








Instance Segmentation and Tracking with Cosine Embeddings and Recurrent Hourglass Networks

Christian Payer¹(✉) , Darko Štern² , Thomas Neff¹ , Horst Bischof¹ ,
and Martin Urschler^{2,3} 

¹ Institute of Computer Graphics and Vision, Graz University of Technology,
Graz, Austria

`christian.payer@icg.tugraz.at`

² Ludwig Boltzmann Institute for Clinical Forensic Imaging, Graz, Austria

³ BioTechMed-Graz, Graz, Austria

Abstract. Different to semantic segmentation, instance segmentation assigns unique labels to each individual instance of the same class. In this work, we propose a novel recurrent fully convolutional network architecture for tracking such instance segmentations over time. The network architecture incorporates convolutional gated recurrent units (ConvGRU) into a stacked hourglass network to utilize temporal video information. Furthermore, we train the network with a novel embedding loss based on cosine similarities, such that the network predicts unique embeddings for every instance throughout videos. Afterwards, these embeddings are clustered among subsequent video frames to create the final tracked instance segmentations. We evaluate the recurrent hourglass network by segmenting left ventricles in MR videos of the heart, where it outperforms a network that does not incorporate video information. Furthermore, we show applicability of the cosine embedding loss for segmenting leaf instances on still images of plants. Finally, we evaluate the framework for instance segmentation and tracking on six datasets of the ISBI celltracking challenge, where it shows state-of-the-art performance.

Keywords: Cell · Tracking · Segmentation · Instances · Recurrent Video · Embeddings

1 Introduction

Instance segmentation plays an important role in biomedical imaging tasks like cell migration, but also in computer vision based tasks like scene understanding. It is considerably more difficult than semantic segmentation (e.g., [10]), since instance segmentation does not only assign class labels to pixels, but also distinguishes between instances within each class, e.g., each individual person on an image from a surveillance camera is assigned a unique ID.

This work was supported by the Austrian Science Fund (FWF): P28078-N33.

© Springer Nature Switzerland AG 2018

A. F. Frangi et al. (Eds.): MICCAI 2018, LNCS 11071, pp. 3–11, 2018.

https://doi.org/10.1007/978-3-030-00934-2_1

Mainly due to the high performance of the U-Net [12], semantic segmentation has been successfully used as a first step in medical instance segmentation tasks, e.g., cell tracking. However, for instances to be separated as connected components during postprocessing, borders of instances have to be treated with special care. In the computer vision community, many methods for instance segmentation have in common that they solely segment one instance at a time. In [4], all instances are first detected and independently segmented, while in [11], recurrent networks are used to memorize which instances were already segmented. Segmenting solely one instance at a time can be problematic when hundreds of instances are visible in the image, as often is the case with e.g., cell instance segmentation. Recent methods are segmenting each instance simultaneously, by predicting embeddings for all pixels at once [5, 8]. These embeddings have similar values within an instance, but differ among instances. In the task of cell segmentation and tracking, temporal information is an important cue to establish coherence between frames, thus preserving instances throughout videos. Despite improvements of instance segmentation using embeddings, to the best of our knowledge, combining them with temporal information for tracking instance segmentations has not been presented.

In this paper, we propose to use recurrent fully convolutional networks for embedding-based instance segmentation and tracking. To memorize temporal information, we integrate convolutional gated recurrent units (ConvGRU [2]) into a stacked hourglass network [9]. Furthermore, we use a novel embedding loss based on cosine similarities, where we exploit the four color map theorem [1], by requiring only neighboring instances to have different embeddings.

2 Instance Segmentation and Tracking

Figure 1 shows our proposed framework on a cell instance segmentation and tracking example. To distinguish cell instances, they are represented as embeddings at different time points. By representing temporal sequences of embeddings in a recurrent hourglass network, a predictor can be learnt from the data, which allows tracking of embeddings also in the case of mitosis events. To finally generate instance segmentations, clustering of the predicted embeddings is performed.

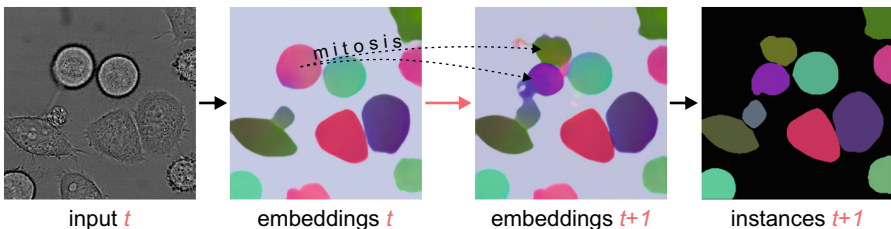


Fig. 1. Overview of our proposed framework showing input image, propagation of cosine embeddings from frame t to frame $t + 1$ (three randomly chosen embedding dimensions as RGB channels), and resulting clustered instances.

2.1 Recurrent Stacked Hourglass Network

We modify the stacked hourglass architecture [9] by integrating ConvGRU [2] to propagate temporal information, as shown in Fig. 2. Differently from the original stacked hourglass network, we use single convolution layers with 3×3 filters and 64 outputs for all blocks in the contracting and expanding paths, while we use ConvGRU with 3×3 filters and 64 outputs in between paths. As proposed by [9], we also stack two hourglasses in a row to improve network predictions. Therefore, we concatenate the output of the first hourglass with the input image to use it as input for the second hourglass. We apply the loss function on the outputs of both hourglasses, while we only use the outputs of the second hourglass for the clustering of embeddings.

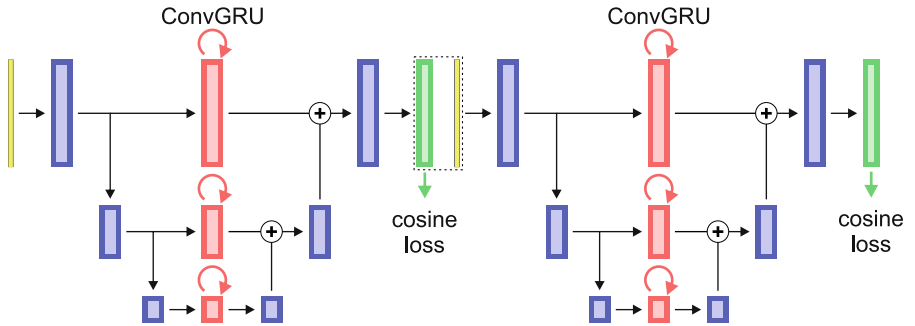


Fig. 2. Overview of the recurrent stacked hourglass network with two hourglasses and three levels. Yellow bars: input; blue boxes: convolutions; red boxes: ConvGRU; dashed black box: concatenation; green boxes: embeddings.

2.2 Cosine Embedding Loss

We let the network predict a d -dimensional embedding vector $\mathbf{e}_p \in \mathbb{R}^d$ for each pixel p of the image. To separate instances $i \in \mathbb{I}$, firstly, embeddings of pixels $p \in \mathbb{S}^{(i)}$ belonging to the same instance i need to be similar, and secondly, embeddings of $\mathbb{S}^{(i)}$ need to be dissimilar to embeddings of pixels $p \in \mathbb{S}^{(j)}$ of other instances $j \neq i$. Here, we treat background as an independent instance. Following from the four color map theorem [1], only neighboring instances need to have different embeddings. Thus, we relax the need of dissimilarity between different instances only to the neighboring ones, i.e., $\mathbb{N}^{(i)} = \bigcup_j \mathbb{S}^{(j)}$ for all instances $j \neq i$ within pixel-wise distance $r_{\mathbb{N}}$ to instance i . This relaxation simplifies the problem by assigning only a limited number of different embeddings to a possibly large number of different instances.

We compare two embeddings with the cosine similarity

$$\cos(\mathbf{e}_1, \mathbf{e}_2) = \frac{\mathbf{e}_1 \cdot \mathbf{e}_2}{\|\mathbf{e}_1\| \|\mathbf{e}_2\|}, \quad (1)$$

which ranges from -1 to 1 , while -1 indicates the vectors have the opposite, 0 orthogonal, and 1 the same direction. We define the cosine embedding loss as

$$L = \frac{1}{|\mathbb{I}|} \sum_{i \in \mathbb{I}} \left(1 - \frac{1}{|\mathbb{S}^{(i)}|} \sum_{p \in \mathbb{S}^{(i)}} \cos(\bar{\mathbf{e}}^{(i)}, \mathbf{e}_p) \right) + \left(\frac{1}{|\mathbb{N}^{(i)}|} \sum_{p \in \mathbb{N}^{(i)}} \cos(\bar{\mathbf{e}}^{(i)}, \mathbf{e}_p)^2 \right), \quad (2)$$

where the mean embedding of instance i is defined as $\bar{\mathbf{e}}^{(i)} = \frac{1}{|\mathbb{S}^{(i)}|} \sum_{p \in \mathbb{S}^{(i)}} \mathbf{e}_p$. By minimizing L , the first term urges embeddings \mathbf{e}_p of pixels $p \in \mathbb{S}^{(i)}$ to have the same direction as the mean $\bar{\mathbf{e}}^{(i)}$, which is the case when $\cos(\bar{\mathbf{e}}^{(i)}, \mathbf{e}_p) \approx 1$, while the second term pushes embeddings \mathbf{e}_p of pixels $p \in \mathbb{N}^{(i)}$ to be orthogonal to the mean $\bar{\mathbf{e}}^{(i)}$, i.e., $\cos(\bar{\mathbf{e}}^{(i)}, \mathbf{e}_p) \approx 0$.

2.3 Clustering of Embeddings

To get the final segmentations from the predicted embeddings, individual groups of embeddings that describe different instances need to be identified. As the number of instances is not known, we perform this grouping with the clustering algorithm HDBSCAN [3] that estimates the number of clusters automatically. For each dataset, two HDBSCAN parameters have to be adjusted: minimal points m_{pts} and minimal cluster size m_{clSize} . To simplify clustering and still be able to detect splitting of instances, we cluster only overlapping pairs of consecutive frames at a time. Since our embedding loss allows same embeddings for different instances that are far apart, we use both image coordinates and value of the embeddings as data points for the clustering algorithm. After identifying the embedding clusters with HDBSCAN and filtering clusters that are smaller than t_{size} , the final segmented instances for each frame pair are obtained.

For merging the segmented instances in overlapping frame pairs, we identify same instances by the highest intersection over union (IoU) between each segmented instance in the overlapping frame. The resulting segmentations are then upsampled back to the original image size, generating the final segmented and tracked instances.

3 Experimental Setup and Results

We train the networks with TensorFlow¹ and perform on-the-fly data augmentation with SimpleITK². We use hourglass networks with seven levels and an input size of 256×256 , while we scale the input images to fit. All recurrent networks are trained on sequences of ten frames. We refer to the supplementary material for individual training and augmentation parameters, as well as individual values of parameter described in Sect. 2.

Left Ventricle Segmentation: To show that our proposed recurrent stacked hourglass network is able to incorporate temporal information, we perform

¹ <https://www.tensorflow.org/>.

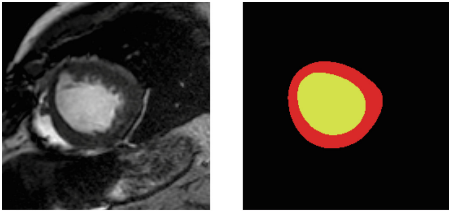
² <http://www.simpleitk.org/>.

semantic segmentation on videos of short-axis MR slices of the heart from the left ventricle segmentation challenge [14]. We compare the recurrent network with a non-recurrent version, where we replace each ConvGRU with a convolution layer to keep the network complexity the same. Since outer slices do not contain parts of the left ventricle, the networks are evaluated on the three central slices that contain both left ventricle myocardium and blood cavity (see Fig. 3a). We train the networks with a softmax cross entropy loss to segment three labels, i.e., background, myocardium, and blood cavity. We use a three-fold cross-validation setup, where we randomly split datasets of 96 patients into three equally sized folds. Table 1a shows the IoU for our internal cross-validation of both recurrent and non-recurrent stacked hourglass networks.

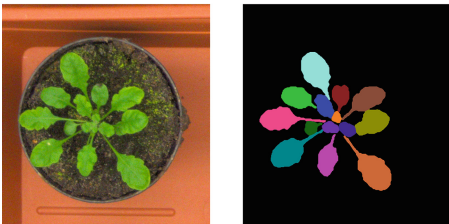
Leaf Instance Segmentation: We show that the cosine embedding loss and the subsequent clustering are suitable for instance segmentation without temporal information, by evaluating on the A1 dataset of the CVPPP challenge for segmenting individual plant leaves [7] (see Fig. 3b). We use the non-recurrent version of the proposed network from the previous experiment to predict embeddings with 32 dimensions. Consequently, the clustering is also performed on single images. As we were not able to provide results on the challenge test set in time before finalizing this paper, we report results of an internal three-fold cross-validation of the 128 training images. In consensus with [13], we report the symmetric best Dice (SBD) and the absolute difference in count ($|\text{DiC}|$) and compare to other methods in Table 1b.

Cell Instance Tracking: As our main experiment, we show applicability of our full framework for instance segmentation and tracking by evaluating six different datasets of cell microscopy videos from the ISBI celltracking challenge [15]. Each celltracking dataset consists of two annotated training videos and two testing videos with image sizes ranging from 512×512 to 1200×1024 and with 48 to 138 frames. We refer to [6] for additional imaging and video parameters. As the instance IDs in groundtruth images are consistent throughout the whole video only for tracking, but not for segmentation, we merge both tracking and segmentation groundtruth for each frame to have consistent instance IDs. Furthermore to learn the background embeddings, we only use the frames on which every cell is segmented. With hyperparameters determined on the two annotated training videos from each dataset, we train the networks for predicting embeddings of size 16 on both videos for our challenge submission.

To compete in the tracking metric of the challenge, the framework is required to identify the parent ID of each cell. As the framework is able to identify splitting cells and to assign new instance IDs (i.e., mitosis as seen on Fig. 1), the parent ID of each newly created instance is determined as the instance with the highest IoU in previous frames. We further postprocess the cells' family tree to be consistent with the evaluation criteria, e.g., an instance ID may not be used after splitting into children. The results in comparison to the top performing methods are presented in Table 2.



(a) Heart MRI input and segmentation.



(b) Plant leaves input and instances.

Fig. 3. Qualitative results of the left ventricle segmentation and the CVPPP leaf instance segmentation. The images on the left side show example inputs, the images on the right side show the predicted segmentations.

Table 1. Quantitative results of the left ventricle segmentation and the CVPPP leaf instance segmentation. Values show mean \pm standard deviation. Note that we report our results for both datasets based on a three-fold cross-validation setup. Thus, they are not directly comparable to other published results. IoU: intersection over union; myo: myocardium; cav: blood cavity; SBD: symmetric best Dice; |DiC|: absolute difference in count.

(a) Quantitative results of the heart MRI left ventricle segmentation.

	IoU _{myo}	IoU _{cav}
non-recurrent	78.3 \pm 9.2	89.1 \pm 7.7
recurrent	79.4 \pm 8.5	89.4 \pm 7.2

(b) Quantitative results of the CVPPP leaf instance segmentation. Values taken from [13].

	SBD	DiC
RIS+CRF	66.6 \pm 8.7	1.1 \pm 0.9
MSU	66.7 \pm 7.6	2.3 \pm 1.6
Nottingham	68.3 \pm 6.3	3.8 \pm 2.0
Wageningen	71.1 \pm 6.2	2.2 \pm 1.6
IPK	74.4 \pm 4.3	2.6 \pm 1.8
IS+RA [11]	84.9 \pm 4.8	0.8 \pm 1.0
Ours	84.5 \pm 5.5	1.5 \pm 1.2

4 Discussion and Conclusion

Up to our knowledge, we are the first to present a method that incorporates temporal information into a network to allow tracking of embeddings for instance segmentation. We perform three experiments to show different aspects of our novel method, i.e., temporal segmentation, instance segmentation, and combined instance segmentation and tracking. Thus, we demonstrate the wide applicability of our approach.

We use the left ventricle segmentation experiment to show that our novel recurrent stacked hourglass network can be used for incorporating temporal information. It can be seen from the results of the experiment that incorporating ConvGRU between contracting and expanding path deeply inside the architecture improves over the baseline stacked hourglass network. Nevertheless, since we simplified the evaluation protocol of the challenge, the results of the experiment should not be directly compared to other reported results. Moreover,

Table 2. Quantitative results of the celltracking datasets for overall performance (OP), segmentation (SEG), and tracking (TRA), as described in [15].

		DIC-HeLa	Fluo-MS-C	Fluo-GOWT1	Fluo-HeLa	PhC-U373	Fluo-SIM+	
OP	1 st	0.864	0.759	0.951	0.942	0.951	0.882	Ours
	2 nd	0.828	0.676	0.914	0.940	0.896	0.878	BGU-IL (1-2)
	3 rd	0.629	0.658	0.902	0.928	0.895	0.874	CUNI-CZ
			5 th 0.631		11 th 0.829	4 th 0.888	9 th 0.810	CVUT-CZ
SEG	1 st	0.814	0.645	0.927	0.903	0.920	0.802	FR-Be-GE
	2 nd	0.776	0.590	0.893	0.893	0.832	0.791	FR-Ro-GE
	3 rd	0.464	0.582	0.887	0.869	0.826	0.781	HD-Har-GE
			5 th 0.496	4 th 0.880	10 th 0.749	5 th 0.793	8 th 0.718	KIT-GE
TRA	1 st	0.915	0.873	0.976	0.991	0.983	0.975	KTH-SE (1-4)
	2 nd	0.881	0.765	0.947	0.987	0.981	0.961	LEID-NL
	3 rd	0.797	0.763	0.925	0.986	0.977	0.957	
					12 th 0.909		10 th 0.902	

benefits of such deep incorporation compared to having recurrent layers on other positions in the network [11] remain to be shown.

This paper also contributes with a novel embedding loss based on cosine similarities. Most of the methods that use embeddings for differentiating between instance segmentations are based on maximizing distances of embeddings in the Euclidean space, e.g., [8]. When using such embedding losses, we observed problems when combining them with recurrent networks, presumably due to unrestricted embedding values. To overcome these problems, we use cosine similarities that normalize embeddings. The only other work that suggests cosine similarities for instance segmentation with embeddings is the unpublished work of [5]. However, compared to their embedding loss that takes all instances into account, our novel loss focuses only on neighboring ones, which can be beneficial for optimization in the case of a large number of instances. We evaluate our novel loss on the CVPPP challenge dedicated to instance segmentation from still images. While waiting for the results of the competition, our method evaluated with three-fold cross-validation shows to be in line with the currently leading method, and has a significant margin to the second best. Moreover, compared to the leading method [11], the architecture of our method is considerably simpler.

In our main experiment for segmentation and tracking of instances, we evaluate our method on the ISBI celltracking challenge, showing large variability in visual appearance, size and number of cells. Our method achieves two first and two second places among the six submitted datasets in the tracking metric. For the dataset DIC-HeLa, having a dense layout of cells as seen in Fig. 1, we outperform all other methods in both tracking and segmentation metrics. On the dataset Fluo-GOWT1 we rank overall second. On the datasets Fluo-HeLa and Fluo-SIM+, which consist of images with small cells, our method does not perform well due to the need to downsample images for the network to process them. When the downsampling results in drastic reduction of cell sizes, our

method fails to create instance segmentations, thus explaining the not satisfying performance also in tracking. To increase the resolution and consequently improve segmentation and tracking, we could split the input image into multiple smaller parts, similarly as done in [12].

In conclusion, our work has shown that embeddings for instance segmentation can be successfully combined with recurrent networks incorporating temporal information to perform instance tracking. In future work, we will investigate the possibility of incorporating the required clustering step inside of a single end-to-end trained network, which could simplify the framework and further improve the segmentation and tracking results.

References

1. Appel, K., Haken, W.: Every planar map is four colorable. *Bull. Am. Math. Soc.* **82**(5), 711–712 (1976)
2. Ballas, N., Yao, L., Pal, C., Courville, A.: Delving deeper into convolutional networks for learning video representations. In: *International Conference on Learning Representations*. CoRR, abs/1511.06432 (2016)
3. Campello, R.J.G.B., Moulavi, D., Zimek, A., Sander, J.: Hierarchical density estimates for data clustering, visualization, and outlier detection. *ACM Trans. Knowl. Discov. Data* **10**(1), 5:1–5:51 (2015)
4. He, K., Gkioxari, G., Dollár, P., Girshick, R.: Mask R-CNN. In: *Proceedings of the International Conference on Computer Vision*, pp. 2980–2988 (2017)
5. Kong, S., Fowlkes, C.: Recurrent pixel embedding for instance grouping. CoRR, abs/1712.08273 (2017)
6. Maška, M., Ulman, V., Svoboda, D., Matula, P., Matula, P., Ederra, C., et al.: A benchmark for comparison of cell tracking algorithms. *Bioinformatics* **30**(11), 1609–1617 (2014)
7. Minervini, M., Fischbach, A., Scharr, H., Tsaftaris, S.A.: Finely-grained annotated datasets for image-based plant phenotyping. *Pattern Recogn. Lett.* **81**, 80–89 (2016)
8. Newell, A., Huang, Z., Deng, J.: Associative embedding: end-to-end learning for joint detection and grouping. In: *Advances in Neural Information Processing Systems*, pp. 2277–2287. Curran Associates, Inc. (2017)
9. Newell, A., Yang, K., Deng, J.: Stacked hourglass networks for human pose estimation. In: Leibe, B., Matas, J., Sebe, N., Welling, M. (eds.) *ECCV 2016*. LNCS, vol. 9912, pp. 483–499. Springer, Cham (2016). https://doi.org/10.1007/978-3-319-46484-8_29
10. Payer, C., Štern, D., Bischof, H., Urschler, M.: Multi-label whole heart segmentation using CNNs and anatomical label configurations. In: Pop, M., et al. (eds.) *STACOM 2017*. LNCS, vol. 10663, pp. 190–198. Springer, Cham (2018). https://doi.org/10.1007/978-3-319-75541-0_20
11. Ren, M., Zemel, R.S.: End-to-end instance segmentation with recurrent attention. In: *Proceedings of the Computer Vision and Pattern Recognition*, pp. 6656–6664 (2017)
12. Ronneberger, O., Fischer, P., Brox, T.: U-Net: convolutional networks for biomedical image segmentation. In: Navab, N., Hornegger, J., Wells, W.M., Frangi, A.F. (eds.) *MICCAI 2015*. LNCS, vol. 9351, pp. 234–241. Springer, Cham (2015). https://doi.org/10.1007/978-3-319-24574-4_28

13. Scharr, H., et al.: Leaf segmentation in plant phenotyping: a collation study. *Mach. Vis. Appl.* **27**(4), 585–606 (2016)
14. Suinesiaputra, A., Cowan, B.R., Al-Agamy, A.O., Elattar, M.A., Ayache, N., et al.: A collaborative resource to build consensus for automated left ventricular segmentation of cardiac MR images. *Med. Image Anal.* **18**(1), 50–62 (2014)
15. Ulman, V., Maška, M., Magnusson, K.E., Ronneberger, O., Haubold, C., Harder, N.: An objective comparison of cell-tracking algorithms. *Nat. Methods* **14**(12), 1141–1152 (2017)