

CELLTRACK R-CNN: A NOVEL END-TO-END DEEP NEURAL NETWORK FOR CELL SEGMENTATION AND TRACKING IN MICROSCOPY IMAGES

Yuqian Chen¹, Yang Song², Chaoyi Zhang¹, Fan Zhang³, Lauren O’Donnell³
Wojciech Chrzanowski^{4,5}, Weidong Cai¹

¹ School of Computer Science, University of Sydney, Australia

² School of Computer Science and Engineering, University of New South Wales, Australia

³ Brigham and Women’s Hospital, Harvard Medical School, USA

⁴ Sydney Pharmacy School, University of Sydney, Australia

⁵ Sydney Nano Institute, University of Sydney, Australia

ABSTRACT

Cell segmentation and tracking in microscopy images are of great significance to new discoveries in biology and medicine. In this study, we propose a novel approach to combine cell segmentation and cell tracking into a unified end-to-end deep learning based framework, where cell detection and segmentation are performed with a current instance segmentation pipeline and cell tracking is implemented by integrating Siamese Network with the pipeline. Besides, tracking performance is improved by incorporating spatial information into the network and fusing spatial and visual prediction. Our approach was evaluated on the DeepCell benchmark dataset. Despite being simple and efficient, our method outperforms state-of-the-art algorithms in terms of both cell segmentation and cell tracking accuracies.

Index Terms— Cell segmentation, cell tracking, deep learning, end-to-end, Siamese Network, spatial information

1. INTRODUCTION

Quantitative analysis of microscopy images provides valuable information for understanding cell structures and cell behaviors, which is of great significance in biology and medicine study. To analyze cell behaviors, it is essential for clinical researches to perform cell segmentation and tracking to identify individual cells and follow them over time [1]. However, the low quality of microscopy images presents special challenges to these tasks.

In recent years, deep learning has proved to be a powerful tool of feature extraction and demonstrated its successful applications in biomedical image analysis. For cell or nuclei segmentation, many studies adopted semantic segmentation architectures such as U-Net to predict foreground and background areas, followed by post-processing procedures [2]. These methods targeted at capturing details but lacked sufficient object-level information, leading to difficulty in

separating individual instances. A few studies used detection based methods and performed segmentation on each detected instance, such as Mask R-CNN [3][4]. These methods demonstrated superior performance in identifying individual instances and no extra post-processing is needed.

The improved detection and segmentation performance of deep learning benefits the subsequent cell tracking task. Most of current tracking methods are conventional algorithms such as overlap intersection-over-union [5] and Viterbi [6], performed independently from detection step. These methods tend to have poor generalization ability and necessitate tuning numerous parameters. Several methods have been proposed to adopt deep learning approaches in cell tracking task [7-13]. In [7], a motion model and a classification neural network were combined for cell tracking. [8] and [9] attempted to achieve joint cell detection and tracking by predicting cell position likelihood and motion map with a neural network, but they were unable to generate segmentation masks. [10], [11] and [12] used deep learning techniques for both cell segmentation and tracking but these processes were executed in sequence with two neural networks trained separately. [13] performed cell segmentation and tracking within a single Recurrent Hourglass Network but showed limited performance and required complex post-processing steps.

In this study, we propose a novel end-to-end neural network framework, dubbed as CellTrack R-CNN, for concurrent cell segmentation and tracking. These tasks are jointly performed by integrating a Siamese tracking branch with the Mask R-CNN pipeline, without any extra post-processing techniques needed. To further enhance tracking performance, spatial information is incorporated into the tracking branch to reach learnable relative position encodings of cell instances, which are then effectively fused with visual features. We evaluated our method on the DeepCell benchmark dataset [10] with state-of-the-art algorithms, and the results demonstrate superior performance of our CellTrack R-CNN in both cell segmentation and tracking.

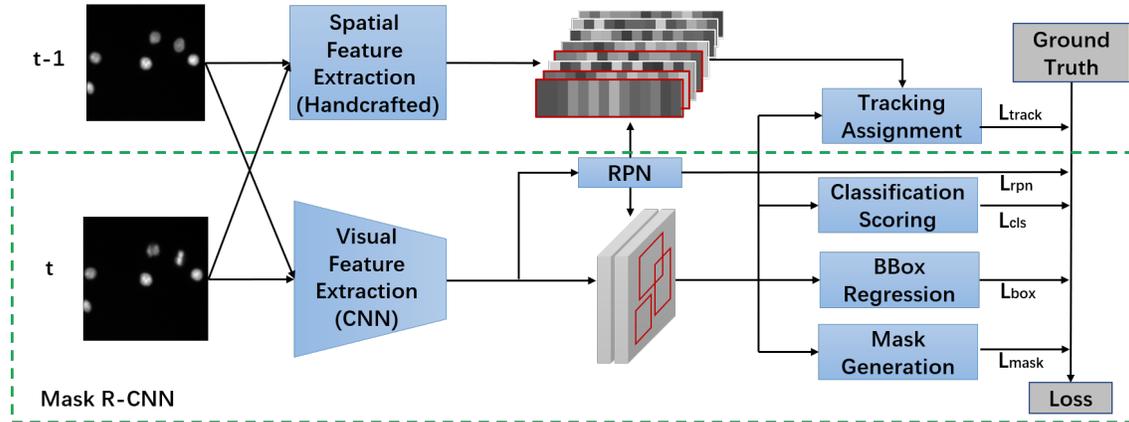


Fig. 1. Overall pipeline of our proposed CellTrack R-CNN. Cell detection and segmentation are performed by Mask R-CNN, the architecture in the green box. A Siamese Network branch is integrated into the Mask R-CNN pipeline for tracking. Visual and spatial features of cell instances in two adjacent frames are extracted and put into the tracking head.

2. METHODS

2.1. Overall Pipeline

The overview of our proposed CellTrack R-CNN framework is shown in Fig. 1. We employ Mask R-CNN [14] as our backbone network for conventional cell detection and segmentation. In addition to the three original branches for predicting classification scores, bounding boxes and segmentation masks, we design a fourth branch with a Siamese Neural Network [15] to achieve cell tracking across frames. Spatial features are extracted in the form of relative position encodings and further aggregated with visual features for better tracking performance. To perform tracking task incorporating temporal information, our framework takes as inputs two adjacent frames to form a training mini-batch.

2.2. Relative Position Encoding

To perform tracking task, we need to identify the same and different instances across frames. In the original architecture of Mask R-CNN, only visual features are extracted with convolutional neural networks (CNN). However, in microscopy images, different cells tend to present similar appearance, making visual information insufficient for cell identification.

Therefore, we propose to incorporate spatial information into our CellTrack R-CNN for cell tracking. Inspired by [16] which studied relative position encoding in NLP tasks, the spatial feature of each cell is presented as relative position between the target cell and its neighboring ones. Specifically, the coordinate differences in x and y dimension between the center of target cell and its nearest n cells are concatenated in sequence as a vector. In this way, a spatial feature vector with length of $2*n$ is generated for each cell instance. Parameter n is set according to the number of cell instances in images. Spatial vectors are normalized before put into tracking head.

2.3. Tracking Branch

Similar to the Mask R-CNN framework, in the first stage, a set of candidate boxes is generated from the region proposal network (RPN). Then the proposed box candidates are used to crop visual features and extract spatial features of the corresponding cell instance. In the second stage, we add a tracking head in parallel to the other three branches. Then the obtained visual and spatial features are fed into the tracking head.

Here we adopt the Siamese Neural Network [15] in the tracking head. Siamese Network is used to calculate the similarity of two inputs. Two different vectors are put into two neural networks with shared weights and output a similarity score. Our tracking branch is composed of two tracking heads for comparing visual and spatial features respectively, as shown in Fig. 2. In each head, the input vector pair is first transformed to two 256-D vectors with a fully-connected layer, then their difference is calculated and finally a similarity score is obtained with another fully-connected layer.

In the first stage, we have obtained features of candidate proposals in frames t and $t-1$. For every candidate in frame t , we calculate its similarity scores with candidates in frame $t-1$ respectively with the Siamese tracking head. During the training process, if the candidate pair is instances of the same cell, the label is set to 1, 0 otherwise. The output similarity scores within $[0, 1]$ are used to calculate losses with a cross entropy function. Tracking head is trained in a way that the same instances have scores close to 1 while different ones close to 0. The tracking loss is defined as the average of visual feature loss and spatial feature loss. Our network is trained in an end-to-end fashion. We add the tracking loss into the total losses of our CellTrack R-CNN, which now becomes $L = L_{rpn} + L_{cls} + L_{box} + L_{mask} + L_{track}$.

For inference, each frame in the testing images is processed in sequence with our pipeline. Both mask segmenta-

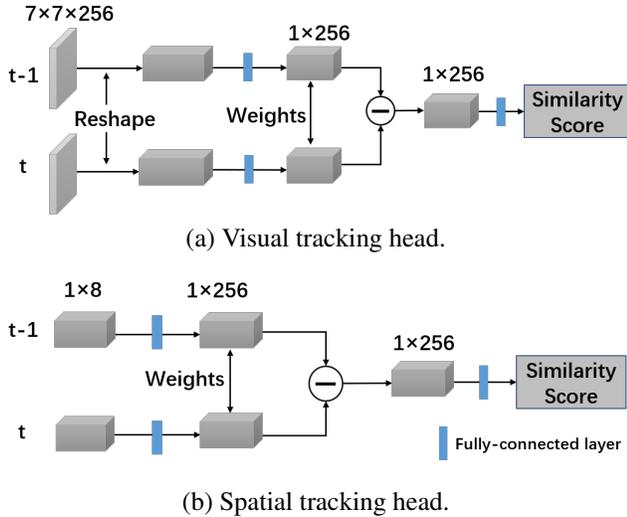


Fig. 2. Architecture of tracking branch

tion and cell tracking are performed on the detected instances. In the tracking task, all instances of the first frame are regarded as new instances of the video and assigned with different tracking ids. To identify the same instance across frames, instances of the current frame are matched to those of the last frame by calculating their visual and spatial similarity scores. After getting the scores, we apply a novel prediction fusion strategy. If the maximum visual and spatial scores correspond to the same cell instance, we assign its tracking id to the target cell. Otherwise, we choose the cell with the maximum overlapping area with the target cell. If the overlapping area is smaller than a threshold α , we regard it as a new instance and assign it a new tracking id. If two target cells match to one cell at the last frame, we regard it as cell mitosis. In this way, our method is able to match cell instances across frames and detect appearing, vanishing cells and cell mitosis. After processing all frames, our method generates a set of instance hypothesis and obtains bounding box, classification score, segmentation mask and tracking id for each instance.

3. EXPERIMENTS AND RESULTS

3.1. Dataset and Evaluation Metrics

We evaluated our method on the public RAW264.7 dataset in the DeepCell database [10]. It is an annotated dataset specific to live-cell imaging, including fluorescence images of the cell nucleus. This dataset provides sufficient image sequences which is important for deep learning-based cell tracking algorithms. Besides, all images in this dataset are fully annotated with detection, segmentation and tracking ground truths. The RAW264.7 dataset contains 13 image sequences and there are 30 images in each of them. We split it into three parts with 8 sequences for training, 2 for validation and 3 for testing.

Table 1. Quantitative comparisons of performance on DeepCell dataset. STH: spatial tracking head. D-T: difference between detection and tracking accuracy.

Method	SEG	DET	TRA	D-T
DeepCell [10]	84.63	95.96	95.46	0.50
FR-Ro-GE [2]	81.82	95.12	94.99	0.13
Ours w/o STH	85.22	96.93	94.19	2.74
Our method	85.26	97.14	97.05	0.09

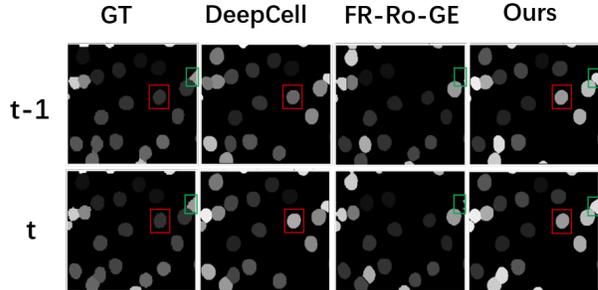


Fig. 3. Visualization of experiment results. Red and green boxes indicate wrong trackings of DeepCell and missed detections of FR-Ro-GE, while our method performed correctly.

For detection and tracking, we adopted graph-based metric [1] for evaluation, which was also used in the DeepCell [10] study. It represents cell lineage in a form of directed acyclic graph (DAG) and generates a tracking score by comparing DAG of generated results and ground truths. Segmentation performance was evaluated by calculating Jaccard Similarity Index of matching objects in generated results and ground truths. These metrics were calculated with publicly available command-line software packages [1]. We compared our method with state-of-the-art methods using these metrics.

3.2. Implementation

We adopted ResNet-50-FPN as the backbone of our CellTrack R-CNN and used the pretrained Mask R-CNN model on MSCOCO Dataset [17]. Our model was trained end-to-end for 40 epochs. We used SGD optimizer with momentum of 0.9 for training with a learning rate of 0.0025. The original images were resized to (950, 800) as input of the neural network. The hyperparameter n in relative position encoding was set to 4 for this dataset because the minimum number of cells in an image was 5. Hyperparameter α during inference was chosen to be 0.1 according to results of the validation set.

3.3. Results and Discussion

In this study, we compared our method with two state-of-the-art studies [2][10]. DeepCell [10] algorithm is the baseline method of this dataset and FR-Ro-GE [2] achieves the best

tracking performance on the ISBI Cell Tracking Challenge [1]. In DeepCell, we also used the segmentation pipeline of Mask R-CNN for better comparison of tracking performance.

We calculated the detection, segmentation and tracking accuracies for comparison. Considering that tracking scores were obtained from detection results, the difference between detection and tracking accuracy was also calculated to clarify the tracking performance. As shown in Table 1, our method achieves the best segmentation, detection and tracking performance. Besides, the minimum difference between detection and tracking score further demonstrates the best performance of our cell tracking method. Visualization of cell segmentation and tracking results are shown in Fig. 3.

An important reason that our method outperforms the others is the fusion of spatial information with visual information in the tracking branch. An ablation study was performed to prove the effectiveness of the spatial tracking head. The results in Table 1 indicate that by incorporating spatial information, the tracking performance shows obvious improvement. In fact, cell tracking is specially challenging in microscopy images due to similar appearance of different cells and severe shape deformations over time, making visual information insufficient to identify individual cells across frames. In this situation, spatial information, such as the relative position distribution of neighboring cells, could make great contributions to cell tracking, which has been confirmed by our study.

In addition to the superior performance, the advantages of our method are its simplicity and efficiency. As we know, most current methods perform cell segmentation and tracking separately and they often require complex and time-consuming post-processing steps. For example, FR-Ro-GE adopted deep learning method only for segmentation and used a separate conventional algorithm for tracking. DeepCell trained two separate neural networks for segmentation and tracking, followed by the conventional tracking algorithm linear programming. However, our method is an end-to-end pipeline. Effective cell tracking is achieved at the price of small increase of model and computation complexity based on an instance segmentation framework.

4. CONCLUSION

To our knowledge, CellTrack R-CNN is the first proposed framework that performs concurrent cell segmentation and tracking within a unified neural network without any post-processing needed. These tasks are fulfilled by integrating Siamese Network with Mask R-CNN pipeline. Spatial information represented as relative position encodings of cell instances is incorporated into the network and proves great benefits to the tracking performance. Our method shows state-of-the-art performance in terms of cell segmentation and tracking. In the future, we will explore more effective strategies for fusing visual and spatial information and investigate performance of our method under various clinical situations.

5. COMPLIANCE WITH ETHICAL STANDARDS

No ethical approval was required for using public dataset.

6. ACKNOWLEDGMENTS

No funding was received for this study. The authors have no relevant financial or non-financial interests to disclose.

7. REFERENCES

- [1] M. Maška, V. Ulman, D. Svoboda, et al., “A benchmark for comparison of cell tracking algorithms,” *Bioinformatics*, vol. 30, no. 11, pp. 1609–1617, 2014.
- [2] O. Ronneberger, P. Fischer, and T. Brox, “U-net: Convolutional networks for biomedical image segmentation,” in *MICCAI*. Springer, 2015, pp. 234–241.
- [3] H. Tsai, J. Gajda, T. Sloan, et al., “Usiigaci: Instance-aware cell tracking in stain-free phase contrast microscopy enabled by machine learning,” *SoftwareX*, vol. 9, pp. 230–237, 2019.
- [4] D. Liu, D. Zhang, Y. Song, C. Zhang, F. Zhang, L. O’Donnell, and W. Cai, “Nuclei segmentation via a deep panoptic model with semantic feature fusion,” in *IJCAI*, 2019, pp. 861–868.
- [5] E. Bochinski, V. Eiselein, and T. Sikora, “High-speed tracking-by-detection without using image information,” in *AVSS*. IEEE, 2017, pp. 1–6.
- [6] K. Magnusson, J. Jaldén, P. M Gilbert, and H. Blau, “Global linking of cell tracks using the viterbi algorithm,” *IEEE transactions on medical imaging*, vol. 34, no. 4, pp. 911–929, 2014.
- [7] T. He, H. Mao, J. Guo, and Z. Yi, “Cell tracking using deep neural networks with multi-task learning,” *Image and Vision Computing*, vol. 60, pp. 142–153, 2017.
- [8] J. Hayashida and R. Bise, “Cell tracking with deep learning for cell detection and motion estimation in low-frame-rate,” in *MICCAI*. Springer, 2019, pp. 397–405.
- [9] J. Hayashida, K. Nishimura, and R. Bise, “Mpm: Joint representation of motion and position map for cell tracking,” in *CVPR*, 2020, pp. 3823–3832.
- [10] E. Moen, E. Borba, G. Miller, et al., “Accurate cell tracking and lineage construction in live-cell imaging experiments with deep learning,” *bioRxiv*, p. 803205, 2019.
- [11] J. Lugagne, H. Lin, and M. Dunlop, “Delta: Automated cell segmentation, tracking, and lineage reconstruction using deep learning,” *PLoS computational biology*, vol. 16, no. 4, pp. e1007673, 2020.
- [12] A. Panteli, D. Gupta, N. de Bruin, and E. Gavves, “Siamese tracking of cell behaviour patterns,” in *MIDL*, 2020.
- [13] C. Payer, D. Štern, T. Neff, et al., “Instance segmentation and tracking with cosine embeddings and recurrent hourglass networks,” in *MICCAI*. Springer, 2018, pp. 3–11.
- [14] K. He, G. Gkioxari, P. Dollár, and R. Girshick, “Mask r-cnn,” in *ICCV*, 2017, pp. 2961–2969.
- [15] S. Chopra, R. Hadsell, and Y. LeCun, “Learning a similarity metric discriminatively, with application to face verification,” in *CVPR’05*. IEEE, 2005, vol. 1, pp. 539–546.
- [16] P. Shaw, J. Uszkoreit, and A. Vaswani, “Self-attention with relative position representations,” in *NAACL-HLT (2)*, 2018.
- [17] T. Lin, M. Maire, S. Belongie, et al., “Microsoft coco: Common objects in context,” in *ECCV*. Springer, 2014, pp. 740–755.