# Star-convex Polyhedra for 3D Object Detection and Segmentation in Microscopy

Martin Weigert[1,2,3,⋆]    Uwe Schmidt[2,3,⋆]    Robert Haase[2,3]    Ko Sugawara[4,5]    Gene Myers[2,3]

[1]Institute of Bioengineering, École Polytechnique Fédérale de Lausanne (EPFL), Switzerland
[2]Max Planck Institute of Molecular Cell Biology and Genetics (MPI-CBG), Dresden, Germany
[3]Center for Systems Biology Dresden (CSBD), Germany
[4]Institut de Génomique Fonctionnelle de Lyon (IGFL), École Normale Supérieure de Lyon, France
[5]Centre National de la Recherche Scientifique (CNRS), Paris, France

## Abstract

*Accurate detection and segmentation of cell nuclei in volumetric (3D) fluorescence microscopy datasets is an important step in many biomedical research projects. Although many automated methods for these tasks exist, they often struggle for images with low signal-to-noise ratios and/or dense packing of nuclei. It was recently shown for 2D microscopy images that these issues can be alleviated by training a neural network to directly predict a suitable shape representation (star-convex polygon) for cell nuclei. In this paper, we adopt and extend this approach to 3D volumes by using star-convex polyhedra to represent cell nuclei and similar shapes. To that end, we overcome the challenges of 1) finding parameter-efficient star-convex polyhedra representations that can faithfully describe cell nuclei shapes, 2) adapting to anisotropic voxel sizes often found in fluorescence microscopy datasets, and 3) efficiently computing intersections between pairs of star-convex polyhedra (required for non-maximum suppression). Although our approach is quite general, since star-convex polyhedra include common shapes like bounding boxes and spheres as special cases, our focus is on accurate detection and segmentation of cell nuclei. Finally, we demonstrate on two challenging datasets that our approach (STARDIST-3D) leads to superior results when compared to classical and deep learning based methods.*

## 1. Introduction

Detection and segmentation of cell nuclei in volumetric (3D) fluorescence microscopy images is a ubiquitous problem in developmental biology and often constitutes the first step when studying cellular expression patterns, or when tracing cell lineages in developing organisms [17, 22]. The task of nuclei *detection* is to roughly locate all individual nuclei inside a 3D volume, *e.g.* by enumerating their center points or bounding boxes. On the other hand, *semantic segmentation* aims to label each pixel with a semantic class (*e.g.*, nucleus or background), but is not concerned with discerning individual nuclei. Finally, *instance segmentation* is more ambitious since it combines these tasks by seeking a separate label mask for each individual nucleus. As modern microscopes produce increasingly large 3D datasets, many automated instance segmentation methods have been proposed over the years [17]. These include classical thresholding approaches with pixel-grouping via connected component, morphological methods based on the watershed transform [4, 16, 7], and optimization via graph-cuts [5]. More recently, methods based on deep learning have been shown to vastly improve results for natural and biomedical images alike [11, 24, 25].

In general, deep learning based instance segmentation can be roughly categorized into *1)* methods that first perform semantic segmentation followed by grouping of pixels into distinct objects (*e.g.* U-Net [9, 6]), and *2)* methods that first predict axis-aligned bounding boxes of individual objects with a subsequent semantic segmentation step for each found object (*e.g.* [11, 27, 26]). Despite the advances made by these methods, they often still underperform due to the low signal-to-noise ratios and dense packing of nuclei in typical fluorescence microscopy datasets. In particular, methods of category *1)* are prone to erroneously fuse touching nuclei, and those of category *2)* may fail to discern objects that are poorly approximated with bounding boxes.

These problems have recently been highlighted by Schmidt *et al.* [21] for the case of 2D fluorescence microscopy images. To alleviate these issues, [21] proposed a method called STARDIST, which uses a neural network that directly predicts an appropriate shape representation (star-convex polygons) for cell nuclei and demonstrated improved results. Concretely, for every pixel inside an object
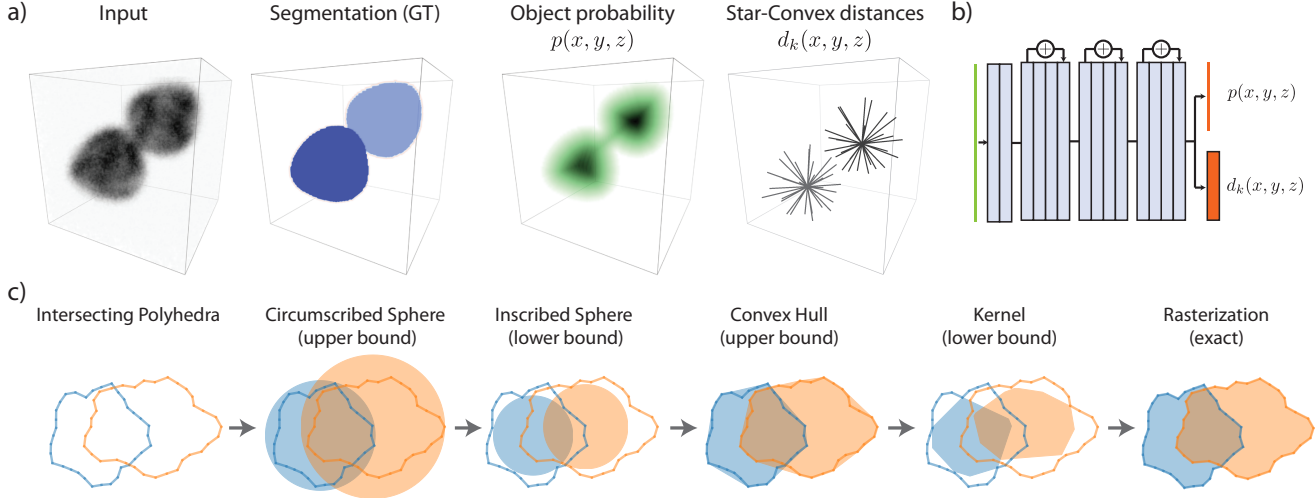
---
⋆Equal contribution

Figure 1: a) The proposed STARDIST-3D method is trained to densely predict object probabilities $p$ and radial distances $d_k$ to object boundaries. b) Schematic of our CNN architecture based on ResNet [12]. c) During non-maximum suppression we use successively tighter bounds to efficiently determine if the intersection volume of two star-convex polyhedra is above a given threshold (only shown in 2D here).

(nucleus) they predict the distance to the object boundary along several radial directions, thus defining a star-convex polygon. Furthermore, they also predict an object probability to determine which pixels are part of cell nuclei and thus are allowed to vote for an object shape. Since every pixel is predicting a polygon to represent the entire shape of the object it belongs to, they perform non-maximum suppression to prune redundant shapes that likely represent the same object. Note that [21] sits somewhere in between object detection and instance segmentation because the predicted shapes are of relatively high fidelity, but are not pixel-accurate.

In this paper, we adopt and extend the approach of [21] to the case of 3D volumes and use *star-convex polyhedra* as shape representations for cell nuclei and similar shapes. We directly predict the polyhedra parameters densely for each pixel and then use non-maximum suppression (NMS) to prune the highly redundant set of obtained polyhedron shapes to ideally retain only one predicted shape for each true object in the image. Please see Fig. 1 for an overview of our approach. Note that we keep the benefits of [21], first and foremost to accurately disambiguate densely packed objects in images with low signal-to-noise ratios. Furthermore, star-convex polygons/polyhedra are a superset of convex shapes in 2D/3D and thus include common shapes like bounding boxes and circles/spheres as special cases.

**Contributions** The extension of [21] from 2D to 3D is challenging and our main contribution in this paper. First, computing the intersection of two star-convex polyhedra (as required for NMS) efficiently is non-trivial (see Section 2.3 and Fig. 1c), but highly necessary to make this approach practical for large 3D volumes. Second, while [21] used 32 radial directions to represent 2D nuclei shapes, a naive ex-

tension to 3D would require $32^2 = 1024$ directions which is not feasible due to the excessive amount of computation and memory required for large 3D volumes. We show that a more judicious selection of radial directions (Section 2 and Fig. 1a) enables faithful shape representations with as little as 64 values. Third, microscopy volumes are commonly acquired with anisotropic voxel sizes that result in squeezed nuclei shapes along the axial (Z) direction. We find that it is critical to adapt the star-convex representation to account for this anisotropy of the data to achieve good results (Sections 2 and 3). Finally, we demonstrate on two challenging datasets that our proposed method (STARDIST-3D) leads to superior results when compared to a classical watershed method and U-Net baselines.

## 2. Method

### 2.1. Star-convex polyhedra

We describe the 3D shape of a single object (cell nucleus) with a *star-convex polyhedron*. Concretely, for each pixel inside an object we compute the distances $d_k$ to the object boundary along a fixed set of $n$ *unit rays* $\vec{r}_k$. To obtain a faithful shape model, we use rays that are approximately evenly distributed on an ellipsoid representative of the objects in a dataset. To that end, we first compute the points $(x_k, y_k, z_k)_{k=0...n-1}$ of a spherical *Fibonacci lattice* [10]

$$z_k = -1 + \frac{2k}{n-1},$$

$$y_k = \sqrt{1 - z_k^2} \sin \left[ 2\pi(1 - \varphi^{-1})k \right],$$

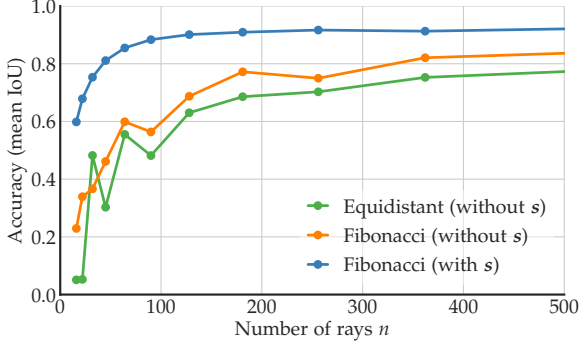$$x_k = \sqrt{1 - z_k^2} \cos \left[ 2\pi(1 - \varphi^{-1})k \right],$$

2

Figure 2: Reconstruction accuracy (mean intersection over union) of ground-truth instances when using different unit rays (Equidistant/Fibonacci) and anisotropy factors $s$ (for dataset PARHYALE).

where $\varphi = \frac{1+\sqrt{5}}{2}$ denotes the golden ratio. To account for anisotropy of the data we generate intermediate, anisotropically scaled vectors $\vec{u}_k = \left(\frac{x_k}{s_x}, \frac{y_k}{s_y}, \frac{z_k}{s_z}\right)$. The respective *anisotropy factor* $\vec{s} = (s_x, s_y, s_z)$ is calculated as the median bounding box size of all objects in the training images. The final unit rays $\vec{r}_k$ are then computed via normalization $\vec{r}_k = \frac{\vec{u_k}}{|\vec{u_k}|}$. The surface of a star-convex polyhedron represented by the distances $d_k$ is then given by its vertices $d_k \cdot \vec{r}_k$ and triangulation induced by the convex hull facets of the unit rays $\vec{r}_k$ (which is a convex set by definition). We generally find that a sufficiently accurate reconstruction of the labeled 3D cell nuclei in our ground-truth (GT) images can be obtained with as few as 64 rays. Fig. 2 shows the reconstruction fidelity for a dataset with highly anisotropic images (PARHYALE, *cf.* Section 3) and highlights the importance of using an appropriate anisotropy factor $\vec{s}$. Note that $\vec{s}$ is automatically computed from the GT images and does not have to be chosen manually. Furthermore, Fig. 2 shows that our ray definition (Fibonacci) is more accurate than using equidistant (polar/azimuthal) distributed rays.

## 2.2. Model

Following [21], we use a convolutional neural network (CNN) to densely predict the star-convex polyhedron representation and a value that indicates how likely a pixel is part of an object. Concretely, for each pixel $(x, y, z)$, the CNN is trained to predict the $n$ *radial distances* $\{d_k(x, y, z)\}_{k=0...n-1}$ to the object boundary as defined above and additionally an *object probability* $p(x, y, z)$ defined as the (normalized) Euclidean distance to the nearest background pixel (Fig. 1a). To save computation and memory we predict at a *grid* of lower spatial resolution than the input image, since a dense (*i.e.*, per input pixel) output is often not necessary (this is similar to the concept of bounding box *anchors* in object detection approaches [19, 18]).

We use a slightly modified 3D variant of ResNet [12] as

a neural network backbone[1] to predict both the radial distances and object probabilities (Fig. 1b). In particular, we use residual blocks with 3 convolution layers of kernel size $3 \times 3 \times 3$. Similar to [12], we start with two convolution layers of kernel sizes $7 \times 7 \times 7$ and $3 \times 3 \times 3$, but without strides to avoid downsampling. This is followed by $m$ residual blocks, where each block only performs downsampling if the spatial resolution is still higher than the prediction grid (see above); we double the number of convolution filters after each downsampling. After the last residual block, we use a single-channel convolution layer with *sigmoid* activation to output the per-pixel[2] object probabilities $p$. The last residual block is additionally connected to an $n$-channel convolution layer to output the radial distances $d_k$. Our code based on Keras/TensorFlow [8, 1] and documentation is available at https://github.com/mpicbg-csbd/stardist.

**Training** Given the pixel-wise object probabilities and distances of the prediction $(\hat{p}, \hat{d}_k)$ and ground-truth $(p, d_k)$, we minimize the following loss function (averaged over all pixels) during training:

$$L(p, \hat{p}, d_k, \hat{d}_k) = L_{obj}(p, \hat{p}) + \lambda_d L_{dist}(p, \hat{p}, d_k, \hat{d}_k). \quad (1)$$

For the object loss $L_{obj}$ we use standard *binary cross-entropy*

$$L_{obj}(p, \hat{p}) = -p \log \hat{p} - (1 - p) \log(1 - \hat{p}). \quad (2)$$

For the distance loss $L_{dist}$ we use the *mean absolute error* weighted by the object probability (active only on object pixels, *i.e.* $p > 0$) and add a regularization term (active only on background pixels, *i.e.* $p = 0$):

$$L_{dist}(p, \hat{p}, d_k, \hat{d}_k) = p \cdot \mathbb{1}_{p>0} \cdot \frac{1}{n} \sum_k |d_k - \hat{d}_k| +$$
$$\lambda_{reg} \cdot \mathbb{1}_{p=0} \cdot \frac{1}{n} \sum_k |\hat{d}_k|. \quad (3)$$

This specific form was chosen to promote increased accuracy for points closer to the object centers (which eventually will be used as polyhedra center candidates).

**Prediction** After the CNN predicts the radial distances $\hat{d}_k$ and object probabilities $\hat{p}$, we collect a set of object *candidates* by only considering radial distances at pixels with object probabilities above a reasonably high threshold, *i.e.* we only retain shapes that very likely belong to an object. Since the set of object candidates is highly redundant, we use non-maximum suppression (NMS) to obtain only one

---

[1]We find that using a U-Net [9] backbone leads to very similar results.
[2]To improve readability we will drop from now on the explicit pixel coordinate $(x, y, z)$ for both $p(x, y, z)$ and $d_k(x, y, z)$.
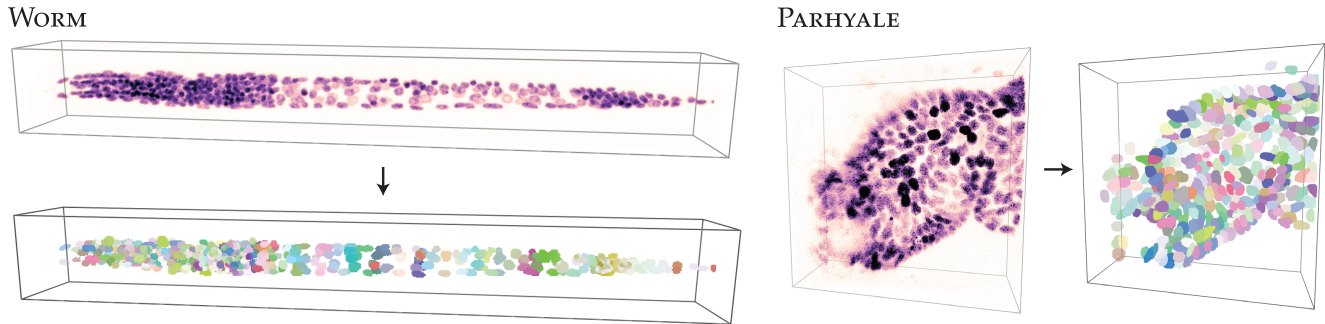
Figure 3: Datasets used in our experiments. Shown are raw input images (purple) and associated ground-truth instance segmentation labels (colored) for a single volume of the WORM (left) and PARHYALE (right) datasets.

shape for every actual object in the image, as is common in object detection (*e.g.*, [11]). Thereby, the object candidate with the highest object probability suppresses all other remaining candidates if they overlap substantially. This process is repeated until there are no further candidates to be suppressed. All remaining (*i.e.* not suppressed) candidates yield the final set of predicted object shapes.

### 2.3. Efficient non-maximum suppression

The NMS step requires to assess the pairwise overlap of a potentially large set of polyhedron candidates ($> 10^4$). Unfortunately, computing the exact intersection volume between two star-convex polyhedra *efficiently* is non-trivial (in contrast to *convex* polyhedra). To address this issue, we employ a filtering scheme that computes as needed successively tighter upper and lower bounds for the overlap of polyhedron pairs (*cf.* Fig. 1c). Concretely, we compute the intersection volume of the respective *i)* bounding spheres (upper bound), *ii)* inscribed spheres (lower bound), *iii)* convex hulls (upper bound), and *iv) kernels*[3] (lower bound). Note that *iii)* and *iv)* involve only the intersection of convex polyhedra and can thus be computed efficiently [3]. If a computed bound is already sufficient to decide whether a candidate should be suppressed or not, no further computation is carried out. Otherwise, we eventually perform the expensive but exact intersection computation by rasterization of both polyhedra. We find that this NMS filtering scheme leads to a noticeable reduction in runtime that makes STARDIST-3D practical (*e.g.* 9 s for a stack of size $1141 \times 140 \times 140$ with 12000 initial candidates).

### 3. Experiments

We consider two qualitatively different datasets (Fig. 3) to validate the efficacy of our approach:

WORM   A subset of 28 images used in Long *et al.* [15], showing DAPI-stained nuclei of the first larval stage

---
[3]The (convex) set of all points that can serve as center of the star-convex polyhedron.

(L1) of *C. elegans* (Fig. 3 left). Stacks are of average size $1157 \times 140 \times 140$ pixels with semi-automatically annotated cell nucleus instances (15148 in total) that underwent subsequent manual curation. We randomly choose 18/3/7 images for training/validation/testing. Note that the images have (near) isotropic resolution.

PARHYALE   A subset of recording #04 of Alwes *et al.* [2], showing *Parhyale hawaiensis* expressing Histone-EGFP (Fig. 3 right). It contains 6 images of $512 \times 512 \times 34$ pixels with manually annotated nucleus instances (1738 in total). We randomly choose 3/1/2 images for training/validation/testing. In contrast to WORM, the images are highly anisotropic in the axial direction. This dataset is more challenging due its substantially lower signal-to-noise ratio (*cf.* Fig. 6). Furthermore, it contains much fewer labeled training images, more akin to what is typical in many biological datasets.

### 3.1. Methods and Evaluation

We compare our proposed STARDIST-3D approach against two kinds of methods (IFT-Watershed [16] and 3D U-Net [9]) that are commonly used for segmentation of fluorescence microscopy images. First, a classical watershed-based method that does not use machine learning. Second, a variant of the popular U-Net with and without more sophisticated postprocessing.

To evaluate the performance of all methods, we use $accuracy(\tau) = \frac{TP}{TP+FN+FP}$ for several overlap thresholds $\tau$. $TP$ are true positives, which are pairs of predicted and ground-truth nuclei having an *intersection over union* (IoU) value $\geq \tau$. $FP$ are false positives (unmatched predicted instances) and $FN$ are false negatives (unmatched ground-truth instances). We use the Hungarian method [14] (implementation from *SciPy* [13]) to compute an optimal matching whereby a single predicted nucleus cannot be assigned to multiple ground-truth instances (and vice versa). Note that a suitable value of $\tau$ depends on the biological application. For example, one would likely use a smaller $\tau < 0.5$
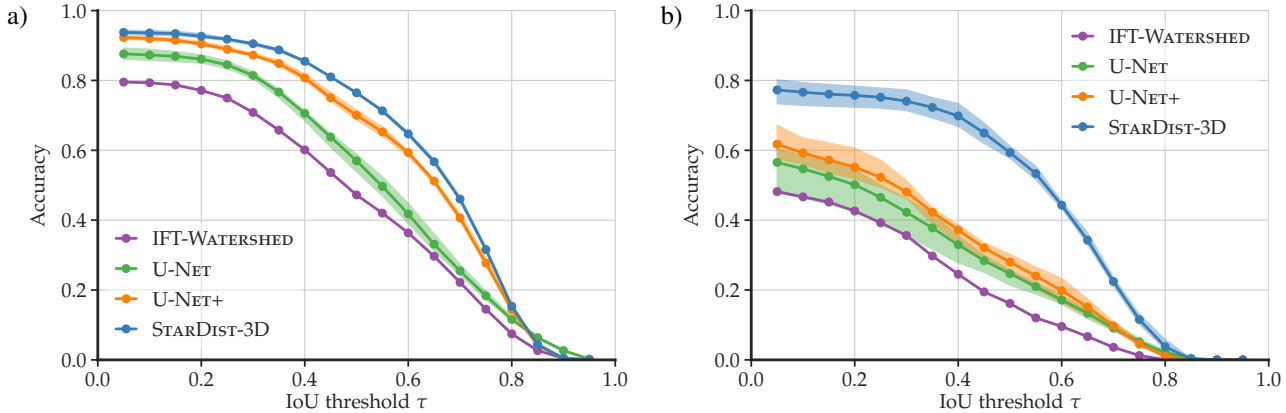
Figure 4: Accuracy for several IoU thresholds $\tau$ for datasets a) WORM and b) PARHYALE. We show the average performance over 5 independent trials for all trained models (shaded regions indicate best and worst result).

for the purpose of counting objects, whereas intensity measurements inside each object would require more accurate shapes and thus demand a higher value of $\tau$.

STARDIST-3D We use STARDIST-3D as explained in Section 2.2 with $n = 96$ radial directions and $m = 3$ residual blocks that start with 32 convolution filters. We predict at a grid half the spatial resolution of the input image, except for the anisotropic Z axis of PARHYALE. We use automatically computed anisotropy factors (*cf.* Section 2.1) of $\vec{s} = (1, 1, 1)$ for WORM and $\vec{s} = (1, 1, 7.1)$ for PARHYALE. We use weights $\lambda_d = 0.1$ and $\lambda_{reg} = 10^{-4}$ for the loss function in Eq. (1).

IFT-WATERSHED The IFT-Watershed [16] is an efficient combination of maxima detection and 3D watershed segmentation. It represents an advanced classical image segmentation method that we know is being used in practice. Concretely, we use the *Interactive Watershed* plugin[4] in *Fiji* [20] and perform extensive parameter tuning (such as Gaussian filter size during preprocessing and maxima detection thresholds) using the training images of each dataset.

U-NET We train a 3D U-Net [9] to classify each pixel into *background*, *nucleus*, and also *nucleus boundary*, as this helps substantially to separate touching nuclei [6]. We threshold the predicted *nucleus* probabilities and group pixels in each connected component to obtain individual nuclei.

U-NET+ We use the same trained 3D U-Net as above, but apply more sophisticated postprocessing. Concretely, we observe improved performance by thresholding the *nucleus* probabilities to obtain seed regions that we

grow (via 3D watershed [23]) until they reach pixels with *background* probability above a second threshold.

We apply random data augmentations during training, including flips, axis-aligned rotations, elastic deformations, intensity rescaling, and noise. After training, thresholds for all methods (as described above) are tuned on validation images to optimize accuracy averaged over $\tau \in \{0.3, 0.5, 0.7\}$.

### 3.2. Results

The results in Table 1 and Fig. 4 show that STARDIST-3D consistently outperforms all other methods that we compared to (note that we report the average result over 5 trials for all trained models). The performance gap between STARDIST-3D and the other methods is especially striking for dataset PARHYALE, which may be explained by STARDIST-3D's shape model being especially helpful to disambiguate between neighboring nuclei in these challenging low-SNR images. In Fig. 6 we show lateral (XY) and axial (XZ) views of segmentation results for both datasets. Here, IFT-WATERSHED often produces imperfect boundaries and erroneous splits, particularly for dataset PARHYALE. This is expected, as the watershed operation uses the input intensities alone without leveraging extracted features. U-NET tends to under-segment the image, generally producing object instances that are too small, as the use of a single threshold for the *nucleus* class leads to a trade-off between object size and avoidance of falsely merged objects. In contrast, U-NET+ exhibits slight over-segmentation, since a larger first threshold produces more (but smaller) objects that are then grown to yield the final instances. Finally, STARDIST-3D produces superior segmentations, although it can sometimes fail to detect some nuclei (especially for dataset PARHYALE). As an additional visualization, we show a 3D rendering of STARDIST-3D segmentation results for both datasets in Fig. 7.

---

[4]https://imagej.net/Interactive_Watershed

| Threshold $\tau$ | 0.10 | 0.20 | 0.30 | 0.40 | 0.50 | 0.60 | 0.70 | 0.80 | 0.90 |
|---|---|---|---|---|---|---|---|---|---|
| | | | | WORM | | | | | |
| IFT-WATERSHED | 0.794 | 0.771 | 0.708 | 0.601 | 0.472 | 0.364 | 0.222 | 0.074 | 0.005 |
| U-NET | 0.873 | 0.861 | 0.814 | 0.706 | 0.570 | 0.418 | 0.255 | 0.116 | **0.027** |
| U-NET+ | 0.920 | 0.905 | 0.872 | 0.807 | 0.700 | 0.593 | 0.406 | 0.144 | 0.005 |
| STARDIST-3D | **0.936** | **0.926** | **0.905** | **0.855** | **0.765** | **0.647** | **0.460** | **0.154** | 0.004 |
| | | | | PARHYALE | | | | | |
| IFT-WATERSHED | 0.467 | 0.426 | 0.356 | 0.245 | 0.161 | 0.096 | 0.036 | 0.000 | 0.000 |
| U-NET | 0.547 | 0.501 | 0.423 | 0.330 | 0.247 | 0.171 | 0.091 | 0.021 | 0.000 |
| U-NET+ | 0.592 | 0.552 | 0.481 | 0.372 | 0.280 | 0.198 | 0.097 | 0.010 | 0.000 |
| STARDIST-3D | **0.766** | **0.757** | **0.741** | **0.698** | **0.593** | **0.443** | **0.224** | **0.038** | 0.000 |

Table 1: Accuracy (average over 5 independent trials for trained models) for several IoU thresholds $\tau$ for datasets WORM and PARHYALE.

Note that we find (not shown) that the accuracy of STARDIST-3D would drop dramatically (for example, from 0.593 to 0.291 for $\tau = 0.5$) if we did not adapt the radial directions to account for the anisotropy of the nuclei shapes (Section 2.1) for PARHYALE. While STARDIST-3D's lead is less pronounced for dataset WORM, this may be due to the higher signal quality of the input images and also the general abundance of labeled cell nuclei available for training and validation (11387 in total). In Fig. 5, we investigate how STARDIST-3D and the other trained models cope with less annotated data by randomly selecting only a partial 3D image slice from each training and validation stack. Interestingly, we find that with only 4.15% of the training and validation data (472 instances in total), STARDIST-3D can for $\tau = 0.5$ reach the same performance (accuracy of 0.7) as U-NET+ with access to 100% of the training and validation data.

## 4. Discussion

We presented STARDIST-3D, an extension of [21] to detect and segment cell nuclei in volumetric fluorescence microscopy images, even when they exhibit substantial anisotropy. Our method outperformed strong watershed and U-Net baselines, yet is easy to train and use, and due to our star-convex polyhedra parameterization and efficient intersection implementation fast enough to process typical large 3D volumes. Furthermore, STARDIST-3D should be generally applicable to segment objects whose shapes are well-represented with star-convex polyhedra.
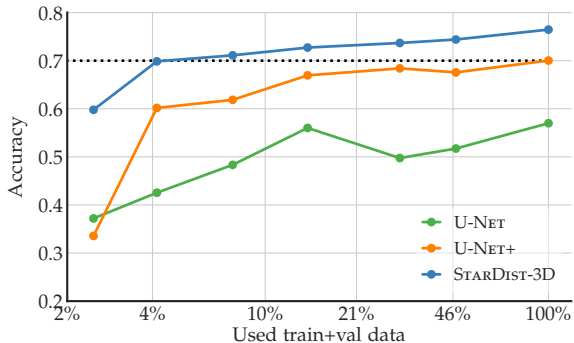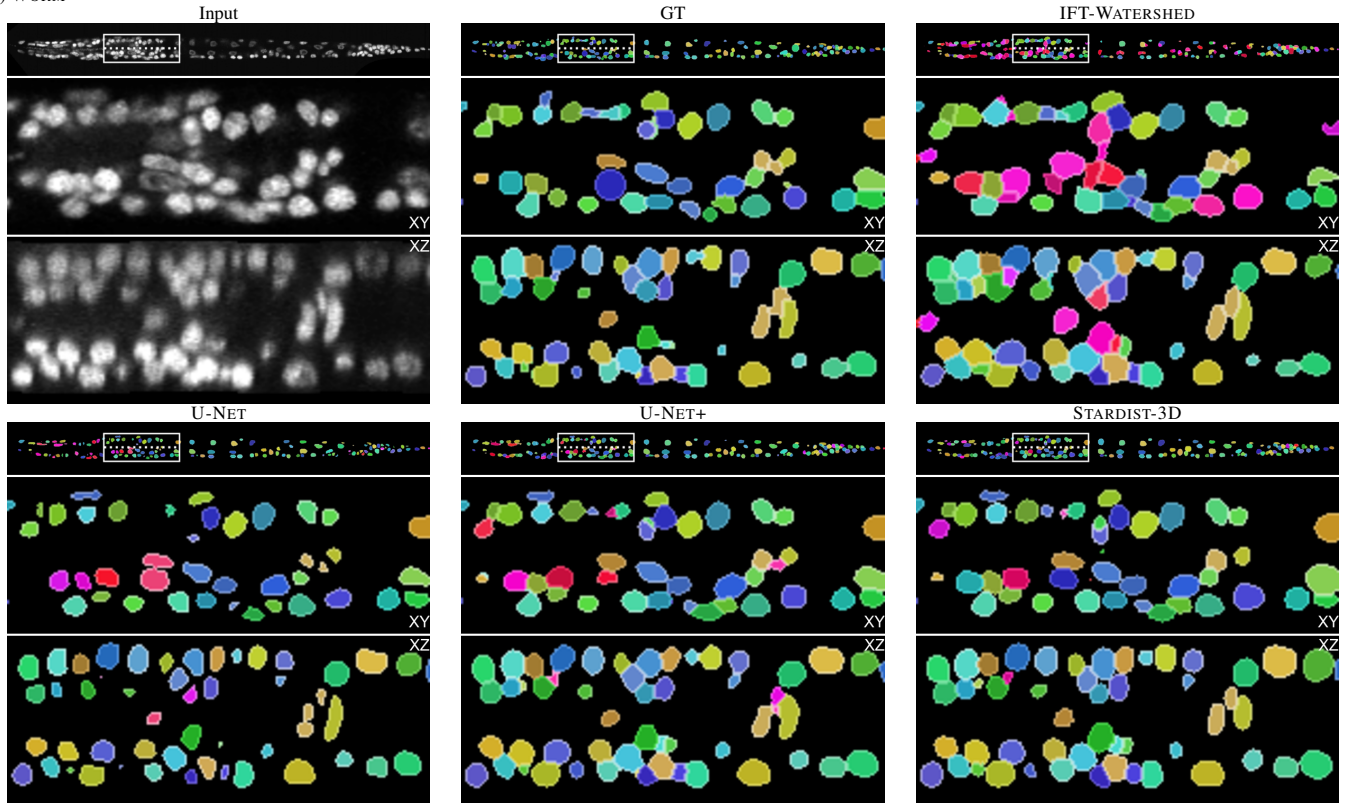
### Acknowledgments

Figure 5: Test accuracy ($\tau = 0.5$) of different methods when using only a fraction of all available training/validation volumes (for dataset WORM).

## References

[1] M. Abadi, P. Barham, J. Chen, Z. Chen, A. Davis, J. Dean, M. Devin, S. Ghemawat, G. Irving, M. Isard, et al. Tensorflow: A system for large-scale machine learning. In *12th USENIX Symposium on Operating Systems Design and Implementation (OSDI 16)*, pages 265–283, 2016.

[2] F. Alwes, C. Enjolras, and M. Averof. Live imaging reveals the progenitors and cell dynamics of limb regeneration. *Elife*, 5, 2016.

[3] C. B. Barber, D. P. Dobkin, D. P. Dobkin, and H. Huhdanpaa. The quickhull algorithm for convex hulls. *ACM Transactions on Mathematical Software (TOMS)*, 22(4):469–483, 1996.

[4] S. Beucher and F. Meyer. The morphological approach to segmentation : The watershed transformation. *Mathematics of Morphology in Image Processing*, pages 433–482, 1993.

[5] Y. Boykov and G. Funka-Lea. Graph cuts and efficient ND image segmentation. *International Journal of Computer Vision*, 70(2), 2006.
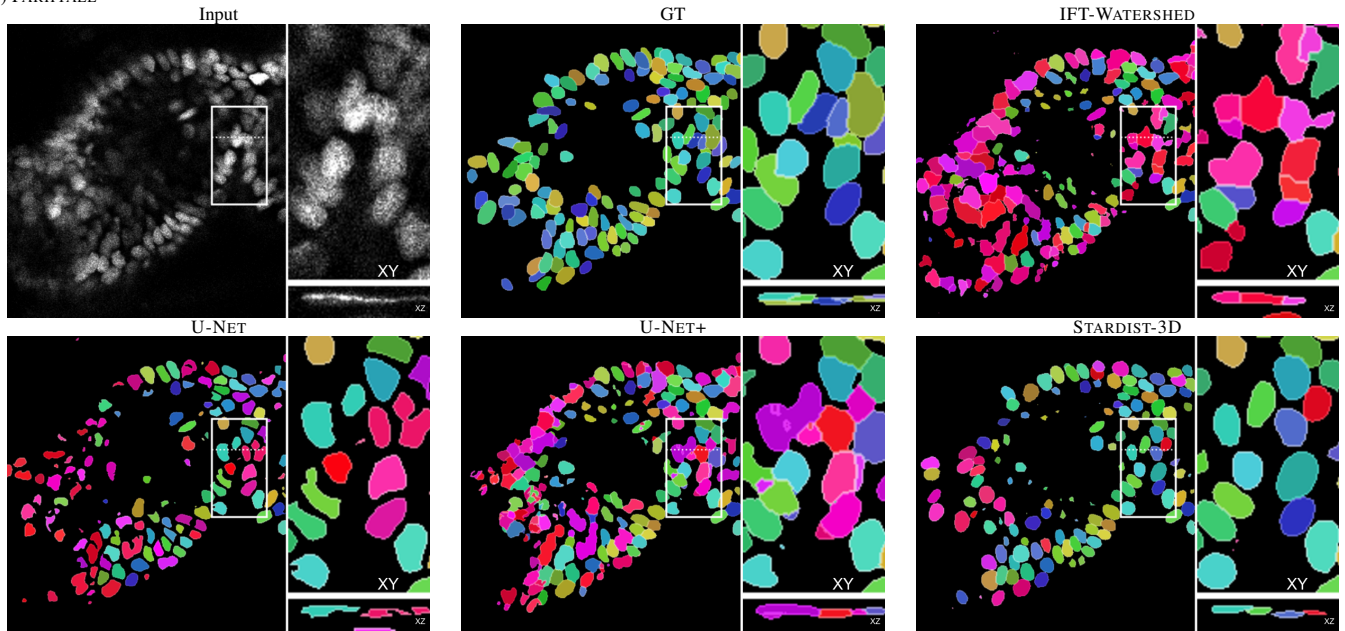
Figure 6: Example results ($\tau = 0.5$) of all methods for both datasets. Colors denote nucleus identities, *i.e.* correct predictions ($TP$) have the same color as the ground-truth (GT). Incorrect predictions ($FP$) are shown in red hues. False negatives ($FN$) are not highlighted. For each inset we show lateral (XY) and axial (XZ, indicated by dotted line) views.
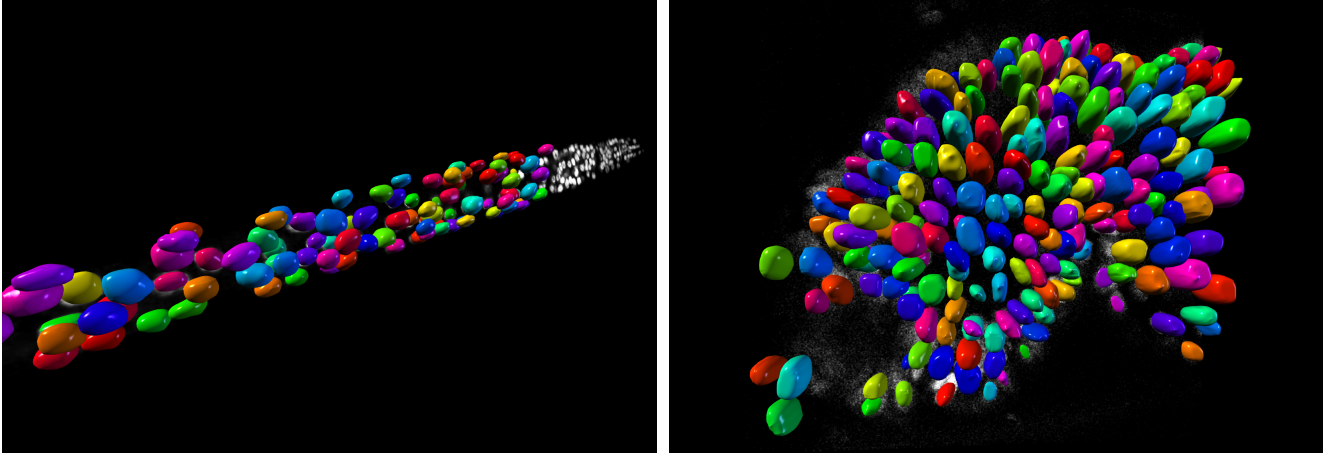
Figure 7: Example results of our STARDIST-3D approach for both datasets WORM (left) and PARHYALE (right). Each instance of a predicted cell nucleus is assigned a random color (not all shown for WORM). 3D rendering via *Paintera* (https://github.com/saalfeldlab/paintera).

[6] J. C. Caicedo, J. Roth, A. Goodman, T. Becker, K. W. Karhohs, M. Broisin, C. Molnar, C. McQuin, S. Singh, F. J. Theis, and A. E. Carpenter. Evaluation of deep learning strategies for nucleus segmentation in fluorescence images. *Cytometry Part A*, 95(9):952–965, 2019.

[7] J. Cheng, J. C. Rajapakse, et al. Segmentation of clustered nuclei with shape markers and marking function. *IEEE Trans. on Biomedical Eng.*, 56(3), 2009.

[8] F. Chollet et al. Keras. https://keras.io, 2015.

[9] Ö. Çiçek, A. Abdulkadir, S. S. Lienkamp, T. Brox, and O. Ronneberger. 3D U-Net: learning dense volumetric segmentation from sparse annotation. In *MICCAI*, 2016.

[10] Á. González. Measurement of areas on a sphere using Fibonacci and latitude–longitude lattices. *Mathematical Geosciences*, 42(1), 2010.

[11] K. He, G. Gkioxari, P. Dollár, and R. Girshick. Mask R-CNN. In *ICCV*, 2017.

[12] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *CVPR*, 2016.

[13] E. Jones, T. Oliphant, P. Peterson, et al. SciPy: Open source scientific tools for Python. http://www.scipy.org, 2001.

[14] H. W. Kuhn. The Hungarian method for the assignment problem. *Naval Research Logistics*, 2(1-2):83–97, 1955.

[15] F. Long, H. Peng, X. Liu, S. K. Kim, and E. Myers. A 3D digital atlas of C. elegans and its application to single-cell analyses. *Nature Methods*, 6(9), 2009.

[16] R. A. Lotufo, A. X. Falcao, and F. A. Zampirolli. IFT-Watershed from gray-scale marker. In *XV Brazilian Symp. on Comp. Graphics and Image Processing*, Oct 2002.

[17] E. Meijering. Cell segmentation: 50 years down the road. *IEEE Signal Processing Magazine*, 29(5), 2012.

[18] J. Redmon, S. Divvala, R. Girshick, and A. Farhadi. You only look once: Unified, real-time object detection. In *CVPR*, 2016.

[19] S. Ren, K. He, R. Girshick, and J. Sun. Faster R-CNN: Towards real-time object detection with region proposal networks. In *NIPS*, 2015.

[20] J. Schindelin, I. Arganda-Carreras, E. Frise, V. Kaynig, M. Longair, T. Pietzsch, S. Preibisch, C. Rueden, S. Saalfeld, B. Schmid, J.-Y. Tinevez, D. J. White, V. Hartenstein, K. Eliceiri, P. Tomancak, and A. Cardona. Fiji: an open-source platform for biological-image analysis. *Nature Methods*, 9, 2012.

[21] U. Schmidt, M. Weigert, C. Broaddus, and G. Myers. Cell detection with star-convex polygons. In *MICCAI*, 2018.

[22] V. Ulman, M. Maška, K. E. Magnusson, O. Ronneberger, C. Haubold, N. Harder, P. Matula, P. Matula, D. Svoboda, M. Radojevic, et al. An objective comparison of cell-tracking algorithms. *Nature Methods*, 14(12), 2017.

[23] S. van der Walt, J. L. Schönberger, J. Nunez-Iglesias, F. Boulogne, J. D. Warner, N. Yager, E. Gouillart, T. Yu, and the scikit-image contributors. scikit-image: image processing in Python. *PeerJ*, 2:e453, 6 2014.

[24] D. A. Van Valen, T. Kudo, K. M. Lane, D. N. Macklin, N. T. Quach, M. M. DeFelice, I. Maayan, Y. Tanouchi, E. A. Ashley, and M. W. Covert. Deep learning automates the quantitative analysis of individual cells in live-cell imaging experiments. *PLoS Computational Biology*, 12(11), 2016.

[25] W. Xie, J. A. Noble, and A. Zisserman. Microscopy cell counting and detection with fully convolutional regression networks. *Computer methods in biomechanics and biomedical engineering: Imaging & Visualization*, 6(3), 2018.

[26] Z. Xu, Z. Wu, and J. Feng. CFUN: Combining Faster R-CNN and U-net network for efficient whole heart segmentation. *arXiv:1812.04914*, 2018.

[27] Z. Zhao, L. Yang, H. Zheng, I. H. Guldner, S. Zhang, and D. Z. Chen. Deep learning based instance segmentation in 3D biomedical images using weak annotation. In *MICCAI*, 2018.