

Appendix A. Self-Supervised Training

Architecture. Our self-supervised training requires a network (here referred to as mini U-Net) with a field of view (FoV) smaller than the expected cell diameter. Since our analyzed datasets contain cells with diameters as small as 20 pixels wide, we use a U-Net architecture [20] with an FoV of only 16×16 . To increase the model’s capabilities without expanding the FoV, we include additional 1×1 convolutions. Each U-Net block is composed of a series of $[3 \times 3, 1 \times 1, 1 \times 1, 3 \times 3]$ valid convolution layers with ReLU activations. We use a downsampling factor of 2×2 , a depth of 1 and constant upsampling layers. In the first layer, we use 64 feature maps and increase it by a factor of 3 after each block.

Training. We train the mini U-Net on batches of 8 randomly chosen images with size 252×252 pixels. We use the Adam optimizer [11] with an initial learning rate of $4e^{-5}$ and train for 50 epochs, reducing the learning rate by a factor of 10 after epochs 20 and 30. In our pairwise loss, defined in Equation 5,

$$\mathcal{L} = \sum_{i,j \in P} \sigma \left(d(i,j) - \hat{d}(i,j) \right) + \lambda_{\text{reg}} \|r(i)\|_2,$$

we use $\sigma(\delta) = \left(1 + \exp\left(-\frac{\|\delta\|_2^2}{\tau}\right) \right)^{-1}$, $\tau = 10$, $\lambda_{\text{reg}} = 1e^{-5}$ and reduce the amount of coordinate pairs to P to reduce the GPU memory footprint. We obtain P by first sampling \mathcal{P}_1 as 20% of all pixels. For every sample in $p_1 \in \mathcal{P}_1$ we then sample $p_2 \in \mathcal{P}_2$, a random coordinate within radius $\kappa = 10$ of p_1 . In our loss we sample $i, j \in P = \mathcal{P}_1 \times \mathcal{P}_2$.

Appendix B. Supervised Training

Architecture. We use the same architecture as the mini U-Net (see Appendix A) but increase the depth of the network to 3 which expands the network’s field of view (FoV).

In conclusion, each U-Net block is composed of a series of $[3 \times 3, 1 \times 1, 1 \times 1, 3 \times 3]$ valid convolution layers with ReLU activations. We use a downsampling factor of 2×2 , a depth of 3 and constant upsampling layers. In the first layer, we use 64 feature maps and increase it by a factor of 3 after each block.

Training. All models are trained with identical training setups. We optimize the loss (see $\mathcal{L}_{\text{STARDIST}}$) with the Adam optimizer with learning rate $1e^{-5}$ for 200 epochs, reducing the learning rate by a factor of 10 at epochs 30, 80 and 160. We use batches of 8 images with size 252×252 pixels, sampling pairs of patches within radius $\kappa = 10$ (see Equation 5), and set the loss temperature $\tau = 10$.

Appendix D. SIMULATED Dataset

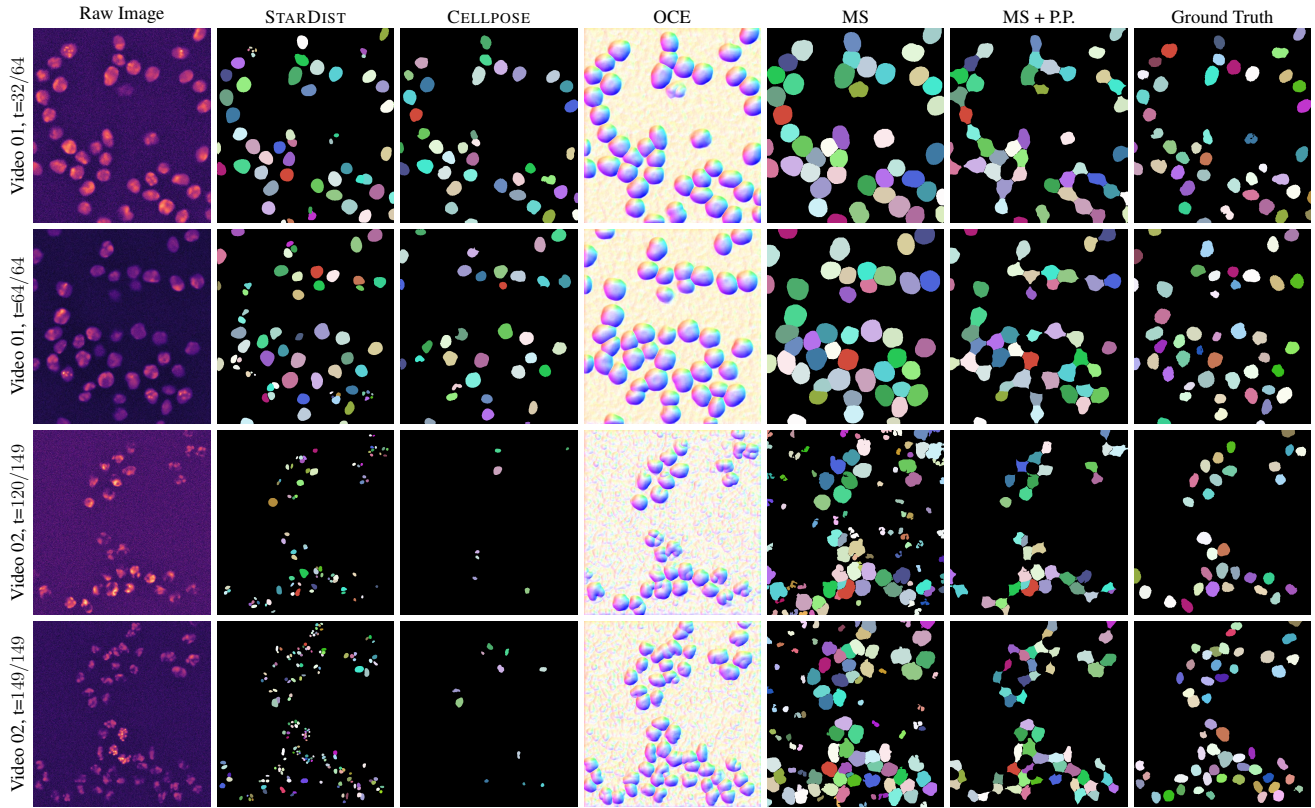


Figure 8. **Qualitative results on the SIMULATED dataset in the Cell Tracking Challenge [26].** The SIMULATED dataset comprises of two time-lapse videos (Videos 01 and 02) which contain 64 and 149 image frames respectively. Shown here are individual raw images from the two videos (first column), predicted instance segmentations obtained using the pre-trained baseline models STARDIST (second column) and CELLPOSE (third column), dense prediction of Object-Centric Embeddings (OCEs) obtained using CELLULUS (fourth column), intermediate instance segmentations obtained by applying mean-shift (MS) clustering on the dense OCEs (fifth column), these intermediate instance segmentations are further post-processed (sixth column, see more details in the Section 4.1 - *Scale Informs All Parameter Choices*) and the Ground Truth Instance Segmentation available for evaluation purposes (seventh column). Video 02 in SIMULATED (see last two rows) contains cells with visible granules which cause over-segmentation for both the evaluated, pre-trained baseline models.