# OPTIMAL TRANSPORT UNLOCKS END-TO-END LEARN-ING FOR SINGLE-MOLECULE LOCALIZATION

## Anonymous authors

Paper under double-blind review

## **ABSTRACT**

Single-molecule localization microscopy (SMLM) surpasses the diffraction limit by detecting and localizing individual fluorophores — fluorescent molecules stained onto the observed specimen — over time to reconstruct super-resolved images. Conventional SMLM requires non-overlapping emitting fluorophores, leading to long acquisition times that hinders live-cell imaging. Although recent deep-learning approaches can handle denser emissions, they rely on non-maximum suppression (NMS) layers, which are non-differentiable and may discard true positives with their local fusion strategy. In this presentation, we reformulate the SMLM training objective as a set-matching problem, deriving an optimal-transport loss that eliminates the need for NMS during inference and enables end-to-end training. Additionally, we propose an iterative neural network that integrates knowledge of the microscope's optical system inside our model. Experiments on synthetic benchmarks and real biological data show that both our new loss function and architecture surpass the state of the art at moderate and high emitter densities. Code and data are provided in the supplementary material.

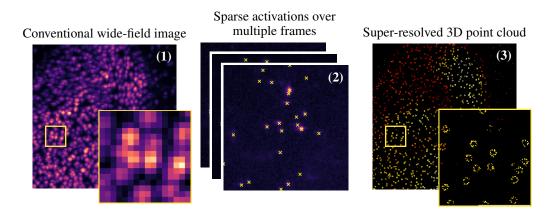


Figure 1: Illustration of the SMLM principle using our method. Data (Fei et al., 2025) show Nup96 in human bone cancer (U2-OS) cells. (1) A conventional wide-field microscope would record an image with limit resolution of  $\sim 200\,\mathrm{nm}$ . (2) Instead, SMLM captures many frames where only a sparse subset of fluorophores actively emit in each one. These can be detected and localized with sub-pixel precision. (3) The union of all detections is rendered as a 3D point cloud (color encodes depth), producing a super-resolved representation of the specimen.

## 1 Introduction

Fluorescence microscopy is a cornerstone tool of biological research that records photon emissions from *fluorophores* (fluorescent molecules) stained onto a specimen to reveal its structure. However, light diffraction restricts the raw image resolution to approximately half the wavelength of light, blurring features smaller than  $\sim 200\,\mathrm{nm}$  in practice (McCutchen, 1967; Schermelleh et al., 2019).

Multiple techniques have been developed to surpass the diffraction limit (Hell & Wichmann, 1994; Gustafsson, 2000; Dertinger et al., 2009; Laine et al., 2023), collectively called *super-resolution microscopy* methods. Among them, *single molecule localization microscopy* (SMLM) takes advantage of the stochastic flickering of fluorophores over a long sequence of images (Betzig et al., 1991). Compared to conventional fluorescence microscopy, the laser power is tuned to achieve a low density of simultaneously active fluorophores such that, with high probability, no two emitters occupy the same diffraction-limited area at the same time (Lelek et al., 2021b). Each emitter can then be individually detected and precisely located. Compiling detections across all frames yields a point cloud representation of the underlying specimen (Rust et al., 2006), effectively achieving super-resolution. Figure 1 illustrates this method. For additional details, see the review by Lelek et al. (2021b).

However, the low-density constraint inherent to this approach limits the number of active fluorophores that can be captured in a single frame, requiring thousands of frames to reconstruct a complete specimen, which hinders live-cell imaging and the observation of dynamic processes (Heilemann et al., 2008). Consequently, high-density setups are desirable, but overlapping fluorophores within the same diffraction-limited area make accurate detection and localization difficult. Deep learning methods have shown success at handling higher densities, but to the best of our knowledge, top models (Sage et al., 2019; Fei et al., 2025) all employ a non-maximum suppression (NMS) layer (Girshick et al., 2014), which prevents end-to-end learning and may discard valid detections in a close neighbourhood, which becomes more and more relevant as density increases.

In this paper, we frame supervised learning for SMLM as a set matching problem, from which we naturally derived a new loss function based on optimal transport (Peyré et al., 2019) that unlocks end-to-end learning. Furthermore, we propose a novel iterative neural network architecture that leverages a reconstruction of the expected frame given the current estimated set of fluorophores, introducing knowledge of the microscope's optic system into the model. We demonstrate that both our loss function and architecture choices improve the state of the art at both low- and high-density regimes on synthetic benchmarks and real data. Our code is available as a supplementary material and will be released as an open-source software package.

## 2 Related Work

Single-molecule localization microscopy. SMLM has been enabled by the development of photoactivable and photoswitchable fluorophores, which allows individual molecules to be precisely localized (Betzig et al., 2006; Hess et al., 2006). Early tools perform detection by locating local maxima and localizing with a Gaussian estimation of the point-spread function (PSF) (Patterson et al., 2010; Rust et al., 2006). Shortly after, the introduction of asymmetry along the z-axis of the PSF enabled 3D localisation (Huang et al., 2008); the most common setup is to introduce astigmatism in the microscope optics, which we employ in this work. To improve localisation accuracy, PSF models has transitioned from being theory-derived to experimentally derived (Babcock et al., 2012). This requires a pre-calibration step using specially designed fluorescent beads, which are imaged to capture how a single point of light appears at different locations. Note that while this calibration phase can be resource- and time-consuming, recent works propose live estimation of the PSF (Liu et al., 2024). 3D-DAOSTORM (Babcock et al., 2012) is a widely used classical method that uses experimentally derived PSFs, and we use it as a baseline in our comparisons.

SMLM delivers excellent resolution  $(10-20\,\mathrm{nm})$  and low phototoxicity, but its main drawback is slow acquisition speed (Lelek et al., 2021b). Methods such as SIM (Gustafsson, 2000), SOFI (Dertinger et al., 2009), and eSRRF (Laine et al., 2023) trade speed for reduced resolution, while STED (Hell & Wichmann, 1994) offers faster imaging at the cost of higher phototoxicity. MINFLUX (Balzarotti et al., 2017) is a promising young technique with comparable resolution to SMLM, but suffers from a small field of view and requires costly, ultra-stable microscopes (Scheiderer et al., 2024). Therefore, SMLM offers an attractive middle ground for biologists among all super-resolution methods, which explains why enhancing performance for high-density setups is of major interest (Lelek et al., 2021a).

Deep-learning methods have been widely applied to fluorescence microscopy (Nehme et al., 2018; Ouyang et al., 2018; Cachia et al., 2023; Li et al., 2023; Mentagui et al., 2024; Fei et al., 2025). Deep SMLM tools include DECODE (Speiser et al., 2021), which is currently the best on the EPFL

SMLM challenge (Sage et al., 2019), a popular benchmark for SMLM tools. More recently, Lite-Loc (Fei et al., 2025) has slightly improved on DECODE's architecture with additional technical refinements. We use those two methods for our comparison in our benchmarks.

**Optimal transport for set matching.** Optimal transport (Peyré et al., 2019; Villani, 2021) has become a fundamental tool for set matching in deep learning. Recent works in object detection (Carion et al., 2020; Zhu et al., 2020; Zhang et al., 2022; Li et al., 2022) have demonstrated success in predicting sets of variable and unknown size using loss functions derived from optimal transport theory, with recent works employing entropic regularization (Cuturi, 2013) to achieve fully differentiable pipelines (Zareapoor et al., 2025). By framing SMLM as a set matching problem, we draw a direct connection to this line of work — substituting fluorophores for objects — enabling the design of an end-to-end training procedure.

**Iterative refinement network.** Iterative refinement with neural networks has proven effective for tasks that benefit from sequential solution improvement (Carreira et al., 2016; Yu et al., 2023). In computer vision, Putzky & Welling (2017) have applied this approach to inverse problems such as image denoising, super-resolution, and inpainting, while Hur & Roth (2019) proposed iterative optical-flow refinement using a feedback loop with a rewarping operator. Because the physics of SMLM is well understood (Etheridge et al., 2022), we show that an accurate simulator of the microscope's physics can provide similar visual feedback, enabling progressive refinement of the solution.

# 3 METHOD

## 3.1 PROBLEM FORMULATION

In this section, we first introduce the image formation model for SMLM and formulate the corresponding inverse problem as a set matching task. We then present a differentiable loss function and an iterative refinement architecture that explicitly leverages the image formation process.

**Image formation model** An *activation* is defined as an emission event from a fluorophore within a given frame (a single fluorophore may produce several activations accross multiple frames). Throughout this work, an activation is represented by a 4D vector  $\mathbf{x} = (x, y, z, n)$ , where (x, y) denote the 2D coordinates in the camera frame (with the origin at the top-left corner), z represents the axial coordinate relative to the focal plane, and n is the photon count. Given N activations within a frame, we denote the complete set as  $\mathcal{X} = \{x_i\}_{1 \le i \le N}$ .

Diffraction within the microscope's optical system is modeled by a convolution with a kernel called the *point spread function* (PSF) (Rossmann, 1969). It can be thought of as the image of a single point source. We represent the PSF as a function  $\mathbf{P}:\mathbb{R}^3\longmapsto\mathbb{R}^{H\times W}$ , that outputs the normalized  $H\times W$  image resulting from the diffraction of a single point source given its 3D coordinates. To ensure photon count independence, the output image is normalized to sum to unity in the focal plane, i.e. for z=0. Given the set of activations  $\mathcal X$  in a frame, the observed  $H\times W$  image, denoted by  $\mathbf{H}(\mathcal X)$ , is formed as a weighted sum of PSFs, where the weights correspond to the photon count n for each activation:

$$\mathbf{H}(\mathcal{X}) = \sum_{(x,y,z,n)\in\mathcal{X}} n\mathbf{P}(x,y,z). \tag{1}$$

The dependence of the PSF on depth z enables 3D localization of activations from the observed image, see (Ovesný et al., 2014). Following Babcock & Zhuang (2017), we assume that the PSF is pre-calibrated on synthetic fluorescent beads and implemented as a collection of 3D splines. This approach is a standard tool in SMLM used in many works (Ries, 2020; Li et al., 2020; Speiser et al., 2021; Etheridge et al., 2022); see Babcock & Zhuang (2017) for further details.

**Noise model.** We adopt the noise model of Sage et al. (2019), which combines shot noise (modeled by a Poisson distribution), amplification noise (modeled by a Gamma distribution only for EM-CCD camera) and readout noise (modeled by a normal distribution). In-depth description of all camera parameters is available in Appendix A.1. The noise for each camera sensor being independent and identically distributed (Fazel & Wester, 2022), it is applied independently to all pixels of  $\mathbf{H}(\mathcal{X})$ . Then, we denote by  $\mathbb P$  the distribution of images  $\mathbf y$  produced by a set  $\mathcal X$  of fluorophores under the noise model such that

$$y \sim \mathbb{P}(\mathcal{X})$$
. (2)

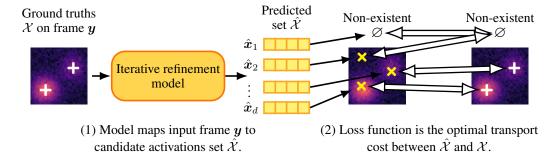


Figure 2: Illustration of our loss function for end-to-end training. (1) Given a simulated image and its ground truth activations (see Section 3.1), our model (see Section 3.3) predicts d candidate activations, each with a detection score quantifying the plausibility of its existence. (2) We solve a regularized optimal transport problem — conceptually similar to a bi-matching between ground truths and predictions — over a cost involving both localisation and detection tasks. Our loss function is the optimal cost yielded by this solution.

**Risk minimization formulation for set matching.** Because no ground truth is available in most scientific imaging applications, supervised-learning models for SMLM have to be trained with a simulator, which is able to generate realistic y from  $\mathbb{P}(\mathcal{X})$  given sets  $\mathcal{X}$  of activations from a distribution  $\mathcal{D}$ . Our approach consists of training a neural network  $f_{\theta}$  which directly predicts a set of activations given an observation y, by minimizing the risk

$$\theta^* = \operatorname*{arg\,min}_{\theta \in \Theta} \mathbb{E}_{\mathcal{X} \sim \mathcal{D}, \boldsymbol{y} \sim \mathbb{P}(\mathcal{X})} \left[ \mathcal{L}(f_{\theta}(\boldsymbol{y}), \mathcal{X}) \right]. \tag{3}$$

Such a formulation raises two major challenges: we need to design a differentiable loss function  $\mathcal{L}$  and an architecture  $f_{\theta}$  that are appropriate to the context of SMLM.

#### 3.2 OPTIMAL TRANSPORT LOSS FUNCTION

We argue that framing SMLM as a supervised-learning problem leads to a set matching formulation, for which optimal transport theory is a natural fit. To the best of our knowledge, however, this framework has not yet been applied to SMLM. Figure 2 provides an overview of our method.

Let  $\mathcal{X} = \{x_i\}_{1 \leq i \leq N}$  be the ground truth set of activations. The size of this set, N, is unknown and varies between frames. We simulate an acquisition  $y \sim \mathbb{P}(\mathcal{X})$  and aim to retrieve  $\mathcal{X}$  from y.

Given y, our neural network  $f_{\theta}$  outputs a set of candidate activations  $\hat{\mathcal{X}} = \{\hat{x}_i\}_{1 \leq i \leq d}$  of fixed size d, each associated with a detection score  $\hat{s}_i$  in (0,1) gathered in a set  $\hat{\mathcal{S}} = \{\hat{s}_i\}_{1 \leq i \leq d}$ . The network architecture is detailed in Section 3.3. The number of candidates d is a constant hyperparameter that fixes the maximum possible number of detectable activations, so it should be set large enough to ensure that the true number of fluorophores N will always be less than or equal to d, see Section 4.

We first define L, a squared cost matrix of size  $d \times d$ , whose components are:

$$\forall 1 \le i, j \le d, L_{i,j} = \begin{cases} (\hat{\boldsymbol{x}}_i - \boldsymbol{x}_j)^T \boldsymbol{\Sigma}^{-1} (\hat{\boldsymbol{x}}_i - \boldsymbol{x}_j) + \log \det (\boldsymbol{\Sigma}) & \text{if } j \le N, \\ 0 & \text{otherwise,} \end{cases}$$
(4)

where  $\Sigma = \operatorname{diag}(\sigma_x^2, \sigma_y^2, \sigma_z^2, \sigma_n^2)$  is a diagonal weighting matrix. Quadratic costs are a natural and principled choice for regression tasks. Extending this formulation to the negative log-likelihood of a multivariate normal distribution allows to learn  $\Sigma$  end-to-end, which can be viewed as an automatic weighting strategy that balances the difficulty of predicting each dimension, similar to the homoscedastic uncertainty weighting method proposed by Kendall et al. (2018).

Similarly, we define D, another  $d \times d$  cost matrix whose components are:

$$\forall 1 \le i, j \le d, D_{i,j} = \begin{cases} -\log(s_i) & \text{if } j \le N, \\ -\log(1 - s_i) & \text{otherwise.} \end{cases}$$
 (5)

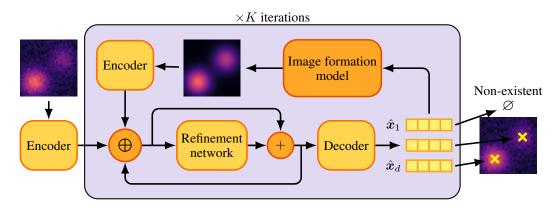


Figure 3: Illustration of our iterative refinement model. Within a classic encoder-decoder architecture, we leverage prior knowledge about the known image formation model (not learned) to simulate the expected frame given the current latent representation. This feedback is used to iteratively refine the model's inner latent representation for K steps. The encoders are identical. + and  $\oplus$  respectively denotes element-wise addition and concatenation.

The binary cross-entropy cost is a natural choice for detection tasks. It favors a high score  $s_i$  when  $\hat{x}_i$  is paired with an element of  $\mathcal{X}$  and low score otherwise, hence promoting good detection. Finally, we define the total cost matrix C = L + D, which integrates both localization and detection tasks.

Considering the initial set matching problem, the optimal solution  $(\hat{\mathcal{X}}^*, \hat{\mathcal{S}}^*)$  given a target  $\mathcal{X}$  would consist of N elements, each identical to one element of  $\mathcal{X}$  and with detection scores close to 1. The remaining d-N elements have detection scores close to 0. Naturally, one would like to compare each candidate in  $\hat{\mathcal{X}}^*$  to its nearest counterpart in  $\mathcal{X}$  and minimize a loss function over these pairs. This can be achieved by solving an optimal-transport problem over C—conceptually creating a bipartite matching between the predictions and the ground truths—where the minimal cost accounts for all pairwise contributions. Therefore, we would ideally like our loss function to be the optimal-transport cost with respect to C, i.e. solve:

$$\min_{\mathbf{\Gamma} \in \mathcal{B}} \langle \mathbf{\Gamma} | \mathbf{C} \rangle_{\mathcal{F}} \quad \text{where } \mathcal{B} = \left\{ \mathbf{\Gamma} \in \mathbb{R}_{+}^{d \times d} \mid \mathbf{\Gamma} \mathbf{1}_{d} = \mathbf{\Gamma}^{\top} \mathbf{1}_{d} = \mathbf{1}_{d} \right\}, \tag{6}$$

and  $\langle .|. \rangle_{\mathcal{F}}$  is the Frobenius inner product. However, while the Hungarian algorithm can exactly solve this problem in  $\mathcal{O}(d^3)$  (Kuhn, 1955), its algorithmic step is non-differentiable, which prevents end-to-end learning. We circumvent this issue by finding  $\Gamma$  through the entropy-regularized optimal transport problem, see (Cuturi, 2013), and therefore define our loss function as follows:

$$\mathcal{L}(\hat{\mathcal{X}}, \hat{\mathcal{S}}, \mathcal{X}) = \langle \mathbf{\Gamma}^* | \mathbf{C} \rangle_{\mathcal{F}}, \quad \text{where } \mathbf{\Gamma}^* = \underset{\mathbf{\Gamma} \in \mathcal{B}}{\arg \min} \langle \mathbf{\Gamma} | \mathbf{C} \rangle_{\mathcal{F}} - \epsilon H(\mathbf{\Gamma}), \tag{7}$$

H is the Shannon entropy and  $\epsilon$  the entropic regularization parameter. A good approximation of  $\Gamma^*$  in Eq. (7) can be found efficiently with a few iterations of the Sinkhorn algorithm, whose steps are differentiable with respect to the elements of C, enabling its use within a deep learning framework, see (Genevay et al., 2018; Mialon et al., 2021).

### 3.3 Iterative refinement scheme

To solve Eq. (3), we investigate architectures that explicitly leverage the image formation process. To this end, we adopt an iterative architecture, an idea that has proven successful for optical flow estimation (Hur & Roth, 2019). At each iteration the network produces a set of candidate activations, turns those proposals into a simulated image, which is then used as feedback to refine the next proposals. This iterative method is illustrated in Figure 3.

Concretely, let  $\boldsymbol{y}$  be the input frame of size  $H \times W$ . An encoder  $\boldsymbol{E} : \mathbb{R}^{H \times W} \longmapsto \mathbb{R}^{C \times H \times W}$  maps  $\boldsymbol{y}$  to a latent representation  $\boldsymbol{z}^{(0)}$ , where C is a hyperparameter controlling the dimension of the latent space. A decoder  $\boldsymbol{D}$  then maps latent variables to a set of candidate activations  $\hat{\mathcal{X}} = \{\hat{\boldsymbol{x}}_i\}_{1 \leq i \leq d}$  and a corresponding set of detection scores  $\hat{\mathcal{S}} = \{\hat{s}_i\}_{1 \leq i \leq d}$ .

# Algorithm 1 Iterative refinement architecture

Require: input frame  $y \in \mathbb{R}^{H \times W}$ , encoder E, decoder D, refinement module R, camera model  $\mathbb{P}(\cdot)$ , number of iterations  $K \in \mathbb{N}$ **Ensure:** final proposals  $(\hat{\mathcal{X}}^{(K)}, \hat{\mathcal{S}}^{(K)})$ 1:  $\boldsymbol{z}^{(0)} \leftarrow \boldsymbol{E}(\boldsymbol{y})$ ⊳ encode original frame 2:  $(\hat{\mathcal{X}}^{(0)}, \hat{\mathcal{S}}^{(0)}) \leftarrow D(z^{(0)})$ 3: **for** k = 0 to K - 1 **do**  $\hat{oldsymbol{y}}^{(k)} \leftarrow \mathbb{E}[\mathbb{P}(\hat{\mathcal{X}}^{(k)}, \hat{\mathcal{S}}^{(k)})]$ ▶ simulate reconstruction from current proposals  $\hat{m{z}}^{(k)} \leftarrow m{E}(\hat{m{y}}^{(k)})$ ▷ encode reconstruction  $z^{(k+1)} \leftarrow z^{(k)} + R(z^{(k)}, \hat{z}^{(k)}, z^{(0)})$  ▷ refine latent iteratively  $(\hat{\mathcal{X}}^{(k+1)}, \hat{\mathcal{S}}^{(k+1)}) \leftarrow D(\boldsymbol{z}^{(k+1)})$ 8: end for 9: return  $(\hat{\mathcal{X}}^{(K)}, \hat{\mathcal{S}}^{(K)})$ 

Given  $(\hat{\mathcal{X}}, \hat{\mathcal{S}})$ , we compute a reconstructed frame  $\hat{y} = \mathbb{E}[\hat{y}|\hat{\mathcal{X}}, \hat{\mathcal{S}}]$ .  $\hat{y}$  is the expected image produced by the current proposal set, and thus provides a visual summary of what the model's output currently explain in the SMLM frame. Comparing the reconstructed image  $\hat{y}$  to the original frame y supplies informative feedback, which helps the model correct errors and refine the candidates set over iterations.

Concretely, we define an iterative refinement operator  $\mathbf{R}: \mathbb{R}^{3 \times C \times H \times W} \longmapsto \mathbb{R}^{C \times H \times W}$  which produces a residual update of the latent representation given it's current estimate, the representation of the simulated frame, and the encoded original frame. Algorithm 1 shows how this proposal is updated successively over K steps.

During training, the final decoded output  $(\hat{\mathcal{X}}^{(K)}, \hat{\mathcal{S}}^{(K)})$  is used as input for our loss function, see Section 3.2. During inference, candidate activations in  $\hat{\mathcal{X}}^{(K)}$  are filtered based on their associated detection scores in  $\hat{\mathcal{S}}^{(K)}$ : only activations exceeding a pre-defined threshold  $\tau \in (0,1)$  are retained; others are discarded.

Details about the decoder architecture and the computation of  $\hat{y}$  given  $\hat{\mathcal{X}}$  and  $\hat{\mathcal{S}}$  can be found in Appendix A.2.

# 4 EXPERIMENTS

**Implementation and training details.** We construct a synthetic target activations set  $\mathcal{X}$  from  $\mathcal{D}$  by uniformly sampling between 10 and 30 activations per frame, assigning each activation independent coordinates that are uniformly distributed across all dimensions. This guarantees that the network cannot learn any specific prior about the activation distribution.

The encoder E is a two-layer U-Net (Ronneberger et al., 2015) with an internal channel width of 48, mapping the input image to a latent image with C=96 channels. The iterative refinement stage is implemented with a similar two-layer U-Net (also 48 channels for the first level). The decoder D is a three-layer convolutional neural network with a number of predicted candidates d equal to HW/4, see Appendix A.2 for additional details. The resulting network contains  $\sim 3$  millions learnable parameters. Empirically, performance improvement stops after three or more refinement iterations; we thus use K=2 in our experiments.

Following (Speiser et al., 2021), we augment y by the previous and the next frame into a tensor  $\bar{y}$  of size  $3 \times H \times W$ . Including these two frames improve the model performance as it provides additional context about the processed frame.

For training, we use AdamW (Loshchilov & Hutter, 2017) for 100,000 steps with a batch size of 128 on a NVIDIA-H100 gpu, taking approximately 20h. The iterative architecture incurs a higher computational burden than single-pass models like DECODE or LiteLoc; further details about our model computational footprint for training and inference is available in Appendix A.3.

At inference, while a threshold  $\tau = 0.5$  yields good results, we find better performance by calibrating  $\tau$  by maximizing the  $E_{3D}$  metric (defined in Section 4) on a separate synthetic dataset.

Density	SNR	Method	Precision ↑	Recall ↑	Jaccard ↑	$RMSE_{lat} \downarrow$	$RMSE_{ax} \downarrow$	E <sub>3D</sub> ↑
0.2	High	3D-DAOSTORM DECODE LiteLoc Ours	$\begin{array}{c} 0.964 \\ 0.961 \pm 0.003 \\ 0.996 \pm 0.002 \\ \textbf{0.998} \pm \textbf{0.002} \end{array}$	0.919 $0.998 \pm 0.001$ $0.987 \pm 0.001$ $0.978 \pm 0.016$	$\begin{array}{c} 0.914 \\ 0.959 \pm 0.003 \\ \textbf{0.983} \pm \textbf{0.001} \\ 0.980 \pm 0.010 \end{array}$	$11.9$ $8.8 \pm 0.1$ $9.0 \pm 0.1$ <b>7.5</b> $\pm$ <b>0.4</b>	$16.9$ $10.7 \pm 0.1$ $11.7 \pm 0.1$ $10.0 \pm 3.5$	$\begin{array}{c} 0.821 \\ 0.895 \pm 0.003 \\ 0.912 \pm 0.001 \\ \textbf{0.920} \pm \textbf{0.007} \end{array}$
	Low	3D-DAOSTORM DECODE LiteLoc Ours	$\begin{array}{c} 0.978 \\ 0.918 \pm 0.002 \\ \textbf{0.995} \pm \textbf{0.001} \\ 0.985 \pm 0.001 \end{array}$	$0.835$ $0.978 \pm 0.001$ $0.939 \pm 0.001$ $0.961 \pm 0.001$	$\begin{array}{c} 0.833 \\ 0.903 \pm 0.002 \\ 0.934 \pm 0.001 \\ \textbf{0.947} \pm \textbf{0.001} \end{array}$	$\begin{array}{c} 19.3 \\ 20.5 \pm 0.1 \\ \textbf{17.0} \pm \textbf{0.1} \\ 18.8 \pm 0.2 \end{array}$	$\begin{array}{c} 29.8 \\ 26.2 \pm 0.1 \\ 25.0 \pm 0.4 \\ \textbf{24.5} \pm \textbf{0.1} \end{array}$	$0.685 \\ 0.757 \pm 0.001 \\ 0.798 \pm 0.001 \\ \mathbf{0.802 \pm 0.002}$
2.0	High	3D-DAOSTORM DECODE LiteLoc Ours	$\begin{array}{c} 0.914 \\ 0.923 \pm 0.003 \\ \textbf{0.993} \pm \textbf{0.001} \\ 0.992 \pm 0.002 \end{array}$	$0.678$ $0.946 \pm 0.002$ $0.863 \pm 0.002$ $0.895 \pm 0.011$	$\begin{array}{c} 0.643 \\ 0.876 \pm 0.004 \\ 0.858 \pm 0.001 \\ \textbf{0.883} \pm \textbf{0.017} \end{array}$	$56.8$ $32.2 \pm 0.3$ $30.7 \pm 0.2$ <b>24.8</b> $\pm$ <b>0.6</b>	$76.6$ $33.0 \pm 0.4$ $36.0 \pm 0.3$ $28.4 \pm 0.5$	$\begin{array}{c} 0.373 \\ 0.706 \pm 0.004 \\ 0.699 \pm 0.001 \\ \textbf{0.750} \pm \textbf{0.004} \end{array}$
	Low	3D-DAOSTORM DECODE LiteLoc Ours	$\begin{array}{c} 0.914 \\ 0.859 \pm 0.034 \\ \textbf{0.992} \pm \textbf{0.001} \\ 0.973 \pm 0.003 \end{array}$	$0.496$ $0.874 \pm 0.006$ $0.729 \pm 0.002$ $0.812 \pm 0.007$	$\begin{array}{c} 0.475 \\ 0.756 \pm 0.027 \\ 0.725 \pm 0.001 \\ \textbf{0.794} \pm \textbf{0.005} \end{array}$	$74.4$ $56.4 \pm 0.3$ $46.2 \pm 0.4$ $48.4 \pm 0.6$	$\begin{array}{c} 120.0 \\ 65.3 & \pm 0.4 \\ 63.8 & \pm 0.2 \\ 59.5 & \pm 0.4 \end{array}$	$\begin{array}{c} 0.116 \\ 0.468 \pm 0.008 \\ 0.500 \pm 0.002 \\ \textbf{0.536} \pm \textbf{0.003} \end{array}$
8.0*	High	3D-DAOSTORM DECODE LiteLoc Ours	$\begin{array}{c} 0.914 \\ 0.869 \pm 0.032 \\ 0.992 \pm 0.001 \\ \textbf{0.995} \pm \textbf{0.001} \end{array}$	$0.388$ $0.692 \pm 0.008$ $0.567 \pm 0.004$ $0.603 \pm 0.022$	$0.376$ $0.632 \pm 0.028$ $0.564 \pm 0.003$ $0.600 \pm 0.028$	$79.6$ $60.3 \pm 0.7$ $54.2 \pm 0.8$ <b>43.3</b> $\pm 1.9$	$\begin{array}{c} 130.1 \\ 92.0 \pm 2.7 \\ 71.4 \pm 1.6 \\ \textbf{49.2} \pm \textbf{1.5} \end{array}$	$\begin{array}{c} 0.028 \\ 0.341 \pm 0.027 \\ 0.360 \pm 0.007 \\ \textbf{0.462} \pm \textbf{0.007} \end{array}$
	Low	3D-DAOSTORM DECODE LiteLoc Ours	$\begin{array}{c} 0.899 \\ 0.824 \pm 0.009 \\ \textbf{0.993} \pm \textbf{0.002} \\ 0.990 \pm 0.002 \end{array}$	$0.211$ $0.591 \pm 0.028$ $0.349 \pm 0.002$ $0.432 \pm 0.025$	$0.207$ $0.528 \pm 0.038$ $0.348 \pm 0.002$ $0.433 \pm 0.029$	$85.8$ $74.6 \pm 1.9$ $66.6 \pm 0.4$ $58.1 \pm 3.5$	$171.3$ $89.0 \pm 4.4$ $92.5 \pm 1.4$ $72.2 \pm 2.7$	$-0.196 \\ 0.224 \pm 0.022 \\ 0.124 \pm 0.002 \\ 0.244 \pm 0.004$

Table 1: Comparative evaluation of SMLM algorithms on the EPFL 2016 challenge datasets and metrics. For each method, means and standard deviations are estimated over four independent training seeds (3D-DAOSTORM is deterministic). \*Note: the EPFL 2016 challenge does not include a density-8.0 dataset; we have created by temporally binning the density-2.0 datasets.

**Synthetic data.** Because no ground-truth annotations exist for real SMLM acquisitions, we performe the initial evaluation on the open synthetic datasets provided by Sage et al. (2019) on the 2016 EPFL challenge, and adopte their set of metrics.

To evaluate candidate activations in a frame, we first solve a Hungarian assignment between ground-truths and predicted activations. A prediction is considered a *true positive* (TP) if it lies within  $\pm 250$  nm in both x and y directions, and  $\pm 500$  nm in z relative to its matched ground-truth (both thresholds come from the EPFL challenge). Otherwise, predictions (resp. ground truths) are labeled as *false positives* (resp. *false negatives*). Detection performance is quantified by computing *precision, recall* and *Jaccard Index* (*area under the curve* is not commonly employed in this field). Localisation performance is evaluated by computing the *root-mean-square error* (RMSE) for TPs, both for the lateral plan and the axial dimension. A global performance metric called *3D efficiency* (E<sub>3D</sub>) is then defined as:

$$E_{3D} = \frac{E_{ax} + E_{lat}}{2} \quad \text{where} \begin{cases} E_{lat} = 1 - \sqrt{(1 - Jaccard)^2 + \alpha_{lat}^2 RMSE_{lat}^2}, \\ E_{ax} = 1 - \sqrt{(1 - Jaccard)^2 + \alpha_{ax}^2 RMSE_{ax}^2}, \end{cases} \tag{8}$$

 $\alpha_{\rm lat} = 1.0 \, \rm nm^{-1}$  and  $\alpha_{\rm ax} = 0.5 \, \rm nm^{-1}$ , following definitions of the EPFL challenge. All metrics are computed frame by frame and averaged.

Our benchmark includes 3D-DAOSTORM (Babcock et al., 2012), DECODE (Ouyang et al., 2018) and LiteLoc (Fei et al., 2025). All algorithms are evaluated on the open-access EPFL 2016 challenge datasets (Sage et al., 2019), all with astigmatism PSFs. To assess performance in a very high-density regime, we have synthesized a density-8.0 benchmark by temporally binning groups of 4 frames in the original density-2.0 sequences. Note that this procedure boosts the SNR through noise averaging.

Results are reported in Table 1. We observe that while our approach yields lower recall than the other methods, it preserves excellent precision and almost always achieves the lowest RMSE in all spatial dimensions. Most notably, it also outperforms all competitors on the  $E_{3D}$  metric for all densities and SNRs, establishing itself as the most "complete" method in this evaluation.

**Real data.** We have evaluated our method on three publicly available datasets, all of which provide beads for calibrating their astigmatic PSFs. The Tubulin and NPC-Nup107 datasets from Li et al. (2018) depict, respectively, the microtubule network and nuclear pore complexes in U2OS cells.

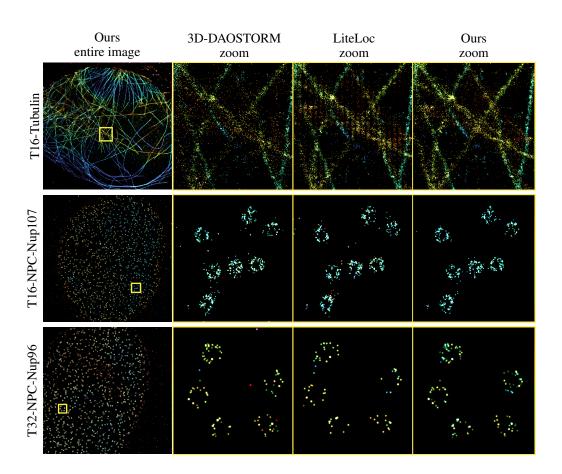


Figure 4: Qualitative comparison of SMLM methods on real data. Although ground truths are unavailable, results show that our approach yields fewer grid-reconstruction artifacts (line 1), improved depth estimation consistency (line 2), and more accurate nuclear pore complex reconstruction (line 3). Refer to the main text for a more thorough discussion.

Dataset	Bin size	Method	FRC (nm) ↓	RSP ↑
Tubulin Li et al. (2018)	×1	LiteLoc Ours	$29.7 \pm 0.3$ $31.9 \pm 0.2$	<b>0.708</b> 0.692
	×16	LiteLoc Ours	$63.0 \pm 0.8$ $58.1 \pm 1.1$	0.649 <b>0.672</b>
NPC-Nup107 Li et al. (2018)	×1	LiteLoc Ours	$19.3 \pm 0.3$ $18.8 \pm 0.4$	<b>0.696</b> 0.686
,	×16	LiteLoc Ours	$25.9 \pm 0.3$ $22.1 \pm 0.1$	0.682 <b>0.684</b>
NPC-Nup96 Fei et al. (2025)	×1	LiteLoc Ours	$\mathbf{29.8 \pm 0.1}$ $31.7 \pm 0.1$	<b>0.713</b> 0.693
	$\times 32$	LiteLoc Ours	$71.5 \pm 0.5$ $44.2 \pm 0.4$	0.671 <b>0.689</b>

Table 2: Quantitative results on real datasets. We temporally binned them to simulate very high-density setups. Our method consistently scored first in those denser regimes.

The NPC-Nup96 dataset from Fei et al. (2025) also features nuclear pore complexes in the same cell line. All datasets were acquired with conventional SMLM activation densities; therefore, to test our method's robustness at higher densities, we applied 16-frame temporal binning to Tubulin and

Iterative arch.	OT loss func.	Jaccard ↑	$RMSE_{vol} \downarrow$	E <sub>3D</sub> ↑
×	×	$0.876 \pm 0.004$	$47.9 \pm 0.5$	$0.705 \pm 0.004$
×	<b>✓</b>	$0.867 \pm 0.004$	$39.6 \pm 0.3$	$0.740 \pm 0.002$
<b>✓</b>	×	$0.883 \pm 0.007$	$41.6 \pm 0.2$	$0.725 \pm 0.003$
<b>~</b>	<b>✓</b>	$0.883 \pm 0.007$	$39.2 \pm 0.5$	$0.750 \pm 0.003$

Table 3: Ablation study of our different modules over EPFL synthetic with high SNR and a density of 2.0. ★ means we used DECODE's original loss function or model architecture.

NPC-Nup107 and 32-frame binning to NPC-Nup96. We refer to the temporally-binned versions as T16-Tubulin, T16-NPC-Nup107, and T32-NPC-Nup96. Note that this approach is an imperfect proxy for truly high-density imaging, as it improves the signal-to-noise ratio via noise averaging.

Figure 4 compares 3D SMLM reconstructions, rendered with SMAP (Ries, 2020), for 3D-DAOSTORM (Babcock et al., 2012), LiteLoc (Fei et al., 2025) and our method. We omit DE-CODE because its output is similar to LiteLoc. On the T16-Tubulin dataset, our algorithm yields a higher-fidelity reconstruction than 3D-DAOSTORM and eliminates the artefacts that appear with LiteLoc. On the T16-NPC-Nup107 dataset, all methods recover comparable structures; however, our method delivers more consistent depth estimates (as indicated by colors), whereas 3D-DAOSTORM and LiteLoc exhibit spatially varying detections. On the T32-NPC-Nup96 dataset, our approach reconstructed NPC's structures with clear greater fidelity than LiteLoc and 3D-DAOSTORM.

Quantitatively, the absence of ground-truth data prevents the use of the metrics introduced in section 4. To evaluate the resolution and fidelity of a reconstructed super-resolution image, we adopted two widely used metrics: Fourier ring correlation (FRC) (Banterle et al., 2013) and the resolution-scaled Pearson's coefficient (RSP) Culley et al. (2018). FRC reconstructs two super-resolution images by splitting localisations into two subsets, computing their Fourier transforms, and then measuring the correlation of their spatial frequency signals against each other. The resulting curve provides an estimate of the spatial frequency at which signal can no longer be distinguished from noise (Banterle et al., 2013). RSP is defined as the Pearson correlation coefficient between the reconstructed super-resolution image and a reference image, typically the mean of all raw wide-field frames. Values close to one indicate strong agreement between the reconstruction and the reference.

Results with these metrics on real datasets are reported in Table 2. In dense-activations regimes, our approach consistently yields lower FRC and higher RSP values than other methods, confirming the visual improvements illustrated in Figure 4.

**Ablation study.** We have conducted an ablation study on synthetic data to validate the effectiveness of our loss function and our iterative architecture. Results are reported in Table 3. It can be seen that the loss function drives most of the improvement, with our iterative architecture providing a modest boost. Given the additional memory and compute overhead of our architecture, a lightweight variant that retains only the optimal loss function can be considered for deployment scenarios with constrained resources.

### 5 DISCUSSION AND CONCLUDING REMARKS

We have presented a novel deep-learning SMLM method that surpasses existing methods in medium and high-density regimes, all without the need for handcrafted layers. By enabling faster data acquisition, our approach extend SMLM's temporal resolution, allowing more accurate observation of rapid biological processes. Furthemore, the integration of optimal transport theory to SMLM could open a path to new localization algorithms.

The main limitation of our method is the longer training and inference times that result from its iterative design. However, training is a one-time cost per experimental setup, and inference remains fast enough ( $\sim 200 \, \mathrm{fps}$  on a modern GPU) to let biologists run multiple experiments sequentially with minimal delay. Another limitation that nearly all top-performing methods share is the dependence for precise PSF calibration (Lelek et al., 2021a). Future work could focus on robust methods invariant to PSF variations, or even pursue blind SMLM super-resolution without sacrificing precision.

# ETHICS STATEMENT

In this work, we explore new model architecture and loss function to improve single-molecule localization microscopy (SMLM) reconstruction. We do not anticipate ethical or societal harms: the work is computational only and all biological data used are public and were used according to their licenses. We believe that by improving SMLM reconstruction and releasing our code openly, this work can broadly benefit biological research and make advanced tools accessible to communities worldwide.

## REPRODUCIBILITY STATEMENT

The project repository includes all requirements to reproduce our results. We provide the full source code and model implementation, all datasets are publicly available, training and evaluation scripts are provided with all hyperparameters set (including random seeds), and python environment specification are supplied, making all experiments from section 4 reproducible.

## REFERENCES

- Hazen Babcock, Yaron M Sigal, and Xiaowei Zhuang. A high-density 3d localization algorithm for stochastic optical reconstruction microscopy. *Optical nanoscopy*, 1(1):6, 2012.
- Hazen P. Babcock and Xiaowei Zhuang. Analyzing Single Molecule Localization Microscopy Data Using Cubic Splines. *Scientific Reports*, 7(1):552, 2017.
- Francisco Balzarotti, Yvan Eilers, Klaus C Gwosch, Arvid H Gynnå, Volker Westphal, Fernando D Stefani, Johan Elf, and Stefan W Hell. Nanometer resolution imaging and tracking of fluorescent molecules with minimal photon fluxes. *Science*, 355(6325):606–612, 2017.
- Niccolo Banterle, Khanh Huy Bui, Edward A Lemke, and Martin Beck. Fourier ring correlation as a resolution criterion for super-resolution microscopy. *Journal of structural biology*, 183(3): 363–367, 2013.
- EJKTDJSRL Betzig, Jay K Trautman, TD Harris, JS Weiner, and RL Kostelak. Breaking the diffraction barrier: optical microscopy on a nanometric scale. *Science*, 251(5000):1468–1470, 1991.
- Eric Betzig, George H Patterson, Rachid Sougrat, O Wolf Lindwasser, Scott Olenych, Juan S Bonifacino, Michael W Davidson, Jennifer Lippincott-Schwartz, and Harald F Hess. Imaging intracellular fluorescent proteins at nanometer resolution. *science*, 313(5793):1642–1645, 2006.
- Mayeul Cachia, Vasiliki Stergiopoulou, Luca Calatroni, Sébastien Schaub, and Laure Blanc-Féraud. Fluorescence image deconvolution microscopy via generative adversarial learning (fluogan). *Inverse Problems*, 39(5):054006, 2023.
- Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-End Object Detection with Transformers, May 2020.
- Joao Carreira, Pulkit Agrawal, Katerina Fragkiadaki, and Jitendra Malik. Human pose estimation with iterative error feedback. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 4733–4742, 2016.
- Siân Culley, David Albrecht, Caron Jacobs, Pedro Matos Pereira, Christophe Leterrier, Jason Mercer, and Ricardo Henriques. Quantitative mapping and minimization of super-resolution optical imaging artifacts. *Nature methods*, 15(4):263–266, 2018.
- Marco Cuturi. Sinkhorn distances: Lightspeed computation of optimal transport. *Advances in neural information processing systems*, 26, 2013.
- Thomas Dertinger, Ryan Colyer, Gopal Iyer, Shimon Weiss, and Jörg Enderlein. Fast, background-free, 3d super-resolution optical fluctuation imaging (sofi). *Proceedings of the National Academy of Sciences*, 106(52):22287–22292, 2009.
  - Thomas J Etheridge, Antony M Carr, and Alex D Herbert. Gdsc smlm: Single-molecule localisation microscopy software for imagej. *Wellcome open research*, 7:241, 2022.

- Mohamadreza Fazel and Michael J Wester. Analysis of super-resolution single molecule localization
   microscopy data: A tutorial. *AIP advances*, 12(1), 2022.
  - Yue Fei, Shuang Fu, Wei Shi, Ke Fang, Ruixiong Wang, Tianlun Zhang, and Yiming Li. Scalable and lightweight deep learning for efficient high accuracy single-molecule localization microscopy. *Nature Communications*, 16(1):7217, 2025.
    - Aude Genevay, Gabriel Peyré, and Marco Cuturi. Learning generative models with sinkhorn divergences. In *International Conference on Artificial Intelligence and Statistics*, pp. 1608–1617. PMLR, 2018.
    - Ross Girshick, Jeff Donahue, Trevor Darrell, and Jitendra Malik. Rich feature hierarchies for accurate object detection and semantic segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 580–587, 2014.
    - Mats GL Gustafsson. Surpassing the lateral resolution limit by a factor of two using structured illumination microscopy. *Journal of microscopy*, 198(2):82–87, 2000.
    - Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770–778, 2016.
    - Mike Heilemann, Sebastian Van De Linde, Mark Schüttpelz, Robert Kasper, Britta Seefeldt, Anindita Mukherjee, Philip Tinnefeld, and Markus Sauer. Subdiffraction-resolution fluorescence imaging with conventional fluorescent probes. *Angewandte Chemie-International Edition*, 47(33), 2008.
    - Stefan W Hell and Jan Wichmann. Breaking the diffraction resolution limit by stimulated emission: stimulated-emission-depletion fluorescence microscopy. *Optics letters*, 19(11):780–782, 1994.
    - Samuel T Hess, Thanu PK Girirajan, and Michael D Mason. Ultra-high resolution imaging by fluorescence photoactivation localization microscopy. *Biophysical journal*, 91(11):4258–4272, 2006.
    - Bo Huang, Sara A Jones, Boerries Brandenburg, and Xiaowei Zhuang. Whole-cell 3d storm reveals interactions between cellular structures with nanometer-scale resolution. *Nature methods*, 5(12): 1047–1052, 2008.
    - Junhwa Hur and Stefan Roth. Iterative Residual Refinement for Joint Optical Flow and Occlusion Estimation, 2019.
    - Alex Kendall, Yarin Gal, and Roberto Cipolla. Multi-task learning using uncertainty to weigh losses for scene geometry and semantics. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 7482–7491, 2018.
    - Harold W Kuhn. The hungarian method for the assignment problem. *Naval research logistics quarterly*, 2(1-2):83–97, 1955.
    - Romain F Laine, Hannah S Heil, Simao Coelho, Jonathon Nixon-Abell, Angélique Jimenez, Theresa Wiesner, Damián Martínez, Tommaso Galgani, Louise Régnier, Aki Stubb, et al. High-fidelity 3d live-cell nanoscopy through data-driven enhanced super-resolution radial fluctuation. *Nature Methods*, 20(12):1949–1956, 2023.
    - Mickaël Lelek, Melina T Gyparaki, Gerti Beliu, Florian Schueder, Juliette Griffié, Suliana Manley, Ralf Jungmann, Markus Sauer, Melike Lakadamyali, and Christophe Zimmer. Single-molecule localization microscopy. *Nature reviews methods primers*, 1(1):39, 2021a.
    - Mickaël Lelek, Melina T. Gyparaki, Gerti Beliu, Florian Schueder, Juliette Griffié, Suliana Manley, Ralf Jungmann, Markus Sauer, Melike Lakadamyali, and Christophe Zimmer. Single-molecule localization microscopy. *Nature Reviews Methods Primers*, 1:1–27, 2021b.
    - Feng Li, Hao Zhang, Shilong Liu, Jian Guo, Lionel M Ni, and Lei Zhang. Dn-detr: Accelerate detr training by introducing query denoising. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 13619–13627, 2022.

- Xinyang Li, Xiaowan Hu, Xingye Chen, Jiaqi Fan, Zhifeng Zhao, Jiamin Wu, Haoqian Wang, and Qionghai Dai. Spatial redundancy transformer for self-supervised fluorescence image denoising. *Nature Computational Science*, 3(12):1067–1080, 2023.
  - Yiming Li, Markus Mund, Philipp Hoess, Joran Deschamps, Ulf Matti, Bianca Nijmeijer, Vilma Jimenez Sabinina, Jan Ellenberg, Ingmar Schoen, and Jonas Ries. Real-time 3d single-molecule localization using experimental point spread functions. *Nature methods*, 15(5):367–369, 2018.
  - Yiming Li, Elena Buglakova, Yongdeng Zhang, Jervis Vermal Thevathasan, Joerg Bewersdorf, and Jonas Ries. Accurate 4pi single-molecule localization using an experimental psf model. *Optics Letters*, 45(13):3765–3768, 2020.
  - Sheng Liu, Jianwei Chen, Jonas Hellgoth, Lucas-Raphael Müller, Boris Ferdman, Christian Karras, Dafei Xiao, Keith A Lidke, Rainer Heintzmann, Yoav Shechtman, et al. Universal inverse modeling of point spread functions for smlm localization and microscope characterization. *Nature Methods*, 21(6):1082–1093, 2024.
  - Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. arXiv preprint arXiv:1711.05101, 2017.
  - Charles W. McCutchen. Superresolution in Microscopy and the Abbe Resolution Limit. *JOSA*, 57 (10):1190–1192, 1967.
  - Hamza Mentagui, Luca Calatroni, Sébastien Schaub, and Laure Blanc-Féraud. Physics-inspired generative adversarial modelling for fluctuation-based super-resolution microscopy. In 2024 IEEE International Symposium on Biomedical Imaging (ISBI), pp. 1–4. IEEE, 2024.
  - Grégoire Mialon, Dexiong Chen, Alexandre d'Aspremont, and Julien Mairal. A trainable optimal transport embedding for feature aggregation and its relationship to attention. In *ICLR*, 2021.
  - Elias Nehme, Lucien E Weiss, Tomer Michaeli, and Yoav Shechtman. Deep-storm: super-resolution single-molecule microscopy by deep learning. *Optica*, 5(4):458–464, 2018.
  - Wei Ouyang, Andrey Aristov, Mickaël Lelek, Xian Hao, and Christophe Zimmer. Deep learning massively accelerates super-resolution localization microscopy. *Nature biotechnology*, 36(5): 460–468, 2018.
  - Martin Ovesný, Pavel Křížek, Josef Borkovec, Zdeněk Švindrych, and Guy M. Hagen. ThunderSTORM: a comprehensive ImageJ plug-in for PALM and STORM data analysis and superresolution imaging. *Bioinformatics*, 30(16):2389–2390, August 2014.
  - George Patterson, Michael Davidson, Suliana Manley, and Jennifer Lippincott-Schwartz. Superresolution imaging using single-molecule localization. *Annual review of physical chemistry*, 61(1): 345–367, 2010.
  - Gabriel Peyré, Marco Cuturi, et al. Computational optimal transport: With applications to data science. *Foundations and Trends® in Machine Learning*, 11(5-6):355–607, 2019.
  - Patrick Putzky and Max Welling. Recurrent Inference Machines for Solving Inverse Problems, 2017.
  - Jonas Ries. Smap: a modular super-resolution microscopy analysis platform for smlm data. *Nature Methods*, 17(9):870–872, 2020.
  - Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *Medical image computing and computer-assisted intervention—MICCAI 2015: 18th international conference, Munich, Germany, October 5-9, 2015, proceedings, part III 18*, pp. 234–241. Springer, 2015.
  - Kurt Rossmann. Point spread-function, line spread-function, and modulation transfer function: tools for the study of imaging systems. *Radiology*, 93(2):257–272, 1969.
  - Michael J Rust, Mark Bates, and Xiaowei Zhuang. Sub-diffraction-limit imaging by stochastic optical reconstruction microscopy (storm). *Nature methods*, 3(10):793–796, 2006.

Daniel Sage, Thanh-An Pham, Hazen Babcock, Tomas Lukes, Thomas Pengo, Jerry Chao, Ramraj Velmuruga, Alex Herbert, Anurag Agrawal, Silvia Colabrese, Ann Wheeler, Anna Archetti, Bernd Rieger, Raimund Ober, Guy M. Hagen, Jean-Baptiste Sibarita, Jonas Ries, Ricardo Henriques, Michael Unser, and Seamus Holden. Super-resolution fight club: Assessment of 2D & 3D single-molecule localization microscopy software. *Nature methods*, 16(5):387–395, 2019.

Lukas Scheiderer, Zach Marin, and Jonas Ries. Minflux—molecular resolution with minimal photons. *arXiv preprint arXiv:2410.15902*, 2024.

Lothar Schermelleh, Alexia Ferrand, Thomas Huser, Christian Eggeling, Markus Sauer, Oliver Biehlmaier, and Gregor PC Drummen. Super-resolution microscopy demystified. *Nature cell biology*, 21(1):72–84, 2019.

Artur Speiser, Lucas-Raphael Müller, Philipp Hoess, Ulf Matti, Christopher J. Obara, Wesley R. Legant, Anna Kreshuk, Jakob H. Macke, Jonas Ries, and Srinivas C. Turaga. Deep learning enables fast and dense single-molecule localization with high accuracy. *Nature Methods*, 18(9): 1082–1090, 2021.

Cédric Villani. Topics in optimal transportation, volume 58. American Mathematical Soc., 2021.

Changdong Yu, Xiaojun Bi, and Yiwei Fan. Deep learning for fluid velocity field estimation: A review. *Ocean Engineering*, 271:113693, 2023.

Masoumeh Zareapoor, Pourya Shamsolmoali, Huiyu Zhou, Yue Lu, and Salvador García. Fractional Correspondence Framework in Detection Transformer, 2025.

Hao Zhang, Feng Li, Shilong Liu, Lei Zhang, Hang Su, Jun Zhu, Lionel M Ni, and Heung-Yeung Shum. Dino: Detr with improved denoising anchor boxes for end-to-end object detection. *arXiv* preprint arXiv:2203.03605, 2022.

Xizhou Zhu, Weijie Su, Lewei Lu, Bin Li, Xiaogang Wang, and Jifeng Dai. Deformable detr: Deformable transformers for end-to-end object detection. *arXiv preprint arXiv:2010.04159*, 2020.

# A APPENDIX

#### A.1 IMAGE FORMATION PROCESS

SMLM experimental setups typically employ either Electron-Multiplying CCD (EM-CCD) or scientific CMOS (sCMOS) cameras. Their sensors converts incident photons into a digital intensity value (ADU) through a sequence of physical processes, each of which introduces noise.

Let n be the incident photon count on the camera sensor. Initially, photon detection is modeled as a Poisson process — known as *shot noise* — with a mean proportional to n and the quantum efficiency (QE), and an offset known as the spurious charge (c):

$$n_{e,1} \sim \mathcal{P}\left(\text{QE} \times n + c\right).$$
 (9)

EM-CCD cameras introduce an additional amplification stage, modeled as a Gamma distribution with parameters  $n_{e,1}$  and the electromagnetic gain (EM):

$$n_{e,2} \sim \Gamma(n_{e,1}, \text{EM})$$
 for EM-CCD, or  $n_{e,2} = n_{e,1}$  for sCMOS. (10)

Subsequently, read noise is modeled by a normal distribution with mean  $n_{e,2}$  and standard deviation  $\sigma_R$ :

$$n_{e,3} \sim \mathcal{N}(n_{e,2}, \sigma_R). \tag{11}$$

Finally, the analog-to-digital conversion process yields the observed ADU, scaled by the electrons per ADU ( $e_{ADU}$ ) and offset by a baseline (B):

$$y = \min\left(\left\lfloor \frac{n_{e,3}}{e_{\text{ADU}}} \right\rfloor + B , 65535\right)$$
 (12)

Table 4 presents the parameters for two cameras commonly used in SMLM. Evolve Delta 512 is used by in the Tubulin and NPC-Nup107 datasets (Li et al., 2018), while Dhyana 400BSI V3 is used in the NPC-Nup96 dataset (Fei et al., 2025).

702
703
704
705
706
707
708
709

Parameter	Evolve Delta 512	Dhyana 400BSI V3
Camera type	EMCCD	sCMOS
Quantum efficiency (QE)	0.90	0.95
Spurious charge (c)	0.002	0.002
EM gain (EM)	300	_
Readout noise $(\sigma_R)$	74.4	1.535
Electrons per $\overrightarrow{ADU}$ ( $e_{ADU}$ )	45	0.7471
ADU baseline (B)	100	100

Table 4: Parameters for two typical cameras used in SMLM.

## A.2 ARCHITECTURE DETAILS

**Decoder architecture.** The decoder is responsible for mapping a  $C \times H \times W$  tensor to a  $d \times 5$  matrix, therefore its architecture is, and foremost, driven by the choice of d. In DECODE, Speiser et al. (2021) picked  $d = H \times W$  and predicts one candidate per pixel in the original frame  $\boldsymbol{y}$ . This led to an enormous amount of predictions, even at extremely high densities. Conversely, object detection using transformer architectures like DETR (Carion et al., 2020) typically predicts 100 candidates, which is an order of magnitude larger than the expected number of objects to detect in a normal object detection task. We picked a middle ground and chose  $d = H \times W/4$ . We found that one candidate per four original pixels offers an effective balance between coverage and computational cost.

Therefore, our decoder is formally defined as  $D: \mathbb{R}^{C \times H \times W} \longmapsto \mathbb{R}^{5 \times H/2 \times W/2}$ , mapping a latent variable to a  $H/2 \times W/2$  map with 5 channels, where each pixel is an activation prototype.

Consider a single pixel i of D's output, and let  $(\tilde{x}_i, \tilde{y}_i)$  be its 2D coordinates in the camera coordinate system. The five elements output for this pixel encode the characteristics of the underlying candidate activation: the detection score  $\hat{s}_i$ , the relative lateral coordinates  $(\Delta \hat{x}_i, \Delta \hat{y}_i)$ , the depth  $\hat{z}_i$  and the number of emitted photons  $\hat{n}_i$ . The absolute lateral coordinates  $(\hat{x}_i, \hat{y}_i)$  are reconstructed by summing  $(\Delta \hat{x}_i, \Delta \hat{y}_i)$  with  $(\tilde{x}_i, \tilde{y}_i)$ . The magnitude of the relative coordinate offsets  $(\Delta \hat{x}, \Delta \hat{y})$  predicted by the decoder is set to three times the pixel size. This extended range permits multiple activations to be mapped within a single pixel area, as neighbouring activations can contribute to their surroundings.

Finally, the output is formatted into a candidate set  $\hat{\mathcal{X}} = \{(\hat{x}_i, \hat{y}_i, \hat{z}_i, \hat{n}_i)\}_{1 \leq i \leq d}$  and a detection scores set  $\mathcal{S} = \{\hat{s}_i\}_{1 \leq i \leq d}$ . We integrate this reconstruction process into the decoder, meaning  $D(z) = (\hat{\mathcal{X}}, \hat{\mathcal{S}})$ .

The practical implementation of the decoder consists sequentially of a 2×2 max pooling layer, followed by a single residual block (He et al., 2016) and a linear projection.

**Differentiable simulation within our model.** During inference, our algorithm selects a subset of candidate detections by thresholding their confidence scores. However, this operation is non-differentiable, preventing direct gradient propagation during training. To mimic this behaviour while retaining differentiability, we replace it by a soft weighting that scales the photon count of each candidate by its detection confidence. For each candidate  $x_i$  in  $\hat{\mathcal{X}}$ , the network outputs the 3D coordinates  $(\hat{x}_i, \hat{y}_i, \hat{z}_i)$ , the raw photon count  $\hat{n}_i$ , and a detection confidence  $\hat{s}_i \in (0, 1)$ . We choose to modulate the photon count by the confidence, producing the weighted activation

$$\tilde{\boldsymbol{x}}_i = (\hat{x}_i, \, \hat{y}_i, \, \hat{z}_i, \, \hat{s}_i \, \hat{n}_i),$$

and the set of all such activations is denoted  $\tilde{\mathcal{X}} = \{\tilde{x}_i\}_{i=1}^d$ . This causes activations with low detection scores to have a number of emitted photons near zero, making them almost non-existent, while keeping almost untouched activations with a detection score close to one, mimicking the effect of a hard threshold while remaining fully differentiable.

After derivation, the expected image  $\hat{y}$  is obtained by:

$$\hat{\mathbf{y}} = \mathbb{E}[\hat{\mathbf{y}}|\tilde{\mathcal{X}}] = \frac{\text{QE} \times \text{EM}}{e_{\text{ADU}}}\mathbf{H}(\tilde{\mathcal{X}}) + B,$$
(13)

Subpart	Multiply-Accumulate operations	Parameters
Encoder	1.03 GMac	1.05 M
Decoder	512.56 MMac	499.4 k
Residual Network	1.87 GMac	1.26 M
Renderer	94.64 MMac	0

Table 5: Multiply-Accumulate operations and number of parameters of our model subparts.

and with EM = 1 for sCMOS camera.  $\tilde{y}$  is an end-to-end differentiable approximation of the reconstructed output, and can be used inside our iterative refinement scheme.

# A.3 COMPUTATIONAL FOOTPRINT

 Training is performed using an NVIDIA H100 GPU using the AdamW optimizer with a learning rate of  $4 \times 10^{-4}$ , a weight decay of 0.01, and a cosine annealing scheduler. We chose a batch size of 128 to maximize GPU usage, filling all 80 GB of VRAM. It can be lowered using smaller batch sizes or gradient accumulation.

We trained for 14 hours 100 epochs of 1024 steps each, totaling approximately 100,000 steps. Excellent results ( $E_{3D} \ge 0.72$  on EPFL's density=2.0 and high SNR dataset) are achieved after only 20 minutes of training, at around 2000 steps.

During inference, a batch size of 16 produces a peak VRAM usage of  $8.7\,\mathrm{GB}$  and processes 2500 64x64 images in  $30\,\mathrm{s}$ , or  $12\,\mathrm{ms/frame}$ .

Table 5 shows an overview of the computational resources for each subpart of our model.