

XLFM-Former: Rapid 3D Reconstruction for Light Field Microscopy

Anonymous ICCV submission

Paper ID 8604

Abstract

001 *Light field microscopy (LFM) has become an emerging*
002 *tool in neuroscience for large-scale neural imaging*
003 *in vivo, with XLFM (eXtended Light Field Microscopy)*
004 *notable for its single-exposure volumetric imaging, broad*
005 *field of view, and high temporal resolution. However,*
006 *3D reconstruction of neural activity from multi-view light*
007 *field data is slow and computationally expensive, hindering*
008 *real-time brain dynamics observation and broader applica-*
009 *tions. Deep learning could provide faster and more*
010 *accurate solutions, but advancements are hampered by a*
011 *lack of datasets, benchmarks, and effective architectures.*
012 *We present eXtended Light Field Microscopy-Transformer*
013 *(XLFM-Former), a transformer-based framework for real-*
014 *time (> 30 volumes/second) and high-quality XLFM recon-*
015 *struction. It uses a Swin Transformer encoder for hierar-*
016 *chical feature extraction and progressive upsampling for cross-*
017 *level fusion. Key innovations include: (1) Masked View*
018 *Modeling-Light Field (MVM-LF) pretraining for global*
019 *view representation; (2) Physics-guided Point Spread Func-*
020 *tion (PSF) loss for optical consistency. When tested on*
021 *the unique XLFM-zebrafish dataset, XLFM-Former outper-*
022 *forms current methods, achieving 7.7% higher PSNR while*
023 *maintaining fine neural structures. To support future re-*
024 *search, we release the benchmark dataset and evaluation*
025 *metrics. Code and dataset available at: xxx.*

026 1. Introduction

027 Light Field Microscopy (LFM) has emerged as a crucial
028 technique for rapid volumetric imaging of nervous systems
029 [19, 20, 36]. Notably, eXtended Light Field Microscopy
030 (XLFM) [8], due to its graceful balance between speed,
031 scale and resolution, is considered one of the most suitable
032 LFM techniques for large-scale neural activity recording
033 in several model organisms, including fish and mouse [2].
034 XLFM offers several advantages: 1) XLFM enables single-
035 exposure acquisition of complete light field information at
036 100 Hz, whereas conventional microscopy techniques (e.g.,
037 two-photon microscopy [11], light-sheet microscopy [17])

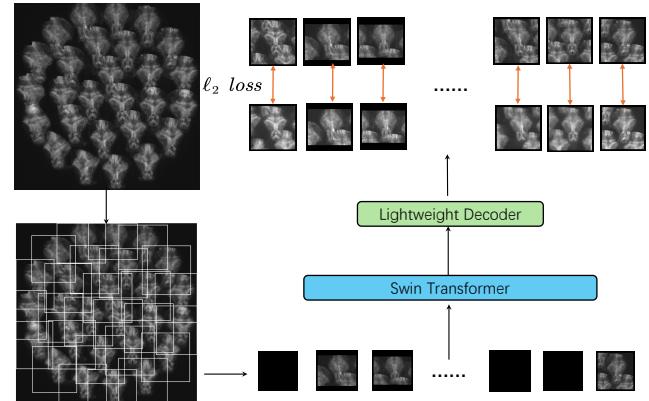


Figure 1. Our pretraining pipeline for XLFM. The raw light field acquired from the microscope is separated into 27 distinct viewpoints based on physical coordinates. With a 70% probability, we randomly mask a subset of these viewpoints and task the model with reconstructing them. The training is supervised by an ℓ_2 loss comparing the predicted and ground-truth views.

038 require sequential layer-by-layer scanning, making it chal-
039 lengeing to capture sub-second large-scale neural dynamics
040 simultaneously. 2) The XLFM system incorporates a point
041 spread function (PSF) that is approximately spatially invari-
042 ant. Consequently, the reconstruction of volumes through
043 3D deconvolution is free from artifacts. 3) The rapid speed
044 of XLFM allows real-time observation of large-scale popu-
045 lation neural activity *in vivo*. Integrating volumetric imag-
046 ing with optogenetic manipulation [4, 35] will enable opti-
047 cal brain-machine interface, namely closed-loop optical in-
048 terrogation of brain-wide activity in both immobilized [28]
049 and freely behaving animals [5, 8].

050 However, the distinctive optical sampling method em-
051 ployed by XLFM causes a speed bottleneck in 3D recon-
052 struction, not data acquisition. Traditional approaches, such
053 as iterative Richardson-Lucy deconvolution, are inherently
054 slow, necessitating the development of more efficient re-
055 construction algorithms that can integrate both the physics
056 priors and the expressiveness inherent in neural network
057 models. Despite its great potential, existing deep learning

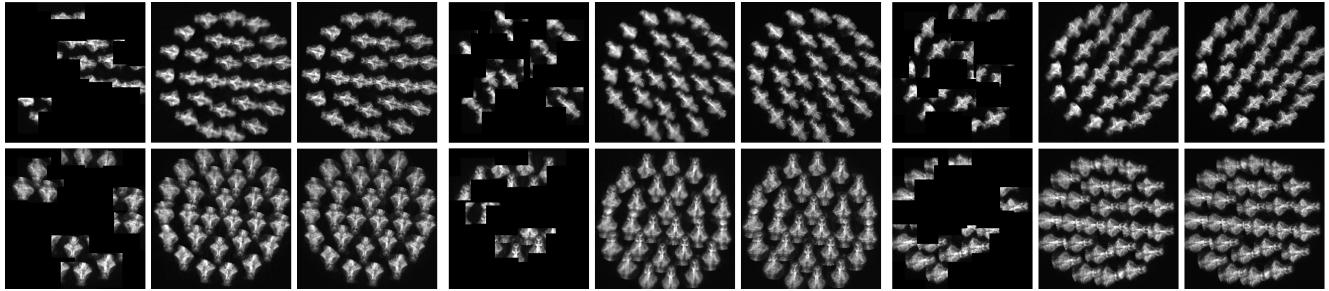


Figure 2. **Zebrafish multi-view light field images used for pretraining an encoder model.** For each triplet, we show the masked image (left), our MVM-LF regenerated image (middle), and the ground-truth (right).

approaches still face several challenges in XLFM reconstruction tasks. Firstly, high-quality annotation is costly. While XLFM sampling data is readily available, generating high-quality 3D reconstruction annotations (e.g., via Richardson-Lucy deconvolution [24]) is highly expensive, limiting the scalability of supervised learning methods. Secondly, multi-view geometric modeling remains insufficient. Although XLFM captures light field information using a microlens array, existing approaches struggle to fully exploit multi-view data for enhancing 3D reconstruction quality. Additionally, the lack of physical-optical priors poses another challenge. Most existing methods rely purely on data-driven approaches, neglecting the physical constraints of microscopic optical imaging, which may lead to reconstruction results that deviate from optical imaging principles.

To address these issues, we propose eXtended Light Field Microscopy-Transformer (XLFM-Former), an end-to-end Transformer framework specifically optimized for XLFM tasks. XLFM-Former employs hierarchical feature encoding and a window attention mechanism to enhance its modeling capability. The XLFM-Former encoder, constructed based on the window attention mechanism, extracts semantic features at four different resolutions and scales, ensuring effective modeling of both local and global information. Subsequently, a CNN-based decoder is employed for feature fusion to restore high-resolution 3D structural information. The decoder is followed by a reconstruction head, which directly computes the final 3D reconstruction results. To further enhance model performance, we introduce a self-supervised pretraining strategy for the Transformer encoder. Specifically, since XLFM acquires light field data from 27 different viewpoints through the microlens array, we propose a novel self-supervised learning task—Masked View Modeling-Light Field (MVM-LF)—to improve the model’s understanding of multi-view light field information. During pre-training, 70% of the viewpoint data is randomly masked, and the model is provided with only 30% of the viewpoints, requiring it to predict the content of the masked viewpoints (See Figure 1 for an

overview). This encourages the model to learn the underlying geometric relationships within the light field data. By leveraging the intrinsic structural priors of XLFM, this approach reduces reliance on costly annotations while significantly improving model generalization. To get a qualitative sense of our reconstruction task, see Figure 2 and 3.

Beyond data-driven learning, we introduce a PSF-based physical constraint loss to ensure that the reconstructed 3D structures conform to the optical system’s physical principles. The Point Spread Function (PSF) [1], derived from the microscope’s optical system, characterizes the diffraction and spreading behavior of a point source within the imaging system. Unlike conventional loss functions (e.g., L1/L2, SSIM), the PSF loss function forward-projects the 3D reconstruction results back into the original light field, ensuring that model predictions adhere to the physical imaging process. This enables XLFM-Former to effectively learn the optical characteristics of light field data, improving reconstruction quality while minimizing artifacts that violate optical principles. Experimental results demonstrate that XLFM-Former outperforms existing state-of-the-art frameworks, achieving superior reconstruction quality on XLFM tasks.

Our main contributions are summarized as follows:

- We establish a standardized framework for the XLFM reconstruction task and construct the first benchmark dataset for XLFM reconstruction. For the first time, we present a large-scale XLFM dataset comprising 22,581 images, including three free-swimming zebrafish, seven immobilized zebrafish, and six test zebrafish, captured at varying acquisition rates (10 fps / 1 fps). This dataset establishes a standardized and reproducible benchmark, facilitating the evaluation and advancement of XLFM reconstruction methodologies.
- We design the eXtended Light Field Microscopy-Transformer (XLFM-Former) framework, integrating Swin Transformer’s hierarchical modeling capabilities to effectively capture both local and global information in XLFM data while adapting to light field viewpoint modeling tasks.

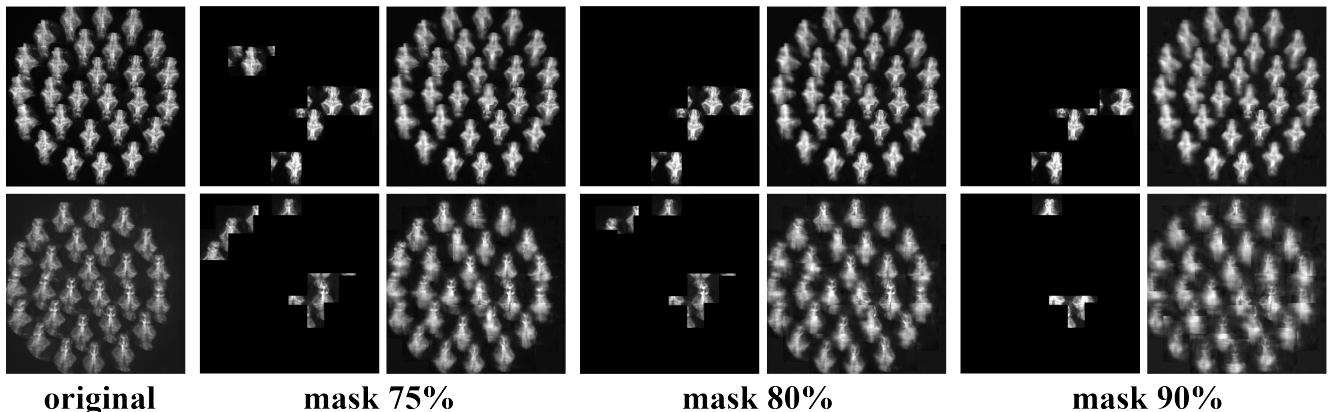


Figure 3. **Regeneration of XLFM light field images via MVM-LF.** The model can still accurately predict the view under appropriate occlusion, indicating that it has learned the global view relationship. Excessive occlusion (90%) causes prediction to crash, indicating that MVM-LF requires a reasonable occlusion ratio to balance information loss and network learning ability.

- Based on XLFM-Former, we propose the Masked View Modeling-Light Field (MVM-LF) self-supervised learning task, which leverages light field viewpoint relationships for masked prediction, improving the model’s geometric understanding at a low cost.
- We introduce a PSF-based physical prior loss function, which explicitly models the optical imaging process of light field microscopy, guiding the network to learn reconstruction results that adhere to the optical system’s characteristics, thereby enhancing physical consistency and generalization capability.

2. Related Work

2.1. Unsupervised Pretraining Methods for Light Field Microscopy

In the computer vision community, unsupervised pretraining methods have gained widespread attention [3, 13, 14, 21, 25]. For example, Masked Autoencoders [16] learn by randomly masking parts of the input to help the model understand spatial relationships and global context. Contrastive learning [6, 7, 9, 18, 30, 33] creates positive and negative pairs of samples to bring similar samples closer and push dissimilar ones apart. However, in the context of LFM, unsupervised methods are still underexplored. A recent approach [34], Masked LF Modeling (MLFM), introduces a self-supervised pre-training scheme to enhance Light Field Super-Resolution (LFSSR). This method uses a transformer-based structure, LFormer, to learn inter-view correlations. While this approach significantly improves performance, its reliance on random masking does not fully capture the interdependencies between views, which are critical for high-quality super-resolution. Although both approaches are applied to light fields, our proposed method is fundamentally different. First, we focus on the specific task

of XLFM reconstruction. Second, our approach is based on viewpoint reconstruction, whereas theirs relies on random pixel masking.

2.2. 3D Reconstruction in XLFM

Recent work [10, 26, 31] has demonstrated the potential of deep learning in addressing computational bottlenecks in XLFM. A recent approach [31] combines two neural networks, SLNet and XLFMNet, for real-time sparse 3D volumetric reconstruction in light field microscopy. SLNet extracts the spatio-temporally sparse components from image sequences, while XLFMNet performs high-fidelity 3D reconstruction. Another recent approach [26] proposes using a conditional normalizing flow architecture for fast 3D reconstruction of neural activity in immobilized zebrafish. However, like XLFMNet, this method remains constrained by its sparsity-driven approach, reconstructing only neural signals while disregarding complete biological morphology. Unlike these prior methods, XLFM-Former is designed for full-volume imaging, reconstructing not only neural activity but also entire volumetric structures. By leveraging a Swin Transformer backbone for hierarchical feature extraction and MVM-LF pretraining to learn global context dependencies, XLFM-Former provides a comprehensive reconstruction of biological samples. This distinction is critical in applications where both functional (neural signals) and anatomical (morphological structures) information are necessary for deeper biological insights.

A recent end-to-end approach [10] combines differentiable simulations of optical systems with deep learning-based reconstruction networks for high-performance computational imaging. The key insight is that global information is crucial for such problems, which is achieved by using Fourier-Nets, a shallow neural network architecture based on global kernel Fourier convolution. However, mapping

205 to the Fourier domain results in a substantial increase in
206 memory usage, requiring multiple GPUs for large-volume
207 reconstruction. This limits the method's scalability and ap-
208 plicability to more general imaging tasks. This method is
209 particularly expensive in terms of video memory and is not
210 suitable for XLFM reconstruction because the final output
211 of XLFM exceeds 100 million pixels and cannot be made
212 into patches due to system design issues.

213 To achieve such global information extraction: 1) we use
214 the Swin Transformer [22] as a feature extraction module
215 and then apply a CNN-based decoder for feature fusion. Us-
216 ing self-attention is more efficient than convolution mapped
217 to the Fourier domain. 2) To force the network to under-
218 stand the dependencies between different views, we propose
219 a proxy task. By masking 70% of the views, we force the
220 network to reconstruct the masked views, enabling unsuper-
221 vised pretraining.

222 3. XLFM-Zebrafish Dataset

223 3.1. Data Collection

224 To construct the XLFM-Zebrafish Dataset, we utilized
225 an advanced XLFM system designed to capture high-
226 resolution volumetric neural activity in zebrafish. The data
227 collection process was carefully structured to ensure di-
228 versity in motion states, imaging conditions, and biolog-
229 ical variability. For free-swimming zebrafish, we recorded
230 neural activity in an unconstrained environment, allowing
231 for the study of brain-wide dynamics during naturalistic
232 behaviors. A real-time tracking system was employed to
233 continuously adjust the imaging field of view, ensuring
234 that the zebrafish remained within the microscope's focal
235 range. Additionally, motion artifacts caused by rapid move-
236 ment were mitigated through dual-color fluorescence im-
237 aging and adaptive filtering techniques. In contrast, fixed ze-
238 brafish were embedded in a stabilizing medium to facili-
239 tate high-precision 3D structural reconstruction. This set-
240 ting eliminated motion-induced distortions, enabling the ex-
241 traction of detailed neural architecture. The immobilized
242 specimens were further divided into different groups for
243 training, validation, and testing purposes, ensuring a struc-
244 tured dataset for benchmarking reconstruction algorithms.
245 To capture both rapid neural dynamics and long-term activ-
246 ity trends, we employed multiple imaging conditions that
247 varied in temporal resolution and sampling strategies. This
248 approach allowed us to balance high-fidelity reconstruction
249 with the need for extended observation periods.

250 3.2. Dataset Statistics

251 To construct a high-quality XLFM-Zebrafish Dataset, we
252 collected zebrafish in different motion states and set various
253 sampling conditions to ensure diversity and applicability.
254 This dataset is the first standardized XLFM zebrafish 3D re-

255 construction dataset, designed to evaluate the performance
256 of deep learning models in XLFM-3D reconstruction tasks.
257 The XLFM-Zebrafish Dataset consists of two categories of
258 zebrafish data: Free-swimming Zebrafish, includes 3 indi-
259 vidual zebrafish, used for studying dynamic neural activity
260 and analyzing the impact of motion blur on 3D reconstruc-
261 tion. Fixed Zebrafish, includes 13 individual zebrafish, suit-
262 able for high-precision 3D structural reconstruction, with
263 7 individuals for training and validation, and 6 for test-
264 ing. Additionally, we introduced two different sampling
265 rates: 10fps (High sampling rate): Suitable for temporal neu-
266 ral activity modeling and high-precision light field recon-
267 struction. 1fps (Low sampling rate): Used for long-term
268 dynamic tracking and reconstruction stability analysis un-
269 der low frame rates. The dataset comprises 22,581 light
270 field images. The detailed statistics are presented in the
271 Supplementary Section 7.

272 4. Methodology

273 4.1. XLFM-Former

274 Figure 4 illustrates the overall architecture of XLFM-
275 Former. We describe the details of encoder and decoder
276 in this subsection.

277 **XLFM-Former Encoder:** The input to XLFM-Former
278 is a raw light field captured by the XLFM system. The
279 original light field data is represented as a 2D array of sub-
280 aperture views with a resolution of $N \times N$, where N denotes
281 the number of views along both horizontal and vertical di-
282 rections. Each sub-aperture view contains spatial informa-
283 tion about the observed scene, forming a multi-view represen-
284 tation. To align the data with the optical coordinate sys-
285 tem and extract relevant information, a cropping operation
286 is applied to each sub-aperture view of cropped light field
287 size $\mathcal{X} \in \mathbb{R}^{H \times W \times D \times S}$. H and W denote the height and
288 width of each cropped sub-aperture view. D represents the
289 depth dimension obtained from the cropping process, which
290 corresponds to the axial resolution. S refers to the number
291 of channels, which is set to $S = 1$ in XLFM light field imag-
292 ing, producing a 3D view token representation. To facili-
293 tate self-attention computation in the transformer-based en-
294 coder, the patch partitioning layer divides the input volume
295 into a sequence of non-overlapping 3D View Tokens, each
296 of size (H', W', D') . The resulting tokenized feature map
297 has dimensions: $\frac{H}{H'} \times \frac{W}{W'} \times \frac{D}{D'}$, where each token is pro-
298 jected into a C -dimensional embedding space using a linear
299 embedding layer. This patch-wise tokenization allows the
300 transformer to process spatial and depth information effi-
301 ciently, enabling effective modeling of local and global de-
302 pendencies in the light field data. The encoder consists of
303 four hierarchical stages, where each stage contains two con-
304secutive Swin Transformer blocks. At each stage, a patch
305 merging layer is applied to downsample the resolution while

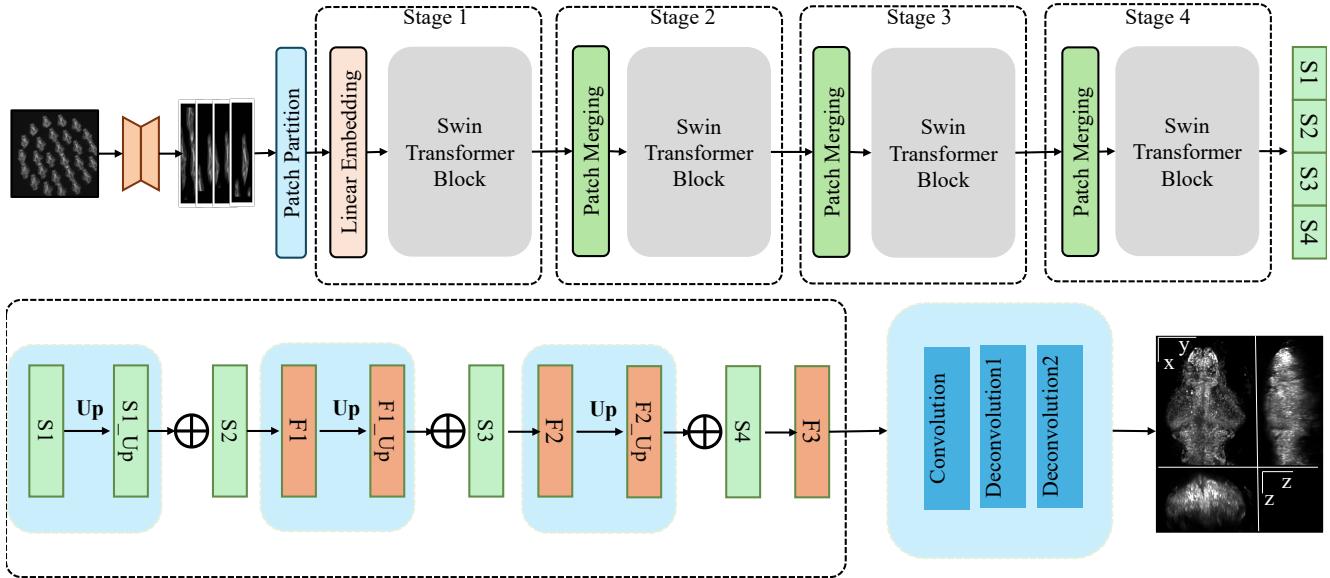


Figure 4. Overview of the Swin-XLFM architecture.

increasing feature dimensionality. Given an input feature map $S \in \mathbb{R}^{h \times w \times d \times c}$ at stage T , the patch merging operation first groups adjacent patches of size $2 \times 2 \times 2$ and concatenates their features where $S' \in \mathbb{R}^{\frac{h}{2} \times \frac{w}{2} \times \frac{d}{2} \times 8c}$ represents the concatenated feature map. A linear projection layer is then applied to transform the feature dimension to $2c$, resulting in: $S_{T+1} = \mathbf{W}S' + \mathbf{b}$, where \mathbf{W} is the learnable weight matrix, and \mathbf{b} is the bias term. Subsequently, the four feature vectors S_i , $i \in \{1, 2, 3\}$ corresponding to the four stages are sent to the decoder for feature fusion.

Decoder: The decoder reconstructs the final high-resolution 3D structure by progressively integrating multi-scale features extracted from the encoder. This process consists of two key components: progressive upsampling and cross-level feature fusion. The decoder takes the hierarchical feature maps from the encoder, denoted as S_i , $i \in \{1, 2, 3\}$, and reconstructs the high-resolution output through a sequence of deconvolutional operations. The highest-level feature map S_1 serves as the starting point, and at each stage, an upsampling operation is applied to gradually recover spatial details. Fused feature maps F_i , $i \in \{1, 2, 3\}$ are obtained by adding features of different scales:

$$\hat{F}_0 = S_1, \quad \hat{F}_i = \text{Up}(\hat{F}_{i-1}) + S_i, \quad i \in \{1, 2, 3\}. \quad (1)$$

The final reconstructed 3D structure is obtained by applying a reconstruction head, which maps the fused feature F_3 to the target 3D volume: First, a convolution is performed to match the number of channels:

$$\hat{F} = \text{Conv}_{1 \times 1}(F_{\text{fusion}}), \quad (2)$$

Then two deconvolutions are performed to match the target resolution:

$$\mathcal{V}_{\text{pred}} = \text{Deconv}_2 \left(\text{Deconv}_1 \left(\text{Conv}_{1 \times 1}(\hat{F}) \right) \right). \quad (3)$$

4.2. Physics-Guided PSF Reprojection Loss

In order to enforce that the reconstructed volume adheres to the underlying optical physics of the XLFM system, we introduce a physics-guided loss based on the system's Point Spread Function (PSF). The PSF, denoted as h , characterizes the optical blurring inherent in the imaging process.

Given a volumetric reconstruction \mathcal{V} (either the predicted volume $\mathcal{V}_{\text{pred}}$ or the ground truth \mathcal{V}_{GT}), we simulate the corresponding raw light field measurement by convolving the volume with the PSF:

$$\mathbf{I}_{\text{sim}} = h * \mathcal{V}, \quad (4)$$

where $*$ denotes convolution.

Accordingly, we compute the simulated measurements for both the predicted volume and the ground truth:

$$\mathbf{I}_{\text{pred}} = h * \mathcal{V}_{\text{pred}}, \quad \mathbf{I}_{\text{GT}} = h * \mathcal{V}_{\text{GT}}. \quad (5)$$

The Physics-Guided PSF Reprojection Loss is then defined as the mean squared error (MSE) between these two simulated measurements:

$$\mathcal{L}_{\text{PSF}} = \|\mathbf{I}_{\text{pred}} - \mathbf{I}_{\text{GT}}\|_2^2 = \|h * \mathcal{V}_{\text{pred}} - h * \mathcal{V}_{\text{GT}}\|_2^2. \quad (6)$$

By minimizing \mathcal{L}_{PSF} , we encourage the network to generate 3D reconstructions that not only approximate the ground truth in the volumetric domain but also conform to the physical imaging process modeled by the PSF.

361 4.3. Masked View Modeling for Light Fields (MVM- 362 LF)

363 To enhance self-supervised learning and inter-view modeling
364 in XLFM, we propose Masked View Modeling for
365 Light Fields (MVM-LF) as a pretraining strategy, enabling
366 the model to reconstruct missing views and capture global
367 scene structures.

368 **Pretrained Encoder:** The encoder architecture and input
369 representation in MVM-LF are identical to those used
370 in XLFM reconstruction. This consistency ensures that the
371 learned features during pretraining are directly transferable
372 to the supervised reconstruction task. The pretrained encoder,
373 denoted as f_θ , serves as the initialization for the
374 XLFM reconstruction model.

375 **Lightweight Decoder:** Inspired by self-supervised
376 masked reconstruction frameworks, we adopt a lightweight
377 decoder consisting of a series of convolutional layers. The
378 decoder is responsible for predicting the missing views during
379 pretraining. Once training is completed, the decoder is
380 discarded, and only the pretrained encoder is retained for
381 fine-tuning in the XLFM reconstruction task.

382 **Masking Strategy:** The core principle of MVM-LF is
383 to randomly mask a proportion r_m of the input views and
384 force the network to reconstruct them based solely on the
385 unmasked views. This proxy task compels the model to
386 learn a joint representation of the global structure and inter-
387 view dependencies. Formally, given a set of sub-aperture
388 views:

$$389 \quad \mathcal{U} = \{U_1, U_2, \dots, U_{N_u}\}, \quad (7)$$

390 we define the masked subset as:

$$391 \quad \mathcal{U}_{\text{mask}} = \{U_i \mid i \in \mathcal{M}\}, \quad (8)$$

392 where \mathcal{M} is a randomly sampled index set satisfying $|\mathcal{M}| =$
393 $r_m N_u$ with $r_m = 0.7$ (i.e., 70% of the views are masked).
394 The network is trained to reconstruct the missing views as:

$$395 \quad \hat{\mathcal{U}}_{\text{mask}} = f_\theta(\mathcal{U} \setminus \mathcal{U}_{\text{mask}}). \quad (9)$$

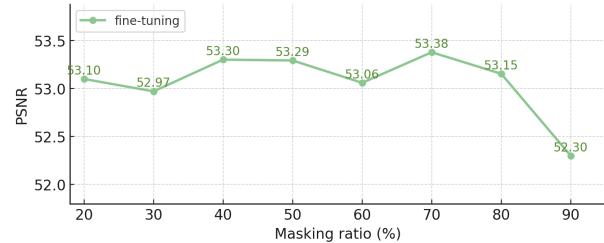
396 The loss function for MVM-LF is defined as the mean
397 squared error (MSE) between the predicted and ground
398 truth masked views:

$$399 \quad \mathcal{L}_{\text{MVM-LF}} = \sum_{U_i \in \mathcal{U}_{\text{mask}}} \|U_i - \hat{U}_i\|_2^2. \quad (10)$$

401 4.4. Loss Function

402 For the pre-training task, we only use ℓ_2 loss. For the XLFM
403 reconstruction task, we complete the reconstruction by min-
404 imizing the following loss combination.

$$405 \quad \begin{aligned} \mathcal{L}_{\text{total}} &= \frac{1}{\lambda_1} \mathcal{L}_{\text{MS_SSIM}} + \frac{1}{\lambda_2} \mathcal{L}_{\text{Edge}} + \frac{1}{\lambda_3} \mathcal{L}_{\text{PSNR}} \\ &+ \frac{1}{\lambda_4} \mathcal{L}_{\text{MSE}} + \frac{1}{\lambda_5} \mathcal{L}_{\text{PSF}}. \end{aligned} \quad (11)$$



406 **Figure 5. Masking ratio.**

407 The detailed loss function presented in the Supplementary
408 Section 8.

409 5. Experiments

410 5.1. Implementation Details

411 For the MVM-LF task, we employ a batch size of 8 to fa-
412 cilitate stable training dynamics. To enhance convergence
413 and mitigate the risk of the model becoming trapped in lo-
414 cal optima, we utilize the ReduceLROnPlateau learning rate
415 scheduler, with an initial learning rate set to 1e-4. The train-
416 ing process is conducted for 250 epochs to ensure robust
417 feature learning. All experiments are performed on a dis-
418 tributed computing setup with four NVIDIA A100-80GB
419 SMX4 GPUs. For the XLFM reconstruction task, the train-
420 ing configuration remains largely consistent, with the pri-
421 mary exception that the batch size is set to 1, aligning with
422 the requirements of volumetric reconstruction. The detailed
423 experimental setup presented in the Supplementary Section
424 9.

425 5.2. Main Results

426 We evaluate our approach against state-of-the-art archi-
427 tectures, including ConvNeXt, ViT, PVT, EfficientNet,
428 ResNet-50/101, and U-Net (Table 1). Our method consis-
429 tently outperforms all baselines across all evaluation met-
430 rics, achieving the highest PSNR (54.04 dB) and SSIM
431 (0.9944), significantly surpassing ConvNeXt (50.16 dB,
432 0.9876) and other transformer- and CNN-based models.
433 Figure 6 further highlights the qualitative advantages of our
434 model, producing sharper, more structurally accurate recon-
435 structions, while competing methods exhibit blurring, dis-
436 tortions, and detail loss. These results demonstrate that
437 physics-guided constraints and hierarchical transformer-
438 based modeling are key to achieving superior XLFM recon-
439 struction, both numerically and perceptually.

440 5.3. Ablation Study

441 **Masking ratio:** As shown in Figure 5, we analyze the im-
442 pact of different masking ratios in MVM-LF pretraining
443 on XLFM reconstruction. The results indicate that mod-
erate masking (50–70%) achieves optimal performance,

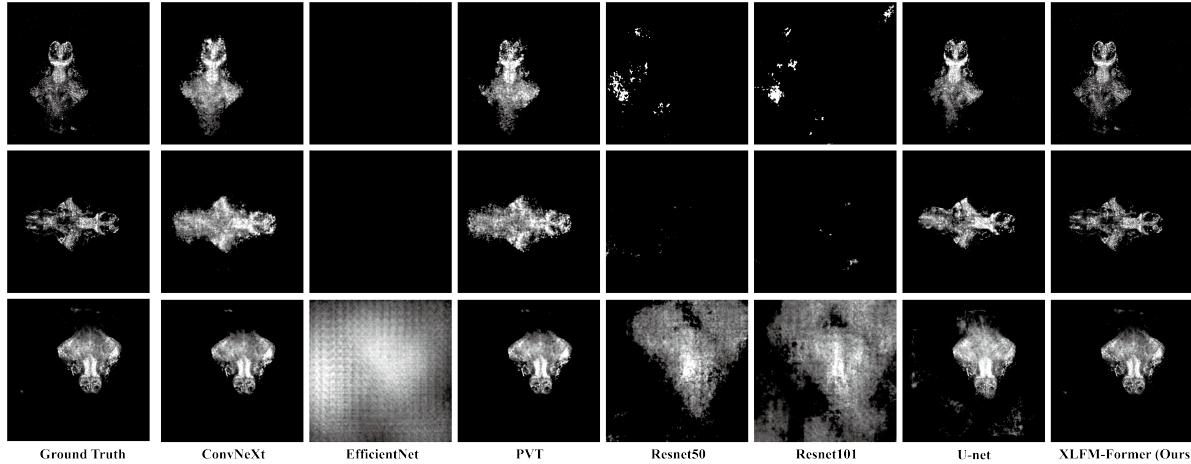


Figure 6. Compare with the state-of-the-art architecture on XLFM-Zebrafish Dataset.

Table 1. Comparison of Methods on XLFM-Zebrafish Dataset. The best results are highlighted in **bold**, while the second-best are underlined.

Method	# 1		# 2		# 3		# 4		# 5		# 6		Avg.	
	PSNR↑	SSIM↑												
ConvNeXt [23]	49.48	<u>0.9851</u>	53.88	0.9867	44.87	0.9833	51.38	0.9882	51.52	0.9892	49.79	<u>0.9935</u>	50.16	<u>0.9876</u>
ViT [12]	49.38	0.9842	52.67	0.9895	<u>45.29</u>	<u>0.9834</u>	51.09	0.9888	51.35	0.9906	45.90	0.9893	49.28	0.9876
PVT [32]	47.21	0.9804	47.93	0.9760	44.50	0.9807	49.46	0.9851	48.32	0.9841	46.60	0.9910	47.34	0.9829
EfficientNet [29]	45.04	0.9550	54.68	0.9851	42.13	0.9541	49.56	0.9801	48.63	0.9772	27.16	0.7264	44.53	0.9296
ResNet-50 [15]	46.46	0.9688	54.89	0.9851	41.46	0.9388	49.47	0.9790	48.82	0.9786	39.98	0.9304	46.85	0.9634
ResNet-101 [15]	47.20	0.9728	54.90	0.9851	41.33	0.9266	49.47	0.9787	49.09	0.9800	39.50	0.8893	46.91	0.9554
U-Net [27]	48.81	0.9807	<u>57.23</u>	<u>0.9928</u>	44.41	0.9808	<u>52.61</u>	<u>0.9908</u>	<u>52.06</u>	<u>0.9904</u>	41.47	0.9725	49.43	0.9847
Ours	53.97	0.9930	59.83	0.9963	49.31	0.9910	54.55	0.9951	54.65	0.9955	51.95	0.9956	54.04	0.9944

with PSNR peaking at 53.38 dB at 70% masking. Lower masking ratios (20–30%) provide insufficient representation learning, leading to suboptimal fine-tuning results, while excessive masking (90%) reduces PSNR to 52.30 dB, indicating difficulty in reconstructing missing views with limited context. These findings highlight the importance of balancing information removal and reconstruction difficulty, and we adopt 70% masking as the default setting for maximal pretraining efficiency.

Efficacy of Pre-training: We assess data efficiency by comparing models trained from scratch and those with pre-training under varying labeled data proportions (Figure 7). Pretraining provides a significant boost, especially in low-data regimes, with a PSNR of 51.92 dB at 10% labeled data, surpassing the 50.73 dB of training from scratch. While the performance gap narrows as more labeled data is available, the pretrained model consistently outperforms the scratch-trained counterpart, even at 80% labeled data. These results confirm that MVM-LF pretraining enhances feature generalization, making the model more data-efficient and robust across different data availability scenarios.

Efficacy of different components: We conduct an ab-

lation study to assess the impact of PSF loss and MVM-LF pretraining (Table 2). The baseline model achieves 52.14 dB PSNR, while adding PSF loss improves reconstruction fidelity to 52.96 dB. MVM-LF pretraining further enhances global view dependency learning, reaching 53.38 dB. Integrating both components into the full model yields the highest performance (54.04 dB PSNR, 0.9944 SSIM), confirming that combining physics-based constraints with self-supervised pretraining results in the most effective XLFM reconstruction.

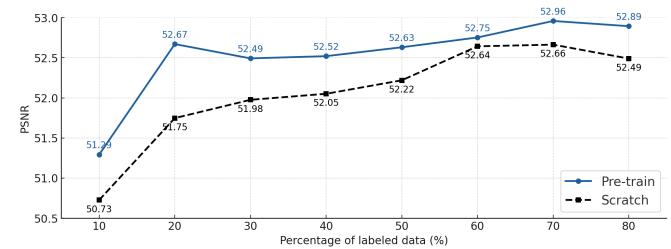


Figure 7. Data-efficient performance on XLFM-Zebrafish Test Dataset.

Table 2. Proposed components. 'Baseline' refers to the standard XLFM-Former model. 'PSF' indicates the inclusion of the physics-guided PSF loss. 'MVM-LF' corresponds to the incorporation of MVM-LF pretraining. 'Full' represents the complete model with both PSF loss and MVM-LF pretraining. The best results are highlighted in **bold**.

Method	PSNR↑	SSIM↑
baseline	52.1417	0.9924
PSF	52.9560	0.9931
MVM-LF	53.3772	0.9938
Full	54.0435	0.9944

Efficacy of Pretraining Strategies: We compare MVM-LF pretraining with alternative methods, including ImageNet-based initialization and pixel-level masked pretraining (Table 3). Training from scratch (Baseline) achieves 52.14 dB PSNR, while ImageNet-1k/22k pretraining provides only marginal improvements (52.70 dB / 52.38 dB), indicating that conventional pretraining is suboptimal for XLFM data. Random-masked pretraining performs slightly better (52.97 dB PSNR), but MVM-LF pretraining achieves the best results (54.04 dB PSNR, 0.9944 SSIM), demonstrating its superior ability to model multi-view dependencies. These findings highlight the importance of task-specific pretraining in optimizing XLFM reconstruction quality.

Generalization Under Reduced Input Views: We evaluate the pretrained model's ability to infer missing views by progressively reducing available inputs (Table 4). The scratch-trained model achieves 52.14 dB PSNR with full views, whereas the pretrained model surpasses it even with only 60% of views, reaching 52.54 dB PSNR. The best performance (53.26 dB PSNR) occurs at 80% input views, confirming strong generalization. These results demonstrate that MVM-LF pretraining enables robust multi-view reconstruction, allowing high-fidelity reconstruction even with incomplete input, making it highly adaptable to real-world imaging constraints.

6. Conclusion

We present XLFM-Former, a transformer-based framework for high-speed, high-fidelity XLFM volumetric reconstruction. It integrates a hierarchical Swin Transformer encoder and a physics-guided PSF loss to capture multi-scale spatial features while ensuring optical consistency. We further propose MVM-LF pretraining, a self-supervised strategy that improves global view dependency learning and missing view inference. Experiments on the XLFM-Zebrafish Dataset show that XLFM-Former surpasses SOTA models in PSNR and SSIM. Ablation studies validate the impact of

Table 3. Pretraining strategies. 'Baseline' refers to training from scratch. 'ImageNet 1k' and 'ImageNet 22k' indicate models initialized with ImageNet-1k and ImageNet-22k pretrained weights, respectively. 'Random Masked' represents pretraining with a randomly masked pixel reconstruction strategy. 'Ours' corresponds to the proposed MVM-LF pretraining. The best results are highlighted in **bold**.

Method	PSNR↑	SSIM↑
baseline	52.1417	0.9924
ImageNet 1k	52.6976	0.9931
ImageNet 22k	52.3770	0.9923
Random masked	52.9697	0.9934
Ours	54.0435	0.9944

Table 4. Missing views in XLFM reconstruction. The table reports PSNR and SSIM performance when the number of available input views is progressively reduced (from 90% to 60%). The baseline ('Scratch') represents a model trained from scratch with all input views. The results demonstrate that even when only 60% of the views are provided, the pretrained model surpasses the scratch-trained model using full input views, confirming its ability to infer missing views.

Input View Ratio	PSNR↑	SSIM↑
Scratch	52.1417	0.9924
90%	52.9650	0.9933
80%	53.2565	0.9936
70%	52.6687	0.9928
60%	52.5372	0.9928

PSF constraints, loss components, and self-supervised pre-training, demonstrating the synergy between physics priors and learned features. Even with reduced input views, MVM-LF pretraining enables robust reconstruction, making XLFM-Former a scalable solution for neuroscience and biomedical imaging.

In the fields of medical imaging and neuroscience, labeled data is often scarce and expensive to obtain, making self-supervised learning an essential direction for future research. Meanwhile, the computer vision community has developed a range of powerful and efficient pretraining strategies and backbone architectures, which have significantly advanced neural imaging techniques. At the same time, novel methodologies emerging from neuroscience have inspired innovations in deep learning model design. We hope that our work serves as a bridge between these two communities, fostering closer interdisciplinary collaboration and encouraging further exploration of self-supervised pretraining for complex biological imaging tasks.

531 **References**

- [1] Ernst Abbe. Beiträge zur theorie des mikroskops und der mikroskopischen wahrnehmung. *Archiv für mikroskopische Anatomie*, 9(1):413–468, 1873. 2
- [2] Lu Bai, Lin Cong, Ziqi Shi, Yuchen Zhao, Yujie Zhang, Bin Lu, Jing Zhang, Zhi-Qi Xiong, Ninglong Xu, Yu Mu, et al. Volumetric voltage imaging of neuronal populations in the mouse brain by confocal light-field microscopy. *Nature Methods*, 21(11):2160–2170, 2024. 1
- [3] Amir Bar, Xin Wang, Vadim Kantorov, Colorado J Reed, Roei Herzig, Gal Chechik, Anna Rohrbach, Trevor Darrell, and Amir Globerson. Detreg: Unsupervised pretraining with region priors for object detection. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 14605–14615, 2022. 3
- [4] Edward S Boyden, Feng Zhang, Ernst Bamberg, Georg Nagel, and Karl Deisseroth. Millisecond-timescale, genetically targeted optical control of neural activity. *Nature neuroscience*, 8(9):1263–1268, 2005. 1
- [5] Yuming Chai, Kexin Qi, Yubin Wu, Daguang Li, Guodong Tan, Yuqi Guo, Jun Chu, Yu Mu, Chen Shen, and Quan Wen. All-optical interrogation of brain-wide activity in freely swimming larval zebrafish. *iScience*, 27(1):108385, 2024. 1
- [6] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *International conference on machine learning*, pages 1597–1607. PMLR, 2020. 3
- [7] Ching-Yao Chuang, Joshua Robinson, Yen-Chen Lin, Antonio Torralba, and Stefanie Jegelka. Debiased contrastive learning. *Advances in neural information processing systems*, 33:8765–8775, 2020. 3
- [8] Lin Cong, Zeguan Wang, Yuming Chai, Wei Hang, Chunfeng Shang, Wenbin Yang, Lu Bai, Jiulin Du, Kai Wang, and Quan Wen. Rapid whole brain imaging of neural activity in freely behaving larval zebrafish (*danio rerio*). *elife*, 6:e28158, 2017. 1
- [9] Jiequan Cui, Zhisheng Zhong, Shu Liu, Bei Yu, and Jiaya Jia. Parametric contrastive learning. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 715–724, 2021. 3
- [10] Diptodip Deb, Zhenfei Jiao, Ruth Sims, Alex Chen, Michael Broxton, Misha B Ahrens, Kaspar Podgorski, and Srinivas C Turaga. Fouriernets enable the design of highly non-local optical encoders for computational imaging. *Advances in Neural Information Processing Systems*, 35:25224–25236, 2022. 3
- [11] Winfried Denk, James H Strickler, and Watt W Webb. Two-photon laser scanning fluorescence microscopy. *Science*, 248(4951):73–76, 1990. 1
- [12] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations*. 7
- [13] Dumitru Erhan, Aaron Courville, Yoshua Bengio, and Pascal Vincent. Why does unsupervised pre-training help deep learning? In *Proceedings of the thirteenth international conference on artificial intelligence and statistics*, pages 201–208. JMLR Workshop and Conference Proceedings, 2010. 3
- [14] Dengpan Fu, Dongdong Chen, Jianmin Bao, Hao Yang, Lu Yuan, Lei Zhang, Houqiang Li, and Dong Chen. Unsupervised pre-training for person re-identification. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 14750–14759, 2021. 3
- [15] 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*, pages 770–778, 2016. 7
- [16] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked autoencoders are scalable vision learners. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 16000–16009, 2022. 3
- [17] Jan Huisken, Jim Swoger, Filippo Del Bene, Joachim Wittbrodt, and Ernst HK Stelzer. Optical sectioning deep inside live embryos by selective plane illumination microscopy. *Science*, 305(5686):1007–1009, 2004. 1
- [18] Prannay Khosla, Piotr Teterwak, Chen Wang, Aaron Sarna, Yonglong Tian, Phillip Isola, Aaron Maschinot, Ce Liu, and Dilip Krishnan. Supervised contrastive learning. *Advances in neural information processing systems*, 33:18661–18673, 2020. 3
- [19] Marc Levoy, Ren Ng, Andrew Adams, Matthew Footer, and Mark Horowitz. Light field microscopy. In *Acm siggraph 2006 papers*, pages 924–934. 2006. 1
- [20] Haoyu Li, Changliang Guo, Deborah Kim-Holzapfel, Weiyi Li, Yelena Altshuller, Bryce Schroeder, Wenhao Liu, Yizhi Meng, Jarrod B French, Ken-Ichi Takamaru, et al. Fast, volumetric live-cell imaging using high-resolution light-field microscopy. *Biomedical optics express*, 10(1):29–49, 2018. 1
- [21] Hao Liu and Pieter Abbeel. Behavior from the void: Unsupervised active pre-training. *Advances in Neural Information Processing Systems*, 34:18459–18473, 2021. 3
- [22] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 10012–10022, 2021. 4
- [23] Zhuang Liu, Hanzi Mao, Chao-Yuan Wu, Christoph Feichtenhofer, Trevor Darrell, and Saining Xie. A convnet for the 2020s. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 11976–11986, 2022. 7
- [24] Leon B Lucy. An iterative technique for the rectification of observed distributions. *Astronomical Journal*, Vol. 79, p. 745 (1974), 79:745, 1974. 2
- [25] Oscar Manas, Alexandre Lacoste, Xavier Giró-i Nieto, David Vazquez, and Pau Rodriguez. Seasonal contrast: Unsupervised pre-training from uncurated remote sensing data. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9414–9423, 2021. 3

- 645 [26] Josué Page Vizcaíno, Panagiotis Symvoulidis, Zeguan Wang,
646 Jonas Jelten, Paolo Favaro, Edward S Boyden, and To-
647 bias Lasser. Fast light-field 3d microscopy with out-of-
648 distribution detection and adaptation through conditional
649 normalizing flows. *Biomedical optics express*, 15(2):1219–
650 1232, 2024. 3
- 651 [27] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-
652 net: Convolutional networks for biomedical image segmen-
653 tation. In *Medical image computing and computer-assisted*
654 *intervention-MICCAI 2015: 18th international conference,*
655 *Munich, Germany, October 5-9, 2015, proceedings, part III*
656 *18*, pages 234–241. Springer, 2015. 7
- 657 [28] C. F. Shang, Y. F. Wang, M. T. Zhao, Q. X. Fan, S. Zhao,
658 Y. Qian, S. J. Xu, Y. Mu, J. Hao, and J. L. Du. Real-time
659 analysis of large-scale neuronal imaging enables closed-loop
660 investigation of neural dynamics. *Nat Neurosci*, 27(5):1014–
661 1018, 2024. 1
- 662 [29] Mingxing Tan and Quoc Le. Efficientnet: Rethinking model
663 scaling for convolutional neural networks. In *International*
664 *conference on machine learning*, pages 6105–6114. PMLR,
665 2019. 7
- 666 [30] Yonglong Tian, Chen Sun, Ben Poole, Dilip Krishnan,
667 Cordelia Schmid, and Phillip Isola. What makes for good
668 views for contrastive learning? *Advances in neural informa-*
669 *tion processing systems*, 33:6827–6839, 2020. 3
- 670 [31] Josue Page Vizcaino, Zeguan Wang, Panagiotis Symvoulidis,
671 Paolo Favaro, Burcu Guner-Ataman, Edward S Boyden,
672 and Tobias Lasser. Real-time light field 3d microscopy
673 via sparsity-driven learned deconvolution. In *2021 IEEE*
674 *International Conference on Computational Photography*
675 (*ICCP*), pages 1–11. IEEE, 2021. 3
- 676 [32] Wenhui Wang, Enze Xie, Xiang Li, Deng-Ping Fan, Kaitao
677 Song, Ding Liang, Tong Lu, Ping Luo, and Ling Shao.
678 Pyramid vision transformer: A versatile backbone for dense
679 prediction without convolutions. In *Proceedings of the*
680 *IEEE/CVF international conference on computer vision*,
681 pages 568–578, 2021. 7
- 682 [33] Xiao Wang and Guo-Jun Qi. Contrastive learning with
683 stronger augmentations. *IEEE transactions on pattern anal-*
684 *ysis and machine intelligence*, 45(5):5549–5560, 2022. 3
- 685 [34] Da Yang, Hao Sheng, Sizhe Wang, Shuai Wang, Zhang
686 Xiong, and Wei Ke. Boosting light field spatial super-
687 resolution via masked light field modeling. *IEEE Transac-*
688 *tions on Computational Imaging*, 2024. 3
- 689 [35] Feng Zhang, Li-Ping Wang, Edward S Boyden, and Karl
690 Deisseroth. Channelrhodopsin-2 and optical control of ex-
691 citable cells. *Nature methods*, 3(10):785–792, 2006. 1
- 692 [36] Zhenkun Zhang, Lu Bai, Lin Cong, Peng Yu, Tianlei Zhang,
693 Wanzhuo Shi, Funing Li, Jiulin Du, and Kai Wang. Imaging
694 volumetric dynamics at high speed in mouse and zebrafish
695 brain with confocal light field microscopy. *Nature biotech-*
696 *nology*, 39(1):74–83, 2021. 1

XLFM-Former: Rapid 3D Reconstruction for Light Field Microscopy

Supplementary Material

697

7. Dataset Statistics Details

698 Table 5 summarizes the statistics of the XLFM-Zebrafish
 699 dataset, which consists of image data from zebrafish in both
 700 free-swimming and fixed conditions. The dataset is cate-
 701 gorized into Pre-training Set, Training/Validation Set, and Test
 702 Set to facilitate different stages of model development. The
 703 dataset is carefully designed to ensure diversity in motion
 704 complexity, viewpoint variations, and temporal resolutions.
 705 By pretraining on large-scale, complex free-swimming ze-
 706 brafish data, the model gains a stronger ability to general-
 707 ize and better reconstruct simpler fixed zebrafish data, lead-
 708 ing to improved performance. This structured dataset and
 709 training methodology provide a standardized benchmark for
 710 evaluating XLFM-Former and contribute to the advance-
 711 ment of XLFM-based 3D reconstruction techniques.

712 **Pre-training Set (Free-swimming Zebrafish):** The
 713 Pre-training Set comprises three free-swimming zebrafish
 714 (m_1, m_2, m_3), totaling 20,123 images, all captured at a 10
 715 fps sampling rate. Free-swimming zebrafish exhibit highly
 716 dynamic and complex motion, leading to significant vari-
 717 ations in viewpoint and pose. This complexity presents a
 718 challenge for 3D reconstruction but also offers an opportu-
 719 nity for the model to learn richer geometric structures. To
 720 leverage this complexity, we employ unsupervised pretrain-
 721 ing on this large-scale free-swimming dataset before train-
 722 ing on the fixed zebrafish dataset. By learning from diverse,
 723 naturally occurring light field transformations, the model
 724 develops a robust understanding of light field geometry and
 725 depth relationships. This approach significantly improves
 726 performance when fine-tuned on simpler fixed zebrafish
 727 data, demonstrating the benefits of pretraining on large, di-
 728 verse datasets before supervised learning on smaller, more
 729 controlled datasets.

730 **Training/Validation Set (Fixed Zebrafish):** The Train-
 731 ing/Validation Set contains seven fixed zebrafish (f_1-f_7)
 732 with a total of 1,761 images. Most samples were collected
 733 at 10 fps, while some (e.g., f_3) were acquired at 1 fps
 734 to introduce variations in temporal resolution. Compared
 735 to free-swimming zebrafish, the fixed zebrafish dataset
 736 presents a more structured and constrained setting, making
 737 it an ideal target for supervised training once the model has
 738 been pretrained on more complex free-swimming data.

739 **Test Set (Fixed Zebrafish):** The Test Set consists of six
 740 fixed zebrafish (t_1-t_6) with a total of 1,011 images. Some
 741 samples (e.g., t_2, t_4, t_5) were acquired at 1 fps, allowing
 742 a comprehensive evaluation of XLFM-Former's reconstruc-
 743 tion performance under different sampling conditions.

Table 5. The XLFM-Zebrafish dataset statistics.

Dataset Name	Number of Images	Sampling Rate (fps)
Free-swimming Zebrafish (Pre-training Set)		
moving_fish1 (m_1)	4000	10
moving_fish2 (m_2)	7332	10
moving_fish3 (m_3)	8791	10
Fixed Zebrafish (Training/Validation Set)		
fixed_fish1 (f_1)	240	10
fixed_fish2 (f_2)	117	10
fixed_fish3 (f_3)	318	1
fixed_fish4 (f_4)	314	10
fixed_fish5 (f_5)	374	10
fixed_fish6 (f_6)	214	10
fixed_fish7 (f_7)	184	10
Fixed Zebrafish (Test Set)		
test_fixed_fish1 (t_1)	300	10
test_fixed_fish2 (t_2)	41	1
test_fixed_fish3 (t_3)	334	10
test_fixed_fish4 (t_4)	61	1
test_fixed_fish5 (t_5)	61	1
test_fixed_fish6 (t_6)	214	10

8. Loss Function Details

In our proposed XLFM-Former framework, we employ different loss functions tailored for **pretraining** and **reconstruction tasks** to ensure robust and high-quality 3D volume generation.

Pretraining Loss: For the self-supervised pretraining task, where we use Masked View Modeling-Light Field (MVM-LF), we adopt a simple ℓ_2 loss to enforce the consistency between the predicted and ground truth light field views:

$$\mathcal{L}_{\text{pretrain}} = \|\hat{I} - I\|_2^2, \quad (12)$$

where \hat{I} represents the predicted light field views, and I denotes the original (ground truth) views before masking. The choice of ℓ_2 loss is motivated by its stability in regression tasks and its ability to ensure smooth reconstructions during pretraining.

XLFM Reconstruction Loss: For the final **3D reconstruction task**, we minimize the following composite loss function:

$$\begin{aligned} \mathcal{L}_{\text{total}} = & \frac{1}{\lambda_1} \mathcal{L}_{\text{MS_SSIM}} + \frac{1}{\lambda_2} \mathcal{L}_{\text{Edge}} + \frac{1}{\lambda_3} \mathcal{L}_{\text{PSNR}} \\ & + \frac{1}{\lambda_4} \mathcal{L}_{\text{MSE}} + \frac{1}{\lambda_5} \mathcal{L}_{\text{PSF}}. \end{aligned} \quad (13)$$

Each term in the loss function contributes to a different aspect of the 3D reconstruction quality, ensuring sharpness,

accuracy, and optical consistency. Below, we provide a detailed breakdown of each component (PSF loss has been described in Section 4.2.):

Multi-Scale Structural Similarity (MS-SSIM) Loss:

Structural Similarity Index (SSIM) is widely used to measure the perceptual similarity between images. We employ a multi-scale SSIM (MS-SSIM) loss to capture both local and global structural fidelity:

$$\mathcal{L}_{\text{MS_SSIM}} = 1 - \text{MS-SSIM}(\hat{V}, V), \quad (14)$$

where \hat{V} and V represent the predicted and ground truth 3D volumes. MS-SSIM helps preserve structural details and enhances the perceptual quality of the reconstruction.

Edge-Aware Loss: To enhance edge sharpness and suppress blurring, we define the edge-preserving loss as a weighted combination of edge loss and multi-scale MSE loss:

$$\begin{aligned} \mathcal{L}_{\text{Edge_Aware}} &= \frac{1}{\lambda_{\text{edge}}} \left(\|\nabla_x \hat{V} - \nabla_x V\|_1 + \|\nabla_y \hat{V} - \nabla_y V\|_1 \right. \\ &\quad \left. + \|\nabla_z \hat{V} - \nabla_z V\|_1 \right) \\ &\quad + \frac{1}{\lambda_{\text{mse}}} \sum_{s=1}^S \|\hat{V}^{(s)} - V^{(s)}\|_2^2. \end{aligned} \quad (15)$$

Here, $\mathcal{L}_{\text{edges_loss}}$ ensures sharpness by penalizing gradient differences along spatial dimensions, while $\mathcal{L}_{\text{multi_scale_MSE}}$ enforces consistency across multiple resolutions, improving global structure and fine details.

Peak Signal-to-Noise Ratio (PSNR) Loss: PSNR is a widely used metric for measuring signal fidelity. We define a loss function that penalizes low PSNR values:

$$\mathcal{L}_{\text{PSNR}} = -\text{PSNR}(\hat{V}, V), \quad (16)$$

where a higher PSNR corresponds to higher-quality reconstructions. By minimizing this loss, we encourage the model to reduce noise and artifacts.

Mean Squared Error (MSE) Loss: The MSE loss ensures pixel-wise intensity similarity between the predicted and ground truth volumes:

$$\mathcal{L}_{\text{MSE}} = \|\hat{V} - V\|_2^2. \quad (17)$$

While MSE is commonly used in image restoration, it is prone to blurring. Thus, we use it in combination with perceptual losses (e.g., MS-SSIM and edge loss) to balance fine details and overall similarity.

8.1. Ablation Study on Loss Functions

To evaluate the impact of different loss components on the overall performance of our model, we conduct an ablation study by removing individual loss terms and measuring the

Table 6. **Loss function.** ‘Full’ denotes the complete model with all loss terms. ‘w.o. ms_ssimm’ removes the multi-scale SSIM loss. ‘w.o. Edge_Aware’ omits the edge-aware loss. ‘w.o. PSNR’ excludes the PSNR loss. ‘w.o. MSE Loss’ removes the MSE loss. The best results are highlighted in **bold**.

Method	PSNR↑	SSIM↑
w.o. $\mathcal{L}_{\text{MS_SSIM}}$	53.2787	0.9937
w.o. $\mathcal{L}_{\text{Edge_Aware}}$	52.1870	0.9922
w.o. $\mathcal{L}_{\text{PSNR}}$	53.0741	0.9935
w.o. \mathcal{L}_{MSE}	53.2521	0.9937
Full	54.0435	0.9944

performance in terms of PSNR and SSIM. The results are summarized in Table 6.

From the table, we observe that the complete model (*Full*) achieves the highest PSNR of **54.0435** and SSIM of **0.9944**, demonstrating the effectiveness of integrating all loss terms.

Effect of Multi-Scale SSIM Loss: Removing the multi-scale SSIM loss (*w.o. $\mathcal{L}_{\text{MS_SSIM}}$*) results in a performance drop, reducing PSNR to 53.2787 and SSIM to 0.9937. This highlights the importance of SSIM-based perceptual loss in preserving structural information.

Effect of Edge-Aware Loss: When the edge-aware loss is removed (*w.o. $\mathcal{L}_{\text{Edge_Aware}}$*), the PSNR decreases significantly to 52.1870, while SSIM drops to 0.9922. This indicates that the edge-aware term plays a crucial role in maintaining sharp details and edge consistency.

Effect of PSNR Loss: Excluding the PSNR-based loss (*w.o. $\mathcal{L}_{\text{PSNR}}$*) results in a minor degradation in performance, with PSNR reducing to 53.0741 and SSIM to 0.9935. This suggests that optimizing directly for PSNR provides some benefits but is not the dominant factor.

Effect of MSE Loss: When the MSE loss is removed (*w.o. \mathcal{L}_{MSE}*), the PSNR slightly drops to 53.2521, while SSIM remains at 0.9937. This indicates that MSE contributes to pixel-wise accuracy but is less critical than the other loss terms.

9. Training Configuration

For the MVM-LF task, we adopt a batch size of 8 to ensure stable training dynamics while maintaining an efficient balance between computational cost and convergence stability. To further enhance the training process and prevent the model from becoming trapped in local optima, we employ the ReduceLROnPlateau learning rate scheduler. The initial learning rate is set to 1e-4, and the scheduler dynamically adjusts the learning rate based on validation loss fluctuations, reducing it when no improvement is observed for a predefined number of epochs. This adaptive learning rate strategy helps in maintaining a steady convergence while

Table 7. Data augmentation techniques applied during training.
Each transformation is applied with a probability of 0.5.

Augmentation Type	Probability
Random Flip	0.5
Random Rotation	0.5
Random Gaussian Noise&Blur	0.5
Random Brightness &Contrast	0.5

846 preventing premature stagnation in suboptimal solutions.

847 To achieve robust feature learning and ensure general-
848 ization across diverse data distributions, we train the model
849 for 250 epochs. Given the complexity of the task and the
850 high-dimensional nature of the input data, prolonged train-
851 ing allows the model to capture intricate spatial and struc-
852 tural information effectively. All experiments are conducted
853 in a distributed computing environment equipped with four
854 NVIDIA A100-80GB SMX4 GPUs, leveraging multi-GPU
855 parallelism to accelerate training and optimize resource util-
856 ization.

857 For the XLFM reconstruction task, the training setup re-
858 mains largely consistent with the MVM-LF configuration.
859 However, a notable difference is the use of a batch size of
860 1, which aligns with the requirements of volumetric recon-
861 struction. Given that volumetric data often involves higher
862 memory footprints due to its three-dimensional repres-
863 entation, a smaller batch size ensures that computations re-
864 main feasible within available GPU memory constraints.
865 All experiments are implemented using PyTorch Lightning,
866 a high-level deep learning framework that simplifies train-
867 ing and enhances reproducibility.

868 9.1. Data Augmentation Strategy

869 To improve the model’s robustness and generalization ca-
870 pability, we apply a series of data augmentation techniques
871 during training. These augmentations are applied with a
872 probability of 0.5, as summarized in Table 7.

873 These augmentation techniques enhance the model’s
874 ability to handle variations in real-world data, reducing
875 overfitting and improving generalization performance.

876 9.2. Model Architecture

877 For all experiments, we adopt the Swin Transformer frame-
878 work in its tiny configuration. The Swin Transformer is a
879 hierarchical vision transformer that efficiently models long-
880 range dependencies while maintaining computational effi-
881 ciency. The tiny variant provides a lightweight architecture
882 suitable for our training setup, balancing performance and
883 efficiency.

9.3. Extended Training on Additional Datasets 884

To further explore the scalability and generalization capa-
885 bility of XLFM-Former, we conducted large-scale training
886 on an extended dataset. This additional training aimed to
887 improve the model’s ability to reconstruct fine-grained de-
888 tails while enhancing robustness across diverse volumetric
889 data.

The large-scale training process required substantial
890 computational resources, consuming approximately 1344
891 A100-80GB GPU hours. This extensive training allowed
892 the model to refine its feature representation and leverage
893 a broader data distribution, leading to improved reconstruc-
894 tion quality.

To demonstrate the effectiveness of this large-scale training,
895 we provide a visualized demo showcasing the en-
896 hanced view synthesis capability. The qualitative results
897 highlight the model’s ability to generate high-fidelity recon-
898 structions with improved structural consistency and percep-
899 tual quality, further validating the benefits of extended train-
900 ing.