# MASKTWINS: DUAL-FORM COMPLEMENTARY MASK-ING FOR DOMAIN-ADAPTIVE IMAGE SEGMENTATION

Anonymous authors

Paper under double-blind review

# ABSTRACT

Recent works have correlated Masked Image Modeling (MIM) with consistency regularization in unsupervised domain adaptation. However, they merely treat masking as a special form of deformation on the input images and neglect the theoretical analysis, which leads to a superficial understanding of masked reconstruction and insufficient exploitation of its potential in enhancing feature extraction and representation learning. In this paper, we reframe masked reconstruction as a sparse signal reconstruction problem and theoretically prove that the dual form of complementary masks possesses superior capabilities in extracting domain-agnostic image features. Based on this compelling insight, we propose MaskTwins, a simple yet effective learning strategy that integrates masked reconstruction directly into the main training pipeline. MaskTwins uncovers intrinsic structural patterns that persist across disparate domains by enforcing consistency between predictions of images masked in complementary ways, enabling domain generalization in an end-to-end manner. Extensive experiments verify the superiority of MaskTwins over baseline methods in natural and biological image segmentation. These results demonstrate the significant advantages of MaskTwins in extracting domain-invariant features without the need for separate pre-training, offering a new paradigm for domain-adaptive segmentation.

# 027 028 029 030

031

004

010 011

012

013

014

015

016

017

018

019

021

025

026

# 1 INTRODUCTION

032 Inspired by Masked Language Modeling (MLM) (Devlin, 2018; Brown, 2020) in natural language 033 processing, Masked Image Modeling (MIM) (Bao et al., 2022; He et al., 2022; Xie et al., 2022b) has 034 achieved remarkable success in self-supervised visual representation learning. MIM learns semantic representations by deliberately obscuring parts of the input and then reconstructing the missing 035 information based on the unmasked parts, e.g., normalized pixels (He et al., 2022; Xie et al., 2022b), HOG feature (Wei et al., 2022), discrete tokens (Bao et al., 2022; Dong et al., 2023), deep features 037 (Zhou et al., 2021; Dong et al., 2022) or frequencies (Xie et al., 2022a; Liu et al., 2023). Their success stems from the ability to learn robust, generalizable features despite incomplete or corrupted data, as masked reconstruction techniques simulate real-world visual occlusions and distortions, 040 enhancing model comprehension of visual concepts. 041

Analogously, consistency regularization in unsupervised domain adaptive segmentation learns 042 domain-invariant features by enforcing consistency between the predictions of transformed images 043 and their original counterparts. In unsupervised domain adaptation (UDA), consistency regulariza-044 tion based methods (Choi et al., 2019; Araslanov & Roth, 2021; Melas-Kyriazi & Manrai, 2021) typically utilize a variety of augmentations, like affine transformations, color jittering and cutout 046 (DeVries, 2017), expecting the learned feature to be invariant to a certain group of transformations 047 on the inputs. Focusing excessively on selecting the most appropriate parameters and perturbation 048 functions makes them depart from the simple principle of consistency. Recently, MIC (Hoyer et al., 2023) uses masked image consistency to learn context relations. However, it considers masking as merely an image deformation and neglect the theoretical analysis, which results in a cursory under-051 standing of masked reconstruction and a failure to fully harness its benefits for feature extraction and representation learning. Moreover, the learning from single masked context is limited and the 052 effectiveness of single-branch masked consistency is largely contingent upon the accuracy of the pseudo-labels generated, whose incorrectness will lead to noisy training and poor generalization.

054 In this paper, we propose a novel perspective on masked reconstruction by reframing it as a sparse 055 signal reconstruction problem and utilize it to design an effective strategy for domain-adaptive seg-056 mentation. Our theoretical analysis reveals that the dual form of complementary masks possesses su-057 perior image feature extraction capabilities. This insight is grounded in the principles of compressed 058 sensing (Donoho, 2006), suggesting that complementary masks can provide a more comprehensive sampling of the input space. Building upon this theoretical foundation, we introduce MaskTwins, a simple yet effective learning strategy for domain-adaptive segmentation. MaskTwins leverages the 060 consistency constraints of complementary masks to extract domain-invariant features. Furthermore, 061 it employs Adaptive Instance Normalization (AdaIN) (Huang & Belongie, 2017) to adjust feature 062 statistics between source and target domains, enhancing adaptability across diverse data sources. 063 This approach not only advances the theoretical understanding of masked reconstruction but also 064 provides a practical framework for improving performance on domain-adaptive vision tasks. 065

- 066 Our contributions can be summarized as follows:
  - 1. We provide a theoretical foundation for masked reconstruction by reframing it as a sparse signal reconstruction problem, offering new insights into the effectiveness of complementary masks. This perspective bridges the gap between masked image modeling and signal processing theory, potentially opening new avenues for future research.
  - 2. We propose MaskTwins, a novel learning strategy that enforces consistency between predictions of dual-form complementary masked images without introducing extra learnable parameters. Therefore, this approach is computationally efficient and can be easily integrated into existing architectures.
  - 3. We demonstrate the superiority of our approach through extensive experiments, showing significant improvements over baseline methods in both natural and biological image segmentation. Our results indicate that MaskTwins can enhance model robustness and adaptability across diverse domains, providing a more conceptual guidance for masked consistency learning in vision tasks.
- 080 081

067

068

069

071

073

075

076

077

078

079

082 083

084

# 2 RELATED WORKS

# 2.1 UNSUPERVISED DOMAIN ADAPTATION

**UDA in natural image segmentation** Unsupervised domain adaptation (UDA) addresses the critical problem of performance degradation in target domains through the effective exploitation of both 087 labeled source domain data and unlabeled target domain data. By bridging the domain gaps, UDA 088 has emerged as a versatile solution to enhance model robustness in various computational domains, demonstrating promising results on various computer vision tasks such as natural image semantic segmentation (Tsai et al., 2018; Mei et al., 2020; Jiang et al., 2022) and medical image segmenta-091 tion (Bermúdez-Chacón et al., 2018; Liu et al., 2020a; Wu et al., 2021). UDA solutions are broadly categorized into three groups: statistical moment alignment (Chen et al., 2019; Liu et al., 2020b), 092 adversarial learning (Tsai et al., 2018; Luo et al., 2021; Zheng & Yang, 2022) and self-training (Zou 093 et al., 2018; Mei et al., 2020; Zhao et al., 2023). Methods based on statistical moment alignment aim 094 to minimize the domain discrepancy employing an appropriate statistical distance function such as 095 entropy minimization (Chen et al., 2019) and Wasserstein distance (Liu et al., 2020b). Adversarial 096 training methods achieve domain invariant feature extraction with a GAN framework (Goodfellow et al., 2014). To overcome the challenges of instability in adversarial learning, Zheng & Yang (2022) 098 adaptively refine the distribution of training data by aggregating the weak models. In self-training, 099 pseudo labels (Lee et al., 2013) are created for the unlabeled target domain using confidence thresh-100 olds (Zou et al., 2018; 2019; Mei et al., 2020), pseudo-label prototypes (Zhang et al., 2019a; 2021; 101 Jiang et al., 2022) or uncertainty (Zheng & Yang, 2021). Recently, Hoyer et al. (2023) and Yang 102 et al. (2024) explore context relations while Zhao et al. (2023) learn pixel-wise representations to boost the quality of pseudo-labels. 103

104

UDA in biological image segmentation For the segmentation of biological images, domain adaptation is receiving increasing attention due to the lack of manually annotated data. Specially, the 3D volumes of microscopy image datasets allow the additional consideration of the consistency of adjacent sections. For example, Huang et al. (2022b) take the inter-slice information into account

and Sun et al. (2023) construct an intricately-designed network that captures long-range sectional variations within structures and effectively discriminates by adaptively aggregating diverse components. Yin et al. (2023) performs domain alignment in the feature space and incorporates the prototype representation into feature alignment. Different from these UDA methods, our proposed method integrates the context relationships by enforcing complementary masked consistency without introducing extra learnable parameters. The dual-form masked image consistency enables the learning of complementary clues, which further boosts the extraction of doamain-invariant features and increases the robustness of networks across different segmentation tasks.

116 117

118

# 2.2 MASKED IMAGE MODELING

119 Masked Image modeling (MIM) (Bao et al., 2022; Wei et al., 2022; He et al., 2022) methods are 120 showing great promise in visual self-supervised representation learning for their ability to learn robust and generalizable features from incomplete or corrupted input data, enhancing the models' 121 comprehension of visual concepts. Many target signals have been conceived for the masked recon-122 struction, encompassing raw pixels (He et al., 2022; Xie et al., 2022b), HOG features (Wei et al., 123 2022), discrete visual tokens (Bao et al., 2022; Dong et al., 2023), frequencies (Xie et al., 2022a; 124 Liu et al., 2023) and deep features (Zhou et al., 2021; Dong et al., 2022). Recently, Wang et al. 125 (2023) further explore the reconstruction process at multiple scales while Kong & Zhang (2023) 126 interprete MIM in a unified framework. However, these works mainly treat masked reconstruction 127 as a pre-training strategy but neglect its potential for downstream tasks related to domain general-128 ization. Hoyer et al. (2023) preliminarily explore the masked target image in the UDA setting and 129 conclude that masked image consistency substantially boosts UDA performance through additional 130 context clues. Shin et al. (2024) superficially perform complementary masking for RGB-Thermal 131 segmentation. Yet, a thorough theoretical foundation for the effectiveness of masked images in domain adaptation remains to be established. In this work, we introduce a novel reconceptualiza-132 tion of the masked reconstruction as a sparse signal reconstruction problem and refine the theory of 133 complementary masks. By surpassing the constraints of domain-specific customization, MaskTwins 134 employs a strategic complementary masking technique on the input data, ensuring a more holistic 135 and nuanced understanding of the intrinsical data patterns. 136

137 138

139 140

141

# 3 Method

# 3.1 OVERVIEW

142 The MaskTwins framework for unsupervised domain adaptation (UDA) in semantic segmentation 143 is detailed in Figure 1. The objective is to train a neural network  $f_{\theta}$  that effectively generalizes to the target domain, given a labeled source domain dataset  $X^S = \{(x_i^S, y_i^S)\}_{i=1}^{N_S} \subseteq \mathcal{D}^S$  and an unlabeled target domain dataset  $X^T = \{x_j^T\}_{j=1}^{N_T} \subseteq \mathcal{D}^T$ . The framework operates by generating two 144 145 146 complementary masked versions of each target image  $x_i^T$ , denoted as  $D \odot x_i^T$  and  $(1 - D) \odot x_i^T$ , 147 where D is a binary mask. A teacher model  $f_{\phi}$ , updated via the Exponential Moving Average (EMÅ) 148 of the student parameters, generates pseudo-labels for the target domain. The student's predictions, 149 together with the pseudo-labels from the teacher model, are used to compute the target-domain 150 losses, while a supervised loss is computed using the labeled source data. This iterative process 151 adapts the model to the target domain, leveraging both the supervised source information and the 152 unsupervised adaptation to the target domain.

153

154 Motivation Consistency regularization (Choi et al., 2019; Araslanov & Roth, 2021; Melas-Kyriazi 155 & Manrai, 2021) is a common technique in UDA. It typically leverages a rich set of augmentations, 156 like affine transformations, cutout (DeVries, 2017), and color jittering in images. Nonetheless, their 157 success heavily depends on the accuracy of the pseudo-labels generated, whose incorrectness will 158 lead to noisy training and poor generalization. Inspired by MIC (Hoyer et al., 2023), we expect the performance of masked consistency in UDA. We further take insights from the paradigm of 159 masked reconstruction (Bao et al., 2022; He et al., 2022; Xie et al., 2022b) and propose the theory of 160 complementary masks to support the application of masked consistency for domain-adaptive image 161 segmentation.



Figure 1: The overall framework of MaskTwins. Given the labeled source data, we calculate the segmentation prediction  $P^S$  with the network  $f_{\theta}$ , supervised by basic segmentation loss  $\mathcal{L}_{sup}^S$ . For the target domain, we obtain the predictions of complementary masked target images, constrained by the pseudo-labels  $P_T$  that are generated based on the unmasked image by an exponential moving average (EMA) teacher  $f_{\phi}$ . "//" on " $\rightarrow$ " means stop gradient. Furthermore, MaskTwins proposes the complementary masked loss between dual-form complementary masked images for deep consistency learning.

3.2 THEORETICAL ANALYSIS OF COMPLEMENTARY MASKING

To provide a formal foundation for the complementary masking strategy in MaskTwins, we present a theoretical analysis addressing the properties of masked training in visual tasks. This analysis focuses on information preservation, generalization bounds, and feature consistency. Detailed proofs of all results are provided in Appendix E.

**Definition 1** (Complementary Mask). Let  $D \in \{0, 1\}^{H \times W}$  be a binary matrix, where each element  $D_{ij} \sim Bernoulli(0.5)$ . The complementary mask pair is defined as (D, 1 - D), where 1 is the all-ones matrix of size  $H \times W$ .

**Definition 2** (Random Mask). Let  $R \in \{0,1\}^{H \times W}$  be a binary matrix where each element  $R_{ij} \sim$ Bernoulli(0.5) independently. The random mask pair is defined as  $(R_1, R_2)$ , where  $R_1$  and  $R_2$  are independent random masks.

Assumption 1 (Visual Data Model). The input image  $X \in \mathbb{R}^{H \times W \times C}$  is generated by the model X = S + E + N, where S represents a sparse signal component, E represents environmental factors, and  $N \sim \mathcal{N}(0, \sigma^2 I)$  is additive Gaussian noise.

Assumption 2 (Feature Extraction Framework). We consider a feature extraction framework with
 the objective function:

214

192 193

194

 $\mathcal{L}(f) = \mathbb{E}_X[\ell(f(X_1), f(X_2))],\tag{1}$ 

215 where  $f : \mathbb{R}^{H \times W \times C} \to \mathbb{R}^k$  is the feature extraction function, and  $\ell : \mathbb{R}^k \times \mathbb{R}^k \to \mathbb{R}$  is the loss function.

Theorem 1 (Information Preservation). For any input image X, define the information preservation metric  $IP(X_1, X_2) = \frac{\langle f(X_1), f(X_2) \rangle}{\|f(X)\|^2}$ . Then:

$$\mathbb{E}[IP(D \odot X, (1-D) \odot X)] \ge \mathbb{E}[IP(R_1 \odot X, R_2 \odot X)]$$
<sup>(2)</sup>

$$Var(IP(D \odot X, (1-D) \odot X)) \le Var(IP(R_1 \odot X, R_2 \odot X)),$$
(3)

where  $\odot$  denotes element-wise multiplication.

**Theorem 2** (Generalization Bound). Assume  $\ell$  is L-Lipschitz and f is  $\beta$ -smooth. For any  $\delta \in (0, 1)$ , with probability at least  $1 - \delta$ :

$$|\mathcal{L}(f) - \hat{\mathcal{L}}_n(f)| \le C_1 L\beta B\left(\sqrt{\frac{1}{n}} + \sqrt{\frac{\log(1/\delta)}{n}}\right) \quad (Complementary) \tag{4}$$

$$|\mathcal{L}(f) - \hat{\mathcal{L}}_n(f)| \le C_2 L\beta B\left(\sqrt{\frac{1}{n}} + \sqrt{\frac{\log(1/\delta)}{n}} + \sqrt{\frac{HWC}{n}}\right) \quad (Random), \tag{5}$$

where  $B = \sup_{X \in \mathcal{X}} ||X||_F$ , and  $C_1$ ,  $C_2$  are constants.

**Theorem 3** (Feature Consistency). Define the feature consistency error as  $FCE(X_1, X_2) = ||f(X_1) - f(X_2)||_2$ . Then for any  $\delta \in (0, 1)$ , with probability at least  $1 - \delta$ :

$$FCE(D \odot X, (1-D) \odot X) \le C_1 \sigma \sqrt{k \log(HWC/\delta)} \quad (Complementary)$$
(6)

$$FCE(R_1 \odot X, R_2 \odot X) \le C_2 \left( \sigma \sqrt{k \log(HWC/\delta)} + \|E\|_F \sqrt{\frac{k \log(HWC/\delta)}{HWC}} \right) \quad (Random)$$
(7)

where  $C_1$ ,  $C_2$  are constants.

**Remark 1.** The theoretical results demonstrate the advantages of complementary masking. Specifically, complementary masks offer better information preservation, tighter generalization bounds, and improved feature consistency, compared to random masking. These properties are critical for extracting domain-invariant features, which are essential in cross-domain tasks such as domain adaptation.

#### 3.3 MASKTWINS: COMPLEMENTARY MASKED LEARNING

Building upon the theoretical framework, we now describe the core complementary masked learning approach in MaskTwins. This strategy employs patch-wise binary masks to generate dual complementary views of the target images. Specifically, for each target image  $x_j^T$ , a binary mask D is sampled from a Bernoulli distribution:

$$D_{mb+1:(m+1)b} \sim \text{Bernoulli}(1-r), \tag{8}$$

$$nb+1:(n+1)b$$

where r is the mask ratio, b is the patch size, and m and n are patch indices. The dual-form complementary masked images are then obtained by element-wise multiplication:

$$X_{cm}^{T} = \{X_{D}^{T}, X_{1-D}^{T}\} = \{D \odot X^{T}, (1-D) \odot X^{T}\}.$$
(9)

These complementary views encourage the model to extract robust, domain-invariant features by enforcing consistency learning upon masked images. To effectively learn from dual-form complementary contexts, we introduce two kinds of consistency losses. First, we constrain the consistent prediction of complementary masked images, which enables the network to integrate the dual-form clues. The complementary masked loss is accordingly defined as:

$$\mathcal{L}_{cm}^{T} = \mathbb{E}[\mathcal{L}_{ce}(p_{j,D}^{T}, p_{j,1-D}^{T})], \qquad (10)$$

where  $p_{j,D}^T$  and  $p_{j,1-D}^T$  are the predictions for the complementary masked images. Intended to encourage successful masked reconstruction for both masked views, we also define a masked consistency learning loss:

$$\mathcal{L}_{cl}^{T} = \mathbb{E}[\lambda \times \mathcal{L}_{ce}(p_{j,D}^{T}, \hat{y}_{j}^{T}) + (1 - \lambda) \times \mathcal{L}_{ce}(p_{j,1-D}^{T}, \hat{y}_{j}^{T})],$$
(11)

where  $\hat{y}_j^T$  are the pseudo-labels,  $\lambda$  defaults to 0.5 to ensure balanced learning from the complementary masks. Since there is no ground truth available for the target domain, a teacher model  $f_{\phi}$ predicts the pseudo-label for the unmasked target image:

$$\hat{\mu}_j^T = [c = \operatorname{argmax} f_\phi(x_j^T)], \tag{12}$$

where c is one category and the pseudo-label is converted into a one-hot categorical form via the Iverson bracket [·].

The parameters of the teacher network  $f_{\phi}$  are updated using an Exponential Moving Average (EMA) of the parameters of the student network  $f_{\theta}$  (Tarvainen & Valpola, 2017):

$$\phi_{t+1} \leftarrow \alpha \phi_t + (1 - \alpha)\theta_t, \tag{13}$$

where t denotes a training step and  $\alpha$  is the EMA decay rate. The teacher model averages the weights of previous student models over time, leading to temporally stable and reliable predictions on the target domain. Moreover, with access to the complete information from the original target images, it can issue guidance for the adaptation process and provide high-quality pseudo-labels which are then used in conjunction with our complementary masking approach to enhance the masked reconstruction.

This complementary masking strategy ensures that the model learns from diverse, yet consistent, views of the target domain, promoting robust generalization to the target domain. The next section details the overall model architecture and training process, which integrates these complementary masking principles.

290 291

303

311

316

317

3.4 MODEL ARCHITECTURE AND TRANING STRATEGY

The MaskTwins architecture consists of a shared encoder and segmentation head for both the source and target domains. To mitigate domain shift, we employ an Adaptive Instance Normalization (AdaIN) (Huang & Belongie, 2017) module in the shallow layers of the network, which aligns feature distributions between the two domains.

During training, we apply the complementary masks to the target domain images and enforce consistency between the predictions of these masked versions. This encourages the model to learn invariant representations that generalize well to the target domain. Our training strategy integrates supervised learning on the source domain with self-training and consistency regularization on the target domain.

The supervised loss on the source domain is defined as:

$$\mathcal{L}_{sup}^{S} = \mathbb{E}[\mathcal{L}_{ce}(p_{i}^{S}, y_{i}^{S})] = \mathbb{E}[-y_{i}^{S}\log(p_{i}^{S})],$$
(14)

where  $p_i^S = f_\theta(x_i^S)$  is the source prediction of the network  $f_\theta$ .

By integrating these components - complementary masking, consistency regularization, and self-training with a teacher model - MaskTwins effectively leverages the complementary information from masked inputs, promoting robust feature learning and improved generalization to the target domain.

The overall loss function that encapsulates our training strategy is formulated as:

$$\mathcal{L}_{total} = \mathcal{L}_{sup}^{S} + \mathcal{L}_{cl}^{T} + \lambda_{cm} \mathcal{L}_{cm}^{T},$$
(15)

where  $\mathcal{L}_{sup}^{S}$  is the supervised loss on the source domain,  $\mathcal{L}_{cl}^{T}$  is the masked consistency learning loss on the target domain,  $\mathcal{L}_{cm}^{T}$  is the complementary masked loss, and  $\lambda_{cm}$  is the weight for the complementary masked loss. We summarize the pipeline of MaskTwins in Algorithm 1 in Appendix B.

4 EXPERIMENTS

# 318 4.1 IMPLEMENTATION DETAILS 319

Datasets To demonstrate the versatility of MaskTwins, we conduct experiments spanning six distinct datasets: SYNTHIA (Ros et al., 2016) and Cityscapes (Cordts et al., 2016) are natural datasets, VNC III (Gerhard et al., 2013), Lucchi (Lucchi et al., 2013), MitoEM (Wei et al., 2020) and WASP-SYN (Li et al., 2024) are biological datasets. The details of the datasets and the task-specific implementation on these datasets can be found in Appendix C.

Table 1: Comparison results with previous UDA methods on SYNTHIA→Cityscapes. "SW" stands
for *sidewalk*, "TL" for *traffic light*, "TS" for *traffic sign*, "Veg." for *vegetation*, "PR" for *person*.
We present pre-class IoU and mean IoU (mIoU), averaged across 13 categories. The competitors
include DAFormer (Hoyer et al., 2022a), CAMix (Zhou et al., 2022b), HRDA (Hoyer et al., 2022b),
MIC (Hoyer et al., 2023), etc. More details are shown in Appendix A.

330	Method	Road	SW	Build	TL	TS	Veg.	Sky	PR	Rider	Car	Bus	Motor	Bike	mIoU
331	SIBAN	82.5	24.0	79.4	16.5	12.7	79.2	82.8	58.3	18.0	79.3	25.3	17.6	25.9	46.3
332	DADA	89.2	44.8	81.4	8.6	11.1	81.8	84.0	54.7	19.3	79.7	40.7	14.0	38.8	49.8
333	BDL	86.0	46.7	80.3	14.1	11.6	79.2	81.3	54.1	27.9	73.7	42.2	25.7	45.3	51.4
334	APODA	86.4	41.3	79.3	22.6	17.3	80.3	81.6	56.9	21.0	84.1	49.1	24.6	45.7	53.1
335	SIM	83.0	44.0	80.3	17.1	15.8	80.5	81.8	59.9	33.1	70.2	37.3	28.5	45.8	52.1
336	FDA	79.3	35.0	73.2	19.9	24.0	61.7	82.6	61.4	31.1	83.9	40.8	38.4	51.1	52.5
227	LSE	82.9	43.1	78.1	9.1	14.4	77.0	83.5	58.1	25.9	71.9	38.0	29.4	31.2	49.4
337	CCM	79.6	36.4	80.6	22.4	14.9	81.8	77.4	56.8	25.9	80.7	45.3	29.9	52.0	52.9
338	LDR	85.1	44.5	81.0	16.4	15.2	80.1	84.8	59.4	31.9	73.2	41.0	32.6	44.7	53.1
339	CD-SAM	82.5	42.2	81.3	18.3	15.9	80.6	83.5	61.4	33.2	72.9	39.3	26.6	43.9	52.4
340	CLAN	82.7	37.2	81.5	17.1	13.1	81.2	83.3	55.5	22.1	76.6	30.1	23.5	30.7	48.8
341	ASA	<u>91.2</u>	48.5	80.4	5.5	5.2	79.5	83.6	56.4	21.9	80.3	36.2	20.0	32.9	49.3
342	DAST	87.1	44.5	82.3	13.9	13.1	81.6	86.0	60.3	25.1	83.1	40.1	24.4	40.5	52.5
042	UncerDA	79.4	34.6	83.5	32.1	26.9	78.8	79.6	66.6	30.3	86.1	36.6	19.5	56.9	54.6
343	RPLR	81.5	36.7	78.6	20.7	23.6	79.1	83.4	57.6	30.4	78.5	38.3	24.7	48.4	52.4
344	UACR	85.5	42.5	83.0	20.9	25.5	82.5	88.0	63.2	31.8	86.5	41.2	25.9	50.7	55.9
345	DACS	80.6	25.1	81.9	22.7	24.0	83.7	90.8	67.6	38.3	82.9	38.9	28.5	47.6	54.8
346	ProDA	87.8	45.7	84.6	54.6	37.0	88.1	84.4	74.2	24.3	88.2	51.1	40.5	45.6	62.0
347	DAFormer	84.5	40.7	88.4	55.0	54.6	86.0	89.8	73.2	48.2	87.2	53.2	53.9	61.7	67.4
2/19	CAMix	87.4	47.5	88.8	55.2	55.4	87.0	91.7	72.0	49.3	86.9	57.0	57.5	63.6	69.2
0.40	HRDA	85.2	47.7	88.8	65.7	60.9	85.3	92.9	/9.4	52.8	89.0	$\frac{64.7}{1.0}$	63.9	64.9	/2.4
349	MIC	86.6	<u>50.5</u>	89.3	66.7	63.4	8/.1	94.6	81.0	<u>58.9</u>	90.1	61.9	6/.1	<u>64.3</u>	/4.0
350	Ours	96.0	/0.1	89.5	00.8	62.1	89.1	<u>94.3</u>	81.5	59.7	90.5	00.0	07.7	03.0	/6./

351 352

353 **MaskTwins parameters** MaskTwins uses the square mask for 2D domain adaptation and the cube 354 mask for 3D respectively. The complementary masks have equal loss weight and the same mask ra-355 tio r = 0.5. The mask patch size is fixed for each task, approximately 1/16 of the input size. For 356 SYNTHIA  $\rightarrow$  Cityscapes, we use a patch size b = 64, a loss weight  $\lambda_{cm} = 0.01$ , and common color 357 augmentation (brightness, contrast, saturation, hue, and blur) following the parameters of Hoyer 358 et al. (2022a), Hoyer et al. (2022b) and Tranheden et al. (2021). For mitochondria semantic segmentation, we use a patch size b = 32, a loss weight  $\lambda_{cm} = 0.01$ , a pseudo-label threshold  $\delta = 0.7$ , 359 and random augmentation including flip, transpose, rotate, resize and elastic transformation. For 360 synapse detection, the point annotations (3D coordinates) are transformed into voxel cubes with a 361 size of  $3 \times 3 \times 3$  to be used as the training target. We use a patch size b = 6, a loss weight  $\lambda_{cm} = 0.1$ 362 Empirically, we set the threshold  $\delta_{pre} = 0.75$  for the pre-synapse,  $\delta_{post} = 0.65$  for the post-synapse 363 by default. The experiments are conducted on  $8 \times RTX$  3090 GPU.

364 365 366

# 4.2 NATURAL IMAGE SEMANTIC SEGMENTATION

367 First, we compare MaskTwins with previous UDA methods on SYNTHIA→Cityscapes in Table 1. 368 It can be seen that MaskTwins outperforms the previously state-of-the-art method by a significant 369 margin of +2.7 mIoU and remains competitive in segmenting almost all classes, which verifies the 370 effectiveness of the dual form of complementary masks on target images. Classes that most profit 371 from our method are sidewalk, road, vegetation, bus, and rider. Particularly, sidewalk owns the 372 lowest UDA performance over 13 categories, meaning that it is the most difficult to adapt for previ-373 ous methods. Here, contextual relationships seem to be crucial for achieving successful adaptation. 374 However, we increase the IoU of the *sidewalk* by +19.6 from 50.5 to 70.1 IoU. Additionally, our 375 performance improvement on road is +4.8 from 91.2 to 96.0 IoU, probably because of its strong correlation with *sidewalk*. For some classes, our method increases the performance by a smaller 376 margin or causes a minor reduction, probably because the small objectives lead MaskTwins to mis-377 understand the complementary masked regions. In Figure 2, we visualize the segmentation results



Figure 2: Qualitative segmentation results on SYNTHIA→Cityscapes. MaskTwins improves the segmentation of classes such as *sidewalk*, *road*, *bus* and *rider*.



Figure 3: Qualitative comparison of MaskTwins with previous methods on VNC III→Lucchi Subset2 (row 1) and MitoEM-H→MitoEM-R (row 2). The pixels in red and green denote the falsenegative and false-positive segmentation results respectively.

and the comparison with previous strong methods HRDA (Hoyer et al., 2022b), MIC (Hoyer et al., 2023) and the ground truth. While previous methods are confused by illumination as well as crossings and fail to distinguish *sidewalk* from *road*, MaskTwins enables a more robust recognition of these categories. We can conclude that the complementary masking significantly enhances semantic segmentation, particularly for large or complex objects, where it effectively preserves structure and enables accurate segmentation despite obstacles.

# 4.3 MITOCHONDRIA SEMANTIC SEGMENTATION

We conduct quantitative comparison results of our approach with multiple UDA baselines on the Lucchi and MitoEM datasets to demonstrate the superiority of our approach. As listed in Table 2, MaskTwins achieves the new state-of-the-art results in all cases, which corroborates the effective-ness of the proposed complementary masking strategy. Specifically, MaskTwins enhances the IoU of VNC III-Lucchi(Subset1) and Lucchi(Subset2) to 75.0% and 78.6%, outperforming the state-of-the-art methods by 3.2% and 3.2%. On the MitoEM dataset with a larger structure discrepancy, our method consistently has remarkable improvements by +2.1% IoU and +1.3% IoU respectively. It is noticeable that the mitochondria in MitoEM-H exhibit denser and more intricate distributions compared to those in MitoEM-R, rendering the domain adaptation from MitoEM-R to MitoEM-H more challenging than the reverse. Despite this, MaskTwins surpasses CAFA (Yin et al., 2023) by a significant margin on the benchmark of MitoEM-R→MitoEM-H. It demonstrates that the proposed strategy can strengthen the generalization capacity of the learned model and adapt it to the chal-lenging and diverse target domain. In Figure 3, we further qualitatively compare MaskTwins with other competitive methods including DAMT-Net (Peng et al., 2020), DA-VSN (Guan et al., 2021), DA-ISC (Huang et al., 2022b), and CAFA (Yin et al., 2023). The results highlighted by yellow boxes reveal that MaskTwins shows better adaptability while other methods fail to handle hard cases with large domain gap. By leveraging the complementary masked context, our method manages to separate mitochondria correctly from the background and delivers more fine-grained results on the target domain. This indicates that MaskTwins is adept at extracting robust features of segmented objectives, thereby achieving effective adaptation from the source domain to the target domain.

132	Table 2: Quantitative comparisons on the Lucchi and MitoEM datasets (metrics in %). "Oracle"
133	denotes the model is trained on target with groundtruth labels, while "NoAdapt" represents the
134	model pretrained on source is directly applied in target for inference without any adaptation strategy.
135	The results of Oracle, NoAdapt, UALR, DAMT-Net, DA-VSN and DA-ISC are adopted from Huang
136	et al. (2022b)

Methods	VNC I	$II \rightarrow L\iota$	icchi (Su	bset1)	VNC III→Lucchi (Subset2)					
	mAP	F1	MCC	IoU	mAP	F1	MCC	IoU		
Oracle	-	92.7	86.5	86.5	-	93.9	-	88.6		
NoAdapt	-	57.3	40.3	40.3	-	61.3	-	44.3		
Advent (Vu et al., 2019a)	78.9	74.8	73.3	59.7	90.5	82.8	81.8	70.7		
UALR (Wu et al., 2021)	80.2	72.5	71.2	57.0	87.2	78.8	77.7	65.2		
DAMT-Net (Peng et al., 2020)	-	74.7	60.0	60.0	-	81.3	-	68.7		
DA-VSN (Guan et al., 2021)	82.8	75.2	73.9	60.3	91.3	83.1	82.2	71.1		
DA-ISC (Huang et al., 2022b)	89.5	81.3	80.5	68.7	92.4	85.2	84.5	74.3		
CAFA(Yin et al., 2023)	91.1	83.4	82.8	71.8	94.8	85.8	85.4	75.4		
MaskTwins(Ours)	92.4	85.6	85.1	75.0	95.2	87.9	87.4	78.6		
Methods	Mito	EM-R-	→MitoEl	M-H	Mito	EM-H-	→MitoEl	M-R		
Methods	Mito   mAP	EM-R- F1	$\rightarrow$ MitoEN MCC	M-H IoU	Mito   mAP	EM-H- F1	→MitoEl MCC	M-R IoU		
Methods Oracle	Mitc   mAP   97.0	EM-R- F1 91.6	$\rightarrow MitoEN \\ MCC \\ 91.2$	M-H IoU 84.5	Mito   mAP   98.2	EM-H- F1 93.2	→MitoEl MCC 92.9	M-R IoU 87.3		
Methods Oracle NoAdapt	Mito   mAP   97.0   74.6	DEM-R- F1 91.6 56.8	→MitoEM MCC 91.2 59.2	M-H IoU 84.5 39.6	Mito mAP 98.2 88.5	EM-H- F1 93.2 76.5	→MitoEl MCC 92.9 76.8	M-R IoU 87.3 61.9		
Methods Oracle NoAdapt Advent (Vu et al., 2019a)	Mito mAP 97.0 74.6 89.7	EM-R- F1 91.6 56.8 82.0	→MitoEN MCC 91.2 59.2 81.3	M-H IoU 84.5 39.6 69.6	Mito   mAP   98.2   88.5   93.5	EM-H- F1 93.2 76.5 85.4	→MitoEl MCC 92.9 76.8 84.8	M-R IoU 87.3 61.9 74.6		
Methods Oracle NoAdapt Advent (Vu et al., 2019a) UALR (Wu et al., 2021)	Mito   mAP   97.0 74.6   89.7 90.7	EM-R- F1 91.6 56.8 82.0 83.8	→MitoEM MCC 91.2 59.2 81.3 83.2	M-H IoU 84.5 39.6 69.6 72.2	Mito mAP 98.2 88.5 93.5 92.6	EM-H- F1 93.2 76.5 85.4 86.3	→MitoEl MCC 92.9 76.8 84.8 85.5	M-R IoU 87.3 61.9 74.6 75.9		
Methods Oracle NoAdapt Advent (Vu et al., 2019a) UALR (Wu et al., 2021) DAMT-Net (Peng et al., 2020)	Mito   mAP   97.0   74.6   89.7   90.7   92.1	EM-R- F1 91.6 56.8 82.0 83.8 84.4	→MitoEM MCC 91.2 59.2 81.3 83.2 83.7	M-H IoU 84.5 39.6 69.6 72.2 73.0	Mito Mito 98.2 88.5 93.5 92.6 94.8	EM-H- F1 93.2 76.5 85.4 86.3 86.0	→MitoEl MCC 92.9 76.8 84.8 85.5 85.7	M-R IoU 87.3 61.9 74.6 75.9 75.4		
Methods Oracle NoAdapt Advent (Vu et al., 2019a) UALR (Wu et al., 2021) DAMT-Net (Peng et al., 2020) DA-VSN (Guan et al., 2021)	Mite mAP 97.0 74.6 89.7 90.7 92.1 91.6	EM-R- F1 91.6 56.8 82.0 83.8 84.4 83.3	→MitoEM MCC 91.2 59.2 81.3 83.2 83.7 82.6	M-H IoU 84.5 39.6 69.6 72.2 73.0 71.4	Mito mAP 98.2 88.5 93.5 92.6 94.8 94.5	EM-H- F1 93.2 76.5 85.4 86.3 86.0 86.7	→MitoEl MCC 92.9 76.8 84.8 85.5 85.7 86.3	M-R IoU 87.3 61.9 74.6 75.9 75.4 76.5		
Methods Oracle NoAdapt Advent (Vu et al., 2019a) UALR (Wu et al., 2021) DAMT-Net (Peng et al., 2020) DA-VSN (Guan et al., 2021) DA-ISC (Huang et al., 2022b)	Mite mAP 97.0 74.6 89.7 90.7 92.1 91.6 92.6	EM-R- F1 91.6 56.8 82.0 83.8 84.4 83.3 85.6	→MitoEM MCC 91.2 59.2 81.3 83.2 83.7 82.6 84.9	M-H IoU 84.5 39.6 69.6 72.2 73.0 71.4 74.8	Mito mAP 98.2 88.5 93.5 92.6 94.8 94.5 96.8	EM-H- F1 93.2 76.5 85.4 86.3 86.0 86.7 88.5	<ul> <li>→MitoEl</li> <li>MCC</li> <li>92.9</li> <li>76.8</li> <li>84.8</li> <li>85.5</li> <li>85.7</li> <li>86.3</li> <li>88.3</li> </ul>	M-R IoU 87.3 61.9 74.6 75.9 75.4 76.5 79.4		
Methods Oracle NoAdapt Advent (Vu et al., 2019a) UALR (Wu et al., 2021) DAMT-Net (Peng et al., 2020) DA-VSN (Guan et al., 2021) DA-ISC (Huang et al., 2022b) CAFA (Yin et al., 2023)	Mite mAP 97.0 74.6 89.7 90.7 92.1 91.6 92.6 92.8	EM-R- F1 91.6 56.8 82.0 83.8 84.4 83.3 85.6 86.6	→MitoEN MCC 91.2 59.2 81.3 83.2 83.7 82.6 84.9 86.0	M-H IoU 84.5 39.6 69.6 72.2 73.0 71.4 74.8 76.3	Mito mAP 98.2 88.5 93.5 92.6 94.8 94.5 96.8 96.8	EM-H- F1 93.2 76.5 85.4 86.3 86.0 86.7 88.5 89.2	<ul> <li>→MitoEl</li> <li>MCC</li> <li>92.9</li> <li>76.8</li> <li>84.8</li> <li>85.5</li> <li>85.7</li> <li>86.3</li> <li>88.3</li> <li>88.9</li> </ul>	M-R IoU 87.3 61.9 74.6 75.9 75.4 76.5 79.4 80.6		

# 4.4 SYNAPSE DETECTION

We also evaluate the effectiveness of our proposed method on 3D synapse detection. This task aims to pinpoint the positions of pre-synaptic and post-synaptic sites in the 3D space, as well as to deter-mine the connectivity between them, specifically identifying the IDs of the pre-synapses to which the post-synapses are linked. For a more vivid depiction of the detection outcomes, we visualize the 3D results of pre- and post-synapse detection in Appendix D.3. Following Chen et al. (2024), we convert the task of synapse detection into a segmentation task. Since there is few prior works on this new challenge, we re-implement SSNS-Net (Huang et al., 2022a), AdaSyn (Chen et al., 2024) and MIC (Hoyer et al., 2023) strictly following their experimental implementations and make a fair comparison. Table 3 shows that MaskTwins achieves the highest F1-score with an outstand-ing gain of 2.02% totally, 3.13% on post-synapse. Due to its high density and the one-to-many synapse connectivity problem, post-synapses are more difficult to identify. Other methods perform poorly on post-synapse detection. However, MaskTwins can learn robust features and capture more post-synapses correctly with the help of complementary masks. 

Table 3: Comparison result on the WASPSYN Challenge. The F1-score is the average of  $F1_{pre}$  and  $F1_{post}$ . Table 4: Effect of the mask patch size on MitoEM-H $\rightarrow$ MitoEM-R.

80 81	Method	F1 <sub>pre</sub> F1 <sub>post</sub> F1-score		F1-score	Patch size(b)	IoU(%)	
2	SSNS-Net	0.7201	0.3072	0.5137	16	81.13	
3	AdaSyn	0.7846	0.3136	0.5491	32	81.88	
4	MIC	0.7823	0.3599	0.5711	64	81.08	
35	MaskTwins(Ours)	0.7914	0.3912	0.5913	128	81.44	

r	٦	6	1
		ŝ	3
	-	-	



Figure 4: Ablation study on mask type and mask ratio with special attention paid on the metrics of F1, MCC and IoU. Bars in blue and gray represents using complementary masks and random masks.

4.5 ABLATION STUDY

501 **Patch size** Table 4 shows the effect of the mask patch size b on MitoEM-H $\rightarrow$ MitoEM-R with a 502 input size of 512. By gradually increasing the mask patch size, we observe that the best performance 503 is achieved when b = 32, i.e. 1/16 of the input size. Patches that are either larger or smaller exhibit 504 varying degrees of performance reduction. This is likely because patches that are too large may ex-505 cessively cover the foreground while those that are too small tend to apply an overly dense masking, potentially hindering the complementary learning of contextual information. On the contrary, by 506 concentrating on context-rich areas using appropriate mask patch size, the model can better utilize 507 the spatial relations within the image, leading to improved performance in unsupervised domain 508 adaptation. Therefore, we use a mask patch size of 1/16 in all experiments. 509

510 Mask type and mask ratio We evaluate the effectiveness of complementary masks on MitoEM-511  $H \rightarrow MitoEM-R$ , compared with random masks. To do this, we systematically alter the mask ratios 512 and specifically explore the combinations of [r, 1-r] with a mask ratio  $r \in \{0.1, 0.2, 0.3, 0.4, 0.5\}$ . 513 For instance, a mask ratio of 0.3 implies a corresponding mask ratio of 0.7 in the dual-form com-514 plementary masks, and we maintain the same experimental settings in the control group of random 515 masks for a fair comparison. Upon observation, while F1, MCC, and IoU consistently show similar 516 trends, mAP stands out as an outlier. Consequently, we place greater emphasis on the metrics of F1, 517 MCC, and IoU for a more reliable assessment. As shown in Figure 4, there is a noticeable decline in 518 performance with the increase of the mask ratio in the experiments utilizing complementary masks. This decrease is attributed to the asymmetrical dual contexts, which may disrupt the bidirectional 519 training process. The best performance on the target domain is attained when employing comple-520 mentary masks with a mask ratio of 0.5. In contrast, we find that the performance of random masks 521 exhibits a fluctuation as the mask ratio changes, according with the characteristic of randomness. 522 Notably, a mask ratio of 0.1 yields marginally better results, which relies on single lightly masked 523 image to give a relatively accurate prediction but discards the mutual learning of dual masks. 524

525 526

527

496

497 498 499

500

# 5 CONCLUSION

In this work, we present a novel perspective on masked reconstruction by reinterpreting it as a 528 sparse signal reconstruction problem and theoretically prove the effectiveness of the dual form of 529 complementary masks. Based on this theoretical foundation, we propose MaskTwins, an effective 530 strategy that utilizes complementary masks to simultaneously enhance the robust feature extraction 531 for domain-adaptive segmentation. Our MaskTwins has demonstrated remarkable superiority over 532 the state-of-the-art methods across a diverse range of domain adaptation scenarios, spanning from 533 natural to biological imaging and from 2D to 3D modalities. For instance, MaskTwins respectively 534 achieves significant performance improvements by +2.7% and +2.5% on SYNTHIA→Cityscapes and biological datasets. Since MaskTwins performs masked image consistency without extra an-536 notations, it offers a flexible technique that can be seamlessly incorporated with other methods to 537 further facilitate the learning of domain-invariant features, ensuring the cross-domain knowledge adaptation process. In the future, we will continue to explore the potential of MaskTwins in a 538 broader spectrum of visual recognition challenges, including but not limited to domain-adaptive video segmentation and image classification.

#### 540 REFERENCES 541

547

553

554

561

562

567

569

581

584

585

542	Nikita Araslanov and Stefan Roth. Self-supervised augmentation consistency for adapting seman-
543	tic segmentation. In Proceedings of the IEEE/CVF conference on computer vision and pattern
544	recognition, pp. 15384–15394, 2021.

- Hangbo Bao, Li Dong, Songhao Piao, and Furu Wei. BEiT: BERT pre-training of image transform-546 ers. In International Conference on Learning Representations, 2022.
- Shai Ben-David, John Blitzer, Koby Crammer, and Fernando Pereira. Analysis of representations 548 for domain adaptation. Advances in neural information processing systems, 19, 2006. 549
- 550 Shai Ben-David, John Blitzer, Koby Crammer, Alex Kulesza, Fernando Pereira, and Jennifer Wort-551 man Vaughan. A theory of learning from different domains. Machine learning, 79:151–175, 552 2010.
  - Róger Bermúdez-Chacón, Pablo Márquez-Neila, Mathieu Salzmann, and Pascal Fua. A domainadaptive two-stream u-net for electron microscopy image segmentation. In ISBI. IEEE, 2018.
- 556 Tom B Brown. Language models are few-shot learners. arXiv preprint arXiv:2005.14165, 2020.
- 558 Minghao Chen, Hongyang Xue, and Deng Cai. Domain adaptation for semantic segmentation with maximum squares loss. In ICCV, 2019. 559
  - Qi Chen, Wei Huang, Yueyi Zhang, and Zhiwei Xiong. Domain adaptive synapse detection with weak point annotations. In ISBI. IEEE, 2024.
- 563 Jaehoon Choi, Taekyung Kim, and Changick Kim. Self-ensembling with gan-based data augmentation for domain adaptation in semantic segmentation. In Proceedings of the IEEE/CVF interna-564 tional conference on computer vision, pp. 6830-6840, 2019. 565
- 566 Marius Cordts, Mohamed Omran, Sebastian Ramos, Timo Rehfeld, Markus Enzweiler, Rodrigo Benenson, Uwe Franke, Stefan Roth, and Bernt Schiele. The cityscapes dataset for semantic 568 urban scene understanding. In CVPR, 2016.
- Jacob Devlin. Bert: Pre-training of deep bidirectional transformers for language understanding. 570 arXiv preprint arXiv:1810.04805, 2018. 571
- 572 Terrance DeVries. Improved regularization of convolutional neural networks with cutout. arXiv 573 preprint arXiv:1708.04552, 2017. 574
- Xiaoyi Dong, Jianmin Bao, Ting Zhang, Dongdong Chen, Weiming Zhang, Lu Yuan, Dong Chen, 575 Fang Wen, and Nenghai Yu. Bootstrapped masked autoencoders for vision bert pretraining. In 576 ECCV, 2022. 577
- 578 Xiaoyi Dong, Jianmin Bao, Ting Zhang, Dongdong Chen, Weiming Zhang, Lu Yuan, Dong Chen, 579 Fang Wen, Nenghai Yu, and Baining Guo. Peco: Perceptual codebook for bert pre-training of 580 vision transformers. In AAAI, 2023.
- David L Donoho. Compressed sensing. IEEE Transactions on Information Theory, 52(4):1289-582 1306, 2006. 583
  - Alexey Dosovitskiy. An image is worth 16x16 words: Transformers for image recognition at scale. arXiv preprint arXiv:2010.11929, 2020.
- Stephan Gerhard, Jan Funke, Julien Martel, Albert Cardona, and Richard Fetter. Segmented 587 anisotropic sstem dataset of neural tissue. figshare, pp. 0-0, 2013. 588
- 589 Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, 590 Aaron Courville, and Yoshua Bengio. Generative adversarial nets. Advances in neural information processing systems, 27, 2014. 592
- Dayan Guan, Jiaxing Huang, Aoran Xiao, and Shijian Lu. Domain adaptive video segmentation via temporal consistency regularization. In ICCV, 2021.

594 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In CVPR, 2016. 596 Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked 597 autoencoders are scalable vision learners. In CVPR, 2022. 598 Judy Hoffman, Dequan Wang, Fisher Yu, and Trevor Darrell. Fcns in the wild: Pixel-level adversar-600 ial and constraint-based adaptation. arXiv preprint arXiv:1612.02649, 2016. 601 Lukas Hoyer, Dengxin Dai, and Luc Van Gool. Daformer: Improving network architectures and 602 training strategies for domain-adaptive semantic segmentation. In CVPR, 2022a. 603 604 Lukas Hoyer, Dengxin Dai, and Luc Van Gool. Hrda: Context-aware high-resolution domain-605 adaptive semantic segmentation. In ECCV, 2022b. Lukas Hoyer, Dengxin Dai, Haoran Wang, and Luc Van Gool. Mic: Masked image consistency for 607 context-enhanced domain adaptation. In CVPR, 2023. 608 Wei Huang, Chang Chen, Zhiwei Xiong, Yueyi Zhang, Xuejin Chen, Xiaoyan Sun, and Feng Wu. 609 Semi-supervised neuron segmentation via reinforced consistency learning. IEEE Transactions on 610 Medical Imaging, 41(11):3016–3028, 2022a. 611 612 Wei Huang, Xiaoyu Liu, Zhen Cheng, Yueyi Zhang, and Zhiwei Xiong. Domain adaptive mitochon-613 dria segmentation via enforcing inter-section consistency. In MICCAI, 2022b. 614 Xun Huang and Serge Belongie. Arbitrary style transfer in real-time with adaptive instance normal-615 ization. In Proceedings of the IEEE International Conference on Computer Vision, 2017. 616 617 Zhengkai Jiang, Yuxi Li, Ceyuan Yang, Peng Gao, Yabiao Wang, Ying Tai, and Chengjie Wang. 618 Prototypical contrast adaptation for domain adaptive semantic segmentation. In ECCV, pp. 36– 54. Springer, 2022. 619 620 Ying Jin, Ximei Wang, Mingsheng Long, and Jianmin Wang. Minimum class confusion for versatile 621 domain adaptation. In ECCV, 2020. 622 Diederik P Kingma. Adam: A method for stochastic optimization. arXiv preprint arXiv:1412.6980, 623 2014. 624 625 Xiangwen Kong and Xiangyu Zhang. Understanding masked image modeling via learning occlusion 626 invariant feature. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern 627 Recognition, 2023. 628 Dong-Hyun Lee et al. Pseudo-label: The simple and efficient semi-supervised learning method for 629 deep neural networks. In ICML, 2013. 630 631 Kisuk Lee, Jonathan Zung, Peter Li, Viren Jain, and H Sebastian Seung. Superhuman accuracy on the snemi3d connectomics challenge. arXiv preprint arXiv:1706.00120, 2017. 632 633 Guangrui Li, Guoliang Kang, Wu Liu, Yunchao Wei, and Yi Yang. Content-consistent matching for 634 domain adaptive semantic segmentation. In ECCV, 2020. 635 Jing Li, Kang Zhou, Shenhan Qian, Wen Li, Lixin Duan, and Shenghua Gao. Feature re-636 representation and reliable pseudo label retraining for cross-domain semantic segmentation. *IEEE* 637 Transactions on Pattern Analysis and Machine Intelligence, 46(3):1682–1694, 2022. 638 639 Yicong Li, Wanhua Li, Qi Chen, Wei Huang, Yuda Zou, Xin Xiao, Kazunori Shinomiya, Pat Gunn, 640 Nishika Gupta, Alexey Polilov, et al. Waspsyn: A challenge for domain adaptive synapse detec-641 tion in microwasp brain connectomes. *IEEE Transactions on Medical Imaging*, 2024. 642 Yunsheng Li, Lu Yuan, and Nuno Vasconcelos. Bidirectional learning for domain adaptation of 643 semantic segmentation. In CVPR, 2019. 644 645 Dongnan Liu, Donghao Zhang, Yang Song, Fan Zhang, Lauren O'Donnell, Heng Huang, Mei Chen, and Weidong Cai. Pdam: A panoptic-level feature alignment framework for unsupervised domain 646 adaptive instance segmentation in microscopy images. *IEEE Transactions on Medical Imaging*, 647 40(1):154-165, 2020a.

648 649 650	Hao Liu, Xinghua Jiang, Xin Li, Antai Guo, Yiqing Hu, Deqiang Jiang, and Bo Ren. The devil is in the frequency: Geminated gestalt autoencoder for self-supervised visual pre-training. In <i>AAAI</i> , 2023.
651	
652	Xiaofeng Liu, Yuzhuo Han, Song Bai, Yi Ge, Tianxing Wang, Xu Han, Site Li, Jane You, and Jun
653 654	In AAAI, 2020b.
655	
656	domain adaptation. Advances in neural information processing systems, 31, 2018.
657 658	I Loshchilov. Decoupled weight decay regularization. arXiv preprint arXiv:1711.05101, 2017.
659	
660	Aurelien Lucchi, Yunpeng Li, and Pascal Fua. Learning for structured prediction using approximate subgradient descent with working sets. In <i>CVPR</i> , 2013.
660	V 'L D' L' T. C. L. '. V J.V' V. C'. 'C' Constant batt
663	rawei Luo, Ping Liu, Tao Guan, Junqing Yu, and Yi Yang. Significance-aware information bottle- neck for domain adaptive semantic segmentation. In <i>ICCV</i> , 2019.
664	Versi Lue Dine Lie Liene Zhang Tre Cren Luncing Versud Vi Veng, Catagory level adverse
665	rawel Luo, Ping Liu, Liang Zheng, Tao Guan, Junqing Tu, and Yi Yang. Category-level adversar-
666 667	Analysis and Machine Intelligence, 44(8):3940–3956, 2021.
668	Brian W Matthews Comparison of the predicted and observed secondary structure of the phage
669	lysozyme. <i>Biochimica et Biophysica Acta (BBA)-Protein Structure</i> , 405(2):442–451, 1975.
670	Ke Mei, Chuang Zhu, Jiagi Zou, and Shanghang Zhang. Instance adaptive self-training for unsuper-
671 672	vised domain adaptation. In ECCV, 2020.
673	Luke Melas-Kyriazi and Ariun K Manrai Pixmatch: Unsupervised domain adaptation via pixel-
674 675	wise consistency training. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and</i>
676	<i>Futern Recognition</i> , pp. 12455–12445, 2021.
677 678 679	Jialin Peng, Jiajin Yi, and Zhimin Yuan. Unsupervised mitochondria segmentation in em images via domain adaptive multi-task learning. <i>IEEE Journal of Selected Topics in Signal Processing</i> , 14 (6):1199–1209, 2020.
680	Xingchao Peng, Ben Usman, Neela Kaushik, Judy Hoffman, Deguan Wang, and Kate Saenko
681	Visda: The visual domain adaptation challenge. <i>arXiv preprint arXiv:1710.06924</i> , 2017.
682	Harsh Rangwani Sumukh K Aithal Mayank Mishra Arihant Jain and Venkatesh Babu Radhakr-
683 684	ishnan. A closer look at smoothness in domain adversarial training. In <i>ICML</i> , 2022.
685	German Ros, Laura Sellart, Joanna Materzynska, David Vazquez, and Antonio M Lopez. The
686	synthia dataset: A large collection of synthetic images for semantic segmentation of urban scenes.
687	In CVPR, 2016.
688	
689	Ukcheol Shin, Kyunghyun Lee, In So Kweon, and Jean Oh. Complementary random masking for rgb-thermal semantic segmentation. In <i>ICRA</i> , 2024.
601	
602	M Naseer Subhani and Mohsen Ali. Learning from scale-invariant examples for domain adaptation
693	in semantic segmentation. In ECCV, 2020.
694	Rui Sun, Huayu Mai, Naisong Luo, Tianzhu Zhang, Zhiwei Xiong, and Feng Wu. Structure-
695	decoupled adaptive part alignment network for domain adaptive mitochondria segmentation. In
696	MICCAI, 2023.
~ ~ ~	

- Antti Tarvainen and Harri Valpola. Mean teachers are better role models: Weight-averaged consistency targets improve semi-supervised deep learning results. Advances in neural information processing systems, 30, 2017.
- 701 Wilhelm Tranheden, Viktor Olsson, Juliano Pinto, and Lennart Svensson. Dacs: Domain adaptation via cross-domain mixed sampling. In *WACV*, 2021.

702 703 704	Yi-Hsuan Tsai, Wei-Chih Hung, Samuel Schulter, Kihyuk Sohn, Ming-Hsuan Yang, and Manmohan Chandraker. Learning to adapt structured output space for semantic segmentation. In <i>CVPR</i> , 2018.
705 706	Tuan-Hung Vu, Himalaya Jain, Maxime Bucher, Matthieu Cord, and Patrick Pérez. Advent: Adversarial entropy minimization for domain adaptation in semantic segmentation. In <i>CVPR</i> , 2019a.
707 708	Tuan-Hung Vu, Himalaya Jain, Maxime Bucher, Matthieu Cord, and Patrick Pérez. Dada: Depth- aware domain adaptation in semantic segmentation. In <i>ICCV</i> , 2019b.
709 710 711	Haoqing Wang, Yehui Tang, Yunhe Wang, Jianyuan Guo, Zhi-Hong Deng, and Kai Han. Masked image modeling with local multi-scale reconstruction. In <i>CVPR</i> , 2023.
712 713	Yuxi Wang, Junran Peng, and ZhaoXiang Zhang. Uncertainty-aware pseudo label refinery for do- main adaptive semantic segmentation. In <i>ICCV</i> , 2021.
714 715 716 717	Zhonghao Wang, Mo Yu, Yunchao Wei, Rogerio Feris, Jinjun Xiong, Wen-mei Hwu, Thomas S Huang, and Honghui Shi. Differential treatment for stuff and things: A simple unsupervised domain adaptation method for semantic segmentation. In <i>CVPR</i> , 2020.
718 719	Chen Wei, Haoqi Fan, Saining Xie, Chao-Yuan Wu, Alan Yuille, and Christoph Feichtenhofer. Masked feature prediction for self-supervised visual pre-training. In <i>CVPR</i> , 2022.
720 721 722 723	Donglai Wei, Zudi Lin, Daniel Franco-Barranco, Nils Wendt, Xingyu Liu, Wenjie Yin, Xin Huang, Aarush Gupta, Won-Dong Jang, Xueying Wang, et al. Mitoem dataset: Large-scale 3d mitochondria instance segmentation from em images. In <i>MICCAI</i> , 2020.
724 725	Siqi Wu, Chang Chen, Zhiwei Xiong, Xuejin Chen, and Xiaoyan Sun. Uncertainty-aware label rectification for domain adaptive mitochondria segmentation. In <i>MICCAI</i> , 2021.
726 727 728	Enze Xie, Wenhai Wang, Zhiding Yu, Anima Anandkumar, Jose M Alvarez, and Ping Luo. Seg- former: Simple and efficient design for semantic segmentation with transformers. <i>Advances in</i> <i>neural information processing systems</i> , 34:12077–12090, 2021.
729 730 731 732	Jiahao Xie, Wei Li, Xiaohang Zhan, Ziwei Liu, Yew Soon Ong, and Chen Change Loy. Masked frequency modeling for self-supervised visual pre-training. <i>arXiv preprint arXiv:2206.07706</i> , 2022a.
733 734	Zhenda Xie, Zheng Zhang, Yue Cao, Yutong Lin, Jianmin Bao, Zhuliang Yao, Qi Dai, and Han Hu. Simmim: A simple framework for masked image modeling. In <i>CVPR</i> , 2022b.
735 736 737	Tongkun Xu, Weihua Chen, Pichao Wang, Fan Wang, Hao Li, and Rong Jin. Cdtrans: Cross-domain transformer for unsupervised domain adaptation. <i>arXiv preprint arXiv:2109.06165</i> , 2021.
738 739 740	Jihan Yang, Ruijia Xu, Ruiyu Li, Xiaojuan Qi, Xiaoyong Shen, Guanbin Li, and Liang Lin. An adversarial perturbation oriented domain adaptation approach for semantic segmentation. In <i>AAAI</i> , 2020a.
741 742 743	Jinyu Yang, Weizhi An, Sheng Wang, Xinliang Zhu, Chaochao Yan, and Junzhou Huang. Label- driven reconstruction for domain adaptation in semantic segmentation. In <i>ECCV</i> , 2020b.
744 745	Jinyu Yang, Weizhi An, Chaochao Yan, Peilin Zhao, and Junzhou Huang. Context-aware domain adaptation in semantic segmentation. In <i>WACV</i> , 2021.
746 747 748	Jinyu Yang, Jingjing Liu, Ning Xu, and Junzhou Huang. Tvt: Transferable vision transformer for unsupervised domain adaptation. In <i>WACV</i> , 2023.
749 750 751	Senqiao Yang, Jiarui Wu, Jiaming Liu, Xiaoqi Li, Qizhe Zhang, Mingjie Pan, Yulu Gan, Zehui Chen, and Shanghang Zhang. Exploring sparse visual prompt for domain adaptive dense prediction. In <i>AAAI</i> , 2024.
752 753	Yanchao Yang and Stefano Soatto. Fda: Fourier domain adaptation for semantic segmentation. In <i>CVPR</i> , 2020.
755	Dan Yin, Wei Huang, Zhiwei Xiong, and Xuejin Chen. Class-aware feature alignment for domain adaptative mitochondria segmentation. In <i>MICCAI</i> . Springer, 2023.

- 756 Fei Yu, Mo Zhang, Hexin Dong, Sheng Hu, Bin Dong, and Li Zhang. Dast: Unsupervised domain 757 adaptation in semantic segmentation based on discriminator attention and self-training. In AAAI, 758 2021. 759 Marvin Zhang, Sergey Levine, and Chelsea Finn. Memo: Test time robustness via adaptation and 760 augmentation. Advances in neural information processing systems, 35:38629–38642, 2022. 761 762 Pan Zhang, Bo Zhang, Ting Zhang, Dong Chen, Yong Wang, and Fang Wen. Prototypical pseudo label denoising and target structure learning for domain adaptive semantic segmentation. In CVPR, 763 764 2021. 765 Qiming Zhang, Jing Zhang, Wei Liu, and Dacheng Tao. Category anchor-guided unsupervised do-766 main adaptation for semantic segmentation. Advances in neural information processing systems, 767 32, 2019a. 768 Yuchen Zhang, Tianle Liu, Mingsheng Long, and Michael Jordan. Bridging theory and algorithm 769 for domain adaptation. In ICML, 2019b. 770 771 Dong Zhao, Shuang Wang, Qi Zang, Dou Quan, Xiutiao Ye, Rui Yang, and Licheng Jiao. Learning 772 pseudo-relations for cross-domain semantic segmentation. In ICCV, 2023. 773 Zhedong Zheng and Yi Yang. Rectifying pseudo label learning via uncertainty estimation for domain 774 adaptive semantic segmentation. International Journal of Computer Vision, 129(4):1106–1120, 775 2021. 776 Zhedong Zheng and Yi Yang. Adaptive boosting for domain adaptation: Toward robust predictions 777 in scene segmentation. IEEE Transactions on Image Processing, 31:5371-5382, 2022. 778 779 Jinghao Zhou, Chen Wei, Huiyu Wang, Wei Shen, Cihang Xie, Alan Yuille, and Tao Kong. ibot: 780 Image bert pre-training with online tokenizer. arXiv preprint arXiv:2111.07832, 2021. 781 Qianyu Zhou, Zhengyang Feng, Qiqi Gu, Guangliang Cheng, Xuequan Lu, Jianping Shi, and 782 Lizhuang Ma. Uncertainty-aware consistency regularization for cross-domain semantic segmen-783 tation. Computer Vision and Image Understanding, 221:103448, 2022a. 784 785 Qianyu Zhou, Zhengyang Feng, Qiqi Gu, Jiangmiao Pang, Guangliang Cheng, Xuequan Lu, Jianping Shi, and Lizhuang Ma. Context-aware mixup for domain adaptive semantic segmentation. 786 *IEEE Transactions on Circuits and Systems for Video Technology*, 33(2):804–817, 2022b. 787 788 Wei Zhou, Yukang Wang, Jiajia Chu, Jiehua Yang, Xiang Bai, and Yongchao Xu. Affinity space 789 adaptation for semantic segmentation across domains. IEEE Transactions on Image Processing, 790 30:2549-2561, 2020. 791 Yang Zou, Zhiding Yu, BVK Kumar, and Jinsong Wang. Unsupervised domain adaptation for 792 semantic segmentation via class-balanced self-training. In ECCV, 2018. 793 794 Yang Zou, Zhiding Yu, Xiaofeng Liu, BVK Kumar, and Jinsong Wang. Confidence regularized 795 self-training. In CVPR, 2019. 796 797 EXTENDED RELATED WORKS ON UDA А 798 799 Adversarial learning Hoffman et al. (2016) are the first to apply the adversarial approach for 800 UDA on semantic segmentation to encourage domain-invariant alignment globally. SIBAN (Luo 801 et al., 2019) employs a significance-aware information bottleneck (SIB) before the adversarial fea-802 ture adaptation to extract latent representations in semantic segmentation tasks. FDA (Yang & 803 Soatto, 2020) performs spectral transfer by swapping the low-frequency component of the spectrum 804 of one with the other. APODA (Yang et al., 2020a) explicitly trains a domain-invariant classifier by
- of one with the other. APODA (Yang et al., 2020a) explicitly trains a domain-invariant classifier by
  generating and defensing against point-wise feature space adversarial perturbations, in order to adapt
  the representations of the tail classes or small objects for semantic segmentation. SIMWang et al.
  (2020) and CLAN (Luo et al., 2021) apply category-level alignment to minimize the discrepancy
  between the source and target distributions. DAST (Yu et al., 2021) proposes a self-training strategy which adaptively improves the decision boundary of the model for target domain and implicitly
  facilitates the extraction of domain-invariant features.

810 **Pseudo-label self-training** DADA (Vu et al., 2019b) introduces a novel depth-aware adaptation 811 scheme while BDL (Li et al., 2019) proposes a novel bidirectional learning framework for domain 812 adaptation of segmentation. LSE (Subhani & Ali, 2020) exploits scale-invarince property of the 813 model to generate pseudo-labels. DACS (Tranheden et al., 2021) mixes images from the two do-814 mains along with the corresponding labels and pseudo labels to perform Cross-domain mixed Sampling. Some generative methods try to acquire target-like synthetic images by content-consistent 815 matching (CCM) (Li et al., 2020) or label-driven reconstruction (LDR) Yang et al. (2020b). To 816 improve the quality of pseudo labels, UncerDA (Wang et al., 2021) provides an uncertainty-aware 817 pseudo label assignment strategy while RPLR (Li et al., 2022) retrains the networks using selected 818 reliable pseudo labels. Many works focus on consistency regularization to capture contextual re-819 lations, such as CD-SAM (Yang et al., 2021), UACR (Zhou et al., 2022a), CAMix (Zhou et al., 820 2022b), HRDA (Hoyer et al., 2022b) and MIC (Hoyer et al., 2023). Researchers also conducted ex-821 tensive attempts, including affinity in ASA (Zhou et al., 2020), representative prototypes in ProDA 822 (Zhang et al., 2021), and Transformer architecture in DAFormer (Hoyer et al., 2022a). In conclu-823 sion, consistency-based methods try to learn domain-variant feature to enhance the robustness of the 824 model, which align with the constrained entropy minimization perspective of MEMO (Zhang et al., 825 2022).

- 826
- 827

838 839

840 841 842

843 844 845

846

847

848

849

828 **Theory for UDA** The theoretical works (Ben-David et al., 2006; 2010; Zhang et al., 2019b) 829 provide fundamental insights into UDA, especially concerning domain discrepancy and theoretical bounds. Specifically, they study margin bounds for classification tasks at the distribution level, 830 while we focus on segmentation tasks and the theory of Masked Image Modeling and compressed 831 sensing at the image level. We have analyzed the information preservation, generalization bounds 832 and feature consistency to demonstrate the effectiveness of complementary masking. Zhang et al. 833 (2019b) discuss generalization bounds based on empirical Rademacher complexity, building upon 834 the domain adaptation theories presented in previous works such as those by Ben-David et al. (2006; 835 2010). We preliminary observe that there exists deeper connections between these works and ours. 836 Hopefully, we will make further theoretical analysis in the future work. 837

# **B** MASKTWINS TRAINING PROCEDURE

We provide the overall training procedure of MaskTwins for image segmentation in Algorithm 1.

# Algorithm 1 MaskTwins Algorithm

- **Input:** Source domain  $\mathcal{D}_S$ , Target domain  $\mathcal{D}_T$ , student model  $f_{\theta}$ , teacher model  $f_{\phi}$ , the total iteration number N.
  - 1: Initialize network parameter  $\theta$  with ImageNet pre-trained parameters. Initialize teacher network  $\phi$  randomly.
- 850 2: for iteration = 1 to N do 851  $\begin{array}{c} x^{S}, y^{S} \sim \mathcal{D}_{S}. \\ x^{T} \sim \mathcal{D}_{T}. \end{array}$ 3: 852 4:  $\begin{array}{l} \overset{\sim}{p}^{S} \leftarrow f_{\theta}(x^{S}), \\ \hat{y}^{T} \leftarrow \operatorname{argmax} f_{\phi}(x^{T}). \end{array}$ 853 5: 854 6: 855  $X_D^T, X_{1-D}^T \leftarrow$  Patch-wise complementary masking by Eq. 8 and 9. 7:  $p_D^T \leftarrow f_\theta(x_D^T), p_{1-D}^T \leftarrow f_\theta(x_{1-D}^T).$ 856 8:  $\mathcal{L}_{total} \leftarrow$  Total loss by Eq. 15. 9: Compute  $\nabla_{\theta} \mathcal{L}_{\text{total}}$  by back-propagation. 858 10: 859 11: Perform stochastic gradient descent on  $\theta$ . 12: Update teacher network  $\phi$  with  $\theta$ . 860 13: end for 861 14: return  $f_{\theta}$ . 862
- 863

# 864 C EXPERIMENTAL DETAILS

# C.1 NATURAL IMAGE SEMANTIC SEGMENTATION

Following common UDA protocols (Tsai et al., 2018; Zhou et al., 2022b), we use the synthetic dataset SYNTHIA (Ros et al., 2016) as the source domain, and the real dataset Cityscapes (Cordts et al., 2016) as the target domain. SYNTHIA is a synthetic dataset composed of 9,400 annotated images with the resolution of  $1280 \times 960$ , while Cityscapes consists of 2,975 training and 500 validation real-world images.

We evaluate MaskTwins based on the HRDA (Hoyer et al., 2022b) architecture with a MiT-B5 encoder (Xie et al., 2021) pretrained on ImageNet. To be specific, we follow the DAFormer (Hoyer et al., 2022a) self-training strategy and training parameters, i.e. AdamW (Loshchilov, 2017) with a learning rate of  $6 \times 10^{-5}$  for the encoder and  $6 \times 10^{-4}$  for the decoder, 40k training iterations, a batch size of 2, linear learning rate warmup, a loss weight  $\lambda_{st} = 1$ , an EMA factor  $\alpha = 0.999$  and DACS (Tranheden et al., 2021) data augmentation.

879

866

867

# 880 C.2 MITOCHONDRIA SEMANTIC SEGMENTATION

We evaluate the proposed method on three challenging EM datasets for 2D domain adaptive mi-882 tochondria segmentation tasks: VNC III (Gerhard et al., 2013), Lucchi (Lucchi et al., 2013) and 883 MitoEM (Wei et al., 2020) dataset. VNC III consists of 20 sections of size  $1024 \times 1024$ . The 884 training subset (Subset1) and the test subset (Subset2) of Lucchi each contain 165 images, with a 885 resolution of  $1024 \times 768$  pixels. MitoEM dataset can be divided into MitoEM-R(Rat) and MitoEM-886 H(Human). Each volume contains 1000 images of size  $4096 \times 4096$ , with the first 500 images anno-887 tated. Following Huang et al. (2022b), four widely used metrics are used for evaluation, i.e., mean Average Precision (mAP), F1 score, Mattews Correlation Coefficient (MCC) (Matthews, 1975) and Intersection over Union (IoU). 889

We use a five-stage U-Net following Huang et al. (2022b) and Yin et al. (2023). During training, we randomly crop the original EM section into  $512 \times 512$  with random augmentation including flip, transpose, rotate, resize and elastic transformation. All models are trained for 200k iterations with a batch size of 2. We use the Adam optimizer (Kingma, 2014) with  $\beta_1 = 0.9$  and  $\beta_2 = 0.999$ . The learning rate is set at  $1 \times 10^{-4}$  and has a polynomial decay with a power of 0.9.

895 896

897

# C.3 SYNAPSE DETECTION

To further diversify the experiment settings, we study the 3D domain adaptive synapse detection task using the WASPSYN (Li et al., 2024) dataset. The WASPSYN dataset includes 14 image chunks from different brain regions of Megaphragma viggianii, and five of them have point annotations. Specifically, we take the first one as the source data, and the remaining four chunks are considered target data.

The experiments are performed based on 3D ResUNet following Lee et al. (2017). Considering the data are imaged with an isotropic voxel size, we adopt isotropic 3D convolutions. Specifically, we set the kernel size for the initial embedding layer to be  $5 \times 5 \times 5$ , whereas the convolutional layers subsequently utilize a default kernel size of  $3 \times 3 \times 3$ . In the training process, we use a crop size of  $96 \times 96 \times 96$  with a batch size of 4 and train for 200k iterations. We use an Adam optimizer with a base learning rate of 0.0001 and a linear warming up in the first 1000 iterations.

- 909
- 910
- 911
- 912
- 913 914
- 915
- 916
- 917

# 918 D MORE RESULTS

# 920 D.1 ABLATION STUDY EXTENSION

We fully ablate the used components on MitoEM-H  $\rightarrow$  MitoEM-R in Tables 5 and 6. Upon observation, mAP stands out as an outlier while F1, MCC, and IoU consistently show similar trends. So we pay more attention to the latter three metrics, especially the IoU.

As shown in Table 5, adding the consistency loss  $L_{cl}$  improves the performance over the supervised loss. Further, we separately incorporate the randomly masked loss  $L_{rm}$  and the complementary masked loss  $L_{cm}$ . The results indicate that both losses contribute to performance improvements, with our proposed complementary masking strategy being more effective than the random masking strategy. The results with  $L_{rm}$  and  $L_{cm}$  have been visually shown in Figure 4 in the main paper.

930 The EMA teacher realizes a temporal ensemble of previous student models, which increases the 931 robustness and temporal stability of pseudo-labels. It is a common strategy used in semi-supervised 932 learning and UDA. In our work, we adopt the EMA teacher to keep consistent with previous meth-933 ods, such as CAMix (Zhou et al., 2022b), MIC (Hoyer et al., 2023), DAFormer (Hoyer et al., 2022a), 934 etc. Table 6 show that both EMA and AdaIN contribute to performance improvements, with EMA 935 having a more significant impact. In Table 5, adding the complementary masked loss to the exist-936 ing consistency loss yields a notable improvement (from 80.64 to 81.88 in IoU). Therefore, while we use some well-constructed modules, the main performance improvement comes from the key 937 contribution of complementary masking. 938

Table 5: Ablation study of each loss component on MitoEM-H  $\rightarrow$  MitoEM-R. The mean and standard deviation are computed over 3 random seeds.  $L_{sup}$  = supervised loss,  $L_{cl}$  = consistency loss,  $L_{cm}$  = complementary masked loss,  $L_{rm}$  = randomly masked loss, with a mask ratio of 0.5.

	mAP	F1	MCC	IoU
$L_{sup}$	$96.38 \pm 0.18$	$88.95 \pm 0.07$	$88.60{\scriptstyle~\pm 0.07}$	$80.11{\scriptstyle~\pm 0.11}$
$L_{sup} + L_{cl}$	$96.60{\scriptstyle~\pm 0.35}$	$89.27{\scriptstyle~\pm 0.10}$	$88.94{\scriptstyle~\pm 0.12}$	$80.64 \pm 0.17$
$L_{sup} + L_{cl} + L_{rm}$	$96.84{\scriptstyle~\pm 0.10}$	$89.80{\scriptstyle~\pm 0.03}$	$89.45{\scriptstyle~\pm 0.04}$	$81.49{\scriptstyle~\pm 0.04}$
$L_{sup} + L_{cl} + L_{cm}$	$96.87{\scriptstyle~\pm 0.06}$	$90.03{\scriptstyle~\pm 0.06}$	$89.66 \pm 0.04$	$81.88 \pm 0.09$

Table 6: Ablation study of EMA and AdaIN on MitoEM-H  $\rightarrow$  MitoEM-R. The mean and standard deviation are computed over 3 random seeds.

	mAP	F1	MCC	IoU
Ours w/o AdaIN & EMA	$96.74{\scriptstyle~\pm 0.08}$	$89.61{\scriptstyle~\pm 0.05}$	$89.23 \pm 0.04$	$81.18 \pm 0.08$
Ours w/o EMA	$96.85{\scriptstyle~\pm 0.15}$	$89.74{\scriptstyle~\pm 0.02}$	$89.38{\scriptstyle~\pm 0.04}$	$81.40{\scriptstyle~\pm 0.04}$
Ours w/o AdaIN	$96.89{\scriptstyle~\pm 0.14}$	$89.88{\scriptstyle~\pm 0.05}$	$89.52 \pm 0.03$	$81.63{\scriptstyle~\pm 0.08}$
Ours	$96.87{\scriptstyle~\pm 0.06}$	$90.03{\scriptstyle~\pm 0.06}$	$89.66 \pm 0.04$	$81.88 \pm 0.09$

965 966

939 940 941

942

943

955

967

968

969

970

#### 972 D.2 CLASSIFICATION TASKS

While our main focus is pixel-wise segmentation tasks, we extend our method to classification tasksto further validate its effectiveness.

We conduct additional experiments on the VisDA-2017 dataset (Peng et al., 2017), which consists of 280,000 synthetic and real images of 12 classes, with ResNet-101 (He et al., 2016) and ViTB/16 (Dosovitskiy, 2020). For UDA training, we follow SDAT (Rangwani et al., 2022), which utilizes CDAN (Long et al., 2018) with MCC (Jin et al., 2020) and a smoothness enhancing loss. We use the same training parameters, i.e. SGD with a learning rate of 0.002, a batch size of 32, and a smoothness parameter of 0.02. We use a patch size b=64, a mask ratio r=0.5, a loss weight  $\lambda_{cm} = 0.01$ .

As shown in Tables 7 and 8, our method improves the UDA performance by +0.3 and +0.4 percent points when used with a ViT and ResNet network, respectively. The improvement is consistent over almost all classes.

Table 7: Image classification accuracy in % on VisDA-2017 for UDA with ViT-B/16. "Bcycl" stands for *bicycle*, "PR" for *person*, "Sktb" for *skateboard*. The competitors include TVT (Yang et al., 2023), CDTrans (Xu et al., 2021), SDAT (Rangwani et al., 2022), and MIC (Hoyer et al., 2023). The results are adopted from Hoyer et al. (2023).

	Method		Plane	Bcycl	Bus	Car	Horse	e Knife	Motor	PR	Plant	Sktb	Train	Truck	Mean
_	TVT		92.9	85.6	77.5	60.5	93.6	98.2	89.3	76.4	93.6	92.0	91.7	55.7	83.9
	CDTrans		97.1	90.5	82.4	77.5	96.6	96.1	93.6	88.6	97.9	86.9	90.3	62.8	88.4
	SDAT		98.4	90.9	85.4	82.1	98.5	97.6	96.3	86.1	96.2	96.7	92.9	56.8	89.8
	SDAT w/ MAE		97.1	88.4	80.9	75.3	95.4	97.9	94.3	85.5	95.8	91.0	93.0	65.4	88.4
	MIC		<u>99.0</u>	<u>93.3</u>	<u>86.5</u>	<u>87.6</u>	98.9	<u>99.0</u>	97.2	89.8	98.9	<u>98.9</u>	96.5	<u>68.0</u>	92.8
	Ours	ĺ	99.1	95.0	86.6	89.0	<u>98.8</u>	99.3	<u>96.8</u>	<u>88.3</u>	<u>98.8</u>	99.1	97.2	<b>69.7</b>	93.1

Table 8: Image classification accuracy in % on VisDA-2017 for UDA with ResNet-101. "Bcycl" stands for *bicycle*, "PR" for *person*, "Sktb" for *skateboard*. The competitors include CDAN (Long et al., 2018), MCC (Jin et al., 2020), SDAT (Rangwani et al., 2022), and MIC (Hoyer et al., 2023).
The results are adopted from Hoyer et al. (2023).

-	Method	Plane	Bcycl	Bus	Car	Horse	e Knife	Motor	PR	Plant	Sktb	Train	Truck	Mean
	CDAN	85.2	66.9	83.0	50.8	84.2	74.9	88.1	74.5	83.4	76.0	81.9	38.0	73.9
	MCC	88.1	80.3	80.5	71.5	90.1	93.2	85.0	71.6	89.4	73.8	85.0	36.9	78.8
	SDAT	95.8	85.5	76.9	69.0	93.5	97.4	88.5	78.2	93.1	91.6	86.3	55.3	84.3
	MIC	96.7	88.5	84.2	74.3	96.0	96.3	90.2	81.2	<b>94.3</b>	95.4	88.9	56.6	86.9
	Ours	96.9	88.8	<u>81.8</u>	77.1	96.4	<u>97.2</u>	90.3	83.8	<u>93.3</u>	<u>94.8</u>	90.2	57.4	87.3

#### D.3 VISUALIZATION RESULTS

We visualize the segmentation results of MaskTwins and qualitatively compare with the state-of-art methods on SYNTHIA→Cityscapes in Figure 5 and mitochondria datasets in Figure 8. We also provide more visualization for synapse detection on the WASPSYN dataset in Figure 6 and 7. 



Figure 5: More segmentation results on SYNTHIA→Cityscapes.



Figure 6: Visualization of the volume in the WASPSYN dataset. Left column to right column: sections from X-Y, X-Z, and Y-Z plane.



Figure 7: An example of visualization of the detection results of pre-synapse (left) and post-synapse (right). Dots and lines: magenta-true positive, yellowfalse negative, and cyan-false positive. 



Figure 8: More segmentation results on VNC III→Lucchi Subset1 (row 1 and 2), VNC III→Lucchi Subset2 (row 3 and 4), MitoEM-R→MitoEM-H (row 5 and 6) and MitoEM-H→MitoEM-R (row 7 and 8). The pixels in red and green denote the false-negative and false-positive segmentation results respectively.

# <sup>1134</sup> E THEORY PROOFS

1141

1145

1151

1152

1154

1155

1159

1160

1161 1162

1136 E.1 COMPLEMENTARY MASKING THEORY: MEAN AND VARIANCE ANALYSIS

**1138 Definition 1: Complementary Mask** Let  $D \in \{0,1\}^d$  be a random binary vector where each 1139 element  $D_i$  is independently drawn from Bernoulli(0.5). The **complementary mask** is 1 - D, 1140 where 1 is the vector of ones in  $\mathbb{R}^d$ .

**Definition 2: Random Masks** Let  $D_1, D_2 \in \{0, 1\}^d$  be independent random binary vectors where each element  $D_{ki}$  (for k = 1, 2) is independently drawn from Bernoulli(0.5). These are the **random masks**.

1146 E.2 INFORMATION PRESERVATION METRIC

1148 Given a deterministic vector  $x \in \mathbb{R}^d$ , we define masked versions of x as:

1149 1150 - For complementary masks:

 $x_1 = D \odot x, \quad x_2 = (1 - D) \odot x$ 

- For random masks:

 $x_1 = D_1 \odot x, \quad x_2 = D_2 \odot x$ 

1156 where  $\odot$  denotes element-wise (Hadamard) product.

1158 Define the information preservation (IP) metric as:

$\operatorname{IP}(x_1, x_2) =$	_	$\langle x_1, x_2 \rangle$
	_	$  x  ^2$

1163 E.3 MEAN AND VARIANCE COMPUTATIONS

1164 1165 E.3.1 Complementary Masks

1166 Mean: 1167

1168 For complementary masks, note that for each coordinate *i*:

 $D_i(1 - D_i) = 0$ 

 $\langle x_1, x_2 \rangle = \sum_{i=1}^{d} D_i (1 - D_i) x_i^2 = 0$ 

 $\mathbf{IP}(x_1, x_2) = \frac{0}{\|x\|^2} = 0$ 

 $\mathbb{E}[\mathrm{IP}(x_1, x_2)] = 0$ 

1171 because  $D_i$  is either 0 or 1.

<sup>1172</sup> Therefore, the inner product:

1174

1169

1170

1177

1179

1176

1178 Thus,

1180

1181 1182 and

1183 1184

1185 Variance:

1186 Since  $IP(x_1, x_2) = 0$  almost surely, 1187

 $Var(IP(x_1, x_2)) = 0$ 

1188 E.3.2 RANDOM MASKS

1190 Mean:

1193

1194

1197

1198 1199

1202

1205 1206

For random masks:

$$\langle x_1, x_2 \rangle = \sum_{i=1}^d D_{1i} D_{2i} x_i^2$$

Since  $D_{1i}$ ,  $D_{2i}$  are independent Bernoulli(0.5), we have:

# $\mathbb{E}[D_{1i}D_{2i}] = \left(\frac{1}{2}\right)\left(\frac{1}{2}\right) = \frac{1}{4}$

1200 Therefore,

$\mathbb{E}[\langle x_1, x_2 \rangle] = \frac{1}{4} \ x\ ^2$
$\mathbb{E}[\mathrm{IP}(x_1, x_2)] = \frac{1}{4}$

1203 1204 and

# 1207 Variance:

1208 Compute  $Var(D_{1i}D_{2i})$ : 1209

1210 1211

1212

1213

1214 1215

$$\operatorname{Var}(D_{1i}D_{2i}) = \frac{1}{4} - \left(\frac{1}{4}\right)^2 = \frac{3}{16}$$

Then,

$$\operatorname{Var}(\langle x_1, x_2 \rangle) = \sum_{i=1}^d \frac{3}{16} x_i^4 = \frac{3}{16} \sum_{i=1}^d x_i^4$$

 $\operatorname{Var}(\operatorname{IP}(x_1, x_2)) = \frac{3}{16} \frac{\sum_{i=1}^d x_i^4}{(\|x\|^2)^2}$ 

1216 1217 Thus,

1218 1219

Theorem 4 (Consistency Bound for Feature Learning). Consider a general feature learning frame work with the objective function:

 $\mathcal{L}(f) = \mathbb{E}_x \left[ \ell \left( f(x_1), f(x_2) \right) \right],$ 

1234 1235

1239

where  $f : \mathbb{R}^d \to \mathbb{R}^k$  is the feature extraction function,  $\ell : \mathbb{R}^k \times \mathbb{R}^k \to \mathbb{R}$  is the loss function, and ( $x_1, x_2$ ) is a sample pair generated from input data x after applying masks or transformations.

1238 Assume:

1240 (a) The loss function  $\ell$  is L-Lipschitz continuous with respect to both arguments, i.e., for any 1241  $a, b, c, d \in \mathbb{R}^k$ ,

$$|\ell(a,b) - \ell(c,d)| \le L \left( \|a - c\|_2 + \|b - d\|_2 \right).$$

(c) The input data x takes values in a compact subset  $\mathcal{X} \subset \mathbb{R}^d$ , and  $\sup_{x \in \mathcal{X}} ||x||_2 \leq B$ .

(b) The feature extraction function f is  $\beta$ -Lipschitz continuous (or  $\beta$ -smooth), i.e., for any

 $||f(x) - f(y)||_2 \le \beta ||x - y||_2.$ 

Then, for any  $\delta \in (0, 1)$ , with probability at least  $1 - \delta$ , the following holds:

#### (i) For complementary masks:

$$|\mathcal{L}(f) - \hat{\mathcal{L}}_n(f)| \le 4L\beta B\left(\sqrt{\frac{2}{n}} + \sqrt{\frac{\log(2/\delta)}{n}}\right),$$

where  $\hat{\mathcal{L}}_n(f) = \frac{1}{n} \sum_{i=1}^n \ell(f(x_{1i}), f(x_{2i}))$  is the empirical risk computed on n samples.

# (ii) For random masks:

 $x, y \in \mathbb{R}^d$ ,

$$|\mathcal{L}(f) - \hat{\mathcal{L}}_n(f)| \le 4L\beta B\left(\sqrt{\frac{2}{n}} + \sqrt{\frac{\log(2/\delta)}{n}}\right) + 2L\beta B\sqrt{\frac{d}{n}}.$$

*Proof.* We will prove the bounds for both complementary masks and random masks separately.

#### 1262 Case (i): Complementary Masks

# 1263 Step 1: Define the Function Class

1265 Let  $\mathcal{F} = \{x \mapsto \ell(f(Dx), f((I-D)x)) : f \text{ is } \beta\text{-Lipschitz}\}$ , where D is a deterministic mask operator (for complementary masks).

#### 1267 Step 2: Bounding the Rademacher Complexity

Consider the empirical Rademacher complexity of  $\mathcal{F}$ :

$$\hat{\mathfrak{R}}_n(\mathcal{F}) = \mathbb{E}_{\boldsymbol{\sigma}}\left[\sup_{f\in\mathcal{F}}\frac{1}{n}\sum_{i=1}^n \sigma_i \ell\left(f(Dx_i), f((I-D)x_i)\right)\right],$$

where  $\sigma = (\sigma_1, \dots, \sigma_n)$  are independent Rademacher random variables (i.e.,  $\mathbb{P}(\sigma_i = +1) = \mathbb{P}(\sigma_i = -1) = 1/2$ ).

<sup>1275</sup> Using the Lipschitz property of  $\ell$  and f, we have:

$$\hat{\mathfrak{R}}_{n}(\mathcal{F}) \leq L\mathbb{E}_{\boldsymbol{\sigma}} \left[ \sup_{f} \frac{1}{n} \sum_{i=1}^{n} \sigma_{i} \left( \|f(Dx_{i}) - f(0)\|_{2} + \|f((I-D)x_{i}) - f(0)\|_{2} \right) \right]$$
$$\leq L\mathbb{E}_{\boldsymbol{\sigma}} \left[ \frac{1}{n} \sum_{i=1}^{n} |\sigma_{i}| \left( \|f(Dx_{i}) - f(0)\|_{2} + \|f((I-D)x_{i}) - f(0)\|_{2} \right) \right]$$

1282  
1283  
1284  
1284  
1285  

$$\leq 2L\beta \mathbb{E}_{\sigma} \left[ \frac{1}{n} \sum_{i=1}^{n} |\sigma_i| \|x_i\|_2 \right]$$

1285  
1286 
$$= 2L\beta \frac{1}{n} \sum_{i=1}^{n} ||x_i||_2 \mathbb{E}_{\sigma_i}[|\sigma_i|]$$
1287

1288  
1289 
$$= 2L\beta \frac{1}{n} \sum_{i=1} \|x_i\|_2 \cdot \mathbb{E}_{\sigma_i}[1]$$
1290

 $= 2L\beta \frac{1}{2} \sum_{i=1}^{n} ||x_i||_2$ 

$$1292 \qquad \qquad -2Lp n \sum_{i=1}^{n-1}$$

 $\begin{array}{ll} 1293 \\ 1294 \end{array} \leq 2L\beta B, \end{array}$ 

since  $||x_i||_2 \leq B$ . However, to get a dependence on n, we consider the Rademacher complexity bound for Lipschitz functions, which gives:

$$\hat{\mathfrak{R}}_n(\mathcal{F}) \leq rac{2Leta B}{\sqrt{n}}$$

# 1300 Step 3: Apply Concentration Inequalities

By McDiarmid's inequality, since changing one sample affects the empirical loss by at most  $\frac{2L\beta B}{n}$ , we have for any t > 0:

$$\mathbb{P}\left(|\mathcal{L}(f) - \hat{\mathcal{L}}_n(f)| \ge \mathbb{E}\left[|\mathcal{L}(f) - \hat{\mathcal{L}}_n(f)|\right] + t\right) \le 2\exp\left(-\frac{2nt^2}{(2L\beta B)^2}\right).$$

Setting  $t = L\beta B \sqrt{\frac{2\log(2/\delta)}{n}}$ , we get with probability at least  $1 - \delta$ :

$$|\mathcal{L}(f) - \hat{\mathcal{L}}_n(f)| \le \mathbb{E}\left[|\mathcal{L}(f) - \hat{\mathcal{L}}_n(f)|\right] + L\beta B\sqrt{\frac{2\log(2/\delta)}{n}}$$

# 13121313Step 4: Combine the Bounds

<sup>1314</sup> Using symmetrization and the bound on  $\hat{\mathfrak{R}}_n(\mathcal{F})$ , we have:

$$\mathbb{E}\left[\left|\mathcal{L}(f) - \hat{\mathcal{L}}_n(f)\right|\right] \le 2\hat{\Re}_n(\mathcal{F}) \le \frac{4L\beta B}{\sqrt{n}}.$$

1319 Therefore, combining the above, we have:

$$|\mathcal{L}(f) - \hat{\mathcal{L}}_n(f)| \le \frac{4L\beta B}{\sqrt{n}} + L\beta B\sqrt{\frac{2\log(2/\delta)}{n}} = 4L\beta B\left(\sqrt{\frac{1}{n}} + \sqrt{\frac{\log(2/\delta)}{n}}\right).$$

#### 1325 Case (ii): Random Masks

# 1326 Step 1: Modify the Function Class

1327 Let  $\mathcal{F}_{rand} = \{x \mapsto \ell(f(R_1x), f(R_2x)) : f \text{ is } \beta\text{-Lipschitz}, R_1, R_2 \text{ are random masks}\}.$ 

#### 1329 Step 2: Bounding the Rademacher Complexity

1330 Similarly, we consider:

$$\hat{\mathfrak{R}}_{n}(\mathcal{F}_{\mathrm{rand}}) = \mathbb{E}_{\boldsymbol{\sigma},R_{1},R_{2}}\left[\sup_{f} \frac{1}{n} \sum_{i=1}^{n} \sigma_{i} \ell\left(f(R_{1i}x_{i}), f(R_{2i}x_{i})\right)\right].$$

Again, using Lipschitz properties, we have:

$$\hat{\mathfrak{R}}_{n}(\mathcal{F}_{\text{rand}}) \leq L\mathbb{E}_{\boldsymbol{\sigma},R_{1},R_{2}}\left[\sup_{f} \frac{1}{n} \sum_{i=1}^{n} \sigma_{i} \left(\|f(R_{1i}x_{i}) - f(0)\|_{2} + \|f(R_{2i}x_{i}) - f(0)\|_{2}\right)\right]$$

1341 Since f is  $\beta$ -Lipschitz and  $||x_i||_2 \leq B$ , we have:

$$||f(R_{1i}x_i) - f(0)||_2 \le \beta ||R_{1i}x_i - 0||_2.$$

Given that  $R_{1i}$  is a random mask (e.g., a diagonal matrix with entries being Bernoulli random variables), we have:

1348  
1349 
$$\mathbb{E}_{R_{1i}}\left[\|R_{1i}x_i\|_2^2\right] = \sum_{j=1}^d \mathbb{E}[(R_{1i})_{jj}^2]x_{ij}^2 = \frac{d}{d}\|x_i\|_2^2 = \|x_i\|_2^2,$$

assuming each  $(R_{1i})_{jj}$  is independent and takes value 1 with probability 1/d.

1352 Therefore,

$$\mathbb{E}_{R_{1i}}\left[\|f(R_{1i}x_i) - f(0)\|_2\right] \le \beta \mathbb{E}_{R_{1i}}\left[\|R_{1i}x_i\|_2\right] \le \beta \sqrt{\mathbb{E}_{R_{1i}}\left[\|R_{1i}x_i\|_2^2\right]} \le \beta \frac{B}{\sqrt{d}}$$

1356 Similarly for  $R_{2i}$ .

1358 Therefore,

# 

# Step 3: Apply Concentration Inequalities

Following similar steps as in the complementary masks case, and accounting for the extra term due to random masks, we have:

 $\hat{\mathfrak{R}}_n(\mathcal{F}_{\mathrm{rand}}) \leq 2L\beta \frac{B}{\sqrt{d}}.$ 

$$|\mathcal{L}(f) - \hat{\mathcal{L}}_n(f)| \le 4L\beta B\sqrt{\frac{2}{n}} + L\beta B\sqrt{\frac{2\log(2/\delta)}{n}} + 2L\beta B\frac{1}{\sqrt{d}}.$$

Since  $\frac{1}{\sqrt{d}} \le \sqrt{\frac{d}{n}}$  for  $d \le n$ , we can write:

$$|\mathcal{L}(f) - \hat{\mathcal{L}}_n(f)| \le 4L\beta B\left(\sqrt{\frac{2}{n}} + \sqrt{\frac{\log(2/\delta)}{n}}\right) + 2L\beta B\sqrt{\frac{d}{n}}.$$

1376 This completes the proof.

**Theorem 5** (Signal Recovery Guarantee). Let  $x \in \mathbb{R}^d$  be a signal generated from the sparse linear model:

 $x = Mz + \xi,$ 

where:

•  $M \in \mathbb{R}^{d \times n}$  is a known measurement matrix (dictionary),

•  $z \in \mathbb{R}^n$  is a k-sparse vector (i.e.,  $||z||_0 \le k$ ),

•  $\xi \sim \mathcal{N}(0, \sigma^2 I_d)$  is additive Gaussian noise.

Suppose we have two masking matrices  $R_1, R_2 \in \mathbb{R}^{m \times d}$  representing partial observations of x:

- For complementary masks,  $R_1$  and  $R_2$  satisfy  $R_1R_2^{\top} = 0$  and  $R_1^{\top}R_1 + R_2^{\top}R_2 = I_d$ , i.e., they partition the indices of x without overlap and cover all entries.
- For random masks,  $R_1$  and  $R_2$  select entries independently at random.

Define the aggregated observation  $y \in \mathbb{R}^{2m}$  as:

$$y = \begin{pmatrix} y_1 \\ y_2 \end{pmatrix} = \begin{pmatrix} R_1 x \\ R_2 x \end{pmatrix} = \begin{pmatrix} R_1 M \\ R_2 M \end{pmatrix} z + \begin{pmatrix} R_1 \xi \\ R_2 \xi \end{pmatrix} = Az + \eta_z$$

where  $A \in \mathbb{R}^{2m \times n}$  is the effective measurement matrix, and  $\eta \in \mathbb{R}^{2m}$  is the aggregated noise.

Assume that A satisfies the Restricted Isometry Property (RIP) of order 2k with constant  $\delta_{2k} < \delta^*$ for some  $\delta^* < 1$ .

1402 Let  $\hat{z}$  be the solution to the basis pursuit denoising problem: 

 $\hat{z} = \arg\min_{u \in \mathbb{R}^n} \|u\|_1$  subject to  $\|y - Au\|_2 \le \epsilon$ ,

where  $\epsilon \geq \|\eta\|_2$ . Then, for any  $\delta > 0$ , with probability at least  $1 - \delta$ ,  $\|\hat{z} - z\|_2 \le C\sigma \sqrt{\frac{k\log(n/\delta)}{m}},$ where C > 0 is a constant depending only on the RIP constant  $\delta_{2k}$ . Moreover, when  $R_1$  and  $R_2$  are complementary masks that together cover all entries of x without overlap, and m = d/2, the recovery error achieves the bound:  $\|\hat{z} - z\|_2 \le C_1 \sigma \sqrt{\frac{k \log(n/\delta)}{d}},$ where  $C_1 > 0$  is a constant depending only on  $\delta_{2k}$ . *Proof.* We will establish an upper bound on the estimation error  $\|\hat{z} - z\|_2$  under the given assumption tions. Step 1: Formulating the Observations The observations are:  $y_1 = R_1 x = R_1 (Mz + \xi) = R_1 Mz + R_1 \xi,$  $y_2 = R_2 x = R_2 (Mz + \xi) = R_2 Mz + R_2 \xi.$ By stacking  $y_1$  and  $y_2$ , we have:  $y = \begin{pmatrix} y_1 \\ y_2 \end{pmatrix} = \begin{pmatrix} R_1 M \\ R_2 M \end{pmatrix} z + \begin{pmatrix} R_1 \xi \\ R_2 \xi \end{pmatrix} = Az + \eta,$ where  $A = \begin{pmatrix} R_1 M \\ R_2 M \end{pmatrix} \in \mathbb{R}^{2m \times n}$  and  $\eta = \begin{pmatrix} R_1 \xi \\ R_2 \xi \end{pmatrix} \in \mathbb{R}^{2m}$ . Step 2: Recovering z via Basis Pursuit Denoising We consider the optimization problem:  $\hat{z} = \arg\min_{u \in \mathbb{R}^n} \|u\|_1$  subject to  $\|y - Au\|_2 \le \epsilon$ , with  $\epsilon \geq \|\eta\|_2$ . Our goal is to bound  $\|\hat{z} - z\|_2$ . Step 3: Applying Compressed Sensing Recovery Guarantees Since A satisfies the RIP of order 2k with constant  $\delta_{2k} < \delta^*$ , standard compressed sensing results (e.g., Candès et al. (2006)) imply that:  $\|\hat{z} - z\|_2 \le C_0 \frac{\|\eta\|_2}{\sqrt{m}},$ where  $C_0 > 0$  depends only on  $\delta_{2k}$ . **Step 4: Bounding**  $\|\eta\|_2$ The noise vector  $\eta$  consists of 2m components, each being either  $\xi_i$  or zero. Since  $\xi \sim \mathcal{N}(0, \sigma^2 I_d)$ , each nonzero entry of  $\eta$  is  $\mathcal{N}(0, \sigma^2)$ . Therefore,  $\|\eta\|_2^2$  is the sum of 2m independent  $\sigma^2 \chi_1^2$  random variables, where  $\chi_1^2$  denotes a chi-squared distribution with one degree of freedom. Using concentration inequalities for chi-squared distributions (see, e.g., Laurent & Massart (2000)), for any t > 0:  $\Pr\left(\|\eta\|_2^2 \ge 2m\sigma^2(1+2\sqrt{t/(2m)}+2t/(2m))\right) \le e^{-t}.$ 

1458 Setting  $t = m \log(n/\delta)$ , we obtain:

$$\Pr\left(\|\eta\|_2^2 \ge 2m\sigma^2\left(1 + 2\sqrt{\frac{\log(n/\delta)}{m}} + \frac{2\log(n/\delta)}{m}\right)\right) \le \left(\frac{\delta}{n}\right)$$

1461 1462 1463

1464

1465 1466

1471 1472 1473

1478 1479

1482 1483

1460

For sufficiently large m, the terms involving 1/m become negligible, and we have, with probability at least  $1 - \delta$ :

m

$$\|\eta\|_2 \le C_1 \sigma \sqrt{m \log\left(\frac{n}{\delta}\right)},$$

where  $C_1 > 0$  is a constant.

# 1468 Step 5: Final Estimation Error Bound

1470 Substituting the bound on  $\|\eta\|_2$  into the recovery guarantee:

$$\|\hat{z} - z\|_2 \le C_0 \frac{C_1 \sigma \sqrt{m \log(n/\delta)}}{\sqrt{m}} = C \sigma \sqrt{\log\left(\frac{n}{\delta}\right)},$$

1474 where  $C = C_0 C_1$ .

To incorporate the sparsity k, we consider the number of possible supports of size k, which is  $\binom{n}{k}$ . Applying a union bound over all supports, we have:

$$\Pr\left(\|\hat{z} - z\|_2 \le C\sigma\sqrt{\log\left(\frac{n}{\delta}\right)}\right) \ge 1 - \delta.$$

1480 Noting that  $\log {n \choose k} \le k \log(n/k)$ , we refine the bound:

$$\|\hat{z} - z\|_2 \le C\sigma \sqrt{k \log\left(\frac{n}{k\delta}\right)} \le C'\sigma \sqrt{\frac{k \log(n/\delta)}{m}}$$

1484 1485 where C' > 0 is a constant.

#### 1486 Step 6: Special Case with Complementary Masks

When  $R_1$  and  $R_2$  are complementary and m = d/2, substituting m = d/2 yields:

$$\|\hat{z} - z\|_2 \le C' \sigma \sqrt{\frac{2k \log(n/\delta)}{d}} = C_1 \sigma \sqrt{\frac{k \log(n/\delta)}{d}}.$$

1490 1491 1492

1487

1488 1489

1493 Remark 3 (Advantages of Complementary Masks). Complementary masks offer significant advan-1494 tages in compressive sensing applications, enhancing both the theoretical foundations and practical 1495 implementations. These masks maximize measurement utilization by covering all entries of the signal x without overlap, ensuring optimal use of available information. This comprehensive coverage 1496 leads to improved Restricted Isometry Property (RIP) constants for the measurement matrix A, re-1497 sulting in tighter recovery bounds. The non-overlapping nature of complementary masks also plays 1498 a crucial role in minimizing noise influence, as it prevents noise accumulation and effectively reduces 1499  $\|\eta\|_2$ . A key benefit is the improved recovery accuracy, where the error bound scales inversely with 1500 the dimensionality d of x, leading to enhanced recovery performance. Furthermore, the structured 1501 nature of these masks contributes to algorithmic efficiency, facilitating faster and more effective 1502 computation in practical recovery algorithms. Collectively, these properties make complementary 1503 masks a powerful tool in compressive sensing, offering a balanced approach that enhances both 1504 theoretical guarantees and practical performance. 1505

#### 1506 1507 REFERENCES

Candès, E. J., Romberg, J., & Tao, T. (2006). Stable signal recovery from incomplete and inaccurate measurements. *Communications on Pure and Applied Mathematics*, 59(8), 1207–1223.

Laurent, B., & Massart, P. (2000). Adaptive estimation of a quadratic functional by model selection.
 *Annals of Statistics*, 28(5), 1302–1338.

# <sup>1512</sup> F APPLICATIONS AND EXTENSIONS

# 1514 F.1 SELF-SUPERVISED LEARNING

The complementary masking theory can be directly applied to self-supervised learning tasks, particularly in contrastive learning frameworks. Here, we present a corollary that demonstrates how our theory can be used to analyze the performance of contrastive learning algorithms.

**Corollary 6** (Contrastive Learning with Complementary Masks). *Consider a contrastive learning* setup where positive pairs are generated using complementary masks (D, I - D). Let  $f_{\theta} : \mathbb{R}^d \to \mathbb{R}^k$ be the encoder network parameterized by  $\theta$ , and let the contrastive loss be defined as:

1522 1523 1524

1525

1535

$$\mathcal{L}(\theta) = -\mathbb{E}_x \left[ \log \frac{e^{sim(f_\theta(Dx), f_\theta((I-D)x))/\tau}}{\sum_{j=1}^N e^{sim(f_\theta(Dx), f_\theta((I-D)x_j))/\tau}} \right]$$

where  $sim(\cdot, \cdot)$  is the cosine similarity and  $\tau$  is a temperature parameter. Then, under the assumptions of Theorem 2, with probability at least  $1 - \delta$ :

$$|\mathcal{L}(\theta) - \hat{\mathcal{L}}_n(\theta)| \le O\left(\frac{L\beta B}{\tau}\left(\sqrt{\frac{1}{n}} + \sqrt{\frac{\log(1/\delta)}{n}}\right)\right)$$

where  $\hat{\mathcal{L}}_n(\theta)$  is the empirical loss on *n* samples, *L* is the Lipschitz constant of the loss function,  $\beta$  is the smoothness parameter of  $f_{\theta}$ , and *B* is the bound on the input norm.

**Proof.** The proof follows directly from Theorem 2 by observing that the contrastive loss is Lipschitz continuous with respect to the encoder outputs, and the encoder network is assumed to be  $\beta$ -smooth. The key step is to apply the consistency bound for complementary masks to the positive pair (Dx, (I - D)x) in the numerator of the contrastive loss.

This corollary provides a theoretical justification for using complementary masks in contrastive learning algorithms. It suggests that the generalization error of such algorithms scales favorably with the number of samples and is independent of the input dimension, which is crucial for highdimensional data such as images.

# 1545 F.2 EXTENSION TO MULTI-VIEW DATA

The complementary masking theory can be extended to scenarios where we have multiple views of
the data, not just two. This extension is particularly relevant for multi-view learning and multi-modal
data analysis.

**Theorem 7** (Multi-View Complementary Masking). Let  $x \in \mathbb{R}^d$  be the input data, and consider K complementary masks  $D_1, \ldots, D_K$  such that  $\sum_{i=1}^K D_i = I$ . Define the multi-view information preservation metric as:

1554

1557

$$MIP(x_1,\ldots,x_K) = \frac{1}{K(K-1)} \sum_{i \neq j} \frac{\langle x_i, x_j \rangle}{\|x\|^2}$$

where  $x_i = D_i x$ . Then:

559 *1.* 
$$\mathbb{E}[MIP(x_1, \dots, x_K)] = \frac{1}{K^2}$$

1560  
1561 2. 
$$Var(MIP(x_1, \dots, x_K)) \le \frac{K-1}{K^3} \frac{\sum_{i=1}^d x_i^4}{\|x\|^4}$$

1562

**1563** *Proof.* (Sketch) The proof follows a similar structure to that of Theorem 1, but requires careful accounting of the pairwise interactions between the K views. The key insight is that the complementary nature of the masks ensures that the expected overlap between any two views is  $1/K^2$  of the total information. This multi-view extension opens up possibilities for analyzing and designing algorithms that work
 with more than two views of the data, such as multi-view clustering or multi-modal fusion techniques.

# 1570 G CONCLUSION

The complementary masking theory presented in this paper provides a rigorous framework for analyzing information preservation in masked data representations. The key advantages of complementary masks over random masks include:

1576 1. Tighter generalization bounds in feature learning tasks.

<sup>1577</sup> 2. More robust signal recovery guarantees, especially in the presence of strong signals.

15781579 3. Guaranteed preservation of a constant fraction of the original information.

These theoretical results have immediate implications for the design and analysis of self-supervised
 learning algorithms, particularly in contrastive learning setups. They also provide insights into why
 certain masking strategies might outperform others in practice.