

MCAE: Masked Contrastive Autoencoder for Face Anti-Spoofing

Tianyi Zheng
Shanghai Jiao Tong University
tyzheng@sjtu.edu.cn

Abstract

Face anti-spoofing (FAS) method performs well under the intra-domain setups. But cross-domain performance of the model is not satisfying. Domain generalization method has been used to align the feature from different domain extracted by convolutional neural network (CNN) backbone. However, the improvement is limited. Recently, the Vision Transformer (ViT) model has performed well on various visual tasks. But ViT model relies heavily on pre-training of large-scale dataset, which cannot be satisfied by existing FAS datasets. In this paper, taking the FAS task as an example, we propose Masked Contrastive Autoencoder (MCAE) method to solve this problem using only limited data. Meanwhile in order for a feature extractor to extract common features in live samples from different domains, we combine Masked Image Model (MIM) with supervised contrastive learning to train our model. Some intriguing design principles are summarized for performing MIM pre-training for downstream tasks. We also provide insightful analysis for our method from an information theory perspective. Experimental results show our approach has good performance on extensive public datasets and outperforms the state-of-the-art methods.

1. Introduction

Face recognition (FR) techniques offers a simple yet convenient way for identity authentication applications, such as mobile access control and electronic payments. Though face biometric systems are widely used, with the emergence of various presentation attacks, critical concerns about security risk on face recognition systems are increasing. An unprotected face recognition system might be fooled by merely presenting artifacts like a photograph or video in front of the camera. Therefore, how to strengthen the face recognition system from a variety of presentation attacks promotes the techniques of face anti-spoofing (FAS).

As an important research topic, a series of face anti-spoofing (FAS) methods have been proposed, from hand-

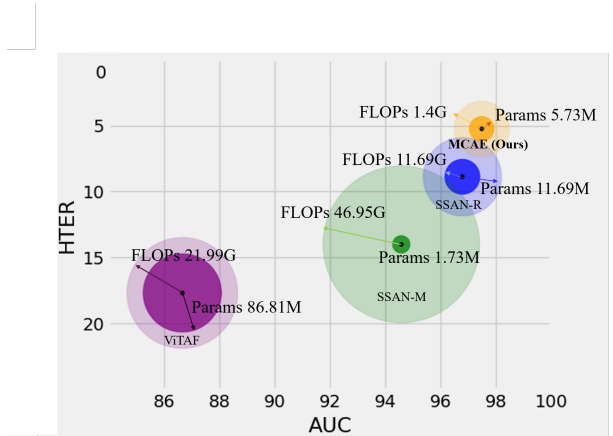


Figure 1. The parameters and FLOPs of different model. Dark circles indicate the number of parameters in the models, light circles indicate FLOPs of the model. ViTAF uses ViT-Base as their backbone; SSAN-M uses MADDFG; SSAN-R uses Resnet-18 and our method uses ViT-Tiny.

craft feature based to deep representation based methods. Although the previous FAS methods have achieved promising performance in intra-domain scenarios, they may suffer dramatic degradation when encounter unseen domains. To alleviate this issue, researchers have proposed various approaches [6, 21, 26, 30, 32, 38, 40] to improve the generalizability of FAS in cross-domain settings. However, these methods are almost based on common convolutional neural network (CNN), which lack taking advantage of subtle properties of global and local image statistics in FAS [40]. Recently, the modern architecture Vision Transformer (ViT) [11] has received increasing attention in the vision community. With their in-built local patchifying and global self-attention mechanisms, ViTs may be potentially better-suited to FAS over their CNN counterparts. Most recently, CNN equipped with attention modules [28] and sophisticated designed ViT variants [13, 18, 20, 24] have been introduced into FAS and obtained promising performance. However, whether a vanilla ViT without extra training samples from upstream tasks can achieve competitive cross-domain generalization has not been explored thus far.

Despite the advantages, directly training ViT as well as CNN with a binary classification model is prone to overfit to the biases introduced by the data come from different domains. The most discriminative spoofing cues can dramatically change or even disappear across various domains, making the learned features vulnerable in the unseen scenarios. To alleviate this problem, we propose to use the recently emerged Masked Image Modeling (MIM) to learn facial representation for FAS. MIM is mostly built upon the Vision Transformer (ViT), which suggests that self-supervised visual representations can be done by masking input image parts while requiring the target model to recover the missing contents. With the masked-reconstruct mechanisms, MIM pretraining forces the ViT to learn maximize the mutual information between local cues and global features, thus constraining the model from easily overfitting to local discriminative spoofing bias. However, it is non-trivial to apply MIM pretrain methods like MAE [16] to FAS. Previous works [11, 16] have demonstrated that the performance of ViT relies heavily on pre-training of large-scale datasets like JFT-300M (300 million images) [34] and ImageNet21K (14 million images) [9]. There is little work to study how to pre-train ViT on small data sets, especially for MIM based pre-training. To fill this gap, in this paper we take MAE as an example and investigate how this pre-training method perform on small FAS datasets. Specifically, some intriguing findings reveal that MIM pre-training on small data has distinctly different design principles than pre-training on large-scale natural images.

It is worth noting that, although MIM pretraining can prevent the model from over-fitting to local local discriminative spoofing bias, as the reconstruction quality continues to improve during pretraining, the network also inevitably learns domain-specific global features like color distortion, illuminations, and image resolutions. These redundant features can significantly reduce the cross-domain generalizability of the pre-trained model in downstream FAS tasks. To mitigate this situation and make it easier to transfer pre-trained features to downstream FAS tasks, we propose to incorporate contrastive learning into the MIM pre-training framework to capture more domain invariant liveness information and suppress domain-specific one. Following the above design principles, we propose a Masked Contrastive Autoencoder training framework for face anti-spoofing, called MACE. As shown in Figure 1, our work demonstrates that even with a lightweight vanilla ViT, superior cross-domain generalizability can be achieved in FAS tasks without using extra training data. We also provide insightful analysis on why the proposed MCAE outperforms the related methods. The main contributions of this work are:

- We introduce the modern architecture Vision Transformer into face anti-spoofing and propose a simple

yet effective framework for robust cross-domain FAS.

- We conducted a systematic study on how to apply the MIM pre-training on small datasets, by taking the FAS task as an example. By incorporating contrastive learning into the MIM pre-training, we propose a Masked Contrastive Autoencoder training framework for face anti-spoofing, called MACE. We also summarize empirical design principles for performing MIM pre-training for downstream cross-domain FAS, which are significantly different from the known experience with large-scale natural image pre-training.
- Extensive experiments on the widely-used benchmark datasets demonstrate the superior cross-domain generalizability of our proposed MCAE in FAS tasks even with a lightweight vanilla ViT.

2. Related Work

2.1. Face anti-spoofing (FAS).

The traditional features used are often hand-crafted features such as LBP [12], HOG [23] and SIFT [29]. Recently, CNN and Vision Transformer [11] are used as the backbone to extract features to distinguish the spoof samples from live samples. In order to get a model that has good performance on cross-domain datasets. Domain adaptation (DA) and domain generation (DG) method are widely used in FAS. In those methods, the algorithms are designed to pull the live samples close in feature space. PatchNet [37] combine the information in capturing device and presenting materials to improve the generality of their model. The method SSAN [40] is designed to split the representation into content and style ones with different supervision. ViTAF method [19] uses adaptive transformers as backbone and has good performance on cross-domain task. But this model requires large additional datasets for supervised pre-training. How to get a model with good generalization performance using limited datasets is still a problem.

2.2. Masked Image Model (MIM).

Recently, the MIM become popular in self-supervised vision tasks thanks to the introduction of ViT model. Masked prediction was inspired by the success of Masked Language Model (MLM) like BERT [10]. Training such models requires two stages namely pre-training and fine-tuning. During the pre-train stage, the MIM mask some patches of the original image and predict the features of those masked patches, such as pixel feature [16], discrete token [2] and HOG feature [41]. Then in the fine-tuning stage the pre-trained model are used as the feature extractor to downstream tasks. Because the model has learned a good representation of the relative images, the feature extracted

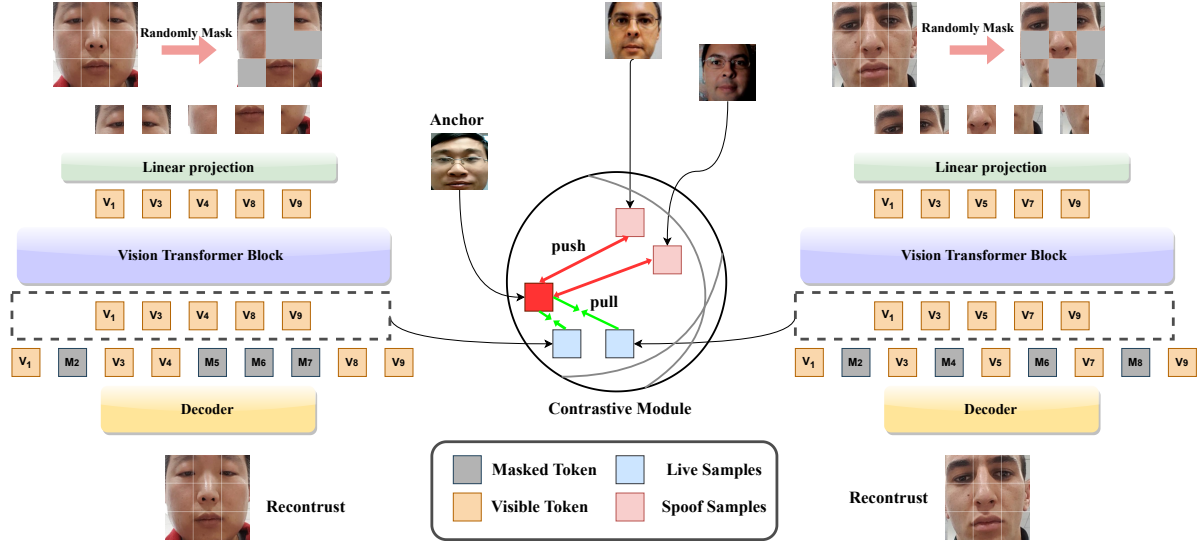


Figure 2. Overview of our MCAE method. We randomly choose some tokens and send them to the encoder. The output of the encoder are used to do reconstruction task with masked tokens and get the aggregate token. The blue square represent features of live sample and the pink square represent features of spoof sample. The ViT blocks have the same weights.

by it can be very useful to downstream tasks. But it should be noticed that the great success of MIM is based on large dataset for pre-training.

2.3. Contrastive learning.

The contrastive learning method in self-supervised learning are based on pretext tasks. The key idea of contrastive learning method is pull the positive samples close and meanwhile push away the negative samples in feature space. There are different ways to construct positive and negative samples. For example, the SimCLR [5] considers the same image with different data augmentation as positive samples, and all of other images as negative samples. The BYOL [14] even does not need to construct negative samples anymore. It uses an online encoder to predict the output of another momentum encoder MoCo [17].

The above methods are self-supervised and no label information is need. In order to find the similarity between samples with the same label, the SupContrast [22] and WACL [4] use all the samples with the same label to compute contrastive loss during pre-training stage. The method proposed by [15] combines contrastive loss with classification loss in the fine-tuning stage. However, none of the above methods combine the contrast task with reconstruction task during the pre-training stage.

3. Proposed Method

3.1. Overview

An overview of our method is show in Figure 2. Our method consists of two stages, including pre-training stage

and fine-tuning stage. We describe our method in Algorithm 1 of the pre-training stage. In the fine-tuning stage, we only keep the encoder to extract features for FAS tasks while discarding the decoder.

3.2. Masked Autoencoders

MIM is a simple but effective way to train a general feature extractor for different vision tasks such as classification, object detection and semantic segmentation.

In the pre-training stage, we use ViT as our encoder. First of all, we divided the input image into non-overlapping patches and then projected each of them into tokens T_i so that we can get a token sequence $\{T_i\}$ where $i = 1, \dots, n$ and n is the sequence length. Then we randomly choose a subset of the token sequence $\{T_{v_i}\}$ to keep, and the other tokens will be masked and denoted as $\{T_{m_i}\}$. The $\{T_{v_i}\}$ will be fed into encoder G_θ to get the latent representation of the original image. we use aggregation to represent the image's features, denoted as T_{agg} . The decoder D_θ combine $\{T_{v_i}\}$ with $\{T_{m_i}\}$ to reconstruct the original input image. The indicate function $\mathbb{1}_{mask}(i)$ indicates whether token T_i is masked is defined as Equation (1):

$$\mathbb{1}_{mask}(i) = \begin{cases} 1, & i \in T_m \\ 0, & i \notin T_m \end{cases}. \quad (1)$$

The loss function \mathcal{L}_{rec} are defined as Equation (2):

$$\mathcal{L}_{rec} = \frac{1}{n} \sum_{i=1}^n \|D_\theta(G_\theta(T_v), T_m) - T_i\|_2^2 \mathbb{1}_{mask}(i). \quad (2)$$

During constructing MIM pre-training for small FAS datasets, we find some intriguing design principles which are different from previous experience [16] with MIM pre-training on large-scale natural datasets. Specifically, unlike a lightweight decoder used in MAE [16] for large-scale natural datasets pre-training, the decoder size matters for the representation of FAS tasks. We also find that compare to nature images a higher proportion of the input image (85%), yields a better self-supervisory task for FAS data. Moreover, introducing some high level semantic information can help the self-supervised model better benefit the downstream task with small size datasets. In the following, we present a practice of introducing supervised contrast learning in MIM pre-training to capture more domain invariant and task-related information for FAS tasks. We hope that these found design principles can inspire researchers to design better MIM pre-training strategies for a variety of downstream tasks, especially on small data sets.

3.3. Information Theoretic Analysis

In this subsection, we give an information theoretic analysis of our method and explain why it's useful to combine reconstruction task with contrastive task. Intuitively, we want to reconstruct all of the T_i of the original image, which is equal to maximize the mutual information between the $G_\theta(T_v)$ and T_i . Therefore, we will find the correlation between the reconstruction task and contrastive module.

Based on the definition of the mutual information [8], the mutual information between the T_v and T_i is given in Equation (3):

$$\begin{aligned} \mathcal{I}(T_i; G_\theta(T_v)) &= H(T_i) - H(T_i | G_\theta(T_v)) \\ &= H(G_\theta(T_v)) - H(G_\theta(T_v) | T_i). \end{aligned} \quad (3)$$

By the definition of the conditional entropy, we have:

$$H(T_i | G_\theta(T_v)) = \mathbb{E}_{P_{T_i, G_\theta(T_v)}} [-\log P(T_i | G_\theta(T_v))].$$

But in practice it's very difficult to get distribution of $P(T_i | G_\theta(T_v))$ directly. The most common way to approximate this distribution is using another distribution $Q(x)$ instead of it and maximize the lower bound of KL divergence between them [1]:

$$\begin{aligned} \mathcal{I}(G_\theta(T_v); T_i) &= H(T_i) - H(T_i | G_\theta(T_v)) \\ &= H(T_i) + \mathbb{E}_{P_{T_i, G_\theta(T_v)}} [\log P(T_i | G_\theta(T_v))] \\ &= H(T_i) + \mathbb{E}_{P_{T_i, G_\theta(T_v)}} [\log Q(T_i | G_\theta(T_v))] \\ &\quad + \underbrace{D_{KL}(P(T_i | G_\theta(T_v)) \| Q(T_i | G_\theta(T_v)))}_{\geq 0} \\ &\geq \mathbb{E}_{P_{T_i, G_\theta(T_v)}} [\log Q(T_i | G_\theta(T_v))]. \end{aligned} \quad (4)$$

The distribution $Q(x)$ can be chosen arbitrarily. So we can use Gaussian distribution with σI diagonal matrix as

$Q(x)$ [27] i.e. $Q(x) \sim \mathcal{N}(T_i | G_\theta(T_v), \sigma I)$. Therefore, the maximize problem can be convert to the minimal problem in Equation (5):

$$\min \mathbb{E}_{P_{T_i, T_v}} [\|D_\theta(G_\theta(T_v), T_m) - T_i\|_2^2]. \quad (5)$$

In the reconstruction task, we want to minimize the Equation (2). In fact Equation (5) is very similar to the Equation (2), the only difference between them is whether to compute loss value on unmasked tokens.

Next, we will proof that T_{agg} can be used as the representation of the original image. Based on the assumption proposed by [33], we know that the reconstruct image and input image are both redundant for the task-relevant information i.e. there exists an ϵ s.t. $\mathcal{I}(T_i; T_{agg} | D_\theta(G_\theta(T_{v_i}), T_m)) \leq \epsilon$.

Theorem 1. The self-supervised learned aggregate token contains all the task-relevant information [35] in the token sequence $\{T_i\}$ where $i = 1, \dots, n$ with a potential loss ϵ .

$$\mathcal{I}(T_i; G_\theta(T_{v_i})) - \epsilon \leq \mathcal{I}(T_{agg}; G_\theta(T_{v_i})) \leq \mathcal{I}(T_i; G_\theta(T_{v_i})).$$

The proofs are provided in the appendix. By Thm 1, we can get the result that T_{agg} can be used as the original image's task-relevant representation. The ideal situation is that T_{agg} is the sufficient statistic for estimating the $\{T_i\}$. In that case, $G_\theta(T_v)$, T_{agg} and $\{T_i\}$ form a Markov chain $\{T_i\} \leftrightarrow T_{agg} \leftrightarrow G_\theta(T_v)$ and T_{agg} can represent $\{T_i\}$ without any information loss.

Contrastive learning focuses on learning common features between instances of the same class and distinguishing differences between instances of different classes. Meanwhile, when we already have an encoder that can extract input features well, the next step is to find the common information of the samples from different domains with the same label. Since the reconstruction task can promote the model to learn a good task-relevant representation of the input image, we can add the contrastive module to maximize the mutual information between positive samples.

3.4. Supervised Contrastive Module

Different datasets of face images are collected by different capture devices in different scenes, and have different resolutions. Thus we want to find the common information about live samples without any additional features due to domain difference. However, the common contrastive learning methods rely on an extensive dictionary, which is difficult to achieve on small datasets. Meanwhile, the datasets of face anti-spoofing is different from the nature image dataset. Since the original datasets of the face are video, there are high similarities between adjacent frames in the video. Therefore, treating images between adjacent frames as negative samples is unreasonable.

According to the previous analysis, we already have an encoder that can extract task-relevant features of the input

image well. So we can aggregate all tokens directly as the representation of the input image. The contrastive module aims to narrow the live samples, which common method used in the previous is triplet loss. However, triplet loss not only narrows the distance between each live sample but also the distance between each spoof sample. We think it is unreasonable to narrow the distance between spoof samples because spoof face images come in many forms and do not have consistent features. Our aim is also to find domain-invariant features between live samples, so we use the domain label and live label as our supervised information.

Therefore, we design our supervised contrastive loss based on the previous work [4, 15, 22]. We special design weighting factors for different samples for the face anti-spoofing task. We have two types of positive samples; one is the positive samples between the same domain, and the other is the positive samples between different domains. We give more significant weight to positive samples from different domains. The supervised contrastive loss is defined in Equation (6):

$$L_{con} = - \mathbb{E} \left[\sum_{j=1}^N \mathbb{1}_{i \neq j} (1 - \mathbb{1}_{y_i \neq y_j}) \log \frac{\lambda_l \exp(s_{i,j}/\tau)}{\lambda_l \exp(s_{i,j}/\tau) + \sum_{k=1}^N \mathbb{1}_{y_i \neq y_k} \exp(s_{i,k}/\tau)} \right]. \quad (6)$$

In the above Equation (6), N is the mini-batch size; y_i and y_j mean the label of sample i and sample j ; τ is the temperature parameter; $\mathbb{1}_{y_i \neq y_j}$ is the indicator function of whether sample i and sample j have the same label; s_{ij} is the cosine similarity between sample i and sample j ; λ_l are the weighting factors. Every sample in a mini-batch is used as an anchor once.

It is worth noting that we should add our supervised contrastive module after the mutual information between T_{agg} and $G_\theta(T_v)$ tends to converge. Only in this case can we consider that the features of the original image are well extracted. If we add modules too early, the features used to contrast are meaningless, and we will illustrate this phenomenon in the next section.

4. Experiments

4.1. Experimental Setups

4.1.1 Experiment Datasets.

We evaluate our proposed methods on cross-dataset testing based on four public datasets CASIA-MFSD(C) [45], Replay-Attack(I) [7], MSU-MFSD(M) [42] and OULU-NPU(O) [3]. Since each datasets are sampled by various devices in different scenarios, there are large differences be-

Algorithm 1 MCAE Pre-training stage

Input: Face images from different domain; A relatively small hyperparameter ϵ ; Various weight parameters λ_l

Parameter: Encoder G_θ and decoder D_θ .

Output: Trained encoder G_θ .

- 1: Choose a subset T_v from all of the token T_i as visible token. The rest of T_i are masked which are denoted as T_m . The mask ratio is 85%.
 - 2: Input visible token T_v to encoder to get the feature representation $G_\theta(T_v)$.
 - 3: **for** $i = 1, \dots, N_{train}$ **do**
 - 4: Input T_v and T_m to decoder to get reconstruct image $D_\theta(G_\theta(T_v), T_m)$.
 - 5: **if** $|G_\theta(T_v) - T_m|^2 < \epsilon$ **then**
 - 6: Compute reconstruct loss \mathcal{L}_{rec} .
 - 7: **else**
 - 8: Aggregate feature token T_{agg} .
 - 9: Compute reconstruct loss \mathcal{L}_{rec} and use T_{agg} to compute supervised contrastive loss L_{con} .
 - 10: **end if**
 - 11: Update parameters θ of encoder G_θ and decoder D_θ .
 - 12: **end for**
 - 13: Keep encoder E_θ and discard decoder D_θ .
 - 14: **return** Trained Encoder G_θ
-

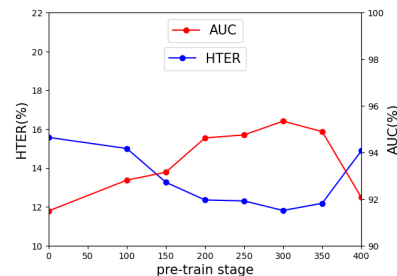


Figure 3. The procedure with the supervised contrastive module. When we add supervised contrastive module after the model are pre-trained a period of time, the performance of the model are effectively improved. The unit of abscissa is epoch.

tween those datasets. Experiments in such a setting can be a good evaluation of the generalization ability of the model.

4.1.2 Implementation Details.

We use MTCNN [44] to detect faces in each dataset, then crop and resize each face image into $256 \times 256 \times 3$. We use ViT-Tiny as our backbone whose embedding dimension is 192 and the patch size of each image is 16×16 . We only use random resized cropping as our data augmentation method. We use the same evaluation metric as previous work [31], i.e. the Half Total Error Rate (HTER) and the Area Under Curve (AUC).

Method	O&C&I to M		O&M&I to C		O&C&M to I		I&C&M to O	
	HTER(%)	AUC(%)	HTER(%)	AUC(%)	HTER(%)	AUC(%)	HTER(%)	AUC(%)
MADDG ([31])	17.69	88.06	24.50	84.51	22.19	84.99	27.98	80.02
NAS-FAS ([43])	16.85	90.42	15.21	92.64	11.63	96.98	13.16	94.18
D^2 AM ([6])	12.70	95.66	20.98	85.58	15.43	91.22	15.27	90.87
SSDG-R ([21])	7.38	97.17	10.44	95.94	11.71	96.59	15.61	91.54
ANRL ([25])	10.83	96.75	17.83	89.26	16.03	91.04	15.67	91.90
DRDG ([26])	12.43	95.81	12.43	95.81	19.05	88.79	15.63	91.75
SSAN-R ([40])	6.67	98.75	10.00	96.67	8.88	96.79	13.72	93.63
PatchNet ([37])	7.10	98.46	11.33	94.58	13.40	95.67	11.82	95.07
ViTAF [†] ([19])	4.75	98.79	15.70	92.76	17.68	86.66	16.46	90.37
MCAE (Ours)	3.81	99.10	10.00	96.71	5.25	97.49	11.81	95.34

Table 1. Comparison results between our MCAE method and state-of-the-art methods on cross-dataset testing. ViTAF[†] denote the ViT-Base model pre-trained by Imagenet dataset. Our MCAE method achieves the best performance in each setting.

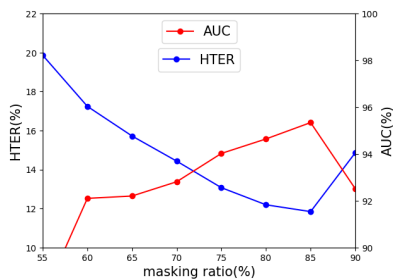


Figure 4. The influence of masking ratio. Because of the strong structural similarity of face images, the best masking ratio 85% which is larger than that of natural images.

4.2. Experiment Result

4.2.1 Cross-domain experiments.

Following previous works, we use Leave-One-Out (LOO) setting to do our cross-domain experiments. We pre-train and fine-tune our model on three datasets and test it on the rest one. The cross-domain result compared to other state-of-the-art method are as shown in Table 1. We observed that our CMAE method achieves the best performance. The result shows that our method has good generalization ability. Although the dataset size of FAS is limited, we still get a model with good generalization performance through pre-training and contrastive learning without any other datasets.

4.2.2 Experiments on Limited Source Domains.

To further evaluate the generalization ability of our model, we do cross-domain experiment based on limited source domain data. Following previous work, we use M and I as our source domain data to train our model. Then test our model in dataset C and O. The result are shown in Table 2, Our method achieves the best performance. The result proves that our method has good generalization ability.

Method	M&I to C		M&I to O	
	HTER(%)	AUC(%)	HTER(%)	AUC(%)
IDA ([42])	45.16	58.80	54.52	42.17
MADDG ([31])	41.02	64.33	39.35	65.10
SSDG-M ([21])	31.89	71.29	36.01	66.88
DR-MD-Net ([39])	31.67	75.23	34.02	72.65
ANRL ([25])	31.06	72.12	30.73	74.10
SSAN-M ([40])	30.00	76.20	29.44	76.62
MCAE (Ours)	29.89	77.65	21.32	87.35

Table 2. The result of the limited source domain experiments. Even though the dataset is limited, our method still get a good performance on each setting.

4.3. Effect Analyses of pre-training stage

In this subsection, we explore the impact of some essential modules in the pre-training stage because our model has some unique designs for the face anti-spoofing tasks. The face image has substantial structural similarity, which is very different from the nature images.

4.3.1 Pre-training schedules.

As shown in Figure 3, adding a supervised contrastive module after the model has been pre-trained at a specific time has good performance. The result is consistent with what we analyzed earlier. When the encoder does not converge, it is meaningless to use the features at this time as contrastive features. In addition, the contrastive loss also takes a particular training time to converge, so we cannot add this module too late.

4.3.2 Mask ratio.

The mask ratio determines the number of visible patches processed by the encoder in pre-training stage. Figure 4 show the result of experiment in different mask ratio. As we can see in Figure 4. The best mask ratio for face anti-spoofing tasks is 85% which is larger than the nature image.

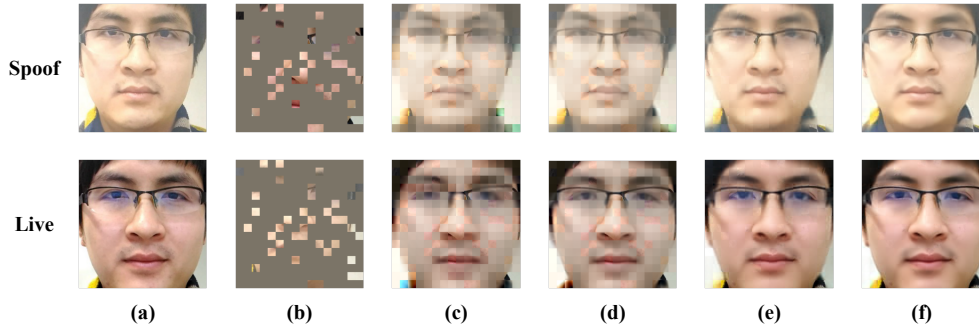


Figure 5. The reconstruct result: (a): Original image. (b): Masked image. (c): Reconstructed image with 1 layer decoder. (d): Reconstructed image with 2 layer decoder. (e): Reconstructed image with 4 layer decoder. (f): Reconstructed image with 8 layer decoder. The reconstructed image with 4 layer and 8 layer decoder distinguish between attack and real sample.

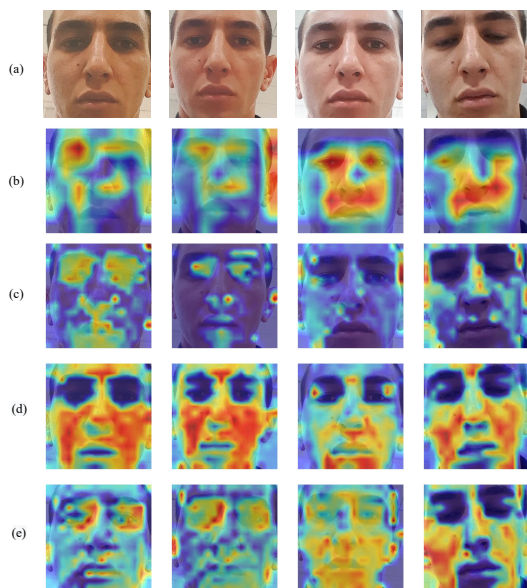


Figure 6. Attention visualization of different model. (a): Original images. (b): Resnet-18. (c): ViT without pretrain. (d): Pre-trained ViT without supervised contrastive module. (e): Pre-trained ViT with supervised contrastive module. The pre-trained ViT model pay more attention about the whole image, while other model care the local feature.

Decoder		Evaluation Metric	
width	depth	HTER(%)	AUC(%)
48	1	17.38	89.40
192	2	17.21	90.84
384	4	13.35	92.56
512	8	11.81	95.34
768	10	14.01	92.30

Table 3. Effect of the docoder size. The performance of model is sensitive to the decoder structure, the decoder with 8 layers has the best performance.

Because of the substantial structural similarity of the face images, masking more of the input images is a useful way to promote the encoder to learn a better representation of the face images.

4.3.3 Decoder size.

Decoder is a key module in pre-training stage. We find that the performance of model is sensitive to the decoder structure. Unlike using lightweight decoders in nature images, we design our decoder especially for FAS tasks. The results of experiment with different decoder size are shown in Table 3. Because in FAS tasks, the detail of the image is crucial to classifier. If the size of the decoder is too small, the encoder cannot learn a good representation of the masked image which may negatively affect the aggregate of the sample with the same label. Meanwhile, the decoder will be discarded in the fine-tuning stage, so we must limit the size of decoder to avoid the phenomenon that too many useful information are in decoder instead of encoder.

4.4. Ablation Studies

4.4.1 The effect of Pre-training.

Table 4 shows the importance of our special design for the pre-training stage. The results are very unsatisfying when we use the ViT model without pre-training to do face anti-spoofing tasks. Meanwhile, even though we use the Imagenet dataset to pre-train our model, the improvement is also limited. Since the images of the face are quite different from the natural image. The features extracted by the encoder trained with the Imagenet dataset may not be suitable for face anti-spoofing tasks.

4.4.2 Supervised contrastive module.

The results in Table 5 show that the contrastive module is beneficial in improving the model’s generalization ability.

Method	O&C&I to M		O&M&I to C		O&C&M to I		I&C&M to O	
	HTER(%)	AUC(%)	HTER(%)	AUC(%)	HTER(%)	AUC(%)	HTER(%)	AUC(%)
ViT-Tiny [†]	10.48	94.80	26.78	79.21	26.67	74.93	24.77	81.47
ViT-Tiny*	8.57	95.84	22.11	84.32	17.33	85.64	22.34	82.92
ViT-Tiny	3.81	99.10	10.00	96.71	5.25	97.49	11.81	95.34

Table 4. Ablation result for model pre-training. ViT-Tiny[†] denoted the ViT-Tiny model without pre-training. ViT-Tiny* denoted the ViT-Tiny model pre-trained on imagenet dataset. The results show that our training method are the best.

Method	O&C&I to M		O&M&I to C		O&C&M to I		I&C&M to O	
	HTER(%)	AUC(%)	HTER(%)	AUC(%)	HTER(%)	AUC(%)	HTER(%)	AUC(%)
w contrastive	3.81	99.10	10.00	96.71	5.25	97.49	11.81	95.34
w/o λ_k parameters	5.71	97.33	13.22	93.23	8.08	96.82	13.05	93.24
w/o contrastive	6.19	96.98	16.67	90.79	12.67	94.97	13.19	92.92

Table 5. Ablation result on supervised contrastive module. With the contrastive module, the performance is improved.

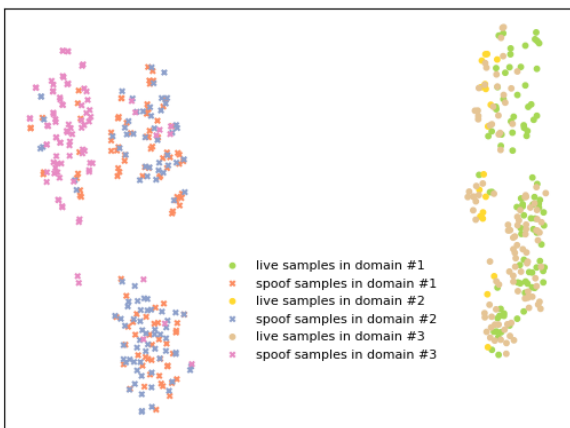


Figure 7. Visualization of the feature distribution of pre-trained MCAE model. Because of the introduction of label information in the pre-training stage, live samples are well differentiated from spoof samples.

If the contrastive module is removed, the model’s performance is unsatisfactory. Meanwhile, the weighting factor brings a significant boost to our model performance. Because this module brings task-related information to the model, it makes sense for small dataset-based tasks.

4.5. Visualization and Analysis

4.5.1 Reconstruct Visualization.

To further explore the role of decoder, we visualize the reconstruction results based on different decoders size. As shown in Figure 5, we can see that the larger decoder can better reconstruct the original image. Meanwhile, we find that when the size of decoder is too small, the reconstructed image of the live samples and spoof samples are very similar in details. In that case, features of live samples cannot be distinguished well.

4.5.2 Attention Visualization.

In order to find the area where the model focuses on, we choose Grad-CAM [46] to describe the activation maps on the original images. We compare the attention visualization with different model in Figure 6. The pre-trained ViT model focus on the whole area of the face image. If we do not pre-train the ViT model, the model pay attention to different location of the input image. Meanwhile, we also compare ViT model with CNN model, we find that the CNN model pay attention to the local feature which is very different from pre-trained ViT model.

4.5.3 t-SNE.

We visualize the t-SNE [36] in Figure 7 to analyze our MCAE model feature space. We can observe that all of the live samples from different domains are aggregated. This phenomenon indicates that our method is effective to push all of the live samples close. Furthermore, we find that the distance between live samples and spoof samples is large, and it is easy to distinguish them from the feature space.

5. Conclusion

In this paper, we propose a novel model Masked Contrastive Autoencoder for cross-domain face anti-spoofing task. Meanwhile in order for a feature extractor to extract common features in live samples from different domains, we combine Masked Image Model (MIM) with supervised contrastive learning to train our model. By taking the FAS tasks as an example, we find the MIM pre-training has potential to improve downstream performance on small size dataset. Some intriguing design principles are summarized for performing MIM pre-training for downstream tasks, which are different from previous experience with MIM pre-training on large-scale natural datasets. We hope that these found design principles can inspire researchers to design better MIM pre-training strategies for a variety of downstream tasks, especially on small data sets.

References

- [1] David Barber Felix Agakov. The im algorithm: a variational approach to information maximization. *Advances in neural information processing systems*, 16(320):201, 2004. 4
- [2] Hangbo Bao, Li Dong, and Furu Wei. Beit: Bert pre-training of image transformers. *arXiv preprint arXiv:2106.08254*, 2021. 2
- [3] Zinelabinde Boulkenafet, Jukka Komulainen, Lei Li, Xiaoyi Feng, and Abdenour Hadid. Oulu-npu: A mobile face presentation attack database with real-world variations. In *2017 12th IEEE international conference on automatic face & gesture recognition (FG 2017)*, pages 612–618. IEEE, 2017. 5
- [4] Long Chen, Fei Wang, Ruijing Yang, Fei Xie, Wenjing Wang, Cai Xu, Wei Zhao, and Ziyu Guan. Representation learning from noisy user-tagged data for sentiment classification. *International Journal of Machine Learning and Cybernetics*, pages 1–16, 2022. 3, 5
- [5] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *International conference on machine learning*, pages 1597–1607. PMLR, 2020. 3
- [6] Zhihong Chen, Taiping Yao, Kekai Sheng, Shouhong Ding, Ying Tai, Jilin Li, Feiyue Huang, and Xinyu Jin. Generalizable representation learning for mixture domain face anti-spoofing. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 1132–1139, 2021. 1, 6
- [7] Ivana Chingovska, André Anjos, and Sébastien Marcel. On the effectiveness of local binary patterns in face anti-spoofing. In *2012 BIOSIG-proceedings of the international conference of biometrics special interest group (BIOSIG)*, pages 1–7. IEEE, 2012. 5
- [8] Thomas M. Cover and Joy A. Thomas. *Elements of information theory (2. ed.)*. Wiley, 2006. 4
- [9] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR 2009)*, 20-25 June 2009, Miami, Florida, USA, pages 248–255. IEEE Computer Society, 2009. 2
- [10] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018. 2
- [11] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020. 1, 2
- [12] Tiago de Freitas Pereira, André Anjos, José Mario De Martino, and Sébastien Marcel. Lbp- top based countermeasure against face spoofing attacks. In *Asian Conference on Computer Vision*, pages 121–132. Springer, 2012. 2
- [13] Anjith George and Sébastien Marcel. On the effectiveness of vision transformers for zero-shot face anti-spoofing. In *International IEEE Joint Conference on Biometrics, IJCB 2021, Shenzhen, China, August 4-7, 2021*, pages 1–8. IEEE, 2021. 1
- [14] Jean-Bastien Grill, Florian Strub, Florent Altché, Corentin Tallec, Pierre Richemond, Elena Buchatskaya, Carl Doersch, Bernardo Avila Pires, Zhaohan Guo, Mohammad Gheshlaghi Azar, et al. Bootstrap your own latent-a new approach to self-supervised learning. *Advances in neural information processing systems*, 33:21271–21284, 2020. 3
- [15] Beliz Gunel, Jingfei Du, Alexis Conneau, and Ves Stoyanov. Supervised contrastive learning for pre-trained language model fine-tuning. *arXiv preprint arXiv:2011.01403*, 2020. 3, 5
- [16] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked autoencoders are scalable vision learners. *arXiv:2111.06377*, 2021. 2, 4
- [17] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for unsupervised visual representation learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 9729–9738, 2020. 3
- [18] Hsin-Ping Huang, Deqing Sun, Yaojie Liu, Wen-Sheng Chu, Taihong Xiao, Jinwei Yuan, Hartwig Adam, and Ming-Hsuan Yang. Adaptive transformers for robust few-shot cross-domain face anti-spoofing. *CoRR*, abs/2203.12175, 2022. 1
- [19] Hsin-Ping Huang, Deqing Sun, Yaojie Liu, Wen-Sheng Chu, Taihong Xiao, Jinwei Yuan, Hartwig Adam, and Ming-Hsuan Yang. Adaptive transformers for robust few-shot cross-domain face anti-spoofing. *arXiv preprint arXiv:2203.12175*, 2022. 2, 6
- [20] Yao-Hui Huang, Jun-Wei Hsieh, Ming-Ching Chang, Lipeng Ke, Siwei Lyu, and Arpita Samanta Santra. Multi-teacher single-student visual transformer with multi-level attention for face spoofing detection. In *32nd British Machine Vision Conference 2021, BMVC 2021, Online, November 22-25, 2021*, page 125. BMVA Press, 2021. 1
- [21] Yunpei Jia, Jie Zhang, Shiguang Shan, and Xilin Chen. Single-side domain generalization for face anti-spoofing. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8484–8493, 2020. 1, 6
- [22] Prannay Khosla, Piotr Teterwak, Chen Wang, Aaron Sarna, Yonglong Tian, Phillip Isola, Aaron Maschinot, Ce Liu, and Dilip Krishnan. Supervised contrastive learning. *Advances in Neural Information Processing Systems*, 33:18661–18673, 2020. 3, 5
- [23] Jukka Komulainen, Abdenour Hadid, and Matti Pietikäinen. Context based face anti-spoofing. In *2013 IEEE Sixth International Conference on Biometrics: Theory, Applications and Systems (BTAS)*, pages 1–8. IEEE, 2013. 2
- [24] Ajian Liu and Yanyan Liang. Ma-vit: Modality-agnostic vision transformers for face anti-spoofing. In Luc De Raedt, editor, *Proceedings of the Thirty-First International Joint Conference on Artificial Intelligence, IJCAI 2022, Vienna, Austria, 23-29 July 2022*, pages 1180–1186. ijcai.org, 2022. 1

- [25] Shubao Liu, Ke-Yue Zhang, Taiping Yao, Mingwei Bi, Shouhong Ding, Jilin Li, Feiyue Huang, and Lizhuang Ma. Adaptive normalized representation learning for generalizable face anti-spoofing. In *Proceedings of the 29th ACM International Conference on Multimedia*, pages 1469–1477, 2021. 6
- [26] Shubao Liu, Ke-Yue Zhang, Taiping Yao, Kekai Sheng, Shouhong Ding, Ying Tai, Jilin Li, Yuan Xie, and Lizhuang Ma. Dual reweighting domain generalization for face presentation attack detection. *arXiv preprint arXiv:2106.16128*, 2021. 1, 6
- [27] Son T Ly, Bai Lin, Hung Q Vo, Dragan Maric, Badri Roysam, and Hien V Nguyen. Student collaboration improves self-supervised learning: Dual-loss adaptive masked autoencoder for brain cell image analysis. *arXiv preprint arXiv:2205.05194*, 2022. 4
- [28] Zuheng Ming, Zitong Yu, Musab Qassem Al-Ghadi, Muriel Visani, Muhammad Muzzamil Luqman, and Jean-Christophe Burie. Vitranspad: Video transformer using convolution and self-attention for face presentation attack detection. *CoRR*, abs/2203.01562, 2022. 1
- [29] Keyurkumar Patel, Hu Han, and Anil K Jain. Secure face unlock: Spoof detection on smartphones. *IEEE transactions on information forensics and security*, 11(10):2268–2283, 2016. 2
- [30] Yunxiao Qin, Chenxu Zhao, Xiangyu Zhu, Zezheng Wang, Zitong Yu, Tianyu Fu, Feng Zhou, Jingping Shi, and Zhen Lei. Learning meta model for zero- and few-shot face anti-spoofing. In *The Thirty-Fourth AAAI Conference on Artificial Intelligence, AAAI 2020, The Thirty-Second Innovative Applications of Artificial Intelligence Conference, IAAI 2020, The Tenth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2020, New York, NY, USA, February 7-12, 2020*, pages 11916–11923. AAAI Press, 2020. 1
- [31] Rui Shao, Xiangyuan Lan, Jiawei Li, and Pong C Yuen. Multi-adversarial discriminative deep domain generalization for face presentation attack detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10023–10031, 2019. 5, 6
- [32] Rui Shao, Xiangyuan Lan, and Pong C. Yuen. Regularized fine-grained meta face anti-spoofing. In *The Thirty-Fourth AAAI Conference on Artificial Intelligence, AAAI 2020, The Thirty-Second Innovative Applications of Artificial Intelligence Conference, IAAI 2020, The Tenth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2020, New York, NY, USA, February 7-12, 2020*, pages 11974–11981. AAAI Press, 2020. 1
- [33] Karthik Sridharan and Sham M. Kakade. An information theoretic framework for multi-view learning. In Rocco A. Servedio and Tong Zhang, editors, *21st Annual Conference on Learning Theory - COLT 2008, Helsinki, Finland, July 9-12, 2008*, pages 403–414. Omnipress, 2008. 4
- [34] Chen Sun, Abhinav Shrivastava, Saurabh Singh, and Abhinav Gupta. Revisiting unreasonable effectiveness of data in deep learning era. In *IEEE International Conference on Computer Vision, ICCV 2017, Venice, Italy, October 22-29, 2017*, pages 843–852. IEEE Computer Society, 2017. 2
- [35] Yao-Hung Hubert Tsai, Yue Wu, Ruslan Salakhutdinov, and Louis-Philippe Morency. Self-supervised learning from a multi-view perspective. In *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. OpenReview.net, 2021. 4
- [36] Laurens Van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. *Journal of machine learning research*, 9(11), 2008. 8
- [37] Chien-Yi Wang, Yu-Ding Lu, Shang-Ta Yang, and Shang-Hong Lai. Patchnet: A simple face anti-spoofing framework via fine-grained patch recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 20281–20290, June 2022. 2, 6
- [38] Guoqing Wang, Hu Han, Shiguang Shan, and Xilin Chen. Cross-domain face presentation attack detection via multi-domain disentangled representation learning. In *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2020, Seattle, WA, USA, June 13-19, 2020*, pages 6677–6686. Computer Vision Foundation / IEEE, 2020. 1
- [39] Guoqing Wang, Hu Han, Shiguang Shan, and Xilin Chen. Cross-domain face presentation attack detection via multi-domain disentangled representation learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6678–6687, 2020. 6
- [40] Zhuo Wang, Zezheng Wang, Zitong Yu, Weihong Deng, Jiahong Li, Tingting Gao, and Zhongyuan Wang. Domain generalization via shuffled style assembly for face anti-spoofing. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4123–4133, June 2022. 1, 2, 6
- [41] Chen Wei, Haoqi Fan, Saining Xie, Chao-Yuan Wu, Alan Yuille, and Christoph Feichtenhofer. Masked feature prediction for self-supervised visual pre-training. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14668–14678, 2022. 2
- [42] Di Wen, Hu Han, and Anil K Jain. Face spoof detection with image distortion analysis. *IEEE Transactions on Information Forensics and Security*, 10(4):746–761, 2015. 5, 6
- [43] Zitong Yu, Jun Wan, Yunxiao Qin, Xiaobai Li, Stan Z Li, and Guoying Zhao. Nas-fas: Static-dynamic central difference network search for face anti-spoofing. *IEEE transactions on pattern analysis and machine intelligence*, 43(9):3005–3023, 2020. 6
- [44] Kaipeng Zhang, Zhanpeng Zhang, Zhifeng Li, and Yu Qiao. Joint face detection and alignment using multitask cascaded convolutional networks. *IEEE signal processing letters*, 23(10):1499–1503, 2016. 5
- [45] Zhiwei Zhang, Junjie Yan, Sifei Liu, Zhen Lei, Dong Yi, and Stan Z Li. A face antispoofing database with diverse attacks. In *2012 5th IAPR international conference on Biometrics (ICB)*, pages 26–31. IEEE, 2012. 5
- [46] Bolei Zhou, Aditya Khosla, Agata Lapedriza, Aude Oliva, and Antonio Torralba. Learning deep features for discriminative localization. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2921–2929, 2016. 8