

# JigsawGAN: Auxiliary Learning for Solving Jigsaw Puzzles with Generative Adversarial Networks

Ru Li, *Student Member, IEEE*, Shuaicheng Liu, *Member, IEEE*, Guangfu Wang, Guanghui Liu, *Senior Member, IEEE*, Bing Zeng, *Fellow, IEEE*

**Abstract**—The paper proposes a solution based on Generative Adversarial Network (GAN) for solving jigsaw puzzles. The problem assumes that an image is divided into equal square pieces, and asks to recover the image according to information provided by the pieces. Conventional jigsaw puzzle solvers often determine the relationships based on the boundaries of pieces, which ignore the important semantic information. In this paper, we propose JigsawGAN, a GAN-based auxiliary learning method for solving jigsaw puzzles with unpaired images (with no prior knowledge of the initial images). We design a multi-task pipeline that includes, (1) a classification branch to classify jigsaw permutations, and (2) a GAN branch to recover features to images in correct orders. The classification branch is constrained by the pseudo-labels generated according to the shuffled pieces. The GAN branch concentrates on the image semantic information, where the generator produces the natural images to fool the discriminator, while the discriminator distinguishes whether a given image belongs to the synthesized or the real target domain. These two branches are connected by a flow-based warp module that is applied to warp features to correct the order according to the classification results. The proposed method can solve jigsaw puzzles more efficiently by utilizing both semantic information and boundary information simultaneously. Qualitative and quantitative comparisons against several representative jigsaw puzzle solvers demonstrate the superiority of our method.

**Index Terms**—Solving jigsaw puzzles, generative adversarial networks, auxiliary learning

## I. INTRODUCTION

Solving the jigsaw puzzle is a challenging problem that involves research in computer science, mathematics and engineering. It abstracts a range of computational problems that a set of unordered fragments should be organized into their original combination. Numerous applications are subsequently raised, such as reassembling archaeological artifacts [1]–[4], recovering shredded documents or photographs [5]–[7] and genome biology [8]. Standard jigsaw puzzles are made by

dividing images into interlocking patterns of pieces. With the pieces separated from each other and mixed randomly, the challenge is to reassemble these pieces into the original image. The degree of difficulty is determined by the number of pieces, the shape of the pieces and the graphical composition of the picture itself. The automatic solution of puzzles, without having any information on the underlying image, is NP-complete [9], [10]. This task relies on computer vision algorithms, such as contour or feature detection [11]. Recent development of deep learning (DL) opens bright perspectives for finding better reorganizations more efficiently.

The first algorithm that attempted to automatically solve general puzzles was introduced by Freeman and Gardner [14]. The original approach was designed to solve puzzles with 9 pieces by only considering the geometric shape of the pieces. Various puzzle reconstruction methods are then introduced, which obey a basic operation that the unassigned piece should link to an existing piece by finding its best neighbor according to some affinity functions [1], [15], [16] or matching contours [6], [17], [18]. These algorithms are slow because they keep iterating themselves until a proper solution is produced, as long as all parts are put in place. Moreover, the ‘best neighbors’ may be false-positives due to the occasionally unobvious continuation between true neighbors. Some subsequent works are proposed to avoid checking all possible permutations of piece placements [19]–[22]. However, they mainly concentrate on the boundary information (four boundaries of each piece), whereas ignoring the semantic information that is useful for image understanding. Recently, some methods based on convolutional neural networks (CNN) have been introduced for predicting the permutations by neural optimization [23], [24] or the relative position of a fragment with respect to another [12], [25], [26]. Similar to previous methods, these algorithms do not consider the semantic content information and only rely on learning the position of each piece.

In this paper, we propose JigsawGAN, a GAN-based auxiliary learning pipeline that combines the boundary information of pieces and the semantic information of the global images to solve jigsaw puzzles, which can better simulate the procedure that how humans solve jigsaw puzzles. Auxiliary learning approaches aim to maximize the prediction of a primary task by supervising the model to additionally learn a secondary task [27], [28]. In this work, we address the problem of permutation classification by simultaneously optimizing the GAN network as an auxiliary task. The proposed pipeline consists of two branches, a classification branch (Fig. 2 (b)) and a GAN branch (Fig. 2 (d), (e) and a discriminator). The

Manuscript received January 12, 2021; revised June 24, 2021 and September 21, 2021; accepted September 30, 2021. Date of publication December 7, 2021; date of current version December 15, 2021. This work was supported in part by the National Natural Science Foundation of China (NSFC) under Grant 61872067, Grant 62071097, Grant 62031009, and Grant 61720106004, in part by the “111” Projects under Grant B17008, in part by Sichuan Science and Technology Program under Grants 2019YFH0016. The associate editor coordinating the review of this manuscript and approving it for publication was Dr. Charles Kervrann. (Corresponding authors: Guanghui Liu and Shuaicheng Liu.)

Ru Li, Shuaicheng Liu, Guanghui Liu and Bing Zeng are with School of Information and Communication Engineering, University of Electronic Science and Technology of China, Chengdu 611731, China (e-mail: guanghuiliu@uestc.edu.cn; liushuaicheng@uestc.edu.cn)

Guangfu Wang is with Megvii Technology, Chengdu 610000, China.

Digital Object Identifier 10.1109/TIP.2021.3120052

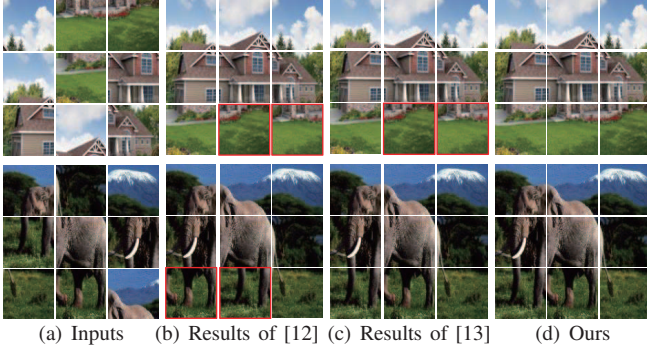


Fig. 1. Comparisons with two recent jigsaw puzzle solvers. (a) Input images. (b) Results of Paumard *et al.*'s method [12]. (c) Results of Huoyan *et al.*'s method [13]. (d) Our results. Paumard *et al.*'s method [12] cannot handle these two cases and Huoyan *et al.*'s method [13] fails to solve the house example. The proposed JigsawGAN works well on these cases due to the concentration of boundary and semantic information simultaneously.

classification branch classifies the jigsaw permutations, and the GAN branch recovers features to images with correct sequences. The two branches are connected by the encoder (Fig. 2 (a)) and a flow-based warp module (Fig. 2 (c)).

The flow-based warp module is applied to warp encoder features to the predicted positions according to the classification results. Moreover, it contributes to the gradient back propagation from the GAN branch to the classification branch. The GAN branch consists of the decoder and the discriminator. The former one recovers the features to the images and the latter one aims to classify that an image belongs to the generated dataset or the real dataset. The discriminator can be regarded as the process that humans judge whether the reorganization of pieces is correct. The judgment process is essential to equip the pipeline with the capability of high-level image understanding, together with low-level boundary clues. The boundary loss and the GAN loss are applied for constraining the boundaries of pieces and distinguishing the semantic information, respectively. The JigsawGAN is mainly designed for solving  $3 \times 3$  puzzles and existing algorithms for  $3 \times 3$  puzzles can benefit from our method. We introduce reference labels to guide the convergence of the classification network. The GAN branch can push the encoder to generate more informative features, thereby obtaining higher classification accuracy. Our method is able to improve the reorganization accuracy of existing methods if we select their results as our reference labels. Detailed descriptions and experiments are presented in Section IV-E4.

Fig. 1 shows the comparisons with two recent methods [12], [13]. Paumard *et al.* proposed a CNN-based method to detect the neighbor pieces and applied the shortest path optimization to recover the images [12]. Huoyan *et al.* applied the graph connection Laplacian algorithm to determine the boundary relationships [13]. Fig. 1 exhibits the reassembled pieces according to the predicted labels. It includes two hard examples, in which three grassland pieces in the house example and two leg pieces in the elephant example can be easily confused. Paumard *et al.*'s method [12] cannot handle the two cases and Huoyan *et al.*'s method [13] fails to solve the house example.

In contrast, our method works well on these cases due to the concentration on the boundary information and semantic clues.

Overall, the main contributions are summarized as follows:

- We propose JigsawGAN, a GAN-based architecture to solve jigsaw puzzles, in which both the boundary information and semantic clues are fused for the inference.
- We introduce a flow-based warp module to reorganize feature pieces, which also ensures the gradient from the GAN branch can be back-propagated to the classification branch during the training process.
- We provide quantitative and qualitative comparisons with several typical jigsaw puzzle solvers to demonstrate the superiority of the proposed method.

## II. RELATED WORKS

Introduced by Freeman and Gardner [14], the puzzle solving problems have been studied by many researchers, even though Demaine *et al.* discovered that puzzle assembling is NP-hard if the dissimilarity is unreliable [10].

### A. Solving Square-piece Puzzles

Most jigsaw solvers assume that the input includes equal-sized squared pieces. The first work was proposed by [29], where a greedy algorithm and a benchmark were proposed. They presented a probabilistic solver to achieve approximated puzzle reconstruction. Pomeranz *et al.* evaluated some compatibility metrics and proposed a new compatibility metric to predict the probability that two given parts are neighbors [16]. Gallagher *et al.* divided the squared jigsaw puzzle problems into three types [19]. ‘Type 1’ puzzle means the orientation of each jigsaw piece is known, and only the location of each piece is unknown. ‘Type 2’ puzzle is defined as a non-overlapping square-piece jigsaw puzzle with unknown dimension, unknown piece rotation, and unknown piece position. Meanwhile, the global geometry and position of every jigsaw piece in ‘Type 3’ puzzle are known, and only the orientation of each piece is unknown. Gallagher *et al.* solved ‘Type 1’ and ‘Type 2’ puzzle problems through the minimum spanning tree (MST) algorithm constrained by geometric consistency between pieces. A Markov Random Field (MRF)-based algorithm was also proposed to solve the ‘Type 3’ puzzle.

Some subsequent works were proposed to solve the ‘Type 1’ puzzle problem. Yu *et al.* applied linear programming to exploit all pairwise matches, and computed the position of each component [30]. There are also some methods aimed to solve the complicated ‘Type 2’ puzzle problem. Son *et al.* introduced loop constraints for assembling non-overlapping jigsaw puzzles where the rotation and the position of each piece are unknown [20]. Huoyan *et al.* utilized graph connection Laplacian to recover the rotations of the pieces when both shuffles and rotations are unknown [13]. Some variants were investigated subsequently, including handling puzzles with missing pieces [31] and eroded boundaries [12]. Paikin *et al.* calculated dissimilarity between pieces and then proposed a greedy algorithm for handling puzzles of unknown size and with missing entries [31]. Bridger *et al.* proposed a GAN-based architecture to recover the eroded regions and

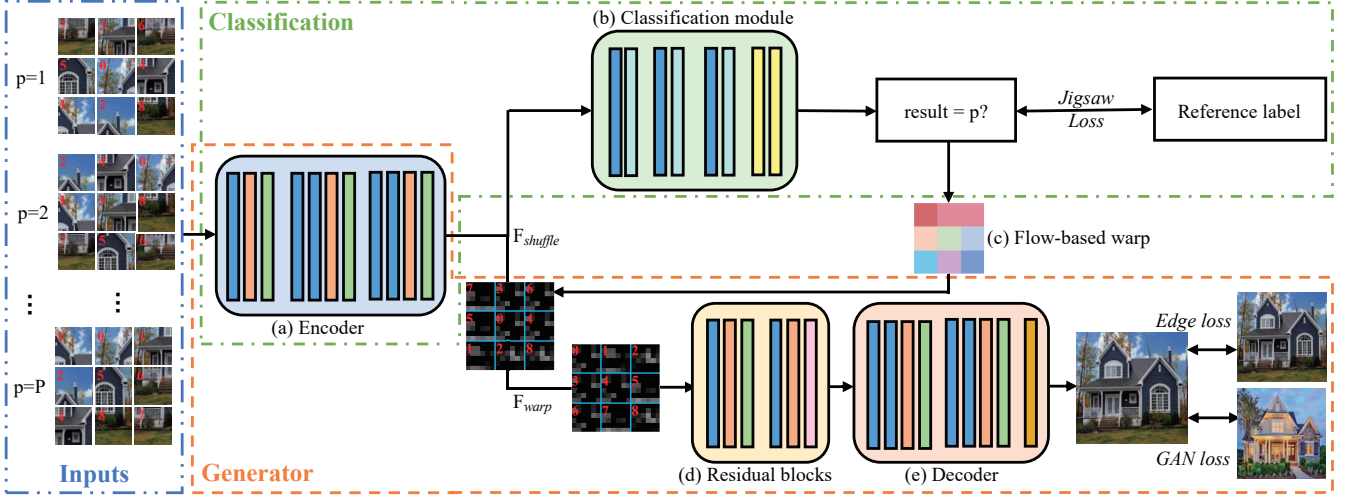


Fig. 2. The pipeline exhibits the architecture of the proposed JigsawGAN, which consists of two streams that operate different functions. **The top stream** classifies the jigsaw permutations. In turn, **the bottom stream** recovers features to images. The two streams are connected by the encoder and a flow-based warp module. The encoder extracts features of input shuffled pieces. The shuffled features are then combined, called  $F_{shuffle}$ , to be fed to the classification module to predict the jigsaw permutation. The classification results provides information to generate flow fields that can reassemble the combined feature  $F_{shuffle}$  to a new combined feature  $F_{warp}$ . The warped feature  $F_{warp}$  is then sent to the residual blocks and the decoder to recover natural image.

reassembled the images using greedy algorithms [26]. Our method concentrates on solving the ‘Type 1’ non-overlapping square-piece jigsaw puzzles.

Current DL-based methods are mainly designed for solving 9 pieces with strong robustness compared to traditional algorithms. Dery *et al.* [24] applied a pre-trained VGG model as feature extractor and corrected the order using a pointer network [32]. Paumard *et al.* utilized a CNN-based method to detect the neighbor pieces and used the shortest path optimization to recover the eroded images [12]. These methods obey the conventional rule to solve jigsaw puzzles, which consists of two steps. The first step is applying the neural network to determine the relationships of pieces, and the second step is using optimization techniques to reorganize the pieces. Their performances are limited by the first step due to the aperture problem based only on the local pieces. We propose a GAN-based approach to include global semantic information apart from local pieces for a better accuracy.

### B. Pre-training Jigsaw Puzzle Solvers

Numerous self-supervised methods consider solving jigsaw puzzles as pre-text tasks. This learning strategy is a recent variation of the unsupervised learning theme that transfers the pre-trained network parameters on jigsaw puzzle tasks to other visual recognition tasks [33]. These methods assume that a rich universal representation has been captured in the pre-trained model, which is useful to be fine-tuned with the task-specific data using various strategies. Noroozi *et al.* introduced a context-free network (CFN) to separate the pieces in the convolutional process. Its main architecture focused on a subset of possible permutations including all the image tiles and solved a classification problem [34]. Santa *et al.* proposed to handle the whole set by approximating the permutation matrix and solving a bi-level optimization problem to recover the right ordering [35]. The above methods tackled the problem

by dealing with the separate pieces and then finding a way to recombine them. Carlucci *et al.* proposed JiGen to train a jigsaw classifier and a object classifier simultaneously. They focused on domain generalization tasks by considering that the jigsaw puzzle solver can improve semantic understanding [36]. Du *et al.* combined the jigsaw puzzle and the progressive training to optimize the fine-grained classification by learning which granularities are the most discriminative and how to fuse information cross multi-granularity [37]. However, solving jigsaw puzzle tasks in these methods are supervised, which depend heavily on the training data. In this paper, we propose a GAN-based auxiliary learning method to solve jigsaw puzzles, where the ground truth of jigsaw puzzle task is unavailable. Auxiliary learning in our method means that we formulate the permutation classification as the primary task with the secondary goal is to optimize the GAN network.

## III. METHOD

We propose a GAN-based architecture for solving the jigsaw puzzles with unpaired datasets. The architecture includes two streams operating different functions. The first stream defines the jigsaw puzzle solving as a classification task to judge which permutation the shuffled input belongs to. The second stream is composed of the generator  $G$  and the discriminator  $D$ . The generator  $G$  learns the mapping function between different domains, while the discriminator  $D$  aims to optimize the generator by distinguishing target domain images from the generated ones. Specifically, we generate a wide diversity of shuffled images  $\{x_i\}_{i=1,\dots,N} \in X$  as the source domain data, and a collection of natural images  $\{y_j\}_{j=1,\dots,M} \in Y$ , as the target domain data. The data distributions of the two domains are denoted as  $x \sim p_{\text{data}}(x)$  and  $y \sim p_{\text{data}}(y)$ , respectively.



TABLE I

LAYER CONFIGURATIONS OF THE ARCHITECTURE.  $B$  REPRESENTS THE BATCH SIZE,  $H_p$  AND  $W_p$  ARE THE HEIGHT AND WIDTH OF INPUT PIECES. THE ENCODER EXTRACTS FEATURES OF INPUT SHUFFLED PIECES. THE SHUFFLED FEATURES ARE THEN COMBINED ( $F_{shuffle}$ ) WITH SIZE  $H_f$  AND  $W_f$ . THE WARPED FEATURE  $F_{warp}$  IS SENT TO THE RESIDUAL BLOCKS AND THE DECODER TO RECOVER NATURAL IMAGE.

(a) Encoder				
Input: $[B \times n \times n, 3, H_p, W_p]$				
Conv			BN	Activation
Kernel	Stride	Channel	Channel	
7	1	64	64	ReLU
3	2	128	-	-
3	1	128	128	ReLU
3	2	256	-	-
3	1	256	256	ReLU
$F_{shuffle}: [B, 256, H_f, W_f]$				
(c) Residual blocks and decoder				
$F_{warp}: [B, 256, H_f, W_f]$				
Conv			BN	Activation
Kernel	Stride	Channel	Channel	
3	2	256	256	ReLU
3	1	256	256	ES
3	1/2	128	-	-
3	1	128	128	ReLU
3	1/2	64	-	-
3	1	64	64	ReLU
7	1	3	-	-
Generator output: $[B, 3, H, W]$				

(b) Classification module					
$F_{shuffle}: [B, 256, H_f, W_f]$					
Conv			Activation	Max pooling	FC
Kernel	Stride	Channel		Stride	Channel
3	2	256	ReLU	2	-
3	2	384	ReLU	2	-
3	2	384	ReLU	2	-
-	-	-	-	-	4096
-	-	-	-	-	$P$
Classification output: $B \times P$					
(d) Discriminator					
Input: $[B, 3, H, W]$					
Conv			BN	Activation	
Kernel	Stride	Channel	Channel		
3	1	32	-	LReLU	
3	2	64	-	LReLU	
3	1	64	64	LReLU	
3	2	128	-	LReLU	
3	1	128	128	LReLU	
3	1	256	256	LReLU	
3	1	1	-	-	
Output: $[B, 1, H/4, W/4]$					

### A. Network Architecture

We present the classification network  $C$ , generator  $G$  in Fig. 2. where  $C$ ,  $G$  and the discriminator  $D$  are simultaneously optimized during the training process. Their detailed layer configurations are shown in Table I. The inputs are obtained by decomposing the source images into  $n \times n$  pieces, which are then shuffled and reassigned to one of the  $n^2$  grid positions to generate  $n^2!$  combinations.  $n^2!$  is a very large number and scarcely possible to be considered as classification tasks if  $n > 2$ , for example,  $(3^2)! = 362,880$ . Generally,  $P$  elements are selected from  $n^2!$  combinations according to the maximal Hamming distance [34], and we assign an index to each entry. In our implementation, solving the jigsaw puzzle is considered as a classification task and each permutation corresponds to a classification label. The number of permutations indicates the classification categories. The ultimate target is to predict correct permutation of the shuffled input image.

The shuffled input image is first divided into  $n \times n$  pieces, and the pieces are sent to the network in parallel. We choose the discrete pieces as inputs in order to prevent the cross-influence between the boundaries of pieces. The encoder (Fig. 2 (a)) extracts features of these shuffled pieces. The shuffled features are then combined, called  $F_{shuffle}$ , to be sent to the classification module (Fig. 2 (b)) to predict the jigsaw permutation. The classification module is constrained by the pseudo-labels, called reference labels. During the training, we find that an unsupervised classification network converges only with difficulty. To solve the problem, we generate the reference labels to guide  $C$ . The classification results provide information to generate flow fields (Fig. 2 (c)) that can reassemble the combined feature  $F_{shuffle}$  to a new combined feature  $F_{warp}$ . The warped feature  $F_{warp}$  is then sent to the residual blocks (Fig. 2 (d)) and the decoder (Fig. 2 (e)) to recover natural images. The decoder can recover a perfect

reorganized image if  $F_{warp}$  is reassembled with the correct order according to the classification label.

When the real permutation labels and the correct natural images are unavailable, we adopt the GAN architecture to achieve the unsupervised optimization. The vanilla GAN consists of two components, a generator  $G$  and a discriminator  $D$ , where  $G$  is responsible for capturing the data distribution while  $D$  tries to distinguish whether a sample comes from the real data or the generator. This framework corresponds to a min-max two-player game, and introduces a powerful way to estimate the target distribution. The GAN branch provides global constraints to enable the network to concentrate on semantic clues. The decoder and the discriminator as well as extra losses enable the encoder to generate more informative features and further improve the classification network.

**Classification Network:** The classification network consists of two parts: the convolutional blocks (the encoder) and the classification module. The encoder extracts useful signals for downstream transformation. The classification module aims to distinguish different permutations, which includes convolutional layers, max pooling layers and fully connected layers.

**Generator:** The generator network is composed of three parts: the encoder, eight residual blocks and the decoder. The encoder of the classification network also extracts features for the generator. Afterward, eight residual blocks with identical layouts are adopted to construct the content and the manifold features. The decoder consists of two identical transposed convolutional blocks and a final convolutional layer.

**Discriminator:** The discriminator network complements the generator and aims to classify each image as real or fake. The discriminator network includes several convolutional blocks. Such a simple discriminator uses fewer parameters and can work on images of arbitrary sizes.

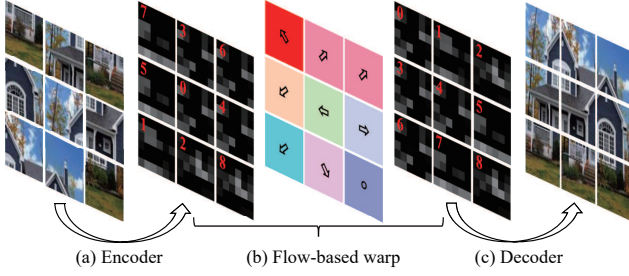


Fig. 3. The flow-based warp module can warp the shuffled features to ordered features according to the classification results.

### B. Flow-based Warp Module

The GAN branch cannot recover the extracted features directly because the features are shuffled. We introduce a flow-based warp module to reassemble the shuffled encoder features and ensure the gradient can be back-propagated correctly during the feature recombination. In detail, we define the flow fields to represent dense pixel correspondences between the shuffled encoder features ( $F_{shuffle}$ ) and the warped features ( $F_{warp}$ ). The optical flows are constructed according to the predicted labels, which are used to warp  $F_{shuffle}$  into  $F_{warp}$ . The size of the flow-based warp module is defined as  $H_f \times W_f \times 2$ , where the first  $H_f \times W_f$  channel controls the horizontal shift of the features and the last  $H_f \times W_f$  channel calculates the vertical shift of the features. The discrete encoder features are combined before reassembling by the flow-based warp module.  $H_f$  and  $W_f$  are the height and width of combined features. The flow-based warp module not only rearranges the shuffled features but also guarantees the gradient back propagation. Note that, some shift operations can also recombine the shuffled features, but block the gradient back propagation. As a result, the semantic information from the GAN branch cannot be delivered to the classification branch, thereby resulting in lower classification accuracy.

The decoder can recover perfect images if the classification labels are accurate, as shown in Fig. 3. Some wrong classification can be tolerated when the network is initialized because the GAN optimization can correct the mistakes gradually. GAN network can detect the reorganization error and rectify it according to the semantic information acquired from the real dataset. The correct information will be transferred to the encoder and the classification module through the gradient back propagation, and can guide the classification network to classify the input correctly on following iterations.

### C. Loss Function

We design our objective function to include the following three losses: (1) the jigsaw loss  $L_{jigsaw}(C)$ , which optimizes the classification network to recognize the correct permutations; (2) the adversarial loss  $L_{GAN}(G, D)$ , which drives the generator network to achieve the desired transformation; (3) the boundary loss  $L_{boundary}(G, D)$ , which pushes the decoder of  $G$  to recover clear images, and further constrains

the encoder to generate useful features that are helpful for the classification network. The full objective function is:

$$L(G, D, C) = L_{jigsaw}(C) + L_{GAN}(G, D) + L_{boundary}(G, D) \quad (1)$$

1) *Jigsaw Loss*: We consider the jigsaw puzzle solving as a classification task and  $P$  is the classification category. Kullback-Leibler (KL) divergence can measure the similarity between the predicted distribution and the target distribution. The KL divergence becomes smaller when two permutations tend to be similar. We aim to minimize the KL divergence between the reassembled result ( $p_{predict}$ ) and the ground truth ( $p_{real}$ ), which can be described as follows:

$$\arg \min KL(p_{predict}, p_{real}) \quad (2)$$

A simple classification network is first proposed to recognize the correct permutation. We minimize the following jigsaw loss to optimize the classification network:

$$L_{jigsaw}(C) = \mathbb{E}_{x \sim p_{data}(x)} [CE(C(x), p)], \quad (3)$$

where  $p$  is the probability distribution of the real data,  $C(x)$  is the probability distribution of the predicted data which indicates the probability that the result belongs to each category.  $p$  and  $C(x)$  are defined as matrix with size  $B \times P$ , where  $B$  is the batch size. CE is cross-entropy loss which is defined as:

$$CE(C(x), p) = - \sum p \cdot \log(C(x)) \quad (4)$$

However, the real jigsaw labels are unavailable for unsupervised tasks. Directly training  $C$  without jigsaw labels is impossible. The classification network tends to assign labels randomly, which achieves 30-40% classification accuracy for three-classification tasks and 20-30% for four-classification tasks according to our experiments. Note that, three-classification tasks mean the dataset contains three categories and the model will predict the most likely category that one input belongs to. Pseudo-labels are introduced to achieve better classification performance. The pseudo-labels, called reference labels, are generated according to the shuffled input images. The reference labels can constrain  $C$  when permutation indexes are unavailable. We apply the following five steps to obtain the reference labels:

- 1) Cutting four boundaries from each piece, and the width of the boundary is determined by a hyper-parameter  $pix$ .
- 2) Calculating the PSNR values between the top boundary of a specific piece and bottom boundaries of other pieces to obtain a vector with size  $1 \times n^2$ . Then, a  $n^2 \times n^2$  matrix can be obtained when computes all top-bottom relationships, as shown in Fig. 4.
- 3) Applying the same way to get another  $n^2 \times n^2$  matrix to indicate the left-right relationships.
- 4) There are  $n \times (n-1)$  correct top-bottom boundary pairs and  $n \times (n-1)$  correct left-right boundary pairs, so we adopt a greedy algorithm [38] to select  $n \times (n-1)$  maximum values from two matrices, respectively.
- 5) A minimum spanning tree (MST) algorithm [19] is then applied to assemble the pieces, and the reference permutation  $p_{ref}$  is returned.

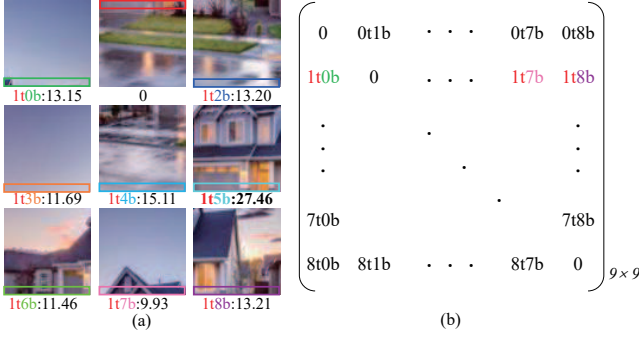


Fig. 4. In this case, we calculate the PSNR values between the top boundary of the top-middle piece and the bottom boundaries of other pieces. Corresponding PSNR values are displayed in (a). These PSNR values compose the second row of the  $9 \times 9$  top-bottom relationship metric, as shown in (b).

TABLE II  
EXPERIMENTS TO SHOW THE INFLUENCE OF THE HYPER-PARAMETER  $pix$ .

	$pix = 3$	$pix = 2$	$pix = 1$
PSNR	20.1	21.5	24.7
Accuracy	56.5%	60.3%	66.4%

Empirically,  $pix$  is set to 1, which means we cut a row of pixels of each boundary. Table II shows the corresponding PSNR values and reorganization accuracy when we select different  $pix$ . A larger  $pix$  may cut a wider boundary strip. However, the redundant information could decrease the PSNR values and the reorganization accuracy. We select  $pix = 1$  when we obtain the boundary strips, which are more sensitive to the boundary discontinuity.

After the reference labels are obtained, focal (FL) loss [39] is applied as the jigsaw loss. FL loss ensures the minimization of the KL divergence whilst increasing the entropy of the predicted distribution, which can prevent the model from becoming overconfident. FL loss can down-weight easy examples and focus training on hard negatives. We replace the CE loss conventionally used with the FL loss to improve the network calibration. FL loss is defined as:

$$FL(C(x), p) = - \sum (1 - C(x) \cdot p)^\gamma \cdot \log(C(x)) \quad (5)$$

where  $(1 - C(x) \cdot p)^\gamma$  is a modulating factor to optimize the imbalance of the dataset and  $\gamma$  is a tunable focusing parameter that smoothly adjusts the rate at which easy examples are down-weighted.  $\gamma$  is set to 2 in our implementations. Finally, the jigsaw loss can be described as follows:

$$L_{jigsaw}(C) = \mathbb{E}_{x \sim p_{data}(x)} [FL(C(x), p_{ref})] \quad (6)$$

2) *Adversarial Loss*: The performance of the classification module  $C$  is limited by the reference labels. After the encoder module generating distinctive features, we add a decoder module to recover the features to natural images. We do not train an autoencoder network that directly recovers the shuffled inputs because it is useless for the encoder to generate more informative features. The proposed flow-based warp module is able to warp the shuffled features ( $F_{shuffled}$ ) to warped features ( $F_{warp}$ ) according to the classification results. The reorganized features  $F_{warp}$  are then sent to the GAN branch

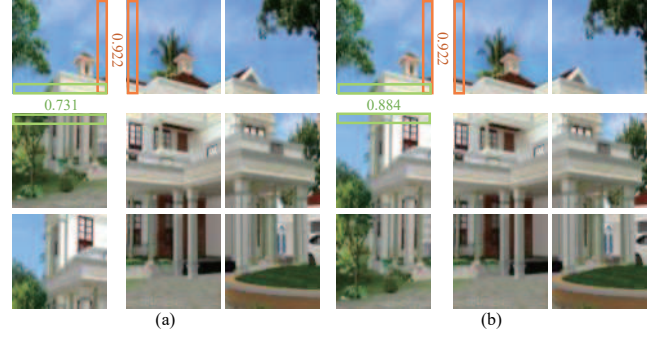


Fig. 5. (a) An incorrect reorganized result that may confuse the discriminator. (b) The correct reassemble result. The numbers indicate the SSIM values between two adjacent boundaries.

to recover the images. The GAN branch is auxiliary to the classification branch and trained to help the encoder module generate useful features that are aware of the semantic information, and further improve the classification performance.

As in classic GAN networks, the adversarial loss is used to constrain the results of  $G$  to look like target domain images. In our task, adversarial loss pushes  $G$  to generate natural images in the absence of corresponding ground truth. Meanwhile,  $D$  aims to distinguish whether a given image belongs to the synthesized or the real target set, which tries to classify an image into two categories: the generated image  $G(x)$  and the real image  $y$ , as formulated in Eq. 7.

$$L_{GAN} = \mathbb{E}_{y \sim p_{data}(y)} [\log D(y)] + \mathbb{E}_{x \sim p_{data}(x)} [\log(1 - D(G(x)))] \quad (7)$$

3) *Boundary Loss*: The adversarial loss ensures that the generated image looks similar to the target domain. It is inadequate for a smooth transaction, especially at the boundaries of pieces. For example, in Fig. 5 (a), sorting adjacent top-bottom pieces in house images mistakenly may confuse the discriminator, such that the discriminator considers the imperfect image as correct image and cannot work as we desired. Therefore, it is essential to enforce a more strict constraint to guarantee the semantic consistency of the generated images. To achieve this, we add penalties at the boundaries of pieces to punish the incorrect combination. We first cut the decoder result into  $n \times n$  pieces and calculate the SSIM of adjacent boundaries. For example, as for  $3 \times 3$  puzzle, the top-left piece needs to calculate the SSIM of its right boundary with the left boundary of the top-middle piece (orange boxes in Fig. 5). Similarly, the SSIM value between its bottom boundary and the top boundary of the middle-left piece is computed (green boxes in Fig. 5). We obtain  $n \times (n - 1)$  values to represent the top-bottom relationships and  $n \times (n - 1)$  values to indicate the left-right relationships. Two average scores of these two relationships are calculated to design the boundary loss, noted as  $SSIM_{tb}$  and  $SSIM_{lr}$ , respectively. As such, the boundary loss is defined as:

$$L_{boundary}(G, D) = \mathbb{E}_{x \sim p_{data}(x)} [(1 - SSIM_{tb}(G(x)))] + \mathbb{E}_{x \sim p_{data}(x)} [(1 - SSIM_{lr}(G(x)))] \quad (8)$$

Matching boundaries generate large SSIM values. The boundary loss provides more penalties to wrong placements.





Fig. 6. Examples generated by the proposed JigsawGAN. The top two rows show the results when the hyper-parameter  $n = 2$ , the middle two rows present the results when the hyper-parameter  $n = 3$ , and the bottom two rows are results when  $n = 4$ . Different styles can be efficiently reassembled by JigsawGAN.

#### IV. EXPERIMENTS

We first illustrate the datasets and implementation details in Section IV-A and then conduct comprehensive experiments to verify the effectiveness of the proposed method, including quantitative comparisons (Section IV-B), qualitative comparisons (Section IV-C), computational times (Section IV-D) and ablation studies (Section IV-E). Specifically, we first compare our method with several representative jigsaw puzzle works, including a classic jigsaw puzzle solver [19], a linear programming-based method [30], a loop constraints-based method [20], a graph connection Laplacian-based method [13], and then compare with a recent DL-based method which relies on the shortest path optimization [12]. Some comparison methods are designed for handling the ‘Type 1’ and ‘Type 2’ puzzles simultaneously, which can be easily applied for our ‘Type 1’ task for solving  $3 \times 3$  pieces. Next, ablation studies are performed to illustrate the importance of different components.

##### A. Datasets and Implementation Details

1) *Data Collection*: The dataset includes 7,639 images, among which 5,156 images originate from the PACS

dataset [40] and 2,483 images are from our own collection. The PACS dataset includes many pictures with similar contents, so we gathered images from movies or Internet that are more distinguishable to increase the diversity. Our dataset covers 4 object categories (House, Person, Elephant and Guitar), and each of them can be divided into 4 domains (Photo, Art paintings, Cartoon and Sketches). Each category is divided into three subsets: a jigsaw set (40%) to generate the inputs, a real dataset (40%) to help the discriminator to distinguish the generated images and real images, and a test set (20%) to evaluate the performance of different methods. They are randomly selected from the overall dataset. The shuffled fragments are then prepared according to [36]: the input image is divided into  $n \times n$  pieces and the size of each piece is set to  $24 \times 24$ . The permutation  $P$  is set to 1,000 in our implementation. The hyper-parameter  $n$  and the permutation  $P$  are important factors for the network performance. We perform ablation studies on these two factors in Section IV-E3.

2) *Training Details*: We implement JigsawGAN in PyTorch, and all the experiments are performed on an NVIDIA RTX 2080Ti GPU with 100 epochs. For each iteration, every

TABLE III  
QUANTITATIVE COMPARISONS WITH SEVERAL REPRESENTATIVE JIGSAW PUZZLE SOLVERS IN TERMS OF REORGANIZATION ACCURACY.

	Gallagher2012 [19]	Son2014 [20]	Yu2015 [30]	Huroyan2020 [13]	Paumard2020 [12]	Ours
House	64.5%	70.3%	73.8%	76.4%	61.2%	<b>80.2%</b>
Person	63.8%	68.2%	70.6%	73.7%	63.6%	<b>79.9%</b>
Guitar	59.4%	64.5%	68.3%	71.1%	60.4%	<b>77.4%</b>
Elephant	60.5%	66%	71.4%	73.1%	58.2%	<b>78.5%</b>
Mean	62.1%	67.3%	71.0%	74.6%	60.9%	<b>79.0%</b>

input image is divided into  $n \times n$  pieces, which are sent to the network in parallel. We choose the discrete pieces as inputs in order to prevent the cross-influence between the boundaries of pieces. Current architecture can handle a maximum of 16 pieces. If we want to deal with images with more pieces, the classification network should be deepened with more convolutional blocks to get a larger respective field. Adam optimizer is applied with the learning rate of  $2.0 \times 10^{-4}$  for the generator, the discriminator and the classification network. The entire training process costs 4 hours on average.

### B. Quantitative Comparisons

Table III lists the reorganization accuracy of the aforementioned methods and the JigsawGAN when  $n = 3$ . The reorganization accuracy means the proportion of correct reassembled images in all test images, and the scores reported in Table III are the average results over the test sets. Some comparison methods [13], [19], [20], [30] obey the basic rules to solve jigsaw puzzles: (1) detecting boundaries to determine the relationship between the pieces, and (2) applying different optimization methods to reassemble the images. Gallagher *et al.* applied the Mahalanobis Gradient Compatibility (MGC) measurement to determine the pieces' relationship and a greedy algorithm to reorganize the pieces [19]. A considerable improvement for [19] was proposed by [20], by adding loop constraints to the piece reorganization process. Yu *et al.* combined the advantages of greedy methods and loop propagation algorithms to introduce a linear programming-based solver [30]. Huroyan *et al.* applied the graph connection Laplacian to better understand the reconstruction mechanism [13]. However, these methods are limited by the boundary detection step, which will lead to the randomness of their results. Paumard *et al.* [12] applied the shortest path optimization algorithm to reorganize the pieces, whose boundary relationship is predicted by a neural network. Their pruning strategy for the shortest path algorithm may affect the performance. Compared to these methods, our JigsawGAN achieves the best scores, which improves the performance against [12], [19], [20] with a relatively large margin and has more than 4 point improvements compared with [13], [30].

### C. Qualitative Comparisons

In this section, we first show some reorganization results generated by JigsawGAN in Fig. 6. Our method can reassemble the inputs and generate high-quality results. Then, qualitative comparisons with aforementioned methods are presented in Fig. 7, which shows the reassembled pieces according to

corresponding predicted labels. The top two rows show the results when the hyper-parameter  $n = 2$ , the middle two rows present the results when  $n = 3$ , and the bottom two rows are results when  $n = 4$ . In the first elephant example, the top two pieces can be easily confused if the methods do not consider the semantic content information. Results of [19], [30] and [12] fail to solve the case. The guitar example has weak boundary constraints between the pieces, which leads to the failure of [20] and [12]. The strategy of Son *et al.* aims to recover the complete shape from pieces based on a dissimilarity metric [20]. The two musical note pieces have a small contribution to construct the guitar structure and confuse Son *et al.*'s method. Paumard *et al.* reassemble pieces with boundary erosion, which directly determine the relationships of pieces through a neural network and ignore the useful boundary information [12]. In the first house example, results of [19], [20] and [30] fail to solve the tree and roof pieces, while the result of [12] cannot distinguish two road pieces. In the second elephant example, they all fail to discriminate the leg pieces. We noted that puzzles with low percentages of recovery by these algorithms have large portions of pieces with the same uniform texture and color. The global semantic information and the strict boundary detection are indispensable to obtain better reassemble results. The failure of the last two house examples of other methods when  $n = 4$  is caused by the similar texture and the unobvious boundary. Note that, the results of [12] are the same as the inputs when  $n = 4$  because their method can only handle puzzles when  $n \leq 3$ . In comparison, by considering the global information (GAN loss) and the boundary information (boundary loss) simultaneously, our method recovers these cases well.

### D. Computational Times

Computing efficiency is also important for evaluating the reassemble performance. We conduct comparisons of computational times with the aforementioned methods. The results for solving  $3 \times 3$  puzzles on the test set are reported in Table IV. Gallagher *et al.*'s method [19] is the fastest algorithm and our method is a bit slower than [19]. However, introducing a useful network with a little more running time is advisable in exchange for better results. Methods [30] and [13] are slower than JigsawGAN due to the linear programming algorithm and the optimized iteration, respectively. As for [20], most of the time is spent in the pairwise matching, the unoptimized merging, the trimming and the filling steps. Paumard *et al.*'s method [12] is the slowest due to the complicated network architecture and the reassemble graph. Although our method takes approximately 4 hours on the training process, the



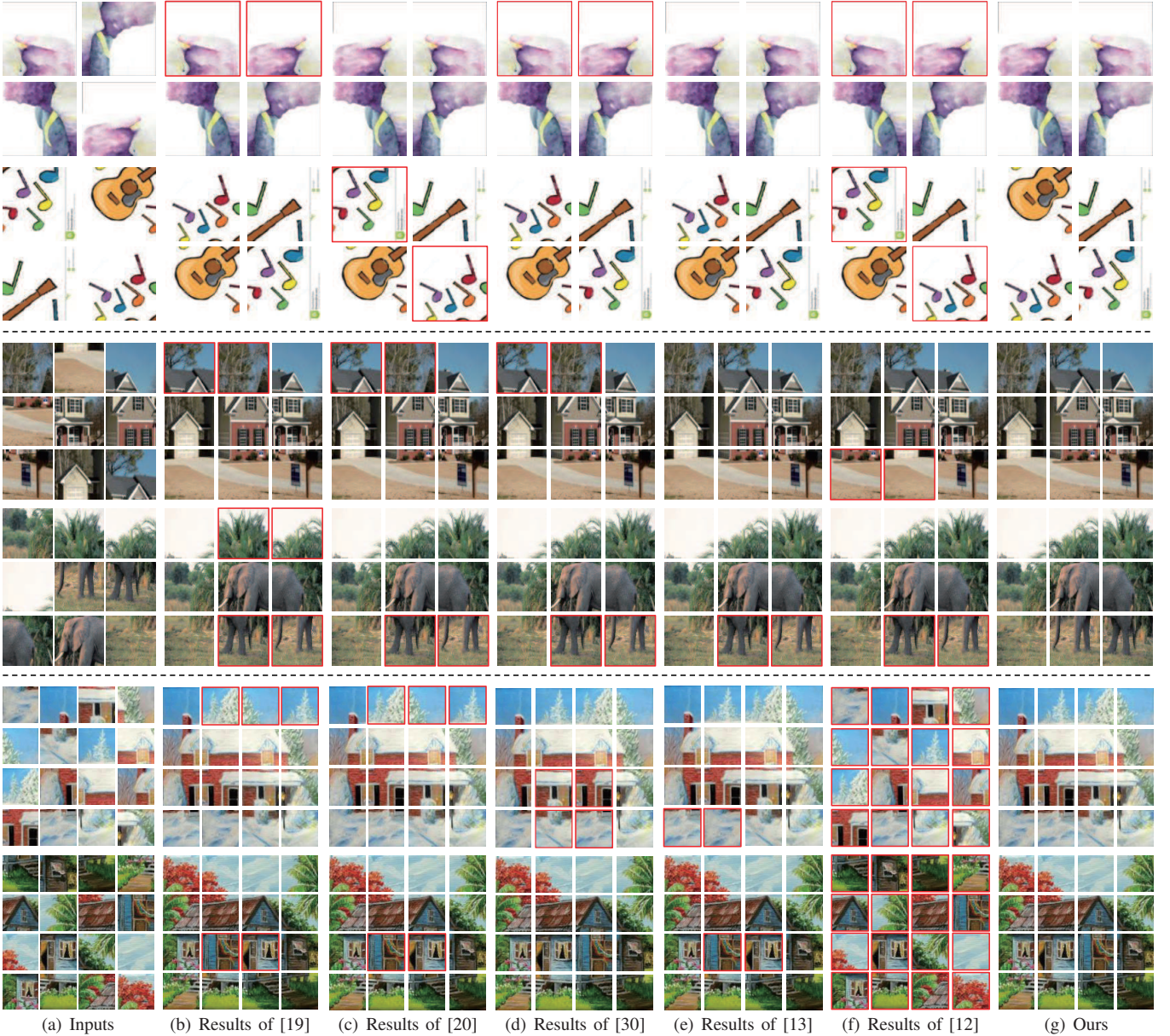


Fig. 7. Qualitative comparisons with Gallagher *et al.*'s method [19], Son *et al.*'s method [20], Yu *et al.*'s method [30], Huroyan *et al.*'s method [13] and Paumard *et al.*'s method [12]. The top two rows show the results when the hyper-parameter  $n = 2$ , the middle two rows present the results when  $n = 3$ , and the bottom two rows are results when  $n = 4$ . The proposed method assembles both the boundary and the semantic information in order to generate more accurate results. The results of [12] are the same as the inputs when  $n = 4$  because their method can only handle puzzles when  $n \leq 3$ .

TABLE IV  
COMPUTATIONAL TIMES OF JIGSAWGAN AND THE COMPARISON METHODS WHEN SOLVING  $3 \times 3$  PUZZLES.

Methods	Gallagher2012 [19]	Son2014 [20]	Yu2015 [30]	Huroyan2020 [13]	Paumard2020 [12]	Ours
Time (s)	0.274	1.121	0.525	0.496	3.054	0.293

PSNR-based reference label and the GAN network are not used in the test process, which take less running time on the test set. The simple yet effective classification network can provide satisfactory reassemble results.

#### E. Ablation Studies

1) *Ablation Study of Loss Terms*: We perform the ablation study on the variants of the loss function to understand how

these main modules contribute to the final results. Table V displays the ablation results, which demonstrates that each component contributes to the objective function. We have illustrated the importance of the jigsaw loss  $L_{jigsaw}$  in Section III-C1. Directly training the classification network without the jigsaw loss is impossible. The network tends to assign results randomly, achieving 30-40% classification accuracy for three-class tasks and 20-30% for four-class tasks. Table V

TABLE V  
ABLATION EXPERIMENTS OF DIFFERENT LOSS TERMS.

	$p_{ref}$	w/o. $L_{GAN}$	w/o. $L_{boundary}$	Ours
Accuracy	66.4%	68.1%	75.8%	79.0%

TABLE VI  
ABLATION EXPERIMENTS OF DIFFERENT DATASETS. P, H, G AND E MEANS THE PERSON, HOUSE, GUITAR AND ELEPHANT DATASETS, RESPECTIVELY. ‘H.G.E AND P’ REPRESENTS THAT WE TRAIN OUR NETWORK ON HOUSE, GUITAR AND ELEPHANT DATASETS, AND THEN TEST ON THE PERSON DATASET.

Accuracy	P and P	H.P.G.E and P	H.G.E and P
	79.9%	72.4%	66.9%
Accuracy	H and H	H.P.G.E and H	P.G.E and H
	80.2%	74.6%	69.7%
Accuracy	G and G	H.P.G.E and G	H.P.E and G
	77.4%	69.8%	63.1%
Accuracy	E and E	H.P.G.E and E	H.P.G and E
	78.5%	72.2%	65.5%

further shows the importance of the adversarial loss  $L_{GAN}$  and the boundary loss  $L_{boundary}$ . As shown in Table V, removing  $L_{GAN}$  degrades the results substantially, so as removing the  $L_{boundary}$ . The third column shows the reorganization accuracy without  $L_{GAN}$ , which indicates the discriminator and corresponding loss are removed. The accuracy of without  $L_{GAN}$  is apparently inferior to the final result. The fourth column displays the result without  $L_{boundary}$ , which means the decoder is only constrained by the adversarial loss. The result is also not as good as our final result. We also conduct the experiment that removes both  $L_{GAN}$  and  $L_{boundary}$ , which represents the GAN branch is not involved, where the final classification accuracy only depends on the reference labels  $p_{ref}$ . The result in the second column is worse than the final result. We conclude that all three loss terms are critical.

2) *Ablation Study of Different Training Sets*: This ablation study aims to demonstrate the effectiveness of the semantic information provided by the GAN branch. In Table VI, ‘P’ is the ‘Person’ dataset, ‘H’ means the ‘House’ dataset, ‘G’ represents the ‘Guitar’ dataset and ‘E’ indicates the ‘Elephant’ dataset. For each item in Table VI, the capital letter before ‘and’ means we train our network on corresponding datasets, and the capital letter after ‘and’ means we test the model on corresponding datasets. For example, ‘H.G.E and P’ represents that we train the network on House, Guitar and Elephant datasets, and then test on the Person dataset. Table VI exhibits four groups of experiments according to different datasets.

We first select ‘P’ as the test dataset and set ‘P’, ‘H.P.G.E’ and ‘H.G.E’ as the training dataset, respectively. When the training and the test dataset belong to the same category, the network can learn the semantic information and obtain the best performance (79.9%). Then, if we set all datasets as the training set, the performance is inferior to the result when the network is trained on a single dataset because the mixed semantic information will influence the judgment of the discriminator (72.4%). To explore the effectiveness of the semantic information deeply, we further conduct the

TABLE VII  
THE REORGANIZATION ACCURACY WHEN SELECT DIFFERENT  $n$  AND  $P$ .

	$P = 10$	$P = 100$	$P = 1000$
$n = 2$	91.9%	87.0%	83.2%
$n = 3$	85.6%	83.3%	79.0%

experiment when selecting ‘H.G.E’ as the training dataset while the test set is ‘P’. The result shows that the network can hardly obtain the correct semantic information if the training set and the test set were less correlated. Therefore, the result (66.9%) is similar to the result when we removing the GAN branch (66.4%). For example, as for the ‘Person’ category, if the generated image exchanges the arm and leg pieces, the GAN network can judge that it is irrational according to the semantic information acquired from the real ‘Person’ dataset, and give more penalties to the case. Otherwise, if the real dataset contains images coming from other datasets, the composition of ‘Person’ cannot be detected accurately. The trend of other three datasets is the same as the Person dataset (row 3 to row 8), which demonstrate the semantic information provided by the GAN branch is important.

3) *Ablation Study of Grids and Permutations*: The selection of different  $n$  and  $P$  have significant influences on the reorganization accuracy. As shown in Table VII, increasing the permutation  $P$  obviously degrades the accuracy, so as increasing the hyper-parameter  $n$ . Increasing  $P$  makes the classification more complicated, while increasing  $n$  improves the difficulty of detecting the relationships of pieces. With the increase of  $P$ , the permutations become close to each other and their features tend to be similar. It is challenging for the classification network to recognize them correctly, and therefore affecting the classification accuracy.

4) *Ablation Study of Reference Labels*: The proposed method can be considered as an improvement technique of existing methods, which can improve the reorganization accuracy of them if we select their results as our reference labels. As described in Section II, the jigsaw puzzle tasks can be specifically divided into three categories: (1) conventional methods that solve hundreds of pieces; (2) DL-based methods which generally handle  $3 \times 3$  pieces; (3) self-supervised learning methods that consider solving jigsaw puzzles as pre-text tasks. The first category methods concentrate on obtaining better performance in detecting the neighbor pieces, whereas ignoring the accuracy of reassembling the images correctly. If we consider assigning each piece to the correct place as a perfect reconstruction, the proportion of correct reconstruction images in all test images of their methods only occupies 50%-60% approximately. The proposed method belongs to the second category and considers solving jigsaw puzzles as a classification task. We define *homogeneous algorithms* to indicate the methods that can be applied to solve  $3 \times 3$  puzzles, including the second category method [12] and some first category methods [13], [19], [20], [30]. The proposed JigsawGAN is compatible with homogeneous methods to create considerable improvement, which is benefited from the GAN branch and corresponding losses.

TABLE VIII  
THE IMPROVEMENT WHEN SELECT DIFFERENT METHODS AS OUR REFERENCE LABELS.

	Gallagher2012 [19]	Son2014 [20]	Yu2015 [30]	Huroyan2020 [13]	Paumard2020 [12]
Original accuracy	62.1%	67.3%	71.0%	74.6%	60.9%
Ours	75.6%	78.4%	79.2%	81.0%	69.4%
Improvements	+13.5%	+11.1%	+7.8%	+6.4%	+8.5%

Table VIII shows the improvements of aforementioned methods if we choose their results as our reference labels, respectively. As for our results, higher reference label accuracy leads to higher reorganization accuracy because the reference labels can better guide the network to learn useful information, such as achieving 81.0% accuracy for method [13]. The improvements of [19] and [20] are significant, while the improvements of [30] and [13] are relatively smaller because there is an upper bound of the network. As for reference labels with lower accuracy, the GAN branch can promote them significantly. Moreover, it is reasonable that the improvement for [12] is small because [12] is mainly designed for their proposed dataset with eroded boundaries. Overall, no matter how the reference labels are obtained, it can be promoted by utilizing our proposed architecture.

#### F. Discussion

Solving the jigsaw puzzle with deep learning methods is a developing research area, especially when the inputs include many pieces. Conventional DL-based methods concentrate on solving puzzles with  $3 \times 3$  pieces. The upper limit of the proposed architecture is determined by the respective field. Current architecture can handle a maximum of 16 pieces. If we want to deal with images with more pieces, the classification network should be deepened with more convolutional blocks to get a larger respective field. Moreover, solving non-square puzzles is still a challenging research problem. The proposed method in the current version is constrained by the reference label and the boundary loss, which are designed to optimize straight boundaries. Applying a metric designed for matching non-square pieces to compute the boundary loss and update the reference label may be useful for solving arbitrary-size puzzles. We consider them as future works.

#### V. CONCLUSION

We have proposed JigsawGAN, a GAN-based auxiliary learning method for solving jigsaw puzzles when the prior knowledge of the initial images is unavailable. The proposed method can apply the boundary information of pieces and the semantic information of generated images to solve jigsaw puzzles more accurately. The architecture contains the classification branch and the GAN branch, which are connected by an encoder and a flow-based warp module. The GAN branch drives the encoder to generate more information features and further improves the classification branch. GAN loss and a novel boundary loss are introduced to constrain the network to focus on the semantic information and the boundary information, respectively. We have conducted comprehensive experiments to demonstrate the effectiveness of JigsawGAN.

#### REFERENCES

- [1] C. Papaodysseus, T. Panagopoulos, M. Exarhos, C. Triantafyllou, D. Fragoulis, and C. Doumas, "Contour-shape based reconstruction of fragmented, 1600 bc wall paintings," *IEEE Trans. on Signal Processing*, vol. 50, no. 6, pp. 1277–1288, 2002.
- [2] B. J. Brown, C. Toler-Franklin, D. Nehab, M. Burns, D. Dobkin, A. Vlachopoulos, C. Doumas, S. Rusinkiewicz, and T. Weyrich, "A system for high-volume acquisition and matching of fresco fragments: Reassembling theran wall paintings," *ACM Trans. Graph.*, vol. 27, no. 3, pp. 1–9, 2008.
- [3] C. Toler-Franklin, B. Brown, T. Weyrich, T. Funkhouser, and S. Rusinkiewicz, "Multi-feature matching of fresco fragments," *ACM Trans. Graph.*, vol. 29, no. 6, pp. 1–12, 2010.
- [4] R. Pintus, K. Pal, Y. Yang, T. Weyrich, E. Gobbetti, and H. Rushmeier, "A survey of geometric analysis in cultural heritage," *Computer Graphics Forum*, vol. 35, no. 1, pp. 4–31, 2016.
- [5] M. A. Marques and C. O. Freitas, "Reconstructing strip-shredded documents using color as feature matching," in *ACM Symposium on Applied Computing*, pp. 893–894, 2009.
- [6] H. Liu, S. Cao, and S. Yan, "Automated assembly of shredded pieces from multiple photos," *IEEE Trans. on Multimedia*, vol. 13, no. 5, pp. 1154–1162, 2011.
- [7] A. Deever and A. Gallagher, "Semi-automatic assembly of real cross-cut shredded documents," in *Proc. ICIP*, pp. 233–236, 2012.
- [8] W. Marande and G. Burger, "Mitochondrial dna as a genomic jigsaw puzzle," *Science*, vol. 318, no. 5849, pp. 415–415, 2007.
- [9] T. Altman, "Solving the jigsaw puzzle problem in linear time," *Applied Artificial Intelligence an International Journal*, vol. 3, no. 4, pp. 453–462, 1989.
- [10] E. D. Demaine and M. L. Demaine, "Jigsaw puzzles, edge matching, and polyomino packing: Connections and complexity," *Graphs and Combinatorics*, vol. 23, no. 1, pp. 195–208, 2007.
- [11] V. Sivapriya, M. Senthilvel, and K. Varghese, "Automatic reassembly of fragments for restoration of heritage site structures," in *International Symposium on Automation and Robotics in Construction*, vol. 35, pp. 1–7, 2018.
- [12] M.-M. Paumard, D. Picard, and H. Tabia, "Deepzle: Solving visual jigsaw puzzles with deep learning and shortest path optimization," *IEEE Trans. on Image Processing*, vol. 29, pp. 3569–3581, 2020.
- [13] V. Huroyan, G. Lerman, and H.-T. Wu, "Solving jigsaw puzzles by the graph connection laplacian," *SIAM Journal on Imaging Sciences*, vol. 13, no. 4, pp. 1717–1753, 2020.
- [14] H. Freeman and L. Garder, "Apictorial jigsaw puzzles: The computer solution of a problem in pattern recognition," *IEEE Trans. on Electronic Computers*, no. 2, pp. 118–127, 1964.
- [15] M. G. Chung, M. M. Fleck, and D. A. Forsyth, "Jigsaw puzzle solver using shape and color," in *International Conference on Signal Processing*, vol. 2, pp. 877–880, 1998.
- [16] D. Pomeranz, M. Shemesh, and O. Ben-Shahar, "A fully automated greedy square jigsaw puzzle solver," in *Proc. CVPR*, pp. 9–16, 2011.
- [17] G. M. Radack and N. I. Badler, "Jigsaw puzzle matching using a boundary-centered polar encoding," *Computer Graphics and Image Processing*, vol. 19, no. 1, pp. 1–17, 1982.
- [18] W. Kong and B. B. Kimia, "On solving 2d and 3d puzzles using curve matching," in *Proc. CVPR*, 2001.
- [19] A. C. Gallagher, "Jigsaw puzzles with pieces of unknown orientation," in *Proc. CVPR*, pp. 382–389, 2012.
- [20] K. Son, J. Hays, and D. B. Cooper, "Solving square jigsaw puzzles with loop constraints," in *Proc. ECCV*, pp. 32–46, 2014.
- [21] L. Chen, D. Cao, and Y. Liu, "A new intelligent jigsaw puzzle algorithm base on mixed similarity and symbol matrix," *International Journal of Pattern Recognition and Artificial Intelligence*, vol. 32, no. 02, p. 1859001, 2018.



- [22] D. Mondal, Y. Wang, and S. Durocher, "Robust solvers for square jigsaw puzzles," in *International Conference on Computer and Robot Vision*, pp. 249–256, 2013.
- [23] D. Sholomon, O. E. David, and N. S. Netanyahu, "Dnn-buddies: A deep neural network-based estimation metric for the jigsaw puzzle problem," in *International Conference on Artificial Neural Networks*, pp. 170–178, 2016.
- [24] L. Dery, R. Mengistu, and O. Awe, "Neural combinatorial optimization for solving jigsaw puzzles: A step towards unsupervised pre-training," 2017.
- [25] M.-M. Paumard, D. Picard, and H. Tabia, "Image reassembly combining deep learning and shortest path problem," in *Proc. ECCV*, pp. 153–167, 2018.
- [26] D. Bridger, D. Danon, and A. Tal, "Solving jigsaw puzzles with eroded boundaries," in *Proc. CVPR*, pp. 3526–3535, 2020.
- [27] A. Valada, N. Radwan, and W. Burgard, "Deep auxiliary learning for visual localization and odometry," in *IEEE International Conference on Robotics and Automation*, pp. 6939–6946, IEEE, 2018.
- [28] S. Lyu, J. Cheng, X. Wu, L. Cui, H. Chen, and C. Miao, "Auxiliary learning for relation extraction," *IEEE Trans. on Emerging Topics in Computational Intelligence*, 2020.
- [29] T. S. Cho, S. Avidan, and W. T. Freeman, "A probabilistic image jigsaw puzzle solver," in *Proc. CVPR*, pp. 183–190, 2010.
- [30] R. Yu, C. Russell, and L. Agapito, "Solving jigsaw puzzles with linear programming," *arXiv preprint arXiv:1511.04472*, 2015.
- [31] G. Paikin and A. Tal, "Solving multiple square jigsaw puzzles with missing pieces," in *Proc. CVPR*, pp. 4832–4839, 2015.
- [32] O. Vinyals, M. Fortunato, and N. Jaitly, "Pointer networks," in *Proc. NIPS*, vol. 28, pp. 2692–2700, 2015.
- [33] C. Doersch, A. Gupta, and A. A. Efros, "Unsupervised visual representation learning by context prediction," in *Proc. ICCV*, pp. 1422–1430, 2015.
- [34] M. Noroozi and P. Favaro, "Unsupervised learning of visual representations by solving jigsaw puzzles," in *Proc. ECCV*, pp. 69–84, 2016.
- [35] R. Santa Cruz, B. Fernando, A. Cherian, and S. Gould, "DeepPermNet: Visual permutation learning," in *Proc. CVPR*, pp. 3949–3957, 2017.
- [36] F. M. Carlucci, A. D'Innocente, S. Bucci, B. Caputo, and T. Tommasi, "Domain generalization by solving jigsaw puzzles," in *Proc. CVPR*, pp. 2229–2238, 2019.
- [37] R. Du, D. Chang, A. K. Bhunia, J. Xie, Z. Ma, Y.-Z. Song, and J. Guo, "Fine-grained visual classification via progressive multi-granularity training of jigsaw patches," in *Proc. ECCV*, pp. 153–168, 2020.
- [38] M.-M. Paumard, D. Picard, and H. Tabia, "Jigsaw puzzle solving using local feature co-occurrences in deep neural networks," in *Proc. ICIP*, pp. 1018–1022, 2018.
- [39] T.-Y. Lin, P. Goyal, R. Girshick, K. He, and P. Dollár, "Focal loss for dense object detection," in *Proc. ICCV*, pp. 2980–2988, 2017.
- [40] D. Li, Y. Yang, Y.-Z. Song, and T. M. Hospedales, "Deeper, broader and artier domain generalization," in *Proc. ICCV*, pp. 5542–5550, 2017.



**Ru Li** (Student Member, IEEE) received the B.E. degree in electronic information engineering from the China University of Petroleum, Qingdao, China, in 2016. She is currently pursuing the Ph.D. degree with the School of Information and Communication Engineering, University of Electronic Science and Technology of China, Chengdu, China. From July 2019 to October 2019, she visited the University of Oxford, where she worked on image enhancement. Her research interests include image processing and computer vision.



**Shuaicheng Liu** (Member, IEEE) received the B.E. degree from Sichuan University, Chengdu, China, in 2008, and the M.Sc. and Ph.D. degrees from the National University of Singapore, Singapore, in 2010 and 2014, respectively. In 2014, he joined the University of Electronic Science and Technology of China, Chengdu, where he is currently an Associate Professor with the School of Information and Communication Engineering, Institute of Image Processing. His research interests include computer vision and computer graphics.



**Guangfu Wang** received the B.E. and M.Sc. degrees from the Department of Automation, University of Electronic Science and Technology of China, Chengdu, China, in 2013 and 2016, respectively. In 2016, he worked at Alibaba Group as an Algorithm Engineer. Now he is a Researcher of Megvii. His research interests include computer vision and deep learning.



**Guanghui Liu** (Senior Member, IEEE) received the M.Sc. and Ph.D. degrees in Electronic Engineering from the University of Electronic Science and Technology of China (UESTC), Chengdu, China, in 2002 and 2005, respectively. In 2005, he joined Samsung Electronics, Seoul, South Korea, as a Senior Engineer. In 2009, he became an Associate Professor with the School of Electronics Engineering, UESTC, where he has been a Full Professor since 2014 and is currently with the School of Information and Communication Engineering. His general research interests include multimedia, remote sensing, and wireless communication. In these areas, he has authored over ten papers in refereed journals or conferences, and holds more than 60 patents (six U.S. granted patents). He served as the Publication Chair for the IEEE ISPACS-2010 and the IEEE VCIP-2016. In 2015, he received the Natural Science Award and the Science and Technology Progress Award from the Ministry of Education of China.



**Bing Zeng** (Fellow, IEEE) received the B.E. and M.Sc. degrees in electronic engineering from the University of Electronic Science and Technology of China (UESTC), Chengdu, China, in 1983 and 1986, respectively, and the Ph.D. degree in electrical engineering from the Tampere University of Technology, Tampere, Finland, in 1991. He worked as a Post-Doctoral Fellow with the University of Toronto from September 1991 to July 1992, and as a Researcher with Concordia University from August 1992 to January 1993. Then, he joined The Hong Kong University of Science and Technology (HKUST). After 20 years of service, he returned to UESTC in Summer 2013, through Chinas 1000-Talent-Scheme. At UESTC, he leads the Institute of Image Processing to work on image and video processing, 3D and multiview video technology, and visual big data. During his tenure with HKUST and UESTC, he has supervised more than 30 master's and Ph.D. students, received over 20 research grants, filed eight international patents, and published more than 250 articles. He was a recipient of the 2nd Class Natural Science Award (the first recipient) from the Ministry of Education of China in 2014, and was elected as a Fellow of the IEEE in 2016, for contributions to image and video coding. He was the General Co-Chair of the IEEE VCIP-2016, Chengdu, in November 2016. He is currently on the Editorial Board of Journal of Visual Communication and Image Representation and serves as the General Co-Chair for PCM-2017. He served as an Associate Editor for the IEEE TRANSACTIONS ON CIRCUITS AND SYSTEMS FOR VIDEO TECHNOLOGY (TCSVT) for eight years, and received the Best Associate Editor Award in 2011.