# EDGE2VEC: A HIGH QUALITY EMBEDDING
# FOR THE JIGSAW PUZZLE PROBLEM

**Daniel Rika**
Dept. of Computer Science
Bar-Ilan University
Israel, Ramat-Gan 52900
danielrika@gmail.com

**Dror Sholomon**
Dept. of Computer Science
Bar-Ilan University
Israel, Ramat-Gan 52900
dror.sholomon@gmail.com

**Eli (Omid) David**
Dept. of Computer Science
Bar-Ilan University
Israel, Ramat-Gan 52900
mail@elidavid.com

**Nathan S. Netanyahu**
Dept. of Computer Science
Bar-Ilan University
Israel, Ramat-Gan 52900
nathan@cs.biu.ac.il

## ABSTRACT

Pairwise compatibility measure (CM) is a key component in solving the jigsaw puzzle problem (JPP) and many of its recently proposed variants. With the rapid rise of deep neural networks (DNNs), a trade-off between performance (*i.e.*, accuracy) and computational efficiency has become a very significant issue. Whereas an end-to-end DNN-based CM model exhibits high performance, it becomes virtually infeasible on very large puzzles, due to its highly intensive computation. On the other hand, exploiting the concept of embeddings to alleviate significantly the computational efficiency, has resulted in degraded performance, according to recent studies. This paper derives an advanced CM model (based on modified embeddings and a new loss function, called *hard batch triplet loss*) for closing the above gap between speed and accuracy; namely a CM model that achieves SOTA results in terms of performance and efficiency combined. We evaluated our newly derived CM on three commonly used datasets, and obtained a reconstruction improvement of 5.8% and 19.5% for so-called Type-1 and Type-2 problem variants, respectively, compared to best known results due to previous CMs.

## 1 Introduction

The jigsaw puzzle problem (JPP) is an ongoing challenging task, which has been studied for decades. Although the problem plays a role in various *real-world* applications, e.g., broken pottery, ancient frescoes, shredded documents, etc., most of the effort has focused on its *synthetic* version of putting together a "perfect" digital image cropped into $N$ square non-overlapping pieces, each of size $S \times S$ pixels. The synthetic puzzle variant has served as a common testbed for many researchers, and has evolved dramatically over the years. In principle, reconstructing a puzzle requires two main phases: (1) Computation of all pairwise compatibility measures (CMs) between pieces, and (2) employment of a reconstruction algorithm based on the first phase. The above framework is not only reasonable, but has proven immensely effective in various problem instances, involving, e.g., puzzles with tens of thousands of pieces, puzzles with missing pieces and even multiple puzzles with mixed pieces. This paper focuses mainly on the first phase, i.e., finding the CM, which directly affects the quality of the resulting reconstruction.

One of the first CMs is the *sum square differences* (SSD) introduced by Cho et al. [1], which computes the sum of squared (pixel value) differences along the common edge between adjacent pieces. Pomeranz et al. [2] later proposed the *prediction-based compatibility* (PBC) measure, which tries to quantify how well a left piece edge can "predict" an adjacent right piece edge and vice versa. Specifically, they considered a 2-pixel width edge for each adjacent piece of a pair in question, and assessed the prediction by the $(L_{3/10})^{1/16}$ normalization. Next, Gallagher [3] proposed the

*Mahalanobis gradient compatibility* (MGC), which compares the difference between color gradients, rather than the difference between RGB values. Paikin et al. [4] suggested an asymmetric CM based on the $L_1$ metric (referred to as $L_1$ CM). The above traditional CMs rely solely on content information of piece edge, and they tend to be computed easily.

Under slightly more realistic conditions, however, of degraded puzzles with eroded piece boundaries, all of the above CMs become virtually useless, since edge pixels along the boundary of adjacent pair pieces lack essential (chromatic) information. A whole new set of CMs, based on *deep neural network* (DNN) architectures, has emerged, to tackle, e.g., eroded piece boundaries. These DNN-based CMs work as an end-to-end (E2E) function $g : \mathbb{R}^{S \times 2S \times C} \to \mathbb{R}$, as it inputs a pair of pieces of size $S \times S$ and $C$ color channels, and outputs a compatibility score. Since these DNN-based E2E functions arrive at a compatibility score based on the *entire* pair of pieces, they are pretty robust to pieces with eroded boundaries and possibly other forms of degradation. Paixão et al. [5] used SqueezeNet [6] and MobileNetV2 [7] as DNN-based E2E CMs to achieve state of the art reconstruction for *strip-cut* shredded documents. Rika et al. [8] designed a DNN-based E2E CM for reconstructing Portuguese tile panels with hundreds of tiles. Bridger et al. [9] used a GAN-based CM to solve puzzles with eroded boundaries, by first filling the gap between a pair of pieces by the *generator*, and then providing a compatibility score by the *discriminator*.

In terms of complexity, there are $4N(N-1)$ and $16N(N-1)$ pairwise combinations for the so-called *Type-1* and *Type-2* variants, respectively, of an $N$-piece puzzle. (Type-1 refers to unknown piece location and known orientation, whereas Type-2 refers to unknown piece location and orientation.) Obviously, using a DNN-based E2E CM, which requires millions (or even billions) of MACs for a single inference, could become extremely expensive computationally, which only worsens as the number of pieces gets larger.

To alleviate the above computational burden, a more efficient scheme was proposed by Rika *et al.* [10], called TEN, which works as follows. In contrast to a typical DNN-based E2E model, which computes the CM of a given pair explicitly (*i.e.*, by intensive processing of all pairwise pairs), TEN first represents the given pieces $k_l, k_r \in \mathbb{R}^{S \times S \times C}$ in latent space $z_l, z_r \in \mathbb{R}^d$, using two embedding DNN-based models, $f_{left}$ and $f_{right}$. These latent representations extract the most relevant information of a tile with respect to its boundary, *i.e.* $f_{left}$ extracts information from the left piece with respect to its right edge, and $f_{right}$ extracts information from the right piece with respect to its left edge. After obtaining all of the $4N$ boundary representations by $f_{left}$ and $f_{right}$, TEN employs a simple and much faster distance measure between embedding vectors, $\mathcal{D} : \mathbb{R}^d \times \mathbb{R}^d \to \mathbb{R}$, to compute all of the $16N(N-1)$ pairwise CMs. (Euclidean distance was reported to work best.) Despite of its significantly improved efficiency (relatively to a DNN-based E2E), TEN was even slightly more accurate than classical CMs. However, its lower accuracy compared to that of a powerful DNN-based E2E, poses a substantial gap in performance to be closed, limiting TEN to a mere concept.

This paper introduces Edge2Vec, a fast, robust, and accurate CM scheme. As opposed to TEN, Edge2Vec eliminates the notion of twin embedding networks, by using the same model to extract features from both the left and right pieces due to a simple *horizontal flipping* technique. Furthermore, we demonstrate the superior performance of Edge2Vec relatively to a DNN-based E2E CM in the domain of *eroded boundaries*, which serves as a testbed for real-world settings. The main contributions of our work are as follows:

- Closed entirely the performance gap from TEN to a DNN-based E2E CM, without compromising computational efficiency.

- Proposed a grouping technique to reduce footprint from 9.6M parameters to only 2.1M (*i.e.* $\times 4.6$ less parameters) without affecting the performance.

- Proposed a new variant of triplet loss, called *hard-batch-triplet loss* for high quality edge representation.

## 2 Edge2Vec Architecture

We show how to improve TEN in two aspects. First, we eliminate the redundancy of twin embedding models (instead of using a single one), and secondly we propose a technique to handle efficiently the last projection layer that contains the vast majority of the parameters of the model.

### 2.1 From Twins to a Single Embedding Network

Given a pair of pieces $(k_l, k_r)$, where $k_l, k_r \in \mathbb{R}^{S \times S \times C}$ are the left and right piece in the pair, respectively, TEN first embeds $k_l, k_r$ in latent vectors $z_l, z_r \in \mathbb{R}^d$, and then computes the Euclidean distance between these embeddings, as the CM of the pieces. Since $z_l$ represents the information of $k_l$ with respect to its right edge, and $z_r$ represents the information of $k_r$ with respect to its left edge, two different models, $f_{left}, f_{right}$ were proposed. We suspect that
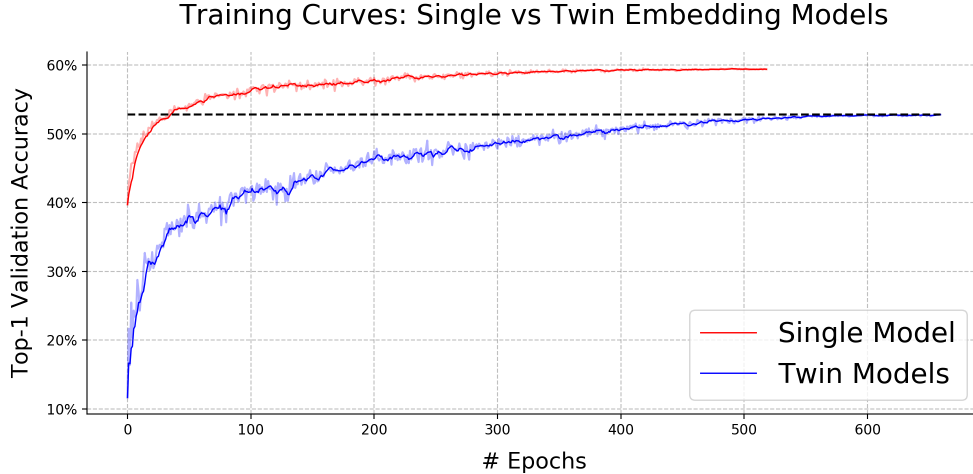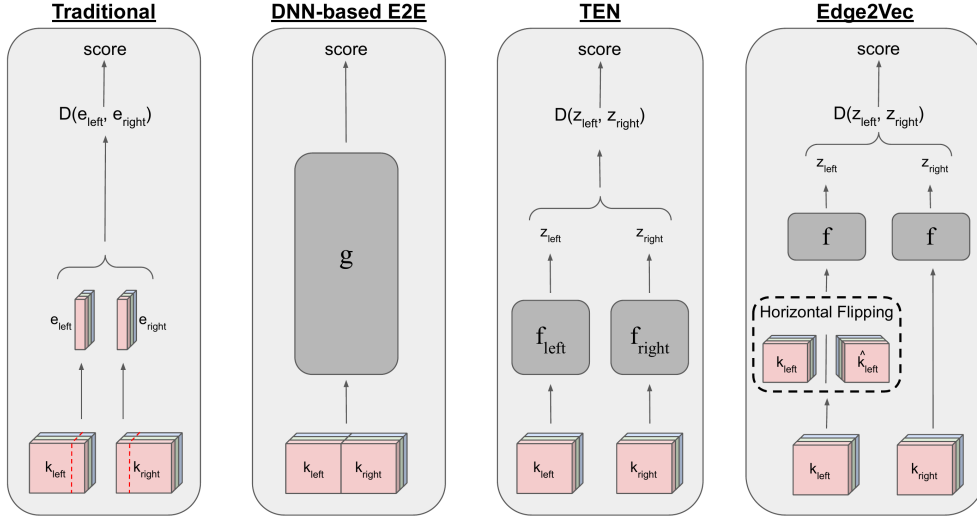
Figure 1: Plots of Top-1 validation accuracy; single embedding model takes only 38 epochs to reach TEN's highest level, obtained after 600 epochs (see black dashed line); also, single model improves Top-1 validation accuracy by 6.7%, while using only half the number of TEN's weights (*i.e.*, 2.5M parameters instead of 5.1M).

training and employing such twin models might prove too hard and unstable a task, since each model needs to adjust millions of parameters for generating embedding vectors of the same distribution as the vectors of its companion model with extreme precision. Instead, applying *horizontal flipping* (*i.e.*, mirroring) to $k_l$ and $k_r$, thereby getting $\hat{k}_l$ and $\hat{k}_r$, respectively, should imply naturally that $f_{left}(k_l) = f_{right}(\hat{k}_l)$ and $f_{right}(k_r) = f_{left}(\hat{k}_r)$. For example, $\hat{k}_l$ contains the same information as $k_l$, only that due to the flipping, the right edge of $k_l$ becomes the left edge of $\hat{k}_l$. (A similar observation holds for $k_r$.) Thus, the redundancy of TEN's $f_{left}$ and $f_{right}$ models should be clear. Instead, one can represent both left and right edges of a pair of pieces, due to *horizontal flipping*, by a single model, *e.g.*, $z_l = f_{right}(\hat{k}_l)$, $z_r = f_{right}(k_r)$. Shifting away from twin embedding networks to a single network, according to the above explanation, is illustrated in Figure 2b. We will refer to this proposed scheme as Edge2Vec.
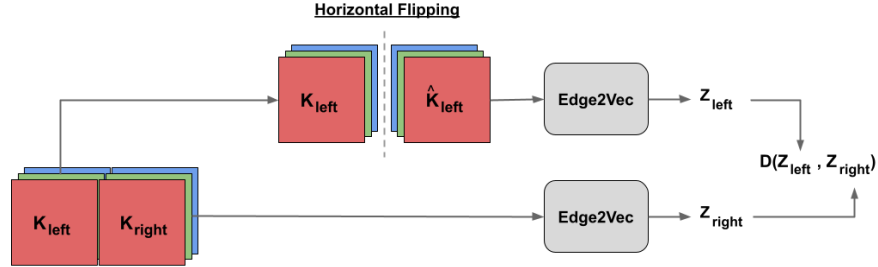
To justify the use of a single model instead of twins, we trained Edge2Vec on the same configurations and datasets described in [10]. Training plots of Top-1 validation (to be defined) for both models are depicted in Figure 1. In addition to its faster convergence – Edge2Vec needed 15.8 times less epochs to reach TEN's highest Top-1 level – it also improved TEN's Top-1 accuracy on the test set by 5.7% and 5.8% for Type-1 and Type-2, respectively.

## 2.2   Grouping

As previously mentioned, Edge2Vec originates from TEN's architecture, which consists of four convolutional layers corresponding to $64, 128, 256$, and $512$ output channels, using a kernel size of $3 \times 3$ and the nonlinear ReLU activation function [11] after each layer. Max pooling was applied after the second and third convolutional layers to reduce spatial dimensionality, and a fully-connected (FC) layer was finally applied to get the embedding vector. It is well known that the largest number of connections (or weights associated with them) lie between the last convolutional layer and the embedding vector. Specifically, this number of parameters is given by $C_o * S_o^2 * d$, where $d$ is the embedding dimension, $C_o = 512$ is the number of output channels, and $S_o = 7$ output spatial size (assuming $28 \times 28$ input images). To mitigate such parameter explosion, we found that splitting the FC layer into $G$ groups of smaller FC sub-layers (for certain $G$ values) could downsize the number of parameters to $(C_o * S_o^2 * d)/G$, without suffering any degradation in performance. Before flattening the output of the last convolutional layer $x$, we split it along the channel axis into $G$ sub-channels, $\{x_1, ..., x_G\}$. Having defined also $G$ FC sub-layers: $FC_i : \mathbb{R}^{(C_o/G)*S_o^2} \to \mathbb{R}^{d/G}$, we then assign each $x_i$ to its own $FC_i$ sub-layer, to obtain $G$ sub-embeddings $\{z_1, ..., z_G\}$ whose concatenation provides the final embedding vector $z$. Namely, instead of a single FC layer, we use $G$ FC sub-layers, each of which is associated with only $1/G^2$ the number of parameters. Thus, the above grouping decreases the total number of parameters associated with the full FC layer by a factor of $G$.

3

(a)



(b)

Figure 2: Illustration of pairwise CM schemes: (a) Block diagrams of traditional, E2E, TEN, and Edge2Vec schemes, where $\mathcal{D}$ denotes any distance metric, *e.g.*, Euclidean distance, in our case ; (b) *horizontal flipping* used by Edge2Vec.

# 3 Training Edge2Vec

In order to train Edge2Vec, we started with the *triplet-loss*:

$$\mathcal{L}_{\text{Triplet}} = \max(0, \mathcal{D}(z_a, z_p) - \mathcal{D}(z_a, z_n) + \gamma) \tag{1}$$

where $\mathcal{D}$ denotes Euclidean distance, and $\gamma$ is a hyperparameter set to $1$. We discuss here the choice of a negative embedding $z_n$, which competes with the positive embedding $z_p$ for the anchor $z_a$. During training, a triplet $(z_a, z_p, z_n)$ is picked randomly from the same image (which is also selected randomly from a pool of training images). It is reasonable to assume that picking $z_n$ from the same image would improve the model's discrimination capability, since a given image tends to be composed of certain color and texture distributions. However, this is only one possibility for selecting $z_n$; other strategies might push the model even harder to extract better embedding features, by picking more challenging $z_n$'s from other training set images, to distinguish from.

In addition, during each training iteration, a batch of $B$ triplet samples are picked randomly and stacked together to get more informative gradients for the next optimization step of the model. Hence, we propose an efficient way which takes advantage of batching for finding harder $z_n$'s. Our proposed selection of *hard batch triplets* (HBT) is explained, in detail, in the following subsection.

## 3.1 Hard Batch Triplet (HBT) Selection

Let us denote each triplet in the batch $(z_a^b, z_p^b, z_n^b)$, where $b \in \{1, .., B\}$ is the batch index. Instead of looking at each triplet $(z_a^b, z_p^b, z_n^b)$ in the batch separately, we can combine the triplets to obtain a more diverse set of negative
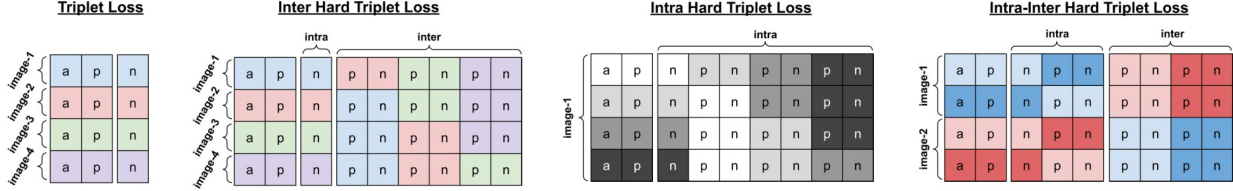
4

Figure 3: Intra-inter triplet batch with 4 samples. Left to right: Original triplet: every anchor has its own negative; Inter hard batch triplet: each triplet is sampled from a different image in the training set; Intra hard batch triplet: all batch triplets are sampled from same image; Intra-inter hard batch triplet: intra-inter ratio of 1:1

candidates, by sharing both the positives and negatives between all of them. In other words, every anchor $z_a^b$ will have $2B$ candidates $\tilde{\mathcal{Z}}^b = \{z_p^1, z_p^2, ..., z_p^B, z_n^1, z_n^2, ..., z_n^B\}$, where $z_p^b$ is the only positive embedding, and all other $2B - 1$ embeddings are negative candidates. Thus, instead of using $z_n^b$ as the negative for anchor $z_a^b$, we can calculate the distances between $z_a^b$ to all of the $2B - 1$ negative candidates in $\tilde{\mathcal{Z}}^b$, and choose the "hardest" one in the batch. Put formally, for every anchor $z_a^b$ in the batch, its new negative is defined as

$$z_{n*}^b = \underset{z \in \tilde{\mathcal{Z}}^b \setminus \{z_p^b\}}{\operatorname{argmin}} \mathcal{D}(z_a^b, z) \tag{2}$$

for which we then apply the regular *triplet loss*:

$$\mathcal{L}_{\text{HBT}}(z_a^b, \tilde{\mathcal{Z}}^b) = \max(0, \mathcal{D}(z_a^b, z_p^b) - \mathcal{D}(z_a^b, z_{n*}^b) + \gamma) \tag{3}$$

Experimenting with our *HBT loss function* during training, we observed initially a boost in Top-1 performance. However, this improvement quickly tapered off on the validation set, as a result of *overfitting*. This was attributed to the phenomenon where the embedding vectors were dominated by only a few features with very large values, such that most of the features were irrelevant in computing the Euclidean distance between embedding vectors. To force the model to generate more balanced embeddings, we used the $L_2$ regularization described below.

## 3.2 $L_2$ Embedding Regularization

Realizing that the HBT loss function should provide at least the same results as the standard triplet loss, since it selects the original negative $z_n^b$ or a harder one in the batch, we just needed to overcome the overfitting mentioned in the previous subsection. This was done by encouraging Edge2Vec to take advantage of as many features as possible, while avoiding very large values, due to the following $L_2$ regularization:

$$\mathcal{R}_{L_2} = \sqrt{\frac{1}{3Bd} \sum_{b=1}^{B} \sum_{i=1}^{d} (z_a^b[i])^2 + (z_p^b[i])^2 + (z_n^b[i])^2} \tag{4}$$

The final loss function used was thus

$$\mathcal{L} = \mathcal{L}_{\text{HBT}} + \lambda \mathcal{R}_{L_2} \tag{5}$$

where $\lambda = 1$ was found to work best in our case.

## 3.3 Intra- and Inter-Image Sampling

The proposed set of HBTs is composed of $B$ triplets, each from a different image (unless $B$ is greater than the number of images in the training set). Thus, the hard negatives per batch sample are *inter* hard negatives, while the only negative from the same image as the anchor ($z_a^b$), is the one from the original triplet, *i.e.*, $z_n^b$. However, it is more likely that distinguishing negatives from other images (*i.e.*, *inter negatives*) will be less difficult than negatives sampled from the anchor's image (*i.e.*, *intra negatives*). To maintain a wide range of possibilities, we control the ratio between the numbers of intra and inter samples in a batch, which directly reflects on the ratio between the numbers of intra and inter hard negatives. In principle, we experienced that an intra-inter sampling ratio of 1:1 works the best, since *intra sampling* forces rather discriminative embeddings, while *inter sampling* keeps the training more stable due to a more diverse set of images in the batch. Figure 3 illustrates all three combinations of HBTs and also a regular triplet for reference.

Note that using intra sampling increases the chance that some negatives in the batch are the positives of other anchors. Hence, before applying HBT loss to the batch, we need to detect all such instances and remove those potential hard negatives for any relevant anchors in the batch.

# 4 Experimental Results

We have experimented deliberately with the more challenging problem of degraded puzzles containing pieces with *eroded boundaries*, as synthetic/perfect puzzles can be solved rather easily. In particular, while classical CM methods perform very well on the latter type of puzzles, they fail miserably on the former.

Specifically, we experimented with eroded 1-pixel frame piece boundaries of $28 \times 28$ pieces. Thus, for every piece, we set the outer frame pixels to 0, and retained the original content of the remaining $26 \times 26$ pixels. The core CNN of our model is identical to TEN's, while the last FC layer was modified according to the *grouping* architecture explained previously. For our datasets, we used the 800 training images of DIV2K [12] downscaled by a factor of 2 (due to bicubic interpolation), which is the same training set used for TEN. For additional evaluation and to avoid bias as much as possible, we augmented our test set by adding two commonly used datasets. In addition to the 100 images from DIV2K's validation set, we added 100 images from the PIRM dataset [13] and 20 images from the MIT dataset [1], Altogether, we obtained a test set of 220 puzzles, each containing 300–1000 pieces. For training, we used the Adam optimizer [14], a learning rate of 1E-4 (with a decay factor of 0.9, in case the loss did not decrease for 5 epochs), a batch size of 1024, and an epoch of 5000 iterations. Also, we used Euclidean distance as our distance metric $\mathcal{D}$ and embeddings of size $d = 320$. To ensure symmetric pairwise CMs, we apply post-processing according to Eqs. 6, 7, as was proposed in [10]. The training was done on a machine with 64 CPU cores, 128GB of RAM and 8 GPUs, each with 11GB memory.

$$C'(k_i, k_j) = \frac{C(k_i, k_j) - min(C(k_i, *))}{max(C(k_i, *)) - min(C(k_i, *))} \tag{6}$$

$$C''(k_i, k_j) = C''(k_j, k_i) = \frac{C'(k_i, k_j) + C'(k_j, k_i)}{2} \tag{7}$$

## 4.1 Grouping: Smaller while Better

As noted before, most of the model's parameters (or connection weights) are associated with the final FC layer. Specifically, the number of weights, mapping the spatial CNN features onto a vector $z$ in latent space $\mathbb{R}^d$, depends on the embedding size $d$. Since no nonlinear function is applied after the projection layer, each component of $z$ is a weighted sum of all the spatial features. We suggest that even if all of the channels in the final layer of spatial features are essential for our embedding $z$, not all of them necessarily contribute to all of $z$'s components. Moreover, if not all of the channels should be associated essentially with certain components of $z$, they could actually prove detrimental to the quality of these components. Motivated by this insight, we conducted experiments to assess the extent to which the dependencies of each component in $z$ on all channel features $C$ could be down scaled to only a subset of $C/G$ channels. Having experimented with $C = 512$, $d = 120$, and $G \in \{1, 2, 4, 8\}$, we noticed that grouping does not degrade the model's performance; to the contrary, it only resulted in slightly better results (see Table 1). Moreover, splitting Edge2Vec into 16 groups decreased the total number of parameters from 9.6M to only 2.1M (*i.e.*, a reduction by a factor of almost 4.6), without any degradation in performance.

## 4.2 HBT Loss and $L_2$ Regularization

As noted before, we should not have observed a worse performance of our model using the HBT loss function (instead of the standard triplet loss), as its basic purpose is to find only harder negatives. To explain the overfitting phenomenon encountered initially, we examined the following parameters of embeddings generated by Edge2Vec: (1) The maximal absolute value of an embedding feature, and (2) the fraction of very large embedding components with respect to the embedding size $d$. Collecting the various embedding values from the validation set, we found that on average, 22.5% of the absolute values of an embedding's components were in the range 2–8, while all the rest were considerably smaller. And since Euclidean distance is used as the distance metric $\mathcal{D}$, small sets of very large elements could well dominate the outcome of the CM computation. So although HBT loss should force the model to distinguish successfully hard negatives, the learning process has been impaired due to a few dominant sets of values, which affect mostly the gradients' update, while ignoring the remaining 77.5% embedding features containing smaller values. Thus, to mitigate the use of very large values in the embedding vector, we added the $L_2$ regularization of Eq. 4 to the final loss function, and controlled its influence by the parameter $\lambda$. Figure 4b depicts nicely the positive impact of $L_2$ regularization on the Top-1 accuracy for the test set, which improved from 58.6% to 64.5% for $\lambda = 0$ and $\lambda = 1$, respectively.

To demonstrate the effect of $L_2$ regularization on the robustness, we gradually masked out the embeddings from the highest absolute value to the lowest, until 50% of an embedding vector were removed. Next, for each level of masking we calculated the Top-1 accuracy on the testset and measured the relative accuracy of the masked embeddings compared to the original unmasked embeddings (which is the original Top-1 accuracy on the testset). Figure 4a depicts the two

Table 1: Top-1 validation accuracy of all ablation experiments ordered "chronologically" from top to bottom; each ablation study contains best hyper parameters found in previous experiments, *e.g.*, if $G$=8 was found for best performance (under *grouping*), all subsequent experiments were conducted with $G$=8.

| Grouping | | | |
|---|---|---|---|
| $G$=1 | $G$=2 | $G$=4 | $G$=8 |
| 65.3% | 65.7% | 65.9% | 66.3% |

| Embedding Size | | | |
|---|---|---|---|
| $d$=40 | $d$=80 | $d$=160 | $d$=320 |
| 64.6% | 66.3% | 67.1% | 67.5% |

| Batch-Size | | | | |
|---|---|---|---|---|
| $B$=64 | $B$=128 | $B$=256 | $B$=512 | $B$=1024 |
| 66.6% | 67.1% | 67.5% | 68.6% | 69.6% |

| Intra-Inter Hard Batch Triplets | | |
|---|---|---|
| Inter | Intra-Inter 1:1 | Intra |
| 69.6% | 72.9% | 69.8% |

Top-1 accuracies curves, with and without the $L_2$ regularization. (All of the other training hyperparameters were identical for the two experiments.) The above plots imply that Edge2Vec spreads more evenly the embedding features with $L_2$ regularization than without it. (The latter confines most of the embedding energy to only a few features.)
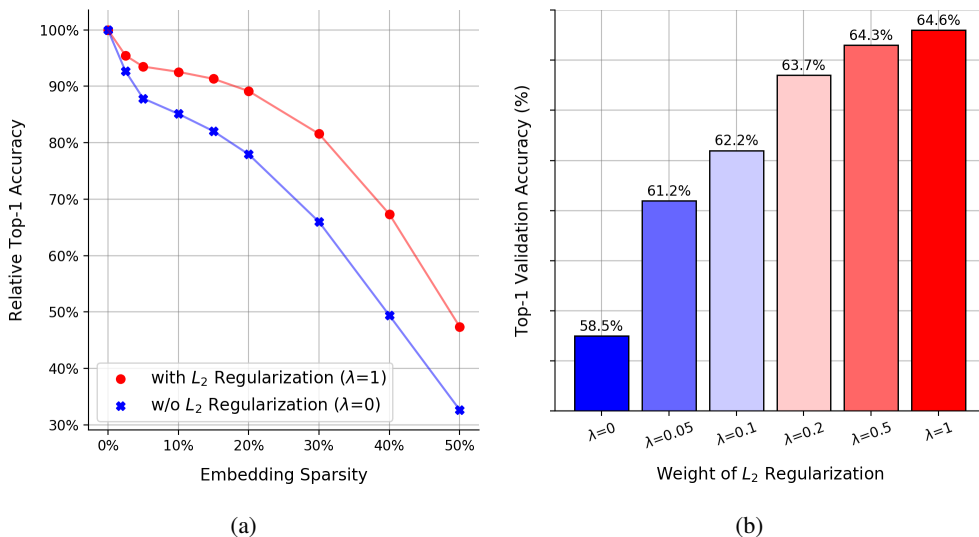


(a)                                                    (b)

Figure 4: Impact of $L_2$ regularization: (a) Influence of embedding features spread more evenly, *i.e.*, model trained with $L_2$ regularization preserves almost 90% of original Top-1 accuracy with 20% most significant features deleted, whereas model without $L_2$ regularization drops to $\approx 80\%$ of its original Top-1 accuracy with same amount of sparsity; (b) impact of $\lambda$ on final performance: 6% improvement in Top-1 validation accuracy from no regularization ($\lambda = 0$) to full regularization ($\lambda = 1$).

Figure 5 demonstrates the combined effect of our HBT loss and $L_2$ regularization on an 864-piece puzzle from the DIV2K dataset. Specifically, it demonstrates the dramatic improvement, due to the use of HBT, of the discrimination (in terms of CM scores) between correct neighboring pieces (shown in yellow on the matrix super diagonal, with respect to corresponding anchors on the main diagonal) versus all other candidates (shown in dark).
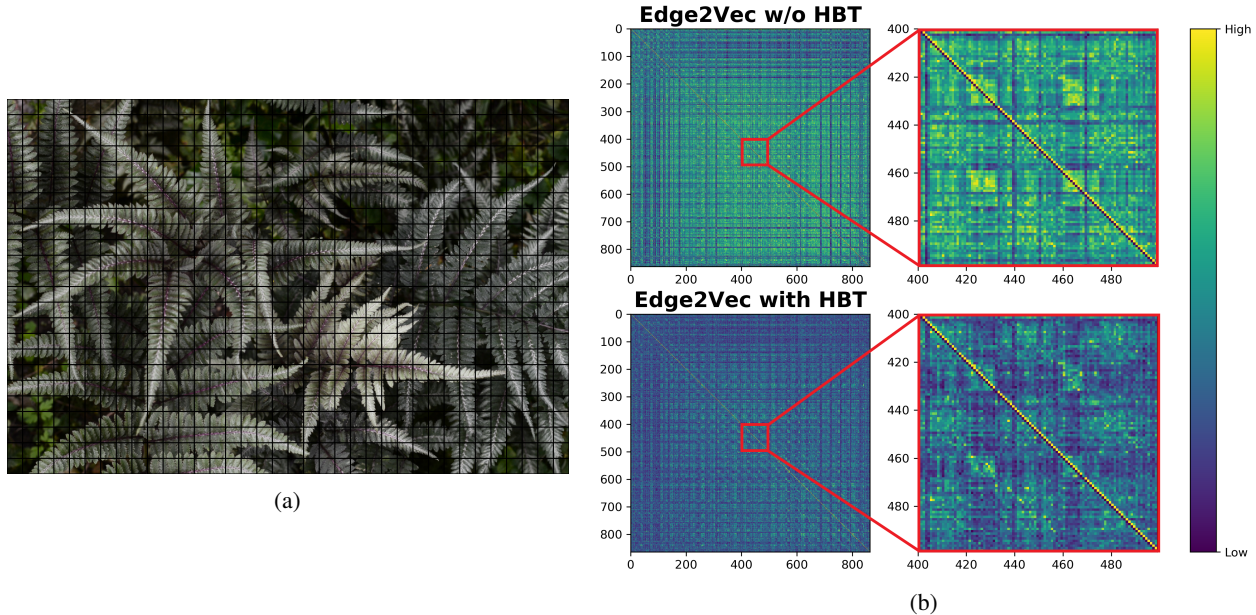
Figure 5: Distance map visualization: (a) 864-piece puzzle from DIV2K testset; (b) distance map of Edge2Vec without (top)/with (bottom) HBT loss; rows represent an anchor piece $N$, and columns represent the examined piece with respect to *right* edge of the anchor; unless anchor is located at the rightmost column of puzzle, its correct neighboring piece on the *right* is $N + 1$. Namely, compatibility scores along super diagonal of CM matrix should be largest (in yellow), while the rest should be (much) smaller (purple); $100 \times 100$ zoom-in maps clearly demonstrate enhanced discrimination with HBT loss.

## 4.3 Dimensionality and Batch Size

**Dimensionality.** The core of this work is to represent a piece with respect to one of its edges by a vector $z$ in the latent space $\mathbb{R}^d$. The vector dimensionality $d$ is the number of features that Edge2Vec works with, and could accidentally become a performance bottleneck when setting it too low. Thus, during training we experimented with various $d$ values, *i.e.*, $d \in \{40, 80, 160, 320\}$. As indicated by Table 1, it is desirable to set $d \geq 160$; the best performance (through our experiments) was obtained for $d = 320$.

**Batch-size.** It is well known that a larger batch size could assist in the training, as the gradients become more stable. In our case, the batch size highly affects our proposed HBT loss function, since more samples in the batch can potentially generate harder negatives as explained before. Hence, we tested Edge2Vec on a variety of batch sizes as shown in Table 1. The results reveal that typically a larger batch size causes Edge2Vec to extract more discriminative features, which results in better performance.

## 4.4 Performance Evaluation

After picking the best hyperparameters for our Edge2Vec model, according to the empirical ablation study reported, we ran a comparative performance evaluation against previous CM models. Before analyzing, in more detail, our experimental evaluation, we note that a CM model could be characterized by the following three main attributes: (1) Top-1 accuracy, (2) *footprint* and computational complexity (relevant to all DNN-based methods), and (3) actual inference time. Figure 6 below clearly indicates that the best trade-off with respect to the above criteria is achieved by Edge2Vec, compared to all of the classical and DNN-based CMs evaluated.

### 4.4.1 Top-1 Accuracy

Recall, for an $N$-piece puzzle, a piece edge has $N - 1$ or $4(N - 1)$ matching candidates for Type-1 or Type-2, respectively. We say that a piece edge is Top-1 if its most compatible edge candidate (according to a CM model in question) is its very adjacent edge in the original puzzle, *i.e.*, its "ground truth" neighboring edge. Following the above Top-1 definition, Top-1 accuracy is the percentage of piece edges that are Top-1 according to a specific CM algorithm. (For a set of multiple puzzles, we can regard Top-1 accuracy as the average Top-1 accuracy over all puzzles in the

dataset.) The higher the Top-1 accuracy for a given CM module, the better its performance. Ideally, Top-1 accuracy should approach $100\%$.

We compared our Edge2Vec to various classical and DNN-based E2E CMs, and found that our method achieved higher Top-1 accuracies than those obtained by the previous embedding-based TEN-L (by 12.2% and 15.5% for Type-1 and Type-2, respectively). Moreover, it also outperformed the highly-intensive DNN-based E2E method, as summarized in Table 4.

### 4.4.2 Footprint and Computational Complexity

The notion of *zero emission* has been gaining recently more attention, as to the operation volume of NN-based models, in terms of footprint and computational complexity. Thus, we compared naturally Edge2Vec to previous NN-based CMs also with respect to these terms. Specifically, footprint is measured by a model's number of parameters, and computational complexity is measured by the number of *multiply–accumulates* (MACs) per a single pair of pieces, or per an entire $N$-piece puzzle (for a total of $16N^2$ pairs). The number of parameters and the number of MACs per single pair are provided typically by a library in use, whereas the number of MACs for an $N$-piece puzzle is implied by the nature of a model in question (see below).

**Number of MACs for CM via DNN-based E2E:** Given a pair of pieces, each of size $28 \times 28$ pixels, the input size of a DNN-based E2E CM is (28, 56, 3), and the number of MACs is just the number of multiplications and additions during a forward pass. Thus, the CM complexity of a DNN-based E2E for an $N$-piece puzzle is the complexity per a single pair multiplied by $16N^2$ (*i.e.*, multiplied by the number of pair combinations).

**Number of MACs for CM via DNN-based with Embeddings:** To compute, in this case, the CM of a single pair of pieces $(k_{right}, k_{left})$, we first represent them in latent space, due to some NN-based model (*e.g.*, TEN, TEN-L or Edge2Vec), obtaining $(z_{right}, z_{left})$. We then calculate the Euclidean distance between these embeddings. Since the execution of NN-based models for each pair piece is much more intensive than the computation of the above Euclidean distance, we consider, essentially, only the number of MACs involving the projection of the two pieces to embedding space via an embedding-based model in question. In other words, this should be multiplied, essentially, just by the total number of embeddings, to obtain the (approximate) number of MACs for an entire $N$-piece puzzle. Exploiting the notion of our embedding-based models, we need only consider, therefore, $4N$ embeddings with respect to all (potential) boundaries of each piece, *i.e.*, multiply by a factor of $8N$ embedding operations (as each piece embedding needs to be considered twice, for the piece itself and for its flipped version). See Table 2 for specific results. (As noted, the number of MACs neglects the $16N^2$ Euclidean distance operations, as their execution is considerably faster than the embedding computations themselves.)

Table 2: Footprint and computational complexity. Our Edge2Vec CM has a competitive footprint and extremely low complexity; note the linear dependence of # MACs wrt $N$ for embedding-based models, in contrast to its quadratic dependence wrt $N$ for DNN-based E2E; although TEN-L is substantially bigger than other NN-based CMs (being an ensemble of TEN), it remains very efficient wrt DNN-based E2E CMs, for the entire puzzle, due to embedding.

| Model | # Params (M) | # MACs (G) | |
| --- | --- | --- | --- |
| | | Single pair | $N$-piece puzzle |
| **TEN** | 5.1 | 0.35 | $1.4N$ |
| **TEN-L** | 20.4 | 1.4 | $5.6N$ |
| **DNN-based E2E** | 1.6 | 0.35 | $5.6N^2$ |
| **Edge2Vec (ours)** | 2.1 | 0.35 | $1.4N$ |

### 4.4.3 Inference Time

Edge2Vec retains, essentially, TEN's faster running times, due to our novel use of embeddings as an intermediate step in the CM computation pipeline. See Table 3 for wall-clock times of all NN-based CMs, which compute all $16N^2$ compatibility scores for different puzzle sizes. The running time of Edge2Vec is between that of TEN and TEN-L, *i.e.*, much faster than that of the intensive DNN-based E2E architecture. (TEN's slightly faster inference times could be attributed to its smaller embedding dimension, *i.e.*, 40, instead of 320 used by Edge2Vec.)

Table 3: Wall-clock inference times (**in seconds**) of our proposed Edge2Vec, compared to previous embedding-based CMs and DNN-based E2E; Edge2Vec remains extremely efficient compared to E2E mode.

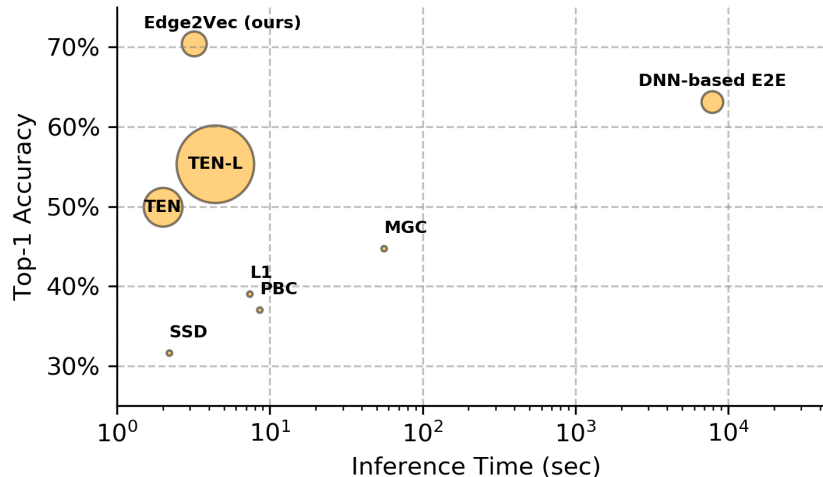| # Pieces | TEN | Edge2Vec (ours) | TEN-L | DNN-based E2E |
|---|---|---|---|---|
| **100** | $7.0 \times 10^{-2}$ | $1.4 \times 10^{-1}$ | $1.8 \times 10^{-1}$ | $3.2 \times 10^{1}$ |
| **200** | $1.4 \times 10^{-1}$ | $2.6 \times 10^{-1}$ | $3.9 \times 10^{-1}$ | $1.2 \times 10^{2}$ |
| **400** | $3.6 \times 10^{-1}$ | $5.2 \times 10^{-1}$ | $8.9 \times 10^{-1}$ | $4.9 \times 10^{2}$ |
| **800** | $1.0 \times 10^{0}$ | $1.5 \times 10^{0}$ | $2.2 \times 10^{0}$ | $1.9 \times 10^{3}$ |
| **1600** | $3.2 \times 10^{0}$ | $4.2 \times 10^{0}$ | $6.2 \times 10^{0}$ | $7.9 \times 10^{3}$ |
| **3200** | $1.1 \times 10^{1}$ | $1.4 \times 10^{1}$ | $1.9 \times 10^{1}$ | $3.1 \times 10^{4}$ |



Figure 6: Trade-off illustration for all evaluated CMs on puzzles with *eroded boundaries*, regarding: (1) Average Top-1 accuracy ($y$-axis) for Type-2 on the 3 testsets; (2) model footprint (depicted by bubble size); and (3) inference (wall-clock) time ($x$-axis) for computing all $16N(N-1)$ scores of 1,600-piece puzzle; although TEN-like model has much larger footprint than E2E (because of larger last fully-connected layer and use of two identical models (rather than one)), it is much faster because of $8N$ executions (instead of $16N^2$); tiny bubbles depict classical CM schemes, which lack footprint. Our Edge2Vec CM provides best trade-off wrt above criteria.

### 4.4.4 Reconstruction Results

Although Top-1 accuracy is a useful metric for evaluating the quality of a CM algorithm, we recall that CM's functional role is to assist in the eventual reconstruction of a given puzzle, as accurately as possible. Hence, in addition to Top-1 accuracy, we also conducted a reconstruction comparison between all of the CMs evaluated. For completeness, we chose two reconstruction algorithms: Gallagher's reconstruction algorithm [3] and a genetic algorithm (GA)-based solver [8]. We decided to work with these two methods, since Gallagher's solver is publicly available and commonly used in many other jigsaw puzzle studies. However, since Gallagher's solver is a greedy algorithm, which could easily get trapped in local optima, we also tested the CMs in conjunction with the above GA-based solver, which is a global optimization method.

In employing Gallagher's greedy reconstruction algorithm to supplement any of the evaluated CMs, we used his rescaling technique for the pairwise CMs

$$\hat{C}(k_i, k_j) = \frac{C(k_i, k_j)}{Second(C(k_i, *))} \tag{8}$$

where $Second(C(k_i, *))$ denotes the second most compatible value for anchor $k_i$, before introducing their scores into the reconstruction phase. The goal of this operation, is to give a degree of confidence to the preliminary pairwise CM values, which can help a lot in a greedy reconstruction algorithm like Gallagher's. Table 4 contains all of the reconstruction results, measured by the *neighbor accuracy* criterion, for Type-1 and Type-2, using Gallagher's algorithm and the GA-based solver. It can be seen that our proposed Edge2Vec model improved the reconstruction results using

Table 4: Comparative evaluation of various CMs for Type-1 (top) and Type-2 (bottom), with respect to Top-1 accuracy (left) and neighbor accuracy due to Gallagher's reconstruction (middle) and GA-based solver (right); average accuracy (over 10 runs for each puzzle) reported for latter reconstruction. Edge2Vec exhibits superior performance to all other methods (in nearly all cases); best results appear in bold.

**TYPE-1**

| Method | Top-1 Accuracy | | | Gallagher's solver | | | GA-based solver | | |
|---|---|---|---|---|---|---|---|---|---|
| | DIV2K | PIRM | MIT | DIV2K | PIRM | MIT | DIV2K | PIRM | MIT |
| SSD | 37.4% | 41% | 42% | 28.7% | 32.4% | 34.5% | 27.9% | 34% | 32.4% |
| PBC | 46.4% | 49.1% | 43.7% | 43.7% | 47.1% | 36.4% | 49.3% | 53.2% | 44.7% |
| $L_1$ | 47% | 49.8% | 46.6% | 42.2% | 46% | 41.3% | 43.6% | 50.9% | 43.2% |
| MGC | 56.3% | 59.5% | 53.5% | 54.7% | 60.2% | 49.4% | 64.4% | 71.4% | 58.2% |
| TEN | 59.4% | 61% | 55.2% | 52.9% | 53.1% | 44.9% | 67.5% | 71.5% | 61.2% |
| TEN-L | 64.6% | 65.4% | 59.1% | 60.7% | 59.8% | 53.4% | 72.9% | 76.8% | 65.1% |
| DNN-based E2E | 71.4% | 72.6% | 69.4% | 75.4% | 76.1% | 69% | 82.5% | 90.8% | **87.3%** |
| Edge2Vec | **76.4%** | **77.8%** | **71.3%** | **80.2%** | **82.3%** | **75.3%** | **87.4%** | **92.1%** | 86.5% |

**TYPE-2**

| Method | Top-1 Accuracy | | | Gallagher's solver | | | GA-based solver | | |
|---|---|---|---|---|---|---|---|---|---|
| | DIV2K | PIRM | MIT | DIV2K | PIRM | MIT | DIV2K | PIRM | MIT |
| SSD | 29% | 32.2% | 33.5% | 7.3% | 9.8% | 10.3% | 16.3% | 18.4% | 18.3% |
| PBC | 37.5% | 39.3% | 34.1% | 13.8% | 16.5% | 9.5% | 26.1% | 30.5% | 20.9% |
| $L_1$ | 38.7% | 40.7% | 37.6% | 12.2% | 15% | 10.7% | 25.6% | 28.2% | 22% |
| MGC | 45.4% | 47.5% | 41.1% | 19.7% | 22.7% | 12.8% | 40.2% | 47.5% | 32.3% |
| TEN | 50.5% | 52.5% | 46.6% | 21.8% | 23.4% | 16.6% | 46.9% | 50.9% | 39.6% |
| TEN-L | 56.8% | 57.5% | 51.5% | 30.2% | 33.6% | 22.7% | 55.5% | 59.4% | 49.1% |
| DNN-based E2E | 64.1% | 64.2% | 60.9% | 48.5% | 49.9% | 39% | 73.3% | 79.7% | 74.5% |
| Edge2Vec | **71.4%** | **74.2%** | **66.9%** | **65.2%** | **70.7%** | **60%** | **82.3%** | **89%** | **81.2%** |

DNN-based E2E with both reconstruction algorithms. Specifically, using Edge2Vec with Gallagher's solver achieved an average improvement across all three test sets of 5.8% and 19.5% for Type-1 and Type-2, respectively. Using Edge2Vec in conjunction with the GA-based solver improved the average reconstruction rates of Type-1 and Type-2 by 1.8% and 8.3%, respectively.

### 4.4.5 Comparison to State-of-the-Art Model

for the completeness of the experiments, we compared Edge2Vec to the previous SOTA method by Bridger et al [9]. As their model requires $64 \times 64$ pieces, we trained our Edge2Vec on the same piece size with similar erosion (1 eroded pixel of a $28 \times 28$ piece is equivalent to about 2 eroded pixels of a $64 \times 64$ piece). We then used Gallagher's open source solver [3] for comparative evaluation, due to a lack of code availability of Paikin and Tal's solver [4]. In any case, [3]'s solver should be inferior compared to the latter, so our relative results could only improve. We also compared our wall-clock run-times per puzzle size against those provided by [9]. (for CM computation and reconstruction). See at Table 5 a reconstruction comparison and inference time between the results in [9] and ours. In summary, our Edge2Vec is accurate and also way more efficient than the method proposed by [9].

## 5   Conclusions

In this paper we proposed a new CM scheme called Edge2Vec. Our newly derived CM uses embeddings in an efficient yet highly-accurate manner. In particular, we demonstrated how the use of a *single embedding model* combined with our newly proposed *hard batch triplet loss* function results in superior performance relatively to that of previous traditional CMs, as well as previous DNN-based E2E CMs, in terms of both Top-1 accuracy and reconstruction neighbor accuracy. To the best of our knowledge, Edge2Vec currently exhibits the best trade-off between precision, efficiency, and footprint.

Table 5: Comparison to state-of-the-art Model. As can be seen, our Edge2Vec CM combined with Gallagher's solver [3] was able to achieve superior performance in all metrics compared to the previously state-of-the-art method [9], which is a GAN-based E2E CM combined with Paikin and Tal's solver [4]. It is worth to mentioned that our proposed Edge2Vec was able to produced these results with less powerful reconstruction algorithm, while better can be done by only using other reconstruction algorithm.

| Dataset | Neighbor accuracy (%) | | Perfect reconstruction | | Running time (sec) | |
|---|---|---|---|---|---|---|
| | [9] | Edge2Vec | [9] | Edge2Vec | [9] | Edge2Vec |
| MIT (70 pieces) | 84.6 | 92.4 (+7.8) | 4 | 10 (+6) | 62 (x149) | 0.415 |
| McGill (88 pieces) | 76.9 | 83.7 (+6.8) | 7 | 7 (+0) | 98 (x196) | 0.498 |
| BGU (150 pieces) | 76.3 | 91.5 (+15.2) | 2 | 5 (+3) | 286 (x293) | 0.976 |



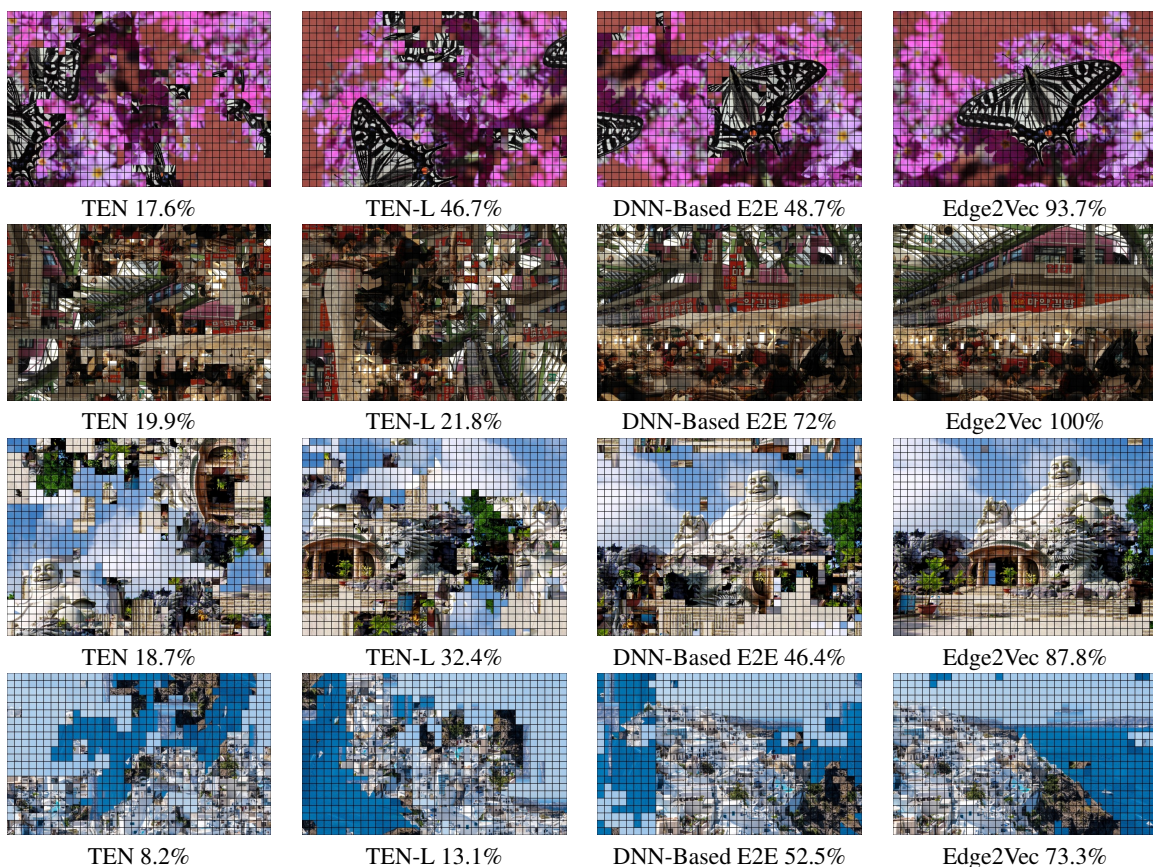| TEN 17.6% | TEN-L 46.7% | DNN-Based E2E 48.7% | Edge2Vec 93.7% |
| TEN 19.9% | TEN-L 21.8% | DNN-Based E2E 72% | Edge2Vec 100% |
| TEN 18.7% | TEN-L 32.4% | DNN-Based E2E 46.4% | Edge2Vec 87.8% |
| TEN 8.2% | TEN-L 13.1% | DNN-Based E2E 52.5% | Edge2Vec 73.3% |

Figure 7: Illustration of reconstruction samples due to Gallagher's solver for Type-2 variant; Edge2Vec exhibits superior performance to previous DNN-based embedding and E2E methods.

# References

[1] T. S. Cho, S. Avidan, and W. T. Freeman. A probabilistic image jigsaw puzzle solver. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 183–190, 2010.

[2] D. Pomeranz, M. Shemesh, and O. Ben-Shahar. A fully automated greedy square jigsaw puzzle solver. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 9–16, 2011.

[3] A. C. Gallagher. Jigsaw puzzles with pieces of unknown orientation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 382–389, 2012.

[4] G. Paikin and A. Tal. Solving multiple square jigsaw puzzles with missing pieces. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4832–4839, 2015.

[5] T. M. Paixao, R. F. Berriel, M. C. Boeres, C. Badue, A. F. De Souza, and T. Oliveira-Santos. A deep learning-based compatibility score for reconstruction of strip-shredded text documents. In *Proceedings of the IEEE 31st SIBGRAPI Conference on Graphics, Patterns and Images*, pages 87–94, 2018.

[6] N. I. Forrest, H. Song, W.M. Matthew, A. Khalid, J.D. William, and K. Kurt. Squeezenet: AlexNet-level accuracy with 50x fewer parameters and <0.5MB model size. *arXiv:1602.07360*, 2015.

[7] M. Sandler, A. Howard, M. Zhu, A. Zhmoginov, and L. Chen. MobileNetV2: Inverted residuals and linear bottlenecks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4510–4520, 2018.

[8] D. Rika, D. Sholomon, O. E. David, and N. S. Netanyahu. A novel hybrid scheme using genetic algorithms and deep learning for the reconstruction of Portuguese tile panels. In *Proceedings of the ACM Conference on Genetic and Evolutionary Computation*, pages 1319–1327, 2019.

[9] D. Bridger, D. Danon, and A. Tal. Solving jigsaw puzzles with eroded boundaries. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3526–3535, 2020.

[10] D. Rika, D. Sholomon, O. E. David, and N. S. Netanyahu. TEN: Twin Embedding Networks for the Jigsaw Puzzle Problem with Eroded Boundaries. *https://doi.org/10.48550/arxiv.2203.06488*, 2022.

[11] V. Nair and G.E. Hinton. Rectified linear units improve restricted boltzmann machines. In *Proceedings of the 27th International Conference on Machine Learning*, pages 807–814, 2010.

[12] E. Agustsson and R. Timofte. NTIRE 2017 challenge on single image super-resolution: Dataset and study. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, July 2017.

[13] Y. Blau, R. Mechrez, R. Timofte, T. Michaeli, and L. Zelnik-Manor. 2018 PIRM challenge on perceptual image super-resolution. *https://doi.org/10.48550/arxiv.1809.07517*, 2018.

[14] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. *arXiv:1412.6980*, 2017.