

# Quasi-Conformal Transformer Network

Han Zhang\*, Qiguang Chen\*, Yuchen Guo\*, and Lok Ming Lui\*

**Abstract.** Information is not generally distributed uniformly in an image domain. Thus, to make the convolutional neural network focus more on those important, some deformation on convolution windows or feature maps should be applied. Besides, the topology of an image should be preserved from the ideas for defining the convolution operation. However, controlling topology is hard and not convenient for existing methods since they all use vector representation for displacements. In this paper, we proposed Quasi-Conformal Transformer Network using Beltrami representation, which is a strong representation to control the bijectivity and the degree of geometric deformations. Together with our Beltrami Solver Net(BSN), we proposed an end-to-end learnable network, which advantages other works on its power to control the geometric deformation of the feature maps.

**Key words.** deformable convolution, deformable pooling, disturbance-invariant, bijectivity

**1. Introduction.** The information is not usually distributed uniformly in images. The target objects in images may differ in pose, position, and scale. Some human-caused errors may also bring challenges to image processing. To solve this problem, some feature representation methods that are invariant to transformation like SIFT (scale-invariant feature transform) [24] are proposed. With the development of computing power, the convolutional neural network (CNN) [17], with its transformation-invariant convolution windows, becomes the most popular and practical model that altered the landscapes of the computer vision community. However, since the convolution window and the pooling window are both defined to be in the fixed size and shape, some works accounting for adaptive convolution are proposed.

The adaptive convolution is mainly defined through two approaches. The spatial transformer (STN) [11] proposed the learnable spatial transformation on feature maps to focus only on regions that contain the most information. However, STN is not convenient to assign non-rigid transformations in the network. Though it may become possible with the thin-plate spline(TSP) method, the topology of the transformation is not guaranteed to be preserved. Another way is direct to define the deformable convolution(DeformConv) [6], where the displacement vectors are assigned directly on convolution windows. Since this work use vector representation, they are easy to result in a messy deformation like that in Figure.1 which contains self-intersections and failed to keep the topology of an image.

However, consideration for topology is a very important advantage of a convolutional neural network over a fully-connected neural network, which ignores it[17]. Reviewing the traditional sliding window methods, like the idea of Sobel operator [13], Prewitt operator [26], *et al.* , is to compute the gradient of the intensity function of images. Since the gradient computation will only involve neighbors, the operator window should always encircle a continuous region. More than that, if the displacement vectors are not restricted properly, overfitting in the process will easily occur. From these perspectives, the topology of images and feature maps should be preserved. In some cases, even the degree of the freedom of deformation for

---

\*Department of Mathematics, The Chinese University of Hong Kong, Shatin, Hong Kong (hanzhang001@cuhk.edu.hk,ycguo@math.cuhk.edu.hk, im.ki@link.cuhk.edu.hk, lmlui@math.cuhk.edu.hk).

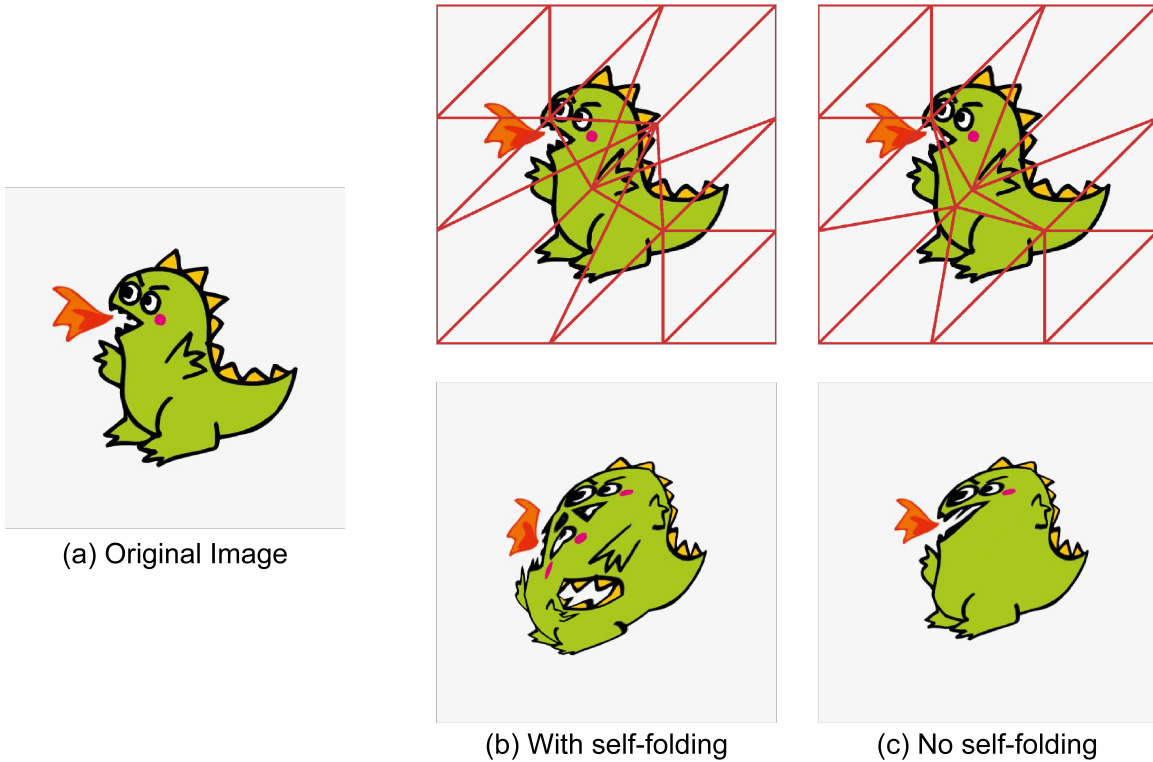


Figure 1: Example for images who failed to preserve the topology during processing. The dinosaur are warped to gain an extra mouth by fold its tail inside.

40 either feature map or convolution window should be controlled to reduce overfitting.

41 However, bijectivity is very important for tasks that want to preserve the topology and  
 42 avoid self-intersection. If the transformation mapping failed to preserve the original topology,  
 43 the transformed image would easily lose its original semantic meaning. The idea can be clearly  
 44 explained by Figure.1, the dinosaur got another mouth in its body which is wrongly mapped  
 45 by its tail. Thus, to assign reasonable spatial attention to the input images, we propose our  
 46 quasi-conformal transformer network in this paper. Compared to the previous work[11], our  
 47 model can produce pixel-wise transformation, which outperforms the TPS variant of spatial  
 48 transformer network that failed to produce topology-preserving mapping.

49 Our model is based on the quasi-conformal theory and uses Beltrami coefficient as the  
 50 mapping representation instead of control points or vector field. By restricting the Beltrami  
 51 coefficient to be less than 1, the associated mapping can be guaranteed to be bijective. Thus,  
 52 our quasi-conformal transformer network is composited of two modules (Figure.7). The coeffi-  
 53 cient estimator predicts the Beltrami coefficients that represent the mapping and the Beltrami  
 54 solver network that transfers the Beltrami representation into vector representation that is  
 55 convenient to do spatial transformation pixel-wisely. Individual modules will be introduced  
 56 and discussed in more detail in Section.4.

57 To evaluate the capacity of the quasi-conformal transformer network, we test the clas-

58 sification results on images with different kinds of deformation. On such deformed images,  
59 conventional standard rectangle convolution can not accomplish the tasks well. However, due  
60 to the spatial attention assigned by the quasi-conformal transformer, our model can obtain  
61 much better results than a network without spatial transformers. Compared to the thin plate  
62 spline variant of spatial transformer network, our method can also acquire a higher accuracy  
63 rate and are easier to optimize because of the proper topology constraint.

64 To sum up, our contributions are:

- 65 • We introduced Beltrami representation for a pixel-wise spatial transformer in neural  
66 network, which is convenient to control the topologic property of the transformation  
67 mapping.
- 68 • By controlling Beltrami coefficients, the mappings produced are guaranteed to be  
69 topology-preserving. By experiments results, such mapping help with classification  
70 results.
- 71 • We can achieve deformable convolution and deformable pooling with quasi-conformal  
72 transformer network, which is difficult by using the pixel-wise thin plate spline variant  
73 of spatial transformer network as it usually failed to preserve the topology of the  
74 transformation mapping.

## 75 2. Related Work.

76 **2.1. Computational Quasi-Conformal.** Computational conformal is a powerful tool to  
77 control the geometric variation and topology in image science[16, 22] and surface processing[20,  
78 8]. Benefitting from the Beltrami representation, the mapping between two different domains  
79 can preserve good geometric properties like bijectivity and smoothness. Through controlling  
80 the Beltrami coefficients with such representation of mappings. Driven by the motivation to  
81 preserve different geometric information, ways of parameterization methods are proposed [9, 3,  
82 2]. Such convenient representation are also popular and succeed in computational fabrication  
83 community [29, 5, 25]. With the capability to handle large deformations, the quasi-conformal  
84 method also succeed in registration for images [16, 22] and surfaces [4] and segmentation with  
85 topology- and convexity prior [33, 28].

86 **2.2. Deep Learning.** Neural network models [18] are widely employed and greatly suc-  
87 cessed in different fields like image science and natural language processing. With convolution  
88 neural network[17], translation invariants become available for learning methods and intro-  
89 duce the development of image science to the next stage. Classification is the first vision task  
90 that benefits from deep layers[15]. Inspired by the success, other tasks adopt deep learning  
91 models to solve the problems[]. Besides, the great success of the application of deep learning,  
92 the model itself is also developing. Max-pooling layers are introduced into the neural network  
93 model for dimensionality reduction by [30]. [27] increase the depth of the network by using  
94 very small ( $3 \times 3$ ) convolution windows to enhance the non-linearity of the network. With  
95 shortcut connections, ResNet [10] overcome the vanishing problem that prevents models from  
96 learning.

97 Jaderberg *et al.* [11] proposed the spatial transformer network, which is the first work  
98 to introduce spatial transformation to assign attention to feature map. Warping a feature  
99 map and doing convolution in the normal way is mathematically equivalent to defining some

100 special convolutions learned from data. Jeon and Kim [12] proposed active convolution that  
 101 can sample the locations of the convolution with some displacements that are shared over  
 102 different spatial locations. However, the information is not uniformly distributed in the feature  
 103 map. Not all pixels contribute equally to the final results[23]. Motivated by this, Qi *et al.*  
 104 [6, 34] designed deformable convolution whose convolution differs not only among different  
 105 maps but even on different spatial locations of the feature map.

106 **3. Mathematical Background.** In this chapter, we introduce some fundamental geometry  
 107 concepts and theories that are related to our model. In brief words, a quasi-conformal mapping  
 108 is mainly used for our registration map. Besides, since we need to compress our registration  
 109 map as discussed in the introduction, Fourier compression for Beltrami representation will  
 110 also be introduced.

### 111 3.1. Conformal Maps.

112 **Definition 3.1.** (*Quasi-conformal mape*). A conformal map is a map  $f : \mathbb{C} \rightarrow \mathbb{C}$  that  
 113 satisfying the Beltrami equation

$$114 \quad (3.1) \quad \frac{\partial f}{\partial \bar{z}} = \mu(z) \frac{\partial f}{\partial z}$$

115 for some complex-valued function named as Beltrami coefficient  $\mu$  satisfying  $\|\mu\|_\infty < 1$  and  
 116  $\frac{\partial f}{\partial \bar{z}}$  is non-vanishing almost everywhere. The complex partial derivatives are given by

$$117 \quad (3.2) \quad \frac{\partial f}{\partial z} := \frac{1}{2} \left( \frac{\partial f}{\partial x} - i \frac{\partial f}{\partial y} \right) \quad \text{and} \quad \frac{\partial f}{\partial \bar{z}} := \frac{1}{2} \left( \frac{\partial f}{\partial x} + i \frac{\partial f}{\partial y} \right)$$

118  $\mu$  is the Beltrami representation, which is also called the Beltrami coefficient, of the quasi-  
 119 conformal map  $f$ . It's worthy to mention that  $\mu$  is a measure of non-conformality. Particularly,  
 120 for a point  $p$ , suppose  $\mu(p) = 0$ . Then, the associated quasi-conformal map  $f$  is conformal  
 121 around a small neighborhood of  $p$ . In this case, Equation 3.1 is the Cauchy-Riemann equation.  
 122 This can also illustrate that conformality analysis of a quasi-conformal map  $f$  can be simplified  
 123 into the analysis of its associated Beltrami coefficient  $\mu$ . Infinitesimally, such a map  $f$  can be  
 124 rewritten as follows in a local neighborhood around a point  $p$ :

$$125 \quad (3.3) \quad \begin{aligned} f(z) &= f(p) + f_z(p)z + f_{\bar{z}}(p)\bar{z} \\ &= f(p) + f_z(p)(z + \mu(p)\bar{z}) \end{aligned}$$

126 This further enhanced our discussion before that  $f$  is conformal when  $\mu(p) = 0$ . To explain  
 127 for the equation above,  $f(p)$  is a translation, while  $f_z(p)$  is a dilation. Since both of them  
 128 are conformal, all the non-conformality of  $f$  is brought by  $D(z) = z + \mu(p)\bar{z}$ . Hence, the  
 129 Beltrami coefficient  $\mu$  actually encode the conformality of  $f$ . Analyzing quasi-conformal  $f$  is  
 130 equivalent to that for its associated Beltrami coefficient  $\mu$ . To be detail, the angle of maximal  
 131 magnification is  $\arg(\mu(p))/2$  with magnifying factor  $1 + |\mu(p)|$ ; for the maximal shrinking is  
 132 the orthogonal angle  $(\arg(\mu(p)) - \pi)/2$  with shirking factor  $1 - |\mu(p)|$ .

133 The maximal quasi-conformal dilation of  $f$  is given by

$$134 \quad (3.4) \quad K = \frac{1 + \|\mu\|_\infty}{1 - \|\mu\|_\infty}$$

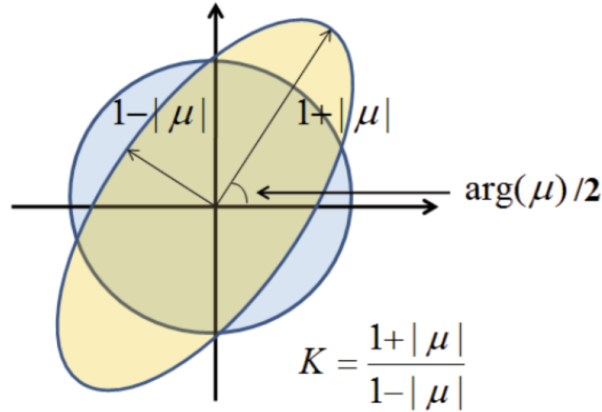


Figure 2: Illustration of how the Beltrami coefficient measures the conformality distortion of a quasi-conformal map

135 Figure 2 illustrated the geometry of quasi-conformal map.

136 Another important relationship between a map and its Beltrami coefficients is the diffeomorphism property. By a norm constraint on  $\mu$ , the bijectivity of  $f$  can be preserved which is explained as following theory.

139 **Theorem 3.2.** *If  $f : \mathbb{C} \rightarrow \mathbb{C}$  is a  $C^1$  map. Define*

140 (3.5) 
$$\mu = \frac{\partial f}{\partial \bar{z}} / \frac{\partial f}{\partial z}$$

141 *If  $\mu$  satisfies  $\|\mu_f\|_\infty < 1$ , then  $f$  is bijective.*

142 For two quasi-conformal maps  $f, g$ , the Beltrami coefficient of their composition can be expressed in terms of their individual Beltrami coefficients  $\mu_f, \mu_g$  directly according to the following theory.

145 **Theorem 3.3.** *For two quasi-conformal maps  $f : \Omega \subset \mathbb{C} \rightarrow f(\Omega)$  and  $g : f(\Omega) \rightarrow \mathbb{C}$ , whose Beltrami coefficients are  $\mu_f, \mu_g$  respectively. The Beltrami coefficient of the composited function  $g \circ f$  is clearly defined as*

148 (3.6) 
$$\mu_{g \circ f} = \frac{\mu_f + (\overline{f_z}/f_z)(\mu_g \circ f)}{1 + (\overline{f_z}/f_z)\overline{\mu_f}(\mu_g \circ f)}$$

149 **Theorem 3.4.** *Suppose  $f : M_1 \rightarrow M_2$  and  $g : M_2 \rightarrow M_3$  are two quasi-conformal maps. Write the associated Beltrami coefficient as  $\mu_{f^{-1}}$  and  $\mu_g$  respectively. Suppose  $\mu_{f^{-1}} = \mu_g$ . Then, the Beltrami coefficient of  $g \circ f$  is equal to 0. In other words,  $g \circ f : M_1 \rightarrow M_3$  is a conformal map.*

153 As we almost have everything for computing Beltrami coefficients from a given quasi-conformal map, we also need the converse. That's to say, given complex function  $\mu$ , we can solve out it associated quasi-conformal map  $f$  if  $\|\mu\|_\infty < 1$ .

156 Given a Beltrami coefficient  $\mu$  which is less than 1. Denote the corresponding quasi-  
157 conformal map  $f : \mathbb{C} \rightarrow \mathbb{C}$  as  $f = u + iv$ , we have

$$158 \quad (3.7) \quad \mu(f) = \frac{(u_x - v_y) + \sqrt{-1}(v_x + u_y)}{(u_x + v_y) + \sqrt{-1}(v_x - u_y)}$$

159 Rewrite  $\mu = \rho + i\tau$ . From the Beltrami Equation 3.1, we are able to use one pair of  $v_x, v_y$  or  
160  $u_x, u_y$  to describe the partial derivatives of the other pairs with linear combinations as:

$$161 \quad (3.8) \quad \begin{aligned} v_y &= \alpha_1 u_x + \alpha_2 u_y; & \text{and} & & -u_y &= \alpha_1 v_x + \alpha_2 v_y \\ -v_x &= \alpha_2 u_x + \alpha_3 u_y & & & u_x &= \alpha_2 v_x + \alpha_3 v_y \end{aligned}$$

162 where  $\alpha_1 = \frac{(\rho-1)^2 + \tau^2}{1 - \rho^2 - \tau^2}$ ;  $\alpha_2 = -\frac{2\tau}{1 - \rho^2 - \tau^2}$ ;  $\alpha_3 = \frac{(1+\rho)^2 + \tau^2}{1 - \rho^2 - \tau^2}$ . Since  $\nabla \cdot \begin{pmatrix} -v_y \\ v_x \end{pmatrix} = 0$  and  $\nabla \cdot$

163  $\begin{pmatrix} -u_y \\ u_x \end{pmatrix} = 0$ , we have

$$164 \quad (3.9) \quad \nabla \cdot \left( A \begin{pmatrix} u_x \\ u_y \end{pmatrix} \right) = 0 \quad \text{and} \quad \nabla \cdot \left( A \begin{pmatrix} v_x \\ v_y \end{pmatrix} \right) = 0$$

165 where  $A = \begin{pmatrix} \alpha_1 & \alpha_2 \\ \alpha_2 & \alpha_3 \end{pmatrix}$  and is easily checked to be symmetric positive definite. Equation 3.9  
166 is called the generalized *Laplace equation*. Solving the equation, so one can obtain everything  
167 for  $f$ . Note that here the information for  $\alpha_1, \alpha_2$  and  $\alpha_3$  are given by the Beltrami coefficient  
168  $\mu$ . In a word, when we need to find an optimal quasi-conformal map, we can convert the  
169 problem to find an optimal complex-valued function alternatively.

170 Alternatively, quasi-conformal maps can also be defined between two Riemann surfaces  $\mathcal{M}$   
171 and  $\mathcal{N}$  by *Beltrami differential* instead of Beltrami coefficient. A *Beltrami differential*  $\mu(z) \frac{dz}{z}$   
172 on  $\mathcal{M}$  is an assignment to each chart  $(U_\alpha, \phi_\alpha)$  of an  $L^\infty$  complex-valued function  $\mu_\alpha$  defined  
173 on the local parameter  $z_\alpha$  such that

$$174 \quad (3.10) \quad \mu_\alpha(z_\alpha) \frac{dz_\alpha}{dz_\alpha} = \mu_\beta(z_\beta) \frac{dz_\beta}{dz_\beta}$$

175 on the domain also covered by another chart  $(U_\beta, \psi_\beta)$ , where  $\frac{dz_\beta}{dz_\alpha} = \frac{d}{dz_\alpha} \phi_{\alpha\beta}$  and  $\phi_{\alpha\beta} = \phi_\beta \circ \phi_\alpha^{-1}$ .  
176 For a better illustration, we give figure 3.

177 **3.2. Fourier Approximation for Beltrami Representation.** [21] shows that Fourier ap-  
178 proximation for Beltrami representation can easily preserve the mapping information while  
179 that for the representation of coordinate functions fails to enforce the homeomorphism. Com-  
180 pared with the latter method which requires that the Jacobian of the coordinate functions  
181 has to be greater than 0, the Fourier approximation for Beltrami representation requires only  
182 that the supreme norm must be less than 1, which is easier to be satisfied.

183 In the discrete case, an  $N \times N$  Beltrami coefficient  $\mu$  can be separated into two images  
184 representing the real and imaginary parts,  $\mu_r$  and  $\mu_i$ , respectively. The DFT of  $\mu_r$  is

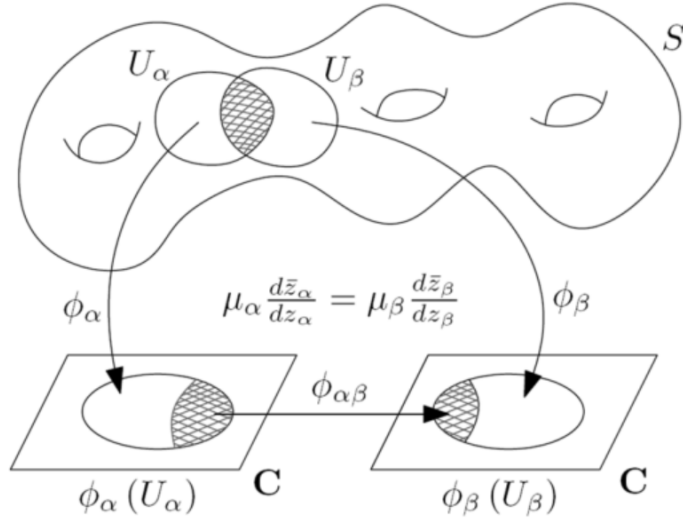


Figure 3: Beltrami differential on general Riemann surfaces

$$185 \quad (3.11) \quad \hat{\mu}_r(m, n) = \frac{1}{N^2} \sum_{k=0}^{N-1} \sum_{l=0}^{N-1} \mu_r(k, l) e^{-\sqrt{-1} \frac{2\pi km}{N}} e^{-\sqrt{-1} \frac{2\pi ln}{N}}$$

186 This is equivalent to

$$187 \quad (3.12) \quad \hat{\mu}_r = U \mu_r U$$

188 where  $U_{kl} = \frac{1}{N} e^{-\sqrt{-1} \frac{2\pi kl}{N}}$ ,  $0 \leq k, l \leq N-1$ . The inverse DFT of  $\hat{\mu}_r$  is

$$189 \quad (3.13) \quad \mu_r(p, q) = \sum_{m=0}^{N-1} \sum_{n=0}^{N-1} \hat{\mu}_r(m, n) e^{\sqrt{-1} \frac{2\pi pm}{N}} e^{\sqrt{-1} \frac{2\pi qn}{N}}$$

190 which can be rewritten as

$$191 \quad (3.14) \quad \mu_r = (NU^*) \hat{\mu}_r (NU^*)$$

192 When using Fourier coefficients  $\hat{\mu}_r, \hat{\mu}_i$  to approximate Beltrami coefficients, we keep only a  
 193 small fraction of the frequency components which acts as the low-frequency components. This

194 is equivalent to saying that only a small fraction of the frequency components can well capture  
 195 the majority of deformation. Motivated by this idea, we propose the Beltrami Solver Network  
 196 (BSNet) which takes a Beltrami coefficient as input and uses the Fourier approximation to  
 197 represent the global information of the input.

198 **3.3. Numerical Implementation of LBS.** Given the Beltrami Coefficient  $\mu$ , we can re-  
 199 construct the corresponding quasi-conformal mapping  $f$  by solving (3.9)

200 In discrete case, the parameter domain  $D$  is a mesh grid. The restriction of  $f$  on each  
 201 triangular face  $T$  is linear and can be written as

$$202 \quad (3.15) \quad f|_T(x, y) = \begin{bmatrix} u|_T(x, y) \\ v|_T(x, y) \end{bmatrix} = \begin{bmatrix} a_T x + b_T y + r_T \\ c_T x + d_T y + s_T \end{bmatrix}$$

203 Obviously, on each face, we have induced partial derivatives from the face-wise linear  
 204 assumption.

205 Hence, the partial derivatives of  $f$  at each face  $T$  can be denoted as  $D_x f(T) = a_T + ic_T$   
 206 and  $D_y f(T) = b_T + id_T$ . Now the gradient  $\nabla_T f := (D_x f(T), D_y f(T))^t$  on  $T$  can be computed  
 207 by solving

$$208 \quad (3.16) \quad \begin{pmatrix} v_1 - v_0 \\ v_2 - v_0 \end{pmatrix} \nabla_T f_i = \begin{pmatrix} f_i(\vec{v}_1) - f_i(\vec{v}_0) \\ f_i(\vec{v}_2) - f_i(\vec{v}_0) \end{pmatrix}$$

209 where  $[\vec{v}_0, \vec{v}_1]$  and  $[\vec{v}_0, \vec{v}_2]$  are two edges on  $T$ .

210 Besides,  $\mu$  is a face-based function. Denote the face-based function  $\alpha_i$  on face  $T$  by  $\alpha_i^T$ ,  
 211 where  $i = 1, 2, 3$ . From 3.8, we have

$$212 \quad (3.17) \quad \begin{aligned} -d_T &= \alpha_1^T a_T + \alpha_2^T b_T \\ c_T &= \alpha_2^T a_T + \alpha_3^T b_T \end{aligned}$$

213 and

$$214 \quad (3.18) \quad \begin{aligned} -b_T &= \alpha_1^T c_T + \alpha_2^T d_T \\ a_T &= \alpha_2^T c_T + \alpha_3^T d_T \end{aligned}$$

215 Let  $T = [\vec{v}_i, \vec{v}_j, \vec{v}_k]$  and  $w_I = f(\vec{v}_I)$ , where  $I = i, j, k$ . Suppose  $v_I = g_I + ih_I$  and  
 216  $w_I = s_I + it_I$  ( $I = i, j, k$ ). Using (3.16), According to mapping (3.16), for each face  $T$ , we  
 217 have

$$218 \quad (3.19) \quad \begin{bmatrix} a_T & b_T \\ c_T & d_T \end{bmatrix} \begin{bmatrix} g_j - g_i & g_k - g_i \\ h_j - h_i & h_k - h_i \end{bmatrix} = \begin{bmatrix} s_j - s_i & s_k - s_i \\ t_j - t_i & t_k - t_i \end{bmatrix}$$

219 Thus,

$$220 \quad (3.20) \quad \begin{bmatrix} a_T & b_T \\ c_T & d_T \end{bmatrix} = \frac{1}{2 \cdot \text{Area}(T)} \begin{bmatrix} s_j - s_i & s_k - s_i \\ t_j - t_i & t_k - t_i \end{bmatrix} \begin{bmatrix} h_k - h_i & g_i - g_k \\ h_i - h_j & g_j - g_i \end{bmatrix}$$

$$221 \quad (3.21) \quad = \begin{bmatrix} A_i^T s_i + A_j^T s_j + A_k^T s_k & B_i^T s_i + B_j^T s_j + B_k^T s_k \\ A_i^T t_i + A_j^T t_j + A_k^T t_k & B_i^T t_i + B_j^T t_j + B_k^T t_k \end{bmatrix}$$



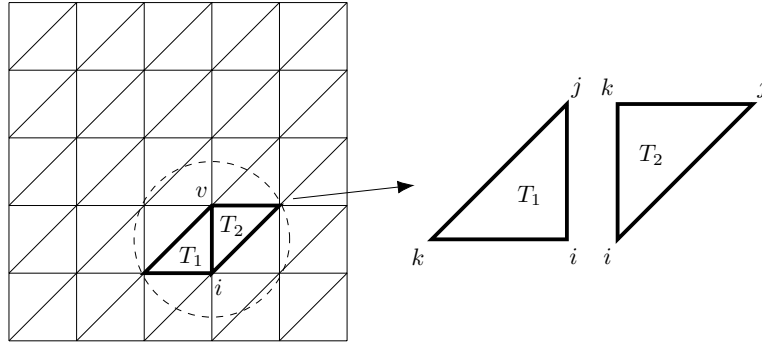


Figure 4: Illustration of the derivation of the coefficient of vertex  $v$

222 where

$$\begin{aligned}
 223 \quad (3.22) \quad & A_i^T = (h_j - h_k) / 2 \cdot \text{Area}(T); & B_i^T &= (g_k - g_j) / 2 \cdot \text{Area}(T) \\
 & A_j^T = (h_k - h_i) / 2 \cdot \text{Area}(T); & B_j^T &= (g_i - g_k) / 2 \cdot \text{Area}(T) \\
 & A_k^T = (h_i - h_j) / 2 \cdot \text{Area}(T); & B_k^T &= (g_j - g_i) / 2 \cdot \text{Area}(T)
 \end{aligned}$$

224 For each vertex  $v_i$ , let  $N_i$  be the collection of neighborhood faces attached to  $v_i$ . By careful  
 225 checking, one can observe that

$$226 \quad (3.23) \quad \sum_{T \in N_i} A_i^T b_T = \sum_{T \in N_i} B_i^T a_T; \quad \sum_{T \in N_i} A_i^T d_T = \sum_{T \in N_i} B_i^T c_T;$$

227 Substituting Equations (3.17) and (3.18) into (3.23), we obtain the following equations

$$228 \quad (3.24) \quad \sum_{T \in N_i} (A_i^T [\alpha_1^T a_T + \alpha_2^T b_T] + B_i^T [\alpha_2^T a_T + \alpha_3^T b_T]) = 0$$

$$229 \quad (3.25) \quad \sum_{T \in N_i} (A_i^T [\alpha_1^T c_T + \alpha_2^T d_T] + B_i^T [\alpha_2^T c_T + \alpha_3^T d_T]) = 0$$

230 Replacing  $a_T$ ,  $b_T$ ,  $c_T$  and  $d_T$  with their corresponding expressions, we derive the following  
 231 coefficient for the central vertex  $i$  of  $N_i$

$$232 \quad (3.26) \quad c_i = \sum_{T \in N_i} [\alpha_1^T (A_i^T)^2 + 2\alpha_2^T A_i^T B_i^T + \alpha_3^T (B_i^T)^2]$$

233 For two incident triangular faces  $T_1, T_2 \in N_i$ , the edge  $e$  between  $T_1$  and  $T_2$  connect the  
 234 central vertex  $i$  and another vertex  $v$ , as shown in Figure 4. In this case, the index of vertex  
 235  $v$  in  $T_1$  is  $j$ , and that in  $T_2$  is  $k$ . The coefficient of vertex  $v$  can then be written as follow

$$\begin{aligned}
236 \quad (3.27) \quad c_v = & \alpha_1^{T_1} A_i^{T_1} A_j^{T_1} + \alpha_2^{T_1} (A_i^{T_1} B_j^{T_1} + A_j^{T_1} B_i^{T_1}) + \alpha_3^{T_1} B_i^{T_1} B_j^{T_1} + \\
& \alpha_1^{T_2} A_i^{T_2} A_k^{T_2} + \alpha_2^{T_2} (A_i^{T_2} B_k^{T_2} + A_k^{T_2} B_i^{T_2}) + \alpha_3^{T_2} B_i^{T_2} B_k^{T_2}
\end{aligned}$$

237 According to Equation (3.26), (3.27), (3.24) and (3.25), for a vertex  $i$ , we can write down  
238 the following equations

$$\begin{aligned}
239 \quad (3.28) \quad c_i s_i + \sum_{v \in V_i} c_v s_v &= 0 \\
c_i t_i + \sum_{v \in V_i} c_v t_v &= 0
\end{aligned}$$

240 where  $V_i$  is the set of adjacent vertices of vertex  $i$ .

241 For an  $N \times N$  mesh grid, we have

$$\begin{aligned}
242 \quad (3.29) \quad C_s \mathbf{s} &= \mathbf{0} \\
C_t \mathbf{t} &= \mathbf{0}
\end{aligned}$$

243 where  $\mathbf{s}$  and  $\mathbf{t}$  are the  $N \times N$  dimensional coordinate vectors, in which the entries of  
244 boundary constraints set to their true values.  $C_s$  and  $C_t$  are the same  $N^2 \times N^2$  sparse matrix,  
245 each row of which contains  $c_i$  and  $c_v$  for a vertex in the mesh grid. The only difference between  
246  $C_s$  and  $C_t$  is that the rows that correspond to the boundary constraints of the two coordinates  
247 are set to 0.

248 Solving this linear system with the boundary constraints, we can obtain the corresponding  
249 mapping  $f$  given a Beltrami coefficient  $\mu_f$ .

250 **3.4. Convolutional Neural Network.** Deep learning is a branch of Machine Learning.  
251 With the development of computation power, deep learning become popular and helped mul-  
252 tiple fields like computer vision, natural language processing, financial modeling, etc. Imi-  
253 tating the way human beings learn from experiences, supervised deep learning methods learn  
254 to solve a task by doing regression from data. A deep learning model, which is generally a  
255 neural network, consists of three main parts: architecture, optimization method, and the loss  
256 function.

257 For the architecture, the neural network is made of neurons that mimic the functionality of  
258 human brain. For each neuron, the main components include the input feature  $x_1, x_2, \dots, x_n$ ,  
259 thier correspondending weighting  $\omega_1, \omega_2, \dots, \omega_n$ , the transfer function which is simply sum-  
260 mation in most cases, the bias  $b$  and the activation function  $\phi$ . When a group of feature passes  
261 through this neuron, it will output a signal  $y$ , which is computed by

$$262 \quad (3.30) \quad y = \phi \left( \sum_{i=1}^n x_i \omega_i + b \right)$$

263 The output  $y$  of this neuron will be the input of the neurons connected to this one in the  
264 after layer. Through such a process, the signals are transmitted layer by layer, and features  
265 are processed to produce the final output. In recent years, the neural network become deeper  
266 and deeper. This may be explained as more layers would enrich the levels of features and

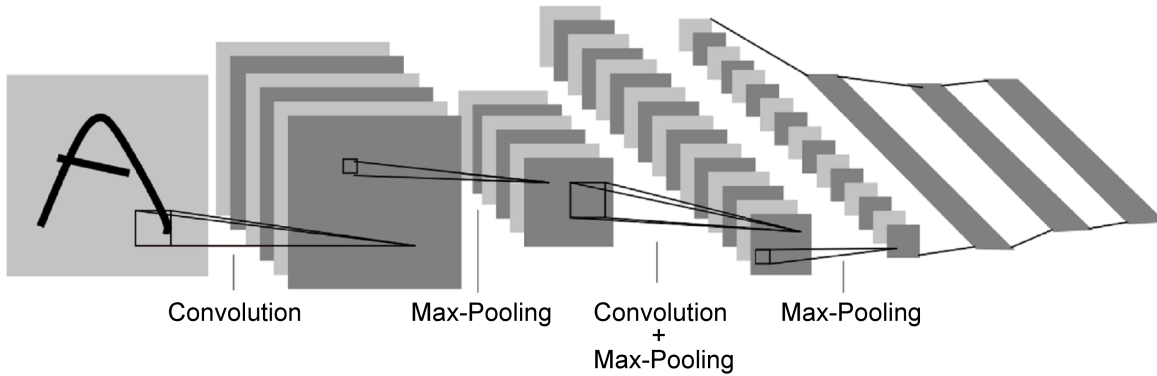


Figure 5: Convolution Neural Network

267 leads to a better result. However, a network model with too many layers is hard to train  
 268 and get to convergence. Through deep layers can bring more non-linearity and flexibility to  
 269 approximate the function that describes the problem to be solved, too many parameters may  
 270 result in vanishing and exploding gradient problems [10]. For tasks like image analysis, pixel  
 271 signals are defined as associated with their neighbors. In other words, to extract the features  
 272 like edges, shapes, textures, and objects, one should not only consider this pixel only but  
 273 also should account for its relation with others nearby. Besides such a need for consideration  
 274 on local, the location of some objects may be different spatially in the image domain. The  
 275 convolutional neural network, which is made of layers of trainable filters, is to solve such tasks.  
 276 In a convolution neural network, the neuron's input feature is encircled by a moving window  
 277 while the output will be located at the place accordingly (see Figure.5).

278 During the training of a deep neural network model, the data are input to generate predic-  
 279 tions. To evaluate how well the model is trained, we need to compare the difference between  
 280 the prediction and the ground truth. Such a difference is not always the Euclidean distance  
 281 between them. For example, for a classification task, cross-entropy loss is used more often  
 282 due to its capacity on measuring the disorder and unpredictability of a system. When we  
 283 design the model to output a probability distribution  $q(x)$  where the input should belong,  
 284 using cross-entropy loss can enhance the data cluster into groups and approach the actual  
 285 distribution  $p(x)$ .

$$286 \quad (3.31) \quad E(p, q) = - \sum p(x) \log(q(x))$$

287 After the loss function and architecture are defined, the distance in the output space  
 288 and the function to be regressed are fixed. The problem that remained is to optimize the  
 289 whole model to obtain a small difference between the predicted and the target. To reduce the  
 290 difference between the output of the network and the ground truth, the weights  $\omega_i$  for each  
 291 neuron  $k$  should be updated by gradient descent. For the error defined according to some loss  
 292 function  $E(Y, \hat{Y})$  where  $Y$  is the label and the prediction  $\hat{Y} = N(X; \Omega)$ , where  $N$  is the neural  
 293 network with weight parameters  $\Omega$ .

294 The gradient for each  $\omega$  is computed based on the chain-rule. This greatly reduce the

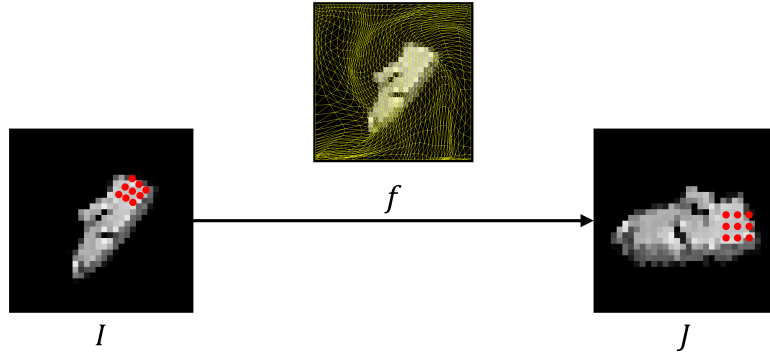


Figure 6: Illustration for deformable convolution via deformable feature map: do regular convolution on the deformed feature map  $J = I \circ f$  is equivalent to do deformable convolution on the original feature map  $I$

295 descent computation between layers. To make it clear, suppose the we have  $N + 1$  layers in  
 296 total, number the first(input layer) as 0 and the last(output layer) as  $N$ . Denote  $k$ -th layer  
 297 has  $M_k$  nodes. The weight parameter from node  $i$  in  $m - 1$  layer to node  $j$  in  $m$  layer is  $\omega_{ij}^m$ .  
 298 The The descent for weights from the last hidden layer  $N - 1$  to the output layer  $N$  is

$$299 \quad (3.32) \quad \frac{dError}{d\omega_{ij}^m} = \frac{dError}{dx_j^m} \frac{dx_j^m}{d\omega_{ij}^m} \quad \text{for } i = 1, \dots, M_{m-1}; j = 1, \dots, M_m$$

300 where  $x_j^m = h_j^m \left( \sum_{i=1}^{M_{m-1}} x_i^{m-1} \omega_{ij}^m \right)$  is the activated neuron value for node  $j$  in the  $m$ -th layer.

301 Then descent for middle hidden layer  $k$  is

$$302 \quad (3.33) \quad \frac{dError}{d\omega_{ij}^k} = \frac{dError}{dx_j^k} \frac{dx_j^k}{d\omega_{ij}^k} \quad \text{for } i = 1, \dots, M_{k-1}; j = 1, \dots, M_k$$

303 where  $x_j^k = h_j^k \left( \sum_{i=1}^{M_{k-1}} x_i^{k-1} \omega_{ij}^k \right)$  is the activated neuron value for node  $j$  in the  $k$ -th layer.

304 **4. Methodology.** In this section, we will introduce our *Quasi-Conformal transformer*,  
 305 which can deform the feature maps through a mapping learned from training data. To make it  
 306 convenient, let's assume our feature map is a one-channel feature map, which is  $I : \mathbb{R}^{H \times W} \rightarrow$   
 307  $\mathbb{R}$ , where  $H, W$  denotes the height and width respectively. For multi-channel images, we can  
 308 do the same as that for one channel for each channel of it. The pipeline is like this: 1) Input  
 309 a feature map  $I$  into *Beltrami Generator*, which is a simple and small network, to obtain  
 310 the Beltrami coefficient  $\mu$  of a Quasi-conformal mapping. 2) Input the Beltrami coefficient  $\mu$   
 311 into a pre-trained *Beltrami Solver Network* (BSnet) to acquire the associated Quasi-conformal  
 312 mapping  $f_\mu$ . 3) Spatially transform the feature map  $I$  with mapping  $f_\mu$  as  $I \circ f_\mu$ . 3) Input the  
 313 deformed image  $I \circ f_\mu$  to a classifier or segmentation network according to the particular task.  
 314 The architecture for our Quasi-Conformal transformer network is illustrated as Figure.7.

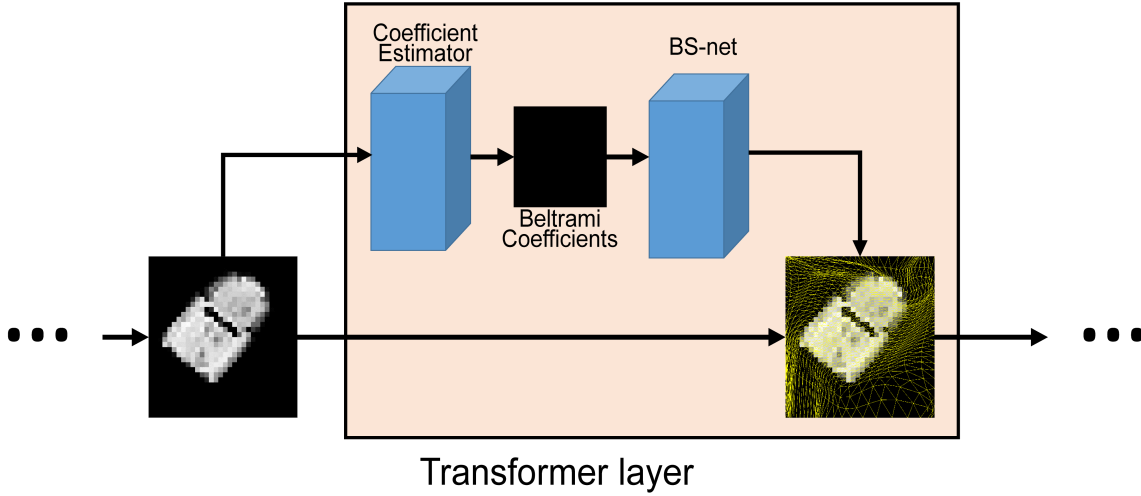


Figure 7: Quasi-Conformal transformer network

315 The whole model is an end-to-end neural network model without manual interference.  
 316 Except for the BS-net, which is to solve for the quasi-conformal mapping given the Beltrami  
 317 coefficients and is fixed as long as it’s pre-trained, every parameter in the network is learnable.  
 318 In the following, we will introduce the grid parameterization of images, the Beltrami coefficient  
 319 generator, Beltrami solver network, and spatial transformation in detail. Besides these, we  
 320 will discuss some training tricks for the overall Quasi-conformal transformer network at the  
 321 end of this section.

322 **4.1. Beltrami Solver Network.** To make the ideal feasible, the first component required is  
 323 a network to convert Beltrami coefficients to their corresponding mappings which is proposed  
 324 by [1]. In the previous work, [21], the conversion from Beltrami coefficients to mappings is  
 325 achieved by using Linear Beltrami Solver (LBS), in which a sparse linear system is required to  
 326 be solved. In this paper, we are going to use a neural network to approximate the mappings  
 327 given their corresponding Beltrami Coefficients. There are two benefits. On the one hand,  
 328 once trained, a well-designed neural network gives predictions much faster than solving the  
 329 sparse linear system. On the other hand, the neural network can backpropagate errors from  
 330 its output to input, which makes it possible to use a trained neural network as a component  
 331 when training another network to solve a complicated task.

332 From above we notice that a single Beltrami coefficient  $\mu_{ij}$  represents the distortion of  
 333 a local region, where  $i$  and  $j$  represent the indices of a triangle in the spatial position. The  
 334 global distortion depends on the entire Beltrami coefficient  $\mu$ , since the mapping is obtained  
 335 by solving the linear system. It is natural that we can use cascaded convolutional and down-  
 336 sampling layers to extract global information from the Beltrami coefficients, which can then  
 337 be used to predict the mappings. The network structure should be similar to U-Net.

338 However, we learn the prior knowledge from [21] that Beltrami representation can be  
 339 easily compressed by Fourier approximation, which means that low-frequency components  
 340 hold most of the global distortion information. With this prior knowledge, we can further

341 simplify the network structure. Compared with the convolutional layer, Fourier Transform  
 342 has no trainable parameter and is much faster. And the generated low-frequency component  
 343 has only two channels, much less than the deep features extracted by neural networks. As a  
 344 result, we use Fourier Transform to extract global information.

345 Once we use the Fourier approximation, we have to think about how to combine spatial  
 346 operations, such as convolution and interpolation, and the Fourier coefficients on the frequency  
 347 domain. Directly performing such operations on the frequency domain easily leads to poor  
 348 performance. It is critical to introduce a layer for transforming the frequency features to the  
 349 spatial domain.

350 Considering that the size of the deep spatial features or Fourier coefficient mentioned  
 351 above is different from the input images and the Beltrami coefficient, we cannot use (3.12)  
 352 and (3.14) directly.

353 In order to tackle this problem, we propose the Domain Transform Layer (DTL) which  
 354 imitates the computation of (3.12) and (3.14). This layer can be formulated

$$355 \quad (4.1) \quad \begin{aligned} \hat{\mu} &= M\mu N = (\mu^T M^T)^T N \\ \mu &= M\hat{\mu} N = (\hat{\mu}^T M^T)^T N \end{aligned}$$

356 where  $M$  and  $N$  are trainable complex matrices. This layer can be implemented by  
 357 stacking two  $1 \times 1$  convolution layers, together with some permutation operations.

358 As shown in subsection 4.1, given an input feature map of shape  $(H, W, C)$ , we can permute  
 359 the feature map to be  $(H, C, W)$ . Then we perform  $K$  kernels  $1 \times 1$  convolution. The shape  
 360 of the resulting feature map should be  $(H, C, K)$ . These operations are equivalent to matrix  
 361 multiplication of a  $H \times W$  matrix and a  $W \times K$  matrix.

362 From (4.1) we notice that (4.1) can be implemented by stacking two matrix multiplication  
 363 blocks. The detail structure is shown in Figure 8. In our experiments,  $H = W = K = L = 14$ .

364 With DTL, the features in the spatial domain can be transformed to a proper domain  
 365 where spatial operations work. We can then perform convolution and upsampling on the  
 366 features to obtain the mappings.

367 However, during the experiments, we found that the output mappings of the network are  
 368 quite similar to their ground truth. But the outputs of the network have fewer details. In  
 369 order to remedy this problem, we introduce a second path to improve the local details of the  
 370 distortion, which are discarded in the computation of approximated Beltrami representation  
 371 in the first path.

372 As in Figure 10, the second path is the upper path. In this skip path, convolution and  
 373 a downsampling are performed on the input  $\mu$ , after which two more convolution layers are  
 374 performed and the output features are concatenated to the output from the first path. Ex-  
 375 perimental results in the following sections show the necessity of this skip path.

376 We trained this network in an unsupervised setting. As mentioned in subsection 4.1, a  
 377 single value in the Beltrami coefficient represents the distortion of a triangular face. The shape  
 378 of the Beltrami coefficient describing an  $N \times N$  image should be  $(N - 1) \times [2 \times (N - 1)]$ . Due  
 379 to this reason, LBS takes an  $(N - 1) \times 2 \times (N - 1)$  dimensional vector as its input. However,  
 380 in this paper, we wish to ensure the consistency between the input and output of our model.  
 381 To cope with this conflict, we adopt the following method.

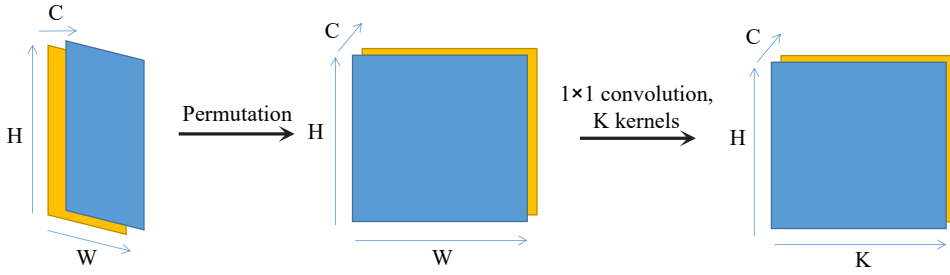


Figure 8: Matrix multiplication in the domain transform layer

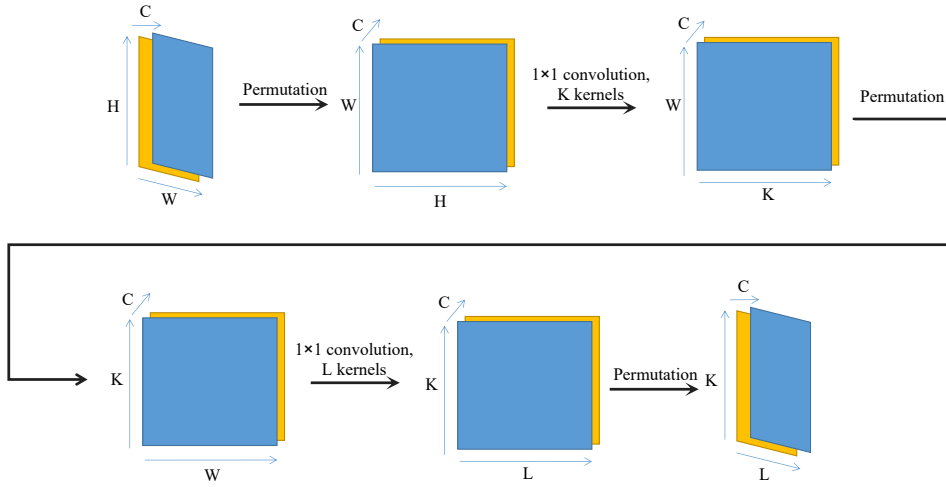


Figure 9: Domain Transform Layer

382 Countless  $N \times N$  Beltrami coefficients  $\mu_{sqr}$  are generated by stacking pairs of images in  
 383 the ILSVRC2012 dataset, which are augmented with some data augmentation tricks, such  
 384 as random crop and flipping.  $\mu_{sqr}$  serves as the input of BSNet, representing the Beltrami  
 385 coefficients of the triangles with odd indices in each row. An  $(N - 1) \times [2 \times (N - 1)]$  Beltrami  
 386 coefficient  $\mu_{rect}$  is then obtained by removing the last row and column of  $\mu_{sqr}$  and interpo-  
 387 lating the Beltrami coefficients representing the triangles with even indices with the values  
 388 representing the surrounding faces. the linear system in the LBS can then be retrieved with  
 389  $\mu_{rect}$ , the coefficients of which are then used to compute the loss.

390 In Quasi-conformal geometry, every vertex with index  $i$  in a mesh satisfies Equations  
 391 (3.29). It is natural to regard Equations (3.29) as a loss function when training a neural  
 392 network to approximate the ground truth mapping of a given Beltrami coefficient  $\mu$ .

393 The loss function can be formulated as follow

394 (4.2) 
$$L_{Lapla} = \frac{1}{2N^2} (\|C_s \mathbf{s}\|_1 + \|C_t \mathbf{t}\|_1)$$

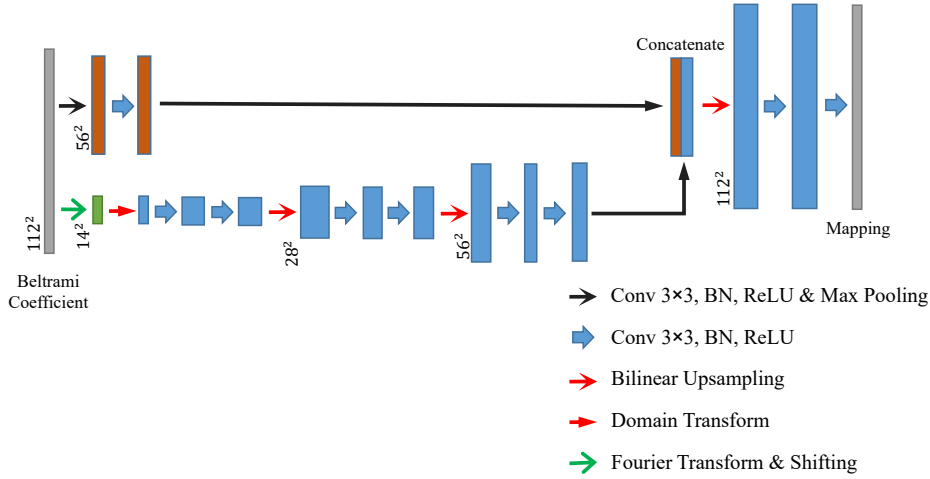


Figure 10: The architecture of the Beltrami Solver Network (BSNet).

395 (4.3)

$$L_{BSNet} = \gamma L_{Lapla}$$

396 Note that each row in  $C_s$  and  $C_t$  represents the relationship between a certain pixel and the  
 397 pixels adjacent to it. In the context of a triangular mesh, a pixel is adjacent to at most six  
 398 pixels. So each row in  $C_s$  and  $C_t$  has at most seven nonzero elements. Although the two  
 399  $N^2 \times N^2$  matrices  $C_s$  and  $C_t$  are sparse, both matrices can be rewritten as two dense arrays in  
 400 the implementation of our method, which means that the computation of  $L_{BSNet}$  is memory  
 401 saving and efficient.

402 **4.2. Beltrami Coefficient Estimator.** Another component in the proposed algorithm is  
 403 the Estimator Network which takes the source and target images as its input and generates  
 404 the corresponding Beltrami coefficient representing the distortion of the input images. The  
 405 network is based on U-Net with an activation added to its output. Its framework is shown in  
 406 Figure 11.

407 The estimator performs convolution and pooling to extract the deep spatial features of  
 408 the two input images. These deep spatial features are then up-sampled to be a two-channel  
 409 image  $\tilde{\mu}$  from which the norm and angle of  $\tilde{\mu}$  can be computed. Notice that when we train  
 410 Estimator, we assume that a BSNet has been trained so that for any  $\mu$  satisfies  $\|\mu\|_2 < 1$ , it  
 411 can convert the  $\mu$  to its corresponding mapping. Naturally,  $\mu$  generated by Estimator should  
 412 also satisfy this condition. In order to ensure  $\|\mu\|_2 < 1$ , we add a  $Tanh$  activation which takes  
 413  $\tilde{\mu}$  as its input.

$$414 \quad \|\mu\|_2 = Tanh(\|\tilde{\mu}\|_2)$$

415 where

$$416 \quad Tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$





$$428 \quad (4.8) \quad L_{Estimator} = \alpha L_F + \beta L_\mu + \eta L_{smooth}$$

429 where  $\alpha$ ,  $\beta$ , and  $\eta$  are hyper-parameters that give different weighting to the corresponding  
430 terms.

431 **4.3. Grid Parameterization and Spatial Transformation.** To make it convenient to per-  
432 form the pixel-wise transformation on images, we need to parametrize the image domain with  
433 coordinates systems. Notice that overall in this paper, the pixel would be referred to not  
434 only for elements in images but also for feature maps in the hidden layers of a neural net-  
435 work model. Also, we don't necessarily distinguish between image and feature map in this  
436 work. For a image with a height of  $H$  and width of  $W$ , to parameterize it within the domain  
437  $[0, 1]^2$ , we discretize the whole image domain by a grid  $G = \{(x_i, y_j)\} = \{(ih_x, jh_y) : i =$   
438  $1, 2, \dots, H; j = 1, 2, \dots, W\}$ , where  $h_x = \frac{1}{H+1}$  and  $h_y = \frac{1}{W+1}$ . In case the feature map contains  
439 multiple channels, the pixels for each channel located at the same position are parameterized  
440 by the same point in the grid.

441 In our implementation, the target image will always be parametrized with such a regular  
442 and normalized grid defined above. Then what we need to deform the image is the source  
443 coordinates on the input image. More cumulatively speaking, for the target coordinate  $(x_i^t, y_j^t)$ ,  
444 its corresponding pixel is the point whose pixel coordinate is  $(x_i^s, y_j^s)$  and lie in the input feature  
445 map. To summary, the output of the coefficient estimator and BSnet should be the coordinates  
446 of the source grid  $G_s = \{(x_i^s, y_j^s)\}$  which sampled on the input image.

447 Through such a setting, we can simply use the embedding function `grid_sample()` in  
448 `PyTorch`, which is also the standard parameterization method used in spatial transformer  
449 network[11] and texture mapping in computer graphics [7].

450 When the source coordinates are determined, generally its not exactly the same to any of  
451 the coordinates that parameterized the input image. Thus, interpolation for the resampling on  
452 points of source coordinates is must. The method for this interpolation should be differentiable  
453 to enable back-propagation. Associated with an input feature map  $I$  defined on  $(x, y) \in$   
454  $\{(x_p, y_l) : p = 1, \dots, H'; l = 1, \dots, W'\}$  and the output mapping of coefficient estimator  
455 and BSnet  $f : (x^t, y^t) \rightarrow (x^s, y^s)$ , we do resampling for  $J = I \circ f$  which is defined on  
456  $(x^s, y^s) \in \{(x_i^t, y_j^t) : i = 1, \dots, H; j = 1, \dots, W\}$ . We can write the resampling as:

$$457 \quad (4.9) \quad J(x_i^t, y_j^t) = \sum_{p=1}^{H'} \sum_{l=1}^{W'} I(x_p, y_l) k(x_i^s - x_p; \Phi_x) k(y_j^s - y_l; \Phi_y)$$

where  $i = 1, \dots, H$  and  $j = 1, \dots, W$

458 where  $\Phi_x$  and  $\Phi_y$  are the parameters of a generic sampling kernel  $k()$ ,  $G_0 = \{(x_p, y_l)\}$  is  
459 a regular and normalized grid that parameterized the image domain of  $I$ . The interpolation  
460 would be identical for each channel of the image. Let's fix the resampling method as bilinear,

461 then the resampling method can be written as

$$462 \quad (4.10) \quad J(x_i^t, y_j^t) = \sum_{p=1}^{H'} \sum_{l=1}^{W'} I(x_p, y_l) \max(0, 1 - \frac{|x_i^s - x_p|}{h_x}) \max(0, 1 - \frac{|y_j^s - y_l|}{h_y})$$

where  $i = 1, \dots, H$  and  $j = 1, \dots, W$

463 Thus, the derivative of the bilinear interpolation for output feature map  $J$  at  $(x_i^t, y_j^t)$  with  
464 respect to the input feature map and the source coordinates are given respectively by

$$465 \quad (4.11) \quad \frac{\partial J}{\partial I(x_p, y_l)}(x_i^s, y_j^s) = \sum_{p=1}^{H'} \sum_{l=1}^{W'} \max(0, 1 - \frac{|x_i^s - x_p|}{h_x}) \max(0, 1 - \frac{|y_j^s - y_l|}{h_y})$$

$$466 \quad (4.12) \quad \frac{\partial J}{\partial x^s}(x_i^s, y_j^s) = \sum_{p=1}^{H'} \sum_{l=1}^{W'} I(x_p, y_l) \max(0, 1 - \frac{|y_j^s - y_l|}{h_y}) h(x_i^s)$$

467 where

$$468 \quad (4.13) \quad h(x_i^s) = \begin{cases} 0 & \text{if } |x_p - x_i^s| \geq h_x \\ \frac{1}{h_x} & \text{if } x_p \geq x_i^s > x_p - h_x \\ -\frac{1}{h_x} & \text{if } x_p < x_i^s < x_p + h_x \end{cases}$$

469 and similarly to 4.12 and 4.13 for  $\frac{\partial J}{\partial y^s}$ .

470 Through such a differentiable resampling method, the gradients in the neural network  
471 model can pass backward and enable the updating of the parameters. With the derivative  
472  $\frac{\partial J}{\partial x^s}$  and  $\frac{\partial J}{\partial y^s}$ , the transformer layer can be trained. Through the gradients with respect to the  
473 input feature map in Equation.4.11, the gradients are able to pass to the previous layer that  
474 outputs it and leads to an end-to-end trainable neural network model.

475 **4.4. Quasi-Conformal Transformer Network.** The combination of the Coefficient Esti-  
476 mator, Beltrami Solver network, and image warper forms up the Quasi-conformal transformer  
477 layer(see Figure. 7). Coefficient Estimator takes the input feature map  $I$  and predicts the  
478 Beltrami coefficients  $\mu$  that is the quasi-conformal representation for the desired mapping  $f$ .  
479 The predicted Beltrami coefficients are solved into the mapping by the pre-trained BS-net.  
480 Directly by the mapping which is the source coordinates  $(x^s, y^s) = f(x^t, y^t)$  that lie on the  
481 input feature map, we use image warper to do resampling and obtain the transformed map.  
482 In this model, except for the Beltrami Solver Network, which should remain static after it ac-  
483 quired enough pretraining, every module in the quasi-conformal transformer layer is trainable.  
484 This enables a quasi-conformal transformer to be inserted in any position of the convolutional  
485 neural network and ends up with an end-to-end trainable neural network model.

486 The operation of what a quasi-conformal transformer layer did in a network is explicit  
487 and can be interpreted visually. Like the idea of spatial transformer network [11], our quasi-  
488 conformal transformer network can assign the input feature map a trainable warping that can

489 help the after layers work better and minimize the overall loss function. For example, in a clas-  
 490 sification task with disturbing images, a transformer layer can be inserted before the classifier  
 491 network to restore the disturbance and recover the semantic meaning for a better classifica-  
 492 tion result. Simply adding it to some specific tasks can also make sense. Like some popular  
 493 works aim for learning convolutions that differ from the regular rectangle windows[6][34], our  
 494 quasi-conformal transformer can also learn deformable convolution definition by warping the  
 495 feature map into a deformed one. In this way, doing regular convolution in the deformed map  
 496 is mathematically equivalent to assigning deformable convolution in the original feature map.

497 Beyond what a spatial neural network can do, quasi-conformal transformers are high-  
 498 lighted to learn dense and large transformations that are calculated point-wisely for the input  
 499 feature map. Though it would also be possible for a thin-plate spline (TPS) variant of the  
 500 spatial neural network, that directly predicts the source coordinates, our method can be much  
 501 easier to control the topology and the degree of the deformation benefit the Beltrami repre-  
 502 sentation. Controlling such topologic and geometric properties of a spatial transformation is  
 503 very important in learning a large and dense deformation. To be in more detail, learning a  
 504 mapping that is too flexible and without proper constraint may easily result in overfitting.  
 505 The loss can converge to a very low point easily due to the flexibility of the mapping but failed  
 506 to do prediction correctly in the testing. This will be evaluated thoroughly in Section. 5.3.

507 To achieve this, we carefully designed a penalty term that can constrain the mapping  
 508 to be a diffeomorphism and enhance its property of topology-preserving. Denote the quasi-  
 509 conformal transformer layer as  $N_{qct}$  and the other parts of the network as  $V$ , such a penalty  
 510 term is written as:

$$\begin{aligned}
 E_{reg}(V, N_{qct}) &= \int |\mu| + \int |\nabla\mu| \\
 (4.14) \quad &= \sum_{i=1}^H \sum_{j=1}^W \mu(x_i^t, y_j^t) + \sum_{i=1}^H \sum_{j=1}^W \nabla\mu(x_i^t, y_j^t)
 \end{aligned}$$

512 The penalty term should be put in the final loss function. Thus, for a particular task, the  
 513 overall cost function is:

$$(4.15) \quad E = E_{task} + \alpha E_{reg}$$

515 where  $E_{reg}$  is defined as Equation.4.14 weighted by  $\alpha$ , while  $E_{task}$  is the penalty term  
 516 associated with the task, e.g. cross-entropy error for classification or mean square error for  
 517 registration.

518 In the work of this paper, the

519 **5. Experimental Results.** In this section, we will thoroughly evaluate the capability of the  
 520 quasi-conformal transformer network and compare it with its most related existing work[11]  
 521 and its thin-plate spline-based variant mentioned in their work. Firstly, like the work of spatial  
 522 transformer, we test the power for localization on distorted versions of MNIST[19] handwriting  
 523 dataset. In this experiment, the images will be transformed affinely. Such deformed images  
 524 will only have characters on a small region of the image domains which makes it challenging  
 525 to do classification. Through our quasi-conformal transformer layer, the convolution will only

526 be performed on regions that contain information. The second experiment is to evaluate the  
527 elastic deformation restoration of the quasi-conformal transformer. To do this, we randomly  
528 deform the images from CIFAR10[14] and employ such deformed CIFAR10 to do supervised  
529 classification. We compare between two networks, one is a simple CNN classifier network  
530 without a quasi-conformal layer while another is the same classifier with the quasi-conformal  
531 transformer layer in its head. Our QCTN can restore the deformed image and result in  
532 better classification accuracy. In the third experiment, we will show experiments on Fashion  
533 MNIST[31] data that can do localization and elastic deformation restoration simultaneously.  
534 Lastly, on the original CIFAR10 dataset without any pre-processing, we put a QC transformer  
535 layer ahead of a classifier and do training without pretraining on the QC transformer layer.  
536 From the deformation by QC transformer on the feature map, we are able to do a learnable  
537 deformable convolution on images. The results promise quasi-conformal layers helped to  
538 improve the classification accuracy.

539 Our method was implemented in *Python* and run on CentOS-7 based central cluster nodes  
540 with a 2.4GHz Intel Xeon E5-2680 CPU, 64GB, and a GeForce GTX 1080 Ti GPU. The  
541 learning rate in the stage of classification is 0.00005. What is worthy to mention is, when we  
542 determine the parameter for learning rate, it's hard for a thin-plate spline spatial transformer  
543 network to converge. A parameter explosion and overfitting can easily occur during its train-  
544 ing. However, for a fair comparison, we take the learning rate to be 0.00005 uniformly in this  
545 work and do an extra clip for gradients especially when implementing the TPS variant for  
546 STN-CNN by 10.

547 **5.1. Localization on Deformed MNIST.** We take MNIST handwriting dataset to illus-  
548 trate the ability of the quasi-conformal transformer network for localizing the characters that  
549 are coarsely standing only in a partial region of the whole image domain. We first affinely  
550 deform the MNIST images. The deformation range is set to enable  $[-\frac{\pi}{3}, \frac{\pi}{3}]$  for rotation angle  $\theta$   
551 and  $[0.2, 0.6]$  for scaling parameter  $s$ . The translation is allowed as long as the characters will  
552 not move out of the image domains. We didn't test that for elastic deformation with MNIST  
553 since the elastic deformation on characters is not visually obvious. However, we do test that  
554 for elastic deformation with Fashion MNIST and will be presented in the later section. We  
555 trained three networks in this part including our quasi-conformal network to compare the  
556 performance. The other two models include the baseline convolutional neural network, which  
557 contains three convolutional layers and two fully-connected layers and the same network with  
558 spatial transformer (ST-CNN) or quasi-conformal transformer (QCT-CNN) layer inserted be-  
559 tween the input and the classifier CNN. The spatial transformer and quasi-conformal trans-  
560 former layer are both pretrained to obtain an initialization. In the process of pretraining, the  
561 transformer layer takes the deformed as input and the original image as the label to update  
562 the weightings in the module.

563 In the stage for the classification where the overall model is trained together, the opti-  
564 mization will run for 100 epochs. We evaluate the classification accuracy rate on the test  
565 dataset. Compared to the baseline convolutional neural network, both models equipped with  
566 transformer acquire better results as the accuracy rate given in Table. 1. From the visual-  
567 ized results shown in Figure. 12, it's obvious that the transformer succeeds in conferring the  
568 feature map and focuses only on regions that encircle the characters. However, our methods

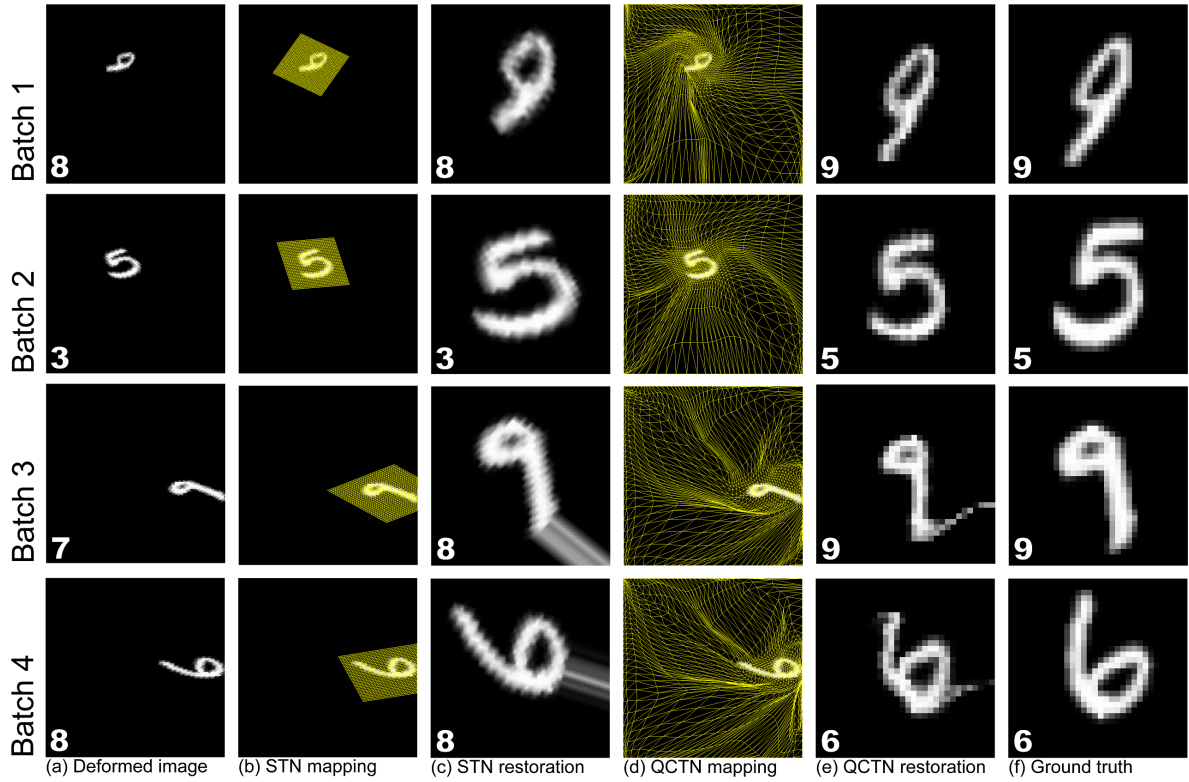


Figure 12: Classification result on affinely deformed MNIST. (a) the deformed images. (b) mapping generated by STN visualized on deformed image. (c) image localized by STN. (d) mapping generated by QCTN visualized on deformed image. (e) image localized by QCTN. (f) ground truth. The characters on the down-left of column (a)(c)(e) indicates the predicted result by baseline CNN, STN and QCTN respectively.

Table 1: The results for classifying affinely transformed MNIST by baseline CNN, STN and QCTN

Method	Train	Test
CNN	83.62	82.73
ST-CNN	94.97	94.90
QCT-CNN	96.45	96.32

569 outperform STN not only by quantity but also by the quality of images transformed. With  
 570 point-wise transformation by our QC transformer, the characters are more clean and sharp  
 571 with little blur. Particularly, when the character lies close to the boundary of the domain,  
 572 STN will generate a large blur because of bilinear interpolation on boundaries. However, our  
 573 method, which is a point-wise learnable mapping, can reduce such advantages efficiently.

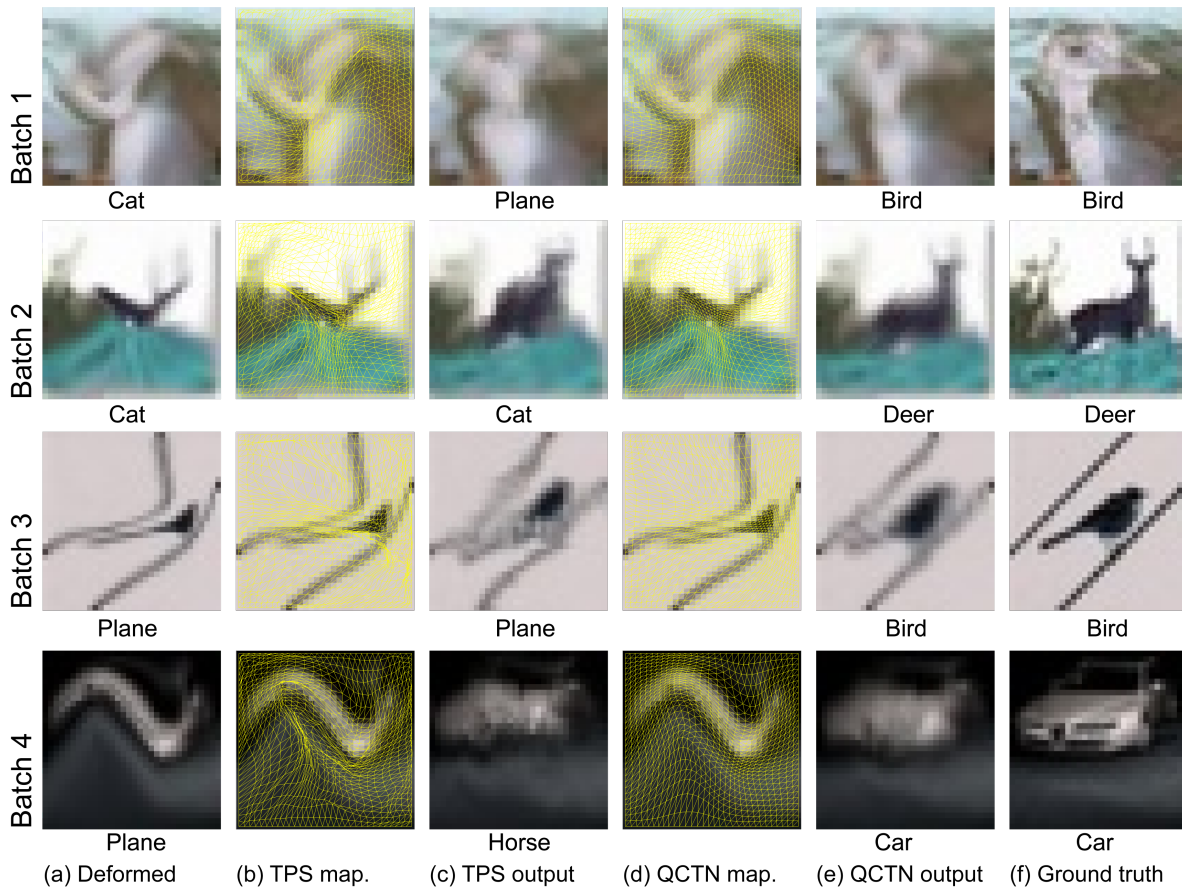


Figure 13: Classification result on CIFAR-10 with large non-rigid deformation. (a) the deformed images. (b) mapping generated by TPS-STN visualized on deformed image. (c) image recovered by TPS-STN. (d) mapping generated by QCTN visualized on deformed image. (e) image recovered by QCTN. (f) ground truth. The class names on the bottom of images in column (a)(c)(e) indicates the predicted class by baseline CNN, TPS-STN and QCTN respectively.

574 **5.2. Elastic Deformed CIFAR-10.** In this section, we perform an experiment to show that  
 575 the quasi-conformal transformer can learn to be invariant to elastic deformation. That's to say,  
 576 given a distorted image that may come from capturing across some uneven surfaces like glasses  
 577 or water, the quasi-conformal transformer layer can restore the image before it goes into the  
 578 classification network. Through such restoration, the original semantic meaning of the image  
 579 should be recovered. We assign the elastic deformation with different scales on the dataset  
 580 CIFAR10, which is a small low-resolution image recognition dataset, to obtain the deformed  
 581 dataset. For comparison, we also implemented a variant of spatial transformer network with  
 582 thin-plate spline transformation, which outputs the mapped coordinates of the control points  
 583 without constraint. The architecture for the localization network that predicts the mapped

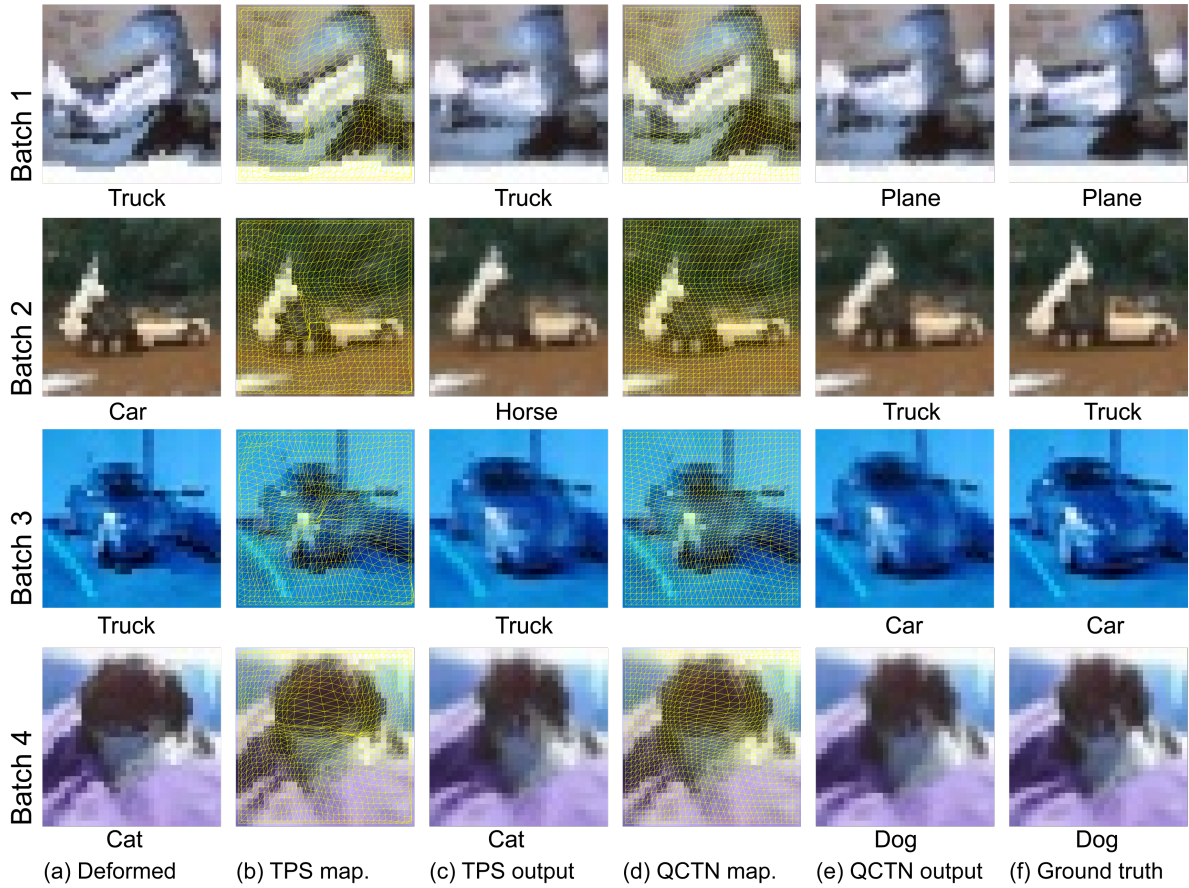


Figure 14: Classification result on CIFAR-10 with small non-rigid deformation. (a) the deformed images. (b) mapping generated by TPS-STN visualized on deformed image. (c) image recovered by TPS-STN. (d) mapping generated by QCTN visualized on deformed image. (e) image recovered by QCTN. (f) ground truth. The class names on the bottom of images in column (a)(c)(e) indicates the predicted class by baseline CNN, TPS-STN and QCTN respectively.

584 coordinates are the same as our coefficient estimator for a fair comparison. Though the  
 585 motivation of thin-plate spline (TPS-STN) and quasi-conformal transformer shares the same  
 586 motivation, the mapping generated by TPS-STN may not be a diffeomorphism and contains  
 587 self-foldings.

588 The parameters in the experiment are set as we discussed at the beginning of this section,  
 589 which is with a learning rate 0.00005 and clips the gradients value for the transformer layer  
 590 with 10. It's worthy to mention that, gradients clipping is necessary for TSP-STN but not  
 591 for our QCTN. Through we testing. The performance is similar for QCTN with or without  
 592 clipping. But for TSP-STN, training without clipping the gradients would result in parameter  
 593 explosion. We take the deep layer aggregation model(DLA,[32]) as the base classifier in this



Table 2: The results for classifying deformed CIFAR-10 with different scales by baseline CNN, TPS-STN and QCTN.

Method	Large Def.		Small Def.	
	Train	Test	Train	Test
CNN	94.55	75.84	99.13	83.06
STN-CNN	96.05	75.67	99.04	83.76
TPS-CNN	95.54	77.37	99.67	84.02
QCT-CNN	96.11	81.41	99.75	85.87

594 experiment.

595 Here we synthesize the deformation in two scales. Figure. 14 presents the small deforma-  
 596 tion, where some minor disturbance occurred. The large deformation is in Figure. 13, where  
 597 the semantic meaning of the image is hard to distinguish. As illustrated in Table.2, the base  
 598 CNN gets the lowest accuracy rate since the distortions bring some noise and degrade the  
 599 images with information loss. For the method with transformer layers, our quasi-conformal  
 600 transformer network obtain a better result than TPS-STN associated with mappings free of  
 601 self-folding. From the visualized figures of the test dataset, images recovered by QCTN are  
 602 more accurate and close to the original images, which helps human beings and neural network  
 603 model to recognize the classes each image belong to.

604 **5.3. Localization and Restoration on FashionMNIST.** The localization  $f_l$  and restora-  
 605 tion  $f_r$  can be composited into a single mapping  $f_c = f_r \circ f_l$  as it's still a diffeomorphism.  
 606 Thus, should be able to be learned in our Quasi-conformal transformer framework. In this  
 607 section, with the deformation consisting of both affine transformation and elastic deformation  
 608 assigned on the dataset FashionMNIST, we present the advantages of our QC-transformer  
 609 network that can solve the localization and the restoration simultaneously in a single QC  
 610 transformer layer. We also do a comparison with the spatial transformer network and its  
 611 thin-plate spline variant as well as the pure classifier without any quasi-conformal or spatial  
 612 transformer layers.

Table 3: The results for classifying deformed Fashion MNIST with different scales by baseline CNN, STN, TPS-STN and QCTN

Method	Large Def.		Small Def.	
	Train	Test	Train	Test
CNN	71.51	68.91	73.77	71.10
STN-CNN	80.86	76.11	82.57	78.83
TPS-CNN	84.29	80.82	86.13	83.59
QCT-CNN	84.18	83.84	86.48	85.79

613 The parameters set for this experiment are just like the previous, gradients clipping by 10

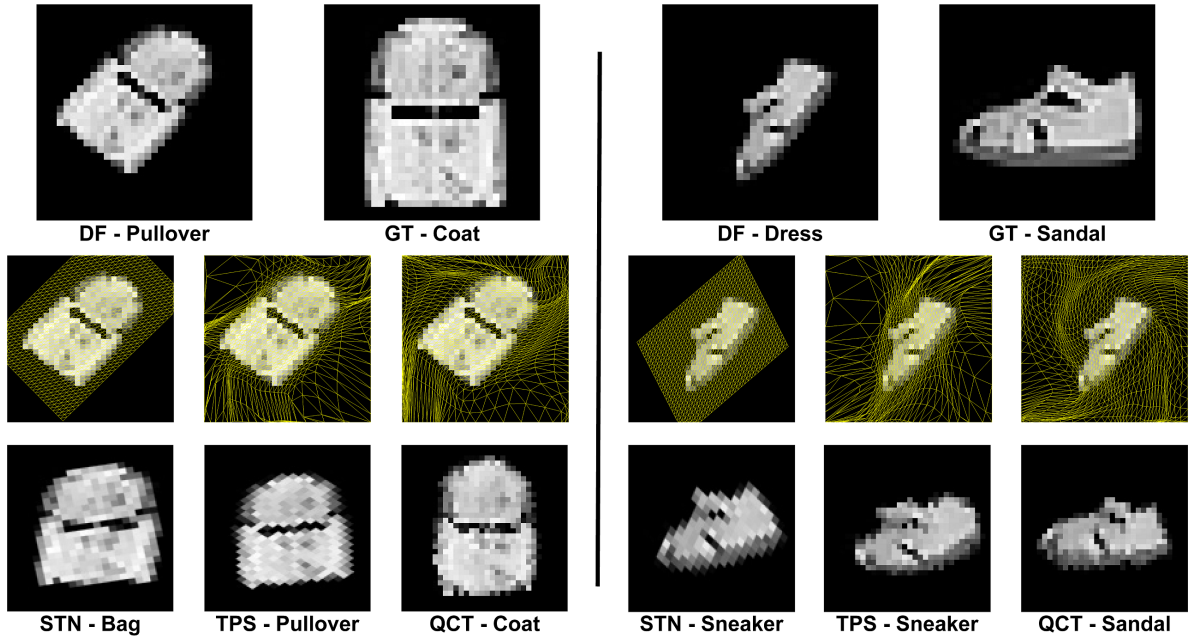


Figure 15: Classification result on Fashion MNIST with small non-rigid&affine deformation. Every batch follow the same illustration rule. As a example, DF-Bag present the deformed image and is predicted as bag by baseline CNN. GT-shirt indicates the ground truth image with the label as shirt. STN-Pullover, TPS-Sandal, QCT-shirt present the image recovered by STN, TPS-STN, QCTN and the predicted results respectively. The mapping on deformed image are visualized above them accordingly.

614 and the learning rate set to be 0.00005. Our methods outperform the other methods including  
 615 the basic STN and its TPS variant as well as the baseline classifier which is the same as that  
 616 in Section. 5.1.

617 As illustrated in Table. 3, our methods did the best on both deformation scales and achieve  
 618 around 3% and 2% higher accuracy than that of TPS-STN. Illustrated in Figure.16 for large  
 619 deformation and Figure.15 for small deformation, QCTN is able to recover the blurred and  
 620 distorted image the best while simultaneously localizing the main objects in the image.

621 **6. Conclusion and Future Work.** In this paper, we introduced the Quasi-conformal trans-  
 622 former network, which can be inserted into any part of a network. QCTN is capable to localize  
 623 the regions that are important for the tasks. Besides, for images that contain distortions that  
 624 may destroy the semantic meaning, QCTN can be trained to restore and recover the features  
 625 for accomplishing tasks. More than that, since QCTN is a self-contained module, it can be  
 626 inserted into any place of a neural network model to form up a new end-to-end trainable  
 627 network model. Through the learnable transformation by QCTN, a deformable convolution  
 628 can be performed by convoluting with regular rectangle on deformed feature map which is  
 629 mathematically equivalent. Through the experiments, QCTN is proven to be able to assign  
 630 not only spatial-invariant but also elastic deformation invariant to a network model. Besides,

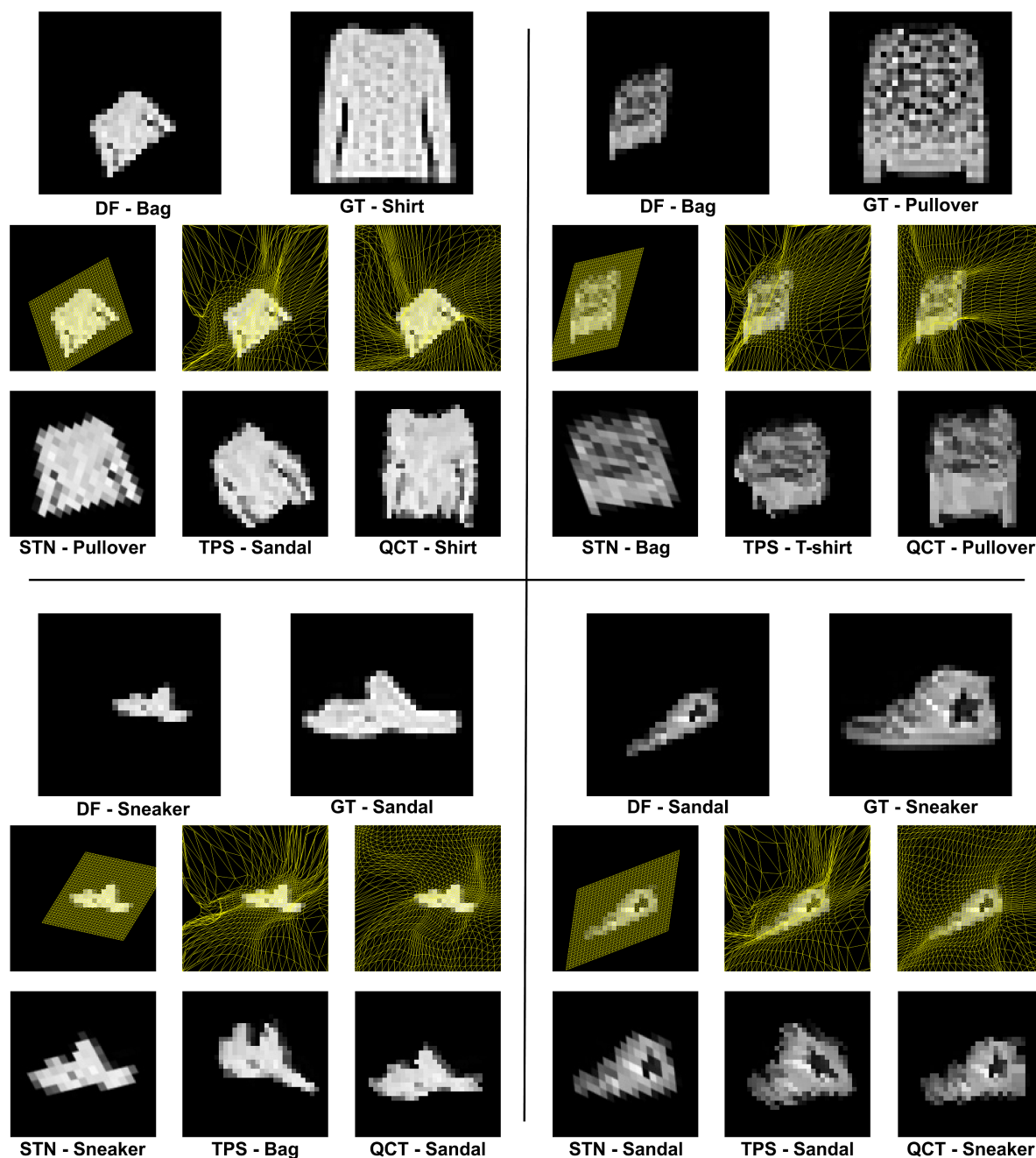


Figure 16: Classification result on Fashion MNIST with large non-rigid&affine deformation. Every batch follow the same illustration rule as Figure.15

631 compared to the similar work spatial transformer network and its thin-plate spline variant, our  
 632 model is a point-wisely learnable transformation whose topology is guaranteed to be preserved

633 by the proper constraint Beltrami coefficients. In our following works, we plan to extend the  
 634 application of the Quasi-conformal transformer network. For example, on registration and  
 635 segmentation. Besides, it would be interesting to assign bijectivity directly on the deformable  
 636 convolution network[6] to see if preserving the topology of the deformed filter can help with  
 637 the performance. Also, designing a bijective transformer for 3D volume deep learning would  
 638 also be valuable, especially for medical images.

639

## REFERENCES

- 640 [1] Q. CHEN, Z. LI, AND L. M. LUI, *A learning framework for diffeomorphic image registration based on*  
 641 *quasi-conformal geometry*, CoRR, abs/2110.10580 (2021), <https://arxiv.org/abs/2110.10580>, <https://arxiv.org/abs/2110.10580>.  
 642
- 643 [2] G. CHOI, Y. LIU, AND L. M. LUI, *Free-boundary conformal parameterization of point clouds*, arXiv  
 644 preprint arXiv:2010.15399, (2020).
- 645 [3] G. P.-T. CHOI, K. T. HO, AND L. M. LUI, *Spherical conformal parameterization of genus-0 point clouds*  
 646 *for meshing*, SIAM Journal on Imaging Sciences, 9 (2016), pp. 1582–1618.
- 647 [4] P. T. CHOI AND L. M. LUI, *Fast disk conformal parameterization of simply-connected open surfaces*,  
 648 Journal of Scientific Computing, 65 (2015), pp. 1065–1090.
- 649 [5] K. CRANE, U. PINKALL, AND P. SCHRÖDER, *Robust fairing via conformal curvature flow*, ACM Trans.  
 650 Graph., 32 (2013).
- 651 [6] J. DAI, H. QI, Y. XIONG, Y. LI, G. ZHANG, H. HU, AND Y. WEI, *Deformable convolutional networks*,  
 652 in iccv, 2017, pp. 764–773.
- 653 [7] J. D. FOLEY, A. VAN DAM, S. K. FEINER, J. F. HUGHES, AND R. L. PHILLIPS, *Introduction to computer*  
 654 *graphics*, vol. 55, Addison-Wesley Reading, 1994.
- 655 [8] X. GU, Y. WANG, T. F. CHAN, P. M. THOMPSON, AND S.-T. YAU, *Genus zero surface conformal*  
 656 *mapping and its application to brain surface mapping*, IEEE Transactions on Medical Imaging, 23  
 657 (2004), pp. 949–958.
- 658 [9] X. GU AND S.-T. YAU, *Global conformal surface parameterization*, in Proceedings of the 2003 Euro-  
 659 graphics/ACM SIGGRAPH symposium on Geometry processing, Eurographics Association, 2003,  
 660 pp. 127–137.
- 661 [10] K. HE, X. ZHANG, S. REN, AND J. SUN, *Deep residual learning for image recognition*, in cvpr, 2016,  
 662 pp. 770–778.
- 663 [11] M. JADERBERG, K. SIMONYAN, A. ZISSERMAN, ET AL., *Spatial transformer networks*, Advances in neural  
 664 information processing systems, 28 (2015), pp. 2017–2025.
- 665 [12] Y. JEON AND J. KIM, *Active convolution: Learning the shape of convolution for image classification*, in  
 666 cvpr, 2017, pp. 4201–4209.
- 667 [13] N. KANOPOULOS, N. VASANTHAVADA, AND R. L. BAKER, *Design of an image edge detection filter using*  
 668 *the sobel operator*, IEEE Journal of solid-state circuits, 23 (1988), pp. 358–367.
- 669 [14] A. KRIZHEVSKY, *Learning multiple layers of features from tiny images*, tech. report, 2009.
- 670 [15] A. KRIZHEVSKY, I. SUTSKEVER, AND G. E. HINTON, *Imagenet classification with deep convolutional*  
 671 *neural networks*, Advances in neural information processing systems, 25 (2012), pp. 1097–1105.
- 672 [16] K. C. LAM AND L. M. LUI, *Landmark-and intensity-based registration with large deformations via quasi-*  
 673 *conformal maps*, SIAM Journal on Imaging Sciences, 7 (2014), pp. 2364–2392.
- 674 [17] Y. LECUN, Y. BENGIO, ET AL., *Convolutional networks for images, speech, and time series*, The hand-  
 675 book of brain theory and neural networks, 3361 (1995), p. 1995.
- 676 [18] Y. LECUN, L. BOTTOU, Y. BENGIO, AND P. HAFFNER, *Gradient-based learning applied to document*  
 677 *recognition*, Proceedings of the IEEE, 86 (1998), pp. 2278–2324.
- 678 [19] Y. LECUN, C. CORTES, AND C. BURGES, *Mnist handwritten digit database*, AT&T Labs [Online]. Available:  
 679 <http://yann.lecun.com/exdb/mnist>, 2 (2010).
- 680 [20] B. LÉVY, S. PETITJEAN, N. RAY, AND J. MAILLOT, *Least squares conformal maps for automatic texture*  
 681 *atlas generation*, ACM transactions on graphics (TOG), 21 (2002), pp. 362–371.
- 682 [21] L. M. LUI, K. C. LAM, T. W. WONG, AND X. GU, *Texture map and video compression using beltrami*

- 683 *representation*, SIAM Journal on Imaging Sciences, 6 (2013), pp. 1880–1902.
- 684 [22] L. M. LUI, K. C. LAM, S.-T. YAU, AND X. GU, *Teichmüller mapping ( $t$ -map) and its applications to*  
685 *landmark matching registration*, SIAM Journal on Imaging Sciences, 7 (2014), pp. 391–426.
- 686 [23] W. LUO, Y. LI, R. URTASUN, AND R. ZEMEL, *Understanding the effective receptive field in deep convo-*  
687 *lutional neural networks*, in Proceedings of the 30th International Conference on Neural Information  
688 Processing Systems, 2016, pp. 4905–4913.
- 689 [24] P. C. NG AND S. HENIKOFF, *Sift: Predicting amino acid changes that affect protein function*, Nucleic  
690 acids research, 31 (2003), pp. 3812–3814.
- 691 [25] J. PANETTA, M. KONAKOVIĆ-LUKOVIĆ, F. ISVORANU, E. BOULEAU, AND M. PAULY, *X-shells: A new*  
692 *class of deployable beam structures*, ACM Transactions on Graphics (TOG), 38 (2019), pp. 1–15.
- 693 [26] J. M. PREWITT, *Object enhancement and extraction*, Picture processing and Psychopictorics, 10 (1970),  
694 pp. 15–19.
- 695 [27] K. SIMONYAN AND A. ZISSERMAN, *Very deep convolutional networks for large-scale image recognition*,  
696 arXiv preprint arXiv:1409.1556, (2014).
- 697 [28] C. Y. SIU, H. L. CHAN, AND R. L. MING LUI, *Image segmentation with partial convexity shape prior*  
698 *using discrete conformality structures*, SIAM Journal on Imaging Sciences, 13 (2020), pp. 2105–2139.
- 699 [29] Y. SOLIMAN, D. SLEPČEV, AND K. CRANE, *Optimal cone singularities for conformal flattening*, *tog*, 37  
700 (2018).
- 701 [30] C. SZEGEDY, W. LIU, Y. JIA, P. SERMANET, S. REED, D. ANGUELOV, D. ERHAN, V. VANHOUCHE,  
702 AND A. RABINOVICH, *Going deeper with convolutions*, in *cvpr*, 2015, pp. 1–9.
- 703 [31] H. XIAO, K. RASUL, AND R. VOLLGRAF, *Fashion-mnist: a novel image dataset for benchmarking machine*  
704 *learning algorithms*, CoRR, abs/1708.07747 (2017), <http://arxiv.org/abs/1708.07747>, <https://arxiv.org/abs/1708.07747>.
- 705 [32] F. YU, D. WANG, E. SHELFHAMER, AND T. DARRELL, *Deep layer aggregation*, in *cvpr*, 2018, pp. 2403–  
706 2412.
- 707 [33] D. ZHANG AND L. M. LUI, *Topology-preserving 3d image segmentation based on hyperelastic regulariza-*  
708 *tion*, Journal of Scientific Computing, 87 (2021), pp. 1–33.
- 709 [34] X. ZHU, H. HU, S. LIN, AND J. DAI, *Deformable convnets v2: More deformable, better results*, in *cvpr*,  
710 2019, pp. 9308–9316.
- 711