
Harmformer: Harmonic Networks Meet Transformers for Continuous Roto-Translation Equivariance

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Convolutional Neural Networks exhibit inherent equivariance to image translation,
2 leading to efficient parameter and data usage, faster learning, and improved ro-
3 bustness. The concept of translation equivariant networks has been successfully
4 extended to rotation transformation using group convolution for discrete rotation
5 groups and harmonic functions for the continuous rotation group encompassing
6 360° . We explore the compatibility of the Self-Attention mechanism with full rota-
7 tion equivariance, in contrast to previous studies that focused on discrete rotation.
8 We introduce the Harmformer, a harmonic transformer with a convolutional stem
9 that achieves equivariance for both translation and continuous rotation. Accom-
10 panied by an end-to-end equivariance proof, the Harmformer not only outperforms
11 previous equivariant transformers, but also demonstrates inherent stability under
12 any continuous rotation, even without seeing rotated samples during training.

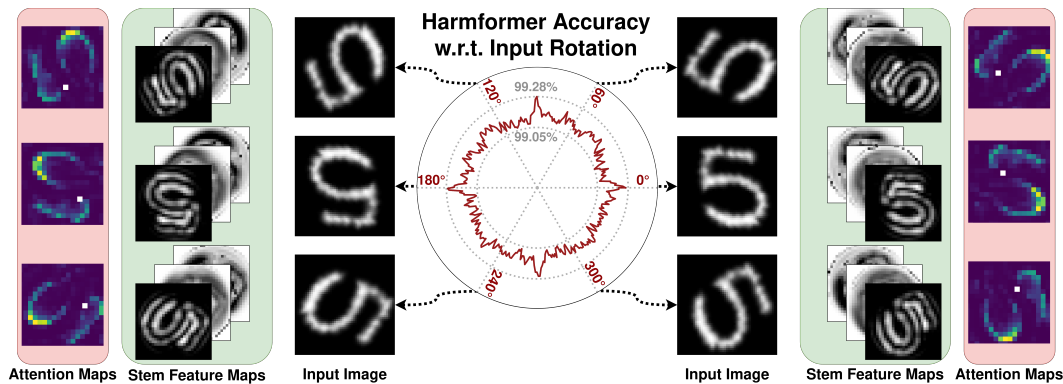


Figure 1: Equivariance of the Harmformer feature and attention maps in response to rotation of the input image: While the maps themselves are rotated, the magnitudes in the maps remain the same.

13 1 Introduction

14 A key strength that positions Convolutional Neural Networks (CNNs) [1] as a superior architecture
15 for computer vision tasks is the weight sharing across the spatial domain. This design ensures that
16 CNN feature maps retain their values as the input is translated, only being shifted according to
17 the input. Formally known as translation equivariance, this property provides CNNs with inherent
18 robustness and efficiency in managing translations. Equivariance can be extended to other groups
19 of transformations, such as rotation, scaling, or mirroring. The advantage of equivariant models is

20 that they ensure a tractable response of the model to the transformation of the input. As a result, the
21 model can eliminate the effects of the transformations and produce predictions that are invariant to
22 them. For instance, to achieve translation invariance in conventional CNNs, the feature maps are
23 commonly aggregated by global average pooling before the classification layer.

24 Group Equivariant Convolutional Neural Networks (G-CNNs) [2] show that CNNs can be modified
25 to become equivariant to any discrete transformation group, such as rotation by a discrete set of
26 angles. An extension to continuous rotation and translation group was introduced by Worrall et al.
27 [3]. The authors proposed Harmonic Networks (H-Nets), which restrict the convolution filters to
28 a family of harmonic functions ideal for expressing full rotation equivariance. Both approaches
29 improve the generalization and efficiency of training for the chosen group, similarly to how CNNs
30 benefit from translation equivariance. For example, rotation equivariant networks are well suited
31 for object detection in aerial imagery because such images lack natural orientation and equivariant
32 networks inherently accommodate all rotations. Beyond aerial imagery [4, 5], equivariant CNNs are
33 effective in many other applications, such as microscopy [6, 7], histology [8], and remote sensing [9].

34 With the adoption of transformer architectures in computer vision, the Self-Attention mechanism
35 has also been integrated into equivariant networks [10, 11, 12]. Equivariant transformers are gaining
36 importance especially in domains such as graph-based structures (e.g. molecules) [13, 14, 15, 16],
37 vector fields [17], manifolds [18], and generic geometric data [19, 20]. In the 2D domain, Romero
38 and Cordonnier [21] proposed a transformer equivariant to discrete rotation and translation groups by
39 using the principle of G-CNNs in the positional encoding of the Self-Attention (SA). The formulation
40 was further improved by Xu et al. [22]. In both cases, the computational complexity of the equivariant
41 SA increases quadratically with the number of angles in the considered rotation group, which limits
42 the model angular resolution. Equivariance to continuous rotation presents a versatile solution.

43 In this paper, we introduce Harmformer, the first vision transformer capable of achieving continuous
44 2D roto-translation equivariance. The name is derived from circular harmonics [23] which provide
45 the equivariance property preserved throughout the architecture. To ensure computational efficiency,
46 our network starts with an equivariant convolutional stem based on Harmonic networks [3], where
47 we redesign the key components, such as activations, normalization layers, and introduce equivariant
48 residual connections. The stem output is divided into equivariant patches, which are then passed to
49 the transformer. Alongside a novel self-attention SA mechanism, we introduce layer normalization
50 and linear layers to guarantee end-to-end equivariance. The equivariance property allows Harmformer
51 to remove the effect of roto-translation just before classification, preserving all relevant information
52 at earlier stages (see Fig. 1).

53 Through experimental validation, we show that Harmformer surpasses all previous discrete equivariant
54 transformers [21, 22] on established benchmarks [24, 25, 26]. It also outperforms earlier invariant
55 models [27, 28, 29] on classification tasks where the model is trained solely on non-rotated data.

56 2 Related Work

57 We review the three foundational concepts from prior research that Harmformer builds upon: the
58 SA mechanism, equivariant convolution networks, and transformers with a convolutional stem stage.
59 Additionally, we discuss other equivariant transformer architectures.

60 **Visual Self-Attention** The well-known SA mechanism originates from natural language processing
61 [30] and is widely used in computer vision since the publication of the Visual Transformer (ViT)
62 [31]. Transformers, unlike CNNs, exhibit larger model capacities but require substantial amounts
63 of data and have quadratic complexity with respect to input size. Transformers closely related to
64 Harmformer include CoAtNet [32] and, more specifically, ViT_p [33]. These architectures begin with
65 a convolution stem to downscale the input and thereby reduce the computational complexity of the
66 subsequent application of SA. However, these architectures are not equivariant to roto-translation.

67 **Equivariant Convolutions** Since the publishing of the G-CNNs [2], the concept of equivariant
68 convolutional networks has expanded across various modalities and transformation groups. In 2D,
69 these transformations include rotation [3, 34], scaling [35, 36, 37], and general $E(2)$ transformations
70 [38]. In 3D, applications cover $SO(3)$ transformations in volumetric data [39, 40] and point clouds
71 [41], as well as spherical CNNs [42]. Equivariant networks are also applied to graphs [43] and
72 non-Euclidean manifolds [44]. Harmformer builds on and extends the H-Nets published by Worrall

73 et al. [3], which are purely convolutional networks equivariant to continuous 360° rotation. In our
 74 implementation of H-Nets, we incorporate the improvements introduced in H-NeXt [27].

75 **Equivariant Transformers** As previously mentioned, equivariant networks have integrated the SA
 76 mechanism in various domains, including 3D graphs and point clouds using irreducible representations
 77 [10, 15, 14], operations on Lie algebras [12], and general geometric data using geometric algebras
 78 [19, 20]. Particularly relevant to our work are the planar roto-translation equivariant transformers,
 79 such as Group Equivariant Self-Attention Networks (GSA-Nets) [21] and $E(2)$ -Equivariant Vision
 80 Transformer (GE-ViT), which reformulate relative positional encoding to construct equivariant
 81 transformers. However, GSA-Nets and GE-ViT operate only on discrete rotation groups such as
 82 $\{k\frac{\pi}{2} \mid k \in \mathbb{Z}\}$, where finer angular sampling substantially increases the computational complexity.

83 3 On Equivariance in Vision Transformers

84 We analyze the roto-translation equivariance of the ViT architecture, a well-known representative
 85 of vision transformers. First, we formalize the notion of equivariance. Intuitively, a function f
 86 is equivariant to a transformation a_g if the transformation and the function commute, $f(a_g(x)) =$
 87 $a_g(f(x))$. For example, processing a rotated input image has the same effect as directly rotating the
 88 features of the unrotated image. In practice, such a definition would be too restrictive. The function
 89 f (layer or network) typically has a different domain and codomain, so the transformation may act
 90 differently on each. The core idea remains the same: the model response to the input transformation is
 91 predictable. To formally define equivariance, we draw upon the seminal work of Cohen and Welling
 92 [2] or the more recent one formulated by Weiler et al. [45].

93 **Definition 3.1** (Equivariance). A function $f : X \rightarrow Y$ (a whole network or a single layer) is called
 94 group equivariant with respect to a group G if for every element g in G , represented by a linear map
 95 $a_g : X \rightarrow X$, there exists a corresponding linear map $b_g : Y \rightarrow Y$ such that the following holds:

$$f(a_g(x)) = b_g(f(x)) \quad \text{for all } x \in X \text{ and } g \in G. \quad (1)$$

96 A composition $f_2(f_1(x))$ of two equivariant functions $f_1 : X \rightarrow Y$ and $f_2 : Y \rightarrow Z$ is equivariant.
 97 We call invariance a special case of equivariance when b_g is the identity for all g in G .

98 **Self-Attention** A key mechanism that distinguishes transformers from previous architectures is the
 99 SA layer [30]. Before discussing the properties of SA, let us formally define it.

100 **Definition 3.2** (Self-Attention). Given an input matrix $Y \in \mathbb{R}^{n \times d}$, where each row of Y represents a
 101 feature vector of dimension d , usually called a patch. The matrices Q (queries), K (keys), and V
 102 (values) are computed as linear projections of Y :

$$[Q, K, V] = [YW_q, YW_k, YW_v] \quad W_{q,k,v} \in \mathbb{R}^{d \times d_h}, \quad (2)$$

103 where d_h is the dimension within the SA layer. The output of the self-attention layer, $\text{SA}(Y)$, is a
 104 weighted sum of the vectors in V , where the weights are defined as the softmax-normalized pairwise
 105 similarity scores between the vectors in Q and K :

$$A = \text{softmax}(QK^T / \sqrt{d_h}) \quad A \in \mathbb{R}^{n \times n}, \quad (3)$$

$$\text{SA}(Y) = AV. \quad (4)$$

106 In practice, SA is typically extended to Multi-Head Self-Attention (MSA), in which multiple SA
 107 layers with different embedding matrices $W_{(k,v,q)}$ are computed in parallel and then combined.

108 The construction of SA implies its well-known property, in the literature often referred to as permu-
 109 tation invariance [46]. According to Def. 3.2, it is more accurate to call the SA layer permutation
 110 equivariant rather than permutation invariant. For if we change the order of the rows in Y , the $\text{SA}(Y)$
 111 remains the same except for the same change in the order of its output rows.

112 **What makes ViT non-equivariant?** As rotation and translation are special cases of permutation,
 113 the permutation equivariance of SA might suggest the roto-translation equivariance of the whole
 114 ViT. However, the permutation-equivariance of SA holds at the patch level and not at the pixel level,
 115 where translation or rotation takes place. In the initial stage of ViT, before the first SA layer, the
 116 image is split into n non-overlapping patches of fixed size, typically 16×16 pixels. These are linearly
 117 transformed and flattened to form the rows of the input matrix Y of the first SA layer. This patch-wise

118 operation breaks the direct rotation (or translation) equivariance of ViT at the image level, because
 119 for an image rotation a_g in Eq. (1) corresponding to an angle g from the rotation group G , there is no
 120 b acting on the patches that can be expressed as a rotation b_g by g ; the same holds for translation.

121 A solution typically used by previous equivariant approaches [21, 22] is to consider pixel-level
 122 “patches” of size 1×1 pixel. Then the image rotation is equivalent to patch-level permutation, and
 123 the corresponding transformer model remains equivariant, assuming that interpolation errors and
 124 boundary effects are minimal. This approach, however, has two major drawbacks. As seen from
 125 Eq. (3), the SA has a quadratic complexity with respect to the number of patches n , so operating
 126 on a pixel grid (as opposed to 16×16) incurs an almost 10^5 penalty factor in memory requirements
 127 and correspondingly increases the processing time. To mitigate this, GSA-Nets and GE-ViT reduce
 128 complexity by using local SA [47] that restricts the attention field to the 7×7 neighborhood of the
 129 patch. The second drawback is that the local self-attention in the first layers is not very informative,
 130 because nearby pixels are usually highly correlated.

131 **Position Encoding** During construction of the input matrix Y for the first SA layer, the patches are
 132 also given absolute position encoding that provides information about their locations. This breaks
 133 equivariance as the patches of a transformed image will receive different encoding compared to their
 134 counterparts in the original image. Equivariant transformers [21, 22] replace the absolute encoding
 135 with circular relative encoding, similar to iRPE introduced by Wu et al. [48].

136 In Harmformer, we address these challenges with a convolutional stem stage that initially reduces
 137 spatial dimensions and extracts high-level features. Subsequently, we create 1×1 patches from these
 138 high-level features and process them by the SA layers. To maintain spatial correspondence among
 139 the patches while ensuring equivariance, Harmformer also uses circular relative position encoding.

140 4 Harmonic Convolutions and Equivariance to Continuous Roto-Translation

141 To understand the equivariance property of Harmformer, it is essential to understand the concept of
 142 harmonic convolutions introduced in H-Nets [3], as they are employed in the stem and affect the
 143 subsequent transformer layers. The main difference from the traditional CNNs is that the convolution
 144 filters based on circular harmonic functions are specifically designed to encode rotational symmetries.
 145 The filters are defined as follows.

146 **Definition 4.1** (Harmonic Filter). A harmonic filter $W_m : \mathbb{R}^2 \rightarrow \mathbb{C}$ parameterized by a rotation order
 147 m is given by:

$$W_m(r, \theta) = R(r)e^{i(m\theta + \beta)}, \quad (5)$$

148 where (r, θ) are polar coordinates. Here, $R : \mathbb{R} \rightarrow \mathbb{R}$ is a learnable radial function and $\beta \in \mathbb{R}$ is a
 149 learnable phase shift. The rotation order m is a parameter that determines the filter symmetry.

150 As translation equivariance is inherently provided by convolution, we will focus solely on rotation in
 151 the following discussion and denote the rotation operator by $[\cdot]^\alpha$, where α is the angle of rotation. Let
 152 us look in detail at how the rotation applied to the input affects feature maps generated by harmonic
 153 convolution. The H-Nets features are represented as complex values in polar form.

154 **Lemma 4.1** (Harmonic Convolution Property). Let I be an input image and W_{m_1} a harmonic filter.
 155 Under image rotation by angle α , convolution of I with W_{m_1} is given by:

$$[I]^\alpha \otimes W_{m_1} = e^{im_1\alpha} [I \otimes W_{m_1}]^\alpha. \quad (6)$$

156 This equation shows that rotating the image only results in a phase shift of the feature values, while
 157 the spatial coordinates are rotated accordingly. This property also holds for subsequent convolution
 158 layers. If the first feature map is denoted as $F_{m_1}([I]^\alpha)$, then convolution with another harmonic filter
 159 W_{m_2} is given by:

$$F_{m_1}([I]^\alpha) \otimes W_{m_2} = e^{i(m_1+m_2)\alpha} [F_{m_1}(I) \otimes W_{m_2}]^\alpha. \quad (7)$$

160 The authors of H-Nets also construct activation, batch normalization, and pooling layers that preserve
 161 this property. As a result, their classifier can be independent of input rotation and translation. To
 162 remove the influence of rotation, they extract only the magnitude from the last feature map and
 163 discard the phase. To aggregate spatial information, they use global average pooling. Note that for

164 tasks where the orientation of the object is relevant, phase can be used as no information is lost due
 165 to equivariance.

166 To unify the equivariance property within the Harmformer architecture, we define Harmonic Equivari-
 167 ance (HE), which is motivated by Lemma 4.1 and satisfies the general definition of equivariance
 168 (Def. 3.1). HE describes how features transform with respect to the rotation of an input image. By
 169 showing that each Harmformer layer satisfies HE, we establish the relationship between the features
 170 and the rotation of an input throughout the model.

171 **Definition 4.2** (Harmonic Equivariance – HE). A layer $F_m(\cdot)$ associated with a rotation order m is
 172 said to be HE, if for any rotation by angle α and admissible input I , it is transformed as follows:

$$F_m([I]^\alpha) = e^{im\alpha} [F_m(I)]^\alpha. \quad (8)$$

173 Here $[F_m(I)]^\alpha$ are features obtained from an unrotated input I and then rotated. The phase is shifted
 174 by a multiple of the rotation angle, where the factor is given by the rotation order of the layer. The
 175 process is illustrated in Fig. 3a.

176 5 Harmformer Architecture

177 The architecture of Harmformer is shown in Figure 2 and its layers will be discussed one by one.
 178 HE (Def. 4.2) of each layer is proved in Appendix A, demonstrating the end-to-end continuous
 179 rotation and translation equivariance. The architecture begins with a stem stage based on H-Nets,
 180 which we have further improved by refining activation and normalization layers and incorporating
 181 residual connections. The stem is followed by an equivariant encoder tailored to maintain HE, and
 182 the last component is a classifier, which takes the HE output of the encoder and computes an invariant
 183 representation for classification.

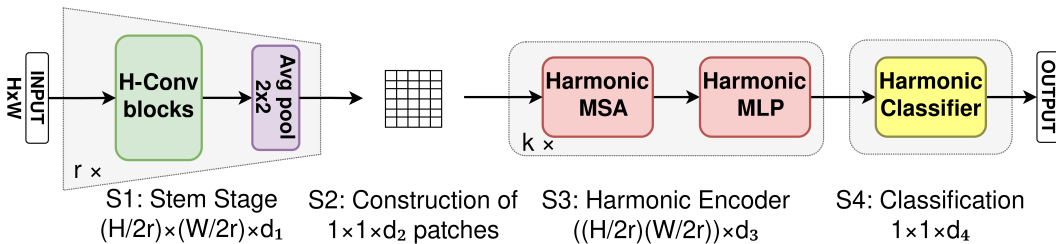


Figure 2: Overview of the Harmformer architecture, divided into four stages: S1 - downscaling the input, S2 - constructing patches from feature maps, S3 - Harmonic Encoder, and S4 - Classifier.

184 5.1 Harmformer: S1 Stem Stage

185 The main role is to prepare features for the Harmonic Encoder (S3) so that they are HE and have
 186 lower spatial resolution to keep the computational complexity of SA manageable, as discussed in Sec.
 187 3. To this end, we design the stage to comprise r iterations of H-Conv blocks, followed by average
 188 pooling, as shown in Fig. 2. Each iteration increases the number of channels while decreasing the
 189 spatial dimension.

190 The stage starts with an input that formally satisfies HE for the rotation order $m = 0$, expressed as
 191 $[I]^\alpha = e^{i0\alpha} [I]^\alpha$, followed by the first H-Conv block shown in Figure 3b.

192 **Rotation Order Streams** The HE and the definition of Harmonic Convolution have already been
 193 detailed in Lemma 4.1 and Def. 4.1. An important aspect that remains to be addressed is the selection
 194 of rotation orders for the harmonic filters. In our initial convolution with the input image (often called
 195 lifting convolution), we employ harmonic filters of rotation orders $-1, 0$, and 1 . This setup produces
 196 three streams of feature maps, each corresponding to one of these rotation orders.

197 Our experiments, along with the results reported in [27, 3], indicate that generating feature maps of
 198 higher rotation orders does not significantly improve performance but increases the computational
 199 complexity. Based on this evidence, we limit rotation orders to $-1, 0$, and 1 .

200 Most layers process these streams independently and those that interact across streams are indicated
 201 in the diagrams by a "spoon" symbol, as in the case of Harmonic Convolution in Figure 3b. Streams

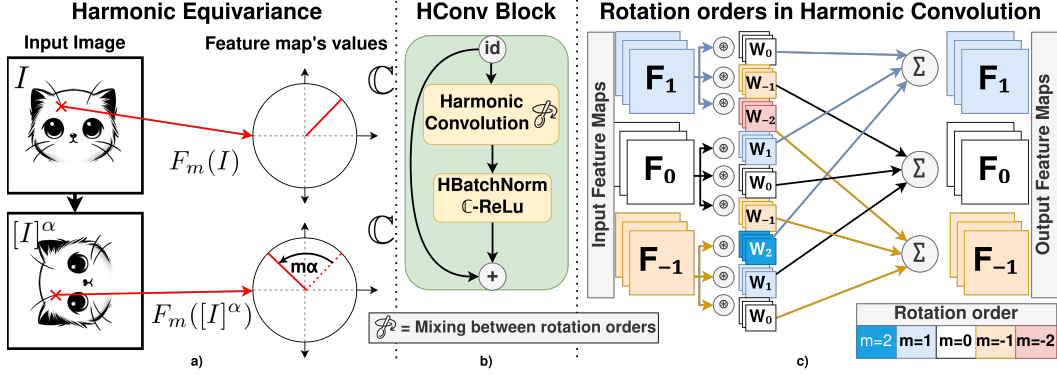


Figure 3: (a) Phase shift of HE feature values when the input is rotated; (b) Harmonic Convolution (H-Conv) Block of the stem stage; (c) Interaction of harmonic filters W_m with feature maps F_m within the Harmonic Convolution layer of the H-Conv Block, where m is the rotation order.

202 in the Harmonic Convolution block are mixed similarly as in H-Nets. The proposed mixing strategy
 203 is shown in Figure 3c and follows from the Harmonic Convolution property in (7), which states that

$$F_m^{out} = \sum_{m=m_1+m_2} F_{m_1}^{in} \otimes W_{m_2}, \quad (9)$$

204 where m , m_1 , and m_2 are the rotation orders of the output, input, and harmonic filter, respectively.

205 **Layers Operating on Magnitude** Because rotation affects only the phase of the features leaving the
 206 magnitude untouched, element-wise functions, such as normalization or activation, operating only
 207 on magnitudes preserve the HE property. In contrast with previous H-Nets [27, 3], we restrict the
 208 codomain of every element-wise function f transforming magnitudes to non-negative numbers, $f : \mathbb{R} \rightarrow \mathbb{R}_0^+$,
 209 since negative magnitudes inadvertently flip the phase, thus violating the HE property. This
 210 consideration leads us to propose a novel normalization fused together with activation (HBatchNorm
 211 and C-ReLu), detailed in Appendix A.4. Restricting the codomain and fusing the normalization with
 212 the activation has a positive impact on performance, as shown in Ablation B.1.

213 **Residual Connection** The final stem element is the residual connection, previously unused in
 214 H-Nets. Residual connections are also used within our encoder blocks. As in standard CNNs, they
 215 improve gradient flow and reduce training time. With respect to rotation orders, they process streams
 216 independently, thus preserving HE according to the following lemma:

217 **Lemma 5.1** (HE of Residual Connections). A residual connection between feature maps of the same
 218 rotation order, $F_m(I)$ and $F'_m(I)$, preserves HE property:

$$F'_m([I]^\alpha) + F_m([I]^\alpha) = e^{im\alpha} [(F'_m(I) + F_m(I))]^\alpha. \quad (10)$$

219 5.2 Harmformer: S2 Construction of the Patches

220 To integrate the stem output with the encoder, the final stem feature maps are divided into 1×1 -sized
 221 patches, as illustrated in Figure 4a. The patches are constructed separately for all three streams of
 222 rotation orders. The resulting stack of patches then comprises three matrices $F_{-1}, F_0, F_1 \in \mathbb{C}^{(h \cdot w) \times d}$,
 223 each representing a single rotation order, where h , w , and d denote the height, width, and number
 224 of channels of the last feature maps, respectively. We keep this notation for encoder feature maps
 225 (patches), as they correspond to the stem feature maps, just reshuffled.

226 Neglecting small interpolation errors, the spatial transformation of the input translates only into
 227 a permutation of the stack of patches F_m as discussed in Sec. 3. For clarity, we use a discrete
 228 representation but it should be noted that the encoder can be modeled using a functional framework,
 229 as shown by Romero and Cordonnier [21].

230 Before feeding the SA with patches, transformer networks typically apply a linear projection to adjust
 231 the dimension d . We use a linear layer that processes the patches independently with respect to their
 232 order of rotation to preserve HE:

233 **Lemma 5.2** (HE of Linear Layer). A linear layer applied to a HE feature map $F_m([I]^\alpha) \in \mathbb{C}^{(hw) \times d_{in}}$
 234 preserves the rotation order m . Formally, we have:

$$F_m([I]^\alpha)W = e^{mi\alpha} [F_m(I)W]^\alpha, \quad (11)$$

235 where $W \in \mathbb{C}^{d_{in} \times d_{out}}$ represents a shared weight matrix applied independently over all spatial
 236 positions of the input feature map.

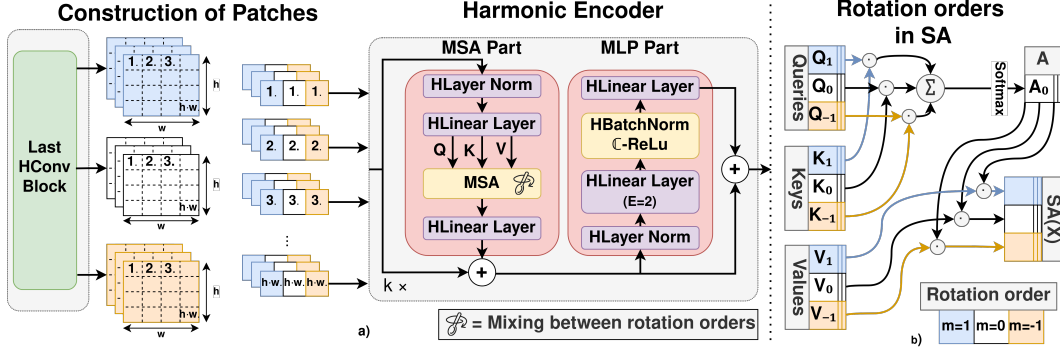


Figure 4: a) Construction of the patches (colors represent rotation orders) and the Harmonic Encoder structure. b) Diagram depicting the interaction of SA mechanisms across different rotation orders.

237 5.3 Harmformer: S3 Harmonic Encoder

238 This section outlines our encoder, which is designed to preserve the HE property. The encoder is
 239 organized into several k blocks, each containing Multi-Head Self-Attention (MSA) and Multi-Layer
 240 Perceptron (MLP) components, as shown in Figure 4a. Along with the layers presented in the previous
 241 sections, we propose a SA mechanism and a layer normalization, both of which preserve the HE.

242 As the following lemma shows, the layer normalization can be adapted to satisfy HE by operating
 243 independently on the streams of rotation orders.

244 **Lemma 5.3** (HE of Layer Norm). A feature map $F_m([I]^\alpha) \in \mathbb{C}^{(hw) \times d}$ with a rotation order m
 245 preserves HE when normalized by its mean and standard deviation:

$$\frac{F_m([I]^\alpha) - \mu}{\sigma + \epsilon} = e^{im\alpha} \frac{[F_m(I)]^\alpha - \mu}{\sigma + \epsilon}, \quad (12)$$

246 where μ, σ are the sample means and standard deviations of the original feature maps computed over
 247 their spatial dimensions, respectively, and ϵ is a small constant added for numerical stability.

248 **Self-Attention** The essential components of the encoder are MSA layers. The proposed MSA mixes
 249 features with different rotation orders. In the first step, queries, keys and values are generated for
 250 each rotation order $-1, 0, \text{ and } 1$ independently, which preserves HE as follows from Lemma 5.2. We
 251 split the SA calculation (Eq. (3),(4)) into two operations: dot product and matrix multiplication, and
 252 demonstrate their properties by the following lemmas.

253 **Lemma 5.4** (Dot product subtracts rotation orders). Consider two HE feature maps $Q_{m_1}([I]^\alpha) \in$
 254 $\mathbb{C}^{(hw) \times d}$ and $K_{m_2}([I]^\alpha) \in \mathbb{C}^{(hw) \times d}$ that represent queries and keys, respectively. The dot product
 255 of these feature maps is HE and has the rotation order $m_1 - m_2$. Formally, we have:

$$Q_{m_1}([I]^\alpha) \overline{K_{m_2}([I]^\alpha)^T} = e^{i(m_1 - m_2)\alpha} [Q_{m_1}(I) \overline{K_{m_2}(I)^T}]^\alpha, \quad (13)$$

256 where $\overline{K_{m_2}([I]^\alpha)^T}$ denotes the complex conjugate transpose of $K_{m_2}([I]^\alpha)$.

257 **Lemma 5.5** (Matrix multiplication sums rotation orders). Consider a HE feature map $A_{m_1}([I]^\alpha) \in$
 258 $\mathbb{C}^{(hw) \times (hw)}$ representing an attention matrix and HE feature map $V_{m_2}([I]^\alpha) \in \mathbb{C}^{(hw) \times d}$ representing
 259 values. The result of their matrix multiplication is HE with a rotation order $m = m_1 + m_2$:

$$A_{m_1}([I]^\alpha) V_{m_2}([I]^\alpha) = e^{i(m_1 + m_2)\alpha} [A_{m_1}(I) V_{m_2}(I)]^\alpha. \quad (14)$$

260 where $[A_{m_1}(I)]^\alpha$ and $[V_{m_2}(I)]^\alpha$ are feature maps created from unrotated I and rotated afterwards.

261 After these operations, the relative circular encodings [48] are added to the result of the dot product
 262 before it undergoes the softmax activation. The Harmformer softmax operates only on magnitudes
 263 and its codomain is within \mathbb{R}_0^+ to avoid breaking HE.

264 We have shown that the dot product between queries Q_{m_1} and keys K_{m_2} results in a rotation order
 265 of $m = m_1 - m_2$. Similarly, matrix multiplication between the attention matrix A_{m_1} and values
 266 V_{m_2} yields a feature map with a rotation order of $m = m_1 + m_2$. The final task is to combine these
 267 rotation orders to produce output feature maps with the same number of rotation orders as the input
 268 feature maps, i.e. -1, 0, and 1.

269 **Mixing Orders in MSA** Since there are multiple strategies for combining rotation orders, we
 270 explored several of them and provide details on other configurations in Ablation B.2. The optimal
 271 approach, according to our experiments, is shown in Figure 4 and involves:

- 272 1. **Dot Product Calculation:** The dot product is computed only between the same rotation
 273 orders, separately. According to Lemma 5.4, this results in three feature maps with the
 274 rotation order 0 and dimension $\mathbb{C}^{(hw) \times (hw)}$.
- 275 2. **Attention Matrix Formation:** These results are summed to form a single matrix of
 276 rotation order 0. A softmax function is then applied to form a single attention matrix
 277 $A_0 \in \mathbb{C}^{(hw) \times (hw)}$, preserving the rotation order 0, because softmax function operates only
 278 on magnitudes and outputs non-negative numbers.
- 279 3. **Self-Attention Output:** Finally, the self-attention output is produced by matrix multiplica-
 280 tion of the attention matrix (rotation order zero) with the values of each rotation order. This
 281 process results in a triplet of outputs with the target rotation orders $(-1, 0, 1)$.

282 Other layers, such as linear layers and residual connections, have been introduced in previous sections.
 283 By stacking HE layers, the last feature map coming from our encoder will maintain the HE property.

284 5.4 Harmformer: S4 Classification

285 Spatial position and orientation are generally redundant for classification tasks, except when classify-
 286 ing directional objects such as arrows. In the final stage, we remove this redundant information from
 287 the feature maps and produce an invariant feature vector. The feature maps entering the classification
 288 stage form a matrix of the shape $\mathbb{C}^{3 \times n \times d}$. To aggregate over different rotation orders, we keep only
 289 the magnitude, resulting in $\mathbb{R}^{n \times 3d}$. The spatial information is then eliminated by applying global
 290 average pooling over the dimension n (patches), reducing the shape to \mathbb{R}^{3d} . The final feature vector,
 291 which is roto-translation invariant, is processed by a single linear layer for classification.

292 6 Experiments

293 To validate the properties of Harmformer, we conducted experiments on four benchmarks listed in
 294 Table 1. For detailed experimental configurations and an ablation study of architectural modifications,
 295 see Appendices C and B, respectively. Addition segmentation experiment is included in Appendix D.

Table 1: Dataset overview, detailing sizes and whether training and test sets contain rotated images.

Dataset Name	Sample Size	Train/Test/Val. Size	Rot. Train/Test	Ref.	Scenario
mnist-rot-test	$28 \times 28 \times 1$	50k / 10k / 10k	X/✓	[27]	1
cifar-rot-test	$32 \times 32 \times 3$	42k / 10k / 8k	X/✓	[27]	1
rotated MNIST	$28 \times 28 \times 1$	10k / 2k / 50k	✓/✓	[24]	2
PCam	$96 \times 96 \times 3$	262k / 32k / 32k	X/X	[25, 26]	2

296 **Model Architecture** The models are designed to match the number of parameters of the previous
 297 state-of-the-art models while maintaining the same overall architecture. Depending on the benchmark,
 298 the stem stage consists of 2-4 blocks to reduce resolution, followed by 3-4 harmonic encoder
 299 blocks. To ensure that the equivariant properties emerge from the architecture, we avoid any data
 300 augmentation. Consistent with H-NeXt [27], the inputs are initially upscaled by a factor of two to
 301 mitigate interpolation errors.

Table 2: Error on mnist-rot-test

Model	Error	Param.
ResNets-50 [28]	57.6%	
SWN-GCN [28]	8.20%	2.7M
H-Nets [3]	7.11%	33k
H-NeXt [27]	1.30%	28k
Harmformer	0.82%	29.7k

Table 3: Error on cifar10-rot-test

Model	Error	Param.
ResNets-50 [28]	63.90%	
SWN-GCN [28]	49.50%	2.7M
H-NeXt[27]	38.54%	118k
Harmformer	31.41%	118k
Harmformer (Large)	29.29%	217k

Table 4: Error on rotated MNIST

SA model	Error	Param.
GSA-Nets[21]	2.03%	44.7k
GE-ViT[22]	1.99%	45k
Harmformer	1.18%	30k
CNNs model		
G-CNNs[2]	1.69%	73.1k

Table 5: PCam

Error	Param.
15.24%	206k
16.18%	
12.47%	146k
10.88%	141k

Table 6: Average Error

Dataset	Avg %±Std
rotMNIST	1.26±0.055
PCam	14.2±0.986
mnist-rot-test	0.91±0.142
cifar-rot-test	31.9±0.636
cifar-rot-test (Large)	29.7±0.203

302 **Invariance Benchmarks** In the first scenario, we verify the equivariance of Harmformer by training
 303 the model exclusively on upright (non-rotated) data and testing it on randomly rotated data; the first
 304 two datasets in Tab. 1. Since the model is trained only on upright images, any equivariant properties
 305 must arise purely from the model design, not from the training data.

306 We outperform previous methods on both datasets, as shown in Tables 2 and 3. Harmformer
 307 improves the robustness to rotation, as we discuss further in Section B.4, which partially explains
 308 the performance gain on mnist-rot-test. Stability under rotation is less enhanced on cifar-rot-test
 309 (Sec. B.4), and the superior results there are likely due to the higher model capacity of the transformer
 310 architecture, which improves the overall detection performance.

311 **Equivariance Benchmarks** In the second scenario, we compare the performance of the Harmformer
 312 on established equivariance benchmarks for roto-translation where there is no significant distribution
 313 shift between the training and test sets, either containing rotated samples or not. This evaluation
 314 assesses how our method stacks up against previous equivariant transformers that are equivariant to
 315 discrete rotation and translation. Tables 4 and 5 show the top results on the rotated MNIST and PCam
 316 datasets, respectively. Harmformer outperforms previous equivariant transformers and narrows the
 317 performance gap between equivariant transformers and convolution-based models.

318 For completeness, we include the average performance in each benchmark listed in Table 6. The
 319 accuracy on the PCam dataset was slightly unstable, probably due to the characteristics of the dataset.
 320 Methods specifically designed for PCam, such as [8], use extensive augmentation techniques, which
 321 were avoided in our case to ensure unbiased results.

322 7 Conclusion and Future Work

323 The proposed Harmformer is the first transformer model to achieve end-to-end equivariance to
 324 continuous rotation and translation in 2D. This was accomplished by designing an equivariant self-
 325 attention inspired by harmonic convolution. Along with the novel SA, we introduced several layers
 326 specifically tailored for equivariance, including linear layers, layer normalization, batch normalization,
 327 activations, and residual connections. Our model outperforms previous equivariant transformers,
 328 narrowing the performance gap with convolution-based equivariant networks.

329 We hypothesize that the full potential of transformers may not be realized due to the nature of
 330 traditional benchmarks. The 2D equivariant transformers have so far been tested on datasets containing
 331 relatively small images that lack global dependencies. Therefore, future research should explore
 332 the application of equivariant transformers on larger datasets where, similar to ViT, they could
 333 demonstrate their potential. In addition, the proposed model can be extended to other modalities
 334 while maintaining its equivariance properties. For example, the harmonic networks that form the
 335 basis of our approach can also be adapted for 3D applications.

336 **References**

- 337 [1] Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to
 338 document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998. doi: 10.1109/5.726791. URL
 339 <https://ieeexplore.ieee.org/abstract/document/726791>.
- 340 [2] Taco Cohen and Max Welling. Group Equivariant Convolutional Networks. In *Proceedings of The 33rd*
 341 *International Conference on Machine Learning*, pages 2990–2999. PMLR, 2016. doi: 10.48550/arXiv.
 342 1602.07576. URL <https://arxiv.org/pdf/1602.07576>.
- 343 [3] Daniel E Worrall, Stephan J Garbin, Daniyar Turmukhambetov, and Gabriel J Brostow. Harmonic
 344 Networks: Deep Translation and Rotation Equivariance. In *Proceedings of the IEEE Conference on*
 345 *Computer Vision and Pattern Recognition (CVPR)*, pages 5028–5037, 2017. doi: 10.48550/arXiv.1612.
 346 04642. URL [https://openaccess.thecvf.com/content_cvpr_2017/html/Worrall_Harmonic_](https://openaccess.thecvf.com/content_cvpr_2017/html/Worrall_Harmonic_Networks_Deep_CVPR_2017_paper.html)
 347 [Networks_Deep_CVPR_2017_paper.html](https://openaccess.thecvf.com/content_cvpr_2017/html/Worrall_Harmonic_Networks_Deep_CVPR_2017_paper.html).
- 348 [4] Jiaming Han, Jian Ding, Nan Xue, and Gui-Song Xia. Redet: A rotation-equivariant de-
 349 tector for aerial object detection. In *Proceedings of the IEEE/CVF Conference on Com-*
 350 *puter Vision and Pattern Recognition (CVPR)*, pages 2786–2795, June 2021. URL [https:](https://openaccess.thecvf.com/content/CVPR2021/html/Han_ReDet_A_Rotation-Equivariant_Detector_for_Aerial_Object_Detection_CVPR_2021_paper.html)
 351 [//openaccess.thecvf.com/content/CVPR2021/html/Han_ReDet_A_Rotation-Equivariant_](https://openaccess.thecvf.com/content/CVPR2021/html/Han_ReDet_A_Rotation-Equivariant_Detector_for_Aerial_Object_Detection_CVPR_2021_paper.html)
 352 [Detector_for_Aerial_Object_Detection_CVPR_2021_paper.html](https://openaccess.thecvf.com/content/CVPR2021/html/Han_ReDet_A_Rotation-Equivariant_Detector_for_Aerial_Object_Detection_CVPR_2021_paper.html).
- 353 [5] Jian Ding, Nan Xue, Yang Long, Gui-Song Xia, and Qikai Lu. Learning RoI transformer
 354 for oriented object detection in aerial images. In *Proceedings of the IEEE/CVF Con-*
 355 *ference on Computer Vision and Pattern Recognition (CVPR)*, June 2019. URL [https:](https://openaccess.thecvf.com/content_CVPR_2019/html/Ding_Learning_RoI_Transformer_for_Oriented_Object_Detection_in_Aerial_Images_CVPR_2019_paper.html)
 356 [//openaccess.thecvf.com/content_CVPR_2019/html/Ding_Learning_RoI_Transformer_](https://openaccess.thecvf.com/content_CVPR_2019/html/Ding_Learning_RoI_Transformer_for_Oriented_Object_Detection_in_Aerial_Images_CVPR_2019_paper.html)
 357 [for_Oriented_Object_Detection_in_Aerial_Images_CVPR_2019_paper.html](https://openaccess.thecvf.com/content_CVPR_2019/html/Ding_Learning_RoI_Transformer_for_Oriented_Object_Detection_in_Aerial_Images_CVPR_2019_paper.html).
- 358 [6] Benjamin Chidester, Tianming Zhou, Minh N Do, and Jian Ma. Rotation equivariant and invariant neural
 359 networks for microscopy image analysis. *Bioinformatics*, 35(14):i530–i537, 07 2019. ISSN 1367-4803.
 360 doi: 10.1093/bioinformatics/btz353. URL <https://doi.org/10.1093/bioinformatics/btz353>.
- 361 [7] Benjamin Chidester, That-Vinh Ton, Minh-Triet Tran, Jian Ma, and Minh N. Do. Enhanced
 362 rotation-equivariant U-Net for nuclear segmentation. In *Proceedings of the IEEE/CVF Con-*
 363 *ference on Computer Vision and Pattern Recognition (CVPR) Workshops*, June 2019. URL
 364 [https://openaccess.thecvf.com/content_CVPRW_2019/html/CVMI/Chidester_Enhanced_](https://openaccess.thecvf.com/content_CVPRW_2019/html/CVMI/Chidester_Enhanced_Rotation-Equivariant_U-Net_for_Nuclear_Segmentation_CVPRW_2019_paper.html)
 365 [Rotation-Equivariant_U-Net_for_Nuclear_Segmentation_CVPRW_2019_paper.html](https://openaccess.thecvf.com/content_CVPRW_2019/html/CVMI/Chidester_Enhanced_Rotation-Equivariant_U-Net_for_Nuclear_Segmentation_CVPRW_2019_paper.html).
- 366 [8] Simon Graham, David Epstein, and Nasir Rajpoot. Dense steerable filter cnns for exploiting rotational
 367 symmetry in histology images. *IEEE Transactions on Medical Imaging*, 39(12):4124–4136, 2020. doi:
 368 10.1109/TMI.2020.3013246. URL <https://ieeexplore.ieee.org/abstract/document/9153847>.
- 369 [9] Gong Cheng, Junwei Han, Peicheng Zhou, and Dong Xu. Learning rotation-invariant and fisher discrimi-
 370 native convolutional neural networks for object detection. *IEEE Transactions on Image Processing*, 28(1):
 371 265–278, 2019. doi: 10.1109/TIP.2018.2867198. URL [https://ieeexplore.ieee.org/abstract/](https://ieeexplore.ieee.org/abstract/document/8445665/)
 372 [document/8445665/](https://ieeexplore.ieee.org/abstract/document/8445665/).
- 373 [10] Fabian Fuchs, Daniel Worrall, Volker Fischer, and Max Welling. Se(3)-transformers: 3d roto-translation
 374 equivariant attention networks. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin,
 375 editors, *Advances in Neural Information Processing Systems*, volume 33, pages 1970–1981. Curran As-
 376 sociates, Inc., 2020. URL [https://proceedings.neurips.cc/paper_files/paper/2020/file/](https://proceedings.neurips.cc/paper_files/paper/2020/file/15231a7ce4ba789d13b722cc5c955834-Paper.pdf)
 377 [15231a7ce4ba789d13b722cc5c955834-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2020/file/15231a7ce4ba789d13b722cc5c955834-Paper.pdf).
- 378 [11] David W. Romero and Mark Hoogendoorn. Co-attentive equivariant neural networks: Focusing equivari-
 379 ance on transformations co-occurring in data. In *International Conference on Learning Representations*,
 380 2020. URL <https://openreview.net/forum?id=r1g6ogrtDr>.
- 381 [12] Michael J Hutchinson, Charline Le Lan, Sheheryar Zaidi, Emilien Dupont, Yee Whye Teh, and Hyunjik
 382 Kim. Lietransformer: Equivariant self-attention for lie groups. In Marina Meila and Tong Zhang, editors,
 383 *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of*
 384 *Machine Learning Research*, pages 4533–4543. PMLR, 18–24 Jul 2021. URL [https://proceedings.](https://proceedings.mlr.press/v139/hutchinson21a.html)
 385 [mlr.press/v139/hutchinson21a.html](https://proceedings.mlr.press/v139/hutchinson21a.html).
- 386 [13] Iliia Igashov, Hannes Stärk, Clément Vignac, Arne Schneuing, Victor Garcia Satorras, Pascal Frossard,
 387 Max Welling, Michael Bronstein, and Bruno Correia. Equivariant 3d-conditional diffusion model for
 388 molecular linker design. *Nature Machine Intelligence*, 6(4):417–427, Apr 2024. ISSN 2522-5839. doi:
 389 10.1038/s42256-024-00815-9. URL <https://doi.org/10.1038/s42256-024-00815-9>.

- 390 [14] Yi-Lun Liao, Brandon M Wood, Abhishek Das, and Tess Smidt. EquiformerV2: Improved equivariant
391 transformer for scaling to higher-degree representations. In *The Twelfth International Conference on*
392 *Learning Representations*, 2024. URL <https://openreview.net/forum?id=mCOBKZmrzD>.
- 393 [15] Yi-Lun Liao and Tess Smidt. Equiformer: Equivariant graph attention transformer for 3d atomistic
394 graphs. In *The Eleventh International Conference on Learning Representations*, 2023. URL <https://openreview.net/forum?id=KwmPfARgOTD>.
395
- 396 [16] Philipp Thölke and Gianni De Fabritiis. Equivariant transformers for neural network based molecular
397 potentials. In *International Conference on Learning Representations*, 2022. URL <https://openreview.net/forum?id=zNHqZ9wrRB>.
398
- 399 [17] Serge Assaad, Carlton Downey, Rami Al-Rfou', Nigamaa Nayakanti, and Benjamin Sapp. VN-transformer:
400 Rotation-equivariant attention for vector neurons. *Transactions on Machine Learning Research*, 2023.
401 ISSN 2835-8856. URL <https://openreview.net/forum?id=EiX2L4sDPG>.
- 402 [18] Lingshen He, Yiming Dong, Yisen Wang, Dacheng Tao, and Zhouchen Lin. Gauge equivariant trans-
403 former. In M. Ranzato, A. Beygelzimer, Y. Dauphin, P.S. Liang, and J. Wortman Vaughan, editors,
404 *Advances in Neural Information Processing Systems*, volume 34, pages 27331–27343. Curran Asso-
405 ciates, Inc., 2021. URL [https://proceedings.neurips.cc/paper_files/paper/2021/file/](https://proceedings.neurips.cc/paper_files/paper/2021/file/e57c6b956a6521b28495f2886ca0977a-Paper.pdf)
406 [e57c6b956a6521b28495f2886ca0977a-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2021/file/e57c6b956a6521b28495f2886ca0977a-Paper.pdf).
- 407 [19] Pim de Haan, Taco Cohen, and Johann Brehmer. Euclidean, projective, conformal: Choosing a geometric
408 algebra for equivariant transformers. In Sanjoy Dasgupta, Stephan Mandt, and Yingzhen Li, editors,
409 *Proceedings of The 27th International Conference on Artificial Intelligence and Statistics*, volume 238
410 of *Proceedings of Machine Learning Research*, pages 3088–3096. PMLR, 02–04 May 2024. URL
411 <https://proceedings.mlr.press/v238/haan24a.html>.
- 412 [20] Johann Brehmer, Pim de Haan, Sönke Behrends, and Taco S Cohen. Geometric algebra trans-
413 former. In A. Oh, T. Neumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine, editors, *Ad-
414 vances in Neural Information Processing Systems*, volume 36, pages 35472–35496. Curran Asso-
415 ciates, Inc., 2023. URL [https://proceedings.neurips.cc/paper_files/paper/2023/file/](https://proceedings.neurips.cc/paper_files/paper/2023/file/6f6dd92b03ff9be7468a6104611c9187-Paper-Conference.pdf)
416 [6f6dd92b03ff9be7468a6104611c9187-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2023/file/6f6dd92b03ff9be7468a6104611c9187-Paper-Conference.pdf).
- 417 [21] David W. Romero and Jean-Baptiste Cordonnier. Group Equivariant Stand-Alone Self-Attention For
418 Vision. In *International Conference on Learning Representations*, 2021. doi: 10.48550/arXiv.2010.00977.
419 URL <https://openreview.net/forum?id=Jkfyjn0Eo6M>.
- 420 [22] Renjun Xu, Kaifan Yang, Ke Liu, and Fengxiang He. $E(2)$ -Equivariant Vision Transformer. In Robin J.
421 Evans and Ilya Shpitser, editors, *Proceedings of the Thirty-Ninth Conference on Uncertainty in Artificial*
422 *Intelligence*, volume 216 of *Proceedings of Machine Learning Research*, pages 2356–2366. PMLR, 31
423 Jul–04 Aug 2023. doi: 10.48550/arXiv.2306.06722. URL [https://proceedings.mlr.press/v216/](https://proceedings.mlr.press/v216/xu23b.html)
424 [xu23b.html](https://proceedings.mlr.press/v216/xu23b.html).
- 425 [23] William T Freeman et al. The design and use of steerable filters. *IEEE Transactions on Pattern analysis*
426 *and machine intelligence*, 13(9):891–906, 1991.
- 427 [24] Hugo Larochelle, Dumitru Erhan, Aaron Courville, James Bergstra, and Yoshua Bengio. An empirical
428 evaluation of deep architectures on problems with many factors of variation. In *Proceedings of the 24th*
429 *International Conference on Machine Learning*, pages 473–480. Association for Computing Machinery,
430 2007. ISBN 978-1-59593-793-3. doi: 10.1145/1273496.1273556. URL [https://dl.acm.org/doi/](https://dl.acm.org/doi/abs/10.1145/1273496.1273556)
431 [abs/10.1145/1273496.1273556](https://dl.acm.org/doi/abs/10.1145/1273496.1273556).
- 432 [25] Babak Ehteshami Bejnordi, Mitko Veta, Paul Johannes van Diest, Bram van Ginneken, Nico Karssemeijer,
433 Geert Litjens, Jeroen A. W. M. van der Laak, , and the CAMELYON16 Consortium. Diagnostic Assessment
434 of Deep Learning Algorithms for Detection of Lymph Node Metastases in Women With Breast Cancer.
435 *JAMA*, 318(22):2199–2210, 12 2017. ISSN 0098-7484. doi: 10.1001/jama.2017.14585. URL <https://doi.org/10.1001/jama.2017.14585>.
436 [//doi.org/10.1001/jama.2017.14585](https://doi.org/10.1001/jama.2017.14585).
- 437 [26] Bastiaan S. Veeling, Jasper Linmans, Jim Winkens, Taco Cohen, and Max Welling. Rotation equivariant
438 cnns for digital pathology. In *Medical Image Computing and Computer Assisted Intervention – MICCAI*
439 *2018: 21st International Conference, Granada, Spain, September 16-20, 2018, Proceedings, Part II*,
440 page 210–218, Berlin, Heidelberg, 2018. Springer-Verlag. ISBN 978-3-030-00933-5. doi: 10.1007/
441 [978-3-030-00934-2_24](https://doi.org/10.1007/978-3-030-00934-2_24). URL https://doi.org/10.1007/978-3-030-00934-2_24.
- 442 [27] Tomáš Karella, Filip Šroubek, Jan Blažek, Jan Flusser, and Václav Košík. H-NeXt: The next step
443 towards roto-translation invariant networks. In *34th British Machine Vision Conference 2023, BMVC 2023,*
444 *Aberdeen, UK, November 20-24, 2023*. BMVA, 2023. URL [https://papers.bmvc2023.org/0578.](https://papers.bmvc2023.org/0578.pdf)
445 [pdf](https://papers.bmvc2023.org/0578.pdf).

- 446 [28] Sungwon Hwang, Hyungtae Lim, and Hyun Myung. Equivariance-bridged SO (2)-invariant representation
447 learning using graph convolutional network. In *The 32nd British Machine Vision Conference (BMVC*
448 *2021)*. The British Machine Vision Association, 2021. doi: 10.48550/arXiv.2106.09996. URL <https://www.bmvc2021-virtualconference.com/assets/papers/0218.pdf>.
449
- 450 [29] Renata Khasanova and Pascal Frossard. Graph-based isometry invariant representation learning. In *Proceed-*
451 *ings of the 34th International Conference on Machine Learning*, volume 70, pages 1847–1856. PMLR, 2017.
452 doi: 10.48550/arXiv.1703.00356. URL <http://proceedings.mlr.press/v70/khasanova17a.html>.
- 453 [30] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz
454 Kaiser, and Illia Polosukhin. Attention is all you need. In I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach,
455 R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems*,
456 volume 30. Curran Associates, Inc., 2017. URL [https://proceedings.neurips.cc/paper_files/](https://proceedings.neurips.cc/paper_files/paper/2017/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf)
457 [paper/2017/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2017/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf).
- 458 [31] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas
459 Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit,
460 and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In
461 *International Conference on Learning Representations*, 2021. URL [https://openreview.net/forum?](https://openreview.net/forum?id=YicbFdNTTy)
462 [id=YicbFdNTTy](https://openreview.net/forum?id=YicbFdNTTy).
- 463 [32] Zihang Dai, Hanxiao Liu, Quoc V Le, and Mingxing Tan. Coatnet: Marrying convolution and attention
464 for all data sizes. In M. Ranzato, A. Beygelzimer, Y. Dauphin, P.S. Liang, and J. Wortman Vaughan,
465 editors, *Advances in Neural Information Processing Systems*, volume 34, pages 3965–3977. Curran As-
466 sociates, Inc., 2021. URL [https://proceedings.neurips.cc/paper_files/paper/2021/file/](https://proceedings.neurips.cc/paper_files/paper/2021/file/20568692db622456cc42a2e853ca21f8-Paper.pdf)
467 [20568692db622456cc42a2e853ca21f8-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2021/file/20568692db622456cc42a2e853ca21f8-Paper.pdf).
- 468 [33] Tete Xiao, Mannat Singh, Eric Mintun, Trevor Darrell, Piotr Dollar, and Ross Girshick. Early convolutions
469 help transformers see better. In M. Ranzato, A. Beygelzimer, Y. Dauphin, P.S. Liang, and J. Wortman
470 Vaughan, editors, *Advances in Neural Information Processing Systems*, volume 34, pages 30392–30400.
471 Curran Associates, Inc., 2021. URL [https://proceedings.neurips.cc/paper_files/paper/](https://proceedings.neurips.cc/paper_files/paper/2021/file/ff1418e8cc993fe8abcfe3ce2003e5c5-Paper.pdf)
472 [2021/file/ff1418e8cc993fe8abcfe3ce2003e5c5-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2021/file/ff1418e8cc993fe8abcfe3ce2003e5c5-Paper.pdf).
- 473 [34] Maurice Weiler, Fred A. Hamprecht, and Martin Storath. Learning steerable filters for rotation equiv-
474 ariant cnns. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*
475 *(CVPR)*, June 2018. URL [https://openaccess.thecvf.com/content_cvpr_2018/html/Weiler_](https://openaccess.thecvf.com/content_cvpr_2018/html/Weiler_Learning_Steerable_Filters_CVPR_2018_paper.html)
476 [Learning_Steerable_Filters_CVPR_2018_paper.html](https://openaccess.thecvf.com/content_cvpr_2018/html/Weiler_Learning_Steerable_Filters_CVPR_2018_paper.html).
- 477 [35] Daniel Worrall and Max Welling. Deep scale-spaces: Equivariance over scale. In H. Wallach, H. Larochelle,
478 A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing*
479 *Systems*, volume 32. Curran Associates, Inc., 2019. URL [https://proceedings.neurips.cc/paper_](https://proceedings.neurips.cc/paper_files/paper/2019/file/f04cd7399b2b0128970efb6d20b5c551-Paper.pdf)
480 [files/paper/2019/file/f04cd7399b2b0128970efb6d20b5c551-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2019/file/f04cd7399b2b0128970efb6d20b5c551-Paper.pdf).
- 481 [36] Ivan Sosnovik, Michał Szmaja, and Arnold Smeulders. Scale-equivariant steerable networks. In *Inter-*
482 *national Conference on Learning Representations*, 2020. URL [https://openreview.net/forum?id=](https://openreview.net/forum?id=HJgpugrKPS)
483 [HJgpugrKPS](https://openreview.net/forum?id=HJgpugrKPS).
- 484 [37] Md Ashiqur Rahman and Raymond A. Yeh. Truly scale-equivariant deep nets with fourier lay-
485 ers. In A. Oh, T. Neumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine, editors, *Ad-*
486 *vances in Neural Information Processing Systems*, volume 36, pages 6092–6104. Curran Asso-
487 ciates, Inc., 2023. URL [https://proceedings.neurips.cc/paper_files/paper/2023/file/](https://proceedings.neurips.cc/paper_files/paper/2023/file/1343edb2739a61a6e20bd8764e814b50-Paper-Conference.pdf)
488 [1343edb2739a61a6e20bd8764e814b50-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2023/file/1343edb2739a61a6e20bd8764e814b50-Paper-Conference.pdf).
- 489 [38] Maurice Weiler and Gabriele Cesa. General E(2)-equivariant steerable CNNs. In H. Wallach, H. Larochelle,
490 A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing*
491 *Systems*, volume 32. Curran Associates, Inc., 2019. URL [https://proceedings.neurips.cc/paper_](https://proceedings.neurips.cc/paper_files/paper/2019/file/45d6637b718d0f24a237069fe41b0db4-Paper.pdf)
492 [files/paper/2019/file/45d6637b718d0f24a237069fe41b0db4-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2019/file/45d6637b718d0f24a237069fe41b0db4-Paper.pdf).
- 493 [39] Maurice Weiler, Mario Geiger, Max Welling, Wouter Boomsma, and Taco S Cohen. 3D Steerable CNNs:
494 Learning rotationally equivariant features in volumetric data. In S. Bengio, H. Wallach, H. Larochelle,
495 K. Grauman, N. Cesa-Bianchi, and R. Garnett, editors, *Advances in Neural Information Processing Systems*,
496 volume 31. Curran Associates, Inc., 2018. URL [https://proceedings.neurips.cc/paper_files/](https://proceedings.neurips.cc/paper_files/paper/2018/file/488e4104520c6aab692863cc1dba45af-Paper.pdf)
497 [paper/2018/file/488e4104520c6aab692863cc1dba45af-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2018/file/488e4104520c6aab692863cc1dba45af-Paper.pdf).
- 498 [40] Daniel Worrall and Gabriel Brostow. CubeNet: Equivariance to 3D rotation and transla-
499 tion. In *Proceedings of the European Conference on Computer Vision (ECCV)*, September
500 2018. URL [https://openaccess.thecvf.com/content_ECCV_2018/html/Daniel_Worrall_](https://openaccess.thecvf.com/content_ECCV_2018/html/Daniel_Worrall_CubeNet_Equivariance_to_ECCV_2018_paper.html)
501 [CubeNet_Equivariance_to_ECCV_2018_paper.html](https://openaccess.thecvf.com/content_ECCV_2018/html/Daniel_Worrall_CubeNet_Equivariance_to_ECCV_2018_paper.html).

- 502 [41] Erik J Bekkers, Sharvaree Vadgama, Rob Hesselink, Putri A Van der Linden, and David W. Romero.
503 Fast, expressive $SE(n)$ equivariant networks through weight-sharing in position-orientation space. In *The*
504 *Twelfth International Conference on Learning Representations*, 2024. URL [https://openreview.net/](https://openreview.net/forum?id=dPHLbUqGbr)
505 [forum?id=dPHLbUqGbr](https://openreview.net/forum?id=dPHLbUqGbr).
- 506 [42] Taco S. Cohen, Mario Geiger, Jonas Köhler, and Max Welling. Spherical CNNs. In *International Confer-*
507 *ence on Learning Representations*, 2018. URL <https://openreview.net/forum?id=Hkbd5xZRb>.
- 508 [43] Víctor García Satorras, Emiel Hooeboom, and Max Welling. E(n) equivariant graph neural networks.
509 In Marina Meila and Tong Zhang, editors, *Proceedings of the 38th International Conference on Machine*
510 *Learning*, volume 139 of *Proceedings of Machine Learning Research*, pages 9323–9332. PMLR, 18–24
511 Jul 2021. URL <https://proceedings.mlr.press/v139/satorras21a.html>.
- 512 [44] Maurice Weiler, Patrick Forré, Erik Verlinde, and Max Welling. Coordinate independent convolutional
513 networks – isometry and gauge equivariant convolutions on Riemannian manifolds, 2021.
- 514 [45] Maurice Weiler, Patrick Forré, Erik Verlinde, and Max Welling. *Equivariant and Coordinate Inde-*
515 *pendent Convolutional Networks*. 2023. URL [https://maurice-weiler.gitlab.io/cnn_book/](https://maurice-weiler.gitlab.io/cnn_book/EquivariantAndCoordinateIndependentCNNs.pdf)
516 [EquivariantAndCoordinateIndependentCNNs.pdf](https://maurice-weiler.gitlab.io/cnn_book/EquivariantAndCoordinateIndependentCNNs.pdf).
- 517 [46] Yang Liu, Yao Zhang, Yixin Wang, Feng Hou, Jin Yuan, Jiang Tian, Yang Zhang, Zhongchao Shi, Jianping
518 Fan, and Zhiqiang He. A survey of visual transformers. *IEEE Transactions on Neural Networks and*
519 *Learning Systems*, pages 1–21, 2023. doi: 10.1109/TNNLS.2022.3227717. URL [https://ieeexplore.](https://ieeexplore.ieee.org/abstract/document/10088164)
520 [ieee.org/abstract/document/10088164](https://ieeexplore.ieee.org/abstract/document/10088164).
- 521 [47] Prajit Ramachandran, Niki Parmar, Ashish Vaswani, Irwan Bello, Anselm Levskaya, and Jon Shlens.
522 Stand-alone self-attention in vision models. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché-Buc,
523 E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 32. Cur-
524 ran Associates, Inc., 2019. URL [https://proceedings.neurips.cc/paper_files/paper/2019/](https://proceedings.neurips.cc/paper_files/paper/2019/file/3416a75f4cea9109507cacd8e2f2aefc-Paper.pdf)
525 [file/3416a75f4cea9109507cacd8e2f2aefc-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2019/file/3416a75f4cea9109507cacd8e2f2aefc-Paper.pdf).
- 526 [48] Kan Wu, Houwen Peng, Minghao Chen, Jianlong Fu, and Hongyang Chao. Rethinking and improving
527 relative position encoding for vision transformer. In *Proceedings of the IEEE/CVF International Confer-*
528 *ence on Computer Vision (ICCV)*, pages 10033–10041, 2021. doi: 10.48550/arXiv.2107.14222. URL
529 [https://openaccess.thecvf.com/content/ICCV2021/html/Wu_Rethinking_and_Improving_](https://openaccess.thecvf.com/content/ICCV2021/html/Wu_Rethinking_and_Improving_Relative_Position-Encoding_for_Vision_Transformer_ICCV_2021_paper.html)
530 [Relative_Position-Encoding_for_Vision_Transformer_ICCV_2021_paper.html](https://openaccess.thecvf.com/content/ICCV2021/html/Wu_Rethinking_and_Improving_Relative_Position-Encoding_for_Vision_Transformer_ICCV_2021_paper.html).
- 531 [49] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing
532 internal covariate shift. In Francis Bach and David Blei, editors, *Proceedings of the 32nd International*
533 *Conference on Machine Learning*, volume 37 of *Proceedings of Machine Learning Research*, pages 448–
534 456, Lille, France, 07–09 Jul 2015. PMLR. URL [https://proceedings.mlr.press/v37/ioffe15.](https://proceedings.mlr.press/v37/ioffe15.html)
535 [html](https://proceedings.mlr.press/v37/ioffe15.html).
- 536 [50] J. Staal, M.D. Abramoff, M. Niemeijer, M.A. Viergever, and B. van Ginneken. Ridge-based vessel
537 segmentation in color images of the retina. *IEEE Transactions on Medical Imaging*, 23(4):501–509, 2004.
538 doi: 10.1109/TMI.2004.825627.
- 539 [51] Erik J. Bekkers, Maxime W. Lafarge, Mitko Veta, Koen A. J. Eppenhof, Josien P. W. Pluim, and Remco
540 Duits. Roto-translation covariant convolutional networks for medical image analysis. In *Medical Image*
541 *Computing and Computer Assisted Intervention – MICCAI 2018: 21st International Conference, Granada,*
542 *Spain, September 16-20, 2018, Proceedings, Part I*, page 440–448, Berlin, Heidelberg, 2018. Springer-
543 Verlag. ISBN 978-3-030-00927-4. doi: 10.1007/978-3-030-00928-1_50. URL [https://doi.org/10.](https://doi.org/10.1007/978-3-030-00928-1_50)
544 [1007/978-3-030-00928-1_50](https://doi.org/10.1007/978-3-030-00928-1_50).
- 545 [52] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical
546 image segmentation. In *Medical image computing and computer-assisted intervention–MICCAI 2015: 18th*
547 *international conference, Munich, Germany, October 5-9, 2015, proceedings, part III 18*, pages 234–241.
548 Springer, 2015.
- 549 [53] Wentao Liu, Huihua Yang, Tong Tian, Zhiwei Cao, Xipeng Pan, Weijin Xu, Yang Jin, and Feng Gao. Full-
550 resolution network and dual-threshold iteration for retinal vessel and coronary angiograph segmentation.
551 *IEEE Journal of Biomedical and Health Informatics*, 26(9):4623–4634, 2022. doi: 10.1109/JBHI.2022.
552 3188710.
- 553 [54] Risi Kondor, Zhen Lin, and Shubhendu Trivedi. Clebsch–gordan nets: a fully fourier space spherical
554 convolutional neural network. *Advances in Neural Information Processing Systems*, 31, 2018.

555 A Proofs of Harmformer Equivariance

556 In this section, we systematically formulate the proofs of the HE property (Definition 4.2) for each
 557 layer of the Harmformer. The central concept of the architecture is the handling of three streams
 558 corresponding to different rotation orders. By oversimplifying the interactions of rotation orders, two
 559 main properties of the harmonic function should be highlighted:

- 560 1. Feature maps with the same rotation order can be summed:

$$e^{im\alpha} F_1 + e^{im\alpha} F_2 = e^{im\alpha} (F_1 + F_2)$$

- 561 2. Multiplication of feature maps results in the sum of their rotation orders:

$$e^{im_1\alpha} F_1 \cdot e^{im_2\alpha} F_2 = e^{i(m_1+m_2)\alpha} (F_1 \cdot F_2)$$

562 Interactions between these streams occur in layers harmonic convolution or Multi-Head Attention
 563 (MSA). Other layers, such as layer normalization, process feature maps of the different rotation order
 564 independently. Additionally, some operations, such as batch normalization and activation functions,
 565 operate solely on the magnitudes of complex numbers leaving the phase untouched.

566 A.1 Equivariance of Harmonic Convolutions

567 Note that the proof of H-Conv equivariance was originally formulated by Worrall et al. [3]. For the
 568 sake of completeness, we have provided a highly simplified version of these proofs. However, we
 569 encourage readers to read the more comprehensive work on G-steerable convolution kernels and the
 570 theory of steerable equivariant convolution networks in [45] (Chapters 4-5), which provides a broader
 571 perspective and demonstrates the equivalence with G-CNNs.

572 **Lemma A.1** (Rotation of a Harmonic Filter). When the coordinates of a harmonic filter are rotated
 573 by an angle α , it only changes by a factor $e^{im\alpha}$, where m is the rotation order of the harmonic filter
 574 and \mathcal{R}_α is a corresponding 2D rotation matrix.

Proof.

$$\begin{aligned} W_m(\mathcal{R}_\alpha^{-1}\mathbf{x}) &\equiv \tilde{W}_m(r, \theta - \alpha) = R(r) \cdot e^{-im(\theta - \alpha)} \\ &= e^{im\alpha} \tilde{W}_m(r, \theta) \equiv e^{im\alpha} W_m(\mathbf{x}), \end{aligned} \quad (15)$$

575 where \mathbf{x} is the spatial coordinates. □

576 Let us denote an input image I that is rotated by the angle α and translated by vector \mathbf{t} as

$$[I]_t^\alpha(\mathbf{x}) \equiv I(\mathcal{R}_\alpha^{-1}\mathbf{x} + \mathbf{t}). \quad (16)$$

577 **Theorem A.1** (Harmonic convolution sums the rotation orders). When an input image I is rotated by
 578 α and translated by \mathbf{t} , the output of a multiple successive harmonic convolution is given by:

$$[W_{m_1} \circledast W_{m_2} \circledast \dots \circledast [I]_t^\alpha](\mathbf{x}) = e^{i(m_1+m_2+\dots)\alpha} [W_{m_1} \circledast W_{m_2} \circledast \dots \circledast I](\mathcal{R}_\alpha\mathbf{x} + \mathbf{t}) \quad (17)$$

579 *Proof.* We start with the very first harmonic convolution.

$$[W_{m_1} \circledast [I]_t^\alpha](\mathbf{x}) = \int_{\mathbb{R}^2} W_{m_1}(\mathbf{z}) [I]_t^\alpha(\mathbf{x} - \mathbf{z}) d\mathbf{z} \quad (18)$$

$$= \int_{\mathbb{R}^2} W_{m_1}(\mathbf{z}) I(\mathcal{R}_\alpha(\mathbf{x} - \mathbf{z}) + \mathbf{t}) d\mathbf{z} \quad (Eq.16) \quad (19)$$

$$= \int_{\mathbb{R}^2} W_{m_1}(\mathcal{R}_\alpha^{-1}\mathbf{z}') I(\mathcal{R}_\alpha\mathbf{x} - \mathbf{z}' + \mathbf{t}) d\mathbf{z}' \quad (\mathbf{z}' = \mathcal{R}_\alpha\mathbf{z})^1 \quad (20)$$

$$= e^{im_1\alpha} \int_{\mathbb{R}^2} W_{m_1}(\mathbf{z}') I((\mathcal{R}_\alpha\mathbf{x} + \mathbf{t}) - \mathbf{z}') d\mathbf{z}' \quad (\text{Lemma A.1}) \quad (21)$$

$$= e^{im_1\alpha} [W_{m_1} \circledast I](\mathcal{R}_\alpha\mathbf{x} + \mathbf{t}) \quad (22)$$

580 Denote the first feature map as F , if we roto-translate the input image:

$$[F]_t^\alpha(\mathbf{x}) \equiv e^{im_1\alpha} F(\mathcal{R}_\alpha \mathbf{x} + \mathbf{t}) \quad (23)$$

581 The following harmonic convolution is given by a similar equation.

$$[W_{m_2} \otimes [F]_t^\alpha](\mathbf{x}) = \int_{\mathbb{R}^2} W_{m_2}(\mathbf{z}) [F]_t^\alpha(\mathbf{x} - \mathbf{z}) d\mathbf{z} \quad (24)$$

$$= e^{im_1\alpha} \int_{\mathbb{R}^2} W_{m_2}(\mathbf{z}) F(\mathcal{R}_\alpha \mathbf{x} - \mathcal{R}_\alpha \mathbf{z} + \mathbf{t}) d\mathbf{z} \quad (25)$$

$$= e^{i(m_1+m_2)\alpha} \int_{\mathbb{R}^2} W_{m_2}(\mathbf{z}') F((\mathcal{R}_\alpha \mathbf{x} + \mathbf{t}) - \mathbf{z}') d\mathbf{z}' \quad (\mathbf{z}' = \mathcal{R}_\alpha \mathbf{z}) \quad (26)$$

$$= e^{i(m_1+m_2)\alpha} [W_{m_2} \otimes F](\mathcal{R}_\alpha \mathbf{x} + \mathbf{t}) \quad (27)$$

582 Accordingly for all following harmonic convolution layers. \square

583 A.2 Layers Operating on Magnitudes

584 This section describes the original H-Nets layers that operate on magnitudes as formulated by Worrall
585 et al. [3], and introduces our proposed enhancements.

Definition A.1 (C-ReLU).

$$\mathbb{C}\text{-ReLU}_b(X e^{i\theta}) = \text{ReLU}(X + b) e^{i\theta}, \quad (28)$$

586 where $X e^{i\theta}$ represents a complex number in exponential form, and $b \in \mathbb{R}$ is a learnable bias parameter
587 of the activation function.

Definition A.2 (Harmformer C-ReLU).

$$\mathbb{C}\text{-ReLU}_{a,b}(X e^{i\theta}) = \text{ReLU}(a \cdot X + b) e^{i\theta}, \quad (29)$$

588 where $X e^{i\theta}$ is a complex number in exponential form, and $a, b \in \mathbb{R}$ are learnable parameters of the
589 activation function.

590 Definition A.1 uses only the bias parameter b . Such an activation function cannot zero out higher
591 values while leaving lower values unaffected. To allow this, our enhanced C-ReLU also incorporates
592 a multiplication by the parameter a .

593 A.3 Complex Batch Normalization in Harmonic Networks

594 In H-Nets, batch normalization is adapted from its traditional definition. The layer standardizes only
595 the magnitudes of the complex numbers, leaving the phase components unaffected. The C-BN can be
596 formally defined as follows:

Definition A.3 (C-BN).

$$\mathbb{C}\text{-BN}_{\gamma,\beta}(X e^{i\theta}) = \left(\gamma \left(\frac{X - \mu}{\sqrt{\sigma^2 + \epsilon}} \right) + \beta \right) e^{i\theta} = \text{BN}_{\gamma,\beta}(X) e^{i\theta}, \quad (30)$$

597 where $X e^{i\theta}$ represents a complex number in exponential form, and $\gamma, \beta \in \mathbb{R}$ are learnable scaling
598 and shifting parameters, respectively. Here, μ and σ denote the running sample mean and variance,
599 which are estimated during the training phase and fixed during inference.

600 However, this formulation can produce negative magnitudes, thus inverting the phase and violating
601 HE. Therefore in Harmformer we instead use a batch normalization integrated with an activation
602 function, that can be defined as:

Definition A.4 (Harmformer HBatchNorm + C-ReLU).

$$\mathbb{C}\text{-BN-ReLU}_{a,b}(X e^{i\theta}) = \text{ReLU} \left(a \left(\frac{X - \mu}{\sqrt{\sigma^2 + \epsilon}} \right) + b \right) e^{i\theta}, \quad (31)$$

603 where $X e^{i\theta}$ represents a complex number in exponential form, and $a, b \in \mathbb{R}$ are learnable scaling
604 and shifting parameters, respectively. Here, μ and σ denote the running sample mean and variance,
605 which are estimated during the training phase and fixed during inference.

606 Our formulation uses the ReLU function, which maps \mathbb{R} to \mathbb{R}^+ , to ensure that changes in magnitudes
607 are always positive. Additionally, by integrating the scaling and shifting parameters a and b into
608 batch normalization, the number of learnable parameters is reduced. See section B.1 for a comparison
609 of different normalization layers.

610 A.4 Discrete Representation

611 The normalization layers are the last from the stem stage, as can be seen in Figure 2. For the sake
612 of clarity, we will make the transition to discrete space and focus only on rotation for the following
613 layers, as mentioned in Section 5.2. Suppose that the feature maps (stack of patches) $F_m(I) \in \mathbb{C}^{n \times d}$,
614 extracted from the input image I , transforms under a rotation of the input as follows

$$F_m([I]^\alpha) = e^{im\alpha} [F_m(I)]^\alpha, \quad (32)$$

615 where m is the rotation order and α is the rotation angle of the image I . Here n is the number of
616 patches and d is the dimension of each patch.

617 This property implies that $[\cdot]^\alpha$ is a linear operator, thus for the feature maps the following applies:

$$[F_{m_1}(I)]^\alpha + [F_{m_2}(I)]^\alpha = [F_{m_1}(I) + F_{m_2}(I)]^\alpha \quad (33)$$

$$[F_{m_1}(I)]^\alpha \cdot [F_{m_2}(I)]^\alpha = [(F_{m_1}(I) \cdot F_{m_2}(I))]^\alpha \quad (34)$$

618 A.5 Residual Connection

619 **Lemma A.2** (HE of Residual Connections (Lemma 5.1)). A residual connection between feature
620 maps of the same rotation order, $F_m(I)$ and $F'_m(I)$, preserves HE property:

$$F'_m([I]^\alpha) + F_m([I]^\alpha) = e^{im\alpha} [(F'_m(I) + F_m(I))]^\alpha. \quad (35)$$

621 *Proof.* By the properties of harmonic equivariance, we have:

$$F'_m([I]^\alpha) + F_m([I]^\alpha) = (e^{im\alpha} [F'_m(I)]^\alpha) + (e^{im\alpha} [F_m(I)]^\alpha) \quad (36)$$

$$= e^{im\alpha} ([F'_m(I)]^\alpha + [F_m(I)]^\alpha) \quad (37)$$

622 Since $[\cdot]^\alpha$ is a linear operator, we can combine the rotated feature maps:

$$= e^{im\alpha} [(F'_m(I) + F_m(I))]^\alpha \quad (38)$$

623 \square

624 A.6 Linear Layers

625 **Lemma A.3** (HE of Linear Layer (Lemma 5.2)). A linear layer applied to a HE feature map
626 $F_m([I]^\alpha) \in \mathbb{C}^{(hw) \times d_{in}}$ preserves the rotation order m . Formally, we have:

$$F_m([I]^\alpha)W = e^{mi\alpha} [F_m(I)W]^\alpha, \quad (39)$$

627 where $W \in \mathbb{C}^{d_{in} \times d_{out}}$ represents a shared weight matrix applied independently over all spatial
628 positions of the input feature map.

629 *Proof.* As the matrix W has a rotation order $m = 0$ because it doesn't change under an input rotation.
630 Then the property comes trivially from Eq (34). \square

631 A.7 Multi-Head Self-Attention

632 **Lemma A.4** (HE of Layer Norm (Lemma 5.3)). A feature map $F_m([I]^\alpha) \in \mathbb{C}^{(hw) \times d}$ with a rotation
633 order m preserves HE when normalized by its mean and standard deviation:

$$\frac{F_m([I]^\alpha) - \mu}{\sigma + \epsilon} = e^{im\alpha} \frac{[F_m(I)]^\alpha - \mu}{\sigma + \epsilon}, \quad (40)$$

634 where μ, σ are the sample means and standard deviations of the original feature maps computed over
635 their spatial dimensions, respectively, and ϵ is a small constant added for numerical stability.

Proof.

$$\frac{F_m([I]^\alpha) - \hat{\mu}}{\hat{\sigma} + \epsilon} = \frac{F_m([I]^\alpha) - \frac{\sum F_m([I]^\alpha)}{h \cdot w}}{\hat{\sigma} + \epsilon} \quad (41)$$

636 By the properties of HE, we express $F_m([I]^\alpha) = e^{im\alpha} F_m(I)$. Sigma of does not change under
637 rotation, because its equal to standard deviation of magnitude in complex numbers.

$$= \frac{e^{im\alpha} [F_m(I)]^\alpha - \frac{\sum e^{im\alpha} [F_m(I)]^\alpha}{h \cdot w}}{\sigma + \epsilon} \quad (42)$$

$$= e^{im\alpha} \frac{[F_m(I)]^\alpha - \frac{\sum [F_m(I)]^\alpha}{h \cdot w}}{\sigma + \epsilon} \quad (43)$$

638 The sum of the feature map is invariant to its rotation.

$$= e^{im\alpha} \frac{[F_m(I)]^\alpha - \mu}{\sigma + \epsilon} \quad (44)$$

639

□

640 **Lemma A.5** (Dot product subtracts rotation orders (Lemma 5.4)). Consider two HE feature maps
641 $Q_{m_1}([I]^\alpha) \in \mathbb{C}^{(hw) \times d}$ and $K_{m_2}([I]^\alpha) \in \mathbb{C}^{(hw) \times d}$ that represent queries and keys, respectively. The
642 dot product of these feature maps is HE and has the rotation order $m_1 - m_2$. Formally, we have:

$$Q_{m_1}([I]^\alpha) \overline{K_{m_2}([I]^\alpha)^T} = e^{i(m_1 - m_2)\alpha} \left[Q_{m_1}(I) \overline{K_{m_2}(I)^T} \right]^\alpha, \quad (45)$$

643 where $\overline{K_{m_2}([I]^\alpha)^T}$ denotes the complex conjugate transpose of $K_{m_2}([I]^\alpha)$.

644 *Proof.* By the properties of harmonic equivariance, we express:

$$\begin{aligned} Q_{m_1}([I]^\alpha) &= e^{im_1\alpha} [Q_{m_1}(I)]^\alpha, \\ K_{m_2}([I]^\alpha) &= e^{im_2\alpha} [K_{m_2}(I)]^\alpha. \end{aligned}$$

645 Taking the complex conjugate transpose of $K_{m_2}([I]^\alpha)$, we obtain:

$$\overline{K_{m_2}([I]^\alpha)^T} = e^{-im_2\alpha} \overline{[K_{m_2}(I)^T]^\alpha} = e^{-im_2\alpha} \left[\overline{K_{m_2}(I)^T} \right]^\alpha.$$

646 As is derived from the commutativity of the scalar multiplication with the matrix multiplication.

$$\begin{aligned} Q_{m_1}([I]^\alpha) \overline{K_{m_2}([I]^\alpha)^T} &= e^{im_1\alpha} [Q_{m_1}(I)]^\alpha e^{-im_2\alpha} \left[\overline{K_{m_2}(I)^T} \right]^\alpha \\ &= e^{i(m_1 - m_2)\alpha} [Q_{m_1}(I)]^\alpha \left[\overline{K_{m_2}(I)^T} \right]^\alpha \\ &= e^{i(m_1 - m_2)\alpha} \left[Q_{m_1}(I) \overline{K_{m_2}(I)^T} \right]^\alpha. \end{aligned}$$

647 This shows that the dot product result is also HE with a rotation order of $m_1 - m_2$. □

648 **Lemma A.6** (Matrix multiplication sums rotation orders (Lemma 5.5)). Consider a HE feature
649 map $A_{m_1}([I]^\alpha) \in \mathbb{C}^{(hw) \times (hw)}$ representing an attention matrix, and HE feature map $V_{m_2}([I]^\alpha) \in$
650 $\mathbb{C}^{(hw) \times d}$ representing values. The result of their matrix multiplication is HE with a rotation order
651 $m = m_1 + m_2$:

$$A_{m_1}([I]^\alpha) V_{m_2}([I]^\alpha) = e^{i(m_1 + m_2)\alpha} [A_{m_1}(I) V_{m_2}(I)]^\alpha. \quad (46)$$

652 where $[A_{m_1}(I)]^\alpha, [V_{m_2}(I)]^\alpha$ are feature maps created from unrotated I and rotated afterwards.

653 *Proof.* Proof is analogical to Lemma 5.4 By the properties of harmonic equivariance, the HE feature
654 maps $A_{m_1}([I]^\alpha)$ and $V_{m_2}([I]^\alpha)$ can be represented as:

$$\begin{aligned} A_{m_1}([I]^\alpha) &= e^{im_1\alpha} [A_{m_1}(I)]^\alpha, \\ V_{m_2}([I]^\alpha) &= e^{im_2\alpha} [V_{m_2}(I)]^\alpha. \end{aligned}$$

655 Multiplying these matrices, we find:

$$\begin{aligned}
 A_{m_1}([I]^\alpha)V_{m_2}([I]^\alpha) &= (e^{im_1\alpha}[A_{m_1}(I)]^\alpha)(e^{im_2\alpha}[V_{m_2}(I)]^\alpha) \\
 &= e^{im_1\alpha}e^{im_2\alpha}[A_{m_1}(I)]^\alpha[V_{m_2}(I)]^\alpha \\
 &= e^{i(m_1+m_2)\alpha}[A_{m_1}(I)]^\alpha[V_{m_2}(I)]^\alpha. \\
 &= e^{i(m_1+m_2)\alpha}[A_{m_1}(I)V_{m_2}(I)]^\alpha.
 \end{aligned}$$

656 This confirms that the product is HE and preserves the combined rotation order of $m_1 + m_2$. \square

657 B Ablation Study and Additional Experiments

658 In addition to the experiments in the main text that compare the Harmformer to other methods, we
 659 include an ablation study that demonstrates its rotational robustness and explores other architectural
 660 choices. We have omitted the PCam benchmark due to computational constraints, as its training time
 661 was extremely long.

662 B.1 Ablation: Normalization Layers in Stem Stage (S1)

663 In Section A.4, we propose a modification of batch normalization by integrating it with an activation
 664 function. Specifically, we first apply batch normalization to the feature magnitudes, followed by
 665 a ReLU activation on these normalized values. This approach is more consistent with the original
 666 purpose of batch normalization as formulated by Ioffe and Szegedy [49], which is to standardize the
 667 distribution of activations across layers.

668 To test this novel normalization approach, we replaced our normalization layers in the Harmformer H-
 669 Conv block (see Figure 3b) with the original H-Nets batch normalization [3] followed by a \mathbb{C} -ReLU.
 670 We also evaluated how our proposed layer normalization used in the encoder block would perform in
 671 the H-Conv block.

672 The results depicted in Figure 5 show that our proposed normalization (blue bar) outperforms the
 673 original H-Nets normalization layer (red bar) across all three benchmarks, significantly reducing
 674 variance across different runs. It also exceeds the performance of layer normalization (yellow bar)
 675 in the rotated MNIST and mnist-rot-test, although it slightly underperforms in the cifar-rot-test. In
 676 addition, the layer normalization was more computationally expensive according to our experiments.

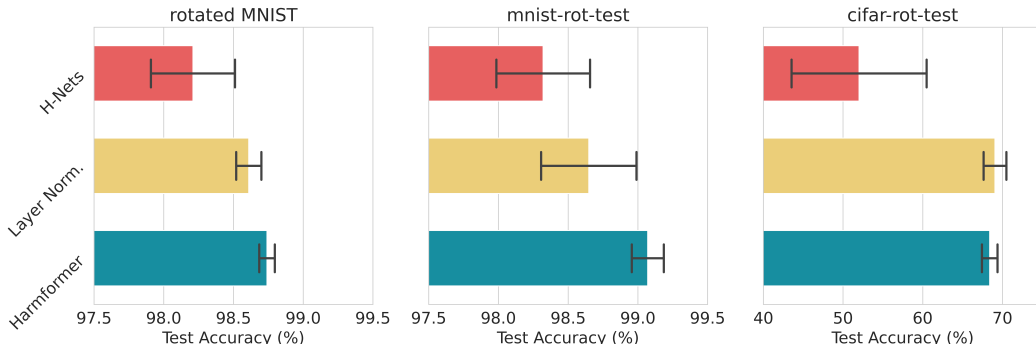


Figure 5: Ablation study on different normalization layers. The rows represent different normalization layers in the H-Conv block. Each plot is aggregated from 5 different runs. The error bars represent the standard deviation.

677 B.2 Ablation: Mixing Rotation Orders in Self-Attention Mechanism

678 As mentioned in the main text (Section 5), determining how queries, keys, and values should
 679 interact based on their rotation order is not intuitive. Therefore, we have extensively tested various
 680 configurations and listed the most promising ones in this section. In choosing the final solution, in
 681 addition to performance, we focused on the principles that the number of streams should not increase
 682 and that the method should not require extensive computation.

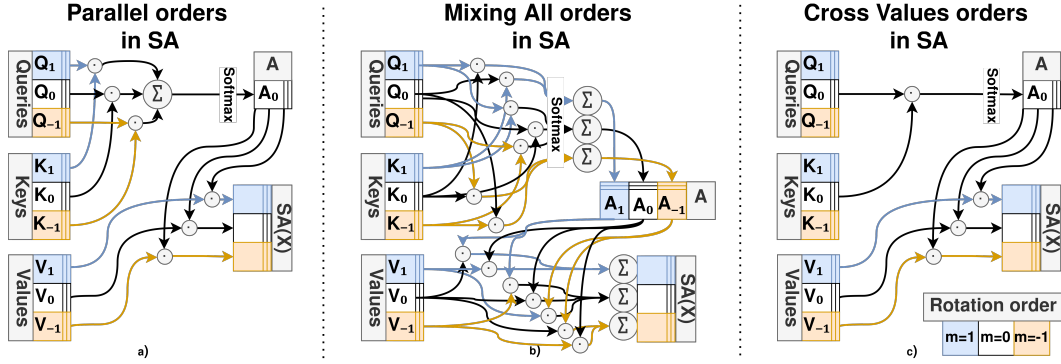


Figure 6: Ablation study on mixing rotation orders within the SA mechanism. a) The principle used in Harmformer b) Mixing all possible combinations of queries, keys and values c) Cross Values a method of mixing only values of different rotation orders.

683 In Lemmas 5.4 and 5.5, we demonstrate that the dot product subtracts the rotation orders and matrix
 684 multiplication sums the rotation orders. Based on this, we propose several configurations, illustrated
 685 in Figure 6. Apart from those mentioned here, we investigated learnable weights for each rotation
 686 order combination and different placements of softmax or combinations of these configurations
 687 together, but none yielded significant improvements. The final configuration used in Harmformer is
 688 shown in Figure 6a. The configuration in Figure 6b allows all possible combinations to produce the
 689 three streams (1, 0, -1). The last configuration, Cross Values, illustrated in Figure 6c, uses higher
 690 rotation orders only for values. Similarly, we tested Cross Keys and Cross Queries only for keys and
 691 queries, respectively.

692 Figure 7 presents the performance of these configurations on our benchmarks. The only configuration
 693 surpassing Harmformer (Figure 6a) was Mixing All (Figure 6b) in the case of rotated MNIST and
 694 mnist-rot-test. Since the performance difference was minor and the computational demands were
 695 significantly higher, we did not use Mixing All in our final architecture.

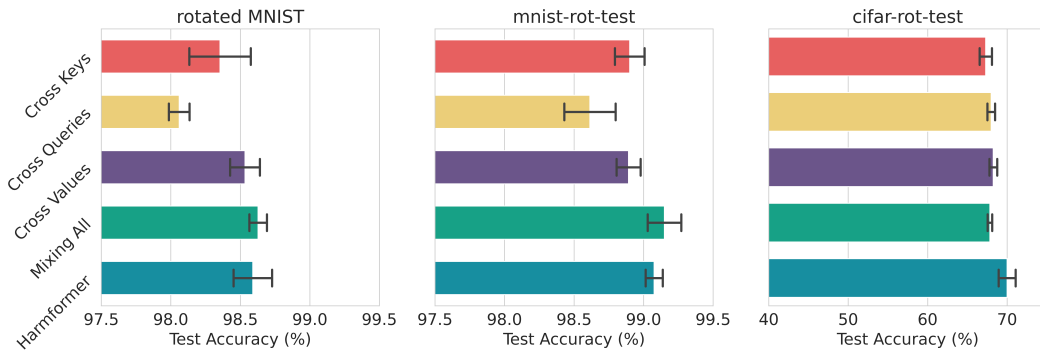


Figure 7: Ablation study on different SA mixing configurations. Each plot is aggregated from 5 different runs. The error bars represent the standard deviation.

696 B.3 Ablation: Relative Positional Encoding (RPE)

697 In Harmformer, we use relative circular encoding similar to those published in iRPE [48]. The
 698 encoding is added immediately after calculating the dot product between keys and queries. RPE sig-
 699 nificantly improves performance, as demonstrated in Figure 8, which shows Harmformer performance
 700 with and without RPE.

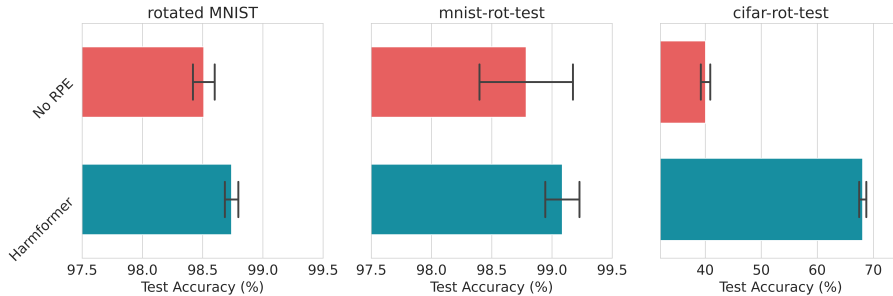


Figure 8: Ablation study on the use of relative circular position encoding (RPE) in the Harmformer. The error bars represent the standard deviation.

701 **B.4 Experiments: Stability of Classification w.r.t. Input Rotation**

702 In these experiments, we investigate the influence of input interpolation errors on performance,
 703 following previous invariant models [29, 28, 27]. Stability is primarily examined using the invariance
 704 benchmarks mnist-rot-test and cifar-rot-test, where the training data consists of non-rotated images,
 705 while the test data consists of randomly rotated images. Note that this implies that the training images
 706 consist of original sharp images, but the test images contain images with interpolation errors. In
 707 contrast, the rotated MNIST dataset contains rotated images in both the training and test sets, resulting
 708 in interpolation errors in both sets.

709 The test accuracy with respect to the input rotation is shown in Figure 9. Since the rotated MNIST
 710 dataset does not contain all images rotated by all angles, we use the original MNIST dataset [1]
 711 for this experiment. As a result, the test set, and therefore its accuracy, is different from that described in
 712 Section 6 of the main text.

713 For the mnist-rot-test, we observe very small oscillations, almost the same as for rotated MNIST. The
 714 accuracy reaches maxima at 0°, 90°, 180°, and 270°, where there is no interpolation effect. For the
 715 cifar-rot-test, the oscillations are more significant due to the low resolution of the dataset relative
 716 to the complexity of the objects, with minima at 45°, 135°, 225°, and 315°, where the interpolation
 717 errors are greatest.

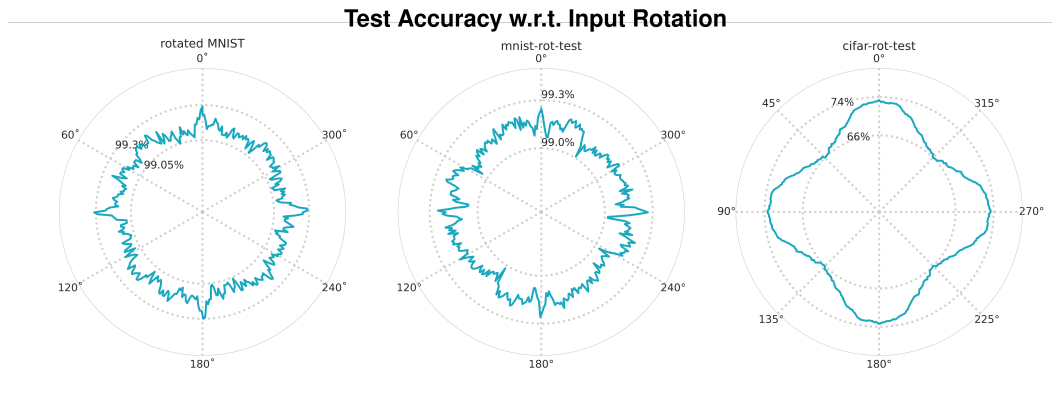


Figure 9: Classification stability with respect to input image rotation. The angle in the circular plots represents the rotation angle of the input image, and the radius represents the test accuracy on specified benchmarks.

718 For comparison with previous invariant models, we have included Table 7. For the mnist-rot-test,
 719 there is a significant improvement in Δ , which represents the difference in accuracy between the
 720 interpolation-free and interpolation-affected images. However, for the cifar-rot-test, the gap between
 721 these cases remains almost the same. For the MNIST datasets, the results are almost the same whether
 722 training with rotated or unrotated data. This leads to the hypothesis that if the resolution of the dataset

Table 7: Test accuracy comparison of the Harmformer and H-NeXt [27] at specific angles of rotation of inputs.

Model	mnist-rot-test			cifar-rot-test		
	0°	45°	Δ	0°	45°	Δ
H-NeXt [27]	98.9%	97.8%	1.1%	64.5%	57.4%	7.1%
Harmformer	99.2%	99.1%	0.1%	73.4%	66.1%	7.3%

723 matches the complexity of the recognition task, the Harmformer should not suffer from interpolation
 724 errors. However, this hypothesis would require further testing on large-scale datasets, which is beyond
 725 the scope of this paper.

726 B.5 Experiments: Evaluating the Role of Harmonic Convolutions

727 To investigate the importance of the convolutional stem and the encoder, we conducted an experiment
 728 using a minimal stem stage with an enlarged attentive field. The purpose of this setup was to ensure
 729 that the recognition was not due to convolution alone. This configuration contained only three
 730 convolutional layers and a single pooling layer.

731 The results, as shown in Table 8, indicate that the model performance remains within the expected
 732 error range despite the simplified convolutional stem. This suggested that the encoder plays a crucial
 733 role in the final classification. Notably, the inclusion of a single pooling layer significantly enhances
 734 the complexity of subsequent attention mechanisms. Due to the increased GPU RAM requirements,
 this configuration was exclusively tested on the rotated MNIST dataset.

Table 8: Performance comparison of the Harmformer architecture with a shallow stem stage on the rotated MNIST dataset.

Model	SA Input Shape	Test Error	Params.
Shallow Stem Stage	$32 \times 32 \times 16$	1.29%	40k
Harmformer	$16 \times 16 \times 16$	$1.26\% \pm 0.055$	30k

735

736 C Experimental Setup

737 C.1 Compute Resource

738 Each experiment was run on a single GPU within our shared, small but diverse cluster comprising 17
 739 GPUs. The cluster includes Tesla P100, V100, and A100 models, NVIDIA GeForce RTX 2080 Ti,
 740 3080, 4090, RTX A5000, and a Quadro P5000. Despite its limited size, our setup allowed for flexible
 741 and scalable computation using various GPU configurations. To provide a better overview, Table 9
 742 lists the epoch training time across each experiment on the NVIDIA GTX 4090.

Table 9: Training time of one epoch across different benchmarks on the NVIDIA GTX 4090.

GPU Model	mnist-rot-test	cifar-rot-test	rotated MNIST	PCam
Epoch Training time (mm:ss)	01:02	02:18	00:16	37:40
Number of Training Samples	52k	42k	10k	262k
Batch Size	32	32	32	8

743 C.2 Computation Complexity w.r.t. Non-Equivariant Convolution and SA Mechanism

744 In general, equivariant networks usually due to their properties impose higher computation complexity
 745 than their classical counterparts. For example, a single classical convolution has complexity $O(N^2 \cdot$
 746 $n^2)$, where $N \times N$ is the spatial dimension of the output feature map and $n \times n$ is the size of

747 the filter. In contrast, the original G-CNN equivariant to rotation and translation has complexity
 748 $O(N^2 \cdot n^2 \cdot |\theta|^2)$, where θ is the number of elements in the rotation group. Thus, a G-CNN equivariant
 749 to 90-degree rotations and translation would have $|\theta| = 4$.

750 **Harmformer Convolution** In Harmformer stem stage, we use convolution layers similar to H-Nets.
 751 This has a complexity of $O(N^2 \cdot n^2 \cdot |o|^2)$, where $|o|$ is the number of rotation orders of the input
 752 and output feature maps.

753 **Harmformer SA mechanism** Classical global SA mechanism has a complexity of $O(N^2 \cdot d + N \cdot d^2)$,
 754 where N is the number of patches and each patch has dimension d . Our SA mechanism, as shown in
 755 Figure 4b, adds multiplication by rotation orders o for matrix multiplication and dot product, resulting
 756 in a complexity of $O(o \cdot N^2 \cdot d + o \cdot N \cdot d^2)$.

757 **Additional Computational Considerations** Harmformer operates in the complex domain, where
 758 each multiplication requires four times and each addition requires two times more operations than
 759 their real counterparts. Additionally, the computational load increases due to upscaling the input and
 760 using large convolution kernels, as recommended in H-NeXt [27]. These factors also contribute to
 761 the overall complexity of Harmformer.

762 C.3 Configurations of Experiments

763 This subsection details the specific configurations of the Harmformer architecture used in the experi-
 764 ments described in Section 6. For convenience, Figure 10, which depicts the complete Harmformer
 765 architecture, is included. The parameters for each dataset are enumerated in the following tables:
 766 Tables 12 and 13 for the MNIST-rot-test and rotated MNIST datasets; Tables 14 and 15 for the
 767 CIFAR-rot-test dataset; and Tables 10 and 11 for the PCam dataset.

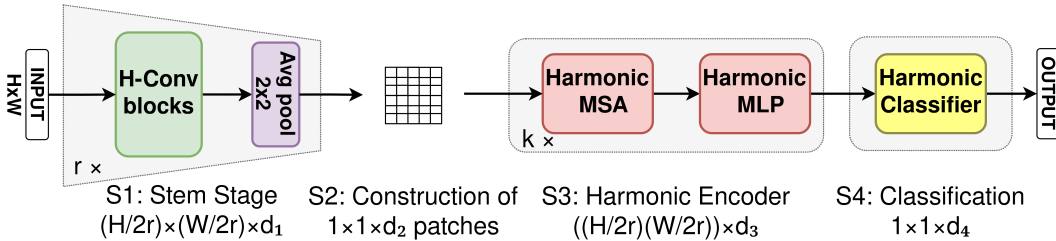


Figure 10: Reproduced from Figure 2 to detail the number of parameters for each experiment.

Table 10: PCam: Architecture

Parameter	Value
Number of S1 Blocks (r)	4
Convolution per S1 Block	2
Number of S1 Channels (d_1)	[4, 8, 16, 32]
S1 Channels Dropout	[0, 0.2, 0.3, 0.4]
Number of S3 Encoders (k)	4
Number of S3 Heads	4
Shape of S3 Patches (d_3)	8
MSA&MLP Dropout	0.4

Table 11: PCam: Training Settings

Parameter	Value
Epochs	100
Batch Size	8
Learning Rate	0.0007
Label Smoothing	0.1
Scheduler	Cosine
Optimizer	AdamW
Weight Decay	0.01
Runs	5
Input Padding	0

Table 12: MNIST datasets: Architecture

Parameter	Value
Number of S1 Blocks (r)	2
Convolution per S1 Block	2
Number of S1 Channels (d_1)	[8, 16]
S1 Channels Dropout	[0, 0]
Number of S3 Encoders (k)	3
Number of S3 Heads	1
Shape of S3 Patches (d_3)	16
MSA&MLP Dropout	0.1

Table 13: MNIST datasets: Training Settings

Parameter	Value
Epochs	100
Batch Size	32
Learning Rate	0.007
Label Smoothing	0.1
Scheduler	Reduce LR on Plateau
Optimizer	AdamW
Weight Decay	0.01
Runs	5
Input Padding	2

Table 14: cifar-rot-test: Architecture

Parameter	Value
Number of S1 Blocks (r)	2
Convolution per S1 Block	3
Number of S1 Channels (d_1)	[8, 16]
S1 Channels Dropout	[0, 0.1]
Number of S3 Encoders (k)	4
Number of S3 Heads	4
Shape of S3 Patches (d_3)	8
MSA&MLP Dropout	0.2

Table 15: cifar-rot-test: Training Settings

Parameter	Value
Epochs	200
Batch Size	32
Learning Rate	0.007
Label Smoothing	0.1
Scheduler	Cosine
Optimizer	AdamW
Weight Decay	0.01
Runs	5
Input Padding	0

768 **D Segmentation experiment: Retina blood vessel segmentation**

769 To demonstrate the generalizability and scalability of our architecture beyond classification tasks,
 770 we introduce Harmformer for retinal blood vessel segmentation using the DRIVE dataset [50]. The
 771 DRIVE is binary segmentation task, where the goal is to extract retinal blood vessels from an RGB
 772 image.

773 The dataset contains 262,080 samples for training and 65,520 samples for validation, similar to the
 774 settings of [51]. Each sample consists of an input image of size $3 \times 64 \times 64$ and a target segmentation
 775 mask of size 64×64 . These samples were generated from 17 training images and 3 validation images,
 776 each of which is 768×584 pixels and represents a different patient.

777 To use Harmformer as an image-to-image model, we adopt a U-Net [52] architecture in Fig 11a.
 778 Unlike our classification models (Section 6), this model processes the images at their original
 779 resolution, without any upscaling before they enter the network. For the output, we use only the
 780 magnitude of the final feature maps. To merge the hidden features (channels) into a single output
 781 layer, we apply a standard 2D convolution layer at the end.

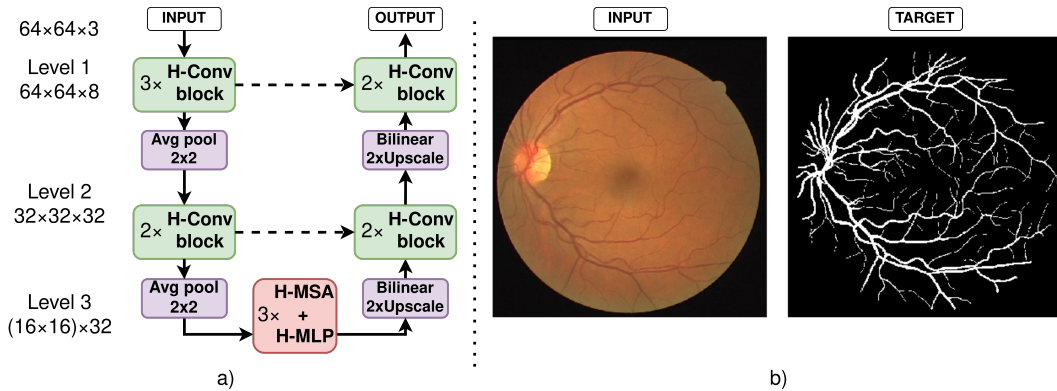


Figure 11: (a) Diagram of the Harmformer architecture for image segmentation. (b) Example of an image from the DRIVE dataset: the RGB image and the target segmentation mask.

782 We trained the U-Net Harmformer for 20 epochs with the AdamW optimizer, a learning rate of
 783 0.001 and 64 batch size. For augmentation, we used horizontal and vertical flipping, color jitter, and
 784 auto-contrast. We ran 4 different experiments with different seeds.

785 The results are shown in Table 16, using the area under the receiver operating characteristic curve
 786 (AUC) as the evaluation metric. For completeness, we have also included the performance of G-CNNs
 787 and the current state-of-the-art model FR-UNet [53]. As expected, these results are consistent with the
 788 findings in the paper, with Harmformer slightly underperforming compared to equivariant convolution
 789 architectures. Nevertheless, we show that our architecture is versatile and can also be applied to
 790 non-classification tasks.

Table 16: AUC for DRIVE segmentation [50].

Model	AUC	Equivariant model
Harmformer	0.9746 ± 0.0002	✓
G-CNNs [51]	0.9784 ± 0.0001	✓
FR-UNet [53]	0.9889	✗

791 E Differences Between 2D and 3D Equivariant Transformers

792 While 2D equivariant transformers [21, 22] have been relatively understudied, 3D equivariant trans-
793 formers [17, 14, 15, 18, 12, 10] have received more attention. In this section, we aim to highlight
794 the key differences that make the 2D case unique, and compare Harmformer with the most closely
795 related $SE(3)$ -Transformer, which operates in 3D but also uses steerable basis representations.

796 An important distinction lies in the nature of the input data, which directly influences the transformer
797 architecture. While 2D datasets typically consist of dense pixels with highly correlated neighborhoods,
798 3D equivariant datasets, often represented as graphs or point clouds, tend to be sparse. In 3D,
799 neighboring elements can vary significantly; for example, in molecular graphs [14, 15], atoms can
800 fulfill entirely different roles within the structure.

801 **Patches** The properties of the input data determine how to prepare patches in an equivariant manner.
802 In the case of the $SE(3)$ -Transformer, each node of the graph can be directly treated as a patch,
803 eliminating the need for a stem stage. For Harmformer, on the other hand, it is necessary to aggregate
804 low-level correlated data into a higher-level representation. Additionally, the classical (16×16)
805 ViT [30] grid cannot be used, as discussed in Section 3. Therefore, we employ a convolutional stem
806 stage, where the convolution kernels are expressed using circular harmonics to maintain equivariance.

807 This reliance on harmonic representations is a common feature between Harmformer and the $SE(3)$ -
808 Transformer. While Harmformer uses circular harmonics, the $SE(3)$ -Transformer uses spherical
809 harmonics. Both approaches leverage steerable basis functions [23], which are widely used in
810 equivariant networks [39, 34, 54]. These steerable bases change predictably under rotation, allowing
811 the effects of rotation to be effectively neutralized—via phase shifts in circular harmonics and via
812 the Wigner-D matrix in spherical harmonics. It is important to note that the use of steerable bases
813 predates both transformers, as shown in [23, 39, 54].

814 **Queries, Keys, and Values** In Harmformer, queries (Q), keys (K), and values (V) are generated
815 independently from individual patches through a linear layer, we proposed in Section 5.2. In contrast,
816 the $SE(3)$ -Transformer creates them by applying convolutions across points (i.e., patches), using
817 steerable spheres that aggregate information from the local neighborhood.

818 **Attention** The $SE(3)$ -Transformer focuses exclusively on invariant attention (type-0) and applies
819 only local attention. In contrast, Harmformer explores multiple strategies for mixing attention
820 and values of various orders (types), while performing global attention across the entire image.
821 Additionally, Harmformer introduces an equivariant layer normalization at the beginning of the
822 attention layer, while the $SE(3)$ -Transformer does not use any layer normalization. Other minor
823 distinctions include Harmformer’s use of an improved activation function and relative embeddings.

824 **NeurIPS Paper Checklist**

825 **1. Claims**

826 Question: Do the main claims made in the abstract and introduction accurately reflect the
827 paper’s contributions and scope?

828 Answer: [Yes]

829 Justification: All claims are supported by either theoretical proofs or experiments.

830 **2. Limitations**

831 Question: Does the paper discuss the limitations of the work performed by the authors?

832 Answer: [Yes]

833 Justification: All the limitations are discussed in the last section of the paper and all
834 assumptions are clearly stated.

835 **3. Theory Assumptions and Proofs**

836 Question: For each theoretical result, does the paper provide the full set of assumptions and
837 a complete (and correct) proof?

838 Answer: [Yes]

839 Justification: All theoretical statements are accompanied by proofs. The proofs are included
840 in the Appendix.

841 **4. Experimental Result Reproducibility**

842 Question: Does the paper fully disclose all the information needed to reproduce the main ex-
843 perimental results of the paper to the extent that it affects the main claims and/or conclusions
844 of the paper (regardless of whether the code and data are provided or not)?

845 Answer: [Yes]

846 Justification: The code is included in the submission, we only use open source datasets. All
847 trained models are also included and their settings are clearly stated in the appendix.

848 **5. Open access to data and code**

849 Question: Does the paper provide open access to the data and code, with sufficient instruc-
850 tions to faithfully reproduce the main experimental results, as described in supplemental
851 material?

852 Answer: [Yes]

853 Justification: Yes, we have included our source code with the submission. If the submission
854 is accepted, we will publish the code on github. The code includes preparation, download
855 scripts for all datasets, and the experimental settings.

856 **6. Experimental Setting/Details**

857 Question: Does the paper specify all the training and test details (e.g., data splits, hyper-
858 parameters, how they were chosen, type of optimizer, etc.) necessary to understand the
859 results?

860 Answer: [Yes]

861 Justification: This information is listed in the Appendix and is also included in the code
862 submission.

863 **7. Experiment Statistical Significance**

864 Question: Does the paper report error bars suitably and correctly defined or other appropriate
865 information about the statistical significance of the experiments?

866 Answer: [Yes]

867 Justification: We show the best results along with the average and standard deviation of at
868 least 5 runs with the performance of our models.

869 **8. Experiments Compute Resources**

870 Question: For each experiment, does the paper provide sufficient information on the com-
871 puter resources (type of compute workers, memory, time of execution) needed to reproduce
872 the experiments?

873 Answer: [Yes]
874 Justification: We enlisted the compute resources in appendix.

875 **9. Code Of Ethics**

876 Question: Does the research conducted in the paper conform, in every respect, with the
877 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

878 Answer: [Yes]
879 Justification: To the best of our knowledge, we have not violated any part of the NeurIPS
880 Code of Ethics.

881 **10. Broader Impacts**

882 Question: Does the paper discuss both potential positive societal impacts and negative
883 societal impacts of the work performed?

884 Answer: [NA]
885 Justification: Our paper focuses on representation learning and model robustness to transfor-
886 mations. We believe that there is no direct path to negative applications.

887 **11. Safeguards**

888 Question: Does the paper describe safeguards that have been put in place for responsible
889 release of data or models that have a high risk for misuse (e.g., pretrained language models,
890 image generators, or scraped datasets)?

891 Answer: [NA]
892 Justification: This paper poses no such risks that we are aware of.

893 **12. Licenses for existing assets**

894 Question: Are the creators or original owners of assets (e.g., code, data, models), used in
895 the paper, properly credited and are the license and terms of use explicitly mentioned and
896 properly respected?

897 Answer: [Yes]
898 Justification: We provide a citation to the paper presenting the benchmarks, models, and
899 the code included in the supplemental materials contains references to the libraries or
900 repositories from which we drew inspiration.

901 **13. New Assets**

902 Question: Are new assets introduced in the paper well documented and is the documentation
903 provided alongside the assets?

904 Answer: [NA]
905 Justification: This paper does not release new assets.

906 **14. Crowdsourcing and Research with Human Subjects**

907 Question: For crowdsourcing experiments and research with human subjects, does the paper
908 include the full text of instructions given to participants and screenshots, if applicable, as
909 well as details about compensation (if any)?

910 Answer: [NA]
911 Justification: This paper does not involve crowdsourcing nor research with human subjects.

912 **15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human
913 Subjects**

914 Question: Does the paper describe potential risks incurred by study participants, whether
915 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)
916 approvals (or an equivalent approval/review based on the requirements of your country or
917 institution) were obtained?

918 Answer: [NA]
919 Justification: This paper does not involve crowdsourcing nor research with human subjects.