

Manuscript version: Author's Accepted Manuscript

The version presented in WRAP is the author's accepted manuscript and may differ from the published version or Version of Record.

Persistent WRAP URL:

<http://wrap.warwick.ac.uk/140052>

How to cite:

Please refer to published version for the most recent bibliographic citation information. If a published version is known of, the repository item page linked to above, will contain details on accessing it.

Copyright and reuse:

The Warwick Research Archive Portal (WRAP) makes this work by researchers of the University of Warwick available open access under the following conditions.

Copyright © and all moral rights to the version of the paper presented here belong to the individual author(s) and/or other copyright owners. To the extent reasonable and practicable the material made available in WRAP has been checked for eligibility before being made available.

Copies of full items can be used for personal research or study, educational, or not-for-profit purposes without prior permission or charge. Provided that the authors, title and full bibliographic details are credited, a hyperlink and/or URL is given for the original metadata page and the content is not changed in any way.

Publisher's statement:

Please refer to the repository item page, publisher's statement section, for further information.

For more information, please contact the WRAP Team at: wrap@warwick.ac.uk.

Dense Steerable Filter CNNs for Exploiting Rotational Symmetry in Histology Images

Simon Graham, David Epstein and Nasir Rajpoot

Abstract—Histology images are inherently symmetric under rotation, where each orientation is equally as likely to appear. However, this rotational symmetry is not widely utilised as prior knowledge in modern Convolutional Neural Networks (CNNs), resulting in *data hungry* models that learn independent features at each orientation. Allowing CNNs to be rotation-equivariant removes the necessity to learn this set of transformations from the data and instead frees up model capacity, allowing more discriminative features to be learned. This reduction in the number of required parameters also reduces the risk of overfitting. In this paper, we propose Dense Steerable Filter CNNs (DSF-CNNs) that use group convolutions with multiple rotated copies of each filter in a densely connected framework. Each filter is defined as a linear combination of steerable basis filters, enabling exact rotation and decreasing the number of trainable parameters compared to standard filters. We also provide the first in-depth comparison of different rotation-equivariant CNNs for histology image analysis and demonstrate the advantage of encoding rotational symmetry into modern architectures. We show that DSF-CNNs achieve state-of-the-art performance, with significantly fewer parameters, when applied to three different tasks in the area of computational pathology: breast tumour classification, colon gland segmentation and multi-tissue nuclear segmentation.

Index Terms—Rotation-equivariance, steerable filters, deep learning, computational pathology.

THE recent advances in the analysis of Haematoxylin & Eosin (H&E) stained whole-slide images (WSIs) can largely be attributed to the rise of digital slide scanning [1]. In particular, Convolutional Neural Networks (CNNs) leverage the prior knowledge that images have translational symmetry and utilise a weight sharing strategy, which guarantees that a translation of the input will result in a proportional translation of the features. This property, known as *translation equiv-*

Copyright (c) 2019 IEEE. Personal use of this material is permitted. However, permission to use this material for any other purposes must be obtained from the IEEE by sending a request to pubpermissions@ieee.org.

S. Graham and N. Rajpoot are part of the PathLAKE digital pathology consortium, which is funded from the Data to Early Diagnosis and Precision Medicine strand of the government's Industrial Strategy Challenge Fund, managed and delivered by UK Research and Innovation (UKRI). DE and NR are also supported in part by the UK Medical Research Council (MRC) through award MR/P015476/1.

S.Graham is with the Mathematics for Real-World Systems Centre for Doctoral Training, University of Warwick, UK (email: s.graham.1@warwick.ac.uk).

N.Rajpoot is with the Department of Computer Science, University of Warwick, UK (email: n.m.rajpoot@warwick.ac.uk).

D.Epstein is with the Mathematics Institute, University of Warwick, UK. (email: david.epstein@warwick.ac.uk)

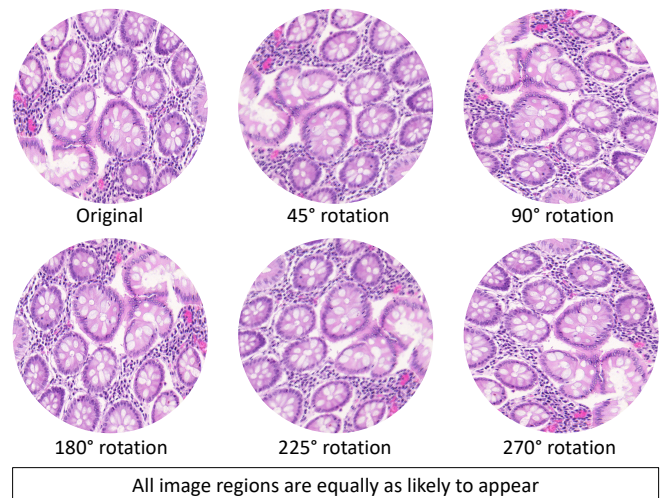


Fig. 1. Cropped circular regions from a whole-slide image. Each orientation is equally as likely to appear.

ariance, is an inherent property of the CNN and removes the need to learn features at all spatial locations, significantly reducing the number of learnable parameters. In certain image analysis applications, where there is no global orientation, it is desirable to extend this property of equivariance beyond translation to also rotation. One such example is the field of computational pathology (CPath) where important image features can appear at any orientation (Fig. 1). Therefore, we should be able to learn those features, regardless of their orientation. In the absence of rotation-equivariance, data augmentation is typically used, where multiple rotated copies of the WSI patches are usually introduced to the network during the training process. However, the augmentation strategy requires many more parameters in order to learn weights of different orientations. Instead, encoding rotational symmetry as a prior knowledge into current deep learning architectures by enforcing rotation-equivariance requires fewer parameters and leads to an overall superior discriminative ability. Also, rotation-equivariant CNNs typically converge quicker because the network does not need to spend time learning different filter orientations.

CPath is ripe ground for the utilisation of rotation-equivariant models, yet most models fail to incorporate this prior knowledge into the CNN architectures. Inspired by recent developments in the study of rotation-equivariant CNNs [2]–[5], we propose Dense Steerable Filter based CNNs (DSF-

CNNs) that integrate steerable filters [6] with group convolution [2] and a densely connected framework [7]. Dense connectivity enables efficient gradient propagation, encourages feature re-use and consequently leads to superior performance. Each filter is defined as a linear combination of circular harmonic basis filters, enabling exact rotation and significantly reducing the number of parameters compared to standard filters. The main contributions of this work are listed as follows:

- A Dense Steerable Filter CNN that achieves rotation-equivariance by integrating steerable filter group convolutions within a densely connected network.
- The first thorough comparison of multiple rotation-equivariant for CPath.
- We demonstrate state-of-the-art performance across multiple histology image datasets with far fewer parameters.

I. RELATED WORK

A. CNNs for translation equivariance

Images can contain numerous symmetries and therefore patterns may appear at various spatial positions and orientations. Recent methods [8] have shown that these symmetries can be detected, yet in this work we focus on how symmetries can be leveraged as a *prior knowledge* to increase the performance of image recognition algorithms. Pioneered by LeCun *et al.* in 1994 [9], CNNs inherently incorporate translation symmetry in images and achieve translation equivariance by re-using filters at all spatial locations. Therefore, a shift of the input leads to a proportional shift of the filter responses. This design drastically reduces the number of required parameters because features do not need to be learned independently at each location. Since the increase in computing power and the development of algorithms that assist network optimisation [10] CNNs have become deeper [7], [11], leading to current state-of-the-art performance in numerous image recognition tasks [12], [13]. As a result of the success of deep learning, CNNs have since been widely used in CPath for various tasks including: gland segmentation [14], [15]; nucleus segmentation [16]–[18]; mitosis detection [19]; cancer type prediction [20] and cancer grading [21], [22]. Yet, unlike translation, CNNs do not behave well with respect to rotation because this symmetry is not built into the network architecture.

B. Exploiting rotational symmetry

Rotating the data: It is well known that histology images have no global orientation and therefore standard practice is to apply rotation augmentation to the training data [23]. This improves performance, but requires many parameters and is therefore prone to overfitting. Also, there is no guarantee that CNNs trained with rotation augmentation will learn an equivariant representation and generalise to data with small rotations [24]. To reduce the variance of predictions of multiple orientations, test-time augmentation (TTA) can be used [25]. However, with TTA inference time scales linearly with the number of augmented copies. TI-Pooling [26] utilises multiple rotated copies of the input in a twin network architecture, where a pooling operation over orientations is performed to

find the optimal canonical instance of the input images for training. However, like TTA, TI-Pooling is computationally expensive.

Rotating the filters: Cohen & Welling [2] pioneered group equivariant CNNs (*G*-CNNs), where the convolution was generalised to share weights over additional symmetry groups beyond translation. However, they limited the filter transformation to 90° rotations and horizontal/vertical flips to ensure exact transformations on the 2D pixel grid. Veeling *et al.* [27] showed that these *G*-CNNs can be used to improve the performance of metastasis detection in breast histology images. Furthermore, Linmans *et al.* [28] and Graham *et al.* [29] extended the application of the *G*-CNNs proposed by Cohen & Welling to pixel-based segmentation in histology images, highlighting an improved performance over conventional CNNs. The symmetries of a square grid are limited to integer translations extended by the dihedral group of order 8 (4 reflections and 4 rotations). To counter the limitation of working with square grids in the *G*-CNN, Hoogeboom *et al.* [30] used hexagonal filters. However, this strategy requires images to be resampled on a hexagonal lattice, which is an additional overhead. Instead of using exact filter rotations, Bekkers *et al.* [31] and Lafarge *et al.* [5] applied *G*-CNNs to several medical imaging tasks by rotating filters with bilinear interpolation. Therefore, this method was not restricted to rotations by multiples of 90°, but may introduce interpolation artefacts. Oriented response networks [32] use active rotating filters during the convolution that explicitly encodes location and orientation information within the feature maps.

The aforementioned methods carry forward the feature maps for each orientation throughout the network. Instead, Marcos *et al.* [4] converted the output of multiple convolutions with rotated filter copies to a vector field by considering the magnitude and angle of the highest scoring orientation at every spatial location, leading to more compact models. To help overcome the issue of inexact filter rotation, the method only considered parameters at the centre of each filter and therefore required larger filters and consequently more parameters.

Rotating the feature maps: Dieleman *et al.* proposed a method similar to the *G*-CNN, but instead of rotating the filters, the feature maps were rotated. This design choice has no effect on the equivariance, yet any rotation that is not a multiple of 90° may suffer from interpolation artefacts.

Steerable filters: CNNs that encode rotation-equivariance are typically only equivariant to *discrete* rotations. Cohen & Welling [33] first proposed steerable CNNs and described a general mathematical theory that applies to both continuous and discrete groups. To achieve full 360° equivariance, Worrall *et al.* [34] used the concept of steerable filters [6] and constrained the weights to be complex circular harmonics. Cheng *et al.* [35] propose a rotation-equivariant CNN, named RotDCF, that decomposes filters over joint steerable bases across the space and the group geometry simultaneously. Weiler *et al.* [3] learned steerable filters as a linear combination of atomic basis filters, which enabled exact filter rotation within *G*-CNNs. Then, these steerable filters were used within the group convolution to enable the network to be equivariant to rotation. Weiler & Cesa [36] then performed an extensive

comparison of rotation equivariant models using steerable filters. Our method builds on the approach proposed by Weiler *et al.* [3], by incorporating steerable filter group convolutions into a densely connected framework for superior performance.

II. MATHEMATICAL FRAMEWORK

In this section we present the key mathematical concepts used in our framework. We first describe images, filters and feature maps as functions. We introduce steerable filters and describe the group-convolution (G -convolution) operation with these filters. This operation leads to G -equivariance. Below, we deal with a single filter at a time, although the method actually needs a whole filter bank to be used. We follow the method described by Weiler *et al.* [3], but we use a slightly different formulation. We encourage readers to read both approaches.

A. Images and feature maps as functions

We model an image as a map $f : \mathbb{C} \cong \mathbb{R}^2 \rightarrow \mathbb{R}$ with compact support¹. Let \mathcal{F} be the vector space over \mathbb{R} of all $f : \mathbb{C} \rightarrow \mathbb{R}$, with compact support, and let $\mathcal{F}_{\mathbb{C}}$ be the vector space over \mathbb{C} of all functions $f : \mathbb{C} \rightarrow \mathbb{C}$ with compact support.

We denote by $\text{SE}(2)$ the group of isometries of the plane, omitting reflections. Each element of $\text{SE}(2)$ can be written in the form $z \mapsto e^{i\theta}z + b$, where $z, b \in \mathbb{C}$ and $\theta \in \mathbb{R}$. If $g \in \text{SE}(2)$ and $f \in \mathcal{F}$, we define $g.f \in \mathcal{F}$ by:

$$(g.f)(z) = f(g^{-1}(z)) \text{ for } z \in \mathbb{C}. \quad (1)$$

The same definition is used for $g.f : \mathbb{C} \rightarrow \mathbb{C}$ when $f \in \mathcal{F}_{\mathbb{C}}$.

B. Steerable functions and filters:

The additive group of real numbers \mathbb{R} acts on \mathbb{C} by rotations keeping 0 fixed. By (1), it acts linearly on \mathcal{F} (and on $\mathcal{F}_{\mathbb{C}}$):

$$f^\theta(z) = f(e^{-i\theta}z) \text{ for } f \in \mathcal{F}, \theta \in \mathbb{R}.$$

We define $V(f) \subset \mathcal{F}_{\mathbb{C}}$ to be the complex vector subspace spanned by the orbit $\{f^\theta \mid \theta \in \mathbb{R}\}$. If $V(f)$ is a finite dimensional vector space, we say that f is *steerable*.

Theorem: A necessary and sufficient condition for $\psi \in \mathcal{F}_{\mathbb{C}}$ to be steerable is that there should exist an integer $A \geq 0$, and radial profile functions $R_k : [0, \infty) \rightarrow \mathbb{C}$ for $k \in \mathbb{Z}$ and $-A \leq k \leq A$, such that, in polar coordinates:

$$\psi(r, \varphi) = \sum_{k=-A}^A R_k(r) e^{ik\varphi}, \quad (2)$$

where some or all of the radial profile functions R_k may be identically zero. To ensure that ψ has compact support, each R_k is assumed to have compact support.

If ψ satisfies (2), then $V(\psi)$ is clearly finite dimensional. The reverse implication takes a bit longer to argue, but easily follows from standard theorems in Group Representation Theory².

Fig. 2 is a graphical representation of basis harmonic filters that appear in (2).

¹The *support* of f is the smallest closed subset of \mathbb{C} containing $\{z \in \mathbb{C} \mid f(z) \neq 0\}$.

²For full mathematical rigour, the theorem requires the additional hypothesis that, for each r , ψ is a continuous function of φ . See also [37] for more technical details.

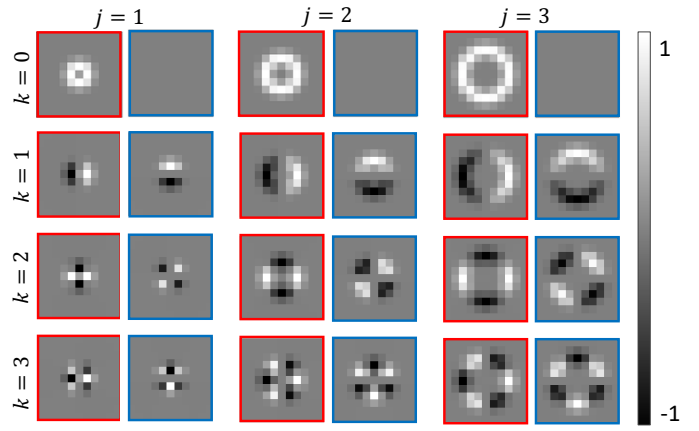


Fig. 2. Example circular harmonic basis filters sampled on the 11×11 square grid. Red and blue borders denote the real and imaginary parts respectively. Each pair of images comes from a single term $R_k(r)e^{ik\theta}$ in (2). In this Fig., the particular radial profile functions R_k are all Gaussians, as they are in our proposed model. These Gaussians have mean/mode/max at j . The integer k specifies the frequency.

Real Version: In practice we will work with steerable real-valued filters. Since a real-valued steerable filter ψ is also a complex-valued steerable filter, we can apply (2) to obtain, in the same notation:

$$\psi(r, \varphi) = \text{Re} \left(\sum_{k=-A}^A R_k(r) e^{ik\varphi} \right).$$

Now $\text{Re}(z) = (z + \bar{z})/2$. It follows that we can write instead (but the radial profiles change):

$$\psi(r, \varphi) = \text{Re} \left(\sum_{k=0}^A R_k(r) e^{ik\varphi} \right) \quad (3)$$

where $R_0 : [0, \infty) \rightarrow \mathbb{R}$ and, for $k > 0$, $R_k : [0, \infty) \rightarrow \mathbb{C}$.

C. Feature maps modelled on a group:

Following the pioneering work of Cohen and Welling [2] and of Weiler *et al.* [3], we explain the changes to the architecture of CNNs, required to express rotation equivariance.

We fix an integer $n > 0$. We use the symbol $\rho_{u, \theta}$ to denote the euclidean transformation given by

$$\rho_{u, \theta}(z) = e^{i\theta}z + u, \quad (4)$$

where $u \in \mathbb{C}$ and $\theta = 2\pi s/n$, for some integer s with $0 \leq s < n$. Let $G \subset \text{E}(2)$ be the subgroup of all such transformations.

Let U be a group, with two subgroups U_1 and U_2 . U is said to be a *semidirect product* of U_1 with U_2 , denoted by $U_1 \rtimes U_2$, if there are projections $p_1 : U \rightarrow U_1$ and $U \rightarrow U_2$ —this means that $p_1|_{U_1}$ and $p_2|_{U_2}$ are both identity maps—such that p_2 is a homomorphism with kernel U_1 , and $p_1 \times p_2 : U \rightarrow U_1 \times U_2$ is a bijection, but, in general, not an isomorphism of groups. The importance of this concept in the study of equivariant CNNs was first pointed out in [2], and there is a systematic study [36].

G has two important subgroups, namely

$$C_n = \{\rho_{0, \theta} \mid \theta = 2\pi s/n, 0 \leq s < n\}, \quad (5)$$

a cyclic subgroup of order n consisting of all rotations in G keeping $0 \in \mathbb{C}$ fixed and

$$T = \{\rho_{u,0} \mid u \in \mathbb{C}\} \cong \mathbb{C},$$

consisting of all translations of \mathbb{C} . We define the group

$$C'_n = \{\theta \mid \theta = 2\pi s/n, 0 \leq s < n\}, \quad (6)$$

with group law addition mod 2π . Clearly, $C_n \cong C'_n$. We also use $\{e\} \cong C_1$ to denote the trivial group with one element.

The bijection

$$\Pi : G \rightarrow \mathbb{C} \times C'_n \text{ defined by } \Pi(\rho_{u,\theta}) = (u, \theta) \quad (7)$$

gives G the semidirect product structure $G = T \rtimes C'_n$. We impose on $\mathbb{C} \times C'_n$ a product metric that is the same as the usual Euclidean metric on \mathbb{C} , and is any convenient fixed metric on the finite discrete space C'_n . The bijection Π is then used to impose a metric on G , so that Π becomes an isometry. Π does not preserve the group structure, unless $n = 1$.

As a metric space G is the disjoint union of the n right cosets

$$\mathbb{C}_\theta = T\rho_{0,\theta} = \{\rho_{u,\theta} \mid u \in \mathbb{C}\} \subset G \text{ for } \theta \in C'_n, \quad (8)$$

such that each coset is isometric to \mathbb{C} .

A G -feature map is defined to be a function $f : G \rightarrow \mathbb{R}$, with compact support.

D. \mathcal{G} -convolutions:

We generalize the concept of a convolution to a G -convolution, that maps one G -feature map to another.

We give the definition of \mathcal{G} -convolution, where \mathcal{G}^3 is a group with a measure $\mu_{\mathcal{G}}$ —this means that, given $f : \mathcal{G} \rightarrow \mathbb{R}$, we can form the integral denoted by $\int_{g \in \mathcal{G}} f(g) d\mu_{\mathcal{G}}$ or $\int_{g \in \mathcal{G}} f(g) dg$. We will stick to the *unimodular* case, which is general enough for all cases of interest in this paper. The word *unimodular* means that we can change the dummy variable g in the integral to g^{-1} , or gh or hg ($h \in \mathcal{G}$ constant), without changing the value of the integral.

Given maps $f : \mathcal{G} \rightarrow \mathbb{R}$ and $\psi : \mathcal{G} \rightarrow \mathbb{R}$, we define their \mathcal{G} -convolution $(f *_G \psi) : \mathcal{G} \rightarrow \mathbb{R}$ by

$$\begin{aligned} (f *_G \psi)(g) &= \int_{h \in \mathcal{G}} f(gh^{-1})\psi(h) dh \\ &= \int_{h \in \mathcal{G}} f(h)\psi(h^{-1}g) dh \text{ for } g \in \mathcal{G}. \end{aligned} \quad (9)$$

The first equality is a definition, whereas the second follows by a change of variable.

\mathcal{G} -convolution is automatically \mathcal{G} -equivariant. To see this, note that, for any $\alpha \in \mathcal{G}$,

$$\begin{aligned} (\rho_\alpha(f) *_G \psi)(g) &= \int f(\alpha^{-1}gh^{-1})\psi(h) dh \\ &= (f *_G \psi)(\alpha^{-1}g) = (\rho_\alpha(f *_G \psi))(g). \end{aligned}$$

It follows that

$$\rho_\alpha(f) *_G \psi = \rho_\alpha(f *_G \psi). \quad (10)$$

³We use \mathcal{G} instead of G because we have reserved the name G for the particular group defined in Subsection II-C and \mathcal{G} denotes an arbitrary group.

E. Hidden layer G -convolutions and G -filters

By a G -filter, we mean a function $G \rightarrow \mathbb{R}$. Formally this is the same as a G -feature map. However, in an implementation of these ideas, a G -feature map will turn out to be a discrete object, specified by a collection of matrices, whereas a G -filter retains its identity as a function. This is what enables exact rotation of a G -filter by an arbitrary angle.

In order to define G -convolutions, we need a measure on the space G , as described for \mathcal{G} in Subsection II-D. The measure μ_G on G is given by using the usual euclidean (area) measure on each $\mathbb{C}_\theta \cong \mathbb{C}$. Note that (G, μ_G) is *unimodular* (term defined in Subsection II-D) because rotation is measure preserving on the plane. Integration of a function $f : G \rightarrow \mathbb{R}$, with respect to μ_G , is carried out by first integrating each of the n functions $f|_{\mathbb{C}_\theta} \cong \mathbb{C} \rightarrow \mathbb{R}$ and adding the n resulting terms.

We now define an “*atomic steerable planar filter*”, which is not learned, but defined and does not change during training (see (13)). Instead our network learns the complex coefficients used in a complex linear combination of the atomic steerable planar filters.

For each non-negative integer j , we define $\tau_j : [0, \infty) \rightarrow \mathbb{R}$ to be a Gaussian, with mode at j , as

$$\tau_j(r) = \exp(-|r - j|^2/2\sigma^2) \text{ for } j \geq 0, r \geq 0. \quad (11)$$

Let j and k be non-negative integers. By a *atomic steerable planar filter*, we mean a map $\psi_{jk} : \mathbb{C} \rightarrow \mathbb{C}$ defined by

$$\psi_{jk}(u) = \tau_j(|u|)e^{ik \arg(u)}. \quad (12)$$

If, in addition, $\lambda \in C'_n$, we define the *atomic steerable G -filter* $\psi_{jk\lambda} : G \rightarrow \mathbb{R}$ by

$$\psi_{jk\lambda}(\rho_{u,\theta}) = \begin{cases} 0 & \text{if } \lambda \neq \theta \\ \tau_j(|u|)e^{ik(\arg(u)-\theta)} & \text{if } \lambda = \theta. \end{cases} \quad (13)$$

From (12)

$$\psi_{jk\lambda}(\rho_{u,\theta}) = e^{-ik\theta}\psi_{jk}(u) \text{ if } \theta = \lambda, \quad (14)$$

which is ψ_{jk} rotated by angle θ .

Any finite complex linear combination of atomic steerable G -filters, $\sum_{j,k,\lambda} w_{jk\lambda}\psi_{jk\lambda}$, is again a steerable G -filter. In our framework, we plan to convolve each G -feature map with the real part of such a sum. By (9) the result of such a convolution is another G -feature map. The complex numbers $w_{jk\lambda}$ are weights in the network, determined by the network during training and each $w_{jk\lambda}$ gives rise to two real weights. We will initially restrict to a single term in the finite sum, in order to keep the formulas uncluttered, and then add them together.

Let $f : G \rightarrow \mathbb{R}$ be a G -feature map. From (9), we have the formula

$$(f *_G \text{Re}(w_{jk\lambda}\psi_{jk\lambda}))(\rho_{z,\theta}) = \int_{\rho_{u,\varphi} \in G} f(\rho_{u,\varphi}) \cdot \text{Re}(w_{jk\lambda}\psi_{jk\lambda}(\rho_{v,\beta})) d\mu_G, \quad (15)$$

where $\rho_{v,\beta} = \rho_{u,\varphi}^{-1} \rho_{z,\theta}$, so that $v = e^{-i\varphi}(z-u)$ and $\beta = \theta - \varphi$. From (12) and (13),

$$\psi_{jk\lambda}(\rho_{v,\beta}) = \begin{cases} 0 & \text{if } \lambda \neq \beta = \theta - \varphi \\ e^{-ik\varphi} \cdot \psi_{jk}(z-u) & \text{if } \lambda = \beta = \theta - \varphi. \end{cases} \quad (16)$$

Writing $f_\varphi(u) = f(\rho_{u,\varphi})$, we obtain from (15) and (16)

$$\begin{aligned} & (f *_G \text{Re}(w_{jk\lambda} \psi_{jk\lambda}))(\rho_{z,\theta}) \\ &= \text{Re} \left(w_{jk\lambda} \cdot e^{-ik(\theta-\lambda)} \cdot (f_{\theta-\lambda} * \psi_{jk}) \right) (z) \\ &= \left(f_{\theta-\lambda} * \text{Re}(w_{jk\lambda} \cdot e^{-ik(\theta-\lambda)} \psi_{jk}) \right) (z). \end{aligned} \quad (17)$$

If we add over $\lambda \in C'_n$, then we can substitute $\varphi = \theta - \lambda$ and add over $\varphi \in C'_n$, since θ is fixed in (17). Adding over j, k and φ , we obtain

$$\begin{aligned} & \left(f *_G \text{Re} \left(\sum_{jk\lambda} w_{jk\lambda} \psi_{jk\lambda} \right) \right) (\rho_{z,\theta}) \\ &= \sum_{jk\varphi} \left(f_\varphi * \text{Re} \left(w_{jk(\theta-\varphi)} \cdot e^{-ik\varphi} \psi_{jk} \right) \right) (z) \end{aligned} \quad (18)$$

which recovers the same result as (10) in [3]. We have ignored the fact that there are normally many channels (G -feature maps) in the domain and many channels in the range. Each pair (channel in domain, channel in base) needs its own G -filter, so each such pair gives rise to different weights.

F. The input layer G -convolution

The input to network is an image that can be thought of as a map $f : \mathbb{C} \rightarrow \mathbb{R}$, which we compose with $P : G \rightarrow \mathbb{C}$ given by $P(\rho_{u,\theta}) = u$, to obtain $f \circ P : G \rightarrow \mathbb{R}$. By (17), we have

$$\begin{aligned} & ((f \circ P) *_G \text{Re}(w \psi_{jk\lambda}))(\rho_{z,\theta}) = \\ & \text{Re}((w_{jk\lambda} \cdot e^{ik\lambda}) \cdot e^{-ik\theta} \cdot (f * \psi_{jk})(z)) \end{aligned}$$

Since $w_{jk\lambda}$ is a complex scalar that the network has to estimate, λ adds no new information and we dispense with it. We then sum over all terms, obtaining a simplified version of (18).

$$\begin{aligned} & \left((f \circ P) *_G \text{Re} \left(\sum_{jk} w_{jk} \psi_{jk} \right) \right) (\rho_{z,\theta}) \\ &= \left(f * \text{Re} \left(\sum_{jk} w_{jk} \cdot e^{-ik\theta} \cdot \psi_{jk} \right) \right) (z). \end{aligned} \quad (19)$$

This gives a principled derivation of Equation (8) in [3]. In particular, our proof of G -equivariance (see (10)) works equally well for input layer and hidden layer G -convolutions. See Fig. 3(b) for a graphical illustration of the method.

G. Sampling and the discrete case

The above formulas assume that the functions involved are continuous. But a computer is a finite machine, so we need to work with discrete data, and this involves sampling.

Sampling planar steerable filters: In the computer, a planar feature map is represented by a matrix, not by a continuous

function. According to (18) and (19), we need to convolve this matrix with the real part of a complex linear combination of atomic planar filters, ψ_{jk} . Now ψ_{jk} is a function, not a matrix—this is exactly what allows rotation of the filter through an arbitrary angle. On the other hand, convolution with a matrix requires a matrix, not a function. We therefore have to sample the atomic filters ψ_{jk} , and their rotations through angles $2\pi s/n$ for $0 \leq s < n$, at the integer points $a + ib$, where a and b are integers. We then perform a weighted linear combination of the sampled filters and apply (18) or (19). As the Nyquist Sampling Theorem suggests, for a fixed size of steerable filter, aliasing may occur unless one bounds the frequencies used from above. In line with Weiler & Cesa [36], we use frequencies up to $k = 0, 2, 3, 2$ for $j = 0, 1, 2, 3$ in all 7×7 steerable basis filters. Using larger filters enables higher frequencies before aliasing, yet leads to an increase in computation time and may lead to overfitting.

Sampling G -filters: As in the case of planar convolution just discussed, our formulas need to be reinterpreted when the various component pieces of a hidden layer G -convolution are formulated as arrays of dimension 3 or higher, rather than as functions. For example a G -feature map has been defined as a function $G \rightarrow \mathbb{R}$, and we need to explain how a function on the continuous group G is represented in the computer by n matrices.

As shown in (8), G as a metric space is the disjoint union $\bigcup_{\theta \in C'_n} \mathbb{C}_\theta$ of n copies of \mathbb{C} , with its usual euclidean metric. For each $\theta \in C'_n$ (see (8)) we define

$$\mathbb{Z}_\theta = \{\rho_{a+ib,\theta} \mid a, b \in \mathbb{Z}\} \subset \mathbb{C}_\theta. \quad (20)$$

Each point of \mathbb{C}_θ is within a distance $1/\sqrt{2}$ of some point in the lattice \mathbb{Z}_θ . It is therefore reasonable to use, as a G -feature map,

$$f : \bigcup_{\theta \in C'_n} \mathbb{Z}_\theta \rightarrow \mathbb{R}. \quad (21)$$

Analogously to the notation just before (17), we write $f_\theta = f|_{\mathbb{Z}_\theta}$. The domain is infinite, but since f is assumed to have compact support, we need only record the values of f at a finite number of elements of G . In this way, a G -feature map is replaced by n real matrices all of the same size.

We have also defined a G -filter as a function $G \rightarrow \mathbb{R}$. This is also sampled on $\bigcup_{\theta \in C'_n} \mathbb{Z}_\theta$. When learning the complex coefficients $w_{jk\lambda}$ that appear in (15), the values of j and k are limited for the reasons just explained for the planar situation, namely to avoid aliasing and overfitting.

III. DENSE STEERABLE FILTER CNN

A. Network architecture

The main building blocks of our proposed rotation-equivariant DSF-CNN⁴ are: an input layer G -convolution layer; steerable filter G -dense-blocks and a G -pooling layer. Below, we build on the theoretical explanation in Section II to describe the separate components of our proposed approach.

⁴Model code: <https://github.com/simongraham/dsf-cnn>

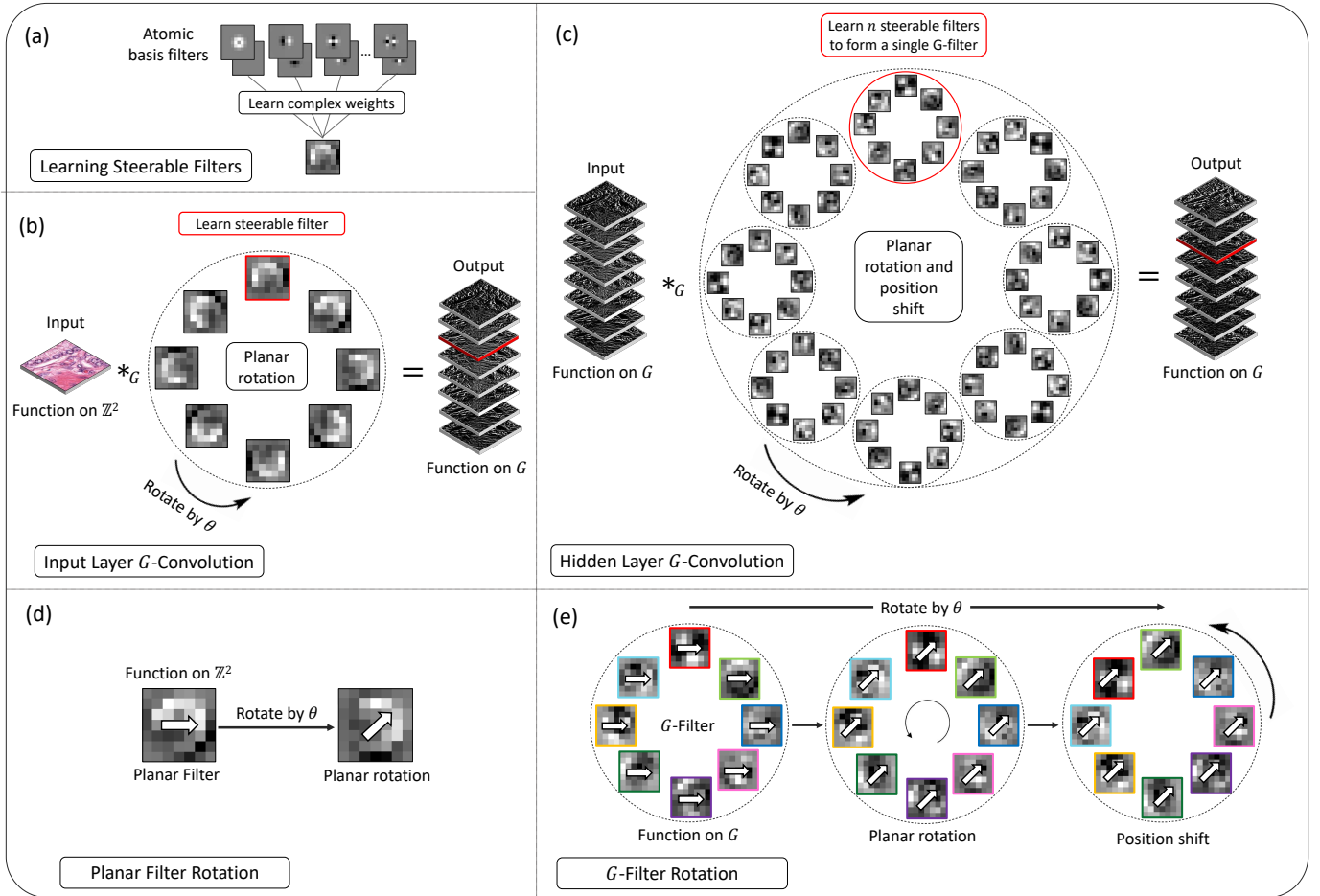


Fig. 3. Overview of the two types of G -convolution used in our approach with 8 filter orientations- best viewed in colour. a) Generation of steerable filters by linearly combining a series of atomic basis filters. b) Illustration of the input layer G -convolution, mapping an image $f : \mathbb{C} \rightarrow \mathbb{R}$ to a G -feature map $h : G \rightarrow \mathbb{R}$. A single steerable planar filter, learned by the network, is rotated n times and each rotated filter is convolved with the planar input f . This gives n planar feature maps, which combine to give a single G -feature map h . The image f is convolved with the red bordered planar filter to give the red bordered planar feature map in the stack on the right. c) Illustration of the hidden layer G -convolution, mapping a G -feature map $f : G \rightarrow \mathbb{R}$ to a G -feature map $h : G \rightarrow \mathbb{R}$. The network learns a single steerable G -filter, which consists of n planar filters, displayed by placing them all in the same circle. Then, a single G -filter is rotated n times and each rotated G -filter is convolved with the input G -feature map f to generate a total of n planar feature maps or a single G -feature map. The convolution between the input f and the red circled G -filter gives the red bordered planar feature map on the right. d) demonstrates rotation of a planar filter, as used in the input layer G -convolution and e) demonstrates rotation of a G -filter, used in the hidden layer G -convolution. It can be seen from e) that G -filters undergo an additional position shift, in line with the group action. In both d) and e), $\theta = \frac{\pi}{4}$.

Input Layer G -convolution: Up to the G -pooling operation, all convolutions within our network are steerable G -convolutions, as described in Section II-E. Therefore, we pre-define a set of circular harmonic basis filters using (2) and sample the filters on the square grid, as can be seen in Fig. 2. Then, we learn how to linearly combine these atomic basis filters to generate steerable filters and consider only the real part for our convolution filter, as shown in (3). This can be visualised in Fig. 3a. The input layer steerable G -convolution maps an image $f : \mathbb{C} \rightarrow \mathbb{R}$ to some G -feature map $h : G \rightarrow \mathbb{R}$. Each G -feature map is determined by its restriction h_θ to each coset $\mathbb{C}_\theta \cong \mathbb{C}$. Specifically, we create n rotated copies of each steerable filter and independently convolve the filters with the input to produce n feature maps (or a single G -feature map). Planar rotation of each filter is performed using (14) and can be observed in Fig. 3d. The input layer G -convolution is demonstrated in Fig. 3b, where the convolution between

the input and the steerable filter bordered in red produces the output also bordered in red. Now, when the input is rotated by an angle $\frac{2\pi s}{n}$, with integers $0 \leq s < n$, and the input layer G -convolution is performed, the feature maps undergo a planar rotation by angle $\frac{2\pi s}{n}$, but in addition shift s positions.

G -dense-blocks: To enable efficient gradient propagation, encourage feature re-use and to improve overall performance, we use dense connectivity [7] between G -convolutions in hidden layers of the network. Each hidden layer steerable G -convolution maps a G -feature map $f : G \rightarrow \mathbb{R}$ to some G -feature map $h : G \rightarrow \mathbb{R}$. We can explain this mapping in terms of the restrictions of f and h to cosets. Because the input to the hidden layer G -convolution is now a function on G , we must similarly ensure that our filters give a function on G . We rotate each G -filter to give n rotated copies and perform a convolution between the input G -feature map f and each filter orientation to produce n feature maps (or a single G -feature

map h). When rotating these G -filters, an additional position shift must be performed, in line with the associated group action. In Fig. 3c, $n = 8$ steerable planar filters are generated as shown by the red circle, forming a single G -filter. This G -filter is convolved with the input G -feature map to generate the output with the red border. We can see that each G -filter, consists of 8 planar filters that individually rotate and shift position as the entire G -filter is rotated. This rotation can be seen in Fig. 3e, where the arrows show the orientation of each planar filter and the coloured borders are used to help visualise the position of each planar filter in the G -filter.

For each G -dense-block, the feature-maps of all preceding layers are concatenated to the input before performing the G -convolution. This increases the number of connections between layers, strengthening feature propagation. Specifically, each G -dense-block consists of k units. Each unit contains a 7×7 G -convolution followed by a 5×5 G -convolution that produce 14 and 6 orientation dependent feature maps respectively. After k units, the G -dense-block concludes by applying a final 5×5 G -convolution.

G -pooling: At the output of the network, we transform each G -feature map f to a planar feature map, by taking the pointwise maximum of the n planar feature maps f_θ that constitute f . This operation ensures that the output of G -pooling is *invariant* to rotation of the input.

G -Batch-Normalisation: Batch normalisation (BN) involves two trainable parameters that scale and shift the normalised output. Standard BN is applied to the output of all feature maps and therefore learned BN parameters are typically different for each planar feature map in the group G . However, when the input is rotated, BN parameters will not transform in accordance with the input and therefore standard BN is not rotation-equivariant. Instead, after each G -convolution, we use a group-equivariant batch normalisation that aggregates moments per group rather than spatial feature map. This is essential to ensure rotation-equivariance throughout the network.

Classification: For our classification DSF-CNN, we initially perform the input layer steerable G -convolution followed by a hidden layer G -convolution. We then use 4 G -dense-blocks, where each block consists of 3,4,5 and 6 dense units. After every G -convolution layer we use a group-equivariant batch normalisation that aggregates moments per group rather than spatial feature map and ReLU non-linearity. Before every G -dense-block, we perform spatial max-pooling to decrease the dimensions of the feature maps. After the final G -dense-block, we perform G -pooling and then apply 3 1×1 classical convolution operations to get the final output.

Segmentation: We extend our DSF-CNN to the task of segmentation by up-sampling feature maps after the final G -dense-block in the aforementioned classification CNN. Specifically, we up-sample by a factor of 2 with bilinear interpolation and then utilise a G -dense-block. This is repeated until the spatial dimensions of the original image are regained. From the deepest layer of the up-sampling branch, each dense-block contain 4, 3 and 2 units. In line with U-Net [38], we also use skip connections to propagate information from the encoder to the decoder. After the feature maps have been up-sampled, we

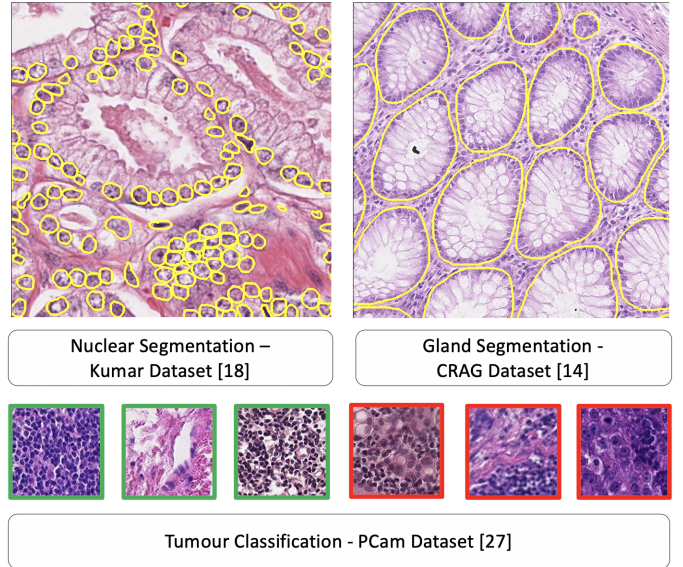


Fig. 4. Image regions from the three datasets. For nuclear segmentation, gland segmentation and tumour classification, we use the Kumar [18], CRAG [14] and PCam [27] datasets. Yellow boundaries show the pathologist annotation, while green and red borders denote non-tumour and tumour image patches.

use a single hidden layer G -convolution, which is followed by G -pooling such that the resulting feature map is a function on \mathbb{C} . Finally we use 2 1×1 classical convolutions to obtain the output, where we segment both the object and the contour to help separate touching instances. For nuclear segmentation, we additionally predict the eroded nuclei masks which are used as markers in marker-controlled watershed.

IV. EXPERIMENTS AND RESULTS

A. Experimental overview

Recently, there has been a growing number of proposed CNNs that achieve rotation-equivariance [2]–[4], [31], [34], yet there is lack of comprehensive evaluation of the various methods for the analysis of histopathology images. We perform a thorough comparison of various rotation-equivariant CNNs and demonstrate the effectiveness of the proposed model. Specifically, we compare a baseline CNN with H-Nets [34], VF-CNNs [4], G -CNNs with standard filters [2], [31] and G -CNNs with steerable filters [3] and assess the impact of increasing the number of filter rotations in each model. For a thorough analysis, each method is applied to the task of breast tumour classification and then the best performing models are applied to the tasks of nucleus and gland segmentation. After gaining an insight into the performance of the different rotation-equivariant models, we then compare our proposed Dense Steerable Filter CNN with the state-of-the-art methods on each of the three datasets used in our experiments.

B. The three datasets

We use the following three publicly available histology image datasets:

Breast tumour classification: PCam [27] is a dataset of 327K

image patches of size 96×96 pixels at $10 \times$ extracted from the Camelyon16 dataset [39], containing 400 H&E stained breast WSIs. Each image patch was labelled as tumour if the central region (32×32) contained at least one tumour pixel as given by the original annotation [39].

Multi-tissue nucleus segmentation: The Kumar [18] dataset contains 30 $1,000 \times 1,000$ image tiles from seven organs (6 breast, 6 liver, 6 kidney, 6 prostate, 2 bladder, 2 colon and 2 stomach) of The Cancer Genome Atlas (TCGA) database acquired at $40 \times$ magnification. Within each image, the boundary of each nucleus is fully annotated.

Colorectal gland segmentation: The CRAG dataset [14] consists of 213 H&E images mostly of size $1,512 \times 1,516$ pixels taken from 38 WSIs acquired at $20 \times$ of colorectal adenocarcinoma (CRA) patients. It is split into 173 training images and 40 test images with different cancer grades with pixel-based gland annotation.

C. Evaluation metrics

Here we describe the metrics used for evaluation. For tumour classification, we calculated the area under the receiver operating characteristic curve (AUC) to assess the binary classification performance. For gland segmentation, we employed the same quantitative measures that were used in the GlaS challenge [40]. These metrics consist of F_1 , DICE and Hausdorff distance at the object level and assess the quality of instance segmentation. For nuclear segmentation, we report the binary DICE and panoptic quality (PQ). Here, the binary DICE assesses the ability of the method to distinguish nuclei from the background, whereas PQ provides insight into the quality of instance segmentation.

D. Comparative analysis of rotation-equivariant models

Baseline models: For the task of breast tumour classification, we implement a baseline CNN for comparison with the aforementioned rotation-equivariant models. The model consists of a series of convolution, batch normalisation, non-linear and spatial pooling operations, which are then followed by three 1×1 convolutions to obtain the final output, denoting the probability of an input patch being tumour.

For the tasks of gland and nuclear segmentation we leverage the fully convolutional neural network architecture, which allows us to use the same model architecture, irrespective of the input size. The encoder of the baseline segmentation model uses the same architecture as the baseline classification CNN. Then a series of up-sampling and convolution operations are used to regain the spatial dimensions of the original image. In line with U-Net, we use skip connections to incorporate features from the encoder, but utilise summation as opposed to concatenation. In line with our proposed model described in Subsection III-A, at the output of the network we perform segmentation of the object and the contour and additionally predict the eroded masks for nuclear segmentation.

Rotation-equivariant models: To assess the performance of various rotation-equivariant approaches, we modify the baseline models, but keep the fundamental architecture the same. The main difference between different models is how the

filters are rotated, how many filter orientations are considered and how the convolution operation is performed.

Aside from H-Nets, each rotation-equivariant model considers 4, 8 and 12 filter orientations. H-Nets encode full 360° equivariance within the model and therefore filters do not need to be explicitly rotated. When applying rotation to a filter with an angle that is a multiple of $\frac{\pi}{2}$, the rotation is *exact* because the output can still be represented on the square grid. However, any other rotation may give interpolation artefacts and therefore may have negative implications for rotation-equivariance. Therefore, in line with Marcos *et al.* [4] and Lafarge *et al.* [5], for both the VF-CNN and standard G -CNN, we apply circular masking to the filters when using the groups C_8 and C_{12} . However, this masking still leads to inevitable interpolation artefacts in the centre of the filter. Steerable filters as defined by (2) do not suffer from interpolation artefacts and, therefore, circular masking is not needed.

In all comparative experiments for rotation-equivariance, we fix each filter to be of size 7×7 . We used a larger filter than typically used in modern CNNs because this size ensures that we can construct a good basis set for steerable filter generation, with reasonable frequency content and reduced aliasing.

For fair comparison, we ensure that the number of parameters is similar between different models. For both standard and steerable G -CNNs, the number of parameters increases with the size of the group, if we fix the number of filters in each layer. This is because one feature map is produced per orientation of the filter, which increases the number of required filters in the subsequent layer. To maintain the same number of parameters as the baseline CNN, we divide the number of filters in each layer of the standard G -CNN by \sqrt{n} , where n is the number of orientations in the group. Steerable G -CNNs learn k parameters (or $k/2$ complex parameters) for each filter, where typically $k < K^2$. Therefore, the number of filters in each layer of a steerable G -CNN should be divided by $\frac{k\sqrt{n}}{K^2}$. Instead of carrying forward all orientations throughout the network, VF-CNNs collapse the orientation dependent feature maps to two feature maps, representing magnitude and angle. Therefore, the VF-CNN requires more filters in the next layer, but the number of parameters stays constant irrespective of the size of the group. To ensure the same number of parameters as the baseline CNN, for all group sizes we divide the number of filters in each layer of VF-CNNs by $\frac{4}{3}$. Each H-Net filter is constrained to be a complex circular harmonic, parameterised by N radial terms and a single phase offset term. Also, the number of parameters is dependent on the maximum frequency m of the filters. Specifically, in H-Nets frequencies in the range $[-m, m]$ are considered, equating to a total of $M = 2m + 1$ frequency terms. Therefore, to ensure a similar number of parameters as the standard CNN, we multiply the number of filters in each layer of a H-Net by $\frac{K^2}{M \cdot (N+1)}$.

In all models, we down-sample with max-pooling, but for VF-CNNs and H-Nets we use a modified pooling strategy, based on the magnitude of the feature maps. Similarly, when using both VF-CNNs and H-Nets, we do not incorporate the angle information when using batch normalisation (BN) and non-linear activation functions; otherwise the angles may change important information about relative and global orien-

tations. For G -CNNs, we use a modified BN that aggregates moments per group rather than spatial feature map.

To verify our implementations of the various rotation-equivariant networks, we cross-checked the performance of each model against reported benchmarks on the rotated MNIST dataset [41] before applying them to the histology datasets. These results are summarised in Table A2.

E. Quantitative results

Tumour classification: We report comparative results of different rotation-equivariant models on the PCam dataset at the top of Table I. We observe that H-Nets do not perform as well as the baseline CNN for the task of tumour classification. Despite this, we observe that we are able to increase the performance when incorporating higher frequency filters in the network, but the performance is still not comparable to conventional CNNs. This may suggest that constraining the filters in this way may not be optimal for detecting complex features in histology. VF-CNNs marginally outperform the conventional CNN, where we observe that increasing the number of filter rotations leads to a slight improvement in performance. When we utilise the group convolution, with filter rotation as performed by Bekkers *et al.* [31] and Lafarge *et al.* [5], we see an improved performance when using up to 8 filter orientations. This gain in performance can be attributed to incorporating our prior knowledge of rotational symmetry into the network. To ensure that we maintain a similar number of parameters, we need to reduce the number of feature maps at each layer when the size of the group is increased. This may explain the drop in performance when using 12 filter orientations. When using steerable filters, but with no filter rotation, we observe an improved performance over conventional CNNs, highlighting the benefit of learning a linear combination of basis filters, rather than standard filters. Then, as we increase the size of the group to 4 and 8 orientations we see an improvement in the performance. We also observe that using steerable filters rather than standard filters within the G -convolution gives a better result.

At the bottom of Table I we compare the performance of our proposed DSF-CNN with the $p4m$ -DenseNet [27], which is the top performing method that was proposed with the introduction of the PCam dataset. This approach integrates the use of G -convolutions on, as proposed by Cohen & Welling [2], into a densely connected CNN [7]. Here, the network uses filter rotations by multiples of 90° and also uses reflections. This is denoted by D_4 , which is the dihedral group containing 4 rotation and 4 reflection symmetries. In addition, we compare results to the commonly used ResNet-34 [11], ResNet-50 [11], DenseNet-121 [7] and DenseNet-169 [7]. Despite the small amount of parameters, we observe that our method achieves the best performance with an AUC of 0.975, which is a promising improvement over the previous state-of-the-art.

Gland segmentation: We compare the performance of the different rotation-equivariant models for gland segmentation on the CRAG dataset in the top part of Table II. For this experiment, when comparing different rotation-equivariant ap-

proaches, we choose to only assess the performance of conventional CNNs, standard G -CNNs and steerable G -CNNs. This is because our previous experiment on breast tumour classification indicates that G -CNNs are capable of achieving a superior result over competing rotation-equivariant approaches. Similar to our observations for breast tumour classification, we see that increasing the group size within the group convolution leads to an increase in performance, but the best performance is achieved when using 8 filter orientations. For this task, using steerable filters in the group convolution led to the best performance.

TABLE I

TUMOUR CLASSIFICATION RESULTS ON THE PCAM DATASET [27]. TOP: COMPARISON OF DIFFERENT ROTATION-EQUIVARIANT MODELS WITH A SIMILAR PARAMETER BUDGET. BOTTOM: COMPARISON OF PROPOSED APPROACH WITH STATE-OF-THE-ART. THE SUPERScript ASSOCIATED WITH H-NET DENOTES THE MAXIMUM FREQUENCY USED.

Method	Group	Parameters	AUC
CNN	{e}	564K	0.947
H-Net ¹ [34]	SO(2)	553K	0.934
H-Net ² [34]	SO(2)	542K	0.939
VF-CNN [4]	C_4	556K	0.949
VF-CNN [4]	C_8	556K	0.951
VF-CNN [4]	C_{12}	556K	0.953
G -CNN [2]	C_4	561K	0.964
G -CNN [5], [31]	C_8	557K	0.968
G -CNN [5], [31]	C_{12}	557K	0.962
Steerable G -CNN [3]	{e}	553K	0.963
Steerable G -CNN [3]	C_4	546K	0.969
Steerable G -CNN [3]	C_8	565K	0.971
Steerable G -CNN [3]	C_{12}	545K	0.969
ResNet-34 [11]	{e}	21.3M	0.942
ResNet-50 [11]	{e}	23.5M	0.948
DenseNet-121 [7]	{e}	7.8M	0.921
DenseNet-169 [7]	{e}	13.3M	0.920
$p4m$ -DenseNet* [27]	D_4	119K	0.963
DSF-CNN (Ours)	C_8	2.2M	0.975

TABLE II

GLAND SEGMENTATION RESULTS ON THE CRAG [14] DATASET. TOP: COMPARISON OF DIFFERENT ROTATION-EQUIVARIANT MODELS WITH A SIMILAR PARAMETER BUDGET. BOTTOM: COMPARISON OF PROPOSED APPROACH WITH STATE-OF-THE-ART.

Method	Group	Params	Obj F ₁	Obj Dice	Obj Haus ↓
CNN	{e}	984K	0.793	0.809	246.0
G -CNN [2]	C_4	982K	0.833	0.856	170.4
G -CNN [5], [31]	C_8	988K	0.837	0.866	157.4
G -CNN [5], [31]	C_{12}	979K	0.818	0.834	192.2
Steerable G -CNN [3]	{e}	981K	0.811	0.848	175.9
Steerable G -CNN [3]	C_4	984K	0.837	0.869	164.8
Steerable G -CNN [3]	C_8	989K	0.861	0.888	139.5
Steerable G -CNN [3]	C_{12}	976K	0.855	0.870	156.2
FCN8 [38]	{e}	134.3M	0.796	0.835	199.5
U-Net [38]	{e}	37.0M	0.827	0.844	196.9
MILD-Net [14]	{e}	83.3M	0.869	0.883	146.2
Rota-Net [29]	C_4	71.3M	0.869	0.887	144.2
DSF-CNN (Ours)	C_8	3.7M	0.874	0.891	138.4

In the bottom part of Table II, we compare our proposed approach with MILD-Net [14] and Rota-Net [29], which are top-performing gland segmentation methods and therefore can be appropriately used for performance benchmarking. Like

TABLE III

NUCLEAR SEGMENTATION RESULTS ON THE KUMAR [18] DATASET. TOP: COMPARISON OF DIFFERENT ROTATION-EQUIVARIANT MODELS WITH A SIMILAR PARAMETER BUDGET. BOTTOM: COMPARISON OF PROPOSED APPROACH WITH STATE-OF-THE-ART.

Method	Group	Params	B-Dice	PQ
CNN	{e}	984K	0.767	0.447
G-CNN [2]	C_4	982K	0.793	0.490
G-CNN [5], [31]	C_8	988K	0.811	0.519
G-CNN [5], [31]	C_{12}	979K	0.814	0.534
Steerable G-CNN [3]	{e}	981K	0.791	0.510
Steerable G-CNN [3]	C_4	984K	0.809	0.542
Steerable G-CNN [3]	C_8	989K	0.818	0.543
Steerable G-CNN [3]	C_{12}	976K	0.820	0.558
FCN8 [42]	{e}	134.3M	0.797	0.312
SegNet [43]	{e}	29.4M	0.811	0.407
U-Net [38]	{e}	37.0M	0.758	0.478
Mask-RCNN [44]	{e}	40.1M	0.760	0.509
DIST [17]	{e}	9.2M	0.789	0.443
Micro-Net [45]	{e}	192.6M	0.797	0.519
CIA-Net [46]	{e}	22.0M	0.818	0.577
HoVer-Net [16]	{e}	54.7M	0.826	0.597
DSF-CNN (Ours)	C_8	3.7M	0.826	0.600

the $p4m$ -DenseNet, Rota-Net makes use of the standard G -convolution, but is limited to only 90° filter rotations. In addition, we compare with FCN8 and U-Net as they are two widely used CNNs for segmentation. We observe that our DSF-CNN achieves the best performance with a fraction of the parameter budget. Notably, our model has around 20 times fewer parameters than Rota-Net and MILD-Net.

Nuclear segmentation: We report the comparative results of different rotation-equivariance methods for nuclear segmentation on the Kumar dataset in the top part of Table III. Similar to above, we compare conventional CNNs with both standard and steerable G -CNNs. Here, we see that all rotation-equivariant approaches show a significant improvement over standard CNNs and we see an improvement when increasing the number of filter orientations to 12 in all models. Once again, we observe that the steerable G -CNNs for segmentation of nuclei are superior to standard G -CNNs that use bilinear interpolation during filter rotation.

We evaluate the performance of our proposed method with several state-of-the-art approaches in the bottom part of Table III. In particular, HoVer-Net [16], CIA-Net [46], Micro-Net [45] and DIST [17] have been purpose-built for the task of nuclear segmentation and, therefore, provide a competitive benchmark. The proposed DSF-CNN once again achieves the best performance compared to other methods for both binary DICE and panoptic quality, on par with the state-of-the-art HoVer-Net method, while requiring a fraction of the parameter count.

F. Visual results

In Fig. 5 we visualise the features and the corresponding outputs as we rotate the input with angle increments of $\frac{\pi}{4}$ (8 in total) for both the baseline CNN and C_8 -steerable G-CNN. Specifically, we analyse the properties of both CNNs trained for the tasks of gland and nuclear segmentation. To observe the feature map transformation with rotation of the input, we

analyse two sets of feature maps in both CNNs: *Feature Map A* at the output of the 2nd convolution and *Feature Map B* at the output of the convolution after the final up-sampling operation. Similarly, we observe how the output probability map transforms when the input is rotated.

To analyse this, we feed each image orientation into the network to obtain a set of feature maps and output probability maps. Then, after rotating features and probability maps back to their original orientation, we compute the pixel-wise variance map of the features and the output to see how they change with rotation of the input. G -CNN feature maps are a function on G and therefore we visualise a single planar feature map within the group. For the rotation-equivariant model, we observe that there is a near-negligible variance between the features of each input orientation. On the other hand, there is much higher variance between the features of standard CNNs after input rotation. This implies that the rotation-equivariant CNN successfully learns an equivariant feature representation. Also, there is a lower variance between the predictions of multiple input orientations for the rotation-equivariant CNN as compared to the standard CNN. Thus, the rotation-equivariant CNN behaves as expected with rotation of the input, which is a particularly desirable property when training CNNs with histology image data. It must be noted that features learned by conventional CNNs are highly complex and it is very difficult to infer the relationship between learned features and input rotation. Nonetheless, we demonstrate that rotation-equivariant CNNs have a predictable transformation with input rotation, making them more stable than conventional CNNs.

G. Implementation and training details

We implemented our framework with the open source software library TensorFlow version 1.12.0 [47] on a workstation equipped with two NVIDIA GeForce 1080 Ti GPUs. During training, data augmentation including flip, rotation, Gaussian blur and median blur was applied. For breast tumour classification, we fed the original patches of size 96×96 into the network. For gland and nuclear segmentation, we used patches of size 448×448 and 256×256 respectively. For tumour classification, we trained our model using a batch size of 32 and then used a batch size of 8 for both gland and nucleus segmentation. We used cross-entropy loss for all tumour classification and gland segmentation models, whereas we used a combination of weighted cross-entropy and dice loss for nuclear segmentation. For all models, we trained using Adam optimisation with an initial learning rate of 10^{-3} , that was reduced as training progressed. The network was trained with an RGB input, normalised between 0 and 1.

V. DISCUSSION AND CONCLUSIONS

Conventional CNNs do not behave as expected with rotation of the input, which is a particularly undesirable property in the field of computational pathology, where important features in histology images can appear at any orientation. Instead, rotation-equivariant CNNs build this prior knowledge of rotational symmetry within the network, such that features rotate in accordance with the input without explicitly learning features

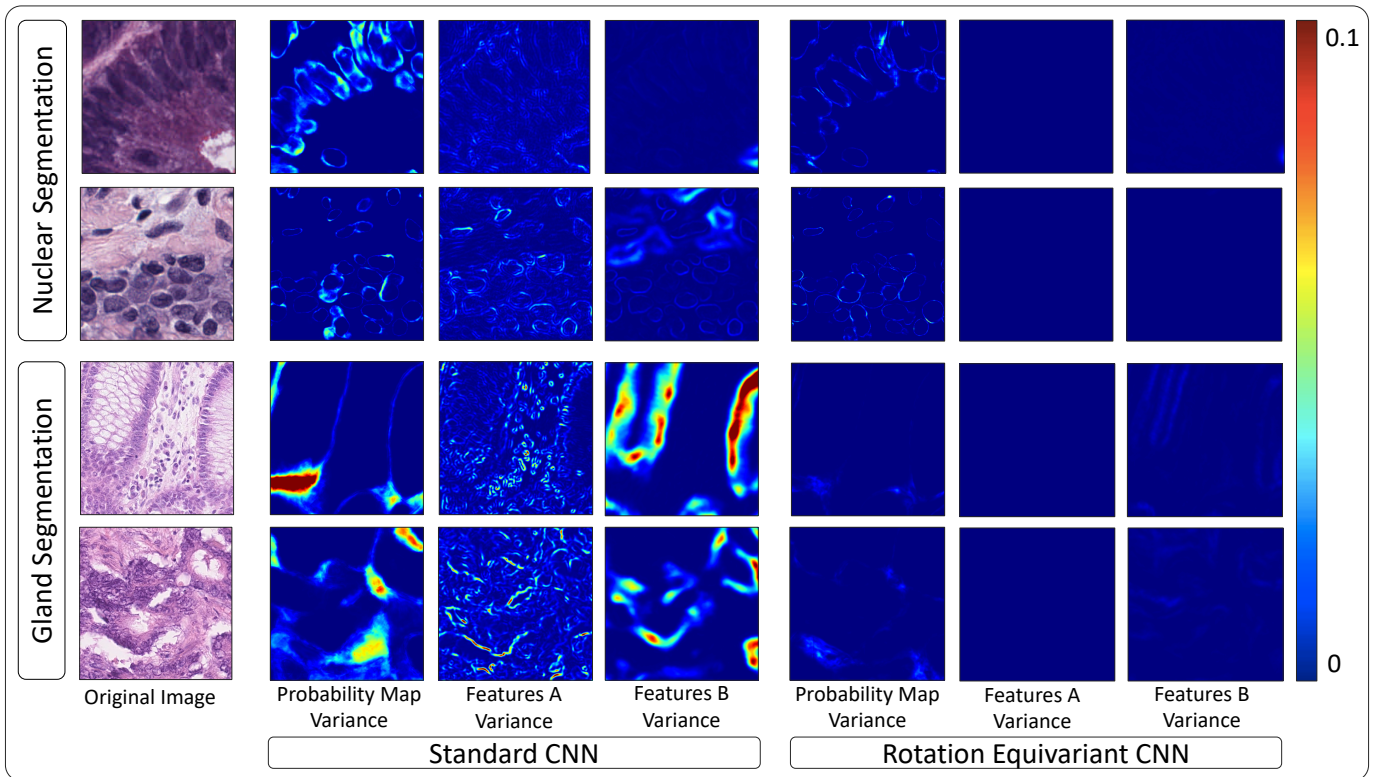


Fig. 5. Variance between the predictions and features for multiple orientations of the input. The original image is rotated with steps of $\frac{\pi}{4}$ to give 8 orientations and each copy is passed through the network to enable variance calculation. Features A and B are located at the beginning and end of the network respectively. The rotation-equivariant CNN we compare with is the C_8 steerable G -CNN.

at various orientations. In this paper, we proposed a densely connected steerable filter CNN that achieves state-of-the-art performance on the three datasets used in our experiments with a fraction of the parameter budget of recent top-performing models. We conducted a thorough comparative analysis of various rotation-equivariant CNNs applied to the tasks of breast tumour classification, gland segmentation and nuclear segmentation. We showed that steerable filter group convolutions gave the best quantitative results on all three tasks, where 8 filter orientations consistently gave a strong performance. We visualised features within a rotation-equivariant model to demonstrate that they rotate with the input and therefore have a higher degree of feature map interpretability. Finally, we showed that rotation-equivariant models give more stable predictions with input rotation than regular CNNs do. In future work, we will consider incorporating additional symmetries into the group convolution, such as mirror and scale symmetries. This will further increase the interpretability of feature maps and may lead to an improvement in performance and help direct future research in computational pathology.

APPENDIX

A. VERIFICATION OF BASELINE MODELS

In order to verify our self implemented approaches, we report the performance of each rotation-equivariant model on the rotated MNIST dataset [41] in Table A2, which is typically used for performance benchmarking in this domain. In particular, we report the performance of a conventional

CNN, H-Nets, standard G -CNNs, VF-CNNs and steerable G -CNNs. This was primarily to ensure that we were able to achieve a comparable performance with the reported results in the original papers. In our experiments all CNNs have the same base-level architecture, where we ensured that the models had the same number of layers, the same filter size and a similar number of parameters. Therefore our experiments are not only used for verification, but also to perform a fair head-to-head comparison between models. To maintain a similar number of parameters, we followed the same strategy as described in Section IV-D. In line with our experiments in the paper, for H-Net we apply spatial max-pooling based on the magnitudes, as opposed to average-pooling, which is used in the original paper.

TABLE A1

PERFORMANCE OF OUR BASELINE MODELS ON ROTATED MNIST DATASET [41]. THE SUPERScript ASSOCIATED WITH H-NET DENOTES THE MAXIMUM FREQUENCY USED.

Method	Group	Parameters	Error
CNN	$\{e\}$	416K	2.001
H-Net ¹ [34]	SO(2)	418K	1.371
H-Net ² [34]	SO(2)	414K	1.352
G -CNN [2]	C_4	413K	0.976
G -CNN [5], [31]	C_8	407K	0.962
G -CNN [5], [31]	C_{12}	411K	0.940
VF-CNN [4]	C_8	418K	1.202
VF-CNN [4]	C_{12}	418K	1.172
Steerable G -CNN [3]	C_8	416K	0.820
Steerable G -CNN [3]	C_{12}	424K	0.809

We observe that all rotation-equivariant CNNs achieve a greater performance than the conventional CNN, where the best performance is achieved by the C_{12} steerable G -CNN. Interestingly, we observe a significant boost in performance for our C_4 G -CNN and H-Net implementations, compared to the originally published results. These models have the same number of layers as the original implementations, but are wider to ensure a similar number of parameters between competing models. Note, we also add $2 \times 1 \times 1$ convolutions after obtaining the invariant map (after G -pooling or computing the magnitude of the complex feature maps), which may have also contributed to the increase in performance. If we use the same architecture used by Weiler *et al.* for the C_{12} steerable G -CNN, then we obtain an error of 0.709, which is very close to the original result. However, this implementation uses around $3.3M$ parameters, which is nearly $8 \times$ the amount that we use in our comparative experiments in Table A2.

TABLE A2

DESCRIPTION OF MATHEMATICAL SYMBOLS USED THROUGHOUT THE PAPER.

Symbol	Description
\mathbb{R}	Set of real numbers
\mathbb{C}	Set of complex numbers
\mathbb{Z}	Set of integers
\mathcal{F}	Real vector space of functions $\mathbb{C} \rightarrow \mathbb{R}$
$\mathcal{F}_{\mathbb{C}}$	Complex vector space of functions $\mathbb{C} \rightarrow \mathbb{C}$
Re	Real part of complex number
$E(2)$	Euclidean group
$SE(2)$	Special euclidean group (no reflections)
$SO(2)$	Special orthogonal group (no reflections)
$\{e\}$	Trivial group containing only the identity on page 3
n	A positive integer, fixed throughout this paper
D_n	Dihedral group of n rotations of \mathbb{C} , fixing 0 and flips
C_n	Cyclic group of n rotations of \mathbb{C} , fixing 0
C'_n	$\{2\pi s/n \mid 0 \leq s < n\}$ group law is addition mod 2π
\mathcal{G}	An arbitrary group
G	Group as defined in Subsection II-C
r	radius in polar coordinates
ψ	a filter
λ, β, θ	usually elements of C'_n , sometimes arbitrary angles
R_k	Radial profile of atomic steerable filters

REFERENCES

- [1] D. R. Snead, Y.-W. Tsang, A. Meskiri, P. K. Kimani, R. Crossman, N. M. Rajpoot, E. Blessing, K. Chen, K. Gopalakrishnan, P. Matthews *et al.*, "Validation of digital pathology imaging for primary histopathological diagnosis," *Histopathology*, vol. 68, no. 7, pp. 1063–1072, 2016.
- [2] T. Cohen and M. Welling, "Group equivariant convolutional networks," in *International conference on machine learning*, 2016, pp. 2990–2999.
- [3] M. Weiler, F. A. Hamprecht, and M. Storath, "Learning steerable filters for rotation equivariant cnns," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018.
- [4] D. Marcos, M. Volpi, N. Komodakis, and D. Tuia, "Rotation equivariant vector field networks," in *Proceedings of the IEEE International Conference on Computer Vision*, 2017, pp. 5048–5057.
- [5] M. W. Lafarge, E. J. Bekkers, J. P. Pluim, R. Duits, and M. Veta, "Roto-translation equivariant convolutional networks: Application to histopathology image analysis," *arXiv preprint arXiv:2002.08725*, 2020.
- [6] W. T. Freeman and E. H. Adelson, "The design and use of steerable filters," *IEEE Transactions on Pattern Analysis & Machine Intelligence*, no. 9, pp. 891–906, 1991.
- [7] G. Huang, Z. Liu, L. van der Maaten, and K. Q. Weinberger, "Densely Connected Convolutional Networks," *ArXiv e-prints*, p. arXiv:1608.06993, Aug. 2016.
- [8] W. Ke, J. Chen, J. Jiao, G. Zhao, and Q. Ye, "Srn: side-output residual network for object symmetry detection in the wild," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 1068–1076.
- [9] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner, "Gradient-based learning applied to document recognition," *Proceedings of the IEEE*, vol. 86, no. 11, pp. 2278–2324, 1998.
- [10] S. Ioffe and C. Szegedy, "Batch normalization: Accelerating deep network training by reducing internal covariate shift," *arXiv preprint arXiv:1502.03167*, 2015.
- [11] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.
- [12] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "ImageNet: A Large-Scale Hierarchical Image Database," in *CVPR09*, 2009.
- [13] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick, "Microsoft coco: Common objects in context," in *European conference on computer vision*. Springer, 2014, pp. 740–755.
- [14] S. Graham, H. Chen, J. Gamper, Q. Dou, P.-A. Heng, D. Snead, Y. W. Tsang, and N. Rajpoot, "Mild-net: Minimal information loss dilated network for gland instance segmentation in colon histology images," *Medical image analysis*, vol. 52, pp. 199–211, 2019.
- [15] H. Chen, X. Qi, L. Yu, and P.-A. Heng, "Dcan: deep contour-aware networks for accurate gland segmentation," in *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, 2016, pp. 2487–2496.
- [16] S. Graham, Q. D. Vu, S. E. A. Raza, A. Azam, Y. W. Tsang, J. T. Kwak, and N. Rajpoot, "Hover-net: Simultaneous segmentation and classification of nuclei in multi-tissue histology images," *Medical Image Analysis*, vol. 58, p. 101563, 2019.
- [17] P. Naylor, M. Laé, F. Reyat, and T. Walter, "Segmentation of nuclei in histopathology images by deep regression of the distance map," *IEEE transactions on medical imaging*, vol. 38, no. 2, pp. 448–459, 2018.
- [18] N. Kumar, R. Verma, S. Sharma, S. Bhargava, A. Vahadane, and A. Sethi, "A dataset and a technique for generalized nuclear segmentation for computational pathology," *IEEE transactions on medical imaging*, vol. 36, no. 7, pp. 1550–1560, 2017.
- [19] S. U. Akram, T. Qaiser, S. Graham, J. Kannala, J. Heikkilä, and N. Rajpoot, "Leveraging unlabeled whole-slide-images for mitosis detection," in *Computational Pathology and Ophthalmic Medical Image Analysis*. Springer, 2018, pp. 69–77.
- [20] S. Graham, M. Shaban, T. Qaiser, N. A. Koohbanani, S. A. Khurram, and N. Rajpoot, "Classification of lung cancer histology images using patch-level summary statistics," in *Medical Imaging 2018: Digital Pathology*, vol. 10581. International Society for Optics and Photonics, 2018, p. 1058119.
- [21] E. Arvaniti, K. S. Fricker, M. Moret, N. Rupp, T. Hermanns, C. Fankhauser, N. Wey, P. J. Wild, J. H. Rueschoff, and M. Claassen, "Automated gleason grading of prostate cancer tissue microarrays via deep learning," *Scientific reports*, vol. 8, no. 1, pp. 1–11, 2018.
- [22] M. Shaban, R. Awan, M. M. Fraz, A. Azam, Y. Tsang, D. Snead, and N. M. Rajpoot, "Context-aware convolutional neural network for grading of colorectal cancer histology images," *IEEE Transactions on Medical Imaging*, pp. 1–1, 2020.
- [23] D. Tellez, G. Litjens, P. Bandi, W. Bulten, J.-M. Bokhorst, F. Ciompi, and J. van der Laak, "Quantifying the effects of data augmentation and stain color normalization in convolutional neural networks for computational pathology," *arXiv preprint arXiv:1902.06543*, 2019.
- [24] A. Azulay and Y. Weiss, "Why do deep convolutional networks generalize so poorly to small image transformations?" *arXiv preprint arXiv:1805.12177*, 2018.
- [25] N. Moshkov, B. Mathe, A. Kertesz-Farkas, R. Hollandi, and P. Horvath, "Test-time augmentation for deep learning-based cell segmentation on microscopy images," *bioRxiv*, p. 814962, 2019.
- [26] D. Laptev, N. Savinov, J. M. Buhmann, and M. Pollefeys, "Ti-pooling: transformation-invariant pooling for feature learning in convolutional neural networks," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 289–297.
- [27] B. S. Veeling, J. Linmans, J. Winkens, T. Cohen, and M. Welling, "Rotation equivariant cnns for digital pathology," in *International Conference on Medical image computing and computer-assisted intervention*. Springer, 2018, pp. 210–218.
- [28] J. Linmans, J. Winkens, B. S. Veeling, T. S. Cohen, and M. Welling, "Sample efficient semantic segmentation using rotation equivariant convolutional networks," *arXiv preprint arXiv:1807.00583*, 2018.

- [29] S. Graham, D. Epstein, and N. Rajpoot, "Rota-net: Rotation equivariant network for simultaneous gland and lumen segmentation in colon histology images," in *European Congress on Digital Pathology*. Springer, 2019, pp. 109–116.
- [30] E. Hoogeboom, J. W. Peters, T. S. Cohen, and M. Welling, "Hexaconv," *arXiv preprint arXiv:1803.02108*, 2018.
- [31] E. J. Bekkers, M. W. Lafarge, M. Veta, K. A. Eppenhof, J. P. Pluim, and R. Duits, "Roto-translation covariant convolutional networks for medical image analysis," in *International Conference on Medical Image Computing and Computer-Assisted Intervention*. Springer, 2018, pp. 440–448.
- [32] Y. Zhou, Q. Ye, Q. Qiu, and J. Jiao, "Oriented response networks," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 519–528.
- [33] T. S. Cohen and M. Welling, "Steerable cnns," *arXiv preprint arXiv:1612.08498*, 2016.
- [34] D. E. Worrall, S. J. Garbin, D. Turmukhambetov, and G. J. Brostow, "Harmonic networks: Deep translation and rotation equivariance," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 5028–5037.
- [35] X. Cheng, Q. Qiu, R. Calderbank, and G. Sapiro, "Rotdcf: Decomposition of convolutional filters for rotation-equivariant deep networks," in *International Conference on Learning Representations 2019 (ICLR'19)*, 2019.
- [36] M. Weiler and G. Cesa, "General e (2)-equivariant steerable cnns," in *Advances in Neural Information Processing Systems*, 2019, pp. 14 334–14 345.
- [37] student (<https://math.stackexchange.com/users/20150/student>), "Every measurable homomorphism from \mathbb{R}^n to \mathbb{C}^* is exponential." Mathematics Stack Exchange, uRL:<https://math.stackexchange.com/q/442980> (version: 2013-07-13). [Online]. Available: <https://math.stackexchange.com/q/442980>
- [38] O. Ronneberger, P. Fischer, and T. Brox, "U-net: Convolutional networks for biomedical image segmentation," in *International Conference on Medical image computing and computer-assisted intervention*. Springer, 2015, pp. 234–241.
- [39] B. E. Bejnordi, M. Veta, P. J. Van Diest, B. Van Ginneken, N. Karssemeijer, G. Litjens, J. A. Van Der Laak, M. Hermsen, Q. F. Manson, M. Balkenhol *et al.*, "Diagnostic assessment of deep learning algorithms for detection of lymph node metastases in women with breast cancer," *Jama*, vol. 318, no. 22, pp. 2199–2210, 2017.
- [40] K. Sirinukunwattana, J. P. Pluim, H. Chen, X. Qi, P.-A. Heng, Y. B. Guo, L. Y. Wang, B. J. Matuszewski, E. Bruni, U. Sanchez *et al.*, "Gland segmentation in colon histology images: The glas challenge contest," *Medical image analysis*, vol. 35, pp. 489–502, 2017.
- [41] H. Larochelle, D. Erhan, A. Courville, J. Bergstra, and Y. Bengio, "An empirical evaluation of deep architectures on problems with many factors of variation," in *Proceedings of the 24th international conference on Machine learning*, 2007, pp. 473–480.
- [42] J. Long, E. Shelhamer, and T. Darrell, "Fully convolutional networks for semantic segmentation," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 3431–3440.
- [43] V. Badrinarayanan, A. Kendall, and R. Cipolla, "Segnet: A deep convolutional encoder-decoder architecture for image segmentation," *IEEE transactions on pattern analysis and machine intelligence*, vol. 39, no. 12, pp. 2481–2495, 2017.
- [44] K. He, G. Gkioxari, P. Dollár, and R. Girshick, "Mask R-CNN," *ArXiv e-prints*, p. arXiv:1703.06870, Mar. 2017.
- [45] S. E. A. Raza, L. Cheung, M. Shaban, S. Graham, D. Epstein, S. Pelenaris, M. Khan, and N. M. Rajpoot, "Micro-net: A unified model for segmentation of various objects in microscopy images," *Medical image analysis*, vol. 52, pp. 160–173, 2019.
- [46] Y. Zhou, O. F. Onder, Q. Dou, E. Tsougenis, H. Chen, and P.-A. Heng, "Cia-net: Robust nuclei instance segmentation with contour-aware information aggregation," *arXiv preprint arXiv:1903.05358*, 2019.
- [47] M. Abadi, P. Barham, J. Chen, Z. Chen, A. Davis, J. Dean, M. Devin, S. Ghemawat, G. Irving, M. Isard *et al.*, "Tensorflow: A system for large-scale machine learning." in *OSDI*, vol. 16, 2016, pp. 265–283.