

B-cos Alignment for Inherently Interpretable CNNs and Vision Transformers

Moritz Böhle, Navdeppal Singh, Mario Fritz, and Bernt Schiele *Fellow, IEEE*

Abstract—We present a new direction for increasing the interpretability of deep neural networks (DNNs) by promoting weight-input alignment during training. For this, we propose to replace the linear transformations in DNNs by our novel B-cos transformation. As we show, a sequence (network) of such transformations induces a single linear transformation that faithfully summarises the full model computations. Moreover, the B-cos transformation is designed such that the weights align with relevant signals during optimisation. As a result, those induced linear transformations become highly interpretable and highlight task-relevant features. Importantly, the B-cos transformation is designed to be compatible with existing architectures and we show that it can easily be integrated into virtually all of the latest state of the art models for computer vision—*e.g.* ResNets, DenseNets, ConvNext models, as well as Vision Transformers—by combining the B-cos-based explanations with normalisation and attention layers, all whilst maintaining similar accuracy on ImageNet. Finally, we show that the resulting explanations are of high visual quality and perform well under quantitative interpretability metrics.

Index Terms—Interpretable Deep Learning, XAI, Convolutional Neural Networks, Vision Transformers, Interpretability

1 INTRODUCTION

WHILE deep neural networks (DNNs) are highly successful in a wide range of tasks, explaining their decisions remains an open research problem [1]. The difficulty here lies in the fact that such explanations need to faithfully summarise the internal model computations *and* present them in a human-interpretable manner. *E.g.*, it is well known that piece-wise linear (*i.e.* ReLU-based [2]) models are accurately summarised by a linear transformation for every input [3]. However, despite providing an accurate summary, these piece-wise linear transformations are generally not intuitively interpretable for humans and tend to perform poorly under quantitative interpretability metrics, cf. [4], [5]. Recent work thus aimed to improve the explanations’ interpretability, often focusing on their visual quality [6]. However, gains in the visual quality of the explanations could come at the cost of their model-faithfulness [6].

Instead of optimising the explanation method, we aim to design the DNNs to inherently provide an explanation that fulfills the aforementioned requirements—the resulting explanations constitute both a faithful summary and have a clear interpretation for humans. To achieve this, we propose the **B-cos transformation** as a drop-in replacement for linear transformations. As such, the B-cos transformation can easily be integrated into a wide range of existing DNN architectures and we show that the resulting models provide high-quality explanations for their decisions, see Fig. 1.

To ensure that these explanations constitute a faithful summary of the models, we design the B-cos transformation

as an input-dependent linear transformation. Importantly, any sequence of such transformations therefore induces a single linear transformation that faithfully summarises the entire sequence. To make the induced linear transformations interpretable, the B-cos transformation is designed such that its weights align with task-relevant input signals during optimisation. The summarising linear transformation thus becomes easily interpretable for humans: it is a direct reflection of the weights the model has learnt during training and specifically reflects those weights that best align with a given input. In short, we make the following contributions:
(1) We introduce the B-cos transformation to improve the interpretability of DNNs. By promoting weight-input alignment, these transformations are explicitly designed to yield explanations that highlight task-relevant input patterns.
(2) Specifically, the B-cos transformation is designed such that any sequence of B-cos transformations can be faithfully summarised by a single linear transformation. We show that this allows to explain not only the models’ outputs, but also representations in intermediate network layers.
(3) We demonstrate that a plain B-cos convolutional neural network without any additional non-linearities, normalisation layers or regularisation schemes achieves competitive performance on CIFAR10 [7], demonstrating the modelling capacity of DNNs that are solely based on the B-cos transformation. In this context, we show that the parameter B gives fine-grained control over the increase in weight alignment and thus the interpretability of the B-cos networks.
(4) We analyse how to integrate normalisation layers into B-cos models to take advantage of the optimisation benefits of those layers without sacrificing interpretability.
(5) We analyse how to combine the proposed B-cos framework with attention-based models (*e.g.* Vision Transformers (ViTs) [8]), which highlights the generality of our approach and shows that it extends beyond convolutional networks.
(6) Finally, we show in a wide range of experiments that B-cos DNNs achieve similar accuracy as their conven-

• Moritz Böhle is the corresponding author. Moritz Böhle, Navdeppal Singh, and Bernt Schiele are with the Department of Computer Vision and Machine Learning, Max Planck Institute for Informatics, Saarland Informatics Campus, Saarbrücken, 66123, Germany. E-mail: {mboehle,nsingh,schiele}@mpi-inf.mpg.de

• Mario Fritz is with the CISPA Helmholtz Center for Information Security, Saarbrücken 66123, Germany. E-mail: fritz@cispa.de.

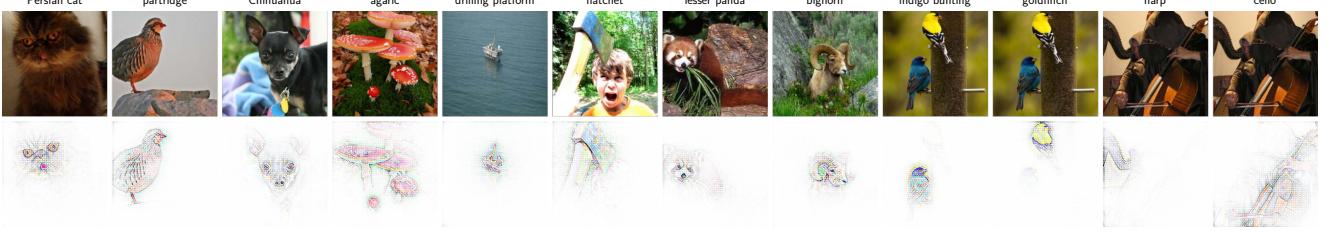


Fig. 1: **Top:** Inputs \mathbf{x}_i to a B-cos DenseNet-121. **Bottom:** B-cos explanation for class c (c : image label). Specifically, we visualise the c -th row of $\mathbf{W}_{1 \rightarrow L}(\mathbf{x}_i)$ as applied by the model, see Eq. (14); no masking of the original image is used. For the last 2 images, we also show the explanation for the 2nd most likely class. For details on visualising $\mathbf{W}_{1 \rightarrow L}(\mathbf{x}_i)$, see Sec. 4.

tional counterparts. Specifically, we evaluate the B-cos versions of various commonly used DNNs such as VGGs [9], ResNets [10], DenseNets [11], ResNeXt [12] and ConvNeXt models [13], as well as ViTs [8]. Our strongest models achieve >80% top-1 accuracy on ImageNet, all while providing highly detailed explanations for their decisions. Specifically, the B-cos explanations outperform other explanation methods across all tested architectures, both under quantitative metrics as well as under qualitative inspection.

As such, this work extends [14] by integrating normalisation and attention layers into B-cos DNNs and provides a significantly broader experimental evaluation of the proposed framework. The code and all of the trained models are publicly available at github.com/B-cos/B-cos-v2.

2 RELATED WORK

Approaches for understanding DNNs typically focus on explaining individual model decisions *post-hoc*, i.e., they are designed to work on any pre-trained DNN. Examples of this include perturbation-based, [15], [16], [17], activation-based, [18], [19], or backpropagation-based explanations, [5], [20], [21], [22], [23], [24], [25], [26]. In order to obtain explanations for the B-cos networks, we also rely on a backpropagation-based approach. In contrast to post-hoc explanation methods, however, we design the B-cos networks to be explainable under this particular form of backpropagation and the resulting explanations are thus model-inherent.

The design of such *inherently interpretable* models has gained increased attention recently. Examples include prototype-based networks [27], BagNets [28], and CoDA Nets [29]. Similar to the BagNets and the CoDA Nets, our B-cos networks can be faithfully summarised by a single linear transformation. Moreover, similar to [29], we rely on a structurally induced alignment pressure to make those transformations interpretable. In contrast to those works, however, our method is specifically designed to be compatible with existing neural network architectures, which allows us to improve the interpretability of a wide range of DNNs.

Weight-input alignment. It has been observed that adversarial training promotes alignment [30] and recent studies suggest that this could increase interpretability via gradient-based explanations [31], [32]. Further, [33] introduce a loss to increase alignment. Instead of relying on loss-based model regularisation, the increase in alignment in B-cos networks is based on *architectural constraints* that promote weight-input alignment during model optimisation. Finally, [34] note that

the weights of a (piece-wise) linear classifier cannot be expected to align with relevant input signals in the presence of distractors (e.g. the background in object classification datasets). To overcome this, [34] propose to learn the relevant input signals for a given pre-trained model from data in order to explain said model. In contrast, by optimising the weights to align with the relevant signals already during training, for B-cos models the weights themselves can directly be used to visualise the relevant input signals.

Non-linear transformations. While the linear transformation is the default operation for most neural network architectures, many non-linear transformations have been investigated [35], [36], [37], [38], [39], [40]. Most similar to our work are [36], [37], [38], which assess transformations that emphasise the cosine similarity (i.e., ‘alignment’) between weights and inputs to improve model performance. In fact, we found that amongst other transformations, [37] evaluates a non-linear transformation that is equivalent to our B-cos operator with $B=2$. In contrast to [37], we explicitly introduce this non-linear transformation to increase model interpretability and show that the resulting models scale to large-scale classification problems.

The effect of biases on performance and interpretability. While ‘Input \times Gradient’ (IxG) [5] accurately summarises a piece-wise *linear* model given by¹ $y(\mathbf{x}) = \mathbf{W}_{\mathbf{x}}^T \mathbf{x}$, most piece-wise linear models in fact additionally contain bias terms and compute $y(\mathbf{x}) = \mathbf{W}_{\mathbf{x}}^T \mathbf{x} + \mathbf{b}(\mathbf{x})$. As bias terms are not accounted for in the explanations $\mathbf{W}_{\mathbf{x}}$, they do not provide a *full* description of the model [25], [41]. To overcome this [25], [41] propose to adapt the explanation methods to reflect the bias terms in the explanations, while [42], [43] show that, for some tasks, removing the bias terms altogether only has a minor effect on the model performance but can significantly improve the models’ interpretability. Similarly, we show that including bias terms in the models only provides minor improvements in model accuracy. Hence, we design the B-cos DNNs to be bias-free. This ensures that the linear model summaries constitute *complete* explanations [26].

Normalisation layers. To design bias-free, yet easily optimisable and performant networks, we adapt the normalisation layers in B-cos DNNs such that they do not add biases

1. Here, we write $\mathbf{W}(\mathbf{x}) = \mathbf{W}_{\mathbf{x}}$ to emphasise the fact that piece-wise linear models compute a linear transformation of the input. Thus, the gradient with respect to the input is given by the linear matrix $\mathbf{W}_{\mathbf{x}}$.

to the model computations at inference time, similar to the modified BatchNorm in [42].

Explaining attention-based models. As the name exemplifies, attention is often thought to give insight into what a model ‘pays attention to’ for its prediction. As such, various methods for using attention to understand the model output have been proposed, such as visualising the attention of single attention heads [44]. However, especially in deep layers the information can be highly distributed and it is unclear whether a given token still represents its original position in the input [45], [46], which complicates the interpretation of high attention scores deep in the neural network [45], [47].

Therefore, [46] proposed ‘attention rollout’, which summarises the various attention maps throughout the layers. However, this summary still only includes the attention layers and neglects all other network components [47]. In response, various improvements over attention rollout have been proposed, such as GradSAM [48] or an LRP-based explanation method [49], that were designed to more accurately reflect the computations of *all* model components.

Instead of providing post-hoc explanations for ViT models, in this work we show that attention layers are inherently compatible with the B-cos-based explanations, as they also compute dynamic linear transformations of their inputs. As such, they can be combined seamlessly with B-cos layers to yield B-cos ViTs, which achieve similar accuracies as their conventional counterparts. As the resulting explanations inherently comprise the multi-layer perceptrons as well as the attention layers, they faithfully summarise the entire model and thus do not require any post-hoc explanations.

3 B-COS NEURAL NETWORKS

In this section, we introduce the B-cos transformation as a replacement for the linear units in DNNs, which are (almost) “at the heart of every deep network” [50], and discuss how this can increase the interpretability of DNNs.

For this, we first introduce the B-cos transformation as a variation of the linear transformation (3.1) and highlight its most important properties. Then, we show how to construct B-cos Networks (3.2) and how to faithfully summarise the network computations to obtain explanations for their outputs (3.2.1). Subsequently, we discuss how the B-cos transformation—combined with the binary cross entropy (BCE) loss—affects the parameter optima of the models (3.2.2). Specifically, by inducing alignment pressure, the B-cos transformation aligns the model weights with task-relevant patterns in the input. Finally, in Sec. 3.3 we integrate the B-cos transformation into conventional DNNs by using it as a drop-in replacement for the ubiquitously used linear units and discuss how to combine B-cos layers with normalisation and attention layers without sacrificing the interpretability of B-cos networks.

3.1 The B-cos transformation

Typically, the individual ‘neurons’ in a DNN compute the dot product between their weights w and an input x :

$$f(x; w) = w^T x = \|w\| \|x\| c(x, w), \quad (1)$$

$$\text{with } c(x, w) = \cos(\angle(x, w)). \quad (2)$$

Here, $\angle(x, w)$ returns the angle between the vectors x and w . As we seek to improve the interpretability of DNNs by promoting weight-input alignment during optimisation, we propose to explicitly emphasise the alignment term in the linear transformation, *i.e.* the impact of the cosine function. This yields the **B-cos transformation**:

$$\text{B-cos}(x; w) = \underbrace{\|\hat{w}\|}_{=1} \|x\| |c(x, \hat{w})|^B \times \text{sgn}(c(x, \hat{w})). \quad (3)$$

Here, the hat-operator scales \hat{w} to unit norm, B is a scalar, and sgn the sign function. The B-cos transformation thus rescales the output of a linear transformation and equals:

$$\text{B-cos}(x; w) = \hat{w}^T x \times |\cos(x, \hat{w})|^{B-1}; \quad (4)$$

for a derivation of the equivalence, see supplement (Sec. D).

Note that this only introduces *minor changes* (highlighted in blue) with respect to Eq. (1); *e.g.*, for $B=1$, Eq. (4) is equivalent to a linear transformation with \hat{w} . However, albeit small, these changes are important for three reasons.

First, they bound the output of B-cos neurons:

$$\|\hat{w}\| = 1 \Rightarrow \text{B-cos}(x; w) \leq \|x\|. \quad (5)$$

As becomes clear from Eq. (3), equality in Eq. (5) is only achieved if x and w are collinear, *i.e.*, *aligned*.

Secondly, by increasing the exponent B , the output for badly aligned weight-input pairs is further suppressed,

$$B \gg 1 \wedge |c(x, \hat{w})| < 1 \Rightarrow \text{B-cos}(x; w) \ll \|x\|, \quad (6)$$

and the respective B-cos unit only yields outputs close to its maximum (*i.e.* $\|x\|$) for a small range of angular deviations from x . Together, these two properties can significantly change the optimal parameters. To illustrate this, we show in Fig. 2 how increasing B affects a simple linear classification problem. In particular, note how increasing B narrows the ‘response window’ (red area) of the B-cos transformation and suppresses the influence of the distractors (grey circles). In contrast to a linear classifier, which has a highly task-dependent optimum (cf. first row in Fig. 2), the B-cos transformation thus classifies based on *cosine similarity* between weights and inputs. As such, the optimal weights constitute ‘angular prototypes’ and lend themselves well for explaining the model prediction: a sample is confidently classified as the red class if it ‘looks like’ the weight vector.

Lastly, the B-cos transformation maintains an important property of the linear transformation: specifically, sequences of B-cos transformations are faithfully summarised by a single linear transformation (Eq. (14)). As a result, the explanations for individual neurons in a DNN can easily be combined to yield explanations for the full model, as we discuss in Sec. 3.2.1. Importantly, these explanations will align with discriminative patterns when optimising a B-cos network for classification (Sec. 3.2.2), making them easily interpretable for humans.

3.2 Simple (convolutional) B-cos networks

In this section, we discuss how to construct simple (convolutional) DNNs based on the B-cos transformation. Then, we show how to summarise the network outputs by a single

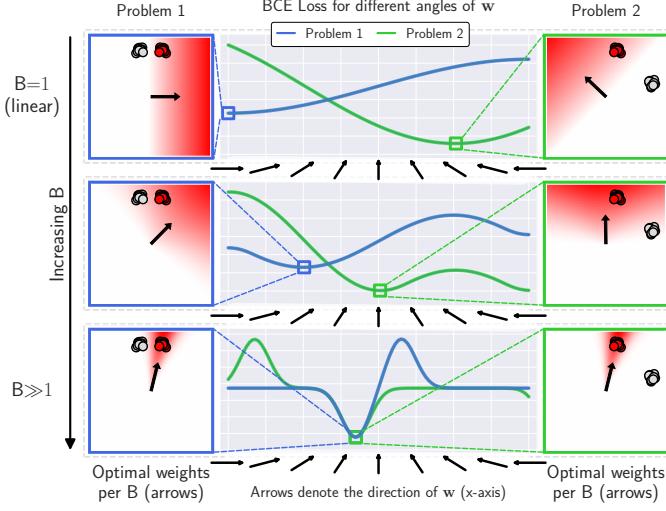


Fig. 2: **Col. 2:** BCE loss for different angles of \mathbf{w} for B-cos classifiers (Eq. (3)) with different values of B (rows) for two classification problems. **Cols. 1+3:** Visualisation of the classification problems and the corresponding optimal weights (arrows) per B . For $B=1$ (first row) the weights \mathbf{w} represent the decision boundary of a linear classifier. Although the red cluster is the same in both cases, the optimal weight vectors differ significantly (compare within row). In contrast, for higher values of B the weights converge to the same optimum in both tasks (see last row). The opacity of the red shading shows the strength of the positive activation of the B-cos transformation for a sample at a given position.

linear transformation and, finally, why this transformation aligns with discriminative patterns in classification tasks.

B-cos networks. The B-cos transformation is designed to replace the linear transformation and can be used in exactly the same way. E.g., consider a *conventional* fully connected multi-layer neural network $\mathbf{f}(\mathbf{x}; \theta)$ of L layers, defined as

$$\mathbf{f}(\mathbf{x}; \theta) = \mathbf{l}_L \circ \mathbf{l}_{L-1} \circ \dots \circ \mathbf{l}_2 \circ \mathbf{l}_1(\mathbf{x}), \quad (7)$$

with \mathbf{l}_j denoting layer j with parameters \mathbf{w}_j^k for neuron k in layer j , and θ the collection of all model parameters. In such a model, each layer \mathbf{l}_j typically computes

$$\mathbf{l}_j(\mathbf{a}_j; \mathbf{W}_j) = \phi(\mathbf{W}_j \mathbf{a}_j), \quad (8)$$

with \mathbf{a}_j the input to layer j , ϕ a non-linear activation function (e.g., ReLU), and the row k of \mathbf{W}_j given by the weight vector \mathbf{w}_j^k of the k -th neuron in that layer. Note that the non-linear activation function ϕ is *required* to be able to model non-linear relationships with multiple layers.

A corresponding B-cos network \mathbf{f}^* with layers \mathbf{l}_j^* can be formulated in exactly the same way as

$$\mathbf{f}^*(\mathbf{x}; \theta) = \mathbf{l}_L^* \circ \mathbf{l}_{L-1}^* \circ \dots \circ \mathbf{l}_2^* \circ \mathbf{l}_1^*(\mathbf{x}), \quad (9)$$

with the only difference being that every dot product (here between rows of \mathbf{W}_j and inputs \mathbf{a}_j) is replaced by the B-cos transformation in Eq. (4). In matrix form, this equates to

$$\mathbf{l}_j^*(\mathbf{a}_j; \mathbf{W}_j) = (\widehat{\mathbf{W}}_j \mathbf{a}_j) \times |c(\mathbf{a}_j; \widehat{\mathbf{W}}_j)|^{B-1}. \quad (10)$$

Here, the power, absolute value, and \times operators are applied element-wise, $c(\mathbf{a}_j; \widehat{\mathbf{W}}_j)$ computes the cosine similarity between input \mathbf{a}_j and the rows of $\widehat{\mathbf{W}}_j$, and the hat operator scales the rows of $\widehat{\mathbf{W}}_j$ to unit norm. Finally, as for $B>1$ the transformation \mathbf{l}_j^* is already *non-linear*, a non-linear ϕ is not required for modelling non-linear relationships.

The above discussion readily generalises to convolutional neural networks (CNNs): in CNNs, we replace the linear transformations computed by the convolutional kernels by B-cos. Further, although we assumed a plain multi-layer network without add-ons such as skip connections, we show in Sec. 5 that the benefits of B-cos also transfer to more advanced architectures (Sec. 3.3).

3.2.1 Computing explanations for B-cos Networks

As can be seen by rewriting Eq. (10), a B-cos layer effectively computes an input-dependent linear transformation:

$$\mathbf{l}_j^*(\mathbf{a}_j; \mathbf{W}_j) = \widetilde{\mathbf{W}}_j(\mathbf{a}_j) \mathbf{a}_j, \quad (11)$$

$$\text{with } \widetilde{\mathbf{W}}_j(\mathbf{a}_j) = |c(\mathbf{a}_j; \widehat{\mathbf{W}}_j)|^{B-1} \odot \widehat{\mathbf{W}}_j. \quad (12)$$

Here, \odot scales the rows of the matrix to its right by the scalar entries of the vector to its left. Hence, the output of a B-cos Network, see Eq. (9), is effectively calculated as

$$\mathbf{f}^*(\mathbf{x}; \theta) = \widetilde{\mathbf{W}}_L(\mathbf{a}_L) \widetilde{\mathbf{W}}_{L-1}(\mathbf{a}_{L-1}) \dots \widetilde{\mathbf{W}}_1(\mathbf{a}_1 = \mathbf{x}) \mathbf{x}. \quad (13)$$

As multiple linear transformations in sequence can be collapsed to a single one, $\mathbf{f}^*(\mathbf{x}; \theta)$ can be written as

$$\mathbf{f}^*(\mathbf{x}; \theta) = \mathbf{W}_{1 \rightarrow L}(\mathbf{x}) \mathbf{x}, \quad (14)$$

$$\text{with } \mathbf{W}_{1 \rightarrow L}(\mathbf{x}) = \prod_{j=1}^L \widetilde{\mathbf{W}}_j(\mathbf{a}_j). \quad (15)$$

Thus, $\mathbf{W}_{1 \rightarrow L}(\mathbf{x})$ faithfully summarises the network computations (Eq. (9)) by a single linear transformation (Eq. (14)).

To explain an activation (e.g., the class logit), we can now either directly visualise the corresponding row in $\mathbf{W}_{1 \rightarrow L}$, see Figs. 1 and 12, or the *contributions* according to $\mathbf{W}_{1 \rightarrow L}$ coming from individual input dimensions. We use the resulting spatial contributions maps to quantitatively evaluate the explanations. In detail, the input contributions $\mathbf{s}_j^l(\mathbf{x})$ to neuron j in layer l for an input \mathbf{x} are given by the

$$\text{contribution map } \mathbf{s}_j^l(\mathbf{x}) = [\mathbf{W}_{1 \rightarrow l}(\mathbf{x})]_j^T \odot \mathbf{x}, \quad (16)$$

with $[\mathbf{W}_{1 \rightarrow l}]_j$ denoting the j th row in matrix $\mathbf{W}_{1 \rightarrow l}$; as such, the contribution from a single pixel location (x, y) is given by $\sum_c [\mathbf{s}_j^l(\mathbf{x})]_{(x,y,c)}$ with c the color channels.

3.2.2 Optimising B-cos Networks for classification

In the following, we discuss why the linear transformations $\mathbf{W}_{1 \rightarrow L}$ (Eq. (14)) can be expected to align with relevant input patterns for models based on B-cos transformations.

For this, first note that the output of each neuron—and thus of each layer—is bounded, cf. Eqs. (5) and (10). Since the output of a B-cos Network is computed as a sequence of such bounded transformations, see Eq. (13), the output of the network as a whole is also bounded. Secondly, note that a B-cos Network only achieves its upper bound for a given input if every unit achieves its upper bound. Importantly, the individual units can only achieve their maxima by aligning with their inputs (cf. Eq. (5)). Hence, optimising

a B-cos Network to maximise its output over a set of inputs will optimise the model weights to align with those inputs.

To take advantage of this for B-cos classification models, we train them with the binary cross entropy (BCE) loss

$$\mathcal{L}(\mathbf{x}_i, \mathbf{y}_i) = \text{BCE}(\sigma(\mathbf{f}^*(\mathbf{x}_i; \theta)/T + \mathbf{b}), \mathbf{y}_i) \quad (17)$$

for input \mathbf{x}_i and the respective one-hot encoding of the label \mathbf{y}_i ; see below (target encoding) for an additional discussion. Here, σ denotes the sigmoid function, \mathbf{b} and T fixed bias and scaling terms, and θ the model parameters. Note that we choose the BCE loss as it directly entails maximisation. Specifically, in order to reduce the BCE loss, the network is optimised to maximise the (negative) class logit for the correct (incorrect) classes. As discussed in the previous paragraph, this will optimise the weights in each layer of the network to align with their inputs. In particular, they will need to align with class-specific input patterns such that these result in large outputs for the respective class logits.

Finally, note that increasing B allows to specifically reduce the output of badly aligned weights in each layer (cf. Eq. (5)). This decreases the layer's output strength and thus the output of the network as a whole for badly aligned weights, which increases the alignment pressure during optimisation (thus, higher $B \rightarrow$ higher alignment, see Fig. 4).

Target encoding and logit bias b . As discussed, BCE optimises the models to maximise the absolute magnitude of both the positive (target class) as well as the negative logits (other classes). To encourage the models to find *positive* evidence (*i.e.* contributions as defined in Eq. (16)) for the target class, rather than negative evidence for the other classes, we set the logit bias \mathbf{b} to $\log(1/(C-1))$ for C classes, which corresponds to a default class probability of $1/C$ for each class. As a result, the magnitude of $\mathbf{f}^*(\mathbf{x}_i; \theta)$ can be small for non-target classes without incurring a large loss, while it has to be large for the target class to achieve a class confidence close to 1 and thus reduce the BCE loss.

Interestingly, in our experiments we found that models with normalisation layers (see Sec. 3.3.2) can still exhibit a tendency to focus on negative evidence, especially when the number of classes C is large (*e.g.* on ImageNet), see Sec. 5.4. To overcome this, for our ImageNet experiments we additionally modify the target encoding \mathbf{y}_i to be $1/C$ for all non-target classes, which encourages the models to produce weights that are orthogonal to the input for non-target classes, *i.e.* the optimum can only be reached if $[\mathbf{f}^*(\mathbf{x}_i; \theta)]_c = 0$ for all c that are not the target class.

3.3 Advanced B-cos Networks

To test the generality of our approach, we investigate how integrating B-cos transformations into conventional DNNs affects their classification performance and interpretability. To do this, in the following we discuss how to combine B-cos layers and their explanations with other commonly used model components, such as activation functions (Sec. 3.3.1), normalisation (Sec. 3.3.2), and attention layers (Sec. 3.3.3).

Note that whenever converting a linear / convolutional layer in a given network to the corresponding B-cos layer, **we remove all bias terms** to ensure that the model output is given by a dynamic linear transformation as in Eq. (14). For a discussion on the impact of biases, see Sec. 5.4.

3.3.1 Activation functions in B-cos networks

Since B-cos transformations themselves are non-linear, B-cos DNNs do not require activation functions to model non-linear relationships (see Sec. 3.2). Hence, most of our models in Sec. 5 do not employ any activation functions. Nonetheless, combining B-cos networks with activation functions can still be beneficial. In particular, as we discuss in Sec. 5.1, adding an additional non-linearity allows us to study the effect of the parameter B in Eq. (3) over a wide range of values, including $B=1$, and thus to directly compare them to conventional piece-wise linear models².

While there are many potential non-linearities to choose from, in this work, we specifically explore the option of combining the B-cos transformation with MaxOut [51]. For this, we model every neuron by 2 B-cos transformations of which only the maximum is forwarded to the next layer:

$$\text{mB-cos}(\mathbf{x}) = \max_{i \in \{1,2\}} \{\text{B-cos}(\mathbf{x}; \mathbf{w}_i)\} . \quad (18)$$

We do so for several reasons. First, this operation maintains the models' dynamic linearity and is thus compatible with the B-cos explanations. Secondly, as discussed above, for $B=1$ the resulting model is a piece-wise linear model, which allows us to interpolate between conventional piece-wise linear models and the proposed B-cos models. Lastly, while the ReLU [2] operation also fulfills these properties, we noticed that B-cos models without any normalisation layers were much easier to optimise with MaxOut instead of ReLU, which we attribute to avoiding 'dying neurons' [51].

3.3.2 Adapting normalisation layers for B-cos Networks

In this section we discuss how normalisation layers can be integrated into B-cos Networks. Specifically, we note that these layers (*e.g.* BatchNorm [52], LayerNorm [53], PositionNorm [54], or InstanceNorm [55]) are (1) non-linear operations and (2) can pose a challenge to the *completeness* of the explanations. By interpreting them as dynamic linear functions and removing bias terms, we can nonetheless seamlessly integrate them into B-cos explanations. First, however, let us define the normalisation layers.

Definition. Normalisation layers have been shown to benefit neural network training by whitening the input \mathbf{x} via

$$\text{*Norm}(\mathbf{x}, \mathcal{X}; \gamma, \beta) = \frac{\mathbf{x} - \langle \mathcal{X} \rangle_*}{\sqrt{\text{var}_*(\mathcal{X})}} \times \gamma + \beta , \quad (19)$$

with $\mathcal{X} \in \mathbb{R}^{b \times c \times h \times w}$ a batch of b inputs, each with c channels and a width w and a height h ; further, $\mathbf{x} \in \mathbb{R}^c$ is a single representation vector from a specific input at a given location. As indicated by the placeholder $*$, the individual normalisation layers differ mainly in the choice of dimensions over which the mean $\langle \cdot \rangle$ and the variance var are computed. Finally, note that at inference time, the mean and variance terms are sometimes replaced by running estimates of those values, most commonly done so in BatchNorm [52].

2. For $B=1$ linear and B-cos transformations are equivalent (Eq. (4)).

To avoid ambiguities³ and facilitate the comparison of different normalisation schemes, we define the normalisation layers by the respective choice of dimensions they use:

$$\text{BatchNorm}(\mathbf{x}, \mathcal{X}; \gamma, \beta) := \text{BHW-Norm}(\mathbf{x}, \mathcal{X}; \gamma, \beta) \quad (20)$$

$$\text{LayerNorm}(\mathbf{x}, \mathcal{X}; \gamma, \beta) := \text{CHW-Norm}(\mathbf{x}, \mathcal{X}; \gamma, \beta) \quad (21)$$

$$\text{InstanceNorm}(\mathbf{x}, \mathcal{X}; \gamma, \beta) := \text{HW-Norm}(\mathbf{x}, \mathcal{X}; \gamma, \beta) \quad (22)$$

$$\text{PositionNorm}(\mathbf{x}, \mathcal{X}; \gamma, \beta) := C - \text{Norm}(\mathbf{x}, \mathcal{X}; \gamma, \beta) \quad (23)$$

i.e. BatchNorm estimates mean and variance for each channel independently (no ‘C’) over all spatial dimensions (HW) and inputs in the batch (B), whereas LayerNorm computes mean and variance across all activations (CHW) in a single input, but does not use other inputs to estimate those values.

By defining the normalisation layers as above, an additional candidate naturally emerges, which we call *AllNorm*:

$$\text{AllNorm}(\mathbf{x}, \mathcal{X}; \gamma, \beta) := \text{BCHW-Norm}(\mathbf{x}, \mathcal{X}; \gamma, \beta). \quad (24)$$

Dynamic linear normalisation. Normalising the input by the variance of course constitutes a highly non-linear operation. To nonetheless integrate this operation into the B-cos-based dynamic linear explanations (cf. Eq. (13)), we interpret the normalisation itself to be a dynamic linear function with:

$$\mathbf{w}_{*\text{Norm}}(\mathbf{x}, \mathcal{X}) = \gamma \times \sqrt{\text{var}_*^{-1}(\mathcal{X})} \quad (25)$$

$$\text{s.t. } * \text{Norm}(\mathbf{x}, \mathcal{X}; \gamma, \beta) = \mathbf{w}_{*\text{Norm}} \times (\mathbf{x} - \langle \mathcal{X} \rangle_*) + \beta, \quad (26)$$

Explanation completeness. As becomes clear from Eq. (19), normalisation layers introduce additive biases into the model prediction (i.e. β and, if running estimates are used, by the running mean estimation of $\langle \mathcal{X} \rangle_*$). The model output is thus not given by a dynamic linear transformation anymore (cf. Eq. (14)), but instead by an *affine* transformation $\mathbf{f}^*(\mathbf{x}; \theta) = \mathbf{W}(\mathbf{x})\mathbf{x} + \mathbf{b}(\mathbf{x})$ with an input-dependent bias $\mathbf{b}(\mathbf{x})$.

The contribution maps in Eq. (16) would thus not be *complete*, as they do not take the bias into account. This, however, is undesirable for model-faithful explanations [25], [26], as the bias terms can play a significant role in the model decision (see Sec. 5.4).

To alleviate this, we use bias-free (BF) normalisation [42], and remove the *learnt*⁴ biases β :

$$* \text{Norm}^{\text{BF}}(\mathbf{x}, \mathcal{X}; \gamma, \beta) = \frac{\mathbf{x} - \langle \mathcal{X} \rangle_*}{\sqrt{\text{var}_*(\mathcal{X})}} \times \gamma + 0. \quad (27)$$

Moreover, for normalization layers that replace the mean computation $\langle \mathcal{X} \rangle_*$ by a running mean estimate at inference time (i.e. BatchNorm), we additionally remove the centering term $\langle \mathcal{X} \rangle_*$, since otherwise external biases would be introduced at inference time.

3.3.3 B-cos ViTs: combining B-cos and attention layers

Thus far, our discussion focused on designing interpretable convolutional neural networks (CNNs), which have long dominated computer vision research. Recently, however,

3. Note that LayerNorm has in fact been interpreted differently by different authors (e.g. [13] vs. [54], [56]), i.e. either as Eq. (21) or Eq. (23).

4. Note that the mean subtraction does not introduce *external* biases that would be unaccounted for in the explanations, as this operation can be modelled by a linear transformation of the input.

CNNs are often surpassed by transformers [44], which—if the current development is any indication—will replace CNNs for ever more tasks and domains. In the following, we therefore investigate how to design B-cos transformers.

For this, first note that the core difference between vision transformers (ViTs) and CNNs is the use of attention layers and additive positional encodings. The remaining ViT components (tokenisation, MLP blocks, normalisation layers, classification head) have a direct convolutional counterpart.

In particular, the tokenisation module is typically given by a $K \times K$ convolution with a stride and kernel size of K , the MLP blocks effectively perform 1x1 convolutions (cf. [13]), the normalisation layers can be expressed as PositionNorm in CNNs (Eq. (23)) and the classification head is either given by a single or multiple linear layers.

Therefore, in the following, we discuss the attention mechanism and the positional embeddings in more detail.

Attention. Interestingly, the attention operation itself is already dynamic linear and attention thus lends itself well to be integrated into the linear summary according to Eq. (14):

$$\text{Att}(\mathbf{X}; \mathbf{Q}, \mathbf{K}, \mathbf{V}) = \underbrace{\text{softmax}(\mathbf{X}^T \mathbf{Q}^T \mathbf{K} \mathbf{X})}_{\text{Attention matrix } \mathbf{A}(\mathbf{X})} \underbrace{\mathbf{V} \mathbf{X}}_{\text{Value}(\mathbf{X})} \quad (28)$$

$$= \underbrace{\mathbf{A}(\mathbf{X}) \mathbf{V}}_{\mathbf{W}(\mathbf{X})} \mathbf{X} = \mathbf{W}(\mathbf{X}) \mathbf{X}. \quad (29)$$

Here, \mathbf{Q} , \mathbf{K} and \mathbf{V} are the query, key, and value transformation matrices and \mathbf{X} denotes the matrix of input tokens to the attention layer. Softmax is computed column-wise.

In multi-head self-attention (MSA), H attention heads are used in parallel. Their concatenated outputs are then linearly projected by a matrix \mathbf{U} via

$$\text{MSA}(\mathbf{X}) = \mathbf{U} [\mathbf{W}_1(\mathbf{X}) \mathbf{X}, \dots, \mathbf{W}_H(\mathbf{X}) \mathbf{X}], \quad (30)$$

which is still dynamic linear. While there are multiple linear operations involved in MSA (query, key, value, and projection computations) and therefore various options for introducing B-cos layers, we empirically found not replacing⁵ the query, key, and value computations to yield the best results, which we attribute to a potentially higher compatibility with the softmax function. Hence, to adapt the attention layers to the B-cos framework, we only replace the linear projection via \mathbf{U} by a B-cos transformation in our experiments.

Positional encoding. In contrast to CNNs, which possess a strong inductive bias regarding spatial relations (local connectivity), transformers are invariant with respect to the token order and thus lack such a ‘locality bias’. To nevertheless leverage spatial information, it is common practice to break the symmetry between tokens by adding a (learnt) embedding \mathbf{E} to the input tokens \mathbf{X} such that $\mathbf{X}' = \mathbf{X} + \mathbf{E}$.

While explanations solely with respect to \mathbf{X} are not thus complete (cf. Sec. 3.3.2), since part of the model output will be based on contributions from \mathbf{E} , we experimentally find that the positional embeddings do not negatively impact the interpretability of the ViTs. How to design ViT models without the need for such external biases represents an interesting direction for future work.

5. I.e., these transformations effectively use a B-cos layer with $B=1$.

4 EXPERIMENTAL SETTING

Datasets. We evaluate the accuracies of a wide range of B-cos networks on the CIFAR-10 [7] and the ImageNet [57] datasets. We use the same datasets for the qualitative and quantitative evaluations of the model-inherent explanations.

CIFAR10 models. For the CIFAR10 experiments, we develop a simple fully-convolutional B-cos DNN, consisting of 9 convolutions, each with a kernel size of 3, followed by a global pooling operation. We evaluate a network without additional non-linearities as well as with MaxOut units, see Sec. 3.3.1. Additionally, to study the normalisation layers and the effect of the bias terms (cf. Eq. (27)), we evaluate various B-cos ResNets. For this we adapt a conventional ResNet-56 [10] by removing all activation functions and bias parameters and replacing all linear / convolutional layers by their corresponding B-cos version. Additionally, we replace the conventional batch normalisation layers by the (modified) normalisation layers as described in Sec. 3.3.2.

ImageNet models. For the ImageNet experiments, we rely on the publicly available [58] implementations of a wide range of CNNs (VGGs [9], ResNets [10], DenseNets [11], ResNext [12] and ConvNext models [13]). Further, we evaluate B-cos versions of the Vision Transformers (ViTs) [8], specifically the version from [59], as well as ViT_C models as in [60], *i.e.* ViTs with a shallow stem of four convolutional layers. We adapt all of those architectures to B-cos networks as described in Sec. 3.3 and train them according to standard protocols. Specifically, we train most CNNs with the default 90 epochs training paradigm as found in the torchvision library [61]; additionally, we evaluate various models when employing a more sophisticated training recipe [62] of 600 epochs, based on the ConvNext [13] training protocol. Finally, for the ViT models, we follow the protocol proposed by [59], which has shown strong performance for ViTs with only 90 epochs of training. For more details on the training procedure, see supplement (Sec. C).

Image encoding. We add three additional channels and encode images as $[r, g, b, 1-r, 1-g, 1-b]$, with $r, g, b \in [0, 1]$ the red, green, and blue color channels. On the one hand, this reduces a bias towards bright regions in the image⁶ [29]. On the other hand, colors with the *same angle in the original encoding*—*i.e.*, $[r_1, g_1, b_1] \propto [r_2, g_2, b_2]$ —are *unambiguously encoded by their angles under the new encoding*. Therefore, the linear transformation $\mathbf{W}_{1 \rightarrow l}$ can be decoded into colors just based on the angles of each pixel, see Fig. 1. For a detailed discussion, see supplement (Sec. D).

Evaluating explanations. To compare explanations for the model decisions and evaluate their faithfulness, we employ the *grid pointing game* [29]. That means we evaluate the trained models on a synthetic 3x3 grid of images of different classes and for each of the corresponding class logits measure how much positive attribution an explanation method assigns to the correct location in the grid, see Fig. 3 for an example. Note that due to the global attention in ViTs such grids are much further out of distribution at every layer of the model than they are for CNNs, for which the

6. Note that the model output is given by a weighted sum over the input (Eq. (14)). Since black pixels are encoded as zero in the conventional encoding, they cannot contribute to the class logits (cf. Eq. (16)).

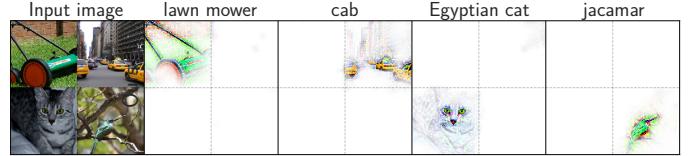


Fig. 3: 2x2 pointing game example. **Column 1:** input image. **Columns 2–5:** explanations for individual class logits.

locally applied convolutional kernels are not affected by the synthetic nature of the input grids. The grid pointing game, however, relies on the assumption that the model extracts similar features in the synthetic setting and *locally* correctly classifies the subimages. To ensure that this assumption (at least approximately) holds, we therefore apply the ViTs in a sliding window fashion to the synthetic image grids, *i.e.* we apply the ViTs to patches of size 224x224 extracted at a stride of 112 from the image grid. As this significantly increases the computational cost, we use 2x2 grids for ViTs.

Following [29], we construct 500 grids from the most confidently and correctly classified images and compare the model-inherent contribution maps (Eq. (16)) against several commonly employed post-hoc explanation methods under two settings. First, we evaluate all methods on B-cos networks to investigate which method provides the best explanation *for the same model*. Secondly, we further evaluate the post-hoc methods on pre-trained versions of the original models. This allows to compare explanations *between different models* and assess the *explainability gain* obtained by converting conventional models to B-cos networks. Lastly, all attribution maps (except for GradCAM, which is of much lower resolution to begin with) are smoothed by a 15x15 (3x3) kernel to better account for negative attributions in the localisation metric for ImageNet (CIFAR-10) images, which is small with respect to the image size of 224x224 (32x32).

Visualisations details. For generating the visualisations of the linear transforms for individual neurons n in layer l (cf. Figs. 1 and 12), we proceed as follows. First, we select all pixel locations (x, y) that positively contribute to the respective activation (*e.g.*, class logit) as computed by Eq. (16); *i.e.*, $\{(x, y) \text{ s.t. } \sum_c [\mathbf{s}_n^l(\mathbf{x})]_{(x,y,c)} > 0\}$ with c the 6 color channels (see image encoding). Then, we normalise the weights of each color channel such that the corresponding weights (*e.g.*, for r and $1-r$) sum to 1. Note that this normalisation maintains the angle for each color channel pair (*i.e.*, r and $1-r$), but produces values in the allowed range $r, g, b \in [0, 1]$. These normalised weights can then directly be visualised as color images. The opacity of a pixel is set to $\min(\|\mathbf{w}_{(x,y)}\|_2 / p_{99.9}, 1)$, with $p_{99.9}$ the 99.9th percentile over the weight norms $\|\mathbf{w}_{(x,y)}\|_2$ across all (x, y) .

5 RESULTS

In this section, we analyse the performance and interpretability of B-cos models. For this, in Sec. 5.1 we show results of ‘simple’ B-cos models without advanced architectural elements such as skip connections. In this context, we investigate how the B parameter influences B-cos models in terms of performance and interpretability. Thereafter, in Sec. 5.2, we present quantitative results of the *advanced* B-cos

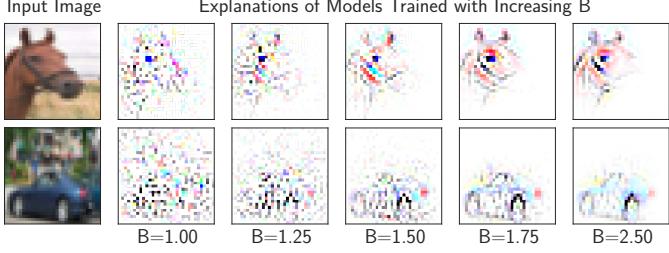


Fig. 4: **Col. 1:** Input images. **Cols. 2-6:** Explanations for different classes c (top: ‘horse’; bottom: ‘car’) of models trained with increasing B . For higher B , the model-inherent linear explanations $[\mathbf{W}_{1 \rightarrow l}(\mathbf{x}_i)]_c$ increasingly align with discriminative input patterns, thus becoming more interpretable.

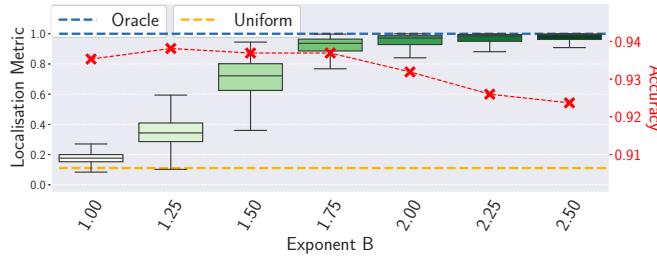


Fig. 5: Accuracy (crosses) and localisation (box plots) results for a B-cos network trained with different B . While decreasing accuracy, larger B significantly improve localisation.

models, i.e., B-cos models based on common DNN architectures (cf. Sec. 3.3). In Sec. 5.3, we visualise and qualitatively discuss explanations for individual neurons of the advanced B-cos models, and finally, in Sec. 5.4, we discuss the impact of the bias terms in the normalisation layers.

5.1 Simple B-cos models

In the following, we discuss the experimental results of simple B-cos models evaluated on the CIFAR-10 dataset.

B	MaxOut B-cos networks						plain	
	1.00	1.25	1.50	1.75	2.00	2.25	2.50	2.00
Accuracy (%)	93.5	93.8	93.7	93.7	93.2	92.6	92.4	91.5

TABLE 1: **CIFAR-10** accuracy of a B-cos model without additional non-linearities (**plain**) and for B-cos models with MaxOut (Eq. (18)) and increasing values for B (left to right).

Accuracy. In Tab. 1, we present the test accuracies of various B-cos models trained on CIFAR-10. We show that a plain B-cos model ($B=2$) without any add-ons (ReLU, batch norm, etc.) can achieve competitive⁷ performance. By modelling each neuron via 2 MaxOut units (Eq. (18)), the performance can be increased and the resulting model ($B=2$) performs on par with a ResNet-56 (achieving 93.0%, see [10]). Further, we see that an increase in the parameter B leads to a decline in performance from 93.8% for $B=1.25$ to 92.4% for $B=2.5$. Notably, despite its simple design, our strongest model with $B=1.25$ performs similarly to the strongest ResNet model (93.6%) reported in [10].

7. A ResNet-20 achieves 91.2% [10] with the same data augmentation.

Model interpretability. As discussed in Sec. 3.2.2, we expect an increase in B to increase the alignment pressure on the weights during optimisation and thus influence the models’ optima, similar to the single unit case in Fig. 2. This is indeed what we observe. For example, in Fig. 4, we visualise $[\mathbf{W}_{1 \rightarrow l}(\mathbf{x}_i)]_{y_i}$ (see Eq. (14)) for different samples i from the CIFAR-10 test set. For higher values of B , the weight alignment increases notably from piece-wise linear models ($B=1$) to B-cos models with higher B ($B=2.5$). Importantly, this does not only lead to an increase in the visual quality of the explanations, but also to quantifiable gains in model interpretability. In particular, as we show in Fig. 5, the spatial contribution maps defined by $\mathbf{W}_{1 \rightarrow l}(\mathbf{x}_i)$ (see Eq. (16)) of models with larger B values score significantly higher in the localisation metric (see Sec. 4).

5.2 Advanced B-cos models

B-cos CNNs – Accuracy. In Tabs. 2 and 3, we present the top-1 accuracies of the B-cos models on the ImageNet validation set for the 90 epoch and 600 epoch training paradigms respectively. For comparison, we also show the difference to the accuracy of the respective baseline models as reported in [61] (Δ^1). As can be observed, the B-cos models are highly competitive, achieving accuracies on par with the baselines for some models (ResNet-50, $\Delta^1 = -0.2$), and a worst-case drop of $\Delta^1 = -2.0$ (ResNext-50-32x4d) under the 90 epoch training paradigm (Tab. 2). Interestingly, we find that the differences to the baseline models vanish when training the baselines without bias terms in the convolution and normalisation layers, *i.e.* as is done in B-cos models; we denote the difference to the baselines without biases by Δ^2 .

When employing a long training protocol with additional data augmentation as in [13], the performance of the B-cos models can be further improved: *e.g.*, the B-cos ResNet-152 increases its top-1 accuracy by 3.2pp (compare Tabs. 2 and 3) and achieves a top-1 accuracy of 80.2%, which is an unprecedented performance for models that inherently provide such detailed explanations (see Sec. 5.3). Nonetheless, especially for the ConvNext models, we observe a significant performance drop with respect to the baseline numbers reported in [61] (up to 5.0pp for ConvNext-Tiny). Note, however, that the training protocol has been meticulously optimised for ConvNext models in [13]; as such, we expect that the performance of B-cos models will significantly improve if a similar effort as in [13] is undertaken for optimising the model architecture and training procedure.

B-cos CNNs – Interpretability. Apart from the accuracy, in Tabs. 2 and 3 we further provide the mean localisation scores obtained in the grid pointing game (cf. Fig. 3 and Sec. 4) of the model-inherent explanations according to Eq. (16). Additionally, we show the difference to one of the most popular post-hoc attribution methods (IntGrad) evaluated on the B-cos models ($\Delta_{\text{IntG}}^{\text{B-cos}}$) and the baselines ($\Delta_{\text{IntG}}^{\text{basel.}}$).

Across all models, we find that the model-inherent explanations achieve significantly higher localisation scores than IntGrad; this difference is even more pronounced when comparing to the IntGrad explanations for the baselines.

To provide more detail, in Fig. 6 we show the localisation results for two specific models (ResNet-152, DenseNet-201) for a wide range of post-hoc explanation methods:

B-cos Model	Accuracy			Localisation		
	%	Δ^1	Δ^2	%	$\Delta_{\text{IntG}}^{\text{B-cos}}$	$\Delta_{\text{IntG}}^{\text{basel.}}$
VGG-11	69.0	-1.4	-0.6	85.6	+26.4	+64.5
ResNet-18	68.7	-1.1	+0.3	86.9	+34.0	+64.0
ResNet-34	72.1	-1.2	-0.1	89.0	+36.9	+65.2
ResNet-50	75.9	-0.2	+2.4	90.2	+32.7	+63.3
ResNet-101	76.3	-1.1	+6.2	91.8	+33.9	+64.4
ResNet-152	76.6	-1.7	+3.2	91.3	+33.6	+63.9
DenseNet-121	73.6	-0.8	+0.2	92.1	+41.2	+69.1
DenseNet-161	76.6	-0.5	+5.2	93.4	+47.1	+68.5
DenseNet-169	75.0	-0.6	+0.3	91.8	+49.3	+67.7
DenseNet-201	75.6	-1.3	+0.3	93.0	+44.8	+68.8
ResNext-50-32x4d	75.6	-2.0	+0.6	91.2	+28.2	+64.1

TABLE 2: **Left:** Top-1 accuracies on ImageNet for a standard 90 epoch training protocol [61]. Δ^1 : Accuracy difference with respect to the numbers reported in [61]. Δ^2 : Accuracy difference with respect to baselines trained without biases. This setting puts the baselines and the B-cos models on an equal footing, as B-cos models do not use biases either. **Right:** Additionally, we show localisation scores of the model-inherent explanations (Eq. (16)) as well as localisation scores of IntGrad evaluated on the B-cos models ($\Delta_{\text{IntG}}^{\text{B-cos}}$) and the pre-trained baselines ($\Delta_{\text{IntG}}^{\text{basel.}}$) as obtained from [61]. We observe significant localisation gains for all models.

B-cos Model	Accuracy		Localisation		
	%	Δ^1	%	$\Delta_{\text{IntG}}^{\text{B-cos}}$	$\Delta_{\text{IntG}}^{\text{basel.}}$
ResNet-50	79.5	-1.4	86.2	+41.0	+58.6
ResNet-152	80.2	-2.1	85.3	+44.2	+55.5
ConvNext-Tiny	77.5	-5.0	69.6	+37.7	+46.5
ConvNext-Base	79.7	-4.4	81.4	+35.2	+54.6

TABLE 3: **Left:** ImageNet accuracies for the ConvNext-inspired training protocol [62]. Δ^1 : Accuracy difference with respect to the numbers reported in [62]. While we observe a significant drop with respect to the ConvNext models, we note that the training protocol has been optimised explicitly for those architectures [13] and that B-cos models do not employ biases (cf. Tab. 2) to ensure $y(x) = W(x)x$. **Right:** Additionally, we show the localisation scores of the model-inherent explanations (Eq. (16)) as well as the localisation scores of IntGrad evaluated on the B-cos models ($\Delta_{\text{IntG}}^{\text{B-cos}}$) and the pre-trained baselines ($\Delta_{\text{IntG}}^{\text{basel.}}$) as obtained from [61]. We observe significant localisation gains for all models.

the ‘vanilla’ gradient (Grad) [63], Guided Backpropagation (GB) [21], Input \times Gradient (IxG) [5], Integrated Gradients (IntGrad) [26], DeepLIFT [5], GradCAM [24], and LIME [17]. In particular, we plot the localisation scores for the B-cos versions of those models (rows 1+2) and the respective baselines (rows 3+4).

In comparison to the post-hoc explanation methods, we observe the model-inherent explanations (denoted as ‘Ours’, Eq. (16)) to consistently and significantly outperform all other methods, both when comparing to post-hoc explanations for the *same* model (rows 1+2), as well as when comparing explanations *across* models (for ease of comparison, we repeat the B-cos results in rows 3+4). Note that while GradCAM also shows strong performance, it explains only a fraction of the entire model [64] and thus yields much coarser attribution maps, see also Fig. 9. Interestingly, we also find that the other post-hoc explanations consistently

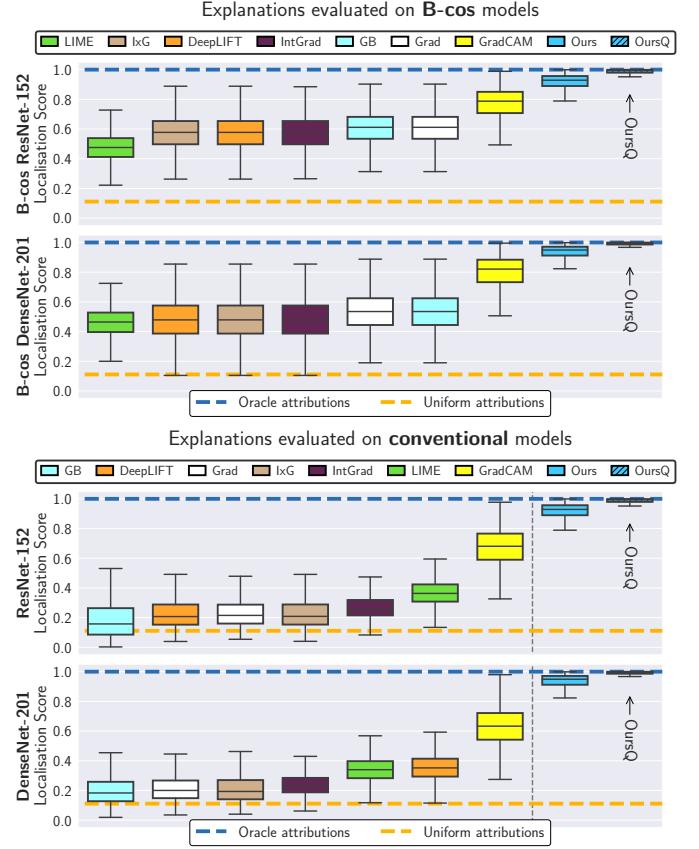


Fig. 6: Localisation results of the model-inherent contribution maps (‘Ours’, Eq. (16)) and various post-hoc explanation methods for a B-cos ResNet-152 and a B-cos DenseNet-201 (rows 1+2) and their conventional counterparts (rows 3+4); we further show the results of the top 0.025 percentile of pixels in the model-inherent explanations (‘OursQ’), i.e. roughly 11.3k pixels, cf. Fig. 7. For an easier cross-model comparison, we repeat the results of the B-cos models of rows 1+2 in rows 3+4 (Ours + OursQ). As can be seen, the model-inherent explanations significantly outperform post-hoc explanation methods when evaluating those methods on the B-cos models (rows 1+2) and on the conventional counterparts (rows 3+4).

perform better for B-cos models than for the baseline models (compare rows 1+3 and 2+4).

Finally, we additionally evaluate how well the most highly contributing pixels according to $W(x)$ fare in the localisation metric, which we denote as OursQuantile (OursQ) in Fig. 6, see also Fig. 7. We find that when using only the top-n contributing pixels in the model-inherent explanations, the localisation score can be improved even further.

B-cos ViTs – Accuracy. In Tab. 4, we report the top-1 accuracies of B-cos ViT and ViT_C models of various commonly used sizes—Tiny (Ti), Small (S), Base (B), and Large (L)—and the difference (Δ^3) to the respective baseline models which we optimised under the same training protocol [59].

For both the original ViTs and the ones with a shallow convolutional stem (ViT_C), we observe significant performance drops for smaller model sizes (Ti+S). For larger models (B+L), however, the difference between the B-cos models and the baselines vanish, see, e.g., the ViT_C -B ($\Delta^3 = +0.3$)

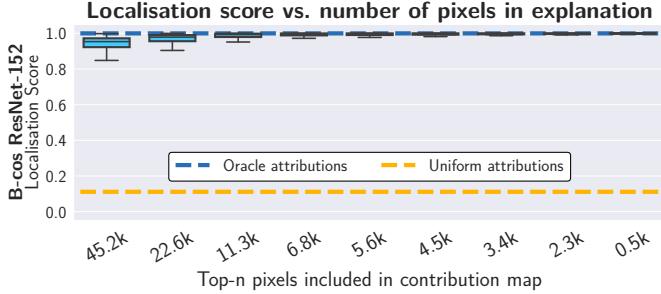


Fig. 7: Top-n pixel localisation. The localisation score significantly improves when only including the n most strongly contributing pixels of the model-inherent explanations in the localisation metric. While highlighting small contributions coming from images belonging to other classes (see Fig. 3) might be model-faithful, this is penalised in the metric. Nonetheless, we find that the correct region is consistently highlighted by the most strongly contributing pixels.

B-cos Model	Accuracy		Localisation			
	%	Δ^3	%	$\Delta_{\text{Rollout}}^*$	$\Delta_{\text{IntGrad}}^{\text{B-cos}}$	$\Delta_{\text{IntGrad}}^{\text{basel.}}$
ViT-Ti	60.0	-10.3	69.8	+44.8	+12.2	+30.0
ViT-S	69.2	-5.2	72.1	+47.1	+15.0	+30.9
ViT-B	74.4	-0.9	74.6	+49.6	+19.4	+32.5
ViT-L	75.1	-0.7	74.0	+49.0	+17.7	+32.2
ViT _C -Ti	67.3	-5.3	86.8	+61.8	+30.3	+46.4
ViT _C -S	74.5	-1.2	87.4	+62.4	+30.4	+45.9
ViT _C -B	77.1	+0.3	86.5	+61.5	+28.6	+44.5
ViT _C -L	77.8	-0.1	87.3	+62.3	+29.9	+45.3

TABLE 4: **Left:** Top-1 accuracies on ImageNet for the 90 epoch Simple-ViT training protocol [59]. Δ^3 : Accuracy difference with respect to training the baseline models according to the same protocol. While we observe a significant drop for smaller ViT and ViT_C models, the difference to the baseline models vanishes for larger ViTs: e.g., the B-cos ViT_C-B and ViT_C-L models perform on par with the baselines. **Right:** Additionally, we show the localisation scores of the model-inherent explanations (Eq. (16)) as well as the differences to the scores obtained via Attention Rollout ($\Delta_{\text{Rollout}}^*$) and IntGrad evaluated on the B-cos models ($\Delta_{\text{IntG}}^{\text{B-cos}}$) and the baselines ($\Delta_{\text{IntG}}^{\text{basel.}}$). We observe significant localisation gains when using the B-cos explanations.

and the ViT_C-L ($\Delta^3 = -0.1$) models. Further, and in line with [60], we find that a shallow convolutional stem leads to significant performance gains; as we discuss next, such a stem can also significantly improve model interpretability.

B-cos ViTs – Interpretability. To assess the B-cos ViTs’ interpretability, we compare the model-inherent explanations (Eq. (16)) to IntGrad evaluated on the B-cos ($\Delta_{\text{IntG}}^{\text{B-cos}}$) and the baseline models ($\Delta_{\text{IntG}}^{\text{basel.}}$); additionally, we report the difference to a commonly used attention-based explanation, namely Attention Rollout [46] ($\Delta_{\text{Rollout}}^*$). Since Rollout is inherently a class-agnostic method⁸, it on average yields the same localisation as uniformly distributed importance values for all models, i.e. a mean localisation score of 25%.

As can be seen in Tab. 4, the B-cos ViT and ViT_C mod-

8. Attention Rollout computes the product of the average (across heads) attention matrices of all layers and is thus not class-specific.

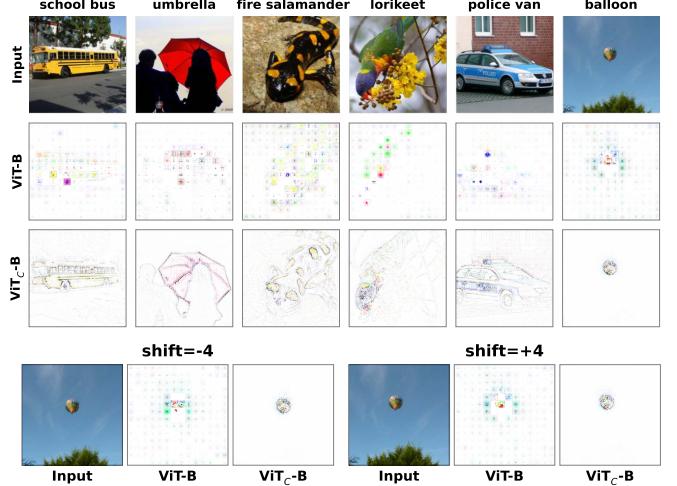


Fig. 8: Top: Comparison between explanations of B-cos ViT-B models without (row 2) and with (row 3) convolutional stem, i.e. between ViT and ViT_C models. By adding just 4 convolutional layers for embedding the image patches, the explanations become much more coherent and less sparse. **Bottom:** For a single image, we qualitatively visualise how the explanations of the different ViT models perform under diagonal shifts $s \in -4, +4$ of the image on the top right. While both models generally show consistent explanations, the inductive bias due to the convolutional stem leads to more stable behaviour and explanations of the ViT_C models.

els show significant gains in localisation compared to the respective baseline models. Interestingly, adding a convolutional stem of just four convolutional layers yields notable improvements in localisation between the respective B-cos models, i.e. between B-cos ViT and B-cos ViT_C models. This suggests that ViT_C models perform better because the transformers are applied to more meaningful representations, which avoids mixing visually similar, yet semantically unrelated, patches in the first few global attention operations.

These quantitative improvements in interpretability between ViT and ViT_C models can also be observed qualitatively, see Fig. 8. Specifically, in contrast to the explanations of the ViTs without convolutional stem (row 2), the B-cos ViT_C explanations are visually more coherent and less sparse, making them more easily interpretable for humans.

5.3 Qualitative evaluation of explanations

The following results are based on the DenseNet-121, cf. Tab. 2, we found similar results for other B-cos models.

Every activation in a B-cos model is the result of a sequence of B-cos transforms. Hence, the intermediate activations in any layer l can also be explained via the corresponding linear transform $\mathbf{W}_{1 \rightarrow l}(\mathbf{x})$, see Eq. (14).

For example, in Fig. 1, we visualise the linear transforms of the respective *class logits* for various input images. Given the alignment pressure during optimisation, these linear transforms align with class-discriminatory patterns in the input and thus actually resemble the class objects.

Similarly, in Figs. 10 and 12, we visualise explanations for *intermediate features*. These features are selected from the most highly contributing feature directions $\hat{\mathbf{a}}_{(x,y)}^l \in \mathbb{R}^{d_l}$

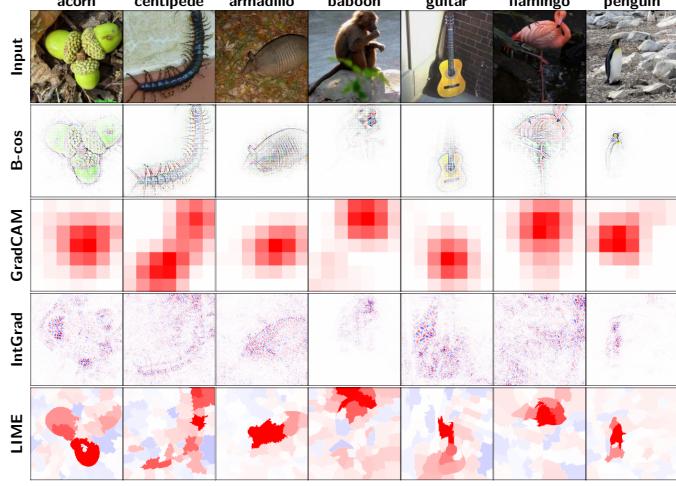


Fig. 9: Comparison between B-cos explanations and various commonly used post-hoc methods on a B-cos DenseNet-121. For all post-hoc explanations, positive, negative, and zero attributions are shown in red, blue, and white respectively.



Fig. 10: Explanations for highly contributing feature directions in various layers. In early layers the features seem to encode low-level concepts (e.g., eyes in layer 39) and represent more high-level concepts in later layers (e.g., layers 89 and 121), such as neck braces or parachutes, see also Fig. 12.

according to the linear transformation $\mathbf{W}_{l \rightarrow L}(\mathbf{x})$, cf. Eq. (14); here, $\hat{\mathbf{a}}$ denotes that $\hat{\mathbf{a}}$ is of unit norm, d_l is the number of features in layer l , and (x, y) denotes the position at which this feature is extracted from the activation map.

Specifically, in Fig. 10, we show explanations for some images that most strongly activate those feature directions across the validation set. We find that features in early layers seem to represent low-level concepts, and become more complex in later layers.

Fig. 12 shows additional results for features in layer 87. We observed that some features are highly specific to certain concepts, such as bike wheels (feature 89), red bird beaks (feature 109), or faces with headgear (feature 123). Importantly, these features do not just learn to align with simple,

fixed patterns—instead, they represent semantic concepts and are robust to changes in colour, size, and pose.

Lastly, in Fig. 11, we show explanations of the two most likely classes for images for which the model produces predictions with high uncertainty; additionally, we show the Δ -Explanation, i.e., the difference in contribution maps for the two classes, see Eq. (16). By means of the model-inherent linear mappings $\mathbf{W}_{1 \rightarrow L}$, the model can provide a human-interpretable explanation for its uncertainty: there are indeed features in each of those images that provide evidence for both of the predicted classes.

5.4 Explicit and implicit model biases

As discussed in Sec. 3.3.2, we design the B-cos CNNs to be bias-free by fixing the *explicit* bias parameters in all layers to be zero. In this section, we aim to motivate this choice further. Additionally, we show that even when doing so, the models can learn to counteract this and add *implicit* biases by exploiting reliable input features such as image edges.

Explicit biases. As discussed by [25] for piece-wise linear networks, the bias parameters of DNNs are often not accounted for in input-level attributions, despite the fact that they can play a critical role in the model prediction. Similarly, the contributions s_c (Eq. (16)) from the input to the c -th class logit are not complete explanations if biases are used within B-cos models. Instead, the output $y(\mathbf{x})$ is

$$y(\mathbf{x}) = \mathbf{W}(\mathbf{x})\mathbf{x} + b(\mathbf{x}) , \quad (31)$$

with $b(\mathbf{x}) \in \mathbb{R}^c$ subsuming all the contributions to the class logits that are not accounted for by the dynamic linear mapping $\mathbf{W}(\mathbf{x})$. As such, B-cos DNNs face the same problem as described by [25]: while explanations based on $\mathbf{W}(\mathbf{x})$ neglect the impact of the contributions from bias terms $b(\mathbf{x})$, in Fig. 13 we show that the bias term can play a critical role in the model decision. In particular, we find that for some models (*e.g.* see the models trained with InstanceNorm or



Fig. 11: **Col. 1:** Ambiguous input image. **Cols. 2+3:** Explanations for the two most likely classes under a B-cos DenseNet-121. **Col. 4:** Difference of contribution maps for the two class logits, i.e., $\Delta(c_1, c_2) = s_{c_1}^L(\mathbf{x}) - s_{c_2}^L(\mathbf{x})$, see Eq. (16); positive values shown in orange (c_1), negative values in blue (c_2).

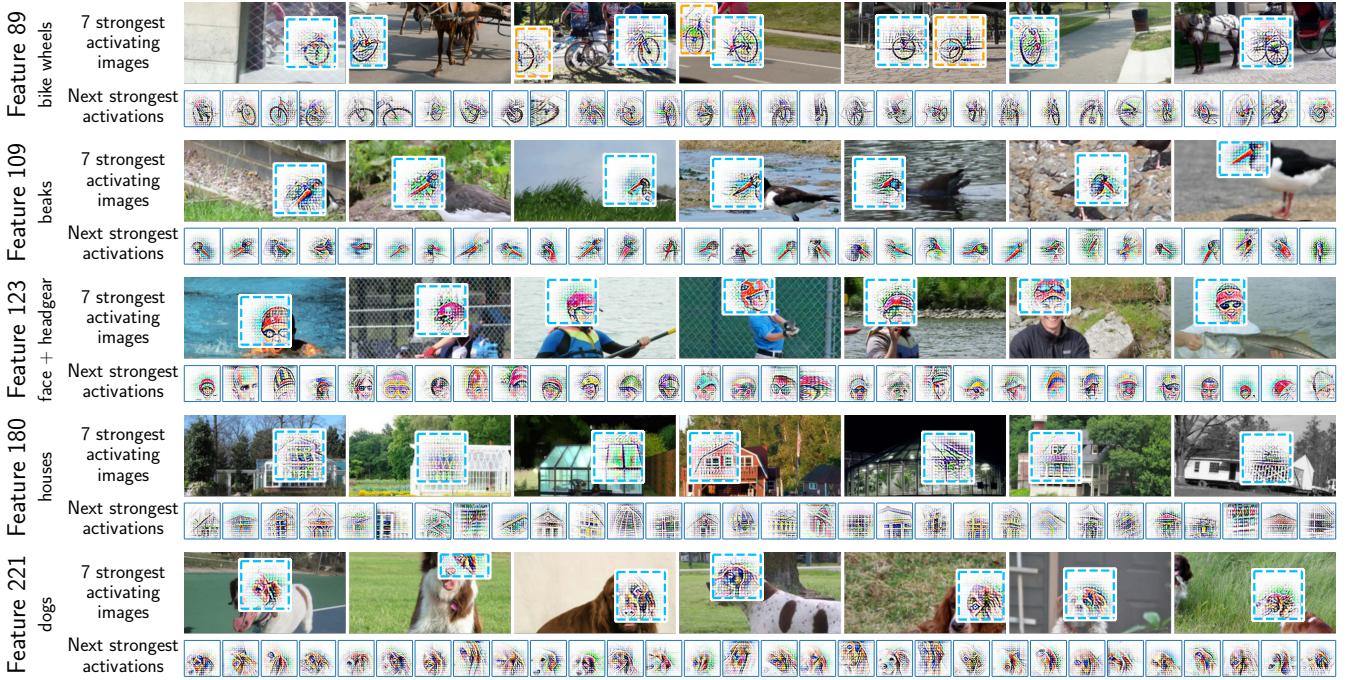


Fig. 12: Explanations of 5 highly contributing feature directions in layer 87 of a DenseNet-121. For each direction, we provide its index number n and a concept description. Further, we show the 7 most activating images for each direction (top row per feature), in which we visualise the explanation for the highest (**blue squares**) activation; i.e., visualise the 72×72 center patch of the weighting $[\mathbf{W}_{1 \rightarrow l}(\mathbf{x})]_n$ for feature n . For some images in the first row, we additionally show the explanation for the 2nd highest activation (**orange squares**). Lastly, we show the explanations of the highest activations (corresponding to the blue squares) for the next 30 images to highlight the features’ specificity.

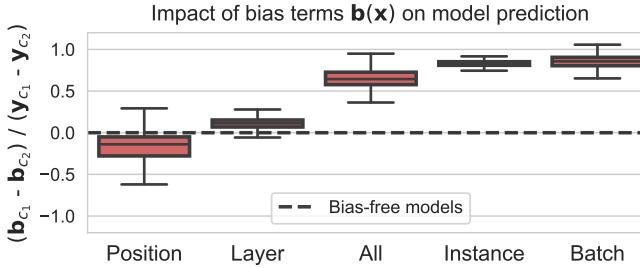


Fig. 13: **Ratio of bias differences to logit differences** for the top-2 predictions (c_1, c_2) of ResNet-20 models on the CIFAR10 test set, trained with different normalisation layers. For some models, the bias term explains most of the difference and thus plays a critical role in the model predictions (e.g., for BatchNorm and InstanceNorm). On the other hand, for the model trained with LayerNorm, the bias term does not seem to be decisive, whereas for PositionNorm it even seems to favour the second highest prediction. To avoid not adequately representing the potentially critical bias $\mathbf{b}(\mathbf{x})$ in the explanations, we ensure $\mathbf{b}(\mathbf{x})=0$ for B-cos CNNs.

BatchNorm), the bias terms make up most of the difference between the top-2 predictions of the model. To ensure that the explanations of B-cos CNNs are complete, we set all bias parameters of the model to zero, see Sec. 3.3.2.

Implicit biases. Even when no *explicit* bias parameters are used, we found that some models learnt to use image edges for computing biases to add to the class logits. While these biases are correctly reflected in the explanations (see Fig. 14,

row 2), this behaviour makes the model explanations less human-interpretable. Interestingly, these models seem to have learnt to solve the optimisation task in Eq. (17) in an unintended manner. Specifically, all classes receive highly positive contributions from image corners or edges (the top edge in Fig. 14), thus yielding high class logits for all classes. To still obtain low scores for the non-target classes, the models then seem to use features of the recognised object in the image to add negative contributions to those classes.

To corroborate this, we additionally visualise the explanations for the mean-corrected logits $\mathbf{y} - \langle \mathbf{y}_c \rangle_c$ in row 3 of Fig. 14; note that given the dynamic linearity of the B-cos models, this is equivalent to computing mean-corrected explanations $\mathbf{W}(\mathbf{x}) - \langle [\mathbf{W}(\mathbf{x})]_c \rangle_c$. As expected, we find that those modified explanations \mathbf{E}' exhibit a high degree of specificity and are thus more easily human-interpretable.

We hypothesise that this behaviour is rooted in the formulation of the optimisation task itself: in particular, the BCE loss ‘asks’ the model to produce highly negative logits for all classes but one, which might bias the model towards using object features to compute negative contributions.

To overcome this, and to incentivise the models to use object features primarily to compute *positive contributions* for the classes this feature belongs to, we propose to change the optimisation task itself, as described in Sec. 3.2.2. As we show in the last row of Fig. 14, this indeed has the intended effect and results in explanations that are highly focused on the class objects in the image.

We believe this to highlight the importance of a holistic approach towards interpretable deep neural networks. In particular, these results show that both the model design as

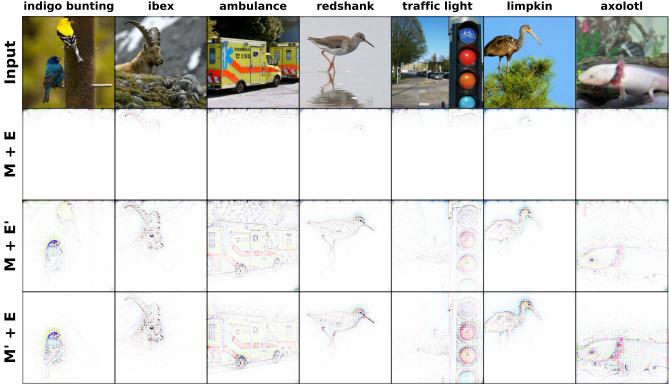


Fig. 14: ‘**You get what you optimise for.**’ We found that for some models M (here, ResNet-18), the explanations E seem to highlight image edges for all samples, implying that most *positive* contributions for all classes come from those edges, see **row 2**. However, the explanations for the difference from the mean logit (i.e., $y_c - \langle y_j \rangle_j$)—denoted by E'—are still qualitatively convincing, see **row 3**, suggesting that the models learn to add a positive bias to all classes, which vanishes in the mean-corrected explanations. To alleviate this, we propose to optimise the models differently (yielding M', see **row 4**) instead of changing the explanations: specifically, by changing the target encoding, as described in Sec. 3.2.2, we encourage the models to focus on positive evidence.

well as the optimisation procedure can significantly impact the model behaviour as well as the resulting explanations, which complicates the development of post-hoc explanations methods that do not take those aspects into account.

6 CONCLUSION

We presented a novel approach for endowing deep neural networks with a high degree of *inherent interpretability*. In particular, we developed the B-cos transformation as a modification of the linear transformation to increase weight-input alignment during optimisation and showed that this can significantly increase interpretability. Importantly, the B-cos transformations can be used as a drop-in replacement for the ubiquitously used linear transformations in conventional DNNs whilst only incurring minor drops in classification accuracy. As such, our approach can increase the interpretability of a wide range of DNNs at a low cost and thus holds great potential to have a significant impact on the deep learning community. In particular, it shows that strong performance and interpretability need not be at odds. Moreover, we demonstrate that by structurally constraining *how* the neural networks are to solve an optimisation task—in the case of B-cos networks via *alignment*—allows for extracting explanations that faithfully reflect the underlying model. We believe this to be an important step on the road towards interpretable deep learning, which is an essential ingredient for building trust in DNN-based decisions, especially in safety-critical situations.

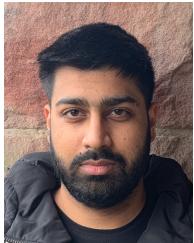
REFERENCES

- [1] W. Samek, G. Montavon, S. Lapuschkin, C. J. Anders, and K.-R. Müller, “Explaining Deep Neural Networks and Beyond: A Review of Methods and Applications,” *Proceedings of the IEEE*, vol. 109, no. 3, 2021.
- [2] V. Nair and G. E. Hinton, “Rectified Linear Units Improve Restricted Boltzmann Machines,” in *International Conference on Machine Learning (ICML)*, 2010.
- [3] G. F. Montúfar, R. Pascanu, K. Cho, and Y. Bengio, “On the Number of Linear Regions of Deep Neural Networks,” in *Advances in Neural Information Processing Systems (NeurIPS)*, 2014.
- [4] M. Yang and B. Kim, “Benchmarking Attribution Methods with Relative Feature Importance,” *CoRR*, vol. abs/1907.09701, 2019.
- [5] A. Shrikumar, P. Greenside, and A. Kundaje, “Learning Important Features Through Propagating Activation Differences,” in *International Conference on Machine Learning (ICML)*, 2017.
- [6] J. Adebayo, J. Gilmer, M. Muelly, I. J. Goodfellow, M. Hardt, and B. Kim, “Sanity Checks for Saliency Maps,” in *Advances in Neural Information Processing Systems (NeurIPS)*, 2018.
- [7] A. Krizhevsky, “Learning multiple layers of features from tiny images,” University of Toronto, Tech. Rep., 2009.
- [8] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly, J. Uszkoreit, and N. Houlsby, “An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale,” in *International Conference on Learning Representations (ICLR)*, 2021.
- [9] K. Simonyan and A. Zisserman, “Very Deep Convolutional Networks for Large-Scale Image Recognition,” in *International Conference on Learning Representations (ICLR)*, 2015.
- [10] 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 (CVPR)*, 2016.
- [11] G. Huang, Z. Liu, L. Van Der Maaten, and K. Q. Weinberger, “Densely Connected Convolutional Networks,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017.
- [12] S. Xie, R. Girshick, P. Dollár, Z. Tu, and K. He, “Aggregated residual transformations for deep neural networks,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017.
- [13] Z. Liu, H. Mao, C.-Y. Wu, C. Feichtenhofer, T. Darrell, and S. Xie, “A convnet for the 2020s,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022.
- [14] M. Böhle, M. Fritz, and B. Schiele, “B-cos Networks: Alignment is All we Need for Interpretability,” in *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022.
- [15] S. M. Lundberg and S. Lee, “A Unified Approach to Interpreting Model Predictions,” in *Advances in Neural Information Processing Systems (NeurIPS)*, 2017.
- [16] V.Petsiuk, A. Das, and K. Saenko, “RISE: Randomized Input Sampling for Explanation of Black-box Models,” in *British Machine Vision Conference (BMVC)*, 2018.
- [17] M. T. Ribeiro, S. Singh, and C. Guestrin, ““Why Should I Trust You?”: Explaining the Predictions of Any Classifier,” in *International Conference on Knowledge Discovery and Data Mining (SIGKDD)*, 2016.
- [18] B. Kim, M. Wattenberg, J. Gilmer, C. J. Cai, J. Wexler, F. B. Viégas, and R. Sayres, “Interpretability Beyond Feature Attribution: Quantitative Testing with Concept Activation Vectors (TCAV),” in *International Conference on Machine Learning (ICML)*, 2018.
- [19] S. Das, P. Xu, Z. Dai, A. Endert, and L. Ren, “Interpreting Deep Neural Networks through Prototype Factorization,” in *International Conference on Data Mining Workshops (ICDMW)*, 2020.
- [20] K. Simonyan, A. Vedaldi, and A. Zisserman, “Deep Inside Convolutional Networks: Visualising Image Classification Models and Saliency Maps,” in *International Conference on Learning Representations (ICLR)*, Workshop, 2014.
- [21] J. T. Springenberg, A. Dosovitskiy, T. Brox, and M. A. Riedmiller, “Striving for Simplicity: The All Convolutional Net,” in *International Conference on Learning Representations (ICLR)*, Workshop, 2015.
- [22] B. Zhou, A. Khosla, Á. Lapedriza, A. Oliva, and A. Torralba, “Learning Deep Features for Discriminative Localization,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.
- [23] S. Bach, A. Binder, G. Montavon, F. Klauschen, K.-R. Müller, and W. Samek, “On Pixel-Wise Explanations for Non-Linear Classifier Decisions by Layer-Wise Relevance Propagation,” *PLoS ONE*, 2015.
- [24] R. R. Selvaraju, M. Cogswell, A. Das, R. Vedantam, D. Parikh, and D. Batra, “Grad-CAM: Visual Explanations from Deep Networks

- via Gradient-Based Localization," in *International Conference on Computer Vision (ICCV)*, 2017.
- [25] S. Srinivas and F. Fleuret, "Full-Gradient Representation for Neural Network Visualization," in *Advances in Neural Information Processing Systems (NeurIPS)*, 2019.
- [26] M. Sundararajan, A. Taly, and Q. Yan, "Axiomatic Attribution for Deep Networks," in *International Conference on Machine Learning (ICML)*, 2017.
- [27] C. Chen, O. Li, D. Tao, A. Barnett, C. Rudin, and J. Su, "This Looks Like That: Deep Learning for Interpretable Image Recognition," in *Advances in Neural Information Processing Systems (NeurIPS)*, 2019.
- [28] W. Brendel and M. Bethge, "Approximating CNNs with Bag-of-local-Features models works surprisingly well on ImageNet," in *International Conference on Learning Representations (ICLR)*, 2019.
- [29] M. Böhle, M. Fritz, and B. Schiele, "Convolutional Dynamic Alignment Networks for Interpretable Classifications," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021.
- [30] D. Tsipras, S. Santurkar, L. Engstrom, A. Turner, and A. Madry, "Robustness May Be at Odds with Accuracy," in *International Conference on Learning Representations (ICLR)*, 2019.
- [31] H. Shah, P. Jain, and P. Netrapalli, "Do Input Gradients Highlight Discriminative Features?" *CoRR*, vol. abs/2102.12781, 2021.
- [32] B. Kim, J. Seo, and T. Jeon, "Bridging Adversarial Robustness and Gradient Interpretability," *Safe Machine Learning workshop at International Conference on Learning Representations (ICLR)*, 2019.
- [33] S. Srinivas and F. Fleuret, "Rethinking the Role of Gradient-based Attribution Methods for Model Interpretability," in *International Conference on Learning Representations (ICLR)*, 2021.
- [34] P.-J. Kindermans, K. T. Schütt, M. Alber, K.-R. Müller, D. Erhan, B. Kim, and S. Dähne, "Learning how to explain neural networks: PatternNet and PatternAttribution," in *International Conference on Learning Representations (ICLR)*, 2018.
- [35] G. Zoumpourlis, A. Doumanoglou, N. Vretos, and P. Daras, "Non-linear convolution filters for CNN-based learning," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017.
- [36] W. Liu, Y. Zhang, X. Li, Z. Liu, B. Dai, T. Zhao, and L. Song, "Deep Hyperspherical Learning," in *Advances in Neural Information Processing Systems (NeurIPS)*, 2017.
- [37] W. Liu, Z. Liu, Z. Yu, B. Dai, R. Lin, Y. Wang, J. M. Rehg, and L. Song, "Decoupled Networks," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018.
- [38] C. Luo, J. Zhan, X. Xue, L. Wang, R. Ren, and Q. Yang, "Cosine normalization: Using cosine similarity instead of dot product in neural networks," in *International Conference on Artificial Neural Networks (ICANN)*, 2018.
- [39] K. Ghiasi-Shirazi, "Generalizing the convolution operator in convolutional neural networks," *Neural Processing Letters*, vol. 50, no. 3, 2019.
- [40] C. Wang, J. Yang, L. Xie, and J. Yuan, "Kervolutional neural networks," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019.
- [41] S. Wang, T. Zhou, and J. Bilmes, "Bias also matters: Bias attribution for deep neural network explanation," in *International Conference on Machine Learning (ICML)*, 2019.
- [42] S. Mohan, Z. Kadkhodaie, E. P. Simoncelli, and C. Fernandez-Granda, "Robust And Interpretable Blind Image Denoising Via Bias-Free Convolutional Neural Networks," in *International Conference on Learning Representations (ICLR)*, 2020.
- [43] R. Hesse, S. Schaub-Meyer, and S. Roth, "Fast Axiomatic Attribution for Neural Networks," in *Advances in Neural Information Processing Systems (NeurIPS)*, vol. 34, 2021.
- [44] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, "Attention is all you need," *Advances in Neural Information Processing Systems (NeurIPS)*, 2017.
- [45] S. Serrano and N. A. Smith, "Is Attention Interpretable?" in *Proceedings of the Annual Meeting of the Association for Computational Linguistics (ACL)*, 2019.
- [46] S. Abnar and W. Zuidema, "Quantifying Attention Flow in Transformers," in *Proceedings of the Annual Meeting of the Association for Computational Linguistics (ACL)*, 2020.
- [47] J. Bastings and K. Filippova, "The elephant in the interpretability room: Why use attention as explanation when we have saliency methods?" in *Proceedings of the Third BlackboxNLP Workshop on Analyzing and Interpreting Neural Networks for NLP*, 2020.
- [48] O. Barkan, E. Hauon, A. Caciularu, O. Katz, I. Malkiel, O. Armstrong, and N. Koenigstein, "Grad-SAM: Explaining Transformers via Gradient Self-Attention Maps," in *Proceedings of the International Conference on Information and Knowledge Management (CIKM)*, 2021.
- [49] H. Chefer, S. Gur, and L. Wolf, "Transformer interpretability beyond attention visualization," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021.
- [50] P. Ramachandran, B. Zoph, and Q. V. Le, "Searching for Activation Functions," in *International Conference on Learning Representations (ICLR) Workshop*, 2018.
- [51] I. Goodfellow, D. Warde-Farley, M. Mirza, A. Courville, and Y. Bengio, "Maxout networks," in *International Conference on Machine Learning (ICML)*, 2013.
- [52] S. Ioffe and C. Szegedy, "Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift," in *International Conference on Machine Learning (ICML)*, 2015.
- [53] J. L. Ba, J. R. Kiros, and G. E. Hinton, "Layer normalization," *CoRR*, vol. abs/1607.06450, 2016.
- [54] B. Li, F. Wu, K. Q. Weinberger, and S. Belongie, "Positional normalization," *Advances in Neural Information Processing Systems*, vol. 32, 2019.
- [55] D. Ulyanov, A. Vedaldi, and V. Lempitsky, "Instance normalization: The missing ingredient for fast stylization," *CoRR*, vol. abs/1607.08022, 2016.
- [56] Y. Wu and K. He, "Group normalization," in *Proceedings of the European conference on computer vision (ECCV)*, 2018.
- [57] J. Deng, W. Dong, R. Socher, L. Li, K. Li, and L. Fei-Fei, "ImageNet: A large-scale hierarchical image database," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2009.
- [58] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga, A. Desmaison, A. Kopf, E. Yang, Z. DeVito, M. Raison, A. Tejani, S. Chilamkurthy, B. Steiner, L. Fang, J. Bai, and S. Chintala, "PyTorch: An Imperative Style, High-Performance Deep Learning Library," in *Advances in Neural Information Processing Systems (NeurIPS)*, 2019.
- [59] L. Beyer, X. Zhai, and A. Kolesnikov, "Better plain ViT baselines for ImageNet-1k," *CoRR*, vol. abs/2205.01580, 2022.
- [60] T. Xiao, M. Singh, E. Mintun, T. Darrell, P. Dollár, and R. Girshick, "Early convolutions help transformers see better," *Advances in Neural Information Processing Systems (NeurIPS)*, vol. 34, 2021.
- [61] "Torchvision library, pretrained models," <https://pytorch.org/vision/stable/models.html>, accessed: 2021-11-11.
- [62] "How to Train State-Of-The-Art Models Using TorchVision's Latest Primitives," <https://pytorch.org/blog/how-to-train-state-of-the-art-models-using-torchvision-latest-primitives/>, 2021, accessed: 2023-05-23.
- [63] D. Baehrens, T. Schroeter, S. Harmeling, M. Kawanabe, K. Hansen, and K.-R. Müller, "How to explain individual classification decisions," *The Journal of Machine Learning Research (JMLR)*, 2010.
- [64] S. Rao, M. Böhle, and B. Schiele, "Towards better understanding attribution methods," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022.



Moritz Böhle is a Ph.D. student in Computer Science at the Max Planck Institute for Informatics, working with Prof. Dr. Bernt Schiele and Prof. Dr. Mario Fritz. He graduated with a bachelor's degree in physics in 2016 from the Freie Universität Berlin and obtained his Master's degree in computational neuroscience in 2019 from the Technische Universität Berlin. His research focuses on understanding the 'decision process' in deep neural networks and designing inherently interpretable neural network models.



Navdeppal Singh is a Master's student pursuing Computer Science at Saarland University, where he also obtained his Bachelor's degree in Computer Science in 2021. For his bachelor thesis, he worked on the robustness evaluation of inherently interpretable neural network models. He's interested in computer vision, interpretable and robust neural networks, and neural search.



Mario Fritz Mario Fritz is faculty member at the CISPA Helmholtz Center for Information Security, honorary professor at the Saarland University, and a fellow of the European Laboratory of Learning and Intelligent Systems (ELLIS). Before, he was senior researcher at the Max Planck Institute for Informatics, PostDoc at UC Berkeley and International Computer Science Institute. He studied computer science at the university Erlangen/Nuremberg and obtained his PhD from TU Darmstadt. His current work is

centered around Trustworthy Information Processing with a focus on the intersection of AI & Machine Learning with Security & Privacy. He is associate editor of IEEE TPAMI, and a leading scientist of the Helmholtz Medical Security, Privacy, and AI Research Center, where he is coordinating projects on privacy and federated learning in health. He has over 100 publications, including 80 in top-tier journals (IJCV, TPAMI) and conferences (NeurIPS, AAAI, IJCAI, ICLR, NDSS, USENIX Security, CCS, S&P, CVPR, ICCV, ECCV).



Bernt Schiele has been Max Planck Director at MPI for Informatics and Professor at Saarland University since 2010. He studied computer science at the University of Karlsruhe, Germany. He worked on his master thesis in the field of robotics in Grenoble, France, where he also obtained the "diplome d'études approfondies d'informatique". In 1994 he worked in the field of multi-modal human-computer interfaces at Carnegie Mellon University, Pittsburgh, PA, USA in the group of Alex Waibel. In 1997 he

obtained his PhD from INP Grenoble, France under the supervision of Prof. James L. Crowley in the field of computer vision. The title of his thesis was "Object Recognition using Multidimensional Receptive Field Histograms". Between 1997 and 2000 he was postdoctoral associate and Visiting Assistant Professor with the group of Prof. Alex Pentland at the Media Laboratory of the Massachusetts Institute of Technology, Cambridge, MA, USA. From 1999 until 2004 he was Assistant Professor at the Swiss Federal Institute of Technology in Zurich (ETH Zurich). Between 2004 and 2010 he was Full Professor at the computer science department of TU Darmstadt. He is fellow of IEEE, ACM, ELLIS, and IAPR.

B-cos Alignment for Inherently Interpretable CNNs and Vision Transformers

Supplementary Material

Moritz Böhle
 MPI for Informatics
 Saarland Informatics Campus

Navdeepal Singh
 MPI for Informatics
 Saarland Informatics Campus

Mario Fritz
 CISPA Helmholtz Center
 for Information Security

Bernt Schiele
 MPI for Informatics
 Saarland Informatics Campus

Table of Contents

In this supplement to our work on B-cos DNNs, we provide:

(A) Additional qualitative results	2
In this section, we show additional <i>qualitative</i> results of the model-inherent explanations. Specifically, we show and discuss explanations for various model architectures evaluated on the same set of images.	
(B) Additional quantitative results	3
In this section, we show additional <i>quantitative</i> results. In particular, we present the localisation metric results for two additional B-cos networks as well as those of the pre-trained conventional DNNs.	
(C) Implementation details	3
In this section, we describe the training and evaluation procedure for our experiments in more detail.	
(D) Additional derivations and discussions	5
In this section, we provide a short derivation for Eq. (3) → Eqs. (4)+(9). Further, we provide a more detailed explanation of the relevance of the image encoding for visualising the linear transformations $\mathbf{W}_{1 \rightarrow l}(\mathbf{x})$.	

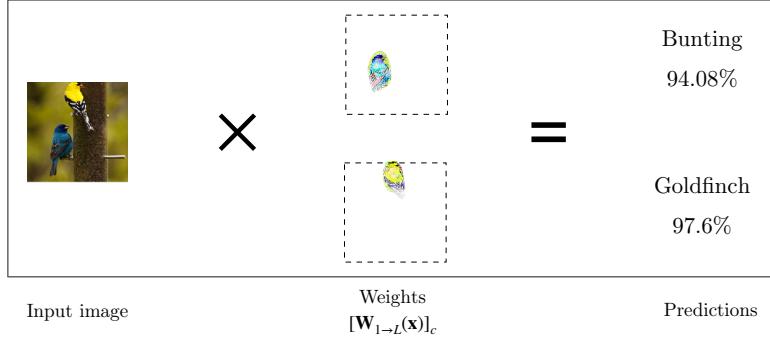


Fig. A1. Illustration of the computations of a B-cos network. For a given input image (left), the model computes an *input-dependent* linear transform $\mathbf{W}_{1 \rightarrow L}(\mathbf{x})$ (center). The scalar product between the input and the weights $[\mathbf{W}_{1 \rightarrow L}(\mathbf{x})]_c$ for class c (row c of $\mathbf{W}_{1 \rightarrow L}(\mathbf{x})$), yields the class logits for the respective class. To obtain class probabilities (right), we apply the sigmoid function. Since the B-cos networks are trained with the BCE loss, they produce probabilities *per class* and *not a probability distribution over classes*. Thus, the probabilities do not sum to 1. For illustration purposes, we only visualise the positive contributions according to $\mathbf{W}_{1 \rightarrow L}(\mathbf{x})$.

A. Additional qualitative examples

In Fig. A1, we illustrate how the linear mappings $\mathbf{W}_{1 \rightarrow L}(\mathbf{x})$ are used to compute the outputs of B-cos networks. In particular, with this we would like to highlight that these linear mappings do not only constitute qualitatively convincing visualisations. Instead, they are in fact based on the actual linear transformation matrix that the model effectively applies to the input to compute its outputs and thus constitute an accurate summary of the model computations.

Comparison between model explanations. In Figs. A2 and A3, we compare explanations for different network architec-

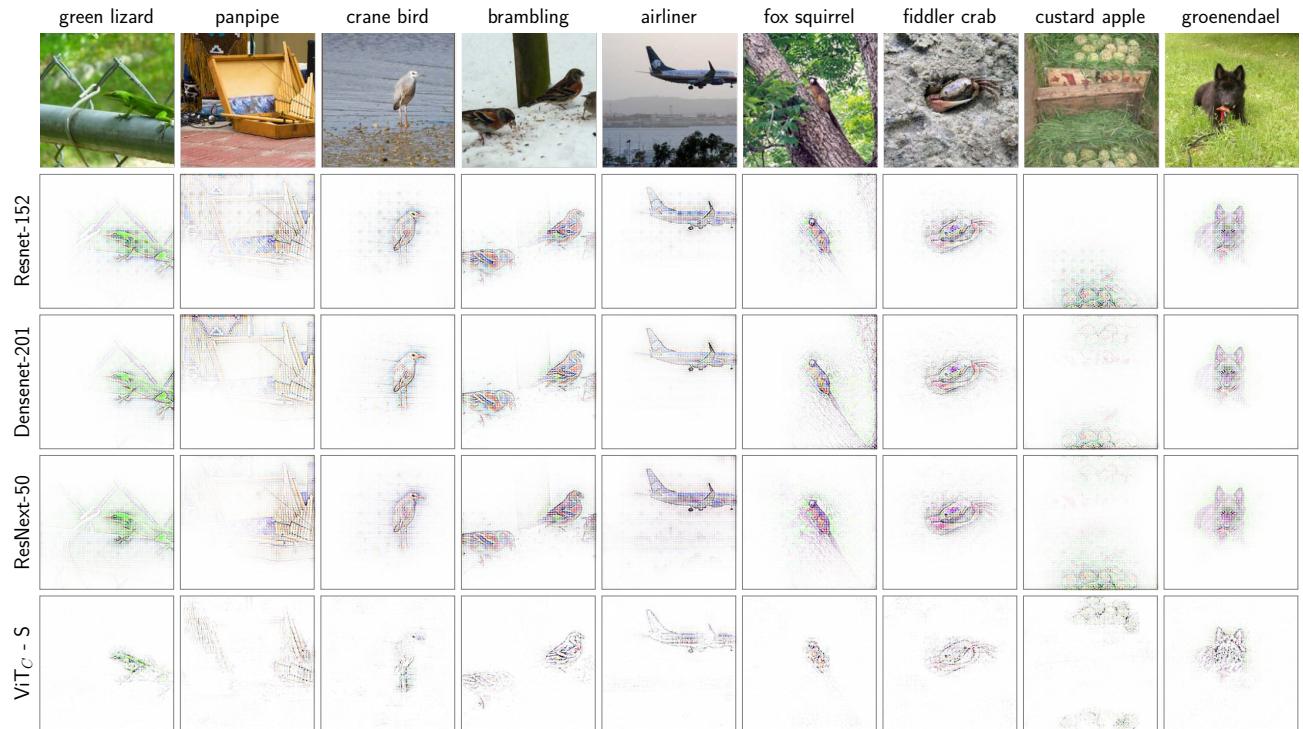


Fig. A2. **Additional model explanations** based on $\mathbf{W}(\mathbf{x})$ for various B-cos architectures, see also Fig. A3. While the explanations are generally consistent between models, architecture-specific details can be observed. E.g., ResNet-based architectures exhibit a grid-pattern in the explanations, and we find that the ViT models generally produce sparser explanations.

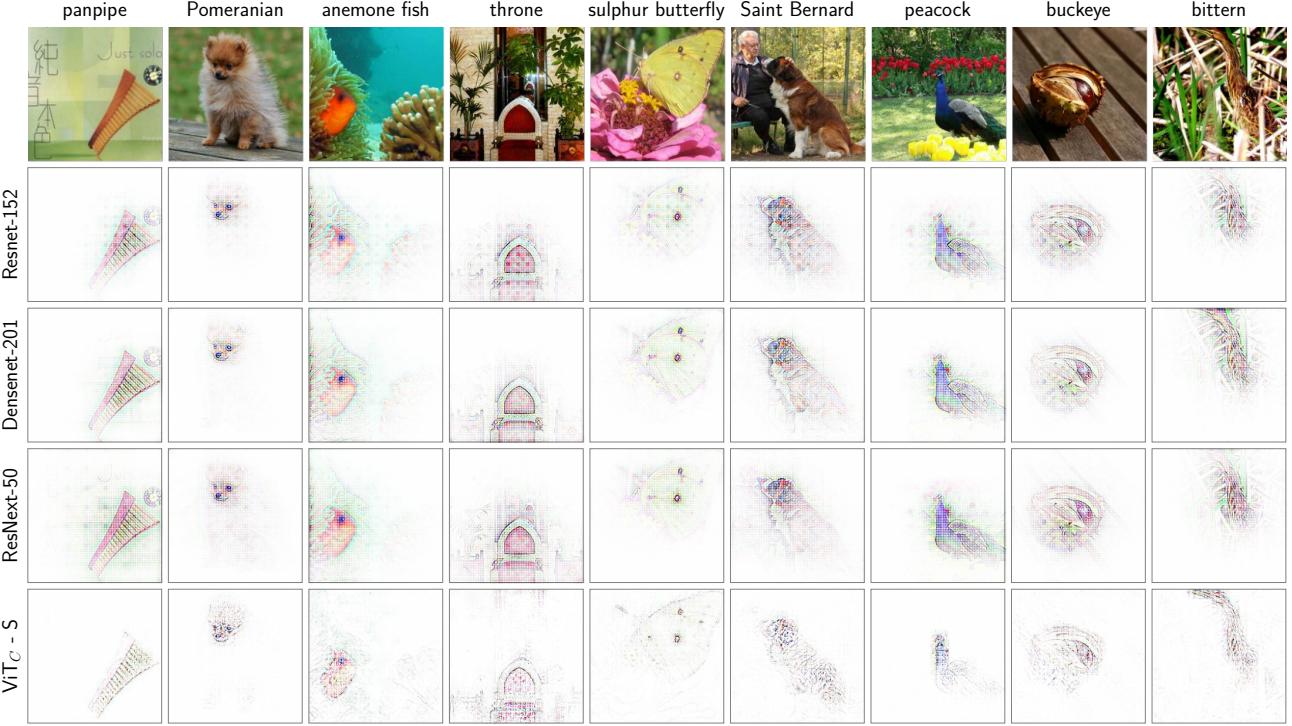


Fig. A3. Model explanations based on $\mathbf{W}(\mathbf{x})$ for various B-cos architectures, see also Fig. A2. While the explanations are generally consistent between models, architecture-specific details can be observed. E.g., ResNet-based architectures exhibit a grid-pattern in the explanations, and we find that the ViT models generally produce sparser explanations.

tures on the same set of images. We find that all models generally seem to focus on the same regions in the image. However, we also observe architecture-specific details: *e.g.*, we find that the vision transformer yield much sparser explanations than the convolutional models and ResNet models (ResNet-152 and ResNext-50 in the figures) exhibit grid-pattern artefacts, which we attribute to the downsampling mechanism in those models. Considering that information is not processed as locally in the ViT_C models as in the convolutional models, these qualitative results provide further evidence for the model-faithfulness of the explanations based on $\mathbf{W}(\mathbf{x})$.

B. Additional quantitative evaluations

In Fig. B1, we present the localisation results of the grid pointing game [S5] for two additional model architectures, see also Fig. 6 in the main paper.

In particular, in Fig. B1 (left), we show that of all methods, the best explanation for the B-cos models is given by the model-inherent linear transformations $\mathbf{W}_{1 \rightarrow L}(\mathbf{x})$. Moreover, in Fig. B1 (right) we can see that the B-cos explanations yield a significant interpretability gain over the respective baselines: specifically, we see that no method explains the baseline models better than the model-inherent linear transformations explain the respective B-cos network.

Finally, as discussed in the main paper (see Fig. 7), the results of the B-cos explanations can be further improved by including only the top-n pixels (here $\approx 11.3k$) in the explanations, which we denote by OursQ in the figure.

C. Implementation details

In the following, we provide further implementation details regarding the training and evaluation procedure (C.1) and the post-hoc attribution methods (C.2).

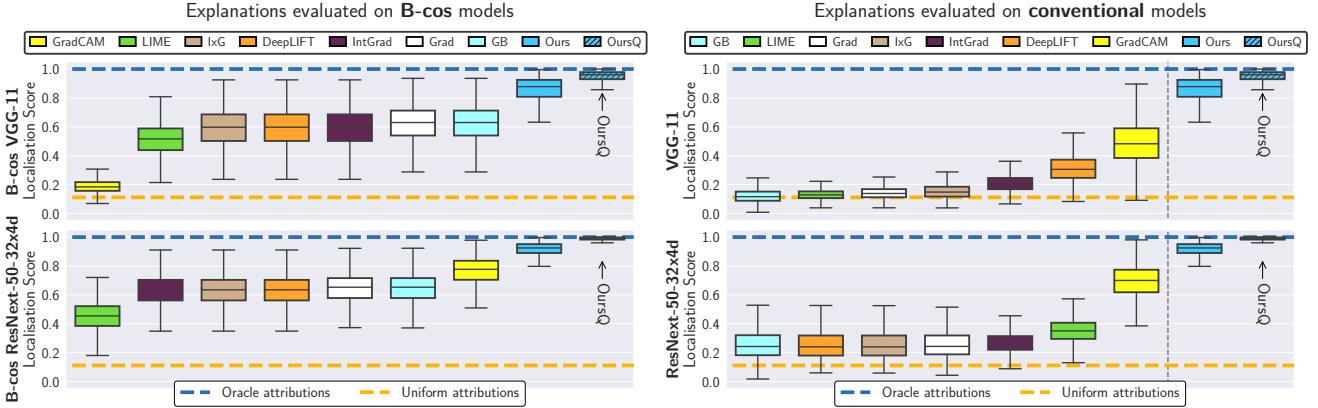


Fig. B1. Additional localisation results. We show localisation results of two additional B-cos models (left) and the conventional counterparts (right). These results are consistent with those shown in Fig. 6 in the main paper. Specifically, we find the model-inherent explanations (Ours) of the B-cos models to yield close to optimal localisation results, especially when only including the top 11.3k pixels in the explanations (OursQ). Additionally, and in accordance with the results reported in the main paper for a DenseNet and a ResNet model, we find that the B-cos models are better explained by post-hoc explanation methods than the conventional models, except for GradCAM on VGG. For VGG, the gradient computation of GradCAM seems to be significantly affected by the non-linear aspects of the B-cos transformations (note that VGG is the only model using a multi-layer perceptron for the classification head).

C.1. Training and evaluation procedure

C.1.1 CIFAR10

Architecture. For our simple B-cos models trained on CIFAR10, we used a 9-layer architecture with the following specifications: kernel size $k = [3, 3, 3, 3, 3, 3, 3, 3, 1]$, stride $s = [1, 1, 2, 1, 1, 2, 1, 1, 1]$, padding $p = [1, 1, 1, 1, 1, 1, 1, 1, 0]$, and output channels $o = [64, 64, 128, 128, 128, 256, 256, 256, 10]$ for layers $l = [1, 2, 3, 4, 5, 6, 7, 8, 9]$ respectively.

When increasing the parameter B , we observed the input signal to decay strongly over the network layers, which resulted in zero outputs and hindered training. To overcome this, we scaled all layer outputs with a fixed scalar γ to improve signal propagation, which we set such that $\log_{10} \gamma = 1.75 + B/10$. To counteract the artificial upscaling of the signal at the network output, we down-scaled the network output by a fixed constant T for each B , such that $\log_{10} T = [8.9, 8.125, 7.35, 6.757, 4.8, 4.525, 4.25]$ for $B = [1, 1.25, 1.5, 1.75, 2, 2.25, 2.5]$ respectively.

Furthermore, for our B-cos ResNet models for CIFAR10 we converted the ResNet-{20,32,44,56} [S7] models to their corresponding B-cos version. Specifically, to evaluate the applicability of different normalisation layers, we replaced the conventionally used Batch-Norm layers by Batch-, Layer-, Instance-, Position-, and All-Norm layers (Eqs. (20)–(24)). As we report in Fig. 13 in the main paper, for some of the normalisation layers (*e.g.* All-, Instance-, and Batch-Norm), the bias terms can play a significant role in the classification decision. Therefore, for the ImageNet experiments, we exclusively use *bias-free* normalisation layers. Furthermore, we switch the order of the the global pooling layer and last linear classifier layer, *i.e.* we first do a linear classification (using a convolution) before doing a global (average) pooling, in order to more easily evaluate those models in the localisation metric, for which the size of the input image changes. Note that for models with a linear classification head, this yields functionally equivalent models; for the B-cos models, however, the ordering matters, since the classification head is a non-linear (*i.e.* B-cos) layer.

Training. We trained all our CIFAR10 models for 100 epochs with the Adam optimiser [S9], an initial learning rate of 1×10^{-3} , and a batch size of 64. Further, we used a cosine learning rate schedule and decayed the learning rate to 1×10^{-5} for our simple B-cos models and to 0 for our B-cos ResNets over 100 epochs. For augmenting the data, we applied horizontal flipping and padded random cropping. We used a bias term of $\mathbf{b} = \log(0.1/0.9)$, which yields a uniform probability distribution for zero inputs ($[\mathbf{f}(\mathbf{x} = \mathbf{0})]_i = [\sigma(\mathbf{W}_{1 \rightarrow L} \mathbf{0} + \mathbf{b})]_i = 0.1 \forall i$).

C.1.2 ImageNet

Training. Our training recipe is based on the `torchvision` (<https://github.com/pytorch/vision/>) reference recipes. For our ResNets, DenseNets, and VGGs, we train for 90 epochs with the Adam optimiser using a learning rate of

1×10^{-3} and a batch size of 256 (distributed over 4 GPUs with each having batch size of 64). Following `torchvision`, we also train our ResNeXt model with the same recipe as for the ResNets but for 100 epochs instead of 90. As mentioned in Sec. 3.2.2 in the main paper, we change the target encoding \mathbf{y}_i for non-target classes to be $1/C$. For the learning rate schedule, we first use a linear warm-up for 5 epochs, then we use a cosine annealing learning schedule, decaying the learning rate to 0. To avoid divergence issues we employ gradient clipping. However, we empirically noticed that the magnitude of the gradient for B-cos models tends to be significantly lower than that of the gradients for standard networks. To avoid manually tuning the clipping threshold, we opt for employing Adaptive Gradient Clipping (AGC) [S4]. For data augmentation, we use random cropping and then resizing with bilinear interpolation to size 224 along with random horizontal flips with a probability of 0.5.

Inspired by `torchvision`'s new longer training recipes¹, we additionally train a set of ConvNeXt-{Tiny, Base}, DenseNet121, and ResNet{50, 152} models for 600 epochs with the Adam optimiser using a learning rate of 1×10^{-3} and a batch size of 1024 (distributed over 8 GPUs with each having a batch size of 128). We keep an exponential moving average (EMA) of the weights of the model with a decay rate of 0.99998 updated after every 32 training steps. We use the same learning rate scheduler as mentioned above and AGC for gradient clipping. For data augmentation, we use random cropping and then resizing to size 176 with bilinear interpolation, MixUp [S17] with $\alpha = 0.2$, CutMix [S16] with $\alpha = 1.0$, Repeated Augmentations [S2, S8] with 4 repetitions, Random Erasing [S6, S18] with a probability of 0.1, horizontal flipping with a probability of 0.5, and finally TrivialAugment [S10] which is parameter-free.

We trained our ViTs from scratch following the 90 epoch protocol described by [S3]. Given that all the B-cos layer weights are normalised to unit norm, we employ Adam instead of AdamW, since no weight decay is necessary. Furthermore, as mentioned previously we employ AGC over standard gradient clipping. We used an effective batch size of 1024 (distributed over 8 GPUs) for all but our largest ViTs (ViT-L models and the Simple ViT-B model), for which we use a batch size of 512 due to memory constraints. We observed these larger ViT models to diverge at the beginning of training and therefore opted to increase the length of the warm-up period from the 10 000 steps described by [S3] to 50 000 steps.

For all further details, we kindly refer the reader to the code at <https://github.com/B-cos/B-cos-v2>, which we make available along with this submission.

Evaluation. We evaluated all networks on the full images of the ImageNet validation set, after rescaling the smaller dimension (height / width) to 256 using bilinear interpolation. For networks with EMA weights, we evaluate both the EMA weights non-EMA weights and report the best performing one.

C.2. Attribution methods

We compare the model-inherent explanations, given by the linear transformation $\mathbf{W}_{1 \rightarrow L}(\mathbf{x})$, against the following post-hoc attribution methods: the vanilla gradient (Grad, [S1]), ‘Input×Gradient’ (IxG, [S13]), Guided Backpropagation’ (GB, [S14]), Integrated Gradients (IntGrad, [S15]), DeepLIFT ([S13]), GradCAM ([S12]), and LIME ([S11]).

For all methods except LIME, we rely on the captum library (github.com/pytorch/captum). For IntGrad, we set `n_steps` = 50 for integrating over the gradients. For LIME, we used the official implementation available at github.com/marotcr/lime, with 500 samples and the default values for the kernel size ($k = 4$) and evaluate the weights assigned to the individual image segments, see Fig. 9 in the main paper for a visualisation.

C.2.1 Localisation metric

We evaluated all attribution methods on the grid pointing game [S5]. For this, we constructed 500 3×3 grid images. For an example of a 2×2 grid, see Fig. 3 in the main paper. As was done in [S5], we sorted the images according to the models’ classification confidence for each class and then sampled a random set of classes for each multi-image. For each of the sampled classes, we then included the most confidently classified image in the grid that had not already been used in a previous grid image.

D. Additional derivations and discussions

D.1. The B-cos transformation as a scaled linear transformation

In the following, we provide additional details on how to express the B-cos transformation as a scaled linear transformation and in matrix form.

¹<https://pytorch.org/blog/how-to-train-state-of-the-art-models-using-torchvision-latest-primitives/>

As described in Eq. (3) in the main paper, the B-cos transformation is given by

$$B\text{-cos}(\mathbf{x}; \mathbf{w}) = \|\widehat{\mathbf{w}}\| \|\mathbf{x}\| \times |c(\mathbf{x}, \widehat{\mathbf{w}})|^B \times \text{sgn}(c(\mathbf{x}, \widehat{\mathbf{w}})) , \quad (\text{D.1})$$

$$\text{with } c(\mathbf{x}, \mathbf{w}) = \cos(\angle(\mathbf{x}, \mathbf{w})) , \quad (\text{D.2})$$

$$\widehat{\mathbf{w}} = \mathbf{w}/\|\mathbf{w}\| , \quad (\text{D.3})$$

$\angle(\mathbf{x}, \mathbf{w})$ returning the angle between \mathbf{x} and \mathbf{w} , and sgn the sign function. Note that the sign function can be expressed as $\text{sgn}(a)=a/|a|$ for $|a| \neq 0$ and zero otherwise. Hence, Eq. (D.1) can be expressed as

$$B\text{-cos}(\mathbf{x}; \mathbf{w}) = \|\widehat{\mathbf{w}}\| \|\mathbf{x}\| \times |c(\mathbf{x}, \widehat{\mathbf{w}})|^B \times \text{sgn}(c(\mathbf{x}, \widehat{\mathbf{w}})) \quad (\text{D.4})$$

$$(\text{replace sgn}) = \|\widehat{\mathbf{w}}\| \|\mathbf{x}\| \times |c(\mathbf{x}, \widehat{\mathbf{w}})|^B \times c(\mathbf{x}, \widehat{\mathbf{w}})/|c(\mathbf{x}, \widehat{\mathbf{w}})| \quad (\text{D.5})$$

$$(\text{combine cos terms}) = \|\widehat{\mathbf{w}}\| \|\mathbf{x}\| \times |c(\mathbf{x}, \widehat{\mathbf{w}})|^{B-1} \times c(\mathbf{x}, \widehat{\mathbf{w}}) \quad (\text{D.6})$$

$$(\text{reorder}) = \|\widehat{\mathbf{w}}\| \|\mathbf{x}\| \times c(\mathbf{x}, \widehat{\mathbf{w}}) \times |c(\mathbf{x}, \widehat{\mathbf{w}})|^{B-1} \quad (\text{D.7})$$

$$(\text{write first three factors as linear transform}) = \widehat{\mathbf{w}}^T \mathbf{x} \times |c(\mathbf{x}, \widehat{\mathbf{w}})|^{B-1} . \quad (\text{D.8})$$

For clarity, we marked the changes between lines in the above equations in red.

From Eq. (D.8) it becomes clear that a B-cos transformation simply computes a rescaled linear transform, as in Eq. (4) in the main paper. Thus, multiple units in parallel (i.e., a layer \mathbf{l}^* of B-cos units) can easily be expressed in matrix form via

$$\mathbf{l}^*(\mathbf{x}) = |c(\mathbf{x}, \widehat{\mathbf{W}})|^{B-1} \times \widehat{\mathbf{W}} \mathbf{x} . \quad (\text{D.9})$$

Here, the \times , \cos , and absolute value operators are applied element-wise and the rows of $\widehat{\mathbf{W}}$ are given by $\widehat{\mathbf{w}}_n$ of the individual units n . Hence, the output of each unit (entry in output vector \mathbf{l}^*) is the down-scaled linear transformation from Eq. (D.8). Note that Eq. (D.9) is the same as Eq. (10) in the main paper.

D.2. On the relevance of image encoding for the visualisations

As we describe in the main paper, we encode image pixels as $[r, g, b, 1-r, 1-g, 1-b]$. This has two important advantages.

On the one hand, as argued by [S5], this overcomes a bias towards bright pixels. For this, note that the model output is computed as a linear transformation of the input \mathbf{x} . As such, the contribution to the output per pixel is given by the weighted input strength. In particular, a specific pixel location (i, j) with color channels c contributes $\sum_c w_{(i,j,c)} x_{(i,j,c)}$ to the output. Under the conventional encoding—i.e., $[r, g, b]$ —, a black pixel is encoded by $x_{(i,j,c)}=0$ for $c \in \{1, 2, 3\}$ and can therefore not contribute to the model output. Since we train the model to maximise its outputs (binary cross entropy loss, see Sec. 3.2.2 in the main paper), the network will preferentially encode bright pixels, as these can produce higher contributions for maximising the output than dark pixels. In contrast, under the new encoding dark and bright pixels have the same amount of ‘signal’ that can be weighted, i.e., $\sum_c x_{(i,j,c)}=3 \forall (i, j)$.

Moreover, this encoding allows to unambiguously infer the color of a pixel solely based on the angle of the pixel vector $[r, g, b, 1-r, 1-g, 1-b]$. To contrast this with the original encoding, consider a pixel that is (almost) completely black and given by $[r, g, b]$ with $r=0, g=0, b=0.001$. This pixel has the same angle as a red pixel, given by $r=1, g=0, b=0$. Thus, these two colors cannot be disambiguated based on their angle. By adding the three additional color channels $[1-r, 1-g, 1-b]$, each color is uniquely encoded by the direction of the pixel vector. Finally, note that the B-cos transformation induces an alignment pressure on the weights, i.e., the model weights are optimised such that $\mathbf{W}_{1 \rightarrow L}$ points in the same direction as (important features in) the input. Consequently, the weights will reproduce the *angles* of the pixels, but there is no constraint on their *norm*. Since the angle is sufficient for inferring the color, we can nevertheless decode the angles of the weight vectors into RGB colors.

References

- [S1] David Baehrens, Timon Schroeter, Stefan Harmeling, Motoaki Kawanabe, Katja Hansen, and Klaus-Robert Müller. How to explain individual classification decisions. *The Journal of Machine Learning Research (JMLR)*, 2010. 5
- [S2] Maxim Berman, Hervé Jégou, Andrea Vedaldi, Iasonas Kokkinos, and Matthijs Douze. MultiGrain: a unified image embedding for classes and instances. *CoRR*, abs/1902.05509, 2019. 5
- [S3] Lucas Beyer, Xiaohua Zhai, and Alexander Kolesnikov. Better plain ViT baselines for ImageNet-1k. *CoRR*,

abs/2205.01580, 2022. 5

- [S4] Andy Brock, Soham De, Samuel L Smith, and Karen Simonyan. High-Performance Large-Scale Image Recognition Without Normalization. In *International Conference on Machine Learning (ICML)*, 2021. 5
- [S5] Moritz Böhle, Mario Fritz, and Bernt Schiele. Convolutional Dynamic Alignment Networks for Interpretable Classifications. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. 3, 5, 6
- [S6] Terrance DeVries and Graham W. Taylor. Improved Regularization of Convolutional Neural Networks with Cutout. *CoRR*, abs/1708.04552, 2017. 5
- [S7] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016. 4
- [S8] Elad Hoffer, Tal Ben-Nun, Itay Hubara, Niv Giladi, Torsten Hoefer, and Daniel Soudry. Augment Your Batch: Improving Generalization Through Instance Repetition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020. 5
- [S9] Diederik P. Kingma and Jimmy Ba. Adam: A Method for Stochastic Optimization. In *International Conference on Learning Representations (ICLR)*, 2015. 4
- [S10] Samuel G. Müller and Frank Hutter. TrivialAugment: Tuning-free Yet State-of-the-Art Data Augmentation. In *International Conference on Computer Vision (ICCV)*, 2021. 5
- [S11] Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. "Why Should I Trust You?": Explaining the Predictions of Any Classifier. In *International Conference on Knowledge Discovery and Data Mining (SIGKDD)*, 2016. 5
- [S12] Ramprasaath R. Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh, and Dhruv Batra. Grad-CAM: Visual Explanations from Deep Networks via Gradient-Based Localization. In *International Conference on Computer Vision (ICCV)*, 2017. 5
- [S13] Avanti Shrikumar, Peyton Greenside, and Anshul Kundaje. Learning Important Features Through Propagating Activation Differences. In *International Conference on Machine Learning (ICML)*, 2017. 5
- [S14] Jost Tobias Springenberg, Alexey Dosovitskiy, Thomas Brox, and Martin A. Riedmiller. Striving for Simplicity: The All Convolutional Net. In *International Conference on Learning Representations (ICLR), Workshop*, 2015. 5
- [S15] Mukund Sundararajan, Ankur Taly, and Qiqi Yan. Axiomatic Attribution for Deep Networks. In *International Conference on Machine Learning (ICML)*, 2017. 5
- [S16] Sangdoo Yun, Dongyoon Han, Seong Joon Oh, Sanghyuk Chun, Junsuk Choe, and Youngjoon Yoo. CutMix: Regularization Strategy to Train Strong Classifiers with Localizable Features. In *International Conference on Computer Vision (ICCV)*, 2019. 5
- [S17] Hongyi Zhang, Moustapha Cissé, Yann N. Dauphin, and David Lopez-Paz. mixup: Beyond Empirical Risk Minimization. In *International Conference on Learning Representations (ICLR)*, 2018. 5
- [S18] Zhun Zhong, Liang Zheng, Guoliang Kang, Shaozi Li, and Yi Yang. Random Erasing Data Augmentation. In *The Thirty-Fourth AAAI Conference on Artificial Intelligence AAAI*, 2020. 5