

---

# Concept Whitening for Interpretable Image Recognition

---

Zhi Chen<sup>1</sup> Yijie Bei<sup>2</sup> Cynthia Rudin<sup>1,2</sup>

## Abstract

What does a neural network encode about a concept as we traverse through the layers? Interpretability in machine learning is undoubtedly important, but the calculations of neural networks are very challenging to understand. Attempts to see inside their hidden layers can either be misleading, unusable, or rely on the latent space to possess properties that it may not have. In this work, rather than attempting to analyze a neural network posthoc, we introduce a mechanism, called *concept whitening* (CW), to alter a given layer of the network to allow us to better understand the computation leading up to that layer. When a concept whitening module is added to a CNN, the axes of the latent space can be aligned with concepts of interest. By experiment, we show that CW can provide us a much clearer understanding for how the network gradually learns concepts over layers without hurting predictive performance.

## 1. Introduction

An important practical challenge that arises with neural networks is the fact that the units within their hidden (intermediate, convolutional) layers are not usually semantically understandable. This is particularly true with computer vision applications, where an expanding body of research has focused centrally on explaining the calculations of neural networks and other black box models. Some of the core questions considered in these posthoc analyses of neural networks include: “What concept does a unit in a hidden layer of a trained neural network represent?” or “Does this unit in the network represent a concept that a human might understand?”

The questions listed above are important, but it is not clear that they would naturally have satisfactory answers when

performing posthoc analysis on a pre-trained neural network. In fact, there are several reasons why various types of posthoc analyses would not answer these questions. Efforts to interpret individual nodes of pre-trained neural networks (e.g., Zhou et al., 2018a; 2014) have shown that some fraction of nodes can be identified to be aligned with some high-level semantic meaning, but these special nodes do not provably contain the network’s full information about the concepts. That is, the nodes are not “pure,” and information about the concept could be scattered throughout the network. Concept-vector methods also (Kim et al., 2017; Zhou et al., 2018b; Ghorbani et al., 2019) have been used to analyze pre-trained neural networks. Here, vectors in the latent space are chosen to align with pre-defined or automatically-discovered concepts. While concept-vectors are more promising, they still make the assumption that the latent space of a neural network admits a posthoc analysis of a specific form. In particular, they assume that the latent space places members of each concept in one easy-to-classify portion of latent space. Since the latent space was not explicitly constructed to have this property, there is no reason to believe it holds.

Ideally, we would want a neural network whose latent space *tells us* how it is disentangling concepts, without needing to resort to extra classifiers like concept-vector methods (Kim et al., 2017; Ghorbani et al., 2019), without surveys to humans (Zhou et al., 2014), and without other manipulations that rely on whether the geometry of a latent space serendipitously admits analysis of concepts. Rather than having to rely on assumptions that the latent space admits disentanglement, we would prefer to constrain the latent space directly. We might even wish that the concepts align themselves along the axes of the latent space, so that each point in the latent space has an interpretation in terms of known concepts.

Let us discuss how one would go about imposing such constraints on the latent space. In particular, we introduce the possibility of what we call *concept whitening*. Concept whitening (CW) is a module inserted into a neural network. It constrains the latent space to represent target concepts and also provides a straightforward means to extract them. It does not force the concepts to be learned as an intermediate step, rather *it imposes the latent space to be aligned along the concepts*. For instance, let us say that, using CW on a lower layer of the network, the concept “airplane” is

---

<sup>1</sup>Department of Computer Science, Duke University, USA <sup>2</sup>Department of Electrical and Computer Engineering, Duke University, USA. Correspondence to: Zhi Chen <zhi.chen1@duke.edu>.

represented along one axis. By examining the images along this axis, we can find the lower-level abstraction that the network is using for the complex concept “airplane,” which might be white or silver objects with blue backgrounds. In the lower layers of a standard neural network we cannot necessarily find this abstraction, because the abstraction of “airplane” might be spread throughout latent space rather than along an “airplane” axis. By looking at images along the airplane axis at each layer, we see how the network gradually represents airplanes with an increasing level of sophistication and complexity.

Concept whitening could be used to replace a plain batch normalization step in a CNN backbone, because it combines batch whitening with an extra step involving a rotation matrix. Batch whitening usually provides helpful properties to latent spaces, but our goal requires the whitening to take place with respect to concepts; the use of the rotation matrix to align the concepts with the axes is the key to interpretability through disentangled concepts. Whitening decorrelates and normalizes each axis (i.e., transforms the post-convolution latent space so that the covariance matrix between channels is the identity). Exploiting the property that a whitening transformation remains valid after applying arbitrary rotation, the rotation matrix strategically matches the concepts to the axes.

The concepts used in CW do not need to be the labels in the classification problem, they can be learned from an auxiliary dataset in which concepts are labeled. The concepts do not need to be labeled in the dataset involved in the main classification task (though they could be), and the main classification labels do not need to be available in the auxiliary concept dataset.

Through qualitative and quantitative experiments, we illustrate how concept whitening applied to the various layers of the neural network illuminates its internal calculations. We verify the interpretability and pureness of concepts in the disentangled latent space. Importantly for practice, we show that by replacing the batch normalization layer in pretrained state-of-the-art models with a CW module, the resulting neural network can achieve accuracy on par with the corresponding original black box neural network on large datasets, and it can do this within one additional epoch of further training. Thus, with fairly minimal effort, one can make a small modification to a neural network architecture (adding a CW module), and in return be able to easily visualize how the network is learning all of the different concepts at any chosen layer.

CW can show us how a concept is represented at a given layer of the network. What we find is that at lower layers, since a complex concept cannot be represented by the network, it often creates lower-level *abstract concepts*. For example, an airplane at an early layer is represented by an

abstract concept that is white or gray objects on a blue background. A bed is represented by an abstract concept that seems to be characterized by warm colors (orange, yellow). In that sense, the CW layer can help us to discover new concepts that can be formally defined and built on.

## 2. Related work

There are several large and rapidly expanding bodies of relevant literature.

### Interpretability and explainability of neural networks:

There have been two schools of thought on improving the interpretability of neural networks: (1) learning an inherently *interpretable* model; (2) providing post-hoc *explanations* for an exist neural network. CW falls within the first type, though it only enlightens what the network is doing, rather than providing a full understanding of the network’s computations. To provide a full explanation of each computation would lead to more constraints and thus a loss in flexibility, whereas CW allows more flexibility in exchange for more general types of explanations. The vast majority of current works on neural networks are of the second type, explainability. A problem with the terminology is that “explanation” methods are often summary statistics of performance (e.g., local approximations, general trends on node activation) rather than actual explanations of the model’s calculations. For instance, if a node is found to activate when a certain concept is present in an image, it does not mean that all information (or even the majority of information) about this concept is involved with that particular node.

Saliency-based methods are the most common form of post-hoc explanations for neural networks (Zeiler and Fergus, 2014; Simonyan et al., 2013; Smilkov et al., 2017; Selvaraju et al., 2017). These methods assign importance weights to each pixel of the input image to show the importance of each pixel to the image’s predicted class. Saliency maps are problematic for well-known reasons: they often provide highlighting of edges in images, regardless of the class. Thus, very similar explanations are given for multiple classes, and often none of them are useful explanations. Saliency methods can be unreliable and fragile (e.g., Adebayo et al., 2018).

Other work provides explanations of how the network’s latent features operate. Some measure the alignment of an individual internal unit or filter of trained neural networks to a predefined concept and find some units have relatively strong alignment to that concept (Zhou et al., 2018a; 2014). While some units (i.e., filters) may align nicely with predefined concepts, the concept can be represented diffusely through many units (the concept representation by individual nodes is impure); this is because the network was not trained to have concepts expressed purely through individ-

ual nodes. To address this weakness, several concept-based post-hoc explanation approaches have recently been proposed that do not rely on the concept aligning with individual units (Kim et al., 2017; Zhou et al., 2018b; Ghorbani et al., 2019; Yeh et al., 2019). Instead of analyzing individual units, these methods try to learn a linear combination of them to represent a predefined concept (Kim et al., 2017) or to automatically discover concepts by clustering patches and defining the clusters as new concepts (Ghorbani et al., 2019). Although these methods are promising, they are based on assumptions of the latent space that may not hold. For instance, these methods assume that a classifier (usually a linear classifier) exists on the latent space such that the concept is correctly classified. Since the network was not trained so that this assumption holds, it may not hold. More importantly, since the latent space is not shaped explicitly to handle this kind of concept-based explanation, unit vectors (directions) in the latent space may not represent concepts purely. We will give an example in the next section to show why latent spaces built without constraints may not achieve concept separation.

CW avoids these problems because it shapes the latent space through training. In that sense, CW is closer to work on inherently interpretable neural networks, though its use-case is in the spirit of concept vectors, in that it is useful for providing important directions in the latent space.

There are emerging works trying to build inherently interpretable neural networks. Like CW, they alter the network structure to encourage different forms of interpretability. For example, neural networks have been designed to perform case-based reasoning (Chen et al., 2018; Li et al., 2018), to incorporate logical or grammatical structures (Li et al., 2017; Granmo et al., 2019; Wu and Song, 2019), to do classification based on hard attention (Mnih et al., 2014; Ba et al., 2014; Sermanet et al., 2014; Elsayed et al., 2019), or to do image recognition by decomposing the components of images (Saralajew et al., 2019). These models all have different forms of interpretability than we consider (understanding how the latent spaces of each layer learn a known set of concepts). One work that is somewhat similar to ours is that of Bouchacourt and Denoyer (2019), who develop a concept-based deep learning method that is inherently interpretable, but it relies on specific properties of textual data that do not readily transfer to image data. Zhang et al. (2018) add losses to the filters to encourage them to detect object parts that can also be viewed as concepts. However, unlike in CW, the method of Zhang et al. (2018) works only for object parts while CW works for any type of concept, such as objects, colors, textures, etc. Adel et al. (2018) transform the density of the current latent representation in an invertible way by normalizing flows and maximizing the mutual information between the transformed representation and side information provided by human users. Although

side information could also include concepts, Adel et al. (2018) query for the side information by active learning which is different from ours.

**Whitening and orthogonality:** Whitening is a linear transformation that transforms the covariance matrix of random input vectors to be the identity matrix. It is a classical preprocessing step in data science. In the realm of deep learning, batch normalization (Ioffe and Szegedy, 2015), which is widely used in many state-of-the-art neural network architectures, retains the standardization part of whitening but not the decorrelation. Earlier attempts (Desjardins et al., 2015; Luo, 2017) whiten by periodically estimating the whitening matrix, which leads to instability in training. Other methods (Cogswell et al., 2015) perform whitening by adding a decorrelation loss. By observing that SVD is differentiable, Huang et al. (2018b; 2019) develop ZCA whitening, supported directly in back-propagation. Siarohin et al. (2018) also propose a differentiable whitening block, but it is based on Cholesky whitening. The whitening part of our CW module borrows techniques from IterNorm (Huang et al., 2019) because it is differentiable and accelerated. CW is different from previous methods because its whitening matrix *is multiplied by an orthogonal matrix and maximizes the activation of known concepts along the latent space axes*.

In the field of deep learning, many initial works about incorporating orthogonality constraints are targeted for RNNs (Vorontsov et al., 2017; Mhammedi et al., 2017; Wisdom et al., 2016), since orthogonality could help avoid vanishing gradients or exploding gradients in RNNs. Other work explores ways to learn orthogonal weights or representations for all types of neural networks (not just RNNs) (Harandi and Fernando, 2016; Huang et al., 2018a; Lezcano-Casado and Martínez-Rubio, 2019; Lezama et al., 2018). For example, Lezama et al. (2018) use special loss functions to force orthogonality. The optimization algorithms used in the above methods are all different from ours. For CW, we optimize the orthogonal matrix by Cayley-transform-based curvilinear search algorithms proposed by Wen and Yin (2013). While Vorontsov et al. (2017) also use a Cayley transform, they do it with a fixed learning rate that does not work effectively in our setting. More importantly, the goal of doing optimization with orthogonality constraints in all these works are completely different from ours. None of them try to align columns of the orthogonal matrix with any type of concept.

### 3. Methodology

Suppose  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n \in \mathcal{X}$  are samples in our dataset and  $y_1, y_2, \dots, y_n \in \mathcal{Y}$  are their labels. From the latent space  $\mathcal{Z}$  defined by a hidden layer, a DNN classifier  $f : \mathcal{X} \rightarrow \mathcal{Y}$  can be divided into two parts, a feature extractor  $\Phi : \mathcal{X} \rightarrow \mathcal{Z}$ , with parameters  $\theta$ , and a classifier  $g : \mathcal{Z} \rightarrow \mathcal{Y}$ , parameterized by

$\omega$ . Then  $\mathbf{z} = \Phi(\mathbf{x}; \theta)$  is the latent representation of the input  $\mathbf{x}$  and  $f(\mathbf{x}) = g(\Phi(\mathbf{x}; \theta); \omega)$  is the predicted label. Suppose we are interested in  $k$  concepts called  $c_1, c_2, \dots, c_k$ . We can then pre-define  $k$  auxiliary datasets  $\mathbf{X}_{c_1}, \mathbf{X}_{c_2}, \dots, \mathbf{X}_{c_k}$  such that samples in  $\mathbf{X}_{c_j}$  are the most representative samples of concept  $c_j$ . Our goal is to learn  $\Phi$  and  $g$  simultaneously, such that (a) the classifier  $g(\Phi(\cdot; \theta); \omega)$  can predict the label accurately; (b) the  $j^{\text{th}}$  dimension  $z_j$  of the latent representation  $\mathbf{z}$  aligns with concept  $c_j$ . In other words, samples in  $\mathbf{X}_{c_j}$  should have larger values of  $z_j$  than other samples. Conversely, samples not in  $\mathbf{X}_{c_j}$  should have smaller values of  $z_j$ .

### 3.1. Standard Neural Networks May Not Achieve Concept Separation

Some posthoc explanation methods have looked at unit vectors in the direction of data where a concept is exhibited to measure how different concepts contribute to a classification task (Zhou et al., 2018b). Other methods consider directional derivatives towards data exhibiting the concept (Kim et al., 2017), for the same reason. There are important reasons why these types of approaches may not work.

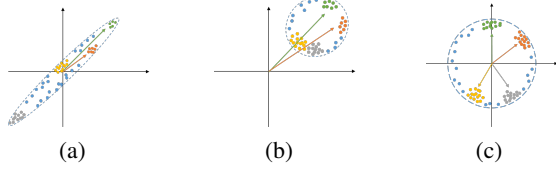


Figure 1. Possible data distributions in the latent space. In (a) the data are standardized but not decorrelated; In (b) the data are not mean centered; In (c) the data are whitened. In both (a) and (b), unit vectors are not valid for representing concepts.

Consider, for instance, an elongated latent space similar to that illustrated in Figure 3.1(a). Here, two unit vectors pointing to different groups of data (perhaps exhibiting two separate concepts) may have a large inner product, suggesting that they may be part of the same concept, when in fact, they may be not be similar at all. Worse, a unit vector in the yellow direction appears to indicate that the red concept has more extreme values of the yellow concept than the members of the yellow concept itself. Thus, even if the latent space is standardized, multiple unrelated concepts can still be found by traversing towards the same general direction, as shown in Figure 3.1(a). For the same reason, taking derivatives towards the parts of the space where various concepts tend to appear may yield similar derivatives for very different concepts. This is true even if these concepts are not co-located in latent space.

If the latent space is not mean-centered, that alone could cause problems for posthoc methods that compute directions

towards concepts. Consider, for instance, a case where all points in the latent space are far from the origin. In that case, *all* concept directions point towards the same part of the space: the part where the data lies (see Figure 3.1)(b).

At the very least, if we are to examine directions in the latent space to look for known concepts, the latent space should be mean-centered, and its axes should be aligned with these concepts.

### 3.2. Why Whitening Could Work

Fortunately, due to properties of high dimensional geometry (Blum et al., 2016), if our data are whitened (standardized and decorrelated), there is some hope that the concepts can be fully represented by unit vectors. Two reasons for this are shown by the following theorems, which are variations of standard results from high-dimensional geometry.

**Theorem 3.1.** (Variation of well-known theorem) *If a random vector  $\mathbf{x} \in \mathbb{R}^d$  is sampled from a  $d$ -dimensional spherical Gaussian distribution  $\mathcal{N}(\mathbf{0}_d, I_d)$ . Then for all  $\epsilon > 0$ ,*

$$Pr\left(\left|\frac{\|\mathbf{x}\|}{\sqrt{d-1}} - 1\right| > \epsilon\right) \leq \frac{4}{\epsilon\sqrt{d-1}} e^{-\frac{\epsilon^2(d-1)}{4}}.$$

This theorem means that if the dimension  $d$  of the latent space is large enough, which is reasonable for DNNs, then as long as the data follow a spherical Gaussian distribution, we have that almost all data are distributed near the surface of the sphere with radius  $\sqrt{d-1}$ . In other words, since all data points have approximately the same distance to the origin, they are distinguished only by their directions. Therefore, as long as all samples representing the concept are near each other, and as long as non-concept samples are not nearby, then a unit vector fully characterizes the location of that concept in latent space.

**Theorem 3.2.** (Variation of well-known theorem) *Suppose two unit vectors  $\mathbf{u}, \mathbf{v} \in \mathbb{R}^d$  are randomly drawn from a  $d$ -dimensional unit sphere. Then for all  $\epsilon > 0$ , the angle between them obeys*

$$Pr(|\mathbf{u} \cdot \mathbf{v}| < \epsilon) \geq 1 - \frac{4}{\epsilon\sqrt{d-2}} e^{-\frac{\epsilon^2(d-2)}{2}}.$$

This theorem indicates that two random unit vectors have high probability to be nearly orthogonal when  $d$  is large. Assuming that drawing images from different uncorrelated concepts is similar to drawing randomly distributed data, the unit vectors pointing to different concepts may be nearly orthogonal to each other. In that case, we may be able to find an orthonormal basis  $\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_d$  and use the first  $k$ -axes to represent the concepts of interest.

The above theorems provide the reasons why concept whitening might work: by aligning the samples representing



a concept in the *same* direction, and aligning each concept in *its own* direction, unit vectors are able to fully characterize that concept in latent space.

### 3.3. Concept Whitening Module

Let  $\mathbf{Z}_{d \times n}$  be the latent representation matrix of  $n$  samples, in which each column  $\mathbf{z}_i \in \mathbb{R}^d$  contains the latent features of the  $i^{th}$  sample. Our Concept Whitening module (CW) consists of two parts, whitening and orthogonal transformation. The whitening transformation  $\psi$  decorrelates and standardizes the data by

$$\psi(\mathbf{Z}) = \mathbf{W}(\mathbf{Z} - \mu \mathbf{1}_{n \times 1}^T) \quad (1)$$

where  $\mu = \frac{1}{n} \sum_{i=1}^n \mathbf{z}_i$  is the sample mean and  $\mathbf{W}_{d \times d}$  is the whitening matrix that obeys  $\mathbf{W}^T \mathbf{W} = \Sigma^{-1}$ . Here,  $\Sigma_{d \times d} = \frac{1}{n} (\mathbf{Z} - \mu \mathbf{1}^T)(\mathbf{Z} - \mu \mathbf{1}^T)^T$  is the covariance matrix. The whitening matrix  $\mathbf{W}$  is not unique and can be calculated in many ways such as ZCA whitening and Cholesky decomposition. Another important property of the whitening matrix is that it is rotation free; suppose  $\mathbf{Q}$  is an orthogonal matrix, then

$$\mathbf{W}' = \mathbf{Q}^T \mathbf{W} \quad (2)$$

is also a valid whitening matrix. In our module, after whitening the latent space to endow it with the properties discussed above, we still need to rotate the samples in their latent space such that the data from concept  $c_j$ , namely  $\mathbf{X}_{c_j}$ , are highly activated on the  $j^{th}$  axis. Specifically, we need to find an orthogonal matrix  $\mathbf{Q}_{d \times d}$  whose column  $\mathbf{q}_j$  is the  $j^{th}$  axis, by optimizing the following objective:

$$\begin{aligned} \max_{\mathbf{q}_1, \mathbf{q}_2, \dots, \mathbf{q}_k} \sum_{j=1}^k \frac{1}{n_j} \mathbf{q}_j^T \psi(\mathbf{Z}_{c_j}) \mathbf{1}_{n_j \times 1} \\ \text{s.t. } \mathbf{Q}^T \mathbf{Q} = \mathbf{I}_d \end{aligned} \quad (3)$$

where  $\mathbf{Z}_{c_j}$  is a  $d \times n_j$  matrix denoting the latent representation of  $\mathbf{X}_{c_j}$  and  $c_1, c_2, \dots, c_k$  are concepts of interest. An optimization problem with an orthogonality constraint like this can be solved by gradient-based approaches on the Stiefel manifold (e.g., the method of Wen and Yin, 2013).

This whole procedure constitutes CW, and can be done for any given layer of a neural network as part of the training of the network.

### 3.4. Optimization and Implementation Detail

Whitening has not (to our knowledge) been previously applied to align the latent space to concepts. In the past, whitening has been used to speed up back-propagation. The specific whitening problem for speeding up back-propagation is different from that for concept alignment—the rotation matrix is not present in other work on whitening, nor is the notion of a concept—however, we can leverage some

of the optimization tools used in that work on whitening (Huang et al., 2019; 2018a; Siarohin et al., 2018). Specifically, we adapt ideas underlying the IterNorm algorithm (Huang et al., 2019), which employs Newton’s iterations to approximate ZCA whitening, to the problem studied here. Let us now describe how this is done.

The whitening matrix in ZCA is

$$\mathbf{W} = \mathbf{D} \mathbf{\Lambda}^{-\frac{1}{2}} \mathbf{D}^T \quad (4)$$

where  $\mathbf{\Lambda}_{d \times d}$  and  $\mathbf{D}_{d \times d}$  are the eigenvalue diagonal matrix and eigenvector matrix given by the eigenvalue decomposition of the covariance matrix,  $\Sigma = \mathbf{D} \mathbf{\Lambda} \mathbf{D}^T$ . Like other normalization methods, we calculate a  $\mu$  and  $\mathbf{W}$  for each mini-batch of data, and average them together to form the model used in testing.

As mentioned in Section 3.3, the challenging part for CW is that we also need to learn an orthogonal matrix by solving an optimization problem. To do this, we will optimize the objective while strictly maintaining the matrix to be orthogonal by performing gradient descent with a curvilinear search on the Stiefel manifold (Wen and Yin, 2013) and adjust it to deal with mini-batch data.

During training, our procedure must handle two types of data: data for calculating the main objective and the data representing the predefined concepts. The model is optimized by alternating optimization: the mini-batches of the main dataset and the auxiliary concept dataset are fed to the network, and the following two objectives are optimized in turns. The first objective is:

$$\min_{\theta, \omega, \mathbf{W}, \mu} \frac{1}{n} \sum_{i=1}^n \text{CrsEnt}(g(\mathbf{Q}^T \psi(\Phi(\mathbf{x}_i; \theta); \mathbf{W}, \mu); \omega), y_i) \quad (5)$$

where  $\Phi$  and  $g$  are layers before and after the CW module parameterized by  $\theta$  and  $\omega$  respectively.  $\psi$  is a whitening transformation parameterized by sample mean  $\mu$  and whitening matrix  $\mathbf{W}$ .  $\mathbf{Q}$  is the orthogonal matrix and  $\mathbf{Q}^T \psi$  together form the CW module (which is also a valid whitening transformation). The second objective is

$$\begin{aligned} \max_{\mathbf{q}_1, \mathbf{q}_2, \dots, \mathbf{q}_k} \sum_{j=1}^k \frac{1}{n_j} \sum_{\mathbf{x}_i^{(c_j)} \in \mathbf{X}_{c_j}} \mathbf{q}_j^T \psi(\Phi(\mathbf{x}_i^{(c_j)}; \theta); \mathbf{W}, \mu) \\ \text{s.t. } \mathbf{Q}^T \mathbf{Q} = \mathbf{I}_d. \end{aligned} \quad (6)$$

The orthogonal matrix  $\mathbf{Q}$  is fixed when training for the main objective and the other parameters are fixed when training for  $\mathbf{Q}$ . The optimization problem is a linear programming problem with quadratic constraints (LPQC) which is generally NP-hard. Since directly solving for the optimal solution is intractable, we optimize it by gradient methods on the Stiefel manifold. At each step  $t$ , the orthogonal matrix  $\mathbf{Q}$  is

updated by Cayley transform

$$\mathbf{Q}^{(t+1)} = (I + \frac{\eta}{2}\mathbf{A})^{-1}(I - \frac{\eta}{2}\mathbf{A})\mathbf{Q}^{(t)}$$

where  $\mathbf{A} = \mathbf{G}(\mathbf{Q}^{(t)})^T - \mathbf{Q}^{(t)}\mathbf{G}^T$  is a skew-symmetric matrix,  $\mathbf{G}$  is the gradient of the loss function and  $\eta$  is the learning rate. The optimization procedure is accelerated by curvilinear search on the learning rate at each step (Wen and Yin, 2013). Note that, in the Cayley transform, the stationary points are reached when  $\mathbf{A} = \mathbf{0}$ , which has multiple solutions. Since the solutions are in high-dimensional space, these stationary points are very likely to be saddle points which can be avoided by SGD. Therefore, we use the stochastic gradient calculated by a mini-batch of samples to replace  $\mathbf{G}$  at each step. To accelerate and stabilize the stochastic gradient, we also apply momentum to it during implementation.

In CNNs, a feature map (a channel within one layer, created by a convolution of one filter) contains the information of how activated a part of the image is by a single filter. That filter may be a detector for a specific concept. Let us reshape the feature map into a vector, where each element of the vector represents how much one part of the image is activated by the filter. Thus, if the feature map for one filter is  $h \times w$  then a vector of length  $hw$  contains the activation information for that filter around the whole feature map. We do this reshaping procedure for each filter, which reshapes the output of a convolution layer  $Z_{h \times w \times d \times n}$  into a matrix  $Z_{d \times (hwn)}$ . We then perform CW on the reshaped matrix. After doing this, the resulting matrix is still size  $d \times (hwn)$ . If we reshape this matrix back to its original size as a tensor, one feature map of the tensor now (after training) represents whether a meaningful concept is detected at each location in the image for that layer. Note that, now the output of a filter is a feature map which is a  $h \times w$  matrix but the concept activation score we used in the optimization problem is a scalar. Therefore, we need to get an activation value from the feature map. We may have multiple ways to do it. We try the following calculations to define activation based on the feature map: (a) mean of all feature map values; (b) max of all feature map values; (c) mean of all positive feature map values; (d) mean of down-sampled feature map obtained by max pooling. We use (d) in our experiments since it is good at capturing both high-level and low-level concepts. Detailed analysis and experiments about the choice of different activation calculations are discussed in Appendix A.

Let us discuss some aspects of practical implementation. The CW module can substitute for other normalization modules such as BatchNorm in an hidden layer of the CNN. Because both whitening and orthogonal optimization require relatively higher computational cost, one can leverage a pretrained model as a warm start. To do this, we might leverage a pretrained model (for the same main objective)

that does not use CW, and replace a BatchNorm layer in that network with a CW layer. The model usually converges in one epoch (one pass over the data) if a pretrained model is used.

## 4. Experiments

In this section, we first show that after replacing one batch norm (BN) layer with our CW module, the accuracy of image recognition is still on par with the original model (4.1). After that, we visualize the concept basis we learn and show that the axes are aligned with the concepts assigned to them. Specifically, we display the images that are most activated along a single axis (2); we then show how two axes interact with each other (4.2.2); and we further show how the same concept evolves in different layers (4.2.3), where we have replaced one layer at a time. Then we measure the purity of our concept axes and compare with other concept-based neural network methods (4.3).

### 4.1. Main Objective Accuracy

We evaluate the image recognition accuracy of the CNNs before and after adding a CW module. We show that simply replacing a BN module with a CW module and training for a single epoch leads to (at most) a small drop in performance on its main objective performance. Specifically, after replacing the BN module with the CW module, we trained popular CNN architectures including VGG16+BN (Simonyan and Zisserman, 2014), ResNet with 18 layers and 50 layers (He et al., 2016) and DenseNet161 (Huang et al., 2017) on the Places365 (Zhou et al., 2017) dataset. The auxiliary concept dataset we used is MS COCO (Lin et al., 2014). Each annotation, e.g. “person,” in MS COCO was used as one concept, and we selected all the images with this annotation (images having “person” in it) as the data representing the concept. In order to limit the total time of the training process, we used pre-trained models for the popular CNN architectures (discussed above) and fine-tuned these models after BN was replaced with CW.

Table 1 shows the average test accuracy on the validation set of Places365 over 5 runs. We randomly selected 3 concepts to learn using CW for each run, and used the average of them to measure accuracy. We repeated this, applying CW to different layers and reported the average accuracy among the layers. The accuracy does not change much among CW applied to the different layers, as shown in Appendix B.

Because we have leveraged a pretrained model, when training with CW, we conduct only one additional epoch of training (one pass over the dataset) for each run. As shown in Table 1, the performance of these models using the CW module is on par with the original model: the difference is within 1% with respect to top-1 and top-5 accuracy. This

Table 1. Top-1 and top-5 test accuracy on Places365 dataset. Our results show that CW does not hurt performance.

	Top-1 acc.		Top-5 acc.	
	Original	+CW	Original	+CW
VGG16-BN	53.6	53.3	84.2	83.8
ResNet18	54.5	53.9	84.6	84.2
ResNet50	54.7	54.9	85.1	85.2
DenseNet161	55.3	55.5	85.2	85.6

means in practice, if a pretrained model (using BN) exists, one can simply replace the BN module with a CW module and train it for one epoch, in which case, the pretrained black-box model can be turned into a more interpretable model that is approximately equally accurate.

## 4.2. Visualizing the Concept Basis

In order to demonstrate the interpretability benefits of models equipped with CW modules, we visualize the concept basis in the CW module and validate that the axes are aligned with their assigned concepts. In detail, (a) we check the most activated images on these axes; (b) we look at how images are distributed in a 2D slice of the latent space; (c) we show how realizations of the same concept change if we apply CW on different layers. All of these experiments were done on a ResNet18 equipped with CW trained on Places365.

### 4.2.1. TOP-10 IMAGES ACTIVATED

We sort all validation samples by their activation values (discussed in Section 3.4) to show how much they are related to the concept. Figure 2 shows the images that have the top-10 largest activations along three different concept’s axes. Note that all these concepts are trained together using one CW module.

From Figure 2(b) we can see that all the top activated images have the same semantic meaning if the CW module is located a higher layer (namely the 16th layer). Figure 2(a) shows that when the CW module is applied to a lower layer (namely the 2nd layer), it tends to capture low level information such as color or texture characteristic of these concepts. For instance, the top activated images on the “airplane” axis generally have a blue background with a white or gray object in the middle, which also happens in real airplanes images. It is reasonable that the lower layer CW module cannot extract complete information about high-level concepts such as “airplane” since the model complexity of the first two layers is limited.

In that sense, the CW layer has *discovered* an abstraction of a more complex concept; namely it has discovered that the blue images with white objects are primitive representations

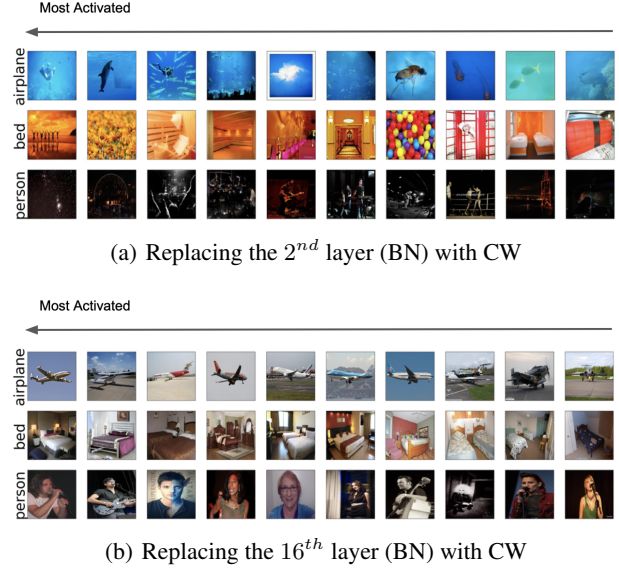


Figure 2. Top-10 Image activated on axes representing different concepts.

of the “airplane” concept. Similarly, the network seems to have discovered that warm colors is a lower-level abstraction of the “bedroom” concept, and that dark background with vertical light is an abstraction of the “person” concept.

Interestingly, when different definitions of activation are used (namely the options discussed in Section 3.4), the abstract concepts discovered by the network often look different. Some of these are shown in Appendix A.2.

### 4.2.2. 2D-REPRESENTATION SPACE VISUALIZATION

Let us consider whether joint information about different concepts is captured by the latent space of CW. To investigate how the data are distributed in the new latent space, we pick a 2D slice of the latent space, which means we select two axes  $\mathbf{q}_i$  and  $\mathbf{q}_j$  and look at the subspace they form.

The data’s joint distribution on the two axes is shown in Figure 3. To visualize the joint distribution, we first compute the activations of all validation data on the two axes, then divide the latent space into a  $50 \times 50$  grid of blocks, where the maximum and minimum activation value are the top and bottom of the grid. For the grid shown in Figure 3 (a), we randomly select one image that falls into each block, and display the image in its corresponding block. If there is no image in the block, the block remains black. From Figure 3 (a), we observe that the axes are not only aligned with their assigned concepts, they also incorporate joint information. For example, a “person in bed” has high activation on both the “person” axis and “bed” axis.

We also include a 2D histogram of the number of images that



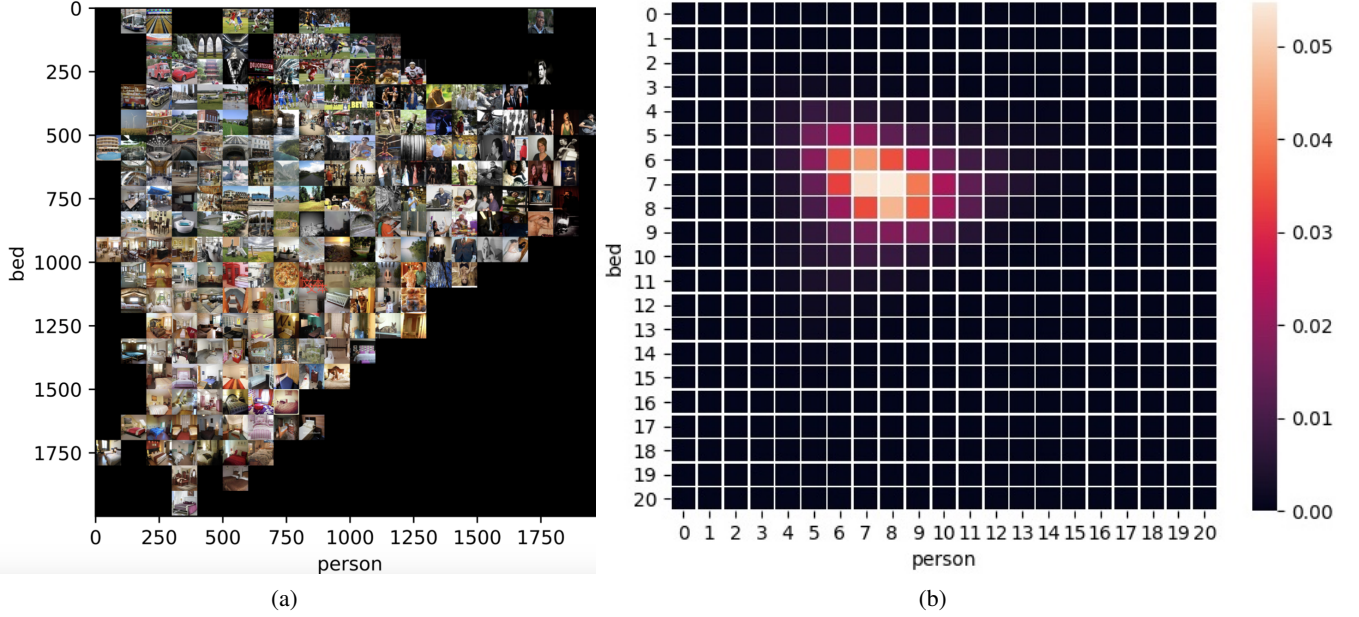


Figure 3. Joint distribution of the bed-person subspace. The bounding box given by projected values in the subspace is evenly divided into  $20 \times 20$  blocks. (a) Plotting a random test image fall into each block; (b) Density map of test image representation

fall into each block. As shown in Figure 3(b), most images are distributed near the center which agrees with Theorem 3.2: the samples’ feature vector has high probability to be nearly orthogonal to the axis we pick, and consequently they have near 0 activation on the axis itself. Note that the 2D histogram does not contradict Theorem 3.1 since the samples that have low activations on these two axes might have high activations on the unplotted axes. Therefore, they can still have distances to the origin that are similar to those of the highly activated samples.

#### 4.2.3. TRAJECTORY OF CONCEPTS IN DIFFERENT LAYERS

Although our objective is the same when we apply the CW module to different layers in the same CNN, the latent space we get might be different. This is because different layers might have different levels of semantic meaning. Because of this, it might be interesting to track how the representation of a single image will change as the CW module is applied to different layers of the CNN.

In order to better understand the latent representation, we plot a 2D slice of the latent space. Unlike in Section 4.2.2, here, a point in the plot is not specified by the activation values themselves but by their rankings. For example, the point (0.7, 0.1) means the point is a 70% quantile in the first axis and 10% quantile in the second axis. We use the percentage instead of using the value, because as mentioned in Section 4.2.2, most points are near the center of the plot,

so the rankings spread the values for plotting purposes.

Figure 4 shows the 2D representation plot of two representative images. Each point in the plot corresponds to the percentile rank representation of the image when the CW module is applied to different layers. The points are connected by arrows according to the depth of the layer. These plots confirm that the abstract concepts learned in the lower layers tend to capture lower-level meaning (such as colors or shapes) while the higher layers capture high-level meaning (such as types of objects). For example, in the left image in Figure 4 (a), the bed is blue, where blue is typical low level information about the “airplane” class but not about the “bed” class since bedrooms are usually warm colors. Therefore, in lower layers, the bed image has higher ranking in the “airplane” axis than the “bed” axis. However, when CW is applied to deeper layers, high level information is available, and thus the image becomes highly ranked on the “bed” axis and lower on the “airplane” axis.

In Figure 4 (b), traversing through the networks’ layers, the image of a sunset does not have the typical blue coloring of a sky. Its warm colors put it high on the “bedroom” concept for the second layer, and low on the “airplane” concept. However, as we look at higher layers, where the network can represent more sophisticated concepts, we see the image’s rank grow on the “airplane” concept (perhaps the network uses the presence of skies to detect airplanes), and decrease on the “bed” concept. From there, as we increase layers, the “airplane” concept decreases slightly (perhaps because



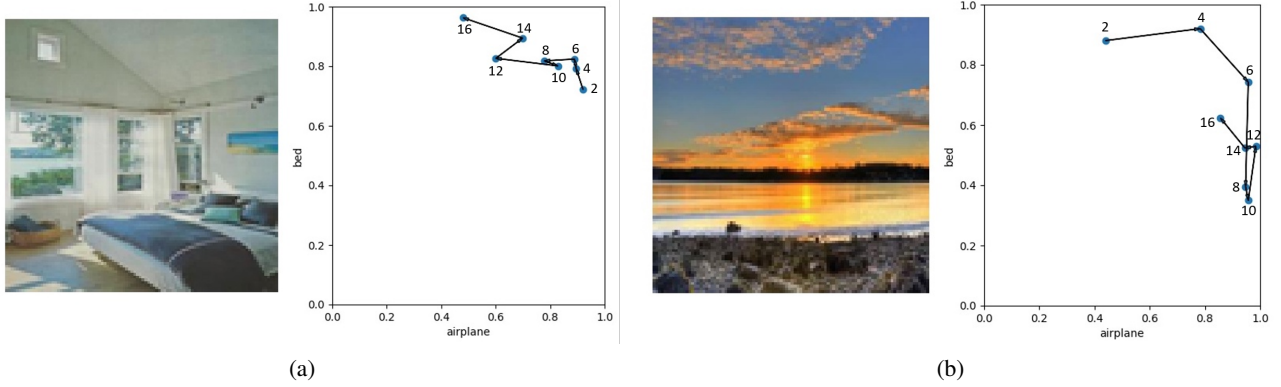


Figure 4. 2D representation plot of two representative images. Each point in the right trajectory plot correspond to the percentile rank for the activation values of on each axis. The number on the points means the layer depth of the CW module. The trajectory shows how the percentile rank of the left image changes when CW is applied to different layers.

there is no airplane in the image), and the “bed” concept increases slightly.

### 4.3. Quantitative Comparison of Interpretability

In this subsection, we measure the interpretability of the latent space quantitatively. To quantitatively define the interpretability with respect to concepts, we measure the purity of the concepts we learned with CW and compare with other concept-based methods. The purity is measured by the AUC calculated from the activation values. Specifically, we choose 10 concepts to learn at the same time. Each concept dataset is divided into a training set and testing set. After training the CW module using the training set, we get the testing samples’ activation values on the 10 concept axes. For each concept axis, we assign samples of this concept to the label 1 while giving other samples label 0. In this way, we can calculate the AUC score of the latent space with respect to each concept. The AUC score measures whether the samples belongs to a concept are ranked higher than other samples. Thus, the AUC score indicates the purity of the concept axis.

We compare the concept purity measured by AUC with the concept vectors learned by TCAV (Kim et al., 2017), IBD (Zhou et al., 2018b) and filters in standard CNNs (Zhou et al., 2014). For TCAV and IBD, since these methods already find concept vectors, we use the samples’ projections on the vectors to measure the AUC score. For filters in standard CNNs, we measure the AUC score for all filters and choose the best one to compare with our method, separately for each concept (denoted “Best Filter”). As shown in Figure 5, we compare these methods across the different layers. The concepts learned in the CW module are generally purer than those of other methods. This results from CW’s whitening of the latent space and optimization of the

loss function, as illustrated in Sections 3.1 and 3.2.

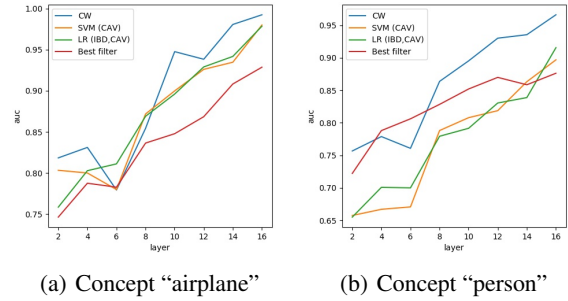


Figure 5. Concept purity measured by AUC score. Concept purity of CW module is compared to other posthoc methods on different layers.

We also compare the correlation of axes in the latent space before and after the CW module is applied. For comparison with posthoc methods like TCAV and IBD, we measure the output of their BN modules in the pretrained model, because the output of these layers are mean centered and normalized, which, as we discussed, are important properties for concept vectors. Shown by the absolute correlation coefficients plotted in Figure 6, the axes still have relatively strong correlation after passing through the BN module. (If CW were applied instead of BN, they would instead be decorrelated). This result reflects why purity of concepts is important; when the axes are pure, the signal of one concept can be concentrated only on its axis, while in standard CNNs, the concept could be distributed throughout the latent space.

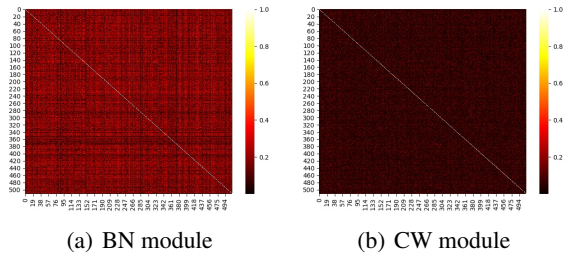


Figure 6. Absolute correlation coefficient of every feature pair in the 16<sup>th</sup> layer. (a) when the 16<sup>th</sup> layer is a BN module; (b) when 16<sup>th</sup> layer is a CW module.

## 5. Conclusion and Future Work

Concept whitening is a module placed at the bottleneck of a CNN, to force the latent space to be disentangled, and to align the axes of the latent space with the predefined concepts. By building an inherently interpretable CNN with concept whitening, we can gain intuition about how the network gradually learns the target concepts over the layers without harming the main objective’s performance.

There are many avenues for possible future work. Since CW modules are useful for helping humans to define primitive abstract concepts, such as those we have seen the network use at early layers, it would be interesting to automatically detect and quantify these new concepts, in the spirit of Ghorbani et al. (2019). Also the requirement of CW to completely decorrelate the outputs of all the filters might be too strong for some tasks. This is because concepts might be highly correlated in practice such as “airplane” and “sky”. In this case, we may want to soften our definition of CW. We could define several general topics that are uncorrelated, and use multiple correlated filters to represent concepts within each general topic. In this scenario, instead of forcing the gram matrix to be the identity matrix, we could make it block diagonal. The orthogonal basis would become a set of orthogonal subspaces.

## References

- Adebayo, J., Gilmer, J., Muelly, M., Goodfellow, I., Hardt, M., and Kim, B. (2018). Sanity checks for saliency maps. In *Advances in Neural Information Processing Systems*, pages 9505–9515.
- Adel, T., Ghahramani, Z., and Weller, A. (2018). Discovering interpretable representations for both deep generative and discriminative models. In *International Conference on Machine Learning*, pages 50–59.
- Ba, J., Mnih, V., and Kavukcuoglu, K. (2014). Multiple object recognition with visual attention. *arXiv preprint arXiv:1412.7755*.
- Blum, A., Hopcroft, J., and Kannan, R. (2016). Foundations of data science. *Vorabversion eines Lehrbuchs*.
- Bouchacourt, D. and Denoyer, L. (2019). Educe: Explaining model decisions through unsupervised concepts extraction. *arXiv preprint arXiv:1905.11852*.
- Chen, C., Li, O., Tao, C., Barnett, A. J., Su, J., and Rudin, C. (2018). This looks like that: deep learning for interpretable image recognition. *arXiv preprint arXiv:1806.10574*.
- Cogswell, M., Ahmed, F., Girshick, R., Zitnick, L., and Batra, D. (2015). Reducing overfitting in deep networks by decorrelating representations. *arXiv preprint arXiv:1511.06068*.
- Desjardins, G., Simonyan, K., Pascanu, R., et al. (2015). Natural neural networks. In *Advances in Neural Information Processing Systems*, pages 2071–2079.
- Elsayed, G., Kornblith, S., and Le, Q. V. (2019). Saccader: Improving accuracy of hard attention models for vision. In *Advances in Neural Information Processing Systems*, pages 700–712.
- Ghorbani, A., Wexler, J., Zou, J. Y., and Kim, B. (2019). Towards automatic concept-based explanations. In *Advances in Neural Information Processing Systems*, pages 9273–9282.
- Granmo, O.-C., Glimsdal, S., Jiao, L., Goodwin, M., Omlin, C. W., and Berge, G. T. (2019). The convolutional tsetlin machine. *arXiv preprint arXiv:1905.09688*.
- Harandi, M. and Fernando, B. (2016). Generalized back-propagation, étude de cas: Orthogonality. *arXiv preprint arXiv:1611.05927*.
- He, K., Zhang, X., Ren, S., and Sun, J. (2016). Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778.
- Huang, G., Liu, Z., Van Der Maaten, L., and Weinberger, K. Q. (2017). Densely connected convolutional networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4700–4708.
- Huang, L., Liu, X., Lang, B., Yu, A. W., Wang, Y., and Li, B. (2018a). Orthogonal weight normalization: Solution to optimization over multiple dependent stiefel manifolds in deep neural networks. In *Thirty-Second AAAI Conference on Artificial Intelligence*.

- Huang, L., Yang, D., Lang, B., and Deng, J. (2018b). Decorrelated batch normalization. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 791–800.
- Huang, L., Zhou, Y., Zhu, F., Liu, L., and Shao, L. (2019). Iterative normalization: Beyond standardization towards efficient whitening. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4874–4883.
- Ioffe, S. and Szegedy, C. (2015). Batch normalization: Accelerating deep network training by reducing internal covariate shift. *arXiv preprint arXiv:1502.03167*.
- Kim, B., Wattenberg, M., Gilmer, J., Cai, C., Wexler, J., Viegas, F., and Sayres, R. (2017). Interpretability beyond feature attribution: Quantitative testing with concept activation vectors (tcav). *arXiv preprint arXiv:1711.11279*.
- Lezama, J., Qiu, Q., Musé, P., and Sapiro, G. (2018). Ole: Orthogonal low-rank embedding-a plug and play geometric loss for deep learning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 8109–8118.
- Lezcano-Casado, M. and Martínez-Rubio, D. (2019). Cheap orthogonal constraints in neural networks: A simple parametrization of the orthogonal and unitary group. *arXiv preprint arXiv:1901.08428*.
- Li, O., Liu, H., Chen, C., and Rudin, C. (2018). Deep learning for case-based reasoning through prototypes: A neural network that explains its predictions. In *Thirty-Second AAAI Conference on Artificial Intelligence*.
- Li, X., Song, X., and Wu, T. (2017). Aognets: Compositional grammatical architectures for deep learning. *arXiv preprint arXiv:1711.05847*.
- Lin, T.-Y., Maire, M., Belongie, S., Hays, J., Perona, P., Ramanan, D., Dollár, P., and Zitnick, C. L. (2014). Microsoft coco: Common objects in context. In *European conference on computer vision*, pages 740–755. Springer.
- Luo, P. (2017). Learning deep architectures via generalized whitened neural networks. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages 2238–2246. JMLR. org.
- Mhammedi, Z., Hellicar, A., Rahman, A., and Bailey, J. (2017). Efficient orthogonal parametrisation of recurrent neural networks using householder reflections. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages 2401–2409. JMLR. org.
- Mnih, V., Heess, N., Graves, A., et al. (2014). Recurrent models of visual attention. In *Advances in neural information processing systems*, pages 2204–2212.
- Saralajew, S., Holdijk, L., Rees, M., Asan, E., and Villmann, T. (2019). Classification-by-components: Probabilistic modeling of reasoning over a set of components. In *Advances in Neural Information Processing Systems*, pages 2788–2799.
- Selvaraju, R. R., Cogswell, M., Das, A., Vedantam, R., Parikh, D., and Batra, D. (2017). Grad-cam: Visual explanations from deep networks via gradient-based localization. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 618–626.
- Sermanet, P., Frome, A., and Real, E. (2014). Attention for fine-grained categorization. *arXiv preprint arXiv:1412.7054*.
- Siarohin, A., Sangineto, E., and Sebe, N. (2018). Whitening and coloring batch transform for gans. *arXiv preprint arXiv:1806.00420*.
- Simonyan, K., Vedaldi, A., and Zisserman, A. (2013). Deep inside convolutional networks: Visualising image classification models and saliency maps. *arXiv preprint arXiv:1312.6034*.
- Simonyan, K. and Zisserman, A. (2014). Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*.
- Smilkov, D., Thorat, N., Kim, B., Viégas, F., and Wattenberg, M. (2017). Smoothgrad: removing noise by adding noise. *arXiv preprint arXiv:1706.03825*.
- Vorontsov, E., Trabelsi, C., Kadoury, S., and Pal, C. (2017). On orthogonality and learning recurrent networks with long term dependencies. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages 3570–3578. JMLR. org.
- Wen, Z. and Yin, W. (2013). A feasible method for optimization with orthogonality constraints. *Mathematical Programming*, 142(1-2):397–434.
- Wisdom, S., Powers, T., Hershey, J., Le Roux, J., and Atlas, L. (2016). Full-capacity unitary recurrent neural networks. In *Advances in Neural Information Processing Systems*, pages 4880–4888.
- Wu, T. and Song, X. (2019). Towards interpretable object detection by unfolding latent structures. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 6033–6043.
- Yeh, C.-K., Kim, B., Arik, S. O., Li, C.-L., Ravikumar, P., and Pfister, T. (2019). On concept-based explanations in deep neural networks. *arXiv preprint arXiv:1910.07969*.

- Zeiler, M. D. and Fergus, R. (2014). Visualizing and understanding convolutional networks. In *European conference on computer vision*, pages 818–833. Springer.
- Zhang, Q., Nian Wu, Y., and Zhu, S.-C. (2018). Interpretable convolutional neural networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 8827–8836.
- Zhou, B., Bau, D., Oliva, A., and Torralba, A. (2018a). Interpreting deep visual representations via network dissection. *IEEE transactions on pattern analysis and machine intelligence*.
- Zhou, B., Khosla, A., Lapedriza, A., Oliva, A., and Torralba, A. (2014). Object detectors emerge in deep scene cnns. *arXiv preprint arXiv:1412.6856*.
- Zhou, B., Lapedriza, A., Khosla, A., Oliva, A., and Torralba, A. (2017). Places: A 10 million image database for scene recognition. *IEEE transactions on pattern analysis and machine intelligence*, 40(6):1452–1464.
- Zhou, B., Sun, Y., Bau, D., and Torralba, A. (2018b). Interpretable basis decomposition for visual explanation. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 119–134.



## A. Concept Activation Calculation and Concept Activation Comparison Experiments

### A.1. Calculations of Concept Activation Based on Feature Maps

The output of a single filter is a  $h \times w$  feature map. However, a scalar is needed to quantify how much a sample is activated on a concept, which is used in both optimization and evaluation. Based on a feature map, multiple reasonable ways exists to calculate the concept activation.

Specifically, we try the following calculations to produce an activation value:

- Mean of all feature map values
- Max of all feature map values
- Mean of all positive feature map values
- Mean of down-sampled feature map obtained by max pooling.

Figure 7 shows these four methods of calculating the activation through demonstration. Among them, the mean of values is more suitable for capturing low-level concepts since they are distributed throughout the feature map. For high-level concepts, the max value and mean of positive values are more powerful: they can capture high-level concepts such as objects, since objects usually occur just in one location, not repeatedly throughout an image. The mean of max-pooled values is a combination of the previous types and is capable of representing both high-level and low-level concepts. Intuitively, the mean of max pooled values is more similar to the max function when applied to higher layers and more similar to the mean function when applied to lower layers. This is because, for higher layers, the mean is taken of only a few values, simply because higher layers are smaller in size. Thus, the max is the dominant calculation. In contrast, for lower layers, which are much larger, the max's are taken over a relatively small number of elements (local regions), and then the mean is taken over all of the local regions. Hence the mean is the dominant calculation for lower layers.

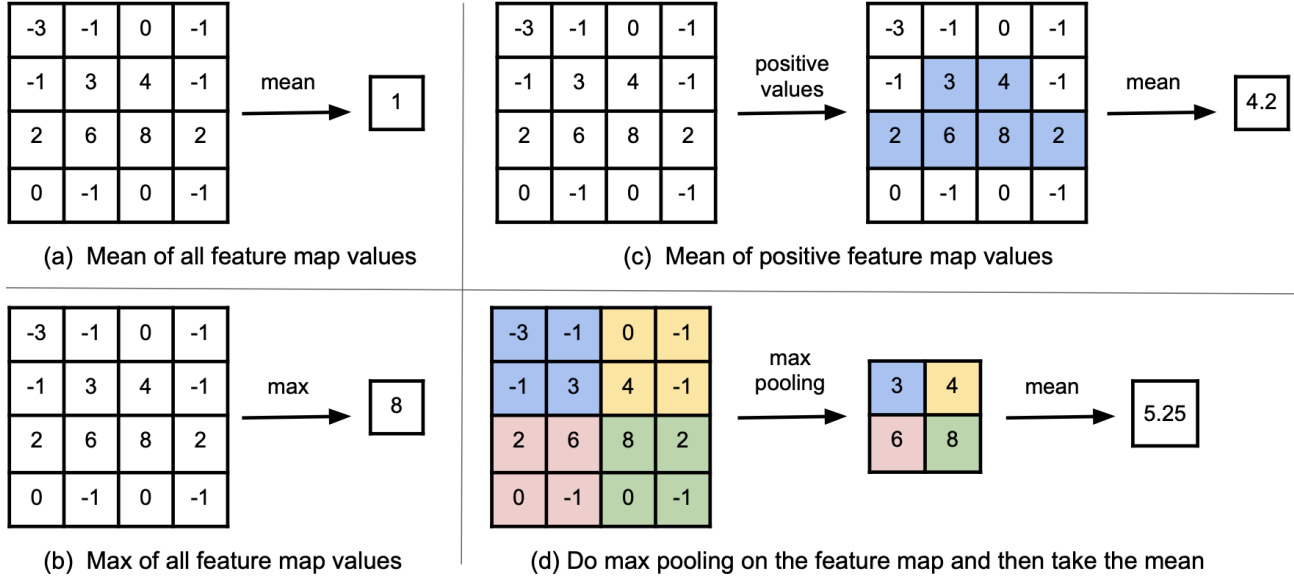


Figure 7. Four methods of calculating concept activation based on the feature map.

### A.2. Top-10 Activated Images Based on Different Calculations

Figure 8 shows the top-10 activated images under the four different calculations for concept activation. The CNN architecture, dataset and the depth of the CW module are the same as before. The figures show that when concept activation is calculated in different ways, the most activated images may look different and the network even may discover completely different abstract concepts. For example, when CW is applied to the  $2^{nd}$  layer, the network discovered the abstraction of the concept “bed” to be warm colors when the activation was the mean of feature map values, while the abstract concepts seems to involve boundaries of colors if activation is calculated as the max value. Also if the activation is calculated as the mean of all

values, the “person” concept gives rise to an abstract concept involving dense texture, while under the mean of max-pooled values, the “person” concept is abstracted to a dark background with vertical lights. This difference in the discovered abstract concepts could be explained by the fact that these calculation methods focus on different locations within the image: the mean value focuses on the whole image while the max value only looks at one place within the image.

### A.3. Concept AUC Based on Different Activation Calculations

Table 2 shows concept AUC when different concept activation definitions are used. The definition and calculation of concept AUC is the same as in Section 4.3. The dataset and CNN architecture are also the same. To compare these concept activations’ capability to capture both high-level concepts and low-level concepts, we apply CW to the  $2^{nd}$  and  $16^{th}$  layers of ResNet18. Table 2 indicates that in the  $2^{nd}$  layer, the max value of the feature map performs poorly in AUC than the other calculation methods. In contrast, in the  $16^{th}$  layer, the mean performs poorly compared to the other methods. The max-pool-mean method performs well on both layers. This result matches our intuitive reasoning that the max-pool-mean combines the advantages of mean and max. It is suitable for capturing both low-level concepts and high-level concepts.

Table 2. Concept AUC obtained by different calculations of concept activation. Max-pool-mean performs well when CW applied to both low and high layers while others all have shortage.

	AUC-“airplane”		AUC-“bed”		AUC-“person”	
	$2^{nd}$ layer	$16^{th}$ layer	$2^{nd}$ layer	$16^{th}$ layer	$2^{nd}$ layer	$16^{th}$ layer
Mean	0.820	0.981	0.687	0.853	0.714	0.918
Max	0.716	0.992	0.589	0.904	0.759	0.969
Positive-mean	0.798	0.992	0.614	0.924	0.757	0.968
Max-pool-mean	0.818	0.993	0.692	0.906	0.757	0.966

## B. Main Objective Accuracy when CW is Applied to Different Layers

As mentioned in Section 4.1, we measure the main objective accuracy when CW applied to different layers. Tables 3 through 6 show the layer-wise test accuracy of different CNN architectures. The dataset and CNN architectures are the same as in Section 4.1. Results in Tables 3 through 6 indicate that no matter which layer we apply CW, accuracy is not substantially impacted.

Table 3. Top-1 and top-5 accuracy of VGG16-CW on Places365 dataset. Our results indicate that the choice of layer to apply CW does not have a practical impact on accuracy.

CW layer	Top-1 acc.	Top-5 acc.
1 <sup>nd</sup>	53.2	83.8
2 <sup>th</sup>	53.3	83.8
3 <sup>th</sup>	53.4	83.8
4 <sup>th</sup>	53.4	83.9
5 <sup>th</sup>	53.2	83.9
6 <sup>th</sup>	53.3	83.8
7 <sup>th</sup>	53.5	83.8
8 <sup>rd</sup>	53.3	83.9
9 <sup>nd</sup>	53.4	83.8
10 <sup>th</sup>	53.2	83.8
11 <sup>nd</sup>	53.2	83.9
12 <sup>th</sup>	53.3	83.7

Table 4. Top-1 and top-5 accuracy of ResNet50-CW on Places365 dataset. Our results indicate that the choice of layer to apply CW does not have a practical impact on accuracy.

CW layer	Top-1 acc.	Top-5 acc.
2 <sup>nd</sup>	55.2	85.4
5 <sup>th</sup>	55.3	85.5
8 <sup>th</sup>	55.3	85.5
11 <sup>th</sup>	55.2	85.5
14 <sup>th</sup>	55.3	85.5
17 <sup>th</sup>	54.8	85.2
20 <sup>th</sup>	54.7	85.0
23 <sup>rd</sup>	54.8	85.0
26 <sup>nd</sup>	54.7	85.0
29 <sup>th</sup>	54.8	85.0
32 <sup>nd</sup>	54.8	85.1
35 <sup>th</sup>	54.7	85.0
38 <sup>th</sup>	54.8	85.1
41 <sup>st</sup>	54.6	85.0
44 <sup>th</sup>	54.7	84.9
47 <sup>th</sup>	54.6	85.0

Table 5. Top-1 and top-5 accuracy of DenseNet161-CW on Places365 dataset. Our results indicate that the choice of layer to apply CW does not have a practical impact on accuracy.

CW layer	Top-1 acc.	Top-5 acc.
14 <sup>th</sup>	55.6	85.7
39 <sup>th</sup>	55.5	85.5
88 <sup>nd</sup>	55.5	85.6
161 <sup>th</sup>	55.5	85.6

Table 6. Top-1 and top-5 accuracy of ResNet18-CW on Places365 dataset. Our results indicate that the choice of layer to apply CW does not have a practical impact on accuracy.

CW layer	Top-1 acc.	Top-5 acc.
2 <sup>nd</sup>	53.9	84.2
4 <sup>th</sup>	54.0	84.5
6 <sup>th</sup>	54.0	84.3
8 <sup>th</sup>	54.0	84.2
10 <sup>th</sup>	54.0	84.3
12 <sup>th</sup>	53.9	84.1
14 <sup>th</sup>	53.7	83.9
16 <sup>th</sup>	53.5	83.8

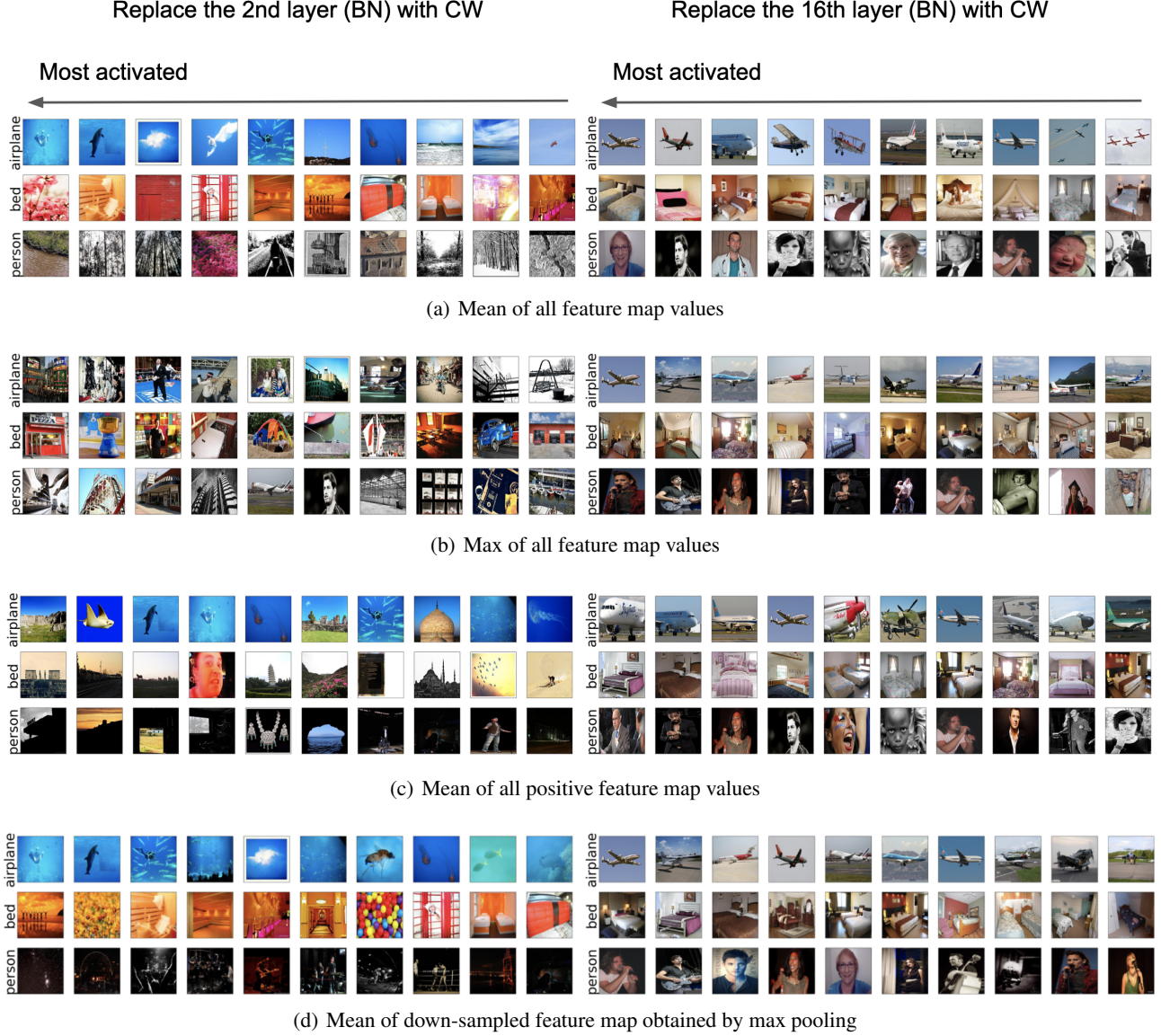


Figure 8. Top-10 activated images obtained by different calculations of concept activation. Depending on which choice (max, mean, mean of positives, mean of max pool values), different abstract concepts are generated. For instance, for the person class, using the mean calculation on the 2nd layer (top left), the abstract concept is a dense texture. For the person class using the mean of max-pooled values (bottom left), the abstract concept is dark background with vertical lights. For the bed concept with the mean of max-pooled values (bottom left), the abstract concept is warm colors, whereas for the max calculation (second row left) the abstract concept is boundaries of different colors. These concepts could later be formalized, if desired, to create better or more interpretable classifiers in the future.