

Exploring the Spatial Dynamics of In-Distribution and Out-of-Distribution Data in Logit Space

Anonymous authors
Paper under double-blind review

Abstract

Out-of-distribution (OOD) data pose a significant challenge to deep learning (DL) classifiers, prompting extensive research into their effective detection methods. Current state-of-the-art OOD detection methods employ a scoring technique designed to assign lower scores to OOD samples compared to in-distribution (ID) ones. Nevertheless, these approaches lack foresight into the configuration of OOD and ID data within the latent space. Instead, they make an implicit assumption about their inherent separation or force a separation post-training by utilizing selected OOD data. As a result, most OOD detection methods result in complicated and hard-to-validate scoring techniques. This study conducts a thorough analysis of the logit embedding landscape, revealing that both ID and OOD data exhibit a distinct trend. Specifically, we demonstrate that OOD data tend to reside near the center of the logit space. In contrast, ID data tend to be situated farther from the center, predominantly in the positive regions of the logit space, thus forming class-wise clusters along the orthogonal axes that span the logit space. This study highlights the critical role of the DL classifier in differentiating between ID and OOD logits.

1 Introduction

Deep learning (DL) classification models perform well at generalizing from large datasets, achieving superior classification accuracy compared to many alternatives. They deliver highly accurate predictions when the test data aligns with the training data's distribution. However, current DL models are unable to handle out-of-distribution (OOD) data. This limits their application in critical fields such as biomedicine, finance, and autonomous systems.

For instance, consider the scenario in biomedicine where DL models classify bacteria from genome sequences. In such cases, it is crucial to account for the presence of novel bacteria types, which can be considered as OOD instances. Failure to account for these could lead to misclassifying these novel bacteria as known types, potentially resulting in flawed diagnostics or misleading scientific conclusions (Ren et al., 2019). This example underscores the need for more robust DL models to effectively identify and handle OOD data. Neglecting these novel entities may result in their misclassification as known types (Ren et al., 2019).

Recent OOD detection methods predominantly operate under the assumption that a classifier, when trained on ID data, intrinsically maps the logits of OOD samples to a distinct spatial location within the logit landscape, divergent from those of ID instances. Thus, differentiating OOD instances from ID data typically involves assigning high likelihood values to the logit (or softmax) location of the ID samples (Vyas et al., 2018; Lee et al., 2018; Sun et al., 2022; Gomes et al., 2022; Liu et al., 2020; Komini & Girdzijauskas, 2024).

Nevertheless, these strategies lack foundational awareness regarding the specific locational distribution of OOD samples in the embedding space. Consequently, these techniques attempt complicated and computationally intensive density estimations of the ID logits, categorizing those samples that fall beneath a certain likelihood threshold as OOD. Furthermore, the scalability of these methods in relation to the number of ID classes is limited, as they require robust performance in density mapping across all ID classes.

An inaccuracy in the mapping of even a single ID class may lead to a substantial error rate, where the methods could erroneously classify ID instances as OOD. Instead, our study demonstrates that a well-trained DL classifier, incorporating nonlinearities that suppress negative values (e.g., ReLU), systematically maps ID data into well-defined, class-specific clusters with a predictable structure. These ID clusters are situated along orthogonal axes within the positively constrained logit space and are notably separated from the logit space's center. Additionally, we reveal that OOD data are not arbitrarily scattered in the logit space but rather centrally positioned.

Although previous research has explored the separation of OODs and IDs in the logit space (Lee et al., 2018; Liu et al., 2020; Choi et al., 2024; Katz-Samuels et al., 2022b), this work demonstrates a configuration of OODs and IDs logits. The noted positioning of OOD and ID logits lays the groundwork for the possible creation of a binary classifier (OOD from ID), which could lead to simpler yet more effective OOD detection models. This study presents the following contributions:

1. An investigation into the spatial allocation of ID data within the logit space.
2. An empirical validation of the observed ID and OOD logit allocation over many models.

2 Method

2.1 Understanding ID and OOD data in DL models

In exploring ID and OOD data, it is crucial to delineate their distinctions in relation to DL classifiers.

The training dataset is regarded as the optimal empirical representation of ID data. Within this dataset, ID data tend to aggregate into class-specific clusters based on discriminative features corresponding to each class. The exact parameterization of this feature space remains unknown; however, the annotated empirical ID dataset serves as a practical surrogate. Assuming, these discriminative features constitute a multimodal distribution where each mode corresponds to one of the target classes. For ID data classification, DL models are engineered to exploit these discriminative features distribution and the annotated labels to effectively map the ID data into respective class-specific clusters within the logit vector space.

Whenever we encounter data whose features deviate significantly from this defined distribution of discriminative ID features, they are considered OOD. The degree to which the features of OOD data diverge from this distribution determines their shift with the ID features. Consequently, the more distant the OOD features are from the ID distribution, the lesser their resemblance to the ID discriminative features, rendering them less perceptible from the DL model.

More concretely, a DL classifier trained to differentiate cats from dogs will encounter difficulties with photos of horses and wolves. However, because wolves' features are more similar to dogs, wolves represent closer OOD data compared to horses, which are farther from both trained categories.

2.2 Analyzing ID and OOD logits in DL models

Assuming that DL models operate as spatial-invariance pattern-matching mechanisms operating over the distribution of discriminative features, these models produce high positive logits when the data features exhibit strong coherence with the feature representations parameterized during training. To do so, DL models rely on convolutions and matrix multiplications, which are composed of dot-product computations between the model weights and the input data. During training, the objective is to maximize the dot-product values between the model weights and the ID training data.

Furthermore, we assume that the feature representations of OOD data reside in regions significantly distant from the distribution of discriminative features associated with ID data. Consequently, these OOD features exhibit a distributional shift relative to the feature representation parameterized by the DL model. As a result, the dot-products for the OOD data with the model weights remain low, leading to their logits being centrally concentrated both before and after training. Simultaneously, the logits associated with ID data are driven towards higher positive values.

To visually represent the empirical distributions of ID and OOD logits before and after the training phase, we employed a binary classification model using a multilayer perceptron (MLP) model (see Appendix A and fig. 1a).

To minimize any initial biases, weights of DL models are initialized using a centered Gaussian distribution (Glorot & Bengio, 2010; He et al., 2015a), while the biases are set to zero. Furthermore, considering that both ID and OOD data originate from different distributions than the model's initial weight distribution, it is reasonable to assume that they are **distributionally shifted** from the model's initialized weights. This **shift exemplifies as reduced co-variability** between the model's weight parameters and both OOD and ID data samples (see Corollary 1), **consequently leading to negligible covariance**.

Corollary 1. *Assuming the features of the ID and OOD data, as well as the initialized model weights, originate from three distinct distributions with minimal overlap—it follows that the expected dot-product value between the data (ID and OOD) and the DL model's weights is minimal.*

Proof. Given the distribution shift between the data vector \hat{x} and the model's weights ω , their covariance is expected to be minimal, i.e.,

$$|\text{Cov}(\hat{x}, \omega)| < c, \quad \text{where } c \text{ is a sufficiently small positive constant.}$$

Define the covariance between model weights ω and data \hat{x} as

$$|\text{Cov}(\hat{x}, \omega)| = |\mathbb{E}[\langle \hat{x}, \omega \rangle] - \mathbb{E}[\hat{x}] \mathbb{E}[\omega]| < c.$$

Considering the DL model are initialized from a centralized distribution (Glorot & Bengio, 2010; He et al., 2015a), the initial expectation of the model weights is zero ($\mathbb{E}[\omega] = 0$), it follows that

$$|\mathbb{E}[\langle \hat{x}, \omega \rangle]| < c.$$

□

Because DL models fundamentally rely on the dot-product operations between inputs and weights, this result suggests that the logits of both ID and OOD data are centered around zero in the logit space prior to model training (see fig. 1b).

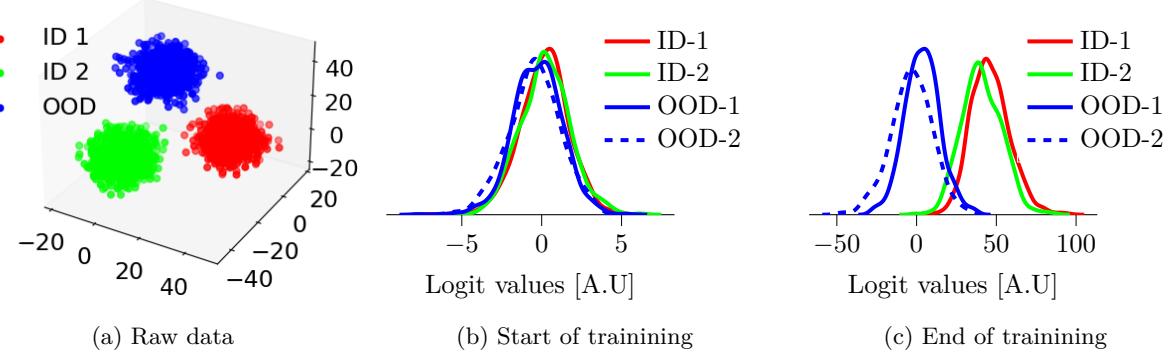


Figure 1: Figure 1a shows raw data sampled from a multimodal Gaussian distribution, utilized as training data for a simple MLP binary classifier depicted in Appendix A. In this figure, red and green points denote ID classes for binary classification, and blue points represent OOD data. Figure 1b and Figure 1c demonstrate kernel density estimations (KDE) across logit cells for both OOD and ID data before and after model training, respectively. In both figures, 'OOD-1' and 'OOD-2' refer to KDEs for OOD data within the first and second logits, while 'ID-1' and 'ID-2' represent KDEs for ID class one data in the first logit cell and ID class two data in the second logit cell, respectively.

2.3 Allocation of ID logits towards positive regions

Training a DL classifier involves utilizing the cross-entropy loss, (*i.e.*, $\mathbf{H}(\mathbf{Y}, \hat{\mathbf{Y}})$), to encourage the prediction ($\hat{\mathbf{Y}}$) to closely align with the ground truth (\mathbf{Y}). When employing one-hot encoding for both $\hat{\mathbf{Y}}$ and \mathbf{Y} , the training objective simplifies to:

$$\mathbf{H}(\mathbf{Y}, \hat{\mathbf{Y}}) = - \sum_i \mathbf{Y}(i) \log(\hat{\mathbf{Y}}(i)) = \underbrace{-\mathbf{Y}(j) \log(\hat{\mathbf{Y}}(j))}_{\mathbf{Y}(j)=1} - \sum_{i,i \neq j} \underbrace{\mathbf{Y}(i) \log(\hat{\mathbf{Y}}(i))}_{\mathbf{Y}(i)=0} = -\log(\hat{\mathbf{Y}}(j)).$$

Eventually, the minimization cross-entropy loss (*i.e.*, $\min[\mathbf{H}(\mathbf{Y}, \hat{\mathbf{Y}})]$) equivilates to maximum likelihood estimation (MLE) (*i.e.*, $\min[-\log(\hat{\mathbf{Y}}(j))]$).

As training progresses, the softmax layer aims to generate a response close to one for the cell corresponding to the correct class (*i.e.*, $\hat{\mathbf{Y}}(j) \rightarrow 1$). Additionally, owing to the inherent property that the softmax output is confined within a simplex (*i.e.*, $\hat{\mathbf{Y}}(j)^\dagger + \sum_{i,i \neq j} \hat{\mathbf{Y}}(i)^\dagger = 1$), the remaining cells are pushed towards values close to zero (*i.e.*, $\hat{\mathbf{Y}}(i)_{i \neq j} \rightarrow 0$). Hence, optimization can be conceptualized as the maximization of the softmax cell corresponding to the correct class and the simultaneous minimization of cells associated with incorrect classes.

This pattern of maximization-minimization is also observed in other classification losses (*i.e.*, Support Vector Machine (Tang, 2015) and Kullback-Leibler divergence (Cui et al., 2024)), which are commonly employed in training DL classification models. This maximization-minimization optimization extends from softmax cells directly to the respective logit cells, as softmax maintains the order of logits. In particular, the logit cell linked to the correct class tries to attain large positive values (see fig. 1c).

However, when suppressing the negative values in an activation layer, the minimization process results in logit values near zero rather than approaching negative values of high magnitudes. Therefore, ID data are projected toward the positive regions of the logit space (see Appendix B for an extended discussion). Given that the logits reach their high positive value for the correct logit cell indicated by the one-hot encoding and approach zero for all other categories, it is evident that the logits for ID samples cluster by class along orthogonal axes within the logit space.

2.4 Central allocation of OOD logits

Since DL classifiers (*i.e.*, ResNet, DenseNet, ViT) are trained using maximum likelihood estimation and their architectures rely on discrete convolutions, which itself relies on dot-products, the training process inherently seeks to maximize the dot-product between the **ID data discriminative features** and the model's parameters (see Appendix B). Furthermore, maximizing the dot-product between two vectors enhances their linear association. Correspondingly, a pronounced linear relationship between two vectors typically results in a higher co-variability between these two vectors. As a result, maximizing the dot-product enables a high degree of covariance and, by extension, cross-correlation. Hence, one can safely assume that there is a positive relationship between the dot-product of two vectors (*i.e.*, $\langle \hat{\mathbf{x}}, \omega \rangle$) and their covariance, given that both metrics assess the degree of alignment between the two vectors.

Prior to training, the initial distributions of the model's weights and any given data are disparate. This disparity typically results in both the covariance and the dot-product between the model's weights and any given data being close to zero due to the lack of any established relationship. Thus, an untrained DL model generally steers any input data towards the center of the logit space (see corollary 1 and fig. 1b).

After the onset of training, the aim is to align the model's weights with the **discriminative features of the ID data**, maximizing the dot-product and **increases** their co-variability, culminating in stronger activation of the logit cell that encodes the correct class (see fig. 1c and Theorem 1 in Appendix B). **Since the discriminative features of OOD and ID data derive from fundamentally different distributions, they inherently exhibit a certain level of distributional shift. This shift is also reflected in the distribution of the OOD features in relation to the model's parameters since the latter are trained to align with the ID discriminative features.** This **shift** implies that the covariance between OOD and **parameters** will likely remain minimal, even post-

train. As a result, their expected dot-product tends to yield smaller magnitudes. Therefore, OOD data tend to remain centered within the logit space even after training (see figs. 1b and 1c).

3 Related Work

Although the differentiation between OOD and ID logits has been extensively studied (Wang et al., 2022; Liu et al., 2020; Hsu et al., 2020; Lee et al., 2018; Ren et al., 2021), existing methods do not anticipate their actual configuration. Conventional OOD detection methods predominantly classify data by first identifying ID samples and subsequently labeling all other samples as OOD by default. A recent empirical investigation has not only highlighted the transferability of ID training strategies to OOD detection but also identified a tangible correlation between the robustness of ID training protocols and OOD detection efficacy (Wenzel et al., 2022). This study suggests that refining ID training methods could unlock potential pathways for enhancing OOD detection. Another study examines the influence of pre-trained Vision Transformers (ViT) (Vaswani et al., 2023) on ImageNet and reports notable improvements in OOD detection performance (Dosovitskiy et al., 2021). Parallel to these observations, another line of research incorporates outlier data — surrogates for OOD samples — within the training phase. This is achieved through an auxiliary loss term that sharpens the contrast between ID and outlier inputs, potentially strengthening OOD detection (Katz-Samuels et al., 2022a; Hendrycks et al., 2019; Wang et al., 2023; Du et al., 2022; Ming et al., 2022). Complementing these approaches, there has been a significant effort to restrict the classification of ID data into a hyperspherical embedding, which intrinsically helps OOD detection(Ming et al., 2023).

Another line of research assumes an inherent separation between OOD and ID logits and tries to devise scoring techniques using solely ID logits or softmax output. The OOD detection works by telling as OOD anything that is not ID. The earliest work on this front assumes clustering of ID logits into a multimodal Gaussian distribution and then tries to utilize Mahalanobis distance (Lee et al., 2018; Ren et al., 2021). More advanced methods try to upgrade the Mahalonobis distance with geometric information using Fisher Information matrix (Gomes et al., 2022) Other works try to perform a data drive density estimation using energy-based models (Liu et al., 2020).

Another promising research demonstrates the utility of enhanced Hopfield networks in amplifying the distinction between ID and OOD data Hofmann et al. (2024). Similarly, another proactive work tries to increase the separability between OOD and ID using kernel principal component analysis on the OOD and ID embeddings (Fang et al., 2024). Last but not least, Zhang et al. (2024) tries to learn the shape of ID feature space using an online expectation maximization, which enhances the detection of OOD post-train.

4 Results

In these experiments, we demonstrate that the OOD logits remain near the center of the logit space both before and after training. In contrast, ID logits consistently gravitate towards clusters around class-specific areas in the positive regions of the logit space. Furthermore, we show that these ID clusters align with the orthogonal axis that spans the logit space embeddings.

OOD vs ID during training: In fig. 2, we empirically illustrate the distribution of ID and OOD logits before and after training. Additionally, we present the evolution of these distributions throughout the training process. To do so, we employed Resnet-9 (He et al., 2015b) with CIFAR-100 (Krizhevsky et al., a) as the ID data and CIFAR-10 (Krizhevsky et al., b) as the OOD data. *Additionally, for the correctly classified ID data, we categorize ‘ID-in’ when the logit values reach their maximum in the cell corresponding to the correct class, as indicated by the one-hot encoding of that class. Conversely, we categorize it as ‘ID-out’ when this condition is not met.*

We represent the empirical distributions of the logit outputs for both ID and OOD samples via kernel density estimation (KDE) (Bishop, 2006). At the beginning of training, one can notice that the densities for both OOD and ID logits are concentrated near zero (see figs. 2a to 2d). While OOD and ID-out logits maintain their central tendency around zero (see fig. 2c) the ID-in logits exhibit a shift towards higher positive values (see fig. 2a). Analyzing the peak (i.e., mode) of each KDE plot (i.e., ID-in, ID-out, and OOD in fig. 2c), it

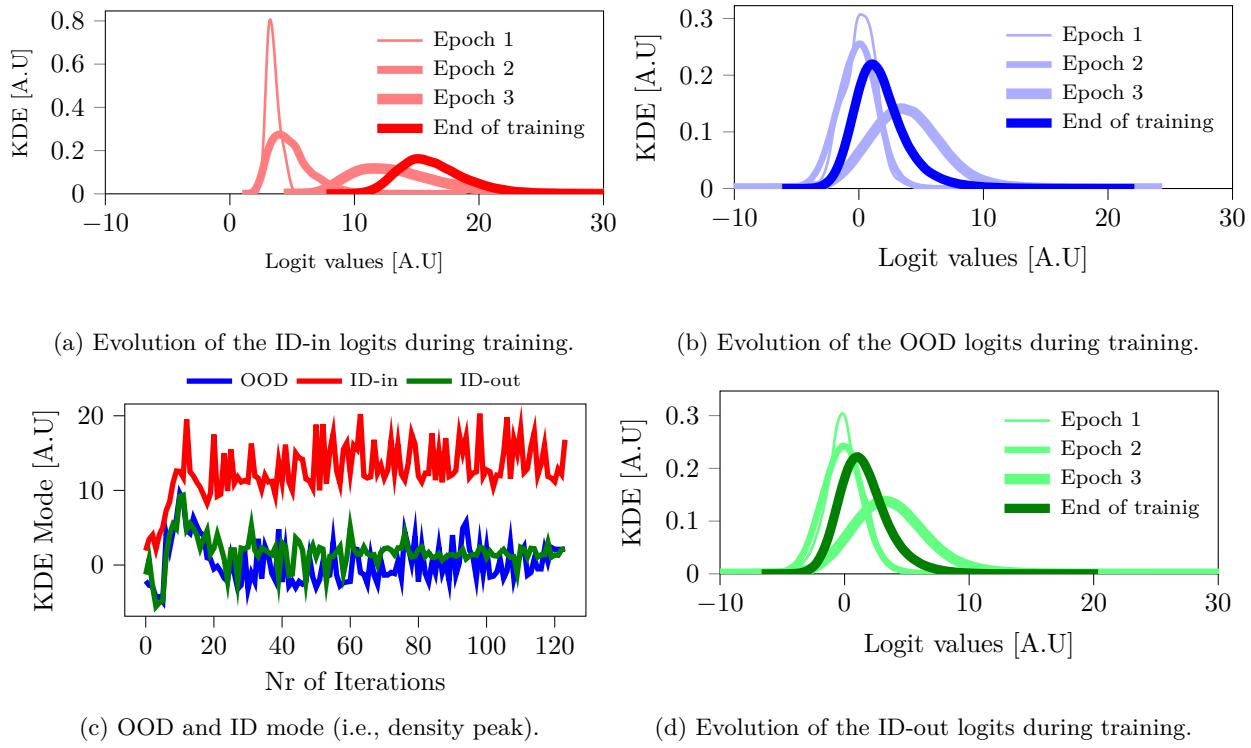


Figure 2: Figure 2a presents the density plot across various epochs for the aggregation of ID-in across all logits, while fig. 2b displays the density plot across different epochs for the aggregation of OOD across logits. Similarly, fig. 2d shows the density plot over different epochs for the aggregation of ID-out across all logits. Since the KDE plots are limited to the first three and the final epochs, we included fig. 2c to provide a comprehensive view of the entire trajectory, featuring the peak (i.e, mode) of the density plot for every epoch.

is evident that ID-in trends towards positive values over time as anticipated by Theorem 1 in Appendix B. Furthermore, the ID-out and OOD logits remain centrally positioned, aligning with our analytical predictions. In addition to the density plots in figs. 2a, 2b and 2d, which illustrate the aggregation of ID-in, ID-out [from training data along with OOD](#) across all logit cells, see Appendix C for detailed visualization of density plots on individual logit cells for a more in-depth analysis.

Effect of activation function: To further understand the configuration of ID and OOD logits, we investigate the impact of various activation functions on a ResNet-34 model. Specifically, we empirically demonstrated this impact by utilizing a selection of activation functions known for their inherent suppression of negative values, including Celu (Barron, 2017), Elu (Clevert et al., 2016), Gelu (Hendrycks & Gimpel, 2023), Selu (Klambauer et al., 2017), Silu (Elfwing et al., 2017), Relu (Hein et al., 2019), Leaky-Relu (Maas et al., 2013), and Mish (Misra, 2020). A ResNet-34 model was trained on the SVHN dataset (Goodfellow et al., 2014) (i.e., ID data), utilizing each activation function. Simultaneously, the CIFAR-10 dataset was used as OOD data.

A ResNet-34 model was trained on the SVHN dataset (Goodfellow et al., 2014), (i.e., ID data), utilizing each activation function (i.e., Celu, Elu, Gelu, Mish, Selu, Silu, ReLU, Leaky ReLU). The model is trained using stochastic gradient descent (SGD) with a cyclical learning rate starting at $lr = 10^{-3}$ with a cosine annealing operation with a periodicity of 200. Furthermore, the momentum is 0.9 while the weight decay $5 * 10^{-4}$. A batch size of 256 is applied for both test and train data. No regularization is applied to the training process, while the training data are augmented with random flipping and cropping.

One can notice that ID-in logits maintain a tendency towards high positive values across all the activation functions (see fig. 3). On the other hand, ID-out and OOD logits are predominantly centralized around

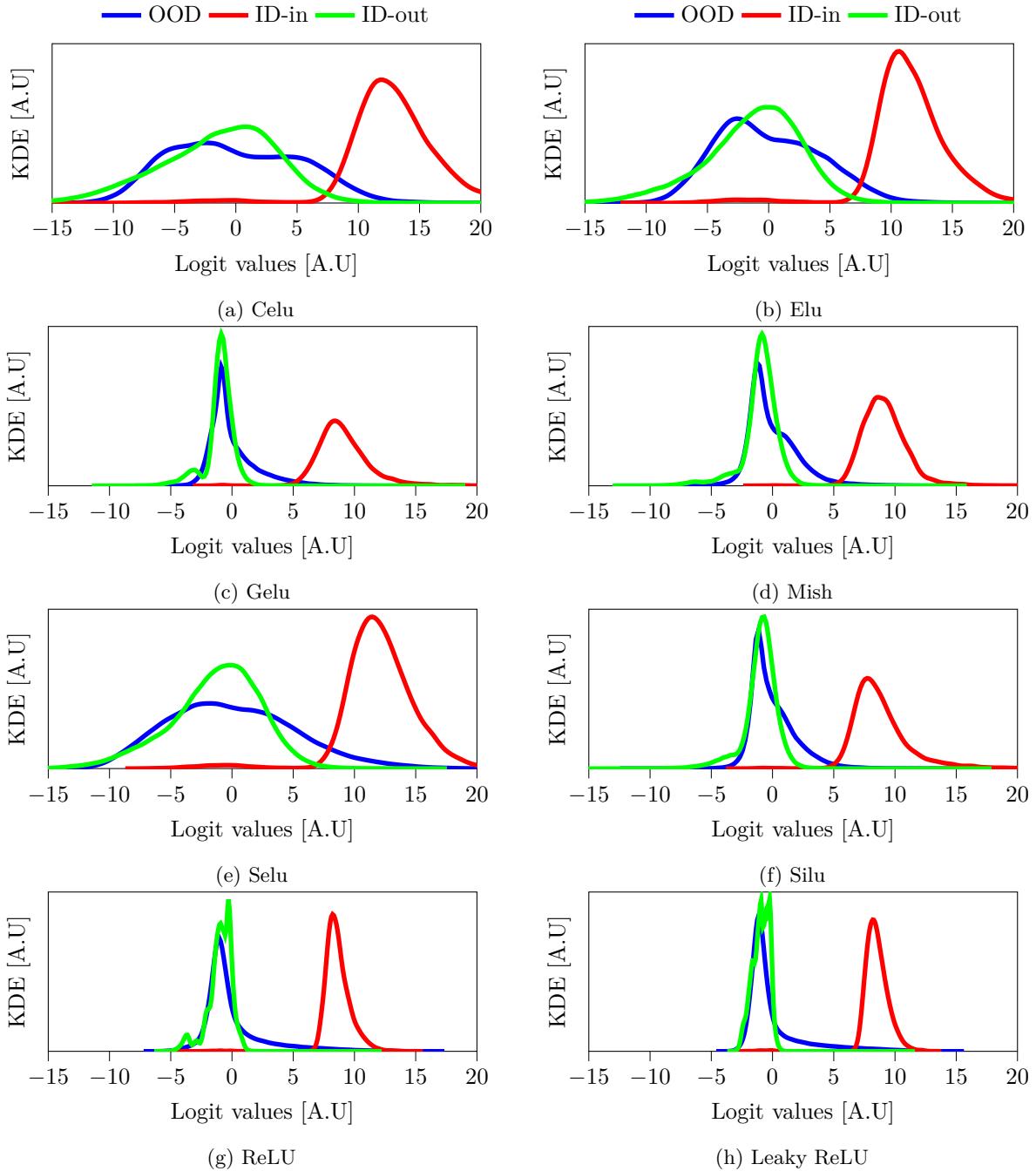


Figure 3: An analysis of the density over logits across eight distinct activation functions that suppress negative values is presented. The ResNet-34 architecture is utilized and trained on the SVHN dataset as the ID data, while the OOD includes CIFAR-10.

zero (see fig. 3). Consequently, despite the application of varying non-linearities, the relative configuration of ID and OOD logits remains similar. In addition to the density plots depicted in fig. 3, which show the distribution of aggregated ID-in, ID-out, and OOD across all logit cells, figs. 10 to 17 in Appendix D provides a detailed visualization of density plots on individual logit cells using various activation functions.

Effect of dropout: Dropout is a pivotal component in enhancing the generalizability of contemporary deep learning methodologies Srivastava et al. (2014). This technique forces the model to yield accurate

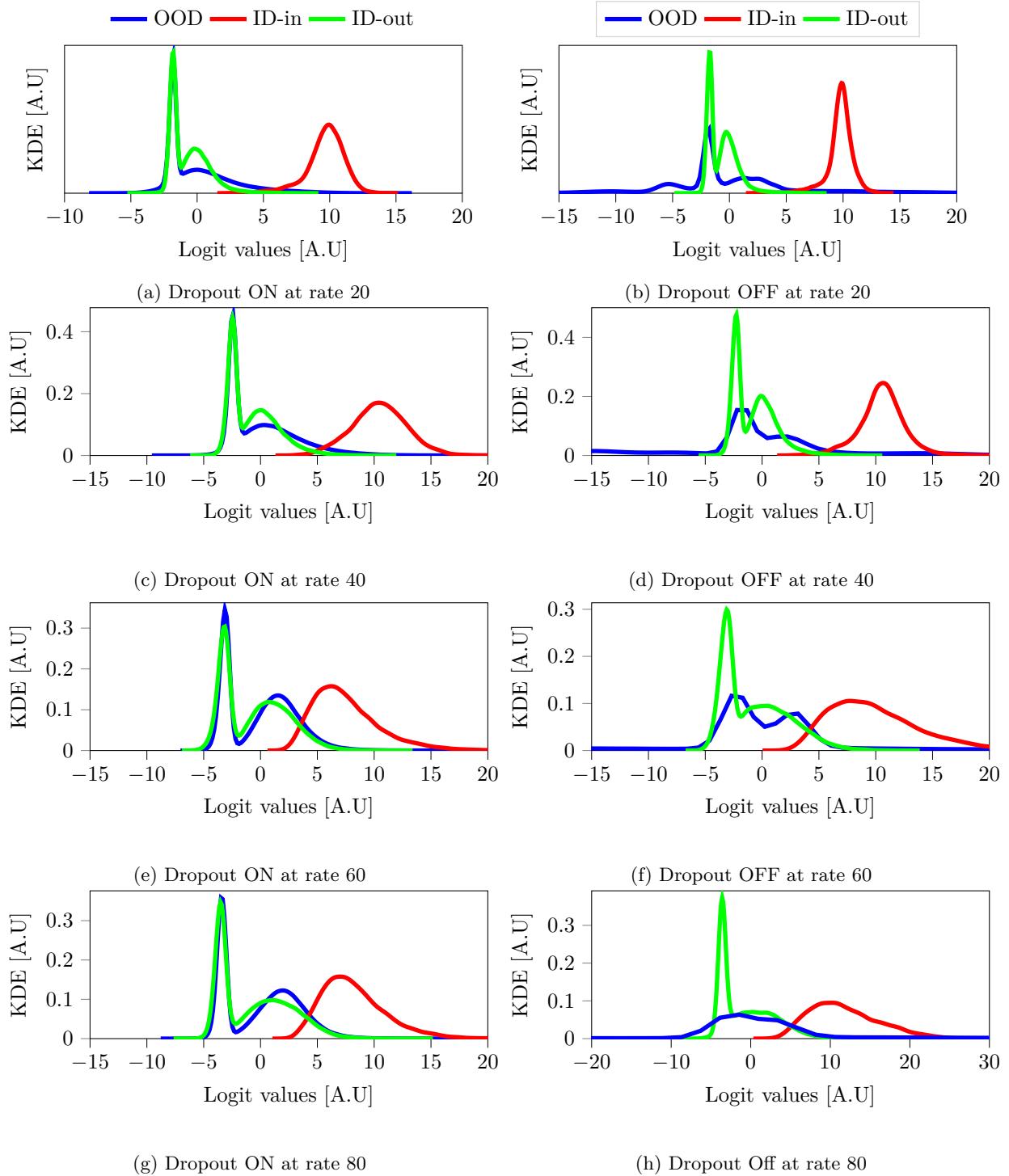


Figure 4: An analysis of the density over logits across eight distinct dropout rates is presented. The ResNet-34 architecture is utilized and trained on the CIFAR-10 dataset as the ID data, while the OOD includes $\{D/\text{CIFAR-10}\}$.

classifications, leveraging a randomly selected subset of its parameters by deactivating a proportion of the gradients during the training phase. Although training involves only a subset of parameters, the entire set is employed during inference. However, the influence of dropout on the differentiation of ID and OOD data has

been relatively understudied, as the main focus of dropout has been optimizing the accuracy of ID datasets rather than maximizing the discriminative capability between OOD and ID datasets. Instead, our study focuses on investigating this discrimination. It reveals that OOD vs. ID behaves differently when using a partial subset of parameters (i.e., dropout ON) versus the entire set (i.e., dropout OFF).

In our experimentation, we employed a ResNet-34 architecture configured with varying dropout ratios of 20%, 40%, 60%, and 80% applied in each convolution layer. The dataset ensemble comprised $\{D\} = \text{SVHN, CIFAR-100, CIFAR-10, Tiny ImageNet (Deng et al., 2009), iSUN (Xu et al., 2015), LSUN (Yu et al., 2016)}$, where the model was trained specifically on CIFAR-10 and utilized $\{D/\text{CIFAR-10}\}$ as the OOD dataset.

One can notice that the anticipated pattern of ID and OOD persist across all levels of dropout utilized. However, an increase in the dropout rate results in an enhanced overlap between ID and OOD predictions (see fig. 4). Moreover, this overlap is more pronounced when the complete set of weights is engaged during inference (i.e., dropout OFF) compared to when only a randomly selected subset is used (i.e., dropout ON) (refer to fig. 4).

Because the overlap between OOD and IID logits increases as the dropout rate rises (see fig. 4), it suggests that applying dropout to the convolutional layers has a detrimental effect on the generalization of the ID discriminative features. Moreover, this empirical observation indicates the unintended introduction of a spurious correlation between the model weights and OOD features.

In addition to the density plots shown in fig. 4, which illustrate the aggregated distribution of ID-in, ID-out, and OOD across all logit cells, a more detailed visualization of density plots on individual logit cells can be found in figs. 18 to 25 of Appendix D.

Experiments on different classifiers: The analysis of ID and OOD logits has been expanded across various DL classifier models. Our study examines various iterations of DenseNet (Huang et al., 2018), specifically versions 121, 161, 169, and 201, as well as ResNet (He et al., 2015b), encompassing versions 18, 34, 50, 101, and 152. Furthermore, the utilized experimental dataset comprises $\{D\} = \{\text{SVHN, CIFAR-100, CIFAR-10, Tiny ImageNet, iSUN, LSUN}\}$.

Densenet and ResNet models are trained using SGD with a cyclical learning rate starting at $lr = 10^{-3}$ with a cosine annealing operation with a periodicity of 200. Furthermore, the momentum is 0.9 while the weight decay $5 * 10^{-4}$. A batch size of 256 is applied for both test and train data, while the number of epochs is 200. ReLU is utilized as an activation function for every layer. No regularization is applied to the training process, while the training data are augmented with random flipping and cropping. Each version of Densenet and ResNet undergoes separate training on CIFAR-10 and SVHN as ID datasets.

When CIFAR-10 is utilized as ID, the remaining datasets are employed as OOD data, specifically $\{D\}$ without CIFAR-10 (i.e., $\{D\}/\text{CIFAR-10}$) is utilized as OOD. Similarly, when SVHN is utilized as ID, the remaining datasets are employed as OOD data, specifically $\{D\}$ without SVHN (i.e., $\{D\}/\text{SVHN}$) is utilized as OOD.

Observations indicate that the ID-in logits consistently tend toward higher positive values across various versions of DenseNet (see fig. 5) and ResNet (see fig. 6). Contrarily, ID-out and OOD logits tend to be concentrated around zero (as shown in fig. 3).

Notice that the density plots in figs. 5 and 6 demonstrate the distribution of the ID-in, ID-out, and OOD for all logit cells aggregated together. For thorough visual representations on a per-logit-cell basis, across all different versions of DensetNet and ResNet see figs. 26 to 43 in Appendix D.

Experiments on different vision transformers: Contrary to traditional convolutional neural networks (e.g., DenseNet, ResNet), which process image patches exclusively on a spatial level, vision transformers (ViT) incorporate an additional component of interleaved processing among patches through attention mechanism (Dosovitskiy et al., 2021). To examine the effects of this interleaved processing on the arrangement of OOD and ID logits, we carried out experiments with various ViT configurations, including the base (ViT-B) and large (ViT-L) models, each with two different patch sizes: 16x16 and 32x32 pixels.

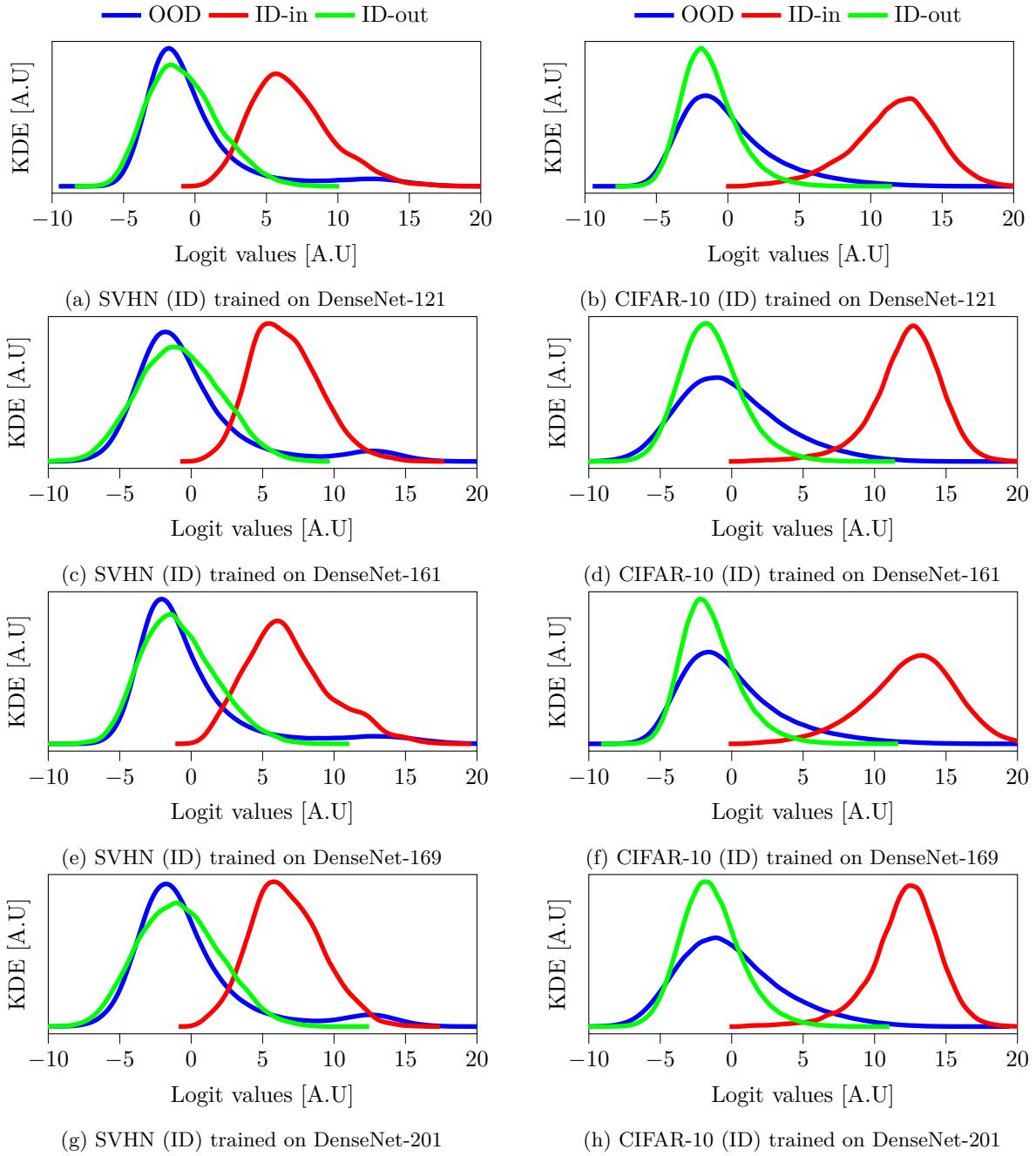


Figure 5: Logit densities across various DenseNet architectures trained on SVHN and CIFAR-10.

Furthermore, the utilized experimental dataset comprises $\{D\} = \{\text{SVHN}, \text{CIFAR-100}, \text{CIFAR-10}, \text{Tiny ImageNet}, \text{iSUN}, \text{LSUN}\}$. Each model undergoes separate training on CIFAR-10 and SVHN as ID datasets. The remaining datasets are employed as OOD data, specifically $\{D\}/\text{CIFAR-10}$ and $\{D\}/\text{SVHN}$.

In fig. 7, one can notice that for all versions of the ViT, ID-in logits converge towards higher positive values as expected. Contrarily, the logits for both the ID-out and OOD samples predominantly cluster around the center of the logit space. Therefore, the intertwined processing among patches in ViT does not alter the anticipated configuration of OOD and ID logits, as it is inherently composed of a dot-product operation. Observe that the density plots shown in fig. 7 depict the spread of ID-in, ID-out, and OOD aggregations over

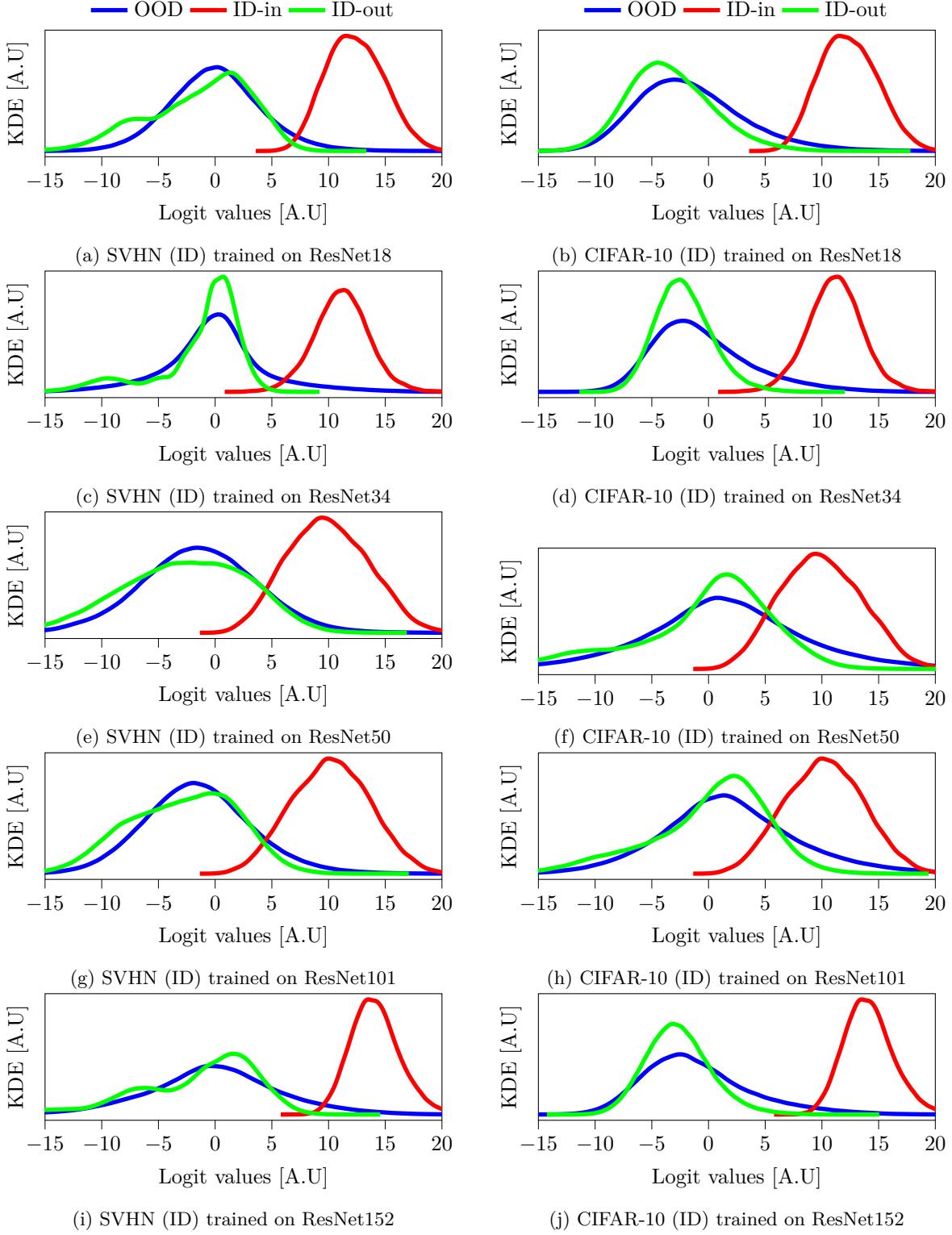


Figure 6: Logit densities across various ResNet architectures trained on SVHN and CIFAR-10.

all logit cells. For a detailed visual analysis of each logit cell, refer to the various versions of ViT illustrated in figs. 44 to 51 in Appendix D.

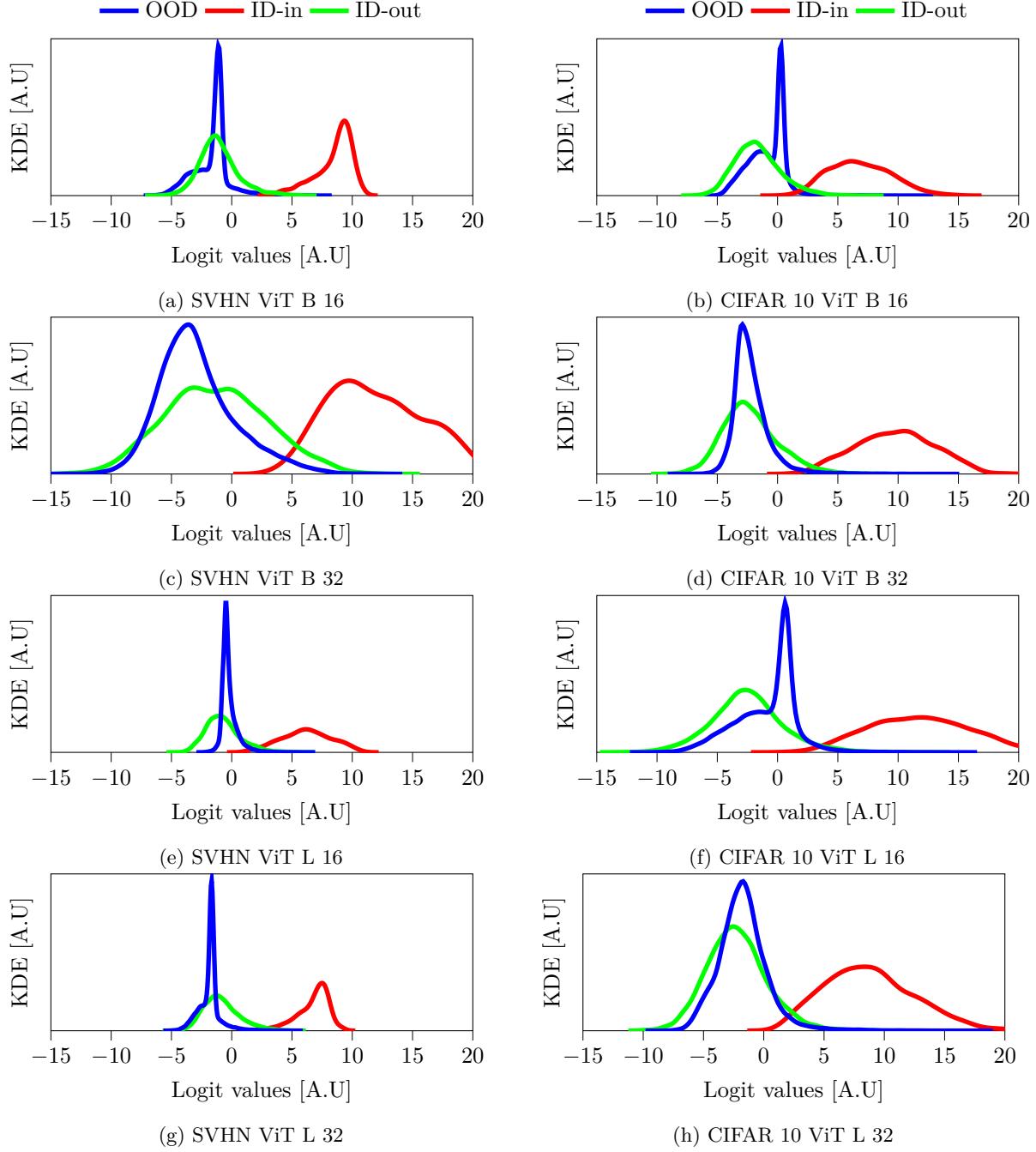


Figure 7: An analysis of the density over logits across distinct ViT architecture trained on the SVHN and CIFAR-10 dataset as the ID data, while the OOD includes $\{D\}/\text{SVHN}$ and $\{D\}/\text{CIFAR-10}$ respectively.

5 Conclusion

While current research on OOD detection focuses on developing new methods that naturally give higher scores to ID data and, by default, lower scores to OOD samples, this study concentrates on analyzing the differences between OOD and ID logit distributions.

Specifically, we demonstrated the anticipated configuration of OODs and IDs logits, i.e., that ID logits are clustered by class towards the positive region of the logit space, aligning with the orthogonal axis that spans

this space. Additionally, OOD logits remain consistently distinct from ID logits, clustering around the center of the logit space.

This behavior of OOD and ID logits is consistent across various architectures (i.e., convolutional neural network models and vision transformers) and activation functions tested on a set of large and diverse OOD data. However, elevated dropout rates on the convolutional layers have been identified as a significant factor in increasing overlap between OOD and ID samples.

As a future direction, the observed patterns within OOD, ID-in, and ID-out logits indicate the potential for a novel approach that leverages ID-out logits as proxies for OOD instances. This approach will facilitate the development of a binary classifier neural network designed to differentiate between OOD and ID samples, employing ID-out logits as representative proxies for OOD instances. Consequently, this method addresses OOD detection as a straightforward classification challenge, thereby mitigating the need for threshold-based discrimination methods.

An additional crucial application of the observed logit configuration is the detection of IID data shifts. Since IID values are typically oriented towards positive values along the corresponding axis, this characteristic can be utilized to develop a more accurate and scalable approximation of the Wasserstein distance. Consequently, enabling a more sensitive metric to detect shifts toward the center of the logit space in the IID test data.

References

- Jonathan T. Barron. Continuously differentiable exponential linear units, 2017.
- Christopher M. Bishop. *Pattern Recognition and Machine Learning (Information Science and Statistics)*. Springer-Verlag, Berlin, Heidelberg, 2006. ISBN 0387310738.
- Caroline Choi, Fahim Tajwar, Yoonho Lee, Huaxiu Yao, Ananya Kumar, and Chelsea Finn. Conservative prediction via data-driven confidence minimization, 2024. URL <https://arxiv.org/abs/2306.04974>.
- Djork-Arné Clevert, Thomas Unterthiner, and Sepp Hochreiter. Fast and accurate deep network learning by exponential linear units (elus), 2016.
- Jiequan Cui, Zhuotao Tian, Zhisheng Zhong, Xiaojuan Qi, Bei Yu, and Hanwang Zhang. Decoupled kullback-leibler divergence loss, 2024. URL <https://arxiv.org/abs/2305.13948>.
- Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pp. 248–255. Ieee, 2009.
- Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations*, 2021. URL <https://openreview.net/forum?id=YicbFdNTTy>.
- Xuefeng Du, Zhaoning Wang, Mu Cai, and Yixuan Li. Vos: Learning what you don’t know by virtual outlier synthesis, 2022.
- Stefan Elfwing, Eiji Uchibe, and Kenji Doya. Sigmoid-weighted linear units for neural network function approximation in reinforcement learning, 2017.
- Kun Fang, Qinghua Tao, Kexin Lv, Mingzhen He, Xiaolin Huang, and Jie Yang. Kernel pca for out-of-distribution detection, 2024. URL <https://arxiv.org/abs/2402.02949>.
- Xavier Glorot and Yoshua Bengio. Understanding the difficulty of training deep feedforward neural networks. In Yee Whye Teh and D. Mike Titterington (eds.), *AISTATS*, volume 9 of *JMLR Proceedings*, pp. 249–256. JMLR.org, 2010. URL <http://dblp.uni-trier.de/db/journals/jmlr/jmlrp9.html#GlorotB10>.
- Eduardo Dadalto Camara Gomes, Florence Alberge, Pierre Duhamel, and Pablo Piantanida. Igeod: An information geometry approach to out-of-distribution detection, 2022.

- Ian J. Goodfellow, Yaroslav Bulatov, Julian Ibarz, Sacha Arnoud, and Vinay Shet. Multi-digit number recognition from street view imagery using deep convolutional neural networks. 2014.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Delving deep into rectifiers: Surpassing human-level performance on imagenet classification, 2015a.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition, 2015b.
- Matthias Hein, Maksym Andriushchenko, and Julian Bitterwolf. Why relu networks yield high-confidence predictions far away from the training data and how to mitigate the problem, 2019.
- Dan Hendrycks and Kevin Gimpel. Gaussian error linear units (gelus), 2023.
- Dan Hendrycks, Mantas Mazeika, and Thomas Dietterich. Deep anomaly detection with outlier exposure, 2019.
- Claus Hofmann, Simon Schmid, Bernhard Lehner, Daniel Klotz, and Sepp Hochreiter. Energy-based hopfield boosting for out-of-distribution detection, 2024. URL <https://arxiv.org/abs/2405.08766>.
- Yen-Chang Hsu, Yilin Shen, Hongxia Jin, and Zsolt Kira. Generalized odin: Detecting out-of-distribution image without learning from out-of-distribution data, 2020.
- Gao Huang, Zhuang Liu, Laurens van der Maaten, and Kilian Q. Weinberger. Densely connected convolutional networks, 2018.
- Julian Katz-Samuels, Julia Nakhleh, Robert Nowak, and Yixuan Li. Training ood detectors in their natural habitats, 2022a.
- Julian Katz-Samuels, Julia Nakhleh, Robert Nowak, and Yixuan Li. Training ood detectors in their natural habitats, 2022b. URL <https://arxiv.org/abs/2202.03299>.
- Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization, 2017.
- Günter Klambauer, Thomas Unterthiner, Andreas Mayr, and Sepp Hochreiter. Self-normalizing neural networks, 2017.
- Vangjush Komini and Sarunas Girdzijauskas. Integrating logit space embeddings for reliable out-of-distribution detection, 2024.
- Alex Krizhevsky, Vinod Nair, and Geoffrey Hinton. Cifar-100 (canadian institute for advanced research). a.
- Alex Krizhevsky, Vinod Nair, and Geoffrey Hinton. Cifar-10 (canadian institute for advanced research). b.
- Kimin Lee, Kibok Lee, Honglak Lee, and Jinwoo Shin. A simple unified framework for detecting out-of-distribution samples and adversarial attacks, 2018.
- Seongwoo Lim, Won Jo, Joohyung Lee, and Jaesik Choi. Pathwise explanation of ReLU neural networks. In Sanjoy Dasgupta, Stephan Mandt, and Yingzhen Li (eds.), *Proceedings of The 27th International Conference on Artificial Intelligence and Statistics*, volume 238 of *Proceedings of Machine Learning Research*, pp. 4645–4653. PMLR, 02–04 May 2024. URL <https://proceedings.mlr.press/v238/lim24a.html>.
- Weitang Liu, Xiaoyun Wang, John D. Owens, and Yixuan Li. Energy-based out-of-distribution detection, 2020.
- Andrew L Maas, Awni Y Hannun, Andrew Y Ng, et al. Rectifier nonlinearities improve neural network acoustic models. In *Proc. icml*, volume 30, pp. 3. Atlanta, GA, 2013.
- Yifei Ming, Ying Fan, and Yixuan Li. Poem: Out-of-distribution detection with posterior sampling, 2022.
- Yifei Ming, Yiyou Sun, Ousmane Dia, and Yixuan Li. How to exploit hyperspherical embeddings for out-of-distribution detection?, 2023.

- Diganta Misra. Mish: A self regularized non-monotonic activation function, 2020.
- Jie Ren, Peter J. Liu, Emily Fertig, Jasper Snoek, Ryan Poplin, Mark A. DePristo, Joshua V. Dillon, and Balaji Lakshminarayanan. Likelihood ratios for out-of-distribution detection, 2019.
- Jie Ren, Stanislav Fort, Jeremiah Liu, Abhijit Guha Roy, Shreyas Padhy, and Balaji Lakshminarayanan. A simple fix to mahalanobis distance for improving near-ood detection, 2021.
- Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. Dropout: A simple way to prevent neural networks from overfitting. *Journal of Machine Learning Research*, 15(56):1929–1958, 2014. URL <http://jmlr.org/papers/v15/srivastava14a.html>.
- Yiyou Sun, Yifei Ming, Xiaojin Zhu, and Yixuan Li. Out-of-distribution detection with deep nearest neighbors, 2022.
- Yichuan Tang. Deep learning using linear support vector machines, 2015. URL <https://arxiv.org/abs/1306.0239>.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need, 2023.
- Apoorv Vyas, Nataraj Jammalamadaka, Xia Zhu, Dipankar Das, Bharat Kaul, and Theodore L. Willke. Out-of-distribution detection using an ensemble of self supervised leave-out classifiers, 2018.
- Haoqi Wang, Zhizhong Li, Litong Feng, and Wayne Zhang. Vim: Out-of-distribution with virtual-logit matching, 2022. URL <https://arxiv.org/abs/2203.10807>.
- Qizhou Wang, Junjie Ye, Feng Liu, Quanyu Dai, Marcus Kalander, Tongliang Liu, Jianye Hao, and Bo Han. Out-of-distribution detection with implicit outlier transformation, 2023.
- Florian Wenzel, Andrea Dittadi, Peter Vincent Gehler, Carl-Johann Simon-Gabriel, Max Horn, Dominik Zietlow, David Kernert, Chris Russell, Thomas Brox, Bernt Schiele, Bernhard Schölkopf, and Francesco Locatello. Assaying out-of-distribution generalization in transfer learning, 2022.
- Pingmei Xu, Krista A Ehinger, Yinda Zhang, Adam Finkelstein, Sanjeev R. Kulkarni, and Jianxiong Xiao. Turkergaze: Crowdsourcing saliency with webcam based eye tracking, 2015.
- Fisher Yu, Ari Seff, Yinda Zhang, Shuran Song, Thomas Funkhouser, and Jianxiong Xiao. Lsun: Construction of a large-scale image dataset using deep learning with humans in the loop, 2016.
- Yonggang Zhang, Jie Lu, Bo Peng, Zhen Fang, and Yiu ming Cheung. Learning to shape in-distribution feature space for out-of-distribution detection. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024. URL <https://openreview.net/forum?id=1Du3mMP5YN>.

A Toy example

The training configuration for the model outlined in table 1 includes a batch size of 64, a learning rate of 0.001, and 30 training epochs. To combat overfitting, a dropout rate of 0.8 is employed.

Table 1: Architecture of the MLP model.

Layer Type	Output Size	Additional Information
Linear	2048	in_features=3
ReLU	2048	-
Dropout	2048	p=0.8
Linear	2048	in_features=2048
ReLU	2048	-
Dropout	2048	p=0.8
Linear	2	in_features=2048

B In distribution positioning in the logit space during training

Theorem 1. *In the training process of a deep learning classifier utilizing an activation function that suppresses negative values, the logit corresponding to the true class (i.e., ID-in denoted by $\hat{L}(j)$), attain big positive magnitude ($\hat{L}(j) \rightarrow +\infty$). Simultaneously, the logits representing the incorrect classes (i.e., ID-out denoted by $\hat{L}(i)$ for $i \neq j$) converge towards minimal magnitude values ($\hat{L}(i)_{i \neq j} \rightarrow 0$).*

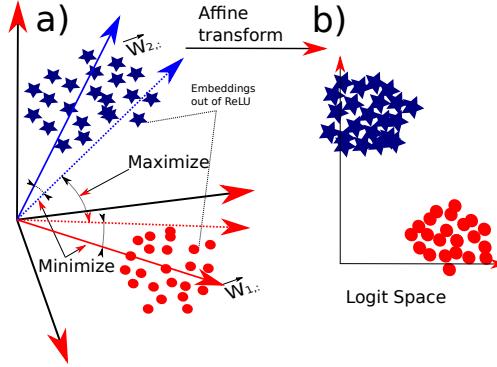


Figure 8: This toy example shows the separation of ID in a binary classification task. Figure a) contains the embeddings (E) rectified with a ReLU. Figure b) shows the linear separation of class-wise clustering of ID data logits (\hat{L}). The smaller the angle between \vec{E} and $\vec{W}_{1,:}$, the higher the dot-product $\langle W_{1,i}, E_i \rangle$ Figure a); thus the more distanced from the center the ID logits are (Figure b). The bigger the angle between (\vec{E}) and $\vec{W}_{2,:}$, the higher the dot-product $\langle \vec{W}_{2,i}, \vec{E}_i \rangle$ (see, fig. 8 a), the more compact the ID logits are.

Proof. To establish the constraint towards zero for the logit cells not corresponding to the correct class (i.e., ID-out $\hat{L}(i)_{i \neq j} \rightarrow 0$), it is crucial to acknowledge that the predecessor latent space ($\hat{E}(i)$) is confined to positive values as the negative values are suppressed (see, fig. 8.a). The layer preceding the softmax constitutes a linear transformation of the data from high-dimensional embeddings (\hat{E}) to the logit space ($\hat{L} = \hat{E} \times W$, where \times denotes matrix multiplication) with dimensions aligning with the number of specified classes (see, fig. 8.b). Since the optimizer seeks maximum response for the logit cell $\hat{L}[i]$ (i.e., ID-in), it aims to maximize the dot-product $\arg \max_{W[i,:]} (\hat{E}[:,], W[i,:])^1$, s.t : $\hat{E}[:] \geq 0$.

Considering the embeddings $\hat{E}[:]$ and $W[i,:]$ as vectors in the vector space (see, fig. 8a), maximizing $\langle \vec{E}[:,], \vec{W}[i,:] \rangle$ results in the minimization of the angle between $\vec{E}[:]$ and $\vec{W}[i,:]$ (i.e., $\min \angle(\vec{W}[i,:], \vec{E}[:])$) while

¹ $\langle \cdot, \cdot \rangle$ indicates the dot-product

ensuring the former always remains in the positive regions. The optimization aims to maintain the direction of the vector $\vec{W}[i,:]$ akin to the cluster of vectors $\vec{E}[:]$, specifically within the positive regions (see, fig. 8.a).

Moreover, the optimization aims to achieve a minimum response for every other logit cell $\hat{L}[j \neq i]$ that does not correspond to the correct class, expressed as $\arg \min_{W[j \neq i,:]} \langle \hat{E}[:, W[j \neq i,:]] \rangle$, subject to the constraint $\hat{E}[:] \geq 0$. In essence, it seeks to maximize the angle between $\vec{W}[j \neq i,:]$ and the cluster of vector data $\vec{E}[:]$ (*i.e.*, $\max \angle(\vec{W}[j \neq i,:], \vec{E}[:])$) (see, fig. 8.a).

Thus, the clusters associated with different classes try to attain maximum angular separation from one another, leading the parameter vectors $\vec{W}[i,:]$ to align accordingly. Given that all vectors $\vec{E}[:]$ are angularly separated within the positive region, the maximum angle between these two vectors approaches perpendicularity (see Lemma 1) asymptotically. Consequently, the minimized logit values ($\arg \min(\vec{W}[j \neq i,:], \vec{E}[:]) \approx 0$) would asymptotically approach zero during training.

Consequently, the asymptotic behavior of the data configuration in the logit space compels the data points to form compact clusters far from the center of the space, corresponding to their respective classes. This process leads to the minimization of interclass distances and the maximization of intraclass distances. \square

Lemma 1. *In the positive region of a high-dimensional space, the maximum angle that two vectors can attain is perpendicular.*

Proof. One way to establish this lemma involves employing the concept of cosine similarity. Let us consider two arbitrary vectors in an N-dimensional space, denoted as X and Y , where $X, Y \in \mathbb{R}^N$. The cosine similarity between these vectors is defined as:

$$\cos_{\text{sim}}(X, Y) = \frac{\sum_{i=1}^N X_i Y_i}{\|X\| \|Y\|} \quad (1)$$

Given that both vectors reside in the positive region of the vector space, meaning that each component of the vectors satisfies $X_i \geq 0, Y_i \geq 0 \forall i \in [1,..,N]$, it is evident that any two vectors in the positive region cannot yield a negative value for the cosine similarity. This is because the numerator, representing the dot-product of the vectors, comprises products of non-negative components. Consequently, the numerator cannot be negative. Therefore, the minimum value that $\cos_{\text{sim}}(X, Y)$ can attain is zero, corresponding to perpendicular vectors. \square

Theorem 1 does not imply that all weights in the model must be non-negative after training. It's important to note that the activation function suppresses negative values, leading to nullified gradients for these outputs that do not contribute to the update process during training. Typically, such negative outputs originate from dot-products involving negative weights, suggesting that these negative weights might not consistently be updated during training sessions. Furthermore, since the predictive response predominantly depends on pathways utilizing positive weights (Lim et al., 2024), these positive weights are more frequently adjusted during the optimization process as outlined in Theorem 1.

C Experimentation on CIFAR-100 (ID) vs CIFAR-10 (OOD)

Resnet-9 is trained using stochastic gradient descent (SGD) with a learning rate starting at $lr = 10^{-1}$. The batch size is 256, and the number of epochs is 200. The learning rate is decimated every quarter of epochs. Within each quarter, the learning rate is scheduled using a 1cycle learning rate. ReLU is utilized as an activation function for every layer. No regularization is applied to the training process, while the training data are augmented with random flipping and cropping. In figs. 9a and 9b, we present the distributions of logit values for ID samples, with the former displaying densities corresponding to the ID-in and the latter for the ID-out. OOD densities are depicted for each logit in fig. 9c.

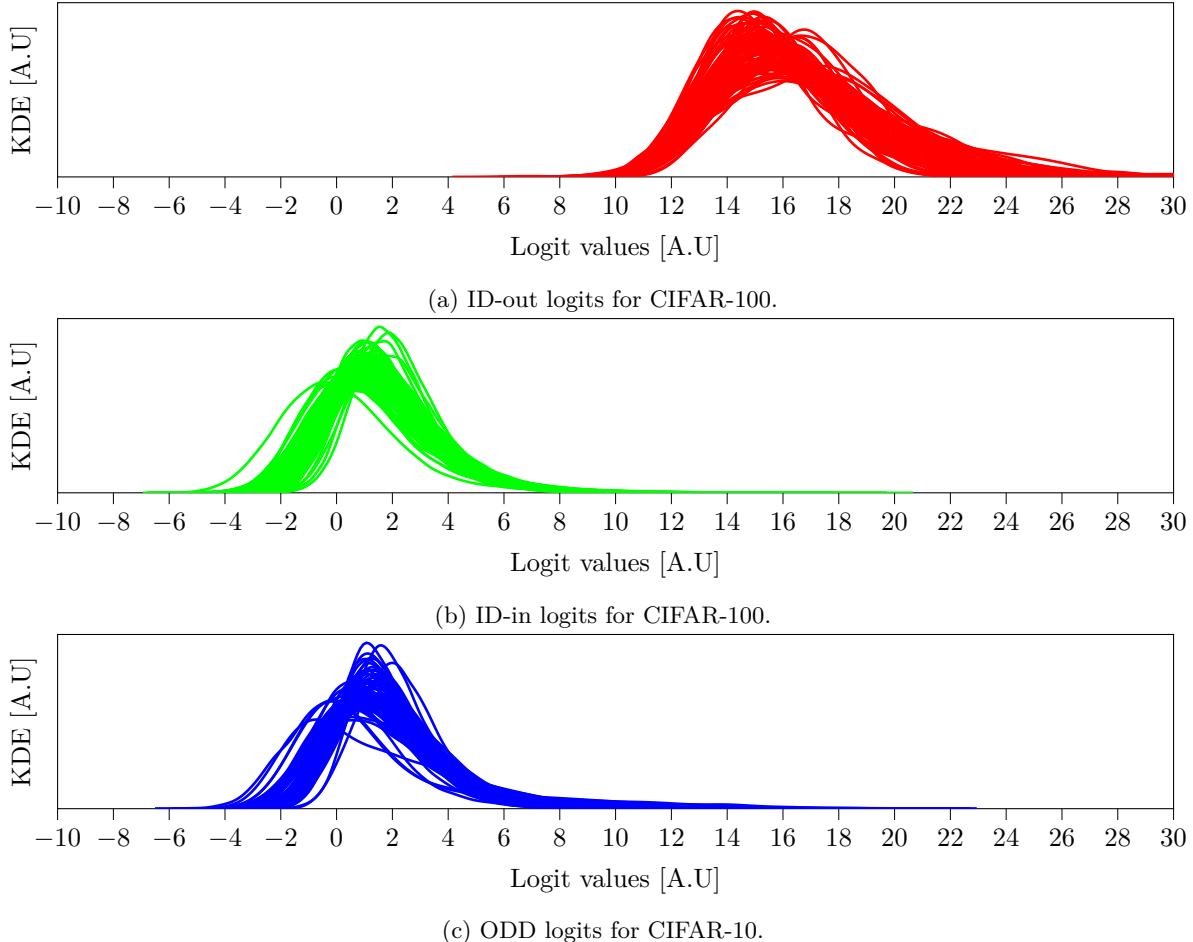


Figure 9: KDE response CIFAR-100 (ID) vs CIFAR-10 (OOD) while using Resnet-9 with ReLU activation function.

D Detailed visualization of density plots on individual logit cells

Figures 10 to 17 showcase a detailed visualization of the ID and OOD logits for each cell across different types of activation functions.

$$\text{Relu: } f(x) = \max(0, x)$$

$$\text{Celu: } f(x) = \max(0, x) + \min(0, \alpha(e^{x/\alpha} - 1))$$

$$\text{Elu: } f(x) = \begin{cases} x & \text{if } x > 0 \\ \alpha(e^x - 1) & \text{if } x \leq 0 \end{cases}$$

$$\text{GELU: } f(x) = x\Phi(x)$$

where $\Phi(x)$ is the cumulative distribution function of the standard Gaussian distribution:

$$\Phi(x) = \frac{1}{2} \left[1 + \operatorname{erf} \left(\frac{x}{\sqrt{2}} \right) \right]$$

$$\text{Selu: } f(x) = \lambda \begin{cases} x & \text{if } x > 0 \\ \alpha e^x - \alpha & \text{if } x \leq 0 \end{cases}$$

$$\text{Silu: } f(x) = \frac{x}{1 + e^{-x}}$$

$$\text{Leaky-Relu: } f(x) = \begin{cases} x & \text{if } x > 0 \\ \alpha x & \text{if } x \leq 0 \end{cases}$$

$$\text{Mish: } f(x) = x \tanh(\ln(1 + e^x))$$

Similarly, figs. 18 to 25 provide a comprehensive visualization of the ID and OOD logits for each cell, delineating the impacts of various dropout rates.

Moreover, figs. 26 to 43 offer a detailed visualization of the ID and OOD logits for each cell across various versions of Densenet and Resnet.

Last but not least figs. 44 to 51 showcase a detailed visualization of the ID and OOD logits for each cell across different variants of Vision Transformers for comprehensive comparative analysis.

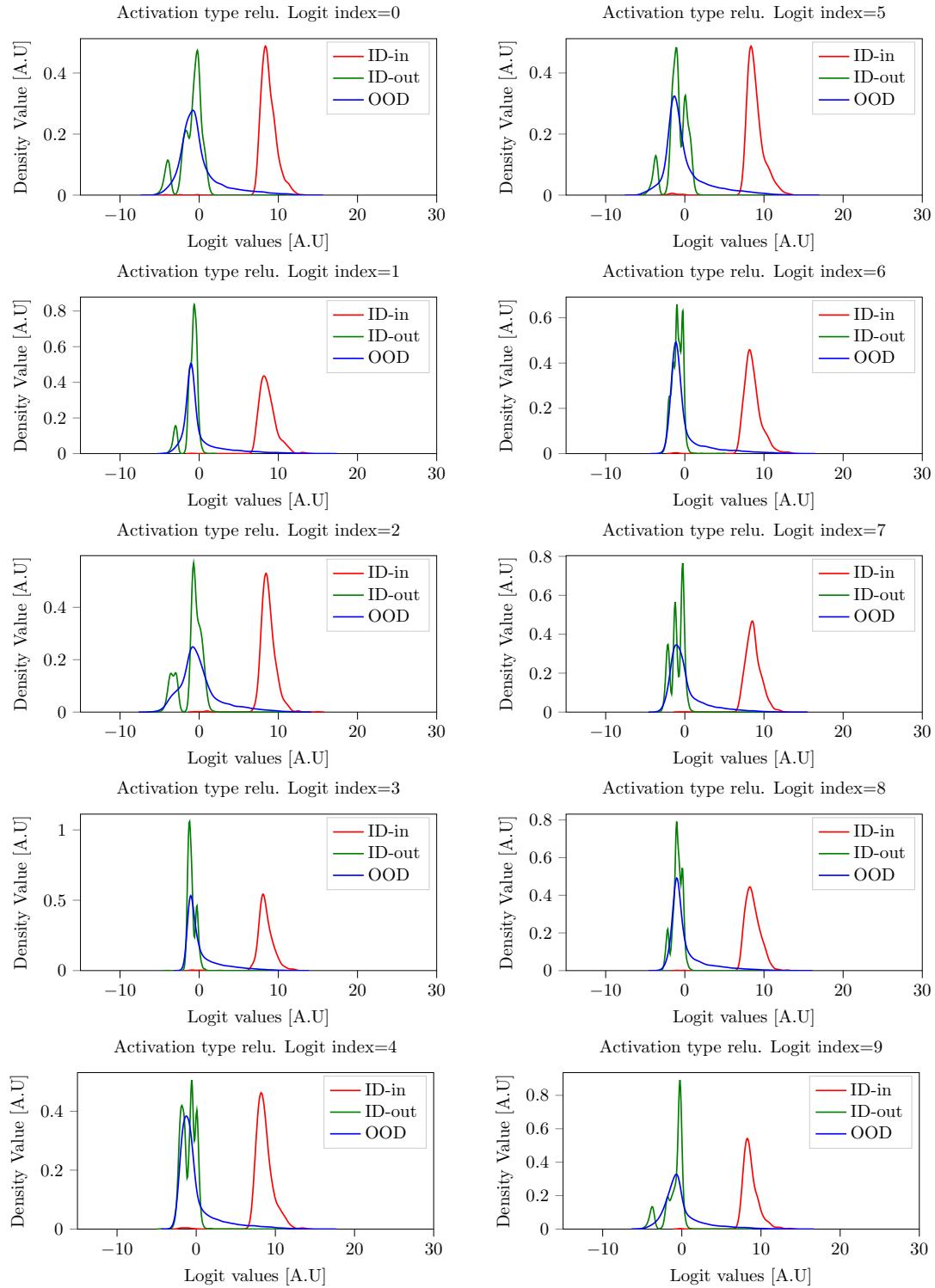


Figure 10: Densities over each logit cell from a Resnet-34 classifier with ReLU activation.

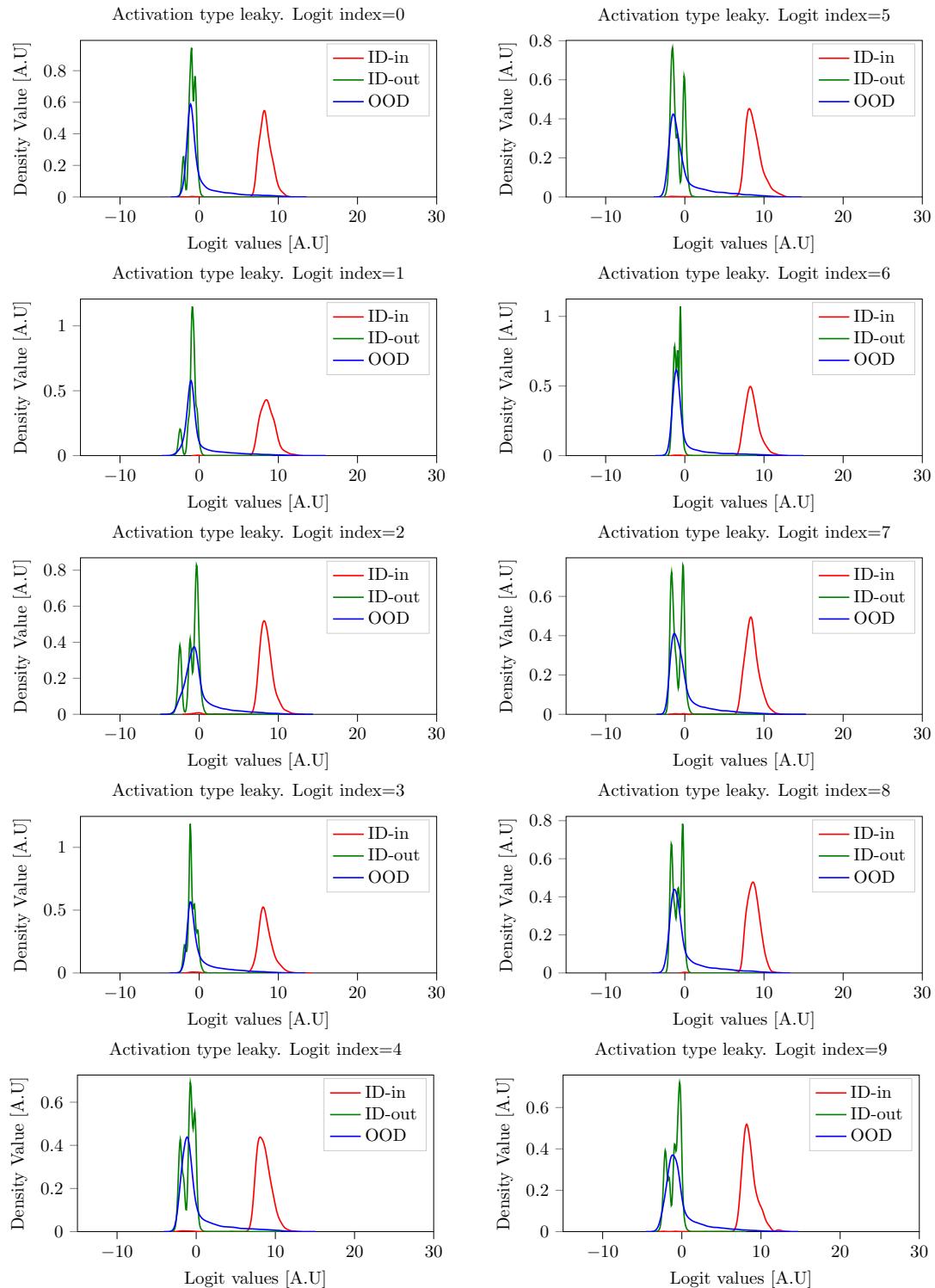


Figure 11: Densities over each logit cell from a Resnet-34 classifier with Leaky Relu activation.

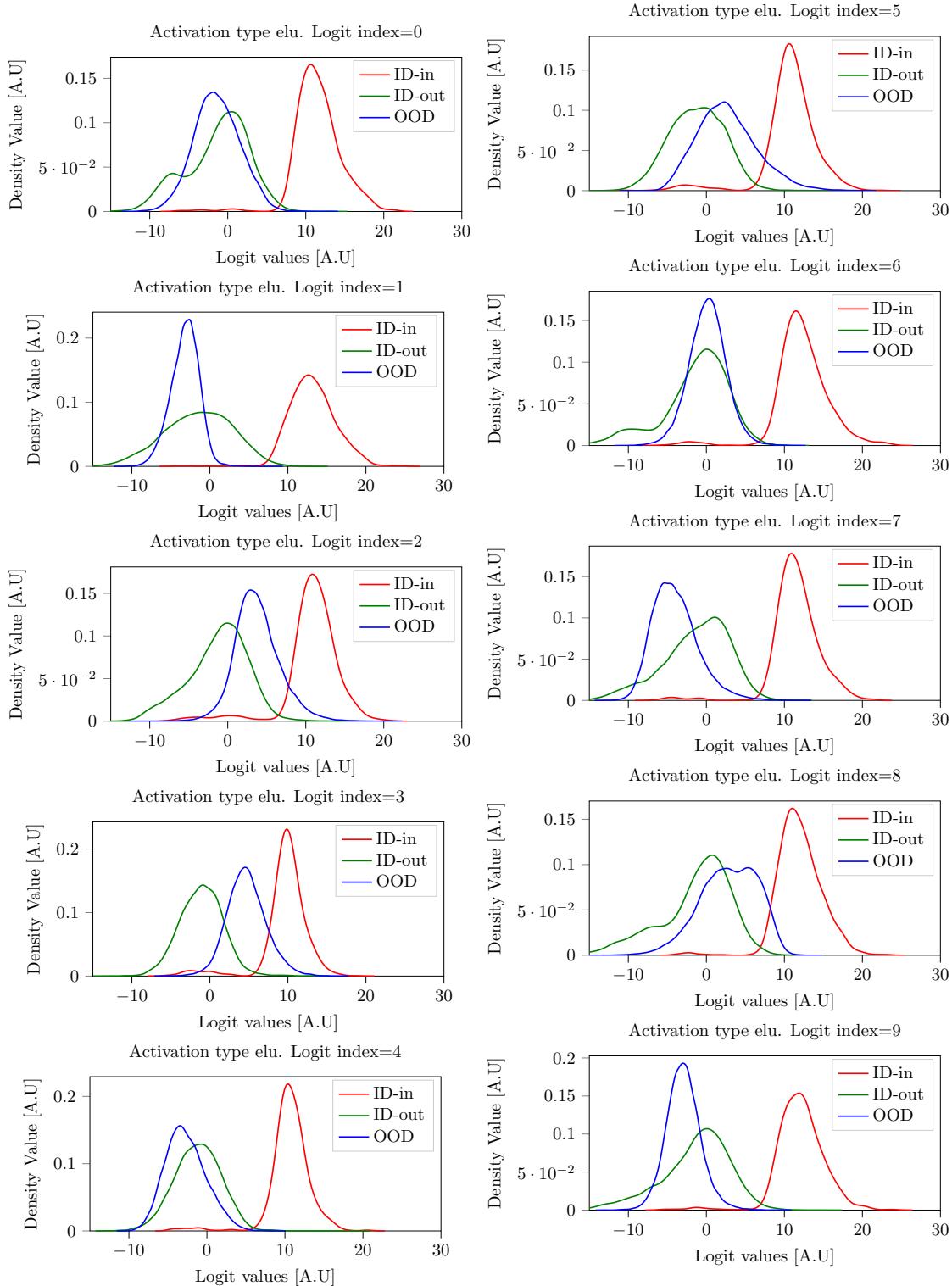


Figure 12: Densities over each logit cell from a Resnet-34 classifier with Elu activation.

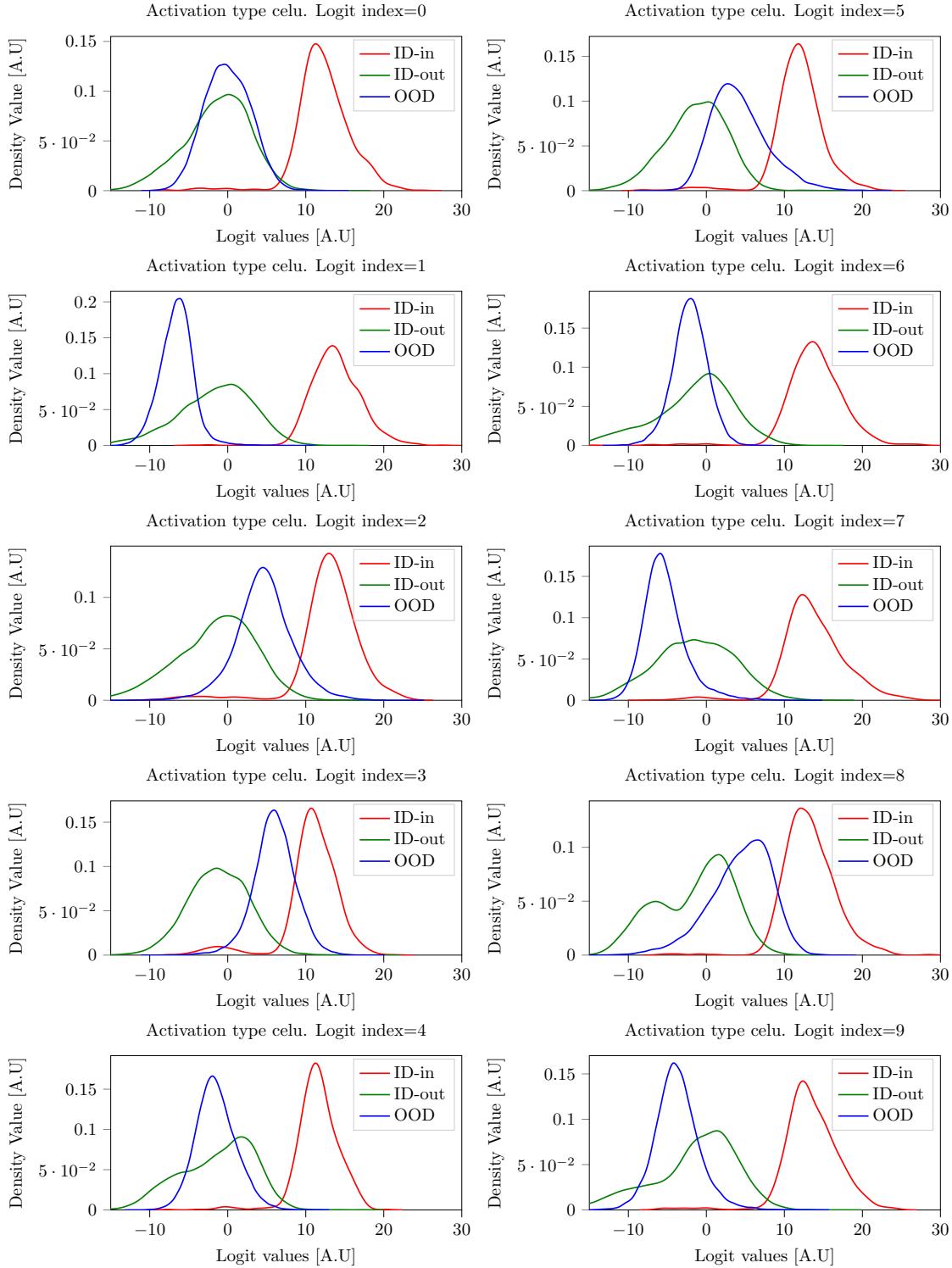


Figure 13: Densities over each logit cell from a Resnet-34 classifier with Celu activation.

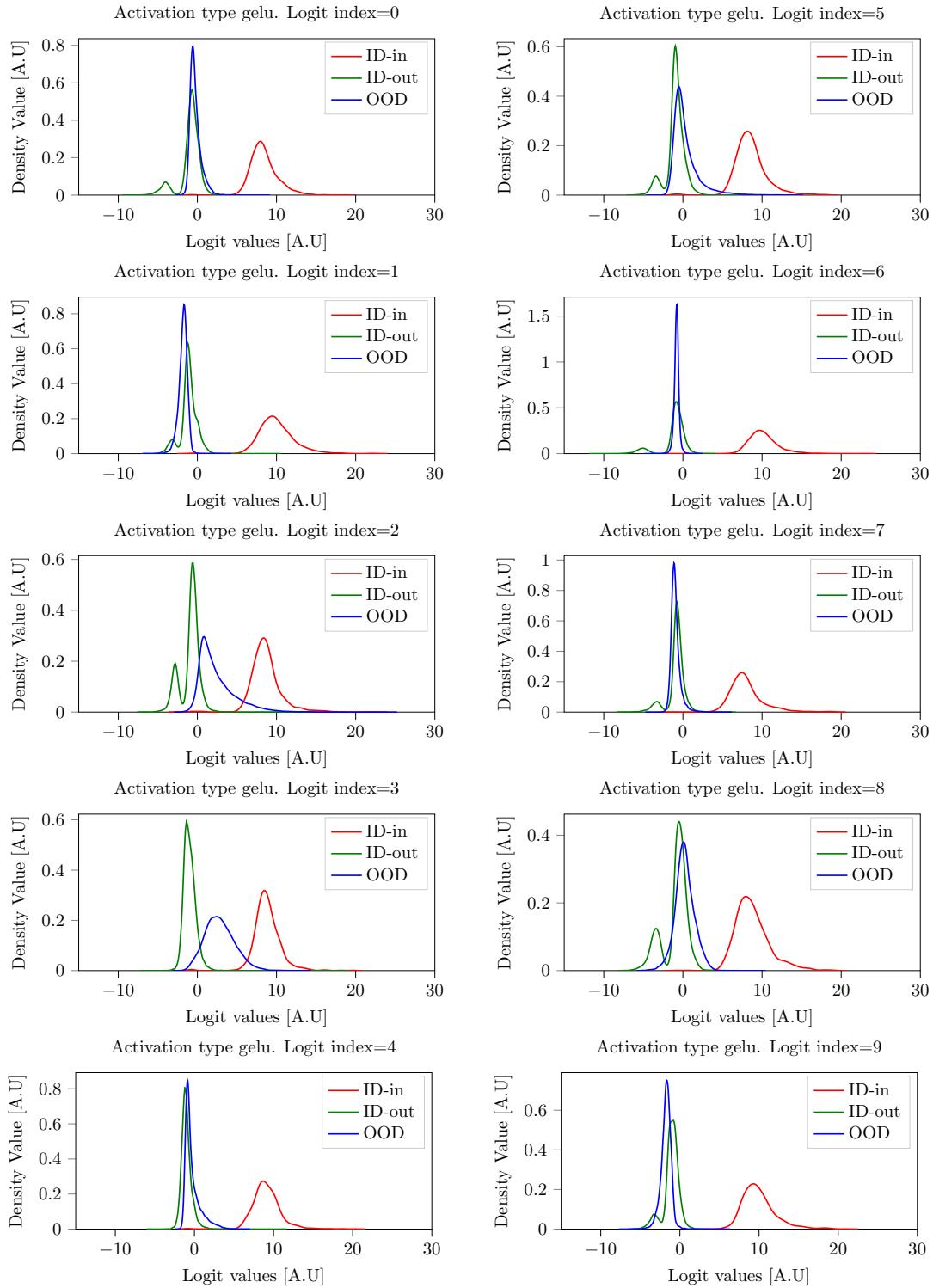


Figure 14: Densities over each logit cell from a Resnet-34 classifier with Gelu activation.

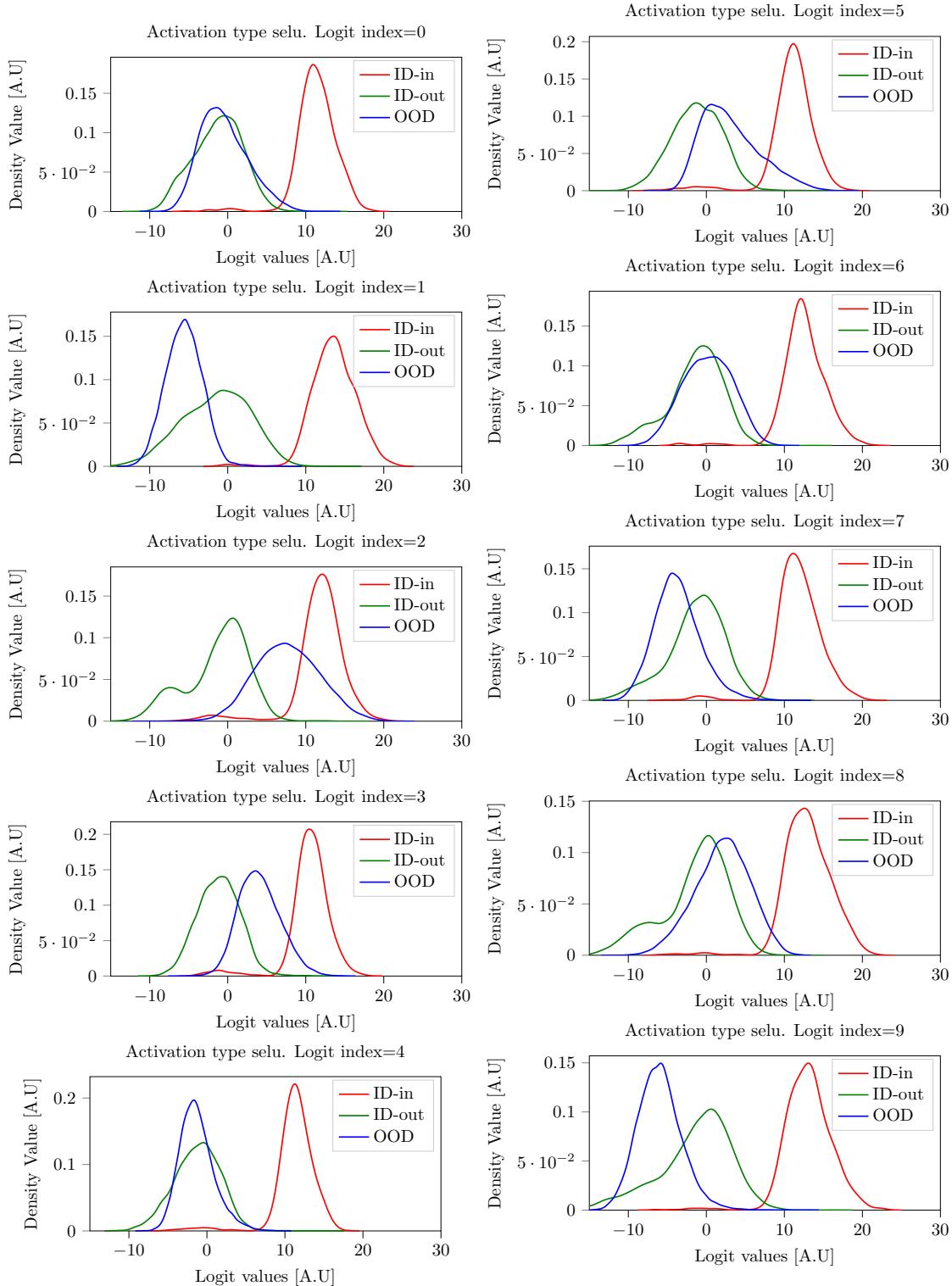


Figure 15: Densities over each logit cell from a Resnet-34 classifier with Selu activation.

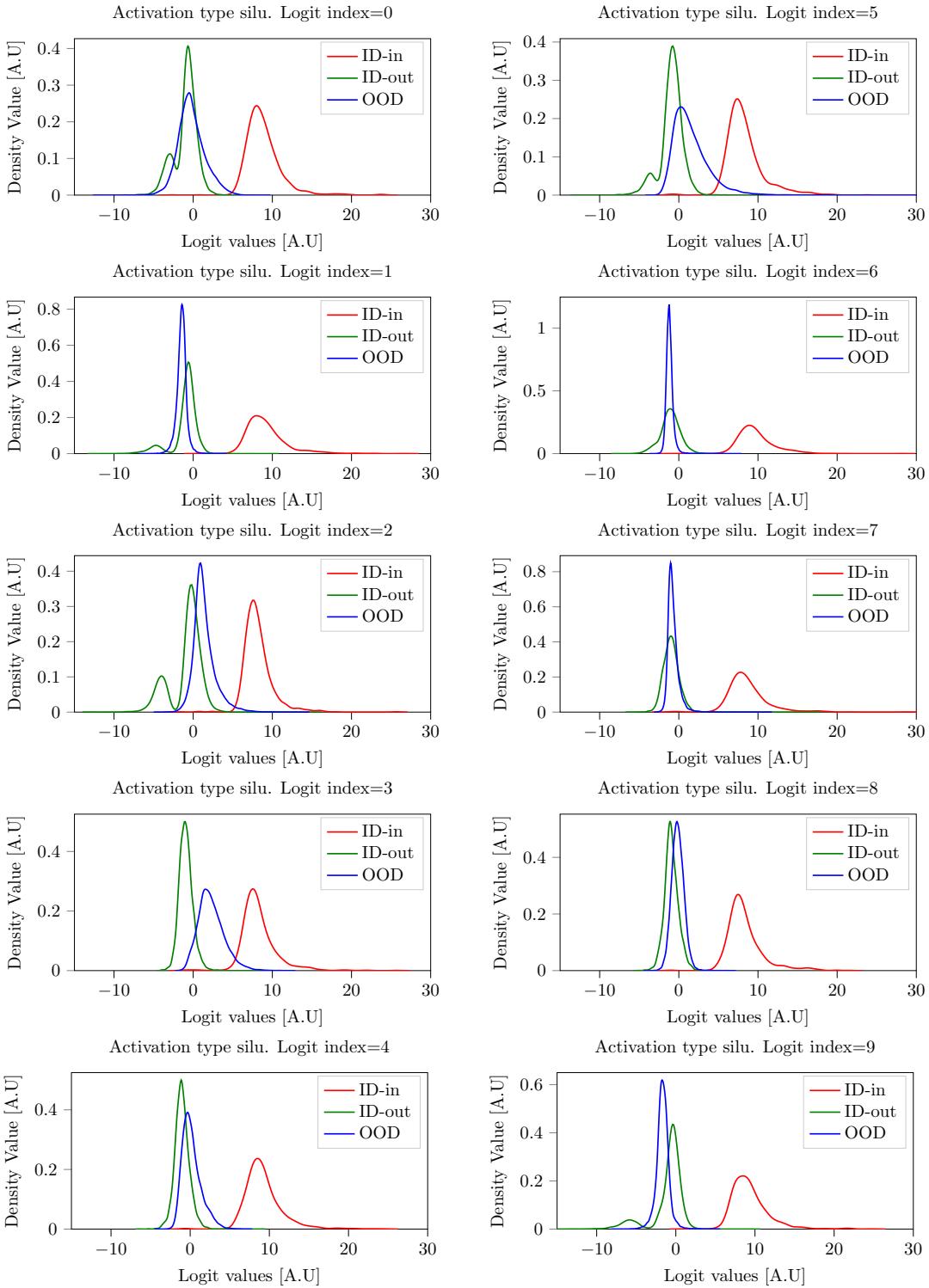


Figure 16: Densities over each logit cell from a Resnet-34 classifier with Silu activation.

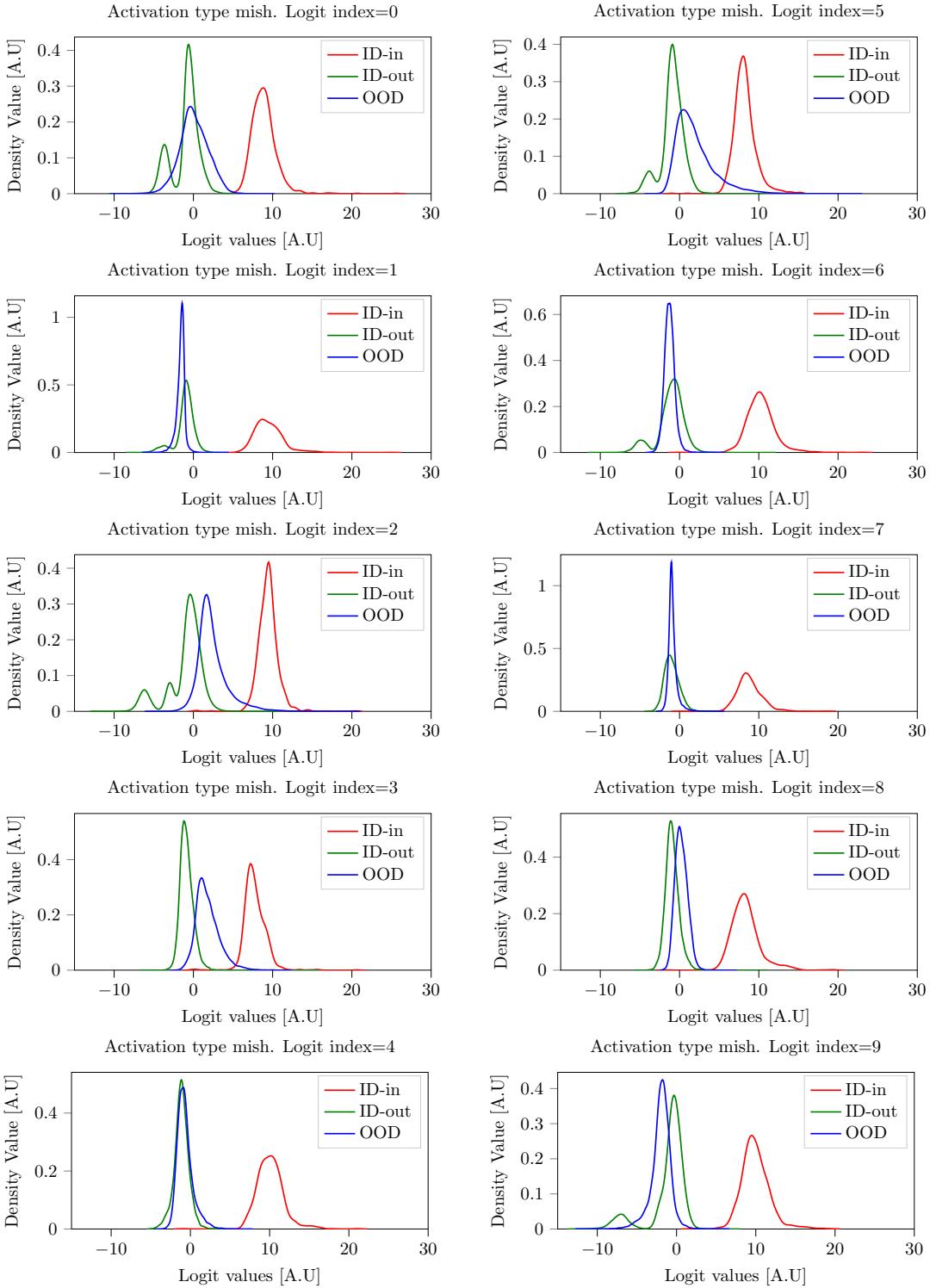


Figure 17: Densities over each logit cell from a Resnet-34 classifier with Mish activation.

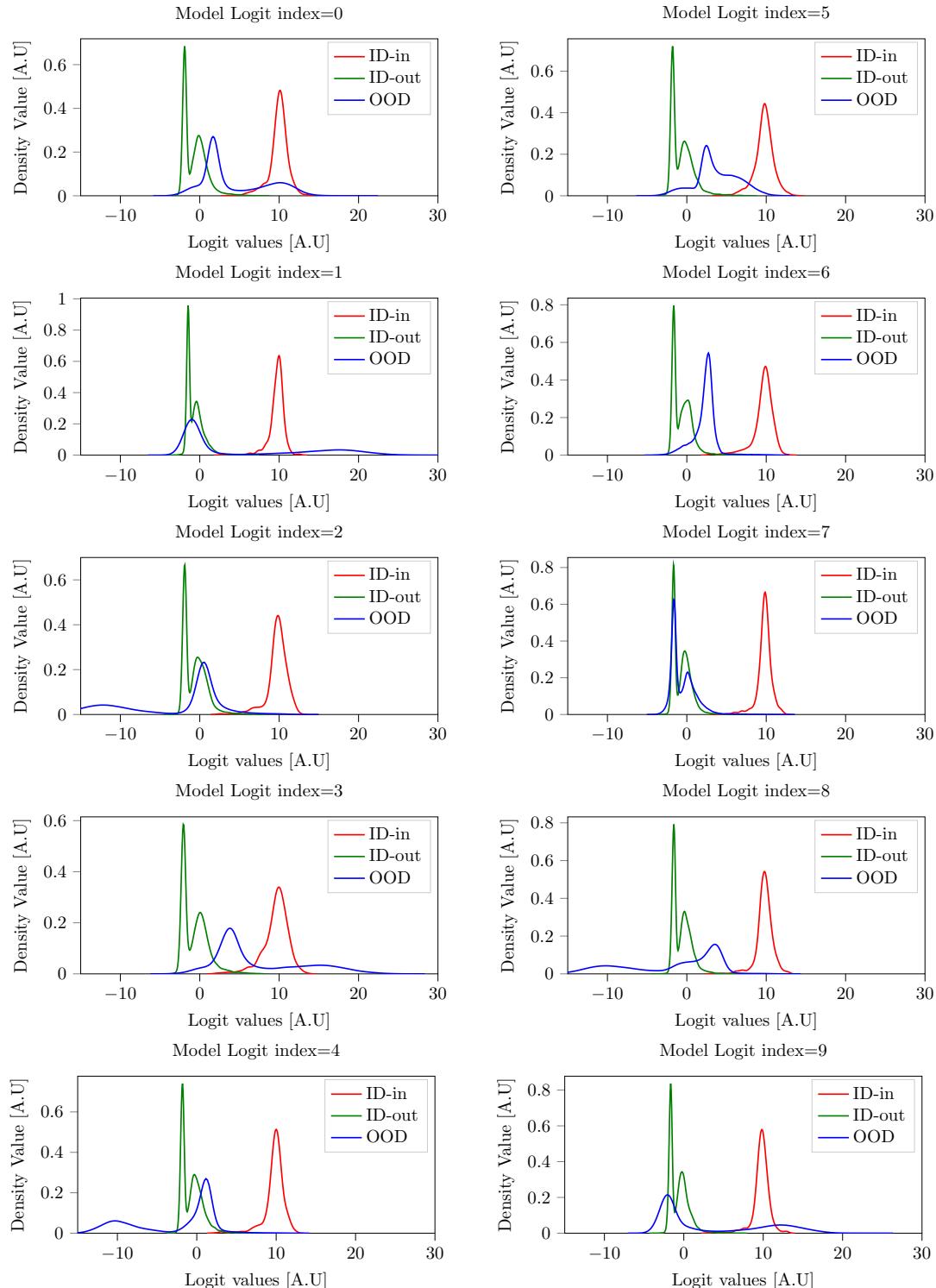


Figure 18: Densities over each logit cell from a Resnet-34 with a dropout 20% which is deactivated post train.

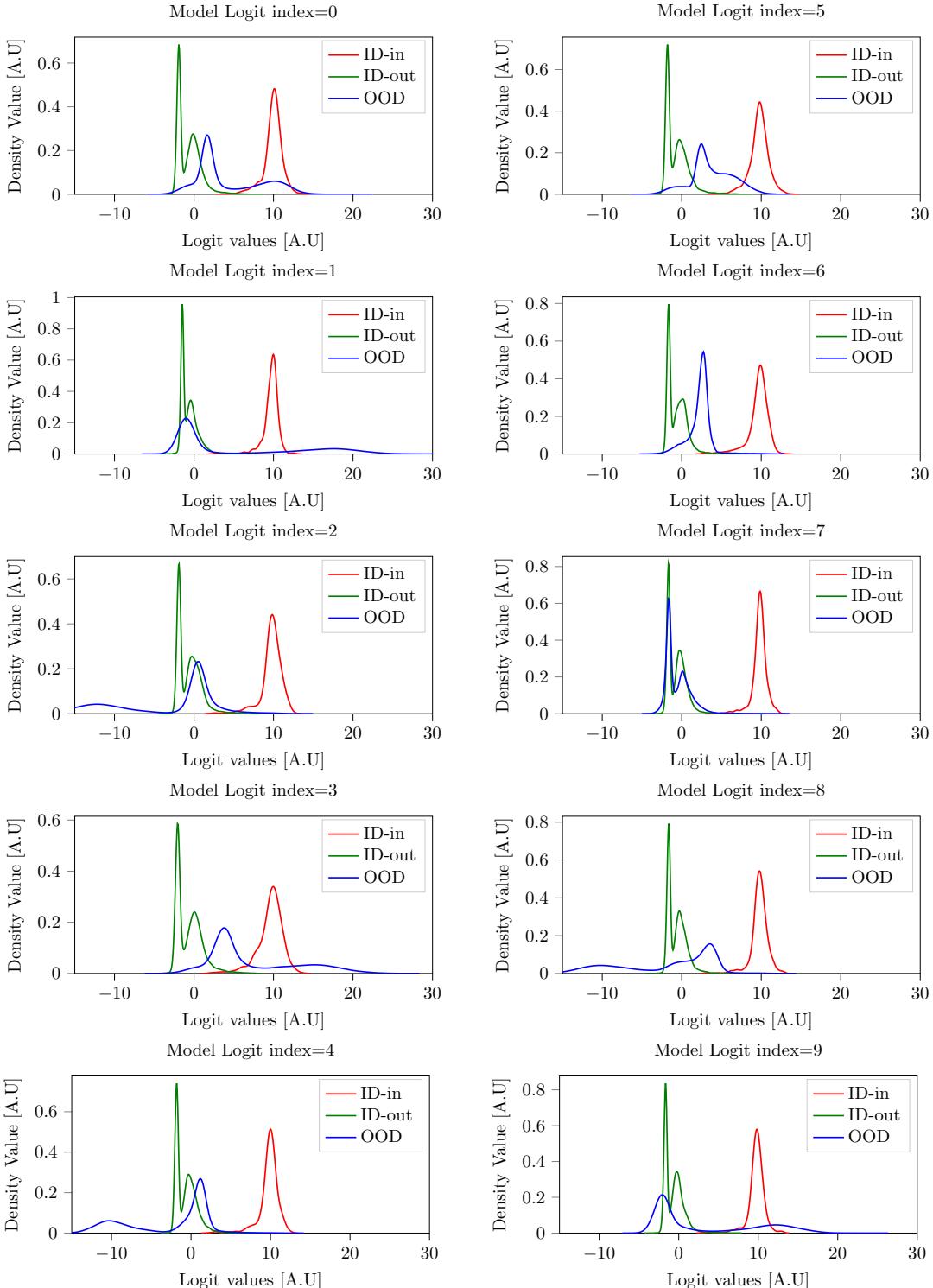


Figure 19: Densities over each logit cell from a Resnet-34 with a dropout 40% which is deactivated post train.

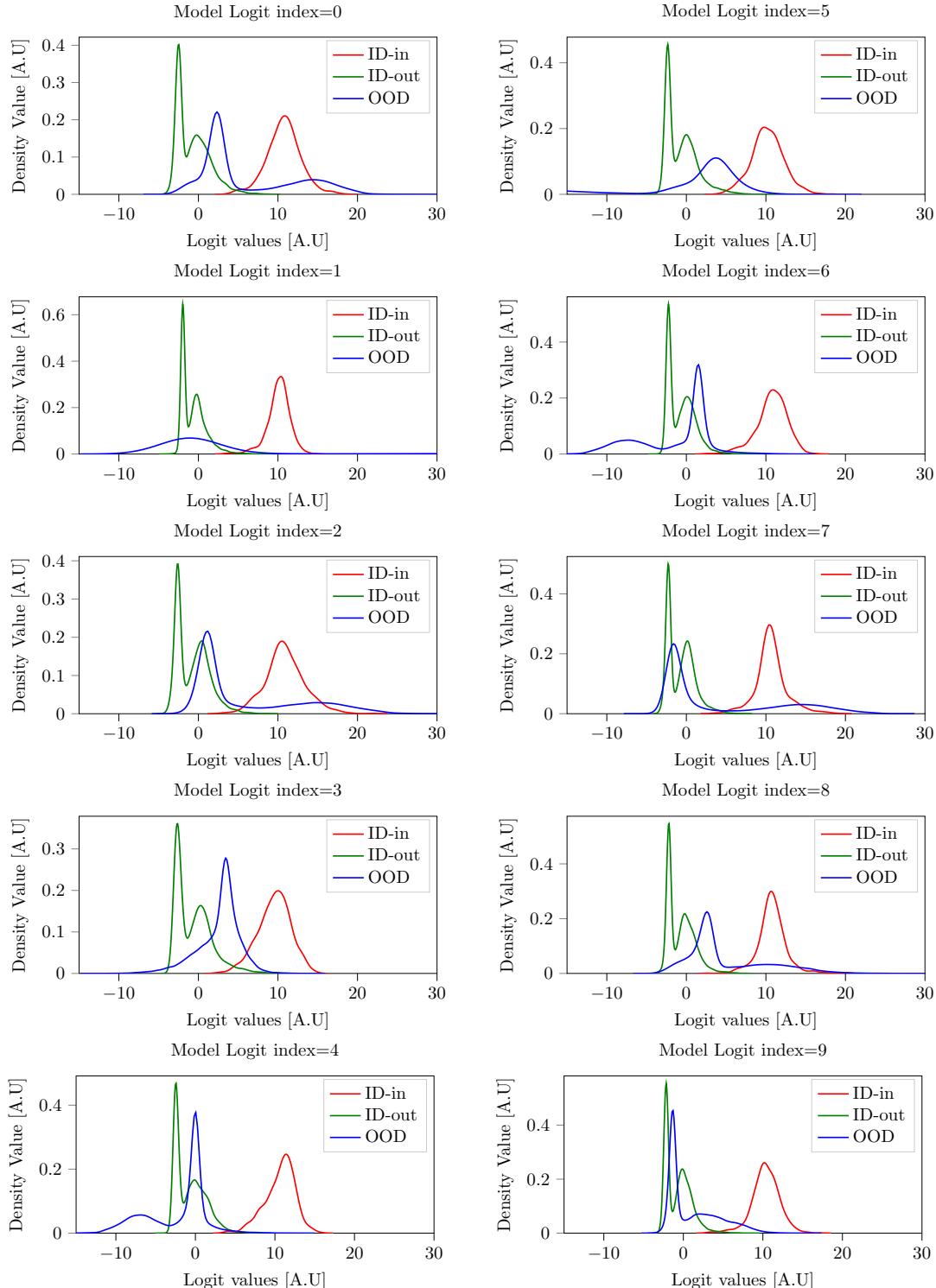


Figure 20: Densities over each logit cell from a Resnet-34 with a dropout 60% which is deactivated post train.

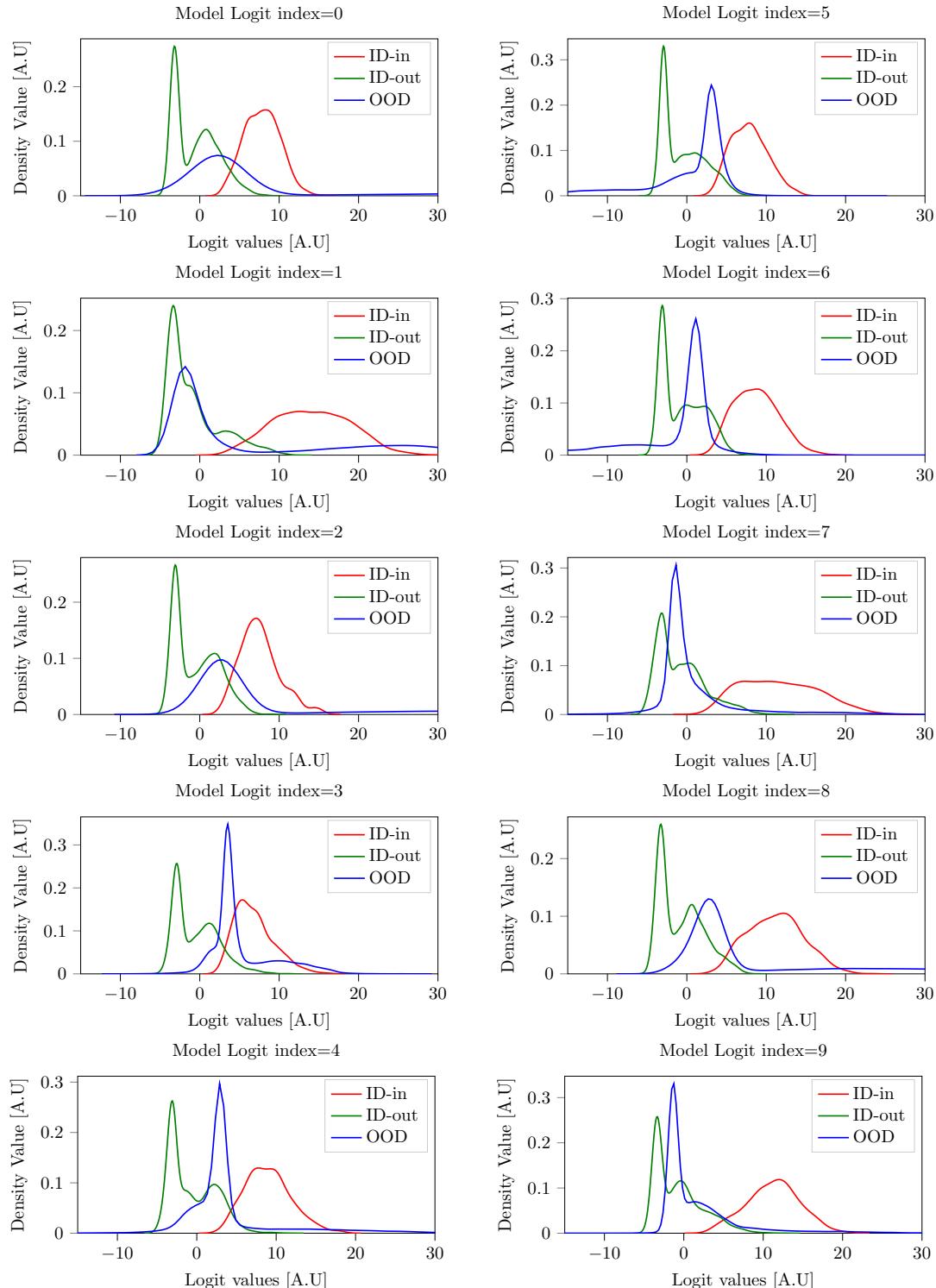


Figure 21: Densities over each logit cell from a Resnet-34 with a dropout 80% which is deactivated post train.

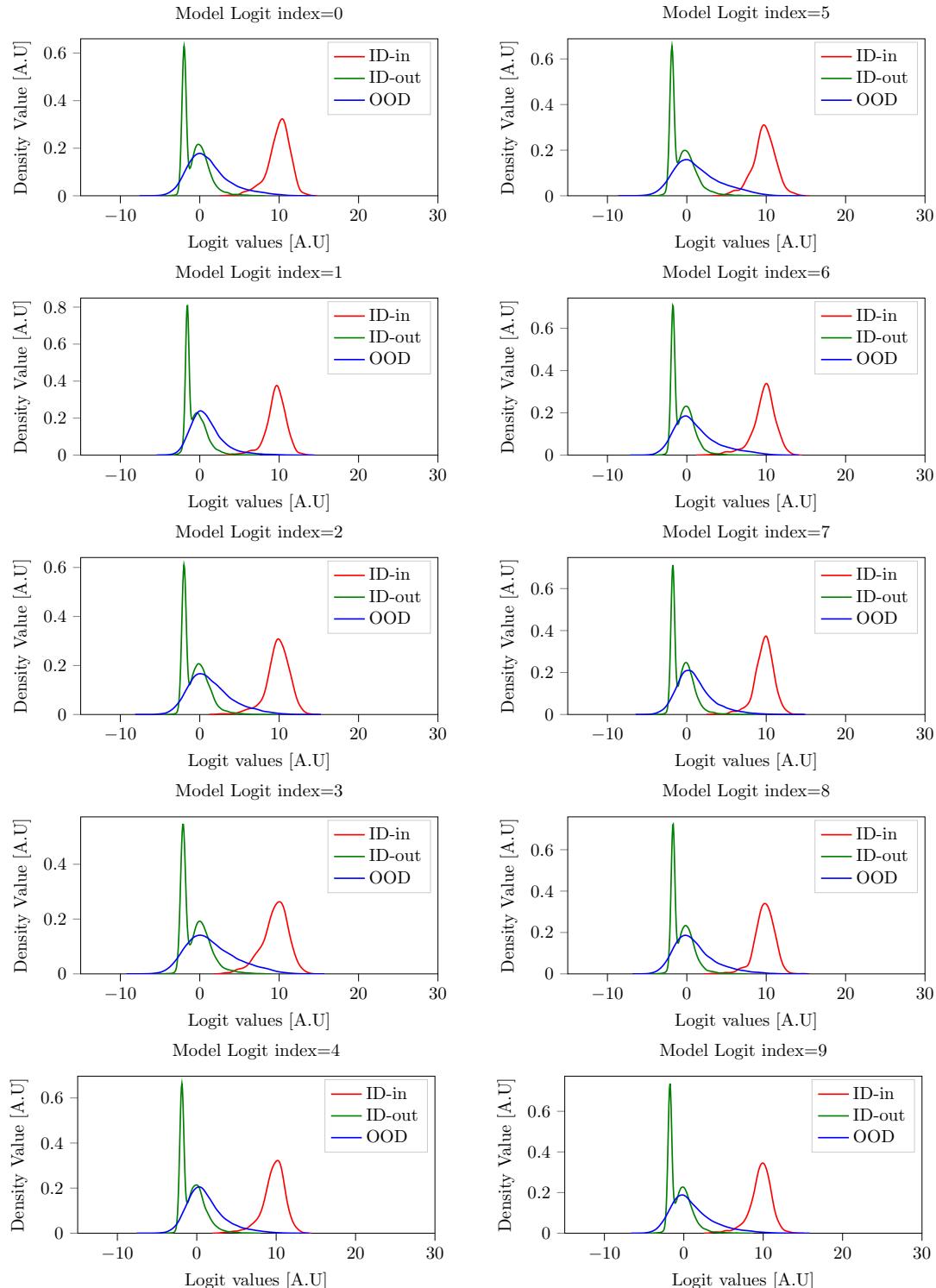


Figure 22: Densities over each logit cell from a Resnet-34 with a dropout 20%, which remains activated post train.

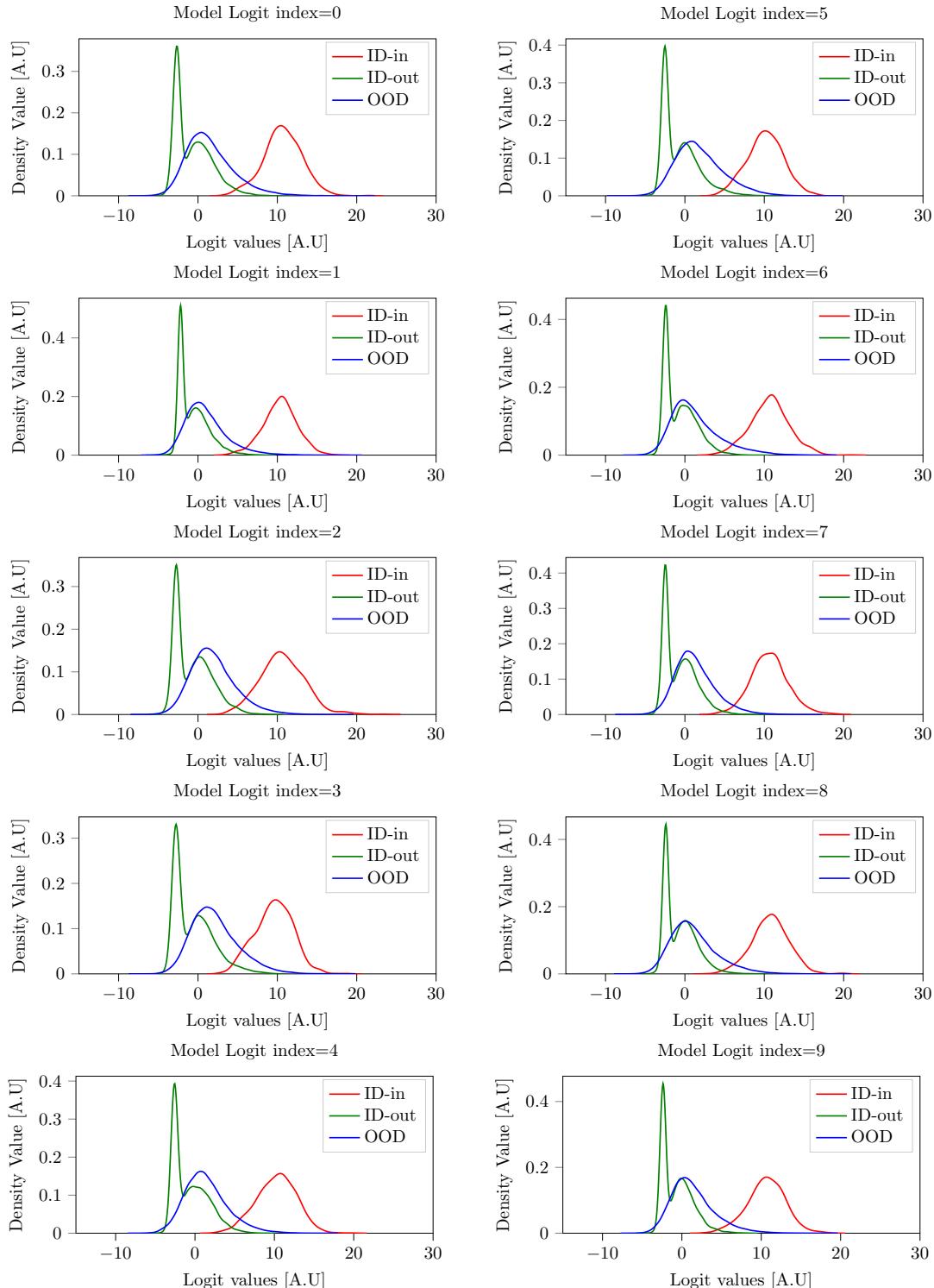


Figure 23: Densities over each logit cell from a Resnet-34 with a dropout 40%, which remains activated post train.

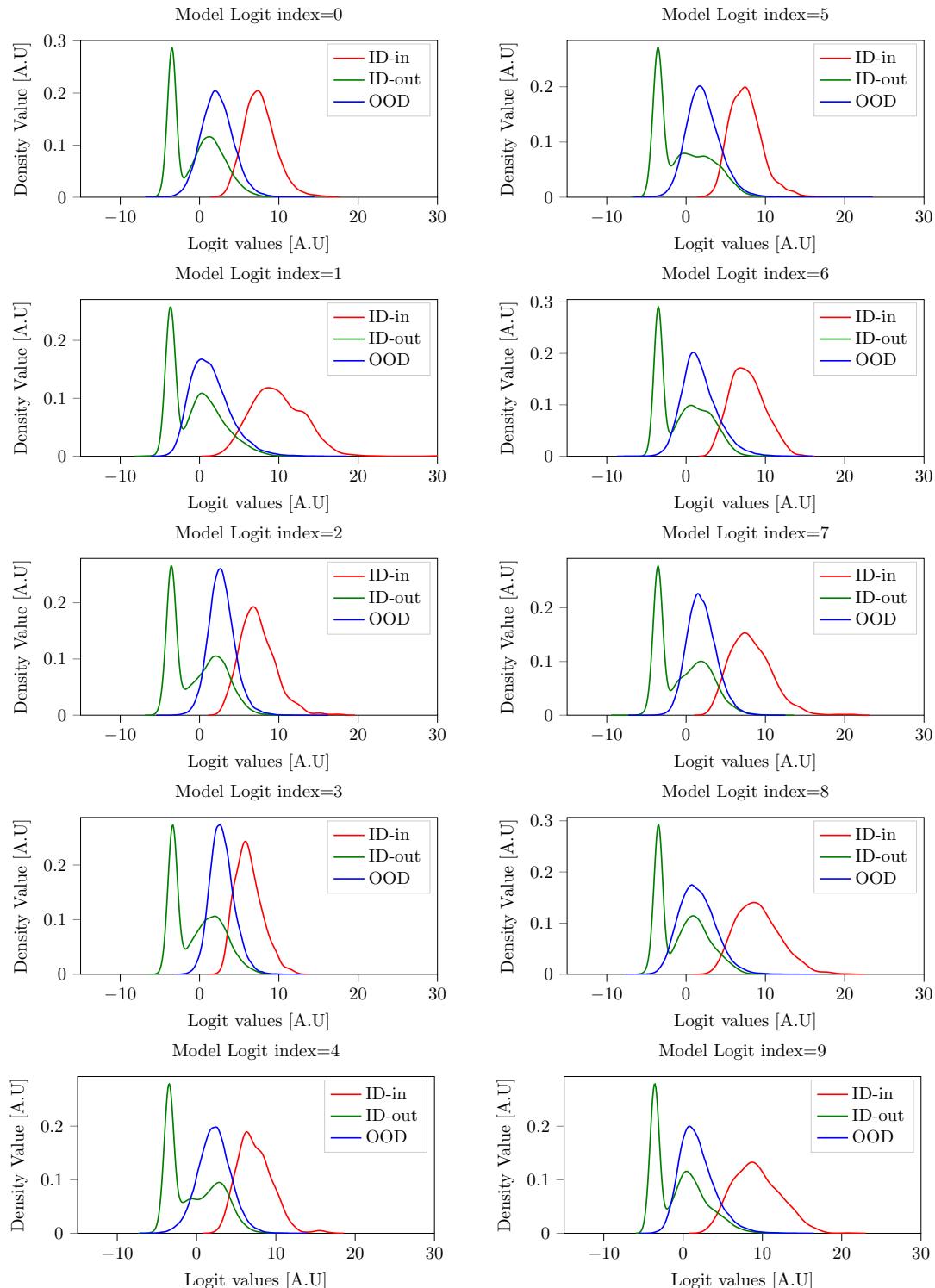


Figure 24: Densities over each logit cell from a Resnet-34 with a dropout 60%, which remains activated post train.

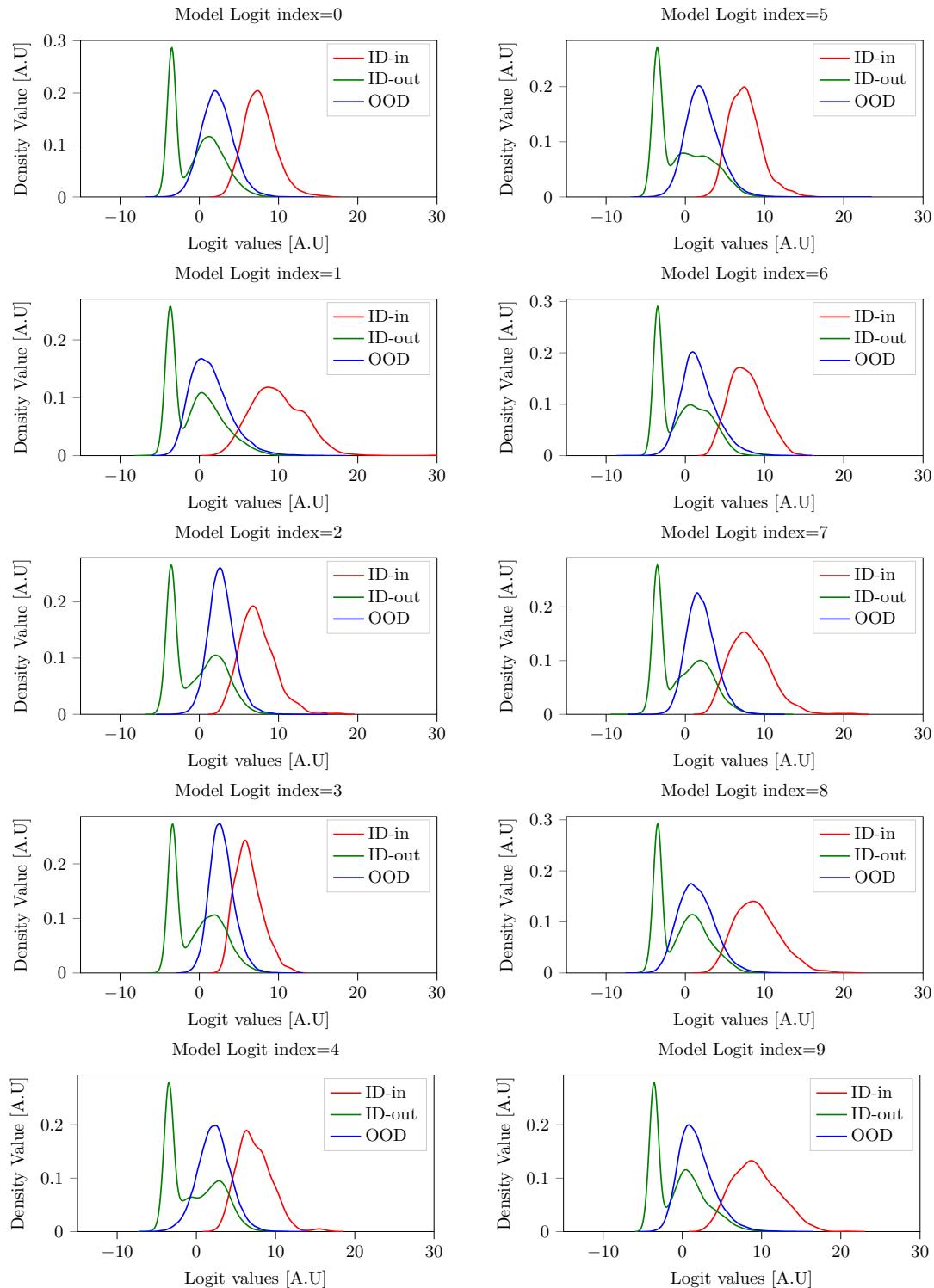


Figure 25: Densities over each logit cell from a Resnet-34 with a dropout 80%, which remains activated post train.

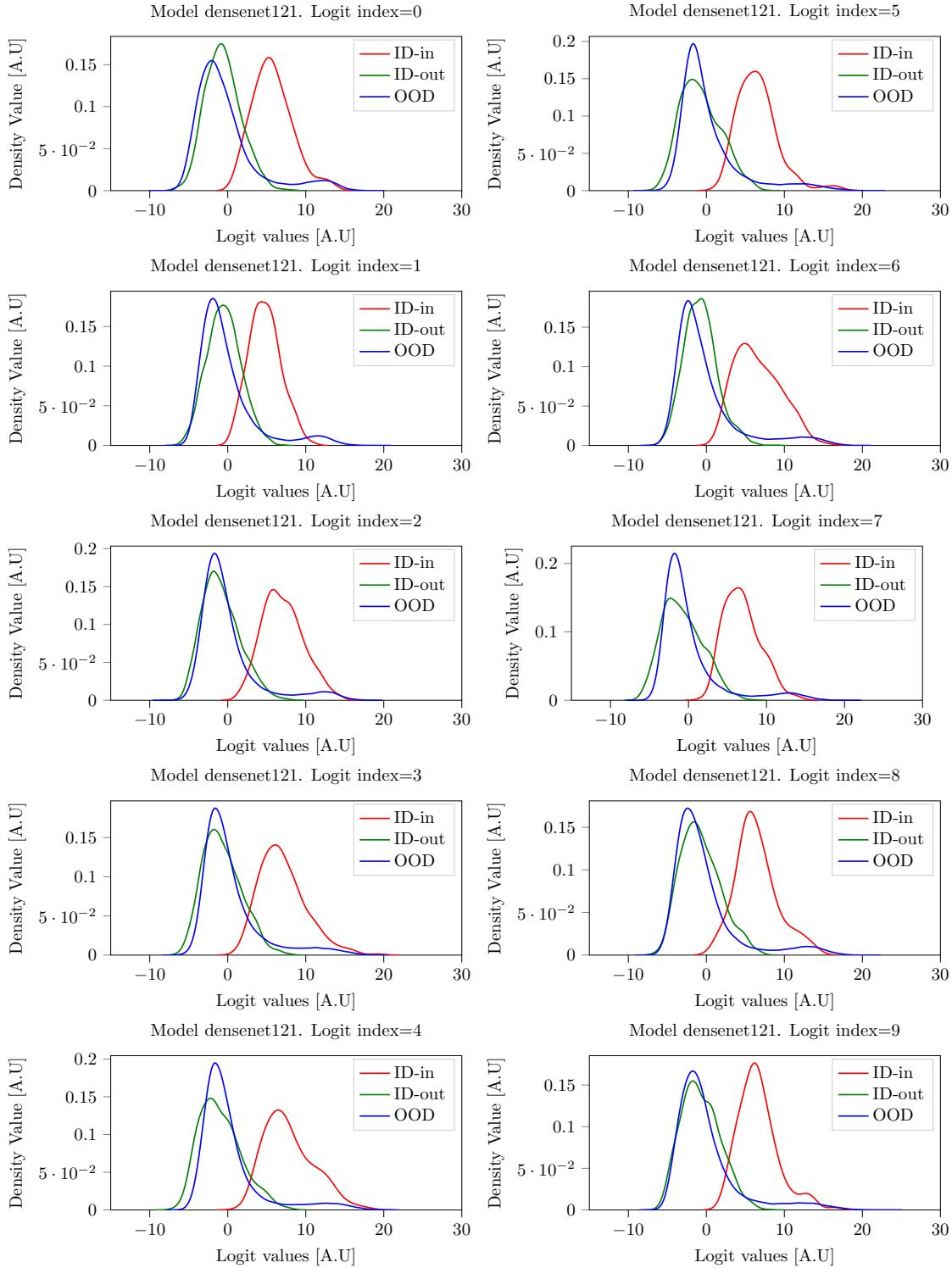


Figure 26: Logit cell densities for SVHN as ID with Densenet121.

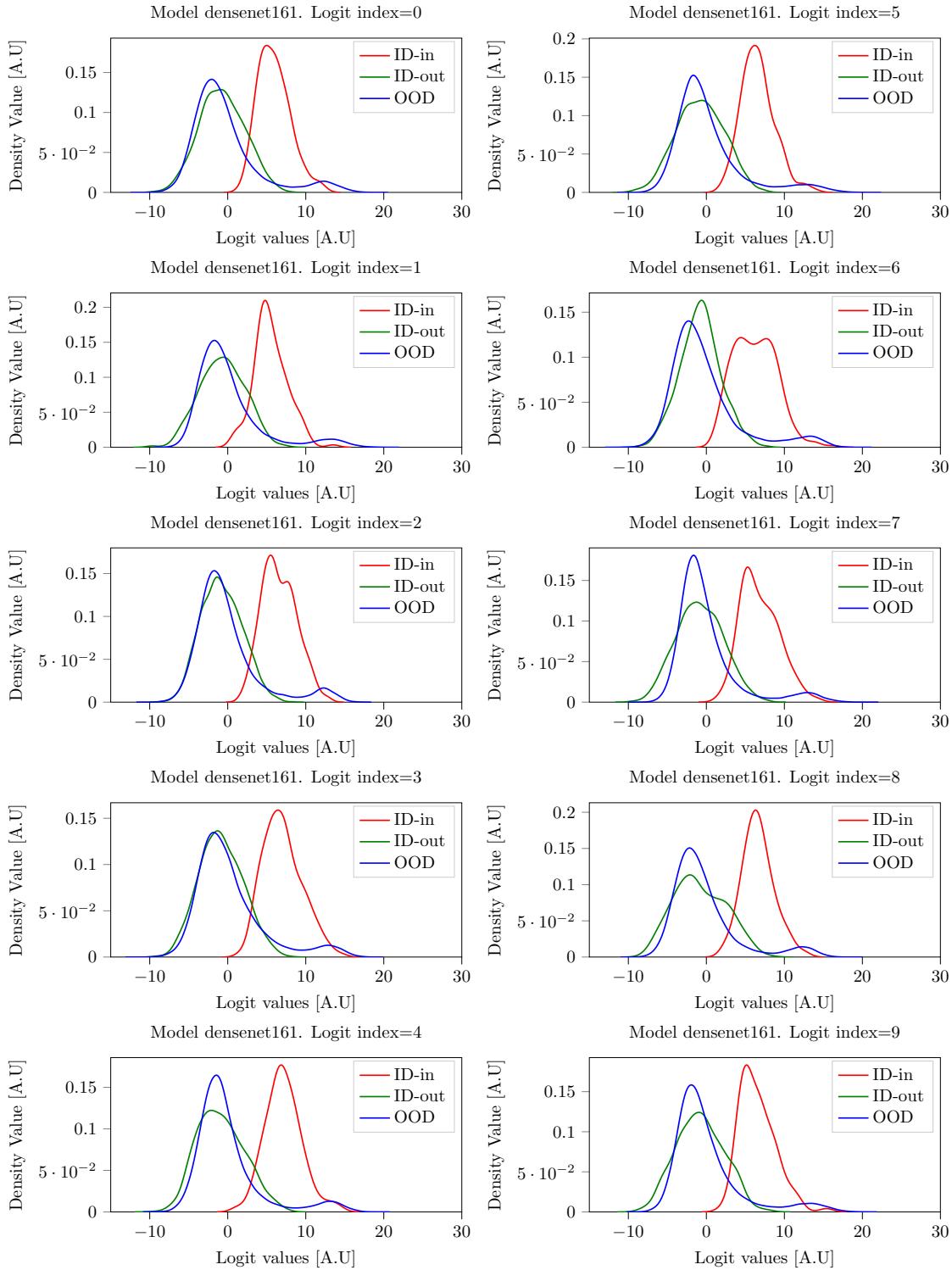


Figure 27: Logit cell densities for SVHN as ID with Densenet161.

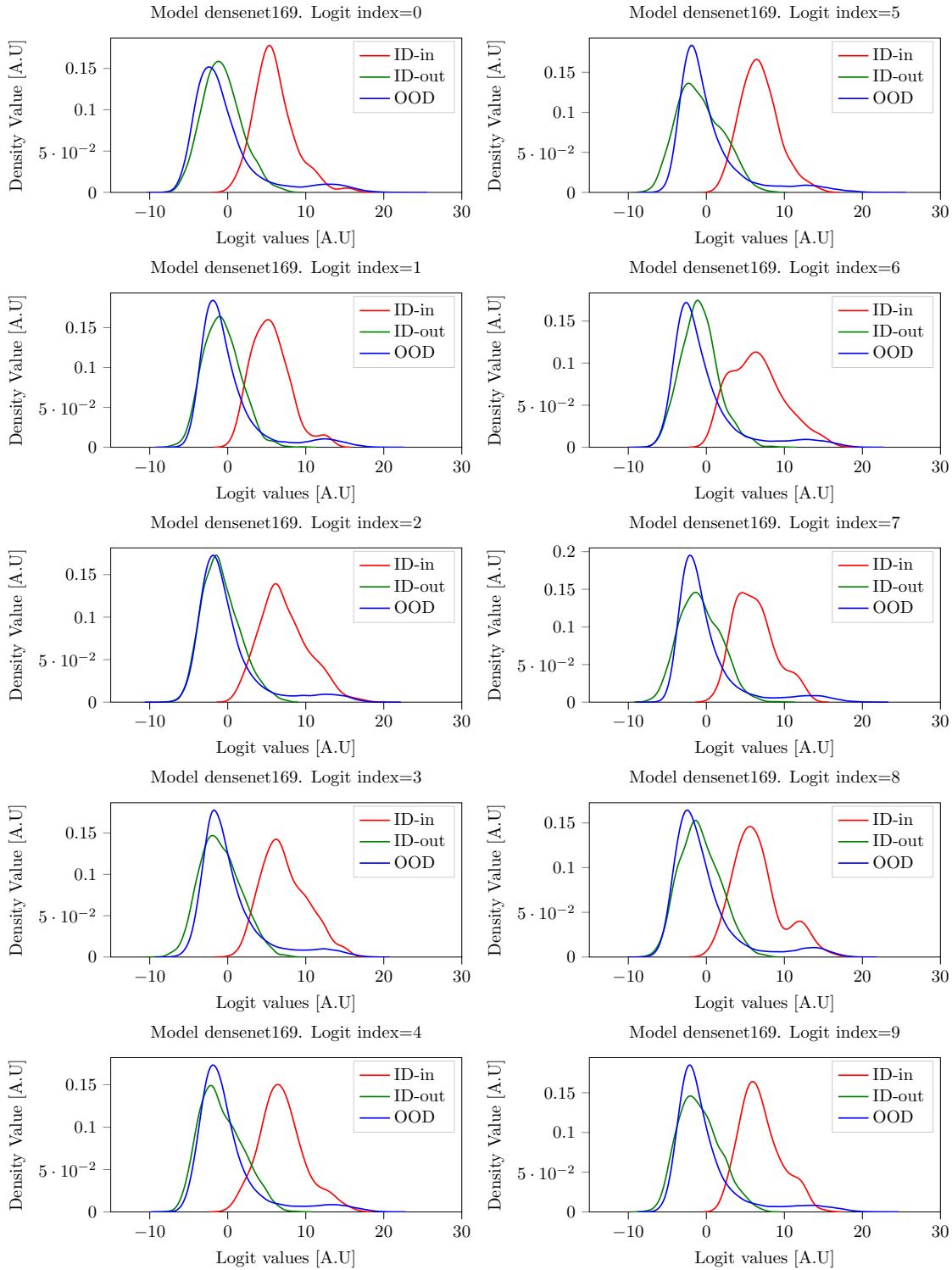


Figure 28: Logit cell densities for SVHN as ID with Densenet169.

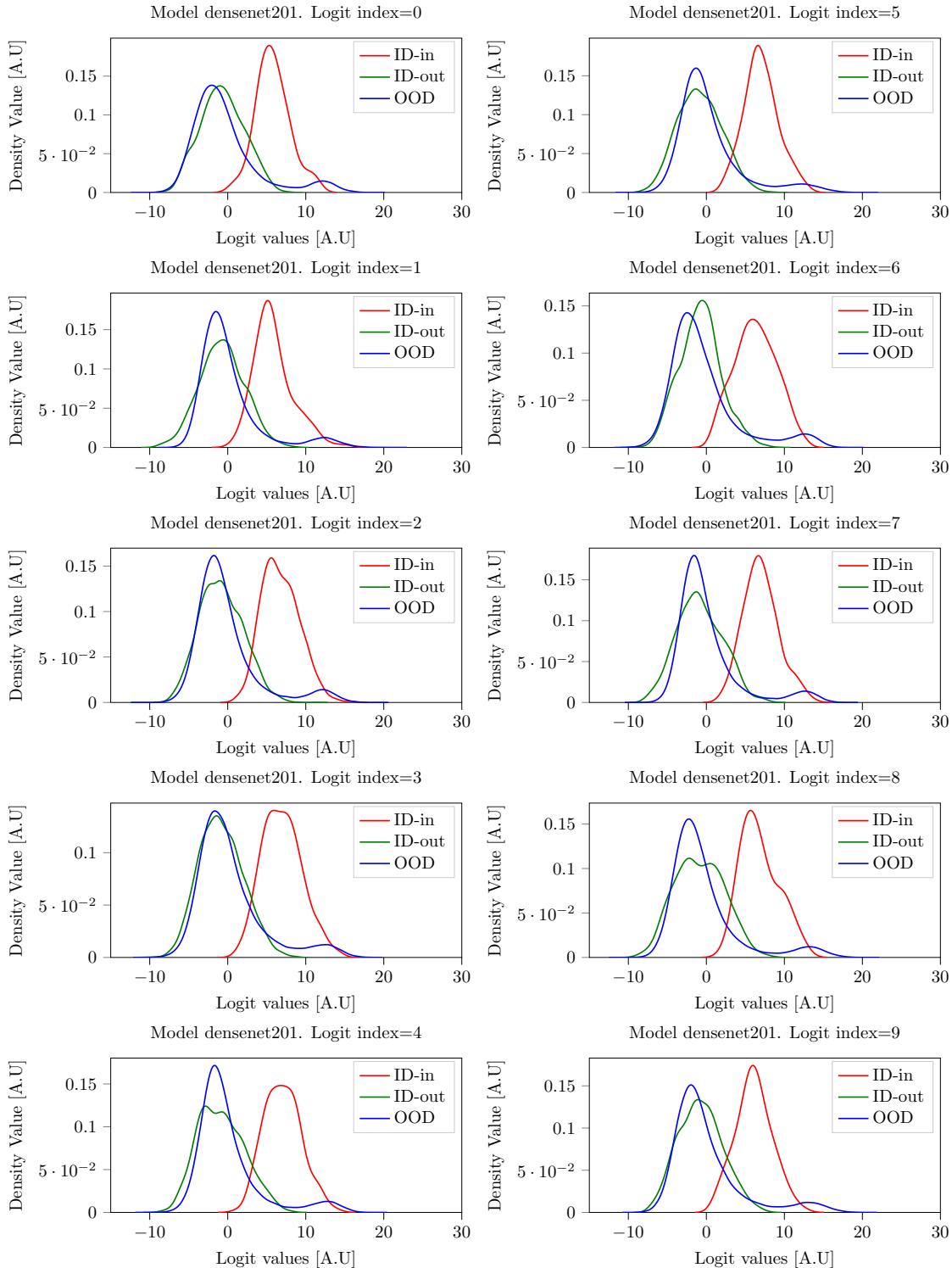


Figure 29: Logit cell densities for SVHN as ID with Densenet201.

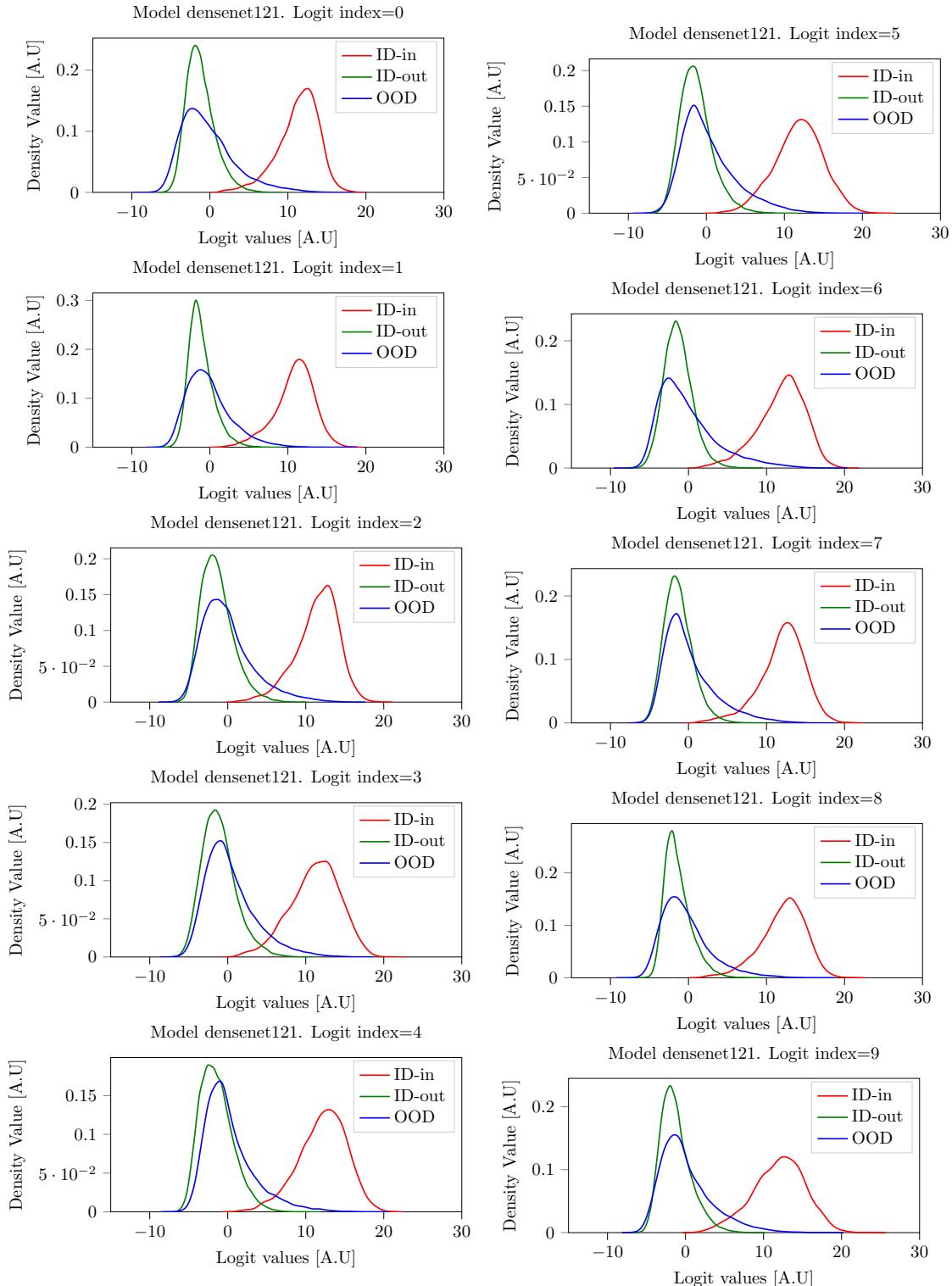


Figure 30: Logit cell densities for CIFAR-10 as ID with Densenet121.

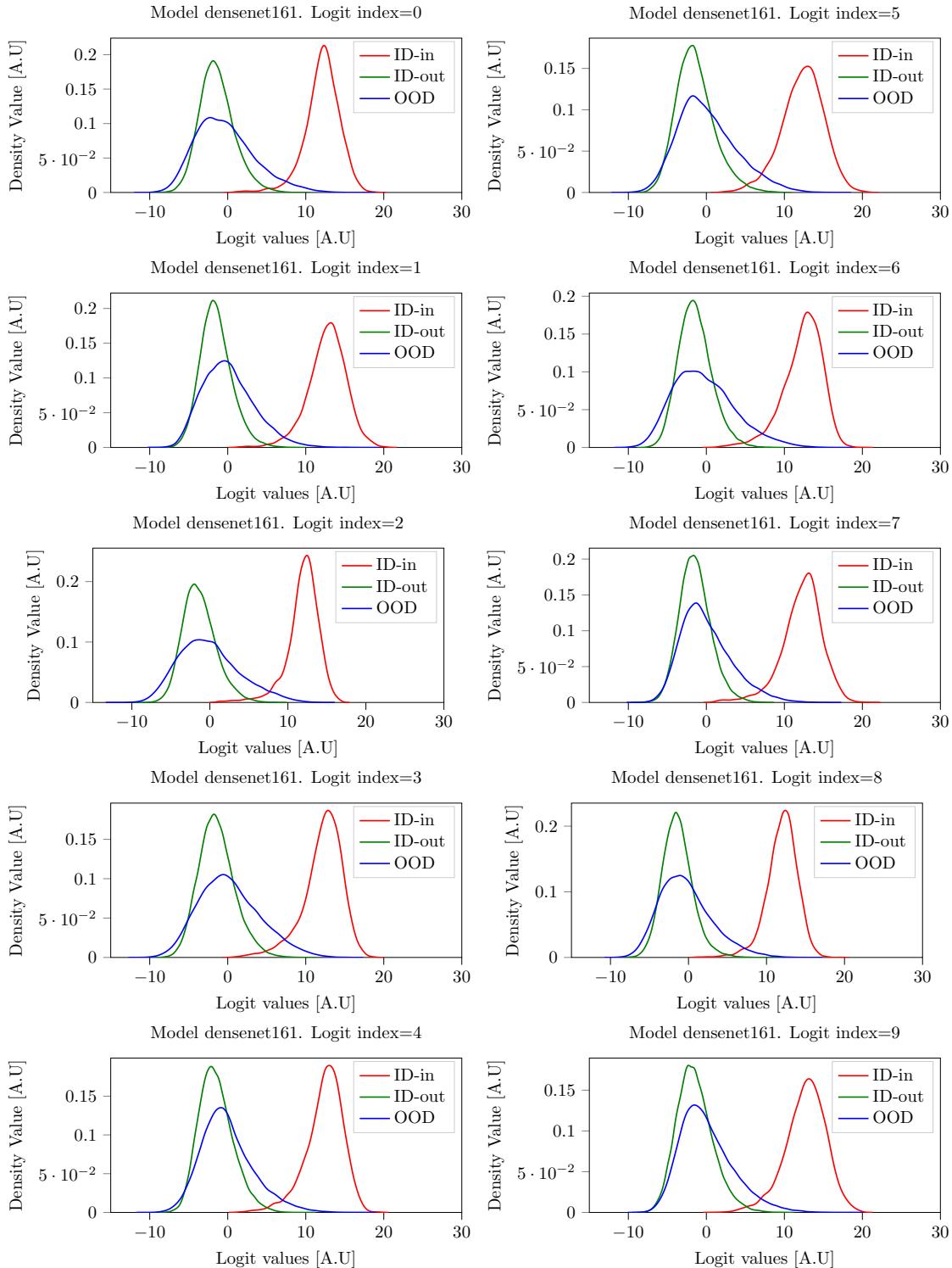


Figure 31: Logit cell densities for CIFAR-10 as ID with Densenet161.

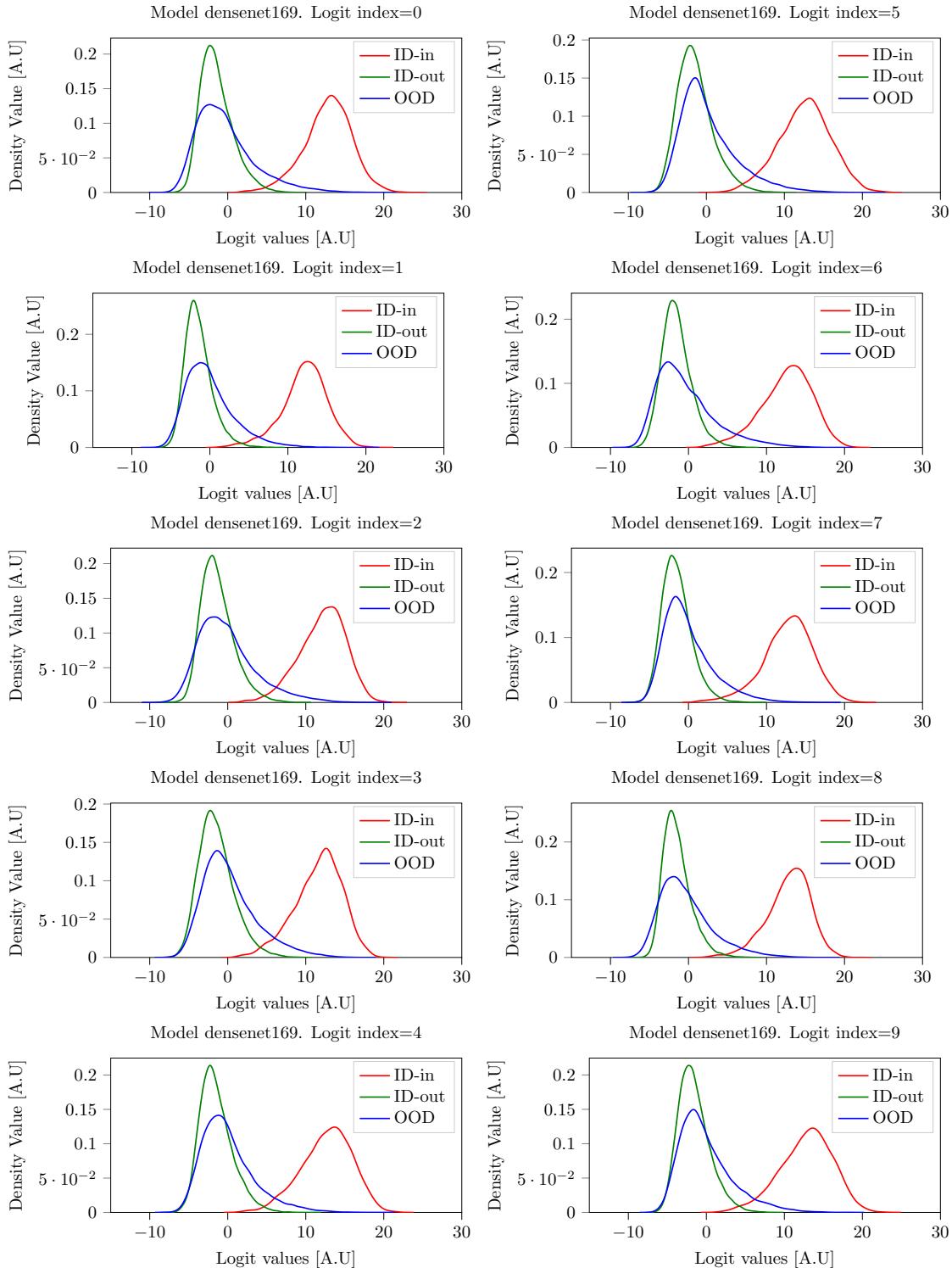


Figure 32: Logit cell densities for CIFAR-10 as ID with Densenet169.

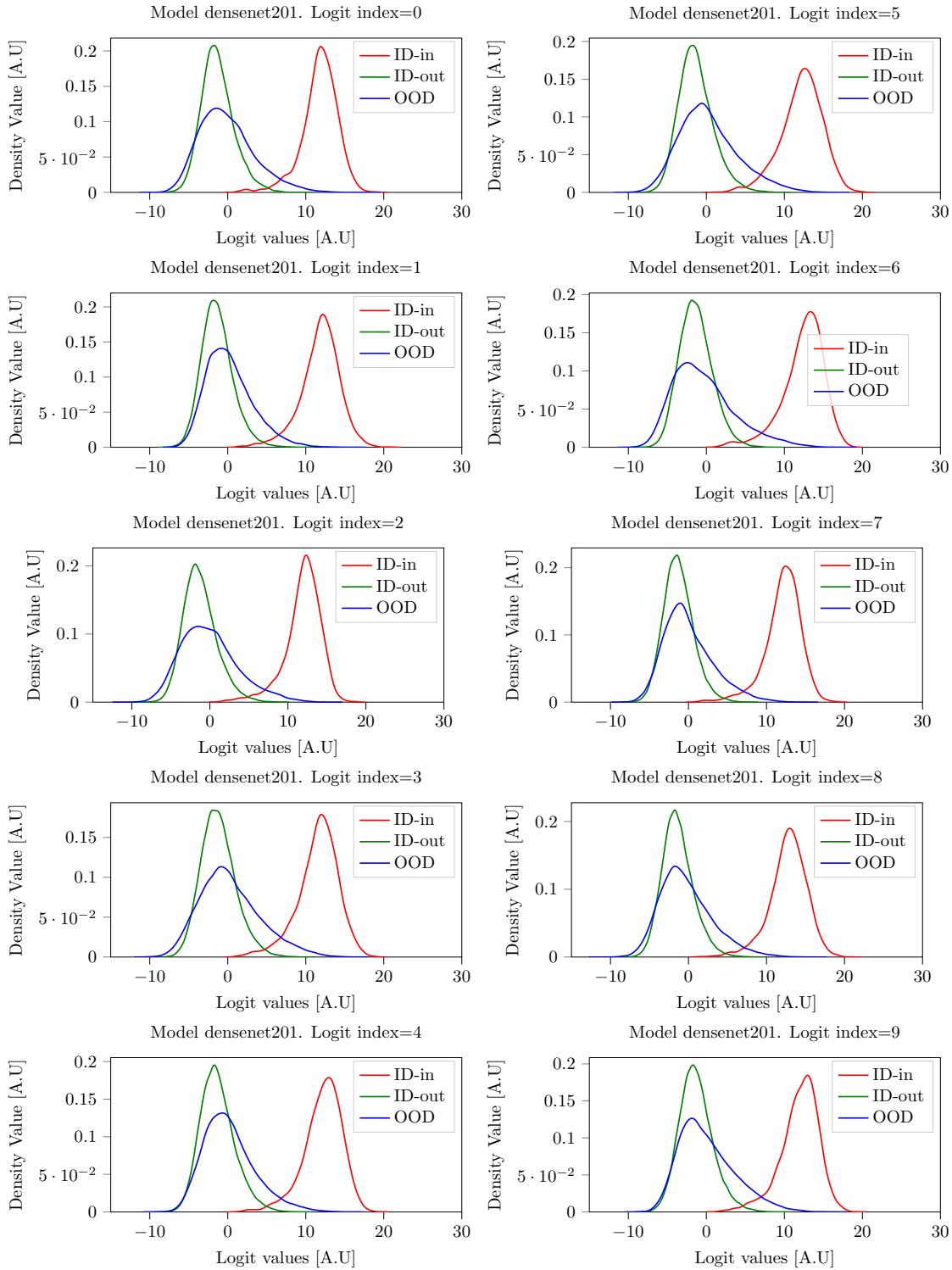


Figure 33: Logit cell densities for CIFAR-10 as ID with Densenet201.

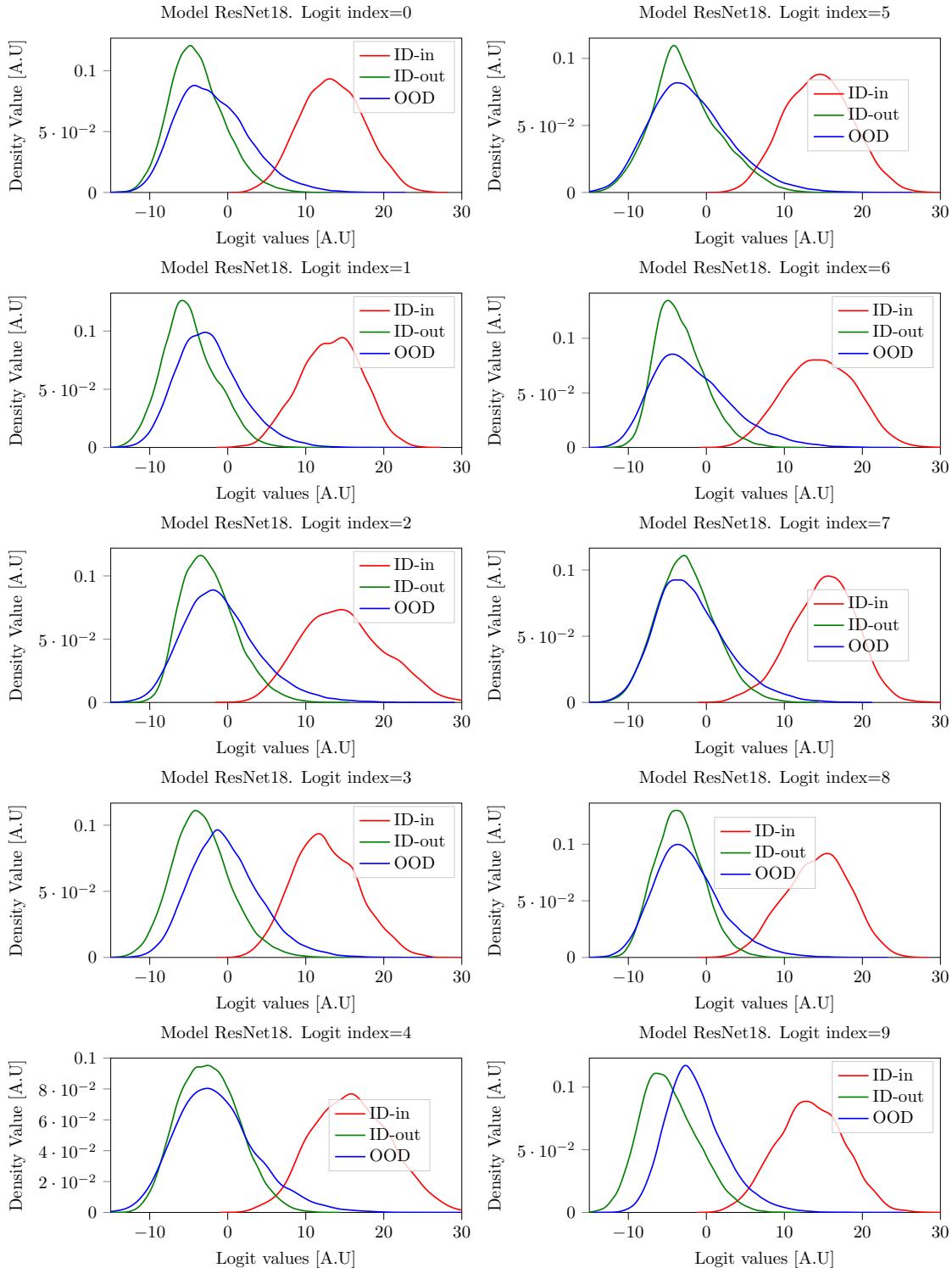


Figure 34: Logit cell densities for CIFAR-10 as ID with ResNet18.

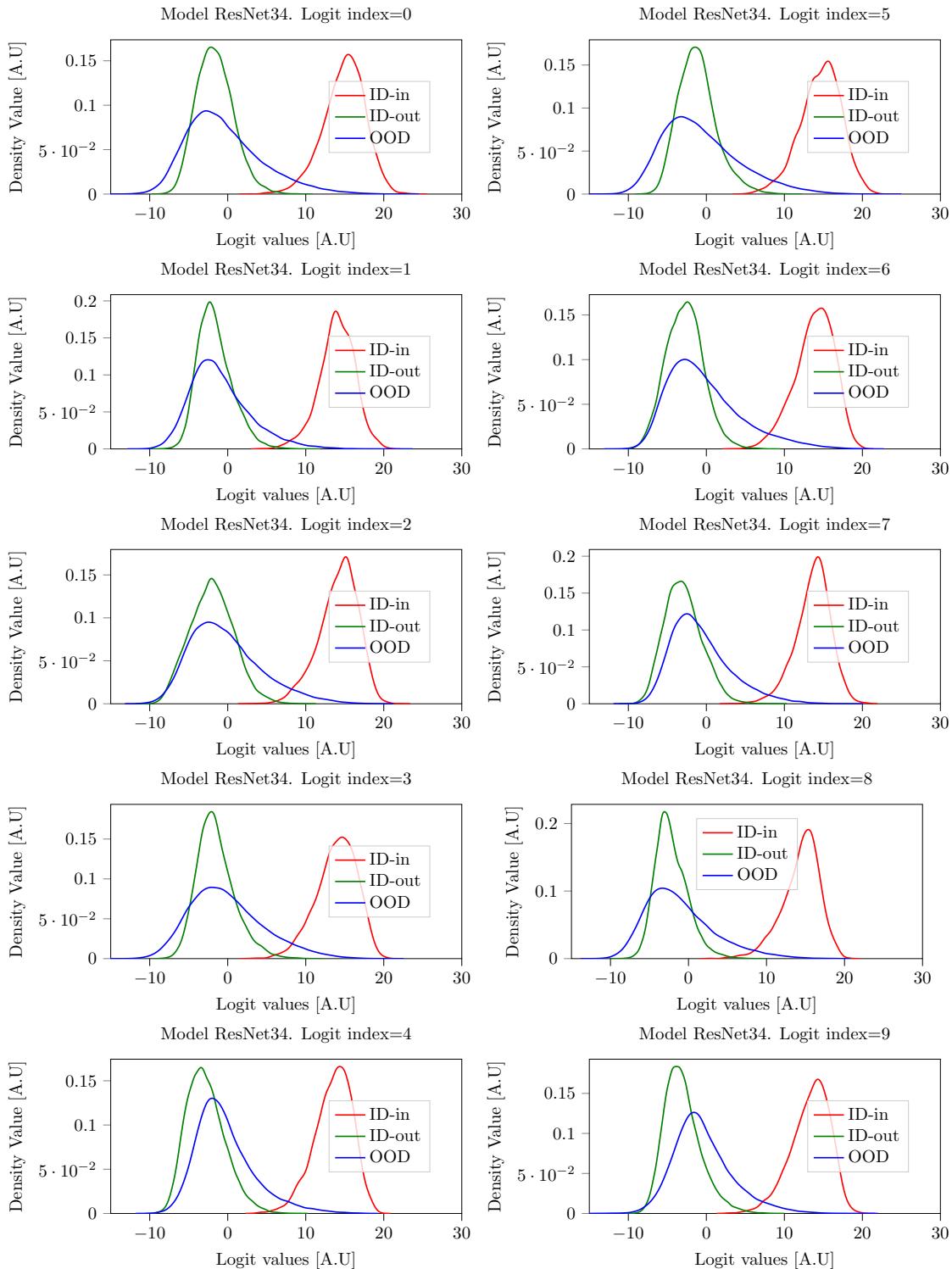


Figure 35: Logit cell densities for CIFAR-10 as ID with ResNet34.

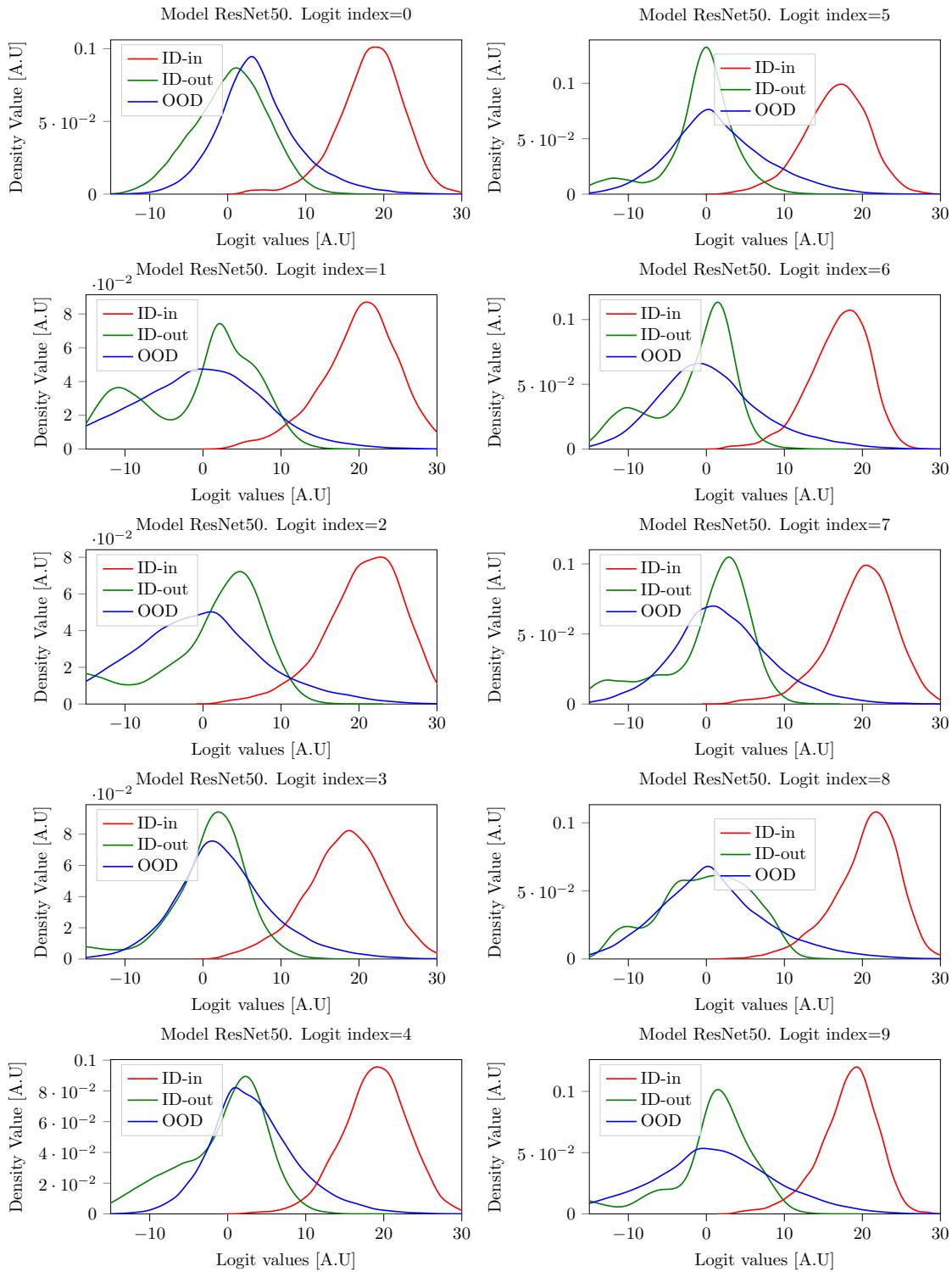


Figure 36: Logit cell densities for CIFAR-10 as ID with ResNet50.

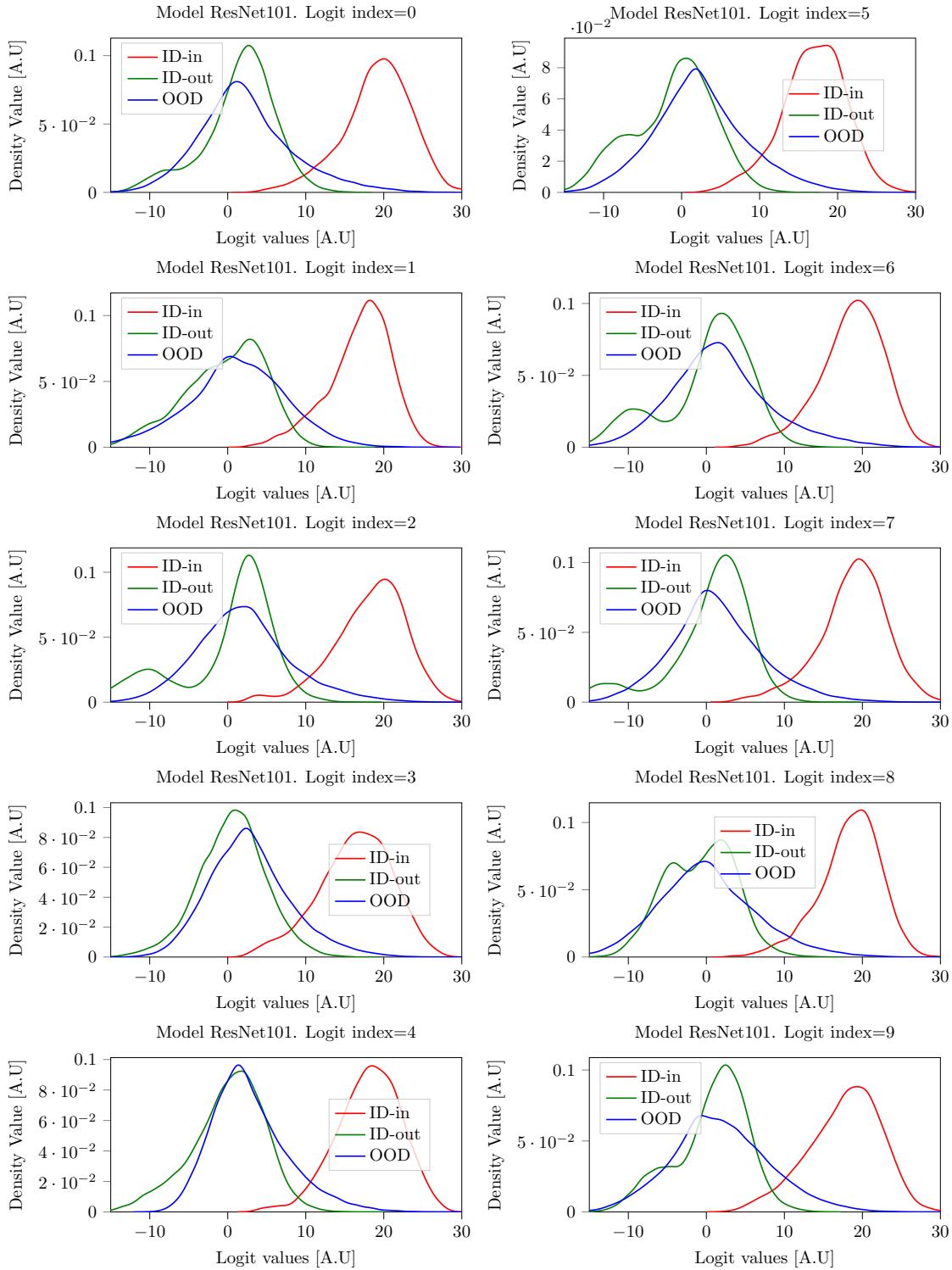


Figure 37: Logit cell densities for CIFAR-10 as ID with ResNet101.

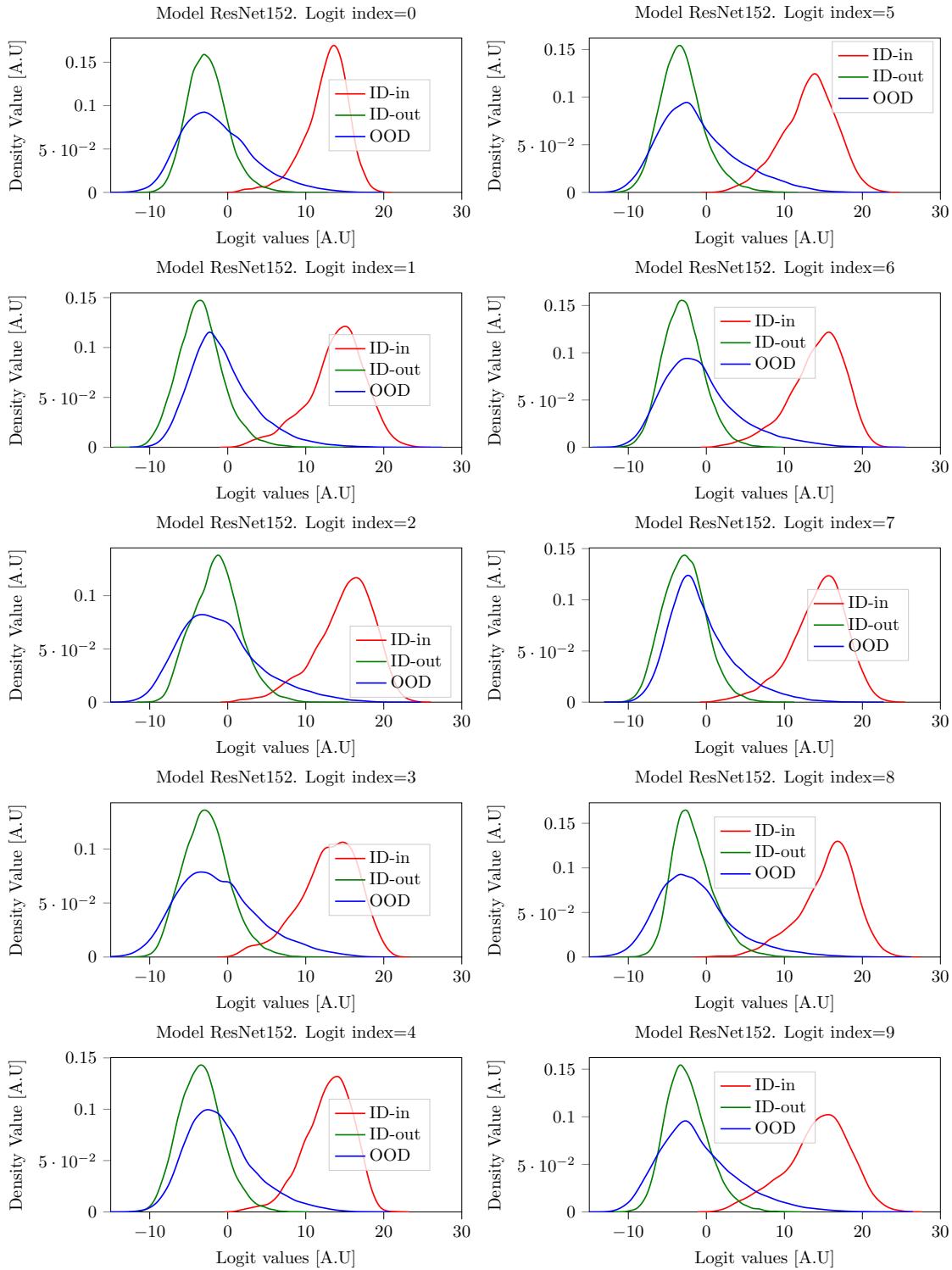


Figure 38: Logit cell densities for CIFAR-10 as ID with ResNet152.

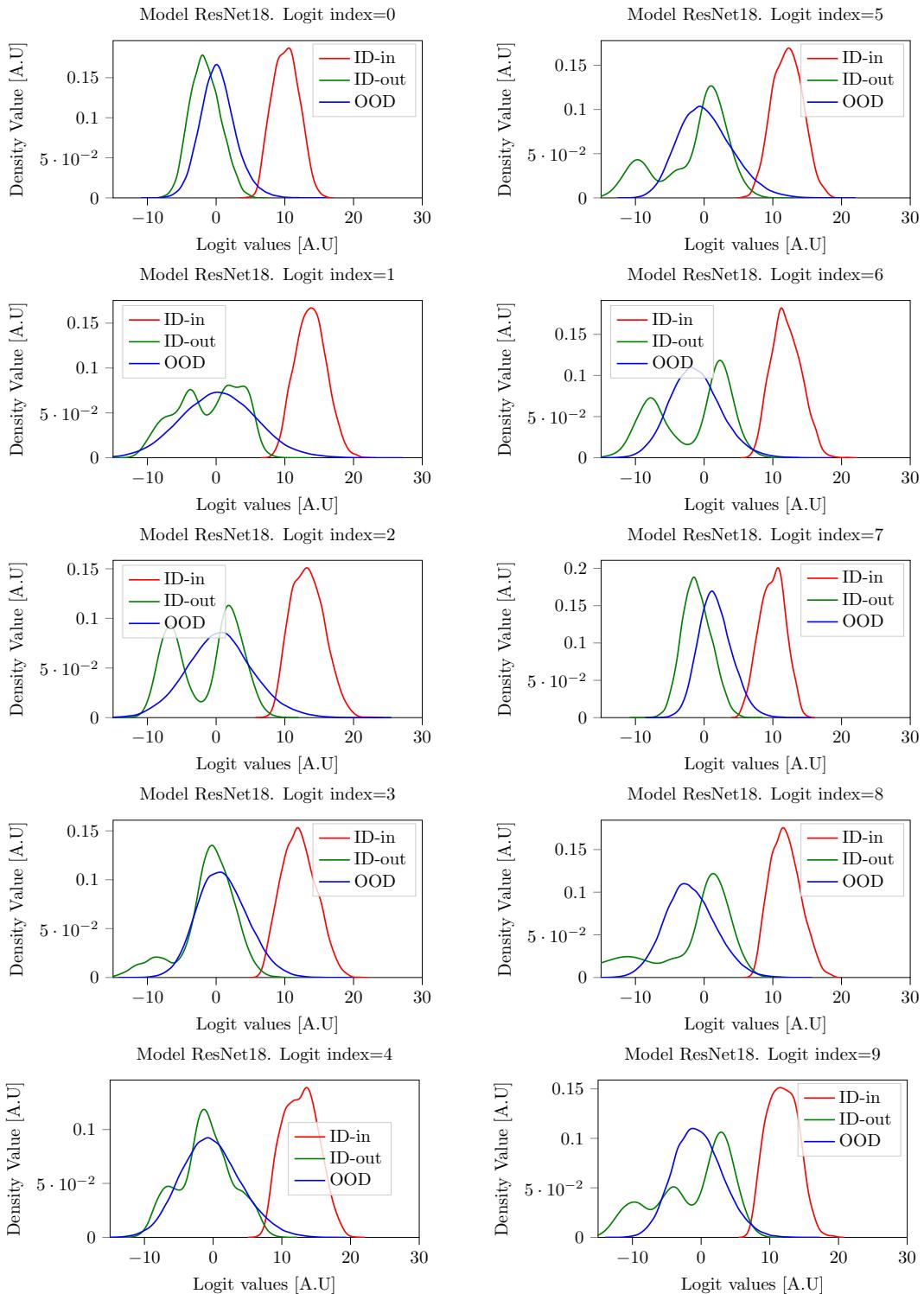


Figure 39: Logit cell densities for SVHN as ID with ResNet18.

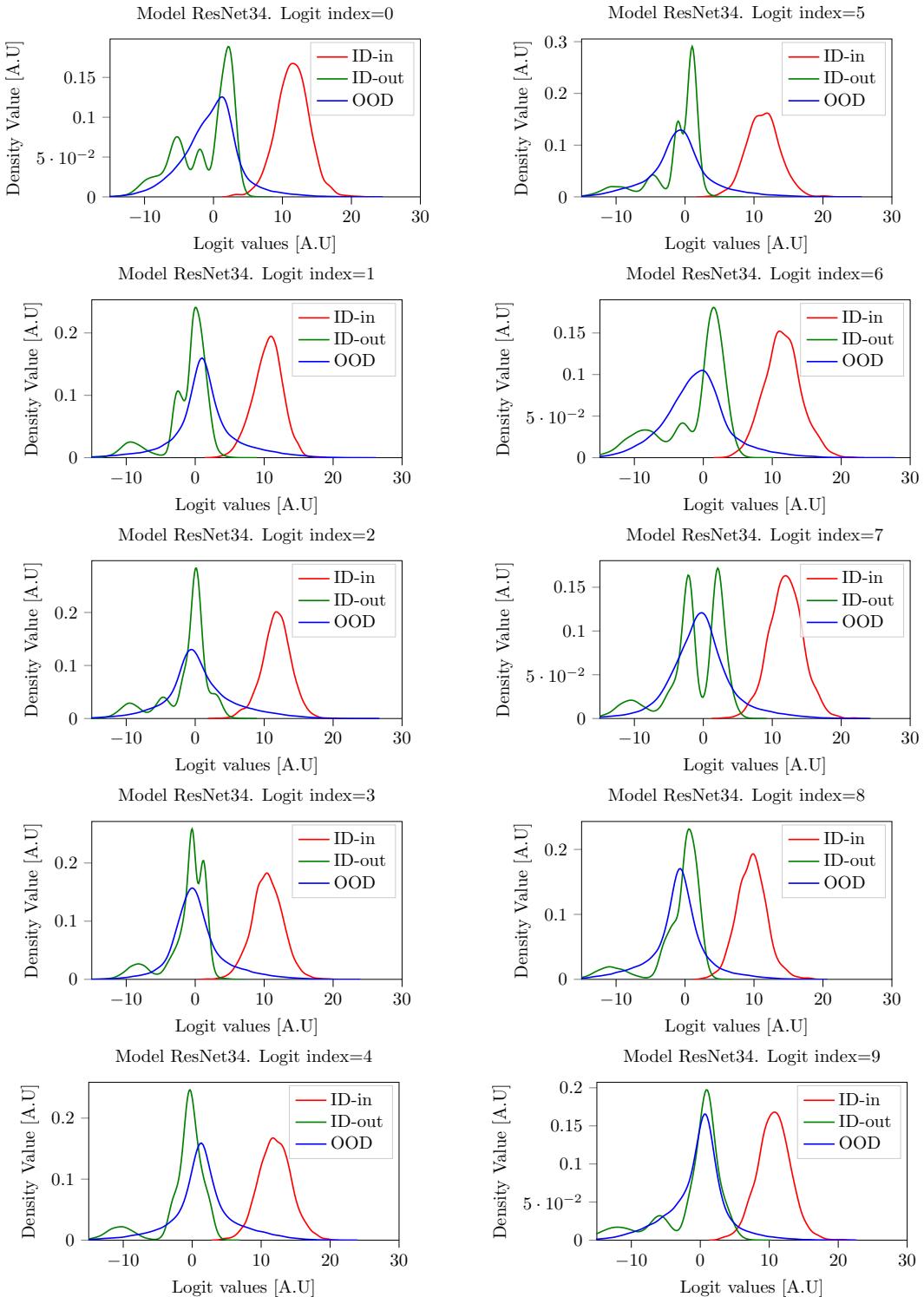


Figure 40: Logit cell densities for SVHN as ID with ResNet34.

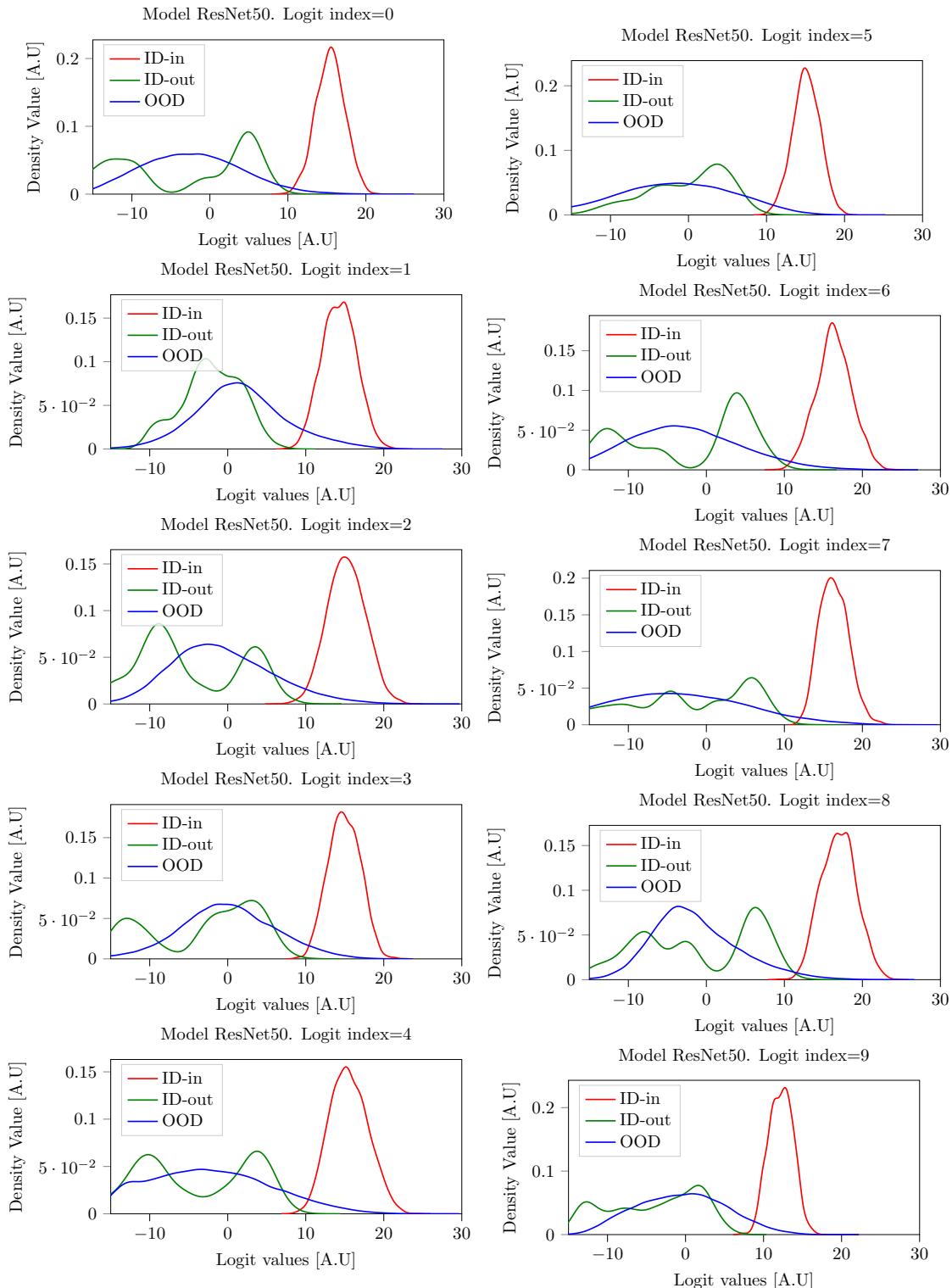


Figure 41: Logit cell densities for SVHN as ID with ResNet50.

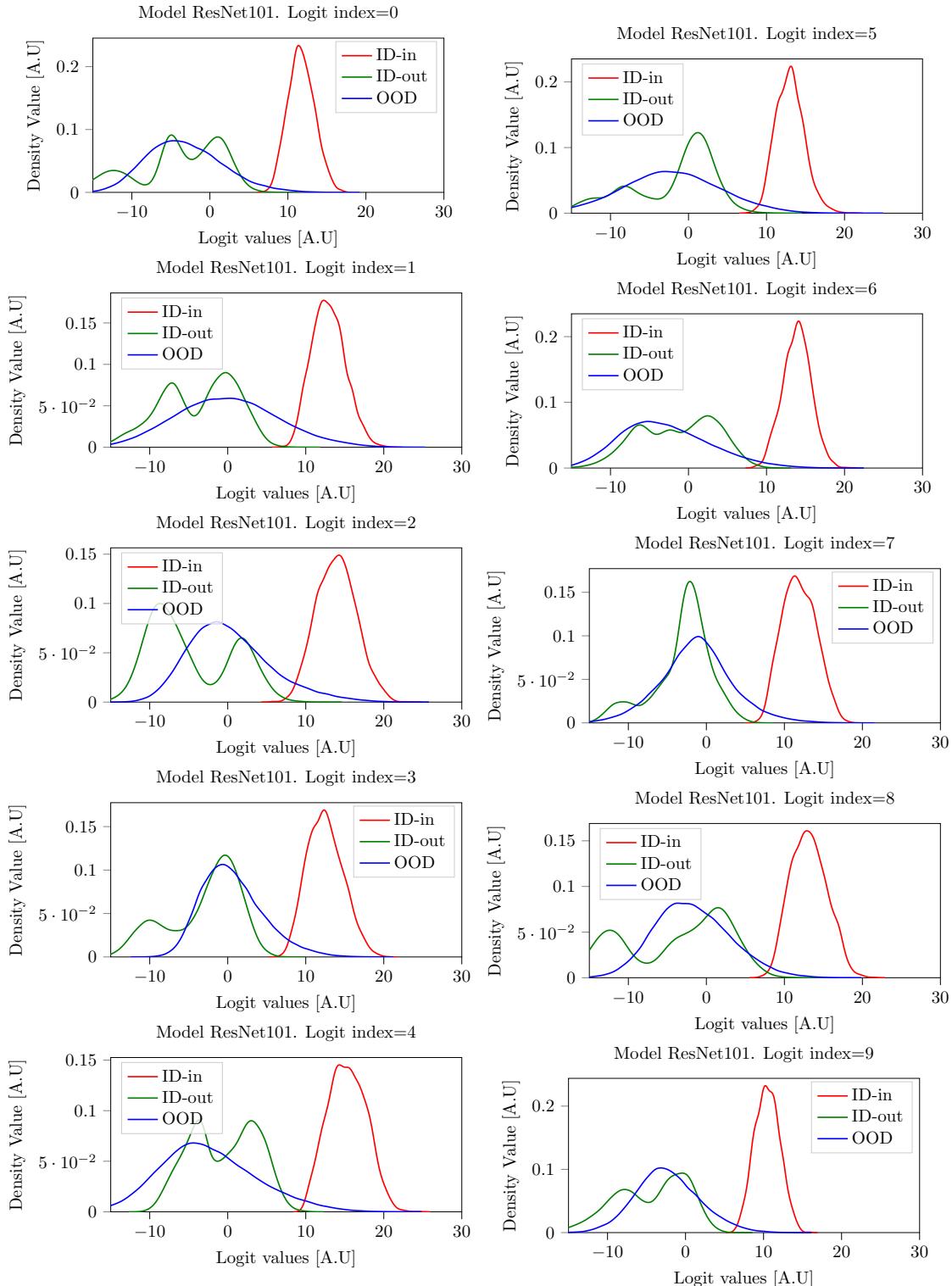


Figure 42: Logit cell densities for SVHN as ID with ResNet101.

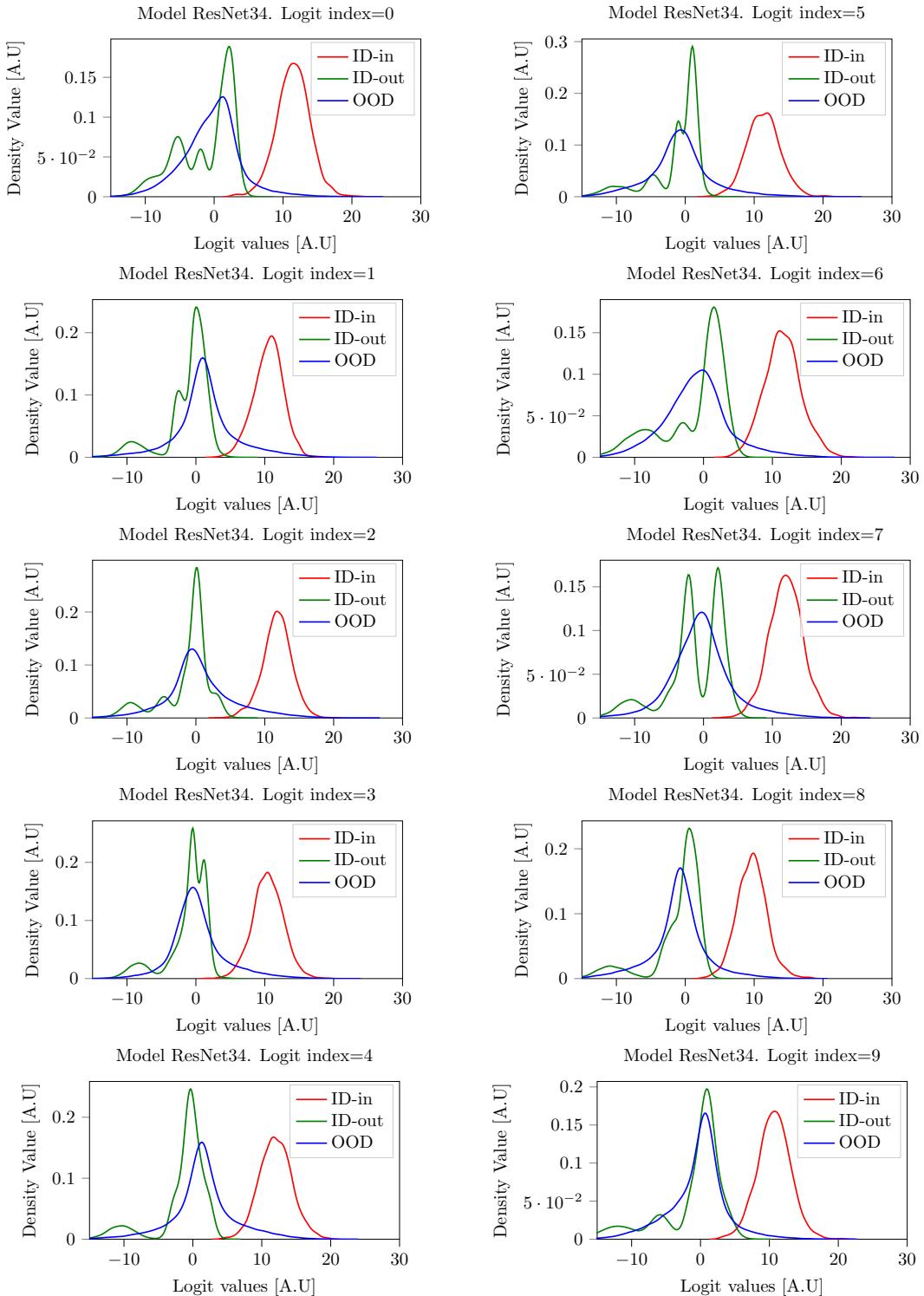


Figure 43: Logit cell densities for SVHN as ID with ResNet152.

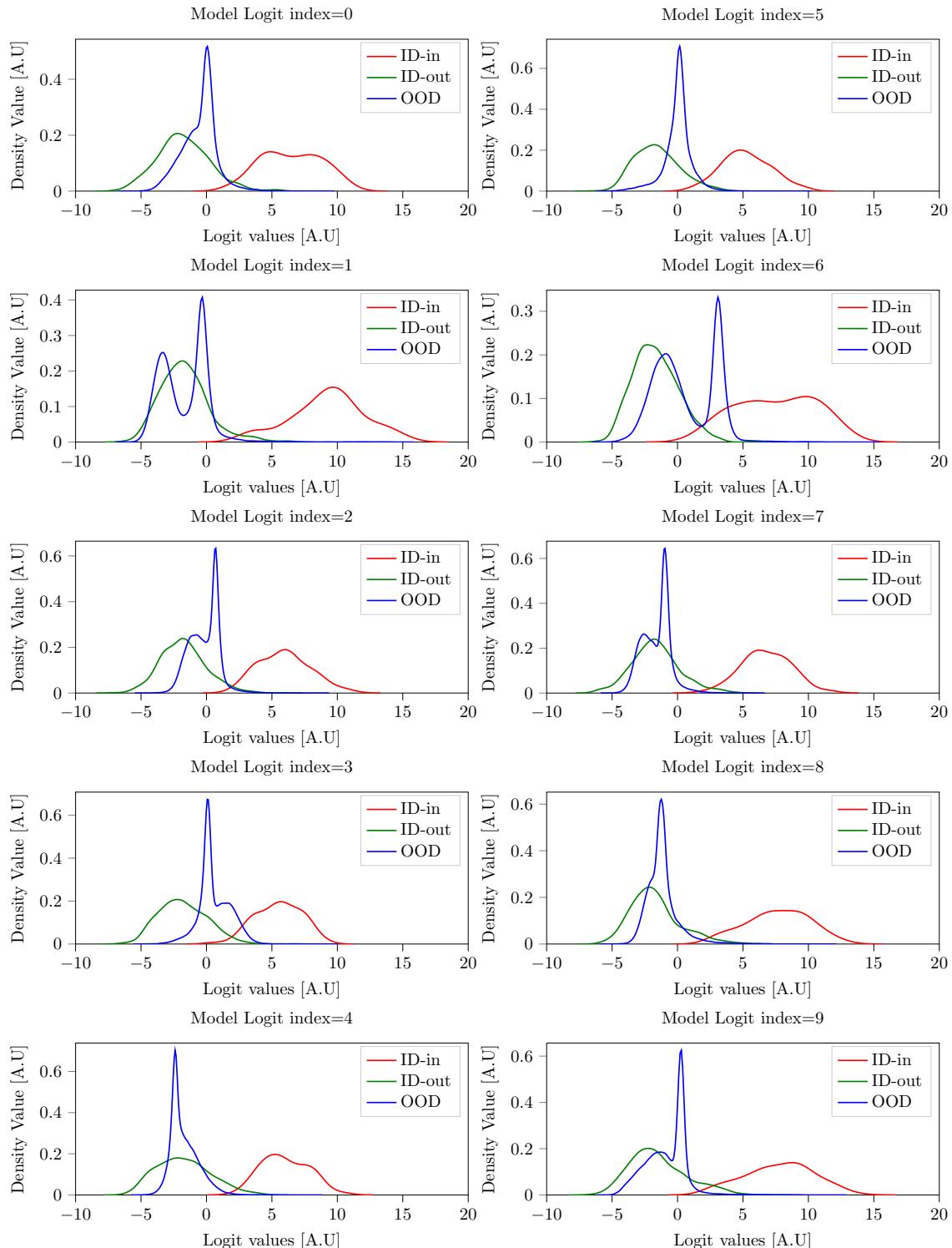


Figure 44: Logit cell densities for CIFAR-10 as ID with ViT-B-16.

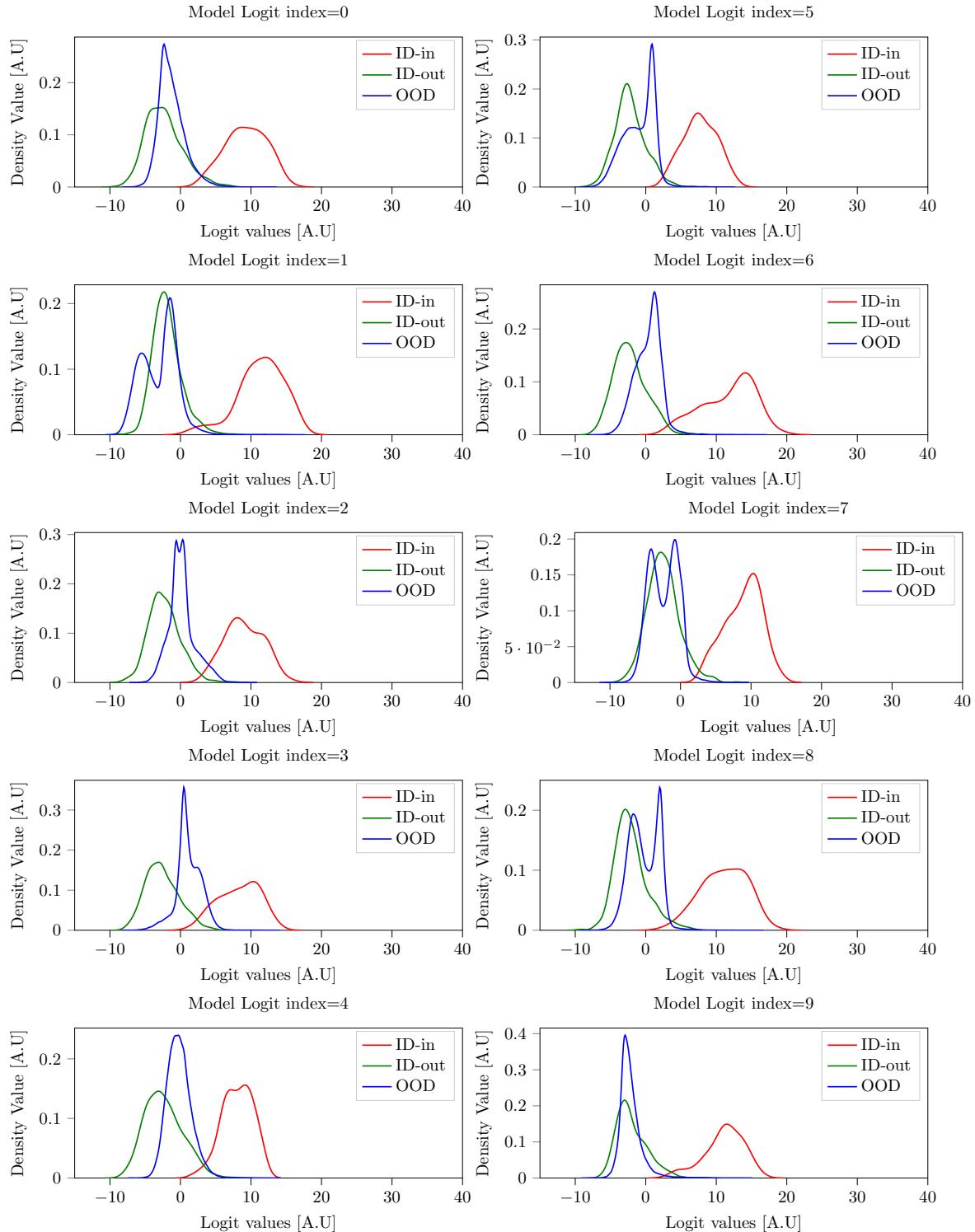


Figure 45: Logit cell densities for CIFAR-10 as ID with ViT-B-32.

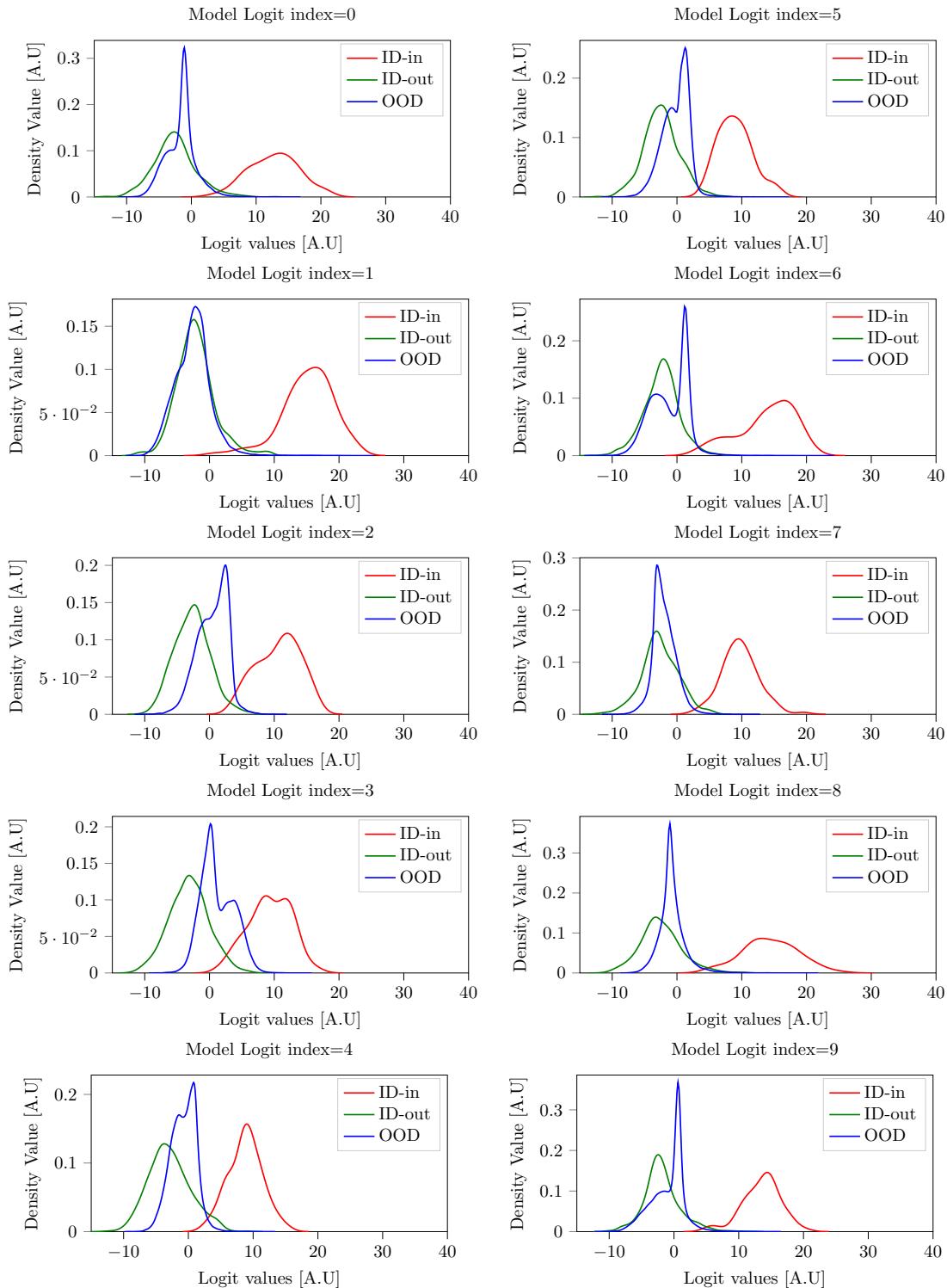


Figure 46: Logit cell densities for CIFAR-10 as ID with ViT-L-16.

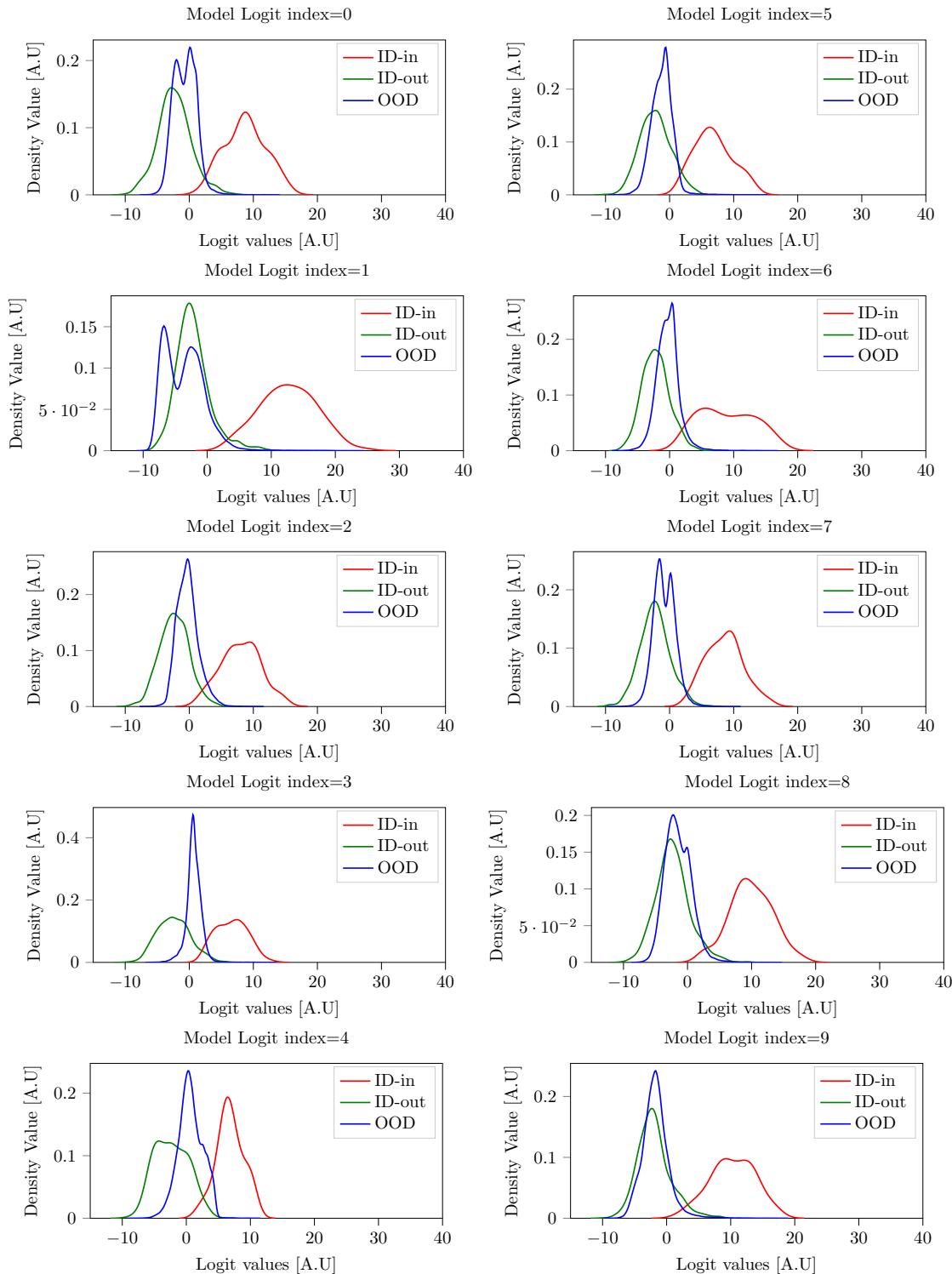


Figure 47: Logit cell densities for CIFAR-10 as ID with ViT-L-32.

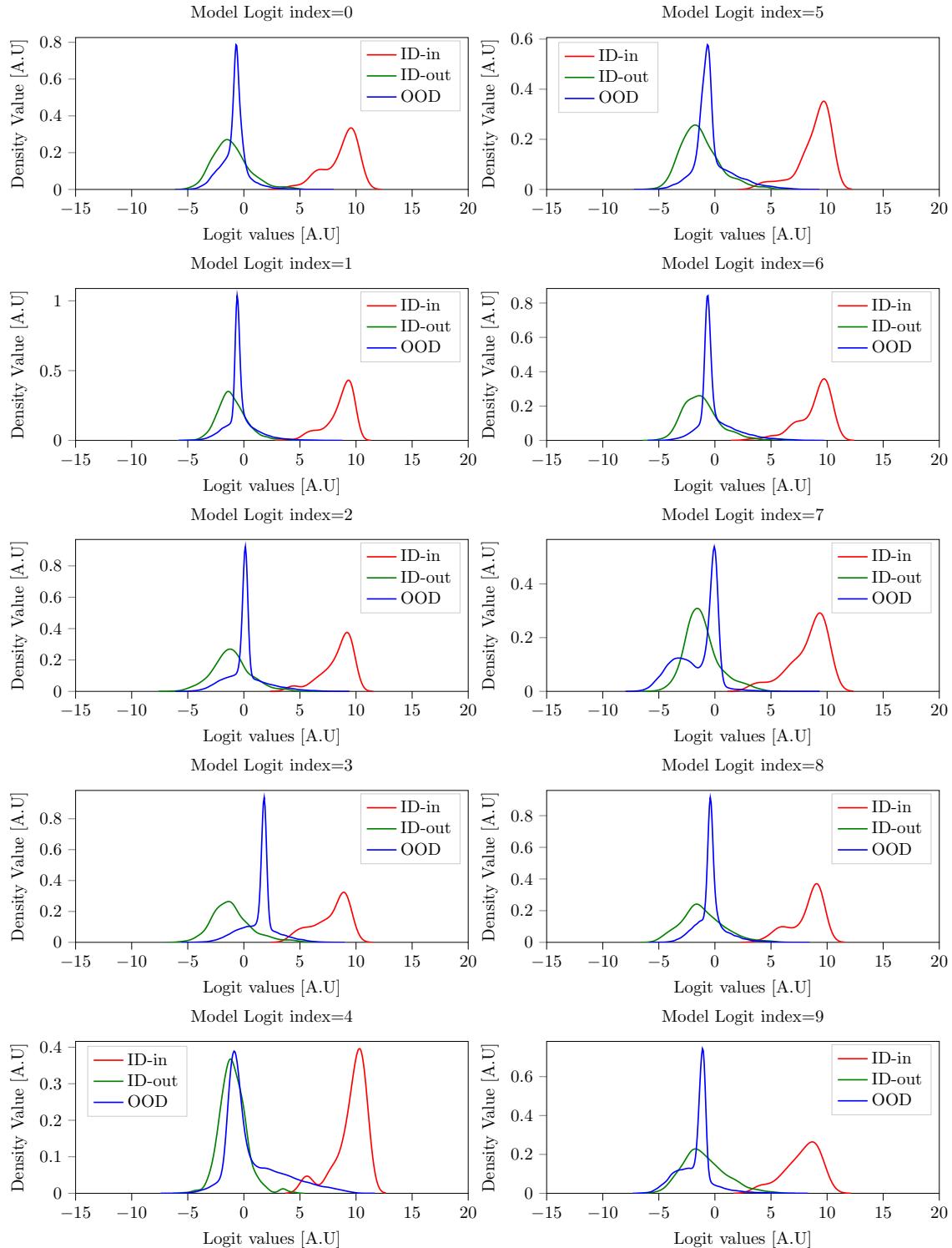


Figure 48: Logit cell densities for SVHN as ID with ViT-B-16.

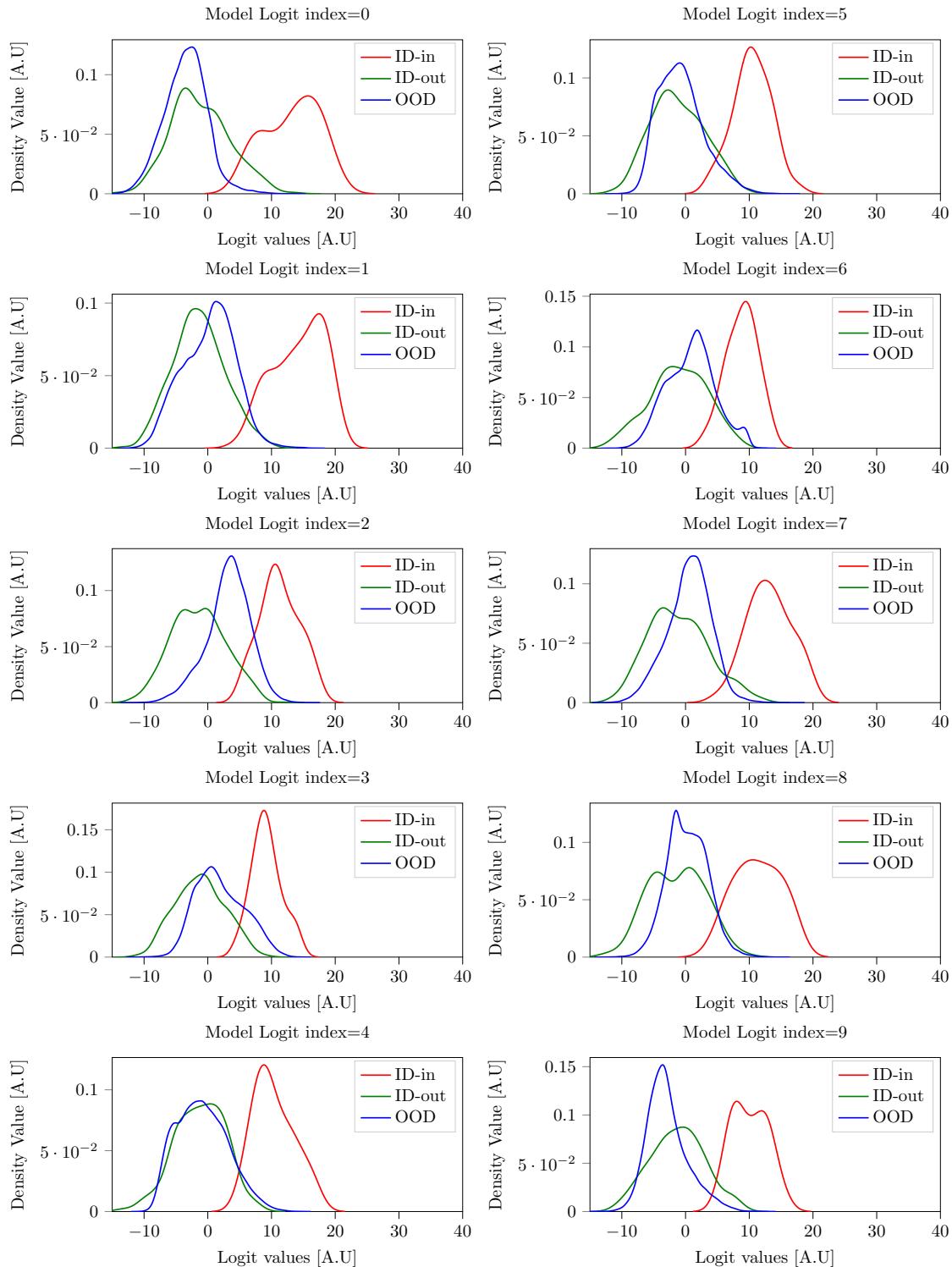


Figure 49: Logit cell densities for SVHN as ID with ViT-B-32.

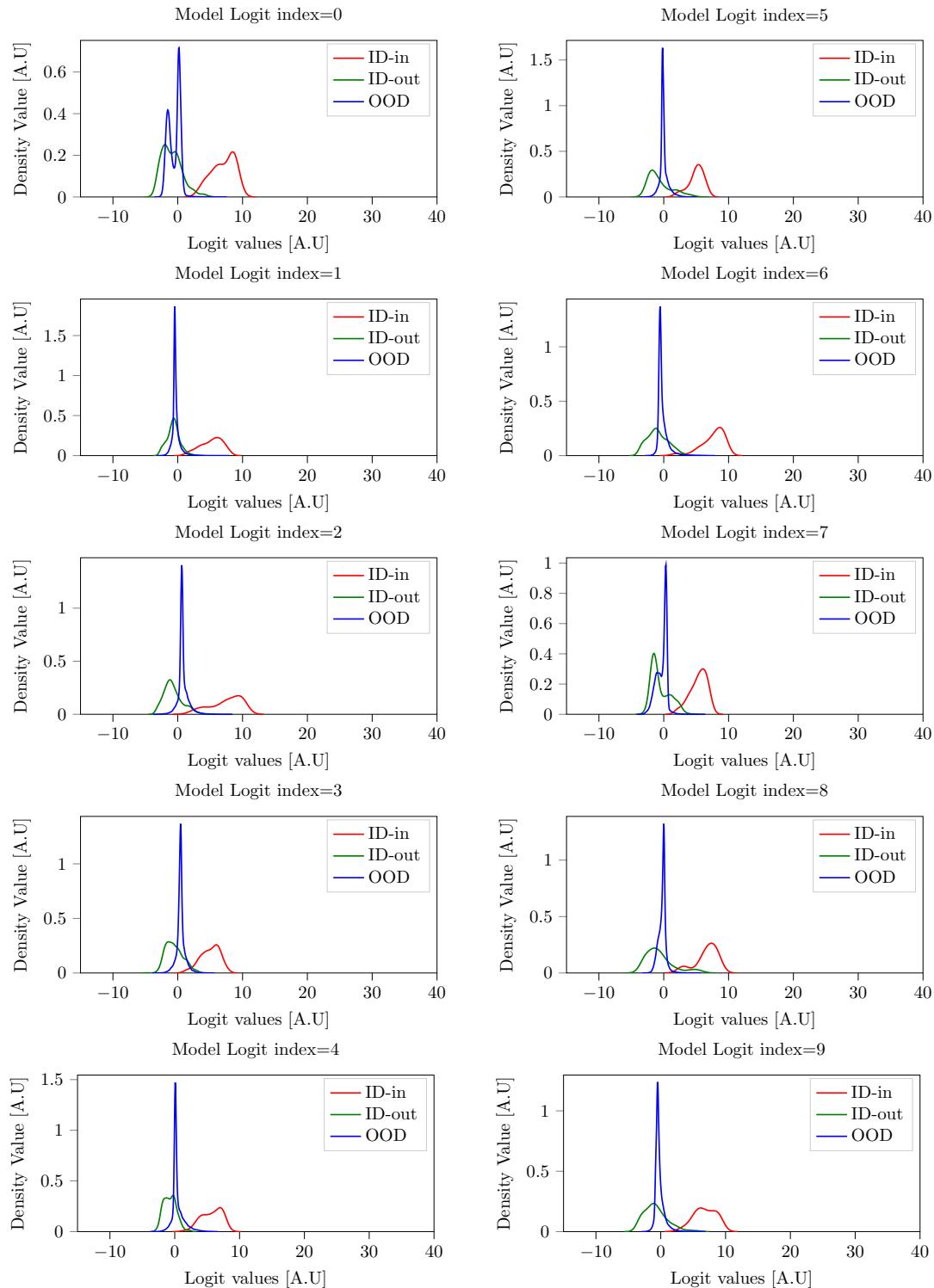


Figure 50: Logit cell densities for SVHN as ID with ViT-L-16.

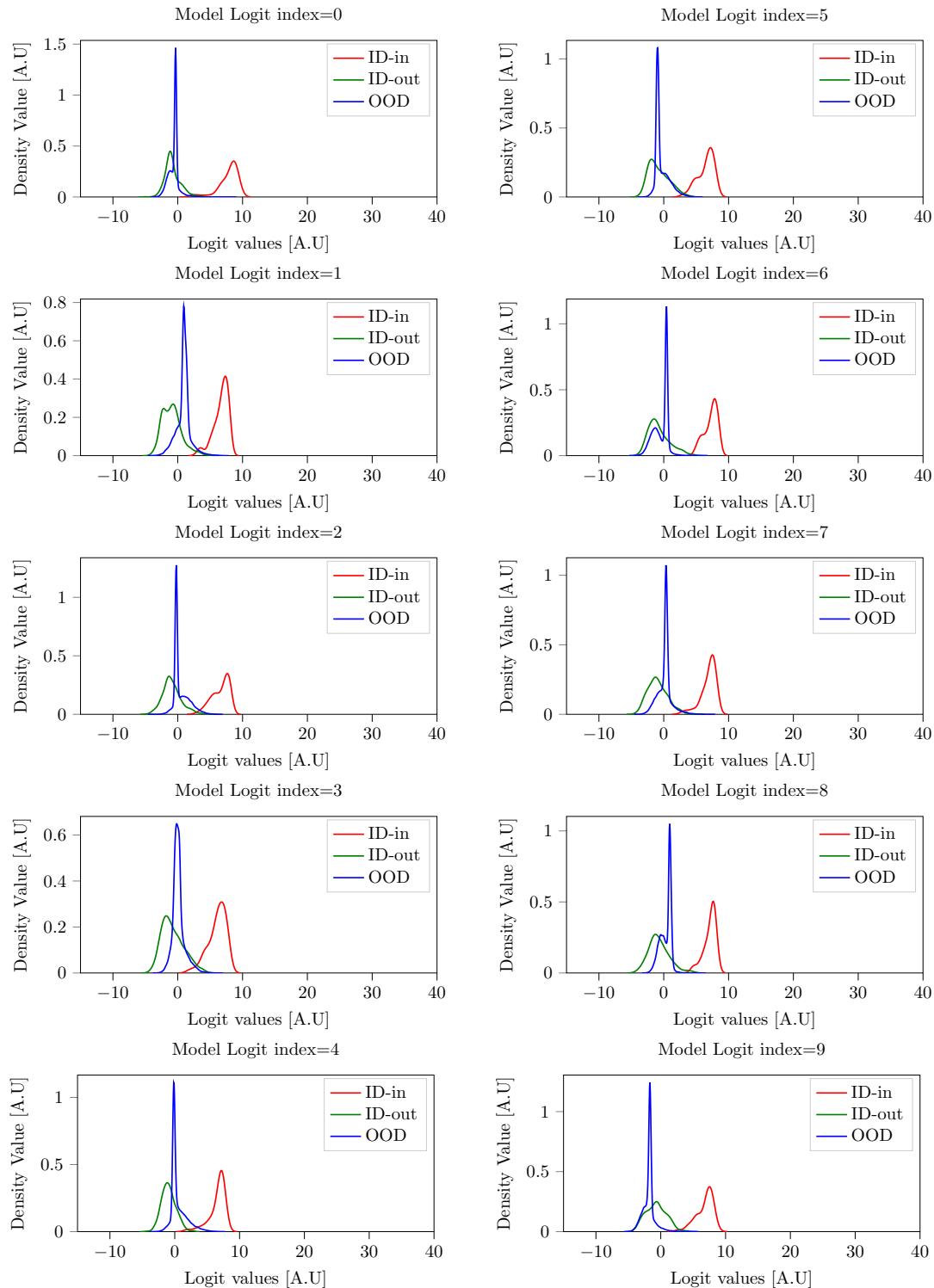


Figure 51: Logit cell densities for SVHN as ID with ViT-L-32.

E Additional experimentation on grayscale image

The classifier model, which is used for this purpose, consists of three convolutional layers followed by two fully connected layers (see table 2). This model is then trained using the Adam optimizer (Kingma & Ba, 2017) via a learning rate of $lr = 10^{-4}$ with weight decay $w_{decay} = 10^{-6}$ and with $\beta_1 = 0.8$ and $\beta_2 = 0.999$. A batch size of 256 is applied for both test and train data. No augmentation or regularization is applied to the training process. ReLU activation is utilized at every layer of the network. For a detailed visual analysis of each logit cell, refer to figs. 52 and 53.

Table 2: A convolutional neural network model for the experiment on the fashion-MNIST and MNIST datasets.

Layer (Type)	Matrix	Nr Parameters
Conv2d-1	[64,28,28]	640
BatchNorm2d-2	[64,28,28]	128
ReLU	[64,28,28]	0
Conv2d-3	[128,14,14]	73,856
BatchNorm2d-4	[128,14,14]	256
ReLU	[128,14,14]	0
Conv2d-5	[256,5,5]	295,168
BatchNorm2d-6	[256,5,5]	512
ReLU	[256,5,5]	0
Linear-7	[16]	16,400
ReLU	[16]	0
Linear-8	[10]	170

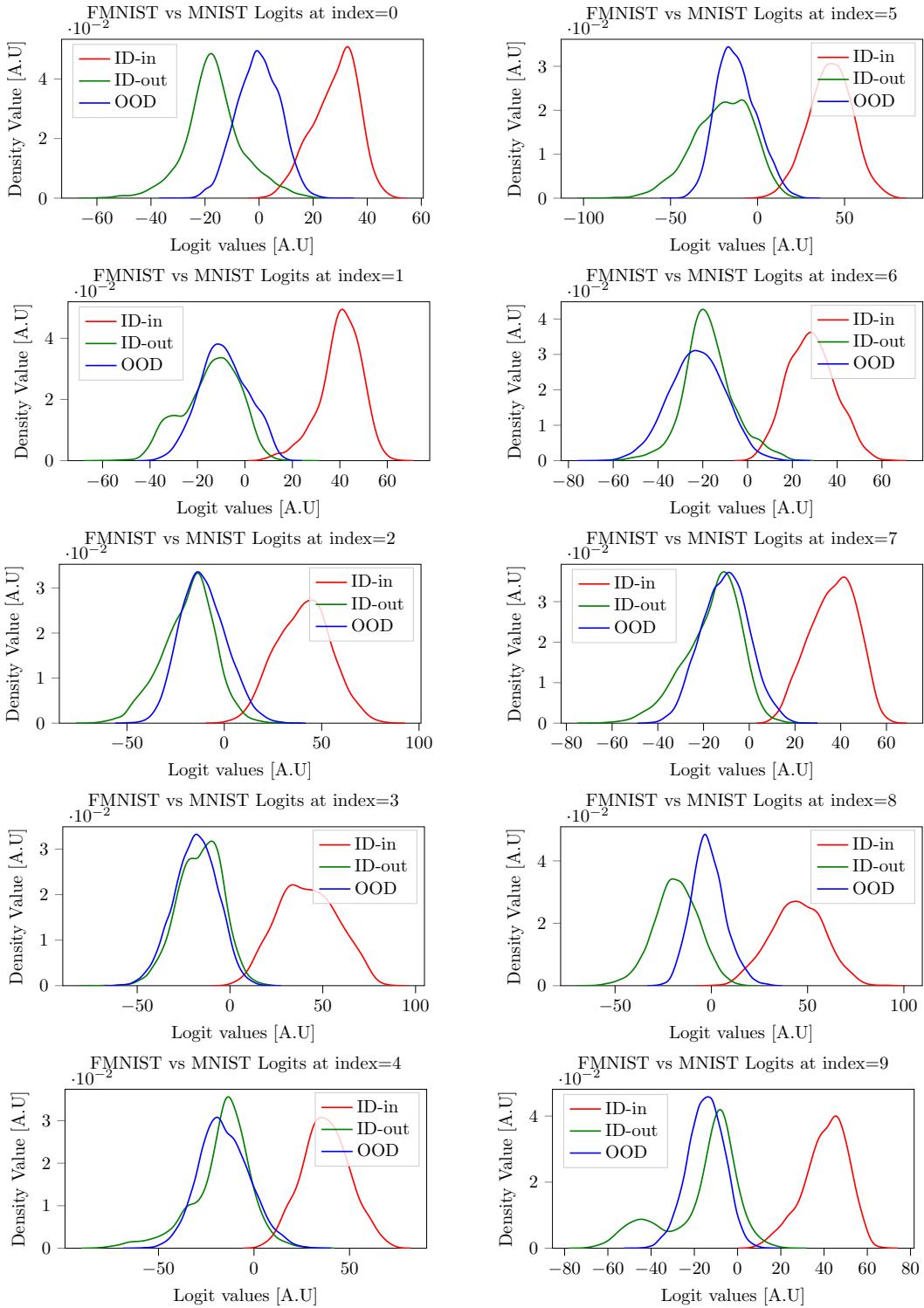


Figure 52: Logit cell densities for FMNIST as ID and MNIST as OOD with the model in table 2.

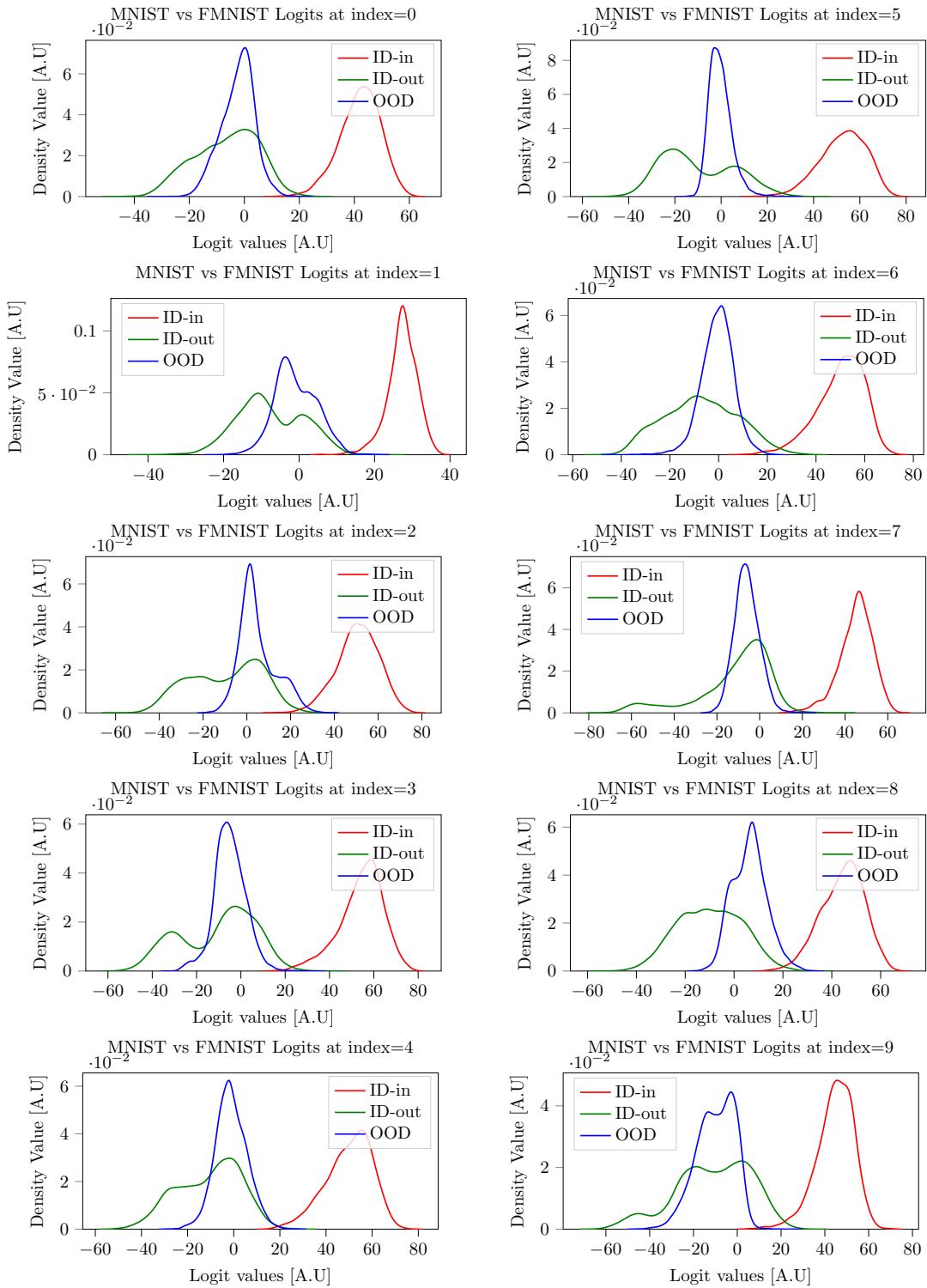


Figure 53: Logit cell densities for MNIST as ID and FMNIST as OOD with the model in table 2.