

# On the Dynamics & Transferability of Latent Generalization during Memorization

Anonymous authors

Paper under double-blind review

## Abstract

Deep networks have been known to have extraordinary generalization abilities, via mechanisms that aren't yet well understood. It is also known that upon shuffling labels in the training data to varying degrees, deep networks, trained with standard methods, can still achieve perfect or high accuracy on this corrupted training data. This phenomenon is called *memorization*, and typically comes at the cost of poorer generalization to true labels. Recent work has demonstrated, surprisingly, that the internal representations of such models retain significantly better latent generalization abilities than is directly apparent from the model. In particular, it has been shown that such latent generalization can be recovered via simple probes (called MASC probes) on the layer-wise representations of the model. However, several basic questions about this phenomenon of latent generalization remain poorly understood: (1) What is the origin and dynamics over training of latent generalization during memorization? Specifically, is it the case that model generalization and latent generalization use largely the same underlying mechanisms? (2) Is the specific nature of the probe critical for our ability to extract latent generalization from the model's layerwise outputs? (3) Does there exist a way to immediately transfer latent generalization to model generalization by suitably modifying model weights directly? On the one hand, this question is conceptually important because it establishes conclusively that the latent generalization manifested by the probe is also within reach of the model, with exactly the information that the model was provided during training, namely the corrupted training data. On the other hand – and more pragmatically – it also suggests the possibility of "repairing" a trained model that has memorized, without requiring expensive retraining from scratch. To address (1), we track the training dynamics, empirically, and find that latent generalization abilities largely peak early in training, with model generalization, suggesting a common origin for both. However, while model generalization degrades steeply over training thereafter, latent generalization falls more modestly & plateaus at a higher level over epochs of training. These experiments lend circumstantial evidence to the hypothesis that latent generalization uses largely similar mechanisms as those that underlie the model's generalization in the early phases of training. To investigate (2), we examine the MASC probe and show that it is a quadratic classifier. The question in (2) thus becomes whether the quadratic nature of the MASC probe underlies its remarkable effectiveness in extracting latent generalization. If this were so, a linear probe constructed along these lines would not be as effective. To investigate this, we designed a new linear probe for this setting, and find, surprisingly, that it has superior generalization performance in comparison to the quadratic probe, in most cases. This suggests that the quadratic nature of the probe is not critical in extracting latent generalization. Importantly, the effectiveness of the linear probe enables us to answer (3) in the affirmative. Specifically, using this new linear probe, we devise a way to transfer the latent generalization present in last-layer representations to the model by directly modifying the model weights. This immediately endows such models with improved generalization, i.e. without additional training. Our findings provide a more detailed account of the rich dynamics of latent generalization during memorization, provide clarifying insight on the specific role of the probe in latent generalization, as well as demonstrate the means to leverage this understanding to directly transfer this generalization to the model.

# 1 Introduction

Overparameterized deep neural networks have seen widespread deployment in many fields, due to their remarkable generalization abilities. However, we still don’t have a clear understanding of the mechanisms underlying their ability to generalize so well to unseen data. It has also been shown (Zhang et al., 2017; 2021) that overparameterized deep networks are capable of achieving high or even perfect training accuracy on datasets, wherein a subset of training data have their labels randomly shuffled. Such models typically have poor generalization performance, i.e. poorer accuracy on test data with correct labels – a phenomenon that has been called *memorization*. It is known (Arpit et al., 2017) that during training, models trained with such corrupted datasets exhibit better generalization during the initial phases of training; however generalization progressively deteriorates as training accuracy improves subsequently.

A recent study (Ketha & Ramaswamy, 2025) has shown that while deep networks trained on datasets having corrupted labels tend to exhibit poor generalization, their intermediate-layer representations retain a surprising degree of latent generalization ability. This ability can be recovered from such trained networks by using a simple probe – the Minimum Angle Subspace Classifier (MASC) – that utilizes the subspace geometry of the corrupted training dataset representations, to this end. Their findings suggest that generalizable features are present in the layer-wise representations of such networks, even when the model fails to utilize them sufficiently. However, the origin and evolution of this latent generalization ability during training are not well understood. In particular, it is unclear whether model generalization and latent generalization arise from the same underlying mechanisms, and whether the specific nature of the probe (MASC) is critical for extracting latent generalization from layer-wise representations. Moreover, it is not even clear whether this latent generalization can be directly transferred to the model. That is, whether it is, in principle, possible to directly modify model weights, so it overtly acquires this latent generalization performance, is not known. Here, we address these questions.

Our main contributions are listed below.

- To investigate the origin and training dynamics of latent generalization, we analyze models trained on standard datasets with varying levels of label corruption and use MASC (Ketha & Ramaswamy, 2025) to characterize how their latent generalization ability evolves over the course of training. We find that as the model exhibits a peak in its test accuracy early in training (Arpit et al., 2017), the MASC test accuracy at all layers also tends to largely peak concurrently with that of the model, albeit at different levels. Following this, the evolution of test accuracies between the model & MASC diverge, with the model showing a marked decline in test accuracies over further epochs of training. In contrast, the MASC test accuracies decline more modestly & tend to plateau higher, which manifests in the improved generalization ability at the end of training, as reported in (Ketha & Ramaswamy, 2025).
- We observe that MASC is a non-linear classifier; and in particular, prove that it is a quadratic classifier. This brings up the possibility that the improved generalization performance of the probe (i.e. MASC) is attributable to the effectiveness of the quadratic nature of the probe itself and may not easily be decodable, e.g. by a linear probe. To address this point, we introduce a simple linear alternative – Vector Linear Probe Intermediate-layer Classifier (VeLPIC). Surprisingly, we find that VeLPIC almost always achieves superior latent generalization performance in comparison to MASC, especially for higher corruption degrees. This establishes that latent generalization during memorization is linearly decodable from layerwise representations.
- By leveraging the linear probe (VeLPIC), we devise a way to directly modify the pre-softmax weights of such deep networks, that immediately transfers to the model, the latent generalization performance of VeLPIC (as applied to the last layer). Notably, this is without requiring additional training. This demonstrates that latent generalization present in layerwise representations can be transferred directly to enhance the model generalization of deep network models, in the memorization regime.

## 2 Related work

In influential work, (Zhang et al., 2017; 2021) showed that deep networks can achieve perfect training accuracy even with randomly shuffled labels, accompanied by poor generalization. In follow-up work, (Arpit et al., 2017) find that in the memorization regime, networks learn simple patterns first during training. Their work provides a detailed account of the early dynamics of training. More recently, (Ketha & Ramaswamy, 2025) in fact show that in spite of the fall in model generalization later on in training, the layerwise representations of the model retain significant latent generalization ability. (Arpit et al., 2017) also show that regularization can help models resist memorization in the label noise case.

Analyzing intermediate representations in deep networks has been previously explored using kernel-PCA (Montavon et al., 2011) and linear classifier probes (Alain & Bengio, 2018). Notably, (Alain & Bengio, 2018) state that they deliberately did not probe deep networks in the memorization setting since they thought that such probes would inevitably overfit. On the contrary, (Ketha & Ramaswamy, 2025) demonstrate that probes on deep networks in the memorization setting, can have enhanced generalization. (Stephenson et al., 2021) show evidence suggesting that memorization occurs in the later layers. Li et al. (2020) show that in the memorization regime, there is substantial deviation from initial weights.

Several training paradigms have been proposed to enhance generalization performance when learning from corrupted datasets. For example, MentorNet (Jiang et al., 2018) introduces a framework wherein a mentor network guides the learning process of a student network by guiding the student model to focus on likely clean labels. Likewise, Co-Teaching (Han et al., 2018) trained two peer networks simultaneously, each selecting small-loss examples to update its counterpart. Early-Learning Regularization (ELR) (Liu et al., 2020) augmented the training objective with a regularization term, towards this end.

Saxe et al. (2013) offer theoretical explanations on generalization for deep linear networks and Lampinen & Ganguli (2018) offer theoretical explanations in the memorization regime. Methodologies such as Canonical Correlation Analysis (Raghu et al., 2017; Morcos et al., 2018) and Centered Kernel Alignment (Kornblith et al., 2019) have been used to characterize training dynamics and network similarity. Representational geometry and structural metrics provide further insights into learned representation properties (Chung et al., 2016; Cohen et al., 2020; Sussillo & Abbott, 2009; Farrell et al., 2019; Bakry et al., 2015; Cayco-Gajic & Silver, 2019; Yosinski et al., 2014).

## 3 Experimental setup

We demonstrate results for the same set of models and datasets as presented in (Ketha & Ramaswamy, 2025). Specifically, we use Multi-Layer Perceptrons (MLPs) trained on the MNIST (Deng, 2012) and CIFAR-10 (Krizhevsky, 2009) datasets; Convolutional Neural Networks (CNNs) trained on MNIST, Fashion-MNIST (Xiao et al., 2017), and CIFAR-10; and AlexNet (Krizhevsky et al., 2012) trained on Tiny ImageNet dataset (Moustafa, 2017). Additional details on the model architectures and training are available in the Appendix Section A.

Each model was trained under two distinct schemes: (i) using training data with true labels, referred to as “generalized models,” and (ii) using training data with labels randomly shuffled to varying degrees (referred to as “memorized models” Zhang et al. (2021)). Similar to Ketha & Ramaswamy (2025), we train the aforementioned models using corruption degrees of 0%, 20%, 40%, 60%, 80%, and 100%. Training with a corruption degree  $c$  implies that, with probability  $c$ , the label of a training datapoint is changed with a randomly selected label drawn uniformly from the set of possible classes. This may result in the label remaining the same after the change as well. All models were trained either until achieving high training accuracy (99% or 100%) or for a maximum of 500 epochs, whichever occurred first.

To study the dynamics of the training process, we conducted the experiments on model checkpoints saved at various stages of training. Specifically, we began with the randomly initialized model (corresponding to epoch 0), followed by checkpoints saved at every second epoch up to the 20th epoch. Beyond epoch 20, results are shown at intervals of five epochs for the MLP and CNN models, and at intervals of ten epochs

for the AlexNet model. The reported results are averaged over three independent training runs, with shaded regions in the plots indicating the range across instances.

Ketha & Ramaswamy (2025) investigate the organization of class-conditional subspaces using the training data at various layers of deep networks. These subspaces are estimated via Principal Components Analysis (PCA), specifically, ensuring that they pass through the origin. To probe the layerwise geometry without relying on subsequent layers, they propose a new probe – the Minimum Angle Subspace Classifier (MASC). For a given test input, MASC projects the layer output onto each class-specific subspace, and computes the angles between the original and projected vectors, for each subspace. The label predicted by MASC corresponds to the class whose subspace yields the projected vector with the smallest such angle. We provide a detailed summary of the working of MASC in the Appendix Section B. We have used 99% as the percentage of variance explained by the principal components that form the class-specific subspaces used by MASC, similar to experiments conducted in Ketha & Ramaswamy (2025).

## 4 Training dynamics of latent generalization using MASC

As shown in Ketha & Ramaswamy (2025), for models trained with corrupted labels, there exists at least one layer where MASC exhibits better generalization than the corresponding trained model. However, the origin & evolution of this latent generalization across training isn’t well understood.

Here, we empirically study the behavior of latent generalization, as manifested by MASC, during training. MASC testing accuracy during training for MLP trained on MNIST, CNN trained on Fashion-MNIST and AlexNet trained on Tiny ImageNet are shown in Figure 1. Results with 0% and 100% corruption degrees as well as the results for additional models i.e. MLP trained on CIFAR-10, CNN trained on MNIST, CNN trained on CIFAR-10 for various corruption degrees are presented in Figure 4 and Figure 5, respectively in the Appendix Section C.

Our findings indicate the presence of two distinct phases in the training process, separated by the point at which the model achieves peak test accuracy. For various non-zero degrees of corruption, in most cases (except those with 100% degrees of corruption), the MASC test accuracy largely follows the rise in the model’s test accuracy up to this peak. However, beyond the peak, while the model’s test accuracy declines significantly, the drop in MASC accuracy is less steep and plateaus at a higher level over the epochs.

For non-zero corruption degrees (except those with 100% corruption), in most cases, for MLP, MASC accuracy on later layers performed better than MASC accuracy on early layers, whereas for CNNs MASC accuracy on early layers performed better.

Our results represent progress in clarifying the origin & evolution of latent generalization by MASC, during training. In particular, given that model generalization & latent generalization show a concurrent initial rise, it suggests the possibility of common mechanisms that drive both in the early phases of training. The subsequent divergence between model generalization & latent generalization is an intriguing phenomenon, whose mechanisms merit future investigation.

## 5 Non-linearity of MASC

Classically (Alain & Bengio, 2018), linear probes have been used to probe layers of deep networks. However, (Ketha & Ramaswamy, 2025) do not use the standard linear probe from (Alain & Bengio, 2018). Below, we prove that MASC (Ketha & Ramaswamy, 2025) is in fact a classifier that is quadratic in the layerwise output of the layer that it is applied to.

**Proposition 1.** *MASC is a quadratic classifier.*

*Proof.* Let  $\mathbf{x}_l$  denote the output of the layer  $l$  of the deep network when it is given input  $\mathbf{x}$ . Let  $\mathbf{p}_1^c, \mathbf{p}_2^c, \dots, \mathbf{p}_k^c$  be a basis<sup>1</sup> of the subspace  $\mathcal{S}_c$  corresponding to class  $c$ . Let  $\mathbf{x}_l^c$  be the projection of  $\mathbf{x}_l$  on  $\mathcal{S}_c$ . We have

$$\mathbf{x}_l^c = (\mathbf{x}_l \cdot \mathbf{p}_1^c) \mathbf{p}_1^c + \dots + (\mathbf{x}_l \cdot \mathbf{p}_k^c) \mathbf{p}_k^c \quad (1)$$

<sup>1</sup>which is typically estimated via PCA, where  $k$  is the number of principal components.

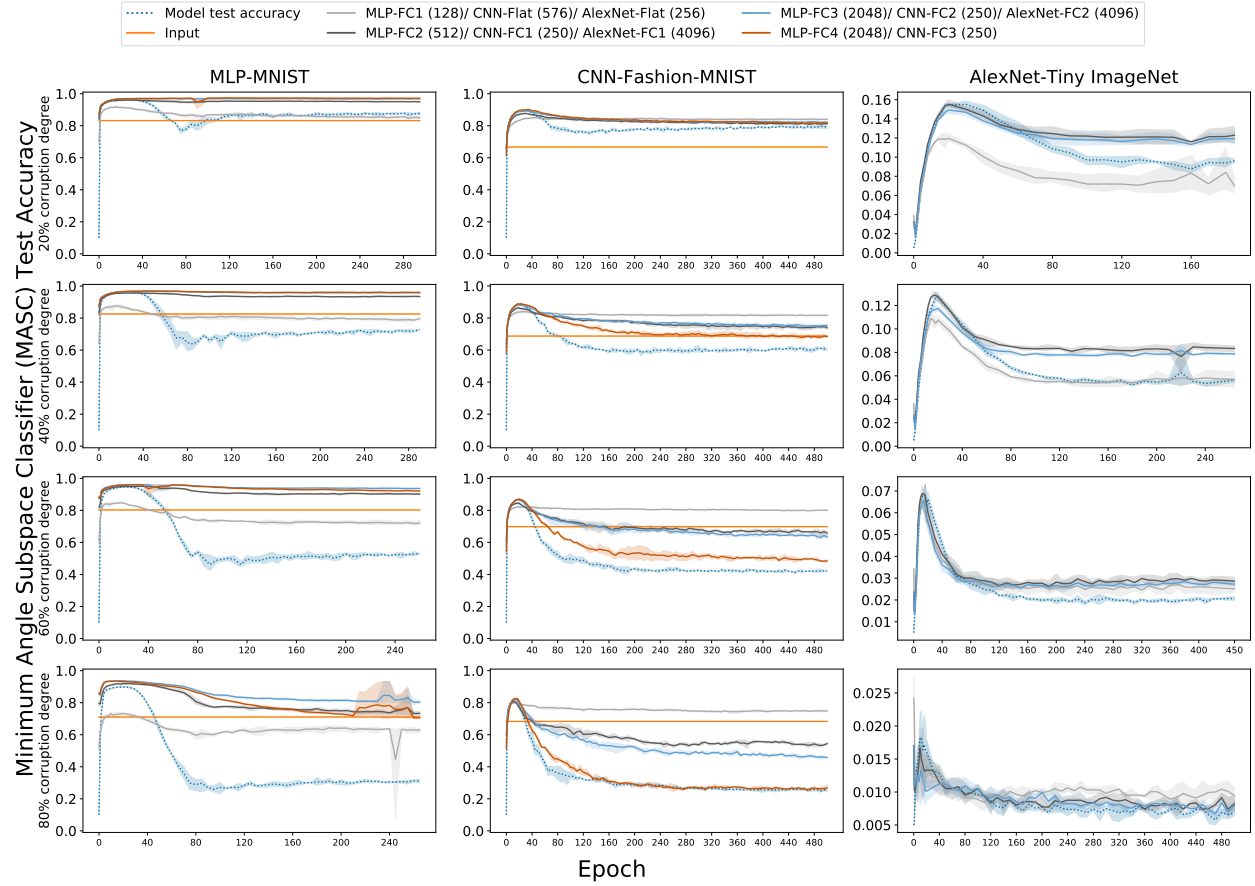


Figure 1: Minimum Angle Subspace Classifier (MASC) test accuracy over epochs of training for multiple models/datasets, where test data is projected onto class-specific subspaces constructed at each epoch from corrupted training data with the indicated label corruption degree. The plots display MASC accuracy across different layers of the network. For reference, the evolution of test accuracy of the corresponding model (blue dotted line) over epochs of training is also shown. FC denotes fully connected layers with *ReLU* activation, and Flat refers to the flatten layer without *ReLU*.

Now, MASC on layer  $l$  predicts the label of  $\mathbf{x}$  as

$$\arg \max_c (\mathbf{x}_l \cdot \mathbf{x}_l^c) = \arg \max_c ((\mathbf{x}_l \cdot \mathbf{p}_1^c)^2 + \dots + (\mathbf{x}_l \cdot \mathbf{p}_k^c)^2) \quad (2)$$

which is quadratic in  $\mathbf{x}_l$ . This establishes that MASC is a quadratic classifier.  $\square$

## 6 Vector Linear Probe Intermediate-layer Classifier (VeLPIC): A new linear probe

Given that MASC is inherently a non-linear classifier as proved above, a natural question is if its extraordinary ability to decode generalization from hidden representations of memorized networks is a consequence of its non-linearity. In other words, if the quadratic nature of MASC is indeed responsible for its effectiveness, then a corresponding linear probe would be expected to perform substantially worse. It raises the question of whether the latent generalization reported in (Ketha & Ramaswamy, 2025) is linearly decodable – with comparable performance – from the layerwise representations of the network.

To investigate this, we build a linear probe analogous to MASC. We sought to retain the same broad idea, namely determine an instance of a mathematical object per class and measure closeness of the layerwise

output of an incoming datapoint to these objects with the prediction corresponding to the class whose object was closest in this sense. In contrast to (Alain & Bengio, 2018), where parameters of their linear probe are learned iteratively by minimizing a cross-entropy loss, we seek to determine the linear probe parameters directly via the geometry of the class-conditional training data. We choose to simply use a vector as this mathematical object and measure closeness in the angle sense. We call this probe the Vector Linear Probe Intermediate-layer Classifier (VeLPIC). As we discuss subsequently, we find, surprisingly, that this choice is significantly more effective than MASC, in most cases. Secondly, we show that we can use the parameters of the probe as applied to the last layer, to modify the model weights to immediately confer the corresponding generalization to the model.

We now discuss how the vector corresponding to each class in VeLPIC is constructed. Each class vector is determined using only the top principal component from PCA run on augmented<sup>2</sup> class-conditional corrupted training data. However, the first principal component can manifest in two opposite directions (i.e. the vector or its negative). This is important here<sup>3</sup> because incoming data vectors can be “close” to this class vector, even though their angles are obtuse and closer to 180°. VeLPIC resolves this directional issue by aligning the class vector based on the sign of the projection of the training data mean; if the mean of the training data projected on this principal component is negative, the direction of the principal component is flipped to obtain the class vector; otherwise, it is retained as is.

Formally, for a given test data point  $\mathbf{x}$ , let  $\mathbf{x}_l$  denote its activation at layer  $l$  obtained in the forward pass of  $\mathbf{x}$  through the deep network until the output of layer  $l$ . For layer  $l$ , let  $\{\mathcal{P}_m\}_{m=1}^M$  be the top principal component vectors, one each per class, of the class-conditional corrupted training data and  $\{T_m\}_{m=1}^M$  be its corresponding<sup>4</sup> projection means, where  $M$  is the number of classes. Let  $\{\mathcal{V}_m\}_{m=1}^M$  be unit vectors representing VeLPIC class vectors. VeLPIC uses  $\{\mathcal{V}_m\}_{m=1}^M$  to predict the label of  $\mathbf{x}_l$  based on its maximum projection among these class vectors<sup>5</sup>, as outlined in Algorithm 1.

---

#### Algorithm 1 Vector Linear Probe Intermediate-layer Classifier (VeLPIC)

---

**Input:** Principal component vectors  $\{\mathcal{P}_m\}_{m=1}^M$ , projection training means  $\{T_m\}_{m=1}^M$ , layer  $l$  output  $\mathbf{x}_l$ , class labels  $\{C_m\}_{m=1}^M$ .

**Output:** Predicted label  $y(\mathbf{x}_l)$ .

```

1: for each class  $m = 1, \dots, M$  do
2:   if  $T_m < 0$  then
3:      $\mathcal{V}_m \leftarrow -\mathcal{P}_m$ 
4:   else
5:      $\mathcal{V}_m \leftarrow \mathcal{P}_m$ 
6:   end if
7: end for
8: for each class  $m = 1, \dots, M$  do
9:    $\mathbf{x}_{lm} \leftarrow$  Projection of  $\mathbf{x}_l$  onto  $\mathcal{V}_m$ 
10: end for
11:  $y(\mathbf{x}_l) \leftarrow C_j$  where  $j = \arg \max_m \mathbf{x}_{lm}$ 
12: Return:  $y(\mathbf{x}_l)$ 

```

---

### 6.1 Training dynamics of the linear probe

Here, we examine if a linear probe (i.e. VeLPIC) can decode latent generalization with performance comparable to MASC. To this end, we tracked the performance of VeLPIC, during training.

---

<sup>2</sup>We augment class training data points with their negative, so as to obtain a 1-D subspace, rather than a 1-D affine space, along the lines of the subspace construction procedure for MASC.

<sup>3</sup>Observe that this isn’t an issue with MASC, since it is a quadratic classifier.

<sup>4</sup>i.e.  $T_i$  is the mean of projecting training data points on  $\mathcal{P}_i$ .

<sup>5</sup>This is equivalent to minimum angle to the VeLPIC class vectors.

VeLPIC test accuracy during training for MLP-MNIST, CNN-Fashion-MNIST and AlexNet-Tiny ImageNet are shown in Figure 2. Results with 0% and 100% corruption degrees as well as the results with additional models are shown in Figure 6 and 7, respectively in Appendix Section D.

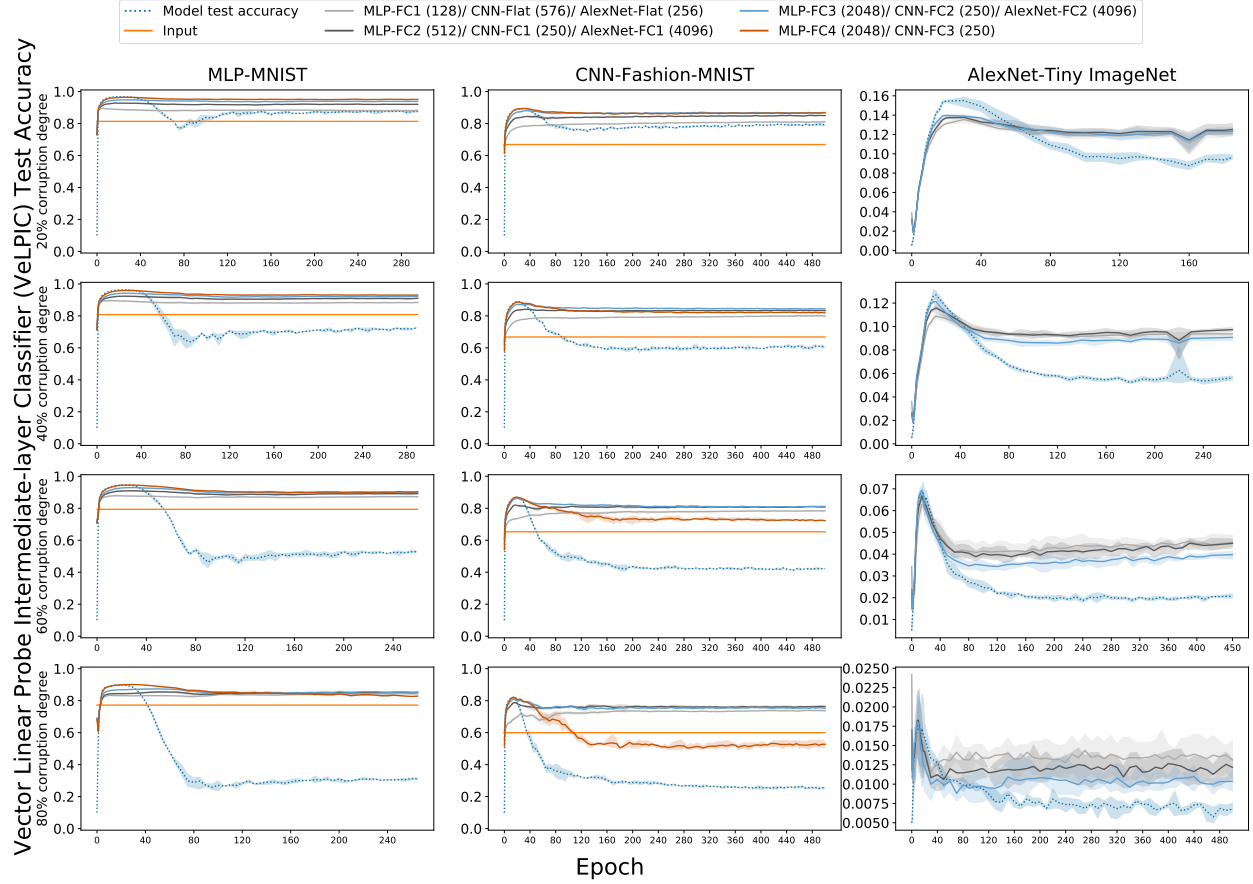


Figure 2: Vector Linear Probe Intermediate-layer Classifier (VeLPIC) test accuracy during training of the network, where test data is projected onto class vectors constructed at each epoch from training data with the indicated label corruption degrees. The plots display VeLPIC accuracy across different layers of the network for various model-dataset combinations. For reference, the test accuracy of the models (blue dotted line) over epochs of training is also shown. FC denotes fully connected layers with *ReLU* activation, and Flat refers to the flatten layer without *ReLU*.

Unexpectedly, not only does VeLPIC not perform worse than MASC, we find indeed that VeLPIC almost always significantly outperforms MASC, especially with higher corruption degrees. In fact, for representations from many layers, VeLPIC is able to extract significantly better latent generalization performance than MASC and our results show that, for these layers, VeLPIC’s performance plateaus at significantly higher levels than MASC. The difference between VeLPIC test accuracy and MASC test accuracy are shown in Figure 8, 9 & 10 in the Appendix Section D.1. This indicates that the probe’s quadratic nature is not essential in extracting latent generalization.

## 7 Transferring latent generalization to model generalization

Given that latent generalization often exceeds the models’ generalization ability, a natural question is whether this hidden latent generalization can be transferred directly to the model by appropriately modifying its weights. An affirmative answer to this question has implications that the phenomenon of latent generalization

is not mainly driven by the extraordinary effectiveness of the probe (MASC) used in prior work (Ketha & Ramaswamy, 2025). Conceptually, this question is significant because it demonstrates that the probe’s latent generalization is inherently accessible to the model using only the corrupted data it was trained on. Yet, the model during the act of standard training does not end up generalizing to the same extent, for reasons that we don’t understand well. Practically, it raises the possibility of repairing memorized model using the latent representations, without the cost of retraining from scratch.

Here, we ask if the latent generalization in models that memorize, can be directly transferred to the model, in order to immediately improve its generalization. To this end, it turns out that the class vectors of VeLPIC applied to the last layer can be directly substituted in the pre-softmax layer of the model as an intervention that transfers VeLPIC’s generalization performance to the model, without further training. We elaborate below on how this is so.

Consider a model whose last layer (i.e. the layer preceding the pre-softmax layer) consists of  $d$  units. Let  $\mathbf{v}_j \in \mathbb{R}^d$  be the VeLPIC class vector for class  $j$ . The new pre-softmax weight matrix  $\mathbf{W}_{\text{pre-softmax}} \in \mathbb{R}^{M \times d}$  is constructed as:

$$\mathbf{W}_{\text{pre-softmax}} = ([\mathbf{v}_1 \quad \mathbf{v}_2 \quad \cdots \quad \mathbf{v}_M])^\top \quad (3)$$

This weight matrix  $\mathbf{W}_{\text{pre-softmax}}$  replaces the original pre-softmax weights, and all biases are set to zero. It is straightforward to see that this substitution results in the model making the same predictions as VeLPIC applied to the last layer. While this substitution is fairly straightforward, it needs the linear probe, i.e. VeLPIC, for it to work. In particular, we note that it is unclear if transferability can happen with MASC alone and therefore the development of VeLPIC was an enabling factor in helping us answer this question.

During model training, we replace the pre-softmax weights with VeLPIC vectors, as indicated above and evaluate the model’s performance on the test dataset at each epoch. Figure 3 presents these results for MLP-MNIST, CNN-Fashion-MNIST, and AlexNet-Tiny ImageNet. Additional results, including models with 0% and 100% corruption levels, and other model-dataset pairs are presented in Figure 11 and Figure 12, respectively in the Appendix Section E.

We observe that the weight intervention that replaces pre-softmax weights with the VeLPIC vectors leads to an immediate & significant improvement in generalization performance in every epoch of the latter phase of training, matching that of the linear probe, & notably without any further training. This establishes that the latent generalization in memorized models can be directly harnessed to enhance their test performance, even in the presence of label noise.

## 8 Discussion

The notion of memorization, where deep networks are able to perfectly learn noisy data at the expense of generalization has posed a challenge to traditional notions of generalization from Statistical Learning Theory (Zhang et al., 2017; 2021). Recent work (Ketha & Ramaswamy, 2025) demonstrating improved latent generalization in such models is an interesting new development in our understanding of memorization and the nature of representations that drive it. Our goal here was to take a deeper dive into this phenomenon, to investigate the origin and dynamics of latent generalization. While the dynamics of memorization and generalization early in training have seen detailed empirical investigation (Arpit et al., 2017), the phenomenon of fall in model generalization in the later phase of training is more poorly understood. We showed that early-on in training, latent generalization and the model’s generalization closely follow each other, suggesting common mechanisms that contribute to both. However, later in training, there is a divergence, with the model retaining significant latent generalization ability, while sacrificing overt model generalization to a greater degree. After showing that MASC (Ketha & Ramaswamy, 2025) is a quadratic classifier, we built a new linear probe (VeLPIC) and found, unexpectedly, that it has better latent generalization performance in comparison to MASC, in most cases. Indeed, while (Ketha & Ramaswamy, 2025) show that MASC applied to at least one layer outperforms the model at the end of training, with respect to generalization, with VeLPIC, we find that, in most cases, all layers’ latent generalization outperform model generalization. This implies that the latent generalization effect during memorization is more pronounced and more widely present among layer representations than previously reported in (Ketha & Ramaswamy, 2025). We were



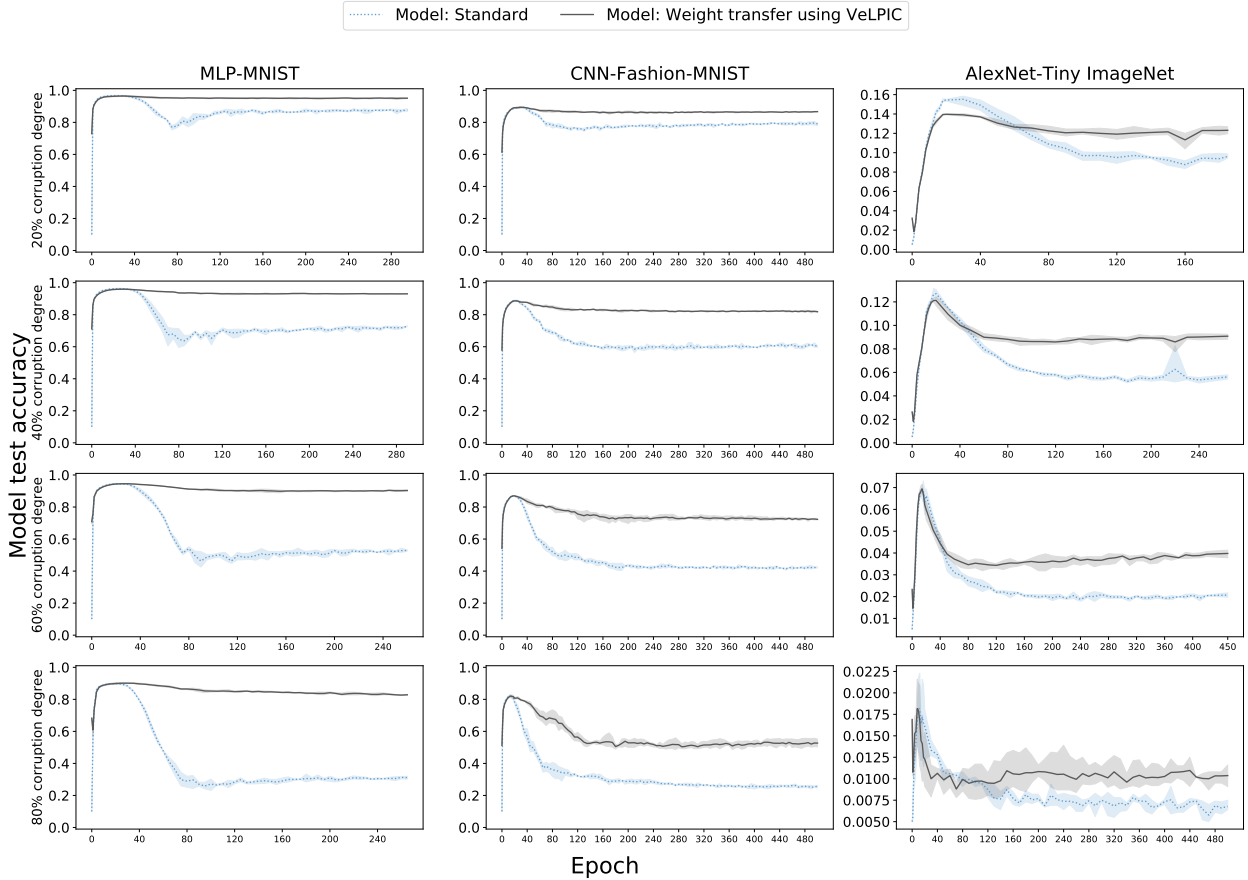


Figure 3: Model test accuracy when the weight intervention is applied to the epoch in question during training. The test accuracy of the model with standard training without weight intervention (blue dotted line) is overlaid for comparison.

also interested in examining if the latent generalization could readily be translated to model generalization by directly modifying model weights. We utilized the linear probe to derive a new set of model pre-softmax weights to make this so.

This work brings up multiple new directions for investigation. While we have made some progress, the detailed mechanisms governing latent generalization during memorization remain to be investigated. It is also an open question, whether there exist other probes that can extract better latent generalization from layerwise representations, in comparison to MASC and VeLPIC. Next, it is unclear if latent generalization from representations of layers other than the last layer can be transferred towards model generalization. This can be useful to do, in cases where early or middle layers exhibit better latent generalization than the last layer. More generally, in light of these results, whether an understanding of generalization in the memorization regime can inform a better understanding of generalization for models trained with uncorrupted labels is a worthwhile direction for future investigation.

In closing, our results highlight the rich role of representations in driving generalization during memorization and how their understanding can be utilized to directly improve model generalization.

## References

- Guillaume Alain and Yoshua Bengio. Understanding intermediate layers using linear classifier probes. 2018. *arXiv preprint arXiv:1610.01644*, 2018.
- Devansh Arpit, Stanisław Jastrzębski, Nicolas Ballas, David Krueger, Emmanuel Bengio, Maxinder S Kanwal, Tegan Maharaj, Asja Fischer, Aaron Courville, Yoshua Bengio, et al. A closer look at memorization in deep networks. In *International conference on machine learning*, pp. 233–242. PMLR, 2017.
- Amr Bakry, Mohamed Elhoseiny, Tarek El-Gaaly, and Ahmed Elgammal. Digging deep into the layers of cnns: In search of how cnns achieve view invariance. *arXiv preprint arXiv:1508.01983*, 2015.
- N Alex Cayco-Gajic and R Angus Silver. Re-evaluating circuit mechanisms underlying pattern separation. *Neuron*, 101(4):584–602, 2019.
- SueYeon Chung, Daniel D Lee, and Haim Sompolsky. Linear readout of object manifolds. *Physical Review E*, 93(6):060301, 2016.
- Uri Cohen, SueYeon Chung, Daniel D Lee, and Haim Sompolsky. Separability and geometry of object manifolds in deep neural networks. *Nature communications*, 11(1):746, 2020.
- Li Deng. The mnist database of handwritten digit images for machine learning research. *IEEE Signal Processing Magazine*, 29(6):141–142, 2012.
- Matthew S Farrell, Stefano Recanatesi, Guillaume Lajoie, and Eric Shea-Brown. Dynamic compression and expansion in a classifying recurrent network. *bioRxiv*, pp. 564476, 2019.
- Bo Han, Quanming Yao, Xingrui Yu, Gang Niu, Miao Xu, Weihua Hu, Ivor Tsang, and Masashi Sugiyama. Co-teaching: Robust training of deep neural networks with extremely noisy labels. *Advances in neural information processing systems*, 31, 2018.
- Lu Jiang, Zhengyuan Zhou, Thomas Leung, Li-Jia Li, and Li Fei-Fei. Mentornet: Learning data-driven curriculum for very deep neural networks on corrupted labels. In *International conference on machine learning*, pp. 2304–2313. PMLR, 2018.
- Simran Ketha and Venkatakrishnan Ramaswamy. Decoding generalization from memorization in deep neural networks. *arXiv preprint arXiv:2501.14687*, 2025.
- Diederik P Kingma. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- Simon Kornblith, Mohammad Norouzi, Honglak Lee, and Geoffrey Hinton. Similarity of neural network representations revisited. In *International conference on machine learning*, pp. 3519–3529. PMLR, 2019.
- Alex Krizhevsky. Learning multiple layers of features from tiny images. Technical report, 2009.
- Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. *Advances in neural information processing systems*, 25, 2012.
- Andrew K Lampinen and Surya Ganguli. An analytic theory of generalization dynamics and transfer learning in deep linear networks. *arXiv preprint arXiv:1809.10374*, 2018.
- Mingchen Li, Mahdi Soltanolkotabi, and Samet Oymak. Gradient descent with early stopping is provably robust to label noise for overparameterized neural networks. In *International conference on artificial intelligence and statistics*, pp. 4313–4324. PMLR, 2020.
- Sheng Liu, Jonathan Niles-Weed, Narges Razavian, and Carlos Fernandez-Granda. Early-learning regularization prevents memorization of noisy labels. *Advances in neural information processing systems*, 33: 20331–20342, 2020.
- Grégoire Montavon, Mikio L Braun, and Klaus-Robert Müller. Kernel analysis of deep networks. *Journal of Machine Learning Research*, 12(9), 2011.

- Ari Morcos, Maithra Raghu, and Samy Bengio. Insights on representational similarity in neural networks with canonical correlation. *Advances in neural information processing systems*, 31, 2018.
- Mohammed Ali Moustafa. Tiny imagenet, 2017. URL <https://kaggle.com/competitions/tiny-imagenet>.
- Ning Qian. On the momentum term in gradient descent learning algorithms. *Neural networks*, 12(1):145–151, 1999.
- Maithra Raghu, Justin Gilmer, Jason Yosinski, and Jascha Sohl-Dickstein. Svcca: Singular vector canonical correlation analysis for deep learning dynamics and interpretability. *Advances in neural information processing systems*, 30, 2017.
- Andrew M Saxe, James L McClelland, and Surya Ganguli. Exact solutions to the nonlinear dynamics of learning in deep linear neural networks. *arXiv preprint arXiv:1312.6120*, 2013.
- Cory Stephenson, Abhinav Ganesh, Yue Hui, Hanlin Tang, SueYeon Chung, et al. On the geometry of generalization and memorization in deep neural networks. In *International Conference on Learning Representations*, 2021.
- David Sussillo and Larry F Abbott. Generating coherent patterns of activity from chaotic neural networks. *Neuron*, 63(4):544–557, 2009.
- Lina M Tran, Adam Santoro, Lulu Liu, Sheena A Josselyn, Blake A Richards, and Paul W Frankland. Adult neurogenesis acts as a neural regularizer. *Proceedings of the National Academy of Sciences*, 119(45): e2206704119, 2022.
- Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms, 2017.
- Jason Yosinski, Jeff Clune, Yoshua Bengio, and Hod Lipson. How transferable are features in deep neural networks? *Advances in neural information processing systems*, 27, 2014.
- Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. Understanding deep learning requires rethinking generalization. iclr 2017. *arXiv preprint arXiv:1611.03530*, 2017.
- Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. Understanding deep learning (still) requires rethinking generalization. *Communications of the ACM*, 64(3):107–115, 2021.

## A Model architectures and training details

**MLP Model.** The MLP architecture consists of four hidden layers with 128, 512, 2048, and 2048 units, respectively. Each layer is followed by a *ReLU* activation, and a *softmax* layer is used for classification. Models were trained SGD Qian (1999) with a learning rate of  $1 \times 10^{-3}$  and momentum 0.9. A batch size of 32 was used across all experiments. Input dataset was normalized by dividing pixel values by 255.

**CNN Model.** The CNN model<sup>6</sup> is composed of three convolutional blocks, each containing two convolutional layers followed by a max pooling layer. The convolutional layers use 16, 32, and 64 filters, respectively, with kernel size  $3 \times 3$  and stride 1. The max pooling layers have a kernel size of  $2 \times 2$  and stride 1. These blocks are followed by three fully connected layers with 250 units each. ReLU activation is used after all layers except pooling, and softmax is used at the output for classification. The CNN was trained using Adam optimizer Kingma (2014) with a learning rate of 0.0002. For MNIST and Fashion-MNIST, a batch size of 32 was used, while for CIFAR-10, a batch size of 128 was used. Input data was normalized by subtracting the mean and dividing by the standard deviation of each channel.

The experiments were conducted on servers and workstations equipped with NVIDIA GeForce RTX 3080, RTX 3090, Tesla V100, and Tesla A100 GPUs. The server runs on Rocky Linux 8.10 (Green Obsidian), while the workstation uses Ubuntu 20.04.3 LTS. Memory requirements varied depending on the specific experiments and models. All model implementations were developed in Python using the PyTorch library, with `torch.manual_seed` set to 42 to ensure reproducibility. Accuracy served as the primary evaluation metric throughout this work.

## B Minimum Angle Subspace Classifier (MASC)

We summarize below the Minimum Angle Subspace Classifier (MASC) from (Ketha & Ramaswamy, 2025), in order to keep the exposition here largely self-contained.

For a given deep network, MASC leverages the class-specific geometric structure of network’s latent representations. For an input data point  $\mathbf{x}$ , let its activation vector at layer  $l$  be denoted by  $\mathbf{x}_l$ . The objective is to classify  $\mathbf{x}_l$  by leveraging a set of class-conditional subspaces,  $\{S_k\}_{k=1}^K$ , estimated from a training dataset  $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^m$ . To predict the class label  $y(\mathbf{x}_l)$ , MASC Algorithm 2 (reproduced verbatim from (Ketha & Ramaswamy, 2025)), assigns  $\mathbf{x}_l$  to the class whose training subspace forms the smallest angle with it.

The class-conditional subspaces  $\{S_k\}_{k=1}^K$  are estimated from the training dataset  $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^m$ , where each  $\mathbf{x}_i \in \mathbb{R}^d$  is paired with a label  $y_i \in \{C_k\}_{k=1}^K$ . For a given layer  $l$ , these subspaces are constructed following Algorithms 3 and 4 (reproduced verbatim from (Ketha & Ramaswamy, 2025)). In practice, each subspace  $S_k$  is represented by its principal components, which provide a compact basis for capturing the underlying class-conditional structure.

---

**Algorithm 2 Minimum Angle Subspace Classifier (MASC)** (reproduced verbatim from (Ketha & Ramaswamy, 2025))

---

- 1: **Input:** Training subspaces  $\{S_k\}_{k=1}^K$ , layer output data point  $\mathbf{x}_l$  from layer  $l$  when input  $\mathbf{x}$  is passed through the network and classes  $\{C_k\}_{k=1}^K$ .
  - 2: **Output:** MASC prediction class label  $y(\mathbf{x}_l)$  according to layer  $l$ .
  - 3: **for** each class  $C_k$  **do**
  - 4:    $\mathbf{x}_{lk} \leftarrow$  compute the projection of  $\mathbf{x}_l$  onto subspace  $S_k$ .
  - 5:   Compute the angle  $\theta(\mathbf{x}_l, \mathbf{x}_{lk})$  between  $\mathbf{x}_l$  and  $\mathbf{x}_{lk}$
  - 6: **end for**
  - 7: Assign the label  $y(\mathbf{x}_l) = C_k$  where  $k = \arg \min_k \theta(\mathbf{x}_l, \mathbf{x}_{lk})$
- 

<sup>6</sup>The convolution network were implemented following the design principals outlined in (Tran et al., 2022).

**Algorithm 3 Subspaces Estimator for MASC**

(reproduced verbatim from (Ketha &amp; Ramaswamy, 2025))

- 
- 1: **Input:** Training dataset  $\mathcal{D}\{(\mathbf{x}_i, y_i)\}_{i=1}^m \in \mathbb{R}^d \times \mathbb{R}$ , where each  $\mathbf{x}_i \in \mathbb{R}^d$  and  $y_i \in \{C_k\}_{k=1}^K$  are input-label pairs, neural network, and layer  $l$ .
  - 2: **Output:** Subspaces  $\{S_k\}_{k=1}^K$  for classes  $K$  and given layer  $l$ .
  - 3:  $\mathcal{D}_l = \phi$
  - 4: **for** each input pair  $(\mathbf{x}_i, y_i)$  in  $\mathcal{D}$  **do**
  - 5:   Pass  $\mathbf{x}_i$  through the network layers to obtain the output of layer  $l$ , denoted as  $\mathbf{x}_l \in \mathbb{R}^{ld}$ .
  - 6:    $\mathcal{D}_l = \mathcal{D}_l \cup \{\mathbf{x}_l\}$
  - 7: **end for**
  - 8: Estimated subspaces  $\{S_k\}_{k=1}^K \leftarrow \text{PCA-Based Subspace Estimation}(\mathcal{D}_l)$
  - 9: **Return:** Subspaces  $\{S_k\}_{k=1}^K$
- 

**Algorithm 4 PCA-Based Subspace Estimation**

(reproduced verbatim from (Ketha &amp; Ramaswamy, 2025))

- 
- 1: **Input:** Layer output  $\mathcal{D}_l = \{(\mathbf{x}_i, y_i)\}_{i=1}^m$ , where  $\mathbf{x}_l \in \mathbb{R}^{ld}$  and  $y_i \in \{C_k\}_{k=1}^K$ .
  - 2: **Output:** Subspaces  $\{S_k\}_{k=1}^K$  for classes  $K$ .
  - 3:  $\mathcal{D}_{\text{new}} \leftarrow \mathcal{D}_l$
  - 4: **for** each data point  $\mathbf{x}_l$  in  $\mathcal{D}_l$  **do**
  - 5:    $\mathcal{D}_{\text{new}} \leftarrow \mathcal{D}_{\text{new}} \cup \{-\mathbf{x}_l\}$
  - 6: **end for**
  - 7: **for** each class  $C_k$  in  $C_K$  **do**
  - 8:   Extract the subset of data  $\mathcal{D}_{\text{new},k} = \{\mathbf{x}_l \mid y_i = k\}$
  - 9:   Apply PCA to  $\mathcal{D}_{\text{new},k}$  to calculate the PCA components
  - 10:   The span of the PCA components defines the subspace  $S_k$
  - 11: **end for**
  - 12: **Return:** Subspaces  $\{S_k\}_{k=1}^K$
-

## C Training dynamics of latent generalization using MASC

MASC testing accuracy during training for MLP trained on MNIST, CNN trained on Fashion-MNIST and AlexNet trained on Tiny ImageNet with 0% and 100% corruption degrees are shown in Figure 4. In the absence of label corruption, the MASC accuracy of the final fully connected layers (MLP-FC4 (2048 units) and CNN-FC3 (250 units)) closely matched the corresponding model test accuracies. Interestingly, in certain cases, such as AlexNet trained on Tiny ImageNet (FC1 and FC2, each with 4096 units) and MLP-CIFAR-10 (FC3 with 2048 units), the MASC accuracy even surpassed the model’s test accuracy. Results with additional models i.e. MLP trained on CIFAR-10, CNN trained on MNIST, CNN trained on CIFAR-10 for various corruption degrees are shown in Figure 5.

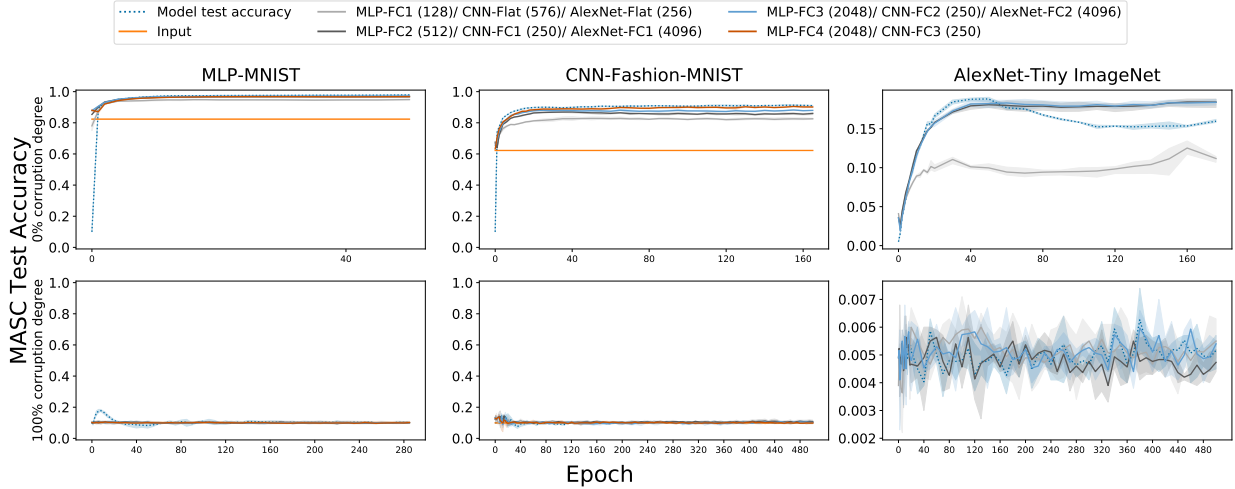


Figure 4: Minimum Angle Subspace Classifier (MASC) Test accuracy for 0% and 100% corruption degrees during training of the network, where test data is projected onto class-specific subspaces constructed from training data with the indicated label corruption degrees. The plots display MASC accuracy across different layers of the network for various model–dataset combinations. For reference, the test accuracy of the models (dotted line) is also shown. Each row corresponds to a specific corruption degree, while columns represent different models, as labeled. FC denotes fully connected layers with *ReLU* activation, and Flat refers to the flatten layer without *ReLU*.

## D Training dynamics of the linear probe: VeLPIC

A linear probe – VeLPIC – test accuracy during training for MLP-MNIST, CNN-Fashion-MNIST and AlexNet-Tiny ImageNet with 0% and 100% corruption degrees are shown in Figure 6. Results with additional models i.e. MLP-CIFAR-10, CNN-MNIST, CNN-CIFAR-10 for various corruption degrees are shown in Figure 7.

### D.1 Difference between VeLPIC and MASC

Here, we present the difference between test accuracy of VeLPIC and MASC during training and for different layer of the networks. For MLP-MNIST, CNN-Fashion-MNIST and AlexNet-Tiny ImageNet, these results are shown in Figure 8 and Figure 9. Results with additional models i.e. MLP-CIFAR-10, CNN-MNIST, CNN-CIFAR-10 for various corruption degrees are shown in Figure 10.

## E Transferring latent generalization to model generalization

For MLP-MNIST, CNN-Fashion-MNIST and AlexNet-Tiny ImageNet with 0% and 100% corruption degrees, model test accuracy during training when we replace the pre-softmax weights with VeLPIC vectors are

shown in Figure 11. Results with additional models i.e. MLP-CIFAR-10, CNN-MNIST, CNN-CIFAR-10 for various corruption degrees are shown in Figure 12. Model corrupted training accuracy for all models-dataset-corruption are plotted in Figure 13 and Figure 14.

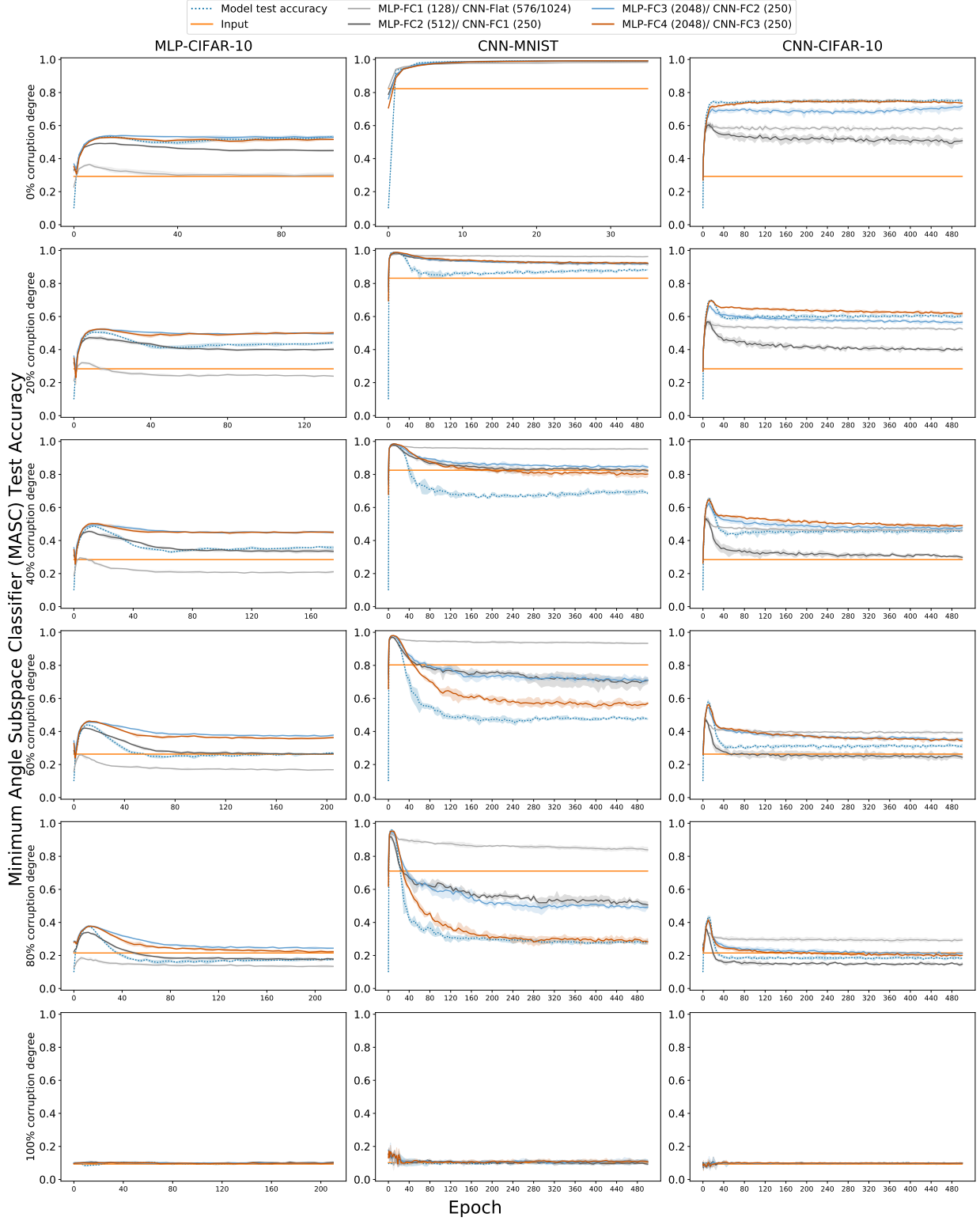


Figure 5: MASC accuracy during training of the network, where test data is projected onto class-specific subspaces constructed from training data with the indicated label corruption degrees. The plots display MASC accuracy across different layers of the network for various model-dataset combinations. For reference, the test accuracy of the models (dotted line) is also shown. Each row corresponds to a specific corruption degree, while columns represent different layer models, as labeled. FC denotes fully connected layers with *ReLU* activation, and Flat refers to the flatten layer without *ReLU*.



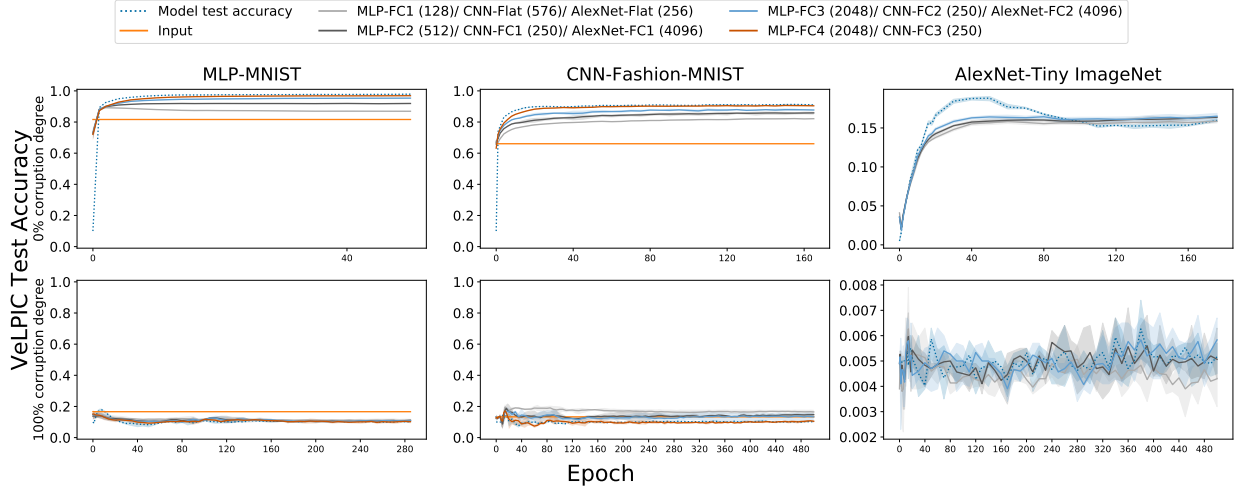


Figure 6: Vector Linear Probe Intermediate-layer Classifier (VeLPIC) test accuracy for 0% and 100% corruption degrees during training of the network, where test data is projected onto class vectors constructed at each epoch from training data with the indicated label corruption degrees. The plots display VeLPIC accuracy across different layers of the network for various model-dataset combinations. For reference, the test accuracy of the models (blue dotted line) over epochs of training is also shown. FC denotes fully connected layers with *ReLU* activation, and Flat refers to the flatten layer without *ReLU*.

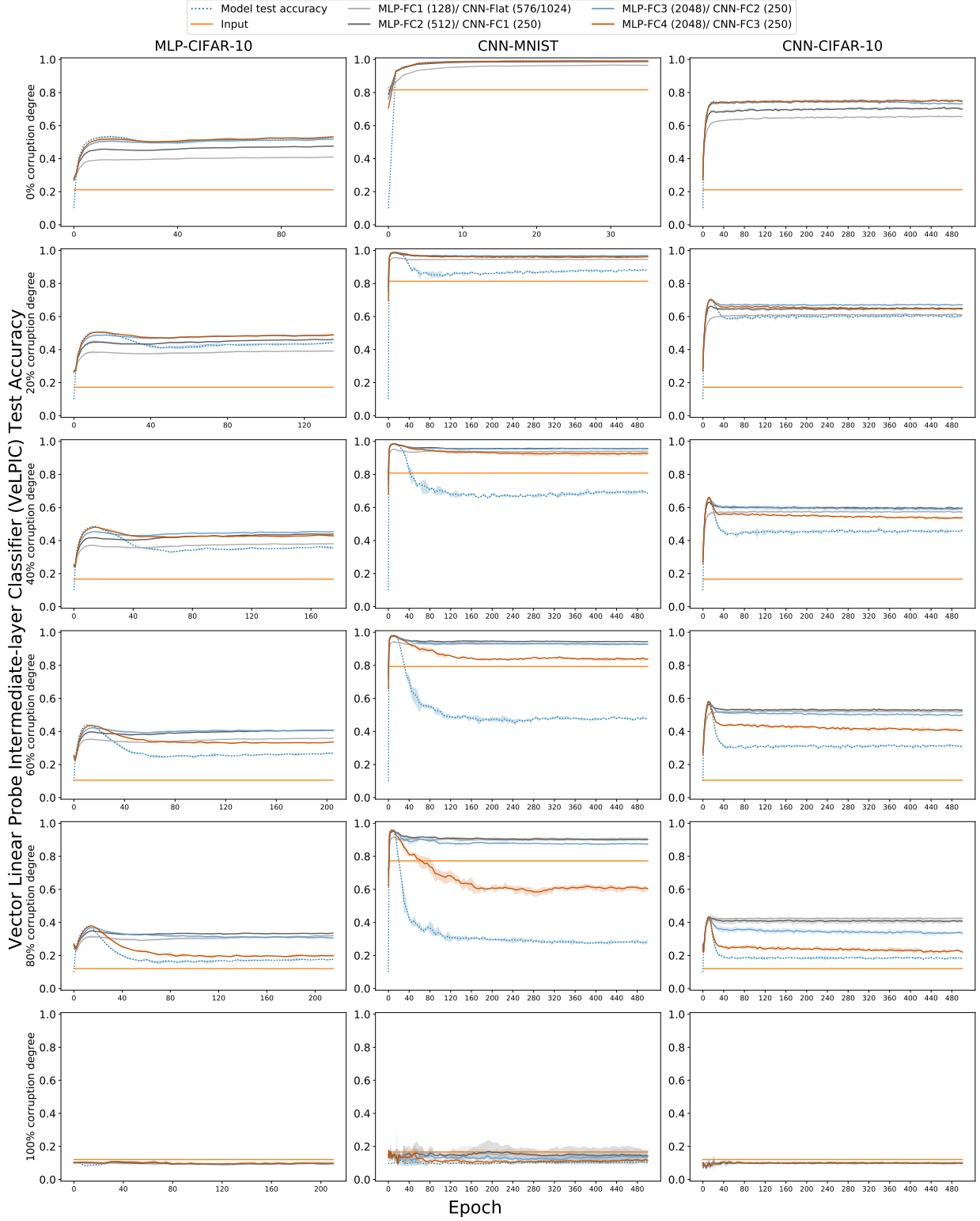


Figure 7: Vector Linear Probe Intermediate-layer Classifier (VeLPIC) test accuracy during training of the network, where test data is projected onto class vectors constructed at each epoch from training data with the indicated label corruption degrees. The plots display VeLPIC accuracy across different layers of the network for various model-dataset combinations. For reference, the test accuracy of the models (blue dotted line) over epochs of training is also shown. FC denotes fully connected layers with *ReLU* activation, and Flat refers to the flatten layer without *ReLU*.

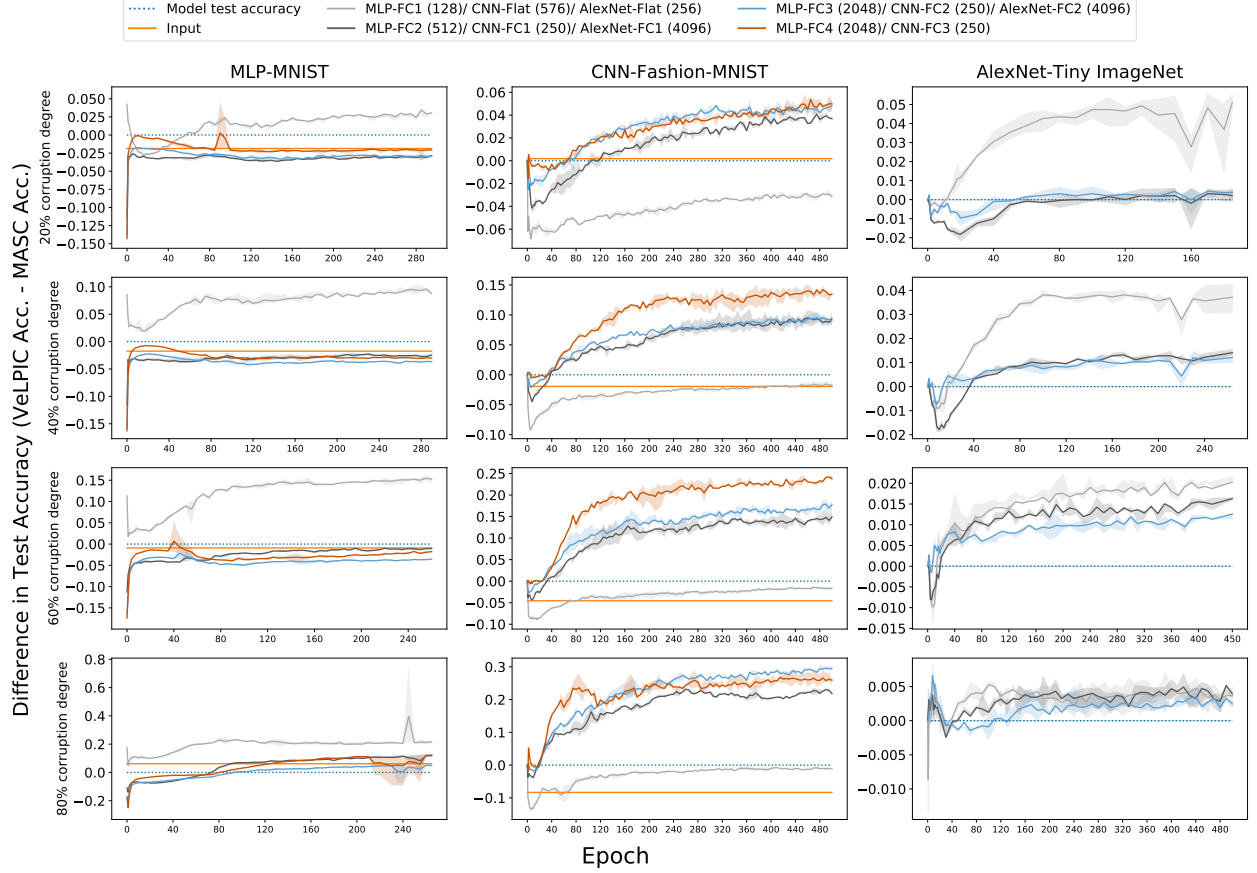


Figure 8: Difference in test accuracy (VeLPIC Accuracy - MASC Accuracy) during training of the network, where test data is projected onto class vectors constructed at each epoch from training data with the indicated label corruption degrees. The plots display difference in accuracy across different layers of the network for various model-dataset combinations. For reference, the test accuracy of the models (blue dotted line) over epochs of training is also shown, which would be 0.

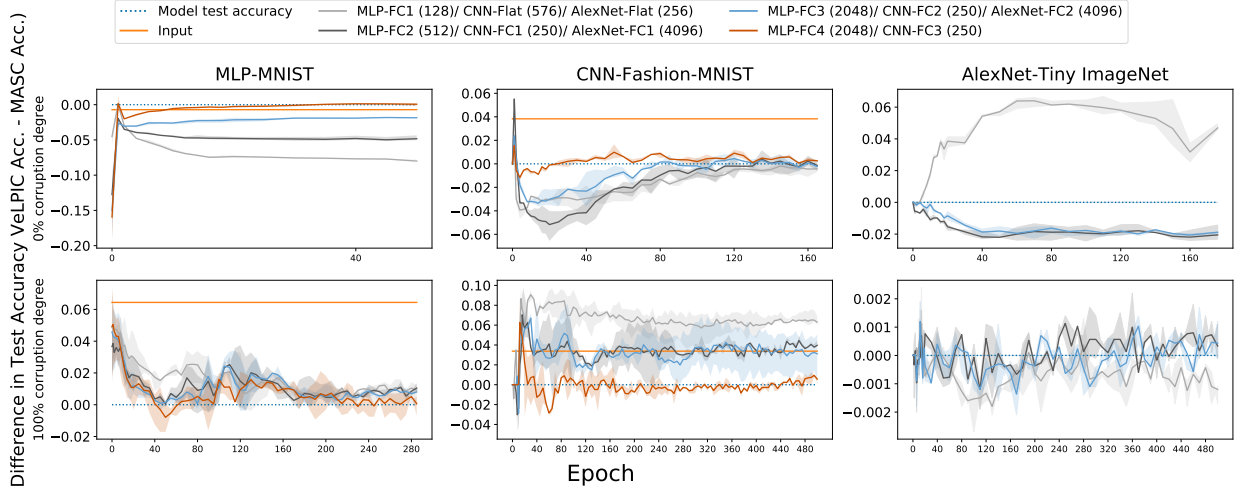


Figure 9: Difference in test accuracy (VeLPIC Accuracy - MASC Accuracy) during training of the network, where test data is projected onto class vectors constructed at each epoch from training data with the indicated label corruption degrees. The plots display difference in accuracy across different layers of the network for various model-dataset combinations. For reference, the test accuracy of the models (blue dotted line) over epochs of training is also shown, which would be 0.

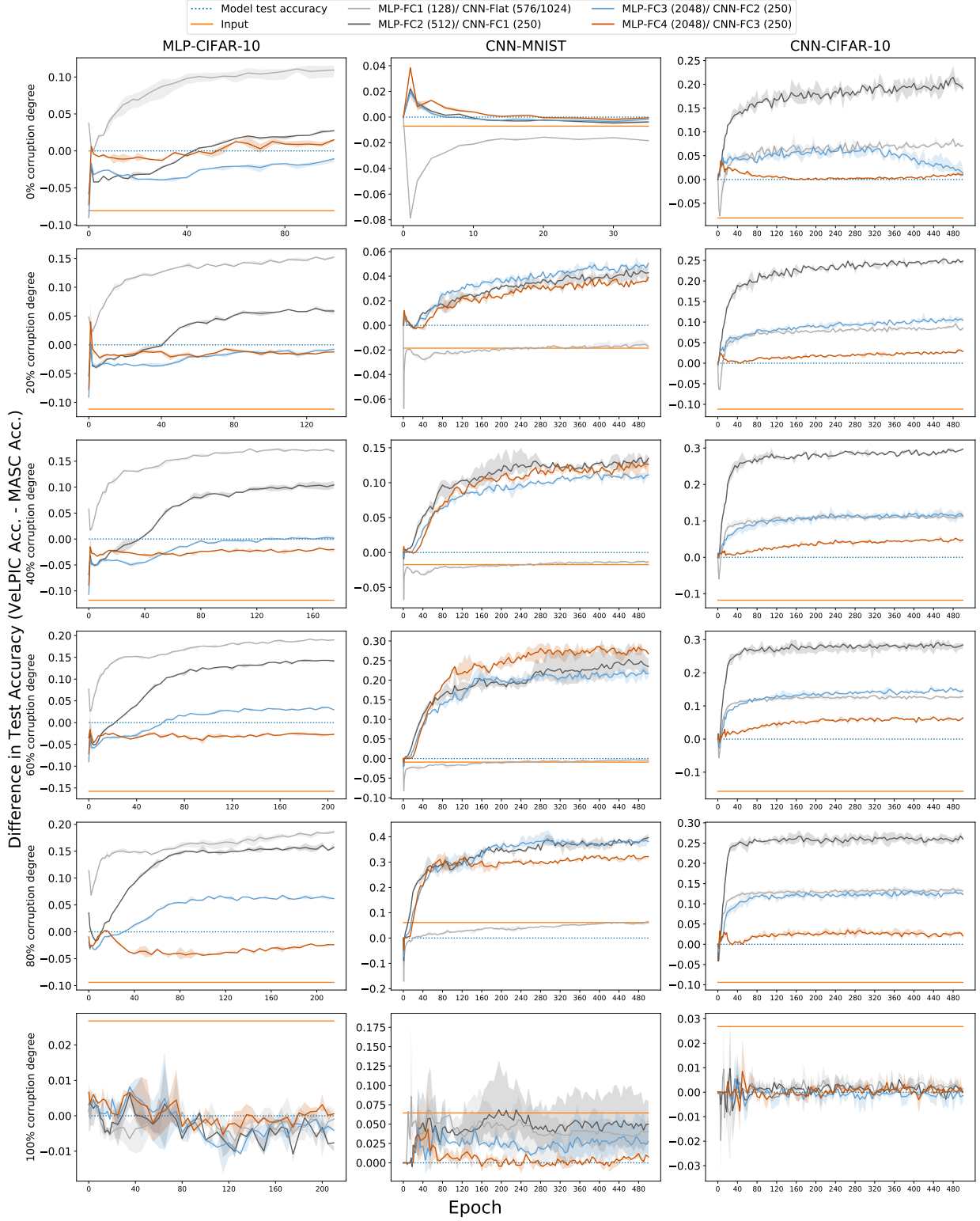


Figure 10: Difference in test accuracy (VeLPIC Accuracy - MASC Accuracy) during training of the network, where test data is projected onto class vectors constructed at each epoch from training data with the indicated label corruption degrees. The plots display difference in accuracy across different layers of the network for various model-dataset combinations. For reference, the test accuracy of the models (blue dotted line) over epochs of training is also shown, which would be 0.

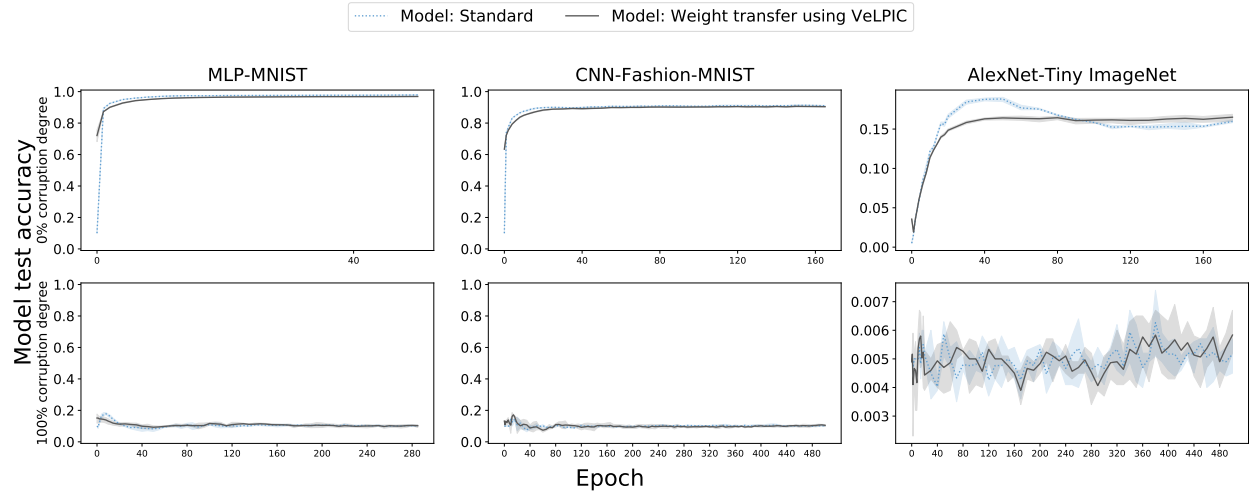


Figure 11: Comparing model test accuracy with VeLPIC transferred accuracy when the weight intervention is applied to the model at the epoch in question during training for corruption degrees 0% and 100%. The test accuracy of the model with standard training without weight intervention (blue dotted line) is overlaid for comparison.

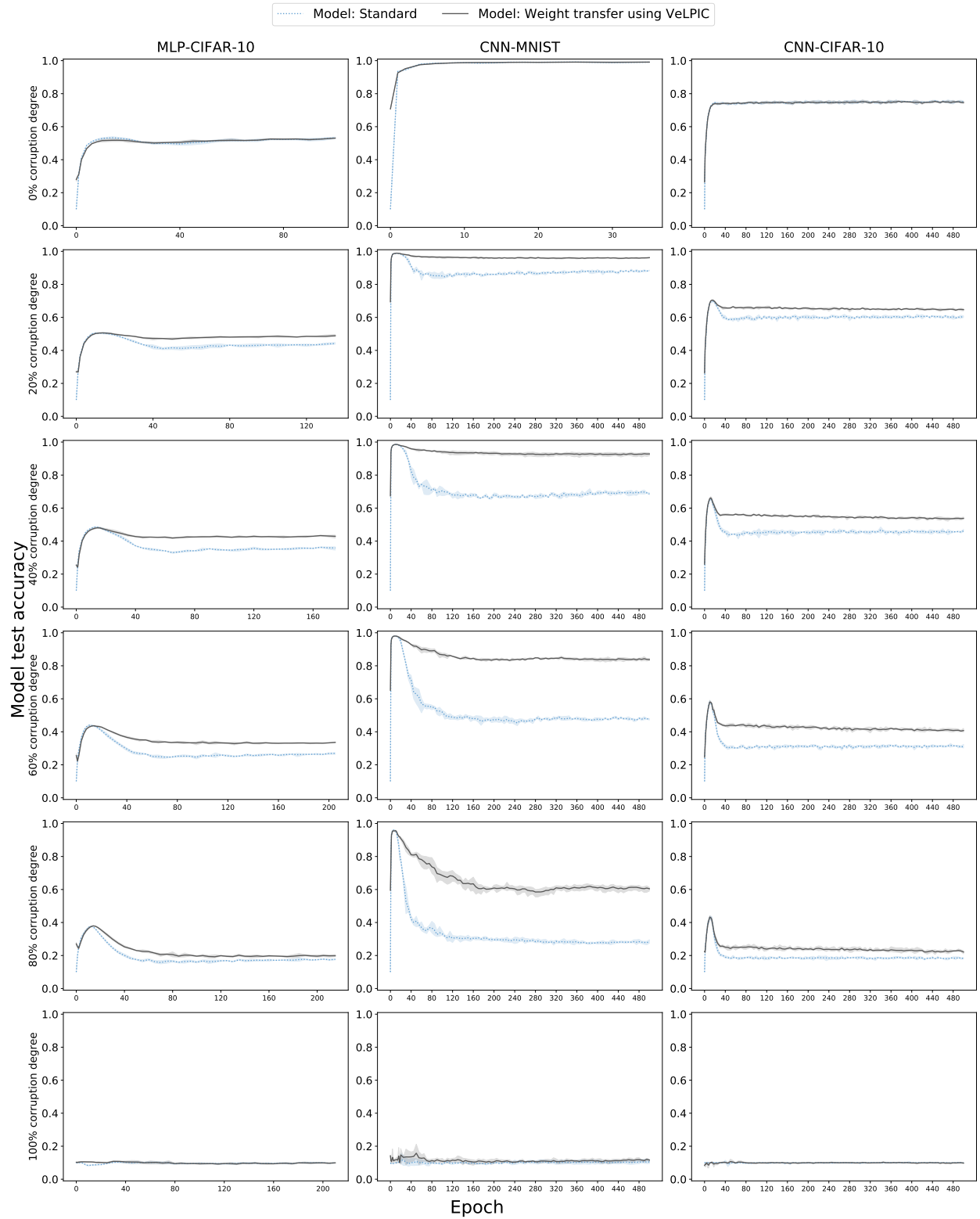


Figure 12: Comparing model test accuracy with VeLPIC transferred accuracy when the weight intervention is applied to the epoch in question during training. The test accuracy of the model with standard training without weight intervention (blue dotted line) is overlaid for comparison.

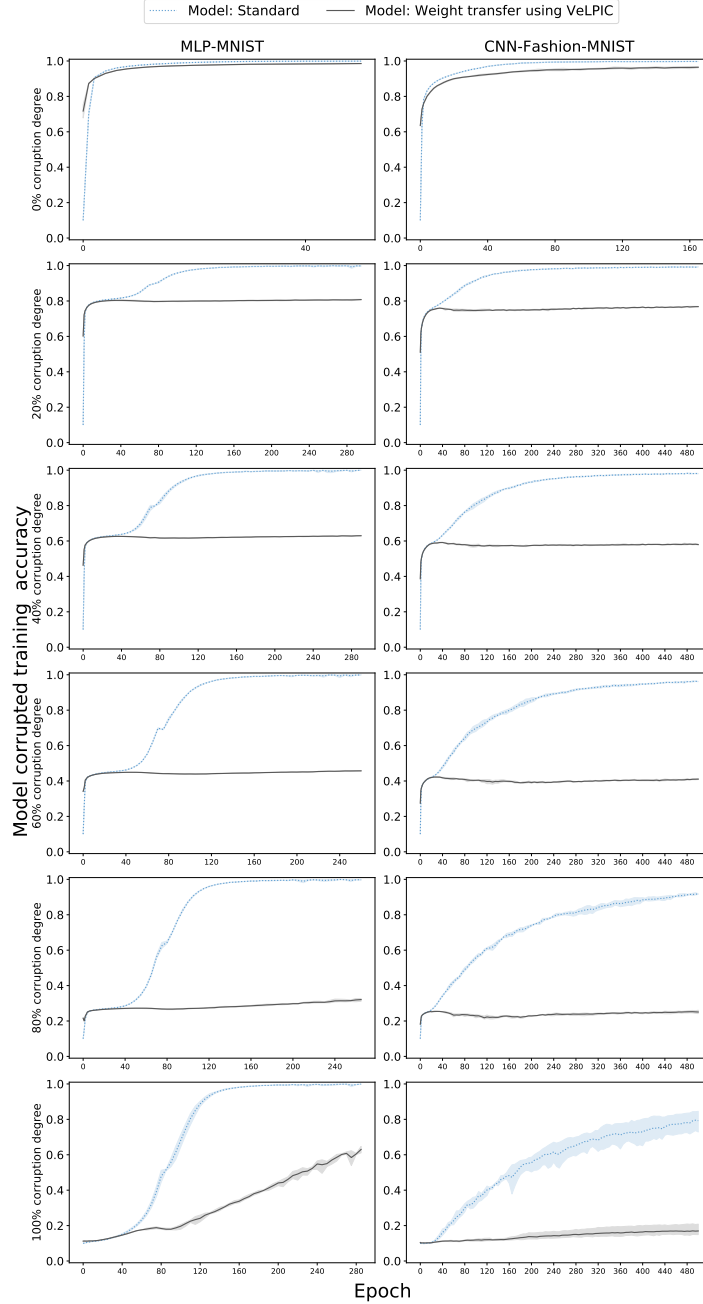


Figure 13: Model train accuracy on corrupted dataset when the VeLPIC weight intervention is applied to pre-softmax weights at the epoch in question during training. The training accuracy on corrupted dataset of the model with standard training without weight intervention (blue dotted line) is overlaid for comparison. Observe that, except for 100% corruption degree, the transferred training accuracy tends to saturate at a level largely consistent with the fraction of true training labels in the corrupted dataset.



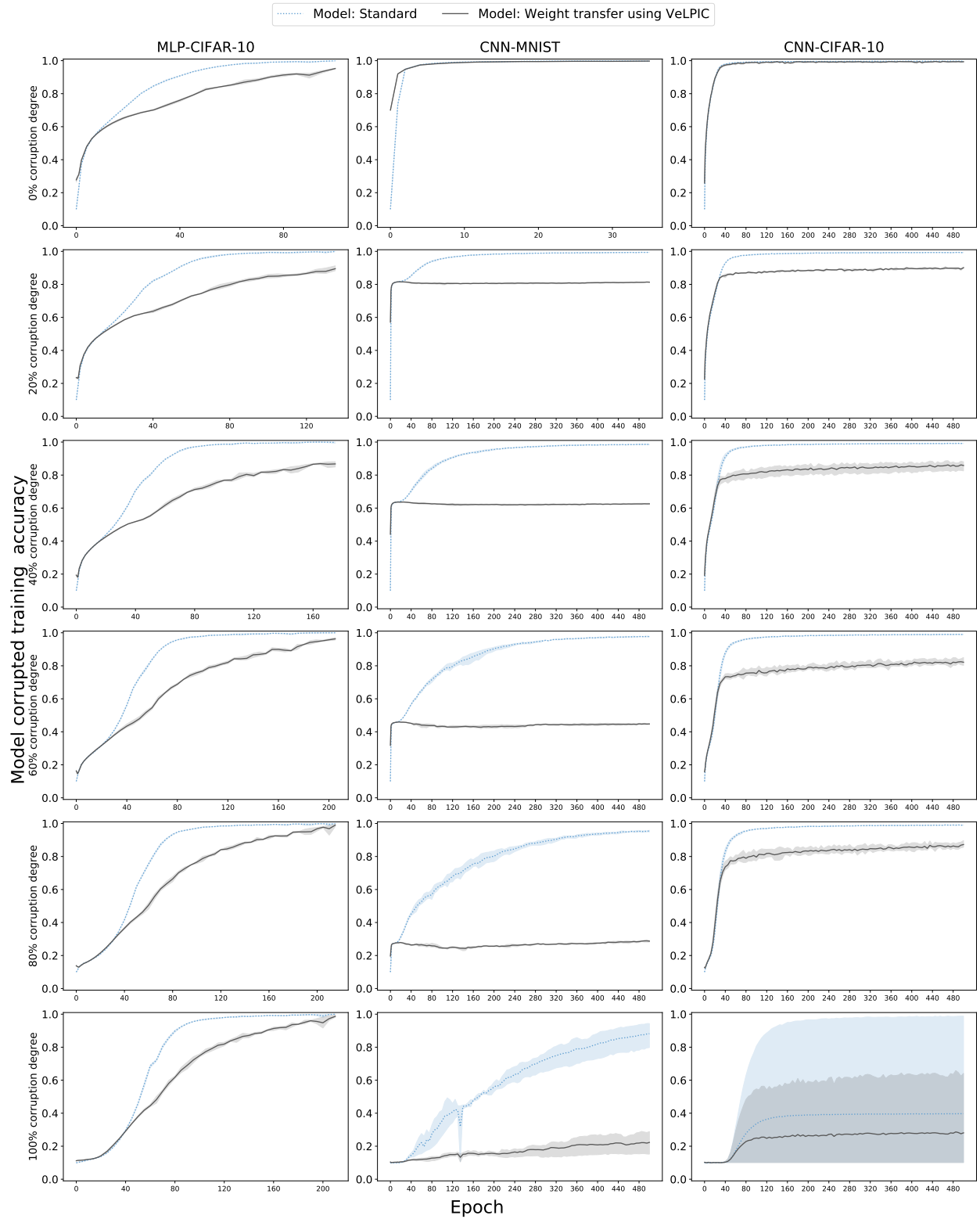


Figure 14: Model train accuracy on corrupted dataset when the VeLPIC weight intervention is applied to pre-softmax weights at the epoch in question during training. The training accuracy on corrupted dataset of the model with standard training without weight intervention (blue dotted line) is overlaid for comparison.