

# Characterizing Datapoints via Second-Split Forgetting

Pratyush Maini<sup>1</sup>   Saurabh Garg<sup>1</sup>   Zachary C. Lipton<sup>1</sup>   J. Zico Kolter<sup>1,2</sup>  
 Carnegie Mellon University<sup>1</sup>   Bosch Center for AI<sup>2</sup>  
 {pratyushmaini,zlipton}@cmu.edu; {sgarg2, zkolter}@cs.cmu.edu

## Abstract

Researchers investigating example hardness have increasingly focused on the dynamics by which neural networks learn and forget examples throughout training. Popular metrics derived from these dynamics include (i) the epoch at which examples are first correctly classified; (ii) the number of times their predictions flip during training; and (iii) whether their prediction flips if they are held out. However, these metrics do not distinguish among examples that are hard for distinct reasons, such as membership in a rare subpopulation, being mislabeled, or belonging to a complex subpopulation. In this paper, we propose *second-split forgetting time* (SSFT), a complementary metric that tracks the epoch (if any) after which an original training example is forgotten as the network is fine-tuned on a randomly held out partition of the data. Across multiple benchmark datasets and modalities, we demonstrate that *mislabeled* examples are forgotten quickly, and *seemingly rare* examples are forgotten comparatively slowly. By contrast, metrics only considering the first split learning dynamics struggle to differentiate the two. At large learning rates, SSFT tends to be robust across architectures, optimizers, and random seeds. From a practical standpoint, the SSFT can (i) help to identify mislabeled samples, the removal of which improves generalization; and (ii) provide insights about failure modes. Through theoretical analysis addressing overparameterized linear models, we provide insights into how the observed phenomena may arise.<sup>1</sup>

## 1 Introduction

A growing literature has investigated metrics for characterizing the difficulty of training examples, driven by such diverse motivations as (i) deriving insights for how to reconcile the ability of deep neural networks to generalize [30] with their ability to memorize noise [15, 48]; (ii) identifying potentially mislabeled examples; and (iii) identifying notably challenging or rare sub-populations of examples. Some of these efforts have turned towards learning dynamics, with researchers noting that neural networks tend to learn cleanly labeled examples before mislabeled examples [17, 18, 33], and more generally tend to learn *simpler* patterns sooner—for several intuitive notions of simplicity [19, 35, 43]. Broadly, works in this area tend to characterize examples as belonging either to *prototypical groups* or *memorized exceptions* [7, 16, 25]. Adapting these intuitions to real datasets, Feldman [15] propose rating the degree to which an example is memorized based on whether its predicted class flips when it is excluded from the training set. These, and other works [8, 21, 35, 43, 47] have proposed many metrics for characterizing example difficulty with Carlini et al. [7] comparing five such metrics. However, while many of these works distinguish some notion of *easy* versus *hard* samples, they seldom (i) offer tools for distinguishing among different types of hard examples; (ii) explain theoretically why these metrics might be useful for distinguishing easy versus hard samples. Moreover, existing metrics tend to give similar scores to examples that are difficult for distinct reasons, e.g. membership in rare, complex, or mislabeled sub-populations.

<sup>1</sup>Code for reproducing our experiments can be found at <https://github.com/pratyushmaini/ssft>.

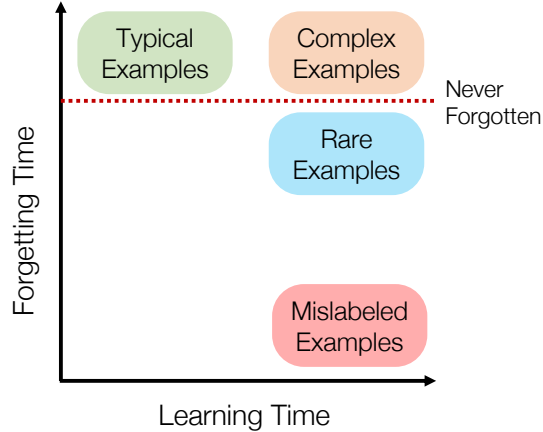


Figure 1: Overview of example separation offered by the unified view of learning and forgetting time.

In this paper, we propose to additionally consider a new metric, Second-Split Forgetting Time (SSFT), calculated based on the forgetting dynamics that arise as training examples are forgotten when a neural network continues to train on a second, randomly held out data partition. SSFT is defined as the fine-tuning epoch after which a first-split training example is no longer classified correctly. We find that SSFT identifies mislabeled examples remarkably well but does little to separate out under- versus over-represented subpopulations. Conversely, metrics based on the (first-split) training dynamics are more discriminative for separating these populations but less useful for detecting mislabeled examples. We leverage the complementarity of first- and second-split metrics, showing that by jointly visualizing the two, we can produce a richer characterization of the training examples.

In our experiments, we operationalize several notions of hard examples, namely: (i) **mislabeled** examples, for which the original label has been flipped to a randomly chosen incorrect label; (ii) **rare** examples, which belong to underrepresented subpopulations; and (iii) **complex** examples, which belong to subpopulations for which the classification task is more difficult (details in Section 3.2). We perform specific ablation studies with datasets complicated by just one type of hard example (Section 4.3), and show how SSFT can help to distinguish among these categories of examples. We observe that during second-split training, neural networks (i) first forget mislabeled examples from the first split; (ii) only slowly begin to forget *rare* examples (e.g., from underrepresented subpopulations) unique to the first training set; and (iii) do not forget complex examples.

This separation of hard example types has multiple practical applications. **First**, we can use the method to identify noisy labels: On CIFAR-10 with 10% added class noise, SSFT achieves 0.94 AUC for identifying mislabeled samples, while the first-split metrics range in AUC between 0.58 to 0.90. **Second**, the method can also help improve generalization in noisy data settings: while the removal of hard examples according to first-split learning time degrades the performance of the classifier, the removal of hard examples according to SSFT can actually *improve* generalization. This is especially beneficial when e.g., training on synthetic data (produced by a generative model) or mislabeled data. **Third**, we show how SSFT can identify failure modes of machine learning models. For example, in a simplified task classifying between horses and airplanes in the CIFAR-10 dataset, we find that training examples containing horses with sky backgrounds and airplanes with green backgrounds are among the earliest forgotten—indicating that the model relies on the background as a spurious feature. **Last**, we also add that our metric is robust across multiple seeds, and the earliest forgotten examples are robust across architectures. Across multiple optimizers, SSFT distinguishes mislabeled samples, whereas first-split metrics appear more sensitive to the choice of optimizer.

Finally, we investigate second-split dynamics theoretically, analyzing overparametrized linear models [46]. We introduce notions of mislabeled, rare, and complex examples appropriate to this toy model. Our analysis shows that mislabeled examples from the first split are forgotten quickly during second-split training whereas rare examples are not. However, as we train for a long time, rare examples from the first split are eventually forgotten as the model converges to the minimum norm solution on the second split while predictions on complex examples remain accurate with high probability.

## 2 Related Work

**Example Hardness.** Several recent works quantify example hardness with various training-time metrics. Many of these metrics are based on first-split learning dynamics [8, 25, 27, 35, 43]. Other works have resorted to properties of deep networks such as compression ability [21] and prediction depth [5]. Carlini et al. [7] study metrics centered around model training such as confidence, ensemble agreement, adversarial robustness, holdout retraining, and accuracy under privacy-preserving training. Closest in spirit to the SSFT studied in our paper are efforts in [7, 47]. Crucially, Carlini et al. [7] study the KL divergence of the prediction vector after fine-tuning on a held-out set at a low learning rate, and do not draw any direct inference of the separation offered by their metric. Focusing on (first-split) forgetting dynamics, Toneva et al. [47] defined a metric based on the number of forgetting events during training and identified sets of *unforgettable* examples that are never misclassified once learned. In our work, we find complementary benefits of analysis based on first- and second-split dynamics.

**Memorization of Data Points.** In order to capture the memorization ability of deep networks, their ability to memorize noise (or randomly labeled samples) has been studied in recent work [3, 48]. As opposed to the memorization of rare examples, the memorization of noisy samples hurts generalization and makes the classifier boundary more complex [15]. On the contrary, a recent line of works has argued how memorization of (atypical) data points is important for achieving optimal generalization performance when data is sampled from long-tailed distributions [6, 11, 15].

**Simplicity Bias.** Another line of work argues that neural networks have a bias toward learning simple features [43], and often do not learn complex features even when the complex feature is more predictive of the true label than the simple features. This suggests that models end up memorizing (through noise) the few samples in the dataset that contain the complex feature alone, and utilize the simple feature for correctly predicting the other training examples [1, 32].

**Label Noise.** Large-scale machine learning datasets are typically labeled with the help of human labelers [12] to facilitate supervised learning. It has been shown that a significant fraction of these labels are erroneous in common machine learning datasets [39]. Learning under noisy labels is a long-studied problem [2, 26, 31, 37]. Various recent methods have also attempted to identify label noise [10, 23, 38, 40]. While the focus of our work is not to propose a new method in this long line of work, we show that the view of forgetting time naturally distills out examples with noisy labels. Future work may benefit by augmenting our metric with SOTA methods for label noise identification.

## 3 Method

The primary goal of our work is to *characterize* the hardness of different datapoints in a given dataset. Suppose we have a dataset  $\mathcal{S}_A = \{\mathbf{x}_i, \mathbf{y}_i\}^n$  such that  $(\mathbf{x}_i, \mathbf{y}_i) \sim \mathcal{D}$ . For the purpose of characterization, we augment each datapoint  $(\mathbf{x}_i, \mathbf{y}_i) \in \mathcal{S}_A$  with parameters  $(\text{fslt}_i, \text{ssft}_i)$  where  $\text{fslt}_i$  quantifies the first-split learning time (FSLT), and  $\text{ssft}_i$  quantifies the second-split forgetting time (SSFT) of the sample. To obtain these parameters, we next describe our proposed procedure.

**Procedure** We train a model  $f$  on  $\mathcal{S}$  to minimize the empirical risk:  $\mathcal{L}(\mathcal{S}; f) = \sum_i \ell(f(\mathbf{x}_i), \mathbf{y}_i)$ . We use  $f_A$  to denote a model  $f$  (initialized with random weights) trained on  $\mathcal{S}_A$  until convergence (100% accuracy on  $\mathcal{S}_A$ ). We then train a model initialized with  $f_A$  on a held-out split  $\mathcal{S}_B \sim \mathcal{D}^n$  until convergence. We denote this model with  $f_{A \rightarrow B}$ . To obtain parameters  $(\text{fslt}_i, \text{ssft}_i)$ , we track per-example predictions  $(\hat{\mathbf{y}}_i^t)$  at the end of every epoch ( $t^{\text{th}}$ ) of training. Unless specified otherwise, we train the model with cross-entropy loss using Stochastic Gradient Descent (SGD).

**Definition 1** (First-Split Learning Time). For  $\{\mathbf{x}_i, \mathbf{y}_i\} \in \mathcal{S}_A$ , learning time is defined as the earliest epoch during the training of a classifier  $f$  on  $\mathcal{S}_A$  after which it is always classified correctly, i.e.,

$$\text{fslt}_i = \underset{t^*}{\operatorname{argmin}} (\hat{\mathbf{y}}_{i,(A)}^t = \mathbf{y}_i \quad \forall t \geq t^*) \quad \forall \{\mathbf{x}_i, \mathbf{y}_i\} \in \mathcal{S}_A. \quad (1)$$

**Definition 2** (Second-Split Forgetting Time). Let  $\hat{\mathbf{y}}_{i,(A \rightarrow B)}^t$  to denote the prediction of sample  $\{\mathbf{x}_i, \mathbf{y}_i\} \in \mathcal{S}_A$  after training  $f_{(A \rightarrow B)}$  for  $t$  epochs on  $\mathcal{S}_B$ . Then, for  $\{\mathbf{x}_i, \mathbf{y}_i\} \in \mathcal{S}_A$  forgetting time is defined as the earliest epoch after which it is never classified correctly, i.e.,

$$\text{ssft}_i = \underset{t^*}{\operatorname{argmin}} (\hat{\mathbf{y}}_{i,(A \rightarrow B)}^t \neq \mathbf{y}_i \quad \forall t \geq t^*) \quad \forall \{\mathbf{x}_i, \mathbf{y}_i\} \in \mathcal{S}_A. \quad (2)$$

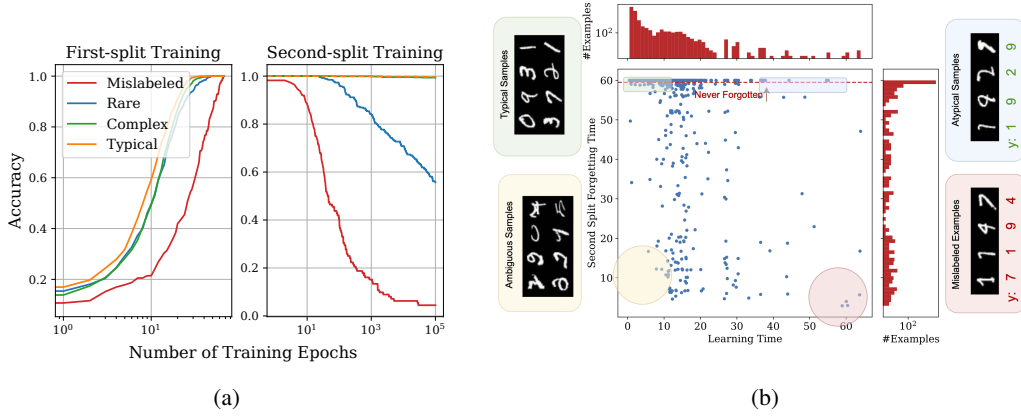


Figure 2: Rate of Learning and Forgetting of examples for different groups in the synthetic dataset. While first-split training is not able to distinguish between rare and complex examples, second-split training succeeds in distinguishing them. Additionally, second-split training separates mislabeled examples from the rest relatively better than first-split training. (b) Visualization of first-split learning and second-split forgetting times when training LeNet model on the MNIST dataset.

### 3.1 Baseline Methods

We provide a brief description of metrics for example hardness considered in recent comparisons [25].

**Number of Forgetting Events:** ( $n_f$ ). An example  $(\mathbf{x}_i, \mathbf{y}_i) \in \mathcal{S}$  undergoes a forgetting event when the accuracy on the example decreases between two consecutive updates. Toneva et al. [47] analyzed the total number of such events  $n_f$  during the training of a neural network to identify hard examples.

**Cumulative Learning Accuracy:** ( $\text{acc}_l$ ). Jiang et al. [25] suggest that rather than using the learning time (Definition 1), using the number of epochs during training when a machine learning model correctly classifies a given sample is a more stable metric for predicting example hardness.

**Cumulative Learning Confidence:** ( $\text{conf}_l$ ). Similar to  $\text{acc}_l$ ,  $\text{conf}_l$  measures the cumulative softmax confidence of the model towards the correct class over the course of training.

### 3.2 Example Characterization

We characterize example hardness via three sources of learning difficulty: **(i) Mislabeled Examples:** We refer to mislabeled examples as those datapoints whose label has been flipped to an incorrect label uniformly at random. **(ii) Rare Examples:** We assume that rare examples belong to sub-populations of the original distribution that have a low probability of occurrence. In particular, there exist  $O(1)$  examples from such sub-populations in a given dataset. In Section 4.3 we describe how we operationalize this notion in the case of the CIFAR-100 dataset. **(iii) Complex Examples:** These constitute samples that are drawn from sub-groups in the dataset that require either (1) a hypothesis class of high complexity; or (2) higher sample complexity to be learnt relative to examples from rest of the dataset. We leave the definition of complex samples mathematically imprecise, but with the same intuitive sense as in prior work [3, 43]. For instance, in a dataset composed of the union of MNIST and CIFAR-10 images, we would consider the subpopulation of CIFAR-10 images to be more *complex*.

## 4 Empirical Investigation of First- and Second-Split Training Dynamics

### 4.1 Experimental Setup

**Datasets** We show results on a variety of image classification datasets—MNIST [13], CIFAR-10 [29], and Imagenette [22]. For experiments in the language domain, we use the SST-2 dataset [45]. For each of the datasets, we split the training set into two equal partitions ( $\mathcal{S}_A, \mathcal{S}_B$ ). For experiments

Sentences in SST-2 dataset with smallest forgetting time	Label
The director explores all three sides of his story with a sensitivity and an inquisitiveness reminiscent of Truffaut	Neg
Beneath the film's obvious determination to shock at any cost lies considerable skill and determination , backed by sheer nerve	Neg
This is a fragmented film, once a good idea that was followed by the bad idea to turn it into a movie	Pos
The holiday message of the 37-minute Santa vs. the Snowman leaves a lot to be desired.	Pos
Epps has neither the charisma nor the natural affability that has made Tucker a star	Pos
The bottom line is the piece works brilliantly	Neg
Alternative medicine obviously has its merits ... but Ayurveda does the field no favors	Pos
What could have easily become a cold, calculated exercise in postmodern pastiche winds up a powerful and deeply moving	Neg
example of melodramatic moviemaking	
Lacks depth	Pos
Certain to be distasteful to children and adults alike , Eight Crazy Nights is a total misfire	Pos

Table 1: First-split sentences that were forgotten by the 3rd epoch of second-split training of a BERT-base model on the SST-2 dataset. Notice that all of these forgotten examples are mislabeled.

with mislabeled examples, we simulate mislabeled examples by randomly selecting a subset of 10% examples from both the partitions and changing their label to an incorrect class.

**Training Details** Unless otherwise specified, we train a ResNet-9 model [4] using SGD optimizer with weight decay  $5e-4$  and momentum 0.9. We use the cyclic learning rate schedule [44] with a peak learning rate of 0.1 at the 10th epoch. We train for a maximum of 100 epochs or until we have 5 epochs of 100% training accuracy. We first train on  $\mathcal{S}_A$ , and then using the pre-initialized weights from stage 1, train on  $\mathcal{S}_B$  with the same learning parameters. All experiments can be performed on a single RTX2080 Ti. Complete hyperparameter details are available in Appendix B.1.

## 4.2 Learning-Forgetting Spectrum for various datasets

**Synthetic Dataset** We consider data  $(\mathbf{x}, \mathbf{y})$  sampled from a mixture of multiple distributions  $\mathcal{D}_g$ , s.t.  $\mathbf{x} \in \mathbb{R}^d$ .  $\mathcal{D}_g$  denotes the  $g^{\text{th}}$  group and has a sampling frequency of  $\pi_g$ . Each group  $\mathcal{D}_g \equiv (\mathcal{X}_g, \{\mathbf{y}_g\})$ , i.e., the true label for all the samples drawn from a given group is the same, and the examples in each group are non-overlapping. Each group is parametrized by a set of  $k \ll d$  unique indices  $\mathcal{I}_g \subset [d]$  such that  $\mathcal{I}_i \cap \mathcal{I}_j = \emptyset$  for  $i \neq j$ . The discriminative characteristic of each group is the vector  $\mathbf{u}_g$ , such that,  $[\mathbf{u}_g]_i = 1$  if  $i \in \mathcal{I}_g$  else 0  $\forall i \in [d]$ . Then for any sample  $(\mathbf{x}, \mathbf{y}) \in \mathcal{S}$ :

$$P(\mathbf{x} \in \mathcal{X}_g) = \pi_g; \quad \mathbf{x} | \mathcal{X}_g \sim \mathcal{N}(0, \sigma^2 \mathbf{I}_d) + \mu_g.$$

For our simulation, we consider a 10 class-classification problem, with  $\mu_g = 5$  for typical groups, and  $\mu_g = 4$  for complex groups (higher signal to noise ratio). For any sample drawn from a rare group, we have  $O(1)$  samples from that group in the entire dataset ( $\mathcal{S}_A \cup \mathcal{S}_B$ ). Mislabeled samples are only generated from the majority typical groups. In Figure 2a, we show the rate of learning and forgetting of examples from each of these categories. We note that in the second-split training, the mislabeled examples are quickly forgotten, and the complex examples are never forgotten. The rare examples are forgotten slowly. In Section 5 we will theoretically justify the observations in the synthetic dataset and show that the rare examples are expected to be forgotten as we train for an infinite time.

**Image Domain** In Figure 2b, we show representative examples in the four quadrants of the learning-forgetting spectrum. More specifically, we find that the examples forgotten fastest and learned last are mislabeled. And the ones learned early and never forgotten once learned are characteristic simple examples of the MNIST dataset. Examples in the first and third quadrant are seemingly atypical and ambiguous respectively. Similar visualizations for other image datasets can be found in Appendix B.2.

**Other Modalities** The forgetting and learning dynamics occur broadly across modalities apart from images. We repeat the same problem setup on the SST-2 [45] dataset for sentiment classification. We fine-tune a pre-trained BERT-base model [14] successively on two disjoint splits of the dataset. In Table 1, we provide a list of the earliest forgotten samples when we train a BERT model on the second split of SST-2 dataset. The results suggest that SSFT is able to identify mislabeled samples.

## 4.3 Ablation Experiments

We design specific experimental setups to capture the three notions of hardness as defined in Section 3.



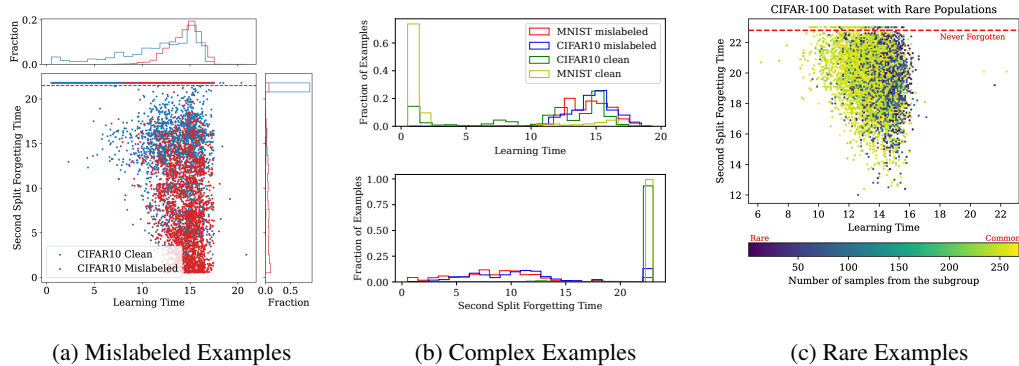


Figure 3: Ablation experiments to distinguish the learning and forgetting dynamics for specific types of hard examples. (a) Mislabeled samples may be learned as slowly as a high fraction of typical samples, but they are forgotten much faster. (b) FSLT distinguishes complex (CIFAR10 clean) and simple (MNIST clean) examples, but SSFT does not. On the contrary, FSLT can not distinguish clean and mislabeled examples of CIFAR10, while the SSFT can. (c) FSLT is able to distinguish examples based on the sub-group frequency, however, SSFT has a low correlation with the sub-group frequency.

**Mislabeled Examples** We sample 10% datapoints from both the first and second split of the CIFAR-10 dataset, and randomly change their label to an incorrect label. Figure 3a shows the learning-forgetting spectrum for the dataset. In the adjoining density histograms, note that a large fraction of the mislabeled and correctly labeled examples are learned at the same time. However, during second-split training, the mislabeled examples are forgotten quickly whereas a large fraction of the clean examples are never forgotten, allowing SSFT to succeed in distinguishing mislabeled samples.

**Complex Examples** We generate a joint dataset that contains the union of both MNIST and CIFAR-10 examples. This is motivated by work in simplicity bias [43] that argue that neural networks learn simpler features first. We also add 10% labeled noise to each of the datasets in the union to understand the learning and forgetting time relationship of a sample that is complex or mislabeled together. In Figure 3b, we show the FSLT and SSFT for MNIST and CIFAR-10 samples. We note that a high fraction of the CIFAR-10 (complex) samples learn at the same speed as the mislabeled samples. However, when looking at the SSFT, we are able to draw a strong separation between the mislabeled samples and complex samples. This indicates that the complexity of a sample has low correlation with its tendency to be forgotten once learnt, but a high correlation with being learned slowly.

**Rare Examples** The CIFAR-100 [29] dataset is a 100-class classification task. The dataset contains 20 superclasses, each containing 5 subclasses. We create a 20-class classification dataset with long tails simulated through the 5 sub-classes within each superclass. More specifically, the number of examples in each subgroup for a given superclass is given by {500, 250, 125, 64, 32} respectively (exponentially decaying with a factor of 2). This is done to simulate the hypothesis of dataset subgroups following a Zipf distribution [49] as argued for by Feldman [15]. This dataset is further divided into two equal splits to analyze the learning-forgetting dynamics. In order to remove any other effects of example hardness (either within a subgroup, or among subgroups), we randomize both the chosen subset of examples and the ordering of the majority and minority groups between each superclass, by training the model on 20 such random splits and aggregating learning and forgetting statistics over these runs. In Figure 3c, we show a scatter plot for the FSLT and SSFT, colored by the frequency of the group a particular example belongs to. We observe that FSLT strongly correlates with the size of the subgroup, whereas the SSFT has a very low correlation with the rareness of a sample.

We provide further ablations to show that FSLT is able to identify hard and rare examples, but SSFT shows nearly no discriminative power at finding the two in Appendix C.

#### 4.4 Dataset Cleansing

**Identifying Label Noise** We present AUC scores for detection of label noise via various popular methods in example difficulty literature, across various datasets in Table 2. We note that (i) cumulative predictions over the course of training help stabilize both the learning time and forgetting time metrics;

Method →	fslt	acc <sub>l</sub>	ssft (Ours)	acc <sub>f</sub> (Ours)	conf <sub>l</sub>	n <sub>f</sub>	Joint (Ours)
Imagenette	0.834	0.912	0.931	0.941	0.786	0.781	<b>0.957</b>
CIFAR10	0.740	0.900	0.938	0.941	0.947	0.580	<b>0.958</b>
MNIST	0.973	<b>0.998</b>	0.997	<b>0.998</b>	0.965	0.377	<b>0.998</b>
CIFAR100	0.700	0.899	0.865	0.885	0.860	0.300	<b>0.926</b>
EMNIST	0.987	0.990	0.987	0.989	0.984	0.386	<b>0.997</b>

Table 2: AUC for identification of label noise using various metrics for example hardness across different datasets. Across all datasets, our **ssft** metric outperforms alternative baselines. We introduce **acc<sub>f</sub>** as the cumulative accuracy on the second-split training, inspired by previous work that suggests using cumulative accuracies helps make first-split learning time more stable [25]. All other notations are described in Section 3. In the case of the Joint method, we select new prediction ranks based on the combined learning and forgetting ranks, further improving over the **ssft** metric alone.

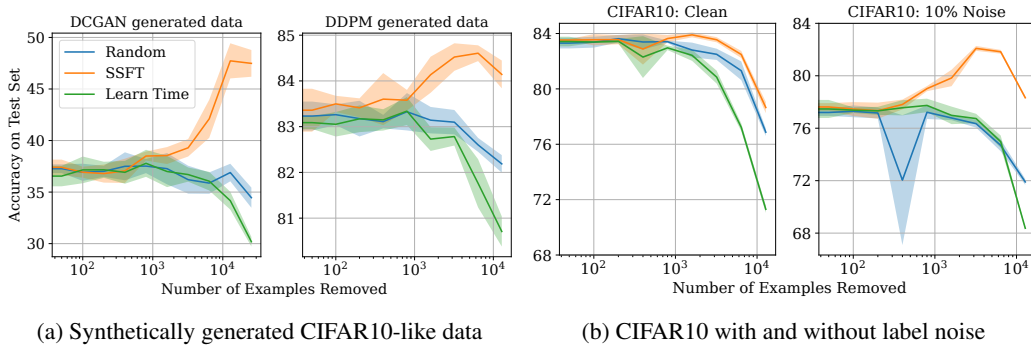


Figure 4: Accuracy on CIFAR-10 test set after removing a varying number of samples from the training set based on (i) random choice, (ii) examples with the lowest SSFT, and (ii) the highest FSLT. Removing examples based on SSFT helps improve the generalization on the original test set.

(ii) for simple datasets such as MNIST with few ambiguous images, all of the baseline methods have very high AUC (greater than 0.99) in finding noisy inputs. However, in datasets such as CIFAR-10 and Imagenette, we find that second-split forgetting metrics do better than first-split training metrics. Finally, we also compare the use of both forgetting and learning time to find noisy samples, and we find a small improvement in the results of just using the forgetting time. While we do not make explicit comparisons with other state of art methods dedicated to finding label noise, our results suggest that augmenting second split forgetting time information may help improve their results. As also observed in recent work [25], we find that the number of forgetting events (**n<sub>f</sub>**) [47] is an unreliable indicator of mislabeled samples. We hypothesize that this is because of the fact that mislabeled examples may often be (first) learnt very late, hence their count of total forgetting events is also low.

**Cleaning synthetically generated datasets** Generative models are capable of mimicking the distribution of a given dataset. We generate synthetic datasets of CIFAR10-like samples using (i) DDPM (denoising diffusion model [20]); and (ii) DCGAN (Deep Convolutional GAN [41]). In both cases, we assign pseudo-labels using the BiT model [28] as in prior work [36]. We collect a sample of 50,000 training examples and record the generalization performance on CIFAR-10 as we remove ‘hard’ samples, as evaluated by various metrics. In Figure 4, we can see that removing the most easily forgotten examples can benefit by up to 10% generalization accuracy on the clean test set of CIFAR-10. In case of the synthetic data generated using DDPM, the gains in generalization performance are under 2%. We hypothesize that this is because the samples generated by DDPM are more representative of the typical distribution of CIFAR-10 than those generated by DCGAN.

**Note:** The ability to train on a second split allows SSFT the *unique* opportunity to train on a clean split of CIFAR-10 in order to assess the alignment of the synthetic samples with the oracle samples. As a result, the SSFT is much more effective in filtering out ambiguous first-split synthetic examples.

#### 4.5 Evaluating Example Utility

Recent works [16, 47] have argued for removing a large fraction of the less memorized examples, and keeping the memorized ones. We will analyze the change in model generalization upon removing varying sizes of examples from the training set, as ranked by lowest SSFT and highest FSLT (Figure 4). In the presence of noisy examples, removing samples based on the SSFT helps improve generalization, whereas FSLT does not do much better than random. We draw the following inferences:

**FSLT finds important samples** As we remove more samples from the dataset, the accuracy of the model trained after samples are removed based on the highest FSLT is significantly lower than random guessing. This suggests that the utility of these samples is higher than random samples. Put in line with the hypothesis of memorization of rare example as proposed in [15], we see that empirically, the examples that are slow to learn are important for the model’s test set generalization.

**SSFT removes pathological samples** On the contrary, removing examples based on the SSFT helps improve model generalization (especially when there is label noise). Even in the setting when there is no label noise, in contrast to FSLT, we find that removing examples that were easily forgotten has a lower negative impact on the model’s generalization as opposed to removing random samples. This suggests that the examples that are forgotten in the early epochs of second-split training hurt a model’s generalization, and may not be characteristic samples of their particular class.

**Practitioner’s view** From the AUC numbers in Table 2, it may appear that removing examples via learning-based metrics such as learning time and cumulative learning accuracy also provides a high rate of removal of noisy samples. However, when we observe the example utility graphs in Figure 4, we draw the inference that the examples that are learned late, are often important examples (such as rare memorized examples). However, even when SSFT fails to capture the correct noisy examples, it still removes unimportant samples and does not hurt generalization. Similar graphs for other metrics described in Table 2 can be found in Appendix B.

#### 4.6 Characterizing Potential Failure Modes

Recent works have attempted to train classifiers on datasets that contain spurious features [24, 42] (example Waterbirds, CelebA [34] dataset). However, a fundamental challenge is to first identify the spurious correlation that the classifier may be relying on. Only then can recent methods be trained to remove the reliance on spurious patterns. We train a ResNet-9 model to classify CIFAR-10 images of horses and airplanes. In Figure 5, we observe that the model forgets planes with green backgrounds and horses with blue backgrounds. This suggests that the model relied on the background as a spurious feature. By analyzing the forgotten examples we can further investigate the examples that the classifier fails to generalize to.



Figure 5: By inspecting the earliest-forgotten examples, we can gain insights into potential failure modes. This model quickly forgets planes with green backgrounds and horses on blue backgrounds.

**Stability of SSFT** We note that SSFT is stable across multiple seeds (Pearson correlation of 0.81), and across architectures (Pearson correlation of 0.63). While the overall correlation for samples ranked by SSFT may be low across architectures, the top-ranked examples have a high correlation (0.85), suggesting the most forgotten examples are consistent across architectures. In contrast, FSLT has a Pearson correlation of 0.52 across seeds. Most interestingly, the learning time metric is brittle to the choice of hyperparameters. As shown by Jiang et al. [25], when using Adam optimizer, examples of different hardness get learned together. In our experiments, we observe the same phenomenon during learning, however, SSFT is robust to the choice of the optimizer. Detailed results in Appendix C.1.

**Limitations** One limitation of the proposed metric is that it is brittle to the choice of the learning rate for the second-split training. If we use a very small learning rate, then overparametrized deep models are capable of learning the new dataset without forgetting examples from the first split. Alternately, if we use a very large learning rate, the model may diverge and undergo catastrophic forgetting. However, under ‘reasonable’ choices of learning rate (like that for first-split training), we find SSFT is robust. We provide a detailed analysis of the same in Appendix C.1.



## 5 Theoretical Results

Through our theoretical analysis, we will characterize the forgetting dynamics of mislabeled, rare and complex examples in a simplified version of the framework used for our synthetic experiments in Figure 2a. Recall, our setup contains two dataset splits  $\mathcal{S}_A, \mathcal{S}_B$ , where we train on the first split until achieving perfect accuracy on all training points, and then with these weights train on  $\mathcal{S}_B$  for infinite time. In particular, we will prove that both mislabeled and rare examples are forgotten upon training for infinite time, with mislabeled examples being forgotten much faster. Further, we will show that complex examples from the first split do not get forgotten if not continually trained on. We assume in our analysis that  $\mathcal{S}_B$  has no mislabeled or rare examples, and  $\mathcal{S}_A$  contains one example of each type.

We consider a dataset  $\mathcal{S} = \{\mathbf{x}_i, \mathbf{y}_i\}^n$  such that  $(\mathbf{x}_i, \mathbf{y}_i) \in \mathcal{X} \times \mathcal{Y}$ , and  $\mathbf{x}_i = \boldsymbol{\mu}_g + \mathbf{z}_i$  where  $\mathbf{z}_i \sim \mathcal{N}(0, \sigma^2 \mathbf{I}_d)$ , and  $\|\boldsymbol{\mu}_g\|_2^2 = k\mu^2$  (as in Section 4.2). Let  $\mathbf{w} \in \mathbb{R}^d$  represent the weight vector of an overparametrized linear model. We analyze the learning and forgetting dynamics by minimizing the empirical risk:  $\mathcal{L}(\mathcal{S}; \mathbf{w}) = \sum_i \ell(\mathbf{w}^\top \mathbf{y}_i \mathbf{x}_i)$ , where  $\ell$  is the exponential loss. Following Chatterji and Long [9], we make the following assumptions about the problem setup:

(A.1) The failure probability satisfies  $0 \leq \delta \leq 1/C$ ,

(A.2) The number of samples satisfies  $n \geq C \log(1/\delta)$ ,

(A.3) The input dimension  $d \geq C \max\{n^2 \log(n/\delta), n(k\mu^2/\sigma^2)\}$ , and  $k\mu^2/\sigma^2 \geq C \log(n/\delta)$ ,

where  $k\mu^2/\sigma^2$  represents the signal to noise ratio and  $C$  is a large constant. Now we formalize the notions of rare, mislabeled and complex examples for our theoretical analysis.

**Definition 1** (Rare Examples,  $\mathcal{R}$  [15]). *Consider a dataset  $\mathcal{S}$  sampled from a mixture of distributions  $\{\mathcal{D}_1, \dots, \mathcal{D}_N\}$  with frequency  $\{\pi_1, \dots, \pi_N\}$  respectively. Let  $\mathcal{R} \subseteq \mathcal{S}$  be the set of rare examples. Then, for all  $(\mathbf{x}_i, \mathbf{y}_i) \in \mathcal{R}$ , if  $(\mathbf{x}_i, \mathbf{y}_i) \sim \mathcal{D}_j$ , then there are  $O(1)$  samples from  $\mathcal{D}_j$  in  $\mathcal{S}$ .*

**Definition 2** (Mislabeled Examples,  $\mathcal{M}$ ). *Consider a  $k$  class classification problem with  $\mathcal{Y} = \{1, 2, \dots, k\}$ . Let  $\mathcal{M} \subset \mathcal{S}$  be the set of mislabeled examples. Then for any  $(\mathbf{x}, \mathbf{y}) \sim \mathcal{D}$ , a corresponding mislabeled example is given by  $(\mathbf{x}, \tilde{\mathbf{y}}) \in \mathcal{M}$  such that  $\tilde{\mathbf{y}} \in \mathcal{Y} \setminus \{\mathbf{y}\}$ .<sup>2</sup>*

**Definition 3** (Complex Examples,  $\mathcal{C}$ ). *Let  $\mathcal{C} \subset \mathcal{S}$  be the set of examples sampled from complex distributions. Let  $(\mathbf{x}_i, \mathbf{y}_i) \in \mathcal{C}$  such that  $(\mathbf{x}_i, \mathbf{y}_i) \sim \mathcal{D}_g$  (complex group), then  $\mu_g = \frac{\mu_t}{\lambda}$ ,  $\lambda > 1$  where  $\mu_t$  is the coordinate-wise mean for samples drawn from any simple distribution  $\mathcal{D}_t$  (Section 4.2).*

**Optimization** We perform gradient descent with fixed learning rate  $\eta$ ,

$$\mathbf{w}(t+1) = \mathbf{w}(t) - \eta \nabla \mathcal{L}(\mathbf{w}(t)) = \mathbf{w}(t) - \eta \sum_i \ell'(\mathbf{w}^\top \mathbf{y}_i \mathbf{x}_i) \cdot \mathbf{y}_i \mathbf{x}_i. \quad (3)$$

**Solution dynamics** For sufficiently small learning rate  $\eta$ , and (bounded) starting point  $\mathbf{w}(0)$ , Soudry et al. [46] showed that:

$$\mathbf{w}(t) = \hat{\mathbf{w}} \log t + \rho(t), \quad (4)$$

where  $\rho(t)$  is a bounded residual term, and  $\hat{\mathbf{w}}$  is the solution to the hard margin SVM:

$$\hat{\mathbf{w}} = \underset{\mathbf{w} \in \mathbb{R}^d}{\operatorname{argmin}} \|\mathbf{w}\|_2^2 \quad \text{s.t.} \quad \mathbf{w}^\top \mathbf{y}_i \mathbf{x}_i \geq 1, \quad (5)$$

### 5.1 First-split Learning

For stage 1, we consider that we train the model for a maximum of  $T$  epochs (until we achieve 100% accuracy on the first training dataset  $\mathcal{S}_A$ ). This means that the learned weight vectors are close to, but have not converged to the max margin solution. The solution at the end of  $t$  epochs is given by  $\mathbf{w}_A(t)$ . At sufficiently large  $T$ , we have:

$$\begin{aligned} \mathbf{w}_A(T) &= \hat{\mathbf{w}}_A \log T + \rho_A(T) \\ \mathbf{w}_A(T)^\top \mathbf{y}_i \mathbf{x}_i &\geq 1 \quad \forall (\mathbf{x}_i, \mathbf{y}_i) \in \mathcal{S}_A \end{aligned} \quad (6)$$

<sup>2</sup>For binary classification,  $\mathcal{Y} = \{-1, +1\}$ . The labels are reversed for mislabeled examples.

## 5.2 Second-split Forgetting

We initialize the weights for second stage of training with  $\mathbf{w}_A(T)$  from first training stage, and then train on  $\mathcal{S}_B$ . We provide the formal theorem statement and complete proofs in Appendix A, but provide informal theorem statements and an intuitive proof sketch below:

**Theorem 1** (Asymptotic Forgetting (informal)). *For sufficiently small learning rate, given datasets  $\mathcal{S}_A, \mathcal{S}_B \sim \mathcal{D}^n$ . After training for  $T' \rightarrow \infty$  epochs, the following hold with high probability:*

1. *Mislabeled and Rare examples from  $\mathcal{S}_A$  are forgotten.*
2. *Complex examples from  $\mathcal{S}_A$  are not forgotten.*

*Proof Sketch.* We use the result from Soudry et al. [46] that for any bounded initialization, when trained on a separable data, the model converges to the same min-norm solution. As a result, we can ignore the impact of  $\mathcal{S}_A$  at infinite time training. Then, we use generalization bounds from Chatterji and Long [9] to argue about the accuracy on mislabeled and complex examples. For the case of rare examples, we show that the probability of correct model prediction can be approximated by a Gaussian CDF with mean 0 and  $\mathcal{O}(1/\sqrt{n})$  variance.

**Theorem 2** (Intermediate-Time Forgetting (informal)). *For sufficiently small learning rate, given two datasets  $\mathcal{S}_A, \mathcal{S}_B \sim \mathcal{D}^n$ . For a model initialized with weights,  $\mathbf{w}_B(0) = \mathbf{w}_A(T)$  and trained for  $T' = f(T)$  epochs, the following hold with high probability:*

1. *Mislabeled examples from  $\mathcal{S}_A$  are no longer incorrectly predicted.*
2. *Rare examples from  $\mathcal{S}_A$  are not forgotten.*

*Proof Sketch.*  $\mathcal{S}_B$  contains examples from the same majority distributions as  $\mathcal{S}_A$ . The mislabeled example also belongs to one of these distributions, but has the opposite label. However,  $\mathcal{S}_B$  does not have samples from rare groups found in  $\mathcal{S}_A$ . Using representer theorem, we decompose the model updates into a weighted sum of each training data point in  $\mathcal{S}_B$ . Then, we analyze the change in prediction on rare and mislabeled examples, which is a dot product of the weight update with  $\mathbf{x}_m$  or  $\mathbf{x}_r$ . Per our assumptions, the mean of each group  $\mu_g$  is orthogonal to the other. As a result, the rare example finds negligible coupling with any example in  $\mathcal{S}_B$ , and the variance of its prediction keeps increasing due to the noise term contributed in the model weights from each example in  $\mathcal{S}_B$ . On the contrary, the mislabeled examples have a strong coupling with all the examples in its group. Due to its incorrect label, the mean of its predictions moves towards the correct label, with variance increasing at a similar rate. The final step is to jointly analyze the rate of change of prediction of both the examples, and find an optimal time  $T'$  when the prediction on the mislabeled example is flipped and the rare example still retains its prediction with high probability.

## 6 Conclusion

While many prior works investigate training time dynamics to characterize the hardness of examples, we enrich this literature with a complementary lens focused on the second-split forgetting time. We demonstrate the potential of SSFT to distinguish among rare, mislabeled, and complex examples; and also show the differences in the example properties captured by first-split and second-split metrics.

Our work opens new lines of inquiry in future work that can utilize the separation of hard examples. First, we expect state of art methods in label noise identification to benefit by augmenting our approach. Further, we believe our ablations showing that complex, noisy, and mislabeled samples may all be learned slowly inspire future work that can unite different takes on the memorization-generalization research—early learning, simplicity bias, and singleton memorization.

## Acknowledgements

We would like to thank Aakash Lahoti and Jeremy Cohen for their insightful comments on this work. SG acknowledges Amazon Graduate Fellowship and JP Morgan PhD Fellowship for their support. ZL acknowledges Amazon AI, Salesforce Research, Facebook, UPMC, Abridge, the PwC Center, the Block Center, the Center for Machine Learning and Health, and the CMU Software Engineering Institute (SEI) via Department of Defense contract FA8702-15-D-0002, for their generous support.

## References

- [1] Zeyuan Allen-Zhu and Yuanzhi Li. Towards understanding ensemble, knowledge distillation and self-distillation in deep learning. *arXiv preprint arXiv:2012.09816*, 2020.
- [2] Dana Angluin and Philip Laird. Learning from noisy examples. *Machine Learning*, 2(4):343–370, 1988.
- [3] 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*, pages 233–242. PMLR, 2017.
- [4] Woonhyuk Baek. Torchskelton. <https://github.com/wbaek/torchskelton>, 2019.
- [5] Robert John Nicholas Baldock, Hartmut Maennel, and Behnam Neyshabur. Deep learning through the lens of example difficulty. In A. Beygelzimer, Y. Dauphin, P. Liang, and J. Wortman Vaughan, editors, *Advances in Neural Information Processing Systems*, 2021. URL <https://openreview.net/forum?id=WWRBHhH158K>.
- [6] Gavin Brown, Mark Bun, Vitaly Feldman, Adam Smith, and Kunal Talwar. When is memorization of irrelevant training data necessary for high-accuracy learning? In *Proceedings of the 53rd Annual ACM SIGACT Symposium on Theory of Computing*, pages 123–132, 2021.
- [7] Nicholas Carlini, Úlfar Erlingsson, and Nicolas Papernot. Distribution density, tails, and outliers in machine learning: Metrics and applications. *ArXiv*, abs/1910.13427, 2019.
- [8] Satrajit Chatterjee. Coherent gradients: An approach to understanding generalization in gradient descent-based optimization. *arXiv preprint arXiv:2002.10657*, 2020.
- [9] Niladri S Chatterji and Philip M Long. Finite-sample analysis of interpolating linear classifiers in the overparameterized regime. *J. Mach. Learn. Res.*, 22:129–1, 2021.
- [10] Pengfei Chen, Ben Ben Liao, Guangyong Chen, and Shengyu Zhang. Understanding and utilizing deep neural networks trained with noisy labels. In *International Conference on Machine Learning*, pages 1062–1070. PMLR, 2019.
- [11] Chen Cheng, John Duchi, and Rohith Kuditipudi. Memorize to generalize: on the necessity of interpolation in high dimensional linear regression. *arXiv preprint arXiv:2202.09889*, 2022.
- [12] 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*, pages 248–255. Ieee, 2009.
- [13] Li Deng. The mnist database of handwritten digit images for machine learning research. *IEEE Signal Processing Magazine*, 29(6):141–142, 2012.
- [14] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- [15] Vitaly Feldman. Does learning require memorization? a short tale about a long tail. *Proceedings of the 52nd Annual ACM SIGACT Symposium on Theory of Computing*, 2020.
- [16] Vitaly Feldman and Chiyuan Zhang. What neural networks memorize and why: Discovering the long tail via influence estimation. *Advances in Neural Information Processing Systems*, 33: 2881–2891, 2020.
- [17] Jonathan Frankle, David J Schwab, and Ari S Morcos. The early phase of neural network training. *arXiv preprint arXiv:2002.10365*, 2020.
- [18] Saurabh Garg, Sivaraman Balakrishnan, Zico Kolter, and Zachary Lipton. RATT: Leveraging unlabeled data to guarantee generalization. *arXiv preprint arXiv:2105.00303*, 2021. URL <https://arxiv.org/abs/2105.00303>.

- [19] Guy Hach Cohen, Leshem Choshen, and Daphna Weinshall. Let’s agree to agree: Neural networks share classification order on real datasets. In *International Conference on Machine Learning*, pages 3950–3960. PMLR, 2020.
- [20] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in Neural Information Processing Systems*, 33:6840–6851, 2020.
- [21] Sara Hooker, Aaron Courville, Gregory Clark, Yann Dauphin, and Andrea Frome. What do compressed deep neural networks forget? *arXiv preprint arXiv:1911.05248*, 2019.
- [22] Jeremy Howard. Imagenette. URL <https://github.com/fastai/imagenette/>.
- [23] Jinchi Huang, Lie Qu, Rongfei Jia, and Binqiang Zhao. O2u-net: A simple noisy label detection approach for deep neural networks. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3326–3334, 2019.
- [24] Badr Youbi Idrissi, Martin Arjovsky, Mohammad Pezeshki, and David Lopez-Paz. Simple data balancing achieves competitive worst-group-accuracy. *arXiv preprint arXiv:2110.14503*, 2021.
- [25] Ziheng Jiang, Chiyuan Zhang, Kunal Talwar, and Michael C Mozer. Characterizing structural regularities of labeled data in overparameterized models. *arXiv preprint arXiv:2002.03206*, 2020.
- [26] Ishan Jindal, Matthew Nokleby, and Xuewen Chen. Learning deep networks from noisy labels with dropout regularization. In *2016 IEEE 16th International Conference on Data Mining (ICDM)*, pages 967–972. IEEE, 2016.
- [27] Gal Kaplun, Nikhil Ghosh, Saurabh Garg, Boaz Barak, and Preetum Nakkiran. Deconstructing distributions: A pointwise framework of learning. *arXiv preprint arXiv:2202.09931*, 2022.
- [28] Alexander Kolesnikov, Lucas Beyer, Xiaohua Zhai, Joan Puigcerver, Jessica Yung, Sylvain Gelly, and Neil Houlsby. Big transfer (bit): General visual representation learning. In *European conference on computer vision*, pages 491–507. Springer, 2020.
- [29] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. *Citeseer*, 2009.
- [30] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. *Advances in neural information processing systems*, 25, 2012.
- [31] Junnan Li, Richard Socher, and Steven C.H. Hoi. Dividemix: Learning with noisy labels as semi-supervised learning. In *International Conference on Learning Representations*, 2020. URL <https://openreview.net/forum?id=HJgExaVtwr>.
- [32] Yuanzhi Li, Colin Wei, and Tengyu Ma. Towards explaining the regularization effect of initial large learning rate in training neural networks. *Advances in Neural Information Processing Systems*, 32, 2019.
- [33] 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.
- [34] Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Deep learning face attributes in the wild. In *ICCV*, pages 3730–3738. IEEE Computer Society, 2015. ISBN 978-1-4673-8391-2. URL <http://dblp.uni-trier.de/db/conf/iccv/iccv2015.html#LiuLWT15>.
- [35] Karttikeya Mangalam and Vinay Uday Prabhu. Do deep neural networks learn shallow learnable examples first? In *ICML 2019 Workshop on Identifying and Understanding Deep Learning Phenomena*, 2019. URL <https://openreview.net/forum?id=HkxHv4rn24>.
- [36] Preetum Nakkiran, Gal Kaplun, Yamini Bansal, Tristan Yang, Boaz Barak, and Ilya Sutskever. Deep double descent: Where bigger models and more data hurt. *Journal of Statistical Mechanics: Theory and Experiment*, 2021(12):124003, 2021.

- [37] Nagarajan Natarajan, Inderjit S Dhillon, Pradeep K Ravikumar, and Ambuj Tewari. Learning with noisy labels. *Advances in neural information processing systems*, 26, 2013.
- [38] Curtis Northcutt, Lu Jiang, and Isaac Chuang. Confident learning: Estimating uncertainty in dataset labels. *Journal of Artificial Intelligence Research*, 70:1373–1411, 2021.
- [39] Curtis G Northcutt, Anish Athalye, and Jonas Mueller. Pervasive label errors in test sets destabilize machine learning benchmarks. *arXiv preprint arXiv:2103.14749*, 2021.
- [40] Geoff Pleiss, Tianyi Zhang, Ethan Elenberg, and Kilian Q Weinberger. Identifying mislabeled data using the area under the margin ranking. *Advances in Neural Information Processing Systems*, 33:17044–17056, 2020.
- [41] Alec Radford, Luke Metz, and Soumith Chintala. Unsupervised representation learning with deep convolutional generative adversarial networks. *arXiv preprint arXiv:1511.06434*, 2015.
- [42] Shiori Sagawa, Pang Wei Koh, Tatsunori B Hashimoto, and Percy Liang. Distributionally robust neural networks for group shifts: On the importance of regularization for worst-case generalization. *arXiv preprint arXiv:1911.08731*, 2019.
- [43] Harshay Shah, Kaustav Tamuly, Aditi Raghunathan, Prateek Jain, and Praneeth Netrapalli. The pitfalls of simplicity bias in neural networks. *Advances in Neural Information Processing Systems*, 33:9573–9585, 2020.
- [44] Leslie N Smith. A disciplined approach to neural network hyper-parameters: Part 1—learning rate, batch size, momentum, and weight decay. *arXiv preprint arXiv:1803.09820*, 2018.
- [45] Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D. Manning, Andrew Ng, and Christopher Potts. Recursive deep models for semantic compositionality over a sentiment treebank. In *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing*, pages 1631–1642, Seattle, Washington, USA, October 2013. Association for Computational Linguistics. URL <https://www.aclweb.org/anthology/D13-1170>.
- [46] Daniel Soudry, Elad Hoffer, Mor Shpigel Nacson, Suriya Gunasekar, and Nathan Srebro. The implicit bias of gradient descent on separable data. *The Journal of Machine Learning Research*, 19(1):2822–2878, 2018.
- [47] Mariya Toneva, Alessandro Sordoni, Remi Tachet des Combes, Adam Trischler, Yoshua Bengio, and Geoffrey J Gordon. An empirical study of example forgetting during deep neural network learning. *arXiv preprint arXiv:1812.05159*, 2018.
- [48] 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.
- [49] George Kingsley Zipf. *The psycho-biology of language: An introduction to dynamic philology*. Routledge, 2013.



# Characterizing Datapoints via Second-Split Forgetting

## Supplementary Material

### A Theoretical Results

#### A.1 Preliminaries

Let  $\mathbf{w} \in \mathbb{R}^d$  represent the weight vector of overparametrized linear model. We analyze the learning and forgetting dynamics by minimizing the empirical risk:  $\mathcal{L}(\mathcal{S}; \mathbf{w}) = \sum_i \ell(\mathbf{w}^\top \mathbf{y}_i \mathbf{x}_i)$ . We consider the exponential loss  $\ell(z) = \exp(-z)$  for our analysis. For completeness, we rewrite the definitions and preliminaries from the main paper below.

**Data Generating Process** We restate the data generating process as detailed in the synthetic experiment in Section 4.2. We consider data  $(\mathbf{x}, \mathbf{y})$  sampled from a mixture of multiple distributions  $\mathcal{D}_g$ , s.t.  $\mathbf{x} \in \mathbb{R}^d$ .  $\mathcal{D}_g$  denotes the  $g^{\text{th}}$  group and has a sampling frequency of  $\pi_g$ . Each group  $\mathcal{D}_g$  is a distribution over  $(\mathcal{X}_g \times \{\mathbf{y}_g\})$ , i.e., the true label for all the samples drawn from a given group is the same, and the examples in each group are non-overlapping. Each group is parametrized by a set of  $k \ll d$  unique indices  $\mathcal{I}_g \subset [d]$  such that  $\mathcal{I}_i \cap \mathcal{I}_j = \emptyset$  for  $i \neq j$ . The discriminative characteristic of each group is the vector  $\mathbf{u}_g$ , such that,  $[\mathbf{u}_g]_i = 1$  if  $i \in \mathcal{I}_g$  else 0  $\forall i \in [d]$ . In the following discussion, we will refer to  $\mu_g$  as the coordinate-wise mean vector for the group  $g$ , such that  $\mu_g = \mu_g \mathbf{u}_g$ . Then for any sample  $(\mathbf{x}, \mathbf{y}) \in \mathcal{S}$ :

$$P(\mathbf{x} \in \mathcal{X}_g) = \pi_g; \quad \mathbf{x} | \mathcal{X}_g \sim \mathcal{N}(0, \sigma^2 \mathbf{I}_d) + \mu_g. \quad (7)$$

**Definition 1** (Rare Examples,  $\mathcal{R}$  [15]). Consider a dataset  $\mathcal{S}$  sampled from a mixture of distributions  $\{\mathcal{D}_1, \dots, \mathcal{D}_N\}$  with frequency  $\{\pi_1, \dots, \pi_N\}$  respectively. Let  $\mathcal{R} \subseteq \mathcal{S}$  be the set of rare examples. Then, for all  $(\mathbf{x}_i, \mathbf{y}_i) \in \mathcal{R}$ , if  $(\mathbf{x}_i, \mathbf{y}_i) \sim \mathcal{D}_j$ , then there are  $O(1)$  samples from  $\mathcal{D}_j$  in  $\mathcal{S}$ .

**Definition 2** (Misabeled Examples,  $\mathcal{M}$ ). Consider a  $k$  class classification problem with  $\mathcal{Y} = \{1, 2, \dots, k\}$ . Let  $\mathcal{M} \subset \mathcal{S}$  be the set of mislabeled examples. Then for any  $(\mathbf{x}, \mathbf{y}) \sim \mathcal{D}$ , a corresponding mislabeled example is given by  $(\mathbf{x}, \tilde{\mathbf{y}}) \in \mathcal{M}$  such that  $\tilde{\mathbf{y}} \in \mathcal{Y} \setminus \{\mathbf{y}\}$ .<sup>1</sup>

**Definition 3** (Complex Examples,  $\mathcal{C}$ ). Let  $\mathcal{C} \subset \mathcal{S}$  be the set of examples sampled from complex distributions. Let  $(\mathbf{x}_i, \mathbf{y}_i) \in \mathcal{C}$  such that  $(\mathbf{x}_i, \mathbf{y}_i) \sim \mathcal{D}_g$  (complex group), then  $\mu_g = \frac{\mu_t}{\lambda}$ ,  $\lambda > 1$  where  $\mu_t$  is the coordinate-wise mean for samples drawn from any simple distribution  $\mathcal{D}_t$  (Section 4.2).

The implication of the aforementioned characterization is that complex distributions have a lower signal-to-noise ratio (SNR) as compared to simple distributions. We assume the sample complexity required to estimate the distribution as a proxy for the complexity of the distribution. In this regard, having a low SNR increases the complexity.

Recall that  $\mathcal{S}_A$  and  $\mathcal{S}_B$  denote the first and second training splits of our dataset. In our theoretical framework, we only consider two majority distributions in the second split dataset ( $\mathcal{S}_B$ ). Let us call them  $\mathcal{D}_1, \mathcal{D}_2$ . Therefore, both  $\mathcal{S}_A$  and  $\mathcal{S}_B$  contain  $O(n)$  samples from  $\mathcal{D}_1, \mathcal{D}_2$ . In the first split dataset ( $\mathcal{S}_A$ ), we consider the presence of another distribution  $\mathcal{D}_r$  that constitutes the rare example (only one sample  $(\mathbf{x}_r, \mathbf{y}_r)$ ). The mislabeled example (only one sample  $(\mathbf{x}_m, \mathbf{y}_m)$ ) belongs to one of the majority distributions, and we will assume without loss of generality that this is from distribution  $\mathcal{D}_1$ . To understand the population accuracy in the case of complex distributions, we will draw a simple analogy in the case when the majority distributions  $\mathcal{D}_1, \mathcal{D}_2$  occur from complex distributions as defined below. Based on equation 7, without loss of generality, we will assume that dimensions  $\{1 \dots k\}, \{k + 1 \dots 2k\}, \{2k + 1 \dots 3k\}$  are the predictive dimensions for the majority group 1, 2 and that for the rare example from dataset 1. We make these assumptions to simplify the theoretical exposition. However, our results can be observed even after relaxing them at the expense of more book-keeping.

Based on Chatterji and Long [9], we make the following assumptions about the problem setup:

(A.1) The failure probability satisfies  $0 \leq \delta \leq 1/C$ ,

(A.2) The number of samples satisfies  $n \geq C \log(1/\delta)$ ,

<sup>1</sup>For binary classification,  $\mathcal{Y} = \{-1, +1\}$ . The labels are reversed for mislabeled examples.

(A.3) The input dimension  $d \geq C \max\{n^2 \log(n/\delta), nk\mu^2/\sigma^2\}$ , and  $k\mu^2/\sigma^2 \geq C \log(n/\delta)$ , where  $k\mu^2/\sigma^2$  represents the signal to noise ratio.

**Optimization** We perform gradient descent with fixed learning rate  $\eta$ ,

$$\mathbf{w}^{(t+1)} = \mathbf{w}(t) - \eta \nabla \mathcal{L}(\mathbf{w}(t)) = \mathbf{w}(t) - \eta \sum_i \ell'(\mathbf{w}^\top \mathbf{y}_i \mathbf{x}_i) \cdot \mathbf{y}_i \mathbf{x}_i. \quad (8)$$

**Asymptotic Solution** [46] For sufficiently small learning rate  $\eta$ , and (bounded) starting point  $\mathbf{w}(0)$ ,

$$\mathbf{w}(t) = \hat{\mathbf{w}} \log t + \rho(t), \quad (9)$$

where  $\hat{\mathbf{w}}$  is the solution to the hard margin SVM:

$$\hat{\mathbf{w}} = \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^d} \|\mathbf{w}\|_2^2 \quad s.t. \quad \mathbf{w}^\top \mathbf{y}_i \mathbf{x}_i \geq 1, \quad (10)$$

## A.2 Learning Stage

For the stage 1, we consider that we train the model for a maximum of  $T$  epochs (until we achieve 100% accuracy on the first training dataset  $\mathcal{S}_A$  with margin greater than 1). This means that the learned weight vectors are close to, but have not converged to the max margin solution. The solution at the end of  $T$  epochs is given by  $\mathbf{w}_A(t)$ . At sufficiently large  $T$ , we have:

$$\begin{aligned} \mathbf{w}_A(T) &= \hat{\mathbf{w}}_A \log T + \rho_A(T) \\ \mathbf{w}_A(T)^\top \mathbf{y}_i \mathbf{x}_i &\geq 1 \quad \forall (\mathbf{x}_i, \mathbf{y}_i) \in \mathcal{S}_A \end{aligned} \quad (11)$$

**Lemma 1** (Bounded Weights). *With probability at least  $1 - \delta$ , there exists a bounded time  $t$  beyond which the model classifies all training points correctly.*

*Proof.* From Lemma 5, we know that the dataset is separable with probability at least  $1 - \delta$ . This means that the max-margin solution (or the SVM solution) for the dataset classifies all training points correctly. From the analysis in Soudry et al. [46], we know that:

$$\begin{aligned} \lim_{t \rightarrow \infty} \frac{\mathbf{w}_A(t)}{\|\mathbf{w}_A(t)\|_2} &= \frac{\hat{\mathbf{w}}_A}{\|\hat{\mathbf{w}}_A\|_2}, \\ \lim_{t \rightarrow \infty} \mathbf{y}_i \mathbf{x}_i^\top \frac{\mathbf{w}_A(t)}{\|\mathbf{w}_A(t)\|_2} &= \mathbf{y}_i \mathbf{x}_i^\top \frac{\hat{\mathbf{w}}_A}{\|\hat{\mathbf{w}}_A\|_2} \end{aligned} \quad (12)$$

Then it directly follows that, for  $\epsilon > 0$ , there exists bounded time  $T > 0$  such that  $\forall t > T$ ,

$$\left| \mathbf{y}_i \mathbf{x}_i^\top \frac{\mathbf{w}_A(t)}{\|\mathbf{w}_A(t)\|_2} - \mathbf{y}_i \mathbf{x}_i^\top \frac{\hat{\mathbf{w}}_A}{\|\hat{\mathbf{w}}_A\|_2} \right| < \epsilon$$

This concludes that there exists a time  $T$  for which the sign of the prediction of both the max-margin solution and the learnt solution will be the same, implying correct prediction for every example in the training set for a bounded time solution.  $\square$

## A.3 Forgetting Stage

We initialize the weights for second stage of training with  $\mathbf{w}_A(T)$  from first training stage, and then train on  $\mathcal{S}_B$  to minimize the empirical loss using gradient descent (Equation 3). Assume that the dataset is balanced in the class labels,  $|\mathcal{S}_B| = n$ . Also, recall that  $\mathcal{S}_B$  does not contain mislabeled or rare examples from the same sub-group as in  $\mathcal{S}_A$ . For analyzing example forgetting in an asymptotic setting, we will directly borrow results from the analysis made by Chatterji and Long [9]. They prove a stronger result for the case where the dataset contains a fraction  $\eta$  of mislabeled examples. However, we will use the setting without label noise. The asymptotic result then builds on to the main theorem of the paper on intermediate time forgetting (Theorem 4).

**Transformations for Equivalence to Chatterji and Long [9]** In our setup, we consider a group structure where each distribution has a mean vector that is orthogonal to the others. We show the equivalence of the same to the data model studied in prior work [9].

Let us define  $\mu_1, \mu_2 \in \mathbb{R}^d$  as follows:

$$[\mu_1]_j = \begin{cases} \mu & \text{if } j \in \{1 \dots k\} \\ 0 & \text{o.w.} \end{cases}$$

Similarly, define  $\mu_2, \mu_3$  as well. For the equivalence condition, we are only concerned about the dataset split  $\mathcal{S}_B$  which is where the generalization bounds hold. Further, let  $z \sim \mathcal{N}(0, \sigma^2 \mathbf{I}_d)$ . Then,  $\mathbf{x}|\mathcal{X}_g \sim z + \mu_g$ .

We can now shift and rescale the axes such that the new origin is located at  $(\mu_1 + \mu_2)/2$ . Then, define  $\mu = (\mu_1 - \mu_2)/2$ . This results in the simplification that the mean of the examples sampled from  $\mathcal{D}_1$  is  $\mu$  and that from  $\mathcal{D}_2$  is  $-\mu$ . We can hence express  $\mathbf{x} = \mathbf{y}\mu + \mathbf{z}$ . This directly follows their model assumption where  $\mathbf{z} \in \mathbb{R}^d$  has each marginal sampled from a mean zero subgaussian distribution with subgaussian norm at most 1. In our case, each marginal is a Gaussian random variable with variance  $\sigma^2$ . Once again, we can rescale the axes such that  $\tilde{\mu} = \mu/\sigma$ . Now, our data model directly follows the data model discussed in prior work [9].

#### A.4 Asymptotic Analysis

First, we analyze infinite-time training case. We will extend the result from Chatterji and Long [9], Soudry et al. [46] to show that (i) mislabeled and rare examples are forgotten when the model is trained for long; and (ii) complex examples are not forgotten. First, recall the result used in Subsection A.2 for any bounded initialization for weights  $\mathbf{w}_B(0)$ ,

$$\mathbf{w}_B(t) = \hat{\mathbf{w}}_B \log t + \rho_B(t). \quad (13)$$

We will first provide a formal version of Theorem 1 which was informally stated in the main paper.

**Theorem 3** (Asymptotic Forgetting). *Under assumptions A.1, A.2, A.3, with probability  $1-\delta$ , fine-tuning for  $t \rightarrow \infty$  iterations on the second dataset  $\mathcal{S}_B$  produces a max-margin classifier  $\hat{\mathbf{w}}_B$  such that*

$$\begin{aligned} \mathbb{P}_{(\mathbf{x}_m, \mathbf{y}_m) \in \mathcal{S}_A} [\text{sign}(\hat{\mathbf{w}}_B \cdot \mathbf{x}_m) = \mathbf{y}_m] &\leq \exp(-c\|\tilde{\mu}\|_2^2/d), \\ \Phi(-1/C) &\leq \mathbb{P}_{(\mathbf{x}_r, \mathbf{y}_r) \in \mathcal{S}_A} [\text{sign}(\hat{\mathbf{w}}_B \cdot \mathbf{x}_r) \neq \mathbf{y}_r] \leq \Phi(1/C), \\ \mathbb{P}_{(\mathbf{x}_c, \mathbf{y}_c) \in \mathcal{S}_A} [\text{sign}(\hat{\mathbf{w}}_B \cdot \mathbf{x}_c) \neq \mathbf{y}_c] &\leq \exp(-c\|\tilde{\mu}/\lambda\|_2^2/d), \end{aligned}$$

for some absolute constant  $c > 0$  and  $(\mathbf{x}_m, \mathbf{y}_m) \in \mathcal{M}$ ,  $(\mathbf{x}_r, \mathbf{y}_r) \in \mathcal{R}$ ,  $(\mathbf{x}_c, \mathbf{y}_c) \in \mathcal{C}$  respectively.

The theorem implies that the probability that mislabeled examples from  $\mathcal{S}_A$  are classified with the given (incorrect) label tends to 0 if  $\|\mu\|_2 = \theta(d^\beta)$  for any  $\beta \in (1/4, 1/2]$ . Note that in our case, we have  $k$  dimensions of signal. Therefore, as long as  $k/d$  is a constant fraction, as the input dimensions  $d \rightarrow \infty$  the above holds. Examples in  $\mathcal{S}_A$  from distributions absent in  $\mathcal{S}_B$  are randomly classified.

*Intuition.* In the asymptotic case, the initial weights from first split training  $\mathbf{w}_A(T) = \mathcal{O}(\log T) \ll \mathbf{w}_B(t)$ , where  $t \rightarrow \infty$  in the limit of infinite training. As a consequence, for any bounded initialization, the model weights converge to the minimum norm solution (SVM) solution from Soudry et al. [46]. We use results from [9] who study the setting of a binary classification problem with noisy label fraction  $\eta$ . In our case, since the second-split is assumed to have only clean samples,  $\eta = 0$ . The final step is to adapt our data generating process to the format used in Chatterji and Long [9].

*Proof.* Recall the transformations described in Section A.3. Let  $(\mathbf{x}_m, \mathbf{y}_m), (\mathbf{x}_r, \mathbf{y}_r), (\mathbf{x}_c, \mathbf{y}_c)$  represent any point from  $\mathcal{S}_A$  which belongs to mislabeled set  $\mathcal{M}$ , rare set  $\mathcal{R}$  and complex set  $\mathcal{C}$ . The important thing to note in the analysis that follows is that each of these examples is independent of the samples in  $\mathcal{S}_B$ . Hence, the probability of correctly predicting on them is same as that of correctly predicting on a population sample, in the limit of infinite training (when initialization does not matter and all models converge to the same solution).

$$\mathbb{P}_{(\mathbf{x}, \mathbf{y}) \sim \mathcal{D}} [\text{sign}(\mathbf{w} \cdot \mathbf{x}) \neq \mathbf{y}] = \mathbb{P}_{(\mathbf{x}, \mathbf{y}) \sim \mathcal{D}} [(\mathbf{y}\mathbf{w} \cdot \mathbf{x}) < 0],$$

We will now analyze the probability of correct prediction of mislabeled, rare, and complex examples separately.

**Mislabeled Examples** Since  $\mathcal{S}_B$  is separable, in the limit of infinite-training time, the classifier correctly predicts all the examples in the dataset. We will use this fact to show that it assigns the opposite label to the mislabeled sample in set  $\mathcal{S}_A$  with high probability. Then, we denote ‘failure’ as the event that the mislabeled samples is still predicted with label  $\mathbf{y}_m$  at infinite-time training.

$$\begin{aligned}\mathbb{P}_{(\mathbf{x}, \mathbf{y}) \sim \mathcal{D}_1} [(\mathbf{y}\mathbf{w} \cdot \mathbf{x}) < 0] &= 1 - \mathbb{P}[(\mathbf{y}_m \mathbf{w} \cdot \mathbf{x}_m) < 0] \\ &= \mathbb{P}[(\mathbf{y}_m \mathbf{w} \cdot \mathbf{x}_m) > 0] \\ &= \mathbb{P}[\text{sign}(\mathbf{w} \cdot \mathbf{x}_m) = \mathbf{y}_m],\end{aligned}$$

Then, we can directly borrow the result from Chatterji and Long [9] (Theorem 4) to find that

$$\begin{aligned}\mathbb{P}_{(\mathbf{x}, \mathbf{y}) \sim \mathcal{D}_1} [(\mathbf{y}\mathbf{w} \cdot \mathbf{x}) > 0] &\leq \exp\left(-c \frac{\|\tilde{\boldsymbol{\mu}}\|^4}{d}\right) \\ \mathbb{P}[\text{sign}(\mathbf{w} \cdot \mathbf{x}_m) = \mathbf{y}_m] &\leq \exp\left(-c \frac{\|\tilde{\boldsymbol{\mu}}\|^4}{d}\right),\end{aligned}\tag{14}$$

**Rare Examples** Without loss of generality, we may assume that the correct label for the rare example  $\mathbf{y}_r = 1$ .

$$\begin{aligned}\mathbb{P}[(\mathbf{y}_r \mathbf{w} \cdot \mathbf{x}_r) < 0] &= \mathbb{P}[(\mathbf{y}_r \mathbf{w} \cdot (\mathbf{x}_r - \boldsymbol{\mu}_r) < -\mathbf{y}_r \mathbf{w} \cdot \boldsymbol{\mu}_r] \\ &= \mathbb{P}[(\mathbf{w} \cdot (\mathbf{x}_r - \boldsymbol{\mu}_r) < -\mathbf{w} \cdot \boldsymbol{\mu}_r], \quad (\text{since } \mathbf{y}_r = 1) \\ &= \mathbb{P}[(\mathbf{w} \cdot \mathbf{z}_r < -\mathbf{w} \cdot \boldsymbol{\mu}_r] \\ &= \Phi\left(\frac{-\mathbf{w} \cdot \boldsymbol{\mu}_r}{\sigma \|\mathbf{w}\|_2}\right),\end{aligned}\tag{15}$$

where  $\Phi$  is the Gaussian CDF. In the last step we use the fact that  $\mathbf{z}_r \sim \mathcal{N}(0, \sigma^2 \mathbf{I}_d)$ . Therefore its dot product with the vector  $\mathbf{w}$  results with a summation of  $d$  Gaussian vectors, each with mean 0 and variance  $\sigma^2 [\mathbf{w}]_i^2 \forall i \in [d]$ . Now what remains is to prove that the value at which we want to calculate the CDF is close to 0, so that the probability of the misclassification is close to 0.5.

Recall from the analysis in Soudry et al. [46] that asymptotically the model converges to the max-margin separator that is comprised of a weighted sum of the support vectors of the dataset (let us call this set  $\mathcal{V}_B$ ). This means that the final weights of the model  $\mathbf{w}$  is a combination of datapoints from the first two majority distributions  $\mathcal{D}_1, \mathcal{D}_2$ .

$$\mathbf{w} = \sum_{i \in \mathcal{V}_B} \alpha_i \mathbf{x}_i,\tag{16}$$

$$\boldsymbol{\mu}_r \cdot \mathbf{w} = \sum_{i \in \mathcal{V}_B} \alpha_i (\boldsymbol{\mu}_r \cdot \mathbf{x}_i)\tag{17}$$

Since  $\boldsymbol{\mu}_r \cdot \boldsymbol{\mu}_i = 0$ , we have that  $(\boldsymbol{\mu}_r \cdot \mathbf{x}_i) = \boldsymbol{\mu}_r \cdot \mathbf{z}_i + 0$ , which is a mean zero random variable.

$$\boldsymbol{\mu}_r \cdot \mathbf{w} = \sum_{i \in \mathcal{V}_B} \alpha_i (\boldsymbol{\mu}_r \cdot \mathbf{z}_i),\tag{18}$$

$$= \sum_{i \in \mathcal{V}_B} \alpha_i \mu \sum_{2k+1 \leq j \leq 3k} [\mathbf{z}_i]_j,\tag{19}$$

$$= \xi,\tag{20}$$

$$(21)$$

where  $\xi \sim \mathcal{N}(0, k\mu^2 \sigma^2 \sum_{i \in \mathcal{V}_B} \alpha_i^2)$ . Also,  $\|\mathbf{w}\|_2^2 = \sum_{i \in \mathcal{V}_B} \alpha_i^2 \mathbf{x}_i^2$ . However,  $\mathbf{x}_i = \boldsymbol{\mu}_i + \mathbf{z}_i$ . From Lemma 4, we know that with probability greater than  $1 - \delta/6$ , for every example,  $\frac{d\sigma^2}{2} \leq \|\mathbf{z}_i\|_2^2 \leq \frac{3d\sigma^2}{2}$ . By Young’s inequality for products:

$$\begin{aligned}
\|\mathbf{z}_i\|_2^2 &= \|\mathbf{x}_i - \boldsymbol{\mu}_i\|_2^2, \\
&\leq 2\|\mathbf{x}_i\|_2^2 + 2\|\boldsymbol{\mu}_i\|_2^2, \\
\|\mathbf{x}_i\|_2^2 &\geq \frac{1}{2}\|\mathbf{z}_i\|_2^2 - \|\boldsymbol{\mu}_i\|_2^2, \\
&\geq \frac{d\sigma^2}{4} - \|\boldsymbol{\mu}_i\|_2^2, \\
&\geq \frac{d\sigma^2}{8} \quad (\text{since } \frac{k\mu^2}{\sigma^2} < d/nC \text{ for large } C), \\
\|\mathbf{w}\|_2^2 &= \sum_{i \in \mathcal{V}_B} \alpha_i^2 \mathbf{x}_i^2 \geq \frac{d\sigma^2}{8} \sum_{i \in \mathcal{V}_B} \alpha_i^2.
\end{aligned}$$

Therefore,

$$\begin{aligned}
\eta &= \frac{-\xi}{\sigma \|\mathbf{w}\|_2} \sim \mathcal{N}\left(0, \frac{k\mu^2}{d\sigma^2/8}\right), \\
\eta &\sim \mathcal{N}\left(0, \frac{1}{nC}\right), \quad (\text{since } \frac{k\mu^2}{\sigma^2} < d/nC \text{ for large } C), \\
\mathbb{P}[(\mathbf{y}_r \mathbf{w} \cdot \mathbf{x}_r) < 0] &= \Phi(\eta)
\end{aligned} \tag{22}$$

Using Gaussian tail bound on  $\eta$ , with probability at least  $1 - \delta$ ,  $\eta \leq \sqrt{\frac{\log(1/\delta)}{nC}}$ . Therefore,

$$\mathbb{P}[(\mathbf{y}_r \mathbf{w} \cdot \mathbf{x}_r) < 0] \leq \Phi\left(\sqrt{\frac{\log(1/\delta)}{nC}}\right) \leq \Phi\left(\frac{1}{C}\right) \approx 0.5.$$

In the last step, we used the assumption (A.2) that  $n \geq C \log(1/\delta)$  for some large constant  $C$ .

*Remark:* This analysis highlights that the meaning of being ‘forgotten’ for rare examples is to predict randomly. However, in case of mislabeled examples, the predicted label of the example approaches its ‘true’ label, irrespective of its training label.

**Complex Examples** The analysis for complex examples follows directly from the analysis in the case of mislabeled examples. The only difference is the magnitude of the signal to noise ratio within the group. Recall that complex examples are also sampled from majority groups. Therefore, the probability that the complex example  $(\mathbf{x}, \mathbf{y})$  is misclassified at the end of training for infinite time is given by:

$$\mathbb{P}[\mathbf{y} \mathbf{x}^\top \mathbf{w} < 0] \leq \exp\left(-c \frac{\|\tilde{\boldsymbol{\mu}}/\lambda\|^4}{d}\right) \tag{23}$$

This approaches 0, indicating perfect classification of complex examples from  $\mathcal{S}_A$ . Note that this is a complimentary case where the second split only has examples from the complex distribution.

□

## A.5 Intermediate Time Analysis

From the analysis in Section A.4, we find that all the mislabeled and rare examples are forgotten by the time we train for  $t \rightarrow \infty$  iterations. Since examples from complex subgroups are not forgotten even at infinite time training, we skip analysis for those examples in this subsection. Our goal is to show that there exists a time  $T'$  such that with high probability, the model forgets all the mislabeled examples, but still correctly classifies all the rare examples. To show this, we will track the accuracy of inputs in the first data split, as we train on the examples in the second split.

The model output for any sample  $(\mathbf{x}_i, \mathbf{y}_i) \in \mathcal{S}_A$  is given by  $\mathbf{w}(t)^\top \mathbf{x}_i$ , and the prediction is considered to be correct if  $\text{sign}(\mathbf{w}(t)^\top \mathbf{x}_i) = \mathbf{y}_i$ , or if  $\mathbf{w}(t)^\top \mathbf{x}_i \mathbf{y}_i > 0$ . From hereon, we will use the notation  $\mathbf{a}_i^t = \mathbf{w}(t)^\top \mathbf{x}_i \mathbf{y}_i$ .



Recall that we considered that there is one example from both the mislabeled and rare example category in  $\mathcal{S}_A$ . We will denote these data points by  $(\mathbf{x}_m, \mathbf{y}_m) \in \mathcal{M}$ , and  $(\mathbf{x}_r, \mathbf{y}_r) \in \mathcal{R}$  respectively. Without loss of generality, we will assume that  $(\mathbf{x}_m, -\mathbf{y}_m)$  was sampled from  $\mathcal{D}_1$  in the mixture of distributions  $\mathcal{D}$ . All the examples in the second split  $\mathcal{S}_B$  sampled from the same distribution are given by  $\mathcal{S}_{B,1} \subset \mathcal{S}_B$ . The remaining examples are in the set denoted by  $\mathcal{S}_{B,-1} \subset \mathcal{S}_B$ .

From the learning time training dynamics, we know that the initialization of the weights for the second round of training is given by:

$$\mathbf{w}_B(0) = \hat{\mathbf{w}}_A \log T + \rho_A(T). \quad (24)$$

Now, from the representer theorem, we can decompose the model weights at any iteration of second-split training as follows,

$$\mathbf{w}_B(t) = \hat{\mathbf{w}}_A \log T + \rho_A(T) + \sum_{j \in \mathcal{S}_B} \beta_j \mathbf{y}_j \mathbf{x}_j. \quad (25)$$

Note that we introduce an additional term  $\mathbf{y}_j$  in the decomposition using representer theorem, but this can be done without loss of generality since  $\mathbf{y}_j \in \{-1, +1\}$ . This helps us in ensuring that each  $\beta_j$  is non-negative as shown in Lemma 6.

From Lemma 8, we know that there exists a bounded time  $T'$  when the mislabeled example prediction flips. We will denote  $\mu$  as the coordinate-wise mean for the  $k$ -signal dimensions in the vector  $\mu_g$  in the discussion that follows. Let us define  $\Delta_t = \frac{\sum_{j \in \mathcal{S}_{B,1}} \beta_j(t)}{\sum_{j \in \mathcal{S}_B} \beta_j(t)}$ , and  $\Delta = \max_t \Delta_t$ . For a symmetric distribution with two majority subgroups with the opposite label, we would expect this value to be close to 0.5. Now, we present formal version of Theorem 2 from the main paper.

**Theorem 4** (Intermediate-Time Forgetting). *Under the distribution outlined in Appendix A.1 with assumptions A.1, A.2, A.3, whenever  $\mathcal{S}_A$  is separable, when fine-tuning on the second dataset split  $\mathcal{S}_B$  with sufficiently small learning rate, there exists some bounded time  $T'$  when*

1.  $\mathbb{P}_{(\mathbf{x}_m, \mathbf{y}_m) \in \mathcal{S}_A} [\mathbf{y}_m \neq \mathbf{w}(T') \cdot \mathbf{x}_m] \geq 1 - c_0 \exp\left(\frac{-k^2 \mu^2 \Delta^2}{c \sigma^2 d}\right) - c_1 \exp(-cd)$
2.  $\mathbb{P}_{(\mathbf{x}_r, \mathbf{y}_r) \in \mathcal{S}_A} [\mathbf{y}_r = \mathbf{w}(T') \cdot \mathbf{x}_r] \geq 1 - c_0 \exp\left(\frac{-k^2 \mu^2 \Delta^2}{c \sigma^2 d}\right) - c_1 \exp(-cd)$

for absolute constants  $c_0, c_1, c$  and  $(\mathbf{x}_m, \mathbf{y}_m) \in \mathcal{M}, (\mathbf{x}_r, \mathbf{y}_r) \in \mathcal{R}$  respectively.

If the fraction of dimensions that contain the signal,  $k/d$ , is considered fixed then both the first and second term above exponentially decay with a factor of  $d$ . This suggests that increasing overparametrization leads to a higher likelihood of the phenomenon of intermediate time forgetting—there exists an epoch when the prediction of mislabeled example is flipped but the rare examples are still correctly predicted with high probability.

In what follows, the key idea is to show that at a time when the prediction of the mislabeled example is incorrect with high probability, the predictions of the rare examples is correct with high probability.

*Proof.* From Lemma 7, we know that  $(\mathbf{x}_m, \mathbf{y}_m)$  and  $(\mathbf{x}_r, \mathbf{y}_r)$  are support vectors for  $\mathcal{S}_A$ . Hence,  $\mathbf{y}_m \mathbf{x}_m^\top \hat{\mathbf{w}}_A = \mathbf{y}_r \mathbf{x}_r^\top \hat{\mathbf{w}}_A = 1$ . We will use this fact to find the distribution for  $\mathbf{a}_m^t, \mathbf{a}_r^t$ .

Let  $[\mathbf{x}_j]_i$  denote the  $i^{\text{th}}$  dimension of the  $j^{\text{th}}$  datapoint in the second split. Recall that the second split comprises of two majority subgroups.  $k$  dimensions of the input vector contain the true signal for class prediction. The dimensions  $\{1 \dots k\}, \{k+1 \dots 2k\}, \{2k+1 \dots 3k\}$  are the predictive dimensions for the majority group 1, 2 and that for the rare example from dataset  $\mathcal{S}_A$ . Therefore,  $[\mathbf{x}_j]_i = \mu + [\mathbf{z}_j]_i$ , if  $i$  is in the predictive dimensions, otherwise  $[\mathbf{x}_j]_i = [\mathbf{z}_j]_i$  where  $[\mathbf{z}_j]_i \sim \mathcal{N}(0, \sigma^2)$ . To make notation simple, we will refer to  $\beta_j(t)$  by using  $\beta_j$ .

*Remark:* We do not perform input transformation to prove the following results.

**Mislabeled Example** The prediction on the mislabeled point times the given label can be written as:

$$\begin{aligned} \mathbf{a}_m^t &= \mathbf{y}_m \mathbf{x}_m^\top \mathbf{w}_B(t) = \mathbf{y}_m \mathbf{x}_m^\top \left( \hat{\mathbf{w}}_A \log T + \rho_A(T) + \sum_{j \in \mathcal{S}_B} \beta_j \mathbf{y}_j \mathbf{x}_j \right), \\ &= \log T + c_m + \sum_{j \in \mathcal{S}_B} \beta_j (\mathbf{y}_m \mathbf{x}_m^\top \mathbf{y}_j \mathbf{x}_j) \quad (\text{using Lemma 7}), \end{aligned} \quad (26)$$

where  $c_m = \mathbf{y}_m \mathbf{x}_m^\top \rho_A(T)$  is a residual term that does not continue grow with  $T$  (Theorem 9 [46]).

Without loss of generality, assume that the mislabeled example  $\mathbf{x}_m$  belongs to majority group 1 (true label = 1), and was originally labeled such that  $\mathbf{y}_m = -1$  in the first-split dataset.

$$\begin{aligned} \sum_{j \in \mathcal{S}_B} \beta_j (\mathbf{y}_m \mathbf{x}_m^\top \mathbf{y}_j \mathbf{x}_j) &= \sum_{j \in \mathcal{S}_{B,1}} \beta_j (\mathbf{y}_m \mathbf{x}_m^\top \mathbf{y}_j \mathbf{x}_j) + \sum_{j \in \mathcal{S}_{B,-1}} \beta_j (\mathbf{y}_m \mathbf{x}_m^\top \mathbf{y}_j \mathbf{x}_j) \\ &= - \sum_{j \in \mathcal{S}_{B,1}} \beta_j (\mathbf{x}_m^\top \mathbf{x}_j) + \sum_{j \in \mathcal{S}_{B,-1}} \beta_j (\mathbf{x}_m^\top \mathbf{x}_j). \end{aligned} \quad (27)$$

Now we can use the distribution properties of the dataset to further simplify per dimension and aggregate the sum across all examples.

$$\sum_{j \in \mathcal{S}_{B,1}} \beta_j (\mathbf{x}_m^\top \mathbf{x}_j) = \mathbf{x}_m^\top \sum_{j \in \mathcal{S}_{B,1}} \beta_j \mathbf{x}_j = \mathbf{x}_m^\top \mathbf{x}_{\mathcal{S}_{B,1}}, \quad (28)$$

where  $[\mathbf{x}_{\mathcal{S}_{B,1}}]_i = \mu \mathbb{1}(i \in \{1 \dots k\}) \sum_{j \in \mathcal{S}_{B,1}} \beta_j + [\mathbf{z}_{\mathcal{S}_{B,1}}]_i$ , where  $[\mathbf{z}_{\mathcal{S}_{B,1}}]_i \sim \mathcal{N}(0, \sigma^2 \sum_{j \in \mathcal{S}_{B,1}} \beta_j^2)$ . Now, we will add up the dot product dimension wise. Let us call  $B_1 = \sum_{j \in \mathcal{S}_{B,1}} \beta_j$  and  $B_1^v = \sum_{j \in \mathcal{S}_{B,1}} \beta_j^2$ . Then  $[\mathbf{x}_{\mathcal{S}_{B,1}}]_i \sim \mathcal{N}(\mu \mathbb{1}(i \leq k) \cdot B_1, \sigma^2 B_1^v)$ .

$$\begin{aligned} \mathbf{x}_m^\top \mathbf{x}_{\mathcal{S}_{B,1}} &= \mu \sum_k [\mathbf{z}_{\mathcal{S}_{B,1}}]_i + \mu B_1 \sum_k [\mathbf{z}_m]_i + k \cdot \mu^2 \cdot B_1 + \sum_d [\mathbf{z}_{\mathcal{S}_{B,1}}]_i \cdot [\mathbf{z}_m]_i \\ &= \mu \cdot \alpha_1 + \sum_d [\mathbf{z}_{\mathcal{S}_{B,1}}]_i \cdot [\mathbf{z}_m]_i \\ \mathbf{x}_m^\top \mathbf{x}_{\mathcal{S}_{B,-1}} &= \mu \cdot \alpha_2 + \sum_d [\mathbf{z}_{\mathcal{S}_{B,-1}}]_i \cdot [\mathbf{z}_m]_i \end{aligned} \quad (29)$$

where  $\alpha_1 \sim \mathcal{N}(k \cdot \mu \cdot B_1, k \cdot \sigma^2 (B_1^v + B_1^2))$  and  $\alpha_2 \sim \mathcal{N}(0, k \cdot \sigma^2 (B_2^v + B_2^2))$ . Notice that in the first step, the first three terms are independent of each other and can be added directly to obtain a new random variable using the independence condition, however, the last term is dependent on the first two. In the final step, we also add the terms for the dot product corresponding to the opposite group.

This gives us the overall expression as follows:

$$\begin{aligned} \sum_{j \in \mathcal{S}_B} \beta_j (\mathbf{y}_m \mathbf{x}_m^\top \mathbf{y}_j \mathbf{x}_j) &= \mu \cdot (\alpha_2 - \alpha_1) + \sum_d ([\mathbf{z}_{\mathcal{S}_{B,-1}}]_i - [\mathbf{z}_{\mathcal{S}_{B,1}}]_i) \cdot [\mathbf{z}_m]_i, \\ &= \mu \cdot \alpha_m + \sum_d [\mathbf{z}_{\mathcal{S}_B}]_i \cdot [\mathbf{z}_m]_i - k \cdot \mu^2 \cdot B_1, \end{aligned} \quad (30)$$

where  $\alpha_m \sim \mathcal{N}(0, k \cdot \sigma^2 (B^v + B_1^2 + B_2^2))$  and  $[\mathbf{z}_{\mathcal{S}_B}]_i \sim \mathcal{N}(0, \sigma^2 B^v)$ , since  $(B_1^v + B_2^v = \sum_{j \in \mathcal{S}_B} \beta_j^2 = B^v)$ .

$$\mathbf{a}_m^t = \log T + c_m + \xi_m - k \cdot \mu^2 \cdot B_1; \text{ s.t. } \xi_m = \mu \cdot \alpha_m + \sum_d [\mathbf{z}_{\mathcal{S}_B}]_i \cdot [\mathbf{z}_m]_j. \quad (31)$$

**Rare Example** Following the analysis of the mislabeled example, we can similarly find an expression for  $\mathbf{a}_r^t$  on the rare example as follows:

$$\begin{aligned}\mathbf{a}_r^t &= \mathbf{y}_r \mathbf{x}_r^\top \mathbf{w}_B(t) = \mathbf{y}_r \mathbf{x}_r^\top (\hat{\mathbf{w}}_A \log T + \rho_A(T) + \sum_{j \in \mathcal{S}_B} \beta_j \mathbf{y}_j \mathbf{x}_j), \\ &= \log T + c_r + \sum_{j \in \mathcal{S}_B} \beta_j (\mathbf{y}_r \mathbf{x}_r^\top \mathbf{y}_j \mathbf{x}_j) \quad (\text{using Lemma 7}),\end{aligned}\tag{32}$$

Following the same procedure as for mislabeled example, we get the overall expression as follows:

$$\sum_{j \in \mathcal{S}_B} \beta_j (\mathbf{y}_r \mathbf{x}_r^\top \mathbf{y}_j \mathbf{x}_j) = \mu \cdot \alpha_r + \sum_d [\mathbf{z}_{\mathcal{S}_B}]_i \cdot [\mathbf{z}_r]_i,\tag{33}$$

where  $\alpha_r \sim \mathcal{N}(0, k \cdot \sigma^2 (B^v + B_1^2 + B_2^2))$  and  $[\mathbf{z}_{\mathcal{S}_B}]_i \sim \mathcal{N}(0, \sigma^2 B^v)$ .

$$\mathbf{a}_r^t = \log T + c_r + \xi_r, \text{ s.t. } \xi_r = \mu \cdot \alpha_r + \sum_d [\mathbf{z}_{\mathcal{S}_B}]_i \cdot [\mathbf{z}_r]_i.\tag{34}$$

**Combining both cases** Recall that  $\beta_j > 0$  for all  $j$ . Therefore,  $B^2 = (\sum_{j \in \mathcal{S}_B} \beta_j)^2 > B^v = \sum_{j \in \mathcal{S}_B} \beta_j^2$ . Based on the analysis of Soudry et al. [46], we know that  $\mathbf{w}_B(t)$  grows as fast as  $O(\log t)$ . Therefore, both  $B_1 = O(\log t)$ , and  $B = O(\log t)$  (must grow at most as fast as that). From the problem definition,  $B_1 \leq \Delta B$ , where  $\Delta \in [0, 1]$ .

Our goal now is to find an epoch  $t$  during the second-split training when the rare example is correctly classified with high probability, and the mislabeled example is incorrectly classified with high probability. The maxima is achieved when  $\mathbf{a}_r^t \approx -\mathbf{a}_m^t$ . We assume that the step sizes are sufficiently small such that we can achieve this condition. Moreover,  $c_r, c_m \ll \log T$ . Hence, the condition is to find  $t$  such that  $\xi_r + \log T > 0$  and  $\xi_m + \log T < k \cdot \mu^2 \cdot B_1$  with high probability. By symmetry, we must find  $\mathbb{P}(|\xi_m| < k \cdot \mu^2 \cdot B_1/2)$ .

$$\mathbb{P}\left(|\xi_m| > k \cdot \mu^2 \cdot \frac{B_1}{2}\right) \leq \mathbb{P}\left(|\xi_1| > k \cdot \mu^2 \cdot \frac{B_1}{4}\right) + \mathbb{P}\left(|\xi_2| > k \cdot \mu^2 \cdot \frac{B_1}{4}\right),\tag{35}$$

where  $\xi_1 = \mu \cdot \alpha_m \sim \mathcal{N}(0, \mu^2 \sigma^2 k (B^v + B_1^2 + B_2^2))$  and  $\xi_2 = \sum_d [\mathbf{z}_{\mathcal{S}_B}]_i \cdot [\mathbf{z}_r]_i$ . For the first term, we can directly use the Gaussian tail bound using Chernoff method.

$$\begin{aligned}\mathbb{P}\left(|\xi_1| > k \cdot \mu^2 \cdot \frac{B_1}{4}\right) &\leq 2 \exp \frac{-k^2 \mu^4 B_1^2}{32 \mu^2 \sigma^2 d (B^v + B_1^2 + B_2^2)}, \\ &\leq 2 \exp \frac{-k^2 \mu^4 B_1^2}{32 \mu^2 \sigma^2 d (2B^2)}, \\ &\leq 2 \exp \frac{-k \mu^2 \Delta^2}{c_1 \sigma^2}.\end{aligned}\tag{36}$$

Now, for the second term, using Lemma 3 we have that.

$$\begin{aligned}\mathbb{P}\left(|\xi_2| > k \cdot \mu^2 \cdot \frac{B_1}{4}\right) &\leq 2 \exp \frac{-k^2 \mu^4 B_1^2}{c_2 \sigma^4 d B^v} + c_1 \exp(-cd), \\ &\leq 2 \exp \frac{-k^2 \mu^2 \Delta}{c_2 \sigma^2 d} + c_1 \exp(-cd),\end{aligned}\tag{37}$$

since  $\frac{\mu}{\sigma} > 1$ . Finally, combining both the cases we have that

$$\mathbb{P}\left(|\xi_m| > k \cdot \mu^2 \cdot \frac{B_1}{2}\right) \leq c_0 \exp \frac{-k^2 \mu^2 \Delta^2}{c \sigma^2 d} + c_1 \exp(-cd),\tag{38}$$

This shows that as the dimensionality of the dataset increases, the likelihood of prediction on mislabeled examples being flipped while the rare examples retain their prediction increases exponentially. This concludes the proof of Theorem 4.  $\square$

## A.6 Concentration Inequalities and Additional Lemmas

To make our work self-contained, we supplement the reader with additional Lemmas and Theorems that are helpful tools for proving the theorems in this work. We restate versions of the Hoeffding and Bernstein inequalities as in Chatterji and Long [9].

**Lemma 2** (Soudry et al. [46], Theorem 3). *For any linearly separable dataset  $\mathcal{S}_A$  and for all small enough step-sizes  $\alpha$ , we have*

$$\frac{\mathbf{w}_A}{\|\mathbf{w}_A\|} = \lim_{t \rightarrow \infty} \frac{\mathbf{w}^{(t)}}{\|\mathbf{w}^{(t)}\|}.$$

**Theorem 5** (General Hoeffding's Inequality). *Let  $\theta_1, \dots, \theta_m$  be independent mean-zero sub-Gaussian random variables and  $a = (a_1, \dots, a_m) \in \mathbb{R}^m$ . Then, for every  $t > 0$ , we have*

$$\mathbb{P} \left[ \left| \sum_{i=1}^m a_i \theta_i \right| \geq t \right] \leq 2 \exp \left( \frac{-c_0 t^2}{K^2 \|a\|_2^2} \right),$$

where  $K = \max_i \|\theta_i\|_{\psi_2}$  and  $c$  is an absolute constant.

In our case, since we deal with Gaussian random variables,  $\|\theta\|_{\psi_2}$  (sub-Gaussian norm) is same as the variance of the random variable. That is,  $\theta \sim \mathcal{N}(0, \sigma^2) \implies \|\theta\|_{\psi_2} = K = \sigma$ .

**Theorem 6** (Bernstein Inequality). *For independent mean-zero sub-exponential random variables  $\theta_1, \dots, \theta_m$ , for every  $t > 0$ , we have*

$$\mathbb{P} \left[ \left| \sum_{i=1}^m \theta_i \right| \geq t \right] \leq 2 \exp \left( -c_1 \min \left\{ \frac{t^2}{\sum_{i=1}^m \|\theta_i\|_{\psi_1}^2}, \frac{t}{\max_i \|\theta_i\|_{\psi_1}} \right\} \right),$$

where  $c_1$  is an absolute constant.

In our case, let  $\mathbf{x}_i \sim \mathcal{N}(0, \sigma^2)$  be independent Gaussian random variables, then for  $Z = \sum_i^d X_i^2$ ,

$$\mathbb{P}(|Z - d\sigma^2| \geq t) \leq 2 \exp \left( -\frac{c_1}{\sigma^2} \min \left\{ \frac{t^2}{d\sigma^2}, t \right\} \right).$$

**Lemma 3** (Gaussian Product). *Let  $\mathbf{z}_1 \sim \mathcal{N}(0, \sigma_1^2 \mathbf{I}_d)$ ,  $\mathbf{z}_2 \sim \mathcal{N}(0, \sigma_2^2 \mathbf{I}_d)$  be independent multivariate Gaussian random variables. Then,*

$$\mathbb{P} [|\mathbf{z}_1 \cdot \mathbf{z}_2| > t] \leq 2 \exp \left( \frac{-c t^2}{d\sigma_1^2 \sigma_2^2} \right) + 2 \exp(-c_0 d)$$

*Proof.* The proof of this lemma is a simplified version of the proof for Lemma 20 in Chatterji and Long [9]. First, consider  $\mathbf{z}_2$  to be fixed, and only  $\mathbf{z}_1$  to be random. Then from Theorem 5 we have

$$\mathbb{P} \left[ \left| \sum_{i=1}^d [\mathbf{z}_2]_i \cdot [\mathbf{z}_1]_i \right| \geq t \right] \leq 2 \exp \left( \frac{-c t^2}{\sigma_1^2 \|\mathbf{z}_2\|_2^2} \right).$$

Also, adapting Theorem 6 such that  $Z = \sum_i^d X_i^2 = \|\mathbf{z}_2\|_2^2$ , and setting  $t = d\sigma^2$

$$\begin{aligned} \mathbb{P} \left( \left| \|\mathbf{z}_2\|_2^2 - d\sigma_2^2 \right| \geq t \right) &\leq 2 \exp \left( -\frac{c_0}{\sigma_2^2} \min \left\{ \frac{t^2}{d\sigma_2^2}, t \right\} \right), \\ \mathbb{P} \left( \|\mathbf{z}_2\|_2^2 \geq 2d\sigma_2^2 \right) &\leq 2 \exp(-c_0 d). \end{aligned}$$

Coming back to the initial problem and again considering both  $\mathbf{z}_1, \mathbf{z}_2$  to be random variables, we have

$$\mathbb{P} [|\mathbf{z}_1 \cdot \mathbf{z}_2| \geq t] \leq \mathbb{P} [|\mathbf{z}_1 \cdot \mathbf{z}_2| \geq t \mid \|\mathbf{z}_2\|_2^2 \leq 2d\sigma_2^2] + \mathbb{P} [\|\mathbf{z}_2\|_2^2 > 2d\sigma_2^2], \quad (39)$$

$$\leq 2 \exp \left( \frac{-c t^2}{d\sigma_1^2 \sigma_2^2} \right) + 2 \exp(-c_0 d). \quad (40)$$

□

**Corollary 1** (Chatterji and Long [9], Lemma 20). *There is a  $c \geq 1$  such that, for all large enough  $C$ , with probability at least  $1 - \delta/6$ , for all  $i \neq j \in [n]$ ,*

$$|\mathbf{z}_i \cdot \mathbf{z}_j| < c \left( \sigma^2 \sqrt{d \log(n/\delta)} \right).$$

**Lemma 4** (Gaussian Square). *There is a  $c \geq 1$  such that, for all large enough  $C$ , with probability at least  $1 - \delta/6$ , for all  $k \in [n]$ ,*

$$\frac{d\sigma^2}{2} \leq \|\mathbf{z}_k\|_2^2 \leq \frac{3d\sigma^2}{2}.$$

*Proof.* Recall that  $\mathbf{x}_k = \boldsymbol{\mu}_k + \mathbf{z}_k$ , where  $\mathbf{z}_k \sim \mathcal{N}(0, \sigma^2 \mathbf{I}_d)$ . Then,

$$\|\mathbf{z}_k\|_2^2 = \sum_i^d [\mathbf{z}_k]_i^2 = Z$$

Adapting Theorem 6, and setting  $t = \lambda d\sigma^2$  with  $0 < \lambda < 1$

$$\begin{aligned} \mathbb{P} \left( \left| \|\mathbf{z}_2\|_2^2 - d\sigma^2 \right| \geq t \right) &\leq 2 \exp \left( -\frac{c_1}{\sigma^2} \min \left\{ \frac{t^2}{d\sigma^2}, t \right\} \right), \\ \mathbb{P} \left( \left| \|\mathbf{z}_2\|_2^2 - d\sigma^2 \right| \geq \lambda d\sigma^2 \right) &= 2 \exp \left( -\frac{c_1}{\sigma^2} \min \left\{ \lambda^2 d\sigma^2, \lambda d\sigma^2 \right\} \right), \\ &= 2 \exp \left( -c_1 \lambda^2 d \right). \quad (\text{since } 0 < \lambda < 1) \end{aligned}$$

Recall that  $d \geq C \log(n/\delta)$ . We can set  $\lambda = 1/2$  so that,  $\|\mathbf{z}_2\|_2^2 > d\sigma^2/2$  with probability at least  $1 - \delta/6n$  (for large enough  $C$ ). Taking a union bound over all examples gives us the desired result.

Note that we can get closer to  $d\sigma^2$  by choosing an appropriately higher value of  $C$ .

□

**Lemma 5** (Dataset Separability). *With probability at least  $1 - \delta$  over random samples of dataset  $\mathcal{S}_A$ , samples  $(\mathbf{x}_1, \mathbf{y}_1), \dots, (\mathbf{x}_n, \mathbf{y}_n)$  are linearly separable.*

*Proof.* We will show that there exists a set of weights that with high probability correctly classify each example in the dataset. Let  $\mathbf{x}_i = \boldsymbol{\mu}_i + \mathbf{z}_i$  for every data point in  $\mathcal{S}$ . From assumptions, our dataset  $(\mathcal{S}_A)$  contains  $\mathbf{z}_r, \mathbf{z}_m$  that belong to rare and mislabeled groups. Rest all points are denoted by  $\mathbf{z}_k$ . Consider the classifier  $\mathbf{w} = \sum_i^n \mathbf{y}_i \mathbf{x}_i$ . Then,

**Case 1: Rare Example**  $(\mathbf{x}_r, \mathbf{y}_r)$

$$\begin{aligned} \mathbf{y}_r \mathbf{w} \cdot \mathbf{x}_r &= \sum_j \mathbf{y}_j \mathbf{y}_r \mathbf{x}_j \cdot \mathbf{x}_r, \\ &= \mathbf{x}_r \cdot \mathbf{x}_r + \sum_{i \neq r} \mathbf{y}_r \mathbf{y}_i \mathbf{x}_i \cdot \mathbf{x}_r, \\ &= \boldsymbol{\mu}_r \cdot \boldsymbol{\mu}_r + \mathbf{z}_i \cdot \mathbf{z}_i + \sum_{r \neq j} \mathbf{y}_r \mathbf{y}_j \mathbf{z}_r \cdot \mathbf{z}_j, \quad (\text{since } \boldsymbol{\mu}_r \cdot \boldsymbol{\mu}_i = 0 \ \forall i \in \{1, 2\}) \\ &= k\mu^2 + \mathbf{z}_i \cdot \mathbf{z}_i + \sum_{r \neq j} \mathbf{y}_r \mathbf{y}_j \mathbf{z}_r \cdot \mathbf{z}_j, \\ &\geq 0 + d\sigma^2/2 - c_1 n \sqrt{d\sigma^2 \log(n/\delta)} \quad (\text{using Lemma 4 and Corollary 1}) \\ &> 0 \end{aligned}$$

for  $d \geq Cn^2 \log(n/\delta)$ .



**Case 2: Misabeled Example**  $(\mathbf{x}_m, \mathbf{y}_m)$  Without loss of generality, assume that the mislabeled example is sampled from the distribution  $\mathcal{D}_1$ , and the set of all correctly labeled examples in the dataset  $\mathcal{S}_A$  from this distribution be  $\mathcal{S}_{A,1}$ .

$$\begin{aligned}
\mathbf{y}_m \mathbf{w} \cdot \mathbf{x}_m &= \sum_j \mathbf{y}_j \mathbf{y}_m \mathbf{x}_j \cdot \mathbf{x}_m, \\
&= \mathbf{x}_m \cdot \mathbf{x}_m + \sum_{i \neq j} \mathbf{y}_m \mathbf{y}_j \mathbf{x}_m \cdot \mathbf{x}_j, \\
&= \boldsymbol{\mu}_m \cdot \boldsymbol{\mu}_m + \mathbf{z}_i \cdot \mathbf{z}_i - \sum_{j \in \mathcal{S}_{A,1}} \boldsymbol{\mu}_m \cdot \boldsymbol{\mu}_m + \sum_{m \neq j} \mathbf{y}_m \mathbf{y}_j \mathbf{z}_m \cdot \mathbf{z}_j, \\
&\quad \text{(since } \boldsymbol{\mu}_m \cdot \boldsymbol{\mu}_1 = 1, \boldsymbol{\mu}_m \cdot \boldsymbol{\mu}_2 = 0, \boldsymbol{\mu}_m \cdot \boldsymbol{\mu}_r = 0) \\
&\geq 0 + d\sigma^2/2 - n_1 k \mu^2 - c_1 n \sqrt{d \log(n/\delta)}, \\
&\quad \text{(if } n_1 = |\mathcal{S}_{A,1}|, \text{ using Lemma 4 and Corollary 1)} \\
&\geq d\sigma^2/2 - nk\mu^2 - c_1 n \sigma^2 \sqrt{d \log(n/\delta)}, \\
&> 0
\end{aligned}$$

for  $d \geq C \max\{n^2 \log(n/\delta), nk\mu^2/\sigma^2\}$ .

**Case 3: Majority Example**  $(\mathbf{x}_i, \mathbf{y}_i)$  The case of majority examples directly follows from the case of mislabeled examples. Rather than having a negative summation over the mean vector for  $n$  examples (in line 3), we will have a positive summation because the true label will match the label of the rest of the examples in the subset  $\mathcal{S}_{A,1}$ . This will make the expected value of the dot product even larger.

*Remark:* Since the first split dataset  $\mathcal{S}_A$  is separable, it directly follows that the second split dataset  $\mathcal{S}_B$  is also separable since it does not have any mislabeled and rare examples.  $\square$

## A.7 Lemmas for Theorem 4

**Lemma 6** (Sign of  $\beta$ ).  $\beta_j \geq 0$  for all  $j$  in Equation 25.

*Proof.* Analyzing the steps of gradient descent, we have:

$$\begin{aligned}
\dot{\mathbf{w}}(t) &= -\nabla \mathcal{L}(\mathbf{w}(t)) \\
&= \sum_{j \in \mathcal{S}_B} \exp(-\mathbf{y}_j \mathbf{x}_j^\top \mathbf{w}(t)) (\mathbf{y}_j \mathbf{x}_j^\top) \\
\mathbf{w}(t) - \mathbf{w}(0) &= \sum_{j \in \mathcal{S}_B} \left( \mathbf{y}_j \mathbf{x}_j^\top \int_0^t \underbrace{\exp(-\mathbf{y}_j \mathbf{x}_j^\top \mathbf{w}(t))}_{\beta_j(t)} dt \right), \tag{41}
\end{aligned}$$

Hence,  $\beta_j \geq 0 \forall j$ .  $\square$

**Lemma 7** (Support Vectors). If dataset  $\mathcal{S}_A$  is separable, then  $(\mathbf{x}_m, \mathbf{y}_m)$  and  $(\mathbf{x}_r, \mathbf{y}_r)$  are support vectors for  $\mathcal{S}_A$ .

*Proof.* We will prove by contradiction. From our assumption,  $(\mathbf{x}_m, \mathbf{y}_m)$  and  $(\mathbf{x}_r, \mathbf{y}_r)$  are the only mislabeled and rare examples in the first-split  $\mathcal{S}_A$  from their respective sub-groups. Let us assume that they are not support vectors. Then, we can directly follow from the Asymptotic Analysis in Subsection A.4 that the probability of correct classification of the rare example is 0.5 and for the mislabeled example approaches 0 as the model is trained for infinite time. But we know that the model achieves 100% accuracy on the training set  $\mathcal{S}_A$  with weights  $\mathbf{w}_A(T)$ . Hence, this is a contradiction, and  $(\mathbf{x}_m, \mathbf{y}_m)$ ,  $(\mathbf{x}_r, \mathbf{y}_r)$  must be support vectors for the original classification problem.  $\square$

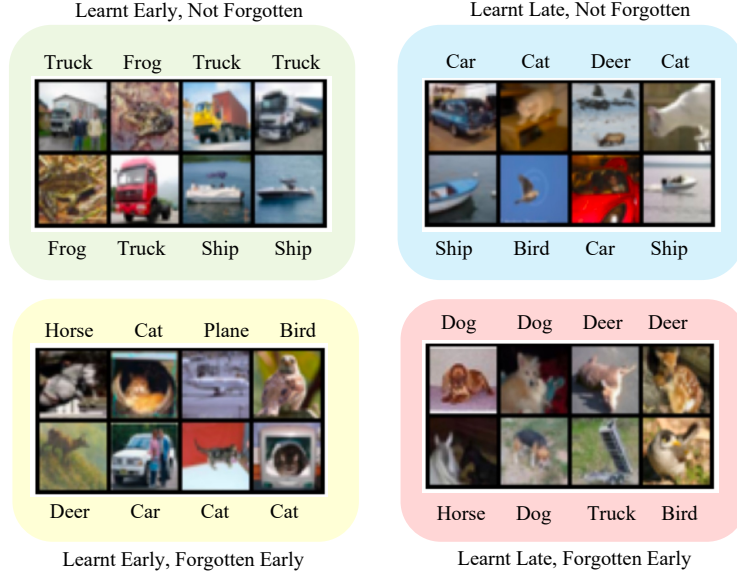


Figure 6: Examples from the CIFAR-10 dataset grouped based on their learning and forgetting time.

**Lemma 8** (Bounded Time Prediction). *With probability at least  $1 - \delta$ , there exists a bounded time  $T' = O(T)$  when the mislabeled examples are incorrectly classified with high probability.*

*Proof.* In Subsection A.4 we have shown that at infinite time training, the model misclassifies mislabeled examples with high probability. Then, the existence of bounded time weights for which the prediction of mislabeled examples is flipped directly follows from the proof of Proposition 1 applied to the result of Theorem 3.  $\square$

## B Experimental Results

### B.1 Experimental Setup

**Architectures** We perform experiments using four different model architectures—LeNet, ResNet-9 [4], ResNet-50, and Bert-base-cased [14]. Comparisons with model architectures are used in analysis of stability of the SSFT metric. For other numbers reported in tables and plots, we use the ResNet-9 model, unless otherwise stated.

**Optimizer** We experiment with three different learning rate scheduling strategies—cyclic learning rate schedule, cosine learning rate, and step decay learning rate. We test for two values of peak learning rate—0.1 and 0.01. All the model are trained using the SGD optimizer with weight decay  $5e-4$  and momentum 0.9, apart from the comparison with optimizers in Appendix C.1 where we also experiment with the Adam optimizer.

**Training Procedure** We train for a maximum of 100 epochs or until we have 5 epochs of 100% training accuracy. We first train on  $\mathcal{S}_A$ , and then using the pre-initialized weights from stage 1, train on  $\mathcal{S}_B$  with the same learning parameters. All experiments can be performed on a single RTX2080 Ti.

### B.2 Image Datasets

In the main paper we present visualizations of training examples from the MNIST dataset based on which quadrant they lie on in the learning-forgetting graph. Here, we complement our findings by showing visualizations for the CIFAR-10 dataset. We note that CIFAR-10 dataset provide many different types of visibility patterns within the same class. Hence, examples may be learnt late due to belonging to a rare visibility pattern. In Figure 6, we see that the examples that were learnt earliest and

never forgotten have similar visibility patterns—for instance all the trucks have a similar perspective. On the contrary, as we move to the first quadrant with examples that were learnt late but never forgotten, we see that all the examples are true to their semantic class, but these visibility patterns occur rarely. Finally, we also analyze the visualizations based on examples that were forgotten during the course of second-split training. While in the case of MNIST dataset, SSFT was able to remove the mislabeled examples well, we see that CIFAR-10 offers more challenges because examples may be ambiguous because of other reasons and may be forgotten owing to the model using spurious features.

## C Ablation Studies

We detail the experimental setup used to conduct our ablation studies directed towards understanding the learning and forgetting dynamics of rare and complex examples respectively.

**Rare Examples** The experiments to show the rate of learning for rare examples are inspired by the singleton hypothesis as proposed by Feldman [15]. The hypothesis was inspired by a long-tailed distribution of visibility patterns in the person and bus category of the PASCAL dataset. For example, the dataset contains many buses with the front visible, but very few buses that were captured from the rear or the side, and even fewer buses whose view is occluded by the presence of other objects in front of them. (Refer to Figure 1 in their work for more details.) In our work, we first attempted to leverage the same training set-up with the provided visibility patterns. However, we noted that there wasn’t a strong correlation between the frequency of an example’s visibility pattern, and the rate at which it was learnt. We hypothesize that this is because there are other factors of example hardness that may make an example be learnt slowly or fast (such as complexity, as detailed in the next paragraph). This can lead to an example being learnt fast if it has a simple pattern yet occurs rarely. Especially when there are only  $O(1)$  samples from a given sub-group (based on the visibility pattern), we can not make any claims based on singleton correlation alone.

Hence, in order to distill the frequency of occurrence of an example with other confounders that may influence its training-time, we created a long-tailed dataset from the CIFAR-100 dataset. CIFAR-100 is a dataset of 100 object classes, which can be further grouped into 20 super-classes. For instance, examples from categories *maple*, *oak*, *palm*, *pine*, *willow* all belong to the ‘superclass’ of *trees*. Similar division of 5 sub-classes is provided in the datasets for each of the superclasses. Each class contains 500 training examples and the overall dataset has 50,000 training examples.

As a first step towards creating a long-tailed dataset, we assign a fixed frequency ordering within the subgroups of a superclass. The most frequent subgroup has 500 examples in the training set, for the next most frequent subgroup, we randomly select 250 examples from the training set, and so on until the last sub group with 31 examples in the training set. This means that there are exactly 20 sub-groups in the final dataset with  $\{500, 250, 125, 62, 31\}$  examples respectively. Irrespective of the class number, the task is to predict the corresponding superclass, that is, we reduce the problem to a 20-class classification problem. However, we track the learning and forgetting dynamics of examples from each of the 100 sub-groups separately, based on their group frequency. To remove any other confounders of example hardness, we (i) randomize the group frequency ordering of the sub-groups within a superclass (in case some classes are harder to learn than the others); and (ii) randomize the examples that were selected based on the group size (in case some examples were ambiguous or hard). We further split the dataset into two IID partitions to analyze the learning time and SSFT, and average the results over 20 random runs of the experiment. Experimental results are detailed in the main paper.

**Complex Examples** Prior works advocating for, and understanding the simplicity bias [43] have operationalized the notion of simplicity via the complexity of hypothesis class required to learn the distribution that a complex example may be sampled from. In particular, Shah et al. [43] construct a synthetic dataset with MNIST and CIFAR-10 images vertically stacked on top of each other—with the part with MNIST images corresponding to the part of the combined image with *simpler* features, and the part with CIFAR-10 images corresponding to the part with *complex* features. They show that the model almost completely relies on the part of the image containing the MNIST digit even when it is less predictive of the true label. Inspired by this argument about the simplicity of features, we create a dataset that has the union of images from the MNIST and the CIFAR-10 dataset. More specifically, we select classes from the MNIST dataset corresponding to digits  $\{0, 1, 2, 3\}$ , and classes from the CIFAR-10 dataset corresponding to  $\{\text{horses, airplanes, dog, frog}\}$  and label them

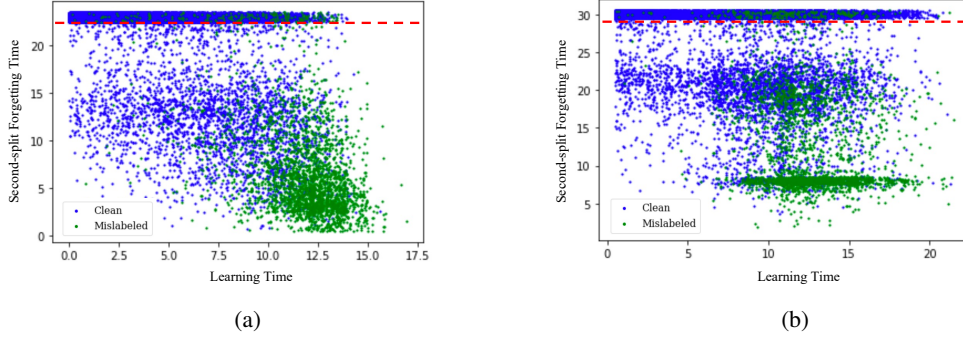


Figure 7: FSLT (First-split learning time) is able to provide some degree of separation between mislabeled and clean samples when trained with the SGD optimizer (left), but fails when the model is trained using Adam (right) on the CIFAR-10 dataset.

from  $\{0, 1, 2, 3\}$ . This means that the model associates the label 0 to both the digit 0 and airplane class. The attempt of this experiment is to draw the link between the simplicity bias and the rate of learning. Experimental results are provided in the main paper.

### C.1 Stability of our metric

**Stability across architectures** The forgetting of examples is a property of both the dataset and the model architecture. As a result, we find that just like the learning time, the forgetting time has a lower correlation between architectures. The average pearson correlation between the ResNet-9 and ResNet-50 models is 0.62 in case of the CIFAR10 dataset. However, we note that the most forgotten examples generalize across datasets. That is, the average pearson correlation between the bottom 10% examples of the dataset is 0.87. This highlights how the forgetting metric is good for finding misaligned examples in the dataset, since they are not a property of the model architecture. We suspect that among the examples that are infrequently forgotten, the model capacity and other inductive biases of the model architecture may have a role in driving the average pearson correlation low.

**Stability across optimizers** Jiang et al. [25] showed that changing the learning optimizer from SGD to Adam can lead to a significant change in the learning rate of examples from different levels of hardness (based on their regularity metric). More specifically, they find that examples with a low consistency score (closely correlated with learning speed) also get learnt fast when using the Adam optimizer. This suggests that using an optimizer like Adam at training time may have an impact on the ability of learning time based metrics to separate examples. In Figure 7, we contrast the ability of forgetting and learning time based metrics for identifying label noise when using the SGD and Adam optimizers. When using an optimizer such as SGD, the mislabeled samples are learnt slower than a large fraction of the training examples, and the learning time metric offers some degree of separation between the clean and mislabeled examples. However, when we use the Adam optimizer, it results in joint learning of a large fraction of both mislabeled and clean samples. Hence, offering a very low degree of separation. However, under the same training procedure, the SSFT still allows us to distinguish between the mislabeled and clean samples.

**Stability across seeds and learning rates** The pearson correlation for stability across seeds for the forgetting time metric is 0.83. This is higher than the corresponding learning time based metric (correlation 0.56). However, one of the drawbacks of our proposed metric is that the SSFT requires the use of an appropriate learning rate that allows the examples to be forgotten slowly. We provide more information about the same in the main paper.

**Stability Across Learning Rate Schedules** We experiment with three different learning rate schedules for the second-split training—triangular, cosine and linear. In triangular learning rate, we increase the learning rate from 0 to the maximum set value linearly over the first 10 epochs, and then decay it back to 0 until we reach the last epoch (maximum of 100 epochs). In case of the linear schedule, we increase the learning rate from 0 to the maximum set value linearly over the course of 100 epochs. The intuition behind using a linear learning rate schedule was to be able to set higher

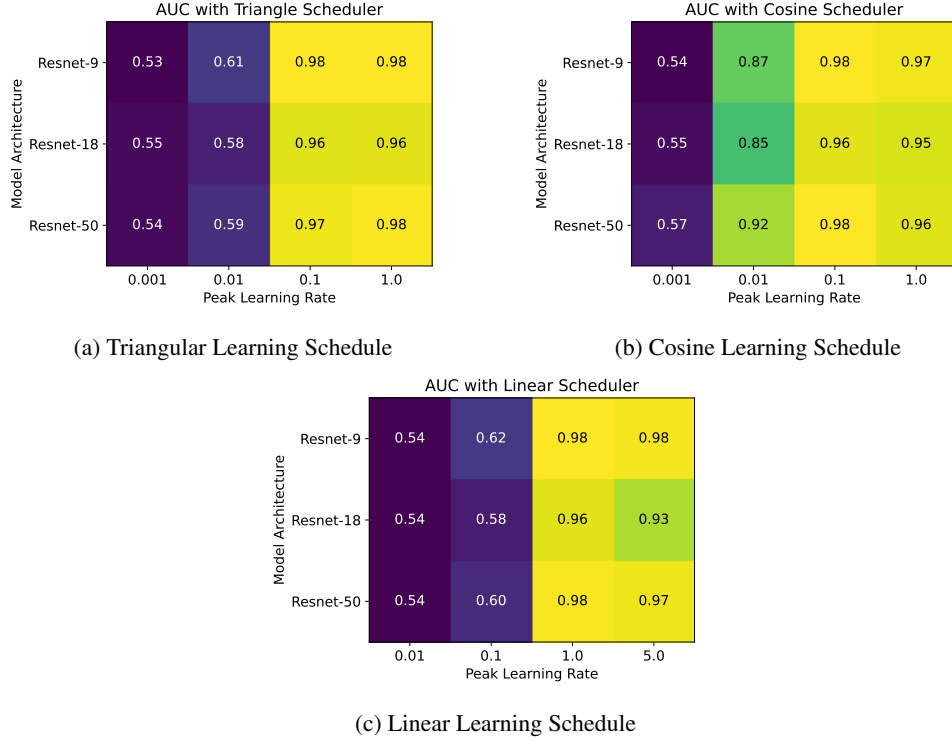


Figure 8: We present a heatmap for the AUC of detection of mislabeled examples using the SSFT metric under various learning rates, architecture sizes and learning rate schedules for the SGD optimizer. The experiment was performed on the CIFAR-10 dataset with 10% label noise, and the forgetting times were averaged over 5 seeds before using them for AUC calculation.

peak learning rates so that the model eventually forget all the examples in the first split and we can create a better ordering between samples based on forgetting time (as opposed to the setting where only a small fraction of examples are ever forgotten).

The results of the combined analysis across hyperparameters such as architecture, learning rate and learning rate schedule are presented in the heatmap of AUC of mislabeled example detection in Figure 8. The experiment was performed on the CIFAR-10 dataset with 10% label noise, and the forgetting times were averaged over 5 seeds before using them for AUC calculation. We can see that uniformly across architectures and learning rate schedules having a very low learning rate makes nearly all examples indistinguishable based on forgetting time. This is because all the models are sufficiently overparametrized to memorize all the examples in the dataset. Hence, when using a very small learning rate the optimization step moves the model weights insignificantly and we do not forget many mislabeled samples from the first split. On the other side, having a large learning rate helps achieve a strong separation between mislabeled and clean examples which also shows up in the form of high AUC values in the figure.

## C.2 Impact of Sampling Frequency of Mislabeled Examples

In the synthetic experiment performed in Section 4.2, we assumed that mislabeled examples occur from the majority subgroups. As a result, we observed that they get forgotten quickly during second-split training. However, in this section we aim to understand the impact of sampling frequency on the forgetting time of mislabeled examples. More specifically, we now assume that mislabeled samples occur in rare subgroups in the synthetic setup. We find that the learning curves of the mislabeled example stays the same as before, but the forgetting time for mislabeled examples closely approaches that of rare examples. This is because there is very little signal for the model to learn the opposite class during second split training since the example occurs only  $O(1)$  times. The learning and forgetting curves pertaining to the same experiment are presented in Figure 9. In contrast with the forgetting



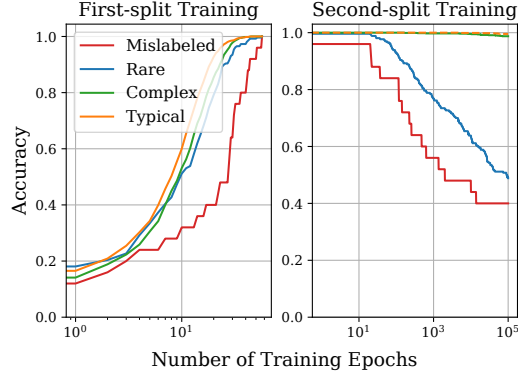


Figure 9: Learning and forgetting curves for mislabeled and rare examples when the mislabeled examples are drawn from rare subgroups in the synthetic setup described in Section 4.2.

curve in Section 4.2 where the mislabeled examples are quickly forgotten and their prediction is flipped, we find that when the subgroup corresponding to the mislabeled examples is infrequent, their forgetting time closely corresponds to that of rare examples; and on aggregate their predictions do not get flipped in the epochs that the model was trained for.

### C.3 Mislabeled Example Detection

In this section we provide additional details about the experimental setting for the results presented in Table 2. In case of the CIFAR-100 dataset, we reduce the learning rate for second-split training by a factor of 10, and use batch size of 128. While all the other training procedures in this paper used a cyclic learning rate, for the case of CIFAR-100, we use warm-up based multi-step decay learning rate schedule.<sup>2</sup> The model used for training was ResNet-18, and 10% label noise was added. The training setting for EMNIST is identical to that of the MNIST dataset. We use the first 10 classes of the dataset to make it a 10-class classification problem.

<sup>2</sup>We follow the code in <https://github.com/weiaicunzai/pytorch-cifar100>