

402 A Analyzing and Visualizing the Results of the Reconstruction Optimization

403 The analysis of the results of the various reconstruction losses Eqs. (6), (12) and (15), involve
404 verifying and checking which of the training samples were reconstructed. In this section we provide
405 further details on our method for analyzing the reconstruction results, and how we measure the quality
406 of our reconstructions.

407 A.1 Analyzing the Results of the Reconstruction Optimization

408 In order to match between samples from the training set and the outputs of the reconstruction
409 algorithm (the so-called "candidates") we follow the same protocol of Haim et al. [2022]. Note that
410 before training our models, we subtract the mean image from the given training set. Therefore the
411 training samples are d -dimensional objects where each entry is in $[-1, 1]$.

412 First, for each training sample we compute the distance to all the candidates using a normalized L_2
413 score:

$$d(\mathbf{x}, \mathbf{y}) = \left\| \frac{\mathbf{x} - \mu_{\mathbf{x}}}{\sigma_{\mathbf{x}}} - \frac{\mathbf{y} - \mu_{\mathbf{y}}}{\sigma_{\mathbf{y}}} \right\|_2^2 \quad (16)$$

414 Where $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$ are a training sample or an output candidate from the reconstruction algorithm,
415 $\mu_{\mathbf{x}} = \frac{1}{d} \sum_{i=1}^d \mathbf{x}(i)$ is the mean of \mathbf{x} and $\sigma_{\mathbf{x}} = \sqrt{\frac{1}{d-1} \sum_{i=1}^d (\mathbf{x}(i) - \mu_{\mathbf{x}})^2}$ is the standard deviation
416 of \mathbf{x} (and the same goes for $\mathbf{y}, \mu_{\mathbf{y}}, \sigma_{\mathbf{y}}$).

417 Second, for each training sample, we take C candidates with the smallest distance according to
418 Eq. (16). C is determined by finding the first candidate whose distance is larger than B times the
419 distance to the closest nearest neighbour (where B is a hyperparameter). Namely, for a training sample
420 \mathbf{x} , the nearest neighbour is \mathbf{y}_1 with a distance $d(\mathbf{x}, \mathbf{y}_1)$, then C is determined by finding a candidate
421 \mathbf{y}_{C+1} whose distance is $d(\mathbf{x}, \mathbf{y}_{C+1}) > B \cdot d(\mathbf{x}, \mathbf{y}_1)$, and for all $j \leq C$, $d(\mathbf{x}, \mathbf{y}_j) \leq B \cdot d(\mathbf{x}, \mathbf{y}_1)$. B
422 was chosen heuristically to be $B = 1.1$ for MLPs, and $B = 1.5$ for convolutional models. The C
423 candidates are then summed to create the reconstructed sample $\hat{\mathbf{x}} = \frac{1}{C} \sum_{j=1}^C \mathbf{y}_j$. In general, we can
424 also take only $C = 1$ candidate, namely just one nearest neighbour per training sample, but choosing
425 more candidates improve the visual quality of the reconstructed samples.

426 Third, the reconstructed sample $\hat{\mathbf{x}}$ is scaled to an image in $[0, 1]$ by adding the training set mean and
427 linearly "stretching" the minimal and maximal values of the result to $[0, 1]$. Finally, we compute
428 the SSIM between the training sample \mathbf{x} and the reconstructed sample $\hat{\mathbf{x}}$ to measure the quality of
429 reconstruction.

430 A.2 Deciding whether a Reconstruction is "Good"

431 Here we justify our selection for SSIM=0.4 as the threshold for what we consider as a "good"
432 reconstruction. In general, the problem of deciding whether a reconstruction is the correct match to a
433 given sample, or whether a reconstruction is a "good" reconstruction is equivalent to the problem of
434 comparing between images. No "synthetic" metric (like SSIM, l_2 etc.) will be aligned with human
435 perception. A common metric for this purpose is LPIPS Zhang et al. [2018] that uses a classifier
436 trained on Imagenet Deng et al. [2009], but since CIFAR images are much smaller than Imagenet
437 images (32×32 vs. 224×224) it is not clear that this metric will be better than SSIM.

438 As a simple rule of thumb, we use SSIM>0.4 for deciding that a given reconstruction is "good".
439 To justify, we plot the best reconstructions (in terms of SSIM) in Fig. 8. Note that almost all
440 samples with SSIM>0.4 are also visually similar (for a human). Also note that some of the samples
441 with SSIM<0.4 are visually similar, so in this sense we are "missing" some good reconstructions.
442 In general, determining whether a candidate output of a reconstruction algorithm is a match to a
443 training sample is an open question and a problem in all other works for data reconstruction, see for
444 example Carlini et al. [2023] that derived a heuristic for reconstructed samples from a generative
445 model. This cannot be dealt in the scope of this paper, and is an interesting future direction for our
446 work.



Figure 8: Justifying the threshold of SSIM= 0.4 as good rule-of-thumb for a threshold for a “good” reconstruction. The SSIM values are shown above each train-reconstruction pair. Note that samples with SSIM> 0.4 (blue) are visually similar. Also some of the samples with SSIM< 0.4 (red) are similar. In general deciding whether a reconstruction is “good” is an open question beyond the scope of this paper.

447 B Implementation Details

448 **Further Training Details.** The models that were reconstructed in the main part of the paper were
 449 trained with learning rates of 0.01 for binary classifiers (both MLP and convolutional), and 0.5 in the
 450 case of multi-class classifier (Section 4). The models were trained with full batch gradient descent
 451 for 10^6 epochs, to guarantee convergence to a KKT point of Eq. (1) or a local minima of Eq. (13).
 452 We note that Haim et al. [2022] observed that models trained with SGD can also be reconstructed.
 453 The experiment in Appendix G (large models with many samples) also uses SGD and results with
 454 similar conclusion, that some models trained with SGD can be reconstructed. In general, exploring
 455 reconstruction from models trained with SGD is an interesting direction for future works.

456 **Runtime and Hardware.** Runtime of a single reconstruction run (specific choice of hyperparamet-
 457 ers) from a model D-1000-1000-1 takes about 20 minutes on a GPU Tesla V-100 32GB or NVIDIA
 458 Ampere Tesla A40 48GB.

459 **Hyperparameters of the Reconstruction Algorithm.** Note that the reconstruction loss contains the
 460 derivative of a model with ReLU layers, which is flat and not-continuous. Thus, taking the derivative
 461 of the reconstruction loss results in a zero function. To address this issue we follow a solution
 462 presented in Haim et al. [2022]. Namely, given a trained model, we replace in the backward phase
 463 of backpropagation the ReLU function with the derivative of a softplus function (or SmoothReLU)
 464 $f(x) = \alpha \log(1 + e^{-x})$, where α is a hyperparameter of the reconstruction scheme. The functionality
 465 of the model itself does not change, as in the forward phase the function remains a ReLU. Only
 466 the backward function is replaced with a smoother version of the derivative of ReLU which is
 467 $f'(x) = \alpha \sigma(x) = \frac{\alpha}{1+e^{-x}}$ (here σ is the Sigmoid function). To find good reconstructions we run the
 468 algorithm multiple times (typically 100 times) with random search over the hyperparameters (using

469 the Weights & Biases framework Biewald [2020]). The exact parameters for the hyperparameters
 470 search are:

- 471 • Learning rate: log-uniform in $[10^{-5}, 1]$
- 472 • σ_x : log-uniform in $[10^{-6}, 0.1]$
- 473 • λ_{\min} : uniform in $[0.01, 0.5]$
- 474 • α : uniform in $[10, 500]$

475 B.1 Small Initialisation

476 Models whose first layer was initialized with a small (non-standard) initialization plays several roles
 477 in our paper. These include models that were trained following the approach in Haim et al. [2022],
 478 as in Section 4, or comparison of such models to models trained with weight-decay, as discussed
 479 in Section 6.1. The models whose results appear in Fig. 2 and Fig. 5 were initialized with a scale of
 480 10^{-3} . After submitting the paper we noticed that the initialization used in Haim et al. [2022] was
 481 in fact smaller - 10^{-4} . In order to make a fair comparison, we re-run the baselines shown in Fig. 5
 482 (red-dashed lines). As seen in Fig. 9, the corrected initialisation increase the number of reconstructed
 483 samples in the case of the binary classifier (Fig. 9a) and decrease in the case of multiclass classifier
 484 (Fig. 9b). In both cases, this does not change the main claim in Section 6.1, that using weight decay
 485 terms during training changes the reconstructability and in some cases dramatically increase the
 486 number of samples that are vulnerable to reconstruction. We will make sure to update Fig. 5 with the
 487 corrected version.

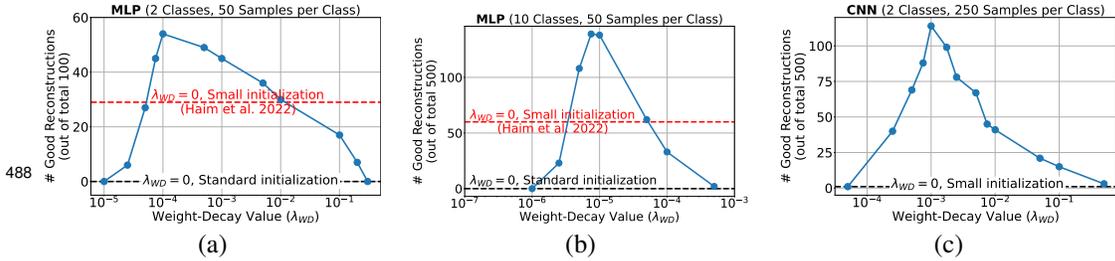


Figure 9: Corrected version of Fig. 5. The correction of the baselines in red did not affect the claims in Section 6.1.

489 C Experiments with Different Number of Classes and Fixed Training Set Size



Figure 10: Experiments of reconstruction from models trained on a fixed training set size (500 samples) for different number of classes. Number of "good" reconstruction is shown for each model.

490 To complete the experiment shown in Fig. 3, we also perform experiments on models trained on
 491 various number of classes ($C \in \{2, 3, 4, 5, 10\}$) and with a fixed training set size of 500
 492 samples (distributed equally between classes), see Fig. 10. It can be seen that as the number of classes
 493 increases, also does the number of good reconstructions, where for 10 classes there are more than 6
 494 times good reconstructions than for 2 classes. Also, the quality of the reconstructions improves as the
 495 number of classes increase, which is depicted by an overall higher SSIM score. We also note, that
 496 the number of good reconstructions in Fig. 10 is very similar to the number of good reconstructions
 497 from Fig. 3 for 50 samples per class. We hypothesize that although the number of training samples
 498 increases, the number of "support vectors" (i.e samples on the margin which can be reconstructed)
 499 that are required for successfully interpolating the entire dataset does not change by much.

500 **D General Losses - More Results**

501 Following the discussion in Section 5 and Fig. 4, Figures 11, 12, 13 present visualizations of training
502 samples and their reconstructions from models trained with L_2 , $L_{2.5}$ and Huber loss, respectively.



Figure 11: **Reconstruction using L_2 loss.** Training samples (red) and their best reconstructions (blue) using an MLP classifier that was trained on 300 CIFAR10 images using an L_2 regression loss, as described in Section 5 and Fig. 4.



Figure 12: **Reconstruction using $L_{2.5}$ loss.** Training samples (red) and their best reconstructions (blue) using an MLP classifier that was trained on 300 CIFAR10 images using an $L_{2.5}$ regression loss, as described in Section 5 and Fig. 4.



Figure 13: **Reconstruction using Huber loss.** Training samples (red) and their best reconstructions (blue) using an MLP classifier that was trained on 300 CIFAR10 images using Huber loss, as described in Section 5 and Fig. 4.

503 E Further Analysis of Weight Decay

504 By looking at the exact distribution of reconstruction quality to the distance from the margin, we
 505 observe that weight-decay (for some values) results in more training samples being on the margin of
 506 the trained classifier, thus being more vulnerable to our reconstruction scheme.

507 This observation is shown in Fig. 14 where we show the scatter plots for all the experiments from Fig. 5
 508 (a). We also provide the train and test errors for each model. It seems that the test error does not
 509 change significantly. However, an interesting observation is that reconstruction is possible even for
 510 models with non-zero training errors, i.e. models that do not interpolate the data, for which the
 511 assumptions of Lyu and Li [2019] do not hold.

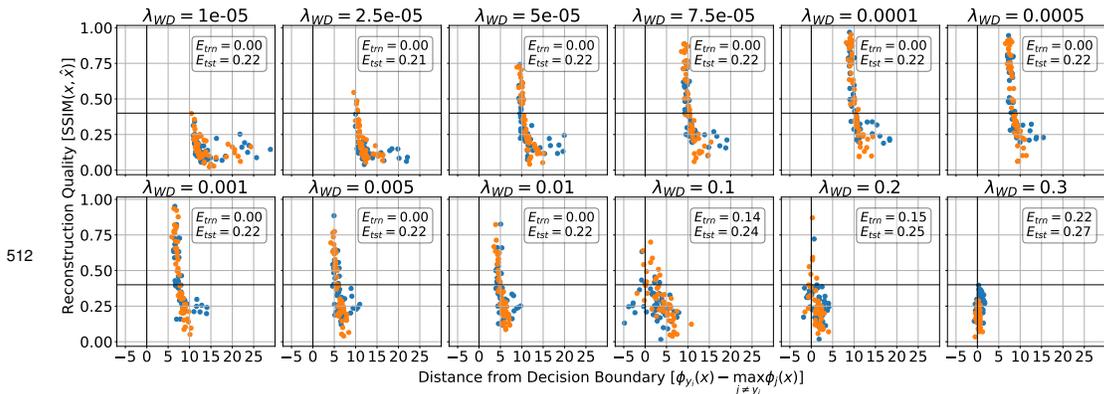


Figure 14: Scatter plots of the 12 experiments from Fig. 5 (a). Each plot is model trained with a different value of weight decay on 2 classes with 50 samples in each class. Certain values of weight decay make the model more susceptible to our reconstruction scheme.

513 F Convolutional Neural Networks - Ablations and Observations

514 In this section we provide more results and visualizations to the experiments on convolutional neural
 515 network in Section 6.1.

516 In Fig. 15 we show ablations for the choice of the kernel-size (k) and number of output channels
 517 (C_{out}) for models with architecture CONV(kernel-size= k , output-channels= C_{out})-1000-1. All models
 518 were trained on 500 images (250 images per class) from the CIFAR10 dataset, with weight-decay

519 term $\lambda_{WD}=0.001$. As can be seen, for such convolutional models we are able to reconstruct samples
 520 for a wide range of choices.

521 Note that the full summary of reconstruction quality versus the distance from the decision boundary
 522 for the model whose reconstructed samples are shown in Fig. 6, is shown in Fig. 15 for kernel-size 3
 523 (first row) and number of output channels 32 (third column).

524 **Further analysis of Fig. 15.** As expected for models with less parameters, the reconstructability
 525 decreases as the number of output channels decrease. An interesting phenomenon is observed for
 526 varying the kernel size: for a fixed number of output channel, as the kernel size increases,
 527 the susceptibility of the model to our reconstruction scheme decreases. However, as the kernel size
 528 approaches 32 (the full resolution of the input image), the reconstructability increases once again.
 529 On the one hand it is expected, since for kernel-size=32 the model is essentially an MLP, albeit
 530 with smaller hidden dimension than usual (at most 64 here, whereas the typical model used in the
 531 paper had 1000). On the other hand, it is not clear why for some intermediate values of kernel size
 532 (in between 3 and 32) the reconstructability decreases dramatically (for many models there are no
 533 reconstructed samples at all). This observation is an interesting research direction for future works.

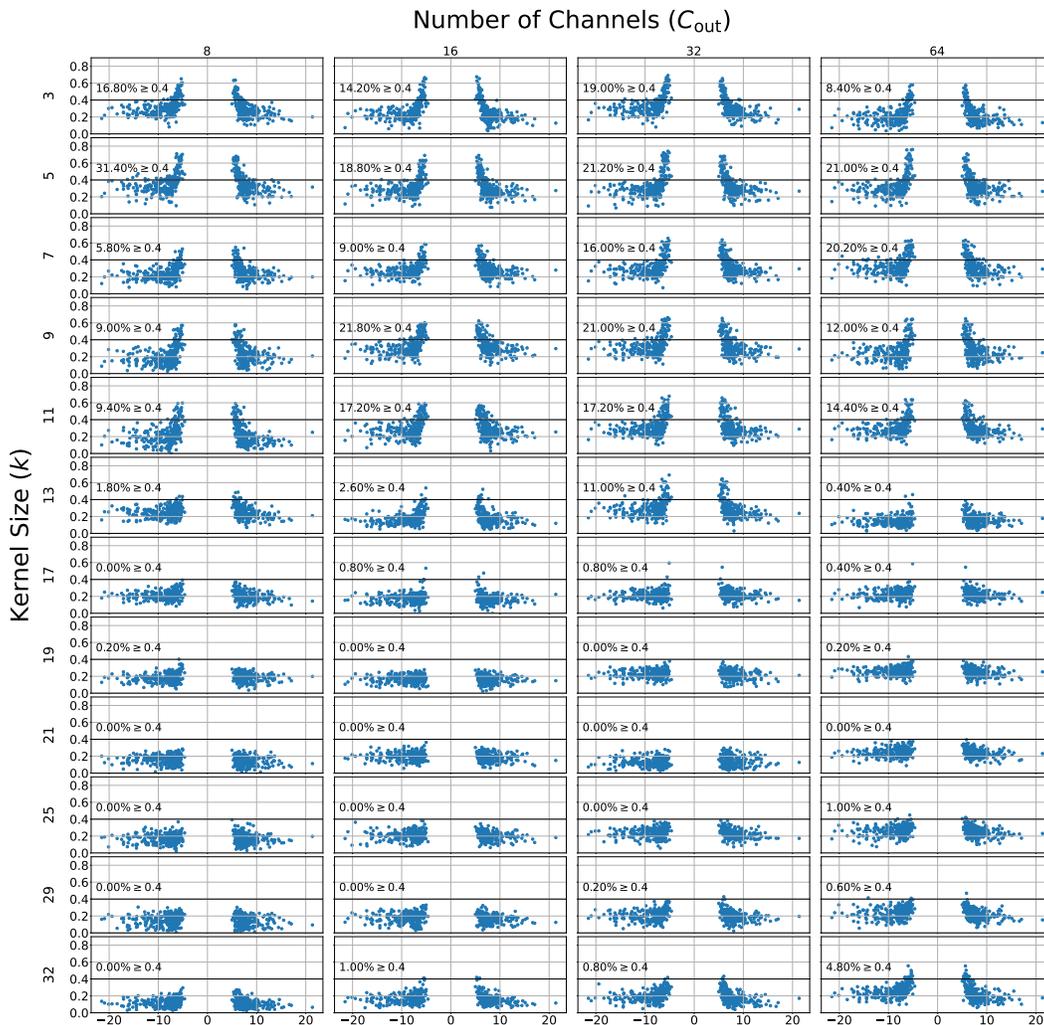


Figure 15: Ablating the choice of the kernel size and output-channels for reconstruction from neural binary classifiers with architecture CONV(kernel-size= k ,output-channels= C_{out})-1000-1.

534 **Visualizing Kernels.** In Haim et al. [2022], it was shown that some of the training samples can
 535 be found in the first layer of the trained MLPs, by reshaping and visualizing the weights of the first
 536 fully-connected layer. As opposed to MLPs, in the case of a model whose first layer is a convolution
 537 layer, this is not possible. For completeness, in Fig. 16 we visualize all 32 kernels of the Conv layer.
 538 Obviously, full images of shape $3 \times 32 \times 32$ cannot be found in kernels of shape $3 \times 3 \times 3$, which makes
 539 reconstruction from such models (with convolution first layer) even more interesting.

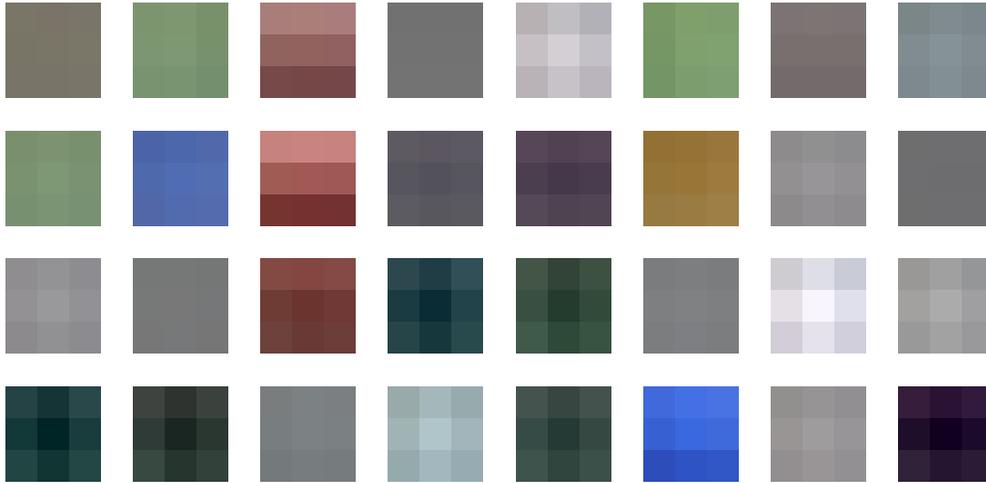


Figure 16: The kernels of the model whose reconstructions are shown in Fig. 6, displayed as RGB images.

540 G Reconstruction From a Larger Number of Samples

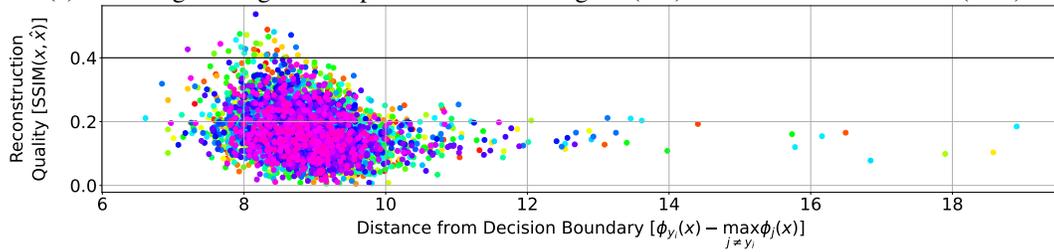
541 One of the major limitations of Haim et al. [2022] is that they reconstruct from models that trained on
 542 a relatively small number of samples. Specifically, in their largest experiment, a model is trained with
 543 only 1,000 samples. Here we take a step further, and apply our reconstruction scheme for a model
 544 trained on 5,000 data samples.

545 To this end, we trained a 3-layer MLP, where the number of neurons in each hidden layer is 10,000.
 546 Note that the size of the hidden layer is 10 times larger than in any other model we used. Increasing
 547 the number of neurons seems to be one of the major reasons for which we are able to reconstruct
 548 from such large datasets, although we believe it could be done with smaller models, which we leave
 549 for future research. We used the CIFAR100 dataset, with 50 samples in each class, for a total of 5000
 550 samples.

551 In Fig. 17a we give the best reconstructions of the model. Note that although there is a degradation in
 552 the quality of the reconstruction w.r.t a model trained on less samples, it is still clear that our scheme
 553 can reconstruct some of the training samples to some extent. In Fig. 17b we show a scatter plot of the
 554 SSIM score w.r.t the distance from the boundary, similar to Fig. 3a. Although most of the samples
 555 are on or close to the margin, only a few dozens achieve an $SSIM > 0.4$. This may indicate that there
 556 is a potential for much more images to reconstruct, and possibly with better quality.



(a) Full Images. Original samples from the training set (*red*) and reconstructed results (*blue*)



(b) Scatter plot (similar to Fig. 3) .

Figure 17: Reconstruction from a model trained on 50 images per class from the CIFAR100 dataset (100 classes, total of 5000 datapoints). The model is a 3-layer MLP with 10000 neurons in each layer.