
Appendix for “Residual Alignment: Uncovering the Mechanisms of Residual Networks”

Anonymous Author(s)

Affiliation

Address

email

1 **1 (RA2+3+4) Imply (RA1)**

2 **Theorem 3.1.** *In a pre-activation ResNet, assuming the Jacobian linearizations are exact and satisfy*
3 *(RA2+3+4), then (RA1) holds for the intermediate representations.*

4 *Proof.* In a pre-activation ResNet,

$$h_{i+1} = h_i + \mathcal{F}(h_i; \mathcal{W}_i).$$

5 Since Jacobian linearizations are exact, we have:

$$h_{i+1} = (I + J_i)h_i.$$

6 Recall, the singular value decomposition of J_i is given by

$$J_i = U_i S_i V_i^\top,$$

7 where U_i and V_i are the respective left and right singular vectors, and S_i is the singular value matrix.

8 Invoking (RA2),

$$J_i = U S_i U^\top,$$

9 and therefore

$$h_{i+1} = (I + U S_i U^\top)h_i.$$

10 Applying recursively the above equality leads to

$$h_k = \left(\prod_{i=1}^{k-1} (I + U S_i U^\top) \right) h_1 = U \left(\prod_{i=1}^{k-1} (I + S_i) \right) U^\top h_1. \quad (1)$$

11 For binary classification, (RA3) implies the Jacobians are rank 1 and therefore

$$h_k = U_{k,1} \left(\prod_{i=1}^{k-1} (1 + S_{i,1}) \right) U_{k,1}^\top h_1.$$

12 According to (RA4),

$$S_{i,1} = \frac{1}{i}$$

13 and

$$\prod_{i=1}^{k-1} (1 + S_{i,1}) = \prod_{i=1}^{k-1} (1 + 1/i) = k.$$

14 Substituting the above into Equation (1), we obtain

$$h_k = k U_{i,1} U_{i,1}^\top h_1.$$

15 This proves that the intermediate representations of a given input are *equispaced* on a *line* embedded
16 in high dimensional space, i.e., (RA1).

17

□

18 **2 Unconstrained Jacobians Model Leads to RA**

19 We start by providing motivation for the unconstrained Jacobians problem introduced in the main text.
 20 Assume a training example, x , is situated next to a point on the classification boundary, denoted by
 21 x_{mid} , satisfying $f(x_{\text{mid}}; \mathcal{W}) = 0$. Performing a Taylor expansion of the logits of x around x_{mid} yields

$$\begin{aligned} f(x; \mathcal{W}) &= f(x_{\text{mid}}; \mathcal{W}) + \left. \frac{\partial f(x; \mathcal{W})}{\partial x} \right|_{x_{\text{mid}}} \Delta_x + h(x, x_{\text{mid}}) \\ &= \left. \frac{\partial f(x; \mathcal{W})}{\partial x} \right|_{x_{\text{mid}}} \Delta_x + h(x, x_{\text{mid}}), \end{aligned} \quad (2)$$

22 where $\Delta_x = x - x_{\text{mid}}$ and $h(x, x_{\text{mid}})$ accounts for the approximation error, which is $\mathcal{O}(\|\Delta_x\|_2^2)$ and
 23 assumed to be negligible in our analysis.

24 Recall the loss associated with training a ResNet:

$$\text{minimize}_{\{\mathcal{W}_i\}_{i=1}^{L+1}} \frac{1}{N} \sum_{n=1}^N \mathcal{L}(f(x_n; \mathcal{W}), y_n) + \frac{\lambda}{2} \|\mathcal{W}\|_2^2. \quad (3)$$

25 Substituting Equation (2) into the above, neglecting the approximation error, and considering only
 26 the objective associated with the training sample x and its label y , we get

$$\mathcal{L} \left(\left. \frac{\partial f(x; \mathcal{W})}{\partial x} \right|_{x_{\text{mid}}} \Delta_x, y \right) + \frac{\lambda}{2} \|\mathcal{W}\|_2^2.$$

27 By using the chain rule, we can then obtain the following:¹

$$\mathcal{L} \left(\frac{\partial \mathcal{C}(h_{L+1}; \mathcal{W}_{L+1})}{\partial h_{L+1}} \prod_{i=1}^L \left(I + \frac{\partial \mathcal{F}(h_i; \mathcal{W}_i)}{\partial h_i} \right) \Delta_x, y \right) + \frac{\lambda}{2} \sum_{i=1}^L \|\mathcal{W}_i\|_F^2 + \frac{\lambda}{2} \|\mathcal{W}_{L+1}\|_F^2, \quad (4)$$

28 where all the Jacobians are evaluated at the point x_{mid} . This naturally gives rise to the following
 29 definition, as introduced in the main text.

30 **Definition 3.2 (Unconstrained Jacobians Model).** *Given a fixed input $\Delta_x \in \mathbb{R}^D$ and its label*
 31 *$y \in \{+1, -1\}$, find matrices $J_i \in \mathbb{R}^{D \times D}$, $1 \leq i \leq L$, and vector $w \in \mathbb{R}^D$ that*

$$\text{minimize}_{w, \{J_i\}_{i=1}^L} \mathcal{L}(w^\top \prod_{i=1}^L (I + J_i) \Delta_x, y) + \frac{\lambda}{2} \sum_{i=1}^L \|J_i\|_F^2 + \frac{\lambda}{2} \|w\|_2^2.$$

32 In the main text, we stated the following theorem.

33 **Theorem 3.3.** *There exists a global optimum of the Unconstrained Jacobians Model where the top*
 34 *Jacobian singular vectors are aligned (RA2), all Jacobians are rank one, analogous to (RA3), and*
 35 *the top Jacobian singular values are equal, analogous to (RA4).*

36 *Proof.* Throughout the proof, we assume, without loss of generality, that the label is $y = 1$. Using
 37 the cyclic property of the trace, the logit, $w^\top \prod_{i=1}^L (I + J_i) \Delta_x$, equals

$$\text{tr} \left\{ \Delta_x w^\top \prod_{i=1}^L (I + J_i) \right\}.$$

38 Denoting the power set of all natural numbers between 1 and L by $\mathcal{P}(L)$, the above can be expressed
 39 as follows:

$$\sum_{s \in \mathcal{P}(L)} \text{tr} \left\{ \Delta_x w^\top \prod_{i \in s} J_i \right\}.$$

40 Each element in the above summation can be upper bounded through the following theorem (a
 41 generalization of Von Neumann's trace inequality [Mirsky, 1975] to the product of more than two
 42 real matrices).

¹For simplicity, we ignore the input transformation that maps the input Δ_x to the first representation $h_1 \in \mathbb{R}^D$ by simply assuming $\Delta_x \in \mathbb{R}^D$.

43 **Theorem 2** ([Miranda and Thompson, 1993]). *Let A_1, \dots, A_m be matrices with real entries. Take*
44 *the singular values of A_j to be $s_1(A_j) \geq \dots \geq s_n(A_j)$, for $j = 1, \dots, m$, and denote $S_j =$*
45 *$\text{diag}(s_1(A_j), \dots, s_n(A_j))$. Then, as the matrices P_1, \dots, P_m range over all possible rotations, i.e.,*
46 *the special orthogonal group $\text{SO}(n)$,*

$$\begin{aligned} & \sup_{P_1, \dots, P_m \in \text{SO}(n)} \text{tr}(A_1 P_1 \dots A_m P_m) \\ &= \sum_{i=1}^{n-1} \prod_{j=1}^m s_i(A_j) + [\text{sign det}(A_1 \dots A_m)] \prod_{j=1}^m s_n(A_j). \end{aligned}$$

47 *Moreover, assuming $\text{sign det}(A_1 \dots A_m) = 1$,*

$$\sup_{P_1, \dots, P_m \in \text{SO}(n)} \text{tr}(A_1 P_1 \dots A_m P_m) = \text{tr} \left\{ \prod_{i=1}^m S_i \right\}.$$

48 We will continue our proof using contradiction. Suppose all existing global optima of the uncon-
49 strained Jacobians problem consist of Jacobians that do not have aligning singular vectors, or do
50 not have equal singular values, or are not rank 1. Then take any solution $\{J_i\}_{i=1}^L$ and w . Using the
51 singular value decomposition, we have

$$J_i = U_i S_i V_i^\top, \quad \text{for } i = 1, \dots, L,$$

52 and

$$\Delta_x w^\top = U_{L+1} S_{L+1} V_{L+1}^\top.$$

53 Then Theorem 2 implies

$$\begin{aligned} \sum_{s \in \mathcal{P}(L)} \text{tr} \left\{ \Delta_x w^\top \prod_{i \in s} J_i \right\} &\leq \sum_{s \in \mathcal{P}(L)} \text{tr} \left\{ S_{L+1} \prod_{i \in s} S_i \right\} \\ &= \text{tr} \left\{ S_{L+1} \prod_{i=1}^L (I + S_i) \right\}. \end{aligned}$$

54 For all $s \in \mathcal{P}(L)$, the inequality becomes equality once the singular vectors of all the Jacobians align
55 with those of $\Delta_x w^\top$ and once the vector w is chosen to be proportional to Δ_x (so that the matrix
56 $\Delta_x w^\top$ is symmetric). The implication of the steps thus far is that one can increase, or at least keep
57 constant the logit, and consequently reduce, or at least keep constant the loss by simply aligning the
58 singular vectors of the Jacobians. In addition, since the regularization term $\|J_i\|_F^2 = \text{tr}\{S_i^2\}$, this
59 change of Jacobians does not affect the regularization terms.

60 Notice that $\Delta_x w^\top$ is a rank one matrix and so S_{L+1} has a single non-zero diagonal entry. Furthermore,
61 the matrices S_i , for $1 \leq i \leq L$, are all diagonal. As such, we can zero out all their other diagonal
62 entries and leave a single non-zero entry at the location that S_{L+1} has one, which does not affect the
63 logits but reduces the regularization terms.

64 Using the inequality of arithmetic and geometric means on this only non-zero entry s_i of every
65 diagonal matrix S_i gives

$$s_{L+1} \prod_{i=1}^L (1 + s_i) \leq s_{L+1} \left(1 + \frac{1}{L} \sum_{i=1}^L s_i \right)^L.$$

66 The implication of the above inequality is that, once the singular vectors of the Jacobians are aligned,
67 one can further increase the logits and reduce the loss by averaging all the top singular values, s_i
68 for $1 \leq i \leq L + 1$, and forcing them to be equal. Furthermore, since $\|J_i\|_F^2 = \text{tr}\{S_i^2\} = s_i^2$ is
69 convex, by Jensen's inequality, averaging the singular values only decreases the value of the Jacobian
70 regularization.

71 All in all, we obtain higher, or at least no lower logit, and lower, or at least no higher loss when
72 all singular vectors are aligned, all top singular values are equal and all other singular values are
73 zero, which contradicts the statement that no global optima of the unconstrained Jacobians problem
74 satisfies all of these conditions. \square

75 **3 Additional Empirical Evidence**

76 **3.1 (RA1)**

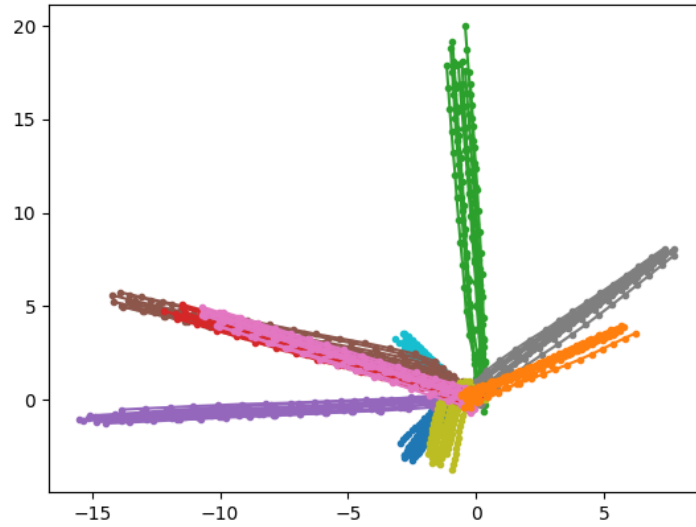


Figure 1: Fully-connected ResNet34 (Type 1 model) trained on MNIST.

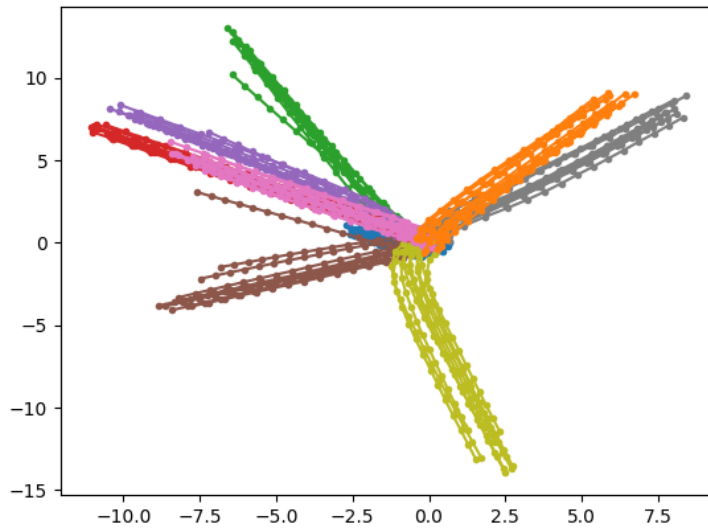


Figure 2: Fully-connected ResNet34 (Type 1 model) trained on FashionMNIST.

77

78

79

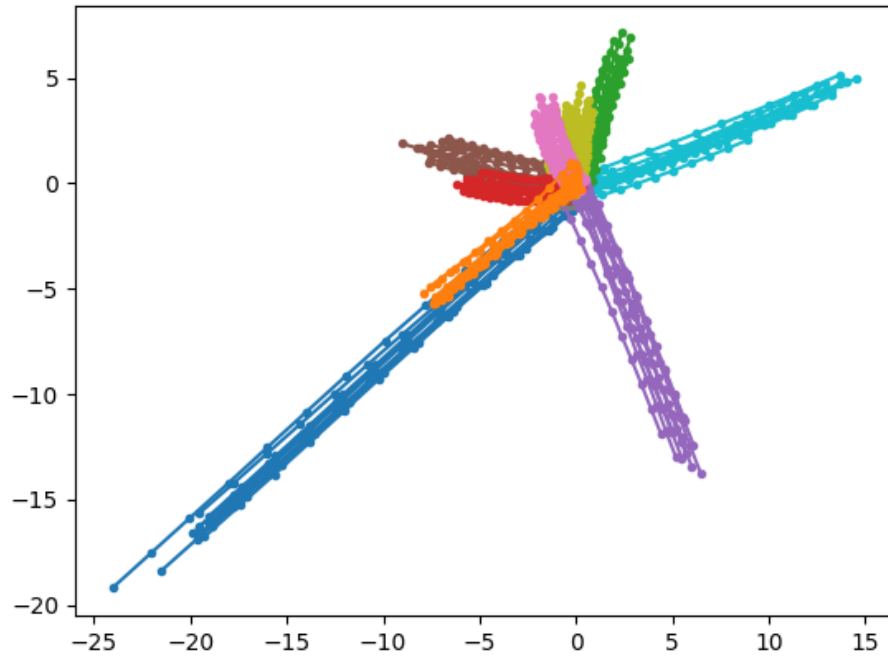


Figure 3: Fully-connected ResNet34 (Type 1 model) trained on CIFAR10.

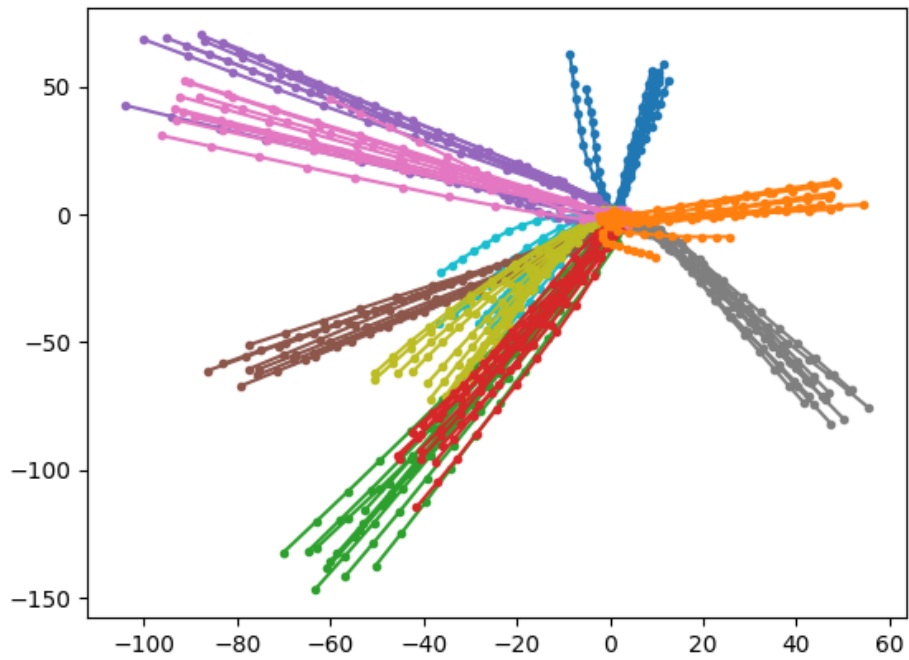


Figure 4: Convolutional ResNet34 (Type 2 model) trained on MNIST.

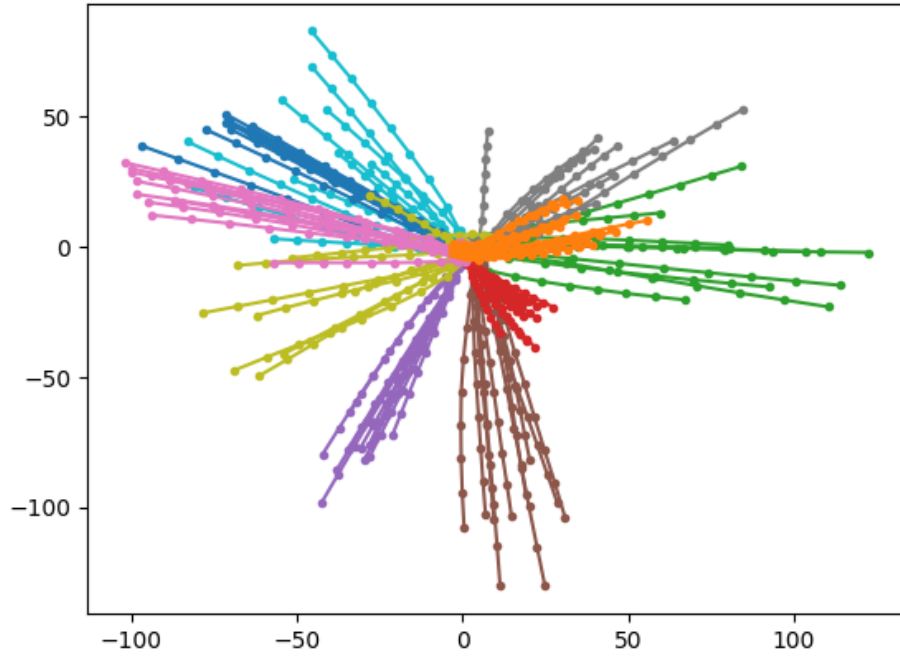


Figure 5: Convolutional ResNet34 (Type 2 model) trained on FashionMNIST.

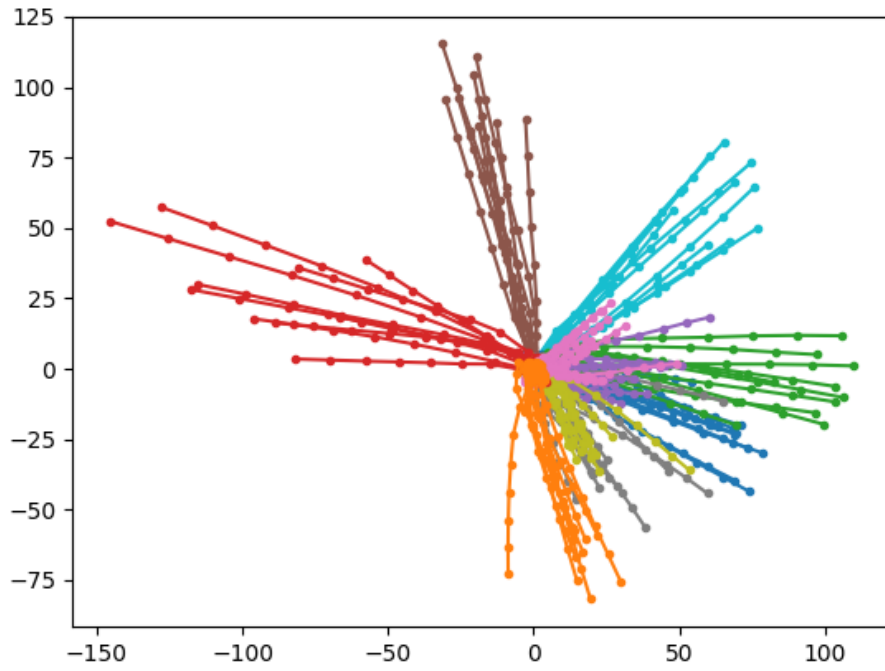


Figure 6: Convolutional ResNet34 (Type 2 model) trained on CIFAR10.

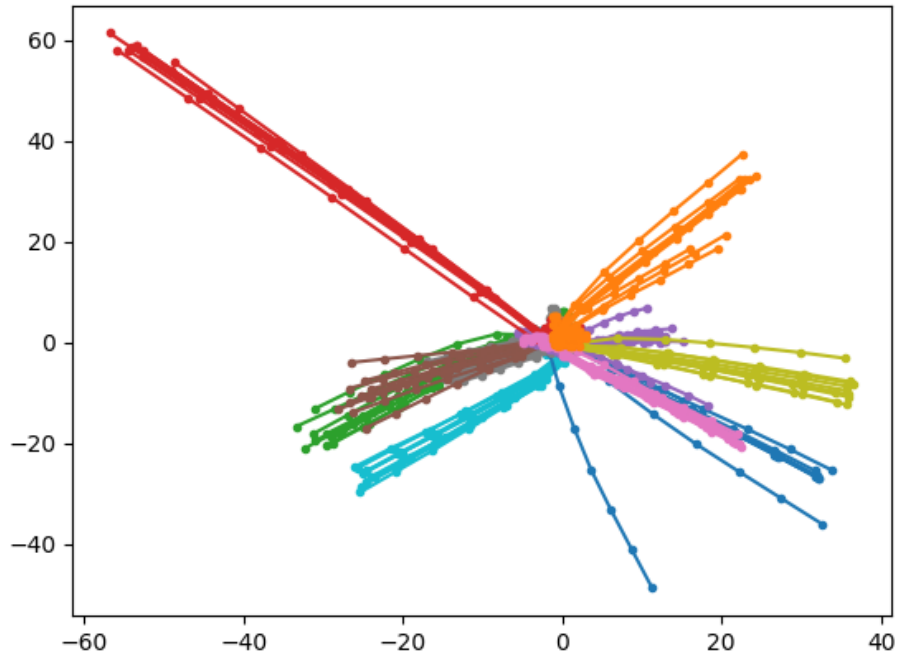


Figure 7: Convolutional ResNet34 with downsampling (Type 3 model) trained on MNIST.

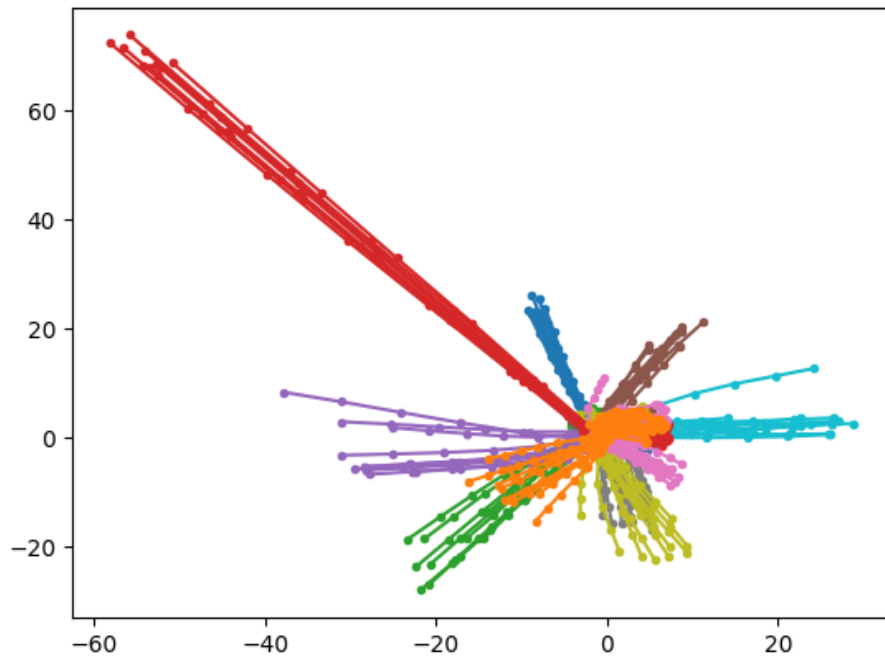


Figure 8: Convolutional ResNet34 with downsampling (Type 3 model) trained on FashionMNIST.

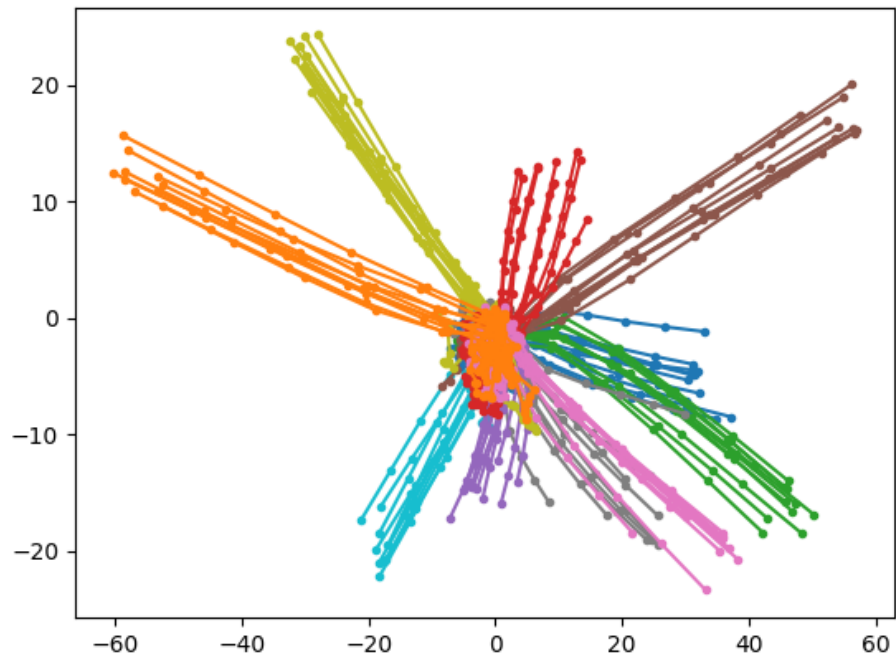


Figure 9: Convolutional ResNet34 with downsampling (Type 3 model) trained on CIFAR10.

80 **3.2** (RA2) : $U_{j,K}^\top J_i V_{j,K}$

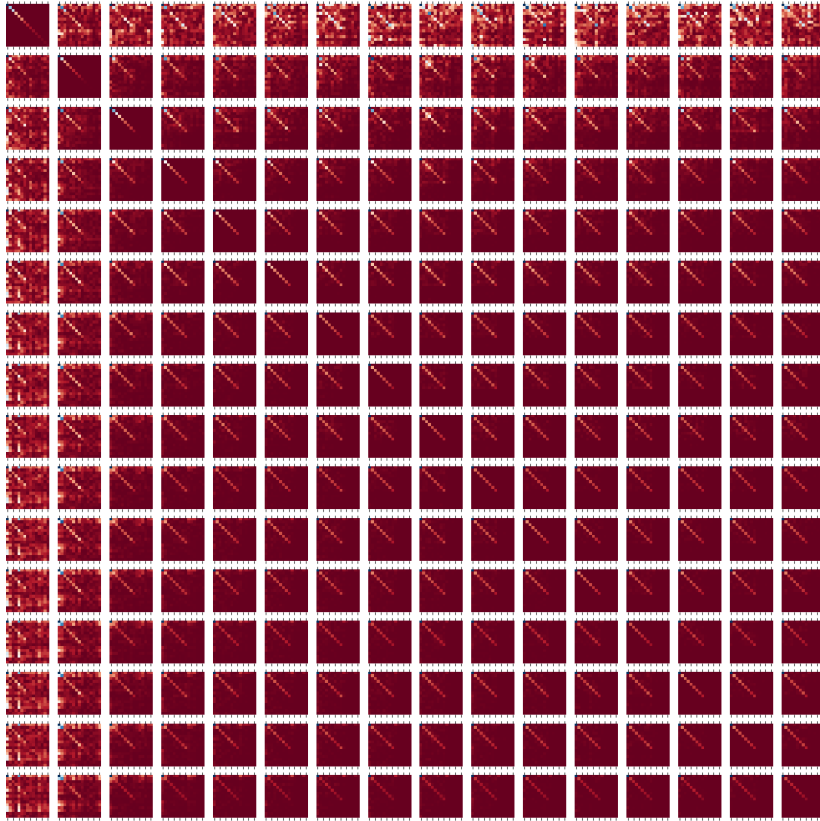


Figure 10: Fully-connected ResNet34 (Type 1 model) trained on MNIST.

81

82

83

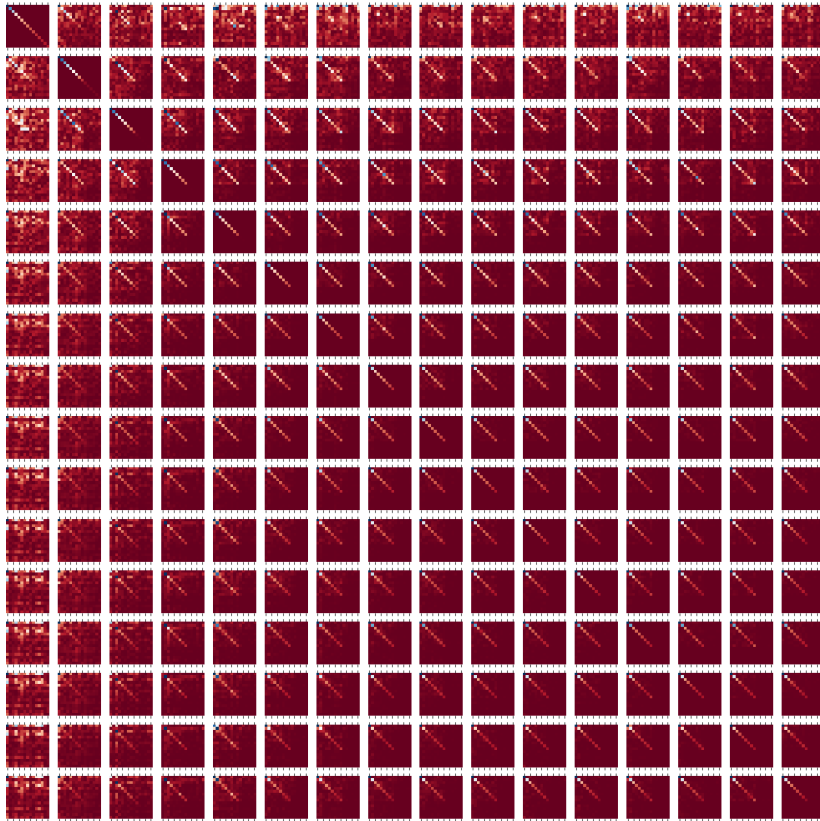


Figure 11: Fully-connected ResNet34 (Type 1 model) trained on FashionMNIST.

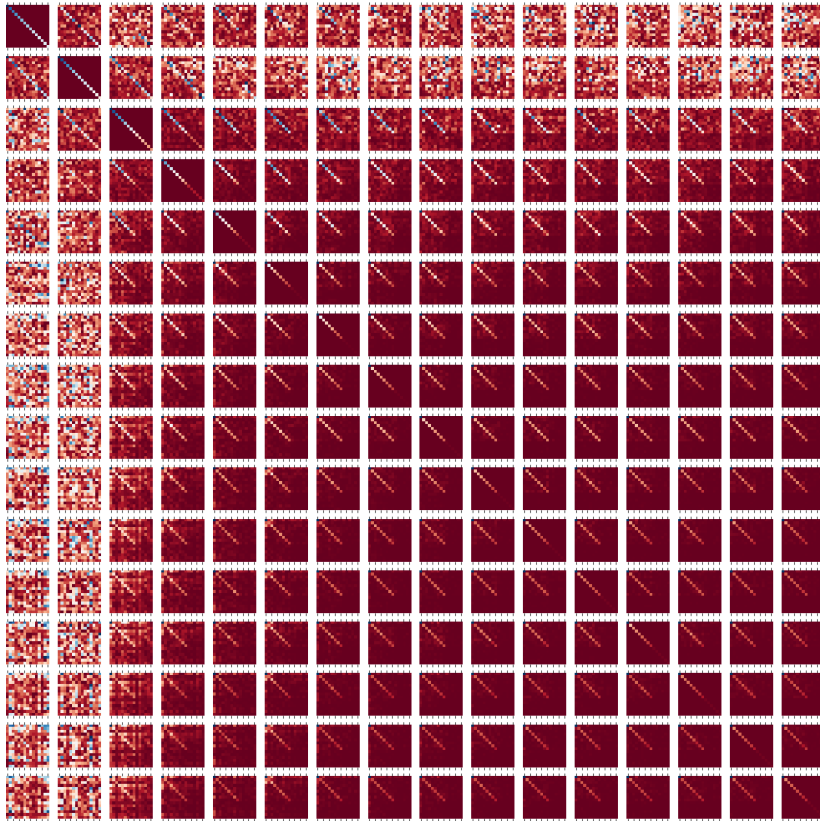


Figure 12: Fully-connected ResNet34 (Type 1 model) trained on CIFAR10.

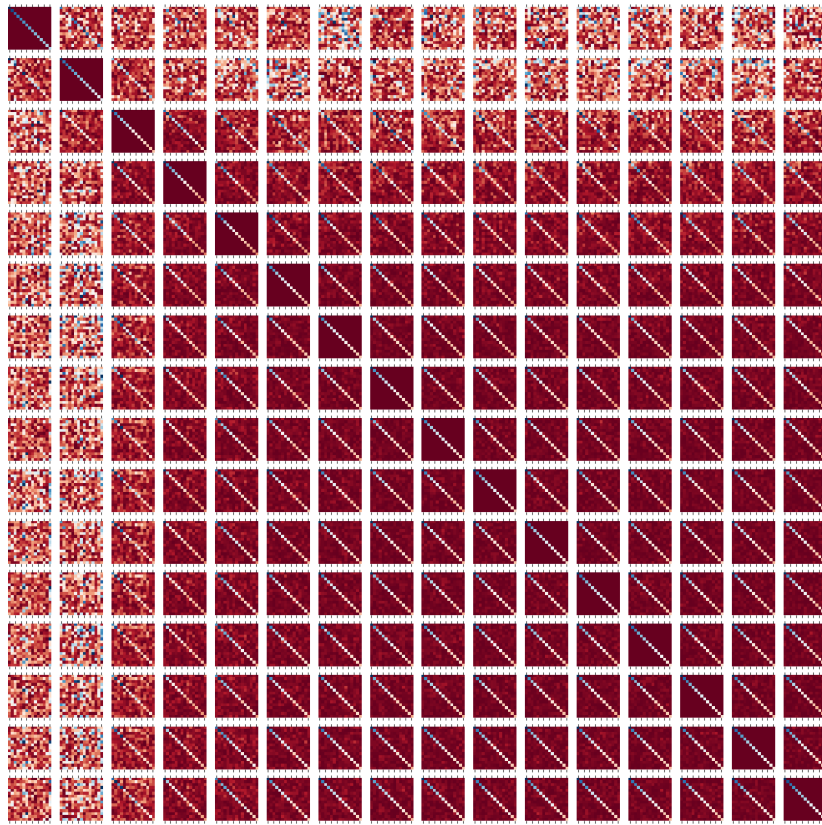


Figure 13: Fully-connected ResNet34 (Type 1 model) trained on CIFAR100.

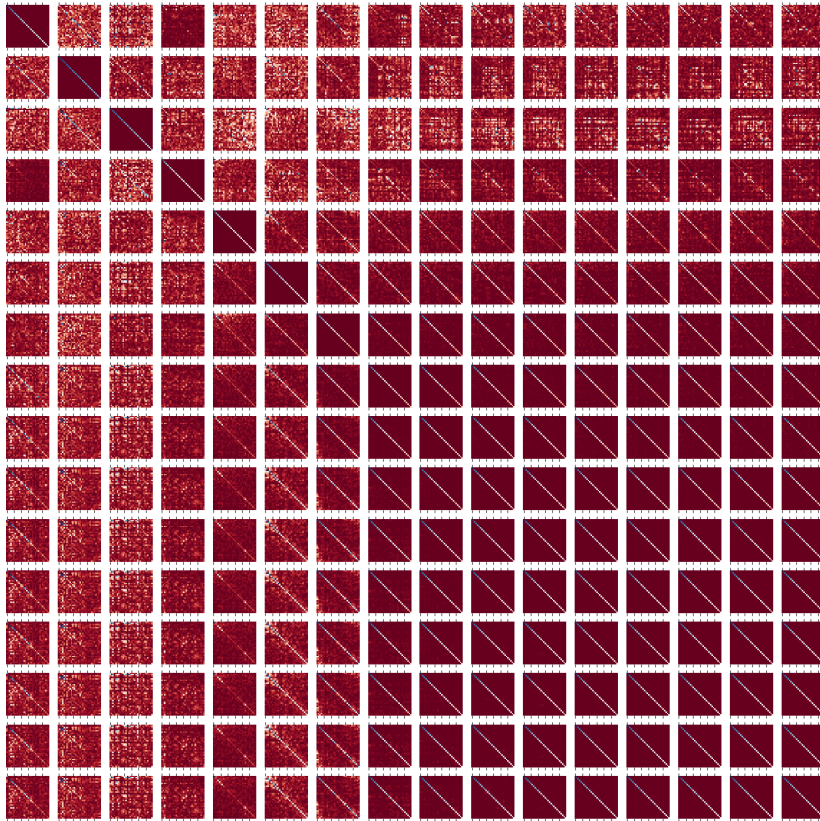


Figure 14: Convolutional ResNet34 (Type 2 model) trained on MNIST.

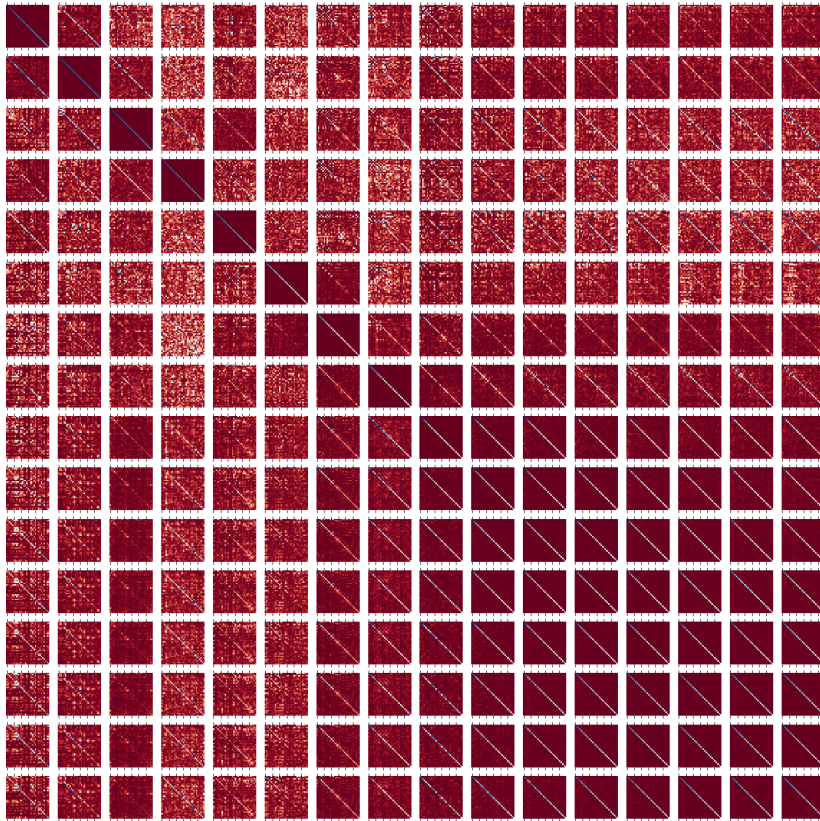


Figure 15: Convolutional ResNet34 (Type 2 model) trained on FashionMNIST.

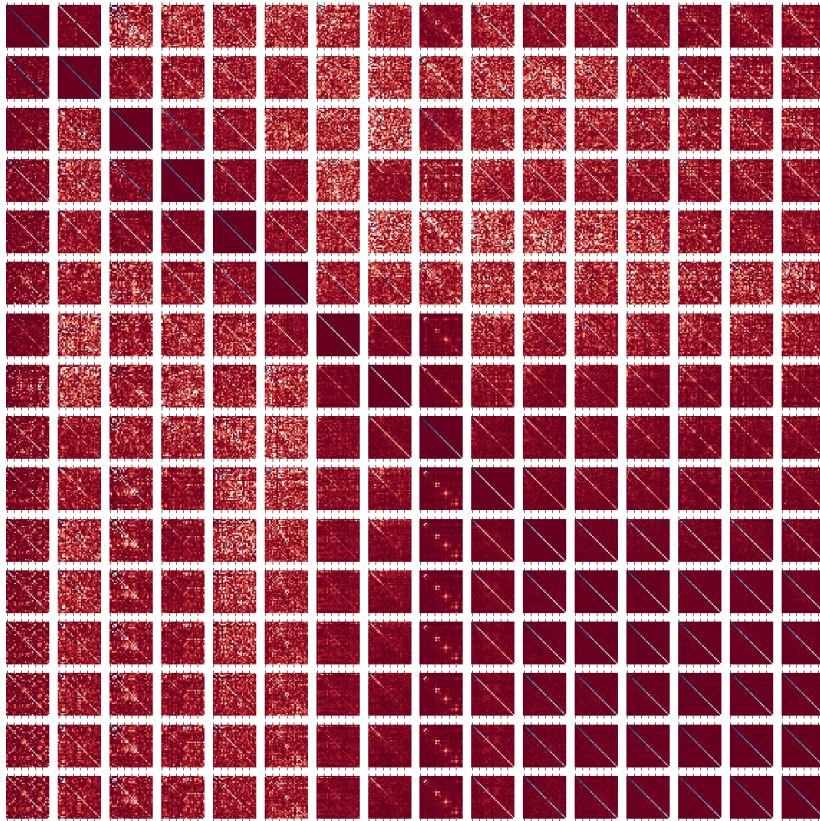


Figure 16: Convolutional ResNet34 (Type 2 model) trained on CIFAR10.

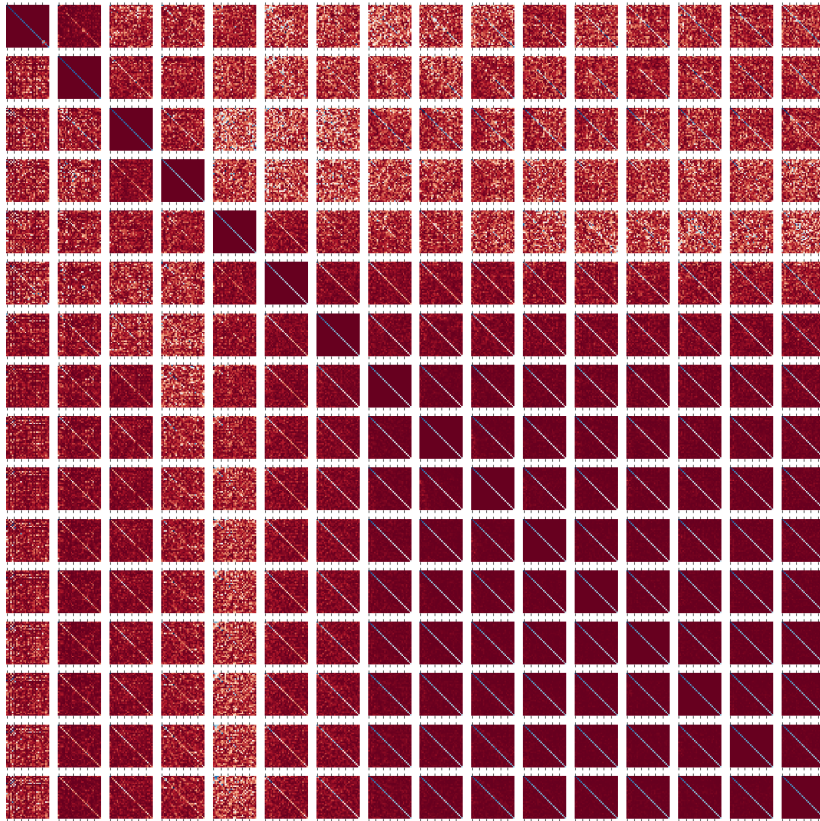


Figure 17: Convolutional ResNet34 (Type 2 model) trained on CIFAR100.

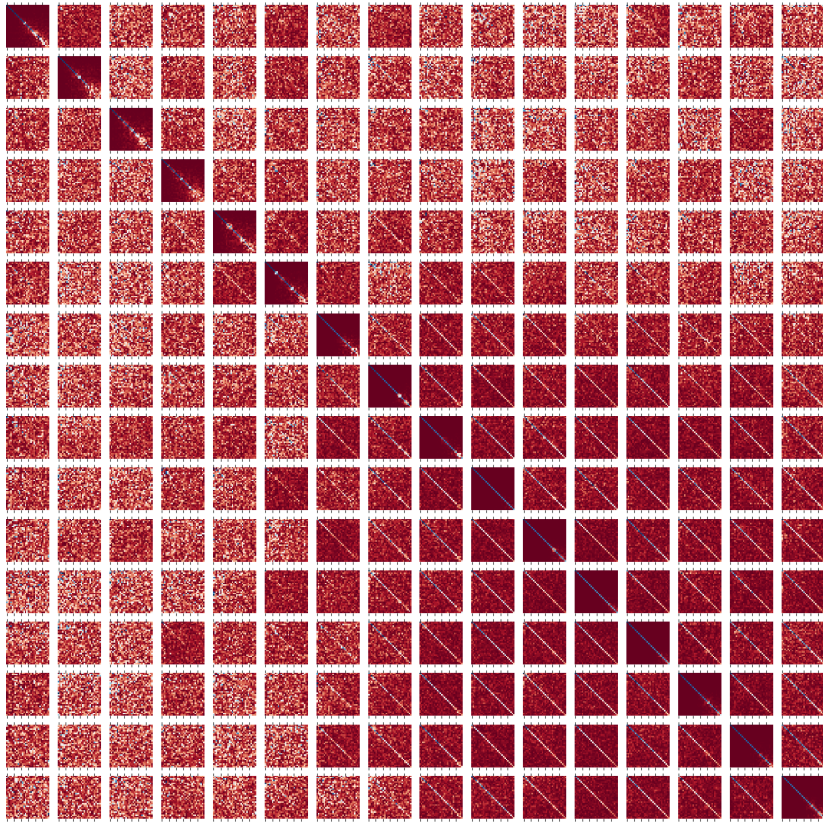


Figure 18: Convolutional ResNet34 (Type 2 model) trained on ImageNette.

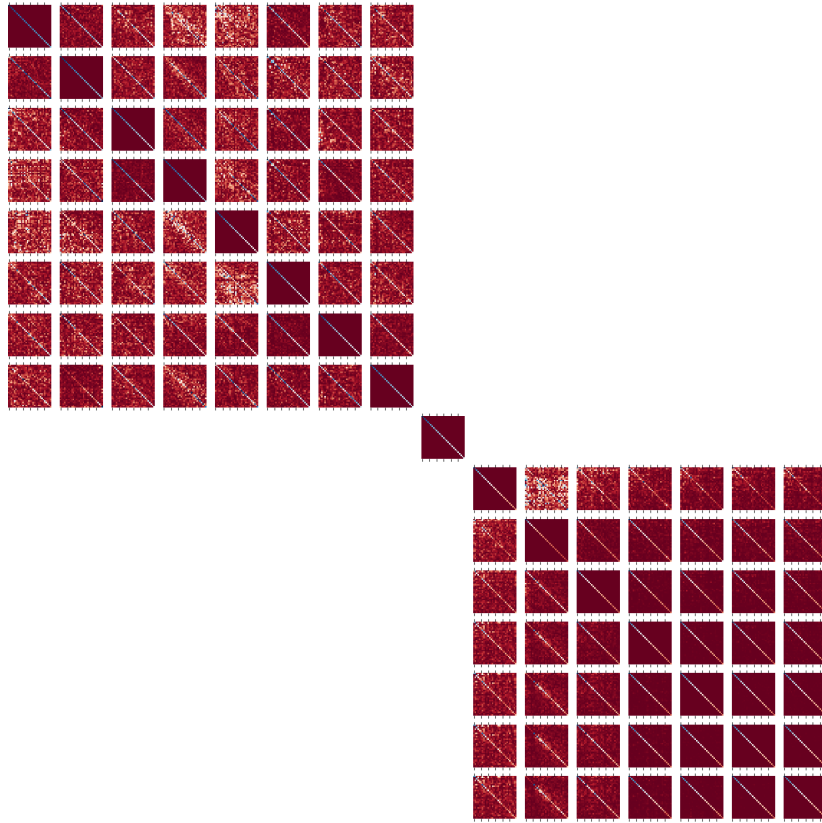


Figure 19: Convolutional ResNet34 with downsampling (Type 3 model) trained on MNIST.

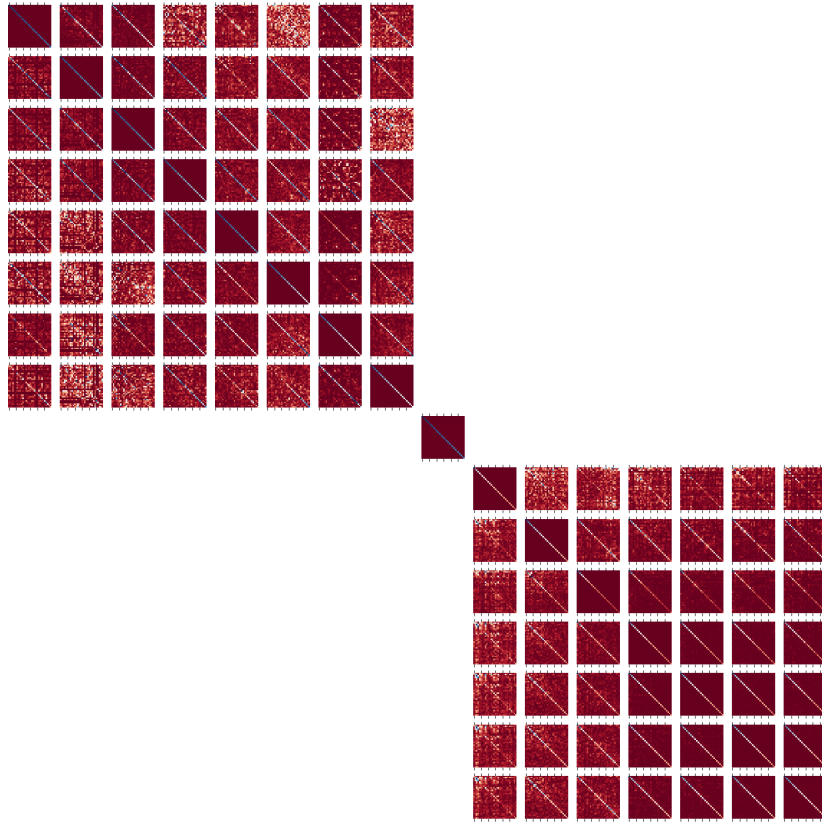


Figure 20: Convolutional ResNet34 with downsampling (Type 3 model) trained on FashionMNIST.

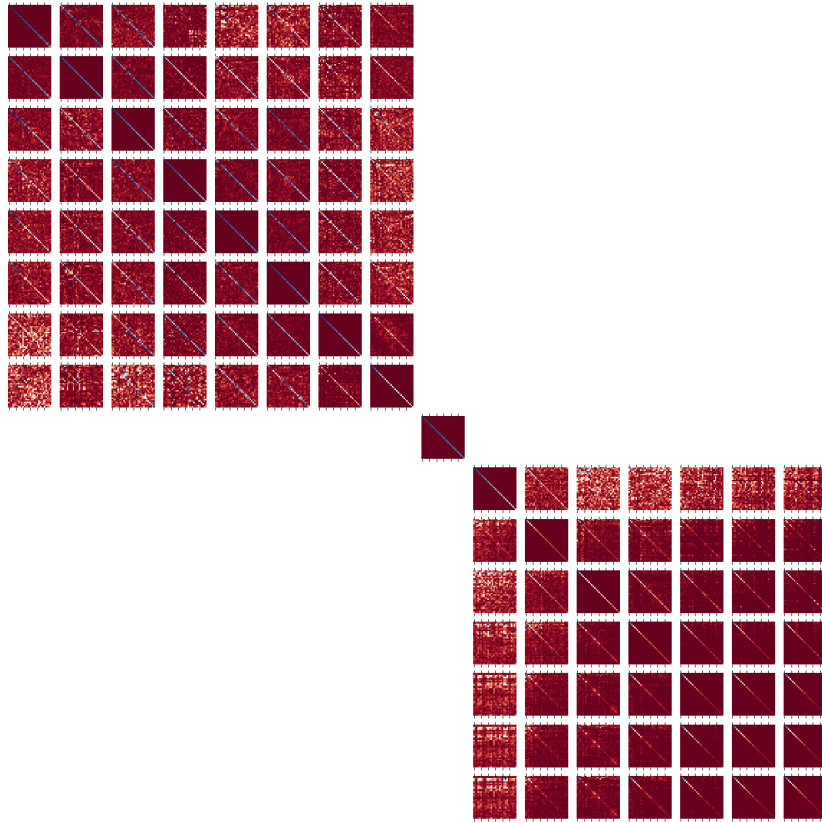


Figure 21: Convolutional ResNet34 with downsampling (Type 3 model) trained on CIFAR10.

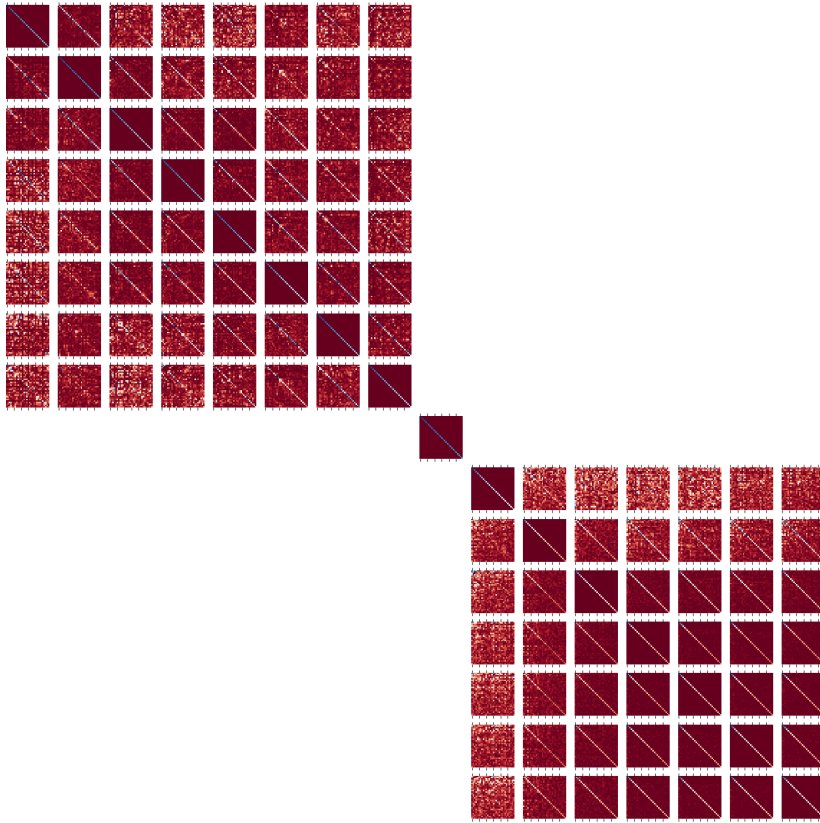


Figure 22: Convolutional ResNet34 with downsampling (Type 3 model) trained on CIFAR100.

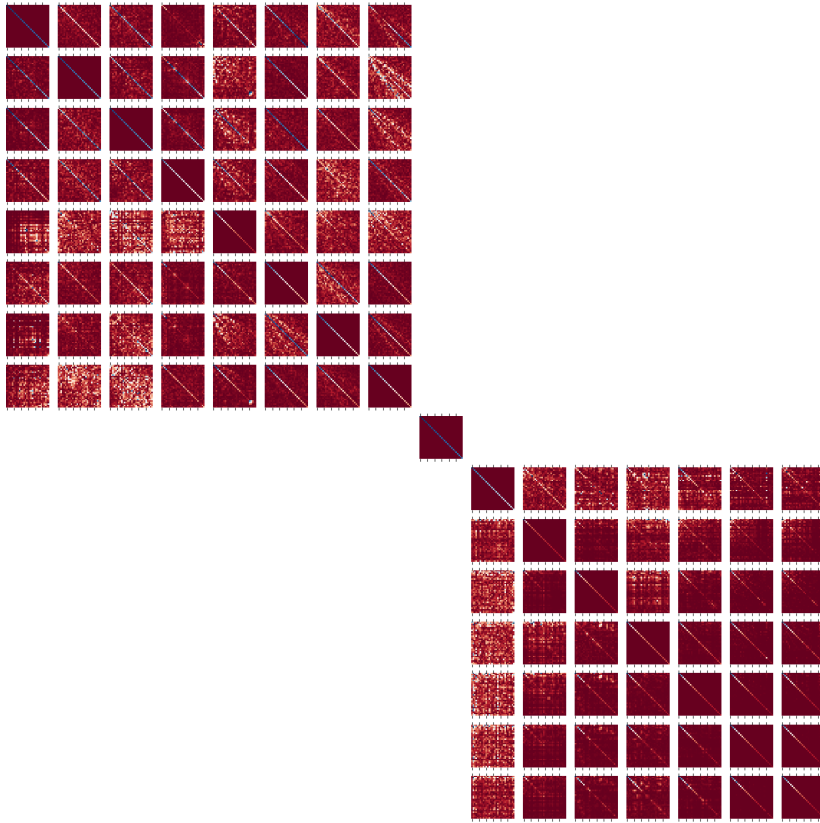


Figure 23: Convolutional ResNet34 with downsampling (Type 3 model) trained on ImageNette.

84 **3.3** (RA2) : $V_{j,K}^\top J_i U_{j,K}$

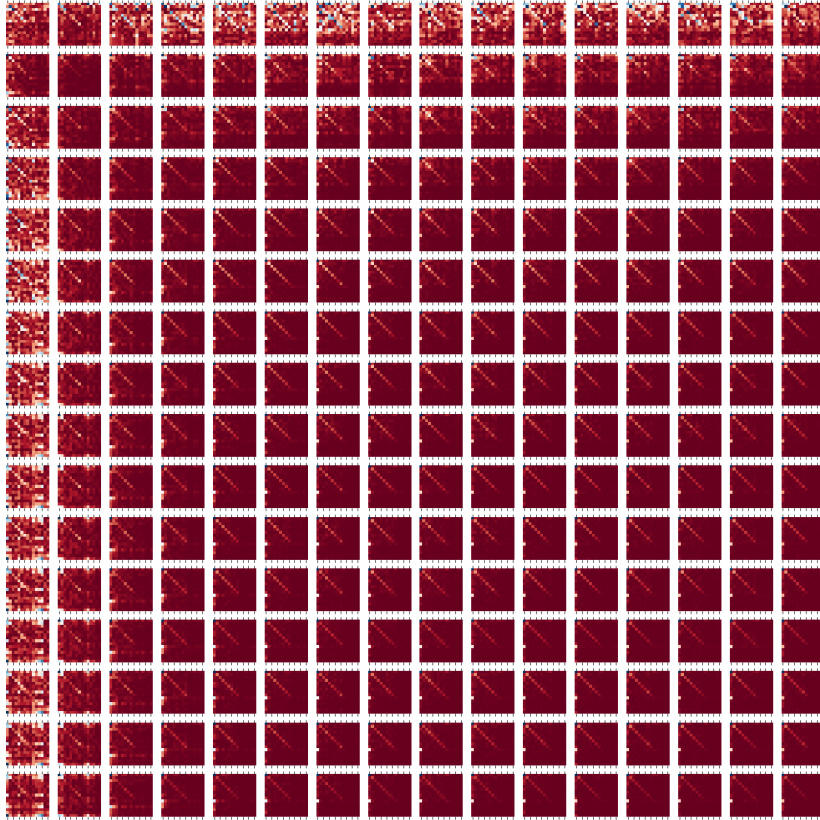


Figure 24: Fully-connected ResNet34 (Type 1 model) trained on MNIST.

85

86

87

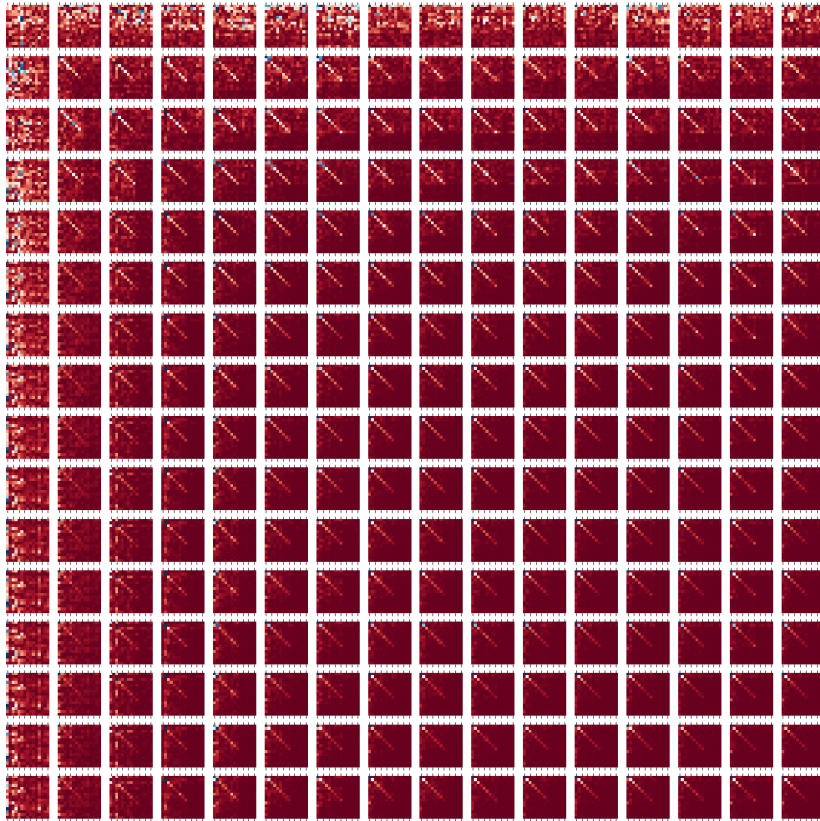


Figure 25: Fully-connected ResNet34 (Type 1 model) trained on FashionMNIST.

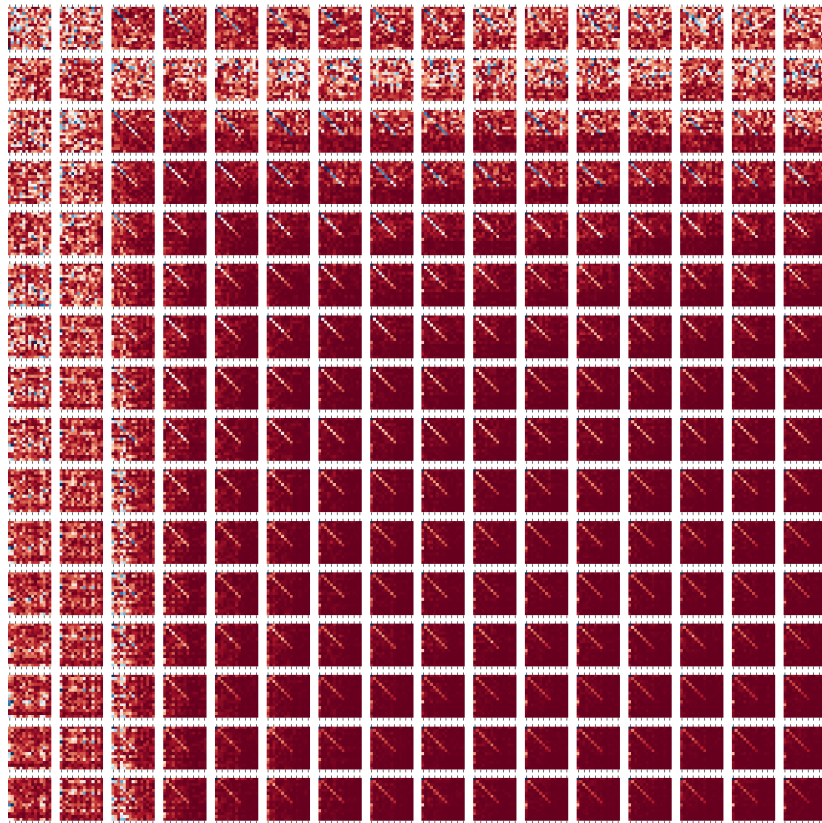


Figure 26: Fully-connected ResNet34 (Type 1 model) trained on CIFAR10.

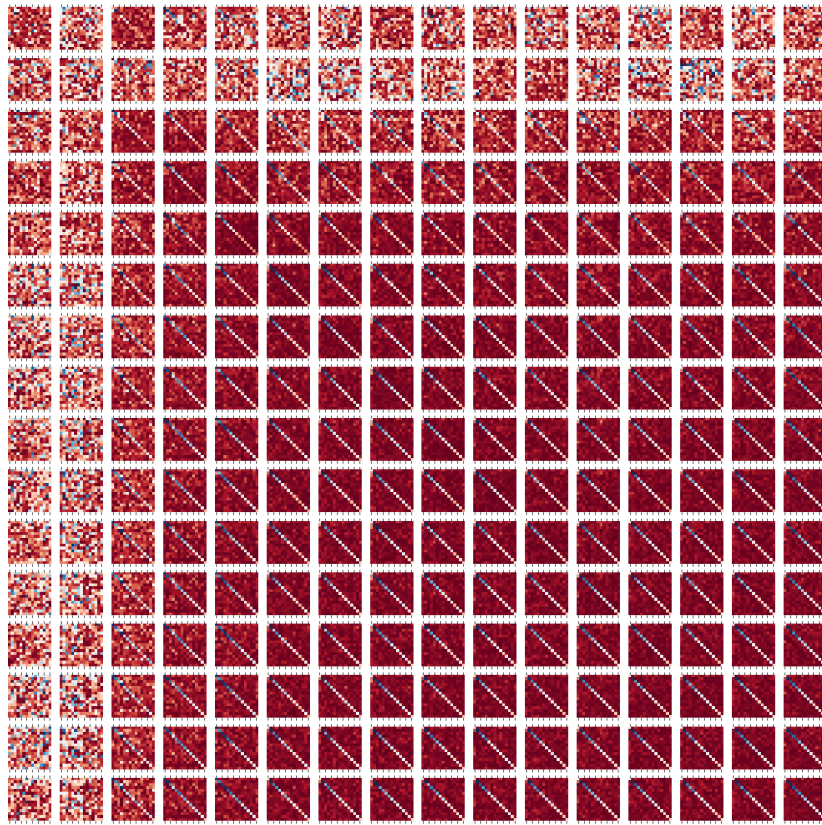


Figure 27: Fully-connected ResNet34 (Type 1 model) trained on CIFAR100.

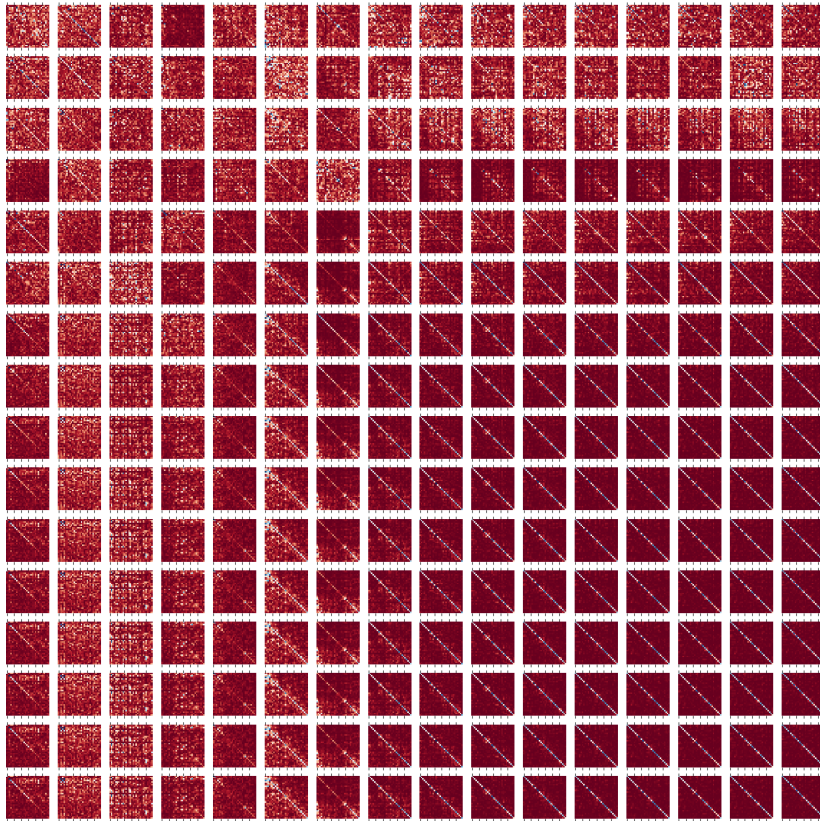


Figure 28: Convolutional ResNet34 (Type 2 model) trained on MNIST.

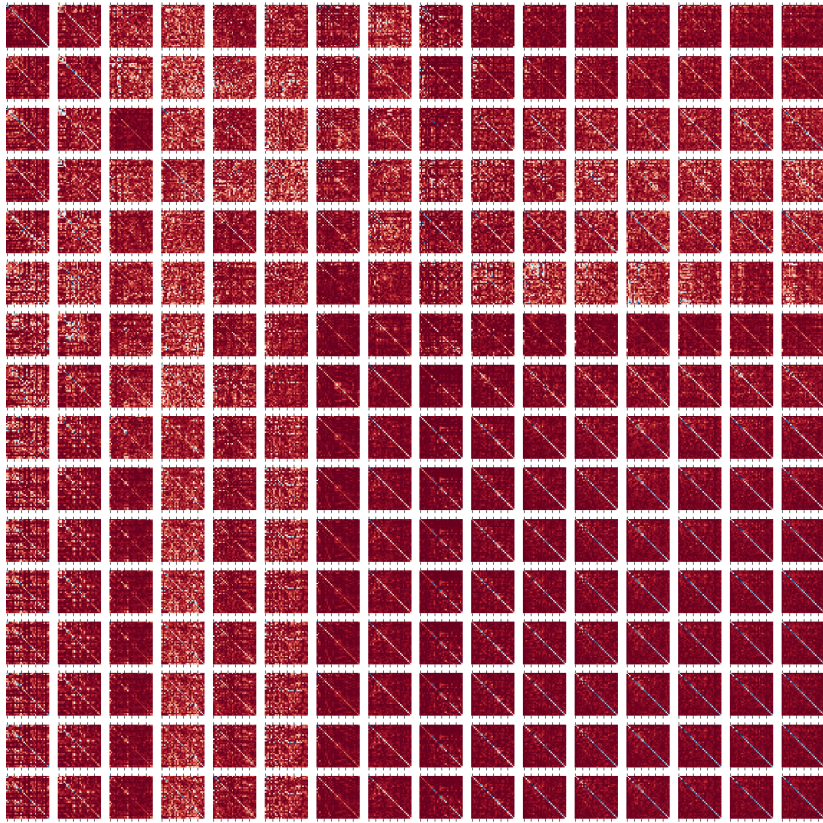


Figure 29: Convolutional ResNet34 (Type 2 model) trained on FashionMNIST.

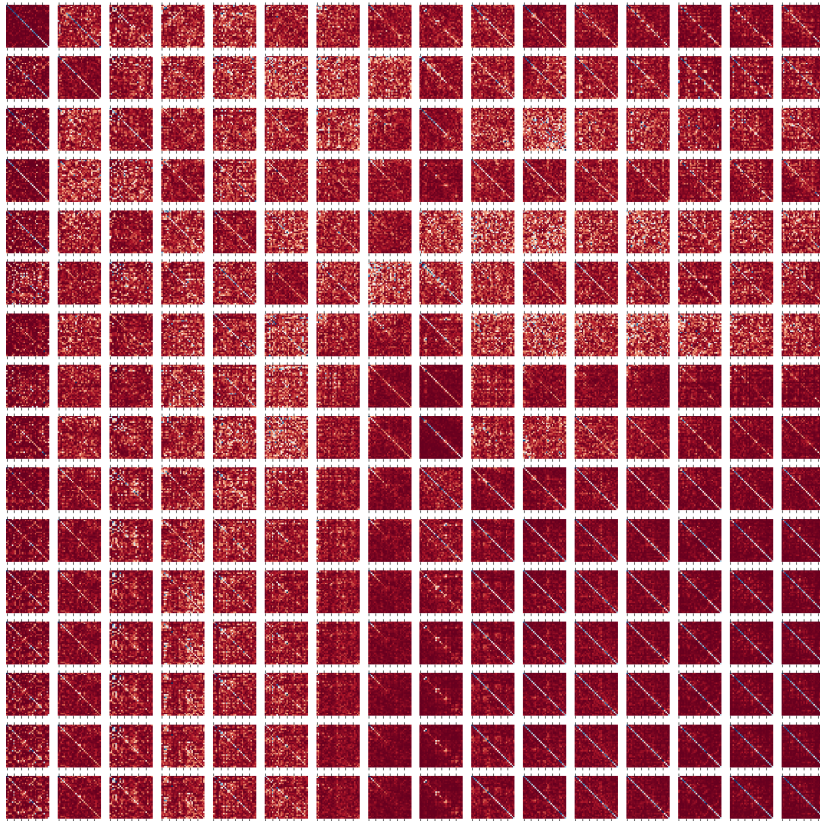


Figure 30: Convolutional ResNet34 (Type 2 model) trained on CIFAR10.

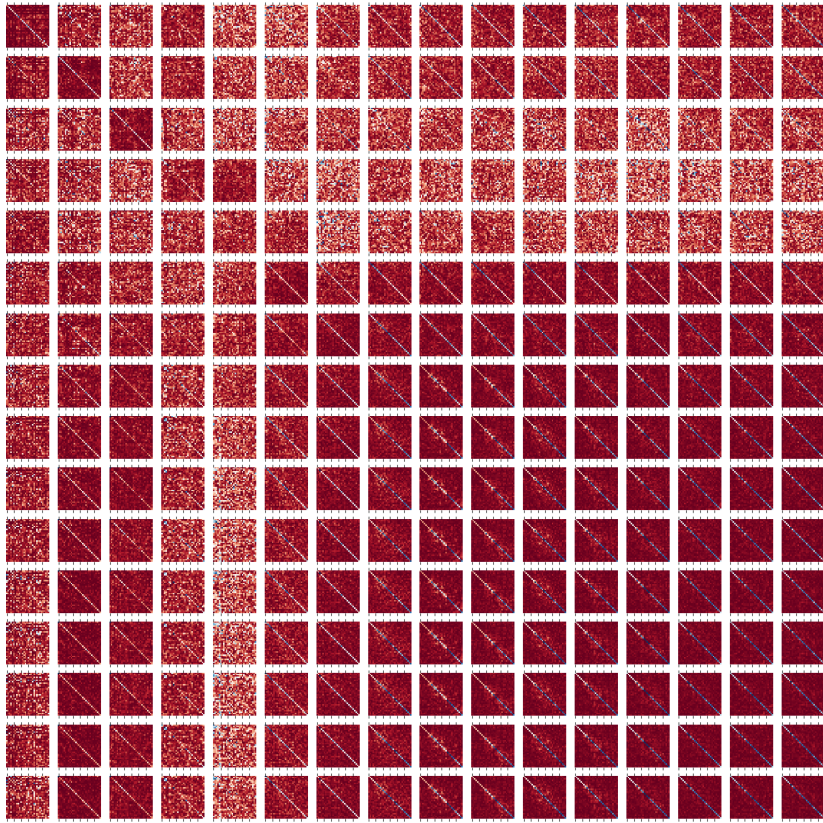


Figure 31: Convolutional ResNet34 (Type 2 model) trained on CIFAR100.

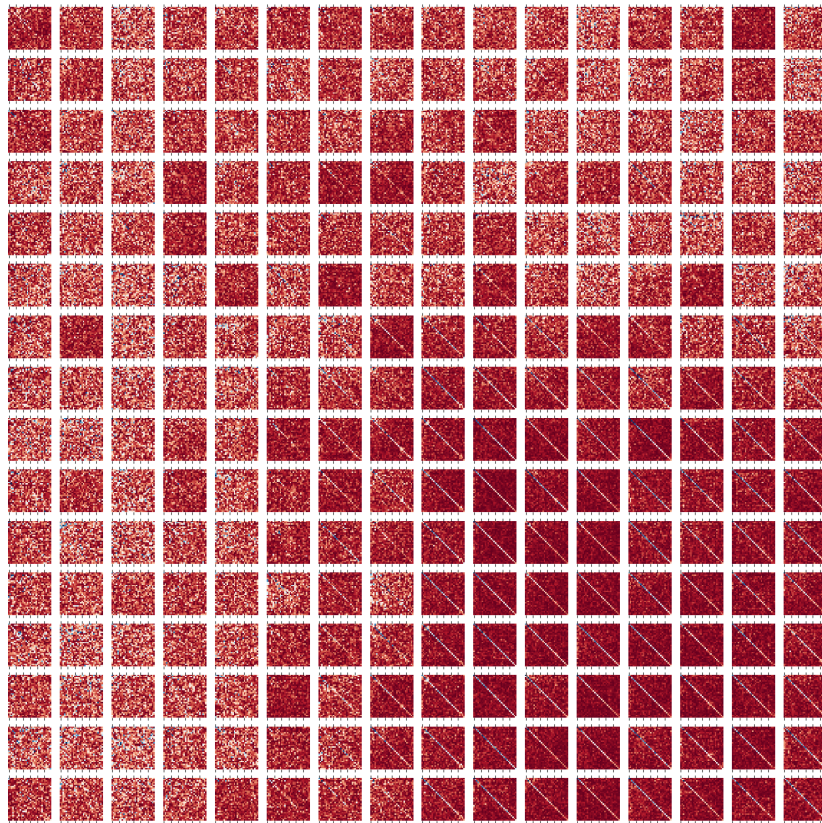


Figure 32: Convolutional ResNet34 (Type 2 model) trained on ImageNette.

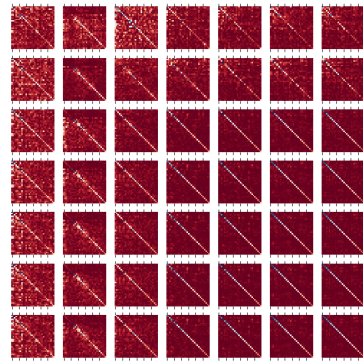
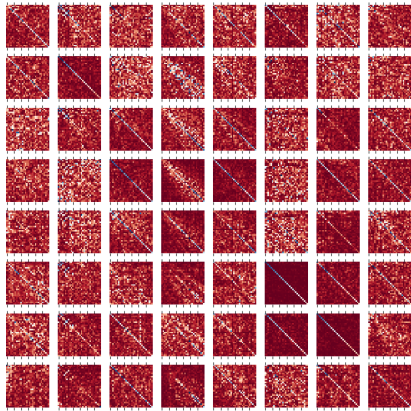


Figure 33: Convolutional ResNet34 with downsampling (Type 3 model) trained on MNIST.

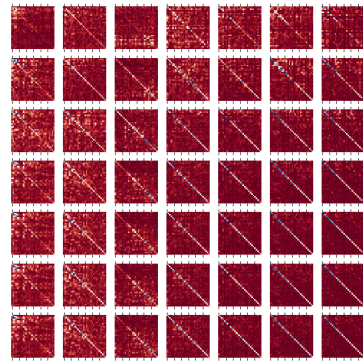
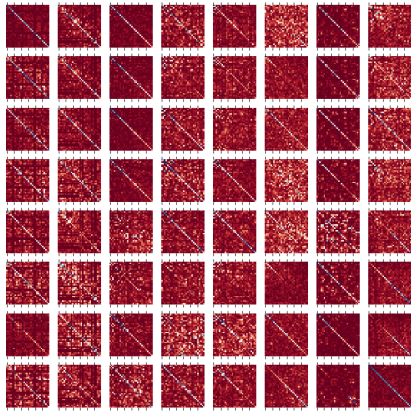


Figure 34: Convolutional ResNet34 with downsampling (Type 3 model) trained on FashionMNIST.

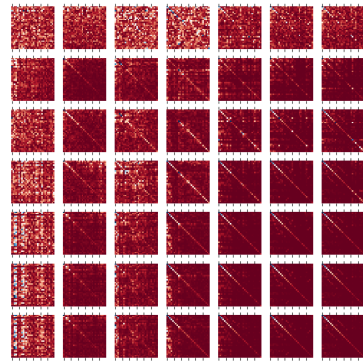
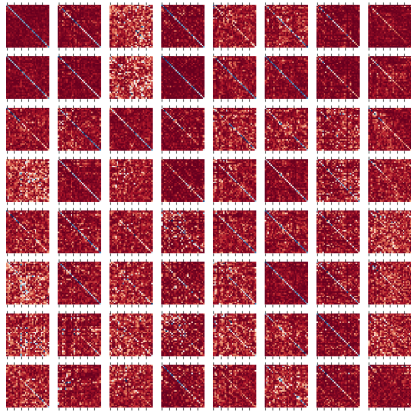


Figure 35: Convolutional ResNet34 with downsampling (Type 3 model) trained on CIFAR10.

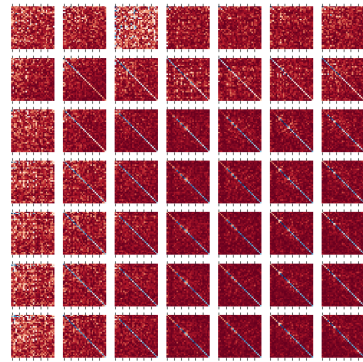
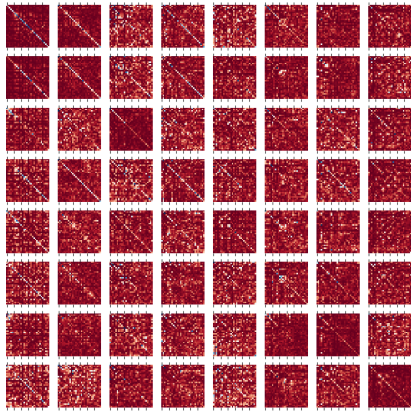


Figure 36: Convolutional ResNet34 with downsampling (Type 3 model) trained on CIFAR100.

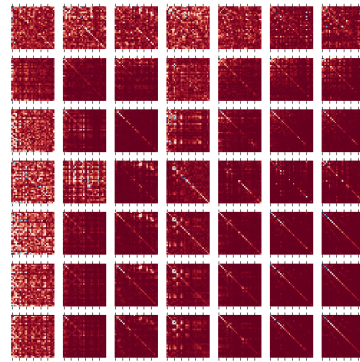
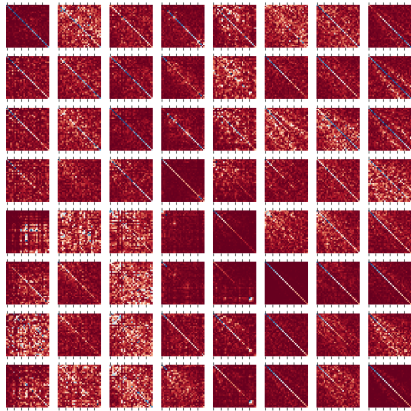


Figure 37: Convolutional ResNet34 with downsampling (Type 3 model) trained on ImageNette.

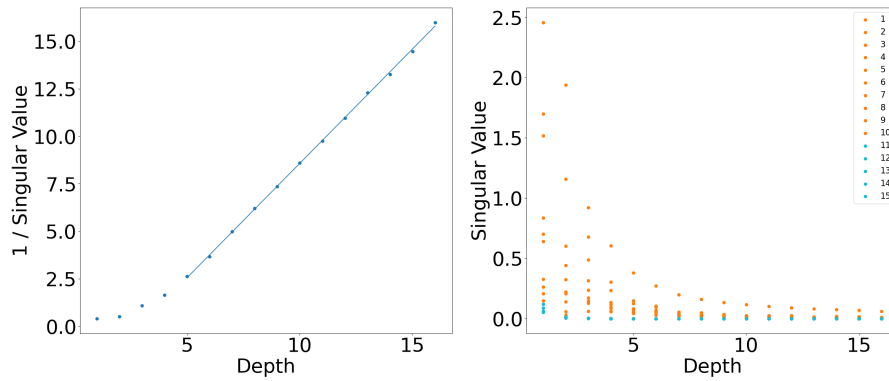


Figure 38: Fully-connected ResNet34 (Type 1 model) trained on MNIST.

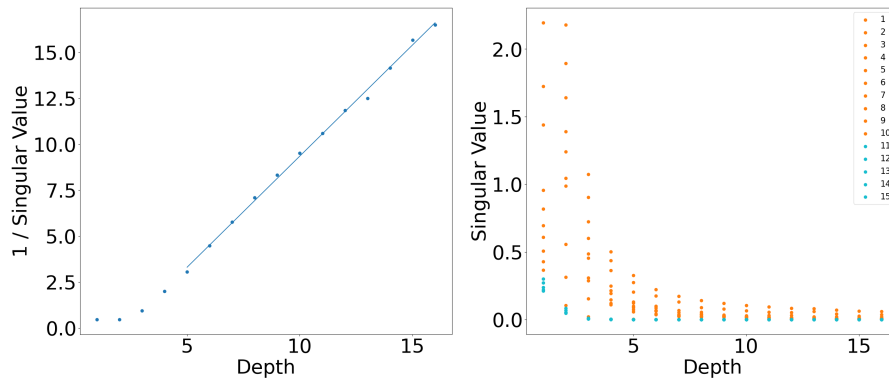


Figure 39: Fully-connected ResNet34 (Type 1 model) trained on FashionMNIST.

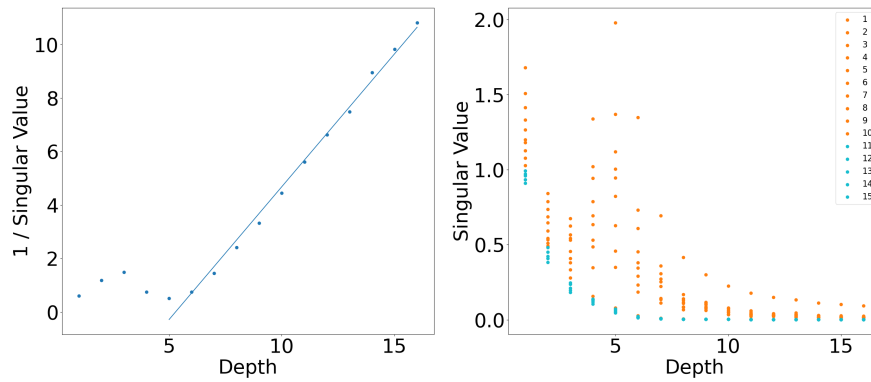


Figure 40: Fully-connected ResNet34 (Type 1 model) trained on CIFAR10.

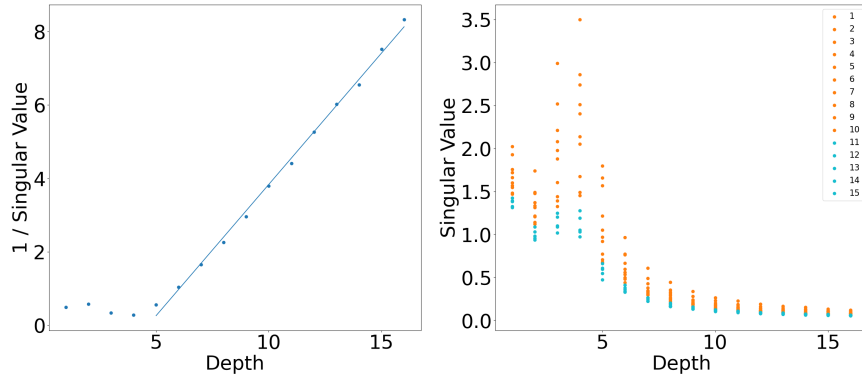


Figure 41: Fully-connected ResNet34 (Type 1 model) trained on CIFAR100.

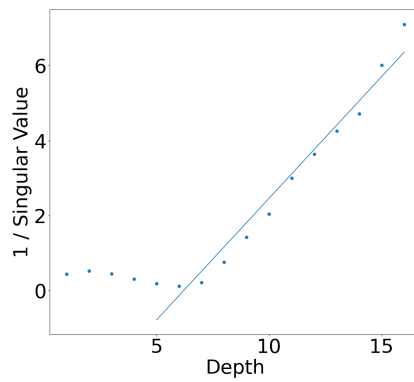


Figure 42: Convolutional ResNet34 (Type 2 model) trained on MNIST.

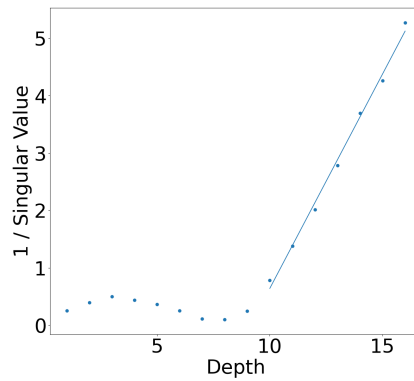


Figure 43: Convolutional ResNet34 (Type 2 model) trained on FashionMNIST.

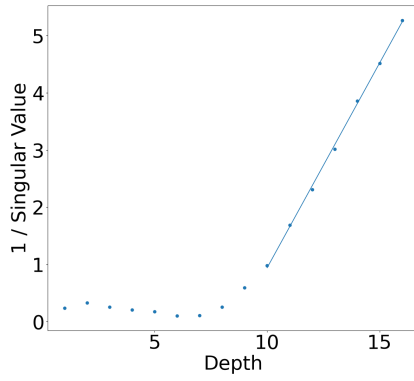


Figure 44: Convolutional ResNet34 (Type 2 model) trained on CIFAR10.

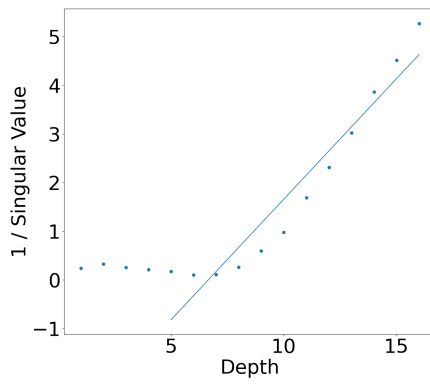


Figure 45: Convolutional ResNet34 (Type 2 model) trained on CIFAR100.

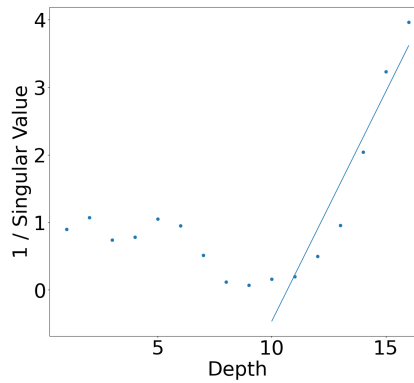


Figure 46: Convolutional ResNet34 (Type 2 model) trained on ImageNette.

91 **4 Broader Impacts**

92 This work presents foundational research and we do not foresee any potential negative societal
93 impacts.

94 **References**

95 Hector Miranda and Robert C Thompson. A trace inequality with a subtracted term. *Linear algebra*
96 *and its applications*, 185:165–172, 1993.

97 Leon Mirsky. A trace inequality of john von neumann. *Monatshefte für mathematik*, 79(4):303–306,
98 1975.