
Emergence in non-neural models: grokking modular arithmetic via average gradient outer product

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Neural networks trained to solve modular arithmetic tasks exhibit *grokking*, the
2 phenomenon where the test accuracy improves only long after the model achieves
3 100% training accuracy in the training process. It is often taken as an example
4 of “emergence”, where model ability manifests sharply through a phase transi-
5 tion. In this work, we show that the phenomenon of grokking is not specific
6 to neural networks nor to gradient descent-based optimization. Specifically, we
7 show that grokking occurs when learning modular arithmetic with Recursive Fea-
8 ture Machines (RFM), an iterative algorithm that uses the Average Gradient Outer
9 Product (AGOP) to enable task-specific feature learning with kernel machines.
10 We show that RFM and, furthermore, neural networks that solve modular arith-
11 metic learn block-circulant features transformations which implement the previ-
12 ously proposed Fourier multiplication algorithm.

13 1 Introduction

14 In recent years the idea of “emergence” has become an important narrative in machine learning.
15 While there is no broad agreement on the definition (Rogers & Luccioni, 2023), it is often argued
16 that “skills” emerge during the training process once certain data size, compute, or model size thresh-
17 olds are achieved (Wei et al., 2022; Arora & Goyal, 2023). Furthermore, these skills are believed to
18 appear rapidly, exhibiting sharp and seemingly unpredictable improvements in performance at these
19 thresholds. One of the simplest and most striking examples supporting this idea is “grokking” mod-
20 ular arithmetic (Power et al., 2022; Nanda et al., 2023). A neural network trained to predict modular
21 addition or another arithmetic operation on a fixed data set rapidly transitions from near-zero to per-
22 fect (100%) test accuracy at a certain point in the optimization process. Surprisingly, this transition
23 point occurs long after perfect *training accuracy* is achieved. Not only is this contradictory to the
24 traditional wisdom regarding overfitting but, as we will show, some aspects of grokking do not fit
25 neatly with our modern understanding of “benign overfitting” Bartlett et al. (2021); Belkin (2021).

26 Despite a large amount of recent work on emergence and, specifically, grokking, (see, e.g., (Power
27 et al., 2022; Liu et al., 2023; Nanda et al., 2023; Thilak et al., 2022; Furuta et al., 2024; Miller et al.,
28 2024)), the nature or even existence of the emergent phenomena remains contested. For example,
29 the recent paper Schaeffer et al. (2023) suggests that the rapid emergence of skills may be a “mirage”
30 due to the mismatch between the discontinuous metrics used for evaluation, such as accuracy, and
31 the continuous loss used in training. The authors argue that, in contrast to accuracy, the test (or
32 validation) loss or some other suitably chosen metric may decrease gradually throughout training
33 and thus provide a useful measure of progress. Another possible progress measure is the training
34 loss. As SGD-type optimization algorithms generally result in a gradual decrease of the training
35 loss, one may posit that skills appear once the training loss falls below a certain threshold in the

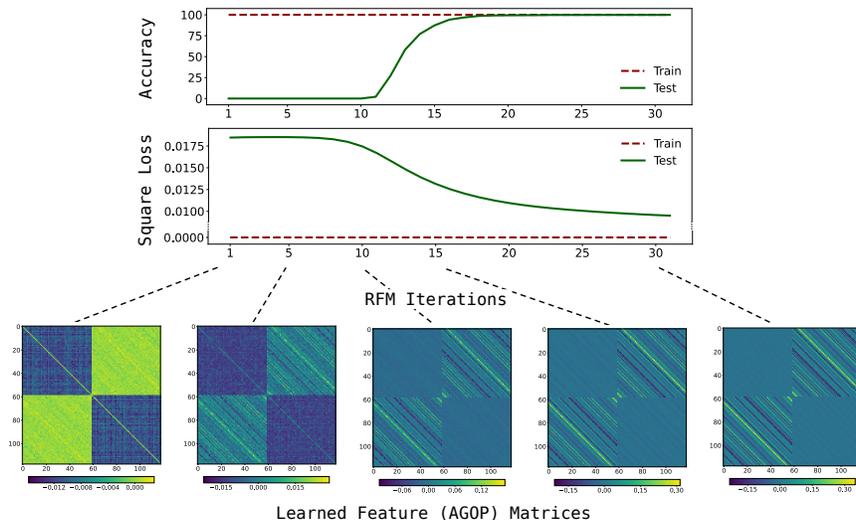


Figure 1: Recursive Feature Machines grok the modular arithmetic task $f^*(x, y) = (x + y) \bmod 59$.

36 optimization process. Indeed, such a conjecture is in the spirit of classical generalization theory,
 37 which considers the training loss to be a useful proxy for the test performance Mohri et al. (2018).

38 In this work, we show that sharp emergence in modular arithmetic arises entirely from feature learn-
 39 ing, independently of other aspects of modeling and training, and is not predicted by the standard
 40 measures of progress. We then clarify the nature of feature learning leading to the emergence of
 41 skills in modular arithmetic. We discuss these contributions in further detail below.

42 **Summary of the contributions.** We demonstrate empirically that grokking modular arithmetic:
 43 (1) is not specific to neural networks; (2) is not tied to gradient-based optimization methods; (3) is
 44 not predicted by training or test loss¹, let alone accuracy.

45 Specifically, we show grokking for Recursive Feature Machines (RFM) (Radhakrishnan et al.,
 46 2024a), an algorithm that iteratively uses the Average Gradient Outer Product (AGOP) to enable
 47 task-specific feature learning in general machine learning models. In this work, we use RFM to en-
 48 able feature learning in kernel machines, which are a class of predictors with no native mechanism
 49 for feature learning. In this setting, RFM iterates between three steps: (i) training a kernel machine,
 50 f , to fit training data; (ii) computing the AGOP matrix of f , M , over the training data to extract
 51 task-relevant features; and (iii) transforming input data, x , using the learned features via the map
 52 $x \rightarrow M^{s/2}x$ for a matrix power $s > 0$ (see Section 2 for details).

53 In Fig. 1 we give a representative example of RFM grokking modular addition, despite not using any
 54 gradient-based optimization methods and achieving perfect (numerically zero) training loss at every
 55 iteration. We see that during the first few iterations both the test loss and and test accuracy remain at
 56 the constant (random) level. Around iteration 10 the test loss starts improving and, a few iterations
 57 later, test accuracy quickly transitions to 100%. We also observe that even early in the iteration,
 58 structure emerges in AGOP feature matrices (see Fig. 1). The gradual appearance of structure in
 59 these feature matrices is striking given that the training loss is identically zero at every iteration and
 60 that the test loss does not significantly change until iteration 8. The striped patterns observed in
 61 feature matrices correspond to matrices whose sub-blocks are circulant with entries that are constant
 62 along the “long” diagonals which wrap around the matrix.² Such *circulant feature matrices* are key
 63 to learning modular arithmetic. In Section 3 we demonstrate that standard kernel machines using
 64 *random* circulant features easily learn modular operations. As these random circulant matrices are
 65 generic, we argue that no additional structure is required to solve modular arithmetic.

¹We note that for neural networks trained by SGD, emergence cannot be decoupled from training loss, as non-zero loss is required for training to occur at all.

²Feature sub-matrices may also be constant on anti-diagonals. We also refer to these matrices as circulant.

66 To demonstrate that the feature matrices evolve toward this structure (including for multiplication
67 and division under an appropriate re-ordering of the input coordinates), we introduce two “hidden
68 progress measures” (Barak et al., 2022): (1) *Circulant deviation*, which measures constancy of
69 the diagonals of a matrix, and (2) *AGOP alignment*, which measures similarity between the feature
70 matrix at iteration t and the AGOP of a fully trained model. We will see that both of these measures
71 show gradual (initially nearly linear) progress toward a model that generalizes.

72 We further argue that emergence in fully connected neural networks trained on modular arithmetic
73 identified in prior work (Gromov, 2023; Liu et al., 2022) is analogous to that for RFM and is exhib-
74 ited through the AGOP (see Section 4). By visualizing covariances of network weights, we observe
75 that these models also learn block-circulant features to grok modular arithmetic. We demonstrate
76 that these features are highly correlated with the AGOP of neural networks, corroborating prior ob-
77 servations from Radhakrishnan et al. (2024a). Furthermore, paralleling our observations for RFM,
78 our progress measures indicate gradual progress toward a generalizing solution during neural net-
79 work training. Finally we demonstrate that training neural networks on data transformed by random
80 block-circulant matrices dramatically decreases training time needed to learn modular arithmetic.

81 Why are these learned block circulant features effective for modular arithmetic? We provide support-
82 ing theoretical evidence that circulant features result in kernel machines implementing the Fourier
83 Multiplication Algorithm (FMA) for modular arithmetic (see Section 5). For the case of neural net-
84 works, several prior works have argued empirically and theoretically that neural networks learn to
85 implement the FMA to solve modular arithmetic (Nanda et al., 2023; Varma et al., 2023; Morwani
86 et al., 2024). While kernel RFM and neural networks utilize different classes of predictive models,
87 our results suggest that they discover similar algorithms for implementing modular arithmetic.

88 By decoupling feature learning from predictor training, our results provide evidence for emergent
89 properties of machine learning models arising purely as a consequence of their ability to learn fea-
90 tures. We hope our work will help isolate the underlying mechanisms of emergence and shed light
91 on the key practical concern of how, when, and why these seemingly unpredictable transitions occur.

92 **Paper outline.** Section 2 reviews preliminary concepts. In Section 3, we demonstrate emergence
93 with RFM and show AGOP features consist of circulant blocks. Section 4, shows that neural network
94 features are circulant and are captured by the AGOP. In Section 5, we prove that kernel machines
95 learn the FMA with circulant features. We provide a discussion and conclude in Section 6.

96 2 Preliminaries

97 **Learning modular arithmetic.** Let $\mathbb{Z}_p = \mathbb{Z}/p\mathbb{Z}$ denote the field of integers modulo a prime p and
98 let $\mathbb{Z}_p^* = \mathbb{Z}_p \setminus \{0\}$. We learn modular functions $f^*(a, b) = g(a, b) \bmod p$ where $f^* : \mathbb{Z}_p \times \mathbb{Z}_p \rightarrow \mathbb{Z}_p$,
99 $a, b \in \mathbb{Z}_p$, and $g : \mathbb{Z} \times \mathbb{Z} \rightarrow \mathbb{Z}$ is an arithmetic operation on a and b , e.g. $g(a, b) = a + b$. Note that
100 there are p^2 discrete input pairs (a, b) for all modular operations except for $f^*(a, b) = (a \div b) \bmod p$,
101 which has $p(p - 1)$ inputs as the denominator cannot be 0.

102 To train models on modular arithmetic tasks, we construct input-label pairs by one-hot encoding the
103 input and label integers. Specifically, for every pair $a, b \in \mathbb{Z}_p$, we write the input as $e_a \oplus e_b \in \mathbb{R}^{2p}$
104 and the output as $e_{f^*(a,b)} \in \mathbb{R}^p$, where $e_i \in \mathbb{R}^p$ is the i -th standard basis vector in p dimensions and
105 \oplus is concatenation. The training dataset consists of a random subset of $n = r \times N$ input/label pairs,
106 where r is the *training fraction* and $N = p^2$ or $p(p - 1)$ is the number of possible discrete inputs.

107 **Circulant matrices.** The features that RFMs and neural networks learn in order to solve modular
108 arithmetic contain blocks of *circulant matrices*, which are defined as follows. Let $\sigma : \mathbb{R}^p \rightarrow \mathbb{R}^p$
109 be the cyclic permutation which acts on a vector $u \in \mathbb{R}^p$ by shifting its coordinates by one cell
110 to the right: $[\sigma(u)]_j = u_{j-1 \bmod p}$, for $j \in [p]$. We write the ℓ -fold composition of this map
111 $\sigma^\ell(u) \in \mathbb{R}^p$ with entries $[\sigma^\ell(u)]_j = u_{j-\ell \bmod p}$. A circulant matrix $C \in \mathbb{R}^{p \times p}$ is determined by
112 a vector $c = [c_0, \dots, c_{p-1}] \in \mathbb{R}^p$, and has rows (in order from first to last): $c, \sigma(c), \dots, \sigma^{p-1}(c)$.
113 Feature matrices may also have have constant anti-diagonals (so-called Hankel matrices). To ease
114 terminology, we will use the word circulant to refer to both Hankel and circulant matrices.

115 **Average Gradient Outer Product (AGOP).** The AGOP matrix, which will be central to our dis-
116 cussion, is defined as follows.

117 **Definition 2.1 (AGOP).** Given a predictor $f : \mathbb{R}^d \rightarrow \mathbb{R}^c$ with c outputs, $f(x) \equiv$
118 $[f_0(x), \dots, f_{c-1}(x)]$, let $\frac{\partial f(x')}{\partial x} \in \mathbb{R}^{d \times c}$ be the Jacobian (transposed) of f evaluated at some point

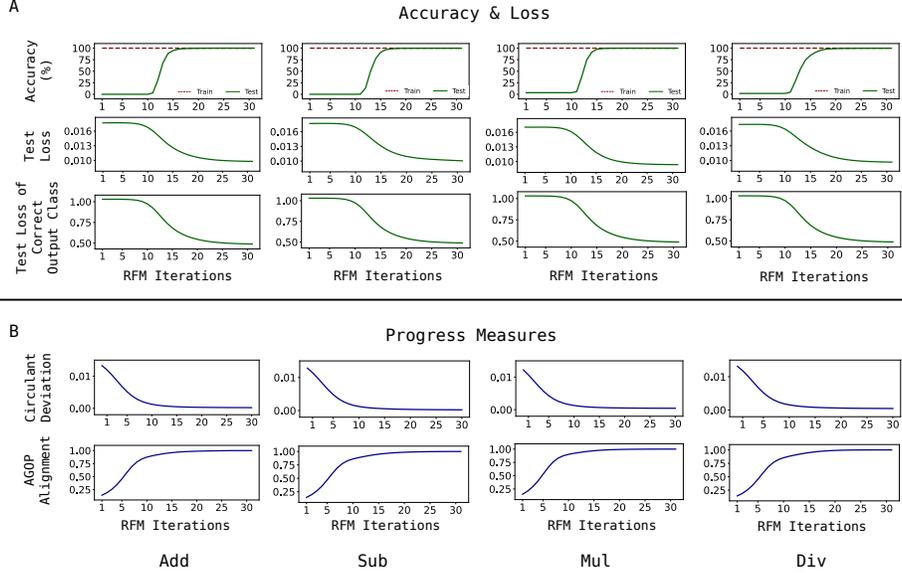


Figure 2: RFM with the quadratic kernel on modular arithmetic with modulus $p = 61$ trained for 30 iterations. (A) Test accuracy, test loss (mean squared error) over all output coordinates, and test loss of the correct class output coordinate do not change in the first 8 iterations and then, sharply transition. (B) Circulant deviation and AGOP alignment show gradual progress towards generalizing solutions despite accuracy and loss metrics not changing in the initial iterations. For multiplication (Mul) and division (Div), circulant deviation is measured with respect to the feature sub-matrices after reordering by the discrete logarithm.

119 $x' \in \mathbb{R}^d$ with entries $[\frac{\partial f(x')}{\partial x}]_{s,\ell} = \frac{\partial f_\ell(x')}{\partial x_s}$. Then, for f trained on a set of data points $\{x^{(j)}\}_{j=1}^n$,
 120 with $x^{(j)} \in \mathbb{R}^d$, the Average Gradient Outer Product (AGOP), G , is defined as,

$$G(f; \{x^{(j)}\}_{j=1}^n) = \frac{1}{n} \sum_{j=1}^n \frac{\partial f(x^{(j)})}{\partial x} \frac{\partial f(x^{(j)})}{\partial x}^\top \in \mathbb{R}^{d \times d}. \quad (1)$$

121 For simplicity, we omit the dependence on the dataset in the notation. Top eigenvectors of AGOP can
 122 be viewed as the “most relevant” input features, those input directions that influence the output of a
 123 general predictor (for example, a kernel machines or a neural network) the most. As a consequence,
 124 the AGOP can be viewed as a task-specific transformation that can be used to amplify relevant
 125 features and improve sample efficiency of machine learning models.

126 Indeed, a line of prior works (Yuan et al., 2023; Trivedi et al., 2014; Hristache et al., 2001) have used
 127 the AGOP to improve the sample efficiency of predictors trained on multi-index models, a class of
 128 predictive tasks in which the target function depends on a low-rank subspace of the data. Though
 129 the study of AGOP has been motivated by these multi-index examples, we will see that the AGOP
 130 can be used to recover useful features for modular arithmetic that are, in fact, not low-rank.

131 **AGOP and feature learning in neural networks.** Radhakrishnan et al. (2024a) posited that AGOP
 132 was a mechanism through which neural networks learn features. In particular, the authors introduce
 133 the *Neural Feature Ansatz (NFA)* stating that for any layer ℓ of a trained neural network with weights
 134 W_ℓ , the *Neural Feature Matrix (NFM)*, $W_\ell^\top W_\ell$, are highly correlated to the AGOP of the model
 135 computed with respect to the input of layer ℓ . The NFA suggests that neural networks learn features
 136 at each layer by utilizing the AGOP. For more details on the NFA, see Appendix C.

137 **Recursive Feature Machine (RFM).** Importantly, AGOP can be computed for any differentiable
 138 predictor, including those such as kernel machines that have no native feature learning mechanism.
 139 As such, the authors of Radhakrishnan et al. (2024a) developed an algorithm known as RFM, which
 140 iteratively uses the AGOP to extract features. Below, we present the RFM algorithm used in conjunc-
 141 tion with kernel machines. Suppose we are given data samples $(X, y) \in \mathbb{R}^{n \times d} \times \mathbb{R}^n$ where X con-
 142 tains n samples denoted $\{x^{(j)}\}_{j=1}^n$. Given an initial symmetric positive-definite matrix $M_0 \in \mathbb{R}^{d \times d}$,

143 and Mahalanobis kernel $k(\cdot, \cdot; M) : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$, RFM iterates the following steps for $t \in [T]$:

$$\text{Step 1 (Predictor training): } f^{(t)}(x) = k(x, X; M_t)\alpha \text{ with } \alpha = k(X, X; M_t)^{-1}y; \quad (2)$$

$$\text{Step 2 (AGOP update): } M_{t+1} = [G(f^{(t)})]^s; \quad (3)$$

144 where $s > 0$ is a matrix power and $k(X, X; M) \in \mathbb{R}^{n \times n}$ denotes the matrix with entries
 145 $[k(X, X; M)]_{j_1 j_2} = k(x^{(j_1)}, x^{(j_2)}; M)$ for $j_1, j_2 \in [n]$. In this work, we select $s = \frac{1}{2}$ for
 146 all experiments (see Algorithm 1 for complete pseudocode). We use the following two Maha-
 147 lanobis kernels: (1) the quadratic kernel, $k(x, x'; M) = (x^\top M x')^2$; and (2) the Gaussian kernel
 148 $k(x, x'; M) = \exp(-\|x - x'\|_M^2/L)$, where for $z \in \mathbb{R}^d$, $\|z\|_M^2 = z^\top M z$, and L is the bandwidth.

149 3 Emergence with Recursive Feature Machines

150 We now show that RFM exhibits sharp transitions in performance on modular arithmetic tasks (ad-
 151 dition, subtraction, multiplication, and division) due to the emergence of block-circulant features.

152 We will use a modulus of $p = 61$ and train RFM with
 153 quadratic and Gaussian kernel machines (experimental
 154 details are provided in Appendix D). As we solve kernel
 155 ridgeless regression exactly, all iterations of RFM
 156 result in zero training loss and 100% training accuracy.
 157 The top two rows of Fig. 2A show that the first several
 158 iterations of RFM result in near-zero test accuracy and
 159 approximately constant, large test loss. Despite these
 160 standard progress measures initially not changing, con-
 161 tinuing to iterate RFM leads to a dramatic, sharp in-
 162 crease to 100% test accuracy and a corresponding de-
 163 crease in the test loss later in the iteration process.

164 **Sharp transition in loss of correct output coordi-
 165 nate.** It is important to note that our total loss function
 166 is the square loss averaged over $p = 61$ classes.
 167 It is thus plausible that, due to averaging, the near-
 168 constancy of the total square loss over the first few it-
 169 erations conceals steady improvements in the predictions
 170 of the correct class. However, in Fig. 2A (third row) we
 171 show that the test loss for the output coordinate (logit)
 172 of the correct class closely tracks the total test loss.

173 **Emergence of block-circulant features in RFM.** To
 174 understand RFM generalization, we visualize the $2p \times$
 175 $2p$ feature matrix given by the square root of the AGOP
 176 from the final iteration of RFM. We first visualize the
 177 feature matrices for RFM trained on modular addition/-
 178 subtraction in Fig. 3A. Their visually-evident striped
 179 structure suggests a more precise characterization:

180 **Observation 1** (Block-circulant features). *Feature matrix $M^* \in \mathbb{R}^{2p \times 2p}$ at the final iteration of
 181 RFM on modular addition/subtraction is of the form*

$$M^* = \begin{pmatrix} A & C^\top \\ C & A \end{pmatrix}, \quad (4)$$

182 where $A, C \in \mathbb{R}^{p \times p}$, C is an asymmetric circulant matrix, $A = c_1 I + c_2 \mathbf{1}\mathbf{1}^\top$ for scalars c_1, c_2 .

183 Similarly to addition and subtraction, RFM successfully learns multiplication and division. Yet,
 184 in contrast to addition and subtraction, the structure of feature matrices for these tasks, shown in
 185 Fig. 3B, is not at all obvious. Nevertheless, re-ordering the rows and columns of the feature matrices
 186 for these tasks brings out their hidden circulant structure of the form stated in Eq. (4). We show the
 187 effect of re-ordering in Fig. 3C (see also Appendix Fig. 1 for the evolution of re-ordered and original
 188 features during training).

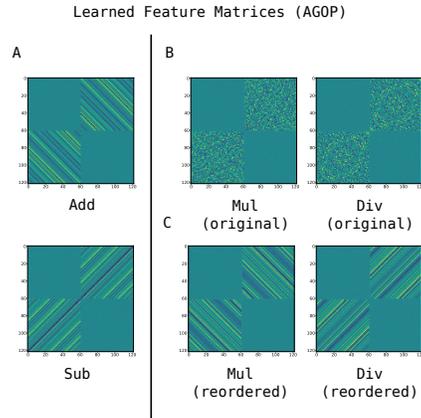


Figure 3: RFM with the quadratic kernel for modular arithmetic with $p = 61$. (A) The square root of the kernel AGOPs for addition (Add), subtraction (Sub) visualized without their diagonals to emphasize the off-diagonal blocks. (B) Square root of the kernel AGOP for multiplication (Mul), division (Div). (C) For Mul and Div, rows and columns of each sub-matrix is re-ordered by the discrete log. base 2.

189 We briefly discuss the reordering procedure below and provide further details in Appendix E. To
 190 reorder, we use the fact of group theory that the multiplicative group \mathbb{Z}_p^* is a cyclic group of order
 191 $p - 1$ (e.g., Koblitz (1994)). By definition of the cyclic group, there exists at least one element
 192 $g \in \mathbb{Z}_p^*$, known as a *generator*, such that $\mathbb{Z}_p^* = \{g^i; i \in \{1, \dots, p - 1\}\}$. As we will see, re-
 193 ordering the rows and columns of the AGOP by powers of a generator reveals circulant structure.
 194 For modular multiplication/division, the map taking g^i to i is known as the *discrete logarithm* base
 195 g (Koblitz, 1994, Ch.3). It is natural to expect block-circulant feature matrices to arise in modular
 196 multiplication/division after reordering by the discrete log as the discrete log converts modular
 197 multiplication/division into modular addition/subtraction. We note the recent work Doshi et al. (2024)
 198 also used the discrete log to reorder coordinates in the context of constructing a solution for solving
 199 modular multiplication with neural networks.

200 **Progress measures.** We propose and examine two measures of feature learning, *circulant deviation*
 201 and *AGOP alignment*.

202 *Circulant deviation.* As the final feature matrices contain circulant sub-blocks, a natural progress
 203 measure for learning modular arithmetic with RFM is how far AGOP feature matrices are from
 204 a block-circulant matrix. For a feature matrix M , let A denote the bottom-left sub-block of M .
 205 We define circulant deviation as the total variance of the (wrapped) diagonals of A normalized by
 206 the norm $\|A\|_F^2$. In particular, let $\mathcal{S} \in \mathbb{R}^{p \times p} \rightarrow \mathbb{R}^{p \times p}$ denote the shift operator, which shifts
 207 the ℓ -th row of the matrix by ℓ positions to the right. Also let $\text{Var}(\mathbf{v}) = \sum_{j=0}^{p-1} (v_j - \mathbb{E}\mathbf{v})^2$ be
 208 the variance of a vector \mathbf{v} . If $A[j]$ denotes the j -th column of A , we define circulant deviation \mathcal{D}
 209 as: $\mathcal{D}(A) = \frac{1}{\|A\|_F^2} \sum_{j=0}^{p-1} \text{Var}(\mathcal{S}(A)[j])$. As circulant matrices are constant along their (wrapped)
 210 diagonals, they have a circulant deviation of 0.

211 We see in Fig. 2B (top row) that circulant deviation exhibits gradual improvement through the course
 212 of training with RFM. We find that for the first 10 iterations, while the training loss is numerically
 213 zero and the test loss does not improve, circulant deviation exhibits gradual, nearly linear, improve-
 214 ment. The improvements in circulant deviation reflect visual improvements in features, as was also
 215 shown in Fig. 1. These curves also provide further support for Observation 1, as the circulant devi-
 216 ation is close to 0 at the end of training.

217 Circulant deviation depends crucially on the observation that for modular arithmetic the feature
 218 matrices contained circulant blocks. For more general tasks, we may not be able to identify such
 219 structure. Thus, we propose a second, more general progress measure, AGOP alignment.

220 *AGOP alignment.* Given two matrices $A, B \in \mathbb{R}^{d \times d}$, let $\rho(A, B)$ denote the standard cosine simi-
 221 larity between these two matrices when vectorized. Specifically, let $\tilde{A}, \tilde{B} \in \mathbb{R}^{d^2}$ denote the vector-
 222 ization of A and B respectively, then $\rho(A, B) = \frac{\langle \tilde{A}, \tilde{B} \rangle}{\|\tilde{A}\| \|\tilde{B}\|}$.

223 If M_t denotes the AGOP at iteration t of RFM (or epoch t of a neural network) and M^* denotes the
 224 final AGOP of the trained RFM (or neural network), then AGOP alignment at iteration t is given by
 225 $\rho(M_t, M^*)$. The same measure of alignment was used in Zhu et al. (2024), except their alignment
 226 was computed with respect to the AGOP of the ground truth model. Note that as modular operations
 227 are discrete, in our setting there is no unique ground truth model for which AGOP can be computed.

228 Like circulant deviation, AGOP alignment exhibits gradual improvement in the regime that test loss
 229 is constant and large (see Fig. 2B, bottom row). Moreover, AGOP alignment is a more general
 230 progress measure since it does not require assumptions on the structure of the AGOP. For instance,
 231 AGOP alignment can be measured without reordering for modular multiplication/division. While
 232 AGOP alignment does not require a specific form of the final features, it is still an *a posteriori*
 233 measurement of progress as it requires access to the features of a fully trained model.

234 **Random circulant features allow standard kernels to generalize.** We conclude this section by
 235 providing further evidence that the form of feature matrices given in Observation 1 is key to enabling
 236 generalization in kernel machines trained to solve modular arithmetic tasks. We now show that a
 237 transformation with a *generic* block-circulant matrix enables kernel machines to learn modular
 238 arithmetic. We generate a random circulant matrix C by first sampling entries of the first column
 239 i.i.d. from the uniform distribution on $[0, 1] \subset \mathbb{R}$ and then shifting the column to generate the
 240 remaining columns of C . We construct M^* in Observation 1 with $c_1 = 1, c_2 = -1/p$. For modular
 241 addition, we transform the input data by mapping $x_{ab} = e_a \oplus e_b$ to $\tilde{x}_{ab} = (M^*)^{\frac{1}{4}} x_{ab}$, and then

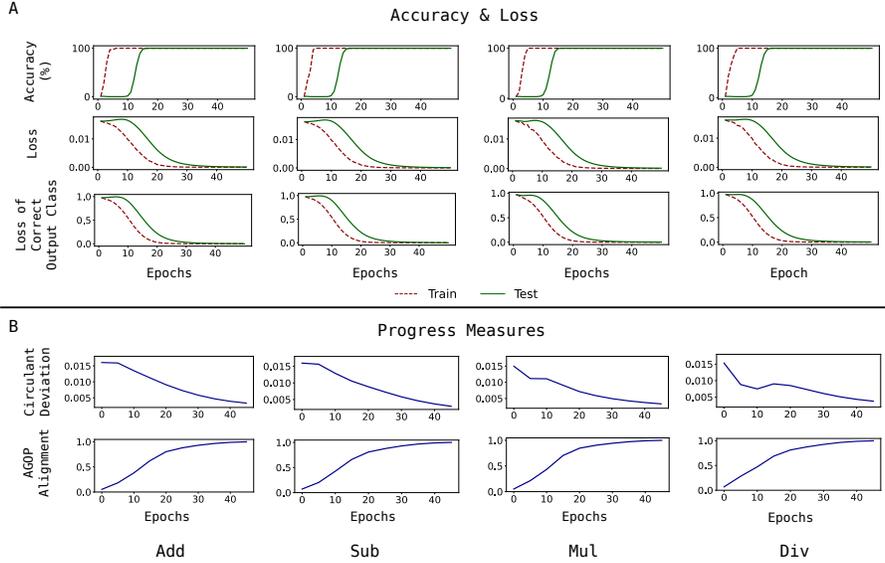


Figure 5: One hidden layer fully-connected networks with quadratic activations trained on modular arithmetic with $p = 61$ trained for 50 epochs with the square loss. (A) Test accuracy, test loss over all outputs, and test loss of the correct class output do not change in the initial iterations. (B) Progress measures for circulant deviation and AGOP alignment. Circulant deviation for Mul and Div are computed after reordering by the discrete logarithm base 2.

242 train on the new data pairs $(\tilde{x}_{ab}, e_{a+b \bmod p})$ for a subset of all possible pairs $(a, b) \in \mathbb{Z}_p^2$. Note that
 243 transforming data with $(M^*)^{\frac{1}{2}}$ is akin to using $s = 1/2$ in RFM.

244 We do the same for modular multiplication after reorder-
 245 ing the random circulant by the discrete logarithm as de-
 246 scribed above. The experiments in Fig. 4 show that stan-
 247 dard kernel machines trained on feature matrices with
 248 random circulant blocks outperform RFM that learns such
 249 features through AGOP. We also find that directly enforc-
 250 ing circulant blocks in the sub-matrices of M_t through-
 251 out RFM iterations accelerates grokking and improves test
 252 loss (see Appendix F, Appendix Fig. 2). These exper-
 253 iments provide direct evidence that the structure in Ob-
 254 servation 1 is key for generalization on modular arith-
 255 metic and, furthermore, *no additional structure* beyond
 256 a generic circulant is required.

257 4 Emergence 258 in neural networks through AGOP

259 We now show that grokking in two-layer neural networks
 260 relies on the same principles as grokking by RFM. Specif-
 261 ically we demonstrate that (1) block-circulant features are
 262 key to neural networks grokking modular arithmetic; and
 263 (2) our measures (circulant deviation and AGOP align-
 264 ment) indicate gradual progress towards generalization,
 265 while standard measures of generalization exhibit sharp
 266 transitions. All experimental details are provided in Ap-
 267 pendix D.

268 **Grokking with neural networks.** We first reproduce grokking with modular arithmetic using fully-
 269 connected networks as identified in prior works (Fig. 5A) (Gromov, 2023). In particular, we train

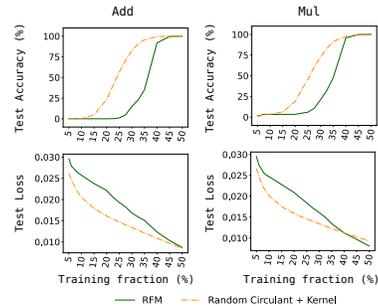


Figure 4: Random circulant features generalize with standard kernels for modular arithmetic. RFM with the Gaussian kernel on addition (Add) and multiplication (Mul) for modulus $p = 61$ is compared to a base Gaussian kernel machine trained on random circulant features (for Mul, the sub-blocks are circulant after re-ordering by the discrete logarithm base 2).

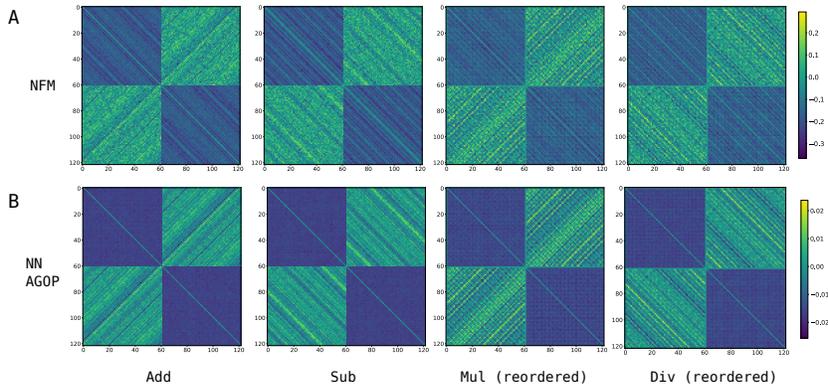


Figure 6: Feature matrices from one hidden layer neural networks with quadratic activations trained on addition, subtraction, multiplication, and division modulo 61. The Pearson correlations between the NFM and square root of the AGOP for each task are 0.955 (Add), 0.942 (Sub), 0.924 (Mul), 0.929 (Div). Mul and Div are shown after reordering by the discrete logarithm base 2.

270 one hidden layer fully connected networks $f : \mathbb{R}^{2p} \rightarrow \mathbb{R}^p$ of the form $f(x) = W_2\phi(W_1x)$ with
 271 quadratic activation $\phi(z) = z^2$ on modulus $p = 61$ data with a training fraction 50%.

272 Consistent with prior work (Gromov, 2023) and analogously to RFMs, neural networks exhibit an
 273 initial training period where the train accuracy reaches 100%, while test accuracy is at 0% and test
 274 loss does not improve (see Fig. 5A). After this point, we see that the accuracy rapidly improves to
 275 achieve perfect generalization. We further verify that the sharp transition in test loss is not an artifact
 276 of averaging the loss over all output coordinates. In the third row of Fig. 5A we show that the test
 277 loss of the individual correct output coordinate closely tracks the total loss.

278 **Emergence of block-circulant features in neural networks.** To understand the features learned
 279 by neural networks we visualize the first layer Neural Feature Matrix, defined as follows.

280 **Definition 4.1.** Given a fully connected network $f(x) = a^\top \phi(W_1x)$, the first layer Neural Feature
 281 Matrix (NFM) is the matrix $W_1^\top W_1 \in \mathbb{R}^{2p \times 2p}$.

282 The NFM is the un-centered covariance of network weights and has been used in prior work in order
 283 to understand the features learned by various neural network architectures at any layer (Radhakrishnan
 284 et al., 2024a; Trockman et al., 2022). Fig. 6A displays the NFM for one hidden layer neural
 285 networks with quadratic activations trained on modular arithmetic tasks. For addition/subtraction,
 286 we find that the NFM exhibits block circulant structure, akin to the feature matrix for RFM. As
 287 described in Section 3 and Appendix E, we reorder the NFM for networks trained on multiplication/
 288 division with respect to a generator for \mathbb{Z}_p^* in order to observe block-circulant structure (see
 289 Appendix Fig. 4A for a comparison of multiplication/division NFMs before and after reordering).
 290 The block-circulant structure in both the NFM and the feature matrix of RFM suggests that the two
 291 models are learning similar sets of features.

292 The work Radhakrishnan et al. (2024a) posited that AGOP is the mechanism through which neural
 293 networks learn features. The authors stated their claim in the form of the Neural Feature Ansatz
 294 (NFA), which states that NFMs are proportional to a matrix power of AGOP through training (see
 295 Eq. (5) for a restatement of the NFA). As such, we additionally compute the square root of the AGOP
 296 to examine the features learned by neural networks trained on modular arithmetic tasks. We visualize
 297 the square root of the AGOPs of these trained models in Fig. 6B and also find that the square root
 298 of the AGOP and the NFM are highly correlated (greater than 0.92), where Pearson correlation is
 299 equal to cosine similarity after centering the inputs to be mean 0. Moreover, we find that the square
 300 root of AGOP of neural networks again exhibits the same structure as stated in Observation 1 (see
 301 Appendix Fig. 4B for a comparison of multiplication/division AGOPs before and after reordering).

302 **Random circulant maps improve generalization of neural networks.** To further establish the
 303 importance and generality of block-circulant features, we demonstrate that training networks on
 304 inputs transformed with a random block-circulant matrix greatly accelerates learning. In Fig. 7, we
 305 compare the performance of neural networks trained on one-hot encoded modulo p integers and the

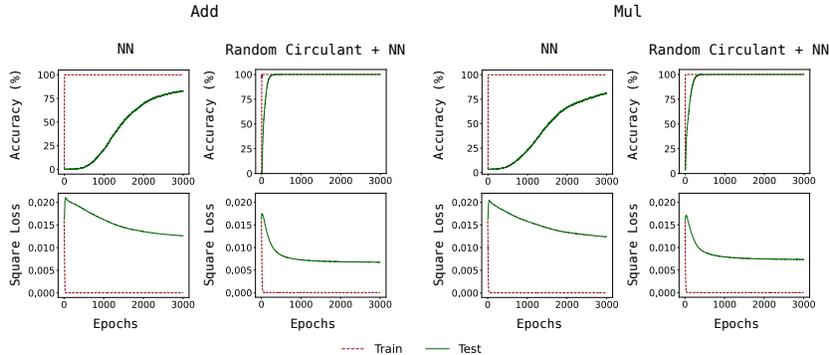


Figure 7: Random circulant features speed up generalization in neural networks for modular arithmetic tasks. We compare one hidden layer MLPs with quadratic activations trained on modular addition and multiplication for $p = 61$ using standard one-hot encodings or those transformed by random circulant matrices (re-ordered by the discrete logarithm for multiplication).

306 same integers transformed with a random block-circulant matrix. At a training fraction of 17.5%,
 307 we find that networks trained on transformed integers achieved 100% test accuracy within several
 308 hundred epochs and exhibit little delayed generalization while networks trained on non-transformed
 309 integers do not achieve 100% test accuracy even within 3000 epochs.

310 **Progress measures.** Given that the square root of the AGOP of neural networks exhibits block-
 311 circulant structure, we can use circulant deviation and AGOP alignment to measure gradual progress
 312 of neural networks toward a generalizing solution. As before, we measure circulant deviation in
 313 the case of multiplication/division after reordering the feature submatrix by a generator of \mathbb{Z}_p^* . In
 314 Fig. 5B, we see that our measures indicate gradual progress in contrast to sharp transitions in the
 315 standard measures of progress shown in Fig. 5A. There is a period of 5-10 epochs where circulant
 316 deviation and AGOP alignment improve while test loss and test accuracy do not. As was the case of
 317 RFM, these metrics reveal gradual progress of neural networks toward generalizing solutions.

318 5 Fourier multiplication algorithm from circulant features

319 We have seen so far that features containing circulant sub-blocks enable generalization for RFMs
 320 and neural networks across modular arithmetic tasks. We now provide theoretical support that shows
 321 how kernel machines equipped with such circulant features learn generalizing solutions. In particu-
 322 lar, we show that there exist block-circulant feature matrices, as in Observation 1, such that kernel
 323 machines equipped with these features and trained on all available data for a given modulus p solve
 324 modular arithmetic through the *Fourier Multiplication Algorithm* (FMA). Notably, the FMA has
 325 been argued both empirically and theoretically in prior works to be the solution found by neural
 326 networks to solve modular arithmetic (Nanda et al., 2023; Zhong et al., 2024).

327 The FMA is a specific solution for implementing modular arithmetic that first represents the data by
 328 its Discrete Fourier Transform (DFT). Intuitively, transforming the data with circulant matrices ex-
 329 tracts the DFT of the one-hot encoded vectors following the well-known fact that circulant matrices
 330 can be diagonalized using the matrix that encodes the DFT (Gray et al., 2006). We state our result
 331 informally here (for more details on the FMA, the precise theorem, and its proof, see Appendix G).

332 **Theorem 5.1** (Circulant features give the FMA). *Training on all of the discrete data for any mod-
 333 ular operation, for each output class $\ell \in \{0, \dots, p - 1\}$, suppose we train a separate quadratic
 334 kernel predictor and particular block-circulant feature matrices M_ℓ (having the structure in Obser-
 335 vation 1). Then, the concatenated predictor given by kernel ridgeless regression on each output is
 336 equivalent to the Fourier Multiplication Algorithm for that modular operation.*

337 Notably, the FMA is defined over all of \mathbb{R}^{2p} , not just on one-hot encoded inputs. Thus, not only do
 338 neural networks and RFM learn similar features, we have established a setting where kernel meth-
 339 ods equipped with block-circulant feature matrices learn the same out-of-domain solution as neural
 340 networks for these tasks. This result is interesting, in part, as the only constraint for generalization
 341 on these tasks is to obtain perfect accuracy on inputs that are standard basis vectors.

343 Most classical analyses of generalization relied on the training loss serving as a proxy for the test loss
 344 and thus a useful measure of generalization. Empirical results of deep learning have upended this
 345 long-standing belief. In many settings, predictors that fit the data exactly can still generalize, thus
 346 invalidating training loss as a predictor of test performance. This has led to the recent developments
 347 in understanding benign overfitting, in neural networks as well as in classical kernel and linear mod-
 348 els Belkin (2021); Bartlett et al. (2021). Since the training loss may not predict generalization, the
 349 common suggestion has been to use the validation loss computed on a separate *validation dataset*.
 350 Emergent phenomena, such as grokking, show that we cannot rely even on validation performance
 351 at intermediate training steps to predict generalization at the end of training. Indeed, validation loss
 352 at a certain iteration may not be indicative of the validation loss itself only a few iterations later.
 353 Further, contrary to Schaeffer et al. (2023), we show these phase transitions in performance are not
 354 generally “a mirage” since, as we observe in this work, they are not always predicted by *a priori*
 355 measures of performance, continuous or discontinuous. Instead, emergence is fully determined by
 356 feature learning, which is difficult to observe without having access to a fully trained model. Indeed,
 357 the progress measures discussed in this work, as well as those suggested in, e.g., Barak et al. (2022);
 358 Nanda et al. (2023); Doshi et al. (2024) can be termed *a posteriori* progress indicators. They all
 359 require either understanding of the algorithm implemented by a generalizing trained model (such as
 360 our circulant deviation, the Fourier gap considered in Barak et al. (2022), or the Inverse Participation
 361 Ratio in Doshi et al. (2024)) or access to such a model (e.g. AGOP alignment).

362 Consider generalizing features for modular multiplication shown in Fig. 3. The original features
 363 shown in panel B of this figure do not have an easily identifiable pattern. In contrast, re-ordered
 364 features in panel C are clearly striped, containing block-circulants. As discussed in Section 3, re-
 365 ordering of features requires understanding that the multiplicative group \mathbb{Z}_p^* is cyclic of order $p - 1$.
 366 While a well-known result, it is far from obvious *a priori*. It is thus plausible that in other settings
 367 hidden feature structures may be hard to identify due to a lack of mathematical insight.

368 **Why is learning modular arithmetic surprising?** The task of learning modular operations is
 369 different from many other statistical machine learning tasks. In continuous ML settings, we typically
 370 posit that the “ground truth” target function is smooth in an appropriate sense. Hence any general
 371 purpose algorithm capable of learning smooth functions (such as, for example, k -nearest neighbors)
 372 should be able to learn the target function given enough data. Primary differences between learning
 373 algorithms are thus in sample and computational efficiency. In contrast, it is unclear what principle
 374 leads to learning modular arithmetic from partial observations. There are many ways to fill in the
 375 missing data and we do not know a simple inductive bias, to guide us toward a solution. Several
 376 recent works argued that margin maximization with respect to certain norms can account for learning
 377 modular arithmetic (Morwani et al., 2024; Lyu et al., 2023; Mohamadi et al., 2024). While the
 378 direction is promising, general underlying principles are not yet clear.

379 **Analyses of grokking.** Recent works (Kumar et al., 2024; Lyu et al., 2023; Mohamadi et al., 2024)
 380 argue that grokking occurs in neural networks through a two phase mechanism that transitions from a
 381 “lazy” regime, with no feature learning, to a “rich” feature learning regime. Our experiments clearly
 382 show that grokking in RFM does not undergo such a transition. For RFM on modular arithmetic
 383 tasks, our progress measures indicate that the features evolve gradually toward the final circulant
 384 matrices, even as test performance initially remains constant (Fig. 2). Grokking in these settings
 385 is entirely due to the gradual feature quality improvement and two-phase grokking does not occur.
 386 Additionally, we have not observed significant evidence of “lazy” to “rich” transition as a mechanism
 387 for grokking in our experiments with neural networks, as most of our measures of feature learning
 388 start improving early on in the training process (improvement in circulant deviation measure is
 389 delayed for addition and subtraction, but not for multiplication and division, while AGOP feature
 390 alignment initially shows near linear improvement for all tasks), see Fig. 5. Our observations for
 391 neural networks are in line with the results in (Doshi et al., 2024; Nanda et al., 2023), where their
 392 proposed progress measures, Inverse Participation Ratio and Gini coefficients of the weights in the
 393 Fourier domain, are shown to increase prior to improvements in test loss and accuracy for modular
 394 arithmetic.

395 **References**

- 396 Navid Ardeshir, Daniel J. Hsu, and Clayton H. Sanford. Intrinsic dimensionality and generalization
397 properties of the r -norm inductive bias. In Gergely Neu and Lorenzo Rosasco (eds.), *Proceedings*
398 *of Thirty Sixth Conference on Learning Theory*, volume 195 of *Proceedings of Machine Learning*
399 *Research*, pp. 3264–3303. PMLR, 12–15 Jul 2023. URL <https://arxiv.org/pdf/2206.05317>.
400
- 401 Sanjeev Arora and Anirudh Goyal. A theory for emergence of complex skills in language models.
402 *arXiv preprint arXiv:2307.15936*, 2023. URL <https://arxiv.org/pdf/2307.15936>.
- 403 Boaz Barak, Benjamin Edelman, Surbhi Goel, Sham Kakade, Eran Malach, and Cyril Zhang. Hid-
404 den progress in deep learning: Sgd learns parities near the computational limit. *Advances in Neu-*
405 *ral Information Processing Systems*, 35:21750–21764, 2022. URL <https://openreview.net/pdf?id=8XWP2ewX-im>.
406
- 407 Peter L Bartlett, Andrea Montanari, and Alexander Rakhlin. Deep learning: a statistical viewpoint.
408 *Acta numerica*, 30:87–201, 2021. URL <https://arxiv.org/pdf/2103.09177>.
- 409 Daniel Beaglehole, Adityanarayanan Radhakrishnan, Parthe Pandit, and Mikhail Belkin. Mech-
410 anism of feature learning in convolutional neural networks. *arXiv preprint arXiv:2309.00570*,
411 2023. URL <https://arxiv.org/pdf/2309.00570>.
- 412 Daniel Beaglehole, Ioannis Mitliagkas, and Atish Agarwala. Feature learning as alignment:
413 a structural property of gradient descent in non-linear neural networks. *arXiv preprint*
414 *arXiv:2402.05271*, 2024a. URL <https://arxiv.org/pdf/2402.05271>.
- 415 Daniel Beaglehole, Peter S ukeny, Marco Mondelli, and Mikhail Belkin. Average gradient outer
416 product as a mechanism for deep neural collapse. *arXiv preprint arXiv:2402.13728*, 2024b. URL
417 <https://arxiv.org/pdf/2402.13728>.
- 418 Mikhail Belkin. Fit without fear: remarkable mathematical phenomena of deep learning through
419 the prism of interpolation. *Acta Numerica*, 30:203–248, 2021. URL <https://arxiv.org/pdf/2105.14368>.
420
- 421 Alexandru Damian, Jason Lee, and Mahdi Soltanolkotabi. Neural networks can learn representations
422 with gradient descent. In *Conference on Learning Theory*, pp. 5413–5452. PMLR, 2022. URL
423 <https://arxiv.org/pdf/2206.15144>.
- 424 Xander Davies, Lauro Langosco, and David Krueger. Unifying grokking and double descent. *ML*
425 *Safety Workshop, 36th Conference on Neural Information Processing Systems (NeurIPS 2022)*,
426 2023. URL <https://arxiv.org/abs/2303.06173>.
- 427 Darshil Doshi, Tianyu He, Aritra Das, and Andrey Gromov. Grokking modular polynomials.
428 *International Conference on Learning Representations (ICLR): BGPT Workshop*, 2024. URL
429 <https://arxiv.org/abs/2406.03495>.
- 430 Hiroki Furuta, Gouki Minegishi, Yusuke Iwasawa, and Yutaka Matsuo. Interpreting grokked
431 transformers in complex modular arithmetic. *arXiv preprint arXiv:2402.16726*, 2024. URL
432 <https://arxiv.org/pdf/2402.16726>.
- 433 Robert M Gray et al. Toeplitz and circulant matrices: A review. *Foundations and Trends® in Com-*
434 *munications and Information Theory*, 2(3):155–239, 2006. URL <https://ee.stanford.edu/~gray/toeplitz.pdf>.
435
- 436 Andrey Gromov. Grokking modular arithmetic. *arXiv preprint arXiv:2301.02679*, 2023. URL
437 <https://arxiv.org/pdf/2301.02679>.

- 438 Suriya Gunasekar, Blake E Woodworth, Srinadh Bhojanapalli, Behnam Neyshabur, and Nati Sre-
439 bro. Implicit regularization in matrix factorization. *Advances in neural information processing*
440 *systems*, 30, 2017.
- 441 Judy Hoffman, Daniel A Roberts, and Sho Yaida. Robust learning with jacobian regularization.
442 *arXiv preprint arXiv:1908.02729*, 5(6):7, 2019.
- 443 Marian Hristache, Anatoli Juditsky, Jorg Polzehl, and Vladimir Spokoiny. Structure adaptive
444 approach for dimension reduction. *Annals of Statistics*, pp. 1537–1566, 2001. URL [https://projecteuclid.org/journals/annals-of-statistics/volume-29/
445 issue-6/Structure-Adaptive-Approach-for-Dimension-Reduction/10.
446 1214/aos/1015345954.full](https://projecteuclid.org/journals/annals-of-statistics/volume-29/issue-6/Structure-Adaptive-Approach-for-Dimension-Reduction/10.1214/aos/1015345954.full).
- 448 Neal Koblitz. *A course in number theory and cryptography*, volume 114. Springer Science &
449 Business Media, 1994.
- 450 Tanishq Kumar, Blake Bordelon, Samuel J. Gershman, and Cengiz Pehlevan. Grokking as the tran-
451 sition from lazy to rich training dynamics. *International Conference on Learning Representations*
452 *(ICLR)*, 2024. URL <https://openreview.net/pdf?id=vt5mnLVIVo>.
- 453 Ziming Liu, Ouail Kitouni, Niklas S Nolte, Eric Michaud, Max Tegmark, and Mike Williams. To-
454 wards understanding grokking: An effective theory of representation learning. *Advances in Neu-
455 ral Information Processing Systems*, 35:34651–34663, 2022.
- 456 Ziming Liu, Eric J. Michaud, and Max Tegmark. Omnigrok: Grokking beyond algorithmic
457 data. *International Conference on Learning Representations (ICLR)*, 2023. URL <https://openreview.net/pdf?id=zDiHoIWa0q1>.
- 459 Kaifeng Lyu, Jikai Jin, Zhiyuan Li, Simon Shaolei Du, Jason D Lee, and Wei Hu. Dichotomy of
460 early and late phase implicit biases can provably induce grokking. In *The Twelfth International
461 Conference on Learning Representations (ICLR)*, 2023. URL [https://openreview.net/
462 forum?id=XsHqr9dEGH](https://openreview.net/forum?id=XsHqr9dEGH).
- 463 Jack Miller, Charles O’Neill, and Thang Bui. Grokking beyond neural networks: An empirical
464 exploration with model complexity. *Transactions on Machine Learning Research (TMLR)*, 2024.
465 URL <https://openreview.net/pdf?id=ux9BrxPC18>.
- 466 Mohamad Amin Mohamadi, Zhiyuan Li, Lei Wu, and Danica J. Sutherland. Why do you
467 grok? a theoretical analysis on grokking modular addition. In *Forty-first International Con-
468 ference on Machine Learning (ICML)*, 2024. URL [https://openreview.net/forum?
469 id=ad5I6No9G1](https://openreview.net/forum?id=ad5I6No9G1).
- 470 Mehryar Mohri, Afshin Rostamizadeh, and Ameet Talwalkar. *Foundations of machine learning*.
471 MIT Press, 2018.
- 472 Ankur Moitra. *Algorithmic aspects of machine learning*. Cambridge University Press, 2018.
- 473 Depen Morwani, Benjamin L. Edelman, Costin-Andrei Oncescu, Rosie Zhao, and Sham Kakade.
474 Feature emergence via margin maximization: case studies in algebraic tasks. *International
475 Conference on Learning Representations (ICLR)*, 2024. URL [https://openreview.net/
476 pdf?id=i9wDX850jR](https://openreview.net/pdf?id=i9wDX850jR).
- 477 Alireza Mousavi-Hosseini, Sejun Park, Manuela Girotti, Ioannis Mitliagkas, and Murat A Erdogdu.
478 Neural networks efficiently learn low-dimensional representations with sgd. *arXiv preprint
479 arXiv:2209.14863*, 2022. URL <https://arxiv.org/pdf/2209.14863>.
- 480 Neel Nanda, Lawrence Chan, Tom Lieberum, Jess Smith, and Jacob Steinhardt. Progress measures
481 for grokking via mechanistic interpretability. *International Conference on Learning Representa-*

- 482 *tions (ICLR)*, 2023. URL <https://openreview.net/pdf?id=9XFSbDPmdW>.
- 483 Suzanna Parkinson, Greg Ongie, and Rebecca Willett. Relu neural networks with linear layers are
484 biased towards single- and multi-index models. *arXiv preprint arXiv:2305.15598*, 2023. URL
485 <https://arxiv.org/pdf/2305.15598>.
- 486 Alethea Power, Yuri Burda, Harri Edwards, Igor Babuschkin, and Vedant Misra. Grokking: Gen-
487 eralization beyond overfitting on small algorithmic datasets. *arXiv preprint arXiv:2201.02177*,
488 2022.
- 489 Adityanarayanan Radhakrishnan, Daniel Beaglehole, Parthe Pandit, and Mikhail Belkin. Mecha-
490 nism of feature learning in deep fully connected networks and kernel machines that recursively
491 learn features. *arXiv preprint arXiv:2212.13881*, 2022.
- 492 Adityanarayanan Radhakrishnan, Daniel Beaglehole, Parthe Pandit, and Mikhail Belkin. Mecha-
493 nism for feature learning in neural networks and backpropagation-free machine learning mod-
494 els. *Science*, 383(6690):1461–1467, 2024a. doi: 10.1126/science.adi5639. URL <https://www.science.org/doi/abs/10.1126/science.adi5639>.
- 496 Adityanarayanan Radhakrishnan, Mikhail Belkin, and Dmitriy Drusvyatskiy. Linear recursive fea-
497 ture machines provably recover low-rank matrices. *arXiv preprint arXiv:2401.04553*, 2024b.
498 URL <https://arxiv.org/pdf/2401.04553>.
- 499 Anna Rogers and Sasha Luccioni. Position: Key claims in llm research have a long tail of footnotes.
500 In *Forty-first International Conference on Machine Learning*, 2023. URL <https://arxiv.org/pdf/2308.07120>.
- 502 Rylan Schaeffer, Brando Miranda, and Sanmi Koyejo. Are emergent abilities of large language
503 models a mirage? In *Thirty-seventh Conference on Neural Information Processing Systems*,
504 2023. URL <https://openreview.net/forum?id=ITw9edRD1D>.
- 505 Vimal Thilak, Etai Littwin, Shuangfei Zhai, Omid Saremi, Roni Paiss, and Joshua Susskind. The
506 slingshot mechanism: An empirical study of adaptive optimizers and the grokking phenomenon.
507 *arXiv preprint arXiv:2206.04817*, 2022. URL <https://arxiv.org/abs/2206.04817>.
- 508 Shubhendu Trivedi, Jialei Wang, Samory Kpotufe, and Gregory Shakhnarovich. A consistent es-
509 timator of the expected gradient outerproduct. In *UAI*, pp. 819–828, 2014. URL <https://www.columbia.edu/~skk2175/Papers/GOP-UAI.pdf>.
- 511 Asher Trockman, Devin Willmott, and J Zico Kolter. Understanding the covariance structure of
512 convolutional filters. *arXiv preprint arXiv:2210.03651*, 2022.
- 513 Vikrant Varma, Rohin Shah, Zachary Kenton, János Kramár, and Ramana Kumar. Explain-
514 ing grokking through circuit efficiency. *International Conference on Learning Representations*
515 (*ICLR*), 2023. URL <https://openreview.net/pdf?id=7Zbg38nA0J>.
- 516 Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yo-
517 gatama, Maarten Bosma, Denny Zhou, Donald Metzler, Ed H. Chi, Tatsunori Hashimoto, Oriol
518 Vinyals, Percy Liang, Jeff Dean, and William Fedus. Emergent abilities of large language models.
519 *Transactions on Machine Learning Research (TMLR)*, 2022. URL <https://openreview.net/pdf?id=yzkSU5zdwD>.
- 521 Gan Yuan, Mingyue Xu, Samory Kpotufe, and Daniel Hsu. Efficient estimation of the central mean
522 subspace via smoothed gradient outer products. *arXiv preprint arXiv:2312.15469*, 2023. URL
523 <https://arxiv.org/pdf/2312.15469>.
- 524 Ziqian Zhong, Ziming Liu, Max Tegmark, and Jacob Andreas. The clock and the pizza: Two
525 stories in mechanistic explanation of neural networks. *Advances in Neural Information Processing*

526 *Systems*, 36, 2024.

527 Libin Zhu, Chaoyue Liu, Adityanarayanan Radhakrishnan, and Mikhail Belkin. Catapults in sgd:
528 spikes in the training loss and their impact on generalization through feature learning. *International Conference on Machine Learning (ICML)*, 235, 2024.
529

Algorithm 1 Recursive Feature Machine (RFM) (Radhakrishnan et al., 2024a)

Require: X, y, k, T, L \triangleright Train data: (X, y) , base kernel: k , iters.: T , matrix power: s , and bandwidth: L
 $M_0 = I_d$
for $t = 0, \dots, T - 1$ **do**
 Solve $\alpha \leftarrow k(X, X; M_t)^{-1}y$ $\triangleright f^{(t)}(x) = k(x, X; M_t)\alpha$
 $M_{t+1} \leftarrow [G(f^{(t)})]^s$
end for
return α, M_{T-1} \triangleright Solution to kernel regression: α , and feature matrix: M_{T-1}

530 A Additional discussion

531 **Low rank learning.** The problem of learning modular arithmetic can be viewed as a type of matrix
532 completion – completing the $p \times p$ matrix (so-called Cayley table) representing modular operations,
533 from partial observations. The best studied matrix completion problem is low rank matrix comple-
534 tion, where the goal is to fill in missing entries of a low rank matrix from observing a subset of the
535 entries (Moitra, 2018, Ch.8). While many specialized algorithms exist, it has been observed that
536 neural networks can recover low rank matrix structures Gunasekar et al. (2017). Notably, in a devel-
537 opment paralleling the results of this paper, low-rank matrix completion can provably be performed
538 by linear RFMs using the same AGOP mechanism Radhakrishnan et al. (2024b).

539 It is thus tempting to posit that grokking modular operations in neural networks or RFM can be
540 explained as a low rank prediction problem. Indeed modular operations can be implemented by an
541 index 4 model, i.e., a function of the form $f = g(Ax)$, where $x \in \mathbb{R}^{2p}$ and A is a rank 4 matrix (see
542 Appendix L for the construction). It is a plausible conjecture as there is strong evidence, empirical
543 and theoretical, that neural networks are capable of learning such multi-index models Damian et al.
544 (2022); Mousavi-Hosseini et al. (2022) as well as low-rank matrix completion. Furthermore, a phe-
545 nomenon similar to grokking was discussed in (Radhakrishnan et al., 2022, Fig. 5, 6) in the context
546 of low rank feature learning for both neural networks and RFM. However, despite the existence of
547 generalizable low rank models, the actual circulant features learned by both Neural Networks and
548 RFM are *not* low rank. Interestingly, this observation mirrors the problem of learning parity func-
549 tions through neural network inspired minimum norm interpolation, which was analyzed in Ardeshir
550 et al. (2023). While single-directional (index one) solutions exist in that setting, the authors show
551 that the minimum norm solutions are all multi-dimensional.

552 **Explanations for deep learning** Finally, this work adds to the growing body of evidence that
553 the AGOP-based mechanisms of feature learning can account for some of the most interesting phe-
554 nomena in deep learning. These include generalization with multi-index models (Parkinson et al.,
555 2023), deep neural collapse (Beaglehole et al., 2024b), and the ability to perform low-rank matrix
556 completion (Radhakrishnan et al., 2024b). Thus, RFM provides a framework that is both practically
557 powerful and serves as a theoretically tractable model of deep learning.

558 B Additional Preliminaries

559 For completeness we replicate the algorithm definition for Recursive Feature Machines (RFM) pro-
560 vided by Radhakrishnan et al. (2024a) in Algorithm 1. This procedure recursively fits a kernel
561 estimator for a chosen base kernel, k , then updates the feature matrix, M , by computing a matrix
562 power of the Average Gradient Outer Product (AGOP) for that estimator. The algorithm termi-
563 nates after a total of T iterations. The final estimator and feature matrix are then returned by the
564 algorithm.

565 C Neural Feature Ansatz

566 While the NFA has been observed generally across depths and architecture types (Radhakrishnan
567 et al., 2024a; Beaglehole et al., 2023, 2024a), we restate this observation for fully-connected net-
568 works with one hidden-layer of the form $f(x) = a^\top \phi(W_1 x)$.

569 **Ansatz 1** (Neural Feature Ansatz for one hidden layer). *For a one hidden-layer neural network f^{NN}*
570 *and a matrix power $\alpha \in (0, 1]$, the following holds:*

$$W_1^\top W_1 \propto G(f^{\text{NN}})^\alpha . \quad (5)$$

571 Note that this statement implies that $W_1^\top W_1$ and $G(f^{\text{NN}})^s$ have a cosine similarity of ± 1 .

572 In this work, we choose $\alpha = \frac{1}{2}$, following the main results in Radhakrishnan et al. (2024a). While
573 the absolute value of the cosine similarity is written in Eq. (5) to be 1, it is typically a high value less
574 than 1, where the exact value depends on choices of initialization, architecture, dataset, and training
575 procedure. For more understanding of these conditions, see Beaglehole et al. (2024a).

576 D Model and training details

577 **Gaussian kernel:** Throughout this work we take bandwidth $L = 2.5$ when using the Mahalanobis
578 Gaussian kernel. We solve ridgeless kernel regression using NumPy on a standard CPU.

579 **Neural networks:** Unless otherwise specified, we train one hidden layer neural networks with
580 quadratic activation functions and no biases in PyTorch on a single A100 GPU. Models are trained
581 using AdamW with hidden width 1024, batch size 32, learning rate of 10^{-3} , weight decay 1.0, and
582 standard PyTorch initialization. All models are trained using the Mean Squared Error loss function
583 (square loss).

584 For the experiments in Appendix Fig. 5, we train one hidden layer neural networks with quadratic
585 activation and no biases on modular addition modulo $p = 61$. We use 40% training fraction, PyTorch
586 standard initialization, hidden width of 512, weight decay 10^{-5} , and AGOP regularizer weight 10^{-3} .
587 Models are trained with vanilla SGD, batch size 128, and learning rate 1.0.

588 E Reordering feature matrices by group generators

589 Our reordering procedure uses the standard fact of group theory that the multiplicative group \mathbb{Z}_p^* is
590 a cyclic group of order $p - 1$ Koblitz (1994). By definition of the cyclic group, there exists at least
591 one element $g \in \mathbb{Z}_p^*$, known as a *generator*, such that $\mathbb{Z}_p^* = \{g^i ; i \in \{1, \dots, p - 1\}\}$.

592 Given a generator $g \in \mathbb{Z}_p^*$, we reorder features according to the map, $\phi_g : \mathbb{Z}_p^* \rightarrow \mathbb{Z}_p^*$, where if $h = g^i$,
593 then $\phi_g(h) = i$. In particular, given a matrix $B \in \mathbb{R}^{p \times p}$, we reorder the bottom right $(p - 1) \times (p - 1)$
594 sub-block of B as follows: we move the entry in coordinate (r, c) with $r, c \in \mathbb{Z}_p^*$ to coordinate
595 $(\phi_g(r), \phi_g(c))$. For example if $g = 2$ in \mathbb{Z}_5^* , then $(2, 3)$ entry of the sub-block would be moved to
596 coordinate $(1, 3)$ since $2^1 = 2$ and $2^3 \bmod 5 = 3$. In the setting of modular multiplication/division,
597 the map ϕ_g defined above is known as the *discrete logarithm* base g (Koblitz, 1994, Ch.3). The
598 discrete logarithm is analogous to the logarithm defined for positive real numbers in the sense that
599 it converts modular multiplication/division into modular addition/subtraction. Lastly, in this setting,
600 we note that we only reorder the bottom $(p - 1) \times (p - 1)$ sub-block of B as the first row and column
601 are 0 (as multiplication by 0 results in 0).

602 Upon re-ordering the $p \times p$ off-diagonal sub-blocks of the feature matrix by the map ϕ_g , the fea-
603 ture matrix of RFM for multiplication/division tasks contains circulant blocks as shown in Fig. 3C.
604 Thus, the reordered feature matrices for these tasks also exhibit the structure in Observation 1. As a
605 remark, we note that there can exist several generators for a cyclic group, and thus far, we have not
606 specified the generator g we use for re-ordering. For example, 2 and 3 are both generators of \mathbb{Z}_5^* since
607 $\{2, 2^2, (2^3 \bmod 5), (2^4 \bmod 5)\} = \{3, (3^2 \bmod 5), (3^3 \bmod 5), (3^4 \bmod 5)\} = \mathbb{Z}_5^*$. Lemma K.1 im-
608 plies that the choice of generator does not matter for observing circulant structure. As a convention,
609 we simply reorder by the smallest generator.

610 F Enforcing circulant structure in RFM

611 We see that the structure in Observation 1 gives generalizing features on modular arithmetic when
612 the circulant C is constructed from the RFM matrix. We observe that enforcing this structure at
613 every iteration, and comparing to the standard RFM model at that iteration, improves test loss and
614 accelerates grokking on e.g. addition (Appendix Fig. 2). The exact procedure to enforce this struc-
615 ture is as follows. We first perform standard RFM to generate feature matrices M_1, \dots, M_T . Then
616 for each iteration of the standard RFM, we construct a new \widetilde{M}_t on which we solve ridgeless kernel
617 regression for a new α and evaluate on the test set. To construct \widetilde{M} , we take $D = \mathbf{diag}(M_t)$ and
618 first let $\widetilde{M} = D^{-1/2} M D^{-1/2}$, to ensure the rows and columns have equal scale. We then reset the
619 top left and bottom right sub-matrices of \widetilde{M} as $I - \frac{1}{p} \mathbf{1}\mathbf{1}^T$, and replace the bottom-left and top-right

620 blocks with C and C^\top , where C is an exactly circulant matrix constructed from M_t . Specifically,
 621 where c is the first column of the bottom-left sub-matrix of M_t , column ℓ of C is equal to $\sigma^\ell(M_t)$.

622 G Fourier multiplication algorithm from circulant features

623 As stated in the main text, using certain circulant matrices, kernel regression will learn the Fourier
 624 Multiplication Algorithm (FMA). We state the FMA for modular addition/subtraction from Nanda
 625 et al. (2023) below. While these prior works write this algorithm in terms of cosines and sines, our
 626 presentation simplifies the statement by using the DFT.

627 **Complex inner product and Discrete Fourier Transform (DFT).** In our theoretical analysis in
 628 Section 5, we will utilize the following notions of complex inner product and DFT. The complex
 629 inner product $\langle \cdot, \cdot \rangle_{\mathbb{C}}$ is a map from $\mathbb{C}^d \times \mathbb{C}^d \rightarrow \mathbb{C}$ of the form

$$\langle u, v \rangle_{\mathbb{C}} = u^\top \bar{v}, \quad (6)$$

630 where \bar{v}_j is the complex conjugate of v_j . Let $i = \sqrt{-1}$ and let $\omega = \exp(\frac{-2\pi i}{d})$. The DFT is the map
 631 $\mathcal{F} : \mathbb{C}^d \rightarrow \mathbb{C}^d$ of the form $\mathcal{F}(u) = Fu$, where $F \in \mathbb{C}^{d \times d}$ is a unitary matrix with $F_{ij} = \frac{1}{\sqrt{d}}\omega^{ij}$. In
 632 matrix form, F is given as

$$F = \frac{1}{\sqrt{d}} \begin{pmatrix} 1 & 1 & 1 & \dots & 1 \\ 1 & \omega & \omega^2 & \dots & \omega^{d-1} \\ 1 & \omega^2 & \omega^4 & \dots & \omega^{2(d-1)} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & \omega^{d-1} & \omega^{2(d-1)} & \dots & \omega^{(d-1)(d-1)} \end{pmatrix}. \quad (7)$$

633 **Fourier Multiplication Algorithm for modular addition/subtraction.** Consider the modular addi-
 634 tion task with $f^*(a, b) = (a + b) \bmod p$. For a given input $x = x_{[1]} \oplus x_{[2]} \in \mathbb{R}^{2p}$, the FMA
 635 generates a value for output class ℓ , $y_{\text{add}}(x; \ell)$, through the following computation:

- 636 1. Compute the Discrete Fourier Transform (DFT) for each digit vector $x_{[1]}$ and $x_{[2]}$, which
 637 we denote $\hat{x}_{[1]} = Fx_{[1]}$ and $\hat{x}_{[2]} = Fx_{[2]}$ where the matrix F is defined in Eq. (7).
- 638 2. Compute the element-wise product $\hat{x}_{[1]} \odot \hat{x}_{[2]}$.
- 639 3. Return $\sqrt{p} \cdot \langle \hat{x}_{[1]} \odot \hat{x}_{[2]}, Fe_\ell \rangle_{\mathbb{C}}$ where e_ℓ denotes ℓ -th standard basis vector and $\langle \cdot, \cdot \rangle_{\mathbb{C}}$
 640 denotes the complex inner product (see Eq. (6)).

641 This algorithmic process can be written concisely in the following equation:

$$y_{\text{add}}(x; \ell) = \sqrt{p} \cdot \langle Fx_{[1]} \odot Fx_{[2]}, Fe_\ell \rangle_{\mathbb{C}}. \quad (8)$$

642 Note that for $x = e_a \oplus e_b$, the second step of the FMA reduces to

$$Fe_a \odot Fe_b = \frac{1}{\sqrt{p}} Fe_{(a+b) \bmod p}. \quad (9)$$

643 Using the fact that F is a unitary matrix, the output of the FMA is given by

$$\sqrt{p} \cdot \left\langle \frac{1}{\sqrt{p}} Fe_{(a+b) \bmod p}, Fe_\ell \right\rangle_{\mathbb{C}} = e_{(a+b) \bmod p}^\top F^\top \bar{F} e_\ell = e_{(a+b) \bmod p}^\top e_\ell = \mathbb{1}_{\{(a+b) \bmod p = \ell\}}. \quad (10)$$

644 Thus, the output of the FMA is a vector $e_{(a+b) \bmod p}$, which is equivalent to modular addition. We
 645 provide an example of this algorithm for $p = 3$ in Appendix J.

646 **Remarks.** We note that our description of the FMA uses all entries of the DFT, referred to as fre-
 647 quencies, while the algorithm as proposed in prior works allows for utilizing a subset of frequencies.
 648 Also note that the FMA for subtraction, written y_{sub} , is similar and given by

$$y_{\text{sub}}(x; \ell) = \sqrt{p} \cdot \langle Fx_{[1]} \odot Fe_{p-\ell-1}, Fx_{[2]} \rangle_{\mathbb{C}}. \quad (11)$$

649 Having described the FMA, we now state our theorem.

650 **Theorem G.1.** Given all of the discrete data $\{(e_a \oplus e_b, e_{(a-b) \bmod p})\}_{a,b=0}^{p-1}$, for each output class
651 $\ell \in \{0, \dots, p-1\}$, suppose we train a separate kernel predictor $f_\ell(x) = k(x, X; M_\ell) \alpha^{(\ell)}$ where
652 $k(\cdot; \cdot; M_\ell)$ is a quadratic kernel with $M_\ell = \begin{pmatrix} 0 & C^\ell \\ (C^\ell)^\top & 0 \end{pmatrix}$ and $C \in \mathbb{R}^{p \times p}$ is a circulant matrix
653 with first row e_1 . When $\alpha^{(\ell)}$ is the solution to kernel ridgeless regression for each ℓ , the kernel pre-
654 predictor $f = [f_0, \dots, f_{p-1}]$ is equivalent to Fourier Multiplication Algorithm for modular subtraction
655 (Eq. (11)).

656 As C is circulant, C^ℓ is also circulant. Hence, each M_ℓ has the structure described in Observation 1,
657 where $A = 0$. Note our construction differs from RFM in that we use a different feature matrix
658 M_ℓ for each output coordinate, rather than a single feature matrix across all output coordinates.
659 Nevertheless, Theorem G.1 provides support for the fact that block-circulant feature matrices can be
660 used to solve modular arithmetic.

661 We provide the proof for Theorem G.1 in Appendix K. The argument for the FMA for addition
662 (Eq. (8)) is identical provided we replace C^ℓ with $C^\ell R$ and $(C^\ell)^\top$ with $(C^\ell R)^\top$ in each M_ℓ , where
663 R is the Hankel matrix that reverses the row order (i.e. ones along the main anti-diagonal, zero's
664 elsewhere), whose first row is e_{p-1} . An analogous result follows for multiplication and division
665 under re-ordering by a group element, as described in Section 3.

666 Our proof uses the well-known fact that circulant matrices can be diagonalized using the DFT matrix
667 (Gray et al., 2006) (see Lemma K.2 for a restatement of this fact). This fundamental relation intu-
668 itively connects circulant features and the FMA. By using kernels with block-circulant Mahalanobis
669 matrices, we effectively represent the one-hot encoded data in terms of their Fourier transforms. We
670 conjecture that this implicit representation is what enables RFM to learn modular arithmetic with
671 more general circulant matrices when training on just a fraction of the discrete data.

672 H Grokking multiple tasks

673 Throughout the main paper, we focused on modular arithmetic settings for a single task. In more
674 general domains such as language, one may expect there to be many ‘‘skills’’ that need to be learned.
675 In such settings, it is possible that these skills are grokked at different rates. While a full discussion
676 is beyond the scope of this work, to illustrate this behavior, we performed additional experiments
677 in here, where we train RFM on a pair of modular arithmetic tasks simultaneously and demonstrate
678 that different tasks are indeed grokked at different points throughout training.

679 We train RFM to simultaneously solve the following two modular polynomial tasks: (1) $x + y \bmod p$
680 ; (2) $x^2 + y^2 \bmod p$ for modulus $p = 61$. We train RFM with the Mahalanobis Gaussian kernel
681 using bandwidth parameter $L = 2.5$. Training data for both tasks is constructed from the same
682 80% training fraction. In addition to concatenating the one-hot encodings for x, y , we also append
683 an extra bit indicating which task to solve (0 indicating task (1) and 1 indicating task (2)). The
684 classification head is shared for both tasks (e.g. output dimension is still \mathbb{R}^p).

685 In Appendix Fig. 3, we observe that there are two sharp transitions in the test loss and test accuracy.
686 By decomposing the loss into the loss per task, we observe that RFM groks task (1) prior to grokking
687 task (2). Overall, these results illustrate that grokking of different tasks can occur at different training
688 iterations.

689 I AGOP regularization and weight decay for grokking modular arithmetic.

690 It has been argued in prior work that weight decay (ℓ_2 regularization on network weights) is neces-
691 sary for grokking to occur when training neural networks for modular arithmetic tasks (Varma et al.,
692 2023; Davies et al., 2023; Nanda et al., 2023). Under the NFA (Eq. (5)), which states that $W_1^\top W_1$ is
693 proportional to a matrix power of $G(f)$, we expect that performing weight decay on the first layer,
694 i.e., penalizing the loss by $\|W_1\|_F^2 = \text{tr}(W_1^\top W_1)$, should behave similarly to penalizing the trace of
695 the AGOP, $\text{tr}(G(f))$, during training.³ To this end, we compare the impact of using (1) no regular-
696 ization; (2) weight decay; and (3) AGOP regularization when training neural networks on modular
697 arithmetic tasks. In Appendix Fig. 5, we find that, akin to weight decay, AGOP regularization leads

³We note this regularizer been used prior work where AGOP is called the Gram matrix of the input-output Jacobian Hoffman et al. (2019).

698 to grokking in cases where using no regularization results in no grokking and poor generalization.
 699 These results provide further evidence that neural networks solve modular arithmetic by using the
 700 AGOP to learn features.

701 J FMA example for $p = 3$

702 We now provide an example of the FMA for $p = 3$. Let $x = e_1 \oplus e_2$. In this case, we expect the
 703 FMA to output the vector e_0 since $(1 + 2) \bmod 3 = 0$. Following the first step of the FMA, we
 704 compute

$$\hat{x}_{[1]} = Fe_1 = \frac{1}{\sqrt{3}}[1, \omega, \omega^2]^\top ; \quad \hat{x}_{[2]} = Fe_2 = \frac{1}{\sqrt{3}}[1, \omega^2, \omega^4]^\top , \quad (12)$$

705 which are the first and second columns of F , respectively. Then their element-wise product is given
 706 by

$$Fe_1 \odot Fe_2 = \frac{1}{3}[1, \omega^3, \omega^6]^\top = \frac{1}{3}[1, 1, 1]^\top = \frac{1}{\sqrt{3}}Fe_0 , \quad (13)$$

707 which is $\frac{1}{\sqrt{3}}$ times the first column of the DFT matrix. Finally, we compute the outputs
 708 $\sqrt{3} \left\langle \frac{1}{\sqrt{3}}Fe_0, Fe_\ell \right\rangle_{\mathbb{C}}$ for each $\ell \in \{0, 1, 2\}$. As F is unitary, $y_{\text{add}}(e_1 \oplus e_2; \ell) = \mathbb{1}_{\{1+2=\ell \bmod 3\}}$, so
 709 that coordinate 0 of the output will have value 1, and all other coordinates have value 0.

710 K Additional results and proofs

711 **Lemma K.1.** *Let $C \in \mathbb{R}^{p \times p}$ with its first row and column entries all equal to 0. Let the $(p -$
 712 $1) \times (p - 1)$ sub-block starting at the second row and column be C^\times . Then, C^\times is either circulant
 713 after re-ordering by any generator q of \mathbb{Z}_p^* , or C^\times is not circulant under re-ordering by any such
 714 generator.*

715 *Proof of Lemma K.1.* We prove the lemma by showing that for any two generators q_1, q_2 of \mathbb{Z}_p^* , if
 716 C^\times is circulant re-ordering with q_1 , then it is also circulant when re-ordering by q_2 .

717 Suppose C^\times is circulant re-ordering with q_1 . Let $i, j \in \{1, \dots, p - 1\}$. Note that by the circulant
 718 assumption, for all $s \in \mathbb{Z}$,

$$C_{q_1^i, q_1^j} = C_{q_1^{i+s}, q_1^{j+s}} , \quad (14)$$

719 where we take each index modulo p .

720 As q_2 is a generator for \mathbb{Z}_p^* , we can access all entries of C^\times by indexing with powers of q_2 . Further,
 721 as q_1 is a generator, we can write $q_2 = q_1^k$, for some power k . Let $a \in \mathbb{Z}$. Then,

$$\begin{aligned} C_{q_2^i, q_2^j} &= C_{q_1^{ki}, q_1^{kj}} \\ &= C_{q_1^{ki+ka}, q_1^{kj+ka}} \\ &= C_{q_1^{k(i+a)}, q_1^{k(j+a)}} \\ &= C_{q_2^{i+a}, q_2^{j+a}} . \end{aligned}$$

722 Therefore, C is constant on the diagonals under re-ordering by q_2 , concluding the proof. \square

723 We next state Lemma K.2, which is used in the proof of Theorem G.1.

724 **Lemma K.2** (See, e.g., Gray et al. (2006)). *Circulant matrices U can be written (diagonalized) as:*

$$U = FD\bar{F}^\top ,$$

725 where F is the DFT matrix, \bar{F}^\top is the element-wise complex conjugate of F^\top (i.e. the Hermitian of
 726 F), and D is a diagonal matrix with diagonal $\sqrt{p} \cdot Fu$, where u is the first row of U .

727 We now present the proof of Theorem G.1, restating the theorem below for the reader's convenience.

728 **Theorem.** Given all of the discrete data $\{(e_a \oplus e_b, e_{(a-b) \bmod p})\}_{a,b=0}^{p-1}$ in modular subtraction
729 task, for each output class $\ell \in \{0, \dots, p-1\}$, we train a separate kernel predictor $f_\ell(x) =$
730 $k(x, X; M_\ell) \alpha^{(\ell)}$. Here $k(\cdot, \cdot; M_\ell)$ is a quadratic kernel with $M_\ell = \begin{pmatrix} 0 & C^\ell \\ (C^\ell)^\top & 0 \end{pmatrix}$ and $C \in \mathbb{R}^{p \times p}$
731 is a circulant matrix with first row e_1 . When $\alpha^{(\ell)}$ is the solution to kernel ridgeless regression for
732 each ℓ , the kernel predictor $f = [f_0, \dots, f_{p-1}]$ is equivalent to Fourier Multiplication Algorithm
733 for modular subtraction (Eq. (11)).

734 *Proof of Theorem G.1.* We present the proof for modular subtraction as the proof for addition fol-
735 lows analogously. We write the standard kernel predictor for class ℓ on input $x = x_{[1]} \oplus x_{[2]} \in \mathbb{R}^{2p}$
736 as,

$$f_\ell(x) = \sum_{a,b=0}^{p-1} \alpha_{a,b}^{(\ell)} k(x, e_a \oplus e_b; M_\ell),$$

737 where we have re-written the index into kernel coefficients for class ℓ , $\alpha^{(\ell)} \in \mathbb{R}^{p \times p}$, so that the
738 coefficients are multi-indexed by the first and second digit. Specifically, now $\alpha_{a,b}^{(\ell)}$ is the kernel
739 coefficient corresponding to the representer $k(\cdot, x)$ for input point $x = e_a \oplus e_b$. Recall we use a
740 quadratic kernel, $k(x, z; M_\ell) = (x^\top M_\ell z)^2$. In this case, the kernel predictor simplifies to,

$$f_\ell(x) = \sum_{a,b=0}^{p-1} \alpha_{a,b}^{(\ell)} \left(x_{[1]}^\top C^\ell e_b + e_a^\top C^\ell x_{[2]} \right)^2.$$

741 Then, the labels for each pair of input digits, written as a matrix $Y^{(\ell)} \in \mathbb{R}^{p \times p}$ for the ℓ -th class
742 where the row and column index the first and second digit respectively, are $Y^{(\ell)} = C^{-\ell}$.

743 For $x = e_{a'} \oplus e_{b'}$, i.e. x in the discrete dataset, we have,

$$\begin{aligned} f_\ell(x) &= \sum_{a,b=0}^{p-1} \alpha_{a,b}^{(\ell)} (\delta_{(a,b'-\ell)} + \delta_{(a',b-\ell)} + 2\delta_{(a,b'-\ell)}\delta_{(a',b-\ell)}) \\ &= e_{b'-\ell}^\top \alpha^{(\ell)} \mathbf{1} + \mathbf{1}^\top \alpha^{(\ell)} e_{a'+\ell} + 2e_{b'-\ell}^\top \alpha^{(\ell)} e_{a'+\ell} \\ &= e_{b'}^\top C^{-\ell} \alpha^{(\ell)} \mathbf{1} + \mathbf{1}^\top \alpha^{(\ell)} C^{-\ell} e_{a'} + 2e_{b'}^\top C^{-\ell} \alpha^{(\ell)} C^{-\ell} e_{a'} \\ &= e_{b'}^\top (C^{-\ell} \alpha \mathbf{1} \mathbf{1}^\top + \mathbf{1} \mathbf{1}^\top \alpha C^{-\ell} + 2C^{-\ell} \alpha C^{-\ell}) e_{a'}, \end{aligned}$$

744 where $\delta_{(u,v)} = \mathbb{1}_{\{u=v\}}$. Let $f_\ell(X) \in \mathbb{R}^{p \times p}$ be the matrix of function values of f_ℓ , where
745 $[f_\ell(X)]_{a,b} = f_\ell(e_a \oplus e_b)$, and, therefore, $f_\ell(e_a \oplus e_b) = e_a^\top f_\ell(X) e_b$. Then, to solve for $\alpha^{(\ell)}$,
746 we need to solve the system of equations for α ,

$$\begin{aligned} f_\ell(X) &= (C^{-\ell} \alpha \mathbf{1} \mathbf{1}^\top + \mathbf{1} \mathbf{1}^\top \alpha C^{-\ell} + 2C^{-\ell} \alpha C^{-\ell})^\top = C^{-\ell} \\ &\iff C^{-\ell} \alpha \mathbf{1} \mathbf{1}^\top + \mathbf{1} \mathbf{1}^\top \alpha C^{-\ell} + 2C^{-\ell} \alpha C^{-\ell} = C^\ell \end{aligned}$$

747 Note, by left-multiplying both sides by $C^{-\ell}$, we see this equation holds iff,

$$C^{-2\ell} \alpha \mathbf{1} \mathbf{1}^\top + \mathbf{1} \mathbf{1}^\top \alpha C^{-\ell} + 2C^{-2\ell} \alpha C^{-\ell} = I.$$

748 Note the solution is unique as the kernel matrix is full rank. We posit the solution α such that
749 $C^{-2\ell} \alpha C^{-\ell} = \frac{1}{2}I + \lambda \mathbf{1} \mathbf{1}^\top$, which is $\alpha = \frac{1}{2}C^{3\ell} + \lambda \mathbf{1} \mathbf{1}^\top$. Then, solving for λ , we require,

$$\mathbf{1} \mathbf{1}^\top + 2p\lambda \mathbf{1} \mathbf{1}^\top + 2\lambda \mathbf{1} \mathbf{1}^\top = 0,$$

750 which implies $\lambda = -\frac{2}{2p+2}$. Substituting this value of λ and simplifying, we see finally that
751 $f_\ell(x) = x_{[1]}^\top C^{-\ell} x_{[2]}$. Therefore, using that circulant matrices are diagonalized by $C = \sqrt{p} F D \bar{F}^\top$
752 (Lemma K.2) and $\bar{F}^\top F = I$, where $D = \text{diag}(F e_1)$, we derive,

$$\begin{aligned} f_\ell(x) &= \sqrt{p} \cdot x_{[1]}^\top F D^{-\ell} \bar{F}^\top x_{[2]} \\ &= \sqrt{p} \cdot x_{[1]}^\top F \text{diag}(F e_{p-\ell-1}) \bar{F}^\top x_{[2]} \\ &= \sqrt{p} \cdot \langle F x_{[1]} \odot F e_{p-\ell-1}, F x_{[2]} \rangle_{\mathbb{C}} \end{aligned}$$

753 which is the output of the FMA on modular subtraction. \square

754 **L Low rank solution to modular arithmetic**

755 **Addition** We present a solution to the modular addition task whose AGOP is low rank, in contrast
 756 to the full rank AGOP recovered by RFM and neural networks.

We define the “encoding” map $\Phi : \mathbb{R}^p \rightarrow \mathbb{C}$ as follows. For a vector $\mathbf{a} = [a_0, \dots, a_{p-1}]$,

$$\Phi(\mathbf{a}) = \sum_{k=0}^{p-1} a_k \exp\left(\frac{k \cdot 2\pi i}{p}\right) .$$

Notice that Φ is a linear map such that $\Phi(\mathbf{e}_k) = \exp\left(\frac{k \cdot 2\pi i}{p}\right)$. Notice also that Φ is partially invertible with the “decoding” map $\Psi : \mathbb{C} \rightarrow \mathbb{R}^p$.

$$\Psi(z) = \widetilde{\max} \left(\left\langle z, \exp\left(\frac{0 \cdot 2\pi i}{p}\right) \right\rangle, \dots, \left\langle z, \exp\left(\frac{(p-1) \cdot 2\pi i}{p}\right) \right\rangle \right) .$$

757 Above $\widetilde{\max}$ is a function that makes all entries zero except for the largest one and the inner product
 758 is the usual inner product in \mathbb{C} considered as \mathbb{R}^2 . Thus

$$\Psi \left(\exp\left(\frac{k \cdot 2\pi i}{p}\right) \right) = \mathbf{e}_k . \tag{15}$$

759 Ψ is a nonlinear map $\mathbb{C} \rightarrow \mathbb{R}^p$. While it is discontinuous but can easily be modified to make it
 760 differentiable.

By slight abuse of notation, we will define $\Phi : \mathbb{R}^p \times \mathbb{R}^p \rightarrow \mathbb{C}^2$ on pairs:

$$\Phi(\mathbf{e}_j, \mathbf{e}_k) = (\Phi(\mathbf{e}_j), \Phi(\mathbf{e}_k)) .$$

761 This is still a linear map but now to \mathbb{C}^2 .

Consider now a quadratic map M on $\mathbb{C}^2 \rightarrow \mathbb{C}$ given by complex multiplication:

$$M(z_1, z_2) = z_1 z_2 .$$

It is clear that the composition $\Psi M \Phi$ implements modular addition

$$\Psi M \Phi(\mathbf{e}_j, \mathbf{e}_k) = \mathbf{e}_{(j+k) \bmod p}$$

762 Furthermore, since Φ is a liner map to a four-dimensional space, the AGOP of the composition
 763 $\Psi M \Phi$ is of rank 4.

Multiplication The construction is for multiplication is very similar with modifications which we sketch below. We first re-order the non-zero coordinates by the discrete logarithm with base equal to a generator of the multiplicative group e_g (see Appendix E), while keeping the order of index 0. Then, we modify Φ to remove index a_0 from the sum for inputs \mathbf{a} . Thus for multiplication,

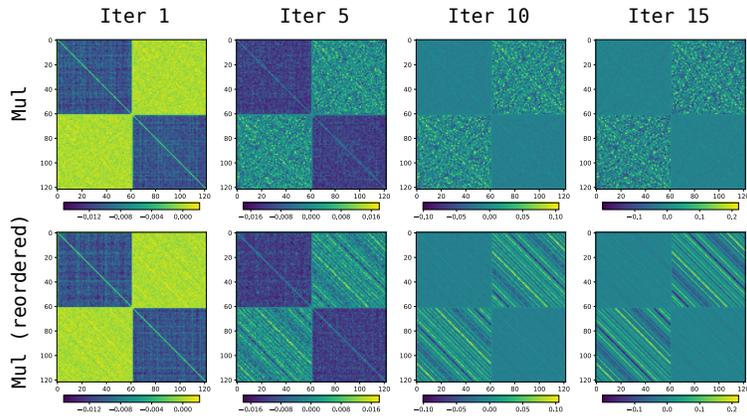
$$\Phi(\mathbf{a}) = \sum_{k=1}^{p-1} a_k \exp\left(\frac{k \cdot 2\pi i}{p-1}\right) ,$$

764 Hence that $\Phi(\mathbf{e}_0) = 0$, $\Phi(\mathbf{e}_g) = \exp\left(\frac{2\pi i}{p-1}\right)$ and $\Phi(\mathbf{e}_{g^k}) = \exp\left(\frac{k \cdot 2\pi i}{p-1}\right)$. We extend Φ to $\mathbb{R}^p \times \mathbb{R}^p$
 765 as in Eq. 15 above. Note that Φ and the re-ordering together are still a linear map of rank 4.

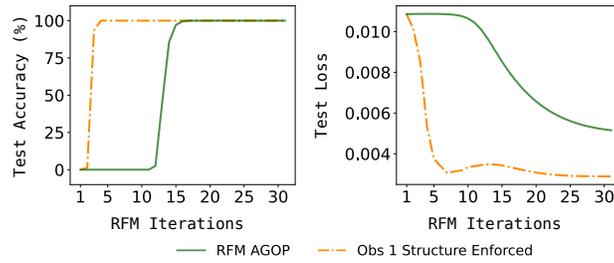
Then, the “decoding” map, $\Psi(z)$, will be modified to return 0, when $z = 0$, and otherwise,

$$\Psi(z) = g^{\widetilde{\max}(\langle z, \exp(\frac{0 \cdot 2\pi i}{p-1}) \rangle, \dots, \langle z, \exp(\frac{(p-2) \cdot 2\pi i}{p-1}) \rangle)} .$$

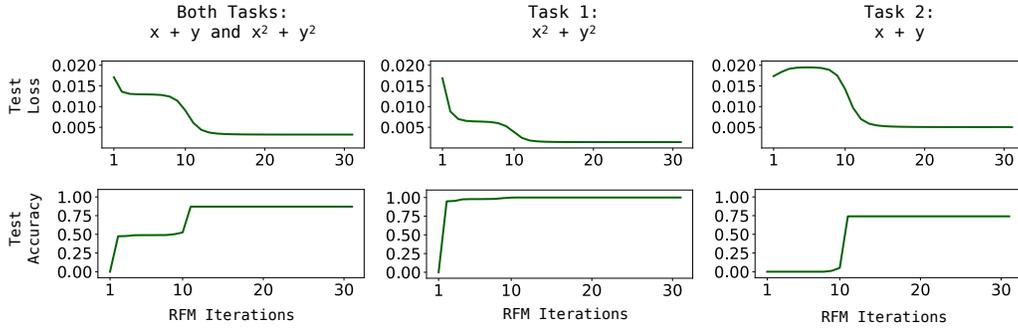
766 M is still defined as above. It is easy to check that the composition of $\Psi M \Phi$ with reordering
 767 implements modular multiplication modulo p and furthermore, the AGOP will also be of rank 4.



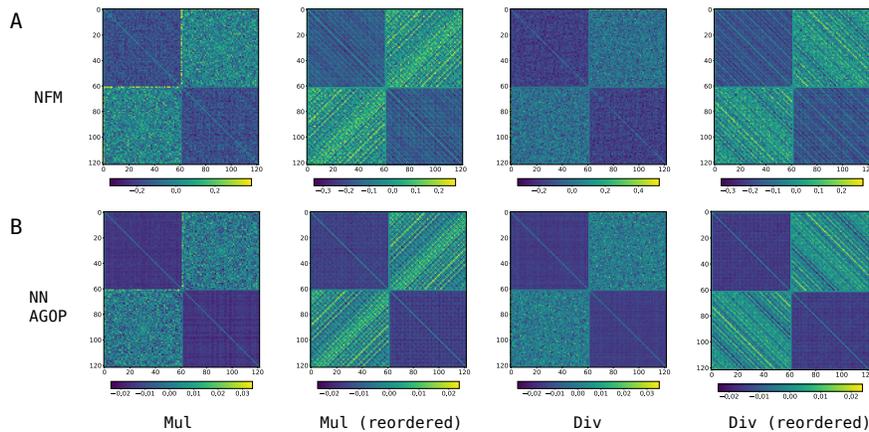
Appendix Figure 1: AGOP evolution for quadratic RFM trained on modular multiplication with $p = 61$ before reordering (top row) and after reordering by the logarithm base 2 (bottom row).



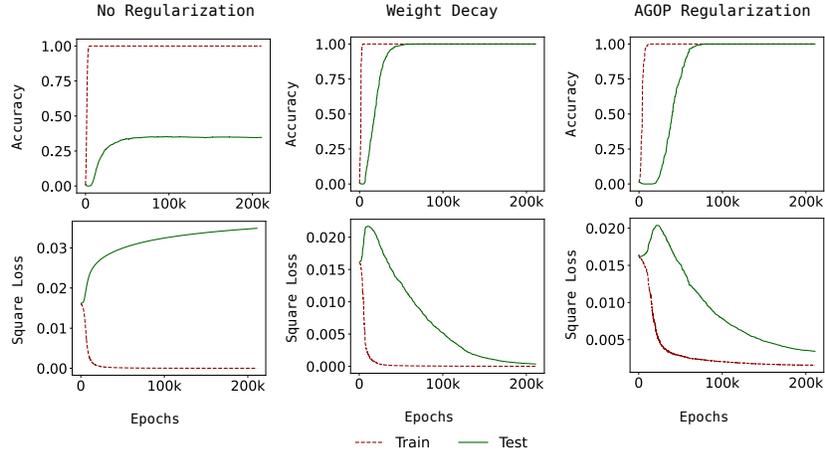
Appendix Figure 2: We train a Gaussian kernel-RFM on $x + y \bmod 97$ and plot test loss and accuracy versus RFM iterations. We also evaluate the performance of the same model upon modifying the M matrix to have exact block-circulant structure stated in Observation 1.



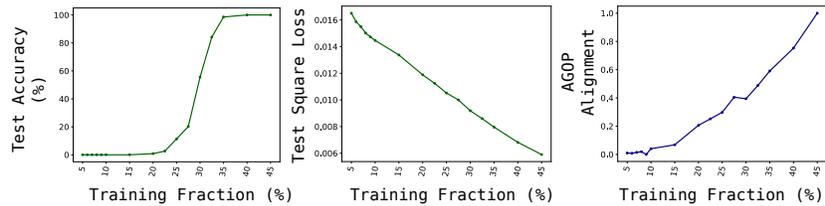
Appendix Figure 3: RFM with the Gaussian kernel trained on two modular arithmetic tasks with modulus $p = 61$. Task 1 is to learn $x^2 + y^2 \bmod p$ and task 2 is to learn $x + y \bmod p$.



Appendix Figure 4: (A) We visualize the neural feature matrix (NFM) from a one hidden layer neural network with quadratic activations trained on modular multiplication and division, before and after reordering by the discrete logarithm. (B) We visualize the square root of the AGOP of the neural network in (A) before and after reordering.



Appendix Figure 5: One hidden layer fully connected networks with quadratic activations trained on modular addition with $p = 61$ with vanilla SGD. Without any regularization the test accuracy does not go to 100% whereas using weight decay or regularizing using the trace of the AGOP result in 100% test accuracy and grokking.



Appendix Figure 6: We train kernel-RFMs for 30 iterations using the Mahalanobis Gaussian kernel for $x + y \bmod 97$. We plot test accuracy, test loss, and AGOP alignment versus percentage of training data used (denoted training fraction). All models reach convergence (i.e., both the test loss and test accuracy no longer change) after 30 iterations. We observe a sharp transition in test accuracy with respect to the training fraction, but we observe gradual change in test loss and AGOP alignment with respect to the training data fraction.