
Composing Global Optimizers to Reasoning Tasks via Algebraic Objects in Neural Nets

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 We prove rich algebraic structures of the solution space for 2-layer neural net-
2 works with quadratic activation and L_2 loss, trained on reasoning tasks in Abelian
3 group (e.g., modular addition). Such a rich structure enables *analytical* construc-
4 tion of global optimal solutions from partial solutions that only satisfy part of the
5 loss, despite its high nonlinearity. We coin the framework as CoGO (*Composing*
6 *Global Optimizers*). Specifically, we show that the weight space over different
7 numbers of hidden nodes of the 2-layer network is equipped with a semi-ring
8 algebraic structure, and the loss function to be optimized consists of *monomial*
9 *potentials*, which are ring homomorphism, allowing partial solutions to be com-
10 posed into global ones by ring addition and multiplication. Our experiments show
11 that around 95% of the solutions obtained by gradient descent match exactly our
12 theoretical constructions. Although the global optimizers constructed only re-
13 quired a small number of hidden nodes, our analysis on gradient dynamics shows
14 that overparameterization asymptotically decouples training dynamics and is ben-
15 efiticial. We further show that training dynamics favors simpler solutions under
16 weight decay, and thus high-order global optimizers such as perfect memorization
17 are unfavorable.

18 1 Introduction

19 Large Language Models (LLMs) have shown impressive results in various disciplines (OpenAI,
20 2024; Anthropic; Team, 2024b,a; Dubey et al., 2024; Jiang et al., 2023), while they also make sur-
21 prising mistakes in basic reasoning tasks (Nezhurina et al., 2024; Berglund et al., 2023). Therefore,
22 it remains an open problem whether it can truly do reasoning tasks. On one hand, existing works
23 demonstrate that the models can learn efficient algorithms (e.g., dynamic programming (Ye et al.,
24 2024) for language structure modeling, etc) and good representations (Jin & Rinard, 2024; Wijmans
25 et al., 2023). Some reports emergent behaviors (Wei et al., 2022) when scaling up with data and
26 model size. On the other hand, many works also show that LLMs cannot self-correct (Huang et al.,
27 2023), and cannot generalize very well beyond the training set for simple tasks (Dziri et al., 2023;
28 Yehudai et al., 2024; Ouellette et al., 2023), let alone complicated planning tasks (Kambhampati
29 et al., 2024; Xie et al., 2024).

30 To understand how the model performs reasoning and further improve its reasoning power, people
31 have been studying simple arithmetic reasoning problems in depth. Modular addition (Nanda et al.,
32 2023; Zhong et al., 2024), i.e., predicting $a + b \pmod d$ given a and b , is a popular one due to its
33 simple and intuitive structure yet surprising behaviors in learning dynamics (e.g., grokking (Power
34 et al., 2022)) and learned representations (e.g., Fourier bases (Zhou et al., 2024)). Most works
35 focus on various metrics to measure the behaviors and extracting interpretable circuits from trained
36 models (Nanda et al., 2023; Varma et al., 2023; Huang et al., 2024). Analytic solutions can be

37 constructed and/or reverse-engineered (Gromov, 2023; Zhong et al., 2024; Nanda et al., 2023) but it
38 is not clear how to construct a systematic framework to explain and generalize the results.

39 In this work, we systematically analyze 2-layer neural networks with quadratic activation and L_2
40 loss on predicting the outcome of group multiplication in Abelian group G , which is an extension
41 of modular addition. We find that global optimizers can be constructed *algebraically* from small
42 partial solutions that are optimal only for parts of the loss. We achieve this by showing that (1) for
43 the 2-layer network, there exists a *semi-ring* structure over the weights space *across different order*
44 (i.e., number of hidden nodes or network width), with specifically defined addition and multipli-
45 cation (Sec. 4.1), and (2) the L_2 loss is a function of *monomial potentials* (MPs), which are ring
46 homomorphisms (Theorem 1) that allow compositions of partial solutions into global ones using
47 ring addition and multiplication.

48 As a result, our theoretical framework, named CoGO (i.e., *Composing Global Optimizers*), success-
49 fully constructs two distinct types of Fourier-based global optimizers of per-frequency order 4 (or
50 “ 2×2 ”) and order 6 (or “ 2×3 ”), and a global optimizer of order d^2 that correspond to perfect mem-
51 orization. Empirically, we demonstrate that around 95% of the solutions obtained from gradient
52 descent (with weight decay) have the predicted structure and match exactly with our theoretical con-
53 struction of order-4 and order-6 solutions. In addition, we also analyze the training dynamics, and
54 show that the dynamics favors low-order global optimizers, since global optimizers algebraically
55 connected by ring multiplication can be proven to also be topologically connected. Therefore, high-
56 order solution like perfect memorization is unfavorable in the dynamics. When the network width
57 goes to infinite, the dynamics of monomial potentials becomes decoupled, demystifying why over-
58 parameterization improves the performance.

59 To our best knowledge, we are the first to discover such algebraic structures inside network training,
60 apply it to analyze solutions to reasoning tasks such as modular additions, and show our theoretical
61 constructions occur in actual gradient descent solutions.

62 2 Related Works

63 **Algebraic structures for machine learning.** Many works leverage symmetry and group structure
64 in deep learning. For example, in geometric deep learning, different forms of symmetry are incor-
65 porated into network architectures (Bronstein et al., 2021). However, they do not open the black
66 box and explore the algebraic structures of the network itself during training.

67 **Expressibility.** Existing works on expressibility (Li et al., 2024; Liu et al., 2022) gives explicit
68 weight construction of neural networks weights (e.g., Transformers) for reasoning tasks like au-
69 tomata, which includes modular addition. However, their works do not discover algebraic structures
70 in the weight space and loss, nor learning dynamics analysis, and it is not clear whether the con-
71 structed weights coincide with the actual solutions found by gradient descent, even in synthetic data.

72 **Fourier Bases in Arithmetic Tasks.** Existing works discovered that pre-trained models use Fourier
73 bases for arithmetic operations (Zhou et al., 2024). This is true even for a simple Transformer, or
74 even a network with one hidden layer (Morwani et al., 2023). Previous works also construct ana-
75 lytic Fourier solutions (Gromov, 2023) for modular addition, but with the additional assumption of
76 infinite width, unaware of the algebraic structures we discover. Existing theoretical work (Morwani
77 et al., 2023) also shows group-theoretical results on algebraic tasks related to finite groups, also for
78 networks with one-hidden layers and quadratic activations. Compared to ours, they use the max-
79 margin framework with a special regularization ($L_{2,3}$ norm) rather than L_2 loss, do not characterize
80 and leverage algebraic structures in the weight space, and do not analyze the training dynamics.

81 3 Decoupling L_2 Loss for reasoning tasks of Abelian group

82 **Basic group theory.** A set G forms a *group*, which means that (1) there exists an operation \cdot (i.e.,
83 “multiplication”): $G \times G \mapsto G$ and it satisfies association: $(g_1 \cdot g_2) \cdot g_3 = g_1 \cdot (g_2 \cdot g_3)$. Often we write
84 $g_1 g_2$ instead of $g_1 \cdot g_2$ for brevity. (2) there exists an identity element $e \in G$ so that $eg = ge = g$,
85 (3) for every group element $g \in G$, there is a unique inverse g^{-1} so that $gg^{-1} = g^{-1}g = e$. In some
86 groups, the multiplication operation is commutative, i.e., $gh = hg$ for any $g, h \in G$. Such groups
87 are called *Abelian group*. Modular addition forms a Abelian (more specifically, cyclic) group by
88 noticing that there exists a mapping $a \mapsto e^{2\pi ai/d}$ and $a+b \pmod d$ is $e^{2\pi ai/d} \cdot e^{2\pi bi/d} = e^{2\pi(a+b)i/d}$.

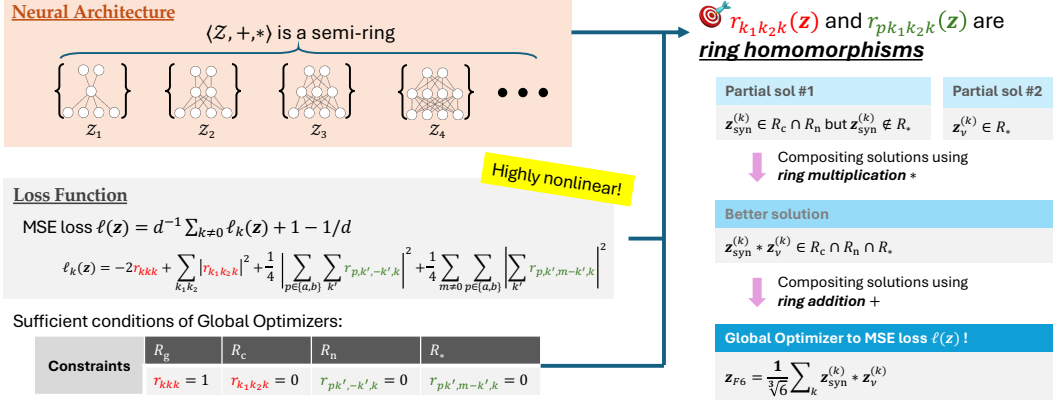


Figure 1: Overview of proposed theoretical framework CoGO. (1) The family of 2-layer neural networks, \mathcal{Z} , form a *semi-ring* algebraic structure (Theorem 2) with ring addition and multiplication (Def. 5). $\mathcal{Z} = \bigcup_{q \geq 0} \mathcal{Z}_q$ where \mathcal{Z}_q is a collection of all weights with order- q (i.e., q hidden nodes). (2) For outcome prediction of Abelian group multiplication, the MSE loss $\ell(\mathbf{z})$ is a function of *monomial potentials* (MPs) $r_{k_1 k_2 k}(\mathbf{z})$ and $r_{p k_1 k_2 k}(\mathbf{z})$ (Theorem 1), which are ring homomorphism (Theorem 3). (3) Thanks to the property of ring homomorphism, global optimizers to MSE loss $\ell(\mathbf{z})$ with quadratic activation can be constructed *algebraically* from partial solutions that only satisfy a subset of constraints (Sec. 5.1) using ring addition and multiplication, instead of running gradient descent. Examples include Fourier solution \mathbf{z}_{F6} (Corollary 2) and $\mathbf{z}_{FA/6}$ (Corollary 4) and perfect memorization solution \mathbf{z}_M (Corollary 3). In Sec. 6, we analyze the role played of MPs in gradient dynamics, showing that the dynamics favors low-order global optimizers (Theorem 6) under weight decay regularization, and the dynamics of MPs become decoupled with infinite width (Theorem 7). In Sec. 7 we show that the gradient descent solutions match exactly with our theoretical construction.

89 **Basic Ring theory.** A set \mathcal{Z} forms a *ring*, if there exists two operations, addition $+$ and multipli-
 90 cation $*$, so that (1) $\langle \mathcal{Z}, + \rangle$ forms an Abelian group, (2) $\langle \mathcal{Z}, * \rangle$ is a monoid (i.e., a group without
 91 inverse), and (3) multiplication distributes with addition (i.e., $a * (b + c) = a * b + a * c$ and
 92 $(b + c) * a = b * a + c * a$). \mathcal{Z} is called *semi-ring* if $\langle \mathcal{Z}, + \rangle$ is a monoid.

93 **Notation.** Let \mathbb{R} be the real field and \mathbb{C} be the complex field. For a complex vector \mathbf{z} , \mathbf{z}^\top is its
 94 transpose, $\bar{\mathbf{z}}$ is its complex conjugate and \mathbf{z}^* its conjugate transpose. For a tensor z_{ijk} , $z_{\cdot jk}$ is a
 95 vector along its first dimension, $z_{i \cdot k}$ along its second dimension, and $z_{ij \cdot}$ along its last dimension.

96 **Problem Setup.** We consider the following 2-layer networks with q hidden nodes, trained with
 97 (projected) ℓ_2 loss on prediction of group multiplication in Abelian group G with $|G| = d$:

$$\ell = \sum_i \left\| P_1^\perp \left(\frac{1}{2d} \mathbf{o}[i] - l[i] \right) \right\|^2, \quad \mathbf{o}[i] = V \sigma(W^\top \mathbf{f}[i]) = \sum_j \mathbf{v}_j \sigma(\mathbf{w}_j^\top \mathbf{f}[i]) \quad (1)$$

98 where $\sigma(x) = x^2$ is the quadratic activation function (Du & Lee, 2018; Allen-Zhu & Li, 2023),
 99 $P_1^\perp = I - \frac{1}{d} \mathbf{1}\mathbf{1}^\top$ is the zero-mean projection matrix, $W = [\mathbf{w}_1, \dots, \mathbf{w}_q] \in \mathbb{R}^{d \times q}$, $V =$
 100 $[\mathbf{v}_1, \dots, \mathbf{v}_q]^\top \in \mathbb{R}^{d \times q}$ are learnable parameters. $\mathbf{f}[i] \in \mathbb{R}^d$ are input embeddings. i is the sam-
 101 ple index. Note that variants of quadratic activation have been used empirically, e.g. squared ReLU
 102 and gated activations (So et al., 2021; Shazeer, 2020; Zhang et al., 2024).

103 **Input and Output.** The input contains the two group elements $g_1[i]$ and $g_2[i]$, encoded as $\mathbf{f}[i] =$
 104 $U_{G_1} \mathbf{e}_{g_1[i]} + U_{G_2} \mathbf{e}_{g_2[i]}$, where U_{G_1} and U_{G_2} are column orthogonal embedding matrices. The output
 105 is the result $g_1[i]g_2[i] \in G$, encoded as the label $l[i] = g_1[i]g_2[i]$ to be predicted. We can extend our
 106 framework to *group action prediction*, in which g_2 may not be a group element but any object (e.g.,
 107 a discrete state in reinforcement learning). See Appendix E for more details.

108 Let $\phi_k = [\phi_k(g)]_{g \in G} \in \mathbb{C}^d$ be the scaled Fourier bases (or more formally, *character function* of the
 109 finite Abelian group G , see Appendix A). Then weight vector \mathbf{w}_j and \mathbf{v}_j can be written as:

$$\mathbf{w}_j = U_{G_1} \sum_{k \neq 0} z_{akj} \phi_k + U_{G_2} \sum_{k \neq 0} z_{bkj} \phi_k, \quad \mathbf{v}_j = \sum_{k \neq 0} z_{ckj} \bar{\phi}_k \quad (2)$$

110 where $\mathbf{z} := \{z_{pkj}\}$ are the complex coefficients, $p \in \{a, b, c\}$, $0 \leq k < d$ and j runs through q
 111 hidden nodes. We exclude $\phi_0 \equiv 1$ because the constant bias term has been filtered out by the top-

112 down gradient from the loss function. Since \mathbf{w}_j and \mathbf{v}_j are all real, the Hermitian constraints holds,
 113 i.e., $\overline{z_{ckj}} = \overline{\phi_k^* \mathbf{v}_j} = \phi_{-k}^* \mathbf{v}_j = z_{c,-k,j}$ (and similar for z_{akj} and z_{bkj}). Leveraging the property of
 114 quadratic activation functions, we can write down the loss function analytically (see Appendix A):

115 **Theorem 1** (Analytic form of L_2 loss with quadratic activation). *The objective of 2-layer MLP*
 116 *network with quadratic activation can be written as $\ell = d^{-1} \sum_{k \neq 0} \ell_k + (d-1)/d$, where*

$$\ell_k = -2r_{kkk} + \sum_{k_1 k_2} |r_{k_1 k_2 k}|^2 + \frac{1}{4} \left| \sum_{p \in \{a,b\}} \sum_{k'} r_{p,k',-k',k} \right|^2 + \frac{1}{4} \sum_{m \neq 0} \sum_{p \in \{a,b\}} \left| \sum_{k'} r_{p,k',m-k',k} \right|^2 \quad (3)$$

117 Here $r_{k_1 k_2 k} := \sum_j z_{ak_1 j} z_{bk_2 j} z_{ckj}$ and $r_{pk_1 k_2 k} := \sum_j z_{pk_1 j} z_{pk_2 j} z_{ckj}$.

118 Note that for cyclic group G , the frequency k is a mod- d integer. For general Abelian group which
 119 can be decomposed into direct sum of cyclic groups according to Fundamental Theorem of Finite
 120 Abelian Groups (Diaconis, 1988), k is a multidimensional frequency index. For convenience, we
 121 define $\phi_{-k} := \overline{\phi_k}$ as the conjugate representation of ϕ_k . Since weights \mathbf{w}_j and \mathbf{v}_j are all real, the
 122 Hermitian constraints holds, i.e., $\overline{z_{ckj}} = \overline{\phi_k^* \mathbf{v}_j} = \phi_{-k}^* \mathbf{v}_j = z_{c,-k,j}$ (and similar for z_{akj} and z_{bkj}).
 123 Therefore, $z_{p,-k,j} = \overline{z_{pkj}}$, $r_{-k,-k,-k} = \overline{r_{kkk}}$ and ℓ is real and can be minimized.

124 Eqn. 3 contains different r terms, which play an important role in determining global optimizers.

125 **Definition 1** (0/1-set). *Let $R := \{r\}$ be a collection of r terms. The weight \mathbf{z} is said to have 0-set*
 126 *R_0 and 1-set R_1 (or 0/1-sets (R_0, R_1)), if $r(\mathbf{z}) = 0$ for all $r \in R_0$ and $r(\mathbf{z}) = 1$ for all $r \in R_1$.*

127 With 0/1-sets, we can characterize rough structures of the global optimizers to the loss:

128 **Lemma 1** (A Sufficient Conditions of Global optimizers of Eqn. 3). *If the weight \mathbf{z} to Eqn. 3 has*
 129 *0-sets $R_c \cup R_n \cup R_*$ and 1-set R_g , i.e.*

$$r_{kkk}(\mathbf{z}) = \mathbb{1}(k \neq 0), \quad r_{k_1 k_2 k}(\mathbf{z}) = 0, \quad r_{pk_1 k_2 k}(\mathbf{z}) = 0 \quad (4)$$

130 *then it is a global optimizer with zero loss $\ell(\mathbf{z}) = 0$. Here $R_g := \{r_{kkk}, k \neq 0\}$, $R_c :=$
 131 $\{r_{k_1 k_2 k}, k_1, k_2, k \text{ not all equal}\}$, $R_n := \{r_{p,k',-k',k}\}$ and $R_* := \{r_{p,k',m-k',k}, m \neq 0\}$.*

132 Lemma 1 provides sufficient conditions since there may exist solutions that achieve global optimum
 133 (e.g., $\sum_{k'} r_{p,k',m-k',k}(\mathbf{z}) = 0$ but $r_{p,k',m-k',k}(\mathbf{z}) \neq 0$). However, as we will see, it already leads
 134 to rich algebraic structure, and serves as a good starting point. Directly finding the global optimizers
 135 using Eqn. 4 can be a bit complicated and highly non-intuitive, due to highly nonlinear structure of
 136 Eqn. 3. However, there are nice structures we can leverage, as we will demonstrate below.

137 4 Beyond Fixed Parameter Space: The Semi-ring structure

138 4.1 The semi-ring structure of the solution space

139 We define the *weight space* $\mathcal{Z}_q = \{\mathbf{z}\}$ to include all the weight matrices with q hidden nodes
 140 (\mathcal{Z}_0 means an empty network), and $\mathcal{Z} = \bigcup_{q \geq 0} \mathcal{Z}_q$ be the solution space of all different number
 141 of hidden nodes. Interestingly, \mathcal{Z} naturally is equipped with a *semi-ring* structure, and each term
 142 of the loss function can effective interact with such a semi-ring structure, yielding provable global
 143 optimizers, including both the Fourier solutions empirically reported in previous works (Zhou et al.,
 144 2024; Gromov, 2023), and the perfect memorization solution (Morwani et al., 2023).

145 To make our argument formal, we start with a few definitions.

146 **Definition 2** (Order of \mathbf{z}). *The order $\text{ord}(\mathbf{z})$ of $\mathbf{z} \in \mathcal{Z}$ is its number of hidden nodes.*

147 **Definition 3** (Scalar multiplication). *$\alpha \mathbf{z} \in \mathcal{Z}$ is element-wise multiplication $[\alpha z_{pkj}]$ of $\mathbf{z} \in \mathcal{Z}$.*

148 **Definition 4** (Identification of \mathcal{Z}). *In \mathcal{Z} , two solutions of the same order that differ only by a per-
 149 mutation along hidden dimension j are considered identical.*

150 For any two solutions $\mathbf{z}_1 := \{z_{pkj}^{(1)}\}$ and $\mathbf{z}_2 := \{z_{pkj}^{(2)}\}$, we can define their operations:

151 **Definition 5** (Addition and Multiplication in \mathcal{Z}). *Define $\mathbf{z} = \mathbf{z}_1 + \mathbf{z}_2$ in which $z_{pk} :=$
 152 $\text{concat}(z_{pk}^{(1)}, z_{pk}^{(2)})$ and $\mathbf{z} = \mathbf{z}_1 * \mathbf{z}_2$, in which $z_{pk} := z_{pk}^{(1)} \otimes z_{pk}^{(2)}$. The addition and multiplication
 153 respect Hermitian constraints and the identity element $\mathbf{1}$ is the 1-order solutions with $\{z_{pk0} = 1\}$.*

154 Note that the multiplication definition is one special case of Khatri–Rao product (Khatri & Rao,
 155 1968). Although the Kronecker product and concatenation are not commutative, thanks to the iden-
 156 tification (Def. 4), it is clear that $z_1 + z_2 = z_2 + z_1$ and $z_1 * z_2 = z_2 * z_1$ and thus both operations
 157 are commutative. Then we can show:

158 **Theorem 2** (Algebraic Structure of \mathcal{Z}). $\langle \mathcal{Z}, +, * \rangle$ is a commutative semi-ring.

159 As we will see, the semi-ring structure of \mathcal{Z} paves the way to construct explicitly global optimizers.

160 4.2 The Monomial Potentials and its connection to semi-ring \mathcal{Z}

161 Now let us study the structure of the loss function Eqn. 3 and how they are related to the semi-ring
 162 structure of \mathcal{Z} . For this, we first define the concept of *monomial potentials*:

163 **Definition 6** (Monomial potential (MP)). $r(z) := \sum_j \prod_{p,k \in \text{id}_x(r)} z_{pkj}$ is called monomial poten-
 164 tial (MP), where $\text{id}_x(r)$ specifies the indices involved in the monomial terms.

165 Following this definition, terms in the loss function (Theorem 1) are examples of MPs.

166 **Observation 1** (Specific MPs). $r_{k_1 k_2 k}(z)$ and $r_{p k_1 k_2 k}(z)$ defined in Theorem 1 are MPs.

167 So what is the relationship between MPs, which are functions that map a weight z to a complex
 168 scalar, and the semi-ring structure of \mathcal{Z} ? The following theorem tells that MPs are *ring homomor-*
 169 *phism*, that is, these mappings respect addition and multiplication:

170 **Theorem 3.** For any monomial potential $r : \mathcal{Z} \mapsto \mathbb{C}$, $r(\mathbf{1}) = 1$, $r(z_1 + z_2) = r(z_1) + r(z_2)$ and
 171 $r(z_1 * z_2) = r(z_1)r(z_2)$ and thus r is a ring homomorphism.

172 **Observation 2.** The order function $\text{ord} : \mathcal{Z} \mapsto \mathbb{N}$ is also a ring homomorphism.

173 Since the loss function $\ell(z)$ depends on the weight z entirely through $r_{k_1 k_2 k}(z)$ and $r_{p k_1 k_2 k}(z)$,
 174 which are MPs, due to the property of ring homomorphism, it is possible to construct a global
 175 optimizer from partial solutions that satisfy only some of the constraints¹:

176 **Lemma 2** (Composing Partial Solutions). If z_1 has 0/1-sets (R_1^-, R_1^+) and z_2 has 0/1-sets
 177 (R_2^-, R_2^+) , then (1) $z_1 * z_2$ has 0/1-sets $(R_1^- \cup R_2^-, R_1^+ \cap R_2^+)$. (2) $z_1 + z_2$ have 0/1-sets
 178 $(R_1^- \cap R_2^-, R_1^+ \cup R_2^+)$.

179 Once we reach 0/1-sets $(R_c \cup R_n \cup R_*, R_g)$, we find a global optimizer. In addition, we also imme-
 180 diately know that there exists infinitely many global optimizers, via ring multiplication (Def. 5):

181 **Definition 7** (Unit). z is called a unit if $r_{kkk}(z) = 1$ for all $k \neq 0$.

182 **Corollary 1.** If z is a global optimizer and y is a unit, then $z * y$ is also a global optimizer.

183 5 Composing Global Optimizers

184 5.1 Constructing Partial Solutions with Polynomials

185 While intuitively one can get global optimizers by manually crafting some partial solutions and
 186 combining, in this section, we provide a more systematic approach to compose global optimizers as
 187 follows. Since \mathcal{Z} enjoys a semi-ring structure, we consider a *polynomial* in \mathcal{Z} in the following form:

$$z = \mathbf{u}^L + c_1 * \mathbf{u}^{L-1} + c_2 * \mathbf{u}^{L-2} + \dots + c_L \quad (5)$$

188 where the *generator* \mathbf{u} and coefficients c_l are order-1 and the power operation \mathbf{u}^l is defined by ring
 189 multiplication. The following construction of a polynomial leads to a partial solution.

190 **Theorem 4** (Construction of partial solutions). Suppose \mathbf{u} has 1-set R_1 , $\Omega_R(\mathbf{u}) := \{r(\mathbf{u}) | r \in$
 191 $R\} \subseteq \mathbb{C}$ is a set of evaluations on R (multiple values counted once), then if $1 \notin \Omega_R$, then the
 192 polynomial solution $\rho_R(\mathbf{u}) := \prod_{s \in \Omega_R(\mathbf{u})} (\mathbf{u} + \hat{s})$ has 0/1-set (R, R_1) up to a scale. Here \hat{s} is any
 193 order-1 weight that satisfies $r(\hat{s}) = -s$ for any $r \in R \cup R_+$. For example, $\hat{s} = -s^{1/3}\mathbf{1}$.

¹Mathematically, the *kernel* $\text{Ker}(r) := \{z : r(z) = 0\}$ of a ring homomorphism r is an *ideal* of the ring,
 and the intersection of ideals are still ideals. For brevity, we omit the formal definitions.

Symbol	$[a, b, c]$	Evaluation on MPs								Maximal polynomial $\rho(\mathbf{u})$	order q	
		R_c			R_n		R_*					
		$\bar{a}bc$	$a\bar{b}c$	$ab\bar{c}$	$\bar{a}ac$	$\bar{b}bc$	aac	bbc	$\bar{a}\bar{a}c$	$\bar{b}\bar{b}c$		
$\mathbf{1}_k$	$[1, 1, 1]$	1	1	1	1	1	1	1	1	1	–	–
$\bar{\mathbf{1}}_k$	$[-1, -1, 1]$	1	1	1	1	1	1	1	1	1	–	–
\mathbf{u}_{one}	$[1, -1, -1]$	1	1	1	-1	-1	-1	-1	-1	-1	$\mathbf{u} + \mathbf{1}$	2
\mathbf{u}_{syn}	$[\omega_3, \omega_3, \omega_3]$	ω_3	ω_3	ω_3	ω_3	ω_3	1	1	$\bar{\omega}_3$	$\bar{\omega}_3$	$\mathbf{u}^2 + \mathbf{u} + \mathbf{1}$	3
\mathbf{u}_{3c}	$[\omega_3, \bar{\omega}_3, 1]$	ω_3	$\bar{\omega}_3$	1	1	1	$\bar{\omega}_3$	ω_3	ω_3	$\bar{\omega}_3$	$\mathbf{u}^2 + \mathbf{u} + \mathbf{1}$	3
\mathbf{u}_{3a}	$[1, \omega_3, \bar{\omega}_3]$	1	ω_3	$\bar{\omega}_3$	$\bar{\omega}_3$	$\bar{\omega}_3$	$\bar{\omega}_3$	ω_3	$\bar{\omega}_3$	1	$\mathbf{u}^2 + \mathbf{u} + \mathbf{1}$	3
\mathbf{u}_{4c}	$[i, -i, 1]$	-1	-1	1	1	1	-1	-1	-1	-1	$\mathbf{u} + \mathbf{1}$	2
\mathbf{u}_{4a}	$[1, i, -i]$	1	-1	-1	-i	-i	-i	i	-i	i	$\mathbf{u}^3 + \mathbf{u}^2 + \mathbf{u} + \mathbf{1}$	4
\mathbf{u}_ν	$[\nu, -\nu, -\bar{\nu}^2]$	ν^2	ν^2	ν^4	$-\bar{\nu}^2$	$-\bar{\nu}^2$	-1	-1	$-\nu^4$	$-\nu^4$	9-th degree	10

Table 1: Exemplar order-1 single frequency generator $\mathbf{u}^{(k)}$ with $r_{kkk}(\mathbf{u}^{(k)}) = 1$. In the single-frequency case, for each MP r we use “ $\bar{a}bc$ ” to represent $r_{-k,k,k}$ and “ $\bar{a}\bar{a}c$ ” to represent $r_{a,-k,-k,k}$, etc. We omit superscript “ (k) ” for clarity and omit conjugate columns (i.e., $\bar{a}bc$ which is conjugate to $ab\bar{c}$). Here, $\omega_3 := e^{2\pi i/3}$ and $\omega_4 := i$ are the 3rd and 4th roots of unity. The constructed solutions are partial, i.e., the evaluation of some MPs yields 1 (red cell) and cannot be the root of the polynomial according to Theorem 4. Note that \mathbf{u}_ν is a general case with $\mathbf{u}_{\nu=1} = \mathbf{u}_{\text{one}}$ and $\mathbf{u}_{\nu=i} = \mathbf{u}_{4c}$.

194 For convenience, we use $\rho(\mathbf{u})$ to represent the *maximal* polynomial, i.e., when $R =$
 195 $\arg \max_{1 \notin \Omega_R(\mathbf{u})} |\Omega_R(\mathbf{u})|$ is the largest subset of MPs with $1 \notin \Omega_R(\mathbf{u})$. Our goal is to find low-order
 196 (partial) solutions, since gradient descent prefers low order solutions (see Theorem 6). Although
 197 there exist high-degree but low-order polynomials, e.g., $\mathbf{u}^9 + \mathbf{1}$, in general, degree L and order q are
 198 correlated, and we can find low-degree ones instead. To achieve that, \mathbf{u} should be properly selected
 199 (e.g., symmetric weights) to create as many duplicate values (but not 1) in R as possible.

200 5.2 Composing Global Solutions

201 We first consider the case that the generator \mathbf{u} is only nonzero at frequency k (and thus $-k$ by
 202 Hermitian constraints), but zero in other frequencies, i.e., $u_{pk'0} = 0$ for $k' \neq \pm k$. Such solutions
 203 correspond to Fourier bases in the original domain. Also, \mathbf{u} has 1-set $R_1 = \{r_{kkk}\}$. This means that
 204 \mathbf{u} can be characterized by three numbers $u_{ak0} = a$, $u_{bk0} = b$, and $u_{ck0} = c$ with $abc = 1$. In this
 205 case, only a subset of monomial potentials (MPs) whose indices only involve a single frequency k
 206 are non-zero (e.g., $r_{k,-k,k} \in R_c$ and $r_{b,-k,k,k} \in R_n$), which makes our construction much easier.

207 Following Theorem 4, we can construct different partial solutions. Some examples are shown in
 208 Table 1, which do not reach the complete set $R_c \cup R_n \cup R_*$ and therefore are not global. Note that
 209 it is possible to create a generator so that all MPs are not 1 (e.g., $\mathbf{u}_{3c} * \mathbf{u}_{4a}$), but then $|\Omega_R(\mathbf{u})|$ will
 210 be too large, producing high-degree polynomials (e.g., $\mathbf{u}_{3c} * \mathbf{u}_{4a}$ gives a 10-th-degree polynomial).

211 However, utilizing these partial solutions, with Lemma 2 we can construct global optimizers:

212 **Corollary 2** (Order-6 global optimizers). *The following “ 3×2 ” Fourier solutions satisfy the suffi-*
 213 *cient condition (Lemma 1) and thus are global optimizers (assuming d is odd):*

$$z_{F6} = \frac{1}{\sqrt[3]{6}} \sum_{k=1}^{(d-1)/2} z_{\text{syn}}^{(k)} * z_\nu^{(k)} * \mathbf{y}_k \quad (6)$$

214 Here $z_{\text{syn}}^{(k)} := \rho(\mathbf{u}_{\text{syn}}^{(k)})$ and $z_\nu^{(k)} := \mathbf{u}_\nu^{(k)} + \mathbf{1}_k$ (i.e., not maximal polynomial), where \mathbf{u}_{syn} and \mathbf{u}_ν
 215 are defined in Table 1. \mathbf{y} is an order-1 unit. As a result, $\text{ord}(z_{F6}) = 3 \cdot 2 \cdot 1 \cdot (d-1)/2 = 3(d-1)$
 216 and each frequency are affiliated with 6 hidden nodes (order-6).

217 **Other solutions.** We may replace \mathbf{u}_{syn} and \mathbf{u}_ν with any other pairs that collectively cover all MPs.
 218 For example, \mathbf{u}_{syn} can be combined with any of $\{\mathbf{u}_{3c}, \mathbf{u}_{3a}, \mathbf{u}_{4a}\}$, and $\mathbf{u}_{\nu=\pm i}$ can be coupled with
 219 \mathbf{u}_{3a} or \mathbf{u}_{4a} , etc. Here we pick one with a small order. Compared to construction from Gromov
 220 (2023), ours is much more concise and does not use infinite-width approximation.

221 **Even d .** For even d , simply replace $(d-1)/2$ with $\lfloor (d-1)/2 \rfloor$ and add an additional order-2 term
 222 $\rho(\mathbf{u}_{\text{one}}) = \mathbf{u}_{\text{one}} + \mathbf{1}$ (Tbl. 1) for the frequency $d/2$. Note that the frequency $k = d/2$ only has r_{kkk} ,
 223 r_{akkk} and r_{bkkk} , and all other conjugate combinations are absent. Thus $\mathbf{u}_{\text{one}}^{(k)} + \mathbf{1}_k$ covers them all.

224 Fig. 2 shows a case with $d = 7$. In this case, each frequency, out of $(d-1)/2 = 3$ total number of
 225 frequencies, is associated with 6 hidden nodes. If we remove the last term in the loss that corresponds
 226 to R_* , then an order-3 solution suffices (i.e. $z_{\text{syn}} = \rho(\mathbf{u}_{\text{syn}})$).

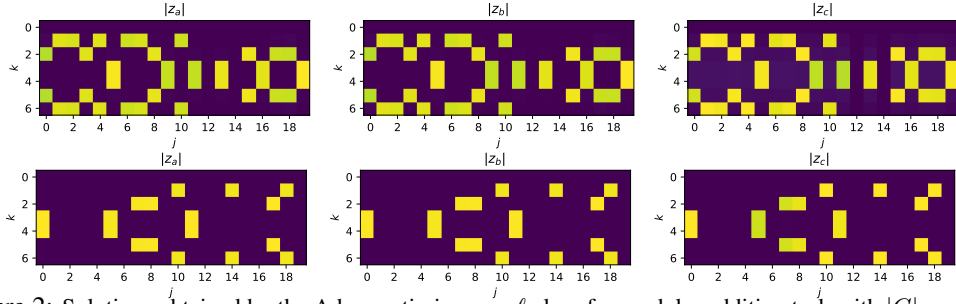


Figure 2: Solutions obtained by the Adam optimizers on ℓ_2 loss for modular addition task with $|G| = d = 7$ and $q = 20$ hidden nodes. **Top:** For each frequency $\pm k$, there are exactly 6 hidden nodes representing such a frequency, consistent with Corollary 1. **Bottom:** Optimizing Eqn. 3 without the last term $\sum_{m \neq 0} \sum_{p \in \{a, b\}} \left| \sum_{k'} r_{p, k', m-k', k} \right|^2$ (equivalently removing the constraint R_*). Now each frequency has exactly 3 hidden nodes, which corresponds to the solution $\mathbf{z}_{\text{syn}} = \boldsymbol{\rho}(\mathbf{u}_{\text{syn}})$ in Tbl. 1.

227 Using polynomials, we can also construct perfect-memorization solutions. For this, we first define
 228 two generators \mathbf{u}_a with $u_{\cdot, k0}^{(\alpha)} = [\omega_d^k, 1, \bar{\omega}_d^k]$ ($k \neq 0$), and \mathbf{u}_b with $u_{\cdot, k0}^{(\beta)} = [1, \omega_d^k, \bar{\omega}_d^k]$ ($k \neq 0$). Here
 229 $\omega_d := e^{2\pi i/d}$ is the d -th root of unity.

230 **Corollary 3** (Perfect Memorization). *We construct two d -order weights \mathbf{z}_a and \mathbf{z}_b :*

$$\mathbf{z}_a = \sum_{j=0}^{d-1} \mathbf{u}_a^j, \quad \mathbf{z}_b = \sum_{j=0}^{d-1} \mathbf{u}_b^j \quad (7)$$

231 Here $\mathbf{z}_a \in R_c(k_1 \neq k) \cap R_n \cap R_*(p = b \text{ or } m \neq k)$, $\mathbf{z}_b \in R_c(k_2 \neq k) \cap R_n \cap R_*(p = a \text{ or } m \neq k)$.
 232 Then $\mathbf{z}_M = d^{-2/3} \mathbf{z}_a * \mathbf{z}_b$ satisfies the sufficient condition (Lemma 1) and is the perfect memorization
 233 solution with $\text{ord}(\mathbf{z}_M) = d^2$:

$$z_{akj_1j_2}^{(M)} = \omega^{kj_1} / \sqrt[3]{d^2}, \quad z_{bkj_1j_2}^{(M)} = \omega^{kj_2} / \sqrt[3]{d^2}, \quad z_{ckj_1j_2}^{(M)} = \omega^{-k(j_1+j_2)} / \sqrt[3]{d^2} \quad (8)$$

234 where each hidden node is indexed by $j = (j_1, j_2)$, $0 \leq j_1, j_2 < d$, $k \neq 0$.

235 To see why this corresponds to perfect memorization, simply apply an inverse Fourier transform for
 236 each hidden node (j_1, j_2) , and the original weights are (zero-mean) delta function located at j_1 , j_2
 237 and $j_1 + j_2$ accordingly.

238 Interestingly, there also exists a lower-order solution, 2×2 , that meets R_c and R_* but not R_n :

239 **Corollary 4** (Order-4 single frequency solution). *Define single frequency order-2 solution \mathbf{z}_ξ :*

$$z_{ak\cdot} = [1, \xi], \quad z_{bk\cdot} = [1, -i\xi], \quad z_{ck\cdot} = [1, i] \quad (9)$$

240 where $|\xi| = 1$. Then the order-4 solution $\mathbf{z}_{F4}^{(k)} := \boldsymbol{\rho}(\mathbf{u}_{\nu=i}^{(k)}) * \mathbf{z}_\xi^{(k)}$ has 0-sets R_c and R_* (but not R_n).

241 While $\mathbf{z}_{F4}^{(k)}$ itself does not satisfy the sufficient condition (Eqn. 4), it is part of a global optimizer
 242 when mixing with \mathbf{z}_{F6} :

243 **Corollary 5** (Mixed order-4/6 global optimizers). *With $\mathbf{z}_{F4}^{(k)}$, there is a global optimizer to Eqn. 3
 244 that does not meet the sufficient condition, i.e., $\sum_{k'} r_{p, k', -k', m} = 0$ but $r_{p, k', -k', m} \neq 0$:*

$$\mathbf{z}_{F4/6} = \frac{1}{\sqrt[3]{6}} \hat{\mathbf{z}}_{F6}^{(k_0)} + \frac{1}{\sqrt[3]{4}} \sum_{k=1, k \neq k_0}^{(d-1)/2} \mathbf{z}_{F4}^{(k)} \quad (10)$$

245 where $\hat{\mathbf{z}}_{F6}^{(k_0)}$ is a perturbation of $\mathbf{z}_{F6}^{(k_0)} := \mathbf{z}_{\text{syn}}^{(k_0)} * \mathbf{z}_{\nu=1}^{(k_0)}$ by adding constant biases to its (c, k) entries
 246 for $k \neq k_0$. The order is lower than \mathbf{z}_{F6} : $\text{ord}(\mathbf{z}_{F4/6}) = 6 + 4 \cdot ((d-1)/2 - 1) = 2d < \text{ord}(\mathbf{z}_{F6})$.

247 **Remarks.** To construct $\hat{\mathbf{z}}_{F6}$, in addition to $\mathbf{z}_{\text{syn}} * \mathbf{z}_{\nu=1}$, we could use other pairs of single frequency
 248 solutions to achieve the same effects. For example, using $\mathbf{z}_{\text{syn}, \alpha\beta} * \mathbf{z}_{\nu=i}$, where $\mathbf{z}_{\text{syn}, \alpha\beta}$ is:

$$z_{ak\cdot} = [1, \omega_3\alpha, \bar{\omega}_3\beta], \quad z_{bk\cdot} = [1, \omega_3\bar{\alpha}, \bar{\omega}_3\bar{\beta}], \quad z_{ck\cdot} = [1, \omega_3, \bar{\omega}_3] \quad (11)$$

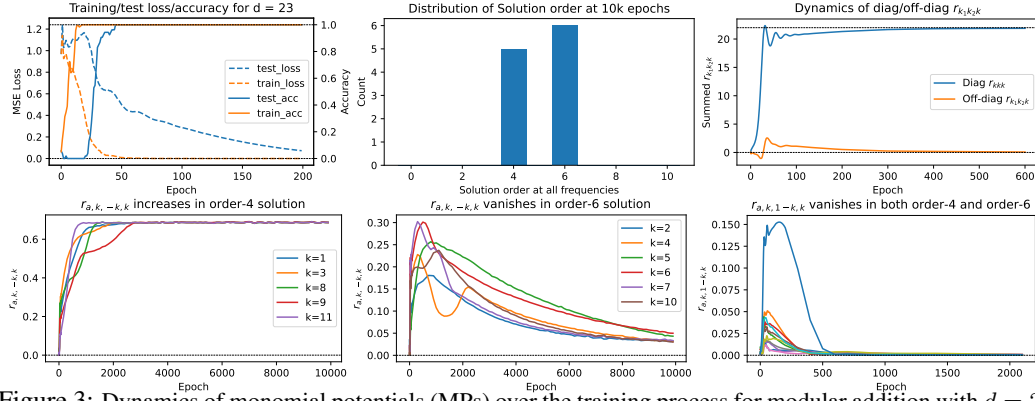


Figure 3: Dynamics of monomial potentials (MPs) over the training process for modular addition with $d = 23$ and $q = 1024$ hidden nodes. **Top Row.** *Left:* Training/test accuracy reaches 100% and loss close to 0. Test accuracy jumps after training reaches 100% (grokking). *Mid:* After 10k epochs, the distribution of solution orders are concentrated at 4 and 6 (Corollary 2 and 4). *Right:* Dynamics of $r_{k_1 k_2 k}$. Summation of diagonal r_{kkk} converges towards $d - 1$ (dotted line) with ripple effects, while off-diagonal $r_{k_1 k_2 k}$ converges towards 0. **Bottom Row.** Dynamics of different MPs. Order-4 and order-6 behave differently on $r_{p,k,-k,k}$, because order-4 does not satisfy the sufficient condition (Lemma 1) but a mixture of order-4 and order-6 (i.e., $\mathcal{Z}_{F4/6}$) is still the global optimizer to the L_2 loss (Corollary 5).

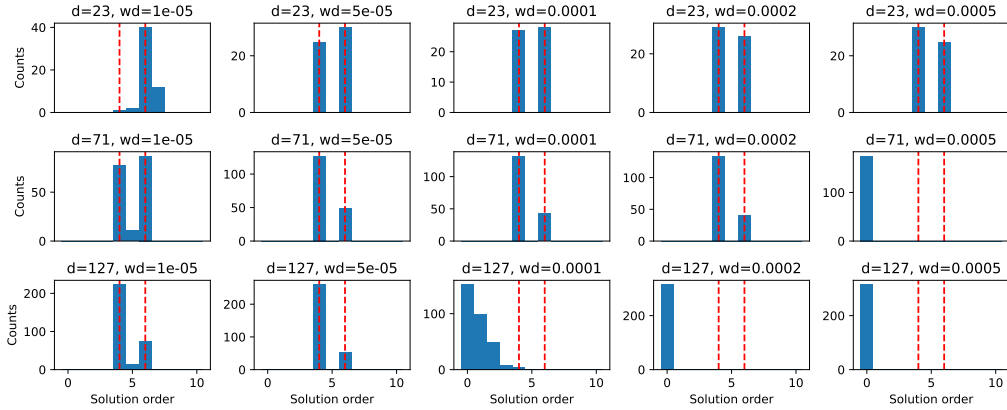


Figure 4: Solution distribution over different weight decay regularization for $q = 512$, trained with 10k epochs with Adams with learning rate 0.01 on modular addition (i.e., predicting $a + b \pmod{d}$) with $d \in \{23, 71, 127\}$. The two red dashed lines correspond to order-4/6 solutions. The histogram is accumulated over 5 random seeds.

249 where $|\alpha| = |\beta| = 1$. $\mathbf{z}_{\text{syn}} = \boldsymbol{\rho}(\mathbf{u}_{\text{syn}})$ is a special case of $\mathbf{z}_{\text{syn},\alpha,\beta}$ when $\alpha = \beta = 1$.

250 Note that multiple per frequency order-6 solutions can be inserted in this construction. Compared
 251 to all order-6 solutions \mathcal{Z}_{F6} , this $\mathcal{Z}_{F4/6}$ mixture solution has a lower order and is perceived in the
 252 experiments (See Fig. 6), in particular when d is large (Tbl. 2), showing a strong preference of
 253 gradient descent towards lower order solutions.

254 6 Gradient dynamics

255 Now we have characterized the structures of global optimizers. One natural question arises: why
 256 does the optimization procedure not converge to the perfect memorization solution \mathbf{z}_M , but to the
 257 Fourier solutions \mathcal{Z}_{F6} and $\mathcal{Z}_{F4/6}$? The answer is given by gradient dynamics.

258 Let $\mathbf{r} = [r_{k_1 k_2 k}, r_{p k_1 k_2 k}] \in \mathbb{C}^{4d^3}$ be a vector of all MPs, and $J := \frac{\partial \mathbf{r}}{\partial \mathbf{z}} \frac{\partial \mathbf{z}}{\partial \mathcal{W}}$ be the Jacobian matrix
 259 of the mapping $\mathbf{r} = \mathbf{r}(\mathbf{z}(\mathcal{W}))$ in which \mathcal{W} is the collection of original weights. Note that when we
 260 take derivatives with respect to \mathbf{r} and apply chain rules, we treat \mathbf{r} and its complex conjugate (e.g.,
 261 r_{kkk} and $r_{-k,-k,-k} = \bar{r}_{kkk}$) as independent variables. Since we run the gradient descent on \mathcal{W} , will
 262 such (indirect) optimization leads to a descent of \mathbf{r} towards the desired targets (Lemma 1)? This is
 263 confirmed by the following theorem:

264 **Theorem 5** (Dynamics of MPs). *The dynamics of MPs satisfies $\dot{\mathbf{r}} = -JJ^* \overline{\nabla_{\mathbf{r}} \ell}$, which has positive
 265 inner product with the negative gradient direction $-\overline{\nabla_{\mathbf{r}} \ell}$.*

d	%not			%non-factorable		error ($\times 10^{-2}$)		solution distribution (%) in factorable ones			
	order-4/6	order-4	order-6	order-4	order-6	$\mathbf{z}_{\nu=i}^{(k)} * \mathbf{z}_{\xi}^{(k)}$	$\mathbf{z}_{\nu=i}^{(k)} * \mathbf{z}_{\text{syn},\alpha\beta}^{(k)}$	$\mathbf{z}_{\nu}^{(k)} * \mathbf{z}_{\text{syn}}^{(k)}$	others		
23	0.0 \pm 0.0	0.00 \pm 0.00	5.71 \pm 5.71	0.05 \pm 0.01	4.80 \pm 0.96	47.07 \pm 1.88	11.31 \pm 1.76	39.80 \pm 2.11	1.82 \pm 1.82		
71	0.0 \pm 0.0	0.00 \pm 0.00	0.00 \pm 0.00	0.03 \pm 0.00	5.02 \pm 0.25	72.57 \pm 0.70	4.00 \pm 1.14	21.14 \pm 2.14	2.29 \pm 1.07		
127	0.0 \pm 0.0	1.50 \pm 0.92	0.00 \pm 0.00	0.26 \pm 0.14	0.93 \pm 0.18	82.96 \pm 0.39	2.25 \pm 0.64	14.13 \pm 0.87	0.66 \pm 0.66		

Table 2: Matches between order-4/6 solutions from gradient descent and those constructed by CoGO. Number of hidden nodes $q = 512$ and weight decay is 5×10^{-5} . Around 95% gradient descent solutions are factorable with very small factorization error (~ 0.04 compared to solution norm on the order of 1). Furthermore, CoGO successfully predicts $\sim 98\%$ of the structure of the empirical solutions, while the remaining 2% are mostly due to insufficient training (i.e., near miss against known theoretical construction). Here \mathbf{z}_{ξ} is defined in Corollary 4, $\mathbf{z}_{\nu} := \mathbf{u}_{\nu} + \mathbf{1}$ is defined in Tbl. 1, and $\mathbf{z}_{\text{syn},\alpha\beta}$ is defined in Eqn. 11. The means and their standard deviations are computed over 5 seeds.

266 Corollary 1 shows that by ring multiplication, we could create infinitely many global optima from a
 267 base one. The following theorem answers which solution gradient dynamics picks.

268 **Theorem 6** (The Occam’s Razer: Preference of low-order solutions). *If $\mathbf{z} = \mathbf{y} * \mathbf{z}'$ and both \mathbf{z} (of
 269 order q) and \mathbf{z}' are global optimal solutions, then there exists a path of zero loss connecting \mathbf{z} and \mathbf{z}'
 270 in the space of \mathcal{Z}_q . As a result, lower-order solutions are preferred if trained with L_2 regularization.*

271 This shows that gradient dynamics with weight decay will pick a lower-order (i.e., simpler) solution,
 272 suggests that perfect memorization may not be not favorable in dynamics. The following theorem
 273 shows that the dynamics also enjoys *asymptotic freedom*:

274 **Theorem 7** (Infinite Width Limits at Initialization). *Considering the modified loss of Eqn. 3 with
 275 only the first two terms: $\tilde{\ell}_k := -2r_{kkk} + \sum_{k_1 k_2} |r_{k_1 k_2 k}|^2$, if the weights are i.i.d Gaussian and
 276 network width $q \rightarrow +\infty$, then JJ^* converge to diagonal and the dynamics of MPs is decoupled.*

277 Intuitively, this means that a large enough network width ($q \rightarrow +\infty$) makes the dynamics much
 278 easier to analyze. On the other hand, the final solution may not require that large q . As analyzed in
 279 Corollary 2, for each frequency, to achieve global optimality, 6 hidden nodes suffice.

280 7 Experiments

281 **Setup.** We train the 2-layer MLP on the modular addition task, which is a special case of outcome
 282 prediction of Abelian group multiplication. We use Adam optimizer with learning rate 0.01, MSE
 283 loss, and train for 10000 epochs with weight decays. We tested on $|G| = d \in \{23, 71, 127\}$. All
 284 data are generated synthetically and training/test split is 90%/10%.

285 **Solution Distributions.** As shown in Fig. 3, we see order-4 and order-6 solutions in each frequency
 286 emerging from well-trained networks on $d = 23$. The mixed solution $\mathbf{z}_{F_{4/6}}$ can be clearly observed
 287 in a small-scale example (Fig. 6). This is also true for larger d (Fig. 4). Although the model is
 288 trained with heavily over-parameterized networks, the final solution order remains constant, which
 289 is consistent with Corollary 1. Large weight decay shifts the distribution to the left (i.e., low-order
 290 solutions) until model collapses (i.e., all weights become zero), consistent with our Theorem 6 that
 291 demonstrates that gradient descent with weight decay favors low-order solutions. Similar conclu-
 292 sions follow for fewer and more overparameterization (Appendix H).

293 **Exact match between theoretical construction and empirical solutions.** A follow-up question
 294 arises: *do the empirical solutions match exactly with our constructions?* After all, distribution
 295 of solution order is a rough metric. For this, we identify all solutions obtained by gradient de-
 296 scent at each frequency, factorize them and compare with theoretical construction up to conjuga-
 297 tion/normalization. To find such a factorization, we use exhaustive search (Appendix H).

298 The answer is yes. Tbl. 2 shows that around 95% of order-4 and order-6 solutions from gradient
 299 descent can be factorized into 2×2 and 2×3 and each component matches our theoretical construc-
 300 tion in Corollary 2 and 4, with minor variations. Furthermore, when d is large, most of the solutions
 301 become order-4, which is consistent with our analysis for mixed solution $\mathbf{z}_{F_{4/6}}$ (Corollary 5) that
 302 one order-6 solution in the form of $\mathbf{z}_{\nu=i} * \mathbf{z}_{\text{syn},\alpha\beta}$ suffices to achieve a global optimizer, with all
 303 other frequencies taking order-4s. In fact, for $d = 127$, the number of order-6 solution taking the
 304 form of $\mathbf{z}_{\nu=i} * \mathbf{z}_{\text{syn},\alpha\beta}$ is $(d - 1)/2 \cdot 2.25\% \approx 1.26$, coinciding with the theoretical results.

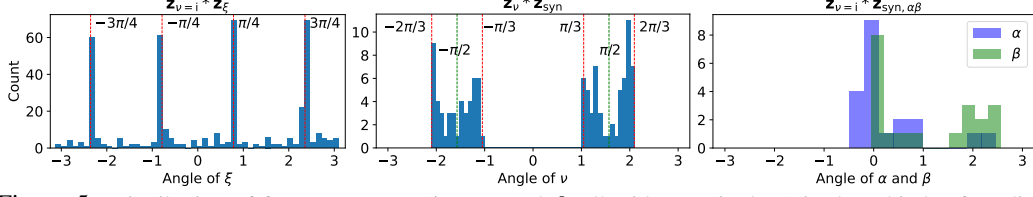


Figure 5: Distribution of free parameters (ξ , ν , α and β , all with magnitude 1) in three kinds of gradient descent solutions identified by CoGO. While any value of these parameters makes a global optimizer, gradient descent dynamics has a particular preference in picking them during optimization.

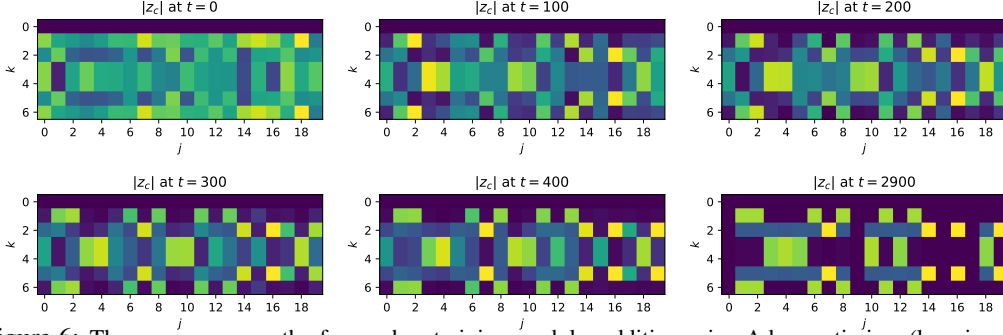


Figure 6: The convergence path of $z_{c..}$ when training modular addition using Adam optimizer (learning rate 0.05, weight decay 0.005). The final solution contains 2 order-6 ($z_{F6}^{(k)}$) and 1 order-4 ($z_{F4}^{(k)}$) solutions. Note that for $z_{c..}$, unlike Fig. 2, each order-6 solution contains a constant bias term to cancel out the artifacts of order-4 solution (Corollary 5). For each hidden node j , once a dominant frequency emerges, others fade away.

305 **Implicit Bias of gradient descent.** Our construction gives other possible solutions (e.g., $z_{3c} * z_{\text{syn}}$)
 306 which are never observed in the gradient solutions. Even for the observed solutions, e.g. $z_{\nu} * z_{\text{syn}}$,
 307 the distribution of free parameters is highly non-uniform (Fig. 5), showing a strong preference of
 308 certain choices. These suggest strong implicit bias in optimization, which we leave for future work.

309 8 Conclusion and future work

310 In this work, we propose CoGO (*Composing Global Optimizers*), a theoretical framework that mod-
 311 els the algebraic structure of global optimizers when training a 2-layer network on reasoning tasks
 312 of Abelian group with L_2 loss. We find that the global optimizers can be algebraically composed by
 313 partial solutions that only fit parts of the loss, using ring operations defined in the weight space of the
 314 2-layer neural networks across different network widths. Under CoGO, we also analyze the training
 315 dynamics, show the benefit of over-parameterization, and the inductive bias towards simpler solu-
 316 tions due to topological connectivity between algebraically linked high-order (i.e., involving more
 317 hidden nodes) and low-order global optimizers. Finally, we show that the gradient descent solutions
 318 exactly match what constructed solutions (e.g. $z_{F4/6}$ and z_{F6} , see Corollary 5 and Corollary 2).

319 **Develop novel training algorithms.** Instead of applying (stochastic) gradient descent to overpa-
 320 rameterized networks, CoGO suggests a completely different path: decompose the loss, find the
 321 MPs, construct low-order solutions and combine them to achieve the final solutions on the fly using
 322 algebraic operations. Such an approach may be more efficient and scalable than gradient descent,
 323 due to its factorable nature. Also, our framework works for losses depending on monomial potentials
 324 (L_2 loss is just one example), which opens a new dimension for loss design.

325 **Putting different widths into the same framework.** Many existing theoretical works study prop-
 326 erties of networks with fixed width. However, CoGO demonstrates that nice mathematical structures
 327 emerge when putting networks of different widths together, an interesting direction to consider.

328 **Grokking.** When learning modular addition, there exists a phase transition from *memorization* to
 329 *generalization* during training, known as *grokking* (Varma et al., 2023; Power et al., 2022), long after
 330 the training performance becomes (almost) perfect. Our work may be expanded to a nonuniformly
 331 distributed training set to study the dynamics of representation learning on grokking.

332 **Extending to other activations.** For other activation than quadratic (e.g., SiLU) with $\sigma(0) = 0$,
 333 with a Taylor expansion, the same framework may still apply (with higher rank MPs).

334 **References**

- 335 Zeyuan Allen-Zhu and Yuanzhi Li. Backward feature correction: How deep learning performs deep
336 (hierarchical) learning. In *The Thirty Sixth Annual Conference on Learning Theory*, pp. 4598–
337 4598. PMLR, 2023.
- 338 Anthropic. The claude 3 model family: Opus, sonnet, haiku. URL <https://www.anthropic.com/news/claude-3-family>.
339
- 340 Lukas Berglund, Meg Tong, Max Kaufmann, Mikita Balesni, Asa Cooper Stickland, Tomasz Kor-
341 bak, and Owain Evans. The reversal curse: Llms trained on” a is b” fail to learn” b is a”. *arXiv*
342 *preprint arXiv:2309.12288*, 2023.
- 343 Michael M Bronstein, Joan Bruna, Taco Cohen, and Petar Veličković. Geometric deep learning:
344 Grids, groups, graphs, geodesics, and gauges. *arXiv preprint arXiv:2104.13478*, 2021.
- 345 Keith Conrad. Characters of finite abelian groups. *Lecture Notes*, 17, 2010.
- 346 Persi Diaconis. Group representations in probability and statistics. *Lecture notes-monograph series*,
347 11:i–192, 1988.
- 348 Simon Du and Jason Lee. On the power of over-parametrization in neural networks with quadratic
349 activation. In *International conference on machine learning*, pp. 1329–1338. PMLR, 2018.
- 350 Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha
351 Letman, Akhil Mathur, Alan Schelten, and et. al. The llama 3 herd of models, 2024. URL
352 <https://arxiv.org/abs/2407.21783>.
- 353 Nouha Dziri, Ximing Lu, Melanie Sclar, Xiang Lorraine Li, Liwei Jiang, Bill Yuchen Lin, Peter
354 West, Chandra Bhagavatula, Ronan Le Bras, Jena D Hwang, et al. Faith and fate: Limits of
355 transformers on compositionality (2023). *arXiv preprint arXiv:2305.18654*, 2023.
- 356 William Fulton and Joe Harris. *Representation theory: a first course*, volume 129. Springer Science
357 & Business Media, 2013.
- 358 Quentin Garrido, Mahmoud Assran, Nicolas Ballas, Adrien Bardes, Laurent Najman, and Yann
359 LeCun. Learning and leveraging world models in visual representation learning. *arXiv preprint*
360 *arXiv:2403.00504*, 2024.
- 361 Andrey Gromov. Grokking modular arithmetic. *arXiv preprint arXiv:2301.02679*, 2023.
- 362 Jie Huang, Xinyun Chen, Swaroop Mishra, Huaixiu Steven Zheng, Adams Wei Yu, Xinying Song,
363 and Denny Zhou. Large language models cannot self-correct reasoning yet. *arXiv preprint*
364 *arXiv:2310.01798*, 2023.
- 365 Yufei Huang, Shengding Hu, Xu Han, Zhiyuan Liu, and Maosong Sun. Unified view of grokking,
366 double descent and emergent abilities: A perspective from circuits competition. *arXiv preprint*
367 *arXiv:2402.15175*, 2024.
- 368 Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chap-
369 lot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier,
370 L  lio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril,
371 Thomas Wang, Timoth  e Lacroix, and William El Sayed. Mistral 7b, 2023. URL <https://arxiv.org/abs/2310.06825>.
372
- 373 Charles Jin and Martin Rinard. Emergent representations of program semantics in language models
374 trained on programs, 2024.
- 375 Subbarao Kambhampati, Karthik Valmeekam, Lin Guan, Mudit Verma, Kaya Stechly, Siddhant
376 Bhambri, Lucas Saldyt, and Anil Murthy. Llms can’t plan, but can help planning in llm-modulo
377 frameworks, 2024. URL <https://arxiv.org/abs/2402.01817>.
- 378 CG Khatri and C Radhakrishna Rao. Solutions to some functional equations and their applications
379 to characterization of probability distributions. *Sankhy  : the Indian journal of statistics, series A*,
380 pp. 167–180, 1968.

- 381 Zhiyuan Li, Hong Liu, Denny Zhou, and Tengyu Ma. Chain of thought empowers transformers to
382 solve inherently serial problems. *ICLR*, 2024.
- 383 Bingbin Liu, Jordan T Ash, Surbhi Goel, Akshay Krishnamurthy, and Cyril Zhang. Transformers
384 learn shortcuts to automata. *arXiv preprint arXiv:2210.10749*, 2022.
- 385 Depen Morwani, Benjamin L Edelman, Costin-Andrei Oncescu, Rosie Zhao, and Sham Kakade.
386 Feature emergence via margin maximization: case studies in algebraic tasks. *arXiv preprint*
387 *arXiv:2311.07568*, 2023.
- 388 Neel Nanda, Lawrence Chan, Tom Lieberum, Jess Smith, and Jacob Steinhardt. Progress measures
389 for grokking via mechanistic interpretability. In *The Eleventh International Conference on Learn-*
390 *ing Representations*, 2023. URL <https://openreview.net/forum?id=9XFSbDPmdW>.
- 391 Marianna Nezhurina, Lucia Cipolina-Kun, Mehdi Cherti, and Jenia Jitsev. Alice in wonderland:
392 Simple tasks showing complete reasoning breakdown in state-of-the-art large language models,
393 2024. URL <https://arxiv.org/abs/2406.02061>.
- 394 OpenAI. Gpt-4 technical report, 2024. URL <https://arxiv.org/abs/2303.08774>.
- 395 Simon Ouellette, Rolf Pfister, and Hansueli Jud. Counting and algorithmic generalization with
396 transformers. *arXiv preprint arXiv:2310.08661*, 2023.
- 397 Alethea Power, Yuri Burda, Harri Edwards, Igor Babuschkin, and Vedant Misra. Grokking: Gen-
398 eralization beyond overfitting on small algorithmic datasets. *arXiv preprint arXiv:2201.02177*,
399 2022.
- 400 Noam Shazeer. Glu variants improve transformer. *arXiv preprint arXiv:2002.05202*, 2020.
- 401 David R. So, Wojciech Manke, Hanxiao Liu, Zihang Dai, Noam Shazeer, and Quoc V. Le. Primer:
402 Searching for efficient transformers for language modeling. *NeurIPS*, 2021. URL [https://](https://arxiv.org/abs/2109.08668)
403 arxiv.org/abs/2109.08668.
- 404 Benjamin Steinberg. Representation theory of finite groups. *Carleton University*, 2009.
- 405 Richard S Sutton. Reinforcement learning: An introduction. *A Bradford Book*, 2018.
- 406 DeepSeek Team. Deepseek-v2: A strong, economical, and efficient mixture-of-experts language
407 model, 2024a. URL <https://arxiv.org/abs/2405.04434>.
- 408 Gemini Team. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of con-
409 text, 2024b. URL <https://arxiv.org/abs/2403.05530>.
- 410 Vikrant Varma, Rohin Shah, Zachary Kenton, János Kramár, and Ramana Kumar. Explaining
411 grokking through circuit efficiency. *arXiv preprint arXiv:2309.02390*, 2023.
- 412 Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yo-
413 gatama, Maarten Bosma, Denny Zhou, Donald Metzler, et al. Emergent abilities of large language
414 models. *TMLR*, 2022.
- 415 Erik Wijmans, Manolis Savva, Irfan Essa, Stefan Lee, Ari S Morcos, and Dhruv Batra. Emergence
416 of maps in the memories of blind navigation agents. *AI Matters*, 9(2):8–14, 2023.
- 417 Jian Xie, Kai Zhang, Jiangjie Chen, Tinghui Zhu, Renze Lou, Yuandong Tian, Yanghua Xiao, and
418 Yu Su. Travelplanner: A benchmark for real-world planning with language agents, 2024.
- 419 Tian Ye, Zicheng Xu, Yuanzhi Li, and Zeyuan Allen-Zhu. Physics of language models: Part 2.1,
420 grade-school math and the hidden reasoning process. *arXiv preprint arXiv:2407.20311*, 2024.
- 421 Gilad Yehudai, Haim Kaplan, Asma Ghandeharioun, Mor Geva, and Amir Globerson. When can
422 transformers count to n? *arXiv preprint arXiv:2407.15160*, 2024.
- 423 Zhengyan Zhang, Yixin Song, Guanghui Yu, Xu Han, Yankai Lin, Chaojun Xiao, Chenyang Song,
424 Zhiyuan Liu, Zeyu Mi, and Maosong Sun. Relu² wins: Discovering efficient activation functions
425 for sparse llms, 2024. URL <https://arxiv.org/abs/2402.03804>.

- 426 Ziqian Zhong, Ziming Liu, Max Tegmark, and Jacob Andreas. The clock and the pizza: Two
427 stories in mechanistic explanation of neural networks. *Advances in Neural Information Processing*
428 *Systems*, 36, 2024.
- 429 Tianyi Zhou, Deqing Fu, Vatsal Sharan, and Robin Jia. Pre-trained large language models use fourier
430 features to compute addition. *arXiv preprint arXiv:2406.03445*, 2024.

431 **A Decoupling L_2 Loss (Proof)**

432 We use the *character function* $\phi : G \rightarrow \mathbb{C}$, which maps a group element g into a complex number.

433 **Lemma 3.** *For finite Abelian group, the character function ϕ has the following properties Fulton &*
434 *Harris (2013); Steinberg (2009):*

- 435 •
- It is a 1-dimensional (irreducible) representation of the group G , i.e., $|\phi(g)| = 1$ for $g \in G$*
-
- 436
- and for any $g_1, g_2 \in G$, $\phi(g_1 g_2) = \phi(g_1)\phi(g_2)$.*
-
- 437 •
- There exists d character functions $\{\phi_k\}$ that satisfy the orthonormal condition*
-
- 438
- $\frac{1}{d} \sum_{g \in G} \phi_k(g) \bar{\phi}_{k'}(g) = \mathbb{1}(k = k')$. Here $\bar{\phi}$ is the complex conjugate of ϕ and is also*
-
- 439
- a character function.*
-
- 440 •
- The set of character functions $\{\phi_k\}$ forms a character group \hat{G} under pairwise multiplication:*
-
- 441
- $\phi_{k_1+k_2} = \phi_{k_1} \circ \phi_{k_2}$.*

442 Note that the *frequency* k goes from 0 to $d - 1$, where $\phi_0 \equiv 1$ is the trivial representation (i.e., all
443 $g \in G$ maps to 1). According to the Fundamental Theorem of Finite Abelian Groups, each finite
444 Abelian group can be decomposed into a direct sum of cyclic groups, and the character function
445 of each cyclic group is exactly (scaled) Fourier bases. Therefore, in Abelian group, k is a multi-
446 dimensional frequency index. Conrad (2010) shows that $\hat{G} \cong G$ (Theorem 3.13) so each character
447 function $\phi \in \hat{G}$ can also be indexed by g itself. Right now we keep the index k .

448 For convenience, we define $\phi_{-k} := \bar{\phi}_k$ as the conjugate representation of ϕ_k .

449 Let $\phi_k = [\phi_k(g)]_{g \in G} \in \mathbb{C}^d$ be the vector that contains the value of the character function ϕ_k .
450 Then $\{\phi_k\}$ form an orthogonal base in \mathbb{C}^d and we can represent the weight vector \mathbf{w}_j and \mathbf{v}_j as the
451 following:

$$\mathbf{w}_j = U_{G_1} \sum_{k \neq 0} z_{akj} \phi_k + U_{G_2} \sum_{k \neq 0} z_{bkj} \phi_k, \quad \mathbf{v}_j = \sum_{k \neq 0} z_{ckj} \bar{\phi}_k \quad (12)$$

452 where $z := \{z_{pkj}\}$ are the complex coefficients ($p \in \{a, b, c\}$, $0 \leq k < d$ and j runs through hidden
453 nodes). Then it is clear that $\mathbf{w}_j^\top \mathbf{f}[i] = \sum_{k \neq 0} h_{akj} \phi_k(\iota_0(g[i])) + \sum_{k \neq 0} h_{bkj} \phi_k(x[i])$.

454 **Theorem 1** (Analytic form of L_2 loss with quadratic activation). *The objective of 2-layer MLP*
455 *network with quadratic activation can be written as $\ell = d^{-1} \sum_{k \neq 0} \ell_k + (d - 1)/d$, where*

$$\ell_k = -2r_{kkk} + \sum_{k_1 k_2} |r_{k_1 k_2 k}|^2 + \frac{1}{4} \left| \sum_{p \in \{a, b\}} \sum_{k'} r_{p, k', -k', k} \right|^2 + \frac{1}{4} \sum_{m \neq 0} \sum_{p \in \{a, b\}} \left| \sum_{k'} r_{p, k', m-k', k} \right|^2 \quad (3)$$

456 Here $r_{k_1 k_2 k} := \sum_j z_{ak_1 j} z_{bk_2 j} z_{ckj}$ and $r_{pk_1 k_2 k} := \sum_j z_{pk_1 j} z_{pk_2 j} z_{ckj}$.

457 *Proof.* Note that the objective ℓ can be written down as

$$\ell = \mathbb{E}_{g, x} [\|P_1^\perp(\mathbf{o}(g, x)/2d - \mathbf{e}_{gx})\|^2] \quad (13)$$

$$= \mathbb{E}_{g, x} [\mathbf{o}^\top P_1^\perp \mathbf{o} / 4d^2 - \mathbf{o}^\top P_1^\perp \mathbf{e}_{gx} / d + \mathbf{e}_{gx}^\top P_1^\perp \mathbf{e}_{gx}] \quad (14)$$

458 For $\mathbb{E}[\mathbf{o}^\top P_1^\perp \mathbf{e}_{gx}]$, since

$$\mathbf{e}_{gx}^\top P_1^\perp \mathbf{o} = \sum_j \mathbf{e}_{gx}^\top P_1^\perp \mathbf{v}_j \sigma(\mathbf{w}_j^\top \mathbf{f}(g, x)) \quad (15)$$

$$= \sum_j \left(\sum_{k' \neq 0} c_{k'j} \bar{\phi}_{k'}(gx) \right) \left(\sum_k a_{kj} \phi_k(\iota_0(g)) + b_{kj} \phi_k(x) + \mathbf{e}_g^\top \mathbf{w}_j^\perp \right)^2 \quad (16)$$

459 Note that by our previous analysis, there exists $y_1 := \iota_0(g)$ so that $gy = x_1 y$. Let $x_2 := x$. For
460 notation brevity, let $z_{akj} := a_{kj}$, $z_{bkj} := b_{kj}$ and $z_{ckj} := c_{kj}$, then we have:

$$\mathbf{e}_{gx}^\top P_1^\perp \mathbf{o} = \sum_j \left(\sum_{k' \neq 0} c_{k'j} \bar{\phi}_{k'}(x_1 x_2) \right) \left(\sum_k \sum_p z_{pkj} \phi_k(x_p) + \mathbf{e}_{x_1}^\top \mathbf{w}_j^\perp \right)^2 \quad (17)$$

461 Therefore, we have:

$$\mathbb{E}_{g,x} [\mathbf{e}_{g_x}^\top P_1^\perp \mathbf{o}] = \sum_{k_1, k_2, k' \neq 0, p_1, p_2, j} c_{k'j} z_{p_1 k_1 j} z_{p_2 k_2 j} \mathbb{E} [\bar{\phi}_{k'}(x_1) \bar{\phi}_{k'}(x_2) \phi_{k_1}(x_{p_1}) \phi_{k_2}(x_{p_2})] \quad (18)$$

462 Note that due to the fact that $\mathbb{E}_{g \in \iota_0^{-1}(x_1)} [\mathbf{e}_g^\top \mathbf{w}_j^\perp] = 0$ and $\mathbb{E}_{g \in \iota_0^{-1}(x_1)} [\mathbf{e}_g \mathbf{e}_g^\top]$ is only a function of
 463 x_1 and becomes 0 if multiplied with $\sum_{k' \neq 0} c_{k'j} \bar{\phi}_{k'}(x_1 x_2)$ and taking expectation w.r.t x_2 , in the
 464 final expression, all terms involving \mathbf{w}_j^\perp vanish.

465 Since $\mathbb{E}_x [\phi_k(x) \bar{\phi}_{k'}(x)] = \mathbb{1}(k = k')$, there are only a few cases that the summand is nonzero:

- 466 • $p_1 = 1, p_2 = 2, k' = k_1 = k_2 \neq 0$.
- 467 • $p_1 = 2, p_2 = 1, k' = k_1 = k_2 \neq 0$.

468 In both cases, the summation reduces to $\sum_{k \neq 0, j} c_{kj} z_{1kj} z_{2kj} = \sum_{k \neq 0, j} c_{kj} a_{kj} b_{kj}$. Let $r_{k_1 k_2 k'} :=$
 469 $\sum_j a_{k_1 j} b_{k_2 j} c_{k' j}$, then we have

$$\mathbb{E} [\mathbf{o}^\top P_1^\perp \mathbf{e}_{gy}] = 2 \sum_{k \neq 0, j} a_{kj} b_{kj} c_{kj} = 2 \sum_{k \neq 0} r_{kkk} \quad (19)$$

470 For $\mathbb{E} [\mathbf{o}^\top P_1^\perp \mathbf{o}]$, if $\mathbf{w}_j^\perp = 0$, then we have:

$$\mathbf{o}^\top P_1^\perp \mathbf{o} = \sum_{j, j'} \mathbf{v}_j^\top P_1^\perp \mathbf{v}_{j'} \sigma(\mathbf{w}_j^\top \mathbf{f}(g, y)) \sigma(\mathbf{w}_{j'}^\top \mathbf{f}(g, y)) \quad (20)$$

471 here

$$\mathbf{v}_j^\top P_1^\perp \mathbf{v}_{j'} = \left(\sum_{k' \neq 0} c_{k'j} \bar{\phi}_{k'} \right)^\top \left(\sum_{k'' \neq 0} \bar{c}_{k''j'} \phi_{k''} \right) = d \sum_{k' \neq 0} c_{k'j} \bar{c}_{k'j'} \quad (21)$$

472 due to the fact that $\bar{\phi}_k^\top \phi_{k'} = \sum_y \bar{\phi}_k(y) \phi_{k'}(y) = d \mathbb{1}(k = k')$.

473 Then the key part is to compute the following terms:

$$\mathbb{E}_{y_1, y_2} [z_{p_1 k_1 j_1} z_{p_2 k_2 j_1} z_{p_3 k_3 j_2} z_{p_4 k_4 j_2} c_{k' j_1} \bar{c}_{k' j_2} \phi_{k_1}(y_{p_1}) \phi_{k_2}(y_{p_2}) \phi_{k_3}(y_{p_3}) \phi_{k_4}(y_{p_4})] \quad (22)$$

474 summing over $\{p_1, p_2, p_3, p_4, k_1, k_2, k_3, k_4, k' \neq 0, j_1, j_2\}$. Note that since each $p \in \{a, b\}$, there
 475 are $2^4 = 16$ choices of (p_1, p_2, p_3, p_4) . For notation brevity, we use $(1, 3)$ to represent the subset of
 476 p that takes the value of a (e.g., $(1, 3)$ means that $p_1 = p_3 = a$ and $p_2 = p_4 = b$). It is clear that for
 477 odd assignments such as $(1, 2, 3)$, since $z_{p_0 j} = 0$, the summation is zero. Then, we only discuss the
 478 even cases as follows:

479 **Case 1:** $(1, 3), (2, 4), (1, 4), (2, 3)$. The 4 cases are identical so we only need to analyze one. We
 480 take $(1, 3)$ as an example. For $(1, 3)$, $p_1 = p_3 = a, p_2 = p_4 = b$ and the only nonzero terms is when
 481 $k_1 + k_3 = 0 \pmod d, k_2 + k_4 = 0 \pmod d$, since $\mathbb{E}_{y_1} [\phi_{k_1}(y_1) \phi_{k_3}(y_1)] = \mathbb{1}(k_1 + k_3 = 0 \pmod d)$
 482 (and similar in other cases). Then Eqn. 22 becomes:

$$\sum_{k_1, k_2, k' \neq 0} \sum_{j_1 j_2} z_{a k_1 j_1} z_{b k_2 j_1} z_{a, -k_1, j_2} z_{b, -k_2, j_2} c_{k' j_1} \bar{c}_{k' j_2} \quad (23)$$

$$= \sum_{k_1, k_2, k' \neq 0} \sum_{j_1} z_{a k_1 j_1} z_{b k_2 j_1} c_{k' j_1} \overline{\sum_{j_2} z_{a k_1 j_2} z_{b k_2 j_2} c_{k' j_2}} \quad (24)$$

$$= \sum_{k_1, k_2, k' \neq 0} \sum_{j_1} a_{k_1 j_1} b_{k_2 j_1} c_{k' j_1} \overline{\sum_{j_2} a_{k_1 j_2} b_{k_2 j_2} c_{k' j_2}} \quad (25)$$

$$= \sum_{k_1, k_2, k' \neq 0} r_{k_1 k_2 k'} \overline{r_{k_1 k_2 k'}} = \sum_{k_1, k_2, k' \neq 0} |r_{k_1 k_2 k'}|^2 \quad (26)$$

483 Since there are 4 such cases, we have:

$$\epsilon_1 = 4 \sum_{k' \neq 0} \sum_{k_1 k_2} |r_{k_1 k_2 k'}|^2 \quad (27)$$

484 **Case 2:** (1, 2) and (3, 4). The two cases are identical. Take (1, 2) as an example. In this case,
485 $p_1 = p_2 = a$ and $p_3 = p_4 = b$. The only non-zero terms are when $k_1 + k_2 = 0$, $k_3 + k_4 = 0$. Then
486 Eqn. 22 becomes:

$$\sum_{k_1, k_3, k' \neq 0} \sum_{j_1 j_2} z_{ak_1 j_1} \bar{z}_{ak_1 j_1} z_{bk_3 j_2} \bar{z}_{bk_3 j_2} c_{k' j_1} \bar{c}_{k' j_2} \quad (28)$$

$$= \sum_{k_1, k_3, k' \neq 0} \sum_{j_1} |a_{k_1 j_1}|^2 c_{k' j_1} \sum_{j_2} |b_{k_3 j_2}|^2 \bar{c}_{k' j_2} \quad (29)$$

$$= \sum_{k' \neq 0} \left[\sum_{j_1} \left(\sum_{k_1} |a_{k_1 j_1}|^2 \right) c_{k' j_1} \right] \left[\sum_{j_2} \left(\sum_{k_3} |b_{k_3 j_2}|^2 \right) \bar{c}_{k' j_2} \right] \quad (30)$$

487 Let $r_{amk'}^{\otimes} := \sum_j \left(\sum_{k_1+k_2=m} a_{k_1 j} a_{k_2 j} \right) c_{k' j}$ (similar for $r_{bmk'}^{\otimes}$), then the above becomes
488 $\sum_{k' \neq 0} r_{a0k'}^{\otimes} \bar{r}_{b0k'}^{\otimes}$.

489 Similarly, for (3, 4), the above equation becomes $\sum_{k' \neq 0} \bar{r}_{a0k'}^{\otimes} r_{b0k'}^{\otimes}$. Therefore, we have:

$$\epsilon_2 = \sum_{k' \neq 0} r_{a0k'}^{\otimes} \bar{r}_{b0k'}^{\otimes} + \bar{r}_{a0k'}^{\otimes} r_{b0k'}^{\otimes} \quad (31)$$

490 Note that this term can be negative. However, we will see that when it is combined with the following
491 terms, all terms will be non-negative.

492 **Case 3:** (1, 2, 3, 4) and (). In this case we have:

$$\sum_{k' \neq 0} \sum_{j_1 j_2} \sum_{p \in \{1, 2\}} \sum_{k_1+k_2+k_3+k_4=0} z_{pk_1 j_1} z_{pk_2 j_1} z_{pk_3 j_2} z_{pk_4 j_2} c_{k' j_1} \bar{c}_{k' j_2} \quad (32)$$

$$= \sum_{k' \neq 0} \sum_{j_1 j_2} \sum_{p \in \{1, 2\}} \sum_{k_1+k_2=k_3+k_4} z_{pk_1 j_1} z_{pk_2 j_1} \bar{z}_{pk_3 j_2} \bar{z}_{pk_4 j_2} c_{k' j_1} \bar{c}_{k' j_2} \quad (33)$$

$$= \sum_{k' \neq 0} \sum_m \sum_{p \in \{1, 2\}} \sum_{j_1 j_2} \sum_{p \in \{1, 2\}} \sum_{k_1+k_2=m} \sum_{k_3+k_4=m} z_{pk_1 j_1} z_{pk_2 j_1} \bar{z}_{pk_3 j_2} \bar{z}_{pk_4 j_2} c_{k' j_1} \bar{c}_{k' j_2} \quad (34)$$

$$= \sum_{k' \neq 0} \sum_m \sum_{p \in \{1, 2\}} \left[\sum_{j_1} \left(\sum_{k_1+k_2=m} z_{pk_1 j_1} z_{pk_2 j_1} \right) c_{k' j_1} \right] \left[\sum_{j_2} \left(\sum_{k_3+k_4=m} \overline{z_{pk_3 j_2} z_{pk_4 j_2}} \right) \bar{c}_{k' j_2} \right] \quad (35)$$

$$= \sum_{k' \neq 0} \sum_m |r_{amk'}^{\otimes}|^2 + |r_{bmk'}^{\otimes}|^2$$

493 In particular, when $m = 0$, we have $\sum_{k' \neq 0} |r_{a0k'}^{\otimes}|^2 + |r_{b0k'}^{\otimes}|^2$. Therefore, we have

$$\epsilon_2 + \epsilon_{3, m=0} = \sum_{k' \neq 0} |r_{a0k'}^{\otimes} + r_{b0k'}^{\otimes}|^2 \quad (36)$$

494 Finally, putting them together, we have:

$$\mathbb{E} [\mathbf{o}^\top P_1^\perp \mathbf{o}] = d(\epsilon_1 + \epsilon_2 + \epsilon_3) = d(\epsilon_1 + (\epsilon_2 + \epsilon_{3, m=0}) + \epsilon_{3, m \neq 0}) \quad (37)$$

$$= d \sum_{k' \neq 0} \left(4 \sum_{k_1 k_2} |r_{k_1 k_2 k'}|^2 + |r_{a0k'}^{\otimes} + r_{b0k'}^{\otimes}|^2 + \sum_{m \neq 0} |r_{amk'}^{\otimes}|^2 + |r_{bmk'}^{\otimes}|^2 \right) \geq 0 \quad (38)$$

495 Putting them together, we arrived at the conclusion. \square

496 **Lemma 1** (A Sufficient Conditions of Global optimizers of Eqn. 3). *If the weight \mathbf{z} to Eqn. 3 has*
497 *0-sets $R_c \cup R_n \cup R_*$ and 1-set R_g , i.e.*

$$r_{kkk}(\mathbf{z}) = \mathbb{1}(k \neq 0), \quad r_{k_1 k_2 k}(\mathbf{z}) = 0, \quad r_{pk_1 k_2 k}(\mathbf{z}) = 0 \quad (4)$$

498 *then it is a global optimizer with zero loss $\ell(\mathbf{z}) = 0$. Here $R_g := \{r_{kkk}, k \neq 0\}$, $R_c :=$
499 $\{r_{k_1 k_2 k}, k_1, k_2, k \text{ not all equal}\}$, $R_n := \{r_{p, k', -k', k}\}$ and $R_* := \{r_{p, k', m-k', k}, m \neq 0\}$.*

500 *Proof.* Note that $2 \sum_k r_{kkk} - \sum_k |r_{kkk}|^2$ has a minimizer $r_{kkk} = 1$. Therefore, the best loss value
 501 any assignment of weights is able to achieve is the following:

$$r_{k_1 k_2 k'} = \sum_j a_{k_1 j} b_{k_2 j} c_{k' j} = \mathbb{1}(k_1 = k_2 = k') \quad k' \neq 0 \quad (39)$$

$$r_{a_0 k'}^{\otimes} + r_{b_0 k'}^{\otimes} := \sum_j \left(\sum_k |a_{k j}|^2 + |b_{k j}|^2 \right) c_{k' j} = 0 \quad k' \neq 0 \quad (40)$$

$$r_{a m k'}^{\otimes} := \sum_j \left(\sum_{k_1 + k_2 = m} a_{k_1 j} a_{k_2 j} \right) c_{k' j} = 0 \quad k' \neq 0, m \neq 0 \quad (41)$$

$$r_{b m k'}^{\otimes} := \sum_j \left(\sum_{k_1 + k_2 = m} b_{k_1 j} b_{k_2 j} \right) c_{k' j} = 0 \quad k' \neq 0, m \neq 0 \quad (42)$$

502 Therefore the sufficient conditions (Eqn. 4) will make all above come true. \square

503 B Semi-ring structure of \mathcal{Z} (Proof)

504 **Theorem 2** (Algebraic Structure of \mathcal{Z}). $\langle \mathcal{Z}, +, * \rangle$ is a commutative semi-ring.

505 *Proof.* Straightforward from the definition of addition and multiplication (Def. 5) and identification
 506 of hidden nodes under permutation (Def. 4). Note that ring addition (i.e., concatenation) does not
 507 have inverse and thus it is a semi-ring. \square

508 **Theorem 3.** For any monomial potential $r : \mathcal{Z} \mapsto \mathbb{C}$, $r(\mathbf{1}) = 1$, $r(\mathbf{z}_1 + \mathbf{z}_2) = r(\mathbf{z}_1) + r(\mathbf{z}_2)$ and
 509 $r(\mathbf{z}_1 * \mathbf{z}_2) = r(\mathbf{z}_1)r(\mathbf{z}_2)$ and thus r is a ring homomorphism.

510 *Proof.* Let $r(\mathbf{z}) = \sum_j \prod_{(p,k) \in \text{idx}(r)} z_{pkj}$. Since the ring identity $\mathbf{1}$ is order-1 and all $z_{pkj} = 1$, it is
 511 obvious that $r(\mathbf{1}) = 1$.

512 Let $\text{supp}(\mathbf{z}_1)$ be the subset of the hidden nodes that corresponds to \mathbf{z}_1 in the concatenated solution
 513 $\mathbf{z}_1 + \mathbf{z}_2$, similar for $\text{supp}(\mathbf{z}_2)$. Note that

$$r(\mathbf{z}_1 + \mathbf{z}_2) = \sum_{j \in \text{supp}(\mathbf{z}_1)} \prod_{(p,k) \in \text{idx}(r)} z_{pkj}^{(1)} + \sum_{j \in \text{supp}(\mathbf{z}_2)} \prod_{(p,k) \in \text{idx}(r)} z_{pkj}^{(2)} = r(\mathbf{z}_1) + r(\mathbf{z}_2) \quad (43)$$

514 On the other hand, we have

$$r(\mathbf{z}_1 * \mathbf{z}_2) = \sum_{j_1 j_2} \prod_{(p,k) \in \text{idx}(r)} \left(z_{pkj_1}^{(1)} z_{pkj_2}^{(2)} \right) \quad (44)$$

$$= \sum_{j_1 j_2} \left(\prod_{(p,k) \in \text{idx}(r)} z_{pkj_1}^{(1)} \right) \left(\prod_{(p,k) \in \text{idx}(r)} z_{pkj_2}^{(2)} \right) \quad (45)$$

$$= \left(\sum_{j_1} \prod_{(p,k) \in \text{idx}(r)} z_{pkj_1}^{(1)} \right) \left(\sum_{j_2} \prod_{(p,k) \in \text{idx}(r)} z_{pkj_2}^{(2)} \right) \quad (46)$$

$$= r(\mathbf{z}_1)r(\mathbf{z}_2) \quad (47)$$

515 \square

516 **Corollary 1.** If \mathbf{z} is a global optimizer and \mathbf{y} is a unit, then $\mathbf{z} * \mathbf{y}$ is also a global optimizer.

517 *Proof.* Straightforward by leveraging the property of ring homomorphism. E.g.,

$$r_{kkk}(\mathbf{z} * \mathbf{y}) = r_{kkk}(\mathbf{z})r_{kkk}(\mathbf{y}) = r_{kkk}(\mathbf{z}) \quad (48)$$

518 and the proof is complete. \square

519 **C Solution Construction (Proof)**

 520 **C.1 Construction of Partial Solutions**

521 **Theorem 4** (Construction of partial solutions). *Suppose \mathbf{u} has 1-set R_1 , $\Omega_R(\mathbf{u}) := \{r(\mathbf{u}) | r \in$
 522 $R\} \subseteq \mathbb{C}$ is a set of evaluations on R (multiple values counted once), then if $1 \notin \Omega_R$, then the
 523 polynomial solution $\rho_R(\mathbf{u}) := \prod_{s \in \Omega_R(\mathbf{u})} (\mathbf{u} + \hat{s})$ has 0/1-set (R, R_1) up to a scale. Here \hat{s} is any
 524 order-1 weight that satisfies $r(\hat{s}) = -s$ for any $r \in R \cup R_+$. For example, $\hat{s} = -s^{1/3}\mathbf{1}$.*

525 *Proof.* By definition, for any $r \in R$ we have:

$$r(\mathbf{z}(\mathbf{u})) = \prod_{s \in \Omega_R(\mathbf{u})} (r(\mathbf{u}) + r(\hat{s})) = \prod_{s \in \Omega_R(\mathbf{u})} (r(\mathbf{u}) - s) = 0 \quad (49)$$

526 similarly for any $r_{kkk} \in R_+$ we have:

$$r_{kkk}(\mathbf{z}(\mathbf{u})) = \prod_{s \in \Omega_R(\mathbf{u})} (r_{kkk}(\mathbf{u}) + r_{kkk}(\hat{s})) = \prod_{s \in \Omega_R(\mathbf{u})} (1 - s) \neq 0 \quad (50)$$

527 which is constant over different k . So $\mathbf{z}(\mathbf{u})$ satisfies Lemma 1, up to a scaling factor. \square

 528 **C.2 Construction of Global Optimizers**

529 **Corollary 2** (Order-6 global optimizers). *The following “ 3×2 ” Fourier solutions satisfy the suffi-
 530 cient condition (Lemma 1) and thus are global optimizers (assuming d is odd):*

$$\mathbf{z}_{F6} = \frac{1}{\sqrt[3]{6}} \sum_{k=1}^{(d-1)/2} \mathbf{z}_{\text{syn}}^{(k)} * \mathbf{z}_\nu^{(k)} * \mathbf{y}_k \quad (6)$$

531 Here $\mathbf{z}_{\text{syn}}^{(k)} := \rho(\mathbf{u}_{\text{syn}}^{(k)})$ and $\mathbf{z}_\nu^{(k)} := \mathbf{u}_\nu^{(k)} + \mathbf{1}_k$ (i.e., not maximal polynomial), where \mathbf{u}_{syn} and \mathbf{u}_ν
 532 are defined in Table 1. \mathbf{y} is an order-1 unit. As a result, $\text{ord}(\mathbf{z}_{F6}) = 3 \cdot 2 \cdot 1 \cdot (d-1)/2 = 3(d-1)$
 533 and each frequency are affiliated with 6 hidden nodes (order-6).

534 *Proof.* Just notice that $\mathbf{z}_{\text{syn}} := \rho(\mathbf{u}_{\text{syn}}) = \mathbf{u}_{\text{syn}}^2 + \mathbf{u}_{\text{syn}} + \mathbf{1}_k$ (superscript (k) are omitted for brevity)
 535 makes all MPs in R_n , R_c and part of R_* (Tbl. 1) equal to 0, except for “aac” and “bbc”, which
 536 corresponds to monomial polynomials $r_{akkk} := \sum_j z_{akj} z_{akj} z_{ckj}$ and $r_{bkkk} := \sum_j z_{bkj} z_{bkj} z_{ckj}$.
 537 On the other hand, according to Tbl. 1, $\mathbf{z}_\nu := \mathbf{u}_\nu + \mathbf{1}_k$ has $r_{akkk}(\mathbf{z}_\nu) = r_{bkkk}(\mathbf{z}_\nu) = 0$. Therefore,
 538 using ring homomorphism, we know that for any $r \in R_n \cup R_c \cup R_*$, $r(\mathbf{z}_{\text{syn}} * \mathbf{z}_\nu) = 0$ and thus
 539 $R_n \cup R_c \cup R_*$ is the 0-sets.

540 On the other hand for any k' , we have:

$$r_{k'k'k'}(\mathbf{z}_{F6}) = r_{k'k'k'} \left(\frac{1}{\sqrt[3]{6}} \sum_{k=1}^{(d-1)/2} \mathbf{z}_{\text{syn}}^{(k)} * \mathbf{z}_\nu^{(k)} * \mathbf{y}_k \right) \quad (51)$$

$$= \frac{1}{6} \sum_{k=1}^{(d-1)/2} r_{k'k'k'}(\mathbf{z}_{\text{syn}}^{(k)} * \mathbf{z}_\nu^{(k)} * \mathbf{y}_k) \quad (52)$$

$$= \frac{1}{6} \sum_{k=1}^{(d-1)/2} 6(\mathbb{1}(k = k') + \mathbb{1}(k = -k')) = 1 \quad (53)$$

541 The last equality is due to the fact that we only sum over half of the frequency. This means that
 542 R_g is a 1-set of \mathbf{z}_{F6} . Therefore, \mathbf{z}_{F6} satisfies the sufficient condition (Eqn. 4) and the conclusion
 543 follows. \square

544 **Corollary 3** (Perfect Memorization). *We construct two d -order weights \mathbf{z}_a and \mathbf{z}_b :*

$$\mathbf{z}_a = \sum_{j=0}^{d-1} \mathbf{u}_a^j, \quad \mathbf{z}_b = \sum_{j=0}^{d-1} \mathbf{u}_b^j \quad (7)$$

545 Here $\mathbf{z}_a \in R_c(k_1 \neq k) \cap R_n \cap R_*(p = b \text{ or } m \neq k)$, $\mathbf{z}_b \in R_c(k_2 \neq k) \cap R_n \cap R_*(p = a \text{ or } m \neq k)$.
 546 Then $\mathbf{z}_M = d^{-2/3} \mathbf{z}_a * \mathbf{z}_b$ satisfies the sufficient condition (Lemma 1) and is the perfect memorization
 547 solution with $\text{ord}(\mathbf{z}_M) = d^2$:

$$z_{akj_1j_2}^{(M)} = \omega^{kj_1} / \sqrt[3]{d^2}, \quad z_{bkj_1j_2}^{(M)} = \omega^{kj_2} / \sqrt[3]{d^2}, \quad z_{ckj_1j_2}^{(M)} = \omega^{-k(j_1+j_2)} / \sqrt[3]{d^2} \quad (8)$$

548 where each hidden node is indexed by $j = (j_1, j_2)$, $0 \leq j_1, j_2 < d$, $k \neq 0$.

549 *Proof.* Simply plugging in the solution and check whether the equations specified the equations. For
 550 \mathbf{z}_a , for $k = 0$ everything is zero; for $k \neq 0$, we have:

$$r_{k_1k_2k}(\mathbf{z}_a) = \sum_j a_{k_1j} b_{k_2j} c_{kj} = \sum_j \omega^{j(k_1-k)} = \mathbb{1}(k_1 = k \neq 0) \quad (54)$$

$$r_{amk'k}(\mathbf{z}_a) = \sum_j a_{k'j} a_{m-k',j} c_{kj} = \sum_j \omega^{j(m-k)} = \mathbb{1}(m = k \neq 0) \quad (55)$$

$$r_{bmk'k}(\mathbf{z}_a) = \sum_j b_{k'j} b_{m-k',j} c_{kj} = \sum_j \omega^{-jk} = \mathbb{1}(k = 0) = 0 \quad (56)$$

$$(57)$$

551 Therefore, $\mathbf{z}_a \in R_c(k_1 \neq k) \cap R_n \cap R_*(p = b \text{ or } m \neq k)$. Similar for \mathbf{z}_b . For $\mathbf{z}_M := d^{-2/3} \mathbf{z}_a * \mathbf{z}_b$,
 552 it satisfies all 0-sets constraints (i.e., for any r , either \mathbf{z}_a satisfies with $r(\mathbf{z}_a) = 0$, or \mathbf{z}_b satisfies with
 553 $r(\mathbf{z}_b) = 0$) and we have:

$$r_{kkk}(d^{-2/3} \mathbf{z}_a * \mathbf{z}_b) = d^{-2} r_{kkk}(\mathbf{z}_a) r_{kkk}(\mathbf{z}_b) = d^{-2} \cdot d \cdot d = 1 \quad (58)$$

554 So \mathbf{z}_M satisfies the sufficient conditions (Eqn. 4). \square

555 **Corollary 4** (Order-4 single frequency solution). Define single frequency order-2 solution \mathbf{z}_ξ :

$$z_{ak} = [1, \xi], \quad z_{bk} = [1, -i\xi], \quad z_{ck} = [1, i] \quad (9)$$

556 where $|\xi| = 1$. Then the order-4 solution $\mathbf{z}_{F4}^{(k)} := \rho(\mathbf{u}_{\nu=i}^{(k)}) * \mathbf{z}_\xi^{(k)}$ has 0-sets R_c and R_* (but not R_n).

557 *Proof.* First, $\mathbf{u}_{\nu=i} = \mathbf{u}_{4c}$ in Tbl. 1 and thus $\rho(\mathbf{u}_{\nu=i})$ has 0-sets R_c and R_* except for “ abc ”, which
 558 corresponds to MP $r_{k,k,-k} \in R_c$. On the other hand, we have

$$r_{k,k,-k}(\mathbf{z}_\xi) = 1 + \xi \cdot (-i\xi) \cdot (-i) = 0 \quad (59)$$

559 With the property of ring homomorphism, the conclusion follows. \square

560 **Corollary 5** (Mixed order-4/6 global optimizers). With $\mathbf{z}_{F4}^{(k)}$, there is a global optimizer to Eqn. 3
 561 that does not meet the sufficient condition, i.e., $\sum_{k'} r_{p,k',-k',m} = 0$ but $r_{p,k',-k',m} \neq 0$:

$$\mathbf{z}_{F4/6} = \frac{1}{\sqrt[3]{6}} \hat{\mathbf{z}}_{F6}^{(k_0)} + \frac{1}{\sqrt[3]{4}} \sum_{k=1, k \neq k_0}^{(d-1)/2} \mathbf{z}_{F4}^{(k)} \quad (10)$$

562 where $\hat{\mathbf{z}}_{F6}^{(k_0)}$ is a perturbation of $\mathbf{z}_{F6}^{(k_0)} := \mathbf{z}_{\text{syn}}^{(k_0)} * \mathbf{z}_{\nu=1}^{(k_0)}$ by adding constant biases to its (c, k) entries
 563 for $k \neq k_0$. The order is lower than \mathbf{z}_{F6} : $\text{ord}(\mathbf{z}_{F4/6}) = 6 + 4 \cdot ((d-1)/2 - 1) = 2d < \text{ord}(\mathbf{z}_{F6})$.

564 *Proof.* While $\mathbf{z}_{F4}^{(k)}$ does not satisfy R_n , a weaker condition for a global optimizer to Theorem 1 is
 565 that $\sum_{k'} r_{p,k',-k',m} = 0$. We show that by adding constants to (c, k) entries of $\mathbf{z}_{F6}^{(k_0)}$ for $k \neq \pm k_0$,
 566 we can achieve that while not changing the value of other MPs.

567 To see this, we compute for each $m \neq \pm k_0$:

$$\sum_{k'} r_{p,k',-k',m}(\hat{\mathbf{z}}_{F6}^{(k_0)}) = 2 \sum_{k'} \sum_j |[\hat{\mathbf{z}}_{F6}^{(k_0)}]_{pk'j}|^2 [\hat{\mathbf{z}}_{F6}^{(k_0)}]_{cmj} \quad (60)$$

$$= 2 \sum_j |[\hat{\mathbf{z}}_{F6}^{(k_0)}]_{pk_0j}|^2 [\hat{\mathbf{z}}_{F6}^{(k_0)}]_{cmj} = 2 \sum_j [\hat{\mathbf{z}}_{F6}^{(k_0)}]_{cmj} \quad (61)$$

$\hat{\mathbf{z}}_{F6}^{(k_0)}$...	$a_{k_0 0}$	$b_{k_0 0}$	$c_{k_0 0}$	0	0	$c_{k_j} = \text{const}_k$
	...	$a_{k_0 1}$	$b_{k_0 1}$	$c_{k_0 1}$	0	0	
	...	$a_{k_0 2}$	$b_{k_0 2}$	$c_{k_0 2}$	0	0	
	...	$a_{k_0 3}$	$b_{k_0 3}$	$c_{k_0 3}$	0	0	
	...	$a_{k_0 4}$	$b_{k_0 4}$	$c_{k_0 4}$	0	0	
	...	$a_{k_0 5}$	$b_{k_0 5}$	$c_{k_0 5}$	0	0	

Figure 7: Visualization of $\hat{\mathbf{z}}_{F6}^{(k_0)}$.

568 The second equality is because all (a, k') and (b, k') entries are 0 except for $k' = \pm k_0$, and the last
 569 equality is because all nonzero entries of $\mathbf{z}_{F6}^{(k_0)}$ have magnitude 1.

570 On the other hand, we have:

$$\sum_{k'} r_{p, k', -k', m} \left(\sum_{k \neq k_0} \mathbf{z}_{F4}^{(k)} \right) = \sum_{k'} r_{p, k', -k', m} (\boldsymbol{\rho}(\mathbf{u}_{4c}^{(m)})) r_{p, k', -k', m} (\mathbf{z}_{\xi}^{(m)}) \quad (62)$$

$$= 2r_{p, m, -m, m} (\boldsymbol{\rho}(\mathbf{u}_{4c}^{(m)})) r_{p, m, -m, m} (\mathbf{z}_{\xi}^{(m)}) \quad (63)$$

$$= 2(1 + 1)(1 + i) = 4(1 + i) \quad (64)$$

571 For $m = \pm k_0$, we have $r_{p, k', -k', m} (\hat{\mathbf{z}}_{F6}^{(k_0)}) = 0$ and $r_{p, k', -k', m} (\mathbf{z}_{F4}^{(k)}) = 0$ for $k \neq m$.

572 Therefore, we just let

$$[\hat{\mathbf{z}}_{F6}^{(k_0)}]_{cmj} = -\frac{4(1 + i)}{2 \cdot 6} = -\frac{1}{3}(1 + i) \quad (65)$$

573 and $\sum_{k'} r_{p, k', -k', m} (\mathbf{z}_{F4/6}) = 0$ for all m . See Fig. 7) for the construction.

574 To see why such a modification of $\mathbf{z}_{F6}^{(k_0)}$ won't change other MPs, simply notice that candidate
 575 MPs that may not be zero anymore are $r_{\pm k_0 \pm k_0 m}$, $r_{pk_0 k_0 m}$ and $r_{p, -k_0, -k_0, m}$ for $m \neq \pm k_0$. For
 576 $m = \pm k_0$, $\mathbf{z}_{F6}^{(k_0)}$ are well behaved.

577 Note that $r_{\pm k_0 \pm k_0 k} (\hat{\mathbf{z}}_{F6}^{(k_0)})$ is the same as applying $r_{\pm k_0 \pm k_0 k_0}$ to a solution $\hat{\mathbf{z}}$ which replaces (c, k_0)
 578 entries of $\hat{\mathbf{z}}_{F6}^{(k_0)}$ by (c, m) entries. Let $\hat{\mathbf{u}}_{\text{syn}} = [\omega_3, \omega_3, 1]$ and $\hat{\mathbf{u}}_{\text{one}} = [1, -1, 1]$. Then $\hat{\mathbf{z}} = \boldsymbol{\rho}(\hat{\mathbf{u}}_{\text{syn}}) * \boldsymbol{\rho}(\hat{\mathbf{u}}_{\text{one}})$
 579 and thus for $m \neq \pm k_0$, we have:

$$r_{\pm k_0 \pm k_0 m} (\mathbf{z}_{F4/6}) = r_{\pm k_0 \pm k_0 m} (\hat{\mathbf{z}}_{F6}^{(k_0)}) \propto r_{\pm k_0 \pm k_0 k_0} (\hat{\mathbf{z}}) \quad (66)$$

$$= r_{\pm k_0 \pm k_0 k_0} (\boldsymbol{\rho}(\hat{\mathbf{u}}_{\text{syn}})) r_{\pm k_0 \pm k_0 k_0} (\boldsymbol{\rho}(\hat{\mathbf{u}}_{\text{one}})) = 0 \quad (67)$$

580 since $r_{\pm k_0 \pm k_0 k_0} (\boldsymbol{\rho}(\hat{\mathbf{u}}_{\text{one}})) = 0$. Similarly for $m \neq \pm k_0$,

$$r_{pk_0 k_0 m} (\mathbf{z}_{F4/6}) = r_{pk_0 k_0 m} (\hat{\mathbf{z}}_{F6}^{(k_0)}) \propto r_{pk_0 k_0 k_0} (\hat{\mathbf{z}}) \quad (68)$$

$$= r_{pk_0 k_0 k_0} (\boldsymbol{\rho}(\hat{\mathbf{u}}_{\text{syn}})) r_{pk_0 k_0 k_0} (\boldsymbol{\rho}(\hat{\mathbf{u}}_{\text{one}})) = 0 \quad (69)$$

581 since $r_{pk_0 k_0 k_0} (\boldsymbol{\rho}(\hat{\mathbf{u}}_{\text{syn}})) = 0$. Similarly for $r_{p, -k_0, -k_0, m}$. \square

582 C.3 Canonical Forms

583 **Definition 8.** A solution \mathbf{z} is called canonical at k_0 , or $\mathbf{z} \in \mathcal{C}_{k_0}$, if $z_{pk_0} = 1$ for all p and $k = \pm k_0$.

584 **Lemma 4** (Canonical Decomposition). Any solution \mathbf{z} with $r_{k_0 k_0 k_0} (\mathbf{z}) \neq 0$ can be decomposed into
 585 $\mathbf{z} = \mathbf{z}' * \mathbf{y}$, where \mathbf{z}' is canonical at k_0 and $\text{ord}(\mathbf{y}) = 1$. Both $r_{k_0 k_0 k_0} (\mathbf{z}') \neq 0$ and $r_{k_0 k_0 k_0} (\mathbf{y}) \neq 0$.

586 *Proof.* Since $r_{k_0 k_0 k_0} (\mathbf{z}) = \sum_j a_{k_0 j} b_{k_0 j} c_{k_0 j} \neq 0$, there must exist some j so that $z_{a_{k_0 j}} z_{b_{k_0 j}} z_{c_{k_0 j}} \neq$
 587 0, which means that $z_{a_{k_0 j}} \neq 0$, $z_{b_{k_0 j}} \neq 0$ and $z_{c_{k_0 j}} \neq 0$. Since the node index j can be permuted,

588 we can let node j be the first node 0 and let $y_{pk0} = z_{pkj}$ and $z'_{pkj'} = z_{pkj} z_{pkj}^{-1}$ for $p \in \{a, b, c\}$ and
 589 $k = \pm k_0$, then \mathbf{z}' is canonical at k_0 and $\text{ord}(\mathbf{y}) = 1$. Finally, by ring homomorphism, since

$$r_{k_0 k_0 k_0}(\mathbf{z}) = r_{k_0 k_0 k_0}(\mathbf{z}') r_{k_0 k_0 k_0}(\mathbf{y}) \neq 0 \quad (70)$$

590 we know that both $r_{k_0 k_0 k_0}(\mathbf{z}') \neq 0$ and $r_{k_0 k_0 k_0}(\mathbf{y}) \neq 0$. \square

591 D Gradient Dynamics (Proof)

592 **Theorem 5** (Dynamics of MPs). *The dynamics of MPs satisfies $\dot{\mathbf{r}} = -JJ^* \overline{\nabla_{\mathbf{r}} \ell}$, which has positive*
 593 *inner product with the negative gradient direction $-\overline{\nabla_{\mathbf{r}} \ell}$.*

594 *Proof.* By gradient descent of \mathcal{W} , we have $\dot{\mathcal{W}} = -\overline{\nabla_{\mathcal{W}} \ell}$. By chain rule, we have:

$$\dot{\mathcal{W}} = -\overline{\nabla_{\mathcal{W}} \ell} = -\overline{J^T \nabla_{\mathbf{r}} \ell} = -J^* \overline{\nabla_{\mathbf{r}} \ell} \quad (71)$$

595 Then the dynamics of $\mathbf{r} = \mathbf{r}(z(\mathcal{W}))$, as driven by the dynamics of \mathcal{W} , is given by

$$\dot{\mathbf{r}} = J \dot{\mathcal{W}} = -JJ^* \overline{\nabla_{\mathbf{r}} \ell} \quad (72)$$

596 To show positive inner product, we have:

$$-\overline{\nabla_{\mathbf{r}} \ell}^* \dot{\mathbf{r}} = \overline{\nabla_{\mathbf{r}} \ell}^* JJ^* \overline{\nabla_{\mathbf{r}} \ell} = \|J^* \overline{\nabla_{\mathbf{r}} \ell}\|_2^2 \geq 0 \quad (73)$$

597 \square

598 **Theorem 6** (The Occam's Razer: Preference of low-order solutions). *If $\mathbf{z} = \mathbf{y} * \mathbf{z}'$ and both \mathbf{z} (of*
 599 *order q) and \mathbf{z}' are global optimal solutions, then there exists a path of zero loss connecting \mathbf{z} and \mathbf{z}'*
 600 *in the space of \mathcal{Z}_q . As a result, lower-order solutions are preferred if trained with L_2 regularization.*

601 *Proof.* Let $\text{ord}(\mathbf{z}) = q$ and $\text{ord}(\mathbf{z}') = q'$. Then $q' | q$. Since both \mathbf{z} and \mathbf{z}' are global optimal. Since
 602 r_{kkk} is ring homomorphism, we know that $r_{kkk}(\mathbf{z}) = r_{kkk}(\mathbf{z}') r_{kkk}(\mathbf{y}) = 1/2d = r_{kkk}(\mathbf{z}')$ and
 603 thus $r_{kkk}(\mathbf{y}) = 1$ for all $k \neq 0$.

604 Let the augmented identity $\mathbf{e} \in \mathcal{Z}_q$ be $e_{pmj} = \mathbb{1}(j = 0)$. Then $r_{kkk}(\mathbf{e}) = 1$ for all $k \neq 0$.

605 We want to construct a path in \mathcal{Z}_q , the space of order- q solutions as follows:

$$\tilde{\mathbf{z}}(t) = \tilde{\mathbf{y}}(t) * \mathbf{z}', \quad 0 \leq t \leq 1 \quad (74)$$

606 in which $\tilde{\mathbf{y}}(0) = \mathbf{e}$, $\tilde{\mathbf{y}}(1) = \mathbf{y}$, and $r_{kkk}(\tilde{\mathbf{y}}(t)) = 1$ for any t . To see why this is possible, pick a
 607 continuous family of trajectories $\hat{\mathbf{y}}(t; \lambda)$ with $\lambda \in [0, 1]$ so that they satisfies

$$\hat{\mathbf{y}}(0; \lambda) = \mathbf{e}, \quad \hat{\mathbf{y}}(1; \lambda) = \mathbf{y}, \quad r_{kkk}(\hat{\mathbf{y}}(t; 0)) \leq 1, \quad r_{kkk}(\hat{\mathbf{y}}(t; 1)) \leq 1 \quad (75)$$

608 which can always be achieved by scaling some trajectory with a factor that depends on λ . Then
 609 by intermediate theorem, there exists $\lambda(t)$ so that $r_{kkk}(\hat{\mathbf{y}}(t; \lambda(t))) = 1$ for some k . Note that for
 610 different frequency k and k' , r_{kkk} and $r_{k'k'k'}$ involves disjoint components of \mathbf{z} so we could find
 611 such a path for all $k \neq 0$.

612 Therefore, for any monomial potential r included in MSE loss (Eqn. 3), we have

$$r(\tilde{\mathbf{z}}(t)) = r(\tilde{\mathbf{y}}(t)) r(\mathbf{z}') = \begin{cases} \text{finite} \cdot 0 = 0 & r \neq r_{kkk} \\ 1 \cdot 1/2d = 1/2d & r = r_{kkk} \end{cases} \quad (76)$$

613 and thus the entire trajectory $\tilde{\mathbf{z}}(t) = \tilde{\mathbf{y}}(t) * \mathbf{z}' \in \mathcal{Z}_q$ connecting \mathbf{z} and $\mathbf{e} * \mathbf{z}'$, which is \mathbf{z}' in the space
 614 of \mathcal{Z}_q , is also globally optimal.

615 To see why weight decay regularization leads to lower-order solution, we could simply compare the
 616 ℓ_2 norm of $\mathbf{z} = \mathbf{y} * \mathbf{z}'$ and $\mathbf{e} * \mathbf{z}'$. At each frequency k , this reduces to the following optimization
 617 problem:

$$\min \sum_j |a_j|^2 + |b_j|^2 + |c_j|^2, \quad \text{s.t.} \sum_j a_j b_j c_j = 1 \quad (77)$$

618 where $a_j := y_{akj}$, $b_j := y_{bkj}$ and $c_j := y_{ckj}$. Since we know that arithmetic mean is no less than
 619 geometric mean:

$$\frac{|a_j|^2 + |b_j|^2 + |c_j|^2}{3} \geq \sqrt[3]{|a_j b_j c_j|^2} \quad (78)$$

620 We have:

$$\sum_j |a_j|^2 + |b_j|^2 + |c_j|^2 \geq 3 \sum_j |a_j b_j c_j|^{2/3} \geq 3 \quad (79)$$

621 The last inequality holds because (1) if any $|a_j b_j c_j| \geq 1$, then it holds, (2) if all $|a_j b_j c_j| < 1$, then
 622 since a^x is a decreasing function for $a < 1$, $\sum_j |a_j b_j c_j|^{2/3} \geq \sum_j |a_j b_j c_j| \geq |\sum_j a_j b_j c_j| = 1$.

623 The minimizer is reached when $|a_j| = |b_j| = |c_j|$. Note that if $a_j b_j c_j$ has any complex phase or
 624 negative, then in order to satisfy $\sum_j a_j b_j c_j = 1$, objective function needs to be larger. So without
 625 loss of generality, we could study $a_j = b_j = c_j = x_j \geq 0$ and the optimization problem becomes

$$\min \sum_j x_j^2, \quad \text{s.t.} \quad \sum_j x_j^3 = 1, \quad x_j \geq 0 \quad (80)$$

626 which has a minimizer at the corners $(1, 0, \dots)$. This corresponds to $a_j = b_j = c_j = \mathbb{1}(j = 0)$,
 627 which is the augmented identity $e \in \mathcal{Z}_q$. \square

628 **Theorem 7** (Infinite Width Limits at Initialization). *Considering the modified loss of Eqn. 3 with
 629 only the first two terms: $\tilde{\ell}_k := -2r_{kkk} + \sum_{k_1 k_2} |r_{k_1 k_2 k}|^2$, if the weights are i.i.d Gaussian and
 630 network width $q \rightarrow +\infty$, then JJ^* converge to diagonal and the dynamics of MPs is decoupled.*

631 *Proof.* Let $\tilde{\ell} := \sum_k \nabla \tilde{\ell}_k$. Let's compute the dynamics of MPs following Theorem 5: $\dot{r} =$
 632 $-JJ^* \nabla_r \tilde{\ell}$.

633 First it is clear that

$$\frac{\partial \tilde{\ell}}{\partial r_{k_1 k_2 k}} = \sum_k \frac{\partial \tilde{\ell}_k}{\partial r_{k_1 k_2 k}} = -2\mathbb{1}(k_1 = k_2 = k) + 2\overline{r_{k_1 k_2 k}} \quad (81)$$

634 So the (k_1, k_2, k) component of $\nabla_r \tilde{\ell}$ only contains $r_{k_1 k_2 k}$.

635 Then we compute $H := JJ^*$ and show that it is asymptotically diagonal. To see this, each compo-
 636 nent of H , i.e., $h_{k_1 k_2 k_3, k'_1 k'_2 k'_3}$ can be computed as the following:

$$h_{k_1 k_2 k_3, k'_1 k'_2 k'_3} = \sum_{pmj} \frac{\partial r_{k_1 k_2 k_3}}{\partial z_{pmj}} \overline{\frac{\partial r_{k'_1 k'_2 k'_3}}{\partial z_{pmj}}} \quad (82)$$

$$= \mathbb{1}(k_1 = k'_1) \sum_j b_{k_2 j} \bar{b}_{k'_2 j} c_{k_3 j} \bar{c}_{k'_3 j} \quad (83)$$

$$+ \mathbb{1}(k_2 = k'_2) \sum_j a_{k_1 j} \bar{a}_{k'_1 j} c_{k_3 j} \bar{c}_{k'_3 j} \quad (84)$$

$$+ \mathbb{1}(k_3 = k'_3) \sum_j a_{k_1 j} \bar{a}_{k'_1 j} b_{k_2 j} \bar{b}_{k'_2 j} \quad (85)$$

637 where $a_{kj} := z_{akj}$, $b_{kj} := z_{bkj}$ and $c_{kj} := z_{ckj}$. Then for component $(k_1 k_2 k_3, k'_1, k'_2, k'_3)$, if any
 638 $k_p \neq k'_p$ for some $p \in \{a, b, c\}$, then the corresponding $z_{pk_p j} \bar{z}_{pk'_p j}$ has random phase for hidden
 639 node j , and $h_{k_1 k_2 k_3, k'_1 k'_2 k'_3} \rightarrow 0$ when $q \rightarrow +\infty$.

640 Combining the two, we know that the dynamics of MPs is decoupled, that is, each $r_{k_1 k_2 k}$ evolves
 641 independently over time. \square

642 **Ripple effects.** While Theorem 7 only holds at initialization, the resulting decoupled MP dynamics,
 643 e.g., $dr_{kkk}/dt = 1 - r_{kkk}$ that leads to $r_{kkk}(t) = 1 - e^{-t}$, already captures the rough shape of
 644 the curve (Fig. 3 top right). To capture its fine structures (e.g., ripples before stabilization), we can
 645 also model the dynamics of the diagonal element in JJ^* . Consider a symmetric 1D case on a fixed

646 frequency k , where all diagonal $r_{kkk} = r_0 - r$ (where $r_0 = 1/2d$) and all off-diagonal $r_{k_1k_2k} = r$,
 647 then

$$\dot{r} = -\dot{r}_{kkk} = \kappa(r_{kkk} - r_0) = -\kappa r, \quad \dot{\kappa} = \alpha(r_0 - r_{kkk}) - (1 - \alpha)r_{k_1k_2k} - c_0 = (2\alpha - 1)r - c_0 \quad (86)$$

648 where $\kappa > 0$ is the diagonal element of JJ^* and α is a coefficient that characterizes the relative
 649 strength of two negative gradient $-\nabla_{r_{kkk}}\ell = r_0 - r_{kkk}$ and $-\nabla_{r_{k_1k_2k}}\ell = -r_{k_1k_2k}$, and c_0 is the
 650 gradient terms caused by asymmetry and/or other frequencies. This yields a second-order ODE that
 651 has complex roots in the characteristic function when $c_0 > 0$.

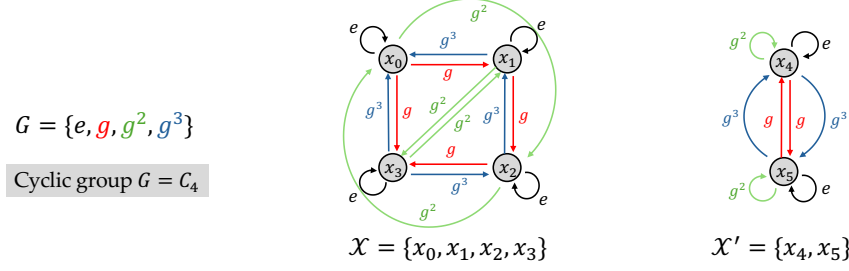


Figure 8: An example case of group action on state set \mathcal{X} , \mathcal{X} can be partitioned into several disjoint components, each is a transitive graph w.r.t the group actions in G .

652 E Extending CoGO to Group Action Prediction

653 While in this work we mainly focus on Abelian group, CoGO can be extended to more general *group*
 654 *action prediction*: given a group element $g \in G$ and the current state $x \in \mathcal{X}$, the goal is to predict
 655 $gx \in \mathcal{X}$, i.e., the next state after action g . Such tasks include modular addition/multiplication in
 656 which the group acts on itself (i.e., $\mathcal{X} = G$), and also includes the transition function in reinforce-
 657 ment learning (Sutton, 2018) and world modeling (Garrido et al., 2024), in which an action changes
 658 the current state to a new one.

659 **Setup.** Consider a state space \mathcal{X} and *group action* $G \times \mathcal{X} \mapsto \mathcal{X}$ where $g \in G$ is a group element
 660 acting on a state $x \in \mathcal{X}$ to get an update state $gx \in \mathcal{X}$. It satisfies two axioms (1) the group identity
 661 maps everything to itself: $ex = x$, and (2) the group action is compatible with group multiplication:
 662 $g(hx) = (gh)x$ for any $g, h \in G$ and $x \in \mathcal{X}$.

663 Equipped with the group action, the state space now can be decoupled into a disjoint of *transitive*
 664 *components*.

665 **Definition 9** (Transitive group action). *A group action is transitive, if for any $x_1, x_2 \in \mathcal{X}$, there*
 666 *exists $g \in G$ so that $gx_1 = x_2$.*

667 Since the group action is compatible with multiplication, \mathcal{X} under G will be partitioned into disjoint
 668 components $\mathcal{X} = \bigcup_l \mathcal{X}_l$ and we can analyze each component separately (Fig. 8).

669 **Transitive Group Action.** For each transitive component \mathcal{X} (dropping l for brevity), under certain
 670 conditions, we could define a *state multiplication* operation (a formal definition in Def. 10 in Ap-
 671 pendix) so that for any group action $gx \in \mathcal{X}$, there is an associated state $x' \in \mathcal{X}$ so that $x' \cdot x = gx$.
 672 Furthermore, under the multiplication, \mathcal{X} itself becomes a group:

673 **Theorem 8** ($\mathcal{X} \cong G/G_{x_0}$). *If the group stabilizer $G_{x_0} := \{g | gx_0 = x_0\}$ is a normal subgroup of*
 674 *G , then \mathcal{X} is isomorphic to the quotient group G/G_{x_0} and thus forms a group.*

675 Moreover, we can prove that for any group element $g \in G$, there exists $x = \iota_0(g) \in \mathcal{X}$ so that for
 676 any state x' , the group action gx' is the same as the state multiplication $x' \cdot x$. Therefore, for group
 677 action prediction tasks, we have (note the difference compared to Eqn. 12):

$$\mathbf{w}_j = U_G \left(P_0 \mathbf{w}_{j,G}^{\parallel} + \mathbf{w}_{j,G}^{\perp} \right) + U_{\mathcal{X}} \mathbf{w}_{j,\mathcal{X}} \quad (87)$$

678 where $\mathbf{w}_{j,G}^{\parallel} \in \mathbb{R}^{|\mathcal{X}|}$ is the “in-graph” component of G , $\mathbf{w}_{j,G}^{\perp} \in \mathbb{R}^{|G|}$ is the “out-of-graph” component
 679 of G , and $P_0 \in \mathbb{R}^{|G| \times |\mathcal{X}|}$ “lifts” from \mathcal{X} to G using ι_0 , i.e., $(P_0)_{gx} = 1$ for $g \in \iota_0^{-1}(x)$, and
 680 $\mathbf{w}_{j,G}^{\perp} \perp P_0 \mathbf{w}_{j,G}^{\parallel}$. Since any g just behaves like $\iota_0(g)$ when acting on \mathcal{X} , our framework can be
 681 applied to characterize the learning of $\mathbf{w}_{j,G}^{\parallel}$. Intuitively, we only learn representation of G ’s element
 682 “module” its kernel G_{x_0} , since element in the kernel is indistinguishable from each other.

683 On the other hand, the behavior of $\mathbf{w}_{j,G}^{\perp}$ will be influenced by g acting on other graphs, and the final
 684 learned representation of a group element g is the direct sum of them.

685 **F Detailed explanation of Sec. E**

686 **Matrix Representation.** Each group element g can be represented by a matrix R_g , i.e., its *matrix*
 687 *representation*, so that it respects the group multiplication (i.e., *homomorphism*): $R_{gh} = R_g R_h$ for
 688 any group elements $g, h \in G$.

689 The dimension of such a representation may differ widely. Some representation can be 1-
 690 dimensional (e.g., for Abelian group), while others can be infinitely dimensional. The *permutation*
 691 *representation* $R_g \in \mathbb{R}^{d \times d}$ maps a one-hot representation $e_x \in \mathbb{R}^d$ of an object \mathcal{X} into its image
 692 $e_{gx} \in \mathbb{R}^d$, also a one-hot representation. Intuitively, $(R_g)_{jk} = 1$ means that it maps the k -th element
 693 into the j -th element.

694 **Lemma 5** (Structure of R_g). *For any $g \in G$, R_g is a permutation matrix.*

695 **Lemma 6** (Summation of R_g). *If the group action is transitive, then $\sum_{g \in G} R_g = \frac{|G|}{d} \mathbf{1}\mathbf{1}^\top$.*

696 **F.1 Transitive Case**

697 To construct the multiplication operation on \mathcal{X} , we first pick reference point $x_0 \in \mathcal{X}$, and establish
 698 a mapping $\iota_0 : G \mapsto \mathcal{X}$: $\iota_0(g) = gx_0$. Note that ι_0 is not necessarily a bijection; in fact we have:

699 **Lemma 7** (Co-set Mapping ι_0). *There is a bijection between $\{\iota_0^{-1}(x)\}_{x \in \mathcal{X}}$ and co-sets $[G : G_{x_0}]$
 700 of group stabilizer $G_{x_0} := \{g \in G | gx_0 = x_0\}$, which is a subgroup of G fixing x_0 .*

701 **Lemma 8** (Uniqueness of Multiplication Mapping). *If G_{x_0} is a normal subgroup, then for all $g_1 \in$
 702 $\iota_0^{-1}(x_1)$ and $g_2 \in \iota_0^{-1}(x_2)$, all $g_1 g_2 G_{x_0}$ correspond to the same coset.*

703 **Definition 10** (The multiplication operator on \mathcal{X}). *When G_{x_0} is a normal subgroup, we define
 704 multiplication on \mathcal{X} : $\mathcal{X} \times \mathcal{X} \mapsto \mathcal{X}$ to be $x_1 x_2 := \iota_0(g_1 g_2 G_{x_0})$ for $x_1 = g_1 x_0$ and $x_2 = g_2 x_0$.
 705 Under this definition, x_0 is the identity element.*

706 **Lemma 9.** *If $g \in \iota_0^{-1}(x)$, then for any $x' \in \mathcal{X}$, $gx' = xx'$.*

707 This means that in terms of group action, the group element g is indistinguishable to x on \mathcal{X} .

708 **F.2 General group action**

709 In this case, R_g can be decomposed into a direct sum of smaller matrices, and all our analysis applies
 710 to each of these small matrices.

711 In the main text, to simplify the notation, we assume that the group action is transitive, i.e., for any
 712 $y, y' \in Y$, there exists $g \in G$ so that $gy = y'$. In the following we will show that for general group
 713 actions, the conclusion still follows.

714 *Group orbit.* For any $x \in \mathcal{X}$, Let $G \cdot y := \{gy | g \in G\} \subseteq Y$ be its *orbit*.

715 **Lemma 10.** *For $y, y' \in G$, either $G \cdot y = G \cdot y'$ (two orbits collapse) or $G \cdot y \cap G \cdot y' \neq \emptyset$ (two
 716 orbits are disjoint). Therefore, orbits form a partition of \mathcal{X} .*

717 Let $X/G := \{G \cdot y | x \in \mathcal{X}\}$ be the collection of all orbits. The following lemma tells that the matrix
 718 representation R_g can be decomposed into a direct sum (i.e., block diagonal matrix) on each orbit.

Lemma 11 (Direct sum decomposition of R_g).

$$R_g = \bigoplus_{Y' \in Y/G} R_g^{Y'} \quad (88)$$

719 *and each $R_g^{Y'} \in \mathbb{R}^{|Y'| \times |Y'|}$ is a permutation matrix with $\sum_g R_g^{Y'} = \frac{|G|}{|Y'|} \mathbf{1}\mathbf{1}^\top$.*

720 *Proof.* By the definition of group orbits, the group action g is closed within each Y' . Therefore, R_g
 721 is a direct sum (i.e., block-diagonal).

722 For each element $x \in \mathcal{X}'$, let's check its destination under G . It is clear that if two group elements
 723 $g, h \in G$ maps \mathcal{X} to the same destination, then

$$gy = hy \iff y = g^{-1}hy \iff g^{-1}h \in G_y \iff h = gG_y \quad (89)$$

724 where G_y is the stabilizer of \mathcal{X} , a subgroup of G . Therefore, g and h map \mathcal{X} to the same destination,
 725 if and only if they are from the same coset of G_y . Therefore, each entry of $\sum_g R_g^{Y'}$ on the column
 726 \mathcal{X} equals to the size of cosets of G_y , which is $|G_y|$. Furthermore, for $y_1, y_2 \in Y'$, since they belong
 727 to the same orbit, there exists g so that $gy_1 = y_2$ and thus for any $g' \in G_{y_1}$, we have

$$g'y_1 = y_1 \iff gg'y_1 = gy_1 = y_2 \iff gg'g^{-1}y_2 = y_2 \iff gg'g^{-1} \in G_{y_2} \quad (90)$$

728 So there exists bijection between G_{y_1} and G_{y_2} . This means that $|G_y|$ is constant for any $x \in \mathcal{X}'$ and
 729 thus all elements in $\sum_g R_g^{Y'}$ are equal to $|G|/|Y'|$ (i.e., the number of the group elements that send
 730 \mathcal{X} out to various destinations in Y' , divided by the possible distinct destinations $|Y'|$, results in the
 731 number of times each destination gets hit). \square

732 G Proofs for the content in Appendix

733 **Lemma 5** (Structure of R_g). *For any $g \in G$, R_g is a permutation matrix.*

734 *Proof.* Since every element needs to have a destination, every column of R_g sums to 1, i.e., $\mathbf{1}^\top R_g =$
 735 $\mathbf{1}^\top$. Then we prove that the mapping $y \mapsto gy$ is a bijection. Suppose there exists y_1, y_2 so that
 736 $gy_1 = gy_2$. Therefore by compatibility we have:

$$g^{-1}(gy_1) = g^{-1}(gy_2) \iff (g^{-1}g)y_1 = (g^{-1}g)y_2 \iff ey_1 = ey_2 \iff y_1 = y_2 \quad (91)$$

737 So any g is a bijective mapping on \mathcal{X} . Since every element of R_g is either 0 or 1, R_g is a permutation
 738 matrix. \square

739 **Lemma 6** (Summation of R_g). *If the group action is transitive, then $\sum_{g \in G} R_g = \frac{|G|}{d} \mathbf{1}\mathbf{1}^\top$.*

740 *Proof.* Simply apply Lemma 11 and notice that for transitive group action, $X/G = \{Y\}$. \square

741 **Lemma 7** (Co-set Mapping ι_0). *There is a bijection between $\{\iota_0^{-1}(x)\}_{x \in \mathcal{X}}$ and co-sets $[G : G_{x_0}]$
 742 of group stabilizer $G_{x_0} := \{g \in G | gx_0 = x_0\}$, which is a subgroup of G fixing x_0 .*

743 *Proof.* First we have

$$\iota_0(g) = \iota_0(h) \iff gy_0 = hy_0 \iff y_0 = g^{-1}hy_0 \iff g^{-1}h \in G_{y_0} \iff h \in gG_{y_0} \quad (92)$$

744 So for any $y = gy_0$, all elements in $\iota_0^{-1}(y)$ are also in gG_{y_0} and vice versa. The bijection is:

$$\iota_0^{-1}(y) \leftrightarrow gG_{y_0}, \quad \text{for } y = gy_0 \quad (93)$$

745 or equivalently,

$$y \leftrightarrow \iota_0(gG_{y_0}) \quad (94)$$

746 \square

747 **Lemma 9.** *If $g \in \iota_0^{-1}(x)$, then for any $x' \in \mathcal{X}$, $gx' = xx'$.*

748 *Proof.* For $g \in \iota_0^{-1}(x)$, we have $gx_0 = x$. For any $x' = hx_0$, we have:

$$gx' = ghx_0 = (gh)x_0 \quad (95)$$

749 On the other hand, by definition, $xx' := \iota_0(ghG_{x_0}) = (gh)x_0$. So for any x' , $gx' = xx'$. \square

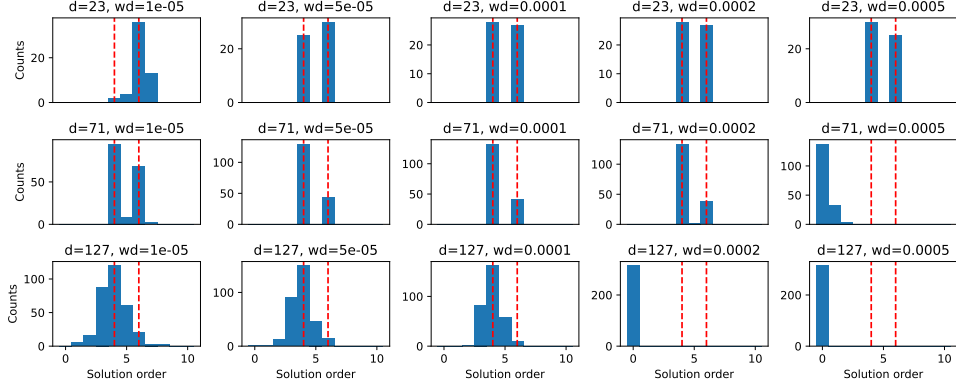


Figure 9: Distribution of solutions with hidden size $q = 256$.

750 **H Additional Experiments**

751 **Algorithm to extract factorization from gradient descent solutions.** Given the solutions obtained
 752 by gradient descent using Adam optimizer, we first compute the corresponding z via the Fourier
 753 transform (that is, Eqn. 12). Here $z = [z_{pkj}]$ is a 3-by- d -by- q tensor. Here $d = |G|$ and q is the
 754 number of hidden nodes in the 2-layer neural networks.

755 Then for each frequency k , we extract the salient components of z by thresholding with a universal
 756 threshold (e.g. 0.05). The number of salient components (e.g., 6 or 4) is the order of the per-
 757 frequency solution.

758 Suppose we now get $z^{(k)}$ for frequency k , which is a 3-by-6 (and thus an order-6) solution. Then
 759 we enumerate all possible permutation of 6 hidden nodes ($6! = 720$ possibilities) to find one permu-
 760 tation τ so that $\|z_{pk\tau(\cdot)} - z_{pk\cdot}^{(1)} \otimes z_{pk\cdot}^{(2)}\|$ is minimized, following ring multiplication defined in Def. 5.
 761 Note that for each permutation, we also need to consider whether $\tilde{\mathbf{1}} := [-1, -1, 1]$ can be applied to
 762 each hidden node j ($\tilde{\mathbf{1}}$ is also defined in Tbl. 1). This is because both $z_1 + z_2$ and $z_1 + \tilde{\mathbf{1}} * z_2$ have
 763 exactly the same values on all monomial potentials (MPs) we consider, due to the fact that $r(\tilde{\mathbf{1}}) = 1$
 764 for any $r \in R_g \cup R_c \cup R_n \cup R_*$. Therefore we call $\tilde{\mathbf{1}}$ “pseudo-1”.

765 For search efficiency, we therefore first consider the permutation τ so that $\|z_{ck\tau(\cdot)} - z_{ck\cdot}^{(1)} \otimes z_{ck\cdot}^{(2)}\|$ is
 766 minimized, since the component c is invariant to the pseudo-1 transformation $\tilde{\mathbf{1}}$, and then for those
 767 eligible τ , we search whether $\tilde{\mathbf{1}}$ should be applied when considering $p \in \{a, b\}$.

768 Once we find such z_1 and z_2 , we convert them into their canonical forms \tilde{z}_1 and \tilde{z}_2 (Def. 8) to
 769 eliminate any possible multiplicative term y so that $z_1 = y * \tilde{z}_1$. We then compare the canonical
 770 forms (up to complex conjugate) with various order-3 and order-2 partial solutions constructed by
 771 CoGO, as detailed in Sec. 5. If their distance is below a certain threshold (e.g., $< 10\%$ of the norm
 772 after normalizing both \hat{z}_1 and \hat{z}_2), then a match is detected.

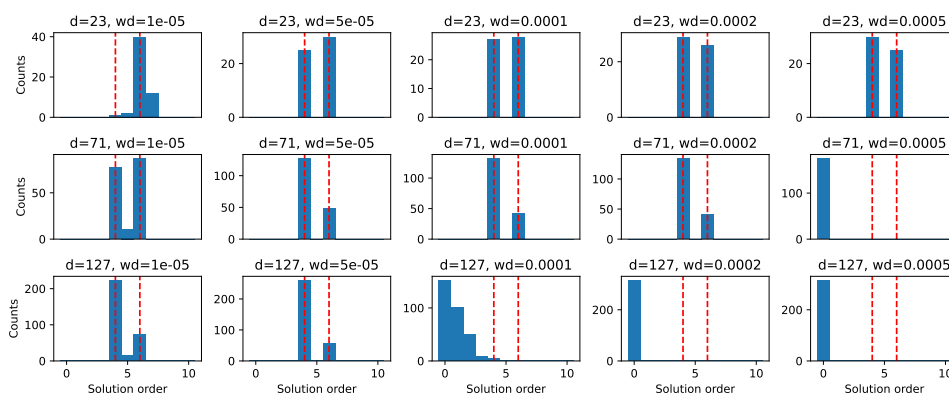


Figure 10: Distribution of solutions with hidden size $q = 512$.

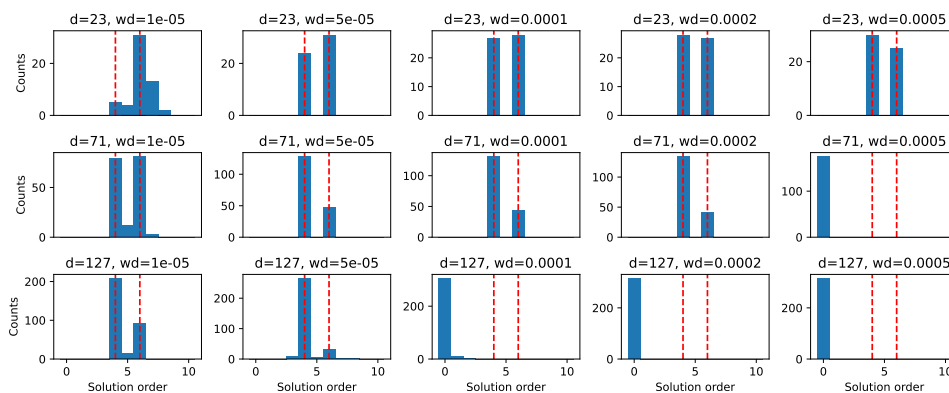


Figure 11: Distribution of solutions with hidden size $q = 1024$.

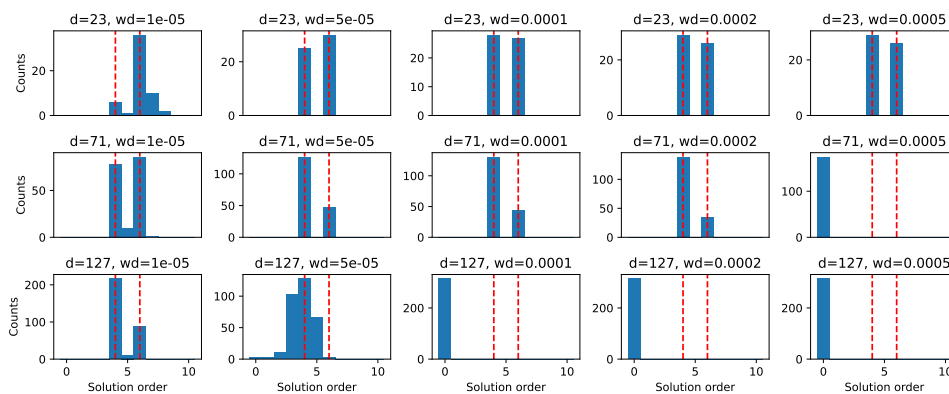


Figure 12: Distribution of solutions with hidden size $q = 2048$.