
Real-Valued Backpropagation is Unsuitable for Complex-Valued Neural Networks

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Recently complex-valued neural networks have received increasing attention due
2 to successful applications in various tasks and the potential advantages of better
3 theoretical properties and richer representational capacity. However, the training
4 dynamics of complex networks compared to real networks remains an open prob-
5 lem. In this paper, we investigate the dynamics of deep complex networks during
6 real-valued backpropagation in the infinite-width limit via neural tangent kernel
7 (NTK). We first extend the Tensor Program to the complex domain, to show that
8 the dynamics of any basic complex network architecture is governed by its NTK
9 under real-valued backpropagation. Then we propose a way to investigate the
10 comparison of training dynamics between complex and real networks by studying
11 their NTKs. As a result, we surprisingly prove that for most complex activation
12 functions, the commonly used real-valued backpropagation reduces the training
13 dynamics of complex networks to that of ordinary real networks, thus eliminating
14 the characteristics of complex-valued neural networks. Finally, we study the results
15 numerically and the experiments verify our theoretical findings.

16 1 Introduction

17 Recently complex-valued neural networks have been successfully applied to various tasks, such
18 as time-series prediction [Wisdom et al., 2016], computer vision [Trabelsi et al., 2018], signal
19 processing [Yao et al., 2020]. Compared to real-valued neural networks, it is shown that complex
20 networks have the potential to provide richer representational capacity [Arjovsky et al., 2016], faster
21 learning [Danilhelka et al., 2016], better motivation and generalization for signal-related tasks [Hirose
22 and Yoshida, 2012, Tygert et al., 2016]. Theoretically, there have been significant advances for
23 complex networks regarding the universal approximation property [Voigtlaender, 2020], critical
24 points [Nitta, 2002], local minima [Nitta, 2013] and separation results [Zhang et al., 2021].

25 However, training deep complex networks has been challenging because of several non-intuitive
26 analytical properties of complex algebra. Firstly, in practice, we often deal with a real-valued cost
27 function, which is non-analytic with respect to complex-valued parameters. Secondly, Liouville’s
28 theorem asserts that every bounded and complex-differentiable function is a constant. Thus almost
29 all activation functions are non-analytic due to the preference for boundedness before the popularity
30 of ReLU [Scardapane et al., 2018]. Moreover, Voigtlaender [2020] has proved that the universal
31 approximation property holds only when using non-holomorphic activation functions.

32 Due to these reasons, different backpropagation algorithms in the complex domain were independently
33 proposed for non-holomorphic networks [Basseley et al., 2021], mostly by optimizing the real and
34 imaginary components separately. Recently, real-valued backpropagation [Nitta, 1997] is widely
35 used due to the convenience of utilizing the real-valued deep learning library [Arjovsky et al., 2016,
36 Trabelsi et al., 2018, Tan et al., 2020]. It optimized the complex network just like a real network by

37 computing partial derivatives of the cost with respect to the real and imaginary parts separately and
38 achieved state-of-the-art performance.

39 As a result, a natural and fundamental problem in complex-valued neural networks arises: Do complex
40 networks have a different inductive bias from real networks during training? Do complex networks
41 trained by gradient descent tend to learn different hypotheses from real networks? However, to the
42 best of the authors’ knowledge, the training dynamics of backpropagation for complex networks
43 compared to real networks remains open.

44 In this paper, we investigate the training dynamics of deep complex networks under real-valued
45 backpropagation from neural tangent kernel (NTK) perspective [Jacot et al., 2018], which captures
46 the optimization behavior of neural networks in the infinite-width limit. Then we provide a way to
47 investigate the comparison of training dynamics between complex and real networks via their NTKs.
48 As a result, we obtain informative results which may guide the algorithm selection and complex
49 network design in practice if people want to take full advantage of complex networks.

50 **Our contributions.** Our main contributions can be summarized as follows:

- 51 • First, we extend the Tensor Program [Yang, 2020] to the complex domain and show that for
52 a complex-valued neural network of any basic architecture in the infinite-width limit, the
53 training dynamics of real-valued backpropagation is determined by kernel gradient descent
54 with a deterministic NTK at initialization.
- 55 • Second, we investigate the comparison of training dynamics between complex and real
56 networks based on their NTKs. We surprisingly prove that the commonly used real-valued
57 backpropagation reduces the training dynamics of complex-valued multi-layer perceptrons
58 (MLPs) to that of ordinary real MLPs, thus eliminating the characteristics of complex-valued
59 neural networks. This result holds for most commonly used complex activation functions,
60 including all split activation functions such as $\mathbb{C}\text{ReLU}$, $\mathbb{C}\text{Sigmoid}$, $\mathbb{C}\text{tanh}$; and part of
61 magnitude-based holomorphic activation functions.
- 62 • Finally, we study the results numerically, and the experiments verifies our findings. Specifi-
63 cally, in several settings with different depths and various activation functions, the NTKs of
64 complex networks converges to the NTKs of real networks as the widths grow.

65 **Organization.** We start with some preliminaries and notations about complex networks and neural
66 tangent kernels in Section 2. In Section 3, we investigate the NTK of complex-valued neural networks
67 of any architecture during real-valued backpropagation in the infinite-width limit. Section 4 firstly
68 presents the NTK of complex MLPs in any depth and then investigates the conditions that training
69 dynamics of complex-valued MLPs reduce to that of ordinary real MLPs. We verified our results
70 empirically in Section 5. Finally we discuss the related works and conclude the paper. Due to the
71 limited space, all proofs are placed in the appendices.

72 2 Preliminaries

73 2.1 Complex-Valued Neural Networks

74 Without loss of generality, we focus on complex-valued neural networks with real-valued output
75 $f_\theta(\mathbf{z}) \in \mathbb{R}^{d_{out}}$ with parameter set $\theta \in \mathbb{C}^p$, input $\mathbf{z} \in \mathbb{C}^d$ and $\mathbf{z} = \mathbf{x} + \mathbf{y}i$ with $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$, and trained
76 by real-valued backpropagation algorithm (also called complex-BP or generalized complex BP, see
77 Appendix B for more details), which is conventional in the literature [Arjovsky et al., 2016, Wisdom
78 et al., 2016, Trabelsi et al., 2018, Zhang et al., 2021, Wu et al., 2021]. Our analysis can be naturally
79 applied to complex-valued output by decomposing the real and imaginary part of output into two
80 functions, or applied to real-valued input by treating its imaginary part as zero.

For an L -hidden layer complex network, we denote the output of last hidden layer as $\mathbf{h}_L \in \mathbb{C}^n$.
Without loss of generality, we consider that the output of a complex network with a linear readout
layer is achieved via

$$f_\theta(\mathbf{z}) = \Re\{\mathbf{W}_{L+1}\mathbf{h}_L\}$$

81 where $\mathbf{W}_{L+1} \in \mathbb{C}^{n \times d_{out}}$ [Wisdom et al., 2016, Zhang et al., 2021]. Note that there are other two
82 common forms of linear readout layer to generate real-valued output, like $f_\theta(\mathbf{z}) = \mathbf{W}_{L+1}\Re\{\mathbf{h}_L\}$

83 where the output weight $\mathbf{W}_{L+1} \in \mathbb{R}^{n \times d_{out}}$ [Wu et al., 2021], and $f_\theta(\mathbf{z}) = \mathbf{W}_{L+1} \begin{bmatrix} \Re(\mathbf{h}_L) \\ \Im(\mathbf{h}_L) \end{bmatrix}$
84 where $\mathbf{W}_{L+1} \in \mathbb{R}^{2n \times d_{out}}$ [Arjovsky et al., 2016]. They can be treated as special cases of our settings.
85 For a complex-valued neural network $f_\theta(\mathbf{z})$, we can always decompose all the complex operations
86 into two-dimensional real-valued operations, denoted as $f_{[\theta_R, \theta_I]}([\mathbf{x}, \mathbf{y}])$ where $\theta_R, \theta_I \in \mathbb{R}^p$ are the
87 real and imaginary parts of all complex parameters respectively. However, it would be erroneous to
88 assume that a complex network is equivalent to an ordinary real-valued neural network, because the
89 operation of complex multiplication limits the degree of freedom [Hirose and Yoshida, 2012].

90 2.2 Neural Tangent Kernel

91 For a real-valued deep neural network $f_{\theta_r}(\mathbf{x}) \in \mathbb{R}^{d_{out}}$ with parameter set $\theta_r \in \mathbb{R}^p$ and input $\mathbf{x} \in \mathbb{R}^d$,
92 its Neural Tangent Kernel (NTK) under gradient descent is defined as

$$\widehat{\Theta}_r(\mathbf{x}, \mathbf{x}') = \langle \nabla_{\theta_r} f_{\theta_r}(\mathbf{x}), \nabla_{\theta_r} f_{\theta_r}(\mathbf{x}') \rangle \quad (1)$$

93 which quantifies the functional gradient descent when taking an infinitely small gradient step on
94 a new observation. In case f_{θ_r} corresponds to an infinite width MLP, Jacot et al. [2018] showed
95 that $\widehat{\Theta}_r(\mathbf{x}, \mathbf{x}')$ converges to a limiting kernel $\dot{\Theta}_r(\mathbf{x}, \mathbf{x}')$ at initialization and remains frozen during
96 training, i.e.,

$$\lim_{n \rightarrow \infty} \widehat{\Theta}_r^t(\mathbf{x}, \mathbf{x}') = \lim_{n \rightarrow \infty} \widehat{\Theta}_r^0(\mathbf{x}, \mathbf{x}') = \dot{\Theta}_r(\mathbf{x}, \mathbf{x}') \quad \forall \text{ training time } t,$$

97 which could give an accurate description of training dynamics with kernel gradient descent trajectory.
98 Thus a infinitely wide neural network is governed by a linear model based on its first order Taylor
99 expansion in the parameter space [Lee et al., 2019].

100 **Tensor Programs.** After the original NTK was derived from multi-layer perceptrons, it is soon
101 extended to a variety of network structures including convolution neural networks (CNTK) [Arora
102 et al., 2019], recurrent neural networks (RNTK) [Yang, 2019, Alemohammad et al., 2020], graph
103 neural networks (GNTK) [Du et al., 2019] and so on. Importantly, Yang [2020], Yang and Littwin
104 [2021] propose NETSOR^T program, a basic form in Tensor Programs series, and prove that for a real-
105 valued neural network of any architecture that can be represented by NETSOR^T program language,
106 its NTK converges to a deterministic limit and stays frozen during training in the infinite-width limit.

107 **Neural Tangent Kernel of complex networks.** For a complex-valued neural network $f_\theta(\mathbf{z})$, also
108 denoted as $f_{[\theta_R, \theta_I]}([\mathbf{x}, \mathbf{y}])$, when it is trained by real-valued backpropagation, the empirical neural
109 tangent kernel is as follows due to that the real and imaginary parts are optimized separately

$$\widehat{\Theta}(\mathbf{z}, \mathbf{z}') = \langle \nabla_{\theta} f_{\theta}(\mathbf{z}), \nabla_{\theta} f_{\theta}(\mathbf{z}') \rangle = \langle \nabla_{\theta_R} f_{\theta_R}(\mathbf{z}), \nabla_{\theta_R} f_{\theta_R}(\mathbf{z}') \rangle + \langle \nabla_{\theta_I} f_{\theta_I}(\mathbf{z}), \nabla_{\theta_I} f_{\theta_I}(\mathbf{z}') \rangle.$$

110 We can always rewrite the NTK as $\Theta([\mathbf{x}, \mathbf{y}], [\mathbf{x}', \mathbf{y}'])$. Note that at each layer of complex networks
111 in both feed-forward and backward procedure, complex matrix multiplication structure leads to
112 numerous interactions and weight sharing between the real and imaginary parts, which makes the
113 analysis of the complex NTK challenging.

114 3 Complex Tensor Program

115 In this section, we show that complex-valued neural networks of any basic architecture also have
116 NTK behavior in the infinite-width limit: the training dynamics of real-valued backpropagation is
117 determined by kernel gradient descent with its NTK at initialization. Specifically, we extend the
118 simplified NETSOR^T program [Yang, 2020] to the complex domain, and propose the basic Complex
119 Tensor Program (CTP). Note that the complex networks are mostly non-holomorphic, thus we do not
120 require the complex tensor program to represent the backward propagation of the complex networks.

121 **Definition 1 (Complex Tensor Program)** *Given an initial set \mathcal{V} of random \mathbb{C}^n vectors and a set \mathcal{W}*
122 *of random $\mathbb{C}^{n \times n}$ complex matrices, a sequence of \mathbb{C}^n vectors is called a complex tensor program if*
123 *they are recursively generated through one of the following ways:*

124 **ComNonlin** *Given $\phi: \mathbb{C}^k \rightarrow \mathbb{C}$ and $\mathbf{z}^1, \dots, \mathbf{z}^k \in \mathbb{C}^n$, generate $\phi(\mathbf{z}^1, \dots, \mathbf{z}^k) \in \mathbb{C}$;*

125 **ComMatMul** *Given $\mathbf{W} \in \mathbb{C}^{n \times n}$ and $\mathbf{z} \in \mathbb{C}^n$, generate $\mathbf{W}\mathbf{z} \in \mathbb{C}^n$.*

126 Obviously, the complex tensor program could represent the forward procedures of all basic complex
 127 network architectures, such as the generic feed-forward full-connected complex networks [Nitta,
 128 2004], the complex-valued recurrent neural networks [Wisdom et al., 2016] and complex-valued
 129 convolutional neural networks [Trabelsi et al., 2018, Tan et al., 2020].

130 3.1 Complex network setup

131 This subsection introduces the settings and assumptions of complex networks considered. Firstly we
 132 introduce the required assumption for activation functions of complex networks, which generalizes
 133 the assumption in Yang [2020] to the complex domain.

134 **Assumption 2** We assume that the complex activation function $\phi : \mathbb{C}^k \rightarrow \mathbb{C}$ used in the complex
 135 networks and its derivative are polynomially-bounded, i.e., $\phi(\mathbf{z})$ satisfies that $|\phi(\mathbf{z})| \leq C\|\mathbf{z}\|^p + c$
 136 for some $C, p, c > 0$ and $\mathbf{z} \in \mathbb{C}^k$; so is its derivative.

137 It is worth mentioning that numerous activation functions have been proposed to deal with complex-
 138 valued representations and basically they satisfy the assumption. For example, the most commonly
 139 used complex sigmoid function $\mathbb{C}\text{Sigmoid}$ [Nitta, 1997, 2004] and complex hyperbolic tangent
 140 function $\mathbb{C}\text{tanh}$ [Nitta, 2002], which apply sigmoid and hyperbolic tangent function respectively to
 141 the real part and imaginary part separately; Moreover, the recently proposed ReLU-based complex
 142 activation functions like $\mathbb{C}\text{ReLU}$ [Trabelsi et al., 2018, Tan et al., 2020], zReLU [Guberman, 2016]
 143 and modReLU [Arjovsky et al., 2016] also satisfy the assumption.

144 **Complex NTK parametrization.** For the complex weights, we initialize each $\mathbf{W} \in \mathcal{W}$ with
 145 $\mathbf{W} = \mathbf{A} + \mathbf{B}i = \frac{\sigma_A}{\sqrt{n}}A + \frac{\sigma_B}{\sqrt{n}}Bi$ where $A_{\alpha\beta}, B_{\alpha\beta} \sim \mathcal{N}(0, 1)$, which we refer to as *complex NTK*
 146 *parametrization*. Without loss of generality, we set the variances of real and imaginary parts of
 147 all layers as σ_A and σ_B respectively in the following paper. It naturally extends the real-valued
 148 NTK parametrization [Jacot et al., 2018, Lee et al., 2019] to complex parameters. Note that this
 149 parametrization is non-vacuous because many previous works [Nitta, 1997, 2004] choose to initialize
 150 the real and imaginary parts separately.

151 **Setup.** Consider a complex-valued neural network $f_\theta(\mathbf{z})$ with complex NTK parametrization, its
 152 feed-forward procedure can be represented by a complex tensor program and complex activation
 153 functions all satisfy Assumption 2. Suppose that there is a multivariate Gaussian $\mathcal{N}_\mathcal{V}$ defined on
 154 $\mathbb{R}^{2|\mathcal{V}|}$ such that the real and imaginary variables of the initial set of vectors \mathcal{V} are sampled like
 155 $\{\Re[q]_\alpha : q \in \mathcal{V}\} \cup \{\Im[q]_\alpha : q \in \mathcal{V}\} \sim \mathcal{N}_\mathcal{V}$ i.i.d. for each coordinate $\alpha \in [n]$. For the output readout
 156 matrix W_{L+1} , we also adopt complex NTK parametrization, and it is sampled independently from
 157 all other parameters and is not used anywhere else in the interior of the network. Without loss of
 158 generality, the network is trained by SGD with batch-size 1 and learning rate 1.

159 3.2 NTK for any complex network

160 **Theorem 3 (Complex NTK at initialization)** Consider a complex-valued neural network $f_\theta(\mathbf{z})$
 161 with above setup, then as its widths go to infinity, its NTK $\hat{\Theta}(\mathbf{z}, \mathbf{z}')$ at initialization converges almost
 162 surely to a deterministic limiting kernel $\hat{\Theta}(\mathbf{z}, \mathbf{z}')$ over any finite set of inputs.

163 **Corollary 4 (Complex NTK during training)** Consider training a complex-valued neural network
 164 $f_\theta(\mathbf{z})$ with above setup. At training time t , denote the input sample as \mathbf{z}_t and the loss function as
 165 $\mathcal{L}_t : \mathbb{R} \rightarrow \mathbb{R}$. Suppose \mathcal{L}_t is continuous for all t . Then as widths approach infinity, for any $\mathbf{z} \in \mathbb{C}^d$
 166 and training time t , $f_t(\mathbf{z})$ converges almost surely to a random variable $\mathring{f}_t(\mathbf{z})$ and

$$\mathring{f}_{t+1}(\mathbf{z}) - \mathring{f}_t(\mathbf{z}) = -\hat{\Theta}(\mathbf{z}, \mathbf{z}_t) \mathcal{L}'_t(\mathring{f}_t(\mathbf{z}_t)), \quad (2)$$

167 where $\hat{\Theta}(\mathbf{z}, \mathbf{z}_t)$ is the limiting NTK of the complex network at initialization.

168 The proofs of Theorem 3 and Corollary 4 are given in Appendix A. Theorem 3 and Corollary 4 show
 169 that for a complex-valued neural network of any architecture trained by real-valued backpropagation,
 170 its NTK at initialization converges to a deterministic limiting kernel in the infinite-width limit, and
 171 the training dynamics is determined by kernel gradient descent with the NTK at initialization.

172 4 Comparison of training dynamics between complex and real networks

173 In this section, we focus on the important problem: when will complex networks have different
 174 inductive bias during training from real networks? The problem can not be solved unless the
 175 training dynamics of complex networks could be captured. Based on the NTK theory obtained in
 176 Section 3, we could provide a way to compare the training dynamics of complex and real networks
 177 by comparing their NTKs. Specifically, we have derived the NTK formula of complex multi-layer
 178 perceptrons (MLPs) and investigated the conditions under which complex MLPs trained by real-
 179 valued backpropagation will reduce to ordinary real MLPs in the infinite-width limit.

180 4.1 Neural Tangent Kernel of complex multi-layer perceptrons

181 This subsection presents the NTK formula of the most generic complex network, i.e., the L -hidden
 182 layer complex MLPs.

183 The network performs the following computation at layer $l \in [1, L]$

$$184 \quad \mathbf{h}_l = \mathbf{s}_l + \mathbf{r}_l i = \phi(\mathbf{W}_l \mathbf{h}_{l-1}) = \phi((\mathbf{A}_l + \mathbf{B}_l i)(\mathbf{s}_{l-1} + \mathbf{r}_{l-1} i)), \quad (3)$$

184 where $\mathbf{W}_l = \mathbf{A}_l + \mathbf{B}_l i$ is the complex weight matrix and $\mathbf{h}_l = \mathbf{s}_l + \mathbf{r}_l i$ with $\mathbf{h}_l \in \mathbb{C}^n$ is the output of
 185 l -th layer. For the first layer, we set $\mathbf{s}_0 = \mathbf{x}$ and $\mathbf{r}_0 = \mathbf{y}$ where $\mathbf{z} = \mathbf{x} + \mathbf{y}i$ and $\mathbf{z} \in \mathbb{C}^d$. Besides, the
 186 output of the complex network with a linear read-out layer is achieved via $f_\theta(\mathbf{x}) = \text{Re}\{\mathbf{W}_{L+1} \mathbf{h}_L\}$
 187 where $\mathbf{W}_{L+1} \in \mathbb{C}^{n \times d_{out}}$. Suppose all activation functions ϕ satisfy the Assumption 2. Complex
 188 NTK parametrization is applied for all complex parameters $\mathbf{W}_1 \in \mathbb{C}^{d \times n}$, $\mathbf{W}_{L+1} \in \mathbb{C}^{n \times d_{out}}$ and
 189 $\mathbf{W}_l \in \mathbb{C}^{n \times n}$ for $l \in [2, L]$.

190 We denote the real part of the l -th hidden layer pre-activation as $\alpha_l(\mathbf{z})$ and the imaginary part as
 191 $\beta_l(\mathbf{z})$. In the feed-forward procedure, we denote the covariance kernel functions between the real
 192 and imaginary part of the pre-activations respectively as

$$\Sigma_\alpha^l(\mathbf{z}, \mathbf{z}') = \mathbb{E}_{\theta \sim \mathcal{N}} [\alpha_l(\mathbf{z})^\top \alpha_l(\mathbf{z}') / n], \quad \Sigma_{\alpha, \beta}^l(\mathbf{z}, \mathbf{z}') = \mathbb{E}_{\theta \sim \mathcal{N}} [\alpha_l(\mathbf{z})^\top \beta_l(\mathbf{z}') / n], \quad (4)$$

$$\Sigma_\beta^l(\mathbf{z}, \mathbf{z}') = \mathbb{E}_{\theta \sim \mathcal{N}} [\beta_l(\mathbf{z})^\top \beta_l(\mathbf{z}') / n], \quad \Sigma_{\beta, \alpha}^l(\mathbf{z}, \mathbf{z}') = \mathbb{E}_{\theta \sim \mathcal{N}} [\beta_l(\mathbf{z})^\top \alpha_l(\mathbf{z}') / n]. \quad (5)$$

193 In the backward procedure, we denote the gradient vector of the real part of the l -th hidden
 194 layer pre-activation as $\delta_\alpha^l(\mathbf{z}) := \sqrt{n} (\nabla_{\alpha_l(\mathbf{z})} f_\theta(\mathbf{z}))$ and that of the imaginary part as $\delta_\beta^l(\mathbf{z}) :=$
 195 $\sqrt{n} (\nabla_{\beta_l(\mathbf{z})} f_\theta(\mathbf{z}))$. Similarly, we denote the covariance kernel functions of the gradient vector
 196 between the real and imaginary part of the pre-activations respectively as

$$\Pi_\alpha^l(\mathbf{z}, \mathbf{z}') = \mathbb{E}_{\theta \sim \mathcal{N}} [\delta_\alpha^l(\mathbf{z})^\top \delta_\alpha^l(\mathbf{z}') / n], \quad \Pi_{\alpha, \beta}^l(\mathbf{z}, \mathbf{z}') = \mathbb{E}_{\theta \sim \mathcal{N}} [\delta_\alpha^l(\mathbf{z})^\top \delta_\beta^l(\mathbf{z}') / n], \quad (6)$$

$$\Pi_\beta^l(\mathbf{z}, \mathbf{z}') = \mathbb{E}_{\theta \sim \mathcal{N}} [\delta_\beta^l(\mathbf{z})^\top \delta_\beta^l(\mathbf{z}') / n], \quad \Pi_{\beta, \alpha}^l(\mathbf{z}, \mathbf{z}') = \mathbb{E}_{\theta \sim \mathcal{N}} [\delta_\beta^l(\mathbf{z})^\top \delta_\alpha^l(\mathbf{z}') / n]. \quad (7)$$

197 **Theorem 5** For a L -hidden layer complex MLP, with all activation functions satisfying Assumption 2,
 198 in the limit as all widths $n \rightarrow \infty$, the empirical NTK at initialization converges to the following
 199 limiting kernel

$$\lim_{n \rightarrow \infty} \hat{\Theta}(\mathbf{z}, \mathbf{z}') = \mathring{\Theta}(\mathbf{z}, \mathbf{z}') = \Theta(\mathbf{z}, \mathbf{z}') \otimes \mathbf{I}_{d_{out}} \quad (8)$$

200 where

$$\Theta(\mathbf{z}, \mathbf{z}') = \sum_{l=1}^L (\Pi_\alpha^l(\mathbf{z}, \mathbf{z}') \Sigma_\alpha^l(\mathbf{z}, \mathbf{z}') + \Pi_\beta^l(\mathbf{z}, \mathbf{z}') \Sigma_\beta^l(\mathbf{z}, \mathbf{z}')) \quad (9)$$

$$+ \Pi_{\alpha, \beta}^l(\mathbf{z}, \mathbf{z}') \Sigma_{\alpha, \beta}^l(\mathbf{z}, \mathbf{z}') + \Pi_{\beta, \alpha}^l(\mathbf{z}, \mathbf{z}') \Sigma_{\beta, \alpha}^l(\mathbf{z}, \mathbf{z}')) + \Sigma_\alpha^{L+1}(\mathbf{z}, \mathbf{z}') \quad (10)$$

201 where the covariance functions $\Sigma_\alpha^l, \Sigma_\beta^l, \Pi_\alpha^l, \Pi_\beta^l$ are defined in Eq. 4-7.

202 The result is proved in the Appendix C, where the detailed recursions of intermediate kernels for the
 203 NTK calculation are also presented.

204 Note that the NTK formula of a complex MLP looks very different from the NTK of a real MLP
 205 given by Jacot et al. [2018] due to the existence of interaction between real and imaginary parts,
 206 which is caused by joint weight sharing in complex matrix multiplication. However, if we go deeper,
 207 does there exist situations that the NTK of a complex MLP will reduce to that of a real MLP?

208 **4.2 Asymptotic equivalence of training dynamics**

209 In this subsection, we provide our main results. Surprisingly, we show that, for commonly used
 210 complex activation functions, complex networks trained by real-valued backpropagation have the
 211 same inductive bias as real networks during training in the infinite-width limit.

212 We first define asymptotic equivalence between neural networks, which represents a perspective to
 213 investigate when will a complex network have different inductive bias from a real network during
 214 training. It also helps if we change the network structure or backpropagation algorithms.

215 **Definition 6** *Two neural networks trained by gradient descent are asymptotic equivalent, if as all
 216 widths go to infinity, their neural tangent kernels Θ converge to the same deterministic limit $\hat{\Theta}$ at
 217 initialization and have the same optimization trajectory during training.*

218 The following theorem is our main result: if we train complex networks with real-valued back
 219 propagation, under very common conditions, the complex MLPs are asymptotic equivalent with real
 220 MLPs, thus they have the same training dynamics.

221 **Theorem 7** *Consider a complex MLP in Eq. (3) and a ordinary real MLP with L hidden layers
 222 trained by real-valued backpropagation. Suppose $\sigma_A = \sigma_B$ at initialization and the activation
 223 functions satisfy Assumption 2. As the widths go to infinity, they are asymptotic equivalent if the
 224 activation functions satisfy one of the following conditions*

225 **Condition 1** *All split activation functions ϕ satisfying $\phi(\alpha, \beta) = \phi_R(\alpha) + \phi_R(\beta)i$*

226 **Condition 2** *A subset of holomorphic activation functions ϕ satisfying $\phi_2(\alpha, \beta) = \phi_1(\beta, -\alpha)$
 227 and $\frac{\partial \phi_1(\alpha, \beta)}{\partial \beta} = \frac{\partial \phi_2(\alpha, \beta)}{\partial \alpha} = 0$*

228 where the general complex activation function is denoted as $\phi(\alpha, \beta) = \phi_1(\alpha, \beta) + \phi_2(\alpha, \beta)i$ with
 229 input pre-activations $\alpha + \beta i$.

230 In the Appendix D the result is proved and we also give the sufficient and necessary conditions.
 231 For simplicity, here we only show the most informative conditions. The key idea of the proof is
 232 transforming asymptotic equivalence into four complex conditions and find the common solutions
 233 based on Rules.F.13 in Appendix F.

234 **Discussion about the Condition 1.** Note that most commonly used complex activation functions
 235 satisfy the Condition 1:

$$\phi(z) = \phi_R(\Re(z)) + \phi_R(\Im(z))i$$

236 like complex sigmoid function [Benvenuto and Piazza, 1992, Nitta, 1997], complex hyperbolic
 237 tangent function [Hirose and Yoshida, 2012], etc. It is also worth mentioning that the recently
 238 proposed ReLU-based complex activation function $\mathbb{C}ReLU$ also satisfies the Condition 1, which has
 239 achieved the best performance in feed-forward complex networks in image processing tasks [Trabelsi
 240 et al., 2018, Tan et al., 2020] among all ReLU-based complex activation functions.

241 **Remark 8** *Because of Liouville’s theorem, the only complex-valued functions that are bounded and
 242 analytic everywhere are constants. Thus in practice, one must choose between boundedness and
 243 analyticity for a complex activation function. Before the popularity of ReLU, almost all activation
 244 functions in the real case were bounded. Consequently, previous works about complex networks
 245 always preferred non-analytic functions to preserve boundedness. Most commonly they applied split
 246 activation functions separately to the real and imaginary parts, as investigated in Bassey et al. [2021]
 247 and Scardapane et al. [2018]. So the condition contains most complex networks in practice.*

248 As a result, the theorem demonstrates that for complex networks with all these common complex
 249 activations, if they are trained by real-valued BP, then these complex networks reduce to real networks
 250 as widths grow, despite the joint interaction weight sharing caused by complex matrix multiplication
 251 structure. Consequently, real-valued backpropagation eliminates the characteristics of most complex
 252 networks. This may guide the selection of training algorithm in practice if people want to take full
 253 advantage of complex networks, and encourage people to use backpropagation algorithms specially
 254 designed for complex networks, like Wirtinger calculus [Adali et al., 2011].

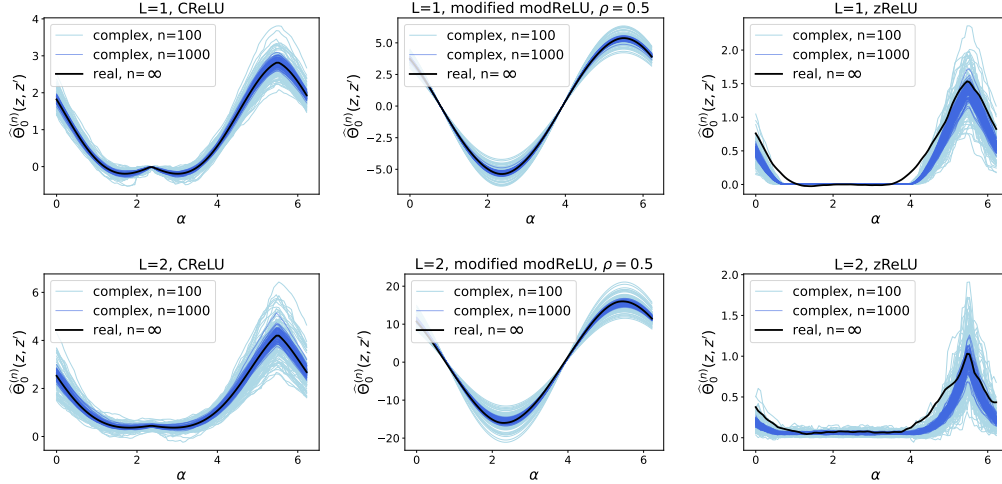


Figure 1: **Convergence of complex NTKs to corresponding real NTKs at initialization.** One input $z = 1 - i$ is fixed and the other input $z' = \cos \alpha + \sin \alpha i$ varies with $\alpha \in [0, 2\pi]$. The black line is the limiting NTKs of real MLPs $\hat{\Theta}_r(z, z')$ while the light blue and blue ones are empirical NTKs of complex MLPs $\hat{\Theta}_0^{(n)}(z, z')$ with width $n = 100$ and 1000 . For each width, $\hat{\Theta}_0^{(n)}(z, z')$ is calculated 100 times randomly, corresponding to 100 lines in the figure.

255 **Discussion about the Condition 2.** Condition 2 is also non-vacuous since it is a subset of Cauchy-
 256 Riemann condition. It includes all the magnitude-based ReLU-type complex activation functions.
 257 For example, we can easily obtain a modified modReLU [Arjovsky et al., 2016] and a modified
 258 phase-based ReLU satisfying Condition 2 as follows

$$\phi(z) = \begin{cases} z & \text{if } |z| \geq \rho, \\ 0 & \text{otherwise.} \end{cases} \quad \phi(z) = \begin{cases} z & \text{if } g(\cos \theta_z) \leq \gamma, \\ 0 & \text{otherwise.} \end{cases}$$

259 where $g(\cos \theta_z)$ can be any function of phase $\cos \theta_z$. Note that, these activation functions are
 260 analytic almost everywhere, and according to previous theory, they enjoy better theoretical results
 261 like separation results [Zhang et al., 2021] and local minima [Wu et al., 2021]. However, our results
 262 indicate that real-valued backpropagation eliminates these advantages of complex networks in the
 263 infinite-width limit, which further illustrates the inappropriateness of real-valued backpropagation.

264 5 Empirical study

265 In this section, we empirically verify the relationship between NTKs of the complex networks and
 266 real networks, and investigate the network widths required for the establishment of our results.
 267 We consider complex-valued MLPs with one or two hidden layers and we use CReLU, modified
 268 modReLU, CSigmoid, Ctanh and zReLU as the activation functions. Note that all these activation
 269 functions satisfy the conditions of our theorem except zReLU. Through the following experiments,
 270 we want to check whether the empirical complex NTKs $\hat{\Theta}_t^{(n)}$ converge to the corresponding real
 271 NTKs $\hat{\Theta}_r$ with different activation functions as the widths n grow.

272 For real networks, the input is the concatenated vector of the real and imaginary parts of the complex-
 273 valued input. Corresponded to the complex-valued fully-connected layer, we use the commonly used
 274 real-valued fully-connected layer without complex matrix multiplication structure. To implement
 275 the corresponding activation functions, complex-valued activation functions are transformed to real-
 276 valued ones in the following way: we divide the pre-activation vector into two half, treat the first
 277 half as real parts and the second as imaginary parts, as the input of $\phi(\alpha, \beta)$. Then we concatenate
 278 the real and imaginary parts after activation. For split activation functions like CReLU, it can just
 279 correspond to real ReLU activation. In NTK initialization, the standard deviations are set as 1 for
 280 complex networks and scaled to $\sqrt{2}$ for real networks. All empirical NTKs of complex networks are
 281 calculated based on the Neural Tangents library [Novak et al., 2019].

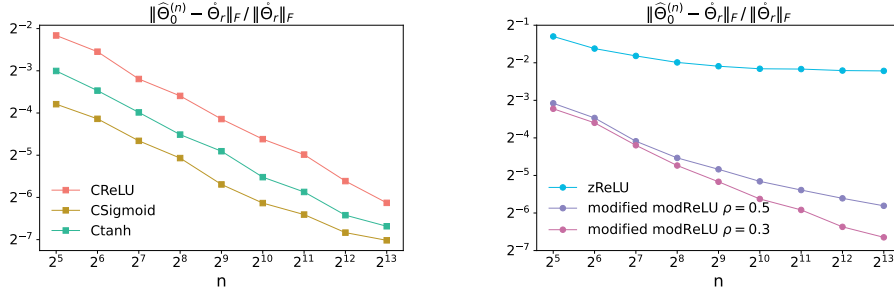


Figure 2: **Convergence of complex NTKs to corresponding real NTKs under various activation functions with a larger range of widths.** **Left:** three kinds of split activation functions satisfying our conditions. **Right:** other complex-valued activation functions where modified modReLU satisfies conditions and zReLU does not.

282 **Verifying asymptotic equivalence at initialization.** The first experiment shows the distribution
 283 of empirical NTKs of complex MLPs $\hat{\Theta}_t^{(n)}(z, z')$ and analytic NTKs of real MLPs $\check{\Theta}_r(z, z')$ with
 284 different z' at initialization on a synthetic dataset. We define $z = 1 - i$ and $z' = \cos \alpha + \sin \alpha i$
 285 for $\alpha \in [0, 2\pi]$. Then we can view $\hat{\Theta}_t^{(n)}(z, z')$ and $\check{\Theta}_r(z, z')$ as functions of α . For the complex
 286 networks, we calculate empirical NTKs for hidden layer width $n = 100, 1000$ at initialization and
 287 hidden layer number $l = 1, 2$. In each case, we calculate $\hat{\Theta}_t^{(n)}(z, z')$ 100 times with different random
 288 NTK initialization. We compare these empirical complex NTKs with corresponding real NTKs. In
 289 the case of CReLU, we calculate $\check{\Theta}_r$ by the analytic form solution of NTK [Cho and Saul, 2009];
 290 in the case of zReLU and modified modReLU, it's hard to get the closed form solution, so we use
 291 the average of a large amount of wide empirical NTKs to approximate $\check{\Theta}_r$. The results are shown
 292 in Figure 1. In the figure we observe that for CReLU and modified modReLU, which satisfy our
 293 conditions, their NTKs $\hat{\Theta}_0^{(n)}$ concentrate to the NTKs of real MLPs $\check{\Theta}_r$ perfectly, and when $n = 1000$,
 294 the convergence is more concentrated and the complex NTKs almost equal to real NTKs; for zReLU
 295 which does not satisfy the conditions, there's a gap between complex and real NTKs. Therefore, the
 296 results verify our results perfectly and demonstrate our results are non-vacuous.

297 **Verifying asymptotic equivalence with more activation functions as widths grow much larger.**
 298 In the first experiment, although for zReLU there is a gap between complex and real networks, the
 299 tendency is still similar. For the second experiment, we do the similar experiment on the same
 300 synthetic dataset, so that we can see what will happen when n becomes much larger. Besides,
 301 we consider more different activation functions which satisfy our conditions including CReLU,
 302 CSigmoid, Ctanh and modified modReLU with different hyper-parameters ρ . We calculate relative
 303 Frobenius norm $\|\hat{\Theta}_0^{(n)}(X, X) - \check{\Theta}_r(X, X)\|_F / \|\check{\Theta}_r(X, X)\|_F$ on set X with widths n ranging from
 304 2^5 to 2^{13} , which measures the difference between complex NTKs $\hat{\Theta}_0^{(n)}$ and real NTKs $\check{\Theta}_r$. Figure 2
 305 shows the result. For all those split activation functions (CReLU, CSigmoid and Ctanh), the tendency
 306 of convergence remains unchanged even when the width n goes to quite a large number 2^{13} . The two
 307 curves of modified modReLU act similarly with that of split activation functions; However, the curve
 308 of zReLU does not converge at all at initialization.

309 **Verifying asymptotic equivalence during training.** The third experiment investigates the conver-
 310 gence of difference between complex NTKs $\hat{\Theta}_t^{(n)}$ and real NTKs $\check{\Theta}_r$ during training as the widths go
 311 to infinity on MNIST [LeCun et al., 1998]. We randomly choose a subset of MNIST as training set
 312 $\mathcal{D} = (X, Y)$ ($|\mathcal{D}| = 128$), and treat the first half of features as real parts, the second half as imaginary
 313 parts. Then we calculate relative Frobenius norm between empirical NTKs of complex networks at
 314 time t and real NTKs at initialization with widths n ranging from 2^5 to 2^{10} at initialization ($t = 0$)
 315 and during training ($t = 1000$). Due to the memory limitations, we cannot try larger n . The learning
 316 rate η is 0.5 for $l = 1$ and 0.2 for $l = 2$. The results are shown in Figure 3. We can see that in all
 317 these cases, the relative Frobenius norm decreases as n goes up, regardless of the training steps and
 318 hidden layer numbers. Therefore, our theoretical results hold during training. More experiments
 319 verifying asymptotic equivalence can be found in Appendix E.

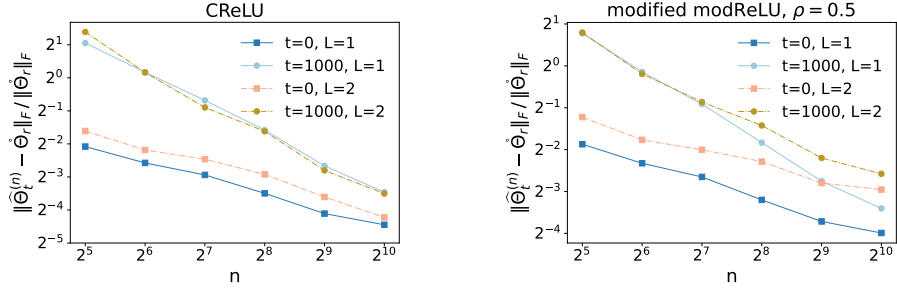


Figure 3: **Convergence of complex NTKs to corresponding real NTKs as widths grow during training.** The Y-axis is the difference between empirical complex NTKs and corresponding real NTK in terms of the relative Frobenius norm. At each point the relative Frobenius norm is calculated 20 times and the mean value is shown. The solid line indicates one hidden layer ($L = 1$) while the dashdot line indicates two hidden layers ($L = 2$). The square marker indicates $t = 0$ (at initialization) while the circle marker indicates $t = 1000$ (during training).

320 **Empirical results.** Overall, in the case of complex activation functions satisfying our theorems,
 321 such as \mathbb{C} ReLU and modified modReLU, complex NTKs $\hat{\Theta}_t^{(n)}$ converges to real NTKs $\hat{\Theta}_r$, quite
 322 well even the widths are about 1000; in the case of complex activation function that does not satisfy our
 323 theorems like zReLU, the convergence from the complex NTK $\hat{\Theta}_0^{(n)}$ to the real NTK $\hat{\Theta}_r$ does not
 324 occur. This validates our theory perfectly and demonstrates that our conditions are non-vacuous.

325 6 Related work

326 Long before the popularity of deep learning, there have been many investigations on complex-
 327 valued neural networks [Hirose, 1992, Benvenuto and Piazza, 1992, Nitta, 1997]. However, It is
 328 always challenging to train complex networks due to analytical properties. The most notable reason
 329 is that almost all cost functions are real-valued and thus non-holomorphic. In order to perform
 330 backpropagation for complex networks, the conventional approach to overcome the limitation is to
 331 use separate derivatives with respect to the real-imaginary parts of a non-analytic function [Nitta,
 332 2004], or split amplitude-phase parts [Hirose, 1992]. Hirose and Yoshida [2012] has shown that split
 333 backpropagation for amplitude-phase parts could achieve better generalization than real networks
 334 on signal processing tasks. Backpropagation based on Wirtinger calculus [Wirtinger, 1927] has
 335 recently received increasing attention [Adali et al., 2011, Bassey et al., 2021], which presents an
 336 elegant alternative and allows keeping all computations in the complex domain. However, real-valued
 337 backpropagation [Nitta, 1997] is widely used recently to train deep complex networks due to the
 338 convenience of utilizing the real-valued deep learning library [Arjovsky et al., 2016, Trabelsi et al.,
 339 2018, Tan et al., 2020] and achieved state-of-the-art performance. Thus we start investigating from
 340 it. Theoretically, there have been important advances for complex networks regarding the universal
 341 approximation property [Voigtlaender, 2020], local minima [Nitta, 2013, Wu et al., 2021], separation
 342 results [Zhang et al., 2021]. However, to our best knowledge, there is still no theoretical analysis of
 343 the training dynamics of complex networks and the equivalence between complex and real networks.

344 7 Conclusion

345 In this paper, we first extend the NTK theory to the complex domain, and propose a way to compare
 346 the training dynamics between complex and real networks based on their NTKs. Surprisingly, we
 347 find that the commonly used real-valued backpropagation reduces the training dynamics of complex-
 348 valued MLPs to that of ordinary real MLPs, thus eliminating the characteristics of complex-valued
 349 neural networks. Empirical study verifies that our results are practical for commonly used complex
 350 activation functions. We hope the proposed method and results could help the design and analysis of
 351 training for complex networks. For future work, developing the analysis for more training algorithms
 352 and complex network structures will be exciting directions. It is also significant to consider other
 353 parametrizations in the future investigation.

354 References

- 355 Tülay Adalı, Peter J Schreier, and Louis L Scharf. Complex-valued signal processing: The proper
356 way to deal with impropriety. *IEEE Transactions on Signal Processing*, 59(11):5101–5125, 2011.
- 357 Sina Alemohammad, Zichao Wang, Randall Balestriero, and Richard Baraniuk. The recurrent neural
358 tangent kernel. In *International Conference on Learning Representations*, 2020.
- 359 Martin Arjovsky, Amar Shah, and Yoshua Bengio. Unitary evolution recurrent neural networks. In
360 *International Conference on Machine Learning*, pages 1120–1128, 2016.
- 361 Sanjeev Arora, Simon S Du, Wei Hu, Zhiyuan Li, Ruslan Salakhutdinov, and Ruosong Wang. On
362 exact computation with an infinitely wide neural net. In *Advances in Neural Information Processing
363 Systems*, pages 8139–8148, 2019.
- 364 Joshua Bassegy, Lijun Qian, and Xianfang Li. A survey of complex-valued neural networks. *arXiv
365 preprint arXiv:2101.12249*, 2021.
- 366 Nevio Benvenuto and Francesco Piazza. On the complex backpropagation algorithm. *IEEE Transac-
367 tions on Signal Processing*, 40(4):967–969, 1992.
- 368 Youngmin Cho and Lawrence Saul. Kernel methods for deep learning. In *Advances in Neural
369 Information Processing Systems*, pages 342–350, 2009.
- 370 Ivo Danihelka, Greg Wayne, Benigno Uribe, Nal Kalchbrenner, and Alex Graves. Associative long
371 short-term memory. In *International Conference on Machine Learning*, pages 1986–1994, 2016.
- 372 Simon S Du, Kangcheng Hou, Russ R Salakhutdinov, Barnabas Poczos, Ruosong Wang, and Keyulu
373 Xu. Graph neural tangent kernel: Fusing graph neural networks with graph kernels. In *Advances
374 in Neural Information Processing Systems*, pages 5724–5734, 2019.
- 375 Nitzan Guberman. On complex valued convolutional neural networks. *arXiv preprint
376 arXiv:1602.09046*, 2016.
- 377 Akira Hirose. Continuous complex-valued back-propagation learning. *Electronics Letters*, 28(20):
378 1854–1855, 1992.
- 379 Akira Hirose and Shotaro Yoshida. Generalization characteristics of complex-valued feedforward
380 neural networks in relation to signal coherence. *IEEE Transactions on Neural Networks and
381 Learning Systems*, 23(4):541–551, 2012.
- 382 Arthur Jacot, Clément Hongler, and Franck Gabriel. Neural tangent kernel: Convergence and
383 generalization in neural networks. In *Advances in Neural Information Processing Systems*, pages
384 8580–8589, 2018.
- 385 Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to
386 document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- 387 Jaehoon Lee, Lechao Xiao, Samuel Schoenholz, Yasaman Bahri, Roman Novak, Jascha Sohl-
388 Dickstein, and Jeffrey Pennington. Wide neural networks of any depth evolve as linear models
389 under gradient descent. In *Advances in Neural Information Processing Systems*, pages 8572–8583,
390 2019.
- 391 Tohru Nitta. An extension of the back-propagation algorithm to complex numbers. *Neural Networks*,
392 10(8):1391–1415, 1997.
- 393 Tohru Nitta. On the critical points of the complex-valued neural network. In *Proceedings of the 9th
394 International Conference on Neural Information Processing, 2002. ICONIP'02.*, volume 3, pages
395 1099–1103, 2002.
- 396 Tohru Nitta. Orthogonality of decision boundaries in complex-valued neural networks. *Neural
397 Computation*, 16(1):73–97, 2004.
- 398 Tohru Nitta. Local minima in hierarchical structures of complex-valued neural networks. *Neural
399 Networks*, 43:1–7, 2013.

- 400 Roman Novak, Lechao Xiao, Jiri Hron, Jaehoon Lee, Alexander A Alemi, Jascha Sohl-Dickstein,
401 and Samuel S Schoenholz. Neural tangents: Fast and easy infinite neural networks in python. In
402 *International Conference on Learning Representations*, 2019.
- 403 Simone Scardapane, Steven Van Vaerenbergh, Amir Hussain, and Aurelio Uncini. Complex-valued
404 neural networks with nonparametric activation functions. *IEEE Transactions on Emerging Topics*
405 *in Computational Intelligence*, 4(2):140–150, 2018.
- 406 Xiaofeng Tan, Ming Li, Peng Zhang, Yan Wu, and Wanying Song. Complex-valued 3-d convolutional
407 neural network for polsar image classification. *IEEE Geoscience and Remote Sensing Letters*, 17
408 (6):1022–1026, 2020.
- 409 Chiheb Trabelsi, Olexa Bilaniuk, Ying Zhang, Dmitriy Serdyuk, Sandeep Subramanian, Joao Felipe
410 Santos, Soroush Mehri, Negar Rostamzadeh, Yoshua Bengio, and Christopher J Pal. Deep complex
411 networks. In *International Conference on Learning Representations*, 2018.
- 412 Mark Tygert, Joan Bruna, Soumith Chintala, Yann LeCun, Serkan Piantino, and Arthur Szlam. A
413 mathematical motivation for complex-valued convolutional networks. *Neural Computation*, 28(5):
414 815–825, 2016.
- 415 Felix Voigtlaender. The universal approximation theorem for complex-valued neural networks. *arXiv*
416 *preprint arXiv:2012.03351*, 2020.
- 417 Wilhelm Wirtinger. Zur formalen theorie der funktionen von mehr komplexen veränderlichen.
418 *Mathematische Annalen*, 97(1):357–375, 1927.
- 419 Scott Wisdom, Thomas Powers, John Hershey, Jonathan Le Roux, and Les Atlas. Full-capacity
420 unitary recurrent neural networks. In *Advances in Neural Information Processing Systems*, pages
421 4880–4888, 2016.
- 422 Jin-Hui Wu, Shao-Qun Zhang, Yuan Jiang, and Zhi-Hua Zhou. Towards theoretical under-
423 standing of flexible transmitter networks via approximation and local minima. *arXiv preprint*
424 *arXiv:2111.06027*, 2021.
- 425 Greg Yang. Scaling limits of wide neural networks with weight sharing: Gaussian process behavior,
426 gradient independence, and neural tangent kernel derivation. *arXiv preprint arXiv:1902.04760*,
427 2019.
- 428 Greg Yang. Tensor programs ii: Neural tangent kernel for any architecture. *arXiv preprint*
429 *arXiv:2006.14548*, 2020.
- 430 Greg Yang and Etai Littwin. Tensor programs iib: Architectural universality of neural tangent kernel
431 training dynamics. In *International Conference on Machine Learning*, pages 11762–11772, 2021.
- 432 Xin Yao, Xiaoran Shi, and Feng Zhou. Human activities classification based on complex-value
433 convolutional neural network. *IEEE Sensors Journal*, 20(13):7169–7180, 2020.
- 434 Shao-Qun Zhang, Wei Gao, and Zhi-Hua Zhou. Towards understanding theoretical advantages of
435 complex-reaction networks. *arXiv preprint arXiv:2108.06711*, 2021.

436 **Checklist**

- 437 1. For all authors...
- 438 (a) Do the main claims made in the abstract and introduction accurately reflect the paper's
439 contributions and scope? [Yes]
- 440 (b) Did you describe the limitations of your work? [Yes] **See Section 3 for assumptions.**
- 441 (c) Did you discuss any potential negative societal impacts of your work? [N/A]
- 442 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
443 them? [Yes]
- 444 2. If you are including theoretical results...
- 445 (a) Did you state the full set of assumptions of all theoretical results? [Yes] **See Section 3**
446 **for assumptions.**
- 447 (b) Did you include complete proofs of all theoretical results? [Yes] **See Appendix.**
- 448 3. If you ran experiments...
- 449 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
450 mental results (either in the supplemental material or as a URL)? [Yes]
- 451 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
452 were chosen)? [Yes]
- 453 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
454 ments multiple times)? [N/A] **The numerical experiments only aim to verify the**
455 **theoretical results.**
- 456 (d) Did you include the total amount of compute and the type of resources used (e.g., type
457 of GPUs, internal cluster, or cloud provider)? [N/A] **The numerical experiments only**
458 **aim to verify the theoretical results.**
- 459 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 460 (a) If your work uses existing assets, did you cite the creators? [Yes] **MNIST.**
- 461 (b) Did you mention the license of the assets? [N/A] **MNIST. GNU General Public**
462 **License v3.0**
- 463 (c) Did you include any new assets either in the supplemental material or as a URL? [No]
- 464 (d) Did you discuss whether and how consent was obtained from people whose data you're
465 using/curating? [N/A]
- 466 (e) Did you discuss whether the data you are using/curating contains personally identifiable
467 information or offensive content? [N/A]
- 468 5. If you used crowdsourcing or conducted research with human subjects...
- 469 (a) Did you include the full text of instructions given to participants and screenshots, if
470 applicable? [N/A]
- 471 (b) Did you describe any potential participant risks, with links to Institutional Review
472 Board (IRB) approvals, if applicable? [N/A]
- 473 (c) Did you include the estimated hourly wage paid to participants and the total amount
474 spent on participant compensation? [N/A]