

Don't Fear Peculiar Activation Functions: EUAF and Beyond

Qianchao Wang^{a,**}, Shijun Zhang^{b,**}, Dong Zeng^c, Zhaoheng Xie^d, Hengtao Guo^e, Tiejong Zeng^a, Feng-Lei Fan^{a,*}

^aCenter of Mathematical Artificial Intelligence, Department of Mathematics, The Chinese University of Hong Kong, Hong Kong, China

^bDepartment of Applied Mathematics, The Hong Kong Polytechnic University, Hong Kong, China

^cDepartment of Biomedical Engineering, Southern Medical University, Guangzhou, China

^dInstitute of Medical Technology, Peking University Health Science Center, Peking University, Beijing, China

^eIndependent Researcher, 708 6th Ave N, Seattle, WA 98109, US,

Abstract

In this paper, we propose a new super-expressive activation function called the Parametric Elementary Universal Activation Function (PEUAF). We demonstrate the effectiveness of PEUAF through systematic and comprehensive experiments on various industrial and image datasets, including CIFAR10, Tiny-ImageNet, and ImageNet. Moreover, we significantly generalize the family of super-expressive activation functions, whose existence has been demonstrated in several recent works by showing that any continuous function can be approximated to any desired accuracy by a fixed-size network with a specific super-expressive activation function. Specifically, our work addresses two major bottlenecks in impeding the development of super-expressive activation functions: the limited identification of super-expressive functions, which raises doubts about their broad applicability, and their often peculiar forms, which lead to skepticism regarding their scalability and practicality in real-world applications.

Keywords: Deep Neural Networks, Approximation Theory, Super-Expressiveness, Parametric Elementary Universal Activation Function (PEUAF), Industrial Applications

1. INTRODUCTION

In recent years, deep learning has achieved significant success in many critical areas (LeCun et al., 2015). A major factor contributing to this success is the development of highly effective nonlinear activation functions, which greatly enhance the information processing capabilities of neural networks. While established options like the Rectified Linear Unit (ReLU) and its variants are widely used (Nair and Hinton, 2010), the fundamental importance of activation functions makes the search for better ones a continuous effort. Researchers are persistently working to design and evaluate various activation functions through both theoretical analysis and empirical studies (Bingham and Miikkulainen, 2022; Apicella et al., 2021; Wang et al., 2024).

In the realm of approximation theory, it has been shown that certain activation functions can empower a neural network with a simple structure to approximate any continuous function with an arbitrarily small error, using a fixed number of neurons (Maiorov and Pinkus,

1999). These functions are termed “super-expressive activation functions” (Yarotsky, 2021). According to research, to achieve super-expressiveness, an activation function should possess both periodic and analytical components (Shen et al., 2022; Yarotsky, 2021). One such example is the elementary universal activation function (EUAF), defined as follows:

$$\text{EUAF}(x) := \begin{cases} |x - 2\lfloor \frac{x+1}{2} \rfloor| & \text{for } x \geq 0, \\ \frac{x}{1+|x|} & \text{for } x < 0, \end{cases}$$

Figure 1 depicts EUAF, an analytical function on $(-\infty, 0)$ and periodic on $[0, \infty)$. The unique and highly desirable property of super-expressiveness allows neural networks to achieve precise approximation accuracy without increasing network complexity. This contrasts with traditional universal approximation methods, where more complex structures and a higher number of neurons are required as the approximation error decreases. By integrating super-expressive activation functions, one can attain the desired approximation accuracy by merely adjusting parameters, thus maintaining a simpler network architecture.

To the best of our knowledge, the development of

*Corresponding author, hitfanfenglei@gmail.com

**Qianchao Wang and Shijun Zhang are co-first authors.

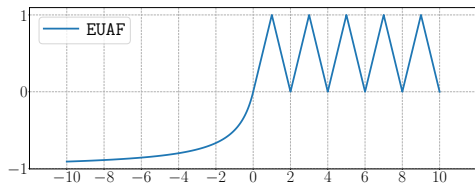


Figure 1: An illustration of EUAF.

super-expressive activation functions faces two technical challenges that hinder their potential value to neural networks: 1) First, only a limited number of super-expressive functions have been identified so far (Maierov and Pinkus, 1999; Shen et al., 2022; Yarotsky, 2021). It is unclear if the super-expressive property can be broadly applied. Additionally, for deep learning practitioners, having a greater variety of activation functions that exhibit learning capabilities is necessary in terms of enriching their armory. Developing more super-expressive functions increases the likelihood of finding their utilities in important applications, as different activation functions differ in their trainability. 2) Second, the practical utility of super-expressive activation functions is questionable. While superior expressiveness can be theoretically established through specialized constructions that demonstrate the existence of an expressive solution (Shen et al., 2021; Yarotsky, 2021), this does not necessarily translate to better practical performance. Furthermore, it is unclear whether gradient-based methods can effectively learn good solutions for networks using these functions.

Compared to commonly used functions like ReLU, sigmoid, and tanh, super-expressive functions usually have peculiar shapes. For example, Figure 1 shows EUAF, which is a typical super-expressive activation function. It has a complex and intimidating form, which makes most practitioners skeptical about its scalability and practicality in real-world applications. If we can demonstrate the practical utility of any super-expressive activation function, it could help resolve the skepticism and bridge the gap between their theoretical elegance and usefulness.

In addressing the first bottleneck, we substantially generalize the scope of EUAF to encompass a large family of functions capable of achieving super-expressiveness. Specifically, an activation function ρ is considered to be super-expressive if it is real analytic within a small interval and a fixed-size ρ -activated network can reproduce a triangle-wave function. To address the second bottleneck, we believe that super-expressive functions can indeed be practically useful. Previous studies (Sitzmann et al., 2020; Ramirez et al.,

2023) successfully applied the periodic function \sin as an activation function within the implicit neural representation. These models have been demonstrated to be suitable for representing complex signals and their derivatives, as well as for solving challenging boundary value problems (Liu et al., 2022a). These studies provide valuable insights into the potential of super-expressive activation functions, since both super-expressive activation functions and \sin share periodicity. Moreover, from the perspective of signal decomposition, normal activation functions like ReLU tend to assist models in identifying the direct component (DC) of a signal (Lee et al., 2024). In contrast, super-expressive activation functions can better handle stationary signals due to their inherent periodicity. This characteristic enhances their ability to manage complex real-world signals more efficiently.

Specifically, we choose EUAF as our representative and investigate a parameterized variant, named PEUAF, which adaptively learns the frequency w on the positive side. Mathematically,

$$\text{PEUAF}(x) := \begin{cases} |wx - 2\lfloor \frac{wx+1}{2} \rfloor| & \text{for } x \geq 0, \\ \frac{x}{1+|x|} & \text{for } x < 0, \end{cases}$$

where w is the trainable parameter representing the frequency on the positive side. PEUAF can adaptively extract the stationary signals with different frequencies. This adaptability allows PEUAF to effectively capture and represent signals with diverse frequency components, which is particularly advantageous in addressing real-world signal complexities. Then, we validate the effectiveness of PEUAF by experimenting with four industrial datasets (1D data) and three image datasets (2D data). For industrial datasets, our tests show that PEUAF surpasses other activation functions in terms of test accuracy, convergence speed, and fault localization ability. For image datasets, we find that combining PEUAF with other activation functions can usually yield better performance than only using a single activation function, although using PEUAF alone cannot achieve satisfactory performance. Thus, PEUAF can serve as a valuable add-on to the network. Our main contributions are as follows:

- We provide a non-trivial generalization of EUAF, showing that a broader family of activation functions can achieve super-expressiveness.
- We bridge the gap between the theoretical elegance and empirical usefulness of super-expressive functions by demonstrating their competitive performance in practical applications through systematic

131 experiments on four industrial datasets and three
132 image datasets including ImageNet.

- 133 • We introduce PEUAF, a parameterized version of
134 EUAF, and demonstrate that PEUAF can be used
135 individually or in conjunction with other well-
136 performing activation functions.

137 2. RELATED WORK

138 In the field of artificial intelligence, deep neural net-
139 works have proven to be highly effective tools. These
140 networks leverage the power of interconnected nodes
141 structured in multiple layers, allowing them to excel in
142 a wide range of complex applications and new domains.
143 At their core, deep neural networks rely on an affine
144 linear transformation followed by a nonlinear activation
145 function. The nonlinear activation function is essential
146 for the successful training of these networks.

147 Later in this section, we will first review conventional
148 activation functions including ReLU and its variants,
149 as well as recent sigmoidal activation functions in Sec-
150 tion 2.1. We will then discuss super-expressive activa-
151 tion functions in Section 2.2.

152 2.1. Conventional Activation Functions

153 In recent years, the Rectified Linear Unit
154 (ReLU (Nair and Hinton, 2010)), defined as
155 $\text{ReLU}(x) = \max(0, x)$, has gained popularity and
156 recognition for its effectiveness in addressing the gra-
157 dient vanishing and explosion issues encountered with
158 Sigmoid and Tanh activation functions. Thus, ReLU
159 has been widely used in the deep learning community
160 such as industrial fault diagnosis (Liu et al., 2024a) and
161 medical image segmentation (Liu et al., 2024b). How-
162 ever, ReLU can suffer from the occurrence of a number
163 of “dead neurons”, which results in information loss
164 and can hurt the neural network’s feature processing
165 ability. To mitigate this issue, several variants of ReLU
166 have been introduced such as Leaky Rectified Linear
167 Unit (LReLU) (Xu et al., 2015), Parametric Rectified
168 Linear Unit (PReLU) (He et al., 2015), Randomized
169 Leaky Rectified Linear Unit (RRReLU) (Xu et al., 2015),
170 Exponential Linear Unit (ELU) (Clevert et al., 2015),
171 Gaussian Error Linear Unit (GELU) (Hendrycks and
172 Gimpel, 2016), and Generalized Linear Unit (GENLU)
173 (Fan et al., 2020). Most recently, Goldenstein et al.
174 (2024) proposed Self-Normalizing ReLU or NeLU
175 to ensure that the prediction model is not affected
176 by the noise level during testing. It has been tested
177 in synthetic data and image de-noising tasks. These

178 variants represents a significant advancement in ac-
179 tivation function design, offering adaptability and
180 potentially better performance. Whereas, their benefits
181 come with the cost of increased model complexity or
182 computation burden and the need for careful tuning
183 and regularization which inspired researchers to create
184 more different activation functions.

185 In addition to these ReLU variants, other kinds of ac-
186 tivation functions have also been developed. For ex-
187 ample, the Swish ($\text{Swish}(x) = x \cdot \text{sigmoid}(\beta x)$) (Ra-
188 machandran et al., 2017) was identified through an au-
189 tomated search using a combination of exhaustive and
190 reinforcement learning as an alternative to ReLU. Its
191 similar shape makes it a reasonable proxy for ReLU in
192 deep learning applications. Mish, defined as $\text{Mish}(x) =$
193 $x \cdot \tanh(\text{softplus}(x))$ (Misra, 2020), exhibits superior em-
194 pirical results compared to ReLU, Swish, and LReLU
195 in CIFAR-10 and ImageNet classification tasks. Frac-
196 tional adaptive linear units FALUs (Zamora et al., 2022)
197 incorporate fractional calculus principles into activation
198 functions, thereby defining a diverse family of activa-
199 tion functions. It has demonstrated enhanced perfor-
200 mance in image classification tasks, improving test ac-
201 curacy. The Seagull activation function, introduced
202 by (Gao and Zhang, 2023), stands out as a customized
203 activation function designed for applications in regres-
204 sion tasks featuring a partially exchangeable target func-
205 tion. It exhibits superiority in addressing the specific
206 demands of regression scenarios.

207 Overall, the above-mentioned activation functions are
208 hard to be generalized across different domains, espe-
209 cially in industrial applications. Another problem is that
210 the lack of theoretical analysis limits the acceptance of
211 these activation functions in spite of their good perfor-
212 mance. Therefore, it is necessary to verify an activation
213 function with a good theoretical guarantee.

214 2.2. Super-Expressive Activation Functions

215 Numerous studies have explored new activation func-
216 tions to make a fixed-size network achieve an arbitrary
217 error, referred to as super-expressive activation func-
218 tions. For example, Maierov and Pinkus (1999) pro-
219 posed an activation function to achieve this goal, but
220 it lacks a closed form and is computationally impracti-
221 cal. Recently, Yarotsky (2021) demonstrated that simple
222 functions such as (\sin, \arcsin) can achieve super-
223 expressiveness, although the relationship between the
224 network size and the dimension was unclear. How-
225 ever, despite the above problems, \sin has been proven
226 to be effective in 3D neural network field, indicating
227 the potential of super-expressiveness in neural networks
228 (Ramirez et al., 2023). Shen et al. (2022) proposed

229 EUAF, showing that a network with EUAF requires 267
 230 only $O(d^2)$ width and $O(1)$ depth to achieve super- 268
 231 expressiveness. The potential of EUAF is demonstrated 269
 232 among simple experiments such as function approximation 270
 233 and Fashion-MNIST classification. They also explored 271
 234 the approximation of a neural network with three 272
 235 hidden layers which is named Floor-Exponential-Step 273
 236 (FLES) networks (Shen et al., 2021). The utilized floor 274
 237 function ($\lfloor x \rfloor$) can be recognized as an activation func- 275
 238 tion with super-expressiveness (Yarotsky, 2021). In a 276
 239 word, these super-expressive activation functions play 277
 240 a theoretically pivotal role in endowing models with 278
 241 the universal approximation property for all continu-
 242 ous functions. However, previous research either lacked
 243 experiments or only included simple ones, leaving it
 244 unknown whether these super-expressive functions are
 245 practically valuable.

246 3. Enriching the Family of Super-expressive Activa- 247 tion Functions

248 In this section, we aim to significantly expand the
 249 scope of EUAF activation function by introducing a
 250 comprehensive collection of activation functions, each
 251 with approximation properties akin to those of EUAF.
 252 For simplicity, let $\mathcal{NN}_\varrho\{N, L; \mathbb{R}^d \rightarrow \mathbb{R}^n\}$ denote the set
 253 of neural networks $\phi : \mathbb{R}^d \rightarrow \mathbb{R}^n$ that can be represented
 254 by ϱ -activated networks, with a maximum width of N
 255 and a maximum depth of L . Let \mathcal{A} represent the set
 256 of all super-expressive activation functions $\varrho : \mathbb{R} \rightarrow \mathbb{R}$,
 257 which satisfy the following conditions:

- 258 • There exists an interval (α, β) with $\alpha < \beta$ where ϱ
 259 is real analytic and non-polynomial on (α, β) .
- 260 • There exists a fixed-size ϱ -activated network ϕ that
 261 can reproduce a triangle-wave function on $[0, \infty)$,
 262 i.e.,

$$\phi(x) = \left| x - 2 \left\lfloor \frac{x+1}{2} \right\rfloor \right| \quad \forall x \in [0, \infty).$$

263 We denote $\overline{\mathcal{A}}$ as the ‘‘closure’’ of \mathcal{A} . This means a func-
 264 tion ϱ is in $\overline{\mathcal{A}}$ if and only if, for any $A > 0$ and $\varepsilon > 0$,
 265 there exists a $\varrho_\varepsilon \in \mathcal{A}$ such that:

$$|\varrho_\varepsilon(x) - \varrho(x)| < \varepsilon \quad \forall x \in [-A, A].$$

Theorem 1. *Given any $\varrho \in \overline{\mathcal{A}}$, the hypothesis
 space*

$$\mathcal{NN}_\varrho\{O(d^2), O(1); \mathbb{R}^d \rightarrow \mathbb{R}\}$$

*is dense in $C([a, b]^d)$ in terms of the supremum
 norm.*

It is crucial to highlight that the constants in the $O(\cdot)$
 notation in Theorem 1 can be explicitly determined and
 depend only on the choice of ϱ . The proof of Theorem 1
 will be provided later in this section.

Before giving the proof, let us provide several exam-
 ples in $\overline{\mathcal{A}}$. The first example, $\varrho_1 \in \overline{\mathcal{A}}$, exhibits an S-
 shape and is defined as follows:

$$\varrho_1 := \begin{cases} \frac{x}{1-x} & \text{for } x \leq 0, \\ \frac{x}{1+x} + \frac{g(x)}{x^2+10} & \text{for } x > 0, \end{cases}$$

where $g(x) = |x - 2\lfloor \frac{x+1}{2} \rfloor|$ for any $x \in \mathbb{R}$.

The second example, $\varrho \in \overline{\mathcal{A}}$, resembles the ReLU
 activation function and is defined as follows:

$$\varrho_2 := \begin{cases} 0 & \text{for } x \leq 0, \\ x + \frac{g(x)}{x+1} & \text{for } x > 0. \end{cases}$$

The third example, $\varrho_1 \in \mathcal{A} \subseteq \overline{\mathcal{A}}$, is defined as fol-
 lows:

$$\varrho_3 := \begin{cases} \frac{2}{\pi} \arcsin(x) & \text{for } -1 \leq x \leq 1, \\ \sin(\frac{\pi}{2}x) & \text{for } |x| > 1. \end{cases}$$

See Figure 2 for visual representations of ϱ_1 , ϱ_2 , and ϱ_3 .

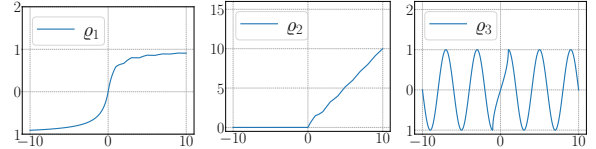


Figure 2: Illustrations of ϱ_1 , ϱ_2 , and ϱ_3 .

Now, we will focus on proving the validity of Theo-
 rem 1. Given any $f \in C([a, b]^d)$ and $\varepsilon > 0$, our goal
 is to construct $\phi \in \mathcal{NN}_\varrho\{O(d^2), O(1); \mathbb{R}^d \rightarrow \mathbb{R}^n\}$ such
 that

$$|\phi(\mathbf{x}) - f(\mathbf{x})| < \varepsilon \quad \forall \mathbf{x} \in [a, b]^d.$$

Several concepts used to establish Theorem 1 can be
 traced back to the research conducted by (Shen et al.,
 2022) and (Yarotsky, 2021). The proof can be divided
 into three main steps as follows.

- The primary objective of the first step is to create
 a neural network that effectively approximates the
 univariate function $f \in C([0, 1])$ within a specific
 ‘‘half’’ interval.

Theorem 2. Given any $f \in C([0, 1])$, $\varrho \in \overline{\mathcal{A}}$, $\varepsilon > 0$, and $K \in \mathbb{N}$, suppose for any $x_1, x_2 \in [0, 1]$, it holds that

$$|f(x_1) - f(x_2)| < \varepsilon/2 \text{ if } |x_1 - x_2| < 1/K. \quad (1)$$

Then there exists $\phi \in \mathcal{NN}_{\varrho}\{\mathcal{O}(1), \mathcal{O}(1); \mathbb{R} \rightarrow \mathbb{R}\}$ such that

$$|\phi(x) - f(x)| < \varepsilon \text{ for any } x \in \bigcup_{k=0}^{K-1} [\frac{2k}{2K}, \frac{2k+1}{2K}].$$

293

- The second step's aim is to utilize the outcome of the first step, Theorem 2, to build a network that effectively approximates the function $f \in C([a, b])$ within the entire interval $[a, b]$.

Theorem 3. Given any $f \in C([a, b])$, $\varrho \in \overline{\mathcal{A}}$, and $\varepsilon > 0$, there exists $\phi \in \mathcal{NN}_{\varrho}\{\mathcal{O}(1), \mathcal{O}(1); \mathbb{R} \rightarrow \mathbb{R}\}$ such that

$$|\phi(x) - f(x)| < \varepsilon \text{ for any } x \in [a, b].$$

298

- The ultimate objective of the final step is to generalize the one-dimensional findings described in Theorem 3 to the multi-dimensional scenario. To achieve this, we will utilize Kolmogorov's superposition theorem (KST) (Kolmogorov, 1957), summarized in Theorem 4. It is important to note that the target function $f \in C([a, b]^d)$ can be appropriately rescaled to facilitate the application of KST.

Theorem 4 (KST). There exist continuous functions $h_{i,j} \in C([0, 1])$ for $i = 0, 1, \dots, 2d$ and $j = 1, 2, \dots, d$ such that any continuous function $f \in C([0, 1]^d)$ can be represented as

$$f(\mathbf{x}) = \sum_{i=0}^{2d} g_i \left(\sum_{j=1}^d h_{i,j}(x_j) \right)$$

for any $\mathbf{x} = (x_1, \dots, x_d) \in [0, 1]^d$, where $g_i : \mathbb{R} \rightarrow \mathbb{R}$ is a continuous function for each $i \in \{0, 1, \dots, 2d\}$.

308

We observe that it is sufficient to demonstrate the case where $\varrho \in \mathcal{A}$ rather than $\varrho \in \overline{\mathcal{A}}$, aided by the following lemma.

309

310

311

Lemma 1 (Proposition 10 of (Zhang et al., 2024)).

Given two functions $\varrho, \tilde{\varrho} : \mathbb{R} \rightarrow \mathbb{R}$ with $\tilde{\varrho} \in C(\mathbb{R})$, suppose for any $M > 0$, there exists $\tilde{\varrho}_\eta \in \mathcal{NN}_{\varrho}\{\tilde{N}, \tilde{L}; \mathbb{R} \rightarrow \mathbb{R}\}$ for each $\eta \in (0, 1)$ such that

$$\tilde{\varrho}_\eta(x) \rightrightarrows \tilde{\varrho}(x) \text{ as } \eta \rightarrow 0^+ \text{ for any } x \in [-M, M].$$

Assuming $\phi_{\tilde{\varrho}} \in \mathcal{NN}_{\tilde{\varrho}}\{N, L; d \rightarrow n\}$, for any $\varepsilon > 0$ and $A > 0$, there exists $\phi_{\varrho} \in \mathcal{NN}_{\varrho}\{\tilde{N} \cdot N, \tilde{L} \cdot L; \mathbb{R}^d \rightarrow \mathbb{R}^n\}$ such that

$$\|\phi_{\varrho} - \phi_{\tilde{\varrho}}\|_{\sup([-A, A]^d)} < \varepsilon.$$

312

Now let's prove the utilized theorems.

312

313

314

315

316

3.1. Proof of Theorem 2

Partition $[0, 1]$ into $2K$ small intervals \mathcal{I}_k and $\tilde{\mathcal{I}}_k$ for $k = 1, 2, \dots, K$, i.e.,

$$\mathcal{I}_k = [\frac{2k-2}{2K}, \frac{2k-1}{2K}] \text{ and } \tilde{\mathcal{I}}_k = [\frac{2k-1}{2K}, \frac{2k}{2K}].$$

317

318

319

Clearly, $[0, 1] = \bigcup_{k=1}^K (\mathcal{I}_k \cup \tilde{\mathcal{I}}_k)$. Let x_k be the right endpoint of \mathcal{I}_k , i.e., $x_k = \frac{2k-1}{2K}$ for $k = 1, 2, \dots, K$.

See an illustration of \mathcal{I}_k , $\tilde{\mathcal{I}}_k$, and x_k in Figure 3 for the case $K = 5$. Our objective is to construct $\phi \in$

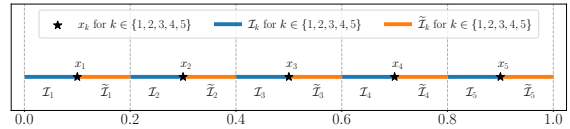


Figure 3: An illustration of \mathcal{I}_k and $\tilde{\mathcal{I}}_k$ for $k \in \{1, 2, \dots, K\}$ with $K = 5$.

320

321

322

323

324

325

326

327

328

$\mathcal{NN}_{\varrho}\{\mathcal{O}(1), \mathcal{O}(1); \mathbb{R} \rightarrow \mathbb{R}\}$ to achieve accurate approximations of f within \mathcal{I}_k for $k = 1, 2, \dots, K$. It is not essential to consider the values of ϕ within $\tilde{\mathcal{I}}_k$ for all k . In other words, our focus is primarily on achieving accurate approximations within one "half" of the interval $[0, 1]$, which is the crucial element in our proof.

Define $\psi(x) := x - \sigma(x)$ for any $x \in \mathbb{R}$, where $\sigma \in \mathcal{NN}_{\varrho}\{\mathcal{O}(1), \mathcal{O}(1); \mathbb{R} \rightarrow \mathbb{R}\}$ with

$$\sigma(x) = \left\lfloor x - 2 \left\lfloor \frac{x+1}{2} \right\rfloor \right\rfloor \text{ for } x \geq 0.$$

329

330

See Figure 4 for an illustration of ψ .

It is easy to verify that

$$\psi(2Kx)/2 + 1 = k \text{ for any } x \in [\frac{2k-2}{2K}, \frac{2k-1}{2K}] = \mathcal{I}_k. \quad (2)$$

331

332

We will make use of the two following lemmas to simplify our proof.

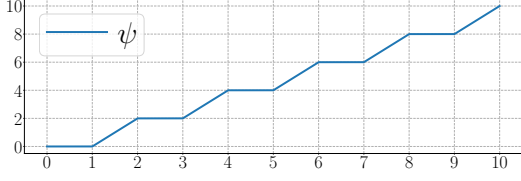


Figure 4: An illustration of ψ on $[0, 10]$.

Lemma 2 (Lemma 23 of Shen et al. (2022)). *Given any rationally independent numbers a_1, a_2, \dots, a_K for any $K \in \mathbb{N}^+$ and an arbitrary periodic function $g : \mathbb{R} \rightarrow \mathbb{R}$ with period T , i.e., $g(x + T) = g(x)$ for any $x \in \mathbb{R}$, assume there exist $x_1, x_2 \in \mathbb{R}$ with $0 < x_2 - x_1 < T$ such that g is continuous on $[x_1, x_2]$. Then the following set*

$$\left\{ (u \cdot g(wa_1) + v, \dots, u \cdot g(wa_K) + v) : u, w, v \in \mathbb{R} \right\}$$

is dense in \mathbb{R}^K provided that

$$\min_{x \in [x_1, x_2]} g(x) < \max_{x \in [x_1, x_2]} g(x).$$

Lemma 3. *Given $K \in \mathbb{N}^+$, suppose ϱ is real analytic and non-polynomial on an interval (α, β) with $\beta > \alpha$. Then there exists $w_0 \in (-\frac{\beta-\alpha}{2K}, \frac{\beta-\alpha}{2K})$ such that $\varrho(\frac{\alpha+\beta}{2} + kw_0)$, for $\{k = 1, 2, \dots, K\}$, are rationally independent.*

Proof. We prove this lemma by contradiction. If it does not hold, then $\varrho(\frac{\alpha+\beta}{2} + kw)$, for $\{k = 1, 2, \dots, K\}$, are rationally dependent for any $w \in (-\frac{\beta-\alpha}{2K}, \frac{\beta-\alpha}{2K}) = \mathcal{I}$. That means, for any $w \in \mathcal{I}$, there exists $\lambda = (\lambda_1, \dots, \lambda_K) \in \mathbb{Q}^K \setminus \{\mathbf{0}\}$ such that $\sum_{k=1}^K \lambda_k \varrho(\frac{\alpha+\beta}{2} + kw) = 0$. We observe that \mathcal{I} is uncountable and $\mathbb{Q}^K \setminus \{\mathbf{0}\}$ is countable. It follows that there exists $\lambda = (\lambda_1, \dots, \lambda_K) \in \mathbb{Q}^K \setminus \{\mathbf{0}\}$ such that $\sum_{k=1}^K \lambda_k \varrho(\frac{\alpha+\beta}{2} + kw) = 0$ for all w in an uncountable subset of \mathcal{I} . Then the real analyticity of ϱ implies $\sum_{k=1}^K \lambda_k \varrho(\frac{\alpha+\beta}{2} + kw) = 0$ for all $w \in \mathcal{I}$. By expanding $\sum_{k=1}^K \lambda_k \varrho(\frac{\alpha+\beta}{2} + kw)$ into the Taylor series at $w = 0$, we get the identity $\sum_{k=1}^K \lambda_k k^m = 0$ for each m with $\frac{d^m \varrho}{dw^m}(\frac{\alpha+\beta}{2}) \neq 0$. Since ϱ is non-polynomial on $(\alpha, \beta) \ni \frac{\alpha+\beta}{2}$, there are infinitely many m with $\frac{d^m \varrho}{dw^m}(\frac{\alpha+\beta}{2}) \neq 0$, implying $\sum_{k=1}^K \lambda_k k^m = 0$. This means $\lambda = (\lambda_1, \dots, \lambda_K) = \mathbf{0}$, a contradiction with $\lambda \in \mathbb{Q}^K \setminus \{\mathbf{0}\}$. So we finish the proof of Lemma 3. \square

Now, let us return to the proof of Theorem 2. We can employ Lemma 3 to produce a collection of rationally independent numbers. Specifically, there exists a

value w_0 such that a_1, a_2, \dots, a_K are linearly independent, where each a_k is defined as $a_k = \varrho(\frac{\alpha+\beta}{2} + kw_0)$.

Next, define

$$g(x) = \left| x - 2 \left\lfloor \frac{x+1}{2} \right\rfloor \right| \quad \text{for } x \in \mathbb{R}.$$

By Lemma 2, there exists $u_1, w_1, v_1 \in \mathbb{R}$ such that

$$\left| u_1 \cdot g(w_1 a_k) + v_1 - f(x_k) \right| < \varepsilon/2 \quad \text{for any } k.$$

Since $\sigma(x) = g(x)$ for any $x \geq 0$ and g is periodic with period 2, we can choose a sufficiently large $m_0 \in \mathbb{N}$ such that

$$\begin{aligned} & \left| u_1 \sigma(w_1 a_k + 2m_0) + v_1 - f(x_k) \right| \\ &= \left| u_1 g(w_1 a_k + 2m_0) + v_1 - f(x_k) \right| \\ &= \left| u_1 g(w_1 a_k) + v_1 - f(x_k) \right| < \varepsilon/2, \end{aligned}$$

for $k = 1, 2, \dots, K$. Define

$$\phi(x) = u_1 \sigma \left(w_1 \varrho \left(\frac{\alpha+\beta}{2} + \left(\frac{\psi(2kx)}{2} - 1 \right) w_0 \right) + 2m_0 \right) + v_1.$$

For any $x \in \mathcal{I}_k$, we have

$$\begin{aligned} \phi(x) &= u_1 \sigma \left(w_1 \varrho \left(\frac{\alpha+\beta}{2} + \left(\frac{\psi(2kx)}{2} - 1 \right) w_0 \right) + 2m_0 \right) + v_1 \\ &= u_1 \sigma \left(w_1 \varrho \left(\frac{\alpha+\beta}{2} + kw_0 \right) + 2m_0 \right) + v_1 \\ &= u_1 \sigma \left(w_1 a_k + 2m_0 \right) + v_1, \end{aligned}$$

implying

$$|\phi(x) - f(x)| \leq \underbrace{|\phi(x) - f(x_k)|}_{< \varepsilon/2} + \underbrace{|f(x_k) - f(x)|}_{< \varepsilon/2 \text{ by (1)}} < \varepsilon.$$

It follows that

$$|\phi(x) - f(x)| < \varepsilon \quad \text{for any } x \in \bigcup_{k=0}^{K-1} \left[\frac{2k}{2K}, \frac{2k+1}{2K} \right].$$

Moreover, we can easily verify $\phi \in \mathcal{NN}_{\varrho}\{\mathcal{O}(1), \mathcal{O}(1); \mathbb{R} \rightarrow \mathbb{R}\}$. So we finish the proof of Theorem 2.

3.2. Proof of Theorem 3 based on Theorem 2.

We claim it suffices to prove the special case $[a, b] = [0, \frac{1}{2}]$ as this simplification readily extends to the broader scenario. To see this, we simply introduce a linear function $\mathcal{L} : [0, \frac{1}{2}] \rightarrow [a, b]$ by defining $\mathcal{L}(x) = 2(b-a)x + a$. The special case implies $f \circ \mathcal{L} : [0, \frac{1}{2}] \rightarrow \mathbb{R}$ can be approximated by a network $\tilde{\phi}$ arbitrarily well. Then $\phi = \tilde{\phi} \circ \mathcal{L}^{-1}$ can approximate $f : [a, b] \rightarrow \mathbb{R}$ well, as desired.

378 We can continuously extend f from $[0, \frac{1}{2}]$ to \mathbb{R} by
 379 setting $f(x) = f(0)$ if $x < 0$ and $f(x) = f(\frac{1}{2})$ if $x > \frac{1}{2}$. It
 380 follows from the uniform continuity of f on $[-1, 2]$ that
 381 there exists $K = K(f, \varepsilon) \in \mathbb{N}^+$ with $K \geq 2$ such that for
 382 any $x_1, x_2 \in [-1, 2]$,

$$|f(x_1) - f(x_2)| < \varepsilon/10 \quad \text{if } |x_1 - x_2| < 1/K.$$

383 For $i = 1, 2, 3, 4$, define

$$f_i(x) := f(x - \frac{i}{4K}) \quad \text{for any } x \in [0, 1].$$

384 Then, for $i = 1, 2, 3, 4$ and $x_1, x_2 \in [0, 1]$, we have

$$|f_i(x_1) - f_i(x_2)| < \varepsilon/10 = \tilde{\varepsilon}/2 \quad \text{if } |x_1 - x_2| < 1/K,$$

385 where $\tilde{\varepsilon} = \varepsilon/5$. For each $i \in \{1, 2, 3, 4\}$, by Theorem 2,
 386 there exists $\phi_i \in \mathcal{NN}_g(\mathcal{O}(1), \mathcal{O}(1); \mathbb{R} \rightarrow \mathbb{R})$ such that

$$|\phi_i(x) - f_i(x)| < \tilde{\varepsilon} = \varepsilon/5 \quad \text{for any } x \in \bigcup_{k=0}^{K-1} [\frac{2k}{2K}, \frac{2k+1}{2K}].$$

387 Define

$$\psi(x) = \sigma(x + 1 - \sigma(x + 1)) \quad \text{for any } x \in \mathbb{R},$$

388 where $\sigma \in \mathcal{NN}_g(\mathcal{O}(1), \mathcal{O}(1); \mathbb{R} \rightarrow \mathbb{R})$ with

$$\sigma(x) = \lfloor x - 2\lfloor \frac{x+1}{2} \rfloor \rfloor \quad \text{for } x \geq 0.$$

389 See an illustration of ψ on $[0, 2K]$ for $K = 5$ in Figure 5.

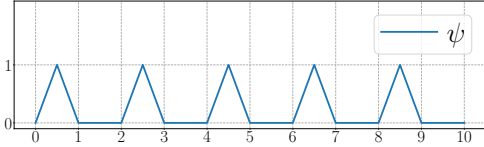


Figure 5: An illustration of ψ on $[0, 2K]$ for $K = 5$.

390 Clearly, $0 \leq \psi(2Kx) \leq 1$ for any $x \in [0, 1]$, from
 391 which we deduce

$$|(\phi_i(x) - f_i(x))\psi(2Kx)| < \varepsilon/5 \quad \forall x \in \bigcup_{k=0}^{K-1} [\frac{2k}{2K}, \frac{2k+1}{2K}].$$

392 Observe that $\psi(y) = 0$ for $y \in \bigcup_{k=0}^{K-1} [2k+1, 2k+2]$,
 393 which implies

$$\psi(2Kx) = 0 \quad \text{for any } x \in \bigcup_{k=0}^{K-1} [\frac{2k+1}{2K}, \frac{2k+2}{2K}].$$

394 Subsequently, by the fact

$$[0, 1] = \left(\bigcup_{k=0}^{K-1} [\frac{2k}{2K}, \frac{2k+1}{2K}] \right) \cup \left(\bigcup_{k=0}^{K-1} [\frac{2k+1}{2K}, \frac{2k+2}{2K}] \right),$$

we have

$$|(\phi_i(x) - f_i(x))\psi(2Kx)| < \varepsilon/5 \quad \text{for any } x \in [0, 1]. \quad (3)$$

396 For each $i \in \{1, 2, 3, 4\}$ and any $z \in [0, \frac{1}{2}] \subseteq [0, 1 - \frac{1}{K}]$
 397 $\subseteq [0, 1 - \frac{i}{4K}]$, we have

$$y_i = z + \frac{i}{4K} \in [\frac{i}{4K}, 1] \subseteq [0, 1].$$

398 By bringing $x = y_i \in [0, 1]$ into Equation (3), we get

$$\begin{aligned} \varepsilon/5 &> |(\phi_i(y_i) - f_i(y_i))\psi(2Ky_i)| \\ &= |\phi_i(y_i)\psi(2Ky_i) - f_i(y_i)\psi(2Ky_i)| \\ &= |\phi_i(z + \frac{i}{4K})\psi(2K(z + \frac{i}{4K})) - f_i(z + \frac{i}{4K})\psi(2K(z + \frac{i}{4K}))| \\ &= |\phi_i(z + \frac{i}{4K})\psi(2Kz + \frac{i}{2}) - f(z)\psi(2Kz + \frac{i}{2})| \end{aligned}$$

399 for any $z \in [0, \frac{1}{2}]$, where the last equality comes from
 400 the fact that $f_i(x) = f(x - \frac{i}{4K})$ for any $x \in [0, 1] \supseteq$
 401 $[\frac{i}{4K}, 1]$. Define

$$\tilde{\phi}(x) := \sum_{i=1}^4 \phi_i(x + \frac{i}{4K})\psi(2Kx + \frac{i}{2}) \quad \text{for any } x \in [0, \frac{1}{2}].$$

402 It is easy to verify that $\sum_{i=1}^4 \psi(x + \frac{i}{2}) = 1$ for any $x \geq 0$
 403 based on the definition of ψ . See Figure 6 for illus-
 404 trations. It follows that $\sum_{i=1}^4 \psi(2Kz + \frac{i}{2}) = 1$ for any
 405 $z \in [0, \frac{1}{2}]$.

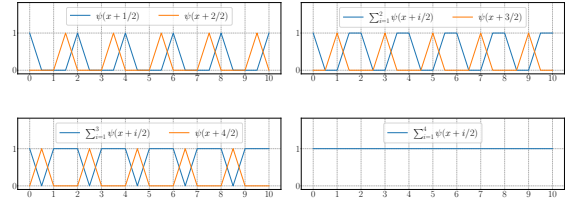


Figure 6: Illustrations of $\sum_{i=1}^4 \psi(x + i/2) = 1$ for any $x \in [0, 10]$.

Hence, for any $z \in [0, \frac{1}{2}]$, we have

$$\begin{aligned} &|\tilde{\phi}(z) - f(z)| \\ &= \left| \sum_{i=1}^4 \phi_i(z + \frac{i}{4K})\psi(2Kz + \frac{i}{2}) - f(z) \sum_{i=1}^4 \psi(2Kz + \frac{i}{2}) \right| \\ &\leq \sum_{i=1}^4 \left| \phi_i(z + \frac{i}{4K})\psi(2Kz + \frac{i}{2}) - f(z)\psi(2Kz + \frac{i}{2}) \right| \\ &< 4 \cdot \frac{\varepsilon}{5} = \frac{4\varepsilon}{5}. \end{aligned}$$

To approximate $(x, y) \mapsto xy$ well, we define

$$\Gamma_\delta(x, y) := \frac{\varrho(x_0 + \delta x + \delta y) - \varrho(x_0 + \delta y) - \varrho(x_0 + \delta x) + \varrho(x_0)}{\delta^2 \varrho''(x_0)}$$

408 for any $x, y \in \mathbb{R}$, where $\rho''(x_0) \neq 0$. Clearly, $\Gamma_\delta(x, y) \rightarrow$ 425 and
 409 xy as $\delta \rightarrow 0$. Then we can define

$$\phi_\delta(x) := \sum_{i=1}^4 \Gamma_\delta\left(\phi_i\left(x + \frac{i}{4K}\right), \psi\left(2Kx + \frac{i}{2}\right)\right) \quad \forall x \in \left[0, \frac{1}{2}\right]. \quad 426$$

410 Clearly, $\phi_\delta \in \mathcal{NN}_\varrho\{\mathcal{O}(1), \mathcal{O}(1); \mathbb{R} \rightarrow \mathbb{R}\}$. Moreover, we
 411 can choose a sufficiently small $\delta_0 > 0$ such that

$$|\phi_{\delta_0}(x) - \tilde{\phi}(x)| < \varepsilon/5 \quad \text{for any } x \in \left[0, \frac{1}{2}\right]. \quad 427$$

412 By defining $\phi := \phi_{\delta_0} \in \mathcal{NN}_\varrho\{\mathcal{O}(1), \mathcal{O}(1); \mathbb{R} \rightarrow \mathbb{R}\}$, we
 413 have

$$|\phi(x) - f(x)| \leq \underbrace{|\phi_{\delta_0}(x) - \tilde{\phi}(x)|}_{< \varepsilon/5} + \underbrace{|\tilde{\phi}(x) - f(x)|}_{< 4\varepsilon/5} < \varepsilon$$

414 for any $x \in \left[0, \frac{1}{2}\right]$. So we finish the proof of Theorem 3.

415 3.3. Proof of Theorem 1 based on Theorem 3 and KST.

416 We can safely assume that $[a, b] = [0, 1]$ since the
 417 general case can be readily extended by incorporating
 418 an affine map such as $\mathcal{L}(a) = (b - a)x + a$. Given any
 419 $f \in C([0, 1]^d)$, by KST, there exist $h_{i,j} \in C([0, 1])$ and
 420 $g_i \in C(\mathbb{R})$ for $i = 0, 1, \dots, 2d$ and $j = 1, 2, \dots, d$ such
 421 that

$$f(\mathbf{x}) = \sum_{i=0}^{2d} g_i\left(\sum_{j=1}^d h_{i,j}(x_j)\right) \quad \forall \mathbf{x} = (x_1, \dots, x_d) \in [0, 1]^d. \quad 431$$

422 Choose a sufficiently large $A > 0$, e.g.,

$$A = 1 + \sup \left\{ \left| \sum_{j=1}^d h_{i,j}(x_j) \right| : i = 0, 1, \dots, 2d, \mathbf{x} \in [0, 1]^d \right\}. \quad 433$$

423 Then for any $\delta > 0$, by Theorem 3, there exist $\psi_{i,j}, \phi_i \in$ 436
 424 $\mathcal{NN}_\varrho\{\mathcal{O}(1), \mathcal{O}(1); \mathbb{R} \rightarrow \mathbb{R}\}$ such that

$$|g_i(t) - \phi_i(t)| < \delta \quad \text{for any } t \in [-A, A] \quad 439$$

$$|h_{i,j}(t) - \psi_{i,j}(t)| < \delta \quad \text{for any } t \in [0, 1],$$

for $i = 0, 1, \dots, 2d$ and $j = 1, 2, \dots, d$. By defining

$$\phi(\mathbf{x}) = \sum_{i=0}^{2d} \phi_i\left(\sum_{j=1}^d \psi_{i,j}(x_j)\right) \quad \forall \mathbf{x} = (x_1, \dots, x_d) \in \mathbb{R}^d.$$

427 we have $\phi \in \mathcal{NN}_\varrho\{\mathcal{O}(d^2), \mathcal{O}(1); \mathbb{R} \rightarrow \mathbb{R}\}$. See an illus-
 428 tration of the architecture of ϕ in Figure 7. Moreover,
 429 by choosing sufficiently small $\delta > 0$, we can conclude
 430 that

$$|\phi(\mathbf{x}) - f(\mathbf{x})| < \varepsilon \quad \forall \mathbf{x} \in [0, 1]^d,$$

which means we finish the proof of Theorem 1.

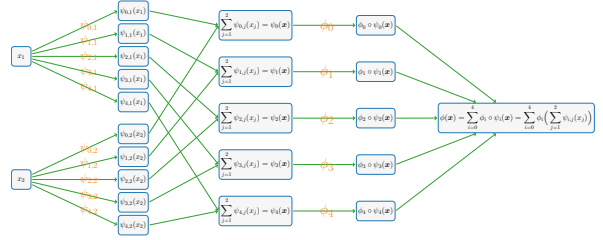


Figure 7: An illustration of the target network realizing ϕ for any $\mathbf{x} \in [a, b]^d$ in the case of $d = 2$. This network contains $(2d + 1)d + (2d + 1) = (d + 1)(2d + 1)$ sub-networks that realize $\psi_{i,j}$ and ϕ_i for $i = 0, 1, \dots, 2d$ and $j = 1, 2, \dots, d$.

432 4. Experimental Results

433 To further validate the efficacy of our activation func-
 434 tions, we evaluate PEUAF against a wide range of base-
 435 line activation functions, including LReLU (Xu et al.,
 436 2015), PReLU (He et al., 2015), Softplus (Zheng
 437 et al., 2015), ELU (Clevert et al., 2015), SELU (Klam-
 438 bauer et al., 2017), ReLU (Nair and Hinton, 2010) and
 439 Swish (Ramachandran et al., 2017). We conduct these

Table 1: A Brief Description of Three Image Datasets and Four Industrial Fault Diagnosis Datasets

Dataset	Description
CIFAR-10	60,000 32×32 resolution RGB images in 10 categories (6,000 images per category)
Tiny ImageNet	100,000 64×64 RGB images in 200 categories (500 for each category)
ImageNet	14,197,122 RGB images over 1,000 categories and 21,841 subcategories
Case Western Reserve University (CWRU)	2,400 vibration signals with 10 types of faults in drive end. Each signals has 1,024 samples
Power Quality Disturbance (PQD)	11200 voltage disturbance signals in 16 types, each disturbance signal at each fault has additive white Gaussian noise
Motor Fault (MF)	6 types of faults and each kind of fault has at least 290 samples (Sun et al., 2023)
Electrical Fault Detection and Classification (EFDC)	12,000 samples with 6 types of faults. Each sample has 6 features including the measured line currents and voltages.

440 comparisons across four industrial signal datasets and
 441 three image datasets (CIFAR-10 (Krizhevsky et al.,
 442 2009), Fashion-MNIST (Xiao et al., 2017) and ImageNet
 443 (Deng et al., 2009)) using three distinct neural
 444 network architectures (LeNet-type (Lecun et al., 1998),
 445 ResNet-18 (He et al., 2016), and VGG-16 (Simonyan
 446 and Zisserman, 2015)).

447 As discussed in (Shen et al., 2022), only a few neurons
 448 with super-expressive activation functions are required
 449 to approximate functions with arbitrary precision to avoid
 450 large generalization errors. However, implementing this
 451 in practical experiments is challenging. Therefore, our
 452 experiments primarily focus on exploring the feature
 453 patterns of PEUAF and determining how it contributes to
 454 improving test accuracy, instead of targeting 100% test
 455 accuracy.

456 4.1. Experimental Setups

457 The datasets used in our experiments are briefly introduced
 458 in Table 1. For each experiment, we train the models
 459 with a batch size of 64 using the “NAdam” optimizer
 460 (Dozat, 2016), with an initial learning rate of 0.01.
 461 The learning rate decays with a factor of 0.2 if the
 462 accuracy change over 5 consecutive epochs is no more
 463 than 1×10^{-4} . We set the number of epochs to 300
 464 to ensure proper convergence. The baseline network
 465 structures employed in our experiments are introduced in
 466 Tables 2 and 3.

Table 2: Baseline A for the CWRU, PQD, and MF datasets.

Layer	Parameters	Activation
1D-Convolution (3×1)	filter size=64, stride = 1	PEUAF
1D-Convolution (3×1)	filter size=64, stride = 1	PEUAF
Batch-normalization (BN)	momentum=0.99, epsilon=0.001	-
Max-pooling	pool size= 3×1 , stride = 1	-
1D-Convolution (3×1)	filter size=64, stride = 1	PEUAF
1D-Convolution (3×1)	filter size=64, stride = 1	PEUAF
Batch-normalization (BN)	momentum=0.99, epsilon=0.001	-
Max-pooling	pool size= 3×1 , stride = 1	-
1D-Convolution (3×1)	filter size=64, stride = 1	PEUAF
1D-Convolution (3×1)	filter size=64, stride = 1	PEUAF
Batch-normalization (BN)	momentum=0.99, epsilon=0.001	-
Global-average-pooling	-	-
Fully connected	size (chosen by tasks)	softmax

467 The most critical hyperparameter is the range of the
 468 adaptive frequency w . To determine this, we conducted
 469 a classification experiment with different w values on
 470 the PQD dataset (A et al.), as illustrated in Figure 8.
 471 The network structure used is Baseline A, a 1D convolutional
 472 neural network. To emphasize the discrepancies in
 473 outcomes, we employ a logarithmic transformation

Table 3: Baseline B for the EFDC dataset.

Layer	Parameters	Activation
1D-Convolution (2×1)	filter size=16, stride = 1	PEUAF
Batch-normalization (BN)	momentum=0.99, epsilon=0.001	-
Max-pooling	pool size= 2×1 , stride = 1	-
1D-Convolution (2×1)	filter size=16, stride = 1	PEUAF
Batch-normalization (BN)	momentum=0.99, epsilon=0.001	-
Max-pooling	pool size= 2×1 , stride = 1	-
Flatten	-	-
Fully connected	size (chosen by tasks)	softmax

474 (log) during the visualization of the loss function. Figure
 475 8 shows the training and validation curves, while Table 4
 476 provides the corresponding test accuracy. The table
 477 reveals two key points: First, when w exceeds 1, the
 478 test accuracy drops significantly, indicating that higher
 479 frequencies pose challenges to the PEUAF’s ability to
 480 effectively extract PQD features. Second, when w lies
 481 in the range of $[0, 1]$, the test accuracy consistently
 482 remains above 98%. Therefore, we reasonably conclude
 483 that the frequency w should be constrained within the
 484 range of $[0, 1]$.

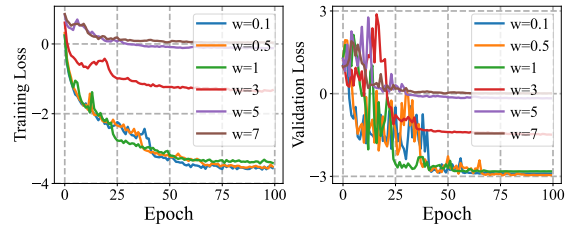


Figure 8: The training and validation loss with different w . (a) the training loss; (b) the validation loss

Table 4: The accuracy of PQD classification with different w .

w	0.1	0.5	1	3	5	7
Accuracy	98.25%	98.30%	98.04%	93.66%	75.26%	53.39%

485 4.2. Analysis Experiments

486 In this section, we conduct experiments to show the
 487 characteristics of PEUAF. For the larger datasets (CWRU,
 488 PQD, and MF), we utilize the Baseline A in Table 2,
 489 while for the EFDC dataset, we use the Baseline B in
 490 Table 3. Baseline B is smaller than Baseline A due to
 491 the smaller size of the EFDC dataset compared to CWRU,
 492 PQD, and MF. Our comparison focuses not only on the
 493 overall performance but also on the convergence behavior
 494 during the training process, fluctuations in the validation
 495 process, and a detailed mechanism analysis.

Table 5: Test accuracy in industrial fault diagnosis datasets.

Model	CWRU	PQD	MF	EFDC	Avg Rank
LReLU	99.58%	98.12%	98.82%	84.24%	4
PReLU	97.08%	97.85%	99.17%	84.50%	7
Softplus	95.00%	97.95%	98.06%	84.37%	8
ELU	100.00%	98.15%	99.70%	83.10%	2
SELU	99.17%	98.12%	99.85%	83.74%	3
ReLU	99.58%	97.89%	99.70%	83.35%	5
Swish	99.16%	98.15%	97.35%	84.63%	6
PEUAF	100.00%	98.17%	100.00%	85.64%	1

497 **Performance.** Table 5 summarizes the performance
498 of several activation functions. All results are the average
499 over three runs. On the EFDC dataset, PEUAF
500 takes the lead by the largest margin, *i.e.*, surpassing
501 the second place Swish by over 1%. On the CWRU,
502 dataset, PEUAF exhibits competitive performance compared
503 to Swish, ReLU, SELU, ELU, and LReLU, while
504 PEUAF outperforms Softplus and PReLU by 2%
505 and 5%, respectively. On the PQD dataset, all activation
506 functions achieve similar test accuracy. Lastly, on
507 the MF dataset, PEUAF shows similar performance with
508 ReLU, SELU, ELU, and PReLU but surpasses Swish
509 and Softplus. Overall, PEUAF proves to be a competent
510 activation function on four industrial fault diagnosis
511 datasets.

512 To further evaluate the effectiveness of PEUAF,
513 we conducted occlusion experiments in two classic
514 datasets: CWRU and PQD. For each dataset, the occluding
515 sizes and strides were set to 100 and 50, respectively.
516 The occluded pixels were all replaced by zeros.
517 Based on the results in Figure 9, we observe that PEUAF
518 outperforms ReLU in locating faults. The experiments
519 reveal two distinct levels of performances: (1) In the
520 PQD dataset, PEUAF and ReLU show similar performance
521 in accurately detecting and localizing faults, as
522 illustrated in Figure 9. This can be attributed to the
523 favorable condition within the PQD dataset, characterized
524 by its low signal-to-noise ratio, contributing to the
525 successful faults localization. However, such ideal
526 conditions are rare in real-world scenarios. (2) In contrast,
527 in the CWRU dataset, PEUAF significantly outperforms
528 ReLU as shown in Figure 9. Despite that both PEUAF
529 and ReLU achieve commendable test accuracy, ReLU
530 tends to capture more holistic features instead of locating
531 the real fault, which makes the outputs less reliable.
532 Conversely, PEUAF effectively locates faults even in
533 the presence of noise interference, offering valuable
534 insights into the timing and severity of fault occurrences,
535 as indicated by the occlusion experiments.

536 **Convergence.** Since PEUAF has a unique shape,

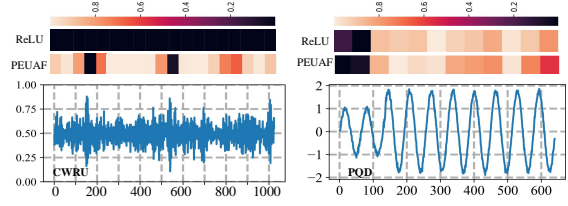


Figure 9: Occlusion experiments of Baselines using PEUAF and ReLU in CWRU and PQD datasets. It is seen that the baseline with PEUAF can better localize the fault when the original signal is noisy.

537 there might be concerns that such an oscillating function
538 could be difficult to optimize. To address this, we
539 compare the training dynamics of PEUAF and ReLU.
540 Figure 10 shows the training and validation curves
541 of PEUAF and ReLU on the CWRU, PQD, and MF
542 datasets. Below are our detailed analyses:

- 543 **1. Convergence speed during training:** The conver-
544 gence rate during the training process is notably
545 influenced by the choice of activation functions and
546 the inherent characteristics of the dataset. All
547 the experiments in Figure 10 consistently demon-
548 strate the superior convergence speed of PEUAF.
549 In dataset with a lower signal-to-noise ratio (such
550 as the PQD dataset), PEUAF and ReLU show sim-
551 ilar convergence speed. In contrast, in noise-free
552 datasets or those with high signal-to-noise ratio,
553 models adopting the PEUAF activation function
554 display significantly faster convergence.
- 555 **2. Convergence effect during training:** The choice
556 of activation functions can impact the convergence
557 effect, particularly in terms of oscillations or fluc-
558 tuations during the training process. In the PQD
559 dataset, the convergence patterns of PEUAF and
560 ReLU are relatively similar, except for some fluc-
561 tuations occurring around the 50th epoch. How-
562 ever, for the MF dataset, noticeable oscillations
563 occur during convergence, particularly within the
564 epoch range between 150 to 200. For the CWRU
565 dataset, the fluctuations happen at around the 50th
566 epoch when using ReLU as the activation function.
567 Therefore, PEUAF helps reduce the oscillation of
568 training losses and improves the training perfor-
569 mance.
- 570 **3. Fluctuation during validation:** In addition to the
571 training process, the effectiveness of PEUAF can
572 also be observed during the validation process.
573 Across all the datasets, PEUAF outperforms ReLU
574 by showing less fluctuation in the validation pro-
575 cess. For the CWRU dataset, both PEUAF and

ReLU exhibit fluctuations at the start of the validation process. However, after a sudden drop in validation loss after approximately 20 epochs, PEUAF shows smaller validation loss fluctuations than ReLU. For the PQD dataset, the validation loss curve for PEUAF and ReLU appear similar, but the amplitude of fluctuations is smaller for PEUAF. The most significant difference in fluctuation patterns is particularly obvious in the MF dataset, where ReLU exhibits high frequency and amplitude of fluctuations. This behavior can potentially be attributed to the fact that, in noise-free data settings, ReLU tends to capture global features initially, rather than precisely pinpointing fine-grained fault details, unlike PEUAF.

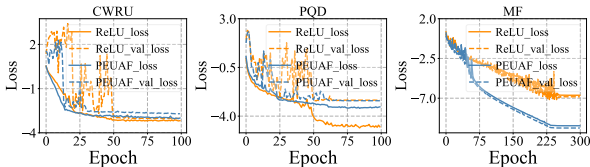


Figure 10: Training dynamics of Baselines using PEUAF and ReLU in three large datasets.

4.3. Comparative Experiments

In this section, we demonstrate combining the super-expressive activation function (PEUAF) with the baseline activation function can enhance the generalization ability of neural networks.

CIFAR-10. In this experiment, we augment the dataset by rotating, shifting, shearing, and horizontally flipping the original images. We mainly focus on the ResNet structure. Table 6 compares ResNet-18a with a mixed activation function to several identical model topologies using ReLU. The mixed activation function achieves a 0.89% error reduction. Tables 7 and 8 separately summarize the test accuracy of ResNet-18 and Vit-B/16 (Dosovitskiy et al., 2020) across various baseline activation functions and mixed activation functions. Notably, the mixed activation function improves the test accuracy, especially in Softplus, which increased by 2.72% and 5.01%.

To further explain the efficacy of mixed activation functions, Figure 11 provides a detailed comparison of the loss and accuracy among ReLU, PEUAF and mixed activation functions. When exclusively applying PEUAF in the CIFAR-10 classification task, both the training convergence and fluctuations are worse than those of ReLU, as shown in the loss curve in Figure 11. However, the ResNet-18 using mixed activation

Table 6: CIFAR-10 Classification error vs the number of parameters, for common compact model architectures vs. ResNet-18a + Mixed ReLU.

Neural Network	#Param	Error%
All-CNN (Springenberg et al., 2014)	1.3M	7.25%
MobileNetV1 (Howard et al., 2017)	3.2M	10.76%
MobileNetV2 (Sandler et al., 2018)	2.24M	7.22%
ShuffleNet 8G (Zhang et al., 2018)	0.91M	7.71%
ShuffleNet 1G (Zhang et al., 2018)	0.24M	8.56%
HENet (Duan et al., 2018)	0.7M	10.16%
ResNet-18a+ReLU	0.27M	8.75%
ResNet-18a+ mixed ReLU	0.27M	7.82%

functions outperforms the models using either ReLU or PEUAF alone. The mixed approach results in smoother loss and accuracy curves during both the training and validation process.

The occlusion experiments further demonstrate that the mixed activation functions can enhance the neural network’s ability to identify essential features. The occlusion sizes and strides are set to 4 and 2, respectively, with occluded pixels replaced by zeros. As in Figure 12, the results provide a clear illustration of this phenomenon. The models using only ReLU or PEUAF successfully identify a multitude of features contributing to the classification. However, they also select too many unnecessary pixel points, recognizing part of the surroundings as the important features for classification. In contrast, the mixed activation function model can accurately locate the critical features while ignoring irrelevant pixels.

Table 7: Comparisons of classification accuracy across several activation functions using ResNet for CIFAR-10.

Activation	Test accuracy
ResNet-18+PEUAF	90.00% / -
ResNet-18+LReLU/Mixed	92.42% / 94.13%
ResNet-18+PReLU/Mixed	92.29% / 94.23%
ResNet-18+Softplus/Mixed	89.28% / 92.09%
ResNet-18+ELU/Mixed	91.09% / 92.11%
ResNet-18+SELU/Mixed	90.47% / 91.32%
ResNet-18+ReLU/Mixed	93.02% / 93.91%
ResNet-18+Swish/Mixed	94.07% / 92.99%
ResNet-34+ReLU/Mixed	93.70% / 94.23%

Tiny-ImageNet. The Tiny-Imagenet dataset is utilized to further demonstrate the expressiveness of PEUAF. The model is trained for 100 epochs with an initial learning rate of 0.1, which decays by an order of magnitude every 30 epochs, using a batch size of 256. Table 9 compares the test accuracy of ResNet-

Table 8: Comparisons of classification accuracy across several activation functions using Vit-B/16 for CIFAR-10.

Activation	Test accuracy
Vit-B/16+PEUAF	90.58% / -
Vit-B/16+LReLU/Mixed	91.15% / 91.31%
Vit-B/16+PReLU/Mixed	90.40% / 90.47%
Vit-B/16+Softplus/Mixed	74.23% / 79.24%
Vit-B/16+ELU/Mixed	89.69% / 89.80%
Vit-B/16+SELU/Mixed	87.26% / 87.37%
Vit-B/16+ReLU/Mixed	89.43% / 89.44%
Vit-B/16+Swish/Mixed	90.66% / 89.65%
Vit-B/16+GELU/Mixed	97.49% / 97.90%

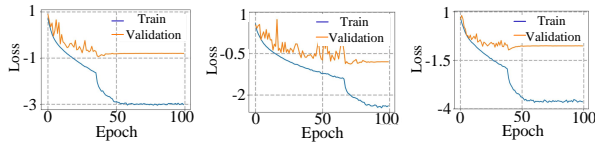


Figure 11: Loss of CIFAR-10 experiments among three activation functions. From left to right are the loss of ResNet-18 using ReLU, PEUAF and mixed activation function.

Table 9: Comparisons of classification accuracy across several activation functions using ResNet-18 for Tiny-ImageNet.

Activation	Test accuracy
ResNet-18+PEUAF	56.86% / -
ResNet-18+LReLU/Mixed	62.39% / 62.29%
ResNet-18+PReLU/Mixed	59.57% / 60.81%
ResNet-18+Softplus/Mixed	56.98% / 57.75%
ResNet-18+ELU/Mixed	59.44% / 59.88%
ResNet-18+SELU/Mixed	59.51% / 59.62%
ResNet-18+ReLU/Mixed	63.40% / 63.42%
ResNet-18+Swish/Mixed	60.76% / 59.53%

18 with several baseline activation functions on Tiny-ImageNet. By replacing the activation functions in the last block, the ResNet-18 with mixed activation functions achieves competitive results, showing slight improvements in the test accuracy across most experiments, except for Swish and PEUAF.

ImageNet. The ImageNet dataset is used to evaluate the effectiveness of PEUAF in large datasets. Due to memory limitations, the model follows the setup from the previous research (Liu et al., 2022b), except for the number of neurons in the first layer and the data enhancement. The neurons of the first layer is reduced to 256 from 512. Table 10 compares the test accuracy between PReLU and the mixed activation. In this large-scale image classification experiment, the drawbacks of using PEUAF alone become apparent, with the accuracy of ResNet-18 with PEUAF being only 63.38%, which is lower than that of PReLU. However, the ResNet-18 with

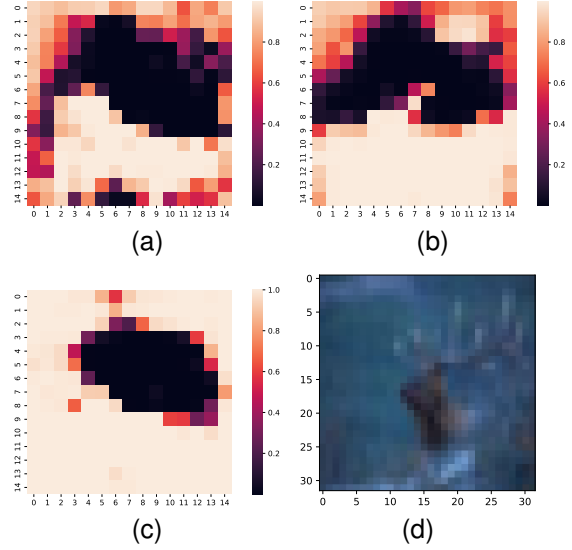


Figure 12: Occlusion experiments on the CIFAR-10 dataset among three activation functions to examine their discriminative ability. (a)~(c) Occlusion experiments using ReLU, PEUAF and mixed activation function. (d) the original figure.

659 mixed activation functions achieves competitive results.

Table 10: Comparisons of classification accuracy across several activation functions using ResNet-18 for ImageNet.

Activation	Test accuracy
ResNet-18+PEUAF	63.38% / -
ResNet-18+LReLU/Mixed	70.65% / 70.96%

660 5. Conclusion and Discussion

This paper provides an in-depth analysis of the characteristics and effectiveness of PEUAF, particularly focusing on its application to industrial and image datasets. By testing the trainable frequency w , we have determined an optimal frequency range for w within the interval $[0, 1]$. To further demonstrate the super-expressiveness of PEUAF, we have conducted experiments using four industrial datasets and three benchmark image datasets. The results indicate that PEUAF surpasses ReLU in terms of convergence speed, oscillation during training, fluctuation during validation, and fault localization ability, especially in industrial datasets with a high signal-to-noise ratio. Additionally, the mixed activation function outperforms the single activation function in most image classification tasks.

Looking ahead, the future of activation function research is promising. The development of PEUAF paves

678 the way for exploring other super-expressive activa-
 679 tion functions that could further enhance neural net-
 680 work performance across various applications. Future
 681 research could focus on expanding the family of super-
 682 expressive activation functions and investigating their
 683 practical utility in more diverse and complex datasets.
 684 Moreover, combining PEUAF with other state-of-the-
 685 art neural network architectures and exploring its bene-
 686 fits in real-world scenarios could yield valuable insights.
 687 The adaptability and effectiveness of PEUAF in hand-
 688 ling stationary signals suggest potential applications
 689 in fields such as signal processing, fault diagnosis, and
 690 time-series analysis.

691 References

692 A, R.M., A, A.C., B, J.B., C, Y.B., A, Y.L., . Open source dataset gen-
 693 erator for power quality disturbances with deep-learning reference
 694 classifiers - sciencedirect. *Electric Power Systems Research* 195.
 695 Apicella, A., Donnarumma, F., Isgrò, F., Prevete, R., 2021. A survey
 696 on modern trainable activation functions. *Neural Networks* 138,
 697 14–32.
 698 Bingham, G., Miikkulainen, R., 2022. Discovering parametric activa-
 699 tion functions. *Neural Networks* 148, 48–65.
 700 Clevert, D.A., Unterthiner, T., Hochreiter, S., 2015. Fast and accurate
 701 deep network learning by exponential linear units (ELUs). *arXiv*
 702 preprint arXiv:1511.07289 .
 703 Deng, J., Dong, W., Socher, R., Li, L.J., Li, K., Fei-Fei, L., 2009. Im-
 704 agenet: A large-scale hierarchical image database, in: 2009 IEEE
 705 conference on computer vision and pattern recognition, IEEE. pp.
 706 248–255.
 707 Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X.,
 708 Unterthiner, T., Dehghani, M., Minderer, M., Heigold, G., Gelly,
 709 S., et al., 2020. An image is worth 16x16 words: Transformers for
 710 image recognition at scale. *arXiv preprint arXiv:2010.11929* .
 711 Dozat, T., 2016. Incorporating nesterov momentum into adam. *ICLR*
 712 workshop.
 713 Duan, J., Zhang, R., Huang, J., Zhu, Q., 2018. The speed improve-
 714 ment by merging batch normalization into previously linear layer
 715 in cnn, in: 2018 International Conference on Audio, Language and
 716 Image Processing (ICALIP), IEEE. pp. 67–72.
 717 Fan, F., Li, M., Teng, Y., Wang, G., 2020. Soft autoencoder and its
 718 wavelet adaptation interpretation. *IEEE Transactions on Computa-*
 719 *tional Imaging* 6, 1245–1257.
 720 Gao, F., Zhang, B., 2023. Data-aware customization of activa-
 721 tion functions reduces neural network error. *arXiv preprint*
 722 *arXiv:2301.06635* .
 723 Goldenstein, N., Sulam, J., Romano, Y., 2024. Pivotal auto-encoder
 724 via self-normalizing relu. *IEEE Transactions on Signal Processing*
 725 , 1–12doi:10.1109/TSP.2024.3418971.
 726 He, K., Zhang, X., Ren, S., Sun, J., 2015. Delving deep into rectifiers:
 727 Surpassing human-level performance on imagenet classification,
 728 in: *Proceedings of the IEEE International Conference on Computer*
 729 *Vision*, pp. 1026–1034.
 730 He, K., Zhang, X., Ren, S., Sun, J., 2016. Deep residual learning
 731 for image recognition, in: *Proceedings of the IEEE conference on*
 732 *computer vision and pattern recognition*, pp. 770–778.
 733 Hendrycks, D., Gimpel, K., 2016. Gaussian error linear units
 734 (GELUs). *arXiv e-prints* , arXiv:1606.08415doi:10.48550/
 735 [arXiv.1606.08415](https://arxiv.org/abs/1606.08415), arXiv:1606.08415.

736 Howard, A.G., Zhu, M., Chen, B., Kalenichenko, D., Wang, W.,
 737 Weyand, T., Andreetto, M., Adam, H., 2017. Mobilenets: Effi-
 738 cient convolutional neural networks for mobile vision applications.
 739 *arXiv preprint arXiv:1704.04861* .
 740 Klambauer, G., Unterthiner, T., Mayr, A., Hochreiter, S., 2017. Self-
 741 normalizing neural networks. *Advances in neural information pro-*
 742 *cessing systems* 30.
 743 Kolmogorov, A.N., 1957. On the representation of continuous func-
 744 tions of many variables by superposition of continuous functions
 745 of one variable and addition. *Doklady Akademii Nauk SSSR* 114,
 746 953–956. URL: <http://mi.mathnet.ru/dan22050>.
 747 Krizhevsky, A., Hinton, G., et al., 2009. Learning multiple layers of
 748 features from tiny images. *Technical Report TR-2009*, Toronto,
 749 ON, Canada.
 750 LeCun, Y., Bengio, Y., Hinton, G., 2015. Deep learning. *nature* 521,
 751 436–444.
 752 Lecun, Y., Bottou, L., Bengio, Y., Haffner, P., 1998. Gradient-based
 753 learning applied to document recognition. *Proceedings of the IEEE*
 754 86, 2278–2324. doi:10.1109/5.726791.
 755 Lee, S., Sim, B., Ye, J.C., 2024. Magnitude and angle dynamics in
 756 training single relu neurons. *Neural Networks* , 106435.
 757 Liu, J., Duan, Z., Liu, H., 2024a. A grid fault diagnosis framework
 758 based on adaptive integrated decomposition and cross-modal at-
 759 tention fusion. *Neural Networks* , 106400.
 760 Liu, Y., Wang, Z., Ma, Q., Shen, H., 2022a. Multistability analy-
 761 sis of delayed recurrent neural networks with a class of piecewise
 762 nonlinear activation functions. *Neural Networks* 152, 80–89.
 763 Liu, Z., Li, S., Wu, D., Liu, Z., Chen, Z., Wu, L., Li, S.Z., 2022b.
 764 Automix: Unveiling the power of mixup for stronger classifiers,
 765 in: *European Conference on Computer Vision*, Springer. pp. 441–
 766 458.
 767 Liu, Z., Lv, Q., Lee, C.H., Shen, L., 2024b. Segmenting medical
 768 images with limited data. *Neural Networks* 177, 106367.
 769 Maiorov, V., Pinkus, A., 1999. Lower bounds for approximation
 770 by MLP neural networks. *Neurocomputing* 25, 81–91. doi:10.
 771 [1016/S0925-2312\(98\)00111-8](https://doi.org/10.1016/S0925-2312(98)00111-8).
 772 Misra, D., 2020. Mish: A self regularized non-monotonic activation
 773 function. *arXiv:1908.08681*.
 774 Nair, V., Hinton, G.E., 2010. Rectified linear units improve restricted
 775 boltzmann machines, in: *Proceedings of the 27th international con-*
 776 *ference on machine learning (ICML-10)*, pp. 807–814.
 777 Ramachandran, P., Zoph, B., Le, Q.V., 2017. Searching for activation
 778 functions. *arXiv preprint arXiv:1710.05941* .
 779 Ramirez, P.Z., De Luigi, L., Sirocchi, D., Cardace, A., Spezialetti, R.,
 780 Ballerini, F., Salti, S., Di Stefano, L., 2023. Deep learning on 3d
 781 neural fields. *arXiv preprint arXiv:2312.13277* .
 782 Sandler, M., Howard, A., Zhu, M., Zhmoginov, A., Chen, L.C., 2018.
 783 Mobilenetv2: Inverted residuals and linear bottlenecks, in: *Pro-*
 784 *ceedings of the IEEE conference on computer vision and pattern*
 785 *recognition*, pp. 4510–4520.
 786 Shen, Z., Yang, H., Zhang, S., 2021. Neural network approximation:
 787 Three hidden layers are enough. *Neural Networks* 141, 160–173.
 788 URL: [https://doi.org/10.1016%2Fj.neunet.2021.](https://doi.org/10.1016%2Fj.neunet.2021.04.011)
 789 [04.011](https://doi.org/10.1016/j.neunet.2021.04.011), doi:10.1016/j.neunet.2021.04.011.
 790 Shen, Z., Yang, H., Zhang, S., 2022. Deep network approxima-
 791 tion: Achieving arbitrary accuracy with fixed number of neu-
 792 rons. *Journal of Machine Learning Research* 23, 1–60. URL:
 793 <http://jmlr.org/papers/v23/21-1404.html>.
 794 Simonyan, K., Zisserman, A., 2015. Very deep convolutional net-
 795 works for large-scale image recognition, in: Bengio, Y., LeCun, Y.
 796 (Eds.), *3rd International Conference on Learning Representations,*
 797 *ICLR 2015*, San Diego, CA, USA, May 7-9, 2015, Conference
 798 Track Proceedings. URL: [http://arxiv.org/abs/1409.](http://arxiv.org/abs/1409.1556)
 799 [1556](https://arxiv.org/abs/1409.1556).
 800 Sitzmann, V., Martel, J., Bergman, A., Lindell, D., Wetzstein, G.,

801 2020. Implicit neural representations with periodic activation functions. *Advances in neural information processing systems* 33, 7462–7473.

802
803

804 Springenberg, J.T., Dosovitskiy, A., Brox, T., Riedmiller, M.A.,
805 2014. Striving for simplicity: The all convolutional net. *CoRR*
806 abs/1412.6806. URL: <https://api.semanticscholar.org/CorpusID:12998557>.

807
808 Sun, Z., Machlev, R., Wang, Q., Belikov, J., Levron, Y., Baimel, D.,
809 2023. A public data-set for synchronous motor electrical faults
810 diagnosis with cnn and lstm reference classifiers. *Energy and AI*
811 14, 100274.

812 Wang, C., Liang, J., Deng, Q., 2024. Dynamics of heterogeneous
813 hopfield neural network with adaptive activation function based on
814 memristor. *Neural Networks* 178, 106408.

815 Xiao, H., Rasul, K., Vollgraf, R., 2017. Fashion-mnist: a novel
816 image dataset for benchmarking machine learning algorithms.
817 [arXiv:1708.07747](https://arxiv.org/abs/1708.07747).

818 Xu, B., Wang, N., Chen, T., Li, M., 2015. Empirical evaluation of
819 rectified activations in convolutional network. *Computer ence* .

820 Yarotsky, D., 2021. Elementary superexpressive activations, in: *International Conference on Machine Learning*, PMLR. pp. 11932–
821 11940.

822
823 Zamora, J., Rhodes, A.D., Nachman, L., 2022. Fractional adaptive
824 linear units, in: *Proceedings of the AAAI Conference on Artificial*
825 *Intelligence*, pp. 8988–8996.

826 Zhang, S., Lu, J., Zhao, H., 2024. Deep network approximation:
827 Beyond ReLU to diverse activation functions. *Journal of Machine*
828 *Learning Research* 25, 1–39. URL: <http://jmlr.org/papers/v25/23-0912.html>.

829
830 Zhang, X., Zhou, X., Lin, M., Sun, J., 2018. Shufflenet: An extremely
831 efficient convolutional neural network for mobile devices, in: *Proceedings of the IEEE conference on computer vision and pattern*
832 *recognition*, pp. 6848–6856.

833
834 Zheng, H., Yang, Z., Liu, W., Liang, J., Li, Y., 2015. Improving
835 deep neural networks using softplus units, in: *2015 International*
836 *Joint Conference on Neural Networks (IJCNN)*, pp. 1–4. doi:10.
837 [1109/IJCNN.2015.7280459](https://doi.org/10.1109/IJCNN.2015.7280459).