

• 专题二:双清论坛“工程科学融合人工智能的关键前沿基础科学问题” •  
DOI: 10.3724/BNSFC-2025-0032

## 工程系统深度神经网络鲁棒性:前沿与展望\*

金龙<sup>1\*\*</sup> 陈良铭<sup>1</sup> 张强强<sup>2\*\*</sup> 李惠<sup>3</sup>

1. 兰州大学 信息科学与工程学院,兰州 730000
2. 兰州大学 土木工程与力学学院,兰州 730000
3. 哈尔滨工业大学 土木工程学院,哈尔滨 150001

**[摘要]** 深度神经网络正为传统工程系统带来一场深刻的范式革命。然而,主流深度神经网络模型固有的脆弱性与工程系统对鲁棒性的极致要求之间形成了尖锐矛盾,即存在“鲁棒性鸿沟”,已成为制约人工智能在工程系统中广泛应用的关键瓶颈之一。目前,尚无针对工程系统深度神经网络鲁棒性的综合性讨论。本文旨在系统性地介绍并分析关于深度神经网络鲁棒性机理、工程系统深度神经网络鲁棒性提升方法的重要研究。首先,在问题发现层面,本文解构了鲁棒性鸿沟问题的内涵与外延;随后,在根源剖析层面,本文介绍了关于深度神经网络鲁棒性缺陷原因分析的理论进展、鲁棒性与多层次网络架构的关系等;进一步地,在现有对策层面,探讨增强工程系统深度神经网络鲁棒性的通用方法,以及针对工业制造、电网系统、自动驾驶等关键工程领域特性的深度神经网络鲁棒性提升方法;最后,在未来展望层面,本文重点讨论内在鲁棒的深度神经网络新架构、新型学习范式、嵌入工程语义约束的新策略等前沿方向,以期构建鲁棒的工程系统深度神经网络提供有价值的参考。

**[关键词]** 人工智能;深度神经网络;工程系统;鲁棒性

以大模型等为代表的深度神经网络,作为新一轮科技与产业变革的核心驱动力,正以前所未有的深度和广度重塑现代工程系统的面貌<sup>[1]</sup>。现代深度神经网络强大的非线性表达能力、对未见场景的泛化性,为工程领域中长期存在的高维复杂问题提供了强有力的解决方案<sup>[2-4]</sup>。这场智能化范式变革深刻地体现在工程系统的全生命周期中,正将生产效率与自动化水平推向新的高度。在智能制造、自动驾驶、智慧医疗等工程系统的多个环节,深度神经网络正从辅助工具逐渐演变为“中枢大脑”,将成为驱动我国工程技术创新的关键。

然而,现代深度神经网络的统计范式固有的脆弱性与工程系统对鲁棒性的严格要求之间存在尖锐的矛盾,本文将其定义为工程系统深度神经网络的“鲁棒性鸿沟”。如图1所示,互联网应用(常规计算机视觉、自然语

言处理任务等)、普通工业应用、安全攸关工程系统对鲁棒性的需求,与深度神经网络在这些应用领域的实际鲁棒性水平之间的差距急剧扩大,进而形成了鲁棒性鸿沟。实践层面,当前深度神经网络极易受到对抗性攻击的影响(即存在内在脆弱性)。攻击者通过在输入数据(如图像、传感器信号甚至物理场景)中精心添加人类难以察觉的微小扰动,便可诱使模型做出灾难性的错误判断。此外,生产实践中广泛存在各类天然噪声,亦会对模型性能造成巨大影响。深度神经网络鲁棒性不足的特性,对于国防安全、智能制造、自动驾驶、医疗影像分析等系统而言,是非常严重的安全隐患。这一特性与安全攸关的工程系统对鲁棒性的严格要求之间,形成了巨大的矛盾,使得在诸多对安全性和鲁棒性有严格要求的领域,深度神经网络很难被广泛地信赖和使用。理论层

收稿日期:2025-07-30; 修回日期:2025-10-30

\* 本文根据国家自然科学基金委员会第409期“双清论坛”讨论的内容整理。

\*\* 通信作者,Email:jinlongsysu@foxmail.com; zhangqq@lzu.edu.cn

本文受到国家自然科学基金项目(62506148, 62476115, U25A20229)、青年教师科研创新能力支持项目(SRICSPYF-BS2025061)和甘肃省科技计划项目(24ZYQA048)的资助。

**引用格式:** 金龙,陈良铭,张强强,等. 工程系统深度神经网络鲁棒性:前沿与展望. 中国科学基金,2026,40(1):118-127.

Jin L, Chen LM, Zhang QQ, et al. Robustness of deep neural networks in engineering systems: Frontiers and outlook. Bulletin of National Natural Science Foundation of China, 2026, 40(1): 118-127. (in Chinese)

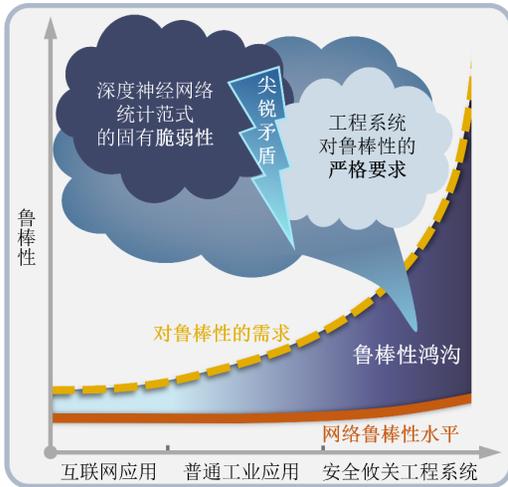


图1 工程系统深度神经网络的鲁棒性鸿沟示意图

Fig.1 Schematic Diagram of the Robustness Gap in Deep Neural Networks for Engineering Systems

面,鲁棒性鸿沟体现在两个方面。(1)理论研究对该鸿沟的进一步印证:现有理论研究揭示了模型鲁棒性与准确性、效率之间存在难以调和的权衡<sup>[5,6]</sup>,这与工程系统同时追求高鲁棒、高精度、高效率的需求相悖;不仅如此,深度神经网络的可证实鲁棒性在数据分布不确定的情况下存在理论上界<sup>[7]</sup>,这与工程系统对鲁棒性的极致追求形成了鲜明反差。(2)工程系统对应的场景超出了当前理论所能刻画的范围:大量鲁棒性理论研究集中于 $L_p$ 范数扰动<sup>[8]</sup>,难以描述工程系统在真实世界中面临物理与语义层面的复杂扰动,如相机丢帧、时间戳不同步、恶劣天气、迷惑性信息等。即便是对于 $L_p$ 范数扰动的分析,也大多要求扰动半径较小<sup>[9]</sup>,与开放世界中无明确限制的扰动形成对比。

鲁棒性鸿沟并非源于单一维度的差距,而是多方面系统性差异的结果。相比于主要任务为常规计算机视觉、自然语言处理等的系统(以下称为通用系统),工程系统在从数据扰动、输出约束、误差容忍、失效后果、适用范围到部署验证的全链路中,均体现出关于鲁棒性的巨大差异。

(1)数据扰动差异:通用系统与工程系统中的扰动来源具有明确差异。通用系统中的深度神经网络鲁棒性研究,其扰动通常是高斯分布、均匀分布等相对确定分布下的噪声,或是根据一定规则设计的对抗性扰动。而工程系统中,深度神经网络必须直面物理世界的高度不确定性:从极端天气、电磁干扰等外部环境扰动,到传感器老化漂移、执行器磨损等内部状态不确定性,这些扰动复杂、动态且难以精确建模。

(2)输出约束差异:通用系统与工程系统中,模型的

输出是否需要满足各类约束也存在差异。通用系统中使用的主流深度神经网络模型通常具有物理不完备性,其输出可能与基本物理定律(如能量守恒、运动学约束)或工程规范相悖。即便采用了物理信息嵌入的方法,也难以穷尽所有应当满足的物理定律与工程规范。

(3)误差容忍差异:通用系统与工程系统可容忍的可靠性量级也具有巨大差异。在通用系统中,95%以上的准确率已是优异表现;但在航空航天、自动驾驶等领域,工程系统要求的准确率常需要高于99.9999%,故障率须被控制在百万分之一甚至更低水平。

(4)失效后果差异:工程系统追求极致可靠性的主要原因在于系统失效的后果通常更为严重。推荐系统的一次错误推送,其代价是用户体验的轻微下降;而航空航天系统的一次感知失误,或电网调度系统的一次错误决策,其后果可能是难以挽回的生命财产损失和重大的社会安全事故。

(5)适用范围差异:通用系统的鲁棒性主要追求性能维持,即在扰动下尽可能保持准确性。然而,在安全攸关的工程系统中,最高优先级是保障系统安全,并在安全性、性价比与准确性之间进行权衡,而非仅仅考虑准确性。例如,一个鲁棒的自动驾驶系统,在识别出超出其安全运行能力的暴风雪时,更合理的行为是安全降级,自动寻找安全地点停车,而非仅仅提升暴风雪中的行驶性能。这种知难而退的灵活性,对维持工程系统鲁棒性十分关键。

(6)部署验证差异:为确保各种鲁棒性增强策略的最终可靠性,工程系统需要一套与通用系统不同的验证保障措施,以确保系统在部署前尽可能安全。传统工程产品需通过严格的形式化验证与安全标准认证方可部署。然而,当前通用系统下神经网络的安全性评估大多依赖基于有限验证集的经验性评估,缺乏数学上可证明的安全边界。

综上所述,鲁棒性鸿沟是一个全链路系统性难题。弥合这一鸿沟,是推动神经网络在关键工程领域实现可信赖应用的核心前提。

近年来,一些综述对神经网络鲁棒性进行了系统总结。然而,目前缺乏对工程系统场景下,神经网络鲁棒性的综述。为促进弥合鲁棒性鸿沟,在上文已进行问题阐述的基础上,本文按照根源分析、现有对策、未来出路的层层递进关系,总结并分析了工程系统神经网络鲁棒性的研究前沿,包括神经网络鲁棒性不足的原因、鲁棒性与网络架构的关系,以及提升工程系统神经网络鲁棒性的通用方法、针对性方法等;在此基础上,本文针对现存挑战,展望了工程系统深度

神经网络鲁棒性领域具有潜力与研究价值的未来关键研究方向。

## 1 深度神经网络鲁棒性不足的根源分析

深度神经网络的鲁棒性是指网络在面临输入数据扰动时,仍能保持其性能(如准确性、稳定性)的能力<sup>[9-13]</sup>。该扰动不仅包括对抗性扰动(即精心设计以使模型性能变差的微小扰动),也包括现实世界中常见的自然扰动,如图像亮度变化、传感器噪声、环境噪声等。理解深度神经网络鲁棒性不足的根源,是构建鲁棒性提升对策的前提。

现代深度神经网络所依赖的统计范式虽取得了令人瞩目的进展,但也带来了天然的、系统性的脆弱倾向。理论上,多种因素共同导致了深度神经网络系统的鲁棒性不足,其中关键的理论解释包括数据的高维特性、模型的特征偏好、鲁棒性与准确性的权衡等。

### 1.1 数据的高维特性

目前,具有挑战性的深度学习任务,其数据通常具有高维特性,为鲁棒性带来严峻挑战。在高维空间中,大量微小扰动会随维度增加而累积,进而显著改变模型的输出<sup>[14]</sup>。实际上,对于任何错误率非零的分类器,均存在对抗样本,且对抗半径与输入维度有关<sup>[15]</sup>,而实现鲁棒性所需的样本量是关于输入维度的多项式级<sup>[16]</sup>。类似地,对抗性攻击的成功率随输入数据的维度增加而呈现多项式级别的增长,表明在高维空间中模型鲁棒性显著不足<sup>[17]</sup>。进一步的研究表明,鲁棒分类器的存在性取决于数据分布在低维子空间的集中程度,且在高维任务中,除非对数据流形做出极强的假设,否则不可避免地会存在对抗样本,造成模型性能急剧下降<sup>[18]</sup>。

### 1.2 模型的特征偏好

除数据层面带来的鲁棒性挑战外,模型层面的隐式特征偏好也是造成深度神经网络鲁棒性不足的重要原因之一。在学习算法、网络架构、优化算法等的共同作用下,深度神经网络的特征具有一些与人类截然不同的隐式偏好<sup>[19,20]</sup>。这些隐式偏好包括对于非鲁棒特征的偏好:在只追求准确率的标准训练下,神经网络会依赖有助于提升准确率但鲁棒性低的特征,而这些特征常常是人类无法感知或认为无关的模式(如特定纹理或像素)<sup>[21]</sup>。在标准训练中,模型倾向于学习不具鲁棒性的特征<sup>[22]</sup>。与标准训练不同,对抗训练会促使学习器转向鲁棒特征,并远离非鲁棒特征<sup>[22]</sup>。与依赖全局形状的人类不同,卷积神经网络模型严重依赖局部纹理线索识别物体<sup>[23]</sup>。

### 1.3 模型架构对鲁棒性的影响

深度神经网络的架构对于其鲁棒性具有关键影响。

按照从宏观到微观的层次,网络架构可分为网络规模、层间连接和层内算子。其中,网络规模指深度、宽度等宏观性质;层间连接指层与层之间是否连接以及如何连接;层内算子指单层对输入施加何种运算。

(1)网络规模:大模型的成功,凸显了增大网络规模对于表达能力和泛化性的重要性。然而,对于鲁棒性,并非网络规模越大越好。更宽且更深的网络通常具有较高的Lipschitz常数,意味着扩大网络规模导致鲁棒性下降是可能的<sup>[24]</sup>。当模型参数不足时,增加宽度会使鲁棒性下降;而在过参数化的情况下,增加宽度有可能增强鲁棒性<sup>[25]</sup>。深度对鲁棒性的影响则更加复杂,主要取决于使用的初始化方式以及训练中参数更新的幅度。具体而言,采用LeCun初始化且训练过程中的参数更新幅度较小的情况下,增加深度有助于提高鲁棒性;然而,若使用神经正切核初始化或He初始化,增加深度反而可能降低鲁棒性<sup>[25]</sup>。

(2)层间连接:经验研究表明,与网络规模这类宏观特性相比,层与层之间的连接方式对模型的鲁棒性有着更加显著的影响<sup>[26]</sup>。通过神经架构搜索结果可以发现,连接较为密集的网络架构通常在鲁棒性方面优于连接较为稀疏的架构<sup>[27]</sup>。在理论方面,基于深度神经网络动力学的研究表明,如果层间连接对应的动力学离散格式具有零稳定性,网络的鲁棒性则更强<sup>[28]</sup>。跨层连接能够有效降低全局字典结构的互协方差,进而提高鲁棒性<sup>[29]</sup>。特别地,通过引入特殊跨层连接以近似隐式欧拉方法,可有效提升鲁棒性<sup>[30]</sup>。此外,残差连接的放置位置也对鲁棒性有着重要影响:在靠近输入的位置增加残差连接,或在靠近输出的位置减少残差连接,均可增强网络的鲁棒性;类似地,增加靠近输入的层宽度,或减少靠近输出的层宽度,也能提升鲁棒性<sup>[24]</sup>。

(3)层内算子:层内算子对鲁棒性的影响颇为显著,并有望从根本上改变网络的鲁棒性。在网络层内设计方面,特定现有算子的搭配有助于提升鲁棒性<sup>[31]</sup>。此外,算子的对称性也起着至关重要的作用:向网络中引入特定的对称性,如使用偶函数作为激活函数,可有效地增强鲁棒性<sup>[32]</sup>。近年来,关于Transformer与卷积神经网络鲁棒性的比较也引起了广泛关注。Transformer更倾向于依赖低频信号,而卷积神经网络更偏向于高频特征,这被认为是Transformer表现出更强鲁棒性的原因之一<sup>[33]</sup>。然而,随后的研究指出,在更公平的设置下,Transformer的鲁棒性优势主要体现在分布外样本的表现上<sup>[34]</sup>。进一步的研究发现,通过对卷积神经网络进行特定调整,如引入图像分块、使用更大卷积核等,其鲁棒性可达到甚至超越Transformer<sup>[35]</sup>。此外,通过利用

$L_\infty$ 距离构建网络基本层,得到的模型具有显著强于常规模型的可认证鲁棒性<sup>[36]</sup>。该思想后续被推广至基于布尔函数的统一Lipschitz网络框架,实现了更优的可认证鲁棒性<sup>[37]</sup>。

### 1.4 鲁棒性与准确性的权衡

除数据和模型层面外,另一个对神经网络鲁棒性的制约因素是鲁棒性和准确性难以兼得,即存在权衡。在相对简单的设定下可证明,任何能完美抵抗有界 $L_\infty$ 范数攻击的模型,必然会在某些干净样本上出错,反之亦然;且经过鲁棒训练所得到的特征与经过标准训练所得到的特征显著不同,但通常与人类感知更一致<sup>[38]</sup>。进一步地,鲁棒性与准确性的权衡可被证明存在,且这种权衡会在数据不平衡的情况下加剧<sup>[39]</sup>。此外,对于数据可分离的二分类问题,只有在网络规模随数据维度指数级增长时,鲁棒性与准确性的权衡才有望消失,即存在维数灾难<sup>[40]</sup>。该结果被后续研究拓展并相互印证:简单地扩大模型规模并不能高效解决鲁棒性与准确性的权衡问题,且当参数量超过某个阈值后,鲁棒性的提升速度会急剧减缓<sup>[5]</sup>。

## 2 工程系统深度神经网络鲁棒性的提升对策

在神经网络鲁棒性不足根源分析研究的指导下,研究者们面向工程系统深度神经网络鲁棒性的提升问题,提出了多种对策。为清晰地讨论这些对策,本节构建了一个提升工程系统人工智能鲁棒性的方法学框架(图2)。本节将首先介绍提升神经网络鲁棒性的

通用对策,这些方法是构建鲁棒工程系统深度学习神经网络的基础。然而,这些通用方法在应用于工程系统时,不一定能适配工程系统对于鲁棒性的独特需求。为此,本节将讨论针对具体工程场景的专用对策,以体现通用方法与工程系统特有需求融合的必要性,并探讨现有对策的局限性。

### 2.1 提升工程系统深度神经网络鲁棒性的通用对策

本节从对抗训练、以数据为中心的方法、模型架构、正则化方法方面介绍提升工程系统深度神经网络鲁棒性的通用方法。

#### 2.1.1 对抗训练

对抗训练是提升模型鲁棒性最有效的方法之一,其目标是通过解决如下极小极大优化问题,以使模型在最坏情况下仍能保持性能<sup>[41]</sup>:

$$\min_{\theta} \mathbb{E}_{(x,y) \sim \mathcal{D}} \left[ \max_{\|p\|_c \leq \epsilon} \left( \mathcal{F}_{\theta}(x+p), y \right) \right] \quad (1)$$

其中, $\theta$ 为模型参数, $\mathbb{E}[\cdot]$ 代表求期望, $\mathcal{Y}$ 为标签, $\mathcal{D}$ 为数据分布, $\epsilon$ 为扰动半径, $\ell$ 为损失函数。对抗训练的开端是快速梯度符号法(Fast Gradient Sign Method,FGSM),其显著提高了模型的鲁棒性<sup>[14]</sup>。进一步地,FGSM对抗训练的思想被推广至多步投影梯度下降(Projected Gradient Descent,PGD)<sup>[42]</sup>。然而,传统对抗训练的算法每轮需多次前后向传播,通常面临巨大的额外计算开销。为处理这一问题,Shafahi等<sup>[43]</sup>提出在同一次反向传播中同时更新模型参数与生成对抗扰动的“免费对抗训练”方法,仅需很小额外计算代价即可大幅提升鲁棒性。

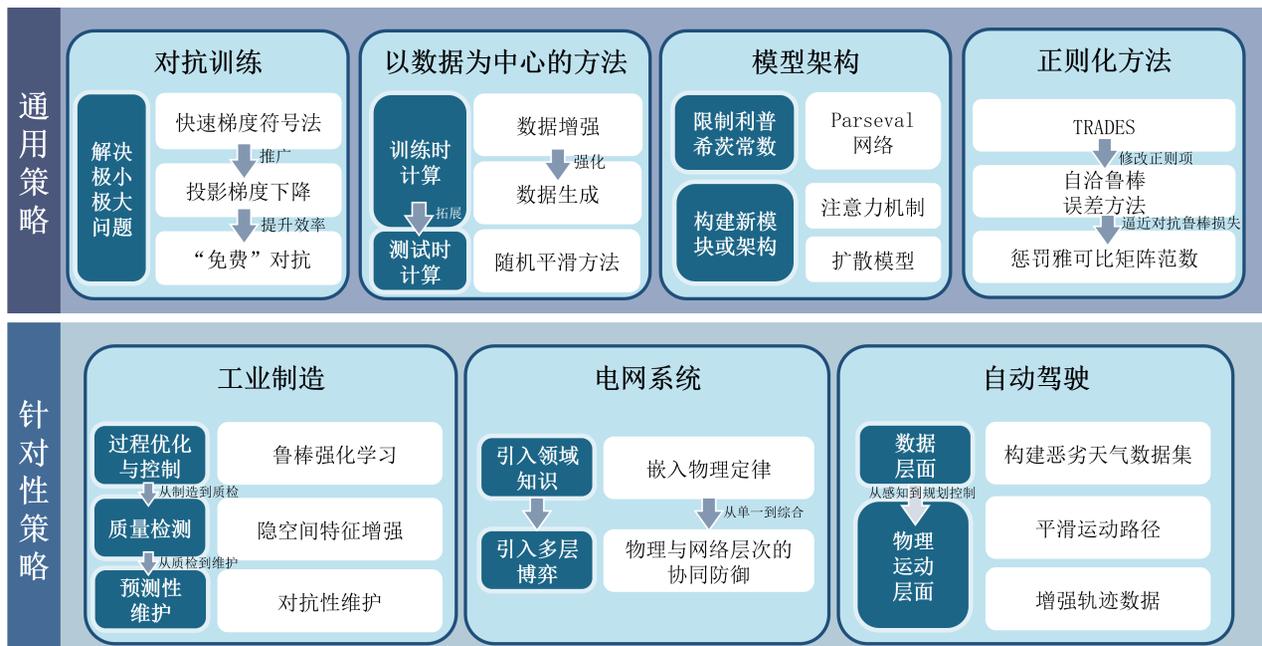


图2 提升工程系统深度神经网络鲁棒性的方法学框架图

Fig.2 Methodological Framework for Enhancing the Robustness of Deep Neural Networks in Engineering Systems

### 2.1.2 以数据为中心的方法

上一小节讨论的对抗训练不仅可看作训练算法,某种意义上也可看作是通过改变数据以提升鲁棒性。除这种依赖训练改变数据的方式外,经典的数据增强并非专门为对抗训练设计,但当与模型权重平均结合时,会显著提高对抗鲁棒性<sup>[44]</sup>。进一步的理论分析表明,最小化使用Mixup数据增强情况下的损失相当于近似地最小化对抗损失的上限<sup>[45]</sup>。除经典的数据增强方法外,通过使用扩散模型、生成对抗网络等神经网络方法生成数据,可进一步提升模型的鲁棒性<sup>[46]</sup>。与在训练时对数据进行操作不同,随机平滑方法在测试时对输入施加不同噪声并进行投票,可获得可证实鲁棒性<sup>[47,48]</sup>。

### 2.1.3 提升鲁棒性的模型架构设计方法

第1.3节表明模型架构对模型鲁棒性具有重要影响;进一步地,本小节介绍利用这种影响提升鲁棒性的方法。由于网络的Lipschitz常数与其鲁棒性直接相关,通过构建恰当的模型架构以限制Lipschitz常数,可控制对抗扰动对输出的影响,进而提升鲁棒性<sup>[49]</sup>。Parseval网络通过促使部分层的Lipschitz常数小于1,实现了对抗准确率提升<sup>[50]</sup>。在神经架构搜索的框架下,通过约束架构参数以降低网络Lipschitz常数,可显著增强网络的对抗鲁棒性<sup>[51]</sup>。

除通过限制Lipschitz常数的技术路线外,构建特定网络模块或架构亦可提升鲁棒性。(1)注意力机制常有助于鲁棒性提升。例如,挤压与激发(Squeeze-and-Excitation, SE)注意力机制、自注意力机制在某些情况下均可提升鲁棒性<sup>[33,52]</sup>。(2)神经常微分方程被认为具有较好的鲁棒性<sup>[53]</sup>。通过采用斜对称动力学保证输出稳定性,所得到的神经常微分方程在不进行对抗训练的情况下取得了较好的鲁棒准确率<sup>[54]</sup>。(3)引入扩散模型有助于提升鲁棒性。通过在图像分类模型前添加一个扩散模型,对抗噪声的影响可被弱化<sup>[55-57]</sup>。进一步地,由一个单一的预训练扩散模型构建分类器,可大幅提升模型的对抗鲁棒性<sup>[58]</sup>。不仅如此,后续理论分析表明该方法具有可证实鲁棒性<sup>[59]</sup>。

### 2.1.4 正则化方法

通过引入关于受扰前后网络输出之间距离的正则项,基于代理损失最小化的权衡启发对抗性防御(Tradeoff-inspired Adversarial Defense via Surrogate-loss Minimization, TRADES)方法在NeurIPS 2018对抗性视觉挑战赛中大幅领先于第二名<sup>[60]</sup>。进一步地,通过将TRADES中的正则项进行修改,自洽鲁棒误差方法有效缓解了鲁棒性与准确性的权衡问题<sup>[61]</sup>。通过构建正则项惩罚输入到输出的雅可比矩阵范数,可显著提升鲁棒

性<sup>[62]</sup>,其原因是该做法可逼近对抗鲁棒性损失<sup>[63]</sup>。

## 2.2 提升工程系统深度神经网络鲁棒性的针对性对策

除第2.1节所述的通用方法外,针对不同工程系统的特性,研究者提出了一些提升特定工程系统鲁棒性的特色方法。

### 2.2.1 工业制造

在工业制造系统中,深度神经网络正被越来越广泛地应用于过程优化与控制、质量检测、预测性维护等。该领域人工智能面临的特有挑战包括:生产环境中的物理扰动多样(如电磁干扰、振动)、传感器数据噪声强烈、对实时性和稳定性要求极高。确保这些深度神经网络具备鲁棒性,对于工业制造系统在系统不确定性和干扰下保持良好性能至关重要。

针对这些挑战,研究者提出了相应的策略。(1)过程优化与控制:虽然深度神经网络在过程优化与控制中正起到越来越重要的作用,但其鲁棒性缺陷为生产过程带来了隐患。Pozdnyakov等<sup>[64]</sup>提出评估并提升过程诊断模型在对抗攻击下鲁棒性的框架,并表明即使是微小的输入扰动也可能严重降低工业过程预测的准确性。在制造过程优化等任务中,原始的强化学习方法往往鲁棒性不足。鲁棒强化学习将安全性和不确定性纳入训练,通过模仿学习和大规模域随机化,可训练出在工业装配基准测试上成功率显著优于传统方法的策略<sup>[65]</sup>,其对抗扰动具有很强的容忍性。

(2)质量检测:在实际的工业生产环境中,诸如烟尘污渍、温度变化、震动摇晃等广泛存在的扰动对质量检测系统的性能造成了挑战。为应对该挑战,一种基于特征提取融合和动态标签分配的缺陷检测方法被提出,显著增强了基准模型在各类扰动下的鲁棒性<sup>[66]</sup>。在数据稀疏、缺陷尺度小的挑战性场景下,通过在隐空间进行特征增强以模拟新数据,制造业中剃须刀外壳质量检测的鲁棒性得到了提升<sup>[67]</sup>。

(3)预测性维护:在工业生产过程中,预测性维护可显著降低设备突发故障带来的危害,并降低周期性维护带来的高成本。为提升预测性维护对数据噪声的鲁棒性,Terziyan等<sup>[68]</sup>引入了对抗性维护的概念,在训练阶段引入复杂攻击,以提升系统抵御恶意行为的能力。通过结合数据插补与保序预测,实现了在不完整数据流下仍能给出有效在线剩余寿命区间,确保维护决策以可靠的置信度为基础<sup>[69]</sup>。

然而,目前提升工业制造系统中深度神经网络鲁棒性的方法仍存在差距。例如,现有的质量检测模型对未知缺陷类型和复合扰动(如光照变化与振动)的鲁棒性依然不足;预测性维护模型在面对数据完整性受损和恶

意攻击时的可靠性不佳。

### 2.2.2 电网系统

现代电网系统越来越依赖深度神经网络完成预测、控制和故障诊断等任务。电网系统作为关键基础设施,深度神经网络对其赋能的核心挑战在于对物理规律(如基尔霍夫定律)的严格遵守和对系统连锁故障的极端敏感性。深度神经网络虽然带来了效率和灵活性,但也引入了新的脆弱性:其在面临对抗性或异常输入时可能被“欺骗”,进而削弱模型在各类任务中的性能,造成能源市场混乱、基础设施故障、人身安全受威胁、停电、用电设备损坏等严重后果。为此,深度神经网络赋能电网系统的应对策略强调与物理知识的有效融合。与通用防御方法不同,电网领域物理知识的引入(如在损失函数中嵌入基尔霍夫电流定律)能使网络在面对常见扰动时依然符合电网物理约束,提升常规防御方法的效果<sup>[70]</sup>。通过引入零和多层Markovian Stackelberg博弈进行网络层防御,并将物理层防御问题转化为一个受安全约束的最优电力流优化问题,实现了物理与网络层次的全链路协同防御<sup>[71]</sup>。

尽管如此,现存差距依然显著。现有的物理约束多为简化模型,难以覆盖电网的全部动态特性。在面对未知类型扰动时,系统的鲁棒性仍是亟待解决的难题。

### 2.2.3 自动驾驶

自动驾驶系统的鲁棒性提升面临场景开放、数据动态的严峻挑战<sup>[72]</sup>。通过贴纸等物理扰动,可使自动驾驶模型在真实道路环境中出现重大错误。例如,将停止标志误识别为限速标志,攻击成功率可达100%<sup>[73]</sup>。除视觉方案外,激光雷达方案的性能在物理对抗攻击下同样会急剧下降。例如,在车辆顶部放置的三维网格对抗物体,可使基于3D点云的激光雷达方案以极高概率对目标车辆视而不见<sup>[74]</sup>。除恶意设计的对抗攻击外,在日光、车灯、手电筒等看似自然光照的特定照射方式下,交通标志识别模型也能出现很高的误判率<sup>[75]</sup>。这表明如果场景足够丰富,自然环境中的光照变化可能产生自然涌现的对抗扰动,对自动驾驶系统造成威胁。除感知外,在自动驾驶的轨迹预测等任务中也观察到了鲁棒性不足的问题<sup>[76]</sup>。

为应对自动驾驶系统鲁棒性不足的缺陷,研究者针对自动驾驶系统的特点提出了一系列解决方案。数据层面,Bijelic等<sup>[77]</sup>构建了包含雾、雨、雪等恶劣天气的多模态数据集,为自动驾驶鲁棒性提升提供了基础。物理运动层面,通过平滑运动路径及增强轨迹数据,可使受攻击时的轨迹预测误差显著降低<sup>[76]</sup>。

然而,现有防御方法大多针对特定类型的攻击或扰

动(如雨、雾),对罕见组合场景难有鲁棒性保证。此外,时序依赖性也可能导致鲁棒性在逐级传递中恶化。这些问题都指向了对更根本性解决方案的需求(见第3节)。

## 3 未来研究议题展望

如第2节所述,尽管当前针对工程系统深度神经网络鲁棒性的研究已取得显著进展,在真实、复杂的工程场景中,这些方法的有效性、泛化性等仍面临严峻挑战,制约深度神经网络在工程系统中广泛应用的鲁棒性鸿沟仍然突出。传统范式,如基于矩阵乘法和加法的网络架构、基于反向传播的对抗训练等,在对数据量和模型参数数量的需求上均已从理论上确认存在维数灾难。为此,未来研究需突破现有范式,探索专为工程系统复杂、高要求场景设计的、具备内在鲁棒性的解决方案。本节重点展望内在鲁棒的深度神经网络新架构、提升鲁棒性的新型学习范式、嵌入工程语义约束的新策略等方向。

### 3.1 内在鲁棒的深度神经网络新架构

常规模型架构虽可通过Lipschitz约束或注意力机制提升鲁棒性,但仍依赖外部防御(如对抗训练),且会牺牲表达能力或泛化性,难以实现内在鲁棒。常规模型架构对数据量和模型参数数量的维数灾难、鲁棒性和准确性的权衡,与矩阵乘法与加法的扰动累积和代数特性密不可分。正如前文所述,目前已有研究采用与乘法不同的方式构建网络基本层,得到显著强于常规模型的可认证鲁棒性,展现了这条道路的希望。然而,这些研究仍然受制于鲁棒性与准确性的权衡。未来研究可尝试探索架构层面具有内在鲁棒性的新型深度神经网络架构,基于不变性、动态系统理论、生物智能启发等的新架构,避免生成不符合物理或工程现实的特征,降低对非鲁棒特征的敏感性。

### 3.2 提升鲁棒性的新型学习范式

除模型架构外,不同的学习范式带来了不同的鲁棒性,且对鲁棒性有至关重要的影响。相比于深度学习等机器学习,采用不同学习范式的生物智能常常鲁棒性更优。这表明了通过构建不同于常规反向传播等的新型学习范式,提升深度神经网络鲁棒性的空间是存在的。例如,受生物神经机制启发,研究者突破误差反向传播算法所获得固定权重造成鲁棒性弱的局限,提出了非固定权重的神经网络电路计算范式,其具有动态联想记忆功能,可突破维数灾难带来的效率瓶颈,为构建鲁棒高效的新型深度神经网络开辟了新方向<sup>[69]</sup>。

### 3.3 嵌入工程语义约束的新策略

工程系统深度神经网络需要满足工业标准、安全规范。与数学上可清晰写出的约束不同,这些标准、规范

既有由数学表达式描述的可形式化部分,也有由自然语言描述的非形式化部分。当前已有一些将神经网络与优化中的约束联系的工作<sup>[70]</sup>。未来,如何将非形式化的语义信息转化为模型可强制遵守的硬约束,具有挑战性,也是值得研究的重要课题。

#### 4 结语

本文围绕工程系统中深度神经网络的鲁棒性挑战、根源、应对策略进行了系统性梳理,总结了提升工程系统深度神经网络鲁棒性的核心挑战,即现代深度神经网络的统计范式与工程系统对鲁棒性的严格要求之间的矛盾。我们梳理了分析深度神经网络脆弱性的理论,以及增强工程系统深度神经网络鲁棒性的通用方法和针对性方法。尽管这些方法取得了一定进展,但它们并未从根本上改变深度神经网络的固有脆弱性倾向。即便鲁棒性取得了一定的提升,但往往需要牺牲其他性能或开销庞大。未来的研究中,真正的突破可能难以来自更精巧的修补,而将源于一场彻底的范式变革。为此,我们呼吁研究力量探索全新的鲁棒深度神经网络范式。本文展望了一些可能的研究方向,旨在为未来面向工程系统的鲁棒深度神经网络探索提供参考,进而推进深度神经网络在安全攸关领域的可信赖应用。

#### 参 考 文 献

- [1] Li YZ, Ding YZ, He SY, et al. Artificial intelligence-based methods for renewable power system operation. *Nature Reviews Electrical Engineering*, 2024, 1(3): 163—179.
- [2] 姜文涛, 高原, 袁姮, 等. 门控机制的图像分类网络. *电子学报*, 2024, 52(7): 2393—2406.  
Jiang WT, Gao Y, Yuan H, et al. Image classification network of gating mechanism. *Acta Electronica Sinica*, 2024, 52(7): 2393—2406. (in Chinese)
- [3] Chen LM, Jin L, Shang MS, et al. Enhancing representation power of deep neural networks with negligible parameter growth for industrial applications. *IEEE Transactions on Systems, Man, and Cybernetics: Systems*, 2024, 54(11): 6837—6848.
- [4] Zonta T, da Costa CA, Zeiser FA, et al. A predictive maintenance model for optimizing production schedule using deep neural networks. *Journal of Manufacturing Systems*, 2022, 62: 450—462.
- [5] Bartoldson BR, Diffenderfer J, Parasyris K, et al. Adversarial robustness limits via scaling-law and human-alignment studies// *Proceedings of the 41st International Conference on Machine Learning*. Vienna, Austria: ACM, 2024: 3046—3072.
- [6] Dohmatob E, Scetbon M. Precise accuracy/robustness tradeoffs in regression; Case of general norms// *International Conference on Machine Learning*. Vienna: PMLR, 2024: 11198—11226.
- [7] Zhang RH, Sun J. Certified robust accuracy of neural networks are bounded due to Bayes errors// *Computer Aided Verification*. Cham: Springer Nature Switzerland, 2024: 352—376.
- [8] Pfrommer S. Safety, robustness, and interpretability in machine learning. Berkeley: University of California, 2025.
- [9] Pelekis S, Koutroubas T, Blika A, et al. Adversarial machine learning: A review of methods, tools, and critical industry sectors. *Artificial Intelligence Review*, 2025, 58(8): 226.
- [10] 桂韬, 奚志恒, 郑锐, 等. 基于深度学习的自然语言处理鲁棒性研究综述. *计算机学报*, 2024, 47(1): 90—112.  
Gui T, Xi ZH, Zheng R, et al. Recent researches of robustness in natural language processing based on deep neural network. *Chinese Journal of Computers*, 2024, 47(1): 90—112. (in Chinese)
- [11] Ghaffari Laleh N, Truhn D, Veldhuizen GP, et al. Adversarial attacks and adversarial robustness in computational pathology. *Nature Communications*, 2022, 13: 5711.
- [12] Zeng ZG, Wang J, Liao XX. Stability analysis of delayed cellular neural networks described using cloning templates. *IEEE Transactions on Circuits and Systems I: Regular Papers*, 2004, 51(11): 2313—2324.
- [13] Zeng ZG, Wang J, Liao XX. Global exponential stability of a general class of recurrent neural networks with time-varying delays. *IEEE Transactions on Circuits and Systems I: Fundamental Theory and Applications*, 2003, 50(10): 1353—1358.
- [14] Goodfellow IJ, Shlens J, Szegedy C. Explaining and harnessing adversarial examples. (2014-12-20)/[2025-07-29]. <https://arxiv.org/pdf/1412.6572.pdf>.
- [15] Gilmer J, Metz L, Faghri F, et al. Adversarial spheres// *International Conference on Learning Representations (ICLR), Workshop Track Proceedings*. Vancouver, BC, Canada: OpenReview.net, 2018. <https://openreview.net/forum?id=SkthILkPf>.
- [16] Schmidt L, Santurkar S, Tsipras D, et al. Adversarially robust generalization requires more data// *Proceedings of the 32nd International Conference on Neural Information Processing Systems*. Montréal, Canada: ACM, 2018: 5019—5031.
- [17] Godfrey C, Kvinge H, Bishoff E, et al. How many dimensions are required to find an adversarial example// *2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*. Vancouver, BC, Canada: IEEE, 2023: 2353—2360.
- [18] Pal A, Sulam J, Vidal R. Adversarial examples might be avoidable; The role of data concentration in adversarial robustness// *Advances in Neural Information Processing Systems*. New Orleans: Curran Associates, 2023: 46989—47015.
- [19] Conwell C, Prince JS, Kay KN, et al. A large-scale examination of inductive biases shaping high-level visual representation in brains and machines. *Nature Communications*, 2024, 15: 9383.
- [20] Feather J, Leclerc G, Mądry A, et al. Model metamers reveal divergent invariances between biological and artificial neural networks. *Nature Neuroscience*, 2023, 26(11): 2017—2034.
- [21] Ilyas A, Santurkar S, Tsipras D, et al. Adversarial examples are not bugs, they are features// *Proceedings of the 33rd Conference on Neural Information Processing Systems*. Vancouver, Canada: Curran Associates, 2019: 125—136.
- [22] Li BH, Li YZ. Adversarial training can provably improve robustness: Theoretical analysis of feature learning process under structured data// *International Conference on Learning Representations (ICLR)*. Singapore: OpenReview.net, 2025. <https://openreview.net/forum?id=in->

- LUnCpDIB.
- [23] Geirhos R, Rubisch P, Michaelis C, et al. ImageNet-trained CNNs are biased towards texture; Increasing shape bias improves accuracy and robustness. (2019-05-06)/[2025-07-29]. <https://openreview.net/forum?id=Bygh9j09KX>.
- [24] Huang HX, Wang YS, Erfani SM, et al. Exploring architectural ingredients of adversarially robust deep neural networks// Advances in Neural Information Processing Systems (NeurIPS). Virtual; Curran Associates, 2021, 34:5545—5559.
- [25] Zhu Z, Liu F, Chrysos G, et al. Robustness in deep learning: The good (width), the bad (depth), and the ugly (initialization)// Proceedings of the 35th International Conference on Neural Information Processing Systems. New Orleans, USA; Curran Associates, 2022:36094—36107.
- [26] Su D, Zhang H, Chen HG, et al. Is robustness the cost of accuracy?—A comprehensive study on the robustness of 18 deep image classification models// Computer Vision—ECCV 2018. Cham; Springer International Publishing, 2018:644—661.
- [27] Guo MH, Yang YZ, Xu R, et al. When NAS meets robustness: In search of robust architectures against adversarial attacks// 2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). Seattle, WA, USA; IEEE, 2020:628—637.
- [28] Chen LM, Jin L, Shang MS. Zero stability well predicts performance of convolutional neural networks// Proceedings of the AAAI Conference on Artificial Intelligence, 2022, 36(6):6268—6277.
- [29] Cazenavette G, Murdock C, Lucey S. Architectural adversarial robustness; The case for deep pursuit// 2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). Nashville, TN, USA; IEEE, 2021:7146—7154.
- [30] Li M, He L, Lin Z. Implicit Euler skip connections; Enhancing adversarial robustness via numerical stability// International Conference on Machine Learning. Vienna; PMLR, 2020:5874—5883.
- [31] Huang SH, Lu ZC, Deb K, et al. Revisiting residual networks for adversarial robustness// 2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). Vancouver, BC, Canada; IEEE, 2023:8202—8211.
- [32] Bombari S, Kiyani S, Mondelli M. Beyond the universal law of robustness; Sharper laws for random features and neural tangent kernels// International Conference on Machine Learning. Hawaii, USA; PMLR, 2023:2738—2776.
- [33] Benz P, Ham S, Zhang C, et al. Adversarial robustness comparison of vision transformer and MLP-Mixer to CNNs. (2021-10-06)/[2025-07-29]. <https://arxiv.org/pdf/2110.02797>.
- [34] Bai Y, Mei J, Yuille AL, et al. Are transformers more robust than CNNs// Advances in Neural Information Processing Systems. Virtual; Curran Associates, 2021, 34:26831—26843.
- [35] Wang ZY, Bai YT, Zhou YY, et al. Can CNNs be more robust than Transformers?// International Conference on Learning Representations (ICLR). Kigali, Rwanda; OpenReview.net, 2023. <https://openreview.net/forum?id=TKIFuQHHECj>.
- [36] Zhang BH, Cai TL, Lu Z, et al. Towards certifying L-infinity robustness using neural networks with L-inf-dist neurons// International Conference on Machine Learning. Virtual; PMLR, 2021:12368—12379.
- [37] Zhang BH, Jiang D, He D, et al. Rethinking Lipschitz neural networks and certified robustness: A boolean function perspective// Proceedings of the 36th International Conference on Neural Information Processing Systems. New Orleans, USA; Curran Associates, 2022, 35:19398—19413.
- [38] Tsipras D, Santurkar S, Engstrom L, et al. Robustness may be at odds with accuracy// International Conference on Learning Representations (ICLR). New Orleans, Louisiana; OpenReview.net, 2019. <https://openreview.net/forum?id=SyxAb30cY7>.
- [39] Dobriban E, Hassani H, Hong D, et al. Provable tradeoffs in adversarially robust classification. IEEE Transactions on Information Theory, 2023, 69(12):7793—7822.
- [40] Li BH, Jin JK, Zhong H, et al. Why robust generalization in deep learning is difficult: Perspective of expressive power. (2022-05-27)/[2026-03-05]. <https://arxiv.org/abs/2205.13863.pdf>.
- [41] 王璐瑶, 曹渊, 刘博涵, 等. 时间序列分类模型的集成对抗训练防御方法. 自动化学报, 2025, 51(1):144—160.  
Wang LY, Cao Y, Liu BH, et al. Ensemble adversarial training defense for time series classification models. Acta Automatica Sinica, 2025, 51(1):144—160. (in Chinese)
- [42] Madry A, Makelov A, Schmidt L, et al. Towards deep learning models resistant to adversarial attacks// International Conference on Learning Representations (ICLR). Vancouver, BC, Canada; OpenReview.net, 2018. <https://openreview.net/forum?id=rJzIBfZAb>.
- [43] Shafahi A, Najibi M, Ghiasi MA, et al. Adversarial training for free!// Advances in Neural Information Processing Systems. Vancouver; Curran Associates, 2019:3358—3369.
- [44] Rebuffi SA, Gowal S, Calian DA, et al. Data augmentation can improve robustness// Proceedings of the 35th International Conference on Neural Information Processing Systems. Virtual; Curran Associates, 2021, 34:29935—29948.
- [45] Zhang LJ, Deng Z, Kawaguchi K, et al. How does mixup help with robustness and generalization?// International Conference on Learning Representations (ICLR). Virtual; OpenReview.net, 2021. <https://openreview.net/forum?id=8yKEo06dKNo>.
- [46] Gowal S, Rebuffi SA, Wiles O, et al. Improving robustness using generated data// Proceedings of the 35th International Conference on Neural Information Processing Systems. Virtual; Curran Associates, 2021, 34:4218—4233.
- [47] Cohen JM, Rosenfeld E, Kolter JZ. Certified adversarial robustness via randomized smoothing// International Conference on Machine Learning (ICML). Long Beach, California; PMLR, 2019:1310—1320.
- [48] 纪守领, 杜天宇, 邓水光, 等. 深度学习模型鲁棒性研究综述. 计算机学报, 2022, 45(1):190—206.  
Ji SL, Du TY, Deng SG, et al. Robustness certification research on deep learning models: A survey. Chinese Journal of Computers, 2022, 45(1):190—206. (in Chinese)
- [49] Zühlke MM, Kudenko D. Adversarial robustness of neural networks from the perspective of Lipschitz Calculus: A survey. ACM Computing Surveys, 2025, 57(6):1—41.
- [50] Cisse M, Bojanowski P, Grave E, et al. Parseval networks; Improving robustness to adversarial examples// International Conference on Machine Learning (ICML). Sydney, Australia; PMLR, 2017:854—863.
- [51] Mok J, Na B, Choe H, et al. AdvRush; Searching for adversarially

- robust neural architectures// 2021 IEEE/CVF International Conference on Computer Vision (ICCV). Montreal, QC, Canada; IEEE, 2022: 12302—12312.
- [52] Xu K, Chen Z, Wang ZY, et al. Toward robust adversarial purification for face recognition under intensity-unknown attacks. *IEEE Transactions on Information Forensics and Security*, 2024, 19: 9550—9565.
- [53] Yan HS, Du JW, Tan VYF, et al. On robustness of neural ordinary differential equations// International Conference on Learning Representations (ICLR). Addis Ababa, Ethiopia; OpenReview.net, 2020. [https://openreview.net/forum?id=HJl8\\_eHYvS](https://openreview.net/forum?id=HJl8_eHYvS).
- [54] Huang YF, Yu YD, Zhang HY, et al. Adversarial robustness of stabilized NeuralODEs might be from obfuscated gradients// Mathematical and Scientific Machine Learning Conference (MSML). PMLR, 2022, 145: 497—515.
- [55] Blau T, Ganz R, Kawar B, et al. Threat model-agnostic adversarial defense using diffusion models. (2022-07-17)/[2026-03-05]. <https://arxiv.org/abs/2207.08089.pdf>.
- [56] Nie WL, Guo B, Huang YJ, et al. Diffusion models for adversarial purification// International Conference on Machine Learning. 2022.
- [57] Wang JY, Lyu ZY, Lin DH, et al. Guided diffusion model for adversarial purification. (2022-05-30)/[2026-03-05]. <https://arxiv.org/abs/2205.14969.pdf>.
- [58] Chen HR, Dong YP, Wang ZY, et al. Robust classification via a single diffusion model// Proceedings of the 41st International Conference on Machine Learning. Vienna, Austria; ACM, 2024: 6643—6665.
- [59] Chen HR, Dong YP, Hao ZK, et al. Diffusion models are certifiably robust classifiers// Advances in Neural Information Processing Systems 37. Vancouver, BC, Canada; Neural Information Processing Systems Foundation, Inc. (NeurIPS), 2024: 50062—50097.
- [60] Zhang HY, Yu YD, Jiao JT, et al. Theoretically principled trade-off between robustness and accuracy// International Conference on Machine Learning (ICML). Long Beach, California; PMLR, 2019: 7472—7482.
- [61] Pang TY, Lin M, Yang X, et al. Robustness and accuracy could be reconcilable by (proper) definition// International Conference on Machine Learning (ICML). PMLR, 2022, 162: 17258—17277.
- [62] Hoffman J, Roberts DA, Yaida S. Robust learning with Jacobian regularization. (2019-08-14)/[2025-07-29]. <https://arxiv.org/abs/1908.06729.pdf>.
- [63] Wu D, Li X. Adversarially robust generalization theory via Jacobian regularization for deep neural networks. (2021-02-23)/[2025-07-29]. <https://arxiv.org/abs/2102.12449.pdf>.
- [64] Pozdnyakov V, Kovalenko A, Makarov I, et al. AADMIP: Adversarial attacks and defenses modeling in industrial processes// International Joint Conference on Artificial Intelligence. Jeju; IJCAI, 2024: 8776—8779.
- [65] Luo JL, Sushkov O, Pevcevič R, et al. Robust multi-modal policies for industrial assembly via reinforcement learning and demonstrations: A large-scale study// Robotics; Science and Systems (RSS). Virtual, 2021. <https://roboticsproceedings.org/rss17/p088.html>.
- [66] Ye QW, Dong YW, Zhang XX, et al. Robustness defect detection: Improving the performance of surface defect detection in interference environment. *Optics and Lasers in Engineering*, 2024, 175: 108035.
- [67] Theodoropoulos S, Dardanis D, Makridis G, et al. Enhancing robustness to novel visual defects through StyleGAN latent space navigation: A manufacturing use case. *Journal of Intelligent Manufacturing*, 2025, 36(5): 3527—3541.
- [68] Terziyan V, Kaikova O. Guardians of reliability, robustness, and resilience: Adversarial maintenance in the era of industry 4.0 and 5.0. *Procedia Computer Science*, 2025, 253: 13—24.
- [69] Wang W, Wang ZQ, Cai ZQ, et al. Robust uncertainty quantification for online remaining useful life prediction with randomly missing and partially faulty sensor data. *Reliability Engineering & System Safety*, 2025, 262: 111177.
- [70] Li WT, Deka D, Wang R, et al. Physics-constrained adversarial training for neural networks in stochastic power grids. *IEEE Transactions on Artificial Intelligence*, 2024, 5(3): 1121—1131.
- [71] Zhang ZM, Huang SW, Chen Y, et al. Cyber-physical coordinated risk mitigation in smart grids based on attack-defense game. *IEEE Transactions on Power Systems*, 2022, 37(1): 530—542.
- [72] 隗寒冰, 吴化腾, 徐进. 考虑驾驶员NMS特征的自动驾驶汽车人机共驾鲁棒横向控制. *机械工程学报*, 2024, 60(16): 280—290.
- Wei HB, Wu HT, Xu J. Robust lateral shared control of autonomous vehicle considering driver NMS characteristics. *Journal of Mechanical Engineering*, 2024, 60(16): 280—290. (in Chinese)
- [73] Eykholt K, Evtimov I, Fernandes E, et al. Robust physical-world attacks on deep learning visual classification// 2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition. Salt Lake City, UT, USA; IEEE, 2018: 1625—1634.
- [74] Tu J, Ren MY, Manivasagam S, et al. Physically realizable adversarial examples for LiDAR object detection// 2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). Seattle, WA, USA; IEEE, 2020: 13713—13722.
- [75] Hsiao TF, Huang BL, Ni ZX, et al. Natural light can also be dangerous: Traffic sign misinterpretation under adversarial natural light attacks// 2024 IEEE/CVF Winter Conference on Applications of Computer Vision (WACV). Waikoloa, HI, USA; IEEE, 2024: 3903—3912.
- [76] Zhang QZ, Hu ST, Sun JC, et al. On adversarial robustness of trajectory prediction for autonomous vehicles// 2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). New Orleans, LA, USA; IEEE, 2022: 15138—15147.
- [77] Bijelic M, Gruber T, Mannan F, et al. Seeing through fog without seeing fog: Deep multimodal sensor fusion in unseen adverse weather// 2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). Seattle, WA, USA; IEEE, 2020: 11679—11689.

## Robustness of Deep Neural Networks in Engineering Systems: Frontiers and Outlook

Long Jin<sup>1\*</sup> Liangming Chen<sup>1</sup> Qiangqiang Zhang<sup>2\*</sup> Hui Li<sup>3</sup>

1. School of Information Science and Engineering, Lanzhou University, Lanzhou 730000, China

2. School of Civil Engineering and Mechanics, Lanzhou University, Lanzhou 730000, China

3. School of Civil Engineering, Harbin Institute of Technology, Harbin 150001, China

**Abstract** Deep neural networks (DNNs) are catalyzing a profound paradigm revolution in traditional engineering systems. However, a sharp contradiction exists between the inherent vulnerability of mainstream DNN models and the extreme robustness requirements of engineering systems. This “robustness gap” has become one of the key bottlenecks restricting the widespread application of artificial intelligence in engineering systems. Currently, there is no comprehensive survey specifically addressing the robustness of DNNs in engineering systems. This survey aims to systematically introduce and analyze important research on the mechanisms of DNN robustness and methods for improving it within engineering systems. First, at the problem-discovery level, this survey deconstructs the connotation and extension of the robustness gap issue. Subsequently, at the root-cause analysis level, it introduces theoretical advancements in analyzing the causes of robustness defects in DNNs and the relationship between robustness and multi-scale network architectures. Furthermore, at the level of existing countermeasures, this survey explores general methods for enhancing the robustness of DNNs in engineering systems, as well as specialized methods for improving DNN robustness tailored to the characteristics of key engineering fields such as industrial manufacturing, power grid systems, and autonomous driving. Finally, from a future outlook perspective, this survey focuses on cutting-edge directions, including novel intrinsically robust DNN architectures, new learning paradigms, and new strategies for embedding engineering semantic constraints. It seeks to provide a valuable reference for constructing robust DNNs for engineering systems.

**Keywords** artificial intelligence; deep neural networks; engineering systems; robustness

**金龙** 兰州大学教授。主要从事神经网络、机器人、智能计算等研究。入选国家级青年人才项目,主持国家自然科学基金项目4项、国家重点研发计划课题等。获甘肃省自然科学奖二等奖、中国自动化学会自然科学奖二等奖等奖励。

**张强强** 兰州大学教授。主要研究先进防护材料智能制造及多尺度计算。入选国家级青年人才项目,主持国家重点研发计划课题、国家自然科学基金项目等。获甘肃省自然科学奖一等奖、甘肃省技术发明奖一等奖、黑龙江省自然科学奖一等奖等。

(责任编辑 贾祖冰 张强)

\* Corresponding Authors, Email: jinlongsysu@foxmail.com; zhangqq@lzu.edu.cn