

Wild Networks: Exposure of 5G Network Infrastructures to Adversarial Examples

Giovanni Apruzzese, Rodion Vladimirov, Aliya Tastemirova, Pavel Laskov
Institute of Information Systems - University of Liechtenstein
 {name.surname}@uni.li

Abstract—Fifth Generation (5G) networks must support billions of heterogeneous devices while guaranteeing optimal Quality of Service (QoS). Such requirements are impossible to meet with human effort alone, and Machine Learning (ML) represents a core asset in 5G. ML, however, is known to be vulnerable to adversarial examples; moreover, as our paper will show, the 5G context is exposed to a yet another type of adversarial ML attacks that cannot be formalized with existing threat models. Proactive assessment of such risks is also challenging due to the lack of ML-powered 5G equipment available for adversarial ML research.

To tackle these problems, we propose a novel adversarial ML threat model that is particularly suited to 5G scenarios, and is agnostic to the precise function solved by ML. In contrast to existing ML threat models, our attacks do not require any compromise of the target 5G system while still being viable due to the QoS guarantees and the open nature of 5G networks. Furthermore, we propose an original framework for realistic ML security assessments based on public data. We proactively evaluate our threat model on 6 applications of ML envisioned in 5G. Our attacks affect both the training and the inference stages, can degrade the performance of state-of-the-art ML systems, and have a lower entry barrier than previous attacks.

Index Terms—Machine Learning, Network Management, 5G Networks, Adversarial Attacks, Cybersecurity

I. INTRODUCTION

THE Fifth Generation (5G) network technology standard represents a revolutionary paradigm in the context of telecommunications. On the one hand, it must provide connectivity to billions of heterogeneous devices. On the other hand, it must ensure content delivery offered by thousands of vendors and assure excellent Quality of Service (QoS) [1]. Such requirements—despite bringing undeniable benefits to the end-users of 5G—represent a serious problem for the *tenants* of the 5G Network Infrastructures (NI). The available resources (e.g., bandwidth or battery capacity) in the 5G NI are limited, and sophisticated management is required to meet the 5G standards. Hence, orchestrating the 5G NI demands timely and precise response to changing environments, and exclusive reliance on hand-crafted or hardcoded methods is impractical. To solve this problem, many works suggest [1, 2], or endorse [3], the integration of Machine Learning (ML) in the 5G NI. For instance, *network slicing*—emblematic of 5G—is focused on dynamic resource allocation and can greatly leverage the automation of ML [4]. Nonetheless, ML automation can only be appreciated with the ‘standalone’ (SA) implementation of 5G [5], whose deployment has just begun [6].¹ However, to

ensure reliable future telecommunication systems, the security of SA 5G must be put under scrutiny, *in advance* [7].

Many studies (e.g., [8, 9]) investigated conventional security aspects in 5G. In contrast, we focus on the specific threat arising from the deployment of ML: *adversarial examples* [7], which can influence the decisions of ML systems via small data manipulations. Some prior works (e.g., [10, 11, 12]) provide evidence that also ML applications envisioned in 5G are susceptible to adversarial examples. Despite their high effectiveness, such examples were always generated by attackers conforming to standard ML threat models, e.g., ‘white-box’ or ‘black-box’. Satisfying the underlying assumptions of all such attacks requires some compromise of the system hosting the target ML ‘box’—which constitutes a high entry barrier in critical infrastructures [13]. We observe that the 5G NI is exposed to a much more subtle variant of adversarial examples which cannot be formalized with prior works.

The major contribution of this paper is the proposal of the new *myopic* threat model, which *complements* existing ML threat models. Our threat model highlights that the 5G paradigm enables attacks that can be launched from the adversary’s legitimately owned devices, without compromising any segment of the 5G NI. In contrast to previous work (e.g., [14, 15, 16, 17]), our myopic attacker—constrained by the 5G context—has less knowledge and capabilities, but can still damage the 5G NI tenants by exploiting the very foundations of the 5G paradigm, such as its open nature and the QoS guarantees [18, 19]. Moreover, our threat model is *agnostic* of the specific ML deployment, allowing coverage of yet to be conceived applications of ML in 5G.

The novelty of our threat model and the nascent rollout of SA 5G demand the respective ML security evaluations that follow the *proactive* approach endorsed by Biggio and Roli [7]: “identifying relevant threats... and simulating the corresponding attacks; devising suitable countermeasures; and repeating this process *before*² system deployment”. However, *any* scientific effort on 5G ML security must face a tough challenge: no real and ML-powered 5G system is currently available for adversarial ML research, and attacking the real 5G NI is not ethical. On the other hand, defeating the ML system on in-house data is not scientifically convincing due to the risk of experimental bias. To solve this dilemma, we propose a *generic framework* for security evaluations of 5G ML components based on open-source data validated by

¹Most operational 5G NI leverage the ‘non-standalone’ (NSA) architecture, a mere an enhancement of the previous 4G NI.

²Emphasized by Biggio and Roli.

the research community. In this framework, we explicitly define the data transformations yielding realizable adversarial perturbations [20] that target ML components of the 5G NI.

Using the proposed framework, we proactively evaluate our threat model through 6 case studies, epitomizing different ML functions in the 5G NI envisioned by the state-of-the-art. We assess myopic attacks at inference and training stage, and also gauge some existing countermeasures. Our findings show that myopic attacks degrade the performance of state-of-the-art ML systems for the 5G NI. Despite being less effective than prior white-/black-box attacks, myopic attacks have a lower entry barrier and can still impact the 5G NI tenants.

Contribution and Organization. In summary, this paper makes the following contributions to the state-of-the-art:

- We propose the new *myopic* threat model. This threat model is tailored for attacks against ML systems in 5G NI, and is agnostic of the specific ML application.
- We present a framework for realistic assessment of adversarial attacks against the 5G NI based on public data.
- We use our framework to proactively evaluate the myopic threat model against 6 state-of-the-art ML prototypes for diverse 5G tasks.

The paper is structured as follows. We introduce the 5G NI, the role of ML, and motivate our paper in §II. We present our myopic threat model in §III. We describe our 5G ML security evaluation framework in §IV. We showcase our case studies in §V. We compare our paper with related work in §VI. We conclude the paper in §VII.

II. BACKGROUND AND MOTIVATION

This paper spans across three broad areas: 5G Networking, Machine Learning, and Cybersecurity. Understanding all such areas is *fundamental* to understand our contribution.

We begin by presenting an overview of the 5G ecosystem (§II-A), and the role of ML in 5G (§II-B). Then, we delve into ML security aspects (§II-C). Finally, we highlight the problems of the state-of-the-art that motivate our paper (§II-D).

A. The 5G Ecosystem

Future generations of mobile networks (5G and beyond) will by far surpass current mobile communication. They will certainly provide more capabilities and enhanced experience for human users in form of better bandwidth and lower latency. They will, however, also support billions of other actors such as self-driving cars, autonomous delivery drones, intelligent sensors, wearable medical devices and the like. All these entities must share—and compete for—the limited available resources (e.g., bandwidth, battery power, latency). This “ecosystem”, and the envisioned technical infrastructure behind it, is depicted in Fig. 1.

The 5G Network Infrastructure (NI) allows users to obtain any remote service, and can be split in two network segments [21]. The *Radio Access Network* (RAN) leverages the New Radio (NR) standard to provide a physical connection between the user equipment (UE) and the air interface devices in the base stations (gNB) whose main task is to forward

user data to the core network³. The gNB integrate specific components, such as distributed units (DU) and centralized units (CU), dedicated to preliminary network management or data preprocessing. The *core network* forwards user data to the Internet and performs more complex orchestration functions, such as the “Access and Mobility Management function” (AMF) which oversees resource optimization and management of UE (e.g., authentication, authorization and billing).

The 5G NI is **open in nature** [18, 22], and accessing its functionalities does not require excessive verifications: in most cases, a valid 5G subscription is sufficient for a UE to be included in the 5G ecosystem [23, 24]. In particular, it is common to group UE into three generic use cases [2]. The Massive Machine Type Communication (MMTC) is characterized by low data rates and very high connection density. The enhanced Mobile Broadband (eMBB) implies low connection density and high average and peak data rates. The Ultra-Reliable and Low-Latency Communication (URLLC) requires extreme reliability and low latency but assumes low data rates.

From the economic perspective, three parties interact in the 5G NI: clients, infrastructure tenants, and service providers. Clients request services from service providers. Infrastructure tenants enable service delivery by connecting clients with the service providers. The business relationships between the 5G NI tenants and the other parties are governed by the **Service Level Agreements** (SLA). Such SLA define the QoS [19] that must be met by the 5G NI, alongside the penalties for failing to meet such requirements. Compared to the past paradigms in which SLA primarily focused on the overall availability (e.g., [25]), in 5G the SLA are expected to have a finer granularity tailored to individual clients or services [26, 27].

Managing such heterogeneous ecosystem while ensuring SLA compliance is hardly feasible via static and human-defined methods [28, 29]. To increase the efficiency of 5G networks, **ML is expected to play a pivotal role in the underlying infrastructure empowering 5G** [1]. This is a significant difference with respect to previous networking paradigms: in 4G (and older) networks, ML was far from being mature and ready for real deployment.

Takeaway. Three characteristics make the (SA) 5G NI a unique setting compared to other network infrastructures: its openness, the fine-grained SLA, the role of ML. Corporate networks may use ML, and—when outsourced—may have strict SLA, but are not open. Older networks (e.g., 4G) are open, but have more relaxed SLA and do not leverage ML.

B. Machine Learning for 5G Networking

Let us briefly explain why ML represents a valuable asset for sustainable management of the 5G NI.

From the viewpoint of the 5G NI tenants, among the main goals of ML in 5G is *meeting the QoS requirements of the SLA*. Failing to reach the agreed levels of QoS violates the SLA and result in penalties for the infrastructure tenants [27, 30].

ML techniques can be deployed by the infrastructure tenants anywhere in their 5G NI (indicated as multiple *ML* boxes in

³NSA 5G also uses eNB-ng, pertaining to 4G, and hence outside our scope.

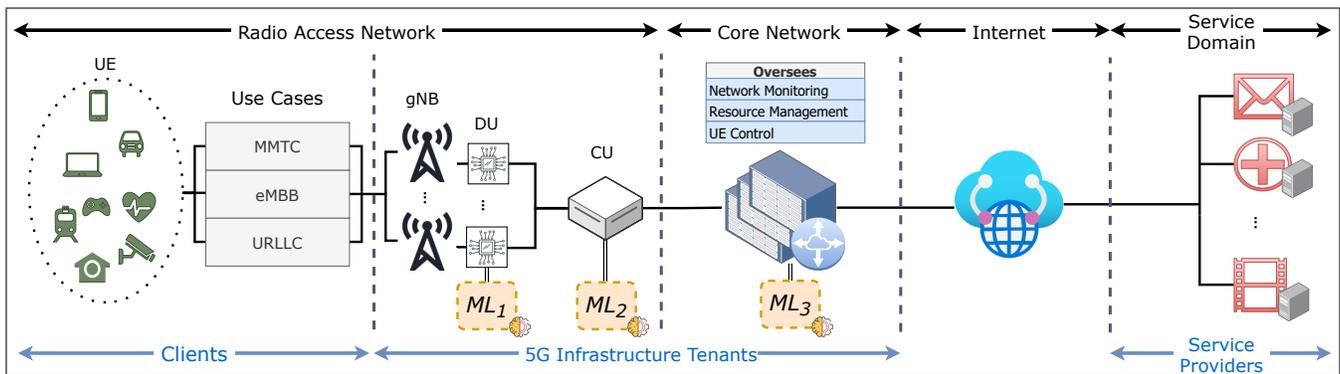


Fig. 1: The 5G Ecosystem. The *clients* transparently use the network infrastructure deployed by the 5G *tenants* to reach the *service providers*.

Fig. 1). Let us showcase four functions envisioned in 5G for which ML proved to be successful by the state-of-the-art.

Network slicing [21] aims at dividing network resources to optimize the delivery of specific groups of services. For example, entertainment video streaming requires high bandwidth but can tolerate temporary packet losses. In contrast, loss of information in eHealth services is unacceptable, but they do not require high bandwidth to operate correctly. Slices may correspond to the three 5G use-cases (eMBB, URLLC and MMTC); a ML approach for assigning UE to such slices was proposed in [31]; but other criteria are also possible [4, 32, 33], including using ML for application-based slices [34, 35]. With respect to Fig. 1, an exemplary ML-based module for network slicing may correspond to the ML_3 box.

Automatic Modulation Recognition (AMR) aims to infer modulation type of a given signal [36]. While AMR was initially proposed for military purposes, it is now an established function of 5G networking [37] and can greatly benefit from ML [38]. By identifying the correct modulation, it is possible to guarantee optimal transmission quality by tuning the transmission parameters accordingly. In Fig. 1, a ML module for AMR can be deployed in the ML_1 box.

Another function of 5G networking is *predicting the Channel Quality Indicator* (CQI). The CQI measures the communication quality received by UE: a UE computes and reports the CQI to its serving gNB; the gNB uses such CQI to ensure optimal transmissions. Communicating the CQI from the UE to the gNB has a high overhead, and frequent mutations in the environment may change the CQI before it is sent to the gNB. To avoid QoS degradation, ML can be used to predict the CQI, either via past reportings [39, 40], or by using other metrics communicated more frequently to the gNB and correlated to the CQI [41, 42]. In Fig. 1, an exemplary ML module for CQI prediction can be deployed in the ML_1 box.

Recently, the utility of deep learning has been demonstrated for *power allocation* in massive Multiple Input-Multiple Output (mMIMO) scenarios [43], a cornerstone of 5G. The idea of mMIMO is to deploy a large number of antennas at gNB, which transmit the same data to UE. The transmission capacity is measured by *spectral efficiency* (SE), which depends on the power allocated to each UE served by a gNB. Such allocation can be computed with ML by analyzing the position of the UEs in a given environment [44]. In Fig. 1, a ML module for power allocation can be placed in the ML_2 box.

These four functions will be the target of our demonstration in §V. The set of functions that exploit ML in 5G is much broader [1], and some are still to be conceived.

C. Security of Machine Learning

Security analysis require the notion of a *threat model* describing the viewpoint of the attacker w.r.t. the target system: depending on the *knowledge* and *capability*, the attacker follows a specific *strategy* to reach the intended *goal* [7].

By focusing on ML security, the so-called ‘adversarial attacks’ aim to adversely effect the target ML system via some data perturbation, i.e., *adversarial examples* [7]. Even imperceptible perturbations (e.g., a single pixel [45], or few extra bytes [46]) can compromise the decisions of ML systems.

Threat models focused on ML security hence revolve around such adversarial examples. An adversary may have a targeted or an indiscriminate *goal* [47], e.g., affecting specific examples, or degrading the overall performance of the ML system. An attacker may have variable degrees of *knowledge* on three key ML elements [48]: the ML model, M ; the training data, T ; and the feature set, F . On the other hand, the *capability* defines how the attacker can interact—possibly under some external constraints [49]—with such elements: e.g., manipulate T , affect F , or use M as an ‘oracle’ [50]. Finally, the attack *strategy* depends on the previous assumptions: e.g., introducing ‘backdoors’ in T [51], or creating a surrogate M and transferring the successful examples to the real M [52].

Classical ML attack scenarios are often expressed via the notion of a ‘box’. A *white-box* attacker has full knowledge of M and F enabling the generation of optimal perturbations [10]. In a *black-box* setting [53], the attacker has no knowledge of M , T or F , but can query M (possibly subject to some query budget [54]). In gray-box settings, the attacker may have some knowledge for optimizing the querying strategy [55]. In *no-box* settings, the attacker cannot interact with the system, but knows T [56]. A further distinction is made according to the attacker’s impact on T . If the attacker has write-access (direct or indirect) to T , then she can affect the *training* stage of M [57], otherwise she is limited to attacks at *inference* stage [58].

Despite an explosive interest to security of ML no universal countermeasures against adversarial examples have been identified so far. Some approaches (e.g., *adversarial training* [59])

or *feature removal* [60]) require foreseeing the exact form of adversarial examples. Other techniques are applicable only when the entire feature space is completely modifiable by the attacker, or to image-related data (e.g., *certified defenses* [61]). Regardless, *any* defense may degrade performance in the absence of adversarial attacks [62].

D. Motivation, Challenges and Scope

Due to its relevance in SA 5G, ML represents an attractive attack target, and adversarial examples require a dedicated treatment w.r.t. traditional security analyses [7]. However, related researches in 5G ML are immature, in particular: (i) *existing ML threat models are inadequate*; (ii) *realistic security assessments are difficult*. We aim to overcome both of these shortcomings which are now discussed in more detail.

The assumptions of classical adversarial scenarios are hard to meet in security sensitive environments. In the 5G NI context, they essentially imply a full security breach. For instance, having complete knowledge of a *trained* ML model for a white-box attack requires that an adversary gains access to a respective 5G NI component, which is well protected by conventional security mechanisms. The same holds true for black-box attacks since the output of ML models may only be visible from within the 5G NI. The training data for such critical components is likewise well protected and cannot be freely manipulated. Although meeting the assumptions of classical ML threat model is in principle possible, compromising the 5G NI is difficult and costly, and attackers may opt for different strategies. Due to the open nature of 5G, an attacker in possession of UE has legitimate access to the 5G RAN. We hence focus on how such an attacker can leverage adversarial examples to inflict damage to the 5G NI tenants by targeting their ML systems. Such offensive strategies are not covered by the state-of-the-art and require formalizing specific threat model—which must be proactively evaluated for deployments of ML in high-risk scenarios [7, 63].

Realistic assessments of adversarial examples require reproduction of the physical constraints binding the attacker [20, 49]. This is particularly challenging in the 5G context: due to the early stage of SA 5G [6], all such systems are protected by NDA⁴ and no ML-equipment is currently available for research. Furthermore, there is a lack of publicly available data for reproducing state-of-the-art ML systems [64]. Finally, practical evaluations must also consider countermeasures and their potential degradation to the baseline performance.

To overcome these shortcomings, we propose the novel *myopic* threat model (§III) which is complementary to existing ML threat models and is specifically tailored to the 5G paradigm. We also present a new framework for security assessment of 5G ML (§IV), which explains how to leverage existing public data to craft realizable adversarial perturbations against state-of-the-art ML components envisioned in 5G. Finally, we apply our framework to evaluate our threat model via 6 case studies (§V) considering different adversarial scenarios.

This work is focused on attacks against ML in SA 5G. Security issues not related to ML, or pertaining to different ML

applications, are beyond the scope of this paper. We stress that we are not the first to consider adversarial ML attacks against the 5G NI. We directly compare our work with the state-of-the-art at the end of this paper (§VI) because the differences of our efforts w.r.t. previous research can only be appreciated after thoroughly understanding our major contributions.

III. MYOPIC THREAT MODEL

Our primary contribution is the “myopic threat model”, describing adversarial ML attacks against the 5G NI which are feasible to stage due to their low cost. Indeed, the corresponding adversarial examples can be generated from the UE of an end-user (i.e. the attacker), without any access to the 5G NI. Before explaining how this is possible, let us introduce the notation for adversarial ML attacks.

Let M be a (trained) ML model that analyzes a set of features F that describe some data samples; let x be a given sample, let F_x be the feature representation of x . Assume that M can correctly predict the ground truth of F_x as $M(F_x)$. In an adversarial attack, a small perturbation ϵ causes M to predict a wrong output on the input F_x . We distinguish two scenarios: (i) *Feature-space Perturbations* (FsP), if ϵ is applied to the feature representation of x , thus resulting in $F_x + \epsilon$; or (ii) *Problem-space Perturbations* (PsP), if ϵ is applied via the process that generates x , thus resulting in $x + \epsilon$. The attack aims to finding such perturbation ϵ and can be formalized as:

$$\text{find } \epsilon \text{ s.t. } \begin{cases} M(F_x + \epsilon) \neq M(F_x) & \text{FsP} \\ M(F_{x+\epsilon}) \neq M(F_x) & \text{PsP} \end{cases} \quad (1)$$

It is implicitly assumed that $F_{x+\epsilon}$ and $F_x + \epsilon$ are associated with the same ground truth as F_x . Also note that the problem- and feature-space can overlap if $F_{x+\epsilon} = x + \epsilon$, meaning that FsP and PsP can be equivalent [49].

A. Definition of the Myopic Threat Model

To substantiate the claim that our threat model is realistic, let us interpret the abstract setting described by Eq. 1 using the general characteristics of 5G networking. This scenario is depicted in Fig. 2. The attacker operates as a client within a 5G NI and hence is part of a heterogeneous environment with many UEs. The attacker uses her UE to interact with the 5G NI (owned by its *tenants*). The results of such interactions ‘enter’ the 5G NI in the form of raw-data x , which is subject to some preprocessing. These operations yield F_x , which is passed as input to a ML model, M , that contributes to a given network function, N (e.g., network slicing). Hence, the prediction of the ML model, $M(F_x)$, is sent to N . However, besides $M(F_x)$, the network function N may use *additional input*, I , that does not depend on x , e.g., actions from other UEs, or state of the 5G NI. The output of the network function, $N(M+I)$, is not directly visible from outside the 5G NI, but all the UEs in the environment (including the attacker’s) may be affected by it. For instance, the resources allocated via network slicing depend on, and affect, the entire environment.

These assumptions describe a *myopic* attacker, whose ‘sight’ does not go beyond her UE.

⁴We interviewed multiple telcos which confirmed this fact.

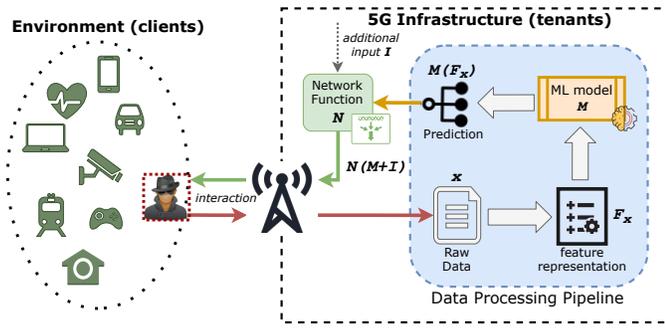


Fig. 2: Illustration of our myopic threat model.

We can now formally define the myopic threat model according to the four criteria described in §II-C:

Goal. The attacker intends to *inflict damage to the 5G NI tenants* via adversarial perturbations against the ML model M used by the network function N .

Knowledge. The attacker has *limited knowledge of the target system*. She only knows a subset of features $\mathcal{F} \subseteq \mathcal{F}$ used by M . The attacker does not know the exact implementation of either N or M , and can observe neither $M(F_x)$ nor I , and not even $N(M+I)$. Finally, the attacker has no information about the training set T used to develop M .

Capability. The attacker has *no control on the 5G NI, but has full control of her UE*. She can freely manipulate the behaviour of her UE to influence her interactions with the 5G NI and thus can consciously affect a subset of features $\bar{\mathcal{F}} \subseteq \mathcal{F}$ via PsP. The attacker cannot modify the raw-data x after it is acquired by the 5G NI and has no access to the data processing pipeline.

Strategy. The attacker *guesses* a PsP affecting some features within her knowledge and capabilities, $\bar{\mathcal{F}}$, analyzed by M .

Let us make three important remarks.

- When the PsP is translated into the *feature space*, such PsP will affect the features $\bar{\mathcal{F}} \supseteq \mathcal{F}$; i.e., the PsP may *inadvertently* affect features beyond the attacker’s knowledge.
- Our threat model is *agnostic* of the task solved by M and is applicable to a broad range of 5G use-cases.
- The PsP can be generated by manipulating the UE’s behaviour *at any layer of the protocol stack*. For instance, an attacker may use some app in a different way, modify an open-source app to affect the network layer, or manipulate the regular functionality of the UE’s operating system or its hardware (e.g., the well-known *flashing* [65]).⁵

In our demonstration (§V) we will investigate six ML case-studies, highlighting: (i) the roles⁶ played by \mathcal{F} , $\bar{\mathcal{F}}$, \mathcal{F} , $\bar{\mathcal{F}}$; and (ii) the effects of PsP on such sets.⁷

B. Comparison with Previous Threat Models

Let us elucidate how the myopic threat model is *complementary* to the conventional ML threat models; Table I summarizes

⁵Such changes, specially at lower layers, *may only possible* via modifying the respective functionality of the UE.

⁶To the best of our knowledge, we are the first to clearly distinguish these four feature sets. Such complexity is necessary to simulate all the imperceptible effects that can occur within the 5G NI.

⁷In §IV we propose an equivalent transformation to PsP that can be leveraged by research endeavours to replicate physically realizable perturbations.

this discussion. The key difference lies in the semantics of the ‘box’ considered by the respective attack scenarios.

Our ‘box’. In the myopic threat model, the ‘box’ is the entire 5G NI (i.e., the dotted square in Fig. 2). This is a complex system made of thousands of components to which the attacker has no internal access. The attacker’s UE interact only with the 5G RAN. The received feedback is the response of the 5G NI to the entire environment, and hence too complex to be usable. Furthermore, the attacker cannot directly interact with the ML component M . Due to such limited knowledge, essentially restricted to \mathcal{F} and further aggravated by the PsP constraint, the attacker is confined to a rough guessing strategy.

Conventional ‘box’. In contrast, in traditional ML threat models the ‘box’ is the ML component itself (e.g., the orange box in Fig. 2). Hence, the attacker can directly observe the impact of her actions to the predictions of the ML component (as in gray/black-box scenarios). Moreover, having complete knowledge of M (in white-box attacks), or of both F and T (for no-box attacks), enables creation of ML components that are an exact match of the real M . In all such scenarios, the attacker can leverage her knowledge and capabilities to optimize her perturbations. Even the ‘physical’ attacks introduced in [66] differ from the scenario envisioned in our threat model: they account for the separation between physical input and its data representation; however, the prediction of M corresponds to the output of the ‘box’. In contrast, our threat model also separates the ML’s predictions from the feedback received by an attacker *by two layers of indirection*: the network function $N(M+I)$, and the final feedback of the 5G NI—which is not immediately usable by our attacker due to its complexity.

TABLE I: Myopic threat model vs existing ‘box’ threat models.

	White box	Gray box	Black box	No box	Myopic
Available Knowledge	M, F	\mathcal{F}	\times	F, T	\mathcal{F}
Optimal Perturb.	✓	✓	✓	✓	\times
ML prediction $M(F_x)$	✓	✓	✓	\times	\times

Finally, we note that the myopic attacker is *less powerful*—and hence *more realistic*—than attackers typically assumed in cryptoanalysis settings (irrespective of the existence of ML). For instance, the Dolev-Yao threat model [67] assumes attackers with complete control of both ends of the communication channel (i.e., the UE and the 5G NI); in contrast, our myopic attacker only controls one end (i.e., her own UE).

C. Viability of Myopic Attacks in 5G

Let us explain why the ‘guessing’ strategy of the myopic attacker is particularly viable in 5G. To this end, we must connect our threat model (§III-A) with the unique characteristics of the 5G NI (§II-A), emphasizing the role played by ML.

We note that real attackers are not interested in crafting the ‘smallest’ perturbation that results in a successful adversarial example⁸, and are not bound to any self-imposed ‘magnitude’ constraint [69]. Indeed, real attackers operate with

⁸Crafting the ‘minimal’ perturbation is the typical assumption in adversarial ML literature focusing on computer vision [68].

a cost/benefit mindset [70]: if their goal is reached at an ‘affordable’ cost, then any strategy is viable. In the case of a myopic attacker, such *goal* is “cause damage to the 5G NI tenants by targeting ML with adversarial examples.”

Malfunctioning of ML in 5G is likely to cause QoS degradation in the environment (e.g., poor connection) or damage the 5G NI equipment (e.g., battery depletion). Since QoS is tied to SLA in the 5G ecosystem, any QoS degradation—even that of the attacker’s own UE—may be a reason for filing a complaint with the respective regulatory authority, leading to financial damages for the 5G NI tenants [71]. The open nature of the 5G NI further aggravates the problem: the attacker can legitimately introduce a large number of UEs into the environment and thus amplify the attack impact at low cost.

The combination of binding SLAs and the openness of the 5G NI makes myopic attacks both feasible and harmful⁹, especially given the lower entry barrier compared to other threats to 5G. To execute a myopic attack, no compromise of the UE or the NI components is needed.

Finally, the assumption that the attacker has full control of the UE is also well-founded and typical in security analyses. Some manipulations may require bypassing the basic security mechanisms of an UE; however, the attacker must learn how to break such mechanisms only *once* for all her UEs (assuming the same brand).¹⁰

Takeaway. Existing ML threat models can be invalidated by denying an attacker internal access to the 5G NI. In contrast, a myopic attacker can inflict damage to the 5G NI tenants by merely changing the behavior of her UE. Such novel threat requires a dedicated treatment and proactive evaluation.

IV. 5G ML SECURITY EVALUATION FRAMEWORK

To set up the stage for security evaluations of ML in 5G, we present a *high-level framework* focused on *fair* and *realistic* assessments of adversarial ML threat models¹¹.

Various previous works investigated ML methods that can empower the 5G NI (cf. §II-B). However, reproducing such methods on a in-house and closed environment—and showing the effectiveness of ad-hoc adversarial perturbations—introduces a substantial experimental bias. To mitigate such bias and enable a *fair* assessment, our framework is based on open-source data. As a main contribution, our framework ensures the evaluation of realistic adversarial ML scenarios in 5G—a challenging task, given the current state-of-the-art.

A. Suitable Public Data

Using publicly available data for security evaluations requires such data to meet four criteria. Specifically, a given public dataset X is *suitable* if: (i) it is validated by the research community, i.e., created by adopting state-of-the-art methodologies; (ii) it contains the ground truth information;

(iii) it complies with 5G; and (iv) it allows creation of realistic adversarial examples. Let us elucidate the last two criteria.

Compliance with 5G. Despite the existence of several tools to perform 5G simulations [64], only few state-of-the-art works release (e.g., [44]), or are fully built on (e.g., [73]), public data. Most prior works on 5G ML do not disclose any dataset (e.g., [32, 33]), preventing accurate replication of their solutions. Robustness of such ML systems *can be assessed on other data* provided that such data relates to the same 5G task. This may require to *adapt* a given dataset, e.g., by removing some samples or features. Such intuition can lead to discovery of ‘novel’ datasets usable for 5G ML evaluations—despite the fact that they were released when 5G was not yet defined or before its rollout began.

Realistic Perturbations. Relying on pre-collected data clearly makes it impossible to operate in the problem space. This challenge impairs the simulation of realistic adversarial examples which assume physically realizable perturbations [20]. Our solution to this challenge is to leverage perturbations applied in the *raw-data space*, which can be made semantically equivalent to PsP. This technique is presented in the following section.

B. Raw-data Space Perturbations

Let us illustrate our intuition by recalling Fig. 2. Here, we can see that the 5G NI ‘receives’ a raw-data sample, x , which is translated into its feature representation, F_x , and then forwarded to the ML model, M . Such operations are invisible to the ML model. Hence, if a perturbation ϵ is applied directly to x , then its effects carry over the entire preprocessing pipeline until $F_{x+\epsilon}$ is created and analyzed by M . Through such *Raw-data space Perturbations* (RsP) it is possible to retroactively simulate PsP *in a research environment* by using existing datasets.

To be semantically equivalent to a PsP, the RsP must consider the influence of a real attacker on the data generation process and anticipate the effects of such influence on the corresponding raw-data. Specifically, the following criteria must be met:

- The perturbation must reflect the *capabilities* assumed in the threat model. For instance, some values simply cannot be influenced by an attacker.
- After its application, the *integrity* of the perturbed raw-data must be preserved. For instance, some values depend on others, and such dependencies must be updated.
- The perturbation must abide by the *constraints* of the data generation process. For instance, the payload of a TCP packet must have between 0 and 1500 bytes.

An RsP that meets these criteria is equivalent to a PsP and hence suitable for realistic security assessments. Otherwise, such RsP either violates the assumptions of the threat model or results in a ‘trivial’ adversarial example that would be rejected by the data processing pipeline before reaching M .

Finally, depending on the specific use-cases, an RsP may have different *intensity* as long as the underlying constraints are met. For instance, a perturbation may entail “sending some extra bytes”, but the amount of such extra bytes can vary (e.g., [46]).

⁹We interviewed several telcos, which acknowledged such risk.

¹⁰All the attacks in our demonstration (§V) leverage well-known data manipulation techniques.

¹¹Our framework complements those that do not assume ML (e.g., [72]).

Takeaway. To evaluate realistic attacks via public data, the ‘mapping’ between the data generation process and the actual data contained in the dataset must be taken into account. This requires: a dataset, X , containing raw-data for the 5G NI (cf. x in Fig. 2); and anticipating a real attacker’s effects on such raw-data by means of an RsP.

C. Workflow

The proposed framework is aimed at *any* research on ML security in the 5G context. Hence it is agnostic of the specific purpose and functionality of the ML component, the format of the dataset, and of the threat model itself.

Our framework has three inputs: a dataset X meeting the criteria in §IV-A; the specifics to devise a state-of-the-art ML model M with X ; and the details to create RsP on X .

We provide a high-level schematic of our framework in Fig. 3. After acquiring a given dataset X (and adapting X to comply with 5G, if required), the first step ① is applying any sensible RsP on a pre-defined subset of X , yielding the adversarial subset of raw-data \mathcal{A} . Such RsP must meet all the criteria listed in §IV-B, and further verifications may be necessary. Then, ② both X and \mathcal{A} are subject to the *same* feature extraction operations, resulting in F_X and $F_{\mathcal{A}}$ respectively. Next, ③ F_X is split into a training T and a validation V partition¹², used to train ④ a ML model M and assess ⑤ its baseline performance. Finally, ⑥ the adversarial examples in $F_{\mathcal{A}}$ are used to ‘attack’ M . To simulate *inference* stage attacks, the performance of M is assessed on $F_{\mathcal{A}}$; otherwise, for attacks at *training* stage, M is re-trained using both T and $F_{\mathcal{A}}$, and its performance is assessed again on V .

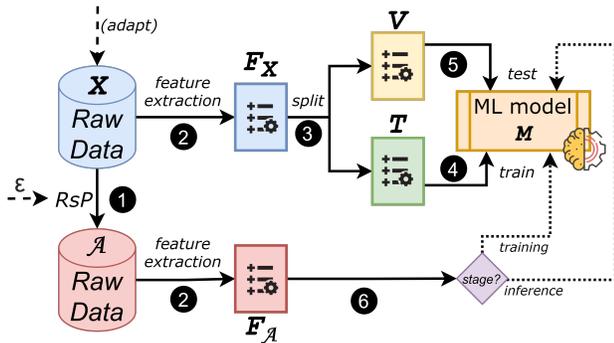


Fig. 3: Workflow of the proposed 5G ML security evaluation framework.

Proactive security evaluations must also consider *countermeasures*. As stated in §II-C, existing defenses may have limited applicability, and can induce *performance degradation* in the absence of adversarial attacks. Nobody would be interested in a defense that protects against an (uncommon) attack at the expense of an unusable ML system in (routine) operation. The workflow in §IV-C can account for the effects of a countermeasure by repeating the same steps for a ‘hardened’ version of M , which we denote as \hat{M} . We define the *tradeoff* of a countermeasure on M as:

$$\mathbb{T}(M) = P_M / P_{\hat{M}} \quad (2)$$

¹²We consider T and V to be in their feature representation.

where P_M (resp. $P_{\hat{M}}$) is the performance of M (resp. \hat{M}) on the same validation data V , according to a given performance metric P (e.g., Accuracy, F1-score). An effective defense should mitigate an attack while achieving \mathbb{T} close to 1.

V. PROACTIVE EVALUATION OF MYOPIC ATTACKS

As a proactive demonstration, we apply the proposed evaluation framework to assess myopic attacks against ML components for 5G NI. Our intention is to showcase myopic attacks in a broad array of ML deployment scenarios in SA 5G—to the extent this is possible given the current state-of-the-art. To the best of our knowledge, our paper provides the largest evaluation of adversarial attacks against ML systems envisioned in the 5G NI.

Overview. Our evaluation spans over 6 Case Studies (CS), each having a different scope. Specifically:

- CS1 (§V-A): we assume a myopic attacker controlling *multiple* UEs, and assess the impact of adversarial examples at both training and inference stages; moreover, in this CS all the feature sets ($F, \bar{F}, \mathcal{F}, \bar{\mathcal{F}}$) are disjoint.
- CS2 (§V-B): we compare the efficacy of two well-known defenses for a ML *regressor*; we also vary the *strength* of the attacker by considering different $\bar{\mathcal{F}}$.
- CS3 (§V-C): we consider attacks against an *online* ML system that leverages time-series analysis.
- CS4 (§V-D): we compare the robustness of *shallow* and *deep learning* against a myopic attacker; we also compare the effectiveness of myopic attacks against *white/black-box* attacks performed by past work against the exactly same ML system.
- CS5 (§V-E): we consider a *single-* and *multi-agent* attack scenario targeting a *physical* quality metric.
- CS6 (§V-F): we showcase a ML system that is *secure-by-design* against myopic attacks.

Each CS considers a different public dataset, used to reproduce state-of-the-art ML systems for a networking task of 5G NI. We provide an aggregated discussion in §V-G. Extensive technical details of our CS are provided in Appendix A, which also includes an overview of the chosen datasets (Table IV).

Evaluation Procedure. Each CS represents a unique setup, but the evaluation follows the same procedures, all based on the workflow of §IV-C. First, we elucidate the considered 5G system by summarizing the *data-stream* that pertains to the targeted ML component. Then, we use a public dataset to train a ML model M using a given set of features F and assess the performance of M . We then simulate myopic attacks against M . We explicitly describe the *knowledge* of the attacker by specifying \mathcal{F} , the subset of features the attacker is aware of. We then elucidate the attacker’s *capabilities* by specifying $\bar{\mathcal{F}}$, the subset of features she can consciously influence. Next, we apply some RsP to the raw-data that affects $\bar{\mathcal{F}}$ (and hence $\bar{F} \supseteq \bar{\mathcal{F}}$). Finally, we verify the integrity of the perturbed raw-data to ensure that all dependencies are preserved and the reported values are correct.

Note that a myopic attacker cannot craft an optimal adversarial example and can only guess a desired perturbation (cf. §III-C). To account for a broad range of potential perturba-

tions, we consider an array of RsP targeting each feature in $\bar{\mathcal{F}}$ at different *intensity*.¹³

In our CS, we will attack both ‘baseline’ ML systems, as well as ‘stronger’ ML systems that integrate some defensive mechanism. Because our focus is on ML security, we consider countermeasures against adversarial examples—which is the typical approach in adversarial ML literature (e.g., [55, 74]). Some CS are expanded with comparisons with white/black-box attacks; or with considerations on some protection strategies that do not belong to the ML domain.

A. CS1: Network Slicing

Highlights. The myopic attacker owns *multiple* UEs, and affects both the *inference* and *training* stages. We also assess *defensive distillation* as a countermeasure. Finally, F , \bar{F} , \mathcal{F} and $\bar{\mathcal{F}}$ are all distinct.

Target 5G system. According to the state-of-the-art (e.g., [32, 33, 35]), ML can support 5G network slicing by analyzing Network Flows (NetFlows). NetFlows capture metadata about the communication sessions in a given network. In this CS, ML is used to distinguish *active* (e.g., web-browsing) from *passive* (e.g., an automated update-check) communications. The intuition is that UE involved in active communications should be assigned to slices with higher importance.

The data-stream of this CS can be modeled as follows. The UEs of the entire environment communicate (in the form of network packets) with the 5G NI, which forwards such data to the service domain. Hence, all the raw-packets passing through the 5G NI are captured by the 5G NI and then exported to NetFlows, which are sent to a dedicated ML component that predicts whether such NetFlows belong to passive or active communications. Such predictions are then further elaborated within the 5G NI to apply the respective slicing policies. The ML component is trained on trusted NetFlows, and is periodically updated with new data to prevent performance degradation due to distribution shifts [75].

Dataset and Baseline. For this CS, we use the CTU13 [76] dataset, containing multiple traces of *real* network traffic. The data in CTU13 comes in the form of raw packet captures (PCAP), enabling the application of RsP¹⁴. The creators of CTU13 provide information allowing to distinguish active (e.g., a human user behaviour) from background (e.g., some passive or scheduled tasks) communications. To the best of our knowledge, we are the first to consider such property of CTU13 to attack corresponding ML classifiers through RsP.

To devise the ML system, we use the PCAP traces as basis from which we extract (and label) the corresponding NetFlows by following the exact procedure explained in CTU13 documentation. For each PCAP trace we obtain a set of NetFlows, used to devise a Random Forest (RF) binary classifier (we devise one classifier per PCAP trace) which analyzes

the most common NetFlow fields (we report the complete F in Table V found in Appendix A-A). Because the majority of CTU13 contains *background* traffic, we use both Accuracy (Acc) and F1-score ($F1$) as performance metrics: $Acc=0.99$ and $F1=0.81$ (the $F1$ focuses on *active* connections).

Attacks and Defense. The assumed attacker knows that the target system uses a ML component analyzing NetFlows for slice allocation. NetFlows summarize communications between two endpoints with unique IP addresses, such as the *duration* (Dur), the *ports*, or the *packets* (Pkt) and *bytes* (Byt) exchanged. The attacker infers such information, hence $\mathcal{F}=(IPs,Dur,Ports,Pkt,Byt)$. However, the attacker knows that she can only influence a subset of \mathcal{F} : the IP is assigned by the 5G NI, the Dur depends on the NetFlow appliance managed by the 5G NI, and the (low) *port* is managed by the service providers; these features are beyond attacker’s control. The attacker can only consciously influence $\bar{\mathcal{F}}=(Pkts,Byt)$ by sending more packets from her UE or adding junk payloads; doing this, however, will also affect other fields in some uncontrollable way. We assume that the attacker owns 6 UE, corresponding to $\sim 5\%$ of the (internal) hosts in CTU13. Through these actions, the attacker can affect the target ML system both at *inference* stage (e.g., an ‘active’ communication that gets assigned to a ‘background’ slice), and at *training* stage (therefore inducing poisoning attacks). The latter is due to the well-known fact that ML systems must be continuously updated to prevent concept-drift [77], meaning that some RsP may be inadvertently included in the training data used to update the ML component.

To craft RsP, we extract the raw traffic of the 6 UE owned by the attacker from the raw PCAP traces and perturb the sent packets. Doing so will induce modifications in the $Pkts$ and Byt features that fall within $\bar{\mathcal{F}}$, but also other features are affected when the raw packets are transformed into NetFlows. Specifically, $\bar{F}=(Dur,SrcByt,DstByt)$, because F distinguishes between source and destination bytes, and some packets will be included in other NetFlows; moreover, in the case of TCP connections such actions will also elicit minimal changes to the responses of the contacted host.

For *inference-stage* attacks, we submit the myopic NetFlows to the baseline RF: our objective is assessing the performance degradation of the 6 UEs owned by the attacker. For *training-stage* attacks, we inject the myopic NetFlows into the training set by randomly replacing the original NetFlows of the 6 myopic UEs with their myopic variants; we stress that such procedure does not violate our assumptions because the myopic UEs are trusted by the 5G NI (otherwise, such UEs would not receive any connectivity). We *do not* manipulate any sample generated by a non-myopic UE (which correspond to 95% of CTU13). We then re-train the RF on such ‘poisoned’ dataset and re-evaluate it on the original testing partition: our objective, here, is assessing the impact on the entire environment. The process is repeated for increasing replacement ratios (25%, 50%, 75%, 90%) of myopic NetFlows, which correspond to just 1–5% of the training data.

To investigate a countermeasure against our attacks, we apply the variant of the *defensive distillation* technique suitable for RF, as reported in [78]. The tradeoff \mathbb{T} of this countermea-

¹³Due to the novelty of our threat model, we will consider defenses that are ‘generally’ applicable. For instance, we will not consider *certified defenses* because they are tailored for adversarial attacks on images and because they assume attackers that are bound by a given perturbation magnitude—both of which are assumptions that do not pertain to our threat model.

¹⁴The CTU13 also contains the processed NetFlows, which we do not consider because it is not raw-data, and hence not valid for RsP.

sure in the absence of attacks, as measured by the Accuracy, is 0.97, showing a slight increase in the baseline performance. We provide more technical details in Appendix A-A, including the detailed workflow for the RsP.

Results. We report the results of our attacks in Fig. 4.

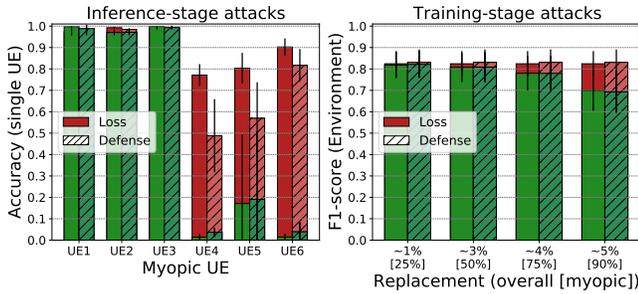


Fig. 4: CS1: Myopic attacks at inference (left) and training (right) stage. The metrics are shown for the baseline *RF* (solid bars) and its hardened counterpart (bars with oblique lines). In both graphs, the red part shows the performance degradation induced by the attack; the black line shows the standard deviation over the multiple considered PCAP traces. Inference stage attacks focus on the effects on each myopic UE (denoted by each pair of bars). Training stage attacks focus on the entire Environment, and the performance is reported by varying replacements of NetFlows **only** of the myopic UE (corresponding at most to 5% of the overall traffic of the Environment).

Attacks at inference stage have notably different effects on the myopic UEs. Some UEs are only slightly affected (UE1–UE3), but others are strongly affected (UE4–UE6). This interesting phenomenon showcases the limited capability of the myopic attacker, who cannot craft optimal and successful perturbations: some UE may be a ‘useless’ vector of RsP, but others can be very powerful—a result which can be amplified by introducing more UEs in the environment. While these attacks only influence the predictions for the attacker’s UEs, we observe that such predictions are used as basis for resource allocation, and hence the entire network can be adversely affected by these attacks.

Attacks at training stage start to impact the whole network after at least 50% of the attacker’s UEs NetFlows are replaced with their myopic variants. Recall, however, the attacker’s UEs represent at most 5% of the entire network population, hence such impact is not negligible.

Defensive distillation is not effective against these attacks; however, it slightly improves the baseline accuracy ($\mathbb{T}=0.97$).

Attacking the considered system via white-/black-box attacks is hardly feasible. For example, the attacker must be aware of the *exact* NetFlow exporter (different tools yield different data [79]), which is confidential information belonging to the 5G NI tenants. At the same time, the attacker cannot leverage the feedback of M because its predictions (passive or active NetFlows) are further elaborated by the 5G NI before applying the slicing policies.

B. CS2: CQI Prediction

Highlights. We compare the effectiveness of two popular countermeasures, *adversarial training* and *feature removal*, applied to a ML *regressor*. We also consider myopic attackers having different knowledge and capabilities, by varying \mathcal{F} .

Target 5G system. This CS considers the 5G task of CQI prediction. The ML components must infer the CQI based on measurements computed by the gNB or other measurements reported more often (once every ms) by the UE [41, 42], such as those related to the Radio Resource Control (RRC) protocol.

The data-stream envisioned in this CS assumes UE that communicate channel quality metrics to the gNB, which integrates a ML component that analyzes such metrics and estimates the CQI. The data received by the gNB is considered to be trusted, because implementing security mechanisms would increase the overhead and therefore introduce delays that would defeat the entire purpose of using ML. The ML component is trained on data collected by the 5G NI tenants and conforming to diverse types of UE [41].

Dataset and Baseline. We use the *ElasticMon* dataset [80] to reproduce the state-of-the-art approach in [41]. Released in 2020, *ElasticMon* contains 5G synthetic raw-data denoting the periodic reportings of a UE to its gNB. The creators of *ElasticMon* considered all the characteristics of 5G in their network environment. To the best of our knowledge, we are the first to consider this dataset in adversarial scenarios.

Our baseline ML model is a *RF* regressor, which we develop by following the exact instructions of [41]. We use the raw-data in *ElasticMon* and follow the same preprocessing steps, obtaining their same feature set F (reported in Table VI in Appendix A-B). Our baseline *RF* achieves similar performance as in [41], as measured via Root Mean Squared Error ($RMSE=0.22$) and accuracy ($Acc=0.95$).

Attacks and Defense. An attacker aware that the gNB predicts the CQI on the basis of RRC reportings can expect that the ML system analyzes the Resource Signal Reference Power (RSRP); or, alternatively, the transmitted packets or bytes. Hence, $\mathcal{F}=(RSRP,Pkt,Byt)$. The *Byt* or *Pkt* can be easily influenced (as explained in §V-A). The RSRP is computed directly by the UE, and the myopic attacker (who physically owns her UE) is fully able to control the reported RSRP (e.g., [81]). Indeed, the RSRP measures the strength of a signal received by a UE from all the surrounding gNBs, with the assumption that the UE will connect to the gNB with the best RSRP. However, an attacker can force her UE to connect to a gNB with a suboptimal RSRP, meaning that the 5G NI will receive an RSRP with different value.¹⁵

Based on \mathcal{F} , we consider four attack scenarios: $\overline{\mathcal{F}}_1=(RSRP)$; $\overline{\mathcal{F}}_2=(Byt)$; $\overline{\mathcal{F}}_3=(Pkt)$; $\overline{\mathcal{F}}_4=(Pkt,Byt)$. We simulate these scenarios as follows. For $\overline{\mathcal{F}}_1$, the RSRP is replaced with a randomly chosen RSRP value in *ElasticMon*. For $\overline{\mathcal{F}}_{2-4}$ we generate RsP at 7 increasing intensity levels, where the *Pkt* (and/or *Byt*) are incrementally increased as a function of their standard deviation across the *ElasticMon* dataset. We always verify the integrity of the myopic raw-data, from which we obtain their feature representation. All attacks occur at inference-stage.

We consider two well-known countermeasures, *adversarial training* and *feature removal*, requiring to foresee the patterns of adversarial attacks. We apply both countermeasures by assuming correct anticipation of each attack scenario (as

¹⁵To the best of our knowledge, there is no way for the 5G NI to prevent such occurrence, as the procedure is carried out on the UE.

in [59,60]). Additional details such as \bar{F} as well as the application of the RsP and the defenses are in Appendix A-B.

Results. We report the results of our attacks in Fig. 5; we only show adversarial training in Fig. 5 because applying feature removal *always defused* the attack (no degradation).

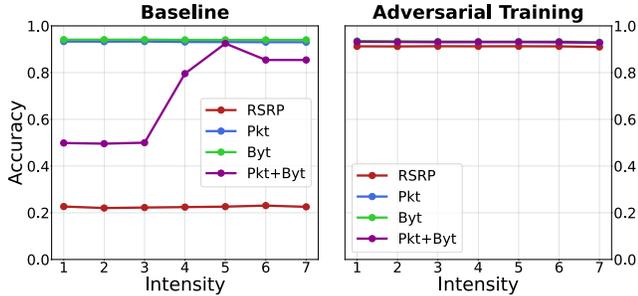


Fig. 5: CS2: results for the baseline RF (left), and the \widehat{RF} hardened through adversarial training (right); feature removal always nullified the attack. The y-axis indicates the accuracy and the x-axis the intensity of each attack (represented with a colored line); for $\bar{\mathcal{F}}_1$ the results are reported for 7 different random RSRP replacements.

By focusing on the baseline RF (left graph), we observe that RsP to $\bar{\mathcal{F}}_1$ (red line) lead to significant performance degradation, whereas RsP to $\bar{\mathcal{F}}_{2,3}$ (blue and green lines) have no impact; it is almost counterintuitive that RsP to $\bar{\mathcal{F}}_4$ (purple line) have a higher impact for lower intensities. Adversarial training (right graph) protects against all our attacks at little cost (its tradeoff shown in Table II is always negligible)—a success shared also by feature removal.

TABLE II: CS2 Tradeoff (lower is better, $\mathbb{T}=1$ is no change).

Defense	$\bar{\mathcal{F}}_1$	$\bar{\mathcal{F}}_2$	$\bar{\mathcal{F}}_3$	$\bar{\mathcal{F}}_4$
Adv. Tra.	1.01	1.01	1.01	1.01
Fea. Rem.	1.00	1.01	1.01	1.01

Launching adversarial attacks conforming to well-known threat models against the envisioned ML system requires many resources. For instance, obtaining detailed information on M may be prohibitive, because such M is embedded in the gNB, and hence ‘difficult’ to reach—unless the attacker already compromised the 5G NI in some way.

C. CS3: CQI Prediction (online)

Highlight. We showcase myopic attacks against *online* ML using real 5G network traffic data.

Target 5G system. This CS also focuses on CQI prediction (as in CS2), but here the prediction is made by analyzing historical CQI reportings via time-series analyses [39]. The data-stream is similar to the one described in CS2, the only difference being the information communicated by the UE to the gNB (as well as the frequency of such communications).

Dataset and Baseline. We are inspired by the recent work in [39] which uses the Irish 5G dataset [82]. To the best of our knowledge we are the first to evaluate such dataset in adversarial scenarios. The Irish 5G contains 5G raw-data metrics (including the CQI) collected by the gNB of a major 5G mobile operator and describing 20 minutes of reportings sampled every second. It contains many traces, focused on a

different UE mobility pattern, ‘static’ or ‘driving’. We remove those traces for which the CQI is not provided. The remaining traces refer to activities of the UE, such as ‘streaming’ or ‘download’. For the ‘static’ mobility pattern, we consider ‘download’ because the CQI of ‘streaming’ never changes; for the ‘driving’ mobility pattern, we consider ‘streaming’ because for the ‘download’ activities the behaviour was too irregular and we never obtained appreciable performance.

As done in [39], our baseline is a Long Short Term Memory ($LSTM$) regressor. Each trace has a dedicated $LSTM$, which consider only the CQI and corresponding timestamp (we do not apply any preprocessing). We use the first half of the trace to pre-train the $LSTM$. Then, the $LSTM$ becomes operational and predicts each next CQI value by using the last 30 reportings. The $LSTM$ is updated in an online fashion when a new sample is received. Hence $F=(\text{last } 30 \text{ CQI})$.

Attacks. A myopic attacker can expect that the 5G NI uses online ML to predict the CQI by using the past history, sampled every second. The attacker cannot know the exact length of such history, hence $\mathcal{F}=(\text{some past CQI})$. However, she can be certain that the most recent value is included in such history, hence: $\bar{\mathcal{F}}=(\text{previous CQI})$. The CQI is computed (and reported) by the UE, so the myopic attacker can influence it arbitrarily [81]. Such myopic attacker can affect the ML predictions if the CQI reported by her UE to the gNB is different from the actual one. All subsequent predictions made by considering the history with the myopic CQI will be affected when the ML model will use the ‘myopic’ history to update itself. We consider two attack scenarios when applying the RsP. In the first scenario, the UE reports that CQI=0, which is the lowest possible value. In the second scenario, the CQI is spoofed with a value within ± 3 of the actual one. In both scenarios, the fake CQI is sent once every minute, therefore $\bar{F}=(\text{at most one CQI})$ in both cases. An attacker can theoretically send a fake CQI continuously, but such attempts can hardly be considered as adversarial examples. Because the $LSTM$ are operational for 10 minutes, our RsP will only affect $\sim 1\%$ of the overall sampled data in both scenarios.

Results. The $LSTM$ accumulate errors over time, which can be measured via Cumulative Root Mean Squared Error (CRMSE). We show in Fig. 6 the difference of such CRMSE for the two considered attack scenarios with respect to clean data (denoted as *differential CRMSE*), during the 10 minutes of the $LSTM$ operation. We use a full (dotted) line to denote the first (second) attack scenario, whereas red (blue) lines refer to the ‘driving’ (‘static’) mobility pattern.

We can see that the first attack (full lines) induces a significantly higher differential CRMSE. In comparison, the second attack (dotted lines) is less effective. We find it interesting that, in both scenarios, the malicious CQI occur only 10 times—representing $\sim 1\%$ of the overall sampled data. Even such a small fraction of malicious samples has a major effect, because each sample influences many future predictions. We report in Appendix A-C the complete time-series of this CS.

Protecting similar ML systems against such myopic attacks is not trivial. The target is an $LSTM$ regressor, and its online nature make it difficult to identify adversarial ML countermeasures that do not impair baseline performance.

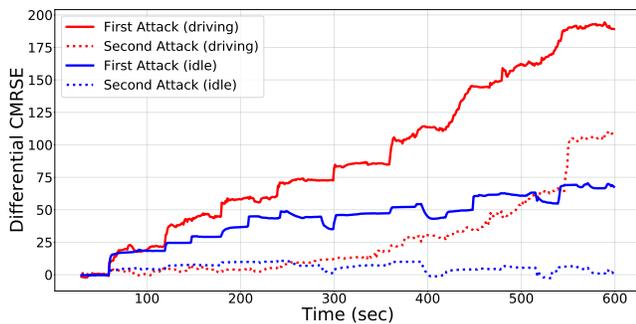


Fig. 6: CS3: differential CRMSE of myopic w.r.t. regular behaviors.

D. CS4: Automatic Modulation Recognition

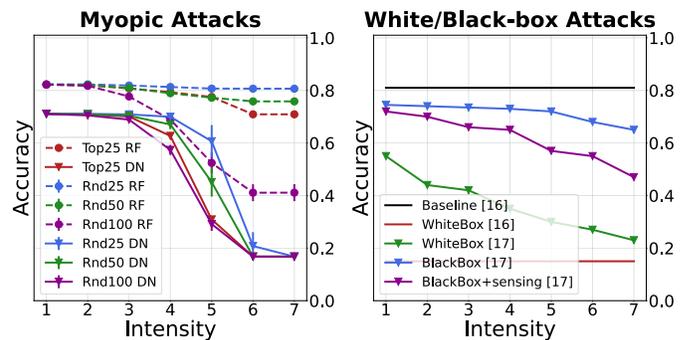
Highlights. The impact of myopic attacks on *deep* and *shallow* learning is assessed. We also consider *white/black-box* attacks targeting exactly the same ML system [16, 17].

Target 5G system. The ML component focuses on AMR, a task extensively studied by research in 5G (cf. §II-B). Hence, the data-stream assumes gNB that capture the radio-signals belonging to a given UE, and measure the In-Phase and Quadrature components of such signals. Afterwards, such measurements are sent to a dedicated ML model which must infer the modulation of the corresponding signal. The data used to develop the ML model belongs to the 5G NI tenants, and can be generated either via controlled simulations, or by direct acquisition of physical radio-signals from real UEs.

Dataset and Baselines. We rely on the `RML2016.10a` dataset [36]. This well-known dataset includes 220K signals collected for signals at 10 different Signal-to-Noise Ratio (SNR), denoting 11 modulation schemes (i.e., the labels). It is created via the GNU-Radio toolkit, which is appreciated in 5G simulations [64]. The dataset was released in 2016, and *to make it compliant with 5G* we consider only digital modulations. Each signals is described by a vector of 128 pairs of In-Phase/Quadrature (I/Q) measurements.

As it was done in previous work (e.g. [10, 16]), we do not apply any preprocessing to the data, hence the features correspond exactly to the raw physical measurements provided in `RML2016.10a`, thus $F=(256 \text{ I/Q measurements})$. To compare shallow with deep learning, we train two multi-class baselines: one is a ‘shallow’ *RF*, the other is exactly the same Deep Neural Network (*DN*) as in [17]. We compute accuracy of both baselines: $Acc_{RF}=0.82$, $Acc_{DN}=0.72$.

Attacks. A *myopic* attacker can easily expect that ML systems for AMR analyze I/Q measurements, but cannot reasonably know the exact composition of F , hence $\mathcal{F}=\bar{\mathcal{F}}=\bar{F}$ =(some I/Q measurements). Such attacker can artificially generate some noise [10] and affect ML at inference stage. However, the I/Q measurements are computed at the receiving end (e.g., gNB), which is not accessible. The myopic attacker is hence limited to random and imprecise perturbations, which we simulate by considering four different scenarios. Three involve RsP of randomly chosen measurements: either 25, 50 or 100 (10, 20 or 40% of F , respectively); whereas in the fourth—a worst-case—the RsP affect exactly the 25 most significant measurements for classification. Hence: $\bar{\mathcal{F}}_{1-3}=(25/50/100 \text{ rnd I/Q})$ and $\bar{\mathcal{F}}_4=(25 \text{ top I/Q})$. For each



(a) Deep vs Shallow learning.

(b) Attacks in [16, 17].

Fig. 7: CS4: Myopic Attacks vs White/Black-box Attacks.

scenario, we craft RsP at 7 increasing intensity levels and attack the baselines. Due to the randomness of $\bar{\mathcal{F}}_{1-3}$, we repeat them 20 times. Additional details are in Appendix A-D.

Let us describe the adversaries considered by related work, all targeting a DN using `RML2016.10a`. In [16] a *white-box* attacker has complete knowledge and can freely apply any perturbation to the input data. The setting in [17] is more constrained: here, both a *white-box* and a *black-box* attackers want to affect a specific victim (e.g., the owner of a different UE). Hence, they can only apply a perturbation synchronized with the victim’s communication channel; such channel can be approximated by sensing the spectrum. All these attackers (in both [17] and [16]) know the entire feature set, and their perturbations affect *all* measurements of each signal; by using our notation, for such attackers $\mathcal{F}=\bar{\mathcal{F}}=\bar{F}=F$. The attacks in [16] apply a fixed perturbation, whereas those in [17] also consider perturbations at 7 increasing intensities.

Results. We analyse the accuracy degradation attained by myopic attacks against our *DN* and *RF* (Fig. 7a, the standard deviation of our 20 trials is denoted as a vertical bar on each marker), and compare it with that of attacks in the related work (Fig. 7b). These results focus on signals with SNR=10db, which is common in the respective literature.

Shallow learning appears to be more robust than deep learning: in Fig. 7a the *DN* (full lines) is more affected by the myopic attacks than the *RF* (dotted lines), although both baselines are defeated with random RsP at high intensity. Prior work only targeted deep learning, so a fair comparison against such work must focus on the *DN* results in Fig. 7a. We can see that the attack in [16] (red line in Fig. 7b) is devastating but it also assumes an extremely powerful attacker. Conversely, the constrained attacker in [17] is much less successful in the black-box setting (blue line in Fig. 7b); she becomes more successful if she can obtain more information about the victim, either by sensing the spectrum (purple line), or in the white-box setting (green line).

E. CS5: Power Allocation in massive MIMO

Highlights. We assess myopic examples in a *single-* and *multi-agent* adversarial setting; the attacks target a *physical* quality metric in 5G networks, the spectral efficiency (SE). We also compare our attacks with a recent work [83]. Finally,

an intriguing property of this CS is that the attacker’s goal is reached through *unsuccessful* adversarial examples.

Target 5G system. The ML component must estimate the power allocation in mMIMO 5G networks. According to the state-of-the-art, the allocated power can be computed by ML as a function of the position of multiple UE with respect to their serving gNB [44]. Hence, we consider a system in which the data-stream envisions UEs that, after obtaining their geographical location [84], communicates such information to the gNB. The 5G NI uses such information to determine the actual *distance*¹⁶ of all the UEs with respect to their serving gNB. All such distances are then provided to a dedicated ML component deployed within the 5G NI, which estimates how much power should be allocated to each gNB in order to support the attached UEs. Such estimates are finally used by the 5G NI for proper resource management. The excellent results of [44] show that proficient ML models for power allocation can be trained through entirely simulated data.

Dataset and Baseline. We use the recent PA-mMIMO dataset [44] generated by a simulation of 20 UE served by 4 gNBs. Its 335K raw-data samples describe the positions of each UE (x/y coordinate pairs) and the corresponding actual allocated power, making it suitable for 5G experiments [44, 83]. We replicate the state-of-the-art technique proposed in [44]: UE report their locations (acquired via GPS) to the 5G NI which uses deep learning to allocate the power of 4 gNBs to 20 UE. We use the source code provided by [44] to devise our Deep Network (DN) baseline whose feature set is $F=(20\ x/y\ \text{pairs})$. The quality of each predicted power allocation vector is estimated my means of the physical SE metric.

Attacks. A myopic attacker can expect the usage of location-based information for power allocation. She cannot reasonably know *how many* UEs are served by the system, but she knows that at least her UE is included among them. In this case, the attacker can spoof her geographical position (e.g., [86, 87]), affecting the DN at the inference stage; hence, $\bar{F}=\bar{F}=\mathcal{F}=(1\ x/y\ \text{pair})$. An attacker may also collude with a partner in a multi-agent setting by synchronizing their attacks, increasing their impact on the whole system; hence, $\bar{F}=\bar{F}=\mathcal{F}=(2\ x/y\ \text{pairs})$. Such scenarios are depicted in Fig. 8, showing an mMIMO network with 4 cells and 4 gNBs (filled dots) each serving 5 UEs (empty dots), for a total of 20 UEs. For the single-agent setting, we consider an attacker whose UE is served by gNB1; for the multi-agent setting, we consider an additional attacker whose UE is served by gNB3. Myopic attackers can fake their location, but they cannot determine the exact values that maximizes their impact. They may, however, expect that greater (faked) distances from the gNB cause more power to be allocated to their UE. We simulate such behavior by “moving” the attackers’ UE away from the gNB. This is done via RsP that alter the x/y coordinates of the myopic UE(s), so that the distance from the serving gNB increases in 8 steps in a range of [0-300] within the same cell (shown with a dotted trajectory in Fig. 8).

Results. We show the impact of the attack on the SE in Fig. 9. Because the attacking UEs never truly move the correct

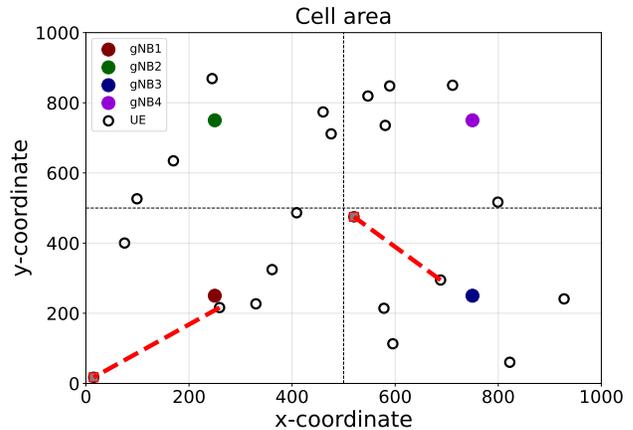


Fig. 8: CS5: network topology and adversarial movement.

power allocation should not change and the corresponding SE should remain the same (i.e., at 100%). In the single-agent setting (solid lines), we can see from the left plot that the SE of *all* UEs in the gNB1 cell is affected: the attacker (UE5) gains about 3% SE at the expense of decreasing the SE of other UEs by 3%–20%. As a consequence, the 5G NI allocates more power to serve the attacker’s UE5 and less for the other UEs, preventing optimal QoS. This setting also slightly affects the UEs in the neighboring cells, as shown in the right plot (with green, blue, purple solid lines). In the multi-agent setting (dotted lines), the impact of an attack is more complex. The *average* SE of a cell, shown in the right plot in dotted lines, can decrease as well as slightly increase as a consequence of an attack. What happens is that the attackers’ UEs (in gNB1 and gNB3) gain more SE than the other UEs in their cells. This leads to an overall degradation of SE in the whole network (confirmed by gNB2 and gNB4 which have only benign UEs).

This CS showcases another intriguing property of the myopic threat model: the attacker can damage the 5G NI even if all adversarial examples are unsuccessful, from the ML point of view. Indeed, the generated perturbations elicit the *correct* response from the DN (as opposed to the wrong response expected from adversarial examples). However, the system QoS is still damaged because the power allocated to serve each UE in the network is incorrect. This happens because the DN approximates power allocation as a multidimensional function depending on the positions of *all* UEs in the network. Due to the novelty of this finding it is not immediately clear how such attacks could be prevented at the ML level.

We now compare our myopic attacks against a very recent work that targets exactly the same DN with white/black-box attacks. In [83], the attacker is very powerful: she controls and can freely “move” *all* 20 UEs of the system. Furthermore, the capability to query the DN (black-box) by manipulating the entire environment, or the full knowledge (white-box) can be exploited to find the optimal attack vector causing an incorrect prediction. Specifically, by moving some UEs *within centimeters*, the SE can be decreased up to 60% and 20% in white- and black-box settings, respectively. We found it surprising that our ‘primitive’ myopic attacker controlling one UEs could cause a similar decrease in SE (e.g., UE3 in the

¹⁶Such distance can be measured by hardcoding the geographical position of the gNB, or by using a dedicated localization service [85].

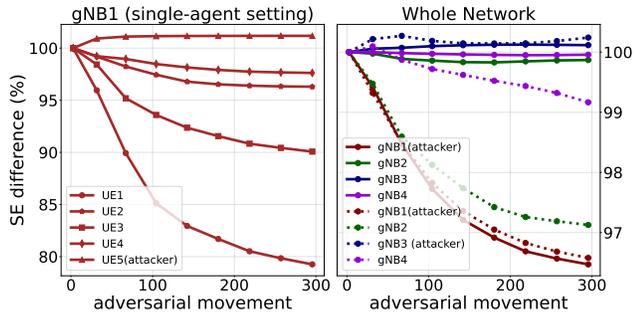


Fig. 9: CS5: attack impact on a single cell and whole network, measured as the difference (in %) of the SE as a function of the adversarial movement. The left plot shows the effects of the single-agent setting on the specific cell served by the closest gNB1; the right plot shows the effects of both the single- and multi-agent setting (full and dotted lines, respectively) on the whole system, where the values are averaged among all UEs in each cell.

left plot of Fig. 9) as a black-box attacker having access to the 5G NI and controlling all 20 UEs.

F. CS6: A secure by design ML system against myopic attacks

Highlights. We showcase a state-of-the-art ML system that *cannot* be affected by a myopic attacker.

Target 5G system. The ML component focuses on network slicing but—differently from CS1 (§V-A)—the objective is to determine the slice assignment by using KPIs *contained in the 5G NI* [31]. Therefore, this CS entails a data-stream in which UEs, served by a given gNB, leverage the connectivity of the 5G NI. However, the decision to allocate more (or less) resources for such gNB (and corresponding UEs) depends on KPI related to the current state of the 5G NI (e.g., the time of the day, or the packet delay budget). Such slicing policies can be implemented by training a deep learning classifier on data provided by the 5G NI tenants, specifying which ‘slice’ should be chosen according to the status of the 5G NI [31].

Dataset and Baseline. We use the DeepSlice dataset [31], released in 2019 by following and characteristics of the three 5G use cases: eMBB, MMTTC, URLLC (which are the three possible ‘slices’). It was also used in [73] to propose a defensive mechanism against DoS attacks. We replicate the same multi-class *DN* as in [31] using the same feature set F (reported in Table VII in Appendix A-E). Our *DN* matches the perfect Accuracy of [31].

Attacks. A myopic attacker that is aware of [31] may want to target a similar ML system to disrupt resource allocation. Such attacker can infer the entire feature set of [31], hence $\mathcal{F}=F$. However, the only features in \mathcal{F} that a myopic attacker can influence are those related to time (e.g., *hour* and *day*): the other features (e.g., *PktDelayBudget*) depend on the 5G specifications and can only be modified by the 5G NI tenants, and are not controllable by the myopic attacker. Hence, $\overline{\mathcal{F}}=\overline{F}=(hour, day)$. We apply the RsP by considering any combination of the *day* and *hour*.

Results. Not a single myopic example was successful against the baseline *DN*. This is an example of a “secure-

by-design” system against myopic attacks¹⁷, because the features relevant for classification (e.g., *PktDelayBudget*) can be influenced only with direct access to the 5G NI. This can be possible for an ‘insider’, but such assumption would violate our threat model. Regardless, for those interested in such ‘insider myopic attacks’, we showcase these attempts in Appendix A-E, which also includes additional details as well as the evaluation of ML-specific countermeasures.

G. Discussion

Let us finalize our evaluation by discussing some crucial remarks from a practical point of view.

Are myopic attacks undefeatable? The attacks in CS1, CS5 and CS3 are successful and ML countermeasures are difficult to apply or not very effective. However, as shown by CS6, some ML systems are secure-by-design against some forms of myopic attacks. Moreover, CS2 clearly shows that some well-known adversarial ML countermeasures can mitigate or completely defuse our myopic attacks. However, to prevent the attacks in CS2, it is necessary to *predict* either the features attacked (for feature removal) or the specific RsP (for adversarial training): doing this requires a proactive approach, which is the one endorsed by our paper. Finally, we hope that our paper will inspire future work aimed at *specifically* countering myopic attacks (with a favorable tradeoff).

Are myopic attacks more dangerous than existing attacks against ML? The white/black-box attacks of [83] and those in [16, 17] (cf. CS5 and CS4) are more disruptive, but they also require a higher resource investment to be staged; for instance, our myopic attacks only require a rough knowledge of the feature set, which can be obtained by reading technical reports or scientific papers. It is hence crucial that both circumstances are taken into account when testing ML components deployed in critical infrastructures.

Can myopic attacks target different network infrastructures? Yes, but—to be viable—only if such infrastructures (i) use ML, (ii) are open, and (iii) are subject to granular SLA (cf. §III-C). If a network is not open, then an attacker cannot easily leverage multiple UEs to stage myopic attacks; without granular SLA, the damage potential of myopic attacks is decreased; and without ML it is simply impossible to conceive attacks based on adversarial examples.

Do the CS represent feasible adversarial attacks? Yes. Our CS are based on the predictions of ML components to adversarial examples, and are adversarial attacks by definition (cf. Eq. 1). Such examples are created via RsP manipulations that are possible by attackers who physically own their UE(s). To provide a broad assessment our RsP have varying intensity, because real attackers are not interested in crafting the ‘minimal’ perturbation. RsP with low intensity may have little effect, but higher intensities can be disruptive (e.g., Fig. 9).

Can myopic attacks be blocked via non-ML protection mechanisms? Yes. For instance, the attacks in CS5 would be

¹⁷We stress that the notion of “secure-by-design” should be contextualized. In our case, we use such notion only to refer to systems that cannot be affected by myopic attacks (which are, by definition, launched from outside the 5G NI). Of course, the ML system in CS6 can still be attacked via other forms of adversarial attacks—if their requirements are met by the adversary.

defused by ensuring the correctness of the position reported by UE. However, such mechanisms are not cheap to implement¹⁸, can be broken (e.g., [89]), and may induce overheads that would nullify the advantages provided by ML in the 5G NI.

Is the real 5G NI endangered by myopic attacks? Yes. In our CS, we attack ML prototypes based on state-of-the-art techniques for 5G networking. The deployed 5G NI will behave differently, since its ML components are trained on different (private) datasets, and the processing pipelines may include additional (proprietary) mechanisms—all of which are at the discretion of the 5G NI tenants and/or under NDA and hence not available for research. However, we interviewed several 5G telcos who acknowledged that our CS are indeed feasible and represent a threat that must be taken into account.

Are our findings generalizable? We acknowledge that, with the exception of CS4, we only attack a single ML model per case study—hence, we do not claim that “every ML model deployed in practice is vulnerable to myopic attacks”. However, all our CS revolve around ML methods proposed by the state-of-the-art, and hence represent those that are more likely to be deployed in the real 5G NI (because they provide the highest performance). Nonetheless, we hope that our paper will inspire future works that will consider different ML models: perhaps, some ML techniques will provide a slightly lower baseline performance, but are naturally robust against myopic perturbations. We believe that such an outcome would be *valuable* for real ML deployments.

VI. RELATED WORK

Limited attention has been given to adversarial examples in 5G networking. Some papers provide a broad overview of the security risks of ML-powered wireless communications [4, 90, 91]; others propose ML-based security mechanisms in 5G (e.g., [73, 92]). An exhaustive survey is [93]. None of these works evaluate or propose original threat models.

Some papers have little in common with 5G networking. The authors of [94, 95, 96] consider attacks against cyber detectors. Similarly, [97, 98, 99, 100] target ML systems for computer vision and autonomous driving. All these tasks—despite being strongly linked with 5G [101]—are unrelated to networking functions. Finally, some orthogonal studies consider the usage of ML as an *offensive* mechanism [75, 102, 103].

Let us directly compare our paper with closely related work, summarized in Table III. For each paper, we report the threat model (○ and ● denote white- and black-box attacks); whether the perturbations adhere to some constraints; the assessment of defenses; and the usage of public data.

Some papers consider white-box attackers with complete knowledge of the target ML model that can freely apply any perturbation (without providing any justification) [11, 16, 83]. Other works consider more realistic scenarios where the perturbations are subject to physical constraints (e.g. [10, 17, 109, 110, 111]), despite considering attackers with full knowledge of the ML system or that can observe its output.

To be viable, all feedback-based strategies pose an additional requirement: the attacker must be certain that the

TABLE III: Prior works on adversarial ML in 5G networking.

Paper	Year	Attacker	Constr. Perturb.	Defenses	Public Data
[104]	2018	○ / ●	✓	✗	1
[12]	2018	●	✓	✗	✗
[15]	2019	○ / ●	✓	✗	✗
[105]	2019	○ / ●	✓	✓	2
[106]	2019	○	✓	✗	1
[107]	2019	○	✓	✓	1
[108]	2019	●	✓	✓	✗
[14]	2019	●	✓	✓	✗
[75]	2019	●	✓	✓	✗
[11]	2019	○	✗	✓	✗
[10]	2019	○	✓	✗	1
[109]	2020	○ / ●	✓	✗	1
[110]	2020	○ / ●	✓	✗	1
[111]	2020	○ / ●	✓	✓	1
[16]	2021	○ / ●	✗	✗	2
[112]	2021	○ / ●	✗	✗	✗
[83]	2021	○ / ●	✗	✗	1
Ours		myopic	✓	✓	6

feedback corresponds exactly to the ML output. By using the notation in §III-A, these attackers must be aware that $N(M+I)=M(F_x)$, and that such $N(M+I)$ corresponds to the feedback of the 5G NI. Obtaining such certainty is tough without access to the infrastructure hosting the target ML systems. Another possibility is if the attacker owned *all* the UEs served by the 5G NI, allowing to monitor the results of all their interactions—which is clearly unfeasible.

The authors of [14] apply specific manipulations to the training set (i.e., trojanning); having such direct access to training data is unlikely in critical environments. More viable poisoning attempts are found in [12, 108] where the attacker has less control and must first ‘sense’ the spectrum. A similar and viable strategy is the jamming attack in [15]: the attacker perpetually inspects the communication channel to replicate (and then disrupt) a ML model. This procedure can only work against ML systems trained on the exact data distribution captured during the sensing activity of the attacker.

All the adversarial attacks against the 5G NI considered by past work are agreeably not impossible; however, launching them *in real 5G contexts* requires information obtainable only through direct access to the 5G NI. This can happen either as a result of an insider threat [113], or via a prolonged and resource intensive APT campaign [114]. We do not claim that our threat model is the only way to realistically attack the 5G NI via adversarial examples.

With respect to existing work, we point out the exposure of the 5G NI to a more affordable but still dangerous adversarial attack strategy, which we formalize with our novel threat model. Moreover, no previous paper has addressed the importance of conducting realistic and generic assessments of adversarial attacks against 5G NI, and only few consider countermeasures. In addition, previous papers usually focus on a single dataset (most notably, RML2016.10a) which sometimes is not publicly released, preventing reproducibility; in contrast, we consider 6 case studies all based on open data, representing the largest assessment of adversarial attacks against the 5G NI. Finally, our original evaluation framework will hopefully facilitate the assessment of adversarial attacks against 5G ML systems—at least until such systems become available for research purposes.

¹⁸They may require stateful analyses of diverse data-sources (e.g., [88]).

VII. CONCLUSIONS

The security of the 5G Network Infrastructure (NI) is of paramount importance due to its critical role in the current and future society. Empowering such infrastructure with ML exposes it to the risk of adversarial examples, which have not received adequate treatment in this context. The 5G paradigm enables a new class of harmful adversarial ML attacks with a low entry barrier, which *cannot* be formalized with existing adversarial ML threat models. Furthermore, such vulnerabilities must be *proactively* assessed, but the early stage of ML in SA 5G makes such evaluations challenging for research.

In this paper, we propose the first threat model that is specific to adversarial ML attacks against the 5G NI. Our ‘myopic’ threat model describes viable attacks that can inflict damage (physical as well as monetary) to the 5G NI tenants. Moreover, we provide an original security evaluation framework based on open source data, enabling realistic assessment of adversarial examples against state-of-the-art ML systems. Both our threat model as well as our framework are *agnostic* of the specific function solved by ML in the 5G NI, and hence cover even yet to be conceived applications of ML in SA 5G.

We apply our framework to evaluate the proposed myopic threat model, and analyze six case studies where we target state-of-the-art ML systems for 5G NI. We show that 5 out of 6 systems can be broken with our myopic attacks which can influence both the training and inference stages; the attacks can also simultaneously affect single devices and the entire environment and can be amplified by multi-agent strategies or acquisition of more UE. All of these circumstances can cause damage to the 5G NI tenants due to SLA violations. Our attacks may have a smaller success rate than prior black-/white-box attacks but do not require any compromise of the 5G NI. Finally, we showcase a ML system which is immune to our attacks—by design.

This paper can inspire many research directions, such as case studies using the real ML elements in SA 5G—when they become available for adversarial ML research purposes. The proposed framework also allows the discovery of existing datasets that are usable for realistic adversarial ML evaluations. Another possibility is estimating the monetary damage of myopic attacks as a consequence of SLA violations.

We have formalized and assessed a new class of attacks against ML systems in SA 5G. Yet, at this point in time, foreseeing the possible incarnations of such systems and evaluating their security risks is difficult—but necessary to ensure their reliability for our society. We hope that our contribution will serve as an important step towards secure ML for 5G networking.

REFERENCES

- [1] M. E. Morocho-Cayamcela, H. Lee, and W. Lim, “Machine learning for 5G/B5G mobile and wireless communications: Potential, limitations, and future directions,” *IEEE Access*, vol. 7, pp. 137 184–137 206, 2019.
- [2] 5G-PPP WG, “5G Architecture,” Tech. Rep. 3.0, June 2019.
- [3] ITU, “Framework for data handling to enable machine learning in future networks IMT-2020,” Tech. Rep. 16, February 2020. [Online]. Available: <https://www.itu.int/en/ITU-T/focusgroups/ml5g/Pages/default.aspx>
- [4] V. P. Kafle, Y. Fukushima, P. Martinez-Julia, and T. Miyazawa, “Consideration on automation of 5G network slicing with machine learning,” in *Proc. IEEE/ITU Kaleidoscope: ML for 5G Future*, 2018, pp. 1–8.
- [5] GSMA, “5G implementation guidelines: SA option 2,” Tech. Rep., 2020.
- [6] “Verizon 5G standalone core,” 2020. [Online]. Available: <https://www.verizon.com/about/news/verizon-5g-standalone-core>
- [7] B. Biggio and F. Roli, “Wild patterns: Ten years after the rise of adversarial machine learning,” *Elsevier Pattern Recogn.*, vol. 84, pp. 317–331, 2018.
- [8] A. Dutta and E. Hammad, “5g security challenges and opportunities: A system approach,” in *Proc. IEEE 5G World Forum*, 2020, pp. 109–114.
- [9] J. Suomalainen, J. Julku, M. Vehkaperä, and H. Posti, “Securing public safety communications on commercial and tactical 5g networks: A survey and future research directions,” *IEEE Journal Commun. Soc.*, 2021.
- [10] B. Flowers, R. M. Buehrer, and W. C. Headley, “Evaluating adversarial evasion attacks in the context of wireless communications,” *IEEE Trans. Inf. Forensics Secur.*, vol. 15, pp. 1102–1113, 2019.
- [11] M. Usama, J. Qadir, M. A. Imran *et al.*, “Adversarial ML attack on self organizing cellular networks,” in *Proc. IEEE UK/China Emerging Tech.*, 2019, pp. 1–5.
- [12] Y. Shi, T. Erpek, Y. E. Sagduyu, and J. H. Li, “Spectrum data poisoning with adversarial deep learning,” in *Proc. IEEE Milit. Commun. Conf.*, 2018, pp. 407–412.
- [13] G. Apruzzese, M. Andreolini, L. Ferretti, M. Marchetti, and M. Colajanni, “Modeling realistic adversarial attacks against network intrusion detection systems,” *ACM Digital Threats: Research and Practice*, 2021.
- [14] K. Davaslioglu and Y. E. Sagduyu, “Trojan attacks on wireless signal classification with adversarial machine learning,” in *Proc. IEEE Int. Symp. Dyn. Spectrum Access Netw.*, 2019, pp. 1–6.
- [15] M. Sadeghi and E. G. Larsson, “Physical adversarial attacks against end-to-end autoencoder communication systems,” *IEEE Commun. Letters*, vol. 23, no. 5, pp. 847–850, 2019.
- [16] M. Usama, I. Ilahi, J. Qadir, R. N. Mitra, and M. K. Marina, “Examining Machine Learning for 5G and beyond through an adversarial lens,” *IEEE Internet Comp.*, vol. 25, no. 2, pp. 26–34, 2021.
- [17] B. Kim, Y. E. Sagduyu, K. Davaslioglu, T. Erpek, and S. Ulukus, “Over-the-Air adversarial attacks on Deep Learning-based Modulation Classifier over Wireless Channels,” in *Proc. IEEE Conf. Inf. Sci. Sys.*, 2020, pp. 1–6.
- [18] S. Komeylian and S. Komeylian, “Deploying an OFDM physical layer security with high rate data for 5G wireless networks,” in *Proc. IEEE Canadian Conf. Elec. Comp. Eng.*, 2020, pp. 1–7.
- [19] H. N. Qureshi, M. Manalastas, S. M. A. Zaidi, A. Imran, and M. O. Al Kalaa, “Service Level Agreements for 5G and Beyond: Overview, Challenges and Enablers of 5G-Healthcare Systems,” *IEEE Access*, 2020.
- [20] L. Tong, B. Li, C. Hajaj, C. Xiao, N. Zhang, and Y. Vorobeychik, “Improving robustness of {ML} classifiers against realizable evasion attacks using conserved features,” in *Proc. USENIX Secur. Symp.*, 2019, pp. 285–302.
- [21] 3G-PP TSG, “Radio access network; NR; NG-RAN description,” Tech. Rep. 16, February 2021.
- [22] Fortinet, “Leveraging a Holistic 5G Security Strategy as Our Digital World Continues to Evolve,” 2021.
- [23] L. Gavrilovska, V. Rakovic, and D. Denkovski, “From Cloud RAN to Open RAN,” *Wirel. Pers. Commun.*, vol. 113, no. 3, pp. 1523–1539, 2020.
- [24] N. Kazemifard and V. Shah-Mansouri, “Minimum delay function placement and resource allocation for Open RAN (O-RAN) 5G networks,” *Elsevier Comp. Netw.*, 2021.
- [25] L. Mastroeni and M. Naldi, “Violation of service availability targets in service level agreements,” in *Proc. IEEE Conf. Comp. Sci. Inf. Sys.*, 2011, pp. 537–540.
- [26] “Improve 5G access performance and differentiate end-to-end SLAs,” Accedian, Tech. Rep., 2021.
- [27] A. Papageorgiou, A. Fernández-Fernández, L. Ochoa-Aday, M. S. Peláez, and M. S. Siddiqui, “SLA Management Procedures in 5G Slicing-based Systems,” in *Proc. IEEE Europ. Conf. Netw. Comm.*, 2020, pp. 7–11.
- [28] T. K. Rodrigues, K. Suto, H. Nishiyama, J. Liu, and N. Kato, “Machine learning meets computation and communication control in evolving edge and cloud: Challenges and future perspective,” *IEEE Commun. Surv. Tut.*, vol. 22, no. 1, pp. 38–67, 2019.

- [29] C. Benzaid and T. Taleb, "Ai-driven zero touch network and service management in 5g and beyond: Challenges and research directions," *IEEE Network*, vol. 34, no. 2, pp. 186–194, 2020.
- [30] TowerXchange, "TowerXchange Asia Dossier 2019," Tech. Rep., 2019.
- [31] A. Thantharate, R. Paropkari, V. Walunj, and C. Beard, "Deepslice: A deep learning approach towards an efficient and reliable network slicing in 5g networks," in *Proc. IEEE Ubiquitous Comput., Elect. & Mobile Commun. Conf.*, 2019, pp. 762–767.
- [32] R. Li, Z. Zhao, Q. Sun, I. Chih-Lin, C. Yang, X. Chen, M. Zhao, and H. Zhang, "Deep reinforcement learning for resource management in network slicing," *IEEE Access*, vol. 6, pp. 74429–74441, 2018.
- [33] E. Coronado and R. Riggio, "Flow-Based Network Slicing: Mapping the Future Mobile Radio Access Networks," in *Proc. IEEE Int. Conf. Comp. Commun. Netw.*, 2019, pp. 1–9.
- [34] L.-V. Le, B.-S. P. Lin, L.-P. Tung, and D. Sinh, "SDN/NFV, machine learning, and big data driven network slicing for 5G," in *IEEE 5G World Forum*, 2018, pp. 20–25.
- [35] L.-V. Le, B.-S. Lin, and S. Do, "Applying big data, machine learning, and SDN/NFV for 5G early-stage traffic classification and network QoS control," *T. Netw. Commun.*, vol. 6, no. 2, p. 36, 2018.
- [36] T. J. O'shea and N. West, "Radio machine learning dataset generation with gnu radio," in *Proc. GNU Radio Conf.*, 2016.
- [37] A. P. Hermawan, R. R. Ginanjar, D.-S. Kim, and J.-M. Lee, "CNN-based automatic modulation classification for beyond 5G communications," *IEEE Commun. Lett.*, vol. 24, no. 5, pp. 1038–1041, 2020.
- [38] M. A. Hazar, N. Odabasioglu, T. Ensari, Y. Kavurucu, and O. Sayan, "Performance analysis and improvement of machine learning algorithms for automatic modulation recognition over rayleigh fading channels," *Neural Comp. Appl.*, vol. 29, no. 9, pp. 351–360, 2018.
- [39] B. Y. L. Kimura, J. Almeida *et al.*, "Deep learning in beyond 5g networks with image-based time-series representation," *arXiv:2104.08584*, 2021.
- [40] C. Parera, A. E. Redondi, M. Cesana, Q. Liao, and I. Malanchini, "Transfer learning for channel quality prediction," in *Proc. IEEE Int. Symp. Measur. Netw.*, 2019, pp. 1–6.
- [41] X. Vasilakos, N. Nikaein, D. H. Lorenz, B. Koksall, and N. Ferdosian, "Integrated methodology to cognitive network & slice management in virtualized 5g networks," *arXiv:2005.04830*, 2020.
- [42] R. Ul Mustafa, S. Ferlin, C. Esteve Rothenberg, D. Raca, and J. J. Quinlan, "A Supervised Machine Learning Approach for DASH Video QoE Prediction in 5G Networks," in *Proc. ACM Symp. QoS Secur. Wireless Mobile Netw.*, 2020, pp. 57–64.
- [43] T. Van Chien, T. N. Canh, E. Björnson, and E. G. Larsson, "Power control in cellular massive mimo with varying user activity: A deep learning solution," *IEEE T. Wireless Commun.*, vol. 19, no. 9, pp. 5732–5748, 2020.
- [44] L. Sanguinetti, A. Zappone, and M. Debbah, "Deep learning power allocation in massive mimo," in *Proc. IEEE Asilomar Conf. Sig. Syst. Comp.*, 2018, pp. 1257–1261.
- [45] J. Su, D. V. Vargas, and K. Sakurai, "One pixel attack for fooling deep neural networks," *IEEE T. Evol. Comput.*, 2019.
- [46] G. Apruzzese, M. Colajanni, and M. Marchetti, "Evaluating the effectiveness of adversarial attacks against botnet detectors," in *Proc. IEEE Int. Symp. Netw. Comput. Appl.*, Oct. 2019, pp. 1–8.
- [47] M. Barreno, B. Nelson, A. D. Joseph, and J. D. Tygar, "The security of machine learning," *Machine Learning*, vol. 81, no. 2, pp. 121–148, 2010.
- [48] N. Šrđić and P. Laskov, "Practical evasion of a learning-based classifier: A case study," in *Proceedings of the 2014 IEEE Symposium on Security and Privacy*, 2014, pp. 197–211.
- [49] F. Pierazzi, F. Pendlebury, J. Cortellazzi, and L. Cavallaro, "Intriguing properties of adversarial ml attacks in the problem space," in *IEEE Symp. Secur. Privacy*, 2020.
- [50] N. Papernot, P. McDaniel, I. Goodfellow, S. Jha, Z. B. Celik, and A. Swami, "Practical black-box attacks against machine learning," in *Proc. ACM Conf. Comput. Commun. Secur.*, 2017, pp. 506–519.
- [51] E. Quiring and K. Rieck, "Backdooring and poisoning neural networks with image-scaling attacks," in *IEEE SP Workshops*, 2020.
- [52] A. Demontis, M. Melis, M. Pintor, M. Jagielski, B. Biggio, A. Oprea, C. Nita-Rotaru, and F. Roli, "Why do adversarial attacks transfer? Explaining transferability of evasion and poisoning attacks," in *Proc. USENIX Secur. Symp.*, 2019, pp. 321–338.
- [53] H. Dang, Y. Huang, and E.-C. Chang, "Evading classifiers by morphing in the dark," in *Proc. ACM SIGSAC Conf. Comput. Commun. Secur.*, 2017, pp. 119–133.
- [54] A. Ilyas, L. Engstrom, A. Athalye, and J. Lin, "Black-box adversarial attacks with limited queries and information," in *ICML*, 2018, pp. 2137–2146.
- [55] G. Apruzzese, M. Andreolini, M. Marchetti, A. Venturi, and M. Colajanni, "Deep reinforcement adversarial learning against botnet evasion attacks," *IEEE T. Netw. Serv. Manag.*, vol. 17, no. 4, pp. 1975–1987, 2020.
- [56] P.-Y. Chen, H. Zhang, Y. Sharma, J. Yi, and C.-J. Hsieh, "Zoo: Zeroth order optimization based black-box attacks to deep neural networks without training substitute models," in *Proc. ACM Workshop Artif. Intel. Secur.*, 2017, pp. 15–26.
- [57] B. Biggio, B. Nelson, and P. Laskov, "Poisoning attacks against support vector machines," in *Proc. Int. Conf. Machin. Learning*. Omnipress, 2012, pp. 1467–1474.
- [58] B. Biggio, I. Corona, D. Maiorca, B. Nelson, N. Šrđić, P. Laskov, G. Giacinto, and F. Roli, "Evasion attacks against machine learning at test time," in *Proc. Europ. Conf. Mach. Learn. and Knowl. Discov. Databases*, 2013, pp. 387–402.
- [59] F. Tramèr, A. Kurakin, N. Papernot, I. Goodfellow, D. Boneh, and P. McDaniel, "Ensemble adversarial training: Attacks and defenses," in *Proc. Int. Conf. Learning Representations*, 2018.
- [60] C. Smutz and A. Stavrou, "Malicious PDF detection using metadata and structural features," in *Proc. Ann. Comp. Secur. Appl. Conf.*, 2012, pp. 239–248.
- [61] M. Lecuyer, V. Atlidakis, R. Geambasu, D. Hsu, and S. Jana, "Certified robustness to adversarial examples with differential privacy," in *Proc. IEEE Symp. Secur. Privacy*, 2019, pp. 656–672.
- [62] W. He, J. Wei, X. Chen, N. Carlini, and D. Song, "Adversarial example defense: Ensembles of weak defenses are not strong," in *USENIX Workshop Offensive Techn.*, 2017.
- [63] "Laying down harmonised rules on artificial intelligence," European Commission, Tech. Rep., 2021.
- [64] L. Bonati, M. Polese, S. D'Oro, S. Basagni, and T. Melodia, "Open, programmable, and virtualized 5G networks: State-of-the-art and the road ahead," *Elsevier Comp. Netw.*, vol. 182, p. 107516, 2020.
- [65] Z. Zhang, Y. Wang, J. Jing, Q. Wang, and L. Lei, "Once root always a threat: Analyzing the security threats of android permission system," in *Australasian Conference on Information Security and Privacy*, 2014, pp. 354–369.
- [66] A. Kurakin, I. Goodfellow, and S. Bengio, "Adversarial examples in the physical world," in *ICLR Workshop*, 2017.
- [67] J. Zhang, Q. Wang, L. Yang, and T. Feng, "Formal verification of 5G-EAP-TLS authentication protocol," in *Proc. IEEE Int. Conf. Data Sci. Cyberspace*, 2019, pp. 503–509.
- [68] N. Papernot, P. McDaniel, A. Sinha, and M. Wellman, "Sok: Security and privacy in machine learning," in *Proc. IEEE Europ. Symp. Secur. Privacy*, Apr. 2018, pp. 399–414.
- [69] N. Carlini, A. Athalye, N. Papernot, W. Brendel, J. Rauber, D. Tsipras, I. J. Goodfellow, A. Madry, and A. Kurakin, "On evaluating adversarial robustness," *Computing Research Repository*, vol. abs/1902.06705, 2019.
- [70] K. S. Wilson and M. A. Kiy, "Some fundamental cybersecurity concepts," *IEEE Access*, 2014.
- [71] V. Sciancalepore, K. Samdanis, X. Costa-Perez, D. Bega, M. Gramaglia, and A. Banchs, "Mobile traffic forecasting for maximizing 5g network slicing resource utilization," in *Proc. IEEE Conf. Comp. Commun.*, 2017, pp. 1–9.
- [72] S. Hussain, O. Chowdhury, S. Mehnaz, and E. Bertino, "LTEInspector: A systematic approach for adversarial testing of 4G LTE," in *Proc. NDSS*, 2018.
- [73] A. Thantharate, R. Paropkari, V. Walunj, C. Beard, and P. Kankariya, "Secure5G: A deep learning framework towards a secure network slicing in 5G and beyond," in *Proc. IEEE Comput. Commun. Workshop Conf.*, 2020, pp. 852–857.
- [74] N. Carlini, "Poisoning the Unlabeled Dataset of Semi-Supervised Learning," in *USENIX Secur. Symp.*, 2021.
- [75] Y. Sagduyu, Y. Shi, and T. Erpek, "Adversarial deep learning for over-the-air spectrum poisoning attacks," *IEEE Trans. Mobil. Comput.*, 2019.
- [76] S. Garcia, M. Grill, J. Stiborek, and A. Zunino, "An empirical comparison of botnet detection methods," *Elsevier Comput. Secur.*, vol. 45, pp. 100–123, 2014.
- [77] A. Imran, A. Zoha, and A. Abu-Dayya, "Challenges in 5G: how to empower SON with big data for enabling 5G," *IEEE Network*, vol. 28, no. 6, pp. 27–33, 2014.
- [78] G. Apruzzese, M. Andreolini, M. Colajanni, and M. Marchetti, "Hardening random forest cyber detectors against adversarial attacks," *IEEE T. Emerging Topics Comp. Int.*, vol. 4, no. 4, pp. 427–439, 2020.

- [79] G. Vormayr, J. Fabini, and T. Zseby, "Why are my flows different? a tutorial on flow exporters," *IEEE Communications Surveys & Tutorials*, vol. 22, no. 3, pp. 2064–2103, 2020.
- [80] X. Vasilakos, B. Köksal, D. H. Izaldi, N. Nikaein, R. Schmidt, N. Ferdosian, R. F. Sari, and R.-G. Cheng, "ElasticSDK: A monitoring software development kit for enabling data-driven management and control in 5G," in *Proc. IEEE Netw. Op. Manag. Symp.*, 2020, pp. 1–7.
- [81] M. Lichtman, R. Rao, V. Marojevic, J. Reed, and R. P. Jover, "5G NR jamming, spoofing, and sniffing: Threat assessment and mitigation," in *Proc. IEEE Int. Conf. Commun.*, 2018, pp. 1–6.
- [82] D. Raca, D. Leahy, C. J. Sreenan, and J. J. Quinlan, "Beyond throughput, the next generation: A 5G dataset with channel and context metrics," in *Proc. ACM Multim. Syst. Conf.*, 2020, p. 303–308.
- [83] B. Manoj, M. Sadeghi, and E. G. Larsson, "Adversarial attacks on Deep Learning-based power allocation in a massive MIMO network," *arXiv:2101.12090*, 2021.
- [84] A. Yudnikov, D. Zabirow, V. Lyashev, and M. Kirichenko, "Doppler spread impact minimization via terminal sensors usage," in *Proc. IEEE Int. Black Sea Conference on Communications and Networking*, 2020, pp. 1–5.
- [85] M. Schüngel, S. Dietrich, D. Ginthör, S.-P. Chen, and M. Kuhn, "Heterogeneous synchronization in converged wired and wireless time-sensitive networks," in *Proc. IEEE Int. Conf. Factory Communication Systems*, 2021, pp. 67–74.
- [86] K. C. Zeng, S. Liu, Y. Shu, D. Wang, H. Li, Y. Dou, G. Wang, and Y. Yang, "All your {GPS} are belong to us: Towards stealthy manipulation of road navigation systems," in *Proc. USENIX Secur. Symp.*, 2018, pp. 1527–1544.
- [87] G. Falco, "Cybersecurity principles for space systems," *Journal of Aerospace Information Systems*, vol. 16, no. 2, pp. 61–70, 2019.
- [88] K. Jansen, M. Schäfer, D. Moser, V. Lenders, C. Pöpper, and J. Schmitt, "Crowd-gps-sec: Leveraging crowdsourcing to detect and localize gps spoofing attacks," in *IEEE Symp. Secur. Privacy*, 2018.
- [89] N. Lakshmanan, N. Budhdev, M. S. Kang, M. C. Chan, and J. Han, "A stealthy location identification attack exploiting carrier aggregation in cellular networks," in *USENIX Security Symposium*, 2021.
- [90] Y. E. Sagduyu, Y. Shi, T. Erpek, W. Headley, B. Flowers, G. Stantchev, and Z. Lu, "When wireless security meets machine learning: Motivation, challenges, and research directions," *arXiv:2001.08883*, 2020.
- [91] I. Ahmad, S. Shahabuddin, T. Kumar, E. Harjula, M. Meisel, M. Juntti, T. Sauter, and M. Ylianttila, "Challenges of ai in wireless networks for iot," *IEEE Ind. Elec. Magazin*, 2020.
- [92] N. Wang, W. Li, A. Alipour-Fanid, L. Jiao, M. Dabaghchian, and K. Zeng, "Pilot contamination attack detection for 5g mmwave grant-free iot networks," *IEEE Trans. Inf. Forensics Secur.*, vol. 16, pp. 658–670, 2020.
- [93] J. Suomalainen, A. Juhola, S. Shahabuddin, A. Mämmelä, and I. Ahmad, "Machine Learning threatens 5G security," *IEEE Access*, vol. 8, pp. 190 822–190 842, 2020.
- [94] M. Usama, J. Qadir, and A. Al-Fuqaha, "Adversarial attacks on cognitive self-organizing networks: The challenge and the way forward," in *Proc. IEEE Conf. Local Comput. Netw. Workshops*, 2018, pp. 90–97.
- [95] O. Ibitoye, O. Shafiq, and A. Matrawy, "Analyzing adversarial attacks against deep learning for intrusion detection in iot networks," in *Proc. IEEE Global Commun. Conf. IEEE*, 2019, pp. 1–6.
- [96] M. A. S. Monge, A. H. González, B. L. Fernández, D. M. Vidal, G. R. García, and J. M. Vidal, "Traffic-flow analysis for source-side DDoS recognition in 5G environments," *Elsevier J. Netw. Comp. Appl.*, vol. 136, pp. 114–131, 2019.
- [97] F. Wu, L. Xiao, W. Yang, and J. Zhu, "Defense against adversarial attacks in traffic sign images identification based on 5g," *Springer J. Wireless Commun. Netw.*, vol. 2020, no. 1, pp. 1–15, 2020.
- [98] G. Li, K. Ota, M. Dong, J. Wu, and J. Li, "DeSVig: Decentralized swift vigilance against adversarial attacks in industrial artificial intelligence systems," *IEEE T. Ind. Inf.*, vol. 16, no. 5, pp. 3267–3277, 2019.
- [99] L. Pajola, L. Pasa, and M. Conti, "Threat is in the air: Machine learning for wireless network applications," in *Proc. ACM Workshop Wireless Secur. Machin. Learn.*, 2019, pp. 16–21.
- [100] J. Qiu, L. Du, Y. Chen, Z. Tian, X. Du, and M. Guizani, "Artificial intelligence security in 5g networks: Adversarial examples for estimating a travel time task," *IEEE Vehic. Tech. Magazine*, vol. 15, no. 3, pp. 95–100, 2020.
- [101] Y. Arjoun and S. Faruque, "Artificial intelligence for 5G wireless systems: Opportunities, challenges, and future research direction," in *Proc. IEEE Ann. Comp. Commun. Conf. Workshop*, 2020, pp. 1023–1028.
- [102] Z. Luo, S. Zhao, Z. Lu, J. Xu, and Y. Sagduyu, "When attackers meet AI: Learning-empowered attacks in cooperative spectrum sensing," *IEEE Trans. Mobile Comp.*, 2020.
- [103] T. Erpek, Y. E. Sagduyu, and Y. Shi, "Deep learning for launching and mitigating wireless jamming attacks," *IEEE Trans. Cogn. Commun. Netw.*, vol. 5, no. 1, pp. 2–14, 2018.
- [104] M. Sadeghi and E. G. Larsson, "Adversarial attacks on deep-learning based radio signal classification," *IEEE Wireless Commun. Letters*, vol. 8, no. 1, pp. 213–216, 2018.
- [105] S. Kokalj-Filipovic, R. Miller, and J. Morman, "Targeted adversarial examples against RF deep classifiers," in *Proc. ACM Workshop Wireless Secur. Machin. Learn.*, 2019, pp. 6–11.
- [106] S. Bair, M. DelVecchio, B. Flowers, A. J. Michaels, and W. C. Headley, "On the limitations of targeted adversarial evasion attacks against deep learning enabled modulation recognition," in *Proc. ACM Workshop Wireless Secur. Machin. Learn.*, 2019, pp. 25–30.
- [107] S. Kokalj-Filipovic, R. Miller, and G. Vanhoy, "Adversarial examples in RF deep learning: Detection and physical robustness," in *Proc. IEEE Glob. Conf. Sig. Inf. Proc.*, 2019, pp. 1–5.
- [108] Y. E. Sagduyu, Y. Shi, and T. Erpek, "IoT network security from the perspective of adversarial deep learning," in *Proc. IEEE Ann. Int. Conf. Sensing. Commun. Netw.*, 2019, pp. 1–9.
- [109] F. Restuccia, S. D'Oro, A. Al-Shawabka, B. C. Rendon, K. Chowdhury, S. Ioannidis, and T. Melodia, "Generalized wireless adversarial deep learning," in *Proc. ACM Workshop Wireless Secur. Machin. Learn.*, 2020, pp. 49–54.
- [110] B. Kim, Y. E. Sagduyu, T. Erpek, K. Davaslioglu, and S. Ulukus, "Adversarial attacks with multiple antennas against deep learning-based modulation classifiers," in *Proc. IEEE Globecom*, 2020, pp. 1–6.
- [111] M. Z. Hameed, A. György, and D. Gündüz, "The best defense is a good offense: Adversarial attacks to avoid modulation detection," *IEEE T. Inf. Forensics Secur.*, vol. 16, pp. 1074–1087, 2020.
- [112] B. Kim, Y. E. Sagduyu, T. Erpek, and S. Ulukus, "Adversarial Attacks on Deep Learning Based mmWave Beam Prediction in 5G and Beyond," *arXiv:2103.13989*, 2021.
- [113] C. Joshi, J. R. Aliaga, and D. R. Insua, "Insider threat modeling: An adversarial risk analysis approach," *IEEE Trans. Inf. Forensics Secur.*, vol. 16, pp. 1131–1142, 2020.
- [114] R. Brewer, "Advanced persistent threats: minimising the damage," *Elsevier Netw. Secur.*, vol. 2014, no. 4, pp. 5–9, 2014.
- [115] M. Stevanovic and J. M. Pedersen, "An analysis of network traffic classification for botnet detection," in *Proc. IEEE Int. Conf. Cyber Situat. Awar., Data Analyt., Assessment*, Jun. 2015, pp. 1–8.
- [116] M. Wen, B. Ye, E. Basar, Q. Li, and F. Ji, "Enhanced orthogonal frequency division multiplexing with index modulation," *IEEE T. Wireless Commun.*, vol. 16, no. 7, pp. 4786–4801, 2017.



and his main expertise lies in the analysis of Network Intrusions, Phishing, and Adversarial Attacks.

Giovanni Apruzzese is a Post-Doctoral researcher within the Institute of Information Systems at the University of Liechtenstein since 2020. He received the PhD Degree and the Master's Degree in Computer Engineering (summa cum laude) in 2020 and 2016 respectively at the University of Modena, Italy. In 2019 he spent 6 months as a Visiting Researcher at Dartmouth College (Hanover, NH, USA) under the supervision of Prof. VS Subrahmanian. His research interests involve all aspects of big data security analytics with a focus on machine learning,



Rodion Vladimirov is a PhD student within the Institute of Information Systems at the University of Liechtenstein since 2021. He received the BSc. Degree in Mathematical Methods in Economics at Ural State Technical University in Russia (2015) and the MSc. Degree in Finance at the University of Liechtenstein (2020). His main research is focused on the security analysis of AI-systems in 5G networks and includes the investigation of attacks against Machine Learning components in critical infrastructures and the design of appropriate countermeasures.



Aliya Tastemirova received the Master’s Degree in Information Systems at the University of Liechtenstein in 2021. Within the same university, she worked as a Student Assistant at the Institute of Information Systems from 2020 to 2021. She is now a Software Developer at Odoo.



Pavel Laskov is Full Professor at the University of Liechtenstein and head of the Hilti Chair of Data and Application Security. He received PhD in computer science at the University of Delaware in 2001 and held research and teaching positions at the Fraunhofer Institute FIRST, University of Tuebingen and Huawei European Research Center. His research is focused on the development of techniques for detection and mitigation of security incidents, especially using custom-built AI techniques. As one of the pioneers of research on AI security, Pavel

Laskov co-designed the first proof-of-concept attacks against mainstream AI algorithms such as neural networks and Support Vector Machines.

APPENDIX A EXPERIMENTAL TESTBED

The datasets chosen for our CS are directly related to the exemplary network functions in 5G infrastructures described in §II-B. An overview of our CS and corresponding datasets is provided in Table IV. For each dataset, we report whether it contains real or synthetic data, its size in samples (after preprocessing), how many features are extracted to devise the corresponding ML methods and how many classes are considered (some CS address a regression problem denoted as “Regr”). Detailed descriptions and motivations for these datasets are provided in each CS.

TABLE IV: Summary of the datasets of our Case Studies.

CS#	Dataset Information					5G Network ML Function
	Name	Origin	Size	Classes	Features	
CS1	CTU13 [76]	Real	5.5M	2	13	Slicing
CS2	ElasticMon [80]	Synt	27K	Regr.	16	CQI Pred.
CS3	Irish 5G [82]	Real	2.4K	Regr.	2	CQI Pred.
CS4	RML2016 [36]	Synt	120K	6	256	AMR
CS5	PA-mMIMO [44]	Synt	335K	Regr.	20	Pow. Alloc.
CS6	DeepSlice [31]	Synt	63K	3	8	Slicing

All our experiments are performed on a machine equipped with an Intel Xeon W-2195 CPU with 36 cores, 256GB RAM, 2TB SSD NVMe, and Nvidia Titan RTX GPU. The implementation leverages Python3 and popular ML libraries: scikit-learn, Keras, and Tensorflow. More details on each CS are found in the remainder of this appendix.

A. CSI: dataset, preprocessing, RsP

This CS is based on the CTU13 dataset [76], collected in a large university campus (almost 300 internal hosts), and provided both as raw packet captures (PCAP) as well as NetFlows generated via Argus¹⁹. Despite being collected in 2013, the considered network is large, has a high bandwidth, and NetFlows are still widely employed today and can be used for slicing operations in 5G NI (as seen in [32, 33]). The CTU13 dataset is popular in the intrusion detection community for the development of ML botnet detectors (e.g., [115]), a task unrelated to the scope of this paper. However, we observe that CTU13 provides ground truth on 4 distinct classes of communications: *active* (e.g., a human user behaviour), *background* (e.g., some passive and automatic communications), *botnet* and *CnC*. Hence *the benign part* of CTU13 (in its 2 classes) can be used for the considered 5G slicing application proposed by the state-of-the-art. In this CS, we thus only consider the benign traffic of CTU, and exclude any malicious communication. The vast majority of NetFlows (over 95%) are for background.

We consider an attacker that controls 6 UEs (i.e., hosts) that perform both *normal* and *background* communications. Multiple PCAP traces are contained in CTU13, each captured at a different date. For our experiments, we use 8 out of the 12 provided traces, choosing those that contain a significant amount of traffic generated by the 6 UEs ‘owned’ by the attacker. We consider each trace independently, i.e., we use each trace to: extract its NetFlows, label them, and devise a specific ML classifier trained and tested on the same trace. Such workflow allows to assess the effects of attacks in different circumstances, as each trace can be considered as a separate scenario.

For this CS, we cannot directly operate on the provided NetFlows because they are not raw-data, and are not valid for RsP. We are forced to operate on the PCAP traces. Each PCAP trace of CTU13 is used to generate NetFlows exactly as done by CTU13 creators. Next, we label the generated NetFlows by using the exact procedure described in the ground truth information; we verify the correctness of our labelling scheme by cross-checking our labels with the original NetFlows in CTU13, *and we obtain an exact matching*. We preprocess our NetFlows by removing missing values and computing additional features related to the *IP addresses* and *Ports*. We cannot use the exact IP addresses to perform the classification, as they would lead to overfitting. Therefore, we differentiate between *internal* and *external* IPs, whereas the Ports are categorized on the basis of the IANA guidelines (as done in [55]). Hence, each NetFlow is described by the set of features F in Table V.

We partition our NetFlows in T and V with a 80:20 split (common in NetFlow analyses [55]). We considered many ML classifiers, but the *RF* outperformed the others so we use this algorithm for our baselines.

To create the RsP, we take the raw PCAP trace and extract all raw traffic of the 6 controlled UEs of the attacker. We then consider every packet originating by such 6 UEs, and we increase its payload by appending small chunks [0-300] of

¹⁹<https://openargus.org/>

TABLE V: Features F of CS1. The affected \bar{F} are in gray.

#	Feature name	Type
1,2	<i>Src/Dst IP type</i>	Bool
3,4	<i>Src/Dst port type</i>	Cat
5	<i>Flow Direction</i>	Bool
6	<i>Connection state</i>	Cat
7	<i>Duration (Dur)</i>	Num
8,9	<i>Src/Dst ToS</i>	Num
10	<i>SrcBytes (SrcByt)</i>	Num
11	<i>DstBytes (DstByt)</i>	Num
12	<i>TotBytes (Byt)</i>	Num
13	<i>TotPkts (Pkt)</i>	Num

random bytes; for TCP packets (which are stateful), we also ensure that the three-way handshake is carried out; we also ensure that the new packets will have the correct checksum; finally, we verify the integrity of the resulting PCAP trace and discard any inconsistent packets. This ‘myopic’ PCAP trace is then subject to the same procedure as the original PCAP trace: we extract, label and preprocess the NetFlows. In Table V, we denote with a light-gray background the features ‘consciously attacked’ representing $\bar{\mathcal{F}}$, and with a darker background the features \bar{F} that are influenced as a consequence of $\bar{\mathcal{F}}$.

For the attacks at inference stage, we submit the myopic NetFlows to the trained RF . For attacks at training-stage, we randomly take a portion (e.g., 25%) of the NetFlows in T that involve the UEs controlled by the attacker, and replace it with as many myopic NetFlows. We then retrain the RF using such ‘poisoned’ T and test it again on V to assess the effects on the whole environment.

All these operations are performed for each of the 8 PCAP traces. The results shown in Fig. 4 report the average performance of all these experiments.

B. CS2: dataset, preprocessing, RsP

The `ElasticMon` dataset is created by leveraging the Mosaic5G FlexRAN controller²⁰. The samples capture information related to the MAC, RRC and PDCP. `ElasticMon` is provided either as raw-data or as preprocessed features. We use the raw-data as basis, to which we apply the same preprocessing as in [41], and use the same 90:10 split for creating T and V . After preprocessing, we obtain the set F of 16 features used by the RF regressor²¹, which we report in Table VI. Here, a light-gray background denotes the features known and targeted by the myopic attacker ($\bar{\mathcal{F}}$); while a dark-gray background are the features \bar{F} altered as a consequence of $\bar{\mathcal{F}}$. All these features are numerical. We recall the 4 attack scenarios: $\bar{\mathcal{F}}_1=(RSRP)$; $\bar{\mathcal{F}}_2=(Byt)$; $\bar{\mathcal{F}}_3=(Pkt)$; $\bar{\mathcal{F}}_4=(Pkt,Byt)$.

In particular, we observe that the RSRP is linked to the RSRQ, so we must also change the latter when doing the RsP for $\bar{\mathcal{F}}_1$. Similarly, increasing the packets ($pktRx$) will also change their inter arrival time ($pktRxAiat$), which must be updated when applying the RsP for $\bar{\mathcal{F}}_2$ and $\bar{\mathcal{F}}_4$.

When determining the intensity of each RsP, we proceed as follows. For $\bar{\mathcal{F}}_1$, we change the RSRP randomly with another

TABLE VI: Features F of CS2. The affected \bar{F} are in gray.

#	Name	#	Name
1	RSRQ	9	totTbsUL
2	RSRP	10	pktRxSn
3	PHR	11	pktRx
4	totPrbDL	12	pktRxByt
5	totPduDL	13	pktTxSn
6	totTbsDL	14	pktRxAiat
7	totPrbUL	15	pktTxAiat
8	totPduUL	16	SFN

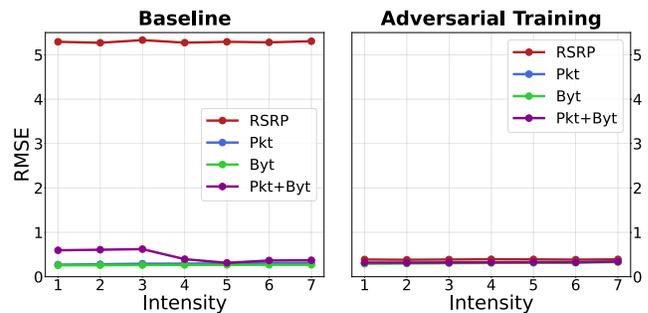
value in `ElasticMon`; we do this 7 times. For the remaining attacks (which target the $pktRx$ and $pktRxByt$), we increase the corresponding values of each sample as a function of their standard deviation in 7 intensity steps \mathcal{I} . Specifically, we use:

$$\mathcal{I} \in (\{0.1, 0.2, 0.5, 1, 2, 5, 10\} \times std_f)$$

where std_f is the standard deviation of the feature f in `ElasticMon`. Let us explain our workflow by using $\bar{\mathcal{F}}_2=pktRx$ as an example. We take the raw-data in `ElasticMon`, from which we compute the standard deviation of the $pktRx$ metric std_{pktRx} ; next, for each intensity value $i \in \mathcal{I}$, we multiply std_{pktRx} by the considered intensity i ; then, we add this value to the $pktRx$ of each sample in the `ElasticMon`. This creates multiple (7) sets of adversarial raw-data, each associated to a given intensity value $i \in \mathcal{I}$ targeting the features of $\bar{\mathcal{F}}_2$. We pre-process each of these sets with the same procedure as in [41], where we also update the derived $pktRxAiat$ feature. We then submit these myopic sets to the baseline RF regressor at inference stage.

We now explain our application of the two countermeasures. We apply *feature removal* 4 times, each time by re-training the baseline RF on the same T but without considering the features influenced by each of the 4 myopic attacks. We also apply *adversarial training* 4 times, by isolating a small portion (5%) of T , denoted T_{aug} , on which we apply the RsP of each specific myopic scenario at all intensity levels. We then re-train the baseline RF on the original T as well as all the ‘perturbed’ variants of T_{aug} .

For completeness, we present in Fig. 10 the efficacy of our attacks as measured via $RMSE$ (lower is better) instead of accuracy as in Fig. 5, shown in the main paper.

Fig. 10: CS2: Attacks and Defense (Baseline $RMSE=0.25$)

C. CS3: time series

We report in Figs. 11 the time series of the predicted CQI. The blue lines correspond to predictions on clean data, while

²⁰https://mosaic5g.io/apidocs/flexran/flexran_spec_v2.2.3.html

²¹To compute accuracy, we round each predicted CQI to the nearest integer.

red lines correspond to adversarial behavior of the two attacks considered in the paper. The time series at the top (bottom) of Figs. 11 correspond to the ‘static’ (‘driving’) mobility pattern.

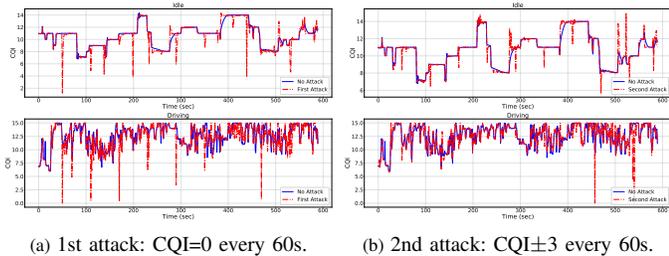


Fig. 11: CS3: predicted CQI with and without myopic examples.

D. CS4 Dataset, RsP, related work

We considered many ‘shallow’ ML classifiers, but the *RF* outperformed the others so we use this algorithm as shallow-learning baseline, whereas the deep-learning baseline uses the same *DN* architecture as prior work. To devise our baselines, we take the *RML2016.10a* dataset, filter the modulations pertaining to 5G (as done in [16], specifically: BPSK, QPSK, 8PSK, CPFSK, GFSK, WBFM), and partition it with the same 50:50 split as in [17] into T and V . We use the exact *DN* architecture as in previous work [17, 36], while the *RF* is finetuned to ensure optimal classification performance.

For each of the 4 considered attack scenarios, we apply the *RsP* in the same way as in CS2 (described in Appendix A-B). We first determine the measurements to perturb: randomly chosen (either 25, 50 or 100 for \bar{F}_1 , \bar{F}_2 , \bar{F}_3 , respectively) or by using the top-25 most significant in the worst-case (\bar{F}_4). Then, we increase the corresponding measurements of each sample as a function of their standard deviation, in 7 intensity steps as done in CS2. We ensure that all such increments fall within acceptable ranges. Finally, we note that *I/Q* measurements are orthogonal [116], hence $\bar{F}=\bar{F}$.

Because the attacks in \bar{F}_{1-3} involve a lot of randomness, we repeat these attacks 20 times (each time by choosing different measurements). In our results (in Fig. 7a) we report the average values and their standard deviation (as a black vertical line on each marker). For the attacks of related work [16, 17], we report the values directly as shown in the corresponding paper because they target a very similar system: the only difference is that [17] does not remove the analog modulations (but their *DN* obtains the same baseline performance as ours), whereas [16] uses only the samples with *SNR*=10db to train and test their *DN*, which explains the slight discrepancy in baseline performance—reported as a black line in Fig. 7b.

E. CS6: details and insider myopic attacks

The *DeepSlice* dataset uses the 5G specifications to describe KPI representing the main 5G use-cases (eMBB, MMTc, URLLC), each denoting a slice and, hence, a label.

We do not perform any pre-processing to *DeepSlice*. We report in Table VII the feature set F describing each sample and analyzed by the baseline *DN*. We partition the dataset by using the same split as [31] of 90:10 for T and V ,

matching their perfect accuracy: $Acc=1.00$. A true myopic attacker can only affect the *day* and *hour*, but such attempts are not successful against the *DN*, as shown in §V-F.

‘Insider’ myopic Attacks. An ‘insider’ myopic attacker that can operate from within the 5G NI and tamper with the gNB configurations. She can thus affect the *PktLossRate* or the *PktDelayBudget* associated to each slice, thus $\bar{F}=(PktDelayBudget, PktLossRate)$. We design 3 attack scenarios: $\bar{F}_1=(PktDelayBudget)$, $\bar{F}_2=(PktLossRate)$, $\bar{F}_3=(PktDelayBudget, PktLossRate)$; in all these cases, $\bar{F}=\bar{F}$, shown in gray in Table VII.

We craft the *RsP* in the same way as described in Appendix A-B. All the affected features are independent.

We counter such attacks with *adversarial training* and *feature removal*, applied in the same way as in CS2 (§V-B).

TABLE VII: F of CS6. Gray rows denote the insider \bar{F} .

#	Feature name	Type
1	UseCase	Cat
2	UEcategory	Cat
3	Technology	Cat
4	Day	Num
5	Hour	Num
6	GuaranteedBitRate	Bool
7	PacketLossRate	Num
8	PacketDelayBudget	Num

Results. We report the results of these insider attacks in Fig. 12, reporting accuracy as a function of *RsP* intensity. We use full lines for attacks against the baseline *DN*, and dotted lines for attacks against the *DN* hardened via adversarial training; as in CS3, feature removal always defused the attacks. The tradeoff of the countermeasures is shown in Table VIII. We observe that \bar{F}_2 has no effect, but \bar{F}_1 and \bar{F}_3 can decrease the accuracy even at low intensities. Feature removal induces an unfavorable tradeoff; in contrast, adversarial training preserves baseline performance but its protection is low.

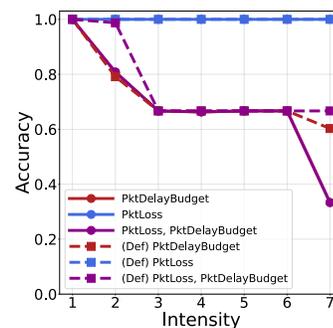


Fig. 12: CS6: Attack and Defense (Base: $Acc=1.00$)

TABLE VIII: CS6-insider: T (lower is better, $T=1$ is no change).

Defense	\bar{F}_1	\bar{F}_2	\bar{F}_3
Adv. Tra.	1.00	1.00	1.00
Fea. Rem.	1.50	1.00	1.50