

A multimodal and perturbation-aware learning approach for robust traffic classification

Idio Guarino ^{a,*}, Giampaolo Bovenzi ^b, Alfredo Nascita ^b, Domenico Ciunzo ^b,
Damiano Carra ^a, Antonio Pescapè ^b

^a University of Verona, Italy

^b University of Napoli Federico II, Italy

ARTICLE INFO

Keywords:

Deep learning
Encrypted traffic
Multimodal techniques
Robust design
Traffic classification

ABSTRACT

Traffic Classification (TC) is pivotal for network management, cybersecurity, and Quality of Experience (QoE) monitoring. However, while Deep Learning (DL) has significantly advanced TC, most existing works assume static, idealized conditions, overlooking key challenges of real-world deployments—such as traffic variability, routing asymmetries, out-of-order packet arrivals, and partial visibility at the Vantage Points (VPs). This motivates the need for robustness evaluations under such scenarios.

In this work, we investigate the robustness of *state-of-the-art* (SOTA) TC models under realistic, yet controlled, perturbation scenarios. Specifically, we introduce novel, model-agnostic traffic perturbations—simulating time jitter, retransmissions, and partial visibility—to reflect conditions commonly encountered in live network traffic. We evaluate our approach on three public datasets—i.e., VPN-16, MIRAGE-19, and MIRAGE-24—and show how MIMETIC-ENHANCED, a multimodal model, tends to outperform two representative single-modal counterparts both in terms of TC effectiveness on clean traffic and robustness under perturbations. Nonetheless, our analysis also reveals that multimodal models remain vulnerable under specific perturbation settings. To address this limitation, we propose a *model-agnostic perturbation-aware training framework* based on Supervised Data Augmentation (Aug) and Contrastive Learning (CL)—considering both self-supervised and supervised variants. Unlike architecture-specific solutions, our approach operates at the *learning strategy level*, allowing it to be seamlessly applied to diverse classifiers without requiring structural modifications. Adopting MIMETIC-ENHANCED as a primary multimodal case study, we integrate the proposed strategies into its two-stage training pipeline. Experimental results demonstrate that perturbation-aware training not only improves TC effectiveness on clean (i.e., unperturbed) traffic—particularly when applied across both training stages—but also significantly strengthens the model's robustness under diverse and realistic perturbation scenarios. Furthermore, we investigate Out-of-Distribution (OOD) detection, model calibration, and TC effectiveness in low-data regimes. Finally, we explicitly demonstrate the framework's *generalizability* by validating it on other SOTA architectures, spanning both single- and multi-modal approaches.

1. Introduction

TC plays a central role in network operations, supporting traffic monitoring and management, cybersecurity enforcement, and QoE monitoring. The field is well-established, with a rich body of research and numerous comprehensive surveys, e.g., [1,2].

Over the years, TC has evolved from traditional Machine Learning (ML) with handcrafted features to DL models that extract patterns directly from raw network data. This shift unfolded in two waves [3]. The first wave, starting in the early 2000s, used ML with per-packet (e.g.,

packet size, inter-arrival time) and per-flow statistics (e.g., total bytes, number of packets, ports) to classify a small set of applications—often with high accuracy and deployability (e.g., supporting early classification by relying only on a few initial packets) [4–7]. However, these methods struggled to keep pace with rising traffic complexity driven by (a) encryption, (b) proliferation of mobile devices, and (c) diversity of the tasks/activities performed by services running over the Internet [8,9]. The second wave has embraced DL, leveraging raw payloads or statistical features in end-to-end architectures to model traffic in an “expert-free” fashion [8,10–14]. In this context, more sophisticated

* Corresponding author.

E-mail addresses: idio.guarino@univr.it (I. Guarino), giampaolo.bovenzi@unina.it (G. Bovenzi), alfredo.nascita@unina.it (A. Nascita), domenico.ciunzo@unina.it (D. Ciunzo), damiano.carra@univr.it (D. Carra), pescapè@unina.it (A. Pescapè).

<https://doi.org/10.1016/j.comnet.2026.112090>

Received 29 July 2025; Received in revised form 30 January 2026; Accepted 4 February 2026

Available online 25 February 2026

1389-1286/© 2026 The Author(s). Published by Elsevier B.V. This is an open access article under the CC BY license (<http://creativecommons.org/licenses/by/4.0/>).

efforts have explored *multimodal DL architectures* [9,15,16], using different (but complementary) facets of traffic data, such as packet payloads and statistical features developed in the first wave [8,10–14]. Indeed, recent literature demonstrates that combining heterogeneous views effectively tackles challenging operational conditions, such as mitigating the impact of noisy annotations [17].

Yet, a key limitation remains: *most SOTA models are trained and tested in static, closed-world settings with clean data from limited VPs*. In contrast, real-world deployments—especially at intermediate network locations such as ISP routers or enterprise middleboxes—are characterized by noisy measurements, asymmetric flows, and partial visibility. This gap highlights *the need to evaluate and improve model robustness* [3] under realistic (perturbed) conditions—especially for encrypted traffic.

Accordingly, this paper tackles a critical shortcoming in today’s traffic classifiers: *their fragility under real-world, dynamic network conditions*. We propose a framework to systematically evaluate model robustness by introducing “black-box” perturbations—transformations applied directly to traffic traces that simulate common network phenomena such as jitter, retransmissions, asymmetric paths, and delayed monitoring. These perturbations are model-agnostic and reflect challenges faced by classifiers operating beyond the client edge, where access to clean data and model internals is precluded. Our emphasis is on *early TC* [4], where decisions must be made based on the initial packets of a flow—this is crucial for maximizing the reactivity of, and minimizing the latency of, network management tools [4].

To improve robustness under such conditions, we thoroughly explore *perturbation-aware training*. These techniques help the model learn invariant representations that are resilient to timing distortions and visibility gaps. While our experiments use MIMETIC-ENHANCED [16]—a multimodal traffic classifier that fuses packet sequences and payloads—the proposed training pipeline *is broadly applicable to any multimodal TC model*. Accordingly, the **contributions** of this work are:

- We propose a taxonomy of *realistic network perturbations*—spanning timing, retransmissions, directional asymmetries, and changes in *VP*—and their systematic application to evaluate the robustness of *DL*-based encrypted traffic classifiers *without relying on model internals*.
- We systematically benchmark the *robustness* of *SOTA* multimodal *TC* models, investigating their susceptibility to network-induced perturbations. While these models perform well on clean data [16], we investigate the potential benefit of multimodality (i.e., the fusion of packet-sequence features and payload content) in close-to-real-world *TC*. By analyzing how different input modalities respond to these conditions, we assess whether (and to what extent) multimodal *DL* can improve resilience by compensating for modality-specific weaknesses.
- Using MIMETIC-ENHANCED as a primary case of study, we investigate *perturbation-aware* training pipelines by injecting realistic traffic transformations into the learning process. Our approach integrates *both Aug* and *CL*—spanning self-supervised and supervised variants—into the *pre-training* and *fine-tuning* stages of model training. This design helps models learn robust representations, which are essential in real-world network environments.
- We assess whether traffic transformations can be leveraged not only to improve model robustness but also to enhance generalization across network conditions. In this respect, we extend the evaluation beyond accuracy—analyzing latent space structure, model calibration, out-of-distribution detection, and sample-scarcity conditions—to capture *TC* effectiveness in practical deployments.
- Finally, we demonstrate that the proposed learning framework is effectively *model-agnostic*. By explicitly validating its *broad generalizability* on representative classifiers—spanning both single-modal [11, 13] and multimodal [15,18] paradigms—we prove that it enhances robustness regardless of the specific model design.

Table 1

List of acronyms and abbreviations in alphabetical order.

Acronym	Definition
AI	Artificial Intelligence
Aug	Supervised Data Augmentation
AUROC	Area Under the Receiver Operating Characteristic Curve
CCE	Categorical Cross-Entropy
CL	Contrastive Learning
DBI	Davies-Bouldin Index
DHiding	Direction Hiding
DHidingloss	Direction Hiding with Packet Loss
DIR	Packet Direction
DL	Deep Learning
ECE	Expected Calibration Error
FPR	False Positive Rate
FT	Fine-Tuning
GenAI	Generative Artificial Intelligence
IAT	Inter Arrival Time
ID	In-Distribution
LSTM	Long Short-Term Memory
ML	Machine Learning
NT-Xent	Normalized Temperature-scaled Cross-Entropy
OOD	Out-of-Distribution
PAY	initial raw bytes transmitted by the application
PHiding	Packets Hiding
PL	Transport-level Payload Length
PS	Packet Size
PSQ	sequence of header fields from the initial packets
PT	Pre-Training
QoE	Quality of Experience
RTO	Retransmission Timeout
RTT	Round-Trip Time
SDU	Service Data Unit
SIL	Silhouette score
SimCon	Self-supervised Contrastive Learning
SOTA	state-of-the-art
SupCon	Supervised Contrastive Learning
TC	Traffic Classification
TCP	Transmission Control Protocol
TCPWIN	TCP Window Size
TJitter	Directional Time Jittering
TLS	Transport Layer Security
TO	Traffic Object
TPR	True Positive Rate
VAE	Variational Autoencoder
VP	Vantage Point
VPC	Vantage Point Centralization
VPCloss	Vantage Point Centralization with Packet Loss
XAI	eXplainable Artificial Intelligence

Our study is validated on *three public datasets*: VPN-16 [19], MIRAGE-19 [8], and MIRAGE-24 [20], showing practical relevance and broad applicability.

The remainder of this paper is organized as follows. *Section 2* reviews related work on *TC* robustness. *Section 3* introduces our multimodal architecture, defines the proposed perturbations, and details the perturbation-aware learning methodologies. *Section 4* outlines the experimental setup. Results are discussed in *Section 5*, and conclusions follow in *Section 6*.

Tables 1 and *2* list the acronyms and symbols used in this paper, together with their definitions, to aid readability.

2. Related works

Recent advances in *TC* have applied *DL* to curated datasets, yielding high accuracy under *static, idealized conditions*. Yet these models often assume full visibility, balanced classes, and stable patterns—assumptions that break down in real-world deployments, especially with encrypted or incomplete traffic. This section surveys how prior work addresses robustness to input perturbations, focusing on how *TC* models respond to deviations from ideal conditions. *Table 3* summarizes the landscape, also detailing the type of *Traffic Object (TO)* used (e.g., flows vs. biflows),

Table 2
List of symbols and their definition.

Symbol	Definition
α	RTT fraction attributed to Client-to-VP path
β	Temperature parameter for contrastive loss
B	Mini-batch of samples
$C \rightarrow S$	Client-to-Server traffic direction
$C \leftarrow S$	Server-to-Client traffic direction
δ	Percentage of retained training data (Data-Constrained)
e	Embedding dimension size ($e = 10$)
$F_\phi(\cdot)$	Fusion function parameterized by ϕ
$f_\theta(\cdot)$	Feature extractor (backbone) parameterized by θ
$f_\varphi(\cdot)$	Classifier head parameterized by φ
\mathcal{L}	Generic loss function (e.g., \mathcal{L}_{CCE})
LP	Number of lost initial packets
L_{\min}, L_{\max}	Minimum and maximum number of retransmissions
M	Number of modalities ($M = 2$)
N	Total number of samples in the dataset
N_b	Number of initial bytes in PAY ($N_b = 576$)
N_p	Number of initial packets in PSQ ($N_p = 10$)
p	Packet loss probability
$p(x)$	Soft-output probability distribution
$\hat{p}(x)$	Predicted confidence
r	Retransmission rate
$sim(\cdot)$	Similarity function (e.g., cosine similarity)
σ_d	Standard deviation for directional time jitter
\mathcal{T}	Predefined set of transformations
T_{out}	Retransmission Timeout duration
$\tau(\cdot)$	A single transformation function sampled from \mathcal{T}
X	Set of original training samples $\{x_i\}_{i=1}^N$
\hat{X}	Set of transformed training samples
\bar{X}	Final augmented training set ($X \cup \hat{X}$)
x	Generic input sample (biflow)
\tilde{x}	Transformed input sample
y, \hat{y}	True class label and predicted label
z	Latent feature vector ($z = f_\theta(x)$)
\emptyset	Empty set (indicating no transformations, $\mathcal{T} = \emptyset$)
\uparrow, \downarrow	Direction of improvement (Higher/Lower is better)

Input type(s) adopted (e.g., raw bytes, packet sequences), **Datasets**, and **TC Task** considered.

We categorize prior efforts by the **Perturbation Type**: Section 2.1 covers methods that introduce deliberate, adversarial modifications, while Section 2.2 focuses on techniques that simulate network fluctuations or changes in operating conditions. Finally, Section 2.3 positions our work with respect to the existing state of the art.

2.1. Robustness to intentional manipulations

A first line of research investigates **intentional perturbations** (“I” in column **Perturbation Type**) and applies adversarial ML [39] with the goal of degrading model performance through carefully crafted input perturbations. These adversarial examples exploit the vulnerability of **Artificial Intelligence (AI)** models to small input manipulations that can lead to incorrect predictions. Adversarial attacks are typically categorized based on the attacker’s knowledge of the “target” model, typically falling into *three categories*: (i) *white-box* attacks assume full access to the model’s architecture and gradients; (ii) *black-box* attacks operate without internal model access, relying only on observed outputs; (iii) *hybrid* attacks exploit access to a “surrogate” model during the perturbation generation phase and apply the resulting perturbations in a model-agnostic fashion.

In the *white-box* setting, Sadeghzadeh et al. [22] craft universal perturbations (named “pad”, “payload”, and “burst” attacks) targeting payloads and packet sequences, while Jin and Apostolaki [36] *PANTS* framework uses satisfiability modulo theory solvers to generate valid yet misleading inputs.

Black-box strategies—e.g., *Tantra* [23] and time-based evasion Zhang et al. [24]—reshape temporal patterns to evade malicious traffic detection without altering original content. *Prism* approach [27] goes further by injecting asymmetric perturbations that mask traffic signatures. Such

perturbations are driven by the temporal transitions of the application behaviour.

Other works adopt a *hybrid strategy* where adversarial perturbations are generated with white-box access to a surrogate model and then directly reused against the target model under black-box conditions. Lu et al. [37] introduce *TSAF*, a framework that uses white-box access to craft time series-based adversarial perturbations via gradient optimization. Then, these are stored as generic patterns and injected into clean traffic at runtime through a masking mechanism. Ding et al. [31] propose a method for generating universal, transferable adversarial perturbations by optimizing a loss function based on internal network activations. Crafted offline with white-box access, these perturbations generalize across models and datasets, enabling efficient, model-agnostic attacks with high success rates and low overhead. Xie et al. [28] introduce *Cactus*, a client-side defense mechanism that applies protocol-compliant perturbations—modifying packet size, timing, and direction—to obfuscate encrypted **Transmission Control Protocol (TCP)** traffic. Although not tailored to a specific model, these transformations are intentionally designed to disrupt DL-based traffic analysis, effectively reducing classification accuracy while preserving application semantics.

Collectively, these studies highlight the fragility of traffic classifiers when exposed to adversarial manipulations and motivate the need for robust training techniques. Works such as [22,28,31,37] primarily assess the robustness of state-of-the-art classifiers, revealing their susceptibility to carefully crafted perturbations, even when such perturbations are limited or protocol-compliant (column **Base Classifier Robustness**). Other contributions, including [23,24,27,36], extend the analysis by proposing enhanced models or defense mechanisms and evaluating their behavior under adversarial conditions (i.e., perturbed inputs). They evaluate how these strengthened designs behave when exposed to perturbed inputs (subcolumn **Ⓢ**), thus providing a more complete picture of robustness under attack. These efforts offer a more holistic view of robustness, combining attack scenarios with mitigation strategies. Still, *adversarial approaches differ in goals and assumptions from studies focused on natural, deployment-induced perturbations—a gap this paper aims to bridge*.

2.2. Robustness to network changes

Beyond adversarial attacks, a substantial line of work examines how network traffic classifiers respond to **non-adversarial perturbations** (“NC” in column **Perturbation Type**)—i.e., variations that *occur naturally in operational settings*. These perturbations, while not malicious, can degrade classification performance or, conversely, serve as a form of data augmentation to improve model generalization and robustness in the face of real-world variability. These include jitter, delay, packet loss, reordering, and incomplete sessions, which can all affect the structure and timing of traffic flows. Unlike adversarial inputs, these perturbations are not crafted to induce misclassification, but they nonetheless pose challenges for models trained under idealized conditions.

Some studies focus on training-time perturbations and assess their impact on the **TC** effectiveness, that is, the model’s performance with unperturbed input. These perturbations are applied during training to promote invariance to minor input variations and to reduce overfitting to specific traffic characteristics.

For instance, Guarino et al. [25] compare contrastive, transfer, and meta-learning strategies, showing that contrastive learning leads to more transferable representations, even when labeled data is limited. Wang et al. [30] benchmark 18 augmentation strategies for **Inter Arrival Time (IAT)**, **Transport-level Payload Length (PL)**, and **Packet Direction (DIR)**, finding that structural changes in **TOs** (e.g., masking, reordering) yield a stronger generalization than amplitude-based perturbations. Along similar lines, Yang et al. [32] propose *TrafCL*, a framework that combines **CL** with session-level augmentations to produce robust representations that remain effective even when sessions are incomplete or labeled samples are scarce. Horowicz et al. [35] explore

Table 3

Overview of related work. The last row summarizes our proposal.

Reference	Year	Perturbation Type	Base Classifier Robustness	Strengthened Design		Classification Details					
				● _{2b}	● ₁	Early	MM	TO	Input	Datasets	Task
Liu et al. [21]	2019	NC	●	●	●	○	○	BF	F	Private	TC
Sadeghzadeh et al. [22]	2021	I	●	○	○	●	○	BF	B, PS	VPN-16	TTC
Sharon et al. [23]	2022	I	●	○	●	○	○	F	PS	CIC-IDS2017 Kitsune	AC
Zhang et al. [24]	2022	I	●	○	●	○	○	BF	F	CIC-IDS2018	AC
Guarino et al. [25]	2023	NC	○	●	○	●	○	BF	PS	MIRAGE-19 AppClassNet VPN-16	TC
Xie et al. [26]	2023	NC	●	○	●	○	○	BF	PS	CICDoHBrw-2020	TC
Li et al. [27]	2023	I	●	○	●	●	○	BF	PS, B, F	USTC2016 VPN-16	ATC
Xie et al. [28]	2024	I	●	○	○	○	○	BF	PS	Private	ATC, WF
Deng et al. [29]	2024	NC	○	○	●	—	○	BF	PS, B, F	SJTU-AN21 VPN-16	TC
Wang et al. [30]	2024	NC	○	●	○	●	○	BF	PS	ISCXTor2016 MIRAGE-19 MIRAGE-2022	TC
Ding et al. [31]	2024	I	●	○	○	●	○	P	B	Private VPN-16 USTC2016 CIC-IoT2023	TC
Yang et al. [32]	2024	NC	○	●	○	—	○	BF	B, PS	CICDoHBrw-2020 RAT-PN	TC
Hajaj et al. [33]	2024	NC	○	●	○	○	○	F	PS	USTC-VPN UC-Quic-Davis FLASH	SC ATC
Jiang et al. [34]	2024	NC	●	○	●	●	●	BF	B, PS	PHAM21 REN19 Private	TC
Horowicz et al. [35]	2024	NC	○	●	○	○	○	F	B	UC-Quic-Davis VPN-16 + ISCXTor2016	TC
Jin and Apostolaki [36]	2025	I	●	○	●	○	○	BF	F	VPN-16 UTMobileNetTraffic2021	ENC ATC
Lu et al. [37]	2025	I	●	○	○	○	○	BF	F, PS	VCAML CIC-IDS2017	QoE-I TC
Wang et al. [38]	2025	NC	●	○	●	●	●	F	B, PS	USTC2016 Malicious-TLS	AC
This work	2025	NC	●	●	●	●	●	BF	B, PS	VPN-16 MIRAGE-19 MIRAGE-24	TC

Perturbation Type: Network Changes (NC), Intentional (I). **Base Classifier Robustness:** Present (●), Absent (○).

Strengthened Design: TC Effectiveness on “unperturbed” traffic (●_{2b}), Robustness to “perturbed” traffic (●₁).

Early Classification (e.g., $N_p \leq 20$): Present (●), Absent (○), Not inferable from the original work (—). **Multi-modal architecture** (MM): Present (●), Absent (○).

Traffic Object (TO): Biflow (BF), Flow (F), Network Channel (NC). **Traffic Input:** Packet Sequences (PS), Raw Bytes (B), Statistical Features (F).

Eval. Task: Traffic Classification (TC), Attack Classification (AC), App Traffic Classification (ATC), Traffic Type Classification (TTC), Website Fingerprint (WF), Encapsulation (ENC), Service Classification (SC), QoE Inference (QoE-I).

generalization in a few-shot learning setting using CL. Their approach, based on FlowPic-based inputs, demonstrates that domain-specific perturbations, such as Packet Loss or Round-Trip Time (RTT) changes, outperform generic computer vision augmentations (e.g., rotate, flip, and color jitter). Differently, Hajaj et al. [33] devise a Long Short-Term Memory (LSTM)-based data augmentation technique to produce synthetic traffic that enriches the training dataset and improves the effectiveness of encrypted traffic classifiers on unperturbed inputs.

Other works evaluate robustness under perturbed inputs, without measuring improvement in TC effectiveness. Deng et al. [29] introduce AN-Net, a classifier exposed to two types of synthetic noise: (i) irrelevant packet noise, simulating traffic from unrelated flows, and (ii) per-packet attribute noise, reflecting fluctuations in network environments, leading to discrepancies between training and test conditions. Additionally, they propose a fusion mechanism that combines attributes from different modalities to mitigate the effects of per-packet attribute noise and improve the robustness of the learned representations. Some studies evaluate models across both perturbed and unperturbed data, offering a more comprehensive robustness evaluation. Focusing on Transport Layer Security (TLS) traffic, Xie et al. [26] propose Rosetta, a self-supervised training strategy with TCP-aware augmentations. Their approach aims to teach models the implicit semantics of encrypted traffic flows to obtain robust generalization, validated across diverse network setups by replaying traffic across diverse network setups with hosts in different geographic locations. Jiang et al. [34] design a packet permutation-aware model that handles out-of-order packet arrival within flows and uses

payload information to enrich the input representation. Their evaluation on both clean and perturbed inputs shows consistent improvements under noise. Conversely, Liu et al. [21] present a more complete evaluation. Their work addresses feature drift by selecting more stable traffic representations and evaluates model robustness under three dimensions: (i) unperturbed test data, (ii) perturbed test data, and (iii) training-time augmentation. This threefold analysis provides one of the most complete assessments of robustness in current TC literature.

Finally, a recent trend explores the use of Generative Artificial Intelligence (GenAI) to improve the quality of learned representations in complex environments. For instance, Wang et al. [38] propose ER-CMGI, a framework combining Variational Autoencoders (VAEs) and Diffusion Models. By generating consistent cross-modal features (e.g., aligning packet lengths with byte sequences), this architectural approach achieves high classification performance on complex traffic patterns. However, such methods rely on computationally-intensive components during inference to reconstruct or align features, aiming to solve robustness through architectural complexity rather than training dynamics.

Collectively, these works highlight the diversity of strategies used to account for real-world traffic variability, ranging from data augmentations to architectural adaptations.

2.3. Literature gap and positioning

While prior work has explored both intentional attacks and network perturbations, important gaps remain, particularly concerning network

-specific transformations (“NC” in the column **Perturbation type**). We address these challenges by introducing *five* categories of traffic perturbations that *emulate realistic network conditions without requiring internal model access*. Our black-box approach highlights model blind spots that clean inputs may conceal, aligning with the goals of adversarial learning. Nonetheless, unlike typical adversarial ML methods (“I” in the column **Perturbation type**) that depend on white-box assumptions or latent feature manipulations, our perturbations operate directly on network traffic and reflect plausible variations, avoiding artificial techniques borrowed from unrelated domains like computer vision [25,30,35].

Only a few prior works, e.g., Horowicz et al. [35], employ comparable black-box, network-level perturbations to assess the robustness of TC models. However, *state-of-the-art approaches typically restrict their analysis to a single dimension/analysis*: they either evaluate how baseline models behave under perturbed inputs (**Base Classifier Robustness**), or they assess robust models (**Strengthened Design**) in terms of *effectiveness* on clean data (sub-column (A)) or their *robustness* on perturbed data (sub-column (B)). When robust models are designed in the same context [35], *a limited number of perturbations and a smaller set of design strategies are explored*. In contrast, our work offers a *comprehensive robustness evaluation by covering all these scenarios using three different representative datasets*. Equally important, our model combines byte-level content and packet sequence patterns (“B” and “PS” in the column **Input**, respectively) *using a multimodal architecture (MM)*, addressing the relevant challenge of defining network perturbations for traffic payload. To the best of our knowledge, alongside Jiang et al. [34], only the recent generative framework by Wang et al. [38] adopts a comparable multimodal design. However, while the latter pursues feature alignment through computationally intensive diffusion models, our framework prioritizes deployment efficiency, embedding robustness directly into the training dynamics. Furthermore, although their evaluation covers traffic with intrinsic perturbations (e.g., jittering and packet loss), such data is also employed during training [38]. Consequently, their analysis does not assess robustness against *unseen* or varying network conditions.

Furthermore, *we explicitly support early TC (Early)*, which remains under-explored in the literature [22,25,27,30,31,34]. In this way, *our work offers a holistic perspective on robustness in TC and fills several key gaps in the existing literature*. The last row of Table 3 summarizes the characteristics of our approach to facilitate direct comparison with the literature.

3. Methodology

Section 3.1 describes the multimodal DL classifier used for TC. Then, Section 3.2 introduces the set of perturbation scenarios considered in this work. Last, Section 3.3 describes the learning approaches used to leverage the defined perturbations.

3.1. Multi-modal traffic classification

TC seeks to assign meaningful labels to entities observed in network traffic, which we refer to as TOs, corresponding to logically grouped sets of packets, as seen from a VP. In this study, as the fundamental TO, we treat bidirectional flows—or *biflows*—each consisting of all packets exchanged between a pair of endpoints, identified by a the tuple: $\langle IP_{src}, port_{src}, IP_{dst}, port_{dst}, L4_{proto} \rangle$, where the source and destination are interchangeable [40].

Our goal is to classify each biflow based solely on its early packets—a strategy known as *early classification*. By making decisions at the outset of a flow, the system can promptly adapt to changing network conditions and enforce traffic management policies, which is vital in real-time scenarios. A biflow can be characterized through *multiple modalities*, each offering a *distinct perspective* for classification [15].

The two most direct inputs are:

1. the **initial raw bytes transmitted by the application (PAY)**—namely, the transport-layer **Service Data Unit (SDU)**, and

Table 4

Impact of each perturbation on PAY and PSQ traffic inputs.

	PAY	PSQ
TJitter	⊗	⊙
RTOfaft	⊗	⊙
RTObef	⊗	⊙
PHiding	⊙	⊙
DHiding	⊙	⊙
DHiding _{loss}	⊙	⊙
VPC	⊗	⊙
VPC _{loss}	⊗	⊙

⊙: Affected.
⊗: Possibly Affected.
⊘: Unaffected.

2. the **sequence of header fields from the initial packets (PSQ)**, such as **Packet Size (PS)**, **DIR** and **IAT**.

Formally, a generic DL classifier can be represented by a function $f_{\theta, \phi}: \mathcal{X} \rightarrow \mathcal{C}$ that transforms input data into output labels. Specifically, $f_{\theta, \phi}(\cdot)$ can be decomposed into two parts, the model backbone $f_{\theta}: \mathcal{X} \rightarrow \mathcal{Z}$ —the *feature extractor*—and the model head $f_{\phi}: \mathcal{Z} \rightarrow \mathcal{C}$ —the *classifier*. The feature extraction process maps input data x into a latent representation $z = f_{\theta}(x) \in \mathcal{Z}$, whereas the classifier maps the latent representation to a per-class probability distribution $p(x) = f_{\phi}(z) \in [0, 1]^{|\mathcal{C}|}$ (i.e., *soft-output*), typically via a “softmax” activation in the final layer, where each component $p_c(x)$ indicates the confidence of the model that the input x belongs to class c . The predicted label is obtained as $\hat{y}(x) = \arg \max_{1, \dots, |\mathcal{C}|} p_i(x)$.

In a standard *supervised learning* setting, model training seeks to optimize parameters (θ, ϕ) to minimize a classification loss $\mathcal{L}(y, \hat{y})$, where $y \in \mathcal{C}$ is true class label. The typical loss is the *Categorical Cross-Entropy (CCE)*: $\mathcal{L}_{CCE}(y, \hat{y}) = -\sum_{i=1}^{|\mathcal{C}|} \mathbf{1}_{\{y=i\}} \log(p_i)$, which compares the soft-output provided by the model with the true label, penalizing lower probabilities assigned to the true class. Since the classifier $f_{\phi}(\cdot)$ is typically a single linear layer, minimizing the CCE loss implicitly encourages the feature extractor $f_{\theta}(\cdot)$ to structure the latent space so that classes become linearly separable, thereby improving accuracy and generalization.

In a multimodal setting, the input data x consists of M distinct modalities x_1, x_2, \dots, x_M , derived from the same biflow. Each modality x_m is processed by a dedicated feature extractor $f_{\theta_m}(\cdot)$ —i.e., a *branch*—, resulting in a latent representation $z_m = f_{\theta_m}(x_m)$, for $m = 1, \dots, M$.

The resulting feature vectors are then combined through a fusion function F_{ϕ} modeling interactions across modalities and producing a joint (latent) representation $\bar{z} = F(z_1, z_2, \dots, z_M)$, which is passed to the classifier $f_{\phi}(\cdot)$. This enables an integrated decision-making process based on complementary modality-specific features.

Notably, *different modalities can be affected differently by network changes*, due to their heterogeneous nature. For example, during congestion, **IATs** (captured by **PSQ**) may increase, while raw bytes (from **PAY**) remain largely unchanged. Accordingly, we hypothesize that the *synergistic combination of these modalities through multimodal learning yields a more comprehensive and robust representation*, enhancing resilience to variations in network conditions.

3.2. Perturbation scenarios overview

To evaluate the robustness of traffic classifiers, we designed traffic perturbation scenarios that reflect realistic deployment conditions. While traffic classifiers are typically trained on traffic captured in controlled environments near the network edges, real-world monitoring points—such as middleboxes—may operate in contexts where the view of traffic is biased or distorted (e.g., middleboxes are mainly deployed in Tier2/ISP networks [41]) due to factors such as asymmetric routing, delayed observations, packet loss, or timing variations [42,43].

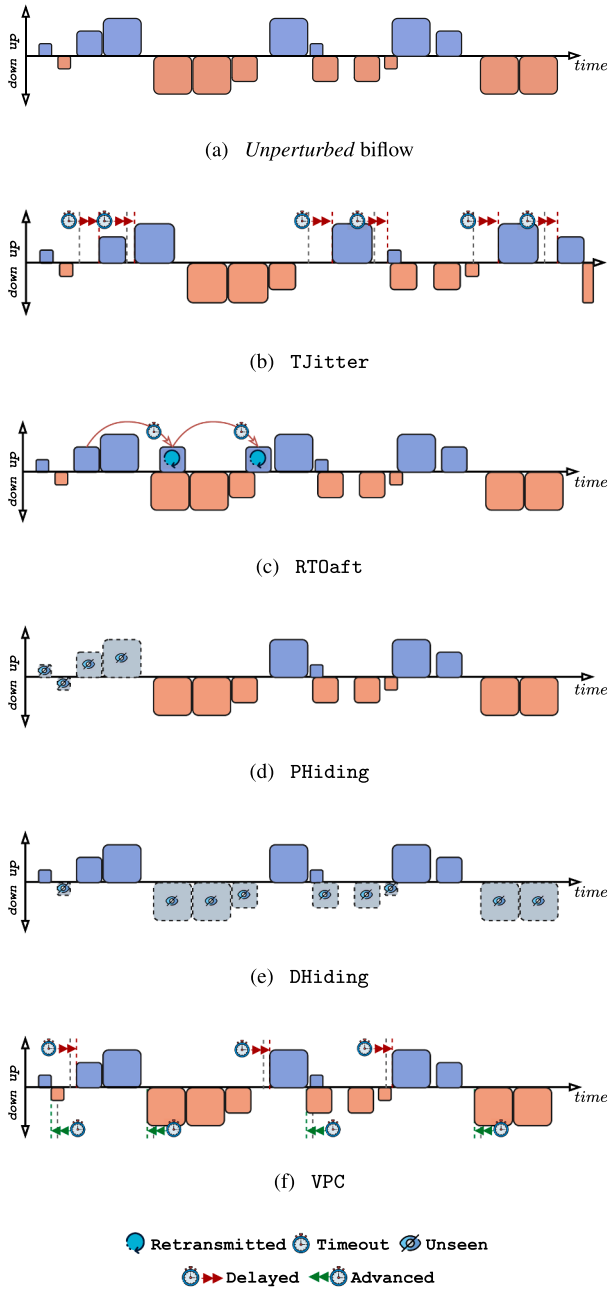


Fig. 1. Visual effects of **TJitter** (b), **RT0aft** (c), **PHiding** (d), **DHiding** (e), and **VPC** (f) on a biflow (a). The *x*-axis represents time, while the *y*-axis indicates the traffic direction (e.g., upstream or downstream). Each box corresponds to a packet, and its size reflects the actual payload length.

Specifically, we introduce a set of perturbations that selectively alter the packet sequence and/or the payload content observed at the **VP** to simulate these effects. Because our goal is to enable early biflow-level **TC** using two widely-adopted inputs—**PAY** and **PSQ**—we systematically explain each perturbation: its underlying motivation, how it reshapes the biflow structure, and which traffic inputs it influences (see Table 4).

Importantly, although these perturbations emulate real-world effects, they are applied synthetically to traffic originally captured in controlled environments. To support reproducibility, we include pseudocode for each perturbation in Appendix A. In the following, we describe the perturbations, whose effects are illustrated in Fig. 1.

Directional Time Jittering (TJitter) (Fig. 1b): since network conditions are inherently dynamic and subject to transient variations, packet timing—particularly **IAT**—can be strongly affected by network environmental changes.

To assess model robustness to timing fluctuations, **TJitter** selectively alters the **IAT** of packets of a given biflow based on their direction (i.e., upstream or downstream). Specifically, a small delta-time is added (*delaying*) to the **IAT** of packets traveling in a given direction, with these values sampled from direction-specific probability distributions, allowing for asymmetric timing distortions.

This transformation shares similarities with the jittering proposed in [44], but differs by explicitly modeling distinct delays per direction, better reflecting the real-world condition of asymmetric path routing. For simplicity, to model jitter in the **IAT** of the i th packet in a biflow, IAT_i , we sample two independent terms $\lambda_{a,i} \sim \mathcal{N}^+(\mu, \sigma_{\text{DIR}_i})$ and $\lambda_{b,i} \sim \mathcal{N}^+(\mu, \sigma_{\text{DIR}_i})$, with $\text{DIR}_i \in \{\text{upstream}, \text{downstream}\}$. The perturbed **IAT**, denoted as $\widehat{\text{IAT}}_i$, is computed as:

$$\widehat{\text{IAT}}_i = \begin{cases} \text{IAT}_i + |\lambda_{a,i} - \lambda_{b,i}|, & \text{if } \text{DIR}_i = \text{DIR}_{i-1} \\ \text{IAT}_i + |\lambda_{a,i} + \lambda_{b,i}|, & \text{if } \text{DIR}_i \neq \text{DIR}_{i-1} \end{cases} \quad (1)$$

which introduces jitter based on the direction of the current packet, DIR_i , and that of the preceding one, DIR_{i-1} . Notably, since **TJitter** affects only packet timing, it does not alter the packet content (i.e., the **PAY** input) or the arrival order of packets as observed from the **VP**.

Retransmission Timeout (RTO) (Fig. 1c) emulates the retransmission of packets within a biflow, where both the packets to be retransmitted and the number of retransmissions are randomly selected. Packet retransmissions are modeled through a **retransmission rate** r , with each retransmitted packet being re-sent after a preset **Timeout**.

In this context, we devised two realistic sub-scenarios to model different packet retransmission conditions: (i) **RT0aft**, where the **VP** observes both the original and retransmitted packets—i.e., the loss occurs *after* the **VP**; (ii) **RT0bef**, where the **VP** observes only retransmitted copies—i.e., the original packet was lost *before* the **VP**.

The proposed transformation is a generalization of the simplified version of **RT0bef** suggested by Xie et al. [26]. Notably, although retransmissions do not alter the content of packets, they can indirectly affect the construction of the **PAY** input since it is built by concatenating the payloads of the first packets until N_b bytes are reached, and the number of packets required depends on their payload sizes.

As a result, retransmissions can substantially affect **PAY**, particularly when they involve the earliest packets. This effect is pronounced in the **RT0bef** setting, where the **VP** misses the original transmissions. If the initial packet is lost before reaching the **VP** and only the retransmission is observed, its payload may not contribute to **PAY** as it would under lossless conditions—leading to a distorted view of the biflow. In addition, our **RTO** transformation—similarly to what is proposed in [26] and adopted in [30,35]—does not account for the potential impact on the **IAT** of the first response packet when the **RTO** involves a request packet. While such a retransmission could delay the server’s response and thus alter the observed **IAT**, we intentionally avoid modeling this effect to maintain consistency with prior work and prevent the introduction of speculative timing artifacts¹ Moreover, **RT0bef** shares some similarities to the packet loss applied by Yang et al. [32], where they randomly remove 10% to 50% packets from a given biflow. However, our proposal goes beyond the simple packet loss by also implementing a retransmission mechanism.

Packets Hiding (PHiding) (Fig. 1d) models the scenario where the **VP** misses a biflow’s initial packets and must classify it based on the later

¹ We do not explicitly model request/response semantics in a biflow and preserve all original packet **IATs** for simplicity.

ones (*late-start classification*), for example, due to a start-up delay in the monitoring system or classifier.

The perturbation affects both **PSQ** and **PAY**, although with different intensity. Intuitively, **PSQ** may be less sensitive to missing initial packets because the classifier still sees a slightly shifted packet sequence that remains close to the original. **PAY**, instead, could be heavily affected, as it relies on the biflow's first bytes. The problem is amplified in **TLS** traffic, where initial packets carry plaintext metadata (e.g., *Server Name Indication*, *Cipher Suites*) while later packets are usually encrypted [16, 45].

Similar to **PHiding**, Yang et al. [32] apply packet loss, but our approach discards only the first packets of a biflow, starting from a random position.

Direction Hiding (DHiding) (Fig. 1e) models the scenario where the forward and backward paths between Client and Server are asymmetric, so only one direction crosses the observation point. This condition becomes more likely as this point is placed deeper within the network, i.e., farther from the endpoints. We also consider a variant, **Direction Hiding with Packet Loss (DHiding_{loss})**, which assumes that even packets in the visible direction may be missing. Each packet is dropped independently with probability p , accounting for situations where packets in the same direction may follow different paths, possibly due to routing asymmetries or load balancing.

Assuming bidirectional communication where both endpoints generate traffic, both versions primarily affect **PSQ** by altering the observed packet sequence, making it unidirectional or even incomplete. Conversely, their impact on **PAY** mainly depends on the direction and the size of the earliest packets of the biflow, which are typically used to extract **PAY**.

Vantage Point Centralization (VPC) (Fig. 1f) models the case where a **VP** is placed far from the endpoints but still observes traffic in both directions. In such conditions, packets exchanged exhibit different timing characteristics due to asymmetric delays, primarily influenced by the **RTT**. We assume that the **RTT** between the Client and the Server can be decomposed into two segments [46]: *from-client-to-VP* and *from-VP-to-server*. We introduce a parameter $\alpha \in [0, 1]$ to represent the proportion of **RTT** attributed to the Client-to-VP path:

$$RTT = RTT_{C \rightarrow V} + RTT_{S \rightarrow V} \quad (2)$$

where $RTT_{C \rightarrow V} = \alpha \cdot RTT$ and $RTT_{S \rightarrow V} = (1 - \alpha) \cdot RTT$, and α represents the relative position of the **VP** along the Client-to-Server path, with $RTT_{C \rightarrow V}$ (resp. $RTT_{S \rightarrow V}$) denoting the **RTT** between Client and **VP** (resp. Server and **VP**).

Based on this decomposition, the **IAT** between consecutive packets sent in opposite directions depends on the total **RTT** and reflects the asymmetry between upstream and downstream paths. Specifically, consider two consecutive packets with opposite directions: a packet p_c sent by the Client and a packet p_s sent by the Server. Whether p_c precedes p_s or p_s precedes p_c , the **IATs** of the later-arriving packet can be computed as:

$$IAT_{p_s} = T_{p_c} + T_{p_s}^k + T_{p_s} + (1 - \alpha) \cdot RTT \quad (3)$$

$$IAT_{p_c} = T_{p_s} + T_{p_c}^k + T_{p_c} + \alpha \cdot RTT \quad (4)$$

Here, IAT_p is the **IAT** of the packet p , T_p^k denotes the thinking time at the sender before transmitting packet p , and T_p denotes the transmission time for the packet p , which typically depends on its **PS**.

In our experimental setup, traffic is usually captured near the client, so we assume $\alpha \approx 0$. To model a **VP** farther from the endpoints, we vary α and recompute the **IATs** to introduce the desired timing asymmetry. Notably, changes in the **VP** position do not affect the **IATs** between packets sent in the same direction, as the **RTT** components cancel out in the difference.

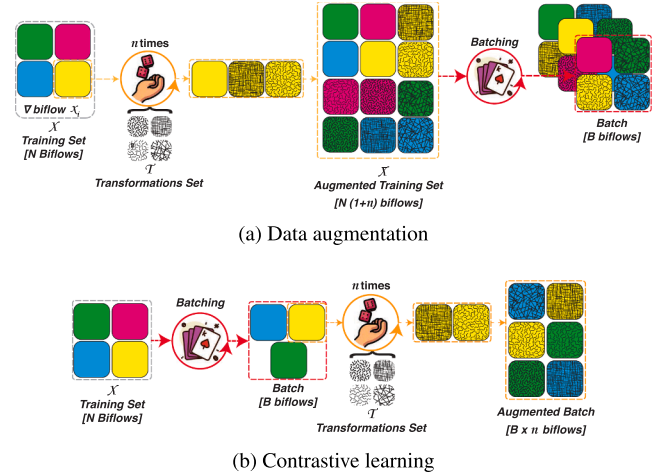


Fig. 2. Transforming procedure in **Aug** (a) and **CL** (b). In **Aug**, transformations are applied before training to create a fixed augmented training set. Batches are then formed and may include both original and transformed samples. In **CL**, transformations are applied “on the fly” during training. Each batch contains only transformed samples. N is the number of original samples. B is the batch size. T is the number of independent transformations per sample.

We further defined the **Vantage Point Centralization with Packet Loss (VPC_{loss})** variant with each packet independently dropped with probability p . This models the increased likelihood of partial visibility when the **VP** is located deeper in the network, due to routing asymmetries, and load balancing.

Regarding input perturbation, **VPC** affects only **PSQ**, as it alters **IATs**, without modifying packet order or payload content, whereas **VPC_{loss}** may also impact **PAY** when early packets are lost. These scenarios help assess whether a model—typically trained on traffic captured near the Client—remains effective when deployed in more central, challenging network positions.

3.3. Perturbation-aware learning approaches

Herein, we describe the two strategies for integrating perturbation-based transformations into the training process: *supervised data augmentation* and *contrastive learning*. Both apply controlled transformations to original samples, but with different objectives. The former directly trains on transformed samples labeled as their corresponding originals to improve generalization [47], while the latter exploits transformations to create semantically-similar pairs, enabling the model to learn robust representations via instance discrimination [48].

Supervised Data Augmentation (Aug): is typically used to enrich supervised learning by increasing sample diversity or mitigating class imbalance. Since **DL** models require large datasets to generalize, **Aug** is particularly beneficial when labeled samples are limited. Hence, we applied perturbations to training samples to generate realistic variants and used within standard supervised learning.

The underlying idea is that, by exposing the model to a wider range of plausible input variations, **Aug** helps it learn invariant representations, i.e., feature embeddings $z = f_{\theta}(x)$ that remain consistent across different perturbations of inputs belonging to the same class. Hence, since the classifier $f_{\phi}(\cdot)$ is linear, the model is implicitly encouraged to cluster perturbed samples of the same class in the latent space, thereby improving inter-class separation generalization. Unlike augmentation techniques used in domains such as computer vision (e.g., rotations or flipping), we *do not apply generic, domain-agnostic transformations*. We instead apply domain-specific perturbations tailored to network biflows, using the realistic transformations defined in Section 3.2 to generate semantically plausible traffic variants.

Formally, given a set of training samples $\mathcal{X} \equiv \{x_i\}_{i=1}^N$, Aug creates n perturbed samples for each x_i by independently applying n transformations sampled from a predefined set \mathcal{T} . We adopt a *pre-augmented* strategy (see Fig. 2a) where all perturbations are applied before training—i.e., prior to mini-batches construction—thus decoupling transformation sampling from the training loop.

For each x_i , we randomly select n transformations $\tau_1(\cdot), \tau_2(\cdot), \dots, \tau_n(\cdot) \in \mathcal{T}$ and apply them separately to obtain the augmented versions, resulting in a *transformed* training set:

$$\hat{\mathcal{X}} \equiv \{\hat{x}_k = \tau_j(x_i) | x_i \in \mathcal{X}, j \in \{1, \dots, n\}\} \quad (5)$$

with $|\hat{\mathcal{X}}| = n \cdot N$. The final *augmented* training set is obtained by combining both the original and augmented samples:

$$\bar{\mathcal{X}} = \mathcal{X} \cup \hat{\mathcal{X}} \quad (6)$$

This ensures that the training set remains representative of the original distribution while also promoting generalization through increased sample diversity.

Contrastive Learning (CL) is a representation learning approach that explicitly regularizes the latent space by teaching models to recognize similarity and dissimilarity among samples. Accordingly, it optimizes a feature extractor by comparing the projections of samples in the latent space using a contrastive loss to bring similar samples closer and push dissimilar ones apart. This is achieved by selecting a reference (ref. to as *anchor*) representation and comparing it to a *positive* (viz. similar) and one or more *negative* (viz. dissimilar) samples. CL methods can be either self-supervised or supervised, depending on how similarity is defined and whether labels are used.

Self-supervised Contrastive Learning (SimCon) [48] can be seen as a branch of unsupervised learning, as it does not involve any manual labeling [49]. It assumes each sample is only similar to its own augmented versions and dissimilar to all others. Unlike Aug, in CL the augmentations are applied independently at each training step, i.e., new pairs are generated for every batch (see Fig. 2b), exposing the model to constantly varying transformations and improving robustness and generalization.

Formally, given a mini-batch $B = \{x_i\}_{i=1}^B \sim \mathcal{X}$, for each $x_i \in B$, two transformed versions $\tilde{x}_i = \tau_1(x_i)$ and $\tilde{x}_j = \tau_2(x_i)$ are obtained by applying two distinctive transformations $\tau_1(\cdot), \tau_2(\cdot) \sim \mathcal{T}$. This results in an “augmented” batch $\tilde{B} = \{\tilde{x}_1, \tilde{x}_2, \dots, \tilde{x}_{2B}\}$ containing two transformed versions of each original input sample x_i . Therefore, for each anchor \tilde{x}_i , its corresponding counterpart \tilde{x}_j (i.e., the one originated from the same original sample), is treated as a positive, while the remaining $2 \cdot (B - 1)$ samples $\tilde{x}_k \in \tilde{B} \setminus \{\tilde{x}_i, \tilde{x}_j\}$ are considered as negative.

The learning process penalizes cases where the latent representation of positive pairs ($f_\theta(\tilde{x}_i), f_\theta(\tilde{x}_j)$) are far apart, and those of negative pairs ($f_\theta(\tilde{x}_i), f_\theta(\tilde{x}_k)$) are too close. This is typically performed by minimizing the *Normalized Temperature-scaled Cross-Entropy (NT-Xent)* loss function [48] (ref. to as $\mathcal{L}_{\text{SimCon}}$ in the following), which aims to maximize the mutual information between the anchor and positive samples. Formally, this loss function is defined as:

$$\mathcal{L}_{\text{SimCon}}(\tilde{x}_i) = -\ln \frac{e^{\text{sim}(f_\theta(\tilde{x}_i), f_\theta(\tilde{x}_j))/\beta}}{\sum_{k=1}^{2N} \mathbf{1}_{[k \neq i]} e^{\text{sim}(f_\theta(\tilde{x}_i), f_\theta(\tilde{x}_k))/\beta}} \quad (7)$$

where $\text{sim}(\cdot)$ is a similarity function (e.g., cosine similarity) and $\beta \in \mathbf{R}^+$ is a temperature parameter.

It is worth noting that the original framework includes a *projection head* that maps the latent representations into a “contrastive space” where $\mathcal{L}_{\text{SimCon}}$ is applied. In our setting, $\mathcal{L}_{\text{SimCon}}$ is directly applied to the latent representations as CL approaches are fully integrated into the supervised pipeline, and the resulting representations are directly used for the classification task.

Supervised Contrastive Learning (SupCon) [50] aims to overcome the limitations of SimCon by incorporating class labels directly into the learning process. In SimCon, samples from the same class as the anchor

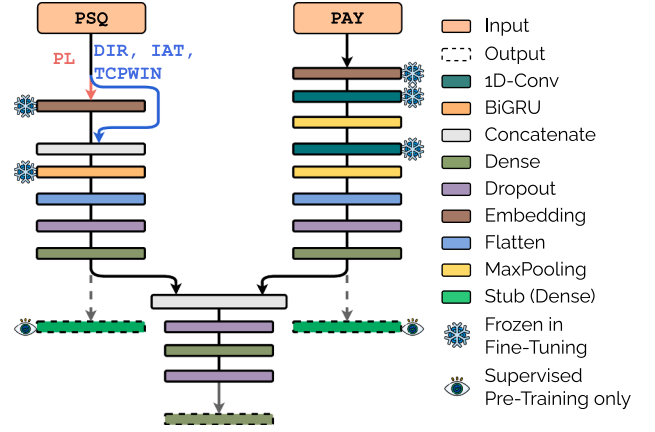


Fig. 3. MIMETIC-ENHANCED classifier. Colors indicate the layer type.

are treated as negatives, thus preventing the model from learning the relationship between samples of the same class. Thus, SupCon refines the latent space by ensuring that samples of the same class are grouped together while samples from different classes are well-separated.

In this case, the learning process is driven by a *SupCon* loss (ref. to as $\mathcal{L}_{\text{SupCon}}$ in the following), which is an extension of the classical $\mathcal{L}_{\text{SimCon}}$. Formally, this loss is defined as:

$$\mathcal{L}_{\text{SupCon}}(\tilde{x}_i) = -\frac{1}{|\mathcal{P}(i)|} \sum_{\tilde{x}_j \in \mathcal{P}(i)} \ln \frac{e^{\text{sim}(f_\theta(\tilde{x}_i), f_\theta(\tilde{x}_j))/\beta}}{\sum_{k=1}^{2N} \mathbf{1}_{[k \neq i]} e^{\text{sim}(f_\theta(\tilde{x}_i), f_\theta(\tilde{x}_k))/\beta}} \quad (8)$$

where $\mathcal{P}(i) \equiv \{\tilde{x}_k \in \tilde{B} \setminus \{\tilde{x}_i\} : y(\tilde{x}_k) = y(\tilde{x}_i)\}$ —with $y(\cdot)$ denoting the actual label—, $\text{sim}(\cdot)$ is a similarity function (e.g., cosine similarity) and $\beta \in \mathbf{R}^+$ is a temperature parameter. In both SimCon and SupCon, the final loss is computed symmetrically across all positive pairs (\tilde{x}_i, \tilde{x}_j) and (\tilde{x}_j, \tilde{x}_i) within the batch.

4. Experimental setup

In this section, we describe the DL architecture (Section 4.1) and the proposed training procedures (Section 4.2). Then, Section 4.3 reports the datasets employed for the experiments, whereas Section 4.4 outlines the evaluation procedure and the metrics.

4.1. Reference multi-modal architecture

We use MIMETIC-ENHANCED [16] as our baseline model. The architecture and training methodology are described below.

Architecture Overview: Fig. 3 depicts the architecture of MIMETIC-ENHANCED, a multimodal DL model that combines two distinct traffic views (i.e., $M=2$). Each modality is processed by a distinct branch, optimized to extract salient features from its corresponding traffic input. The first branch processes a sequence of $N_h=4$ header fields extracted from the first $N_p=10$ packets (viz. PSQ). These fields include the PL, DIR—represented as $\{0=\text{upstream}, 1=\text{downstream}\}$, TCP Window Size (TCPWIN), which is set to zero for UDP packets, and IAT. The second branch handles the first $N_b=576$ bytes of payload data (viz. PAY), treated as a flat byte sequence. Each input is first passed through a *trainable embedding layer* that maps raw input elements into a vector of dimension $e=10$. This results in an embedding matrix $E \in \mathbf{R}^{N \times e}$, where N is equal to N_p and N_b for PSQ and PAY, respectively. Notably, for PSQ, only the sequence of PL is processed with the embedding layer. The resulting vectors are concatenated with the remaining header fields to form the full input to the PSQ branch. In other terms, after the embedding layer, the PSQ data has a resulting size $N_p \times (e + N_h - 1)$.

The PSQ branch comprises a Bidirectional Gated Recurrent Unit layer (BiGRU, with 64 units—returning sequences, followed by a 256-unit dense layer. The PAY branch is composed of two 1D convolutional

Table 5

Overview of the training variants used to integrate **Aug** and **CL** (i.e., **SimCon** or **SupCon**) into the two-phase training procedure of MIMETIC-ENHANCED. Each configuration applies perturbation-based strategies—either individually or in combination—across the **PT** and **FT** phases. The baseline configuration is highlighted in blue.

	Pre-Training					Fine-Tuning				
	M	S	O	T	\mathcal{L}	M	S	O	T	\mathcal{L}
Base	⊗	⊗	⊗	⊗	\mathcal{L}_{CCE}	⊗	⊗	⊗	⊗	\mathcal{L}_{CCE}
Aug _{PT}	⊗	⊗	⊗	⊗	\mathcal{L}_{CCE}	⊗	⊗	⊗	⊗	\mathcal{L}_{CCE}
SimCon _{PT}	⊗	⊗	⊗	⊗	\mathcal{L}_{SimCon}	⊗	⊗	⊗	⊗	\mathcal{L}_{CCE}
SupCon _{PT}	⊗	⊗	⊗	⊗	\mathcal{L}_{SupCon}	⊗	⊗	⊗	⊗	\mathcal{L}_{CCE}
Aug _{FT}	⊗	⊗	⊗	⊗	\mathcal{L}_{CCE}	⊗	⊗	⊗	⊗	\mathcal{L}_{CCE}
Aug _{PT} → Aug _{FT}	⊗	⊗	⊗	⊗	\mathcal{L}_{CCE}	⊗	⊗	⊗	⊗	\mathcal{L}_{CCE}
SimCon _{PT} → Aug _{FT}	⊗	⊗	⊗	⊗	\mathcal{L}_{SimCon}	⊗	⊗	⊗	⊗	\mathcal{L}_{CCE}
SupCon _{PT} → Aug _{FT}	⊗	⊗	⊗	⊗	\mathcal{L}_{SupCon}	⊗	⊗	⊗	⊗	\mathcal{L}_{CCE}

Architectural Layers (⊗): Per-modality Branches (M), Final Layers (S), Trained (⊗), Not Trained (⊗).—Objective Loss (\mathcal{L}): CCE (\mathcal{L}_{CCE}), NT-Xent/SimCon (\mathcal{L}_{SimCon}), SupCon (\mathcal{L}_{SupCon}).—Training Data (⊗): Original (O), Transformed (T), Used (⊗), Not Used (⊗).

layers (1D-Conv, with 16 and 32 filters, kernel size of 25, and stride 1, each followed by a max-pooling layer with window size of 3 and unit stride, and a dense layer with 256 neurons. To model cross-modality correlations, outputs from both branches are concatenated and passed through a shared dense layer with 128 neurons, followed by a softmax-activated dense layer for classification. All layers use ReLU activations. To regularize training and reduce overfitting, a dropout—with a rate of 0.2—is applied after each dense layer and at the end of both modality branches.

Details of the training process: MIMETIC-ENHANCED is trained using a *two-stage procedure*. In the **Pre-Training (PT)** stage, each single-modality branch is independently trained using a *softmax stub* to guide the branch-specific learning. This is followed by a **Fine-Tuning (FT)** stage, where the entire architecture is optimized after freezing the embedding, convolutional, and BiGRU layers (*). Both training phases minimize a **CCE** loss using the SGD optimizer (batch size of 64), combined with a custom *ReduceLROnPlateau* scheduler that monitors the validation loss and decreases the learning rate by a factor of 3 after 20 consecutive epochs without improvement.²

4.2. Advanced training for multi-modal traffic classification

We integrated **Aug** and **CL** into the “supervised” training pipeline of MIMETIC-ENHANCED to assess their impact on both **TC** effectiveness and robustness. As summarized in **Table 5**, we devised several training variants by applying **Aug** and **CL** either individually or jointly across the different phases of the training process (cf. **Section 4.1**).

For each variant, the table details the setup adopted during the **PT** and **FT** phases, in terms of:

1. the architectural layers involved (⊗)—i.e., per-modality branches (M) or final layers (S);
2. the training data used (⊗)—i.e., original (O) or transformed (T); and the objective loss function used to optimize the model—i.e., \mathcal{L}_{CCE} , \mathcal{L}_{SimCon} , or \mathcal{L}_{SupCon} .

The baseline configuration (ref. to as **Base** in the following) relies on original data (O) and the **CCE** in both **PT** and **FT**. We then considered configurations where transformed data (T) is used only during **PT** (i.e., Aug_{PT}, SimCon_{PT}, and SupCon_{PT}) to investigate whether training only the

early layers (M) improves representation learning. The classifier is always trained solely on original data, so it is never exposed to transformed inputs and is not explicitly trained to recognize them. We also assess the Aug_{FT} configuration where **Aug** is applied only during **FT**. In this setup, **PT** uses only original samples, while **FT** is performed on the augmented training set (O + T). This approach isolates the effect of exposing only the model’s final layers (S) to transformed data, allowing us to assess whether this selective exposure improves generalization by refining the classifier’s decision boundaries without altering feature extraction. We do not apply SimCon or SupCon during **FT**, since this phase specifically aims to refine the classifier using a supervised objective, whereas **CL** focuses on shaping the latent space rather than directly optimizing class separation.

In addition, we explore combined configurations in which **Aug** or **CL** is applied in both **PT** and **FT**. More in detail, Aug_{PT} → Aug_{FT} uses original and transformed samples (O + T) in both training phases, enabling us to assess the impact of consistent exposure to increased data variability throughout the entire learning process. Conversely, SimCon_{PT} → Aug_{FT} and SupCon_{PT} → Aug_{FT} adopt a hybrid strategy: the model is pre-trained with **CL** on transformed data (T) and then fine-tuned on an augmented dataset (O + T) using standard cross-entropy.

These setups aim to evaluate whether combining robust representation learning—enabled by **CL**—with augmented **FT** improves generalization by leveraging both latent space regularization and classifier adaptation to data variability. For contrastive variants, no output stub is used during **PT**, and the loss is applied directly to the latent representations. Additionally, before the final supervised **FT**, an intermediate contrastive **FT** step updates only the shared layers and the final dense layer (S), while all remaining layers (M) remain frozen.

In our setup, we consistently generated $n = 2$ transformed versions of each sample during all training stages involving **Aug** or **CL**. In both (training) strategies, perturbations are sampled independently for every sample: both the transformation type (e.g., RTOaf t) and its hyperparameters (e.g., retransmission probability, $r\%$) are drawn from predefined distributions. This prevents the model from overfitting to fixed, deterministic transformations and encourages the learning of invariant features across diverse and realistic variations. Conversely, during robustness evaluation, the same transformation is applied deterministically to all test samples under fixed hyperparameter settings. This procedure is repeated across multiple configurations to enable controlled sensitivity analysis and to ensure reproducible results.

4.3. Datasets description

In this work, we employ three publicly available datasets, namely VPN-16 [19], MIRAGE-19 [8], and MIRAGE-24 [20].

VPN-16 [19] was collected at the Canadian Institute for Cybersecurity and provided in raw PCAP format. It includes human-generated traffic spanning various types, with app-related information collected from both regular and VPN-encapsulated sessions. The dataset comprises 6 traffic types and 15 apps, providing the ground truth at the trace level for the encapsulation type, service type, and app.

Given that VPN encapsulation can alter apps’ traffic patterns, we consider both the app and the encapsulation type as joint labels. **Fig. 4a** depicts the biflow distribution across all (app, encapsulation) pairs. Since some classes contain only a few dozen biflows, we retained only those with at least 100 biflows (17 classes, highlighted in yellow) to ensure sufficient representation.

MIRAGE-19 [8] contains per-flow mobile traffic logs of 40 Android apps collected at the ARCLAB laboratory of the University of Napoli Federico II between May 2017 and May 2019. The data was collected by employing 3 rooted Android devices³, which were used by ≈ 300 volunteers,

² The initial learning rate is set to 0.1, and training stops early if the learning rate falls below 0.0001. To avoid overfitting, we use a validation-based early-stopping technique, where we set the patience to 20 epochs without setting the min_delta parameter.

³ These devices were running CyanogenMod v13.0 custom firmware, equivalent to Android 6.0.1

where N is the total number of samples. **SIL** ranges in $[-1, 1]$, with 1 being the best value—indicating high compactness and clear separation—and -1 the worst, corresponding to samples assigned to the wrong cluster.

Model Calibration: we assess model calibration—i.e., how closely predicted confidence matches actual accuracy—using reliability diagrams [53] and the **Expected Calibration Error (ECE)** [54]. A model is *perfectly calibrated* if, whenever it predicts a confidence level \hat{p} , the true probability of correctness matches that confidence—that is, $\Pr\{\hat{y} = y \mid \hat{p}\} = \hat{p}$, where y and \hat{y} are the true and predicted labels, respectively.

For each test sample n , we denote its true label as $y_{(n)}$, the predicted label as $\hat{y}_{(n)} \triangleq \arg \max_{1, \dots, |C|} p_i(n)$, and the predicted confidence as $\hat{p}_{(n)} \triangleq \max_{1, \dots, |C|} p_i(n)$.

To estimate calibration, we partition predictions into M bins, evenly spaced across the confidence range from $[1/|C|, 1]$. For bin B_m , containing samples with confidence \hat{p} in interval $I_m \triangleq (\frac{m-1}{M}, \frac{m}{M})$, we compute the “expected” accuracy:

$$\text{acc}(B_m) = \frac{1}{|B_m|} \sum_{n \in B_m} 1(\hat{y}_{(n)} = y_{(n)}) \quad (13)$$

and the center-bin confidence $\hat{p}_m \triangleq \frac{m}{M} - \frac{1}{2M}$. A reliability diagram plots these two quantities for each bin and deviations from the diagonal indicate *over-* or *under-confidence*.

We summarize (mis)calibration using the **ECE**, defined as $E_{\hat{p}}\{|P\{\hat{\ell} = \ell \mid \hat{p}\} - \hat{p}|\}$, and approximated as:

$$\text{ECE} \approx \sum_{m=1}^M \frac{|B_m|}{N} |\text{acc}(B_m) - \text{conf}(B_m)| \quad (14)$$

where N is the total number of test samples and $\text{conf}(B_m) = \frac{1}{|B_m|} \sum_{n \in B_m} \hat{p}_{(n)}$ represents the averaged confidence within m th bin. This metric succinctly captures how far the model’s confidence diverges from reality.

Out-of-Distribution (OOD) Detection: to evaluate the model’s ability to identify samples that do not belong to any of the known classes (i.e., seen during the training), we assess its performance in detecting **OOD** traffic. This capability is fundamental in open-set recognition, where the model is required to appropriately classify inputs from previously unseen classes as unknown [55], and is particularly relevant in open-world settings such as network traffic analysis, where novel or anomalous patterns may emerge.

In our case, the (anomaly) *scoring function* corresponds to the negated predicted confidence associated with the inferred class—i.e., $s(x) = -\hat{p}(x)$ where $\hat{p}(x) = \max_{c \in C} p_c(x)$. This score is justified since higher **OOD** scores (closer to zero) correspond to a lower confidence of the classifier for the known classes.

To measure **OOD** capabilities, we use the **Area Under the Receiver Operating Characteristic Curve (AUROC)** to discriminate between positive (viz. out-of-distribution) and negative (viz. in-distribution) samples. From a geometric standpoint, **AUROC** is defined as the area of the ROC curve in the TPR-FPR space:

$$\text{AUROC} = \int_0^1 \text{TPR}(\text{FPR}) d\text{FPR} \quad (15)$$

where $\text{TPR}(\text{FPR})$ denotes the **True Positive Rate (TPR)** as a function of the **False Positive Rate (FPR)**. The **AUROC** span over $[0, 1]$, with 1 indicating the perfect separation between **In-Distribution (ID)** and **OOD** samples, while 0.5 corresponds to random guessing.

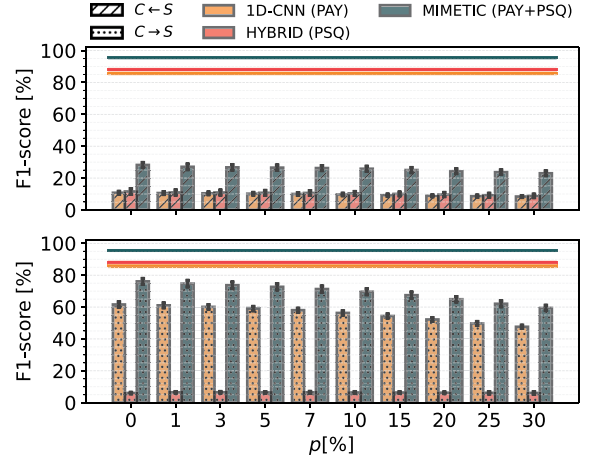


Fig. 5. Comparison of HYBRID, 1D-CNN, and MIMETIC-ENHANCED under **DHiding** ($p=0\%$) and **DHiding_{loss}** ($p>0\%$) with packet-loss rate $p \in \{0, 1, 3, 5, 7, 10, 15, 20, 25, 30\}$. Each model is identified by its traffic type (i.e., PAY or PSQ), with performance on unperturbed samples shown by a solid line. Results are distinct for Client-to-Server ($C \rightarrow S$) and Server-to-Client ($C \leftarrow S$) directions of the biflow observed at the **VP**. Values refer to **MIRAGE-24** and are reported as *mean* \pm *std.dev.* over 10 folds.

5. Experimental results

In this section, we report our experimental evaluation⁷ We start by comparing single-modal and multi-modal baseline architectures across unperturbed and perturbed scenarios (Section 5.1). Then, we evaluate the baseline multi-modal architecture by assessing: (i) its ability to generalize to perturbed inputs (Section 5.2); (ii) whether our training strategies improve generalization on clean (unperturbed) data (Section 5.3) and increase robustness to perturbations (Section 5.4). Next, we evaluate model calibration—how well predicted confidence matches actual accuracy (Section 5.5). We assess performance on **OOD** samples (Section 5.6) and analyze the impact of limited labeled data (Section 5.7). We further investigate the framework’s generalizability across distinct model architectures (Section 5.8) and dissect the impact of individual perturbation types on overall **TC** performance (Section 5.9). Finally, we assess the framework’s deployability (Section 5.10), offering insights into computational complexity and real-world feasibility Fig. 5.

5.1. Comparing the effectiveness and robustness of single- and multi-modal baselines

We first investigate the robustness of multimodal architectures compared to state-of-the-art single-modal models under limited traffic visibility. As a case study, we compare MIMETIC-ENHANCED with *two* representative *single-modal architectures*:

1. a 1D-CNN [11] fed with **PAY** input, and
2. HYBRID [13], a hybrid architecture composed by a 2D-CNN followed by a RNN and fed with **PSQ** input.

We focus on one representative perturbation from Section 3.2, noting that results are consistent across other cases but omitted for brevity.

Section 5.1 shows the **F1-score** for **DHidings** scenarios when the **VP** is traversed solely by either the Client-to-Server ($C \rightarrow S$) or Server-to-Client ($C \leftarrow S$) path, across increasing packet-loss rates p . The results highlight *two main advantages* of MIMETIC-ENHANCED: (i) it outperforms

⁷ Implementation uses PyTorch, numpy, and pandas, with scikit-learn for metrics (custom **ECE**) and the official **SupCon** code for contrastive losses [50]. Experiments ran on an Ubuntu 22.04 server with Intel i9-10900X, 256GB RAM, and 2xNVIDIA RTX A5000 (24GB) GPUs.

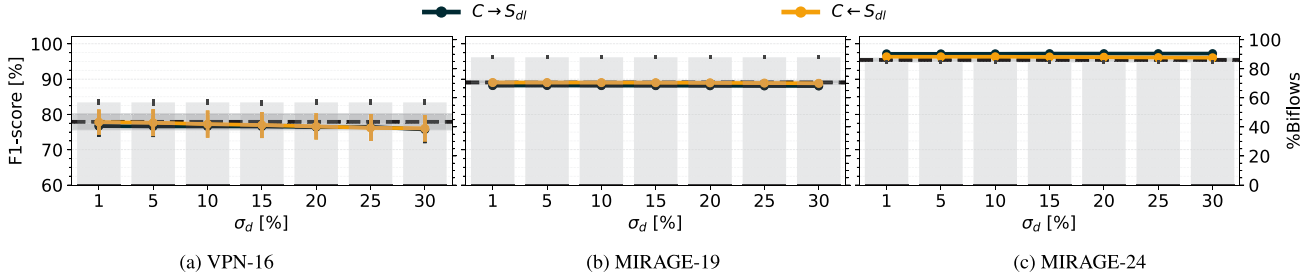


Fig. 6. Performance evaluation of Base under **TJitter** as the standard deviation σ_d varies in $\{1, 5, 10, 15, 20, 25, 30\}\%$, with σ_d expressed as a percentage of the maximum **IAT** per dataset. Results are reported separately for scenarios where a delay is added only in the Client-to-Server ($C \rightarrow S_{dl}$) or Server-to-Client ($C \leftarrow S_{dl}$) direction. The dashed-black line reports the performance on the unperturbed test samples. Values are reported as $mean \pm std.dev.$ over 10 folds. Gray bars show the share of test samples on which performance has been computed.

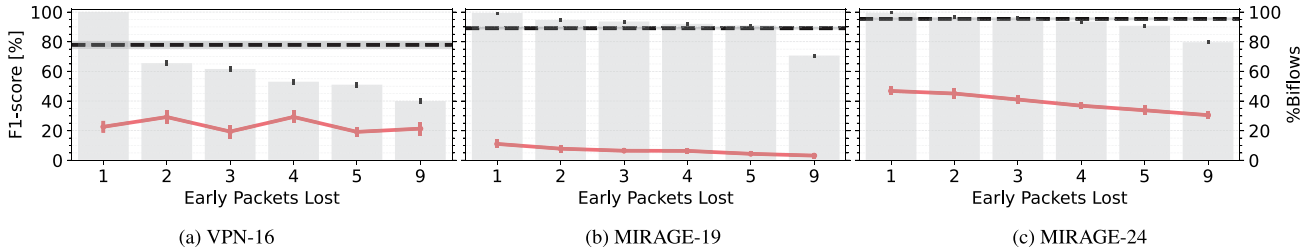


Fig. 7. Performance evaluation of Base under **PHiding** as the number of lost initial packets in the biflow varies in $\{1, 2, 3, 4, 5, 9\}$. The dashed-black line reports the performance on the unperturbed test samples. Values are reported as $mean \pm std.dev.$ over 10 folds. Gray bars show the share of test samples on which performance has been computed.

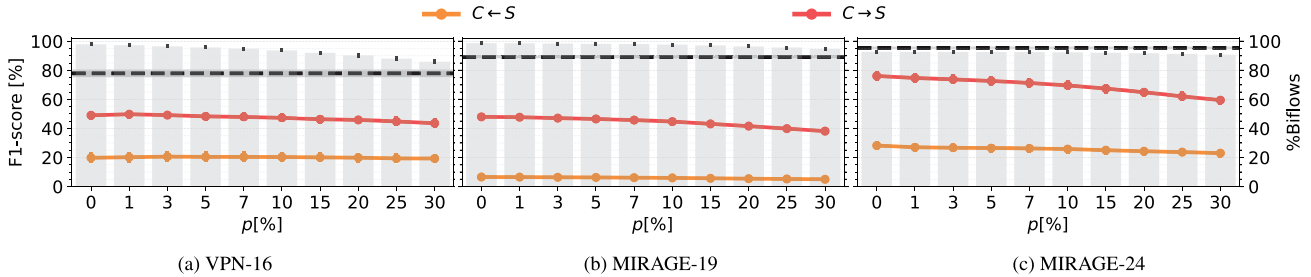


Fig. 8. Performance evaluation of Base under **DHiding** ($p=0\%$) and **DHiding_{loss}** ($p>0\%$) with packet-loss rate $p \in \{0, 1, 3, 5, 7, 10, 15, 20, 25, 30\}\%$. Results are reported separately for scenarios where only the Client-to-Server ($C \rightarrow S$) or Server-to-Client ($C \leftarrow S$) direction of the biflow is visible at the **VP**. The dashed-black line reports the performance on the unperturbed test samples. Values are reported as $mean \pm std.dev.$ over 10 folds. Gray bars show the share of test samples on which performance has been computed.

single-modal models on unperturbed data—up to a $\approx 10\%$ F1-score gain over 1D-CNN—demonstrating the value of combining complementary inputs like **PAY** and **PSQ**; (ii) it maintains greater robustness under perturbations, with smaller performance drops as packet loss (p) increases. Among single-modal models, 1D-CNN and HYBRID behave similarly in the $C \leftarrow S$ scenario but diverge when $C \rightarrow S$, where 1D-CNN shows better resilience and a degradation trend closer to MIMETIC-ENHANCED.

Takeaway:

Multimodal architectures outperform single-modal models by delivering higher classification accuracy in both ideal and perturbed settings, providing enhanced robustness when conditions degrade.

5.2. Baseline sensitivity to perturbations

Herein, we assess the generalization ability of the baseline model (viz. Base)—trained on “clean” data—when exposed to perturbed test samples. To this end, in **Figs. 6–10**, we evaluated its performance—in terms of F1-score—under each of the perturbed scenarios described

in **Section 3.2**. For each scenario, we defined a set of parameter values and applied the corresponding variants uniformly to all test samples. We also compare the resulting performance against that obtained on the unperturbed samples (highlighted with a dashed-black line). This analysis is useful for evaluating the sensitivity of the model to such controlled perturbations.

At a high level, we note that the model performance varies greatly depending on the type of perturbation and is almost consistent across all the datasets. On the one hand, perturbations that solely affect packet timing—such as **TJitter** (**Fig. 6**) and **VPC** (**Fig. 10**)—have a limited impact on classification performance, as they do not alter the sequence or structure of the biflows. Under **TJitter**, performance degradation remains limited as σ_d increases—e.g., up to -1% on MIRAGE-19. Interestingly, on MIRAGE-24, these perturbations lead to an improvement of up to $+1.8\%$, particularly when the delay is introduced in the forward ($C \rightarrow S_{dl}$) direction. This suggests the robustness of the model against timing noise and that it likely leverages alternative, more stable traffic features. Similar findings are observed for **VPC** on MIRAGE-19 and MIRAGE-24 (**Fig. 10b** and **c**). Notably, performance remains stable across the entire variability range $\alpha \in [0, 1]$, which represents the fraction of **RTT** attributed to the Client-to-VP path. In contrast, the performance drop is more pronounced on VPN-16, indicating that the

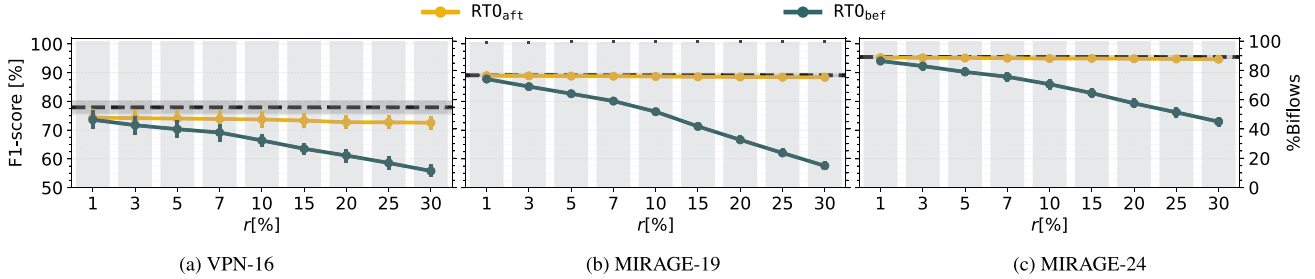


Fig. 9. Performance evaluation of Base under $RT0_{aft}$ and $RT0_{bef}$ with retransmission rate $r \in \{1, 5, 10, 20, 30\}\%$ and timeout $T_{out} = 1s$. The dashed-black line reports the performance on the unperturbed test samples. Values are reported as $mean \pm std.dev.$ over 10 folds. Gray bars show the share of test samples on which performance has been computed.

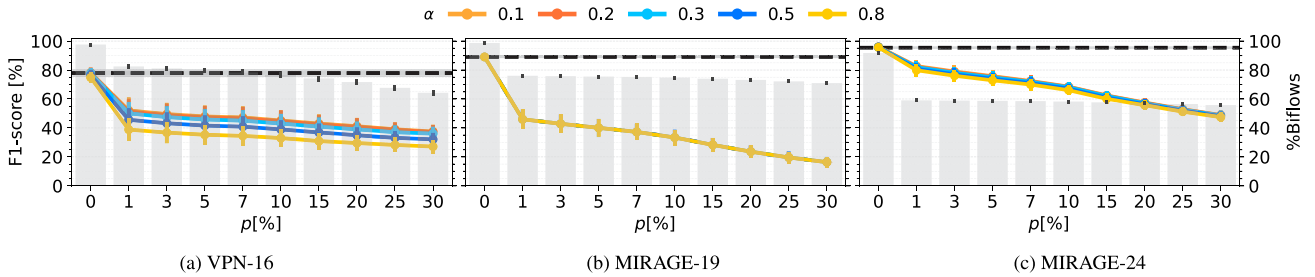


Fig. 10. Performance evaluation of Base under VPC ($p=0\%$) and VPC_{loss} ($p>0\%$) with the packet-loss rate $p \in \{0, 1, 3, 5, 7, 10, 15, 20, 25, 30\}\%$ and $\alpha \in \{0.1, 0.2, 0.3, 0.5, 0.8\}$. The dashed-black line reports the performance on the unperturbed test samples. Values are reported as $mean \pm std.dev.$ over 10 folds. Gray bars show the share of test samples on which performance has been computed.

model may be more sensitive to variations in the relative position of the VP w.r.t. the endpoints on this dataset.

On the other hand, perturbations affecting biflow visibility through packet reordering, hiding, or duplication—e.g., PHiding, DHiding, $RT0_{aft}$, and VPC_{loss} —result in more pronounced performance degradation, indicating that the model is more sensitive to structural than temporal variations.

Specifically, performance significantly drops even when only the first two packets of the biflow are lost (i.e., unseen) by the VP (Fig. 7). This holds especially on MIRAGE-19, where the drop is up to -81% . Conversely, on VPN-16 and MIRAGE-24, the drop—while still substantial—is smaller and increases more gradually as the number of unseen packets increases. These findings highlight that the model heavily relies on the earliest packets of the biflow, from which the PAY input is typically extracted. This is particularly critical in TLS traffic, where the initial packets often carry plaintext information (e.g., SNI, Cipher Suite, etc.) [16,45].

Similar findings are also observed when the VP is crossed only by one of the two paths connecting the client and the server (Fig. 8). Here, performance worsens mainly when the VP is crossed only by the backward path ($C \leftarrow S$), with drops up to -82% on MIRAGE-19. In this case, the impact remains relatively stable as the packet loss rate p increases up to 30%. This range was selected to rigorously cover the three severity classes identified in [56]: low ($< 5\%$), middle (5–20%), and high ($> 20\%$). When only the forward path ($C \rightarrow S$) is visible to the VP, a performance drop is observed across all datasets. The degradation is relatively stable on VPN-16, while it worsens more sharply with p on the others—e.g., from -19% to -36% on MIRAGE-24.

Finally, under $RT0$ scenarios (Fig. 9), we observe that the impact strongly depends on where the loss occurs. In the $RT0_{aft}$ case, performance remains almost unchanged (i.e., $< 1\%$), even as the retransmission probability r increases across the considered range [1%, 30%] on MIRAGE-19 and MIRAGE-24. This is likely because all copies of a packet traverse the VP, preserving the expected packet sequence—especially for the initial packets of the biflow, which are particularly important for classification. Conversely, on VPN-16, a more noticeable degradation is observed, with F1-score dropping from -3.6% to -5.4% as r increases.

When losses occur before the VP ($RT0_{bef}$), performance drops significantly as r varies within [1%, 30%]. For instance, the F1-score drops to up to -32% on MIRAGE-19 as r reaches the maximum severity of 30%—a value consistent with retransmission rates observed for small flows in congested networks [57]. This finding is likely driven by the increased probability of observing out-of-order packets, since retransmitted packets may be observed at the VP in positions that differ from their expected order in the unperturbed biflow.

Takeaway:

The baseline model handles time-based perturbations fairly well, meaning it does not depend much on precise inter-packet timing. But structural perturbations—like dropping, hiding, or reordering packets—hurt performance more, especially if they affect the first few packets or if the VP sees traffic from only one direction. This shows the model is sensitive to how much of the biflow it can observe.

5.3. Improving effectiveness with augmentation and contrastive learning

Herein, we test whether the training strategies from Section 4.2 improve TC effectiveness under clean (unperturbed) conditions. We compare each variant to the baseline (Base) using both (i) TC performance and (ii) the structure of the learned latent space. For each variant, Table 6 reports the results in terms of macro F1-score (i), DBI (ii), and SIL (ii) on all datasets.

On the one hand, focusing on TC performance (i.e., F1-score), we observe that employing Aug, SupCon, or a combination of both in at least one of the training phases consistently improves results over Base.

On both MIRAGE-19 and MIRAGE-24, applying either Aug or SupCon in a single phase (viz. Aug_{PT} , Aug_{FT} , and $SupCon_{PT}$) consistently leads to performance improvements. Notably, using Aug only during FT (Aug_{FT}) typically results in greater improvements than applying it solely during PT (Aug_{PT}). The best performance is achieved when Aug is employed in both phases—e.g., $+1.2\%$ on MIRAGE-19. Interestingly, apply-

Table 6

Performance of Aug_{PT}, SimCon_{PT}, SupCon_{PT}, Aug_{PT} → Aug_{FT}, SimCon_{PT} → Aug_{FT}, and SupCon_{PT} → Aug_{FT} compared to Base on VPN-16, MIRAGE-19, and MIRAGE-24. Results are reported as *mean ± std.dev.* over 10 folds. ↓ (resp. ↑): lower (resp. higher) is better. Baseline (Base) performance is shown in blue. For other variants, green (resp. red) highlights improvements (resp. degradations) w.r.t. Base, with darker shades indicating larger changes.

	VPN-16			MIRAGE-19			MIRAGE-24		
	F1-score [↑]	DBI [↓]	SIL [↑]	F1-score [↑]	DBI [↓]	SIL [↑]	F1-score [↑]	DBI [↓]	SIL [↑]
Base	77.91 ± 2.25	1.84 ± 0.09	0.13 ± 0.01	89.06 ± 0.44	1.92 ± 0.03	0.14 ± 0.00	95.42 ± 0.48	1.36 ± 0.04	0.27 ± 0.01
Aug _{PT}	79.82 ± 1.34	1.74 ± 0.05	0.14 ± 0.01	89.90 ± 0.39	1.78 ± 0.02	0.18 ± 0.01	96.16 ± 0.57	1.27 ± 0.03	0.33 ± 0.01
SimCon _{PT}	74.29 ± 1.06	2.16 ± 0.11	0.07 ± 0.01	88.72 ± 0.67	2.89 ± 0.08	-0.04 ± 0.01	92.40 ± 0.46	2.15 ± 0.1	0.05 ± 0.01
SupCon _{PT}	78.81 ± 1.62	1.51 ± 0.06	0.19 ± 0.01	88.95 ± 0.47	1.83 ± 0.31	0.08 ± 0.08	95.67 ± 0.55	1.11 ± 0.06	0.36 ± 0.02
Aug _{FT}	78.65 ± 1.21	1.91 ± 0.08	0.12 ± 0.01	90.20 ± 0.30	2.26 ± 0.03	0.09 ± 0.01	96.23 ± 0.59	1.61 ± 0.05	0.22 ± 0.01
Aug _{PT} → Aug _{FT}	80.16 ± 1.93	1.70 ± 0.08	0.14 ± 0.01	90.26 ± 0.26	1.75 ± 0.03	0.19 ± 0.01	96.57 ± 0.49	1.18 ± 0.04	0.35 ± 0.01
SimCon _{PT} → Aug _{FT}	77.47 ± 2.48	2.13 ± 0.11	0.07 ± 0.02	90.01 ± 0.42	3.17 ± 0.09	-0.05 ± 0.01	95.31 ± 0.54	2.54 ± 0.13	-0.00 ± 0.02
SupCon _{PT} → Aug _{FT}	80.46 ± 1.16	1.47 ± 0.06	0.20 ± 0.01	90.08 ± 0.36	1.79 ± 0.08	0.12 ± 0.04	96.22 ± 0.41	1.01 ± 0.06	0.39 ± 0.02

ing Aug only during FT (Aug_{FT}) provides results comparable to those obtained when combining SupCon with Aug (SupCon_{PT} → Aug_{FT}). Although these improvements may appear limited, it is important to note that the baseline already performs well under clean conditions (i.e., > 89%), leaving little room for further improvements.

A different trend is observed on VPN-16, where the baseline performance is substantially lower (i.e., 77.9%). Here, the adoption of Aug and SupCon leads to more pronounced improvements (i.e., up to +2.6%). The best results are achieved when such strategies are applied consistently across both training phases (Aug_{PT} → Aug_{FT} and SupCon_{PT} → Aug_{FT}), confirming the importance of combining these techniques in more challenging scenarios.

On the other hand, analysis of the latent space geometry shows that employing SupCon during PT enhances both the DBI and SIL, indicating better intra-class compactness and inter-class separation. These improvements are further amplified when Aug is applied during FT, suggesting that SupCon effectively structures the latent space while Aug helps refine class decision boundaries. Conversely, applying Aug in PT leads to only minor improvements, while using it in FT alone substantially worsens the representations learned during PT.

Interestingly, using SimCon during PT (SimCon_{PT} and SimCon_{PT} → Aug_{FT}) often results in marginal improvements or even performance drops on all datasets. This is likely due to its unsupervised nature, which promotes invariance across samples rather than class-discriminative features. These limitations are also reflected in the latent space geometry, underscoring its reduced effectiveness.

Takeaway:

Training strategies based on Aug and SupCon lead to consistent improvements in terms of both classification performance and latent space structure under clean (i.e., unperturbed) conditions. These gains are largest when the techniques are applied across both training phases—especially when combining SupCon in PT with Aug in FT. In contrast, SimCon consistently performs worse across all metrics.

5.4. Improving robustness with augmentation and contrastive learning

Herein, we investigate whether incorporating perturbation-based transformations into the training process—i.e., leveraging Aug and CL—can improve robustness to perturbed traffic conditions. Accordingly, we compare the performance of the best models identified in Section 5.3 (i.e., Aug_{PT}, SupCon_{PT}, Aug_{FT}, Aug_{PT} → Aug_{FT}, and SupCon_{PT} → Aug_{FT}) in terms of F1-score across the PHiding (Fig. 11), DHidings (Fig. 12), and RTOs (Fig. 13) conditions. We omit TJitter as it has a negligible impact on performance, and VPC for brevity, since its results are comparable to those of RTObef. We report results across all

datasets and parameter settings, with the baseline (Base) performance under the same perturbations included for comparison. A solid black line shows the Base model's accuracy on unperturbed samples, serving as an upper bound.

Overall, we observe that the models trained leveraging Aug and SupCon exhibit significantly higher robustness—i.e., a lower gap against the upper-bound—compared to Base across all considered scenarios, with SupCon_{PT} → Aug_{FT} and Aug_{PT} → Aug_{FT} consistently showing the highest robustness to perturbations.

Focusing on PktHid (Fig. 11), we observe that on MIRAGE-19 and VPN-16, Base shows a drastic drop when the first two packets of the biflow are lost (e.g., up to -81% of F1-score on MIRAGE-19). Conversely, both SupCon_{PT} → Aug_{FT} and Aug_{PT} → Aug_{FT} undergo a much smaller degradation (i.e., < 15% compared to the unperturbed case). Although the degradation increases as more initial packets are lost, these models consistently outperform Base under such perturbation. Notably, on MIRAGE-24, the performance drop remains relatively small (i.e., < 6%) even when the first four packets are lost.

Examining the DHidings cases (Fig. 12), we observe that the Aug- and SupCon-based approaches exhibit much lower sensitivity to both the observed directional flow (i.e., C→S or C←S) and the packet-loss probability p , compared to Base. For instance, Base experiences a substantial performance drop of up to -84% of F1-score when the VP observed only the $C \leftarrow S$ traffic with $p=30\%$ on the MIRAGE-19. In contrast, with Aug_{PT} → Aug_{FT} and SupCon_{PT} → Aug_{FT}, this drop is limited to at most -20% and is further reduced to -9% on MIRAGE-24.

Finally, looking at RTOs (Fig. 13), we observe that while Base shows distinct behaviors between the RTOaft and RTObef scenarios, this does not hold across the new variants. Specifically, under RTOaft, Base exhibits very limited performance degradation as the retransmission rate r increases (e.g., < 1% on both MIRAGE-19 and MIRAGE-24). Conversely, the Aug- and CL-based variants, in particular Aug_{PT} → Aug_{FT} and SupCon_{PT} → Aug_{FT}, not only show negligible degradation but even achieve slight improvements of up to +2% on VPN-16 in some cases. Lastly, in the case of RTObef, Base shows a performance drop of up to -22% on MIRAGE-24 when $r=30\%$. In contrast, Aug_{PT} → Aug_{FT} and SupCon_{PT} → Aug_{FT} show a significantly lower drop of only -1% on the same dataset.

Takeaway:

While Base suffers substantial performance degradation under perturbations and limited visibility, integrating perturbation-based transformations during the training process greatly improves model robustness across all the considered scenarios. Notably, the Aug_{PT} → Aug_{FT} and SupCon_{PT} → Aug_{FT} variants consistently exhibit the highest robustness across diverse scenarios and datasets.

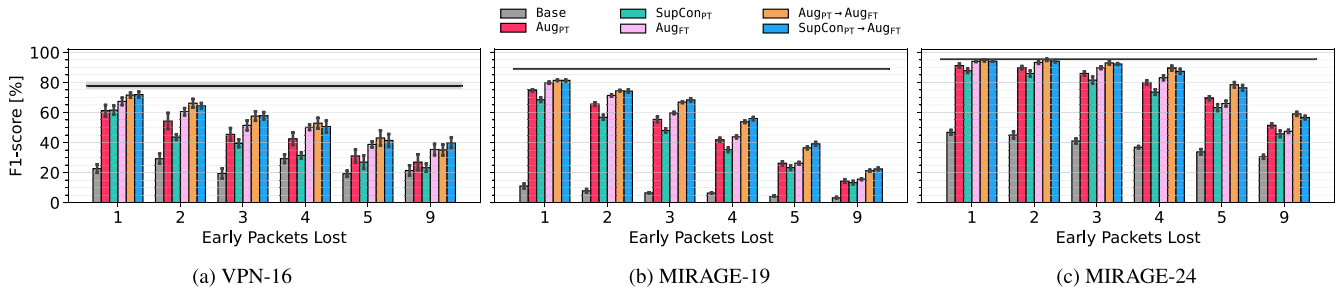


Fig. 11. Comparison of Base, Aug_{PT} , $SupCon_{PT}$, Aug_{FT} , $Aug_{PT} \rightarrow Aug_{FT}$, and $SupCon_{PT} \rightarrow Aug_{FT}$ under $PHiding$ as the number of lost initial packets in the biflow varies in $\{1, 2, 3, 4, 5, 9\}$. The black line reports the performance on the unperturbed test samples. Values are reported as $mean \pm std.dev.$ over 10 folds.

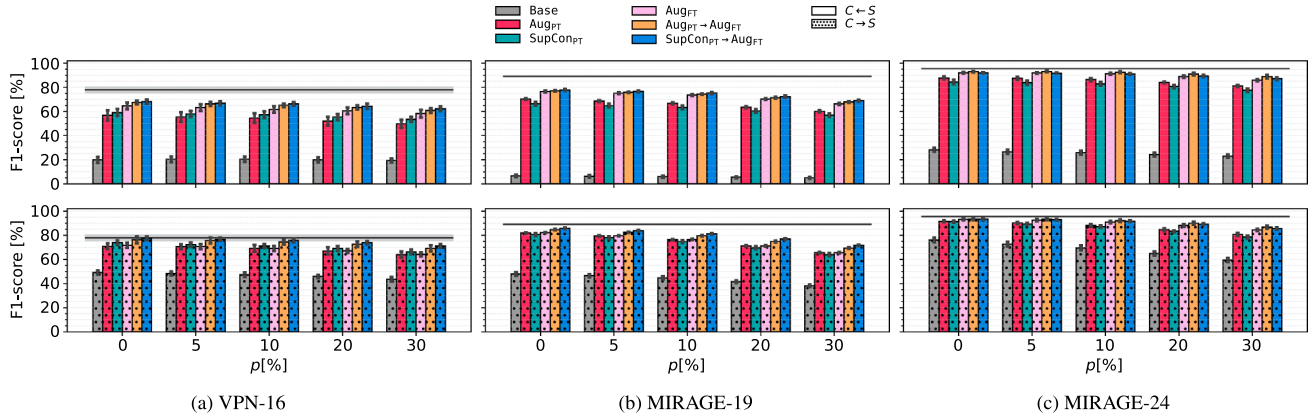


Fig. 12. Comparison of Base, Aug_{PT} , $SupCon_{PT}$, Aug_{FT} , $Aug_{PT} \rightarrow Aug_{FT}$, and $SupCon_{PT} \rightarrow Aug_{FT}$ under $DHiding$ ($p=0\%$) and $DHiding_{loss}$ ($p>0\%$) with packet-loss rate $p \in \{0, 5, 10, 20, 30\}$. Results are reported separately for scenarios where only the Client-to-Server ($C \rightarrow S$) or Server-to-Client ($C \leftarrow S$) direction of the biflow is visible at the VP. The black line reports the performance on the unperturbed test samples. Values are reported as $mean \pm std.dev.$ over 10 folds.

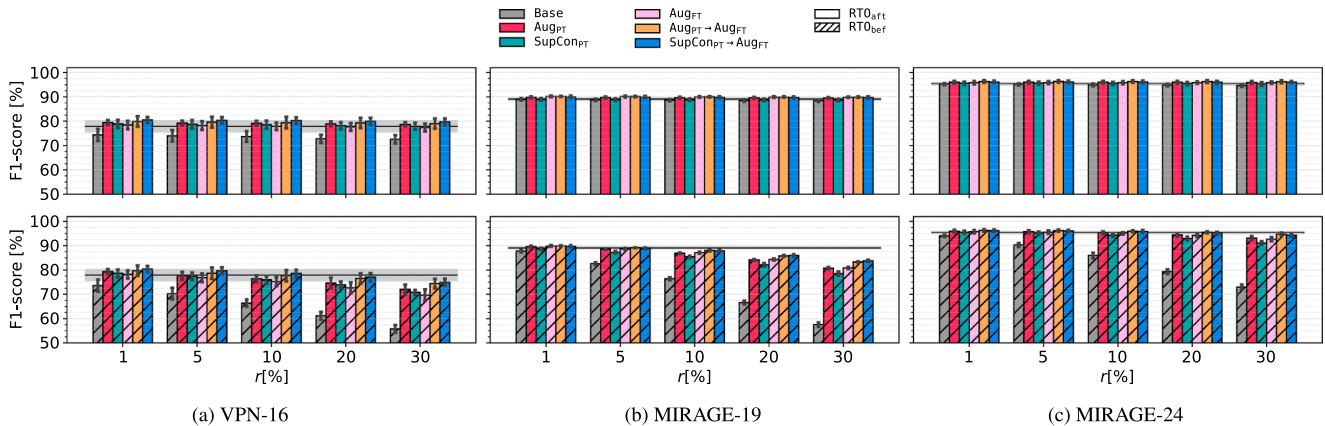


Fig. 13. Comparison of Base, Aug_{PT} , $SupCon_{PT}$, Aug_{FT} , $Aug_{PT} \rightarrow Aug_{FT}$, and $SupCon_{PT} \rightarrow Aug_{FT}$ under RTO scenarios with retransmission rate $r \in \{1, 5, 10, 20, 30\}$ and timeout $T_{outr} = 1s$. The black line reports the performance on the unperturbed test samples. Values are reported as $mean \pm std.dev.$ over 10 folds.

5.5. Assessing model reliability

In this section, we analyze the calibration of the models using the ECE metric. To this end, Table 7 reports the results obtained for Base, Aug_{PT} , $SupCon_{PT}$, Aug_{FT} , $Aug_{PT} \rightarrow Aug_{FT}$, and $SupCon_{PT} \rightarrow Aug_{FT}$ across all datasets.

Overall, we note that calibration varies significantly across datasets and models. Specifically, on MIRAGE-19 and MIRAGE-24, all models achieve good calibration, with ECE values typically below 3%. In particular, on MIRAGE-24, the ECE consistently remains below 1%, indicating well-calibrated models in this scenario, including the baseline. Conversely, calibration on VPN-16 is generally worse, with ECE values

around 6%–7% for all models. This is consistent with the lower F1-score of Base on VPN-16 compared to the other datasets, even in the clean setting, as previously discussed. Comparing the different strategies, we note that incorporating perturbations solely in PT consistently improves calibration (i.e., Aug_{PT} or $SupCon_{PT}$). For instance, $SupCon_{PT}$ achieves the lowest ECE on MIRAGE-19 and MIRAGE-24, outperforming Base and all other variants. Conversely, using transformations only in the FT phase (i.e., Aug_{FT}) consistently worsens calibration across all datasets compared to Base. This is likely due to applying perturbations during FT, which introduces input noise without significantly altering the feature space, as only the final layers—operating on fixed features from the frozen branch layers—are updated. Interestingly, for $SupCon_{PT} \rightarrow Aug_{FT}$,

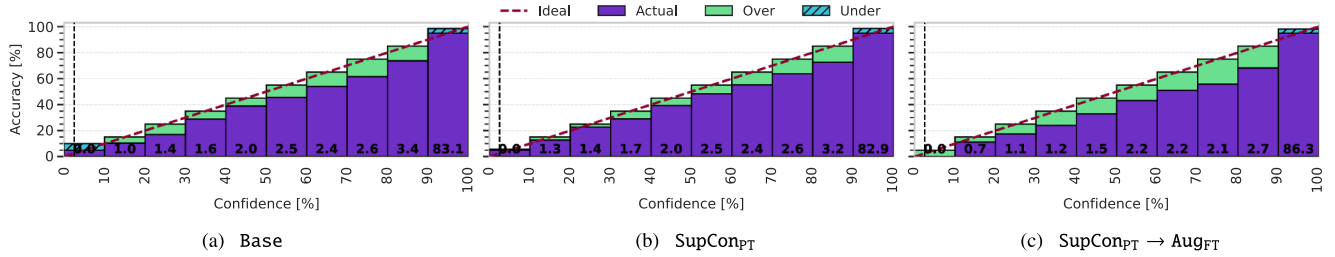


Fig. 14. Reliability diagrams of Base (a), SupCon_{PT} (b), and SupCon_{PT} → Aug_{FT} (c) on MIRAGE-19. Confidence is divided into 10 bins, and it is $\geq \frac{1}{|C|}$ (vertical dashed line), with $|C|$ being the number of classes. Over and under gaps represent an over-confident (optimistic) and under-confident (pessimistic) miscalibration pattern, respectively. The number at the bottom of each bar reports the share of samples within the corresponding bin.

Table 7

Calibration (ECE) of Aug_{PT}, Aug_{FT}, SupCon, Aug_{PT} → Aug_{FT}, and SupCon_{PT} → Aug_{FT} compared with Base on VPN-16, MIRAGE-19, and MIRAGE-24. Results are reported as *mean ± std.dev.* over 10 folds. ECE is reported as a percentage (lower is better). Baseline (Base) performance is shown in blue. For the new variants, green (resp. red) highlights improvements (resp. degradations) w.r.t. Base, with darker shades indicating larger changes.

	VPN-16	MIRAGE-19	MIRAGE-24
Base	5.75 ± 2.05	2.48 ± 0.34	0.75 ± 0.12
Aug _{PT}	5.78 ± 1.46	2.53 ± 0.38	0.64 ± 0.20
SupCon _{PT}	5.53 ± 0.87	2.28 ± 0.41	0.63 ± 0.11
Aug _{FT}	7.05 ± 1.09	4.06 ± 0.23	0.93 ± 0.16
Aug _{PT} → Aug _{FT}	6.53 ± 1.00	3.29 ± 0.34	0.82 ± 0.17
SupCon _{PT} → Aug _{FT}	6.47 ± 1.95	3.20 ± 0.23	0.81 ± 0.15

the benefits of SupCon in PT are partially offset by the use of augmentation in FT. As a result, SupCon_{PT} → Aug_{FT} shows worse calibration than SupCon_{PT} alone, although it still outperforms Aug_{PT} → Aug_{FT}.

In Fig. 14, we report the reliability diagrams for Base, SupCon_{PT}, and SupCon_{PT} → Aug_{FT} on MIRAGE-19. For Base, the sample distribution across bins is highly skewed, with the majority of predictions falling in the last bin (i.e., $\hat{p} > 90\%$). In this range, the model exhibits slight under-confidence. For lower-confidence bins, containing $< 20\%$ of samples, overconfidence increases with \hat{p} .

For SupCon_{PT}, the distribution is nearly identical compared to Base, but over-confidence is reduced, especially in the range $[0, 30]\%$. For SupCon_{PT} → Aug_{FT}, predicted confidence increases overall, with samples shifting toward high-confidence bins. Specifically, overconfidence increases in bins with confidence $< 90\%$, but its impact is reduced as more samples shift into the last bin. Since the behavior in this bin remains similar to Base, the share of samples with reliable predictions increases, leading to improved F1-score (see Section 5.3).

Takeaway:

Calibration varies with the training strategy but shows consistent patterns across datasets. PT with SupCon consistently improves model calibration. Conversely, FT with Aug increases overconfidence in low-confidence predictions. For SupCon_{PT} → Aug_{FT}, this effect is mitigated by a shift toward higher-confidence predictions, improving both reliability and classification performance.

5.6. Benefits in OOD setups

Herein, we evaluate the ability of models to detect OOD traffic samples—i.e., samples from classes not seen during training.

Accordingly, Table 8 compares the AUROC of Base with the best variants from Section 5.3 on the MIRAGE-19 and MIRAGE-24 datasets. To test generalization to truly unseen classes, we discarded overlapping

Table 8

OOD detection in terms of AUROC (%). When training on MIRAGE-19 (resp. MIRAGE-24), test samples from MIRAGE-19 (resp. MIRAGE-24) are considered ID, whereas those from MIRAGE-24 (resp. MIRAGE-19) are treated as OOD. AUROC (higher is better) is reported as *mean ± std.dev.* over 10 folds. Baseline (Base) performance is shown in blue, while the best-performing approach in each configuration (i.e., column) is highlighted in green.

	MIRAGE-19 (ID)	MIRAGE-24 (ID)
Base	85.33 ± 1.17	91.97 ± 1.31
Aug _{PT}	91.46 ± 1.24	94.40 ± 0.53
SupCon _{PT}	86.82 ± 1.76	93.46 ± 2.07
Aug _{FT}	90.42 ± 1.32	92.78 ± 0.83
Aug _{PT} → Aug _{FT}	91.25 ± 0.78	94.99 ± 0.44
SupCon _{PT} → Aug _{FT}	91.03 ± 0.92	95.27 ± 0.61
	MIRAGE-24 (OOD)	MIRAGE-19 (OOD)

classes between the training and test sets (i.e., Messenger, Skype, and Telegram). We focused on MIRAGE-19 and MIRAGE-24, which share the same collection setup, ensuring that observed performance differences are due to intrinsic traffic characteristics rather than artifacts from differing acquisition conditions. Specifically, when training on MIRAGE-19 (resp. MIRAGE-24), test samples from MIRAGE-19 (resp. MIRAGE-24) are treated as ID, while those from the other dataset are considered OOD.

The results show that models trained with Aug or CL consistently achieve higher AUROC values compared to the baseline, indicating better OOD detection capabilities.

When MIRAGE-24 is considered as OOD, the greatest improvement (i.e., +6% of AUROC) is obtained by using only augmentation during PT (i.e., Aug_{PT}). In the reverse scenario, using augmentation or SupCon in PT (i.e., Aug_{PT} or SupCon_{PT}) already provides benefits (e.g., +2.74% with Aug_{PT}), which are further enhanced by incorporating augmentation also in FT (e.g., +3.3% with SupCon_{PT} → Aug_{FT}).

Notably, although MIRAGE-24 contains only 20 classes compared to the 40 in MIRAGE-19, the baseline model already exhibits significantly higher OOD detection capabilities when trained on MIRAGE-24. This suggests that OOD detection performance is influenced not only by the number of classes but also by the intrinsic characteristics of the dataset.

Takeaway:

Incorporating perturbation-based transformations into the learning process, especially during both the pre-training and fine-tuning phases, substantially improves also OOD detection capabilities of the models.

5.7. Advantages in data-constrained training

In this section, we analyze the ability of models to generalize when trained with a limited amount of data. This scenario is particularly useful in real-world contexts where data scarcity often occurs due to privacy issues, collection costs, or other limitations. Hence, we systematically

Table 9

Performance of Base, Aug_{PT} → Aug_{FT}, and SupCon_{PT} → Aug_{FT} on MIRAGE-19, VPN-16, and MIRAGE-24 when reducing the training set size. δ represents the percentage of the training set used for training. Results are reported as *mean* \pm *std.dev.* over 10 folds. F1-score is reported as a percentage. \downarrow (resp. \uparrow): lower (resp. higher) is better.

δ [%]	Config.	VPN-16			MIRAGE-19			MIRAGE-24		
		F1-score [\uparrow]	DBI [\downarrow]	SIL [\uparrow]	F1-score [\uparrow]	DBI [\downarrow]	SIL [\uparrow]	F1-score [\uparrow]	DBI [\downarrow]	SIL [\uparrow]
100	Base	77.91 \pm 2.25	1.84 \pm 0.09	0.13 \pm 0.01	89.06 \pm 0.44	1.92 \pm 0.03	0.14 \pm 0.00	95.42 \pm 0.48	1.36 \pm 0.04	0.27 \pm 0.01
30	Base	69.30 \pm 2.18	2.31 \pm 0.12	0.09 \pm 0.01	82.89 \pm 0.69	2.29 \pm 0.03	0.10 \pm 0.01	91.54 \pm 1.11	1.53 \pm 0.06	0.21 \pm 0.01
	Aug _{PT} → Aug _{FT}	70.10 \pm 2.25	2.14 \pm 0.11	0.09 \pm 0.01	84.60 \pm 0.51	2.04 \pm 0.03	0.14 \pm 0.00	93.73 \pm 0.80	1.34 \pm 0.06	0.29 \pm 0.02
	SupCon _{PT} → Aug _{FT}	72.34 \pm 2.17	1.97 \pm 0.10	0.12 \pm 0.01	84.03 \pm 0.51	1.72 \pm 0.02	0.20 \pm 0.01	93.28 \pm 1.12	1.11 \pm 0.06	0.37 \pm 0.01
10	Base	57.73 \pm 3.18	2.94 \pm 0.13	0.02 \pm 0.01	73.34 \pm 0.91	2.69 \pm 0.05	0.07 \pm 0.01	83.89 \pm 1.60	1.77 \pm 0.06	0.16 \pm 0.01
	Aug _{PT} → Aug _{FT}	57.46 \pm 2.34	2.80 \pm 0.10	0.03 \pm 0.01	76.00 \pm 0.50	2.47 \pm 0.04	0.10 \pm 0.00	88.07 \pm 1.43	1.58 \pm 0.05	0.22 \pm 0.01
	SupCon _{PT} → Aug _{FT}	58.73 \pm 2.90	2.93 \pm 0.19	0.01 \pm 0.02	75.84 \pm 0.40	2.25 \pm 0.05	0.12 \pm 0.01	87.66 \pm 1.80	1.36 \pm 0.07	0.30 \pm 0.01
5	Base	50.09 \pm 1.85	3.52 \pm 0.34	-0.02 \pm 0.01	64.87 \pm 1.84	3.15 \pm 0.06	0.04 \pm 0.01	77.97 \pm 2.29	1.99 \pm 0.05	0.13 \pm 0.01
	Aug _{PT} → Aug _{FT}	49.26 \pm 3.01	3.42 \pm 0.26	-0.02 \pm 0.01	68.45 \pm 0.99	2.88 \pm 0.06	0.07 \pm 0.01	81.05 \pm 1.83	1.82 \pm 0.03	0.17 \pm 0.02
	SupCon _{PT} → Aug _{FT}	45.97 \pm 6.98	3.82 \pm 0.47	-0.06 \pm 0.02	67.96 \pm 1.99	2.70 \pm 0.10	0.06 \pm 0.02	82.72 \pm 2.06	1.62 \pm 0.07	0.24 \pm 0.01

reduce the training set and evaluate the impact on models performance. The primary focus is to assess whether **Aug** and **CL** can still yield benefits under these challenging conditions.

To this end, Table 9 reports the performance of Base, Aug_{PT} → Aug_{FT}, and SupCon_{PT} → Aug_{FT} in terms of F1-score, DBI, and SIL across all datasets when reducing (viz. downsampling) the original training set to $\delta \in \{5, 10, 30\}$ %. The downsampling was performed in a stratified manner to maintain the relative class proportions within each dataset. For comparison, we also report the performance of Base using the full training set ($\delta = 100$ %).

As expected, reducing the training set leads to a significant performance degradation across all datasets, especially in terms of F1-score (e.g., -24% for Base at $\delta = 5$ % on MIRAGE-19). Nevertheless, Aug_{PT} → Aug_{FT} and SupCon_{PT} → Aug_{FT} variants consistently outperform Base, especially as δ decreases on MIRAGE-19 and MIRAGE-24. For instance, SupCon_{PT} → Aug_{FT} outperforms Base of up to +4.8% of F1-score when using only 5% of MIRAGE-24, along with lower DBI (i.e., 1.6 vs. 2) and higher SIL (i.e., 0.24 vs. 0.13) values. Notably, as δ increases to 30%, both Aug_{PT} → Aug_{FT} and SupCon_{PT} → Aug_{FT} nearly recover the performance gap with Base with the entire training set (e.g., ≈ -2 % of F1-score on MIRAGE-24), with SupCon_{PT} → Aug_{FT} constantly achieving lower DBI and higher SIL values.

These results confirm that Aug_{PT} → Aug_{FT} and SupCon_{PT} → Aug_{FT} perform better under limited data conditions, showing improved performance on MIRAGE-19 and MIRAGE-24, along with better generalization despite fewer training samples. Conversely, on VPN-16, these trends are less marked, and in some cases, Base slightly outperforms both approaches in terms of F1-score—e.g., 50% for Base vs. 49% for Aug_{PT} → Aug_{FT} when $\delta = 5$ %. However, Aug_{PT} → Aug_{FT} and SupCon_{PT} → Aug_{FT} still achieve better DBI and SIL values.

Takeaway:

Our findings reveal that while downsampling inevitably reduces performance, particularly in terms of F1-score, the use of Aug and SupCon strategies can partially mitigate these effects on MIRAGE-19 and MIRAGE-24. However, this does not hold for VPN-16, where, in some cases, the baseline outperforms the other approaches in F1-score, despite a worse latent space regularization in terms of both DBI and SIL.

5.8. Assessing framework's generalizability

In this section, we investigate whether the proposed perturbation-aware learning approach can serve as a general, “model-agnostic” mechanism to robustify existing traffic classifiers, regardless of their underlying DL architecture.

Table 10

Performance (F1-score) of single-modal—1D-CNN (PAY) and HYBRID (PSQ)—and multi-modal—APPNET, MIMETIC, and MIMETIC-ENHANCED (PAY and PSQ)—architectures on MIRAGE-24. The results compare the standard training procedure (Base) with the proposed Aug strategies under both clean traffic and critical perturbation scenarios. The results for MIMETIC and MIMETIC-ENHANCED refer to the Aug_{PT} → Aug_{FT} strategy, while other (single-modal) models to the Aug_{PT} one. Green highlighted the best configuration for each architecture. Results are reported as *mean* \pm *std.dev.* over 10 folds.

Arch.	Var.	Clean	⊙ Test Perturbation ⊙		
			PHiding (LP = 5)	DHiding _{loss} (C ← S, p = 30%)	RTObef (r = 30%)
1D-CNN	Base	85.78 \pm 1.10	12.16 \pm 1.36	8.39 \pm 0.60	59.88 \pm 0.89
	Aug _{PT}	87.30 \pm 1.29	32.62 \pm 1.21	64.39 \pm 1.01	80.32 \pm 0.77
HYBRID	Base	88.09 \pm 0.80	21.32 \pm 1.25	8.82 \pm 1.11	8.74 \pm 0.80
	Aug _{PT}	92.08 \pm 0.77	48.86 \pm 3.83	56.91 \pm 5.58	63.35 \pm 4.34
APPNET	Base	92.30 \pm 0.69	14.50 \pm 1.95	8.16 \pm 1.13	66.92 \pm 1.02
	Aug _{PT}	92.67 \pm 0.76	41.13 \pm 3.27	71.90 \pm 2.36	87.00 \pm 0.86
MIMETIC	Base	93.28 \pm 0.38	18.53 \pm 1.30	8.86 \pm 0.64	57.45 \pm 0.80
	Aug _{PT} → Aug _{FT}	95.28 \pm 0.44	58.47 \pm 1.50	81.49 \pm 1.10	88.80 \pm 0.62
MIMETIC-ENHANCED	Base	95.42 \pm 0.48	33.76 \pm 1.82	22.99 \pm 1.49	72.94 \pm 1.04
	Aug _{PT} → Aug _{FT}	96.57 \pm 0.49	78.49 \pm 1.71	88.87 \pm 1.42	94.76 \pm 0.57

Clean: unperturbed traffic samples. LP: number of lost initial packets in the bilflow. p: packet-loss rate. C ← S: only Server-to-Client direction of the bilflow is visible at the VP. r: retransmission rate (%) with $T_{out} = 1$ s.

To this end, we extended our analysis to a diverse set of state-of-the-art baselines, applying the proposed training framework to both single-modal and external multimodal architectures. Specifically, we selected 1D-CNN [11] and HYBRID [13] as representative single-modal baselines operating on PAY and PSQ inputs, respectively. To validate the framework's transferability to different multimodal architectures, we also included MIMETIC [15] and APPNET [18], leveraging both traffic inputs.⁸

Table 10 details the comparative results on the MIRAGE-24 dataset. For each model, we report the F1-score obtained without integrating perturbations (viz. Base) against that achieved using the proposed perturbation-aware strategies. Specifically, we applied Aug_{PT} for 1D-CNN, HYBRID, and APPNET, whereas we employed Aug_{PT} → Aug_{FT} for MIMETIC and MIMETIC-ENHANCED. The evaluation considers both clean traffic—monitoring potential degradation on nominal data—and critical perturbation scenarios, corresponding to specific parameter settings that induce critical degradations.

The results highlight that integrating perturbations into the learning process yields consistent and substantial improvements in both scenarios. On the one hand, the proposed approach enhances the models' capability to discern classes in nominal (viz. clean) scenarios, particularly for architectures leveraging only PSQ input (i.e., +4% in F1-score). On the other hand, it improves model robustness when classifying perturbed traffic characterized by packet retransmission (i.e., RTObef) or partial visibility (i.e., PktHid and DirHidLoss). Finally,

⁸ APPNET considers only PL as PSQ input.

Table 11

Performance (F1-score) of MIMETIC-ENHANCED ($\text{Aug}_{PT} \rightarrow \text{Aug}_{FT}$) trained with individual perturbation types on MIRAGE-24 on both clean traffic and under critical perturbation scenarios. The *Baseline* ($\mathcal{T} = \emptyset$) and the *Complete Model* ($\mathcal{T} = \text{ALL}$) are included as lower and upper bounds, respectively. For each column, the top two performing configurations are highlighted in green, with darker shades indicating higher values. Entries on the main diagonal represent the performance when the model is evaluated against the same perturbation “seen” during training. Results are reported as *mean* \pm *std.dev.* over 10 folds.

		⊕ Test Perturbation ⊕					
\mathcal{T}		Clean	TJitter ($C \leftarrow S_{dl}, \sigma_d = 30\%$)	PHiding ($LP = 5$)	DHiding _{loss} ($C \leftarrow S, p = 30\%$)	RTObef ($r = 30\%$)	VPC _{loss} ($\alpha = 0.8, p = 30\%$)
\emptyset		95.42 \pm 0.48	96.07 \pm 0.53	33.76 \pm 1.82	22.99 \pm 1.49	72.94 \pm 1.04	71.98 \pm 0.50
⊕ Train Perturb.	TJitter	95.31 \pm 0.50	96.32 \pm 0.60	33.30 \pm 1.46	23.41 \pm 1.45	73.51 \pm 0.91	72.30 \pm 1.56
	PHidings	96.26 \pm 0.55	96.62 \pm 0.72	73.96 \pm 0.82	69.66 \pm 1.25	88.20 \pm 1.47	91.02 \pm 1.08
	DHidings	96.81 \pm 0.41	97.32 \pm 0.52	75.14 \pm 1.10	92.41 \pm 0.73	92.28 \pm 0.66	92.97 \pm 1.05
	RTOs	96.02 \pm 0.65	96.69 \pm 0.63	57.59 \pm 1.79	63.71 \pm 3.32	94.36 \pm 0.66	90.35 \pm 1.55
	VPCs	95.63 \pm 0.55	96.03 \pm 0.65	53.91 \pm 2.11	43.86 \pm 1.90	80.91 \pm 0.82	84.87 \pm 1.68
⊕ ALL		96.57 \pm 0.49	97.39 \pm 0.49	78.49 \pm 1.71	88.87 \pm 1.42	94.76 \pm 0.57	94.38 \pm 0.92

while multi-modal architectures consistently outperform single-modal ones, our perturbation-based strategy significantly bridges this gap on clean traffic. Notably, HYBRID performs almost on par with APPNET (with a difference of $< 0.5\%$ in F1-score), despite relying solely on the PSQ modality. Overall, among the multi-modal architectures considered, MIMETIC-ENHANCED proves to be the best performing model.

Takeaway:

Integrating perturbations into the training loop yields consistent improvements on both clean and perturbed scenarios across diverse architectures without requiring structural modifications. This confirms the model-agnostic nature of our approach. Moreover, while multi-modal architectures consistently outperform single-modal ones, our perturbation-based strategy significantly bridges this gap on clean traffic.

5.9. Assessing the contribution of specific perturbations

In this section, we disentangle the specific contribution of each perturbation type (defined in Section 3.2) to the model’s overall performance. In contrast to the comprehensive perturbation-aware strategy—which aggregates all transformations into a unified pipeline—we adopted a targeted ablation approach. Specifically, we trained distinct model variants, each exposed to only one specific class of perturbation during the learning phase (e.g., leveraging only DHiding-based augmentation).

Table 11 reports the results of the $\text{Aug}_{PT} \rightarrow \text{Aug}_{FT}$ variant on MIRAGE-24. For each perturbation type considered during training (see column \mathcal{T}), we report the F1-score on both clean traffic and under critical perturbation scenarios, corresponding to the parameter settings that induced the maximum degradation. For comparison, we also report the performance of the Base model—trained without perturbations ($\mathcal{T} = \emptyset$)—and the model obtained by considering all perturbations ($\mathcal{T} = \text{ALL}$).

As expected, across all scenarios, integrating a specific perturbation into the training loop confers robustness against that specific perturbation (see main diagonal). Remarkably, none of the perturbations introduced during training compromises the model’s ability to classify unperturbed samples. In fact, except for perturbations that mainly alter packet timing (i.e., TJitter and VPC), those modifying the biflow structure yield a slight improvement on unperturbed traffic compared to Base (up to +1.4% F1-score). Notably, DHidings lead to consistent

improvements in both nominal and perturbed scenarios. This suggests that integrating these perturbation types enhances robustness against other categories as well, effectively reducing the model’s dependence on complete packet sequences.

Consequently, by exposing the network to structural degradations such as DHidings, the model is compelled to reconstruct its classification logic from partial information. This indicates that the learned patterns are not merely a template match of the input stream, but rather represent a deeper abstraction of the communication flow. Lastly, the combined integration of all perturbations (viz. ALL) into the learning process yields the best performance across nearly all scenarios.

It is worth noting that, although this study adopts thresholds commonly used in the literature to limit PAY and PSQ inputs—i.e., $N_p = 10$ packets and $N_b = 576$ bytes, respectively—, such limits may vary in real-world scenarios based on operational needs. A robust classifier should not solely depend on the exact presence of specific initial units. Our evaluation of the PHiding perturbation serves as a proxy for assessing this sensitivity. By hiding the first LP packets, we effectively simulate a scenario where the “optimal” early window is either unavailable or shifted. Our results show that while the Base model experiences a significant drop when the initial handshake is missed (e.g., a drop of up to 81% in F1-score on MIRAGE-19), the proposed perturbation-aware strategies (e.g., $\text{Aug}_{PT} \rightarrow \text{Aug}_{FT}$) maintain stability. This confirms that our framework decreases the model’s sensitivity to the specific thresholds adopted. Additionally, the multimodal approach ensures that if a stricter byte threshold truncates the PAY modality, the PSQ modality effectively compensates for the missing data.

Takeaway:

Our analysis highlights a hierarchy in the effectiveness of domain-specific augmentations: while timing-based perturbations mainly enhance local robustness, structural changes, such as DHidings, enable the model to learn deeper and more invariant representations of traffic flow. This suggests that encouraging the model to reconstruct missing information is a more powerful learning signal than simply exposing it to variations in packet timing.

5.10. Computational complexity and real-world feasibility

In this section, we evaluate the practical deployability of our framework through a comprehensive analysis of the computational overhead during the learning phase (viz. training) and the computational costs, in

Table 12

Training time [HH : mm] of Aug_{PT}, SimCon_{PT}, SupCon_{PT}, Aug_{FT}, Aug_{PT} → Aug_{FT}, SimCon_{PT} → Aug_{FT}, and SupCon_{PT} → Aug_{FT} compared with Base on MIRAGE-19, VPN-16, and MIRAGE-24. Results are reported as *mean* ± *std.dev.* over 10 folds.

	VPN-16	MIRAGE-19	MIRAGE-24
Base	00:03 ± 00:00	00:21 ± 00:01	00:17 ± 00:02
Aug _{PT}	00:30 ± 00:04	06:42 ± 00:53	05:40 ± 00:50
SimCon _{PT}	02:12 ± 00:16	28:08 ± 01:51	21:17 ± 01:12
SupCon _{PT}	01:40 ± 00:10	24:18 ± 00:44	14:21 ± 01:18
Aug _{FT}	00:14 ± 00:02	05:27 ± 00:50	04:30 ± 00:59
Aug _{PT} → Aug _{FT}	00:39 ± 00:03	08:23 ± 01:11	06:59 ± 00:59
SimCon _{PT} → Aug _{FT}	02:24 ± 00:16	34:55 ± 02:52	25:11 ± 01:37
SupCon _{PT} → Aug _{FT}	01:51 ± 00:10	27:35 ± 00:56	16:17 ± 01:22

terms of memory and inference time, during the operational phase (viz. inference).

Training Overhead Analysis. We analyze the training overhead introduced by the strategies defined in Section 4.2. While these strategies enable the parallel training of individual branches during PT, we report sequential training times as the worst-case training time—i.e., measured by individually training the branches one after another during PT, followed by the subsequent FT phase.

Table 12 reports the training time of all evaluated methods across the three datasets. As expected, variants incorporating transformations in PT (Aug_{PT}, SimCon_{PT}, SupCon_{PT}), in FT (Aug_{FT}), or in both (Aug_{PT} → Aug_{FT}, SimCon_{PT} → Aug_{FT}, SupCon_{PT} → Aug_{FT}) exhibit significantly longer training times than Base, which completes in under 30 minutes overall.

Among variants leveraging transformations solely during PT, SimCon_{PT} shows the highest time (e.g., ≈80× on MIRAGE-19), followed by SupCon_{PT}. This overhead is partially due to the contrastive loss computation, which requires evaluating the similarity for all $B \cdot (B - 1)$ pairs in each training batch—introducing a quadratic computational cost with respect to batch size B . Notably, despite involving the same number of comparisons, SupCon_{PT} exhibits shorter training times compared to SimCon_{PT}, likely because its supervised loss benefits from a more structured optimization objective guided by class labels. Conversely, applying Supervised Data Augmentation (Aug) either during PT (Aug_{PT}) or FT (Aug_{FT}) remains relatively efficient (e.g., up to 19× on MIRAGE-19), and is consistently faster than the former variants.

Lastly, leveraging transformations in both training phases inevitably leads to a further increase in the total training time. However, the results suggest that starting from richer pre-trained branches—via Aug or SupCon—can facilitate optimization during FT, resulting in only a moderate additional overhead when employing Aug during the subsequent FT (e.g., an additional 10× for SupCon_{PT} → Aug_{FT} on MIRAGE-19). This is further supported by the fact that, when starting from pre-trained branches with SimCon, the FT based on Aug requires a significantly higher additional training time (e.g., an additional 19× on MIRAGE-19).

Since FT with Aug is consistent across all variants, this overhead does not stem from the FT setup itself. Instead, as earlier results show, it reflects the lower quality of the representation learned during PT with SimCon—leading to poorer classification, worse calibration, and less efficient FT.

Takeaway:

Training models with transformations adds a substantial time overhead compared to the Base model. While contrastive loss—especially SimCon—results in longer training times, a simpler Aug strategy offers a more efficient alternative. Higher-quality pre-trained models also improve FT efficiency, whereas weaker ones slow down the optimization process.

Inference Feasibility and Edge Deployment. Unlike the training phase, the operational phase must satisfy strict latency and memory constraints, typically imposed by resource-limited edge devices. In this regard, it is worth noting that the additional overhead induced by the proposed training strategies occurs only during the “offline” learning phase (one-time cost). Since these strategies act solely on the model weights, leaving the architectural backbone invariant, *the inference complexity remains unchanged compared to that of the Base model.*

To validate real-world deployability, Table 13 summarizes the computational performance of MIMETIC-ENHANCED across all datasets. Specifically, we report:

1. the number of trainable parameters (TP),
2. the number of FLOPs⁹ (FP),
3. the memory size on disk (DS), and
4. the inference time (iTime).

Furthermore, we analyze the impact of post-training quantization on these metrics. We compare the reference full-precision model (FP32) on both CPU and GPU against reduced-precision configurations, namely Half-Precision¹⁰ (FP16) on both architectures and Dynamic Quantization¹¹ (INT8) on CPU.

Overall, our findings highlight that the full-precision model is computationally lightweight (≈10.8 MFLOPs) and has an inference latency on CPU that is comparable to—and in some cases lower than—that of the GPU (i.e., 1.20 vs. 1.25 ms). This indicates that the model is latency-bound rather than compute-bound, making expensive GPU accelerators unnecessary for single-flow processing. In terms of memory optimization, dynamic quantization (INT8) proves to be a suitable strategy for CPU-based edge devices, reducing the memory footprint by ≈3.5× compared to the highest precision model (FP32), albeit with a slight latency overhead. As expected, FP16 emulation on standard CPUs results in a significant latency increase (≈+5 ms) due to the lack of native hardware support for half-precision arithmetic.

Finally, we assessed whether the efficiency gains provided by quantization come at the expense of classification effectiveness and robustness. Focusing on the MIMETIC-ENHANCED architecture trained with Aug_{PT} → Aug_{FT}, Table 14 compares the full-precision (FP32) model against its quantized (INT8) counterpart on MIRAGE-24. The evaluation covers both nominal conditions and specific critical perturbation scenarios.

As reported, the performance gap between the two models is negligible across all tested conditions. This confirms that the robust feature representations learned via our perturbation-aware strategy are resilient to aggressive numerical compression, validating the viability of INT8 deployment on edge devices without compromising reliability.

Takeaway:

Mimetic-Enhanced is computationally lightweight (≈1.20 ms/flow on standard CPUs), rendering dedicated GPU accelerators unnecessary. Furthermore, Dynamic Quantization (INT8) effectively reduces memory usage by ≈70% while preserving classification effectiveness, even under severe network perturbations, confirming suitability for resource-constrained edge devices.

⁹ We used the thop Python library <https://pypi.org/project/thop/>

¹⁰ We leveraged the native implementation provided by PyTorch.

¹¹ Leveraging the native PyTorch implementation, we selectively quantized the Dense and BiGRU layers, as they account for the vast majority of the model parameters.

Table 13

Computational performance of MIMETIC-ENHANCED on VPN-16, MIRAGE-19, and MIRAGE-24. The results compare the full-precision model (FP32) against reduced-precision configurations on both CPU and GPU—i.e., Half-Precision (FP16) and Dynamic Quantization (INT8). Inference time (iTime) is reported as *mean ± std.dev.* over 1,000 runs.

Dev.	Prec.	VPN-16				MIRAGE-19				MIRAGE-24			
		TP[K]	FP[M]	DS[MB]	iTime[ms]	TP[K]	FP[M]	DS[MB]	iTime[ms]	TP[K]	FP[M]	DS[MB]	iTime[ms]
GPU	FP32	877	10.79	5.90	1.24 ± 0.01	880	10.80	3.52	1.25 ± 0.01	878	10.8	3.47	1.25 ± 0.01
	FP16	877	10.79	2.96	1.27 ± 0.01	880	10.80	1.77	1.27 ± 0.01	878	10.8	1.74	1.27 ± 0.01
CPU	FP32	877	10.79	5.90	1.20 ± 0.03	880	10.80	3.52	1.20 ± 0.02	878	10.8	3.47	1.20 ± 0.03
	FP16	877	10.79	2.96	6.26 ± 0.04	880	10.80	1.77	6.25 ± 0.07	878	10.8	1.74	6.25 ± 0.05
	INT8	16.9	8.51	3.44	1.25 ± 0.03	16.9	8.51	1.01	1.24 ± 0.02	16.9	8.51	0.99	1.27 ± 0.03

TP: Trainable Parameters, FP: FLOPs, DS: Memory Size (on Disk), iTime: inference time.

Table 14

Performance (F1-score) of MIMETIC-ENHANCED ($A_{\text{AugPT}} \rightarrow A_{\text{AugFT}}$) under post-training quantization on MIRAGE-24. The full-precision model (FP32) is compared against the quantized model (INT8) across clean traffic and critical perturbation scenarios. Results are reported as *mean ± std.dev.* over 10 folds.

⊕ Test Perturbation ⊖				
Prec.	Clean	PHiding	DHiding _{loss}	RTObef
		(LP=5)	($C \leftarrow S, p = 30\%$)	($r = 30\%$)
FP32	96.57 ± 0.49	78.49 ± 1.71	88.87 ± 1.32	94.76 ± 0.57
INT8	96.58 ± 0.48	78.44 ± 1.72	88.86 ± 1.40	94.77 ± 0.53

Clean: unperturbed traffic samples. LP: number of lost initial packets in the biflow. p: packet-loss rate (%). r: retransmission rate (%) with $T_{\text{out}} = 1$ s. $C \leftarrow S$: only Server-to-Client direction of the biflow is visible at the VP.

6. Conclusion

Ensuring robustness in TC is essential for maintaining accurate and resilient network management in the face of evolving traffic patterns and the inherent uncertainty of real-world deployments.

In this work, we analyzed the impact of temporal and structural perturbations on TC performance. Through an extensive evaluation on three public datasets—i.e., VPN-16, MIRAGE-19, and MIRAGE-24—we showed that MIMETIC-ENHANCED, a state-of-the-art multimodal architecture, offers improved classification performance on both clean (i.e., *effectiveness*) and perturbed (i.e., *robustness*) traffic when compared to single-modal approaches. However, while resilient to time-based perturbations, MIMETIC-ENHANCED is *not robust to structural perturbations* involving packet reordering, delayed monitoring, and asymmetric paths, especially when its view on traffic is limited (e.g., up to -81% F1-score when just the first two packets of the biflow are missed on MIRAGE-19). This suggests a *strong dependence on the visibility and completeness of the biflow* observed at the VP.

To mitigate these vulnerabilities, we integrated *perturbation-aware training* into the MIMETIC-ENHANCED pipeline by leveraging Supervised Data Augmentation (Aug) and Contrastive Learning (CL) (SimCon and SupCon) during its Pre-Training (PT) and Fine-Tuning (FT) phases. Our analysis revealed that Aug and SupCon consistently enhance TC performance and robustness when applied across both phases (i.e., $A_{\text{AugPT}} \rightarrow A_{\text{AugFT}}$ and $\text{SupCon}_{\text{PT}} \rightarrow \text{SupCon}_{\text{FT}}$), while SimCon-based strategies fail to yield similar gains.

Beyond gains in effectiveness and robustness, *our analysis highlights further benefits*. In particular, using Aug during PT (or FT) *enhances OOD detection*—e.g., yielding an improvement of up to $+6\%$ AUROC when training on MIRAGE-19 and testing on MIRAGE-24. Similarly, employing SupCon in PT improves model calibration, with reductions of up to -0.2 ECE on MIRAGE-19. However, these benefits come at the cost of a slight degradation in calibration performance when applying Aug during FT, as it tends to increase model overconfidence. This effect is well-mitigated by $\text{SupCon}_{\text{PT}} \rightarrow \text{Aug}_{\text{FT}}$ and $\text{Aug}_{\text{PT}} \rightarrow \text{Aug}_{\text{FT}}$ strategies, which offer the best trade-off between classification performance and reliability.

We further investigated the impact of these strategies under *constrained training data conditions*. Our finding reveals that, while data scarcity inevitably reduces effectiveness, the adoption of Aug and SupCon can partially mitigate its effects in almost all cases.

Although perturbation-aware strategies significantly increase training time, Aug-based methods offer the best trade-off between effectiveness, robustness improvement, and computational overhead. In particular, while CL-based approaches (e.g., $\text{SupCon}_{\text{PT}}$) introduce a training overhead of over $80\times$, the simpler Aug strategies (e.g., Aug_{PT}) achieve comparable or even better performance with less than $30\times$ additional training time. Crucially, this computational burden represents a *one-time offline cost*. It does not affect the operational phase—as the model architecture and inference latency remain unchanged—but significantly enhances generalizability and robustness against real-world network fluctuations.

Finally, we demonstrate that the proposed *learning framework* is effectively *model-agnostic*. By explicitly validating its *broad generalizability* on representative classifiers—spanning both single-modal and multimodal paradigms—we prove that it enhances robustness regardless of the specific model design or training scheme.

In summary, our findings support the adoption of *multimodal architectures combined with perturbation-aware training strategies* as a practical, general-purpose solution to build robust, calibrated, and generalizable traffic classifiers, prioritizing deployment efficiency over architectural complexity.

Future directions will account for

1. adopting **eXplainable Artificial Intelligence (XAI)** for analyzing and improving model robustness;
2. extending robustness analysis to few-shot class-incremental learning scenarios;
3. evaluating model robustness in real-perturbed scenarios; and
4. comparing our traditional perturbation-aware training with Generative AI-based data augmentation

Data availability

Data will be made available on request.

CRedit authorship contribution statement

Idio Guarino: Writing – review & editing, Writing – original draft, Visualization, Validation, Software, Methodology, Investigation, Data curation, Conceptualization; **Giampaolo Bovenzi:** Writing – review & editing, Writing – original draft, Software, Methodology, Data curation, Conceptualization; **Alfredo Nascita:** Writing – review & editing, Writing – original draft, Visualization, Validation, Investigation; **Domenico Ciunzio:** Writing – review & editing, Writing – original draft, Validation, Supervision, Methodology, Investigation; **Damiano Carra:** Writing – review & editing, Supervision, Resources, Conceptualization; **Antonio Pescapè:** Writing – review & editing.

Algorithm 1 DHiding/DHiding_{loss}.

Input: Biflow BF of N packets; number of output packets N_p ; number of output bytes N_b ; maximum drop probability p_{\max} ; magnitude a

Output: PAY: N_b payload bytes; PSQ: headers of up to N_p packets

Note: EXTRPAY and EXTRPSQ return the PAY and PSQ modalities from the selected packets, applying padding if needed;

- 1: Identify unique directions D in BF
- 2: **if** $|D| = 1$ **then**
- 3: **return** (EXTRPAY(BF, N_b), EXTRPSQ(BF, N_p)) \triangleright **Skip:** one-way flow
- 4: **end if**
- 5: Sample direction: $d_{tx} \sim \mathcal{U}\{0, 1\}$ \triangleright Direction to transmit
- 6: Sample $\hat{p} \sim \mathcal{U}(0, p_{\max})$; $p \leftarrow a \cdot \hat{p}$
- 7: $I \leftarrow \{i \in [0, N] \mid \text{DIR}_i = d_{tx}\}$
- 8: For each $i \in I$, sample $z_i \sim \mathcal{U}(0, 1)$
- 9: $I' \leftarrow \{i \in I \mid z_i \geq p\}$ \triangleright Keep only non-dropped packets
- 10: **if** $|I'| = 0$ **then**
- 11: **return** (EXTRPAY(BF, N_b), EXTRPSQ(BF, N_p)) \triangleright **Skip:** no drop all
- 12: **end if**
- 13: PAY \leftarrow EXTRPAY(BF, N_b, I')
- 14: PSQ \leftarrow EXTRPSQ(BF, N_p, I')
- 15: **return** (PAY, PSQ)

Declaration of interests

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

Acknowledgements

This work is partially supported by the Italian PNRR MUR ‘‘Centro Nazionale HPC, Big Data e Quantum Computing, Spoke9 - Digital Society & Smart Cities’’; in part by the ‘‘xInternet’’ Project—funded by MUR—within the PRIN 2022 program (D.D.104-02/02/2022); in part by the European Union (EU) under the Italian National Recovery and Resilience Plan (NRRP) of NextGenerationEU, partnership on ‘‘Telecommunications of the Future’’ (PE00000001 - program ‘‘RESTART’’); and in part by the EU and MUR under the NRRP through the project ‘‘CSI-Future (Joint Communication and Sensing: CSI-Based Sensing for Future Wireless Networks)’’, PRIN 2022 PNRR (CUP B53D23023920001).’’ This manuscript reflects only the authors’ views and opinions and the Ministry cannot be considered responsible for them.

Appendix A. Pseudo-code

The present Appendix provides the pseudo-code of the transformations described in Section 3.2, including their input parameters, internal logic, and returned modalities.

Direction Hiding with Packet Loss (DHiding_{loss}) (Algorithm 1) simulates asymmetric paths where only one directional flow (i.e., upstream or downstream) crosses the VP, and packets may be dropped with probability $p \sim a \cdot \mathcal{U}(0, p_{\max})$. Setting $p_{\max} = 0$ disables packet dropping and corresponds to the basic **Direction Hiding** (DHiding).

Packets Hiding (PHiding) (Algorithm 2) simulates the scenario in which the VP starts observing the biflow with delay. The first visible packet is randomly selected from the initial portion of the flow, with the starting index $f \sim \mathcal{U}(1, \lfloor r \cdot T \rfloor)$, where $r \sim a \cdot \mathcal{U}(0, r_{\max})$.

Retransmission Timeout (RT0) (Algorithm 3) simulates packet retransmission within a biflow, where the packets to be retransmitted and the number of retransmissions are chosen randomly. Packets may be

Algorithm 2 Packets Hiding (PHiding).

Input: Biflow BF of N packets; number of output packets N_p ; number of output bytes N_b ; maximum drop ratio r_{\max} ; magnitude parameter a

Output: PAY: N_b payload bytes; PSQ: headers of up to N_p packets

Note: EXTRPAY and EXTRPSQ return the PAY and PSQ modalities from the selected packets, applying padding if needed; EXTRPSQ also recomputes IATs.

- 1: $T \leftarrow \min(N, N_p)$
- 2: **if** $T = 1$ **then** \triangleright Single packet
- 3: **return** (EXTRPAY(BF, N_b), EXTRPSQ(BF, N_p)) \triangleright **Skip:** single packet
- 4: **end if**
- 5: Sample $r \sim \mathcal{U}(0, r_{\max})$; $r \leftarrow a \cdot r$
- 6: Sample $f \sim \mathcal{U}\{1, \lfloor r \cdot T \rfloor\}$ \triangleright Index of first visible packet
- 7: $I \leftarrow \{f, f + 1, \dots, \min(f + N_p - 1, N - 1)\}$ \triangleright Up to N_p packets, within bounds
- 8: PAY \leftarrow EXTRPAY(BF, N_b, I)
- 9: PSQ \leftarrow EXTRPSQ(BF, N_p, I)
- 10: **return** (PAY, PSQ)

Algorithm 3 RT0aft/ RT0bef.

Input: Biflow BF of N packets; number of output packets N_p ; number of output bytes N_b ; max retransmission probability λ_{\max} ; magnitude parameter a ; timeout T_{out} ; min and max retransmissions L_{\min}, L_{\max} ; RT0aft flag drop_a

Output: PAY: N_b payload bytes; PSQ: headers of up to N_p packets

Note: EXTRPAY and EXTRPSQ return the PAY and PSQ modalities from the selected packets, applying padding if needed; EXTRPSQ also recomputes IATs.

- 1: Sample $\lambda \sim a \cdot \mathcal{U}(0, \lambda_{\max})$ \triangleright Retransmission rate
- 2: $T \leftarrow \min(N, N_p)$ \triangleright Maximum packets to consider
- 3: Sample $r_i \sim \mathcal{U}(0, 1)$ for each $i \in [0, T - 1]$
- 4: $P \leftarrow []$ \triangleright List of selected packets and retransmissions
- 5: **for** $i = 0$ to $T - 1$ **do**
- 6: $t_e^i \leftarrow \sum_{k=0}^{i-1} \text{IAT}_k$ \triangleright Elapsed time from start
- 7: **if** $r_i < \lambda$ **then** \triangleright Retransmission condition
- 8: **if** drop_a **then** \triangleright Add original packet (RT0aft)
- 9: Add $\langle i, t_e^i \rangle$ to P
- 10: **end if**
- 11: Sample $L \sim \mathcal{U}(L_{\min}, L_{\max})$ \triangleright # retransmissions
- 12: **for** $r = 1$ to L **do**
- 13: Add $\langle i, t_e^i + T_{out} \cdot r \rangle$ to P
- 14: **end for**
- 15: **else** \triangleright No retransmission
- 16: Add $\langle i, t_e^i \rangle$ to P
- 17: **end if**
- 18: **end for**
- 19: $P \leftarrow \text{sort}(P)$ \triangleright Sort packets by elapsed time
- 20: PAY \leftarrow EXTRPAY(BF, N_b, P)
- 21: PSQ \leftarrow EXTRPSQ(BF, N_p, P)
- 22: **return** (PAY, PSQ)

retransmitted with probability $\lambda \sim a \cdot \mathcal{U}(0, \lambda_{\max})$. Each selected packet may be retransmitted $L \sim \mathcal{U}(L_{\min}, L_{\max})$ times. Each selected packet is retransmitted after a preset timeout T_{out} . Two scenarios are considered: (i) RT0aft, where the original packet and all retransmissions are observed (loss after VP), and (ii) RT0bef, where only retransmissions are observed (original packet lost before VP).

Vantage Point Centralization (VPC) (Algorithm 4) simulates a scenario where both forward and backward paths cross the VP, but the VP is not

Algorithm 4 Vantage Point Centralization (VPC)/Vantage Point Centralization with Packet Loss (VPC_{loss}).

Input: Biflow BF of N packets; number of output packets N_p ; number of output bytes N_b ; α_{\max} ; maximum drop probability p_{\max} ; magnitude parameter a
Output: PAY: N_p payload bytes; PSQ: headers of up to N_p packets
Note: EXTRPAY and EXTRPSQ return the PAY and PSQ modalities from the selected packets, applying padding if needed;

```

1:  $R\hat{T}T \leftarrow \text{ESTIMATERTT}(BF)$   $\triangleright$  RTT estimation
2: Sample  $\hat{p} \sim \mathcal{U}(0, p_{\max})$ ;  $p \leftarrow a \cdot \hat{p}$   $\triangleright$  Packet drop rate
3: Sample  $\alpha \sim \mathcal{U}(0, \alpha_{\max})$ ;  $\triangleright$  RTT lag
4:  $T \leftarrow \min(N, N_p)$ 
5: if  $T = 1$  then
6:   return (EXTRPAY( $BF, N_b$ ), EXTRPSQ( $BF, N_p$ ))  $\triangleright$  Skip: single packet
7: end if
8: Sample  $p_i \sim \mathcal{U}(0, 1)$  for each  $i \in [0, T - 1]$ 
9:  $I \leftarrow \{i \mid i \in [0, T - 1] \text{ and } p_i \geq \hat{p}\}$   $\triangleright$  Indices of visible packets
10: if  $|I| = 0$  then
11:   return (EXTRPAY( $BF, N_b$ ), EXTRPSQ( $BF, N_p$ ))  $\triangleright$  Skip: no visible packets
12: end if
13: for  $i \in I \setminus \{\max(I)\}$  do  $\triangleright$  Iterate across the selected packets of  $BF$ 
14:   if  $DIR_i \neq DIR_{i+1}$  then
15:      $IAT_i \leftarrow IAT_i + R\hat{T}T \times \alpha \times (1 - 2 \cdot 1(DIR_i = du))$   $\triangleright$  Adjust IAT of  $i$ th packet of  $BF$ 
16:   end if
17: end for
18: PAY  $\leftarrow$  EXTRPAY( $BF, N_b, I$ )
19: PSQ  $\leftarrow$  EXTRPSQ( $BF, N_p, I$ )
20: return (PAY, PSQ)

```

Table A.1

IATs and jitter contributions observed at the VP.

Interval	Direction(s)	IAT _{<i>i</i>} (no jitter)	Jitter term
$t_1 - t_0$	C→S → C→S (cv)	$\Delta_1 + \delta_1^{cv} - \delta_0^{cv}$	$\lambda_1^{cv} - \lambda_0^{cv}$
$t_2 - t_1$	C→S → S→C (vs, sv)	$\Delta_2 + \delta_2^{sv} + \delta_2^{vs}$	$\lambda_2^{sv} + \lambda_2^{vs}$
$t_3 - t_2$	S→C → S→C (sv)	$\Delta_3 + \delta_3^{sv} - \delta_2^{sv}$	$\lambda_3^{sv} - \lambda_2^{sv}$
$t_4 - t_3$	S→C → C→S (vc, cv)	$\Delta_4 + \delta_4^{cv} + \delta_4^{vc}$	$\lambda_4^{cv} + \lambda_4^{vc}$

located near the endpoints. Under the assumption that traffic is captured close to the client, the transformation adjusts the IATs of packets sent in opposite directions to emulate the timing distortions observed by a VP placed farther away from the client (viz. *centralized*).

The IAT of each packet is adjusted according to the estimated RTT and a lag factor $\alpha \in [0, 1]$, which is sampled from $\mathcal{U}(0, \alpha_{\max})$.

The RTT is estimated as the time between the SYN and SYN+ACK packets for TCP biflows or the delay between the pair of packets with minimum cumulative size in opposite directions for UDP biflows. The Vantage Point Centralization with Packet Loss (VPC_{loss}) variant additionally introduces packet loss with probability $p \sim a \cdot \mathcal{U}(0, p_{\max})$. Setting $p_{\max} = 0$ disables packet dropping and corresponds to the basic VPC.

Directional Time Jittering (TJitter) (Algorithm 5) introduces delay jitter by modifying the IAT of packet data based on traffic direction (i.e., 0 = Client-to-VP and 1 = Server-to-VP). For each specified direction $d \in D$, jitter values are sampled from the positive half of a normal distribution— $\mathcal{N}^+(u, \sigma_d)$, where σ_d is the standard deviation for that direction. For the i th packet of a biflow, jitter is applied according to its direction, DIR_i , and influences its IAT, IAT_i , depending on whether the previous packet had the same direction or a different one.

To obtain this transformation, we provide a *packet timing model at a vantage point*. In particular, we model a sequence of packets exchanged

Algorithm 5 Directional Time Jittering (TJitter).

Input: Biflow BF of N packets; number of output packets N_p ; number of output bytes N_b ; list of traffic directions to jitter D ; mean jitter u ; per-direction standard deviation $\Delta = \{\sigma_d \mid d \in D\}$.
Output: PAY: N_b payload bytes; PSQ: headers of up to N_p packets
Note: EXTRPAY and EXTRPSQ return the PAY and PSQ modalities from the selected packets, applying padding if needed;

```

1: for  $d$  in  $D$  do
2:   for  $i \in [1, N_p - 1]$  do  $\triangleright$  Iterate from the 2nd to the  $N_p^{th}$  packet of  $BF$ 
3:     if  $DIR_i = d$  then
4:       Sample  $\delta_a^i, \delta_b^i \sim \mathcal{N}^+(u, \sigma_d)$   $\triangleright$  Per-direction jitter values
5:       if  $DIR_{i-1} = DIR_i$  then  $\triangleright$  Adjust IAT of  $i$ th packet of  $BF$ 
6:          $IAT_i \leftarrow \min(IAT_i + |\delta_a^i - \delta_b^i|, IAT_{\max})$ 
7:       else  $\triangleright$  Adjust IAT of  $i$ th packet of  $BF$ 
8:          $IAT_i \leftarrow \min(IAT_i + |\delta_a^i + \delta_b^i|, IAT_{\max})$ 
9:       end if
10:     end if
11:   end for
12: end for
13: PAY  $\leftarrow$  EXTRPAY( $BF, N_b$ )
14: PSQ  $\leftarrow$  EXTRPSQ( $BF, N_p$ )
15: return (PAY, PSQ)

```

between a Client (C) and a Server (S), as observed from a Vantage Point (V) located between them. Each packet is associated with:

- Δ_i : host-side thinking time before sending the i th packet,
- δ_i^{xy} : base propagation delay from x to y , where $xy \in \{cv, vs, sv, vc\}$,
- λ_i^{xy} : network-induced jitter on path xy for the i th packet,
- t_i : arrival time at the VP of the i th packet,
- $IAT_i = t_i - t_{i-1}$: inter-arrival time between packets $i - 1$ and i at VP.

The following equations define the arrival times of five example packets, sent in this order: C→S, C→S, S→C, S→C, and C→S. This sequence encompasses all relevant client/server interleaving, assuming a request/response communication pattern, as appropriate for a model restricted to application-level packets.

Accordingly, the arrival time t_i of each packet is composed as follows.

$$\begin{aligned}
t_0 &= \delta_0^{cv} + \lambda_0^{cv} \\
t_1 &= \Delta_1 + \delta_1^{cv} + \lambda_1^{cv} \\
t_2 &= t_1 + \delta_2^{vs} + \lambda_2^{vs} + \Delta_2 + \delta_2^{sv} + \lambda_2^{sv} \\
t_3 &= t_1 + \delta_2^{vs} + \lambda_2^{vs} + \Delta_2 + \Delta_3 + \delta_3^{sv} + \lambda_3^{sv} \\
t_4 &= t_3 + \delta_4^{vc} + \lambda_4^{vc} + \delta_4^{cv} + \lambda_4^{cv} + \Delta_4
\end{aligned}$$

Following these equations, the jittered IATs of packets results in Table A.1. Let DIR_i denote the direction of the i th packet in the biflow. The perturbed IAT, \widehat{IAT}_i , is computed as follows:

- **If** the preceding packet had the same direction as the current one ($DIR_i = DIR_{i-1}$):

$$\widehat{IAT}_i = \Delta_i + \delta_i^{DIR_i} - \delta_{i-1}^{DIR_i} + \lambda_i^{DIR_i} - \lambda_{i-1}^{DIR_i} = IAT_i + \lambda_i^{DIR_i} - \lambda_{i-1}^{DIR_i}$$
- **If** the preceding packet had the opposite direction as the current one ($DIR_i \neq DIR_{i-1}$, direction switch):

$$\widehat{IAT}_i = \Delta_i + \delta_i^{DIR_i} + \delta_{i-1}^{DIR_{i-1}} + \lambda_i^{DIR_i} + \lambda_{i-1}^{DIR_{i-1}} = IAT_i + \lambda_i^{DIR_i} + \lambda_{i-1}^{DIR_{i-1}}$$

References

- [1] T.T.T. Nguyen, G.J. Armitage, A survey of techniques for internet traffic classification using machine learning, IEEE Commun. Surv. Tutorials 10 (1–4) (2008) 56–76.

- [2] F. Pacheco, E. Exposito, M. Gineste, C. Baudoin, J. Aguilar, Towards the deployment of machine learning solutions in network traffic classification: a systematic survey, *IEEE Commun. Surv. Tutorials* (2018) 1–1.
- [3] G. Aceto, D. Ciuonzo, A. Montieri, V. Persico, A. Pescapé, AI-powered internet traffic classification: past, present, and future, *IEEE Commun. Mag.* 62 (9) (2024) 168–175.
- [4] L. Bernaille, R. Teixeira, I. Akodkenou, A. Soule, K. Salamatian, Traffic classification on the fly, *ACM SIGCOMM Comput. Commun. Rev.* 36 (2) (2006) 23–26.
- [5] M. Crotti, M. Dusi, F. Gringoli, L. Salgarelli, Traffic classification through simple statistical fingerprinting, *ACM SIGCOMM Comput. Commun. Rev.* 37 (1) (2007) 5–16.
- [6] A. Dainotti, F. Gargiulo, L.I. Kuncheva, A. Pescapé, C. Sansone, Identification of traffic flows hiding behind TCP port 80, in: *IEEE Int. Conf. Commun. (ICC)*, 2010, pp. 1–6.
- [7] P. M. Santiago del Rio, D. Rossi, F. Gringoli, L. Nava, L. Salgarelli, J. Aracil, Wire-Speed statistical classification of network traffic on commodity hardware, in: *Proceedings of the ACM Internet Measurement Conference (IMC)*, 2012, pp. 65–72.
- [8] G. Aceto, D. Ciuonzo, A. Montieri, V. Persico, A. Pescapé, MIRAGE: mobile-app traffic capture and ground-truth creation, in: *IEEE 4th International Conference on Computing, Communications and Security (ICCCS)*, 2019, pp. 1–8.
- [9] I. Guarino, G. Aceto, D. Ciuonzo, A. Montieri, V. Persico, A. Pescapé, Contextual counters and multimodal deep learning for activity-level traffic classification of mobile communication apps during COVID-19 pandemic, *Elsevier Comput. Netw.* 219 (2022) 109452.
- [10] Z. Wang, The applications of deep learning on traffic identification, *BlackHat USA* (2015).
- [11] W. Wang, M. Zhu, J. Wang, X. Zeng, Z. Yang, End-to-end encrypted traffic classification with one-dimensional convolution neural networks, in: *IEEE International Conference on Intelligence and Security Informatics (ISI)*, 2017, pp. 43–48.
- [12] G. Aceto, D. Ciuonzo, A. Montieri, A. Pescapé, Mobile encrypted traffic classification using deep learning: experimental evaluation, lessons learned, and challenges, *IEEE Trans. Netw. Serv. Manage.* 16 (2) (2019) 445–458.
- [13] M. Lopez-Martin, B. Carro, A. Sanchez-Esguevillas, J. Lloret, Network traffic classifier with convolutional and recurrent neural networks for internet of things, *IEEE Access* 5 (2017) 18042–18050.
- [14] M. Lotfollahi, M.J. Siavoshani, R.S.H. Zade, M. Saberian, Deep packet: a novel approach for encrypted traffic classification using deep learning, *Springer Soft Comput.* 24 (3) (2020).
- [15] G. Aceto, D. Ciuonzo, A. Montieri, A. Pescapé, MIMETIC: Mobile encrypted traffic classification using multimodal deep learning, *Elsevier Comput. Netw.* 165 (2019) 106944.
- [16] A. Nascita, A. Montieri, G. Aceto, D. Ciuonzo, V. Persico, A. Pescapé, XAI meets mobile traffic classification: understanding and improving multimodal deep learning architectures, *IEEE Trans. Netw. Serv. Manage.* 18 (4) (2021) 4225–4246.
- [17] Q. Yuan, G. Gou, Y. Zhu, Y. Wang, MMCo: Using multimodal deep learning to detect malicious traffic with noisy labels, *Front. Comput. Sci.* 18 (1) (2023) 181809.
- [18] X. Wang, S. Chen, J. Su, Automatic mobile app identification from encrypted traffic with hybrid neural networks, *IEEE Access* 8 (2020) 182065–182077.
- [19] G. Draper-Gil, A.H. Lashkari, M.S.I. Mamun, A.A. Ghorbani, Characterization of encrypted and VPN traffic using time-related features, in: *Proceedings of the Springer International Conference on Information Systems Security and Privacy (ICISSP)*, 2016, pp. 407–414.
- [20] I. Guarino, D. Ciuonzo, A. Montieri, A. Pescapé, Mirage-app×act-2024: a novel dataset for mobile app and activity traffic analysis, in: *IEEE 20th International Conference on Wireless and Mobile Computing, Networking and Communications (WiMob)*, 2024, pp. 663–666.
- [21] Z. Liu, R. Wang, N. Japkowicz, Y. Cai, D. Tang, X. Cai, Mobile app traffic flow feature extraction and selection for improving classification robustness, *Elsevier J. Netw. Comput. Appl.* 125 (2019) 190–208.
- [22] A.M. Sadeghzadeh, S. Shiravi, R. Jalili, Adversarial network traffic: towards evaluating the robustness of deep-learning-based network traffic classification, *IEEE Trans. Netw. Serv. Manage.* 18 (2) (2021) 1962–1976.
- [23] Y. Sharon, D. Berend, Y. Liu, A. Shabtai, Y. Elovici, TANTRA: timing-based adversarial network traffic reshaping attack, *IEEE Trans. Inf. Forensics Secur.* 17 (2022) 3225–3237.
- [24] C. Zhang, X. Costa-Pérez, P. Patras, Adversarial attacks against deep learning-based network intrusion detection systems and defense mechanisms, *IEEE/ACM Trans. Networking* 30 (3) (2022) 1294–1311.
- [25] I. Guarino, C. Wang, A. Finamore, A. Pescapé, D. Rossi, Many or few samples?: comparing transfer, contrastive and meta-learning in encrypted traffic classification, in: *IEEE 7th Network Traffic Measurement and Analysis Conference (TMA)*, 2023, pp. 1–10.
- [26] R. Xie, J. Cao, E. Dong, M. Xu, K. Sun, Q. Li, L. Shen, M. Zhang, Rosetta: enabling robust TLS encrypted traffic classification in diverse network environments with TCP-aware traffic augmentation, in: *Proceedings of the 32nd USENIX Security Symposium*, 2023, pp. 625–642.
- [27] W. Li, X.-Y. Zhang, H. Bao, B. Yang, Z. Li, H. Shi, Q. Wang, Prism: real-time privacy protection against temporal network traffic analyzers, *IEEE Trans. Inf. Forensics Secur.* 18 (2023) 2524–2537.
- [28] R. Xie, J. Cao, Y. Zhu, Y. Zhang, Y. He, H. Peng, Y. Wang, M. Xu, K. Sun, E. Dong, Q. Li, M. Zhang, J. Li, Cactus: obfuscating bidirectional encrypted TCP traffic at client side, *IEEE Trans. Inf. Forensics Secur.* 19 (2024) 7659–7673.
- [29] X. Deng, Y. Wang, Z. Xue, AN-Net: an anti-noise network for anonymous traffic classification, in: *Proceedings of the ACM Web Conference (WWW)*, 2024, p. 4417–4428.
- [30] C. Wang, A. Finamore, P. Michiardi, M. Gallo, D. Rossi, Data augmentation for traffic classification, in: *Springer Nature Passive and Active Measurement (PAM)*, 2024, pp. 159–186.
- [31] R. Ding, L. Sun, W. Zang, L. Dai, Z. Ding, B. Xu, Towards universal and transferable adversarial attacks against network traffic classification, *Elsevier Comput. Netw.* (2024) 110790.
- [32] X. Yang, S. Ruan, J.L. Yinliang, Y. Bo Sun, TrafCL: robust encrypted malicious traffic detection via contrastive learning, in: *Proceedings of the 33rd ACM International Conference on Information & Knowledge Management (CIKM)*, 2024, pp. 2910–2919.
- [33] C. Hajaj, P. Aharon, R. Dubin, A. Dvir, The art of time-bending: data augmentation and early prediction for efficient traffic classification, *Elsevier Expert Syst. Appl.* 252 (2024) 124166.
- [34] Y. Jiang, X. Wang, Y. Lai, Y. Wang, A packet sequence permutation-aware approach to robust network traffic classification, *IEEE Netw. Lett.* 6 (3) (2024) 203–207.
- [35] E. Horowicz, T. Shapira, Y. Shavitt, Self-supervised traffic classification: flow embedding and few-shot solutions, *IEEE Trans. Netw. Serv. Manage.* 21 (3) (2024) 3054–3067.
- [36] M. Jin, M. Apostolaki, Robustifying ML-powered network classifiers with PANTS, in: *Proceedings of the 34th USENIX Security Symposium*, 2025, pp. 625–642.
- [37] H. Lu, J. Liu, J. Peng, J. Lu, Adversarial attacks based on time-series features for traffic detection, *Elsevier Comput. Secur.* 148 (2025) 104175.
- [38] X. Wang, Q. Yuan, W. Yu, Q. Meng, S. Lu, W. He, C. Gu, Y. Wang, Entropy-regulated cross-modal generative fusion for multimodal network intrusion detection, *Elsevier Inf. Fusion* 126 (2026) 103581.
- [39] X. Wang, J. Li, X. Kuang, Y.-a. Tan, J. Li, The security of machine learning in an adversarial setting: a survey, *J. Parallel Distrib. Comput.* 130 (2019) 12–23.
- [40] A. Dainotti, A. Pescapé, K.C. Claffy, Issues and future directions in traffic classification, *IEEE Netw.* 26 (1) (2012) 35–40.
- [41] S. Huang, F. Cuadrado, S. Uhlig, Middleboxes in the internet: a HTTP perspective, in: *IEEE Network Traffic Measurement and Analysis Conference (TMA)*, 2017, pp. 1–9.
- [42] G. Wan, S. Liu, F. Bronzino, N. Feamster, Z. Durumeric, CATO: end-to-end optimization of ML-based traffic analysis pipelines, in: *Proceedings of the 22nd USENIX Symposium on Networked Systems Design and Implementation (NSDI 25)*, 2025, pp. 1523–1540.
- [43] A. D’Alconzo, I. Drago, A. Morichetta, M. Mellia, P. Casas, A survey on big data for network traffic monitoring and analysis, *IEEE Trans. Netw. Serv. Manage.* 16 (3) (2019) 800–813.
- [44] F. Cerasuolo, A. Nascita, G. Bovenzi, G. Aceto, D. Ciuonzo, A. Pescapé, D. Rossi, MEMENTO: a novel approach for class incremental learning of encrypted traffic, *Elsevier Comput. Netw.* 245 (2024) 110374.
- [45] F. Cerasuolo, I. Guarino, V. Spadari, G. Aceto, A. Pescapé, XAI For interpretable multimodal architectures with contextual input in mobile network traffic classification, in: *IEEE/IFIP Networking Conference (IFIP Networking)*, 2024, pp. 757–762.
- [46] X. Chen, H. Kim, J.M. Aman, W. Chang, M. Lee, J. Rexford, Measuring TCP round-trip time in the data plane, in: *Proceedings of the ACM Workshop on Secure Programmable Network Infrastructure (SPIN ’20)*, Association for Computing Machinery, 2020, pp. 35–41.
- [47] C. Shorten, T.M. Khoshgoftaar, A survey on image data augmentation for deep learning, *Springer J. Big Data* 6 (1) (2019) 1–48.
- [48] T. Chen, S. Kornblith, M. Norouzi, G. Hinton, A simple framework for contrastive learning of visual representations, in: *Proceedings of the 37th ACM International Conference on Machine Learning (CML’20)*, 2020, pp. 1597–1607.
- [49] X. Liu, F. Zhang, Z. Hou, L. Mian, Z. Wang, J. Zhang, J. Tang, Self-supervised learning: generative or contrastive, *IEEE Trans. Knowl. Data Eng.* 35 (1) (2023) 857–876.
- [50] P. Khosla, P. Teterwak, C. Wang, A. Sarna, Y. Tian, P. Isola, A. Maschinot, C. Liu, D. Krishnan, Supervised contrastive learning, in: *Proceedings of the 34th International Conference on Neural Information Processing Systems (NIPS ’20)*, 33, Curran Associates Inc., 2020, pp. 18661–18673.
- [51] D.L. Davies, D.W. Bouldin, A cluster separation measure, *IEEE Trans. Pattern Anal. Mach. Intell. PAMI-1* (2) (1979) 224–227.
- [52] P.J. Rousseeuw, Silhouettes: a graphical aid to the interpretation and validation of cluster analysis, *Elsevier J. Comput. Appl. Math.* 20 (1987) 53–65.
- [53] M. Kull, M. Perello-Nieto, M. Kängsepp, T.S. Filho, H. Song, P. Flach, Beyond temperature scaling: obtaining well-calibrated multiclass probabilities with dirichlet calibration, in: *Proc. of the 33rd International Conference on Neural Information Processing Systems (NeurIPS)*, 2019, pp. 12316–12326.
- [54] M.P. Naeini, G.F. Cooper, M. Hauskrecht, Obtaining well calibrated probabilities using bayesian binning, in: *Proceedings of the ACM 29th AAAI Conference on Artificial Intelligence (AAAI’15)*, 2015, p. 2901–2907.
- [55] L. Yang, A. Finamore, F. Jun, D. Rossi, Deep learning and zero-day traffic classification: lessons learned from a commercial-grade dataset, *IEEE Trans. Netw. Serv. Manage.* 18 (4) (2021) 4103–4118.
- [56] H. Wu, Y. Liu, S. Ni, G. Cheng, X. Hu, Lossdetection: real-time packet loss monitoring system for sampled traffic data, *IEEE Trans. Netw. Serv. Manage.* 20 (1) (2023) 30–45.
- [57] M.A. Ahmed, B. Bensaou, Enhancing TCP via hysteresis switching: theoretical analysis and empirical evaluation, *IEEE/ACM Trans. Networking* 31 (6) (2023) 2614–2623.