

# Magmaw: Modality-Agnostic Adversarial Attacks on Machine Learning-Based Wireless Communication Systems

Jung-Woo Chang\*, Ke Sun\*, Nasimeh Heydaribeni\*, Seira Hidano<sup>†</sup>, Xinyu Zhang\*, Farinaz Koushanfar\*

\*University of California San Diego <sup>†</sup>KDDI Research, Inc.

{juc023, kesun, nheydaribeni}@ucsd.edu, se-hidano@kddi.com, {xyzhang, farinaz}@ucsd.edu

**Abstract**—Machine Learning (ML) has been instrumental in enabling joint transceiver optimization by merging all physical layer blocks of the end-to-end wireless communication systems. Although there have been a number of adversarial attacks on ML-based wireless systems, the existing methods do not provide a comprehensive view including multi-modality of the source data, common physical layer components, and wireless domain constraints. This paper proposes Magmaw, the first black-box attack methodology capable of generating universal adversarial perturbations for any multimodal signal transmitted over a wireless channel. We further introduce new objectives for adversarial attacks on ML-based downstream applications. The resilience of the attack to the existing widely-used defense methods of adversarial training and perturbation signal subtraction is experimentally verified. For proof-of-concept evaluation, we build a real-time wireless attack platform using a software-defined radio system. Experimental results demonstrate that Magmaw causes significant performance degradation even in the presence of the defense mechanisms. Surprisingly, Magmaw is also effective against encrypted communication channels and conventional communications.

## 1. Introduction

Next-generation (NextG) wireless networks promise to support ultra-reliable and low-latency communication for rapidly evolving wireless devices [21]. The emerging networks are thus challenged to establish new features, such as adaptive coding and enhanced modulation schemes to overcome rapidly changing channel conditions and to achieve more efficient use of spectrum [11], [53], [84]. Machine Learning (ML) overcomes this barrier by revolutionizing the entire wireless network protocol stack [56].

Recent research [11] introduces an end-to-end communication system using deep neural networks (DNNs) for both transmitter and receiver, termed joint source-channel coding (JSCC). This ML-based approach optimizes source and channel coding together in a cross-layer framework to handle diverse channel conditions. To cope with multipath fading channels, the JSCC-encoded data can be further modulated into continuous signal waveforms by the orthogonal frequency division multiplexing (OFDM) [82]. ML-based JSCC model can convey the semantic information of each

modality more accurately than traditional communication systems [71], [75], [78], [82]. These advantages are recognized by the industry leaders and the Third Generation Partnership Project (3GPP), which are in the process of adopting ML in 5G and beyond [1], [27], [58], [70], [72].

DNN-based models have shown to be vulnerable to adversarial attacks [63]. Such attacks were first identified in the image domain [14], [22], [25], [45] and were later extended to other modalities. The susceptibility of the models to adversarial perturbations raises serious concerns for the safety of ML adoption in NextG [7], [33], [43]. To be specific, an unexpected high semantic error of the receiver due to small perturbations can cause significant disruptions or threaten customer safety in quality-sensitive applications, e.g., remote surgery [4] and autonomous vehicles [26], [49]. Emerging applications based on multimodal data, such as virtual reality [39] and robots [48] would suffer even more from the corruption of multiple inputs.

Many recent systems [5], [7], [28], [38], [43], [46], [57] aim at crafting adversarial signals on end-to-end wireless communication systems, but they make unrealistic assumptions about the attacker’s capability. For example, although JSCC has a modality-specific structure to maximize the coding efficiency, they assume that only a single modality (e.g., one-hot vector or image) is wirelessly communicated [5], [7], [28], [38], [46], [57]. They also assume that an adversary knows which modality is sent by the transmitter.

In practice, the above assumptions are not valid for the following reasons: 1) the transmitter typically incorporates data from all modalities into the data blocks and then sends them to the receiver; 2) if the adversary wants to recognize the modality of the signal, it needs to have access to the target ML model that carries out JSCC, and this is not always feasible, and 3) even if the adversary can detect the modality, high latency occurs until adversarial perturbations are generated and added to the victim signal. As a result, the realistic attacks must be modality-invariant.

**Our Contributions.** In this paper, we propose Magmaw, a black-box attack framework that generates universal multimodal perturbations on ML-based wireless communication systems. We show, for the first time, that modulated multimodal data can be exploited by adversaries, resulting in failure to restore the original data as well as subversion of

downstream services. We consider examples of downstream services such as video analytics, which analyzes human activity with spatio-temporal features of continuous images, and audio-visual event recognition (AVE) that predicts the event label based on representations over multiple input modalities. We assume OFDM is used as the underlying modulation scheme, which is commonly used in modern wireless communication standards (4G LTE [3], IEEE 802.11 family [29], 5G NR [40]). OFDM divides a data stream into multiple sub-streams, which are transmitted in parallel. We present a comprehensive framework that models the attack objective as an optimization problem and integrates the features of heterogeneous wireless networks as the domain constraints. We incorporate a transformation mechanism in the objective function so that our Universal Adversarial Perturbation (UAP) is trained to be robust to various transformations occurring in the physical layer. Therefore, Magmaw can operate over a hardware-implemented attack platform without requiring synchronization with the sender and the receiver and directly injects adversarial signals into the receiver’s antenna. Unlike previous studies [5], [7], [57] that only considered ideal wireless channels with additive white Gaussian noise (AWGN), Magmaw works in realistic and challenging multipath fading channels. We assume an adversary who does not have access to victim ML models and produces the adversarial perturbations by transferring them from a set of surrogate models for multimodal data. We implement Magmaw on the USRP software-defined radio platform and demonstrate its attack feasibility against multimodal wireless transmission. As case studies to understand transferability of our adversarial examples, we construct two target scenarios: (1) encryption-based secure communication that prevents adversaries from eavesdropping on radio signals, and (2) conventional wireless communication that is not based on ML techniques. We show that Magmaw causes high semantic error at the receiver in both scenarios, proving that our attack is transferable to unseen systems.

Table 1 shows Magmaw comparison with the existing wireless attacks. It includes the key elements to create practically feasible attacks. In the following, we describe in more details the challenges and our solutions.

**Input-Agnostic Perturbations.** We assume a black-box attack setting where the adversary lacks prior knowledge about the data’s modality and the exact channel model. The diversity of models for different source distributions makes crafting UAPs against multimodal wireless signals challenging. For instance, video transmission models use spatio-temporal prediction to optimize transmission efficiency through a temporal chain of dependency among adjacent data symbols in the OFDM packet. Additionally, the attacker’s ability to adjust its transmit signal pattern effectively depends on knowing the wireless channel matrix between the sender and receiver ( $\mathbf{H}_t$ ). However, since  $\mathbf{H}_t$  varies due to factors like link distance, mobility, and environment, not having this information makes crafting an effective attack challenging.

**Protocol-Agnostic.** Our black-box setting assumes the adversary does not know the modulation constellation mapping, the coding rate for the channel bandwidth, or how

the OFDM system assigns the complex symbols to multiple subcarriers. It is possible to design an attacker that classifies the modulation from the complex-valued baseband samples using modulation recognition [73], but the attack’s real-time constraints prohibits this attack.

Magmaw addresses the above challenges by designing a perturbation generator model (PGM) trained to create input- and protocol-agnostic perturbations on surrogate wireless models. Specifically, we build an ensemble of surrogate models with different modalities, modulation, and coding rates for transfer-based attack. We define the optimization problem as maximizing a weighted loss function over a fixed set of surrogate models. We train the perturbations to be universal to different  $\mathbf{H}_t$  by arbitrarily changing the channel matrix. To address the issue of insufficient information regarding OFDM subcarriers, we concatenate the OFDM with our PGM to generate adversarial multi-carrier signals. We then randomly shuffle the order of frequency-domain complex-valued symbols to train the UAPs that are robust to the unknown distribution of benign multi-carrier signals.

**Synchronization-Free Attacks.** An adversarial wireless device may not be precisely synchronized with legitimate transmitter or receiver in time or frequency domain. We address time and frequency de-synchronization issues between the adversarial device and the legitimate transmitter/receiver using a novel offline training process. Adversarial perturbations are less effective due to time and phase offsets, which the adversary cannot predict. To tackle this, we train the PGM with time shift and phase rotation functions, ensuring UAPs remain effective despite varying offsets.

**Evaluation on Existing Defenses.** Previous studies [7], [28], [43] have evaluated adversarial robustness through several defense strategies at the physical layer: adversarial training (AT), perturbation subtraction (PS). However, their attacks are not physically realizable, as they do not take into account the practical constraints in the wireless domain, as shown in Table 1. A recent study [43] generated adversarial perturbations to compromise the pilot signals that are used to perform downlink channel prediction and indoor localization. However, such adversarial signals are inefficient for perturbing OFDM data symbols containing semantic information from multimodal sources. We generate an adversarial attack capable of targeting both OFDM pilots and data symbols, taking into account the impact of channel transformation, and further evaluate our physically plausible attack against extensive existing defenses.

In summary, we make the following major contributions:

- Introduction of Magmaw, the first black-box universal adversarial attack to subvert ML-based wireless systems. We are the first to demonstrate adversarial attacks over USRP software-defined radio against multimodal ML-based wireless communication systems that include the physical layers of the modern WiFi and cellular systems.
- Our attack is input- and protocol-agnostic, i.e., oblivious to the input modality, channel conditions, constellation, coding rate, and OFDM specifications.

TABLE 1: A comparison of adversarial attacks against ML-based wireless transmission systems.

Attacks	Type	Channel	Black-box ML	HW Demo	Input-Agnostic		Protocol-Agnostic			Sync-Free		Defenses	
					Multimodal	$H_t$	Constellation	Coding Rate	OFDM	Time	Phase	AT	PS
[57]	Offline Attacks	AWGN								✓		□	□
[5]			✓								✓	□	□
[46]		Multipath Fading										□	□
[38]												□	□
[7]	Online Attacks	AWGN	✓								✓	■	■
[28]		Multipath Fading	✓				✓					■	□
[43]				✓	✓					✓	✓	■	□
Ours			✓	✓		✓	✓	✓	✓	✓	✓	■	■

$H_t$ : a channel matrix between the sender and the receiver; ✓: the item is supported by the attack.

■: the attack can compromise the defense; ■: the defense was considered, but the attack is not valid in the real world.; □: not mentioned in the corresponding research.

- Our attack is shown to fool the downstream tasks such as video classification and audio-visual event recognition.
- Experiments verify the resiliency of Magmaw against the well-known defenses of adversarial training and perturbation subtraction. The attack’s effectiveness for encrypted communication channels and non-ML-based wireless communication is also evaluated.

## 2. Background

We provide the background on ML-based wireless communications and state-of-the-art JSCC models. We describe related work about adversarial examples in Appendix B.

### 2.1. ML-based Wireless Communications

Conventional wireless communication systems suffer from dramatic performance degradation due to the cliff-effect in which the receiver is unable to recover the transmitted data if channel conditions are worse than a certain threshold [11]. ML-based JSCC is considered as a key enabler for NextG wireless systems to overcome such limitations [70], [72]. ML-based JSCC aims to train robust JSCC encoder and decoder on wireless channels infused with channel conditions similar to the real physical world. Therefore, ML-based communication systems have shown great potential in efficiently transmitting various modalities of data.

### 2.2. Modality-Specific JSCC Model

Existing JSCC systems [11], [71], [75], [78], [82] adopt modality-specific ML models. That is because each modality needs a specialized way of extracting semantic information for accurate symbol recovery at the receiver. Consequently, the model architectures and parameters have to be tailored to the characteristics of each modality. In this paper, we consider four state-of-the-art ML-based JSCC models for image [82], video [71], speech [75], and text transmission [78]. Figure 1 illustrates the commonly used structures for each modality. The image JSCC model [82] is trained to minimize distortion on a frame-by-frame basis. The video JSCC model [71] leverages spatio-temporal similarities between successive image sequences to remove the redundancy. In order to exploit this property, the video JSCC model adopts

the temporal coding structure  $\sigma$ , which clusters each consecutive sequence of pictures into a group of pictures (GOP). Each frame within the GOP is entered into the video JSCC in coding order rather than display order. This means that the video JSCC encoder compresses frames in a specific order. Specifically, for a total of  $P$  frames included in the GOP, the coding order of each frame is determined by the mapping function  $m_\sigma(t)$ , where  $1 \leq m_\sigma(t) \leq P$ . On the other hand, speech signals contain modality-specific characteristics such as speech rate and tone. Thus, the attention mechanism is utilized for speech JSCC models [75] to identify the essential speech features which helps accurate recovery of the signals at the receiver. The text JSCC model [78] is designed to effectively extract semantic information and cope with semantic distortion based on Transformers [69]. The text features recovered by the receiver are decoded into the text sentence through a greedy decoder [74]. In addition, a cross-entropy loss function [36] is utilized to understand semantic meaning while maximizing system capacity.

## 3. Threat Model

### 3.1. Attack Scenario

To explore the adversarial robustness of the NextG wireless communications, we focus on ML-based wireless systems. The attacker<sup>1</sup> is targeted towards wireless signals created by front-end sources (e.g., IoT devices) that are used to transmit the multimodal source data to back-end user(s).

In order to adopt the widely used wireless standards, we concatenate the ML-based models with OFDM blocks to increase the spectral efficiency and reduce the multipath channel effects [24], [61], [77], [82]. Since multipath fading channels and OFDM blocks can be represented as differentiable layers, ML-based wireless systems are trained end-to-end. We consider a practical scenario where the adversary aims to attack the target wireless channel in real time via a PGM trained in an offline process. Our adversary trains UAPs that are agnostic to any input signals and do not require synchronization with the transmitter/receiver. Since PGM is designed based on DNNs, there may be a delay in generating the UAP. To circumvent such delays, the adversary collects several perturbation signals generated

1. Throughout this paper, we use the terms “attacker” and “adversary” interchangeably.

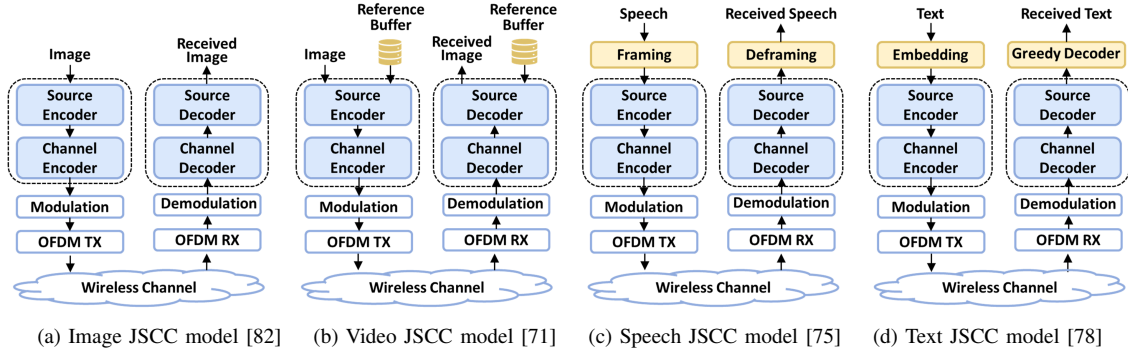


Figure 1: The modality-specific JSCC model for end-to-end wireless communication system.

offline by the PGM. At run time, it transmits these signals which will be superimposed with the legitimate signals at the receiver side. For simplicity, we assume the adversarial device employs a single antenna.

### 3.2. Adversary’s Goal

The goal of Magmaw is to transmit well-crafted wireless perturbations over the target wireless channel such that the legitimate receiver fails to restore the source data. To ensure the stealthiness, Magmaw aims to send adversarial signals with a small magnitude. Therefore, the wireless receiver cannot differentiate between generated adversarial perturbations and natural noise from wireless channels. Following the previous study [7], [57], we use a perturbation-to-signal ratio (PSR) metric to compare the power of the adversarial perturbation at the receiver with the received legitimate signal power. The PSR is set to be [-20,-10] dB [7], [57] so that the adversarial perturbation is not distinguishable with the expected natural noise in the radio channel.

### 3.3. Adversary’s Capability and Knowledge

**Wireless Communication System.** We consider the black-box attack settings where the adversary has no prior knowledge about the ML model architecture/parameters used by the victim transmitter/receiver. We assume that the adversary knows the model family of target ML models (e.g., auto-encoder), how to process each modality, and key physical layer techniques (e.g., OFDM transmitter and receiver). This is a realistic assumption for the following reasons: 1) standard documentation usually describes the core technology and is open to the public, and 2) specialized operations for each modality have already been widely known in the ML community. The adversary trains surrogate ML-based JSCC models using a large amount of publicly available multimodal data.

**Knowledge about Input.** We assume that the adversary does not know the modality, the constellation mapping method, and the number of OFDM symbols encoded by the JSCC model, due to the following reasons: 1) all the application-layer source data, regardless of modalities, need

to multiplex the transmitter radio and wireless channel, 2) the transmitter can adapt several types of modulation techniques according to channel conditions, and 3) the JSCC model varies the coding rate according to the state of the channel. For example, at a high SNR, a model with a low coding rate is used because source data can be transmitted with high throughput. Furthermore, we do not assume that the adversary knows how the transmitter maps OFDM symbol to the subcarriers. At the OFDM transmitter stage, the benign transmitter may assign modulated symbols arbitrarily to the subcarriers rather than in a fixed order. Each subcarrier will transmit a symbol vector with a different distribution.

**Target Wireless Channel.** We consider a practical attack scenario where the adversary cannot have access to the channel state information between the sender and the receiver, i.e.,  $\mathbf{H}_t$ . In addition, we do not assume that the adversary is synchronized with either the transmitter or the receiver. Therefore, the adversary does not know when the transmitter sends the wireless signal, which leads to a random time offset. Consequently, when the perturbation signal arrives at the receiver, there is a phase difference with the victim wireless signal. In addition, we assume that the attacker can determine the carrier frequency used by the targeted channel. The attacker can overhear the victim’s signals by arbitrarily adjusting its waveform bandwidth and carrier frequency using a reconfigurable wireless devices such as a software-defined radio. This allows an attacker to distinguish the target channel from other channels [42].

**Adversary’s Wireless Channel.** The perturbation signal generated by the adversary undergoes a multipath fading channel. We denote the channel matrix for the adversary as  $\mathbf{H}_a$ . According to the Wi-Fi protocol, the receiver periodically sends beacons to wireless devices within the range [8], [43]. An adversary can overhear this transmission and estimate the channel matrix from the receiver to itself. Due to the principle of reciprocity, this channel is the same as  $\mathbf{H}_a$ . In contrast to recent work [43], we relax the assumption that the adversary knows the exact channel matrix between the adversary and the receiver. We make a weaker assumption that the adversary has limited information, i.e., the distribution of the channel between the adversary and the receiver. To realize the UAP, we generate  $N_a$  random

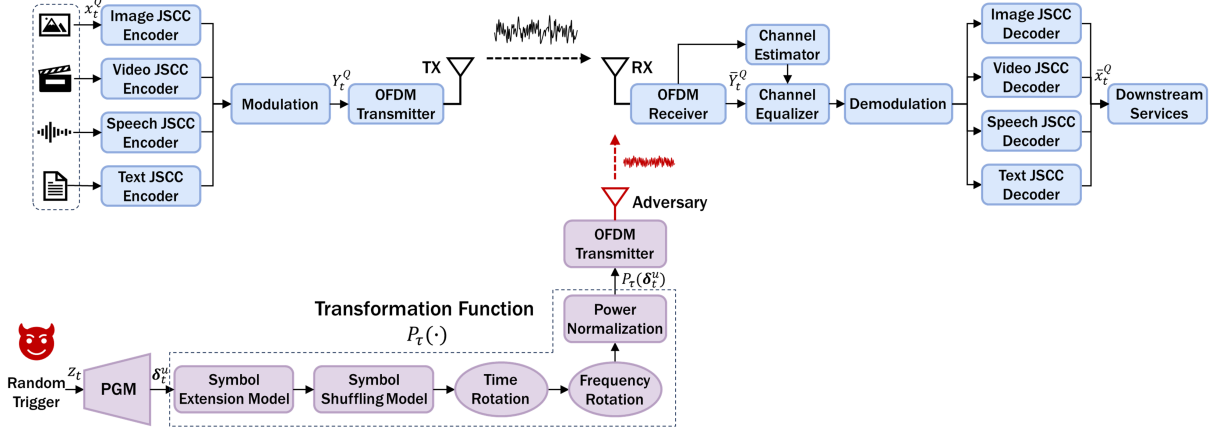


Figure 2: Overview of Magmaw.

samples  $\{\mathbf{H}_a^1, \dots, \mathbf{H}_a^{N_a}\}$  from the distribution. Then, we divide the samples into training and validation samples. We train the UAP with channel information from the training set and determine the best model from the validation set. After training is complete, we evaluate the UAP in real-world scenarios.

#### 4. System Model

Figure 2 illustrates the core processing blocks in the victim communication link along with the Magmaw attacker. **ML-based Transmitter.** We consider OFDM-based JSCC over a multipath fading channel with  $L_t$  paths. The multimodal source data are transmitted using  $N_s$  OFDM symbols with  $L_{fft}$  OFDM subcarriers. Note that  $N_s$  has different values depending on the modality and the coding rate. For channel estimation, the sender transmits a preamble (according to the publicly available wireless communication standards) on the subcarriers. We denote the source data as  $x_t^Q$  with modality  $Q \in \{\mathcal{I}, \mathcal{V}, \mathcal{S}, \mathcal{T}\}$  at time step  $t$ , where  $\mathcal{I}, \mathcal{V}, \mathcal{S}, \mathcal{T}$  denote the image, video, speech, and text, respectively. We describe a JSCC encoder for processing a modality  $Q$  with a given coding rate  $\lambda$  and modulation scheme  $C$  as a function  $E_{Q,C,\lambda}(x_t^Q, \mathcal{B}_t^Q)$ , where  $\mathcal{B}_t^Q$  is the transmitter's reference buffer used for the video JSCC model, as illustrated in Figure 1 (b). We define the reference buffer  $\mathcal{B}_t^Q$  containing the previously decoded frame  $\tilde{x}^{\mathcal{V}}(\cdot)$  as:

$$\mathcal{B}_t^Q = \begin{cases} \{\tilde{x}_{m_\sigma(1)}^{\mathcal{V}}, \dots, \tilde{x}_{m_\sigma(t-1)}^{\mathcal{V}}\}, & \text{if } Q = \mathcal{V}, \\ \emptyset, & \text{if } Q \neq \mathcal{V}. \end{cases} \quad (1)$$

Recall that  $m_\sigma$  is a function that finds the coding order of the  $t$ -th image in the given GOP structure  $\sigma$ .  $\mathcal{B}_t^{\mathcal{V}} = \emptyset$  when  $t=1$ . This is because the first frame is coded by the image JSCC model. Note that  $\tilde{x}^{\mathcal{V}}(\cdot)$  is reconstructed as the output of a video JSCC decoder that takes encoded video sequence  $E_{\mathcal{V},C,\lambda}(x_t^{\mathcal{V}}, \mathcal{B}_t^{\mathcal{V}})$  as input.  $\lambda$  is the coding rate to control the number of symbols per source data.

Then, a constellation mapping method  $M_C(\cdot)$  moves symbols to the nearest constellation points in a finite constel-

lation diagram  $C$ . The modulated symbol,  $Y_t^Q \in \mathbb{C}^{N_s \times N_{fft}}$ , can then be obtained as:

$$Y_t^Q = M_C(E_{Q,C,\lambda}(x_t^Q, \mathcal{B}_t^Q)), \quad (2)$$

Without loss of generality, we assume the target transmitter/receiver uses a single antenna following the 802.11a/g/n Wi-Fi standard [76]. We split  $Y_t^Q$  into a number of signal vectors with dimension of  $N_{fft}$ . Afterwards, an OFDM transmitter allocates divided signals on each subcarrier. Each OFDM symbol passes through an inverse discrete Fourier transform (IFFT), then a cyclic prefix (CP) is added and transmitted to the receiver over a multipath fading channel.

**ML-based Receiver.** The receiver obtains the complex-valued symbols from the channel output by removing the CP and applying FFT with an OFDM receiver. The received signal of the  $k$ -th subcarrier in the  $i$ -th OFDM symbol is given by:

$$\hat{Y}_t^Q[i, k] = \mathbf{H}_t[k] Y_t^Q[i, k] + W[i, k], \quad (3)$$

where  $\mathbf{H}_t \in \mathbb{C}^{N_{fft} \times N_{fft}}$  is the frequency-domain channel matrix, which is a diagonal matrix, and  $W \in \mathbb{C}^{N_s \times N_{fft}}$  is the frequency-domain AWGN matrix.

Given the FFT output of the pilot signals, the channel estimation and equalization are performed to compensate the channel-induced transformation. We adopt a least squares (LS) algorithm to predict the channel state information. After equalizing all of the divided signals with the channel equalizer  $R(\cdot)$ , we quantize the phase and amplitude of the signal on each subcarrier with  $M_C(\cdot)$ . Finally, we employ the JSCC decoder  $D_{Q,C,\lambda}(\cdot)$  to reconstruct an estimate  $\hat{x}_t^Q$  of the original signal. We express the entire process after OFDM receiver as follows:

$$\begin{aligned} \hat{x}_t^Q &= D_{Q,C,\lambda}(M_C(R(\hat{Y}_t^Q)), \hat{\mathcal{B}}_t^Q) \\ &= \mathcal{F}_{Q,C,\lambda}(\hat{Y}_t^Q, \hat{\mathcal{B}}_t^Q), \end{aligned} \quad (4)$$

where  $\hat{\mathcal{B}}_t^Q$  is the receiver's decoded buffer for the video JSCC model.  $\hat{\mathcal{B}}_t^{\mathcal{V}} = \{\hat{x}_{m_\sigma(1)}^{\mathcal{V}}, \dots, \hat{x}_{m_\sigma(t-1)}^{\mathcal{V}}\}$ , where  $\hat{\mathcal{B}}_t^{\mathcal{V}} = \emptyset$

when  $t=1$ .  $\hat{\mathcal{B}}_t^Q = \emptyset$  for other modalities. For simplicity, we denote all processes performed by the receiver after the OFDM receiver as  $\mathcal{F}_{Q,C,\lambda}(\cdot)$ .

## 5. Attack Construction

The framework of Magmaw is illustrated in Figure 2. Our attack methodology follows a hardware-algorithm code-sign to be robust to various distortions as the adversarial signals pass through the wireless channel in real-world scenarios. We formulate the adversaries' objective as an optimization problem, and adapt the Adam [34] optimizer to carefully craft UAPs.

### 5.1. Our Attack Formulation

**General Attack Formulation.** The objective of Magmaw is to find the adversarial perturbation that causes the receiver to fail to decode the wireless signal into the original multimodal data when the adversarial perturbation is superimposed on the receiver. Recall that the adversary does not know the transmitted signal  $Y_t^Q$  beforehand (Section 3). Instead of computing signal-wise perturbations, the adversary performs an offline process to train a single UAP that can target all benign signals. Moreover, the UAP must be agnostic to the input signals generated by all modalities and constellation mapping techniques. We define a UAP reflecting the above three conditions for input-agnostic perturbation as  $\delta^s \in \mathbb{C}^{N_s \times N_{fft}}$ . Therefore, when  $\delta^s$  is injected into the victim channel, the receiver obtains the frequency-domain channel output  $\bar{Y}_t^Q$  as:

$$\bar{Y}_t^Q[i, k] = \mathbf{H}_t[k]Y_t^Q[i, k] + \mathbf{H}_a^l[k]\delta^s[i, k] + W[i, k], \quad (5)$$

where the magnitude of perturbation  $\delta^s$  is bounded by the attacker's power budget  $\epsilon \in \mathbb{R}$ .  $\bar{Y}_t^Q[i, k]$  and  $\delta^s[i, k]$  represent, respectively, the frequency-domain perturbed response and the adversarial perturbation at the  $k$ -th subcarrier of the  $i$ -th OFDM symbol.  $\mathbf{H}_a^l$  is the  $l$ -th sample of the wireless channel between the adversary and the receiver.

Using Equation 4, the receiver then feeds this perturbed signal  $\bar{Y}_t^Q$  to the remaining physical layer components to reconstruct the source data with modality  $Q$  as:

$$\bar{x}_t^Q = \mathcal{F}_{Q,C,\lambda}(\bar{Y}_t^Q, \bar{\mathcal{B}}_t^Q), \quad (6)$$

where  $\bar{\mathcal{B}}_t^Q$  is the perturbed decoded frame buffer to be used in the video JSCC model.  $\bar{\mathcal{B}}_t^V = \{\bar{x}_{m_\sigma(1)}^V, \dots, \bar{x}_{m_\sigma(t-1)}^V\}$ , where  $\bar{\mathcal{B}}_t^V = \emptyset$  when  $t=1$ .  $\bar{\mathcal{B}}_t^Q = \emptyset$  for other modalities.

We aim to find the perturbation signals in a black-box setting, where the attacker has no access to the target DNN models. An effective way to deal with the black-box setup is to exploit the fact that adversarial examples show good transferability between different models [45]. To be specific, adversarial examples trained to fool a surrogate DNN model can be effective to subvert unknown DNN models. By adopting the attack transferability, we first train a surrogate JSCC model for each modality using publicly

available datasets that have different distributions from the target model's training data. Then we find a modality-agnostic adversarial perturbation  $\delta^s$  by solving the following optimization problem:

$$\arg \max_{\delta^s} \left[ \sum_{w \in \Psi^s} \mathcal{L}(w) \right], \text{ s.t. } \|\delta^s\|_2 < \epsilon, \quad (7)$$

where  $\Psi^s$  is a set of all wireless signals that can be created by physical layer elements.  $\mathcal{L}(w)$  is the loss function of ML-based JSCC model when  $w$  is sampled from  $\Psi^s$ .

However, this attack formulation is not suitable for making the UAPs physically realizable for the following reasons. First, having a single  $\delta^s$  as the UAP allows the receiver to estimate the perturbation signal using OFDM pilot signals, resulting in low robustness and persistence of adversarial attacks. Second, the adversary has no prior knowledge of the number of OFDM symbols in the target signal  $Y_t^Q$ , and thus is unable to define  $\delta^s$  as a matrix of the same size as  $Y_t^Q$ . Third, the video JSCC model has a network structure that forms a temporal chain between all video frames within the same GOP, so the model encodes current source data based on previous encoding results. This constructs the inter-frame dependency within a video sequence and it should be considered in generating the UAPs. Fourth, the adversary does not know the distribution of the channel inputs carried by each OFDM subcarrier. Finally, when the perturbation signal overlaps with the benign signal, time or phase offsets may occur.

**Practical Attack Formulation.** To address the above problems, we first construct an ML-based PGM  $G(z_t) = \delta_t^u$  that generates a modality-agnostic adversarial signal by receiving a random trigger  $z_t$  at time step  $t$ . The adversary changes  $z_t$  and injects a new perturbation signal into the target channel each time. Compared with using a single  $\delta^s$  as the UAP, the adversary can create an extremely large set of perturbations, which makes it difficult for the receiver to predict the perturbation signals. The following equation holds for frequency-domain complex-valued symbols at the receiver in the attacker's surrogate wireless systems:

$$\bar{Y}_t^Q[i, k] = \mathbf{H}_t[k]Y_t^Q[i, k] + \mathbf{H}_a^l[k]P_\tau(\delta_t^u)[i, k] + W[i, k], \quad (8)$$

where  $\delta_t^u \in \mathbb{C}^{N_g \times N_{fft}}$  denotes a UAP which contains  $N_g$  information symbols. Since the attacker does not know the number of target symbols,  $N_g$  may not be equal to the  $N_s$ . We define a novel transformation function  $P_\tau$  which enables the PGM-generated wireless signals to faithfully model the distribution of real wireless data. The transformation function consists of several steps: 1) symbol extension model, 2) symbol shuffling model, 3) time rotation, and 4) frequency rotation. The symbol extension model concatenates multiple PGM-generated perturbations such that the symbol-extended perturbations can perturb all OFDM symbols of the target radio signal. The symbol shuffling model makes our attack robust against unknown target symbols by randomly shuffling symbols between the OFDM subcarriers of the adversarial signal. The time and phase rotation function changes the offset of the adversarial signal during offline training

so that the adversarial signals are agnostic to random time and phase shifts in the real world. We also incorporate the power normalization into the transformation to make Magmaw undetectable from natural wireless noise. The wireless properties controlled by the transformation function are parameterized with  $\tau$ . Figure 2 shows all the modules included in the transformation function. With the help of  $P_\tau$ , the PGM is able to be optimized to produce the perturbation signals that are resilient to real-world transformations. In Section 5.2, we explain the internal mechanisms of  $P_\tau$ .

We define an optimization problem to train the PGM that generates a hardware-implementable perturbation signal as:

$$\arg \max_G \mathbb{E}_{z_t \sim p_z} \left[ \sum_{w \in \Psi^u} \mathcal{L}_{rx}(z_t, w) \right],$$

$$\mathcal{L}_{rx}(z_t, w) = \begin{cases} \mathcal{L}_D(x_t^I, \bar{x}_t^I), & \text{if } Q = \mathcal{I}, \\ \sum_{t=m_\sigma(1)}^{m_\sigma(P)} \mathcal{L}_D(x_t^V, \bar{x}_t^V), & \text{if } Q = \mathcal{V}, \\ \mathcal{L}_D(H_F(x_t^S), H_F(\bar{x}_t^S)), & \text{if } Q = \mathcal{S}, \\ \mathcal{L}_{CE}(H_G(x_t^T), H_G(\bar{x}_t^T)), & \text{if } Q = \mathcal{T}, \end{cases} \quad (9)$$

where  $\Psi^u$  is a set containing all radio signals that can be generated by the surrogate ML models. The perturbed signals at the receiver is computed from Equation 6. We use mean-squared error (MSE) loss as the distortion function  $\mathcal{L}_D$ . We train the PGM to maximize distortion on a frame-by-frame basis for the image JSCC model. For the video JSCC model, we consider the inter-frame dependency between adjacent frames as the sum of the distortions over all frames within the GOP. This allows the PGM to adapt to any GOP without the need to reconfigure the attack. As for speech, we transform the speech data into a one-dimensional vector via the deframing function  $H_F$  before the loss is calculated. Since the text JSCC model completes sentence restoration by sequentially finding the probabilities that words will appear with a greedy decoder  $H_G$ , we use a cross-entropy loss function  $\mathcal{L}_{CE}$  between the predicted sentence  $H_G(\bar{x}_t^T)$  and the ground truth sentence  $H_G(x_t^T)$ .

**Downstream Attack Formulation.** Figure 2 illustrates a downstream service appended to the ML-based wireless communication pipeline. As examples of downstream services, we consider two ML-based classifiers: 1) video classification (VC) and 2) audio-visual event recognition (AVR). Let  $F^N(\cdot)$  denote a discriminant function for the receiver's downstream recognition task  $N \in \{\text{VC}, \text{AVR}\}$ . After the receiver demodulates incoming perturbed signals into data, the discriminant function takes the data  $\bar{X}_N$  and outputs a probability distribution over a set  $K_N$  of class labels. Note that the video classifier takes a video clip  $\bar{X}_{\text{VC}} = \{\bar{x}_t^V\}_{t=1}^T$  consisting of  $T$  consecutive frames and the audio-visual model receives  $\bar{X}_{\text{AVR}} = \{\bar{x}_t^I, \bar{x}_t^S\}$  as two inputs. A classifier for task  $N$ ,  $\mathcal{C}^N$ , points  $\bar{X}_N$  to the class with the maximum probability:  $\mathcal{C}^N(\bar{X}_N) = \arg \max_{c \in K_N} F_c^N(\bar{X}_N)$ , where  $F_c^N$  is the probability of the perturbed input belonging to a specific class  $c$ . For downstream task  $N$ , we define the adversarial loss  $\mathcal{L}_{cls}^N$  to train PGM that subverts classifiers:

$$\mathcal{L}_{cls}^N = \max_{c \neq \mathcal{C}^N(\bar{X}_N)} F_c^N(\bar{X}_N) - F_{\mathcal{C}^N(\bar{X}_N)}^N(\bar{X}_N), \quad (10)$$

---

### Algorithm 1 Magmaw

---

**Input:** Dataset  $\mathbb{T}^Q$ , Surrogate JSCC model, Power constraint  $\epsilon$   
**Output:** PGM  $G(\cdot)$   
**for** epoch  $l < \text{MaxIter}$  **do**  
  **for** each modality  $Q \in \{\mathcal{I}, \mathcal{V}, \mathcal{S}, \mathcal{T}\}$  **do**  
    **for** each batch  $\mathbf{B}^Q \in \mathbb{T}^Q$  **do**  
       $C, \lambda \leftarrow$  is sampled uniformly from candidates  
       $\mathbf{H}_t$  is randomly sampled from channel model  
       $\mathbf{H}_a$  is sampled uniformly from training set  
      **if**  $Q = \mathcal{V}$  **then**  
        **for**  $\mathbf{x}_t^V \in \mathbf{B}^V (= \{\mathbf{x}_1^V, \dots, \mathbf{x}_P^V\})$  **do**  
           $Y_t^V \leftarrow$  Equation 2  
           $\mathcal{B}_t^V.append(\bar{x}_t^V)$   
           $z_t \sim \text{Uniform}(0, 1)$   
           $\tau \leftarrow$  uniformly at random  
           $\bar{Y}_t^V[i, k] \leftarrow$  Equation 8  
           $\bar{x}_t^V \leftarrow$  Equation 6  
           $\bar{\mathcal{B}}_t^V.append(\bar{x}_t^V)$   
        **else**  
           $Y_t^Q \leftarrow$  Equation 2  
           $z_t \sim \text{Uniform}(0, 1)$   
           $\tau \leftarrow$  uniformly at random  
           $\bar{Y}_t^Q[i, k] \leftarrow$  Equation 8  
           $\bar{x}_t^Q \leftarrow$  Equation 6  
        Update PGM  $G$  by solving Equation 11  
      **Return:** PGM  $G$

---

where  $\hat{X}_N$  denotes the reconstructed data when there is no attack.  $\hat{X}_{\text{VC}} = \{\hat{x}_t^V\}_{t=1}^T$  and  $\hat{X}_{\text{AVR}} = \{\hat{x}_t^I, \hat{x}_t^S\}$ . The attack succeeds when  $\mathcal{L}_{cls}^N > 0$ . Following the black-box setting, we find UAPs that maximize  $\mathcal{L}_{cls}^N$  for the surrogate model with different architectures from the target model. We then fool the downstream services by transferring the attacks calculated from the surrogate model to the target model.

**Unified Attack Formulation.** Finally, we integrate  $\mathcal{L}_{rx}$  and  $\mathcal{L}_{cls}^N$  into the objective function so that UAPs generated by the PGM can perturb wireless communication and downstream services simultaneously. Specifically, our goal is to solve the following objective function:

$$\arg \max_G \mathbb{E}_{z_t \sim p_z} \left[ \sum_{w \in \Psi^u} [\mathcal{L}_{rx}(z_t, w) + \sum_{N \in \mathcal{N}} \beta_N \mathcal{L}_{cls}^N(z_t, w)] \right], \quad (11)$$

where the parameter  $\beta_N$  weighs the relative importance of each downstream service and  $\mathcal{N} = \{\text{VC}, \text{AVR}\}$ .

In Algorithm 1, we outline the offline training process of Magmaw. Essentially, our goal is to train the PGM  $G$  that generates UAPs to subvert multimodal JSCC models. The ML model used for training is a surrogate model that is different from the target model. We utilize the transformation function  $P_\tau$  to change the outputs of PGM to practically feasible adversarial signals. At each iteration of the training, the algorithm selects a batch from the training dataset  $\mathbb{T}$  with a different distribution from the training dataset of the target model, and then calculates the objective function. We then calculate the gradient of our loss function in Eq. 11, and optimize the PGM by shifting the model parameters in the direction of the gradient. We apply the well-known optimization technique, Adam [34], to find the UAPs.

## 5.2. Design of the Transformation Function

To cope with challenging real-world scenarios, the adversary should craft input-agnostic UAP signals regardless of synchronization with the legitimate receiver. The transformation function  $P_\tau$  helps PGM learn to produce perturbations with a distribution similar to that of adversarial signals that can be realized in the real environment. Therefore, our adversarial signals are agnostic to 1) inconsistency of the number of information symbols between the benign signal and the adversarial signal, 2) unknown symbol allocation across the OFDM subcarriers, 3) time misalignment, and 4) unknown phase rotation. We additionally include a power regularization module in  $P_\tau$  for undetectability. The modules included in the transformation function are shown in Figure 2 and detailed below.

**Symbol Extension Model.** The number of OFDM symbols varies greatly depending on the modality and coding rate. For instance, the number of symbols generated after encoding a video clip is proportional to the number of frames and the resolution. Furthermore, the coding rate of the JSCC encoder determines the amount of information compressed. In an online attack scenario, the modality and coding rate are unknown to the adversary. This leads to the need to make the perturbation signal invariant to the number of OFDM symbols contained in the target signal. As the information about coding rate is publicly available, we assume that the adversary knows the maximum value of  $N_s$ . The adversary concatenates the perturbation signal multiple times through function  $K(\cdot)$  such that  $\mu \cdot N_g$  is equal to the maximum value of  $N_s$ , where  $\mu$  is a parameter to adjust the number of symbols. Therefore, the augmented perturbation signal can perturb all target OFDM symbols without knowing  $N_s$ .

**Symbol Shuffling Model.** Previous works [28], [43] make the assumption that the adversary knows how the target wireless system allocates symbols to each subcarrier. However, this is infeasible in practice, because standard wireless communication systems often randomize the allocation to prevent consecutive repetition of the same symbols. Our adversary aims to make a subcarrier-invariant perturbation that is universally applicable to any symbol distribution of subcarriers. We define a function  $\Gamma(\cdot)$  that randomly shuffles the symbols assigned to the OFDM subcarrier based on  $\zeta$ . Consequently, we train the PGM to generate the adversarial signal that is robust to the unknown symbol distribution across OFDM subcarriers.

**Time and Frequency Rotation.** To create shift-invariant perturbation signals, adversaries can simultaneously inject perturbation signals when a benign sender begins to transmit multimodal signals. However, due to the time misalignment between the sender and the adversary in the physical attack scenario, it is difficult to match the perturbation signal accurately with the victim signal when it arrives at the receiver side. At the OFDM receiver, phase rotation occurs in each OFDM subcarrier due to the time difference  $\Delta t$  between the adversarial and the benign signals. We model the phase shift that occurs in the  $k$ -th subcarrier of each OFDM symbol as  $e^{-j2\pi f_k \Delta t}$ . This phase shift is common in typical OFDM



Figure 3: A scenario where an unknown adversary sends a adversarial signal from behind a wall.

receivers because some of the OFDM symbols are out of time sync [35], [66]. Furthermore, the adversary lacks phase synchronization with the transmitter, resulting in a random phase offset  $\phi$  between them. The phase rotation is independent of the frequency of the subcarrier, and all subcarriers have the same offset,  $e^{j\phi}$ . To overcome the lack of time and phase synchronization, we arbitrarily generate phase shifts through  $e^{-j2\pi f_k \Delta t + j\phi}$  during the offline training process.

**Power Normalization.**  $\mathcal{M}(\cdot)$  is a power normalization function that adjusts the perturbation signal according to  $\epsilon$ , which is the upper bound on the attacker's signal power. We follow the existing power remapping function [7] to preserve the power constraint of the perturbations as follows:

$$\mathcal{M}(\gamma_t^u, \epsilon) = \begin{cases} \sqrt{\epsilon} \frac{\gamma_t^u}{\|\gamma_t^u\|_2}, & \|\gamma_t^u\|_2^2 > \epsilon, \\ \gamma_t^u, & \|\gamma_t^u\|_2^2 \leq \epsilon. \end{cases} \quad (12)$$

where  $\epsilon$  has a different value depending on PSR, which is the ratio of the signal power of the attacker to the signal power of the receiver.  $\gamma_t^u$  is the output of the symbol extension model and symbol shuffling model.

**Transformation Function.** Consequently, we obtain the converted perturbation signal transmitted from the  $k$ -th subcarrier of the  $i$ -th OFDM symbol through the transformation function  $P_{\mu, \zeta, \epsilon, \phi, \Delta t}(\cdot)$  as follows:

$$P_{\mu, \zeta, \epsilon, \phi, \Delta t}(\delta_t^u)[i, k] = \mathcal{M}(\gamma_t^u, \epsilon)[i, k] e^{j\phi} e^{-j2\pi f_k \Delta t}, \quad (13)$$

where  $\gamma_t^u = \Gamma(K(\delta_t^u, \mu), \zeta)$ .

Here, the transformation function is controlled by various parameters  $\mu, \zeta, \epsilon, \phi, \Delta t$ .

## 5.3. Hardware Implementation

Figure 3 shows a scenario in which the adversary sends a perturbation signal behind a thick wall, disrupting communication between the transmitter and the receiver.

**Target Wireless System.** We first implement the ML-based wireless communication system depicted in Figure 2 through USRP B210, a software-defined radio widely used in designing wireless communication system. We drive the



USRP B210 using GNURadio software package [10] that provides a graphical programming interface for configuring transceivers and allows us to model the customized blocks. The transmitter and receiver consist of a USRP B210 and a Linux laptop, respectively, and they communicate through a single antenna, where the carrier frequency is set to 2.4 GHz. The number of cyclic prefixes and subcarriers  $L_{fft}$  is 16 and 64, respectively. Of the 64 subcarriers, 48 are used to carry symbols for ML-based JSCC, 4 of which are used for pilot symbols. We have additionally measured the attack performance at different SNRs. To simulate the signal power of AWGN in the real world, we inject random noise through GNURadio.

**Attack System.** We build an adversarial transmitter using a USRP N310 device with a single antenna and a Linux desktop. We randomly move the antenna to collect 2000 random realizations of the channel  $\{\mathbf{H}_a^l\}_{l=1}^{2000}$  between the adversarial transmitter and receiver. Following the previous work [7], we set the range of PSR to [-20,-10] dB. To perform the black-box attack, we adopt surrogate models with different architectures and parameters from the target wireless communication system and the downstream classifier. We train the PGM offline according to the Algorithm 1. After training is complete, we collect several samples of perturbation signals  $P_\tau(\delta_t^u)$  to facilitate online attacks. In our setup, the attacker feeds random triggers into the PGM to generate a large set of UAPs, and randomly selects UAP adversarial examples from among them. The adversarial transmitter system then loads the stored perturbation signal and transmits it through the OFDM transmitter.

## 6. Attack Evaluation

### 6.1. Experimental Setup

**Victim Models.** We consider four state-of-the-art JSCC models for delivering the multimodal source in the wireless channel. We implemented the JSCC models based on several open source frameworks [69], [78], [82]. Each JSCC model has different weight parameters depending on the variations of the constellation mapping method  $C \in \{\text{QPSK}, 16\text{-QAM}, 64\text{-QAM}\}$  and the coding rate  $\lambda \in \{\frac{1}{6}, \frac{1}{12}\}$ . We train the model with an SNR uniformly distributed in [0,20] dB, to represent a wide range of channel conditions. We also consider scenarios where the receiver applies the demodulated multimodal data to ML-based downstream services, such as video classification and audio-visual event recognition. For the video classification task, we benchmark three state-of-the-art DNN models, namely, I3D [16], SlowFast [23] and TPN [81]. As a benchmark model for audiovisual recognition, we choose the audio-visual networks (AVE) [65] that leverage multi-modality.

**Dataset.** We choose popular multimodal datasets to train and evaluate ML-based JSCC models. For training the image and video JSCC models, we adopt the Vimeo90K dataset [80], which is widely used in evaluating image and video processing tasks. To facilitate efficient training, the

video sequences are cropped to a resolution of  $256 \times 256$ . We then evaluate the image and video JSCC models using the UCF-101 dataset [62]. For the speech JSCC model, we use the speech dataset from Edinburgh DataShare [68], which contains more than 10,000 training data and 800 test data with sampling rate 16 KHz. We truncate the speech sample sequence to have 128 frames with a frame length of 128 after framing. For the text JSCC model, we select the proceedings of the European Parliament, which includes about 2 million sentences and 53 million words. We pre-process the dataset to have sentence lengths between 4 and 30 words. We then split it into training set and test set. We also select widely used datasets as benchmarks to evaluate video classification and audio-visual recognition downstream tasks. We adopt the human action recognition dataset UCF-101 [62] to verify our attack on video classification. For audio-visual recognition, we adopt the audio-visual event dataset [64] which contains 4,143 video clips with 28 event categories.

**Evaluation Metrics.** The goal of JSCC is to minimize the semantic distortion between raw and reconstructed data. We use performance metrics that effectively reflect the semantic information of each modality. For image and video domains, we select the peak signal-to-noise ratio (PSNR) to measure the distortion of the reconstructed frames. In speech domain, mean squared error (MSE) loss reflects the quality of the received speech. For text domain, bilingual evaluation understudy (BLEU) score [9] is widely used to compare the difference between the original sentence and the reconstructed one.

**Baseline Attacks.** We compare Magmaw with state-of-the-art black-box attacks, which are summarized in Table 1. We implement four types of baseline attacks: (1) random attack, (2) single UAP attack [28], (3) multiple UAP attack [7], and (4) RAFA [43]. We design the random attack by generating perturbations that are sampled from a Gaussian distribution. For other attacks [7], [28], [43], we re-implement perturbation signals based on details provided in the papers. Since these studies did not consider the attack on downstream services, we evaluate our attack feasibility in downstream tasks by comparing it to white-box and random attacks.

### 6.2. Experimental Results

We provide two sets of experimental results, where in one, the x-axis is PSR and in the other, it is SNR.

**ML-based Wireless System.** Figure 4 presents the performance of Magmaw in the ML-based wireless transmission systems. We compare the performance of Magmaw to that of the baseline attacks. As shown in Figure 4 (a), Magmaw dramatically deteriorates the performance metrics in the range of all PSRs. Note that “no attack” shows the original performance of the benign model. When applying the adversarial attacks on the image JSCC model, the PSNR drops by up to 8.12dB. For the video JSCC model, PSNR is lowered by 8.43dB on average by Magmaw. We see that the video model is more vulnerable to our adversarial signals than the image JSCC model. The main reason is that the video JSCC model encodes the current frame based on the previously decoded

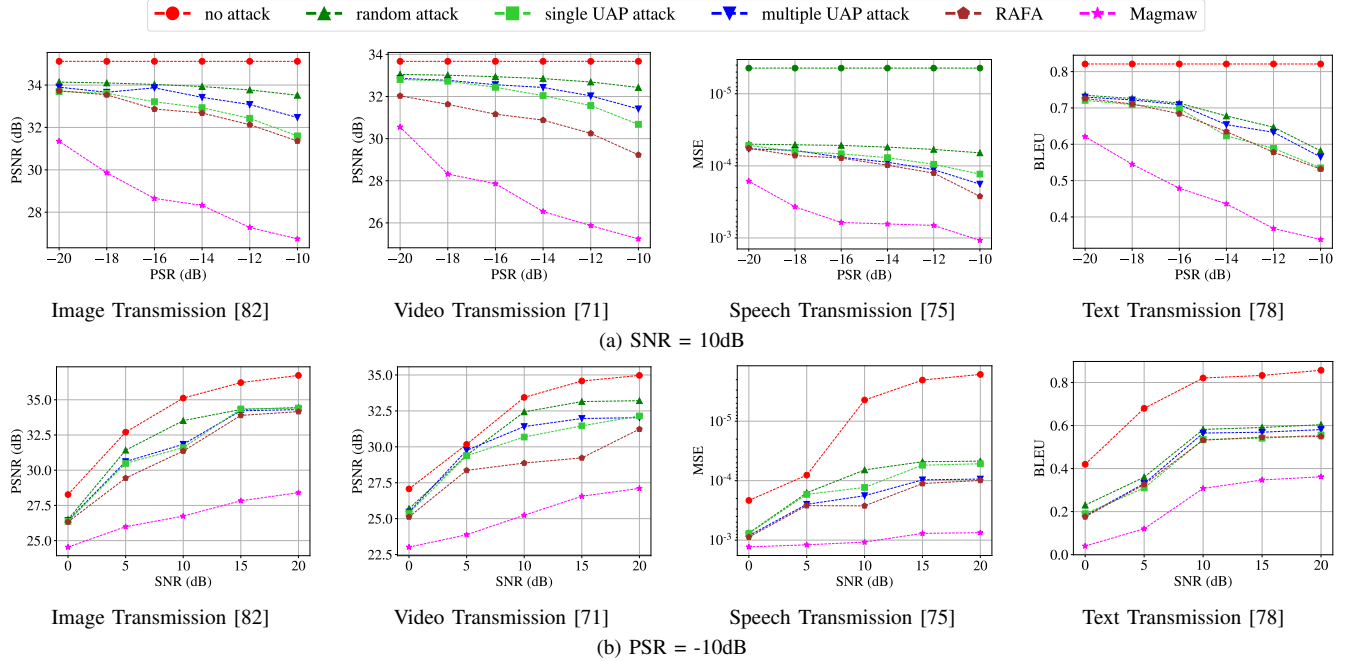


Figure 4: Magmaw on ML-based wireless communication systems.

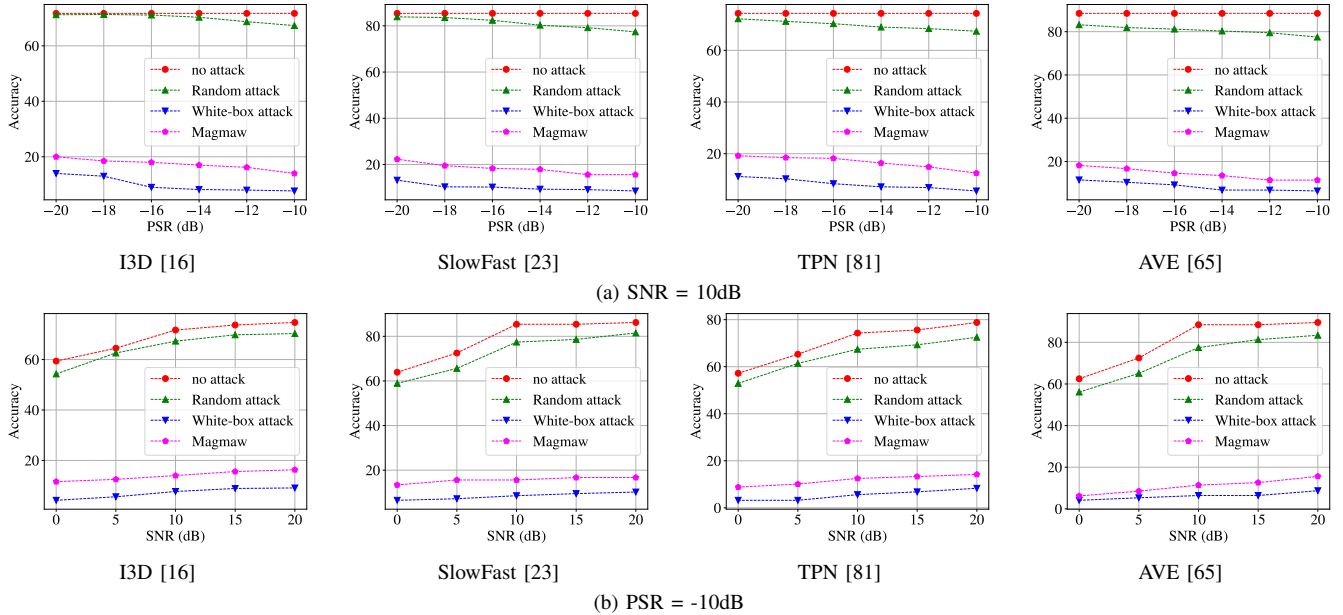


Figure 5: Magmaw on ML-based downstream classification services.

frame, thus propagating the reconstruction distortion to the next frame. For the speech model, we find that Magmaw degrades MSE loss by 4.07 times more than RAFA. We also observe that the BLEU score of the text JSCC model is dropped by up to 0.483 points under Magmaw.

Figure 4 (b) shows the attack performance in the ML-based wireless communication systems for different values of SNR. Here we set the PSR to -10dB. We observe that the highest attack performance can be achieved when the SNR is 0dB. Notably, our attack shows significantly higher attack performance against multimodal wireless transmission com-

pared to the baseline attacks. This is because we train the PGM to generate UAPs that are input- and protocol-agnostic and invariant to time and phase offsets.

In Appendix Figure 12, we further demonstrate the attack results of Magmaw for different constellation mapping methods. We confirm that our attack severely degrades the performance of JSCC models regardless of constellation type. As 64-QAM has slightly higher recovery performance than other modulations (16-QAM, QPSK) in all modalities, we confirm that the higher order of the modulation helps to increase the robustness.

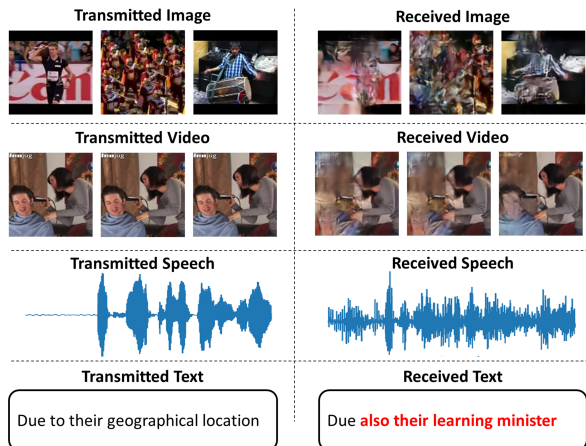


Figure 6: Visualization of attacked multimodal data.

We visualize the attack effect of Magmaw on the multimodal data reconstruction from the receiver in Figure 6. As seen, the JSCC decoder fails to retain semantic information by our perturbation signals. Specifically, the restored images and videos have noise-like artifacts, which dramatically reduce the users’ quality of experience (QoE). Furthermore, the user cannot hear the speaker’s voice in a speech sequence that the speech JSCC decoder failed to recover. The text JSCC decoder generates sentences with incorrect grammar and context, so the user cannot understand the sender’s message.

**Downstream Services.** We evaluate the accuracy of each classifier when Magmaw is directed to a downstream classifier. Then, we provide a comparison with other baseline attacks. Figure 5 shows the attack results for the video classifiers I3D [16], SlowFast [23], and TPN [81] and the audio-visual event classifier AVE [65]. We compare the performance of Magmaw to white-box and random attack scenarios. In the white-box attack scenario, the attacker has complete knowledge of the classification model. Figure 5 (a) presents the accuracy of each baseline for different PSRs. As shown, transmitting randomly sampled perturbations performs very poorly compared to Magmaw. As our attack consistently achieves comparable attack performance compared to the white-box attacks, we confirm that our black-box perturbation signals are successfully transferable to unseen downstream models. Specifically, we achieve an average attack success rate of 81.8% in black box attacks, which is 8.3% lower on average than white-box attacks. In Figure 5 (b), we measure the accuracy of each classifier for different SNRs with the PSR fixed at -10dB. Our experimental results demonstrate that our black-box attacker can subvert downstream services in different SNRs. In particular, we obtain an average attack success rate of 85.2% on the AVE model, which is higher than the video classifiers. Because AVE models receive multiple input modalities, adversarial signals can interfere with larger-dimensional input signals than single-modal classifiers. To analyze the influence of different constellation mapping techniques on the downstream services, we illustrate the attack results on the downstream

---

### Algorithm 2 Adversarial Training against Magmaw

---

**Input:** Dataset  $\mathbb{D}^Q$ , ML-based JSCC model  $\mathcal{J}_{Q,C,\lambda}$ , PGM  $\mathcal{G}$ ,  
**Output:** Robust JSCC model  $\mathcal{J}_{Q,C,\lambda}$   
 $Q \leftarrow$  Modality  
 $C \leftarrow$  Constellation mapping scheme  
 $\lambda \leftarrow$  Coding rate  
Initialize underlying ML-based JSCC model  $\mathcal{J}_{Q,C,\lambda}$   
**for** epoch  $l < \text{MaxIter}$  **do**  
 $\mathbf{H}_t$  is randomly sampled from channel model  
 $\mathbf{H}_a$  is sampled uniformly from training set  
 $\mathbb{B}^{adv} \leftarrow []$   
**for** each batch  $\mathbf{B}^Q \in \mathbb{D}^Q$  **do**  
Train the JSCC model  $\mathcal{J}_{Q,C,\lambda}$  on  $\mathbf{B}^Q$   
 $z_t \sim \text{Uniform}(0, 1)$   
 $\tau_l \leftarrow$  randomly sampled  $\{\mu, \zeta, \epsilon, \phi, \Delta t\}$   
Store  $P_{\tau_l}(\mathcal{G}(z_t))$  in  $\mathbb{B}^{adv}$  for each data in  $\mathbf{B}^Q$   
 $\mathbb{D}^Q.append(\mathbb{D}^Q + \mathbb{B}^{adv})$   
**Return:** Robust JSCC model  $\mathcal{J}_{Q,C,\lambda}$

---

classifiers when different constellation mapping methods are applied to ML-based wireless communication systems in Appendix Figure 13. Although 64-QAM increases accuracy slightly more than other modulations, our input-agnostic attack defeats all modulation techniques.

## 7. Resiliency to Defense Schemes

In this section, we evaluate the performance of the conventional defense mechanisms [7], [43], namely, adversarial training and perturbation subtraction against Magmaw. The performance of the defense algorithms depends on what information the defender knows about the attack formulation. We assume that there is a strong defender who knows the PGM’s model architecture, the channel distribution between the adversary and the receiver, and attack algorithms.

### 7.1. Adversarial Training

The defender aims to obtain a robust ML-based JSCC model for each modality to protect the physical layer from the Magmaw. Since we assume that the defender knows the model architecture of the PGM, adversarial training extends the training dataset to include all adversarial examples and then trains a JSCC model on the augmented dataset. Algorithm 2 shows detailed steps of our adversarial training. We refer to the target JSCC models as  $\mathcal{J}_{Q,C,\lambda}$ , and denote the PGM as  $\mathcal{G}$ , which is identical to the attacker’s model architecture but with different model parameters. The defender trains a ML-based JSCC model by selecting a batch from the training dataset  $\mathbb{D}^Q$  and then generates the adversarial signals controlled by several parameters of the transformation function  $P_\tau$ . We then expand the training dataset to include all adversarial examples and train the model on the augmented training dataset.

**ML-based Wireless System.** We validate Magmaw against the ML-based wireless communication systems, whose resiliency has been improved by adversarial training. As shown in Figure 7, incorporating adversarial examples inside the model training process results in a lower ability to

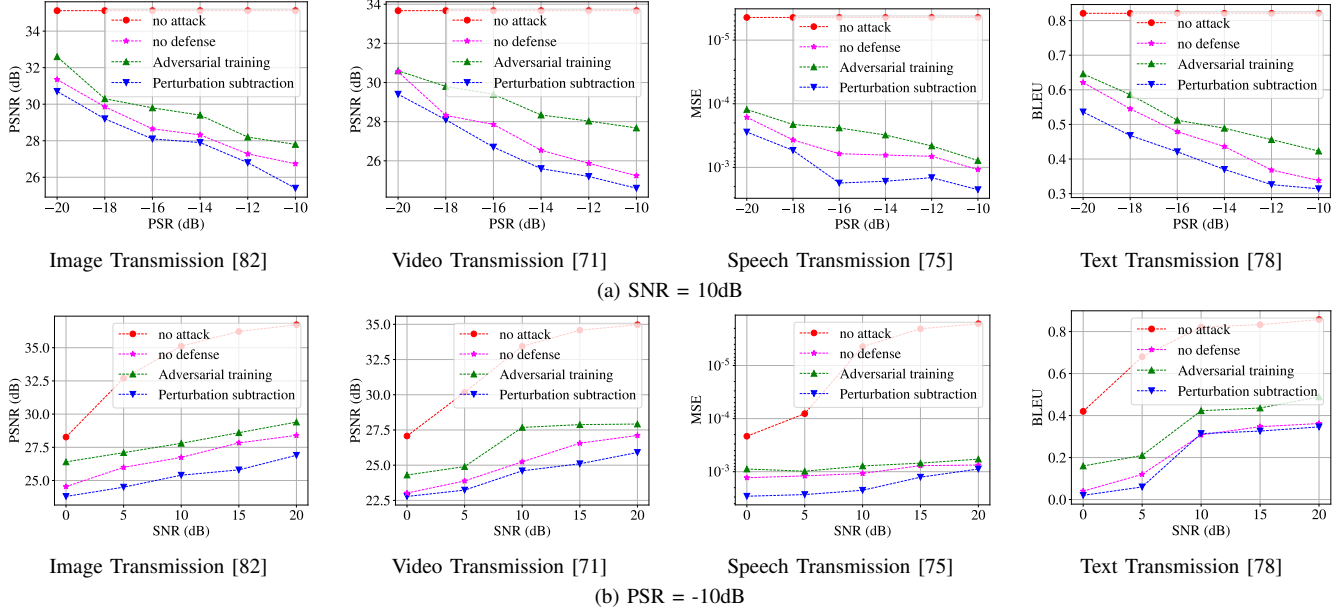


Figure 7: Defense against Magmaw on ML-based wireless communication systems.

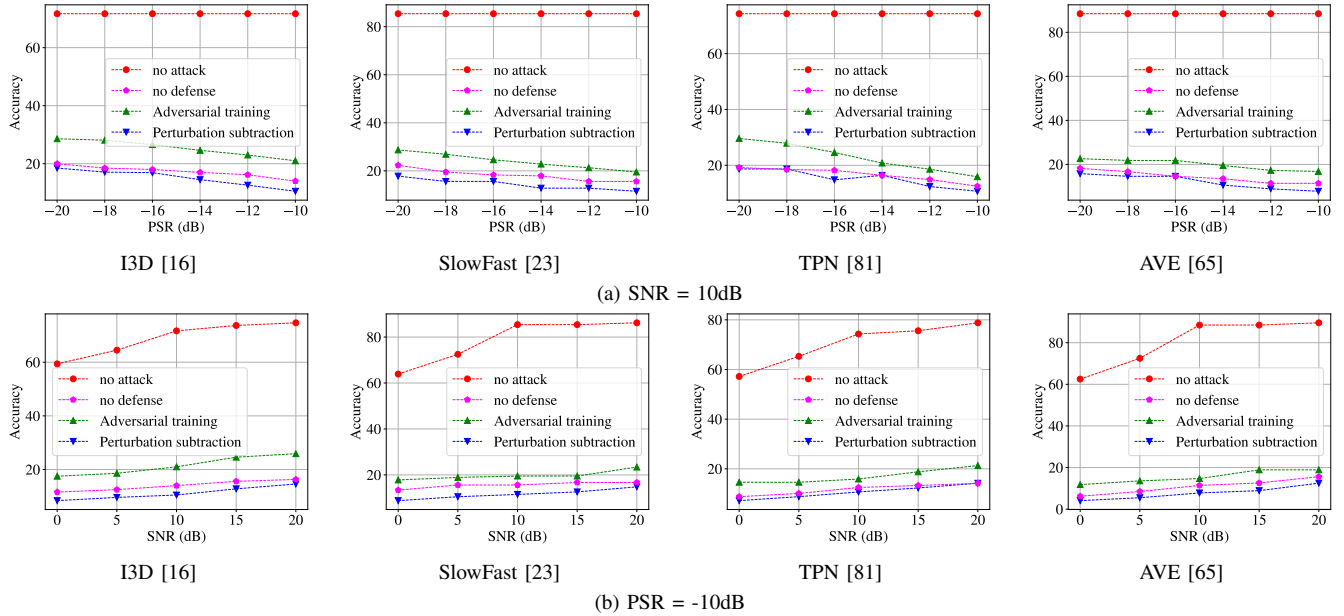


Figure 8: Defense against Magmaw on ML-based downstream classification services.

restore source data even if the underlying victim model is not attacked. Moreover, we observe that adversarial training cannot protect ML-based wireless communication from Magmaw. The reason is that the JSCC model has to be trained on a huge set of perturbations that the defender generates with PGM. Yet it is not feasible for the defender to train JSCC models that are resilient to all possible perturbations. Another reason is that the defender uses a PGM with different parameters from the attacker’s model, so the distribution of adversarial signals generated by the two models is different.

**Downstream services.** Figure 8 shows the classification accuracy of the downstream models trained by adversarial

training. Adversarial training significantly reduces the accuracy of benign models, hindering their applicability. We observe that Magmaw still achieves a high attack success rate even though the benign model undergoes adversarial training. This is because training a model that is universally robust to different types of adversarial signals, while being able to correctly classify multimodal data, is a fundamentally challenging problem.

## 7.2. Perturbation Signal Subtraction

This defense scheme can be performed at the physical layer before the signal is passed through the OFDM receiver.

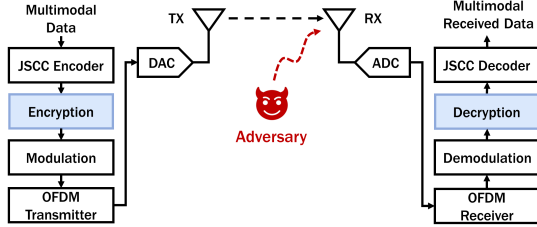


Figure 9: Overview of our attack on secure communication.

Defenders aim to mitigate the effects of adversarial attacks and reconstruct the originally transmitted signal. Thus, the defender takes action on incoming signals that are attacked based on the knowledge of the adversary. As we assume that the defender has knowledge of the adversary’s model architecture, the receiver generates a perturbation signal via the defender’s PGM and then subtracts it from the received wireless signal.

**ML-based Wireless System.** The defense results are summarized in Figure 7. We observe that the source data restored by each JSSC model is more degraded than before the defense. This is because the cancellation of the adversarial signal fails and further amplifies the power of the perturbation. Even if the defender knows the structure of the PGM, the defender cannot generate exactly the same perturbation signal if the model parameters of the PGM are different.

**Downstream services.** As shown in Figure 8, applying perturbation signal subtraction reduces the accuracy of the downstream services by an average of 3.3%. We see that the defender cannot increase the accuracy of the downstream classifier by simply subtracting an estimate of the perturbation. The accuracy of the classifier tends to depend heavily on the quality of the input source.

## 8. Case Studies

### 8.1. Attacks on Encrypted Communication

Encryption schemes are commonly applied in the communication pipeline to protect users’ private data [67]. While the robustness of privacy-preserving communications with ML-based JSSC has not been investigated before, we add the encryption and decryption blocks in our system model to examine the impact of Magmaw on secure communication systems and analyze the vulnerability of encrypted wireless signal (See Figure 9).

**Experiment Setup.** We consider privacy-preserving image transmission [67] with learning with errors (LWE)-based encryption [55], where the transmitter and receiver share a public key, but only the receiver knows the secret key. We set a security level of 192 according to [41]. We do not assume that the adversary has access to the public and the secret keys.

**Attack Results.** Figure 10 (a) shows the performance of the secure communication system with different PSR. As seen, Magmaw lowers the performance of secure image transmission by 2.6dB on average more than RAFA [43].

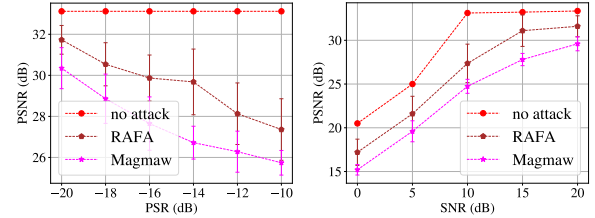


Figure 10: Attack results on secure image communication.

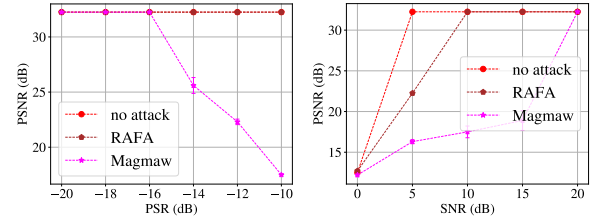


Figure 11: Attack results on conventional image communication.

Furthermore, Figure 10 (b) shows that our attack has an average of 2.4dB higher attack performance compared to RAFA across different SNRs. It can be seen that the OFDM symbols carrying the ciphertext of the image data are vulnerable to our perturbation signal. Our attack is trained to effectively disrupt both the OFDM pilots and data symbols regardless of modality, confirming that it can effectively target radio signals carrying ciphertext as well.

### 8.2. Attacks on Conventional Communication

Current wireless communication standards such as Wi-Fi and 5G follow separate source and channel coding designs that allow independent optimization of each component. We investigate if our UAPs would also be effective for such conventional wireless systems.

**Experiment Setup.** We validate Magmaw against a traditional image transmission system. We adopt the latest BPG as the image codec and LDPC as the channel code, and test with 1/2 LDPC code rates and 16-QAM for OFDM modulation.

**Attack Results.** As shown in Figure 11 (a), the source image cannot be reconstructed when the PSR of Magmaw is greater than -14dB. This is because existing communication systems have a cliff-effect problem in which the source information cannot be recovered if the channel condition deteriorates below the level expected by the channel encoder. However, RAFA can’t subvert existing communications even when the PSR is -10dB. Figure 11 (b) shows that Magmaw is effective even when the SNR is about 10dB higher than where RAFA is effective. This implies that our attack can be transferable to non-ML-based wireless communications.

## 9. Conclusion And Future Work

This paper studies physical-layer adversarial attacks on ML-based wireless systems. We present Magmaw, the first

modality-agnostic attack methodology that generates perturbation signals under the black-box assumption. Magmaw disrupts multimodal data transmission by superimposing a small amount of adversarial signals on victim radio channel, causing the benign receiver to fail to restore the source data. Our results demonstrate that the UAPs generated by Magmaw are feasible in the real world, and can degrade the performance of both target wireless communication and downstream recognition systems simultaneously for multiple modalities. Magmaw can also maintain a high attack success rate by evading defense techniques widely used in the wireless domain.

While Magmaw demonstrates great success in attacking multimodal wireless communications, the perturbation designed in this work needs to be powered by software-defined radios for flexible generation of the perturbation signals. Promising future work is to explore new attack techniques, e.g., intelligent reflecting surfaces [19], [52], [60], to induce rapid phase changes of multimodal radio signals.

## References

- [1] Qualcomm whitepaper vision market-drivers and research directions on the path to 6g. <https://www.qualcomm.com/content/dam/qcomm-martech/dm-assets/documents/qualcomm-whitepaper-vision-market-drivers-and-research-directions-on-the-path-to-6g.pdf>.
- [2] Hadi Abdullah, Muhammad Sajidur Rahman, Washington Garcia, Kevin Warren, Anurag Swarnim Yadav, Tom Shrimpton, and Patrick Traynor. Hear" no evil", see" kenansville": Efficient and transferable black-box attacks on speech recognition and voice identification systems. In *2021 IEEE Symposium on Security and Privacy (SP)*, pages 712–729. IEEE, 2021.
- [3] Najah Abu-Ali, Abd-Elhamid M Taha, Mohamed Salah, and Hossam Hassanein. Uplink scheduling in lte and lte-advanced: Tutorial, survey and evaluation framework. *IEEE Communications surveys & tutorials*, 16(3):1239–1265, 2013.
- [4] Alperen Acemoglu, Jan Kriegelstein, Darwin G Caldwell, Francesco Mora, Luca Guastini, Matteo Trimarchi, Alessandro Vinciguerra, Andrea Luigi Camillo Carobbio, Juljana Hysenbelli, Marco Delsanto, et al. 5g robotic telesurgery: Remote transoral laser microsurgies on a cadaver. *IEEE Transactions on Medical Robotics and Bionics*, 2(4):511–518, 2020.
- [5] Abdullatif Albaseer, Bekir Sait Ciftler, and Mohamed M Abdallah. Performance evaluation of physical attacks against e2e autoencoder over rayleigh fading channel. In *2020 IEEE International Conference on Informatics, IoT, and Enabling Technologies (ICIoT)*, pages 177–182. IEEE, 2020.
- [6] Moustafa Alzantot, Bharathan Balaji, and Mani Srivastava. Did you hear that? adversarial examples against automatic speech recognition. *arXiv preprint arXiv:1801.00554*, 2018.
- [7] Alireza Bahramali, Milad Nasr, Amir Houmansadr, Dennis Goeckel, and Don Towsley. Robust adversarial attacks against dnn-based wireless communication systems. In *Proceedings of the 2021 ACM SIGSAC Conference on Computer and Communications Security*, pages 126–140, 2021.
- [8] Souransu Banerji and Rahul Singha Chowdhury. On ieee 802.11: wireless lan technology. *arXiv preprint arXiv:1307.2661*, 2013.
- [9] Mohamed Ishmael Belghazi, Aristide Baratin, Sai Rajeshwar, Sherjil Ozair, Yoshua Bengio, Aaron Courville, and Devon Hjelm. Mutual information neural estimation. In *International conference on machine learning*, pages 531–540. PMLR, 2018.
- [10] Eric Blossom. Gnu radio: tools for exploring the radio frequency spectrum. *Linux journal*, 2004(122):4, 2004.
- [11] Eirina Boursoulatze, David Burth Kurka, and Deniz Gündüz. Deep joint source-channel coding for wireless image transmission. *IEEE Transactions on Cognitive Communications and Networking*, 5(3):567–579, 2019.
- [12] Yulong Cao, Chaowei Xiao, Benjamin Cyr, Yimeng Zhou, Won Park, Sara Rampazzi, Qi Alfred Chen, Kevin Fu, and Z Morley Mao. Adversarial sensor attack on lidar-based perception in autonomous driving. In *Proceedings of the 2019 ACM SIGSAC conference on computer and communications security*, pages 2267–2281, 2019.
- [13] Yuxin Cao, Xi Xiao, Ruoxi Sun, Derui Wang, Minhui Xue, and Sheng Wen. Stylefool: Fooling video classification systems via style transfer. In *2023 IEEE Symposium on Security and Privacy (SP)*, pages 1631–1648. IEEE, 2023.
- [14] Nicholas Carlini and David Wagner. Towards evaluating the robustness of neural networks. In *2017 IEEE Symposium on Security and Privacy (SP)*, pages 39–57. IEEE, 2017.
- [15] Nicholas Carlini and David Wagner. Audio adversarial examples: Targeted attacks on speech-to-text. In *2018 IEEE security and privacy workshops (SPW)*, pages 1–7. IEEE, 2018.
- [16] Joao Carreira and Andrew Zisserman. Quo vadis, action recognition? a new model and the kinetics dataset. In *proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 6299–6308, 2017.
- [17] Jung-Woo Chang, Mojan Javaheripi, Seira Hidano, and Farinaz Koushanfar. Rovisq: Reduction of video service quality via adversarial attacks on deep learning-based video compression. In *NDSS*, 2023.
- [18] Jianbo Chen, Michael I Jordan, and Martin J Wainwright. Hop-skipjumpattack: A query-efficient decision-based attack. In *2020 IEEE Symposium on Security and Privacy (SP)*, pages 1277–1294. IEEE, 2020.
- [19] Xingyu Chen, Zhengxiong Li, Baicheng Chen, Yi Zhu, Chris Xiaoxuan Lu, Zhengyu Peng, Feng Lin, Wenyao Xu, Kui Ren, and Chunming Qiao. Metawave: Attacking mmwave sensing with metamaterial-enhanced tags. In *The 30th Network and Distributed System Security (NDSS) Symposium*, volume 2023, 2023.
- [20] Yuxuan Chen, Xuejing Yuan, Jiangshan Zhang, Yue Zhao, Shengzhi Zhang, Kai Chen, and Xiaofeng Wang. {Devil's} whisper: A general approach for physical adversarial attacks against commercial black-box speech recognition devices. In *29th USENIX Security Symposium (USENIX Security 20)*, pages 2667–2684, 2020.
- [21] Mostafa Zaman Chowdhury, Md Shahjalal, Shakil Ahmed, and Yeong Min Jang. 6g wireless communication systems: Applications, requirements, technologies, challenges, and research directions. *IEEE Open Journal of the Communications Society*, 1:957–975, 2020.
- [22] Jeremy Cohen, Elan Rosenfeld, and Zico Kolter. Certified adversarial robustness via randomized smoothing. In *international conference on machine learning*, pages 1310–1320. PMLR, 2019.
- [23] Christoph Feichtenhofer, Haoqi Fan, Jitendra Malik, and Kaiming He. Slowfast networks for video recognition. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 6202–6211, 2019.
- [24] Alexander Felix, Sebastian Cammerer, Sebastian Dörner, Jakob Hoydis, and Stephan Ten Brink. Ofdm-autoencoder for end-to-end learning of communications systems. In *2018 IEEE 19th International Workshop on Signal Processing Advances in Wireless Communications (SPAWC)*, pages 1–5. IEEE, 2018.
- [25] Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. *arXiv preprint arXiv:1412.6572*, 2014.
- [26] Jianhua He, Kun Yang, and Hsiao-Hwa Chen. 6g cellular networks and connected autonomous vehicles. *IEEE Network*, 35(4):255–261, 2020.

- [27] Jakob Hoydis, Sebastian Cammerer, Façal Ait Aoudia, Avinash Vem, Nikolaus Binder, Guillermo Marcus, and Alexander Keller. Sionna: An open-source library for next-generation physical layer research. *arXiv preprint arXiv:2203.11854*, 2022.
- [28] Qiyu Hu, Guangyi Zhang, Zhijin Qin, Yunlong Cai, Guanding Yu, and Geoffrey Ye Li. Robust semantic communications with masked vq-vae enabled codebook. *IEEE Transactions on Wireless Communications*, 2023.
- [29] Taewon Hwang, Chenyang Yang, Gang Wu, Shaoqian Li, and Geoffrey Ye Li. Ofdm and its wireless applications: A survey. *IEEE transactions on Vehicular Technology*, 58(4):1673–1694, 2008.
- [30] Jinyuan Jia, Wenjie Qu, and Neil Gong. Multiguard: Provably robust multi-label classification against adversarial examples. *Advances in Neural Information Processing Systems*, 35:10150–10163, 2022.
- [31] Jinyuan Jia, Binghui Wang, Xiaoyu Cao, Hongbin Liu, and Neil Zhenqiang Gong. Almost tight l0-norm certified robustness of top-k predictions against adversarial perturbations. *arXiv preprint arXiv:2011.07633*, 2020.
- [32] Zizhi Jin, Xiaoyu Ji, Yushi Cheng, Bo Yang, Chen Yan, and Wenyuan Xu. Pla-lidar: Physical laser attacks against lidar-based 3d object detection in autonomous vehicle. In *2023 IEEE Symposium on Security and Privacy (SP)*, pages 1822–1839. IEEE, 2023.
- [33] Brian Kim, Yalin E Sagduyu, Kemal Davaslioglu, Tugba Erpek, and Sennur Ulukus. Channel-aware adversarial attacks against deep learning-based wireless signal classifiers. *IEEE Transactions on Wireless Communications*, 21(6):3868–3880, 2021.
- [34] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- [35] Haesoon Lee, Jaeyoung Choi, Dongkyu Kim, and Daesik Hong. Impact of time and frequency misalignments in ofdm based in-band full-duplex systems. In *2017 IEEE Wireless Communications and Networking Conference (WCNC)*, pages 1–6. IEEE, 2017.
- [36] Haoran Li and Wei Lu. Mixed cross entropy loss for neural machine translation. In *International Conference on Machine Learning*, pages 6425–6436. PMLR, 2021.
- [37] Shasha Li, Ajaya Neupane, Sujoy Paul, Chengyu Song, Srikanth V Krishnamurthy, Amit K Roy Chowdhury, and Ananthram Swami. Adversarial perturbations against real-time video classification systems. *arXiv preprint arXiv:1807.00458*, 2018.
- [38] Zeju Li, Xinghan Liu, Guoshun Nan, Jinfei Zhou, Xinchun Lyu, Qimei Cui, and Xiaofeng Tao. Boosting physical layer black-box attacks with semantic adversaries in semantic communications. In *ICC 2023-IEEE International Conference on Communications*, 2023.
- [39] Peng Lin, Qingyang Song, F Richard Yu, Dan Wang, Abbas Jamalipour, and Lei Guo. Wireless virtual reality in beyond 5g systems with the internet of intelligence. *IEEE Wireless Communications*, 28(2):70–77, 2021.
- [40] Xingqin Lin. An overview of 5g advanced evolution in 3gpp release 18. *IEEE Communications Standards Magazine*, 6(3):77–83, 2022.
- [41] Richard Lindner and Chris Peikert. Better key sizes (and attacks) for lwe-based encryption. In *Topics in Cryptology—CT-RSA 2011: The Cryptographers’ Track at the RSA Conference 2011, San Francisco, CA, USA, February 14–18, 2011. Proceedings*, pages 319–339. Springer, 2011.
- [42] Jianwei Liu, Yinghui He, Chaowei Xiao, Jinsong Han, and Kui Ren. Time to think the security of wifi-based behavior recognition systems. *IEEE Transactions on Dependable and Secure Computing*, 2023.
- [43] Zikun Liu, Changming Xu, Emerson Sie, Gagandeep Singh, and Deepak Vasisht. Exploring practical vulnerabilities of machine learning-based wireless systems. In *20th USENIX Symposium on Networked Systems Design and Implementation (NSDI 23)*, pages 1801–1817, 2023.
- [44] Giulio Lovisotto, Henry Turner, Ivo Sluganovic, Martin Strohmeier, and Ivan Martinovic. {SLAP}: Improving physical adversarial examples with {Short-Lived} adversarial perturbations. In *30th USENIX Security Symposium (USENIX Security 21)*, pages 1865–1882, 2021.
- [45] Seyed-Mohsen Moosavi-Dezfooli, Alhussein Fawzi, Omar Fawzi, and Pascal Frossard. Universal adversarial perturbations. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1765–1773, 2017.
- [46] Guoshun Nan, Zhichun Li, Jinli Zhai, Qimei Cui, Gong Chen, Xin Du, Xuefei Zhang, Xiaofeng Tao, Zhu Han, and Tony QS Quek. Physical-layer adversarial robustness for deep learning-based semantic communications. *IEEE journal on selected areas in communications*, 2023.
- [47] Dinh-Luan Nguyen, Sunpreet S Arora, Yuhang Wu, and Hao Yang. Adversarial light projection attacks on face recognition systems: A feasibility study. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition workshops*, pages 814–815, 2020.
- [48] Kuniaki Noda, Hiroaki Arie, Yuki Suga, and Tetsuya Ogata. Multimodal integration learning of robot behavior using deep neural networks. *Robotics and Autonomous Systems*, 62(6):721–736, 2014.
- [49] John Nolan, Kun Qian, and Xinyu Zhang. Ros: passive smart surface for roadside-to-vehicle communication. In *Proceedings of the 2021 ACM SIGCOMM 2021 Conference*, pages 165–178, 2021.
- [50] Patrick O’Reilly, Andreas Bugler, Keshav Bhandari, Max Morrison, and Bryan Pardo. Voiceblock: Privacy through real-time adversarial attacks with audio-to-audio models. *Advances in Neural Information Processing Systems*, 35:30058–30070, 2022.
- [51] Nicolas Papernot, Patrick McDaniel, Ian Goodfellow, Somesh Jha, Z Berkay Celik, and Ananthram Swami. Practical black-box attacks against machine learning. In *Proceedings of the 2017 ACM on Asia conference on computer and communications security*, pages 506–519, 2017.
- [52] Kun Qian, Lulu Yao, Xinyu Zhang, and Tse Nga Ng. Millimirror: 3d printed reflecting surface for millimeter-wave coverage expansion. In *Proceedings of the 28th Annual International Conference on Mobile Computing And Networking*, pages 15–28, 2022.
- [53] Zhijin Qin, Xiaoming Tao, Jianhua Lu, Wen Tong, and Geoffrey Ye Li. Semantic communications: Principles and challenges. *arXiv preprint arXiv:2201.01389*, 2021.
- [54] Erwin Quiring, David Klein, Daniel Arp, Martin Johns, and Konrad Rieck. Adversarial preprocessing: Understanding and preventing {Image-Scaling} attacks in machine learning. In *29th USENIX Security Symposium (USENIX Security 20)*, pages 1363–1380, 2020.
- [55] Oded Regev. On lattices, learning with errors, random linear codes, and cryptography. *Journal of the ACM (JACM)*, 56(6):1–40, 2009.
- [56] Walid Saad, Mehdi Bennis, and Mingzhe Chen. A vision of 6g wireless systems: Applications, trends, technologies, and open research problems. *IEEE network*, 34(3):134–142, 2019.
- [57] Meysam Sadeghi and Erik G Larsson. Physical adversarial attacks against end-to-end autoencoder communication systems. *IEEE Communications Letters*, 23(5):847–850, 2019.
- [58] Sharad Sambhwani, Zdravko Boos, Sidharth Dalmia, Arman Fazeli, Bertram Gunzelmann, Anatoliy Ioffe, Murali Narasimha, Francesco Negro, Laxminarayana Pillutla, and John Zhou. Transitioning to 6g part 1: Radio technologies. *IEEE Wireless Communications*, 29(1):6–8, 2022.
- [59] Takami Sato, Junjie Shen, Ningfei Wang, Yunhan Jia, Xue Lin, and Qi Alfred Chen. Dirty road can attack: Security of deep learning based automated lane centering under {Physical-World} attack. In *30th USENIX Security Symposium (USENIX Security 21)*, pages 3309–3326, 2021.

- [60] Zhambyl Shaikhanov, Fahid Hassan, Hichem Guerboukha, Daniel Mittleman, and Edward Knightly. Metasurface-in-the-middle attack: from theory to experiment. In *Proceedings of the 15th ACM Conference on Security and Privacy in Wireless and Mobile Networks*, pages 257–267, 2022.
- [61] Yulin Shao and Deniz Gunduz. Semantic communications with discrete-time analog transmission: A paper perspective. *IEEE Wireless Communications Letters*, 12(3):510–514, 2022.
- [62] Khurram Soomro, Amir Roshan Zamir, and Mubarak Shah. Ucf101: A dataset of 101 human actions classes from videos in the wild. *arXiv preprint arXiv:1212.0402*, 2012.
- [63] Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian Goodfellow, and Rob Fergus. Intriguing properties of neural networks. *arXiv preprint arXiv:1312.6199*, 2013.
- [64] Yapeng Tian, Jing Shi, Bochen Li, Zhiyao Duan, and Chenliang Xu. Audio-visual event localization in unconstrained videos. In *Proceedings of the European conference on computer vision (ECCV)*, pages 247–263, 2018.
- [65] Yapeng Tian and Chenliang Xu. Can audio-visual integration strengthen robustness under multimodal attacks? In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 5601–5611, 2021.
- [66] David Tse and Pramod Viswanath. *Fundamentals of wireless communication*. Cambridge university press, 2005.
- [67] Tze-Yang Tung and Deniz Gunduz. Deep joint source-channel and encryption coding: Secure semantic communications. *arXiv preprint arXiv:2208.09245*, 2022.
- [68] Cassia Valentini-Botinhao et al. Noisy speech database for training speech enhancement algorithms and its models. *University of Edinburgh. School of Informatics. Centre for Speech Technology Research (CSTR)*, 2017.
- [69] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
- [70] Cheng-Xiang Wang, Marco Di Renzo, Slawomir Stanczak, Sen Wang, and Erik G Larsson. Artificial intelligence enabled wireless networking for 5g and beyond: Recent advances and future challenges. *IEEE Wireless Communications*, 27(1):16–23, 2020.
- [71] Sixian Wang, Jincheng Dai, Zijian Liang, Kai Niu, Zhongwei Si, Chao Dong, Xiaoqi Qin, and Ping Zhang. Wireless deep video semantic transmission. *IEEE Journal on Selected Areas in Communications*, 41(1):214–229, 2022.
- [72] Yang Wang, Zhen Gao, Dezhi Zheng, Sheng Chen, Deniz Gunduz, and H Vincent Poor. Transformer-empowered 6g intelligent networks: From massive mimo processing to semantic communication. *IEEE Wireless Communications*, 2022.
- [73] Yu Wang, Miao Liu, Jie Yang, and Guan Gui. Data-driven deep learning for automatic modulation recognition in cognitive radios. *IEEE Transactions on Vehicular Technology*, 68(4):4074–4077, 2019.
- [74] Sean Welleck, Ilya Kulikov, Stephen Roller, Emily Dinan, Kyunghyun Cho, and Jason Weston. Neural text generation with unlikelihood training. *arXiv preprint arXiv:1908.04319*, 2019.
- [75] Zhenzi Weng and Zhijin Qin. Semantic communication systems for speech transmission. *IEEE Journal on Selected Areas in Communications*, 39(8):2434–2444, 2021.
- [76] IEEE 802.11 working group et al. Wireless lan medium access control (mac) and physical layer (phy) specifications-amendment 2: Enhanced throughput for operation in license-exempt bands above 45 ghz. *IEEE Standard*, 802, 2021.
- [77] Haotian Wu, Yulin Shao, Krystian Mikolajczyk, and Deniz Gunduz. Channel-adaptive wireless image transmission with ofdm. *IEEE Wireless Communications Letters*, 11(11):2400–2404, 2022.
- [78] Huiqiang Xie, Zhijin Qin, Geoffrey Ye Li, and Biing-Hwang Juang. Deep learning enabled semantic communication systems. *IEEE Transactions on Signal Processing*, 69:2663–2675, 2021.
- [79] Shangyu Xie, Han Wang, Yu Kong, and Yuan Hong. Universal 3-dimensional perturbations for black-box attacks on video recognition systems. In *2022 IEEE Symposium on Security and Privacy (SP)*, pages 1390–1407. IEEE, 2022.
- [80] Tianfan Xue, Baian Chen, Jiajun Wu, Donglai Wei, and William T Freeman. Video enhancement with task-oriented flow. *International Journal of Computer Vision*, 127:1106–1125, 2019.
- [81] Ceyuan Yang, Yinghao Xu, Jianping Shi, Bo Dai, and Bolei Zhou. Temporal pyramid network for action recognition. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 591–600, 2020.
- [82] Mingyu Yang, Chenghong Bian, and Hun-Seok Kim. Ofdm-guided deep joint source channel coding for wireless multipath fading channels. *IEEE Transactions on Cognitive Communications and Networking*, 8(2):584–599, 2022.
- [83] Yi Yu, Yufei Wang, Wenhan Yang, Shijian Lu, Yap-Peng Tan, and Alex C Kot. Backdoor attacks against deep image compression via adaptive frequency trigger. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12250–12259, 2023.
- [84] Haiyang Zhang, Nir Shlezinger, Francesco Guidi, Davide Dardari, Mohammadreza F Imani, and Yonina C Eldar. Near-field wireless power transfer for 6g internet of everything mobile networks: Opportunities and challenges. *IEEE Communications Magazine*, 60(3):12–18, 2022.

## Appendix A. Additional Experimental Results

Figure 12 shows the performance of Magmaw on multimodal data transmission when each constellation mapping technique is applied to an ML-based wireless communication system. Figure 13 shows the accuracy of the downstream classifiers when our adversary attacks the wireless communication system to which each constellation mapping technique is applied.

## Appendix B. Related Work

**Adversarial Attacks on Other Domains.** Adversarial attacks have been studied to validate and analyze the robustness of the ML models across multiple domains of computer vision, such as image classification [14], [18], [30], [31], [51], [54], speech recognition [6], [15], [50], human activity recognition [13], [37], [79], neural compression [17], [83], etc. Most of these studies provide an attacker with capabilities to perform a man-in-the-middle attack where she can intercept data in the middle and then injects a small amount of adversarial perturbation. Therefore, the perturbation is not physically feasible and only exposes theoretical vulnerabilities. As the demand for physically feasible adversarial research grows, recent studies [2], [12], [20], [32], [44], [47], [59] define attack methodologies more practically so that adversarial attacks can be realized in the real world. For example, the adversary superimpose light from a projector onto an object, causing the object detector to misclassify the



object [44], [47]. However, compared to wireless attacks, physical attacks in the vision domain are less susceptible to signal distortion and have fewer protocols to challenge.

**Adversarial Attacks on Wireless Systems.** There are several studies [5], [7], [28], [38], [43], [46], [57] that have analyzed the vulnerability of adversarial attacks against ML-based wireless communication systems. The most recent studies [7], [43] have tried to make the attack viable by limiting the adversary’s abilities. Bahramali *et al.* [7] presents a DNN-based attack mechanism to generate UAP signals. However, they make the unrealistic assumption that the target wireless system transmits only one-hot vector messages. Furthermore, they do not take into account the physical layer components, which are traditionally used in end-to-end wireless communication, thus their attacks are not practical in the real world. Liu *et al.* [43] conducts adversarial attacks against ML-based wireless systems for channel prediction and indoor localization. As they aim to contaminate the pilot signals, their perturbations are not suitable for attacking our target system, end-to-end wireless transmission. Furthermore, their attack method does not take into account multi-modality of the source data and different settings of wireless protocols (constellation, coding rate, and OFDM specifications).

In this paper, we present, Magmaw, a black-box attack framework to generate universal multi-modal perturbations that target both OFDM pilots and data symbols. We show, for the first time, that modulated multimodal data can be exploited by adversaries, resulting in subverting downstream services without restoring the original data.

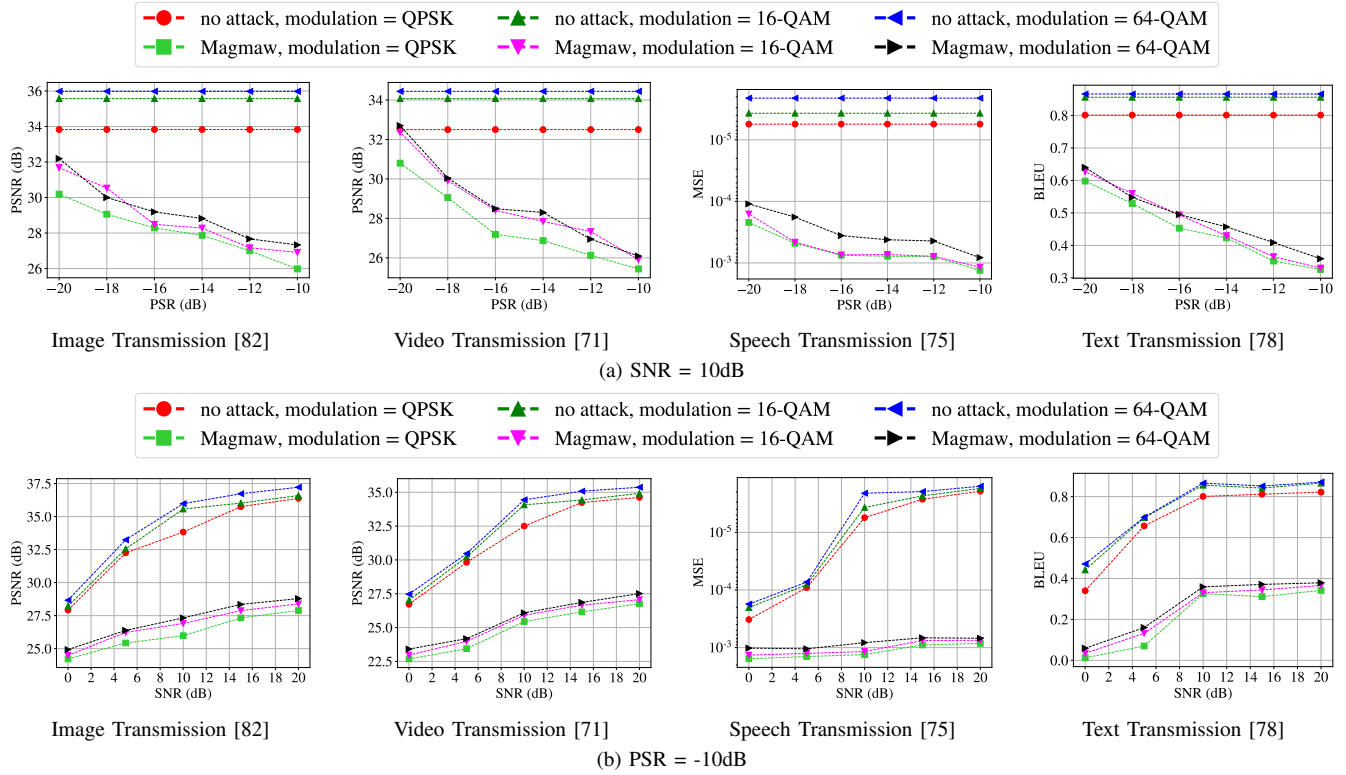


Figure 12: Magmaw on ML-based wireless communication systems with different types of constellation mapping schemes.

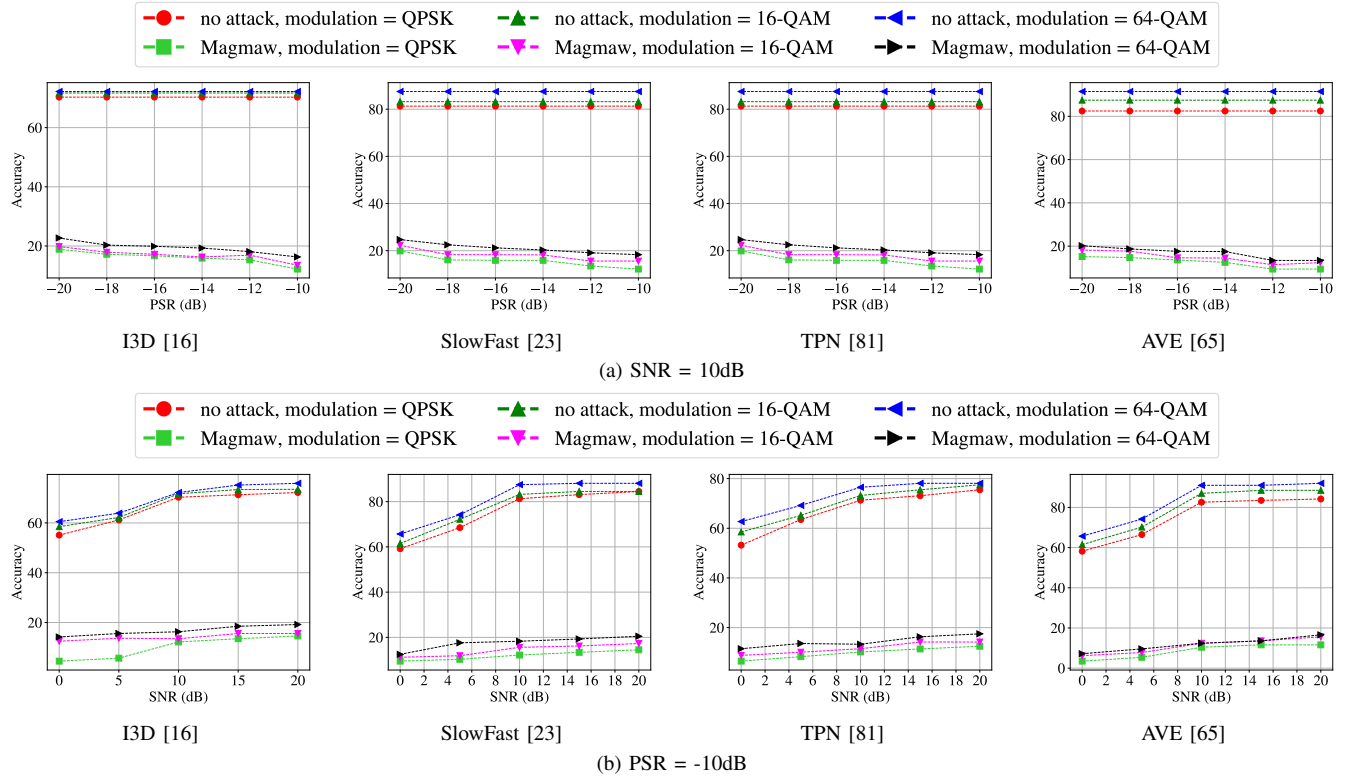


Figure 13: Magmaw on ML-based downstream classification services with different types of constellation mapping schemes.