# UCLA
## UCLA Electronic Theses and Dissertations

**Title**

Application-Tailored Security: Lessons from Theory to Practice

**Permalink**

https://escholarship.org/uc/item/7r71w6j5

**Author**

Karmoose, Mohammed

**Publication Date**

2019

Peer reviewed|Thesis/dissertation

UNIVERSITY OF CALIFORNIA

Los Angeles

Application-Tailored Security:

Lessons from Theory to Practice

A dissertation submitted in partial satisfaction

of the requirements for the degree

Doctor of Philosophy in Electrical and

Computer Engineering

by

Mohammed Karmoose

2019

ABSTRACT OF THE DISSERTATION

Application-Tailored Security:

Lessons from Theory to Practice

by

Mohammed Karmoose

Doctor of Philosophy in Electrical and

Computer Engineering

University of California, Los Angeles, 2019

Professor Christina Fragouli, Chair

With the increase of inter-connected devices, it is of paramount importance to ensure the security of the exchanged information. While cryptographic techniques provide tools to provide confidentiality and data integrity, such techniques may not provide the most efficient solutions. In addition, applications today have challenging performance requirements, with the emergence of time-critical applications such as vehicle networks, as well as resource-limited inter-connected devices such as the devices used in the Internet-of-Things. For such applications, novel security solutions that are application-tailored are needed to meet the performance requirements while adhering to the constraints imposed by the available resources. In this thesis, we adopt this methodology for security design: by understanding the nature of the application, the possible adversaries that may target the communication system, as well as the performance requirements, we design suitable and efficient security solutions. We show this methodology in the context of three different scenarios.

The first scenario is data broadcasting in the context of the index coding problem. We study the problem of providing privacy guarantees against curious clients who are interested in knowing the requests and side information sets of other clients. We first design index codes with higher privacy levels than conventional index codes. We also provide a mechanism, which we call $k$-limited-access schemes, which transforms any index coding technique into

another code with higher privacy guarantees.

The second scenario is in the context of communication systems which relies on millimeter waves. We tackle the problem of secret key establishment. We propose a secret key establishment protocol which allows two communicating parties to establish shared secret keys at very high rates. We showcase the performance of our proposed technique in two different applications: in millimeter wave wireless systems such as 5G networks and IEEE 802.11ay, and vehicle platooning.

The last scenario is in the context of Cyber-Physical Systems. We first argue that, in many situations, an adversary is interested in learning the state vector of the control system. In such cases, a more suitable security metric would be a distortion-based one which leads the adversary to make state estimates that are far from the actual value. We then propose security schemes that require a very small number of secret key bits and still perform well according to the proposed metric: we show that our proposed schemes are in fact optimal for many cases.

The dissertation of Mohammed Karmoose is approved.

Wotao Yin

Danijela Cabric

Suhas Diggavi

Christina Fragouli, Committee Chair

University of California, Los Angeles

2019

*To my parents . . .*

*I am what I am because of you*

TABLE OF CONTENTS

# ACKNOWLEDGMENTS

| | |
|---|---|
| 2009 | B.S. (Electrical and Electronic Engineering), Alexandria University, Egypt |
| 2009 – 2014 | Teaching Assistant, Electrical and Electronic Engineering, Alexandria University, Egypt |
| 2013 | M.S. (Electrical Engineering), Alexandria University, Egypt |
| 2014 | Joined Ph.D program at the ECE department, UCLA (UCLA graduate fellowship) |
| 2016 | Teaching Assistant, ECE Department, UCLA |
| 2017 | Advanced to Ph.D candidacy |
| 2017 | Research Intern, Security and Privacy Research Group, Intel Labs |
| 2018 | Teaching Assistant, ECE Department, UCLA |

PUBLICATIONS

**M. Karmoose**, M. Cardone and C. Fragouli, "Simplifying wireless social caching," in *the 2016 IEEE International Symposium on Information Theory (ISIT), Barcelona*, 2016, pp. 410-414.

**M. Karmoose**, L. Song, M. Cardone and C. Fragouli, "Private broadcasting: an index coding approach", in *the 2017 IEEE International Symposium on Information Theory (ISIT), Aachen*, 2017, pp. 2543-2547.

**M. Karmoose**, L. Song, M. Cardone and C. Fragouli, "Preserving privacy while broadcasting: K-limited-access schemes". in *the 2017 IEEE Information Theory Workshop (ITW), Kaosiung*, 2017, pp. 514-518.

Y. Ezzeldin, **M. Karmoose** and C. Fragouli, "Communication vs distributed computation: an alternative trade-off curve", in *the 2017 IEEE Information Theory Workshop (ITW), Kaosiung*, 2017, pp. 279-283.

**M. Karmoose**, L. Song, M. Cardone and C. Fragouli, "Privacy in index coding: Improved bounds and coding schemes", in *the 2018 IEEE International Symposium on Information Theory (ISIT), Vail*, 2018, pp. 831-835.

**M. Karmoose**, M. Cardone, and C. Fragouli, "Simplifying Wireless Social Caching via Network Coding", in *the IEEE Transactions on Communications* 66, no. 11, 2018, 5512-5525.

**M. Karmoose**, L. Song, M. Cardone and C. Fragouli, "Privacy in index coding: k-limited-access schemes", submitted to *the IEEE Transactions on Information Theory*, 2018.

G. K. Agarwal, **M. Karmoose**, S. Diggavi, C. Fragouli and P. Tabuada, "Distorting an adversary's view in cyber-physical systems", in *the 2018 IEEE Conference on Decision and Control (CDC), Florida*, 2018, pp. 1476-1481.

**M. Karmoose**, G. K. Agarwal, S. Diggavi, C. Fragouli and P. Tabuada, "Distorting an adversary's view in cyber-physical systems", submitted to *the IEEE Transactions on Automatic Control*, 2019.

**M. Karmoose**, C. Fragouli, S. Diggavi, R. Misoczki, L. L. Yang and Z. Zhang, "Using mm-Waves for Secret Key Establishment", in *the IEEE Communications Letters*, 2019.

# CHAPTER 1

# Introduction

Communication systems have been increasingly prevalent in our lives. We are constantly using our connected device to communicate, share information, store data on cloud services, perform financial transactions, and remotely connecting and controlling other devices in our smart homes and vehicles. The emergence of social networks and social media has also influenced us to share a great deal of our daily lives with others.

This exchange of information is performed over a worldwide network of interconnected devices. Typically, information exchange between two devices occur by routing information through other devices and/or sub-networks of devices, most of which are intended to act as intermediate relaying nodes without actually reading or altering the information. This setup creates opportunities for malicious nodes to read and/or alter information in an unauthorized manner. Naturally, this poses a major security concern which could lead to catastrophic consequences if not adequately addressed.

Security vulnerabilities in communication systems create the possibility to be exploited in harmful ways. [Arm] shows the biggest data breaches of the 21st century and the resultant financial losses, which goes to the order of billions of dollars. The recent events of data breaches of Facebook [IF], Apple [Art] and Google [Goo] show that, in addition to corporate-level financial losses, such vulnerabilities can have harmful consequences on individuals [Lar]; in many situations, these consequences can be fatal, as is the case with the recent cyber-attacks on autonomous vehicles [Gar, Gre].

It is due to the importance of Cybersecurity that many initiatives have been established to understand security threats and provide methodological solutions to current information systems. These initiatives include the ETSI Cyber Security Technical Committee [ets], the ISO

information security management standards [iso], and NIST Cybersecurity framework [nis]. While information systems naturally differ in their security requirements, three common security objectives are generally targeted when developing any security solution. These objectives are *confidentiality*, *integrity* and *availability*; commonly referred to as the *CIA Triad*. *Confidentiality* maintains that exchanged information is only observed/consumed by the intended receiver, with no unauthorized access by other parties. *Integrity* ensures the information received by the intended receiver is unchanged with respect to when it was generated. Finally, *availability* is to make sure that the system is continuously delivering the needed amount of information from sources to destination, even in the presence of possible malicious actions.

## 1.1 Role of Cryptography – Can we always rely on it?

An obvious question arise: "can cryptographic tools provide an adequate security solution?" Indeed, there exist cryptographic primitives which promise to provide security guarantees in accordance to the CIA triad. Specifically, cryptography provides encryption schemes such as Integrated Encryption Scheme (IES) [Sti05] and Advanced Encryption Scheme (AES) [Sti05] which ensures that only receivers with the right secret key can decrypt and read the data, thus providing confidentiality. In addition, data integrity can be maintained by the use of authentication mechanisms such as Message Authentication Codes (MAC) [Sti05] and Digital Signature Algorithms (DSA) [Sti05]. While these tools are commonly used to ensure security in many of today's application, I would like to argue that they do not always provide the most efficient solutions. To better illustrate, a more careful look into these techniques is required[1]. These techniques can be generally classified into two categories based on the type of secret keys used for encryption: *symmetric key encryption* (*e.g.*, AES and MAC) in which the same key is used for encryption by the source and decryption by the destination, and *asymmetric key encryption* (*e.g.*, IES and DSA) where different keys are used for

---

[1]Most of the following discussion can be readily found in most modern cryptography textbooks, *e.g.*, [Sti05, KL14]

encryption and decryption. Symmetric-key-based techniques are generally more computationally efficient than their public-key-based counterparts. However, the establishment of shared symmetric keys between communicating parties is usually considered as a burden, especially in the context of dynamic networks which change rapidly. In contrast, the dynamic and efficient distribution of private/public key pairs among communicating nodes is usually possible through the establishment of dedicated security facilities in the network (what is called *public key infrastructure*). In this case, public-key-based schemes are used despite the added penalty in computation efficiency.

This brief exposition of cryptographic techniques immediately shows that, even within the realms of cryptography, there is no "one-size-fits-all" security solution to communication systems; a concept that is already realized in many of today's security initiatives [nis]. More specifically, dynamic networks require security solutions that are scalable, which bodes well with a public-key infrastructure. However, this requires the availability of trusted infrastructure that establishes and distributes keys across nodes, and it uses computationally-heavy encryption and verification techniques. This is in contrast to symmetric-key-based schemes, which are highly efficient yet not as scalable. In the following, I provide two examples which highlight this particular concept.

### 1.1.1   Example 1 - Autonomous Vehicles and Safety Messages

Autonomous vehicles have gained considerable attention from both industry and academic communities. The main objective is to rely on the processing and sensing capabilities of modern vehicles to (semi)-automate the driving process, further enhancing the driving experience of passengers as well as providing benefits in terms of traffic management and road safety [MGL16]. However, in order to ensure the safety of such technologies, autonomous vehicles on the road need to obtain essential information about surrounding vehicles. This is achieved by mandating vehicles to exchange data packets known as *Basic Safety Messages (BSMs)*, which is a part of the IEEE Standard for Wireless Access in Vehicle Environments (WAVE) [IEE16a]. A key property that is imperative for the safety of autonomous driving

3

is that BSMs are received within a 100-ms latency constraint [IEE16a]. The security of BSMs are of paramount importance for the safety of autonomous vehicles, and therefore a security amendment has been proposed for the WAVE standard [IEE16b]; the standard proposes the use of public-key infrastructure for key establishment and the use of digital signatures (namely ECDSA) for message authenticity verification, among other cryptographic primitives.

However, the use of the aforementioned schemes pose a challenge as to adhere to the stringent latency requirement of BSMs. Specifically, typical off-the-shelf WAVE modules which implement the WAVE protocol stack report an approximately 50-ms delay to perform ECDSA verification of one packet [Une]. This means that, in a 100-ms latency window, one such WAVE module is only able to process up to 2 BSM packets from 2 different cars; a situation that is unlikely reflecting of real world environment. Efforts are currently being made to 1) provide more efficient implementations of ECDSA creation/verification schemes, and/or 2) provide dynamic and scalable security architecture which does not rely on the less-efficient public-key-based cryptographic schemes.

### 1.1.2  Example 2 - Vehicle Platooning

Vehicle platooning is another example of autonomous driving in which vehicles form tightly-packed platoons on the road. This promises to provide great fuel savings especially for trucks on highways [AGJ10], as well as better traffic conditions [FN12]. To realize such benefits, platoons are envisioned to have inter-vehicle spacings as small as 10 meters. In such a critical situation, vehicles in a platoon has to be highly adaptable to changes that might suddenly occur on the road which would lead to a change in driving behavior of leading or trailing cars in a platoon. This can only be achieved by equipping vehicles with high sensing capabilities, as well as inter-vehicle communications. Apparently, trust has to be established between platoon vehicles (*e.g.*, by distributing secret keys) in order to prevent unauthorized vehicles from adversely affecting communication. A typical way of establishing such trust is via the use of a public-key infrastructure, as suggested by the security amendment for IEEE

WAVE standard [IEE16b]. However, this trust has to be established while accounting for the dynamic nature of platoons (platoons are expected to form and/or disband frequently to account for changes on the road). Considering that platoons (especially of trucks) are expected to be on highways for most of the time, providing the necessary infrastructure for a public-key infrastructure security solution can be a challenge. Therefore, other mechanisms have to be provided for an efficient trust establishment that does not highly depend on trusted infrastructure – one solution replaces the use of trusted infrastructure and private/public keys with the use of symmetric keys and the concept of *trusted platoon leader* to establish trust between current and new platoon vehicles [KMY19].

### 1.1.3   Example 3 - Wearable Health Monitoring Devices

A typical wearable health monitoring platform consists of an array of sensors connected to a person's body, which collect physical and chemical information about his/her health status in real-time. [GEN16]. The sensor array is coupled with a communication module which enables the transmission of the collected information to a central hub responsible for data collection and analysis – the collection of the sensor array and communication module, along with additional circuitry, is what constitutes a wearable device [HBH17]. As this collected data contains sensitive information about the person's health status in an extended period of time, it is imperative to equip the wearable device with a sufficiently strong encryption schemes [HBH17]. However, a typical wearable sensing device is battery-operated. In addition, due to the size limitation, these devices are usually equipped with small batteries that come with limited power capacities [GEN16]. To enable data collection through different sensors and seamlessly perform data analysis in an extended period of time, power consumption has to be handled very efficiently. With the foreseen increase in the amounts of data collection and transmission, it is not clear how typical cryptographic techniques can be utilized with the required power efficiency in terms of encrypting and validating the data.

## 1.2  Design Methodology

The aforementioned discussion emphasizes that a "one-time-fits-all" solution approach to security is not always fruitful. Although cryptography provides useful general tools for securing communication, an application-based security design is in many cases required. Such a design methodology would require the following:

1. An understanding of the <u>adversaries</u> and <u>threats</u> that may impact the underlying communication system;

2. An understanding of the nature of the application in terms of the <u>available resources</u> and/or <u>performance requirements</u>; and

3. Choosing the right <u>security metric</u> based on which the security system is designed and optimized.

### 1.2.1  Adversaries and Threats

It is important to have a precise understanding of the adversaries and threats against which a security solution should protect. Generally, adversaries can be classified into passive (*i.e.*, can listen to communication but not maliciously tamper with it) and active (*i.e.*, can actively tamper with communication). Although the focus of the thesis is on passive adversaries, more information about the adversary can help provide an effective solution. For example,

1. if the adversary is an outsider to the system, then it may be possible for the legitimate communicating parties to establish a secure communication channel via shared keys, whereas it may not be possible/efficient to do so if the adversary is an honest-but-curious client in the system – this happens to be the case when we discuss our work about data privacy in broadcasting domains in Chapters 2 and 3,

2. the adversary may be eavesdropping communication to extract information, but its end goal is not to learn the communicated data itself, but rather a particular function of it. In this case, it may be easier/more efficient to design a security scheme with the

6

target of preventing the adversary from learning this information, rather than targeting the protection of the data itself – this is a key observation that lead to our proposed security solution for Cyber-physical systems which we talk about in Chapter 4,

3. the capabilities of the adversaries is crucial. An adversary with limited computational power can be prevented by using typical cryptographic techniques, whereas a quantum adversary may not be. In this case, alternative schemes can be used for protection (*e.g.*, physical layer security). While such schemes may be protective against computationally-capable adversaries, they are affected by other properties of the adversaries which are also essential to understand (*e.g.*, the network presence of the adversary in the case of physical layer security schemes) – we touch more on this in discussing security applications for millimeter waves in Chapter 3.

### 1.2.2 Available Resources and Performance Requirements

Understanding the performance requirements is essential for the design of a suitable security scheme. As was discussed with Example 1, using a computationally-expensive scheme for data verification is not suitable to meet the stringent latency requirement of the application. On the other hand, the nature of the application provides the system designer with resources to use in designing the security solution. This was the case with Example 2, where the compact nature of the platoon allows for the use of symmetric key encryption techniques. We also show how this manifests in the application of data privacy in broadcasting domains (Chapters 2 and 3) and in the use of milli-meter waves for communication (Chapter 4).

### 1.2.3 Security Metric

The right choice of a security metric influences the design of a suitable and effective security scheme. The right metric is also dependent on the nature of the application. In cases where an adversary is present with unlimited computational capabilities, an information-theoretic security metric could be a suitable choice. We also discuss a variety of information-theoretic metrics when tackling private information broadcasting in Chapters 2 and 3 (*e.g.*, conditional

entropy of requests and side information, maximal information leakage). Contrary, we show that a distortion-based metric can be useful in applications such as Cyber-physical systems, which we elaborate on more in Chapter 5.

## 1.3  Contributions

The focus of this thesis is to use the previously discussed application-tailored security in designing security solutions for different applications. Namely, in this thesis, I showcase how this methodology can be used in three different scenarios:

1. The first application is in the context of data broadcasting. I consider the problem of index coding where a server delivers data requested by a set of clients over an error-free broadcast channel. By knowing the requests and the side information (*e.g.*, previously known messages) of the clients, the server uses network-coded transmissions to simultaneously deliver the requested data in a small number of transmissions. In this application, we consider an <u>adversary</u> who is a curious client in the system, and the <u>threat</u> it poses is to gain information about the requests and/or side information sets of other clients in the channel. We show how the advesray can gain such information by having access to the broadcast transmission. Since the adversary is also a legitimate client with its own requested data, giving it access to the transmissions is unavoidable and therefore it is unclear how cryptographic approaches are useful in this context. However, we show that the nature of the application shows different <u>resources</u> to be used for security: the unknown requests/side information of other clients, as well as the number of transmissions. We first show that the server can indeed design index codes that tradeoff the privacy of one quantity (the requests or side information) at the expense of the privacy of the other. In addition, we show how we can transform index codes into other codes that provide more privacy for both quantities by using additional transmissions. This is the main topic of Chapters 2 and 3 of this thesis.

2. The second application is in the context of wireless communication systems which use

milli-meter Waves (mmWaves). We consider here the problem of secret key estab-lishment between two communicating parties. With no wired or secured connection between the two parties, key establishment needs to be perfomed over the wireless channel. This creates the possibility of an adversary who has access to the wireless medium to pose the threat of overhearing the secret key being shared between the two parties. With the emergence of quantum computers, it is reasonable to assume that the adversary has unlimited computational power, and therefore traditional key establishment techniques, *e.g.*, Diffie-Hellman (DH), can be broken [SKK18]. However, the nature of the application provides an excellent resource to be used for securing the process. Namely, communication over mmWaves is required to be highly directional in order to overcome the large signal attentuation. This directionality, if used properly, can be used as a countermeasure against an adversary with wireless access. We show how the right secret key generation protocol allows us to generate a significantly higher rate of secret keys than existing key establishment mechanisms. This is the main topic of Chapter 3 of this thesis.

3. The final application is in the context of Cyber-Physical Systems (CPSs). Namely, we consider a typical scenario where an agent is required to transmit information about its current state vector. An adversary has access to the transmission medium and therefore poses the threat to retrieve this information. While general cryptographic techniques may be suitable in some of such scenarios, consider the case where the agent represent a CPS with low computation capabilities such as a small Internet-of-Things (IoT) device. In this case, a computationally-heavy process such as cryptographic encryption technique may not be feasible. We show in this case that a better understanding of the adversary can lead to a better-suited security solution. Specifically, we argue that, in many cases, the adversary is interested in the value of the state vector. However, confusing the adversary into estimating a state vector value that is "far" from the actual value is in many cases sufficient. In this case, we propose a new distortion-based security metric, and we show that more efficient schemes (than typical cryptographic techniques) can be developed to optimize for this metric. This is the main topic of

9

Chapter 4 of this thesis.

## 1.4 Outline of the thesis

The thesis is organized as follows. Chapter 2 discusses the problem of private broadcasting in index coding, where index codes are designed which provide privacy guarantees for the requests and side information sets. Chapter 3 discusses the use of $k$-limited-access schemes to increase the privacy levels of existing index codes. Chapter 4 shows how millimeter waves can be leveraged to provide better secret key establishment solutions. Chapter 4 discusses the use of distortion-based security metrics for cyber-physical systems and how efficient security mechanisms are designed. We finish this thesis by a conclusion.

## 1.5 Notation

Calligraphic letters indicate sets; $|\mathcal{X}|$ is the cardinality of $\mathcal{X}$; by $[m]$ we denote $\{1, 2, \ldots, m\}$ where $m \in \mathbb{Z}^+$; and by $[m_1 : m_2]$ we denote $\{m_1, m_1 + 1, \ldots, m_2\}$ where $m_1, m_2 \in \mathbb{Z}^+$ and $m_2 > m_1$; $2^{[n]}$ and $\binom{[n]}{s}$ are the power set and the set of all possible subsets of $[n]$ of size $s$, respectively; for a sequence $X = \{X_1, \ldots, X_n\}$, $X_{\mathcal{S}}$ is the subsequence of $X$ where only the elements indexed by $\mathcal{S}$ are retained; boldface lower case letters denote vectors and boldface upper case letters indicate matrices; For a matrix $\mathbf{A}$, we denote by $\mathbf{A}'$ the transpose of $\mathbf{A}$; given a vector $\mathbf{b}$, $b_i$ indicates the $i$-th element of $\mathbf{b}$; given matrices $\mathbf{A}$ and $\mathbf{B}$, $\mathbf{B} \subset_k \mathbf{A}$ indicates that $\mathbf{B}$ is formed by a set of $k$ rows of $\mathbf{A}$; $\mathbf{A}_{\mathcal{S}}$ is the submatrix of $\mathbf{A}$ where only the columns indexed by $\mathcal{S}$ are retained; let $\mathbf{x}_i$ be a set of column vectors indexed by $i$, then $\mathbf{x}_a^b = [\mathbf{x}_a' \; \mathbf{x}_{a+1}' \; \cdots \; \mathbf{x}_b']'$ for $b \geq a$ and $a, b \in \mathbb{Z}$; $\text{span}(\mathbf{A})$ is the linear span of the columns of $\mathbf{A}$; $\mathbf{0}_j$ is the all-zero row vector of dimension $j$; $\mathbf{0}_{i \times j}$ is the all-zero matrix of dimension $i \times j$; $\mathbf{1}_j$ denotes a row vector of dimension $j$ of all ones – sometimes alternatively denoted by $\mathbf{I}$ while the size is understood from context; $\mathbf{I}_j$ is the identity matrix of dimension $j$; $\mathbf{e}_i^j$ is the all-zero row vector of length $j$ with a 1 in position $i$; $\mathbf{Pr}(X)$ refers to the probability of event $X$; $H(X|y)$ is the entropy of the random variable $X$, conditioned on the *specific*

realization $y$; $f_X(x)$ denotes the probability density function of a random vector $X$, which can be alternatively denoted by $f(x)$ for brevity; for any random vector $Y$, we denote the mean and covariance matrices of $Y$ by $\mu_Y$ and $R_Y$ respectively; for all $x \in \mathbb{R}$, the floor and ceiling functions are denoted with $\lfloor x \rfloor$ and $\lceil x \rceil$, respectively; $\binom{n}{k} = 0$ if $k < 0$ or $k > n$; finally, logarithms are in base 2.

# CHAPTER 2

# Private Broadcasting: an Index Coding Approach

The first example of our application-tailored security design is in the context of index coding. In the traditional index coding problem, a server employs coding to send messages to $n$ clients within the same broadcast domain. Each client already has some messages as side information and requests a particular unknown message from the server. This setup is an abstraction which captures the essence of many data broadcasting scenarios []. In index coding, all clients learn the coding matrix so that they can decode and retrieve their requested data. Our starting observation is that, learning the coding matrix can pose privacy concerns: it may enable a client to infer information about the requests and side information of other clients. In other words, we are concerned in this application by an adversary who is one of the clients in the index coding setting, but is curious, *i.e.*, it wants to learn information about the requests and side information of other clients. In order to capture the effect of such adversary, our security metric of choice is an information-theoretic one which captures how much information is leaked by the adversary from learning the coding matrix. Since our adversary of concern is a legitimate client (hence it learns the coding matrix), it is not clear how cryptographic tools can be used to provide privacy. However, with a careful examination of the problem, two different resources can be used to provide a security solution. The first of which is the requests and side information sets: index codes can be designed to provide privacy of one quantity at the expense of the other. The main focus of this chapter is to design index codes which provides a trade-off between the privacy of these two quantities. In the following chapter, we discuss the use of another resource: increasing the number of transmissions to attain better privacy guarantees.

## 2.1 Introduction

Consider a set of clients who share the same broadcast domain and wish to download data content from a server. Even though the content that they request may be publicly available, they wish to preserve the anonymity of their requests. For instance, assume that a client requests a video from YouTube related to a particular medical condition. If other clients learn about the identity of that request, this may then violate the privacy of that client. In this chapter, we are interested in studying how to maintain the privacy of clients sharing a broadcast domain.

It is well established that coding across the content messages of the clients is needed to efficiently use the shared broadcast domain, as formalized in index coding [BBJ11]. A typical index coding instance consists of a server with $m$ messages, connected through a broadcast channel to a set of $n$ clients. Each client possesses a subset of the messages as side information and requires a specific new message. The server then uses these side information sets to send coded transmissions, which efficiently deliver the required messages to the clients.

In this chapter we claim that index coding poses a privacy challenge. Consider, for example, that a server transmits $b_1 + b_2$ to satisfy client 1. Since this is a broadcast transmission, other clients observing this transmission will infer that the request of client 1 is either $b_1$ or $b_2$, while the other message must belong to her side information. This suggests that, although the clients can securely convey their requests to the server (e.g., through pairwise keys), a curious client may be able to infer information about the requests and/or side information sets of other clients by learning the encoding matrix used to generate the broadcast transmissions.

The first question we ask is: how much information does the encoding matrix in index coding reveal about the requests and the side information of other users? At a high level, one can think of the request and side information as two shared secrets between each client and the server, where one secret could be used to protect the other. Therefore, as we also show in the chapter, these two aspects exhibit a trade-off: maintaining a certain level of privacy on one aspect limits the amount of privacy level achieved on the other. We also ask: can we

13

design index coding matrices that, for a given number of transmissions, achieve the highest possible level of privacy? How should these matrices be designed and how much privacy can they guarantee?

In this chapter, we take first steps in answering such questions. Our main contributions can be summarized as follows:

1) We propose an information-theoretic metric to characterize the levels of privacy that can be guaranteed. We then provide guidelines for designing encoding matrices and transmission strategies to achieve high privacy levels;

2) We design an encoding matrix and characterize the maximum levels of privacy that it can achieve;

3) We derive universal upper bounds (i.e., which hold independently of the scheme that is used) on the maximum levels of privacy that can be attained;

4) We consider a special case of the problem and we characterize in closed-form the levels of privacy achieved by our scheme, which then we compare to the outer bounds, hence highlighting the privacy trade-off.

**Related Work.** In secure index coding [DSC12], the primal goal is to design strategies such that a passive external eavesdropper – who wiretaps the communication from the server to the clients – cannot learn any information about the messages. Differently, in this work we seek to protect clients' *privacy* against adversaries who wish to learn information about the identity of the requests and side information sets of the clients.

Recently, there has been a lot of effort trying to address privacy concerns in communication setups. For instance, a set of relevant work has considered the problem of protecting privacy of a user against a database. This problem was introduced in [CKG98] and is known as *Private Information Retrieval* (PIR). Specifically, in PIR a client wishes to receive a specific message from a set of (possibly colluding) databases, without revealing the identity of the request. Towards this end, data request and/or storage schemes were designed [TR16, FGH16] and recently the PIR capacity was characterized [SJ16, BU16].

In cryptography, the *Oblivious Transfer* (OT) problem [BCR87] has a close connection

to PIR [MDP14]. Specifically, in OT the goal is to protect both the privacy of the client against the server (i.e., as in PIR, the identity of the request of the client is not revealed to the server) and the privacy of the server against the client (i.e., the client learns only the requested message). OT has also been used as a primitive to build techniques for secure multi-party computation [MDP14].

Different from these works, in this chapter we seek to understand the privacy issues that can arise among clients who share the same broadcast domain. Specifically, we seek to design techniques that guarantee high levels of privacy both in the side information and in the request of a client against another curious client. Given the different problem formulation, the techniques developed to solve the PIR and OT problems do not easily extend to our setup.

**Chapter Organization.** The chapter is organized as follows. In Section 2.2 we define our setup. In Section 2.3 we provide definitions and guidelines on how to design privacy-preserving transmission schemes and we derive fundamental upper bounds. In Section 2.4 we present the design of a privacy-preserving matrix. Based on this matrix, in Section 2.5 we consider a specific scenario for which we propose a transmission scheme and assess its performance. In Section 3.7 we conclude the chapter. Some of the proofs are delegated to the appendices.

## 2.2 Setup

We consider a typical index coding instance, where a set of clients $\mathcal{N} = \{c_{[n]}\}$, with $|\mathcal{N}| = n$, are connected to a server through a shared broadcast channel. The server has a database of messages $\mathcal{M} = \{b_{[m]}\}$, with $|\mathcal{M}| = m$. Each client $c_i, i \in [n]$, is represented by a pair of random variables, namely: (i) $\bar{Q}_i \in [m]$ associated with the index of the message that $c_i$ wishes to download from the server and (ii) $\bar{S}_i \in 2^{[m]}$, associated with the indices of the subset of messages she already has as side information. We indicate with $\bar{q}_i$ and $\bar{\mathcal{S}}_i$ the realizations of $\bar{Q}_i$ and $\bar{S}_i$, respectively, which are chosen uniformly at random from their respective domains. Clearly, $\bar{q}_i \notin \bar{\mathcal{S}}_i$. We assume that the pairs $(\bar{Q}_i, \bar{S}_i)$, $\forall i \in [n]$, are

independent across $i \in [n]$.

**Server Model.** We assume that the server knows the request and the side information of each client, i.e., it is aware of the realizations of the random variables $\bar{Q}_i = \bar{q}_i$ and $\bar{S}_i = \bar{\mathcal{S}}_i$, with $i \in [n]$. Given this, the server seeks to satisfy the requests of the clients through $T$ broadcast transmissions. The server employs linear encoding, i.e., each transmission consists of a linear combination of the $m$ messages, where the coefficients are chosen from a finite field $\mathbb{F}_L$ with $L$ being large enough. This can be mathematically formulated as $\mathbf{A}\mathbf{b} = \mathbf{y}$, where $\mathbf{b} \in \mathbb{F}_L^m$ is the column vector of the $m$ messages, $\mathbf{A} \in \mathbb{F}_L^{T \times m}$ is the encoding matrix used by the server and $\mathbf{y} \in \mathbb{F}_L^T$ is the column vector with linear combinations of the messages.

Therefore, a *transmission scheme* employed by the server consists of the following two components:

i) *Transmission space*: a specific set $\mathcal{A}$ of encoding matrices designed to satisfy the clients and protect their privacy;

ii) *Transmission strategy:* a function that, given $(\bar{q}_{[n]}, \bar{\mathcal{S}}_{[n]})$, determines the encoding matrix $\mathbf{A} \in \mathcal{A}$ to be used. We model the output of the function as a random variable $\mathbf{A}$ where $\mathbf{A} = \hat{\mathbf{A}}$ according to a probability distribution $p_{\mathbf{A}|\bar{Q}_{[n]}, \bar{S}_{[n]}}(\hat{\mathbf{A}}|\bar{q}_{[n]}, \bar{\mathcal{S}}_{[n]})$ that has to be designed.

**Adversary Model.** We assume that some of the clients – referred to as *eavesdroppers* – are malicious. Specifically, the eavesdroppers are non-cooperative clients who, based on the broadcast transmissions they receive, are eager to infer information about the requests and the side information sets of other clients. Since the eavesdroppers do not cooperate, without loss of generality, we can assume that there is only one eavesdropper in the system, namely client $c_n$. In addition, we assume that the eavesdropper $c_n$: (i) is aware of both the transmission scheme employed by the server and the underlying distribution based on which the clients obtain their requests and side information sets; (ii) has infinite computational power; (iii) knows the size of the side information set of each client, i.e., $s_i = |\bar{\mathcal{S}}_i|, i \in [n]$.

This last assumption, which we make to simplify the analysis, provides pessimistic privacy guarantees with respect to a scenario where the eavesdropper does not have this information.

Based on this knowledge, the eavesdropper $c_n$ wishes to infer information about the request and side information of the other clients. Specifically, we denote with $Q_i$ and $S_i$ the random variables, which represent the eavesdropper's estimate of the request and side information of client $c_i$, respectively and we let $p_{Q_i}(q_i)$ and $p_{S_i}(\mathcal{S}_i)$ be the corresponding probability density functions. For ease of notation, in the rest of the chapter, we drop the subscripts from the probability density functions while retaining the arguments. Clearly, $Q_n = \bar{Q}_n$ and $S_n = \bar{S}_n$. Before transmission, the eavesdropper is completely oblivious to $Q_i$ and $S_i$ for $i \in [n-1]$; we model this situation by having $p(q_i|s_i)$ and $p(\mathcal{S}_i|s_i)$ uniformly distributed over $[m]$ and $\binom{[m]}{s_i}$, respectively[1]. Then, by learning the specific encoding matrix $\hat{\mathbf{A}}$ employed by the server, the eavesdropper infers some information about the other clients, which is reflected in the conditional probability distributions $p(q_i|\hat{\mathbf{A}}, s_{[n]}, q_n, \mathcal{S}_n)$ and $p(\mathcal{S}_i|\hat{\mathbf{A}}, s_{[n]}, q_n, \mathcal{S}_n)$.

**Privacy Metric.** We consider the amount of knowledge the eavesdropper has about the variables $Q_i$ and $S_i$ as a privacy metric. In particular, we evaluate how far the uniform distribution is from the conditional distribution that the eavesdropper has after learning the encoding matrix $\hat{\mathbf{A}}$. Let $X \in \{Q_{[n]}, S_{[n]}\}$. Then, inspired by the t-closeness metric for data privacy [LLV07], we consider the *Kullback–Leibler divergence* as a distance metric between the distributions $p(x|\hat{\mathbf{A}}, s_i, q_n, \mathcal{S}_n)$ and $p(x|s_i)$, namely

$$D_{\mathrm{KL}}(p(x|\hat{\mathbf{A}}, s_{[n]}, q_n, \mathcal{S}_n)||p(x|s_i)) = \log(|\mathcal{X}|) - H(X|\hat{\mathbf{A}}, s_{[n]}, q_n, \mathcal{S}_n), \qquad (2.1)$$

where $\mathcal{X}$ is the support of $X$ (note that the entropy used throughout the chapter is conditioned on specific realizations). If $D_{\mathrm{KL}}(p(x|\hat{\mathbf{A}}, s_{[n]}, q_n, \mathcal{S}_n)||p(x|s_i)) = 0$, i.e., $H(X|\hat{\mathbf{A}}, s_{[n]}, q_n, \mathcal{S}_n) = \log(|\mathcal{X}|))$, then the eavesdropper has no knowledge of the variable $X$. Differently, larger values of $D_{\mathrm{KL}}(p(x|\hat{\mathbf{A}}, s_{[n]}, q_n, \mathcal{S}_n)||p(x|s_i))$, i.e., smaller values of $H(X|\hat{\mathbf{A}}, s_{[n]}, q_n, \mathcal{S}_n)$ indicate lower levels of privacy. Therefore, we consider $H(X|\hat{\mathbf{A}}, s_{[n]}, q_n, \mathcal{S}_n)$ as an indication of the

---

[1]In principle, in $p(q_i|s_i)$ and $p(\mathcal{S}_i|s_i)$ we should also have $q_n$, $\mathcal{S}_n$ and $s_{[n]\setminus\{i\}}$ in the conditioning. However, since $(\bar{Q}_i, \bar{S}_i)$, $\forall i \in [n]$, are independent across $i$, we can safely drop this dependence.

level of privacy attained for the variable $X$. We focus on designing transmission schemes with guaranteed levels of privacy regarding three different quantities for each client:

i) Privacy in the request, captured by $H(Q_i|\hat{\mathbf{A}}, s_{[n]}, q_n, \mathcal{S}_n)$;

ii) Privacy in the side information, captured by $H(S_i|\hat{\mathbf{A}}, s_{[n]}, q_n, \mathcal{S}_n)$;

iii) Joint privacy, captured by $H(Q_i, S_i|\hat{\mathbf{A}}, s_{[n]}, q_n, \mathcal{S}_n)$.

Therefore, our goal is to design a transmission scheme which provides privacy guarantees - in terms of the aforementioned metrics - for a given number of transmissions.

## 2.3 Guidelines for Protecting Privacy

Based on the knowledge of $(\bar{Q}_{[n]}, \bar{S}_{[n]})$, the server chooses to use an encoding matrix $\mathbf{A} = \hat{\mathbf{A}}$ such that it satisfies all clients, i.e., it allows each client to decode her request using her side information set.

**Definition 2.3.1.** A $(q,\mathcal{S})$ pair is said to be *decodable in* $\hat{\mathbf{A}}$ if, using $\hat{\mathbf{A}}$ as encoding matrix, message $b_q$ can be decoded knowing $b_\mathcal{S}$.

**Definition 2.3.2.** A $q$ (or $\mathcal{S}$) is said to be *decodable in* $\hat{\mathbf{A}}$ if there exists $\mathcal{S}$ (or $q$) such that $(q,\mathcal{S})$ is decodable in $\hat{\mathbf{A}}$.

In order to design an encoding matrix that satisfies all clients, we rely on the following lemma – a slight variation of [SF15, Lemma 4] – which provides a decodability criterion for $(q,\mathcal{S})$ using a matrix $\hat{\mathbf{A}}$.

**Lemma 2.3.1** (Decodability Criterion). *Let $\hat{\mathbf{A}}$ be the encoding matrix used by the server. Then, the pair $(q,\mathcal{S})$ is decodable in $\hat{\mathbf{A}}$ iff $\hat{\mathbf{A}}_q \notin span(\hat{\mathbf{A}}_{[m]\setminus\{q\cup\mathcal{S}\}})$.*

Lemma 2.3.1 provides a necessary and sufficient algebraic condition on whether a particular $(q, \mathcal{S})$ pair is decodable using a given encoding matrix. The eavesdropper, when trying to infer information about $c_i, i \in [n-1]$, can therefore apply this decodability criterion on

all possible $(q_i, \mathcal{S}_i)$ pairs with $|\mathcal{S}_i| = s_i$, to determine the subset of pairs that are decodable using $\hat{\mathbf{A}}$. In other words, since she knows that the request of client $c_i$ must be satisfied, then the actual $(\bar{q}_i, \bar{\mathcal{S}}_i)$ pair of client $c_i$ must belong to this set of decodable pairs. Thus, the size of the set of decodable pairs with side information sets of size $s_i$ determines the uncertainty that the eavesdropper has regarding the information of client $c_i$ and hence the attained levels of privacy for $c_i$. Therefore, in order to maintain high levels of privacy, it is imperative to design encoding matrices with decodable sets of large sizes.

We next formalize this intuition. Towards this end, we define the following three quantities: (i) $\mathcal{D}(\hat{\mathbf{A}}, s_i)$, i.e., the set of decodable $(q_i, \mathcal{S}_i)$ pairs in $\hat{\mathbf{A}}$ for client $c_i$; (ii) $\mathcal{D}^Q(\hat{\mathbf{A}}, s_i)$, i.e., the set of decodable $q_i$ in $\hat{\mathbf{A}}$ for client $c_i$, and (iii) $\mathcal{D}^S(\hat{\mathbf{A}}, s_i)$, i.e., the set of decodable $\mathcal{S}_i$ in $\hat{\mathbf{A}}$ for client $c_i$. To better understand this notation, consider the following example.

**Example.** Consider $m = 5$, $n = 2$ and $s_1 = 1$. If the server uses $\hat{\mathbf{A}}_1 = \begin{bmatrix} 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 \end{bmatrix}$ as an encoding matrix, then $\mathcal{D}(\hat{\mathbf{A}}_1, 1) = \{(1, i), (3, j)\}$ with $i \in [5]\backslash\{1\}$ and $j \in [5]\backslash\{3\}$, $\mathcal{D}^Q(\hat{\mathbf{A}}_1, 1) = \{1, 3\}$ and $\mathcal{D}^S(\hat{\mathbf{A}}_1, 1) = [5]$. Now, suppose that the server uses $\hat{\mathbf{A}}_2 = \begin{bmatrix} 1 & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 1 & 0 \end{bmatrix}$. Then, $\mathcal{D}(\hat{\mathbf{A}}_2, 1) = \{(1, 2), (2, 1), (3, 4), (4, 3)\}$, $\mathcal{D}^Q(\hat{\mathbf{A}}_2, 1) = \mathcal{D}^S(\hat{\mathbf{A}}_2, 1) = [4]$. Clearly, $|\mathcal{D}(\hat{\mathbf{A}}_1, 1)| > |\mathcal{D}(\hat{\mathbf{A}}_2, 1)|$ and $|\mathcal{D}^S(\hat{\mathbf{A}}_1, 1)| > |\mathcal{D}^S(\hat{\mathbf{A}}_2, 1)|$, but $|\mathcal{D}^Q(\hat{\mathbf{A}}_1, 1)| < |\mathcal{D}^Q(\hat{\mathbf{A}}_2, 1)|$.

With this, we have the following remark that relates the privacy metrics to the sizes of the decodable sets (see Appendix 2.7.1 for details).

**Remark 2.3.2.** *When the eavesdropper observes the encoding matrix $\hat{\mathbf{A}}$, then for all $i \in [n-1]$ and $s_i \in [m-1]$, we have*

$$H(Q_i, S_i | \hat{\mathbf{A}}, s_{[n]}, q_n, \mathcal{S}_n) \leq \log|\mathcal{D}(\hat{\mathbf{A}}, s_i)|, \tag{2.2a}$$

$$H(Q_i | \hat{\mathbf{A}}, s_{[n]}, q_n, \mathcal{S}_n) \leq \log|\mathcal{D}^Q(\hat{\mathbf{A}}, s_i)|, \tag{2.2b}$$

$$H(S_i | \hat{\mathbf{A}}, s_{[n]}, q_n, \mathcal{S}_n) \leq \log|\mathcal{D}^S(\hat{\mathbf{A}}, s_i)|. \tag{2.2c}$$

*Moreover, these bounds are tight iff the corresponding probability distributions are uniform. Namely:*

*i) eq.(2.2a) is tight iff $p(q_i, \mathcal{S}_i | \hat{\mathbf{A}}, s_{[n]}, q_n, \mathcal{S}_n)$ is uniform over $(q_i, \mathcal{S}_i) \in \mathcal{D}(\hat{\mathbf{A}}, s_i)$;*

19

*ii) eq.(2.2b) is tight iff $p(q_i|\hat{\mathbf{A}}, s_{[n]}, q_n, \mathcal{S}_n)$ is uniform over $q_i \in \mathcal{D}^Q(\hat{\mathbf{A}}, s_i)$;*

*iii) eq.(2.2c) is tight iff $p(\mathcal{S}_i|\hat{\mathbf{A}}, s_{[n]}, q_n, \mathcal{S}_n)$ is uniform over $\mathcal{S}_i \in \mathcal{D}^S(\hat{\mathbf{A}}, s_i)$.*

Remark 2.3.2 implies that the sizes of the decodable sets give an upper bound on the corresponding levels of the privacy metrics. Moreover, one can show that the conditions i) to iii) in Remark 2.3.2 hold – and hence bounds (2.2a) to (2.2c) are tight – if $p(\hat{\mathbf{A}}|\bar{q}_{[n]}, \bar{\mathcal{S}}_{[n]})$ in the transmission strategy (described in Section 2.2) is properly designed. For instance, using Bayes' rule, it can be shown – see Appendix 2.7.1 for the details – that condition i) is satisfied iff

$$\sum_{q_{\mathcal{K}}, \mathcal{S}_{\mathcal{K}} \in \prod_{j \in \mathcal{K}} \mathcal{D}(\hat{\mathbf{A}}, s_j)} p(\hat{\mathbf{A}}|q_{[n]}, \mathcal{S}_{[n]}, s_{[n]}), \quad \mathcal{K} = [n-1]\backslash\{i\}$$

is the same for all $(q_i, \mathcal{S}_i) \in \mathcal{D}(\hat{\mathbf{A}}, s_i)$.

From Remark 2.3.2, it follows that the design of privacy-preserving transmission schemes consists of two main steps: (i) designing encoding matrices with large decodable sets and (ii) using transmission strategies which satisfy uniformity conditions and hence achieve maximum levels of privacy.

Based on the result in Remark 2.3.2, we now derive universal upper bounds (i.e., which hold independently of the encoding matrix that the server uses) on the decodable sets and hence on the levels of the privacy metrics. In particular, we have

**Lemma 2.3.3.** *For any $\hat{\mathbf{A}} \in \mathbb{F}_L^{T \times m}$ and $s_i \in [m-1]$, we have*

$$|\mathcal{D}(\hat{\mathbf{A}}, s_i)| \leq T\binom{m}{s_i} =: \mathsf{UB}_{Q,S}, \tag{2.3a}$$

$$|\mathcal{D}^Q(\hat{\mathbf{A}}, s_i)| \leq m =: \mathsf{UB}_Q, \tag{2.3b}$$

$$|\mathcal{D}^S(\hat{\mathbf{A}}, s_i)| \leq \binom{m}{s_i} =: \mathsf{UB}_S. \tag{2.3c}$$

**Proof:** The upper bounds in (2.3b) and (2.3c) simply follow by noticing that the size of a decodable set is upper bounded by the size of the support of the corresponding random variable. We next prove the bound in (2.3a). For a given encoding matrix $\hat{\mathbf{A}} \in \mathbb{F}_L^{T \times m}$, one can

$$\begin{bmatrix} \overbrace{\mathbf{A}_b}^{\text{Seg. 1}} & \overbrace{\mathbf{0}_{\frac{T}{k}\times\ell}}^{\text{Seg. 2}} & \cdots & \overbrace{\mathbf{0}_{\frac{T}{k}\times\ell}}^{\text{Seg. }k} & \overbrace{\phantom{\mathbf{0}_{T\times m-k\ell}}}^{\text{Seg. 0}} \\ \mathbf{0}_{\frac{T}{k}\times\ell} & \mathbf{A}_b & \cdots & \mathbf{0}_{\frac{T}{k}\times\ell} & \\ \vdots & \vdots & \ddots & \vdots & \mathbf{0}_{T\times m-k\ell} \\ \mathbf{0}_{\frac{T}{k}\times\ell} & \mathbf{0}_{\frac{T}{k}\times\ell} & \cdots & \mathbf{A}_b & \end{bmatrix}$$

Figure 2.1: Design of the base matrix $\mathbf{A}^{\text{base}}$ for the achievable scheme.

write $\mathcal{D}(\hat{\mathbf{A}}, s_i) = \sum_{\mathcal{S}_i \in \binom{[m]}{s_i}} \mathcal{N}(\hat{\mathbf{A}}, \mathcal{S}_i)$, where $\mathcal{N}(\hat{\mathbf{A}}, \mathcal{S}_i)$ is the set of requests $q_i \in \mathcal{D}^Q(\hat{\mathbf{A}}, s_i)$ for which the pair $(q_i, \mathcal{S}_i)$ is decodable. According to Lemma 2.3.1, for each $q_i \in \mathcal{N}(\hat{\mathbf{A}}, \mathcal{S}_i)$, $\hat{\mathbf{A}}_{q_i}$ is not in the span of $\hat{\mathbf{A}}_{[m] \setminus \mathcal{S}_i \cup q_i}$. It is therefore straightforward to show that the columns of $\hat{\mathbf{A}}_{\mathcal{N}(\hat{\mathbf{A}}, \mathcal{S}_i)}$ are linearly independent. Thus, $|\mathcal{N}(\hat{\mathbf{A}}, \mathcal{S}_i)| \leq T$ and hence we have $|\mathcal{D}(\hat{\mathbf{A}}, s_i)| \leq T\binom{m}{s_i}$. ∎

## 2.4 Design of a Transmission Space

In this section, we take first steps towards designing a privacy-preserving transmission scheme. Specifically, we design an encoding matrix, referred to as the *base* matrix $\mathbf{A}^{\text{base}}$. Then, we populate the transmission space with the matrices obtained from $\mathbf{A}^{\text{base}}$ by taking all the permutations of its columns. Our design of $\mathbf{A}^{\text{base}}$ is based on the use of Maximum Distance Separable (MDS) codes. A generator matrix of an $[m, T]$ MDS code has the property that any $T \times T$ submatrix is full rank, i.e., any $T$ columns are linearly independent. Such matrices promise to provide large decodable sets. To see this notice that, for a given side information set $\mathcal{S}$ with $|\mathcal{S}| \geq m - T$, all requests in $[m] \setminus \mathcal{S}$ are decodable with $\mathcal{S}$. Therefore, if $\mathbf{B} \in \mathbb{F}_L^{T \times m}$ is a generator matrix of an $[m, T]$ MDS code, then, for all $s \geq m - T$, we have $|\mathcal{D}^Q(\mathbf{B}, s)| = m$ and $|\mathcal{D}(\mathbf{B}, s)| = m\binom{m-1}{s} = O(m^s)$. However, this scheme might require a prohibitively large number of transmissions $T$, especially when $m$ is large and $s$ is small compared to $m$. To achieve high levels of privacy with $T$ that is not that large, we next propose the design of $\mathbf{A}^{\text{base}}$, which is based on a block-MDS as shown in Figure 2.1 and

Figure 2.2: Numerical evaluation of $r_Q$, $r_S$ and $r_{Q,S}$ - $m = 30$ and $s = 3$.

structured as follows:

i) The columns of $\mathbf{A}^{\text{base}}$ are divided in $k+1$ segments, labeled as "Seg. 0 to $k$", where $T$ is a multiple of $k$;

ii) Segments from 1 to $k$ consist of $\ell$ columns, where $\ell \leq \min\{s_{\min} + T/k, \lfloor m/k \rfloor\}$, with $s_{\min} = \min_{i \in [n]} s_i$;

iii) A matrix $\mathbf{A}_b \in \mathbb{F}_L^{\frac{T}{k} \times \ell}$ is constructed as the generator matrix of an $[\ell, T/k]$ MDS code; then, $\mathbf{A}_b$ is repeated $k$ times and positioned in $\mathbf{A}^{\text{base}}$ as shown in Figure 2.1;

iv) The rest of $\mathbf{A}^{\text{base}}$ is filled with zeros.

Note that, for any number of clients $n$ and messages $m$, one can always find values of $k$, $\ell$ and $T$ so that $\mathbf{A}^{\text{base}}$ satisfies all clients (e.g., $k = 1$, $\ell = s_{\min}$ and $T = n$).

We now analyze the performance of our proposed $\mathbf{A}^{\text{base}}$ in terms of the sizes of its decodable sets (see Appendix 2.7.2). These, by means of Remark 2.3.2, provide upper bounds on the levels of privacy that could be attained using $\mathbf{A}^{\text{base}}$.

**Theorem 2.4.1.** *For $\mathbf{A}^{base}$ and any $s_i \in [m - 1]$, we have*

$$H(Q_i, S_i | \hat{\mathbf{A}}, s_{[n]}, q_n, \mathcal{S}_n) \leq \log \left( k\ell \sum_{j=\ell-T/k}^{\ell-1} \binom{\ell - 1}{j} \binom{m - \ell}{s_i - j} \right), \tag{2.4a}$$

$$H(Q_i | \hat{\mathbf{A}}, s_{[n]}, q_n, \mathcal{S}_n) \leq (k\ell). \tag{2.4b}$$

*where the bounds can be achieved by satisfying the uniformity conditions in Remark 2.3.2*

In the next section we study a special scenario in which we use the transmission space here proposed (i.e., populated by the matrices obtained from $\mathbf{A}^{\text{base}}$ by taking all the permutations of its columns) and we design the transmission strategy.

22

## 2.5 Transmission Strategy for a Special Case

In the previous section, we designed a transmission space that consists of all possible matrices obtained by permuting the columns of the matrix $\mathbf{A}^{\text{base}}$. Thus, as discussed in Section 2.2, in order to design a transmission scheme, we need to design a transmission strategy that selects which specific matrix to use according to a probability distribution. However, designing such a transmission strategy that achieves the upper bounds in Remark 2.3.2 is non-trivial. To get an analytical handle on the problem, we take a first step and consider a simplified model: we assume $n = 2$ and an eavesdropper who does *not* have a request. Such a scenario can model a situation where the $n = 2$ clients (the second of which is the eavesdropper) do not have a simultaneous request.

Since only one client needs to be satisfied, then we can use our proposed encoding matrix $\mathbf{A}^{\text{base}}$ with $k = T$ and $\ell \leq \min\{s_1 + 1, \lfloor m/T \rfloor\}$, knowing that the client $c_1$ can always be satisfied by using the appropriate column-permutation of $\mathbf{A}^{\text{base}}$ (i.e., by ensuring that $\mathbf{A}_{q_1}^{\text{base}}$ is non-zero, and all other columns belonging to the same segment of $\mathbf{A}_{q_1}^{\text{base}}$ correspond to messages in $\mathcal{S}_1$). In this case, $\mathbf{A}_b$ is a row vector of arbitrary non-zero values. The following theorem (whose proof can be found in Appendix 2.7.3) then provides analytical guarantees on the attained performance of this scheme.

**Theorem 2.5.1.** *For the scheme described above, we have*

$$H(Q_1, S_1 | \hat{\mathbf{A}}, s_1) = \log T\ell \binom{m - \ell}{s_1 - \ell + 1} =: \mathsf{LB}_{Q,S}, \tag{2.5a}$$

$$H(Q_1 | \hat{\mathbf{A}}, s_1) = \log T\ell =: \mathsf{LB}_Q, \tag{2.5b}$$

$$H(S_1 | \hat{\mathbf{A}}, s_1) = \log T\ell \binom{m - \ell}{s_1 - \ell + 1} - \mathsf{K} =: \mathsf{LB}_S, \tag{2.5c}$$

$$\mathsf{K} = \sum_{i=1}^{T} \binom{T-1}{i-1} \ell^{i-1} \frac{\binom{m-i\ell}{s_1 - i(\ell-1)}}{\binom{m-\ell}{s_1-\ell+1}} \sum_{x=1}^{i} (-1)^{i-x} \binom{i-1}{x-1} \log x,$$

*where $\hat{\mathbf{A}}$ is the column permutation of $\mathbf{A}^{base}$ that is used.*

Note that the two quantities in (2.5a) and (2.5b) meet the upper bounds that follow from Theorem 2.4.1 by applying the conditions in Remark 2.3.2. Moreover, in order to get

23

the bounds in (2.5), we used a transmission strategy for which $p(\hat{\mathbf{A}}|\bar{q}_1, \bar{\mathcal{S}}_1)$ is uniform over all $\hat{\mathbf{A}}$ that satisfy $(\bar{q}_1, \bar{\mathcal{S}}_1)$ for all $(\bar{q}_1, \bar{\mathcal{S}}_1) \in \mathcal{D}(\hat{\mathbf{A}}, s_1)$. This is because, thanks to the special structure of $\mathbf{A}^{\text{base}}$, the number of column-permutations of $\mathbf{A}^{\text{base}}$ that satisfies a given $(\bar{q}_1, \bar{\mathcal{S}}_1)$ is equal for all $(\bar{q}_1, \bar{\mathcal{S}}_1)$.

We next analyze the performance of our scheme. Towards this end, we define the following quantities:

- $\mathsf{G}_{Q,S} := \log(\mathsf{UB}_{Q,S}) - \mathsf{LB}_{Q,S}$, $\mathsf{r}_{Q,S} = 2^{-\mathsf{G}_{Q,S}}$;

- $\mathsf{G}_Q := \log(\mathsf{UB}_Q) - \mathsf{LB}_Q$, $\mathsf{r}_Q = 2^{-\mathsf{G}_Q}$;

- $\mathsf{G}_S := \log(\mathsf{UB}_S) - \mathsf{LB}_S$, $\mathsf{r}_S = 2^{-\mathsf{G}_S}$.

Figure 2.2 shows an example of how the quantities $\mathsf{r}_{Q,S}$, $\mathsf{r}_Q$ and $\mathsf{r}_S$ behave as $\ell$ changes. Note that all these quantities are fractions and hence the maximum level of privacy (y-axis) is 1. Figure 2.2 shows that as $\ell$ increases, higher values of privacy are attained in the requests (i.e., $\mathsf{r}_Q$ increases), but smaller levels of privacy are achieved in the side information (i.e., $\mathsf{r}_S$ decreases). This highlights a trade-off: maintaining a certain level of privacy on one aspect limits the amount of privacy level achieved on the other. It is also noted that increasing $T$ increases the attained values of $\mathsf{r}_Q$ and $\mathsf{r}_S$ for the same value of $\ell$. We believe that the reason such increase does not occur in $\mathsf{r}_{Q,S}$ is because $\mathsf{UB}_{Q,S}$ in (2.3a) is loose.

Next, we assess the performance of our scheme when the parameters of the system grow. We assume that $s_1 = c \cdot m$ and $\ell = b \cdot m + 1$, where $b \leq c \leq \frac{m-1}{m}$. We consider two cases:

**Case I: $c = \frac{m-k_c}{m}$ where $k_c > 0$ is a constant.** In this case, full privacy in the request, side information and their joint can be achieved by using an $[m, T]$ MDS code with $T = k_c$.

**Case II: $c$ and $T$ are constants.** In this case, by choosing $b = 0$, we get $\mathsf{G}_Q = \log \frac{m}{T} = O(\log m)$ and $\mathsf{G}_{Q,S} = \log \frac{\binom{m}{cm}}{T\binom{m-1}{cm}} = \log \frac{1}{T(1-c)} = O(1)$. Also, since conditioning reduces the entropy, we have $H(S_1|\hat{\mathbf{A}}, s_1) \geq$ eq. (2.5a) $-$ eq. (2.5b), which implies $\mathsf{G}_S \leq \log \frac{\binom{m}{cm}}{\binom{m-1}{cm}} = \log \frac{1}{1-c} = O(1)$. This suggests that when $s_1$ grows as a constant fraction of $m$, then with a constant number of transmissions we can have almost perfect side information (and joint)

24

privacy, but very little privacy in the request. However, if we choose $b = c$, then we get $\mathsf{G}_Q \leq \log \frac{1}{Tc} = O(1)$, $\mathsf{G}_{Q,S} = \mathsf{G}_S \leq \log \binom{m}{cm} = O(m \log m)$ since, under these conditions, $\mathsf{K} = 0$ in (2.5c). Thus, in this case almost full privacy is achieved in the request while very little privacy is attained in the side information (and in the joint).

## 2.6 Conclusion

We considered an index coding instance where some clients are malicious: they wish to learn information about the requests and side information of the other clients. We showed how this privacy breach is possible by learning the encoding matrix used by the server. We proposed information-theoretic metrics to model the levels of privacy that can be guaranteed and we designed an encoding matrix for protecting privacy. Then, for a special case of the problem, we derived in closed-form the levels of privacy that our proposed scheme achieves. We showed an inherent trade-off between protecting privacy of either the request or the side information set of the clients.

## 2.7 Appendices

### 2.7.1 Proof of Remark 2.3.2

We prove the result for the upper bound in (2.2a). Given $\hat{\mathbf{A}}$ and $s_i$, the set $\mathcal{D}(\hat{\mathbf{A}}, s_i)$ consists of all possible $(q_i, \mathcal{S}_i)$ pairs that could be the request/side information pair for $c_i$. Therefore, $p(q_i, \mathcal{S}_i | \hat{\mathbf{A}}, s_{[n]}, q_n, \mathcal{S}_n) = 0$ for all $(q_i, \mathcal{S}_i) \notin \mathcal{D}(\hat{\mathbf{A}}, s_i)$. Therefore,

$$H(Q_i, S_i | \hat{\mathbf{A}}, s_{[n]}, q_n, \mathcal{S}_n) = - \sum_{(q_i, \mathcal{S}_i) \in \mathcal{D}(\hat{\mathbf{A}}, s_i)} p(q_i, \mathcal{S}_i | \hat{\mathbf{A}}, s_{[n]}, q_n, \mathcal{S}_n) \log p(q_i, \mathcal{S}_i | \hat{\mathbf{A}}, s_{[n]}, q_n, \mathcal{S}_n) \leq \log |\mathcal{D}(\hat{\mathbf{A}}, s_i)|,$$

thus proving (2.2a). Since $p(q_i, \mathcal{S}_i | \hat{\mathbf{A}}, s_{[n]}, q_n, \mathcal{S}_n) = 0$ for all $(q_i, \mathcal{S}_i) \notin \mathcal{D}(\hat{\mathbf{A}}, s_i)$, then this upper bound is achieved if and only if $p(q_i, \mathcal{S}_i | \hat{\mathbf{A}}, s_{[n]}, q_n, \mathcal{S}_n)$ is uniform over for $(q_i, \mathcal{S}_i) \in \mathcal{D}(\hat{\mathbf{A}}, s_i)$, thus proving the uniformity condition $i$) on (2.2a). Similar arguments can be made to prove (2.2b) and (2.2c).

Next, we show that the uniformity conditions in $i$)-$iii$) imply constraints on the design

of the transmission strategy $p(\hat{\mathbf{A}}|q_{[n]}, \mathcal{S}_{[n]})$. To see this, note that we can write

$$p(q_i, \mathcal{S}_i | \hat{\mathbf{A}}, s_{[n]}, q_n, \mathcal{S}_n) = p(\hat{\mathbf{A}}|q_{\{i,n\}}, \mathcal{S}_{\{i,n\}}, s_{[n]}) \frac{p(q_i, \mathcal{S}_i | s_{[n]}, q_n, \mathcal{S}_n)}{p(\hat{\mathbf{A}}|s_{[n]}, q_n, \mathcal{S}_n)},$$

which follows by applying Bayes' rule. Since the probabilities in the fraction term do not depend on the value of $(q_i, \mathcal{S}_i)$ (note that $p(q_i, \mathcal{S}_i | s_{[n]})$ is uniform), then the uniformity condition i) is satisfied if and only if the term $p(\hat{\mathbf{A}}|q_{\{i,n\}}, \mathcal{S}_{\{i,n\}}, s_{[n]})$ is the same for all $(q_i, \mathcal{S}_i) \in \mathcal{D}(\hat{\mathbf{A}}, s_i)$. We can further write

$$p(\hat{\mathbf{A}}|q_{\{i,n\}}, \mathcal{S}_{\{i,n\}}, s_{[n]}) = \sum_{q_{\mathcal{K}}, \mathcal{S}_{\mathcal{K}} \in \prod_{j \in \mathcal{K}} \mathcal{D}(\hat{\mathbf{A}}, s_j)} p(\hat{\mathbf{A}}|q_{[n]}, \mathcal{S}_{[n]}, s_{[n]}) p(q_{\mathcal{K}}, \mathcal{S}_{\mathcal{K}}|q_i, \mathcal{S}_i, s_{[n]}), \ \ \mathcal{K} = [n-1] \setminus i.$$

Note that the distribution $p(q_{\mathcal{K}}, \mathcal{S}_{\mathcal{K}}|q_i, \mathcal{S}_i, s_{[n]})$ is assumed to be uniform and independent over $i \in [n]$. Therefore, to satisfy the uniformity condition, we must have the summation term on the Righ-Hand Side to be the same for all $(q_i, \mathcal{S}_i) \in \mathcal{D}(\hat{\mathbf{A}}, s_i)$. This therefore imposes constraints on the transmission strategy used by the server. We can similarly show that the uniformity conditions on (2.2b) and (2.2c) also impose constraints on the used transmission strategy.

### 2.7.2 Proof of Theorem 2.4.1

In order to prove Theorem 2.4.1, we need to characterize the quantities $|\mathcal{D}(\hat{\mathbf{A}}, s)|$ and $|\mathcal{D}^Q(\hat{\mathbf{A}}, s)|$, and therefore, using Remark 2.3.2 the result in Theorem 2.4.1 follows.

**Characterizing $|\mathcal{D}^Q(\hat{\mathbf{A}}, s)|$:** One can show that every request $q$ whose corresponding column $\mathbf{A}_q^{\text{base}}$ is non-zero has at least one side information set $\mathcal{S}$ with which $(q, \mathcal{S})$ is decodable in $\mathbf{A}^{\text{base}}$. If this in fact is true, then the result $|\mathcal{D}^Q(\mathbf{A}^{\text{base}}, s)| = k\ell$ follows immediately, since we have $k\ell$ such requests. To prove this statement then, notice that $\ell \leq s_{\min} + T/k$. Then consider a side information set with $|\mathcal{S}| = s_{\min}$ and where all the elements of $\mathcal{S}$ correspond to columns of the same segment as $\mathbf{A}_q^{\text{base}}$. Therefore, the set of all columns of $\mathbf{A}^{\text{base}}$ belonging

to the same segment as $\mathbf{A}_q^{\text{base}}$ and do not belong to $\mathcal{S}$ is of size $\ell - \mathcal{S} = T/k$. They are therefore linearly independent, and $q$ is decodable with $\mathcal{S}$.

**Characterizing $|\mathcal{D}(\hat{\mathbf{A}}, s)|$:** To prove the remaining quantity, notice that we can write $\mathcal{D}(\mathbf{A}^{\text{base}}, s) = \sum_{q \in [m]} \mathcal{N}(\mathbf{A}^{\text{base}}, q)$, where $\mathcal{N}(\mathbf{A}^{\text{base}}, q)$ is the number of side information sets that are decodable with $q$ in $\mathbf{A}^{\text{base}}$. For a given $q$, this quantity is equal to

$$\mathcal{N}(\mathbf{A}^{\text{base}}, q) = \sum_{i=\ell-T/k}^{\ell-1} \binom{\ell-1}{i} \binom{m-\ell}{s-i}, \tag{2.6}$$

for all $q$ with $\mathbf{A}^{\text{base}}$ being non-zero, and 0 otherwise. Since this quantity does not depend on the value of $q$, then the result follows that $\mathcal{D}(\mathbf{A}^{\text{base}}, s) = k\ell \sum_{i=\ell-T/k}^{\ell-1} \binom{\ell-1}{i} \binom{m-\ell}{s-i}$. What remains is to prove (2.6), which we justify as follows: Consider a given $q$ with a non-zero corresponding column in $A^{\text{base}}$, and let $j$ be the index of the segment to which $\mathbf{A}_q^{\text{base}}$ belongs. For a given side information set $\mathcal{S}$, let $i$ be the number of elements in $\mathcal{S}$ whose corresponding columns in $\mathbf{A}^{\text{base}}$ belong to $j$. Then, $(q, \mathcal{S})$ is decodable in $\mathbf{A}^{\text{base}}$ if and only if the elements $\ell - T/k \leq i \leq \ell - 1$; the lower bound is to ensure that the columns of $\mathbf{A}^{\text{base}}$ belonging to segment $j$ that fall outside of $\mathcal{S}$ are linearly independent, and the upper bound is to ensure that $q$ is not in $\mathcal{S}$. The number of subsets $\mathcal{S}$ with $i$ columns in segment $j$ is equal to $\binom{\ell-1}{i} \binom{m-\ell}{s-\ell+1}$. Therefore, by summing over all possible $i$ and multiplying by the number of possible requests we get the expression in (2.6).

### 2.7.3   Proof of Theorem 2.5.1

For this scheme, we can have $p(\hat{\mathbf{A}}|q_1, \mathcal{S}_1) = 1/K$ for all $\hat{\mathbf{A}} \in \mathcal{A}$ for all $(q_1, \mathcal{S}_1) \in \mathcal{D}(\hat{\mathbf{A}}, s_1)$, where $K$ is equal to

$$K = T\binom{s}{\ell-1} \underbrace{\binom{m-\ell}{\ell \ell \cdots \ell}}_{k-1}^{(\text{M})},$$

where the last term is a multinomial coefficient. This is because the number of column-permutations of $\hat{\mathbf{A}}^{\text{base}}$ that satisfies a given $(q_1, \mathcal{S}_1)$ is equal to $K$, independently of the value of $(q_1, \mathcal{S}_1)$. This statement can be justified as follows: for a pair to be decodable, the column of the encoding matrix corresponding to $q$ should be non-zero, and since we have $T$ segments, then there are $T$ possibilities for that column; thus the term $T$ in the

27

expression. Next, all remaining $\ell - 1$ columns of the same segment must correspond to elements in the side information set; thus the term $\binom{s}{\ell-1}$. Finally, among the remaining $m - \ell$ columns, we have to choose $k - 1$ segments, each of length $\ell$; thus the final multinomial term.

**Calculating $H(Q_1, S_1 | \hat{\mathbf{A}}, s_1)$:** Note that by using the transmission strategy described above, we satisfy the uniformity condition of Remark 2.3.2 for (2.2a). Therefore, we have $H(Q_1, S_1 | \hat{\mathbf{A}}, s_1) = \log |\mathcal{D}(\hat{\mathbf{A}}, s_1)| = \log Tl\binom{m-\ell}{s-\ell+1}$. The last equality can be obtained by considering (2.4b) with $k = T$.

**Calculating $H(Q_1 | \hat{\mathbf{A}}, s_1)$:** Using the transmission strategy described above also satisfies the uniformity condition of Remark 2.3.2 for (2.2b). To see this, note that

$$p(q_1 | \hat{\mathbf{A}}, s_1) = \sum_{S_1 : (q_1, S_1) \in \mathcal{D}(\hat{\mathbf{A}}, s_1)} p(q_1, S_1 | \hat{\mathbf{A}}, s_1),$$

where the number of elements in the summation corresponds to the number of subsets $S_1$ that are decodable with $q_1$, which is equal to $\binom{m-\ell}{s-\ell+1}$ irrespective of $q_1$. Therefore, $p(q_1 | \hat{\mathbf{A}}, s_1)$ is uniform over all $q_1 \in \mathcal{D}^Q(\hat{\mathbf{A}}, s_1)$. Thus we have $H(Q_1 | \hat{\mathbf{A}}, s_1) = \log |D^Q(\hat{\mathbf{A}}, s_1)| = \log T\ell$, where the last equality similarly holds by considering (2.4b) with $k = T$.

**Calculating $H(S_1 | \hat{\mathbf{A}}, s_1)$:** Using the transmission strategy above does not satisfy the uniformity condition of Remark 2.3.2 for (2.2c). Therefore, we now seek to quantify the achieved value of $H(S_1 | \hat{\mathbf{A}}, s_1)$.

Note that the used transmission strategy would yield $p(q_1, S_1 | \hat{\mathbf{A}}, s_1) = 1/|\mathcal{D}(\hat{\mathbf{A}}, s_1)|$ for all $(q_1, S_1) \in \mathcal{D}(\hat{\mathbf{A}}, s_1)$ and 0 otherwise. One can then write the marginal $p(S_1 | \hat{\mathbf{A}}, s_1)$ as

$$p(S_1 | \hat{\mathbf{A}}, s) = \sum_{q_1 \in \mathcal{D}^Q(\hat{\mathbf{A}}, s_1)} p(q_1, S_1 | \hat{\mathbf{A}}, s_1) = \frac{N_{\hat{\mathbf{A}}, S_1}}{|\mathcal{D}(\hat{\mathbf{A}}, s_1)|},$$

where $N_{\hat{\mathbf{A}}, S_1}$ is the number of requests $q_1$ that are decodable with $S_1$ in $\hat{\mathbf{A}}$. Therefore, we have

$$H(\mathcal{S}_1|\hat{\mathbf{A}}, s_1) = -\sum_{\mathcal{S}_1 \in \mathcal{D}^S(\hat{\mathbf{A}}, s_1)} \frac{N_{\hat{\mathbf{A}}, \mathcal{S}_1}}{|\mathcal{D}(\hat{\mathbf{A}}, s_1)|} \log \frac{N_{\hat{\mathbf{A}}, \mathcal{S}_1}}{|\mathcal{D}(\hat{\mathbf{A}}, s_1)|}$$

$$= \log|\mathcal{D}(\hat{\mathbf{A}}, s_1)| - \frac{1}{|\mathcal{D}(\hat{\mathbf{A}}, s_1)|} \underbrace{\sum_{\mathcal{S}_1 \in \mathcal{D}^S(\hat{\mathbf{A}}, s_1)} N_{\hat{\mathbf{A}}, \mathcal{S}_1} \log N_{\hat{\mathbf{A}}, \mathcal{S}_1}}_{\bar{N}_t}. \qquad (2.7)$$

Next we calculate $\bar{N}_t$. For a given $\mathcal{S}_1$, let $\ell_j, j \in [T]$ be the number of elements of $\mathcal{S}_1$ for which the corresponding columns in $\hat{\mathbf{A}}$ belong to segment $j$. Then in order for a pair $(q_1, \mathcal{S}_1)$ to be decodable, then $\ell_j$ must be exactly equal to $\ell - 1$, where $j$ corresponds to the segment to which $\hat{\mathbf{A}}_q$ belongs.

Note that $N_{\hat{\mathbf{A}}, \mathcal{S}_1}$ only depends on the values of $\ell_j$, and therefore all subsets $\mathcal{S}_1$ for which $\ell_j, j \in [T]$ are the same will have the same value for $N_{\hat{\mathbf{A}}, \mathcal{S}_1}$. Based on this fact, we can then write

$$\bar{N}_t = \sum_{\ell_1=0}^{\ell} \cdots \sum_{\ell_T=0}^{\ell} \binom{\ell}{\ell_1} \cdots \binom{\ell}{\ell_T} \binom{m - T\ell}{s_1 - \sum_{i=1}^T \ell_i} \left( \sum_{i=1}^T \mathbb{1}_{\{\ell_i = \ell-1\}} \right) \log \left( \sum_{i=1}^T \mathbb{1}_{\{\ell_i = \ell-1\}} \right)$$

$$\stackrel{(a)}{=} \sum_{x=1}^T x \log x \binom{T}{x} \ell^x \underbrace{\left[ \sum_{\substack{\ell_1=0 \\ \ell_1 \neq \ell-1}}^{\ell} \cdots \sum_{\substack{\ell_{T-x}=0 \\ \ell_{T-x} \neq \ell-1}}^{\ell} \binom{\ell}{\ell_1} \cdots \binom{\ell}{\ell_{T-x}} \binom{m - T\ell}{s_1 - x(\ell-1) - \sum_{i=1}^{T-x} \ell_i} \right]}_{C_{s_1, T}(T-x)} \qquad (2.8)$$

where $(a)$ can be justified as follows: note that the possible values to which the term $\sum_{i=1}^T \mathbb{1}_{\{\ell_i = \ell-1\}}$ evaluates are $x \in [T]$ ($x = 0$ is also possible, but trivial). Moreover, it is equal to $x$ if and only if there are exactly $x$ indices from the set $\ell_{[T]}$ which are equal to $\ell - 1$, while the remaining indices can take any value (except $\ell - 1$). Therefore, by means of counting arguments, $\bar{N}_t$ can be expressed as (2.8).

29

Note that we can write

$$C_{s_1,T}(T-x) = \underbrace{\left[\sum_{\substack{\ell_1=0 \\ \ell_1 \neq \ell-1}}^{\ell} \cdots \sum_{\substack{\ell_{T-x}=0 \\ \ell_{T-x} \neq \ell-1}}^{l} \binom{\ell}{\ell_1} \cdots \binom{\ell}{\ell_{T-x}} \binom{m-T\ell}{s_1 - x(\ell-1) - \sum_{i=1}^{T-x} \ell_i}\right]}_{B_{s_1,T}(T-x)}$$

$$\overset{(b)}{=} \left[\sum_{\ell_1=0}^{\ell} \cdots \sum_{\ell_{T-x}=0}^{\ell} \binom{l}{\ell_1} \cdots \binom{\ell}{\ell_{T-x}} \binom{m-T\ell}{s_1 - x(\ell-1) - \sum_{i=1}^{T-x} \ell_i}\right] -$$

$$\sum_{y=1}^{T-x} \binom{T-x}{y} \ell^y \left[\sum_{\substack{\ell_1=0 \\ \ell_1 \neq \ell-1}}^{\ell} \cdots \sum_{\substack{\ell_{T-x-y}=0 \\ \ell_{T-x-y} \neq \ell-1}}^{\ell} \binom{\ell}{\ell_1} \cdots \binom{\ell}{\ell_{T-x-y}} \binom{m-T\ell}{s_1 - (x+y)(\ell-1) - \sum_{i=1}^{T-x-y} \ell_i}\right]$$

$$= B_{s_1,T}(T-x) - \sum_{y=1}^{T-x} \binom{T-x}{y} \ell^y C_{s_1,T}(T-x-y) \tag{2.9}$$

where $(b)$ follows by adding the missing summation terms of $C_{s_1,T}(T-x)$ corresponding to $\ell_i = \ell-1$ and - by means of counting - subtracting them. By noting that $C_{s_1,T}(0) = \binom{m-T\ell}{s_1-T(\ell-1)}$, equation (2.9) then defines a linear recurrence relation on $C_{s_1,T}(T-x)$ which we solve in the following lemma.

**Lemma 2.7.1.** *The solution to the linear recurrence relation in (2.9) is*

$$C_{s_1,T}(T-x) = \sum_{v=0}^{T-x} (-1)^v \ell^v \binom{T-x}{v} B_{s_1,T}(T-x-v) \tag{2.10}$$

*where $B_{s_1,T}(0) = \binom{m-T\ell}{s_1-T(\ell-1)}$.*

**Proof:** We will solve the recurrence relation using strong induction. Specifically, assume that

$$C_{s_1,T}(T-x-y) = \sum_{v=0}^{T-x-y} (-1)^v \ell^v \binom{T-x-y}{v} B_{s_1,T}(T-x-v-y)$$

for $1 \leq y \leq T-x$. Then consider

30

$$\sum_{y=1}^{T-x} \binom{T-x}{y} \ell^y C_{s_1,T}(T-x-y) =$$

$$= \sum_{y=1}^{T-x} \sum_{v=0}^{T-x-y} (-1)^v \ell^{v+y} \binom{T-x}{y} \binom{T-x-y}{v} B_{s_1,T}(T-x-v-y)$$

$$\stackrel{(c)}{=} \sum_{k=1}^{T-x} (-1)^k \ell^k \binom{T-x}{k} B_{s_1,T}(T-x-k) \sum_{v=0}^{k-1} (-1)^{v-k} \frac{\binom{T-x}{k-v}\binom{T-x-k+v}{v}}{\binom{T-x}{k}}$$

$$= \sum_{k=1}^{T-x} (-1)^k \ell^k \binom{T-x}{k} B_{s_1,T}(T-x-k) \sum_{v=0}^{k-1} (-1)^{k-v} \binom{k}{k-v}$$

$$= \sum_{k=1}^{T-x} (-1)^k \ell^k \binom{T-x}{k} B_{s_1,T}(T-x-k) \sum_{v'=1}^{k} (-1)^{v'} \binom{k}{v'}$$

$$= \sum_{k=1}^{T-x} (-1)^k \ell^k \binom{T-x}{k} B_{s_1,T}(T-x-k)(\delta_{k0}-1)$$

$$= -\sum_{k=0}^{T-x} (-1)^k \ell^k \binom{T-x}{k} B_{s_1,T}(T-x-k) + B_{s_1,T}(T-x)$$

where $(c)$ follows by i) changing summation variables as $v+y=k$ and ii) multiplying and dividing by $(-1)^k \binom{T-x}{k}$, and where $\delta_{ij}$ is the Kronecher delta function. Therefore we have

$$C_{s_1,T}(T-x) = \sum_{k=0}^{T-x} (-1)^k \ell^k \binom{T-x}{k} B_{s_1,T}(T-x-k)$$

$$= B_{s_1,T}(T-x) - \sum_{y=1}^{T-x} \binom{T-x}{y} \ell^y C_{s_1,T}(T-x-y)$$

satisfying (2.9), thus completing the proof. ∎

By plugging (2.10) in (2.8), we can further simply (2.8) as follows

$$\bar{N}_t = \sum_{x=1}^{T} x \log x \binom{T}{x} \ell^x \sum_{v=0}^{T-x} (-1)^v \ell^v \binom{T-x}{v} B_{s_1,T}(T-x-v)$$

$$= \sum_{x=1}^{T} \sum_{v=0}^{T-x} x \log x \binom{T}{x} \binom{T-x}{v} \ell^{x+v}(-1)^v B_{s_1,T}(T-x-v)$$

$$= \sum_{x=1}^{T} \sum_{v=0}^{T-x} x \log x \binom{T}{x+v} \binom{x+v}{x} \ell^{x+v}(-1)^v B_{s_1,T}(T-x-v)$$

$$= \sum_{i=1}^{T} \sum_{x=1}^{i} x \log x \binom{T}{i} \binom{i}{x} \ell^i (-1)^{i-x} B_{s_1,T}(T-i)$$

$$= \sum_{i=1}^{T} \binom{T}{i} \ell^i B_{s_1,T}(T-i) \sum_{x=1}^{i} (-1)^{i-x} \binom{i}{x} x \log x$$

$$= \sum_{i=1}^{T} \binom{T}{i} \ell^i B_{s_1,T}(T-i) \sum_{x=1}^{i} (-1)^{i-x} i \binom{i-1}{x-1} \log x$$

$$= T \sum_{i=1}^{T} \binom{T-1}{i-1} \ell^i B_{s_1,T}(T-i) \sum_{x=1}^{i} (-1)^{i-x} \binom{i-1}{x-1} \log x. \tag{2.11}$$

Also, we can write

$$B_{s_1,T}(T-x) = \underbrace{\sum_{\ell_1=0}^{\ell} \cdots \sum_{\ell_{T-x}=0}^{\ell} \sum_{y=0}^{m-T\ell} \binom{\ell}{\ell_1} \cdots \binom{\ell}{\ell_{T-x}} \binom{m-T\ell}{y}}_{\sum_{i=1}^{T-x} l_i + y = s_1 - x(\ell-1)} \overset{(d)}{=} \binom{m-x\ell}{s_1 - x(\ell-1)} \tag{2.12}$$

where $(d)$ follows by using Vandermonde's identity. Using (2.7), (2.11) and (2.12) thus proves the theorem.

# CHAPTER 3

# Privacy in Index Coding: $k$-Limited-Access Schemes

In this chapter, we continue our study of the index coding problem and the associated privacy concerns. We continue with the same setup: we assume an index coding instance with an adversary who is a curious client. The adversary wishes to learn information about the requests and side information sets of other clients. We discusses in the previous chapter how the adversary is able to do so by learning the coding matrix (which it naturally learns since it is a legitimate client). The key idea in the work of previous chapter was to utilize the requests and side information sets as resources for security: we showed how index codes can be designed which provides a trade-off between the privacy of one quantity at the expense of the other.

The focus of this chapter is to utilize another resource to guarantee privacy: the number of transmissions. We start our approach by an intuitive observation: a client would have less information if it has access to a fewer number of transmissions. Therefore, in this chapter, we mitigate the privacy concern by allowing each client to have limited access to the coding matrix. Keeping in mind that the adversary needs to decode its own request, we design coding matrices so that each client needs only to learn some of (and not all) the rows to decode her requested message. Designing such a scheme may require additional transmissions, which comes as a cost for the added privacy level.

First, we show that such an approach indeed increases the level of privacy. We propose to privacy metrics which we use to show that this is indeed the case. Based on this, we propose the use of $k$-limited-access schemes: given an index coding scheme that employs $T$ transmissions, we create a $k$-limited-access scheme with $T_k \geq T$ transmissions, and with the property that each client needs at most $k$ transmissions to decode her message. We derive

upper and lower bounds on $T_k$ for all values of $k$, and develop deterministic designs for these schemes, which are universal, *i.e.*, independent of the coding matrix. We show that our schemes are order-optimal when either $k$ or $n$ is large. Moreover, we propose heuristics that complement the universal schemes for the case when both $n$ and $k$ are small.

## 3.1  Introduction

It is well recognized that broadcasting can offer significant bandwidth savings compared to point-to-point communication [EK11, FLW06], and could be leveraged in several wireless network applications. Use cases include Wi-Fi (cellular) networks where an access point (a base station) is connected to a set of Wi-Fi (cellular) devices through a wireless broadcast channel, and where devices request messages, such as YouTube videos. Another use case has recently emerged in the context of distributed computing [LMY18, EKF], where worker nodes exchange data among themselves to complete computational tasks.

A canonical setup which captures the essence of broadcast channels is the index coding framework [BBJ11]. In an index coding instance, a server is connected to a set of clients through a noiseless broadcast channel. The server has a database that contains a set of messages. Each client: 1) possesses a subset of the messages that she already knows, which is referred to as the *side information set*, and 2) requests a message from the database which is not in her side information set. The server has full knowledge of the requests and side information sets of all clients. A *linear index code* (or *index code* in short)[1] is a linear coding scheme that comprises a set of coded broadcast transmissions which allow each client to decode her requested message using her side information set. The goal is to find an index code which uses the smallest possible number of broadcast transmissions. The key ingredient in designing efficient (*i.e.*, with a small number of transmissions) index codes is the use of coding across messages.

The starting observation of this work is that, using coding over broadcast channels can

---

[1]In this work, we solely focus on linear index codes.

Figure 3.1: An index coding example with 5 messages and 4 clients. Each client wants one message and has another as shown above. The optimal index code consists of sending the two transmissions $\mathbf{b}_1 + \mathbf{b}_2$ and $\mathbf{b}_3 + \mathbf{b}_4$.

cause privacy risks. In particular, a curious client may infer information about the requests and side information sets of other clients, which can be deemed sensitive by their owners. For example, consider a set of clients that use a server to download YouTube videos. Although YouTube videos are publicly available, a client requesting a video about a medical condition may not wish for others to learn her request, or learn what are other videos that she has already downloaded.

To illustrate why coding can create privacy leakage, consider the index coding instance shown in Figure 3.1. A server possesses a set of 5 messages, which we refer to as $\mathbf{b}_1$ to $\mathbf{b}_5$. The server is connected to a set of 4 clients: client 1 wants message $\mathbf{b}_1$ and has as side information message $\mathbf{b}_2$; client 2 wants $\mathbf{b}_2$ and has $\mathbf{b}_1$; client 3 wants $\mathbf{b}_3$ and has $\mathbf{b}_4$; and client 4 wants $\mathbf{b}_4$ and has $\mathbf{b}_3$. In this case, an optimal (*i.e.*, with the minimum number of transmissions) index code consists of sending 2 transmissions, namely $\mathbf{b}_1 + \mathbf{b}_2$ and $\mathbf{b}_3 + \mathbf{b}_4$: it is easy to see that each client can decode the requested message from one of these transmissions using the side information. However, this index code can allow curious clients to violate the privacy of other clients who share the broadcast channel, by learning information that pertains to their requests and/or side information sets. For example, assume that client 4 is curious. Upon learning the two transmissions, client 4 knows that nobody is requesting message $\mathbf{b}_5$.

35

Moreover, she knows that if a client is requesting $\mathbf{b}_1$ or $\mathbf{b}_2$ (similarly, $\mathbf{b}_3$ or $\mathbf{b}_4$), then this client should have the other message as side information in order to decode the requested message.

The solution that we propose to limit this privacy leakage stems from the following observation: it may not be necessary to provide clients with the entire set of broadcast transmissions. Instead, each client can be given access, and learn the coding operations, for only a subset of the transmissions, *i.e.,* the subset that would allow her to decode the message that she requested. Consider again the example in Figure 3.1. The optimal index code consists of two transmissions. However, each client is able to decode her request using exactly one of the two transmissions. Therefore, if each client only learns the coding coefficients for the transmission that she needs, then she will have no knowledge of the content of the other transmission, and thus would have less information about the requests of the other clients. Limiting the access of each client to just one out of the two transmissions was possible for this particular example; however, it is not the case that every index code has this property.

Our approach in this chapter builds on the idea described above. In particular, given an index coding instance that uses $T$ transmissions, we ask: Can we limit the access of each client to at most $k \leq T$ transmissions, while still allowing each client to decode her requested message? In other words, for a given index coding instance, what is the best (in terms of number of transmissions) index code that we can design such that each client is able to decode her request using at most $k$ out of these transmissions? Our work attempts to understand the fundamental relation between limiting the accessibility of clients to the coding matrix and the attained level of privacy. In particular, we propose the use of *k-limited-access schemes*, that transform the coding matrix so as to restrict each client to access at most $k$ rows of the transformed matrix, as opposed to the whole of it. Our contributions include:

- We formalize the intuition that using $k$-limited-access-schemes can indeed increase the attained level of privacy against curious clients. We demonstrate this using two privacy metrics, namely an *entropy-based* metric and the *maximal information leakage*. In both cases, we show that the attained level of privacy is linearly dependent on the value of

$k$, *i.e.,* privacy increases linearly with the number of rows of the coding matrix that we hide.

- We design polynomial time (in the number of clients) universal $k$-limited-access schemes (*i.e.*, that do not depend on the structure of the coding matrix), and require a simple matrix multiplication. We prove that these schemes are order-optimal in some regimes, in particular when either $k$ or $n$ (the number of clients) is large. Interestingly, when $k$ is larger than a threshold, these schemes enable to restrict the amount of access to half of the coding matrix with an overhead of exactly one additional transmission. This result indicates that some privacy-bandwidth trade-off points can be achieved with minimal overhead.

- We propose algorithms that depend on the structure of the coding matrix and show that, when $n$ and $k$ are both small, they provide improved performance with respect to the universal schemes mentioned above. These schemes use a graph-theory representation of the problem, and are optimal for some special instances.

- We provide analytical and numerical performance evaluations of our schemes. We show how our proposed $k$-limited-access schemes provide a bandwidth-privacy trade-off, namely how much bandwidth usage (*i.e.*, number of transmissions) is needed to achieve a certain level of privacy (captured by the value of $k$). We show that our proposed schemes provide a trade-off curve that is close to the lower bound when either $k$ or $n$ is large. In the case where both $n$ and $k$ are small, we show through numerical evaluations that our proposed algorithms give an average performance that is close to the lower bound.

The chapter is organized as follows. Section 3.2 introduces our notation, formulates the problem, and gives a geometric interpretation. Section 3.3 discusses how $k$-limited-access schemes limit the privacy leakage. Section 3.4 shows the construction of $k$-limited-access schemes and proves their order-optimality when either $n$ or $k$ is large. Section 3.5 designs algorithms which are better-suited for cases when both $n$ and $k$ are small. Section 3.6 dis-

cusses related work and Section 3.7 concludes the chapter. Some of the proofs are delegated to the appendices.

## 3.2 Problem Formulation and Geometric Interpretation

**Index Coding.** We consider an index coding instance, where a server has a database $\mathcal{B}$ of $m$ messages $\mathcal{B} = \{\mathbf{b}_\mathcal{M}\}$, where $\mathcal{M} = [m]$ is the set of message indices, and $\mathbf{b}_j \in \mathbb{F}_2^F, j \in \mathcal{M}$, with $F$ being the message size, and where operations are done over the binary field. The server is connected through a broadcast channel to a set of clients $\mathcal{C} = \{c_\mathcal{N}\}$, where $\mathcal{N} = [n]$ is the set of client indices. We assume that $m \geq n$. Each client $c_i, i \in \mathcal{N}$, has a subset of the messages $\{\mathbf{b}_{\mathcal{S}_i}\}$, with $\mathcal{S}_i \subset \mathcal{M}$, as side information and requests a new message $\mathbf{b}_{q_i}$ with $q_i \in \mathcal{M} \setminus \mathcal{S}_i$ that she does not have. We assume that the server employs a *linear code, i.e.,* it designs a set of broadcast transmissions that are linear combinations of the messages in $\mathcal{B}$. The linear index code can be represented as $\mathbf{AB} = \mathbf{Y}$, where $\mathbf{A} \in \mathbb{F}_2^{T \times m}$ is the coding matrix, $\mathbf{B} \in \mathbb{F}_2^{m \times F}$ is the matrix of all the messages and $\mathbf{Y} \in \mathbb{F}_2^{T \times F}$ is the resulting matrix of linear combinations. Upon receiving $\mathbf{Y}$, client $c_i, i \in \mathcal{N}$, employs linear decoding to decode the requested message $\mathbf{b}_{q_i}$.

**Problem Formulation.** In [BBJ11], it was shown that the index coding problem is equivalent to the rank minimization of an $n \times m$ matrix $\mathbf{G} \in \mathbb{F}_2^{n \times m}$, whose $i$-th row $\mathbf{g}_i$, $i \in [n]$, has the following properties: (i) has a 1 in the position $q_i$ (*i.e.,* the index of the message requested by client $c_i$), (ii) has a 0 in the $j$-th position for all $j \in \mathcal{M} \setminus \mathcal{S}_i$, (iii) can have either 0 or 1 in all the remaining positions. For instance, with reference to the example in Figure 3.1, we would have

$$\mathbf{G} = \begin{bmatrix} 1 & \star & 0 & 0 & 0 \\ \star & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & \star & 0 \\ 0 & 0 & \star & 1 & 0 \end{bmatrix},$$

where $\star$ can be either 0 or 1. It was shown in [BBJ11] that finding an optimal linear coding scheme *i.e.,* with minimum number of transmissions) is equivalent to completing $\mathbf{G}$ (*i.e.,*

assign values to the $\star$ components of $\mathbf{G}$) so that it has the minimum possible rank. Once we have completed $\mathbf{G}$, we can use a basis of the row space of $\mathbf{G}$ (of size $T = \text{rank}(\mathbf{G})$) as a coding matrix $\mathbf{A}$. In this case, client $c_i$ can construct $\mathbf{g}_i$ as a linear combination of the rows of $\mathbf{A}$, i.e., $c_i$ performs the decoding operation $\mathbf{d}_i \mathbf{A} \mathbf{B} = \mathbf{d}_i \mathbf{Y}$, where $\mathbf{d}_i \in \mathbb{F}_2^T$ is the decoding row vector of $c_i$ chosen such that $\mathbf{d}_i \mathbf{A} = \mathbf{g}_i$. Finally, client $c_i$ can successfully decode $\mathbf{b}_{q_i}$ by subtracting from $\mathbf{d}_i \mathbf{Y}$ the messages corresponding to the non-zero entries of $\mathbf{g}_i$ (other than the requested message). We remark that any linear index code that satisfies all clients with $T$ transmissions (where $T$ is not necessarily optimal) – and can be obtained by any index code design algorithm [ELH, HE15, CS08] – corresponds to a completion of $\mathbf{G}$ (i.e., given $\mathbf{A} \in \mathbb{F}_2^{T \times m}$, we can create a corresponding $\mathbf{G}$ in polynomial time).

In our problem formulation we assume that we start with a given matrix $\mathbf{G}$ of rank $T$, i.e., we are given $n$ *distinct* vectors that belong to a $T$-dimensional subspace. Using a basis of the row space of the given $\mathbf{G}$, we construct $\mathbf{A} \in \mathbb{F}_2^{T \times m}$. Then, we ask: *Given $n$ distinct vectors $\mathbf{g}_i$, $i \in [n]$, in a $T$-dimensional space, can we find a minimum-size set $\mathcal{A}_k$ with $T_k \geq T$ vectors, such that each $\mathbf{g}_i$ can be expressed as a linear combination of at most $k$ vectors in $\mathcal{A}_k$ (with $1 \leq k \leq T$)?* The vectors in $\mathcal{A}_k$ form the rows of the coding matrix $\mathbf{A}_k$ that we will employ. Then by definition, client $c_i$ will be able to reconstruct $\mathbf{g}_i$ using the matrix $\mathbf{A}_k^{(i)} \subset_k \mathbf{A}_k$. We can equivalently restate the question as follows: *Given a coding matrix $\mathbf{A}$, can we find $\mathbf{P} \in \mathbb{F}_2^{T_k \times T}$, with $T_k$ as small as possible, such that $\mathbf{A}_k = \mathbf{P} \mathbf{A}$ and each row of $\mathbf{G}$ can be reconstructed by combining at most $k$ rows of $\mathbf{A}_k$?* Note that $k = T$ corresponds to the conventional transmission scheme of an index coding problem for which $\mathbf{P} = \mathbf{I}_T$. In the remainder of the chapter we will refer to a scheme that chooses $\mathbf{A}_k$ to be the coding matrix as $k$-limited-access scheme.

**Transmission Protocol.** In order to realize the privacy benefits of using $k$-limited-access schemes – which we will thoroughly illustrate in Section 3.3 – we propose a different transmission protocol for the index coding setup. Figure 3.2 shows both the conventional and the proposed transmission protocols. In the conventional protocol, the server designs a set of $T$ packets, each corresponding to an equation from the set of equations $\mathbf{A} \mathbf{B} = \mathbf{Y}$. As shown in Figure 3.2(a), packet $i \in [T]$ consists of (i) a payload which contains the linear combi-

F − bit messages
F − bit messages

(a) Conventional Transmission Protocol.　　(b) Proposed Transmission Protocol.

Figure 3.2: A comparison between the conventional and the proposed transmission protocols. The proposed transmission protocol incurs a negligible increase in the transmission overhead when both $n$ and $m$ are $o(F)$.

nation $\mathbf{y}_i$ and (ii) a header which contains the coefficients $\mathbf{a}_i$ used to create the equation. In the conventional protocol, the server sends these packets (both headers and payloads) on the broadcast channel to all clients. Our proposed protocol, however, operates differently. Specifically, the server generates packets which correspond to the set of equations $\mathbf{A}_k\mathbf{B} = \mathbf{Y}_k$ in a way that is similar to the conventional protocol. The server then sends *only* the payloads of these packets on the broadcast channel. Differently, the server sends the coefficients corresponding to *only* $\mathbf{A}_k^{(i)} \subset_k \mathbf{A}_k$ to client $c_i$ using a private key or on a dedicated private channel (*e.g.,* the same channel used by $c_i$ to convey her request to the server). Thus, using a $k$-limited-access scheme incurs an extra transmission overhead to privately convey the coding vectors. In particular, the total number of transmitted bits $C_k$ can be upper bounded as $C_k \leq nkm + T_kF$, while the total number of transmitted bits $C$ using a conventional scheme is $C = T(F + m)$. The extra overhead incurred is negligible in comparison to the broadcast transmissions that convey the encoded messages when $n$ and $m$ are both $o(F)$, which is a reasonable assumption for large file sizes (for instance, when sharing YouTube videos).

**Geometric Interpretation.** The geometric interpretation of our problem is depicted in Figure 3.3. An index code $\mathbf{A}$ corresponds to a particular completion of the matrix $\mathbf{G}$. Therefore, the set of row vectors in $\mathbf{G}$ lies in the row span of $\mathbf{A}$ (which is of dimension $T$). We denote this subspace of dimension $T$ by $L$. The problem of finding a matrix $\mathbf{A}_k$ can be

40

Figure 3.3: A geometric interpretation of $k$-limited-access schemes. An index code $\mathbf{A}$ is obtained from a particular filling of the matrix $\mathbf{G}$. Therefore, the collection of row vectors of $\mathbf{G}$ lies in the span of $\mathbf{A}$. Finding $\mathbf{A}_k$ is equivalent to finding a collection of subspaces, each of dimension at most $k$, to cover $\mathbf{G}$. Client $c_i$ is sent a collection of (at most) $k$ rows of $\mathbf{A}_k$; these correspond to one subspace which covers $\mathbf{g}_i$.

interpreted as finding a set of subspaces, each of dimension at most $k$, such that each row vector $\mathbf{g}_i$, $i \in [n]$, is covered by at least one of these smaller subspaces. Once these subspaces are selected, then the rows of $\mathbf{A}_k$ are taken as the union of the basis vectors of all these subspaces. Client $c_i$ is then given the basis vectors of subspace $L_i$, *i.e.*, the one which covers $\mathbf{g}_i$, instead of the whole matrix $\mathbf{A}_k$. Therefore $c_i$ would have perfect knowledge of $L_i$ instead of $L$. Having less information about $L$ naturally translates to less information about the requests of other clients, as we more formally discuss in the next section.

## 3.3  Achieved Privacy Levels

In this section, we investigate and quantify the level of privacy that $k$-limited-access schemes can achieve compared to a conventional index coding scheme (*i.e.,* when each client has access to the entire coding matrix). In what follows, we consider the setup described in the previous section and suppose that client $c_n$ is curious, *i.e.,* by leveraging the (at most) $k$ rows $\mathbf{A}_k^{(n)}$ that she receives, she seeks to infer information about client $c_i, i \in [n-1]$. Specifically,

$Q_{[n]}, S_{[n]}$ → [Index Coding Algorithm] → $A$ → [$k$-limited-access Scheme] → $A_k$ → [Receiver $i$] → $A_k^{(i)}$

Figure 3.4: The procedure of designing an index code and applying $k$-limited-access schemes.

we are interested in quantifying the amount of information that $c_n$ can obtain about $q_i$ (*i.e.*, the identity of the request of $c_i$) as a function of $k$.

We assume that the index coding instance is random, *i.e.*, we consider the requests and side information sets of clients as random variables and denote them as $Q_{[n]}$ and $S_{[n]}$, respectively. The operation of the server is shown in Figure 3.4 and is described as follows:

**Step-1:** The server obtains the information about the requests $Q_{[n]}$ and side information sets $S_{[n]}$ of all clients $c_{[n]}$.

**Step-2:** Based on this information, the server designs an index code $\mathbf{A}$ by means of some index coding algorithm [ELH, HE15, CS08].

**Step-3:** The server then applies the $k$-limited-access scheme to obtain $\mathbf{A}_k = \mathbf{PA}$, where $\mathbf{P}$ is a deterministic mapping from $\mathbf{A}$ to $\mathbf{A}_k$ (see Section 3.4 for the construction of $\mathbf{P}$). This implies that $T_k$ is a deterministic function of $T$ and $k$ (*i.e.*, the parameter of the scheme).

**Step-4:** The server sends $\mathbf{A}_k^{(i)}$ to client $c_i$. If multiple $\mathbf{A}_k^{(i)}$ can be selected, then the server picks and transmits one such matrix uniformly at random, independently of the underlying $\mathbf{A}$ which might have generated this $\mathbf{A}_k$.

We are now interested in quantifying the level of privacy that is achieved by the protocol described above. Towards this end, we use two privacy metrics, namely an *entropy-based* metric and the *maximal information leakage*.

### 3.3.1 Entropy-Based Privacy Metric

The entropy-based privacy metric is inspired by the geometric interpretation of our problem in Figure 3.3. We let $L$ (respectively, $L_n$) be the random variable associated with the

subspace spanned by the $T$ rows of the coding matrix $\mathbf{A}$ (respectively, spanned by the $k$ row vectors of $\mathbf{A}_k^{(n)}$). Client $c_n$ receives the matrix $\mathbf{Y}_k$ and as such she knows $T_k$. Given this, we now define the entropy-based privacy metric and evaluate it for the proposed protocol.

**Definition 3.3.1.** The entropy-based privacy metric is defined as

$$P_k^{(\text{Ent})} = H\left(L | L_n, T_k\right),$$

and quantifies the amount of uncertainty that $c_n$ has about the subspace spanned by the $T$ rows of the index coding matrix $\mathbf{A}$.

Before characterizing $P_k^{(\text{Ent})}$, we state the following lemma, which is proved in Appendix 3.8.1.

**Lemma 3.3.1.** *Given a subspace $L_n \subseteq \mathbb{F}_2^m$ of dimension $k$, let $\mathcal{L}(T, L_n)$ be the set of subspaces $L \subseteq \mathbb{F}_2^m$ of dimension $T \geq k$ where $L_n \subseteq L$. Then $|\mathcal{L}(T, L_n)|$ is equal to*

$$|\mathcal{L}(T, L_n)| = \prod_{\ell=0}^{T-k-1} \frac{2^m - 2^{k+\ell}}{2^T - 2^{k+\ell}}.$$

Assume an index coding setting with $c_n$ observing a particular subspace $L_n = \ell_n$ and a number of transmissions $T_k = t_k$ for the $k$-limited access scheme. Moreover, we consider a stronger adversary (*i.e.*, curious client) and assume that she also knows the specific realization of $T = t$. Given this, we can compute

$$
\begin{aligned}
P_k^{(\text{Ent})} &= H\left(L | L_n = \ell_n, T_k = t_k, T = t\right) \\
&\overset{(a)}{=} H\left(L | L_n = \ell_n, T = t\right) \overset{(b)}{=} \log\left(|\mathcal{L}(t, \ell_n)|\right) \\
&\overset{(c)}{=} \log\left(\prod_{\ell=0}^{t-k-1} \frac{2^m - 2^{k+\ell}}{2^t - 2^{k+\ell}}\right) \overset{m \gg t}{\approx} m(t - k),
\end{aligned}
\tag{3.1}
$$

where: (i) the equality in (a) follows because $T_k$ is a deterministic function of $T$ and $k$, which is the parameter of the scheme (see Step-3); (ii) the equality in (b) follows by assuming that the underlying system maintains a uniform distribution across all feasible $t$-dimensional subspaces of $\mathbb{F}_2^m$; (iii) the equality in (c) follows by virtue of Lemma 3.3.1. We note that when $m \gg t$, then the quantity in (3.1) decreases linearly with $k$, *i.e.*, as intuitively expected, the

less rows of the coding matrix $c_n$ learns, the less she can infer about the subspace spanned by the $T$ rows of the coding matrix $\mathbf{A}$. This suggests that, by increasing $k$, $c_n$ has less uncertainty about $q_i$. Note also that $P_k^{(\text{Ent})}$ is zero when $k = t$; this is because, under this condition, $c_n$ receives the entire index coding matrix, $i.e.$, $L_n = L$, and hence she is able to perfectly reconstruct the subspace spanned by its rows. However, although $P_k^{(\text{Ent})} = 0$ when $k = t$, $c_n$ might still have uncertainty about $q_i$ [KSC17]. Quantifying this uncertainty is an interesting open problem; this uncertainty, in fact, depends on the underlying system, $e.g.$, on the index code used by the server and on the distribution with which the index code matrix is selected.

### 3.3.2 Maximal Information Leakage

The second metric that we consider as our privacy metric is the Maximal Information Leakage (MIL) [IKW16]. Given two discrete random variables $X$ and $Y$ with alphabets $\mathcal{X}$ and $\mathcal{Y}$, the MIL from $X$ to $Y$ is denoted by $\mathcal{L}(X \to Y)$ and defined as

$$\mathcal{L}(X \to Y) = \sup_{S-X-Y} \log \frac{\sum_{y \in \mathcal{Y}} \max_{s \in \mathcal{S}} p_{SY}(s, y)}{\max_{s \in \mathcal{S}} p_S(s)} = \log \sum_{y \in \mathcal{Y}} \max_{x \in \mathcal{X}: p_X(x) > 0} p_{Y|X}(y|x), \qquad (3.2)$$

where the second equality is shown in [IKW16]. The MIL metric captures the amount of information leaked about $X$ through $Y$ to an adversary, who is interested in estimating a (possibly probabilistic) function $S$ of $X$. This is captured by the fact that $S - X - Y$ forms a Markov chain as shown in the expression in (3.2). The metric considers a worst-case such adversary, that is, an adversary who is interested in computing a function $S$ for which the maximum information can be leaked out of $Y$. The result in [IKW16] shows that this quantity depends only on the joint distribution of $X$ and $Y$. The following properties of the MIL are useful [IKW16]:

- (Property 1): If $X - Y - Z$, then $\mathcal{L}(X \to Z) \leq \min\{\mathcal{L}(X \to Y), \mathcal{L}(Y \to Z)\}$,

- (Property 2): $\mathcal{L}(X \to Y) \leq \min\{\log |\mathcal{X}|, \log |\mathcal{Y}|\}$,

- (Property 3): $\mathcal{L}(X \to X) = \log |\{x : p_X(x) > 0\}|$.

To describe how we use the MIL as a privacy metric in our setup, we first need to define what are the corresponding random variables $X$ and $Y$, and then argue that the estimation of client $c_n$ of the requests of other clients forms a Markov chain as required by the MIL definition. To do so, we first define the following sets:

*1)* Given $\mathbf{g}_i$, $\mathbf{A}_k$ and an integer $r$, let $\mathcal{P}(\mathbf{g}_i, \mathbf{A}_k, r)$ be the set of all possible sub-matrices $\mathbf{A}_k^{(i)}$ of $\mathbf{A}_k$ with *exactly* $r$ rows, that client $c_i$ can use to reconstruct the vector $\mathbf{g}_i$:

$$\mathcal{P}(\mathbf{g}_i, \mathbf{A}_k, r) = \left\{ \mathbf{Z} \subset_r \mathbf{A}_k \mid \exists \mathbf{d} \in \mathbb{F}_2^r \text{ s.t. } \mathbf{g}_i = \mathbf{d}\mathbf{Z} \right\},$$

*2)* Given $q_i$, $\mathcal{S}_i$ and $\mathbf{A}_k$, let $\mathcal{T}(q_i, \mathcal{S}_i, \mathbf{A}_k)$ be the set of all possible sub-matrices $\mathbf{A}_k^{(i)}$ of $\mathbf{A}_k$ with the minimum possible number of rows, such that client $c_i$ with side information $\mathcal{S}_i$ can decode $q_i$:

$$\mathcal{T}(q_i, \mathcal{S}_i, \mathbf{A}_k) = \bigcup_{\mathbf{g}_i \in \mathcal{G}(q_i, \mathcal{S}_i)} \mathcal{P}(\mathbf{g}_i, \mathbf{A}_k, r_{\min}),$$

where

$$\mathcal{G}(q_i, \mathcal{S}_i) = \left\{ \mathbf{g} \in \mathbb{F}_2^m \mid g_{q_i} = 1, g_{[m] \setminus \{q_i \cup \mathcal{S}_i\}} = 0 \right\},$$

and

$$r_{\min} = \min \mathcal{R}, \ \mathcal{R} = \left\{ r \in \mathbb{N}^+ : \ \exists \mathbf{g}_i \in \mathcal{G}(q_i, \mathcal{S}_i) \text{ such that } \mathcal{P}(\mathbf{g}_i, \mathbf{A}_k, r) \neq \emptyset \right\}.$$

Since the requests and the side information sets are considered as random variables, then all subsequently generated codes, namely $\mathbf{A}$, $\mathbf{A}_k$ and $\mathbf{A}_k^{(i)}$ can be treated as random variables as well. We denote the corresponding random variables of these quantities as $A$, $A_k$ and $A_k^{(i)}$ respectively. In other words, for a given realization of $Q_{[n]} = q_{[n]}$ and $S_{[n]} = \mathcal{S}_{[n]}$, the corresponding realizations of the aforementioned codes used by the server are $A = \mathbf{A}$, $A_k = \mathbf{A}_k$ and $A_k^{(i)} = \mathbf{A}_k^{(i)}$.

When using conventional index codes (*i.e.,* without $k$-limited-access schemes), client $c_n$ (*i.e.,* the curious client and hence the adversary) would try to infer information about $Q_{[n-1]}$ from observing $A$ and given her information of $Q_n, S_n$. Therefore, one can think of client $c_n$ estimate of $Q_{[n-1]}$ as being a particular estimation function, the input of which is $A$. Differently, after using $k$-limited-access schemes, client $c_n$ would only have observed $A_k^{(n)}$

instead of $A$. Therefore, in the context of MIL, one choice of the variables $X$ and $Y$ is $A$ and $A_k^{(n)}$ respectively. The function $S$ would therefore be client $c_n$'s estimate of $Q_{[n-1]}$ out of $A$. The following proposition shows that this choice of variables $X$, $Y$ and $S$ allows us to use the MIL as a metric.

**Proposition 3.3.2.** *The following Markov chain holds*

$$Q_{[n-1]} - A - A_k - A_k^{(n)}, \tag{3.3}$$

*conditioned on the knowledge of $Q_n, S_n$ in every stage of the chain.*

**Proof:** We have the following:

- $Q_{[n-1]} - A - A_k$ holds since $A_k$ is a deterministic function of $A$ (see also Step-3 of the proposed protocol);

- $A - A_k - A_k^{(n)}$ holds since $p(A_k^{(n)}|A_k, Q_n, S_n) = 1/|\mathcal{T}(Q_n, S_n, A_k)|$, independent of $A$, as described in Step-4 of the proposed protocol.

■

We define $P_k^{(\text{MIL})} = \mathcal{L}\left(A \to A_k^{(n)}|Q_n = q_n, S_n = \mathcal{S}_n\right)$ as our MIL privacy metric[2]. The quantity $P_k^{(\text{MIL})}$ gives the maximum amount of information that $c_n$ can extract about $Q_{[n-1]}$ given the knowledge of $Q_n, S_n$. The following theorem – proved in Appendix 3.8.2 – provides a guarantee on $P_k^{(\text{MIL})}$.

**Theorem 3.3.3.** *Using the MIL, the attained level of privacy against a curious client when $k$-limited-access schemes are used is*

$$P_k^{(MIL)} = O(|\mathcal{S}_n| + mk). \tag{3.4}$$

The quantity in (3.4) characterizes the maximum amount of information that can be leaked to a curious client when $k$-limited-access schemes are used. It is clear that decreasing $k$ would decrease this amount of information; this aligns with the intuition that the less rows

---

[2]We use the notation $\mathcal{L}(X \to Y|Z)$ to denote that the variables $X$ and $Y$ are conditioned on $Z$.

Figure 3.5: This figure shows how the MIL privacy metrics compare for the conventional index coding schemes and the $k$-limited-access schemes. Taking $k = o(T)$ would guarantee privacy gains when using $k$-limited-access schemes.

a server gives to a client, the less information a client would be able to infer about other clients sharing the broadcast domain. In order to shed more light on the benefits of using $k$-limited-access schemes, one could compare the quantity $P_k^{(\text{MIL})}$ with the MIL obtained when $k$-limited-access schemes are not used, *i.e.*, when a client observes the whole matrix $A$. Let this quantity be denoted as $\bar{P}_k^{(\text{MIL})} = \mathcal{L}(A \to A | Q_n = q_n, S_n = \mathcal{S}_n)$. Then we have the following result, which is proved in Appendix 3.8.3.

**Theorem 3.3.4.** *Using the MIL, the attained level of privacy against a curious client for a conventional index coding setup is*

$$\bar{P}_k^{(MIL)} = \Omega\left(mT - T^2\right). \tag{3.5}$$

The results in Theorem 3.3.3 and Theorem 3.3.4 can be interpreted with the help of Figure 3.5. The $k$-limited-access schemes achieve privacy gains as compared to conventional index codes, when the two bounds in (3.4) and (3.5) strictly mismatch. A sufficient (but not necessary) condition for this is to select $k = o(T)$.

## 3.4   Construction of $k$-limited-access Schemes

In this section, we focus on designing $k$-limited-access schemes and assessing their theoretical performance in terms of number of additional transmissions required with respect to a conventional index coding scheme. Recall that we are given a coding matrix $\mathbf{A}$ that requires $T$ transmissions. Then, we seek to construct a matrix $\mathbf{P} \in \mathbb{F}_2^{T_k \times T}$, so that $\mathbf{A}_k = \mathbf{P}\mathbf{A}$, and each client needs to access at most $k$ rows of $\mathbf{A}_k$ to decode her requested message. In particular,

47

we aim at constructing matrices $\mathbf{P}$ with $T_k$ as small as possible. Trivially, $T_k \geq T$. Towards this end, we first derive upper and lower bounds on $T_k$. Our main result is stated in the theorem below.

**Theorem 3.4.1.** *Given an index coding matrix* $\mathbf{A} \in \mathbb{F}_2^{T \times m}$ *with* $T \geq 2$, *it is possible to transform it into* $\mathbf{A}_k = \mathbf{PA}$ *with* $\mathbf{P} \in \mathbb{F}_2^{T_k \times T}$, *such that each client can decode her requested message by combining at most* $k$ *rows of* $\mathbf{A}_k$, *if and only if*

$$T_k \geq \max\left\{T, T^\star\right\}, \quad T^\star = \min\left\{T_k : \sum_{i=1}^{k} \binom{T_k}{i} \geq n\right\}. \tag{3.6}$$

*Moreover, we provide polynomial time (in n) constructions of* $\mathbf{P}$ *such that:*

- *When* $\lceil T/2 \rceil \leq k < T$, *then*

$$T_k \leq \min\left\{n, T+1\right\}; \tag{3.7}$$

- *When* $1 \leq k < \lceil T/2 \rceil$, *then*

$$T_k \leq \min\left\{n, k2^{\lceil \frac{T}{k} \rceil}\right\}. \tag{3.8}$$

**Proof:** The lower bound on $T_k$ in (3.6) is proved in Appendix 3.8.4. In particular, the bound in (3.6) says that, if we are allowed to combine at most $k$ out of the $T_k$ vectors, then we should be able to create a sufficient number of vectors. The two upper bounds on $T_k$ in (3.7) and (3.8) are proved in Section 3.4.1, where we give explicit constructions for $\mathbf{P}$. ∎

We note that, as expected, the smaller the value of $k$ that we require, the larger the value of $T_k$ that we need to use. Trivially, for $k = 1$ we would need $T_k = n$, *i.e.*, the server would need to send uncoded transmissions. Thus, there is a trade-off between the bandwidth – measured as the number $T_k$ of broadcast transmissions – and privacy – captured by the value of $k$ that we require. Interestingly, when $k \geq \lceil T/2 \rceil$, with just one extra transmission, *i.e.*, $T_k = T + 1$, we can restrict the access of each client to at most half of the coding matrix, independently of the coding matrix $\mathbf{A}$. In other words, for this regime, we can achieve a certain level of privacy with minimal overhead. However, as we further reduce the value of $k$, the overhead becomes more significant. Moreover, the results in Theorem 3.4.1 also imply

|  (a) $n = 2^T - 1$ | (b) $n = T^4$ | (c) $n = T^2$ |

Figure 3.6: Bandwidth ($T_k$ on the y-axis) versus privacy ($k$ on the x-axis) trade-off when using the $k$-limited-access schemes in Theorem 3.4.1 for different values of $n$. The plots in this figure are for $T = 20$.

that our constructions are order-optimal in the case of large values of $n$ (when $n = \Theta(2^T))^3$. In addition, when $\lceil T/2 \rceil \leq k < T$, our scheme is at most one transmission away from the optimal number of transmissions, and this is for *any* value of $n$. This is shown in the following lemma, which is proved in Appendix 3.8.4.

**Lemma 3.4.2.** *Consider an index coding setup. We have*

- *When $n = 2^T - 1$ and $\lceil T/2 \rceil \leq k < T$, the bounds in (3.6) and (3.7) coincide, i.e., the provided construction of $\mathbf{P}$ is optimal;*

- *For any value of $n < 2^T - 1$ and $\lceil T/2 \rceil \leq k < T$, the bound in (3.7) is at most one transmission away from the bound in (3.6);*

- *When $n = \Theta(2^T)$ and for any value of $k$, then $T_k = \Theta(k2^{\frac{T}{k}})$, i.e., the provided construction is order-optimal.*

Figure 3.6 shows the trade-off exhibited by our proposed $k$-limited-access schemes between bandwidth usage ($T_k$) and the attained privacy ($k$) - we use $k$ as a proxy to the amount of attained privacy against a curious client (see Section 3.3). The figure shows the

---

[3]Note that $n$ is always $O(2^T)$ (*i.e.*, the number of distinct vectors $\mathbf{g}_i$ for a given $T$ is at most $2^T - 1$). The case of large values of $n$ corresponds to the case where this bound on the number of distinct vectors $\mathbf{g}_i$ is not loose: there is a corresponding lower bound on $n$, *i.e.*, $n = \Omega(2^T)$. Therefore, the case of large values of $n$ corresponds to $n = \Theta(2^T)$.

performance of our constructions in Theorem 3.4.1 (labeled as Scheme-1), as well as the lower bound in (3.6) (labeled as LB) and an upper bound which corresponds to uncoded transmissions (labeled as UB). Figure 3.6(a) confirms the order-optimality of our constructions when $n = 2^T - 1$. In addition, our schemes perform similarly well when $n$ is sufficiently large (and not necessarily equal to $2^T - 1$) as shown in Figure 3.6(b) where $n = T^4$. Finally, Figure 3.6(c) shows the performance for a small value of $n$ ($n = T^2$). The figure shows that our proposed constructions do not perform as well when $n$ and $k$ are small, a case which we study in more details in Section 3.5.

We now conclude this section by giving explicit constructions of the $\mathbf{P}$ matrix and prove the two upper bounds on $T_k$ in (3.7) and (3.8). Our design of $\mathbf{P}$ allows to reconstruct any of the $2^T$ vectors of size $T$. As such our constructions are universal, in the sense that the matrix $\mathbf{P}$ that we construct does not depend on the specific index coding matrix $\mathbf{A}$.

### 3.4.1 Proof of Theorem 3.4.1, Equations (3.7) and (3.8)

Recall that $\mathbf{A}$ is full rank and that the $i$-th row of $\mathbf{G}$ can be expressed as $\mathbf{g}_i = \mathbf{d}_i\mathbf{A}$, where $\mathbf{d}_i \in \mathbb{F}_2^T$ is the coefficients row vector associated with $\mathbf{g}_i$. We next analyze two different cases/regimes, which depend on the value of $k$.

**Case I:** $\lceil T/2 \rceil \leq k < T$. When $n \geq T + 1$, let

$$\mathbf{P} = \begin{bmatrix} \mathbf{I}_T \\ \mathbf{1}_T \end{bmatrix}, \tag{3.9}$$

which results in a matrix $\mathbf{A}_k$ with $T_k = T + 1$, matching the bound in (3.7). We now show that each $\mathbf{g}_i = \mathbf{d}_i\mathbf{A}, i \in [n]$, can be reconstructed by combining up to $k$ vectors of $\mathbf{A}_k$. Let $w(\mathbf{d}_i)$ be the Hamming weight of $\mathbf{d}_i$. If $w(\mathbf{d}_i) \leq \lceil T/2 \rceil$, then we can reconstruct $\mathbf{g}_i$ as $\mathbf{g}_i = [\mathbf{d}_i\ 0]\mathbf{A}_k$, which involves adding $w(\mathbf{d}_i) \leq \lceil T/2 \rceil \leq k$ rows of $\mathbf{A}_k$. Differently, if $w(\mathbf{d}_i) \geq \lceil T/2 \rceil + 1$, then we can reconstruct $\mathbf{g}_i$ as $\mathbf{g}_i = [\bar{\mathbf{d}}_i\ 1]\mathbf{A}_k$, where $\bar{\mathbf{d}}_i$ is the bitwise complement of $\mathbf{d}_i$. In this case, reconstructing $\mathbf{g}_i$ involves adding $T - w(\mathbf{d}_i) + 1 \leq \lfloor T/2 \rfloor \leq k$ rows of $\mathbf{A}_k$.

When $n < T + 1$, then it is sufficient to send $n$ uncoded transmissions, where the $i$-th

transmission satisfies $c_i, i \in [n]$. In this case $c_i$ has access only to the $i$-th transmission, $i.e.,$ $k = 1$. This completes the proof of the upper bound in (3.7).

**Example:** We show how the scheme works via a small example, where $T = 4$ and $k = 2$. In this case, we have

$$\mathbf{P} = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \\ 1 & 1 & 1 & 1 \end{bmatrix}.$$

If $\mathbf{g}_i = \begin{bmatrix} 1 & 1 & 0 & 0 \end{bmatrix} \mathbf{A}$, then it can be reconstructed as $\mathbf{g}_i = \begin{bmatrix} 1 & 1 & 0 & 0 & 0 \end{bmatrix} \mathbf{PA}$ with 2 rows of $\mathbf{PA}$ used in the reconstruction. Differently, if $\mathbf{g}_i = \begin{bmatrix} 1 & 1 & 1 & 0 \end{bmatrix} \mathbf{A}$, then it can be reconstructed as $\mathbf{g}_i = \begin{bmatrix} 0 & 0 & 0 & 1 & 1 \end{bmatrix} \mathbf{PA}$ with again 2 rows of $\mathbf{PA}$ used in the reconstruction.

**Case II:** $1 \leq k < \lceil T/2 \rceil$. Let $Q = \lfloor T / \lceil \frac{T}{k} \rceil \rfloor$ and $T_{\text{rem}} = T - Q \lceil \frac{T}{k} \rceil$. If $k$ divides $T$, then $Q = k$, $T_{\text{rem}} = 0$, otherwise $Q \leq k - 1$ and $T_{\text{rem}} \leq \lceil \frac{T}{k} \rceil$. Then, we can write

$$\mathbf{P} = \begin{bmatrix} \mathbf{Z}_1 & \mathbf{0} & \cdots & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{Z}_2 & \cdots & \mathbf{0} & \mathbf{0} \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ \mathbf{0} & \mathbf{0} & \cdots & \mathbf{Z}_Q & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \cdots & \mathbf{0} & \mathbf{Z}_{Q+1} \end{bmatrix},$$

where, for $i \in [Q]$, the matrix $\mathbf{Z}_i$, of dimension $\lambda_i \times T$, is constructed as follows

$$\mathbf{Z}_i = \begin{bmatrix} \mathbf{0}_{\lambda_i \times (i-1) \lceil \frac{T}{k} \rceil} & \bar{\mathbf{Z}}_i & \mathbf{0}_{\lambda_i \times (Q-i) \lceil \frac{T}{k} \rceil} & \mathbf{0}_{\lambda_i \times T_{\text{rem}}} \end{bmatrix},$$

where $\bar{\mathbf{Z}}_i$, of dimension $\lambda_i \times \lceil \frac{T}{k} \rceil$, has as rows all non-zero vectors of length $\lceil \frac{T}{k} \rceil$. Therefore, $\lambda_i = 2^{\lceil T/k \rceil} - 1$. Similarly, the matrix $\mathbf{Z}_{Q+1}$, of dimension $\lambda_{Q+1} \times T$, is constructed as follows

$$\mathbf{Z}_{Q+1} = \begin{bmatrix} \mathbf{0}_{\lambda_{Q+1} \times Q \lceil \frac{T}{k} \rceil} & \bar{\mathbf{Z}}_{Q+1} \end{bmatrix},$$

51

where $\bar{\mathbf{Z}}_{Q+1}$, of dimension $\lambda_{Q+1} \times T_{\text{rem}}$, has as rows all non-zero vectors of length $T_{\text{rem}}$. Therefore, $\lambda_{Q+1} = 2^{T_{\text{rem}}} - 1$.

In other words, the matrix $\mathbf{P}$ is constructed as a block-diagonal matrix, with the diagonal elements being $\bar{\mathbf{Z}}_i$ for all $i \in [Q+1]$. Therefore, equation (3.8) holds by computing

$$T_k = \sum_{i=1}^{Q+1} \lambda_i = Q\left(2^{\lceil \frac{T}{k} \rceil} - 1\right) + 2^{T_{\text{rem}}} - 1 \le k2^{\lceil \frac{T}{k} \rceil}.$$

What remains is to show that any vector $\mathbf{g}_i, i \in [n]$, can be reconstructed by adding at most $k$ vectors of $\mathbf{P}$. To show this, we prove that any vector $\mathbf{v} \in \mathbb{F}_2^T$ can indeed be constructed with the proposed design of $\mathbf{P}$. We note that we can express the vector $\mathbf{v}$ as $\mathbf{v} = [\mathbf{v}_1 \cdots \mathbf{v}_{Q+1}]$, where $\mathbf{v}_i, i \in [Q]$ are parts of the vector $\mathbf{v}$ each of length $\lceil \frac{T}{k} \rceil$, while $\mathbf{v}_{Q+1}$ is the last part of $\mathbf{v}$ of length $T_{\text{rem}}$. Then, we can write $\mathbf{v} = \sum_{i \in \mathcal{K}(\mathbf{v})} \bar{\mathbf{v}}_i$, where $\bar{\mathbf{v}}_i = \left[\mathbf{0}_{(i-1)\lceil \frac{T}{k} \rceil} \quad \mathbf{v}_i \quad \mathbf{0}_{(Q-i)\lceil \frac{T}{k} \rceil} \mathbf{0}_{T_{\text{rem}}}\right]$ for $i \in [Q]$, $\bar{\mathbf{v}}_{Q+1} = \left[\mathbf{0}_{Q\lceil \frac{T}{k} \rceil} \quad \mathbf{v}_{Q+1}\right]$ and $\mathcal{K}(\mathbf{v}) \subseteq [Q+1]$ is the set of indices for which $\mathbf{v}_i$ is not all-zero. According to the construction of $\mathbf{P}$, for all $i \in \mathcal{K}(\mathbf{v})$, the corresponding vector $\mathbf{v}_i$ is one of the rows in $\mathbf{Z}_i$. The proof concludes by noting that $|\mathcal{K}(\mathbf{v})| \le k$. This is true because, if $k$ does not divide $T$, then $Q \le k - 1$; otherwise, $Q = k$ but $T_{\text{rem}} = 0$ (*i.e.*, $\mathbf{v}_{Q+1}$ does not exist), therefore $\mathcal{K}(\mathbf{v}) \subseteq [k]$. This completes the proof of the upper bound in (3.8).

**Example:** We show how the scheme works via a small example, where $T = 8$ and $k = 3$. For this particular example, we have $Q = \lfloor T/\lceil \frac{T}{k} \rceil \rfloor = 2$ and $T_{\text{rem}} = T - Q\lceil \frac{T}{k} \rceil = 2$. Thus, the idea is that, to reconstruct a vector $\mathbf{v} \in \mathbb{F}_2^8$, we treat $\mathbf{v}$ as $k = 3$ disjoint parts; the first 2 are of length $\lceil \frac{T}{k} \rceil = 3$ and the remaining part is of length $T_{\text{rem}} = 2$. We then construct $\mathbf{P}$ as $k = 3$ disjoint sections, where each section allows us to reconstruct one part of the vector.

(a) $k = 2$          (b) $k = 5$

Figure 3.7: Performance of the scheme in Theorem 3.4.1 (referred to as Scheme-1) for different values of $n$, compared against the lower bound LB in equation (3.6) and the upper bound UB of sending uncoded transmissions - $T = 20$.

Specifically, we construct

$$\bar{\mathbf{Z}}_1 = \bar{\mathbf{Z}}_2 = \begin{bmatrix} 0 & 0 & 1 \\ 0 & 1 & 0 \\ 0 & 1 & 1 \\ 1 & 0 & 0 \\ 1 & 0 & 1 \\ 1 & 1 & 0 \\ 1 & 1 & 1 \end{bmatrix}, \bar{\mathbf{Z}}_3 = \begin{bmatrix} 0 & 1 \\ 1 & 0 \\ 1 & 1 \end{bmatrix}, \qquad \mathbf{P} = \begin{bmatrix} \bar{\mathbf{Z}}_1 & \mathbf{0}_{7\times3} & \mathbf{0}_{7\times2} \\ \mathbf{0}_{7\times3} & \bar{\mathbf{Z}}_2 & \mathbf{0}_{7\times2} \\ \mathbf{0}_{3\times3} & \mathbf{0}_{3\times3} & \bar{\mathbf{Z}}_3 \end{bmatrix}.$$

Any vector $\mathbf{v}$ can be reconstructed by picking at most $k$ vectors out of $\mathbf{P}$, one from each section. For example, let $\mathbf{v} = [0\,1\,0\,0\,1\,1\,1\,0]$. This vector can be reconstructed by adding vectors number 2, 10 and 16 from $\mathbf{P}$.

## 3.5 Constructions for small values of $n$ and $k$

In Section 3.4, we have proved that, independently of the value of $n$, if $k \geq \lceil T/2 \rceil$, then it is sufficient to add one additional transmission to the $T$ transmissions of the conventional index coding scheme. Moreover, the analysis provided in Lemma 3.4.2 showed the order-optimality

of our universal scheme in Theorem 3.4.1 (referred to as Scheme-1) for values of $k < \lceil T/2 \rceil$ when $n$ is large (*i.e.,* exponential in $T$). Figure 3.7 shows the performance of Scheme-1 in Theorem 3.4.1 as a function of the values of $n$ for $T = 20$, with $k = 2$ in Figure 3.7(a) and $k = 5$ in Figure 3.7(b). The performance of Scheme-1 was obtained by averaging over 1000 random index coding instances. In each instance, a code is constructed using the scheme described in Section 3.4.1, and only the rows actually used by the clients $c_{[n]}$ are retained. The performance of the scheme is finally computed by the average number of rows retained in those 1000 iterations. Figure 3.7 shows that our proposed scheme performs well not only for the case of large $n$ (*i.e.,* $n = 2^T - 1$) but also for lower values of $n$. However, Figure 3.7 also suggests that for small values of both $n$ and $k$ (note the left-half of the plot in Figure 3.7(a)), we need to devise schemes that better adapt to the specific values of the index coding matrix $\mathbf{A}$ and vectors $\mathbf{g}_i, i \in [n]$ (recall that Scheme-1 is universal, and hence independent of the value of $\mathbf{A}$). We next propose and analyze the performance of such algorithms.

### 3.5.1 Special Instances

We first represent the problem through a bipartite graph as follows. We assume that the rank of the matrix $\mathbf{G}$ is $T$. Then, there exists a set of $T$ linearly independent vectors in $\mathbf{G}$; without loss of generality, we denote them as $\mathbf{g}_1$ to $\mathbf{g}_T$. Therefore, each vector $\mathbf{g}_{i+T}, i \in [n-T]$, can be expressed as a linear combination of some/all vectors from $\mathbf{g}_{[T]}$; we denote these vectors as the component vectors of $\mathbf{g}_{i+T}$. We can then represent the problem as a bipartite graph $(\mathcal{U} \cup \mathcal{V}, \mathcal{E})$ with $|\mathcal{U}| = T$ and $|\mathcal{V}| = n - T$, where $u_i \in \mathcal{U}$ represents the vector $\mathbf{g}_i$ for $i \in [T]$, $v_j \in \mathcal{V}$ represents the vector $\mathbf{g}_{j+T}$ for $j \in [n-T]$, and an edge exists from node $u_i$ to node $v_j$ if $\mathbf{g}_i$ is one of the component vectors of $\mathbf{g}_{j+T}$. Figure 3.8 shows an example of such graph, where $n = 9$ and $T = 6$. For instance, $v_1$ (*i.e.,* $\mathbf{g}_7$) can be reconstructed by adding $u_i, i \in [4]$ (*i.e.,* $\mathbf{g}_i, i \in [4]$). Given a node $s$ in the graph, we refer to the sets $\mathcal{O}_s$ and $\mathcal{I}_s$ as the *outbound* and *inbound* sets of $s$, respectively: the inbound set contains the nodes which have edges outgoing to node $s$, and the outbound set contains the nodes to which node $s$ has outgoing edges (*i.e.,* the nodes each of which has an incoming edge from $s$). Nodes on either sides of the bipartite graph have either inbound or outbound sets. For instance,

Figure 3.8: Bipartite graph represen-
tation.

Figure 3.9: Optimal representation
when $k = 2$.

with reference to Figure 3.8, $\mathcal{O}_{u_1} = \{v_1, v_2, v_3\}$ and $\mathcal{I}_{v_1} = \{u_1, u_2, u_3, u_4\}$. For this particular example, there exists a scheme with $T_2 = 6$ which can reconstruct any vector with at most $k = 2$ additions. The matrix $\mathbf{A}_2$ which corresponds to this solution consists of the following vectors: $\mathbf{g}_1$, $\mathbf{g}_1 + \mathbf{g}_2$, $\mathbf{g}_1 + \mathbf{g}_2 + \mathbf{g}_3$, $\mathbf{g}_1 + \mathbf{g}_2 + \mathbf{g}_3 + \mathbf{g}_4$, $\mathbf{g}_5$ and $\mathbf{g}_5 + \mathbf{g}_6$. It is not hard to see that each vector in $\mathbf{G}$ can be reconstructed by adding at most 2 vectors in $\mathbf{A}_2$. The vectors in $\mathbf{A}_2$ that are not in $\mathbf{G}$ can be aptly represented as intermediate nodes on the previously described bipartite graph. These intermediate nodes are shown in Figure 3.9 as highlighted nodes. Each added node represents a new vector, which is the sum of the vectors associated to the nodes in its inbound set. We refer to the process of adding these intermediate nodes as creating a *branch*, which is defined next.

**Definition 3.5.1.** Given an ordered set $\mathcal{S} = \{s_1, \cdots, s_S\}$ of nodes, where $s_i$ precedes $s_{i+1}$ for $i \in [S-1]$, a *branch on* $\mathcal{S}$ is a set $\mathcal{S}' = \{s_1', \cdots, s_{S-1}'\}$ of $S-1$ intermediate nodes added to the graph with the following connections: node $s_1'$ has two incoming edges from $s_1$ and $s_2$, and for $i \in [S-1] \setminus \{1\}$, $s_i'$ has two incoming edges from nodes $s_{i-1}'$ and $s_{i+1}$.

For the example in Figure 3.9, we created branches on two ordered sets, $\mathcal{S}_1 = \{u_1, u_2, u_3, u_4\}$ and $\mathcal{S}_2 = \{u_5, u_6\}$. Once the branch is added, we can change the connections of the nodes in $\mathcal{V}$ in accordance to the added vectors. For the example in Figure 3.9, we can replace $u_{[4]}$ in $\mathcal{I}_{v_1}$ with only $s_3$. Using this representation, we have the following lemma.

**Lemma 3.5.1.** *If $\mathcal{O}_{u_{i_T}} \subseteq \mathcal{O}_{u_{i_{T-1}}} \subseteq \cdots \subseteq \mathcal{O}_{u_{i_1}}$ for some permutation $i_1, \cdots, i_T$ of $[T]$, then this instance can be solved by exactly $T$ transmissions for any $k \geq 2$.*

55

**Proof:** One solution of such instance would involve creating a branch on the set $\mathcal{S} = \{u_{i_1}, u_{i_2}, \cdots, u_{i_T}\}$. The scheme used would have the matrix $\mathbf{A}_2$ with its $t$-th row $\mathbf{a}_t = \sum_{\ell=1}^{t} \mathbf{g}_{i_\ell}$ for $t \in [T]$. Note that $\mathbf{g}_{i_1} = \mathbf{a}_1$ and $\mathbf{a}_t + \mathbf{a}_{t-1} = \mathbf{g}_{i_t}$ for all $t \in [T] \setminus \{1\}$. Moreover, for $j \in [n] \setminus [T]$, if $v_{j-T} \in \mathcal{O}_{u_{i_t}}$ for some $i_t$, then $v_{j-T} \in \mathcal{O}_{u_{i_\ell}}$ for all $\ell \leq t$. If we let $t$ be the maximum index for which $v_{j-T} \in \mathcal{O}_{u_{i_t}}$, then we have $\mathcal{I}_{v_{j-T}} = \{u_{i_1}, \cdots, u_{i_t}\}$, and so we get $\mathbf{g}_j = \sum_{\ell=1}^{t} \mathbf{g}_{i_\ell} = \mathbf{a}_t$. This completes the proof. ∎

**Corollary 3.5.2.** *For* $\mathbf{G} \in \mathbb{F}_2^{n \times T}$ *of rank* $T$*, if* $n = T + 1$*, then this instance can be solved in* $T$ *transmissions for any* $k \geq 2$*.*

**Proof:** Without loss of generality, let $\mathbf{g}_{[T]}$ be a set of linearly independent vectors of $\mathbf{G}$. Then, we have $\mathcal{O}_{u_i} = \{v_1\}$ for $i \in \mathcal{I}_{v_1}$ and $\mathcal{O}_{u_j} = \emptyset$ for $j \in [T] \setminus \mathcal{I}_{v_1}$. Thus, from Lemma 3.5.1, this instance can be solved in $T$ transmissions. This completes the proof. ∎

### 3.5.2 Algorithms for General Instances

We here propose two different algorithms, namely Successive Circuit Removing (SCR) and Branch-Search, and analyze their performance.

**Algorithm 1: Successive Circuit Removing (SCR).** Our first proposed algorithm is based on Corollary 3.5.2, which can be interpreted as follows: any matrix $\mathbf{G}$ of $r + 1$ row vectors and rank $r$ can be reconstructed by a corresponding $\mathbf{A}_2$ matrix with $r$ rows. If there does not exist any subset of rows of $\mathbf{G}$ with rank less than $r$, we call $\mathbf{G}$ *a circuit*[4]. Our algorithm works for the case $k = 2^q$, for some integer $q$. We first describe SCR for the case where $q = 1$, and then extend it to general values of $q$. The algorithm works as follows:

1) *Circuit Finding:* find a set of vectors of $\mathbf{G}$ that form a circuit of small size. Denote the size of this circuit as $r + 1$.

2) *Matrix Update:* apply Corollary 3.5.2 to find a set of $r$ vectors that can optimally reconstruct the circuit by adding at most $k = 2$ of them, and add this set to $\mathbf{A}_2$.

3) *Circuit Removing:* update $\mathbf{G}$ by removing the circuit. Repeat the first two steps until the

---

[4]This is in accordance to the definition of a circuit for a matroid[Oxl06].

matrix $\mathbf{G}$ is of size $T' \times T$ and of rank $T'$, where $T' \leq T$. Then, add these vectors to $\mathbf{A}_2$.

Once SCR is executed, the output is a matrix $\mathbf{A}_2$ such that any vector in $\mathbf{G}$ can be reconstructed by adding at most $k = 2$ vectors of $\mathbf{A}_2$. Consider now the case where $q = 2$ (*i.e.*, $k = 4$) for example. In this case, a second application of SCR on the matrix $\mathbf{A}_2$ would yield another matrix, denoted as $\mathbf{A}_4$, such that any row in $\mathbf{A}_2$ can be reconstructed by adding at most 2 vectors of $\mathbf{A}_4$. Therefore, any vector in $\mathbf{G}$ can now be reconstructed by adding at most 4 vectors of $\mathbf{A}_4$. We can therefore extrapolate this idea for a general $q$ by successively applying SCR $q$ times on $\mathbf{G}$ to obtain $\mathbf{A}_k$, with $k = 2^q$.

The following theorem gives a closed form characterization of the best and worst case performance of SCR.

**Theorem 3.5.3.** *Let $T_q^{SCR}$ be the number of vectors in $\mathbf{A}_k$ obtained via SCR. Then, for $k = 2^q$ and integer $q$, we have*

$$\underbrace{f^{Best}(f^{Best}(\cdots f^{Best}(n)))}_{q \text{ times}} \leq T_q^{SCR} \leq \underbrace{f^{Worst}(f^{Worst}(\cdots f^{Worst}(n)))}_{q \text{ times}}, \qquad (3.10)$$

*where $f^{Best}(n) = 2 \left\lfloor \frac{n}{3} \right\rfloor$ and $f^{Worst}(n) = T \left( \left\lfloor \frac{n}{T+1} \right\rfloor + 1 \right)$.*

**Proof:** First we focus on the case $q = 1$. The lower bound in (3.10) corresponds to the best case when the matrix $\mathbf{G}$ can be partitioned into disjoint circuits of size 3. In this case, if SCR finds one such circuit in each iteration, then each circuit is replaced with 2 vectors in $\mathbf{A}_2$ according to Corollary 3.5.2. To obtain the upper bound, note that any collection of $T + 1$ has at most $T$ independent vectors, and therefore contains a circuit of at most size $T + 1$. Therefore, the upper bound corresponds to the case where the matrix $\mathbf{G}$ can be partitioned into circuits of size $T + 1$ and an extra $T$ linearly independent vectors. In that case, the algorithm can go through each of these circuits, adding $T$ vectors to $\mathbf{A}_2$ for each of these circuits, and then add the last $T$ vectors in the last step of the algorithm. Finally, the bounds in (3.10) for a general $q$ can be proven by a successive repetition of the above arguments. ∎

**Algorithm 2: Branch-Search.** A naive approach to determining the optimal matrix $\mathbf{A}_k$ is to consider the whole space $\mathbb{F}_2^T$, loop over all possible subsets of vectors of $\mathbb{F}_2^T$ and, for

every subset, check if it can be used as a matrix $\mathbf{A}_k$. The minimum-size subset which can be used as $\mathbf{A}_k$ is indeed the optimal matrix. However, such algorithm requires in the worst case $O\left(2^{2^T}\right)$ number of operations, which makes it prohibitively slow even for very small values of $T$. Instead, the heuristic that we here propose finds a matrix $\mathbf{A}_k$ more efficiently than the naive search scheme. The main idea behind the heuristic is based on providing a subset $\mathcal{R} \subset \mathbb{F}_2^T$ which is much smaller than $2^T$ and is guaranteed to have at least one solution. The heuristic then searches for a matrix $\mathbf{A}_k$ by looping over all possible subsets of $\mathcal{R}$. Our heuristic therefore consists of two sub-algorithms, namely Branch and Search. Branch takes as input $\mathbf{G}$, and produces as output a set of vectors $\mathcal{R}$ which contains at least one solution $\mathbf{A}_k$. The algorithm works as follows:

1) Find a set of $T$ vectors of $\mathbf{G}$ that are linearly independent. Denote this set as $\mathcal{B}$.
2) Create a bipartite graph representation of $\mathbf{G}$ as discussed in Section 3.5.1, using $\mathcal{B}$ as the independent vectors for $\mathcal{U}$.
3) Pick the dependent node $v_i$ with the highest degree, and split ties arbitrarily. Denote by $\deg(v_i)$ the degree of node $v_i$.
4) Consider the inbound set $\mathcal{I}_{v_i}$, and sort its elements in a descending order according to their degrees. Without loss of generality, assume that this set of ordered independent nodes is $\mathcal{I}_{v_i} = \{u_1,\, u_2,\, \cdots,\, u_{\deg(v_i)}\}$.
5) Create a branch on $\mathcal{I}_{v_i}$. Denote the new branch nodes as $\{u_1^\star,\, u_2^\star,\, \cdots,\, u_{\deg(v_i)}^\star\}$.
6) Update the connections of all dependent nodes in accordance with the constructed branch. This is done as follows: for each node $v_j \in \mathcal{V}$ with $\deg(v_j) \geq k$, if $\mathcal{I}_{v_j} \cap \mathcal{I}_{v_i}$ is of the form $\{u_1,\, u_2,\, \cdots,\, u_\ell\}$ for some $\ell \leq \deg(v_i)$), then replace $\{u_1,\, u_2,\, \cdots,\, u_\ell\}$ in $\mathcal{I}_{v_j}$ with the single node $u_\ell^\star$. Do such replacement for the maximum possible value of $\ell$.
7) Repeat 3) to 6) until all nodes in $\mathcal{V}$ have degree at most $k$.

The output $\mathcal{R}$ is the set of vectors corresponding to all nodes in the graph. The next theorem shows that $\mathcal{R}$ in fact contains one possible $\mathbf{A}_k$, and characterizes the performance of Branch.

**Theorem 3.5.4.** *For a matrix $\mathbf{G}$ of dimension $n \times T$, (a) Branch produces a set $\mathcal{R}$ which*

*contains at least one possible* $\mathbf{A}_k$, *(b) the worst-case time complexity* $t_{Branch}$ *of Branch is* $O(n^2)$, *and (c)* $|\mathcal{R}| \leq (n - T)T$.

**Proof:** To see (a), note that the algorithm terminates when all dependent nodes have a degree of $k$ or less. In every iteration of the algorithm, the degrees of all dependent nodes either remain the same or are reduced. In addition, at least one dependent node is updated and its degree is reduced to 1. Therefore the algorithm is guaranteed to terminate. Since all dependent nodes have degrees $k$ or less, their corresponding vectors can be reconstructed by at most $k$ vectors in $\mathcal{R}$. Therefore, $\mathcal{R}$ contains at least one solution $\mathbf{A}_k$.

To prove (b), the worst-case runtime of Branch corresponds to going over all nodes in $\mathcal{V}$, creating a branch for each one. For the $i$-th node considered by Branch, the algorithm would update the dependencies of all dependent nodes with degrees greater than $k$, which are at most $n - i$ nodes. Therefore $t_{\text{Branch}} = \sum_{i=0}^{n-1}(n - i) = n(n - 1) = O(n^2)$.

To prove (c), note that $|\mathcal{R}|$ is equal to the total number of nodes in all branches created by the algorithm. Therefore we can write $|\mathcal{R}| \leq \sum_{v_i \in \mathcal{V}} \deg(v_i) \leq (n - T)T = O(nT)$. ∎

Let $t_{\text{Search}}$ be the worst-time complexity of the Search step in Branch-Search. Then the worst-case time complexity of Branch-Search is equal to $t_{\text{BS}} = t_{\text{Branch}} + t_{\text{Search}} \leq O(n^2) + 2^{|\mathcal{R}|} = O(n^2) + O(2^{nT}) = O(2^{nT})$, which is exponentially better than the complexity of the naive search. Although our heuristic is still of exponential runtime complexity, we observe from numerical simulations that $|\mathcal{R}|$ is usually much less than $(n - T)T$. Finding more efficient ways of searching through the set $\mathcal{R}$ to find a solution $\mathbf{A}_k$ is an open question.

### 3.5.3   Numerical Evaluation

We here explore the performance of our proposed schemes through numerical evaluations. Specifically, we assess the performance in terms of $T_k$ of SCR and Branch-Search (labeled as BS). We compare their performance against the lower bound in equation (3.6) (labeled as LB), and the upper bound of sending uncoded transmissions (labeled as UB). In particular, we are interested in regimes for which $k < \lceil T/2 \rceil$, because otherwise we know from Theorem 3.4.1 that $T_k = T + 1$. Moreover, we consider values of $n < 2^T - 1$, because

Figure 3.10: Performance comparison for different schemes - $T = 6$, $k = 2$.

if $n = 2^T - 1$ we know from Lemma 3.4.2 that Scheme-1 is order optimal. For SCR, we evaluate its average performance (averaged over 1000 iterations) as well as its upper and lower bounds performance established in Theorem 3.5.3. For Branch-Search, we evaluate its average performance (averaged over 1000 iterations). Figure 3.10 shows the performance of all the aforementioned schemes for $T = 6$ and $k = 2$. As can be seen from Figure 3.10, SCR consistently performs better than uncoded transmissions. In addition, although the current implementation of SCR greedily searches for a small circuit to remove, more sophisticated algorithms for small circuit finding could potentially improve its performance. However, the bounds in (3.10) suggest that the performance of SCR is asymptotically $O(n)$. Branch-Search appears to perform better than other schemes in the average sense. Understanding its asymptotic behavior in the worst-case is an interesting open problem.

## 3.6   Related Work

Index coding was introduced in [BBJ11], where the problem was proven to be NP-hard. Given this, several works have aimed at providing approximate algorithms for the index

coding problem [ELH, BKL10, CS08]. In our work, we were interested in studying the index coding problem from the perspective of private information delivery.

The problem of protecting privacy was initially proposed to enable the disclosure of databases for public access, while maintaining the anonymity of the clients [AP08]. Similar concerns have been raised in the context of *Private Information Retrieval* (PIR), which was introduced in [CKG98] and has received a fair amount of attention [FGH16, CWJ17, SJ18, BU18b, BU18a]. In particular, in PIR the goal is to ensure that no information about the identity of clients' requests is revealed to a set of malicious databases when clients are trying to retrieve information from them. Similarly, the problem of *Oblivious Transfer* was studied [BCR87, MDP14] to establish, by means of cryptographic techniques, two-way private connections between the clients and the server. We note that it is not clear how the use of cryptographic approaches would help in our setup. A curious client, in fact, obtains information about other clients once she learns the transmitted combinations of the messages, *i.e.,* the coding operations. In other words, given that a curious client has also requested data, she needs to learn how the transmitted messages are coded, in order to be able to decode her own requested message.

We were here interested in addressing privacy concerns in broadcast domains. In particular, we analyzed this problem within the index coding framework, as we recently proposed in [KSC17]. This problem differs from secure index coding [DSC12, NPR18], where the goal is to guarantee that an external eavesdropper (with her own side information set) in [DSC12], and each client in [NPR18], does not learn any information about the *content* of the messages other than her requested message. Differently, our goal was to limit the information that a client can learn about the *identities* of the requests of other clients (however, the two approaches could be combined). Note that the techniques developed here can fundamentally differ from those designed for secure index coding. As an extreme example, in fact, the server in our setup can trivially send all the messages that it possesses in an uncoded manner on the broadcast channel. In this case, a curious client will be able to decode all messages, but would still not be able to infer which messages were requested/possessed by other clients, and would learn nothing about their side information. This property is what fun-

damentally contrasts the problem under consideration from the works in [DSC12, NPR18]. Moreover, our approach here has a significant difference with respect to [KSC17]. In fact, while in [KSC17] our goal was to design the coding matrix to guarantee a high-level of privacy, here we assumed that an index coding matrix (that satisfies all clients) was given to us and we developed methods to increase its achieved level of privacy.

The use of $k$-limited-access schemes allows the server to transform an existing index code into a *locally decodable index code* [HL12, NKL18]. Locally decodable index codes allow each client to decode her request using at most $k$ symbols out of the codeword, where $k$ is referred to as the locality of the code. In [HL12], the authors showed that the optimal scalar linear locally decodable index codes with locality 1 are the ones obtained from the coloring of the information graph of the index coding problem. In addition, they provided probabilistic results on the existence (and the impossibility of existence) of locally decodable codes with particular lengths and localities for index coding problems on random graphs. In [NKL18], the authors extended one result in [HL12] where they showed that the optimal *vector* linear locally decodable index codes with locality 1 are obtained from the fractional coloring of the information graph. In addition, they provided a scheme which allows the construction of locally decodable codes for a particular set of index coding instances with special properties, *i.e.*, when certain covering properties are maintained on the side information graph of the index coding problem. Differently from these works, one of the main results of this chapter consisted of providing deterministic constructions/schemes which transform any existing index code into an equivalent code with locality $k$. In addition, our schemes are universal, *i.e.*, they do not depend on the underlying index coding instance.

The solution that we here proposed to limit the privacy leakage is based on finding overcomplete bases. This approach is closely related to compressed sensing and dictionary learning [CN15], where the goal is to learn a dictionary of signals such that other signals can be *sparsely* and *accurately* represented using atoms from this dictionary. These problems seek lossy solutions, *i.e.,* signal reconstruction is not necessarily perfect. This allows a convex optimization formulation of the problem, which can be solved efficiently [RBE10]. In contrast, our problem was concerned with lossless reconstructions, in which case the

optimization problem is no longer convex.

## 3.7 Conclusion

In this chapter, we studied privacy risks in index coding. This problem is motivated by the observation that, since the coding matrix needs to be available to all clients, then some clients may be able to infer the identity of the request and side information of other clients. We proposed the use of $k$-limited-access schemes: these schemes transform the coding matrix so that we can restrict each client to access at most $k$-rows of the transformed matrix as opposed to the whole of it. We explored two privacy metrics, one based on entropy arguments, and the other on the maximal information leakage. Both metrics indicate that the amount of privacy increases with the number of rows that we hide. We then designed polynomial time universal $k$-limited-access schemes, that do not depend on the structure of the index coding matrix $\mathbf{A}$ and proved that they are order-optimal when either $k$ or $n$ is large. For the case where both $k$ and $n$ are small, we proposed algorithms that depend on the structure of the index coding matrix $\mathbf{A}$ and provide improved performance. We overall found that there exists an inherent trade-off between privacy and bandwidth (number of broadcast transmissions), and that in some cases we can achieve significant privacy with minimal overhead.

## 3.8 Appendices

### 3.8.1 Proof of Lemma 3.3.1

The proof is based on simple counting arguments. A subspace $L$ contains all vectors in $L_n$, the number of which is $2^k$. A subspace $L$ therefore consists of a set of $T - k$ linearly independent vectors $\{v_1, \cdots v_{T-k}\}$ that are in $\mathbb{F}_2^m \setminus L_n$, and all linear combinations of $\{v_{[T-k]}\}$ and vectors in $L_n$. We now enumerate the number of ways such a subspace $L$, with $L_n \subseteq L$, can be constructed. We first pick a vector $v_1 \in \mathbb{F}_2^m \setminus L_n$. The total number of possible choices for $v_1$ is equal to $2^m - 2^k$. Once $v_1$ is selected to be in $L$, then all vectors in $v_1 + L_n$ are added to $L$, where $v_1 + L_n$ is the set of vectors obtained by adding $v_1$ to all possible

vectors in $L_n$. Therefore, by picking $v_1$, the total number of vectors of $\mathbb{F}_2^m$ that do not belong to $L$ is now equal to $2^m - 2^{k+1}$, out of which we pick $v_2$. The above process is repeated until all vectors $\{v_{[T-k]}\}$ are selected. Therefore, the total number of such choices becomes $\prod_{\ell=0}^{T-k-1} \left(2^m - 2^{k+\ell}\right)$. In order to compute the total number of subspaces, we need to divide this number by the total number of basis vectors (*i.e.*, linearly independent vectors) used to represent the vectors in $L \setminus L_n$; we denote them by $\{b_1, \cdots, b_{T-k}\}$. The number of vectors in such a basis is $T - k$. Given a subspace $L$, we pick $b_1$ from the set of vectors in $L \setminus L_n$, the number of which is $2^T - 2^k$. Then we pick $b_2$ from the set of vectors $L \setminus (L_n + b_1)$, the number of which is $2^T - 2^{k+1}$. We repeat the previous argument for all $T - k$ vectors. The total number of such basis vectors is therefore equal to $\prod_{\ell=0}^{T-k-1} \left(2^T - 2^{k+\ell}\right)$. Dividing the two quantities therefore proves Lemma 3.3.1.

### 3.8.2  Proof of Theorem 3.3.3

To prove Theorem 3.3.3, we first recall the definition of $\mathcal{G}(q_i, \mathcal{S}_i)$. Given $q_i$ and $\mathcal{S}_i$, $\mathcal{G}(q_i, \mathcal{S}_i)$ is the set which contains all possible $i$-th vectors $\mathbf{g}_i$ of the realization $G$ of the matrix $\mathbf{G}$, namely

$$\mathcal{G}(q_i, \mathcal{S}_i) = \left\{\mathbf{g} \in \mathbb{F}_2^m \mid g_{q_i} = 1, g_{[m] \setminus \{q_i \cup \mathcal{S}_i\}} = 0\right\}.$$

In addition, we define the following set. Given $\mathbf{g}_i$ and an integer $r$, we let $\mathcal{D}(\mathbf{g}_i, r)$ be the set of all possible matrices $\mathbf{A}_k^{(i)}$ of $r$ rows from which $\mathbf{g}_i$ can be reconstructed, namely

$$\mathcal{D}(\mathbf{g}_i, r) = \left\{\mathbf{Z} \in \mathbb{F}_2^{r \times m} \mid \exists \mathbf{d} \in \mathbb{F}_2^r \text{ s.t. } \mathbf{g}_i = \mathbf{d}\mathbf{Z}\right\}.$$

Note that the definition of $\mathcal{D}(\mathbf{g}_i, r)$ is different than that of $\mathcal{P}(\mathbf{g}_i, \mathbf{A}_k, r)$ in that it is not dependent on a specific matrix $\mathbf{A}_k$. Then, we can write

$$P_k^{(\text{MIL})} = \mathcal{L}(A \to A_k^{(n)}|Q_n = q_n, S_n = \mathcal{S}_n) \overset{(a)}{\leq} \log\left|A_k^{(n)}|Q_n = q_n, S_n = \mathcal{S}_n\right|$$

$$\overset{(b)}{=} \log\left|\bigcup_{r=1}^{k} \bigcup_{\mathbf{g}_n \in \mathcal{G}(q_n, \mathcal{S}_n)} \mathcal{D}(\mathbf{g}_n, r)\right|$$

$$\leq \log\left(\sum_{r=1}^{k} \sum_{\mathbf{g}_n \in \mathcal{G}(q_n, \mathcal{S}_n)} |\mathcal{D}(\mathbf{g}_n, r)|\right)$$

$$\overset{(c)}{=} \log\left(2^{|\mathcal{S}_n|} \sum_{r=1}^{k} |\mathcal{D}(\mathbf{g}_n', r)|\right)$$

$$\overset{(d)}{\leq} \log\left(2^{|\mathcal{S}_n|} \sum_{r=1}^{k} \prod_{j=0}^{r-2}(2^m - 2^{j+1})\right)$$

$$\leq \log\left(2^{|\mathcal{S}_n|} k (2^m - 2)^{k-1}\right)$$

$$= O(|\mathcal{S}_n| + mk),$$

where: (i) the equality in (a) follows from Property 2 of the MIL; (ii) the equality in (b) follows by noting that, given $Q_n$ and $S_n$, a possible $A_k^{(n)}$ would belong to $\mathcal{D}(\mathbf{g}_n, r)$ for some $r \in [k]$ and some $\mathbf{g}_n \in \mathcal{G}(Q_n, S_n)$; (iii) the equality in (c) follows by noting that, by symmetry, the number of matrices with $r$ rows from which the vector $\mathbf{g}_i$ can be reconstructed is the same for every possible vector $\mathbf{g}_i \in \mathcal{G}(q_i, \mathcal{S}_i)$. Therefore, the sum over $\mathbf{g}_n$ can be replaced by $\mathcal{D}(\mathbf{g}_n', r) \times |\mathcal{G}(q_n, \mathcal{S}_n)|$ where $\mathbf{g}_n'$ is any arbitrary vector in $\mathcal{G}(q_n, \mathcal{S}_n)$. Based on the structure of the vectors $\mathbf{g}_n \in \mathcal{G}(q_n, \mathcal{S}_n)$, i.e., one in position $q_n$ and zeros in the positions $[m] \setminus \{q_n \cup \mathcal{S}_n\}$, it follows that $|\mathcal{G}(q_n, \mathcal{S}_n)| = 2^{|\mathcal{S}_n|}$; (iv) the inequality in (d) is obtained by counting arguments similar to those in the proof of Lemma 3.3.1. In particular, we enumerate the number of ways we can construct a matrix $\mathbf{A}_k^{(n)}$ with $r$ linearly independent rows, which when linearly combined gives $\mathbf{g}_i$. We first pick a row vector $v_1 \in \mathbb{F}_2^m \setminus \text{Span}(\mathbf{g}_i)$, where $\text{Span}(\mathcal{X})$ of a set of row vectors $\mathcal{X}$ is the row span of these vectors; the number of possible vectors $v_1$ is $2^m - 2$. Then, we pick a second row vector $v_2 \in \mathbb{F}_2^m \setminus \text{Span}(\{\mathbf{g}_i, v_1\})$; the number of possible vectors $v_2$ is $2^m - 2^2$. We repeat this argument for $r - 1$ vectors; the $r$-th vector is then selected so that a linear combination of all $r$ vectors is equal to $\mathbf{g}_i$.

### 3.8.3 Proof of Theorem 3.3.4

We have

$$\bar{P}_k^{(\mathrm{MIL})} = \mathcal{L}(A \to A | Q_n = q_n, S_n = \mathcal{S}_n) \overset{(a)}{=} \log |\{A \; : \; p(A | Q_n = q_n, S_n = \mathcal{S}_n) > 0\}|$$

$$= \log \left| \bigcup_{\mathbf{g} \in \mathcal{G}(q_n, \mathcal{S}_n)} \{A \; : \; \exists \mathbf{d} \in \mathbb{F}_2^T, \mathbf{g} = \mathbf{d}A\} \right|$$

$$\geq \log \left| \{A \; : \; \exists \mathbf{d} \in \mathbb{F}_2^T, \mathbf{g}' = \mathbf{d}A\} \right|$$

$$\overset{(b)}{\geq} \log \left| \{L \subseteq \mathbb{F}_2^m \; : \; \dim(L) = T, \mathbf{g}' \in L\} \right|$$

$$\overset{(c)}{=} \log \prod_{j=1}^{T-1} \left( \frac{2^m - 2^j}{2^T - 2^j} \right)$$

$$\overset{(d)}{\geq} \log \left( \frac{2^m - 2}{2^T - 2} \right)^{T-1} = \Omega \left( mT - T^2 \right),$$

where: (i) the equality in (a) follows from Property 3 of the MIL; (ii) the inequality in (b) follows by letting $L \subseteq \mathbb{F}_2^m$ be a subspace of dimension $\dim(L)$; (iii) the equality in (c) follows by using Lemma 3.3.1 with $k = 1$ (since $\mathbf{g}'$ has only one row) and $t = T$; (iv) the inequality in (d) follows by noting that $\left( \frac{2^m - 2^j}{2^T - 2^j} \right) \geq \left( \frac{2^m - 2}{2^T - 2} \right)$ for $j \in [T-1]$.

### 3.8.4 Proof of Theorem 3.4.1 - Equation (3.6) and Lemma 3.4.2

**Theorem 3.4.1 - Equation** (3.6). Given an index coding matrix $\mathbf{A}$, we denote by $V_{\mathbf{A}} \subseteq \mathbb{F}_2^T$ the subspace formed by the span of the rows of $\mathbf{A}$. It is clear that the dimension of $V_{\mathbf{A}}$ is at most $T$ (exactly $T$ if $\mathbf{A}$ is full rank) and that the $n$ distinct rows of $\mathbf{G}$ lie in $V_{\mathbf{A}}$. Let $\mathbf{a}_i \in \mathbb{F}_2^m, i \in [T_k]$, be the $i$-th row of $\mathbf{A}_k$. Then, the problem of finding a lower bound on the value of $T_k$ can be formulated as follows: *what is a minimum-size set of vectors $\mathcal{A}_k = \{\mathbf{a}_{[T_k]}\}$ such that any row vector of $\mathbf{G}$ can be represented by a linear combination of at most $k$ vectors of $\mathcal{A}_k$?*

A lower bound on $T_k$ can be obtained as follows. Given $\mathcal{A}_k$, there must exist a linear combination of at most $k$ vectors of $\mathcal{A}_k$ that is equal to each of the $n$ distinct row vectors of $\mathbf{G}$. The number of *distinct* non-zero linear combinations of up to $k$ vectors is at most equal

to $\sum_{j=1}^{k} \binom{T_k}{j}$. Thus, we have

$$\sum_{i=1}^{k} \binom{T_k}{i} \geq n. \tag{3.11}$$

Combining this with the fact that $T_k \geq T$ gives precisely the bound in (3.6).

**Lemma 3.4.2.** We now derive the lower bound in Lemma 3.4.2. We first consider the case where $n = 2^T - 1$. From (3.11), we obtain

$$\sum_{i=1}^{k} \binom{T_k}{i} \geq 2^T - 1. \tag{3.12}$$

Since in general $T_k \geq T$, to prove that $T_k \geq T + 1$ for $k < T$, it is sufficient to show that we have a contradiction for $T_k = T$. Indeed, by setting $T_k = T$, the bound in (3.12) becomes

$$\sum_{i=1}^{k} \binom{T}{i} \geq 2^T - 1 = \sum_{i=1}^{T} \binom{T}{i},$$

which clearly is not possible since $k < T$. Hence, $T_k \geq T + 1$ for all $k < T$.

For a general $n$ and $1 \leq k < \lceil T/2 \rceil$, we have

$$k \left( \frac{T_k e}{k} \right)^k \geq k \binom{T_k}{k} \geq \sum_{i=1}^{k} \binom{T_k}{i} \geq n$$

$$\implies T_k \geq \frac{k^{\frac{k-1}{k}}}{e} n^{1/k} = \Omega(kn^{\frac{1}{k}}).$$

Therefore, $T_k = \Omega(k2^{\frac{T}{k}})$ when $n = \Theta(2^T)$. This lower bound, along with the upper bound in equation (3.8) concludes the proof of Lemma 3.4.2.

# CHAPTER 4

# Using mm-Waves for Secret Key Establishment

This chapter shows our second example of application-tailored security solution; namely the use of Millimeter Waves (mmWaves) for ensuring security. The fact that mmWave communication needs to be directional is usually perceived as a challenge; in this chapter we argue that it can be perceived as an opportunity for more resilient security solutions. We are concerned with the problem of efficient secret key sharing among communicating parties. We consider an adversary who passively eavesdrops communication and wishes to learn the secret key being established. Our goal is to establish unconditionally (*e.g.*, regardless of the computational and/or storage capabilities of the adversary) secure keys at a high rate – the reason for these goals will be apparent later in the chapter. Therefore, our security metrics are namely the secret key generation rate and the information-theoretic secrecy of the established keys. We use the directionality in mmWave communication as a resource in our protocol: as will be shown later, our proposed protocol builds on packet erasures, which we show can be induced by the appropriate use of spatial and temporal coding of the transmitted data. We showcase the potential of our approach in two setups: mmWave-based WiFi networks and vehicle platooning. We show that in the first case, we can establish a few hundred secret bits with minimal changes to standard communication protocol; while in both cases, with the right choice of parameters, we can potentially establish keys in the order of tenths of Mbps.

## 4.1 Introduction

Millimeter Wave (mmWave) communications are expected to have significant impact on wireless communication networks such as 5G networks [WHG14], over-60 GHz-WiFi networks [GSC17], autonomous vehicles and vehicle platoons [PSV]. In addition, the inherent directionality of mmWaves can be utilized to establish physical layer secrecy. In fact, inherent properties to the mmWave communication systems (*e.g.*, channel variability, directionality, wider bandwidth allocations) are exploited to reduce the eavesdropping capabilities of Eve [ZWW17].

This work proposes a secret key establishment technique based on mmWave communication. The main motivation stems from the fact that packet erasures can help create secrecy. In fact, a recent line of work [ADD13, SCA16] demonstrated secret key establishment protocols that relied on packet erasures from multi-hop and multipath communication, as well as the use of wiretap codes and beamforming over WiFi. Experiments showed that these protocols yield several Mbps of shared secret keys. This work extends these secret key exchange protocols in a mmWave environment. In fact, the mmWave directional transmissions, if not perfectly aligned, inherently lead to packet losses, and thus it seems a natural host environment for erasure-based key establishment. Differently from existing physical layer secrecy works [HFA18], our proposed protocol: 1) is information-theoretically secure against passive eavesdroppers with limited network presence, 2) allows for very high secret key exchange rates, and 3) is autonomous (*i.e.*, does not require third-party assistance).

Our proposed technique differs from existing cryptographic encryption measures which depend on limited adversarial capabilities: computational capabilities, *e.g.* Diffie-Hellman (DH), or storage capabilities, *e.g.*, Bounded-Storage Model [Dzi06]. In contrast, our proposed scheme establishes secret keys that are information-theoretically secure against eavesdroppers with limited network presence. Moreover, current cryptographic techniques rely on high complexity algorithms to compensate for the low rate of secret key establishment. In this work, we show, through two different scenarios, that our scheme promises secret key generation rates in the order of tens/hundreds of Mbps.

**Main Contributions.** We showcase our approach for two scenarios: (1) over-60 GHz-WiFi networks, where base stations use mmWave antenna arrays for transmissions. First, we propose an analytical model for the instantaneous received Signal-to-Noise-Ratio (SNR), that is inspired from the empirical channel model in [TNM14] and system parameters (*e.g.*, antenna array sizes and beam patterns as described in [GSC17]. We show that, with the right choice of parameters, with minor modification to the standard beamsweeping mechanism, a considerable amount of secret bits (up to a few hundreds) can be established between the base station and mobile devices for virtually no additional transmission cost. In addition, we show that a more invasive secret key establishment protocol achieves few hundred Mbps of secret key generation rates with increased security guarantees.

(2) Vehicular platooning, which is a safety-critical application. We show that, with appropriate choices of code parameters and antenna placement, our technique allows platoons to establish keys with rates up to 166 Mbps – 4 orders of magnitude gain over rates achieved by DH; this allows the use of (otherwise impractical) One-Time Pad (OTP) encryption (an information-theoretically secure encryption technique).

The chapter is organized as follows: Section 4.2 presents our adversary model and background; Section 4.3 discusses the WiFi network application; Section 4.4 discusses the vehicle platooning application and Section 4.5 concludes the chapter.

## 4.2 Model and Background

**System and Adversary Model.** We consider a pair of communicating parties who wish to establish a pairwise key using the scheme in [ADD13, SCA16]. The transmitting, *a.k.a.* Alice (resp. receiving, *a.k.a.* Bob) party is connected to a set of $N_T$ transmitting (resp. $N_R$ receiving) mmWave antenna arrays, each labeled by $t_i$ and situated at location $\mathcal{T}_i, i \in \{1, \cdots, N_T\}$ (resp. $r_j$ and $\mathcal{R}_j, j \in \{1, \cdots, N_R\}$). In addition, the communicating parties wish to communicate secretly in the presence of an eavesdropper (*a.k.a.* Eve), which is equipped with a set of $N_E$ antenna arrays, each label by $e_k$ and situated at location $\mathcal{E}_k, k \in \{1, \cdots, N_E\}$. We assume Eve to be located anywhere within the transmission radius

of the communicating parties, and is passive; therefore the locations of its antennas are unknown. We assume also that Eve has access to the same physical layer technology as the legitimate nodes, has infinite memory as well as unbounded computational capabilities at her disposal, and has perfect knowledge of the protocols. The transmitting power used by each transmitting antenna is denoted by $P_T$, while the noise figure at each receiving antenna array (for both Bob and Eve) is $N_o$. We assume that the available bandwidth is $B$. We assume that each transmitting antenna array is capable of focusing its transmitting energy in desired directions by the use of appropriate beamforming mechanisms. Given a particular beamforming direction, the received SNR at the $j$th receiver (resp. $k$th eavesdropper) antenna from the $i$th transmitter antenna is denoted by $\gamma^{(t_i, r_j)}$ (resp. $\gamma^{(t_i, e_k)}$). Considering the fact that the wireless channel is typically random, then $\gamma^{(t_i, r_j)}$ and $\gamma^{(t_i, e_k)}$ are considered as random variables, with distributions denoted as $f_{\gamma^{(t_i, r_j)}}$ and $f_{\gamma^{(t_i, e_k)}}$.

**mmWave Channel Model and Antenna Patterns.** In mmWaves, transmitters are expected to employ transmit beamforming in order to focus transmission energy in a particular direction in space. However, the radiated energy pattern in space as a result of beamforming strongly relies on 1) the wireless channel between the transmitters and receivers, and 2) the assumed antenna radiation pattern. Therefore, in this work, we strive to employ realistic channel models and antenna patterns in order to give a realistic assessment of our proposed mechanisms. In particular, (1) For over-60 GHz-WiFi cellular networks, we implement the point-to-point 73 GHz outdoor channel model proposed in [TNM14] which takes into account line-of-sight as well as multipath fading signal components. Moreover, in order to take into account the fact that transmitters/receivers that are close by in space exhibit similar channel characteristic, we also implement space consistency between receivers and transmitters, as specified in [Net16]. We also use the standardized antenna radiation pattern proposed in [Net16]. Based on empirical data, we deduce an analytical expression for the received SNR which we describe in the next section.

(2) For vehicular networks, similar models for mmWave channel models are lacking. We developed instead a channel model based on ray tracing, which takes into account reflections off the hood, back and roof of the cars in the platoon. We also used a realistic model for a

vertically-polarized 70 GHz antenna array system.

**Secret Key Protocol [CPF15].** The protocol proposed in [CPF15] allows Alice and Bob (each with possibly multiple antennas) to establish a shared key which is secret from Eve. We here briefly explain how the protocol works, and delegate the details of the protocol to [CPF15, ADD13]. The protocol operates in two rounds of transmission.

<u>Round 1:</u> Alice sends a set of random packets to Bob, who sends feedback of which packets were correctly received. The key idea is that some of these packets are received by Bob while being erased for Eve. If Alice knows a *lower-bound estimate* $N$ of such erased packets, then it can create a shared secret key of size $N$ with Bob in the second round.

<u>Round 2:</u> Assume that $x_1, \cdots, x_M$ are the $M$ packets now shared between Alice and Bob, and Eve knows $M - N$ of them. Then, the secure common key is the concatenation of $y_1, \cdots, y_N$, where the packets $y_i$ are carefully designed (based on MDS codes [SCA16]) linear combinations of $x_1, \cdots, x_M$. Note that Alice does not need to know which $N$ packets are erased; only the number of such packets suffices.

*Example.* In round 1, Alice sends $x_1$ to $x_4$ to Bob, who sends a feedback that packets $x_1$ to $x_3$ were correctly received. Assuming Eve missed two packets $N = 2$, then in round 2, the secret key is the concatenation of $y_1 = x_1 + x_2$ and $y_2 = x_2 + x_3$ which Eve would know nothing about.

Note that this secret is created by knowing that Eve misses at least two packets, but not necessarily knowing *which* two exactly. The security is guaranteed by the fact that the second round does not involve sending the secret itself but rather the packet indices used to create the secrets [SCA16] (*e.g.*, the indices $(1, 2)$ and $(2, 3)$ in the discussed example). In our setting, we make worst-case estimates on how much Eve misses *i.e.*, $N$, and assess how good these estimates are via the *insecure areas* concept as we show next.

**Creating Erasures and Insecure Areas.** The method we follow to enforce erasures in the protocol described previously is based on wiretap codes and directionality. A high-level description of wiretap codes is as follows. The performance of a wiretap code is dictated through two parameters, namely $\text{Th}_1$ and $\text{Th}_2$. Three distinct situations can occur when a

Figure 4.1: Transmission Figure 4.2: Transmission with wiretap codes - Code with wiretap codes - Code Figure 4.3: An illustrative parameters values largely parameters are close in example of the protocol. differ. value.

wiretap-coded packet is received: 1) if it is received with $\gamma \geq \text{Th}_1$ then it is decoded perfectly (*i.e.*, "received"), 2) if $\gamma \leq \text{Th}_2$ then it is completely missed by the receiver (*i.e.*, "erased"), and 3) if $\text{Th}_1 > \gamma > \text{Th}_2$ then partial information can be extracted from the packet. The three aforementioned modes of reception are shown in Figure 4.3. The green area highlights an area in space where a receiver would experience a value of $\gamma \geq \text{Th}_1$ and therefore would decode all transmitted information. The orange region (which typically encloses the green region) highlights an area where $\text{Th}_1 > \gamma > \text{Th}_2$ and therefore a receiver may decode part of the transmitted information. Finally, a receiver outside the green and orange regions (white region) will not be able to infer any information. We assume that packets transmitted in the first round are encoded using wiretap codes. Therefore, Alice hopes that Bob receives these packets while at least some of the packets are erased at Eve; we refer to this event as the *Secret Reception (SecRec) event*. In our setup, we always make the assumption that $N = 1$, *i.e.*, packets from *at least* one transmission/receiver link are erased at Eve while received by Bob. More formally, we denote by a SecRec event that there exists $(t_i, r_j)$ for $(i, j) \in \{1, \cdots, T\} \times \{1, \cdots, R\}$ for which the transmission is correctly received by $r_j$ and is erased by *all* eavesdropper's antennas. The probability of such event is defined as

$$P_{\text{SecRec}} := \Pr\left(\exists (i,j) : \gamma^{(t_i, r_j)} \geq \text{Th}_1, \gamma^{(t_i, e_k)} \leq \text{Th}_2 \; ; \; \forall k\right) \tag{4.1}$$

If $M$ is the number of packets sent on one such link, then the protocol can create shared secret key of size $N = M \cdot L$ where $L$ is the size of one packet. Note that increasing the number of transmitters/receivers (and dispersing them geographically) increase the probability of this estimate of $N$. The protocol fails when the Secure Reception event does not happen, *i.e.*, either Bob does not receive data or Eve receives all packets the Bob does. Consider the example in Figure 4.3 with two transmitting antennas at Alice, and assume that Bob and Eve each has one antenna. If Eve resides in the "not white" region of any link then it would be able to receive the packets transmitted on that link; Secure Reception will not occur with high probability. We therefore consider the protocol to be vulnerable if $P_{\text{SecRec}}$ is not high enough, *i.e.*, $P_{\text{SecRec}} \leq 1 - \delta$ where $\delta$ is our security level. We refer to this situation as the protocol being not $\delta$-secure. The region in space where this occurs (*i.e.*, the probability is not high enough) is referred to as the *Insecure Area (IA)*. Other mechanisms may be needed to protect against eavesdroppers in the IA, and therefore a smaller IA indicates a stronger key agreement mechanism. The choices of $\text{Th}_1$ and $\text{Th}_2$ affect the secret key generation rate as well as the size of the IA. In what follows, we use these two quantities as the performance metrics of our proposed protocol.

**Performance Metrics.** We define the following two performance metrics:

1. *The average secret key rate:* the average number of bits per second established between the communicating parties secretly from the eavesdropper. Given a wiretap code with parameters $\text{Th}_1$ and $\text{Th}_2$, a key generation rate equal to $B\left[\log(1 + \text{Th}_1) - \log(1 + \text{Th}_2)\right]$ can be established between communicating parties while being secure from an eavesdropper, assuming that a secret reception event occurs. Therefore, the average secret key rate is equal to

$$R_{\text{av}} = R_{\max} P_{\text{SecRec}}, \qquad R_{\max} = \left[ \underbrace{B \log_2(1 + \text{Th}_1)}_{\text{Decoding Rate}} - \underbrace{B \log_2(1 + \text{Th}_2)}_{\text{Secrecy Overhead}} \right]. \qquad (4.2)$$

The Decoding Rate component corresponds to the raw data transmission rate achieved between the $(i, j)$-th transmitting/receiving antennas whenever $\gamma^{(t_i, r_j)} \geq \text{Th}_1$. The

Figure 4.4: Antenna sectors [GSC17].

Secrecy Overhead component accounts for the coding overhead due to the use of wiretap codes, thus the difference is the achieved data throughput.

2. *δ-Insecure Area or $IA_\delta$*: we define the set $\mathcal{A} = \{r, \theta : P_{\mathrm{SecRec}} \leq 1 - \delta\}$ as the set of locations in space where the protocol is not δ-secure, where $r, \theta$ are polar coordinates. Therefore, we can define the δ-Insecure Area as

$$IA_\delta := \int_{r,\theta \in \mathcal{A}} r dr d\theta.$$

δ-Insecure Area captures the regions where the likelihood of Eve breaking the secret key establishment mechanism is too high (*i.e.*, at least δ).

Choosing $\mathrm{Th}_1$ and $\mathrm{Th}_2$ gives contradicting effects with respect to the last two objectives. Specifically, when $\mathrm{Th}_1$ and $\mathrm{Th}_2$ are relatively different in value, this results in a relatively larger $P_{\mathrm{SecRec}}$ (therefore a larger insecure area) and larger value of $R_{\max}$. The reverse effect happens when $\mathrm{Th}_1$ and $\mathrm{Th}_2$ are relatively close. We finally note that today a number of practical designs for wiretap codes are emerging, based on polar [MV11], LDPC [TDC07] and lattice codes [LHO], which enable with low complexity to achieve performance curves similar to $R_{\max}$ in (4.2). For this chapter, we will directly use the expression in (4.2) to estimate potential benefits and trade-offs.

## 4.3 Showcase I - IEEE 802.11ay

Our first showcase application is in the context of 60-GHz-based WiFi networks [GSC17]. The IEEE 802.11ay amendment proposes the use of directional communication to cope with the increased signal attenuation that accompanies transmission in the mmWave band.

**Directional Communication.** IEEE 802.11ay proposes the use of *virtual antenna sectors* which discretizes the azimuth angle. Shown in Figure 4.4, a base station sectorizes the azimuth range into $32 - 64$ sectors. Being equipped with up to 3 antenna arrays, each array is responsible from transmission in one-third of these sectors[1]. A mobile device is typically equipped with one antenna array and can have up to 4 sectors. Each device has a set of pre-computed beamforming weights that correspond to transmission in each of the predefined sectors. When a base station wishes to communicate with a mobile device, both communicating parties have to agree on the best sector to use (i.e. best set of beamforming weights to employ) so that received signal strength is maximized. This sector training phase is referred to as *the beamsweeping phase*, and it is split into to sub-phases: 1) a *Sector-Level Sweep (SLS)* phase where both communicating parties agree on the best two sectors to use, and 2) a *Beam Refinement Phase (BRP)* in which the predefined beamforming weights are fine-tuned to further maximize the received signal strength. The SLS phase is also comprised of two sub-phases: the Transmit-SLS for negotiating the best sector to use at the transmitter, and the Receive-SLS for the receiver. We claim that the proposed mechanism for beam training in IEEE 802.11ay creates an excellent opportunity to establish secret keys between mobile devices and WiFi back-end services. For the sake of demonstrating our ideas we only focus on the Transmit-SLS phase, noting that they can be extended to other phases of beamsweeping. We next describe Transmit-SLS:

*1)* The initiator (e.g. base station) sends a sequence of beacon frames, one in each sector. The responder (e.g. mobile device) receives these frames with a quasi-omnidirectional antenna pattern. Each beacon frame is marked with an ID for the used antenna array and sector.

---

[1]Antenna arrays do not cooperatively transmit in the same sector.

*2)* The responder receives the aforementioned frames with varying levels of SNR. It then sends a feedback packet containing the optimal SNR value, and the sector ID of initiator transmitted beacon which was received with this SNR. This feedback packet is transmitted once in every sector of the responder. The initiator receives these frames with a quasi-omnidirectional antenna pattern.

*3)* Upon receiving the feedback packet from the responder, the initiator will be informed of the best sector to use for transmission. The initiator will then send one feedback packet on this sector, containing the optimal SNR value and the ID of the sector used by the responder which was received with this SNR.

*4)* Upon receiving the feedback packet from the initiator, the responder will be informed of the best sector to use for transmission.

**System Parameters.** We assume a WiFi network with $N_T$ base stations (which act as Alice's transmitters) and a mobile device with $N_R$ antennas. Each base station is equipped with mmWave planar antenna arrays with $6 \times 6$ elements, while each mobile device is equipped with a single antenna array. The antenna arrays specifications and radiation patterns follow the standard in [Net16]. As mentioned earlier, we use the channel model proposed in [TNM14] with space consistency as specified in [Net16]. We assume that $P_T/N_oB = -99$ dB and the channel bandwidth is 1 GHz. All transceivers have a noise figure of $-99$ dBm. We assume that base stations have 36 transmission sectors, with the first sector centered at $0°$ with inter-sector separation of $10°$.

### 4.3.1 Analytical Expressions

Empirically, and according to the channel model we use, $\gamma^{(i,j)}$ in dB is normally distributed, i.e., $f_{\gamma^{(i,j)}}(x) \sim \mathcal{N}(\gamma_{\text{av}}^{(i,j)}, \sigma^2)$. The parameters are $\sigma^2 = 24$ and

$$\gamma_{\text{av}}^{(i,j)}(d, \theta) = \gamma_{\text{av}}^{(i,j)}(1, \theta) - 21\log_{10}(d), \tag{4.3}$$

where $d$ and $\theta$ are the distance and the azimuth angle between $i$th transmitter and $j$th receiver ($j$ here refers to an antenna of either Bob or Eve), and where we explicate by

Figure 4.5: $\gamma_{\text{av}}^{(i,j)}$ versus $r$ for $\theta = 0$.



Figure 4.6: $G(\theta)$ versus $\theta$ for $r = 1$.

$\gamma_{\text{av}}^{(i,j)}(d,\theta)$ that $\gamma_{\text{av}}^{(i,j)}$ depends on the factors $d$ and $\theta$. The term $\gamma_{\text{av}}^{(i,j)}(1,\theta)$ corresponds to the average received SNR at a receiver located 1 m away from the $i$th transmitter and with the same azimuth angle as the $j$th receiver. This value is dependent on the transmitted power, receiver noise figure as well as the beampattern of the transmitting antenna array. Empirically, $\gamma_{\text{av}}^{(i,j)}(1,\theta)$ can be approximately modeled as

$$\gamma_{\text{av}}^{(i,j)}(1,\theta) = G(\theta) + \gamma_{\text{init}}^{(i,j)} - 66.8,$$

$$G(\theta) = 10 \log_{10} \left| \left( \sum_{i=1}^{3} 0.33 \, \cos\left(\frac{\theta}{1.8}\right) \cos\left(\frac{2i-1}{2}\pi \sin(\theta)\right) \right)^3 \right|, \quad \theta \in [-180, 180],$$

$$(4.4)$$

where $\gamma_{\text{init}}^{(i,j)} = 10 \log_{10}(\frac{P_T}{N_o B})$ and the subtraction of $-66.8$ dB is to account for the path loss due to 1-m of signal propagation. The empirical expressions discussed here are obtained from Monte-Carlo simulations with 300000 iterations. Figure 4.5 shows $\gamma_{\text{av}}^{(i,j)}$ versus $r$ for $\theta = 0$, both from numerical simulations as well as the expression in (4.3). Figure 4.6 shows $\gamma_{\text{av}}^{(i,j)}$ versus $\theta$ for $r = 1$ both from numerical simulations as well as the expression in (4.4). Figure 4.7 shows the empirical histogram of $\gamma^{(i,j)}$ values in dB, as well as the normal distribution $\mathcal{N}(\gamma_{\text{av}}^{(i,j)}, 24)$ for $r = 0.2$ m, 1 m and 2 m.

Based on the aforementioned assumptions, we can express $P_{\text{SecRec}}$ as

(a) $r = 0.2m$        (b) $r = 0.6m$        (c) $r = 1m$

Figure 4.7: Probability density function for $\gamma_{\text{av}}^{(i,j)}$ - blue line is the empirical distribution and the red line is the analytical expression for the PDF.

$$P_{\text{SecRec}} = \Pr\left(\exists (i,j) : \gamma^{(t_i,r_j)} \geq \text{Th}_1, \gamma^{(t_i,e_k)} \leq \text{Th}_2 \; ; \; \forall k\right) \tag{4.5}$$

$$= 1 - \Pr\left(\forall (i,j) : \overline{\gamma^{(t_i,r_j)} \geq \text{Th}_1, \gamma^{(t_i,e_k)} \leq \text{Th}_2 \; ; \; \forall k}\right)$$

$$\overset{(a)}{=} 1 - \prod_{(i,j)} \Pr\left(\overline{\gamma^{(t_i,r_j)} \geq \text{Th}_1, \gamma^{(t_i,e_k)} \leq \text{Th}_2 \; ; \; \forall k}\right) \tag{4.6}$$

$$= 1 - \prod_{(i,j)} \left(1 - \Pr\left(\gamma^{(t_i,r_j)} \geq \text{Th}_1, \gamma^{(t_i,e_k)} \leq \text{Th}_2 \; ; \; \forall k\right)\right)$$

$$= 1 - \prod_{(i,j)} \left(1 - Q\left(\frac{\text{Th}_1 - \gamma_{\text{av}}^{(t_i,r_j)}}{\sqrt{24}}\right) \prod_{k} \left(1 - Q\left(\frac{\text{Th}_2 - \gamma_{\text{av}}^{(t_i,e_k)}}{\sqrt{24}}\right)\right)\right)$$

$$\tag{4.7}$$

where $(a)$ follows by making an assumption that the variables $\gamma^{(i,j)}, \forall (i,j)$ are independent.

### 4.3.2 Secret Key Establishment Protocols

Incorporating the secret key exchange scheme into the assumed WiFi network can be done in various ways, each with different levels of effectiveness (in terms of the proposed security metrics) as well as its complexity (e.g., how much change in the communication protocol is required to facilitate the scheme). In addition, the performance of the proposed scheme is dictated by the values of $\text{Th}_1$ and $\text{Th}_2$. The thresholds are computed based on the assumed transmission rate and the target secret key rate as follows. Let the beacon frames be transmitted at a rate of $R_T$ Mbps. Setting the Decoding rate in equation (4.2) to this value would

give the corresponding value of $Th_1$. To achieve a target secret key exchange rate of $R_S$, the Secrecy Overhead rate in equation (4.2) should be equal to $R_T - R_S$; this directly gives the corresponding value of $Th_2$.

We now consider two possible secret key establishment protocols:

**1. Beamsweeping-Based (Less Overhead).** In this protocol, the secret key establishment protocol is bootstrapped on top of the T-SLS protocol. Specifically, the protocol includes a chunk of random bits in the beacon frame used by transmitter antennas during the T-SLS phase. Beacon frames are transmitted at a rate of $R_T = 27.5$ Mbps [NCF14]. We assume that a chunk of 1 kbits in the beacon frame is allocated for random bits. In this chunk, we assume that 250 random bits are encoded by a wiretap code and inserted. Therefore we have the Secrecy Overhead to be equal to $750/1k \times 27.5$ Mbps. The corresponding values of $Th_1$ and $Th_2$ can therefore be computed as described earlier. This approach does not require an intrusive change in the existing transmission protocol of the assumed WiFi network. Therefore, as will be shown in the next section, it allows for the establishment of shared secret keys at virtually no additional transmission overhead.

**2. Dedicated Secret Key Exchange Packets (High Secrecy Rate).** In this protocol, dedicated frames are sent for the purpose of secret key establishment. The Decoding Rate as well as the Secrecy Overhead are determined so as the establish a good secret key rate and a small insecure area. These dedicated frames are sent after a legitimate receiver is detected by the transmitter at the T-SLS phase, and therefore an estimate of the receiver's SNR value is known by the transmitter.

### 4.3.3 Performance Evaluation

We assume that a legitimate receiver has $N_R = 1$ one antenna which is located at position $(0,0)$ in space, and receives in an omni-directional way. The transmitter has $N_T$ transmitting antenna array, which are symmetrically distributed around the point $(0,0)$ on a circle with

Figure 4.8: Beamsweeping-Based - $R_{\text{av}}$ for $d = 2$, $N_T = 4$ and $\delta = 0.01$.

Figure 4.9: Beamsweeping-Based - $IA_\delta$ versus $N$ for $d = 2$, $3$, $5$ and $\delta = 0.01$.

radius $d$. We also assume that the eavesdropper is equipped with one antenna which receives in an omni-directional manner (similar to the receiver). When a transmitter antenna array is transmitting data to the legitimate receiver antenna, we assume that the transmitter antenna array beamforms at the location of the receiver's antenna, even during T-SLS phase (*i.e.*, we assume that the legitimate receiver is always aligned with the main direction of a sector). This assumption is reasonable given the locations of transmitters we consider with respect to the legitimate receiver.

**1) Beamsweeping-Based.** Figure 4.8 shows the average number of secret bits when $N_T = 4$ and assuming a Beamsweeping-Based protocol. The value of $R_{\text{av}}$ depends on the probability of Secret Reception, and therefore is dependent on the location of the eavesdropper. Nevertheless, Figure 4.8 shows that an average of 250 bits of secret keys can be achieved against eavesdroppers 5 meters away from the legitimate receiver antenna. Therefore, our simulations suggest T-SLS can automatically provide an average of 250 bits between the transmitter and receiver which are kept secret from eavesdroppers that are 5 meters away from the receiver.

Figure 4.9 shows the insecure area $IA_\delta$ versus $N_T$ for different values of $d$. We consider here $\delta = 0.01$. For $N_T = 4$ and $d = 2$, the insecure region (*i.e.*, the area in which an eavesdroppers renders $P_{\text{SecRec}} \leq 0.99$) is approximately 20 m$^2$. As intuitively expected,

Figure 4.10: $IA_\delta$ versus $Th_2$ for $\delta = 0.01$, 0.1, 0.5, $Th_1 = 7$ dB and $d = 2$.

Figure 4.11: $R_{av}$ for $Th_1 = 7$, $Th_2 = 3$, $\delta = 0.01$ and $d = 2$.

increasing $d$ further increases the insecure area. Moreover, Figure 4.9 suggests that further increasing the number of transmitters does not always decrease the insecure area – the amount of decrease in the insecure area as $N_T$ increases diminishes, with a seeming plateau reached for the case of $d = 2$ at approximately $N_T = 8$.

**2) Dedicated Secret Key Exchange Packets.** The relatively high values of an insecure areas exhibited in Figure 4.9 is due to the particular choices of $Th_1$ and $Th_2$ so as to be compatible with T-SLS frames. On the other hand, if dedicated frames (with specific Decoding Rates and Secrecy Overhead) are to be used, better-performing (in terms of average secret key exchange rate and insecure areas) secret key exchange protocols can be established. Figure 4.10 shows the insecure area versus different choices of $Th_2$ for $d = 2$, $N_T = 4$ and $Th_1 = 7$. The figure suggests that the insecure region can be significantly decreased (approximately to an order of magnitude) by increasing the value of $Th_2$ to 3 dB. In fact, using these particular choices of the thresholds yields a maximum secret key generation rate of 301 Mbps. Figure 4.11 shows how much the average secret key rate is achieved against eavesdroppers in different locations. It is clear that in most of the region, the maximum secret key generation rate is achieved.

## 4.4 Showcase II - Vehicle platooning

Vehicle platooning comprises a set of autonomous cars which drive on the road in a line formation with approximately the same speed and relatively small inter-vehicle distances [PSV].

**Setup and Protocol.** We assume that each car has two mmWave antenna arrays used for transmission and two omni-directional antennas for reception. One pair of transmit/receive antennas (pair-1) is mounted on top of the roof of the car at a height of 0.5 m and the other pair (pair-2) at 1 m. We assume that $N_o = -80$ dBm and $P_T = 30$ dBm.

The secret key protocol works as follows: 1) Pair-1 from the front car sends random packets encoded with a suitable wiretap code to Pair-1 of the back car (Link-1), 2) Pair-2 from the front car sends random packets encoded with a suitable wiretap code to Pair-2 of the back car (Link-2), 3) the front car sends a set of carefully-designed packets as per the protocol in [CPF15] to the back car to establish secret keys.

**Analysis and Discussion.** The preliminary channel model we developed does not account for random channel fading. Therefore, the concept of $\delta$-insecure area becomes a deterministic one, *i.e.*, we only consider $IA_0$. Our key agreement protocol can establish up to 166 Mbps of secret bits, with $IA_0 = 0$. To put this number in perspective, a typical symmetric key exchange algorithm such as (DH-2048), implemented on an off-the-shelf Dedicated-Short-Range-Communication (DSCR) transceiver, gives a key generation rate of approximately 20 kbps; that is, there is a performance gain of approximately 4 orders of magnitude. Table 4.1 shows a comparison between DH-2048 and our proposed secret key establishment protocol.

**Application Example.** We will show next that, thanks to the high rate of secret key generation, our protocol allows for the use of OTP to secure the string stability functionality of vehicle platoons. In order to maintain string stability within the platoon controllers, each car in the platoon exchanges data packets every 100 ms, each of size 60 bytes [PSV], with both the cars in front of and to the back of it. We will show that our suggested key agreement technique can generate enough secret bits which allows the use of OTP to encrypt such messages. Assume that our proposed algorithm for key generation is used every

|  | DH-2048 | Erasure-based mechanism |
|---|---|---|
| Critical resource | Computation power | Bandwidth |
| Secret Key Rate (realistic setup) | 20 kbps | 166 Mbps |
| Complexity of encryption technique | Moderate (AES) | Simple (OTP) |
| Quantum-Vulnerable | Yes | No (Info. theoretically secure) |
| Adversary with high network presence | Resilient | Weak |

Table 4.1: Comparison between DH and proposed mechanism for vehicle platooning.

5 minutes for a duration of 10 ms. Therefore, each two consecutive cars will have an amount of secret keys equal to 10 ms $\times$166 Mbps $\approx$ 200 kB to use for encryption during the next 5 minutes. The total amount of data to be transmitted during the next 5 minutes is equal to 5 min $\times$ 60 B $\approx$ 180 kB $\leq$ 200 kB of secret bits. Therefore, OTP is a practical solution, something rarely achieved in any other kind of security application.

**Discussion.** Comparing our proposed key establishment mechanism that is based on channel erasures, against conventional DH algorithms, we note that our solution is superior in the following aspects: 1) it attains 4 orders of magnitude gains in terms of key generation rates, 2) it allows for using encryption techniques with very low complexity (e.g., OTP) and 3) it is not vulnerable against eavesdroppers with high computational powers (e.g., quantum adversaries). However, it is affected by the availability of a wide transmission bandwidth and the network-presence of adversaries.

## 4.5  Conclusion and Discussion

In this chapter we investigated how the directional nature of mmWave communication can be used to enhance security. We showcased how mmWaves and wiretap codes can enhance

the performance of secret key generation techniques in the context of two applications, over-60 GHz-WiFi networks and vehicle platooning. For both cases, we used/developed channel models with realistic antenna parameters to give realistic assessment of such protocols. For the case of WiFi networks, we empirically developed analytical expressions for the received SNR. We showed that existing T-SLS protocol in IEEE 802.11ay can be used to create a few hundred secret keys at virtually no additional transmission cost, while dedicated protocols in both cases establish very high rates. This work is an initial investigation on the topic. We believe that our results are enticing enough to build a complete system-level implementations of our proposed scheme and analyze its performance in real-world.

# CHAPTER 5

# Distortion based Light-weight Security for Cyber-Physical Systems

(The work in this chapter is based on joint work with Ph.D candidate Gaurav Kumar Agarwal.)

The final example of our application-tailored approach to security is in the context of Cyber-Physical Systems (CPS). In many CPSs, agents can affect the operation of other agents using data that is inter-communicated among them; this data pertains to the control operations of these agents (*e.g.*, state vectors and input vectors). Therefore, an unauthorized and malicious access to this data can hazardously affect the operations of these CPSs. In this work, we consider an adversary who wishes to learn the communicated state vector of a particular control system. Many CPSs are comprised of agents with limited computational and/or energy resources (*e.g.*, IoT devices). In these cases, typical cryptographic tools may not be the most efficient security solutions, and encryption and decryption scheme needs to be devised which are suitable for such control systems. Fortunately, the following observation exists: when designing a security solution, it is (in many cases) sufficient to influence the estimate of the adversary to an estimate that is "far away" from the actual value of the data vector. Based on this observation, we propose a distortion-based security metric which we believe is more appropriate for these applications and is quite frugal in terms of prior requirements on shared keys. In this chapter, we propose distortion-based metrics to protect CPS communication and show that it is possible to confuse adversaries with just a few bits of pre-shared keys. In particular, we will show that a linear dynamical system can communicate its state in a manner that prevents an eavesdropper from accurately learning the state.

## 5.1  Introduction

Wireless networked environments are a natural host for a number of cyber-physical control applications, ranging from autonomous cars and drones, to the Internet-of-Things (IoT), to immersive environments such as augmented reality. It is well recognized that wireless networking is essential to realize the potential of new CPS applications, and is equally well recognized that private and secure exchange of information are necessary and not simply desirable conditions for the CPS ecosystem to thrive. For instance, personal health data in assisted environments, car positions and trajectories, proprietary interests, all need to be protected. This chapter introduces a new approach to secure communication in CPS, that aims to distort an adversary's view of a control system's states. In particular, we will show that a linear dynamical system can securely communicate its state to a trusted party in a manner that prevents a malicious adversary eavesdropping the communication from accurately learning the state.

Our starting observation is that information security measures (cryptographic and information theoretic secrecy), are not well matched to CPS applications as they impose unnecessary requirements, such as protecting all the raw data, and thus can cause high operational costs. Cryptographic methods rely on computational complexity: they require short keys, but high complexity at the communicating nodes (that can be simple sensors in some cases), and can impose a significant overhead on short packet transmissions, therefore increasing delay [WLF16, ZGH13, THM15, KS13]. Information theoretic methods rely on keys: they have low complexity and do not add packet overhead, but require the communicating nodes to share large keys - every communication link needs to use a shared secret key (for a one-time pad) of length equal to the entropy (effectively length) of the transmitted data [Sha49]. These costs accumulate rapidly given that large CPS applications can have dense communication patterns.

Instead, we propose a lightweight approach, that uses small amounts of key and low complexity operations, and builds upon a distortion measure. The following example illustrates

Figure 5.1: Example of drone motion: protection of the most significant bit.

the effect of maximizing distortion[1]. Consider the following simple example of a drone's flying motion, depicted in Fig. 5.1. The drone starts at any position, and moves between adjacent points within the grid. It regularly communicates its location to a legitimate receiver, Bob. A passive eavesdropper, Eve, wishes to infer the drone's locations, and can perfectly overhear all the transmissions the drone makes. We assume the drone and Bob share just one bit of key, that is secret from Eve, and ask: what is the best use we can make of the key?

Using the one bit of shared key to protect the most significant bit (MSB) is not a good solution. The MSB can be protected by XORing a one bit of shared key with the MSB. As shown in Fig. 5.1 the adversary can discover the fake trajectory after a few time steps since this scheme leads to trajectories that do not adhere to the dynamics or environment constraints. In particular, the fake trajectory abruptly moves from the left end of the grid to the right end. At this point, the adversary can learn the real trajectory by flipping back the MSB (we assume that the used scheme is known to everyone). Similar attacks can be made if we use a one-time pad [Sha49] using the same keys over time: as time progresses, more fake trajectories can be discovered and discarded.

Conventional entropy measures also fail to provide insights on how to use the key. For instance, assume we label the 64 squares in Fig. 5.1 sequentially row per row, and consider

---

[1]Although we illustrate our approach for a specific simple example, it extends to protecting general system states.

two cases: in case I, Eve learns that the drone is in one of the neighboring squares $\{1, 2\}$, each with probability 1/2. For case II, Eve knows that the drone is in one of the squares $\{1, 64\}$, again each with probability 1/2. Both cases are equivalent from an information security perspective since in both cases Eve's uncertainty is a set of two equiprobable elements and hence its entropy is 1. However, the security risks in both situations are different. For example, if Eve aims to take a photo of the drone, in the first case she knows where to turn her camera (squares 1 and 2 are close by) while in the second case, she does not (squares 1 and 64 are far apart).

Instead, we propose to use an Euclidean distance distortion measure: how far (in Euclidean space) is Eve's estimate from the actual location. We then propose encoding/decoding schemes which utilize the shared key to maximize this distance. We first consider an "average" distortion measure. Note that if Eve had not received any of the drone transmissions, then the best (adversarial) estimate of the drone's location at any given time is the center point of the confined region in Fig. 5.1. Therefore, a good encryption scheme would strive to maintain Eve's estimate to be as close to the center point as possible; and we achieve the maximum possible distortion, if, after overhearing the drone's transmissions, Eve's best estimate still remains the center point.

The following scheme can achieve this maximum distortion by using exactly one bit of shared secret key. When encoding, the drone either sends its actual trajectory, or a "mirrored" version of it, depending on the value of the secret key. The mirrored trajectory is obtained by reflecting the actual trajectory across a mirroring point in space; in this example, the mirroring point is the center point in Fig. 5.2. Since Eve does not know the value of the shared key, its best estimate of the drone's location - after receiving the drone's transmissions - would be the average location given the trajectory and its mirrored version, which is exactly the center point.

Our results in Section 5.3 extend this idea of mirroring to more general light-weight mappings for dynamical systems in higher dimensional spaces, and theoretically analyze the performance in terms of average distortion for a larger variety of distributions (with certain symmetry conditions). We also discuss a class of systems and controllers for which we can

Figure 5.2: Example of drone motion: mirroring based scheme.

always achieve the perfect distortion with just one bit of key.

The main idea is that many CPS applications can be effectively secured using lightweight approaches which utilize a small amount of keys. To better illustrate this idea, consider a general encryption scheme which uses a $K$-bit key to encrypt the states of a dynamical system. From an abstract point of view, such a scheme hides the true value of the state among a set of $2^K$ states; without knowing the value of the key, an outside observer of the encrypted state cannot resolve the ambiguity among these fake states – we refer to this set as the *ambiguity set*. General encryption (*e.g.*, cryptographic or information-theoretic) schemes aim at increasing the size of the ambiguity set. Differently, in CPS applications, increasing the size of the ambiguity set may not be effective if all of these states are close to each other in a metric space. Distortion based schemes enable to make the most out of a given ambiguity set size.

The main contributions in this chapter are as follows:

• We define security measures that are based on assessing the distortion in the average sense over time and over data.

• We develop a scheme which uses exactly one bit of key and can provide maximum possible distortion (equivalent to Eve with no observations) in some cases. We also discuss the cases where it is not optimal and give an analytical characterization of the attained distortion.

• We then discuss a class of systems and controllers for which we can always guarantee the perfect distortion with just one bit of shared key.

90

- Since for some applications an ambiguity set of size two (corresponding to one bit of key) may not be enough, we also derive an expression of attained distortion when we use larger keys.

## 5.2 System Model

### 5.2.1 System Dynamics

We consider the linear dynamical system,

$$\widetilde{\mathbf{x}}_{t+1} = \mathbf{A}\widetilde{\mathbf{x}}_t + \mathbf{B}\mathbf{u}_t + \mathbf{w}_t, \qquad\qquad \mathbf{y}_t = \mathbf{C}\widetilde{\mathbf{x}}_t + \mathbf{v}_t, \qquad\qquad (5.1)$$

where $\widetilde{\mathbf{x}}_t \in \mathbb{R}^n$ is the state of the system at time $t \in \mathbb{N}$, $\mathbf{u}_t \in \mathbb{R}^m$ is the input to the system at time $t$, $\mathbf{w}_t \in \mathbb{R}^n$ is the process noise, $\mathbf{y}_t$ are the system observations, and $\mathbf{v}_t \in \mathbb{R}^p$ is the observation noise. We denote $\widetilde{\mathbf{x}}_1^T$ by $\widetilde{\mathbf{x}}$, $\mathbf{u}_1^{T-1}$ by $\mathbf{u}$ and $\mathbf{w}_1^{T-1}$ by $\mathbf{w}$. Based on the initial state $\widetilde{\mathbf{x}}_1$ and target state $\widetilde{\mathbf{x}}_T$, the controller computes a sequence of inputs that moves the state from initial state $\widetilde{\mathbf{x}}_1$ to the target state $\widetilde{\mathbf{x}}_T$ in $T$ time instances. We assume that the system uses the obsevations $\mathbf{y}_1^T$ to optimally estimate the states $\widetilde{\mathbf{x}}$. The optimal estimates of $\widetilde{\mathbf{x}}$ made by the system are denoted by $\mathbf{x}$ – in the case of *perfect observation, i.e.,* noiseless and observable systems, then $\mathbf{x} = \widetilde{\mathbf{x}}$.

### 5.2.2 Communication and Adversary Models

At each time instance the system (Alice) transmits information about its state estimate to a legitimate receiver, which is referred to as Bob, via a noiseless link. This situation occurs for example when Bob is remotely monitoring the execution of the system as in Supervisory Control And Data Acquisition (SCADA) systems or in the remote operation of drones.

A malicious receiver, referred to as Eve, is assumed to eavesdrop on the communication between the system and Bob and is able to receive all transmitted signals. The goal of Eve is to make an estimate that is as close to $\mathbf{x}$ as possible: since Bob receives $\mathbf{x}$ and makes control decisions with this information, Eve is interested in $\mathbf{x}$. Eve is assumed to be passive: she

does not actively communicate but is interested in learning the system's states from $t = 1$ to $T$. We assume that the System and Bob have a shared $k$-bit key $K$ which they use to encode/decode the transmitted messages.

### 5.2.3   Inputs and States Random Process Model

We assume that both receivers are only aware of the system model, the matrices $\mathbf{A}, \mathbf{B}, \mathbf{C}$ and the statistics of noises. Therefore, from the perspective of the receivers, the input and output sequences have random distributions which depend on $\mathbf{A}, \mathbf{B}, \mathbf{C}$ and the statistics of the noise. In addition to the process noise $\mathbf{w}$, the joint distribution $f(\mathbf{x}, \mathbf{u}, \mathbf{w})$ depends on *i)* the initial and target states, *ii)* the control law of the system and *iii)* the state estimation process. So, even in noiseless systems, $\mathbf{x}$ and $\mathbf{u}$ possess inherent randomness from a receiver's perspective due to its lack of knowledge about the initial and target states.

### 5.2.4   Encoding Model

The system encodes and transmits packets $\mathbf{z}_1^T$ to ensure that Bob is able to accurately receive $\mathbf{x}_1^T$, the optimal estimates of the system. To do so, the system transmits a packet $\mathbf{z}_t$ at each time step $t$. In this work, we use light-weight memoryless encryption schemes. The $t$-th transmitted packet is a function of only the current state estimate and the shared keys, thus, $\mathbf{z}_t := \mathcal{E}_t(\mathbf{x}_t, K)$, where $\mathcal{E}_t$ is the encoding function used at time $t$. We will denote $\mathbf{z}_1^T$ by $\mathbf{z}$.

### 5.2.5   Bob/Eve Models of Decoding

Bob noiselessly receives the transmitted packets from the system, and decodes them using the shared key. Then, using the decoded information, it generates an estimate of the state of the system at times $t \in [T]$. We require that Bob's estimate is as accurate as Alice's. If we assume that, at time $t \in [T]$, Bob's decoding function is $\Gamma_t(\mathbf{z}_1^t, K)$, then the previous condition is satisfied by ensuring that $\Gamma_t(\mathbf{z}_1^t, K) = \mathbf{x}_t$ for all $t \in [T]$.

Similarly, Eve also receives all transmissions from the system. However, unlike Bob, she

does not have the key $K$. Therefore, Eve's estimate of $\mathbf{x}_t$ is $\hat{\mathbf{x}}_t := \phi_t \left( \mathbf{z}_1^T \right), t \in [T]$, where $\phi_t$ is the decoding function used by Eve at time $t$.

### 5.2.6   Distortion Metrics

We consider a distortion-based security metric which captures how far an estimate is from the actual value. In particular, our analysis is based on the Euclidean distance as our distance metric. However, our analysis can be extended to any $p$-norm, since other norms are just a constant factor away, *i.e.*, $\|\mathbf{x}\|_p \le n^{\frac{1}{p} - \frac{1}{q}} \|\mathbf{x}\|_q$. We assess the performance of Eve as how far its estimate $\hat{\mathbf{x}}$, is from Alice's estimate $\mathbf{x}$. Formally, for a given time instance $t$ and a transmitted codeword $\mathbf{z}_1^T$, we define the following quantity,

$$D(t, \mathbf{z}_1^T) := \mathbb{E}_{\mathbf{x}_t | \mathbf{z}_1^T} \| \mathbf{x}_t - \hat{\mathbf{x}}_t \|^2 \overset{(a)}{=} \mathrm{tr} \left( R_{\mathbf{x}_t | \mathbf{z}_1^T} \right), \tag{5.2}$$

where (5.2) captures the distortion incurred by Eve while estimating $\mathbf{x}_t$ for transmitted symbols $\mathbf{z}_1^T$. Equality in (a) follows because the best (minimizing) estimates of Eve at time $t$ are, $\hat{\mathbf{x}}_t = \phi_t \left( \mathbf{z}_1^T \right) = \mathbb{E} \left[ \mathbf{x}_t | \mathbf{z}_1^T \right]$.

Note that Bob is required to successfully estimate $\mathbf{x}_t$ knowing $\mathbf{z}_1^t$ and the key. Therefore, for a given realization of the key, the encoding function can only map one $\mathbf{x}_t$ and that key realization to each value of $\mathbf{z}_1^T$. Therefore Eve realizes that only trajectories from a particular subset can be the true trajectory for a given $\mathbf{z}_1^T$: those are the ones which correspond to each key realization. Therefore, the expectation in (5.2) is in fact taken over the randomness in the key taking into account posterior probabilities given $\mathbf{z}_1^T$. If Eve does not have observations, the expectation is taken over $X_t$ with prior distribution and we get $D(t, \mathbf{z}_1^T) = \mathrm{tr}(R_{\mathbf{x}_t})$.

As $D(t, \mathbf{z}_1^T)$ is a function of time $t$ and the transmitted sequence $\mathbf{z}_1^T$, we consider an "average case" distortion (denoted by $D_E$) where we take expectation over all possible $\mathbf{z}_1^T$ and average out over time[2].

$$\text{Average} \quad - \quad D_E := \mathbb{E}_{\mathbf{z}_1^T} \left[ \frac{1}{T} \sum_{t=1}^{T} D(t, \mathbf{z}_1^T) \right] \tag{5.3}$$
$$\text{Distortion}$$

---

[2] Another notion of "worst-case" distortion is considered in [AKD18] which is not included in this thesis.

It is worth to note that the definitions of $D_E$ in (5.3) implies that Eve's state estimation must be associated to a time instance. In other words, making a random/constant estimate of the state hoping that it matches the actual state at some time will lead to high distortion values.

### 5.2.7 Design Goals

Our goal is to choose the encoding and decoding functions, $\mathcal{E}_t$ and $\phi_t$, so that Bob can decode loselessly while the distortion is maximized for Eve's estimate. In addition, we seek to achieve this with the minimum amount of shared keys $K$. In absence of any observation by Eve, the distortion will be,

$$D_E^{\max} = \frac{1}{T} \sum_{t=1}^{T} \mathrm{tr}(R_{\mathbf{x}_t}),$$

This will serve as upper bounds as,

$$D_E = \frac{1}{T} \mathbb{E}_{\mathbf{z}_1^T} \sum_{t=1}^{T} \mathrm{tr}(R_{\mathbf{x}_t | \mathbf{z}_1^T}) \overset{(a)}{\leq} \frac{1}{T} \sum_{t=1}^{T} \mathrm{tr}(R_{\mathbf{x}_t}) = D_E^{\max}, \tag{5.4}$$

where (a) follows by noting that the trace of the conditional covariance matrix is a quadratic (convex) function in $\mathbf{z}_1^T$ and therefore we can use Jensen's inequality.

## 5.3  Optimizing Average Distortion $D_E$

In this section, we will first discuss schemes to optimize the Average Distortion ($D_E$). We will initially analyze encoding schemes which use *one* bit of secret key, and characterize their attained level of distortion. We then show that such schemes attain the maximum level of distortion for a family of distributions on $\mathbf{x}$ which exhibit a certain class of symmetry. Later we describe how this analysis extends to the use of multiple keys, as for some application having an ambiguity set of size two might not be enough.

### 5.3.1 Encoding Schemes with 1-bit Shared Secret Key

We now discuss encoding schemes that use one bit of shared key and show how the achieved distortion compares to the upper bound in (5.4). These encoding schemes work as follows:

$$
\mathbf{z}_t = 
\begin{cases}
\mathbf{x}_t & \text{if } K = 0, \\
\alpha_t(\mathbf{x}_t) & \text{if } K = 1,
\end{cases}
\quad \forall t \in [T],
\tag{5.5}
$$

where $K \in \{0, 1\}$ is the shared bit and $\alpha_t(X_t)$ is a transformation of the state vector $\mathbf{x}_t$. We denote by $\alpha_t^{-1}(\mathbf{x}_t)$ the inverse transformation of $\alpha_t$. We will next show the attained distortion of such schemes.

**Theorem 5.3.1** (Proof in Appendix 5.4.1). *The average distortion ($D_E$) attained by using the scheme in (5.5) is,*

$$
\frac{1}{2T} \sum_{t=1}^{T} \mathbb{E}_{\mathbf{x}} \left\{ \frac{f_{\mathbf{x}}(\alpha^{-1}(\mathbf{x}))}{f_{\mathbf{x}}(\mathbf{x}) + f_{\mathbf{x}}(\alpha^{-1}(\mathbf{x}))} \left\| \mathbf{x}_t - \alpha_t^{-1}(\mathbf{x}_t) \right\|^2 \right\},
\tag{5.6}
$$

*where $\alpha^{-1}(\mathbf{x}) := [\alpha_1(\mathbf{x}_1)' \, \alpha_2(\mathbf{x}_2)' \, \cdots \, \alpha_T(\mathbf{x}_T)']'$. Moreover, if the following condition holds,*

$$
f_{\mathbf{x}}(\mathbf{x}) = f_{\mathbf{x}}(\alpha^{-1}(\mathbf{x})), \qquad \qquad \text{for all } \mathbf{x} \in \mathcal{X},
\tag{5.7}
$$

*then the expression simplifies to*

$$
D_E = \frac{1}{4T} \sum_{t=1}^{T} \mathbb{E}_{\mathbf{x}} \left\| \mathbf{x}_t - \alpha_t(\mathbf{x}_t) \right\|^2.
\tag{5.8}
$$

Condition (5.7) implies a general notion of symmetry in the distribution of $f_{\mathbf{x}}(\mathbf{x})$. In the following, we focus on a particular notion of distribution symmetry, for which we show the corresponding choice of $\alpha_t(\mathbf{x}_t)$ and how it can achieve high levels of distortion. Consider a transformation function $\alpha_t(\mathbf{x})$ which reflects a point $\mathbf{x}$ across an affine subspace of dimension $d$, defined by the equations $\mathbf{S}_t \mathbf{x} = \mathbf{b}_t$ where $\mathbf{S}_t \in \mathbb{R}^{d \times n}$ consists of $d \leq n$ orthonormal rows, and $\mathbf{b}_t \in \mathbb{R}^d$; the transformation is $\alpha_t(\mathbf{x}) = (\mathbf{I} - 2\mathbf{S}_t'\mathbf{S}_t)\,\mathbf{x} + 2\mathbf{S}_t'\mathbf{b}_t$. The choice of the dimension $d$ and the subspace $(\mathbf{S}_t, \mathbf{b}_t)$ depend on the properties we would like the encoded trajectories to have. We refer to encoding schemes that are based on this transformation as *mirroring*

Figure 5.3: Mirroring across the line passing through the origin and having a 45° angle with the $X$-axis.

*schemes.* For example, consider $\mathbf{x}_t \in \mathbb{R}^2$ where $\mathbf{S}_t = \frac{1}{\sqrt{2}}[-1 \ 1]$ and $\mathbf{b}_t = 0$. Then $\alpha_t(\mathbf{x}_t)$ corresponds to mirroring across a line that passes through the origin with a 45° angle. This is shown in Fig. 5.3. We are interested in *mirroring schemes* as they are light-weight and can be implemented on low-complexity IoT devices. Moreover, such schemes can provide the maximum distortion level for a class of distributions with what we refer to as *Point Symmetry.*

**Definition 5.3.1 (Point Symmetry).** A random vector $\mathbf{x}$ is said to have Point Symmetry if there exists a point $\mathbf{v}$ for which $f_{\mathbf{x}}(\mathbf{x}) = f_{\mathbf{x}}(2\mathbf{v} - \mathbf{x})$, $\forall \mathbf{x} \in \mathcal{X}$.

**Lemma 5.3.2.** *If $\mathbf{x}$ has Point Symmetry across $\mathbf{v}$, then $\mathbf{v} = \mu_{\mathbf{x}}$.*

 **Proof:** Since $\mathbf{x}$ has Point Symmetry, then

$$f_{\mathbf{x}}(\mathbf{x}) = f_{\mathbf{x}}(2\mathbf{v} - \mathbf{x}) \qquad \Rightarrow \qquad f_{\mathbf{x}}(\mathbf{x}) = f_{2\mathbf{v} - \mathbf{x}}(\mathbf{x})$$

$$\Rightarrow \qquad \mu_{\mathbf{x}} = 2\mathbf{v} - \mu_{\mathbf{x}} \qquad \Rightarrow \qquad \mu_{\mathbf{x}} = \mathbf{v}.$$

$\blacksquare$

 The following result characterizes the performance of the mirroring scheme, and shows that it achieves the maximum distortion for distributions with Point Symmetry.

**Corollary 5.3.3.** *If $\alpha_t(\mathbf{x}_t)$ is based on a **mirroring scheme** along the planes given by $\mathbf{S}_t \mathbf{x} = \mathbf{b}_t$, $t \in [T]$ and the condition (5.7) holds, then (5.8) becomes,*

$$D_E = \frac{1}{T} \sum_{i=1}^{T} tr\left(\mathbf{S}_t R_{\mathbf{x}_t} \mathbf{S}_t' + (\mathbf{b}_t - \mathbf{S}_t \mu_{\mathbf{x}_t})(\mathbf{b}_t - \mathbf{S}_t \mu_{\mathbf{x}_t})'\right). \tag{5.9}$$

*Moreover, if* $\mathbf{x}$ *has Point Symmetry, then* $D_E = \frac{1}{T}\sum_{t=1}^{T} tr(R_{\mathbf{x}_t})$, *the maximum possible distortion.*

**Proof:** If condition (5.7) holds, then by simply plugging the expression of $\alpha_t(\mathbf{x}_t)$ for the mirroring scheme along $\mathbf{S}_t\mathbf{x} = \mathbf{b}_t$ that is $\alpha_t(\mathbf{x}_t) = (\mathbf{I} - 2\mathbf{S}'_t\mathbf{S}_t)\,\mathbf{x}_t + 2\mathbf{S}'_t b_t$ in (5.8) we get (5.9) (Formal proof in Appendix 5.4.1). Choosing $\mathbf{S}_t = I$ and $\mathbf{b}_t = \mu_{\mathbf{x}_t}$ makes $\alpha^{-1}(\mathbf{x}_1^T) = 2\mu_{\mathbf{x}_1^T} - \mathbf{x}_1^T$ which by Point Symmetry satisfies (5.7). Therefore, we get $D^E = \frac{1}{T}\sum_{t=1}^{T} tr(R_{\mathbf{x}_t})$. ∎

Now, we show the implications of our results for mirroring based schemes in the context of a few examples.

**Example 1.** Consider an example where $\mathbf{u}$ is distributed as Gaussian with mean $\mu_{\mathbf{u}}$ and covariance matrix $R_{\mathbf{u}}$. Then for a noiseless system with perfect observation and a zero initial state, $\mathbf{x}_2^T$ is also Gaussian distributed with mean $\mu_{\mathbf{x}_2^T} = \mathbf{Q}\mu_{\mathbf{u}}$ and variance $R_{\mathbf{x}_2^T} = \mathbf{Q}R_{\mathbf{U}}\mathbf{Q}^T$, where

$$
\mathbf{Q} = \begin{bmatrix} \mathbf{B} & 0 & \cdots & 0 \\ \mathbf{AB} & \mathbf{B} & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{A}^{T-2}\mathbf{B} & \cdots & \mathbf{AB} & \mathbf{B} \end{bmatrix}.
$$

A Gaussian random vector has Point Symmetry and therefore, according to Corollary 5.3.3, we can get maximum distortion by setting $\mathbf{b}_t = \mu_{\mathbf{x}_t}$ and $\mathbf{S}_t = \mathbf{I}$.

The next example is based on a Markov-based model for the dynamical system. For this example, the following lemma is useful.

**Lemma 5.3.4.** *Consider the random vector* $\mathbf{x}_1^T$ *where the following conditions hold: 1)* $f_{\mathbf{x}_1}(\mathbf{x}_1)$ *has Point Symmetry, and 2)* $f_{\mathbf{x}_t|\mathbf{x}_1^{t-1}}(\mathbf{x}_t|\mathbf{x}_1^{t-1})$ *has Point Symmetry, then so does* $f_{\mathbf{x}}(\mathbf{x})$, *where* $\mathbf{x} = \mathbf{x}_1^T$ *and* $\mu = [\mu_{\mathbf{x}_1}'\ \mu_{\mathbf{x}_2}'\ \cdots\ \mu_{\mathbf{x}_T}']'$. *Therefore, by virtue of Corollary 5.3.3, mirroring schemes can achieve the maximum distortion.*

Lemma 5.3.4 allows us to characterize the performance of the following example.

97

**Example 2.** Consider the following random walk mobility model. Let $a \in \mathbb{N}^+$, and $\mathbf{x}_t$ be its location at time $t$, then,

$$\mathbf{x}_1 \sim \text{Uni}([-a : a])$$

$$\mathbf{x}_t | \mathbf{x}_{t-1} \sim \text{Uni}([-a : a] \cap \{\mathbf{x}_{t-1} - 1, \mathbf{x}_{t-1}, \mathbf{x}_{t-1} + 1\}).$$

One can see that these distributions satisfy the conditions in Lemma 5.3.4. Therefore, one can set $\mathbf{b}_t = \mu_t = 0$ and $\mathbf{S}_t = 1$, which will achieve maximum distortion of $D_E$.

**Example 3.** Here we provide a numerical example which shows how our mirroring scheme performs for situations where we compute the state distributions using numerical simulations. We consider the quadrotor dynamical system provided in (4) of [KM12]. The quadrotor moves in a 3-dimensional cubed space with a width, length and height of 2 meters, where the origin is the center point of the space. The quadrotor starts its trajectory from an initial point $(-1, y_1, z_1)$ and finishes its trajectory at a target point $(1, y_T, z_T)$ after $T$ time steps, where the points $y_1, z_1, y_T, z_T$ are picked uniformly at random in $[-1, 1]^4$. We assume that $T = 10$ time steps, and that the continuous model in [KM12, (4)] is discretized with a sample time of $T_s = 0.5$ seconds. We assume that the quadrotor encodes and transmits only the states which contain the location information (first three elements of the state vector $\mathbf{x}_t$). The quadrotor is equipped with an LQR controller which designs the input sequence $\mathbf{u}_1^{T-1}$ as the solution of the following problem

$$
\begin{aligned}
\text{minimize} \quad & \|\mathbf{u}\|^2 + 10 \left\|\mathbf{x}_2^{T-1}\right\|^2 \\
\text{subject to} \quad & \mathbf{x}_{t+1} = \mathbf{A}^{\text{quad}}\mathbf{x}_t + \mathbf{B}^{\text{quad}}\mathbf{u}_t, \quad \forall t \in [T-1] \\
& \mathbf{x}_1 = \begin{bmatrix} -1 & y_1 & z_1 & 0 & \cdots & 0 \end{bmatrix}', \\
& \mathbf{x}_T = \begin{bmatrix} 1 & y_T & z_T & 0 & \cdots & 0 \end{bmatrix}',
\end{aligned}
\tag{5.10}
$$

where $\mathbf{A}^{\text{quad}}$ and $\mathbf{B}^{\text{quad}}$ define the quadrotor's discrete-time model. The remaining states of $\mathbf{x}_1$ and $\mathbf{x}_T$ are set to zero to allow the drone to hover at the respective locations. We perform numerical simulation of the aforementioned setup: we run 2 millions iterations, where in each iteration a new initial and target points are picked, and the resultant trajectory is recorded.

Figure 5.4: An illustration of some trajectories. The reflection plane is shown as a dashed-black line. One trajectory (solid-black) is shown along with its mirrored image (dotted-black).

Based on the recorded data, we consider different mirroring schemes and numerically evaluate the attained distortion. To facilitate numerical evaluations, the simulation space is gridded into bins with 0.2 meters of separation, and the location of the drone at each trajectory is approximated to the nearest space bin.

Figure 5.4 shows some of the drone trajectories obtained from our numerical simulation. It is clear that not all trajectories are equiprobable, and therefore the distribution of $\mathbf{x}_t$ is not uniform across all bins in space. Since the motion of the drone is mainly progressive in the positive x-axis direction, reflection across a fixed point results in mirrored trajectories that are progressing in the opposite direction, and therefore are identified to be fake automatically. Therefore, mirroring across a point here is useless: the numerically computed distortion for this scheme is equal to zero.

Next we consider mirroring across the reflection plane shown in Figure 5.4, where $\mathbf{b}_t = 0$ and $\mathbf{S}_t = \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$. As can be seen from the figure, the reflection plane is indeed an axis of symmetry for the distribution of the drones trajectories, and therefore is expected to provide high distortion values. We numerically evaluate the attained distortion using the scheme by using equation (5.6), which evaluates to $D_E = 0.3971$. This is slightly less than

$D_E^{\max} = 0.3979$.

### 5.3.2  Encoding Schemes with $k$-bits Shared Secret Key

The scheme in (5.5) assumes the use of one bit for encryption. However, it is straightforward to extend the scheme when we require a larger ambiguity set. For $k$ bits, we denote the possible values of the shared key as $K \in [0 : 2^k - 1]$. Therefore, the scheme works as follows

$$\mathbf{z}_t(K) = \alpha_t^{(K)}(\mathbf{x}_t), \ \forall t \in [T], \tag{5.11}$$

where $\alpha_t^{(K)}$ is an invertible transformation function used at time $t$ when the value of the key is $K$, and $\alpha_t^{(0)}(x) = \alpha_t^{-(0)}(\mathbf{x}) = \mathbf{x}$. The following theorem shows the achieved value of the distortion in this case, which is a direct extension of Theorem 5.3.1.

**Theorem 5.3.5** (Proof in Appendix 5.4.2). *The average distortion $D_E$ attained by using the scheme in* (5.11) *is*

$$\frac{1}{2^k T} \sum_{t=1}^{T} \mathbb{E}_{\mathbf{x}} \left\{ \frac{\sum_{K=0}^{2^k-1} f_{\mathbf{x}}(\alpha^{-(K)}(\mathbf{x})) \left\| R_t^{(K)} \right\|^2}{\left[ \sum_{K=0}^{2^k-1} f_{\mathbf{x}}(\alpha^{-(K)}(\mathbf{x})) \right]^2 f_{\mathbf{x}}(\mathbf{x})} \right\}, \tag{5.12}$$

*where*

$$R_t^{(K)} = \sum_{\ell=0}^{2^k-1} f_{\mathbf{x}}(\alpha^{-(\ell)}(\mathbf{x})) \left( \alpha_t^{-(\ell)}(\mathbf{x}_t) - \alpha_t^{-(K)}(\mathbf{x}_t) \right),$$

$$\alpha^{-(K)}(\mathbf{x}) := [\alpha_1^{-(K)}(\mathbf{x}_1)' \, \alpha_2^{-(K)}(\mathbf{x}_2)' \, \cdots \, \alpha_T^{-(K)}(\mathbf{x}_T)']'.$$

*Moreover, if the following condition holds,*

$$f_{\mathbf{x}}(\alpha^{-(K)}(\mathbf{x})) = C(\mathbf{x}), \qquad\qquad \forall x \in \mathcal{X}, \forall K \in [0 : 2^k - 1], \tag{5.13}$$

*where $C(\mathbf{x})$ is a constant, then the expression on $D_E$ simplifies to*

$$\frac{1}{2^{3k}T} \sum_{t=1}^{T} \mathbb{E}_{\mathbf{x}} \left[ \sum_{K=0}^{2^k-1} \left\| \sum_{\ell=0}^{2^k-1} \left( \alpha_t^{-(K)}(\mathbf{x}_t) - \alpha_t^{-(\ell)}(\mathbf{x}_t) \right) \right\|^2 \right]. \tag{5.14}$$

100

Using multiple bits of shared keys can provide benefits beyond having a larger ambiguity set. In fact, while we show the optimality of 1-bit mirroring schemes for distributions with point symmetries, using multiple bits of shared key can provide a better distortion for general distributions. For example, it was shown that, for a general finite alphabet: (1) 1-bit schemes are not sufficient to achieve the maximum distortion, and (2) with just 5 bits of shared keys, a scheme achieves more than 97% of the maximum possible distortion [TAF17].

## 5.4 Appendices

### 5.4.1 Proof of Theorem 5.3.1 and Corollary 5.3.3

We start by computing $R_{\mathbf{x}_t | \mathbf{z}_1^T}$. Note that given a sequence of transmitted symbol $\mathbf{z}_1^T$ there are two possible values of sequence of message symbols $\mathbf{x}_1^T$ which are $\mathbf{x}_1^T = \mathbf{z}_1^T$ and $\mathbf{x}_1^T = \tilde{\mathbf{z}}_1^T$, where $\tilde{\mathbf{z}}_t$ is $\alpha_t^{-1}(\mathbf{z}_t)$.

The posterior probability of $\mathbf{x}_t = \mathbf{z}_t$ given $\mathbf{z}_1^T$, *i.e.*, $Pr(\mathbf{x}_t = \mathbf{z}_t | \mathbf{z}_1^T)$ will be equal to $Pr(\mathbf{x}_1^T = \mathbf{z}_1^T | \mathbf{z}_1^T) := p_{\mathbf{z}}$. We note that $p_{\mathbf{z}} = \frac{f(\mathbf{z})}{f(\mathbf{z}) + f(\tilde{\mathbf{z}})}$, where $\tilde{\mathbf{z}} := [\tilde{\mathbf{z}}_1' \ \tilde{\mathbf{z}}_2' \ \cdots \ \tilde{\mathbf{z}}_T']'$. Then, $\mathbb{E}(\mathbf{x}_t | \mathbf{z}_1^T) = p_{\mathbf{z}} \mathbf{z}_t + (1 - p_{\mathbf{z}})(\tilde{\mathbf{z}}_t)$. With this,

$$
\begin{aligned}
R_{\mathbf{x}_t | \mathbf{z}_1^T} &= \mathbb{E}_{\mathbf{x}_t | \mathbf{z}_1^T} \left[ \left( \mathbf{x}_t - \mathbb{E}(\mathbf{x}_t | \mathbf{z}_1^T) \right) \left( \mathbf{x}_t - \mathbb{E}(\mathbf{x}_t | \mathbf{z}_1^T) \right)' \right] \\
&= p_{\mathbf{z}}(1 - p_{\mathbf{z}})^2 (\mathbf{z}_t - \tilde{\mathbf{z}}_t)(\mathbf{z}_t - \tilde{\mathbf{z}}_t)' + (1 - p_{\mathbf{z}}) p_{\mathbf{z}}^2 (\mathbf{z}_t - \tilde{\mathbf{z}}_t)(\mathbf{z}_t - \tilde{\mathbf{z}}_t)' \\
&= p_{\mathbf{z}}(1 - p_{\mathbf{z}})(\mathbf{z}_t - \tilde{\mathbf{z}}_t)(\mathbf{z}_t - \tilde{\mathbf{z}}_t)' \\
D_E &= \frac{1}{T} \mathbb{E}_{\mathbf{z}} \sum_{t=1}^{T} \operatorname{tr} \left( R_{\mathbf{x}_t | \mathbf{z}_1^T} \right) = \frac{1}{T} \mathbb{E}_{\mathbf{z}} \sum_{t=1}^{T} \operatorname{tr} \left( p_{\mathbf{z}}(1 - p_{\mathbf{z}})(\mathbf{z}_t - \tilde{\mathbf{z}}_t)(\mathbf{z}_t - \tilde{\mathbf{z}}_t)' \right) \\
&= \frac{1}{T} \mathbb{E}_{\mathbf{z}} \sum_{t=1}^{T} p_{\mathbf{z}}(1 - p_{\mathbf{z}}) \operatorname{tr} \left( (\mathbf{z}_t - \tilde{\mathbf{z}}_t)(\mathbf{z}_t - \tilde{\mathbf{z}}_t)' \right) \\
&= \frac{1}{T} \mathbb{E}_{\mathbf{z}} \sum_{t=1}^{T} p_{\mathbf{z}}(1 - p_{\mathbf{z}}) \| \mathbf{z}_t - \tilde{\mathbf{z}}_t \|^2 = \frac{1}{T} \mathbb{E}_{\mathbf{z}} \sum_{t=1}^{T} \frac{f_{\mathbf{x}}(\mathbf{z}) f_{\mathbf{x}}(\tilde{\mathbf{z}})}{(f_{\mathbf{x}}(\mathbf{z}) + f_{\mathbf{x}}(\tilde{\mathbf{z}}))^2} \| \mathbf{z}_t - \tilde{\mathbf{z}}_t \|^2.
\end{aligned}
$$

Now, $\mathbf{z}_1^T$ is the transmitted symbols if $\mathbf{x}_1^T = \mathbf{z}_1^T$ and key was zero or if $\{ \mathbf{x}_t = \tilde{\mathbf{z}}_t, \ \forall t \in [T] \}$

and key was one. So $f_{\mathbf{z}}(\mathbf{z}) = \frac{f_{\mathbf{x}}(\mathbf{z}) + f_{\mathbf{x}}(\tilde{\mathbf{z}})}{2}$. Thus $D_E$,

$$
= \frac{1}{T} \mathbb{E}_{\mathbf{z}} \sum_{t=1}^{T} \frac{f_{\mathbf{x}}(\mathbf{z}) f_{\mathbf{x}}(\tilde{\mathbf{z}})}{(f_{\mathbf{x}}(\mathbf{z}) + f_{\mathbf{x}}(\tilde{\mathbf{z}}))^2} \|\mathbf{z}_t - \tilde{\mathbf{z}}_t\|^2 = \frac{1}{T} \int f_{\mathbf{z}}(\mathbf{z}) \sum_{t=1}^{T} \frac{f_{\mathbf{x}}(\mathbf{z}) f_{\mathbf{x}}(\tilde{\mathbf{z}})}{(f_{\mathbf{x}}(\mathbf{z}) + f_{\mathbf{x}}(\tilde{\mathbf{z}}))^2} \|\mathbf{z}_t - \tilde{\mathbf{z}}_t\|^2 d\mathbf{z}
$$

$$
= \frac{1}{2T} \int \sum_{t=1}^{T} \frac{f_{\mathbf{x}}(\mathbf{z}) f_{\mathbf{x}}(\tilde{\mathbf{z}})}{f_{\mathbf{x}}(\mathbf{z}) + f_{\mathbf{x}}(\tilde{\mathbf{z}})} \|\mathbf{z}_t - \tilde{\mathbf{z}}_t\|^2 d\mathbf{z} = \frac{1}{2T} \mathbb{E}_{\mathbf{x}} \sum_{t=1}^{T} \frac{f_{\mathbf{x}}(\tilde{\mathbf{x}})}{f_{\mathbf{x}}(\mathbf{x}) + f_{\mathbf{x}}(\tilde{\mathbf{x}})} \|\mathbf{z}_t - \tilde{\mathbf{z}}_t\|^2
$$

$$
= \frac{1}{2T} \mathbb{E}_{\mathbf{x}} \sum_{t=1}^{T} \frac{f_{\mathbf{x}}(\alpha^{-1}(\mathbf{x}))}{f_{\mathbf{x}}(\mathbf{x}) + f_{\mathbf{x}}(\alpha^{-1}(\mathbf{x}))} \|\mathbf{x}_t - \alpha_t^{-1}(\mathbf{x}_t)\|^2,
$$

which proves (5.6). Again, if we can choose $\mathbf{S}_t$'s, $\mathbf{b}_t$'s where $\alpha_t()$ is mirroring across planes given by $\mathbf{S}_t \mathbf{x} = \mathbf{b}_t$ such that,

$$
f_{\mathbf{x}}(\mathbf{x}) = f_{\mathbf{x}}(\alpha^{-1}(\mathbf{x})), \ \forall \mathbf{x} \in \mathbb{R}^{nT},
$$

the distortion $D_E$ becomes,

$$
D_E = \frac{1}{4T} \mathbb{E}_{\mathbf{x}} \sum_{t=1}^{T} \|\mathbf{x}_t - \alpha_t^{-1}(\mathbf{x}_t)\|^2 \overset{(a)}{=} \frac{1}{T} \sum_{t=1}^{T} \mathbb{E}_{\mathbf{x}_t} \|\mathbf{S}_t \mathbf{x}_t - \mathbf{b}_t\|^2
$$

$$
= \frac{1}{T} \sum_{t=1}^{T} \mathrm{tr} \left( \mathbf{S}_t R_{\mathbf{x}_t} \mathbf{S}_t' + (\mathbf{b}_t - \mathbf{S}_t \mu_{\mathbf{x}_t})(\mathbf{b}_t - \mathbf{S}_t \mu_{\mathbf{x}_t})' \right),
$$

where (a) follows as $\alpha_t(.)$ is mirroring across plane given by $\mathbf{S}_t \mathbf{x} = \mathbf{b}_t$, and thus $\alpha_t(\mathbf{x}) = \alpha_t^{-1}(\mathbf{x}) = (\mathbf{I} - 2\mathbf{S}_t'\mathbf{S}_t)\mathbf{x}_t + 2\mathbf{S}_t'\mathbf{b}_t$. This proves (5.9).

### 5.4.2  Proof of Theorem 5.3.5

Since given $\mathbf{z}$, there are $2^k$ possibilities of $\mathbf{x}_1^T$; $\mathbf{x}_1^T = \alpha^{-1(K)}(\mathbf{z}), K \in [0 : 2^k - 1]$, we start by computing,

$$
p_{\mathbf{z}}^{(K)} := Pr(\mathbf{x}_t = \alpha_t^{-(K)}(\mathbf{z}_t)|\mathbf{z}) = Pr(\mathbf{x} = \alpha^{-(K)}(\mathbf{z})|\mathbf{z})
$$

$$
= \frac{1}{f_{\mathbf{z}}(\mathbf{z})} Pr(\mathbf{z}|\mathbf{x} = \alpha^{-(K)}(\mathbf{z})) f_{\mathbf{x}}(\alpha^{-(K)}(\mathbf{z}))
$$

$$
\overset{(a)}{=} \frac{f_{\mathbf{x}}(\alpha^{-(K)}(\mathbf{z}))}{\sum_{j=0}^{2^k-1} f_{\mathbf{x}}(\alpha^{-(j)}(\mathbf{z}))}, \ K \in [0 : 2^k - 1],
$$

where $(a)$ follows by noting that $Pr(\mathbf{z}|\mathbf{x} = \alpha^{-(K)}(\mathbf{z})|\mathbf{z})$ is equal to the probability of the key being equal to $K$, which is $1/2^k$. Let $S = \sum_{j=0}^{2^k-1} f_{\mathbf{x}}(\alpha^{-(j)}(\mathbf{z}))$. Then $\mathbb{E}(\mathbf{x}_t|\mathbf{z})$ equals

$$\sum_{K=0}^{2^k-1} \alpha^{-(K)}(\mathbf{z}_t)p_{\mathbf{z}}^{(K)} = \frac{1}{S} \sum_{K=0}^{2^k-1} \alpha^{-(K)}(\mathbf{z}_t)f_{\mathbf{x}}(\alpha^{-(K)}(\mathbf{z})).$$

We can then compute $\mathrm{tr}\left(R_{\mathbf{x}_t|\mathbf{z}}\right)$ as,

$$\mathbb{E}_{\mathbf{x}_t|\mathbf{z}}\|\mathbf{x}_t - \mathbb{E}(\mathbf{x}_t|\mathbf{z})\|^2 = \frac{1}{S^3} \sum_{K=0}^{2^k-1} f_{\mathbf{x}}\left(\alpha^{-(K)}(\mathbf{z})\right)\left\|R_t^{(K)}\right\|^2,$$

where $R_t^{(K)} = \sum_{\ell=0}^{2^k-1} f_{\mathbf{x}}(\alpha^{-(\ell)}(\mathbf{x}))\left(\alpha_t^{-(\ell)}(\mathbf{x}_t) - \alpha_t^{-(K)}(\mathbf{x}_t)\right)$. Plugging $tr\left(R_{\mathbf{x}_t|\mathbf{z}}\right)$ in the expression of $D_E$ gives (5.12). Moreover, if condition (5.13) is met, (5.12) simplifies to (5.14).

### 5.4.3   Related Work

Secure data communication where the adversary has unlimited computational power is studied from the lens of information theory, most notably by Shannon [Sha49] and Wyner [Wyn75]. The study of secure communication while using distortion as a measure of security is relatively new and is first studied by Yamamoto [Yam88], where the goal is to maximize the distortion of an eavesdropper's estimate on a message, viewed from an asymptotic (in block length) information-theoretic approach. Schieler and Cuff [SC14] later showed that, in the limit of an infinite block length $n$ code, only $\log(n)$ bits of secret keys are needed to achieve the maximum possible distortion. The idea of using finite block length (and even single-shot) distortion as a performance measure was initiated in [TAF17], where schemes for single shot communication were considered. It demonstrated the exponential benefits for each additional bit of shared key. The schemes examined were for single-shot sensor observations, and not for time-series data, which is the focus of our work in this paper.

Secure communication in control systems is studied in [TGP17b, TGP17a, TSS17, MMS13, CDH16]. Securing the system state from an adversary was explored in [TGP17b, TGP17a], where an asymptotic steady-state analysis was explored. In contrast our work also deals with transients and is not asymptotic. Information-theoretic security was explored in [TSS17],

where the mutual information was used as a privacy measure. Security of the terminal state is considered in [MMS13] where an adversary makes partial noisy measurement of the state trajectory. Differential privacy for control systems was explored in [CDH16], which uses standard statistical indistinguishability which is equally applicable to categorical (non-metric space) data; in our work, we use the estimation error of the adversary in order to quantify privacy, utilizing the fact that CPS data lies in an Euclidean space, as argued earlier.

# CHAPTER 6

# Conclusion and Open Questions

In this thesis, we discussed a general methodology for designing security schemes for different communication systems. The core of this design methodology is that it is application-tailored: the nature of the application influences how the security scheme is designed. In doing so, various aspects of the application affects the design of the scheme, *e.g.*, the nature of the adversary and the assumptions on its capabilities, the performance requirements dictated by the application, the available opportunities and resources, as well as the right choice of a security metric. We considered three different applications for which we applied the aforementioned methodology for designing a security solution. These applications were namely: private data broadcasting in the context of index coding, established secret keys between communicating parties in mmWave systems, and distortion-based security in CPSs. We showed how a good understanding of these aspects allow for the designed of a suitable security solution for the application at hand.

## 6.1 Application 1: Private Data Broadcasting

For the problem of private data broadcasting, we considered an index coding instance where some clients are malicious: they wish to learn information about the requests and side information of the other clients. We first showed how this setup can cause a privacy breach in case there exists an honest-but-curious client. We showed how this breach can be manifested just by learning the encoding matrix used for the index code. To capture the amount of information leakage, we proposed the use of two different metrics: a conditional entropy metric and maximum information leakage. To provide a security solution, we first attempted

to design index codes that are private-aware, *i.e.*, it provides a trade-off between the amount of privacy achieved for the requests and the side information sets of the clients. We then proposed a different scheme for maintaining privacy. The scheme, referred to as $k$-limited-access scheme, transforms an index code into an equivalent one in which each client needs only $k$ transmissions to decode their request. We showed that this scheme is order optimal in some regimes, and for other regimes we provide several heuristic methods to design $k$-limited-access schemes.

Several open questions remain. First, our analysis for private-aware index codes is only for a specific configuration of clients and adversaries. A general analysis for the scheme for arbitrary configurations is still an open question. Second, we showed the order-optimality of our proposed polynomial-time $k$-limited-access schemes for some parameter regimes (namely, when the number of distinct clients is maximum). Whether an optimal polynomial-time scheme exists for general parameter regimes is still an open question. In addition, a rigorous analysis of the provided heuristics for designing $k$-limited-access schemes is also an open question.

Protecting data privacy when using broadcasting domains is a challenging problem. The nature of the broadcast channel allows all parties with access to the channel to receive the broadcast data. This leads to an unavoidable situation where malicious parties get access to data that, if not well protected, can leak sensitive information. Our work was a theoretical characterization of this problem through the index coding model. Other approaches can be attempted using different models for broadcast channels. In addition, a real-world practical implementation of privacy-aware encoding schemes is also an interesting area of research.

## 6.2  Application 2: Secret Key Establishment in mmWave Systems

Next, we discussed the problem of establishing secret keys in mmWave systems. We first remarked that an adversary with a quantum-computing capabilities may be able to break standard number-based encryption mechanisms that are nowadays used. To circumvent this, we proposed the use of a physical-layer-based secret key generation scheme. The proposed

scheme relied on the idea of packet erasures happening on the adversary's side. We proposed the use of mmWave beamforming and wiretap codes to increase the likelihood of such erasures. We should the anticipated benefits of using our scheme in two different scenarios, namely IEEE 802.11ay networks and vehicle platooning. In both cases, results showed that a few hundred Mbps of secret key can be established.

Physical-layer security is an attractive area of research with much potential and many interesting open problems. The emergence of mmWave communication systems has further increased the potential of such security schemes due to the inherent directionality in mmWave communications. However, mmWave communications also come with the challenge of beam coordination in order to establish links between communicating parties. Therefore, developing practical and high performing and reliable physical layer security schemes in mmWave communication systems is a challenge open area of research. Our work studies the problem of secret key establishment using physical layer security techniques. As we remarked earlier, a main open research problem is an actual implementation of the aforementioned scheme. In this work, results were obtained based on numerical simulations and a few theoretical assumptions. For example, ideal wiretap codes were assumed to be used, and the corresponding data communication rates expressions were used. In addition, some practical issues were not taken into account during this initial study, such as packet drop rates, key mismatch, etc. It would be quite an interesting research work to have an actual implementation of the end-to-end system to account for all missing factors in this initial study and provide a realistic assessment for the benefits of such a system.

## 6.3   Application 3: Distortion-based Security in CPSs

The final application is related to security in CPSs. We began noting that, in many applications, an adversary would be interested in learning information about the state vector of the CPSs. In which case, a more suitable goal may be to influence the adversary's estimate of the state vector to "far" from the actual state vector. Based on this observation, we suggested the use of a distortion-based metric, which captures the distance as well as the likelihood of

107

the adversary's estimate in comparison to the actual value. Targeting this metric as our security metric, we proposed the use of mirroring-based schemes which utilizes a small number of secret keys. We showed that, for a specific class of state vector distributions, our proposed scheme is optimal.

Several open research questions remain. The general class of distributions for which our proposed scheme is optimal is still not characterized. In addition, a more ambitious goal is to characterize a general encryption scheme that is optimal for arbitrary distributions. Finally, our proposed scheme works towards optimizing the proposed distortion-based metric; in which case, it sometimes suffices to use exactly one bit of secret key to achieve optimality. However, in other applications, it may be equally important to lead the adversary into having a larger ambiguity set (*i.e.*, a larger set of possible estimates). The use of one-bit of shared key does not appear to achieve such a goal. In addition, it is not clear how several keys of shared keys can be used in a manner that optimizes the distortion-based metric. Therefore, combining these two metrics (distortion-based and ambiguity sets) and developing correspondingly suitable encryption schemes are still open questions.

Distortion-based security in CPS is a relatively recent area of research. The general aim is to secure the system against an adversary, not necessarily by limiting its knowledge about the communicated information (*i.e.*, from an information-theoretic sense), but rather by confusing the adversary into performing a worse inference job. The notion of distortion-based security captures the effectiveness of the adversary's inference task and how far the outcome is from the actual value being inferred. Our work is an initial attempt to devise efficient security schemes from a distortion-based security perspective, where an eavesdropper is interested in estimating the state vector. However, in general, an adversary may be interested in specific functions of the underlying variables (*e.g.*, a function of the state vectors). In this case, designing an efficient (optimal) security scheme is an open question. In addition, in any such application, the evaluation of the attained level of distortion is dependent on the understanding of the random process of the underlying system variables. This underlying random process is dependent on the nature of the CPS application. Understanding this process is also an interesting open question.

# REFERENCES

[ADD13]  Katerina Argyraki, Suhas Diggavi, Melissa Duarte, Christina Fragouli, Marios Gatzianas, and Panagiotis Kostopoulos. "Creating secrets out of erasures." In *Proceedings of the 19th annual international conference on Mobile computing & networking*, pp. 429–440. ACM, 2013.

[AGJ10]  Assad Al Alam, Ather Gattami, and Karl Henrik Johansson. "An experimental study on the fuel reduction potential of heavy duty vehicle platooning." In *13th International IEEE Conference on Intelligent Transportation Systems*, pp. 306–311. IEEE, 2010.

[AKD18]  G. K. Agarwal, M. Karmoose, S. Diggavi, C. Fragouli, and P. Tabuada. "Distorting an Adversary's View in Cyber-Physical Systems." In *IEEE Conference on Decision and Control (CDC)*, pp. 1476–1481, Dec 2018.

[AP08]  Charu C Aggarwal and S Yu Philip. "A general survey of privacy-preserving data mining models and algorithms." In *Privacy-preserving data mining*, pp. 11–52. Springer, 2008.

[Arm]  Taylor Armerding. "The 18 biggest data breaches of the 21st century." https://www.csoonline.com/article/2130877/the-biggest-data-breaches-of-the-21st-century.html.

[Art]  Charles Arthur. "Security leak leaves US Apple iPad owners at risk." https://www.theguardian.com/technology/2010/jun/10/apple-ipad-security-leak?INTCMP=SRCH.

[BBJ11]  Ziv Bar-Yossef, Yitzhak Birk, TS Jayram, and Tomer Kol. "Index coding with side information." *IEEE Transactions on Information Theory*, **57**(3):1479–1494, Mar. 2011.

[BCR87]  G. Brassard, C. Crepeau, and J.-M. Robert. "All-or-nothing disclosure of secrets." *Advances in Cryptology: Proceedings of Crypto '86, Springer-Verlag*, pp. 234–238, 1987.

[BKL10]  Anna Blasiak, Robert Kleinberg, and Eyal Lubetzky. "Index coding via linear programming." *arXiv preprint arXiv:1004.1379*, 2010.

[BU16]  Karim Banawan and Sennur Ulukus. "The capacity of private information retrieval from coded databases." *arXiv:1609.08138*, Sep. 2016.

[BU18a]  Karim Banawan and Sennur Ulukus. "The capacity of private information retrieval from Byzantine and colluding databases." *IEEE Transactions on Information Theory*, 2018.

[BU18b]  Karim Banawan and Sennur Ulukus. "The capacity of private information retrieval from coded databases." *IEEE Transactions on Information Theory*, **64**(3):1945–1956, 2018.

[CDH16]    J. Cortés, G. E. Dullerud, S. Han, J. L. Ny, S. Mitra, and G. J. Pappas. "Differential privacy in control and network systems." In *2016 IEEE 55th Conference on Decision and Control (CDC)*, pp. 4252–4272, Dec 2016.

[CKG98]    Benny Chor, Eyal Kushilevitz, Oded Goldreich, and Madhu Sudan. "Private information retrieval." *Journal of the ACM (JACM)*, **45**(6):965–981, Nov. 1998.

[CN15]     Guangliang Chen and Deanna Needell. "Compressed sensing and dictionary learning." *Finite Frame Theory: A Complete Introduction to Overcompleteness*, **73**:201, January 2015.

[CPF15]    László Czap, Vinod M Prabhakaran, Christina Fragouli, and Suhas N Diggavi. "Secret communication over broadcast erasure channels with state-feedback." *IEEE Transactions on Information Theory*, **61**(9):4788–4808, 2015.

[CS08]     Mohammad Asad R Chaudhry and Alex Sprintson. "Efficient algorithms for index coding." In *IEEE INFOCOM Workshops 2008*, pp. 1–4, 2008.

[CWJ17]    Zhen Chen, Zhiying Wang, and Syed Jafar. "The capacity of private information retrieval with private side information." *arXiv preprint arXiv:1709.03022*, 2017.

[DSC12]    Son Hoang Dau, Vitaly Skachek, and Yeow Meng Chee. "On the security of index coding with side information." *IEEE Transactions on Information Theory*, **58**(6):3975–3988, Jun. 2012.

[Dzi06]    Stefan Dziembowski. "Intrusion-resilience via the bounded-storage model." In *Theory of Cryptography Conference*, pp. 207–224. Springer, 2006.

[EK11]     Abbas El Gamal and Young-Han Kim. *Network information theory*. Cambridge university press, 2011.

[EKF]      Yahya H Ezzeldin, Mohammed Karmoose, and Christina Fragouli. "Communication vs distributed computation: an alternative trade-off curve." In *2017 IEEE Information Theory Workshop (ITW)*, pp. 279–283.

[ELH]      Homa Esfahanizadeh, Farshad Lahouti, and Babak Hassibi. "A matrix completion approach to linear index coding problem." In *2014 IEEE Information Theory Workshop (ITW 2014)*.

[ets]      "ETSI Cyber Security." https://www.etsi.org/technologies/cyber-security.

[FGH16]    Ragnar Freij-Hollanti, Oliver Gnilke, Camilla Hollanti, and David Karpuk. "Private Information Retrieval from Coded Databases with Colluding Servers." *arXiv:1611.02062*, Nov. 2016.

[FLW06]    Christina Fragouli, Jean-Yves Le Boudec, and Jörg Widmer. "Network coding: an instant primer." *ACM SIGCOMM Computer Communication Review*, **36**(1):63–68, 2006.

[FN12]     Pedro Fernandes and Urbano Nunes. "Platooning with IVC-enabled autonomous vehicles: Strategies to mitigate communication delays, improve safety and traffic flow." *IEEE Transactions on Intelligent Transportation Systems*, **13**(1):91–106, 2012.

[Gar]      Simson Garfinkel. "Hackers Are the Real Obstacle for Self-Driving Vehicles." https://www.technologyreview.com/s/608618/hackers-are-the-real-obstacle-for-self-driving-vehicles/.

[GEN16]    Wei Gao, Sam Emaminejad, Hnin Yin Yin Nyein, Samyuktha Challa, Kevin Chen, Austin Peck, Hossain M Fahad, Hiroki Ota, Hiroshi Shiraki, Daisuke Kiriya, et al. "Fully integrated wearable sensor arrays for multiplexed in situ perspiration analysis." *Nature*, **529**(7587):509, 2016.

[Goo]      Dan Goodin. "Google+ shutting down after data leak affecting 500,000 users." https://arstechnica.com/tech-policy/2018/10/google-exposed-non-public-data-for-500k-users-then-kept-it-quiet/.

[Gre]      Andy Greenberg. "Hackers remotely kill a jeep on the highway-with me in it." https://www.wired.com/2015/07/hackers-remotely-kill-jeep-highway/.

[GSC17]    Yasaman Ghasempour, Claudio RCM da Silva, Carlos Cordeiro, and Edward W Knightly. "IEEE 802.11 ay: Next-generation 60 GHz communication for 100 Gb/s Wi-Fi." *IEEE Communications Magazine*, **55**(12):186–192, 2017.

[HBH17]    Anahita Hosseini, Chris Buonocore, Sepideh Hashemzadeh, Hannaneh Hojaiji, Haik Kalantarian, Costas Sideris, Alex Bui, Christine King, and Majid Sarrafzadeh. "Feasibility of a secure wireless sensing smartwatch application for the self-management of pediatric asthma." *Sensors*, **17**(8):1780, 2017.

[HE15]     Xiao Huang and Salim El Rouayheb. "Index coding and network coding via rank minimization." In *IEEE Information Theory Workshop-Fall (ITW)*, pp. 14–18. IEEE, 2015.

[HFA18]    Jehad M Hamamreh, Haji M Furqan, and Huseyin Arslan. "Classifications and Applications of Physical Layer Security Techniques for Confidentiality: A Comprehensive Survey." *IEEE Communications Surveys & Tutorials*, 2018.

[HL12]     Ishay Haviv and Michael Langberg. "On linear index coding for random graphs." In *IEEE International Symposium on Information Theory (ISIT)*, 2012.

[IEE16a]   IEEE. "IEEE Standard for Wireless Access in Vehicular Environments." Standard, 2016.

[IEE16b]   IEEE. "IEEE Standard for Wireless Access in Vehicular Environments–Security Services for Applications and Management Messages." Standard, 2016.

[IF]       Mike Isaac and Sheera Frenkel. "Facebook Security Breach Exposes Accounts of 50 Million Users." https://www.nytimes.com/2018/09/28/technology/facebook-hack-data-breach.html.

[IKW16]  Ibrahim Issa, Sudeep Kamath, and Aaron B Wagner. "An operational measure of information leakage." In *IEEE Annual Conference on Information Science and Systems (CISS)*, pp. 234–239, 2016.

[iso]  "ISO/IEC 27000 family - Information Security Management Systems." https://www.iso.org/isoiec-27001-information-security.html.

[KL14]  Jonathan Katz and Yehuda Lindell. *Introduction to modern cryptography*. CRC press, 2014.

[KM12]  Vijay Kumar and Nathan Michael. "Opportunities and challenges with autonomous micro aerial vehicles." *The International Journal of Robotics Research*, **31**(11):1279–1291, 2012.

[KMY19]  Mohammed Karmoose, Rafael Misoczki, Liuyang Yang, Xiruo Liu, Moreno Ambrosin, and Manoj R Sastry. "Methods and arrangements for vehicle-to-vehicle communications.", February 7 2019. US Patent App. 15/848,785.

[KS13]  P. Koopman and C. Szilagyi. "Integrity in embedded control networks." *IEEE Security Privacy*, **11**(3):61–63, May 2013.

[KSC17]  Mohammed Karmoose, Linqi Song, Martina Cardone, and Christina Fragouli. "Private Broadcasting: an Index Coding Approach." In *IEEE International Symposium on Information Theory (ISIT)*, pp. 2548–2552, June 2017.

[Lar]  Selena Larson. "FDA confirms that St. Jude's cardiac devices can be hacked." https://money.cnn.com/2017/01/09/technology/fda-st-jude-cardiac-hack/.

[LHO]  Jinlong Lu, J Harshan, and Frédérique Oggier. "A USRP implementation of wiretap lattice codes." In *2014 IEEE Information Theory Workshop (ITW)*, pp. 316–320.

[LLV07]  Ninghui Li, Tiancheng Li, and Suresh Venkatasubramanian. "t-closeness: Privacy beyond k-anonymity and l-diversity." In *IEEE 23rd International Conference on Data Engineering*, pp. 106–115, Apr. 2007.

[LMY18]  Songze Li, Mohammad Ali Maddah-Ali, Qian Yu, and A Salman Avestimehr. "A Fundamental Tradeoff Between Computation and Communication in Distributed Computing." *IEEE Transactions on Information Theory*, **64**(1):109–128, 2018.

[MDP14]  Manoj Mishra, Bikash Kumar Dey, Vinod M Prabhakaran, and Suhas Diggavi. "The oblivious transfer capacity of the wiretapped binary erasure channel." In *IEEE International Symposium on Information Theory*, pp. 1539–1543, Jun. 2014.

[MGL16]  Markus Maurer, J Christian Gerdes, Barbara Lenz, Hermann Winner, et al. "Autonomous driving." *Berlin, Germany: Springer Berlin Heidelberg*, **10**:978–3, 2016.

[MMS13]   Waseem A. Malik, Nuno C. Martins, and Ananthram Swami. *LQ Control under Security Constraints*, pp. 101–120. Springer Int. Publishing, Heidelberg, 2013.

[MV11]   Hessam Mahdavifar and Alexander Vardy. "Achieving the secrecy capacity of wiretap channels using polar codes." *IEEE Transactions on Information Theory*, **57**(10):6428–6443, 2011.

[NCF14]   Thomas Nitsche, Carlos Cordeiro, Adriana B Flores, Edward W Knightly, Eldad Perahia, and Joerg C Widmer. "IEEE 802.11 ad: directional 60 GHz communication for multi-Gigabit-per-second Wi-Fi." *IEEE Communications Magazine*, **52**(12):132–141, 2014.

[Net16]   Technical Specification Group Radio Access Network. *Study on channel model for frequency spectrum above 6 GHz*. 3rd Generation Partnership Project, July 2016.

[nis]   "NIST Cybersecurity Framework." https://www.nist.gov/cyberframework.

[NKL18]   Lakshmi Natarajan, Prasad Krishnan, and V Lalitha. "On Locally Decodable Index Codes." In *IEEE International Symposium on Information Theory (ISIT)*, June 2018.

[NPR18]   Varun Narayanan, Vinod M Prabhakaran, Jithin Ravi, Vivek K Mishra, Bikash K Dey, and Nikhil Karamchandani. "Private Index Coding." In *IEEE International Symposium on Information Theory (ISIT)*, pp. 596–600, 2018.

[Oxl06]   James G Oxley. *Matroid theory*, volume 3. Oxford University Press, USA, 2006.

[PSV]   Jeroen Ploeg, Bart TM Scheepers, Ellen Van Nunen, Nathan Van de Wouw, and Henk Nijmeijer. "Design and experimental evaluation of cooperative adaptive cruise control." In *2011 14th International IEEE Conference on Intelligent Transportation Systems (ITSC)*, pp. 260–265.

[RBE10]   Ron Rubinstein, Alfred M Bruckstein, and Michael Elad. "Dictionaries for sparse representation modeling." *Proceedings of the IEEE*, **98**(6):1045–1057, June 2010.

[SC14]   C. Schieler and P. Cuff. "Rate-Distortion Theory for Secrecy Systems." *IEEE Trans.Info. Theory*, **60**(12):7584–7605, Dec 2014.

[SCA16]   Iris Safaka, László Czap, Katerina Argyraki, and Christina Fragouli. "Creating Secrets Out of Packet Erasures." *IEEE Transactions on Information Forensics and Security*, **11**(6):1177–1191, 2016.

[SF15]   Linqi Song and Christina Fragouli. "Content-type coding." In *International Symposium on Network Coding (NetCod)*, pp. 31–35, Jun. 2015.

[Sha49]   Claude E Shannon. "Communication theory of secrecy systems." *Bell system technical journal*, **28**(4):656–715, 1949.

[SJ16]     Hua Sun and Syed A Jafar. "The Capacity of Private Information Retrieval."
           *arXiv:1602.09134*, Feb. 2016.

[SJ18]     Hua Sun and Syed Ali Jafar. "Private Information Retrieval from MDS Coded
           Data With Colluding Servers: Settling a Conjecture by Freij-Hollanti et al." *IEEE
           Transactions on Information Theory*, **64**(2):1000–1022, 2018.

[SKK18]    Jani Suomalainen, Adrian Kotelba, Jari Kreku, and Sami Lehtonen. "Evaluat-
           ing the efficiency of physical and cryptographic security solutions for quantum
           immune iot." *Cryptography*, **2**(1):5, 2018.

[Sti05]    Douglas R Stinson. *Cryptography: theory and practice.* Chapman and Hall/CRC,
           2005.

[TAF17]    C. Tsai, G. K. Agarwal, C. Fragouli, and S. Diggavi. "A distortion based ap-
           proach for protecting inferences." In *2017 IEEE International Symposium on
           Information Theory (ISIT)*, pp. 1913–1917, June 2017.

[TDC07]    Andrew Thangaraj, Souvik Dihidar, A Robert Calderbank, Steven W McLaugh-
           lin, and Jean-Marc Merolla. "Applications of LDPC codes to the wiretap chan-
           nel." *IEEE Transactions on Information Theory*, **53**(8):2933–2945, 2007.

[TGP17a]   Anastasios Tsiamis, Konstantinos Gatsis, and George J. Pappas. "State Es-
           timation with Secrecy against Eavesdroppers." *20th IFAC World Congress*,
           **50**(1):8385 – 8392, 2017.

[TGP17b]   Anastasios Tsiamis, Konstantinos Gatsis, and George J. Pappas. "State-Secrecy
           Codes for Networked Linear Systems." *CoRR*, **abs/1709.04530**, 2017.

[THM15]    W. Trappe, R. Howard, and R. S. Moore. "Low-Energy Security: Limits and
           Opportunities in the Internet of Things." *IEEE Security Privacy*, **13**(1):14–21,
           Jan 2015.

[TNM14]    Timothy A Thomas, Huan Cong Nguyen, George R MacCartney, and Theodore S
           Rappaport. "3D mmWave channel model proposal." In *Vehicular Technology
           Conference (VTC Fall), 2014 IEEE 80th*, pp. 1–6. IEEE, 2014.

[TR16]     Razan Tajeddine and Salim El Rouayheb. "Private Information Retrieval from
           MDS Coded Data in Distributed Storage Systems." *arXiv:1602.01458*, Feb. 2016.

[TSS17]    T. Tanaka, M. Skoglund, H. Sandberg, and K. H. Johansson. "Directed in-
           formation and privacy loss in cloud-based control." In *2017 American Control
           Conference (ACC)*, pp. 1666–1672, May 2017.

[Une]      Unex. "OBU-201U Specification — V2X On-Board Unit, IEEE 1609.x protocol
           stack." datasheet.

[WHG14]  Cheng-Xiang Wang, Fourat Haider, Xiqi Gao, Xiao-Hu You, Yang Yang, Dongfeng Yuan, Hadi M Aggoune, Harald Haas, Simon Fletcher, and Erol Hepsaydir. "Cellular architecture and key technologies for 5G wireless communication networks." *IEEE communications magazine*, **52**(2):122–130, 2014.

[WLF16]  J. Wan, A. B. Lopez, and M. A. Al Faruque. "Exploiting Wireless Channel Randomness to Generate Keys for Automotive Cyber-Physical System Security." In *2016 ACM/IEEE 7th International Conference on Cyber-Physical Systems (IC-CPS)*, pp. 1–10, April 2016.

[Wyn75]  Aaron D Wyner. "The wire-tap channel." *Bell system technical journal*, **54**(8):1355–1387, 1975.

[Yam88]  H. Yamamoto. "A rate-distortion problem for a communication system with a secondary decoder to be hindered." *IEEE Transactions on Information Theory*, **34**(4):835–842, July 1988.

[ZGH13]  B. Zan, M. Gruteser, and F. Hu. "Key Agreement Algorithms for Vehicular Communication Networks Based on Reciprocity and Diversity Theorems." *IEEE Tran. Vehicular Tech.*, **62**(8):4020–4027, Oct 2013.

[ZWW17]  Yongxu Zhu, Lifeng Wang, Kai-Kit Wong, and Robert W Heath. "Secure communications in millimeter wave ad hoc networks." *IEEE Transactions on Wireless Communications*, **16**(5):3205–3217, 2017.