# Chapter 1

# Introduction

The trend of both curiosity and profit-driven human development has caused a surge in the amount of openly accessible raw data. More often than not, the data is generated much faster than it can be processed into something interpretable or useful. In the endeavor of keeping up with the inflow of information, there are two major factors that significantly hinder our progress. First, Moore's law implicitly sets a physical limit to the number of transistors that can be placed on a chip, consequently limiting how powerful and how fast electronic systems can become (barring a paradigm shift in the fundamental design of semiconductors). The second is the Nyquist-Shannon sampling theorem (NST), which limits the range of frequencies a recording device can successfully capture. This study explores the use of compressive sensing (CS)—an emergent sampling theorem that allows reconstruction of signals from much fewer samples than required by the NST—as a viable method for compression, encryption, and/or enhancement. In this framework, the computational burden of encoding/decoding is shifted from the sampling device to the device performing reconstruction, decompression, or

other modes of post-processing. As such, there exist many ways to reconstruct a signal from compressive measurements. CS has found its applications in simple audio signals containing stable frequencies, such as pure tones [1, 2], and dynamic frequencies, such as speech [3–5], images [6–8], and grayscale videos [9, 10]. The formulation of a sensing matrix in CS requires a basis conforming to some uniform uncertainty principle, and most common starting points would be partial discrete cosine transforms (DCT) or partial discrete wavelet transforms (DWT). Recent studies, however, have shown that learned bases perform much better on more complex signals [11–13], i.e., those that would be typically encountered in real life situations. The learning algorithms associated with the construction of these bases range from classical iterative methods, which have long been used in optimization problems, to the more contemporary machine learning methods.

## 1.1 Related literature

In 2004, Candès, Romberg, Tao [14], and Donoho [15] asked the question,

> With the recent breakthroughs in lossy compression technologies, we now know that most of the data we acquire can be thrown away with minimal perceptual loss. Why bother to acquire all the data when we can just directly measure the part that will not be thrown away?

This was eventually answered in many different ways by the same people, ultimately birthing the field which we now know as compressive sensing. The methods in CS apply concepts from time-frequency uncertainty principles [16] and sparse representations, which were studied rigorously by Donoho and Elad [17]. CS can be viewed as a strategic undersampling method: the signal is sampled at random

locations in the real domain, and the ratio of the indices where it is sampled to the size of the signal can be associated with some quasi-frequency which may be lower than the Nyquist rate.

[18] demonstrated the use of deterministic chaos filters to acquire samples instead of random distributions. Sampling using a Gaussian-Logistic map was applied to acoustic signals in [1]. Normally, a deterministic chaotic function will need one or more initialization values as a "seed", and the sequence of numbers produced by different combinations of initial values rapidly diverge from each other. This phenomenon led to investigating the use of compressive sensing as an encryption algorithm. Simultaneous compression and encryption was achieved by [6], and it was found that the initial values were sensitive to perturbations on the order of $10^{-15}$. Their image compression-encryption model via compressive sensing was shown to have a key space on the order of $10^{34}$, making it extremely resistant to brute force and other types of attacks. This was extended in [7] to utilize a higher-dimensional variant of the Lorenz attractor, subsequently expanding the key space to the order of $10^{83}$. In the methods above, sampling was performed in the signal domain (i.e., temporal domain for audio, spatial domain for images), and the reconstruction was performed in the frequency domain with a DCT or similar basis. [2] proposed a method to perform both sampling and reconstruction in the time domain using differential evolution.

Audio signals, compared to images, are much more densely packed with information. Whereas images are not naturally bandlimited and rather, are dependent on the spatial resolution and bit depth of the imaging device, audio size scales proportionally with time and takes on a wider range of values. The accepted frequency range of human hearing is from 20 Hz to 20 kHz, so by the

NST, a sampling frequency of at least 40 kHz is needed to ensure that an audio sample is recorded correctly. Any meaningful audio recording, especially those containing speech, will certainly have a duration of a few seconds up to a few hours, so one cannot straightforwardly apply methodologies used for images or recordings with relatively static frequencies, as the first challenge this would pose for electronic systems is insufficient memory to process the entire signal all at once. Low [3, 4] circumvented this problem by transforming the signal to the modulation domain, i.e., the signal's spectrogram, essentially raising a one-dimensional signal to $N$-dimensions, where the value of $N$ is dependent on the desired spectrogram resolution, number of subbands, and percent overlap between adjacent subbands. In such signals, recordings with an observed noise floor could be easily be denoised, which is an inherent property of CS [19].

Due to the large size of video information as a consequence of its high dimensionality, it is possible, but impractical, to apply image CS techniques on an entire frame-by-frame basis. Correlations between adjacent frames are utilized instead, and can be obtained using dictionary learning [11] or principal components analysis [9]. For the same reason, the application of CS to grayscale videos presuppose the use of machine learning methods. Iliadis [20, 21] came up with two different deep neural network architectures whose inputs and outputs are patches derived from grayscale videos. This idea was utilized in [22] who modified the architecture into a residual network containing several convolutional layers. The original design was targeted towards image reconstruction, but could easily be extended to videos.

In the same vein, neural network methods could also be used in CS of speech. Advances in natural language processing were primarily made using recurrent neural

networks (RNN). In [23], a speech signal was first modeled by their proposed RNN architecture based on a noise-constrained least squares estimate, and final recovery is done via Kalman filtering. A new simple recurrent unit (SRU) network was created in [24] which maps the relation between noisy and clean speech recordings for speech enhancement.

## 1.2 Novelty

This study aims to provide a generalization for applying CS techniques to signals of arbitrary dimensions. Previous studies worked exclusively with either audio or image sequences as the target for CS, and due to the computational demands, the focus of most of the research in the field has been to optimize the computational complexity for real-time applications, and improve signal reconstruction quality. In the establishment of CS methods, two different general frameworks to compressively sample signals arise, namely, one-dimensional CS (1DCS) and two-dimensional CS (2DCS). It is shown that an $N$-dimensional signal can be decomposed into factors of one-dimensional and two-dimensional signals, and can be processed using methods appropriate for each type of signal. Furthermore, it is shown that $N$ is bound not only by the type of signals being worked with, but also the computational power of the decoding/decompressing device. In particular, large values of $N$ are useful in encryption, where a signal is first raised to a high dimension in a certain basis, the sensing matrix is derived from another high-dimensional basis, and the result is cast back to either one or two dimensions to yield the encrypted message.

## 1.3 Thesis overview

The next chapter establishes the relevant mathematical concepts and notation to be used throughout this study, algorithms used in signal reconstruction, and appropriate metrics per type of signal. Chapters 3–5 respectively focus on two-dimensional CS, one-dimensional CS, and $N$-dimensional CS. The reason behind the ordering of Chapters 3 & 4 will become apparent as the usage of spectrograms are introduced. Each of these chapters are self-contained methodologies, results, and discussions to emphasize that the methods can work independently of each other, save for the generalization to $N$-dimensions. The study is concluded and recommendations for future studies are presented in Chapter 6.