

1 Introduction

Any agent, be it a human, an animal or a robot, has to react to it's environment to take advantage of opportunities and to avoid dangers. The transformation of events to reaction can be partitioned into three steps: **(1) physical events** are transformed by sensors into **raw data**, **(2)** Computation transforms the **raw data** into **knowledge** (a representation of the environment), and **(3)** an **action** is chosen based on the acquired **knowledge**.

Rayan : I think we need to edit this a little. The first and second paragraphs feel disconnected from each other. **Yoav :** Rayan. I agree. See if the paragraph I added improves the flow. Feel free to edit.

When detecting and localizing physical there are significant advantages to using several sensors placed at different locations. That is why we humans have *two* eyes and *two* ears. Our two eyes provide our vision with depth, while our two ears enable us localize the direction from which sound is arriving. These abilities are based on *comparing* the signals arriving from the two eyes or the two ears, and the comparing signals requires *communication* between the sensors.

Moving from natural to artificial sensors. One is often interested in detecting and tracking events in larger areas, such as a house or a city. In such cases there is an additional reason to deploying multiple sensors: sensors typically have a finite range. Therefore the number of sensors we need is proportional to the area we wish to cover.

The result of these trends is a rapidly increasing deployment of *sensor networks*. These are composed of a large number of sensors, connected by a communication network.

The design of individual sensors is dominated by considerations of sensitivity and resolution (temporal and spatial). The goal is to detect the smallest, faintest and most transient signals, by exploiting priors on the physical model of signal propagation and information on the location of the sensors. Computation is carried out on each sensor to identify significant events and to correlated them with events identified on other sensors.

These days the leading architecture of reactive systems is wireless sensor networks. Sensor networks consist of large numbers of small independent units, each with sensors, computation and wireless communication. Such systems are constrained by power and communication bandwidth.

One important consequence of these constraints is *pushing computation to the edge*. Instead of communicating the raw information from each sensor to a central computer, each sensor unit locally computes summaries, or sketches, which are shorter and therefore cheaper to communicate. This also reduces the computation load on the units that receive the information.

2 Framework

This proposal combines several related lines of work. To facilitate the exposition, we start by introducing some terminology and notation that will be used throughout. **Figure 1** describes a simple sensor network that tracks a car using the sound waves the car is emitting.

We now expand this simple example into a more general framework. We assume that the network consists of a n sensors and m targets. Sensor i 's state at time τ is denoted $\psi_i(\tau)$. Similarly, the state of target j at time τ is denoted $\theta_j(\tau)$. Here and in the rest of this section, we don't specify the spaces in which ψ_i or θ_j are members. This allows for a general introduction, and more which will be made more specific in later sections. When appropriate, we will denote the combined state of all targets by $\Theta(\tau)$ and the combined state of all sensor by $\Psi(\tau)$. The third state component is the state of physical environment in which the targets and sensors reside. We denote the state of the environment by $\mathbf{E}(\tau)$

The targets generate signals, which we call the *raw* signals. We denote the raw signal generated by target i as $\mathbf{x}_i(\tau)$. We denote the collection of all m signals by $\mathbf{X}(\tau)$. On the receiving end, each sensor i captures a digital signal \mathbf{y}_i . These digital signals are the inputs to the computations we will discuss. As the signals arrive at physically separated sensors, the computation is inherently distributed. The main goal of this proposal is to develop algorithms that achieve desired tasks with minimal communication between the sensors.

Rayan : How are we deciding what variables are in bold and what variables are not? For example, why is $y_i(t)$ not bold in the caption, but $\mathbf{x}_i(\tau)$ in bold?

The transfer function Φ defines the way by which the raw signals \mathbf{X} are transformed into the digitized signals \mathbf{Y} . This function represents both the point transfer function of the physical environment, the analog-to-digital transformation of the physical signal into a discrete time physical signal, and the noise that is added through this process. The transformation is defined by:¹

$$\mathbf{Y} = \Phi(\mathbf{X}, \Theta, \Psi, \mathbf{E})$$

Many of the problems we plan to tackle in this proposal are inverse computation problems. We assume that some aspects of the physical space are known, i.e. we know a subset of $\mathbf{X}, \Theta, \Psi, \mathbf{E}$. Given the digitized signals \mathbf{Y} , our task is to estimate the unknown parts of the physical space. Reliable methods for computing such estimates exist. However, they typically require high communication bandwidth. The goal of this proposal is to find distributed estimation algorithms that achieve good performance while using significantly less communication.

Some specific tasks

We give a few specific examples of tasks. We will elaborate on some of these tasks below

1. **Target Localization:** Figure 1 depicts an archetypal target localization task. In this case the locations of the sensors Ψ and the state of the environment \mathbf{E} are assumed to be fixed. A typical additional assumption, which is represented in the transfer function Φ , is that strongest signal corresponds to the straight line of transmission between the target and the sensors. A common approach to target localization is to estimate the delay between the arrival of the signal at different sensors by using some type of cross correlation []. This calculation is performed. It is well known that placing sensors far from each other provides the most accurate localization. However, achieving this accuracy with bounded communication between the sensors remains a challenge.
2. **Sensor Calibration:** A common situation is that many sensors are installed in an existing environment such as a home, and the state of these sensors Ψ which would typically include location and orientation, is not known. Manual calibration is often labor intensive or impossible. The challenge is to design algorithms through which the network can self-calibrate. We consider two types of calibration, in the easier *active calibration* the system can control the generated signals, while in *blind calibration* the calibration has to be done using signals that are generated by the environment.
3. **Signal reconstruction** Voice based systems such as speakerphones, and voice activated computers, need to reconstruct the speech signal. Microphone arrays are sensor networks where the sensor is a microphone. Accurately reproducing the speech signal when there is more than one speaker is an open challenge. In this problem the location of the sensors is known, the goal is to reproduce the raw signal \mathbf{X} . This task is relatively easy when Θ, Ψ, \mathbf{E} are known and the transfer function Φ is known and simple. It becomes significantly harder when some of Θ, Ψ, \mathbf{E} are unknown or when Φ is complex, such as multi-path radio signals or reverberating audio signals.
4. **Optimizing sensor placement:** The accuracy of target localization depends on the location of the target and of the sensors. While the location of the target is not under our control, the location of the sensors is. Methods for optimizing the locations of the sensors will be described in Section sec:sensor-placement.
5. **Mapping the environment** Sometime the goal of the system is to estimate the environment \mathbf{E} . One example is to use Radar, Sonar or Lidar to create a 2D or 3D representation of the environment for a smart car. Another example, coming from seismology is to use controlled vibration sources and many acceleration sensors to map the subterranean earth. In these settings the locations of the sensors Ψ and the targets Θ (called transmitters in this context) is fixed and known, as is the raw (transmitted) signal \mathbf{X} . The goal is to deduce \mathbf{E} from the collected signal \mathbf{Y} . (Peter, does this make sense to you, can you use this formulation in your sections about tomography and dictionaries?)
6. **Monitoring** In many situations the goal of the sensor network is to track the environment, identify trends and detect anomalies. Motivating examples include: Security systems, systems for monitoring patients or the elderly, highway monitoring and factory floor monitoring. Many of these environments are too complex

¹Note that the transfer function Φ operates on the whole sequences, not just on the sequence at a single time τ . That is because signal propagation takes time, so $\mathbf{Y}(t)$ depends on \mathbf{X} at multiple time points.

to estimate a fully detailed representation. Instead, we suggest building a statistical model which implicitly captures the major degrees of freedom of the environment and the way they relate to major variables such as time of day and day of week. The challenge here is to learn such a model in an unsupervised or weakly supervised way, without heavy use of computational or communication resources. Here, we will propose, and rigorously analyze techniques based on Kernel methods combined with sketching and low dimensional binary embedding techniques. The combination of these tools will simultaneously facilitate performing various desired statistical tasks, while minimizing storage, communication, and computational costs.

3 Target Localization using minimal number of sensors

Yoav : Piya, can you use the notation I defined in the framework section? Also I would like to merge this subsection and the following one, which describes sensing geometry and the goal of minimizing the number of sensors. A sensor network consisting of M sensing units aims to capture information of interest (often described in terms of parameters) regarding the physical environment by acquiring measurements in space (dictated by sensor locations) and in time (dictated by the sampling technique employed at each sensor). In many applications (especially those concerning high-resolution/super-resolution imaging), the goal is to detect certain parameters $\theta_i \in \mathbb{C}^P, i = 1, \dots, K$ from K targets of interest in the environment by acquiring signals emitted by them.

Yoav : This describes a more general framework, using \mathbb{C}^P (does that mean each coordinate is complex?). What is gained from this generality? maybe drop the general notation? Also how does high resolution/super-resolution fit here? If you have worked on such problem, I suggest you devote a paragraph and cite, rather than just mentioning in passing.

As an example, consider a network consisting of active radar units (for example, those mounted on autonomous vehicles) attempting to create a map of the environment. In this case, K can denote the total number of pedestrians, bicyclist's and other cars and $\theta_i \in \mathbb{R}^3$ for the i th target will consist of its location $\mathbf{x}_i = [x_i, y_i]^T$ and velocity (v_i) parameters, i.e.

$$\theta_i = [x_i, y_i, vx_i, vy_i]^T, \quad 1 \leq i \leq K \quad (1)$$

Yoav : Shouldn't K be estimated? Mathematically the space-time measurements collected at the m th sensing element can be described as

$$y_m(t) = \sum_{i=1}^K \phi(\mathbf{d}_m, \theta_i, t) + w_m(t), \quad 1 \leq m \leq M \quad (2)$$

where $w_m(t)$ is the additive noise. Here $\mathbf{d}_m \in \mathbb{R}^3$ denotes the location of the m th sensor and the function $\phi(\cdot)$ characterizes the measurement model (often referred to as the point-spread function in the context of imaging) that depends on the physical laws governing wave propagation, and properties of the medium. Depending on the application and model assumptions, the function $\phi(\cdot)$ can be linear, non-linear, and potentially, even non-convex. However, it can be *partially designed* by choice of sensor locations \mathbf{d}_m . This will be a key enabler towards obtaining compressed sketches of measurements (or reducing the number of sensing units) while preserving the ability to reliably infer the parameter $\theta_i, 1 \leq i \leq K$.

The basic model assumes targets as point sources, but in many situations, they are distributed. **Piya :** Perhaps Peter can help characterize this model, since SONAR deals with such targets.

The main objective is to obtain estimates $\hat{\theta}_i, 1 \leq i \leq K$ of the parameters of interest (θ_i) using *minimal number of measurements/minimizing the number of sensing elements*. These estimates essentially are some appropriate functions of the spatio-temporal measurements $Y_T = \{y_m(t), 1 \leq m \leq M, 1 \leq t \leq T\}$, i.e.,

$$\hat{\theta}_i(T) = g_i(Y_T) \quad (3)$$

In many scenarios, the parameters of interest can be reliably inferred from the *correlation of the measurements*. In other words, the correlation of the measurements act as a sufficient statistic for the parameters to

be inferred. Depending on the application, the correlation matrix can be spatial (when the source signals are stationary), or spatio-temporal (when the temporal dynamics need to be tracked, such as for change-point detection). In these cases, we can effectively summarize the large amount of raw sensor measurements by only retaining and communicating their correlation. **Yoav :** The way I was thinking about it, each sensor has only one signal. In a one scenario, the quantity of interest is the "time delay of arrival" or the time shift of one signal relative to another that would maximize the correlation. Is there anything known about computing this time delay without communicating the whole time series?

Spatial Correlation and Localization: Suppose we compute the spatial correlation between $y_m(t)$ and $y_n(t)$ by averaging over T time samples (the signals are assumed to be stationary over this interval)²

$$r_{m,n}(T) = \frac{1}{T} \sum_{t=1}^T y_m(t) y_n^*(t) \quad (4)$$

We can summarize the self and cross correlation between M time-series measurements (collected at M sensors) using these M^2 correlation values (collected in the form of a correlation matrix R_T). Owing to the geometry of the measurements, these correlation values directly depend on the sensor locations \mathbf{d}_m (via the mapping $\phi(\cdot)$). Hence, it is natural to ask

1. Can we exploit the geometry of the measurement model to further compress the correlation matrix R_T ? What is the role of sensor geometry in this case? We should still be able reliably infer $\theta_i, i = 1, 2, \dots, K$ from such a compressed sketch.
2. How large should M be (in comparison to K) ?

4 Localization of weak sources (Peter)

Yoav : Can the description of SCM be folded into Piya's introduction?

The focus here is detecting weak sources within a sensor network without a fusion center. To observe weak sources, as much information as possible should be used. Thus, at first there is no attempt to reduce the information in the data by sketching or special sensor arrangements. The network could consist of sensors with know location, partially unknown or unknown positions.

The propagation path from a given source location would here represent multiple propagation paths in a non-uniform media. The frequency domain transfer function from a source location to N receivers \mathbf{a} . Assuming K uncorrelated sources of complex amplitude \mathbf{s} at spatial location \mathbf{x}_k , the received signal $\mathbf{y} \in \mathcal{R}^N$ on N receivers is

$$\mathbf{y} = \mathbf{A}\mathbf{s} + \mathbf{n}, \quad (5)$$

where $\mathbf{A} = [\mathbf{a}_1, \dots, \mathbf{a}_K]$ and \mathbf{n} is uncorrelated noise. The sources might be located in the near field and composed of many propagation paths. Examples of many propagation paths from a single source could be waves from 1) a source in a house propagating though the air and though the wall. 2) a cell phone signal with a direct path, a reflected path or refracted path. 3) a car radiating noise though the air and though the ground. Further, the sensors are not placed in a regular order, but where practical and maybe with unknown location. Thus the elements in \mathbf{a}_k are unknown.

Yoav : What is the relationship between a_k and x_k ?

To make observations of weak sources we observe L snapshots assuming stationarity $\mathbf{Y} = [\mathbf{y}_1 \dots \mathbf{y}_L]$. We can here form the sample covariance matrix (SCM)

$$\mathbf{S} = \mathbf{Y}\mathbf{Y}^H / L \quad (6)$$

and form the the normalized SCM $\hat{\mathbf{S}}$ or coherence with elements

$$S_{ji} = \frac{S_{ji}}{\sqrt{S_{ii}S_{jj}}} \quad (7)$$

²Reasonable to do so when the source signals are stationary and emit independent signals. This is the common practice in source localization using antenna arrays. We can also use more sophisticated regularized estimation of correlation.

Yoav : I am confused about the definition of coherence, should it not be the maximum correlation when one signal is shifted relative to the other?

Forming the ensemble mean over multiple snapshots give the cross spectral density matrix $\mathbf{C} \in \mathcal{R}^N \times N$

$$\mathbf{C} = \mathcal{E}[\mathbf{y}\mathbf{y}^H] = \mathbf{A}\mathbf{S}\mathbf{A}^H + \mathbf{N}, \quad (8)$$

The array signal processing literature is ample with processing of this type, especially with the structure of the \mathbf{A} matrix partially known. In this work we will focus on pushing the computations to the sensor nodes and thus only observing part of SCM.

5 Passive narrow-band multi-target localization (Piya)

Yoav : How is this problem different from the problem of finding the best placement for the sensors in order to maximize the accuracy of localization, ignoring communication bandwidth?

With the aim of obtaining a compressive sketch of the correlation matrix (also termed as compressive covariance sensing), we will optimize the design of sensor array (i.e. choice of $\mathbf{d}_m, 1 \leq m \leq M$) by understanding how the array geometry controls the algebraic structure of R_T . One of the main objectives will be to understand how much communication is needed (and between which subset of sensors) to achieve a certain level of accuracy. To illustrate this, we briefly discuss Co-PI Pal's recent work in structured sampler design (e.g., nested, coprime and generalized nested samplers) which utilize the idea of difference sets.

- **Difference set-inspired Designs:** I will review some results in the context of array processing and DOA estimation...(to be filled in).
- **Proposed Research:** Motivated by these results, our goal will be to develop a rigorous framework for further developing the key idea of correlation-aware sensing to a distributed scenario and make it applicable for imaging problems beyond point target localization.

1. **Distributed Sensing:** The idea of difference set inspired sampler design can be actually generalized beyond that of antenna arrays, to acquire *compressive sketches* of the correlation between signals acquired between pairs of sensors. In general, given N sensors, it is natural to think that one needs to compute the correlation between all $\binom{N}{2}$ time series (from all possible sensor-pairs) to construct the overall $N \times N$ correlation matrix R_T . However, using the idea of difference-set sampling, one can only compute cross-correlation values between a much smaller subset of size $\Theta(\sqrt{N})$ of *suitably selected sensor-pairs* and recreate the entire $N \times N$ correlation matrix R_T . In the context of distributed sensing, this automatically means that only these sensors need to communicate and exchange information.

Exploiting Distance-based Redundancies: The key idea behind achieving such reduction is to exploit the redundancies present in the correlation values that naturally result from the physical spatial signal model. A widely used example of such a redundancy is that the correlation $r_{m,n} = E(y_m(t)y_n^*(t))$ between m th and n th sensors is of the following form

$$r_{m,n} \approx f(\mathbf{d}_m - \mathbf{d}_n) \quad (9)$$

In other words, the correlation is spatially only a function of the *inter-sensor distance*, and this is a direct consequence of the functional form of $\phi(\cdot)$. **Piya :** Can give specific examples if needed. This is also referred to as spatial stationarity and it is (exactly or approximately) true for many applications as narrowband and wideband radar ³, super-resolution optical imaging [], mmWave wireless channels [] and so forth. Hence, depending on the inter-sensor distances, many of these $\binom{N}{2}$ correlation values are actually repeated/redundant. Based on this observation, we propose to use a new sketching technique developed by co-PI Pal, called **Generalized Nested Sampling (GNS) to reduce the amount of inter-sensor communication**. Suppose the sensors are located on a uniform grid. In one dimension, (9) implies that the correlation matrix R_T has Toeplitz structure and GNS provides an optimal way to select sensors to sketch such a matrix.

Definition 1 Piya : Definition of GNS goes here..

³In the latter case, this holds at individual frequency bands after splitting the wideband signal into narrow frequency bins using a filter bank

Hence, GNS dictates how to select a subset \mathcal{S}_{GNS} of $M = \Theta(\sqrt{N})$ sensors out of N available sensors. Let $R_{\mathcal{S}} \in \mathbb{C}^{M \times M}$ be the correlation matrix computed by aggregating the signals from these sensors. Then GNS ensures that $R_{\mathcal{S}}$ is a *lossless* sketch of the high-dimensional correlation matrix R_T . **Piya :** To Add (i) Finite sample performance guarantees (ii) two and three-dimensional extension (iii) low-rank extension and (iv) Time-varying model.

2. *Beyond Point Target Localization: Using Priors and Sparsity* In many applications such as camera networks, the quantities of interest are not the low-level measurements acquired at the CCD sensors, but the processed images I_t . In such cases, we need to obtain a compressive sketch of the image $A(I_t)$ via the sketching operator $A(\cdot)$ using low dimensional representation (over unions of subspaces or manifolds). In addition to conventional sparsity and low-rank priors, one can also utilize (partial) knowledge of the prior distribution of the images $I_t \sim \mathcal{D}$. Utilizing these priors can lead to more effective compression for a given level of sparsity. [To be written..]

Piya : These tasks can be further integrated with the binary embedding based sketching ideas proposed by Rayan and Alex. **Rayan :** Agree!

5.1 Graph signal processing approach without a fusion center

Yoav : I think this section can be combined with Piya's sections. Choosing which pairs should communicate is clearly related to their geometric layout.

Here the processing is done locally at each node. A graph signal processing approach was used in Ref [26] for a 5000 element seismic array by processing the whole normalized SCM at once, i.e., using a fusion center. When the coherence $\hat{\mathbf{S}}$ is above a certain threshold at element ij it is likely that a signal is observed and has propagated between nodes i and j , essentially forming an edge between nodes i and j in a graph. When a sufficient set of connected edges are detected in a region of the network a source is likely in that region. Part of the extracted SCM can then be used to localize the source more precisely.

To extract very weak signals with a well estimated and robust SCM is needed. Thus we pass the full time series between local nodes i and j , not the whole array and develop robust signal processing methods[32]. This will represent a lot of communication demand and thus we will only pass signal between neighboring stations. Once a graph edge is formed it could either be communicated further to a wider set of nodes. Robust signal processing methods[32] would entail making the processing insensitive to outliers. Qualitative robustness can be investigated via the influence function. A qualitatively robust estimator is characterized by an IF that is continuous and bounded. Continuity implies that small changes in the observed sample cause only small changes in the estimate. The boundedness implies that a small amount of contamination cannot lead to an unbounded error in the estimate.

5.2 Tomography

Yoav : I believe you are talking here about tomography, or reconstructing the environment. Can you write a paragraph of introduction, what is the problem? What is the desired solution? Sparse modeling assumes that signals can be reconstructed using a few (sparse) vectors, called atoms, from a potentially large set of atoms, called a dictionary. Recent ocean acoustics works utilizing sparse modeling is beamforming[31], matched field processing [9], and geoacoustic inversion [10]. One challenge in sparse modeling is finding the best dictionary for sparsely representing specific signals. Such dictionaries can be composed of wavelets, or the discrete cosine transform (DCT). These predefined dictionaries perform well for many signals. However, using a form of unsupervised machine learning, called dictionary learning, optimal dictionaries can be learned directly from specific data[21]. It has been shown that learned dictionaries outperform generic dictionaries when sufficient signal examples are available. Machine learning, and specifically dictionary learning, have recently obtained compelling results in ocean acoustics citeBianco2017 and seismology[19]. **Rayan :** I might be able to throw in some text about learning fast dictionaries, i.e., dictionaries learned from the data, but that also admit fast transforms like FFTs. Peter, Yoav, what do you think?

Yoav : I don't understand the following three paragraphs, can you give some technical details? Formulas? In current work, we have developed a machine learning-based travel time tomography method called locally sparse travel time tomography (LST)[1]. In LST, small scale local features contained in small rectangular groups of pixels, called patches, in an overall slowness (inverse speed) map are constrained using a sparse

model. Further, the sparsifying dictionary is adapted to the specific slowness data using dictionary learning. Larger scale, or global features spanning the map, are constrained with least-squares regularization. Unlike conventional tomography, in which model features are forced to be exclusively smooth or discontinuous, the LST approach permits smooth and discontinuous local features via dictionary learning.

Whereas many machine learning techniques in geoscience[19], are reliant on large amounts of training data, LST requires none. In LST we adopt the adaptive dictionary learning paradigm from image denoising [7] and medical imaging[25], in which dictionaries are learned directly from patches of the corrupted image. In LST, slowness dictionaries are learned from patches of a least squares regularized inversion, and are then used to reconstruct a sparsity-constrained slowness image. Assuming sufficiently dense ray sampling, the dictionary is initially unknown and is learned in parallel with the inversion. LST obtains high resolution by assuming that small patches of discrete slowness maps are repetitions of few elemental patterns from a dictionary of patterns. These patterns, which are described by the atoms in the dictionaries, are extracted from the data by dictionary learning. The increase in performance for synthetic slownesses relative to competing methods, are demonstrated for ambient noise tomography[1].

Assuming that the travel paths between sensors has been estimated[30, 11] We here propose the future development of machine learning-based tomography methods in ocean acoustics. Such methods will help to more fully-exploit both existing hydrophone and environmental data, as well as very dense sampling from future arrays with many sensors. Such large scale, mobile, and deformable arrays, will use ambient noise processing [30], to obtain very dense and rich data sets. We propose for future work to: (1) further develop a dictionary learning-based travel time tomography [1], accounting for uncertainty in the measurements and physics; (2) formulate the dictionary learning-based approach as CNN via CSC; and (3) apply this CSC tomography framework to data assimilation, to obtain higher-resolution estimates of water column parameters over conventional methods. We further propose to develop (5) acoustic event detection methods that leverage recent advances in machine learning.

6 Statistics on Binarized Sketching of Complex Measurements

Rayan : under construction - need a better title!

Consider an application where a sensor, or a network of sensors monitoring a system aims to detect whether a major change in the status of the system has occurred. In such a setting sensors may collect large amounts of data in an ongoing way. However, communicating all this data across all the sensors to perform this statistical change-detection task may be prohibitively expensive and tremendously wasteful. Similarly, consider the (e.g., anomaly detection) related setting where one may wish to determine whether data being collected at two sensors, or at two sets of sensors is similar (i.e., drawn from the same distribution), without communicating all the data across sensors. The ideal solutions for such tasks would be having each sensor maintain and communicate only one vector whose entries are of low bit-depth, i.e., taking on one of very few values, and then perform a simple computation on the communicated vectors to solve the problem. One goal of this proposal is to devise and study methods for obtaining such low bit-depth sketches in the context of performing statistical tasks, and to study their computational complexity and performance guarantees. To that end, we will first provide some relevant background, both on kernel and compressive statistics [12, 13] as well as quantization and binary embeddings [17, 29, 28, 14].

6.1 Background and Prior Work: Quantization and Binary Embeddings

Rayan : under construction: need to shorten We will rely on extremely coarse, e.g., binary, quantization of the data-sketches. In addition to minimizing storage and computation costs (see e.g., [8, 2]) these methods have the added benefit of being appealing in hardware implementation particularly if they are computationally inexpensive and thus promote speed in hardware devices [17, 20]. A growing body of work, which co-PI Saab has contributed significantly to (e.g., [29, 18, 28, 27, 15, 6]), has focused on signal reconstruction from coarsely quantized measurements. *One important theme that emerges from this line of work is that if one collects more coarsely quantized (even 1-bit) measurements than a critical minimal number, and uses sophisticated quantization schemes, then the extra measurements can be efficiently used in quantization-aware algorithms to rapidly drive the reconstruction error down as a function of the number of measurements.* This theme has held true in a wide range of signal and measurement contexts, including bandlimited functions

[6], finite frame expansions [16], and compressed sensing of approximately sparse vectors [29, 28], low-rank matrices [27], and manifold valued signals [15].

Co-PI Saab has recently extended this observation beyond signal reconstruction, to the context of (Euclidean) distance-preserving binary embeddings [14]. Here the goal is to map points in \mathbb{R}^n to the binary cube $\{\pm 1\}^m$, where $m \ll n$, in such a way that pairwise distances in \mathbb{R}^n can be well-approximated by appropriate functions on $\{\pm 1\}^m \times \{\pm 1\}^m$. To be precise, we now briefly describe this contribution as it is pertinent to our ensuing discussion. In [14], $A : \mathcal{T} \subset \mathbb{R}^n \rightarrow \mathbb{R}^m$ is a *random* Johnson-Lindenstrauss (i.e., distance preserving) map and $\mathcal{T} \subset \mathbb{R}^n$ is a set of finite or infinite cardinality. Moreover, $Q : \mathbb{R}^m \rightarrow \{\pm 1\}^m$ is a noise-shaping quantizer (e.g., a $\Sigma\Delta$ [5], or β [4] quantizer). In their most commonly used form, noise-shaping quantizers act sequentially on the measurements, say \mathbf{y}_i . **Alex : May be able to remove some of the description of 1st order and only focus on general, if need to reduce space.** For example, the simplest noise-shaping quantization scheme, the so-called greedy 1st order $\Sigma\Delta$ scheme, is given, for $i = 1, \dots, m$, by

$$\mathbf{q}_i = \text{sign}(\mathbf{y}_i + \mathbf{u}_{i-1}) \quad (10)$$

$$\mathbf{u}_i = \mathbf{u}_{i-1} + \mathbf{y}_i - \mathbf{q}_i, \quad (11)$$

where the state-variable sequence \mathbf{u}_i is initialized via, e.g., $\mathbf{u}_0 = 0$. In matrix-vector notation, this yields

$$\mathbf{q} = \mathbf{y} - D\mathbf{u} \quad (12)$$

relating the quantization to the measurement vector and the state variables, with D being the $m \times m$ first-order difference matrix. Crucially both for the analysis and for practical implementation, this scheme is *stable*, that is, for dimension independent constants c_1, c_2

$$\|\mathbf{y}\|_\infty \leq c_1 \implies \|\mathbf{u}\|_\infty \leq c_2, \quad (13)$$

More generally, other *stable* noise-shaping quantizers act sequentially on the measurements \mathbf{y}_i , and yield

$$\mathbf{q} = \mathbf{y} - H\mathbf{u},$$

where H is a lower-triangular matrix associated with the scheme. Quantizers of interest include stable r th-order (with $r \geq 1$) $\Sigma\Delta$ schemes [5] where $\mathbf{q} = \mathbf{y} - D^r\mathbf{u}$, and distributed- β encoding schemes [4] where H is block diagonal with identical lower-triangular blocks G , given by $G_{i,i} = 1, G_{i+1,i} = -\beta$, and $G_{i,j} = 0$ otherwise.

With these quantization schemes playing a prominent role, in [14] co-PI Saab constructed approximately isometric (i.e., distance preserving) embeddings between the metric space $(\mathcal{T}, \|\cdot\|_2)$ and the binary cube $\{-1, +1\}^m$ endowed with the pseudometric

$$d_V(\tilde{\mathbf{q}}, \mathbf{q}) := \|V(\tilde{\mathbf{q}} - \mathbf{q})\|_2$$

where V is a carefully constructed matrix. For a matrix $A \in \mathbb{R}^{m \times n}$, and a noise-shaping quantizer Q , as above, the algorithm for computing these embeddings is simply given by

$$\begin{aligned} g : \mathcal{T} &\rightarrow \{\pm 1\}^m \\ x &\mapsto \mathbf{q} = Q(A\mathbf{x}). \end{aligned}$$

In particular, when A is a fast Johnson-Lindenstrauss matrix (e.g., []), the constructed embeddings support fast computation. Despite their highly quantized non-linear nature, these binary embeddings perform as well as linear Johnson-Lindenstrauss methods! Indeed, when \mathcal{T} is finite [14] shows that with high probability and for prescribed distortion α

$$m \gtrsim \frac{\log(|\mathcal{T}|) \log^4 n}{\alpha^2} \implies |d_{\tilde{V}}(g(\mathbf{x}), g(\tilde{\mathbf{x}})) - \|\mathbf{x} - \tilde{\mathbf{x}}\|_2| \leq \alpha \|\mathbf{x} - \tilde{\mathbf{x}}\|_2 + c\eta(m),$$

where $\eta(m) \xrightarrow{m \rightarrow \infty} 0$.

Above, $\eta(m)$ decays polynomially fast in m (when \mathcal{Q} is a $\Sigma\Delta$ quantizer), or exponentially fast (when \mathcal{Q} is a distributed noise shaping quantizer). Additionally, when \mathcal{T} is arbitrary (with possibly infinite cardinality, e.g., a compact manifold) [14] show that with high probability and for prescribed distortion α

$$m \gtrsim \frac{\log^4 n}{\alpha^2} \cdot \frac{\omega(\mathcal{T})^2}{\mathcal{R}(\mathcal{T})^2} \implies |d_{\tilde{\mathcal{V}}}(g(\mathbf{x}), g(\tilde{\mathbf{x}})) - \|\mathbf{x} - \tilde{\mathbf{x}}\|_2| \leq \max(\sqrt{\alpha}, \alpha) \mathcal{R}(\mathcal{T}) + c\eta(m)$$

where $\eta(m)$ is as before and where $\mathcal{R}(\mathcal{T})$ and $\omega(\mathcal{T})$ denote the Euclidean radius of \mathcal{T} and its Gaussian width (which roughly scales with the average radius of \mathcal{T} , so that intrinsically low-dimensional sets have a small Gaussian width).

In short, with very few measurements compared to the ambient dimension of the signals, one can very efficiently (roughly at the cost of a fast Fourier transform) obtain low-dimensional binary sketches of the data. These sketches approximately preserve all pairwise distances in the original set and the distances in the embedded space can be computed efficiently.

While signal reconstruction is not a major focus of this section, the promising results obtained in that context, and also in the context of binary embeddings, lead us to believe that the above techniques can be generalized to other tasks.

6.2 Background and Prior Work: Kernel Statistics

Statistical distances, or accurately measuring distances between distributions, arise in a large number of applications. For sensors, these are important quantities for monitoring and tracking. An example of this would be for acoustic sensors, where each sensor collects the local power spectral density of a signal. Each sensor i now has a set of high dimensional data $\mathbf{Y}_i = \{\mathbf{y}_i(t)\}_{t=t_0}^{t_n} \subset \mathbb{R}^d$, and the question is whether \mathbf{Y}_1 and \mathbf{Y}_2 are distributionally the same, up to a time shift. After constructing a kernel $K : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}_+^d$ to define similarity between any two points, and ignoring communication constraints, a simple statistic to construct for a point $z \in \mathbf{Y}_i$ is a difference of kernel means over the two data sets $\frac{1}{n} \sum_{x \in \mathbf{Y}_1} K(z, x) - \frac{1}{n} \sum_{x \in \mathbf{Y}_2} K(z, x)$. If \mathbf{Y}_1 and \mathbf{Y}_2 came from the same distribution, then this statistic would be unbiased at $z \in \mathbf{Y}_i$, otherwise there would be a bias for particular z . Thus, computing the mean square error over all $z \in \mathbf{Y}_1 \cup \mathbf{Y}_2$ yields a statistic that would be close to 0 if $\mathbf{Y}_1, \mathbf{Y}_2 \sim p$, and would be biased if $\mathbf{Y}_1 \sim p$ and $\mathbf{Y}_2 \sim q$ for $p \neq q$. This is a well studied statistic known as kernel Maximum Mean Discrepancy [12]. In what follows, we will describe the mathematical framework and guarantees established, and propose methods for dealing with communication and computation constraints in this framework through randomized sketching and binary embeddings, as well as how to establish guarantees of convergence under non-i.i.d. sampling situations such as time series data. We will also detail a larger set of sensor problems that can be addressed in this framework.

Generally, the approach described above is measuring the distance between the distributions' *mean embeddings* [23]. A mean embedding of a distribution $\mu_p : \mathbb{R}^d \rightarrow \mathcal{H}$ of a probability distribution p is computed as

$$\mu_p(z) := \mathbb{E}_{x \sim p}[K(z, x)].$$

Effectively, mean embedding is mapping the distribution p to a point in the Reproducing Kernel Hilbert Space \mathcal{H} that is induced by the kernel K .

When K is a *universal* kernel (e.g. Gaussian, linear correlation) [22], then the mean embedding satisfies a key property that $\|\mu_p - \mu_q\|_{\mathcal{H}}$ is bi-Lipschitz with respect to $\|p - q\|_{L^\infty}$ for absolutely continuous distributions p and q . This effectively means that the mean embedding transform maintains the same information as working in \mathbb{R}^d , with the benefit that mean embeddings also satisfy nice statistical convergence properties. One key property is that, if we are only given n finite samples $\mathbf{Y}_1 \sim p$ and $\mathbf{Y}_2 \sim q$ to compute the empirical mean embeddings $\hat{\mu}_X$ and $\hat{\mu}_Y$, and we compute the mean embedding at all $z \in \mathbf{Y}_1 \cup \mathbf{Y}_2$, then $\|\hat{\mu}_{\mathbf{Y}_1} - \hat{\mu}_{\mathbf{Y}_2}\|_2 \rightarrow \|\mu_p - \mu_q\|$ at a rate $O\left(\frac{1}{\sqrt{n}}\right)$. A statistical interpretation of the mean embedding distance is that it computes a mean shift alternative test on the eigenfunctions of K rather than in the original space \mathbb{R}^d , and that two distributions having matching means in the eigenfunction space is equivalent to the distributions having all moments matching in \mathbb{R}^d . A large benefit of the mean embeddings is that this calculation can be done without explicitly computing the eigendecomposition of K .

However, statistics of this type suffer from a number of issues under computation and communication bottlenecks, as they require storing and communicating all points in $\mathbf{Y}_1 \cup \mathbf{Y}_2$. In particular, computing $\hat{\mu}(z)$ at any one point z requires evaluating the kernel at all points in $\mathbf{Y}_1 \cup \mathbf{Y}_2$, which can be prohibitively expensive. This has motivated computational speed ups presented by co-PI Cloninger in [3] of undersampling $z \in S \subset \mathbf{Y}_1 \cup \mathbf{Y}_2$ under the condition that a kernel matrix K can be decomposed as $K \approx RR^T$ for R that can be efficiently accessed. There is a similar vein of computation speed up accomplished through kernel compression via randomized sketching [13]. However, even these approaches still require communication of a number of points between sensors, or passing double precision complex valued summary statistics.

6.3 Framework

Rayan : I like this paragraph :), we should find a home for it. Our goal is to present a complete, theoretically rigorous, framework for performing various statistical, signal processing, and learning tasks from highly quantized data representations. We focus on the case of 1-bit representations as a theoretical extreme case, but emphasize that our methods should apply to more finely quantized data. We will propose, and analyze (a) algorithms for producing quantized sketches of data as well as (b) associated algorithms for performing the afore-mentioned tasks. Our strategy will be to develop these algorithms in tandem; that is, we will propose task-based quantization algorithms and quantization-aware algorithms for performing the tasks. We will strive for methods that support fast computation, and that lend themselves to distributed computing on, e.g., a sensor network.

Suppose signals of interest are represented as N points $x[j] \in \mathbf{X} \subset \mathbb{R}^d$ for $j \in \{1, \dots, N\}$, where \mathbf{X} could be a finite set, i.e., a point cloud, or an infinite set (e.g., a compact manifold, or the set of all sparse vectors). If the data is collected as a time series, then $x[j] = x(t_j)$. Further, assume that the measurement operator, accounting for all the digitization and measurements at all the sensors, is given by $\Phi : \mathbb{R}^d \rightarrow \mathbb{R}^m$. In the case of distributed sensing systems, we can assume each sensor collects/computes a portion, $\Phi_i \mathbf{X} \in \mathbb{R}^{m_i}$, of $\Phi \mathbf{X} \in \mathbb{R}^m$ (so that $\sum_i m_i = m$). The corresponding signal at sensor i , $Y_i = \Phi_i \mathbf{X}$ can then be quantized via $Q(\Phi_i \mathbf{X})$ with a quantization map $Q : \mathbb{R}^{m_i} \rightarrow \{\pm 1\}^{q_i}$ that maps the measurements to bits. The full quantization map of all sensors can be represented as

$$x \mapsto Q(\Phi x) = \begin{bmatrix} Q(\Phi_1 x) \\ \vdots \\ Q(\Phi_k x) \end{bmatrix} \mapsto f(Q(\Phi x)) \approx g(x) \quad (14)$$

6.4 Proposed Work: Binary Embeddings of Nonlinear Similarity Measures

The type of information that can be derived from $f(Q(\Phi x))$ depends on the choice of measurement operator Φ_i at each sensor, as well as the quantization scheme. In [14], Co-PI Saab constructed the quantization scheme f , Φ , and Q so that $\|g(x) - g(y)\| \approx \|x - y\|_2$. For several of the statistical tasks below, we require construction of the quantization scheme to approximate more complex relationships between x, y . We here propose the development of a new set quantization schemes that incorporate nonlinear measurements of the data x . In particular, we aim to design f , Φ , and Q so that $\langle f(Q(\Phi x)), f(Q(\Phi y)) \rangle$ approximates the Gaussian kernel $\langle g(x), g(y) \rangle = e^{-\|x-y\|_2^2/\sigma^2}$. As a first step, we will show that this can be achieved by appropriately quantizing $\cos(2\pi\langle w, \cdot \rangle)$ for random choices of w [24]. To achieve this in practice, each physical sensor could collect one or more linear measurements $\langle \omega, \cdot \rangle$, and then the non-linearity (e.g., cosine) could be implemented as part of the quantizer. Alternatively, the sensor could measure the non-linear function directly prior to quantization.

Rayan : under heavy construction Theoretical considerations: Random Fourier Features (RFFs) were first introduced by Rahimi and Recht [?] to approximate shift-invariant kernels (e.g., Gaussian kernels $k(x, y) = \exp(-\|x - y\|_2^2/2)$). They considered $\Phi(x) = \sqrt{\frac{2}{m}} \cos(Ax + b)$, where A is a Gaussian matrix and b is a uniform random vector over $[0, 2\pi]$ and showed a concentration of measure [] result for

$$\sup_{x, y \in \mathcal{X}} |\langle \Phi(x), \Phi(y) \rangle - k(x, y)| = \sup_{x, y \in \mathcal{X}} \left| \frac{2}{m} \sum \cos(\langle a_i, x \rangle + b_i) \cos(\langle a_i, y \rangle + b_i) - k(x, y) \right|$$

by using an ϵ -net argument. The result was later improved by Sriperumbudur and Szabo [?] using Dudley’s inequality []. Still, the above results’ dependence on the data’s geometric properties (e.g., its intrinsic dimension) was sub-optimal and the application of Gaussian matrices to large data sets could be prohibitively expensive as they do not admit fast transforms. To address the second issue, inspired by [?], Le, Sarlos, and Smola [?, ?] considered a structured random matrix A to speed up RFFs. The structured matrix they proposed satisfied $A = SHG\Pi HB$, where Π is a permutation matrix, H is the Hadamard matrix, and S, G, B are diagonal random matrices. With this construction [?] claimed a concentration result analogous to that of [?] using an ϵ -net argument, and they *conjectured* that H can be replaced by any orthogonal matrix and reported numerical evidence showing that random partial Fourier matrices also gave a comparable results. There is not much theoretical work for this setting.

Questions. This is a great example of nonlinear + structured random measurements. We need to understand the empirical process:

$$\mathbb{P} \left(\sup_{(x,y) \in \mathcal{X} \times \mathcal{Y}} \left| \frac{2}{m} \sum \cos(Ax) \cos(Ay) - k(x,y) \right| \geq t \right),$$

where A is a structured random matrix.

In general, we can study

$$\mathbb{P} \left(\sup_{(f,g) \in \mathcal{F} \times \mathcal{G}} \left| \frac{2}{m} \sum f(a_i)g(a_i) - \mathbb{E}[fg] \right| \geq t \right). \quad (15)$$

The empirical process (15) has been studied by Mendelson [?] when f and g are sub-gaussian, using the generic chaining method.

Understanding (15) is also helpful for us when we incorporate quantization.

Quantize RFFs The first work was introduced by Raginsky and Lazebnik [?] in which they quantized RFFs $Q(\Phi(x)) = Q(\sqrt{\frac{2}{m}} \cos(Ax + b))$ where Q is a dither quantization. They showed that we could quantize and use hamming distance of these quantized data to approximate a (complicated) function of the Gaussian kernel.

Questions. Consider noise-shaping quantization methods for RFFs. Then we need to understand (15). Note that there is also recent work by Re and co-authors [?]

6.5 Proposed Work: Quantized Statistics

Our overarching goal is to reduce the communication constraints by utilizing both sketching and binarization through use of $g(x) = f(Q(\Phi(x)))$. The idea is that, rather than communicate and compute with $x \in \mathbb{R}^d$ where d may be large, one only needs to compute $Q(\Phi x) \in \{\pm 1\}^m$ with $m \ll d$, and communicate $f(Q(\Phi x))$. In what follows, $\mathbb{E}_{x \in X} g(x)$ serves as a proxy for the empirical mean embedding of the data set X . For sensor networks, we can either compare between sensors by embedding each sensors’ data through $g(\mathbf{Y}_i) = f(Q(\Phi_i \mathbf{X}))$ when Φ_i at each sensor are constructed using the same random measurement operator (though there may be differences due to sensor location, obstructions, etc). A separate comparison is to have the sensors work collaboratively across a common set of points by building a concatenated sketching from data across all sensors $g(\mathbf{Y}) = [f(Q(\Phi_1(\mathbf{X}))), \dots, f(Q(\Phi_k(\mathbf{X})))]$. We aim to prove that this low complexity vector still converges to a type of mean embedding, and to provide a rigorous analysis of the statistical power, convergence rates, and minimal detectable separation criteria between the distributions. Below we highlight the benefit of this approach in a number of different sensor problems.

- **Two Sample Testing:** In the context of sensors, the two sample problem can be summarized as follows: each sensor collects a data set $\mathbf{Y}_i \sim p_i$, and the goal is to determine whether the \mathbf{Y}_i were distributed similarly. To address this, we will analyze the binarized statistic

$$\left\| \frac{1}{n} \sum_{x \in \mathbf{Y}_1} g(x) - \frac{1}{m} \sum_{y \in \mathbf{Y}_2} g(y) \right\|,$$

under appropriate norm, and seek to characterize the minimal conditions under which a deviation between p_1 and p_2 can be detected. The approach requires characterizing the types of deviations $p_1 - p_2$ that can

be detected, namely those for which

$$\left\| \int e^{-\|x-y\|^2/\sigma^2} (p_1(y) - p_2(y)) dy \right\| > \epsilon,$$

as well as the rates at which these deviations can be detected. The communication benefit of such a statistic is that the sensors need only transmit the mean of $g(x)$, rather than all the individual points.

Alex : Could expand on this, or reference other comments in proposal and a few papers about this for JL embeddings

Rayan : let's refine this discussion a little – I have some ideas here

- **Change Point Detection:** A variant of the two sample testing problem is change point detection, in which the data is streaming $X(t)$ according to some underlying stochastic process. At some time t^* , the distribution of the stochastic process changes from one distribution to another, and the issue is how quickly after t^* this change can be detected. Unlike the two sample context in which we were testing whether the sensors detected the same distribution, here we can use the sensors collaboratively by constructing the concatenated sketching matrix $g(\mathbf{Y})$. There exists a mean embedding approach to change point analysis [?] which uses the kernel Fisher discriminant ratio and mean embedding to measure the homogeneity between time segments of the process. However, this once again requires storage of all points over the length of the detection window and communication of those points across sensors for both computing the window mean and variance. It also suffers from an inability to begin computing the change point statistic until all points in the window have been collected. We aim to introduce the binarized sketching framework to produce an efficient computation to the kernel mean as in two sample testing, as well as a fast construction of a low rank approximation to the kernel covariance matrix. We will derive it's new limiting distribution under the null model of no change, as well as the consistency under the alternative distribution when a change does occur.

- **Multiple Sensor Common Factor Identification:** A common issue is aggregating multiple sensors to identify and magnify the signal detected by both sensors. Under a linear model, algorithms such as Canonical Correlation Analysis [?] act on multiple streams of simultaneously collected data $\mathbf{y}_i(t)$ to filter noise and recover highly correlated linear projections from two the data sets. Recently, a kernel CCA technique called alternating diffusion [?] has been used to identify common nonlinear effects by building a kernel K_i from each sensing modality and analyzing the product kernel $(K_1 K_2)^t$. Concretely, assume that there exists a hidden manifold \mathcal{M} and two nuisance manifolds \mathcal{N}_1 and \mathcal{N}_2 , and samples $s((x_i, y_i^{(1)}, y_i^{(2)}))$ for $(x_i, y_i^{(1)}, y_i^{(2)}) \in \mathcal{M} \times \mathcal{N}_1 \times \mathcal{N}_2$. Due to sensor location or modality, sensor 1 collects data points $S_1 = s((x_i, y_i^{(1)}, \xi))$ and sensor 2 collecting data points $S_2 = s((x_i, \zeta, y_i^{(2)}))$. Then roughly speaking, for some assumptions on s and for $K_1 : S_1 \times S_1 \rightarrow \mathbb{R}$ and $K_2 : S_2 \times S_2 \rightarrow \mathbb{R}$, Talmon and Wu [?] proved $(K_1 K_2)^t$ is the diffusion kernel on \mathcal{M} only. This is a very beneficial feature, as it means that one can compute kernel statistics for distributions defined on \mathcal{M} only, independent of the nuisance features that may differ between sensors due only to modality or location.

However, this requires communication of all n points across sensors to compute the kernel product. We propose to analyze such approaches under minimal communication constraints by utilizing the low rank binarized decomposition $K_i \approx g(S_i)g(S_i)^*$. As the key feature of alternating diffusion is computing the inner product matrix $g(S_1)^*g(S_2)$, which has dimension $d < n$ and is also low rank. This implies it is possible to sketch $g(S_1)$ (resp. $g(S_2)$) with a small set of vectors v_j (resp. w_k) and only communicate vectors $g(S_1)^*v_j$ (resp. $w_k^*g(S_2)$) such that $\mathbb{E}_{j,k}[g(S_1)^*v_j w_k^*g(S_2)] \approx g(S_1)^*g(S_2)$. The amplification of the common factors observed by both sensors can serve to boost the power of the two sample and change point statistics on the shared observable manifold \mathcal{M} , as well as other compressed statistics that can be computed [13].

- **Rayan :** Maybe put in classification

– **Idea 1:** Training reduced to computing averages μ_j , over the data of say $Q(\cos(Ax + b))$ for each class $x \sim \mathcal{X}_j$. Classification reduced to computing $\mu_y = Q(\cos(Ay + b))$ and finding the closest μ_j to μ_y . In a sense, this is like 2-sample testing, but with a point mass at y . Advantages: Each sensor only keeps the

averages relevant to its own portion of $Q(\cos(Ax))$ for x in training set, and later for y to be classified. Each sensor can compute its own piece of the inner product $\langle \mu_j, \mu_y \rangle$, you are ultimately averaging bits, so memory needed at each sensor grows like $m_i \times \log(P)$ where P is the total number of points and m_i is the size of the sketch at each sensor... The idea is that $m_i \ll d$ and $\log(P) \ll P$ so you save a ton on storage/communication.

- Idea 2: Exploit the hierarchical/distributed nature of the bits produced by binary embeddings. Requires more thought... and clustering: Analogous to the above...

7 Tell me something new

Consider a video system for monitoring a highway. Lets say that the highway is 100 miles long and that we have a 5 cameras per mile, i.e. 500 cameras. Each camera generates a large raw data stream, on the order of 100MB per second. Using video compression this stream can be brought down to XXMB/sec. Totalling XXXMB/sec. A commong approach is to send all of this data to a centralized site where a team of analysts monitor it to identify accidents, road hazards etc.

Several issues should be noted:

- A high bandwidth channel is needed to transmit the information to the center. If this channel is wireless, this transmission will be costly in terms of energy and available bandwidth.
- Manually monitoring 500 video channels a significant task for the analysts. AI methods can be used to aid in that task, our goal is to perform this AI on the camera, rather than in the centralized site.
- Under normal conditions, highway traffic is highly predictable. Only unpredictable situations such as an accident, an object on the road or a speeding car warrant the transmission of high-resolution images.

We propose a general framework for designing sensor networks for anomaly detection called “tell me something new” [?]. In this approach, sensors analyze the data they collect, and send message to the center only when it estimates that this information is significantly different than what the center will predict without the information.

Let $\hat{\theta}_1(t), \hat{\theta}_2(t)$ be the estimates of the target location for each sensor. *In addition, each sensor maintains an estimate of the estimate of the other sensor.* $\hat{\theta}_{1,2}(t), \hat{\theta}_{2,1}$. Each sensor updates it’s estimate of the velocity of the location of the target according to the signal it measures, but it does not update it’s estimate of the other’s estimate. If the two estimate are close to each other $\hat{\theta}_i(t) \approx \hat{\theta}_{i,j}(t)$ then sensor i sends no information to sensor j . On the other hand, if $\hat{\theta}_i(t)$ is far from $\hat{\theta}_{i,j}(t)$, then $\hat{\theta}_i$ is transmitted from sensor i to sensor j . Thus if the target is moving in constant speed, uninterrupted, there is not communication between the sensors.

The basic idea here is that a sensor sends out information only if that information cannot be predicted by the receiver. Similar ideas have been used in arithmetic coding and **Yoav : I think** in $\Sigma\Delta$ encoding.

Recently, PI Freund [?] proposed an asynchronous computation model called “Tell Me Something New” in which each agent broadcasts a message only when the estimate it computes differs from the existing estimates in a statistically significant way.

One important application of sensor networks is to monitor activity and identify anomalies. Examples include: building security systems, factory floors, highway monitoring, health monitoring for the sick or elderly and many others.

On its face, this might seem like an under-constrained impossible problem. However, note that for all of the environments listed above there is a highly repetitive pattern from day to day and from week to week. Add to that the sensors are stationary, and one would expect that most sensors observe highly regular and highly predictable patterns.

The approach we propose in this case is that each sensor creates a model of the characteristics of the signals that it observes during normal operations. It alerts neighboring sensors if it observes something that is abnormal, i.e. a signal that has very low probability according to the model. When several sensors send an alert with a short time window, and when the alerts are consistent with each other, a global alert is sent to the human operators.

8 results from previous grants

8.1 Sampling and quantization theorems for modern data acquisition

Co-PI Saab: *Sampling and quantization theorems for modern data acquisition (DMS-1517204, 08/01/2015–07/31/19, \$160,404)* This grant resulted in 9 published or accepted journal articles [?, ?, ?, ?, ?, ?, ?, ?, ?], and 4 conference papers [?, ?, ?, ?]. **Intellectual Merit:** [?, ?, ?, ?, ?] study quantization of compressed sensing (CS) measurements, dealing with 1-bit scalar quantization, compression of bit-streams, circulant measurement matrices, and low-rank matrices, and manifold-valued signals. [?] develops a state-of-the-art technique, and accompanying theory, for embedding datasets into the binary cube while preserving Euclidean distances. [?] provides a framework for using 1-bit measurements for classification and analyzes the case of two clusters. The papers [?, ?] show that using prior support information and weighted ℓ_1 minimization yield improved recovery from fewer CS measurements. Finally, [?, ?] study phase retrieval from local measurements. The works use and develop tools in random matrix theory, mathematical signal processing, and applied harmonic analysis. **Broader Impacts:** In addition to dissemination of results through multiple seminars and workshop talks, Saab has developed and taught graduate courses on (1) compressed sensing and its applications, (2) applied and computational harmonic analysis, and (3) mathematical methods in data science. He has mentored a UCSD undergraduate (S. Gagniere, currently at UCLA) in research on quantization of CS measurements, and UCSD graduate students (E. Lybrand, A. Nelson, B. Preskitt, J. Liang), and postdocs (A. Ma, T. Huynh), some of whom are co-authors on the works above, while others are taking part in ongoing work. Two of the above mentioned graduate students and postdocs are female, one is an air-force officer, and at least two are members of under-represented groups.

P. Gerstoft: 2014-2018; PLR-1246151 (SIO; 844k); Collaborative Research: Dynamic Response of the Ross Ice Shelf to Wave-induced Vibrations. **Intellectual Merit:** This project (DRIS) investigates the response of ice shelves to ocean wave forcing to (1) infer bulk elastic properties from signal propagation characteristics, and (2) to determine how IG wave and other gravity wave forcing propagates across the Ross Ice Shelf (RIS), and (3) to monitor seasonal variability of the RIS response and icequake activity. In 2014, 16 DRIS broadband stations were installed to acquire data for 2 years. Subsequently, 13 geodetic GPS stations were installed during the 2015 field season recovery of the first year’s seismic data. Joint processing of the 2016 seismic and GPS data are underway. **Broader Impacts:** Supported 2 postdocs, 2 graduate students. NSF Artists and Writers awardee Glenn McClure produced symphonic and choral works from the DRIS data. The Birch Aquarium at Scripps is developing an Antarctic exhibit featuring DRIS results. Publications: [Bromirski et al., 2015; Diez et al., 2015; Bromirski et al., 2017; Chen et al., 2018]. A project web site [Bromirski, 2014], suitable for informing the public, is maintained.

References

- [1] M. J. Bianco and P. Gerstoft. Travel time tomography with adaptive dictionaries. *IEEE Transactions on Computational Imaging*, 4(4):499–511, 2018.
- [2] P. T. Boufounos and R. G. Baraniuk. 1-bit compressive sensing. In *2008 42nd Annual Conference on Information Sciences and Systems*, pages 16–21. IEEE, 2008.
- [3] X. Cheng, A. Cloninger, and R. R. Coifman. Two-sample statistics based on anisotropic kernels. *arXiv preprint arXiv:1709.05006*, 2017.
- [4] E. Chou and C. S. Güntürk. Distributed noise-shaping quantization: I. beta duals of finite frames and near-optimal quantization of random measurements. *Constructive Approximation*, 44(1):1–22, 2016.
- [5] I. Daubechies and R. DeVore. Approximating a bandlimited function using very coarsely quantized data: A family of stable sigma-delta modulators of arbitrary order. *Annals of mathematics*, 158(2):679–710, 2003.
- [6] I. Daubechies and R. Saab. A deterministic analysis of decimation for sigma-delta quantization of bandlimited functions. *IEEE Signal Processing Letters*, 22(11):2093–2096, 2015.
- [7] M. Elad. *Sparse and redundant representations: from theory to applications in signal and image processing*. Springer Science & Business Media, 2010.
- [8] J. Fang, Y. Shen, H. Li, and Z. Ren. Sparse signal recovery from one-bit quantized data: An iterative reweighted algorithm. *Signal Processing*, 102:201–206, 2014.
- [9] K. L. Gemba, S. Nannuru, P. Gerstoft, and W. Hodgkiss. Multi-frequency sparse Bayesian learning for robust matched field processing. *J. Acoust. Soc. Am.*, 141(5):3411–3420, 2017.
- [10] P. Gerstoft, C. F. Mecklenbräuker, W. Seong, and M. Bianco. Introduction to compressive sensing in acoustics, 2018.
- [11] P. Gerstoft, K. G. Sabra, P. Roux, W. Kuperman, and M. C. Fehler. Green’s functions extraction and surface-wave tomography from microseisms in southern california. *Geophysics*, 71(4):SI23–SI31, 2006.
- [12] A. Gretton, K. M. Borgwardt, M. J. Rasch, B. Schölkopf, and A. Smola. A kernel two-sample test. *Journal of Machine Learning Research*, 13(Mar):723–773, 2012.
- [13] R. Gribonval, G. Blanchard, N. Keriven, and Y. Traonmilin. Compressive statistical learning with random feature moments. *arXiv preprint arXiv:1706.07180*, 2017.
- [14] Z. Harchaoui, E. Moulines, and F. R. Bach. Kernel change-point analysis. In *Advances in neural information processing systems*, pages 609–616, 2009.
- [15] D. R. Hardoon, S. Szedmak, and J. Shawe-Taylor. Canonical correlation analysis: An overview with application to learning methods. *Neural computation*, 16(12):2639–2664, 2004.
- [16] T. Huynh and R. Saab. Fast binary embeddings, and quantized compressed sensing with structured matrices. *arXiv preprint arXiv:1801.08639*, 2018.
- [17] M. Iwen, E. Lybrand, A. Nelson, and R. Saab. New algorithms and improved guarantees for one-bit compressed sensing on manifolds. *arXiv preprint arXiv:1902.03726*, 2019.
- [18] M. Iwen and R. Saab. Near-optimal encoding for sigma-delta quantization of finite frame expansions. *Journal of Fourier Analysis and Applications*, 19(6):1255–1273, 2013.

- [19] L. Jacques, J. N. Laska, P. T. Boufounos, and R. G. Baraniuk. Robust 1-bit compressive sensing via binary stable embeddings of sparse vectors. *IEEE Transactions on Information Theory*, 59(4):2082–2102, 2013.
- [20] K. Knudson, R. Saab, and R. Ward. One-bit compressive sensing with norm estimation. *IEEE Transactions on Information Theory*, 62(5):2748–2758, 2016.
- [21] Q. Kong, D. T. Trugman, Z. E. Ross, M. J. Bianco, B. J. Meade, and P. Gerstoft. Machine learning in seismology: Turning data into insights. *Seismological Research Letters*, 90(1):3–14, 2018.
- [22] B. Le, T. W. Rondeau, J. H. Reed, and C. W. Bostian. Analog-to-digital converters. *IEEE Signal Processing Magazine*, 22(6):69–77, 2005.
- [23] R. R. Lederman and R. Talmon. Learning the geometry of common latent variables using alternating-diffusion. *Applied and Computational Harmonic Analysis*, 44(3):509–536, 2018.
- [24] S. Mallat. *A wavelet tour of signal processing*. Elsevier, 1999.
- [25] C. A. Micchelli, Y. Xu, and H. Zhang. Universal kernels. *Journal of Machine Learning Research*, 7(Dec):2651–2667, 2006.
- [26] K. Muandet, K. Fukumizu, B. Sriperumbudur, B. Schölkopf, et al. Kernel mean embedding of distributions: A review and beyond. *Foundations and Trends® in Machine Learning*, 10(1-2):1–141, 2017.
- [27] A. Rahimi and B. Recht. Random features for large-scale kernel machines. In *Advances in neural information processing systems*, pages 1177–1184, 2008.
- [28] S. Ravishanker and Y. Bresler. Mr image reconstruction from highly undersampled k-space data by dictionary learning. *IEEE transactions on medical imaging*, 30(5):1028–1041, 2011.
- [29] N. Riahi and P. Gerstoft. Using graph clustering to locate sources within a dense sensor array. *Signal Processing*, 132:110–120, 2017.
- [30] R. Saab and E. Lybrand. Quantization for low-rank matrix recovery. *Information and Inference: A Journal of the IMA*, 8(1):161–180, 05 2018.
- [31] R. Saab, R. Wang, and Ö. Yilmaz. Quantization of compressive samples with stable and robust recovery. *Applied and Computational Harmonic Analysis*, 44(1):123–143, 2018.
- [32] R. Saab, R. Wang, and O. Yilmaz. From compressed sensing to compressed bit-streams: Practical encoders, tractable decoders. *IEEE Transactions on Information Theory*, 64(9):6098–6114, Sep. 2018.
- [33] K. G. Sabra, P. Gerstoft, P. Roux, W. Kuperman, and M. C. Fehler. Surface wave tomography from microseisms in southern california. *Geophysical Research Letters*, 32(14), 2005.
- [34] R. Talmon and H.-t. Wu. Latent common manifold learning with alternating diffusion: analysis and applications. *Applied and Computational Harmonic Analysis*, 2018.
- [35] A. Xenaki, P. Gerstoft, and K. Mosegaard. Compressive beamforming. *J. Acoust. Soc. Am.*, **136**(1):260–271, 2014.
- [36] A. M. Zoubir, V. Koivunen, E. Ollila, and M. Muma. *Robust Statistics for Signal Processing*. Cambridge University Press, 2018.

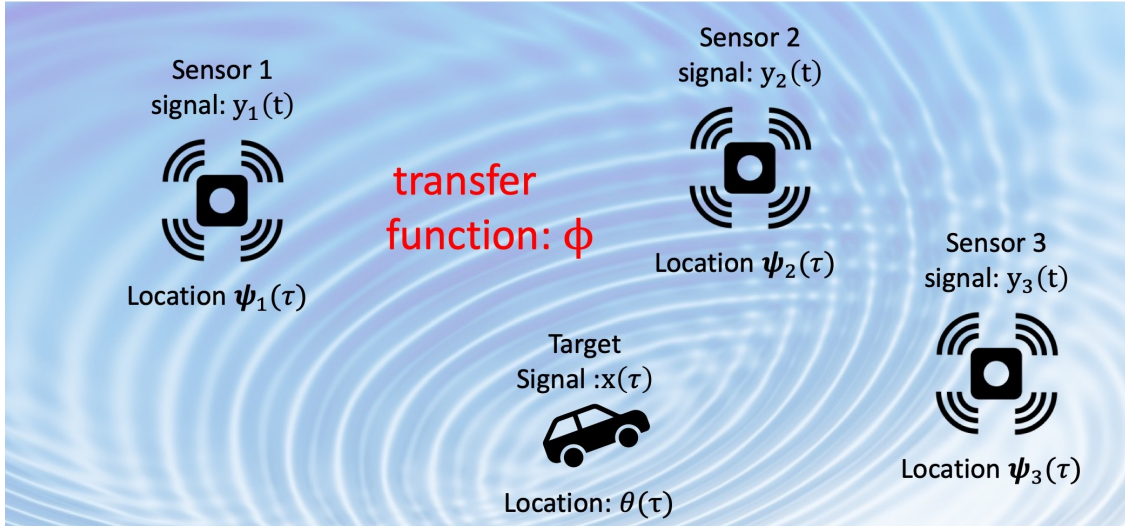


Figure 1: An example of a sensor network. The goal of the network is to track a car. The location of the car as a function of time is $\theta(\tau)$, and the goal of the network to produce $\hat{\theta}(\tau)$. The car emits a sound wave, which we denote by $x(\tau)$. The sound wave travels through the physical environment and arrives at each sensor i , where it is digitized and made into a time sequence $y_i(t)$. The transformation of the signal $x(\tau)$ into $y_i(t)$ is represented by a transfer function ϕ . In this car tracking example $y_i(t) = \phi(x(\tau), \theta(\tau), \psi_i(\tau))$ **Rayan** : ψ is not defined here. We are seeking an inverse transformation that would map the measurements vector $y_1(t), y_2(t), y_3(t)$ to an estimate of the car location $\hat{\theta}(t)$