# Some of My Group Current and Past Research

## Vahid Tarokh

# Current Research Projects

## 1. Current Research Projects:

### *Fundamental Data Science/Math Research*

1- Adaptive Exploitation of Non-Commutative Multimodal Information, Period of performance: 08/17/2015 - 08/16/2020

   **Short Description:** Shannon Information Theory is based on mutual information that is commutative. In practice information between two elements may be non-commutative. In this effort, we extend Shannon's Theory to non-commutative and free random variables.

2- Beyond Shannon Inspired Approach to Limits of Learning, Source and amount of funding: Period of performance: 04/01/18-03/31/19.

   **Short Description:** Shannon's Information Theory assumes knowledge of data generating processes. In contrast, here we are interested in determining how much can be learned from given data without the knowledge of data-generating processes.

3- Physics Inspired Score for Data-Driven Modeling and Prediction, Source and amount of funding: Period of performance:  01/02/17-01/01/21

   We will build a system to automatically model and discover the dynamics of data-generating process. To this end, inspired by physics, we will use Fisher's divergence to produce a new score (in contrast to the traditional log-likelihood) that can be computed and implemented in real time.

### *Brain-Computer Interfacing*

4- MURI: Closed-Loop Multi-sensory Brain-Computer Interface for Enhanced

Decision Accuracy), Period of performance:  6/1/16-5/31/21

**Short Description:**  In this project, we use local field potential (LFP) data collected from brains of Macaque monkeys for predicting their future actions, and human EEG data for Brain-Computer Interfacing.

## *Machine Learning*

5- DANCERS: Dynamic Adaptive Neurally-inspired Control for Efficient RF Surveillance, Period of performance: 07/01/18-03/31/2020.

**Short Description:** We will consider Radio Frequency Activity Data provided to us by DARPA. We will develop algorithms and systems in order to identify transmissions of interest.

6- Decentralized Perception From Online Learning and Semantic Understanding, Period of performance: 01/02/17-01/01/21

**Short Description:** We will develop methods for a group of distributed sensory devices to provide essential information to a decision-maker in order to produce perception of the observed scenario.

## *Optimization and Game Theory*

7- Stochastic Adaptive Game Analytic Multi-Player Optimal Resilient Execution (SAGAMORE), Period of performance: 10/09/18-06/08/2020

We use a Fisher Information Theoretic based approach to prediction of actions of an adversary in a stochastic game setting based on prior history and knowledge.

8- Optimal Online Data-Driven Optimization with Multiple Time-Varying Non-Convex Objectives, Period of performance: 04/1/18-06/07/2020.

**Short Description:** The major goals of this project are *data-driven optimization and control* solutions in distributed settings when the objective functions are time-varying, data generating function is unknown and network challenges between distributed agents may exist.

9- Approximate Computing on Real World Data Using Representation and Coding, Source and amount of funding: Period of performance: 08/1/16-07/31/19

The major goals of this project is to develop approximate computing techniques since the diminishing benefits from CMOS scaling has coincided with an overwhelming increase in the rate of data generation.

## Cybersecurity

Architectures for Security of Internet of Things, Source and amount of funding: DHS (Subcontract from XALGO5), Period of performance: 01/01/2018 - 01/31/2018

We study various methods of detecting malware based from network activity patterns in a hybrid network setting.

# Overview of Vahid Tarokh's Group Past Research

## 2. Past Research

## 2.1. Contributions to Communications:

### 2.1.1. Space-Time Codes:

One of the earliest contributions of our group is to multiple antenna transmission and reception techniques. I worked at AT&T Labs from 1996-2000 where I was fortunate to have worked with several truly outstanding colleagues on this topic. I co-invented and co-established the area of space-time coding by introducing the first codes (TCM-like). I also co-invented space-time block codes. These results have been incorporated into technologies such as LTE, WiMAX, HSPA, UMTS, IEEE 802.16e, and ANSI IS136. Finally, I came up with early papers on differential techniques in space and time, and combined-array processing and space-time coding methods.

Some representative publications are given below.

B.M. Hochwald, T.L. Marzetta, and V. Tarokh. 'Multi-antenna Channel Hardening and its Implication for Rate Feedback and Scheduling', *IEEE Transactions on Information Theory*, vol. 50, No. 9, pp. 1893-1909, Sep. 2004.

V. Tarokh and H. Jafarkhani, 'A Differential Detection Scheme for Transmit Diversity', *IEEE Journal on Selected Areas in Communications*, vol. 18, No. 7, pp. 1169-1174, Jul. 2000.

V. Tarokh, H. Jafarkhani, and A. R. Calderbank, 'Space-Time Block Coding from Orthogonal Designs', *IEEE Transactions on Information Theory*, vol. 45, No. 5, pp. 1456-1467, Jul. 1999.

V. Tarokh V, A. Naguib, N. Seshadri, and A.R. Calderbank, 'Combined Array Processing and Space-Time Coding', *IEEE Transactions on Information Theory*, vol. 45, No.4, pp.~1121-1128, May 1999.

V. Tarokh, N. Seshadri, and A.R. Calderbank, 'Space-time Codes for High Data Rate Wireless Communication: Performance Criterion and Code Construction', *IEEE Transactions on Information Theory*, vol. 44, No. 2, pp. 744-765, Mar. 1998.

### 2.1.2. Contributions to Multicarrier Communications and OFDM:

In 2000 and 2001, my colleagues and I focused on designing powerful codebooks that not only perform near theoretical limits, but their associated signals also have low peak-to-average power ratios in multicarrier communications. Multicarrier communications is the technology of choice for DSL, cable modems, and fourth generation cellular systems. If the peak-to-average power ratio is high, then for a given peak power (determined by the limitation of the power amplifiers or by regulatory issues), the average transmit power is smaller and the range is reduced. Alternatively, for a fixed average transmit power, high peak to average power ratios force power amplifiers with larger linear regions that are much more expensive. Together with my colleague K. Paterson, we

showed in 2000 and 2001 that there exist capacity achieving codes with low peak to average power ratios. I also devised with my colleague H. Jafarkhani methods to compute the peak to average power ratio.

In 2010, I returned to this problem and (jointly with students and postdocs Yongjun Kwak, Maryam Sabbaghian, and Besma Smida) constructed practical signaling schemes that achieve performance near Shannon capacity (comparable to turbo codes) but also have very low peak to average power ratios.

Some representative publications are given below.

V. Tarokh and H. Jafarkhani, "On the Computation and Reduction of the Peak to Average Power Ratio in Multicarrier Communications", *IEEE Transactions on Communications*, vol. 48, No. 1, pp. 37-44, Jan. 2000.

K.G. Paterson and V. Tarokh, "On the Existence and Construction of Good Codes With Low Peak-to-Average Power Ratios", *IEEE Transactions on Information Theory*, vol. 46, No. 6, pp. 1974-1987, Sept. 2000.

C. Roessing and V. Tarokh, ``A Construction of OFDM 16-QAM Sequences Having Low Peak Powers", *IEEE Transactions on Information Theory*, vol. 47, No. 5, pp. 2091-2094, Jul. 2001.

C.V. Chong and V. Tarokh, A Simple Encodable/Decodable OFDM QPSK Code with Low Peak-to-Mean Envelope Power Ratio", *IEEE Transactions on Information Theory*, vol. 47, No. 7, pp. 3025-3029, Nov. 2001.

M. Sabbaghian, Y. Kwak, B. Smida, V. Tarokh, "Near Shannon Limit and Low Peak to Average Power Ratio Turbo Block Coded OFDM," *IEEE Transactions on Communications*, August 2011.

## 2.1.3. Measurement and Modeling of UWB Channels:

Controversy surrounded the practical limitations of Ultra-wideband (UWB) channels around that time with UWB proponents seeming to promise huge data rates. This motivated our measurement campaign (with AT&T Labs Colleagues Saeed Ghassemzadeh, Robert Miller and M.Eng. student Thorvardur Sveinsson) where we produced the largest UWB channel measurement database available at that time. The UWB Path loss Channel Model (proposed by S. Ghassemzadeh and myself) has been adopted by the IEEE 802.15 Study Group for evaluation of vendor proposed ultrawideband files. [Standard document: S.S. Ghassemzadeh and V. Tarokh, "The Ultra Wideband indoor path loss model," *Tech. Report P802.15 02/277r1SG3a, June 2002*]. In addition, our results were used to study coexistence issues between UWB and Wi-Fi Systems.

Some representative publications are given below.

S.S. Ghassemzadeh, L.J. Greenstein, T. Sveinsson, A. Kavcic, and V. Tarokh, ``UWB Delay Profile Models for Residential and Commercial Indoor Environments", *IEEE Transactions on Vehicular Technology*, Jul. 2005.

S.S. Ghassemzadeh, R. Jana, C.W. Rice, W. Turin, and V. Tarokh, "Measurement and Modeling of An Indoor UWB Channel", *IEEE Transactions on Communications*, vol. 52, No. 10, pp. 1786-1796, Oct. 2004.

J. Bellorado, S.S. Ghassemzadeh, L.J. Greenstein, T. Sveinsson and V. Tarokh, "Coexistence of Ultra-Wideband Systems With IEEE-802.11 a Wireless LANs," *IEEE Global Telecommunications Conference* (GLOBECOM), 2003.

## 2.1.4. Distributed, Co-operative and Cognitive Communications:

From 2002-2008, we were interested in distributed communications. We put forth the idea that if a node "C" in a distributed network has a stronger channel to a "transmitting node A" then its intended receiver "B" (assuming that the underlying codebooks are appropriately designed", then "C" can extract the message of transmitter "A" faster than "B". For instance if A is transmitting 1 bit to B and the channel capacity between A and B is 0.1 but for A and C is 1, then A needs to send 10 seconds of coded data to B, while C can figure out the message only in 1s. This implies that C will know the signal that A is transmitted to B for the rest of 9 seconds (ignoring the processing time). This side information can be exploited in various useful manners.

Here (with my former student and postdoc Patrick Mitran and Hideki Ochiai), we envisioned that this "side information" can be used by C to help A in transmission in "opportunistic collaborative communications" (This was *later independently rediscovered* and published by Azarian and Hesham El Gamal too).

In another approach, we (with my former students Natasha Devroye and Patrick Mitran) pushed the alternative idea that C can use this "side information" to send messages to "D" while minimizing the interference (cognitive communications). Our group members (N. Devroye, P. Mitran, M. Vu, S. Yiu, W.Y. Shin) then used and generalized the above ideas to analyze gains of opportunistic cooperative and cognitive channels in a number of papers.

We have also designed a cooperative communications OFDM modem (jointly with Vanu's system on a SBIR grant subcontracted to Harvard). This led to several interesting problem (in particular distributed time and frequency synchronization, Pilot Sequence Design, Synchronization Sequence Design, etc.) that we addressed in various papers (with Peter Parker, Oh-Soon Shin and Lincoln Labs colleague Dan Bliss).

Also during the same period with (Hideki Ochiai, Patrick Mitran, and sabbatical visitor Professor H. Vincent Poor), we introduced a statistical approach and analyzed the performance of distributed beam-forming schemes, where the elements of the array are randomly positioned.

Yet another interesting and extremely challenging problem is that of distributed reverse link joint data and rate control for modern CDMA systems. This is implemented in 1xEVDO standard with 1 bit feedback from the base-station. All reverse link users must adjust their rates and powers with this 1 bit feedback reverse activity bit (RAB) channel. To achieve this, our team (jointly with Raymond Yim and Oh-Soon Shin) came up with an efficient method that has been mathematically proven to achieve the optimum performance subject to short-term and long-term fairness criteria.

In 2009, Behtash Babadi and I turned to exploiting mechanisms found in nature to deal with the problem of distributed dynamic spectrum allocation. Here, we used *Glauber Dynamics* to design an algorithm that not only performed significantly better than benchmark iterative water-filling but also had significantly lower complexity.

Some representative publications are given below.

Patrick Mitran, Hideki Ochiai, and V. Tarokh, "Space-Time Diversity Enhancements Using Collaborative Communications", *IEEE Transactions Information Theory*, vol. 51, No. 6, pp. 2041-2057, Jun. 2005.

Hideki Ochiai, Patrick Mitran, H.V. Poor and V. Tarokh, ``Collaborative Beamforming for Distributed Wireless Ad Hoc Sensor Networks", *IEEE Transactions Signal Processing*, vol. 53, No. 11, pp. 4110-4124, Nov. 2005.

Natasha Devroye, Patrick Mitran and V. Tarokh, "Achievable Rates in Cognitive Radio Channels," *IEEE Transactions Information Theory*, vol. 52, No. 5, May 2006.

Oh-Soon Shin, Albert Chan, H.T. Kung and V. Tarokh, "Design of An OFDM Co-Operative Diversity System," *IEEE Trans. Vehicular Technology*, vol. 56, No. 4, pp. 2203-2215, July 2007.

B. Babadi and V. Tarokh, "GADIA: A Greedy Asynchronous Distributed Interference Avoidance Algorithm," *IEEE Transactions on Information Theory*, vol. 56, No. 12, December 2010.

## 2..1.5. Contributions to Bi-directional Relaying/Coding:

In 2006-2008, we focused on the possibility of taking advantage of the fact that communications may be two-way (bi-directional) to significantly improve on classical communications using coded bi-directional relaying. Our team (with Sang Joon Kim and Patrick Mitran) computed theoretical limits of various existing coded bi-directional protocols. We (jointly with Toshiake Koike-Akino and Peter Popovski) designed unusual constellations for this purpose and showed great improvements over classical network coding.

Some representative publications are given below.

Sang Joon Kim, Patrick Mitran and Vahid Tarokh, "Performance Bounds for Bi-Directional Coded Cooperation Protocols," *IEEE Trans. Info. Theory*, vol. 54, No. 11, pp. 5235-5240, Nov. 2008.

T. Koike-Akino, P. Popovski, and Vahid Tarokh, "Optimized Constellations for Two-Way Wireless Relaying With Physical Network Coding", *IEEE Journal on Selected Areas in Communications*, vol. 27, No. 5, June 2009.

## 2.1.6. Ultimate Limits of Scheduling:

At about the same time, we were also interested in limit capacity of a multiuser wireless system assuming optimal scheduling. Here we (jointly with Dongwoon Bai) used the theory of extermal order statistics to compute these limits asymptotically.

Our results are highly *counter-intuitive*. It is well-known that Claude Shannon computed the capacity of some communications channels in closed form.  For various channels, a closed form expression for channel capacity is not known.  However, we showed using the *Theorem of Fisher-Tippett and Gnendenko* that even for a large class of such channels, the ultimate limits of scheduling can be calculated in closed form. Our tools are technical, and require advance knowledge of extreme value theory.

We have also shown that with a moderate number of users these scheduling limits can be closely approached, and in particular applied these methods to antenna selection schemes.

A representative publication is given below.

Dongwoon Bai, Patrick Mitran, Saeed S. Ghassemzadeh, Robert R. Miller, and Vahid Tarokh, "Rate of Channel Hardening of Antenna Selection Diversity Schemes and Its Implication on Scheduling," *IEEE Trans. Info. Theory*, vol. 55, No. 10, pp. 4353 - 4365, October 2009.

## 2.1.7. Communications and Signal Processing in Nonlinear Regime

Communications devices (e.g. power amplifiers and mixers) are inherently nonlinear. Most designers try to create/use devices that manifest nearly linear behaviors in certain regions. The operations of communications devices are then limited to these regions. Additionally many living communications mechanisms (such as the operation of neurons, etc.) are best modeled as nonlinear channels.

There have been some studies of communications and signal processing in the nonlinear regime, most notably for compensation of amplifier nonlinearities. In spite of this, information theory of nonlinear communications is in a primitive stage. A similar assessment can be made of communications and signal processing (for communications) in this regime.

We have worked on computing the theoretical limits of communications rates (channel capacities) for nonlinear channels. We have established some recent results on the capacity of Gaussian memory-less nonlinear channels. The results indicate that typical compressive nonlinearities limit the system capacity at high signal to noise ratios. As expected, our results reduce to Shannon's famous capacity formula when we approach the linear regime. We have extended these results to memory-less fading channels and MIMO channels. Specifically, we have shown that MIMO capacity gains continue to exist under memory-less nonlinear assumptions.

Many open and extremely challenging problems remain that we may further investigate. These include: the capacities of nonlinear channels with memory, the effect of feed-back, the capacities of channels with time-varying nonlinearities, and nonlinear network information theory.

A representative publication is given below.

M. Sabbaghian, A.I. Sulyman and V. Tarokh, "Analysis of the Impact of Nonlinearity on the Capacity of Communication Channels", *IEEE Transactions on Information Theory*, Vol. 59, No. 11, Nov. 2013.

## 2.2. Contributions to Signal Processing

### 2.2.1. Contributions to Compressed Sampling and Sparse Representation:

During 2008-2011, our team (Mehmet Akckaya and Behtash Babadi, and myself) paid attention to the problems of sparse representation and compressed sampling. These bounds improve on the $L_1$ approach and results/conjectures of Donoho, Candes and Tao, etc. Inspired by coding theory, we have constructed methods that perform close to these bounds with very low complexity.

Additionally, we considered adaptive identification of sparse linear system. We came up with algorithms that perform much better (of the order of 20 dB) and have a lot less complexity than the benchmark RLS algorithm.

Some representative publications are given below.

Behtash Babadi, Nick Kalouptsidis, and Vahid Tarokh, "SPARLS: The Sparse RLS Algorithm," *IEEE Transactions on Signal Processing*, Vol. 58, No. 8, pp. 4013 - 4025, Aug. 2010.

M. Akcakaya and V. Tarokh, "Shannon Theoretic Limits on Noisy Compressive Sampling," *IEEE Trans. on Information Theory*, vol. 56, No. 1, pp. 492-504, Jan. 2010.

Mehmet Akcakaya and Vahid Tarokh, "A Frame Construction and A Universal Distortion Bound for Sparse Representations," *IEEE Trans. Signal Processing*, Vol. 56, No. 6, pp. 2443-2550, June 2008.

### 2.2.2. Applications of Sparse Representation in Coronary MRI:

Former students Mehmet Akcakaya, Yongjun Kwak and Seunghoon Nam became also interested in applying our signal processing results in the area of compressed sensing (described above) to Coronary MRI. Here, we worked closely with Harvard Medical School. This portion of research of students and postdoctoral fellows was mainly supervised by Dr. R. Nezafat (with me a co-supervisor). We have since then continued our collaboration with the Harvard Medical School Radiology department.

Some representative publications are given below.

M. Akcakaya, P. Hu, M. L. Chuang, T. H. Hauser, L. H. Ngo, W. J. Manning, V. Tarokh and R. Nezafat, "Accelerated Non-Contrast Enhanced Pulmonary Vein MRA with Distributed Compressed Sensing," *Journal of Magnetic Resonance Imaging*, vol. 33, No. 5, pp. 1248-1255, May 2011.

M. Akcakaya, S. Nam, P. Hu, M. H. Moghari, L. H. Ngo, V. Tarokh, W. J. Manning and R. Nezafat, "Compressed Sensing with Wavelet Domain Dependencies for Coronary MRI: A Retrospective Study," *IEEE Transactions on Medical Imaging*, Volume: 30 , No. 6,  pp. 1090 - 1099 , May 2011.

## 2.2.3 Signal Processing for Nonlinear Communications

We have also considered signal processing for nonlinear channels.  For example, we have considered the identification of nonlinear channels represented by an unknown but small number of dominant non-linear modes (e.g. Volterra series terms). We have developed a number of recursive online algorithms for adaptive identification of such nonlinear channels. These algorithms provide significant improvement over the conventional algorithms even when restricted to linear regime, both in terms of mean squared error (MSE) and computational complexity. We have observed that gains on the order of 20 dB can be achieved at much lower complexities with an online implementation for both linear and also typical nonlinear channels. Currently, we are investigating various signal processing directions. The first is the problem of constellation design/coding/signal design for nonlinear channels. Previously, we had developed codes with low peak to average power ratios for multicarrier communications. However, these codes were designed to avoid the nonlinear regime of amplifiers. Little is known about constellation and signal design for given nonlinear channels. Additionally, we are investigating equalization of nonlinear channels.

A representative publication is given below.

G. Mileounis, B. Babadi, Nick Kalouptsidis and Vahid Tarokh, "An Adaptive Greedy Algorithm With Application to Nonlinear Communications," *IEEE Transactions on Signal Processing*, Vol. 58, No. 6, pp. 2998-3007, June 2010.

## 2.2.4. Localization in Heterogeneous Media

Localization of objects is a classical problem with applications in navigation, E-911 (emergency services), etc. The Global Positioning System (GPS) and other satellite based geo-location techniques are now in widespread commercial use. Other methods of localization based on Gyros, Wi-Fi Networks have also found some commercial success.

Recently source localization using measurements from spatially distributed sensors has received significant attention in many wireless network applications, such as cellular communication systems, search and rescue operations, environmental monitoring, logistics, inhabitant monitoring, etc.

In *Capsule Endoscopy*, a camera capsule takes pictures of a patient's GI tract and transmits the acquired data. We are interested in localizing this capsule inside the patient's body. Additionally, we anticipate that with the emergence of *medical micro-robots for minimally invasive medicine*, interest in localization of objects inside the human body may only increase in the years to come. However, the human body is not a heterogeneous media. Permittivity, refractive index, and absorption properties of human tissues vary between organs. Additionally, the presence or of absence of body fluids (e.g. in bladder) or other materials can significantly change the media of transmission. This motivates our studies of localization in heterogeneous media. Specifically, Esmaeil Nadimi and I have studied Bayesian approaches to this problem. We are currently investigating extensions to these methods.

Associated with the above is the important problem of *Channel Modeling for Radio Wave Propagation in Live Organs,* which we are also investigating (jointly with Professor K. Pahlavan). Some representative publications are given below.

Esmaeil S. Nadimi and Vahid Tarokh, "Bayesian Source Localization in Networks with Heterogeneous Transmission Medium'', *NAVIGATION: Journal of The Institute of Navigation*, Vol. 59, No. 3, Fall 2012.

Esmaeil S. Nadimi, Victoria Blanes-Vidal, Jakob L.F. Harslund, Mohammad H. Ramezani, Jens Kjeldsen, Per Michael Johansen, David Thiel and Vahid Tarokh, ``In Vivo and In Situ Measurement and Modeling of Intra-Body Effective Complex Permittivity", *Health Technology Letters*, Vol. 2, No. 6, pp. 135-140, Dec. 2015.

Kaveh Pahlavan, Yishuang Geng, David R. Cave, Guanqun Bao, Liang Mi, Emmanuel Agu, Andrew Karellas, Kamran Sayrafian and Vahid Tarokh, ``A Novel Cyber Physical System for 3-D Imaging of the Small Intestine In Vivo", *IEEE Access*, pp. 2730-2742, Dec. 201

## *2.2.4 Design Methods and Signal Processing for Interferometric Arrays:*

Jointly with colleagues at Lincoln Labs, we derived new results enabling robust interferometric image reconstruction in the presence of unknown aperture piston variation via the technique of redundant spacing calibration (RSC). The RSC technique uses redundant measurements of the same interferometric baseline with different pairs of apertures to reveal the piston variation among these pairs. In both optical and radio interferometry, the presence of phase-wrapping ambiguities in the measurements is a fundamental issue that needs to be addressed for reliable image reconstruction.

We showed that these ambiguities affect recently developed RSC phasor-based reconstruction approaches operating on the complex visibilities, as well as traditional phase-based approaches operating on their logarithm. We also derived new sufficient conditions for an interferometric array to be immune to these ambiguities in the sense that their effect can be rendered benign in image reconstruction. This property, which we called *wrap-invariance*, has implications for the reliability of imaging via classical three-baseline phase closures as well as generalized closures. We showed that wrap-invariance is conferred upon arrays whose interferometric graph satisfies a certain cycle-free condition. For cases in which this condition is not satisfied, a simple algorithm was designed for identifying those graph cycles which prevent its satisfaction.

We applied this algorithm to diagnose and corrected a member of a pattern family popular in the literature.

For image reconstruction using these arrays, we also developed a new RSC-based algorithm for prior-less phase recovery in which we generalize the bispectrum to higher order atmosphere-invariants (n-spectra) for improved sensitivity.

Some representative publications are given below.

Binoy G. Kurien, Jonathan B. Ashcom, Vinay N. Shah, Yaron Rachlin and Vahid Tarokh, ``Robust Interferometric Imaging via Prior-less Phase Recovery: Redundant Spacing Calibration with Generalized Closure Phases'', *Monthly Notices of Royal Astronomical Society*, vol. 464, No. 2, pp. 2356-2376, Sept. 2016

Binoy G. Kurien, Vahid Tarokh, Yaron Rachlin, Vinay N. Shah and Jonathan B. Ashcom, ``Resolving Phase Ambiguities in the Calibration of Redundant Interferometric Arrays: Implications for Array Design'', *Monthly Notices of Royal Astronomical Society*, vol. 461, No. 4, pp. 3585-3597, June 2016.

## 2.3.   Biological Signal Processing, Remote Health Monitoring, and Body Area Networks

In recent years, there has been an explosion of interest in remote health monitoring via mobile and electronic devices. Body area networks and their applications are also emerging at a rapid pace. We are actively focusing on these important domains (Please see the course I designed course on these topics "ES155: Biological Signal Processing" in the Teaching Section).

### 2.3.1 EEG Signal Detection and Applications

We have been focusing on analysis of EEG waves, including on *spindle detection* in the waves collected during sleep. Currently we are investigating whether spindles can be used to detect insomnia by examining the detected spindles for healthy and insomniac subjects.

Additionally, we are integrating and programming *devices that could be controlled by the human mind*. In general, we are interested in applications of EEG wave analysis such as brain computer interaction (BCI).

*Epileptic seizure prediction* may be another application of EEG signals. While detection of epileptic seizure signals is well-understood, the ability to predict the onset of seizures remains a challenging problem. With further research in the area of EEG wave analysis, we hope to develop the ability to predict a seizure at least an hour in advance of an episode. We thank Freiburg University for making their database available.

A representative publication is given below.

B. Babadi, S. M. McKinney, V. Tarokh, and J. M. Ellenbogen, "DiBa: A Data-Driven Bayesian Algorithm for Sleep Spindle Detection," *IEEE Transactions on Biomedical Engineering*, vol. 59, No. 2, pp. 483-493, Feb 2012.

### 2..3.2 GAIT Analysis and Applications

We are interested in unsupervised and real-time gait sensor data analysis for remote health monitoring, and in the detection of modes of human motion. Applications of interest include fall detection and early detection of Alzheimer's. We note that falls are a major cause of injury and death in the elderly. Also, connections between Alzheimer's and Gait data have been recently confirmed by various research groups ( A New York Times article on this topic can be found at http://www.nytimes.com/2012/07/17/health/research/signs-of-cognitive-decline-and-alzheimers-are-seen-in-gait.html?pagewanted=all&_r=0).

Jointly, with our colleagues at AT&T Labs-Research, we have developed algorithms and a system based on dedicated shoe insoles, and demonstrated some success in detection and classification of various modes of motion in real time (Short demo of our system can be found at http://people.seas.harvard.edu/~vahid/GAE.avi ).

### 2.3.3 Causality Relations in Alzheimer's Study

Jointly with colleagues at the Department of Radiology of Harvard Medical School, we have been looking at A-beta/Tau accumulation in different parts of brain. Our goal is to determine if there are causality relationships between the accumulations in different parts of brain. This can lead to the discovery of new pathways of Abeta development, etc. Our neurologist colleague (Dr. Jorge Sepulcre) at MGH believes that our algorithms may have found a new pathway for Abeta development that (according to his intuition) seems possible.

A representative publication is given below.

Hamed Farhadi, Yu Xiang, Seongah Jeong, Xiang Li, Ning Guo, Jorge Sepulcre, Vahid Tarokh, and Quanzheng Li, ``Alzheimer's Disease Study: A-beta/Tau Causality Test", *Society of Nuclear Medicine and Molecular Imaging (SNMMI) Annual Meeting*, June 2017 (to appear).

## 2.4   Contributions to Mathematics

We have produced various mathematical for the reason of personal curiosity. Some of these results are described below.

### 2.4.1. Infinite Dimensional Representations of Semi-simple Lie Algebras (1991-1992):

Finite dimensional representations of semi-simple Lie Algebras were totally classified by Cartan and Weyl. For the infinite dimensional representation, it is known that a full classification reduces

to that of torsion free Lie modules. In this line, I showed (jointly with D.J. Britten and F.W. Lemire) that certain torsion free infinite dimensional representations of Lie Algebra $A_n$ do not exist.

A representative publication is given below.

D.J. Britten, F.W. Lemire, and V. Tarokh, ``A Constraint on the Existence of Simple Torsion-Free Lie Modules'', *Proceedings of the American Mathematics Society*, vol. 123, No. 8, pp. 2315-2321, Aug. 1995.

### 2.4.2. Complexity of representation of Lattices by Trellises (1993-1995):

I (and I.F. Blake) computed the relationship between minimal complexity of trellis representation of lattices and their coding gain (equivalently sphere packing properties).

Some representative publications are given below.

V. Tarokh, and I.F. Blake, ``On the Trellis Complexity of The Densest Lattice Packings in $R^n$, *SIAM Journal of Discrete Math*, vol. 9, No. 4, pp. 597-601, Nov. 1996.

V. Tarokh, and I.F. Blake, ``Trellis Complexity Versus The Coding Gain of Lattices: Parts I and II'', *IEEE Transactions on Information Theory*, vol. 42, No. 6, pp. 1796-1816, Nov. 1996.

### 2.4.3. Sequence Design (1999-2003):

We have designed 16-QAM Golay complementary sequences from a multilevel representation of 16-QAM sequences.

A representative publication is given below.

C.V. Chong, R. Venkataramani and V. Tarokh, ``Two Constructions of 16-QAM Golay Complementary Sequences'', *IEEE Transactions on Information Theory*, vol. 49, No. 11, pp. 2953-2959, Nov. 2003.

### 2.4.4. Random matrix theory and properties of pseudorandom matrices produced from binary codes (2010-2012, 2016-2018):

We have considered pseudo-random matrices generated from binary linear block codes whose rows are formed by picking up the codewords of the code at random according to uniform and i.i.d distribution and by changing 0 $\rightarrow$ 1 and 1 $\rightarrow$ -1. We have shown that the spectra of these matrices and their products resemble those of Redmacher's matrices if the underlying codes have large dual distance.

Additionally, we have come up with constructions of pseudo-random binary (0 $\rightarrow$ 1 and 1 $\rightarrow$ -1) matrices with very low Kolmogorov complexity that achieve Wigner's spectrum. These

constructions prove that having semi-circular asymptotic spectrum for symmetric matrices is a weak property.

Some representative publications are given below.

B. Babadi and V. Tarokh, "Spectral Distribution of Product of Pseudorandom Matrices Formed From Binary Block Codes," *IEEE Transactions on Information Theory*, vol. 59, No. 2, Feb. 2013.

B. Babadi and V. Tarokh, "Spectral Distribution of Random Matrices from Binary Linear Block Codes," *IEEE Transactions on Information Theory*, vol. 57, No. 6, pp. 3955 - 3962, June 2011.

Ilya Soloveychik, Yu Xiang, Vahid Tarokh, "Pseudo-Wigner Matrices", IEEE Trans- actions on Information Theory, Vol. 64, No. 4, pp. 3170-3178, April 2018.

Ilya Soloveychik, Yu Xiang, Vahid Tarokh, "Symmetric Pseudo-Random Matrices", IEEE Transactions on Information Theory, Vol. 64, No. 4, pp. 3179-3196, April 2018.

## 2.5 Contributions to Data Science (2015-Present)

### 2.5.1. Representation, Modeling, Prediction and Inference from Data

One of the goals of our present research is developing "domain agnostic methods" for modeling, representation, inference and prediction from data. Many methods of learning from data in different domains are analogous, and involve applying variants of classical algorithms. This gives some hope that a domain agnostic method for vector-time series may be possible. Additionally, most observed natural vector-time series are multi-regime (a.k.a. multi-state) processes, where the dynamics of data slowly varies in a given state, and varies from a state to the other. Our approach involves detecting the state changes, modeling each state originally by a linear model and then bootstrapping it to a non-linear dynamical model (potentially from a dictionary of models), labeling the states (as some states may repeat themselves), and building dictionaries involving state transitions. This must be done both for online and off-line (batch) data. In building the original linear model, we have recently solved a classical problem. Specifically, we considered the problem of model selection for both auto-regression and regression, and proposed a new "information" criterion known as the bridge criterion. The new criterion has the benefits of the two well-known model selection techniques, the Akaike information criterion (AIC) and the Bayesian information criterion (BIC). In the specified case, BIC is known to be consistent, and so is the new criterion. In contrast, in the mis-specified case, the Akaike information criterion is known to be efficient (in the sense that its predictive performance is asymptotically equivalent to the best offered by the candidate models) and the new criterion behaves in a similar manner. Different from the classical criteria, the proposed criterion adaptively achieves either consistency or efficiency depending on the underlying true model. The proofs of his results are extremely difficult, but we can use sequential decision making on the Cohn Polyhedra to provide intuition. Interestingly, as a small byproduct of his work, we produce a very simple algorithm for uniform generation of stable

polynomials of a given degree, a highly non-trivial result. We also generalized Takeuchi's information criterion for general loss functions and nonlinear models. We proved that Generalized TIC (GTIC) can be used in a much broader context, and built a tensor graph of GTIC upon the theano platform and released a python package gtic. Our implementation was applicable for both generalized linear models and single-layer feed-forward neural networks. We then build on the work of Shibata to come up with methods of state change detection in multi-state processes. This algorithm is proved to be strongly consistent (under mild conditions). A simple version of this algorithm with reduced complexity is for instance applied to environmental data, where results of various papers on El Nino, Atlantic Mean Oscillation, etc. are automatically discovered in a domain agnostic manner. Other methods of change detection are also being considered. For bootstrapping from linear to non-linear models in each state assuming that the dynamics can be either described by a Volterra, or is n-times differentiable with continuous derivatives. The results are built into a system that automatically recovers very complicated dynamics and changes with small amount of data. Another ongoing research activity is in rare event prediction. We use Pickands-Balkemade Haan Theorem to build models for extreme events. We have built a system based on this result that can predict. Another line of our research is in localized stationarity tests. Knowing that a vector time series is locally stationary allows us for exploring relationships between various data streams. Additionally, computation of localized stationary windows may be used for change detection. Here we have extended Priestly-Rao's tests from univariate to multivariate, and built on Friedman's test to produce non-parametric local stationarity tests that perform extremely well for limited data. Additionally, we have developed new localized independence and Granger type causality tests based on the modified dynamic time warping (DTW) algorithm. We have also pursued a novel methodology that efficiently predicts outcomes/detect abrupt parameter changes in time-series with underlying time-varying true parameters based on casting statistical inference as expert learning. The experts correspond to the centers of the finitely many epsilon-balls covering the parameter space (assuming that the parameter space closure is compact and applying the Heine-Borel Theorem). Both theoretical analysis and practical experiments show that the proposed method can achieve near-optimal prediction loss under (both abrupt and smooth) unknown variations in data generating processes. We are now focusing on other regimes of vector-time series data. For instance, on the regime when the temporal number of observations is low but the space dimension is large. We are extending on the existing theories on panel data by allowing the use of prior knowledge into our tests. We are also applying our results to various applications such as BCI, approximate computing, environmental data, biological data, etc. Some representative publications are given below.

Stephane Shao, Pierre E. Jacob, Jie Ding, and Vahid Tarokh, " Bayesian Model Comparison with the Hyvarinen Score: Computation and Consistency", Journal of the American Statistical Association, Accepted for Publication.

Jie Ding, Vahid Tarokh, and Yuhong Yang, "Model Selection Techniques: An Overview", IEEE Signal Processing Magazine, Vol. 35, No. 6. pp. 16-34, Nov. 2018.

S. Shahrampour, M. Noshad, J. Ding, and V. Tarokh, "Online Learning for Multi- modal Data Fusion With Application to Object Recognition." IEEE Transactions on Circuits and Systems Ii: Express Briefs Vol. 65, No. 9, pp. 1259-1263, Septem- ber 2018.

Jie Ding, Vahid Tarokh and Yuhong Yang, "Bridging AIC and BIC: A New Crite- rion for Autoregression", IEEE Transaction on Information Theory, Vol. 64, No. 6, pp. 4024-4043, June 2018.

Jie Ding, S. Shahrampour, K. Heal, and V. Tarokh, "Analysis of Multi-state Au- toregressive Models", IEEE Transaction on Signal Processing, Vol. 66, No. 9, pp. 2429-2440, May 2018.

Jie Ding, Yu Xiang, Lu Shen, and Vahid Tarokh, "Multiple Change Point Analysis: Fast Implementation And Strong Consistency", IEEE Transaction on Signal Processing, Vol. 65, No. 17, pp. 4495 - 4510, Sept. 2017.

Qiuyi Han, Jie Ding, Edoardo Airoldi, and Vahid Tarokh, "SLANTS: Sequential Adaptive Nonlinear Modeling of for Vector Time Series", IEEE Transaction on Signal Processing, Vol. 65, No. 19, pp. 4994 - 5005, June 2017.

Shahin Shahrampour, Mohammad Noshad, Vahid Tarokh, "On Sequential Elim- ination Algorithms for Best-Arm Identification in Multi-Armed Bandits", IEEE Transaction on Signal Processing, Vol. 65, No. 16, pp. 4281- 4292, May 2017.

## 2.5.2. Limits of Learning

Currently, we lack an explicit computationally feasible method of determining the information content of data. This is unlike Shannon's information theory of certain channels and sources, where the objectives of transmission and compression are well-defined, and unlimited computational resources are available to the coders and decoders. If we were to make analogous theories to compute the underlying limits of learning from data, we must support multiple objectives, consider the cases where that learners may or may not be computationally limited, and when there may be no clear input-output relationship between different data streams. We first formalize how humans learn information content: First choose a concept from a perceived set of admissible hypotheses/theories, then create a hypothesis/theory, then as more observations arrive, refine or refute the theories, then synthesize more powerful explanations. We formalize the universe of admissible models in a concept space W, equipped with a filtration Vi (i=1,2, …) of ascending complexity that is dense in W. Learning is represented as moving higher in the filtration (refining the model) until either the training data is exhausted or the incremental gain (with respect to the learning objective) becomes small. The choice of concept space and filtration is critical. We will later on investigate dictionary-building methods for this purpose, which could be data- driven or based on prior knowledge (e.g. with help from a tutor/commander). In our formalization, the concept of information content learned from the data at a given time depends on the learning objective and the model class. We note the following important similarities with Shannon's information theory that we are exploring: (1) Tutors/Commanders & Side Information; (2) Learner Limitations & Cascade Channels, etc. We have recently come up with the required formalization and some significant results to calculate the limits of learning. Both for well-specified and mis-specified model selection methods, when the measure of learning is L2 Prediction error, we have come up with asymptotic closed form expressions for limits of learning. For a wide class of

objective functions, by making a connection between expert learning theory and Kolmogorov geometric entropy, we have come up with numerical methods to calculate the limits of learning. Some representative publications are given below.

Jie Ding, Jiawei Zhou and Vahid Tarokh, "Asymptotically Optimal Prediction for Time-Varying Data Generating Processes", IEEE Transactions on Information Theory, Accepted for Publication.

Jie Ding, Enmao Diao, Jiawei Zhou, and Vahid Tarokh, "A Penalized Method for the Predictive Limit of Learning" IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), pp. 4414-4418, April 2018.

Ahamd Beirami, M. Razaviyayn, S, Shahrampour, and Vahid Tarokh, "On Optimal Generalizability in Parametric Learning", 2017 Conference on Neural Information Processing System (NIPS), 2017.

## 2.5.3. Kernel Methods

At the heart of many machine learning problems, kernel methods (such as support vector machine) describe the nonlinear representation of data via mapping the features to a high dimensional feature space. Without recourse to explicit form of the feature maps, one can compute their inner products inexpensively using a kernel function, an idea known as the "kernel trick". However, unfortunately, methods using kernel matrices are not applicable to large-scale machine learning as they incur a massive computational cost on large datasets. This observation motivated researchers to consider kernel approximation using random features and extend the idea to train shallow architectures. A natural concern is then the stochastic oracle from which the features are sampled. Recently, it has been shown that employing data-dependent randomization improves the performance in terms of the required number of random features. Following this line of research, we are concerned with the randomized-feature approach in supervised learning for accurate prediction (generalization). We propose the Energy-based Exploration of Random Features (EERF) algorithm which relies on a data-dependent score function that explores the set of possible features and exploits the promising regions. We theoretically prove that the proposed score function with high probability recovers the spectrum of the best fit within the model class. We have also applied EERF to several practical datasets (e.g. MNIST for digit recognition). Our empirical results verify that our method requires smaller number of random features to achieve a certain generalization error compared to the state-of-the-art while introducing negligible pre-processing overhead and requiring no additional tuning parameters. A representative publication is given below.

Shahin Shahrampour and Vahid Tarokh, "Learning Boun for Greedy Approximation with Explicit Feature Maps from Multiple Kernels", 2018 Conference on Neural Information Processing System (NIPS), 2018.
Shahin Shahrampour, Ahmad Beirami, and Vahid Tarokh, On Data-Dependent Random Features for Improved Generalization in Supervised Learning, Thirty- Second AAAI Conference on Artificial Intelligence (AAAI) , 2018.

## *2.5.4. Inference of Non-Stationary Processes*

We have proposed a new inference procedure for understanding non-stationary processes, under the framework of evolutionary spectra developed by Priestley. Among various frameworks of modeling non-stationary processes, the distinguishing feature of the evolutionary spectra is its focus on the physical meaning of frequency. The classical estimate of the evolutionary spectral density is based on a double-window technique consisting of a short-Fourier transform and a smoothing. However, smoothing is known to suffer from the so-called bias leakage problem. By incorporating Thomson's multitaper method that was originally designed for stationary processes, we propose an improved estimate of the evolutionary spectral density, and analyze its bias/variance/resolution tradeoff. As an application of the new estimate, we further propose a non-parametric rank-based stationarity test, and provide various experimental studies. A representative publication is given below.

Yu Xiang, Jie Ding and Vahid Tarokh, "Estimation for Evolutionary Spectra with Application to Stationarity Test", IEEE Transactions on Signal Processing, Accepted for Publication.

## *2.5.6. High-Dimensional Statistics with Applications to Brain-Computer Interfaces*

A Brain-Computer Interface (BCI) is a mechanism to record data from the brain, process the recorded signal, and if required, provide feedback stimuli to activate regions of the brain. BCIs has potential applications in treating neurological disorders and the development of neuroprosthetics. Design of effective BCIs requires a fundamental understanding of how information is encoded in the brain. We are interested in developing theory and algorithms for high-dimensional statistical inference with applications to brain signal processing. Towards this broader goal, we have developed robust algorithms for classification of local field potentials (LFPs). LFPs are signals recorded from the brain of a human or an animal using microelectrodes. In the experiment of interest (performed by our collaborators at NYU), a macaque monkey is trained to perform memory guided and visually guided saccades. The objective is to use the LFPs to predict where the monkey is looking in each trial and also to decode what type of trial the monkey is doing. This type of inference is fundamental to neuroscience and it is important to develop classification algorithms that perform well over a wide range of problem instances and across different animals. We have shown that the LFP time-series data can be modeled in a non-parametric regression framework, i.e., as functions corrupted by Gaussian noise. We have then argued that the LFP classification problem should be modeled as a composite hypothesis testing problem for a class of functions. This formulation allows us to prove that using minimax function estimators as features for a classifier leads to robust and consistent classification. The theory of Gaussian sequence models allows us to represent minimax function estimators as finite-dimensional objects. We have used this idea to develop new classifiers and have successfully applied them to decode eye movement goals in the memory experiment. Some representative publications are given below.

Taposh Banerjee, John Choi, Bijan Pesaran, Demba Ba and Vahid Tarokh, "Wavelet Shrinkage and Thresholding Based Robust Classification for Brain-Computer Inter- face", IEEE

International Conference on Acoustics, Speech and Signal Processing (ICASSP), pp. 836-840, April 2018.

Taposh Banerjee, John Choi, Bijan Pesaran, Demba Ba, and Vahid Tarokh, "Classi- fication of Local Field Potentials using Gaussian Sequence Model", IEEE Statistical Signal Processing Workshop (SSP), pp. 683-687, June 2018.

## 2.5.7. Quickest Change Detection

With Applications to Neuroscience and Cyber-Physical Systems The core problem in this area is to find algorithms to detect a change in the distribution of a sequence of random variables, in real time, as quickly as possible. The change has to be detected while constraining the false alarms. The work in this area started in the 1930s for statistical quality control in manufacturing, where algorithms were needed to raise an alarm as soon as the quality of the products being manufactured deteriorates. Contemporary applications of change detection include intrusion detection, anomaly detection in cyber-physical systems, detecting the arrival of animals to their habitat, detecting a sudden change in stress on infrastructure, detecting line outages in power grids, etc. The theoretical foundations for this area were laid by Abraham Wald in the 1940s, and by Shiryaev and Kolmogorov in the 1960s. A typical algorithm for change detection is a single-threshold test, in which likelihood ratio of the observations is accumulated over time, and an alarm is raised the first time the accumulated statistic crosses a threshold. We are interested in developing theory and algorithms for change detection with applications to neuroscience and cyber-physical systems. In neuroscience, spike data collected from the brain of primates or mice are used for statistical inference and for brain signal processing. In some applications, it is of interest to detect a change in the firing pattern of neurons. In the experiment of interest, data is collected from the brain of a mouse under a shock experiment. It is believed that a change in firing pattern establishes that the mouse has learned to associate the cue with the shock. We have developed algorithms that can detect a change in firing patterns using just the baseline firing data for training. We have also developed algorithms for detection of changes using multi-modal data in non-stationary environments. The objective here is to use data from physics based sensors (video, EO/IR, etc) and social sensors (Twitter, Instagram) to detect anomalous behavior. This has applications in mission driven real-time tactical military surveillance to provide indications and warnings, e.g., to a soldier in degraded conditions. The anomaly detection problem using multimodal data also has applications for congestion control and event detection in big cities. In these applications, one common theme is the non-stationary of data even in normal conditions. We have developed a theory (mathematical models, optimization formulations, and optimal algorithms) for quickest change detection using multiple modalities in such nonstationary environments. A representative publication is given below.

Taposh Banerjee, Gene Whipps, Prudhvi Gurram, and Vahid Tarokh, "Sequential Event Detection Using Multimodal Data in Non-stationary Environments", 21st International Conference on Information Fusion (FUSION), pp. 1940-1947, July 2018.

## 2.5.8. Multi-Regime Random Fields

We are also interested in multi-regime random fields, which is a multidimensional generalization of the multi-regime random processes. In the latter, one of the most challenging and practically important problems is change detection. Assuming that the data generating process changes from one epoch of time to the other, one wants to find the point of the change which is also the boundary between the regions. The multidimensional change region detection is the problem of splitting the spatial/temporal domain occupied by a random field into a disjoint union of homogeneous regions inside which the data generating model is stable. Due to not being merely indexed by time axis, the notions of cause and effect lose their simple meanings and make the analysis highly non-trivial. The situation may be compared to the difference between one and two dimensional Ising models. To this end, the main focus of our work is on developing novel approaches and tools capable of detecting the regime changes, labeling the epochs between them (since a state may reappear again), describing the process in each state by a compact adaptive nonlinear model, and determining how regimes (states) evolve into each other. One of the most critical challenges in various applications is the scarce number of independent snapshots because of the constant change of the environment and necessity to take prompt decisions. Moreover, despite the fact that more available data increases the estimation precision, processing bigger volumes of measurements requires more computational resources from agents with limited capabilities and affects the system latency. These problems call for development of new techniques to solve real-time tasks with insufficient data. To approach the problem, we exploit the spatial structure of the network. Indeed, the sensors that are close to each other (in either space, time or frequency) may record signals that exhibit similarities. This structure affects the strength of interaction between the agents and dictates the sparsity pattern of the underlying graph making the learning with insufficient samples feasible. One of the most fundamental probabilistic models capturing the essence of the behavior of complex networks is the family of Markov Random Fields (RMF), or Undirected Graphical Models. Assuming that the underlying data generating process is an RMF embedded into a Euclidean space, our goal is to detect different regions of the network so that inside each region the interaction between the agents is similar. It is important to emphasize that in various applications the knowledge of all the edges is usually redundant, but what is more critical is the distribution of the edge parameters over the graph. Remarkably, even the two-dimensional Gaussian case already exhibits all the complexity of the problem, thus for concreteness we focus on it. Below we show that already this seemingly simple model relates to some recent deep mathematical results. The problem of region detection in RMF essentially boils down into model selection in graphical models where instead of detecting all the edges of the graph we want to detect homogeneous regions. The classical results require at least $n = c \log p$ samples to ensure reliable detection of all the edges, where $p$ is the number of vertices and $c$ depends on the graph properties and may be prohibitively large. Using the spatial structure, we develop a novel framework enabling learning the region in the sample starving scenario. The regions are assumed to have reasonably regular boundary and, to make the model class finite, are approximated by lattice polygons. Using rigorous information-theoretic approach, we derive tight necessary sample complexity bounds demonstrating that even bounded number of samples may be enough for consistent recovery of the graph regions. Derivation of such bound requires counting the admissible models represented in our case by polygons with certain parameters (such as perimeter, area, etc.) tiling the plane. This is a very challenging task that has become accessible only recently

due to some prominent developments in the theory of random integer partitions. In our paper work we utilize the Large Deviation Principle to count certain families of lattice polygons and to eventually derive the lower sample complexity bounds. We also propose a simple greedy algorithm capable of efficiently, reliably and quickly partitioning the graph into regions, and rigorously analyze its performance bound. We prove that the proposed method has optimal (up to a logarithmic factor) sample complexity. Some representative publications are given below.

Ilya Soloveychik, Vahid Tarokh, "Region Detection in Markov Random Fields: Gaussian Case", available online arxiv:1802.03848, 2018.

Ilya Soloveychik, Vahid Tarokh. "Large Deviations of Convex Polyominoes", available online arxiv:1802.03849, 2018.