

**BSc project proposals**

**in**

**Algorithms, Data Structures and  
Foundations of Machine Learning**

Peyman Afshani  
Gerth Stølting Brodal  
Kasper Green Larsen  
Chris Schwiegelshohn

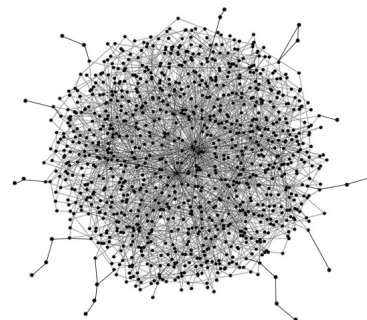
{peyman, gerth, larsen, schwiegelshohn}@cs.au.dk  
Nygaard 3<sup>rd</sup> floor

The following pages contain potential topics for BSc projects, but other projects are also possible. The final topic and direction of the project is settled under guidance by the advisor. In particular the weightening between theory and practical implementations is done on an individual basis, and is often adjusted during the BSc project process. Please do not hesitate to pass by our offices on or send us an e-mail for setting up a meeting to discuss potential BSc projects.

# Connectivity in Graph Sketching

**Problem** Let  $G(V, E)$  be a graph over  $n$  vertices, where  $n$  is known. We process edge insertions and deletion of the form  $(A_{u,v} + 1$  signifying an insertion and  $(A_{u,v} - 1$  signifying a deletion, where  $A$  is the adjacency matrix of  $G$ .

Our goal is to determine whether the graph is connected or not, using as little space as possible.



Graph streaming is a model of computation for very large graphs. Here, we process a sequence of edge insertions and deletions. At any point, we want to be able to query certain properties of the graph. The challenge is to do so in little space, i.e. we want to store roughly  $O(n)$  bits of space, instead of the maximum size of the graph which corresponds to  $O(n^2)$  space.

We focus on the connectivity problem, where we want to decide whether a graph is connected or not. If we are only processing edge insertion, it is very simple to decide whether the graph is connected with union find data structures. Details we leave as an exercise to the interested reader.

If we also process edge deletions, the problem becomes very challenging. Perhaps surprisingly, even in this case,  $O(n \log^3 n)$  bits of space are enough!

The aim of the project is to understand, reproduce and implement the classic algorithm as well as a few heuristics, run them on real world data sets and compare their performance.

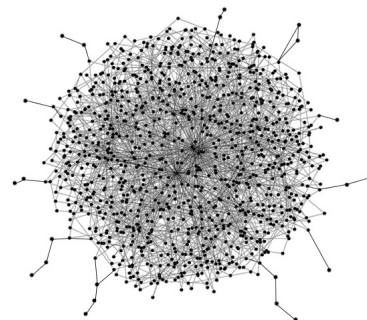
## References

- *Analyzing graph structure via linear measurements.* Kook Jin Ahn, Sudipto Guha, and Andrew McGregor. Proceedings of the Twenty-Third Annual ACM-SIAM Symposium on Discrete Algorithms, SODA 2012, Kyoto, Japan, January 17-19, 2012, pages 459–467.  
DOI: 10.1137/1.9781611973099.40.

**Contact:** Chris Schwiegelshohn

# Two-Pass Min-Cut in Graph Streaming

**Problem** Let  $G(V, E)$  be a graph over  $n$  vertices, where  $n$  is known. We process a sequence of edge insertions. At the end we want to output a global min cut of the graph, i.e. a partition of the graph into two sets  $S$  and  $V \setminus S$  such that  $|E \cap (S \times V \setminus S)|$  is minimized. Our goal is to use as little space and as few passes over the data as possible, while solving this problem exactly.



Graph streaming is a model of computation for very large graphs. Here, we process a sequence of edge insertions and deletions. At any point, we want to be able to query certain properties of the graph. The challenge is to do so in little space, i.e. we want to store roughly  $O(n)$  bits of space, instead of the maximum size of the graph which corresponds to  $O(n^2)$  space.

We focus on the min-cut problem, i.e. a partition of the graph into two sets  $S$  and  $V \setminus S$  such that  $|E \cap (S \times V \setminus S)|$  is minimized. One can show that if we use a single pass over the data, finding the optimal min-cut requires  $\Omega(n^2)$  space, i.e. we cannot do substantially better than simply storing everything. Perhaps surprisingly, two passes are enough to solve this problem in no more than  $O(n \log^2 n)$  space!

The aim of this project will be to give an exposition and literature review on this problem, both for exact and approximate guarantees and compare the different techniques.

## References

- *A Simple Semi-Streaming Algorithm for Global Minimum Cuts*. Sepehr Assadi and Aditi Dudeja, 4th Symposium on Simplicity in Algorithms, SOSA 2021, Virtual Conference, January 11-12, 2021, pages 172–180. [dl.acm.org/doi/10.1137/1.9781611976496.19](https://dl.acm.org/doi/10.1137/1.9781611976496.19).
- *Intractability of min- and max-cut in streaming graphs*, Mariano Zelke, Inf. Process. Lett., 111(3), 145–150, 2011. DOI: 10.1016/j.ipl.2010.10.017.

**Contact:** Chris Schwiegelshohn

# Tracking Frequent Items in Data Streams

Finding frequently occurring items in a dataset is a common task in data mining. In the data streaming literature, this problem is typically referred to as the heavy-hitters problem, which is as follows: a huge stream of items is processed from one end to the other. Each item in the stream can be thought of as an integer and the goal is to report the integers occurring most often while using limited memory. An integer in the stream may represent e.g. a source IP address of a TCP packet. One then may want to find sources with high utilization in a network traffic monitoring application. Another example has each integer representing how many times a concrete web surfer clicked a web ad on a particular site. The goal then becomes to identify the frequent ad clickers in web advertising. Common to these applications is that one wants a solution with tiny memory consumption and processing time.



Simons Institute

The aim of the project is to work with both real-life and synthetic data streams, and to implement, experiment with and compare different solutions to the heavy-hitters problem. An important aspect of several of the solutions to the heavy-hitters problem is that they rely on randomization. Thus a central task is to use basic probability theory to give a theoretical analysis of the solutions. The final report should include a theoretical section covering the basic analysis of different algorithms for the problem, a summary of the implemented algorithms, an experimental evaluation of the algorithms and a discussion of the obtained results.

## References

- Lecture notes by Kasper Green Larsen and from U.C. Berkeley.
- *An Improved Data Stream Summary: The Count-min Sketch and its Applications*. Graham Cormode, S. Muthukrishnan. Journal of Algorithms, 55(1), 58-75, 2005. DOI: 10.1016/j.jalgor.2003.12.001.
- *Finding Repeated Elements*. J. Misra, David Gries. Science of Computer Programming, 2(2), 143-152, 1982. DOI: 10.1016/0167-6423(82)90012-0.

**Contact:** Kasper Green Larsen

# Closest Points

**Problem** Let  $P$  be a set of  $n$  points in the  $\mathbb{R}^d$ , given by their coordinates. The problem is to find two points  $p, q \in P$  such that the Euclidean distance between them is the smallest.

This is a classical computational geometry problem. Clearly, we can find the closest pair by going through all the pairs of points, compute the distance between every pair and then find the minimum among them. Unfortunately, this is a rather slow algorithm which runs in  $\Theta(n^2)$  time.

The first elegant solution came by Shamos and later published together with Hoey in 1975. Their algorithm was a simple divide and conquer algorithm: partition the point set into two equal sizes,  $P_\ell$  and  $P_r$ , e.g., with a vertical line  $h$  and then find the closest pair in each subset recursively. Let  $d$  be the smallest distance that we find in this way. If the closest pair is completely inside one subset, we are done. Otherwise, one point from  $P_\ell$  and another point from  $P_r$  must be the closest points. In particular, we only need to consider the points within distance  $d$  of the line  $h$ . Now, they observed that only a  $O(n)$  pairs need to be considered as each point in  $P_\ell$  or  $P_r$  cannot have too many points of distance  $d$  next to them. This gives the recursion

$$f(n) = 2f(n/2) + O(n)$$

which solves to  $f(n) = O(n \log n)$ .

This solution generalizes to a point set in 3D in a more or less straightforward way and it gives an algorithm with  $O(n \log^2 n)$  running time. The running time comes from a recursion

$$g(n) = 2g(n/2) + f(n)$$

where  $f(n)$  is the time it takes to solve a 2D closest pair problem. If we go with  $f(n) = O(n \log n)$ , then we get  $g(n) = O(n \log^2 n)$ . However, a brilliant solution by Shamos and Bentley showed that this can actually be improved to  $O(n \log n)$ .

The goal of this project is to follow these developments. The tasks are: (1) Implement the  $O(n \log n)$  algorithm for finding the closest pair in a set of  $n$  points in the plane. (2) Implement an  $O(n \log^2 n)$  algorithm for finding the closest pair in a set of  $n$  points in 3D. (3) Implement the algorithm by Bentley and Shamos that runs in  $O(n \log n)$  time to find the closest pair in a set of  $n$  points in 3D.

## References

- *Divide-and-Conquer in Multidimensional Space*. Jon Louis Bentley and Michael Ian Shamos. 8th Annual ACM Symposium on Theory of Computing (STOC), 220-230, 1976. DOI: 10.1145/800113.803652.
- *Closest-Point Problems*. Michael Ian Shamos and Dan Hoey. 16th Annual Symposium on Foundations of Computer Science (FOCS), 151-162, 1975. DOI: 10.1109/SFCS.1975.8.

**Contact:** Peyman Afshani

# Smallest Enclosing Ball

**Problem** Let  $P$  be a set of  $n$  points in  $d$ -dimensions, given by their coordinates. The problem is to find the smallest ball that contains all the points. Thus, the output is the center of the ball and its radius.

This is one of the classical problems in clustering. If we assume  $d$  is a constant, then the problem can be solved in linear time but often the dependencies are exponential in  $d$ . In particular, the best algorithm has the running time of  $O(d^2n + 2^{O(\sqrt{\log n})})$  which is exponential in  $d$ . As a result, when the dimension is large, settling for approximations is a reasonable choice.

In this project, you will implement and compare two simple solutions.

1. This algorithm achieves a  $(1 + \varepsilon)$  factor approximation in time  $O(nd/\varepsilon + 1/\varepsilon^5)$  (Badoiu and Clarkson, 2003).
2. However, there is also a simpler algorithm that achieves a  $3/2$  factor approximation in linear time, in fact, using only one scan of the input (Zarrabi-Zadeh and Chan, 2006).

The goal of this project is to implement these solutions and to compare them.

**Task one.** Implement the algorithms. Run some basic tests to make sure that they are correct.

**Task two.** Run some experiments to measure the running time, and the quality of the solution. Note that one algorithm should be fast but less accurate while the other one should be the opposite.

**Task three.** In your report, explain briefly how you implemented the algorithms and list any difficulties in doing so. You also need to briefly explain how the algorithms work but you don't have to cover the proof of correctness parts. Were the descriptions of the algorithms sufficient for the implementation or did you have to make some non-obvious choices? Describe how you designed your test cases and how you generated the input. Then, try to justify the results of the experiments.

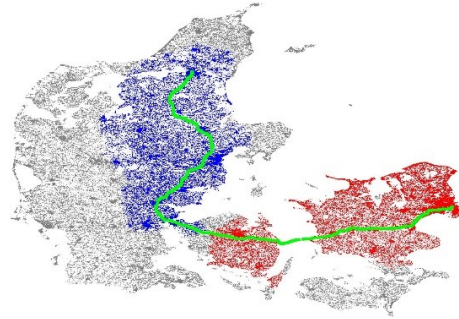
## References

- *Smaller Core-sets for Balls*. Mihai Badoiu and Kenneth L. Clarkson. Proceedings of the Fourteenth Annual ACM-SIAM Symposium on Discrete Algorithms (SODA), 801–802, 2003. <https://dl.acm.org/citation.cfm?id=644240>
- *A Simple Streaming Algorithm for Minimum Enclosing Balls*. Hamid Zarrabi-Zadeh and Timothy Chan. Proc. 18th Canadian Conference on Computational Geometry (CCCG), August, 2006. <http://www.cs.queensu.ca/cccg/papers/cccg36.pdf>

**Contact:** Peyman Afshani

# Shortest Paths Computations on Open Street Map Data

Open Street Map data contains a very detailed description of road networks worldwide and is the underlying data use in many map applications. The data is annotated with both geometric features and logically information of e.g. connection of road segments, speed limits, road types, one-way streets e.t.c.



Martin Jacobsen

The aim of this project is to be able to work with open street map data, in particular to visualize open street map data, convert open street map data to graph representations, implement algorithms for finding shortest paths in road network graphs, and to evaluate experimentally the performance of various algorithms. An important part of the process will be to read the theory behind selected shortest path algorithm, e.g. Dijkstra's algorithm, bidirectional shortest path algorithms, A\*, landmark based algorithms, contraction hierarchies, highway hierarchies, hub labelling, and transit node algorithms. The final report should include a summary of the implemented theory, a description of the implementation, an experimental evaluation of the implemented algorithms, and a discussion of the obtained results — in particular a discussion of space usage versus query time.

The Open Street Map data can also be combined with other data sources, like the Danish height elevation model, to take elevation differences into account for e.g. finding optimal bicycle routes.

## References

- *Route Planning in Transportation Networks*. Hannah Bast, Daniel Delling, Andrew Goldberg, Matthias Müller-Hannemann, Thomas Pajor, Peter Sanders, Dorothea Wagner, Renato F. Werneck. Algorithm Engineering 2016: 19-80. DOI: 10.1007/978-3-319-49487-6\_2.
- Videos with Andrew Goldberg on the topic: Beyond Worst Case Analysis Workshop 2011, NWU 2013.
- Slides from the Algorithm Engineering course, 2017: route.pdf.

**Contact:** Gerth Stølting Brodal.

# **BSc project proposals**

## **in**

# **Computational Complexity and Game Theory**

Ioannis Caragiannis  
Kristoffer Arnsfelt Hansen  
Srikanth Srinivasan

{iannis, arnsfelt, srikanth}@cs.au.dk  
Nygaard 3<sup>rd</sup> floor

The following pages contain potential topics for BSc projects, but other projects are also possible. The final topic and direction of the project is settled under guidance by the advisor. In particular the weighting between theory and practical implementations is done on an individual basis, and is often adjusted during the BSc project process. Please do not hesitate to pass by our offices on or send us an e-mail for setting up a meeting to discuss potential BSc projects.



# Gradient Descent: Theory and Practice

The aim of this project is to study first-order methods in convex optimization. The project should be a mix of theory and practice. The theoretical part of the project consists presenting the mathematical foundations and of analyzing convergence rates and complexity of several of the main methods used for solving large-scale problems. The practical part of the project consists of implementing and comparing these methods.

## References

- To be determined based upon scope of project in terms of theory and practice.

**Contact:** Kristoffer Arnsfelt Hansen

# The Multiplicative Weights Update Method in Game Theory

The idea of the multiplicative weights update method is an important online learning method with applications to many areas, including machine learning and optimization. The goal of this project is to explore the method in the context of game theory. Classic work of Freund and Schapire show that the averages of the strategies produced by the method converge to optimal strategies of zero-sum games. Recent work of Bailey and Piliouras show that while the averages converge to optimal strategies, the sequence may in fact be repelled by optimal strategies. The goal of the project is to present the theoretical background of the multiplicative update method and perform experiments with the method on concrete games.

## References

- *Adaptive game playing using multiplicative weights*. Yoav Freund and Robert E Schapire. Games and Economic Behavior, 29(1-2), 79–103, 1999. DOI: 10.1006/game.1999.0738.
- *Multiplicative Weights Update in Zero-Sum Games*. James P. Bailey and Georgios Piliouras. Proceedings 2018 ACM Conference on Economics and Computation (EC), 321–338, 2018. DOI: 10.1145/3219166.3219235.

**Contact:** Kristoffer Arnsfelt Hansen

# Fair Division Algorithms for Indivisible Items

Fair division is the task of allocating items to agents so that some fairness criterion is satisfied. The problem finds applications in every setting where some kind of allocation of resources is required. In the very popular variations with indivisible items, problem instances consist of a set of  $n$  agents and a collection of  $m$  items. Agents have valuations for the items: agent  $i$  has a valuation  $v_i(j)$  for item  $j$ . Valuations for sets (or bundles) of items are then additive, with the valuation of agent  $i$  for the bundle of items  $S$  being  $v_i(S) = \sum_{j \in S} v_i(j)$ . An allocation of items to agents is simply a partition  $A = (A_1, A_2, \dots, A_n)$  of the items to the agents so that agent  $i$  gets the bundle of items  $A_i$ .

An important fairness criterion is envy-freeness. We say that an allocation  $A$  is envy-free if for every pair of agents  $i$  and  $j$ , it holds that  $v_i(A_i) \geq v_i(A_j)$ . In words, every agent is happy with his/her bundle of items and does not envy the bundle allocated to some other agent. It can be easily seen that envy-freeness is rarely feasible. Indeed, consider an instance with two agents, who have positive valuations for a single item. Then, any allocation that gives the item to one of the agents is not envy-free. For this reason, several relaxations of envy-freeness have been considered in the literature.

Among them, the notion of envy-freeness up to any item (or EFX) is the most appealing. An allocation  $A$  is EFX, if for every pair of agents  $i$  and  $j$ , it holds that  $v_i(A_i) \geq v_i(A_j \setminus \{g\})$ , for any item  $g \in A_j$ . In words, an allocation is EFX if the possible envy of an agent for another is eliminated by removing any item from the bundle of items allocated to the latter. Unfortunately, until now, it is not known whether all instances have EFX instances or not. Fortunately, there are polynomial-time algorithms to compute partial EFX allocations (i.e., by discarding some of the items), simultaneously providing additional efficiency guarantees (e.g., they make sure that all agents get a bundle of items, for which they have a high value). The objective of this project is to study the recent literature on fair division of indivisible items, with a particular focus on the proposed algorithms for computing partial EFX allocations. Part of the work will include the implementation of these algorithms (and of variations of them) and their testing on real and synthetic instances.

## References

- *A Little Charity Guarantees Almost Envy-Freeness*. Bhaskar Ray Chaudhury, Telikepalli Kavitha, Kurt Mehlhorn, and Alkmini Sgouritsa. In *Proceedings of the 31st ACM-SIAM Annual Symposium on Discrete Algorithms (SODA)*, 2658–2672, 2020. DOI: 10.1137/1.9781611975994.162.
- *Envy-Freeness Up to Any Item with High Nash Welfare: The Virtue of Donating Items*. Ioannis Caragiannis, Nick Gravin, and Xin Huang. In *Proceedings of the 20th ACM Conference on Economics and Computation (EC)*, 527–545, 2019. DOI: 10.1145/3328526.3329574.
- *Spliddit: Unleashing Fair Division Algorithms*. Jonathan R. Goldman and Ariel D. Procaccia. *SIGecom Exchanges*, 13(2), 41–46, 2014. DOI: 10.1145/2728732.2728738.

**Contact:** Ioannis Caragiannis

# Algorithmic Aspects of Participatory Budgeting

Participatory budgeting is a relatively recent trend used by municipalities and regional governments worldwide to decide the distribution of funding to public projects. The main idea is to let the citizens express their opinion through voting over spending priorities and then implement the outcome of the vote as a public decision.

The details of voting are rather complicated (compared to a typical government election) and include a series of steps and features. We briefly present three components that are interesting from an algorithmic viewpoint. The first is the selection of the available alternatives and their nature. For example, alternatives could be discrete such as the construction of a bridge or continuous such as the length of bicycle paths. Second, how are the actual preferences of the citizens will be structured in the ballots? This step may result in a loss of information. Third, how should the votes be aggregated into a collective decision? This part is the most crucial algorithmically but strongly depends and interplays with the previous two.

The interdisciplinary field of computational social choice (lying at the intersection of social choice theory and theoretical computer science), which traditionally studies the computational complexity of voting rules, has extensively considered participatory budgeting recently, together with other modern trends of participatory democracy. The focus has been on axiomatic properties, and on quantifying the lack of efficiency of the voting process, due to the loss of information when expressing preferences.

The aim of the project is to survey the recent approaches in participatory budgeting, focusing on theoretical developments and analysis, case studies from its application to municipalities around the world, and available online tools. The project will involve the implementation of representative aggregation methods and experiments/simulations with synthetic voting scenarios.

## References

- Web: [www.participatorybudgeting.org](http://www.participatorybudgeting.org), [pbstanford.org](http://pbstanford.org).
- *Participatory Budgeting: Models and Approaches*. Haris Aziz and Nisarg Shah. In *Pathways Between Social Science and Computational Social Science: Theories, Methods, and Interpretations*, Springer, 2021. Forthcoming (available from Nisarg Shah's homepage).
- *Interactive Democracy*. Markus Brill. In *Proceedings of the 17th International Conference on Autonomous Agents and Multi-Agent Systems (AAMAS)*, 1183–1187, 2018. [dl.acm.org/doi/10.5555/3237383.3237873](https://dl.acm.org/doi/10.5555/3237383.3237873).

**Contact:** Ioannis Caragiannis

# Pseudorandom Generators from Theory in Practice

Generating random bits is a task of fundamental importance in Computer Science. Many real-world algorithms need access to a large number of high quality random

bits in order to perform various computational tasks efficiently and/or securely. Randomness is also useful in computational simulations of complex phenomena such as the climate or the economy.

But where do these algorithms get access to these random bits? It is hard to find sources of ‘true randomness’. Many standard practical ways of obtaining ‘random’ sequences are not random at all: there are standard ways of finding patterns in such sequences (see, e.g. [http://www.pro-technix.com/information/crypto/pages/rfc1750\\_base.html](http://www.pro-technix.com/information/crypto/pages/rfc1750_base.html)).

Given the difficulties involved in obtaining true randomness, it makes sense to think of randomness as a *resource*, much like running time or memory. A standard way to minimize the use of this resource is to use *Pseudorandom Generators* (PRGs).

A PRG is an algorithm that stretches a short, truly random string to a long string that ‘looks random’ to many standard algorithms. The study of the theory of PRGs began in cryptography in the early 1980s, and since then, researchers have developed PRGs for various kinds of applications.

This project is to understand and possibly develop provably correct kinds of PRGs for various computational tasks.

On the theoretical side, you can understand various constructions of PRGs. This will require some elementary notions in probability and algebra, which you can pick up during the course of the project.

On the practical side, you can implement these theoretical constructions on standard statistical tests used to test random sequences in practice.

## References

- Michael Luby, Avi Wigderson: *Pairwise Independence and Derandomization*. Found. Trends Theor. Comput. Sci. 1(4) (2005)
- Salil Vadhan: *Pseudorandomness*: A monograph on PRG constructions. <https://people.seas.harvard.edu/~salil/pseudorandomness/>
- [https://en.wikipedia.org/wiki/Diehard\\_tests](https://en.wikipedia.org/wiki/Diehard_tests): Wikipedia article on standard tests for pseudorandom sequences.

**Contact:** Srikanth Srinivasan

# The Minimum Circuit Size Problem

Consider the following two strings:

$x = 0101010101010101010101010101010101$     $y = 111010001010000110101100011110111101$ .

Which of these is more random? The reasonable answer is that  $y$  is more random than  $x$ , but what does this mean? And more importantly, what does this have to do with Computer Science?

The question of defining and quantifying notions of randomness has been taken up by many mathematicians. Some of the most interesting such notions are *computational*, defined by means of computers and algorithms. The notion of *Kolmogorov complexity* deals with how compactly a given string  $x$  can be compressed by a general algorithm. A string is defined to be random if it cannot be meaningfully compressed by any algorithm.

Unfortunately, the Kolmogorov complexity of a string is itself a difficult quantity to understand, as the Kolmogorov complexity of string is uncomputable. However, to overcome this, we have a number of resource-bounded notions of Kolmogorov complexity, all of which quantify various notions of compressibility. The *Minimum Circuit Size Problem* is an umbrella term for the corresponding computational questions: it is the algorithmic problem of computing the most ‘compressed’ version of a given string  $x$ .

These notions are useful in both theory and practice. In practice, these problems are related to fundamental problems in logic design. In theory, the computational complexity of this problem has been linked to the existence of fundamental objects in cryptography.

The aim of the project is to understand variants of the Minimum Circuit Size Problem, the algorithms we know for these variants, and where we suspect (or can show) that the problem is computationally hard. Understanding these results may require some elementary notions related to probability, algebra, approximation algorithms and basic polynomial-time reductions, which we can pick up as we go along.

## References

- Eric Allender: *The Complexity of Complexity*. Computability and Complexity 2017: 79-94.
- Eric Allender, Lisa Hellerstein, Paul McCabe, Toniann Pitassi, Michael E. Saks: *Minimizing Disjunctive Normal Form Formulas and AC0 Circuits Given a Truth Table*. SIAM J. Comput. 38(1): 63-84 (2008).
- Shuichi Hirahara, Igor Carboni Oliveira, Rahul Santhanam: *NP-hardness of Minimum Circuit Size Problem for OR-AND-MOD Circuits*. Computational Complexity Conference 2018: 5:1-5:31.

**Contact:** Srikanth Srinivasan