



FACULTY OF VETERINARY MEDICINE  
approved by EAEVE



# Studying dynamic processes with networks and trajectories

Robrecht Cannoodt

Thesis submitted in fulfilment of the requirements for the degree of  
Doctor in Computer Science, 2019

Supervisors:

Prof. Dr. Yvan Saeys

Prof. Dr. Katleen De Preter

Vakgroep Toegepaste Wiskunde, Informatica, en Statistiek  
Faculteit Wetenschappen, Universiteit Gent  
Krijgslaan 281 - S2, 9000 Gent



---

## Acknowledgements

---

Hello, here is some text without a meaning. This text should show what a printed text will look like at this place. If you read this text, you will get no information. Really? Is there no information? Is there a difference between this text and some nonsense like “Huardest gefburn”? Kjift – not at all! A blind text like this gives you information about the selected font, how the letters are written and an impression of the look. This text should contain all letters of the alphabet and it should be written in of the original language. There is no need for special content, but the length of words should match the language.



<b>Acknowledgements</b>	<b>iii</b>
<b>1 Introduction</b>	<b>1</b>
1.1 Single cell biology . . . . .	2
1.1.1 Dynamic processes . . . . .	2
1.1.2 Gene regulation . . . . .	2
1.1.3 Single cell transcriptomics . . . . .	2
1.2 Machine learning . . . . .	2
1.2.1 Supervised learning . . . . .	2
1.2.2 Unsupervised learning . . . . .	2
1.2.3 Feature selection . . . . .	2
1.2.4 Trajectory inference . . . . .	2
1.2.5 Network inference . . . . .	2
1.3 Research objectives . . . . .	2
1.4 Outline . . . . .	2
<b>2 A comparison of single-cell trajectory inference methods</b>	<b>3</b>
<b>3 Fast, accurate, and robust single-cell pseudotime</b>	<b>5</b>
<b>4 A toolkit for inferring and interpreting trajectories</b>	<b>7</b>
<b>5 Inferring single cell regulatory networks</b>	<b>9</b>
<b>6 Optimising regulatory networks</b>	<b>11</b>
<b>7 General discussion</b>	<b>21</b>
7.1 Overview and conclusions of the presented work . . . . .	21
7.2 Future research directions . . . . .	21
<b>Samenvatting</b>	<b>23</b>
<b>Summary</b>	<b>25</b>
<b>List of Publications</b>	<b>27</b>



CART Classification And Regression Trees

DNA Deoxyribonucleic Acid

GRN Gene Regulatory Network

IM Importance Measure

ML Machine Learning

mRNA Messenger RNA

NI Network Inference

RF Random Forests

RNA Ribonucleic Acid

TF Transcription Factor





## CHAPTER 1

---

Introduction

---

**Abstract:** Recent developments in single-cell transcriptomics have opened new opportunities for studying dynamic processes in immunology in a high-throughput and unbiased manner. Starting from a mixture of cells in different stages of a developmental process, unsupervised trajectory inference algorithms aim to automatically reconstruct the underlying developmental path that cells are following. In this review, we break down the strategies used by this novel class of methods, and organize their components into a common framework, highlighting several practical advantages and disadvantages of the individual methods. We also give an overview of new insights these methods have already provided regarding the wiring and gene regulation of cell differentiation. As the trajectory inference field is still in its infancy, we propose several future developments which will ultimately lead to a global and data-driven way of studying immune cell differentiation.

Adapted from:

**Cannoodt, R.\***, Saelens, W.\*, and Saeys, Y. Computational methods for trajectory inference from single-cell transcriptomics. *European Journal of Immunology* 46, 11 (2016), 2496–2506. doi:10.1002/eji.201646347.

\* Equal contribution

## 1.1 Single cell biology

### 1.1.1 Dynamic processes

### 1.1.2 Gene regulation

### 1.1.3 Single cell transcriptomics

## 1.2 Machine learning

### 1.2.1 Supervised learning

### 1.2.2 Unsupervised learning

### 1.2.3 Feature selection

### 1.2.4 Trajectory inference

### 1.2.5 Network inference

## 1.3 Research objectives

## 1.4 Outline

---

### A comparison of single-cell trajectory inference methods

---

**Abstract:** Hello, here is some text without a meaning. This text should show what a printed text will look like at this place. If you read this text, you will get no information. Really? Is there no information? Is there a difference between this text and some nonsense like “Huardest gefburn”? Kjift – not at all! A blind text like this gives you information about the selected font, how the letters are written and an impression of the look. This text should contain all letters of the alphabet and it should be written in of the original language. There is no need for special content, but the length of words should match the language.

Adapted from:

Saelens, W.\*, **Cannoodt, R.\***, Todorov, H., and Saeys, Y. A comparison of single-cell trajectory inference methods. *Nature Biotechnology* 37, 5 (2019), 547–554. doi:10.1038/s41587-019-0071-9.

\* Equal contribution

Hello, here is some text without a meaning. This text should show what a printed text will look like at this place. If you read this text, you will get no information. Really? Is there no information? Is there a difference between this text and some nonsense like “Huardest gefburn”? Kjift – not at all! A blind text like this gives you information about the selected font, how the letters are written and an impression of the look. This text should contain all letters of the alphabet and it should be written in of the original language. There is no need for special content, but the length of words should match the language.

---

### Fast, accurate, and robust single-cell pseudotime

---

**Abstract:** Hello, here is some text without a meaning. This text should show what a printed text will look like at this place. If you read this text, you will get no information. Really? Is there no information? Is there a difference between this text and some nonsense like “Huardest gefburn”? Kjift – not at all! A blind text like this gives you information about the selected font, how the letters are written and an impression of the look. This text should contain all letters of the alphabet and it should be written in of the original language. There is no need for special content, but the length of words should match the language.

Adapted from:

**Cannoodt, R.**, . . . , De Preter, K., and Saeys, Y. SCORPIUS improves trajectory inference and identifies novel modules in dendritic cell development. *Journal* vol, issue (2019), page–page. doi:10.1101/079509v2.

Hello, here is some text without a meaning. This text should show what a printed text will look like at this place. If you read this text, you will get no information. Really? Is there no information? Is there a difference between this text and some nonsense like “Huardest gefburn”? Kjift – not at all! A blind text like this gives you information about the selected font, how the letters are written and an impression of the look. This text should contain all letters of the alphabet and it should be written in of the original language. There is no need for special content, but the length of words should match the language.

---

### A toolkit for inferring and interpreting trajectories

---

**Abstract:** Hello, here is some text without a meaning. This text should show what a printed text will look like at this place. If you read this text, you will get no information. Really? Is there no information? Is there a difference between this text and some nonsense like “Huardest gefburn”? Kjift – not at all! A blind text like this gives you information about the selected font, how the letters are written and an impression of the look. This text should contain all letters of the alphabet and it should be written in of the original language. There is no need for special content, but the length of words should match the language.

Adapted from:

**Cannoodt, R.**, Saelens, W., and Saeys, Y. dyno. *Journal* vol, issue (2019), page–page. doi.

Hello, here is some text without a meaning. This text should show what a printed text will look like at this place. If you read this text, you will get no information. Really? Is there no information? Is there a difference between this text and some nonsense like “Huardest gefburn”? Kjift – not at all! A blind text like this gives you information about the selected font, how the letters are written and an impression of the look. This text should contain all letters of the alphabet and it should be written in of the original language. There is no need for special content, but the length of words should match the language.



---

### Inferring single cell regulatory networks

---

**Abstract:** Hello, here is some text without a meaning. This text should show what a printed text will look like at this place. If you read this text, you will get no information. Really? Is there no information? Is there a difference between this text and some nonsense like “Huardest gefburn”? Kjift – not at all! A blind text like this gives you information about the selected font, how the letters are written and an impression of the look. This text should contain all letters of the alphabet and it should be written in of the original language. There is no need for special content, but the length of words should match the language.

Adapted from:

**Cannoodt, R.**, Saelens, W., and Saeys, Y. Inferring Single Cell Regulatory Networks. *Journal* vol, issue (2019), page–page. doi.

Hello, here is some text without a meaning. This text should show what a printed text will look like at this place. If you read this text, you will get no information. Really? Is there no information? Is there a difference between this text and some nonsense like “Huardest gefburn”? Kjift – not at all! A blind text like this gives you information about the selected font, how the letters are written and an impression of the look. This text should contain all letters of the alphabet and it should be written in of the original language. There is no need for special content, but the length of words should match the language.

---

Optimising regulatory networks

---

**Abstract:** Graphlets are small network patterns that can be counted in order to characterise the structure of a network (topology). As part of a topology optimisation process, one could use graphlet counts to iteratively modify a network and keep track of the graphlet counts, in order to achieve certain topological properties. Up until now, however, graphlets were not suited as a metric for performing topology optimisation; when millions of minor changes are made to the network structure it becomes computationally intractable to recalculate all the graphlet counts for each of the edge modifications.

IncGraph is a method for calculating the differences in graphlet counts with respect to the network in its previous state, which is much more efficient than calculating the graphlet occurrences from scratch at every edge modification made. In comparison to static counting approaches, our findings show IncGraph reduces the execution time by several orders of magnitude. The usefulness of this approach was demonstrated by developing a graphlet-based metric to optimise gene regulatory networks. IncGraph is able to quickly quantify the topological impact of small changes to a network, which opens novel research opportunities to study changes in topologies in evolving or online networks, or develop graphlet-based criteria for topology optimisation.

IncGraph is freely available as an open-source R package on CRAN ([incgraph](#)). The development version is also available on GitHub ([rcannood/incgraph](#)).

Adapted from:

**Cannoodt, R.**, Ruyssinck, J., Ramon, J., De Preter, K., and Saeys, Y. IncGraph: Incremental graphlet counting for topology optimisation. *PLOS ONE* 13, 4 (2018), e0195997. doi:10.1371/journal.pone.0195997

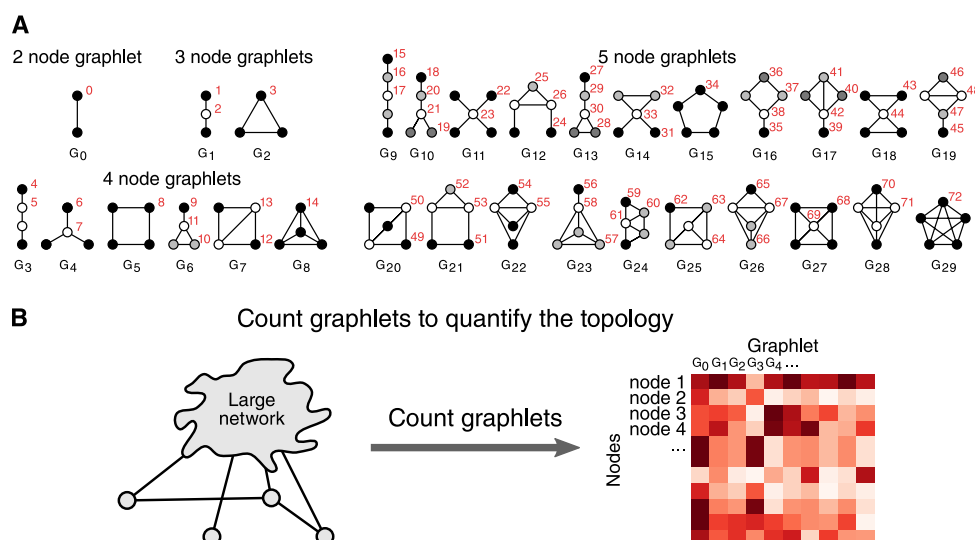
## Introduction

Even the simplest of living organisms already consist of complex biochemical networks which must be able to respond to a variety of stressful conditions in order to survive. An organism can be characterised using numerous interaction networks, including gene regulation, metabolic, signalling, and protein-protein interaction. The advent of high-throughput profiling methods (e.g. microarrays and RNA sequencing) have allowed to observe the molecular contents of a cell, which in turn have enabled computational network inference methods to reverse engineer the biochemical interaction networks [1]. Improving the accuracy of inferred networks has been a long-standing challenge, but the development of ever more sophisticated algorithms and community-wide benchmarking studies have resulted in significant process [2, 3, 4, 5].

Several recent developments involve incorporating topological priors, either to guide the inference process [6] or post-process the network [7]. A common prior is that biological networks are highly modular [8], as they consist of multiple collections of functionally or physically linked molecules [9, 10]. At the lowest level, each module is made up out of biochemical interactions arranged in small topological patterns, which act as fundamental building blocks [11].

Graphlets [12] are one of the tools which could be used to add a topological prior to a biological network, Graphlets are small connected subnetworks which can be counted to identify which low-level topological patterns are present in a network. By comparing how topologically similar a predicted network is to what would be expected of a true biological network, a predicted network can be optimised in order to obtain a better topology.

By counting the number of occurrences of each of the different graphlets (Fig 6.1A) touching a specific node, one can characterise the topology surrounding it. The graphlet counts of a network can be represented as a matrix with one row for each of the nodes and one column for each of the graphlets (Fig 6.1B). An orbit represents a class of isomorphic (i.e. resulting in the same structure) positions of nodes within a graphlet (Fig 6.1A, coloured in red). Both graphlets and orbits have been used extensively for predicting the properties of nodes such as protein functionality [13, 14, 15] and gene oncogenicity [16], by performing network alignment [17, 18] or using them as a similarity measure in machine learning tasks [19, 20].



**Figure 6.1: Graphlet counting in a network characterises its local topologies.** (A) In total, there are 30 different graphlets containing 2 to 5 nodes, ranging from  $G_0$  to  $G_{29}$ . Orbits are an extension of graphlets which also take into account the position of a node within a graphlet. The 73 different orbits are coloured in red. (B) By counting the occurrences of these graphlets in the network, the local topology surrounding a node can be quantified.

In this work, we focus on optimising gene regulatory networks by incorporating a topological

prior as part of a topology optimisation process. We seek to optimise a predicted network by iteratively modifying the network and accepting modifications that lead to topological properties that resemble better those of true biological networks.

However, using graphlets to perform topology optimisation was hitherto not possible. Calculating the graphlet counts using the most state-of-the-art graphlet counting of a moderately sized gene regulatory network already has an execution time of about five seconds (*E. coli*,  $\sim 3000$  genes,  $\sim 10000$  interactions, up to graphlets up to 5 nodes). While this computational time poses no issue for regular static networks, recalculating all graphlet counts for every network modification made as part of a topology optimisation is computationally intractable. For example, performing 100'000 iterations of topology optimisation on a similarly sized network and calculating the topological impact of 10 possible edge modification at each iteration, already results in a computational time of about 12 days. Graphlet-based topology optimisation thus quickly becomes infeasible for larger networks.

When adding or removing an edge to a large network, the number of altered graphlets is much smaller than the total number of graphlets in the network. It could therefore be much more efficient to iterate over and count all the graphlets that have been added or removed as a result of the edge modification, than it is to recalculate the graphlet counts from scratch.

Eppstein et al. introduced data structures and algorithms for updating the counts of size-three[21] and size-four[22] subgraphs in a dynamic setting. The data structures were determined such that the amortised time is  $O(h)$  and  $O(h^2)$ , respectively, where  $h$  is the h-index of the network[23].

We developed IncGraph, an alternative algorithm and implementation for performing incremental counting of graphlets up to size five. We show empirically that IncGraph is several orders of magnitude faster at calculating the differences in graphlet counts in comparison to non-incremental counting approaches. In addition, we demonstrate the practical applicability by developing a graphlet-based optimisation criterion for biological networks.

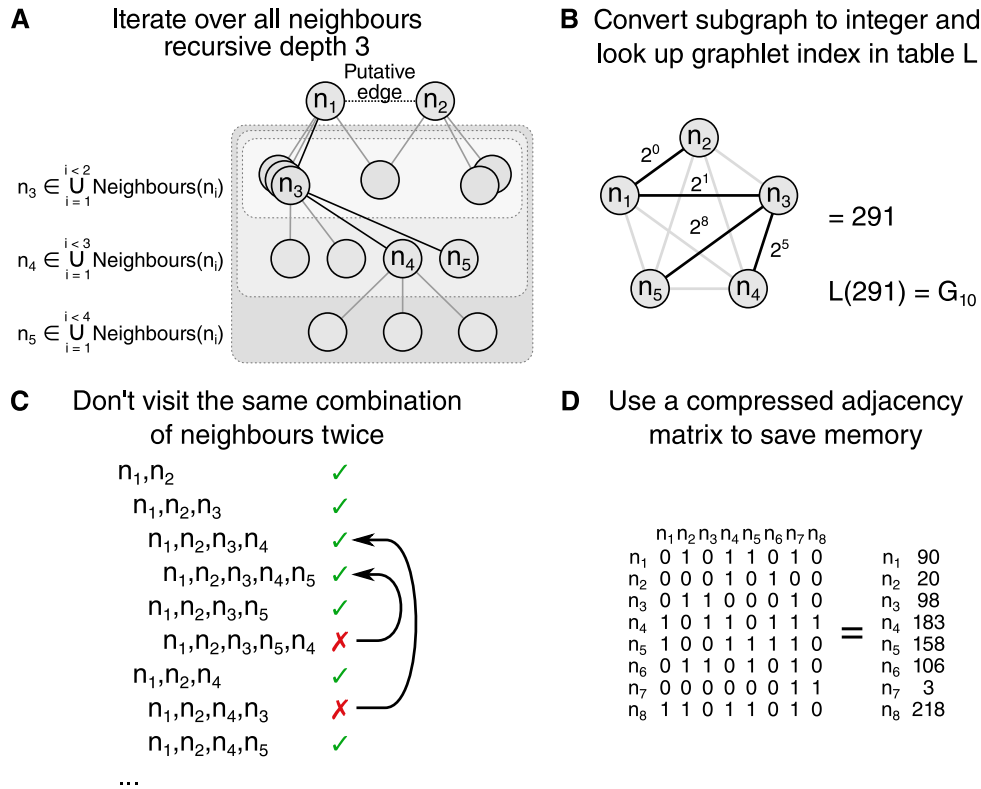
## Materials and methods

Assume a network  $G$  of which the graphlet counts  $C_G$  are known.  $C_G$  is an  $n$ -by- $m$  matrix, with  $n$  the number of vertices in the network,  $m = 73$  is the number of different orbits, and where  $C_G[i, j]$  is the number of times node  $i$  is part of a graphlet at orbit  $O_j$ . Further assume that one edge has either been added or removed from  $G$ , resulting in  $G'$ , at which point the counts  $C_{G'}$  need to be observed. If multiple edges have been modified, the method described below can be repeated for each edge individually.

### Incremental graphlet counting

As stated earlier, recalculating the graphlet counts for each modification made to the network quickly becomes computationally intractable for larger network sizes. However, as the differences in topology between  $G$  and  $G'$  are small, it is instead possible to calculate the differences in graphlet counts  $\Delta_{G,G'}$  instead. This is much more efficient to calculate, as only the neighbourhood of the modified edges needs to be explored.  $C_{G'}$  can thus be calculated as  $C_{G'} = C_G + \Delta_{G,G'}$ . The differences in graphlet counts  $\Delta_{G,G'}$  are calculated by iterating recursively over the neighbours surrounding each of the modified edges (See S1 Pseudocode). Several strategies are used in order to calculate  $\Delta_{G,G'}$  as efficiently as possible (Fig 6.2). (A) The delta matrix is calculated separately for each modified edge. Since the edge already contains two out of five nodes and any modified graphlet is a connected subgraph, the neighbourhood of this edge only needs to be explored up to depth 3 in order to iterate over all modified graphlets. (B) We propose a lookup table to look up the graphlet index of each node of a given subgraph. By weighting each possible edge in a 5-node graphlet, the sum of the edges of a subgraph can be used to easily look up

the corresponding graphlet index. (C) During the recursive iteration of the neighbourhood, the same combination of nodes is never visited twice. (D) As the network can be relatively large, the adjacency matrix is binary compressed in order to save memory. One integer requires 4 bytes and contains the adjacency boolean values of 32 edges, whereas otherwise 32 booleans would require 32 bytes. For example, 1GB of memory is large enough to store a compressed adjacency matrix of 92681 nodes. If the network contains too many nodes to fit a compressed adjacency matrix into the memory, a list of sets containing each node's neighbours is used instead.



**Figure 6.2:** Several strategies are employed in order to reduce time and memory consumption.

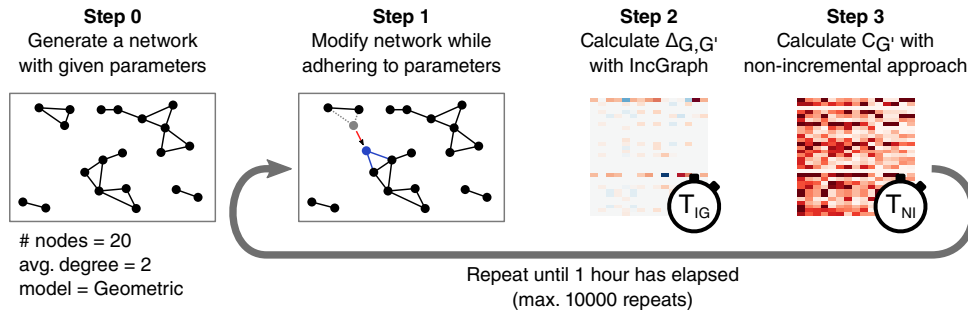
(A) Only the depth 3 neighbourhood of each modified edge needs to be explored in order to have visited all modified five-node graphlets. (B) A lookup table can be used to easily look up the graphlet index of a subgraph, by weighing each edge in a 5-node subgraph by a power of 2. (C) The same combination of five nodes is never repeated, as to avoid counting the same graphlet multiple times. (D) The adjacency matrix of the network is compressed in order to reduce memory usage.

IncGraph supports counting graphlets and orbits of subgraphs up to five nodes in undirected networks. By modifying the lookup table, the method can be easily extended to directed graphlets or larger-node graphlets, or limited to only a selection of graphlets. This allows for variations of the typical graphlets to be studied in an incremental setting.

## Timing experiments

We compared the execution time needed to calculate the graphlet counts in iteratively modified networks between our method and a state-of-the-art non-incremental algorithm, Orca [24]. Orca is a heavily optimised algorithm for counting 5-node graphlets in static networks, and outperforms all other static graphlet counting algorithms by an order of magnitude [24].

The timing experiments were performed by generating networks from 3 different network models, making modifications to those networks while still adhering to the network model, and measuring the execution times taken for both approaches to calculate the new graphlet counts (Fig 6.3). The network models were based on three static network models: Barabási-Albert [25], Erdős-Rényi [26], and Geometric [27]. Pseudo code for these random evolving network models can



**Figure 6.3:** Static network model generators were modified to generate dynamic networks.

Three network models were used: BarabásiAlbert, ErdsRényi, and Geometric. Step 0: a network is generated according to the network model and the given parameters. Step 1: the network is modified such that the result is as likely to have been generated by the network model. Step 2: The time for IncGraph to calculate the differences in graphlet counts is measured ( $T_{IG}$ ). Step 3: The time for the non-incremental approach to calculate the graphlet counts of the new network is measured ( $T_{NI}$ ). Steps 1 to 3 are repeated until all modifications generated at step 0 are processed, or until an execution time threshold has been reached.

be found in S2 Pseudocode, S3 Pseudocode, and S4 Pseudocode respectively. Each model generates an initial network according to the static network model it is based on, and a list of network modifications (removing an edge from or adding an edge to the network). These network modifications are made such that at any given time point in the evolving network, it is likely that the network at its current state could have been generated by the original static network model.

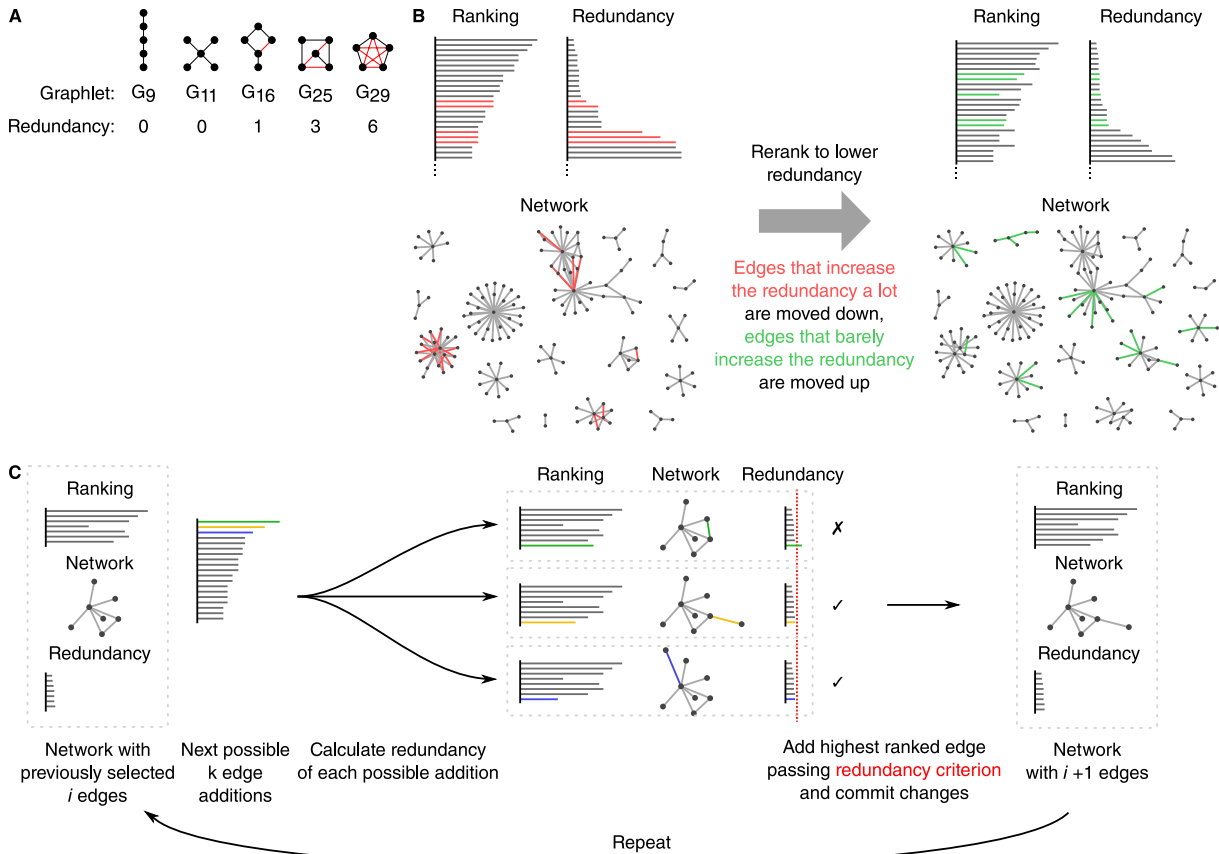
Networks were generated with varying network models, between 1000 and 16000 nodes, node degrees between 2 and 64, and 10000 time points. We measured the time needed to calculate the delta matrix at random time points until 1 hour has passed. All timings experiments were carried out on Intel(R) Xeon(R) CPU E5-2665 @ 2.40GHz processors, with one thread per processor. The generation of networks with higher node counts or degrees was constrained by the execution time of the network generators, not by IncGraph. All data and scripts are made available at [github.com/rcannood/incgraph-scripts](https://github.com/rcannood/incgraph-scripts).

## Gene regulatory network optimisation experiments

We demonstrate the usefulness of IncGraph by using a simple graphlet-based metric to optimise gene regulatory networks. One of the striking differences between real and predicted gene regulatory networks is that the predicted networks contain highly connected subnetworks, which contain high amounts of false positives. We determine a penalty score such that predicted networks containing graphlets with many redundant edges will be penalised in comparison to very sparse networks.

The *redundancy penalty* (Fig 6.4A) of a network is defined as the sum of occurrences of each graphlet multiplied by the redundancy associated with each individual graphlet. The redundancy of a graphlet is the number of edges that can be removed without disconnecting the nodes from one another. By using the redundancy penalty score, we aim to improve the gene regulatory network (Fig 6.4B).

The topology optimisation procedure uses an empty network as initialisation and grows the network by selecting interactions iteratively. Each iteration, the top  $k = 100$  highest ranked interactions that are not currently part of the network are evaluated, and the highest ranked interaction passing the redundancy criterion is selected (Fig 6.4C). This procedure is repeated until a predefined amount of time has passed. As the aim of this experiment is not to obtain the highest performing topology optimisation method, parameter optimisation of  $k$  has not been performed and is considered to be outside the scope of this work.



**Figure 6.4:** Predicted gene regulatory networks of model organisms are optimised to reduce the false positive rate. A) The number of redundant edges in each graphlet are counted. B) The network is optimised in order to obtain a lower redundancy over time. Two networks are shown, one before and one after the optimisation procedure. Edges coloured in red have been removed from the network after optimisation, green edges have been added. C) Starting from an empty network, the interactions are modified by iteratively evaluating the increase in redundancy of the next  $k$  interactions, and adding the first edge for which its redundancy is less than the 90<sup>th</sup> percentile redundancy.

We optimised gene regulatory networks of *E. coli* and *S. cerevisiae*. The predicted networks were generated using the network inference method GENIE3 [28] with default parameters. Gene expression data was obtained from COLOMBOS [29] and GEO [30], respectively. The predicted networks and the optimised versions thereof were compared against respective lists of known gene regulatory interactions [31, 32].

## Results and discussion

The contributions of this work are twofold. Firstly, we propose a new method for incrementally calculating the differences in graphlet counts in changing graphs, and show that it is orders of magnitude faster than non-incremental approaches. Secondly, we demonstrate its applicability by optimising a predicted gene regulatory network in order to reduce the false positive rate therein.

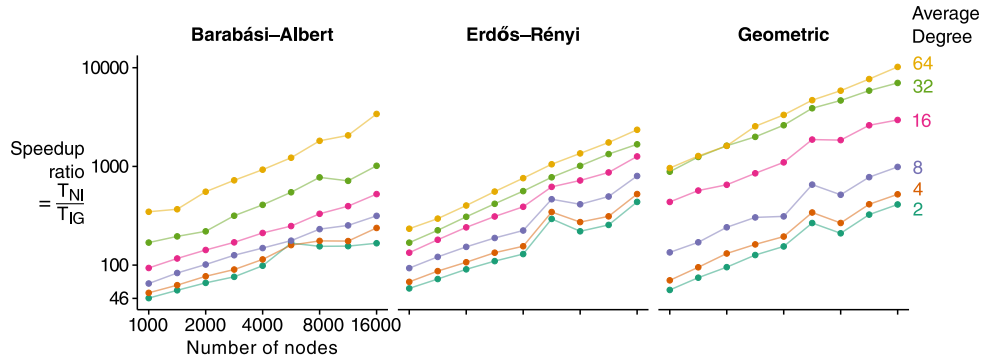
### Execution time is reduced by orders of magnitude

Timing experiments show that IncGraph is significantly faster in calculating the delta matrix in comparison to calculating the graphlet counts from scratch at each iteration (Fig 6.5). The observed speedup ratios between IncGraph and the non-incremental approach Orca ranges from



about  $50\times$  to  $10000\times$ . The speedup ratio increases as the network size increases. For larger networks, IncGraph can thus calculate the delta matrices of 10000 edge modifications while the non-incremental approach calculates one graphlet count matrix.

Surprisingly, IncGraph obtains higher execution times for networks with 5657 nodes than for networks with 8000 nodes. One possible explanation is that the size of the data structures containing those networks are particularly favourable in avoiding cache misses. Confirmation of this explanation, however, would require further investigation.



**Figure 6.5: IncGraph is significantly faster than non-incremental approaches.** For small networks, the execution time of IncGraph  $T_{IG}$  is already 50 times less than that of non-incremental approaches  $T_{NI}$ . This ratio increases even further for networks with higher numbers of nodes or higher average degrees.

Comparing the execution time of IncGraph to the h-index of the networks indicates that the amortised time of IncGraph could be  $O(h^3)$  (S1 Fig). This is in line with the amortised times  $O(h)$  and  $O(h^2)$  of the algorithm described by Eppstein et al.[22] for counting three-size and four-size subgraphs respectively.

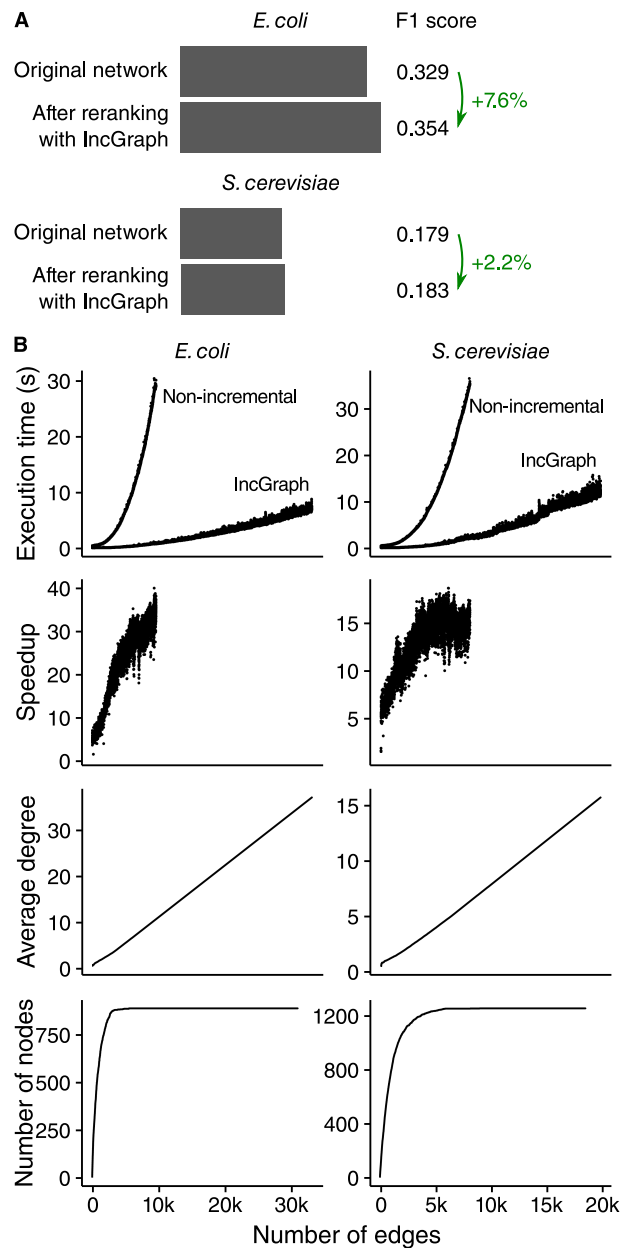
## IncGraph allows for better regulatory network optimisation

We implemented a graphlet-based optimisation algorithm for improving the false positive rate of the predicted gene regulatory networks of *E. coli* and *S. cerevisiae*. After reranking the regulatory interactions, the F1 score of the first 1000 interactions had increased by 7.6% and 2.2% respectively (Fig 6.6A). The obtained speedup of about  $15\text{--}30\times$  (Fig 6.6B) is in line with the experiments on *in silico* networks. Namely, for the *E. coli* network at 9618 interactions and 889 nodes (average degree = 10.8), a speedup of about  $30\times$  was obtained. Similarly, for the *S. cerevisiae* network at 8013 interactions and 1254 nodes (average degree = 6.4), a speedup of about  $15\times$  was obtained. These speedups are in the same order of magnitude of similarly sized networks (1000 nodes and 8000 interactions) generated with a Barabási-Albert model, with a speedup of  $65\times$ . This is to be expected, as such networks share the same scale-free property that gene regulatory networks have.

## Conclusion

Many improvements over the past few years have resulted in efficient graphlet counting algorithms, even for large static networks. However, needing to perform the simplest of tasks tens of thousands of times quickly becomes computationally intractable. As such, recalculating the graphlet counts of a network after each of the many network modification is infeasible.

This study introduces a method for calculating the differences in graphlet (and orbit) counts, which we call incremental graphlet counting or IncGraph for short. We show that IncGraph is at least 10-100 times faster than non-incremental methods for networks of moderate size, and that the speedup ratio increases even further for larger networks. To demonstrate the applicability of IncGraph, we optimised a predicted gene regulatory network by enumerating over the ranked



**Figure 6.6: A simple graphlet-based scoring method improves predicted regulatory networks.**

(A) The F1 score was calculated by calculating the harmonic mean of the AUROC and AUPR scores of the first 1000 interactions. (B) IncGraph is significantly faster than the non-incremental approach. Note that for each interaction added to the network, the graphlet counts of 100 putative interactions were evaluated.

edges and observing the graphlet counts of several candidate edges before deciding which edge to add to the network.

IncGraph enables graphlet-based metrics to characterize online networks, e.g. in order to track certain network patterns over time, as a similarity measure in a machine learning task, or as a criterion in a topology optimisation.

## Supporting information

**S1 Pseudocode.** IncGraph calculates  $\Delta_{G,G'}$  using a strict branch-and-bound strategy.

**S2 Pseudocode.** Pseudo code for generating an evolving Barabási-Albert (BA) network. It first generates a BA network, and then generates  $o$  operations such that at any time point, the network is or very closely resembles a BA network.

**S3 Pseudocode.** Pseudo code for generating an evolving Erdős-Rényi (ER) network. It first generates an ER network, and then generates  $o$  operations such that at any time point, the network is or very closely resembles an ER network.

**S4 Pseudocode.** Pseudo code for generating an evolving geometric network. It first generates a geometric network, and then generates  $o$  operations such that at any time point, the network is or very closely resembles a geometric network.

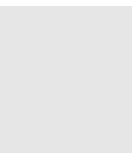
**S1 Fig.** Empirical measurements show a strong relation between the execution time of IncGraph and the h-index cubed of the network it was applied to. This is in line with the findings by Eppstein et al., where counting 3-size subgraphs has an amortised time of  $O(h)$  and counting 4-size subgraphs has an amortised time of  $O(h^2)$ .



Hello, here is some text without a meaning. This text should show what a printed text will look like at this place. If you read this text, you will get no information. Really? Is there no information? Is there a difference between this text and some nonsense like “Huardest gefburn”? Kjift – not at all! A blind text like this gives you information about the selected font, how the letters are written and an impression of the look. This text should contain all letters of the alphabet and it should be written in of the original language. There is no need for special content, but the length of words should match the language.

### **7.1 Overview and conclusions of the presented work**

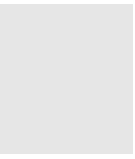
### **7.2 Future research directions**



---

## Samenvatting

---

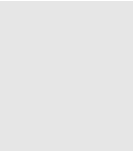




---

## Summary

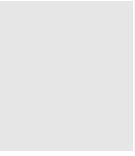
---



---

## List of Publications

---



---

## Bibliography

---

- [1] R. Albert. “Network Inference, Analysis, and Modeling in Systems Biology”. In: *the Plant Cell Online* 19.11 (2007), pp. 3327–3338. ISSN: 1040-4651. DOI: 10.1105/tpc.107.054700.
- [2] Daniel Marbach et al. “Revealing strengths and weaknesses of methods for gene network inference”. In: *Proceedings of the National Academy of Sciences* 107.14 (Apr. 2010), pp. 6286–6291. ISSN: 1091-6490. DOI: 10.1073/pnas.0913357107. URL: <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=2851985%7B%5C%7Dtool=pmcentrez%7B%5C%7Drendertype=abstract>.
- [3] Varun Narendra et al. “A comprehensive assessment of methods for de-novo reverse-engineering of genome-scale regulatory networks”. In: *Genomics* 97.1 (2011), pp. 7–18. ISSN: 08887543. DOI: 10.1016/j.ygeno.2010.10.003. URL: <http://dx.doi.org/10.1016/j.ygeno.2010.10.003>.
- [4] Daniel Marbach et al. “Wisdom of crowds for robust gene network inference”. In: *Nat Meth* 9.8 (Aug. 2012), pp. 796–804. ISSN: 1548-7091. DOI: 10.1038/nmeth.2016. URL: <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3512113%7B%5C%7Dtool=pmcentrez%7B%5C%7Drendertype=abstract>.
- [5] Tarmo Äijö and Richard Bonneau. “Biophysically motivated regulatory network inference: Progress and prospects”. In: *Human Heredity* 81.2 (2017), pp. 62–77. ISSN: 14230062. DOI: 10.1159/000446614.
- [6] Fabrício M. Lopes et al. “A feature selection technique for inference of graphs from their known topological properties: Revealing scale-free gene regulatory networks”. In: *Information Sciences* 272 (2014), pp. 1–15. ISSN: 00200255. DOI: 10.1016/j.ins.2014.02.096.
- [7] Joeri Ruysinck et al. “Netter: re-ranking gene network inference predictions using structural network properties.” In: *BMC Bioinformatics* 17.1 (2016), p. 76. ISSN: 1471-2105. DOI: 10.1186/s12859-016-0913-0. URL: <http://bmcbioinformatics.biomedcentral.com/articles/10.1186/s12859-016-0913-0>.
- [8] Alexander W Rives and Timothy Galitski. “Modular organization of cellular networks.” In: *Proceedings of the National Academy of Sciences of the United States of America* 100.3 (2003), pp. 1128–33. ISSN: 0027-8424. DOI: 10.1073/pnas.0237338100. URL: <http://www.ncbi.nlm.nih.gov/pubmed/12538875%7B%5C%7D5Cnhttp://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=PMC298738>.
- [9] L H Hartwell et al. “From molecular to modular cell biology.” In: *Nature* 402.6761 Suppl (1999), pp. C47–C52. ISSN: 0028-0836. DOI: 10.1038/35011540. arXiv: 05218657199780521865715.

- [10] a Barabasi et al. “Network biology: understanding the cell’s functional organization.” In: *Nature reviews. Genetics* 5.2 (Feb. 2004), pp. 101–13. ISSN: 1471-0056. DOI: 10.1038/nrg1272. URL: <http://www.ncbi.nlm.nih.gov/pubmed/14735121>.
- [11] R Milo et al. “Network motifs: simple building blocks of complex networks.” In: *Science (New York, N.Y.)* 298.2002 (2002), pp. 824–827. ISSN: 00368075. DOI: 10.1126/science.298.5594.824. arXiv: 0908.1143v1. URL: <http://www.sns.ias.edu/%7B%5C%7D7B%7B~%7D%7B%5C%7D7Dtlusty/courses/InfoInBio/Papers/AlonMotifs2002.pdf%20http://www.sns.ias.edu/%7B~%7Dtlusty/courses/InfoInBio/Papers/AlonMotifs2002.pdf>.
- [12] N Przulj et al. “Modeling interactome: scale-free or geometric?” In: *Bioinformatics (Oxford, England)* 20.18 (Dec. 2004), pp. 3508–15. ISSN: 1367-4803. DOI: 10.1093/bioinformatics/bth436. URL: <http://www.ncbi.nlm.nih.gov/pubmed/15284103>.
- [13] Tijana Milenkovi, Nataa Prulj, and Natasa Przulj. “Uncovering biological network function via graphlet degree signatures.” In: *Cancer informatics* 6 (Jan. 2008), pp. 257–73. ISSN: 1176-9351. URL: <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=2623288%7B%5C%7Dtool=pmcentrez%7B%5C%7Drendertype=abstract>.
- [14] Cortnie Guerrero et al. “Characterization of the proteasome interaction network using a QTAX-based tag-team strategy and protein interaction network analysis.” In: *Proceedings of the National Academy of Sciences of the United States of America* 105.36 (2008), pp. 13333–13338. ISSN: 0027-8424. DOI: 10.1073/pnas.0801870105.
- [15] Omkar Singh, Kunal Sawariya, and Polamarasetty Aparoy. “Graphlet signature-based scoring method to estimate protein-ligand binding affinity.” In: *Royal Society open science* 1.4 (2014), p. 140306. ISSN: 2054-5703. DOI: 10.1098/rsos.140306. URL: <http://www.ncbi.nlm.nih.gov/sites/entrez?Db=pubmed%7B%5C%7DDbFrom=pubmed%7B%5C%7DCmd=Link%7B%5C%7DLinkName=pubmed%7B%5C%7Dpubmed%7B%5C%7DLinkReadableName=Related%20Articles%7B%5C%7DIdsFromResult=26064572%7B%5C%7Dordinalpos=3%7B%5C%7Ditool=EntrezSystem2.PEntrez.Pubmed.Pubmed%7B%5C%7DResultsPanel.Pubmed%7B%5C%7DRVDocSum%7B%5C%7D5Cnhttp://www.ncbi>.
- [16] Tijana Milenkovi et al. “Systems-level cancer gene identification from protein interaction network topology applied to melanogenesis-related functional genomics data.” In: *Journal of the Royal Society, Interface* 7 (2010), pp. 423–437. ISSN: 1742-5662. DOI: 10.1098/rsif.2009.0192.
- [17] Oleksii Kuchaiev et al. “Topological network alignment uncovers biological function and phylogeny.” In: *Journal of the Royal Society, Interface / the Royal Society* 7.50 (2010), pp. 1341–1354. ISSN: 1742-5662. DOI: 10.1098/rsif.2010.0063. arXiv: 0810.3280.
- [18] T Milenkovi, H Zhao, and FE Faisal. “Global network alignment in the context of aging”. In: *Proceedings of the International Conference on Bioinformatics, Computational Biology and Biomedical Informatics* (2013), pp. 23–32. URL: <http://dl.acm.org/citation.cfm?id=2508968>.
- [19] Nino Shervashidze et al. “Efficient graphlet kernels for large graph comparison”. In: *AISTATS* 5 (2009), pp. 488–495. URL: <http://machinelearning.wustl.edu/mlpapers/paper%7B%5C%7Dfiles/AISTATS09%7B%5C%7DServashidzeVPMB.pdf>.
- [20] Vladimir Vacic et al. “Graphlet kernels for prediction of functional residues in protein structures.” In: *Journal of computational biology : a journal of computational molecular cell biology* 17.1 (2010), pp. 55–72. ISSN: 1557-8666. DOI: 10.1089/cmb.2009.0029. URL: <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=2921594%7B%5C%7Dtool=pmcentrez%7B%5C%7Drendertype=abstract>.

- [21] David Eppstein and Emma S. Spiro. “The h-index of a graph and its application to dynamic subgraph statistics”. In: *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*. Vol. 5664 LNCS. Springer-Verlag, 2009, pp. 278–289. ISBN: 3642033660. DOI: 10.1007/978-3-642-03367-4\_25. arXiv: 0904.3741. URL: [http://link.springer.com/10.1007/978-3-642-03367-4\\_25](http://link.springer.com/10.1007/978-3-642-03367-4_25).
- [22] David Eppstein et al. “Extended dynamic subgraph statistics using h-index parameterized data structures”. In: *Theoretical Computer Science* 447 (Aug. 2012), pp. 44–52. ISSN: 03043975. DOI: 10.1016/j.tcs.2011.11.034. arXiv: 1009.0783. URL: <https://www.sciencedirect.com/science/article/pii/S0304397511009534>.
- [23] J E Hirsch. “An index to quantify an individual’s scientific research output”. In: *Proceedings of the National Academy of Sciences of the United States of America* 102.46 (Nov. 2005), pp. 16569–72. ISSN: 0027-8424. DOI: 10.1073/pnas.0507655102. arXiv: 0508025 [physics]. URL: <http://www.ncbi.nlm.nih.gov/pubmed/16275915><http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=PMC1283832><http://arxiv.org/abs/physics/0508025><http://dx.doi.org/10.1073/pnas.0507655102>.
- [24] Toma Hoever and Janez Demar. “A combinatorial approach to graphlet counting.” In: *Bioinformatics (Oxford, England)* 30.4 (Feb. 2014), pp. 559–65. ISSN: 1367-4811. DOI: 10.1093/bioinformatics/btt717. URL: <http://www.ncbi.nlm.nih.gov/pubmed/24336411>.
- [25] Réka Albert and Albert Laszlo Barabasi. “Statistical mechanics of complex networks”. In: *Reviews of Modern Physics* 74. January (2002), pp. 47–97. ISSN: 1478-3967. DOI: 10.1088/1478-3967/1/3/006. arXiv: 0106096v1 [arXiv:cond-mat]. URL: <http://dx.doi.org/10.1103/RevModPhys.74.47>.
- [26] P. Erds and A Rényi. “On random graphs”. In: *Publicationes Mathematicae* 6 (1959), pp. 290–297. ISSN: 00029947. DOI: 10.2307/1999405.
- [27] M J B Appel and R P Russo. “The minimum vertex degree of a graph on uniform points in  $[0,1]^d$ ”. In: *Adv. in Appl. Probab.* 29.3 (1997), pp. 582–594. ISSN: 00018678.
- [28] Vân Anh Huynh-Thu et al. “Inferring Regulatory Networks from Expression Data Using Tree-Based Methods”. In: *PLoS ONE* 5.9 (Jan. 2010), e12776. ISSN: 1932-6203. DOI: 10.1371/journal.pone.0012776. URL: <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=2946910><http://www.ncbi.nlm.nih.gov/pmc/articles/PMC2946910/><http://dx.doi.org/10.1371/journal.pone.0012776>.
- [29] Marco Moretto et al. “COLOMBOS v3.0: Leveraging gene expression compendia for cross-species analyses”. In: *Nucleic Acids Research* 44.D1 (2016), pp. D620–D623. ISSN: 13624962. DOI: 10.1093/nar/gkv1251.
- [30] R. Edgar. “Gene Expression Omnibus: NCBI gene expression and hybridization array data repository”. In: *Nucleic Acids Research* 30.1 (2002), pp. 207–210. ISSN: 13624962. DOI: 10.1093/nar/30.1.207. URL: <https://academic.oup.com/nar/article-lookup/doi/10.1093/nar/30.1.207>.
- [31] Socorro Gama-Castro et al. “RegulonDB version 9.0: High-level integration of gene regulation, coexpression, motif clustering and beyond”. In: *Nucleic Acids Research* 44.D1 (2016), pp. D133–D143. ISSN: 13624962. DOI: 10.1093/nar/gkv1156.
- [32] Sisi Ma et al. “De-novo learning of genome-scale regulatory networks in *S. cerevisiae*”. In: *PLoS ONE* 9.9 (2014). ISSN: 19326203. DOI: 10.1371/journal.pone.0106479.