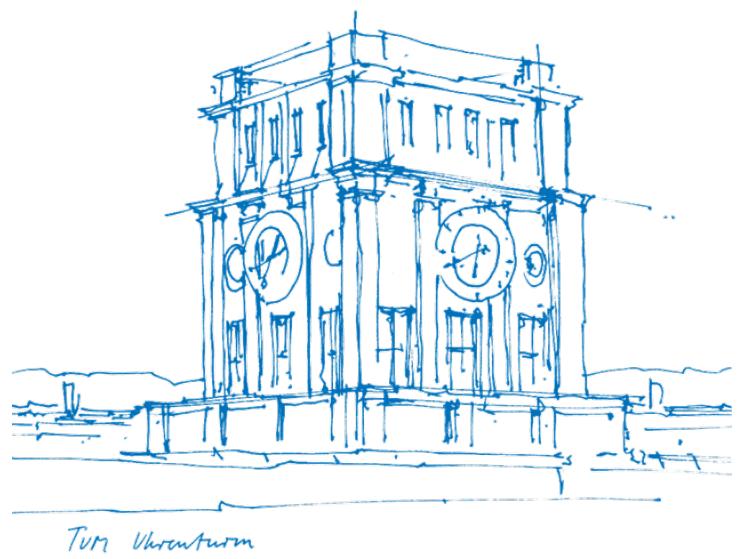


Thesis title

Subtitle of the thesis

Juan David Henao Sanchez



Thesis title

Subtitle of the thesis

Juan David Henao Sanchez

Thesis title

Subtitle of the thesis

Juan David Henao Sanchez

Vollständiger Abdruck der von der Fakultät für Elektrotechnik und Informationstechnik der Technischen Universität München zur Erlangung des akademischen Grades eines

Doktor-Ingenieurs (Dr.-Ing.)

genehmigten Dissertation.

Vorsitzende(r):

Prof. Franz X. Gabelsberger

Prüfer der Dissertation:

1. Prof. Dr. Georg Simon Ohm
2. Prof. James Clerk Maxwell

Die Dissertation wurde am 29.04.2016 bei der Technischen Universität München eingereicht und durch die Fakultät für Elektrotechnik und Informationstechnik am 11.07.2016 angenommen.

To Franz X. Gabelsberger, inventor of the street named after him.

Abstract

The abstract of your thesis goes here.

 Lorem ipsum dolor sit amet, consectetur adipiscing elit. Ut purus elit, vestibulum ut, placerat ac, adipiscing vitae, felis. Curabitur dictum gravida mauris. Nam arcu libero, nonummy eget, consectetur id, vulputate a, magna. Donec vehicula augue eu neque. Pellentesque habitant morbi tristique senectus et netus et malesuada fames ac turpis egestas. Mauris ut leo. Cras viverra metus rhoncus sem. Nulla et lectus vestibulum urna fringilla ultrices. Phasellus eu tellus sit amet tortor gravida placerat. Integer sapien est, iaculis in, pretium quis, viverra ac, nunc. Praesent eget sem vel leo ultrices bibendum. Aenean faucibus. Morbi dolor nulla, malesuada eu, pulvinar at, mollis ac, nulla. Curabitur auctor semper nulla. Donec varius orci eget risus. Duis nibh mi, congue eu, accumsan eleifend, sagittis quis, diam. Duis eget orci sit amet orci dignissim rutrum.

 Nam dui ligula, fringilla a, euismod sodales, sollicitudin vel, wisi. Morbi auctor lorem non justo. Nam lacus libero, pretium at, lobortis vitae, ultricies et, tellus. Donec aliquet, tortor sed accumsan bibendum, erat ligula aliquet magna, vitae ornare odio metus a mi. Morbi ac orci et nisl hendrerit mollis. Suspendisse ut massa. Cras nec ante. Pellentesque a nulla. Cum sociis natoque penatibus et magnis dis parturient montes, nascetur ridiculus mus. Aliquam tincidunt urna. Nulla ullamcorper vestibulum turpis. Pellentesque cursus luctus mauris.

Contents

Abstract	ix
1 Introduction	1
1.1 Pathophysiology of chronical lung diseases	1
1.1.1 Bronchopulmonary Dysplasia (BPD)	1
1.1.2 Asthma	1
1.1.3 Chronic Obstructive Pulmonary Disease (COPD)	1
1.1.4 Idiopathic Pulmonary Fibrosis (IPF)	1
1.2 Computational biology and chronic lung diseases	1
1.2.1 Multi-omics data integration	1
1.2.2 Clinical prediction	1
1.2.3 Systems biology	1
1.3 Aims	1
2 Methodology	3
2.1 Data gathering	4
2.1.1 Mice data	4
2.1.2 Human data	4
2.1.3 Public data	4
2.2 Preprocessing	4
2.2.1 Normalization	4
2.2.2 Data imputation	4
2.2.3 Batch-effect detection	4
2.3 Differential expression analysis	4
2.3.1 Limma	4
2.3.2 DESeq2	4
2.4 Enrichment analysis	4
2.4.1 Gene list functional enrichment analysis	4
2.5 Multi-omics factor analysis (MOFA)	4
2.6 Clinical data correlation	5
2.6.1 Linear regression	5
2.6.2 Binomial regression	5
2.6.3 Ordinal regression	5
2.6.4 Multinomial logistic regression	5
2.6.5 Dirichlet regression	5
2.7 Benchmarking of Lasso models dealing with missing values	5
2.7.1 Knowledge guided multi-level network inference	5
2.7.2 Two-steps based models	5
2.7.3 Inverse covariance based methods	5
2.8 Adult data correlation	5
2.8.1 Random forest	5
2.8.2 t-test, manova, log-reg	5
3 Summary of publications	7
4 Discussion	11

1 Introduction

This is the introduction of the thesis.

1.1 Pathophysiology of chronical lung diseases

1.1.1 Bronchopulmonary Dysplasia (BPD)

Bronchopulmonary Dysplasia

Definition Context Bronchopulmonary Dysplasia (BPD) is the most common complication of prematurity, characterized by the impaired development of the gas exchange area and the bronchial tree. Body Wrap Diagnosis + severity Context (x) The diagnosis of BPD has been dynamically changed over time since it was first reported in 1967 by Northway, Rosen, and Porter. However, the main clinical variables behind its diagnosis remain the basis complemented with new measurements based on recent insights about its onset. Body (x) BPD was first described as a hyaline membrane disease in preterm infants provoked by mechanical ventilation without any ending expiratory pressure and high levels of oxygen supplementation [REF]. In 2001, a workshop organized by the National Institute of Child Health and Human Development (NICHD), the National Heart, Lung, and Blood Institute (NHLBI), and the Office of Rare Diseases (ORD) (NICHD/NHLBI/ORD workshop), improved the diagnosis as well as providing severity classification based on the criteria described in Table XX. Table XX. Diagnostic criteria of Bronchopulmonary Dysplasia* GA < 32 weeks GA >= 32 weeks Time point for assessment 36 weeks PMA of at discharge to home > 28 days and <56 days postnatal age or at discharge to home BPD diagnosis Treatment with more than 21Mild BPD Breathing room air at 36 weeks PMA or at discharge Breathing room air by 56 days postnatal age or at discharge Moderate BPD Need of less than 30Severe BPD Need of more or equal than 30* Table taken and modified from REF GA= gestational age; PMA= postmenstrual age; PPV= positive-pressure ventilation; NCPAP= nasal continuous positive airway pressure.

Although new modifications have been proposed recently, REF. The definition provided by NICHD/NHLBI/ORD workshop in 2001 has been widely used even nowadays to diagnose BPD. Wrap (x) Most of those proposed modifications remained restrictive to oxygen supplementation and mechanical ventilation usage or minimal addition of complementary clinical variables such as image-based readings REF. However, prenatal and postnatal risk variables have been widely studied as influential factors in triggering prematurity and potentially Bronchopulmonary Dysplasia as well. Antenatal and postnatal risk factors Context (x) An effective and efficient diagnosis of premature neonates with BPD has been a central matter of study even in recent years REF. However, detecting risk factors is essential to predict potential susceptibility to BPD before the formal clinical diagnosis. Those risk factors do not include only parameters after neonate delivery (postnatal) but clinical variables during the pregnancy process (prenatal) REF. Body (x) Prematurity is by itself the major risk factor of BPD, which could be divided into gestational age (GA) and the birth weight of the neonate; the risk of developing BPD increases inversely proportionate with the increase in both factors REF. Furthermore, other factors such as intrauterine growth restriction (IUGR), male sex, chorioamnionitis (infection of the placenta and the amniotic fluid), race or ethnicity, and smoking have been defined as risk factors for BPD REF. Likewise, other pre and postnatal variables, including comorbidities, have been evaluated as potential risk factors, including maternal age, antenatal (prenatal) and postnatal corticosteroid administration, surfactant administration, inhaled nitric oxide therapy, patent ductus arteriosus (PDA; a neonatal cardiovascular condition), pulmonary hypertension, intraventricular hemorrhage (IVH; a brain-vascular condition), periventricular leukomalacia (PVL; a type of neonatal brain injury), and retinopathy of prematurity (ROP; a neonatal optical-vascular condition) REF. Wrap (x) The heterogeneity of

clinical variables associated with BPD shed light on the complex molecular mechanisms behind the development and onset of neonatal chronic lung disease. Therefore, a deep molecular understanding, including the set of main biological pathways and molecular entities (genes, proteins and/or metabolites), seems necessary to improve the prevention and diagnosis powered by biomarkers and specific cellular-oriented treatments. Molecular description Context (x) The molecular description of any disease becomes a pivot for efficient diagnosis and well-oriented treatments. In the case of BPD, the heterogeneous clinical landscape encompasses a set of biological pathways that could be grouped into inflammation, oxygen toxicity, growth factor signaling, and extracellular matrix-related processes. Body (x) Inflammation is understood as a body's immunity reaction mainly characterized by tissue swelling REF. In BPD, the inflammation occurs at the lung tissue level by the release of pro-inflammatory proteins called cytokines and the presence of pro-inflammatory cells like neutrophils and monocytes REF. Besides, adaptive immune cells like CD4+ T-cells have been observed to activate T-cells by decreasing the expression of CD62L in infants with BPD REF. Among the risk factors associated previously with BPD, chorioamnionitis has been associated with antenatal immune activation by increasing levels of cytokines such as IL-6, IL-8, and TNF-a in fetal circulation REF. In addition, antenatal lung inflammation impacts a variety of molecular regulatory pathways, such as toll-like receptors 2 and 4 (TLR2 and TLR4), growth factors like TGF-b and CTGF, and mesenchymal structural proteins like bone morphogenetic protein-4 inducing vascular remodeling and alveolar simplification, phenotypes associated to mild BPD onset REF. Likewise, oxygen supplementation could cause inflammation and oxygen toxicity when elevated concentrations (>21

Hyperoxia is a well-known source of oxidative stress in premature neonates by the action of free radicals known as reactive oxygen species (ROS), which causes lung damage directly over epithelial and endothelial cells and proteins, metabolites, and nucleic acids destruction REF. ROS initiates apoptosis (programmed cell death) at the metabolic level through membrane lipids peroxidation, provoking the activation of the protein sphingomyelinase REF. The action of this protein releases large amounts of ceramides, which are the inductors of apoptosis REF. Likewise, cell death by ROS action could be achieved by activating proteins like key caspases and triggering receptor surface death receptors like Fas or mitochondrial cell death pathway activated by Bax proteins REF. On the other hand, ROS can directly oxidate nucleic acids, damaging the double-strain structure and causing cell death by necrosis or apoptosis REF. In BPD, increased levels of NOX1, a protein that produces superoxide radicals (a type of ROS), have been identified as a relevant participant in hyperoxia-induced acute oxidative stress injury, specifically by damaging the alveolar-capillary barrier REF. Both inflammation and hyperoxia-induced injuries affect internal cell functionality. However, extracellular processes could affect normal lung development in key regions such as alveolar and vascular tissues, inducing BPD REF.

The extracellular matrix (ECM) is the complex network of proteins and other molecules which offer structural and functional support to tissues REF. Collagen, the most abundant protein within the interstitial ECM, is composed mainly of collagen types I and III in the developing lung REF. Animal models have shown an increase in collagen 1 by the action of the transforming growth factor TGF-beta provoking thickened collagen fibres, increasing lung rigidity REF. This observation was confirmed through microscopic studies of BPD-diagnosed patients after positive-pressure ventilation REF. On the other hand, elastin, a component of elastic fibres in ECM, has been observed to decrease during impaired lung development, mainly affecting the correct alveolarisation process. Interestingly, the expression of the gene Eln (elastin) is stimulated by the expression of TGF-beta, which upregulation provokes rising levels of collagen 1 REF. Wrap (x) The different efforts to understand the molecular mechanisms behind BPD have diverged in complex and heterogenic biological pathways that potentially hampered normal lung development, causing the most relevant lung phenotypes of BPD, gas exchange area and vascularity affection. However, most of the experiments have been carried out on animal models, delaying the discovery of molecular markers for an effective diagnosis and well-established treatment design. Open problems Context Our increasing knowledge about BPD has improved the early treatment strategies efficiently, increasing the surveillance chances even in the most preterm infants (< 28 weeks PMA) REF. Nevertheless, this fact has opened new challenges, including dealing with comorbidities later in adolescence and/or adulthood REF. Furthermore, the current clinical-based BPD definition is not optimal for early diagnosis and severity prediction. Body

Nowadays, computational-based strategies are efficient and precise for prediction, treatment, and basic understanding-oriented tasks REF. The power of the current computational-based biomedical research is due to the capacity to join diverse types of data, including clinical records, diagnostic images (magnetic resonance image (MRI) or nuclear magnetic resonance (NMR), for example), and molecular data (e.g. gene or protein expression, and metabolite levels). This is possible thanks to the power of current machine learning and statistical techniques REF. Wrap

1.1.2 Asthma

1.1.3 Chronic Obstructive Pulmonary Disease (COPD)

1.1.4 Idiopathic Pulmonary Fibrosis (IPF)

1.2 Computational biology and chronic lung diseases

1.2.1 Multi-omics data integration

1.2.2 Clinical prediction

1.2.3 Systems biology

1.3 Aims

2 Methodology

This is the methodology of the thesis.

2.1 Data gathering

2.1.1 Mice data

Transcriptomics

2.1.2 Human data

Transcriptomics

Metabolomics

2.1.3 Public data

Multi-omics bulk data

Neonatal single-cell transcriptomics

2.2 Preprocessing

2.2.1 Normalization

DESeq2

Pareto scaling

Size-effect

2.2.2 Data imputation

Random-forest

knn-Imputation

2.2.3 Batch-effect detection

Principal component analysis (PCA)

Hierarchical clustering

K-BET

2.3 Differential expression analysis

2.3.1 Limma

2.3.2 DESeq2

2.4 Enrichment analysis

2.4.1 Gene list functional enrichment analysis

2.5 Multi-omics factor analysis (MOFA)

Lorem ipsum dolor sit amet, consectetuer adipiscing elit. Ut purus elit, vestibulum ut, placerat ac, adipiscing vitae, felis. Curabitur dictum gravida mauris. Nam arcu libero, nonummy eget, consectetuer id, vulputate a, magna. Donec vehicula augue eu neque. Pellentesque habitant morbi tristique senectus et netus et malesuada fames ac turpis egestas. Mauris ut leo. Cras viverra metus rhoncus sem. Nulla et lectus vestibulum urna fringilla ultrices. Phasellus eu tellus sit amet tortor gravida placerat. Integer sapien

est, iaculis in, pretium quis, viverra ac, nunc. Praesent eget sem vel leo ultrices bibendum. Aenean faucibus. Morbi dolor nulla, malesuada eu, pulvinar at, mollis ac, nulla. Curabitur auctor semper nulla. Donec varius orci eget risus. Duis nibh mi, congue eu, accumsan eleifend, sagittis quis, diam. Duis eget orci sit amet orci dignissim rutrum.

2.6 Clinical data correlation

2.6.1 Linear regression

2.6.2 Binomial regression

2.6.3 Ordinal regression

2.6.4 Multinomial logistic regression

2.6.5 Dirichlet regression

2.7 Benchmarking of Lasso models dealing with missing values

2.7.1 Knowledge guided multi-level network inference

2.7.2 Two-steps based models

Grouped adaptive Lasso (GALasso)

Stacked adaptive Lasso (SALasso)

2.7.3 Inverse covariance based methods

Convexed conditioned Lasso (CoCoLasso)

Lasso with high missing rate (HMLasso)

2.8 Adult data correlation

2.8.1 Random forest

Imbalanced random forest

Nested cross-validation in random forest

2.8.2 t-test, manova, log-reg

3 Summary of publications

1. Juan Henao, Alida Kindt, Tanja Seegmüller, Kai Foerster, Andreas Flemmer3, Juergen Behr4, Niko-laus Kneidinger, Marion Frankenberger, Fabian Theis, Benjamin Schubert, Markus List, Anne Hilgendorff. Multi-omic signatures relate to the severity of pulmonary outcome in neonates traced into adult disease.

Summary: This project focused on the detection of endotypes behind Bronchopulmonary Dysplasia (BPD) by proteomics, metabolomics, and clinical data integration using a cohort of 55 neonates with and without BPD. The endotypes were detected using Multi-Omics Latent Factor Analysis (MOFA) REF with sensitivity selection. We caught seven latent factors. However, none showed a sign of endotyping discrimination given the combined distribution of latent scores between no BPD and BPD patients in each latent factor. Nevertheless, the biological interpretation of each latent factor allowed us to discover a persistent inflammatory disease component in BPD.

We expanded our analysis by looking for individual molecular features with the potential to be biomarkers of severity using ANOVA with a t-test as a post-hoc method comparing no BPD, mild BPD, and moderate/severe BPD. Acknowledging the clinical heterogeneity signal of BPD cases, we reclassified them into no or moderate/severe BPD using a random forest model trained using oxygen supplementation and mechanical ventilation days (clinical variables used to diagnose BPD). We applied a t-test to identify significant molecular features between no BPD and moderate/severe BPD. We complement our analysis by training different random forest models combining significant molecular features and sets of increasing BPD characterization:

- a) **BPD descriptors:** Oxygen supplementation and mechanical ventilation days.
- b) **Main risk variables:** Gestational age and birth weight.
- c) **Deep clinical phenotyping:** *Main risk variables* and a compendium of clinical measurements encompassing comorbidities, medical interventions, and previously defined MRI-based scores.

The metabolite PC(O-36:5) was detected in both significant analyses and, combined with deep clinical phenotyping, improves the BPD classification along with PC(O-44:5) and gestational age. The protein CCL22 was detected in both significant analyses and improved the BPD classification according to random forest when combined with the main risk variables. Besides, SCGF-alpha, SCGF-beta, and KIR3DL2 were significantly different by ANOVA analysis of no, mild, and moderate/severe BPD comparison.

We traced our significant proteins in an adult chronic lung disease cohort composed of Chronic Obstructive Pulmonary Disease (COPD), Idiopathic Pulmonary Fibrosis (IPF), and healthy donors by ANOVA analysis comparing the three conditions. CCL22 and KIR3DL2 were detected in COPD, while SCGF-beta was significant in COPD and IPF. Those results support the hypothesis regarding the susceptibility of neonates with a BPD diagnosis to develop chronic lung diseases in adulthood.

Contribution: I performed the data pre-processing and all the analyses used in this project. Besides, I created all the data visualization and wrote the first draft of the paper, which was reviewed and edited by Anne Hilgendorff, Markus List, and Tanja Segmuller.

2. Erika Gonzalez Rodriguez1, Juan Henao2, Motaharehsadat Heydarian1, Tina Pritzke1, Alida Kindt3, Anna M. Dmitrieva1, Heiko Adler4, 5, Melanie Markmann6, Valeria Viteri-Alvarez1, Prajakta Oak1,

Markus Koschlig¹, Xin Zhang¹, Kai M. Foerster⁷, Andreas Flemmer⁷, Hamid Hossain^{6,8}, Xavier Pastor², Holger Kirsten⁹, Peter Ahnert⁹, Juergen Behr¹⁰, Tushar J. Desai¹¹, Benjamin Schubert², Anne Hilgendorff^{1,12}. Hyperoxia-induced cell cycle arrest drives long-term impairment of lung development and DNA repair in neonates.

3. Juan David Henao Sanchez^{3,14}, Mustafa Abdo^{1,2}, MD, MSc, Benjamin Schubert^{3,14}, PhD, Markus List⁴, PhD, Henrik Watz^{2,14}, MD, Frauke Pedersen^{1,2,14}, PhD, Alina Bauer^{3,15}, MSc, Dominik Thiele^{5,14}, MSc, Adam M. Chaker^{6,7}, MD, Constanze A. Jakwerth^{7,15}, PhD, Benjamin Waschki^{1,8,14}, MD, Anne Kirsten^{2,14}, MD, Markus Weckmann^{9,14}, PhD, Oliver Fuchs^{9,10,14}, MD, PhD, Gesine Hansen^{11,16}, MD, Matthias V. Kopp^{9,14}, MD, Erika v. Mutius^{12,13,15}, MD, MSc, Inke R. König^{4,14}, PhD, Klaus F. Rabe^{1,14}, MD, PhD, Thomas Bahmer^{1,14}, MD, Carsten B. Schmidt-Weber^{7,15}, PhD, Ulrich M. Zissler^{7,15}, PhD, and the ALLIANCE Study Group*. Cytokines Derived from Nasal Epithelial Lining Fluid in Patients with Asthma.
4. Henao, J. D., Lauber, M., Azevedo, M., Grekova, A., Theis, F., List, M., ... & Schubert, B. (2023). Multi-omics regulatory network inference in the presence of missing data. *Briefings in Bioinformatics*, 24(5), bbad309.

Summary: This project is our contribution to solving one of the most common issues in biological network inferences: the presence of missing data. Here, we extended the previous R package developed by Christoph Ogris, KiMONo (Knowledge guided Multi-Omics Network inference) REF, which uses sparse-group lasso (SGLasso) REF to establish multi-omics edges between molecular features based on a prior network. KiMONo was extended by benchmarking L1-regularized model extensions (sparse models) dealing with missing data and which implementation and code source were performed in R language. We tested SGLasso with imputed data by kNN-imputation (knnS-GLasso), two methods based on inverse covariance matrix (CoCoLasso and HMLasso) REF, and four methods using multi-imputed data (S(A)Lasso and G(A)Lasso) REF. The connection between independent and dependent variables (edge in the network) was established if the beta coefficient was non-zero and the r-squared was larger than 0.1.

We evaluated those sparse models by comparing performance using the same multi-omics information: transcriptomics, CNVs, and methylation. We repeated the comparison using three different datasets extracted from TCGA: Breast invasive carcinoma (TCGA-BRCA), muscle-invasive bladder cancer (TGCA-MIBC), and prostate adenocarcinoma (TGCA-PRAD). We used two prior networks, one extracted from BioGrid for TCGA-BRCA (as the original KiMONo's paper) and one created from FunCoup V5 to evaluate TGCA-MIBC and TGCA-PRAD.

We simulate the most common missing data scenarios, single-omics (transcriptomics), and multi-omics random missingness from 0 to 50%, increasing by 10%. Besides, we added Gaussian noise by adding 0, 0.5, and 1.5 times the standard deviation to the normal distribution. We also evaluated block-missingness cases, i.e., when samples did not match between omics measurements by random experiment removal from 0 to 50%, increasing by 10%. We repeated the experimental setup five times, changing the random seed. We compared the number of nodes, transitivity, the median of R-squared, and the F1 score calculated in two ways, using the whole data (no missingness) as a reference and a network inferred using stability selection fitting 100 random seeds. In addition, we compared the run time per sparse method, and we reached the ability of each method to detect expression quantitative trait methylation (eQTM) regarding a state-of-the-art method, Matrix eQTL.

In general, for small datasets (TCGA-BRCA), the sparse models based on the inverse covariance matrix (CoCoLasso and HMLasso) failed to infer networks even at a minimal missing ratio (10%). However, with larger datasets (TGCA-MIBC and TGCA-PRAD), they could outperform the methods based on multiple imputations (S(A)Lasso and G(A)Lasso) and are computationally more efficient. However, SLasso and knnSGLasso tend to perform well, independent of the size of the dataset. The matrix eQTL and sparse models' performance were quite similar. Nevertheless, the sparse models detected marginally more eQTM-linked genes than Matrix eQTL.

Contribution: I created the benchmark framework to automatize the missing data simulations and run the different methods simultaneously. Manuel Acevedo and Michael Laube implemented the inverse covariance matrix-based methods (CoCoLasso and HMLasso), Anastasiia Grekova implemented the knnSGLasso method, and I implemented the multiple imputation-based methods (S(A)Lasso and G(A)Lasso). I joined all the results and conducted the model performance evaluation, run time comparison, and eQTM analysis. I wrote the manuscript draft along with Benjamin Schubert, Christoph Ogris, and Markus List. I created the result visualization with the help of Christoph Ogris.

5. Caroline Johansson¹, Yvonne Kraus¹, Juan David Henao Sanchez², Kathrin Wolf³, Carola Voss⁴, Tobias Stöger⁴, Friederike Häfner^{1,4}, Sophia Stöcklein⁵, Andreas W. Flemmer⁶, Kai Förster⁶, Anne Hilgendorff, Marie Standl. Impact of maternal exposure to airborne pollutants during pregnancy on pulmonary morbidity in preterm infants

4 Discussion

This is the discussion of the thesis.

A Appendix

Multi-omics regulatory network inference in the presence of missing data

Juan D. Henao , Michael Lauber, Manuel Azevedo, Anastasiia Grekova, Fabian Theis, Markus List, Christoph Ogris[†] and

Benjamin Schubert[†]

Corresponding author: Benjamin Schubert, Member of the German Center for Lung Research (DZL), Helmholtz Zentrum München, Computational Health Department, Ingolstädter Landstraße 1, 85764 Munich, Germany. Telephone: +49 89 3187 43046. E-mail: benjamin.schubert@helmholtz-munich.de

[†]Christoph Ogris and Benjamin Schubert are joint last authors

Abstract

A key problem in systems biology is the discovery of regulatory mechanisms that drive phenotypic behaviour of complex biological systems in the form of multi-level networks. Modern multi-omics profiling techniques probe these fundamental regulatory networks but are often hampered by experimental restrictions leading to missing data or partially measured omics types for subsets of individuals due to cost restrictions. In such scenarios, in which missing data is present, classical computational approaches to infer regulatory networks are limited. In recent years, approaches have been proposed to infer sparse regression models in the presence of missing information. Nevertheless, these methods have not been adopted for regulatory network inference yet. In this study, we integrated regression-based methods that can handle missingness into KiMONo, a Knowledge guided Multi-Omics Network inference approach, and benchmarked their performance on commonly encountered missing data scenarios in single- and multi-omics studies. Overall, two-step approaches that explicitly handle missingness performed best for a wide range of random- and block-missingness scenarios on imbalanced omics-layers dimensions, while methods implicitly handling missingness performed best on balanced omics-layers dimensions. Our results show that robust multi-omics network inference in the presence of missing data with KiMONo is feasible and thus allows users to leverage available multi-omics data to its full extent.

Keywords: multi-omics integration, network inference, data missingness, Lasso model, data imputation

INTRODUCTION

Complex biological systems are organised in multi-level, dynamically controlled networks that regulate and maintain the phenotypic behaviour of individual cells and their response to environmental changes [1]. Uncovering these multi-level networks and systemically understanding the interplay of their elements is a key problem in computational biology. Modern high-throughput multi-omics techniques now enable access to each regulatory network level, even at single-cell resolution [2, 3].

However, combining multi-omics measurements and reconstructing the underlying regulatory network remains challenging [4]. Generally, sparse interaction networks in the form of directed or undirected graphs are constructed from dynamic

interventional omics or large observational data using different classes of statistical methods [4]. Common approaches are either correlation-based [5], use techniques from information theory [6–8], or use (regularised) regression and variable selection frameworks to infer graphical models [9–11]. Most recent methods also integrate prior knowledge [12, 13], such as experimentally determined protein–protein interaction networks, known metabolic pathways, or even predicted miRNA–mRNA interactions [14].

One such recent approach is KiMONo, Knowledge guided Multi-Omics Network inference [15], a two-step prior knowledge-based approach for multi-omics regulatory network inference. In the first step, the framework uses the whole dataset to model each omics element individually, detecting statistical effects between

Juan Henao is a 3rd year PhD candidate at Computational Health Center at Helmholtz Center Munich working on multi-omics and clinical data integration using both, bulk and single-cell data.

Michael Lauber is a PhD Candidate at the Chair of Experimental Bioinformatics at the Technical University Munich. Currently, he is working on an approach for inference of reprogramming transcription factors for trans-differentiation.

Manuel Azevedo is a Master's student at the Technical University of Munich in Mathematics with a focus on Biomathematics and Biostatistics. Currently, he is working as a Student Assistant at Helmholtz Munich, where he is also doing his master's thesis.

Anastasiia Grekova is a Master's student of bioinformatics at the Technical University of Munich and the Ludwig-Maximilians-University Munich, working on multi-omics data integration in Marsico Lab at HMGU.

Fabian Theis is the Head of the Institute of Computational Biology and leading the group for Machine Learning at Helmholtz Center Munich. He also holds the chair of 'Mathematical modelling of biological systems', Department of Mathematics, Technical University of Munich as an Associate Professor.

Markus List obtained his PhD at the University of Southern Denmark and worked as a postdoctoral fellow at the Max Planck Institute for Informatics before starting his group Big Data in BioMedicine at the Technical University of Munich.

Christoph Ogris holds a PostDoc position in the Marsico Lab at Helmholtz-Center Munich. His research focuses on predicting and exploiting multi-modal biological networks to identify disease-specific cross-omic interactions.

Benjamin Schubert obtained his PhD at the University of Tübingen and worked as a postdoctoral fellow at Harvard Medical School and Dana-Farber Cancer Institute USA before starting his group for Translational Immunoinformatics at the Helmholtz Center Munich.

Received: December 8, 2022. Revised: May 6, 2023. Accepted: May 29, 2023

© The Author(s) 2023. Published by Oxford University Press.

This is an Open Access article distributed under the terms of the Creative Commons Attribution License (<https://creativecommons.org/licenses/by/4.0/>), which permits unrestricted reuse, distribution, and reproduction in any medium, provided the original work is properly cited.

them. The framework combines all models in a second step, assembling a multi-omics graph with the input features as nodes linked via edges representing the detected effects.

However, a major drawback of most network inference methods is their inability to handle missing data. Many omics technologies such as mass-spectrometry-based proteomics and metabolomics or single-cell transcriptomics suffer from inherent missingness due to the stochasticity of the biological processes and technical limitations. Additionally, it is often necessary to combine multiple studies that only partially measure the same omics levels to reach sample sizes adequate for network inference, creating patterns of block-wise missingness. Many classical regulatory network inference methods ignore missing data and focus only on analysing complete cases, thus underutilising the collected data and severely limiting the amount of information used. Removing samples with missing features can also lead to biased estimates if the missingness is not completely random [16], potentially affecting the extracted regulatory network. Multiple imputation [17] is another popular approach to deal with missingness, followed by applying any classical network inference method using *ad hoc* rules to harmonise variable selection across multiply-imputed datasets [18]. However, Ganti and Willet [19] demonstrated that such two-step approaches can be sub-optimal and instead require integrated or more general frameworks to handle missing data and variable selection jointly.

In recent years, advances have been made in using sparse graphical models for data with missing information. These approaches can be roughly categorised in Bayesian methods using data augmentation strategies [20], methods using pooled posterior [21], bootstrapped inclusion probabilities [22, 23], methods performing variable selection through stacked [18, 24] or group Lasso integrated multiple imputation methods [25–28], low-rank matrix completion [19, 29], inverse probability weighting [30], Lasso regularised inverse covariance estimation [31–34] and Expectation-Maximisation-based approaches [35, 36]. While most methods address the missingness of individual features, some methods exist that also explicitly model block-missingness [28, 37–39].

Incorporating such approaches in multi-omics network inference is attractive. However, a comprehensive benchmark of existing methods that can handle missing data is lacking. We, therefore, extended KiMONo with various regression-based approaches that integrate and combine prior imputed data [28] and Lasso-regularised inverse covariance estimation methods [33, 34]. We systematically evaluated how these methods handle gradually increasing levels of artificial noise and missingness for regulatory network inference on single- and multi-omics data.

We observed that approaches explicitly handling missingness in a two-step manner performed best over a wide range of random, block-missingness and noise levels in an imbalanced multi-omics dimensional dataset (TCGA-BRCA), while implicit covariance-estimation-based methods performed best in multi-omics with balanced dimensional dataset (TCGA-MIBC and TCGA-PRAD).

METHODS

Knowledge guided Multi-Omics Network (KiMONo) inference

KiMONo, Knowledge guided Multi-Omics Network inference [15], is a two-step inference procedure to detect statistical dependencies between omics features and construct a multi-omics network. Here, the inference complexity is decreased by pre-selecting

Table 1. Inference models included in this benchmark and capable of dealing with missing data

Method	Category	Citation
knnSGLasso	Single imputation + KiMONo	[15, 42]
S(A)Lasso	Stacked multiple imputation	[28]
G(A)Lasso	Grouped multiple imputation	[28]
HMLasso	Inverse covariance estimation	[33]
CoCoLasso	Inverse covariance estimation	[33, 34]
BDCoCoLasso	Inverse covariance estimation	[43]

feature dependencies based on existing prior knowledge of biological mechanisms such as known protein–protein interaction. Based on such prior knowledge, a system of linear multivariate regression models with sparse-group Lasso penalty is constructed [40]:

$$\min_{\beta} \frac{1}{2} \left\| y - \sum_{l=1}^m X^{(l)} \beta^{(l)} \right\|_2^2 + (1 - \alpha) \lambda \sum_{l=1}^m \sqrt{p^{(l)}} \left\| \beta^{(l)} \right\|_2 + \alpha \lambda \|\beta\|_1 \quad (1)$$

where $\alpha = [0, 1]$, and l the group index of m groups (omics layer), being $X^{(l)}$, $\beta^{(l)}$, $p^{(l)}$ the submatrix, coefficients, and length of coefficients for group l , respectively. This allows two types of sparsities: (i) ‘groupwise sparsity’, referring to all groups with at least one non-zero coefficient and (ii) ‘within group sparsity’, referring to all non-zero coefficients within every specific group (omics layer). The framework combines all models in a second step, assembling a multi-omics graph connecting regressed omics elements with its non-zero regressor elements if the regression performance (coefficient of determination) exceeds a predefined threshold.

Regression-based methods for network inference and imputation

We focused on methods with a working R implementation and consistent documentation. These requirements left us with five advanced statistical approaches of three categories (i) stacked and (ii) grouped multiple-imputation, as well as (iii) Lasso-based inverse covariance estimation approaches (Table 1). All mentioned methods have been integrated into the KiMONo framework (<https://github.com/cellmapslab/kimono>). The data underlying this article, including detailed benchmarking results and code are available at Zenodo [41] (<https://doi.org/10.5281/zenodo.7900595>).

kNN-imputation & sparse group Lasso (knnSGLasso)

We implemented a two-step approach that first imputes missing information using nearest neighbour averaging followed by applying the Sparse Group Lasso approach of the classical KiMONo. The kNN-based imputation method [42] implemented in the R package impute v1.46.0 was applied separately to individual omics layers and other covariates. Originally designed for the imputation of gene expression data, the method replaces missing values by averaging non-missing values of its nearest neighbours. If the percentage of missing data allowed for every variable (e.g. 50% per gene (default)) is exceeded, the missing values are imputed using the overall mean per sample. Only samples with missingness <80% (default) were considered for the imputation. The algorithm’s parameters were set to default values: the number of neighbours used in the imputation was set to $k = 10$, and the largest block of variables imputed using the kNN algorithm before

recursively dividing the feature into smaller chunks was set to max = 1500.

Stacked adaptive Lasso (SALasso)

Stacked approaches combine prior D -times multiply imputed datasets by averaging over them, making such approaches applicable to existing sparse regression frameworks:

$$\hat{\theta} = \operatorname{argmin}_{\theta} \left\{ -\frac{1}{n} \sum_{d=1}^D \sum_{i=1}^n o_i \log L(\theta | Y_{d,i}, X_{d,i}) + \lambda \sum_{j=1}^p |\hat{\alpha}_j| \right\} \quad (2)$$

This results in a pooled estimate $\hat{\theta}$ across all imputed datasets, thereby enforcing uniform variable selection across p covariates. Here α_j are so-called adaptive weights used to address estimation and selection inconsistencies that can occur and are defined as $\hat{\alpha}_j = (\beta_j^0 + 1/n)^{-\gamma}$ with some $\gamma > 0$ and β_j^0 initial estimates obtained with (stacked)-Lasso or (stacked)-elastic-net. Following Du et al. [28] γ is set to $\lceil \frac{2v}{1-v} \rceil + 1$ with $v = \log(p)/\log(n)$. As stacking multiple imputed datasets can be seen as artificially increasing the sample size thereby potentially creating oversampling biases, observation weights o_i for each subject can be included to account for the artificial inflation of the dataset with either uniform weights or with weights accounting for varying degrees of missingness per subject. We used the SALasso implementation released in the R package `miselect 0.9.0` [28]. We tested 50 λ values with a `lambda.min.ratio` of 1e-4 in a 5-fold cross-validation with and without adaptive weights (SALasso), while the sample weights o_i were set to be uniform.

Grouped adaptive Lasso (GALasso)

Similar to SALasso, GALasso pools across prior imputed datasets by adding a group Lasso penalty term enforcing consistent variable selection across multiply imputed datasets yielding:

$$\begin{aligned} (\hat{\theta}_1, \dots, \hat{\theta}_D) = & \operatorname{argmin}_{\theta} \left\{ -\frac{1}{n} \sum_{d=1}^D \sum_{i=1}^n \log L(\theta | Y_{d,i}, X_{d,i}) \right. \\ & \left. + \lambda \sum_{j=1}^p \hat{\alpha}_j \sqrt{\sum_{d=1}^D \beta_{d,j}^2} \right\} \end{aligned} \quad (3)$$

with $\hat{\alpha}_j = \left(\sqrt{\sum_{d=1}^D \hat{\beta}_{d,j}^2} + 1/(nD) \right)^{-\gamma}$, where $\hat{\beta}_{d,j}$ is the estimate obtained with $\hat{\alpha}_j = 1$, $\gamma = \lceil 2v/1-v \rceil + 1$, and $v = \log(pD)/\log(nD)$ [28]. Contrary to SALasso, individual sets of parameters θ_d are obtained for each of the D imputed datasets. However, the group-Lasso regularisation term jointly shrinks all $\beta_{d,j}$ to zero, despite their potentially different estimates across the D -imputed datasets enabling a uniform variable selection. We used the GALasso implementation released in the R package `miselect 0.9.0` [28]. We tested 50 λ values with a `lambda.min.ratio` of 1e-4 in a 5-fold cross-validation with and without adaptive weights (GLasso).

Convex conditioned Lasso (CoCoLasso)

CoCoLasso is an inverse covariance estimation method for high-dimensional data with missing values. The main idea is to reformulate the Lasso regression by working with the sample covariance matrix of X , $S = \frac{1}{n}X'X$, and the sample covariance vector of X and y , $\rho = \frac{1}{n}X'y$. With this reformulation β is estimated via S and ρ instead of X and y . Since we are dealing with missing data, CoCoLasso works with the pairwise covariance $S^{pair} := (S_{jk}^{pair})$ and

$\rho^{pair} := (\rho_{jk}^{pair})$ instead of S and ρ , where

$$S_{jk}^{pair} = \frac{1}{n_{jk}} \sum_{i \in I_{jk}} X_{ij} X_{ik} \text{ and } \rho_{jk}^{pair} := \frac{1}{n_{jj}} \sum_{i \in I_{jj}} X_{ij} y_i, \quad (4)$$

where $I_{jk} := \{i : X_{ij} \text{ and } X_{ik} \text{ are observed}\}$, and n_{jk} is the number of elements of I_{jk} . Then the inference problem can be defined as:

$$\beta = \operatorname{argmin}_{\beta} \beta' \Sigma \beta - \rho^{pair} \beta + \lambda \|\beta\|_1, \quad (5)$$

$$\Sigma = \operatorname{argmin}_{\Sigma \geq 0} \|\Sigma - S^{pair}\|_{max}, \quad (6)$$

Since the matrix S^{pair} needs to be positive semidefinite, CoCoLasso applies in (5) the alternating direction method of multipliers algorithm to obtain a positive semidefinite covariance matrix, Σ . Afterwards, CoCoLasso optimises the Lasso objective function (4). We used the CoCoLasso implementation released in the R package `HMLasso 0.0.1` [34] with the following selection of hyperparameters: For λ , we tested 50 values with a `lambda.min.ratio` of 1e-1 in a 5-fold cross-validation.

To improve the computational efficiency of CoCoLasso, BDCoCoLasso was developed by implementing a block coordinate descent strategy [43] over corrupted (with missing data) and uncorrupted (full data) covariates. However, our preliminary results showed the random missingness of 5% and more encompassed already all samples and therefore BDCoCoLasso could not be applied. We, thus, refrain from discussing the performance of BDCoCoLasso.

Lasso with high missing rate (HMLasso)

HMLasso can be seen as an optimally weighted modification of CoCoLasso according to the missing ratio. HMLasso uses the mean imputation method. Instead of X , the mean imputed data variable, Z is used, where $Z_{jk} = X_{jk}$ for an observed element and $Z_{jk} = 0$ otherwise. Followed by that the S^{pair} covariance matrix from CoCoLasso is modified by the covariance matrix of the mean imputed data, defined as $S^{imp} = R S^{pair}$, where $R_{jk} = n_{jk}/n$.

The HMLasso is then formulated as:

$$\beta = \operatorname{argmin}_{\beta} \beta' \Sigma \beta - \rho^{pair} \beta + \lambda \|\beta\|_1, \quad (7)$$

$$\Sigma = \operatorname{argmin}_{\Sigma \geq 0} \|W \odot (\Sigma - S^{pair})\|_F^2, \quad (8)$$

where $W = R_{jk}^\alpha$ is defined as the weight matrix and α is the weight power parameter. We obtain the positive semi-definite matrix by minimising the weighted Frobenius norm in (8). Afterwards, HMLasso optimises the Lasso objective (7). Note that if we set $\alpha = 0$ (non-weighted case), it is equal to CoCoLasso, replacing the Frobenius norm with the max norm. We used the HMLasso implementation released in the R package `HMLasso 0.0.1` [34] with the following selection of hyperparameters: For α , we tested values between 0.5 and 2 with an interval of 0.5. For λ , we tested 50 values with a `lambda.min.ratio` of 1e-1 in a 5-fold cross-validation.

Datasets

For the benchmarking performance evaluation, we collected two datasets consisting of triple-omics data (transcriptomics, copy number variation (CNV), and methylation) from the PanCancer Projects [44] using The Cancer Genome Atlas (TCGA) data portal and the cBioPortal [45]. Both datasets were already pre-processed including normalised and log-transformed gene expression, beta-values for methylation, and linear CNVs.

The first dataset encompassed 871 patients with breast invasive carcinoma (retrieved on 3 July 2022), while the second dataset comprised 414 patients with muscle-invasive bladder cancer (retrieved on 23 August 2022). All samples containing missing information were removed to construct a complete data set as the baseline, thus restricting the data sets to 604 and 404 samples, respectively. Similarly, features with low variance were removed, resulting in 11,530 transcriptomics, 1,366 methylation, and 84 CNV features for the first dataset (TCGA-BRCA from now on), and 20,085 transcriptomics, 15,460 methylation and 24,765 CNV features for the second dataset (TCGA-MIBC from now on), respectively. These two datasets allowed us to study the impact of data dimensionality on performance.

A third dataset was collected from the PanCancer Projects [44] using The Cancer Genome Atlas (TCGA) data portal and the cBioPortal [45] for prostate adenocarcinoma (TCGA-PRAD from now on) (retrieved on 2 March 2023) comparable to TCGA-MIBC to evaluate performance consistency across similar datasets. The third dataset encompassed 491 samples with 20,123 transcriptomics, 15,576 methylation and 24,765 CNV features after low variance removal.

Prior network generation

For the first dataset (TCGA-BRCA), we used the prior network of the original KiMONo publication [15]. Briefly, protein–protein interactions were extracted from the BioGrid interactome (Release 3.5.188) [46], associated with the extracted gene expression information, and linked each CNV and methylation site to its associated gene since both omics layers were already annotated to gene identifiers. The final prior network contained 11,645 nodes (10,848 genes, 84 CNVs and 713 methylation sites).

The construction of the prior network for the second and third datasets (TCGA-MIBC and TCGA-PRAD) followed the same steps as the first prior network construction using FunCoup v5 [47] (built-in September 2020) as the basis. We used interactions with maximum Final Bayesian Score supporting protein–protein interactions, and log-likelihood ratio for physical protein–protein interaction evidence ≥ 99 percentile to extract reliable protein–protein interactions. The final prior network contained 51,318 nodes (18,990 genes, 18,460 CNVs and 13,868 methylation sites).

Network-based multiple imputation

Stacked and grouped adaptive Lasso approaches require multiple imputed data as input. However, multiple imputation methods do not scale well to high-dimensional data with high missingness, and standard implementations such as those offered in the R package MICE, take multiple hours to days to finish. Thus, we developed a novel network-guided multiple imputation by chained equation approach (ngMICE) by utilising KiMONo's prior network. Instead of considering all covariates for imputation, we restrict each imputation attempt to the covariates that are directly linked to the missing covariate in a prior network as other covariates will be removed during network inference by KiMONo and therefore can be neglected. The number of covariates can be further reduced by correlation-based filtering. For missing elements retaining $< k$ covariates for imputation, the top k correlated covariates are used. Once the covariate matrix has been constructed as described, the standard MICE procedure is run. We used the R package MICE 3.14.0 [48], with predictive mean matching (default) as a multiple imputation approach, an absolute Pearson correlation coefficient of 0.1 as the threshold, and $k = 5$ most correlated features in the final regression formula. ngMICE performed similarly to kNN-based imputation in terms

of RMSE across omics types and missingness with slightly worse average performance (Supplementary Figure S1).

Benchmark

We assessed the performance of the selected inference methods in the presence of missing data by simulating three typical scenarios—(i) random missing information in a single omics level and across (ii) multiple levels, as well as (iii) block-missingness structures (Figure 1A). To test the methods' capabilities even further, we decreased the signal-to-noise ratio by systematically adding covariate-specific white noise to the input data. Each experiment was repeated five times for robust performance estimation and corrected for confounding age and sex effects.

Single-omics missing

We selected the transcriptomics level as a single-omics layer to test the different models' capabilities to infer gene regulatory networks with less directly informative co-correlation structures that could be used to impute the missing gene expression information. We removed $m \in \{0, 10, 20, 30, 40\}$ randomly selected entries from the input data. Additionally, we added noise with increasing intensity to the data by drawing from a normal distribution $\epsilon \sim N(0, a\sigma_g^2)$ per gene with a gene-specific variance term estimated from the real data and $a \in \{0, 0.5, 1.5\}$.

Multi-omics missing

To test the models' capabilities to handle more complex co-correlation structures that could potentially be exploited for better imputation, we expanded the single-omics experiment to jointly consider the three available omics types. As before, we randomly removed $m \in \{0, 10, 20, 30, 40\}$ of entries independently per omics layer and added feature-dependent noise to the data as described before while ensuring to bound the beta values of the methylation data to the range between 0 and 1.

Multi-omics block-missing

To test the capabilities of the method to handle block-wise missing information, we removed $m \in \{0, 10, 20, 30, 40\}$ samples per omics layer such that at least two omics-layers still remained per sample. Additionally, we added noise to the remaining samples as described before.

Downsampling

Similarly, to identify the minimal number of samples required to infer reliably the regulatory network, we downsampled the dataset to $k \in \{90, 80, 70, 60\}$ of samples.

Runtime

Using TCGA-BRCA dataset, we tested the runtime of each method on a dedicated machine with an AMD EPYC 7502P 32-Core Processor with 2.5GHz base clock speed and 860 GB RAM using the multi-omics missingness experiment with the same configurations as before. We ran all experiments with ncores=60 (with hyperthreading enabled).

Evaluation metrics

To construct the final network from the individual regressions, we applied a strict filter connecting independent to dependent variables if their beta coefficient was non-zero and their $R^2 > 0.1$.

Prediction metrics

To measure the prediction qualities of each regression model, we recorded the root means squared error (RMSE) and R^2 respectively,

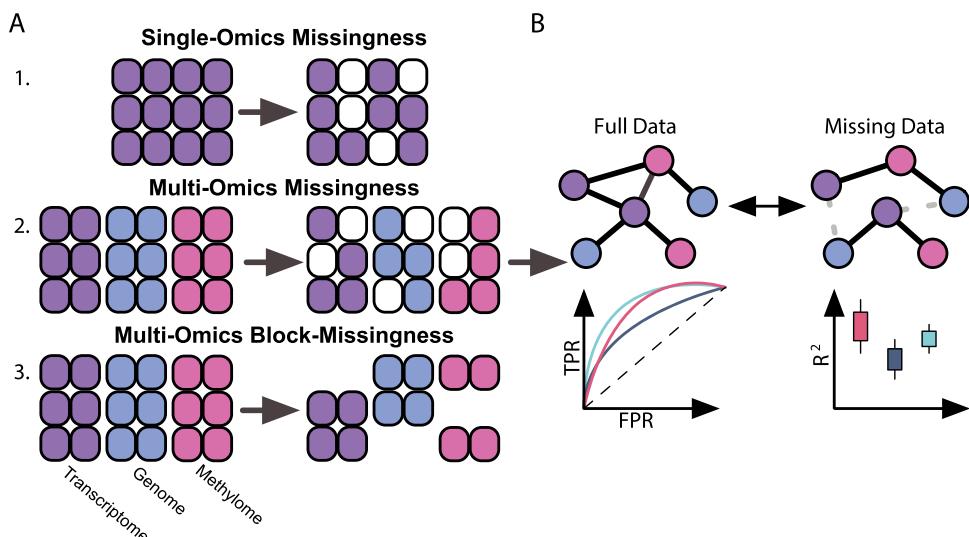


Figure 1. Benchmark schematic. **(A)** Three common missingness scenarios in regulatory network inferences are tested: (1) single-omics, (2) multi-omics random missingness of individual elements, and (3) block-wise missingness in which entire omics layers are missing for an individual. **(B)** We tested five approaches of two categories: (1) Two-step approaches that first impute and then aggregate imputation through and (2) Inverse covariance estimation approaches that implicitly handle missingness during inference. We inferred regulatory networks from full data and data missing for each method and compared the resulting networks with multiple performance metrics.

and compared their distributions based on a Wilcoxon rank-sum test (Figure 1B).

Network reconstruction metrics

To measure the methods' abilities to handle missing data well, we compared the inferred regulatory networks from missing data to their counterpart inferred on full data, and calculated precision, recall, and F1-scores of the recovered network edges (Figure 1B).

Similarly, we inferred a ground truth network with KiMONO using stability selection by repeating the network inference 100 times fitting different random seeds and averaging over the resulting coefficients and R² values before constructing the network, improving the robustness of the graph resulting in a final network consisting of 2,458 nodes (2,182 genes, 63 CNVs and 211 methylations) and 6,554 edges for TCGA-BRCA dataset, 33,652 nodes (12,738 genes, 11,052 CNVs and 9,862 methylations) and 87,229 edges for TCGA-MIBC dataset, and 29,502 nodes (11,251 genes, 9,862 CNVs, and 8,389 methylations), and 75,943 edges for TCGA-PRAD dataset. We compared each stability-selected reference network to all inferred networks on missing data to determine performance differences across the individual methods.

Topological metrics

The interpretation of complex heterogeneous networks and identification of key modules and important network nodes relies on the topological network features. Hence, it is also vital to evaluate if the methods can robustly infer topological structures. Therefore, we use multiple network-based metrics such as node-degree distribution, betweenness centrality, and clustering coefficient to quantify and compare the topological changes of the networks inferred from missing data. Node degree indicates the sparseness of the network, while betweenness centrality indicates how interconnected the network is, and the global clustering coefficient (transitivity) indicates how densely connected neighbouring nodes are.

Quantitative trait analysis

The sparse methods implicitly computed multivariate expression quantitative trait methylation (eQTM) as we mostly restricted

the calculation between gene expression and their corresponding methylation site measurements. Hence, we compared our findings with a state-of-the-art method based on a linear regression approach to eQTM calculation, Matrix eQTL [49]. Our results enabled us to retrieved the genes connected to methylation sites based on topologies constructed using previously chosen thresholding: (i) R² per model >0.1 and (ii) regression coefficient different to zero. For Matrix eQTL, we retrieved the significantly linked genes to eQTMs detected running both, imputed gene expression and imputed methylation sites measurements as inputs for the set of features in the different inferred networks, with a significance threshold per model of 1e-5, and FDR threshold of 0.01.

RESULTS

Most topological features can be conserved in data with missingness

To investigate to what extent topologies of inferred networks are affected by noise and missingness, we computed multiple network properties across all benchmarking scenarios for all methods.

In single-omics scenarios, knnSGLasso demonstrated the most stable topology with increasing missingness rate on TCGA-BRCA data (Figure 2A and B, Supplementary Figures S2 and S3), decomposing into small subnetwork at high missingness rates according to the increase in transitivity (0.054 ± 0.007 , missingness = 0.1; 0.067 ± 0.003 , missingness = 0.5) and reduction in betweenness centrality (0.005 ± 0.007 , missingness = 0.1; 0.002 ± 0.004 , missingness = 0.5; Figure 2A and B). All methods were topologically more stable on TCGA-MIBC data except for CoCoLasso and HMLasso whose transitivity and betweenness centrality dropped at high missingness (transitivity: 0.077 ± 0.002 , betweenness: 0.003 ± 0.002 , missingness = 0.4; transitivity: 0.059 ± 0.004 , betweenness: 0.002 ± 0.002 , missingness = 0.5 for CoCoLasso; Figure 2A and C, Supplementary Figure S4). TCGA-PRAD showed similar topological results to TCGA-MIBC including the dropping of CoCoLasso and HMLasso at a high missingness rate

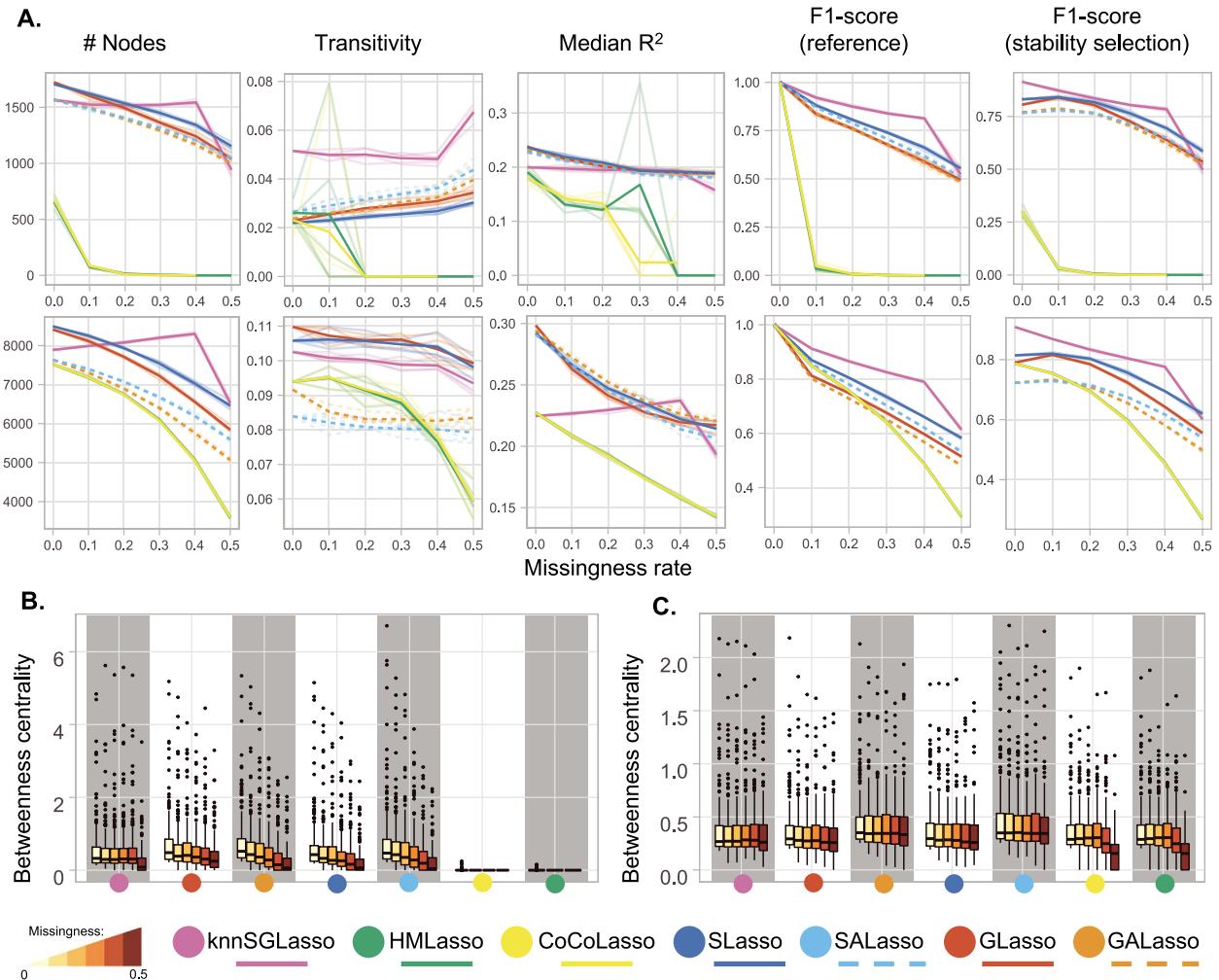


Figure 2. Benchmark results across single-omics experimental setups for both datasets. Transparent lines denote individual runs, while bold lines refer to the average performance. **(A)** TCGA-BRCA dataset (top), TCGA-MIBC dataset (bottom). Performance was evaluated using network size (number of nodes), transitivity (global clustering coefficient), median R^2 , and F1-scores compared to a reference (i.e. the same method applied on the full data and stability selected networks generated with KiMONO). **(B and C)** Illustrating the topological change in TCGA-BRCA **(B)** and TCGA-MIBC **(C)** datasets scaling values for the top 200 nodes with the highest betweenness centrality through increased missingness.

(transitivity: 0.066 ± 0.005 , betweenness: 0.006 ± 0.003 , missingness = 0.5 for CoCoLasso, *Supplementary Figure S8*).

In multi-omics scenarios, knnSGLasso had the most stable topology across increasing missingness rates on TCGA-BRCA. However, in presence of the highest missingness rate, knnSGLasso suffered an abrupt decay in betweenness centrality (*Figure 3A* and *B*). In contrast, knnSGLasso networks became sparser as indicated by the increase in the betweenness centrality and the slight reduction of transitivity with increasing missingness on TCGA-MIBC and TCGA-PRAD data (*Figure 3A* and *C*; *Supplementary Figures S3, S8* and *S9*). S(A)Lasso and G(A)Lasso inferred networks on TCGA-MIBC and TCGA-PRAD data decreased their betweenness and transitivity abruptly compared to networks inferred on TCGA-BRCA data (*Figure 3B* and *C*; *Supplementary Figures S3, S8* and *S9*). CoCoLasso and HMLasso slightly decreased their betweenness centrality for no to medium missingness rates (CoCoLasso and HMLasso: 0.002 ± 0.001 , for missingness = 0/0.3), and dropped in the presence of high missingness (CoCoLasso and HMLasso: 0.001 ± 0.001 , missingness = 0.5; *Figure 3C*; *Supplementary Figure S9*).

knnSGLasso and SALasso presented the most stable topologies in presence of block-missingness on TCGA-BRCA dataset

(*Figure 4B*), although knnSGLasso's full-data topology varied considerably compared with the counterpart networks inferred in presence of missing data (betweenness: 0.006 ± 0.007 and 0.005 ± 0.007 , transitivity: 0.022 ± 0.0 and 0.050 ± 0.006 , for missingness = 0 and missingness = 0.5, respectively; *Figure 4A* and *B*; *Supplementary Figure S3*). S(A)Lasso and G(A)Lasso were less stable, declining slightly in betweenness centrality with increasing missingness rate (*Figure 4B*). knnSGLasso behaved similarly on TCGA-MIBC and TCGA-PRAD data, the betweenness centrality (0.001 ± 0.001 , in both) was drastically different than the rest of its own inferred topologies in presence of missing data (0.004 ± 0.003 and 0.004 ± 0.004 for missingness = 0.1, respectively). The betweenness centrality for the approaches dealing with multiple imputed data dropped abruptly near to zero even at 10% of missingness (*Figure 4C*; *Supplementary Figure S9*). CoCoLasso and HMLasso had the most stable topologies including the full-data inferred network. However, their betweenness and transitivity were lower than knnSGLasso (CoCoLasso and HMLasso: 0.002 ± 0.001 and $0.096 \pm 5e-4$, knnSGLasso: 0.003 ± 0.003 and 0.099 ± 0.001 , respectively for missingness = 0.5; *Figure 4A* and *C*; *Supplementary Figures S3* and *S9*).

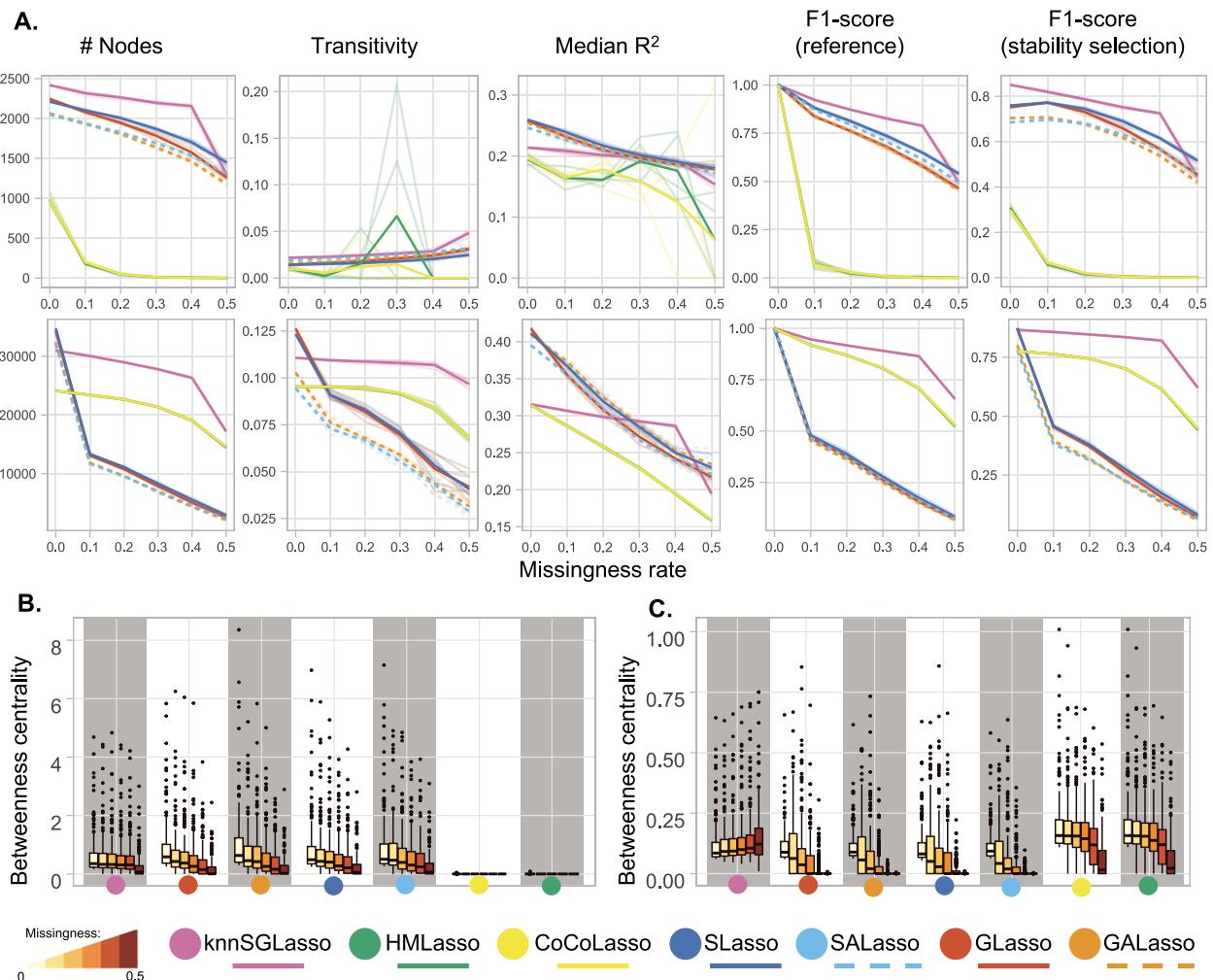


Figure 3. Benchmark results across multi-omics experimental setups for both datasets. Transparent lines denote individual runs, while bold lines refer to the average performance. **(A)** TCGA-BRCA dataset (top), TCGA-MIBC dataset (bottom). Performance was evaluated using network size (number of nodes), transitivity (global clustering coefficient), median R^2 , and F1-scores compared to a reference (i.e. the same method applied on the full data and stability selected networks generated with KiMONo). **(B and C)** Illustrating the topological change in TCGA-BRCA **(B)** and TCGA-MIBC **(C)** datasets scaling values for the top 200 nodes with the highest betweenness centrality through increased missingness.

Across all benchmark settings, the methods showed a decreasing network size with increasing missingness. The constant loss of edges showed a stable ratio indicating no bias towards a specific omics layer (Supplementary Figures S4 and S8). GLasso was less affected in terms of noise, while knnSGLasso produced the most robust results in terms of missingness (Supplementary Figure S2).

kNN imputation-based model performs best for single- and multi-omics data with random missingness

In both the single- and multi-omics setting, knnSGLasso was the best performing method independently of the dataset, reaching F1-scores of 0.921 ± 0.005 and 0.912 ± 0.001 on the single-omics and 0.923 ± 0.002 and 0.947 ± 0.001 on the multi-omics data at 10% missingness (TCGA-BRCA and TCGA-MIBC data, respectively). Followed by SLasso (single-omics: 0.880 ± 0.005 and 0.870 ± 0.002) for single-omics, and CoCoLasso (0.919 ± 0.001) for multi-omics when compared to the reference networks inferred on full data (Figures 2A and 3A; Supplementary Figures S5 and S6).

The performance gradually decreased with increasing missingness and noise levels to a similar degree for all the

methods except for the specific case of multiple imputation-based methods in the multi-omics experiment on TCGA-MIBC data. In this particular case, the F1-score increased to 0.485 ± 0.003 for SLasso, 0.442 ± 0.004 for SALasso, 0.405 ± 0.005 for GLasso and 0.398 ± 0.004 for GALasso, at 50% missingness and a medium noise ($\alpha = 0.5$; Supplementary Figure S6). At higher noise levels, the performance of all methods declined dramatically reaching only F1-scores below 0.15.

For the TCGA-BRCA dataset, HMLasso and CoCoLasso came in last, reaching F1-scores below 0.1 both on single- and multi-omics data even for samples with only 10% missingness and no noise (Figure 2A; Supplementary Figure S6). Both methods even failed to infer networks in 11/15 and 11/15 single-omics experiments at 30 and 40% missingness, as well as in 9/11 and 8/11 multi-omics experiments, at 40 and 50% missingness, respectively.

Similar behaviour could be observed when comparing the inferred networks to the stability selection-based reference network: knnSGLasso performed best with an F1-score of 0.873 ± 0.005 and 0.869 ± 0.008 on single-omics, and 0.820 ± 0.004 and 0.858 ± 0.007 on multi-omics at 10% missingness and no noise in the TCGA-BRCA and TCGA-MIBC dataset, respectively, followed by SLasso (TCGA-BRCA: 0.841 ± 0.007 , TCGA-MIBC: 0.820 ± 0.003),

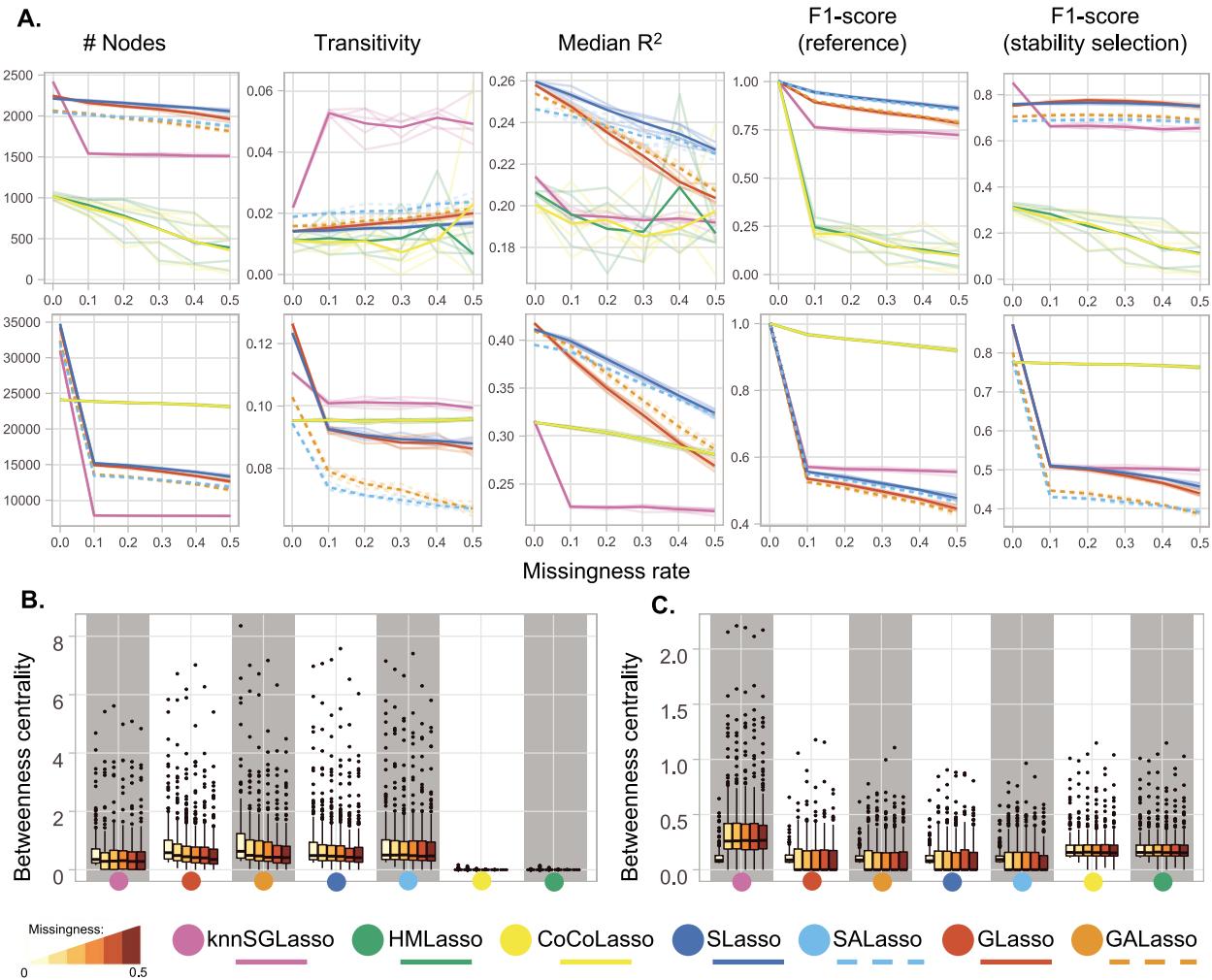


Figure 4. Benchmark results across block-missingness experimental setups for both datasets. Transparent lines denote individual runs, while bold lines refer to the average performance. **(A)** TCGA-BRCA dataset (top), TCGA-MIBC dataset (bottom). Performance was evaluated using network size (number of nodes), transitivity (global clustering coefficient), median R^2 , and F1-scores compared to a reference (i.e. the same method applied on the full data and stability selected networks generated with KIMONO). **(B and C)** Illustrating the topological change in TCGA-BRCA (**B**) and TCGA-MIBC (**C**) datasets scaling values for the top 200 nodes with the highest betweenness centrality through increased missingness.

for single-omics, and SLasso (TCGA-BRCA: 0.772 ± 0.003) and CoCoLasso (TCGA-MIBC: 0.764 ± 0.001) for multi-omics random missingness experiment (Figures 2A and 3A; Supplementary Figure S7). SLasso did marginally outperform knnSGLasso in the TCGA-BRCA dataset in scenarios with high noise and/or high missingness and was almost equivalent at high noise and high missingness (50% missingness, $\alpha = 0.5$) in the multi-omics scenario on the TCGA-MIBC dataset. At the highest noise levels ($\alpha = 1.5$), even in the absence of missingness, none of the methods exceeded F1-scores of 0.18 (Supplementary Figure S7).

TCGA-PRAD showed similar results to TCGA-MIBC in both single- and multi-omics (Supplementary Figure S10). At the single-omics level, knnSGLasso outperformed the rest of the methods in both low- and high-missingness based on F1-scores calculated using full data as reference (0.937 ± 0.001 and 0.912 ± 0.001 at 10% missingness; 0.71 ± 0.017 and 0.613 ± 0.004 at 50% missingness, respectively). Whilst CoCoLasso and HMLasso persisted as worst performing methods in presence of high missing ratio (0.488 ± 0.004 and 0.492 ± 0.004 for TCGA-PRAD, respectively; 0.293 ± 0.004 and 0.295 ± 0.004 for TCGA-MIBC, at 50% missingness) (Supplementary Figure S10). On the other hand, F1-score calculated based on a stable network at a high missingness

ratio showed a slight difference between both datasets, in TCGA-MIBC, SLasso slightly outperformed knnSGLasso (0.619 ± 0.007 and 0.601 ± 0.003 at 50% missingness, respectively). TCGA-PRAD showed the opposite result (0.662 ± 0.002 and 0.702 ± 0.018 at 50% missingness, respectively) (Supplementary Figure S10).

At multi-omics level, knnSGLasso outperformed the rest of the methods followed by CoCoLasso and HMLasso in both, TCGA-MIBC and TCGA-PRAD in low- and high-missingness ratio according with F1-score calculated using whole data as reference (0.947 ± 0.001 , 0.919 ± 0.001 , 0.919 ± 0.001 , respectively at 10% missingness, and 0.656 ± 0.001 , 0.524 ± 0.006 , 0.528 ± 0.005 , respectively, at 50% missingness for TCGA-MIBC; 0.955 ± 0.001 , 0.930 ± 0.001 , 0.930 ± 0.001 , respectively, at 10% missingness, and 0.697 ± 0.017 , 0.593 ± 0.006 , 0.596 ± 0.005 , respectively, at 50% missingness for TCGA-PRAD), and the same measurement using stable network as reference instead (0.858 ± 0.001 , 0.764 ± 0.001 , 0.764 ± 0.001 , respectively, at 10% missingness, and 0.620 ± 0.002 , 0.443 ± 0.005 , 0.447 ± 0.005 , respectively, at 50% missingness for TCGA-MIBC; 0.874 ± 0.001 , 0.778 ± 0.001 , 0.778 ± 0.001 , respectively, at 10% missingness, and 0.663 ± 0.016 , 0.506 ± 0.005 , 0.509 ± 0.005 , respectively at 50% missingness for TCGA-PRAD) (Supplementary Figure S10).

Adaptive weights had a marginal impact on performance with a slight negative effect, irrespective of the dataset. Only at high noise levels did adaptive weights stabilise performance (*Supplementary Figures S6, S7 and S10*).

SLasso and inverse covariance-based approaches perform best for data with block-missingness

When investigating block-missingness, i.e. where entire omics layers are missing for a subset of samples, SLasso performed best when TCGA-BRCA data were used, while CoCoLasso and HMLasso performed best on TCGA-MIBC and TCGA-PRAD data (*Figure 4*; *Supplementary Figures S2 and S10*). SLasso, CoCoLasso and HMLasso reached high consistency with the networks inferred on full data with SLasso achieving F1-scores between 0.945 ± 0.004 at 10% missingness ($\alpha = 0$) and 0.688 ± 0.002 at 50% missingness ($\alpha = 0.5$) on TCGA-BRCA data. CoCoLasso and HMLasso reached F1-scores of 0.967 ± 0.002 at 10% missingness ($\alpha = 0$), and 0.810 ± 0.004 at 50% missingness ($\alpha = 0.5$) respectively on TCGA-MIBC input data (*Figure 4A*; *Supplementary Figure S5 and S6*). The results obtained for TCGA-MIBC were consistently similar to those obtained for TCGA-PRAD with F1-scores of 0.976 ± 0.002 at 10% missingness ($\alpha = 0$), and 0.920 ± 0.004 for TCGA-MIBC and 0.938 ± 0.005 for TCGA-PRAD at 50% missingness ($\alpha = 0$) (*Supplementary Figure S10*).

With higher noise, the performance dropped below an F1-score of 0.2 (*Supplementary Figure S6*). This behaviour could also be observed for GLasso and knnSGLasso, although with generally lower F1-scores. Notably, knnSGLasso networks almost exclusively consisted of gene nodes, while all other methods had a proportional representation of all omics types. The implementation of kNN-imputation used here is not able to handle entire block-missing samples and, consequently, knnSGLasso removes a substantial amount of the features.

Similar behaviour could be observed when comparing the inferred networks on missing data to the stability-selection-based reference. Both SLasso and GLasso outperformed all other methods on TCGA-BRCA data with SLasso reaching the highest F1-scores of 0.763 ± 0.003 (10% missingness, $\alpha = 0$) to 0.655 ± 0.007 (50% missingness, $\alpha = 0.5$), while CoCoLasso and HMLasso outperformed all methods on TCGA-MIBC data reaching F1-scores of 0.773 ± 0.001 (10% missingness, $\alpha = 0$) to 0.728 ± 0.003 (50% missingness, $\alpha = 0.5$; *Figure 4A*). This same pattern was observed on TCGA-PRAD with F1-scores of 0.788 ± 0.001 at 10% missingness ($\alpha = 0$), and 0.763 ± 0.002 for TCGA-MIBC and 0.778 ± 0.002 for TCGA-PRAD at 50% missingness ($\alpha = 0$) (*Supplementary Figure S10*). Adaptive weights (SALasso and GLasso) again had negligible effects and only improved performance markedly at high noise levels as observed in TCGA-MIBC ($\alpha = 1.5$; *Supplementary Figure S7*).

HMLasso and CoCoLasso performed worst on TCGA-BRCA data, reaching average F1-scores of 0.212 ± 0.022 , and 0.246 ± 0.039 at 10% block-missingness, dropping to 0.063 ± 0.033 and 0.073 ± 0.031 respectively with 50% block-missingness and medium noise ($\alpha = 0.5$; *Supplementary Figure S6*). Comparing those methods to the stability-selection-based reference depicted a similar picture (*Supplementary Figure S7*).

knnSGLasso is the least affected by sample size reduction

All methods showed a decrease in concordance already at 10% sample reduction (*Figure 5*). SLasso and knnSGLasso were the most stable method at single-omics level over progressive sample reduction (*Figure 5*). The F1-score of knnSGLasso was

similar to SLasso on TCGA-BRCA data at 10% sample reduction (knnSGLasso: 0.906 ± 0.018 , SLasso: 0.906 ± 0.003). However, at a high sample removal rate (50% sample reduction), SLasso (0.810 ± 0.012) slightly outperformed knnSGLasso (0.796 ± 0.020 ; *Figure 5A*). CoCoLasso and HMLasso reached similar performance as knnSGLasso on TCGA-MIBC and TCGA-PRAD data when 10% sample reduction (CoCoLasso and HMLasso: 0.916 ± 0.005 and 0.949 ± 0.005 , knnSGLasso: 0.917 ± 0.002 and 0.950 ± 0.003 , respectively), but were consistently outperformed by knnSGLasso at higher sample reduction rates (CoCoLasso and HMLasso: 0.782 ± 0.005 and 0.849 ± 0.005 , knnSGLasso: 0.806 ± 0.013 and 0.871 ± 0.003 , respectively; *Figure 5A*; *Supplementary Figure S10*). Similar behaviour but more prominently could be observed when comparing performance based on the stability-selection reference network irrespective of the dataset. knnSGLasso outperformed the rest of the methods followed by SLasso and GLasso. CoCoLasso and HMLasso were among the worst-performing methods (*Figure 5A* and *B*; *Supplementary Figure S10*).

On multi-omics data, the results showed a similar pattern: knnSGLasso consistently outperformed all other methods on TCGA-BRCA data (*Figure 5B*). On TCGA-MIBC and TCGA-PRAD data, CoCoLasso and HMLasso outperformed knnSGLasso at 10% sample reduction (CoCoLasso and HMLasso: 0.959 ± 0.003 and 0.967 ± 0.003 , knnSGLasso: 0.958 ± 0.001 and 0.966 ± 0.002 , respectively), but were overtaken by knnSGLasso with higher reduction rates (CoCoLasso and HMLasso: 0.886 ± 0.005 and 0.907 ± 0.005 , knnSGLasso: 0.893 ± 0.006 and 0.911 ± 0.003 , respectively; *Figure 5B*; *Supplementary Figure S10*). knnSGLasso strongly outperformed all other methods on TCGA-BRCA data when compared to the stability-selection-based reference network (*Figure 5B*). On TCGA-MIBC data, however, SLasso performed best, closely followed by GLasso. CoCoLasso and HMLasso performed worst for no to medium missingness rate, narrowly overtaking SALasso in presence of high missingness rate (CoCoLasso and HMLasso: 0.774 ± 0.001 and 0.756 ± 0.004 ; SALasso 0.776 ± 0.001 and 0.752 ± 0.004 , for 10% missingness and 50% missingness, respectively; *Figure 5B*). Compared with TCGA-MIBC, TCGA-PRAD had similar pattern results being SLasso best performing, closely followed by GLasso in low to high missingness rate (0.877 ± 0.001 and 0.875 ± 0.001 at 10% missingness, 0.854 ± 0.001 and 0.848 ± 0.003 at 50% missingness, respectively). However, at no-missingness, knnSGLasso narrowly overtook SLasso and GLasso (SLasso and GLasso: 0.879 ± 0.0 , knnSGLasso: 0.883 ± 0.0). CoCoLasso and HMLasso performed worst in most missingness rates (*Supplementary Figure S10*).

Inter-layer links possess biological meaning comparable to state-of-the-art methods

We have successfully identified inter-layer links representing potential biological interactions between different molecular abstraction layers. Here, we performed multivariate regression to detect links between gene expression and methylation site measurements, which are expression quantitative trait methylation (eQTM). To evaluate the viability of the eQTM detected by Lasso models dealing with missing data, we compared our results against a state-of-the-art method based on linear regression to eQTM detection, Matrix eQTL [49] (*Figure 5C*). Matrix eQTL detected more eQTM than sparse model-based networks, however, the eQTM-linked genes (number of genes with at least one methylation site significantly related) detected by Matrix eQTL were quite similar but marginally less than those eQTM-linked genes detected by sparse models across the different experimental setups performed in this benchmarking (*Figure 5C*,

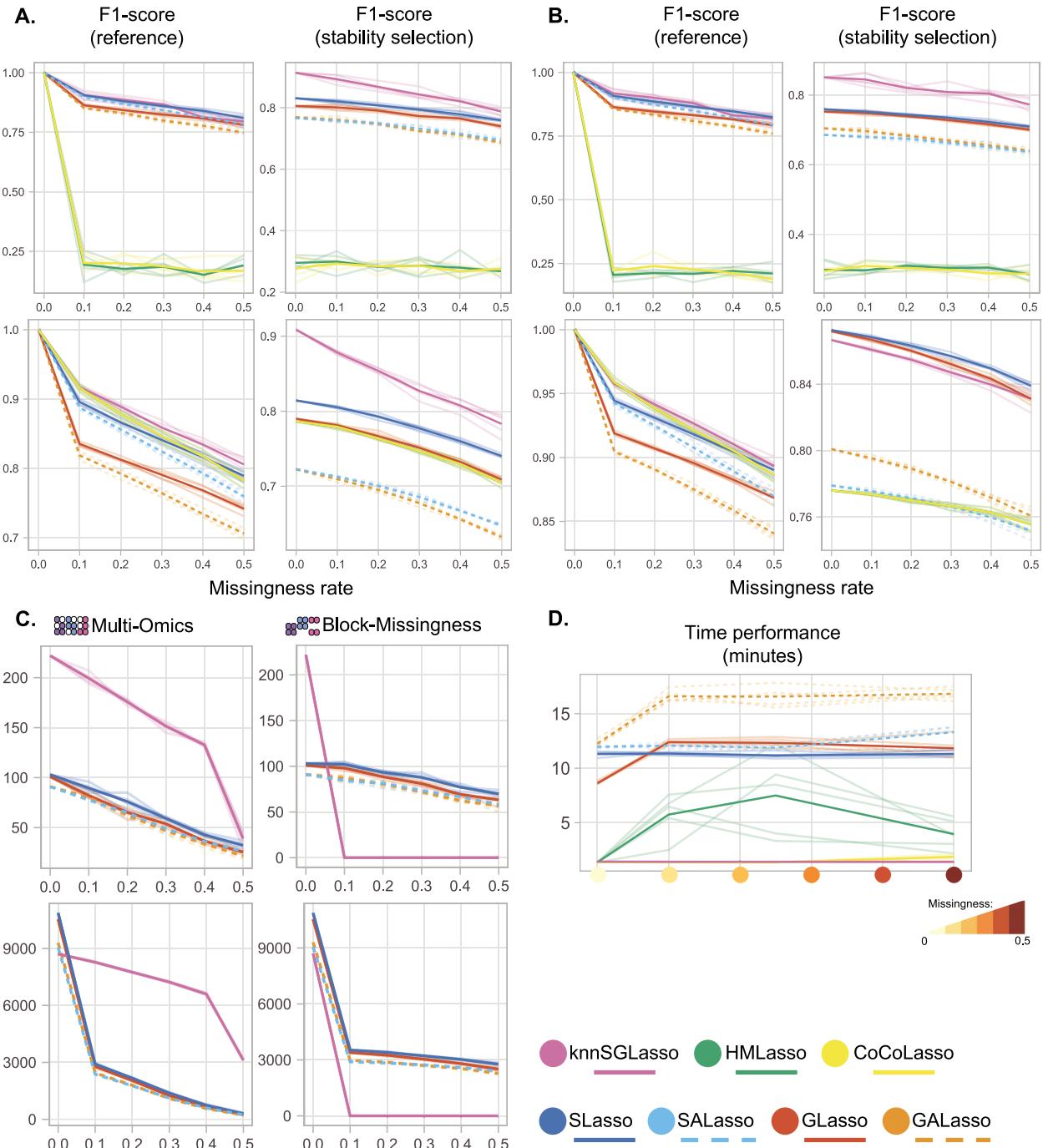


Figure 5. Sample size and inter-layer results across all experimental setups. Transparent lines denote individual runs, while bold lines refer to the average performance. F1-scores for single-omics (**A**) and multi-omics (**B**) on TCGA-BRCA data (top) and TCGA-MIBC data (bottom). (**C**) eQTM-linked genes (number of genes with at least one methylation site significantly related) detected by Matrix eQTL based on the subset of nodes from sparse models inferred networks. (**D**) Runtime evaluation over multi-omics experiments using TCGA-BRCA data.

Supplementary Figure S11). knnSGLasso results detected more eQTM-linked genes at the presence of multi-omics missingness in both datasets at a low missingness rate (199.4 ± 5.03 for TCGA-BRCA and $8,275.6 \pm 29.99$ for TCGA-MIBC, at 10% missingness using Matrix eQTL; 202 ± 4.90 for TCGA-BRCA and $8,286.4 \pm 30.96$ for TCGA-MIBC at 10% missingness using sparse models), as well as at high missingness rate (39.75 ± 8.10 for TCGA-BRCA and $3,117.4 \pm 45.84$ for TCGA-MIBC, at 50% missingness using Matrix eQTL; 41.4 ± 8.41 for TCGA-BRCA and $3,305.8 \pm 44.95$ for TCGA-MIBC, at 50% missingness using sparse models. Except at

whole data performance where SLasso outperformed knnSGLasso followed by GLasso for TCGA-MIBC dataset ($10,875 \pm 0$ for SLasso and $10,540 \pm 0$ for GLasso using Matrix eQTL; $10,881 \pm 0$ for SLasso and $10,547 \pm 0$ for GLasso using sparse models) (Figure 5C, Supplementary Figure S11).

At the presence of block-missingness, Matrix eQTL performance based on SLasso results detected the most eQTM-linked genes at low missingness rate (102 ± 2.92 for TCGA-BRCA and $3,527.2 \pm 12.76$ for TCGA-MIBC, at 10% missingness using Matrix eQTL; 111 ± 2.24 for TCGA-BRCA and $3,552.6 \pm 12.52$ for

TCGA-MIBC, at 10% missingness using sparse models), as well as at high missingness rate (69.8 ± 3.03 for TCGA-BRCA and $2,771 \pm 61.53$ for TCGA-MIBC, at 50% missingness using Matrix eQTL; 83 ± 1.58 for TCGA-BRCA and $2,886.8 \pm 54.41$ for TCGA-MIBC, at 50% missingness using sparse models) (Figure 5C, Supplementary Figure S11).

HMLasso is the fastest approach in datasets with no to medium missingness

knnSGLasso, HMLasso and CoCoLasso had a similar runtime of 82–85 s on the TCGA-BRCA dataset without missing information and gradually increased in runtime with rising missingness (Figure 5D). HMLasso was the fastest approach with an average runtime of 81.843 ± 0.375 s at low to medium missingness levels. Only in scenarios with high missingness, HMLasso was outperformed by knnSGLasso, reaching an average runtime of 97.089 ± 1.327 s. knnSGLasso demonstrated a very consistent runtime across all missingness levels with an average runtime of 97.725 ± 7.995 s. CoCoLasso behaved similarly but was affected by the degree of missingness, reaching maximum average runtimes of 449.6 ± 225.2 s. GALasso and SALasso were the slowest, with average runtimes of 739 ± 15.532 and 919.785 ± 11.833 s on the complete dataset, respectively, of which 241.191 ± 2.975 s were dedicated to imputation. Overall runtime gradually increased for both methods reaching an overall average runtime of $1,148.871 \pm 245.077$ and $1,184.688 \pm 157.794$ s, respectively, of which on average 507.674 ± 58.126 s was spent in imputation. For GALasso, adaptive weights calculations added another 275.363 ± 54.180 s on average to the overall runtime.

DISCUSSION

Due to economic or technical restrictions, missingness of individual values or block-missingness of entire omics layers in a subset of samples is typical for high-throughput multi-omics experiments, rendering multi-omics network inference challenging. In this study, we benchmarked novel regression approaches that can handle missing information across common missingness scenarios in single- and multi-omics experiments and integrated these approaches into KiMONo, a recent approach for network-guided multi-omics network inference.

We observed that the performance of the different methods was dependent on the individual feature size per omics layer. However, some general trends could be identified: kNN-imputation combined with the standard KiMONo approach performed best in most cases except when block-missingness was present. For this particular case, the number of features per omic layer in the dataset was a crucial factor, with SLasso and GLasso being suitable for TCGA-BRCA data, while CoCoLasso and HMLasso were better suited for TCGA-MIBC and TCGA-PRAD data. Both types of methods were also able to handle high degrees of block-missingness consistently well on TCGA-BRCA, TCGA-MIBC and TCGA-PRAD data, respectively, providing an advantage in real-world applications.

One reason why multiple-imputation-based methods provided good performances in the presence of low-dimensional omics-layers such as CNVs in TCGA-BRCA with 84 features, might be due to their utility of the posterior predictive distribution of the missing data and harmonised feature selection over the feature estimates resulting in robust selections more so than point estimates of inverse covariance-based methods. While the original paper of SALasso and GLasso claimed that adaptive weights improved the general performance [28], our results showed a marginal reduction. This discrepancy could be due to the lower

imputation quality of the network-based multiple-imputation approach (ngMICE) we applied, propagating the imputation uncertainty into network inference. Increasing the number of multiple imputations could improve the performance, however, with an increase in runtime.

Multiple imputation in high-dimensional data is generally challenging since existing approaches do not scale to the number of covariates typically encountered in multi-omics studies. Here, dimensionality reduction methods for multiple imputation [50], latent factor models [51] or deep-learning-based approaches [52–54] might improve multiple imputation and therefore network inference quality. However, a benefit of such explicit approaches is their ability to use and adequately address prior imputed datasets often provided by larger consortia.

Implicit methods relying on inverse covariance matrix estimation performed poorly on TCGA-BRCA data but performed generally well on TCGA-MIBC and TCGA-PRAD data, even outperforming knnSGLasso and SLasso/SALasso in presence of block-missingness, indicating that large volumes of information per omic-layer are required to estimate the covariance structure properly with such methods.

Inter-layer links with biological meaning were captured by the different methods across different missingness rates as the number of genes linked to methylation sites analysis demonstrated when compared to state-of-the-art methods such as Matrix eQTL [49] which is dedicated to eQTM detection.

Matrix eQTL detected more eQTM than sparse models-based networks as expected given this method uses inference based on linear regression for all possible combinations of genes-methylation sites as background procedure [49]. However, the number of total genes significantly involved in at least one eQTM was marginally less than those detected by the Lasso inferred network in presence of missing data. Therefore, we proved the biological value of inter-layer links detected across different network inferences and showed the power of the methods to detect new biological insights. Furthermore, knnSGLasso conserved at least twice as many eQTM-linked genes as the rest of the methods including when no-missingness data were used as input in TCGA-BRCA. The absence of advanced methods for multiple imputation, which optimally encompassed high imputation quality and low computational cost were factors that affected SLasso/SALasso and GLasso/GALasso in this task. However, the approach used here (ngMICE) demonstrated to be suitable across increasing missingness rates at the point that SLasso and GLasso outperformed the rest of the methods in certain conditions, such as block missingness in the TCGA-BRCA dataset.

Comparing similar datasets, such as TCGA-MIBC and TCGA-PRAD, and analysing a comparable dataset, such as TCGA-BRCA, through the implementation of different prior networks allowed for the characterisation of the strengths and weaknesses of different sparse models in network inference. We expect these findings to hold in other datasets, such as colorectal cancer from TCGA used by Welz et al. [55] for network inference using KiMONo based on additional omics layers evaluated in the present benchmarking such as proteomics and DNA mutation [55]. This includes the improvement of recent frameworks such as DiffBrainNet, which includes differential expression analysis and prior-knowledge network inference (KiMONo) to study stress response in mouse brain [56].

We note that a true gold standard for evaluating the performance of network inference methods is missing. In its absence, we rely on networks inferred from complete, unperturbed data that nevertheless are likely to contain both false positive and

false negative interactions which may affect the results. Hence, our reference networks are not suited for evaluating methods following different principles for inference as GENIE3 [57], or other nonlinear approaches such as ARACNe and their derivatives [7, 8]. Nevertheless, the methods tested here can be used as direct substitutes for other Lasso-based network inference frameworks such as wgLasso [13], pLasso [58] or PoLoBag [59] with similar expected behaviour as demonstrated here. Also, other linear methods such as WGCNA [5] or Petal [60] could benefit from two-step approaches and consistent feature selection across multiple imputed datasets, similar to G(A)Lasso. A way of such approaches to handle multiple imputed datasets could be through the usage of Rubin's rule after Fisher's Z-transformation, which allows to pool the estimated correlation coefficients between omics features from multiple imputed datasets in a consistent manner [61, 62].

In summary, we found explicit methods to be more robust than methods implicitly handling missingness in most scenarios except on block-missingness. While most methods were tolerant to high levels of missingness, they were strongly affected by noise. While HMLasso was the fastest tested method, knnSGLasso showed the best trade-off between performance and runtime and is thus our recommended approach for handling missingness in KiMONo except for consistently high-dimensional omics-layers data with block-missingness where we recommend HMLasso instead. While we see room for further method improvements, particularly with respect to multiple imputations of high-dimensional single matrices and the robustness of inverse covariance methods, our results show that robust multi-omics network inference in the presence of missingness is feasible with KiMONo and thus allows users to leverage available multi-omics data to their fullest extent.

Key Points

- We extended KiMONo, a multi-omics network inference approach that is using prior network information, to handle data with missingness by integrating advanced Lasso-based inference techniques.
- knnSGLasso outperformed all other methods in the majority of scenarios except in presence of block-missingness where the kNN-imputation step tends to fail. In such cases, we recommend using SLasso on imbalanced and HMLasso on balanced omics layers dimensional data.
- Inverse covariance matrix estimation-based approaches were considerably sensible to the dimensionality of the input data, requiring high-dimensional data per omics layer to reach high performance.

SUPPLEMENTARY DATA

Supplementary data are available online at <https://academic.oup.com/bib/article/24/5/bbad309/7260520>

AUTHORS CONTRIBUTIONS

JDH, MiL, MA and AG implemented the different Lasso approaches. JDH implemented the benchmark framework. JDH and BS analysed the results. BS, CO and MaL conceived and supervised the study. JDH, BS, CO and MaL wrote the original manuscript.

MaL, CO, FJT, BS and MiL reviewed and edited the manuscript. All authors read and approved the final manuscript.

CONFLICT OF INTEREST

FJT reports receiving consulting fees from Roche Diagnostics GmbH and Cellarity Inc., and an ownership interest in Cellarity, Inc. The remaining authors declare no competing interests.

FUNDING

German Centre of Lung Research (DZL). Helmholtz International Lab 'Causal Cell Dynamics' to B.S. Hanns Seidel Foundation to MiL. BMBF (German Federal Ministry of Education and Research) Project TRY-IBD (Grant 01ZX1915B) for C.O. and A.G.

REFERENCES

1. Romero IG, Ruvinsky I, Gilad Y. Comparative studies of gene expression and the evolution of gene regulation. *Nat Rev Genet* 2012;13:505–16.
2. Lee J, Hyeon DY, Hwang D. Single-cell multiomics: technologies and data analysis methods. *Exp Mol Med* 2020;52:1428–42.
3. Li Y, Ma L, Wu D, Chen G. Advances in bulk and single-cell multi-omics approaches for systems biology and precision medicine. *Brief Bioinform* 2021;22:bbab024.
4. Hawe JS, Theis FJ, Heinig M. Inferring interaction networks from multi-omics data. *Front Genet* 2019;10:535.
5. Langfelder P, Horvath S. WGCNA: an R package for weighted correlation network analysis. *BMC Bioinformatics* 2008;9:559.
6. Song L, Langfelder P, Horvath S. Comparison of co-expression measures: mutual information, correlation, and model based indices. *BMC Bioinformatics* 2012;13:328.
7. Margolin AA, Nemenman I, Basso K, et al. ARACNE: an algorithm for the reconstruction of gene regulatory networks in a mammalian cellular context. *BMC Bioinformatics* 2006;7(Suppl 1):S7.
8. Lachmann A, Giorgi FM, Lopez G, Califano A. ARACNe-AP: gene network reverse engineering through adaptive partitioning inference of mutual information. *Bioinformatics* 2016;32:2233–5.
9. Krumsiek J, Suhre K, Illig T, et al. Gaussian graphical modeling reconstructs pathway reactions from high-throughput metabolomics data. *BMC Syst Biol* 2011;5:21.
10. Schäfer J, Strimmer K. An empirical Bayes approach to inferring large-scale gene association networks. *Bioinformatics* 2005;21:754–64.
11. Petralia F, Wang P, Yang J, Tu Z. Integrative random forest for gene regulatory network inference. *Bioinformatics* 2015;31:i197–205.
12. Sass S, Buettner F, Mueller NS, Theis FJ. A modular framework for gene set analysis integrating multilevel omics data. *Nucleic Acids Res* 2013;41:9622–33.
13. Li Y, Jackson SA. Gene network reconstruction by integration of prior biological knowledge. *G3 (Bethesda)* 2015;5:1075–9.
14. List M, Dehghani Amirabad A, Kostka D, Schulz MH. Large-scale inference of competing endogenous RNA networks with sparse partial correlation. *Bioinformatics* 2019;35:i596–604.
15. Ogris C, Hu Y, Arloth J, Müller NS. Versatile knowledge guided network inference method for prioritizing key regulatory factors in multi-omics data. *Sci Rep* 2021;11:6806.
16. Rubin DB. Multiple Imputation for Nonresponse in Surveys. John Wiley & Sons. Hoboken, NJ, USA; 2004.

17. Donders ART, van der Heijden GJMG, Stijnen T, Moons KGM. Review: a gentle introduction to imputation of missing values. *J Clin Epidemiol* 2006; **59**:1087–91.
18. Wood AM, White IR, Royston P. How should variable selection be performed with multiply imputed data? *Stat Med* 2008; **27**: 3227–46.
19. Ganti R, Willett RM. Sparse linear regression with missing data. arXiv [stat.ML]. 2015.
20. Ibrahim JG, Chen M-H, Kim S. Bayesian variable selection for the cox regression model with missing covariates. *Lifetime Data Anal* 2008; **14**:496–520.
21. Yang X, Belin TR, Boscardin WJ. Imputation and variable selection in linear regression models with missing covariates. *Biometrics* 2005; **61**:498–506.
22. Heymans MW, van Buuren S, Knol DL, et al. Variable selection under multiple imputation using the bootstrap in a prognostic study. *BMC Med Res Methodol* 2007; **7**:33.
23. Liu Y, Wang Y, Feng Y, Wall MM. Variable selection and prediction with incomplete high-dimensional data. *Ann Appl Stat* 2016; **10**:418–50.
24. Wan Y, Datta S, Conklin DJ, Kong M. Variable selection models based on multiple imputation with an application for predicting median effective dose and maximum effect. *J Stat Comput Simul* 2015; **85**:1902–16.
25. Chen Q, Wang S. Variable selection for multiply-imputed data with application to dioxin exposure study. *Stat Med* 2013; **32**: 3646–59.
26. Geromini J, Saporta G. Variable selection for multiply-imputed data with penalized generalized estimating equations. *Comput Stat Data Anal* 2017; **110**:103–14.
27. Marino M, Buxton OM, Li Y. Covariate selection for multilevel models with missing data. *Stat* 2017; **6**:31–46.
28. Du J, Boss J, Han P, et al. Variable selection with multiply-imputed datasets: choosing between stacked and grouped methods. *J Comput Graph Stat* 2022; **31**:1–35.
29. Choi Y, Tibshirani R. An investigation of methods for handling missing data with penalized regression. arXiv [stat.AP]. 2013.
30. Johnson BA, Lin DY, Zeng D. Penalized estimating functions and variable selection in semiparametric regression models. *J Am Stat Assoc* 2008; **103**:672–80.
31. Loh P-L, Wainwright MJ. High-dimensional regression with noisy and missing data: provable guarantees with non-convexity. *Adv Neural Inf Process Syst* 2011; **40**:24.
32. Städler N, Bühlmann P. Missing values: sparse inverse covariance estimation and an extension to sparse regression. *Stat Comput* 2012; **22**:219–35.
33. Takada M, Fujisawa H, Nishikawa T. HMLasso: lasso with high missing rate. arXiv [stat.ML]. 2018.
34. Datta A, Zou H. CoCoLasso for high-dimensional error-in-variables regression. *Ann Statistics* 2017; **45**:2400–26.
35. Shen C-W, Chen Y-H. Model selection for generalized estimating equations accommodating dropout missingness. *Biometrics* 2012; **68**:1046–54.
36. Sabbe N, Thas O, Ottoy J-P. EMLasso: logistic lasso with missing data. *Stat Med* 2013; **32**:3143–57.
37. Yu G, Li Q, Shen D, Liu Y. Optimal sparse linear prediction for block-missing multi-modality data without imputation. *J Am Stat Assoc* 2020; **115**:1406–19.
38. Xue F, Qu A. Integrating multisource block-wise missing data in model selection. *J Am Stat Assoc* 2021; **116**: 1914–27.
39. Gentry AE, Kirkpatrick RM, Peterson RE, Webb BT. Missingness adapted group informed clustered (MAGIC)-LASSO: a novel paradigm for phenotype prediction to improve power for genetic loci discovery. *Front Genet* 2023; **14**:1162690.
40. Simon N, Friedman J, Hastie T, Tibshirani R. A sparse-group lasso. *J Comput Graph Stat* 2013; **22**:231–45.
41. Henao JD, Lauber M, Azevedo M, et al. Multi-Omics Regulatory Network Inference in the Presents of Missing Data, 2022.
42. Troyanskaya O, Cantor M, Sherlock G, et al. Missing value estimation methods for DNA microarrays. *Bioinformatics* 2001; **17**:520–5.
43. Escribe C, Lu T, Keller-Baruch J, et al. Block coordinate descent algorithm improves variable selection and estimation in error-in-variables regression. *Genet Epidemiol* 2021; **45**:874–90.
44. Weinstein JN, Collisson EA, Mills GB, et al. The cancer genome atlas pan-cancer analysis project. *Nat Genet* 2013; **45**:1113–20.
45. Gao J, Lindsay J, Watt S, et al. Abstract 5277: the cBioPortal for cancer genomics and its application in precision oncology. *Cancer Res* 2016; **76**:5277–7.
46. Oughtred R, Rust J, Chang C, et al. The BioGRID database: a comprehensive biomedical resource of curated protein, genetic, and chemical interactions. *Protein Sci* 2021; **30**:187–200.
47. Ogris C, Guala D, Sonnhammer ELL. FunCoup 4: new species, data, and visualization. *Nucleic Acids Res* 2018; **46**:D601–7.
48. van Buuren S, Groothuis-Oudshoorn K. Mice: multivariate imputation by chained equations in R. *J Stat Softw* 2011; **45**:1–67.
49. Shabalin AA. Matrix eQTL: ultra fast eQTL analysis via large matrix operations. *Bioinformatics* 2012; **28**(10):1353–8.
50. Hodge DW, Safo SE, Long Q. Multiple imputation using dimension reduction techniques for high-dimensional data. arXiv [stat.ME]. 2019.
51. Argelaguet R, Arnol D, Bredikhin D, et al. MOFA+: a statistical framework for comprehensive integration of multi-modal single-cell data. *Genome Biol* 2020; **21**:111.
52. Qiu YL, Zheng H, Gevaert O. Genomic data imputation with variational auto-encoders. *Gigascience* 2020; **9**.
53. Gayoso A, Steier Z, Lopez R, et al. Joint probabilistic modeling of single-cell multi-omic data with totalVI. *Nat Methods* 2021; **18**: 272–82.
54. Lotfollahi M, Litinetskaya A, Theis FJ. Multigate: single-cell multi-omic data integration. bioRxiv. 2022; 2022.03.16.484643.
55. Welz L, Kakavand N, Hang X, et al. Epithelial X-box binding protein 1 coordinates tumor protein p53-driven DNA damage responses and suppression of intestinal carcinogenesis. *Gastroenterology* 2022; **162**(1):223–237.e11.
56. Gerstner N, Krontira AC, Cruceanu C, et al. DiffBrainNet: differential analyses add new insights into the response to glucocorticoids at the level of genes, networks and brain regions. *Neurobiol Stress* 2022; **21**:100496.
57. Huynh-Thu VA, Irrthum A, Wehenkel L, Geurts P. Inferring regulatory networks from expression data using tree-based methods. *PloS One* 2010; **5**:e12776.
58. Tibshirani R, Friedman J. A pliable lasso. *J Comput Graph Stat* 2020; **29**(1):215–25.
59. Ghosh Roy G, Geard N, Verspoor K, He S. PoLoBag: polynomial lasso bagging for signed gene regulatory network inference from expression data. *Bioinformatics* 2021; **36**(21):5187–93.
60. Petereit J, Smith S, Harris FC, Jr, Schlauch KA. Petal: co-expression network modelling in R. *BMC Syst Biol* 2016; **10** (Suppl 2):51.
61. Marshall A, Altman DG, Holder RL, Royston P. Combining estimates of interest in prognostic modelling studies after multiple imputation: current practice and guidelines. *BMC Med Res Methodol* 2009; **9**:57.
62. Panken AM, Heymans MW. A simple pooling method for variable selection in multiply imputed datasets outperformed complex methods. *BMC Med Res Methodol* 2022; **22**(1):214.