

Towards Causal Analysis of Genetic Factors for Colorectal Cancer

David Miller

causalens

Hana Chockler (✉ hana@causalens.com)

causalens <https://orcid.org/0000-0003-1219-0713>

Andrew Lawrence

causalens

Daniel McNamee

causalens

Nicholas Chia

Mayo Clinic

Brief Communication

Keywords: colorectal cancer, genetic factors, dependencies

Posted Date: October 22nd, 2021

DOI: <https://doi.org/10.21203/rs.3.rs-967255/v1>

License:  This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

Towards Causal Analysis of Genetic Factors for Colorectal Cancer

David W. Miller¹, Hana Chockler^{*1,2}, Andrew R. Lawrence¹, Daniel C. McNamee¹, and
Nicholas Chia^{3,4,5}

¹causaLens, London, UK

¹{*david.miller,hana,andrew,dan*}@causalens.com

²Department of Informatics, King's College London, UK

³Microbiome Program, Center for Individualized Medicine, Mayo Clinic, Rochester, MN,
USA

⁴Division of Surgical Research, Department of Surgery, Mayo Clinic, Rochester, MN, USA

⁵Department of Physiology and Biomedical Engineering, Mayo Clinic, Rochester, MN, USA

³*chia.nicholas@mayo.edu*

Abstract

We present preliminary results and outline future work for the causal analysis of genetic factors influencing the occurrence and severity of colorectal cancer. The findings are based on publicly available datasets. We show that using relatively simple methods we are able to detect meaningful dependencies, reflecting the current biomedical knowledge about the causes of cancer and related conditions (such as infection and sepsis). We also argue that a deeper analysis, taking into account gene regulation, would require a significantly larger dataset.

1 Introduction and Motivation

Colorectal cancer is common (147,950 new cancers annually) and lethal (over 53,000 deaths each year) [1]. In order to combat this deadly disease, screening colonoscopy, where the goal is to identify and remove precursor lesions that may lead to colorectal cancer, has been used for cancer prevention for almost fifty years [2]. The success of this approach rests largely on the recognition that adenomas are causally linked to carcinogenesis. By recognizing the causal lesions and removing them, we effectively prevent cancer from ever forming, leading to an ideal approach to managing an incredibly deadly and common disease.

The power of this causal understanding is clear. However, this approach has limitations. The colorectal cancer incidence rate appears to be increasing in younger individuals (2% annually for ages 50 and under [1]), making the reliance on screening colonoscopy both more challenging and less cost-effective with a broader population to screen. In addition, the only recourse for preventing disease in higher-risk patients has been shorter and shorter screening intervals. Despite these efforts, interval cancers are reported at frequencies of 3.4 – 7.6% [3, 4], even in the face of adequate bowel preparation and seasoned endoscopists [5]. For these reasons, reliable early detection via non-invasive methods, such as blood or stool testing, would significantly improve the long-term prospects of patients with colorectal cancer.

Problem Statement In this paper, we make the first steps towards a causal analysis of the expression of genes and proteins on the manifestation and stage of colorectal cancer. In order to avoid observer bias [6], we adopt a strictly data-driven approach. Therefore, we do not consider any prior knowledge or domain expertise when performing this initial analysis.

*Corresponding author

39 The data is obtained from The Cancer Genome Atlas (TCGA) project, specifically the Gene Expression
40 by RNAseq (IlluminaHiSeq) data and RPPA Protein Expression data from the Colon and Rectal Cancer
41 (COADREAD) cohort. The TCGA COADREAD data was downloaded via the UCSC Xena platform [7].
42 The main challenge in analysing this data lies in the high number of candidate RNA expressions and proteins
43 compared to the orders of magnitude smaller number of patients. Additionally, there is not a perfect
44 correspondence between the RNA and protein data tables, which reduces the total number of patients
45 further.

46 An additional challenge arises from the extreme imbalance in the data, as the dataset does not contain
47 any healthy individuals; therefore, all the observations are from patients with varying severity of colorectal
48 cancer. This extreme imbalance within the dataset can lead to spurious correlations and to missing causal
49 connections [8]; it is inherent to the medical domain, as healthy individuals do not undergo colon resection
50 surgery, and hence any meaningful progress would need to address this imbalance first.

51 A final challenge is due to the inability to intervene. A common method to discover the underlying causal
52 structure of a system is by performing interventional experimentation, often in the form of randomized control
53 trials. These are ubiquitous in the medical domain, particularly around testing for drug efficacy. However,
54 it is not possible (or ethical) to intervene in this scenario to give someone colorectal cancer or change the
55 concentration of specific RNA molecules or proteins in their cells. Therefore, we must perform this analysis
56 using exclusively observational data.

57 **Related Work** There has been a recent rise in the effort to apply AI techniques, and in particular machine
58 learning, to biological data, and specifically to cancer research. The unique challenges posed by the cancer
59 data, namely that it is high-dimensional and low in sample size, as well as the fact that, similarly to
60 other biological data, interventions are not possible or not ethical, motivate different approaches than the
61 traditional causal inference ones.

62 A recent paper on machine learning for cancer research presents a platform for causal inference suitable
63 for this type of data and its application to colorectal cancer, with the goal of determining causal drivers that
64 differentiate between two subtypes of colorectal cancer [9]. Kalantari et al. [10] use inverse reinforcement
65 learning to gain insights about cancer progression from genome data. Farahmand et al. [11] use causal
66 inference to identify transcriptional regulators.

67 2 Background and General Approach

68 The motivation underlying our work is that a principled approach based on *causal counterfactual reason-*
69 *ing* [12] is the correct method for obtaining meaningful results. This approach is based on postulating the
70 existence of an underlying causal model, encompassing the causal dependencies between variables. The pro-
71 cess of causal discovery [13] amounts to discovering this model by performing observations and experiments
72 on the variables. Essentially, the existence of a causal dependency between variables is demonstrated by
73 intervening on the candidate cause and observing the changes in the outcome variable. This approach can
74 be extended to interventions on more than one variable, where we observe the effect of changes in several
75 variables on the outcome. If direct interventions are impossible, as they are in this dataset, we can attempt
76 to partially replace it by the analysis of the existing data, where we search for the records that differ only
77 in the values of candidate causes and examine the value of the outcome. This approach requires a large
78 and sufficiently varied dataset. In particular, the number of features (candidate causes) should be small in
79 comparison to the number of records, to guarantee a sufficient variation in their values. In this section, we
80 outline existing approaches to estimate causality from observations and describe the approach we propose
81 for the analysis of the colorectal cancer dataset.

82 **Intractability of causality** The computational cost of causality-based methods is high and increases
83 dramatically following even a modest increase in the number of variables (p). For our case $p = 20,531$.
84 As we are unable to intervene, the goal is to perform causal discovery on observational-only data, which
85 in general, is an NP-hard problem [14]. A brief overview highlighting the computational complexities of
86 popular methods follows; please see [13] for a thorough review of causal discovery.

87 Constraint-based methods, such as the Peter and Clark (PC) [15] and fast causal inference (FCI) [16, 17]
88 algorithms, utilize conditional independence tests to discover the causal structure. At the worst case (dense
89 causal structure), these scale exponentially in the number of variables ($\mathcal{O}(c^p)$). This does not capture the
90 complexity of the conditional independence tests either, which may be expensive to calculate when trying
91 to identify more complex relationships between random variables than a simple linear one.

92 Score-based methods, which include greedy, e.g., [18], and exhaustive, e.g., [19], search methods, are
93 another key category. These may scale worse than constraint-based methods as some require all permutations
94 to be enumerated, which equals the factorial of the number of variables ($\mathcal{O}(p!)$). Hybrid methods combining
95 constrained-based and score-based methods have also been proposed.

96 Finally, methods utilizing functional causal models, such as LiNGAM [20] and NOTEARS [21], are also
97 popular but they make assumptions on the types of functional dependencies between variables. To allow for
98 faster convergence, simple linear models are often used, which greatly limits the expressiveness of the models
99 and likely misses complex relationships between variables.

100 The functional causal model methods also scale poorly with respect to the number of variables. For
101 example, the costs of the DirectLiNGAM [22] and ICA-LiNGAM [20] algorithms scale as at least $\mathcal{O}(np^3 + p^4)$,
102 where n is the number of samples/observations. If $p \gg n$, as we have in this case as $n = 283$, then twice
103 the number of RNA molecules considered results in sixteen times the computational cost.

104 One cannot blindly apply popular causal discovery methods to this problem. In order to find meaning-
105 ful drivers of protein expression, we must restrict the number of RNA candidates from 20,531 to a more
106 manageable quantity.

107 **Assumptions** Two common assumptions made when performing causal discovery are: (1) Causal Markov
108 Condition and (2) Causal Faithfulness Condition [15]. We make these assumptions in our analysis. Without
109 going into the formal mathematical definitions, when combined it implies “that two [random] variables
110 are directly causally related *if and only if* they are not conditionally independent given any subset of the
111 remaining variables” [19]. This subset can be the empty set. Therefore, a common first step of constraint-
112 based causal discovery is to measure if two random variables are independent given the empty set. We
113 utilize Spearman’s rank correlation as this independence test. It makes no assumption on the underlying
114 probability distributions of the random variables and only assumes a possible monotonic dependency between
115 the variables, which is a natural assumption for this data. The null hypothesis is that the joint distribution
116 factorizes, i.e., $p(x, y) = p(x)p(y)$, which means x and y are independent. Calculating a p-value from the
117 null distribution represents the likelihood of calculating the measured correlation coefficient given that the
118 null hypothesis is true. Therefore, a really small p-value means it is unlikely that the null hypothesis is true
119 and we can assume the two random variables are dependent.

120 **Proposed approach** We performed several types of analysis as detailed below in § 3. We calculated
121 Spearman’s rank correlation between all pairs of RNA and protein expressions focusing on the ataxia-
122 telangiectasia mutated (ATM) protein as our primary example. The ATM protein kinase has been extensively
123 studied for its role in the DNA damage response, and there is increasing evidence that ATM plays an
124 important role in other cellular processes, including carbon metabolism. Carbon metabolism is highly
125 dysregulated in cancer due to the increased need for cellular biomass. Some therapeutic strategies for cancer
126 involve the development of ATM inhibitors. We also computed an intercorrelation between RNA expressions
127 associated with a given protein (specifically, ATM) and examined the resulting clusters. The expectation is
128 that the clusters provide a smaller space to perform more expensive causal discovery techniques. Finally,
129 we analysed the bimodal activations, which infer whether an RNA expression is on or off in a given record,
130 based on the bimodal graph of the correlation with the protein.

131 3 Detailed Methodology and Results

132 3.1 Data characteristics

133 The TCGA Gene Expression dataset contains the log-normalised gene-level transcription estimates of 20,531
134 distinct RNA molecules for 434 distinct individuals, while the Protein Expression dataset contains the RBN-

135 normalized RPPA values of 131 distinct proteins for 464 individuals.

136 We join the datasets by individual, however despite being from the same cohort of patients, there is
137 not a complete overlap between individuals in the two. Only 283 individuals were present in both datasets,
138 reducing the overall amount of data available for analysis by roughly a third.

139 3.2 Bivariate filtering by independence between RNA and protein expression

140 As a first measure to infer relationships between RNA expression and protein expression, we look at the
141 bivariate correlations calculated across individuals as a statistical test to infer independence. We calculate
142 Spearman’s rank correlation [23] and associated significance between all pairs of RNA and protein expressions.
143 Spearman’s rank correlation between two random variables can be measured by calculating the Pearson
144 Correlation [24] between the ranks of the values in the observations of each random variable.

145 For each of the 131 protein expressions in the dataset we look at which RNA molecules expressions exhibit
146 a statistically significant correlation. We then use this significance as a filter for RNA molecules independent
147 to the protein, in order to reduce the number of features to be included in the problem.

148 Transforming from raw values to ranks reduces the impact of outliers and non-linearity on the Pearson’s
149 R correlation metric, making Spearman Rank Correlation a more robust measure for use in noisy datasets.

150 Table 1 shows the percentile of significant RNA relationship counts across proteins, depending on the p-
151 value chosen as a cut-off between dependence and independence. For example, using a p-value cutoff of 0.01
152 (1e-2), the median (50th percentile) protein has 2492 significant relationships amongst the 20,531 possible
153 RNA molecules.

percentile p-value	0.05	0.25	0.50	0.75	0.95
1e-1	3610	5104	6498	7774	10110
1e-2	697	1526	2492	3841	6254
1e-3	118	450	901	2172	4240
1e-4	18	118	318	1193	3038
1e-5	2	30	102	678	2262

Table 1: Percentiles of significant RNA relationship counts, given p-value, across all proteins

154 In statistical analysis it is common practice to use p-values to determine the reliability of a correlation
155 between two variables. In practice, a limit of 0.05 is often used as a cutoff, representing a maximum 5%
156 chance that a relationship is caused by random noise rather than by a true correlation. However, the
157 limitations of this approach are significant and well documented, especially when managing datasets with a
158 large number of variables.

159 As we can see from Table 1, selecting by strength of correlation or significance requires extremely low
160 p-value cutoff to filter to a reasonable numbers of relevant RNA expression variables.

161 The issue here is two-fold. Firstly, and most widely understood, by simple definition we can expect 1%
162 of the relationships that survive the significance cutoff will not be true at all, and will only appear to be
163 correlated by chance, given noise in our dataset. Secondly, and more importantly, an unknown proportion
164 of the remaining 99% relationships will not describe true causal relationships (where expression of the RNA
165 molecule in question drives Protein expression), but instead describe confounded relationships (where both
166 expression of the protein and of the RNA molecule are driven by a third variable). An example of a
167 confounded relationship would be where both the RNA and protein are produced in separate parts of the
168 same physiological pathway, however, it is also possible that both RNA and protein expression are driven
169 by some other, unmeasured, variable.

170 Using a p-value cutoff of 0.01 we found that the expressions of 100 different proteins were significantly
171 correlated with the expressions of more than 1000 different RNA molecules. Conservatively, we can expect
172 990 ‘true’ correlations per protein in this case. While the pathways that include generating proteins may be
173 complex, we cannot claim to have produced a practical range of RNA candidates to investigate in detail.

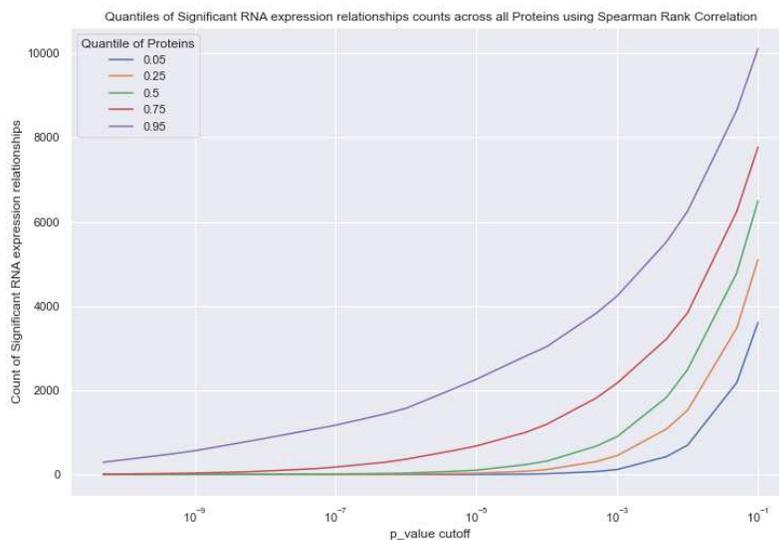


Figure 1: Percentiles of significant RNA relationship counts, given p value, across all proteins

174 A straightforward next step might be to decrease the p-value cutoff further, to limit the set of variables
 175 that we need to consider. However, this risks filtering out true relationships that are have weaker correlations,
 176 either due to a complex relationship or because they are subject to noise or small sample size effects.

177 At this point, simple Bivariate Filtering using expression correlations has reached the limits of its use
 178 and we are forced to consider more sophisticated methods to remove spurious correlations.

179 3.3 Intercorrelation between RNA expressions associated with a given protein

180 To investigate this problem further, we narrowed our focus to the expression of individual proteins. For the
 181 purposes of this communication we discuss results relating to expression of ATM, as mentioned in § 2.

182 For the purposes of illustration, we focus on the intercorrelation between the expressions of 50 RNA
 183 molecules that are most correlated with ATM expression.

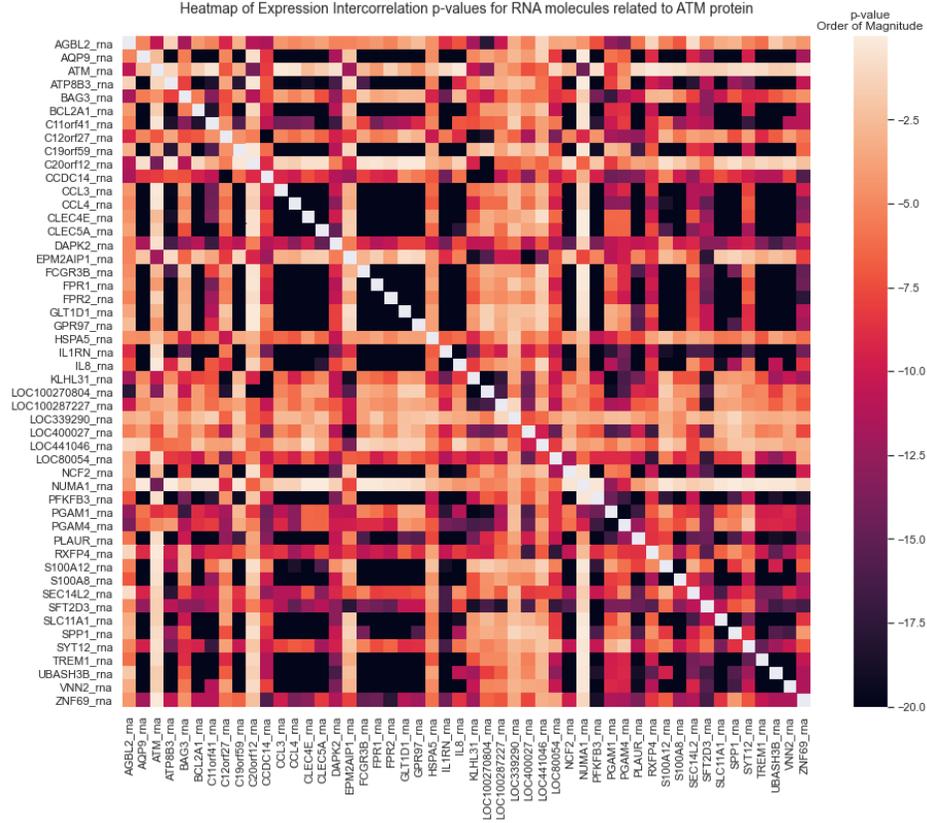


Figure 2: Expression Intercorrelation Heatmap for top 50 RNA molecules related to ATM protein

184 Figure 2 shows a heatmap of the intercorrelation p-values (log scale) for these RNA molecule expressions.
 185 The multiple dark patches show that many RNA molecules' expressions are highly correlated.

186 Interestingly, we note that many of highly inter-correlated RNA molecule expressions are grouped together
 187 into larger blocks. This is due to the fact that the RNA molecule columns are ordered alphabetically, and
 188 molecules with similar names appear to share a similar function (i.e. with stronger intercorrelations).

189 Although we are selecting for RNA expressions that are all strongly correlated with ATM expression, and
 190 therefore should expect some overlap between the top most correlated features, we have placed no bounds
 191 on the sign or the strength of that correlation, only its significance, and the existence of multiple separate
 192 and non-correlating blocks shows that there is a high redundancy in the dataset beyond the RNA molecules'
 193 shared expression relationship with the ATM protein.

194 Further, such blocks of clustered RNA molecules imply that information is being shared within them.
 195 This type of relationship hints at a confounder within each block, or a hidden confounder not included in
 196 the dataset, which is driving the remaining block members' expressions.

197 This shows how the Spearman Rank correlations between RNA molecules' expressions can identify sim-
 198 ilarities between the RNA molecules without any prior data, and implies that such intercorrelations can be
 199 used to identify RNA molecules that share a confounder.

200 3.4 Clustering of RNA expressions for Dimensionality Reduction

201 In order to use these intercorrelations to reduce the dimensionality of the problem, we propose clustering
 202 similar RNA molecules together according to their expressions within the COADREAD Colon and Rectal
 203 Cancer dataset.

204 Agglomerative Clustering, or Hierarchical Clustering, is a method that produces clusters by repeatedly
 205 merging the closest individuals or existing clusters according to some method to determine cluster similarity,

206 based on some distance metric (between individuals) [25]. Agglomerative clustering allows intuitive visual-
 207 isation of clusters, and the ability to dynamically choose the number of clusters to be generated, based on
 208 the value of each new merge.

209 By choosing an appropriate number of clusters, we can group together the RNA molecules that share
 210 extremely similar patterns of expression in the dataset, and take a representative, or an aggregate, from
 211 each. Applying an initial causal analysis on this smaller dataset will then allow us to eliminate entire
 212 clusters thought to be non-causally or spuriously correlated, and focus on smaller clusters of potential causal
 213 drivers. Reducing the dimensionality of the dataset makes the application of causal tools feasible in terms
 214 of computational costs.

215 To produce a cluster hierarchy for the top 50 RNA molecules associated with the ATM protein, we use
 216 the Ward method to determine cluster similarity, based on a measure of Euclidean distance across the RNA
 217 expression per individual within the dataset. Figure 3 shows the dendrogram representing this hierarchy.

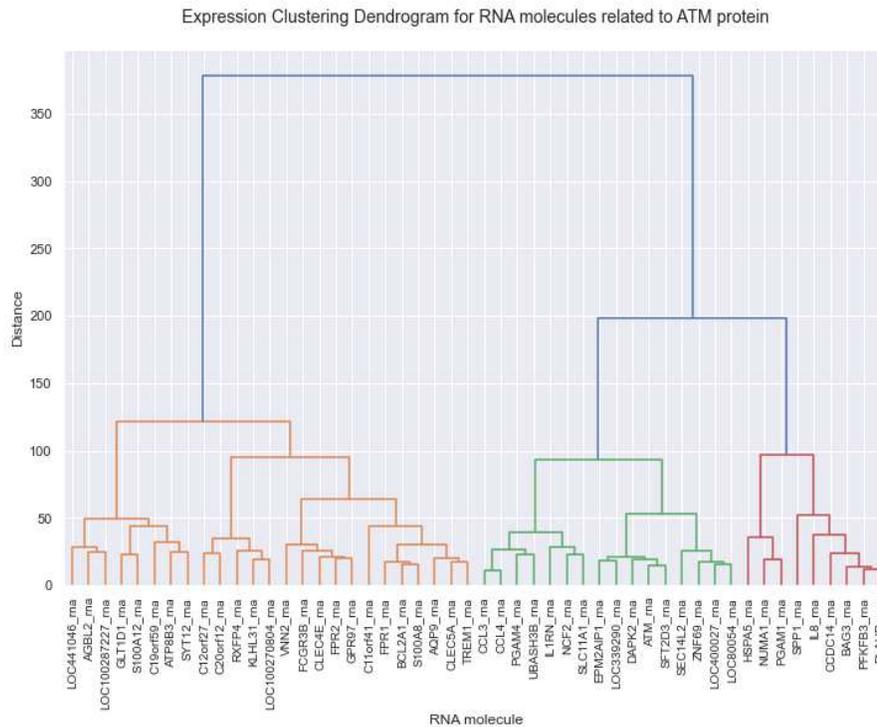


Figure 3: Expression Clustering Dendrogram of top 50 RNA molecules associated with the ATM protein

218 The inverted-U shaped joins on a Dendrogram show merges between clusters, and the height of the legs
 219 of these merges represents the difference between the two clusters. For example, we can see how the CCL3
 220 and CCL4 RNA molecules (green) are considered closer to each other than the nearest neighbouring pair
 221 of RNA molecules, PGAM4 and UBASH3B, and that the algorithm only slightly prefers to merge these
 222 individuals into pairs before merging the pairs into a cluster of 4.

223 As with the heatmap of intercorrelation p-values in figure 2, the dendrogram shows the high redundancy
 224 between the RNA expression features. A significant proportion of the cumulative distance between consec-
 225 utive clusters is concentrated between the final 6 merges (i.e. between the top 7 clusters, compared to the
 226 total of 50 input features).

227 While the cumulative distance is not as direct a measure of information compression as calculating
 228 cumulative variance of features from the eigenvalues of produced in Principal Component Analysis, it is still
 229 a measure of total Ward distances within the cluster hierarchy, which equals the sum of minimum increases
 230 of within-cluster variance from the remaining sub-cluster merges.

231 This density of cumulative distance found within a small number of clusters therefore implies that there
 232 is significant room for compression within the RNA expression variables.

233 More completely, when applied to the entire COADREAD Colon and Rectal Cancer dataset, we find that
234 reducing the 20,500 RNA expression features to 4,100 agglomerative clusters (20% of the dataset’s original
235 size), retains 50.2% of the cumulative distance between consecutive clusters.

236 3.5 Bimodal Activations

237 Bimodal Activations are a widely used transformation in analysis of RNA expression.

238 The transform assumes that globally, across all RNA molecules and all individuals in a dataset, there is
239 a bimodal distribution of expression levels. This is interpreted as defining a cutoff above or below which an
240 RNA molecule is considered ‘active’ or ‘inactive’ respectively. Around the cutoff there is also often assumed
241 to be a grey area of uncertain or noisy expressions that are ignored by the transform, which are to be
242 excluded.

243 We attempted to expand our analysis of Bivariate Filtering to include Bimodal Activations of the RNA
244 expressions in the dataset. We use Point Bi-Serial Correlations to measure the strength of relationship
245 between RNA Bimodal Activations and Protein Expression and, as with the Spearman Rank Correlations,
246 examine the top 50 most correlated RNA molecules for illustration. Point Bi-Serial Correlation measures the
247 relationship between the state of a binary variable and the value of a continuous variable, and is a special
248 case of Pearson Correlation.

249 In order to define the cutoffs for bimodal activation quantitatively we train a Mixed Gaussian Model on the
250 distribution of RNA expressions across all RNA molecules and all individuals in the dataset. By asserting that
251 the Bimodal Activations are described by 3 Gaussian distributions (the distribution of ‘inactive’ expressions,
252 the distribution of ‘active’ individuals, and a distribution for noise between the peaks), we calculate cutoffs as
253 the crossover points of these distributions, with expressions above 6.35 counting as ‘active’ (28%), expressions
254 equal to or below 0.05 counting as ‘inactive’ (57%). Approximately 15% of the data falls between these cutoffs,
255 and is discarded as noise.

256 In addition to the proportion discarded as noise, many RNA features are found to be either ‘active’ or
257 ‘inactive’ across all individuals in the dataset, and therefore have no relationship with protein expression.
258 This results in only 2878 RNA molecules out of the initial 20,000 having a calculable expression correlation
259 with the ATM protein, reducing the available RNA features by 86% before Bivariate Filtering could be
260 applied. This loss of data significantly impacted the analysis, with only 1 out of the top 50 most significant
261 RNA molecules according to Spearman Rank Expression Correlations even having a calculable Point Bi-Serial
262 Correlation using Bimodal Activations.

263 To address this problem, we repeated the analysis without discarding noisy data and use a rule of thumb
264 cutoff expression of 6.00 to divide between the ‘active’ (41%) and ‘inactive’ (59%) expressions. This still
265 results in only 8389 RNA molecules out of the initial 20000 having a calculable expression correlation with
266 the ATM protein, and only 28 out of the top 50 most significant RNA molecules according to Spearman
267 Rank Expression Correlations even having a calculable Point Bi-Serial Correlation.

268 Reducing a continuous variable (RNA expression) to a binary one (RNA Bimodal Activation) clearly
269 eliminates a large amount of information from the dataset, and such transformations should be used with
270 care. In a dataset of this small size, and in a problem where we are searching for relationships amongst a large
271 number of features, this information loss is unmanageable.

272 3.6 Towards Causal Analysis

273 As mentioned in § 2, analysis based upon causal counterfactual reasoning requires a large and sufficiently
274 varied dataset. The COADREAD dataset used in this preliminary analysis does not meet these criteria. A
275 natural next step would be to combine several datasets from the TCGA for different types of cancer in order
276 to provide a larger and more balanced dataset.

277 The analysis performed in § 3.4 identified clusters of RNA molecules which allows us to move forward
278 using causal discovery methods, as referenced in § 2, for this reduced dimensionality problem. We will
279 decompose the problem in two ways: (1) independent data clusters, and (2) a coarse-to-fine approach.
280 The first method decomposes the model into loosely related components, see [26] for a formal definition of
281 decomposable models.

282 For the coarse-to-fine approach, we can learn a high-level causal model where each node represents a
283 cluster. The definition of each cluster node can be a single representative from the cluster, the centroid, or
284 a linear/non-linear combination of all elements, which can be found using classical dimensionality reduction
285 techniques. Once the structure between clusters has been learned, we can zoom in on the local structure
286 within clusters. The dependencies between clusters can also be further refined.

287 4 Conclusions and Future Work

288 Our findings as described in this paper are aligned with the current understanding of the role of different
289 RNA expressions, thus supporting our claim that a screening based on the blood test results is meaningful.
290 We have also shown that the analysis of bimodal activations is not aligned with the current understanding
291 of the roles and connections between RNA expressions, due to the small size of the dataset after a sizeable
292 subset of records is discarded (those where we cannot deduce with any certainty whether the RNA expression
293 is on or off). We deduce that the bimodal activation analysis is applicable only to sufficiently large datasets.

294 Our preliminary analysis is based on a relatively small publicly available dataset, hence the findings are
295 limited.

296 In future work, we plan to perform more complex types of causal analysis and inference on larger datasets,
297 as we predict that they have a potential to uncover previously unknown or not well-understood connections
298 between RNA and occurrence of cancer.

299 Acknowledgements

300 The results published here are in whole or part based upon data generated by the TCGA Research Network:
301 <https://www.cancer.gov/tcga>.

302 References

- 303 [1] Rebecca L. Siegel, Kimberly D. Miller, Ann Goding Sauer, Stacey A. Fedewa, Lynn F. Butterly,
304 Joseph C. Anderson, Andrea Cercek, Robert A. Smith, and Ahmedin Jemal. Colorectal cancer statistics.
305 *CA: A Cancer Journal for Clinicians*, 70:145–164, 2020.
- 306 [2] Wolff W.I. Colonoscopy: history and development. *American Journal of Gastroenterology*, 84, 1989.
- 307 [3] Petrini J.L. Colonoscopy surveillance intervals: we are not there yet. *Gastrointestinal Endoscopy*,
308 85:1271–1272, 2017.
- 309 [4] Singh S., Singh P.P., Murad M.H., Singh H., and Samadder N.J. Prevalence, risk factors, and outcomes
310 of interval colorectal cancers: a systematic review and meta-analysis. *American Journal of Gastroen-*
311 *terology*, 109:1375, 2014.
- 312 [5] Naylor J., Saltzman J.R., Campbell E.J., Perencevich M.L., Jajoo K., and Richter J.M. Impact of physi-
313 cian compliance with colonoscopy surveillance guidelines on interval colorectal cancer. *Gastrointestinal*
314 *Endoscopy*, 85:1263–1270, 2017.
- 315 [6] Kamal Mahtani, Elizabeth A. Spencer, Jon Brassey, and Carl Heneghan. Catalogue of bias: observer
316 bias. *BMJ Evidence-Based Medicine*, 23:23–24, 2018.
- 317 [7] Mary J Goldman, Brian Craft, Mim Hastie, Kristupas Repečka, Fran McDade, Akhil Kamath, Ayan
318 Banerjee, Yunhai Luo, Dave Rogers, Angela N Brooks, et al. Visualizing and interpreting cancer
319 genomics data via the xena platform. *Nature Biotechnology*, 38(6):675–678, 2020.
- 320 [8] Miroslav Kubat and Stan Matwin. Addressing the curse of imbalanced training sets: One-sided selection.
321 In *Proceedings of the Fourteenth International Conference on Machine Learning (ICML)*, pages 179–186,
322 1997.

- 323 [9] Samuel Budd, Arno Blaas, Adrienne Hoarfrost, Kia Khezeli, Krittika D’Silva, Frank Soboczenski, Gra-
324 ham Mackintosh, Nicholas Chia, and John Kalantari. Prototyping crisp: A causal relation and inference
325 search platform applied to colorectal cancer data. In *2021 IEEE 3rd Global Conference on Life Sciences
326 and Technologies (LifeTech)*, pages 517–521, 2021.
- 327 [10] John Kalantari, Heidi Nelson, and Nicholas Chia. The unreasonable effectiveness of inverse reinforcement
328 learning in advancing cancer research. *Proceedings of the AAAI Conference on Artificial Intelligence*,
329 34(01):437–445, 2020.
- 330 [11] Saman Farahmand, Corey O’Connor, Jill A Macoska, and Kouros Zarringhalam. Causal Inference
331 Engine: a platform for directional gene set enrichment analysis and inference of active transcriptional
332 regulators. *Nucleic Acids Research*, 47(22):11563–11573, 2019.
- 333 [12] Judea Pearl. *Causality*. Cambridge University Press, 2 edition, 2009.
- 334 [13] Clark Glymour, Kun Zhang, and Peter Spirtes. Review of causal discovery methods based on graphical
335 models. *Frontiers in Genetics*, 10:524, 2019.
- 336 [14] David Maxwell Chickering, David Heckerman, and Christopher Meek. Large-sample learning of bayesian
337 networks is NP-hard. *Journal of Machine Learning Research*, 5:1287—1330, 2004.
- 338 [15] Peter Spirtes, Clark N Glymour, Richard Scheines, and David Heckerman. *Causation, Prediction, and
339 Search*. MIT press, 2000.
- 340 [16] Peter Spirtes, Christopher Meek, and Thomas Richardson. An algorithm for causal inference in the
341 presence of latent variables and selection bias. *Computation, causation, and discovery*, 21:211–252,
342 1999.
- 343 [17] Jiji Zhang. On the completeness of orientation rules for causal discovery in the presence of latent
344 confounders and selection bias. *Artificial Intelligence*, 172(16-17):1873–1896, 2008.
- 345 [18] David Maxwell Chickering. Optimal structure identification with greedy search. *Journal of Machine
346 Learning Research*, 3:507–554, 2002.
- 347 [19] Ni Y Lu, Kun Zhang, and Changhe Yuan. Improving causal discovery by optimal bayesian network
348 learning. *Proceedings of the AAAI Conference on Artificial Intelligence*, 35(10):8741–8748, May 2021.
- 349 [20] Shohei Shimizu, Patrik O. Hoyer, Aapo Hyvärinen, and Antti Kerminen. A linear non-gaussian acyclic
350 model for causal discovery. *Journal of Machine Learning Research*, 7(72):2003–2030, 2006.
- 351 [21] Xun Zheng, Bryon Aragam, Pradeep K Ravikumar, and Eric P Xing. Dags with no tears: Continuous
352 optimization for structure learning. In *Advances in Neural Information Processing Systems*, volume 31,
353 2018.
- 354 [22] Shohei Shimizu, Takanori Inazumi, Yasuhiro Sogawa, Aapo Hyvärinen, Yoshinobu Kawahara, Takashi
355 Washio, Patrik O Hoyer, and Kenneth Bollen. Directlingam: A direct method for learning a linear
356 non-gaussian structural equation model. *Journal of Machine Learning Research*, 12:1225–1248, 2011.
- 357 [23] Charles Spearman. The proof and measurement of association between two things. *American Journal
358 of Psychology*, 15:72–101, 1904.
- 359 [24] Karl Pearson. Notes on regression and inheritance in the case of two parents. *Proceedings of the Royal
360 Society of London*, 58:240–242, 1895.
- 361 [25] Mirco Nanni. Speeding-up hierarchical agglomerative clustering in presence of expensive metrics. In
362 *Advances in Knowledge Discovery and Data Mining, 9th Pacific-Asia Conference, PAKDD*, volume
363 3518 of *Lecture Notes in Computer Science*, pages 378–387. Springer, 2005.
- 364 [26] Dalal Alrajeh, Hana Chockler, and Joseph Y. Halpern. Combining experts’ causal judgments. *Artif.
365 Intell.*, 288:103355, 2020.