

Model-free estimation of higher-order interactions in single cell expression data

L Del Debbio², A Jansma^{1,2}, A Khamseh^{1,2,3}, C P Ponting¹

¹MRC Human Genetics Unit, Institute of Genetics & Cancer, University of Edinburgh, Edinburgh EH4 2XU, United Kingdom.

²Higgs Centre for Theoretical Physics, School of Physics & Astronomy, University of Edinburgh, Edinburgh EH9 3FD, United Kingdom.

³School of Informatics, University of Edinburgh, Edinburgh EH8 9AB, United Kingdom



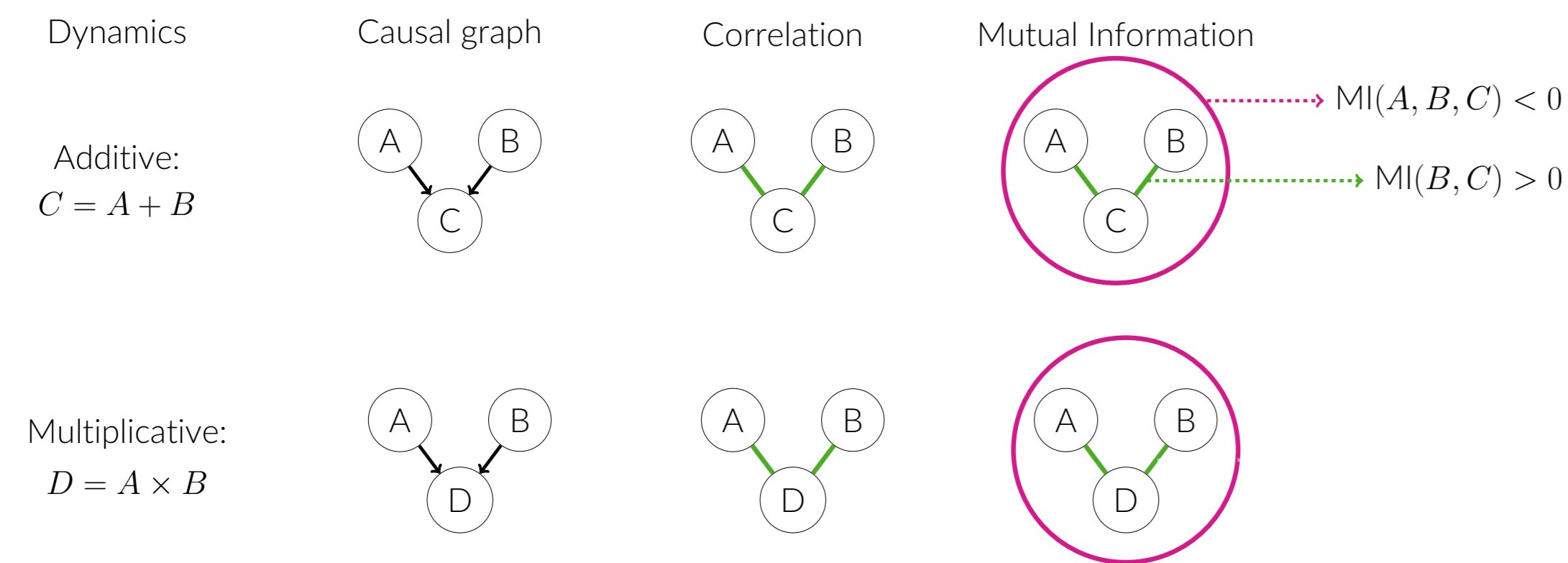
1 Motivation

Causal networks of molecular cell biology can guide our understanding of biological processes, but identifying causal effects from observational data is non-trivial. Existing methods suffer from model-misspecification, ignoring higher-order interactions and dependent variables.

We propose a new way to reconstruct genetic networks that suffers from none of these issues, and generalises to interactions that involve more than 2 genes, potentially revealing higher-order, combinatorial gene regulation. **In four different cell types, we find hundreds of 3-point interactions.**

2 Traditional methods hide dynamics

Consider two binary transcription factors A and B separately affecting a target gene C , and a gene D that is only affected by a bound complex of A and B . Estimating the causal graph, correlation structure, and mutual information gives the networks below:



That is, none of these methods are able to distinguish between these biologically different scenarios - they hide the underlying dynamics.

3 Definition: Model-free interactions

The **effect** I_i of a gene $X_i \in X$ on an outcome Y is the extent to which Y changes when the expression of X_i changes, all other genes (\underline{X}) being fixed:

$$I_i = \frac{\partial Y}{\partial X_i} \Big|_{\underline{X}=0}$$

We say that two genes X_i and X_j **interact** when the level of X_j changes the effect of X_i on Y :

$$I_{ij} = \frac{\partial I_i}{\partial X_j} \Big|_{\underline{X}=0} = \frac{\partial^2 Y}{\partial X_j \partial X_i} \Big|_{\underline{X}=0}$$

Similarly, a third gene X_k might modulate this interaction, which results in a **3-point interaction**:

$$I_{ijk} = \frac{\partial I_{ij}}{\partial X_k} \Big|_{\underline{X}=0} = \frac{\partial^3 Y}{\partial X_k \partial X_j \partial X_i} \Big|_{\underline{X}=0}$$

If we take the most general outcome Y possible - the (log of) the joint distribution of all genes X - we get the definition of interaction from [1]:

Definition 1 (Pairwise interaction between binary genes)

A pair of binary genes $\{X_i, X_j\} \in X$ has a pairwise interaction I_{ij} where

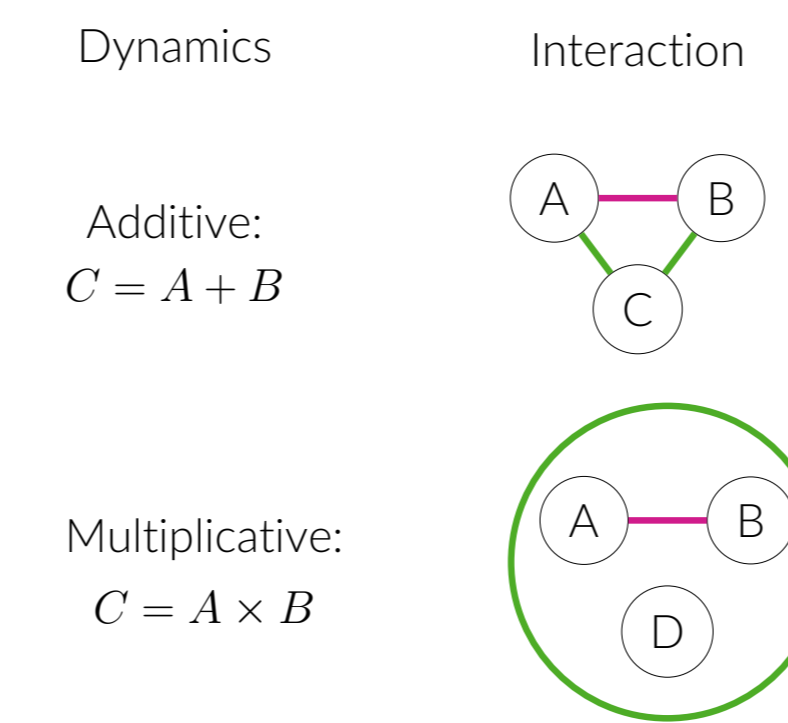
$$I_{ij} = \log \frac{p(X_i = 1, X_j = 1 | \underline{X} = 0) p(X_i = 0, X_j = 0 | \underline{X} = 0)}{p(X_i = 1, X_j = 0 | \underline{X} = 0) p(X_i = 0, X_j = 1 | \underline{X} = 0)}$$

4 Properties of model-free interactions

Definition 1 has the following properties:

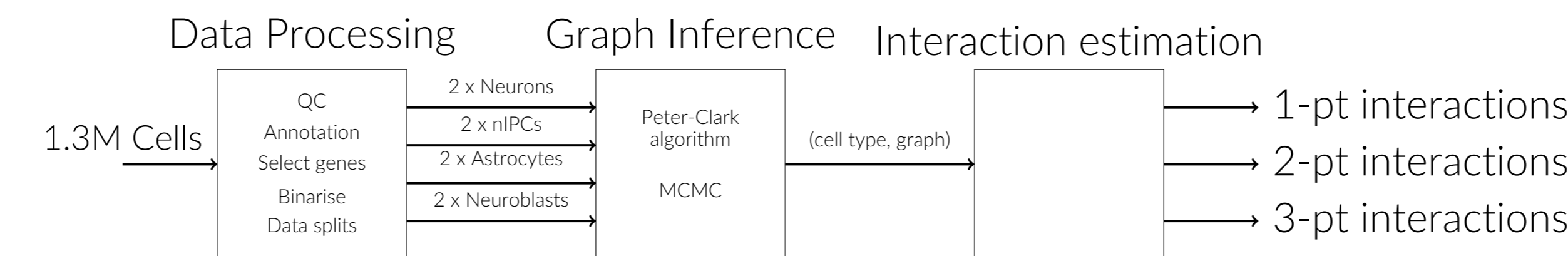
- It is symmetric: $I_{ij} = I_{ji}$.
- Conditionally independent genes do not interact: $X_i \perp\!\!\!\perp X_j | \underline{X} \implies I_{ij} = 0$.
- If $\underline{X} = \emptyset$, the interaction reduces to a log-odds ratio.
- It is model-independent and can be directly estimated from expression data.
- It can be naturally extended to n -point interactions by taking n 'th derivatives of $\log p(X)$.

On the right we show the interactions in the systems from section 2. We can now distinguish the two systems - **D being affected by a bound complex of A and B lead to a 3-point interaction.**

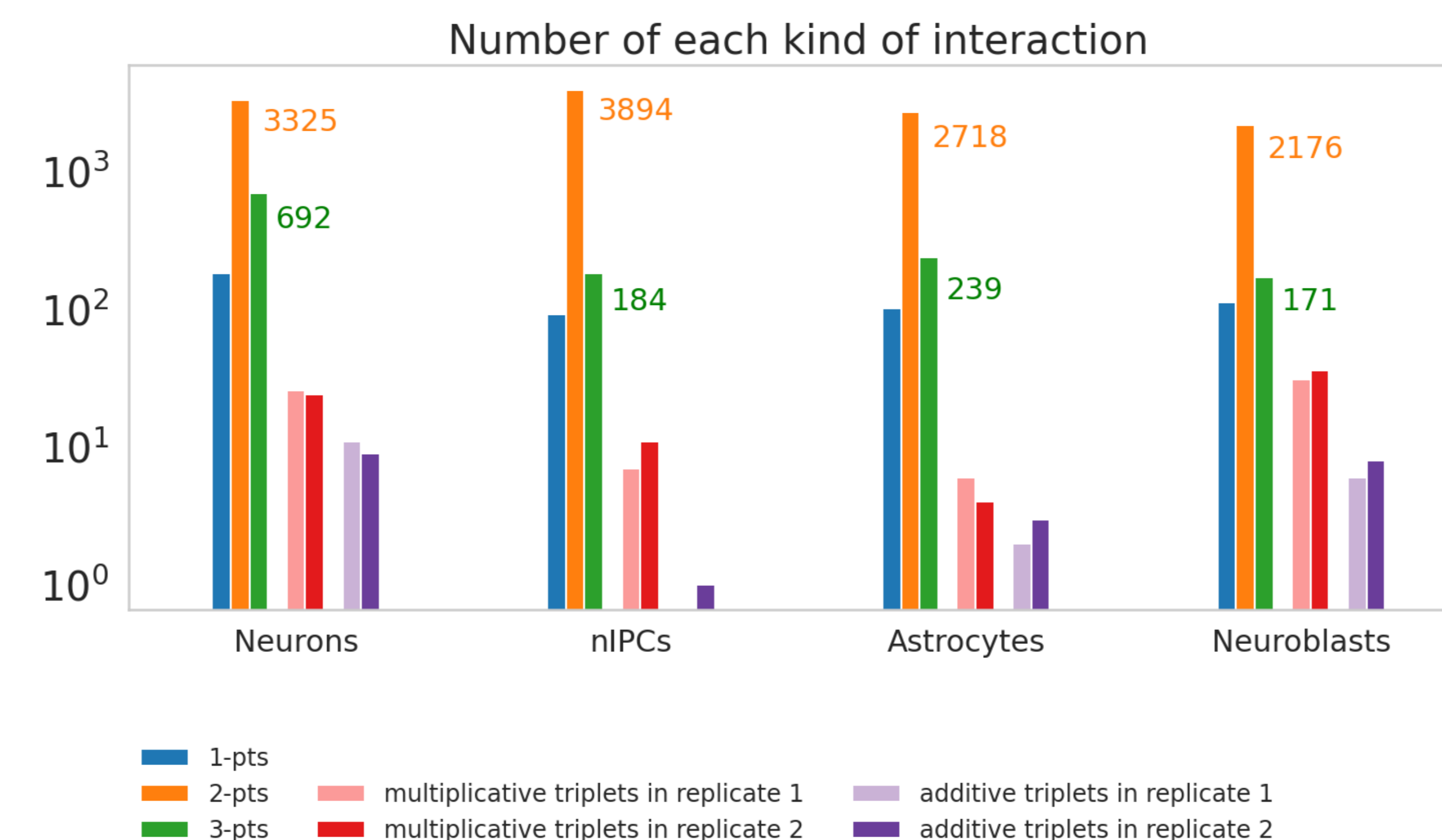


5 Mouse brain cells contain hundreds of 3-point interactions

We consider interactions in a data set of 1.3M embryonic (E18.5) mouse brain cells [2]. From each of four cell types, we construct two biological replicates of 10k cells. To specifically focus on collider triplets, we also estimate the causal graph in each cell type.



Below we show the number of 1-, 2-, and 3-point interactions that are significant in both replicates. Also shown is the number of triplets with an additive or multiplicative interaction pattern in each replicate separately.



1-point interactions predict housekeeping genes better than expression

Given a reference set S of organism-wide mouse housekeeping genes [3], we can compare the enrichment in S of the set E_N of top N expressed genes and the set C_N of top N 1-point interacting genes (self-coupling). Figure 1 shows that this is positive for all N , so 1-point interactions are a better predictor of housekeeping genes than expression.

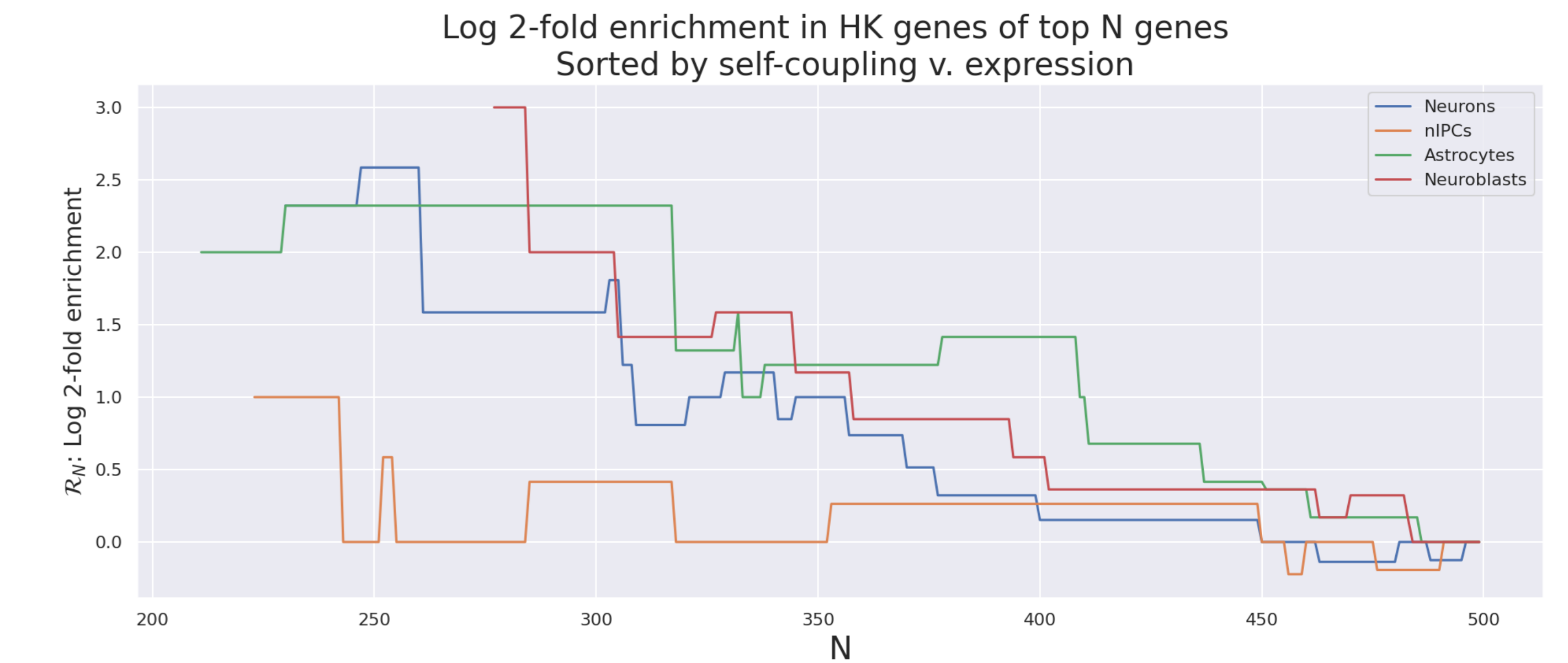


Figure 1: Relative enrichment in housekeeping: $\mathcal{R}_N = \log \frac{|C_N \cap S|}{|E_N \cap S|}$.

6 Conclusion

- Model-free interactions can reveal underlying dynamics in observational data where other methods cannot.
- We claim the existence of these replicated 3-point interactions, and propose that they reflect higher-order gene regulation.**
- The interactions contain biological information.
- We can extend a famous phrase:

Correlation **is not** Causation

is not

Interaction

References

[1] S. V. Beentjes and A. Khamseh, *Higher-order interactions in statistical physics and machine learning: A model-independent solution to the inverse problem at equilibrium*, Physical Review E, 102 (2020), p. 053314.

[2] X. Genomics, *Million cell dataset*, 2017.

[3] B. W. Hounkpe, F. Chenou, F. de Lima, and E. V. De Paula, *HRT Atlas v1.0 database*, Nucleic Acids Research, 49 (2020), pp. D947--D955.