

# A graph feature autoencoder for predicting perturbations to steady state gene expression

C.W. Seitz

June 30, 2022

## 1 Revealing causal relations and predicting gene expression

Suppose that we wanted to predict the response of gene expression to some external perturbation, like an osmotic shock and treatment with a cytokine. In principle, it is *not* possible to infer that response unless a detailed list of biochemical reactions is known apriori - a list that is very rare and likely doesn't exist at all. Luckily, there is another pharmacological problem which can potentially be addressed using this formalism. If the drug of interest is known, and can potentially induce gene expression through unknown mechanisms, those mechanisms can be identified given the right dataset. We do not need to measure detailed interactions between variables, for example, the precise effect of histone modifications and chromatin structure on the rate of transcription. Instead, these interactions can be inferred from the data as long as those variables can be measured with reasonable precision. Ultimately, this framework cannot be used for screening a large panel of potential drugs but it is appropriate for learning mechanisms of say adaptive resistance to drug treatment. We can begin to answer questions like: what are the responsible mechanisms for drug resistance in single cells? In any case, this will require some effort for the learning methods to become interpretable in a biologically significant way. The mechanisms will have more of a qualitative description as well, unless a deep network were to be supplemented by more biophysical characterization. This could be a valuable tool in trying to understand which genetic mechanisms are responsible for the development of resistant states. Deep learning is a key method for integrating information provided by sequencing tools that are becoming readily available.

Another potential application (which is arguably less valuable/feasible) is *prediction* - the ability to predict the effects of the fold-change of one gene on its neighbors in the gene regulatory network. This requires some kind of domain knowledge for selecting the gene of interest and genes affected downstream. It is less feasible because drugs often don't act on genes directly but instead target certain proteins, presenting many unknowns between disruption

of protein-mediated signaling and transcription.