

# Combined effects of protein expression variance and correlation on multicomponent systems

Kyle M. Kovary<sup>1</sup> Mary N. Teruel<sup>1</sup>

<sup>1</sup>Department of Chemical and Systems Biology, Stanford University, Stanford, CA 94305, USA.

## Abstract

Protein expression variation leads to phenotypic variance in otherwise identical cells. This has been demonstrated in cell signaling and differentiation decisions. Additionally coordinated expression of proteins between cells can tune signaling pathways either towards a more binary or analog (correlability) of protein levels. There is some evidence in bulk cell measurements and in bacteria that certain heteromeric subunits or metabolic pathways may be expressed in a coordinated fashion (i.e. operons in bacteria), there has not been a direct measurement of coordinated expression of proteins (independent of TFs) in metazoans (vertebrates?). Here, we measure cell-to-cell variability of relative protein abundance using quantitative proteomics of individual *Xenopus laevis* eggs and show that proteins involved in metabolic pathways or members of heteromeric complexes tend to have high correlations with other members of those pathways/complexes. Our previous work highlighted the fact that correlated expression increases the total variation of a pathway, so one would reason that certain pathways or complexes would need to compensate of this extra source of by reducing the variation of expression of these pathways or complexes. To test this we computed total variance score that took into account both the coefficient of variance and correlations between proteins in a pathway and found that the lower 10% of GO terms were highly enriched for metabolic pathways. When we looked at the relationship between CV and R<sup>2</sup> between proteins in GO terms we found a negative relationship between them, demonstrating that increased correlation needs to come at the expense of decreased variance. Simple molecular models of heteromeric complexes and metabolic pathways demonstrate that this tradeoff can result in higher efficiencies in both function and reduced energy waste. Together, our study argues for a control principle whereby coordinated expression of proteins in a pathway can require lower variance in order to regulate large single cells *Xenopus laevis* eggs. An advantage of this approach is that

**Keywords:** single cell, proteomics, stochasticity

## Introduction

The coordinated expression of proteins is vital across cellular function, from maintaining a dynamic steady state to differentiation of cells into specialized types in different tissues. The ability to regulate coordinated expression has been demonstrated to occur via transcription (transcription factors, chromatin regulation), translation (specialized ribosomes, mRNA structure/modifications), and degradation (E3 ligases). These processes, as well as the molecules they target, are all subject to noise, which has been shown

to lead to differences in cell behavior and decisions in otherwise identical cells. Variability and coordinated expression (correlability) of protein levels leads to increased population level control in binary decisions, and a decreased ability to execute analog signaling. Despite these important insights, a global understanding of the variability of protein expression between single cells and the coordinated expression of proteins is largely an open question. To systematically assess these properties of single cells, we carried out a proteomics experiment in order to regulate large single cells *Xenopus laevis* eggs. An advantage of this approach is that we can get around much of the signal to noise issues that accompany studying single cells due to low sample amounts. Additionally, at this stage of

development transcription is restricted so we are able to gain insights into non-transcriptional control of protein expression. We also utilized isobaric tagging in order to measure 25 individual eggs in 5 mass spectrometry runs. In this study, we were able to measure the relative abundance  $>1000$  proteins across single cells in order to better understand the properties of stochastic and coordinated expression of proteins.

With this dataset we have been able to, for the first time, measured the relationship between protein expression variance and coordinated protein expression on a proteome scale in single cells. We have observed that certain classes of proteins, including protein complexes and metabolic pathways, are expressed in such a way that increased coordinated expression is balanced by decreased variation. By doing this, cells are able to decrease the total variance of a given complex or pathway, an elegant balancing act that allows for finer control of metabolic throughput and controlling the number of potentially formed complexes through stoichiometric control. Though this kind of coordinated expression leads to an increase in variance in a population of cells, it can reduce variation within a cell.

To fetch bibliographic metadata automatically from the web. For example, citing a paper can be as easy as providing its DOI (Clark and Gelfand,

2006) or even just a few keywords (Ahrends et al., 2014; Shi et al., 2017).

## Results

### **Single cell proteomics reveals global protein expression variability and coordinated expression between protein pairs.**

Activated *Xenopus laevis* eggs were collected at 5 time points across the first cell cycle (0, 20, 40, 60, and 80 minutes), with 5 eggs at each time point. Using TMT multiplexing and mass spectrometry, we were able to determine the relative abundance of more than 1300 proteins. Expression of these proteins were largely invariant across the cell cycle, revealing that these highly expressed genes are likely not regulated by cell cycle processes. Additionally, a PCA analysis of these eggs showed no discernible clustering on cell cycle time (Fig S1).

To determine the variance across these proteins, we calculated the coefficient of variation (CV) using all time points collected. This showed a wide range of variation containing multiple distributions (Fig 1B), many of which are consistent with our previous study of variation using targeted mass spectrometry (Fig S2). In order to see if variation was a regulated process, we grouped the protein variation by gene ontology terms (GO terms) and plotted them in ranked order. We saw that in

general, processes xxx were...

Since we were able to measure all of these proteins within single cells, we are able to calculate coordinated expression of protein pairs at single cell resolution. Using Pearson correlation coefficient, we could determine coordinated expression of nearly 2 million protein pairs (Fig 1C) The distribution of correlation coefficients fit a normal distribution, with the majority of protein pairs appearing to not show significant coregulation. However, there appear to be a significant number of protein pairs containing high correlation coefficients, and the heatmap showed a lot of clustering of these highly coregulated pairs (Fig 1D).

### Protein modules display different levels of regulation of noise and co-expression of constituents

To investigate whether coordinated expression goes beyond just protein pairs, we wanted to see if pathways, complexes, and other modules had proteins that with coordinated expression. To answer this question, we grouped proteins by GO terms and calculated the median correlation coefficient across all pairs of proteins within the term. When the ranked media correlation is plotted it's clear that there are a significant number of GO terms with higher than expected levels of coordinate expression.

### Total variance metric reveals certain modules enforce a tradeoff between protein co-expression and noise

#### Variance Sum Law — Dependent Case

The sum of variance for a group of proteins  $X_1, \dots, X_n$  with dependent relationships:

$$var\left(\sum_{i=1}^n X_i\right) = \sum_{i=1}^n var(X_i) + 2 \sum_{i,j:i < j} cov(X_i, X_j)$$

$$\sigma_{total} = \sqrt{var\left(\sum_{i=1}^n X_i\right)}$$

$$\sigma_{norm} = \frac{\sigma_{total}}{n}$$

When calculating the total variance of a pathway or general group of proteins, the total variance is the square root of the sums of the squares of the coefficient of variance. However, if any of the proteins are correlated in expression, an extra term must be added to account for this correlation. This means that if a group of proteins have a positive correlation, then the total variance will be higher than if there was no correlation. In order to study how variance and coordinated expression are related between single cells, we wanted to calculate a normalized total variance metric. Since we were

not interested in the increase of variance due to the number of proteins in a module, we calculated the total variance using the coefficient of variation and correlation coefficients and normalize to the total number of proteins in the module. When total variance is plotted against the number of proteins, there is a clear linear relationship between the two. Once we normalize to the total number of proteins this relationship is lost and we are left with the effects of coefficient of variation and correlation.

### **Low total variance modules are enriched for metabolic pathways and heteromeric protein complexes**

### **Enforcement of a tradeoff between protein co-expression and noise increases efficiency of metabolic pathways and heteromeric protein complexes**

## **Methods**

### **Collection and activation of *Xenopus laevis* eggs:**

*Xenopus* egg extracts were prepared based on modifications of a previous protocol (Tsai et al, 2014). All of the animal protocols used in this manuscript were approved by the Stanford University Administrative Panel on Laboratory Animal Care. To induce egg laying, female *Xenopus laevis* were injected with human chorionic gonadotropin injection

the night before each experiment. To collect the eggs, the frogs were subjected to pelvic massage, and the eggs were collected in 1X Marc's Modified Ringer's (MMR) buffer (0.1 M NaCl, 2 mM KCl, 1 mM MgCl<sub>2</sub>, 2 mM CaCl<sub>2</sub>, 5 mM HEPES, pH 7.8). To remove the jelly coat from the eggs, they were placed in a solution of 2% cysteine in 1× MMR buffer for 4 min and gently agitated, after which they were washed four times with 1× MMR buffer. To activate the cell cycle, eggs were placed in a solution of 0.5 μg/ml of calcium ionophore A23187 (Sigma) and 1X MMR buffer for 3 min, after which they were washed four times with 1× MMR buffer. Single eggs were collected at their respective timepoints and placed into 600uL tubes and snap frozen in liquid nitrogen before being stored at 80°C.

### **Sample preparation for mass spectrometry:**

Single eggs were lysed mechanically by pipetting the egg in 100μL of lysis buffer (100 mM NaCl, 25 mM Tris pH 8.2, Complete EDTA- free protease inhibitor cocktail (Sigma). The lysate was then placed in a 400 uL natural polyethylene microcentrifuge tube (E&K Scientific #485050) and spun at 15,000 g in a right angle centrifuge (Beckman Microfuge E) at 4°C for 5 min. The lipid layer was removed by using a razor blade to cut the tube off just beneath it, and the cytoplasmic

fraction was pipetted into a 1.5-ml protein LoBind tube (Fisher Scientific #13-698-794), being careful to leave the yolk behind. To precipitate the proteins from the cytoplasmic fraction, 1 ml of ice cold acetone was added to each sample and placed at 20°C overnight. To collect precipitated proteins, the samples were centrifuged at 18,000 g for 20 min at 4°C. Acetone was decanted, and the protein pellets were resolubilized in 25 $\mu$ L of 8 M urea. To fully solubilize the protein pellet, the samples were placed in a shaker for 1 h at room temperature. The samples were then diluted to 2 M urea with 50 mM ammonium bicarbonate to a 100 $\mu$ L volume, after which protein concentration was measured in duplicate with a BCA assay by taking two 10 $\mu$ L aliquots of each sample. The proteins in the remaining 80 $\mu$ L of sample volume were reduced with 10 mM TCEP and incubated for 30 min at 37°C, then alkylated with 15 mM iodoacetamide and incubated in the dark at room temperature.

Next, the samples were diluted to 1 M urea with 50 mM ammonium bicarbonate. Trypsin (Promega #V5113) was then added at a ratio of 10 ng trypsin per 1 $\mu$ g protein (no < 500 ng was added to a sample). The trypsin digestion was carried out at 37°C for 12–16 h. To stop the trypsin, formic acid (Fisher #A117-50) was added at a ratio of 3 $\mu$ L per 100 $\mu$ L of sample to bring the pH down to < 3.

Peptides were cleaned up using an Oasis HLB uELUTION plate (Waters), equilibrated, and washed with 0.04% trifluoroacetic acid in water, and eluted in 80% acetonitrile with 0.2% formic acid. All solutions used are HPLC grade. Samples were then lyophilized. To remove any variance produced by phosphorylated peptides, the samples were phosphatase-treated. Peptides were resolubilized in 50 $\mu$ L of 1X NEBuffer 3 (no BSA), and calf intestinal alkaline phosphatase (NEB #M0290S) was added at a ratio 0.25 units per  $\mu$ g of peptide and incubated for 1 h at 37°C. The peptides were cleaned up again according to steps described above. Peptides were resolubilized in 2% acetonitrile and 0.1% formic acid before MS analysis.

### **Mass spectrometry data collection:**

Ask SUMS for details

### **Mass spectrometry data analysis:**

Ask SUMS for details

### **Data processing:**

To minimize the effects of non-biological variance, a correction factor was used to correct for these biases. First, each peptide was normalized by the median across all of the samples. Second, the vector of all peptides for each cell was divided by the median value across all peptides. Before peptides were used to estimate protein abundance, highly

variable peptides needed to be filtered out. Peptides were grouped by their master UniProt accession number, and all peptides whose log transformed values were more than  $2\sigma$  away from the mean were discarded. Next, protein abundances were estimated by taking the median normalized value for each peptide for a protein in each sample. Missing protein levels were imputed using the k-nearest neighbors algorithm, with k being set to 1 and the similarity measure for distance the Gower's distance between the proteome vectors.

## References

Ahrends, R., Ota, A., Kovary, K.M., Kudo, T., Park, B.O., and Teruel, M.N. (2014). Controlling low rates of cell differentiation through noise and ultrahigh feedback. *Science* *344*, 1384–1389.

Clark, J.S., and Gelfand, A.E. (2006). A future for models and data in environmental science. *Trends in Ecology & Evolution* *21*, 375–380.

Shi, Z., Fujii, K., Kovary, K.M., Genuth, N.R., Röst, H.L., Teruel, M.N., and Barna, M. (2017). Heterogeneous ribosomes preferentially translate distinct subpools of mRNAs genome-wide. *Molecular Cell* *67*, 71–83.e7.

R version 3.6.1 (2019-07-05)

Platform: x86\_64-apple-darwin15.6.0 (64-bit)

Running under: macOS Catalina 10.15.3

Matrix products: default

BLAS: /Library/Frameworks/R.framework/Versions/3.6/Resources/lib/libRblas.0.dylib

LAPACK: /Library/Frameworks/R.framework/Versions/3.6/Resources/lib/libRlapack.dylib

locale:

[1] en\_US.UTF-8/en\_US.UTF-8/en\_US.UTF-8/C/en\_US.UTF-8/en\_US.UTF-8

attached base packages:

[1] stats graphics grDevices utils datasets methods base

other attached packages:

[1] knitr\_1.28

loaded via a namespace (and not attached):

[1] Rcpp_1.0.3	lubridate_1.7.4	digest_0.6.24	plyr_1.8.5
[5] R6_2.4.1	jsonlite_1.6.1	formatR_1.7	magrittr_1.5
[9] evaluate_0.14	bibtex_0.4.2.2	httr_1.4.1	rlang_0.4.4
[13] stringi_1.4.6	curl_4.3	xml2_1.2.2	rmarkdown_2.1
[17] tools_3.6.1	stringr_1.4.0	RefManager_1.2.12	xfun_0.12
[21] yaml_2.2.1	compiler_3.6.1	htmltools_0.4.0	