# Bioinformatics 3 V8 – Gene Regulation

Mon, Nov 17, 2014

#### Rates of mRNA transcription and protein translation

#### ARTICLE

doi:10.1038/nature10098

#### Global quantification of mammalian gene expression control

Björn Schwanhäusser<sup>1</sup>, Dorothea Busse<sup>1</sup>, Na Li<sup>1</sup>, Gunnar Dittmar<sup>1</sup>, Johannes Schuchhardt<sup>2</sup>, Jana Wolf<sup>1</sup>, Wei Chen<sup>1</sup> & Matthias Selbach<sup>1</sup>

SILAC: "stable isotope labelling by amino acids in cell culture" means that cells are cultivated in a medium containing heavy stable-isotope versions of essential amino acids.

When non-labelled (i.e. light) cells are transferred to heavy SILAC growth medium, newly synthesized proteins incorporate the heavy label while preexisting proteins remain in the light form.

Schwanhäuser et al. Nature 473, 337 (2011)



Parallel quantification of mRNA and protein turnover and levels. Mouse fibroblasts were pulse-labelled with heavy amino acids (SILAC, left) and the nucleoside 4-thiouridine (4sU, right). Protein and mRNA turnover is quantified by mass spectrometry and next-generation sequencing, respectively.

#### Rates of mRNA transcription and protein translation

84,676 peptide sequences were identified by MS and assigned to 6,445 unique proteins.

5,279 of these proteins were quantified by at least three heavy to light (H/L) peptide ratios

Mass spectra of peptides for two proteins.

Top: high-turnover protein Bottom: low-turnover protein.

Over time, the heavy to light (H/L) ratios increase.

You should understand these spectra!



3

Extract ratio *r* of protein with heavy amino acids ( $P_H$ ) and light amino acids ( $P_L$ ):  $r = \frac{P_H}{P_c}$ 

Assume that proteins labelled with light amino acids decay exponentially with degradation rate constant  $k_{dp}$ :  $P_{i} = P_{n}e^{-k_{dp}t}$ .

Express ( $P_H$ ) as difference between total number of a specific protein  $P_{total}$  and  $P_L$ :

$$P_{H}(t) = P_{total}(t) - P_{L}(t)$$

Assume that  $P_{total}$  doubles during duration of one cell cycle (which lasts  $t_{\infty}$ ):

$$\begin{aligned} P_{H}(t) &= P_{total}(t) - P_{L}(t) = P_{0} 2^{t/t_{cc}} - P_{L}(t), \\ r &= \frac{P_{H}}{P_{L}} = \frac{P_{0}}{P_{L}} 2^{\frac{t}{t_{cc}}} - 1 \\ \frac{P_{H}}{P_{L}} + 1 &= \frac{P_{0}}{P_{L}} 2^{\frac{t}{t_{cc}}} \end{aligned}$$

#### Protein half-lifes and decay rates



Consider *m* intermediate time points:

$$k_{dp} = \frac{\sum_{i=1}^{m} \log_e (r_{t_i} + 1)t_i}{\sum_{i=1}^{m} t_i^2} - \frac{\log_e 2}{t_{cc}},$$
  
$$T_{1/2} = \frac{\log_e 2}{k_{dp}}.$$
 Since then

$$P_L = P_0 e^{-k_{dp}t} = P_0 e^{-k_{dp}\frac{\log e^2}{k_{dp}}} = P_0 e^{\log e^{\frac{1}{2}}} = \frac{1}{2}P_0$$

 $\ln (ratio + 1) = \ln \frac{P_0}{P_L} 2^{\frac{t}{t_{cc}}} = \ln e^{k_{dp}t} + \ln 2^{\frac{t}{t_{cc}}} = k_{dp}t + \ln 2^{\frac{t}{t_{cc}}}$  The same is done to compute mRNA half-lives (not shown).

#### mRNA and protein levels and half-lives



a, b, Histograms of mRNA (blue) and protein (red) half-lives (a) and levels (b).

Proteins were on average 5 times more stable (9h vs. 46h) and 900 times more abundant than mRNAs and spanned a higher dynamic range.

c, d, Although mRNA and protein levels correlated significantly, correlation of halflives was virtually absent

#### Mathematical model of transcription and translation

A widely used minimal description of the dynamics of transcription and translation includes the synthesis and degradation of mRNA and protein, respectively





The mRNA (*R*) is synthesized with a constant rate  $v_{sr}$  and degraded proportional to their numbers with rate constant  $k_{dr}$ .

The protein level (*P*) depends on the number of mRNAs, which are translated with rate constant  $k_{sp}$ .

Protein degradation is characterized by the rate constant  $k_{dp}$ .

The synthesis rates of mRNA and protein are calculated from their measured half lives and levels.

#### Computed transcription and translation rates

Average cellular transcription rates predicted by the model span two orders of magnitude.

The median is about 2 mRNA molecules per hour (b). An extreme example is Mdm2 with more than 500 mRNAs per hour







Calculated translation rate constants are not uniform

#### Maximal translation constant

Abundant proteins are translated about 100 times more efficiently than those of low abundance

Translation rate constants of abundant proteins saturate between approximately 120 and 240 proteins per mRNA per hour.

The maximal translation rate constant in mammals is not known.

The estimated maximal translation rate constant in sea urchin embryos is 140 copies per mRNA per hour, which is surprisingly close to the prediction of this model.



#### gene-regulatory networks

What are gene-regulatory networks (GRNs)?

- networks between genes coding for transcription factors and genes

How does one generate GRNs?

- from co-expression + regulatory information (e.g. presence of TF binding sites)

What can these GRNs be used for?

functional interpretation of exp. data, guide inhibitor design etc.

Limitations of current GRN models:

incomplete in terms of TF-interactions, usually do not account for epigenetic effects and miRNAs

#### How does one generate GRNs?

(1.) "by hand" based on individual experimental observations

(2) Infer GRNs by computational methods from gene expression data (see reference below)

Briefings in Bioinformatics Advance Access published May 21, 2013 BRIEFINGS IN BIOINFORMATICS. page 1 of 17 doi:10.1093/bib/bbt034

#### Supervised, semi-supervised and unsupervised inference of gene regulatory networks

Stefan R. Maetschke, Piyush B. Madhamshettiwar, Melissa J. Davis and Mark A. Ragan

Submitted: I9th January 2013; Received (in nevised form): I5th April 2013

#### **Unsupervised methods**

Unsupervised methods are either based on **correlation** or on **mutual information**.

**Correlation-based network inference methods** assume that correlated expression levels between two genes are indicative of a regulatory interaction.

Correlation coefficients range from -1 to 1.

A **positive** correlation coefficient indicates an **activating interaction**, whereas a **negative** coefficient indicates an **inhibitory interaction**.

The common correlation measure by **Pearson** is defined as

$$corr(X_i, X_j) = \frac{cov(X_i, X_j)}{\sigma(X_i) \cdot \sigma(X_j)}$$

where  $X_i$  and  $X_j$  are the expression levels of genes *i* and *j*, cov(.,.) denotes the covariance, and  $\sigma$  is the standard deviation.

Bioinformatics 3 – WS 14/15

#### Rank-based unsupervised methods

Pearson's correlation measure assumes normally distributed values. This assumption does not necessarily hold for gene expression data.

Therefore rank-based measures are frequently used. The measures by Spearman and Kendall are the most common.

**Spearman's method** is simply Pearson's correlation coefficient for the ranked expression values

Kendall's 
$$\tau$$
 coefficient :  $\tau(X_i, X_j) = \frac{con(X_i^r, X_j^r) - dis(X_i^r, X_j^r)}{\frac{1}{2}n(n-1)}$ 

where  $X_{i}^{r}$  and  $X_{i}^{r}$  are the ranked expression profiles of genes *i* and *j*.

*Con(.)* denotes the number of concordant value pairs (i.e. where the ranks for both elements agree). *dis(.)* is the number of disconcordant value pairs in  $X_i^r$  and  $X_j^r$ . Both profiles are of length *n*.

#### WGCNA

WGCNA is a modification of correlation-based inference methods that **amplifies high correlation coefficients** by raising the absolute value to the power of  $\beta$  ('softpower').

$$w_{ij} = |\mathit{corr}(X_i, X_j)|^{eta}$$

with  $\beta \geq 1$ .

Because softpower is a nonlinear but monotonic transformation of the correlation coefficient, the prediction accuracy measured by AUC will be no different from that of the underlying correlation method itself.

# Unsupervised methods based on mutual information

Relevance networks (RN) introduced by Butte and Kohane measure the **mutual information (MI)** between gene expression profiles to infer interactions.

The MI *I* between discrete variables  $X_i$  and  $X_j$  is defined as

$$I(X_i, X_j) = \sum_{x_i \in X_i} \sum_{x_j \in X_j} p(x_i, x_j) \log\left(\frac{p(x_i, x_j)}{p(x_i)p(x_j)}\right)$$

where  $p(x_i, x_j)$  is the **joint probability distribution** of  $X_i$  and  $X_j$ (both variables fall into given ranges) and  $p(x_i)$  and  $p(x_i)$  are the **marginal probabilities** of the two variables (ignoring the value of the other one).

 $X_i$  and  $X_j$  are required to be discrete variables.

#### Unsupervised methods: Z-score

Z-SCORE is a network inference strategy by Prill et al. that takes advantage of knockout data.

It assumes that a knockout affects directly interacting genes more strongly than others.

The z-score  $z_{ij}$  describes the effect of a knockout of gene *i* in the *k*-th experiment on gene *j* as the normalized deviation of the expression level  $X_{jk}$  of gene *j* for experiment *k* from the average expression  $\mu(X_j)$  of gene *j*:

$$z_{ij} = |rac{X_{jk} - \mu(X_j)}{\sigma(X_j)}|$$

#### supervised inference method: SVM

In contrast to unsupervised methods, e.g. correlation methods, the supervised approach does not directly operate on pairs of expression profiles but on feature vectors that can be constructed in various ways.

E.g. one may use the outer product of two gene expression profiles  $X_i$  and  $X_j$  to construct feature vectors:

$$\mathbf{x} = X_i X_j^T$$

A sample set for the training of the SVM is then composed of feature vectors  $\mathbf{x}_i$ that are labeled  $\mathbf{x}_i = \pm 1$  for gone pairs that interact and  $\mathbf{x}_i = -1$  for these that do

that are labeled  $\gamma_i$  = +1 for gene pairs that interact and  $\gamma_i$  = -1 for those that do not interact.

#### Measure accuracy of GRNs

**Inference methods** (to infer = *dt. aus etwas ableiten/folgern*) aim to recreate the topology of a genetic regulatory network e.g. based on expression data only.

The **accuracy** of a method is assessed by the extent to which the network it infers is similar to the true regulatory network.

We quantify similarity e.g. by the area under the Receiver Operator Characteristic curve (AUC)  $1 \prod_{n=1}^{n}$ 

$$AUC = \frac{1}{2} \sum_{k=1}^{n} (X_k - X_{k-1})(Y_k + Y_{k-1})$$

where  $X_k$  is the false-positive rate and  $Y_k$  is the true positive rate for the *k*-th output in the ranked list of predicted edge weights.

An AUC of 1.0 indicates a perfect prediction, while an AUC of 0.5 indicates a performance no better than random predictions.

#### AUC



Divide data into bins.

Measure value of function Y at midpoint of bin -> factor 0.5

$$AUC = \frac{1}{2} \sum_{k=1}^{n} (X_k - X_{k-1})(Y_k + Y_{k-1})$$

www.wikipedia.org

## Summary

Network inference is a very important active research field.

Inference methods allow to construct the topologies of gene-regulatory networks solely from expression data (unsupervised methods).

Supervised methods show far better performance.

Performance on real data is lower than on synthetic data because regulation in cells is not only due to interaction of TFs with genes,

but also depends on epigenetic effects (DNA methylation, chromatin structure/histone modifications, and miRNAs).

#### **Network Reconstruction**

Experimental data: DNA microarray  $\rightarrow$  expression profiles

**Clustering**  $\rightarrow$  genes that are **regulated simultaneously** 

→ Cause and action??? Are all genes known???

Shown below are 3 different networks that lead to the same expression profiles  $\rightarrow$  combinatorial explosion of number of compatible networks

 $\rightarrow$  static information usually not sufficient



Some formalism may help

→ Bayesian networks (formalized conditional probabilities) but usually too many candidates...

#### **Network Motifs**



*Nature Genetics* **31** (2002) 64

RegulonDB + their own hand-curated findings

 $\rightarrow$  break down network into motifs

- $\rightarrow$  statistical significance of the motifs?
  - $\rightarrow$  behavior of the motifs <=> location in the network?

#### Motif 1: Feed-Forward-Loop



- X = general transcription factor
- Y = specific transcription factor
- Z = effector operon(s)

X and Y **together** regulate Z:

"**coherent**", if X and Y have the **same** effect on Z (activation vs. repression), otherwise "incoherent"

85% of the FFL in E coli are coherent

## FFL dynamics



Delay between X and Y  $\rightarrow$  signal must persist longer than delay  $\rightarrow$  reject transient signal, react only to **persistent** signals  $\rightarrow$  enables fast shutdown

Helps with **decisions** based on **fluctuating signals** 

## Motif 2: Single-Input-Module



Set of operons controlled by a single transcription factor

- same sign
- no additional regulation
- control is usually autoregulatory (70% vs. 50% overall)

Mainly found in genes that code for **parts** of a protein **complex** or metabolic **pathway** (here machinery for arginine biosynthesis)  $\rightarrow$  relative stoichiometries

## **SIM-Dynamics**



If different thresholds exist for each regulated operon:

- $\rightarrow$  first gene that is activated is the last that is deactivated
  - $\rightarrow$  well defined temporal ordering (e.g. flagella synthesis) + stoichiometries

## Motif 3: Densely Overlapping Regulon



Dense layer between groups of transcription factors and operons → much denser than network average (≈ community)

Usually each operon is regulated by a different combination of TFs.

Main "computational" units of the regulation system

Sometimes: same set of TFs for group of operons  $\rightarrow$  "multiple input module"

## **Detection of motifs**

Represent transcriptional network as a connectivity matrix *M* such that  $M_{ij} = 1$  if operon *j* encodes a TF that transcriptionally regulates operon *i* 

and  $M_{ij} = 0$  otherwise.

Scan all  $n \times n$  submatrices of *M* generated by choosing *n* nodes that lie in a connected graph, for n = 3 and n = 4.

Submatrices were enumerated efficiently by recursively searching for nonzero elements.



Connectivity matrix for causal regulation of transcription factor j (row) by transcription factor i (column). Dark fields indicate regulation. (Left) Feed-forward loop motif. TF 2 regulates TFs 3 and 6, and TF 3 again regulates TF 6. (Middle) Single-input multiple-output motif. (Right)

Densely-overlapping region.

Compute a P value for submatrices representing each type of connected subgraph by comparing # of times they appear in real network vs. in random network.

For n = 3, the only significant motif is the feedforward loop. For n = 4, only the overlapping regulation motif is significant. SIMs and multi-input modules were identified by searching for identical rows of *M*. Shen-Orr et al. Nature Gen. 31, 64 (2002)

Bioinformatics 3 – WS 14/15

#### **Motif Statistics**

Structure	Appearances in real network	Appearances in randomized network (mean ± s.d.)	<i>P</i> value
Coherent feedforward loop	34	4.4 ± 3	<i>P</i> < 0.001
Incoherent feedforward loop	6	2.5 ± 2	<i>P</i> ~ 0.03
Operons controlled by SIM (>13 operons)	68	28 ± 7	<i>P</i> < 0.01
Pairs of operons regulated by same two transcription factors	203	57 ± 14	<i>P</i> < 0.001
Nodes that participate in cycles*	0	$0.18 \pm 0.6$	<i>P</i> ~0.8

All motifs are highly overrepresented compared to randomized networks

No cycles  $(X \rightarrow Y \rightarrow Z \rightarrow X)$  were identified, but this was not statistically significant in comparison to to random networks

Bioinformatics 3 – WS 14/15

Shen-Orr et al., Nature Genetics 31 (2002) 64

#### **Network with Motifs**



(flagella and nitrogen systems)

## Summary

#### Today:

- Gene regulation networks have **hierarchies**:
- $\rightarrow$  global "cell states" with specific expression levels
- Network motifs: FFLs, SIMs, DORs are overrepresented
- $\rightarrow$  different functions, different temporal behavior