# Reverse engineering of genotypephenotype map using individual genetic variation

Vitaly V. Gursky<sup>1,2</sup>, Konstantin N. Kozlov<sup>2</sup>, Ivan V. Kulakovskiy<sup>3-5</sup>, Sergey V. Nuzhdin<sup>6</sup>, Maria G. Samsonova<sup>2</sup>

- <sup>1</sup> Ioffe Institute, St. Petersburg, Russia
- <sup>2</sup> Peter the Great St. Petersburg Polytechnic University, St. Petersburg, Russia
- <sup>3</sup> Engelhardt Institute of Molecular Biology, Moscow, Russia
- <sup>4</sup> Vavilov Institute of General Genetics, Moscow, Russia
- <sup>5</sup> Center for Data-Intensive Biomedicine and Biotechnology, Skolkovo Institute of Science and Technology, Moscow, Russia
- <sup>6</sup> University of Southern California, Los Angeles, U.S.A.



(Image by Brooke Wolford; http://misciwriters.com/2016/11/29)



(Image by Brooke Wolford; http://misciwriters.com/2016/11/29)



(Image by Brooke Wolford; http://misciwriters.com/2016/11/29)

Alternative approach:

**Modeling** gene regulatory networks on an 'average' genome and **predicting** how genetic variation effects **molecular phenotype** (gene expression patterns)



(Image by Brooke Wolford; http://misciwriters.com/2016/11/29)

Alternative approach:

**Modeling** gene regulatory networks on an 'average' genome and **predicting** how genetic variation effects **molecular phenotype** (gene expression patterns) Model organisms and processes are helpful as a training ground



#### Gene expression



#### Gene expression



#### Gene expression



**Transcriptional regulation** 



#### Gene regulatory network



### **Reaction-diffusion equations**



$$\frac{\partial u}{\partial t} = R\left(x,t\right) - \lambda_R u(x,t) + D_R \Delta u(x,t)$$

$$\frac{\partial v}{\partial t} = P(x,t) - \lambda_P v(x,t) + D_P \Delta v(x,t)$$

#### Segmentation genes in Drosophila development



## Gap gene network





#### Gap gene network



gene expression

#### Gap gene network



% embryo length

90

35

#### Information about binding sites

Probability of TF binding to DNA according to ChIP-Seq data for *D. melanogaster*:



## **Binding motif**

Processing data:



#### How many binding sites can we expect?

Connection of **PWM-scores** with **binding energy**:

Berg and von Hippel PH (1987) J Mol Biol, 193; Stormo and Fields (1998) Trends Biochem Sci 23.



Typical binding energy histograms (in log scale):

From Sheinman et al. (2012) Rep. Progr. Phys., 75(2)

#### Known enhancers are not always enough

# Quantitative and predictive model of transcriptional control of the *Drosophila melanogaster even skipped* gene

Hilde Janssens<sup>1</sup>, Shuling Hou<sup>2</sup>, Johannes Jaeger<sup>1</sup>, Ah-Ram Kim<sup>1</sup>, Ekaterina Myasnikova<sup>3</sup>, David Sharp<sup>4</sup> & John Reinitz<sup>1</sup>

OPEN ORCESS Freely available online

PLOS GENETICS

nature

# Rearrangements of 2.5 Kilobases of Noncoding DNA from the *Drosophila even-skipped* Locus Define Predictive Rules of Genomic *cis*-Regulatory Logic

Ah-Ram Kim<sup>1,2</sup>, Carlos Martinez<sup>1</sup>, John Ionides<sup>3</sup>, Alexandre F. Ramos<sup>4</sup>, Michael Z. Ludwig<sup>1</sup>, Nobuo Ogawa<sup>5</sup>, David H. Sharp<sup>6</sup>, John Reinitz<sup>1,7</sup>\*

#### Open chromatin regions match binding patterns



Figure 1 DNasel accessibility and *in vivo* DNA binding by transcription factors across the *eve* locus.

### Analysis of gap gene regulatory regions

Improved D. melanogaster major **position matrix motifs by ChiPMunk** software (http://autosome.ru/iDMMPMM/) (Kulakovskiy et al. 2009, 2010). Example of information obtained by the software for Kruppel's regulatory region:



Select binding sites from the **open chromatin regions** (according to DNaseI accessibility data). Total number of sites: from 889 (28 per TF per target gene) to 1419 (44 per TF per target gene).

21

Configurations:

Statistical weights:



$$q_i = Kv_i \exp\left(-\Delta G_i\right) = Kv_i \exp\left(P(S_i) - P(S_{\max})\right)$$

Configurations:

Statistical weights:



$$q_i = Kv_i \exp\left(-\Delta G_i\right) = Kv_i \exp\left(P(S_i) - P(S_{\max})\right)$$

$$W(\sigma) = \prod_{i} q_{i} \prod_{j} \omega_{ij}$$

Configurations:

 $\omega_{12}$ 

 $\omega_{23}$ 

Statistical weights:

$$q_i = Kv_i \exp\left(-\Delta G_i\right) = Kv_i \exp\left(P(S_i) - P(S_{\max})\right)$$





$$W(\sigma)Q(\sigma), \ Q(\sigma) = \prod_{i} \alpha_{i}$$
  
 $\alpha_{i}$  - activation strengths.

Configurations:

 $\omega_{12}$ 

Statistical weights:

$$q_i = Kv_i \exp\left(-\Delta G_i\right) = Kv_i \exp\left(P(S_i) - P(S_{\max})\right)$$



$$W(\sigma)Q(\sigma), \ Q(\sigma) = \prod_{i} \alpha_{i}$$
  
$$\alpha_{i} - \text{activation strengths.}$$



 $\alpha_1$ 

 $\alpha_2$ 

B

Short-range repression:  $q_i \rightarrow \beta q_i$  $\beta$  - repression strengths. Transcriptional activation via thermodynamic approach



$$W(\sigma) = \prod_{i} q_i \prod_{j} \omega_{ij}$$

$$W(\sigma)Q(\sigma), \ Q(\sigma) = \prod_{i} \alpha_{i}$$
  
 $\alpha_{i}$  - activation strengths.

Probability of transcriptional activation of gene 'a':

$$E_a = \frac{Z_{ON}}{Z_{ON} + Z_{OFF}} = \frac{\sum_{\sigma} W(\sigma)Q(\sigma)}{\sum_{\sigma} W(\sigma)Q(\sigma) + \sum_{\sigma} W(\sigma)}$$

He et al. PLoS Comput. Biol., 2010; Kozlov et al. BMC Genomics, 2014, 2015

#### Hybrid dynamical model for gap gene network

Reaction-diffusion equations for mRNA concentrations (  $\mathcal{U}_a$ ) and protein concentrations (  $\mathcal{V}_a$ ):

$$E_a = \frac{Z_{ON}}{Z_{ON} + Z_{OFF}} = \frac{\sum_{\sigma} W(\sigma)Q(\sigma)}{\sum_{\sigma} W(\sigma)Q(\sigma) + \sum_{\sigma} W(\sigma)}$$

$$\frac{\partial u_a}{\partial t} = r_a E_a \left( \mathbf{v}(x,t), \mathbf{V}(x,t) \right) - \lambda_a u_a(x,t) + d_a \frac{\partial^2 u_a}{\partial x^2},$$

$$\frac{\partial v_a}{\partial t} = R_a u_a (t - \tau) - \Lambda_a v_a (x, t) + D_a \frac{\partial^2 v_a}{\partial x^2}$$

 $\mathbf{v}(x,t)$  – vector of gap protein concentrations (Hb, Kr, Gt, Kni).  $\mathbf{V}(x,t)$  – vector of protein concentrations for external regulators (Bcd, Cad, Tll, Hkb). au – delay time (transcription+translation times).

#### Fitting to wild-type expression patterns

The model was fitted to the gap gene expression data for cleavage cycles 13-14A.



#### Cross-validation, negative control, local identifiability analysis



- 1: cross-validation
- 2: original fitting
- 3: fitting to random data

#### Testing on expression in reporter constructs

*In silico* predictions of expression patterns for various reporter constructs (Schroeder et al., 2004; Gallo et al., 2010):



#### Predictions for the regulatory landscape: Many weak binding sites working in concert

We calculate the **regulatory weight** of each binding site as the RMS-difference between the model outputs with and without this site.



## Predictions for the regulatory landscape: Many weak binding sites working in concert

We calculate the **regulatory weight** of each binding site as the RMS-difference between the model outputs with and without this site.



of sites) did not change the patterns qualitatively. This shows that proximate sites demonstrate extensive **compensatory actions**.

#### Direct calculation of compensatory effects



Heatmap of correlations between binding sites:



#### Binding affinity vs. influence on expression

Another type of 'weakness' in the concept of 'weak sites working in concert': Only **weak correlation** between the strength of influence on expression and the binding affinity for binding sites:



#### Simulating evolution of regulatory regions under elevated mutational pressure



#### Dynamics of the number of binding sites



#### Dynamics of the number of binding sites



The dynamics of the population average number of overlapping events



#### Core binding sites



Distributions of site scores for wt conditions (zero generation):



#### Analysis of genetic variation in Drosophila lines

We apply the gene expression model to analyze single nucleotide polymorphisms (**SNPs**) in the regulatory regions of gap genes in a panel of **213 natural sequenced** *D. melanogaster* lines from two populations (Raleigh, NC, and Winters, CA) (Campo et al. 2013).



### Contrasting natural genetic variation with simulated SNPs

Simulation of random point mutations:

1) **One** point mutation per genotype:



#### Contrasting natural genetic variation with simulated SNPs

Simulation of random point mutations:

2) Sets of SNPs per genotype, considering 2 frequency spectra:

2.1) **Neutral spectrum**: the frequency distribution of SNPs extracted from short intron positions of the D. mel genome.

2.2) Population-derived spectrum



## Contrasting natural genetic variation with simulated SNPs

Simulation of random point mutations:

2) Sets of SNPs per genotype, considering 2 frequency spectra:



#### Translating natural genetic variation to gene expression



#### Model simulation of the population genotypes





# Weighted pattern generating potential (wPGP) is better than RMS difference

$$w = \sum_{t,a} f^{a}(t), \quad f^{a} = 0.5 - 0.5 * (\text{reward} - \text{penalty}),$$
  
reward = 
$$\frac{\sum_{i} V_{i}^{a} \min(V_{i}^{a}, v_{i}^{a})}{\sum_{i} (V_{i}^{a})^{2}}, \quad \text{penalty} = \frac{\sum_{i} \max(0, v_{i}^{a} - V_{i}^{a})(V_{\max}^{a} - V_{i}^{a})}{\sum_{i} (V_{\max}^{a} - V_{i}^{a})^{2}},$$

Characteristic	Expression	PGP	CC	1-RMSE
Sensitive to scaling		0.81	0.96	0.81
		0.58	0.96	0.43
Sensitive to shift in basal expression		0.69	0.96	0.62
		0.39	0.96	0.71
Normalization for length of expression domain		0.69	0.96	0.62
		0.69	0.94	0.74
Sensitive to partial pattern		0.56	0.40	0.53
		0.51	0.63	0.62

Figure 3B from Kazemian et al. (2010) Plos Biol. 8: e1000456

#### Sign of SNP influence on expression

SNP's influence on expression can be **sign-alternating**:



#### Specific examples of sign analysis



#### Specific examples of sign analysis



Other examples of alternating sign of SNP influence on expression:



#### SNPs activate (repress) expression via repressing (activating) sites

Distribution of SNPs from the population among activating and repressing binding sites:



#### Combination of influences from multiple SNPs

Mutating a set of *n* binding sites in the regulatory region of a gene:

$$S_i \to \tilde{S}_i, \quad i = 1, \dots, n$$

 $\delta_i$  – scaled difference between a PCM matrix element due to *i* th SNP.

Probability of transcriptional activation for the mutated regulatory region:

$$\tilde{E} = \frac{Z_{ON} + P_n(\delta_1, \dots, \delta_n)}{Z_{ON} + Z_{OFF} + Q_n(\delta_1, \dots, \delta_n)}$$

$$\Delta v_{gen}$$
 = nonlinear function of  $\sum_{SNP} \Delta v_{SNP}$ 

#### SNP influence combination in the population genotypes



#### SNP influence combination in the population genotypes



#### OPEN O ACCESS Freely available online

PLOS GENETICS

#### Data and Theory Point to Mainly Additive Genetic Variance for Complex Traits

William G. Hill<sup>1</sup>\*, Michael E. Goddard<sup>2,3</sup>, Peter M. Visscher<sup>4</sup>

1 Institute of Evolutionary Biology, School of Biological Sciences, University of Edinburgh, Edinburgh, United Kingdom, 2 Faculty of Land and Food Resources, University of Melbourne, Victoria, Australia, 3 Department of Primary Industries, Victoria, Australia, 4 Queensland Institute of Medical Research, Brisbane, Australia

Individual SNPs do not seem to be under purifying selection, but combinations of SNPs do.



Individual SNPs do not seem to be under purifying selection, but combinations of SNPs do.



Only SNP **combinations with strong influence** can be distinguished from random SNPs under the **population-derived spectrum**.

#### RMS and wPGP measures lead to different evolutionary predictions



## Thank you

#### Systems biology and bioinformatics Lab, Peter the Great Polytechnic University

Maria Samsonova Sergey Nuzhdin Konstantin Kozlov Svetlana Surkova Alyona Sokolkova Alexandra Chertkova Anna Igolkina

#### **loffe Institute**

Alexander Samsonov Sergey Rukolaine

#### University of Southern California

Sergey Nuzhdin Paul Marjoram Engelhardt Institute of Molecular Biology,

Ivan Kulakovskiy