# Bayesian Markov models consistently surpass PWMs at predicting regulatory motifs

## 30 years TRANSFAC

### Goettingen
### 8. March 2018

## Johannes Söding, MPI for Biophysical Chemistry, Göttingen

# Methods and research topics of the Söding lab

**Computational metagenomics**
- Fast & deep sequence searching and clustering
- Protein-level assembly
- Large-scale binning & X-assembly
- Gene prediction

**Protein function & structure**
- Coevolution analysis
- Viral metagenomics (Virus-X)
- Functional module discovery in massive metagenomic data

smart algorithms

fast parallelized code

**Bayesian statistical modeling**

**(Post-) Transcriptional regulation**
- Regulatory motif discovery
- RNA-protein binding cooperat.
- PAR-CLIP data analysis
- NGS data analysis

**Systems medicine of complex diseases**
- Risc locus prediction
- Risk variant fine-mapping
- Discovery of drug targets from GWAS & eQTL data

**Single-cell transcriptomics**
- Reconstruction of cellular lineage trees
- Denoising scRNA-seq data

# Why are we interested in transcriptional regulation?

# Genetic causes of common diseases linked to dysregulation of gene networks

(SNP = position in genome with variation in population)

10,000 patients with disease

cases

62% C
38% T

Cytosine increases risk for disease

10,000 healthy people

controls

49% C
51% T

- **≥ 90% of causal SNPs non-coding**

- These SNPs **disrupt transcription factor binding sites** and thereby influence the expression of target genes

# How is an organism encoded in its genome?



embryo

Eve expressed in 7 stripes

1 2 3 4 5 6 7

enhancers

eve

4+6

LacZ expression

4+6

lacZ

Genomically inserted
4+6 enhancer with lacZ reporter

- Genomes contain all information for a single cell to develop into a complex organism and to survive and procreate

- Genomes are molecular programs, which are read by transcription factors binding to specific DNA sequences.

- Transcription rates are the result of complex molecular computations at promoters and enhancers

- We want to understand and predict these molecular computations

# Are we there yet?

## (What I cannot create, I do not understand)



**Transcription factors missing?**

**TF-TF interactions?**

**Weak binding sites missing?**

**Binding site strengths?**

type enhancer

reconsituted 1

reconsituted 2

Activators: Bcd  Hb  Zelda

Repressors: Gt  Kr

We cannot reconstitute even best-studied enhancers with designed sequences

BJ Vincent, J Estrada & AH DePace, Integrative Biol. 2016

# De-novo motif discovery: binding site motifs for TFs enriched in set of sequences

DNA

- **ChIP-seq**
- **SELEX-seq**
- **Protein binding microarrays** (PBMs)
- DNase-seq, FAIRE-seq, ATAC-seq: open chromatin
- CAGE, RACE: transcription start sites
- RNA-seq: co-expressed genes
- Hi-C, ChIA-PET: loops and 3D nuclear structure
- ...

RNA

- PAR-CLIP, ICLIP: RNAs bound by RNA-binding factors
- SELEX-seq: binding motifs
- RNA-seq: 3 'UTRs of co-expressed genes

motif

DNA

# Postion weight matrices (PWMs) assume independence of nucleotides within site

# But how important are correlations among nucleotides in regulatory motifs?

Correlations between neighboring nucleotides:

◆ Shape readout of DNA



bend      kink      minor groove width

◆ Multiple (sequence-dependent) binding modes

◆ Variable spacers between half-sites

◆ Complex combination of motifs at varying distances, e.g. through multiple DNA binding domains, collaborative binding etc.

# Diversity and Complexity in DNA Recognition by Transcription Factors

Gwenael Badis,[1]* Michael F. Berger,[2,3]* Anthony A. Philippakis,[2,3,4]* Shaheynoor Talukder,[1,5]* Andrew R. Gehrke,[2]* Savina A. Jaeger,[2]* Esther T. Chan,[5]* Genita Metzler,[6] Anastasia Vedenko,[7] Xiaoyu Chen,[1] Hanna Kuznetsov,[6] Chi-Fong Wang,[8] David Coburn,[1] Daniel E. Newburger,[2] Quaid Morris,[1,5,9,10] Timothy R. Hughes,[1,5,10]† Martha L. Bulyk[2,3,4,11]†

# Quantitative analysis demonstrates most transcription factors require only simple models of specificity

Y Zhao & G Stormo, *Nature Biotech* 29: 480 – 483 (2011).

…

# Jury remains out on simple models of transcription factor specificity

Q Morris, ML Bulyk, TR Hughes, *Nature Biotech* 29: 483 – 485 (2011).

Protein binding microarray

# Evaluation of methods for modeling transcription factor sequence specificity

Matthew T Weirauch[1,2], Atina Cote[1], Raquel Norel[3], Matti Annala[4], Yue Zhao[5], Todd R Riley[6], Julio Saez-Rodriguez[7], Thomas Cokelaer[7], Anastasia Vedenko[8], Shaheynoor Talukder[1], DREAM5 Consortium[9], Harmen J Bussemaker[6], Quaid D Morris[1,10], Martha L Bulyk[8,11,12], Gustavo Stolovitzky[3], Timothy R Hughes[1,10]

Genomic analyses often involve scanning for potential transcription factor (TF) binding sites using models of the sequence specificity of DNA binding proteins. Many approaches have been developed to model and learn a protein's DNA-binding specificity, but these methods have not been systematically compared. Here we applied 26 such approaches to *in vitro* protein binding microarray data for 66 mouse TFs belonging to various families. For nine TFs, we also scored the resulting motif models on *in vivo* data, and found that the best *in vitro*–derived motifs performed similarly to motifs derived from the *in vivo* data. Our results indicate that simple models based on mononucleotide position weight matrices trained by the best methods perform similarly to more complex models for most TFs examined, but fall short in specific cases (<10% of the TFs examined here). In addition, the best-performing motifs typically have relatively low information content, consistent with widespread degeneracy in eukaryotic TF sequence preferences.

Weirauch et al., Nature Biotechnology (2013)

# The Next Generation of Transcription Factor Binding Site Prediction

Anthony Mathelier*, Wyeth W. Wasserman*   *PLoS Comput Biol* 9, e1003214 (2013)

TFFMs (1st order Markov models) perform >5% better than PWMs in 20% of cases



96 ChIP-seq ENCODE datasets

# Markov Models (MMs) model correlations among nucleotides

k'th order MM: probability depends on *k* previous nucleotides

$$\text{Score}(x_1 \ldots x_W) = \sum_{j=1}^{W} \log \frac{p_j(x_j)}{p_{\text{bg}}(x_j)}$$

0th order, PWM

$$\text{Score}(x_1 \ldots x_W) = \sum_{j=1}^{W} \log \frac{p_j(x_j | x_{j-1})}{p_{\text{bg}}(x_j | x_{j-1})}$$

1st order

$$\text{Score}(x_1 \ldots x_W) = \sum_{j=1}^{W} \log \frac{p_j(x_j | x_{j-1}, x_{j-2})}{p_{\text{bg}}(x_j | x_{j-1}, x_{j-2})}$$

# Markov Models (MMs) model correlations among nucleotides

k'th order MM: probability depends on *k* previous nucleotides

| A | T | C | G | C | T | **A** | $\cdots$

$$p_j(\mathbf{A}) = \frac{\#(\mathbf{A})}{\#\text{ seqs}}$$

0th order, PWM

| A | T | C | G | C | **T** | **A** | $\cdots$

*j*-1 *j*

$$p_j(\mathbf{A}|\mathbf{T}) = \frac{\#(\mathbf{TA})}{\#(\mathbf{T})}$$

1st order

| A | T | C | G | C | **T** | **A** | $\cdots$

*j*-1 *j*

$$p_j(\mathbf{A}|\mathbf{CT}) = \frac{\#(\mathbf{CTA})}{\#(\mathbf{CT})}$$

2nd order

$\vdots$

# Markov Models (MMs) model correlations among nucleotides

*k*'th order MM: probability depends on *k* previous nucleotides

$\cdots$ | A | T | C | G | C | T | **A** | $\cdots$

$p_j(\mathbf{A}) = \dfrac{\#(\mathbf{A}) + 1}{\#\text{ seqs} + 4}$     0th order, PWM

*j*

Pseudo-counts

$\cdots$ | A | T | C | G | C | **T** | **A** | $\cdots$

$p_j(\mathbf{A}|\mathbf{T}) = \dfrac{\#(\mathbf{TA}) + 1}{\#(\mathbf{T}) + 4}$     1st order

*j*-1 *j*

$\cdots$ | A | T | C | G | **C** | **T** | **A** | $\cdots$

$p_j(\mathbf{A}|\mathbf{CT}) = \dfrac{\#(\mathbf{CTA}) + 1}{\#(\mathbf{CT}) + 4}$     2nd order

For order *k* one needs **~100 ×4$^{k+1}$** sequences (!)
to learn probabilities with 10% relative accuracy

# Many higher-order models prune the dependency graph and pool contexts



inhomogeneous parsimonious Markov model

inhomogeneous variable-order Markov model

**Optimization requires comparing very many discrete tree topologies**
**⇒ Slow and challenging to train (model comparison)**
**⇒ Cannot discover motifs de-novo, require pre-aligned motif sequences**

# We use pseudocounts from lower-order!

$\cdots$ $\boxed{A}$ $\boxed{T}$ $\boxed{C}$ $\boxed{G}$ $\boxed{C}$ $\boxed{\textbf{T}}$ $\boxed{\textbf{A}}$ $\cdots$
$j$

$$p_j(\textbf{A}|\textbf{T}) = \frac{\#(\textbf{TA}) + 20\; p_j(\textbf{A})}{\#(\textbf{T}) + 20}$$

$\cdots$ $\boxed{A}$ $\boxed{T}$ $\boxed{C}$ $\boxed{G}$ $\boxed{C}$ $\boxed{T}$ $\boxed{A}$ $\cdots$

$$p_j(\textbf{A}|\textbf{CT}) = \frac{\#(\textbf{CTA}) + 60\; p_j(\textbf{A}|\textbf{T})}{\#(\textbf{CT}) + 60}$$

$\cdots$ $\boxed{A}$ $\boxed{T}$ $\boxed{C}$ $\boxed{G}$ $\boxed{C}$ $\boxed{T}$ $\boxed{A}$ $\cdots$

$$p_j(\textbf{A}|\text{G}\textbf{CT}) = \frac{\#(\text{GCTA}) + 180\; p_j(\textbf{A}|\textbf{CT})}{\#(\text{GCT}) + 180}$$

$\vdots$

Siebert and Soding, NAR 2016

# We use pseudocounts from lower-order!

**If many counts for k-mer**

$\Rightarrow$ counts dominate over pseudocounts

$\Rightarrow$ use maximum likelihood estimate

$\cdots$ | A | T | C | G | C | **T** | **A** | $\cdots$

$j$

$$p_j(\text{A}|\text{T}) = \frac{400 + 20\, p_j(\text{A})}{500 + 20}$$

$\cdots$ | A | T | C | G | C | T | **A** | $\cdots$

$j$

$$p_j(\text{A}|\text{CT}) = \frac{360 + 60 \times 0.8}{400 + 60}$$

$\cdots$ | A | T | C | G | C | T | **A** | $\cdots$

$$p_j(\text{A}|\text{GCT}) = \frac{340 + 180 \times 0.9}{350 + 180}$$

Pseudo-counts

$\vdots$

Siebert and Soding, NAR 2016

# We use pseudocounts from lower-order!

**If few counts for k-mer**

⇒ pseudocounts dominate over counts

⇒ fall back on lower-order estimate

$$p_j(\text{C}|\text{A}) = \frac{50 + 20 \times 0.4}{200 + 20}$$

Pseudo-counts

$$p_j(\text{C}|\text{TA}) = \frac{4 + 60 \times 0.25}{40 + 60}$$

$$p_j(\text{C}|\text{CTA}) = \frac{1 + 180 \times 0.2}{10 + 180}$$

*Bayesian Markov models* only learn parameters for which enough information exists to estimate accurately!

No need for optimizing discrete dependency trees.

# 5th order BaMMs learn binding motifs from ChIP-seq better than PWMs

Increase of partial Area under ROC curve at 5% FPs (pAUC)
for each of 446 ENCODE ChIP-seq datasets (4-fold cross-validated)



100:1 neg:pos,
neg. sequences from
2nd order hom. model

# Gains of 5th order BaMMs over PWMs grow when including flanking nucleotides

Increase of pAUC on 446 ENCODE ChIP-seq sets for +8bp-extended models

# 5th order BaMMs achieve sizeable gains even over 1st order BaMMs

Increase of pAUC on 446 ENCODE ChIP-seq sets for +8bp-extended models



Improvements smaller when training on ChIP-seq data and testing on HT-SELEX data!

Some of the improvement could be due to learning secondary motifs in higher orders

# Klf4 motifs trained on ChIP-seq, tested on EMSA affinities of mutated binding sites

# Klf4 motifs trained on ChIP-seq, tested on EMSA affinities of mutated binding sites

**FoxA2 motifs trained on ChIP-seq tested on EMSA affinities of mutated binding sites**

# Detecting fly narrow-peak TSSs

## (CAGE data from Adelman lab)

Detecting fly broad-peak TSSs

(CAGE data from Adelman lab)

# Pause sites in E. coli
## (data from Landick lab)

# BaMMs are robust to overtraining

# RNA-binding sites from PAR-CLIP: higher order vs. PWM

(data from Cramer lab, Goettingen)

# BaMMmotif server offers 4 tools

https://bammmotif.mpibpc.mpg.de

# How to assess motif models?

- E-values / P-values for enrichment of motif occurrences in input set vs. background set are popular

- But P-values can be very significant for motifs without biological relevance, e.g. input set is large and background set is not 100% realistic

- We need a quality measure that informs us about how well the model will identify binding sites in unseen datasets

- The demands on model specificity depend on the expected ratio of positive to negative sequences! E.g. ChIP-seq: ~1:1, scanning promoters: 1:100

# How to assess motif models?

**Use ROC-like analysis**

# Partial ROC curve?

Relevant range of false positive rate (FPR) depends on expected pos:neg ratio!

# Precision-recall curve?

Relevant range of precision depends on expected pos:neg ratio!

# TP/FP ratio-recall curve!

- Covers entire relevant range of precision or FDR (log scale)
- Different ratios pos:neg simply result in shifted curves

# Summary (main part)

▸ **Higher-order correlations significantly contribute to the binding specificity of most transcription factors**

▸ **Modeling higher correlations in a Bayesian framework can significantly improve predictions on 97% of tested factors, by +36% in pAUC on average. BUT: how much improvements due to learning >1 model in ChIP-seq data?**

▸ **BaMMs are very robust. They never overtrained**

▸ **BaMMs excel in learning complex motif architectures**

▸ **We should move from PWMs to higher-order models**

▸ **BaMMmotif server at https://bammmotif.mpibpc.mpg.de**

**Thank you for listening!**