

Finding regulatory elements shared by a set of genes

Matthieu Defrance - Hélène Touzet - LIFL

January 2006

Introduction

Method

Examples

Conclusion

Overview

- ▶ Motif over-representation in regulatory regions
- ▶ A fast algorithm to extract significant local over-representation
- ▶ Example of Rel/NF- κ B target genes and Muscle specific genes

Biological questions

- ▶ Understanding gene transcriptional regulation in **higher eukaryotes**
- ▶ Detecting Transcription Factors involved in regulatory mechanisms

Over-represented motifs & regulation

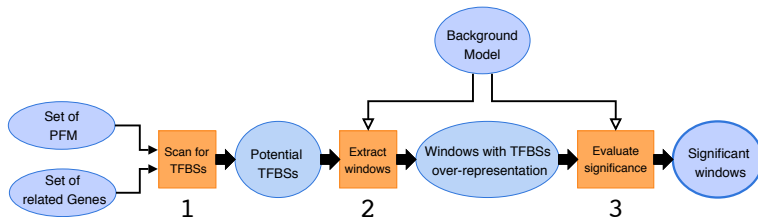
Hypothesis: over-represented motifs are involved in regulation.

- ▶ Working with a set of genes that (are assumed to) share regulatory mechanisms
 - ▶ Functionally related genes
 - ▶ Clusters of genes derived from DNA array analysis
 - ▶ ...
- ▶ A motif can be:
 - ▶ Aligo-nucleotide
 - ▶ Motif whose profile is known: Position Frequency Matrix, HMM, regular expression
- ▶ Need a **background model** to evaluate over-representation
 - ▶ Markov model
 - ▶ Empiric model

Finding motifs over-representation

- ▶ Regulatory motifs are highly degenerated in higher eukaryotes
- ▶ In order to provide accurate predictions we choose to:
 - ▶ Restrict motif search to known profiles
 - ▶ Use motif conservation across multiple species

Finding motifs over-representation: TFM-Explorer



1. **Scan for all potential TFBSs** (exhaustively)
2. **Extract regions** where predicted TFBSs are over-represented
3. **Evaluate significance** of extracted windows (P-value, E-value)

1. Scan for potential TFBSs

Input:

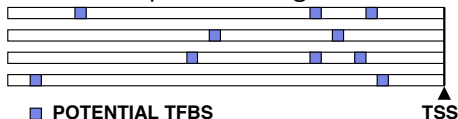
- ▶ Set of regulatory sequences
- ▶ Set of PFM (for example all TRANSFAC matrices)

Output:

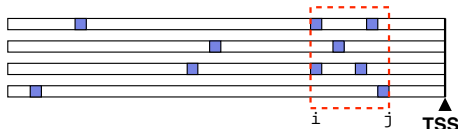
- ▶ Exhaustive list of potential **T**ranscription **F**actor **B**inding **S**ite
 - ▶ Overlapping sites (for a TF) are cut
 - ▶ Forward and reverse strands are scanned
 - ▶ Use an uniform cutting threshold based on P-value

2. Extract windows

- ▶ Sequences are aligned on **T**ranscription **S**tart **S**ite
- ▶ All TFBSs found in the previous stage are used

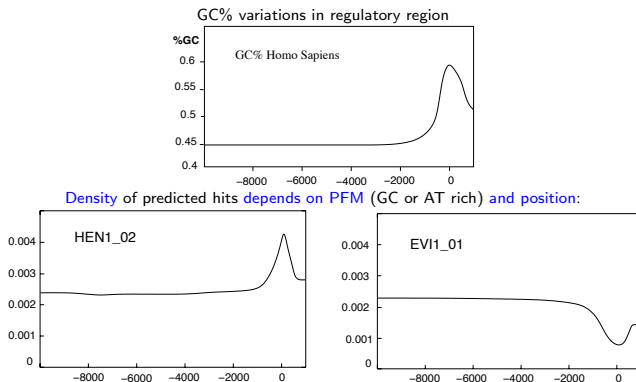


- ▶ Search for common regulatory elements in promoter regions of these genes?



Detect regions where predicted binding sites are **locally over-represented**

2. Extract windows (2) - Background Model?



We need to determine a local distribution of predicted TFBSs for each matrix

-
- Diagram illustrating the potential for Topologically Anomalous Surface States (TASS) in a 2D system. The diagram shows five horizontal lines representing energy levels. Blue squares represent potential TASS. A vertical dashed red line labeled '1' indicates a specific energy level. A black triangle labeled 'TSS' is at the bottom right.

$$s_i = k_i \ln \frac{\lambda_i}{\lambda_i^b} + |E|(\lambda_i^b - \lambda_i)$$

- ▶ E sequences set
- ▶ k_i number of hits at position i
- ▶ λ_i^b Poisson parameter in the background model at position i
- ▶ λ_i Poisson parameter in the expected model at position i

2. Extract windows (3) - Scoring function

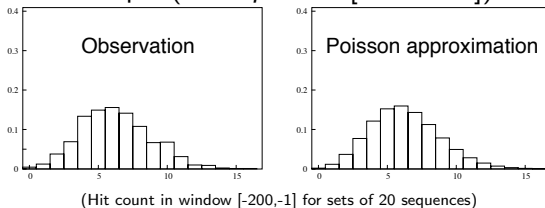
Extension to heterogeneous sequences (multiple organisms)

$$s_i = \left(\frac{\sum_{j=1}^{|E|} \lambda_{i,j}^{\text{green}}}{\sum_{j=1}^{|E|} \lambda_{i,j}^{\text{blue}}} \right)^{k_i} e^{|E| \left(\sum_{j=1}^{|E|} \lambda_{i,j}^{\text{blue}} - \sum_{j=1}^{|E|} \lambda_{i,j}^{\text{green}} \right)}$$

- ▶ E sequences set
- ▶ k_i number of hits at position i
- ▶ $\lambda_{i,j}^b$ Poisson parameter in the background model at position i for sequence j
- ▶ $\lambda_{i,j}$ Poisson parameter in the expected model at position i for sequence j

3. Evaluate significance of window

- ▶ P-value: Probability $P(k, \alpha, \beta, N)$ to observe $X \geq k$ hits in window $[\alpha, \beta]$ for $|E|$ sequences
- ▶ Approximate hit count distribution in a window by a **Poisson** distribution
- ▶ IK3_01 PFM example (*20 sequences* $[-200, -1]$):



3. Evaluate significance of window (2)

Heterogeneous sequences set case

$$P(X \geq k) = 1 - \sum_{z=0}^{k-1} \frac{(|w| \sum_{j=1}^{|E|} \lambda_j)^z}{z!} e^{-|w| \sum_{j=1}^{|E|} \lambda_j}$$

- ▶ E sequences set
- ▶ w window
- ▶ k number of hits in w ($[\alpha, \beta]$)
- ▶ λ_j Poisson parameter of his count distribution in w for the sequence j

3. Evaluate significance of window (3) - Limitations

Is this (Poisson approximation) correct for all matrices?

Fitting Observation/Poisson with χ^2 test			
α (error type I)	[100%, 5%]]5%, 0.5%]]0.5%, 0.0]
% PFM	70%	10%	20%

Application example NF- κ B target genes set

[Karo Gosselin, Corinne Abbadie IBL]

- ▶ NF- κ B control expression of genes involved in
 - ▶ Inflammation and immunity
 - ▶ Stress responses, including apoptosis
- ▶ Gene set compiled from literature data
- ▶ Gene considered as true NF- κ B target when
 - ▶ A NF- κ B was map
 - ▶ The functionality of motif was experimentally validated

NF- κ B target genes set (2)

- ▶ Set of 102 human genes
 - ▶ Promoter sequences retrieved from University California Santa Cruz Genome Browser (region $[-10000 + 1000]$)
 - ▶ List of NF- κ B motif with exact position, exact sequence (for validation)
- ▶ Empiric Background model
- ▶ All TRANSFAC matrices

NF- κ B target genes set (3) - Results of TFM-Explorer

Factor	Matrix ID	Location	Hits	Sequences	P-Value
TATA	V\$TATA_01	[-0074:-0009]	042	042 (41.18%)	3.57e-14
NF-kappaB	V\$NFKB_C	[-0507:-0016]	206	087 (85.29%)	4.89e-14
NF-kappaB	V\$NFKAPPAB65_01	[-0520:-0013]	192	084 (82.35%)	1.93e-13
NF-kappaB	V\$NFKAPPAB_01	[-0227:-0017]	112	075 (73.53%)	2.71e-13
NF-kappaB	V\$NFKB_Q6	[-0230:-0020]	096	069 (67.65%)	3.12e-11
c-Rel	V\$CREL_01	[-0511:-0017]	175	076 (74.51%)	5.37e-11
RREB-1	V\$RREB1_01	[-4382:-3850]	246	089 (87.25%)	1.13e-10
TATA	V\$TATA_C	[-0060:+0042]	039	035 (34.31%)	3.61e-10
NF-AT	V\$NFAT_Q6	[-0251:-0016]	107	066 (64.71%)	4.34e-08
SRY	V\$SRY_01	[-9915:-9736]	110	059 (57.84%)	1.00e-07
CdxA	V\$CDXA_02	[-5849:-5523]	099	050 (49.02%)	2.23e-07

e-value < 1.0

TFM-Scan is able to delineate promoter areas that share relevant over-represented TF binding sites

Muscle data set

[Wasserman]

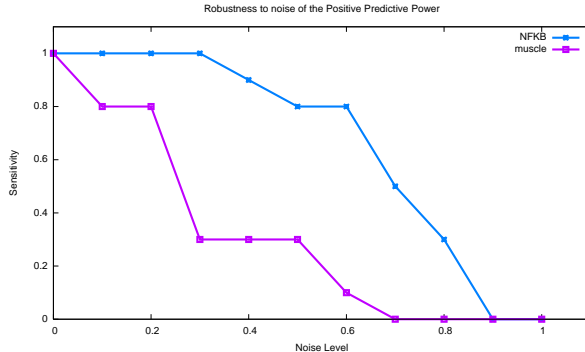
- ▶ Set of 27 genes that have skeletal muscle-specific expression (13 human, 17 mouse and 7 rat genes)
 - ▶ Promoter sequences retrieved from University California Santa Cruz Genome Browser (region $[-2000 + 200]$)
 - ▶ 5 factors that have muscle-specific expression are known: MyF, MEF-2, SRF, TEF (for validation)
- ▶ Empiric Background model
- ▶ All TRANSFAC matrices

Muscle data set (2) - Results

Rank	Transcription Factor (and PSSM)	Score
TFM-Explorer		
1	SRF* ([-0241:-0027])	6.459e-08
2	MyF* ([-0123:-0024])	9.965e-07
3	MEF2* ([-0073:-0026])	1.003e-05
4	p50 ([-0089:-0058])	3.637e-05
5	Hen-1 ([-0489:-0331])	1.498e-04
OTFBS		
1	MYOD.Q6*	3.076e-08
2	AP4.Q5	1.644e-07
3	TAL1BETAE47.01	1.468e-06
4	E47.01	3.599e-06
5	FOXJ2.01	5.798e-06
Toucan		
1	TGIF.01	3.002e-05 (2.29)
2	SRF.C*	7.748e-05 (1.878)
3	E47.02	2.194e-04 (1.426)
4	RFX1.02	2.462e-04 (1.386)
5	LMO2COM.01	3.232e-04 (1.258)
oPOSSUM		
1	MEF2*	1.663e-05
2	SRY*	4.190e-04
3	c-MYB.1	5.022e-04
4	S8	9.329e-04
5	SP1*	1.023e-03
6	Hen-1	1.034e-03

Rank	Transcription Factor (and PSSM)	Score
TFM-Explorer		
1	SRF* ([-0241:-0027])	6.459e-08
2	MyF* ([-0123:-0024])	9.965e-07
3	MEF2* ([-0073:-0026])	1.003e-05
4	p50 ([-0089:-0058])	3.637e-05
5	Ahr-ARNT ([-1342:-1268])	4.505e-05
OTFBS		
1	TAL1BETAITF2.01	3.244e-10
2	TAL1BETAE47.01	5.304e-09
3	YY1.02	1.506e-08
4	TAL1ALPHA47.01	7.534e-08
5	AP4.Q5	3.401e-07
6	MYOD.Q6*	7.808e-07
Toucan		
1	E47.02	5.414e-02 (-0.969)
2	MEF2.01*	1.128e-01 (-1.288)
3	TAL1ALPHA47.01	1.586e-01 (-1.436)
4	MEF2.02*	2.080e-01 (-1.554)
5	MEF2.03*	2.080e-01 (-1.554)
6	CEBP.C	2.196e-01 (-1.577)
oPOSSUM		
1	MEF2*	1.663e-05
2	SRY*	4.190e-04
3	c-MYB.1	5.022e-04
4	S8	9.329e-04
5	SP1*	1.023e-03
6	Hen-1	1.034e-03

Robustness to noise



Web interface

<http://bioinfo.lifl.fr>

Home

Results

Help

Sequences Name

Number of sequences102

Number of matrices243

Region-2000:+0200

DateWed Jan 11 15:16:14 2006

Minimum window size30bp

Maximum window size1000bp

Maximum number of windows to show20

Ratio2.50

Download:[plain text format](#)

Click on a line to get more detailed information

Select two lines and click on **pairwise comparison** to compute correlation between two predictions

Rank	Factor	Matrix ID	Location	Sequences	P-Value	pairwise comparison
1	NF-kappaB (p65)	V\$NFKAPPAB65_01	[-0521:-0019]	72 (70%)	2.46e-21	<input checked="" type="checkbox"/>
2	NF-kappaB	V\$NFKB_C	[-0537:-0020]	76 (74%)	7.51e-20	<input checked="" type="checkbox"/>
3	NF-kappaB	V\$NFKAPPAB_01	[-0536:-0018]	76 (74%)	1.67e-19	<input type="checkbox"/>
4	TATA	V\$TATA_01	[-0059:-0026]	38 (37%)	2.34e-19	<input type="checkbox"/>
5	c-Rel	V\$CREL_01	[-0501:-0020]	66 (64%)	1.64e-15	<input type="checkbox"/>
6	NF-kappaB	V\$NFKB_Q6	[-0537:-0021]	69 (67%)	1.87e-15	<input type="checkbox"/>
7	CdxA	V\$CDXA_01	[-0058:-0018]	38 (37%)	1.97e-15	<input type="checkbox"/>
8	TATA	V\$TATA_C	[-0060:-0015]	38 (37%)	2.46e-14	<input type="checkbox"/>
9	NF-kappaB (p50)	V\$NFKAPPAB50_01	[-0223:-0019]	41 (40%)	2.44e-09	<input type="checkbox"/>
10	NF-AT	V\$NFAT_Q6	[-0223:-0095]	54 (52%)	7.70e-09	<input type="checkbox"/>

Web interface (2)

Full results for matrix V\$NFKAPPAB65_01 and region [-0521:-0019]

Matrix information

Matrix name	V\$NFKAPPAB65_01
Transcription Factor name	NF-kappaB (p65)
Source	TRANSFAC OR JASPAR
Information Content	14.76 bits

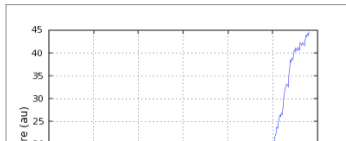
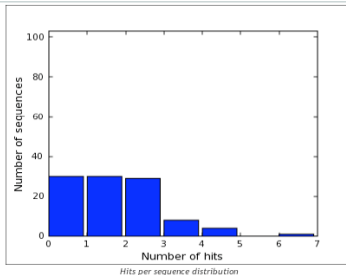


WebLogo representation of the matrix V\$NFKAPPAB65_01 (NF-kappaB (p65))

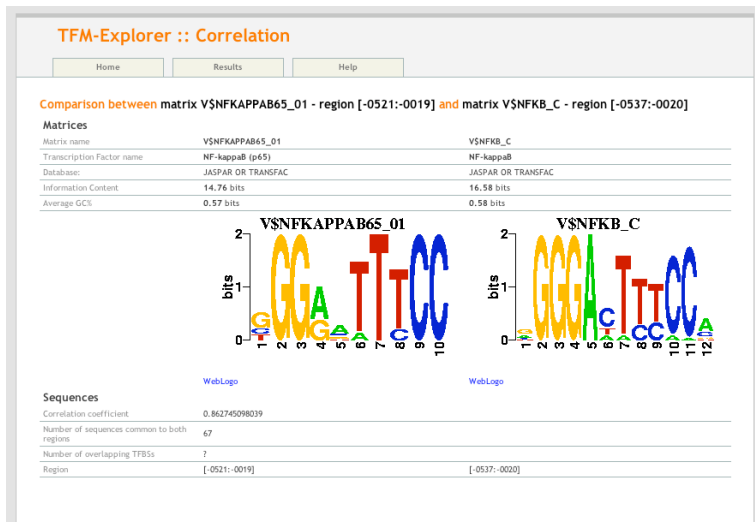
Web interface (3)

Sequences information

Total number of hits in region [-0521:-0019]	134
Total number of sequences	102 (2.61e-03 TFBS per sequences per bp)
Number of sequences having at least one putative TFBS in region [-0521:-0019]	72 (3.70e-03 per sequence per bp)
P-value	2.46e-21
E-value	2.81e-15



Web interface (4)



Conclusion

- ▶ Extract promoter areas that share relevant over-represented TF binding sites
 - ▶ No a priori knowledge of areas size or location is needed
 - ▶ Any kind of TF profile can be used
- ▶ Use regulatory motifs conservation across species
- ▶ Run on the fly