



TUTORIAL T01 ECCB 2014

Analysis of Cis-Regulatory Motifs from High-Throughput Sequence Sets

Session 1 : Web site

Goal and organisation of this session

Goal: introduction to ChIP-seq data analysis with RSAT

- **processing steps:** from reads to peaks. => not covered
- **downstream analyses:**
 - focus on motif analyses

Access

- **Web site**

Steps

1. Retrieving **sequences** from a set of peak coordinates (*fetch-sequences*)
2. Discovering **motifs** from peak sequences (*peak-motifs*)
3. **Visualizing** the sites in the context of genome annotations (*UCSC genome browser*)
4. Measure the **enrichment** of peaks for the expected motif (*matrix-quality*)
5. **A quick tour** of utility tools on RSAT

Morgane Thomas-Chollier

Biological concepts of transcriptional regulation

The diagram illustrates the process of transcriptional regulation. It shows a loop of chromatin where a distal transcription factor binding site (TFBS) interacts with a proximal TFBS through a co-activator complex. This leads to the formation of a transcription initiation complex and subsequent transcription initiation. A legend defines transcription factors as proteins that modulate target genes through DNA cis-regulatory elements.

Transcription factors are proteins that modulate (activate/repress) the expression of **target genes** through the binding on **DNA cis-regulatory elements**

Chromatin accessibility (open/closed) and histone modifications (eg: acetylation) also regulate gene expression

Wasserman et al, Nat Rev Genet, 2004

Morgane Thomas-Chollier

Gene "switched on"

- Active (open) chromatin
- Unmethylated cytosines (white circles)
- Acetylated histones

Gene "switched off"

- Silent (condensed) chromatin
- Methylated cytosines (red circles)
- Deacetylated histones

Transcription Factors / Co-activators

SWI/SNF, HAT, RNA Pol II, HDAC, HMT

Transcription possible

Transcription impeded

Wikipedia

in vivo experimental methods to identify binding sites

ChIP (=Chromatin Immuno-Precipitation)

=> differences in **methods to detect the bound DNA**

- small-scale: PCR / qPCR
- large-scale:
 - microarray = **ChIP-on-chip**
 - sequencing = **ChIP-seq**

Main challenge:

- quality/specificity of the antibodies

The flowchart details the ChIP protocol:

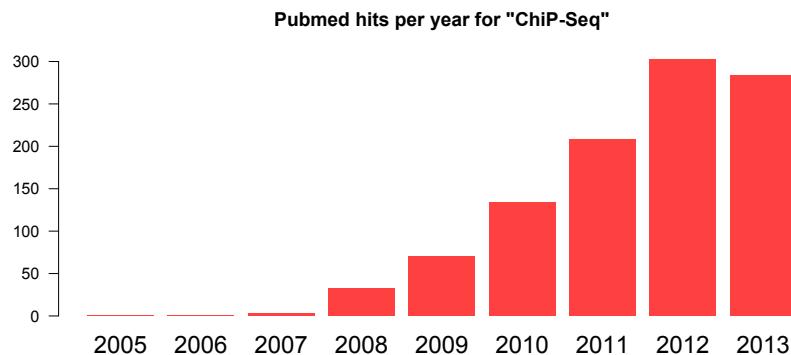
- DNA-protein cross-linking
- Cell lysis
- Sonication or enzyme digestion
- Fragmented chromatin
- Immunoprecipitation with specific antibody
- Immune precipitate (ChIP material)
- DNA purification
- Analysis of bound DNA

Methods for analysis include PCR, qPCR, Microarray, and Sequencing. A sequencing gel electrophoresis image is shown at the bottom right.

Morgane Thomas-Chollier

<http://www.chip-antibodies.com/>

ChIP-seq is a recently-adopted technique !

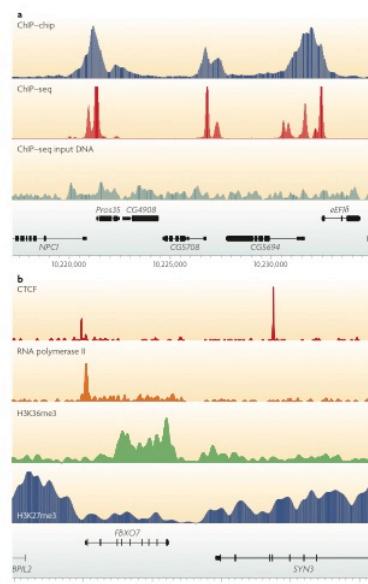


Morgane Thomas-Chollier

ChIP-seq applications

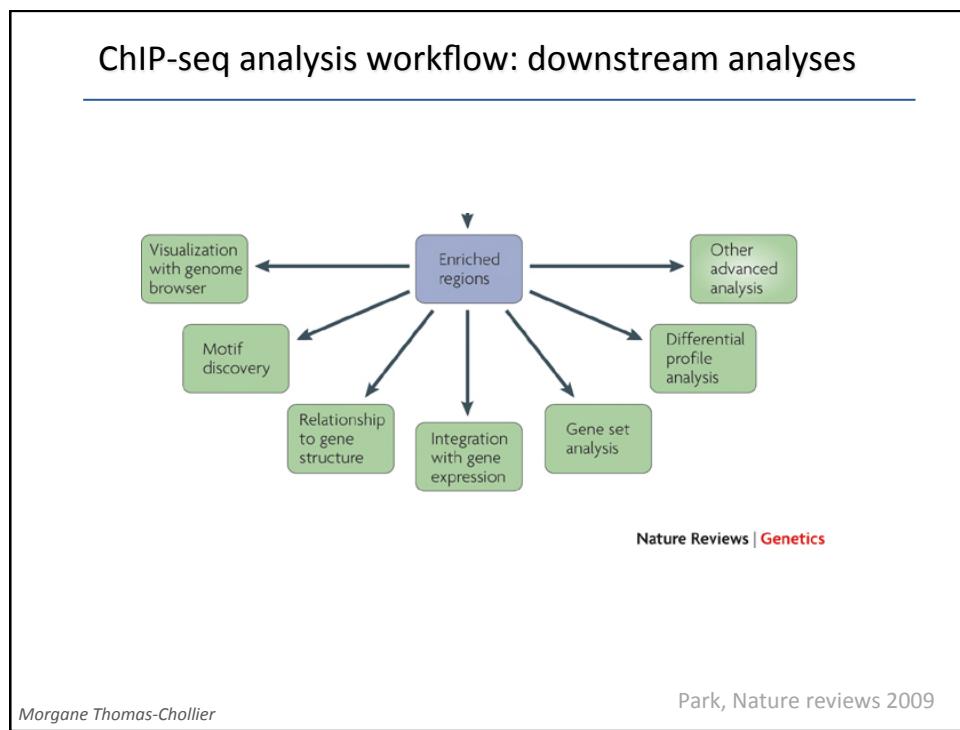
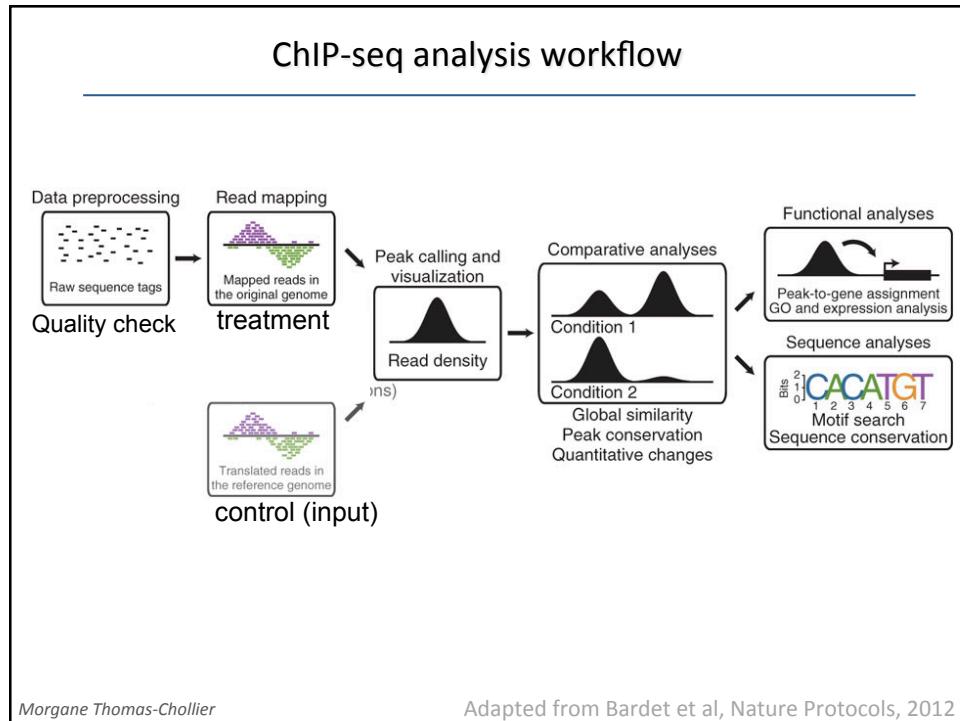
- find **all** regions in the genome bound by
 - a specific **transcription factor**
 - **histones** bearing a specific **modification**
- in a given **experimental condition** (cell type, developmental stage,...)

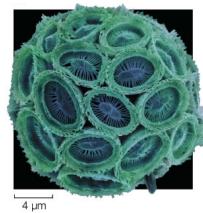
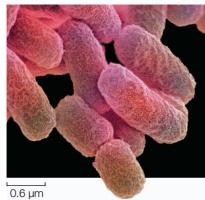
The obtain ChIP-seq **profiles** have **different shapes**, depending on the targeted protein



Park, Nature reviews 2009

Morgane Thomas-Chollier





What is the biological question ?



Morgane Thomas-Chollier

What is the biological question ?

« see if you can find something in the data »

Morgane Thomas-Chollier

What is the biological question ?

~~« see if you can find something in the data »~~

Morgane Thomas-Chollier

What is the biological question ?

- **Where** do a transcription factor (TF) bind ?
 - ✓ In a **specific context** (tissue, developmental stage, mutant)
 - ✓ By **comparison** to another context (WT vs mutant, different time points)

Morgane Thomas-Chollier

What is the biological question ?

- **Where** do a transcription factor (TF) bind ?
 - ✓ In a **specific context** (tissue, developmental stage, mutant)
 - ✓ By **comparison** to another context (WT vs mutant, different time points)

- **How** do a transcription factor (TF) bind ?
 - ✓ Which **binding motif(s)** (can be several for a given TF !!)
 - ✓ Is the **binding** direct to DNA or via **protein-protein** interactions ?
 - ✓ Are there **cofactors** (maybe affecting the motif !!), and if so, identify them

Morgane Thomas-Chollier

What is the biological question ?

- **Where** do a transcription factor (TF) bind ?
 - ✓ In a **specific context** (tissue, developmental stage, mutant)
 - ✓ By **comparison** to another context (WT vs mutant, different time points)

- **How** do a transcription factor (TF) bind ?
 - ✓ Which **binding motif(s)** (can be several for a given TF !!)
 - ✓ Is the **binding** direct to DNA or via **protein-protein** interactions ?
 - ✓ Are there **cofactors** (maybe affecting the motif !!), and if so, identify them

- Which **regulated genes** are directly regulated by a given TF ?

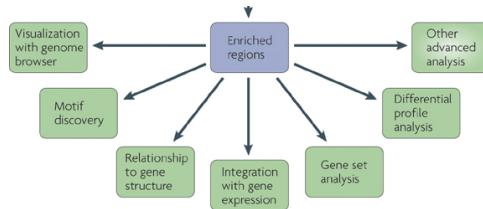
- What are the **targets** of a given TF ?

- Where are the **promoters** (PolII) and **chromatin marks** ?

Morgane Thomas-Chollier

What is the biological question ?

→ Should drive all « downstream » analyses



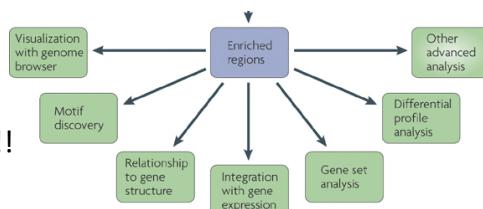
Morgane Thomas-Chollier

Nature Reviews | Genetics

What is the biological question ?

→ Should drive all « downstream » analyses

Will take time
to « do it all » !!!



Morgane Thomas-Chollier

Nature Reviews | Genetics

What is the biological question ?
What can be the following experimental work ?

Morgane Thomas-Chollier

What is the biological question ?
What can be the following experimental work ?

- ➔ cell biology (eg: luciferase assay) ?
- ➔ in vitro assays (eg: EMSA) ?
- ➔ Proteomic (eg: mass spectrometry) ?
- ➔ Transgenics ?
- ➔ Will depend on
 - ✓ the organism
 - ✓ available infrastructure

Morgane Thomas-Chollier

What is the biological question ?

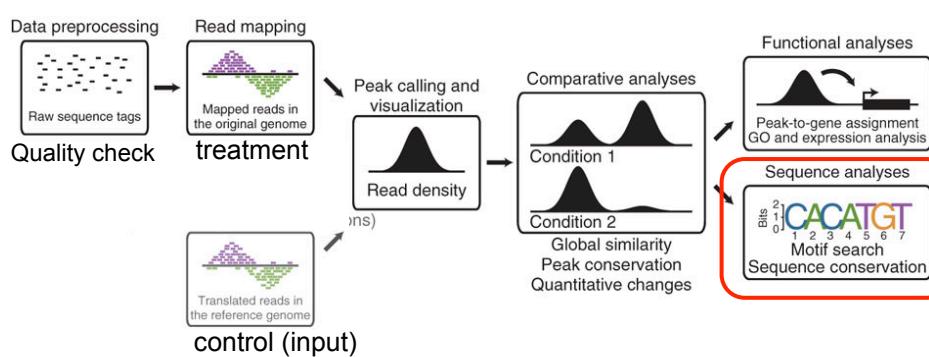
- **Where** do a transcription factor (TF) bind ?
 - ✓ In a **specific context** (tissue, developmental stage, mutant)
 - ✓ By **comparison** to another context (WT vs mutant, different time points)

- **How** do a transcription factor (TF) bind ?
 - ✓ Which **binding motif(s)** (can be several for a given TF !!)
 - ✓ Is the **binding** direct to DNA or via **protein-protein** interactions ?
 - ✓ Are there **cofactors** (maybe affecting the motif !!), and if so, identify them

- Which **regulated genes** are directly regulated by a given TF ?
- What are the **targets** of a given TF ?
- Where are the **promoters** (PolII) and **chromatin marks** ?

Morgane Thomas-Chollier

ChIP-seq analysis workflow



Morgane Thomas-Chollier

Adapted from Bardet et al, Nature Protocols, 2012

Transcription factor specificity

How do TF « know » where to bind DNA ?

TF recognize TFBS with specific DNA sequences

Morgane Thomas-Chollier

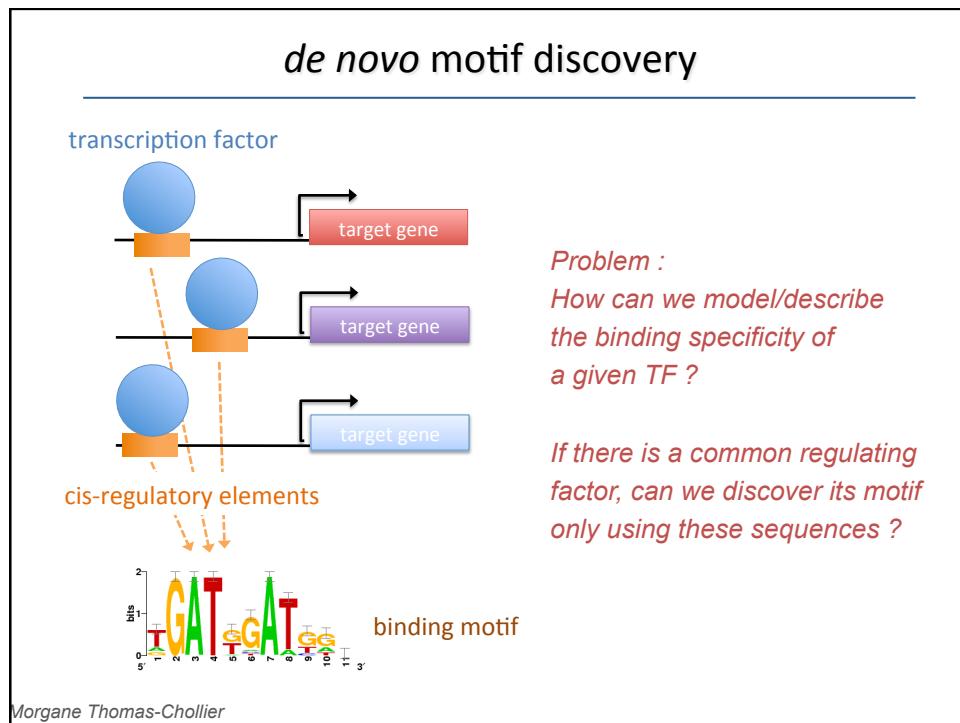
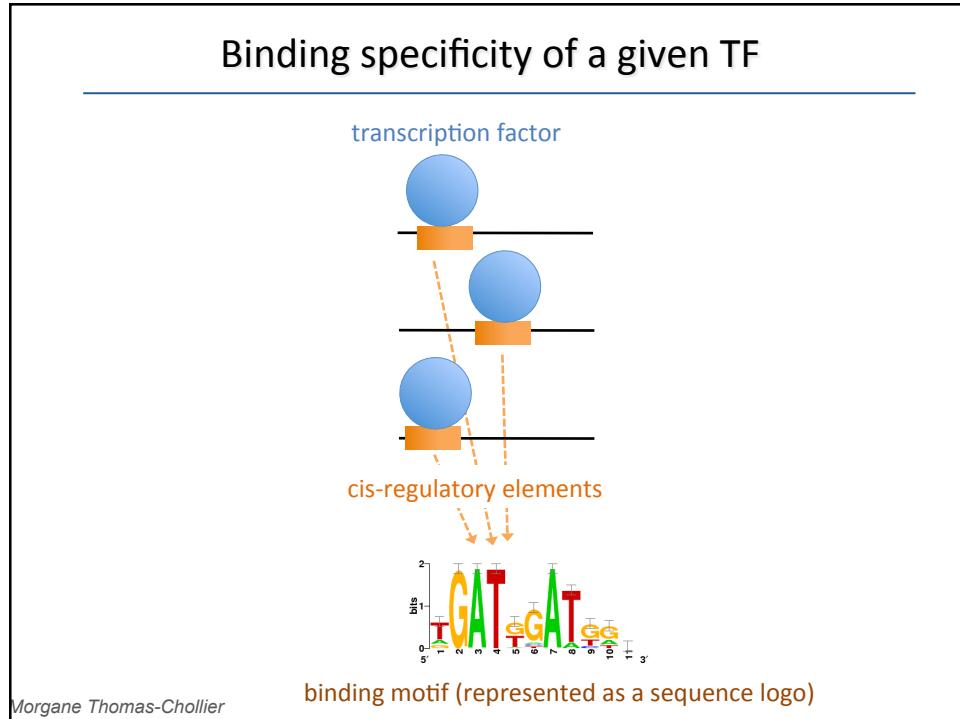
Transcription factor specificity

How do TF « know » where to bind DNA ?

TF recognize TFBS with specific DNA sequences

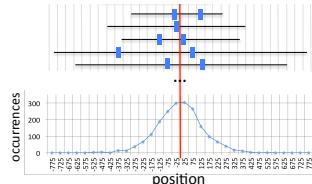
TFBSs are *degenerate*:
a given TF is able to bind DNA on TFBSs with different sequences

Morgane Thomas-Chollier



de novo motif discovery

- Find exceptional motifs based on the sequence only
(*A priori* no knowledge of the motif to look for)
- Criteria of exceptionality:
 - higher/lower frequency than expected by chance
(over-/under-representation)
 - concentration at specific positions relative to some reference coordinate
(positional bias)



Morgane Thomas-Chollier

de novo motif discovery

- Tools already exist for a long time !
 - MEME (1994)
 - RSAT oligo-analysis (1998)
 - AlignACE (2000)
 - Weeder (2001)
 - MotifSampler (2001)

Morgane Thomas-Chollier

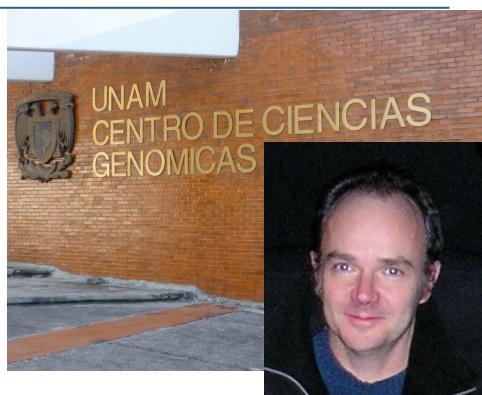
de novo motif discovery

- Tools already exist for a long time !
 - MEME (1994)
 - **RSAT oligo-analysis (1998)**
 - AlignACE (2000)
 - Weeder (2001)
 - MotifSampler (2001)

Morgane Thomas-Chollier

Regulatory Sequence Analysis Tools (**RSAT**)

- Since 1998 (15 years !)
- Initiated in Cuernavaca, Mexico
- yeast cis-regulatory elements



Jacques van Helden



Morgane Thomas-Chollier

RSAT improvements over the years

The screenshot shows two versions of the RSAT website side-by-side. On the left is the version from 2000, featuring a dark blue header with 'RSAT NeAT' and a sidebar with 'Most popular tools' like 'Matrix-wise (quick)', 'Matrix-wise (detailed)', and 'Genomes and genes'. The main content area has sections for 'Tool Map', 'Introduction', 'Tutorials', 'Publications', 'Credits', 'People', 'Data', and 'Download'. On the right is the version from 2012, which has a more modern design with a orange and blue logo, a 'Warnings' box, and a 'Vertebrate genomes' link. It also includes a section for 'Regulatory Sequence Analysis Tools - Web servers' listing servers in Brussels, Uppsala, Paris, and Commissariat à l'Energie Atomique.

Thomas-Chollier, Darbo, Herrmann, Defrance , Thieffry, van Helden *Nature Protocols*, 2012
Thomas-Chollier Defrance, Medina-Rivera, Sand, Herrmann, Thieffry, van Helden *Nucleic Acids Research*, 2012
Medina-Rivera, Abreu-Goodger, Thomas-Chollier, Salgado, Collado-Vides, van Helden *Nucleic Acids Research*, 2011
Sand, Thomas-Chollier, van Helden *Bioinformatics*, 2009
Thomas-Chollier*, Sand*, Turatsinze, Janky, Defrance, Vervisch, van Helden *Nucleic Acids Research*, 2008
Sand, Thomas-Chollier, Vervisch, van Helden *Nature Protocols*, 2008
Thomas-Chollier*, Turatsinze*, Defrance, van Helden *Nature Protocols*, 2008
van Helden, *Nucleic Acids Research*, 2003
van Helden, André, Collado-Vides *Yeast*, 2000

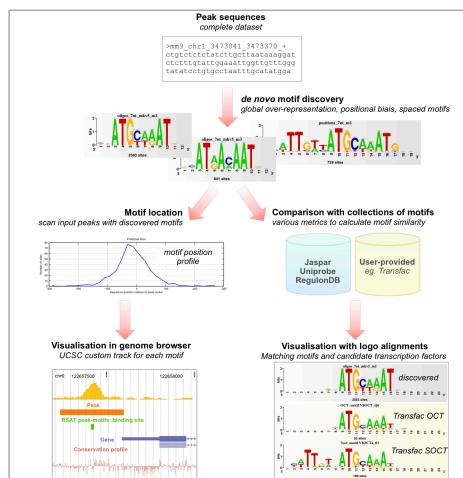
New approaches for ChIP-seq datasets

- Size, size, size**
 - limited numbers of promoters and enhancers
 - ↓
 - dozens of thousands of peaks !!!!!!
- the problem is slightly different**
 - promoters: 200-2000bp from co-regulated genes
 - ↓
 - peaks: 300bp, positional bias
- motif analysis: not just for specialists anymore !**
 - complete user-friendly workflows

Morgane Thomas-Chollier <http://www.genomequest.com/landing-pages/ODI-webinar-web.html>

New approaches for ChIP-seq datasets

- *de novo* motif discovery (*peak-motifs* in RSAT)



Thomas-Chollier et al Nucleic Acids Research, 2012

Morgane Thomas-Chollier

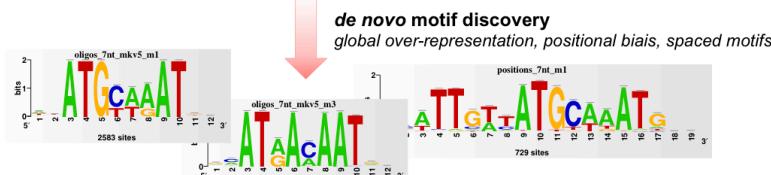
Peaks coordinates

BED

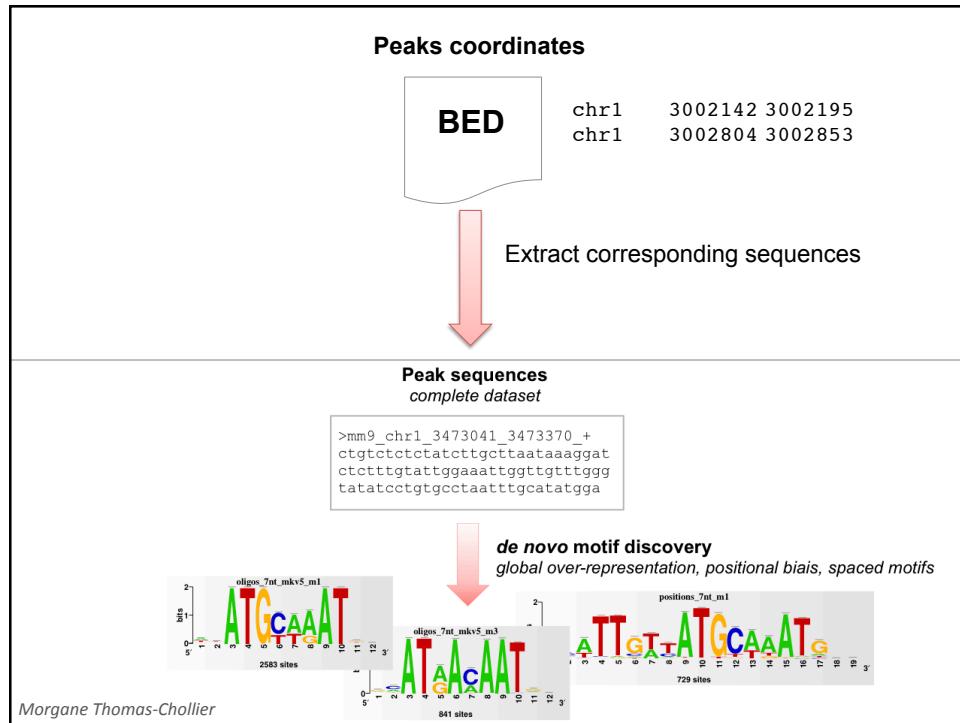
```
chr1    3002142 3002195  
chr1    3002804 3002853
```

Peak sequences complete dataset

```
>mm9_chrl_3473041_3473370_+
ctgtctcttatcttgcttaataaaggat
ctctttgtattggaaattgggtttggg
tatatccgtgcctaattgcataatgg
```



Morgane Thomas-Chollier



Hands on !

- Go to the companion website
- Follow all steps of **Retrieving sequences from your peaks**

Morgane Thomas-Chollier

Using RSAT

The screenshot shows the RSAT homepage. On the left is a sidebar with a tree menu:

- RSAT [NeAT]**
- Regulatory Sequence Analysis Tools**
- Most popular tools**
 - retrieve sequence
 - oligo-analysis (words)
 - matrix-scan (matrices)
 - random sequence
- > view all tools
- New!**
 - Genomes and genes
 - Sequence retrieval
 - Pattern discovery
 - Pattern matching
 - Comparative genomics
 - Conversion/Utilities
 - Drawing
 - Web services
- Help**
 - Map of the tools
 - Introduction
 - Tutorials
 - Course
 - Contact & Forum New!
- Information**
 - Feedback
 - Jacques van Helden

To the right is a dropdown menu titled "Sequence retrieval" with the following options:

- retrieve sequence
- retrieve Ensembl sequence New!
- purge sequence
- convert sequence
- random sequences

Annotations with arrows point to specific elements:

- An arrow points from the "Sequence retrieval" menu to the "2. Run the analysis" link.
- An arrow points from the "Sequence retrieval" menu to the "3. Visualization" link.
- Two orange arrows point to the "Help: tutorials, forum" and "Information: publications,..." links.

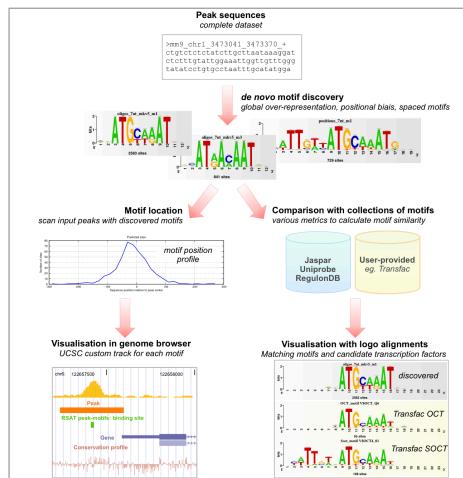
RSAT Web forms

The screenshot shows the "RSA-tools - retrieve sequence" web form. Annotations point to various parts:

- A blue bar at the top contains the text "RSA-tools - retrieve sequence" and "Tool name".
- The text "Returns upstream, downstream or ORF sequences for a list of genes" is labeled "Tool description".
- A note below the input fields reads: "Remark: If you want to retrieve sequences from an organism that is in the Ensembl database, we recommend to use the [retrieve-ensembl-seq](#) program instead."
- The "Organism" dropdown is set to "Saccharomyces cerevisiae".
- The "Single organism" radio button is selected.
- The "Multiple organisms" radio button is available but not selected.
- The "Genes" section includes a "selection" checkbox and a text area for entering gene IDs. Annotations point to this area as "Tool parameters".
- Below the genes input is a "Upload gene list from file" input field with a "Browse..." button. A checkbox "Query contains only IDs (no synonyms)" is also present.
- Tool parameters include "Feature type" (radio buttons for CDS, mRNA, tRNA, rRNA, scRNA), "Sequence type" (radio buttons for upstream, downstream, ORF), and checkboxes for "Prevent overlap with neighbour genes (noorf)", "Mask repeats (only valid for organisms with annotated repeats)", and "Admit imprecise positions".
- The "Sequence format" dropdown is set to "fasta".
- The "Sequence label" dropdown is set to "gene name".
- The "Output" section includes radio buttons for "server", "display", and "email". Annotations point to this section as "Output".
- At the bottom are buttons for "GO", "Reset", "DEMO", and "MANUAL TUTORIAL". The "DEMO" button is highlighted with an orange arrow. Annotations point to the "DEMO" button as "Demo button (fill in the form for test purposes)".
- At the very bottom is a "Help" link.

New approaches for ChIP-seq datasets

- *de novo* motif discovery (**peak-motifs** in RSAT)

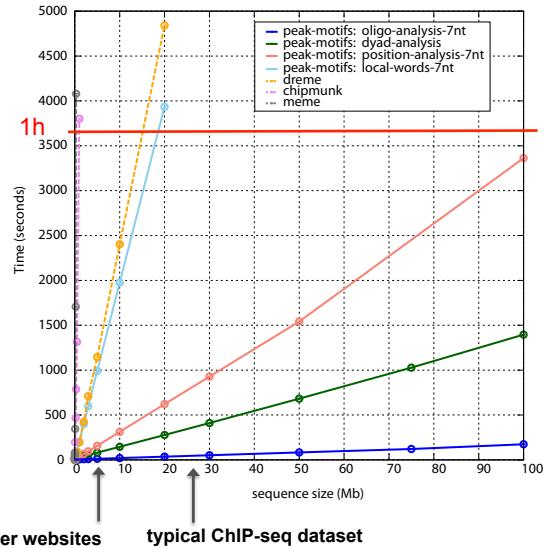


Thomas-Chollier et al Nucleic Acids Research, 2012

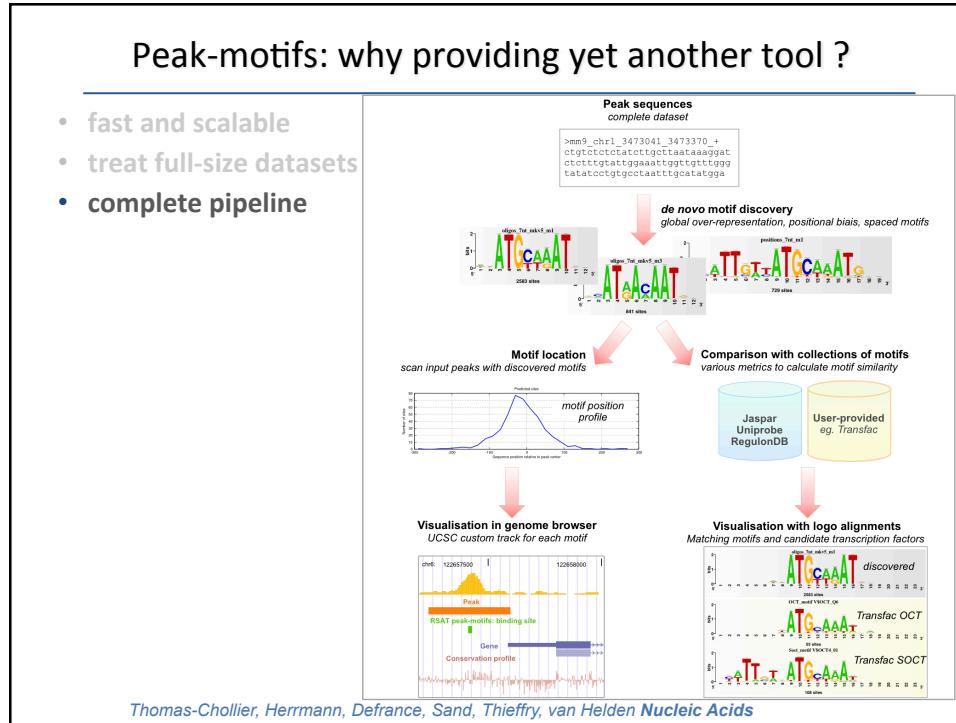
Morgane Thomas-Chollier

Peak-motifs: why providing yet another tool ?

- fast and scalable
 - treat full-size datasets



Thomas-Chollier, Herrmann, Defrance, Sand, Thieffry, van Helden *Nucleic Acids Research*, 2012



Peak-motifs: why providing yet another tool ?

- fast and scalable
- treat full-size datasets
- complete pipeline
- web interface
- accessible to non-specialists
 - Demo buttons
 - Tutorials & Protocols

Thomas-Chollier, Darbo, Herrmann, Defrance, Thieffry, van Helden Nature Protocols, 2012

- HTML report

Morgane Thomas-Chollier

The screenshot shows the 'Peak-motifs' web interface. At the top, it says 'RSA-tools • peak-motifs' and 'Peak-motifs for discovering motifs in massive ChIP-seq peak datasets'. Below this are sections for 'Peak sequences' (with a file input field for 'Peak sequences' and 'Optional: control dataset for differential analysis (test vs control)'), 'Control sequences' (with a file input field for 'Control sequences'), and 'Output' (radio buttons for 'display' or 'email'). At the bottom right are links for 'MANUAL', 'TUTORIAL', and 'ASK A QUESTION'.

Hands on !

- Go to the companion website
- Follow **step 1 of Discovering motifs from peak sequences**

Morgane Thomas-Chollier

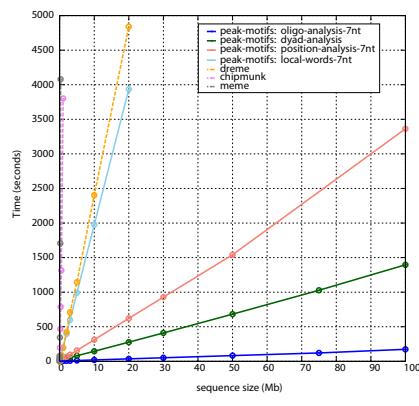
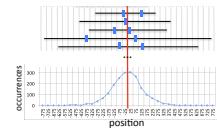
Peak-motifs: why providing yet another tool ?

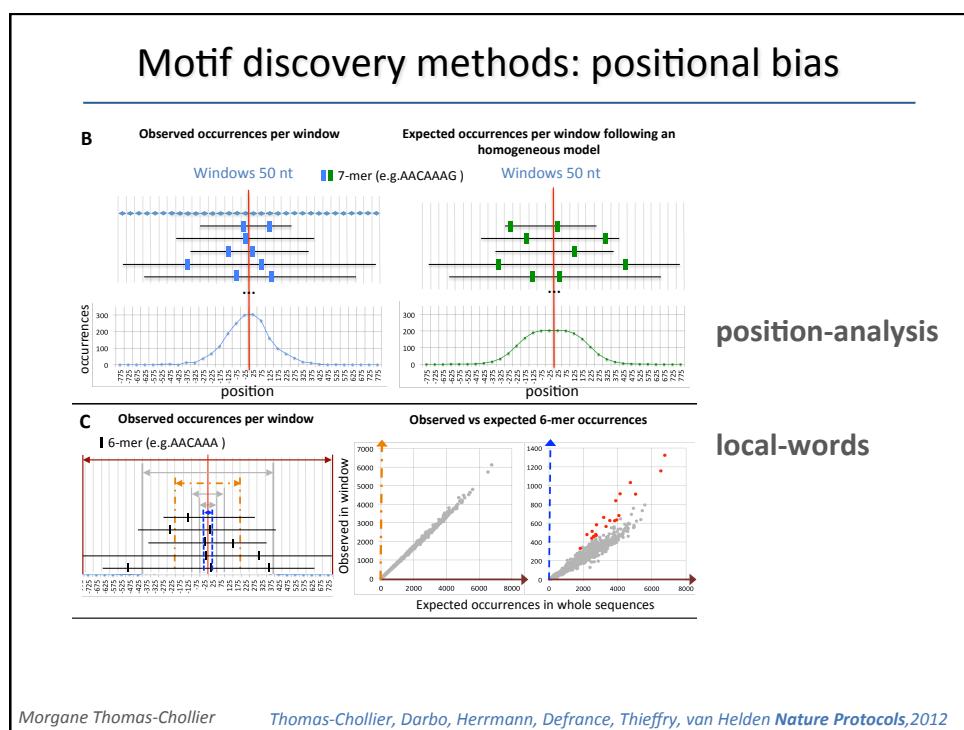
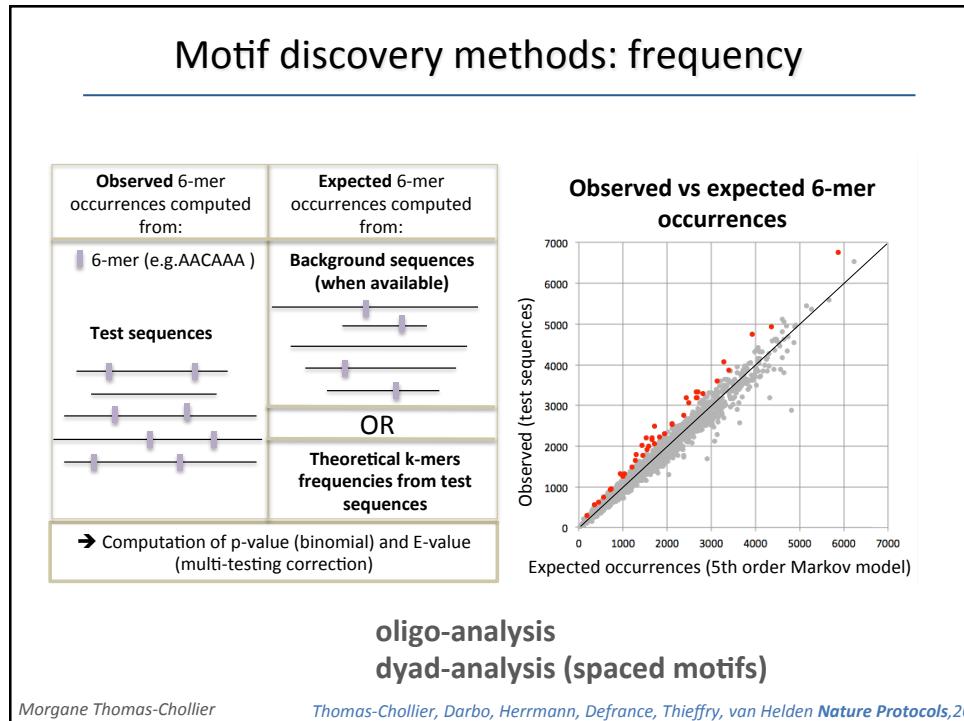
- fast and scalable
- treat full-size datasets
- complete pipeline
- web interface
- accessible to non-specialists
- using 4 complementary algorithms

- Global over-representation
 - **oligo-analysis**
 - **dyad-analysis (spaced motifs)**

- Positional bias

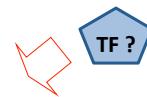
- **position-analysis**
- **local-words**





A motif discovery problem

Motif discovery



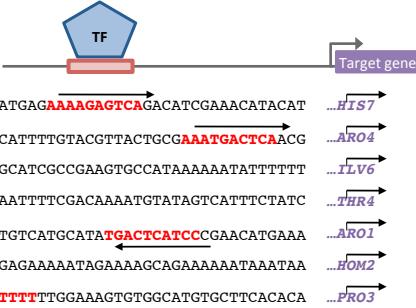
5' - TCTCTCTCACGGCTAATTAGGTGATCATGAAAAAATTGAGAAAAGAGTCAGACATCGAAACATACAT ...*HIS7*
 5' - ATGGCAGAACATCACTTTAACCGTGGCCCAACCGCTGCACCCCTGTGCATTTGACGTTACTGCCAAATGACTCAACG ...*ARO4*
 5' - CACATCCAACGAATCACCTCACCGTTATCGTACACTCACTTCTTCATGCCGAAGTGCATAAAAATTTTTT ...*ILV6*
 5' - TGCGAACAAAAGACTTACAACGAGGAAATAGAAGAAAATGAAAAAATTGACAAAATGATAGTCATTTCTATC ...*THR4*
 5' - ACAAAAGTACCTTCCTGGCCAATCTCACAGATTAATAGTAAATTGTCATGCCATATGACTCATCCGAACATGAAA ...*ARO1*
 5' - ATTGATTGACTCACTTCTGACTACTACCAGTTAACATGTTAGAGAAAATAGAAGACAGAAAAAATAATAAA ...*HOM2*
 5' - GGCACAGTCGGCTTGGTTATCCGGTACTCATCTGACTCTTTGGAAAGTGTGGCATGTGCTTCACACA ...*PRO3*

Problem : If there is a common regulating factor, can we discover its motif (some signal) on the basis of these sequences ONLY ?

- We have a set of sequences
- We suspect that they share some functional signal
- We ignore the transcription factors involved in this regulation.
- We ignore the cis-acting elements

Morgane Thomas-Chollier

Principle: detect unexpected patterns



- Binding sites are represented as “words” = “string”=“k-mer”
 - e.g. **acgtga** is a 6-mer
- Signal is likely to be **more frequent** in the upstream regions of the co-regulated genes than in a random selection of genes
- We will thus detect **over-represented words**

Motif discovery using word counting

Idea:

motifs corresponding to binding sites are generally repeated in the dataset
 → capture this statistical signal

Algorithm

- count occurrences of **all k-mers** in a set of related sequences
 (promoters of co-expressed genes, in ChIP bound regions,...)

Let's take an example (yeast *Saccharomyces cerevisiae*)

- NIT
 - 7 genes expressed under low nitrogen conditions
- MET
 - 10 genes expressed in absence of methionine
- PHO
 - 5 genes expressed under phosphate stress



PHO	MET	NIT
aaaaaa ttttt 51	aaaaaa ttttt 105	aaaaaa ttttt 80
aaaaag ctttt 15	atata atatat 41	cttatc gataag 26
aagaaa tttc 14	gaaaaa ttttc 40	tatata tatata 22
aaaaaa ttttc 13	tatata tatata 40	ataaga tcttat 20
tgccaa ttggca 12	aaaaat attttt 35	aagaaa tttc 20
aaaaat attttt 12	aagaaa tttc 29	aaaaaa ttttc 19
aaatta taattt 12	agaaaa ttttct 28	atatat atatat 19
agaaaa tttct 11	aaaata tatttt 26	agataa ttatct 17
caagaa ttcttg 11	aaaaag cttttt 25	agaaaa tttct 17
aaacgt acgttt 11	agaaaat atttc 24	aaagaa ttcttt 16
aaagaa ttcttt 11	aaataa tttattt 22	aaaaca tgtttt 16
acgtgc gcacgt 10	taaaaa ttttta 21	aaaaag cttttt 15
aataat attatt 10	tgaaaa tttca 21	agaaga tcttct 14
aagaag cttctt 10	ataata tattat 20	tgataa ttatca 14
atataa ttatat 10	atataa tttata 20	atataa ttatat 14

The most frequent oligonucleotides are not informative

- A (too) simple approach would consist in **detecting the most frequent oligonucleotides** (for example hexanucleotides) for each group of upstream sequences.
- This would however lead to deceiving results.
 - In all the sequence sets, the same kind of patterns are selected: **AT-rich hexanucleotides**.

PHO	MET	NIT
aaaaaa tttttt 51	aaaaaa tttttt 105	aaaaaa tttttt 80
aaaaag cttttt 15	atatat atatat 41	cttata gataag 26
aagaaa tttctt 14	gaaaaa tttttc 40	tatata tatata 22
gaaaaa tttttc 13	tatata tatata 40	ataaga tcttat 20
tgccaa ttggca 12	aaaaat attttt 35	aagaaa tttctt 20
aaaaat attttt 12	aagaaa tttctt 29	gaaaaa tttttc 19
aaatta taattt 12	agaaaa ttttct 28	atatat atatat 19
agaaaa ttttct 11	aaaata tatttt 26	agataa ttatct 17
caagaa ttcttg 11	aaaaag cttttt 25	agaaaa ttttct 17
aaacgt acgttt 11	agaaaat atttct 24	aaagaa ttcttt 16
aaagaa ttcttt 11	aaataa tttattt 22	aaaaca tggttt 16
acgtgc gcacgt 10	taaaaa ttttta 21	aaaaag cttttt 15
aataat attatt 10	tgaaaa ttttca 21	agaaga tctttt 14
aagaag cttctt 10	ataata tattat 20	tgataa ttatca 14
atataa ttatat 10	atataa ttatat 20	atataa ttatat 14

A more relevant criterion for over-representation

- The most frequent patterns do not reveal the motifs specifically bound by specific transcription factors.
- They merely **reflect the compositional biases** of upstream sequences.
- A more relevant criterion for over-representation is to detect patterns which **are more frequent** in the upstream sequences of the selected genes (co-regulated) **than the random expectation**.
- The **random expectation** is calculated by counting the frequency of each pattern in the complete set of upstream sequences (all genes of the genome).
 - => **“Background”**

Motif discovery using word counting

Idea:

motifs corresponding to binding sites are generally repeated in the dataset
 → capture this statistical signal

Algorithm

- count occurrences of **all k-mers** in a set of related sequences (promoters of co-expressed genes, in ChIP bound regions,...)
- estimate the **expected number of occurrences** from a background model
 - empirical based on observed k-mer frequencies
 - theoretical background model (Markov Models)

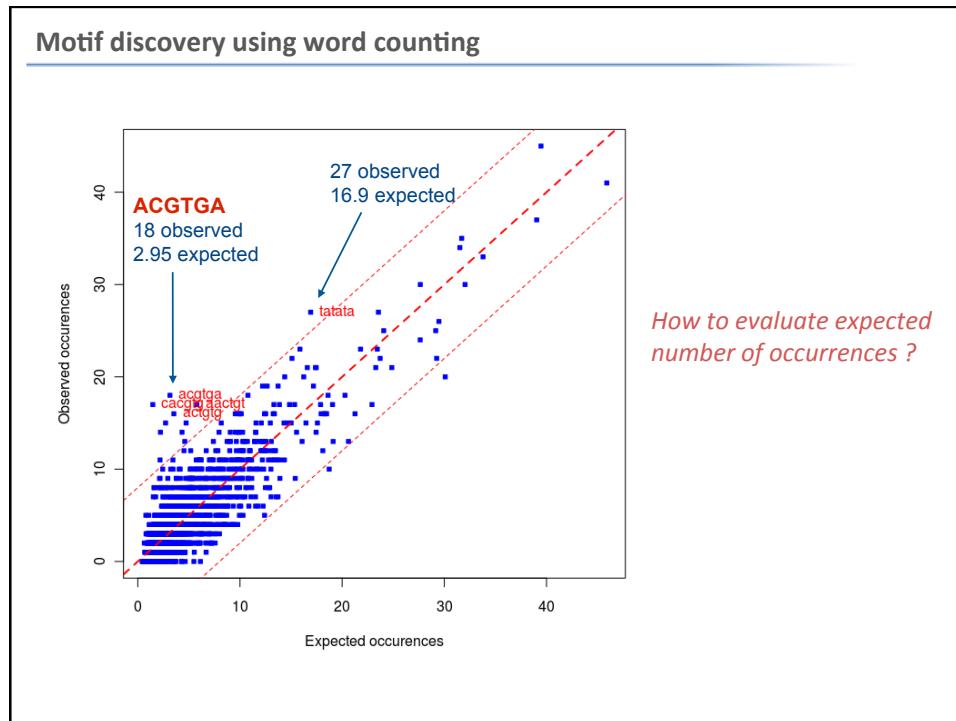
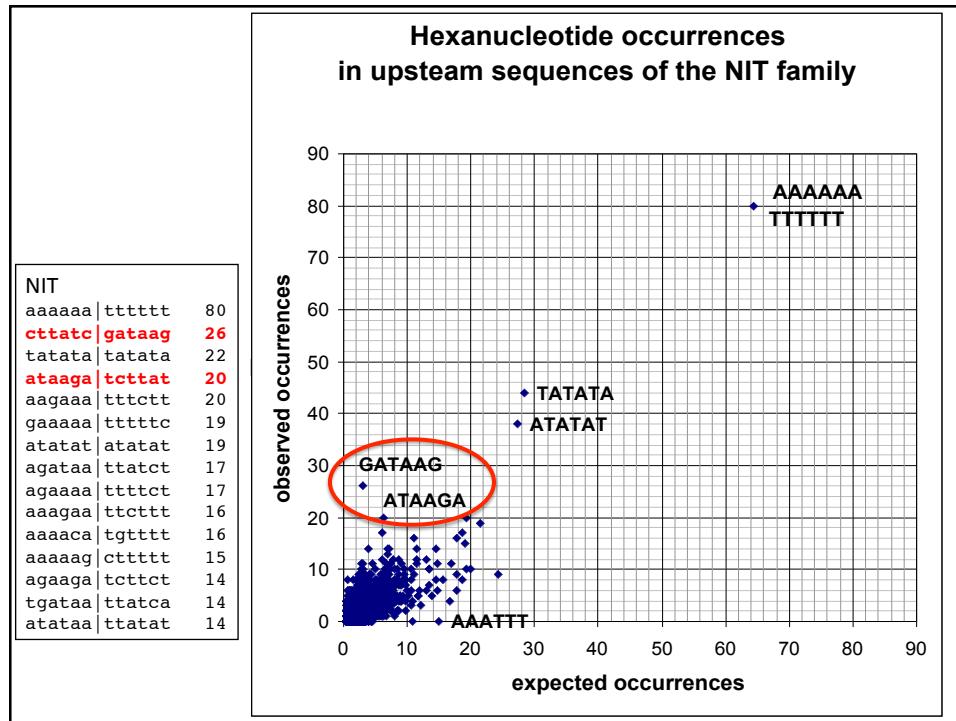
Estimation of word expected frequencies from background sequences



Example:

6nt frequencies in the whole set of 6000 yeast **upstream** sequences

:seq	identifier	observed_freq occ	
aaaaaa	aaaaaa ttttt	0,00510699	14555
aaaaac	aaaaac gtttt	0,00207402	5911
aaaaag	aaaaag ctttt	0,00375191	10693
aaaaat	aaaaat atttt	0,00423577	12072
aaaaca	aaaaca tgttt	0,0019828	5651
aaaacc	aaaacc ggttt	0,00088526	2523
aaaacg	aaaacg cgttt	0,00090105	2568
aaaact	aaaact agttt	0,0014621	4167
aaaaga	aaaaga tcctt	0,00323016	9206
aaaagc	aaaagc gcctt	0,00135824	3871
aaaagg	aaaagg ccctt	0,0017849	5087
aaaagt	aaaagt acttt	0,0019035	5425
aaaaata	aaaata tattt	0,00336805	9599
aaaaatc	aaaatc gattt	0,00131368	3744
aaaaatg	aaaatg cattt	0,00185648	5291
aaaatt	aaaatt aattt	0,00269156	7671
aaacaa	aaacaa tttgt	0,00209999	5985
aaacac	aaacac gttgt	0,00071684	2043
aaacag	aaacag ctgtt	0,00096491	2750
aaacat	aaacat atgtt	0,00108982	3106
aaacca	aaacca tggtt	0,00074421	2121



Motif discovery using word counting

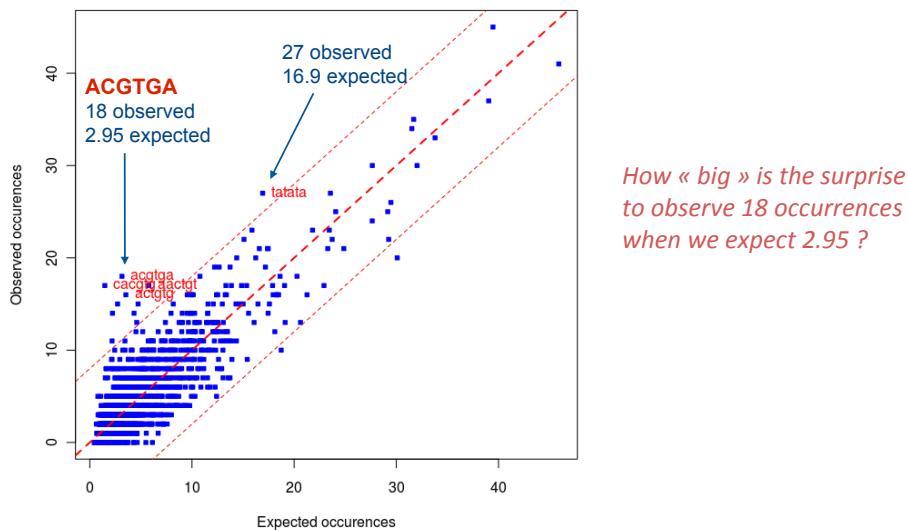
Idea:

motifs corresponding to binding sites are generally repeated in the dataset
 → capture this statistical signal

Algorithm

- count occurrences of **all k-mers** in a set of related sequences (promoters of co-expressed genes, in ChIP bound regions,...)
- estimate the **expected number of occurrences** from a background model
 - empirical based on observed k-mer frequencies
 - theoretical background model (Markov Models)
- **statistical evaluation of the deviation observed** (P-value/E-value)

Statistical evaluation



Statistical evaluation

How « big » is the surprise to observe 18 occurrences when we expect 2.95 ?

- at each position in the sequence, there is a **probability p** that the word starting at this position is ACGTGA
- we consider n positions
- what is the probability that k of these n positions correspond to ACGTGA ?
- **Application :** $p = 3.4e-4$ (intergenic frequencies)
 $n = 9000$ position
 $x = 18$ observed occurrences

$$P(X \geq x) = \sum_{i=x}^n \frac{n!}{i!(n-i)!} p^i (1-p)^{n-i}$$

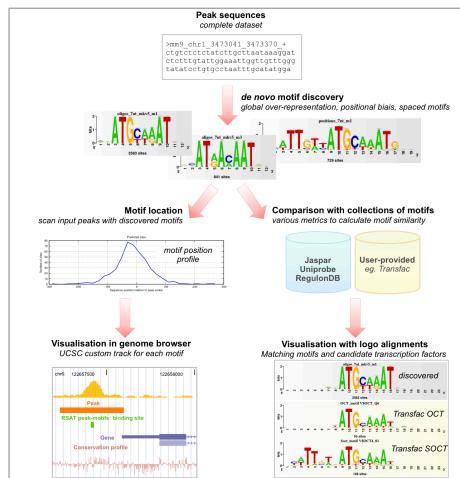
Binomial distribution to measure the “surprise”

Hands on !

- Go to the companion website
- Follow **step 2 of Discovering motifs from peak sequences**

New approaches for ChIP-seq datasets

- *de novo motif discovery* (**peak-motifs** in RSAT)

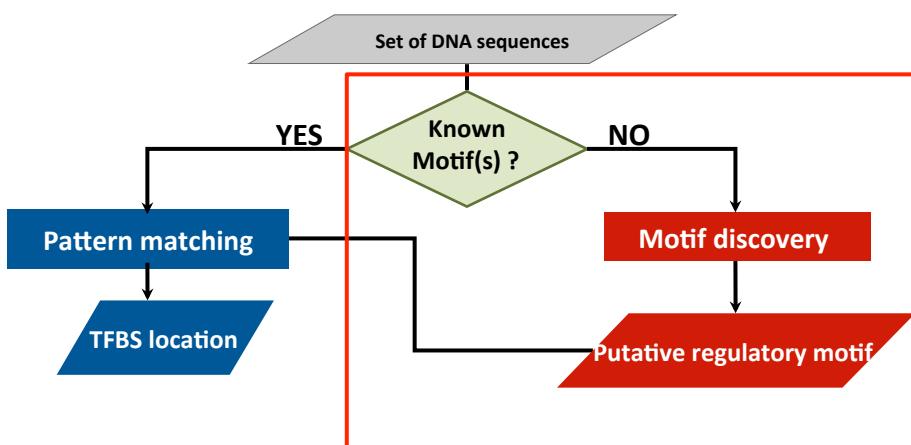


Thomas-Chollier et al Nucleic Acids Research, 2012

Morgane Thomas-Chollier

Motif scanning vs. Motif discovery

eg: ChIP-seq peak sequences



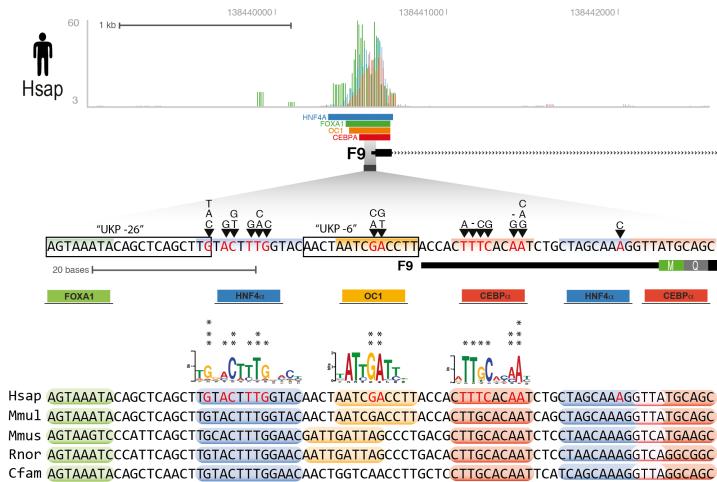
Morgane Thomas-Chollier

Hands on !

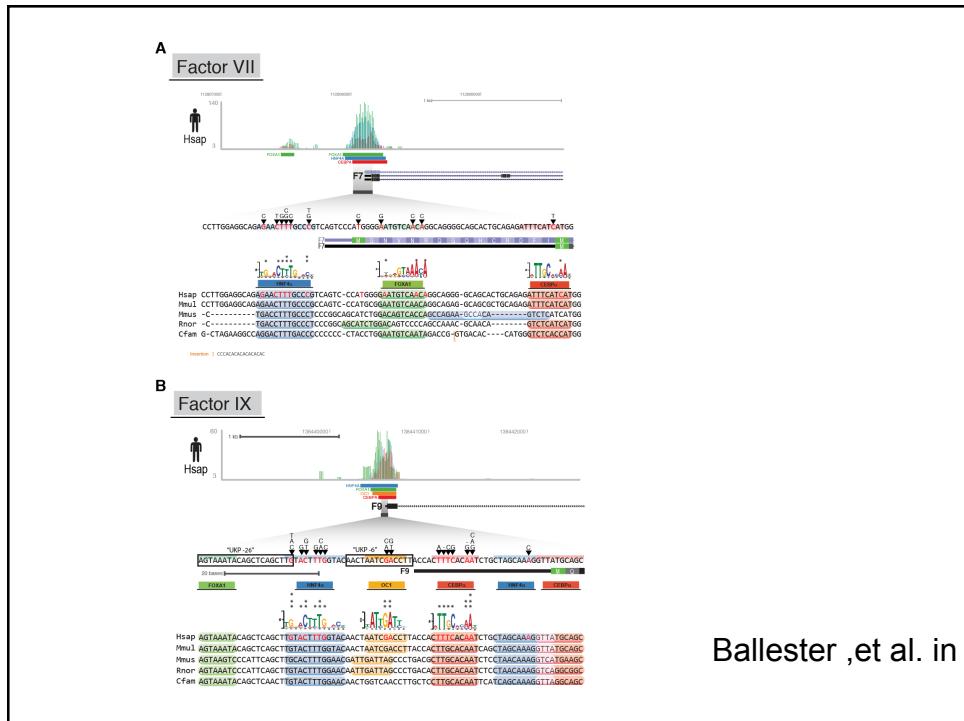
- Go to the companion website
- Follow all steps of **Visualizing the sites in the context of genome annotations**

Morgane Thomas-Chollier

Liver related genes have suffered rare mutations in their regulatory regions leading to human disease



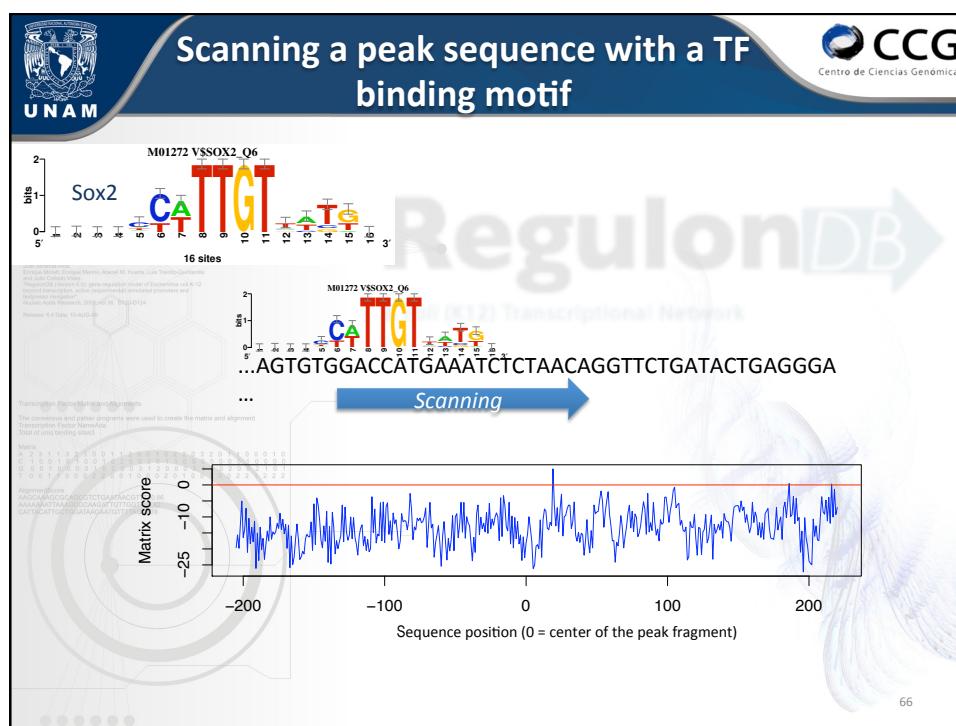
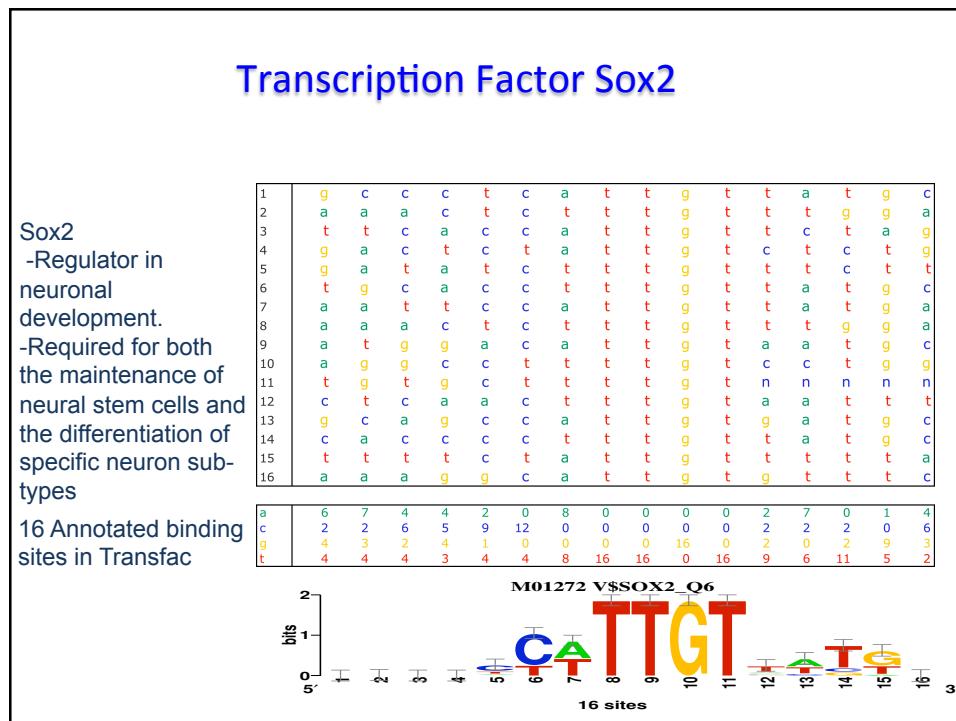
A CpG Mutational Hotspot in a ONECUT Binding Site Accounts for the Prevalent Variant of Hemophilia B
Leyden. Funnell AP, Wilson MD, Ballester B, et al. Am J Hum Genet. 2013

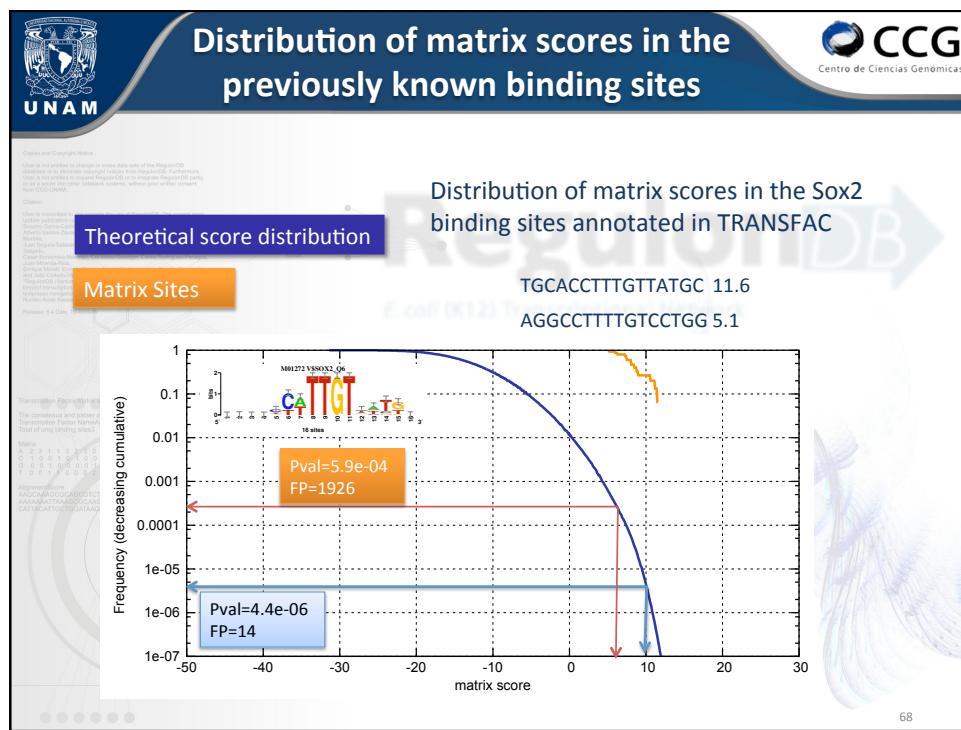
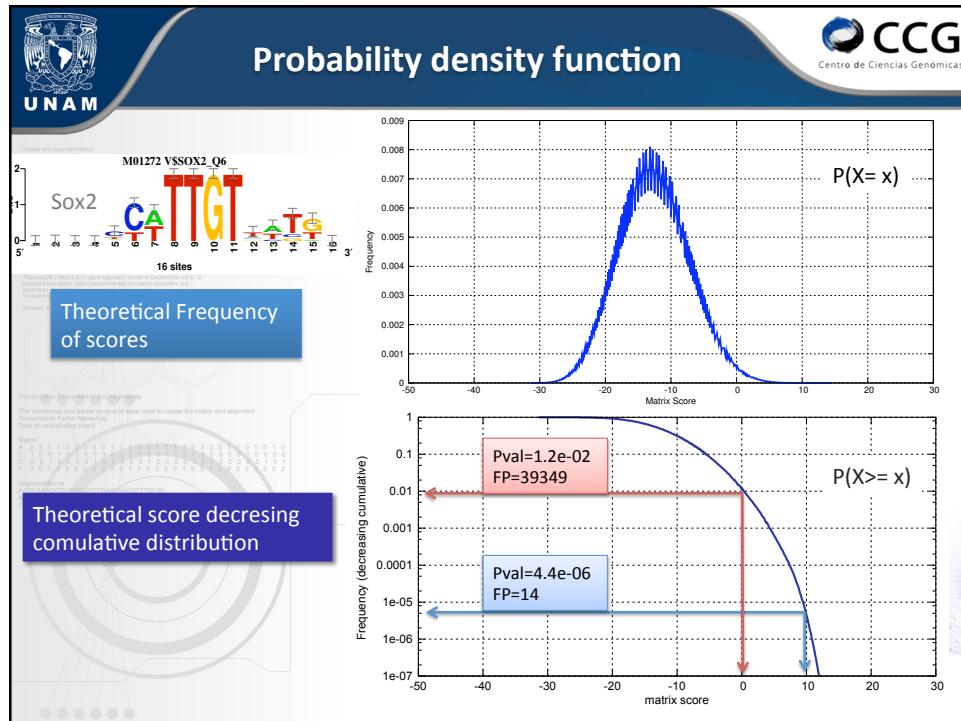


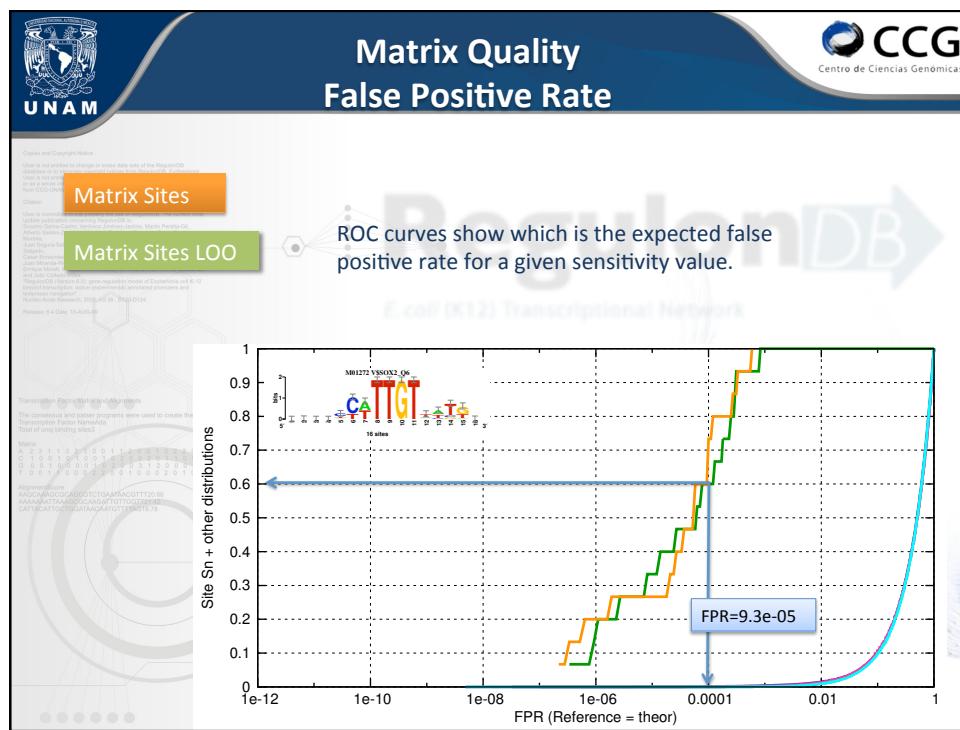
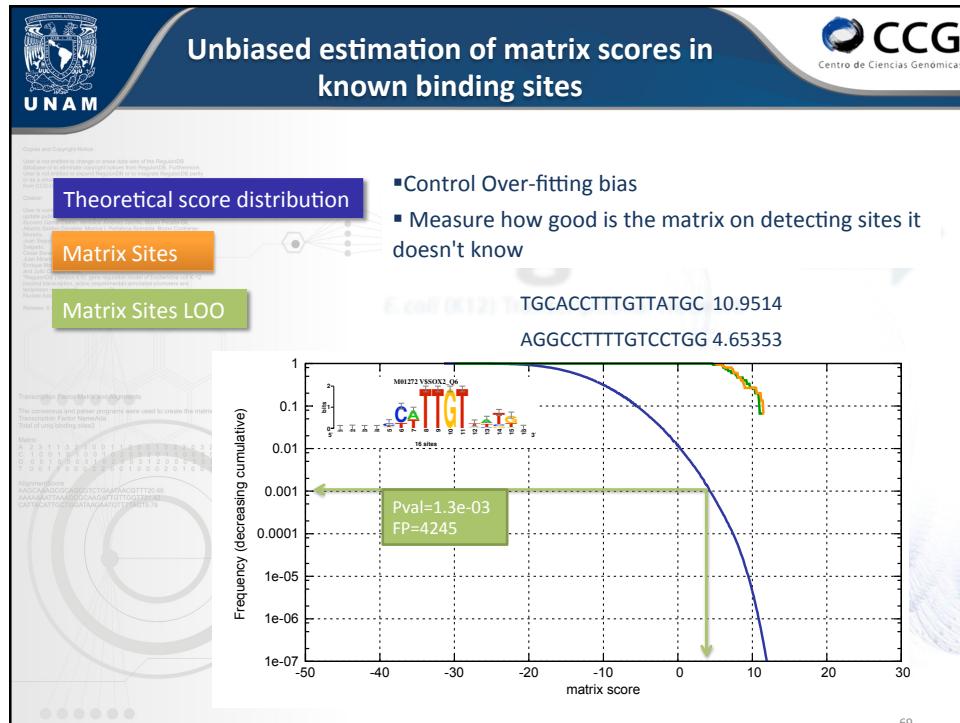
Hands on !

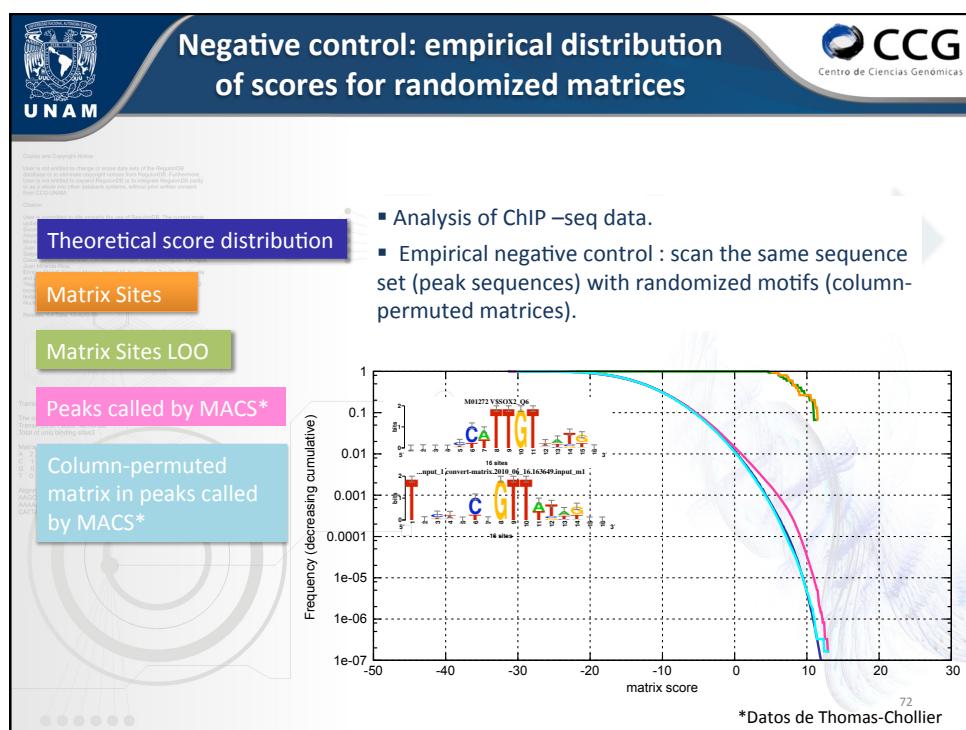
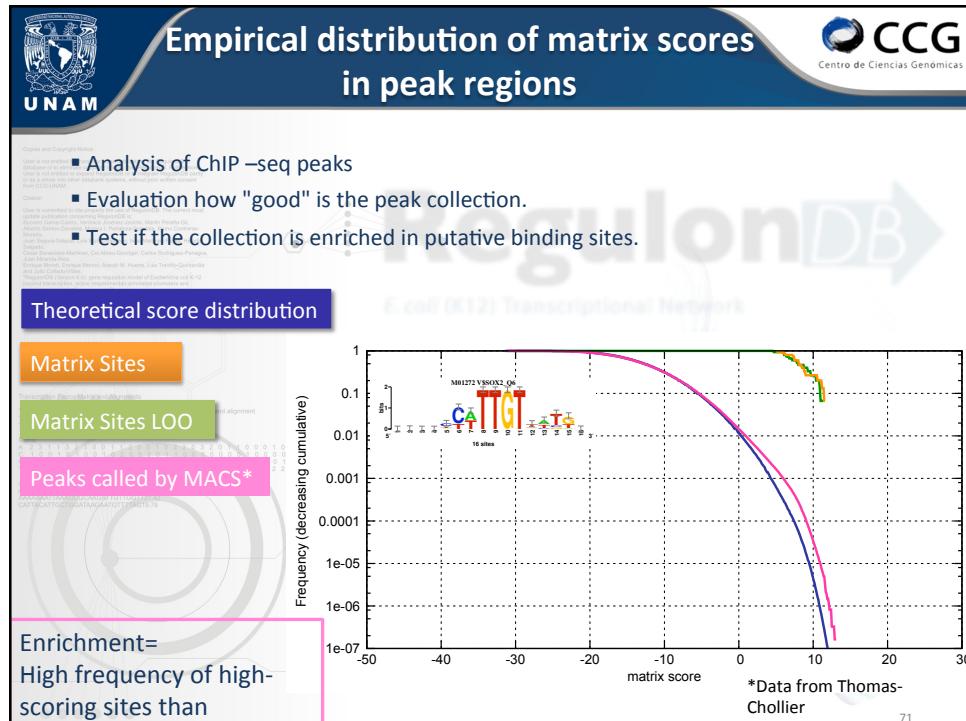
- Go to the companion website
- Follow all steps of **ChIP-seq in bacteria**

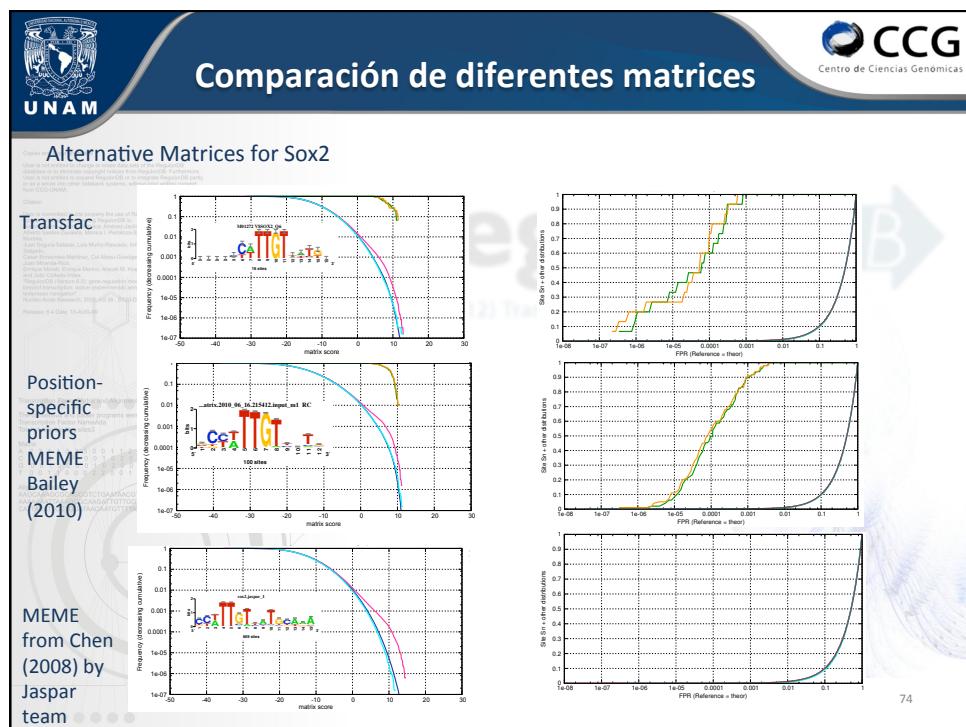
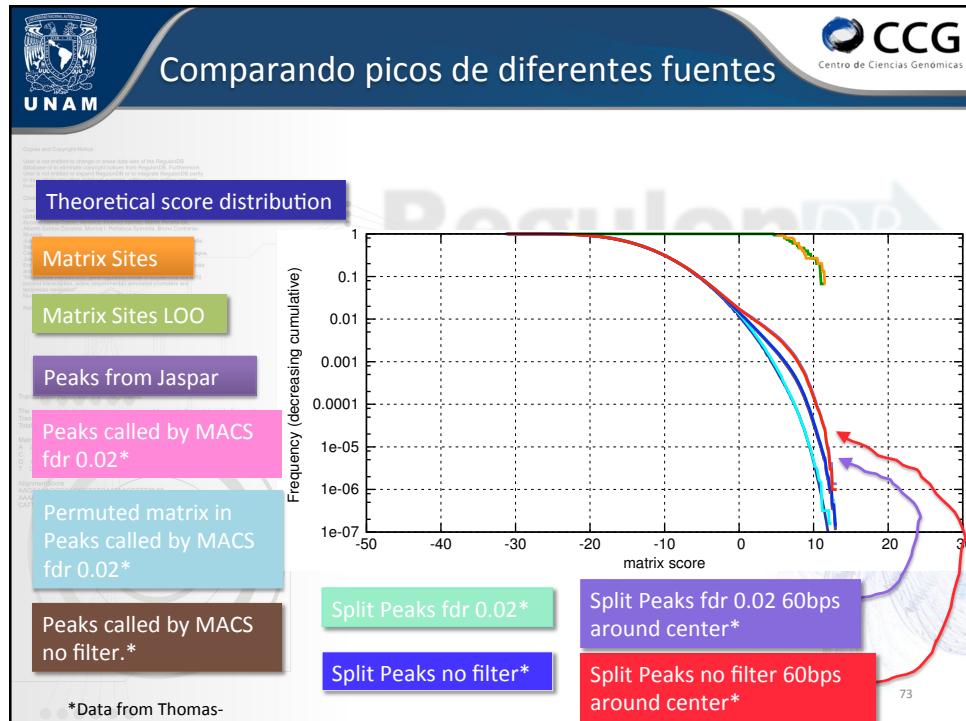
Morgane Thomas-Chollier











Hands on !

- Go to the companion website
- Follow all steps of **Measure the enrichment of your peak for the expected motif**

Morgane Thomas-Chollier