

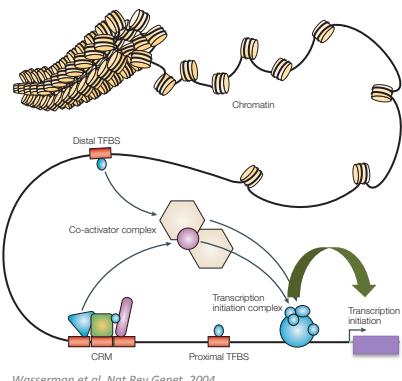
Détection de motifs dans les pics ChIP-seq

M. Thomas-Chollier, M. Defrance,
C. Herrmann, D. Puthier

Ecole de bioinformatique AVIESAN - Roscoff – 5-10 Octobre 2014

1

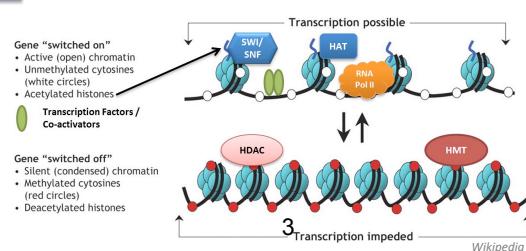
Biological concepts of transcriptional regulation



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

Morgane Thomas-Chollier

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



Goal and organisation of this session

Goal: introduction to motif analysis in ChIP-seq data

- **processing steps:** from reads to peaks. => see previous session
- **downstream analyses:**
 - focus on motif analyses

Tools

- RSAT

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*)

Morgane Thomas-Chollier

2

in vivo experimental methods to identify binding sites

ChIP (=Chromatin Immuno-Precipitation)

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

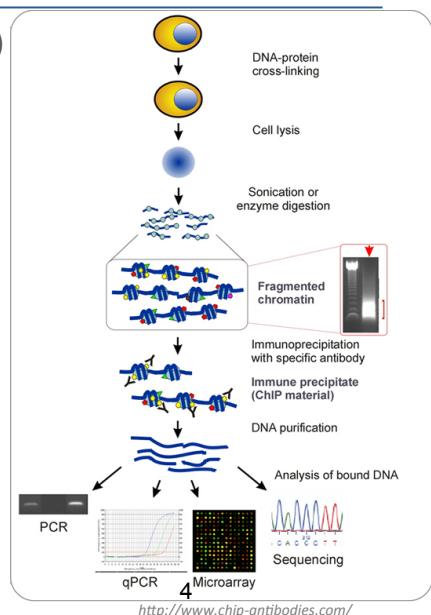
-small-scale: PCR / qPCR

-large-scale:

- microarray = **ChIP-on-chip**
- sequencing = **ChIP-seq**

Main challenge:

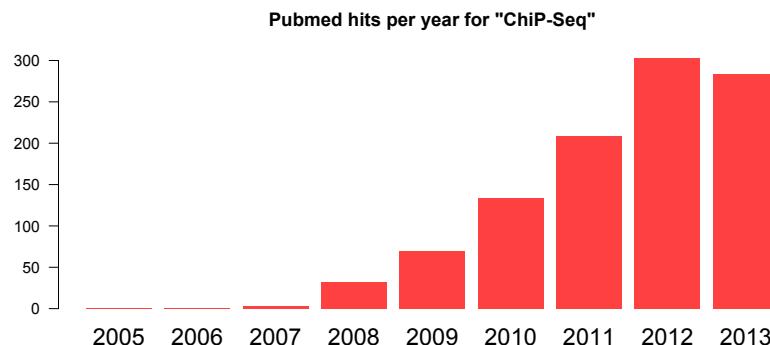
-quality/specificity of the antibodies



Morgane Thomas-Chollier

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

ChIP-seq is a recently-adopted technique !



Morgane Thomas-Chollier

5

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,...)

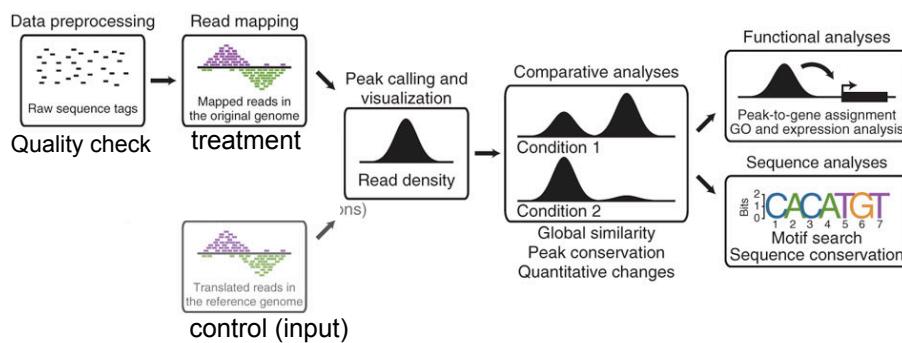
The obtain ChIP-seq **profiles** have **different shapes**, depending on the targeted protein

Morgane Thomas-Chollier



Park, Nature reviews 2009

ChIP-seq analysis workflow

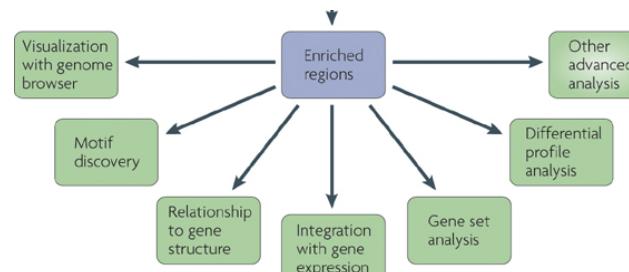


Morgane Thomas-Chollier

Adapted from Bardet et al, Nature Protocols, 2012

7

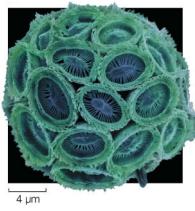
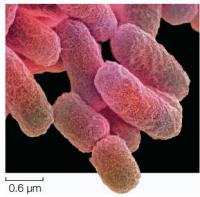
ChIP-seq analysis workflow: downstream analyses



Nature Reviews | Genetics

Morgane Thomas-Chollier

Par8, Nature reviews 2009



What is the biological question ?



Morgane Thomas-Chollier

9

What is the biological question ?

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

Morgane Thomas-Chollier

11

What is the biological question ?

« see if you can find something in the data »

Morgane Thomas-Chollier

10

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

12

What is the biological question ?

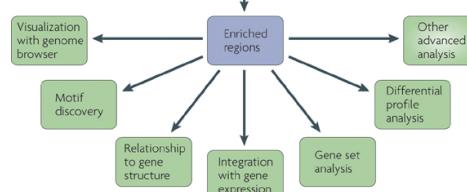
- **Where** do a transcription factor (TF) bind ?
 - ✓ In a **specific context** (tissue, developmental stage, mutant)
 - ✓ By **comparison** to another context (WT vs mutant, different time points)
- **How** do a transcription factor (TF) bind ?
 - ✓ Which **binding motif(s)** (can be several for a 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

13

What is the biological question ?

→ Should drive all « downstream » analyses



Morgane Thomas-Chollier

Nature Reviews | Genetics

15

What is the biological question ?

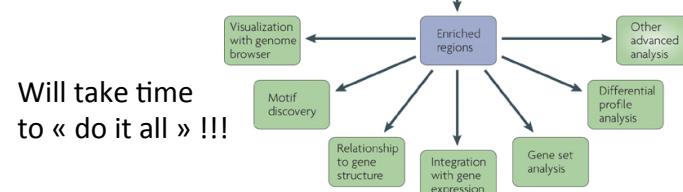
- **Where** do a transcription factor (TF) bind ?
 - ✓ In a **specific context** (tissue, developmental stage, mutant)
 - ✓ By **comparison** to another context (WT vs mutant, different time points)
- **How** do a transcription factor (TF) bind ?
 - ✓ Which **binding motif(s)** (can be several for a 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

14

What is the biological question ?

→ Should drive all « downstream » analyses



Morgane Thomas-Chollier

Nature Reviews | Genetics

16

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

Morgane Thomas-Chollier

17

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

18

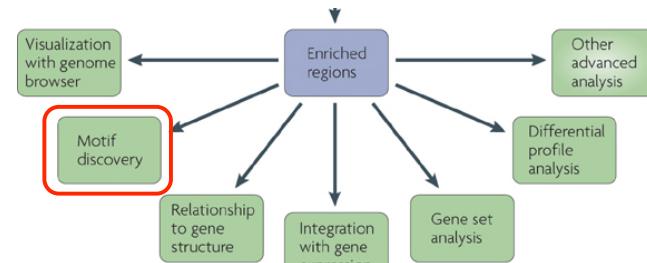
What is the biological question ?

- 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

19

ChIP-seq analysis workflow: downstream analyses



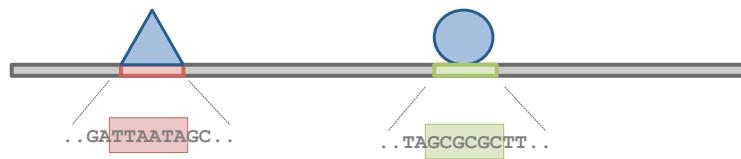
Nature Reviews | Genetics

Morgane Thomas-Chollier

Par20 Nature reviews 2009

Transcription factor specificity

How do TF « know » where to bind DNA ?



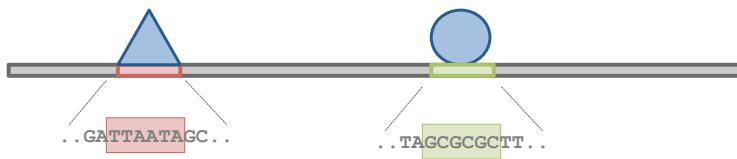
TF recognize TFBS with specific DNA sequences

Morgane Thomas-Chollier

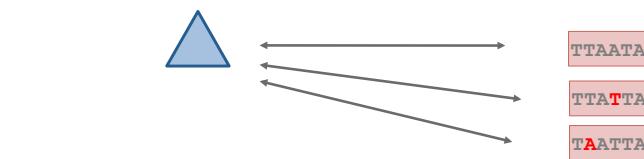
21

Transcription factor specificity

How do TF « know » where to bind DNA ?



TF recognize TFBS with specific DNA sequences

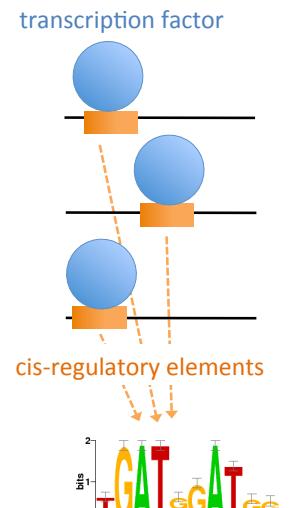


TFBSs are *degenerate*:
a given TF is able to bind DNA on TFBSs with different sequences

Morgane Thomas-Chollier

22

Binding specificity of a given TF

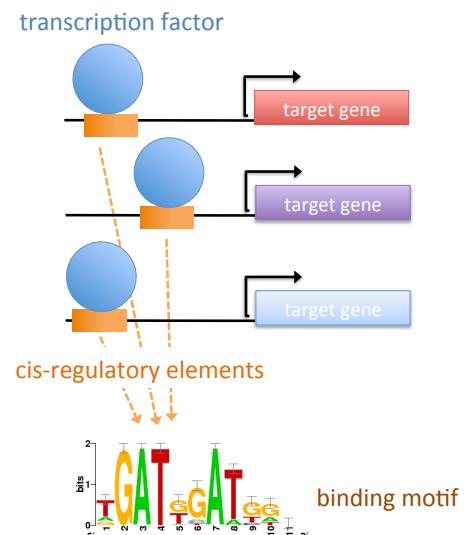


Morgane Thomas-Chollier

binding motif (represented as a sequence logo)

23

de novo motif discovery



Morgane Thomas-Chollier

Problem :
How can we model/describe
the binding specificity of
a given TF ?

If there is a common regulating
factor, can we discover its motif
only using these sequences ?

24

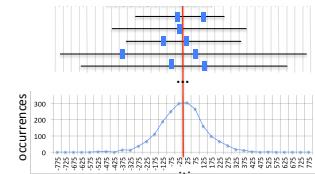
de novo motif discovery

- 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

25

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

27

de novo motif discovery

- Tools already exist for a long time !

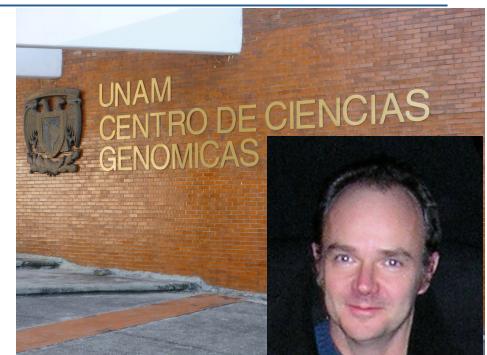
- MEME (1994)
 - RSAT oligo-analysis (1998)
 - AlignACE (2000)
 - Weeder (2001)
 - MotifSampler (2001)

Morgane Thomas-Chollier

26

Regulatory Sequence Analysis Tools (RSAT)

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

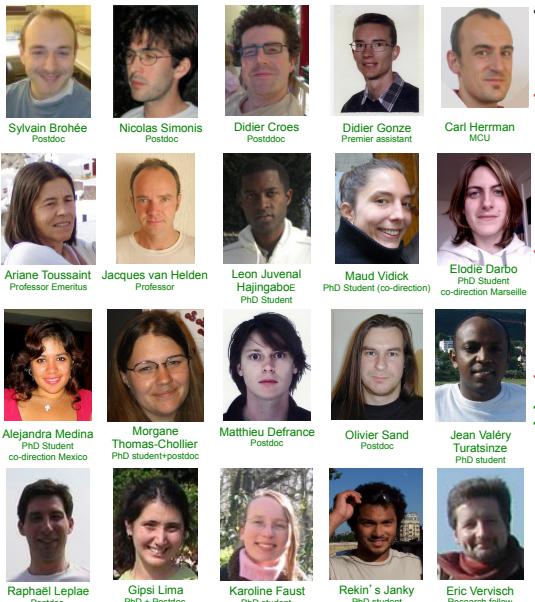


Jacques van Helden

28

Morgane Thomas-Chollier

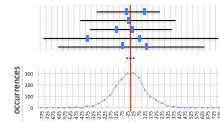
RSAT developers/contributors over the years



29

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



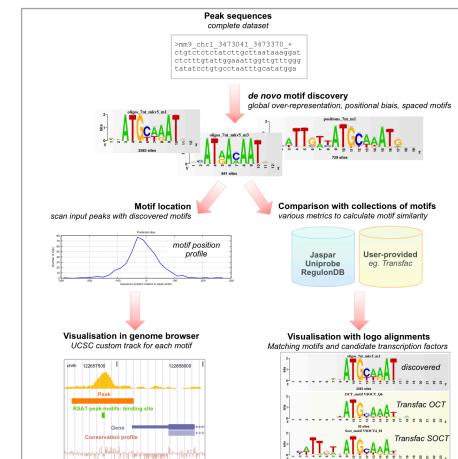
RSAT improvements over the years

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 *Nucleic Acids Research*, 2008
 van Helden, *Nucleic Acids Research*, 2003
 van Helden, André, Collado-Vides Yeast, 2000

30

New approaches for ChIP-seq datasets

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



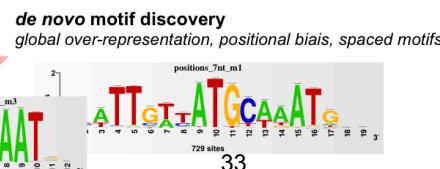
Peaks coordinates

BED

chr1 3002142 3002195
chr1 3002804 3002853

Peak sequences complete dataset

```
>mm9_chrl_3473041_3473370_+
ctgtctcttatcttgcttaataaaaggat
ctctttgtattggaaatttgggtttttgg
tatatatccgtgcctaatttgcataatgga
```



Peaks coordinates

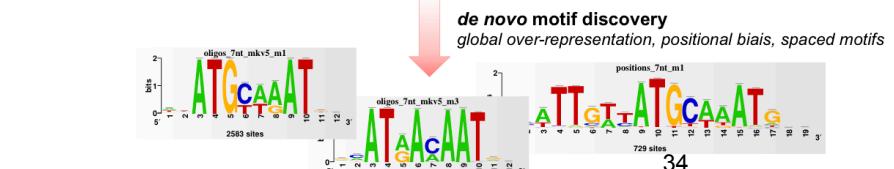
BED

chr1 3002142 3002195
chr1 3002804 3002853

Extract corresponding sequences

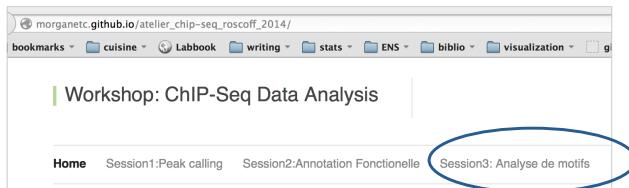
Peak sequences complete dataset

```
>mm9_chrl_3473041_3473370_+
ctgtctcttatcttgcttaataaaaggat
ctctttgtattggaaatttgggtttttgg
tatatatccgtgcctaatttgcataatgga
```



Hands on !

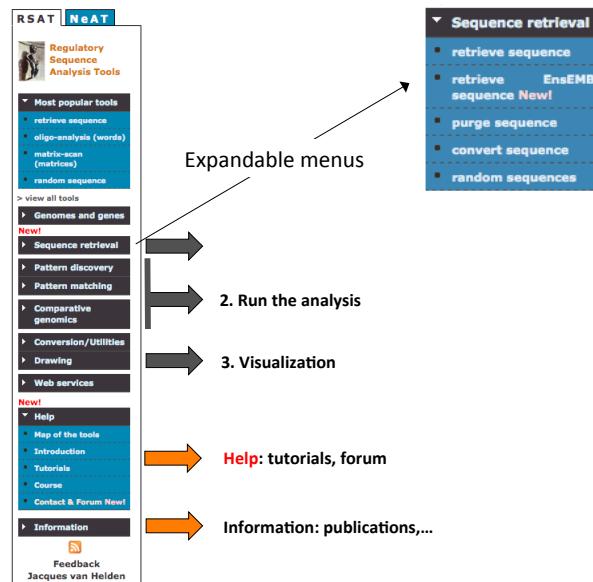
- Go to the companion website
 - <http://ecole-bioinfo-aviesan.sb-roscoff.fr/>
 - Section “chip-seq”



- Follow all steps of **Retrieving sequences from your peaks**

Morgane Thomas-Chollier

Using RSAT



RSAT Web forms

RSA-tools - retrieve sequence

Tool name: RSA-tools - retrieve sequence

Tool description: Returns upstream, downstream or ORF sequences for a list of genes

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](#).

Tool parameters:

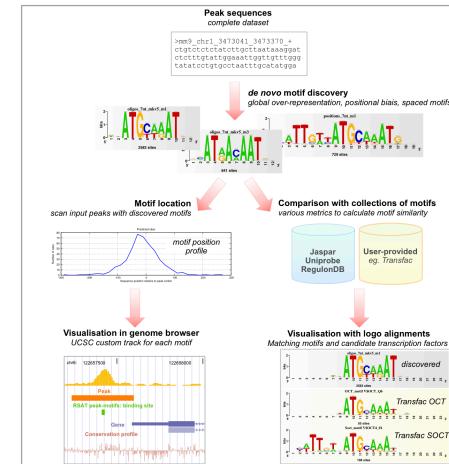
- Single organism: Organism: [Saccharomyces cerevisiae](#)
- Multiple organisms
- Genes: all selection
- Upload gene list from file:
- Query contains only IDs (no synonyms):
- Feature type: CDS, mRNA, tRNA, rRNA, scRNA
- Sequence type: upstream, from default, To default
- Prevent overlap with neighbour genes (noorf):
- Mask repeats (only valid for organisms with annotated repeats):
- Admit imprecise positions:
- Sequence format: fasta
- Sequence label: gene name
- Output: server, display, email

Output:

Go button (launches the analysis) Demo button (fill in the form for test purposes) 37 Help

New approaches for ChIP-seq datasets

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



Thomas-Chollier et al *Nucleic Acids Research*, 2012 38

Morgane Thomas-Chollier

Comparison of tools for ChIP-seq

Program	peak-motifs	ChIPMunk	CompleteMotifs	MEME-ChIP	MICSA	GimmeMotifs
Web interface	yes	yes	yes	yes	no	no
Size limitation	unrestricted (Web site tested with 22 Mb)	100kb (web site)	500kb (web site)	unrestricted, but motif discovery restricted to a few hundred base pairs (clipped to 100bp)	motif discovery restricted to a few hundred base pairs	-
Stand-alone version	yes	yes	yes	yes	yes	yes
Tasks						
peak finding	no	no	no	yes	no	
identification of peak-flanking genes	no	no	yes	no	no	
sequence composition (mono- and di-nucleotides)	yes	no	no	no	no	
motif scanning	yes	yes	yes	yes	yes	
enrichment in motifs from databases	no	yes	no	no	no	
enrichment in discovered motifs	yes	no	no	yes	no	
peak scoring	no	no	yes	yes	no	
motif clustering	no	no	no	no	yes	
comparison discovered motifs / motif DB	yes	no	no	yes	yes	
sequence scanning for site prediction	yes	no	no	yes	no	
positional distribution of sites inside peaks	yes	no	yes	no	yes	
visualization in genome browsers	yes	no	yes	no	no	
Motif discovery algorithms	RSAT motif-analysis RSAT dyad-analysis RSAT position-analysis RSAT local-position-analysis + in stand-alone version: MEME ChIPMunk	ChIPMunk MEME Weeder	MEME DREME	MEME Weeder MotifSampler BioProspector Gadem Improbizer HOMERsuite Trawler MoAn		
Pattern matching algorithms	RSAT matrix-scan-quick	no	patser	MAST + AME (bioProspector)	no	
Motif comparison algorithm	RSAT compare-motifs	no	STAMP	TOMTOM	STAMP	
Motif clustering algorithm						
Comparison of discovered motifs	yes	no	no	no	no	
Motif database comparisons	no	JASPAR TRANSFAC DMPMMB RegulonDB upload your own database	JASPAR TRANSFAC UNIPROBE FLYREG DMPMMB SCPD DMPMMB and many others	JASPAR TRANSFAC UNIPROBE FLYREG DMPMMB SCPD DMPMMB and many others	JASPAR TRANSFAC UNIPROBE FLYREG DMPMMB SCPD DMPMMB and many others	
Motif sizes	variable (multiple word assembly)	user-specified <=35 for MEME <=12 for Weeder <=13 for ChIPMunk		predefined ranges (small, medium, large, extra-large)		predefined ranges (small, medium, large, extra-large)
Multiple motifs	yes	no	yes	yes	yes	yes
Ref (PMID)	This article	20736340 21163795 21480936	20373099 21081511	20373099 21081511	20373099 21081511	20373099 21081511

39

Comparison of tools for ChIP-seq

Program	peak-motifs	ChIPMunk	CompleteMotifs	MEME-ChIP	MICSA	GimmeMotifs
Web interface	yes	yes	yes	yes	no	no
Size limitation	unrestricted (Web site tested with 22 Mb)	100kb (web site)	500kb (web site)	unrestricted, but motif discovery restricted to a few hundred base pairs (clipped to 100bp)	motif discovery restricted to a few hundred base pairs	-
Tasks						
peak finding	no	no	no	yes	yes	no
identification of peak-flanking genes	no	yes	yes	no	no	no
sequence composition (mono- and di-nucleotides)	yes	no	no	no	no	no
enrichment in motifs from databases	no	no	yes	yes	yes	no
enrichment in discovered motifs	yes	no	no	yes	yes	no
peak scoring	no	no	yes	yes	yes	no
motif clustering	no	no	no	no	yes	yes
comparison discovered motifs / motif DB	yes	no	no	yes	yes	yes
sequence scanning for site prediction	yes	no	no	yes	no	no
positional distribution of sites inside peaks	yes	no	yes	yes	yes	yes
visualization in genome browsers	yes	no	yes	no	no	
Motif discovery algorithms	RSAT motif-analysis RSAT dyad-analysis RSAT position-analysis RSAT local-position-analysis + in stand-alone version: MEME ChIPMunk	ChIPMunk MEME Weeder	MEME DREME	MEME Weeder MotifSampler BioProspector Gadem Improbizer HOMERsuite Trawler MoAn		
Pattern matching algorithms	RSAT matrix-scan-quick	no	patser	MAST + AME (bioProspector)	no	
Motif comparison algorithm	RSAT compare-motifs	no	STAMP	TOMTOM	STAMP	
Motif clustering algorithm						
Comparison of discovered motifs	yes	no	no	yes	no	
Motif database comparisons	no	JASPAR TRANSFAC DMPMMB RegulonDB upload your own database	JASPAR TRANSFAC UNIPROBE FLYREG DMPMMB SCPD DMPMMB and many others	JASPAR TRANSFAC UNIPROBE FLYREG DMPMMB SCPD DMPMMB and many others	JASPAR TRANSFAC UNIPROBE FLYREG DMPMMB SCPD DMPMMB and many others	
Motif sizes	variable (multiple word assembly)	user-specified <=35 for MEME <=12 for Weeder <=13 for ChIPMunk		predefined ranges (small, medium, large, extra-large)		predefined ranges (small, medium, large, extra-large)
Multiple motifs	yes	no	yes	yes	yes	yes
Ref (PMID)	This article	20736340 21163795 21480936	20373099 21081511	20373099 21081511	20373099 21081511	20373099 21081511

Thomas-Chollier, Herrmann, DeFrance, Sand, Thieffry, van Helden *Nucleic Acids*

40

Comparison of tools for ChIP-seq

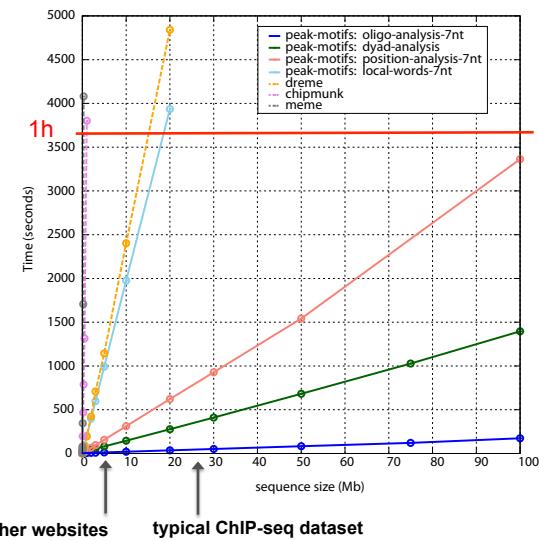
Program	peak-motifs	Chimpunk	CompleteMotifs	MEME-ChIP	MICSA	GimmeMotifs
Web interface	yes	yes	yes	yes	no	no
Size limitation	unrestricted (Web site tested with 22 Mb)	100kb (web site)	500kb (web site)	unrestricted, but motif discovery restricted to 600 peaks clipped to 100bp	motif discovery restricted to a few hundred base pairs	-
peak finding	no	no	no	yes	no	-
annotation of peak-flanking genes	no	no	yes	no	-	-
discovered motifs (mono- and di-nucleotides)	yes	no	no	no	-	-
motif discovery	yes	yes	yes	yes	yes	-
enriched motifs from databases	yes	yes	yes	yes	yes	-
enrichment in discovered motifs	yes	no	no	no	no	-
peak scoring	no	no	no	yes	no	-
motif clustering	no	no	no	no	yes	-
comparison with motifDB / motifDB	yes	no	yes	yes	yes	-
sequence scanning for site presence	yes	no	yes	no	no	-
positional distribution of sites inside peaks	yes	no	yes	no	yes	-
visualization in genome browsers	yes	no	no	no	-	-
Motif discovery algorithms	RSTAT oligo-analysis RSTAT dyad-analysis RSTAT position analysis RSTAT local-word-analysis + in stand-alone version: Chimpunk	Chimpunk MEME Weeder	Chimpunk MEME DREME Weeder	MEME Weeder MotSinger BioProspector Gadem Imprimer HDMotif MotNMF	-	-
Pattern matching algorithms	RSAT matrix-scan-quick	RSAT	RSAT	MAST + AHE (enrichment)	-	-
Motif comparison algorithm	RSAT compare-motifs	RSAT	RSAT	TOMTOM	STAMP	STAMP
Motif clustering algorithm	-	-	-	-	-	YES
Comparisons between discovered motifs	-	-	-	-	-	-
Motif database comparisons	JASPAR UNIPROBE DMMHMM RegulonDB upload your own database	JASPAR TRANSFAC	JASPAR TRANSFAC UNIPROBE DPMEME DPINTERACT DMMHMM and many others	-	-	-
Motif sizes	variable (multiple word assembly)	user-specified	<=25 for RSAT <=10 for Weeder <=13 for Chimpunk	-	predicted ranges (small, medium, large, extra-large)	-
Multiple motifs	yes	yes	yes	yes	yes	yes
Ref (PMID)	This article	20738340	21133395	21400794	20373509	21091311

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

41

Peak-motifs: why providing yet another tool ?

- fast and scalable
- treat full-size datasets

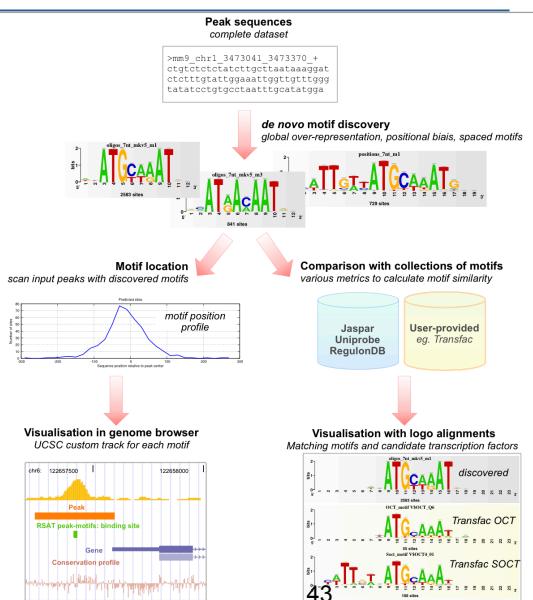


Thomas-Chollier, Herrmann, Defrance, Sand, Thieffry, van Helden *Nucleic Acids Research*, 2012

42

Peak-motifs: why providing yet another tool ?

- fast and scalable
- treat full-size datasets
- complete pipeline



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

43

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
- HTML report

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

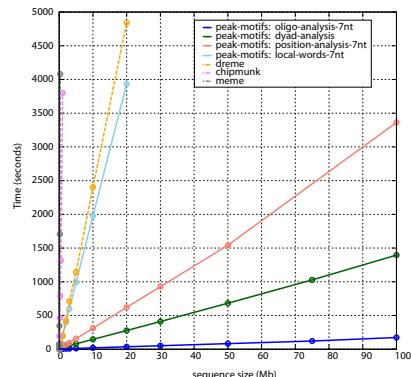
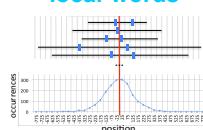
44

Morgane Thomas-Chollier

Peak-motifs: why providing yet another tool ?

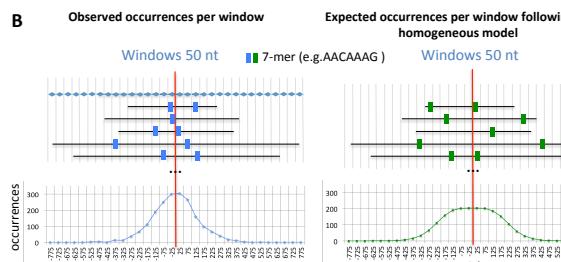
- fast and scalable
- treat full-size datasets
- complete pipeline
- web interface
- accessible to non-specialists
- using 4 complementary algorithms

- Global over-representation
 - oligo-analysis
 - dyad-analysis (spaced motifs)
- Positional bias
 - position-analysis
 - local-words

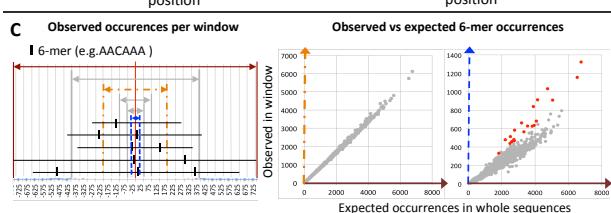


45

Motif discovery methods: positional bias

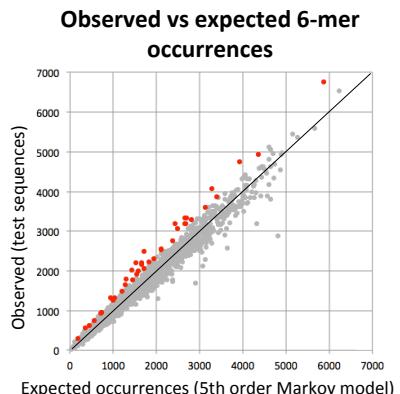
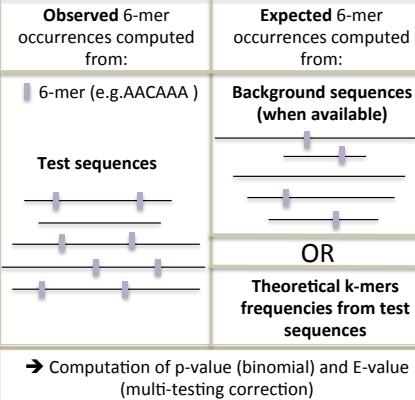


position-analysis



local-words

Motif discovery methods: frequency



oligo-analysis
dyad-analysis (spaced motifs)

46

Hands on !

- Go to the companion website
- Follow all steps of Discovering motifs from peak sequences

Hands on !

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

Morgane Thomas-Chollier

49

Hands on !

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

Morgane Thomas-Chollier

50

Hands on !

- Go to the companion website
- Follow all steps of **Motif analysis with MEME-chip**

Morgane Thomas-Chollier

51

To go further

- The next slides explain step by step the algorithm behind **oligo-analysis**
- **Peak-motifs** : follow this protocol to grasp the detailed tweaking of parameters (send us an email to have free access to the PDF if necessary)
Thomas-Chollier et al. *A complete workflow for the analysis of full-size ChIP-seq (and similar) data sets using peak-motifs*. Nature Protocols 7, 1551–1568 (2012).
- **Matrix-quality** : RSAT program that can be used to evaluate the enrichment of motifs in peaks
Medina-Rivera A, Abreu-Goodger C, Thomas-Chollier M, Salgado H, Collado-Vides J, van Helden J. Theoretical and empirical quality assessment of transcription factor-binding motifs. Nucleic Acids Res. 2011 Feb;39(3):808-24. doi: 10.1093/nar/gkq710. Epub 2010 Oct 4.

Morgane Thomas-Chollier

52

To go further

- Tutorial for ECCB 2014 : <http://rsat.ulb.ac.be/eccb14/>
- Master classes in analysis of cis-regulatory regions (over one week) at Ecole Normale Supérieure every september (contact : mthomas@biologie.ens.fr)

Morgane Thomas-Chollier

53

Motif discovery using word counting

Idea:

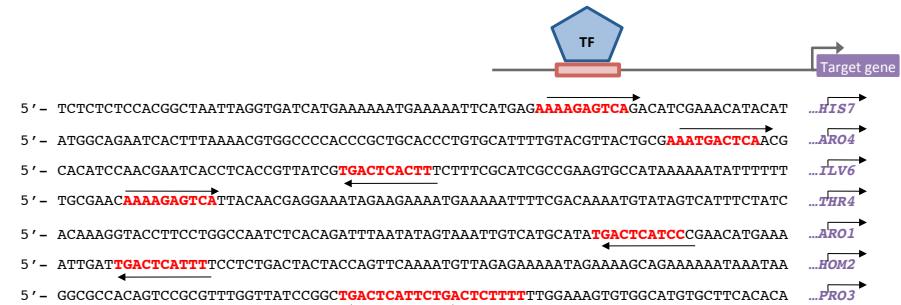
motifs corresponding to binding sites are generally repeated in the dataset
→ capture this statistical signal

Algorithm

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

55

Principle: detect unexpected patterns

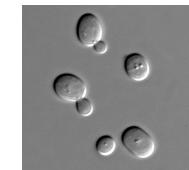


- Binding sites are represented as “words” = “string”=“k-mer”
 - e.g. acgtga is a 6-mer
- Signal is likely to be **more frequent** in the upstream regions of the co-regulated genes than in a random selection of genes
- We will thus detect **over-represented words**

54

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	atatat atatat	41	cttata gataag	26
aagaaa ttttt	14	aaaaaa ttttc	40	tatata tatata	22
gaaaaa ttttc	13	tatata tatata	40	ataaga tcttat	20
tgccaa ttggca	12	aaaaat atttt	35	aagaaa tttctt	20
aaaaat atttt	12	aagaaa ttttt	29	gaaaaa tttttc	19
aaatta taattt	12	agaaaa ttttt	28	atatat atatat	19
agaaaa ttttct	11	aaaata tatttt	26	agataa ttatct	17
caagaa ttcttg	11	aaaaag ctttt	25	agaaaa ttttct	17
aaacgt acgttt	11	agaaaat attttc	24	aaaggaa ttcttt	16
aaagaa ttttt	11	aaataa tttatt	22	aaaaca tgtttt	16
acgtgc gcacgt	10	aaaaaa ttttta	21	aaaaaa cttttt	15
aataat attatt	10	tgaaaa ttttca	21	agaaga tctttt	14
aagaag cttctt	10	ataata tattat	20	tgataa ttatca	14
atataa ttatat	10	atataa tttata	20	atataa tttata	14

56

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		
aaaaaa tttttt	51	
aaaaag cttttt	15	
aagaaa ttttct	14	
gaaaaa tttttc	13	
tgccaa ttggca	12	
aaaaat attttt	12	
aaatata taattt	12	
agaaaa ttttct	11	
caagaa ttcttg	11	
aaacgt acgttt	11	
aaagaa tttcctt	11	
acgtgc gcacgt	10	
ataaat attattt	10	
agaagg cttctt	10	
atataaa ttatat	10	

MET		
aaaaaa tttttt	105	
atatat atatat	41	
gaaaaa tttttc	40	
tatata tatata	40	
aaaaat attttt	35	
aagaaa tttctt	29	
aaaaaa tttttc	28	
atatat atatat	28	
aaaata tatttt	26	
aaaaag cttttt	25	
agaaat tttctt	24	
aaataa tttatt	22	
aaaaaa tttttt	21	
tgaaaa tttca	21	
ataata tattat	20	
atataaa ttatat	20	

NIT		
aaaaaa tttttt	80	
cttatac gataag	26	
tatata tatata	22	
ataaga tcttat	20	
aagaaa tttctt	20	
aaaaaa tttttc	19	
atatat atatat	19	
agataaa ttatct	17	
agaaaa ttttct	17	
aaagaa ttcttt	16	
aaaaca tgtttt	16	
aaaaaa cttttt	15	
agaaga tcttct	14	
tgataa ttatca	14	
atataaa ttatat	14	

57

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”

58

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)

59

Estimation of word expected frequencies from background sequences



Example:

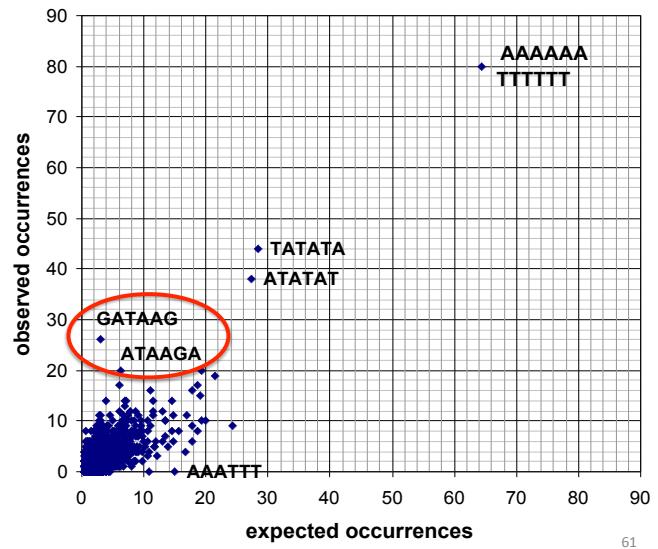
6nt frequencies in the whole set of 6000 yeast **upstream** sequences

: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
aaaaaca	aaaaca tgttt	0,0019828	5651
aaaaacc	aaaacc ggttt	0,00088526	2523
aaaaacg	aaaacg cgttt	0,00090105	2568
aaaaact	aaaact agttt	0,0014621	4167
aaaaga	aaaaga tcctt	0,00323016	9206
aaaagc	aaaagc gcttt	0,00135824	3871
aaaagg	aaaagg ccctt	0,0017849	5087
aaaagt	aaaagt acttt	0,0019035	5425
aaaata	aaaata tattt	0,00336805	9599
aaatac	aaatac gattt	0,00131368	3744
aaaatg	aaaatg cattt	0,00185648	5291
aaaatt	aaaatt atttt	0,00269156	7671
aaacaa	aaacaa ttgtt	0,00209999	5985
aacac	aacac gtgtt	0,00071684	2043
aacacg	aacacg ctgtt	0,00096491	2750
aacat	aacat atgtt	0,00108982	3106
aaacca	aaacca tggtt	0,00074421	2121

60

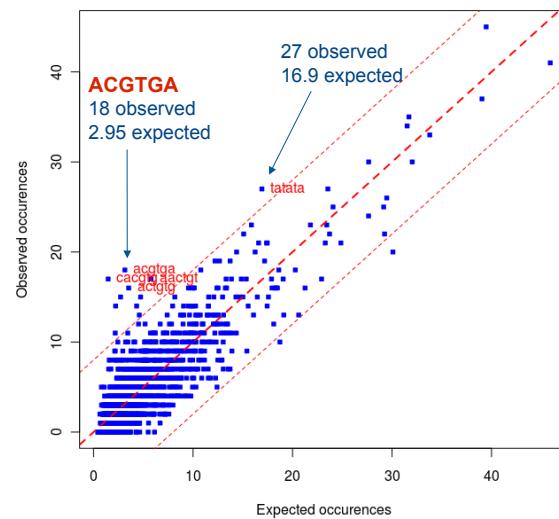
NIT		
aaaaaaa ttttttt	80	
cttatac gataag	26	
tatata tatata	22	
ataaga tcttat	20	
aagaaa tttctt	20	
gaaaaaa tttttc	19	
atataat atatat	19	
agataaa ttatct	17	
agaaaaa ttttct	17	
aaagaaa ttcttt	16	
aaaaaca tggttt	16	
aaaaaaag cttttt	15	
agaaga tcttct	14	
tgataaa ttatca	14	
atataaa ttatat	14	

Hexanucleotide occurrences in upstream sequences of the NIT family



61

Motif discovery using word counting



How to evaluate expected number of occurrences ?

62

Motif discovery using word counting

Idea:

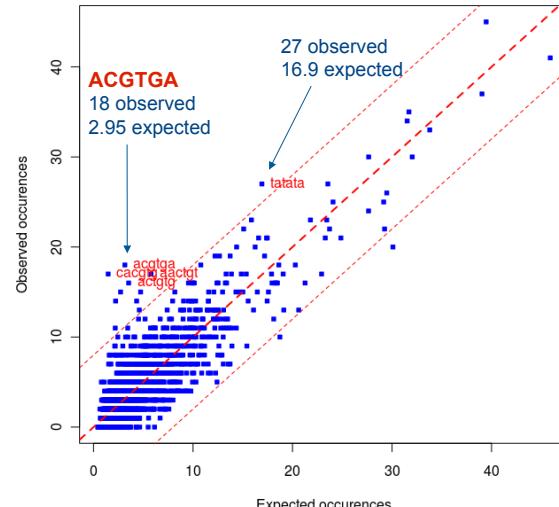
motifs corresponding to binding sites are generally repeated in the dataset
→ capture this statistical signal

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)

63

Statistical evaluation



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

64

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”