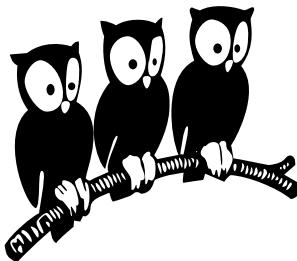


Département de Physique
École Normale Supérieure

Laboratoire de Physique Statistique



THÈSE de DOCTORAT de l'UNIVERSITÉ PARIS 7

Spécialité : Physique Théorique

présentée par

Marc SANTOLINI

pour obtenir le grade de DOCTEUR de l'UNIVERSITÉ PARIS 7

**Analyse computationnelle des éléments cis-régulateurs
dans les génomes d'eucaryotes supérieurs**

Soutenue le ZZ septembre 2013
devant le jury composé de :

M.	Emmanuel BARILLOT	Examinateur
M.	Vincent HAKIM	Directeur de thèse
M.	Pascal MAIRE	Examinateur
M.	Massimo VERGASSOLA	Rapporteur
M.	Martin WEIGT	Rapporteur
M.	Alain ZIDER	Examinateur

these:version du vendredi 5 juillet 2013 à 20 h 42

Remerciements

...

thèse:version du vendredi 5 juillet 2013 à 20 h 42

Table des matières

Liste des figures	vii
Principales abréviations utilisées	ix
Avant-propos	1
Chapitre 1 - Introduction générale.	3
1.1 Le phénotype cellulaire	4
1.2 Les réseaux de régulation génétique	9
1.3 Les interactions protéine-ADN : modèles mathématiques	19
1.4 Les interactions protéine-ADN : mesures expérimentales	26
1.5 Les modules de cis-régulation (CRMs)	36
1.6 Prédiction et validation des CRMs	49
1.7 Bases de données	58
Chapitre 2 - Modèles de fixation des Facteurs de Transcription à l'ADN.	65
2.1 Observations de corrélations au sein des TFBS	66
2.2 Modèles existants permettant de décrire la statistique des TFBS	67
2.3 Modèles de maximum d'entropie	72
2.4 Article	77
2.5 Analyse thermodynamique des modèles	109
2.6 Conclusion et perspectives	112

Liste des figures

Introduction générale.	3
1.1 Le paysage de la différenciation cellulaire	5
1.2 Spécification spatio-temporelle du type cellulaire	7
1.3 Différents exemples de reprogrammation cellulaire	8
1.4 Vision cybernétique du traitement de l'information par la cellule	10
1.5 Un réseau de régulation génétique type	11
1.6 Caractéristiques de l'épigénome	14
1.7 Exemples de motifs dans les réseaux de régulation génétique	15
1.8 Évolution du réseau transcriptionnel : l'exemple de la régulation myogénique. .	17
1.9 Différents états du facteur de transcription	19
1.10 Construction et utilisation du modèle PWM	22
1.11 Étapes d'une expérience de ChIP-on-chip et ChIP-seq	32
1.12 Résolution des expériences ChIP-on-chip et ChIP-seq	33
1.13 Expérience d'empreinte à la DNase I chez la levure : vers une résolution au nucléotide près	35
1.14 Les différents types de CRMs et leurs marques épigénétiques	37
1.15 Différents <i>enhancers</i> conduisent à différents patterns d'expression	39
1.16 Deux modèles d' <i>enhancers</i> : enhanceosome et billboard	41
1.17 L'enhanceosome de l'interferon- β	42
1.18 Flexibilité du code de cis-régulation au cours de l'évolution chez les <i>Drosophiles</i> .	43
1.19 Évolution de la fixation de HNF4 α chez les mammifères	45
1.20 « Shadow enhancer » du gène de segmentation <i>Hunchback</i>	46
1.21 De l' <i>enhancer</i> au super- <i>enhancer</i>	47
1.22 Différentes approches pour la prédiction des CRMs	50
1.23 Méthodes de validation des CRMs par transfection et transgenèse	55
1.24 Impact physiologique de la délétion et de la mutation d'un enhancer	56
1.25 Évolution du coût de séquençage	58
1.26 Distribution des tailles intergéniques et introniques chez différentes espèces .	60

Liste des figures

1.27	Visualisation de données ChIP-seq <i>via</i> le site UCSC	62
1.28	Les différentes données obtenues par le projet ENCODE	63
Modèles de fixation des Facteurs de Transcription à l'ADN.		65
2.1	Différents modèles pour décrire les corrélations entre nucléotides dans les sites de fixation de facteurs de transcription	68
2.2	Illustration d'un système dont on veut maximiser l'entropie	72
2.3	Chaleur spécifique pour différents TFs	111
2.4	Histogrammes des valeurs des champs h et couplages J	113

Principales abréviations utilisées

ARNm	ARN messager
bHLH	<i>basic Helix-Loop-Helix</i>
bp	Paire de base
ChIP	Immunoprécipitation de la chromatine (<i>Chromatin immunoprecipitation</i>)
CRM	Module de cis-régulation (<i>Cis-Regulatory Module</i>)
DHS	Hypersensible à la DNase I (<i>DNaseI-hypersensitive</i>)
ESC	Cellule souche embryonnaire (<i>Embryonic Stem Cell</i>)
ISH	Hybridation <i>in situ</i> (<i>In-Situ Hybridization</i>)
kb	kilobases (1000bp)
MRF	Facteur de régulation myogénique (<i>Myogenic Regulatory Factor</i>)
nt	Nucléotide
PCR	Réaction en chaîne par polymérase (<i>Polymerase Chain Reaction</i>)
PWM	Matrice de poids (<i>Position Weight Matrix</i>)
TF	Facteur de transcription (<i>Transcription Factor</i>)
TFBS	Site de fixation d'un facteur de transcription (<i>Transcription Factor Binding Site</i>)
TSS	Site d'initiation de la transcription (<i>Transcription Start Site</i>)

Avant-propos

Cette thèse se présente sous la forme suivante...

Voici quelques remarques sur la version pdf de ce manuscrit, qui peuvent rendre la lecture plus aisée. Dans la table des matières, la liste des figures et la liste des annexes, les titres sont des liens hypertexte qui pointent vers l'item décrit. Dans la liste des notations utilisées et la bibliographie, ce sont les numéros de page qui sont des liens hypertexte.

these:version du vendredi 5 juillet 2013 à 20 h 42

Avant-propos

Chapitre 1

Introduction générale.

1.1 Le phénotype cellulaire

1.1.1 Qu'est-ce que le phénotype d'une cellule ?

Les organismes vivants sont constitués de cellules de l'ordre de quelques microns, facilement observables à l'aide d'un simple microscope optique. Chaque cellule contient un certain nombre de constituants (gènes, protéines, métabolites...) enclos par une membrane. Il existe des organismes unicellulaires (bactérie, levure) et multicellulaires (mouche, souris, homme). Ce sont ces derniers auxquels nous nous intéressons dans cette thèse. Les cellules qui les constituent sont majoritairement eucaryotes, c'est-à-dire qu'elles possèdent un noyau renfermant le matériel génétique.¹

Bien que possédant toutes le même matériel génétique, les cellules d'un organisme apparaissent d'emblée comme hétérogènes, que ce soit dans leur forme ou dans leurs constituants. Par exemple, chez l'homme, les érythrocytes ou globules rouges présents dans le sang sont des cellules de la forme d'un disque biconcave, dépourvues de noyau et riches en hémostoglobine, tandis que les fibres musculaires squelettiques sont de forme longue et tubulaire, possèdent plusieurs noyaux et expriment actine et myosine.

Cette diversité semble néanmoins limitée. Aussi, parmi les $\sim 6 \cdot 10^{13}$ cellules du corps humain, on peut distinguer ~ 320 différents types cellulaires (?). Bien entendu, ce nombre dépend du seuil de similarité choisi : deux cellules d'un même type n'expriment pas *exactement* le même nombre de molécules. Classiquement, la classification d'un type cellulaire se base sur des propriétés morphologiques observables au microscope ou encore sur l'analyse de molécules présentes à la surface des cellules. Par ailleurs, différents types cellulaires sont associés à différentes fonctions : dans notre exemple la fixation et le transport de l'oxygène dans le cas des globules rouges, la contraction dans le cas des fibres musculaires.

Ces différentes propriétés observables caractérisent le *phénotype* cellulaire (étymologiquement « exhiber un type » en grec). Nous allons le voir, ce phénotype est le résultat de la modulation par des facteurs environnementaux de l'expression génétique qui conditionne le contenu en protéines de la cellule.

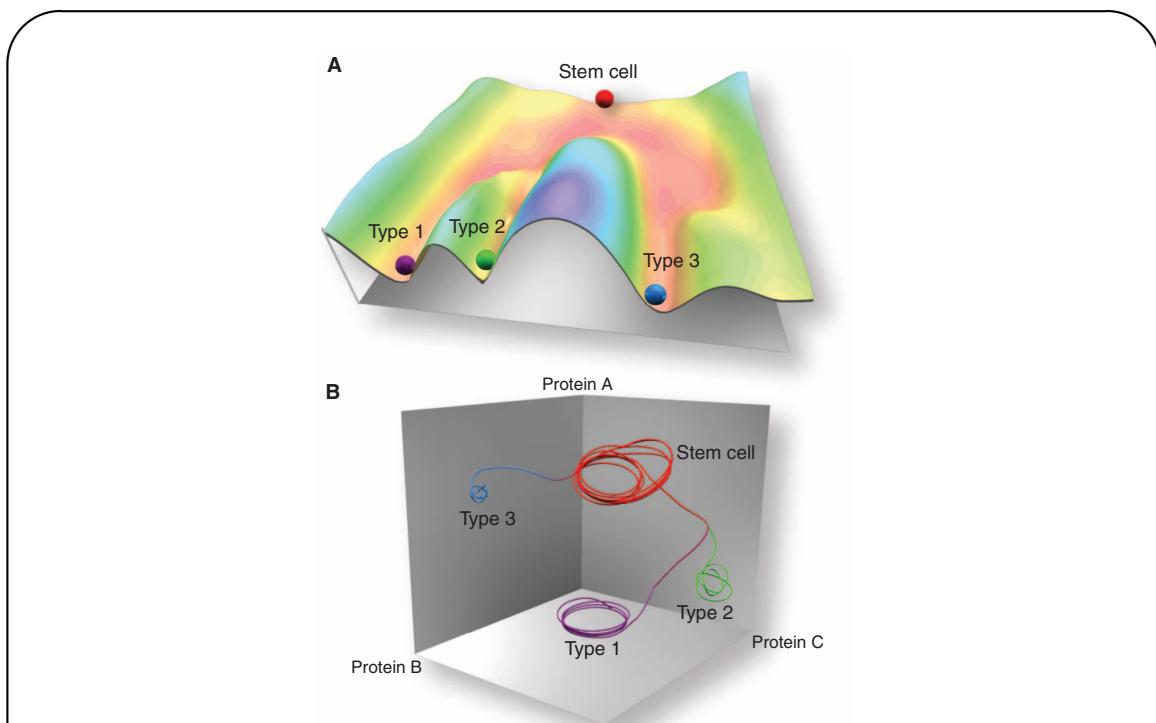


FIGURE 1.1 – Le paysage de la différenciation cellulaire.

Figure tirée de (?). **A.** Paysage épigénétique tel qu'imaginé par Waddington (?) en résonance avec la notion de paysage énergétique en physique. Le développement cellulaire est représenté par une bille dévalant un paysage composé de différentes vallées séparées par des barrières difficilement franchissables, représentant les différents types cellulaires et leur robustesse face aux fluctuations. **B.** Représentation dynamique de l'évolution des états cellulaires. Le phénotype est ici caractérisé par l'expression de trois protéines A, B et C, dont l'évolution dans le temps peut être représentée par une trajectoire dans un espace tridimensionnel. Les états souches et différenciés sont caractérisés par des bassins d'attraction correspondant à différents types cellulaires.

1.1.2 La différenciation cellulaire

L'acquisition d'un phénotype cellulaire particulier au sein d'un organisme est le sujet de la biologie du développement. Cette acquisition passe par différentes étapes de différenciation cellulaire. Schématiquement, au cours du développement d'un organisme, des cellules non différenciées ou souches empruntent un chemin unidirectionnel de différenciation qui

1. Il existe cependant quelques cas connus d'organismes multicellulaires procaryotes, dont les cellules ne possèdent pas de noyau, par exemple chez les bactéries magnétotactiques (?).

Chapitre 1. Introduction générale.

reste最基本的 peu à peu le nombre de types cellulaires qu'elles peuvent potentiellement devenir, passant de l'état souche totipotent à des états pluripotents successifs avant la différenciation finale. Ainsi, la formation des cellules somatiques, qui sont les cellules d'un organisme n'étant ni souches ni germinales (les cellules qui donnent naissance aux gamètes ou cellules sexuelles), est le résultat d'un processus de différenciation initial lors du développement embryonnaire au cours duquel les cellules souches issues de l'œuf donnent naissance à trois couches de tissus distinctes : l'endoderme (feuillet interne), l'ectoderme (feuillet externe) et le mésoderme (feuillet intermédiaire). Des différenciations successives ont ensuite lieu au sein de ces couches pour former divers organes tels que le tube digestif (endoderme), les muscles et les os (mésoderme), ou encore la peau et le système nerveux (ectoderme).

Dans un écrit aujourd'hui célèbre datant de 1957 (?), Waddington proposa une représentation de ces différentes étapes sous la forme d'un paysage épigénétique semblable aux paysages énergétiques dont sont coutumiers les physiciens (fig 1.1A). Dans cette représentation, le processus de différenciation cellulaire est comparé à une bille dévalant une pente et dont la trajectoire suit les multiples embranchements de vallées escarpées, chacune représentant un état de développement différent. Les vallées sont séparées par des pics dont la hauteur reflète la difficulté de passer d'un état à un autre, et les destinations finales possibles de la bille correspondent aux différents types cellulaires.

La notion de trajectoire de différenciation peut être rendue plus parlante en adoptant une représentation de système dynamique. Comme nous l'avons vu en 1.1.1, la cellule contient de nombreux composants : gènes, protéines ou autres métabolites, qui pris dans leur ensemble déterminent à un instant donné l'état cellulaire. Il est ainsi possible de représenter l'état cellulaire à un temps donné comme un point dans un espace de grande dimension dans lequel chaque axe représente l'abondance d'un certain composant (fig 1.1B). De par leur rôle primordial dans la définition de l'état cellulaire, l'expression des protéines (et donc des gènes qui les produisent) domine généralement ces composants, et on parle de « niveau d'expression génétique » pour décrire leur abondance. Les changements d'expression génétique, au cours desquels certains gènes vont être activés et d'autres réprimés, induisent un changement de l'état cellulaire, ce qui se traduit par une trajectoire dans l'espace des états. Ces changements d'expression restreignent finalement l'état cellulaire à une certaine région, définie comme un « attracteur » de la dynamique. Une fois au sein d'un attracteur, l'état cellulaire est robuste aux perturbations du niveau d'expression génétique des différentes composantes. Les attracteurs

peuvent alors être vu comme des types cellulaires distincts correspondant aux différentes vallées de la représentation de Waddington (?).

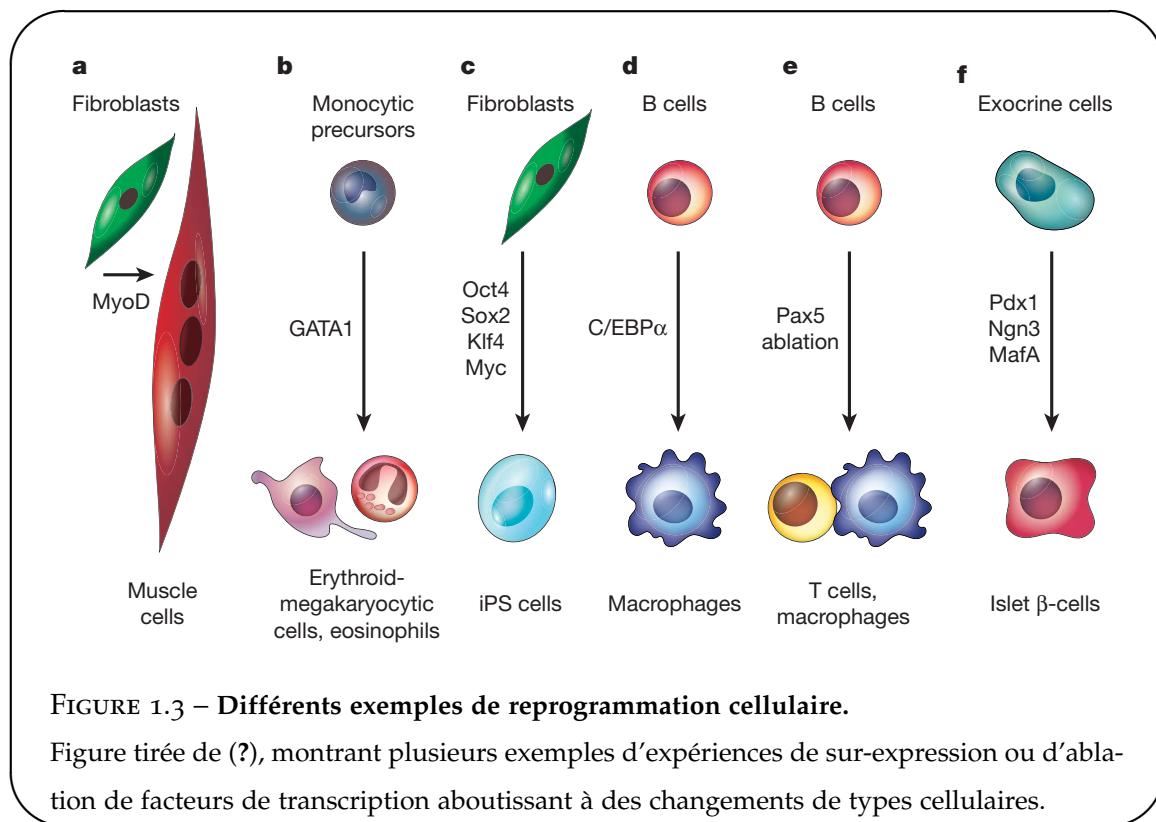
1.1.3 La cellule dans l'organisme : une spécification spatio-temporelle



FIGURE 1.2 – Spécification spatio-temporelle du type cellulaire.

Hybridation *in situ* de l'ARN du gène *Myog*, marqueur de la différenciation des progéniteurs du muscle squelettique, chez des embryons de souris âgés de 9.5, 10.5 et 11.5 jours (de gauche à droite), observés sous un même grossissement de 10. Le motif (*pattern*) de spécification du muscle squelettique est clairement visible au niveau des somites, les futures vertèbres. Images tirées de la base de donnée Embrys (<http://embrys.jp>).

Au sein de l'organisme, la différenciation cellulaire opère à un rythme précis et dans un contexte cellulaire bien défini. Aussi, les trajectoires dans l'espace d'expression génétique que nous avons présentées précédemment sont fonction de l'espace – la position de la cellule dans l'organisme, qui détermine en particulier la concentration des signaux qu'elle reçoit de son environnement – et du temps – le stade de développement de l'organisme –. Il est ainsi possible d'observer chez l'embryon certains motifs ou *patterns* spatio-temporels d'expression génétique correspondant à des organes en formation et révélés par la hybridation *in situ* de l'ARN de certains gènes spécifiques d'un type cellulaire. Par exemple, dans le cas de la formation des muscles squelettiques, le gène de différenciation terminale *Myog* est exprimé chez la souris dès 8 jours embryonnaires au niveau des somites, segments correspondant aux futures vertèbres de la souris adulte, puis commence à être exprimé au niveau des bourgeons de membres à 11.5 jours (voir fig 1.2).



1.1.4 La reprogrammation cellulaire

Depuis plusieurs décennies, différentes expériences ont exhibé la plasticité des états différenciés, élargissant ainsi considérablement la vision classique selon laquelle des cellules souches totipotentes se différencient de manière irréversible en des cellules de moins en moins plastiques, jusqu'à atteindre un état différencié stable. Par exemple, (?) ont montré en 1985 que des programmes d'expression génétique dormants peuvent être exprimés de manière dominante dans des cellules différencierées par la fusion de différents types cellulaires : ainsi, la fusion de cellules musculaires avec des cellules non musculaires permettait l'activation des gènes de type musculaire dans le type cellulaire non musculaire. Puis différents travaux ont montré qu'il était possible de convertir des lignées de cellules différencierées en un autre type cellulaire en introduisant certaines protéines régulatrices de la transcription, ou Facteurs de Transcription (TFs) (??) : on parle alors de transdifférenciation, dont l'un des exemples canoniques est la différenciation de cellules de la peau ou fibroblastes en cellules musculaires par l'introduction du facteur de différenciation myogénique MyoD (voir fig 1.3). Parallèlement, des expériences réalisées chez plusieurs espèces de mammifères ont montré que le trans-

fert de noyaux de cellules différenciées embryonnaires ou adultes dans un oeuf énucléé peut mener à la formation d'un organisme complet, montrant de manière univoque que l'identité des cellules différenciées peut être complètement renversée (?). Enfin, l'avancée la plus récente dans ce domaine a été la démonstration que des cellules somatiques différenciées peuvent être reprogrammées en cellules souches puripotentes par simple introduction d'un « cocktail » de 4 facteurs de transcription : Oct4, Sox2, c-Myc et Klf4 (?) (fig 1.3C).

1.2 Les réseaux de régulation génétique

Afin de pouvoir mieux comprendre les mécanismes de différenciation et de reprogrammation exposés en 1.1, il convient de se plonger dans les mécanismes internes de la cellule qui régissent ses changements d'états.

1.2.1 Vision cybernétique de la cellule

Le paradigme qui règne sur la biologie moléculaire depuis plus d'un demi siècle est celui des réseaux génétiques. L'expression des gènes est en effet régulée par des protéines, les facteurs de transcription, qui sont eux-mêmes issus de l'expression d'autres gènes, créant ainsi un réseau d'interactions entre gènes. Certaines protéines peuvent par ailleurs directement réguler l'activité d'autres protéines, et certains ARNs issus de la transcription de gènes non codants jouent aussi un rôle fondamental dans la régulation de l'activité génétique, le tout formant un réseau complexe d'interactions à différents niveaux. La compréhension de ce réseau et des fonctions qui en résultent forme le socle de la biologie des systèmes. Dans ce cadre, la cellule est vue comme une unité de traitement d'information qui interprète différents signaux reçus en entrée, les traite par un réseau interne de régulation, et réagit en sortie en modifiant son état ou son comportement (fig 1.4). L'intérêt d'une telle description mécanistique est qu'elle permet d'opérer quantifications mathématiques et prédictions, ce qui l'a rendue extrêmement fertile au cours des dernières décennies (?).

1.2.2 Divers modes de régulation

Les modes de régulation qui permettent à la cellule d'interpréter des signaux afin de changer d'état sont nombreux. Nous allons nous concentrer ici sur ceux affectant la production d'ARNs ou de protéines (fig. 1.5).

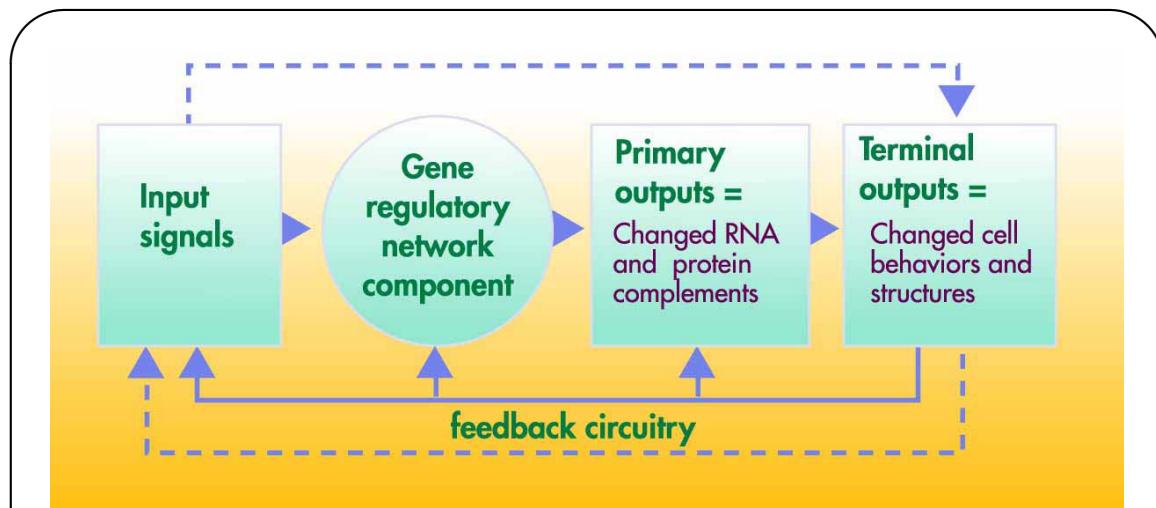


FIGURE 1.4 – Vision cybernétique du traitement de l’information par la cellule.

Figure tirée du programme ‘Genomes to life’ du Département de l’Énergie des États-Unis datant de 2001 (?) schématisant un réseau de régulation cellulaire comme un système de traitement entrée/sortie, possédant trois composantes fondamentales : (1) un système de réception et de transduction des signaux d’entrées qui peuvent être intra- ou extra-cellulaires (plusieurs signaux pouvant affecter un même gène cible), (2) un « composant central » (*core component*) composé du réseau de régulation génétique traitant les signaux, et (3) d’un signal de sortie consistant en l’expression moléculaire des ARNs et protéines des gènes cibles. Le processus résulte en la modification du phénotype de la cellule. Des boucles de régulation (*feedback*) assurent le contrôle et la stabilité des différentes étapes.

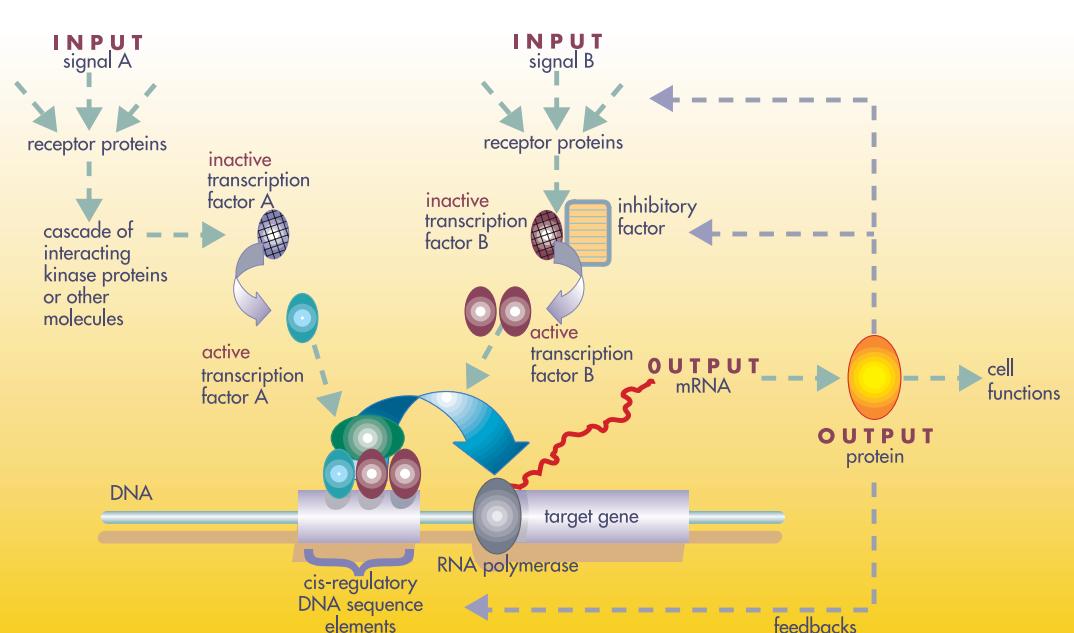


FIGURE 1.5 – Un réseau de régulation génétique type.

Dans cette représentation schématique tirée du rapport du ?, deux voies de signalisation A et B transmettent des signaux d'entrée (qui peuvent être intra ou extra cellulaires) en rendant des facteurs de transcription actifs. Une fois activés, ces derniers interagissent avec des séquences d'ADN proches d'un gène cible en se fixant sur des sites de petite taille ($\sim 10\text{bp}$). Les différents facteurs de transcription interagissent entre eux pour former des complexes occupant des régions de $\sim 1000\text{bp}$ appelées modules de cis-régulation ou CRMs (voir section 1.5). Lorsque les facteurs de transcription sont fixés sur le CRM de leur gène cible, il peuvent activer ou inhiber la transcription d'ARN et donc la production de la protéine correspondante.

- **Régulation génétique**

Tout d'abord, un réseau d'expression génétique est caractérisé par un jeu d'interactions entre différents gènes. Ces interactions se font par l'intermédiaire de protéines régulatrices appelées facteurs de transcription ou TFs, qui sont au nombre de ~ 1400 chez l'homme (?), soit 6% des protéines encodées. Les gènes qui les expriment représentent donc $\sim 3\%$ de l'ensemble des 30,000 gènes connus à ce jour. Pour réguler (activer ou inhiber) la transcription d'un gène cible, les TFs se fixent sur des sites de reconnaissance spécifiques sur l'ADN de $\sim 10\text{bp}$ et interagissent avec la machinerie transcriptionnelle au niveau du promoteur du gène cible. Les TFs peuvent se fixer sur le promoteur même, comme c'est souvent le cas chez la bactérie, ou dans des régions distales allant jusqu'à plusieurs centaines de kb, comme on trouve plus couramment chez les organismes complexes. Par ailleurs, différents TFs peuvent se combiner sur certaines régions de régulation contenant de multiples sites de fixation pour former des complexes protéiques. Ces régions, appelées modules de cis-régulation (CRMs) ou plus communément *enhancers*, sont d'une taille typique de $\sim 1000\text{bp}$ et ont la particularité de conduire à une expression spatio-temporelle très spécifique du gène cible. Ces différents points seront amplement développés en section 1.5.

- **Régulation épigénétique**

Outre la régulation génétique, due à l'action de protéines issues de séquences codantes et se fixant sur des séquences d'ADN – régulation qui est donc entièrement encodée dans le génome et transmise à la descendance –, il existe un autre mode de régulation de la transcription des gènes qui permet notamment d'acquérir une modification d'expression génétique transmise à la descendance sans qu'il y ait modification du code génétique : c'est la régulation épigénétique. Cette régulation passe notamment par la modification des propriétés chimiques de l'ADN et des histones sur lequel il s'enroule pour former la chromatine (fig. 1.6). Ainsi, la méthylation des dimères CpG de l'ADN² au niveau des régions riches en CG, ou îlots CpG, situées près de nombreux promoteurs et habituellement dépourvues de ces marques conduit à une inactivation du gène cible (?). Par ailleurs, la méthylation des histones au niveau des résidus lysines entraîne la fermeture de la chromatine, empêchant l'expression du ou des gène(s) situés à leur niveau, alors que l'acétylation des mêmes lysines entraîne au contraire une ouverture de la chromatine, favorisant ainsi la transcription génétique (?). Ce mode de

2. Les dimères C-G sont appelés CpG, où p caractérise le phosphore liant les deux bases, pour les différencier du CG utilisé pour parler de la statistique en C et G de l'ADN

régulation sera développé plus en détails en section 1.5.1.

- **Régulation post-transcriptionnelle**

Les modifications post-transcriptionnelles affectent les ARNs issus de la transcription des gènes. Ces modifications peuvent être causées par des microARNs ou miRNAs qui sont des ARNs de ~ 23 nts issus d'ARNs se repliant en structure double brin de type « épingles à cheveux » ou *hairpins*. Les miRNAs s'associent à la protéine *Argonaute* du complexe RISC (*RNA-induced silencing complex*) pour entraîner la dégradation spécifique d'ARNms (?). De manière similaire, certains *hairpins* de taille plus importante sont clivés par la protéine Dicer pour former plusieurs petits ARNs de taille similaire aux miRNAs : ce sont les siRNAs (*small interfering RNAs*). Ceux-ci recrutent aussi le complexe protéique RISC et ciblent spécifiquement des ARNm (??). Ce phénomène est connu sous le nom d'interférence ARN (RNAi) et a donné lieu à une méthode aujourd'hui couramment utilisée pour inhiber l'expression d'un gène.

- **Régulation post-traductionnelle**

Les modifications post-traductionnelles affectent les protéines issues de la traduction des ARNs. Elles passent par une modification chimique des protéines, typiquement la phosphorylation, ou comme nous l'avons vu pour la régulation épigénétique, la méthylation ou l'acétylation. Ces modifications peuvent avoir pour effet de changer l'activité de la protéine, que ce soit en modifiant son activité enzymatique ou en déclenchant sa relocalisation nucléaire. Il existe aussi des modifications de structure de la protéine, comme c'est le cas du facteur de transcription *Shavenbaby* chez la Drosophile : dans sa forme native, cette protéine inhibe la transcription de ses gènes cible ; cependant ses résidus terminaux peuvent être clivés par des petits peptides de 11 à 32 acides aminés encodés par le gène *Pri*, rendant la protéine transcriptionnellement active (?).

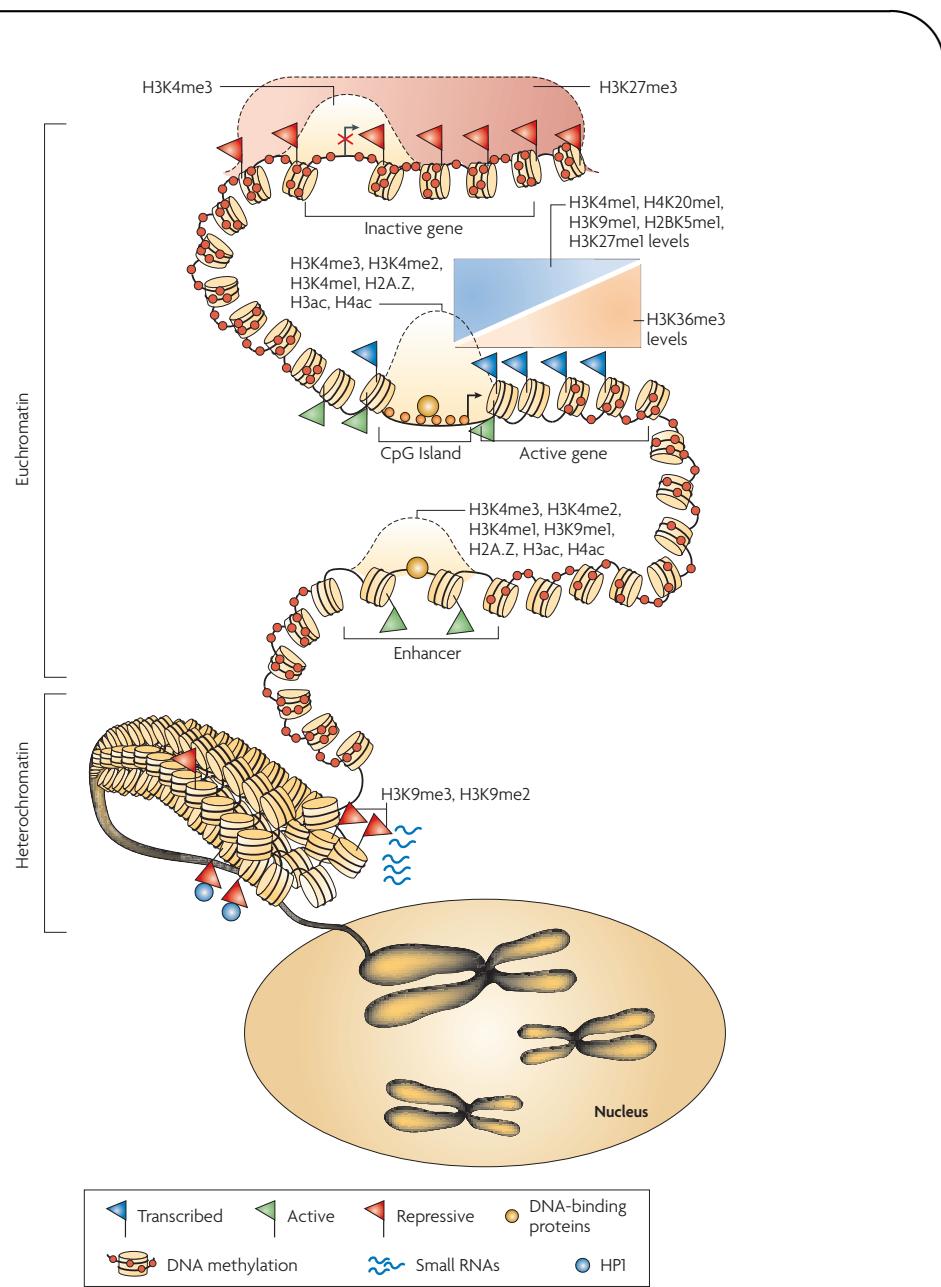
**FIGURE 1.6 – Caractéristiques de l'épigénome.**

Figure tirée de (?). Les chromosomes sont partagés entre régions accessibles d'euchromatine et régions difficilement accessibles d'hétérochromatine. Les régions hétérochromatiques sont marquées par la di- et triméthylation de la lysine 9 de l'histone H3 (H_3K9me2 et H_3K9me3). La méthylation de l'ADN est répandue à travers tout le génome, mais est absente de certaines régions comme les îlots CpG, les promoteurs et les CRMs. La modification $H_3K27me3$ couvre de larges régions englobant des gènes inactifs. Les marques H_3K4me3 , H_3K4me2 , H_3K4me1 et l'acétylation des histones marquent les TSSs des gènes actifs. Les marques H_3K4 , H_3K9 , H_3K27 , H_4K20 et H_2BK5 marquent les régions transcris activement à proximité de la région 5' des gènes (en amont), alors que la marque H_3K36 marque les gènes transcrits dans leur région 3' (en aval).

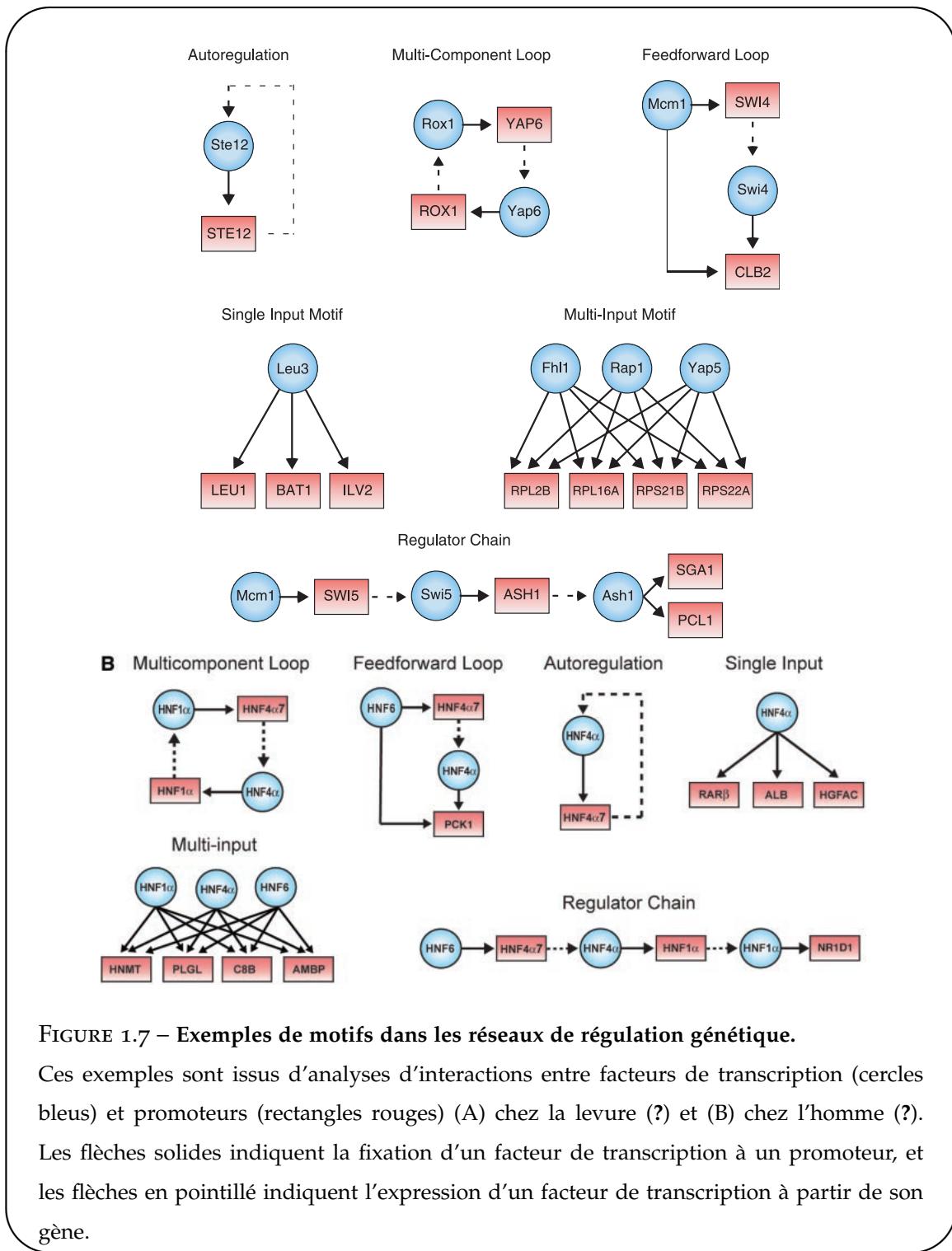


FIGURE 1.7 – Exemples de motifs dans les réseaux de régulation génétique.

Ces exemples sont issus d'analyses d'interactions entre facteurs de transcription (cercles bleus) et promoteurs (rectangles rouges) (A) chez la levure (?) et (B) chez l'homme (?). Les flèches solides indiquent la fixation d'un facteur de transcription à un promoteur, et les flèches en pointillé indiquent l'expression d'un facteur de transcription à partir de son gène.

1.2.3 Câblage du réseau et fonction

Maintenant que nous avons vu la nature des interactions au sein des réseaux génétiques, nous pouvons nous pencher sur leur structure. Notamment, plusieurs études réalisées chez divers organismes de la bactérie à l'homme ont révélé que les réseaux de transcription contiennent un petit ensemble de motifs de régulation récurrents, appelés motifs de réseaux (???) (fig. 1.7). Ces motifs peuvent être vus comme les pièces élémentaires servant à la construction de réseaux fonctionnels. De tels motifs furent d'abord détectés de manière systématique chez la bactérie *Escherichia coli* en remarquant qu'ils apparaissaient dans le réseau de transcription bien plus souvent qu'on ne l'attendrait dans un réseau aléatoire (?). Les mêmes motifs ont ensuite été trouvés chez la levure (??) et chez l'homme (?). La récurrence de ces motifs est liée aux fonctions qu'ils remplissent. Par exemple, la boucle d'autorégulation négative, qui est trouvée chez la moitié des répresseurs d'*Escherichia coli*, possède deux fonctions : l'une est de parvenir rapidement à un état d'équilibre en utilisant un promoteur fort, l'autre est de servir de tampon au bruit d'expression (?). Un autre motif récurrent est la boucle feedforward. Celle-ci consiste en 3 gènes : un régulateur X, qui régule Y, tous deux régulant Z. Dans le cas où des interactions sont des activations et que X et Y sont requis pour activer Z, cette boucle peut servir de tampon au bruit d'expression de X, évitant que des fluctuations de son niveau d'expression n'entraînent par erreur l'activation de Z.

1.2.4 Évolution des réseaux génétiques

L'importance des motifs est rendue plus claire lorsque l'on s'intéresse à l'évolution des réseaux. En effet, au cours de l'évolution, les réseaux de régulation génétique changent : modification des constituants, recâblage du réseau, duplication d'éléments... Néanmoins, certaines modifications sont plus défavorisées du point de vue évolutif que des autres. Ainsi, les motifs tels que les boucles d'autorégulation ou les boucles feedforward, de par leur importance fonctionnelle, auront tendance à être conservés. Pour ce qui est des éléments du réseau, la modification d'un régulateur, par exemple la mutation d'un acide aminé au sein d'un facteur de transcription, aura des conséquences sur l'ensemble des éléments régulés par ce facteur de transcription et pourra donc être fortement délétère. Par contre, la modification d'un site de reconnaissance de ce facteur de transcription sur l'ADN n'aura qu'une portée locale sur la régulation du gène associé.

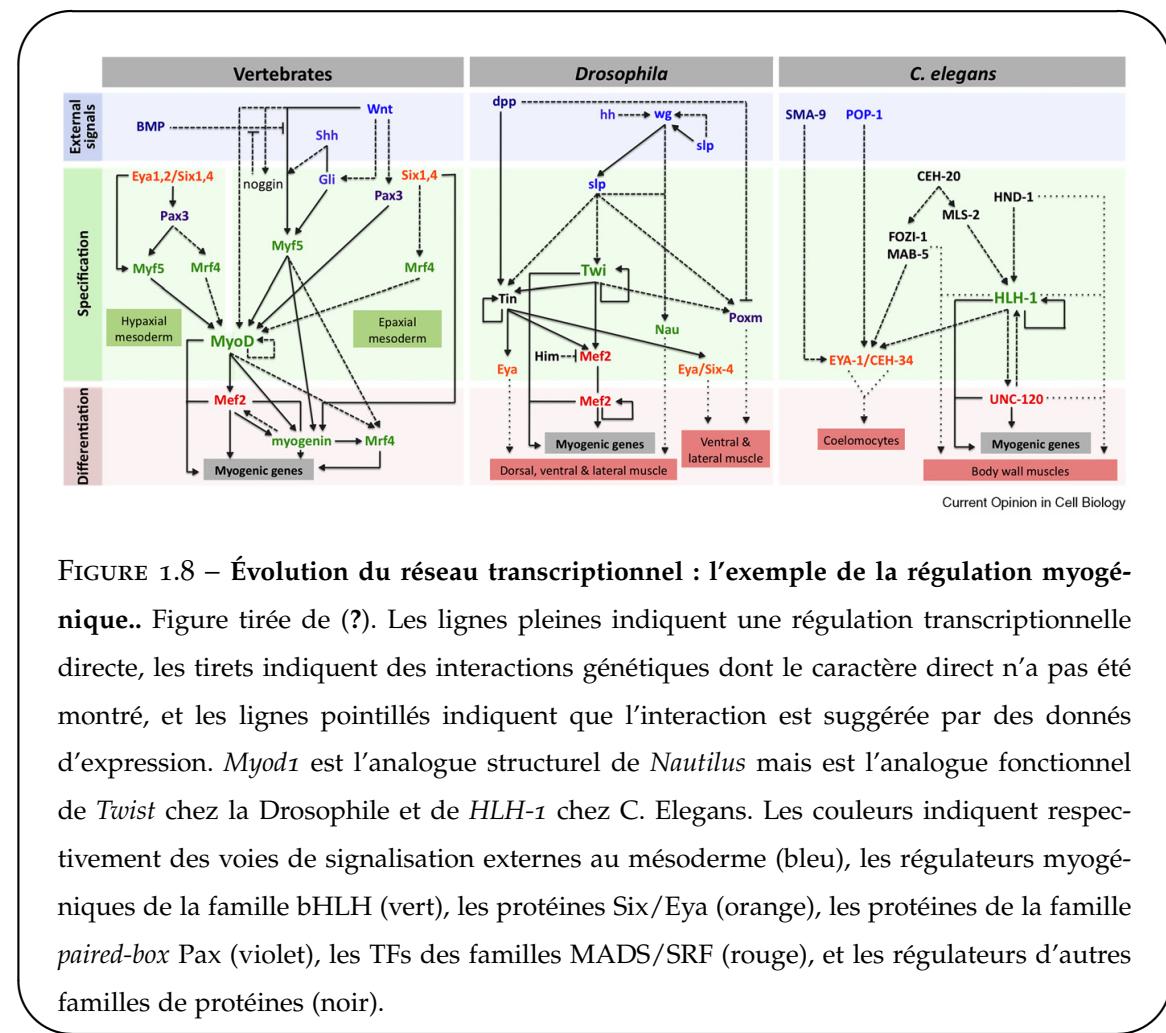


FIGURE 1.8 – Évolution du réseau transcriptionnel : l'exemple de la régulation myogénique.. Figure tirée de (?). Les lignes pleines indiquent une régulation transcriptionnelle directe, les tirets indiquent des interactions génétiques dont le caractère direct n'a pas été montré, et les lignes pointillés indiquent que l'interaction est suggérée par des données d'expression. *Myod1* est l'analogie structurelle de *Nautilus* mais est l'analogie fonctionnelle de *Twist* chez la Drosophile et de *HLH-1* chez *C. Elegans*. Les couleurs indiquent respectivement des voies de signalisation externes au mésoderme (bleu), les régulateurs myogéniques de la famille bHLH (vert), les protéines Six/Eya (orange), les protéines de la famille paired-box Pax (violet), les TFs des familles MADS/SRF (rouge), et les régulateurs d'autres familles de protéines (noir).

À titre d'exemple, prenons le cas du réseau de différenciation du muscle squelettique présenté en figure 1.8, que nous étudierons plus en détail dans le chapitre ?? de ce manuscrit. Au cœur de ce réseau génétique se trouvent les facteurs de régulation myogéniques ou MRFs, des facteurs de transcription de type bHLH qui ont la capacité de convertir des cellules non mesodermiques, c'est-à-dire n'étant pas destinées à devenir des progéniteurs musculaires, en cellules ayant des propriétés musculaires (?). Ces facteurs sont dits « régulateurs maîtres » de la différenciation musculaire. Chez les vertébrés il y a quatre MRFs : *Myf5*, *Mrf4*, *Myod1*, qui ont des rôles redondants dans la spécification des progéniteurs musculaires, et *Myog*, qui conduit à la différenciation terminale. Chez la Drosophile c'est le TF *Twist* qui semble être le principal MRF, mais contrairement aux MRFs des vertébrés, son rôle ne s'arrête pas au contrôle de la différenciation musculaire mais est plus général dans le développement du

Chapitre 1. Introduction générale.

mésoderme (?). C'est cependant le gène *Nautilus* qui possède la séquence d'acides aminés la plus proche de celle des MRFs vertébrés. Ce dernier permet la spécification des progéniteurs myogéniques, et son expression est restreinte au développement musculaire. Néanmoins, les mutants *nautilus* sont viables et son rôle semble mineur comparé aux MRFs vertébrés. Enfin, chez le ver *Caenorhabditis elegans*, c'est l'orthologue de *Myod1*, *hlh-1*, qui tient rôle de MRF.

Malgré ces différences (nombre de MRFs, membre de la famille bHLH tenant ce rôle), on retrouve dans les trois cas une boucle feedforward conservée au niveau de la régulation des cibles des MRFs (fig. 1.8). Ainsi, MyoD régule l'expression de Mef2 et l'activité de MAPK p38 en même temps que l'expression de plusieurs cibles initiales, et par la suite MyoD et phospho-Mef2 co-régulent des gènes plus tardifs. Ce mécanisme permet ainsi de réguler l'aspect temporel de l'expression génétique. Chez la Drosophile, le même motif est observé avec Twist et Mef2 et chez *C. Elegans* avec HLH-1 et le TF UNC-129, de la même famille que Mef2.

Le cœur du réseau est donc conservé dans la forme (topologie), même s'il y a des divergences dans le fond (membres de la famille de TFs impliqués). Néanmoins, les éléments régulateurs en amont, ainsi que les membres périphériques du réseau ont rapidement évolué. Par exemple, chez les vertébrés le TF Pax3 est très en amont dans la hiérarchie génétique et permet l'activation des MRFs et la spécification myogénique, alors que chez la Drosophile son homologue *poxm* est en aval des MRFs et sa perte de fonction n'a que des effets mineurs sur la myogenèse. Par ailleurs, le complexe composé de protéine Six et de leur cofacteur Eya, initialement découvert comme régulateur majeur de la différenciation oculaire chez la Drosophile, est chez les vertébrés un régulateur essentiel situés en amont des MRFs. Chez la Drosophile, il possède aussi un rôle dans la spécification myogénique, mais bien plus en aval que chez les vertébrés. Enfin, chez *C. Elegans* ce complexe est aussi en aval des MRFs mais il participe en plus à la détermination de cellules non myogéniques.

Nous voyons donc que l'évolution d'un réseau génétique possède de multiples facettes : conservation de motifs de réseau fonctionnellement importants (dans notre exemple, la boucle feedforward au cœur du réseau régissant l'aspect temporel de l'expression des cibles), recâblage des interactions pour traiter différents signaux d'entrée... Par ailleurs, il apparaît que plus qu'à des TFs particuliers, c'est à des familles de TFs que nous avons affaire. Aussi un

même rôle au sein du réseau peut-il être rempli par différents membres d'une même famille, comme c'est le cas pour *Myod1* et *Twist*. Ceci s'explique par le fait que les membres d'une même famille partagent des propriétés d'interaction avec l'ADN semblables. Ces interactions sont à la source du fonctionnement du réseau, et nous allons maintenant présenter plus en avant leurs propriétés.

1.3 Les interactions protéine-ADN : modèles mathématiques

Nous l'avons vu, les interactions entre facteurs de transcription et ADN sont une composante essentielle des réseaux génétiques. Les TFs se fixent sur des sites spécifiques de ~ 10 bp dans le voisinage des gènes qu'ils régulent. Trouver ces sites est donc un premier pas vers la reconstruction des réseaux de régulation sous-jacents. Dans cette section nous présentons les modèles d'interactions protéine-ADN qui ont été proposés, et leur application concrète à la recherche de sites de fixation.

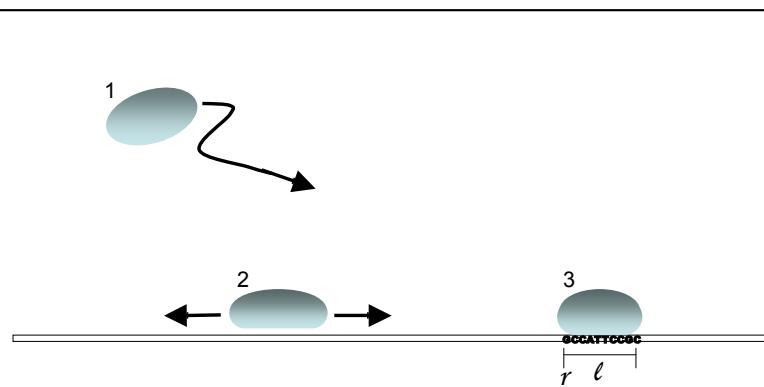


FIGURE 1.9 – Différents états du facteur de transcription. Figure tirée de (?). Lors de sa recherche de site de fixation, le TF peut se trouver dans trois états distincts : (1) un état libre de diffusion tridimensionnelle, (2) un état de diffusion unidimensionnelle sur l'ADN par fixation non spécifique, et (3) un état de fixation spécifique. L'énergie de fixation dépend du site de fixation, de taille l et de coordonnée r .

1.3.1 Modes de recherche du site de fixation par le TF

Un facteur de transcription peut être dans plusieurs états : en diffusion tridimensionnelle, auquel cas il est dit « libre », ou bien fixé sur l'ADN. Dans ce dernier cas, il interagit avec

Chapitre 1. Introduction générale.

l'ADN selon deux modes : une attraction non spécifique d'énergie E_{ns} indépendante de la position sur l'ADN, et une interaction spécifique $E_s(r)$ qui dépend de la séquence de taille $l \sim 10$ à la position r sur l'ADN. L'interaction non spécifique est due à l'interaction électrostatique entre la protéine chargée positivement et l'ADN chargé négativement, alors que l'interaction spécifique implique des liaisons hydrogènes entre le domaine de fixation de la protéine et les nucléotides du site de fixation. La protéine passe d'un mode à l'autre en changeant de conformation. Le facteur de transcription peut ainsi se trouver dans trois états thermodynamiques représentés en figure 1.9 : en diffusion tridimensionnelle libre, fixé non spécifiquement (diffusion unidimensionnelle le long de la structure d'ADN), et fixé spécifiquement sur l'ADN. Ces trois modes contribuent à la cinétique de la recherche d'un site fonctionnel (??). Ainsi, l'attraction non spécifique conduit la protéine à passer à peu près autant de temps fixé sur l'ADN qu'en diffusion libre. La recherche de site de reconnaissance est donc un processus mixte de diffusion unidimensionnelle sur l'ADN et de diffusion tridimensionnelle dans le milieu. Lorsqu'il est fixé sur l'ADN, le facteur diffuse dans un paysage d'énergie E_{ns} plat lorsqu'il est dans sa conformation de fixation non spécifique, ou dans un paysage d'énergie $E_s(r)$ dans sa conformation de fixation spécifique. Cela permet au facteur d'échantillonner les sites de faible énergie $E_s(r)$ tout en évitant d'être bloqué par les barrières de haute énergie en passant en mode de recherche non spécifique. Ce processus s'avère au final très efficace : les temps de recherche sont typiquement inférieurs à une minute, ce qui est petit devant les processus de régulation de la cellule qui se déroulent au mieux sur quelques minutes (??). Il est donc pertinent de décrire l'effet d'un site de fixation sur la régulation d'un gène cible par la probabilité qu'il a de fixer un facteur de transcription à l'équilibre thermodynamique.

1.3.2 Modèle PWM

Présenté en 1987 par Berg et von Hippel (?), le modèle PWM est le modèle le plus simple décrivant l'énergie de fixation spécifique entre un facteur de transcription et un site de fixation sur l'ADN. Ce modèle repose sur plusieurs hypothèses. Tout d'abord, il y a l'hypothèse importante que les sites de fixation des TFs sur l'ADN ont été sélectionnés au cours de l'évolution pour leur propriété de sites de reconnaissance, quelle que soit la concentration du TF dans la cellule. En d'autres termes, le processus de sélection discrimine les sites de fixation sur la seule base de leur énergie de fixation à un TF donné : les sites ayant une énergie de fixation dans une certaine gamme sont retenus, les autres rejettés. Par ailleurs, au sein de cette gamme

 1.3. *Les interactions protéine-ADN : modèles mathématiques*

d'énergie « utile », toutes les séquences sont équiprobables. Enfin, la dernière hypothèse est que chaque nucléotide d'un site de fixation contribue de manière indépendante, c'est-à-dire additive à l'énergie totale du site. Cette hypothèse permet de simplifier le problème en gardant le nombre de paramètres petit.

L'argument de Berg et von Hippel est que ce problème est analogue à celui de physique statistique consistant à déduire les taux d'occupation des niveaux d'énergie de particules indépendantes sachant que l'énergie totale doit avoir une certaine valeur moyenne E . La solution de ce problème est donnée par la formule de Boltzmann reliant énergie et taux d'occupation :

$$f_{i,b} = \exp(-\lambda E_{i,b}) / \mathcal{Z}_i \quad (1.1)$$

où $f_{i,b}$ est la probabilité d'observer la base b à la position i du site de fixation, $E_{i,b}$ est l'énergie associée (en $k_B T$), \mathcal{Z}_i est la fonction partition qui permet de normaliser la distribution à la position i , et λ est un facteur sans dimension, analogue du β de la thermodynamique, et lié au processus de sélection. Dans la suite, nous intégrerons ce facteur à l'énergie.

La connaissance des fréquences des bases permet de définir une autre quantité utile caractérisant la variabilité des séquences de fixation, l'information relative des sites par rapport à une séquence d'ADN aléatoire (?) :

$$\mathcal{I} = \sum_{i=1}^L \sum_{b=A,C,G,T} f_{i,b} \ln \left(\frac{f_{i,b}}{\pi_b} \right) \quad (1.2)$$

où L est la taille du site de fixation et π_b correspond à la probabilité *a priori* d'observer la base b dans le génome. Parce que l'énergie est définie à une constante près, il est usuel de la définir relativement au fond génomique :

$$\tilde{E}_{i,b} = \ln \left(\frac{f_{i,b}}{\pi_b} \right) \quad (1.3)$$

L'énergie totale d'un site S_i est alors

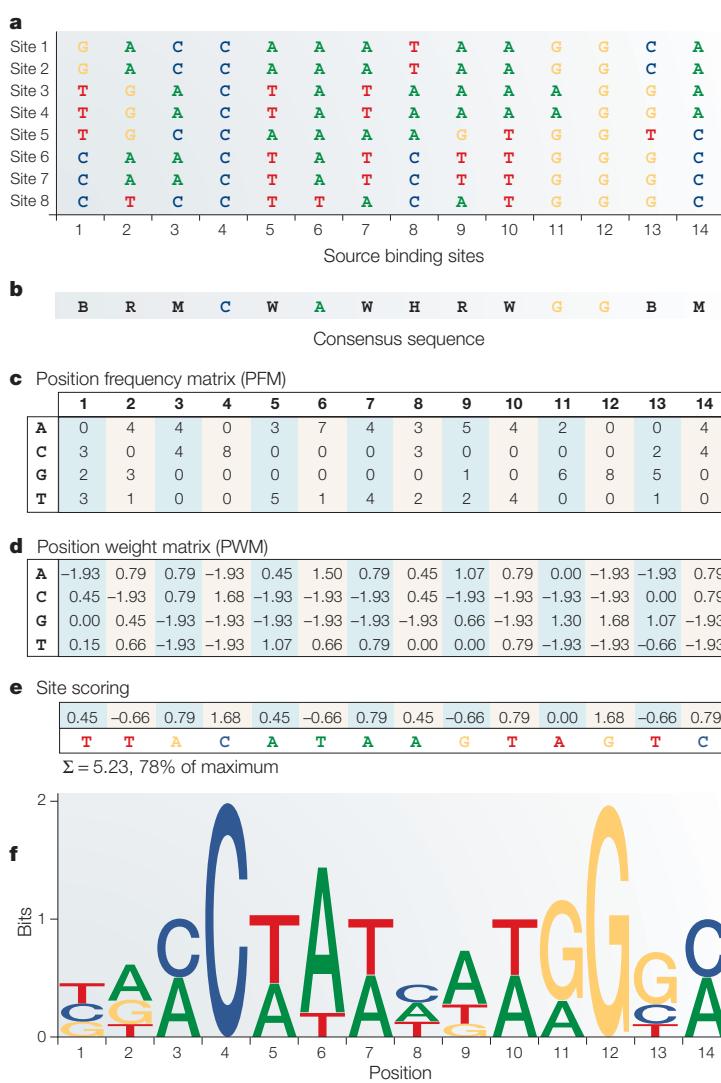


FIGURE 1.10 – Construction et utilisation du modèle PWM. Figure tirée de (?). (a) Supposons connus un certain nombre de sites de fixation d'un facteur de transcription (dans ce cas MEF2). (b) Séquence consensus correspondante utilisant les symboles IUPAC. (c) Une matrice de fréquence est construite, indiquant pour chaque nucléotide sa multiplicité à une position donnée dans le site. (d) La PWM est simplement construite en prenant le logarithme relatif des fréquences PWMs par rapport aux fréquences *a priori* des nucléotides. (e) Le score (ou énergie) d'une séquence d'ADN donnée est calculé en additionnant les poids PWM correspondants. (f) La PWM peut être représentée sous forme de logo (?). Dans cette représentation, la hauteur d'une colonne représente le contenu en information ou information relative moyenne d'une position, et la taille des bases reflète leur fréquence.

$$\begin{aligned}
E &= \sum_{i=1}^L \tilde{E}_{i,b} \\
&= \sum_{i=1}^L \ln \left(\frac{f_{b(i)}}{\pi_b} \right) \\
&= \ln \left(\frac{\prod_{i=1}^L f_{b(i)}}{\prod_{i=1}^L \pi_b} \right) \\
&= \ln \left(\frac{P(S_i|\text{TF})}{P(S_i|\text{fond génomique})} \right)
\end{aligned} \tag{1.4}$$

où $b(i)$ est la base située à la position i du site de fixation. Cette énergie quantifie simplement à quel point la séquence S_i est plus ($E > 0$) ou moins ($E < 0$) probablement un site de fixation (de probabilité $P(S_i|\text{TF})$) qu'un site tiré au hasard dans le génome (de probabilité $P(S_i|\text{fond génomique})$). On parle aussi de *score* de la séquence. L'information relative \mathcal{I} , qui est le score moyen des séquences fixées par le TF, peut alors être vue comme quantifiant à quel point l'ensemble des sites de fixation se distingue d'un ensemble de même taille de sites tirés au hasard.

Avec ces outils en main, il devient alors simple de bâtir un modèle PWM et de l'utiliser pour prédire des séquences fixées (fig. 1.10). Étant donnés des sites de fixation connus, il suffit d'évaluer la fréquence d'occurrence de chaque base à chaque position. La comparaison avec les probabilités génomiques *a priori* d'occurrence permet alors de bâtir une matrice de score, la PWM. Cette matrice peut alors être utilisée pour attribuer un score à une séquence d'ADN en additionnant les scores à chaque position. Finalement, les séquences ayant un score dépassant un certain seuil sont considérées comme des séquences de fixation.

1.3.3 Modèle biophysique

Le modèle PWM est basé sur une hypothèse forte, celle que les sites de fixation ont été sélectionnés sur la base de leur seule affinité ou énergie envers un TF. Néanmoins, à aucun moment n'intervient la concentration du TF dans la cellule, dont dépend pourtant la probabilité de fixation. C'est ce que tente de capturer le modèle biophysique (???).

Considérons l'interaction entre un TF et une séquence d'ADN S_i :



où $TF : S_i$ dénote le complexe entre le TF et le site S_i . La constante d'équilibre de cette réaction s'écrit selon la loi d'action de masse :

$$K_i = \frac{[TF : S_i]}{[TF][S_i]} \quad (1.6)$$

Le site peut être dans deux états : occupé par le TF ou libre. Aussi, la probabilité que le TF soit fixé au site s'écrit simplement

$$P(\text{fixation}|S_i) = \frac{[TF : S_i]}{[TF : S_i] + [S_i]} = \frac{1}{1 + \frac{1}{K_i[TF]}} = \frac{1}{1 + e^{\beta(E_i - \mu)}} \quad (1.7)$$

où $E_i = -kT \ln(K_i)$ est l'énergie libre standard de fixation (souvent notée ΔG), $\mu = kT \ln[TF]$ est le potentiel chimique, k est la constante de Boltzmann, T la température et $\beta = 1/kT$. Ici nous avons considéré qu'il n'y avait qu'un seul site de fixation. De manière générale, le site est en compétition avec le fond génomique, ce qui ajoute une contribution à μ (voir section 1.3.4). À l'instar du modèle PWM, l'énergie E_i est généralement prise comme étant une fonction additive des énergies individuelles des différentes bases du site. Ainsi, lorsque le TF est à faible concentration ($\mu \rightarrow -\infty$), le modèle biophysique écrit en équation 1.7 se réduit au modèle PWM.

1.3.4 Modèle thermodynamique

La description biophysique peut être réécrite en termes thermodynamiques en utilisant des raisonnements simples sur le nombre d'états possibles et leur énergie (et donc poids de Boltzmann) associée. Nous adoptons ici l'approche de (?). On pourra par ailleurs se référer à l'excellente revue (?). Considérons le cas simple d'un seul facteur de transcription interagissant avec un génome de taille $L \gg 1$ ne contenant qu'un seul site fonctionnel, le reste de la séquence étant aléatoire. Nous l'avons vu, l'expérience montre que la protéine se fixe à l'ADN avec une probabilité 1/2. Lorsqu'elle est fixée, elle est à l'équilibre entre le mode spécifique et le mode non spécifique. Nous désirons savoir avec quelle probabilité elle est fixée de manière spécifique. La fonction de partition, énumérant tous les poids de Boltzmann associés aux différents états accessibles au TF fixé s'écrit :

$$\mathcal{Z} = \sum_{r=1}^L e^{-\beta E_s(r)} + L e^{-\beta E_{ns}} \quad (1.8)$$

Notons i la position du site fonctionnel. On peut écrire :

$$\begin{aligned}\mathcal{Z} &= e^{-\beta E_s(i)} + e^{-\beta E_{ns}} + \sum_{r \neq i} e^{-\beta E_s(r)} + (L-1)e^{-\beta E_{ns}} \\ &\simeq e^{-\beta E_i} + \mathcal{Z}_0\end{aligned}\tag{1.9}$$

où Z_0 est la fonction de partition d'une séquence aléatoire, et nous avons introduit l'énergie E_i définie par

$$e^{-\beta E_i} = e^{-\beta E_s(i)} + e^{-\beta E_{ns}}\tag{1.10}$$

Dans le cas d'un site de reconnaissance, $E_{ns} \gg E_s(i)$ de sorte que $E_i \simeq E_s(i)$ (?). La probabilité que le facteur soit fixé sur le site fonctionnel s'écrit finalement :

$$P(\text{fixation spécifique}|E_i) = \frac{e^{-\beta E_i}}{\mathcal{Z}} = \frac{1}{1 + e^{\beta(E_i - F_0)}}\tag{1.11}$$

où $F_0 = -kT \log Z_0$ est l'énergie libre d'une séquence génomique aléatoire. On reconnaît une fonction de Fermi, avec un seuil d'énergie à F_0 : pour $E_i < F_0$, la protéine est essentiellement fixée de manière spécifique à son site de reconnaissance, alors que pour $E_i > F_0$, elle ne distingue plus le site du fond génomique et y est faiblement fixée.

Généralisons à présent au cas de plusieurs facteurs de transcription et sites de reconnaissance. Nous négligeons le recouvrement entre facteurs de transcription fixés sur des sites proches, qui poserait des problèmes stériques et corrèlerait les sites de fixation dans un certain voisinage (la présence d'un TF empêchant la présence d'un autre), et considérons que le nombre de TFs est grand devant le nombre de sites de reconnaissance pour éviter les problèmes de saturation : ainsi, le génome est composé de L séquences indépendantes, chacune pouvant être soit non occupée, soit occupée de manière non spécifique, soit occupée de manière spécifique. Notons μ le potentiel chimique du TF en solution. La fonction de partition totale est le produit des fonctions de partition des sites indépendants,

$$\mathcal{Z}(\mu) = \prod_{r=1}^L \mathcal{Z}(\mu, r)\tag{1.12}$$

où la fonction de partition d'un site s'écrit :

$$\mathcal{Z}(\mu, r) = e^{-\beta\mu} + e^{-\beta E_s(r)} + e^{-\beta E_{ns}}\tag{1.13}$$

Chapitre 1. Introduction générale.

En utilisant à nouveau la définition de E_i en éq. 1.10, la probabilité de fixation d'un site à la position i s'écrit

$$P(\text{fixation spécifique}|E_i) = \frac{e^{-\beta E_i}}{\mathcal{Z}(\mu, i)} = \frac{1}{1 + e^{\beta(E_i - \mu)}} \quad (1.14)$$

La valeur de μ est liée à la fois au nombre de TFs ainsi qu'à la possibilité de se fixer dans le fond génomique. Elle est fixée implicitement par l'équation :

$$n = \sum_{r=1}^L \frac{1}{1 + e^{\beta(E_r - \mu)}} \quad (1.15)$$

qui signifie simplement que le nombre de TFs n dans le système est égal à la somme sur tous les sites de fixation possibles pondérée par la probabilité que le TF y soit fixé. Lorsque $\mu \rightarrow -\infty$ et que la fonction de Fermi peut être approximée par la loi de Boltzmann, l'équation peut s'inverser et l'on trouve (?)

$$\mu = F_0 + kT \log n \quad (1.16)$$

où F_0 est l'énergie libre du fond génomique introduite en éq. 1.11. Ainsi, la prise en compte d'une multiplicité de TFs ajoute un facteur $kT \log n$ au seuil de la fonction de Fermi par rapport au cas d'un seul TF. Par ailleurs, cette approche thermodynamique nous a permis de généraliser le modèle biophysique simple introduit en section 1.3.3.

1.4 Les interactions protéine-ADN : mesures expérimentales

Ces dernières années, des avancées technologiques considérables ont permis d'une part d'établir des modèles de fixation spécifique pour de nombreux TFs, d'autre part de localiser leurs sites de fixation dans le génome. Ces avancées ont eu lieu autant sur le plan *in vitro*, utilisant protéines purifiées et séquences nucléiques artificielles pour déduire l'affinité protéine-ADN, que sur le plan *in vivo*, mesurant l'interaction de la protéine avec l'ADN génomique (?).

1.4.1 Approches *in vitro* : MITOMI, SPR, PBM, CSI, SELEX, et HT-SELEX

- **Approche microfluidique : MITOMI**

En 2007, Maerkl et Quake ont mis au point une technique appelée MITOMI (Mechanically Induced Trapping Of Molecular Interactions) permettant une mesure directe de l'affinité d'un TF à des centaines de séquences d'ADN à la fois (?). Cette technique repose sur l'utilisation d'un système microfluidique composé de chambres dans lesquelles un fluide dont on peut facilement modifier la composition circule dans des canaux d'un diamètre de l'ordre de $1\mu\text{m}$ dont le microenvironnement est finement contrôlé. Le fluide contient des gènes synthétiques codant pour le TF ainsi que du matériel permettant la synthèse de la protéine directement au sein de la chambre, ce qui évite de purifier préalablement le TF. Chaque chambre du système contient des anticorps fixés à la surface permettant de capturer le TF et une certaine concentration d'une séquence d'ADN spécifique contenant une marque fluorescente. Le système contient ainsi des centaines de séquences d'ADN différentes, chacune étant présente à différentes concentrations. Lorsque le TF est fixé par les anticorps, il recrute des séquences d'ADN selon leur affinité. Celles qui ne se fixent pas sont lavées. Au final, les séquences fixées produisent un signal de fluorescence. La comparaison des signaux pour différentes concentrations d'ADN donne accès au rapport des constantes d'équilibre K_{eq} (eq. 1.6). La comparaison avec une séquence référence dont la constante K_{eq} est connue permet alors de déterminer le K_{eq} absolu pour chaque séquence de fixation.

En utilisant 17 systèmes de ce type, ils ont ainsi pu mesurer l'affinité de 4 TFs de type bHLH à 464 séquences d'ADN différentes : les séquences consensus et des séquences ayant une, deux, trois ou quatre mutations. À titre de comparaison, ils ont construit une PWM à partir des séquences contenant une seule mutation, puis ont prédit les énergies attendues des séquences à plusieurs mutations. La prédiction de la PWM s'est avérée bonne dans seulement 56% des cas pour les séquences à deux mutations, 10% pour les séquences à 3 mutations et 0% des cas pour les séquences à 4 mutations, montrant les limites de ce modèle indépendant confronté à des données d'interactions d'ordre supérieur. Un modèle plus raffiné prenant en compte l'énergie d'interaction non spécifique et incluant des interactions entre nucléotides voisins permet néanmoins de rendre compte des valeurs observées (?). Nous reviendrons sur la nécessité de prendre en compte les interactions entre paires de nucléotides lors de

l'interaction spécifique entre TF et ADN dans le chapitre 2.

- **Approche physique : la microscopie SPR**

La méthode de résonance des plasmons de surface (*Surface Plasmon Resonance* ou SPR) est habituellement utilisée pour étudier l'interaction d'une protéine avec un ligand (qui peut être une autre protéine), mais elle peut aussi être utilisée pour mesurer les interactions entre une protéine et quelques centaines de séquences d'ADN différentes (??). Le principe de la microscopie SPR est que l'angle de réflexion de la lumière sur une fine surface d'or, par exemple, dépend de la masse de molécules fixées de l'autre côté de sa surface. Si de l'ADN est lié à la surface, la fixation du TF induit un changement de masse et donc d'angle de reflection lumineuse mesurable au cours du temps. Ainsi, la cinétique de fixation du TF jusqu'à l'atteinte de l'équilibre est accessible. On peut de même étudier la dissociation du TF lors du lavage de la surface. Ces mesures donnent directement accès aux taux d'association k_{on} et de dissociation k_{off} que la simple mesure de la constante d'équilibre $K_{eq} = k_{on}/k_{off}$ ne permet habituellement pas de déterminer.

- **Approches basées sur des puces à ADN : PBM et CSI**

L'analyse de fixation des protéines par puce à ADN (*Protein-Binding Microarray* ou PBM) est une technologie haut débit qui a été développée au cours des 10 dernières années (?). Les puces sont composées de 44,000 puits auxquels sont liés des brins d'ADN. Une puce contient tous les sites de fixation de 8bp possibles ($4^8/2 = 32,768$ séquences en prenant en compte le fait qu'il y a un site sur chacun des deux brins d'ADN) plus deux bases flanquantes (une à chaque extrémité) qu'il est possible de faire varier. Un TF purifié à partir de cellules ou synthétisé *in vitro* est ajouté à la puce, qui est ensuite lavée pour se débarrasser des fixations non spécifiques. La quantité de protéine fixée à un puits donné est déterminée grâce à un anticorps fluorescent contre la protéine. L'enrichissement en protéine est calculé relativement au bruit de fond (anticorps non spécifique par exemple). Il est alors possible d'utiliser ces mesures pour bâtir une PWM du TF (voir par exemple ?).

Une autre méthode utilise aussi des puces à ADN : c'est l'identification de site apparenté (*Cognate Site Identifier* ou CSI) (?). Une différence technique avec les PBMs est que l'ADN est d'abord synthétisé en simple brin puis se replie en double brin pour former le site de fixation, évitant ainsi de devoir générer l'ADN double brin à partir de précurseurs. Par ailleurs, le TF

 1.4. *Les interactions protéine-ADN : mesures expérimentales*

est en compétition avec un marqueur fluorescent qui peut se fixer à l'ADN : il n'est donc pas nécessaire d'utiliser un marquage spécifique sur le TF ou sur un anticorps, ce qui rend la procédure plus généralisable. Finalement, la spécificité du TF est représentée par un « paysage de spécificité » qui encapsule l'information de fluorescence de l'ensemble des variations par rapport à une séquence consensus dans une représentation simple (?).

- **Approche par purification des séquences fixées : SELEX et HT-SELEX**

Mise au point il y a plus de 20 ans, la méthode SELEX (*Systematic Evolution of Ligands by EXponential enrichment*) repose sur la sélection de séquences d'ADN aléatoires par un TF *in vitro* (????). Une bibliothèque de sites de fixation potentiels est d'abord générée en synthétisant des séquences d'ADN aléatoires ou en utilisant des séquences génomiques. Les extrémités de ces séquences contiennent des précurseurs permettant l'amplification exponentielle par PCR. Le TF purifié est ajouté aux sites et les séquences fixées sont séparées des séquences non fixées, par exemple par retard sur gel. Après un cycle de sélection, les séquences récupérées sont encore enrichies en séquences de basse affinité pour le TF, car celles-ci sont simplement initialement bien plus abondantes que les séquences de haute affinité. Afin d'augmenter la proportion de séquence de grande affinité, les séquences filtrées sont amplifiées puis filtrées à nouveau, ceci sur plusieurs cycles. À la fin de ce processus, les séquences sélectionnées sont clonées et séquencées, résultant en un nombre typique de moins de ~ 100 séquences indépendantes (?). Si les séquences initiales sont issues d'ADN génomique, il est possible d'utiliser l'hybridation des séquences à des puces à ADN. La présence de plusieurs cycles de sélection rend néanmoins la détermination des énergies de fixation moins directe qu'avec les techniques précédentes. Une variante de la technique appelée SELEX-SAGE utilise des multimères de sites à la place de sites uniques et permet de réduire le nombre de cycles de sélection et d'augmenter ainsi le nombre de séquences de fixation obtenues (?), permettant de réaliser des modèles plus précis (?).

Depuis la mise au point de la méthode SELEX, des avancées considérables ont été réalisées dans les techniques de séquençage, permettant l'obtention de millions de séquences à la fois : on parle de séquence haut-débit (*high-throughput*) ou encore séquençage massivement parallèle. L'utilisation de ces nouvelles techniques dans l'expérience SELEX a mené à la méthode HT-SELEX (?), aussi appelée Bind-n-Seq (?). Il est alors possible d'estimer un modèle d'énergie à partir des fréquences d'observation des différentes séquences dès le premier cycle (?). Des

Chapitre 1. Introduction générale.

cycles supplémentaires permettent d'obtenir plus d'information sur les séquences les plus spécifiques, notamment sur la présence de contributions non indépendantes à l'énergie, ou de compenser la faible spécificité d'un TF. L'avantage de cette technique est que la taille des sites de fixation n'est pas limitée. Ainsi, avec une nanomole d'ADN ($\sim 10^{15}$ séquences) on peut couvrir l'ensemble des sites de 25bp possibles. Le séquençage haut-débit permet d'en échantillonner $\sim 10^8$, ce qui est largement suffisant pour contraindre des modèles d'énergie indépendants, même pour des TFs ayant des sites de fixations de taille > 15 bp comme c'est souvent le cas chez la bactérie. Cette technique a récemment été poussée encore plus loin (?). En utilisant des protéines marquées, les auteurs ont réalisé un HT-SELEX à partir d'extraits cellulaires, et en ajoutant un code barre aux séquences d'ADN de chaque expérience, ils ont pu analyser les sites de fixation pour plusieurs TFs en parallèle. Ils ont ensuite utilisé cette technique pour obtenir des modèles de spécificité pour 411 TFs humains, la plus grande étude de ce genre réalisée à ce jour (?).

1.4.2 Approche clonale : la technique de simple hybride

Contrairement aux approches précédentes, la technique de simple hybride (*Bacterial one-hybrid* ou B1H) n'est pas purement *in vitro*, au sens où l'interaction protéine-ADN est testée au sein d'une bactérie. Néanmoins, parce que l'interaction n'est pas testée dans son contexte cellulaire d'origine, nous la considérerons comme telle. Cette approche repose sur l'intégration par une bactérie hôte de deux vecteurs d'expression génétique, ou plasmides. Le premier exprime le facteur de transcription d'intérêt fusionné à une sous-unité de l'ARN polymérase (l'appât), c'est la protéine « hybride ». L'autre contient une région de séquence aléatoire représentant un site de fixation potentiel (la proie) en amont d'un promoteur à faible activité. La fixation de cette région par la protéine hybride permet l'activation d'un gène de sélection, généralement *HIS3*, un gène de la levure requis pour la biosynthèse de l'histidine et dont l'homologue bactérien est absent de la souche d'*Escherichia coli* utilisée. La croissance des cellules a lieu dans un milieu ne contenant pas l'histidine. Dans ces conditions, les bactéries n'exprimant pas *HIS3* ne peuvent croître. Ainsi, seules les bactéries au sein desquelles le facteur de transcription se fixe à la proie expriment *HIS3*, croissent et forment des colonies, d'où la notion de gène de sélection. Par ailleurs, la stringence de la sélection peut être modulée en ajoutant au milieu différentes concentrations de 3-amino-triazole (3-AT), un inhibiteur de *HIS3*. De cette façon l'affinité du site de fixation peut être estimée plus finement.

Dans les études de ce type, les sites de fixation présents au sein des colonies sont séquencés individuellement, ce qui permet d'obtenir environ 50 séquences pour une expérience de sélection donnée. Néanmoins, il semble possible d'utiliser les nouvelles technologies de séquençage pour récupérer l'ensemble des sites de fixation des bactéries présentes sur une plaque (?). À l'instar de la méthode HT-SELEX, on obtient des millions de sites, ceux ayant une plus grande affinité étant présents à plusieurs centaines de milliers d'exemplaires, et ceux ayant une faible affinité étant présent en un seul voire aucun exemplaire.

Notons qu'il est aussi possible d'adopter la démarche inverse, c'est-à-dire de partir de quelques sites de fixation présumés fonctionnels mais pour lesquels on ne connaît pas le TF associé. En utilisant une bibliothèque de plasmides codant pour différents TFs hybrides, il est alors possible de déterminer si l'un d'entre eux possède une affinité importante avec les sites testés.

1.4.3 Approches *in vivo* : ChIP-on-chip, ChIP-seq, DNase I

Dans cette section, nous nous intéressons aux techniques permettant d'identifier les sites de fixation d'un facteur de transcription sur le génome. Ces méthodes se basent sur des extraits cellulaires (de 10^4 à 10^8 cellules) qui peuvent provenir d'un tissu homogène (un seul type de cellule) ou hétérogène (plusieurs types de cellules), voire de l'organisme entier si la dissection est impossible (embryon de mouche par exemple). L'information obtenue est donc toujours conditionnée par ce matériau de départ, et l'on n'obtient que les sites *accessibles* étant donnés le type cellulaire et la période de développement étudiés.

- **Immunoprécipitation de la chromatine : ChIP-on-chip et ChIP-seq**

La technique d'immunoprécipitation de la chromatine (ChIP) (fig. 1.11) consiste dans un premier temps à induire la réticulation (*crosslink*) des protéines se liant à l'ADN en traitant les cellules avec de la formaldéhyde. Cette étape permet de transformer les liaisons faibles protéine-ADN en liaisons covalentes. Une fois les protéines fixées, la chromatine est découpée par digestion enzymatique ou en la soumettant à des ultrasons (c'est la sonication), résultant en des fragments de taille variant entre 200 et 600bp. Ces fragments sont ensuite immunoprécipités en présence d'un anticorps spécifique d'un facteur de transcription ou d'un isoforme

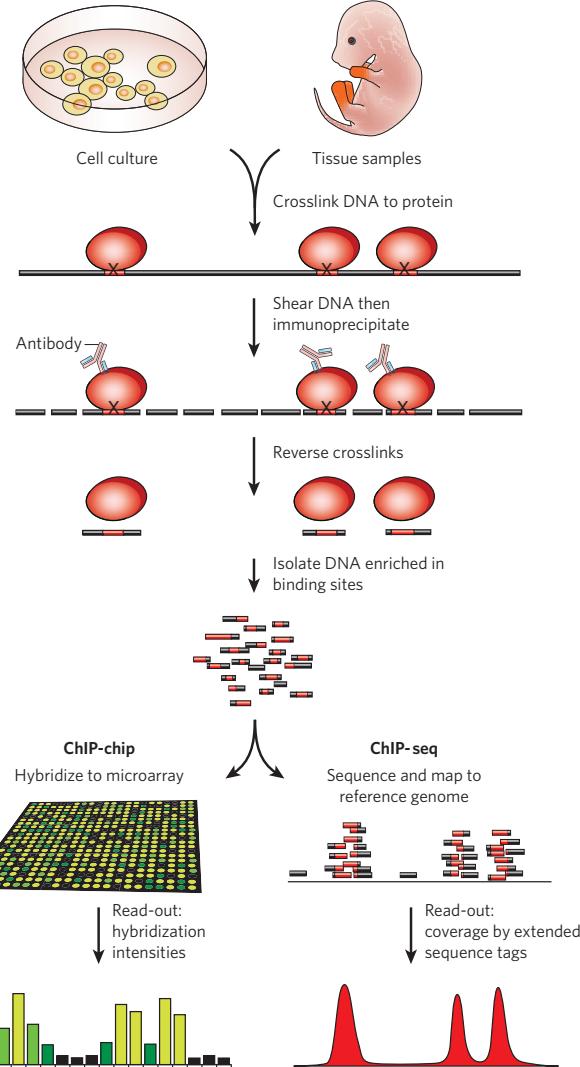


FIGURE 1.11 – Étapes d'une expérience de ChIP-on-chip et ChIP-seq.

Figure tirée de ?. À partir d'extraits cellulaires issus de cultures *in vitro* ou prélevés *in vivo*, plusieurs étapes permettent de récupérer les régions fixées par un TF d'intérêt. D'abord, les protéines, liées de manière non covalente à l'ADN, sont fixées par utilisation d'un agent de réticulation, le formaldéhyde. Puis la chromatine est découpée en fragments de ~ 500bp par utilisation d'ultrasons : on parle de sonication. L'utilisation d'un anticorps spécifique au TF d'intérêt permet de précipiter les fragments d'ADN fixés. Les fragments résultants sont alors soit hybridés à une puce à ADN contenant de nombreuses séquences génomiques (ChIP-on-chip), soit directement séquencés et alignés à un génome de référence (ChIP-seq). La comparaison à un échantillon d'ADN non précipité (« input ») permet alors de définir les régions significativement fixées, avec une résolution plus précise dans le cas du ChIP-seq, où l'ensemble du génome est couvert, par rapport au ChIP-on-chip, pour lequel seulement un sous ensemble du génome est analysé.

1.4. Les interactions protéine-ADN : mesures expérimentales

d'histone (dans le cas d'une étude du paysage épigénétique) d'intérêt, permettant ainsi de récupérer tous les sites de fixation dans le génome. Après purification des fragments précipités, l'échantillon peut être analysé soit par hybridation sur puce (ChIP-on-chip) ou par séquençage haut débit (ChIP-seq).

Dans le cas du ChIP-on-chip, l'échantillon immunoprecipité et l'ADN de départ (*input*) sont marqués avec des colorants fluorescents et hybriderés sur une puce à ADN composée de très nombreux puits contenant des oligonucléotides (courtes séquences d'ADN) correspondant à différentes régions du génome. Dans le meilleur cas, ces oligonucléotides couvrent l'ensemble du génome. Les sites de liaison sont identifiés par l'écart d'intensité entre les signaux de fluorescence des conditions d'immunoprecipitation et d'*input*.

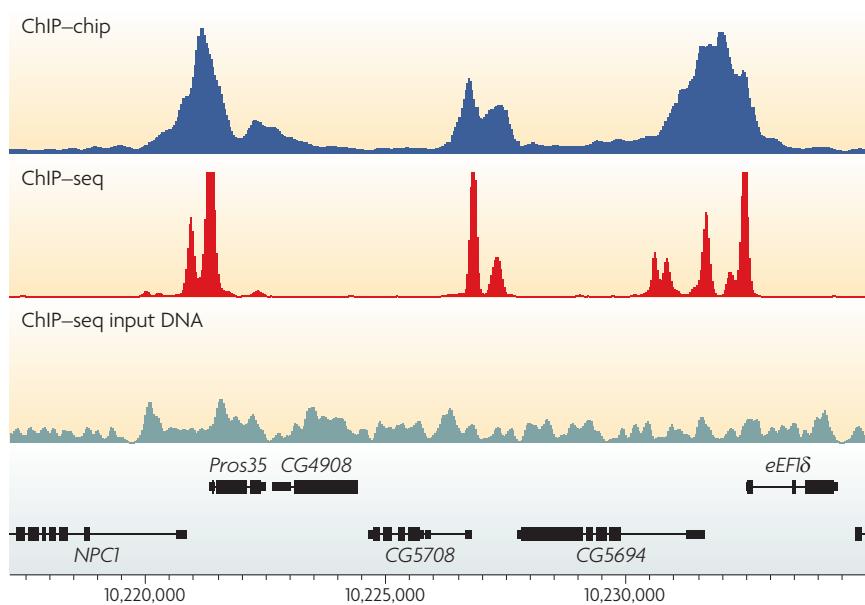


FIGURE 1.12 – Résolution des expériences ChIP-on-chip et ChIP-seq. Figure tirée de (?), montrant les profils de fixation de la protéine Chromator générés à partir d'expériences de ChIP-on-chip (intensité relative par rapport au contrôle, bleu) et de ChIP-seq (densité de séquences, rouge) dans la lignée cellulaire S2 de *Drosophila melanogaster*. On peut noter la plus grande résolution de l'expérience ChIP-seq pour déterminer les sites de liaison. L'ADN utilisé en *input* de l'expérience de ChIP-seq et servant de contrôle est montré en gris, et les gènes du locus indiqués en noir.

Dans le cas du ChIP-seq, l'échantillon immunoprecipité est analysé par séquençage à haut débit, résultant en une librairie de *reads* d'une longueur typique variant entre 27 et 50bp issus des extrémités des séquences. Ces *reads* sont ensuite alignés sur un génome de référence. À chaque position du génome correspond ainsi un certain nombre de séquences précipitées et d'*input*. En comparant ce nombre au nombre moyen dans le locus et à l'*input*, il est possible d'identifier des pics correspondant à la fixation du facteur (voir par exemple le programme d'appel de pics ChIP-seq MACS (?)).

Dans les deux cas, il faut noter que l'on a affaire à la fixation *moyenne* du facteur sur l'ADN dans la population de cellules étudiée. Ainsi, un petit pic peut représenter aussi bien une fixation forte dans un petit sous-ensemble de cellules (par exemple celles qui sont à un certain état d'avancement du cycle cellulaire) qu'une fixation moyenne dans l'ensemble de la population. L'expérience de ChIP-seq offre une résolution bien plus précise ($\leq 100\text{bp}$) que la méthode ChIP-on-chip (fig. 1.12). En effet, dans ce dernier cas la résolution est limitée par le nombre d'oligonucléotides utilisés, qui sont dans le meilleur des cas répartis sur le génome avec 35 – 100 nucléotides d'écart entre deux instances. Pour se comparer à la ChIP-seq, il faudrait que tous les oligonucléotides se superposent à une base près, ce qui demanderait un trop grand nombre de puces.

- **Empreinte à la DNase I (*DNase I footprinting*)**

Contrairement aux techniques précédentes, l'empreinte à la DNase I ne repose pas sur l'étude d'un facteur de transcription précis, mais permet au contraire d'obtenir l'ensemble des sites de fixation dans le génome pour un type cellulaire donné, avec une précision au nucléotide près. Cette méthode repose sur le fait que la fixation stable des facteurs de transcription à l'ADN n'est possible que si la région est pauvre en nucléosomes, les protéines autour desquelles s'enroule l'ADN : on parle de région de chromatine ouverte. Ces régions sont préférentiellement digérées par l'endonucléase DNase I. Étant donné que la majorité de l'ADN est enroulé autour de nucléosomes, les sites hypersensibles à la digestion par DNase I (*DNase I-hypersensitive* ou DHS) correspondent essentiellement à des régions de chromatine ouverte ayant des rôles de régulation génétique : promoteurs, enhancers...

En combinant la technique de DHS avec le séquençage à haut débit, l'expérience de DNase-

1.4. Les interactions protéine-ADN : mesures expérimentales

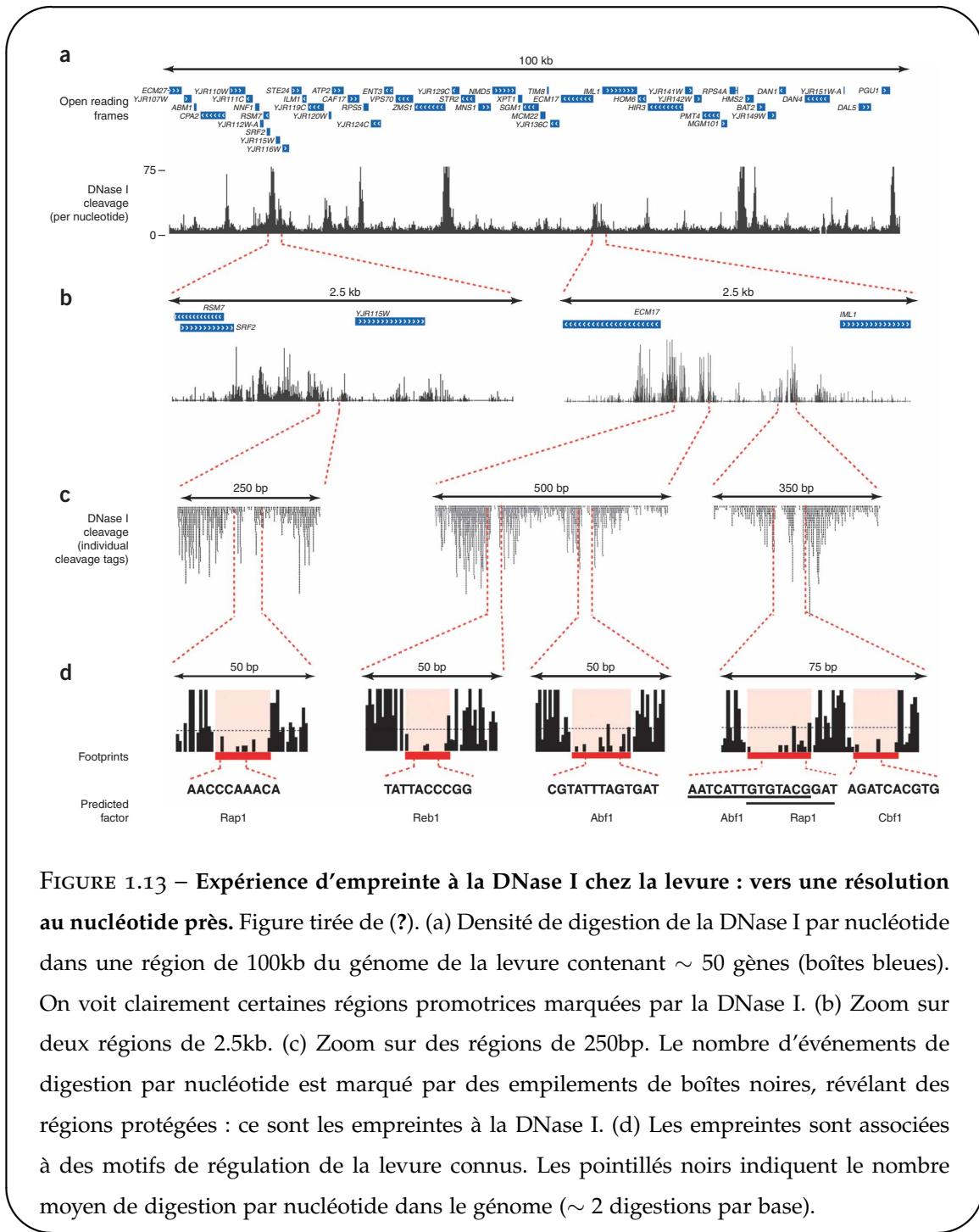


FIGURE 1.13 – Expérience d'empreinte à la DNase I chez la levure : vers une résolution au nucléotide près. Figure tirée de (?). (a) Densité de digestion de la DNase I par nucléotide dans une région de 100kb du génome de la levure contenant ~ 50 gènes (boîtes bleues). On voit clairement certaines régions promotrices marquées par la DNase I. (b) Zoom sur deux régions de 2.5kb. (c) Zoom sur des régions de 250bp. Le nombre d'événements de digestion par nucléotide est marqué par des empilements de boîtes noires, révélant des régions protégées : ce sont les empreintes à la DNase I. (d) Les empreintes sont associées à des motifs de régulation de la levure connus. Les pointillés noirs indiquent le nombre moyen de digestion par nucléotide dans le génome (~ 2 digestions par base).

seq permet d'identifier tous les types de région de régulation à l'échelle du génome (?). Les régions riches en sites de digestion identifient alors les sites DHS. Par ailleurs, au sein d'un site DHS, il y a de petites régions ($\sim 15\text{bp}$) qui sont protégées de la digestion par DNase I : ce sont les empreintes à la DNase I ou *DNase I footprints* (fig. 1.13). Ces empreintes sont dues à la présence de protéines ou de complexes fixés à l'ADN. Cette technique de détection de sites de liaison par empreinte à la DNase I existe depuis 30 ans mais n'a que récemment été porté à l'échelle génomique. En comparant à des données ChIP-seq ou en utilisant des bases de données de motifs de facteurs de transcription, il est possible d'identifier le facteur correspondant dont les sites de fixation sont alors connus au nucléotide près.

1.5 Les modules de cis-régulation (CRMs)

Nous l'avons vu en section 1.2.2, les séquences d'ADN régulant l'expression génétique – CRMs pour *Cis-Regulatory Modules* – jouent un rôle prépondérant au cours du développement des organismes. Ces CRMs assurent en effet l'orchestration de l'expression de gènes spécifiques aux différentes étapes du développement et aux divers types cellulaires. Ils sont au cœur de l'évolution des réseaux génétiques, car ils dictent les interactions entre gènes. De plus, leur altération peut conduire à de nombreuses pathologies, liées pour la plupart à une expression génétique aberrante. Notamment, la majeure partie des variants génétiques qui sont associés de manière significative à une susceptibilité envers une maladie sont situés hors des régions codant pour des protéines, suggérant qu'un certain nombre affectent non pas la forme de la protéine engendrée mais l'expression du gène la produisant en détruisant une activité CRM. Dans cette partie, nous présentons les différents types de CRMs, leur structure, et leur évolution.

1.5.1 Les différents types de CRMs

Selon leur rôle dans la régulation de l'expression génétique, les CRMs peuvent être distingués en trois catégories.

- **Promoteurs**

Les promoteurs permettent la fixation de l'ARN polymérase pour débuter la formation d'un transcrit ARN au site d'initiation de transcription (*Transcription Start Site* ou TSS). Dans

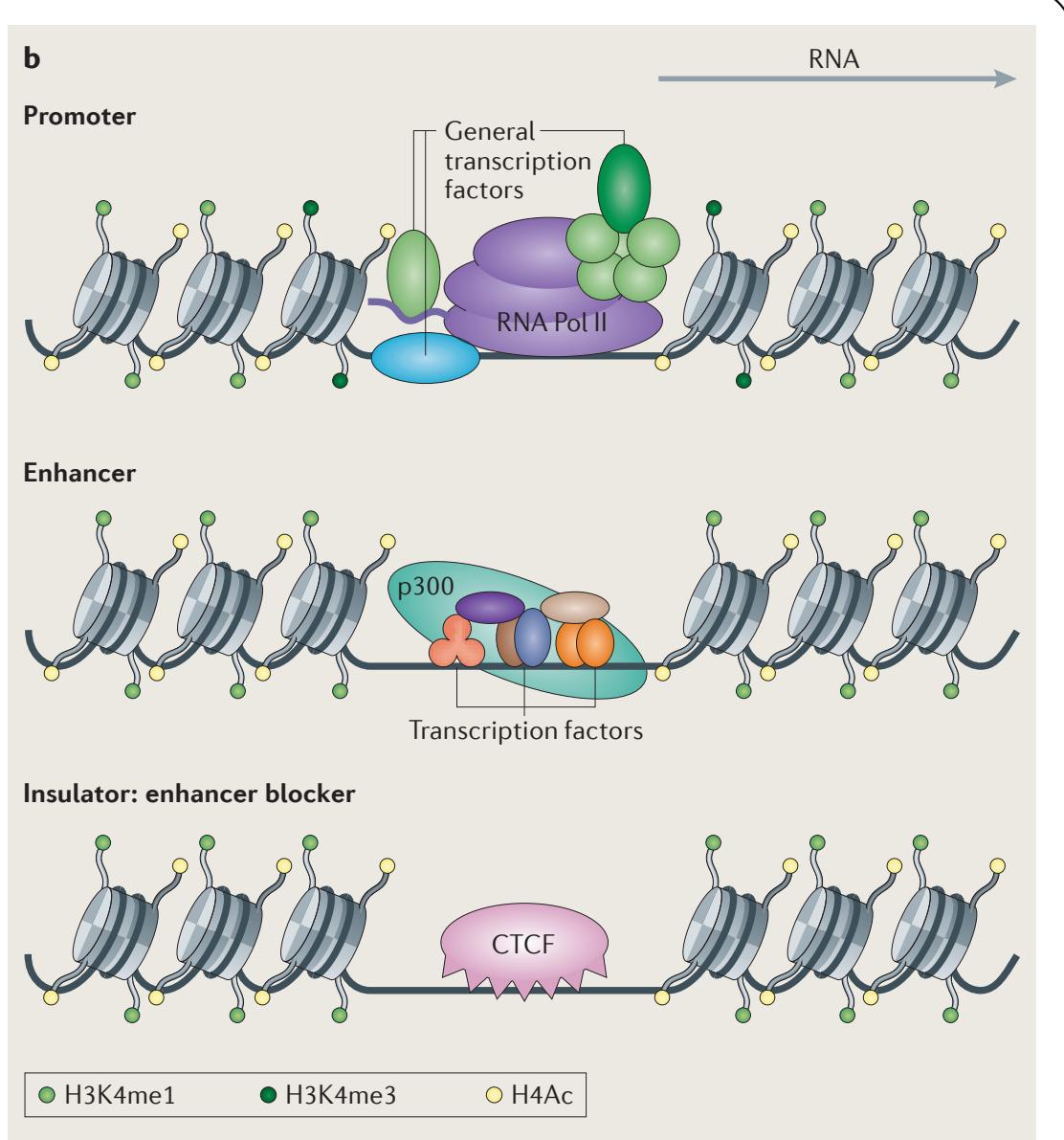


FIGURE 1.14 – Les différents types de CRMs et leurs marques épigénétiques.

Figure tirée de (?). La notion de CRM renvoie à un regroupement de sites de liaison pour un ou plusieurs facteurs de transcription. Les CRMs peuvent être regroupés en plusieurs classes : les promoteurs, les *enhancers/silencers*, et les insulateurs. Les CRMs des différentes classes partagent les marques d'acétylation H3Ac et H4Ac, les promoteurs actifs sont spécifiquement marqués par H3K4me3, et les enhancers et insulateurs par H3K4me1. Les enhancers sont par ailleurs souvent fixés par le co-activateur p300. Enfin, chez les mammifères les insulateurs recrutent CTCF pour bloquer l'activation par les enhancers.

Chapitre 1. Introduction générale.

les promoteurs fixant l'ARN polymérase II (la majorité des promoteurs eucaryotes), des facteurs de transcription généraux se fixent à un cœur de $\sim 100\text{bp}$ autour du TSS afin de faciliter la fixation du complexe de polymérase. Ces coeurs de promoteurs contiennent pour certains des motifs stéréotypés, comme la boîte TATA, et ont un TSS bien déterminé ; néanmoins la plupart des promoteurs des génomes mammifères sont des régions riches en GC et en dinuclétoïdes CpG (les « îlots CpG ») qui ne possèdent pas de boîte TATA et permettent l'initiation de la transcription dans un interval d'environ 100 bases (?). Au niveau épigénétique, les promoteurs actifs sont caractérisés par une région pauvre en nucléosomes en amont du TSS, flanquée de nucléosomes possédant la marque de méthylation H3K4me3.

• *Enhancers et silencers*

Les *enhancers* et *silencers* sont respectivement définis par leur effet positif ou négatif sur l'expression d'un gène cible. Cet effet peut notamment être observé par transfert d'un plasmide contenant l'élément régulateur en amont d'un gène rapporteur dans un animal transgénique ou dans des cultures cellulaires transfectées (voir 1.6.4). Leur activité ne dépend généralement pas de leur position et de leur orientation sur le plasmide. Selon l'environnement cellulaire, une région régulatrice peut être soit *enhancer* soit *silencer*, en fonction de la nature de co-activateurs ou de co-répresseurs des TFs recrutés. Il y a néanmoins relativement peu de *silencers* caractérisés et l'on utilise le terme d'*enhancers* pour désigner de manière générale ces régions régulatrices.

Les *enhancers* peuvent se situer à des distances variables du gène qu'ils régulent (?), pouvant parfois aller jusqu'à 1 Mb comme dans le cas de *Shh* chez la souris (?) (voir fig. 1.24). Les enhancers contiennent de multiples sites de fixations de TFs. Cette multiplicité est requise pour l'activité enhancer, comme cela l'a été montré pour le premier enhancer découvert : celui du virus simien 40 (SV40) (??). Un gène peut par ailleurs posséder plusieurs enhancers distincts conduisant à des expressions spécifiques dans différents tissus, comme cela l'a été montré dans le cas du gène *eve* chez la *Drosophila* (?) ou dans le cas du cluster de gènes de détermination myogénique *Myf5* et *Mrf4* chez les mammifères (?) (fig. 1.15). Ainsi, les différents enhancers d'un même gène peuvent être vus comme autant de points d'entrée d'un réseau de régulation génétique, représentant diverses fonctions logiques et intégrant différentes information spatio-temporelles pour produire en sortie une expression génétique spatio-temporelle finement contrôlée (??).

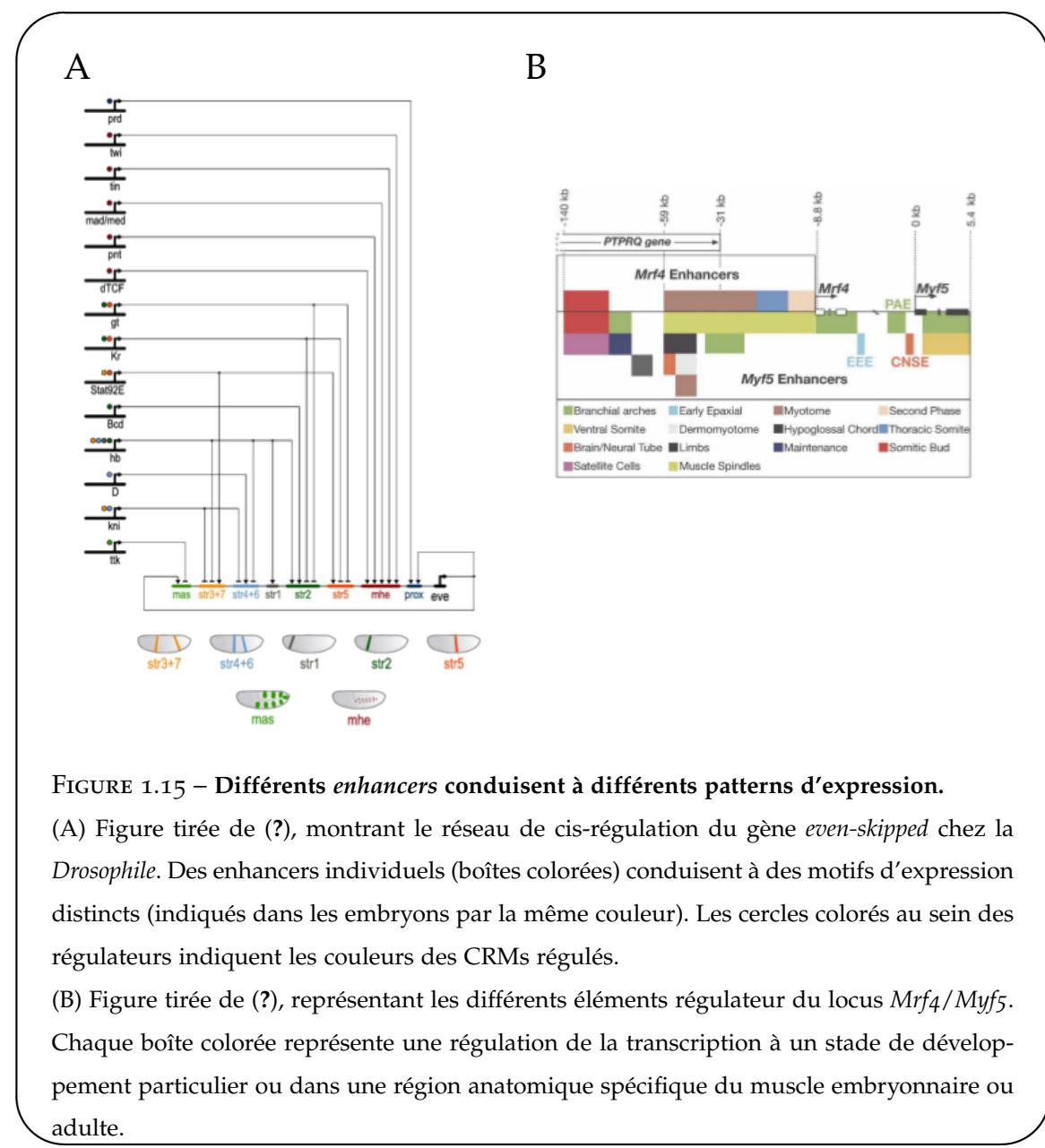


FIGURE 1.15 – Différents *enhancers* conduisent à différents patterns d’expression.

(A) Figure tirée de (?), montrant le réseau de cis-régulation du gène *even-skipped* chez la *Drosophila*. Des enhancers individuels (boîtes colorées) conduisent à des motifs d’expression distincts (indiqués dans les embryons par la même couleur). Les cercles colorés au sein des régulateurs indiquent les couleurs des CRMs régulés.

(B) Figure tirée de (?), représentant les différents éléments régulateur du locus *Mrf4/Myf5*. Chaque boîte colorée représente une régulation de la transcription à un stade de développement particulier ou dans une région anatomique spécifique du muscle embryonnaire ou adulte.

Enfin, comme décrit en fig. 1.14, les enhancers sont associés à de hauts niveaux de marque épigénétique H3K4me1 (?) et sont souvent fixés par le co-activateur p300 (??).

- **Insulateurs**

Les insulateurs sont des CRMs qui restreignent l’effet des enhancers sur leur gène cible (?). Ainsi, certains insulateurs possèdent une activité de blocage d’enhancers. Situés entre un enhancer et un promoteur cible, ces insulateurs bloquent l’activité de l’enhancer, conduisant à une réduction de l’expression du gène cible (?). Chez les mammifères, la fixation de la pro-

Chapitre 1. Introduction générale.

téine CTCF est nécessaire à cette activité de blocage de l'activité enhancer (?), alors que chez la *Drosophila* et plusieurs autres insectes il existe au moins quatre protéines additionnelles qui sont suffisantes à la réalisation de cette activité (?). Par ailleurs, les insulateurs peuvent servir de barrière de protection contre des marques d'hétérochromatine répressives. De tels insulateurs permettent notamment d'éviter les effets de positions – la modification de l'expression d'un gène selon sa position dans le chromosome – lorsqu'ils entourent un gène rapporteur intégré au hasard dans le génome (?). Cette activité passe notamment par le recrutement de USF, protéine qui recrute des enzymes de modification de la chromatine. Un insulateur peut combiner les activités de barrière de protection et de blocage d'enhancer.

De même que les enhancers, les insulateurs peuvent se situer à des distances variables des gènes qu'ils régulent. Il est à noter que la protéine CTCF possède d'autres fonctions que celle d'isolation, et tous les sites de CTCF ne correspondent pas forcément à des insulateurs (?).

1.5.2 Grammaire des enhancers : enhanceosome vs billboard

Nous l'avons vu, les CRMs contiennent en général de multiples sites de liaisons (TFBS) pour un ou plusieurs TFs. On parle de *clustering* (regroupement). Lorsque les TFBS correspondent à plusieurs TFs différents, on parle de CRM hétérotypique, et dans le cas où ils correspondent à un même TF, on parle de CRM homotypique. Cette distinction est surtout utile pour décrire les différentes méthodes de prédiction de CRM, car la plupart des CRMs identifiés chez les Métazoaires sont hétérotypiques (?). L'organisation de ces sites de liaison relève de deux types d'architecture principaux (fig. 1.16).

- **Le modèle “enhanceosome”**

Dans ce modèle, l'architecture des sites de liaison est de prime importance. Le paradigme en est l'enhancer du gène humain interferon- β , sur lequel 8 TFs se fient pour former une surface de reconnaissance continue (?). Les TFBS de cet enhancer se recouvrent les uns les autres, créant au final un complexe de TFs fixés à l'ADN agissant comme une seule unité de régulation (fig. 1.17).

- **Le modèle “billboard”**

La majorité des CRMs adhèrent à ce type d'organisation. L'architecture y est libre : les sites de liaisons n'ont pas de contrainte de nombre, d'ordre, de sens, ou d'espacement (?). De tels CRMs sont propices à une détection informatique basée sur la densité en sites de liaisons

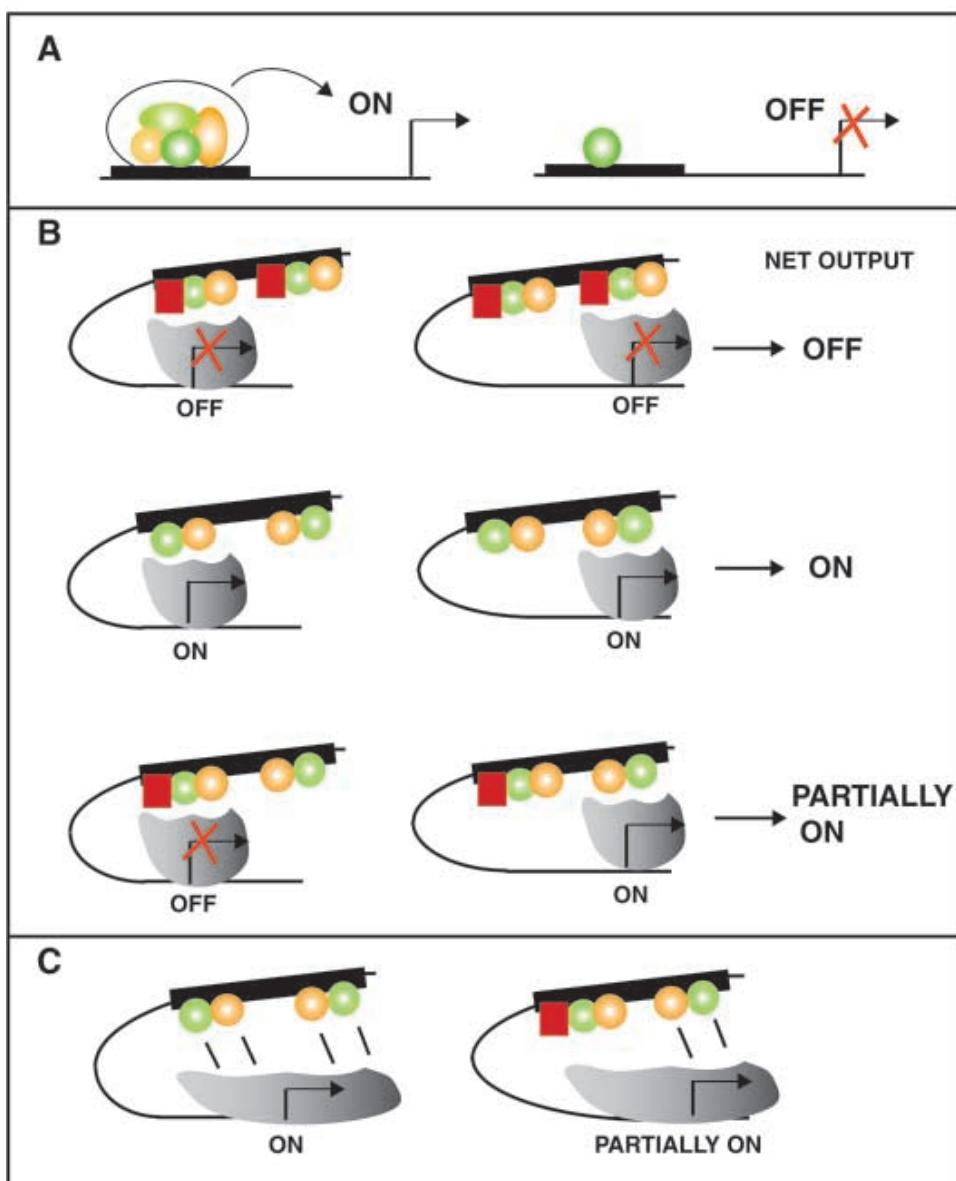
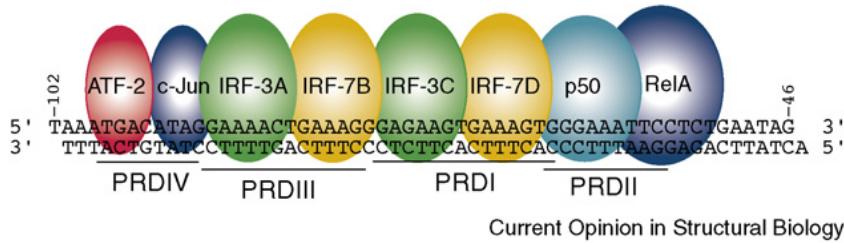


FIGURE 1.16 – Deux modèles d’enhancers : enhanceosome et billboard.

Figure tirée de (?). (A) Dans le modèle enhanceosome, l’enhancer traite l’information des multiples TFs qui le fixent. Un complexe très structuré crée une interface qui recrute la machinerie de transcription basale. L’enhancer peut être vu comme un ordinateur moléculaire qui produit à partir d’entrées multiples un seul signal vers la machinerie de transcription. Le gène cible n’est activé qu’en cas de formation du complexe entier, ce qui fournit un interrupteur binaire on/off seulement activé en cas de stimulus adéquat. La déstabilisation du complexe en changeant par exemple la concentration d’une des protéines permettrait alors d’obtenir une réponse graduelle. (B,C) Modèle d’enhancer « billboard ». Dans ce cas, l’enhancer ne consiste pas en une seule unité de régulation, mais en des sous-unités pouvant contenir différentes informations (répression ou activation par exemple) que la machinerie basale échantillonne soit itérativement (B), soit simultanément (C).



Current Opinion in Structural Biology

FIGURE 1.17 – L’enhanceosome de l’interferon- β .

Figure tirée de (?) représentant le complexe de TFs assemblés sur l’enhanceosome IFN- β , formant une surface de reconnaissance agissant comme une seule unité de régulation.

pour différents TFs.

1.5.3 Évolution des enhancers

La fonction centrale que jouent les enhancers dans la régulation de l’expression génétique laisse à penser que ceux-ci seraient sous sélection et leur séquence serait donc plus conservée que celle des régions non codantes du génome. De fait, la comparaison de séquences non-codantes entre espèces proches s’avère être un mode de détection puissant des régions de régulation (?). Ainsi, l’utilisation de la conservation entre des espèces lointaines comme l’homme et le poisson *Fugu* ou de l’extrême conservation entre des espèces proches comme l’homme, la souris et le rat, permet de détecter des régions ayant une activité enhancer *in vivo* avec un succès proche de 50% (?). À l’instar de la régulation de l’interféron- β , de telles séquences très contraintes obéissent à une logique de type « enhanceosome » où la fonction est intimement liée à la séquence.

Contrastant avec cette vision d’enhancers très contraints, plusieurs études pointent vers une plus grande flexibilité des séquences enhancers (???). Supportant l’idée que la plupart des enhancers se comportent selon le modèle « billboard », la grammaire des sites de fixation dans des séquences orthologues apparaît comme étant loin d’être figée (?). Ainsi, l’enhancer régulant le gène *short gastrulation* (*sog*), bien que présentant chez différentes espèces de *Drosophila* une architecture variable des sites de fixation le composant, conduit à un même motif d’expression (fig. 1.18). Cette idée qu’une panoplie de grammaires conduisent à une même régulation est confortée par les résultats de ? où des enhancers ayant des « entrées » différentes (i.e étant fixés par des TFs différents pendant des durées variables) produisent des « sorties »

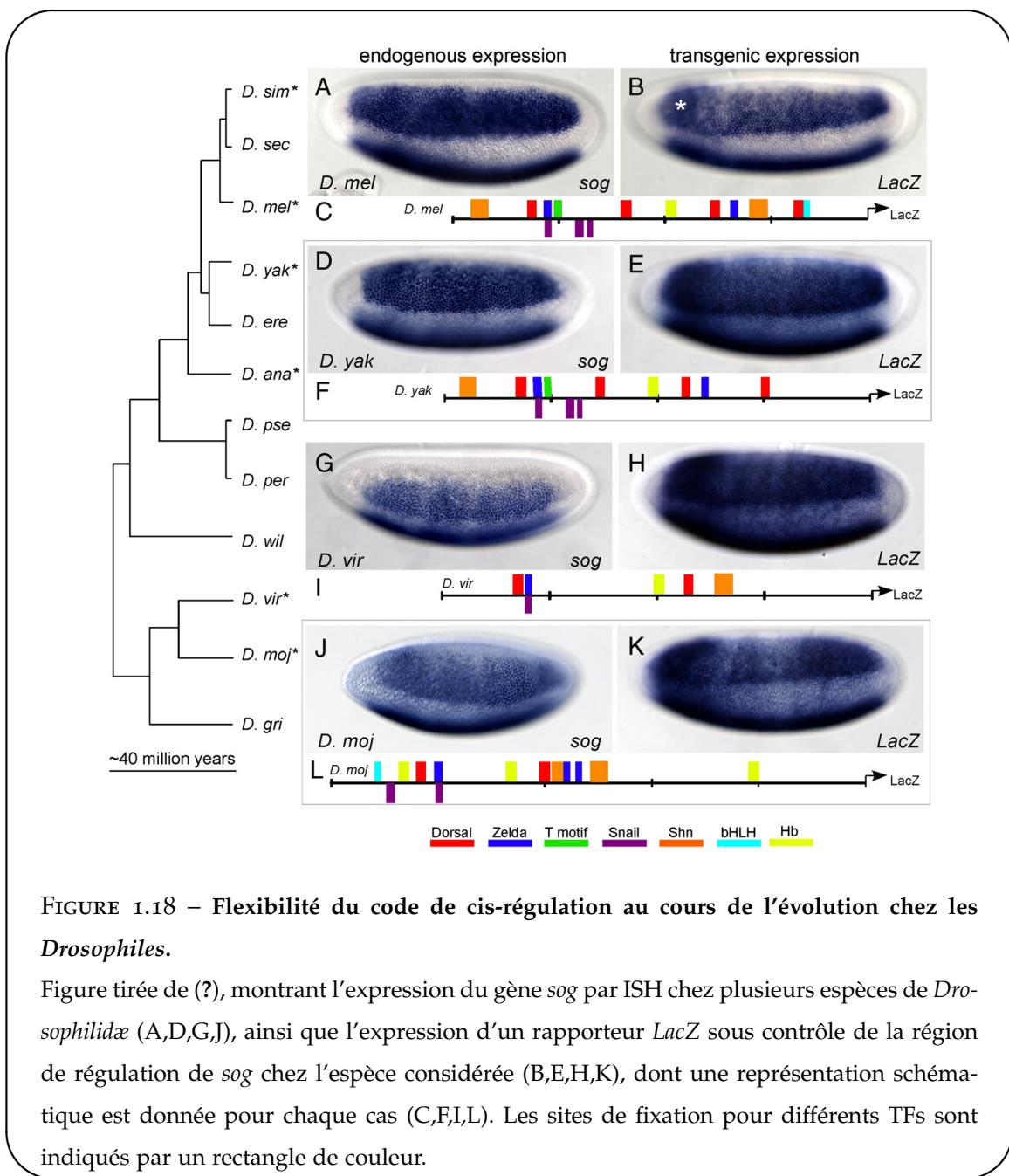


FIGURE 1.18 – Flexibilité du code de cis-régulation au cours de l'évolution chez les *Drosophiles*.

Figure tirée de (?), montrant l'expression du gène *sog* par ISH chez plusieurs espèces de *Drosophilidæ* (A,D,G,J), ainsi que l'expression d'un rapporteur *LacZ* sous contrôle de la région de régulation de *sog* chez l'espèce considérée (B,E,H,K), dont une représentation schématique est donnée pour chaque cas (C,F,I,L). Les sites de fixation pour différents TFs sont indiqués par un rectangle de couleur.

Chapitre 1. Introduction générale.

similaires, dans ce cas une expression spécifique à un tissu donné.

Supportant l'idée d'une flexibilité de la régulation, plusieurs études ont exhibé l'évolution rapide des sites de liaison de TFs dans le génome (?). Une étude de la fixation génomique des facteurs de transcription CEBP α et HNF4 α dans les cellules du foie de 5 espèces de vertébrés (l'homme, deux espèces de souris, le chien et le poulet) a notamment montré que les événements de fixation conservés chez les 5 espèces sont très rares ($\sim 0.3\%$ des pics humains) et correspondent à des régions ultraconservées proches de gènes importants dans la spécification du foie (?). Par ailleurs, lors de la perte de fixation dans l'une des espèce, un gain de fixation proche ($\pm 10\text{kb}$) est observé dans la moitié des cas. Étonnamment, ces changements rapides du câblage du réseau affectent peu l'expression génétique globale (??).

Cette évolution est en grande partie due à une évolution de séquence de fixation. Ainsi, une étude récente a utilisé une souris portant le chromosome 21 de l'homme pour comparer la fixation du facteur HNF4 α dans un contexte murin par rapport au contexte original (?). Le paysage de fixation sur le chromosome 21 exogène a très précisément récapitulé celui observé chez l'homme (fig. 1.19), montrant que le contexte cellulaire est sensiblement le même chez les deux espèces. Par ailleurs, des modifications épigénétiques ainsi que l'expression des ARNm ont pu être récapitulées.

Reste la question du mécanisme permettant cette évolution rapide. Une étude portant sur 7 facteurs de transcription chez les mammifères a montré qu'une proportion importante ($\sim 20\%$) des régions de fixation de ces TFs se situent au sein de différentes familles de transposons (?) (fig. 1.19). Ces transposons, ou éléments transposables, sont des anciens rétrovirus intégrés dans les génomes mammifères ayant la capacité de se dupliquer pour s'intégrer dans une autre région du génome et jouent un rôle fondamental dans l'évolution des génomes (?). Leur accumulation dans le génome a vraisemblablement permis d'obtenir un matériau de base permettant de produire par mutations ponctuelles des éléments de régulation *de novo* (?). Par ailleurs, les transposons peuvent permettre de diffuser par « copier-coller » des éléments de régulation existant. Ainsi, des vagues d'expansion de transposons spécifiques à différentes espèces de mammifères sont à l'origine de la variabilité des régions de fixation observée dans le cas du facteur CTCF (?).

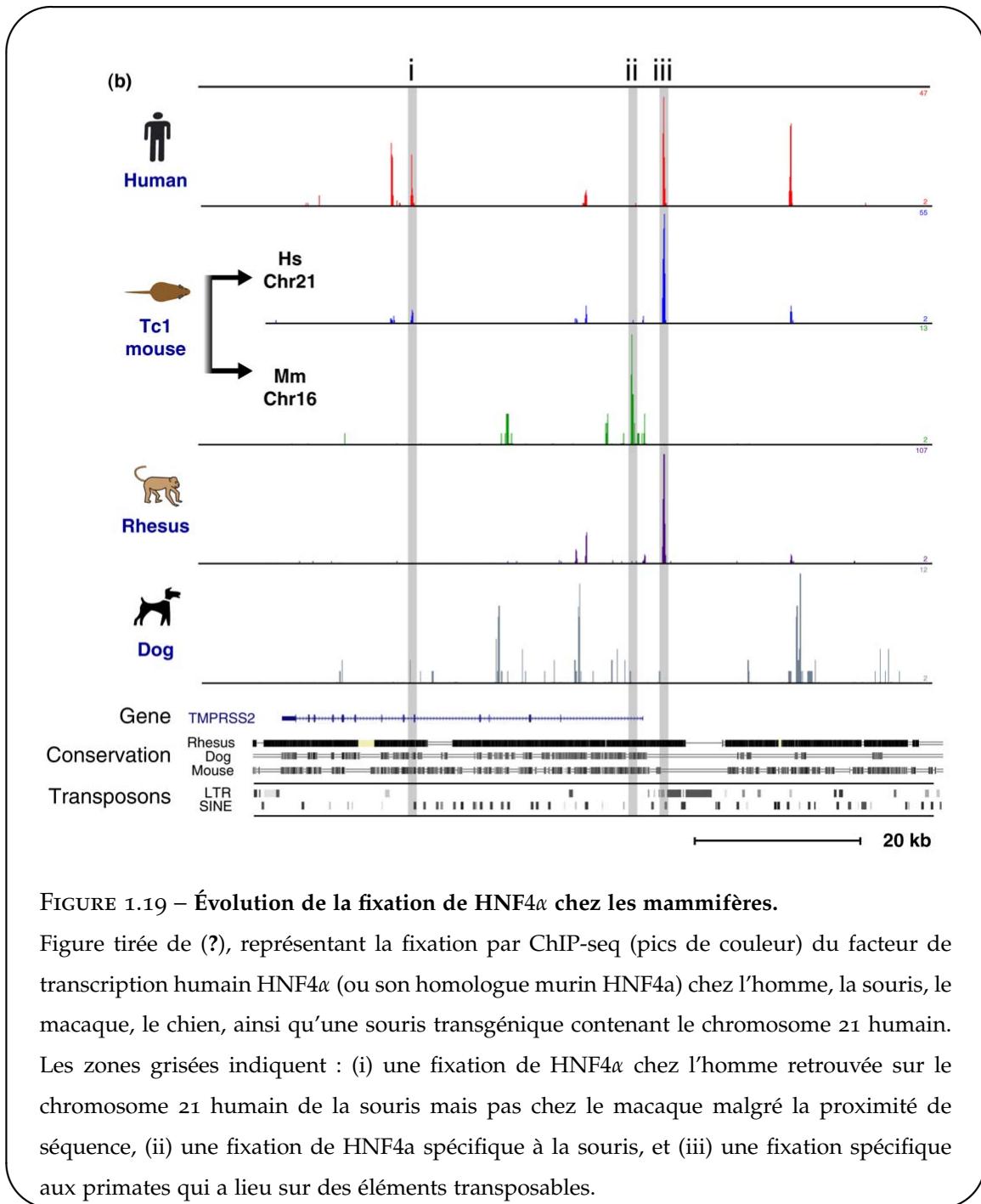


FIGURE 1.19 – Évolution de la fixation de HNF4 α chez les mammifères.

Figure tirée de (?), représentant la fixation par ChIP-seq (pics de couleur) du facteur de transcription humain HNF4 α (ou son homologue murin HNF4a) chez l'homme, la souris, le macaque, le chien, ainsi qu'une souris transgénique contenant le chromosome 21 humain. Les zones grisées indiquent : (i) une fixation de HNF4 α chez l'homme retrouvée sur le chromosome 21 humain de la souris mais pas chez le macaque malgré la proximité de séquence, (ii) une fixation de HNF4a spécifique à la souris, et (iii) une fixation spécifique aux primates qui a lieu sur des éléments transposables.

1.5.4 Les « shadow enhancers »

L'évolution des éléments de cis-régulation est un mécanisme majeur permettant la diversité animale. Néanmoins, de tels changements pourraient compromettre certaines activités génétiques essentielles. Des expériences de ChIP-on-chip ont suggéré que plusieurs gènes de développement actifs lors du développement précoce de l'embryon de Drosophile possèdent des CRMs secondaires, qui conduisent à des motifs d'expression génétique comparables à ceux produits par des CRMs « primaires » plus proximaux (?). L'expression de « shadow enhancer » a été proposée par Michael Levine en 2008 pour décrire ces CRMs redondants et souvent distaux de plusieurs dizaines de kb du gène régulé (?). Il est probable que de tels CRMs soient apparus au cours de l'évolution par duplication du CRM primaire, à l'instar du phénomène de duplication des séquences codant pour des protéines. L'avantage évident que peut conférer la redondance d'un élément de régulation est d'offrir de la robustesse face aux mutations. Par ailleurs, une telle redondance permet de faciliter la divergence et donc la spécialisation des différents CRMs. Ainsi les « shadow enhancers » semblent évoluer plus rapidement que les CRMs primaires auxquels ils sont apparentés (?) pour fournir de nouveaux sites de fixation et conduire à de nouvelles activités de régulation sans bloquer la fonction critique de certains gènes de développement.

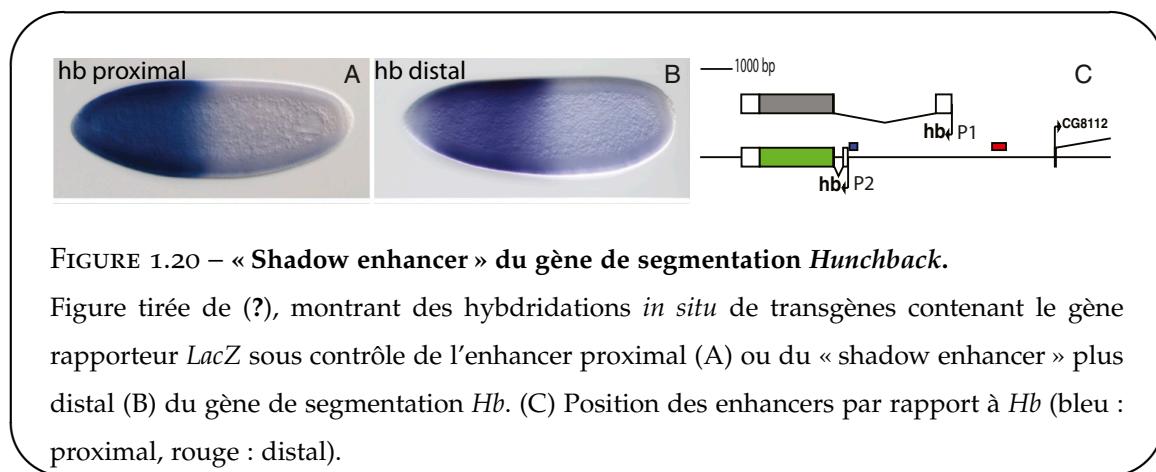


FIGURE 1.20 – « Shadow enhancer » du gène de segmentation *Hunchback*.

Figure tirée de (?), montrant des hybrides in situ de transgènes contenant le gène rapporteur LacZ sous contrôle de l'enhancer proximal (A) ou du « shadow enhancer » plus distal (B) du gène de segmentation *Hb*. (C) Position des enhancers par rapport à *Hb* (bleu : proximal, rouge : distal).

Un exemple mêlant robustesse et divergence est le cas des multiples CRMs régulant le gène *Svb* chez la Drosophile. Chaque CRM est lié à la production d'un motif distinct de trichomes (excroissances de l'épithélium comparables à des poils) sur la larve : ainsi, plusieurs mutations dans ces différents CRMs sont nécessaires pour observer un changement morphologique conséquent (?). Dans ce même système, il a été montré que deux CRMs supplémentaires, des

« shadow enhancers », sont dispensables dans des conditions de température usuelles, mais requis lorsque les embryons se développent dans des conditions de température extrêmes (?).

Par ailleurs, il a été montré que les gènes de segmentation (ou gènes *gap*) de la Drosophile possèdent tous des « shadow enhancers » (fig. 1.20). Leur rôle semble être d'assurer une plus grande précision spatiale du motif d'expression du gène régulé : la perte de l'un des CRMs, proximal aussi bien que « shadow », conduisant à une expression trop restreinte ou trop répandue spatialement selon le cas (?).

1.5.5 Par delà les enhancers : les « super-enhancers »

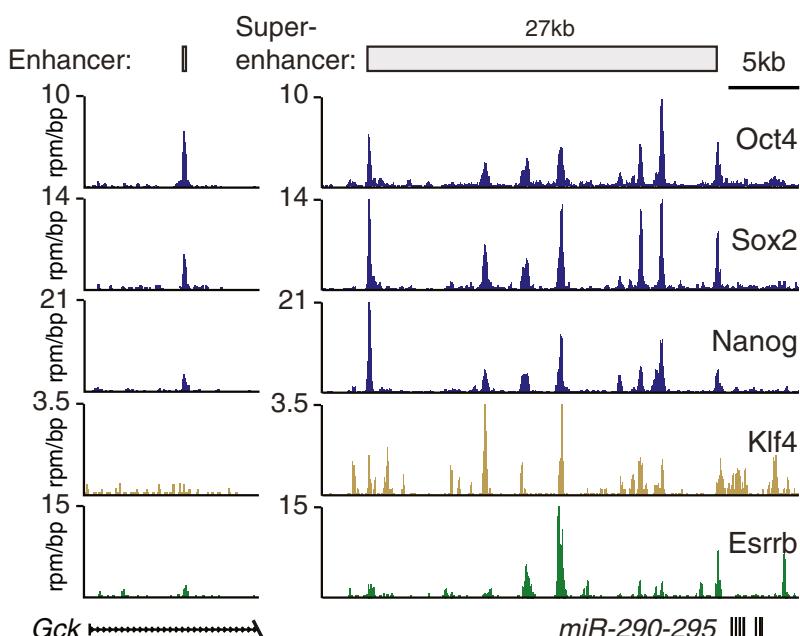


FIGURE 1.21 – De l'enhancer au super-enhancer.

Figure tirée de (?), montrant les profils de ChIP-seq des TFs maîtres Oct4, Sox2, Nanog, Klf4 et Esrrb aux loci de *Gck* et *miR-290-295* dans les cellules souches embryonnaires. Le super-enhancer se distingue du simple enhancer par sa taille (27kb), sa grande concentration en TFs maîtres, notamment Klf4 et Esrrb, et la fixation de la protéine Med1 du complexe Mediator.

Récemment, il a été montré que certains groupements d'enhancers peuvent agir comme une même unité de régulation : on parle de *super-enhancers* (?). Ces régions de taille typique $\sim 10\text{kb}$ (fig. 1.21), sont fixées par des TFs maîtres et sont associées à des gènes encodant

Chapitre 1. Introduction générale.

des régulateurs clés de l'identité cellulaire. Identifiés dans les cellules souches embryonnaires (ESCs), ces ensembles d'enhancers sont fixés par le complexe co-activateur Mediator, qui interagit avec la cohésine pour former un anneau permettant de connecter la région de régulation au promoteur (?). Par ailleurs, les gènes associés aux super-enhancers possèdent un niveau particulièrement élevé d'expression et leur knock-down est associée à une perte de l'état souche des cellules.

Ainsi, ce second niveau d'organisation de la régulation pourrait simplifier la modélisation de la régulation du type cellulaire, en passant de millier de traces de fixation pour différents TFs à quelques centaines de super-enhancers contrôlant les gènes clés de l'identité cellulaire.

1.6 Prédiction et validation des CRMs

1.6.1 Méthodes utilisant la concentration en sites de fixation

Nous l'avons vu, une propriété des CRMs est leur grande concentration en TFBS. Ceci a motivé des approches de prédiction de promoteurs et d'enhancers basées sur leur contenu ou *clustering* en motif (fig. 1.22a). L'avantage de telles approches est qu'elles peuvent être réalisées avec seulement la séquence d'ADN génomique et des modèles de TFs ou motifs (par exemple des PWMs, voir fig. 1.10) représentant les facteurs de transcription impliqués dans le processus étudié. Cependant, les clusters de motifs sont très répandus dans les grands génomes, et sans l'ajout d'informations supplémentaires comme les marques épigénétiques ou l'expression des gènes voisins, ces approches produisent un grand nombre de faux positifs (éléments prédits comme positifs mais étant en réalité négatifs). Par ailleurs, les TFs impliqués ne sont pas toujours connus, et il faut alors apprendre des motifs putatifs à partir de séquences fonctionnelles.

- **Approches utilisant des motifs connus**

L'une des premières investigations basée sur le regroupement de TFBS utilisait 5 motifs connus de la détermination musculaire pour prédire par régression linéaire les CRMs actifs dans le muscle (?). Le taux de validation était relativement bas, autour de 20%. De même, chez *Drosophila melanogaster*, plusieurs études ont utilisé le clustering de motifs pour prédire des CRMs de différents processus développementaux (par ex ?). Ces études ont trouvé de nouveaux enhancers validés expérimentalement (bonne sensibilité) mais avaient des taux de prédiction relativement bas, entre 15 et 30%. L'algorithme *Ahab* (?), utilisant un modèle thermodynamique de fixation des TFs sur les CRMs, a quant à lui réussi à prédire un nombre bien plus important de régions fonctionnelles : $\sim 80\%$ des modules prédits à proximité de 29 gènes de segmentation chez la drosophile ont effectivement récapitulé le motif d'expression du gène associé (?). Ce succès semble notamment être dû au fait que ce modèle thermodynamique, basé sur une prise en compte exhaustive de toutes les segmentations possibles des CRMs en motifs et en ADN « background », permet de donner plus de poids au cas où plusieurs sites de faibles affinité pour un TF se trouvent au sein d'un même module, alors que les autres méthodes utilisent généralement un seuil de probabilité relativement élevé (afin d'éviter les faux positifs) à partir duquel une séquence est considérée comme fixée par un TF. Par ailleurs,

Chapitre 1. Introduction générale.

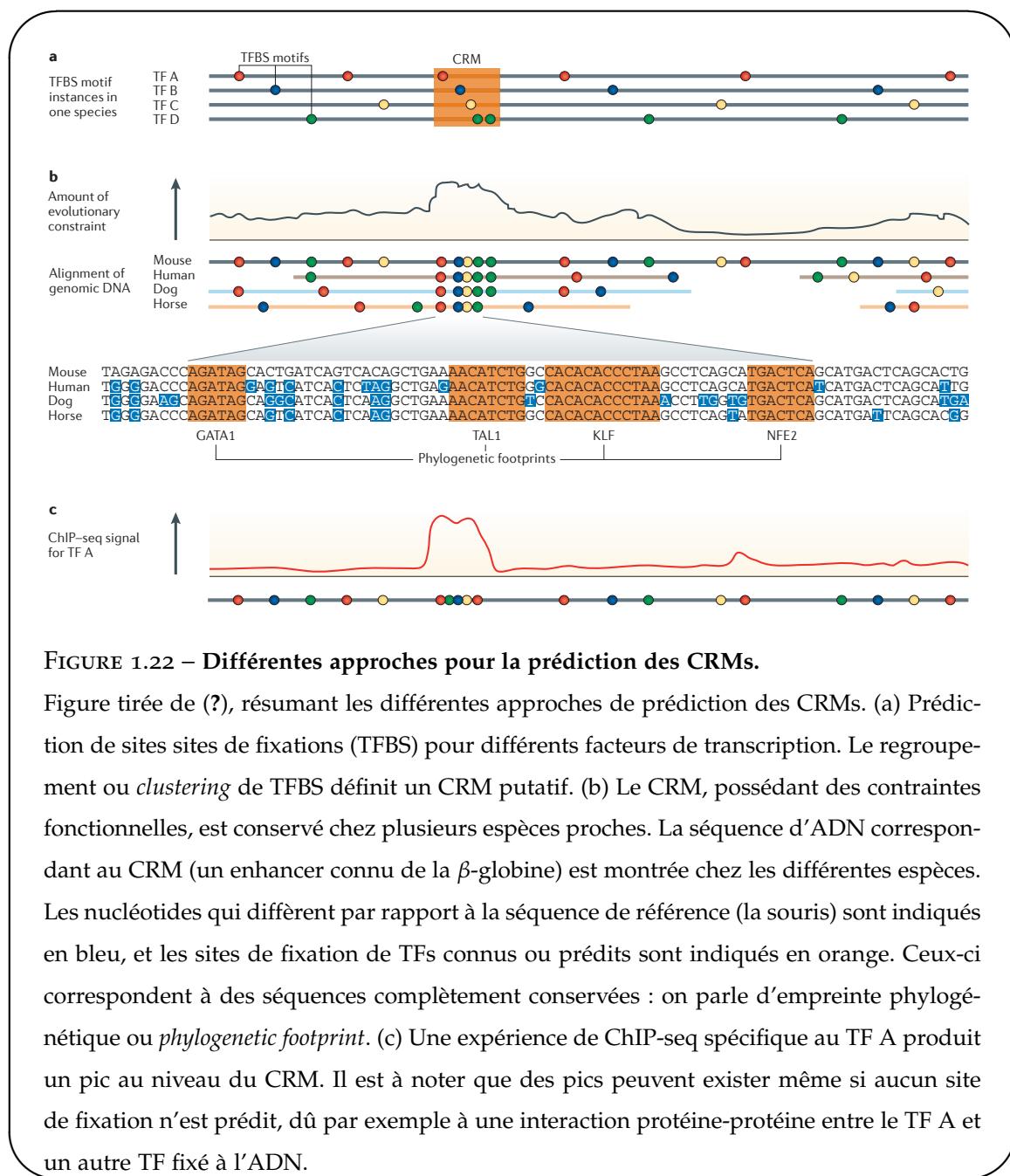


FIGURE 1.22 – Différentes approches pour la prédition des CRMs.

Figure tirée de (?), résumant les différentes approches de prédition des CRMs. (a) Prédiction de sites sites de fixations (TFBS) pour différents facteurs de transcription. Le regroupement ou *clustering* de TFBS définit un CRM putatif. (b) Le CRM, possédant des contraintes fonctionnelles, est conservé chez plusieurs espèces proches. La séquence d'ADN correspondant au CRM (un enhancer connu de la β -globine) est montrée chez les différentes espèces. Les nucléotides qui diffèrent par rapport à la séquence de référence (la souris) sont indiqués en bleu, et les sites de fixation de TFs connus ou prédits sont indiqués en orange. Ceux-ci correspondent à des séquences complètement conservées : on parle d'empreinte phylogénétique ou *phylogenetic footprint*. (c) Une expérience de ChIP-seq spécifique au TF A produit un pic au niveau du CRM. Il est à noter que des pics peuvent exister même si aucun site de fixation n'est prédict, dû par exemple à une interaction protéine-protéine entre le TF A et un autre TF fixé à l'ADN.

cette étude s'est restreinte à un ensemble de gènes connus pour lesquels les régions à proximité riches en TFBS ont *a priori* plus de chances d'être fonctionnelles. De manière générale, plus le domaine de recherche est étendu (par exemple, le génome entier), plus le nombre de faux positifs augmente.

- **Approches *de novo* où les motifs ne sont pas connus**

Lorsque les motifs (PWMs) ne sont pas connus à l'avance, il faut les générer *de novo* à partir de leur surreprésentation dans des CRMs connus. Par exemple, l'algorithme CisModule permet de générer des motifs et des modules simultanément en utilisant un modèle de mélange hiérarchique (?). Lorsqu'il est appliqué aux CRMs musculaires introduits précédemment, il permet de retrouver certains motifs connus et permet de retrouver $\sim 70 - 80\%$ des séquences connues lorsqu'elles sont mélangées avec un nombre similaire de séquences aléatoires. Par ailleurs, l'apprentissage de modèles permettant de discriminer différentes classes de CRMs entre elles plutôt qu'une classe de CRMs par rapport à des séquences aléatoires ou intergéniques peut s'avérer plus fructueux. Ainsi, (?) ont utilisé des motifs connus ainsi que des motifs appris *de novo* avec le programme DME (?) pour leur capacité à discriminer des séquences appartenant à différents jeux de données de régions promotrices pour bâtir un modèle de régression logistique permettant de prédire l'activité tissu-spécifique dans 45 des 56 tissus humains et murins considérés. Il existe aussi plusieurs méthodes qui n'utilisent pas de motifs du type PWM, mais de purs modèles probabilistes tels que des chaînes de Markov d'ordre 5 ou des regroupements de « mots » de k nucléotides ou k -mers selon des critères de distance de Hamming et surreprésentés dans les séquences d'intérêt, par exemple (?). Ces méthodes sont passées en revue dans (?), et elles peuvent atteindre des sensibilités de $\sim 60\%$ pour la prédiction de CRMs mammifères. L'intérêt est que ces études ne présument pas d'un modèle de fixation des TFs à l'ADN. C'est aussi un désavantage, puisqu'elles sont moins informatives quant au réseau génétique sous-jacent et aux mécanismes de régulation impliqués.

1.6.2 Méthodes utilisant la phylogénie

Les approches utilisant la comparaison des génomes de différentes espèces pour prédire des CRMs sont basées sur l'idée que les séquences de régulation sont plus fortement conservées que l'ADN non fonctionnel les entourant. Nous l'avons vu en 1.5.3, une proportion importante de CRMs ne satisfont pas à cette règle. Cette approche ne permet donc d'étudier que

Chapitre 1. Introduction générale.

le sous-ensemble de CRMs qui a subi une forte pression de sélection depuis le dernier ancêtre commun aux espèces considérées et ne donne pas accès aux CRMs apparus récemment au sein d'une espèce.

• Prédiction à partir de la contrainte évolutive seule

L'alignement de séquences non-codantes orthologues fait apparaître des parties très conservées, avec peu de variations dans les séquences sous-jacentes, entourées de séquences accumulant les variations (fig. 1.22b). De telles séquences conservées sont alors interprétées comme ayant été sous sélection, les substitutions délétères ayant été rejetées au cours de l'évolution (?). Par analogie avec les empreintes à la DNAse I, on parle d'empreinte phylogénétique pour caractériser ces courtes séquences très conservées ($\sim 10\text{bp}$), traces de la fixation putative d'un facteur de transcription. Ces empreintes s'avèrent être un indicateur fiable de fonctionnalité (?) et, parce qu'elles ne reposent pas sur des modèles *a priori* de fixation, elles permettent de plus de trouver des motifs de régulation non connus (?). Au niveau de séquences plus longues ($\sim 100\text{bp}$), la contrainte évolutive permet de détecter des CRMs entiers. Ainsi, comme nous l'avons vu en 1.5.3, l'utilisation de la conservation extrême permet d'atteindre 50% de taux validation (?). Néanmoins, lorsque ces contraintes de conservation extrême (par exemple homme-Fugu) sont relâchées, le taux de validation tombe drastiquement, atteignant $\sim 5\%$ (?), montrant la nécessité d'allier le critère de conservation à d'autres données (expression, ChIP...) pour améliorer la prédiction des CRMs.

• Prédiction utilisant la phylogénie et des motifs connus

Une approche pour améliorer les prédictions est de combiner les approches précédentes en utilisant à la fois le *clustering* en TFBS et la contrainte évolutive. À l'échelle du génome entier, cette approche permet de filtrer les résultats pour améliorer le signal de détection chez la Drosophile (?). Du côté des mammifères, en utilisant les motifs de la base de données TRANSFAC et la conservation entre l'homme et la souris, ? ont créé une base de données de modules, PReMods, qui retrouve $\sim 17\%$ de CRMs connus et recoupe 40% des fragments occupés par le co-activateur et marqueur de l'activité enhancer p300. D'autres méthodes se sont concentrées sur des types cellulaires bien définis. Par exemple, la recherche de sites conservés pour des motifs de TF des cellules sanguines connus (?) a permis de définir des CRMs dont 2 ont été testés et validés.

Certains efforts ont par ailleurs été menés pour sortir du cadre d'une conservation de

séquence stricte en modélisant l'évolution d'un CRM fixé par un certain nombre de motifs connus. Par exemple, le modèle MorphMS (?) cherche au sein d'un alignement de deux séquences orthologues des régions prédictes par un modèle d'évolution dérivé d'un ensemble de motifs choisis par l'utilisateur. Une extension de cette approche incorpore le gain et la perte de sites de fixation, mais n'a cependant pas encore été appliquée à l'échelle du génome (?).

- **Approches utilisant la phylogénie pour générer des motifs *de novo***

De même que précédemment, tous les motifs ne sont pas connus et il peut être utile d'avoir recours à de l'apprentissage direct à partir de séquences fonctionnelles connues pour aider à la prédiction. Par exemple, l'algorithme ESPER cherche des patterns (TFBS, %GC, etc) surreprésentés dans des alignements multi-espèces de CRMs connus par rapport à des alignements d'ADN *a priori* non fonctionnel (?). Cette méthode n'est pas restreinte à l'analyse de séquences conservées puisqu'elle peut potentiellement capturer des signatures de changements systématiques. La prédiction de régions de haut potentiel de régulation recouvre presque entièrement les prédictions de PReMods, et le test par transfection de ces régions à proximité de gènes exprimés dans les cellules érythroïdes et possédant un site pour un TF spécifique de l'érythroïde mène à un taux de validation de 50%. Une autre méthode consiste à chercher des mots surreprésentés dans un ensemble d'apprentissage de CRMs connus puis à restreindre les prédictions aux régions conservées (?). Les prédictions réalisées ont toutes été validées chez la Drosophile comme chez la souris.

1.6.3 Méthodes utilisant les marques épigénétiques et de ChIP-seq pour des TFs

- **Prédiction des promoteurs**

La méthode la plus fiable de prédiction d'un promoteur utilise le fait qu'il est toujours localisé au niveau d'un TSS, dont la position peut facilement être obtenue en alignant les séquences de l'ARN du gène correspondant sur le génome (?). Le taux de validation avec cette seule contrainte est très élevé : 91% ont une activité dans au moins un type cellulaire. Par ailleurs, la marque épigénétique H3K4me3 est aussi un indicateur des promoteurs actifs dans le type cellulaire étudié (?) (fig. 1.14).

- **Prédiction des enhancers**

La prédiction des enhancers à partir des marques épigénétiques, comme l'acétylation des histones (?), la méthylation H3K4me1 (?), ou encore la présence du co-activateur p300 (?),

Chapitre 1. Introduction générale.

est très efficace, avec une expression tissu-spécifique dans $\sim 80\%$ des cas (?). Par exemple, ces différentes marques, présentes dans différents tissus, peuvent être utilisées comme autant d'entrées d'un modèle de Markov caché pour produire des prédictions fiables de CRMs tissu-spécifiques chez l'homme (?).

En fait, les prédictions d'activité enhancer à partir de ces marques épigénétiques est plus fiable qu'en utilisant la fixation de facteurs de transcription tissu-spécifiques. Par exemple, sur 63 séquences ADN fixées *in vivo* par le facteur spécifique des cellules sanguines GATA1 chez la souris, seulement la moitié conduisent à une activité après transfection dans des cultures cellulaires (?). Ces enhancers fonctionnels sont par ailleurs plus particulièrement associés à un site de fixation conservé pour GATA1, montrant à nouveau la nécessité de combiner les approches pour améliorer la détection. Un taux de validation similaire a été observé pour le facteur de différenciation myogénique MyoD, avec 40% de régions fixées ayant une activité après transfection en cellules.

L'utilisation de données de fixation pour plusieurs TFs à la fois semble cependant améliorer le pouvoir de prédiction. Ainsi, ? ont étudié la co-fixation de GATA1 avec 4 autres TFs hématopoïétiques dans des mégacaryocytes. En s'intéressant aux gènes à proximité de ces régions, ils en ont découvert plusieurs qui n'étaient pas précédemment connus comme étant important dans l'hématopoïèse. Leur fonction a été testée par knock-down, avec dans 8 cas sur les 9 testés une réduction de la production de globules rouges.

1.6.4 Validation expérimentale

Il existe plusieurs méthodes pour s'assurer de la fonctionnalité d'un CRM prédit.

Tout d'abord, une méthode indirecte donnant du crédit à la prédiction d'un CRM est d'examiner le motif d'expression du gène dont le TSS est le plus proche. Si cette expression reproduit les caractéristiques utilisées pour prédire le CRM (par exemple, s'exprimer dans le muscle pour une prédiction de CRMs utilisant l'abondance de sites de liaison de TFs musculaires), alors cela soutient l'idée (mais ne la démontre pas) que la présence du CRM en est la cause.

Une méthode plus directe permettant de démontrer qu'un fragment d'ADN régule l'expression génétique consiste en une expérience de gain de fonction dans laquelle un plasmide contenant le CRM prédit à proximité d'un gène rapporteur est introduit par transfection *in vitro* en cellule, permettant un suivi quantitatif de l'activité, ou par transgenèse *in vivo* dans

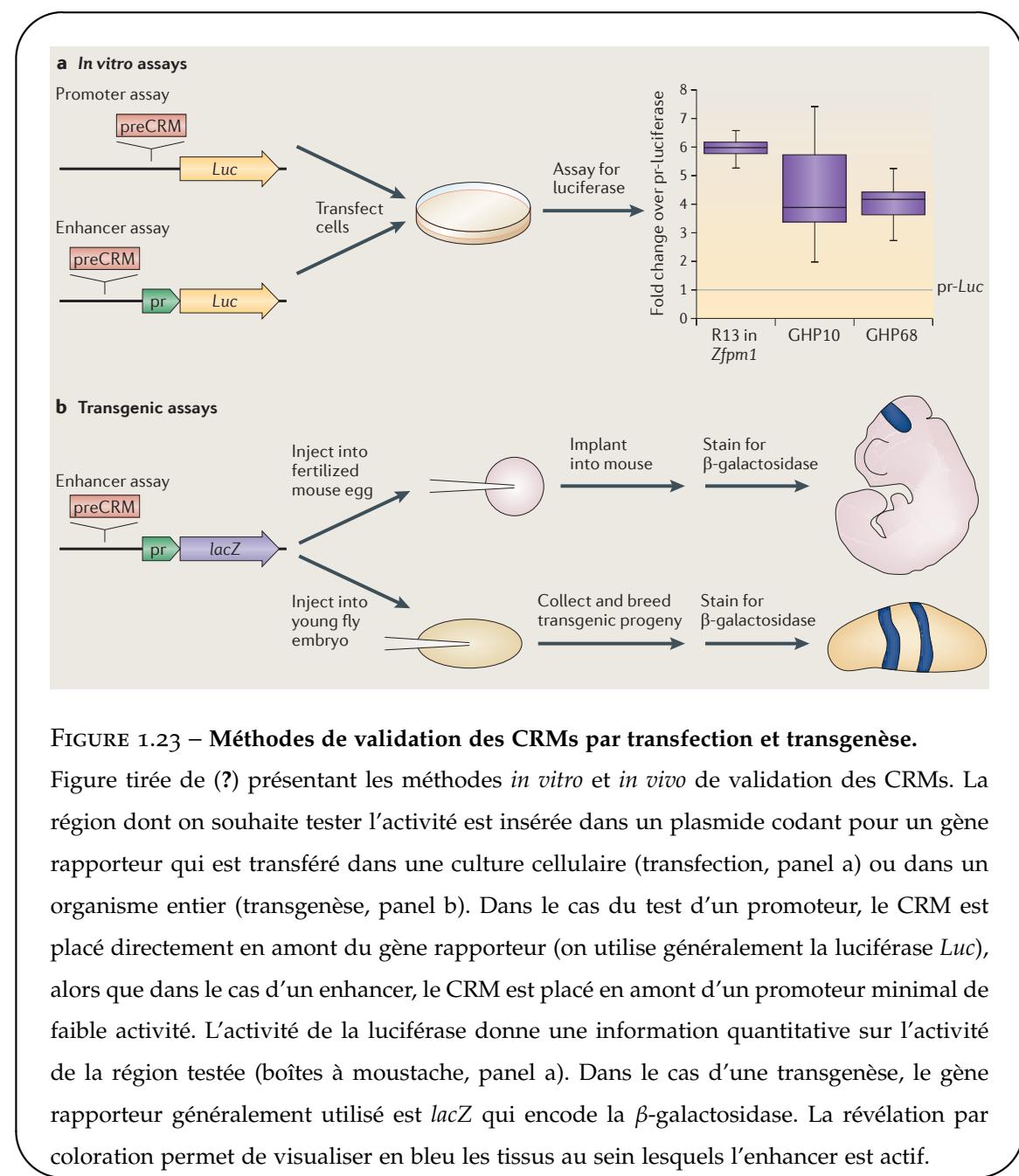


FIGURE 1.23 – Méthodes de validation des CRMs par transfection et transgenèse.

Figure tirée de (?) présentant les méthodes *in vitro* et *in vivo* de validation des CRMs. La région dont on souhaite tester l'activité est insérée dans un plasmide codant pour un gène rapporteur qui est transféré dans une culture cellulaire (transfection, panel a) ou dans un organisme entier (transgenèse, panel b). Dans le cas du test d'un promoteur, le CRM est placé directement en amont du gène rapporteur (on utilise généralement la luciférase *Luc*), alors que dans le cas d'un enhancer, le CRM est placé en amont d'un promoteur minimal de faible activité. L'activité de la luciférase donne une information quantitative sur l'activité de la région testée (boîtes à moustache, panel a). Dans le cas d'une transgenèse, le gène rapporteur généralement utilisé est *lacZ* qui encode la β -galactosidase. La révélation par coloration permet de visualiser en bleu les tissus au sein lesquels l'enhancer est actif.

un organisme, auquel cas le suivi est plus qualitatif mais permet d'établir la spécificité spatio-temporelle (tissu et stade de développement) de l'élément de régulation (fig. 1.23). Ce type d'expérience montre que le CRM prédict est *suffisant* pour reproduire le motif génétique observé. De manière optimale, il faudrait aussi montrer par délétion ciblée de l'élément de régulation au sein du génome que ce dernier est *nécessaire* à l'expression du gène endogène.

1.6.5 Implication des CRMs dans les maladies humaines

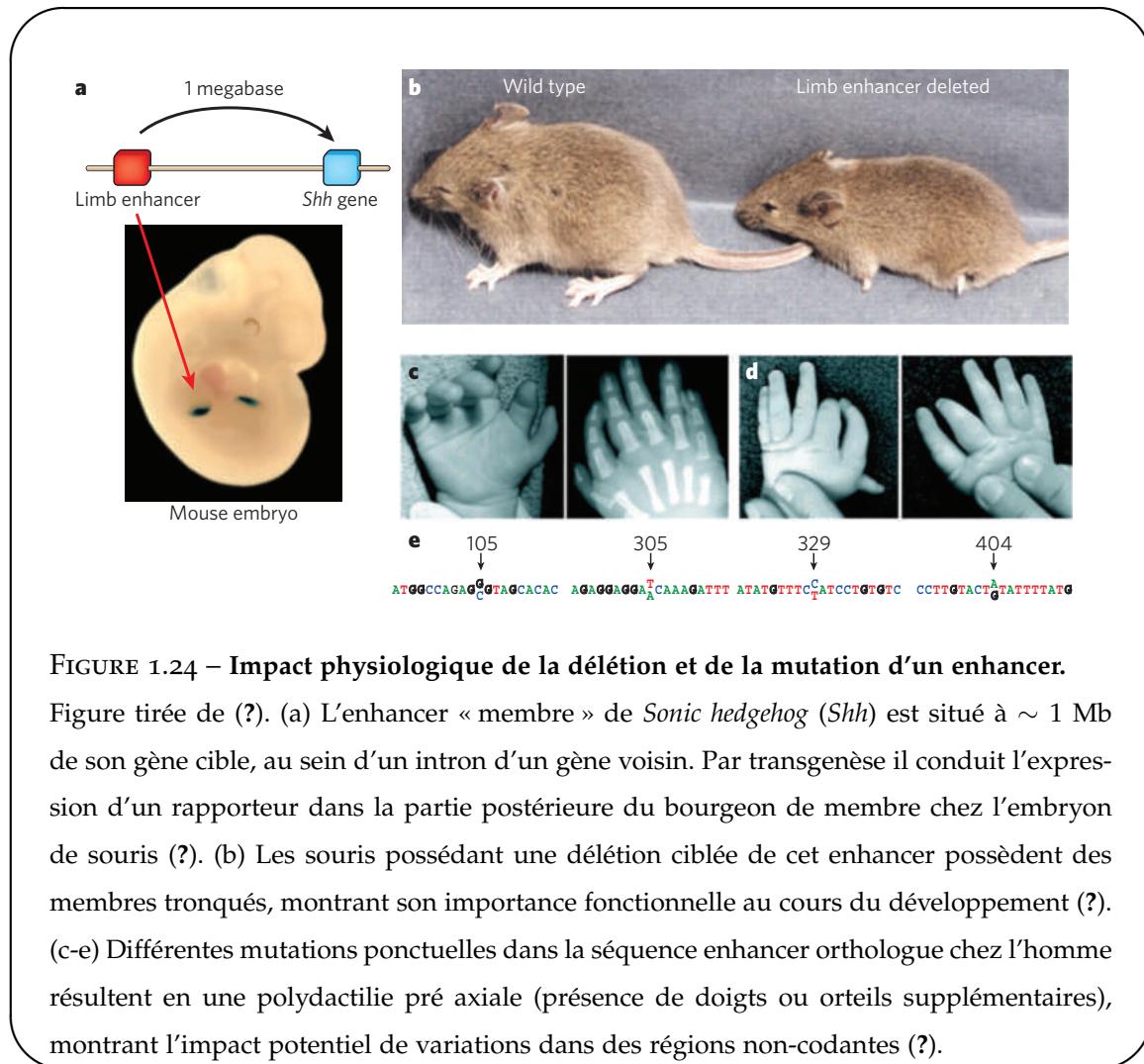


FIGURE 1.24 – Impact physiologique de la délétion et de la mutation d'un enhancer.

Figure tirée de (?). (a) L'enhancer « membre » de *Sonic hedgehog* (*Shh*) est situé à ~ 1 Mb de son gène cible, au sein d'un intron d'un gène voisin. Par transgenèse il conduit l'expression d'un rapporteur dans la partie postérieure du bourgeon de membre chez l'embryon de souris (?). (b) Les souris possédant une délétion ciblée de cet enhancer possèdent des membres tronqués, montrant son importance fonctionnelle au cours du développement (?). (c-e) Différentes mutations ponctuelles dans la séquence enhancer orthologue chez l'homme résultent en une polydactylie pré axiale (présence de doigts ou orteils supplémentaires), montrant l'impact potentiel de variations dans des régions non-codantes (?).

Au cours des dernières décennies, de nombreuses mutations dans les régions codantes des gènes, impliquant des défauts structurels des protéines associées, ont pu être associées à des maladies génétiques. À l'inverse, le rôle des mutations affectant des régions non codantes n'a été que peu exploré, essentiellement du fait de la difficulté d'annoter ces régions correctement afin de définir celles qui pourraient avoir une fonction d'intérêt. Plusieurs études ont cependant pu montrer que des variations affectant des enhancers distaux pouvaient conduire à des pathologies (?).

L'une de ces études concerne l'enhancer spécifique du membre de *Shh* (fig. 1.24). Cet enhancer, initialement décrit chez la souris, se situe à environ 1 Mb de distance de *Shh*, au sein de l'intron d'un gène voisin. Le séquençage de cet enhancer chez plusieurs individus humains a

permis d'associer une douzaine de variations mono-nucléotidiques à la polydactilie pré axiale, c'est-à-dire la présence de doigts ou d'orteils supplémentaires (?). Des études supplémentaires chez la souris ont montré que les variations de séquences observées dans cet enhancer conduisent à une expression ectopique dans la partie antérieure du membre au cours du développement, ce qui est consistant avec la présence de doigts supplémentaires (?). Par ailleurs, la délétion de l'enhancer orthologue de la souris entraîne la troncation des membres (?).

Ainsi, ces résultats montrent l'importance de l'identification des enhancers pour permettre à des études de génétique humaine d'explorer le rôle potentiellement pathologique de mutations dans des régions non codantes fonctionnelles.

1.7 Bases de données

La biologie moderne est caractérisée par l'accumulation de données biologiques qu'il s'agit d'intégrer puis d'interpréter : on parle de biologie intégrative. En particulier, depuis le séquençage du génome humain il y a maintenant plus de dix ans (?), le nombre de génome séquencés n'a cessé d'augmenter, tandis que dans le même temps le prix du séquençage diminuait drastiquement (fig. 1.25). Afin de permettre la gestion et l'utilisation de ces données, de nombreux outils et bases de données ont été mis à disposition (?). Nous évoquons ici ceux qui nous paraissent essentiels du point de vue de la régulation en *cis*.

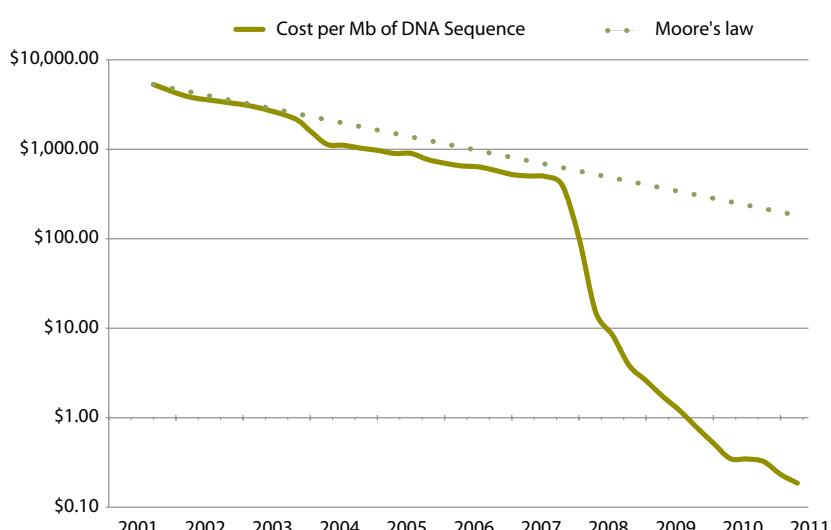


FIGURE 1.25 – Évolution du coût de séquençage.

Figure adaptée de (?), montrant l'évolution du coût du séquençage d'1 Mb d'ADN au cours de la dernière décennie, comparé à une évolution de type « Loi de Moore » où le prix serait diminué de moitié tous les 18 mois.

1.7.1 Obtention de données génomiques

Tout d'abord, les différents génomes séquencés sont à disposition sur des bases de données publiques d'où ils peuvent être téléchargés puis analysés en aval. Parmi les plus généralistes se trouvent la base de donnée de UCSC (UCSC Genome Browser, <http://genome.ucsc.edu>)

et celle de l'EMBL (Ensembl, <http://www.ensembl.org>)³.

Sont à disposition les génomes des différentes espèces séquencées pour les différents assemblages réalisés, des alignements des génomes de différentes espèces deux par deux (*pair-wise alignments*) ou par groupes d'espèces (*multiple alignments*), ainsi qu'un certain nombre d'annotations essentielles à l'analyse de ces génomes : coordonnées des gènes (TSSs, exons, introns avec potentiellement différents transcrits alternatifs), miRNA ou lincRNA, ontologies associées, coordonnées des séquences répétitives (les *repeats*, en partie liés aux éléments transposables abordés en 1.5.3, et qui sont abondants dans les génomes vertébrés), différentes données ChIP-seq, indices de conservation⁴...

Au final, ces différentes données constituent une base de travail fiable et régulièrement mise à jour. Afin de faciliter leur obtention, il est possible d'utiliser le navigateur de tables de UCSC⁵ ou la section BioMart d'Ensembl⁶.

Situé plus en amont, le projet Galaxy (<http://galaxyproject.org>) permet à l'utilisateur de récupérer des données depuis les différentes banques existantes, puis de leur faire subir divers traitements et analyses par divers outils de bioinformatique. Cet outil, qui peut être utilisé sur internet ou bien localement, a l'avantage de permettre la sauvegarde de plans de travail ou *workflows*, successions de commandes utilisées pour traiter une entrée donnée par différents outils stéréotypés et obtenir directement le résultat final, favorisant une approche conviviale orientée utilisateur.

En guise d'exemple, nous montrons en figure 1.26 des statistiques obtenues aisément à partir d'annotations génétiques présentes sur UCSC et traitées avec Galaxy. Ces statistiques sont les distribution de tailles des régions intergéniques et introniques chez plusieurs espèces : la bactérie *Escherichia coli*, la levure *Saccharomyces cerevisiae*, le ver *Caenorhabditis elegans*, la mouche *Drosophila melanogaster*, la souris, le poulet et l'homme.

3. Les données sont accessibles sur les pages de téléchargement, respectivement <http://hgdownload.cse.ucsc.edu/downloads.html> pour UCSC et <http://www.ensembl.org/info/data/ftp/index.html> pour Ensembl

4. Pour le cas de l'assemblage mm9 de la souris, ces annotations sont accessibles à l'adresse suivante : <http://hgdownload.cse.ucsc.edu/goldenPath/mm9/database/>

5. <http://genome.ucsc.edu/cgi-bin/hgTables>

6. <http://www.ensembl.org/biomart/martview>

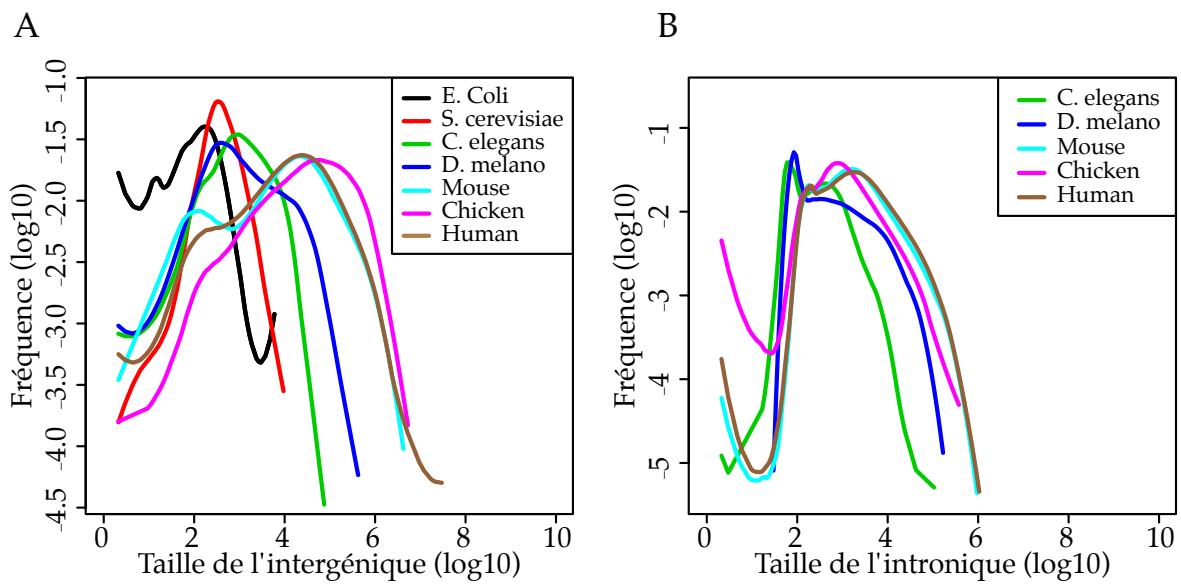


FIGURE 1.26 – Distribution des tailles intergéniques et introniques chez différentes espèces.

Distributions log-log de la taille des régions intergéniques (A) et introniques (B) chez différentes espèces. Les histogrammes sont réalisés avec un intervalle de 0.05 puis lissés avec l'estimateur local LOESS de paramètre $span = 0.3$ (logiciel R). (A) Les régions intergéniques sont définies comme les régions complémentaires aux régions transcrtes (données UCSC), celles-ci étant préalablement fusionnées pour éviter les redondances liées aux multiples transcrits d'un même gène. De la bactérie à l'homme, on observe une inflation de la quantité de génome non codant. (B) Les régions introniques sont définies par le fait qu'elles sont entourées par deux exons d'un même gène. Pour pouvoir être épissés lors de la maturation des preARNm, les introns doivent posséder des sites d'épissage, imposant une borne inférieure à leur taille pour que l'ARNm final soit fonctionnel.

1.7.2 Obtention de données sur les TFs

Nous l'avons vu, les données de fixation des TFs (ChIP-seq, ChIP-on-chip) peuvent être obtenues à partir du site UCSC Genome Browser. Ces données sont aussi généralement accessibles sur le site du NCBI (Gene Expression Omnibus, <http://www.ncbi.nlm.nih.gov/geo/>) via un numéro d'accession donné lors de la publication des données.

De nombreux modèles de TFs ont déjà été bâtis préalablement à l'avènement des données haut-débit de type ChIP-seq, par exemple avec des données SELEX, et il existe des bases de données stockant les PWMs correspondantes : JAPSAR, base de donnée publique⁷, et TRANSFAC, qui marche par abonnement⁸. Il est à noter que ces PWMs ayant souvent été

7. <http://jaspar.cgb.ki.se>

8. <http://www.gene-regulation.com/pub/databases.html>

construites à partir d'un faible nombre de sites de fixations et de données *in vitro*, elles peuvent être relativement inadaptées à l'analyse de données *in vivo*.

1.7.3 Outils de visualisation

Afin d'avoir une idée plus claire des événements de régulation qui se déroulent à un locus donné, il existe plusieurs outils de visualisation des annotations génomiques et épigénétiques, que ce soit sur le site du NCBI (<http://www.ncbi.nlm.nih.gov/gene>), sur Ensembl ou sur UCSC Genome Browser. Ce dernier possède notamment l'avantage qu'il est possible d'importer des données personnelles sous un grand nombre de formats, obtenues à partir de la littérature ou à partir de ses propres travaux. Ainsi, nous présentons en figure 1.27 quelques données de ChIP-seq pour des TFs musculaires et pour des marques épigénétiques, ainsi que des prédictions bioinformatiques de sites de fixation conservés pour les homéoprotéines Six réalisée par nos soins. La visualisation sur UCSC Genome Browser permet de rapidement déterminer le mode de régulation putatif du gène *Chrng* : fixation de Six et MyoD au niveau du promoteur et apparition de marques épigénétiques H3K4me1 et H3Ac sur les histones au cours de la différenciation de progéniteurs musculaires.

Par ailleurs, il existe un outil de visualisation complémentaire de ceux cités : le visualiseur de régions conservées au cours de l'évolution ECR Browser (<http://ecrbrowser.dcode.org>), intégrant de nombreux outils bioinformatiques (?). Ce navigateur permet de visualiser la conservation génomique d'un locus donné chez plusieurs espèces plus ou moins lointaines (par exemple souris, homme, vache, grenouille et poisson zèbre) afin de cibler l'étude de la régulation sur des régions extrêmement conservées. Il est ensuite possible d'analyser les séquences ultraconservées sélectionnées en utilisant les motifs de la base de donnée TRANSFAC via l'outil rVISTA (?). Un exemple d'utilisation de cet outil est donné par la découverte de plusieurs régions de régulation fonctionnelles de l'homéoprotéine Six1 possédant une extrême conservation (?).

1.7.4 Le projet ENCODE

Le projet ENCODE (pour *Encyclopedia of DNA Elements*) est un consortium de groupes de recherche internationaux financés par le NHGRI (*National Human Genome Research Institute*) qui a vu le jour afin de systématiser les méthodes permettant l'annotation des génomes et de

Chapitre 1. Introduction générale.

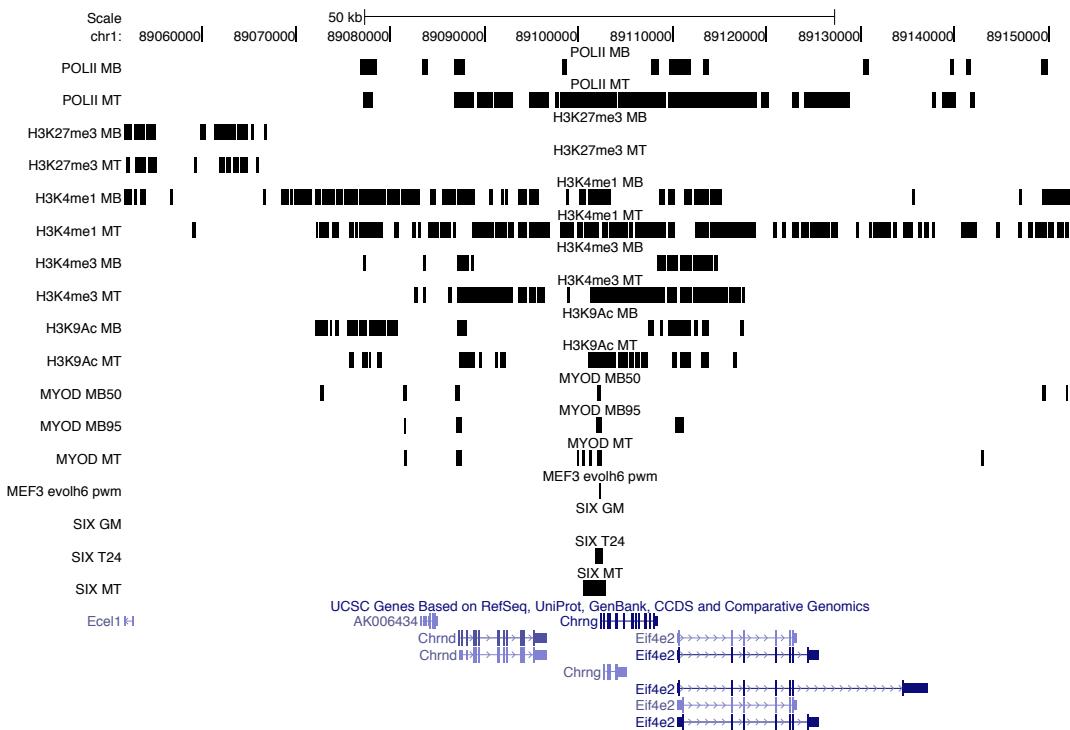


FIGURE 1.27 – Visualisation de données ChIP-seq via le site UCSC.

La visualisation de différentes données ChIP-seq et bioinformatiques (bandes noires) permet de mettre en perspective le cas de la régulation de *Chrng* (en bleu au bas de l'image) lors de la différenciation musculaire. Les données ChIP-seq sont issues de la littérature et les données bioinformatiques (MEF3) ont été obtenues au cours de cette thèse. Les rectangles noirs correspondent aux coordonnées des pics obtenus après avoir appliqué un seuil de filtrage du bruit. Les données de fixation de PolII ainsi que les données de méthylation et d'acétylation des histones (H3K4me3 et H3K9Ac, marques de l'activité transcriptionnelle, voir fig. 1.14), tirées de ?, indiquent que le locus est transcrit lors de la formation de myotubes. Le TF MyoD est fixé au niveau du promoteur de *Chrng* au cours de la prolifération des myoblastes (MB50, 50% de confluence, et MB95, 95% de confluence) et au cours de la différenciation en myotubes MT (?). Le TF Six est co-fixé avec MyoD lors de la différenciation : à T24, 24h après différenciation, et à MT (?). De plus, des analyses bioinformatiques montre l'existence dans cette région d'un site de fixation MEF3 pour la protéine Six conservé chez les vertébrés, corroborant une liaison directe de l'ADN par Six. Prises ensemble, la simple visualisation de ces données suggèrent une régulation par Six et MyoD de *Chrng*.

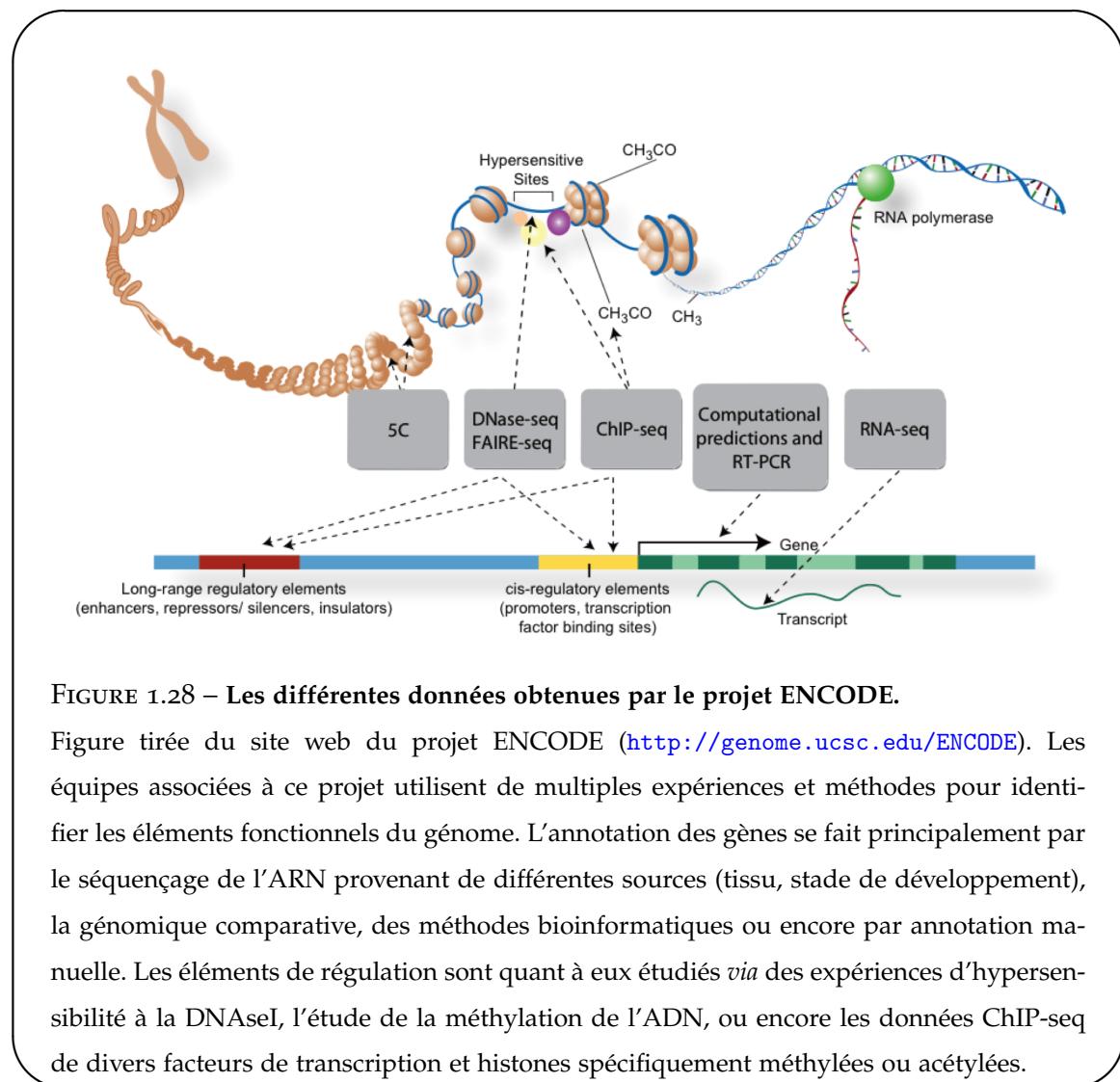


FIGURE 1.28 – Les différentes données obtenues par le projet ENCODE.

Figure tirée du site web du projet ENCODE (<http://genome.ucsc.edu/ENCODE>). Les équipes associées à ce projet utilisent de multiples expériences et méthodes pour identifier les éléments fonctionnels du génome. L'annotation des gènes se fait principalement par le séquençage de l'ARN provenant de différentes sources (tissu, stade de développement), la génomique comparative, des méthodes bioinformatiques ou encore par annotation manuelle. Les éléments de régulation sont quant à eux étudiés *via* des expériences d'hypersensibilité à la DNaseI, l'étude de la méthylation de l'ADN, ou encore les données ChIP-seq de divers facteurs de transcription et histones spécifiquement méthylées ou acétylées.

faciliter l'intégration des nombreuses données obtenues. Son but est de construire une liste exhaustive des éléments fonctionnels du génome humain, qu'ils agissent au niveau de l'ADN, de l'ARN ou des protéines, et des éléments de régulation qui contrôlent l'état cellulaire et l'activité des gènes. Les données sont mises à disposition du public gratuitement sur internet (<http://genome.ucsc.edu/ENCODE/>). À noter que des projets équivalents existent pour d'autres organismes, comme la souris (<http://mouseencode.org>), ou encore le ver *Caenorhabditis elegans* et la mouche *Drosophila melanogaster* (<http://www.modencode.org>).

Totalisant en septembre 2012 plus de 1600 expériences dans plus de 147 types cellulaires, les premières conclusions pointent vers une profusion d'événements de régulation, loin de l'idée d'ADN poubelle (*junk DNA*) : ainsi, 80% du génome est associé à un événement bio-

Chapitre 1. Introduction générale.

chimique associé à de la formation d'ARN ou au remodelage de la chromatine, $\sim 400,000$ régions possèdent un état chromatinien caractéristique des enhancers et $\sim 70,000$ des promoteurs (?). Depuis mai 2013, les données ChIP-seq de 161 TFs couvrant 91 types cellulaires ont été mises à disposition sur UCSC Genome Browser⁹.

9. <http://genome.ucsc.edu/cgi-bin/hgTrackUi?db=wgEncodeAwgTfbsUniform>

Chapitre 2

Modèles de fixation des Facteurs de Transcription à l'ADN.

2.1	Observations de corrélations au sein des TFBS	66
2.2	Modèles existants permettant de décrire la statistique des TFBS	67
2.2.1	Modèle de référence sans corrélations : la PWM	67
2.2.2	Une PWM généralisée : le modèle GWM	69
2.2.3	Réseaux bayésiens	70
2.2.4	Modèles de mélange	71
2.3	Modèles de maximum d'entropie	72
2.3.1	Pourquoi maximiser l'entropie ?	72
2.3.2	Maximisation de l'entropie sous contraintes	74
2.3.3	Application aux sites de fixation	75
2.4	Article	77
2.5	Analyse thermodynamique des modèles	109
2.5.1	Chaleur spécifique	109
2.5.2	Lien avec les valeurs des champs et des couplages	110
2.6	Conclusion et perspectives	112

Introduction du chapitre 2

Dans cette partie, nous nous intéressons à la description de l'interaction entre les facteurs de transcription et leurs sites de reconnaissance sur l'ADN. Pendant longtemps, la qualité de cette description a été limitée par la quantité de données disponibles. Ainsi, les expériences de type SELEX (voir 1.4.1), où des expériences de ChIP au cas par cas permettaient de récupérer de l'ordre de quelques dizaines de sites de fixation pour un TF d'intérêt. Or, le modèle PWM, qui est le modèle le plus simple (en terme de nombre de paramètres) que l'on puisse bâtrir pour décrire l'interaction possède déjà plusieurs dizaines de paramètres – les fréquences des nucléotides à chaque position –.

Ces données ne permettaient donc pas d'explorer plus en avant des modèles plus complexes de fixation incluant par exemple des termes d'interaction entre nucléotides au sein des sites de fixation. Cependant, les avancées récentes en séquençage à haut débit ont permis l'obtention de données très grande échelle, que ce soit *in vivo* par ChIP-seq ou *in vitro* par HT-SELEX (voir 1.4). Le nombre de sites de fixation obtenus est de l'ordre de quelques milliers, ce qui permet de contraindre des modèles de fixation plus complexe que le modèle PWM.

En utilisant des données ChIP-seq pour un grand nombre de facteurs de transcription de la Drosophile et des vertébrés, nous avons contraint différents modèles de fixation incluant implicitement ou explicitement des interactions entre nucléotides. Nous les avons comparés sur leur capacité à décrire les statistiques de fixation TF-ADN observées *in vivo*. Nous présentons préalablement un survol des observations et modèles existant au sujet des corrélations dans les sites de fixations de facteurs de transcription.

2.1 Observations de corrélations au sein des TFBS

Différents travaux ont mis en exergue l'existence de corrélations entre nucléotides au sein des sites de fixation de TFs. Parce que limitées par la quantité de données alors possible d'obtenir, les premières études de ce genre ont centré leur attention sur quelques corrélations importantes pour des cas particuliers. Ainsi, ? ont observé que la protéine Mnt induit des corrélations entre les positions 16 et 17 de ses sites de reconnaissance *in vitro*. Ils ont mesuré expérimentalement la spécificité aux sites de liaisons contenant tous les variants possibles à ces deux positions. Ils ont ainsi observé que la mutation de la base consensus C en position 17 induisait un changement de préférence en position 16 de la base A vers la base C. Par

 2.2. Modèles existants permettant de décrire la statistique des TFBS

ailleurs, ? ont montré que la protéine EGR1 induisait des corrélations au sein d'un triplet de nucléotides central de leur site de reconnaissance. La prise en compte de ces corrélations dans l'énergie de fixation permettait alors d'améliorer la description des données par rapport au modèle additif PWM.

À une plus grande échelle, ? ont utilisé des puces à ADN (technique PBM, cf 1.4.1) pour étudier la fixation *in vitro* de 104 TFs de la souris sur toutes les séquences d'ADN de 10 bp possibles. Pour chaque facteur, plusieurs centaines de séquences de fixation ont ainsi été obtenues. L'étude a révélé l'existence d'une multiplicité de motifs (PWMs) pour la plupart des TFs (seulement 15 étant mieux décrit par un motif unique). Certains motifs reconnaissent notamment des séquences à espacement variable pour lesquelles deux régions spécifiques du site sont séparées par un nombre variable de nucléotides. Enfin, les auteurs ont noté la présence de corrélations fortes dans 19 cas, celles-ci n'étant pas forcément limitées à des dinucléotides mais pouvant impliquer des trinucléotides. Plus récemment, ? ont analysé par HT-SELEX plusieurs centaines de domaines de fixations à l'ADN de TFs humains et de la souris, révélant aussi l'importance d'espacements variables et surtout des corrélations dinucléotidiques entre plus proches voisins.

2.2 Modèles existants permettant de décrire la statistique des TFBS

Différents modèles ont été proposés pour décrire ces corrélations (fig. 2.1). La méthode la plus directe consiste à partir du modèle PWM (fig. 2.1a) et à ajouter des corrélations mutuellement exclusives aux positions les plus corrélées (fig. 2.1b). D'autres méthodes utilisent des structures probabilistes de dépendances sous forme de chaînes de Markov (fig. 2.1c) ou plus généralement de réseau bayésien ou (fig. 2.1d-e). Enfin, une dernière méthode consiste à réaliser des mélanges de modèles afin de capturer des ensembles distincts de corrélations (fig. 2.1f-g).

2.2.1 Modèle de référence sans corrélations : la PWM

Nous l'avons vu, le modèle le plus simple (en termes de nombre de paramètres) décrivant l'interaction entre un TF et son site de reconnaissance sur l'ADN consiste à faire l'hypothèse que les nucléotides contribuent indépendamment à l'énergie de fixation. Cette hypothèse

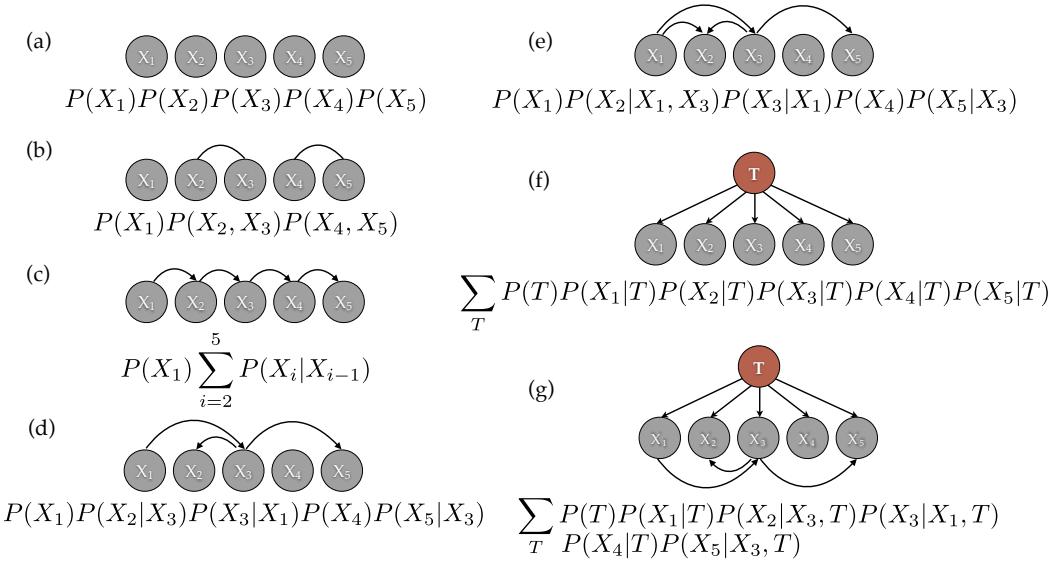


FIGURE 2.1 – Différents modèles pour décrire les corrélations entre nucléotides dans les sites de fixation de facteurs de transcription.

Exemples illustrant différents modèles de fixation sur un site de longueur 5. Pour chaque modèle, la structure du réseau de dépendances sous-jacent est représentée, ainsi que la distribution de probabilité $P(X_1, X_2, X_3, X_4, X_5)$ correspondante, où X_i est une variable aléatoire prenant la valeur (A, C, G ou T) du nucléotide à la position i . Les modèles représentés sont les suivants : (a) PWM (pas de corrélations), (b) GWM (corrélations mutuellement exclusives), (c) chaîne de Markov d'ordre 1 (corrélations entre plus proches voisins), (d) réseau bayésien en arbre (au plus un parent par nœud) ou (e) pas en arbre (le nœud 2 a deux parents), (f) mélange de PWMs et (e) mélange d'arbres à dépendances fixées.

conduit au modèle PWM (section 1.3.2 et fig.2.1a), qui s'écrit¹⁰ :

$$P(X_1, \dots, X_k) = \prod_{i=1}^K P(X_i) \quad (2.1)$$

où $P(X_i)$ est la probabilité marginale d'observer le nucléotide $X \in \{A, C, G, T\}$ à la position i . Un tel modèle possède $3K$ paramètres – 3 paramètres $P(X_i)$ par position, la normalisation des probabilités permettant de fixer le paramètre restant –. Pour une longueur de site typique $K = 10$, le modèle PWM contient 30 paramètres à contraindre, sachant qu'un « modèle »

¹⁰. Comme nous l'avons signalé en 1.3.2, le terme PWM (*Position Weight Matrix*) réfère en fait à la matrice des poids $\log(P(X_i)/\pi_{X_i})$ où π_{X_i} est une distribution neutre indépendante de la position (dite distribution *background*), par exemple calculée sur des régions intergéniques.

complet paramétrant la distribution jointe sans faire d'hypothèse comporterait $4^{10} - 1 \sim 10^6$ paramètres.

2.2.2 Une PWM généralisée : le modèle GWM

Une première méthode permettant de complexifier le modèle PWM consiste à intégrer explicitement des groupes mutuellement exclusif¹¹ de nucléotides corrélés au sein du modèle (fig. 2.1b). Une telle méthode fut d'abord employée par ? pour prendre en compte des corrélations préalablement définies entre nucléotides plus proches voisins. De manière plus générale, ? ont développé un modèle de matrice de poids généralisée (GWM pour *Generalized Weight Matrix*) qui prend en compte de manière systématique les corrélations permettant d'améliorer le modèle indépendant. Pour ce faire, les auteurs utilisent une méthode de Monte-Carlo par chaîne de Markov (MCMC) : des corrélations sont ajoutées ou enlevées au hasard au modèle et acceptées selon la règle de Metropolis-Hastings (?). Cette acceptation est proportionnelle au facteur de Bayes, une quantité qui permet de comparer des modèles possédant des nombres de paramètres différent¹². Ce facteur est défini par le rapport entre la probabilité de générer les données D (les séquences de fixation) avec un modèle M_1 de paramètres θ_1 plutôt qu'avec un autre modèle M_2 de paramètres θ_2 :

$$BF = \frac{P(D|M_1)}{P(D|M_2)} = \frac{\int P(D|\theta_1, M_1)P(\theta_1|M_1)d\theta_1}{\int P(D|\theta_2, M_2)P(\theta_2|M_2)d\theta_2} \quad (2.2)$$

Le modèle final consiste en un ensemble de paramètres décrivant des positions indépendantes et des positions corrélées mutuellement exclusives . En analysant les données TRANSFAC, les auteurs ont noté que dans 25% des cas (22/95) le modèle GWM était significativement meilleur que le modèle PWM (facteur de Bayes supérieur à 6).

Cette méthode a par la suite été utilisée sur des données ChIP-seq pour 4 TFs mammifères – NRSF, STAT1, CTCF et ER – (?). En utilisant les 10% des pics les plus importants comme ensemble d'apprentissage et en se restreignant aux régions de 200bp centrées autour du sommet du pic ChIP, les auteurs ont réalisé un échantillonnage de Gibbs (?) pour obtenir les sites de fixation suivant les hypothèses que (1) chaque pic contient au plus un seul site de fixation (modèle ZOOPS pour *Zero or One Occurrences Per Sequence*), (2) la probabilité *a priori*

11. Les corrélations entre des couples de positions (i,j) et (j,k) ne peuvent être admises au sein du même modèle.

12. Sous certaines approximation, ce facteur peut se rapporter à une différence de valeurs du BIC, introduit dans l'article en 2.4.

d'avoir un site à une certaine position sur la séquence est plus forte autour du sommet du pic, et (3) les sites sont décrits par un modèle GWM. L'étude a révélé l'existence de corrélations fortes limitées aux nucléotides plus proches voisins dans les quatre cas étudiés. Les nucléotides participant aux corrélations se situaient à des positions ayant un faible contenu en information dans le modèle PWM. Enfin, les auteurs ont noté la présence de plusieurs triplets de nucléotides voisins corrélés.

2.2.3 Réseaux bayésiens

Une généralisation du modèle GWM consiste à supprimer la condition d'exclusion mutuelle des paires de nucléotides corrélés en décrivant de manière plus générale le réseau de dépendance entre positions. Une telle description est possible en utilisant le langage des réseaux bayésiens. Les dépendances y sont représentées par un graphe orienté acyclique¹³ G , dont les nœuds sont les variables X_i et les liens représentent les conditionnements d'une variable avec des variables parentes (fig. 2.1e). La probabilité jointe s'écrit :

$$P(X_1, \dots, X_k) = \prod_{i=1}^K P(X_i | P_i^G) \quad (2.3)$$

où P_i^G est l'ensemble (pouvant être vide) des parents de X_i dans G . Le nombre de paramètres peut rapidement devenir grand : si l'on note N_i le nombre de parents de X_i , alors le nombre de paramètres du modèle est $3 \sum_{i=1}^K 4^{N_i}$.

Lorsque les différents nœuds possèdent au plus un parent, on parle d'arbre bayésien (fig. 2.1d). Ce type de réseau bayésien généralise notamment le cas des chaînes de Markov d'ordre 1, où chaque nœud dépend du nœud précédent (fig. 2.1c). Le nombre de paramètres est alors restreint, puisqu'il est au plus de $3 \cdot 4K$.

L'avantage des arbres bayésiens est qu'il existe des algorithmes efficaces permettant de trouver la meilleure structure d'arbre (?). De tels modèles d'arbres ont été utilisés pour décrire les données de 95 TFs de Transfac (?). Dans $\sim 25\%$ des cas (22/95), le modèle d'arbre bayésien s'avère significativement meilleur qu'un modèle PWM, ce qui est du même ordre de grandeur que pour le modèle GWM vu en 2.2.2.

¹³. Un graphe orienté acyclique est un réseau dont les liens sont orientés et au sein duquel il n'est pas possible de revenir à son point de départ en suivant les flèches

2.2.4 Modèles de mélange

Dans les cas précédents, nous avons présenté des modèles capturant des dépendances « locales » entre paires de nucléotides. Néanmoins, il peut exister des dépendances plus largement réparties entre les positions, comme cela a déjà été observé empiriquement (??). De telles corrélations à plus grande échelle peuvent être modélisées en supposant que le facteur de transcription possède plusieurs « modes » de fixation. Ceux-ci peuvent par exemple correspondre à différentes conformations de la protéine sur son site de fixation, chaque configuration possédant ses propres préférences de fixation. Ces modes sont décrits par une variable aléatoire T (le *type* de fixation) de probabilité $P(T)$. Il est ensuite possible de décrire la fixation au sein de chaque mode par l'un des modèles précédents.

- **Mélange de PWMs**

Le cas le plus naturel consiste à utiliser comme modèle de fixation le modèle PWM, c'est-à-dire que dans chaque mode il y a indépendance entre les positions. La probabilité d'observer un site est alors donnée par la somme sur les différents modes de fixation de la probabilité de fixer un site, conditionnée par la probabilité d'être dans ce mode :

$$P(X_1, \dots, X_K) = \sum_{T=1}^N P(T) \prod_{i=1}^K P(X_i|T) \quad (2.4)$$

où N est le nombre de modes de fixation. Ce modèle a plusieurs avantages. D'abord, le nombre de paramètres reste linéaire en K : pour décrire $P(T)$ et les N PWMs il faut $N - 1 + 3KN$ paramètres. Ce nombre reste donc raisonnablement faible devant le nombre de paramètres requis pour complètement décrire les interactions à deux nucléotides, qui croît comme K^2 . Ensuite, le modèle a une interprétation claire qui peut permettre de mettre en exergue un mécanisme biologique sous-jacent.

Ce type de modèle permet de dépasser le modèle PWM dans un nombre substantiel de cas. Ainsi, ? ont montré que ~ 40% des TFs de Transfac (36/95) sont significativement mieux représentés par un mélange de 2 PWMs que par une seule PWM. En utilisant des données *in vitro* plus précises pour 104 TFs de la souris, ? ont montré que ~ 85% (89/104) étaient mieux représenté par une combinaison de PWMs que par une PWM seule, plaidant pour un portée générale de l'existence de « motifs secondaires ».

- **Mélange d'arbres**

De la même manière que les PWMs, il est possible d'étendre les modèles d'arbres en réalisant un mélange d'arbres. Intuitivement, ceci permet de capturer des dépendances additionnelles en gardant un nombre de paramètres linéaire en fonction de la taille du motif. Un tel modèle semble posséder des performances comparables au mélange de PWM, et améliore la description des TFs de Transfac dans $\sim 40\%$ des cas (35/95) (?).

2.3 Modèles de maximum d'entropie

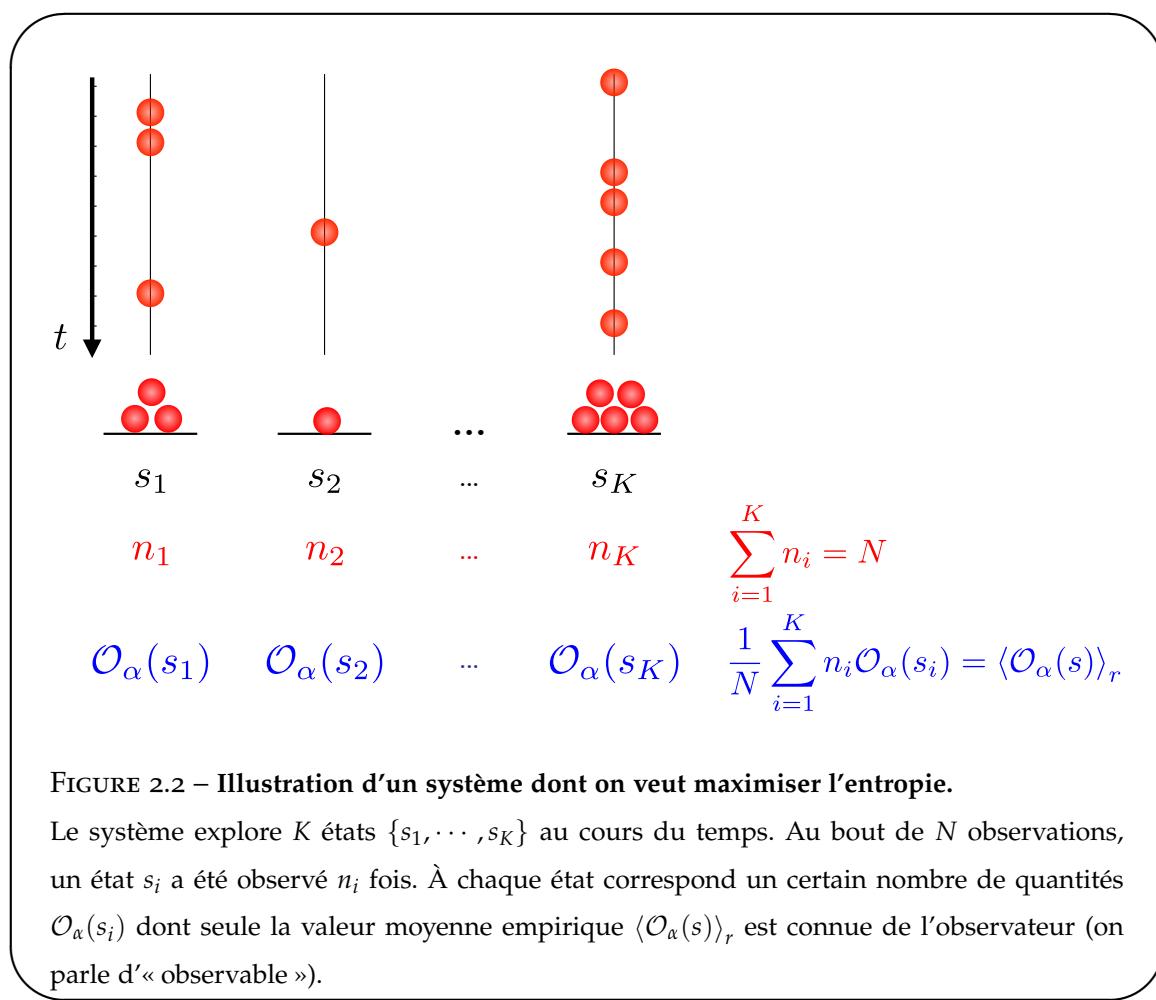


FIGURE 2.2 – Illustration d'un système dont on veut maximiser l'entropie.

Le système explore K états $\{s_1, \dots, s_K\}$ au cours du temps. Au bout de N observations, un état s_i a été observé n_i fois. À chaque état correspond un certain nombre de quantités $\mathcal{O}_\alpha(s_i)$ dont seule la valeur moyenne empirique $\langle \mathcal{O}_\alpha(s) \rangle_r$ est connue de l'observateur (on parle d'« observable »).

2.3.1 Pourquoi maximiser l'entropie ?

Le concept d'entropie remonte aux prémisses de la physique statistique (?). Dans l'essence, il peut être compris de la manière suivante. Supposons qu'un système comporte K états dis-

tincts $\{s_1, \dots, s_K\}$. Au cours du temps, le système explore les différents états (fig. 2.2). Au bout de N observations, chaque état a été observé un nombre n_i de fois. La question sous-jacente au calcul de l'entropie est la suivante : sans connaissance *a priori* sur le système, que puis-je dire de ces n_i ? Prenons l'exemple de la figure 2.2. On a certaines valeurs pour les n_i ($n_1 = 3$, $n_2 = 1$, etc.), et on aimeraient savoir de combien de manières il est possible de réaliser un tel ensemble de valeurs. Notons ce nombre $\mathcal{N}(n_1, \dots, n_K)$. Il est donné par la formule suivante :

$$\begin{aligned}\mathcal{N}(n_1, \dots, n_K) &= \binom{N}{n_1} \binom{N-n_1}{n_2} \cdots \binom{N-\sum_{i=1}^{K-1} n_i}{n_K} \\ &= \frac{N!}{(N-n_1)!n_1!} \times \frac{(N-n_1)!}{(N-n_1-n_2)!n_2!} \times \cdots \times \frac{(N-\sum_{i=1}^{K-1} n_i)!}{0!n_K!}\end{aligned}\quad (2.5)$$

soit

$$\mathcal{N}(n_1, \dots, n_K) = \frac{N!}{n_1!n_2!\cdots n_K!} \quad (2.6)$$

Il convient alors de s'intéresser au logarithme de cette quantité. En effet, dans le cas où les nombres d'observation sont grands $n_i \gg 1$, ceux-ci s'expriment simplement grâce à la formule de Stirling :

$$\log(n!) \xrightarrow{n \rightarrow \infty} n \log(n) - n \quad (2.7)$$

On peut alors écrire

$$\begin{aligned}\log \mathcal{N}(n_1, \dots, n_K) &= N \log(N) - N - \sum_{i=1}^K (n_i \log(n_i) - n_i) \\ &= \sum_{i=1}^K n_i \log\left(\frac{N}{n_i}\right) \\ &= -N \sum_{i=1}^K \frac{n_i}{N} \log\left(\frac{n_i}{N}\right)\end{aligned}\quad (2.8)$$

On note l'apparition des probabilités empiriques $f(s_i) = n_i/N$ d'observer l'état s_i , qui tendent asymptotiquement (dans la limite « thermodynamique » $N \rightarrow \infty$) vers les « vraies » probabilités $P(s_i)$. L'entropie est définie dans cette limite comme étant égale à $1/N \log \mathcal{N}(n_1, \dots, n_K)$, soit

$$S[P] = - \sum_{\{s\}} P(s) \log P(s) \quad (2.9)$$

où $\{s\} = \{s_1, \dots, s_K\}$ dénote l'ensemble des états accessibles. L'idée est alors la suivante : nous souhaitons savoir quels états le système a le plus probablement visité au cours des N transitions. Sans connaissance *a priori* sur le système, il est plus probable que les nombres (n_1, \dots, n_K) obtenus soient ceux qui sont réalisés le plus souvent, c'est-à-dire ceux qui maximisent la quantité $\mathcal{N}(n_1, \dots, n_K)$ et donc au final l'entropie. Par ailleurs, les fluctuations relatives des quantités n_i sont de l'ordre de $1/\sqrt{n_i}$ (?). Ainsi, la solution de maximum d'entropie domine largement les autres solutions possibles dans la limite thermodynamique.

2.3.2 Maximisation de l'entropie sous contraintes

Notons $\mathcal{O}_\alpha(s)$ une quantité attachée à s (fig. 2.2). En thermodynamique, une telle quantité correspond par exemple à l'énergie d'un état. L'observateur n'a lui accès qu'aux valeurs moyennes de telles quantités sous-jacentes. À l'état d'équilibre thermodynamique, l'échantillonnage des états est réalisé au sein de la distribution de probabilité de maximum d'entropie $P(s)$, et les valeurs moyennes calculées avec les fréquences empiriques $f(s)$ doivent donc être compatibles avec les valeurs moyennes calculées avec la distribution de probabilité $P(s)$:

$$\sum_{\{s\}} P(s) \mathcal{O}_\alpha(s) = \sum_{\{s\}} f(s) \mathcal{O}_\alpha(s) \quad (2.10)$$

Nous souhaiterions maintenant connaître la distribution $P(s)$ la moins biaisée (i.e de maximum d'entropie) qui satisfait les contraintes de l'éq. 2.10 imposées par l'observation des données (l'information que possède l'observateur). Ce problème revient à maximiser le Lagrangien suivant :

$$\mathcal{L} = - \sum_{\{s\}} P(s) \log P(s) + \lambda \left(\sum_{\{s\}} P(s) - 1 \right) + \sum_\alpha \beta_\alpha \sum_{\{s\}} (P(s) - f(s)) \mathcal{O}_\alpha(s) \quad (2.11)$$

où les paramètres λ et β_α sont les multiplicateurs de Lagrange correspondant respectivement à la contrainte de normalisation de la distribution de probabilité et aux informations qu'a l'observateur sur certaines valeurs moyennes du système (éq. 2.10). La maximisation de ce Lagrangien est obtenue en annulant la dérivée fonctionnelle par rapport à la distribution de probabilité $P(s)$:

$$\frac{\delta \mathcal{L}}{\delta P(s)} = 0 = -\ln P(s) - 1 + \lambda + \sum_\alpha \beta_\alpha \mathcal{O}_\alpha(s) \quad (2.12)$$

En utilisant la normalisation des probabilités, il est possible de trouver λ , et la solution se met finalement sous la forme

$$P(s) = \frac{1}{\mathcal{Z}} e^{-\mathcal{H}(s)} \quad (2.13)$$

où \mathcal{H} est l'Hamiltonien du système :

$$\mathcal{H} = \sum_{\alpha} \beta_{\alpha} \mathcal{O}_{\alpha}(s) \quad (2.14)$$

et \mathcal{Z} est la fonction de partition permettant la normalisation de la distribution $P(s)$:

$$\mathcal{Z} = \sum_{\{s\}} \exp[-\mathcal{H}(s)]. \quad (2.15)$$

Remarque

Il est possible de montrer que la maximisation de l'entropie, partant des contraintes de l'éq. 2.10 sur les valeurs moyennes pour en arriver à une forme exponentielle de la distribution de probabilité, est le contrepoint d'une maximisation de la vraisemblance partant d'une forme exponentielle pour en arriver aux mêmes contraintes sur les valeurs moyennes (??).

2.3.3 Application aux sites de fixation

- **Corrélations à un point : le modèle PWM**

Dans le cas qui nous intéresse, un état s correspond à une séquence d'ADN appartenant à l'ensemble $\{s\}$ des sites de fixation d'un facteur de transcription. Considérons maintenant l'observable quantifiant la présence du nucléotide a à la position i d'un site :

$$\mathcal{O}_{i,a}(s) = \delta(s_i, a) \quad (2.16)$$

où δ est la fonction de Kronecker qui vaut 1 lorsque le nucléotide à la position i du site s_i vaut a et 0 sinon. De cette définition il suit que la valeur moyenne sur les fréquences empiriques

$$\sum_{\{s\}} f(s) \mathcal{O}_{i,a}(s) = f_{i,a} \quad (2.17)$$

se réduit à la fréquence du nucléotide a à la position i . Notons $h_i(a)$ le multiplicateur de Lagrange correspondant et $\mathcal{A} = \{A, C, G, T\}$. On trouve alors

$$\begin{aligned}\mathcal{H}(s) &= \sum_{i=1}^L \sum_{a \in \mathcal{A}} h_i(a) \delta(s_i, a) \\ &= \sum_{i=1}^L h_i(s_i)\end{aligned}\tag{2.18}$$

Les différentes positions étant indépendantes, la fonction de partition \mathcal{Z} peut par ailleurs se scinder en différentes fonctions de partition par position : $\mathcal{Z} = \prod_{i=1}^L \mathcal{Z}_i$. On obtient au final

$$P(s) = \frac{1}{\mathcal{Z}} e^{-\sum_{i=1}^L h_i(s_i)} = \prod_{i=1}^L \frac{e^{-h_i(s_i)}}{\mathcal{Z}_i}\tag{2.19}$$

On retrouve le modèle PWM introduit dans l'éq. 2.13.

- **Corrélations à deux points : le modèle de Potts**

Il est maintenant relativement direct de complexifier le modèle en ajoutant l'observation des couples d'interaction au sein des sites de fixation :

$$\mathcal{O}_{i,a,j,b}(s) = \delta(s_i, a) \delta(s_j, b)\tag{2.20}$$

La corrélation à deux points entre le nucléotide a en position i et b en position j s'écrit donc

$$\sum_{\{s\}} f(s) \mathcal{O}_{i,a,j,b}(s) = f_{i,a,j,b}\tag{2.21}$$

où $f_{i,a,j,b}$ est la fréquence empirique d'observation de la paire de nucléotide (a, b) aux positions (i, j) . Notons $J_{i,j}(a, b)$ le multiplicateur de Lagrange correspondant. L'Hamiltonien sous les contraintes imposées par les équations 2.17 et 2.21 s'écrit :

$$\begin{aligned}\mathcal{H}(s) &= \sum_{i=1}^L \sum_{a \in \mathcal{A}} h_i(a) \delta(s_i, a) + \sum_{i=1}^{L-1} \sum_{j>i} \sum_{a \in \mathcal{A}} \sum_{b \in \mathcal{A}} J_{i,j}(a, b) \delta(s_i, a) \delta(s_j, b) \\ &= \sum_{i=1}^L h_i(s_i) + \sum_{i=1}^{L-1} \sum_{j>i} J_{i,j}(s_i, s_j)\end{aligned}\tag{2.22}$$

Le modèle de maximum d'entropie est finalement

$$P(s) = \frac{1}{\mathcal{Z}} e^{-\sum_{i=1}^L h_i(s_i) - \sum_{i=1}^{L-1} \sum_{j>i} J_{i,j}(s_i, s_j)}\tag{2.23}$$

On reconnaît le modèle de Potts inhomogène de champs magnétiques locaux h_i et de termes d'interaction $J_{i,j}$ couramment utilisé dans la description des verres de spins (?).

2.4 Article

L'article qui suit décrit l'analyse de données de fixation *in vivo* à grande échelle pour plusieurs TFs drosophiles et mammifères. Différents modèles sont comparés, incluant ou non des dépendances : un modèle PWM, un modèle de mélange de PWMs, et un modèle de Potts.

Beyond position weight matrices: nucleotide correlations in transcription factor binding sites and their description

Marc Santolini, Thierry Mora, and Vincent Hakim
Laboratoire de Physique Statistique, CNRS, Université P. et M. Curie,
Université D. Diderot, École Normale Supérieure, Paris, France.

The identification of transcription factor binding sites (TFBSs) on genomic DNA is of crucial importance for understanding and predicting regulatory elements in gene networks. TFBS motifs are commonly described by Position Weight Matrices (PWMs), in which each DNA base pair independently contributes to the transcription factor (TF) binding, despite mounting evidence of interdependence between base pairs positions. The recent availability of genome-wide data on TF-bound DNA regions offers the possibility to revisit this question in detail for TF binding *in vivo*. Here, we use available fly and mouse ChIPseq data, and show that the independent model generally does not reproduce the observed statistics of TFBS, generalizing previous observations. We further show that TFBS description and predictability can be systematically improved by taking into account pairwise correlations in the TFBS via the principle of maximum entropy. The resulting pairwise interaction model is formally equivalent to the disordered Potts models of statistical mechanics and it generalizes previous approaches to interdependent positions. Its structure allows for co-variation of two or more base pairs, as well as secondary motifs. Although models consisting of mixtures of PWMs also have this last feature, we show that pairwise interaction models outperform them. The significant pairwise interactions are found to be sparse and found dominantly between consecutive base pairs. Finally, the use of a pairwise interaction model for the identification of TFBSs is shown to give significantly different predictions than a model based on independent positions.

Author Summary

Transcription factors are proteins that bind on DNA to regulate several processes such as gene transcription or epigenetic modifications. Being able to predict the Transcription Factor Binding Sites (TFBSs) with accuracy on a genome-wide scale is one of the challenges of modern biology, as it allows for the bottom-up reconstruction of the gene regulatory networks. The description of the TFBSs has been to date mostly limited to a simple model, where the affinity of the protein for DNA, or binding energy, is the sum of independent contributions from uncorrelated amino-acids bound on base pairs. However, structural aspects are of prime importance in proteins and could imply appreciable correlations throughout the observed binding sequences. Using a statistical physics inspired description and high-throughput ChIPseq data for a variety of Drosophilae and mammals TFs, we show that such correlations exist and that accounting for their contribution greatly improves the predictability of genomic TFBSs.

Introduction

Gene regulatory networks are at the basis of our understanding of a cell state and of the dynamics of its response to environmental cues. Central effectors of this regulation are Transcription Factors (TF) that bind on short DNA regulatory sequences and interact with the transcription apparatus or with histone-modifying proteins to alter target gene expressions [1]. The determination of Transcription Factor Binding Sites (TFBSs) on a genome-wide scale is thus of importance and is the focus

of many current experiments [2]. An important feature of TF in eukaryotes is that their binding specificity is moderate and that a given TF is found to bind a variety of different sequences *in vivo* [3]. The collection of binding sequences for a TF-DNA is widely described by a Position Weight Matrix (PWM) which simply gives the probability that a particular base pair stands at a given position in the TFBS. The PWM provides a full statistical description of the TFBS collection when there are no correlations between nucleotides at different positions. Provided that the TF concentration is far from saturation, the PWM description applies exactly at thermodynamic equilibrium in the simple case where the different nucleotides in the TFBS contribute independently to the TF-DNA interaction, such that the total binding energy is the mere sum of the individual contributions [4, 5].

Previous works have reported several cases of correlations between nucleotides at different positions in TFBSs [6–9]. A systematic *in vitro* study of 104 TFs using DNA microarrays revealed a rich picture of binding patterns [10], including the existence of multiple motifs, strong nucleotide position interdependence, and variable spacer motifs, where two small determining regions of the binding site are separated by a variable number of base pairs. Recently, the specificity of several hundred human and mouse DNA-binding domains was investigated using high-throughput SELEX. Correlations between nucleotides were found to be widespread among TFBSs and predominantly located between adjacent flanking bases in the TFBS [9]. The relevance of nucleotide correlations remains however debated [11].

On the modeling side, probabilistic models have been proposed to describe these correlations, either by explicitly identifying mutually exclusive groups of co-varying

nucleotide positions [7, 12, 13], or by assuming a specific and tractable probabilistic structure such as Bayesian networks or Markov chains [9, 14, 15]. However, the extent of nucleotide correlations in TFBSs *in vivo* remains to be assessed, and a systematic and general framework that accounts for the rich landscape of observed TF binding behaviours is yet to be applied in this context. The recent breakthrough in the experimental acquisition of precise, genome-wide TF-bound DNA regions with the ChIPseq technology offers the opportunity to address these two important issues. Using a variety of ChIPseq experiments coming both from fly and mouse, we first show that the independent model generally does not reproduce well the observed TFBS statistics for a majority of TF. This calls for a refinement of the PWM description that accounts for interdependence between nucleotide positions.

The general problem of devising interaction parameters from observed state frequencies has been recently studied in different contexts where large amounts of data have become available. These include describing the probability of coinciding spikes [16, 17] or activation sequences [18, 19] in neural data, the statistics of protein sequences [20, 21], and even the flight directions of birds in large flocks [22]. Maximum-entropy models accounting for pairwise correlations in the least constrained way have been found to provide significant improvement over independent models. The PWM description of TF binding is equivalent to the maximum entropy solely constrained by nucleotide frequencies at each position. Thus, we propose, in the present paper, to refine this model by further constraining pairwise correlations between nucleotide positions. This corresponds to including effective pairwise interactions between nucleotides in an equilibrium thermodynamic model of TF-DNA interaction, as already proposed [23]. When enough data are available, the TFBS statistics and predictability are found to be significantly improved in this refined model. We consider, for comparison, a model that describes the statistics of TFBSs as a statistical mixture of PWMs [14] and generalizes previous proposals [24, 25]. This alternative model can directly capture some higher-order correlations between nucleotides but is found to be outperformed for all considered TF by the pairwise interaction model.

We further show that the pairwise interaction model accounts for the different PWMs appearing in the mixture model by studying its energy landscape: each basin of attraction of a metastable energy minimum in the pairwise interaction model is generally dominantly described by one PWM in the mixture model. Significant pairwise interactions between nucleotides are sparse and found dominantly between consecutive nucleotides, in general qualitative agreement with *in vitro* binding results [9]. The proposed model with pairwise interactions only requires a modest computational effort. When enough data are available, it should thus generally prove worth using the refined description of TFBS that it affords.

Results

The PWM model does not reproduce the TFBS statistics

We first tested how well the usual PWM model reproduced the observed TFBS statistics, *i.e.* how well the frequencies of different TFBSs were retrieved by using only single nucleotide frequencies. For this purpose, we used a collection of ChIPseq data available from the literature [26–28], both from *D. Melanogaster* and from mouse embryonic stem cells (ESC) and a myogenic cell line (C2C12). The TFBSs are short L -mers (we take here $L = 12$), which are determined in each few hundred nucleotides long ChIP-bound region with the help of a model of TF binding. One important consequence and specific features of these data, is that the TFBS collection is not independent of the model used to describe it. Thus, in order to self-consistently determine the collection of binding sites for a given TF from a collection of ChIPseq sequences, we iteratively refined the PWM together with the collection of TFBSs in the ChIPseq data (see Figure and *Methods*). This process ensured that the frequency of different nucleotides at a given position in the considered ensemble of binding sites was exactly accounted by the PWM. We then enquired whether the probability of the different binding sequences in the collection agreed with that predicted by the PWM, as would be the case if the probabilities of observing nucleotides at different positions were independent. Figure 2 displays the results for three different TFs, one from each of the three considered categories: Twi (*Drosophila*), Esrrb (mammals, ESC), and MyoD (mammals, C2C12). For each factor, the ten most frequent sequences in the TFBS collection are shown. For comparison, Figure 2 also displays the probabilities for these sequences as predicted by the PWM built from the TFBS collection. The independent PWM model strongly underestimates the probabilities of the most frequent sequences. Moreover, the PWM model does not correctly predict the frequency order of the sequences and attributes comparable probabilities to these different sequences, in contrast to their observed frequencies.

The relative entropy or Kullback-Leibler divergence (DKL) is a general way to measure the difference between two probability distributions [29]. In order to better quantify the differences between the observed binding sequence frequencies and the PWM frequencies, we computed the DKL between these distributions for all the considered TF, as shown in Figure 2D. For each transcription factor T, part of the differences comes from the finite number $N(T)$ of its observed binding sites. The results are thus compared for each factor T to DKLs between the PWM probabilities and frequencies obtained for artificial sequence samples of size $N(T)$ generated with the same PWM probabilities. For most TFs (22 out of 28), the difference between the observed binding sequence frequencies and the PWM frequencies is signifi-

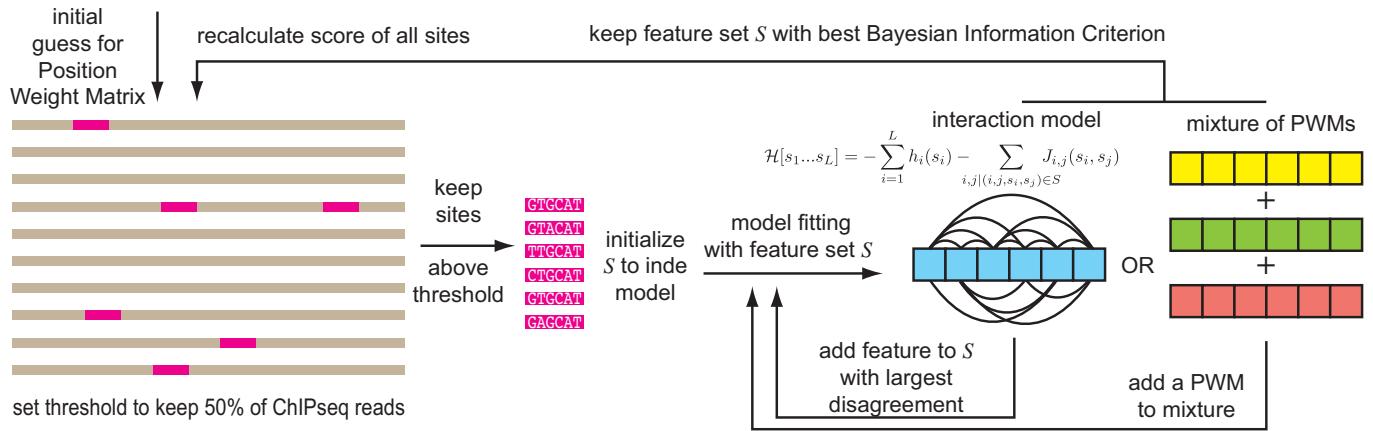


FIG. 1: **Workflow.** An initial Position Weight Matrix (PWM) is used to find a set of binding sites on ChIPseq data. Models are then learned using single-point frequencies (independent), two-point correlations (pairwise) or a mixture of independent models learned on sites clustered by K-Means (mixture) with increasing complexity, *i.e.* increasing number of features in the model. Finally the models with best Bayesian Information Criteria (BIC) are used to predict new sites until convergence to a stable set of sites.

cantly larger than expected from finite size sampling. In the following we focus on these 22 factors for which the PWM description of the TFBSs needs to be refined. It can be noted that the 6 factors for which the PWM description appears satisfactory are predominantly those for which the smallest number of ChIP sequences is available (see Table 1 and Figure S1).

Pairwise interactions in the binding energy improve the TFBS description

The discrepancy between the observed statistics of TFBSs and the statistics predicted by the PWM model calls for a re-evaluation of the PWM main hypothesis, namely the independence of bound nucleotides. As recalled above, the inverse problem of devising interaction parameters from observed frequencies of “words” has been recently studied in different contexts. It has been proposed to include systematically pairwise correlations between the “letters” comprising the words to refine the independent letter description. In the case of a two-letter alphabet, the obtained model is equivalent to the classical Ising model of statistical mechanics[30]. In the present case, the 4-nucleotide alphabet (A,C,G,T) leads to a model equivalent to the so-called inhomogeneous Potts model [30] (hereafter called pairwise interaction model), a generalization of the Ising model to the case where spins assume q values and their fields and interaction parameters depend on the sites considered. In this analogy, nucleotides are spins with $q = 4$ colors.

In practice, the probability of observing a given word

$(s_1 \dots s_L)$ in the dataset is expressed as $P[s_1 \dots s_L] = (1/\mathcal{Z}) \exp(-\mathcal{H}[s_1 \dots s_L])$, where \mathcal{Z} is a normalization constant. \mathcal{H} is formally equivalent to a Hamiltonian in the language of statistical mechanics, and reads:

$$\mathcal{H}[s_1 \dots s_L] = - \sum_{i=1}^L h_i(s_i) - \sum_{i=1}^L \sum_{j < i} J_{i,j}(s_i, s_j), \quad (1)$$

$$s_i \in \{A, C, G, T\}$$

The “magnetic fields” h_i at each site i , along with the interaction parameters J_{ij} between nucleotides at positions i and j , are computed so as to reproduce the frequency of nucleotide usage at each position in the TFBS as well as the pairwise correlations between nucleotides at different positions (see *Methods*). In principle, the number of parameters in the model is sufficient to reproduce the observed values of all pairwise correlations between nucleotides. This however would result in over-fitting the finite-size data with an unrealistically large number of parameters. Therefore, to obtain the model parameters we instead maximized the likelihood that the data was generated by the model with a penalty proportional to the numbers of parameters involved, as provided by the Bayesian Information Criterion (BIC) [31]. Similarly to the procedure followed for the PWM, the pairwise interaction model and the collection of TFBSs for a given factor were iteratively refined together, as schematized in Figure .

Figure 3 shows the improvement in the description of TFBS statistics when using the final pairwise interaction model, for the three factors chosen for illustrative

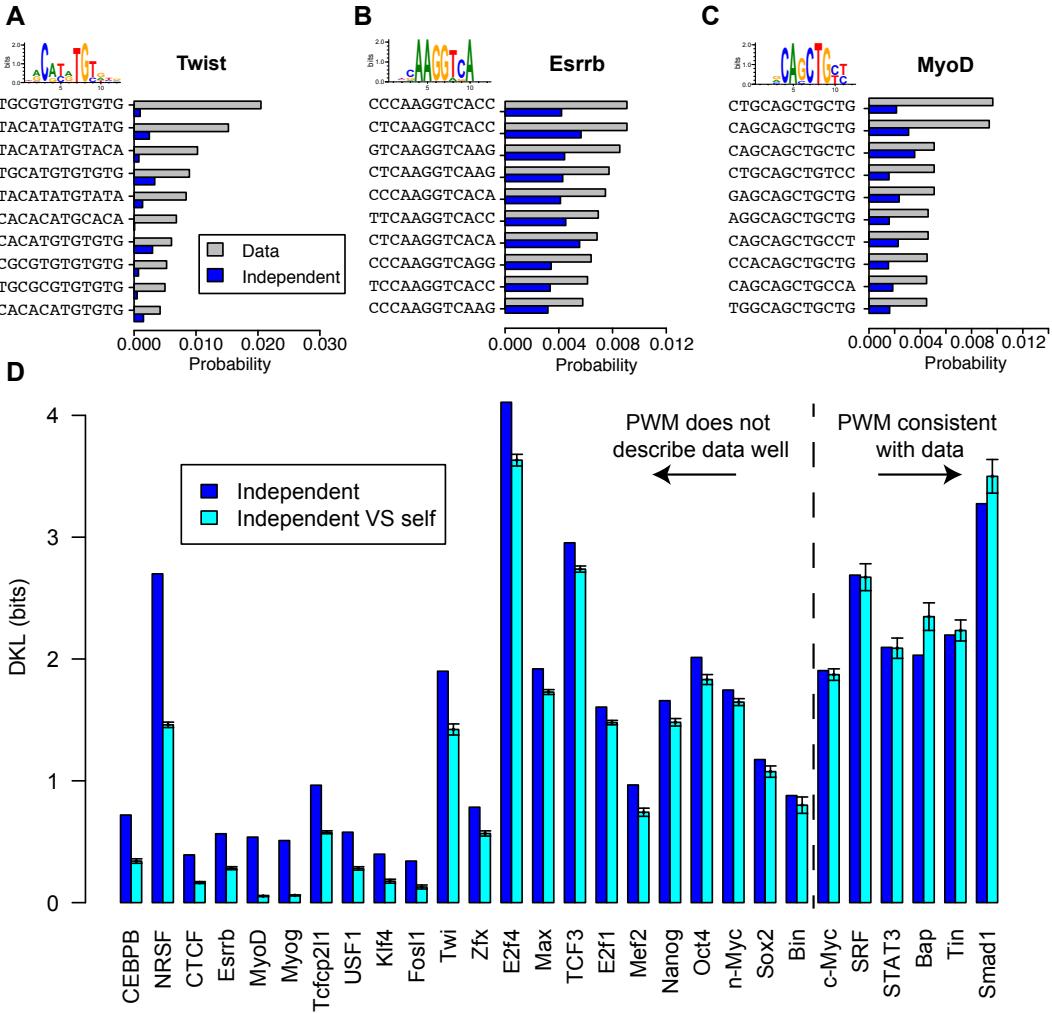


FIG. 2: Observed TFBS frequencies are poorly predicted by a PWM model. The observed frequencies of the most represented binding site sequences for the TF Twist (A), Esrrb (B) and MyoD (C) are shown (gray bars) as well as the probabilities of these sequences as predicted by the PWM model (blue bars). (D) Kullback-Leibler Divergence (DKL) between the observed probability distribution and the independent model distribution (blue). As a control we show the mean (cyan bars) along with two standard deviations of the DKL between the independent model and a finite sample drawn from it (see Methods). A discrepancy between the observed and predicted sequence probabilities is reported for 22 out of 28 factors.

purposes. Where the independent model failed at reproducing the strong amplitude and non-linear decrease in the frequencies of the most over-represented TFBSs, the pairwise interaction model provides a substantial improvement in reproducing the observed statistics. The improvement is most apparent when comparing the frequencies of the ten most observed TFBSs between the model and the ChIPseq data (Figure 3 A, C, E), and is further shown by the statistics of the full collection of TFBSs (Figure 3 B, D, F).

The pairwise model ranks binding sites differently from the PWM

Precise predictions of TFBSs are one important output of ChIPseq data. Moreover, they condition further validation experiments such as gel mobility shift assays or mutageneses. We therefore found it worth assessing the difference in TFBS predictions between pairwise and independent models.

First, we compared the set of ChIP sequences retrieved by the independent and pairwise models model at the cutoff of 50% TPR (True Positive Rate) used in the learning scheme, as shown in Figure 4A. The non overlapping set of ChIPseq sequences (*i.e.* sequences that were picked by one model but not by the other) was found to range

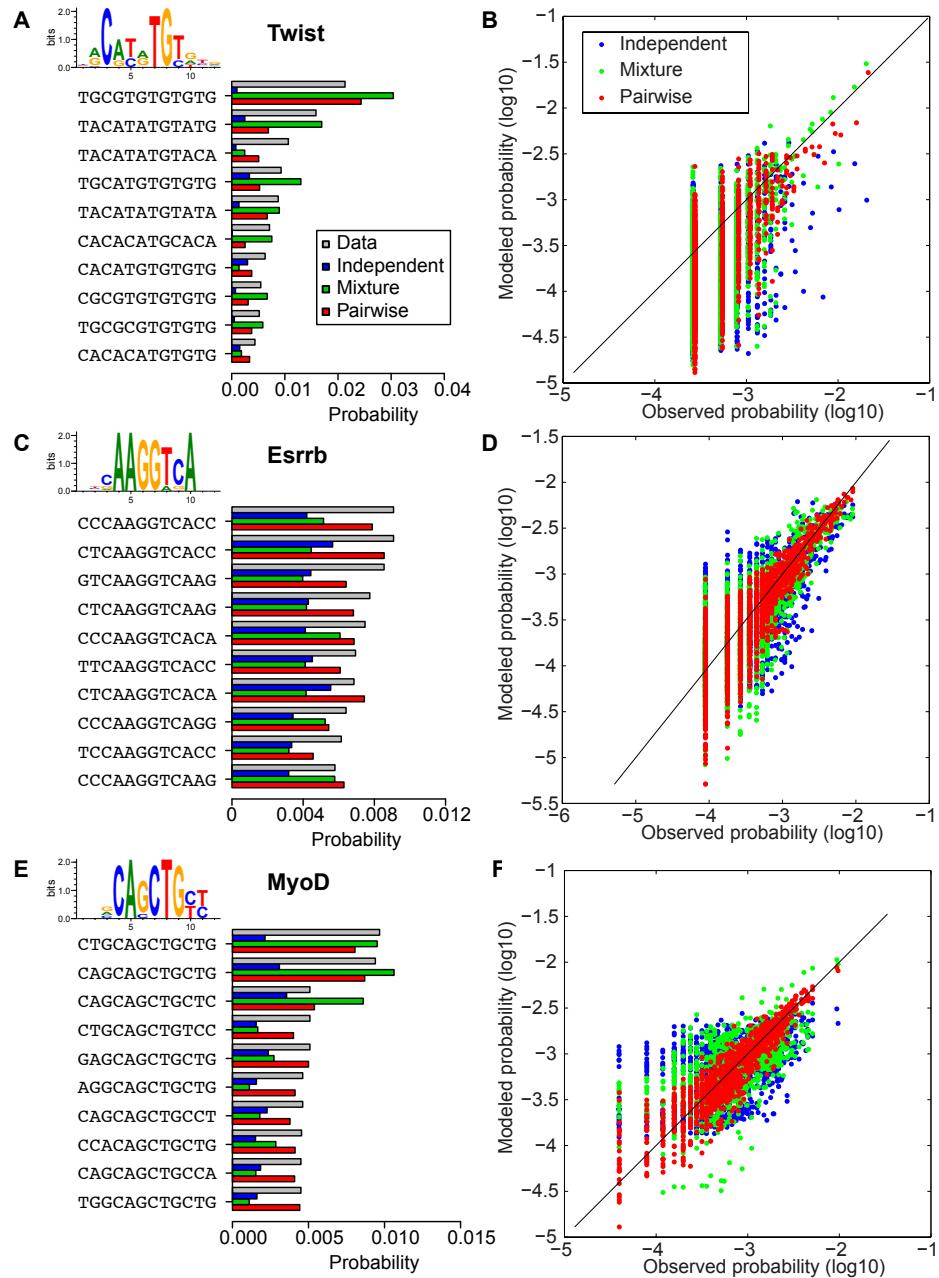


FIG. 3: Models with correlations improve TFBS statistics prediction. The observed frequencies (gray bars) of the most represented TFBSs for Twist (A), Esrrb (B) and MyoD (C) TFs, are shown together with the probabilities of these sequences predicted by the independent energy model (blue bars), the pairwise model taking into account interactions between nucleotides (red bars), and the K-means mixture model (green bars). (B,D,F) show the comparison between frequencies for all binding sequences and predicted sequence probabilities for the three models (same color code). The probability predictions of the pairwise model and to a lesser extent of the mixture model are in much better agreement with the observed frequencies than those of the PWM model.

from a few percent for TF like Esrrb, up to about 15 % for Twist. Thus, even when stemming from the same ChIPseq data, the two models can be learnt from significantly distinct set of sites.

Second, using the set of ChIPseq peaks on which the pairwise model was learned, we looked for the best pre-

dicted sites on each ChIPseq bound fragment using both the pairwise and PWM models (Figure 4B).

The overlap was found to be about 80% on average. The overlap between the sets comprising the two best TFBSs of each ChIPseq was also computed. This resulted in an overlap increase or decrease between the

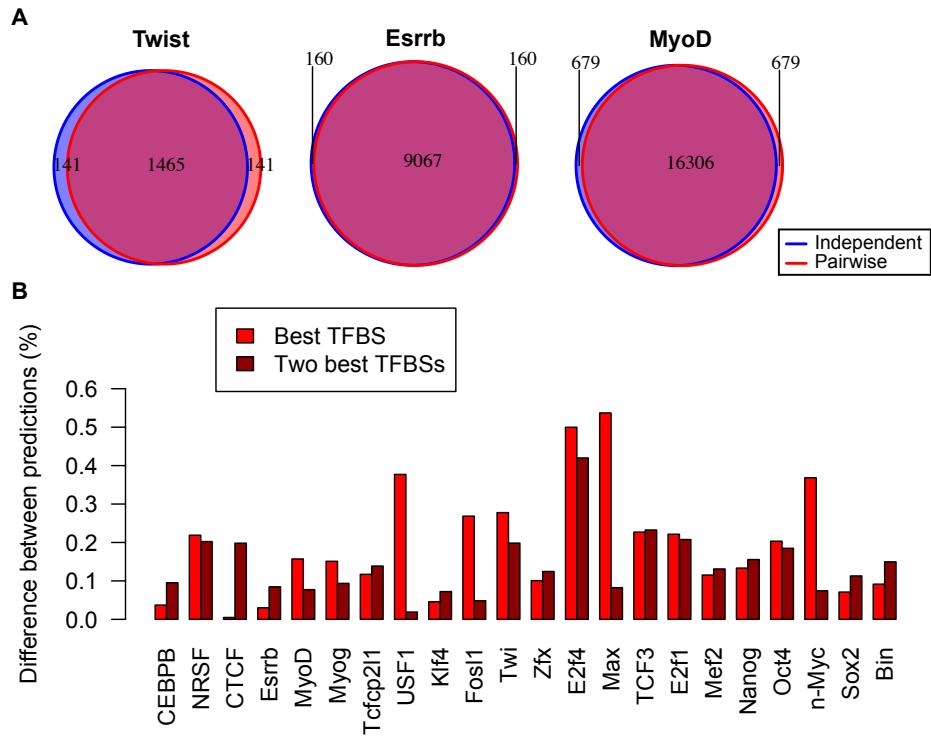


FIG. 4: Overlap between predicted sites. (A) Venn diagrams showing the overlap between the ChIP predicted by the independent (blue) and pairwise (red) models. (B) Difference (one minus the proportion of shared sites) between the best sites predicted by pairwise and PWM models on ChIPseq peaks (light red), and the same quantity when including the next best predicted sites on each peak (dark red). In several cases (*e.g.* Fosl1, Max, n-Myc, Srf, Stat3, Usf1), the difference between predicted sites is much smaller when the two best sites are considered, indicating that the pairwise model and the PWM model rank differently the two best sites in ChIP peaks with multiple bound sites.

prediction of the two models depending on the average of number of binding sites per retrieved ChIPseq fragment. In a few cases (*e.g.* CTCF, Esrrb), the inclusion of the second best TFBS increased the difference between the two models. This generally happened when the ChIPseq fragments were retrieved with typically a single TFBS above threshold (*e.g.* for Esrrb the TFBS specificity was fixed to retrieve 50% of 18453 ChIPseq and about 11000 fragments where found by the two models—see Table I). In these cases, the low specificity TFBSs tended to differ more between the two models than the very specific ones. In several other cases (*e.g.* for Fosl1, Max, n-Myc, USF1), the inclusion of the second best predicted binding sites (Figure 4B) greatly increased the overlap between the two model predictions. This corresponded to cases for which the retrieved fragments contained on average two or more TFBSs above the specificity threshold (Table I). This showed that for these cases the prediction difference between the two models arose predominantly from a different ranking of the best TFBSs.

In conclusion, the TFBS predictions made by the two models can differ significantly both in the rank of ChIPseq fragments and in the rank of binding sites on these fragments.

Comparison with a PWM-mixture model

When described by a PWM, the binding energies of a TF for different nucleotide sequences form a simple energy well with a single minimum at a preferred consensus sequence. Some authors have instead analyzed the binding specificity of transcription factors by introducing multiple preferred sequences [24, 25]. A model of this type that naturally generalizes the PWM description consists of using multiple PWMs [14]. We found it interesting to investigate this approach based on a mixture of PWMs and compare it with the pairwise interaction model to get some insights into potentially important high-order correlations that would not be captured by the pairwise model. As precisely described in *Methods*, an initial mixture of K PWMs was generated by grouping into K clusters the TFBS data for a given TF. Similarly to the pairwise interactions, the number of clusters K was constrained, to avoid over-fitting, by penalizing the corresponding model score using the BIC. For a given TF, the PWM mixture and the collection of TFBSs in the ChIPSeq data were refined iteratively until convergence, usually reached after 10 iterations. The results are shown in Figure 5A for the three representative factors,

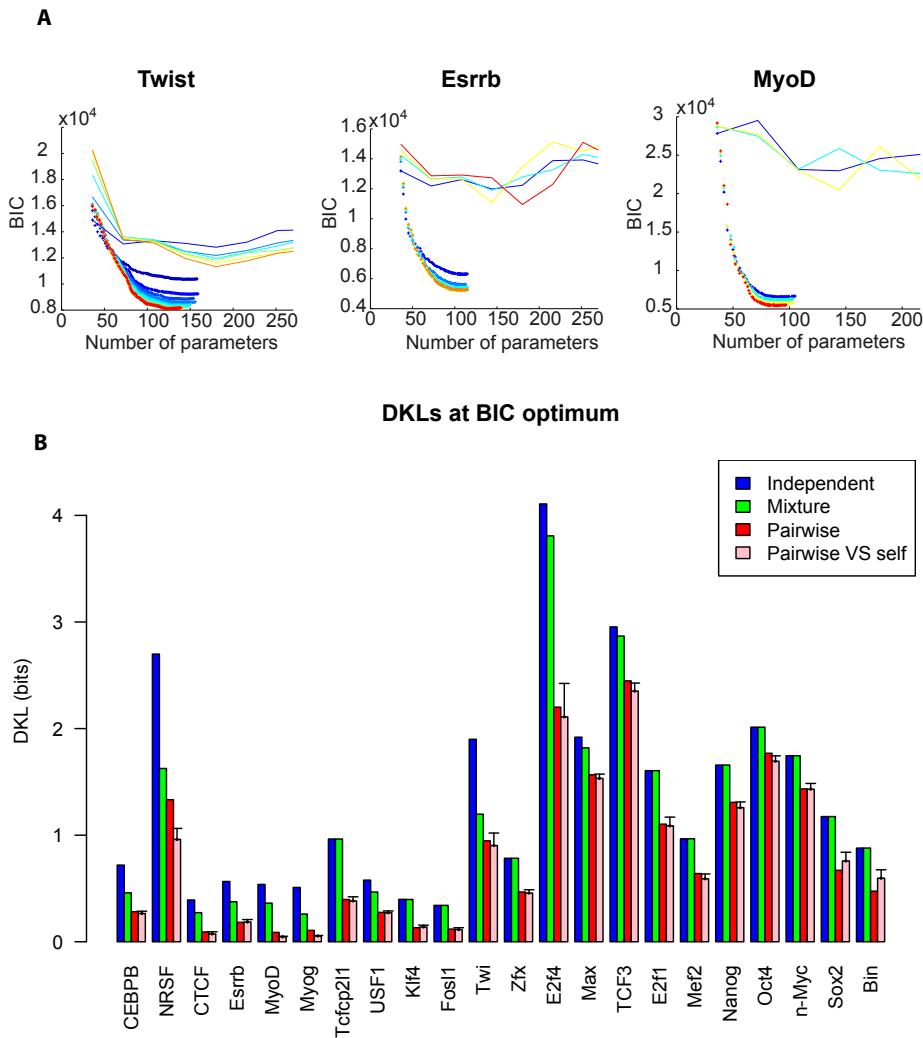


FIG. 5: Model selection. (A) Minimisation of the Bayesian information criterion (BIC, see *Methods*) is used to select the optimal number of model parameters and avoid over-fitting the training set. The evolution of the BIC is shown for the pairwise model (crosses) and the PWM-mixture model (lines). Colors from dark blue to red indicate the number of iterations (see Fig.).

(B) Kullback-Leibler divergences (DKL) between the independent, K-means and pairwise distributions and the observed distribution for the different TFs, for the BIC optimal parameters. We also show the DKL of the pairwise model with a finite-size sample of sequences drawn from it (pink, see *Methods*). Error bars represent two standard deviations.

Twi, Esrrb and MyoD.

The best description of Twi ChIPSeq data is, for instance, provided by a mixture of 5 PWMs, which corresponds to 184 independent parameters. The mixture model yields a significant improvement when compared to the single-PWM model for Twi, and milder ones for Esrrb and MyoD. In the three cases however, it proves inferior to the pairwise model.

More generally, Figure 5B shows the performances of the different models for all studied TFs using the Kullback-Leibler Divergence or DKL between the data distribution $P(s)$ and the models distributions $P_m(s)$. On the one hand, the mixture model improves the de-

scription of the binding data for 12 out of 27 TFs as compared to the single PWM model. The mixture model gives in particular strong improvements in the cases for which the binding sites have a palindromic structure (eg Twi, MyoD, Myog, Max, USF1). This feature often stems from the fact that the TF binds DNA as a dimer, which could give some concreteness to the mixture model: the recruitment of different partners by bHLH factors like MyoD or Myog could indeed lead to a mixture of TFs binding the same sites. On the other hand, the pairwise model clearly outperforms the other models in all cases studied.

As in the PWM case, the finite size of the datasets

leads us to expect fluctuations in the estimation of the DKL. In order to assess the magnitude of these finite-size fluctuations, we computed the average DKL between the best-fitting (pairwise) model and a finite-size artificial sample drawn from its own distribution, as shown in Figure 5B. Values of this DKL that are larger than the one obtained with the real dataset are indicative of overfitting, while the opposite case would suggest that the model is incomplete. In all cases, however, the DKL obtained with this control procedure was within error bars of the value computed with respect to the observed sample, with the exception of NRSF, MyoD, and Myog, as seen in Figure 5B. Thus, the pairwise model is generally the best possible model, insofar as the available dataset allows us to probe.

The metastable states of the pairwise interaction model

In order to more directly relate the pairwise interaction and the mixture models, it is useful to consider the energy landscape of the pairwise interaction model in the space of all possible TFBSSs. By contrast with the simple, single-minimum energy well of the PWM model, the pairwise interaction model has multiple metastable energy minima. The energy landscape of the pairwise interaction model can thus be seen as a collection of energy wells, each centered on its metastable energy minimum. The span of the different energy wells in sequence space can be precisely defined as the basins of attraction of the different metastable minima in an energy minimizing procedure (see *Methods*). This allows one to associate each observed TFBS to a particular energy minimum. This defines basins of attraction that are used to build representative PWMs for each metastable minimum together with a weight—the number of sequences in the basin of attraction—for this energy minimum. We compared each metastable minimum to the PWMs of the mixture model, by calculating the DKL between the PWM computed from the sequences in its basin of attraction and the PWMs of the mixture model. This gave an effective distance which allowed us to associate each metastable state to the nearest PWM of the mixture model.

Using this procedure, we computed the set of PWMs and weights corresponding to the 27 considered TF pairwise interaction models. Examples are shown in Figure 6. In the case of Twi, the PWMs of the pairwise model (“metastable PWMs”) can be clearly associated to the $K = 5$ PWMs of the mixture model. For MyoD, three of the 5 “metastable PWMs” can be clearly assigned to PWMs of the mixture model. The other two have a more spread out representation. The case of Esrrb is similar with one “metastable PWM” in clear correspondence with one PWM of the mixture model, and the other one less clearly so. The correspondence between the two models is shown in Figure S2 for the other TFs for which the mixture model uses more than a single PWM. This

representation allows one to identify some features captured by the pairwise model. For example, in the case of Twist, most of the correlations are coming from the two nucleotides at the center of the motif, which take mainly 3 values among the 16 possible: CA,TG and TA. In the case of MyoD, the representation makes apparent the interdependencies between the two nucleotides following the core E-Box motif, and the restriction to the three main cases of CT, TC and TT.

Properties of the pairwise interactions

The computation of the interaction parameters allows an analysis of some of their properties. In particular, it is interesting to quantify their strengths and measure the typical distance between interacting nucleotides. We address these two questions in turn.

The concept of Direct Information was previously introduced to predict contacts between residues from large-scale correlation data in protein families [33]. We used it here to measure the strength of the pairwise interaction between two nucleotides. Using the previously generated interaction parameters from the pairwise model, we built the Normalized Direct Information (NDI), a quantity which varies from 0 for non-existing interactions, to 1 when interactions are so strong that knowing the amino acid identity at one position entirely determines the amino acid identity at the other position (see *Methods*). Heatmaps displaying the results for the representative Twist, Esrrb and MyoD factors are shown in Figure 7 and in Figures S3 for the other factors. An important observation is that the direct information between different nucleotides is rather weak—usually smaller than 10%—but substantially larger than the direct interaction between nucleotides in the surrounding background (1-3%, see Figure S4). It is interesting to note that such weak pairwise interactions give rise to a substantial improvement in the description of TFBS statistics, similarly to what was previously found in a different context [16]. The pairwise interactions are furthermore observed in Figure 7 to be concentrated on a small subset of all possible interactions. This can be made quantitative by computing the Participation Ratio of the interaction weights, an indicator of the fraction of pairwise interactions that accounts for the observed Direct Information (see *Methods*). Here, typical values of 10 – 20% were found (Figure 7 and Table I), showing that the interactions tend to be concentrated on a few nucleotide pairs.

The interaction weights can also be used to measure the typical distance between interacting nucleotides. To that purpose, we computed the relative weight of the Direct Information as a function of the distance between nucleotides (see *Methods*). Figure 8 A shows box plots that summarize the results for the considered Drosophila and mammalian TFs. Both plots show a clear bias towards nearest-neighbor interactions with a strong peak at $d = 1$, and a rapid decrease for $d \geq 2$. Finally, the

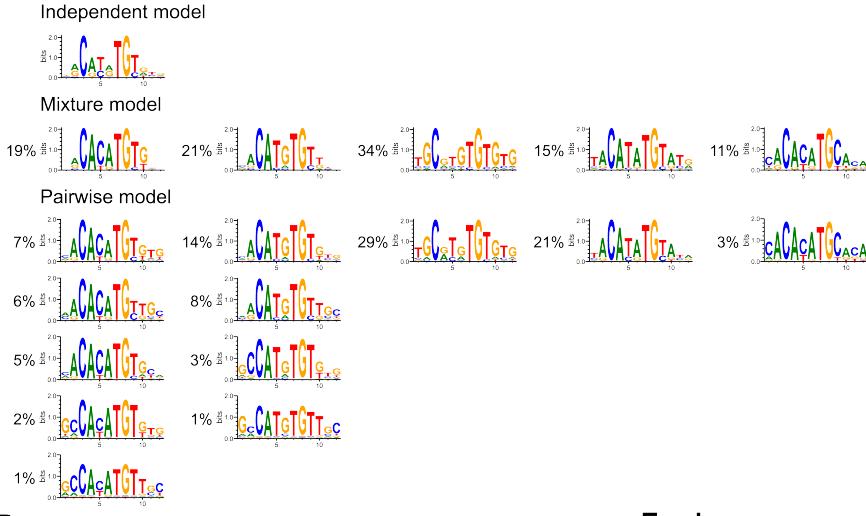
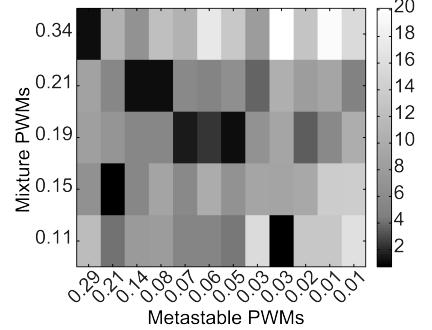
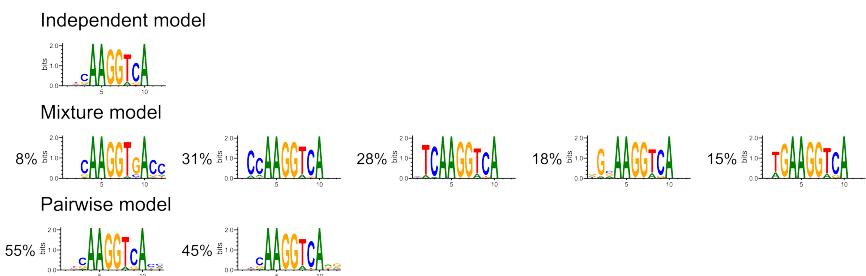
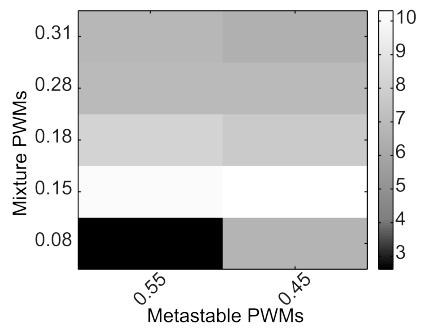
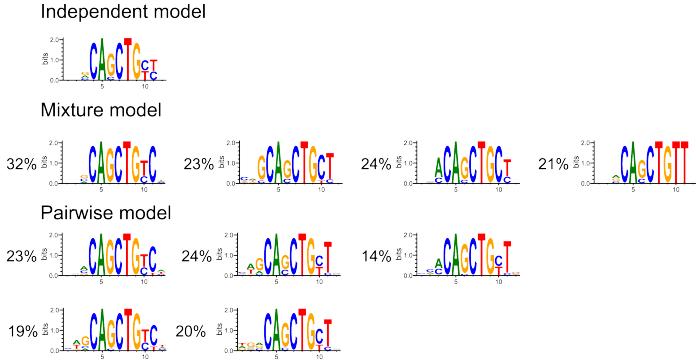
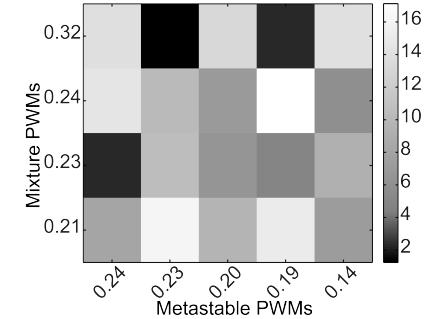
A**Twist****B****Esrrb****C****MyoD**

FIG. 6: Metastable states. The DNA sequence variety described by each model is illustrated using weblogs [32]. Shown are PWMs built from all sites, from the PWM-mixture model, and from the basins of attraction of the pairwise interaction model for Twist (A), Esrrb (B), and MyoD (C). The metastable PWMs are grouped under the mixture PWMs with smallest distance (measured by DKL, in bits). Heatmaps showing the DKLs between metastable PWMs and mixture PWMs are displayed on the right for each factor (minimal DKLs are in black). The proportions of sites used for each logo are also indicated and serve to denote the corresponding PWM.

dominant pair interactions are on average located in the flanking regions of the BS in clear anti-correlation with the most informative nucleotides which are on average in the central region (Figure 8 B). These observations for TF binding *in vivo* agree with similar ones made from a large recent analysis of TF binding *in vitro* [9]. The

fact that for pair correlations to be important, nucleotide variation at a given location is required, may be one way to rationalize them.

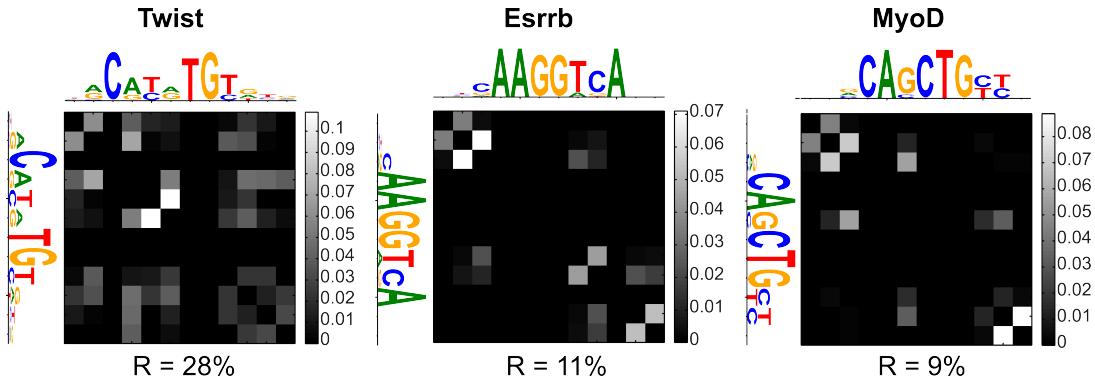


FIG. 7: **Nucleotide pair interactions.** Heat maps showing the values of the Normalized Direct Information between pairs of nucleotides. The matrix is symmetric by definition. PWMs are shown on the side for better visualization of the interacting nucleotides. The participation ratio R is indicated below each heat map.

TABLE I: **Participation Ratios**

Name	Part. Ratio
Bin	0.11
Mef2	0.19
Twi	0.28
E2f1	0.13
Esrrb	0.11
Klf4	0.16
Nanog	0.10
n-Myc	0.09
Oct4	0.24
Sox2	0.12
Tcfcp2l1	0.12
Zfx	0.10
CEPB	0.05
CTCF	0.23
E2f4	0.14
Fosl1	0.09
Max	0.18
MyoD	0.09
Myog	0.09
NRSF	0.27
TCF3	0.19
USF1	0.07

Alternative representation of interactions by Hopfield patterns

Using a simple binary description of neurons, JJ Hopfield suggested, in a classic piece of work [34], that neural memories could be attractors corresponding to patterns arising from pair interactions between neurons. These interaction patterns can be computed in the present case. They offer an alternative way to analyze the patterns of

correlation from the pair-interactions between positions, as already proposed in a mean-field context in [35]. Because the matrix of interactions J_{ij} is symmetric, it can be diagonalized in an orthonormal basis of eigenvectors ξ^k , the Hopfield patterns in the present case, with corresponding real eigenvalues λ_k . These orthonormal eigenvectors correspond to the Hopfield patterns in the present case. The Potts energy (Eq. (1)) for a binding sequence $s_1 \cdots s_L$ can be rewritten in terms of the Hopfield patterns as (see Methods):

$$\mathcal{H} = - \sum_i h_i(s_i) - \frac{1}{2} \sum_{k=1}^{4L} \lambda_k \left(\sum_{i=1}^L \xi_i^k(s_i) \right)^2. \quad (2)$$

Although here the presence of the diagonal h term prevents the patterns to be metastable energy states, they can still be useful to analyze the interaction matrix. This spectral decomposition of the interaction matrix is also similar in spirit to a principal component analysis, and even equivalent in the case of Gaussian variable. One can thus wonder how many patterns are needed to well approximate the full matrix of interactions J . To address this question, one can rank the eigenvalues λ_k in order of decreasing moduli and note J_p the restriction of the interaction matrix generated by the first p eigenvalues and their associated patterns. The full interaction matrix naturally corresponds to J_{48} . Approximate interaction matrices obtained by keeping different numbers of dominant patterns are shown in Figure 9 for the three considered representative factors. Pairs of successive patterns appear to provide the main interaction domains in this representation, as is particularly clear in the case of MyoD. One can see in Figure 9 that J_6 already closely approximates the full interaction matrix, a reflection in the present representation, that the important interactions are concentrated on a few links between pairs of nucleotides.

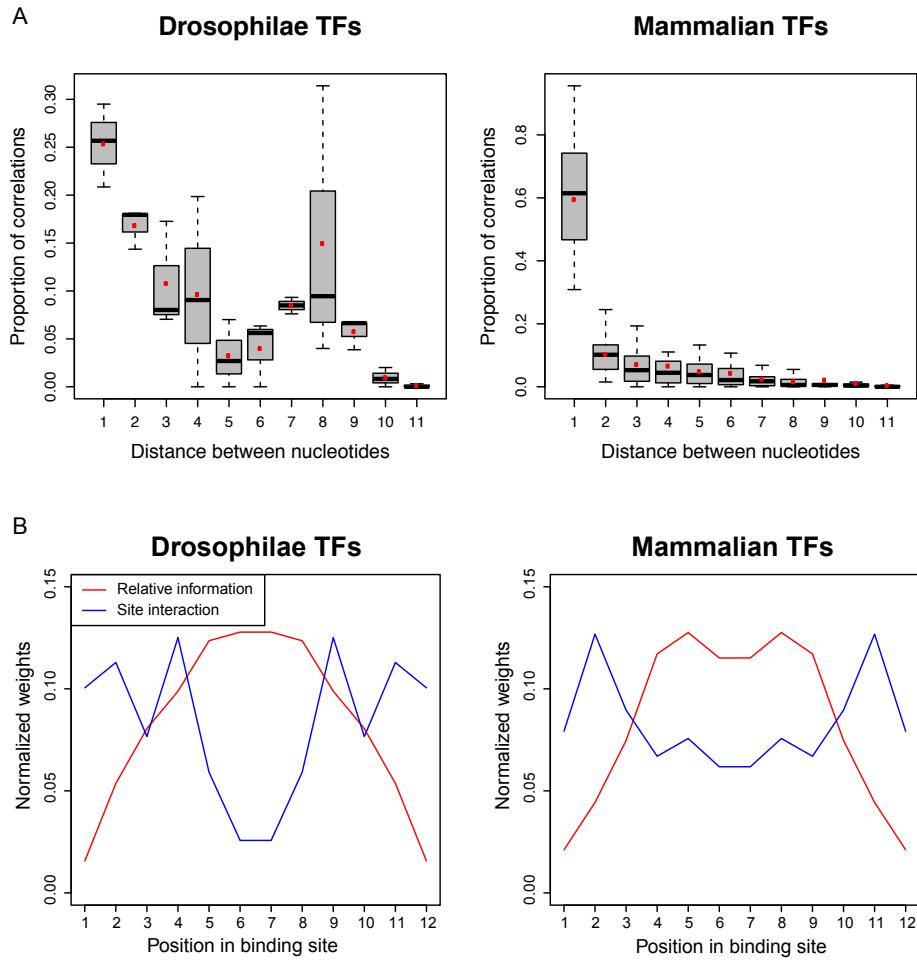


FIG. 8: Properties of the pair interactions. (A) Distances between interacting nucleotides. The box plots show the relative importance of the Normalized Direct Information as a function of the distance between interacting nucleotides. Red dots denote average values. (B) Sum of normalized direct informations in the TFBSSs at a given position, averaged over all considered factors (blue line). The average site information content relative to background as a function of position is also shown (red line). In both quantities, the average over the two TFBSS orientations has been taken.

Discussion

The availability of ChIPseq data for many TFs has led us to revisit the question of nucleotide correlations in TFBSSs. In order to perform a fully consistent analysis of this type of data, we have developed a workflow in which the TFBSS collection and the model describing them are simultaneously obtained and refined together. We have found that when a sufficiently large number of TFBSSs is available, the PWM description does not account well for their statistics. The general presence of correlations that follows from this finding, agrees with previous reports for particular transcription factors [6, 8, 24] and with the conclusions of large scale *in vitro* TF binding studies [9, 10].

In order to refine the PWM description, we have analyzed a model with pairwise interactions [23], and a PWM mixture model [14]. Data overfitting is a concern

for multi-parameter models and has been addressed by putting a penalty on the parameter number using the BIC. While the mixture-model improved in some cases the PWM description, especially for palindromic binding sites, a much more significant and general improvement was found with the pairwise interaction model. The success of the pairwise interaction model agrees with the results of its recent application (however, without the BIC) to high-throughput *in vitro* binding data [23]. It moreover shows that, at least in the case we considered, pairwise interactions are sufficient to account for higher-order correlations, and that an explicit description like the one provided by the PWM-mixture model is not necessary. For example, for Essrb, metastable states arising from nearest-neighbor interactions reproduce a triplet of flanking nucleotides with a variable spacer from the core motif (Figure S5).

Our detailed analysis of the obtained interaction mod-

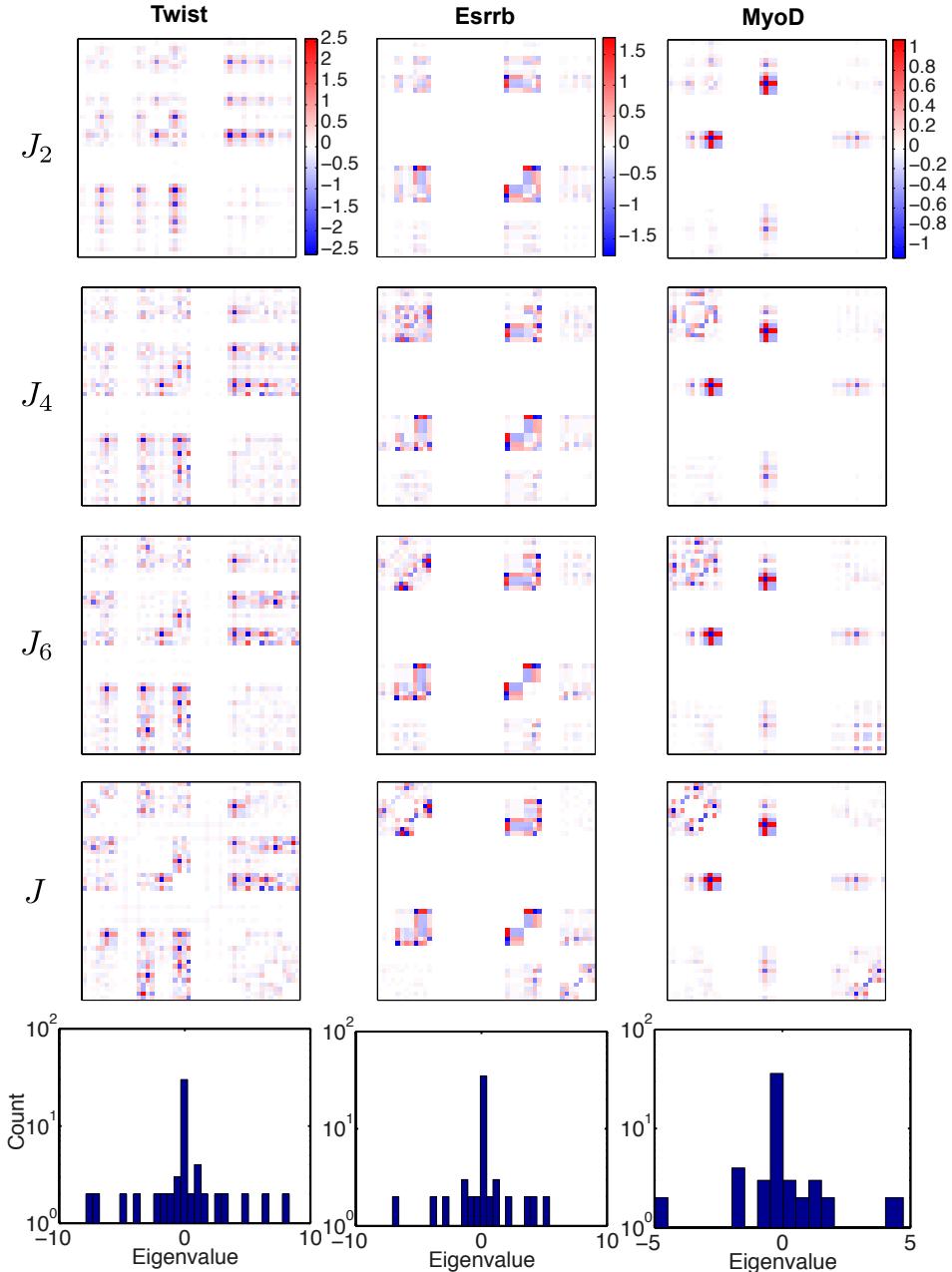


FIG. 9: Representation of interactions by Hopfield patterns. The full interaction matrix J is approximated by a matrix J_p built from the p Hopfield patterns with highest eigenvalue moduli. We show J_2 , J_4 , J_6 and the full matrix J in the basis (i, b) with $i = \{1, \dots, 12\}$ and $b = \{\text{A}, \text{C}, \text{G}, \text{T}\}$. Color bars are shown on the first row for each factor. For MyoD, the correspondence between successive pairs of patterns and distinct interaction domains is seen particularly clearly. In all cases the full interaction matrix is already well approximated by J_6 .

els for different TFs shows that the weights of pairwise interactions are generally weak. The most important are only about 10 % of the PWM weights, but significantly above the interaction weights in the surrounding background DNA (of the order 1-3% by the same measure). Nonetheless, collectively these interactions significantly improve the model description of the TF binding data as found in other examples [16].

We have here obtained the pairwise interaction models based on the principle of maximum entropy, constrained to account for the pair-correlations measured in the data. This approach has already been followed in a variety of biological contexts, from populations of spiking neurons [16, 17] to protein sequences [20] to bird flocks [22]. An interesting feature of these interaction models is their non-convexity, which allows for the existence of many lo-

cal maxima in the probability distribution of sequences, or local minima of energy. This was noted for repertoires of antibodies in a single individual [21], where many of these local states were observed and suggested as possible signatures of past infections. In a very different context, local probability maxima in the probability distribution of retinal spiking patterns was reported and linked to error-correcting properties of the visual system [36]. In the present case of TFBSs, these local minima reflect the multiplicity of binding solutions and resemble the individual PWMs of the mixture model. Pairwise interaction models thus somehow incorporate models of multiple PWMs while outperforming them.

The previously considered case of protein sequences shares many similarities to the statistics of TFBSs, since correlations in protein sequences as in TFBSs reflect both structural and functional constraints. In proteins families, correlations are usually interpreted as resulting from the co-evolution of residues interacting with each other in the protein structure. These effects are hard to distinguish from phylogenetic correlations or other observational biases. Nonetheless, the inference of interaction models from data was successfully used to predict physical contacts between amino-acids in the tertiary structure [37], and to aid molecular dynamics simulations in predicting protein structure [38–40]. In the case of TFBSs, comparison between *in vitro*[9, 10] and *in vivo* binding data may help to disentangle the different possible origins of the found correlations and seems worth pursuing. It appears similarly interesting to study how much of the found pair correlations can be explained on the basis of structural data. Finally, the role of nucleotide interaction in TFBS evolution [41] should be considered and could improve the reconstruction of TFBSs from multi-species comparison [42–44].

Independently of these future prospects, we have found that the TFBSs predicted from ChIP-seq data significantly depended on the model used to extract them. Since the pairwise interaction model and the developed workflow significantly improve TFBS description and require a modest computational effort, they should prove worthy tools in future data analyses.

Materials and Methods

Genome-wide data retrieval

We use both ChIP-on-chip data from *Drosophila Melanogaster* and ChIPseq data from *Mus Musculus*. Data was retrieved from the litterature [26, 27] and from ENCODE data accessible through the UCSC website <http://hgdownload.cse.ucsc.edu/goldenPath/mm9/encodeDCC/wgEncodeCaltechTfbs/>, for a total of 27 TFs. Among them, there are 5 developmental Drosophilae TFs: Bap, Bin, Mef2, Tin and Twi, 11 mammalian stem cells TFs: c-Myc, E2f1, Esrrb, Klf4, Nanog, n-Myc, Oct4, Sox2, Stat3, Tcfcp2l1, Zfx, and 11 factors

involved in mammalian myogenesis: Cebpb, E2f4, Fosl1, Max, MyoD, Myog, Nrsf, Smad1, Srf, Tcf3, Usf1. Overall, there are between 678 and 38292 ChIP peaks, with average size 280bp. DNA sequences were masked for repeats using RepeatMasker [45].

Background models

It is important to discriminate the statistics of the motifs proper from that of the background DNA on which motifs are found. Besides particular nucleotides frequencies, the background DNA can exhibit significant nucleotide correlations, for instance arising from CpG depletion in mammalian genomes (Figure S4). For each ChIPseq data, we used, as background, all sites from both strands of the sequences. This serves to learn independent and pairwise background models which were used as reference models to score the corresponding TFBS models. The position information content in all plotted PWM logos is measured with respect to the nucleotide background frequencies (*i.e.* the independent background model)

Initial PWM refinement

Along with the ChIPseq data for the different factors, we also retrieved corresponding PWMs from the literature [26] or from TRANSFAC database [46]. These initial PWMs were refined according to the following protocol.

Given ChIPseq data (bound regions) for a given TF and an initial PWM of length L ($L = 12$ was taken for all computations in the present paper), we scanned both strands of each bound region and attributed to all observed L -mers a score defined as the ratio between the PWM and background models probabilities. A cutoff was set such that half of the bound regions had at least one predicted TFBS with a score above the cutoff, setting a True Positive Rate (TPR) of 50%. This heuristic criterium overcame the problem of False Positives among the ChIPseq peaks that might have polluted the data. This defined a training set of N L -mers with probability higher than the cutoff. Bound sites were again predicted using the same cutoff. This procedure was repeated until stabilization of the predicted sites to a fixed subset. This resulted in a refined PWM with its associated set of bound sites.

Independent model evaluation

The independent model consist of a matrix of single nucleotide probabilities of size $4 \times L$, where L is the width of the binding site. In a first approximation, the parameters appearing in the matrix can be estimated from a set of binding sites by computing the observed frequency $f_{b,i}$ of

TABLE II: Information about the TFs

Name	N_{chip}	$\Delta_{\text{inde-mixture}}$	$\Delta_{\text{inde-pairwise}}$	$\Delta_{\text{mixture-pairwise}}$	N_{inde}	N_{mixture}	N_{pairwise}
Bap	678	0	12	12	2205	2208	2117
Bin	1857	2	80	81	1300	1298	1228
Mef2	4545	0	161	161	3681	3681	3665
Tin	1791	0	40	40	1333	1333	1310
Twi	3211	182	141	128	3810	3862	3722
c-Myc	3038	0	95	95	2996	2996	2920
E2f1	17367	0	877	877	16625	16625	14915
Esrrb	18453	172	160	167	11243	11333	11275
Klf4	9404	0	97	97	5912	5912	5913
Nanog	8022	0	111	111	6196	6196	6224
n-Myc	6367	0	54	54	6981	6981	6954
Oct4	3147	0	74	74	3187	3187	3079
Smad1	907	0	24	24	690	690	667
Sox2	3523	0	95	95	2306	2306	2293
STAT3	2099	54	58	62	2308	2264	2231
Tcfcp2l1	22406	0	418	418	16691	16691	16649
Zfx	9152	0	203	203	6473	6473	6473
CEBPP	14500	399	337	334	8275	8322	8267
CTCF	32958	360	492	579	17087	17098	17060
E2f4	4132	248	590	517	4643	5146	3879
Fosl1	5981	0	90	90	5088	5088	5039
Max	8751	24	70	81	12531	12495	12386
MyoD	33969	717	679	665	25416	25430	25344
Myog	38292	1116	584	835	29520	29334	29647
NRSF	13756	639	672	488	13183	14363	13440
SRF	2370	1	34	35	2929	2928	2948
TCF3	9453	185	277	257	8528	8690	8775
USF1	8956	11	14	12	8628	8619	8625

For each TF, we show the number N_{chip} of ChIP sequences retrieved, the numbers $\Delta_{\text{inde-pairwise}}$, $\Delta_{\text{inde-mixture}}$, and $\Delta_{\text{pairwise-mixture}}$ of different ChIP sequences used for training between either two models, and the numbers N_{inde} , N_{mixture} , and N_{pairwise} of TFBSSs used to learn each model.

nucleotide b at position i . However, this frequency fluctuates around the “true” probability due to finite sample size, and for example unobserved nucleotides could actually have a low probability of being observed provided that the number of observations be high enough. It is usual to correct for this effect by using the Bayesian pseudo-count approach stemming from Laplace’s rule of succession [3]. The probability to observe nucleotide b at position i is given by:

$$p_{i,b} = \frac{n_{i,b} + \alpha_b}{N + \sum_b \alpha_b} \quad (3)$$

where $n_{i,b}$ is the number of observed b at position i , N is the total number of sites, and α_b ’s are the pseudo-counts, or prior probabilities to observe nucleotide b at position i . The pseudo-counts were all set to 1, however no significant effect was noted when changing this value,

as expected from the large number of observations.

Kullback-Leibler divergence

The Kullback-Leibler divergence is a measure of distance between two probability distributions p and q of a variable s , and is defined as:

$$\text{DKL}(p\|q) = \sum_s p(s) \log \frac{p(s)}{q(s)}. \quad (4)$$

Throughout this paper, when a DKL is calculated between a finite sample and a model distribution, p corresponds to the sites frequencies in the sample, and q to the model distribution. When the DKL is calculated between a PWM of a basin of attraction of a metastable

state and a PWM from the mixture model, p is used for the former, and q for the latter.

Estimation of the fluctuations due to finite sampling: DKL vs self

To estimate whether the description of the data by a model (*e.g.* independent or pairwise) could be improved or was consistent with the finite number N of observed sequences, we computed the ‘self’ DKL between the distribution of a set of N sequences drawn from the model distribution and the model distribution itself. This procedure was repeated 100 times. TFs for which the independent model DKL was smaller than or within two standard deviations of the self DKL were discarded for later analysis.

Derivation of the pairwise interaction model

Information theory offers a principled way to determine the probabilities of a set of states given some mea-

surable constraints. It consists in maximizing a functional known as the entropy[47, 48] over the set of possible probability distributions given the imposed constraints. Here, we wish to determine the probability $P(s)$ of a DNA sequence s of length L , in the set of TFBSs for a transcription factor, given the constraints that the probability distribution P retrieves the one- and two-point correlations observed in a set of bound DNA sequences. We denote by \mathcal{A} the alphabet of possible nucleotides, $\mathcal{A} = \{A, C, G, T\}$ and by s_i the nucleotide at position i in the sequence s so that $s = s_1 \dots s_L$. With these notations, the entropy with the considered constraints translates into the the following functional:

$$\begin{aligned} \mathcal{L} = & - \sum_{\{s\}} P(s) \ln P(s) + \lambda \left(\sum_{\{s\}} P(s) - 1 \right) + \sum_{i=1}^L \sum_{a \in \mathcal{A}} h_i(a) \left(\sum_{\{s\}} P(s) \delta(s_i, a) - P_i(a) \right) \\ & + \sum_{i=1}^{L-1} \sum_{j>i} \sum_{a \in \mathcal{A}} \sum_{a' \in \mathcal{A}} J_{i,j}(a, a') \left(\sum_{\{\sigma\}} P(a) \delta(s_i, a) \delta(s_j, a') - P_{i,j}(a, a') \right), \end{aligned} \quad (5)$$

where $P_i(a)$ (resp. $P_{i,j}(a, a')$) is the probability of having nucleotide a at position i (resp. nucleotides a and a' at position i and j) in the TFBS data set. The function δ denotes the Kronecker δ -function defined by $\delta(a, a') = 1$ if $a = a'$, and 0 otherwise. The first term in Eq. (5) is the entropy of the probability distribution to be found and the other terms are the given constraints along with their Lagrangian multipliers. Maximization of the functional \mathcal{L} is performed in a usual way by setting the functional derivative with respect to the probability distribution P to zero:

$$\frac{\delta \mathcal{L}}{\delta P(s)} = 0 = -\ln P(s) - 1 + \lambda + \sum_{i=1}^L h_i(s_i) + \sum_{i=1}^{L-1} \sum_{j>i} J_{i,j}(s_i, s_j). \quad (6)$$

Finally, using the constraint $\sum_{\{s\}} P(s) = 1$, one finds the probability distribution that maximizes entropy given the constraints that it reproduces the observed one- and two-point correlations:

$$P[s] = \exp[-\mathcal{H}(s)]/\mathcal{Z}, \quad (7)$$

where $\mathcal{H}(s)$ is the inhomogeneous Potts model Hamiltonian,

$$\mathcal{H}[s_1 \dots s_L] = - \sum_{i=1}^L h_i(s_i) - \sum_{i=1}^L \sum_{j<i} J_{i,j}(s_i, s_j), \quad (8)$$

$$s_i \in \{A, C, G, T\}.$$

The normalization constant \mathcal{Z} is the partition function,

$$\mathcal{Z} = \sum_{\{s\}} \exp[-\mathcal{H}(s)]. \quad (9)$$

Gauge fixing

The probability distribution of sequences, as given by Eqs. (7, 8), is invariant under shifts of the local fields $h_i(a)$ and under transformations between the interaction terms $J_{i,j}(a, a')$ and the local fields. In order to uniquely determine \mathcal{H} , this arbitrariness needs to be taken care of by adding further conditions that uniquely fix the in-

teraction parameters, a process known as gauge fixation [20] that we detail here.

a. *Local fields.* Since it amounts to changing the reference energy and is cancelled by the normalization, the probability is invariant with respect to the following global shift of the $h_i(a)$

$$h_i(s_i) \rightarrow \tilde{h}_i(s_i) = h_i(s_i) + \varepsilon_i. \quad (10)$$

We choose to fix this invariance by minimizing the square norm $S_i = \sum_{a \in \mathcal{A}} \tilde{h}_i(a)^2$ of local field terms with respect to the gauge degree of freedom. The corresponding gauge-fixing condition reads

$$\sum_{a \in \mathcal{A}} \tilde{h}_i(a) = 0. \quad (11)$$

This condition can be imposed on any set of fields h_i by using the symmetry (10) and redefining the fields as follows,

$$h_i(s_i) \rightarrow h_i(s_i) - \frac{1}{4} \sum_{a \in \mathcal{A}} h_i(a). \quad (12)$$

b. *Interaction terms.* Another invariance stems from the fact that contributions can be shifted between local fields and interaction energies. Namely, the following change of variables does not affect the probability:

$$J_{ij}(s_i, s_j) \rightarrow \tilde{J}_{ij}(s_i, s_j) = J_{ij}(s_i, s_j) + \psi_i(s_i) + \phi_j(s_j) + C_{i,j}, \quad (13)$$

since the local fields ψ_i and ϕ_j can be redistributed in h and the constant $C_{i,j}$ gives an energy reference for the interacting nucleotides that is cancelled by the normalization process. Again, a gauge condition is obtained by minimizing the square norm $S_{i,j} = \sum_{a,a' \in \mathcal{A}} [\tilde{J}_{ij}(a, a')]^2$ of interaction terms with respect to the gauge degrees of freedom. This yields the conditions:

$$\sum_{a \in \mathcal{A}} \tilde{J}_{i,j}(a, a') = \sum_{a' \in \mathcal{A}} \tilde{J}_{i,j}(a, a') = 0. \quad (14)$$

These can be imposed on a set a of J_{ij} parameters by redefining them as follows:

$$\begin{aligned} J_{ij}(s_i, s_j) &\rightarrow J_{ij}(s_i, s_j) + \frac{1}{16} \sum_{a, a' \in \mathcal{A}} J_{i,j}(a, a') \\ &\quad - \frac{1}{4} \sum_{a \in \mathcal{A}} J_{i,j}(a, s_j) - \frac{1}{4} \sum_{a \in \mathcal{A}} J_{i,j}(s_i, a). \end{aligned} \quad (15)$$

Determination of the pairwise interaction model from the data

The parameters of the inhomogeneous Potts model in Eq. (8), giving the energy of an observed sequence of length L , must be computed from the data. The parameters h and J represent the energy contributions respectively coming from individual nucleotides and from

their interactions. The PWM model is the particular case where all the interaction parameters vanish: $J_{i,j}(a, a') = 0$.

To build the model, we start from the PWM description, characterized by the set of initial $h_i(a) = \log p_{i,a}$ and the interaction parameters J 's set to zero. We add one interaction parameter $J_{i,j}(a, a')$ at a time, corresponding to the pair of nucleotides whose pairwise distribution predicted by the model differs most from data, as estimated by a binomial p -value. We then fit the augmented model to data, use this model to select a new set binding sites from the reads, and repeat the whole procedure. In each of these steps, fitting is performed by a gradient descent algorithm:

$$J \rightarrow J + \epsilon [c_2^{\text{data}} - c_2^{\text{model}}], \quad (16)$$

$$h \rightarrow h + \epsilon [c_1^{\text{data}} - c_1^{\text{model}}], \quad (17)$$

where c_1 and c_2 are matrices of size $4 \times L$ and $4L \times 4L$ respectively corresponding to the single- and two-point frequencies, and superscripts denote whether the matrices are computed from the data or from the model distribution. This algorithm converges to the set of parameters $(\{\tilde{h}_i\}, \tilde{J}_{i,j})$ that match all single marginals and the pairwise marginals of interest. The number of interaction parameters that are being added is controlled by the Bayesian Information Criterion, or BIC (Figure 5). The BIC computes the opposite log-likelihood and adds a penalty proportional to the number of parameters involved. This advertises the over-fitting of a finite dataset with an extravagant number of parameters. The procedure is iterated until minimization of the BIC, yielding the best pairwise model with the full set of parameters $(\{h_i(a)\}, \{J_{i,j}(a, a')\})$. As in the case of the PWM model, we score each sequence using the ratio between the TF and background pairwise models and impose a score cutoff so as to select a set of bound sites yielding 50% TPR, on which a new pairwise model is learned. This process is iterated until convergence to a stable set of bound sites.

BIC computation

Consider a sample $X = (X_1, \dots, X_N)$ of N TFBSs drawn from an unknown distribution function f we wish to estimate. To this extent, several models $\{M_1, \dots, M_m\}$ are proposed, each model M_i having a density g_{M_i} with parameter θ_i of dimension K_i . It is straightforward to see that, as K_i increases, the fit to the observed sample as measured by the likelihood function $g_{M_i}(X|\theta_i)$ increases as well, the limiting case being when f is estimated as the sample distribution. However, such an estimator is inappropriate to account for new, yet unobserved TFBSSs, *i.e.* it is not predictive. Such a case where the number of parameters used to estimate a distribution becomes of the order of the size of the sample is known as overfitting. The BIC allows to overcome overfitting by penalizing high dimension parameters. Using

Bayes Rule, and a uniform a priori distribution on the models, we have

$$P(M_i|X) \propto P(X|M_i). \quad (18)$$

That is, the probability of the model given the data can be inferred from the probability that the data is generated by the model. The latter is obtained by marginalizing the joint distribution of the data and the parameters over the space of parameters Θ :

$$P(X|M_i) = \int_{\Theta} P(X, \theta|M_i) d\theta = \int_{\Theta} g_{M_i}(X|\theta) P(\theta|M_i) d\theta. \quad (19)$$

For a unidimensional parameter θ , the likelihood $g_{M_i}(X|\theta)$ is maximized at some particular $\hat{\theta}_i$ with an uncertainty (or width) proportional to $1/\sqrt{N}$ in the limit of large N . Assuming a broad prior, then for large N the integral is dominated by the likelihood which is concentrated around its maximum. One can then approximate the integral by the area of the region of height the maximum likelihood and of width $1/\sqrt{N}$, that is $g_{M_i}(X, \hat{\theta}_i)/\sqrt{N}$. This result can be retrieved analytically using the method of steepest descent. For a number K_i of parameters, one gets a total volume $g_{M_i}(X, \hat{\theta}_i)/N^{K_i/2}$ [31]. Taking the logarithm yields the BIC condition:

$$BIC_i = -2 \log(P(X|M_i)) \simeq -2 \log(g_{M_i}(X, \hat{\theta}_i)) + K_i \log(N). \quad (20)$$

In the present case, the sample X is the set of observed TFBSSs and the model M_i determines the probability $P_{M_i}(s)$ of belonging to X ,

$$\log(g_{M_i}(X, \hat{\theta}_i)) = \sum_{s \in X} \log[P_{M_i(\hat{\theta}_i)}((s))]. \quad (21)$$

The interpretation of Eq. (20) is clear: adding new parameters improves the fit, but also adds new sources of uncertainty about these parameters due to the finite size of the data. This uncertainty disappears as $N \rightarrow \infty$, since the log-likelihood scales with N while the correction scales with $\log(N)$.

Finally, Eq. (20) is a functional over models, the chosen model M_{BIC} is the one that minimizes it,

$$M_{BIC} = \operatorname{argmin}_{M_i} BIC_i. \quad (22)$$

PWM mixture model

We investigated an approach based on a mixture of PWMs. For that purpose, we used a comparable setup as for the pairwise model. However, instead of adding correlations to a given PWM, new PWMs were added to a mixture model. More precisely, a mixture of K PWMs, with $1 \leq K \leq 10$, was generated by using a K-means algorithm with a Hamming distance metrics on the initial

set of bound sites. This resulted in K clusters, each comprising n_k sites among the initial N sites. A PWM was generated on each of these clusters, with probability distribution \mathcal{P}_k . The mixture model of order K was then defined as [31]:

$$\mathcal{P}[s] = \sum_{k=1}^K p_k \mathcal{P}_k[s], \quad (23)$$

where $p_k = n_k/N$ is the cluster weight. Because a PWM has $3 \times L$ degrees of freedom (L of them being constrained by the summation of nucleotide probabilities to one) and there are $K - 1$ free weight parameters, the number of parameters corresponding to a mixture of order K is $3LK + (K - 1)$. As previously, the model showing minimal BIC score was used for sites detection, a new set of PWMs and weights p_k was generated by clustering the set of detected sites and the procedure was iterated until convergence to a stable set of sites.

Metastable minima of the pairwise interaction model and their basins of attractions

We defined the basins of attraction of a pairwise interaction model energy landscape, in the following fashion. Let s be a site with energy $\mathcal{H}(s)$. We looked for the nucleotides that could be changed to minimize $\mathcal{H}(s)$. If such nucleotides existed, one of them was chosen at random, and its value was updated. One local minimum of the energy landscape, or metastable state, was reached when no such nucleotide could be found. The basin of attraction of a metastable state was then defined as the ensemble of sites that fell to this metastable state when their energy was minimized following the above procedure. We computed metastable states and their basins of attraction for the final set of bound sites obtained with the best pairwise model. A PWM was learned on each basin of attraction, leading to a set of representative PWMs, with different weights representing different proportions of bound sites in their basins.

Computation of the Direct Information

We wanted to build a quantity based solely on direct interactions $J_{i,j}$ between nucleotides, discarding indirect interactions. To this end, we used the interaction parameters obtained from the pairwise model to build the direct dinucleotide probability function:

$$P_{i,j}^d(a, a') = e^{\tilde{h}_i(a) + \tilde{h}_j(a') + J_{i,j}(a, a')} / \mathcal{Z}_{i,j}, \quad (24)$$

where

$$\mathcal{Z}_{i,j} = \sum_{a, a'} e^{\tilde{h}_i(a) + \tilde{h}_j(a') + J_{i,j}(a, a')}.$$

The 8 effective fields \tilde{h}_i and \tilde{h}_j were fully determined by the constraints that the direct probability function matches the observed one-point frequencies:

$$\begin{aligned} \sum_{a'} P_{i,j}^d(a, a') &= P_i(a), & a' \in \{A, C, G, T\}, \\ \sum_a P_{i,j}^d(a, a') &= P_j(a'), & a \in \{A, C, G, T\}. \end{aligned} \quad (25)$$

The normalization of the probabilities $\sum_a P_i(a) = 1$, served to reduce this system to 6 equations. The fields $\tilde{h}_i(a)$, which are determined up to a constant, were fixed by the gauge condition that they vanished for the nucleotide A , $\tilde{h}(A) = 0$. The system was solved using the Levenberg-Marquadt algorithm with $\lambda = 0.005$.

The Direct Information [37] was then defined as:

$$DI_{i,j} = \sum_{a,a'} P_{i,j}^d(a, a') \log_2 \left(\frac{P_{i,j}^d(a, a')}{P_i(a)P_j(a')} \right). \quad (26)$$

As there is no upper bound for this direct information, we built a normalized version of the direct information:

$$NDI_{i,j} = \frac{DI_{i,j}}{\sqrt{S_i S_j}}, \quad (27)$$

where S_i denotes the entropy at position i . Note that $S_i = DI_{i,i}$ so that $NDI_{i,i} = 1$ for this maximally correlated case. On the contrary, independent nucleotides give $NDI_{i,j} = DI_{i,j} = 0$.

Participation Ratio

For each TF, an interaction weight was defined for each pair of nucleotides as

$$w_{i,j} = NDI_{i,j} / \sum_{i \neq j} NDI_{i,j}. \quad (28)$$

Self-interactions have no meaning here and were attributed $w_{i,i} = 0$. Let us note $N = L(L - 1)$ the number of possible interactions. Using our weight, one writes the Participation Ratio as:

$$R = \frac{1}{N \sum_{i \neq j} w_{i,j}^2}. \quad (29)$$

The interpretation is simple: if all weights are equal, $w_{i,j} = 1/N$ and $R = 1$, that is all possible interactions are represented. Conversely, if only one interaction accounts in the distribution budget, then $R = 1/N$, meaning that only one of all possible interactions is represented.

Distance between interactions

The previously defined interaction weights were averaged over all possible pairs of nucleotides at a given distance d of one another, yielding the distance distribution:

$$P(|i - j| = d) = \mathcal{Z}^{-1} \frac{1}{N-d} \sum_{|i-j|=d} w_{i,j}, \quad (30)$$

where

$$\mathcal{Z} = \sum_{d=1}^{N-1} \frac{1}{N-d} \sum_{|i-j|=d} w_{i,j} \quad (31)$$

is a normalization factor. Note that we introduced a correction $1/(N-d)$ to account for finite-size effects, namely the fact that randomly distributed interactions will lead to an overrepresentation of nearest neighbours interactions just because these are more numerous.

Interaction matrix and Hopfield patterns

In the Hamiltonian shown in (1), only $16L(L - 1)/2$ terms appear in the interaction budget: indeed, we forbid self-interactions (already accounted for by the local field h) and do not count the interactions twice. However, we can straightforwardly extend the interaction matrix to a full symmetric matrix $\hat{J}_{(i,a),(j,b)}$ of size $(4L)^2$, with $4L$ -valued indices $(i, a), i \in \{1, \dots, L\}, a \in \mathcal{A}$. The matrix \hat{J} is such that for $i > j$, $\hat{J}_{(i,a),(j,b)} = J_{i,j}(a, b)$ with furthermore $\hat{J}_{(i,a),(i,b)} = 0$ and $\hat{J}_{(i,a),(j,b)} = \hat{J}_{(j,b),(i,a)}$. The energy of a sequence s can then be written with these notations

$$\sum_{1 \leq i < j \leq L} J_{i,j}(s_i, s_j) = \frac{1}{2} \sum_{i=1}^L \sum_{j=1}^L \hat{J}_{(i,s_i),(j,s_j)} = v(s)^\dagger \hat{J} v(s), \quad (32)$$

where in the last equality the \dagger sign denotes vector transposition and we have introduced the $4L$ vector $v(s)$ associated to sequence s , $v(s)_{i,a} = 1$ if $a = s_i$ and $v(s)_{i,a} = 0$ otherwise.

Since the matrix \hat{J} is symmetric, it can be diagonalized in an orthonormal basis of eigenvectors ξ^k , $k = 1, \dots, L$ with real eigenvalues λ_k ,

$$\hat{J} = \sum_k \lambda_k \xi^k \xi^{k\dagger}. \quad (33)$$

Denoting by $\xi_{(i,a)}^k$ the coordinates of the k -th eigenvector then, one can rewrite Eq. (32) as

$$\sum_{1 \leq i < j \leq L} J_{i,j}(s_i, s_j) = \frac{1}{2} \sum_{k=1}^{4L} \lambda_k \left(\sum_{i=1}^L \xi_{(i,s_i)}^k \right)^2. \quad (34)$$

Finally, the full Hamiltonian is given by:

$$\mathcal{H} = - \sum_i h_i(s_i) - \frac{1}{2} \sum_{k=1}^{4L} \lambda_k \left(\sum_{i=1}^L \xi_{(i,s_i)}^k \right)^2. \quad (35)$$

Acknowledgments

We wish to thank PY Bourguignon and I Grosse for stimulating discussions at a preliminary stage of this

work.

-
- [1] F. Spitz and E. E. Furlong, *Nat. Rev. Genet.* **13**, 613 (2012).
- [2] J. A. Stamatoyannopoulos, *Genome Res.* **22**, 1602 (2012).
- [3] W. W. Wasserman and A. Sandelin, *Nat. Rev. Genet.* **5**, 276 (2004).
- [4] O. G. Berg and P. H. von Hippel, *J Mol Biol* **193**, 723 (1987).
- [5] G. D. Stormo and D. S. Fields, *Trends Biochem Sci* **23**, 109 (1998).
- [6] T. K. Man and G. D. Stormo, *Nucleic Acids Res* **29**, 2471 (2001).
- [7] P. V. Benos, M. L. Bulyk, and G. D. Stormo, *Nucleic Acids Res* **30**, 4442 (2002).
- [8] M. L. Bulyk, P. L. F. Johnson, and G. M. Church, *Nucleic Acids Res* **30**, 1255 (2002).
- [9] A. Jolma, J. Yan, T. Whitington, J. Toivonen, K. R. Nitta, P. Rastas, E. Morgunova, M. Enge, M. Taipale, G. Wei, et al., *Cell* **152**, 327 (2013).
- [10] G. Badis, M. F. Berger, A. A. Philippakis, S. Talukder, A. R. Gehrke, S. A. Jaeger, E. T. Chan, G. Metzler, A. Vedenko, X. Chen, et al., *Science* **324**, 1720 (2009), URL <http://eutils.ncbi.nlm.nih.gov/entrez/eutils/elink.fcgi?dbfrom=pubmed&id=19443739&retmode=ref&cmd=prlinks>.
- [11] Y. Zhao and G. D. Stormo, *Nat. Biotechnol.* **29**, 480 (2011).
- [12] Q. Zhou and J. S. Liu, *Bioinformatics* **20**, 909 (2004).
- [13] M. Hu, J. Yu, J. M. G. Taylor, A. M. Chinaiyan, and Z. S. Qin, *Nucleic Acids Res* **38**, 2154 (2010).
- [14] Y. Barash, G. Elidan, N. Friedman, and T. Kaplan, in *Proceedings of the seventh annual international conference on Research in computational molecular biology* (ACM, 2003), pp. 28–37.
- [15] E. Sharon, S. Lubliner, and E. Segal, *PLoS Comput Biol* **4**, e1000154 (2008).
- [16] E. Schneidman, M. J. Berry, R. Segev, and W. Bialek, *Nature* **440**, 1007 (2006), URL <http://www.nature.com/nature/journal/v440/n7087/full/nature04701.html>.
- [17] J. Shlens, G. D. Field, J. L. Gauthier, M. I. Grivich, D. Petrusca, A. Sher, A. M. Litke, and E. J. Chichilnisky, *J Neurosci* **26**, 8254 (2006), URL <http://www.jneurosci.org/cgi/content/full/26/32/8254>.
- [18] Y. Ikegaya, G. Aaron, R. Cossart, D. Aronov, I. Lampl, D. Ferster, and R. Yuste, *Science* **304**, 559 (2004).
- [19] A. Roxin, V. Hakim, and N. Brunel, *J. Neurosci.* **28**, 10734 (2008).
- [20] M. Weigt, R. A. White, H. Szurmant, J. A. Hoch, and T. Hwa, *Proc Natl Acad Sci USA* **106**, 67 (2009), URL <http://www.pnas.org/content/106/1/67.long>.
- [21] T. Mora, A. M. Walczak, W. Bialek, and C. G. Callan, *Proc Natl Acad Sci USA* **107**, 5405 (2010), URL <http://www.pnas.org/content/107/12/5405>.
- [22] W. Bialek, A. Cavagna, I. Giardina, T. Mora, E. Silvestri, M. Viale, and A. M. Walczak, *Proc Natl Acad Sci USA* **109**, 4786 (2012).
- [23] Y. Zhao, S. Ruan, M. Pandey, and G. D. Stormo, *Genetics* **191**, 781 (2012).
- [24] Y. Cao, Z. Yao, D. Sarkar, M. Lawrence, G. J. Sanchez, M. H. Parker, K. L. MacQuarrie, J. Davison, M. T. Morgan, W. L. Ruzzo, et al., *Dev Cell* **18**, 662 (2010).
- [25] S. Heinz, C. Benner, N. Spann, E. Bertolino, Y. C. Lin, P. Laslo, J. X. Cheng, C. Murre, H. Singh, and C. K. Glass, *Mol Cell* **38**, 576 (2010).
- [26] R. Zinzen, C. Girardot, J. Gagneur, M. Braun, and E. Furlong, *Nature* **462**, 65 (2009).
- [27] X. Chen, H. Xu, P. Yuan, F. Fang, M. Huss, V. B. Vega, E. Wong, Y. L. Orlov, W. Zhang, J. Jiang, et al., *Cell* **133**, 1106 (2008).
- [28] I. Dunham and et al., *Nature* **489**, 57 (2012).
- [29] T. Cover and J. Thomas, *Elements of information theory* (Wiley-interscience, 2006).
- [30] R. Baxter, *Exactly solved models in statistical mechanics* (Dover Publications, 2008).
- [31] C. Bishop et al., *Pattern recognition and machine learning* (Springer New York, 2006).
- [32] G. E. Crooks, G. Hon, J.-M. Chandonia, and S. E. Brenner, *Genome Res* **14**, 1188 (2004).
- [33] M. Weigt, R. A. White, H. Szurmant, J. A. Hoch, and T. Hwa, *Proc Natl Acad Sci USA* **106**, 67 (2009).
- [34] J. J. Hopfield, *Proc. Natl. Acad. Sci. U.S.A.* **79**, 2554 (1982).
- [35] S. Cocco, R. Monasson, and V. Sessak, *Phys Rev E Stat Nonlin Soft Matter Phys* **83**, 051123 (2011).
- [36] G. Tkacik, E. Schneidman, M. J. B. II, and W. Bialek, *arXiv q-bio.NC* (2006), 4 pages, 3 figures, [q-bio/0611072v1](http://arxiv.org/abs/q-bio/0611072v1), URL <http://arxiv.org/abs/q-bio/0611072v1>.
- [37] F. Morcos, A. Pagnani, B. Lunt, A. Bertolino, D. S. Marks, C. Sander, R. Zecchina, J. N. Onuchic, T. Hwa, and M. Weigt, *Proc Natl Acad Sci USA* **108**, E1293 (2011).
- [38] J. I. Sulkowska, F. Morcos, M. Weigt, T. Hwa, and J. N. Onuchic, *Proc Natl Acad Sci USA* **109**, 10340 (2012).
- [39] T. A. Hopf, L. J. Colwell, R. Sheridan, B. Rost, C. Sander, and D. S. Marks, *Cell* **149**, 1607 (2012).
- [40] D. S. Marks, L. J. Colwell, R. Sheridan, T. A. Hopf, A. Pagnani, R. Zecchina, and C. Sander, *PLoS ONE* **6**, e28766 (2011).
- [41] M. Lassig, *BMC Bioinformatics* **8 Suppl 6**, S7 (2007).
- [42] A. Moses, D. Chiang, D. Pollard, V. Iyer, and M. Eisen, *Genome biology* **5**, R98 (2004).
- [43] R. Siddharthan, E. Siggia, and E. van Nimwegen, *PLoS Comput Biol* **1**, e67 (2005).
- [44] H. Rouault, K. Mazouni, L. Couturier, V. Hakim, and F. Schweisguth, *Proc Natl Acad Sci U S A* **107**, 14615 (2010).
- [45] Z. Bao and S. R. Eddy, *Genome Res* **12**, 1269 (2002).
- [46] E. Wingender, X. Chen, R. Hehl, H. Karas, I. Lieblich, V. Matys, T. Meinhardt, M. Prüss, I. Reuter, and

- F. Schacherer, Nucleic Acids Res **28**, 316 (2000).
- [47] E. Jaynes, Physical review **108**, 171 (1957).
- [48] C. Shannon, Bell Syst Tech J **27**, 623 (1948).

Supporting Figures

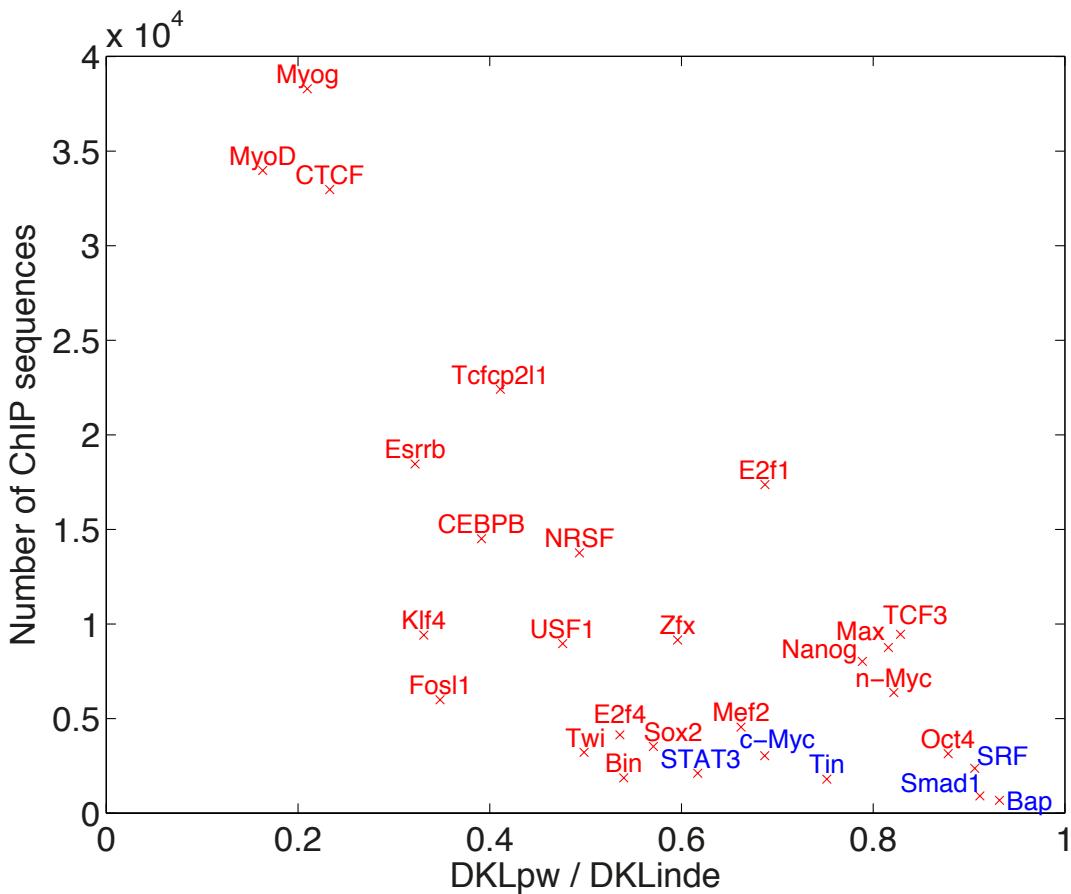
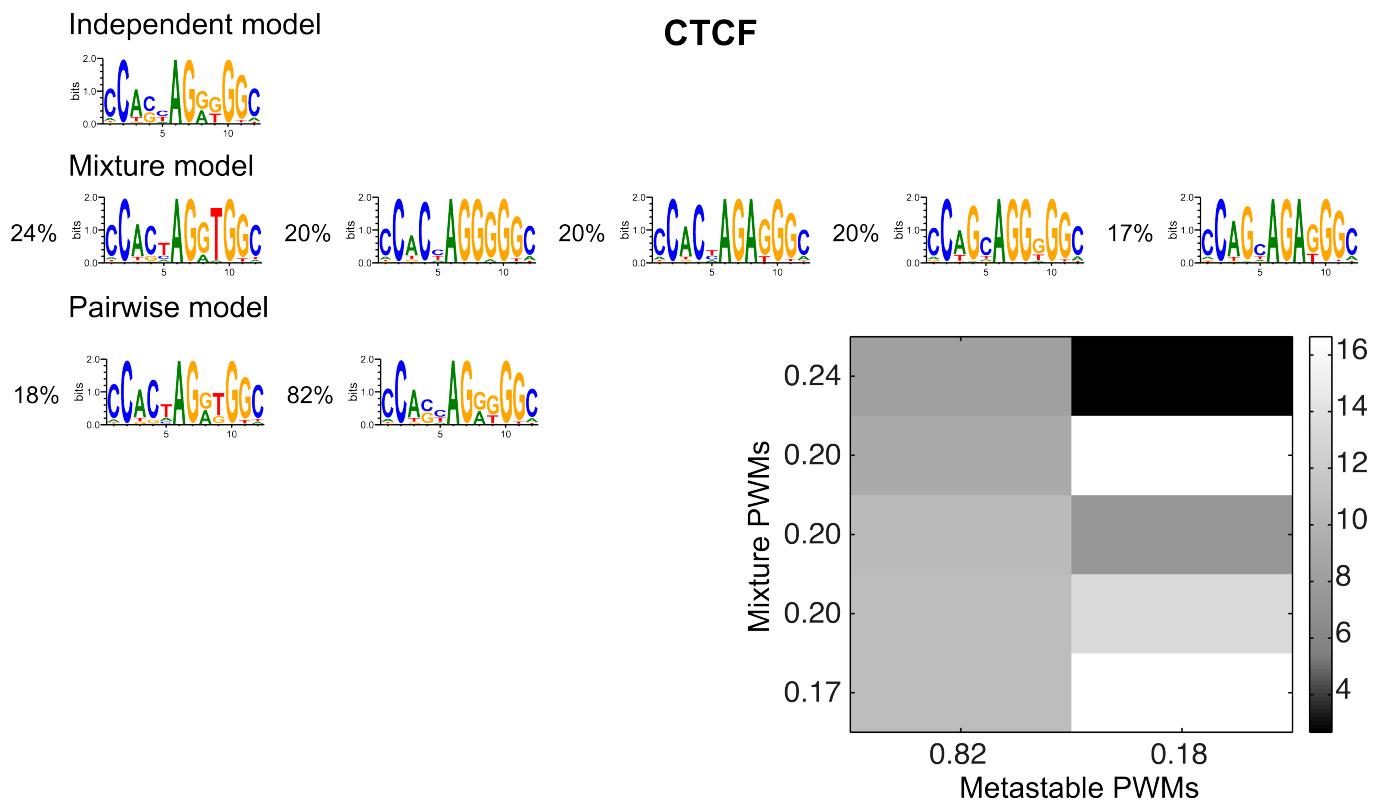
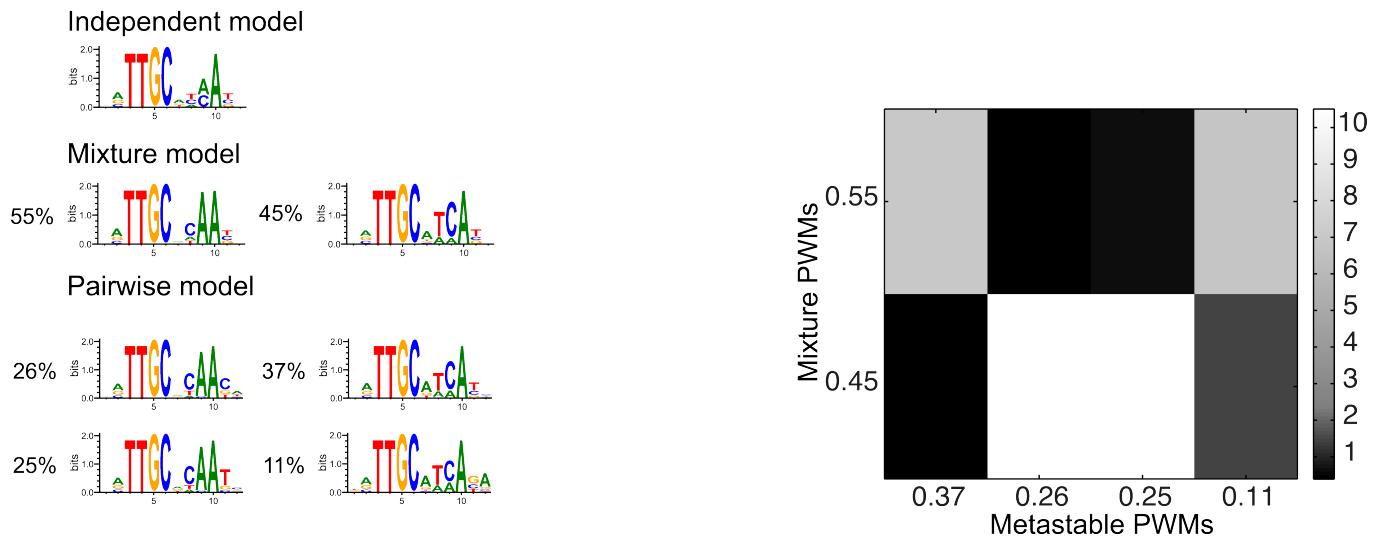
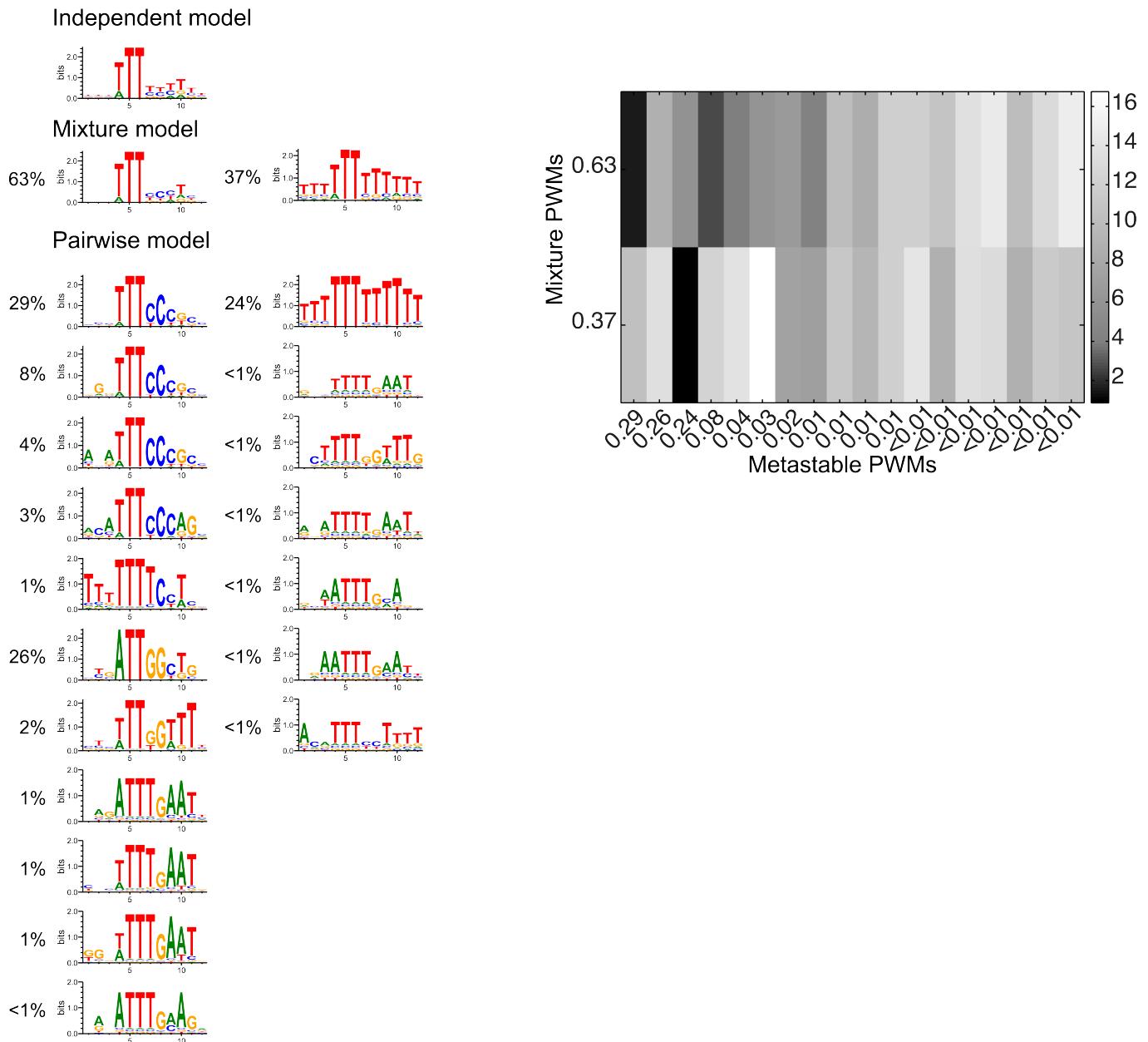
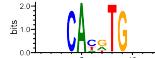
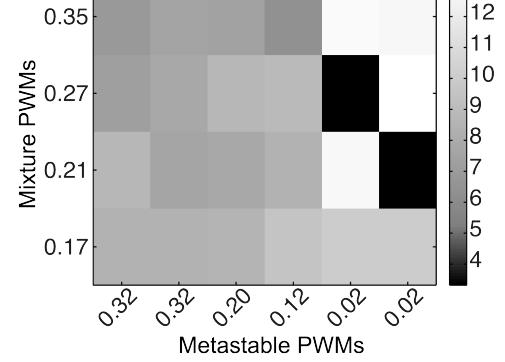
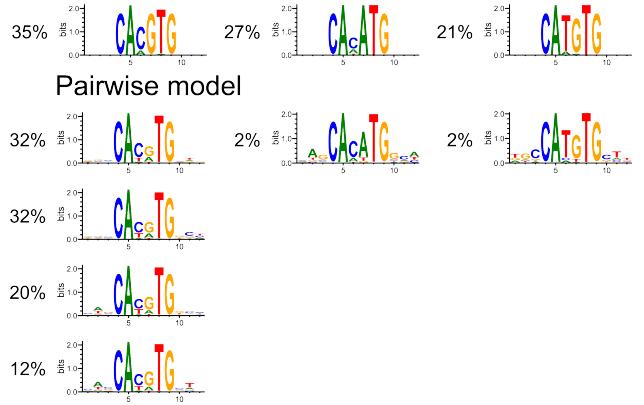
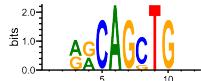
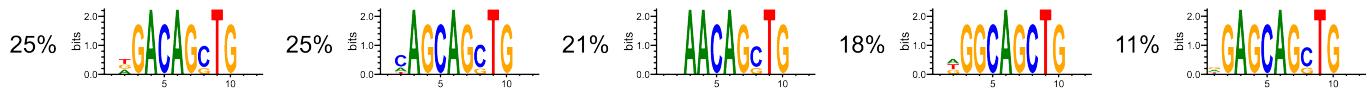
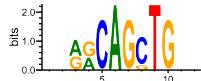
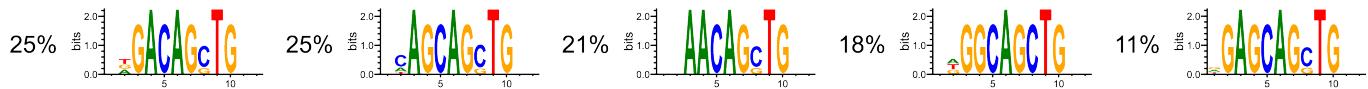
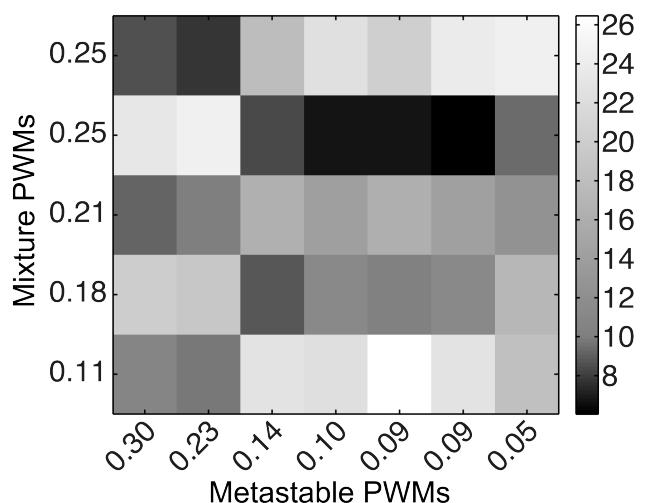
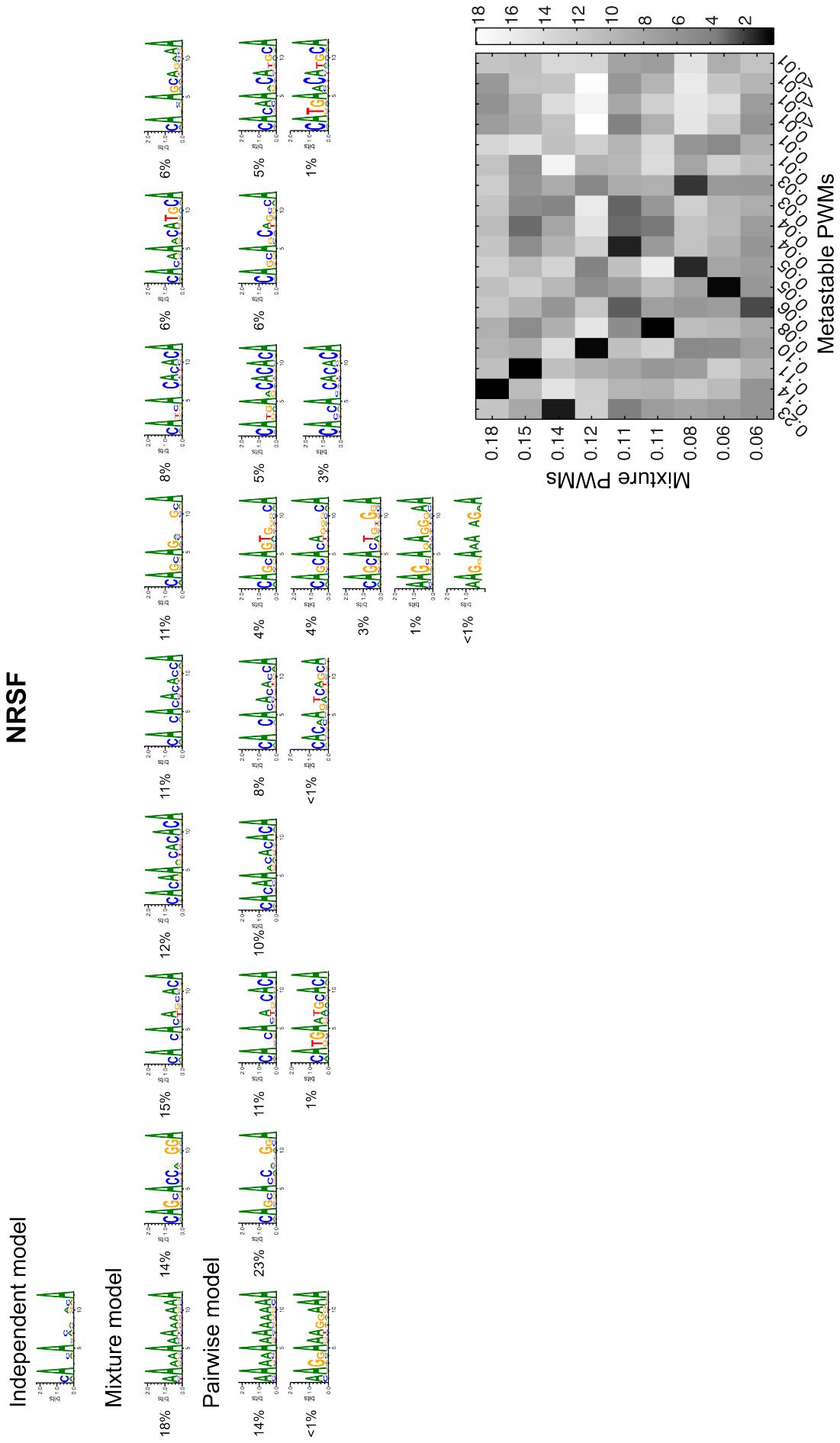


FIG. S1: **Dependence of the fit on the number of ChIP sequences.** For each TF, the number of available ChIP sequences is plotted *vs.* the improvement in the description of its TFBS statistics, provided by the he pairwise model as compared to the PWM independent model. The latter is quantified by the ratio of DKL between the respective model probability distributions and the experimental ones provided by the ChIP data, DKL_{pw}/DKL_{inde} . The improvement afforded by the pairwise model is clearly seen to be correlated to the number of ChIP sequences available.

CEPB

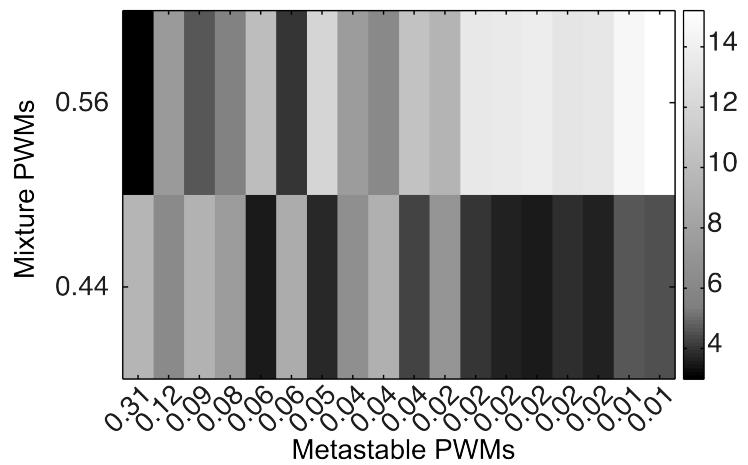
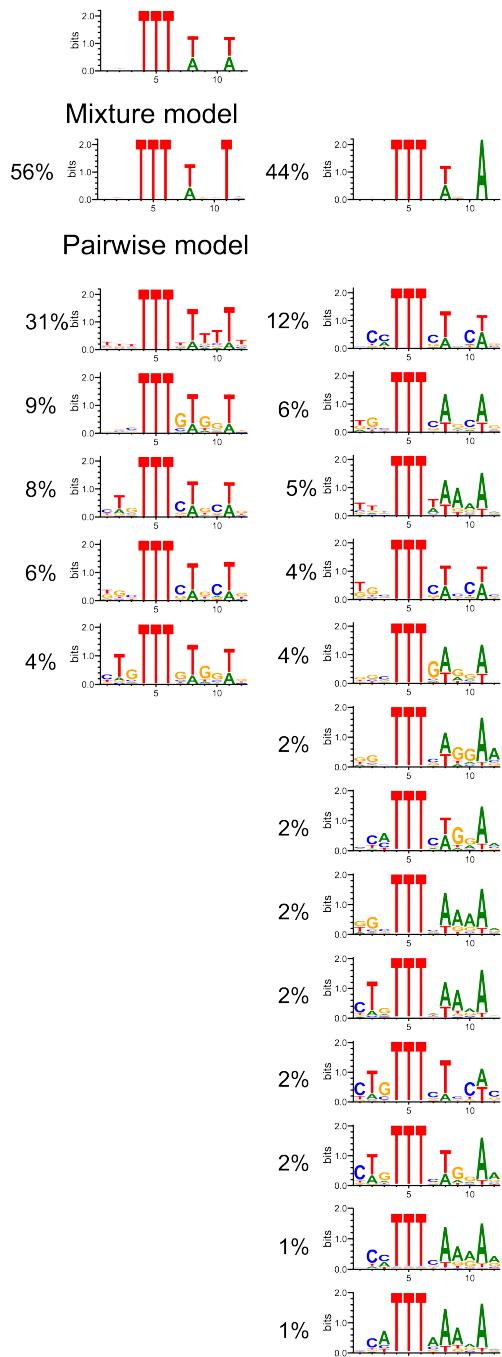
E2f4

Max**Independent model****Mixture model****Independent model****Mixture model****Myog****Independent model****Mixture model****Mixture PWMs**



Tcf3

Independent model



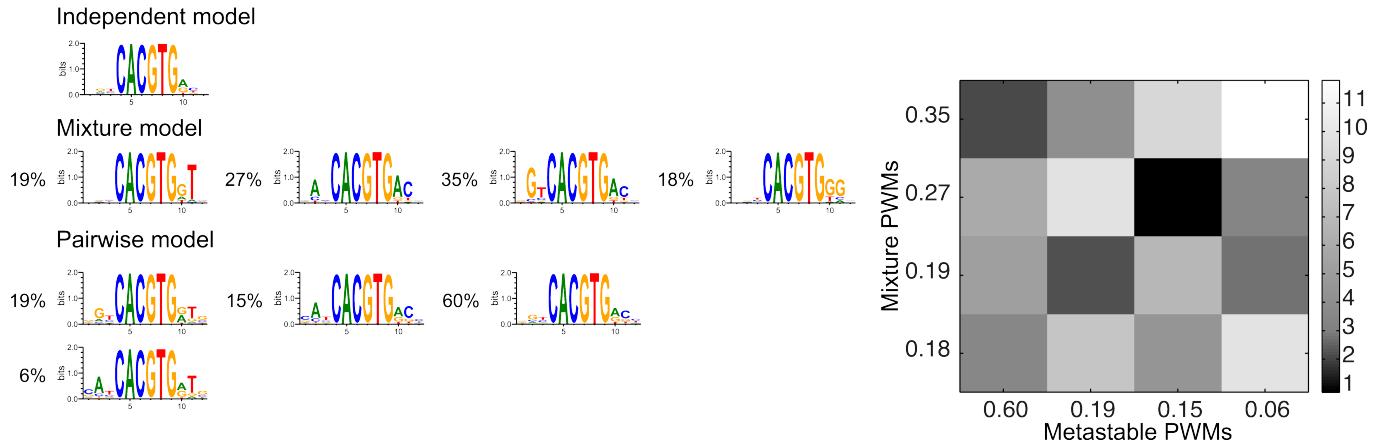
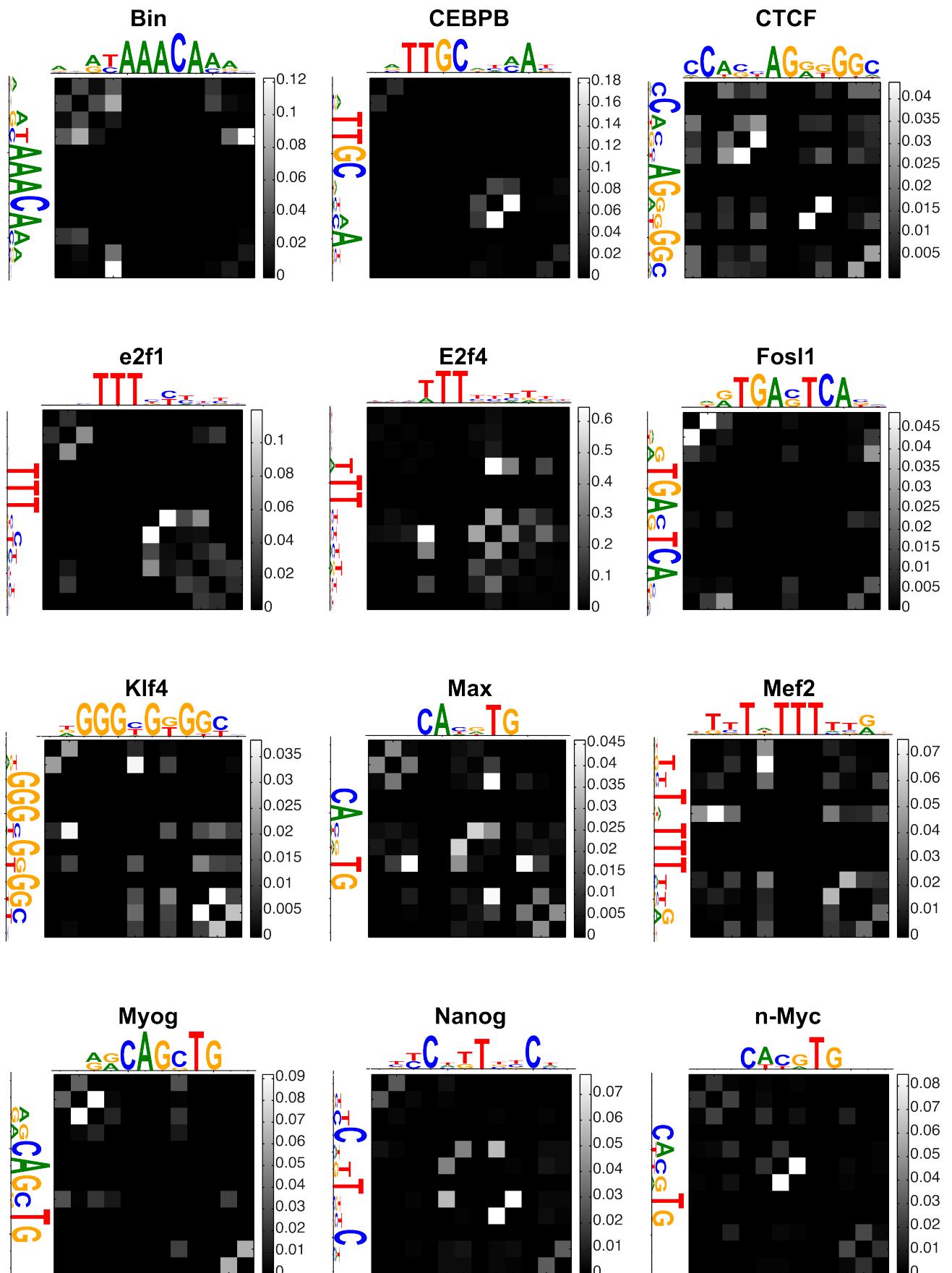
USF1

FIG. S2: Same as Figure 6 of the main text for all considered factors described by a mixture model with two or more PWMs.



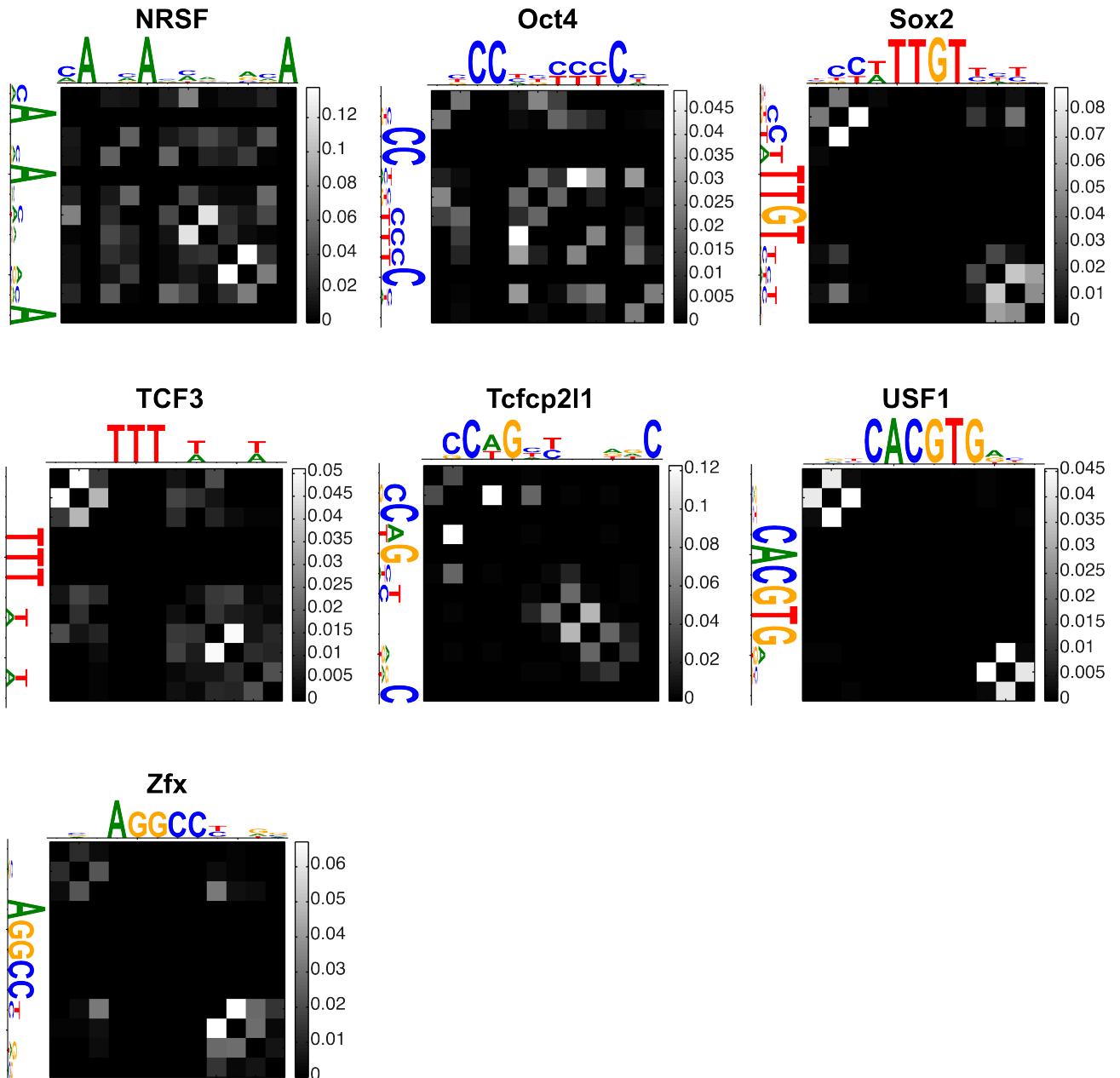


FIG. S3: Same as Figure 7 of the main text for the other considered factors.

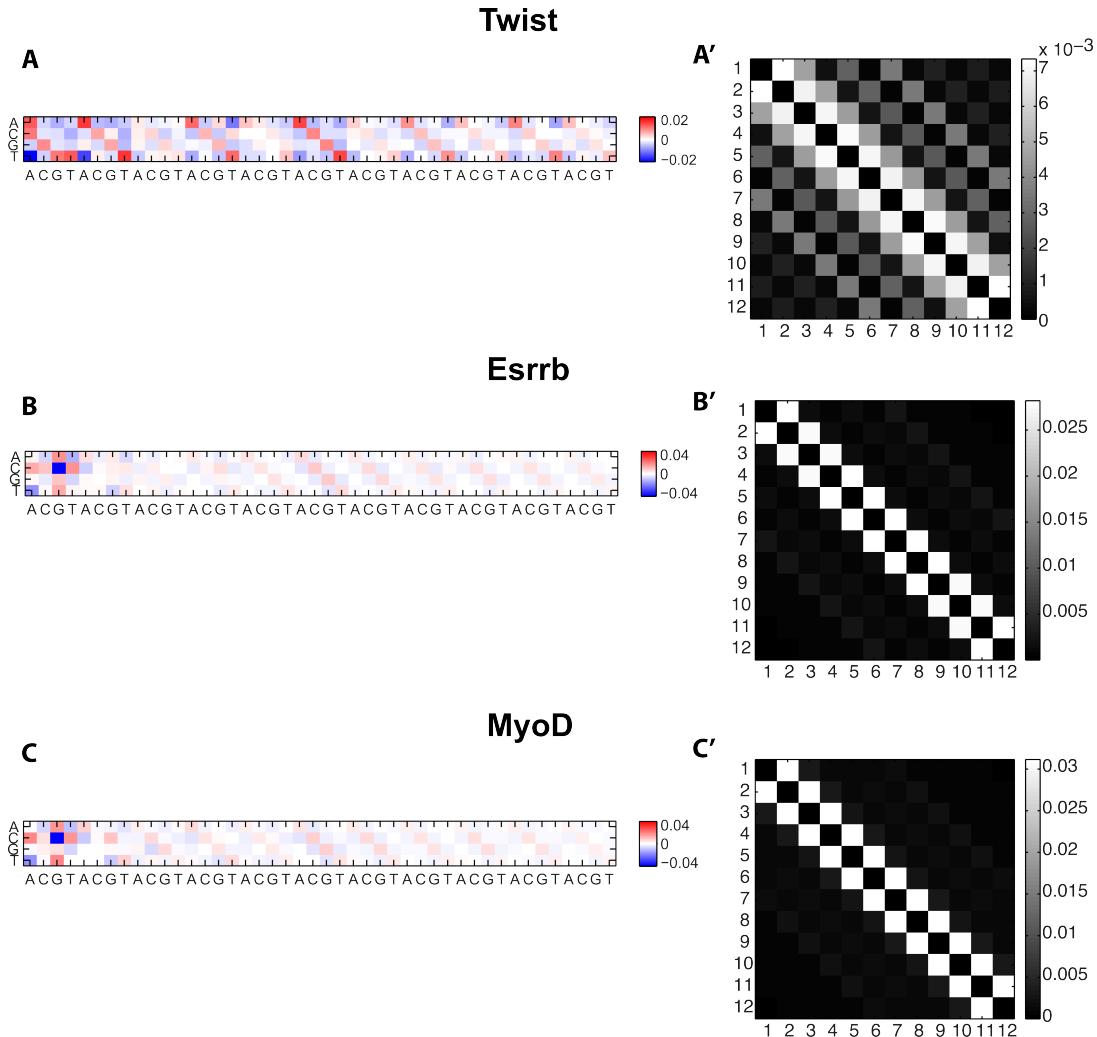


FIG. S4: **Background correlations** (A,B,C) Heat maps showing the correlations between nucleotides in the ChIP data of the 3 factors from the main text. Because of translation invariance, we only show the correlations between a nucleotide (rows) and the next nearest (first four columns) to farthest (last four columns) nucleotides, using the binding site length of $L = 12$. We see in the Drosophila data the appreciable presence of repeated sequences (of type AA, TT, CC, and GG). In the mammalian data sets, we observe the known CpG depletion. (A',B',C') Heat maps showing the values of the Normalized Direct Information between pairs of nucleotides.

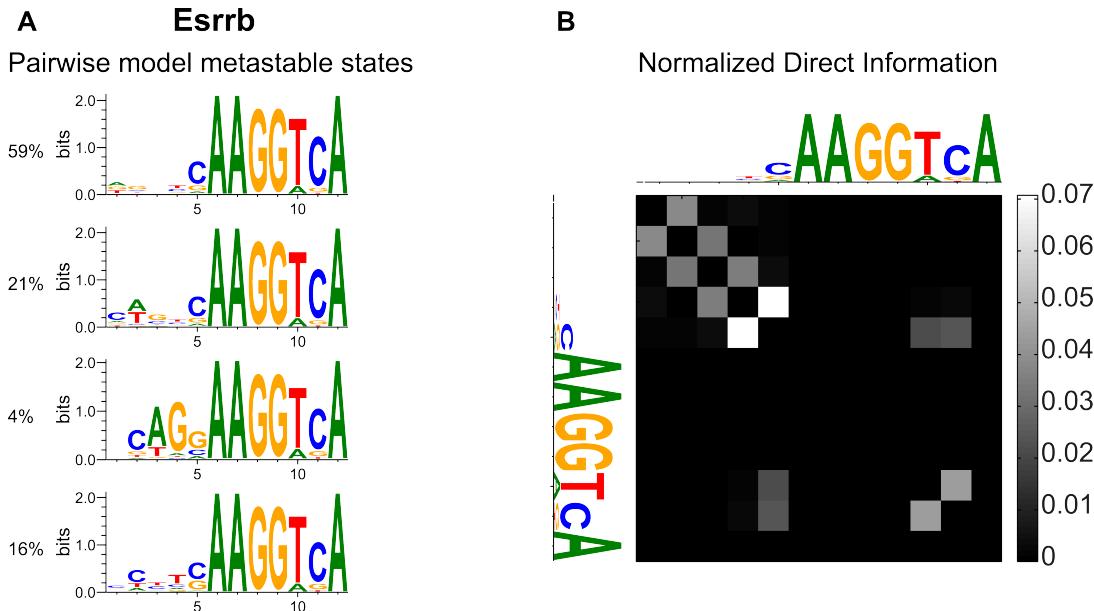


FIG. S5: Variable spacer length We learned a pairwise model for Esrrb including the 4 flanking nucleotides on the left of the main motif. (A) The metastable states of this model show a feature not captured in the main text where binding sites are defined symmetrically around the center of mass of the information content: namely a ‘CAG’ trinucleotide with variable spacer length from the main motif. This feature is apparent in the first 3 logos shown here. (B) The contribution of this trinucleotidic interaction to the Direct Information is captured through strong direct links between the 4 flanking nucleotides, showing that the pairwise model is implicitly able to capture higher order correlations. Logos from the PWM model are surrounding the heatmap for clarity.

2.5 Analyse thermodynamique des modèles

2.5.1 Chaleur spécifique

En plus des résultats présentés dans l'article, nous nous sommes intéressés à une quantité classique de la thermodynamique : la chaleur spécifique ou capacité calorifique. Considérons un modèle décrit par la statistique de Boltzmann à la température inverse $\beta = 1/T$ (on omet la constante de Boltzmann k en l'intégrant à l'énergie) :

$$P(s) = \frac{1}{Z} e^{-\beta E(s)} \quad (2.24)$$

Le cas de l'équation 2.23 correspond au cas particulier $\beta = 1$. Nous voulons voir comment l'amplification ou la diminution globale de l'écart entre les énergies affecte la possibilité du système d'explorer les différents états possibles. À température nulle ($T \rightarrow 0$ ou $\beta \rightarrow \infty$), le système reste dans le niveau fondamental de minimum d'énergie et de probabilité 1, alors qu'à des températures non nulles le système à l'énergie E_0 transite vers un état d'énergie supérieure E_1 avec une probabilité $\propto \exp(-\beta(E_1 - E_0))$. Lorsqu'un paysage énergétique est composé de plusieurs puits d'énergie séparés par des barrières énergétiques importantes, on s'attend à avoir une (ou plusieurs) températures critiques à partir desquelles de fortes différences d'énergie deviennent franchissables. L'énergie moyenne peut alors être significativement affectée, sautant soudainement à une nouvelle valeur du fait du poids des nouveaux états explorés.

La chaleur spécifique permet de caractériser ces sauts soudains d'énergie moyenne lors de la variation de la température, caractéristiques des transitions de phase. Elle mesure simplement la variation de l'énergie moyenne lors d'une variation de température :

$$C(T) = \frac{d\langle E \rangle}{dT} \quad (2.25)$$

où

$$\langle E \rangle = \sum_{\{s\}} E(s) \frac{e^{-\beta E(s)}}{Z} \quad (2.26)$$

Cette chaleur spécifique peut par ailleurs s'écrire sous une forme plus utile :

$$\begin{aligned}
\frac{d\langle E \rangle}{dT} &= -\beta^2 \frac{d\langle E \rangle}{d\beta} \\
&= -\beta^2 \left[\sum_{\{s\}} E(s) \left(-E(s)e^{-\beta E(s)} \right) \frac{1}{Z} + \sum_{\{s\}} E(s)e^{-\beta(E(s))} \left(-\frac{dZ}{d\beta} \frac{1}{Z^2} \right) \right] \\
&= \beta^2 \left[\langle E^2 \rangle - \langle E \rangle^2 \right]
\end{aligned} \tag{2.27}$$

Ainsi, la chaleur spécifique $C(T)$ est directement accessible en regardant les corrélations de l'énergie sur l'ensemble des états du système, ce qui peut se calculer simplement à partir des modèles de fixation. Nous avons calculé la variation de $C(T)$ en fonction de la température pour les modèles indépendant et avec dépendances obtenus en 2.4 pour les différents TFs étudiés. Une température fictive est introduite dans les modèles en multipliant les énergies par β , afin de se placer dans le cadre de l'équation 2.24. Les résultats sont montrés en figure 2.3 (modèle indépendant en bleu, modèle de Potts en rouge). On observe pour la plupart des facteurs l'existence de deux pics de chaleur spécifique pour des températures de l'ordre de $T \sim 10^{-1}$ et $T \sim 5$ (par exemple, $T = 0.4$ et $T = 2.8$ dans le cas du modèle indépendant de Twist). Il y a de légères variations entre les deux modèles : notamment, le premier pic semble renforcé par le modèle de Potts dans plusieurs cas (par exemple, E2f4, NRSF, TCF3 ou Twist). Néanmoins, le nombre de pics (ou de transitions de phases) reste le même.

2.5.2 Lien avec les valeurs des champs et des couplages

Afin de comprendre l'existence des pics de chaleur spécifique et les énergies (températures) associées, il faut revenir aux modèles d'énergie. Lorsque l'on regarde l'histogramme des valeurs absolues de h_i obtenues dans les modèles indépendant des différents TFs étudiés, on trouve plusieurs valeurs typiques autour de 10^{-4} , 1 et 10 (fig. 2.4A). Celles-ci peuvent s'expliquer de la manière suivante. Dans le modèle indépendant, les champs sont simplement le logarithme naturel de la probabilité d'observer un nucléotide a à une position i donnée $h_i(a) = -\log P_i(a)$ (la jauge est choisie telle que $Z_j = 1$). En valeur absolue, les champs h_i proches de 0 ($h_i \sim 10^{-4} - 10^{-3}$) correspondent aux nucléotides très conservés (toujours observés), les valeurs autour de 1 correspondent à des nucléotides dégénérés (i.e également observés : $|\log(1/4)| \sim 1.4$) et les valeurs autour de 10 correspondent aux nucléotides qui ne sont jamais observés, au pseudocount près (pour un pseudocount de 1 et 10^4 séquences, $|\log(10^{-4})| \sim 9.2$). On peut maintenant mieux comprendre les pics de chaleur spécifique. À

2.5. Analyse thermodynamique des modèles

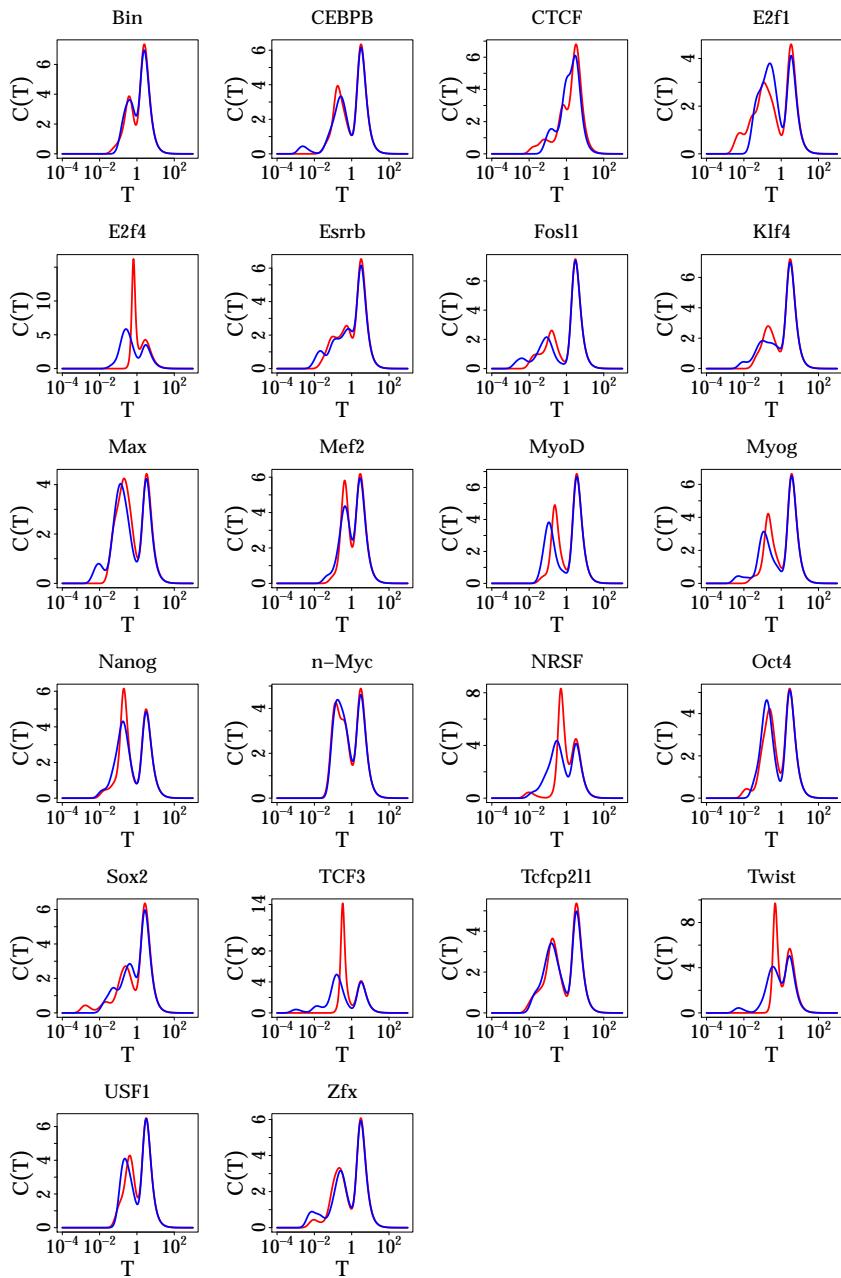


FIGURE 2.3 – Chaleur spécifique pour différents TFs.

La chaleur spécifique (l'équivalent de la capacité calorifique en thermodynamique) $C(T) = d\langle E \rangle / dT$ est tracée en fonction de la température kT (échelle logarithmique) pour les différents TFs considérés. Le modèle indépendant (bleu) et le modèle de Potts avec interactions (rouge) sont comparables dans la plupart des cas.

température nulle, seuls les sites consensus sont accessibles. Lorsque la température se rapproche de 1, les nucléotides dégénérés d'énergie $h_i \sim 1$ deviennent accessibles, augmentant significativement la valeur de l'énergie moyenne (premier pic). Puis, lorsque la température se rapproche de 10, les nucléotides non observés d'énergie $h_i \sim 10$ deviennent à leur tour accessibles, augmentant à nouveau l'énergie moyenne (deuxième pic).

Dans le cas du modèle de Potts (fig. 2.4B), les champs h_i prennent des valeurs proches de celles obtenues avec le modèle indépendant. Par ailleurs, les interactions $J_{i,j}$ sont réparties autour d'un mode centré autour de $J_{i,j} \sim 0.5$, ce qui correspond l'échelle d'énergie du premier pic. Ainsi, le renforcement du premier pic de chaleur spécifique par rapport au cas indépendant observé pour plusieurs TFs de la figure 2.3 peut s'expliquer par l'effet des termes d'interaction $J_{i,j}$.

2.6 Conclusion et perspectives

Nous avons analysé les dépendances au sein des sites de fixation liés *in vivo* pour différents facteurs de transcription Drosophiles et mammifères. Nous avons comparé les performances d'un modèle PWM, d'un modèle de mélange de PWMs, et d'un modèle de Potts, en utilisant un critère bayésien (BIC) pénalisant les modèles à grand nombre de paramètres. Nous avons exhibé l'existence de corrélations faibles dont la prise en compte permet de significativement améliorer la description des données, le modèle de Potts étant significativement supérieur aux deux autres modèles dans la plupart des cas (22/28). Les interactions ont été étudiées systématiquement, montrant notamment une prépondérance des interactions entre plus proches voisins. Nous avons établi une correspondance entre les PWMs du modèle de mélange et les PWMs décrivant les états métastables du paysage énergétique généré par le modèle de Potts. Enfin, nous avons montré que les corrélations pouvaient être groupées en patterns de Hopfield ou « mémoires », et qu'un petit nombre était suffisant à reconstruire le paysage d'interactions.

Une perspective intéressante de ce travail serait de conduire la même analyse sur des données grande échelle obtenues *in vitro* par la méthode HT-SELEX (?). Notamment, certains des facteurs que nous avons étudiés *in vivo* sont représentés dans ces données, et il serait intéressant de voir les différences entre les modèles obtenus. Notamment, retrouve-t-on les mêmes corrélations ? Peut-on exhiber des spécificités de la fixation *in vivo*, où l'on s'attend à

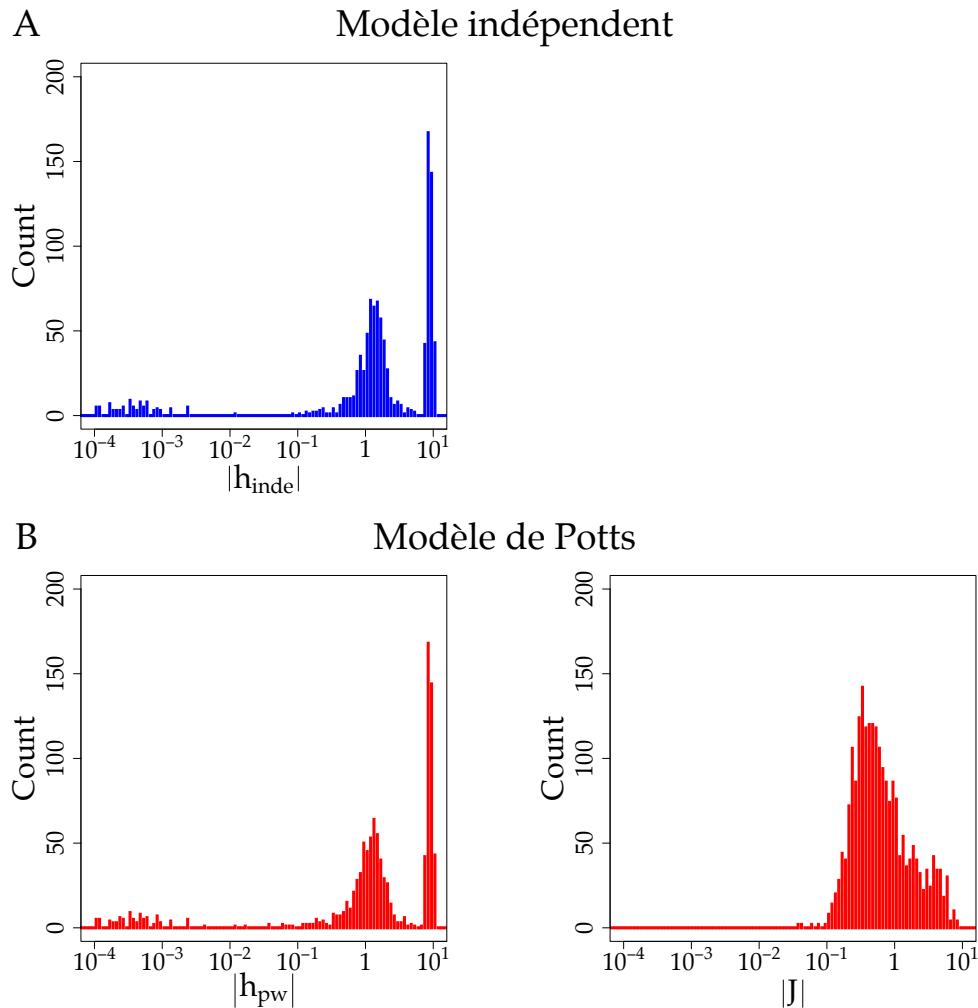


FIGURE 2.4 – Histogrammes des valeurs des champs h et couplages J .

Histogrammes réalisés à partir des valeurs obtenues pour l'ensemble des TFs. Les champs et les couplages sont montrés en valeur absolue sur une échelle logarithmique d'espace-ment 0.05, et les valeurs nulles ne sont pas représentées. (A) Champs h_{inde} dans le modèle indépendant. (B) Champs h_{pw} et couplages J dans le modèle de Potts.

avoir des effets provenant de diverses sources (fixation de nucléosomes, superposition de sites de fixations, ...)? Ces questions feront certainement l'objet d'un prochain travail.