

Introduction

La connaissance du génome dans son intégralité passe par son séquençage et sa cartographie. Cependant, la taille des génomes étant de plusieurs millions de bases, il est nécessaire de coupler les approches de biologie moléculaire et d'informatique, pour pouvoir gérer les quantités immenses de données.

La génomique a ainsi vu le jour, et a permis de caractériser l'organisation du génome et de mettre en évidence l'existence de séquences codantes et de séquences non codantes. Ces dernières sont présentes dans tous les génomes mais leur abondance est corrélée à la complexité des espèces, et donc les mécanismes de régulation des génomes.

Le génome humain serait constitué de 5% seulement de séquences codantes qu'on nomme «gènes». Des programmes bioinformatiques permettent l'analyse des séquences d'ADN et la prédition des gènes. La traduction *in silico* de la séquence nucléotidique en séquence d'acides aminés permet éventuellement de confirmer ou proposer une fonction de la protéine codée par le gène.

1- Séquences codantes

Une séquence codante est une séquence ADN double brin transcrrite en ARN et encadrée de séquences régulatrices. Dans un génome eucaryote, la séquence transcrrite est constituée d'exons et d'introns.

Les séquences régulatrices servent de signaux de début (promoteur, codon initiateur) et de fin de la transcription (codon stop)

Les gènes peuvent donc être prédits par des algorithmes grâce à ces caractéristiques communes aux séquences codantes.

Il faut distinguer deux types de gènes:

* Ceux qui codent pour des protéines et donc sont transcrits en ARN codant qui est l'ARNm, traduit ensuite en protéines

* Ceux qui codent des ARN non codants c'est à dire qui ne seront pas traduits en protéines.

Les ARN non codants sont impliqués dans l'expression du génome et sa régulation. Leur fonction repose sur leur structure et leur séquence:

* L'ARN ribosomique constitue 80% de l'ARN total dans les cellules eucaryotes.

* L'ARN de transfert constitue 15% des ARN dans une cellule eucaryote.

- * Les petits ARN nucléaires et nucléolaires, présents dans le noyau des cellules eucaryotes, interviennent dans l'épissage et la maturation des ARNr
- * Les ARN régulateurs sont impliqués dans les mécanismes d'interférences et régulent l'expression des gènes (transcription et traduction)

2- Sequences régulatrices

La transcription d'un gène est finement régulée dans le temps et l'espace. Un transcrit donné ne sera synthétisé que dans certains types cellulaires à une certaine étape du développement de l'individu et/ou différents facteurs de l'environnement.

Les informations nécessaires à cette régulation sont portées à la fois par des protéines régulatrices (éléments trans), et par des séquences nucléotidiques non transcrrites (éléments cis) placées en amont de la séquence codante.

Un complexe de transcription initie la transcription en interagissant avec le promoteur. Le promoteur est situé en amont du +1 de la transcription. Des séquences particulières permettent de reconnaître une séquence promotrice comme les boîtes TATA et les îlots CpG; une analyse bioinformatique peut révéler dans ces régions, des séquences nucléotidiques conservées et reconnues par les facteurs de transcription.

L'inhibition ou l'activation de la transcription dépend d'autres facteurs et séquences (enhancers et silencers)

Génome- Transcriptome - Protéome : Chez les eucaryotes contrairement aux procaryotes, le génome n'est pas corrélé au transcriptome ni le transcriptome au protéome. Un décalage important est noté entre le nombre de gènes et le nombre de transcrits et de protéines produits par une cellule eucaryote.

Ce décalage révèle l'importance des mécanismes de régulation post-transcriptionnelles et post-traductionnelles dans le niveau d'expression final d'une protéine dans la cellule.

Le mécanisme d'épissage alternatif explique en partie le décalage entre le nombre de gènes et le nombre de transcrits et explique ainsi la présence de séquences introniques dans les gènes eucaryotes.

3- Sequences non codantes

3.1 Sequences répétées

a- Sequences répétées dispersées

Les éléments transposables sont des séquences d'ADN capables d'être transférées, au sein du génome, d'un site donneur à un site cible sous l'effet d'une enzyme spécialisée. Les deux séquences ne présentent aucune homologie et donc le mécanisme d'insertion est une recombinaison illicite.

Les enzymes qui effectuent la transposition sont soit des transposases soit des intégrases (dans le cas de retrovirus par exemple); ces deux enzymes sont des recombinases.

Elles reconnaissent et se lient à des séquences spécifiques aux extrémités du transposon, excisent l'ADN et lient.

Le système transposon/transposase peut être utilisé comme outil pour la transgenèse.

On distingue deux classes d'éléments transposables:

Classe I ou retrotransposons à ARN

Classe II ou transposon à ADN

1- Transposons: transposent via un mode couper-coller, c'est-à-dire que le transposon est excisé de son site et intégré à un autre site sur le même chromosome ou un autre chromosome.

On parle de gène sauteur, car la séquence a changé de place

2- Retrotransposons: constituent 42% du génome humain, la majorité sont cependant inactifs.

Ils transposent par un mode copier-coller, c'est-à-dire que la séquence donneuse reste à sa place, elle est transcrrite en ARN qui est retrotranscrit en ADN, cette séquence est intégrée à un nouvel endroit du génome.

Les séquences L1 sont impliquées dans les cassures de l'ADN et les cancers et l'inactivation du chromosome X chez les mammifères; chez la drosophile, ils se sont substitués aux telomerases pour la maintenance des telomères

b- Séquences répétées en tandem

Les ADN satellites sont des répétitions de séquences nucléotidiques (de quelques paires de bases à plusieurs centaines) qui peuvent s'étendre sur plusieurs milliers de kilobases. Ces séquences, constituant majeur de l'hétérochromatine centromérique et péri-centrique, représentent une fraction importante des génomes eucaryotes. Paradoxalement, la fonction des ADN satellites demeure largement incomprise. En effet, ces séquences ont été longtemps négligées car considérées comme de l'ADN "poubelle" sans fonction cellulaire. Ces séquences jouent un rôle crucial dans la structure des chromosomes et peuvent influencer la régulation génomique. Bien que l'ADN satellite ne code pas pour des protéines, il est essentiel pour la

stabilité chromosomique et la division cellulaire. De plus, des maladies humaines telles que la dystrophie facio-scapulo-humérale ou certains types de cancers, sont associées à des perturbations de l'organisation de satellites

Les répétitions en tandem sont constituées de répétitions adjacentes plus ou moins nombreuses d'un monomère (motif) donné. Certains de ces éléments ont la propriété d'être variables en nombre de répétitions, une conséquence de phénomènes de mutation particuliers (glissement de polymérase et recombinaison inégale). Ces répétitions en tandem peuvent être séparées en trois classes distinctes : les satellites, les minisatellites et les microsatellites.

*Les satellites se dénissent comme des répétitions en tandem possédant un très grand nombre de répétitions, pouvant régulièrement atteindre plusieurs mégabases

*Les minisatellites ou VNTR (variable number tandem repeat), sont des répétitions en tandem d'un motif compris entre dix et une centaine de paires de bases.

* Les microsatellites ou STR (Short tandem repeats) ou SSR (simple sequence repeat), le motif répété est plus court (1-6 nts) long (10-60 nts), il s'agit de minisatellite ou VNTR (variable number tandem repeat)

3.2 Les Pseudogenes

Sont vestiges de genes qui ont perdu leur fonctionnalité. Les pseudogenes donc ressemblent aux genes, ils peuvent être reconnus grâce à des séquences caractéristiques des genes comme les promoteurs et les sites d'épissage

Ils proviennent soit:

* De la duplication de gene ancestral qui existe dans le genome et qui est actif, et ont accumulé beaucoup de mutations, soit:

- Au niveau de la séquence promotrice, ce qui empêche la reconnaissance par le complexe de transcription et donc l'initiation de ce mécanisme

- Ou au niveau de la séquence codante, par l'apparition de codons stop prematurés. Dans ce cas le pseudogene est transcrit, puisque le promoteur est intact, mais le cadre de lecture ne peut pas donner une protéine fonctionnelle (trop courte)

* Soit de la rétrotranscription d'ARNm, ce qui explique que certains pseudogenes sont dépourvus de séquences promotrices, ce qui empêche leur transcription et donc leur expression.

Il a été démontré que certains pseudogenes jouent un rôle essentiel dans la régulation de

l'expression de leurs genes parentaux. Les pseudogenes transcrits peuvent aussi former des ARN intervertis. Ils peuvent aussi réguler les oncogènes et les suppresseurs de tumeur.

4- Chromatine et régulation épigénétique

La régulation de l'expression des gènes eucaryotes implique plusieurs niveaux de contrôle. Un premier niveau dit chromatinien implique la structure de la chromatine.

* **La chromatine** : Dans les cellules eucaryotes, le génome est organisé en une structure complexe constituée d'ADN et de protéines histones, appelée « chromatine ». Cette structure permet la compaction de l'ADN mais aussi la régulation de ses fonctions et de son expression. La chromatine peut adopter deux états qui peuvent être distingués par la présence de modifications post-traductionnelles des histones et méthylation de l'ADN.

- Euchromatine apparaît décondensée au cours de l'interphase, elle correspond aux zones codantes du génome

- Heterochromatine : Région hautement condensée, garde le même état de condensation au cours du cycle cellulaire; elle se compose de

* heterochromatine constitutive: structure permanente et figée, contient peu de gènes, est formée principalement de zones répétées dont les plus grandes régions sont situées à proximité du centromère et des telomères

* heterochromatine facultative : structure dynamique, contient les régions codantes maintenues sous une forme compacte, transcriptionnellement inactive, selon le stade de développement et le type cellulaire

* **La régulation épigénétique** : Chez les eucaryotes, la régulation dite épigénétique implique :

- Les processus par lesquels le génotype, l'ensemble des gènes, engendrent le phénotype, les caractéristiques de l'organisme.

- Les changements dans l'expression des gènes, c'est-à-dire les changements des états de la transcription des gènes. Ces changements sont stables et héréditaires au cours des divisions cellulaires et n'impliquent aucun changement de la séquence des gènes

Les changements épigénétiques sont des caractères reversibles en général et reprogrammables selon le type cellulaire. En effet, elles permettent la différenciation cellulaire en imposant une empreinte génomique caractéristique de chaque type cellulaire. Les gènes sont allumés ou éteints par compaction de la chromatine qui empêche l'accès au promoteur des facteurs de transcription

Le profil de methylation, c'est a dire les positions des methylcytosines (C5) dans le génome, permet d'identifier les gènes ainsi régulés par ces modifications chimiques

Des altérations de cette régulation génétiques dont des épimutations, et sont impliquées dans de nombreuses pathologies et cancers.

* **le code histone** : La région NH₂ terminale des histones est accessible à l'extrémité du nucleosome

Les modifications post traductionnelles des histones s'effectuent dans cette region et entraînent des modifications des états de condensation de la chromatine (euchromatine et heterochromatine) et donc l'accessibilité des complexes protéiques de transcription au promoteur

Il existe une grande diversité des modifications post-traductionnelles des extrémités N-terminales des histones;

Les différentes combinaisons de ces modifications (code histone) agissent de manière coordonnée pour réguler l'expression des gènes. Chaque état transcriptionnel pourrait être associé à un code histone particulier

Les principaux types de modifications des histones sont :

- Acétylation des lysines
- Méthylation des lysines/arginine
- Phosphorylation des serines/tyrosines
- Ubiquitination des lysines
- Sumoylation des lysines

La combinaison des modifications post-traductionnelles des queues des histones constituent le « code histone », un code spécifique reconnu par le système des complexes protéiques « writer » et « reader ». Ce système permet d'interpréter ces profils de modifications et conduit à des processus biologiques spécifiques (régulation de la transcription, réparation de l'ADN, inactivation du chromosome X, organisation du centromère...)

Ces modifications permettent la mise en place de l'empreinte génétique, l'empreinte parentale ainsi que l'inactivation du chromosome X :

* **Empreinte génomique** : Définit de façon stable la répression ou l'expression des gènes dans des lignées cellulaires spécifiques (les deux alleles sont soumis à l'empreinte) en fonction du type cellulaire (différenciation cellulaire), selon les besoins de la cellule et/ou en réponse à des stimuli

* **Lyonisation**: Correspond à l'inactivation du chromosome X. Un chromosome

complet condensé chez la femme connu également sous le nom de corpuscule de Barr
(chromosome X condensé et inactif)

L'inactivation du chromosome X se met en place de manière aléatoire au cours du développement embryonnaire, puis elle se maintient dans les cellules somatiques

* **Empreinte parentale:** Est un mécanisme physiologique conduisant à l'inactivation de l'un des deux alleles parentaux de certains gènes, selon leur origine paternelle ou maternelle

- Empreinte maternelle : Le gène paternel s'exprime et le gène maternel ne s'exprime pas
- Empreinte paternelle : Le gène paternel ne s'exprime pas et le gène maternel s'exprime