

# Chromosomes and Chromosome Biology

The human genome contains large amounts of repetitive DNA (see lecture 1)

Retro-transposon type sequences between 6kb and 8kb long

Contain 2 open reading frames, encoding a reverse transcriptase and an integrase

The LINE element is transcribed by cellular RNA polymerase

The LINE transcript is then copied into DNA by the element-encoded reverse transcriptase

The integrase then inserts this DNA back into the genome at a new site

LINEs copy themselves and insert into new genomic locations as follows

LINEs (long interspersed nuclear elements) - 20% of human genome

Most lines in the human genome are partially deleted and therefore no longer transpose

Some elements of the LINE1 family appear active

Hitchhike on the enzymatic activity of the LINEs

Transcribed by RNA polymerase III and use the LINE encoded enzymes to produce DNA and integrate themselves into the genome

SINES (short interspersed nuclear elements) - 13% of human genome

The most abundant SINES are Alu repeats (so called "cats" they have a recognition site for Alu)

These resemble retroviruses

They encode reverse transcriptase (pol), group specific antigen (gag) and envelope (env) proteins

LTR-like sequences

They are flanked by long terminal repeats (LTR)

So far, we've been seeing retrotransposons - those which originate from reverse transcriptase copied RNA (Or by a DNA based replicative mechanism)

A final class move by a DNA-based cut and paste mechanism

They use transposon-encoded transposase protein

No RNA intermediate

Perfect or imperfect repeats of short sequences occur - 3% of the human genome, occurring on average every 33 kb

One of the most common is the CA repeat, with n between 10 and 100

i.e.: different number of repeats

An individual is likely to have different alleles on each of a pair of chromosomes, allowing us to follow particular regions of the genome by the transmission of SSRs

Individuals are likely to have different spectrums of SSRs if a number of loci are compared

This means that

Many SSRs are polymorphic in the population, meaning that there are different alleles on each repeat locus

This is immensely useful in DNA fingerprinting

The repeats we have encountered so far are all apparently junk

Large tandem arrays, each repeat consist of the 3 rRNA found in the ribosome

Each repeat is separated by a "non-transcribed spacer"

rDNA genes are usually localized to the nucleolus

In some species, the transcription can be visualised as a Christmas tree-like structure

Recombination between arrays can lead to deletions and duplications

The organisation of repeats of DNA coding for histones is far more complex than that for DNA

It can also differ between species

Recombination between arrays can lead to deletions and duplications, and also to new arrangements of histone genes

More complex arrangements (e.g.: histones)

This has been one of the driving forces behind the generation of diversity

Genes are duplicated and then, freed from the constraints of selection, can evolve new functions

Both alpha and beta globin genes have been duplicated during evolution and have evolved slightly different functions

Include specialised genes expressed in the embryo and foetus

Remains of this process can be seen in the degenerate non-functional pseudo-genes "left-over" at both alpha and beta globin loci

Another kind of gene replication is interesting from an evolutionary perspective

Classical examples

The globin gene family

Actual genes

Spacer DNA (non-coding)

Pseudo-genes

These form a large portion of the 11% "misc." part of the human genome

Found in all metazoans (animals)

Control development

Regulate how sets of genes are turned on in different parts of the body

Hox gene clusters

All are clearly related, and it is believed that they have been maintained as a group to facilitate precise control of expression

Very little space left for these!

The actual number in a given organism is a thorny issue, depending on the quality of the sequence and its annotation

These divide fairly quickly (E. Coli > every 20 min under good conditions)

In unicellular organisms, gene numbers are fairly low and the genes are very densely packed in the genome

Their DNA therefore also replicated very quickly, and they have "eliminated waste" from their genome

However, even here there can be up to 5 fold differences in gene densities between related species (of rice and maize)

In metazoans, genes are more numerous and "space out"

In part, this is due to the fact the genes themselves are rather more complex

Almost all budding yeast genes have a single exon

83% of drosophila genes have at least one intron

5% of mammalian genes have over 30 introns

Some are very complex

If we look at the number of exons a gene contains...

We are limited to a set of experimental and comparative tools

Compare with processed mRNAs (obtained from cDNA libraries)

Compare with the primary sequence of proteins

Compare with information derived from simpler organisms

This makes the annotation of metazoan (and in particular mammalian) genes difficult

After DNA replication, chromosomes must be separated during mitosis

To achieve this, chromosomes are coordinated with the mitotic spindle via their centromeres

These are specialised structures on each chromosome that interact with a large number of proteins and the microtubules of the spindle to form a structure known as the kinetochore

If it is a 125bp long sequence (roughly 0.1% of chromosome) with an AT rich core

The best characterised centromeres are those of the budding yeast *Saccharomyces cerevisiae*

The core coordinates a set of protein complexes that secure the spindles to the chromosome

They're much longer (25kb to 110kb - up to 4% of the chromosome)

Those of the fission yeast *Schizosaccharomyces pombe* are less understood

They're composed of a complex arrangement of repeats flanking the inner core bordered by direct repeats

The reason for this dramatic difference is unknown

Since they are extremely repetitive, they haven't yet been fully sequenced

They appear to be composed of long tandem arrays of 171bp alpha satellite repeats

Most poorly characterised are mammalian centromeres

As with yeast, a complex array of proteins are assembled at the centromeres (a very large number of which have been identified)

A functional drosophila centromere was identified and isolated by virtue of its ability to confer mitotic stability of a marker gene

A 420kb drosophila centromere has been described at the sequence level

Sequence analysis reveals a complex organisation of simple sequence repeats interspersed with several different types of transposable elements

These are demarcated by types of transposable elements and their derivatives

Crudely, the chance that two unrelated people have an identical total DNA sequence is 1 in 6 million

Since the diploid genome contains 6 billion bp, this amounts to 6 million differences

This becomes more and more probably for closely related people

And is definite for identical twins

This allows us to use DNA to identify people

It would be a bit tough to sequence a person's entire DNA, though!

We can instead take advantage of polymorphisms associated with SSRs

We take specific primers that bind to unique sequences flanking a polymorphic SSR

We use those in a PCR reaction with the DNA to be identified

Since each chromosome in the pair will have a different allele (length) of SSR, each individual should produce 2 fragments in a PCR reaction

These can easily be resolved on a polyacrylamide gel (capable of base-pair resolution)

Two people compared are unlikely to have the same bands

However, the polymorphisms displayed at each polymorphic region are by themselves very common

Typically each polymorphism will be shared by around 5 - 20% of the population

However, if we look at multiple loci, the unique combination of several polymorphisms in an individual makes the method discriminating enough to be a robust identification tool

The national DNA database in the UK is the SCoP - DNA profiling system - it includes 10 different regions and a sex-specific test

The probability of an unrelated match is 1 in 1 billion

In eukaryotes, the genome is partitioned into chromosomes

Composed of a complex arrangement of DNA and proteins, including a set of functional elements necessary to ensure that DNA is efficiently replicated and partitioned into daughter cells

Replication is initiated at DNA replication origins

At the end of chromosomes, telomeres maintain the integrity of the chromosome

At cell division, the centromere and kinetochore facilitate chromosome segregation

\* Junk DNA is packaged into highly compacted chromatin (heterochromatin) and genes are found in less compact euchromatin

Cytogeneticists have utilised a variety of staining techniques to generate maps based on invariant banding patterns

G-bands are revealed by staining with Giemsa after treating with a protease to digest some of the proteins

Such darkly stained band have low G+C content

They also generally replicate late in the cycle

R-bands, however, are obtained by Giemsa staining after incorporation of the nucleotide analogue Bromodeoxyuridine

Dark bands here are G+C rich

And they usually replicate early

Even today, this cytological analysis can remain useful

For example, some cancers involve translocations or deletions, which modify banding patterns

The structures of some bacterial replication origins are well characterised and understood

The best characterised eukaryotic replication origins are yeast autonomously replicated sequences (ARS) which binds a complex set of ORC proteins (origin recognition complex) that act to control DNA replication

Very few metazoan replication origins have been defined

For example, in early embryos, each chromosome originally replicates from multiple origins

The situation is more complex than in yeast, though

Eventually, as the cell cycle gets longer, the number of origins used per chromosome declines

Similarly, different parts of the chromosome replicate at different times (C-banding analysis)

In eukaryotes, DNA exists as a linear molecule, and information is therefore lost from the end of each chromosome at each DNA replication

The ends of chromosomes are protected by special terminal structures known as telomeres

These consist of arrays of short G-rich repeats

They stabilise the ends of chromosomes by forming a complex T-loop structure that binds several telomere specific proteins

When telomeres become too short, cell death or senescence is triggered