BNF 5106 - Lecture 1 Genetics, Genes, Genetic codes, and Mutations







 Since each parent has 23 pairs of chromosomes, the probability that each parent gives twice the same chromosomes to two different children is (1/2)<sup>23</sup> x (1/2)<sup>23</sup>, or about 1 in 70 000 billion.



A meiosis can thus produce up to 4 gametes each having a recombinant chromosome. Non-recombined regions are called haplotypes; genes on a given haplotype are inherited together

Distribution of meiotic recombination events: talking to your neighbors

Enrique Martinez-Perez<sup>1</sup> and Monica P Colaiácovo<sup>2</sup> Current Opinion in Genetics & Development 2009, 19:105-112





+		
Chromosomes 1	Chromosomes 2	Chromosomes 3
-	Abnormally long haplotype	



## Genetic diseases

#### 1) No genetic disease (but carrying the recessive genes b and g as well as the lethal E gene)

A	b	С	D	E	F	G				
А	В	С	D	X	F	g				
2) Person with the genetic disease associated with gene b										
(and ca	arrying the	e recessiv	e gene g a	s well as t	he lethal l	Egene)				
А	b	С	D	E	F	G				
А	b	С	D	×	F	g				
2) 1 - 44	-1			F						
3) Leth	al mutati	on due to	defective	E gene						
А	b	С	D	×	F	G				
А	в	С	D	K	F	g				
				- C		-				







# Dominant mutations

- A single copy of the mutant gene causes the disease
  - E. g., Hungtiton disease (see below)







Gene = stretch of DNA with a function



- Double helix
- Anti-parallel strands
- C pairs with G (3 hydrogen bonds) and A pairs with T (2 hydrogen bonds)







- speed of replication:
  - ~ 500 bases/sec. in bacteria
  - ~ 50 bases/sec. in eukaryotes



- 100-200 nucleotides in eukaryotes
- Origin of replication
  - 1/genome (chromosome) in bacteria
  - 1/30 000 300 000 bp in eukaryotes





### 4

- Carried out by a single RNA polymerase in eubacteria
- Carried out by three different RNA polymerases in eukaryotes
  - RNA polymerase I = ribosomal RNA (rRNA) genes
  - RNA polymerase II = protein-coding genes
  - RNA polymerase III = small cytoplasmic RNA genes (e.g., 5S rRNA, tRNAs)





# Nucleotide sequences (DNA)

sugar (deoxyribose) + base

TABLE 1.1	One-letter abbreviations for the DNA alphabet				
Symbol	Description				
А	Adenine				
C	Cytosine				
т	Thymine				
G	Guanine				
w	Weak bonds (A, T)				
S	Strong bonds (C, G)				
R	Purines (A, G)				
Y	Pyrimidines (C, T)				
K	Keto (T, G)				
M	Amino (A, C)				
в	C, G, or T				
D	A, G, or T				
н	A, C, or T				
v	A, C, or G				
N	A, C, T, or G				
	No nucleotide (gap symbol)				



- Genome = the entire complement of genetic material carried by an individual
- Transcriptome = the entire set of transcribed sequences produced by the genome
- Proteome = the entire set of proteins encoded by the genome

# Genes and gene structure

- Gene = A sequence of DNA (RNA in some viruses) that is essential for a specific function.
  - Performing the function may not require the gene to be translated or transcribed.
  - protein-coding, RNA-coding, untranscribed genes (replicator genes, recombinator genes, telomeric sequences, segregator genes, etc.)

# Protein-coding genes

- GC and CAAT = binding of RNA polymerase
- TATA = start-point of transcription
- Antisense strand = DNA strand from which the RNA is transcribed
- Sense strand = untranscribed complementary strand (sequence identical to pre-mRNA)





- In vertebrates, intron sizes vary from  $\sim 100~{\rm bp}$  to 100s of Kbs; in contrast, the average exon size is 150 bp
- Human dystrophin gene -> 79 exons spanning over 2.3Mb, the mRNA is 12Kb long; exon/intron = 0.5%
- In other eukaryotes (e.g., yeast, Drosophila, C. elegans, plants), introns are fewer and shorter



# Pseudogenes

- Pseudogene = DNA segment that exhibits a high degree of similarity to a functional gene but which contains defects, such as nonsense and frameshift mutations, that prevent it to be expressed properly (most are not transcribed).
- Composed of processed (originated from the reverse transcription of a mRNA molecule) and unprocessed pseudogenes (originated from duplication of chromosome (DNA) fragments).



- Characteristics of processed pseudogenes:
  - 1) flanked by direct repeats
  - 2) intronless
  - 3) the downstream direct repeat is preceded by a poly-A tail





TABLE 1.3	The universal genetic code										
Codon	Amino acid	Codon	Amino acid	Codon	Amino acid	Codon	Amino acid				
UUU	Phe	UCU	Ser	UAU	Tyr	UGU	Cys				
UUC	Phe	UCC	Ser	UAC	Tyr	UGC	Cys				
UUA	Leu	UCA	Ser	UAA	Stop	UGA	Stop				
UUG	Leu	UCG	Ser	UAG	Stop	UGG	Trp				
CUU	Leu	CCU	Pro	CAU	His	CGU	Arg				
CUC	Leu	CCC	Pro	CAC	His	CGC	Arg				
CUA	Leu	CCA	Pro	CAA	Gln	CGA	Arg				
CUG	Leu	CCG	Pro	CAG	Gln	CGG	Arg				
AUU	Ile	ACU	Thr	AAU	Asn	AGU	Ser				
AUC	Ile	ACC	Thr	AAC	Asn	AGC	Ser				
AUA	lle	ACA	Thr	AAA	Lys	AGA	Arg				
AUG	Met	ACG	Thr	AAG	Lys	AGG	Arg				
GUU	Val	GCU	Ala	GAU	Asp	GGU	Gly				
GUC	Val	GCC	Ala	GAC	Asp	GGC	Gly				
GUA	Val	GCA	Ala	GAA	Glu	GGA	Gly				
GUG	Val	GCG	Ala	GAG	Glu	GGG	Gly				



# The universal genetic code

- Three stop codons: UAA, UAG, and UGA
- 61 sense codons for 20 amino acids
- Degenerate code: 18/20 amino acids are encoded by more than one codon
- The different codons specifying the same amino acid are called <u>synonymous codons</u>
- Some amino acids are encoded by a single codons (Met and Trp), by two codons (e.g., Phe), by three codons (Ile), by four codons (e.g., Val), by six codons (e.g., Leu)
- Synonymous codons often differ by transition (purine to purine, pyrimidine to pyrimidine) and not transversions (purine to pyrimidine and vice-versa)

TABLE 1.4	The verte	brate mitoch	ondrial gene	tic code*			
Codon	Amino acid	Codon	Amino acid	Codon	Amino acid	Codon	Amine acid
UUU	Phe	UCU	Ser	UAU	Tyr	UGU	Cys
UUC	Phe	UCC	Ser	UAC	Tyr	UGC	Cys
UUA	Leu	UCA	Ser	UAA	Stop	UGA	Trp
UUG	Leu	UCG	Ser	UAG	Stop	UGG	Trp
CUU	Leu	CCU	Pro	CAU	His	CGU	Arg
CUC	Leu	CCC	Pro	CAC	His	CGC	Arg
CUA	Leu	CCA	Pro	CAA	Gln	CGA	Arg
CUG	Leu	CCG	Pro	CAG	Gln	CGG	Arg
AUU	Ile	ACU	Thr	AAU	Asn	AGU	Ser
AUC	Ile	ACC	Thr	AAC	Asn	AGC	Ser
AUA	Met	ACA	Thr	AAA	Lys	AGA	Stop
AUG	Met	ACG	Thr	AAG	Lys	AGG	Stop
GUU	Val	GCU	Ala	GAU	Asp	GGU	Gly
GUC	Val	GCC	Ala	GAC	Asp	GGC	Gly
GUA	Val	GCA	Ala	GAA	Glu	GGA	Gly
GUG	Val	GCG	Ala	GAG	Glu	GGG	Gly

#### The vertebrate mitochondrial code



- 4 stop codons
- Two of the codons that specify serine in the universal genetic code are used as termination codons
- Tryptophan and methionine are each encoded by two codons rather than one
- Other exceptions to the "universal genetic code" (codon capture)
  - Mycoplasma uses UGA to code for tryptophan
  - Paramecium and Tetrahymena use UAA and UAG to code for glutamine
  - The yeast *Candida cylindracea* uses the codon CUG to code for serine



generat for a coden in that particular clade is the same as in the In the inference column, OGMP refers to the Organic colournel. The mitschardral encourse is such clade are bleed in



 Mostly due to errors of incorporation during either DNA replication or repair (low rate; maybe due to wrong tautomeric forms). Occur in both somatic and germline cells. Only germline mutations are considered in evolutionary studies.



- Can be classified by the type of change caused:
  - 1) **substitution**, the replacement of one nucleotide by another
  - 2) recombination, the exchange of a sequence by another
  - 3) deletions, the lost of one or more nucleotide
  - 4) **insertions**, the addition of one or more nucleotide
  - 5) inversions, the 180° degree rotation of two or more nucleotides



- Substitution = mutation that has not been eliminated by natural selection (or lost by random genetic drift)
- Synonymous substitution = substitution that does not change the amino acid (at 3rd base of codons and 1st base of Arg and Leu codons)
- Nonsynonymous subs. = substitution that change the amino acid (at all 2nd base of codons and 96% of 1st base of codons)

(a	) Ile	Cys	Ile	Lys	Ala	Leu	Val	Leu	Leu	Thr
	ATA	TGT	ATA	AAG	GCA	CTG	GTC	CTG	TTA	ACA
	ATA	TGT	ATA	AAG	GCA	CTG	GTÅ	CTG	TTA	ACA
	lle	Cys	Ile	Lys	Ala	Leu	Val	Leu	Leu	Thr
(b)	) lle	Cys	lle	Lys	Ala	Asn	Val	Leu	Leu	Thr
	ATA	TGT	ATA	AAG	GCA	AAC	стс 	CTG	TTA	ACA
	ATA	TGT	ATA	AAG	GCA	AAC	ттс	CTG	TTA	ACA
	Ile	Cys	Ile	Lys	Ala	Asn	Phe	Leu	Leu	Thr
(c)	Ile	Cys	Пе	Lys	Ala	Asn	Val	Leu	Leu	Thr
	ATA	TGT	ATA	AAG	GCA	AAC	GTC	СТG	TTA	ACA
	ATA	TGT	ATA	TAG	GCA/	AACG	гссто	GTTA	CA	
	Ile	Cys	Ile	Stop						

# Relative frequencies of different types of substitutions

- Each of the sense codons can mutate to nine other codons by means of a single nucleotide change. (e.g., CCC codon for Pro)
- 61 sense codon x 9 = 549 possible substitutions
- If we assume that these substitutions occur with equal frequency, and that all codons are equally frequent in coding regions, we can compute the expected proportion of the different types of substitutions from the genetic code



TABLE 1.5 Relative frequencies random protein-cod	s of different types of mutational substitutions in a ding sequence					
Substitution	Number	Percent				
Total in all codons	549	100				
Synonymous	134	25				
Nonsynonymous	415	75				
Missense	392	71				
Nonsense	23	4				
Total in first codons	183	100				
Synonymous	8	4				
Nonsynonymous	175	96				
Missense	166	91				
Nonsense	9	5				
Total in second codons	183	100				
Synonymous	0	0				
Nonsynonymous	183	100				
Missense	176	96				
Nonsense	7	4				
Total in third codons	183	100				
Synonymous	126	69				
Nonsynonymous	57	31				
Missense	50	27				
Nonconce	7	4				



- Synonymous substitutions occur almost exclusively at third base of codons (there is some in first base of Leu and Arg codons, but none in second base of codons)
- All substitutions at second base of codons are nonsynonymous
- Most nonsynonymous substitutions are missense substitutions

## Recombination

Two types of homologous recombination:

- 1) Crossingover (reciprocal recombination)
  - Even exchange of homologous sequences between homologous chromosomes
  - Both variants involved in the recombination events are retained but produces new combinations of adjacent sequences
- 2) Gene conversion (nonreciprocal recombination):
  - Uneven replacement of a sequence by another
  - Results in the loss of one of the variant sequence involved in the recombination event



Vertical cut = crossingover; Horizontal cut = gene conversion



By unequal crossingover







By replication slippage (splipped-strand mispairing) In DNA sequences that contain contiguous short repeats Normal pairing during DNA replication Slipped-st and r S-AATC TATA-5-AATCCTAGTATATA-ATATGTGCTTAA-5 ::::: 3'-TTAG G 3-TTAGGATCATATATGTGCTTAA-5 Replication continues inserting TA repeat unit Rep rues after at unit KG TATATAČAČGŽAŤT−3 TATATAČAČGŽAŤT−3 5'-AATC 5-AATCCTAGTATACACGAATT-3 3-TTAG ATAIGTGCTTAA-5 A A TC 3'-TTAGGATCATATATGTGCTTAA-5'



# Huntington disease, a domnant disease, is due to slippage mispairing (Cell 72 : 971-983, 1993)









disease	chromosome 1	sex bias of parent donating severe form	repeated sequence	normal number of copies	いのないのない
fragile X syndrome	X chromosome	maternal	CCG	6-50	fu

fragile X syndrome	X chromosome	maternal	CCG	6-50	premutation = 50-230 full mutation = 230-2,000
spinobulbar muscular atrophy (Kennedy disease)	X chromosome	?	AGC	11–31	40-62
myotonic dystrophy	chromosome 19	maternal	AGC	5-35	premutation = 50-80 full mutation = 80-2,000
Huntington disease	chromosome 4	paternal	CAG AGC	9–37	premutation = 30–38 full mutation = 37–121
spinocerebellar ataxia type 1	chromosome 6	paternal (possibly)	AGC	25-36	43-81
FRAXE	X chromosome	?	CCG	6-25	premutation = 25-200 full mutation = 200 and up
dentatorubral and pallidoluysian atrophy	chromosome 12	paternal (mainly)	AGC	7–23	49-754

Figure 4. Dynamic mutations are now linked with a growing number of human diseases. In addition to fragile X, the list includes myotonic dystrophy and Humington disease, as well as some lies common neurofogenerative and muncular disorders. The exact sequence of mutetidises and direct but the series always includes three mutedides. In addition, the mutation can expand when it is passed from parent to child.

American Scientist 82: 157-163

of copies ated with



### Mutation rates

- Average mutation rate in mammalian nuclear DNA = 3 to 5 x 10<sup>-9</sup> substitutions/site/year
  - Human genome = 3 x 10<sup>9</sup> sites, therefore 1 to 1.66 substitutions/year (in the germ line)
  - This rate varies depending on the genomic regions. E.g., the rate of indels of human microsatellites (tandemly repeated sequences 1 to 8 bases long) is > 10<sup>-3</sup> s/s/y
- In mammals, regions rich in CG dinucleotides (where the C is often methylated) are more prone to substitutions

4

- This rate is also different in different genomes
  - E.g., the rate of synonymous substitutions of the mammalian mitochondrial genome is at least 10 times greater than that of the mammalian nuclear DNA genome
- This rate is also different in different species
  - E. coli = 1 x 10<sup>-10</sup> s/s/y,
  - HIV = 1 x 10<sup>-3</sup> s/s/y

# Patterns of mutation

- The direction of substitutions is often not random.
  - If substitutions were random, we would expect to observe that 66% of the substitutions would be transversions and that 33% of the substitutions would be transitions (8 vrs 4)



# -

- In animal nuclear DNA, transitions have been found to account for 60-70% of all substitutions. Part of this excess is due to the genetics code (transitions are often synonymous substitutions) and to the deamination of methylated cytosines into thymine.
- In animal mitochondrial genomes, the ratio of transitions to transversions varies from 15 to 20.