

Example: HMMs for gene finding, part 1

Ben Langmead



For original Keynote files, email me (ben.langmead@gmail.com)

A human gene

chr11:5246500-5248500 (reverse strand):

```
ATATCTTAGAGGGAGGGCTGAGGGTTTGAAGTCCAACCTCCTAAGCCAGTGCCAGAAGAGCCAAGGACAGGTA
CGGCTGTCATCACTTAGACCTCACCTGTGGAGCCACACCCTAGGGTTGGCCAATCTACTCCCAGGAGCAGG
GAGGGCAGGAGCCAGGGCTGGGCATAAAAGTCAGGGCAGAGCCATCTATTGCTTACATTTGCTTCTGACACA
ACTGTGTTCACTAGCAACCTCAAACAGACACCATGGTGCATCTGACTCCTGAGGAGAAGTCTGCCGTTACTG
CCCTGTGGGGCAAGGTGAACGTGGATGAAGTTGGTGGTGAGGCCCTGGGCAGGTTGGTATCAAGGTTACAAG
ACAGGTTTAAGGAGACCAATAGAACTGGGCATGTGGAGACAGAGAAGACTCTTGGGTTTCTGATAGGCACT
GACTCTCTCTGCCTATTGGTCTATTTTCCCACCCTTAGGCTGCTGGTGGTCTACCCTTGGACCCAGAGGTTT
TTTGAGTCCTTTGGGGATCTGTCCACTCCTGATGCTGTTATGGGCAACCCTAAGGTGAAGGCTCATGGCAAG
AAAGTGCTCGGTGCCTTTAGTGATGGCCTGGCTCACCTGGACAACCTCAAGGGCACCTTTGCCACACTGAGT
GAGCTGCACTGTGACAAGCTGCACGTGGATCCTGAGAACTTCAGGGTGAGTCTATGGGACGCTTGATGTTTT
CTTTCCCCTTCTTTTCTATGGTTAAGTTCATGTCATAGGAAGGGGATAAGTAACAGGGTACAGTTTAGAATG
GGAAACAGACGAATGATTGCATCAGTGTGGAAGTCTCAGGATCGTTTTAGTTTTCTTTTATTTGCTGTTCATA
ACAATTGTTTTCTTTTGTTAATTCTTGCTTTCTTTTTTTTTCTTCTCCGCAATTTTTACTATTATACTTAA
TGCCTTAACATTGTGTATAACAAAAGGAAATATCTCTGAGATACATTAAGTAACTTAAAAAAAAAACTTTACA
CAGTCTGCCTAGTACATTACTATTTGGAATATATGTGTGCTTATTTGCATATTCATAATCTCCCTACTTTAT
TTTTTTTTATTTTTAATTGATACATAATCATTATACATATTTATGGGTTAAAGTGTAATGTTTTAATATGTG
TACACATATTGACCAAATCAGGGTAATTTTGCATTTGTAATTTTAAAAAATGCTTTCTTCTTTTAATATACT
TTTTTGTTTATCTTATTTCTAATACTTTCCCTAATCTCTTTCTTTCAGGGCAATAATGATACAATGTATCAT
GCCTCTTTGCACCATTCTAAAGAATAACAGTGATAATTTCTGGGTTAAGGCAATAGCAATATCTCTGCATAT
AAATATTTCTGCATATAAATTGTAACCTGATGTAAGAGGTTTCATATTGCTAATAGCAGCTACAATCCAGCTA
CCATTCTGCTTTTATTTTATGGTTGGGATAAGGCTGGATTATTCTGAGTCCAAGCTAGGCCCTTTTGCTAAT
CATGTTCATACTCTTATCTTCTCCACAGCTCCTGGGCAACGTGCTGGTCTGTGTGCTGGCCCATCACTT
TGGCAAAGAATTCACCCCACCAGTGCAGGCTGCCTATCAGAAAGTGGTGGCTGGTGTGGCTAATGCCCTGGC
CCACAAGTATCACTAAGCTCGCTTTCTTGCTGTCCAATTTCTATTAAAGGTTCTTTGTTCCCTAAGTCCAA
CTACTAAACTGGGGGATATTATGAAGGGCCTTGAGCATCTGGATTCTGCCTAATAAAAAACATTTATTTTCA
TTGCAATGATGTATTTAAATTATTTCTGAATATTTTACTAAAAAGGGAATGTGGGAGGTCAGTGCATTTAAA
ACATAAAGAAATGAAGAGCTAGTTCAAACCTTGGGAAAATACACTATATCTTAAACTCCATGAAAGAAGGTG
```

A human gene

chr11:5246500-5248500 (reverse strand):

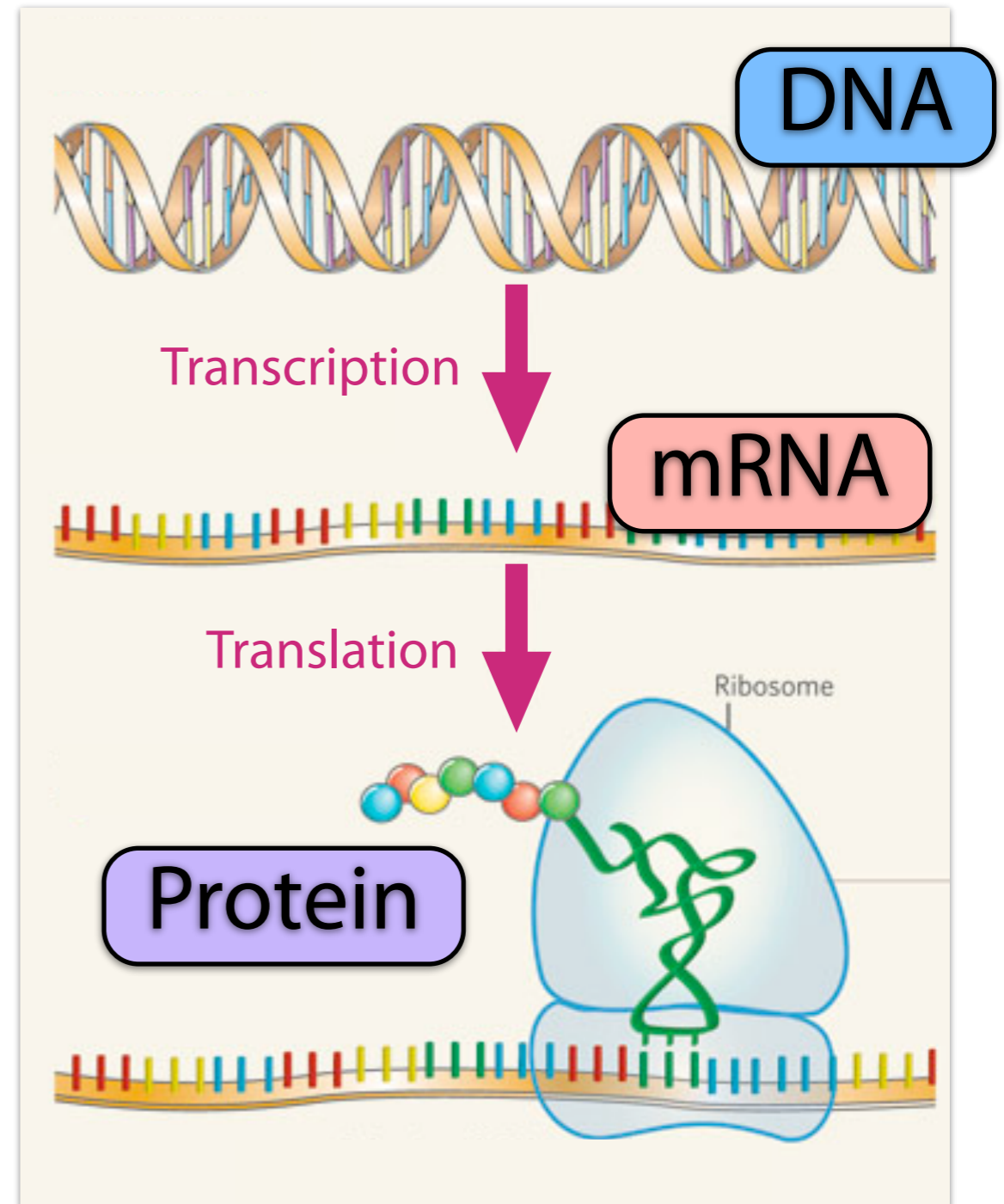
ATATCTTAGAGGGAGGGCTGAGGGTTTGAAGTCCAACCTCCTAAGCCAGTGCCAGAAGAGCCAAGGACAGGTA
CGGCTGTCATCACTTAGACCTCACCCCTGTGGAGCCACACCCTAGGGTTGGCCAATCTACTCCCAGGAGCAGG
GAGGGCAGGAGCCAGGGCTGGGCATAAAAGTCAGGGCAGAGCCATCTATTGCTTACATTTGCTTCTGACACA
ACTGTGTTCACTAGCAACCTCAAACAGACACC**ATGGTGCATCTGACTCCTGAGGAGAAGTCTGCCGTTACTG**
CCCTGTGGGGCAAGGTGAACGTGGATGAAGTTGGTGGTGGAGGCCCTGGGCAGGTTGGTATCAAGGTTACAAG
ACAGGTTTAAGGAGACCAATAGAACTGGGCATGTGGAGACAGAGAAGACTCTTGGGTTTCTGATAGGCACT
GACTCTCTCTGCCTATTGGTCTATTTTCCCACCCTTAG**GCTGCTGGTGGTCTACCCTTGGACCCAGAGGTTT**
TTTGAGTCCTTTGGGGATCTGTCCACTCCTGATGCTGTTATGGGCAACCCTAAGGTGAAGGCTCATGGCAAG
AAAGTGCTCGGTGCCTTTAGTGATGGCCTGGCTCACCTGGACAACCTCAAGGGCACCTTTGCCACACTGAGT
GAGCTGCACTGTGACAAGCTGCACGTGGATCCTGAGAACTTCAGGGTGAGTCTATGGGACGCTTGATGTTTT
CTTTCCCCTTCTTTTCTATGGTTAAGTTCATGTCATAGGAAGGGGATAAGTAACAGGGTACAGTTTAGAATG
GGAAACAGACGAATGATTGCATCAGTGTGGAAGTCTCAGGATCGTTTTAGTTTCTTTTATTTGCTGTTCATA
ACAATTGTTTTCTTTTGTTAATTCTTGCTTTCTTTTTTTTTCTTCTCCGCAATTTTACTATTATACTTAA
TGCCTTAACATTGTGTATAACAAAAGGAAATATCTCTGAGATACATTAAGTAACTTAAAAAAAAAACTTTACA
CAGTCTGCCTAGTACATTACTATTTGGAATATATGTGTGCTTATTTGCATATTCATAATCTCCCTACTTTAT
TTTCTTTTATTTTAAATTGATACATAATCATTATACATATTTATGGGTTAAAGTGTAATGTTTTAATATGTG
TACACATATTGACCAAATCAGGGTAATTTTGCATTTGTAATTTTAAAAAATGCTTTCTTCTTTAATACT
TTTTTGTTTATCTTATTTCTAATACTTTCCCTAATCTCTTTCTTTCAGGGCAATAATGATACAATGTATCAT
GCCTCTTTGCACCATTCTAAAGAATAACAGTGATAATTTCTGGGTTAAGGCAATAGCAATATCTCTGCATAT
AAATATTTCTGCATATAAATTGTAACCTGATGTAAGAGGTTTCATATTGCTAATAGCAGCTACAATCCAGCTA
CCATTCTGCTTTTATTTTATGGTTGGGATAAGGCTGGATTATTCTGAGTCCAAGCTAGGCCCTTTTGCTAAT
CATGTTCATACTCTTATCTTCTCCACAG**CTCCTGGGCAACGTGCTGGTCTGTGTGCTGGCCCATCACTT**
TGGCAAAGAATTCACCCCACCAGTGCAGGCTGCCTATCAGAAAGTGGTGGCTGGTGTGGCTAATGCCCTGGC
CCACAAGTATCACTAAGCTCGCTTTCTTGCTGTCCAATTTCTATTAAAGGTTCTTTGTTCCCTAAGTCCAA
CTACTAAACTGGGGGATATTATGAAGGGCCTTGAGCATCTGGATTCTGCCTAATAAAAAACATTTATTTTCA
TTGCAATGATGTATTTAAATTATTTCTGAATATTTTACTAAAAAGGGAATGTGGGAGGTCAGTGCATTTAAA
ACATAAAGAAATGAAGAGCTAGTTCAAACCTTGGGAAAATACACTATATCTTAAACTCCATGAAAGAAGGTG

Genes

"Central Dogma"

DNA molecules contain information about how to create proteins; this is *transcribed* into [messenger] **RNA** molecules, which, in turn, direct chemical machinery to *translate* the message into a **protein**.

Hunter, Lawrence. "Life and its molecules: A brief introduction." *AI Magazine* 25.1 (2004): 9.



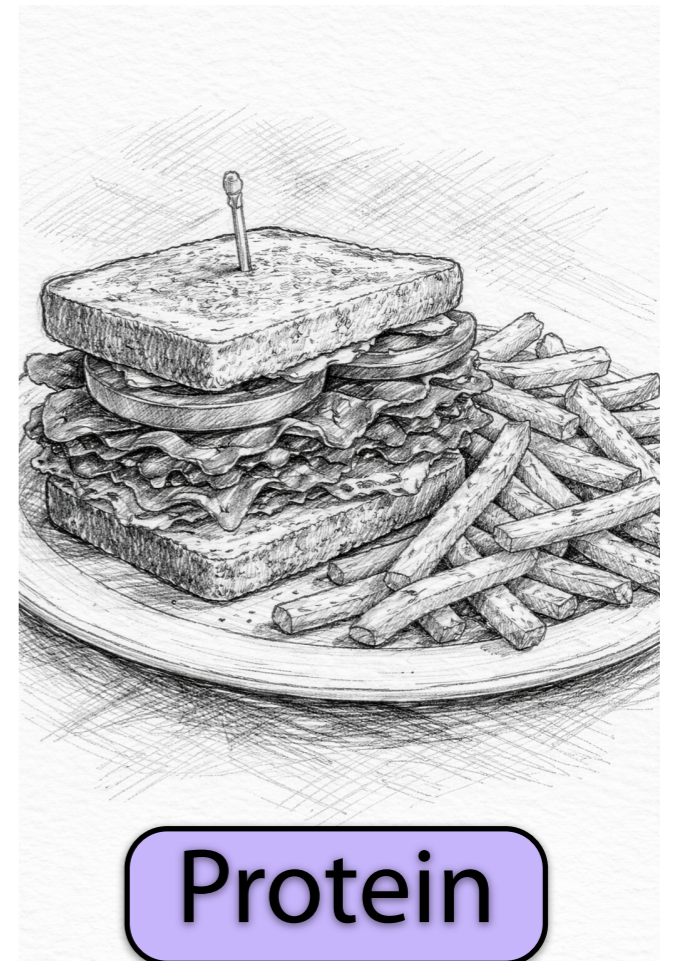
Picture from: Roy H, Ibba M. Molecular biology: sticky end in protein synthesis. *Nature*. 2006 Sep 7;443(7107):41-2.

Genes



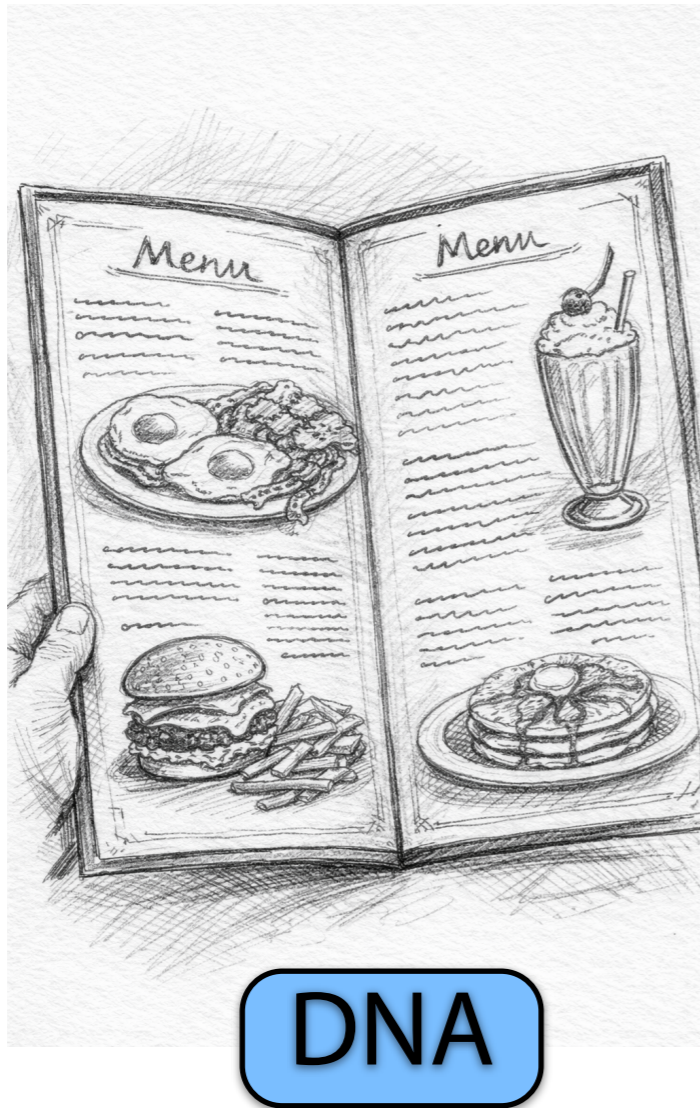
DNA

mRNA



Protein

Genes



Gene :: one item on the menu

Customizations are possible

“Without mayo”

“Add avocado”

etc

Gene finding



What are all
the **orders**
(genes and
their mRNA
versions) that
might come
from the
menu (DNA)



Gene finding

Answer given in Human Genome publications was around 26K—31K; estimate has decreased gradually since

Using this definition, though, do we have agreement on the number of protein-coding genes? The short answer is no. The human genome began with the assumption that our genome contains 100,000 genes. The number slightly decreased as more data became available. The initial human genome catalog was based on a limited set of genes, and when the more complete catalogs were released, it was found that a complete catalog of human genes and transcripts contained 34,214 transcripts.

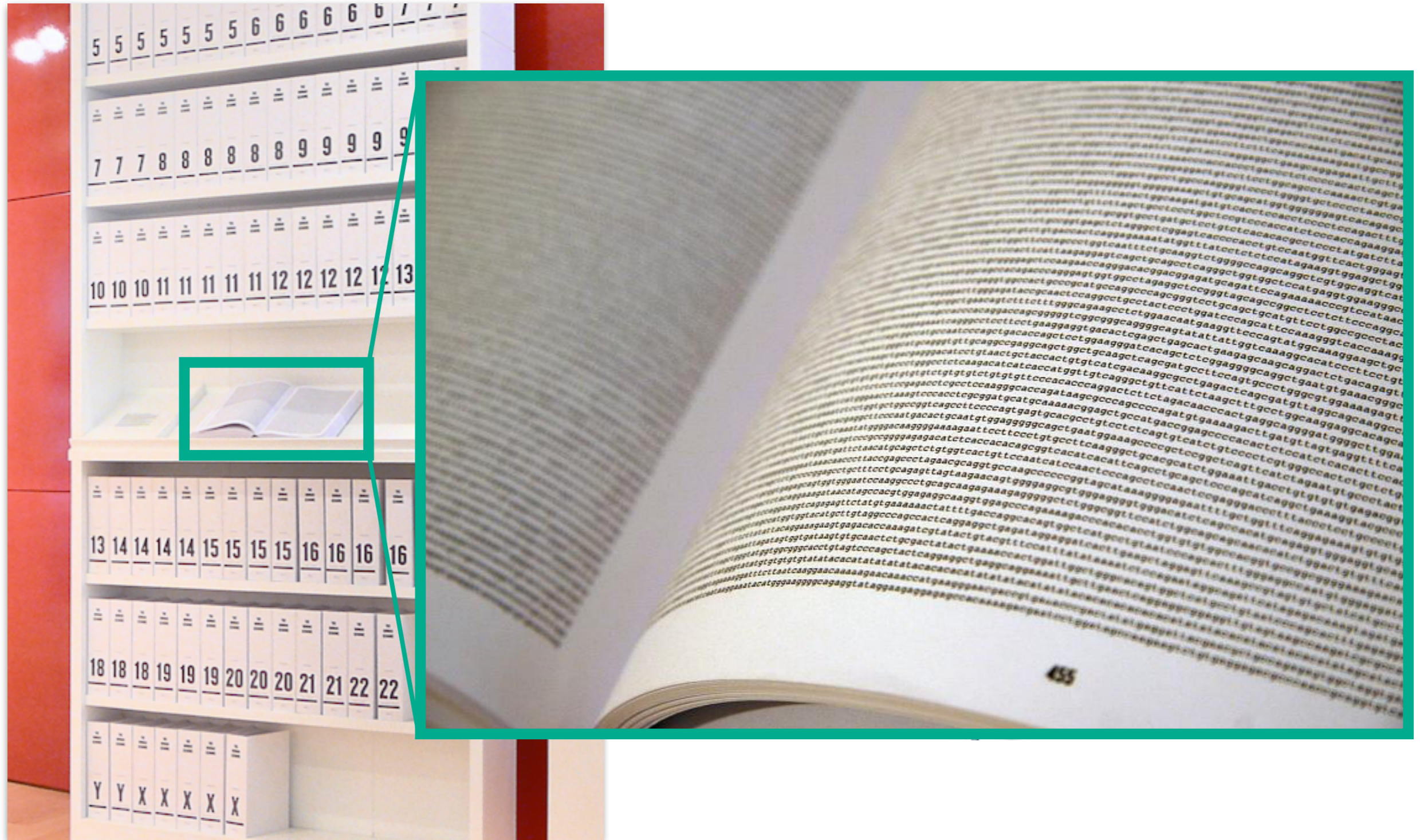
Abstract

CHESS 3 represents an improved human gene catalog based on nearly 10,000 RNA-seq experiments across 54 body sites. It significantly improves current genome annotation by integrating the latest reference data and algorithms, machine learning techniques for noise filtering, and new protein structure prediction methods. CHESS 3 contains 41,356 genes, including 19,839 protein-coding genes and 158,377 transcripts, with 14,863 protein-coding transcripts not in other catalogs. It includes all MANE transcripts and at least one transcript for most RefSeq and GENCODE genes. On the CHM13 human genome, the CHESS 3 catalog contains an additional 129 protein-coding genes. CHESS 3 is available at <http://ccb.jhu.edu/chess>.

Salzberg SL. Open questions: How many genes do we have? *BMC Biol.* 2018 Aug 20;16(1):94.

Varabyou A, Sommer MJ, Erdogdu B, Shinder I, Minkin I, Chao KH, Park S, Heinz J, Pockrandt C, Shumate A, Rincon N, Puiu D, Steinegger M, Salzberg SL & Pertea M 2023. CHESS 3: an improved, comprehensive catalog of human genes and transcripts based on large-scale expression data, phylogenetic analysis, and protein structure. *Genome biology*, 24(1), p.249.

Human genome



https://en.wikipedia.org/wiki/File:Wellcome_genome_bookcase.png

A human gene

chr11:5246500-5248500 (reverse strand):

```
ATATCTTAGAGGGAGGGCTGAGGGTTTGAAGTCCAACCTCCTAAGCCAGTGCCAGAAGAGCCAAGGACAGGTA
CGGCTGTCATCACTTAGACCTCACCTGTGGAGCCACACCCTAGGGTTGGCCAATCTACTCCCAGGAGCAGG
GAGGGCAGGAGCCAGGGCTGGGCATAAAAGTCAGGGCAGAGCCATCTATTGCTTACATTTGCTTCTGACACA
ACTGTGTTCACTAGCAACCTCAAACAGACACCATGGTGCATCTGACTCCTGAGGAGAAGTCTGCCGTTACTG
CCCTGTGGGGCAAGGTGAACGTGGATGAAGTTGGTGGTGGAGCCCTGGGCAGGTTGGTATCAAGGTTACAAG
ACAGGTTTAAGGAGACCAATAGAACTGGGCATGTGGAGACAGAGAAGACTCTTGGGTTTCTGATAGGCACT
GACTCTCTCTGCCTATTGGTCTATTTTCCCACCCTTAGGCTGCTGGTGGTCTACCCTTGGACCCAGAGGTTT
TTTGAGTCCTTTGGGGATCTGTCCACTCCTGATGCTGTTATGGGCAACCCTAAGGTGAAGGCTCATGGCAAG
AAAGTGCTCGGTGCCTTTAGTGATGGCCTGGCTCACCTGGACAACCTCAAGGGCACCTTTGCCACACTGAGT
GAGCTGCACTGTGACAAGCTGCACGTGGATCCTGAGAACTTCAGGGTGAGTCTATGGGACGCTTGATGTTTT
CTTTCCCCTTCTTTTCTATGGTTAAGTTCATGTCATAGGAAGGGGATAAGTAACAGGGTACAGTTTAGAATG
GGAAACAGACGAATGATTGCATCAGTGTGGAAGTCTCAGGATCGTTTTAGTTTTCTTTTATTTGCTGTTCATA
ACAATTGTTTTCTTTTGTTAATTCTTGCTTTCTTTTTTTTTCTTCTCCGCAATTTTTACTATTATACTTAA
TGCCTTAACATTGTGTATAACAAAAGGAAATATCTCTGAGATACATTAAGTAACTTAAAAAAAAAACTTTACA
CAGTCTGCCTAGTACATTACTATTTGGAATATATGTGTGCTTATTTGCATATTCATAATCTCCCTACTTTAT
TTTTTTTTATTTTTAATTGATACATAATCATTATACATATTTATGGGTTAAAGTGTAATGTTTTAATATGTG
TACACATATTGACCAAATCAGGGTAATTTTGCATTTGTAATTTTAAAAAATGCTTTCTTCTTTTAATACT
TTTTTGTTTATCTTATTTCTAATACTTTCCCTAATCTCTTTCTTTCAGGGCAATAATGATACAATGTATCAT
GCCTCTTTGCACCATTCTAAAGAATAACAGTGATAATTTCTGGGTTAAGGCAATAGCAATATCTCTGCATAT
AAATATTTCTGCATATAAATTGTAACCTGATGTAAGAGGTTTCATATTGCTAATAGCAGCTACAATCCAGCTA
CCATTCTGCTTTTATTTTATGGTTGGGATAAGGCTGGATTATTCTGAGTCCAAGCTAGGCCCTTTTGCTAAT
CATGTTCATACTCTTATCTTCTCCACAGCTCCTGGGCAACGTGCTGGTCTGTGTGCTGGCCCATCACTT
TGGCAAAGAATTCACCCCACCAGTGCAGGCTGCCTATCAGAAAGTGGTGGCTGGTGTGGCTAATGCCCTGGC
CCACAAGTATCACTAAGCTCGCTTTCTTGCTGTCCAATTTCTATTAAAGGTTCTTTGTTCCCTAAGTCCAA
CTACTAAACTGGGGGATATTATGAAGGGCCTTGAGCATCTGGATTCTGCCTAATAAAAAACATTTATTTTCA
TTGCAATGATGTATTTAAATTATTTCTGAATATTTTACTAAAAAGGGAATGTGGGAGGTCAGTGCATTTAAA
ACATAAAGAAATGAAGAGCTAGTTCAAACCTTGGGAAAATACACTATATCTTAAACTCCATGAAAGAAGGTG
```

A human gene

chr11:5246500-5248500 (reverse strand):

ATATCTTAGAGGGAGGGCTGAGGGTTTGAAGTCCAACCTCCTAAGCCAGTGCCAGAAGAGCCAAGGACAGGTA
CGGCTGTCATCACTTAGACCTCACCCCTGTGGAGCCACACCCTAGGGTTGGCCAATCTACTCCCAGGAGCAGG
GAGGGCAGGAGCCAGGGCTGGGCATAAAAAGTCAGGGCAGAGCCATCTATTGCTTACATTTGCTTCTGACACA
ACTGTGTTCACTAGCAACCTCAAACAGACACC**ATGGTGCATCTGACTCCTGAGGAGAAGTCTGCCGTTACTG**
CCCTGTGGGGCAAGGTGAACGTGGATGAAGTTGGTGGTGGAGGCCCTGGGCAGGTTGGTATCAAGGTTACAAG
ACAGGTTTAAGGAGACCAATAGAACTGGGCATGTGGAGACAGAGAAGACTCTTGGGTTTCTGATAGGCACT
GACTCTCTCTGCCTATTGGTCTATTTTCCCACCCTTAG**GCTGCTGGTGGTCTACCCTTGGACCCAGAGGTTT**
TTTGAGTCCTTTGGGGATCTGTCCACTCCTGATGCTGTTATGGGCAACCCTAAGGTGAAGGCTCATGGCAAG
AAAGTGCTCGGTGCCTTTAGTGATGGCCTGGCTCACCTGGACAACCTCAAGGGCACCTTTGCCACACTGAGT
GAGCTGCACTGTGACAAGCTGCACGTGGATCCTGAGAACTTCAGGGTGAGTCTATGGGACGCTTGATGTTTT
CTTTCCCCTTCTTTTCTATGGTTAAGTTCATGTCATAGGAAGGGGATAAGTAACAGGGTACAGTTTAGAATG
GGAAACAGACGAATGATTGCATCAGTGTGGAAGTCTCAGGATCGTTTTAGTTTCTTTTATTTGCTGTTCATA
ACAATTGTTTTCTTTTGTTTAATTCTTGCTTTCTTTTTTTTTCTTCTCCGCAATTTTACTATTATACTTAA
TGCCTTAACATT**GTGTATAACAAAAGGAAATATCTCTGAGATACATTAAGTAACTTAAAAAAA**ACTTTACA
CAGTCTGCCTAG**TTCTCCCTACTTTAT**
TTTCTTTTATTT**TGTTTTAATATGTG**
TACACATATTGACCAAATCAGGGTAATTTTGCATTTGTAATTTTAAAAAATGCTTTCTTCTTTTAATATACT
TTTTTGTTTATCTTATTTCTAATACTTTCCCTAATCTTTTCTTTCAGGGCAATAATGATACAATGTATCAT
GCCTCTTTGCACCATTCTAAAGAATAACAGTGATAATTTCTGGGTTAAGGCAATAGCAATATCTCTGCATAT
AAATATTTCTGCATATAAATTGTAACCTGATGTAAGAGGTTTCATATTGCTAATAGCAGCTACAATCCAGCTA
CCATTCTGCTTTTATTTTATGGTTGGGATAAGGCTGGATTATTCTGAGTCCAAGCTAGGCCCTTTTGCTAAT
CATGTTCATACTCTTATCTTCTCCACAG**CTCCTGGGCAACGTGCTGGTCTGTGTGCTGGCCCATCACTT**
TGGCAAAGAATTCACCCACCAGTGCAGGCTGCCTATCAGAAAGTGGTGGCTGGTGTGGCTAATGCCCTGGC
CCACAAGTATCACTAAGCTCGCTTTCTTGCTGTCCAATTTCTATTAAAGGTTCTTTGTTCCCTAAGTCCAA
CTACTAAACTGGGGGATATTATGAAGGGCCTTGAGCATCTGGATTCTGCCTAATAAAAAACATTTATTTTCA
TTGCAATGATGTATTTAAATTATTTCTGAATATTTTACTAAAAAGGGAATGTGGGAGGTCAGTGCATTTAAA
ACATAAAGAAATGAAGAGCTAGTTCAAACCTTGGGAAAATACACTATATCTTAAACTCCATGAAAGAAGGTG

Homo sapiens hemoglobin, beta (HBB)

Genes

Do genes “pop out” of the genomic “background”? 🤔 🤔

Genes

Do genes “pop out” of the genomic “background”? 🤔 🤔

Nucleotides triples
("codons") are translated
into amino acids via the
genetic code

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

Note: Uracil (U) is the RNA
version of Thymine (T)

Genes

Do genes “pop out” of the genomic “background”? 🤔 🤔

Nucleotides triples
("codons") are translated
into amino acids via the
genetic code

CCA → Proline (P)

GGU → Glycine (G)

AAA → Lysine (K)

(etc)

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

Note: Uracil (U) is the RNA
version of Thymine (T)

Genes

Do genes “pop out” of the genomic “background”? 🤔 🤔

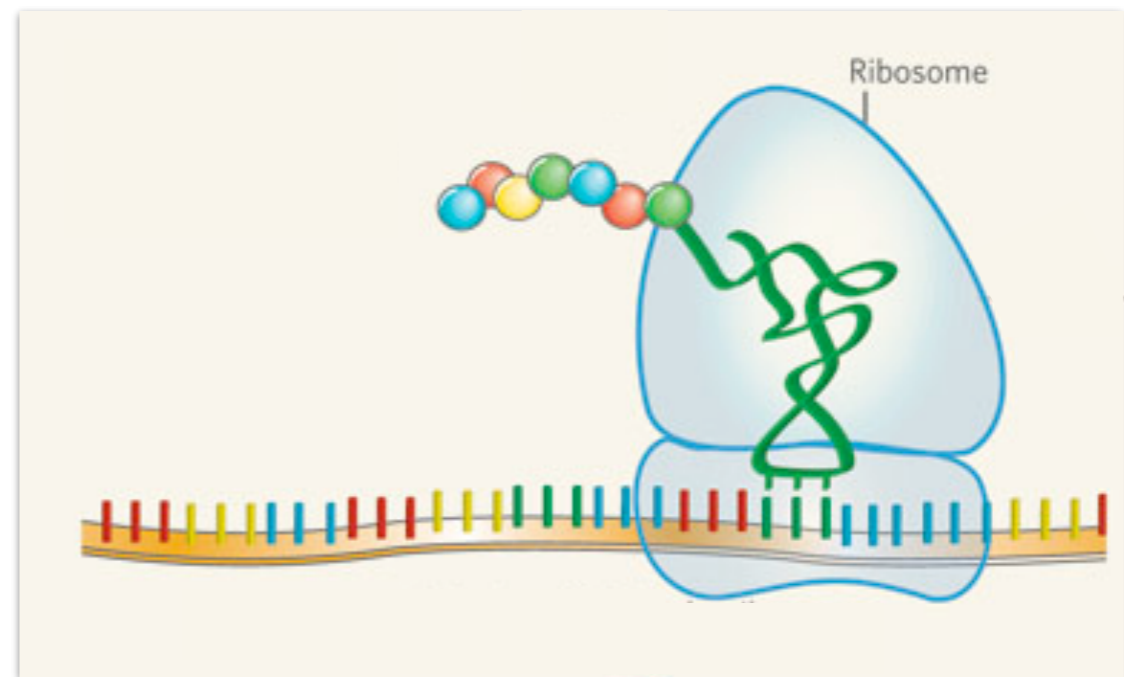
Nucleotides triples
("codons") are translated
into amino acids via the
genetic code

CCA → Proline (P)

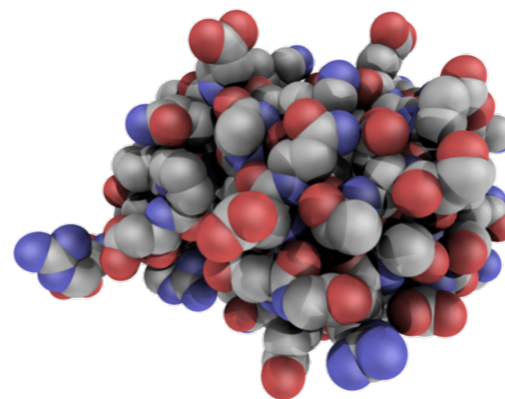
GGU → Glycine (G)

AAA → Lysine (K)

(etc)



Picture from: Roy H, Ibba M. Molecular biology: sticky end in protein synthesis. Nature. 2006 Sep 7;443(7107):41-2.



https://commons.wikimedia.org/wiki/File:Ubiquitin_spheres.png

Genes

Some codons are special

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

Genes

Some codons are special

Stop codons signal for translation to stop

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

Genes

Some codons are special

Stop codons signal for translation to stop

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

Start codon (AUG) signals translation should begin (& codes for Methionine)

Genes

“Splicing” ✂ is a process by which some portions of the mRNA are cut out prior to translation

ATATCTTAGAGGGAGGGCTGAGGGTTTGAAGTCCAACCTCCTAAGCCAGTGCCAGAAGAGCCAAGGACAGGTACGGCTGTCATCACTTAGACCTCACCC
TGTGGAGCCACACCCTAGGGTTGGCCAATCTACTCCCAGGAGCAGGGAGGGCAGGAGCCAGGGCTGGGCATAAAAGTCAGGGCAGAGCCATCTATTGC
TTACATTTGCTTCTGACACAACCTGTGTTCACTAGCAACCTCAAACAGACACC**ATGGTGCATCTGACTCCTGAGGAGAAGTCTGCCGTTACTGCCCTGT**
GGGGCAAGGTGAACGTGGATGAAGTTGGTGGTGGAGGCCCTGGGCAGGTTGGTATCAAGGTTACAAGACAGGTTTAAGGAGACCAATAGAAACTGGGCA
TGTGGAGACAGAGAAGACTCTTGGGTTTCTGATAGGCACTGACTCTCTCTGCCTATTGGTCTATTTTCCCACCCTTAG**GCTGCTGGTGGTCTACCCTT**
GGACCCAGAGGTTCTTTGAGTCCTTTGGGGATCTGTCCACTCCTGATGCTGTTATGGGCAACCCTAAGGTGAAGGCTCATGGCAAGAAAGTGCTCGGT
GCCTTTAGTGATGGCCTGGCTCACCTGGACAACCTCAAGGGCACCTTTGCCACACTGAGTGAGCTGCACTGTGACAAGCTGCACGTGGATCCTGAGAA
CTTCAGGGTGAGTCTATGGGACGCTTGATGTTTTCTTTCCCCTTCTTTTCTATGGTTAAGTTCATGTCATAGGAAGGGGATAAGTAACAGGGTACAGT
TTAGAATGGGAAACAGACGAATGATTGCATCAGTGTGGAAGTCTCAGGATCGTTTTAGTTTTCTTTTATTTGCTGTTTATAACAATTGTTTTCTTTTGT
TTAATTCTTGCTTTCTTTTTTTTTCTTCTCCGCAATTTTTACTATTATACTTAATGCCTTAACATTGTGTATAACAAAAGGAAATATCTCTGAGATAC
ATTAAGTAACTTAAAAAAAAAACTTTACACAGTCTGCCTAGTACATTACTATTTGGAATATATGTGTGCTTATTTGCATATTCATAATCTCCCTACTTT
ATTTTCTTTTATTTTAAATTGATACATAATCATTATACATATTTATGGGTAAAGTGTAATGTTTTAATATGTGTACACATATTGACCAAATCAGGGT
AATTTTGCATTTGTAATTTTAAAAAATGCTTTCTTCTTTTAAATATACTTTTTTGTTTATCTTATTTCTAATACTTTCCCTAATCTCTTTCTTTTCAGGG
CAATAATGATACAATGTATCATGCCTCTTTGCACCATTCTAAAGAATAACAGTGATAATTTCTGGGTTAAGGCAATAGCAATATCTCTGCATATAAAT
ATTTCTGCATATAAATTGTAAGTACTGATGTAAGAGGTTTCATATTGCTAATAGCAGCTACAATCCAGCTACCATTCTGCTTTTATTTTATGGTTGGGATA
AGGCTGGATTATTCTGAGTCCAAGCTAGGCCCTTTTGCTAATCATGTTTACATACCTCTTATCTTCTCCACAG**CTCCTGGGCAACGTGCTGGTCTGTG**
TGCTGGCCCATCACTTTGGCAAAGAATTCACCCACCAGTGCAGGCTGCCTATCAGAAAGTGGTGGCTGGTGTGGCTAATGCCCTGGCCCACAAGTAT
CACTAAGCTCGCTTTCTTGCTGTCCAATTTCTATTAAAGGTTCTTTGTTCCCTAAGTCCAACACTAAACTGGGGGATATTATGAAGGGCCTTGAGC
ATCTGGATTCTGCCTAATAAAAAACATTTATTTTCATTGCAATGATGTATTTAAATTATTTCTGAATATTTTACTAAAAAGGGAATGTGGGAGGTCAG
TGCATTTAAACATAAAGAAATGAAGAGCTAGTTCAAACCTTGGGAAAATACACTATATCTTAAACTCCATGAAAGAAGGTGAGGCTGCAAACAGCTA
ATGCACATTGGCAACAGCCCCTGATGCATATGCCTTATTC

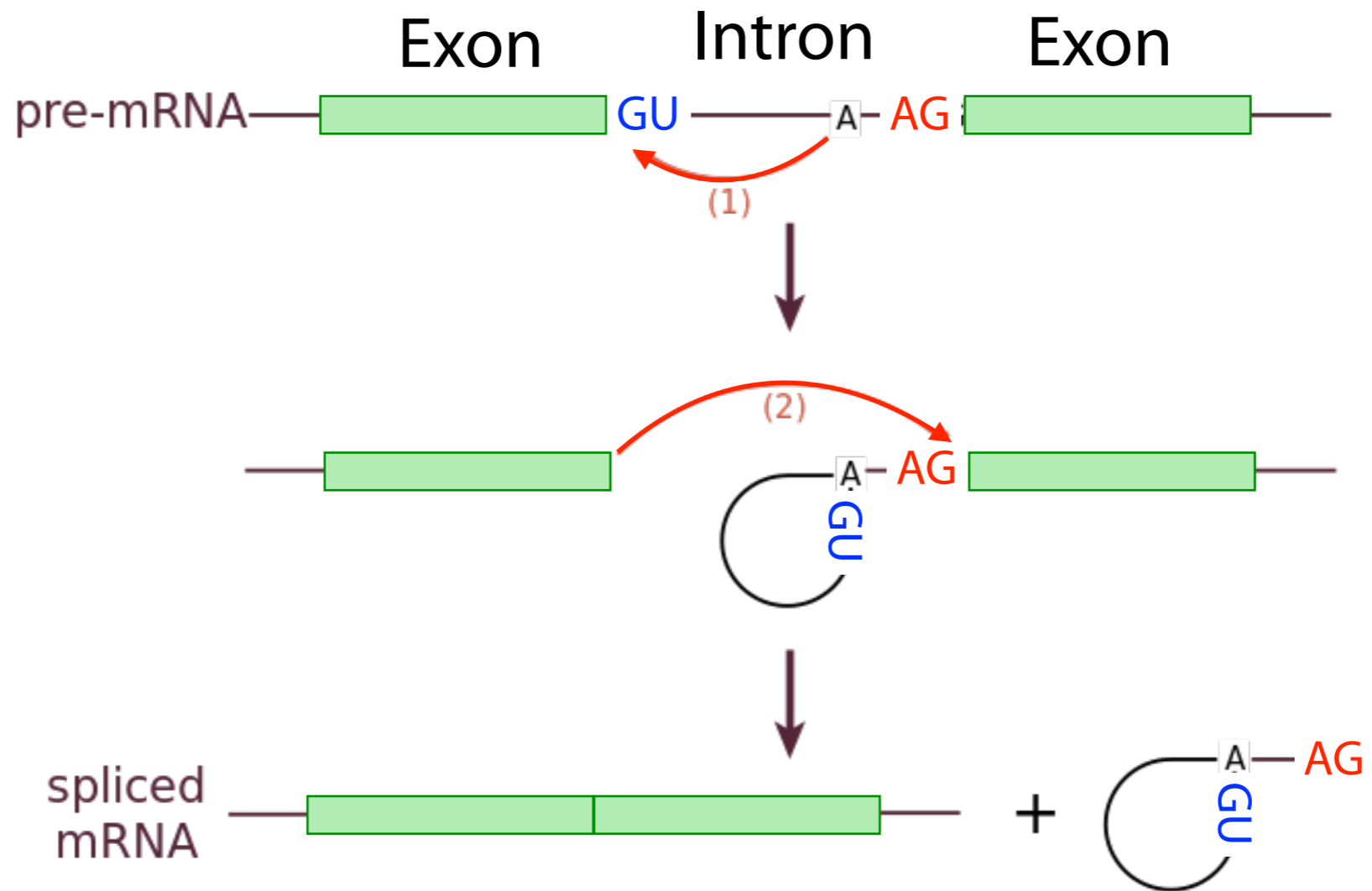
Genes

“Splicing” ✂ is a process by which some portions of the mRNA are cut out prior to translation

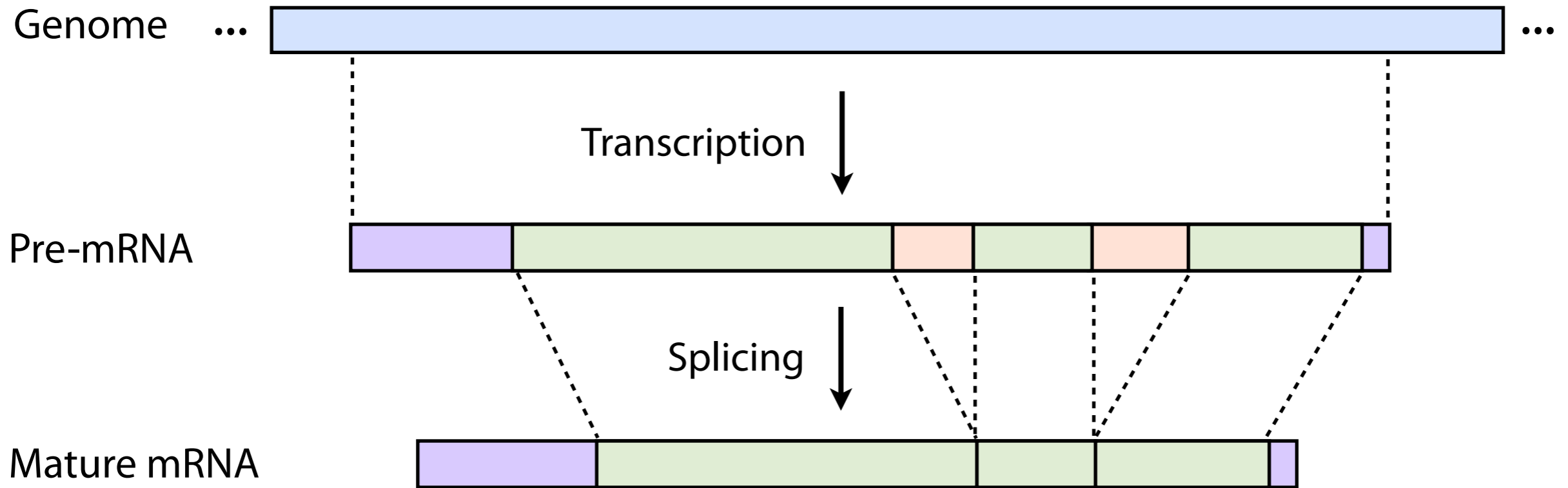
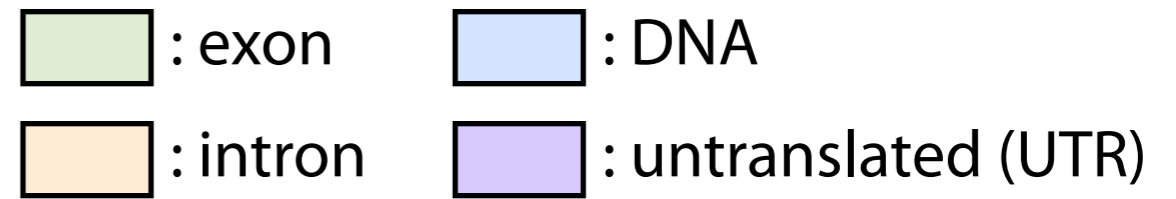
```
ATATCTTAGAGGGAGGGCTGAGGGTTTGAAGTCCAACCTCCTAAGCCAGTGCCAGAAGAGCCAAGGACAGGTACGGCTGTCATCACTTAGACCTCACCC
TGTGGAGCCACACCCTAGGGTTGGCCAATCTACTCCCAGGAGCAGGGAGGGCAGGAGCCAGGGCTGGGCATAAAAAGTCAGGGCAGAGCCATCTATTGC
TTACATTTGCTTCTGACACAACCTGTGTTCACTAGCAACCTCAAACAGACACCATGGTGCATCTGACTCCTGAGGAGAAGTCTGCCGTTACTGCCCTGT
GGGGCAAGGTGAACGTGGATGAAGTTGGTGGTGGAGGCCCTGGGCAGGTTGGTATCAAGGTTACAAGACAGGTTTAAAGGAGACCAATAGAAACTGGGCA
TGTGGAGACAGAGAAGACTCTTGGGTTTCTGATAGGCACTGACTCTCTCTGCCTATTGGTCTATTTTCCCACCCTTAGGCTGCTGGTGGTCTACCCTT
GGACCCAGAGGTTCTTTGAGTCCTTTGGGGATCTGTCCACTCCTGATGCTGTTATGGGCAACCCTAAGGTGAAGGCTCATGGCAAGAAAGTGCTCGGT
GCCTTTAGTGATGGCCTGGCTCACCTGGACAACCTCAAGGGCACCTTTGCCCACTGAGTGAGCTGCACTGTGACAAGCTGCACGTGGATCCTGAGAA
CTTCAGGGTGAGTCTATGGGACGCTTGATGTTTTCTTTCCCCTTCTTTTCTATGGTTAAGTTCATGTCATAGGAAGGGGATAAGTAACAGGGGTACAGT
TTAGAATGGGAAACAGACGAATGATTGCATCAGTGTGGAAGTCTCAGGATCGTTTTAGTTTTCTTTTATTTGCTGTTTATAACAATTGTTTTCTTTTGT
TTAATTCTTGCTTTCTTTTTTTTTCTTCTCCGCAATTTTTACTATTATACTTAATGCCTTAACATTGTGTATAACAAAAGGAAATATCTCTGAGATAC
ATTAAGTAACTTAAAAAAAAAACTTTACACAGTCTGCCTAGTACATTACTATTTGGAATATATGTGTGCTTATTTGCATATTCATAATCTCCCTACTTT
ATTTTCTTTTATTTTAAATTGATACATAATCATTATACATATTTATGGGTAAAGTGTAATGTTTTAATATGTGTACACATATTGACCAAATCAGGGT
AATTTTGCATTTGTAATTTTAAAAAATGCTTTCTTCTTTTAAATATACTTTTTTGTTTATCTTATTTCTAATACTTTCCCTAATCTCTTTCTTTTCAAGG
CAATAATGATACAATGTATCATGCCTCTTTGCACCATTCTAAAGAATAACAGTGATAATTTCTGGGTAAAGGCAATAGCAATATCTCTGCATATAAAT
ATTTCTGCATATAAATTGTAACCTGATGTAAGAGGTTTCATATTGCTAATAGCAGCTACAATCCAGCTACCATTCTGCTTTTATTTTATGGTTGGGATA
AGGCTGGATTATTCTGAGTCCAAGCTAGGCCCTTTTGCTAATCATGTTTACATACCTCTTATCTTCTCCACAGCTCCTGGGCAACGTGCTGGTCTGTG
TGCTGGCCCATCACTTTGGCAAAGAATTCACCCACCAGTGCAAGGCTGCCTATCAGAAAGTGGTGGCTGGTGTGGCTAATGCCCTGGCCCACAAGTAT
CACTAAGCTCGCTTTCTTGCTGTCCAATTTCTATTAAAGGTTCTTTGTTCCCTAAGTCCAACCTACTAAACTGGGGGATATTATGAAGGGCCTTGAGC
ATCTGGATTCTGCCTAATAAAAAACATTTATTTTCATTGCAATGATGTATTTAAATTATTTCTGAATATTTTACTAAAAAGGGAATGTGGGAGGTCAG
TGCATTTAAACATAAAGAAATGAAGAGCTAGTTCAAACCTTGGGAAAATACACTATATCTTAAACTCCATGAAAGAAGGTGAGGCTGCAAACAGCTA
ATGCACATTGGCAACAGCCCCTGATGCATATGCCTTATTC
```

Genes

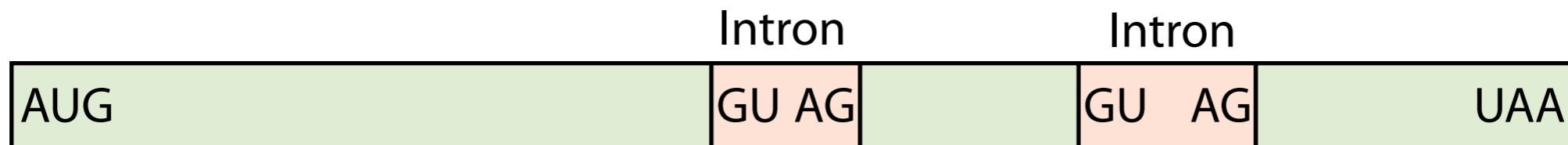
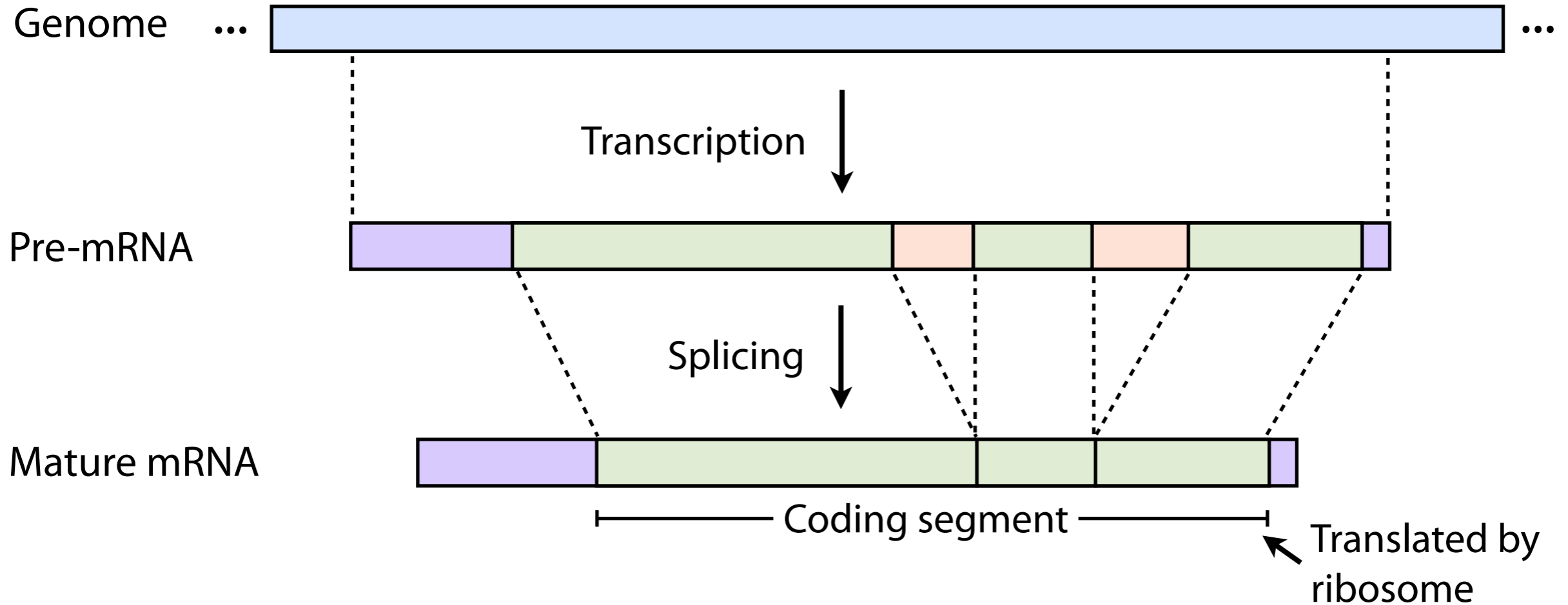
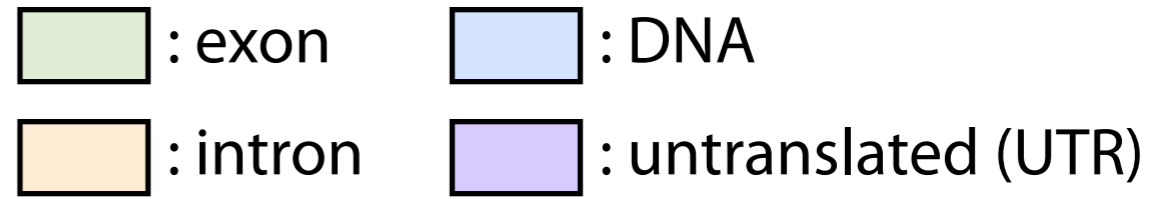
Cutting happens at certain nucleotide patterns: **GU** and **AG**



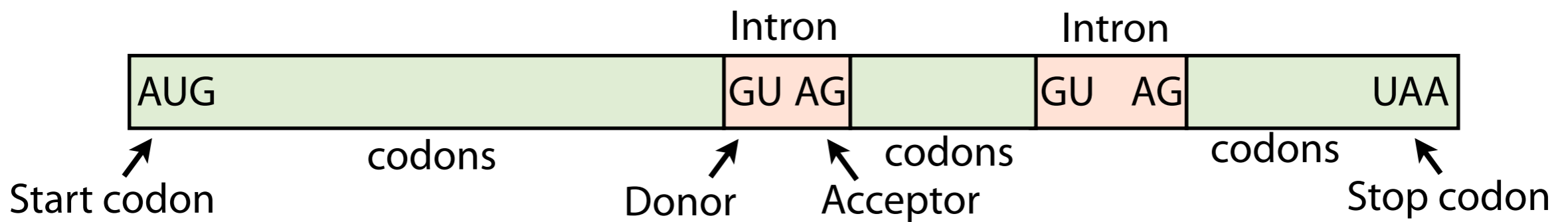
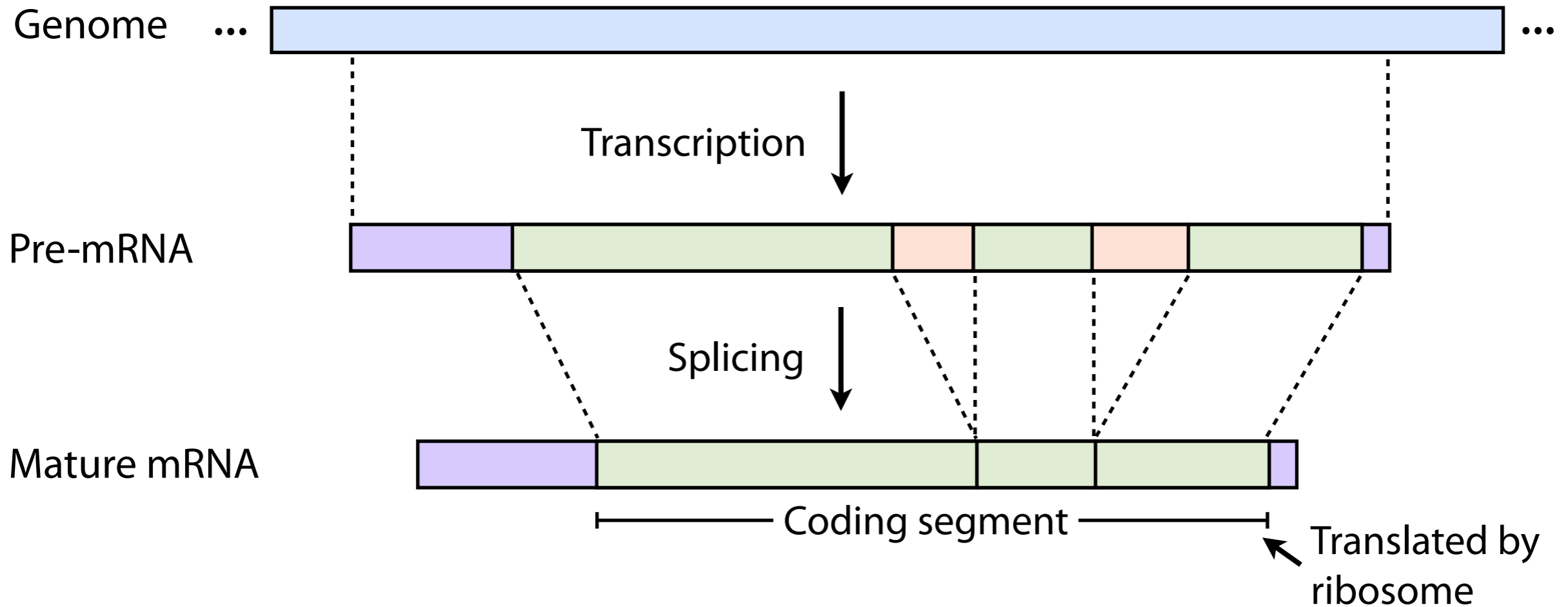
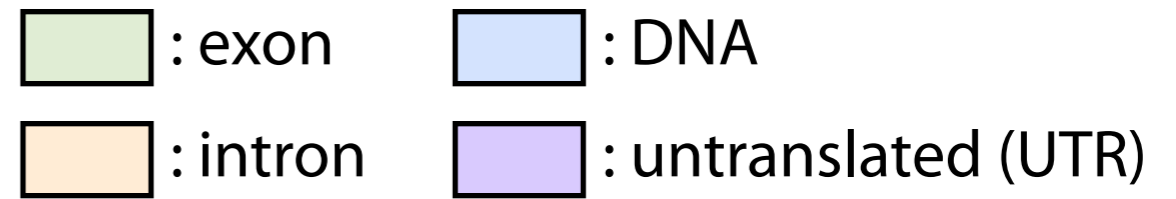
Transcription



Transcription



Gene signals



A human gene

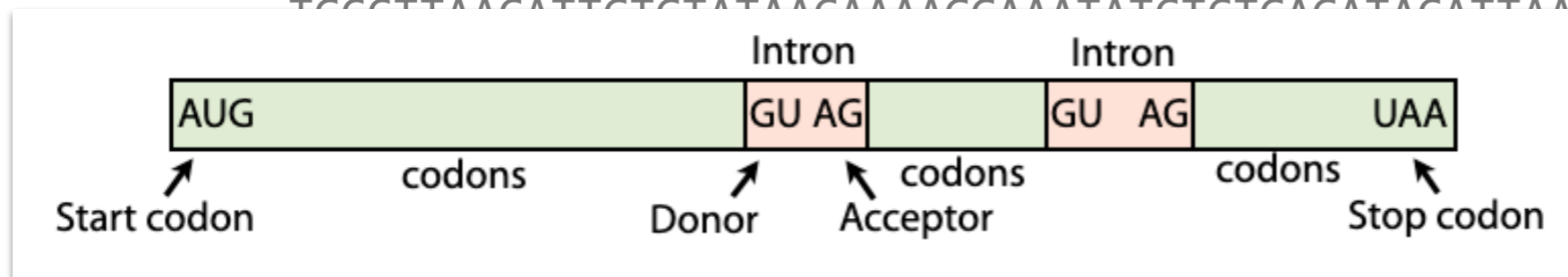
chr11:5246500-5248500 (reverse strand):

ATATCTTAGAGGGAGGGCTGAGGGTTTGAAGTCCAACCTCCTAAGCCAGTGCCAGAAGAGCCAAGGACAGGTA
CGGCTGTCATCACTTAGACCTCACCCCTGTGGAGCCACACCCTAGGGTTGGCCAATCTACTCCCAGGAGCAGG
GAGGGCAGGAGCCAGGGCTGGGCATAAAAAGTCAGGGCAGAGCCATCTATTGCTTACATTTGCTTCTGACACA
ACTGTGTTCACTAGCAACCTCAAACAGACACC**ATGGTGCATCTGACTCCTGAGGAGAAGTCTGCCGTTACTG**
CCCTGTGGGGCAAGGTGAACGTGGATGAAGTTGGTGGTGGAGGCCCTGGGCAGGTTGGTATCAAGGTTACAAG
ACAGGTTTAAGGAGACCAATAGAACTGGGCATGTGGAGACAGAGAAGACTCTTGGGTTTCTGATAGGCACT
GACTCTCTCTGCCTATTGGTCTATTTTCCCACCCTTAG**GCTGCTGGTGGTCTACCCTTGGACCCAGAGGTTT**
TTTGAGTCCTTTGGGGATCTGTCCACTCCTGATGCTGTTATGGGCAACCCTAAGGTGAAGGCTCATGGCAAG
AAAGTGCTCGGTGCCTTTAGTGATGGCCTGGCTCACCTGGACAACCTCAAGGGCACCTTTGCCACACTGAGT
GAGCTGCACTGTGACAAGCTGCACGTGGATCCTGAGAACTTCAGGGTGAGTCTATGGGACGCTTGATGTTTT
CTTTCCCCTTCTTTTCTATGGTTAAGTTCATGTCATAGGAAGGGGATAAGTAACAGGGTACAGTTTAGAATG
GGAAACAGACGAATGATTGCATCAGTGTGGAAGTCTCAGGATCGTTTTAGTTTCTTTTATTTGCTGTTCATA
ACAATTGTTTTCTTTTGTTAATTCTTGCTTTCTTTTTTTTTCTTCTCCGCAATTTTACTATTATACTTAA
TGCCTTAACATTGTGTATAACAAAAGGAAATATCTCTGAGATACATTAAGTAACTTAAAAAAAAAACTTTACA
CAGTCTGCCTAGTACATTACTATTTGGAATATATGTGTGCTTATTTGCATATTCATAATCTCCCTACTTTAT
TTTCTTTTATTTTAAATTGATACATAATCATTATACATATTTATGGGTTAAAGTGTAATGTTTTAATATGTG
TACACATATTGACCAAATCAGGGTAATTTTGCATTTGTAATTTTAAAAAATGCTTTCTTCTTTAATACT
TTTTTGTTTATCTTATTTCTAATACTTTCCCTAATCTCTTTCTTTCAGGGCAATAATGATACAATGTATCAT
GCCTCTTTGCACCATTCTAAAGAATAACAGTGATAATTTCTGGGTTAAGGCAATAGCAATATCTCTGCATAT
AAATATTTCTGCATATAAATTGTAACCTGATGTAAGAGGTTTCATATTGCTAATAGCAGCTACAATCCAGCTA
CCATTCTGCTTTTATTTTATGGTTGGGATAAGGCTGGATTATTCTGAGTCCAAGCTAGGCCCTTTTGCTAAT
CATGTTCATACTCTTATCTTCTCCACAG**CTCCTGGGCAACGTGCTGGTCTGTGTGCTGGCCCATCACTT**
TGGCAAAGAATTCACCCCACCAGTGCAGGCTGCCTATCAGAAAGTGGTGGCTGGTGTGGCTAATGCCCTGGC
CCACAAGTATCACTAAGCTCGCTTTCTTGCTGTCCAATTTCTATTAAAGGTTCTTTGTTCCCTAAGTCCAA
CTACTAAACTGGGGGATATTATGAAGGGCCTTGAGCATCTGGATTCTGCCTAATAAAAAACATTTATTTTCA
TTGCAATGATGTATTTAAATTATTTCTGAATATTTTACTAAAAAGGGAATGTGGGAGGTCAGTGCATTTAAA
ACATAAAGAAATGAAGAGCTAGTTCAAACCTTGGGAAAATACACTATATCTTAAACTCCATGAAAGAAGGTG

A human gene

chr11:5246500-5248500 (reverse strand):

ATATCTTAGAGGGAGGGCTGAGGGTTTGAAGTCCAACCTCCTAAGCCAGTGCCAGAAGAGCCAAGGACAGGTA
CGGCTGTCATCACTTAGACCTCACCCCTGTGGAGCCACACCCTAGGGTTGGCCAATCTACTCCCAGGAGCAGG
GAGGGCAGGAGCCAGGGCTGGGCATAAAAAGTCAGGGCAGAGCCATCTATTGCTTACATTTGCTTCTGACACA
ACTGTGTTCACTAGCAACCTCAAACAGACACC**ATGGTGCATCTGACTCCTGAGGAGAAGTCTGCCGTTACTG**
CCCTGTGGGGCAAGGTGAACGTGGATGAAGTTGGTGGTGGAGGCCCTGGGCAGGTTGGTATCAAGGTTACAAG
ACAGGTTTAAGGAGACCAATAGAACTGGGCATGTGGAGACAGAGAAGACTCTTGGGTTTCTGATAGGCACT
GACTCTCTCTGCCTATTGGTCTATTTTCCCACCCTTAG**GCTGCTGGTGGTCTACCCTTGGACCCAGAGGTTT**
TTTGAGTCCTTTGGGGATCTGTCCACTCCTGATGCTGTTATGGGCAACCCTAAGGTGAAGGCTCATGGCAAG
AAAGTGCTCGGTGCCTTTAGTGATGGCCTGGCTCACCTGGACAACCTCAAGGGCACCTTTGCCACACTGAGT
GAGCTGCACTGTGACAAGCTGCACGTGGATCCTGAGAACTTCAGGGTGAGTCTATGGGACGCTTGATGTTTT
CTTTCCCCTTCTTTTCTATGGTTAAGTTCATGTCATAGGAAGGGGATAAGTAACAGGGTACAGTTTAGAATG
GGAAACAGACGAATGATTGCATCAGTGTGGAAGTCTCAGGATCGTTTTAGTTTCTTTTATTTGCTGTTTCATA
ACAATTGTTTTCTTTTGTTTAATTCTTGCTTTCTTTTTTTTTCTTCTCCGCAATTTTACTATTATACTTAA



GGCTCTTTGCACCAATTC TAAAGAATAACAGTGATAATTTCTGGGTTAAGGCAATAGCAATATCTCTGCATAT
AAATATTTCTGCATATAAATTGTAAGTACTGATGTAAGAGGTTTCATATTGCTAATAGCAGCTACAATCCAGCTA
CCATTCTGCTTTTATTTTATGGTTGGGATAAGGCTGGATTATTCTGAGTCCAAGCTAGGCCCTTTTGCTAAT
CATGTTCATACTCTTATCTTCTCCACAG**CTCCTGGGCAACGTGCTGGTCTGTGTGCTGGCCATCACTT**
TGGCAAAGAATTCACCCACCAGTGCAGGCTGCCTATCAGAAAGTGGTGGCTGGTGTGGCTAATGCCCTGGC
CCACAAGTATCACTAAGCTCGCTTTCTTGCTGTCCAATTTCTATTAAAGGTTCTTTGTTCCCTAAGTCCAA
CTACTAAACTGGGGGATATTATGAAGGGCCTTGAGCATCTGGATTCTGCCTAATAAAAAACATTTATTTTCA
TTGCAATGATGTATTTAAATTATTTCTGAATATTTTACTAAAAAGGGAATGTGGGAGGTCAGTGCATTTAAA
ACATAAAGAAATGAAGAGCTAGTTCAAACCTTGGGAAAATACACTATATCTTAAACTCCATGAAAGAAGGTG

A human gene

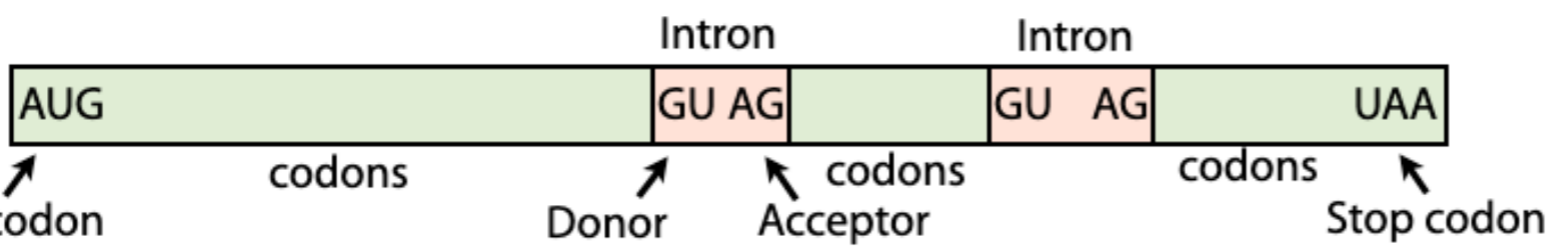
chr11:5246500-5248500 (reverse strand):

```
ATATCTTAGAGGGAGGGGCTGAGGGTTTGAAGTCCAACTCCTAAGCCAGTGGCCAGAAGAGCCAAGGACAGGTA
CGGCTGTCATCACTTAGACCTCACCCCTGTGGAGCAGGATCTACTCCCAGGAGCAGG
GAGGGCAGGAGCCAGGGCTGGGCATAAAAAGTCAGTACATTTGCTTCTGACACA
ACTGTGTTCACTAGCAACCTCAAACAGACACCATGATGTCATCTGACTCCTGAGGAGAAGTCTGCCGTTACTG
CCCTGTGGGGCAAGGTGAACGTGGATGAAGTTGGTGGTGGAGGCCCTGGGCAGGTGTTACAAG
ACAGGTTTAAGGAGACCAATAGAACTGGGCATGTGGAGACAGAGAAGACTCTTGCACACT
GACTCTCTCTGCCTATTGGTCTATTTTCCACCCTTAGAGTGGACCCAGAGGTTTC
TTTGAGTCCTTTGGGGATCTGTCCACTCCTGATGCTGTGTGAAGGCTCATGGCAAG
AAAGTGCTCGGTGCCTTTAGTGATGGCCTGGCTCACCTGGACAACCTCAAGGGCACCTTTGCCACACTGAGT
GAGCTGCACTGTGACAAGCTGCACGTGGATCCTGAGAACTTCAGGGTGTGAGTCTATGGGACGCTTGATGTTTT
CTTTCCCCTTCTTTTCTATGGTTAAGTTCATGTCATAGGAAGGGGATAAGTAACAGGGTACAGTTTAGAATG
GGAAACAGACGAATGATTGCATCAGTGTGGAAGTCTCAGGATCGTTTTAGTTTCTTTTATTTGCTGTTTCATA
ACAATTGTTTTCTTTTGTTTAATTCTTGCTTTCTTTTTTTTTTTCTTCTCCGCAATTTTTACTATTACTTAA
TCCCTTAACATTCTCTATAAGCAAAAGCAAAATATCTCTCAGATAGATTAACTAACTTAAAAAAAACTTTACA
```

Start codon

Donor

Acceptor



Stop codon

```
GCCCTCTTTGCACCAATTCATAAAGAATAACAGTGATAAATTTCTGGGTAAAGGCAATAGCAATATCTCTGCATAT
AAATATTTCTGCATATAAATTGTAACCTGATGTAAGAGGTTTCATATTGCTAATAGCAGCTACAATCCAGCTA
ATTTTATGGTTGGGATAAGGCTGGATTATTCTGAGTCCAAGCTAGGCCCTTTTGCTAAT
CTTATCTTCTCCCAAGCTCCTGGGCAACGTGCTGGTCTGTGTGCTGGCCCATCACTT
ACCCACCAAGTGCAGGCTGCCTATCAGAAAGTGGTGGCTGGTGTGGCTAATGCCCTGGC
CCACAAGTATCACTAACTCGCTTTCTTGCTGTCCAATTTCTATTAAAGGTTCTTTGTTCCCTAAGTCCAA
CTACTAAACTGGGGGATATTATGAAGGGCCTTGAGCATCTGGATTCTGCCTAATAAAAAACATTTATTTTCA
TTGCAATGATGTATTTAAATTATTTCTGAATATTTTACTAAAAAGGGAATGTGGGAGGTCAGTGCATTTAAA
ACATAAAGAAATGAAGAGCTAGTTCAAACCTTGGGAAAATACACTATATCTTAAACTCCATGAAAGAAGGTG
```

A human gene

chr11:5246500-5248500 (reverse strand):

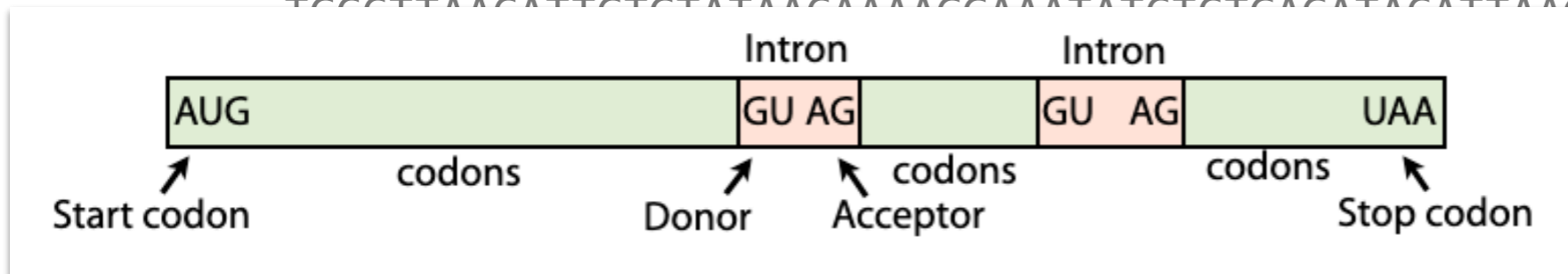
```
ATATCTTAGAGGGAGGGGCTGAGGGTTTGAAGTCCAACTCCTAAGCCAGTGGCCAGAAGAGCCAAGGACAGGTA
CGGCTGTCATCACTTAGACCTCACCCCTGTGGAGCAGGATCTACTCCCAGGAGCAGG
GAGGGCAGGAGCCAGGGCTGGGCATAAAAAGTCAGTACATTTGCTTCTGACACA
ACTGTGTTCACTAGCAACCTCAAACAGACACCATGATGTCATCTGACTCCTGAGGAGAAGTCTGCCGTTACTG
CCCTGTGGGGCAAGGTGAACGTGGATGAAGTTGGTGGTGGAGGCCCTGGGCAGGTGTTACAAG
ACAGGTTTAAGGAGACCAATAGAACTGGGCATGTGGAGACAGAGAAGACTCTTGCACACTGCACT
GACTCTCTCTGCCTATTGGTCTATTTTCCACCCTTAGTGGTGGTGGAGGCCCTGGGCAGGTGTTACAAG
TTTGAGTCCTTTGGGGATCTGTCCACTCCTGATGCTGTGTTGAAGGCTCATGGCAAG
TCACCTGGACAACCTCAAGGGCACCTTTGCCCACTGAGT
TGAGAACTTCAGGGTGTGAGTCTATGGGACGCTTGATGTTTT
TCATAGGAAGGGGATAAGTAACAGGGTACAGTTTAGAATG
GTCTCAGGATCGTTTTAGTTTTCTTTTATTTGCTGTTTCATA
ACAAATGTTTTCTTTTGTTTAATTTCTTTCTTTTTTTTTTTCTTCTCCGCAATTTTTACTATTATACTTAA
TCCCTTAACATTCTCTATAAGCAAAAGCAAAATATCTCTCAGATAGATTAACTAACTTAAAAAAAACTTTACA
ATTCATAATCTCCCTACTTTAT
AAGTGTAATGTTTTAATATGTG
TGCTTTCTTCTTTAATATACT
CAATAATGATACAATGTATCAT
GCAATAGCAATATCTCTGCATAT
AAATATTTCTGCATATAAATTGTAAGTACTGATGTAAGAGGTTTCATATTGCTAATAGCAGCTACAATCCAGCTA
ATTTTATGGTTGGGATAAGGCTGGATTATTCTGAGTCCAAGCTAGGCCCTTTTGCTAAT
CTTATCTTCTCCCAAGCTCCTGGGCAACGTGCTGGTCTGTGTGCTGGCCCATCACTT
ACCCACCAAGTGCAGGCTGCCTATCAGAAAGTGGTGGCTGGTGTGGCTAATGCCCTGGC
CCACAAGTATCACTAACTCGCTTTCTTGCTGTCCAATTTCTATTAAAGGTTCTTTGTTCCCTAAGTCCAA
CTACTAAACTGGGGGATATTATGAAGGGCCTTGAGCATCTGGATTCTGCCTAATAAAAAACATTTATTTTCA
TTGCAATGATGTATTTAAATTATTTCTGAATATTTTACTAAAAAGGGAATGTGGGAGGTCAGTGCATTTAAA
ACATAAAGAAATGAAGAGCTAGTTCAAACCTTGGGAAAATACACTATATCTTAAACTCCATGAAAGAAGGTG
```

Start codon

Donor

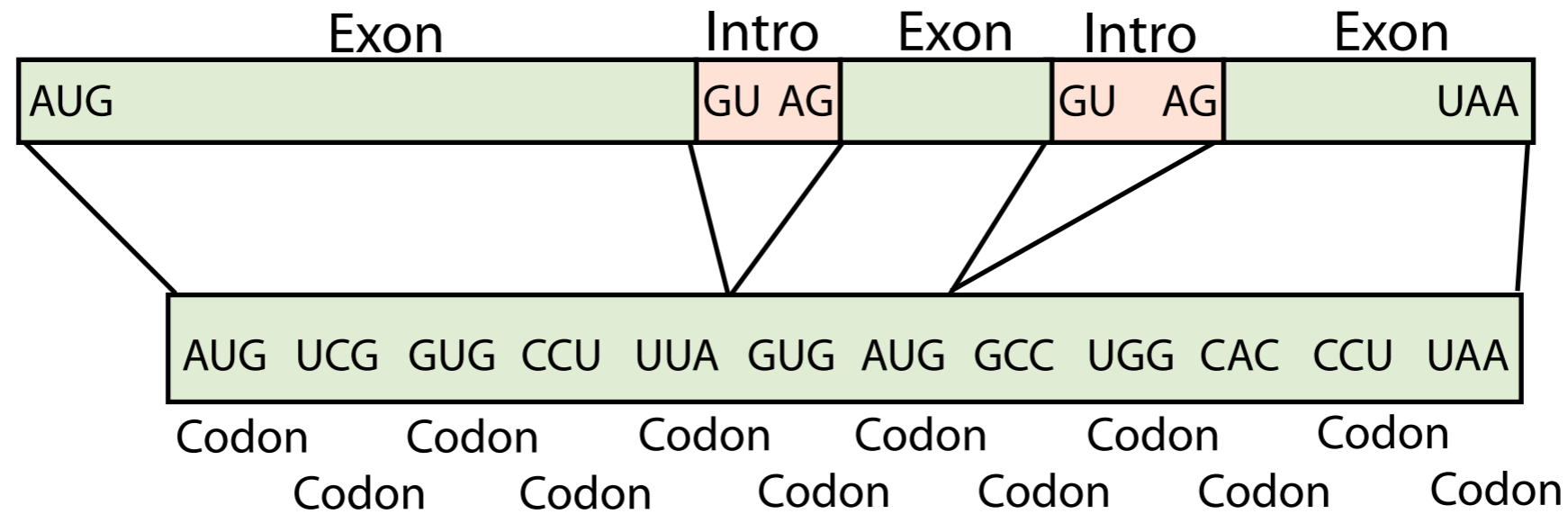
Acceptor

Coding sequence has length 444 - a multiple of 3!



Stop codon

Codons and the genetic code



Let's put these in a sequence model!

We note that signals can be "fuzzy"

Most ATGs **aren't** start codons, most
TAAs **aren't** stop codons, most GTs
aren't donors, etc

Only some donor or acceptors *in a gene*
are involved in splicing

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