

IFT 3290 HIVER 2006

Bioinformatique

Miklós Csűrös

André-Aisenstadt 3149

csuros@iro.umontreal.ca

<http://www.iro.umontreal.ca/~csuros/IFT3290/>

PLAN DE COURS

Préalables :

IFT2010 — Programmation 2

IFT1978 — Probabilités et statistique

Horaire :

cours	mercredi 11 :30–13 :30	Z 215
cours	vendredi 14 :30–16 :30	Z 350

Évaluation : 30% TP, 30% intra, 30% final, 10% présentation

4 travaux : programmation et théorie

DIVERSITÉ ET UNIVERSALITÉ

Diversité de la vie

Universalité au niveau moléculaire : interaction de macromolécules, mécanismes communes (évidence de l'Évolution)

Protéines et acides nucléiques

PROTÉINES

Fonctions :

- structure de la cellule [p.e collagène form des fibres en tissu conjonctif]
- enzymes : catalyseur de réactions spécifiques
- protéines membranaires
- signalisation [p.e. EGF — facteur de croissance pour l'épiderme]
- ...

EXEMPLE : KINESINE

[kinesine]

(moteur moléculaire)

©UIUC Theoretical and Computational Biophysics Group

INFO SUR UNE PROTÉINE

facteur de croissance EGF : voie de signalisation ([www](#))

recherche dans une base de donnée : SWISS-PROT ([www](#))

(Tapez «EGF receptor» au [site Web de SWISS-PROT](#) et suivez le lien vers EGFR_HUMAN.)

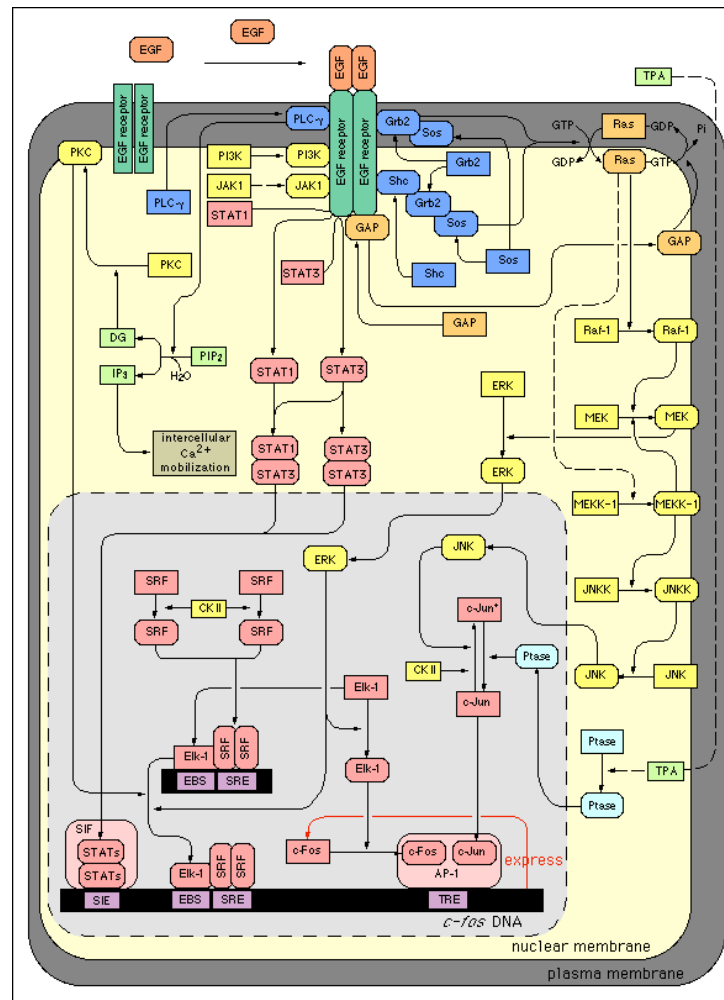
pleine d'information liée : e.g., structure 3D ([www](#))

(Suivez le lien «PDB» dans la ligne «HSSP».)

Séquence

MRPSGTAGAA LLALLAALCP ASRALEEKV CQGTSNKLTQ
LGTfEDHFLS LQRMFNNCEV VLGnLEITYV QRNYDLSFLK
...

VOIE DE SIGNALISATION EGF



EGF DANS SWISS-PROT

NcbProt View of SWISS-PROT: P00533

11/03/ 9:42 AM

ExPASy Home page	Site Map	Search ExPASy	Contact us	SWISS-PROT
Hosted by CBR Canada Mirror sites: Bolivia China Korea Switzerland Taiwan USA				

NcbProt View of SWISS-PROT: P00533 Printer-friendly view | Quick BlastP search

[\[General\]](#) [\[Name and origin\]](#) [\[References\]](#) [\[Comments\]](#) [\[Cross-references\]](#) [\[Keywords\]](#) [\[Features\]](#) [\[Sequence\]](#) [\[Links\]](#)

General information about the entry

Entry name **EGFR_HUMAN**

Primary accession number **P00533**

Secondary accession numbers P06268 Q14225 Q9UMD7

Entered in SWISS-PROT in Release 01, July 1986

Sequence was last modified in Release 35, November 1997

Annotations were last modified in Release 41, June 2002

Name and origin of the protein

Protein name **Epidermal growth factor receptor [Precursor]**

Synonyms **EC 2.7.1.12**

Gene name **EGFR or ERBB1**

From **Eukaryota; Metazoa; Chordata; Craniata; Vertebrata; Euteleostomi; Mammalia; Eutheria; Primates; Catarrhini; Hominidae; Homo.**

References

[1] SEQUENCE FROM NUCLEIC ACID (ISOFORM 1).
MEDLINE=84219729; PubMed=6328312; [NCBI, ExPASy, EBI, Israel, Japan]
Ulirsch A, Coussens L, Havlick J S, Dull T J, Gray A, Tam A W, Lee J, Yarden Y, Libermann T A, Schlessinger J, Downward J, Mayes F L V, White M, Waterfield M D, Seeburg P H.
"Human epidermal growth factor receptor cDNA sequence and aberrant expression of the amplified gene in A431 epidermoid carcinoma cells"; Nature 309:418-425(1984).

[2] SEQUENCE FROM NUCLEIC ACID (ISOFORM 2).
TISSUE=Placenta;
MEDLINE=95382957; PubMed=7654368; [NCBI, ExPASy, EBI, Israel, Japan]
Hekis J Y, Strunk B C, Scoccia B.
"Possible role of variant RNA transcripts in the regulation of epidermal growth factor receptor expression in human placenta"; Mol. Reprod. Dev. 41:149-156(1995).

[3] SEQUENCE FROM NUCLEIC ACID (ISOFORM 2).
TISSUE=Placenta;
MEDLINE=97078686; PubMed=8918811; [NCBI, ExPASy, EBI, Israel, Japan]
Reiter J L, Mailhe N J.
"A 1.8 kb alternative transcript from the human epidermal growth factor receptor gene encodes a truncated form of the receptor."; Nucleic Acids Res. 24:4050-4056(1996).

[4] SEQUENCE FROM NUCLEIC ACID (ISOFORM 2).
TISSUE=Placenta;
MEDLINE=97256547; PubMed=9103388; [NCBI, ExPASy, EBI, Israel, Japan]
Hekis J Y, Garin J, Niederberger C, Scoccia B.
"Expression of a truncated epidermal growth factor receptor-like protein (TEGFR) in ovarian cancer."; Gynecol. Oncol. 65:36-41(1997).

[5] SEQUENCE FROM NUCLEIC ACID (ISOFORMS 3 AND 4).
TISSUE=Placenta;
MEDLINE=21100872; PubMed=11161793; [NCBI, ExPASy, EBI, Israel, Japan]
Reiter J L, Threadgill D W, Eley G D, Strunk K E, Daniels A J, Scheel Sinclair C, Pearsall R S, Green P J, Yee D, Lampland A L, Balasubramaniam S, Crossley T D, Magnuson T R, James C D, Mailhe N J.
"Comparative genomic sequence analysis and isolation of human and mouse alternative EGFR transcripts encoding truncated receptor isoforms."; Genomes 71:1-20(2001).

[6] SEQUENCE OF 575-687 FROM NUCLEIC ACID.
Reiter J L, Threadgill D W, Daniels A J, Scheel C M, Lampland A L, Balasubramaniam S, Crossley T O, Magnuson T R, Mailhe N J.
"Human and mouse alternative EGFR transcripts encoding only the extracellular domain of the receptor."; Submitted (FEB-1999) to the EMBL/GenBank/DBJ databases.

[7] SEQUENCE OF 713-924 FROM NUCLEIC ACID.
MEDLINE=84196372; PubMed=6326261; [NCBI, ExPASy, EBI, Israel, Japan]
Lin C R, Chen W S, Kraiger W, Stolarsky L S, Weber W, Evans R M, Verma I M, Gill G N, Rosenfeld M G.
"Expression cloning of human EGF receptor complementary DNA: gene amplification and three related messenger RNA products in A431 cells."; Science 224:843-848(1984).

[8] SEQUENCE OF 150-962 FROM NUCLEIC ACID.
MEDLINE=8425835; PubMed=6330563; [NCBI, ExPASy, EBI, Israel, Japan]
Xu Y H, Ishii S, Clark A J L, Sullivan M, Wilson R K, Ma D P, Roe B A, Merlino G T, Pastan I.

http://ca.expasy.org/cgi-bin/ncbprotview.pl?P00533

Page 1 of 5

NcbProt View of SWISS-PROT: P00533

11/03/ 9:43 AM

Cross-references

X00588, CAA25240.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
U95089, AAB53063.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
U48722, AAC50802.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
U48723, AAC50804.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
U48724, AAC50796.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
U48725, AAC50797.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
U48726, AAC50798.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
U48727, AAC50799.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
U48728, AAC50800.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
U48729, AAC50801.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF288738, AAG35786.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF288738, AAG35787.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF288738, AAG35788.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF288738, AAG35789.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF288738, AAG35790.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01080.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01081.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01082.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01083.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01084.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01085.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01086.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01087.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01088.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01089.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01090.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01091.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01092.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01093.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01094.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01095.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01096.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01097.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01098.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01099.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01100.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01101.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01102.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01103.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01104.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01105.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01106.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01107.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01108.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01109.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01110.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01111.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01112.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01113.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01114.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01115.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01116.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01117.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01118.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01119.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01120.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01121.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01122.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01123.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01124.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01125.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01126.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01127.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01128.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01129.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01130.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01131.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01132.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01133.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01134.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01135.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01136.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01137.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01138.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01139.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01140.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01141.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01142.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01143.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01144.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01145.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01146.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01147.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01148.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01149.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01150.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01151.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01152.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01153.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01154.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01155.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01156.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01157.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01158.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01159.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01160.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01161.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01162.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01163.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01164.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01165.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01166.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01167.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01168.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01169.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01170.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01171.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01172.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01173.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01174.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01175.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01176.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01177.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01178.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01179.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01180.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01181.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01182.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01183.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01184.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01185.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01186.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01187.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01188.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01189.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01190.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01191.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01192.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01193.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01194.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01195.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01196.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01197.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01198.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01199.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01200.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01201.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01202.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01203.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01204.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01205.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01206.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01207.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01208.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01209.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01210.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01211.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01212.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01213.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01214.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01215.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01216.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01217.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01218.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01219.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01220.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01221.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01222.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01223.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01224.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01225.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01226.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01227.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01228.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01229.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01230.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01231.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01232.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01233.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01234.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01235.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01236.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01237.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01238.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01239.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01240.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01241.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01242.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01243.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01244.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01245.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01246.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01247.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01248.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01249.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01250.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01251.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01252.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01253.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01254.1, --	[EMBL / GenBank / DDBJ] [CoDingSequence]
AF277897, AAK01	

EGF DANS PDB

Structure Explorer - 1FGK

1/10/03 9:22 AM



Structure Explorer - 1FGK



Title Crystal Structure Of The Tyrosine Kinase Domain Of Fibroblast Growth Factor Receptor 1
Classification Phosphotransferase
Compound Mol. Id: 1; Molecule: Fgf Receptor 1; Chain: A, B; Fragment: Tyrosine Kinase Domain, Human Fgfr1 Residues That Possess Ptk Activity; Synonym: Fgfr1K, Fibroblast Growth Factor Receptor 1; E.c: 2.7.1.112; Engineered: Yes; Mutation: L457V, C488A, C584S
Exp. Method X-ray Diffraction



View Structure



[Summary Information](#)

[View Structure](#)

[Download/Display File](#)

[Structural Neighbors](#)

[Geometry](#)

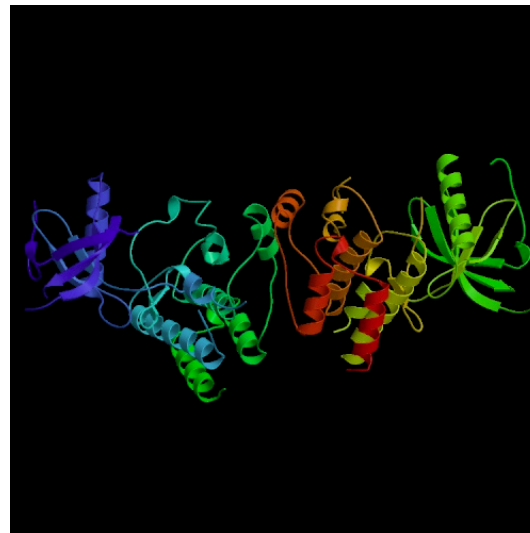
[Other Sources](#)

[Sequence Details](#)

[Structure Factors
\(compressed\)](#)

Explore

[SearchLite](#) [SearchFields](#)



© RCSB

<http://www.rcsb.org/pdb/cgi/structure/1FGK>

Page 1 of 1

SÉQUENCES PROTÉIQUES

Protéine : un ou plusieurs polymères d'acides aminés (polymère est une grosse molécule formée de l'union de molécules plus petites)

Alanine (Ala) Glycine (Gly) Méthionine (Met)

Sérine (Ser) Cystéine (Cys) Histidine (His)

Asparagine (Asn) Thréonine (Thr) Acide aspartique (Asp)

Isoleucine (Ilu) Proline (Pro) Valine (Val)

Acide glutamique (Glu) Lysine (Lys) Glutamine (Gln)

Tryptophane (Trp) Phénylalanine (Phe) Leucine (Leu)

Arginine (Arg) Tyrosine (Tyr)

La séquence détermine la structure

⇒ **Problème de prédiction de structure protéique.**

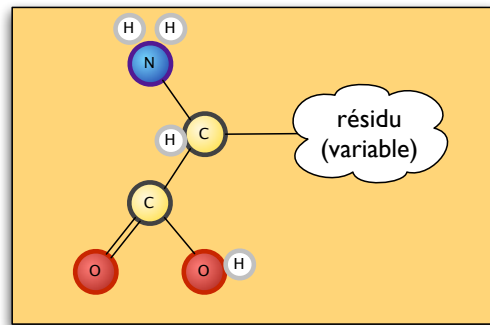
ACIDES AMINÉS

[animation]

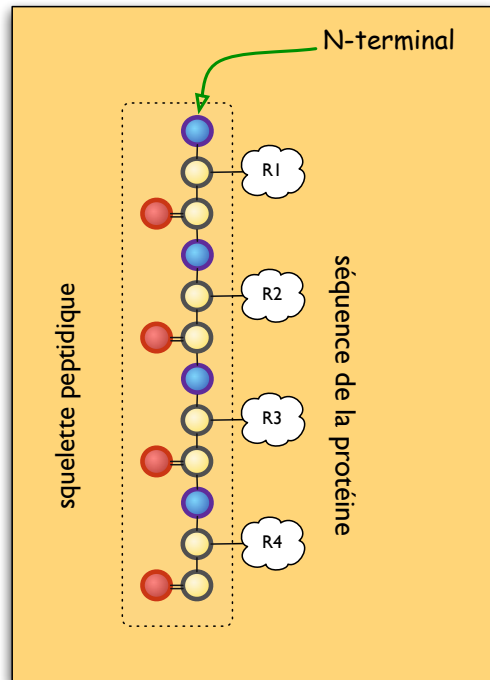
(Phe, Pro, Met, Gly)

Song Tan, Penn State U

ACIDES AMINÉS ET PROTÉINES



Acide aminé



Protéine

AUTRES PROBLÈMES

Comment est la structure établie ?

⇒ Repliement de protéines

La structure détermine la fonction

⇒ Problème de prédiction de fonction protéique.

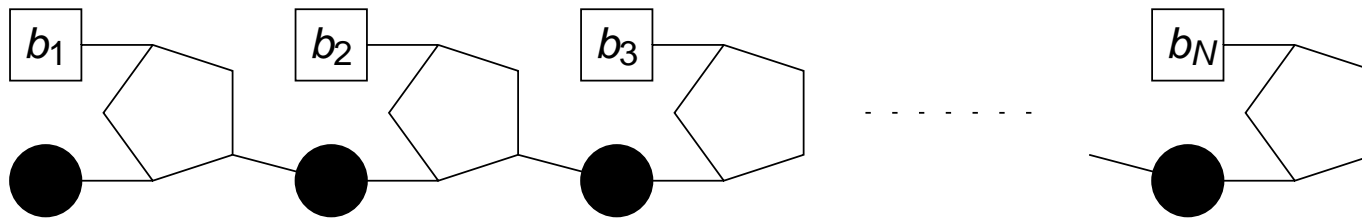
Interaction de protéines

⇒ Analyses de voies métaboliques, chemins de signalisation, réseaux régulateurs

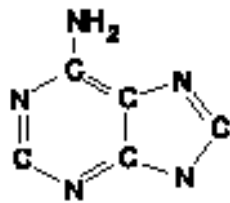
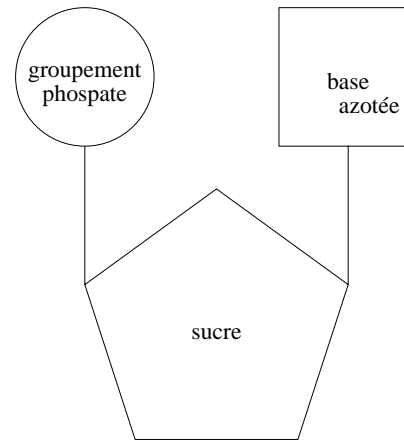
L'ADN ?

Polynucléotide :

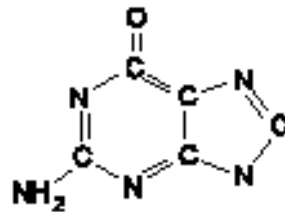
polymère (ie, chaîne) de **nucléotides**



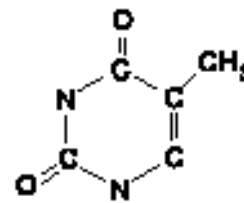
NUCLÉOTIDES



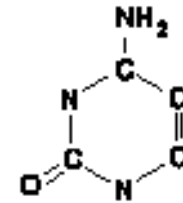
Adénine



Guanine



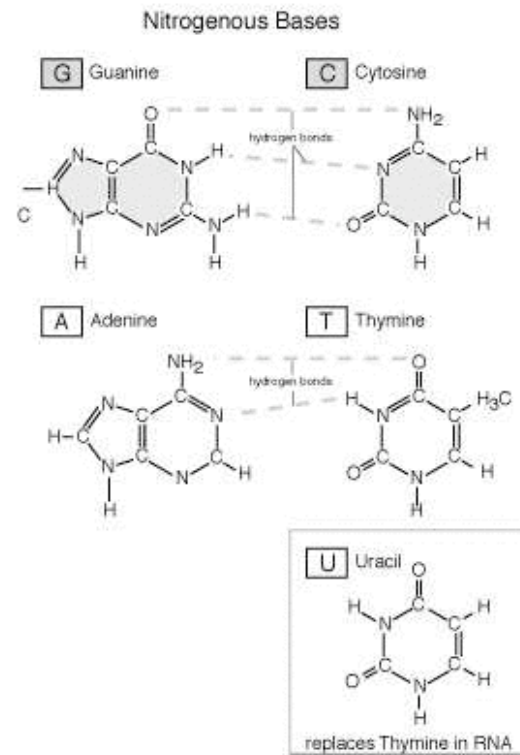
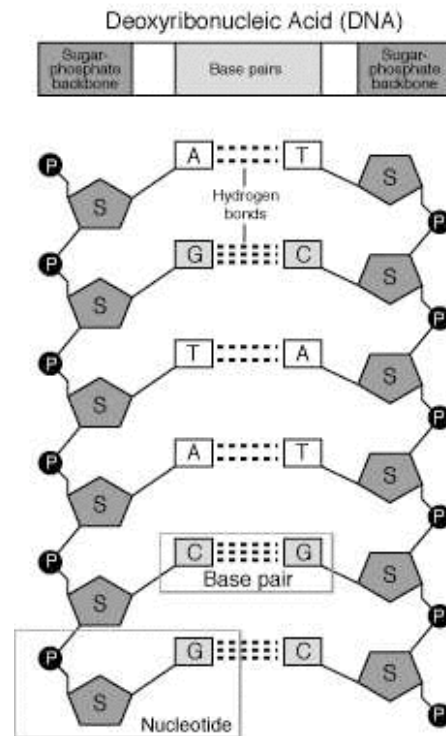
Thymine



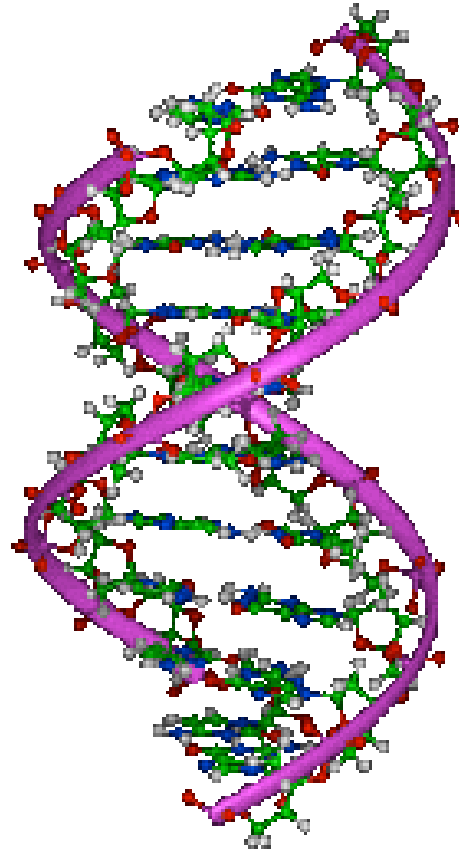
Cytosine

4 bases :

ADN : BASES COMPLÉMENTAIRES



WATSON ET CRICK : DOUBLE HÉLICE



1 pas de l'hélice : 10 paires de base, 3.4 nm (forme usuelle)

equipment, and to Dr. G. E. R. Descom and the captain and officers of R.R.S. *Discovery II* for their part in making the observations.

¹Young, F. B., Gerrard, H., and Jevons, W., *Phil. Mag.*, **46**, 149 (1900).
²Laguarda-Hugos, M. S., *Mem. Soc. Roy. Astr. Soc., Geophys. Surv.*, **8**, 263 (1949).
³Von Arz, H. S., *Wegede Hole Papers in Phys. Oceanog. Meteor.*, **11** (5) (1946).
⁴Ekman, V. W., *Arkiv. Mat. Astron. Fysik. (Stockholm)*, **2**(11) (1955).

MOLECULAR STRUCTURE OF NUCLEIC ACIDS

A Structure for Deoxyribose Nucleic Acid

WE wish to suggest a structure for the salt of deoxyribose nucleic acid (D.N.A.). This structure has novel features which are of considerable biological interest.

A structure for nucleic acid has already been proposed by Pauling and Corey¹. They kindly made their manuscript available to us in advance of publication. Their model consists of three intertwined chains, with the phosphates near the fibre axis, and the bases on the outside. In our opinion, this structure is unsatisfactory for two reasons: (1) We believe that the material which gives the X-ray diagrams is the salt, not the free acid. Without the acidic hydrogen atoms it is not clear what forces would hold the structure together, especially as the negatively charged phosphates near the axis will repel each other. (2) Some of the van der Waals distances appear to be too small.

Another three-chain structure has also been suggested by Fraser (in the press). In his model the phosphates are on the outside and the bases on the inside, linked together by hydrogen bonds. This structure as described is rather ill-defined, and for this reason we shall not comment on it.

We wish to put forward a radically different structure for the salt of deoxyribose nucleic acid. This structure has two helical chains each coiled round the same axis (see diagram). We have made the usual chemical assumptions, namely, that each chain consists of phosphate diester groups joining β -D-deoxy-ribofuranose residues with 3',5' linkages. The two chains (but not their bases) are related by a dyad perpendicular to the fibre axis. Both chains follow right-handed helices, but owing to the dyad the sequence of the atoms in the two chains run in opposite directions. Each chain loosely resembles Furbert's model No. 1; that is, the bases are on the inside of the helix and the phosphates on the outside. The configuration of the sugar and the atoms near it is close to Furbert's 'standard configuration', the sugar being roughly perpendicular to the attached base. There



This figure is purely diagrammatic. The two ribbons represent the two ribbons—sugar and phosphate—groups. The horizontal rungs represent the pairs of bases holding the chains together. The vertical line marks the fibre axis.

is a residue on each chain every 3.4 Å. in the z-direction. We have assumed an angle of 36° between adjacent residues in the same chain, so that the structure repeats after 10 residues on each chain, that is, after 34 Å. The distance of a phosphorus atom from the fibre axis is 10 Å. As the phosphates are on the outside, cations have easy access to them.

The structure is an open one, and its water content is rather high. At lower water contents we would expect the bases to tilt so that the structure could become more compact.

The novel feature of the structure is the manner in which the two chains are held together by the purine and pyrimidine bases. The planes of the bases are perpendicular to the fibre axis. They are joined together in pairs, a single base from one chain being hydrogen-bonded to a single base from the other chain, so that the two lie side by side with identical z-co-ordinates. One of the pair must be a purine and the other a pyrimidine for bonding to occur. The hydrogen bonds are made as follows: purine position 1 to pyrimidine position 1; purine position 6 to pyrimidine position 6.

If it is assumed that the bases only occur in the structure in the most plausible tautomeric forms (that is, with the keto rather than the enol configurations) it is found that only specific pairs of bases can bond together. These pairs are: adenine (purine) with thymine (pyrimidine), and guanine (purine) with cytosine (pyrimidine).

In other words, if an adenine forms one member of a pair on either chain, then on these assumptions

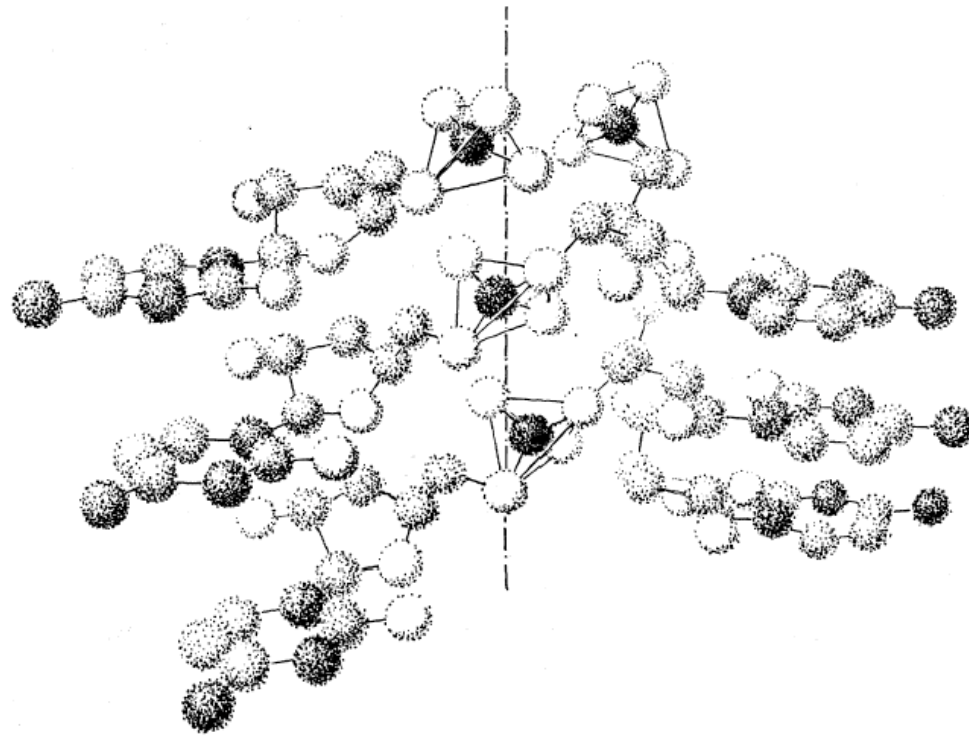
It has not escaped our notice that the specific pairing we have postulated immediately suggests a possible copying mechanism for the genetic material

It has not escaped our notice that the specific pairing we have postulated immediately suggests a possible copying mechanism for the genetic material.

Full details of the structure, including the conditions assumed in building it, together with a set of co-ordinates for the atoms, will be published elsewhere.

We are much indebted to Dr. Jerry Donohue for constant advice and criticism, especially on interatomic distances. We have also been stimulated by a knowledge of the general nature of the unpublished experimental results and ideas of Dr. M. H. F. Wilkins, Dr. R. E. Franklin and their co-workers at

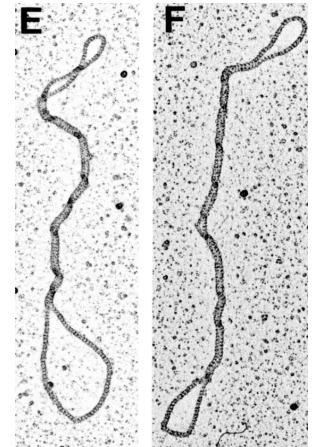
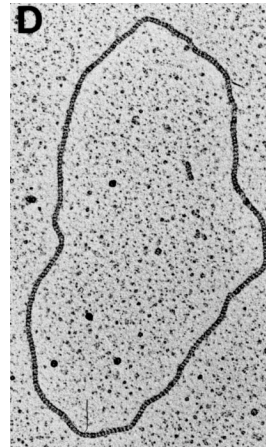
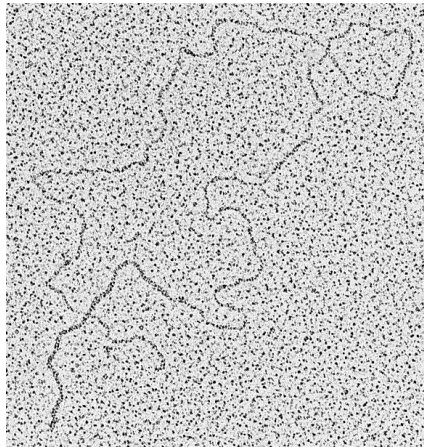
UN AUTRE MODÈLE — LE TRIPLE HÉLICE



Pauling & Corey, *PNAS* 39 : 84 (1953)

ADN - STRUCTURE

la molécule d'ADN peut être
linéaire (nos chromosomes), ou
circulaire (bactéries, organelles)



ADN

stockage de l'ADN : dans le noyau de la cellule, organisé en chromosomes (en eukaryotes)

1. duplication de l'information en ADN : hérédité
2. duplication de l'information en ARN (souvent traduit en protéines)

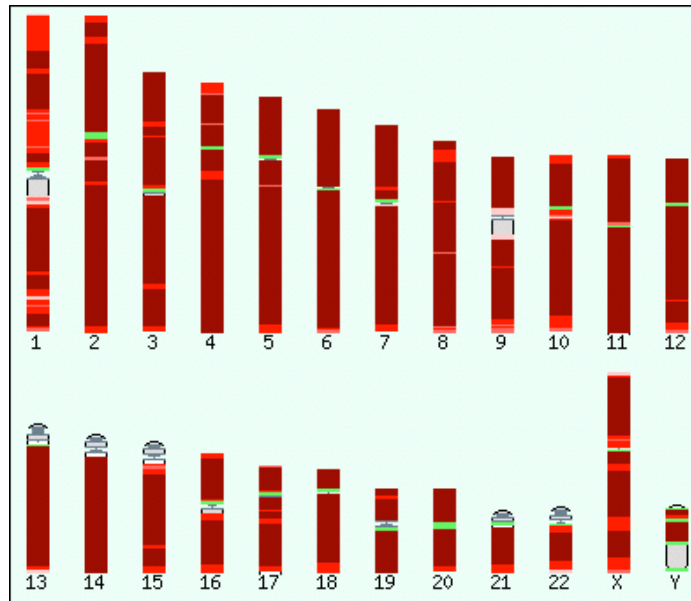
REPLICATION

[animation]

[animation]

Lodish et al., *Molecular Biology of the Cell*, 2002

CHROMOSOMES



22 chromosomes diploïdes et 2 chromosomes sexuels

tailles : $250 \cdot 10^6$ – $23 \cdot 10^6$ paires de base

taille totale du génome humaine : $3 \cdot 10^9$ pb

STRUCTURE DE CHROMOSOMES ET REPLICATION

[animation]

CODE GÉNÉTIQUE

20 acides aminés encodés par 4 nucléotides ?

Encodages par triplets : 64 codons

AAA = Phe	AAG = Phe	AAT = Leu	AAC = Leu
AGA = Ser	AGG = Ser	AGT = Ser	AGC = Ser
ATA = Tyr	ATG = Tyr	ATT = FIN	ATC = FIN
...			

3 triplets d'arrêt, 1 triplet de début (encode aussi Met)

DE L'ADN À PROTÉINE

1. transcription : copie du brin informatif à ARN messenger
(ARN : utilise Uracile au lieu de Thymine)

2. traduction : ARNm à protéine (par le ribosome) : acides aminés fournis par ARN de transfert

Mécanisme universelle !

⇒ Problème de déterminer la séquence de l'ADN.

GÈNES

Gène : unité d'hérédité, se traduit à une protéine (discussion plus raffinée plus tard)

Exons et introns

⇒ Problème de prédiction de gènes.

Expression d'un gène : quantité de ARNm transcrit ⇒ Problème d'analyse de l'expression génique.

COMPARAISONS

Comparaisons de séquences : pour déterminer structure, fonction, ...

⇒ Problème de comparaison de séquences.

Grandes bases de données

⇒ Problème de recherche de séquences.

AUTRES SUJETS

Séquençage de l'ADN

Phylogénies

Calcul moléculaire