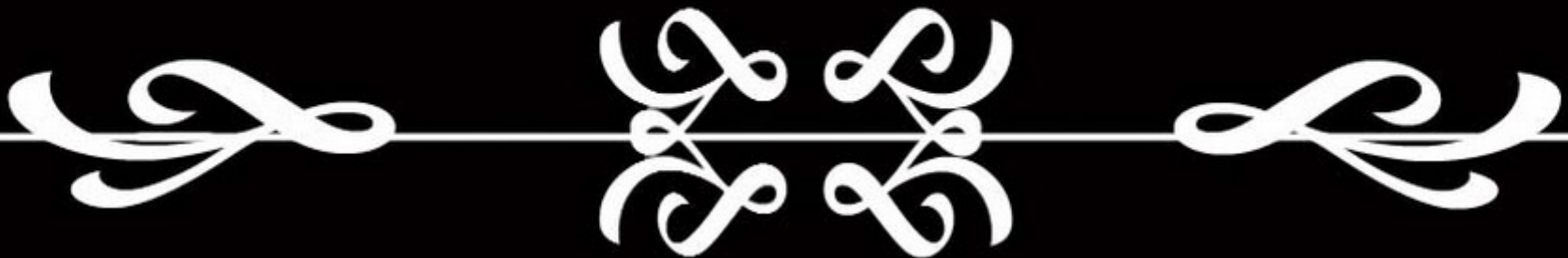




Une brève histoire
– *de la* –
BIOINFORMATIQUE

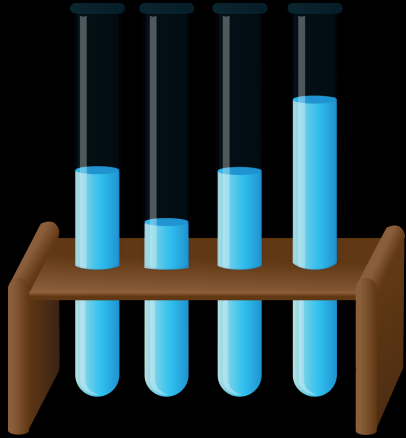


PARTIE I

(1950-1970)

FACTEURS LIMITANTS EN RECHERCHE...

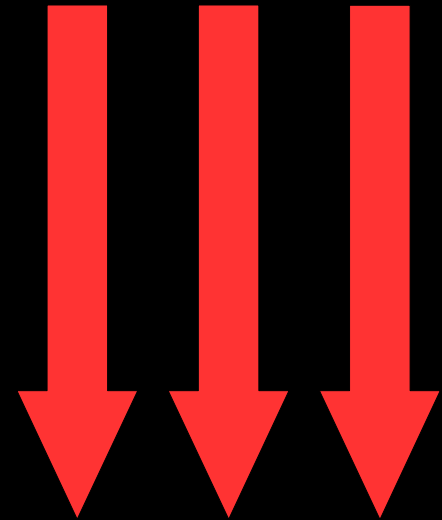
...en 2015



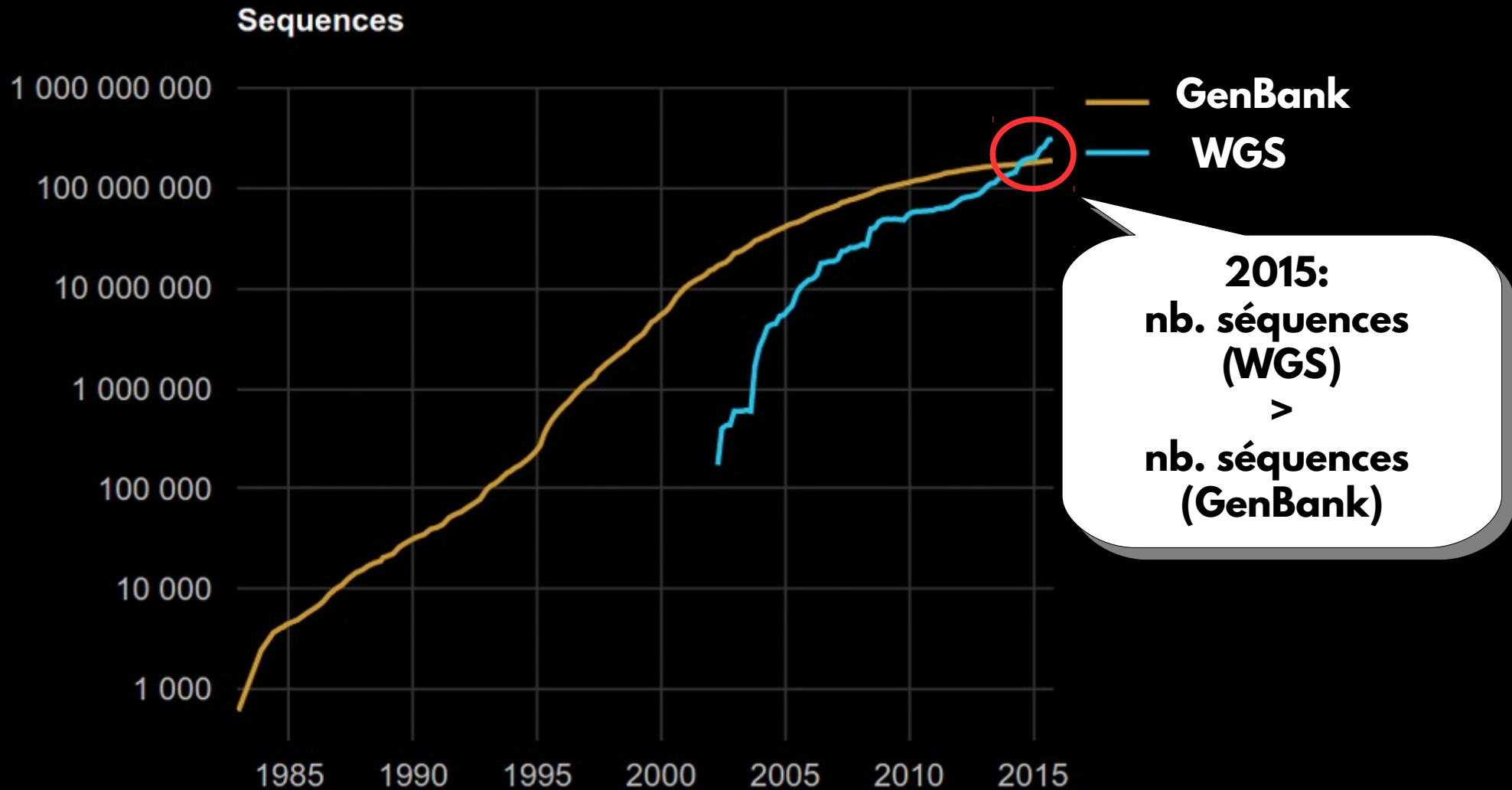
“Générer les données”



“Analyser les données”



CROISSANCE EXPONENTIELLE DE LA QUANTITÉ DE SÉQUENCES D'ADN PUBLIÉES !



PROBLÈMES:

1) Notre capacité de calcul



PROBLÈMES:

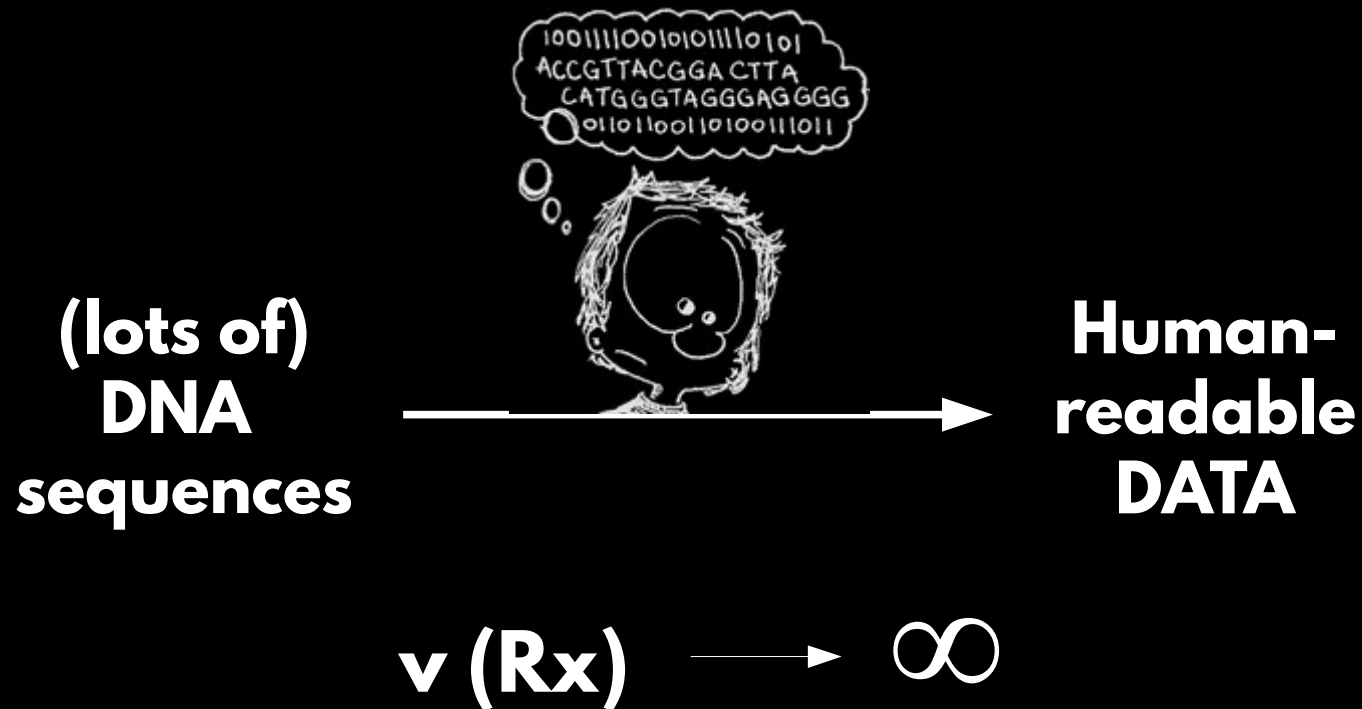
2) Surcharge de données brutes

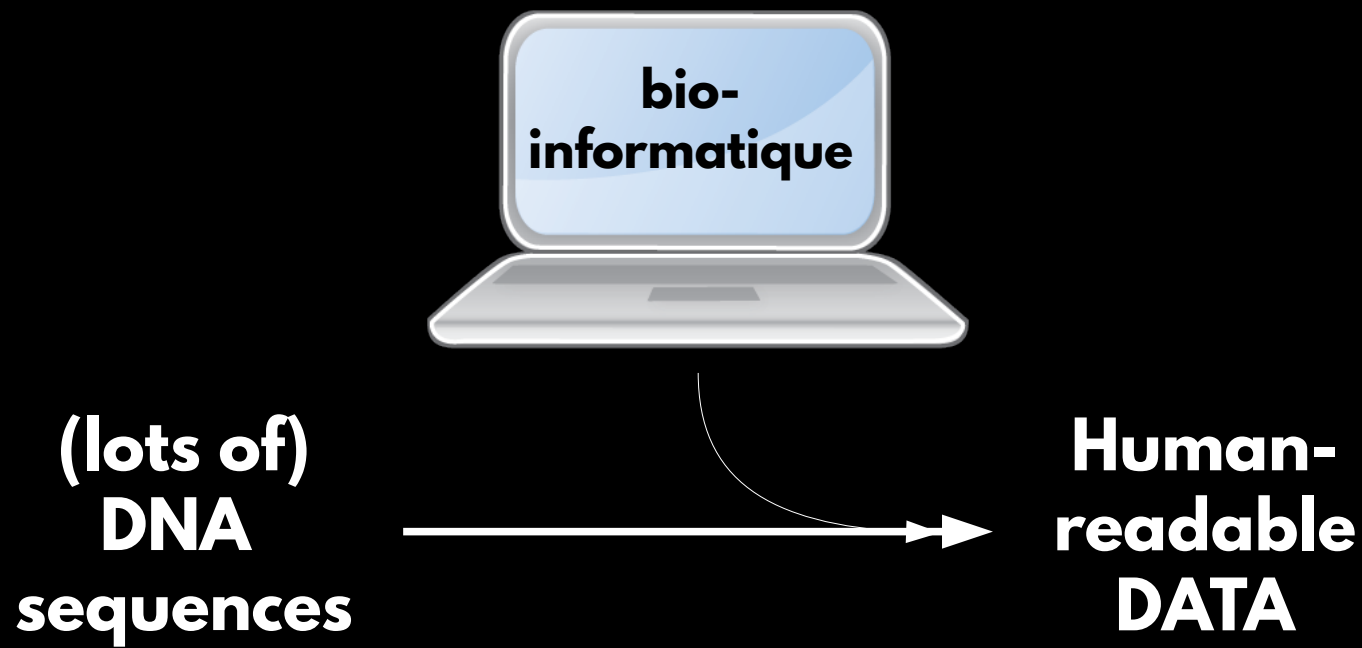


...le cerveau humain ne suffit pas à la tâche !

100111100101011110101
ACCGTTACGGA CTTA
CATGGGTAGGGAGGGG
01101100110100111011

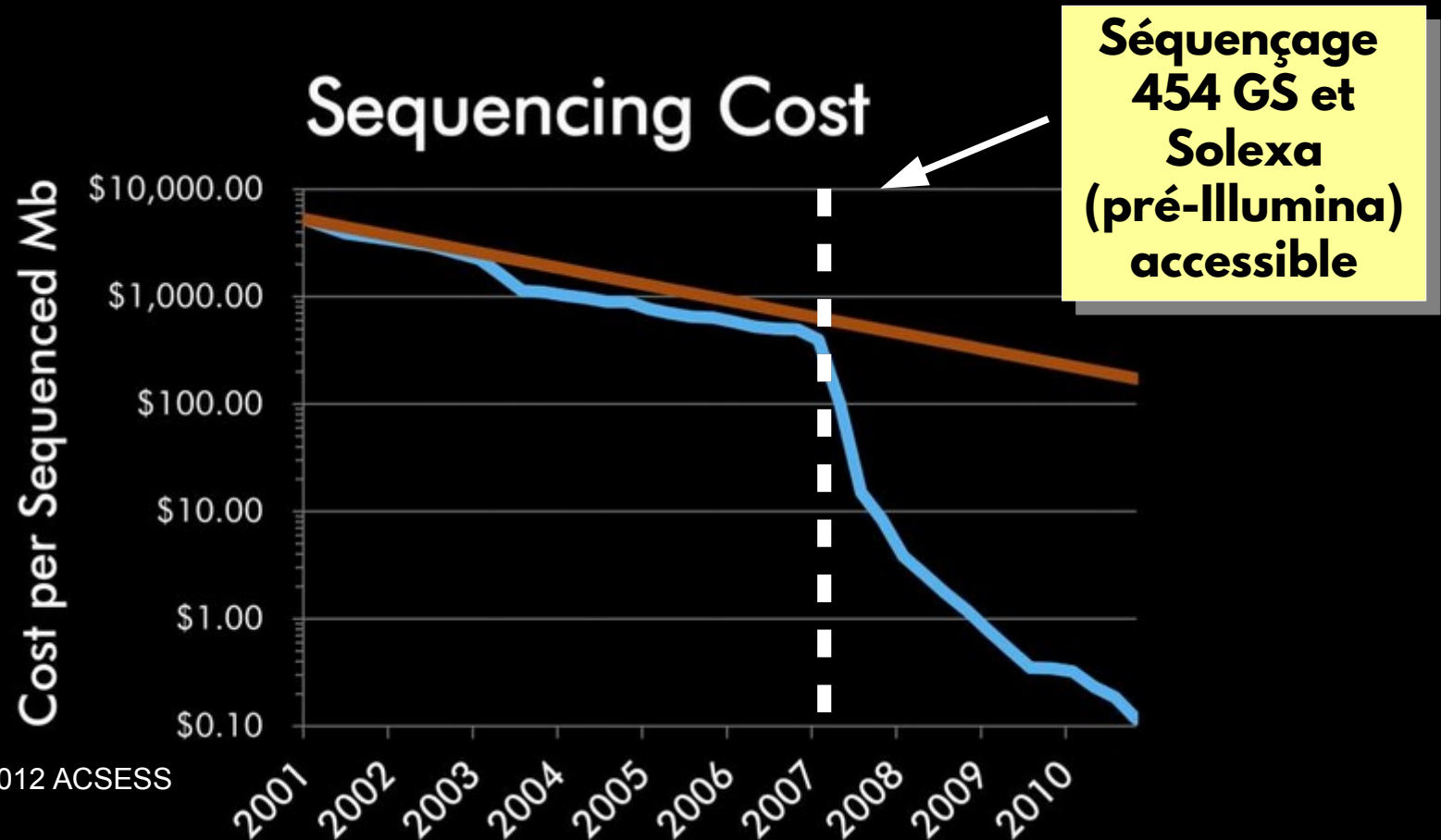






v (Rx) → **Heures,
jours,
semaines...**

TOUT PORTE À CROIRE QUE...



...la bioinformatique serait une discipline récente, venue prêter main forte à l'analyse de milliers de séquences d'ADN.

Est-ce vraiment le cas ? Pour le savoir...



RETOURNONS EN 1950...

1950

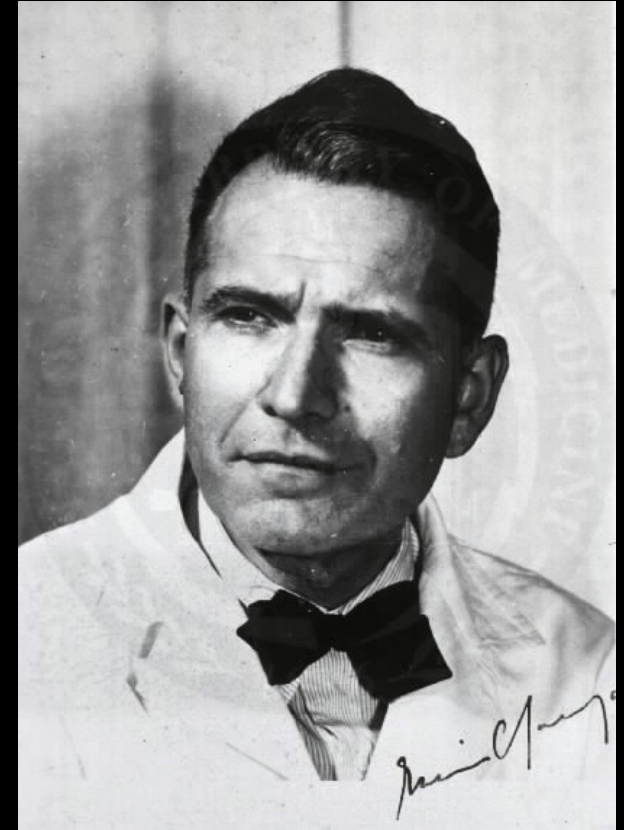
ON SAIT QUE... PEU DE CHOSES SUR LA STRUCTURE DE L'ADN.

- DÉSOXYRIBOSE
- PHOSPHATE
- QUATRE BASES AZOTÉES (A-T-C-G)

$\% A = \% T$

$\% C = \% G$

...et c'est presque tout.



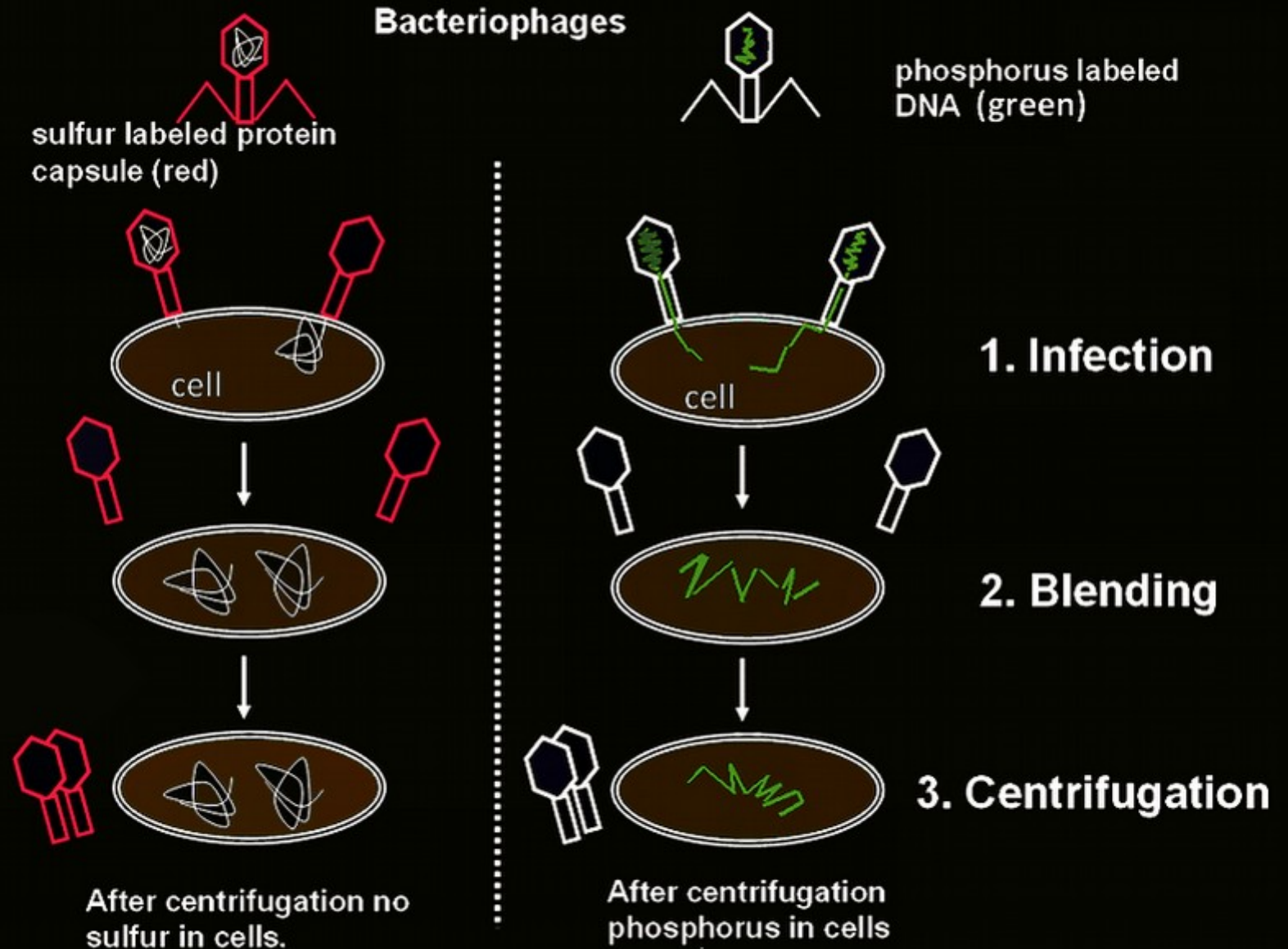
CHARGAFF (1950)

1970
-
-
-
-
1965
-
-
-
-
1960
-
-
-
-
1955
-
-
-
-
1950

1952

ON SAIT QUE... L'ADN PORTE L'INFORMATION GÉNÉTIQUE.

DÉBAT CLOS PAR HERSHEY ET CHASE (1952)



1952

ON SAIT QUE... PEU DE CHOSES SUR LA STRUCTURE DE L'ADN.

IL FAUDRA ATTENDRE :

- **1 an avant de connaître la structure de l'ADN;**
(Watson et Crick 1953)
- **13 ans avant d'en déchiffrer le code;**
(Nirenberg, Leder et al. 1965)
- **25 ans pour séquencer une molécule d'ADN.**
(Gilbert et al. 1976 ; Sanger et al. 1977)

1970

-

-

-

-

1965

-

-

-

-

1960

-

-

-

-

1955

-

-

-

-

1950



1955

MAIS NOUS EN SAVIONS BEAUCOUP PLUS SUR... LES PROTÉINES.

1970

-

-

-

-

1965

-

-

-

-

1960

-

-

-

-

1955

-

-

-

-

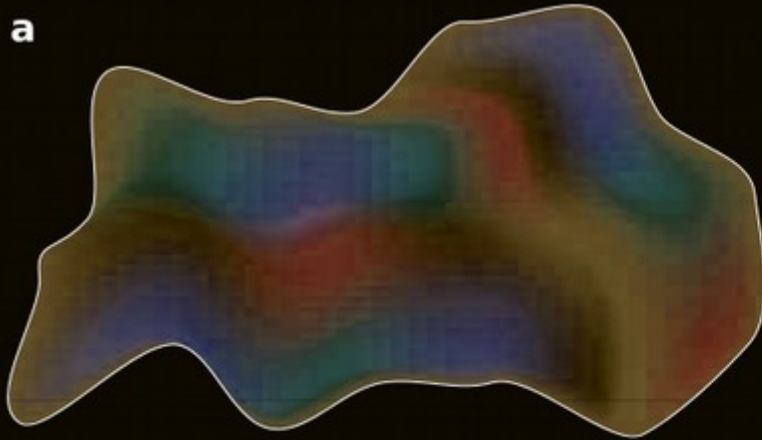
1950



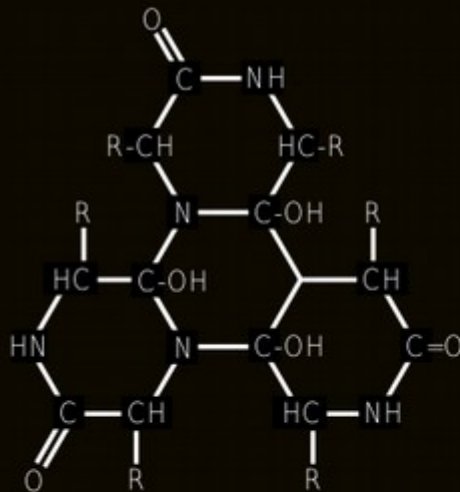
1955

MAIS NOUS EN SAVIONS BEAUCOUP PLUS SUR... LES PROTÉINES.

a



b



c



AVANT 1955:

**TROIS MODÈLES
EN COMPÉTITION**

a) Protéines = colloides amorphes sans structure moléculaire définie.

b) Structure polycyclique.

c) Protéines = simples chaînes répétitives d'acides aminés.

1955

MAIS NOUS EN SAVIONS BEAUCOUP PLUS SUR... LES PROTÉINES.

Découverte de la structure primaire de l'insuline
(Sanger et al. 1955)

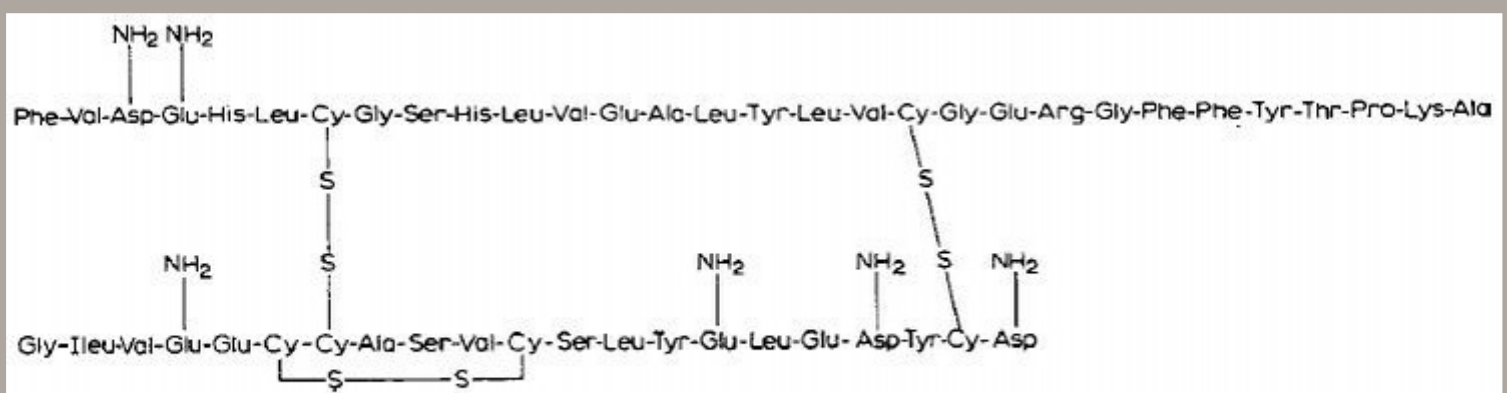


Fig. 2. The structure of insulin.

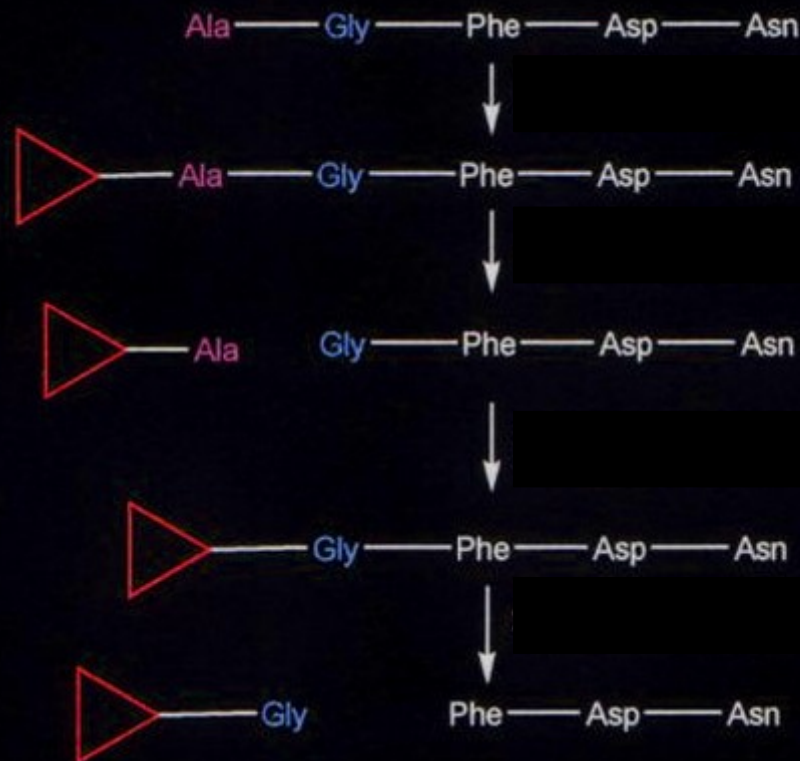
DÉBAT SUR LA NATURE POLYPEPTIDIQUE DES PROTÉINES: CLOS.

1956-1960

OPTIMISATION DU SÉQUENÇAGE DES PROTÉINES.

RÉACTION DE DÉGRADATION D'EDMAN

= Séquençage à 1 a.a. à la fois, à partir du côté N-terminal.

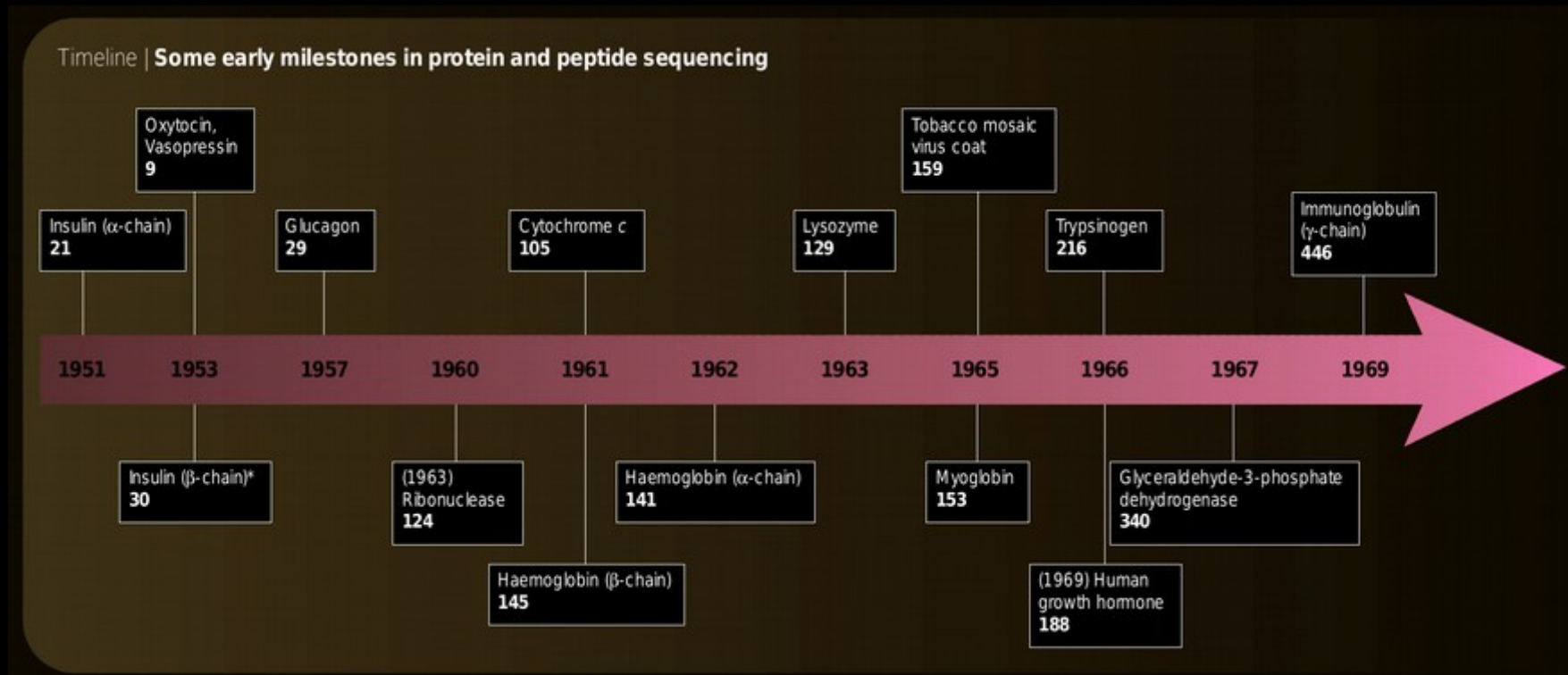


... MAINTENANT SEMI-AUTOMATISÉE !

1960-1965

OPTIMISATION DU SÉQUENÇAGE DES PROTÉINES.

Dégradation d'Edman + Automatisation =

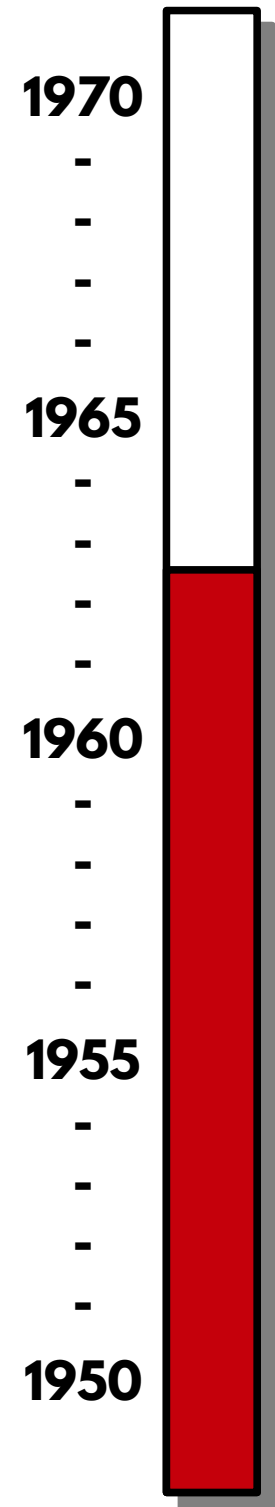


PUBLICATION MASSIVE DE SÉQUENCES DE PROTÉINES.

1960-1965

L'APOGÉE DU “PARADIGME DES PROTÉINES”

STRUCTURE → FONCTION



1960-1965

L'APOGÉE DU “PARADIGME DES PROTÉINES”

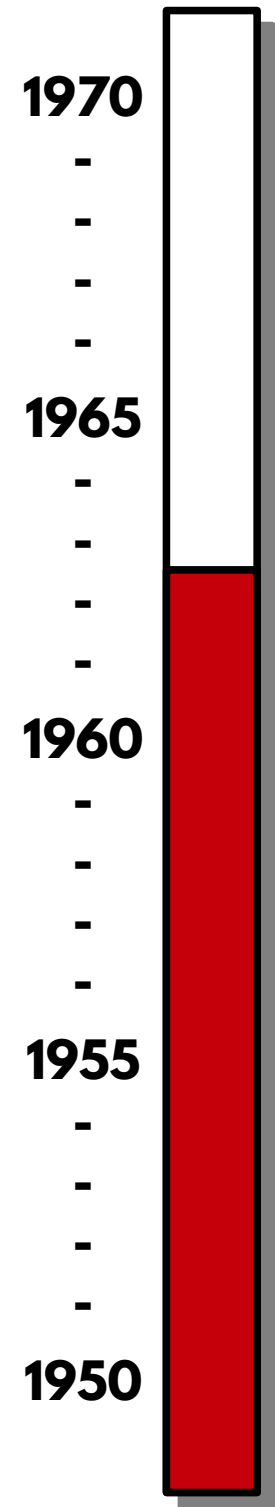
SÉQUENCE



STRUCTURE



FONCTION

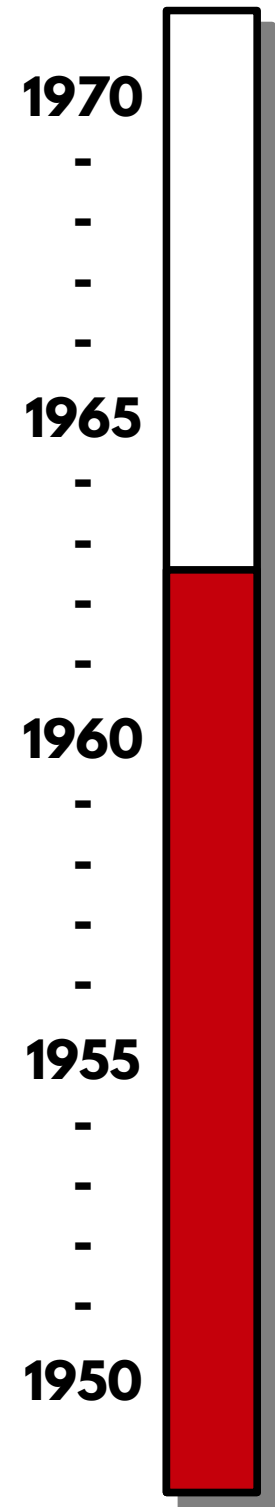


1960-1965

L'APOGÉE DU “PARADIGME DES PROTÉINES”

SÉQUENCE → **STRUCTURE** → **FONCTION**

AUTRE CONSTAT INTÉRESSANT...



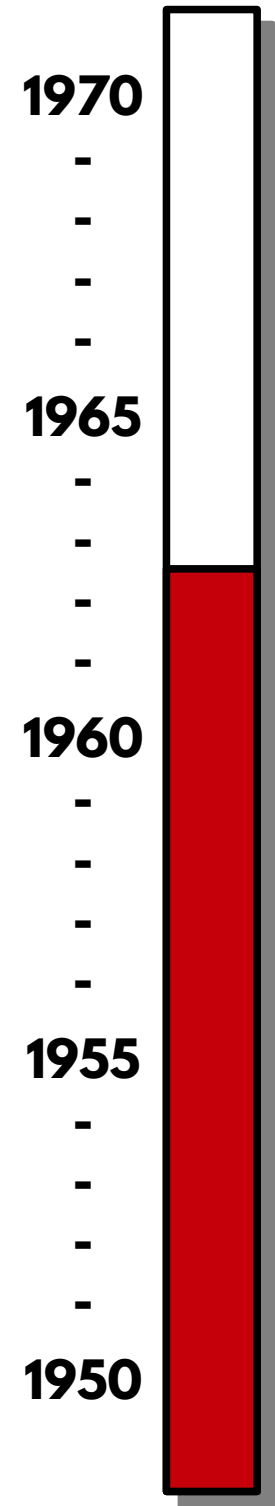
1960-1965

L'APOGÉE DU “PARADIGME DES PROTÉINES”

SÉQUENCE → **STRUCTURE** → **FONCTION**

AUTRE CONSTAT INTÉRESSANT...

Humain >MYOGLOBINE
 NH_2 -MIKTHECATATETHERATINTHEHAT...



1960-1965

L'APOGÉE DU "PARADIGME DES PROTÉINES"

SÉQUENCE → **STRUCTURE** → **FONCTION**

AUTRE CONSTAT INTÉRESSANT...

>MYOGLOBINE

Humain NH_2 -MIKTHECATATETHERATINTHEHAT...

Gorille NH_2 -MIKTHDCATATETHERGTINTHEHAT...

1970

-

-

-

-

1965

-

-

-

-

1960

-

-

-

-

1955

-

-

-

-

1950

1960-1965

L'APOGÉE DU "PARADIGME DES PROTÉINES"

SÉQUENCE → **STRUCTURE** → **FONCTION**

AUTRE CONSTAT INTÉRESSANT...

>MYOGLOBINE

Humain NH_2 -MIKTHECATATETHERATIN THEHAT...

Gorille NH_2 -MIKTHDCATATETHERGTIN THEHAT...

Chat NH_2 -MIKTHDCATATD THERGTIN THEHAT...

1960-1965

L'APOGÉE DU "PARADIGME DES PROTÉINES"

SÉQUENCE → **STRUCTURE** → **FONCTION**

AUTRE CONSTAT INTÉRESSANT...

>MYOGLOBINE

Humain	NH ₂ -MIKTHECATATETHERATIN THEHAT...
Gorille	NH ₂ -MIKTHDCATATETHERGTIN THEHAT...
Chat	NH ₂ -MIKTHDCATATD THERGTIN THEHAT...
Fourmi	NH ₂ -MLKTHDCATGT DTH-RGTIQTH--AT...

1960-1965

L'APOGÉE DU "PARADIGME DES PROTÉINES"

SÉQUENCE → STRUCTURE → FONCTION

AUTRE CONSTAT INTÉRESSANT...

>MYOGLOBINE

Humain	NH ₂ -MIKTHECATATETHERATIN THEHAT...
Gorille	NH ₂ -MIKTHDCATATETHERGTIN THEHAT...
Chat	NH ₂ -MIKTHDCATATD THERGTIN THEHAT...
Fourmi	NH ₂ -MLKTHDCATGT DTH-RGTIQTH--AT...
Pieuvre	NH ₂ -MLKTHDTGTGT DTH-RGTIQTH--AS...

1960-1965

L'APOGÉE DU "PARADIGME DES PROTÉINES"

SÉQUENCE → STRUCTURE → FONCTION

AUTRE CONSTAT INTÉRESSANT...

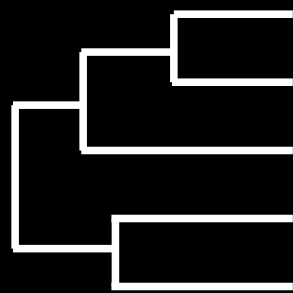
		>MYOGLOBINE
	Humain	NH ₂ -MIKTHECATATETHERATINTHEHAT...
	Gorille	NH ₂ -MIKTHDCATATETHERGTINTHEHAT...
	Chat	NH ₂ -MIKTHDCATATDTHERTGTINTHEHAT...
	Fourmi	NH ₂ -MLKTHDCATGTDTHT-RGTIQTH--AT...
	Pieuvre	NH ₂ -MLKTHDTGTGTDTHT-RGTIQTH--AS...

1960-1965

L'APOGÉE DU "PARADIGME DES PROTÉINES"

SÉQUENCE → STRUCTURE → FONCTION

AUTRE CONSTAT INTÉRESSANT...



Humain

Gorille

Chat

Fourmi

Pieuvre

>MYOGLOBINE

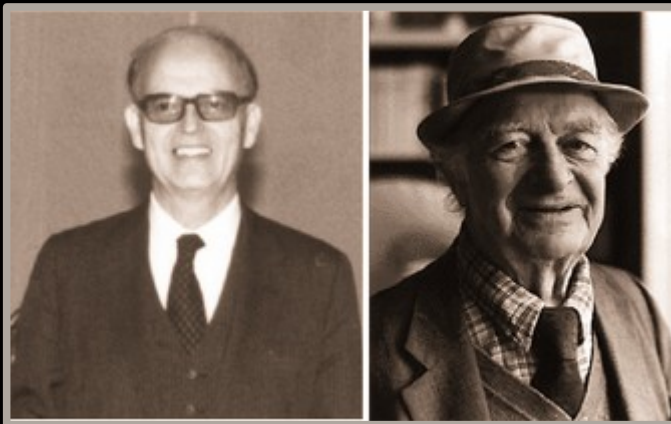
NH₂-MIKTHECATATETHERATIN THEHAT...

NH₂-MIKTHDCATATETHERGTIN THEHAT...

NH₂-MIKTHDCATATD THERGTIN THEHAT...

NH₂-MLKTHDCATGT DTH-RGTIQTH--AT...

NH₂-MLKTHDTGTGT DTH-RGTIQTH--AS...



Emile
ZUCKERKANDL

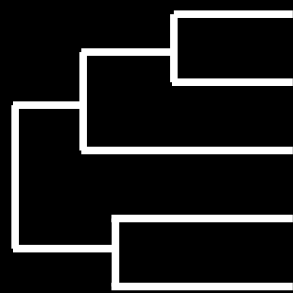
Linus
PAULING

1960-1965

L'APOGÉE DU "PARADIGME DES PROTÉINES"

SÉQUENCE → STRUCTURE → FONCTION

AUTRE CONSTAT INTÉRESSANT...



Humain

Gorille

Chat

Fourmi

Pieuvre

>MYOGLOBINE

NH₂-MIKTHECATATETHERATIN THEHAT...

NH₂-MIKTHDCATATETHERGTIN THEHAT...

NH₂-MIKTHDCATATD THERGTIN THEHAT...

NH₂-MLKTHDCATGT DTH-RGTIQTH--AT...

NH₂-MLKTHDTGTGT DTH-RGTIQTH--AS...



Emile
ZUCKERKANDL



Linus
PAULING

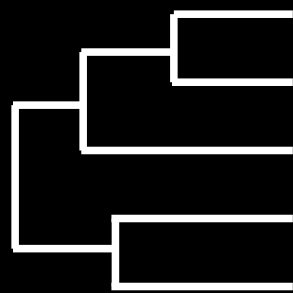
ancestry™

1960-1965

L'APOGÉE DU "PARADIGME DES PROTÉINES"

SÉQUENCE → STRUCTURE → FONCTION

AUTRE CONSTAT INTÉRESSANT...



Humain

Gorille

Chat

Fourmi

Pieuvre

>MYOGLOBINE

NH₂-MIKTHECATATETHERATIN THEHAT...

NH₂-MIKTHDCATATETHERGTIN THEHAT...

NH₂-MIKTHDCATATD THERGTIN THEHAT...

NH₂-MLKTHDCATGT DTH-RGTIQTH--AT...

NH₂-MLKTHDTGTGT DTH-RGTIQTH--AS...



Emile
ZUCKERKANDL



Linus
PAULING

Divergence entre
séquences orthologues

~ =

Histoire évolutive
des protéines
(et des espèces) ?

1960-1965

L'APOGÉE DU “PARADIGME DES PROTÉINES”

DÉBUT DE LA “ PALÉOGÉNÉTIQUE ”

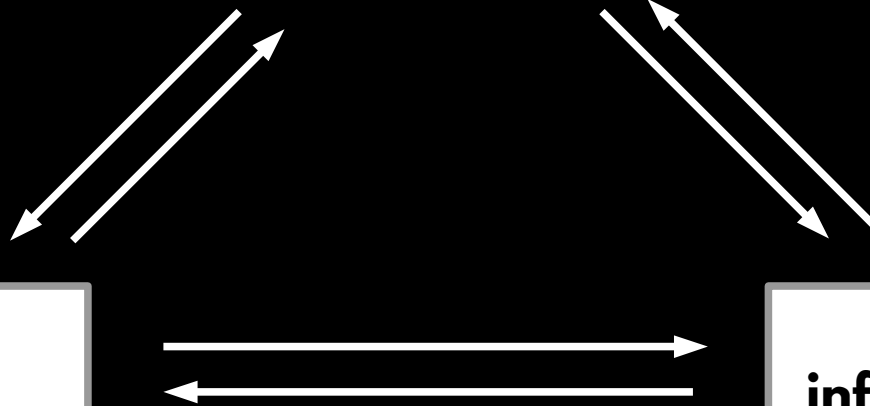
Synthèse entre...

Biochimie des protéines

**Biologie
évolutive**

**Sciences
informatiques**

1970
-
-
-
-
1965
-
-
-
-
1960
-
-
-
-
1955
-
-
-
-
1950



1960-1965

L'APOGÉE DU “PARADIGME DES PROTÉINES”

DÉBUT DE LA “ PALÉOGÉNÉTIQUE ”

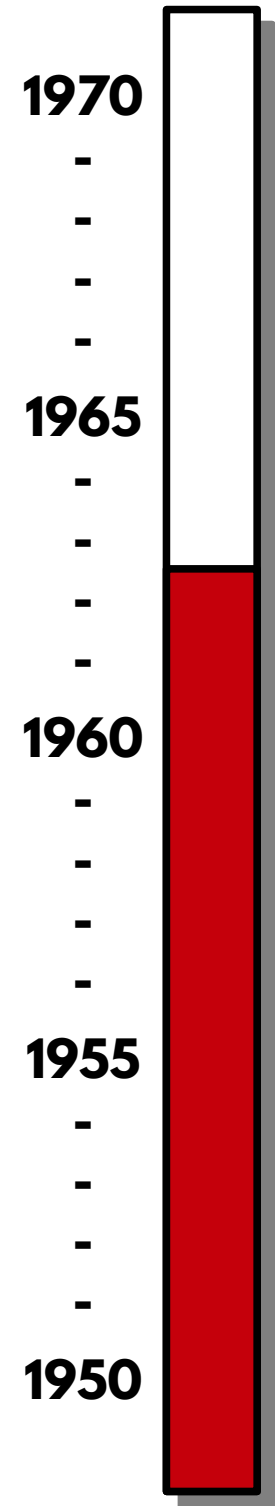
Synthèse entre...

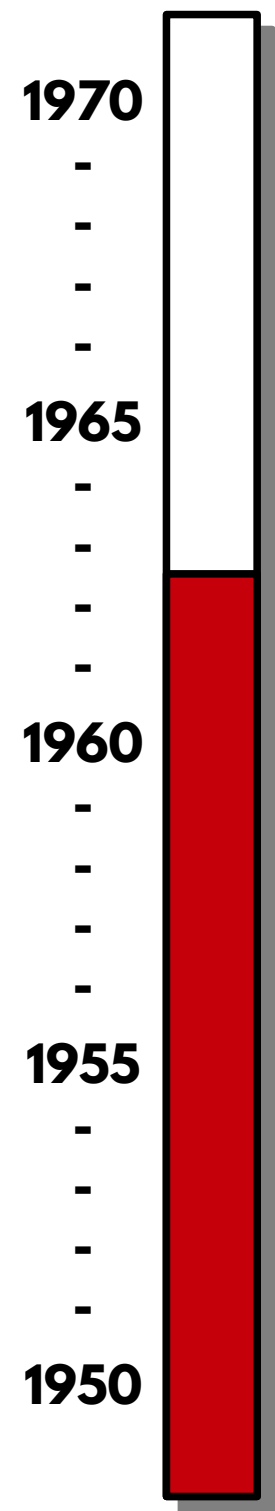
Biochimie des protéines

**Biologie
évolutive**

**Sciences
informatiques**

AU COEUR DE CETTE CHIMIE...





... SE TROUVE UNE FEMME.

1960-1965

GÉNÈSE DE LA “ BIOLOGIE COMPUTATIONNELLE ”

MARGARET OAKLEY DAYHOFF

(1925-1983)

“ ... the mother and father of bioinformatics ”
- David Lipman (NCBI)



- Directrice adjointe du National Biomedical Research Foundation (NBRF)
- A utilisé de manière extensive l'informatique lors de son Ph.D. en électrochimie.
- Voyait le potentiel de l'informatique en **PALÉOGÉNÉTIQUE**.

1970

-

-

-

-

1965

-

-

-

-

1960

-

-

-

-

1955

-

-

-

-

1950

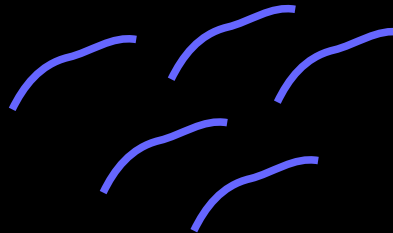
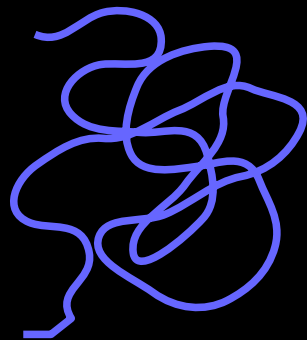
1960-1965

GÉNÈSE DE LA “ BIOLOGIE COMPUTATIONNELLE ”

PRINCIPALES CONTRIBUTIONS DE M.O. DAYHOFF

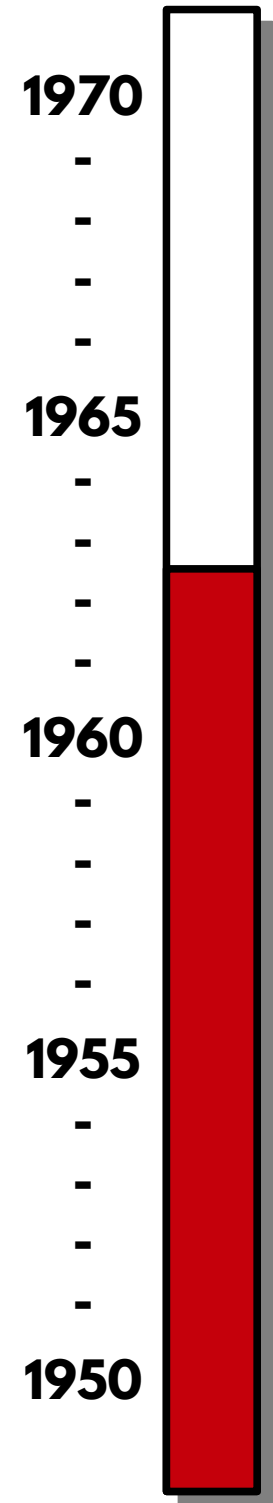
DÉGRADATION D'EDMAN: *PROBLÈME*

- Rendement: 98% = 50 a.a. MAX
- Une protéine > 50 a.a. doit être dégradée en **plus petits peptides**, qui eux seront séquencés.



KTHECAT
ATATETH
ATET
ERAT
HERA
MIK
...

Séquence de la protéine complète ?



1960-1965

GÉNÈSE DE LA “ BIOLOGIE COMPUTATIONNELLE ”

PRINCIPALES CONTRIBUTIONS DE M.O. DAYHOFF

1- Le premier “ outil bioinformatique ” (1962)

COMPROTEIN: A COMPUTER PROGRAM TO AID
PRIMARY PROTEIN STRUCTURE DETERMINATION*

*Margaret Oakley Dayhoff and Robert S. Ledley
National Biomedical Research Foundation
Silver Spring, Maryland*

KTHECAT
ATATETH
ATET
ERAT
HERA
MIK



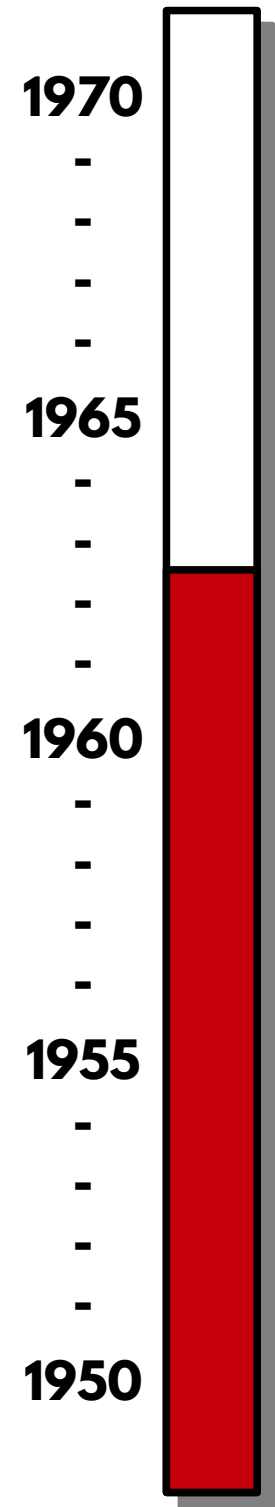
KTHECAT	Peptide
ATATETH	Peptide
ATET	Peptide
ERAT	Peptide
HERA	Peptide
MIK	Peptide

MIKTHECATATETHERAT	Protein

1960-1965

GÉNÈSE DE LA “ BIOLOGIE COMPUTATIONNELLE ”

FACILE ?



1960-1965

GÉNÈSE DE LA “ BIOLOGIE COMPUTATIONNELLE ”

FACILE ?



Ordinateur typique de l'époque (IBM 7090)



1970

-

-

-

-

1965

-

-

-

-

1960

-

-

-

-

1955

-

-

-

-

1950

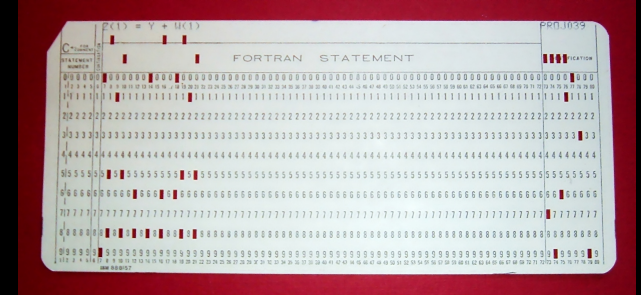
1960-1965

GÉNÈSE DE LA “ BIOLOGIE COMPUTATIONNELLE ”

FACILE ?



Ordinateur typique de l'époque (IBM 7090)



Une ligne de code en 1962 (FORTRAN)



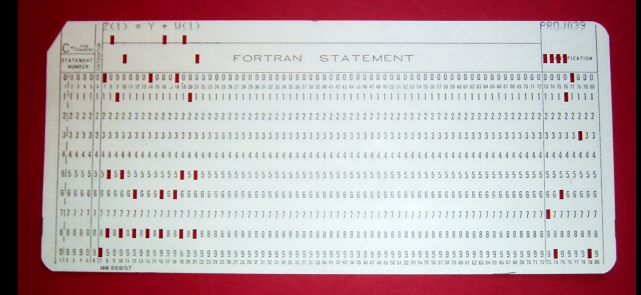
1960-1965

GÉNÈSE DE LA “ BIOLOGIE COMPUTATIONNELLE ”

FACILE ?



Ordinateur typique de l'époque (IBM 7090)



Une ligne de code en 1962 (FORTRAN)



Code source d'un programme

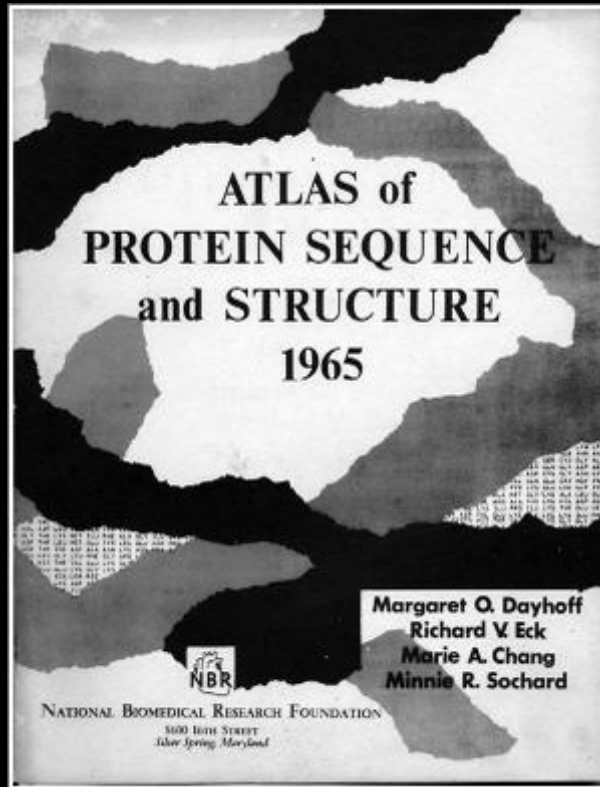


1965

GÉNÈSE DE LA “ BIOLOGIE COMPUTATIONNELLE ”

PRINCIPALES CONTRIBUTIONS DE M.O. DAYHOFF

2- La première “base de données” de séquences de protéines...



BOVINE GROWTH HORMONE

	5	10	15	20	25	30
1	A	F	P	A	M	S
31	E	F	E	R	T	Y
61	P	T	G	K	N	E
91	Q	F	L	S	R	V
121	L	A	L	M	R	E
151	S	D	D	A	L	L
181	C	R	R	F	G	E

COMPOSITION

15 ALA A	11 GLN Q	27 LEU L	13 SER S
13 ARG R	13 GLU E	11 LYS K	12 THR T
6 ASN N	10 GLY G	4 MET M	1 TRP W
10 ASP D	3 HIS H	13 PHE F	6 TYR Y
4 CYS C	7 ILE I	6 PRO P	6 VAL V

MOL. WT. = 21,816

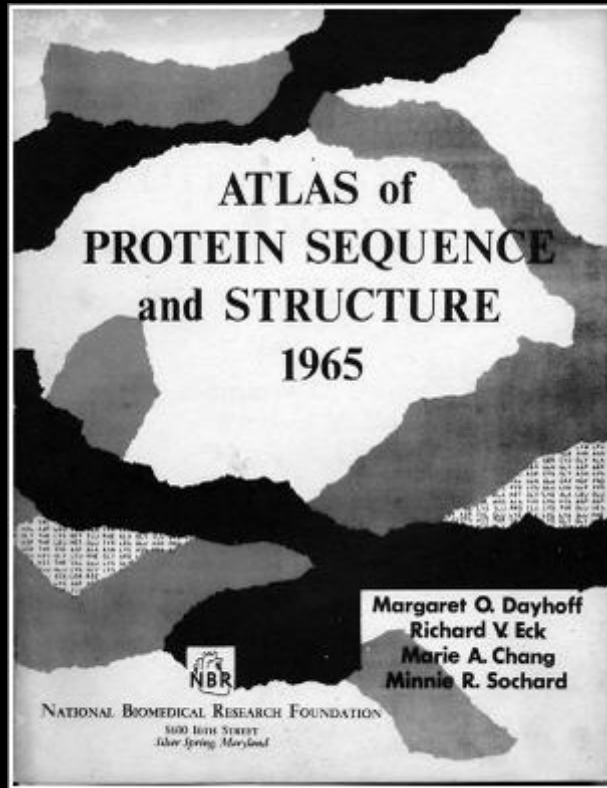
NUMBER OF RESIDUES = 191

1965

GÉNÈSE DE LA “ BIOLOGIE COMPUTATIONNELLE ”

PRINCIPALES CONTRIBUTIONS DE M.O. DAYHOFF

2- La première “base de données” de séquences de protéines...



Au début...

BEAUCOUP de variants interspécifiques
d'une **DIZAIN**e de protéines.

>PROT1	
Humain	NH ₂ -MIKTHECATATETHERATINHE...-COOH
Gorille	NH ₂ -MIKTHDCATATETHERGTINHE...-COOH
Chat	NH ₂ -MIKTHDCATATDTHRGITINHE...-COOH
Fourmi	NH ₂ -MLKTHDCATGTDTHRGITQTHAT...-COOH
Pieuvre	NH ₂ -MLKTHDTGTGTDTHRGITQTHAS...-COOH
...	
S. cerev	NH ₂ -MLKNHDTGTGTDTHRGITQTHAS...-COOH
E. coli	NH ₂ -MLKTHNDTGTGARQTGTIQTTHAS...-COOH

1965

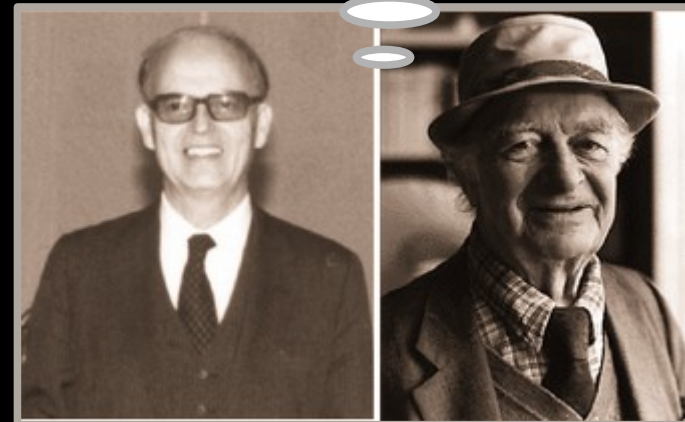
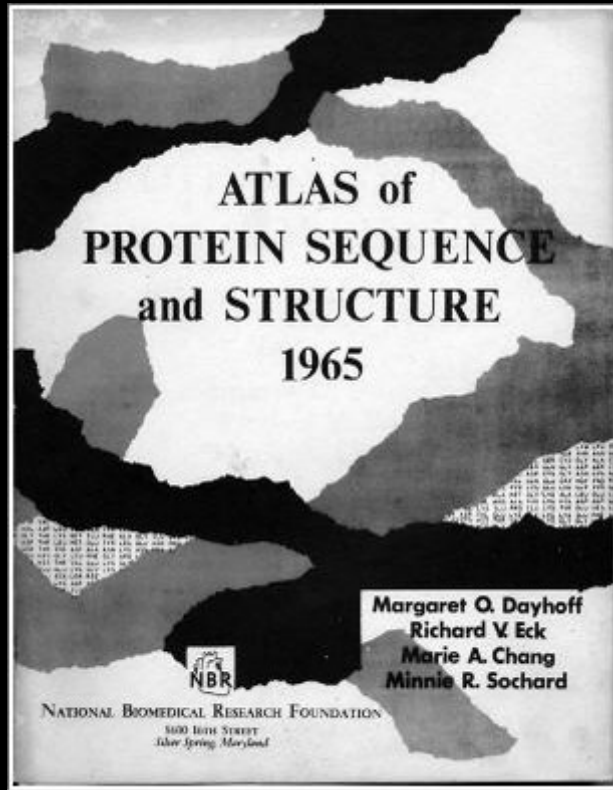
GÉNÈSE DE LA “ BIOLOGIE COMPUTATIONNELLE ”

PRINCIPALES CONTRIBUTIONS DE M.O. DAYHOFF

2- La première “base de données” de séquences de protéines...

Données idéales pour...

ancestry™



1966

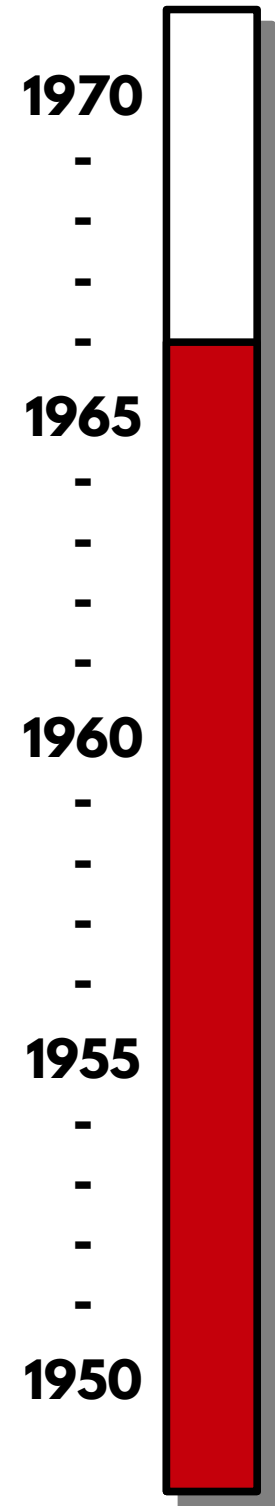
DÉBUTS DE LA PHYLOGÉNIE MOLÉCULAIRE

Exemple...

Protéine 1 AVQH

Protéine 2 ALQH

Protéine 3 AIQH



1966

DÉBUTS DE LA PHYLOGÉNIE MOLÉCULAIRE

Exemple...

Protéine 1 AVQH

Protéine 2 ALQH

Protéine 3 AIQH

Phylogénie ? Trois possibilités :



Quelle hypothèse est la bonne?

Vérifions la distance entre les séquences
À L'AIDE D'ALIGNEMENTS.

1970
-
-
-
-
1965
-
-
-
-
1960
-
-
-
-
1955
-
-
-
-
1950

1966

DÉBUTS DE LA PHYLOGÉNIE MOLÉCULAIRE

PROBLÈME :

Tous ces alignements sont
"ÉQUIVALENTS"

Prot1 AVQH

: ::

Prot2 ALQH

Prot3 AIQH

: ::

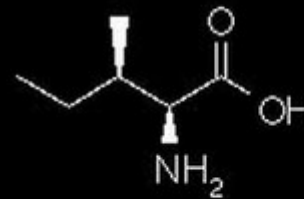
Prot2 ALQH

Prot3 AIQH

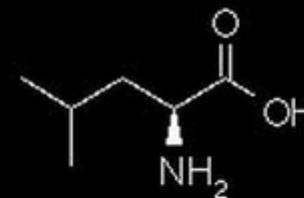
: ::

Prot1 AVQH

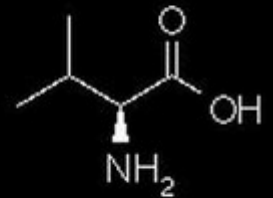
Iso-Leucine



Leucine



Valine

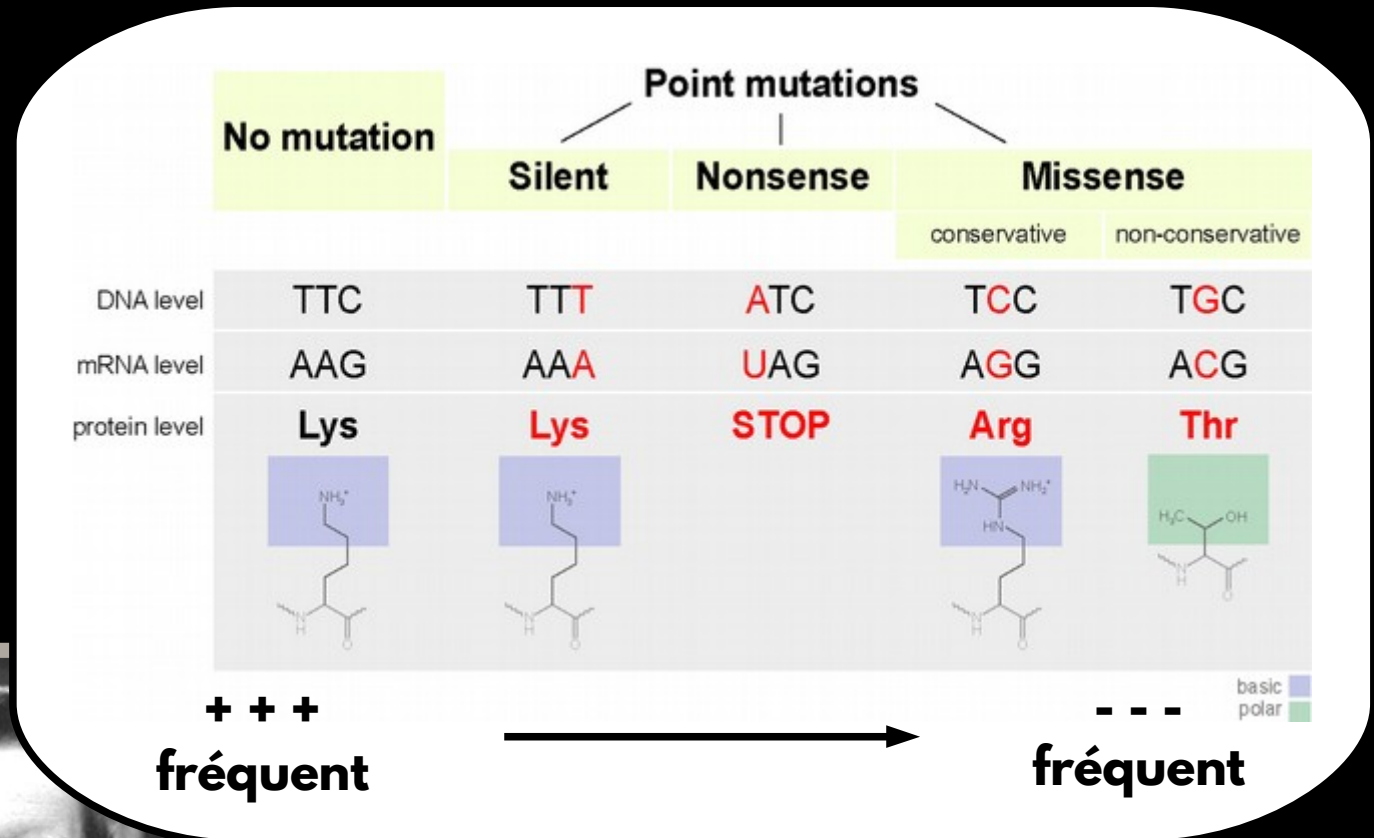


A.A. similaires, mais non-identiques.

**Comment alors quantifier
la distance entre
Prot1, Prot2 et Prot3 ?**

1966

DÉBUTS DE LA PHYLOGÉNIE MOLÉCULAIRE



1966

DÉBUTS DE LA PHYLOGÉNIE MOLÉCULAIRE

Vu que...

Séquence => Structure => Fonction

Certaines mutations ponctuelles
ont plus de chance de
se produire que d'autres.



1970

-

-

-

-

1965

-

-

-

-

1960

-

-

-

-

1955

-

-

-

-

1950

1966

DÉBUTS DE LA PHYLOGÉNIE MOLÉCULAIRE

CECI IMPLIQUE QUE...

**Plus deux séquences comportent
de substitutions rares, plus elles
sont dissimilaires.**

**= POSSIBLE de discriminer des
ALIGNEMENTS ÉQUIVALENTS.**



1970

-

-

-

-

1965

-

-

-

-

1960

-

-

-

-

1955

-

-

-

-

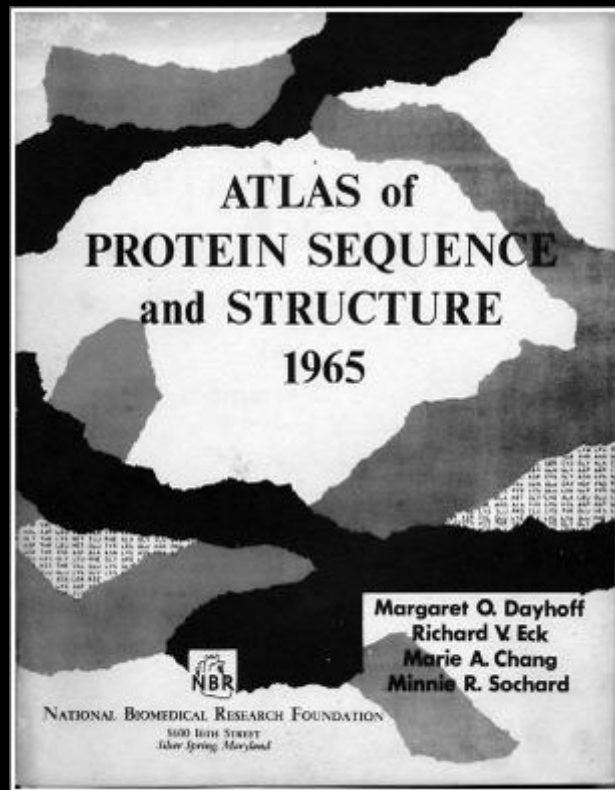
1950

1966

DÉBUTS DE LA PHYLOGÉNIE MOLÉCULAIRE

PRINCIPALES CONTRIBUTIONS DE M.O. DAYHOFF

3- MATRICES DE SUBSTITUTION



PAM (complétée en 1978!)
À partir de **1572** substitutions observées
dans **71** familles de protéines
partagent **>85% d'identité**

a.a. original

	A Ala	R Arg	N Asn	D Asp	C Cys
A	9867	2	9	10	3
R	1	9913	1	0	1
N	4	1	9822	36	0
D	6	0	42	9859	0
C	1	1	0	0	9
Q	3	9	4	5	0

a.a. substitut

$10^4 \times P$ (Arg > Cys)



1966

DÉBUTS DE LA PHYLOGÉNIE MOLÉCULAIRE

RETOURNONS À NOS TROIS PROTÉINES...
Calcul des probabilités de substitution avec PAM1

Prot1 AVQH

: ::

Prot2 ALQH

Prot3 AIQH

: ::

Prot2 ALQH

Prot3 AIQH

: ::

Prot1 AVQH

1970

-

-

-

-

1965

-

-

-

-

1960

-

-

-

-

1955

-

-

-

-

1950

1966

DÉBUTS DE LA PHYLOGÉNIE MOLÉCULAIRE

RETOURNONS À NOS TROIS PROTÉINES...
Calcul des probabilités de substitution avec PAM1

Prot1 AVQH

: : :

Prot2 ALQH

Prot3 AIQH

: : :

Prot2 ALQH

Prot3 AIQH

: : :

Prot1 AVQH

$$P(V > L) = 0.0015$$

$$P(L > V) = 0.0011$$

$$P(I > L) = 0.0022$$

$$P(L > I) = 0.0009$$

$$P(I > V) = 0.0057$$

$$P(V > I) = 0.0033$$

1966

DÉBUTS DE LA PHYLOGÉNIE MOLÉCULAIRE

RETOURNONS À NOS TROIS PROTÉINES...
Calcul des probabilités de substitution avec PAM1

Prot1 AVQH

: : :

Prot2 ALQH

Prot3 AIQH

: : :

Prot2 ALQH

Prot3 AIQH

: : :

Prot1 AVQH

$$P(V > L) = 0.0015$$

$$P(L > V) = 0.0011$$

$$P(I > L) = 0.0022$$

$$P(L > I) = 0.0009$$

$$P(I > V) = 0.0057$$

$$P(V > I) = 0.0033$$

1966

DÉBUTS DE LA PHYLOGÉNIE MOLÉCULAIRE

RETOURNONS À NOS TROIS PROTÉINES...
Calcul des probabilités de substitution avec PAM1

Prot1 AVQH

: : :

Prot2 ALQH

Prot3 AIQH

: : :

Prot2 ALQH

Prot3 AIQH

: : :

Prot1 AVQH

$$P(V > L) = 0.0015$$

$$P(L > V) = 0.0011$$

$$P(I > L) = 0.0022$$

$$P(L > I) = 0.0009$$

$$P(I > V) = 0.0057$$

$$P(V > I) = 0.0033$$

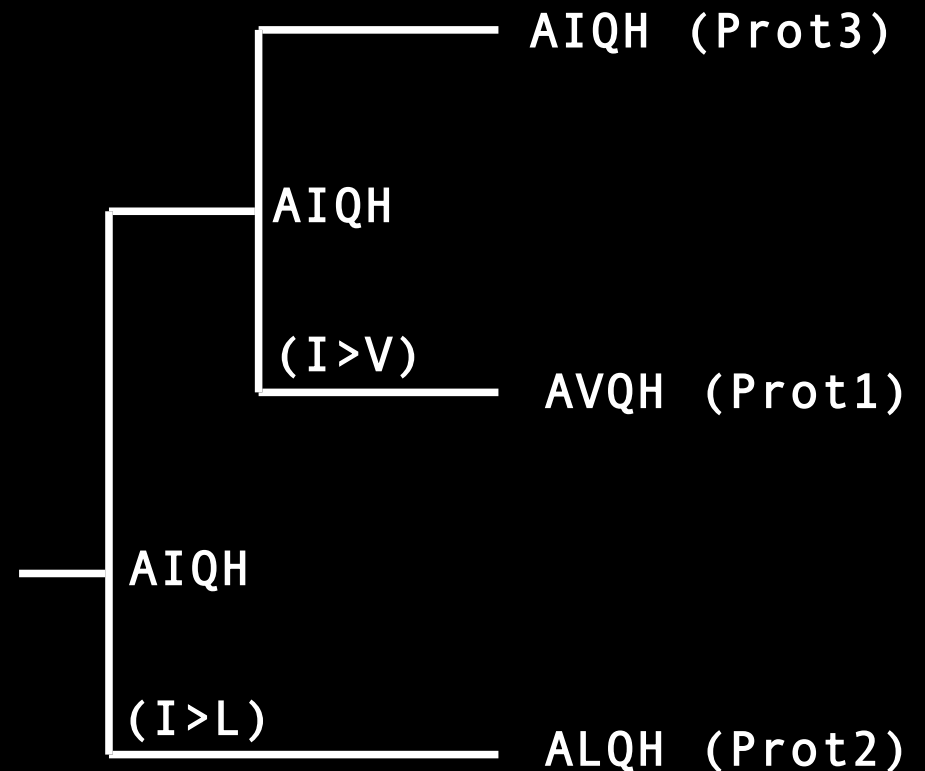
LES PLUS
SIMILAIRES

1966

DÉBUTS DE LA PHYLOGÉNIE MOLÉCULAIRE

RETOURNONS À NOS TROIS PROTÉINES...
Calcul des probabilités de substitution avec PAM1

Prot1	AVQH
	: ::
Prot2	ALQH
	: ::
Prot3	AIQH
	: ::
Prot2	ALQH
	: ::
Prot3	AIQH
	: ::
Prot1	AVQH



Hypothèse la plus PROBABLE





DE RETOUR EN 2015 ...



La session prochaine...

PARTIE II (1971-1990)

- **MINIATURISATION DE L'INFORMATIQUE;**
- **ARRIVÉE DU SÉQUENÇAGE DE L'ADN;**
- **PREMIERS RÉSEAUX D'ORDINATEURS;**
- **DÉBUTS DE L'ÈRE GÉNOMIQUE.**

