

**Modélisation  
et génome**

**ou**

**Bioinformatique**

**ou**

**Biologie  
“Computationnelle”**

# IDE'E D'UNE ECHELLE

## DES ETRES NATURELS.

L'HOMME.
Orang-Outang.
Singe.
QUADRUPEDES.
Ecureuil volant.
Chauvesouris.
Autruche.
OISEAUX.
Oiseaux aquatiques.
Oiseaux amphibies.
Poissons volans.
POISSONS.
Poissons rampans.
Anguilles.
Serpens d'eau.
SERPENS.
Limaces.
Limaçons.
COQUILLAGES.
Vers à tuyau.
Teignes.
INSECTES.
Gallinectes.
Terma, ou Solitaire.
Polypes.

Orties de Mer.
Sensitive.
PLANTES.
Lychens.
Mouffures.
Champignons, Agarics.
Truffes.
Coraux & Coralloides.
Lithophytes.
Amianthe.
Talcs, Gyps, Sélénites.
Ardoises.
PIERRES.
Pierres figurées.
Cristallisations.
SELS.
Vitriols.
METAUX.
DEMI-METAUX.
SOUFRES.
Bitumes.
TERRES.
Terre pure.
EAU.
AIR.
FEU.
Matières plus subtiles.

Charles Bonnet, 1745

**ADN**

**ARN**

**protéine**

**assemblée supramoléculaire**

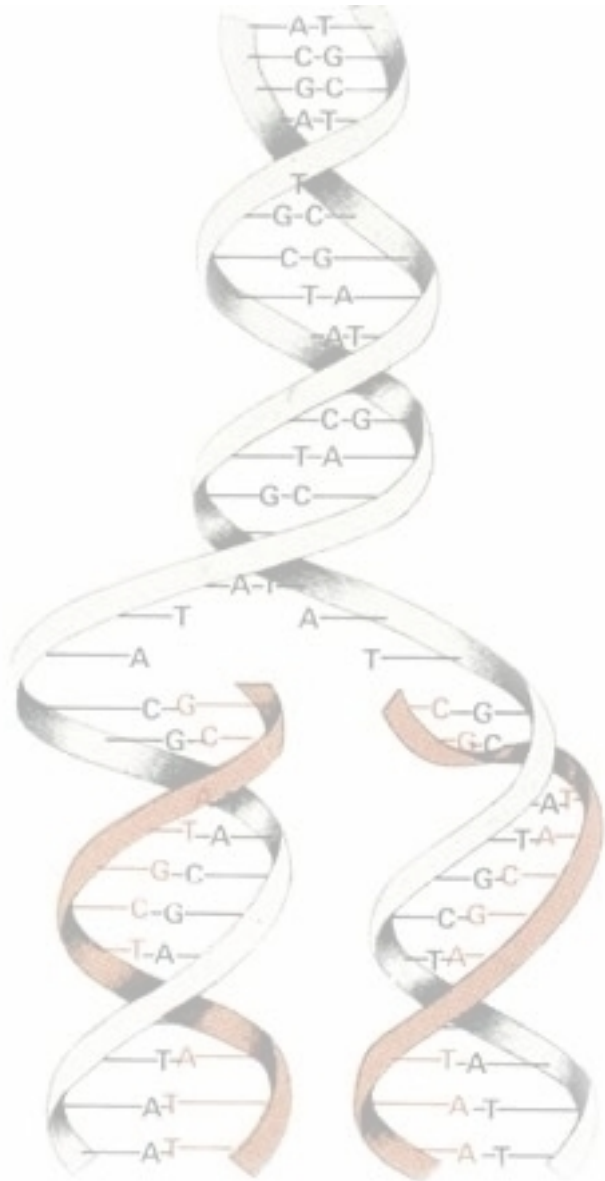
**réseaux et voies**

**cellule**

**tissu**

**organisme**

**population**



# Comprendre la structure et fonction des biomolécules

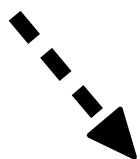
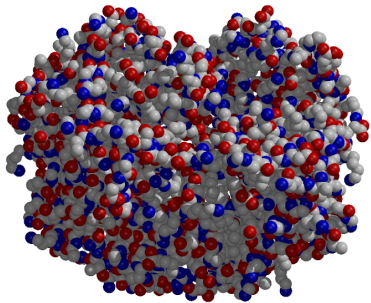
A  
U  
G  
C  
G  
C  
U  
U  
A  
U  
A  
G  
C  
C  
A  
A  
G  
G  
:



**Séquence**



**Structure**

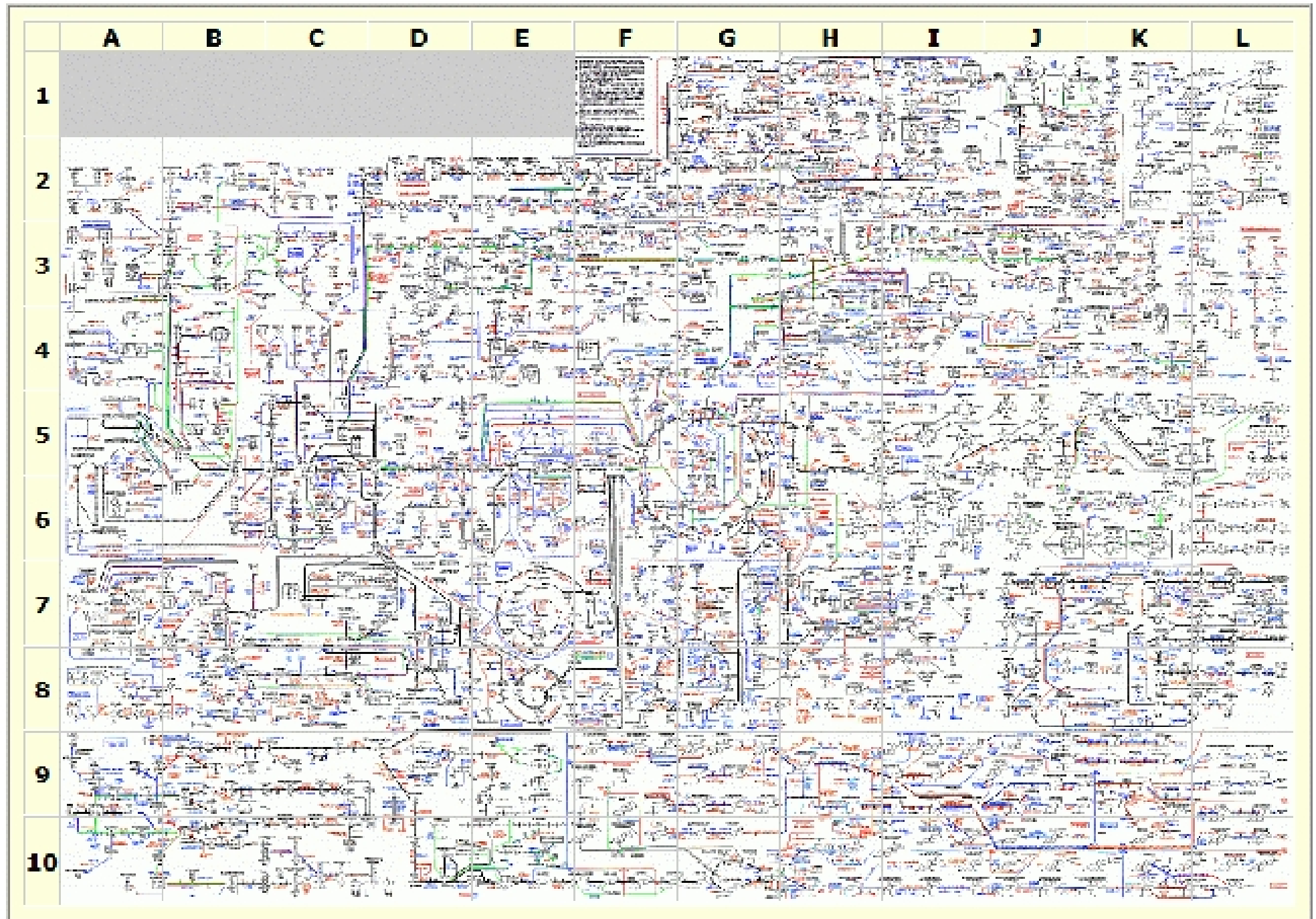


**Fonction**

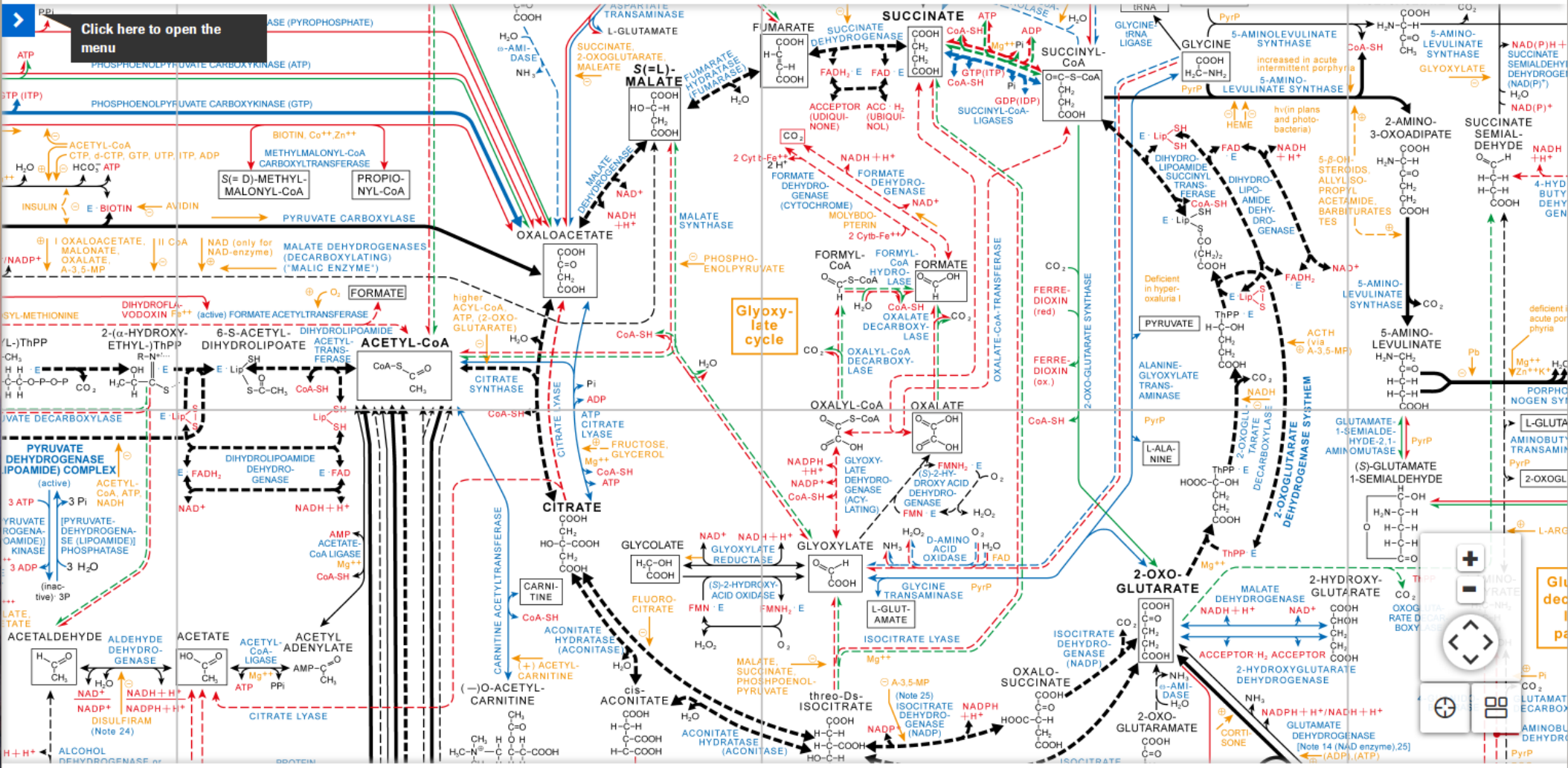


# Comprendre leur intégration dans la cellule (~10<sup>9</sup> macromolécules)

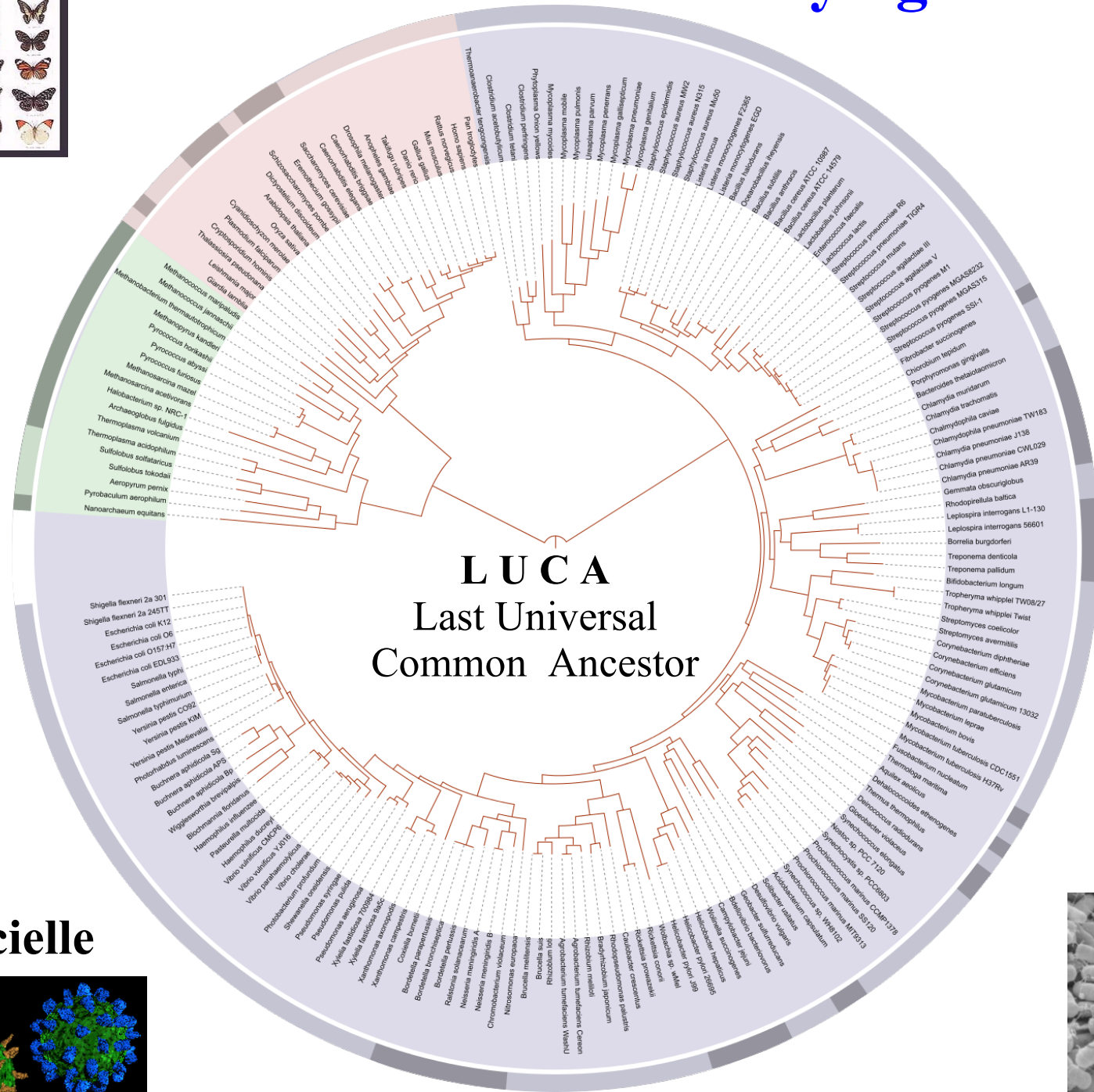
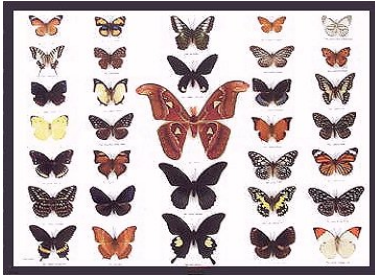
## Biologie “intégrative” ou “systémique” ([web.expasy.org/pathways](http://web.expasy.org/pathways))



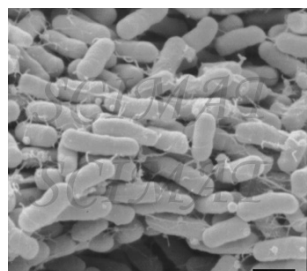
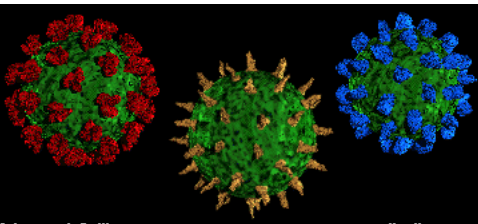
Part 1: Metabolic Pathways Part 2: Cellular and Molecular Processes



# Phylogénie et évolution

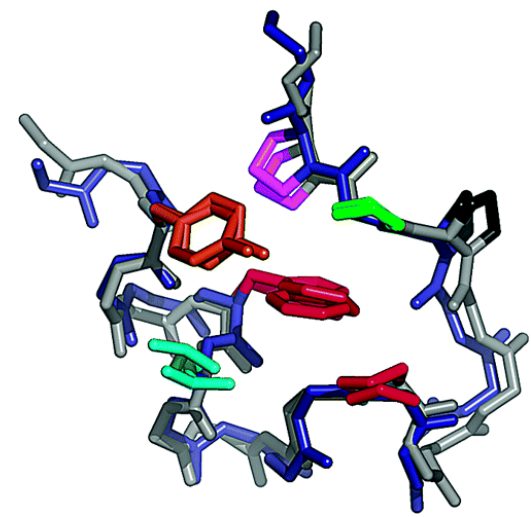
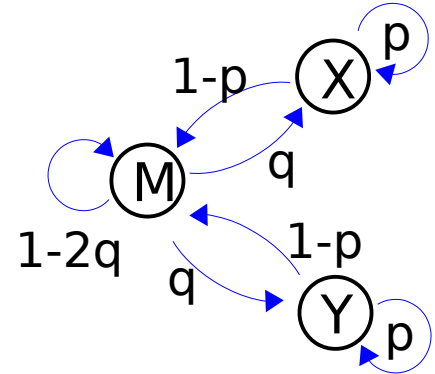
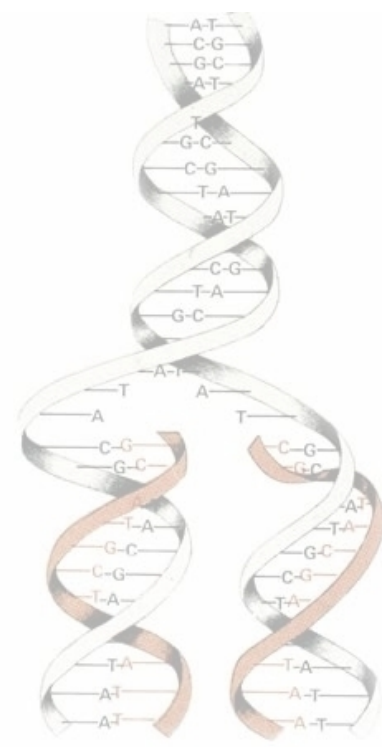
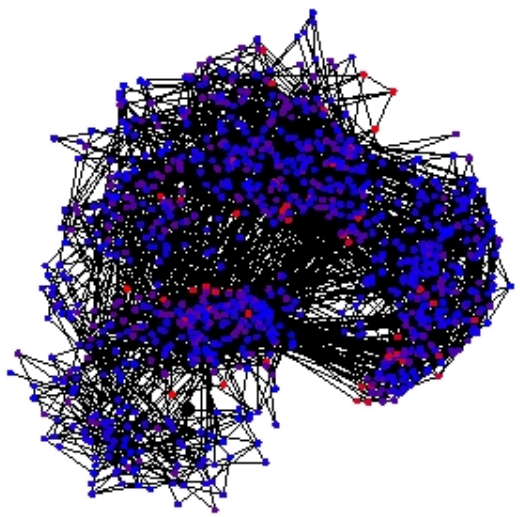


## Vie artificielle



# Les outils de l'informatique, des statistiques, de la chimie, de la physique

		193	195	198		217	ATP	ATP	231	233	235	
YstASP	293	.....	KGKAY.LAQS	PQFNKQQLIV	ADFERVYEIG	PVFR	AE	NSNT	HRHMTEFTGL	DME	MAFEEHY	HEVL 355
ThtASP	190	.....	PGLFYALPOS	PQLFKOQLMV	AGLDRYFQIA	KCFR	DE	DLRA	DRQP	DF	TQL	DLEMSF.VEV EDVL 251
EcoASP	184	.....	KGKFYALPOS	HQLFKQQLMM	SGFDRYYQIV	KCFR	DE	DLRA	DRQP	EF	TQI	DVETSF.MTA PQVR 245
EcoASN	191	QGVDFDKDF	FGKESFLTVS	GQLNGEYIAC	.ALSKIYTFG	PTFR	AE	NSNT	SRHLAEFWML	EPE	VAFAN.L	NDIA 262
YstASN	198	.NTSPTASSY	FGKPTYLTVS	TQLHLEILAL	.SLSRCWTL	PCFR	AE	KSDT	PRHLSEFWML	EVE	MCVNSV	NELT 269
ThtLYS	205	.PFKTYHNAL	DHEFY.LRIS	LELYLKRLLV	GGYEKVF	EIG	RNFR	NEGIDH	.NHNPEFTML	EAY	WAYAD.Y	QDMA 274
EcoLYS	221	.PFITHHNAL	DLDY.LRIA	PELYLKRLV	GGFERV	EIN	RNFR	NEGISV	.RHNPEFTMM	ELY	MAYAD.Y	KOLI 290
YstLYS	284	.PFITHHNDL	DMDY.MRIA	PELFLKQLV	GGLDRV	EIG	RQFR	NEGIDM	.THNPEFTTC	EFY	QAYAD.V	YDLM 353



# Organisation du cours

Thomas Simonson (Dépt Bio)  
thomas.simonson@polytechnique.fr

- 2/12 Comparaisons de séquences, I
- 9/12 Comparaisons de séquences, II
- 16/12 Protéines
- 6/1 Prédiction de structures d'ARN
- 13/1 Phylogénie (T Gaillard)
- 20/1 Analyse statistique du génome
- 27/1 Organisation et séquençage du génome
- 3/2 Réseaux cellulaires
- 10/2 Evolution moléculaire

# Petites classes

Amphi Curie

**Apportez vos ordinateurs pour les PC 2, 3**

thomas.simonson@polytechnique.edu

thomas.gaillard@polytechnique.edu

## Document écrit (chaps 1-7)

Cf aussi [biology.polytechnique.fr/biocomputing](http://biology.polytechnique.fr/biocomputing)  
cf. biblio en fin de cours

## Contrôle écrit

**Langue = français + anglais**

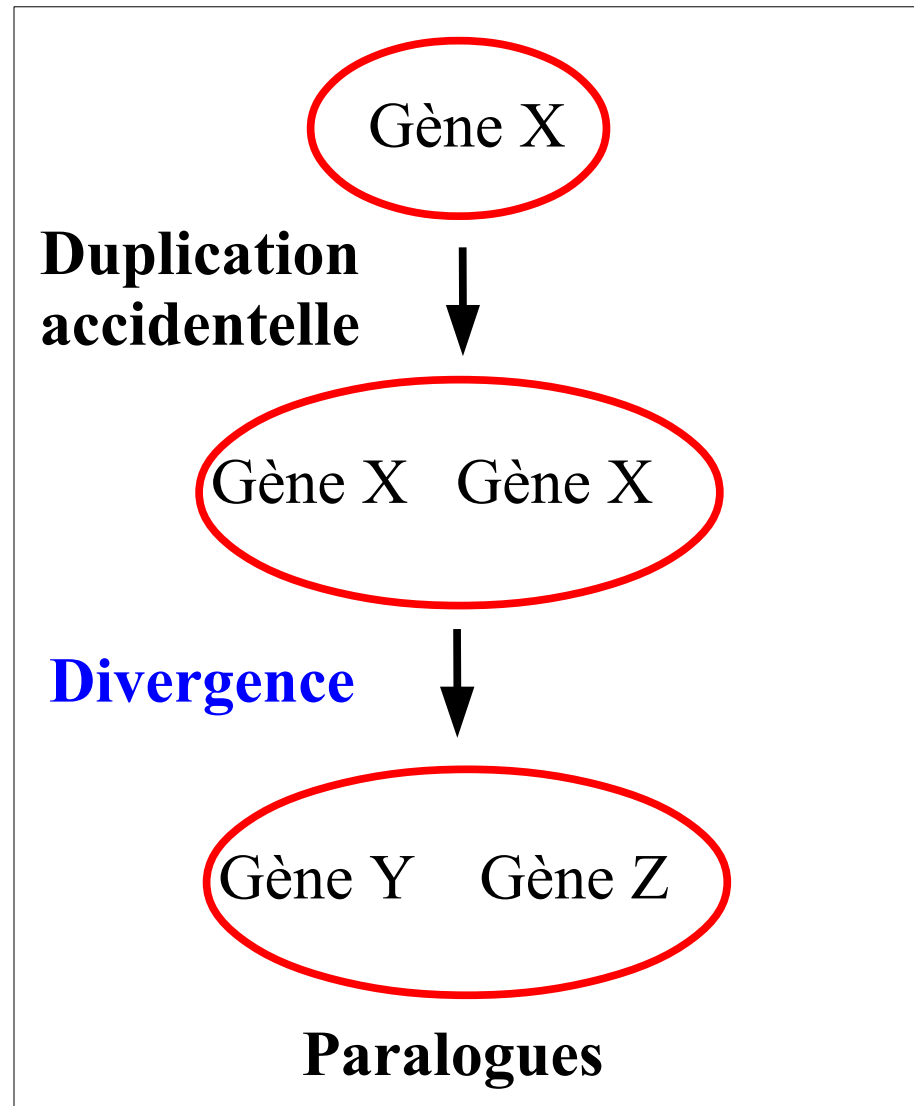
# Comparaison et alignement de séquences. I.

						motif 2								
			193	195	198		217	ATP	ATP	231	233	235		
YstASP	293	.....	KGKAY.LAQS	PQFNKQQLIV	ADFERVYEIG	PVFRAENSNT	HRHMTEFTGL	DMEMAFEEHY	HEVL	355				
ThtASP	190	.....	PGLFYALPQS	PQLFKQMLMV	AGLDRYFQIA	RCFRDEDLRA	.DRQPDFTQL	DLEMSF.VEV	EDVL	251				
EcoASP	184	.....	KGKFYALPQS	PQLFKQLLMM	SGFDRYYQIV	KCFRDEDLRA	.DRQPEFTQL	DVETSF.MTA	PQVR	245				
EcoASN	191	QGKVDFDKDF	FGKESFLTVS	GQLNGETYAC	.ALSKIYTFG	PTFRAENSNT	SRHLAEFWML	EPEVAFAN.L	NDIA	262				
YstASN	198	.NTSPTASSY	FGKPTYLTVS	TQLHLEILAL	.SLSRCWTLS	PCFRAEKSDT	PRHLSEFWML	EVEMCFVNSV	NELT	269				
ThtLYS	205	.PFKTYHNAL	DHEFY.LRIS	LELYLKRLLV	GGYEKVFYIG	RNFRNEGIDH	.NHNPEFTML	EAYWAYAD.Y	QDMA	274				
EcoLYS	221	.PFITHHNAL	DLDMY.LRIA	PELYLKRLVV	GGFERVFEIN	RNFRNEGISV	.RHNPEFTMM	ELYMAYAD.Y	KDLI	290				
YstLYS	284	.PFITHHNDL	DMDMY.MRIA	PELFLKQLVV	GGLDRVYEIG	RQFRNEGIDM	.THNPEFTTC	EFYQAYAD.V	YDLM	353				

- **Comparaisons de séquences: pourquoi faire?**
- **L'alignement de séquences comme modèle d'un processus évolutif**
- **Alignement de séquences: algorithmes**

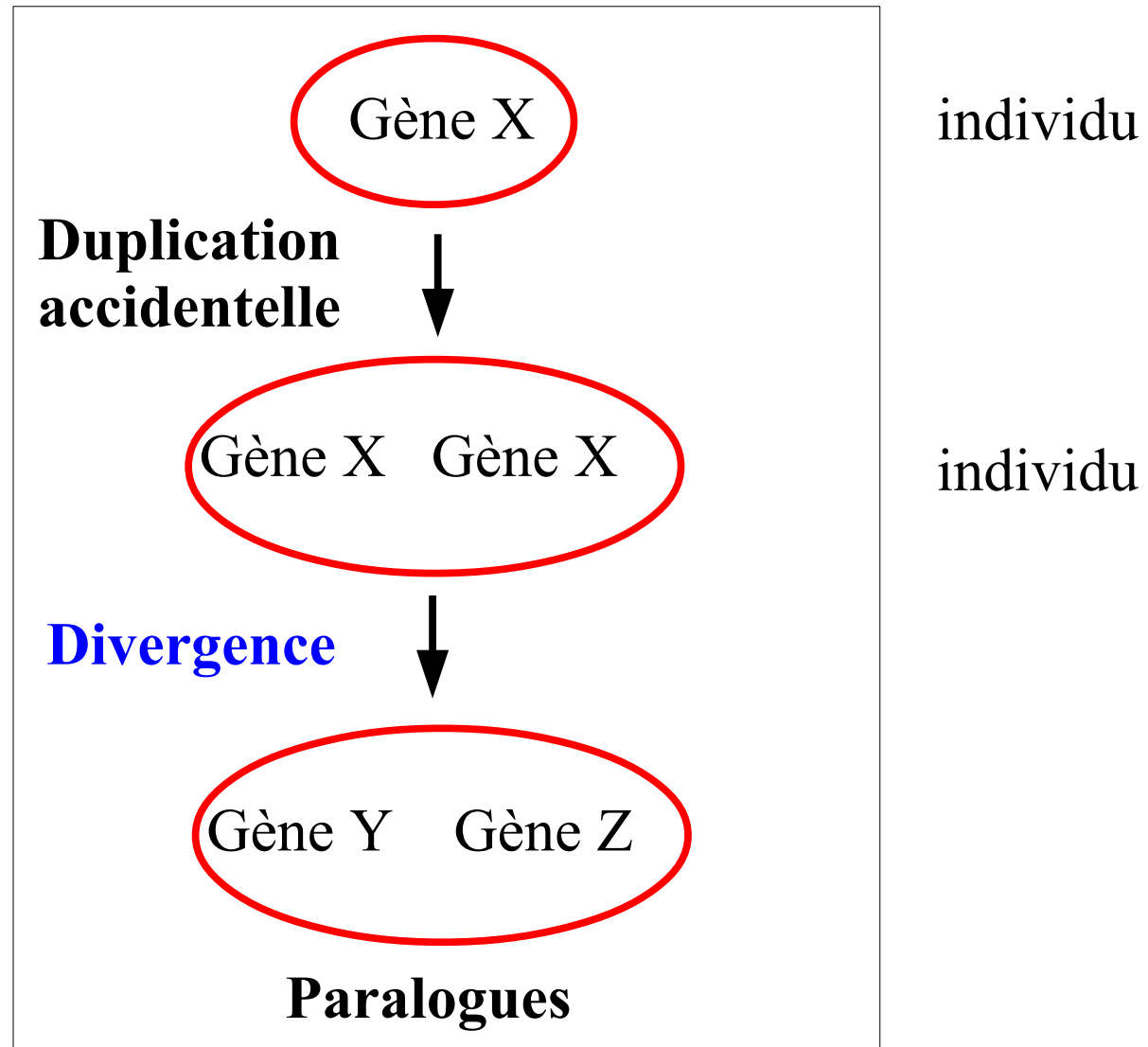


# La duplication conduit à des gènes paralogues



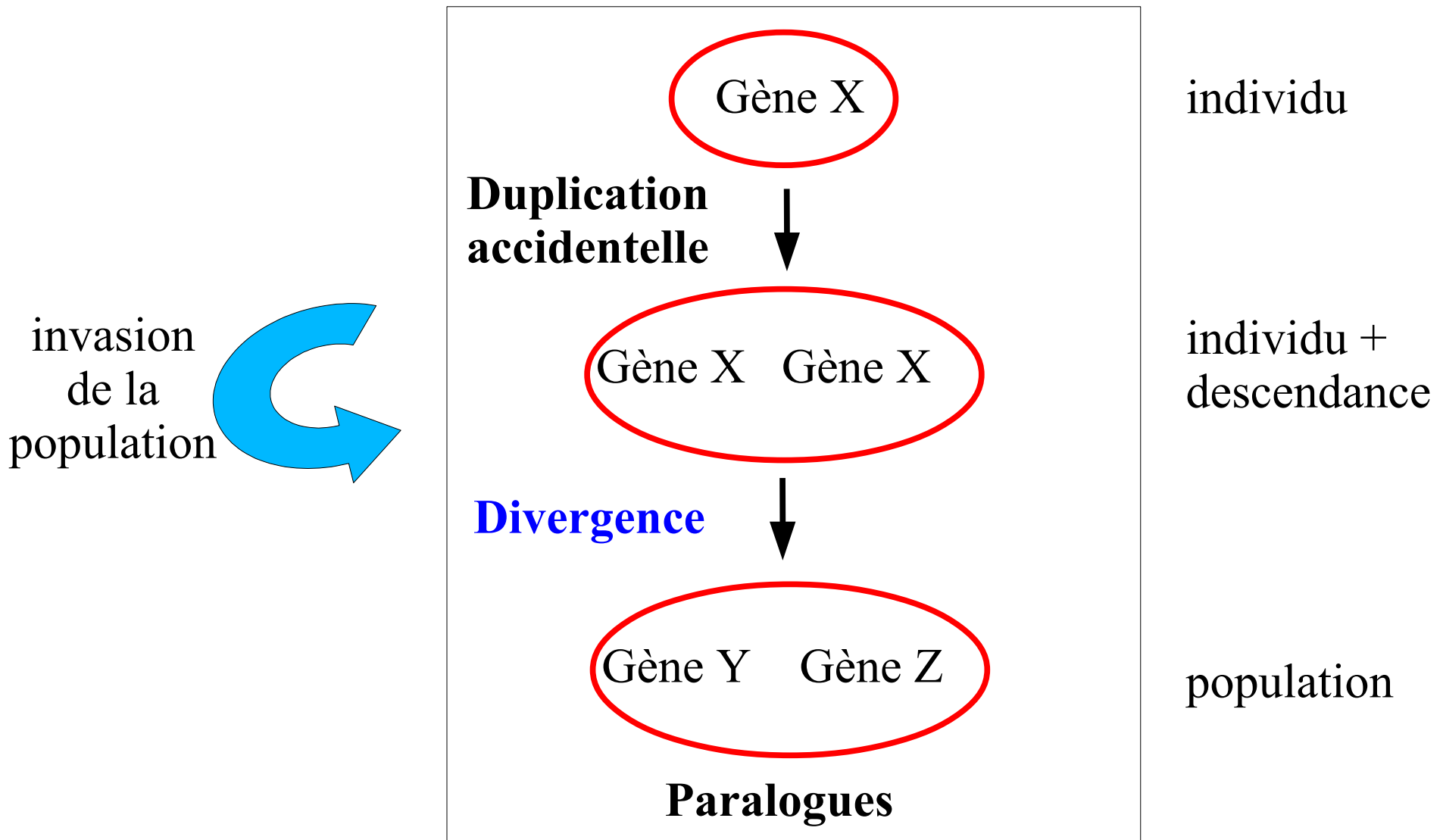
**Les paralogues changent généralement de fonction**

# La duplication conduit à des gènes paralogues



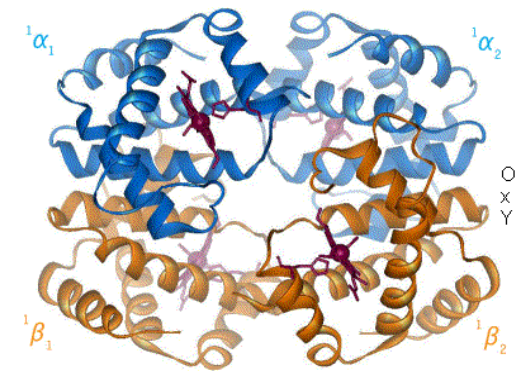
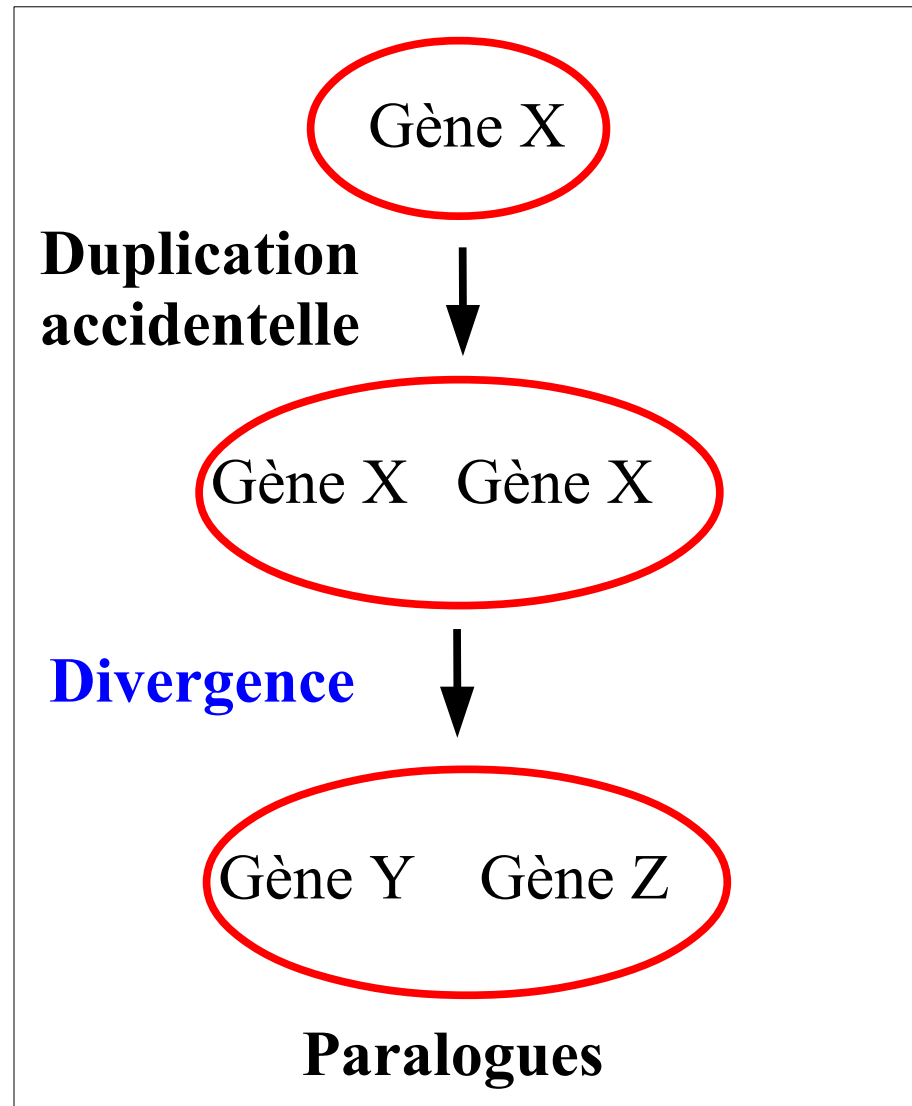
**Les paralogues changent généralement de fonction**

# La duplication conduit à des gènes paralogues



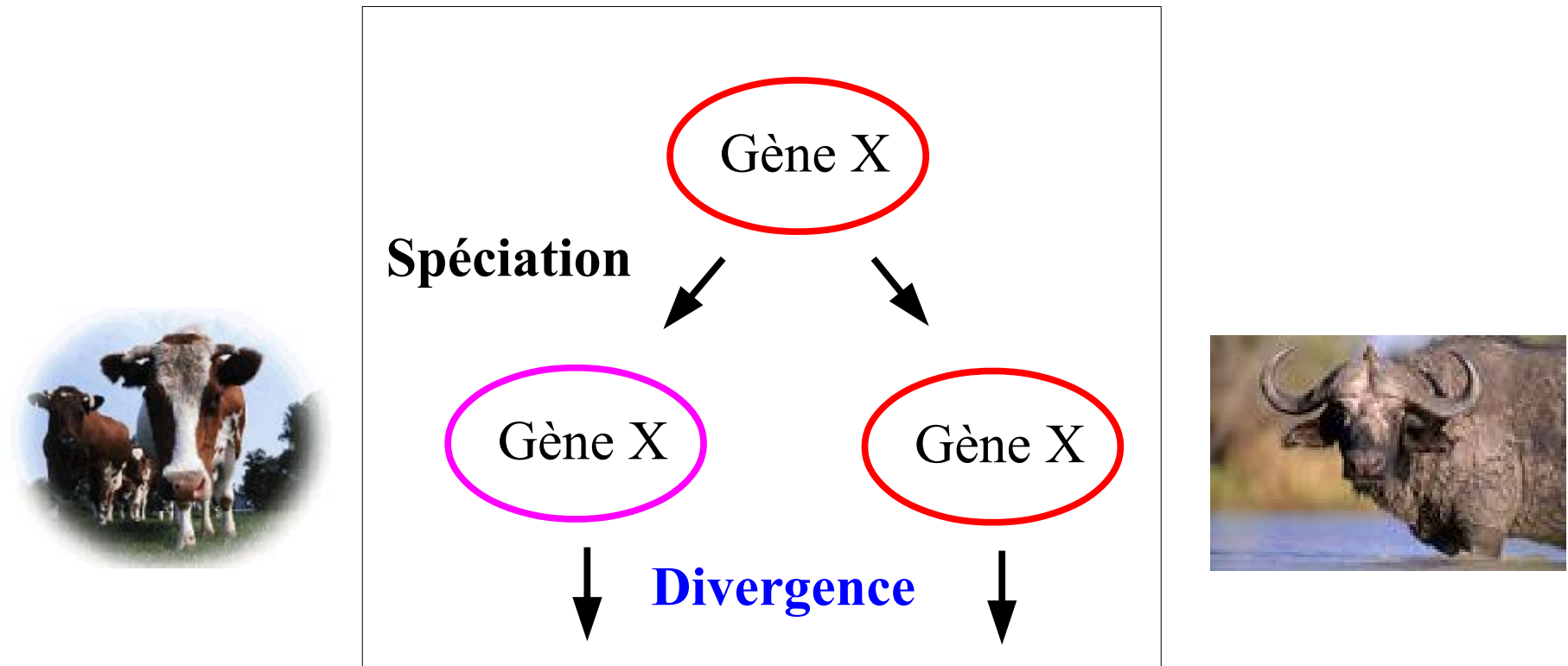
**Les paralogues changent généralement de fonction**

# La duplication conduit à des gènes paralogues



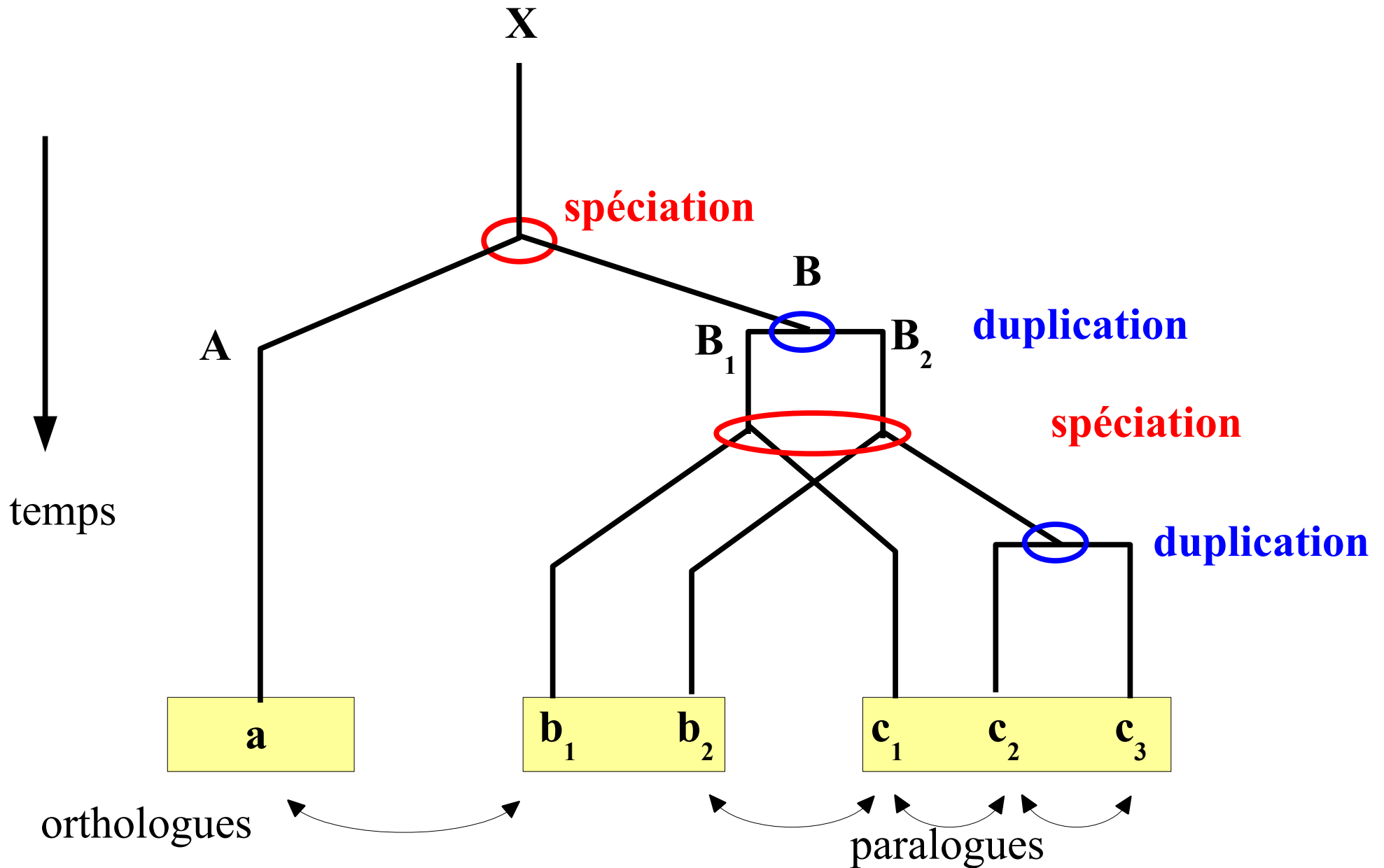
**Les paralogues changent généralement de fonction**

# La spéciation produit des gènes orthologues

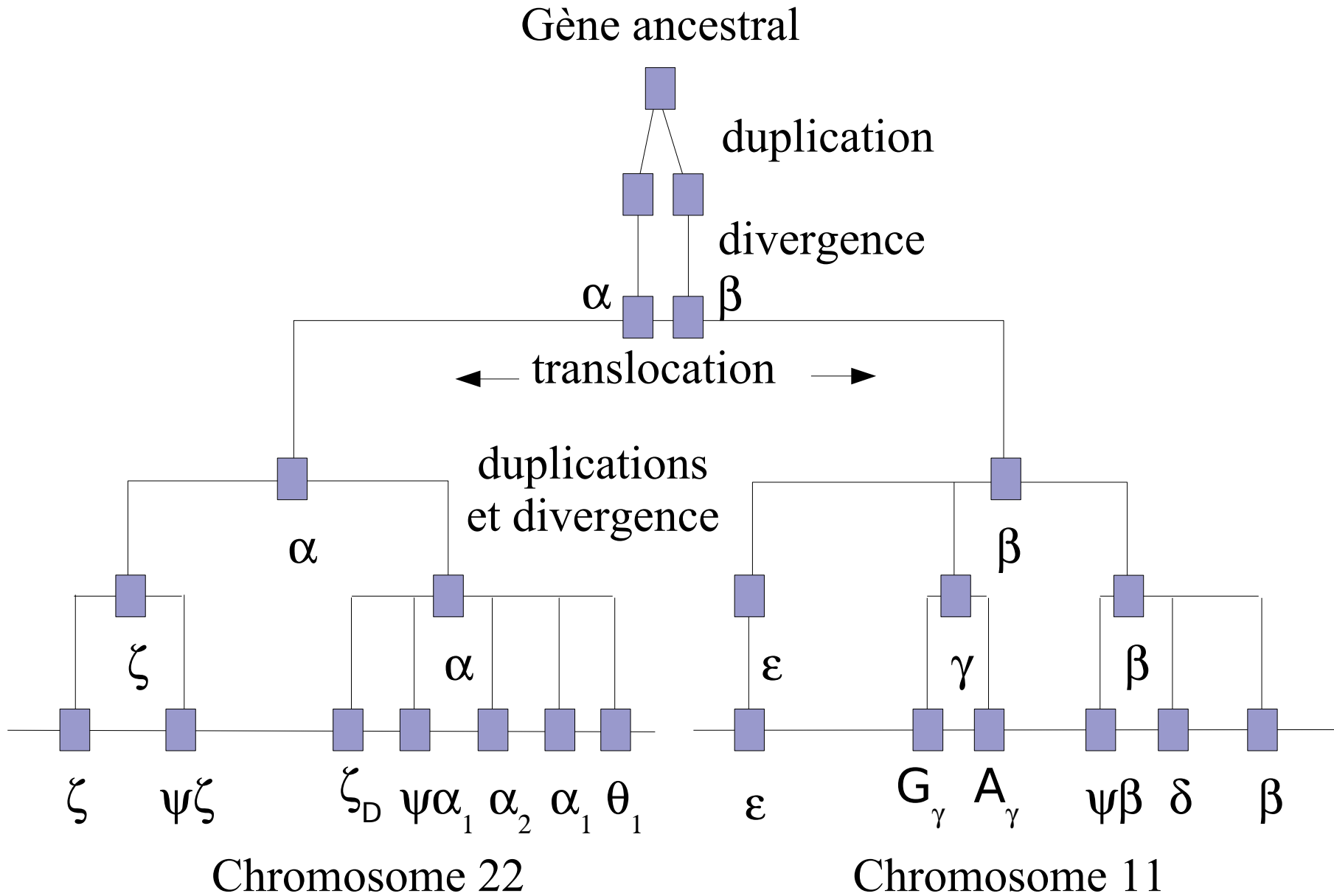


**Les gènes orthologues conservent la fonction**

# Evolution d'un gène par spéciation, duplication et divergence

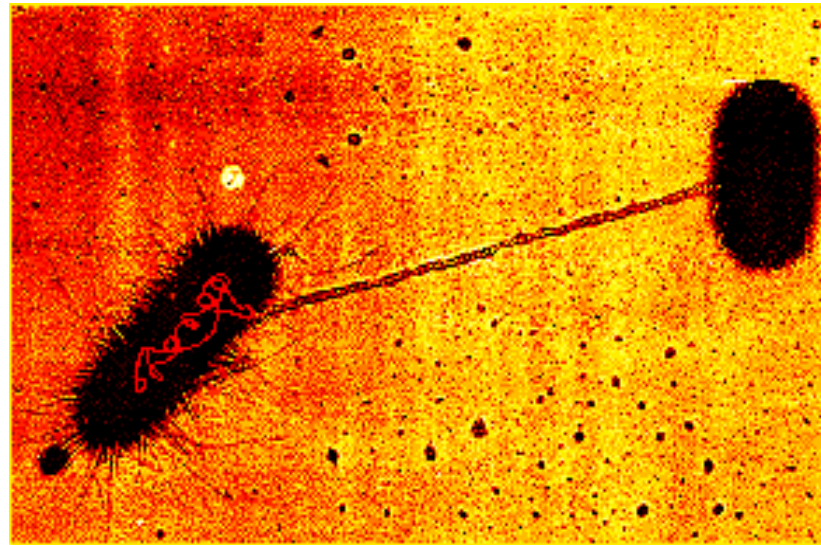


# Gènes et pseudogènes d'hémoglobine



# Transfert horizontal chez les bactéries

Un plasmide en cours de transfert d'une bactérie à une autre



~25% des gènes d'*Escherichia coli* acquis par transfert horizontal

Ce processus ne respecte pas les arbres phylogénétiques...

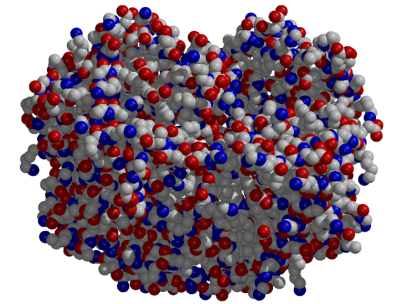
## 2) Les protéines aux séquences similaires ont souvent des structures et des fonctions similaires

K  
L  
H  
G  
G  
P  
M  
L  
D  
S  
D  
Q  
K  
F  
W  
R  
T  
P  
A  
A  
L

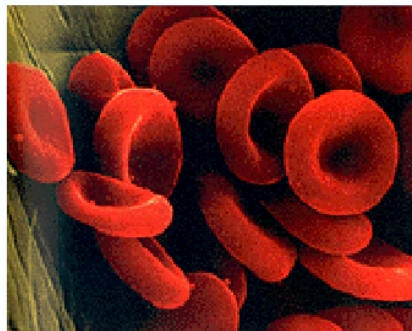
**Séquence**



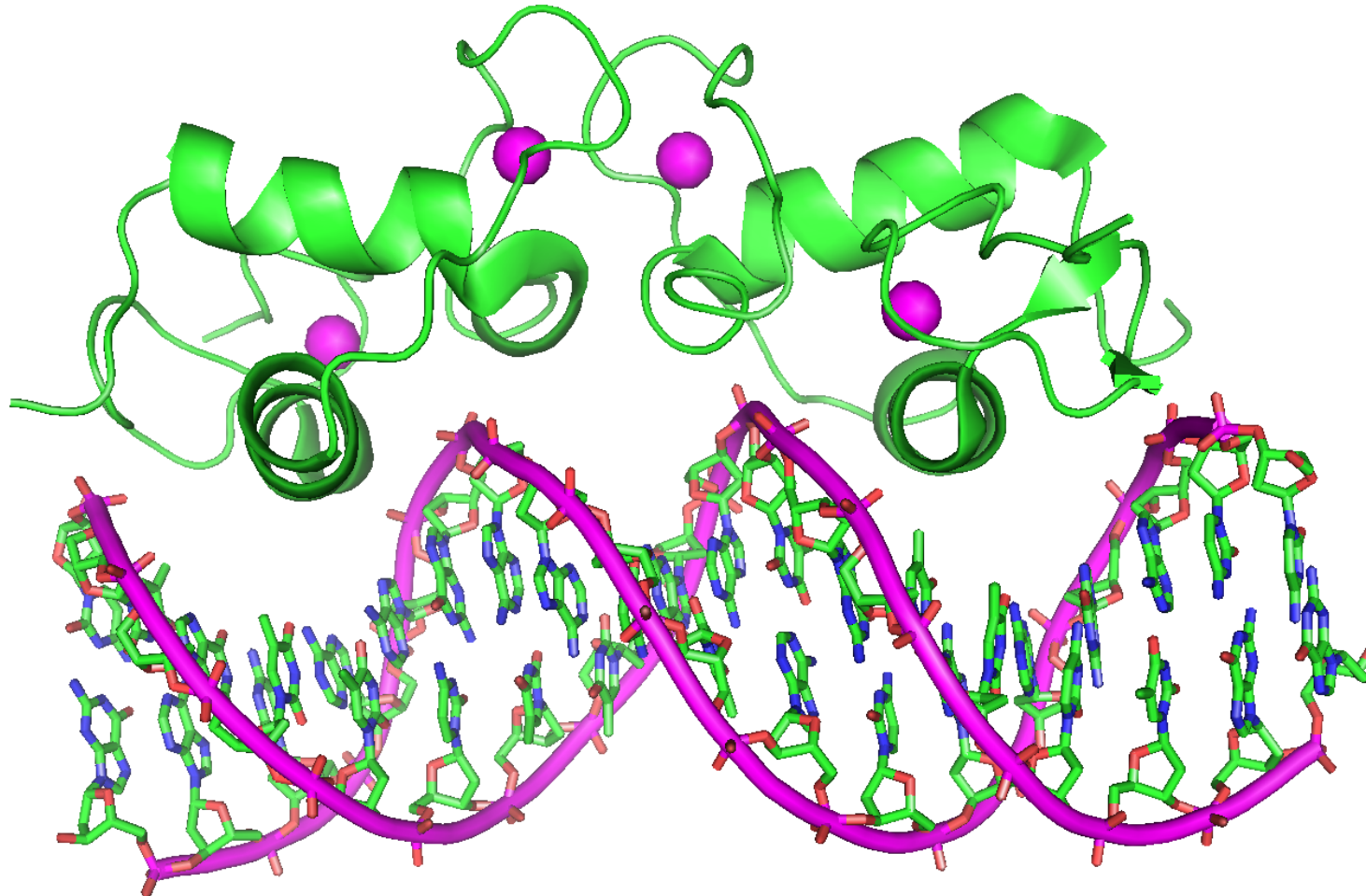
**Structure**



**Fonction**

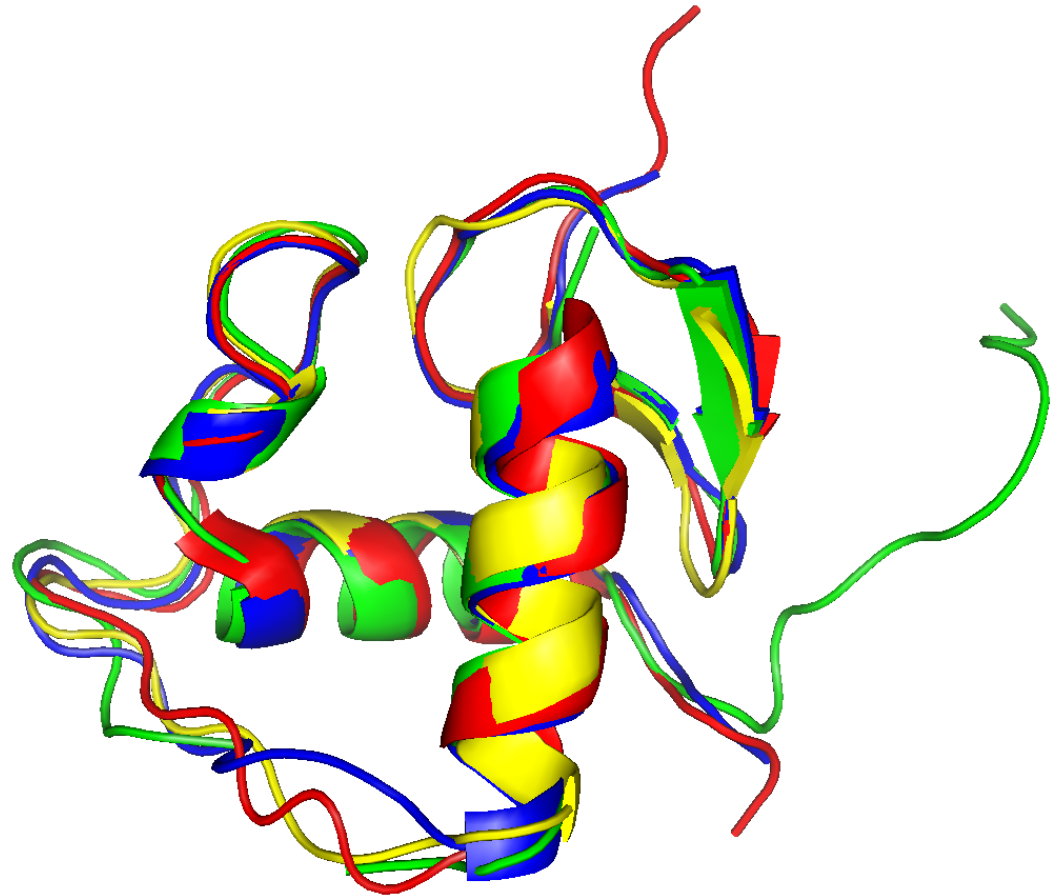


# Récepteur de l'androstérone



# Récepteurs nucléaires

Se fixent sur l'ADN et régulent l'expression génétique sous le contrôle de petits ligands: stéroïdes, vitamines, hormones, ...



androgen

CLICGDEASGAHYGALTCGSCKVFFKRAAEGKQKYL-CASRNDCTIDKFRRKNCPSCLRRLKCYEAGMTLGA

Rev Erb

CKVCGDVASGFHYGVLACEGCKGFFRRSIQQNIQYKRCLKNENCSIVRINRNRCQQCRFKKCLSVGMSRD-

glucocorticoid

CLVCSDEASGCHYGVLTCGCKAFFKRAVEGQHNYL-CKYEGKCIIDKIRRNKCPACRYRKCLQAGMNLEA

retinoic acid

CAICGDRSSGKHYGVYSCEGCKGFFKRTVRKDLTYT-CRDNKDCLIDKRQRNRCQYCRYQKCLAMGM---



# Récepteur de l'androsténone

Alignement avec le récepteur humain de la progesténone

```
PQKTCLICGDEASGAHYGALTCGSCKVFFKRAAEGKQKYL CASRNDCTIDKFRRKNCPSC
PQRVCVICGDEASGCHYGVLTCSCKVFFKRAVEGHHQYLCAGRNDCI VDKIRRKNCPAC
**! *!***** ** *!***** ** *!***** ** *!*****!
```

Alignement avec le récepteur humain de l'hormone thyroïdienne

```
PQKTCLICGDEASGAHYGALTCGSCKVFFKRAAEG--KQKYL CASRNDCTIDKFRRKNCPSC
KDEQCVVCGDKATGYHYRCITCEGCKGFFRRTIQKNLHPTYSCKYDSCCVIDKITRNQCQLC
*! !***! *! * *! * *! * *! * *! * *! * *! * *! * *! * *!
```

Alignement avec la ferrédoxine de la bactérie *Proteus vulgaris*

```
PQKTCLICGDEASGAHYGTLTCGSCKVFFKRAAEGKQKYL CASRNDCTIDKFRRKNCPSC
DQDKCIGCKTCVLACP YGTM EVVSRPVMRKL TALNTIEAFKAEANKCDLCHHRAEG-PAC
* *! * * *! * * * *! * * * *! * *! *! *
```

# **L'alignement comme modèle d'un processus évolutif**

**Hypothèse: deux protéines similaires  
ont toujours un ancêtre commun**

# L'alignement comme hypothèse d'un ancêtre commun “parsimonieux”

T	C	L	I	C	G	D	E	A	S	G	C	H	Y	Récepteur de l'androstérone
<b>T</b>	<b>C</b>	<b>L</b>	<b>V</b>	<b>C</b>	<b>G</b>	<b>D</b>	<b>E</b>	<b>A</b>	<b>T</b>	<b>G</b>	<b>Y</b>	<b>H</b>	<b>Y</b>	Ancêtre commun hypothétique
L	C	V	V	C	G	D	K	A	T	G	Y	H	Y	Récepteur de l'hormone thyroïdienne

Mutations hypothétiques en rouge

**L'hypothèse ci-dessus est-elle vraisemblable?**

# L'alignement comme hypothèse d'un ancêtre commun “parsimonieux”

E|K

T C L I C G D E A S G C H Y

Récepteur de l'androstérone

T C L V C G D **E** A T G Y H Y

Ancêtre commun hypothétique

L C V V C G D K A T G Y H Y

Récepteur de l'hormone thyroïdienne

Mutations hypothétiques en rouge

**L'hypothèse ci-dessus est-elle vraisemblable?**

# L'alignement comme hypothèse d'un ancêtre commun “parsimonieux”

T	C	L	I	C	G	D	E	A	S	G	C	H	Y	Récepteur de l'androstérone
<b>T</b>	<b>C</b>	<b>L</b>	<b>V</b>	<b>C</b>	<b>G</b>	<b>D</b>	<b>E</b>	<b>A</b>	<b>T</b>	<b>G</b>	<b>Y</b>	<b>H</b>	<b>Y</b>	Ancêtre commun hypothétique
L	C	V	V	C	G	D	K	A	T	G	Y	H	Y	Récepteur de l'hormone thyroïdienne

Mutations hypothétiques en rouge

**L'hypothèse ci-dessus est-elle vraisemblable?**

# Les mutations ont des probabilités différentes

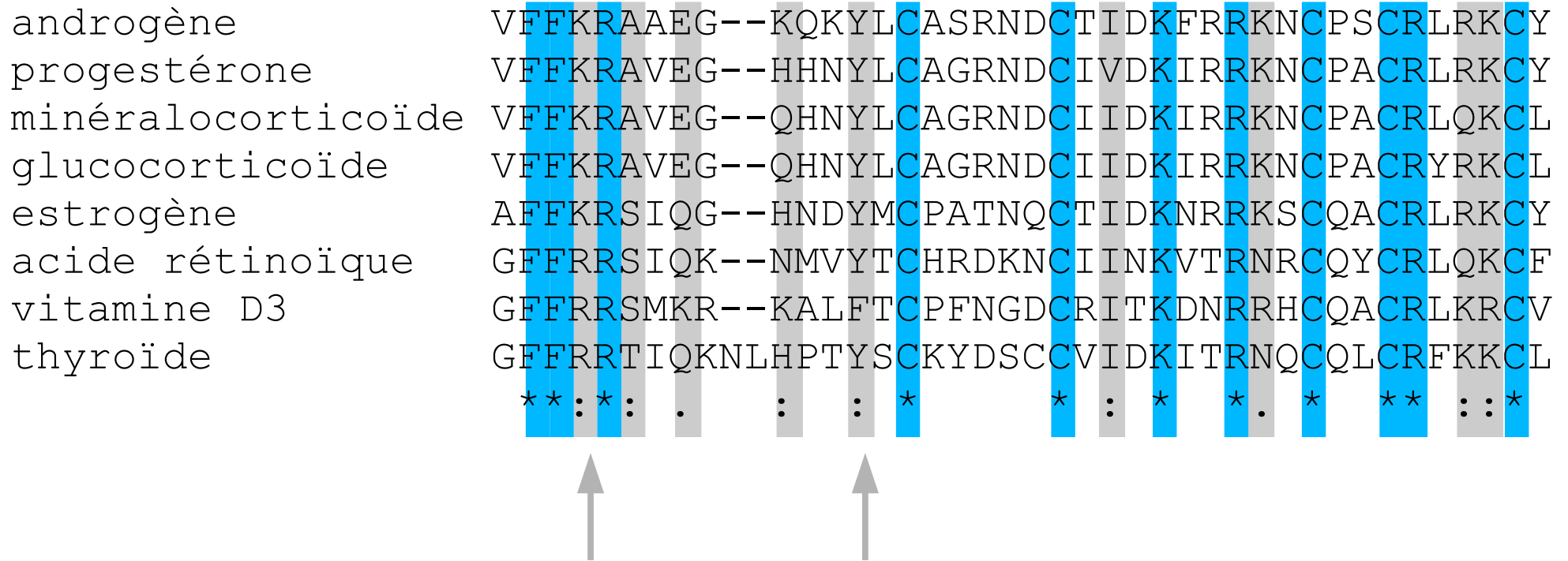
androgène	VFFKRAAEG--KQKYL	CASRNDCTIDKFRRKNC	PSCRLRKCY											
progestérone	VFFKRAVEG--HHNYL	CAGRNDCIIVDKIRRKNC	PACRLRKCY											
minéralocorticoïde	VFFKRAVEG--QHNYL	CAGRNDCIIDKIRRKNC	PACRLQKCL											
glucocorticoïde	VFFKRAVEG--QHNYL	CAGRNDCIIDKIRRKNC	PACRYRKCL											
estrogène	AFFKRSIQG--HNDYMC	PATNQCTIDKNRRKSC	QACRLRKCY											
acide rétinoïque	GFFRRSIQK--NMVYT	CHRDKNCIINKVTRNRC	QYCRLQKCF											
vitamine D3	GFFRRSMKR--KALFTC	PFNGDCRITKDNRRHC	QACRLKRCV											
thyroïde	GFFRRTIQKNLHPTYS	CKYDSCCVIDKITRNQC	QLCRFKKCL											
	** : * :	.	:	:	*	*	:	*	*	.	*	**	::	*

# Les mutations ont des probabilités différentes

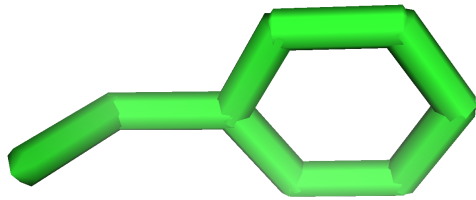
androgène	V	F	F	K	R	A	A	E	G	--	K	Q	K	Y	L	C	A	S	R	N	D	C	T	I	D	K	F	R	R	K	N	C	P	S	C	R	L	R	K	C	Y	
progestérone	V	F	F	K	R	A	V	E	G	--	H	H	N	Y	L	C	A	G	R	N	D	C	I	V	D	K	I	R	R	K	N	C	P	A	C	R	L	R	K	C	Y	
minéralocorticoïde	V	F	F	K	R	A	V	E	G	--	Q	H	N	Y	L	C	A	G	R	N	D	C	I	I	D	K	I	R	R	K	N	C	P	A	C	R	L	Q	K	C	L	
glucocorticoïde	V	F	F	K	R	A	V	E	G	--	Q	H	N	Y	L	C	A	G	R	N	D	C	I	I	D	K	I	R	R	K	N	C	P	A	C	R	Y	R	K	C	L	
estrogène	A	F	F	K	R	S	I	Q	G	--	H	N	D	Y	M	C	P	A	T	N	Q	C	T	I	D	K	N	R	R	K	S	C	Q	A	C	R	L	R	K	C	Y	
acide rétinoïque	G	F	F	R	R	S	I	Q	K	--	N	M	V	Y	T	C	H	R	D	K	N	C	I	I	N	K	V	T	R	N	R	C	Q	Y	C	R	L	Q	K	C	F	
vitamine D3	G	F	F	R	R	S	M	K	R	--	K	A	L	F	T	C	P	F	N	G	D	C	R	I	T	K	D	N	R	R	H	C	Q	A	C	R	L	K	R	C	V	
thyroïde	G	F	F	R	R	T	I	Q	K	N	L	H	P	T	Y	S	C	K	Y	D	S	C	C	V	I	D	K	I	T	R	N	Q	C	Q	L	C	R	F	K	K	C	L
	*	*	:	*	:	.	:	:	*						*		*	:	*	*	.	*		*	:	*	*	.	*		*	*	:	:	*							



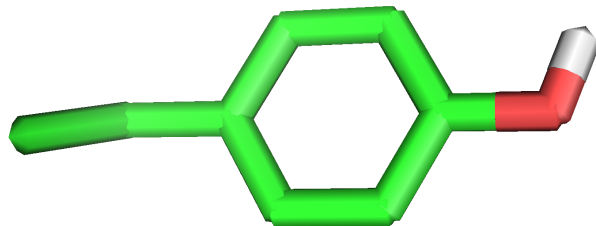
# Les mutations ont des probabilités différentes



# Les mutations ont des probabilités différentes



PHE  
F



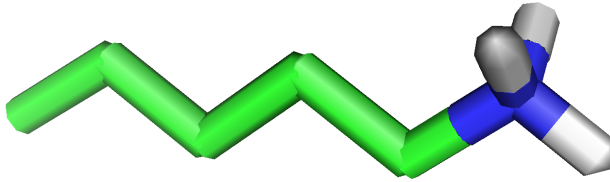
TYR  
Y

Similaires  
ou  
“homologues”

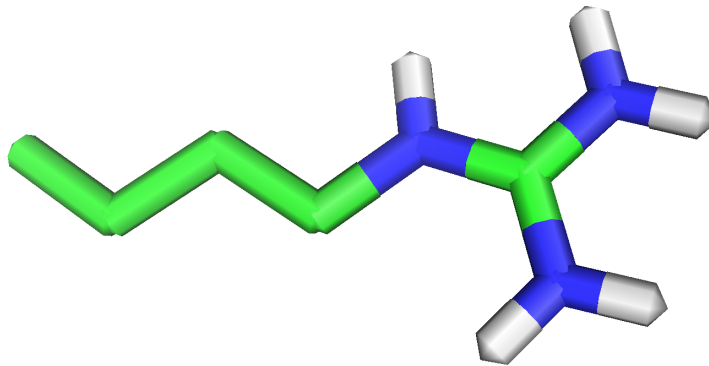
KYL  
NYL  
NYL  
NYL  
DYM  
VYT  
LFT  
TYS  
:

# Les mutations ont des probabilités différentes

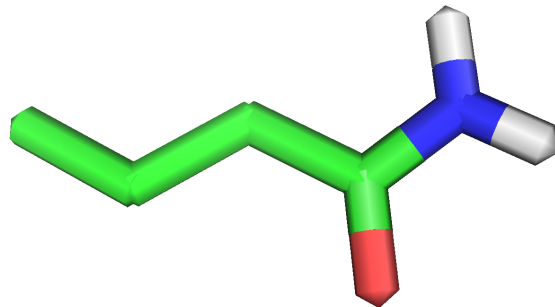
LYS  
K



ARG  
R



GLN  
Q



homologues

R	K	C	Y
R	K	C	Y
Q	K	C	L
R	K	C	L
R	K	C	Y
Q	K	C	F
K	R	C	V
K	K	C	L
:	:	*	

# Un modèle probabiliste d'évolution

$$\begin{aligned} P(x_i, y_j) &= \text{probabilité d'observer } x_i \text{ aligné avec } y_j \\ &= P(\text{"}x_i \text{ muté en } y_j\text{" ou "y}_j \text{ muté en } x_i\text{"}) \end{aligned}$$

séquence x:      T C L I C G D E A S G C H Y  
séquence y:      L C V V C G D K A T G Y H Y

- Probabilités estimées à partir d'alignements tests
- Hypothèse de positions équivalentes et indépendantes

# Récepteur de l'androstérone

Alignement avec le récepteur humain de la progestérone

```
PQKTCLICGDEASGAHYGALTCGSCKVFFKRAAEGKQKYL CASRNDCTIDKFRRKNCPSC
PQRVCVICGDEASGCHYGVLTCSCKVFFKRAVEGHHQYLCAGRNDCIVDKIRRKNCPAC
**! *!***** ** *!***** ** *!***** ** *!*****!
```

Alignement avec le récepteur humain de l'hormone thyroïdienne

```
PQKTCLICGDEASGAHYGALTCGSCKVFFKRAAEG--KQKYL CASRNDCTIDKFRRKNCPSC
KDEQCVVCGDKATGYHYRCITCEGCKGFFRRTIQKNLHPTYSCKYDSCCVIDKITRNQCQLC
*! !***!*! * *! * * *! *! * * * * * * * * * * * * *
```

Alignement avec la ferrédoxine de la bactérie *Proteus vulgaris*

```
PQKTCLICGDEASGAHYGTLTCGSCKVFFKRAAEGKQKYL CASRNDCTIDKFRRKNCPSC
DQDKCIGCKTCVLACP YGTMEVVS RPVMRKL TALNTIEAFKAEANKCDLCHHRAEG-PAC
* *! * * *! * * * * * * * * *! *! *!*
```

Notre explication (hypothèse) est plus ou moins probable: oui, mais...

**Ce n'est pas parce qu'un modèle  
reproduit les données expérimentales  
qu'il est juste**

**Ce n'est pas parce qu'un modèle  
reproduit les données expérimentales  
qu'il est juste**

**Il faut aussi exclure les autres modèles  
ou explications a priori plausibles.**

# Modèle “aléatoire” comme hypothèse concurrente

$Q(x_i, y_j)$  = probabilité d'observer  $x_i$  et  $y_j$   
dans deux protéines indépendantes

$\approx q_{x_i} q_{y_j}$  fréquences naturelles de  $x_i, y_j$

	x	$q_x$	
tryptophane	W	1.3 %	des acides aminés
leucine	L	9.0 %	

W      L  
W    ≠    L

# Vraisemblance d'une colonne dans un alignement

déf.

$$s(x_i, y_j) = \log P(x_i, y_j) / Q(x_i, y_j)$$
$$= \log [ P(x_i, y_j) / q_{x_i} q_{y_j} ]$$

$x_i$

T	C	L	I	C	G	D	E	A	S	G	C	H	Y
L	C	V	V	C	G	D	K	A	T	G	Y	H	Y

$y_i$

$s$  = matrice de “score”



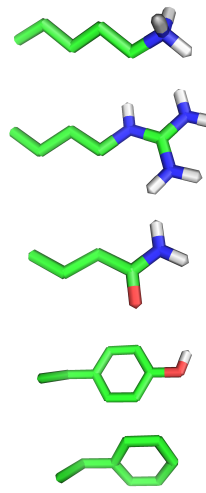
# Une matrice empirique s très utilisée

**BLOSUM62**

	C	S	T	P	A	G	N	D	E	Q	H	R	K	M	I	L	V	F	Y	W
C	9	-1	-1	-3	0	-3	-3	-3	-4	-3	-3	-3	-3	-1	-1	-1	-1	-2	-2	-2
S		4	1	-1	1	0	1	0	0	0	-1	-1	0	-1	-2	-2	-2	-2	-2	-3
T			5	-1	0	-2	0	-1	-1	-1	-2	-1	-1	-1	-1	-1	0	-2	-2	-2
P				7	-1	-2	-2	-1	-1	-1	-2	-2	-1	-2	-3	-3	-2	-4	-3	-4
A					4	0	-2	-2	-1	-1	-2	-1	-1	-1	-1	-1	0	-2	-2	-3
G						6	0	-1	-2	-2	-2	-2	-2	-3	-4	-4	-3	-3	-3	-2
N							6	1	0	0	1	0	0	-2	-3	-3	-3	-3	-2	-4
D								6	2	0	-1	-2	-1	-3	-3	-4	-3	-3	-3	-4
E									5	2	0	0	1	-2	-3	-3	-2	-3	-2	-3
Q										5	0	1	1	0	-3	-2	-2	-3	-1	-2
H											8	0	-1	-2	-3	-3	-3	-1	2	-2
R												5	2	-1	-3	-2	-3	-3	-2	-3
K													5	-1	-3	-2	-2	-3	-2	-3
M														5	1	2	1	0	-1	-1
I															4	2	3	0	-1	-3
L																4	1	0	-1	-2
V																	4	-1	-1	-3
F																		6	3	1
Y																			7	2
W																				11

Pénalité pour gaps: voir plus loin

# Les mutations ont des vraisemblances différentes



	K	R	Q	Y	F
K	<b>6</b>	<b>3</b>	<b>2</b>	-2	-4
R		<b>7</b>	<b>1</b>	-1	-3
Q			<b>7</b>	-1	-4
Y				<b>8</b>	<b>4</b>
F					<b>8</b>

Extrait de la matrice empirique BLOSUM50

# La vraisemblance ou probabilité d'un alignement

$$P(x_i, y_j) = \text{probabilité d'observer } x_i \text{ aligné avec } y_j$$
$$= P(\text{“}x_i \text{ muté en } y_j\text{” ou “}y_j \text{ muté en } x_i\text{”})$$

T	C	L	I	C	G	D	E	A	S	G	C	H	Y
L	C	V	V	-	G	D	K	A	T	G	Y	H	Y

- Probabilités estimées à partir d'alignements tests
- Hypothèse de positions équivalentes
- Hypothèse de positions indépendantes:  $P(\text{alignement}) = \prod_{(i,j)} P(x_i, y_j)$

# Vraisemblance ou “score” d'un alignement

$$\begin{aligned} s(x_i, y_j) &= \log P(x_i, y_j) / Q(x_i, y_j) \\ &= \log [ P(x_i, y_j) / q_{x_i} q_{y_j} ] \end{aligned}$$

$$S(x, y) = \sum_{(i,j)} s(x_i, y_j)$$

T C L I C G D E A S G C H Y  
L C V V - G D K A T G Y H Y

# La probabilité de mutation dépend du temps

$P(x,y) = P(x,y|T)$  = probabilité d'une mutation pendant le temps T

Les cytochromes c de chimpanzee et d'humain sont plus similaires que les cytochromes c d'humain et d'Escherichia coli...

**Si on compare des protéines de mammifère, par exemple, les vrais homologues seront très similaires: scores élevés. Si on compare mammifère et bactérie, de vrais homologues seront nettement moins similaires: scores plus faibles si on utilise la même matrice**

# La probabilité de mutation dépend du temps

$P(x,y) = P(x,y|T)$  = probabilité d'une mutation pendant le temps T

Les cytochromes c de chimpanzee et d'humain sont plus similaires que les cytochromes c d'humain et d'Escherichia coli...

**Une matrice de similarité n'est optimale que pour une échelle de temps, ou de similarité donnée.**

BLOSUM40  $\neq$  BLOSUM50  $\neq$  BLOSUM62  $\neq$  BLOSUM80

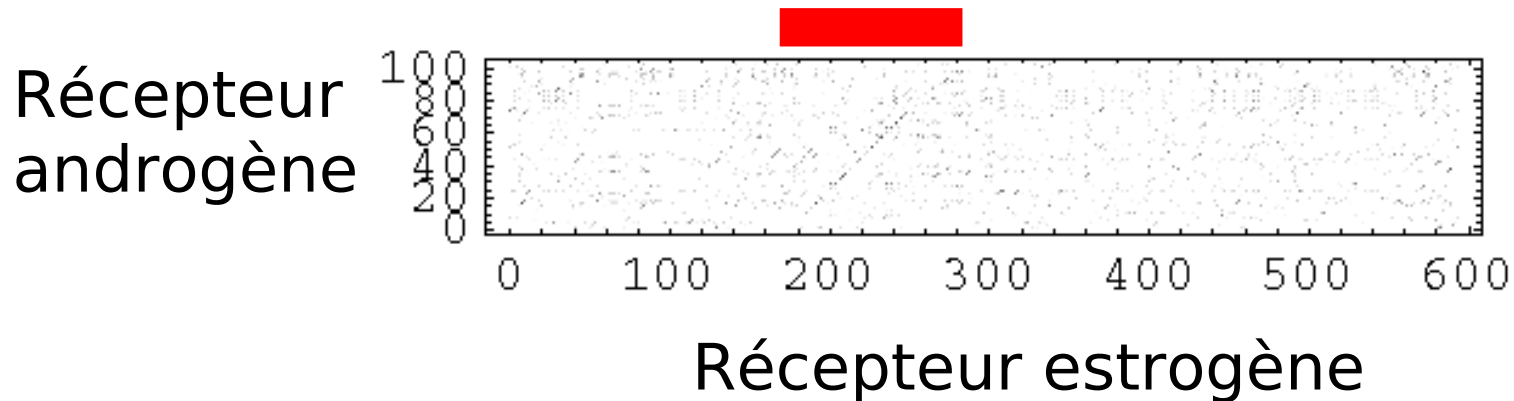
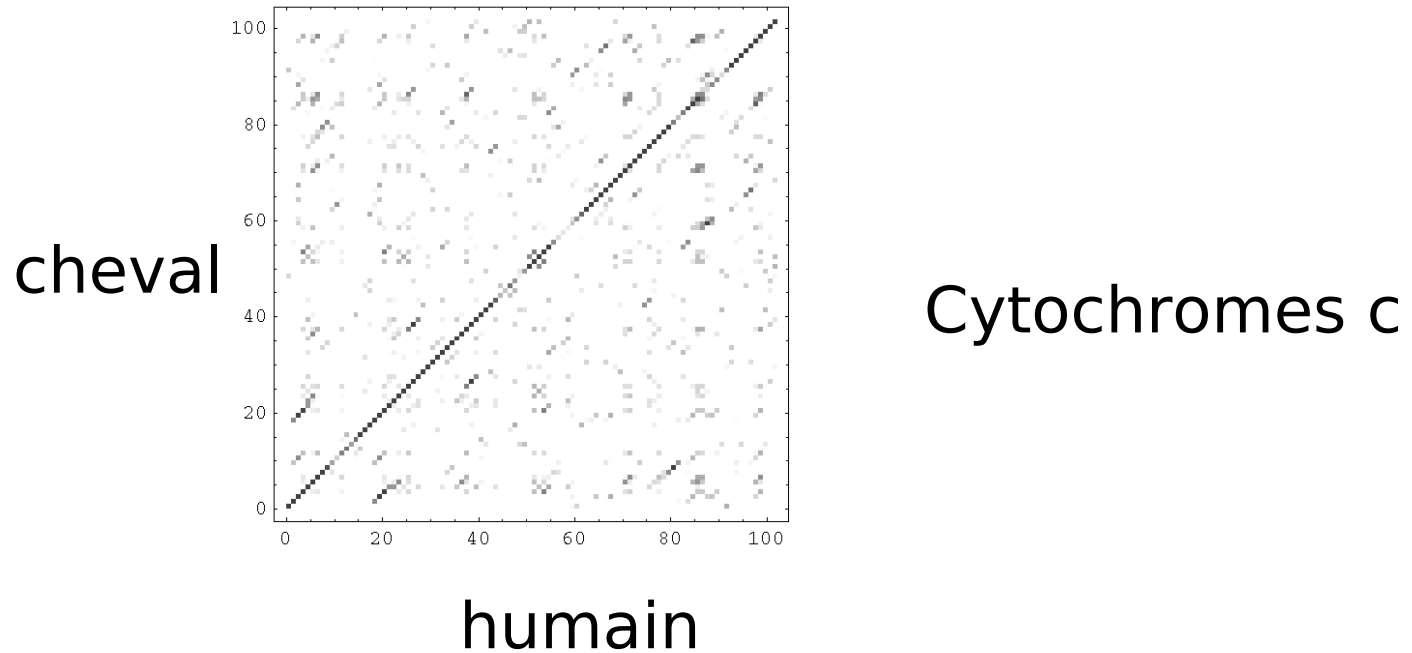
## **En résumé, un alignement correspond à un modèle d'évolution, qui comprend:**

- Hypothèse de divergence “minimale” depuis un ancêtre commun
- Détermination empirique des probabilités de mutations
- Hypothèse de positions équivalentes et indépendantes
- Un modèle de référence (hypothèse “nulle”)

# **L'alignement de séquences: algorithmes**



# Comparaison graphique de deux séquences par une matrice densité



# Comment trouver l'alignement le plus vraisemblable?

Pour deux séquences de longueur n:

$$N \approx \binom{2n}{n} \approx \frac{2^{2n}}{\sqrt{\pi n}} \text{ alignements possibles!}$$

N	1	2	3	4	5	6	7	8	9	....	100
n	2	6	20	70	252	924	3.432	17.160	48.620	....	$10^{59}$

Recherche exhaustive impossible!

# Existence d'un algorithme récursif

Un alignement donné,  $A$ , comprend de nombreux alignements plus petits,  $A_1, A_2, A_3, \dots$

$A$

T	C	L	I	C	G	D	E	A	S	G	C	H	Y
L	C	V	V	-	G	D	K	A	T	G	Y	H	Y

$A_1$

T	C	L	I	C	G	D
L	C	V	V	-	G	D

$A_2$

E	A	S	G	C	H	Y
K	A	T	G	Y	H	Y

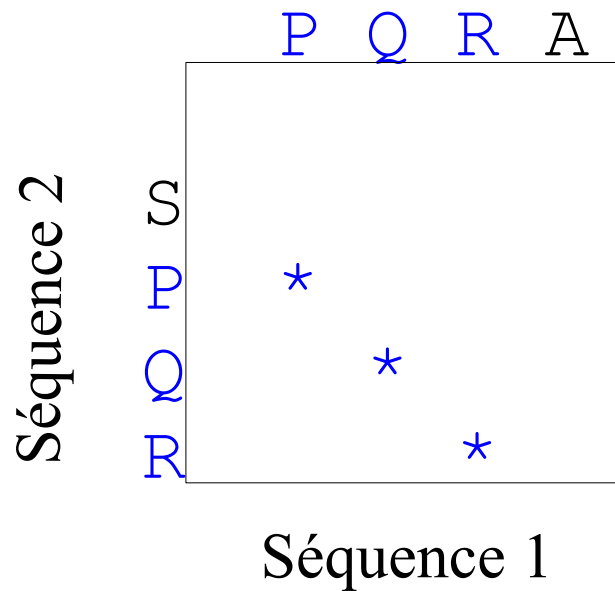
$A_3$

I	C	G	D	E	A
V	-	G	D	K	A

→ Pour évaluer un alignement  $A$ , on peut utiliser les résultats déjà acquis pour des sous-alignements (scores additifs):  
algorithme récursif

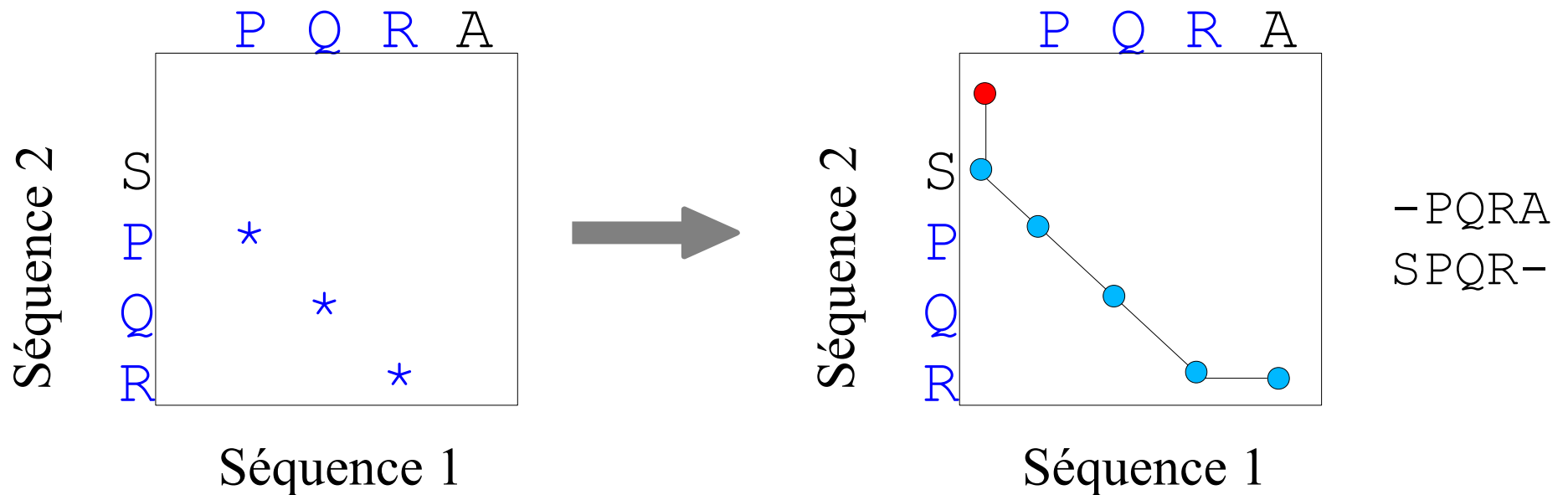
# Une méthode qui s'inspire de la “matrice” déjà vue

On dispose les deux séquences sur les bords d'une matrice:



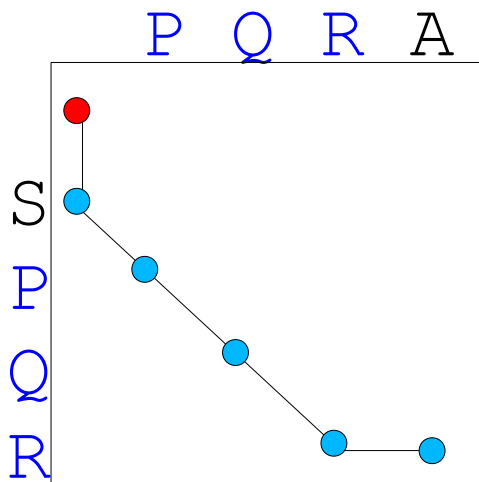
# Une méthode qui s'inspire de la “matrice” déjà vue

On dispose les deux séquences sur les bords d'une matrice.  
**Un alignement correspond à une ligne à travers la matrice:**

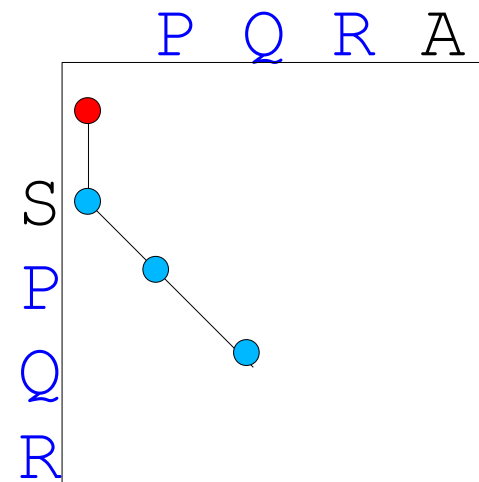


# Une méthode qui s'inspire de la “matrice” déjà vue

On dispose les deux séquences sur les bords d'une matrice.  
**Un alignement correspond à une ligne à travers la matrice:**



-PQRA  
SPQR-

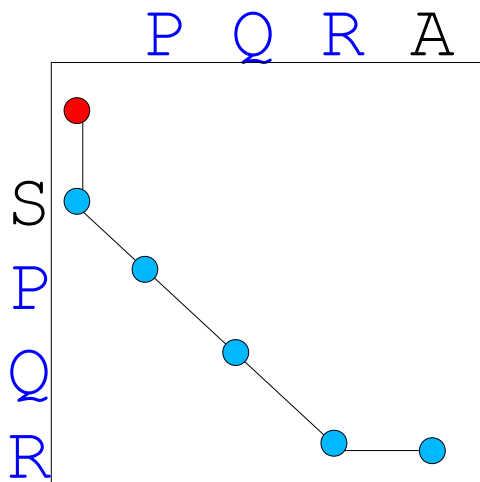


-PQ  
SPQ

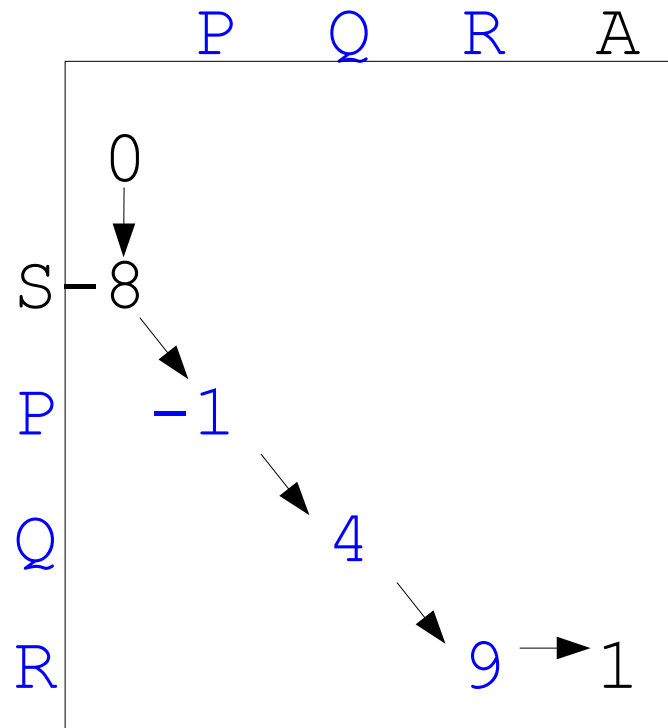
# Une méthode qui s'inspire de la “matrice” déjà vue

Un alignement correspond à une ligne à travers la matrice.

**On annote la table avec les scores:**



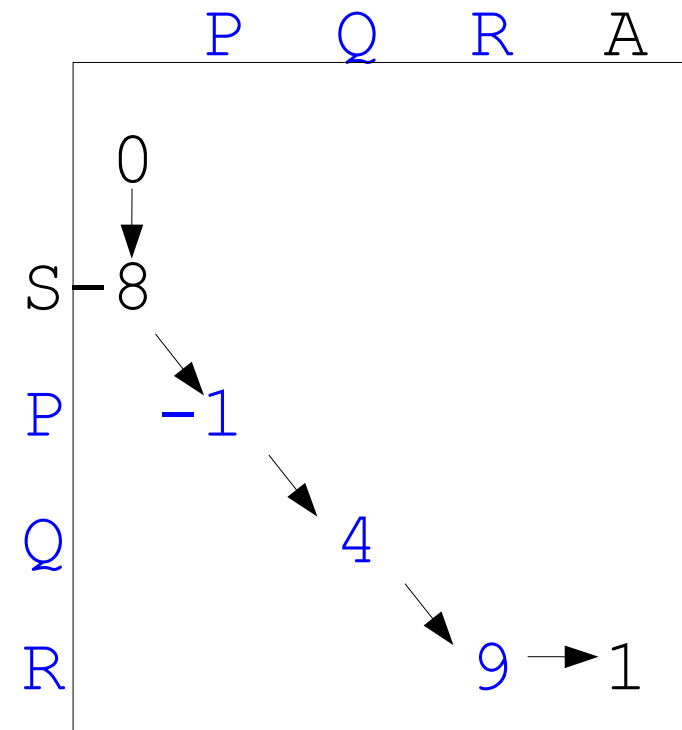
-PQRA  
SPQR-



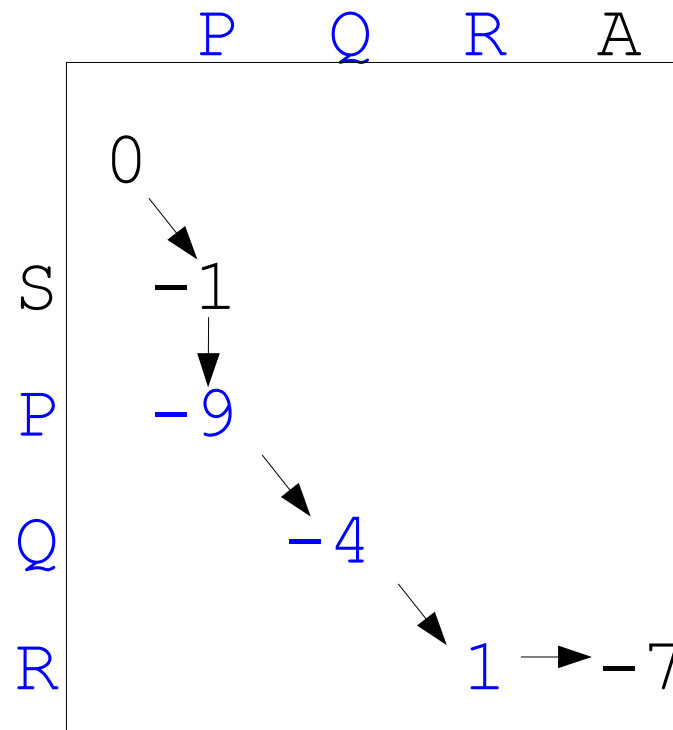
	S	P	Q	R
P	-1	7	-1	-2
Q	0	-1	5	1
R	-1	-2	1	5
A	1	-1	-1	-1

# Une méthode qui s'inspire de la “matrice” déjà vue

L'alignement optimal est en compétition avec beaucoup d'autres alignements possibles:



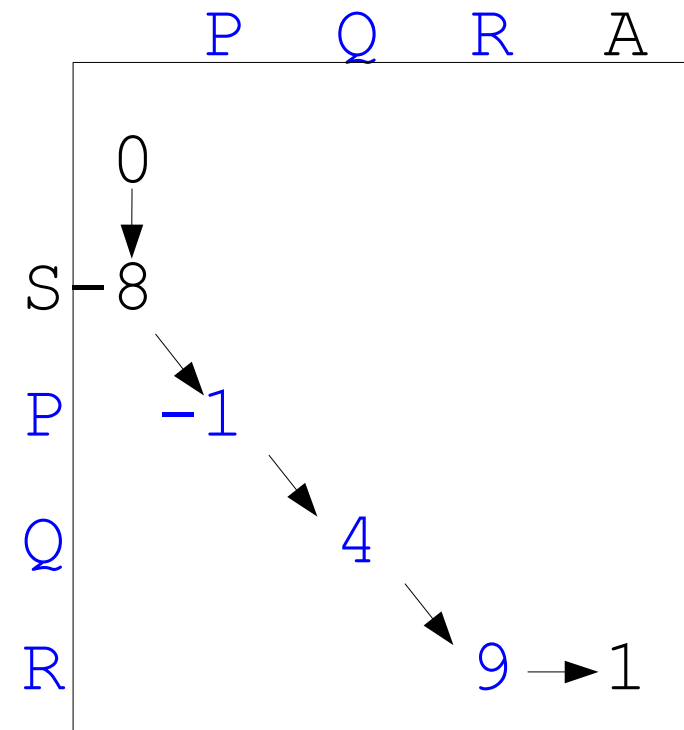
-PQRA  
SPQR-



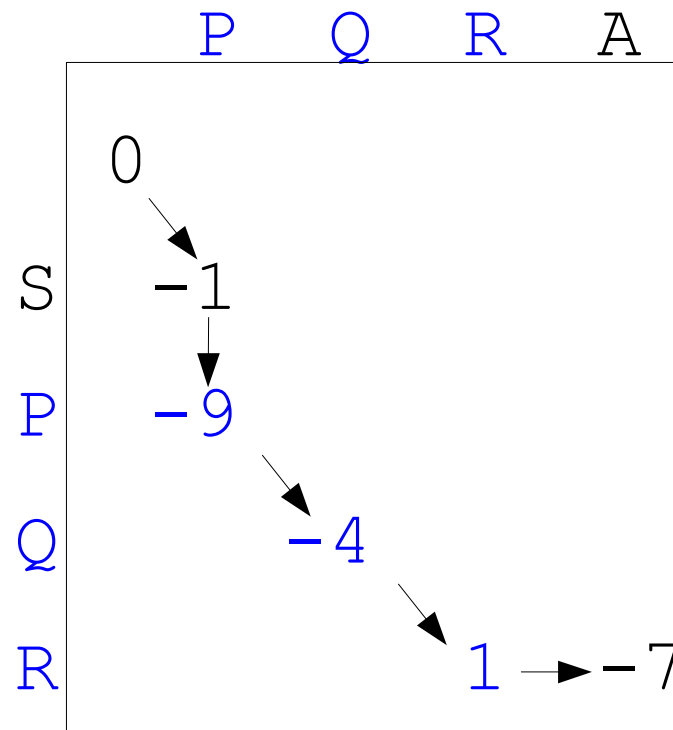
P-QRA  
SPQR-

# Une méthode qui s'inspire de la “matrice” déjà vue

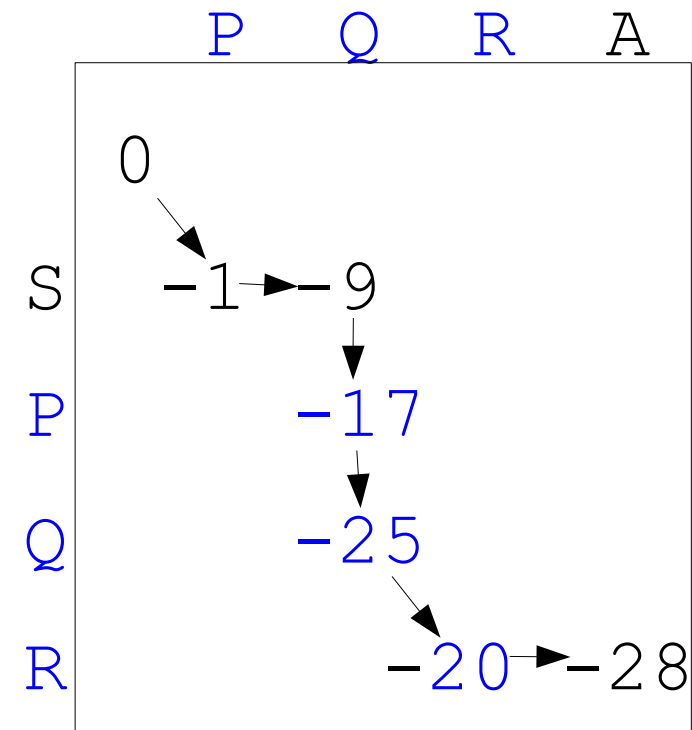
L'alignement optimal est en compétition avec beaucoup d'autres alignements possibles:



-PQRA  
SPQR-



P-QRA  
SPQR-



PQ--RA  
S-PQR-

# Comment trouver le meilleur chemin à travers la table?

Supposons que nous connaissons les scores de trois alignements particuliers:

	P	Q	R	A
S	0			
P				
Q			-1	-9
R			9	?

Chaque score bleu correspond à un alignement incomplet, inconnu à ce stade.

**A partir de ces 3 scores, je peux trouver le score manquant**

# Comment trouver le meilleur chemin à travers la table?

Supposons que nous connaissons les scores de trois alignements particuliers:

	P	Q	R	A
S	0			
P				
Q			-1	-9
R			9	?

**A partir des 3 scores bleus, je peux trouver le score manquant**

3 possibilités:



$$\text{score } -1 - 1 = -2$$



$$\text{score } 9 - 8 = 1$$



$$\text{score } -9 - 8 = -17$$

# Comment trouver le meilleur chemin à travers la table?

Supposons que nous connaissons les scores de trois alignements particuliers:

	P	Q	R	A
S	0			
P				
Q			-1	-9
R			9	<b>1</b>

**A partir des 3 scores bleus, je peux trouver le score manquant**

3 possibilités:



$$\text{score } -1 - 1 = -2$$



$$\text{score } 9 - 8 = 1$$



$$\text{score } -9 - 8 = -17$$

*Si les chemins aboutissant aux cases bleues sont optimaux, le chemin final aussi*

# Algorithme de “programmation dynamique” de Needleman-Wunsch

Pour aligner deux séquences:

$$\begin{array}{ccccccc} X_1 & X_2 & X_3 & X_4 & \dots & X_p \\ Y_1 & Y_2 & Y_3 & Y_4 & Y_5 & \dots & Y_q \end{array}$$

On construira les **alignements optimaux** de toutes les tailles:

alignement optimal de  $x_1, \dots, x_i$  avec  $y_1, \dots, y_j$

pour chaque  $i \leq p, j \leq q$

# Matrice F de scores d'alignements partiels

Soit:  $F(i,j)$  = score du meilleur alignement de  $x_1, \dots, x_i$  avec  $y_1, \dots, y_j$

déf.  
 $A(i,j)$  = l'alignement correspondant

	$Y_1$	$Y_2$	$Y_3$	$Y_4$	$Y_5$	$Y_6$
$X_1$						
$X_2$			<b>F(3,4)</b>			
$X_3$						
$X_4$						
$X_5$						

**Matrice F(i,j)  
de scores optimaux**

# Algorithme de programmation dynamique de Needleman-Wunsch

On va déduire  $F(i,j)$  des scores d'alignements (optimaux) plus petits

	$Y_1$	$Y_2$	$Y_3$	$Y_4$	$Y_5$	$Y_6$
$X_1$						
$X_2$			<b>F(3,4)</b>			
$X_3$						
$X_4$						
$X_5$						

# Calcul récursif de $F(i,j)$

Un alignement de

$x_1 \ x_2 \ \dots \ x_i$

avec

$y_1 \ y_2 \ \dots \ y_j$

peut être obtenu en ajoutant  $x_i$  ou  $y_j$  ou les deux au bout d'un alignement plus petit.

I G A  $x_i$   
L G V  $y_j$

A I G  $x_i$   
G V  $y_j$  -

G A  $x_i$  - -  
S L G V  $y_j$

création ou extension d'un "gap"

création ou extension d'un "gap"

# Calcul récursif de $F(i,j)$

L'alignement **optimal** de

$x_1 \ x_2 \ \dots \ x_i$

avec

$y_1 \ y_2 \ \dots \ y_j$

peut être obtenu en ajoutant  $x_i$  ou  $y_j$  ou les deux au bout d'un alignement **optimal** plus petit.

I G A  $x_i$   
L G V  $y_j$

A I G  $x_i$   
G V  $y_j$  -

G A  $x_i$  - -  
S L G V  $y_j$

# Calcul récursif de $F(i,j)$

$$F(i,j) = \max \left\{ \begin{array}{l} F(i-1,j-1) + s(x_i, y_j) \\ F(i-1,j) - d \\ F(i,j-1) - d \end{array} \right.$$

I	G	A	$x_i$	
L	G	V	$y_j$	

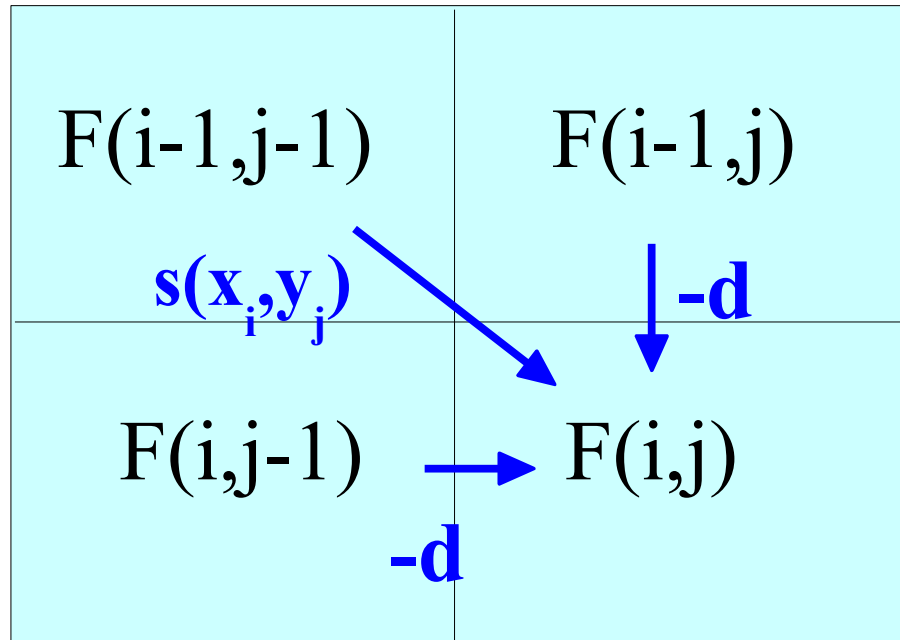
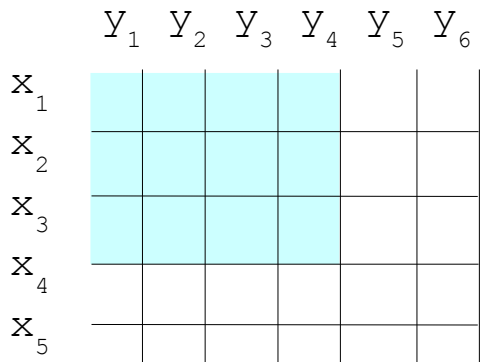
A	I	G	$x_i$	
G	V	$y_j$	-	création ou extension d'un "gap"

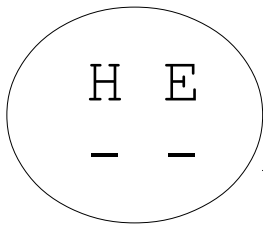
G	A	$x_i$	-	-
S	L	G	V	$y_j$

$s$  est la matrice de score employée (eg, BLOSUM62)

# Calcul récursif de $F(i,j)$



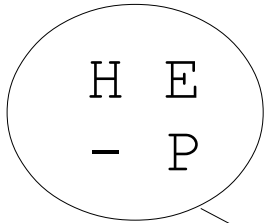
# Initialisation de F



	-	H	E	A	G	A	W
-	0	-d	-2d	-3d	-4d	-5d	-6d
P	-d						
A	-2d						
W	-3d						
H	-4d						
E	-5d						

$$s(x_i, -) = s(-, y_j) = -d = \text{coût d'un gap}$$

# Extension de F



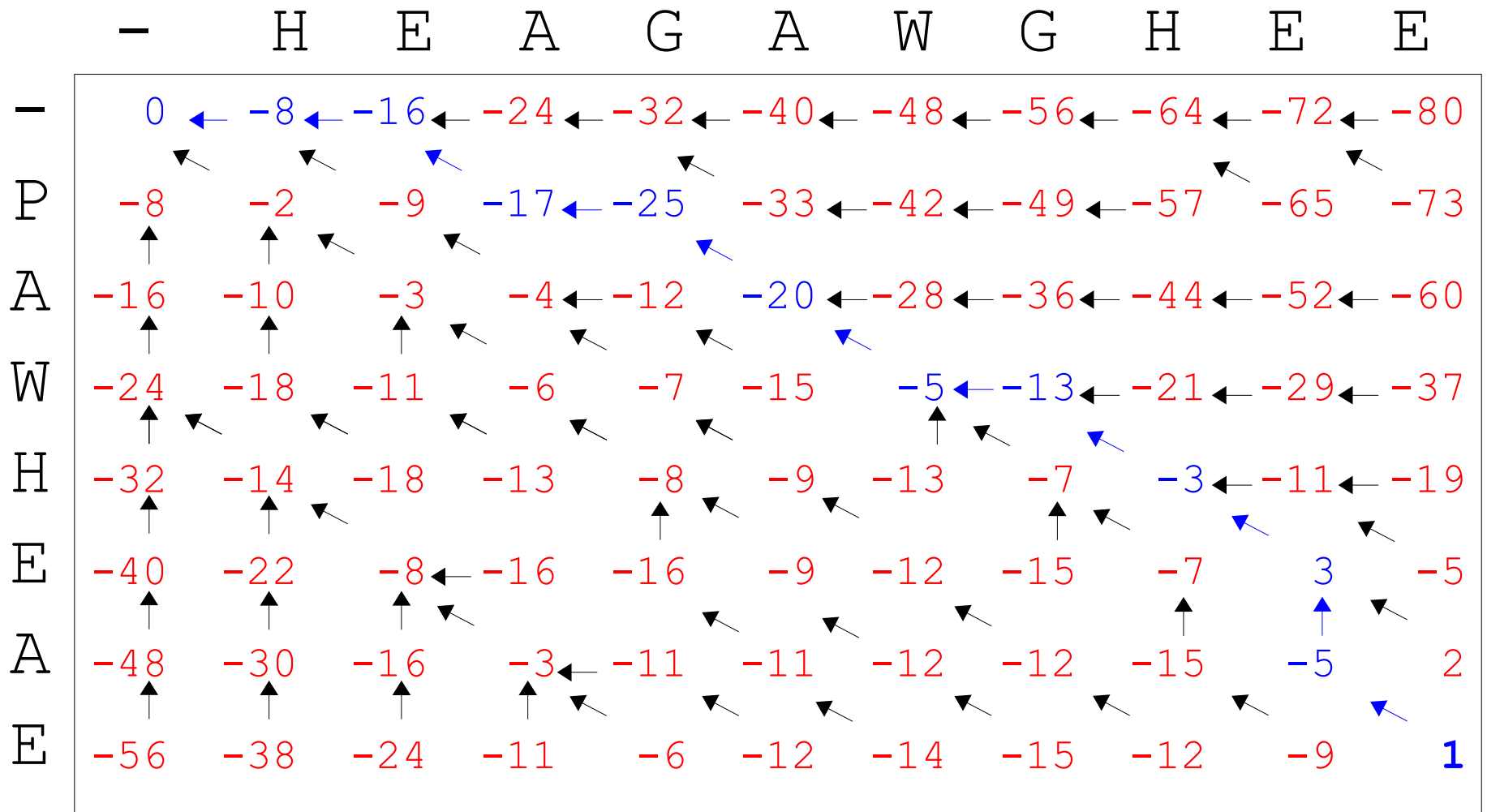
	-	H	E	A	G	A	W
-	0	-8	-16	-24	-32	-40	-48
P	-8	-2	-9				
A							
W	-16						
H	-24						
E	-32						
	-40						

Matrice de similarité

	H	E	A
P	-2	-1	-1
A	-2	-1	5

Coût  $\mu$  pq en mémoire et calcul

# Identification de l'alignement optimal par "tracé inverse"

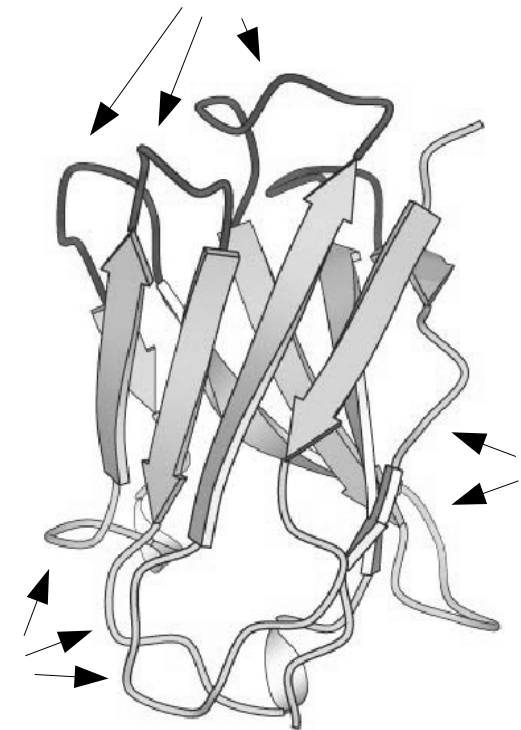


Alignement optimal: HEAGAWGHE-E  
 --P-AW-HEAE

# Il est plus facile d'allonger un gap existant que d'ouvrir un nouveau gap

- Ouvrir un gap:  $-d$
- L'allonger d'un résidu:  $-e$  ( $e < d$ )
- Coût d'un gap de longueur  $m$ :  
$$g(m) = -d - (m-1) e$$

**Eg, boucles**



$\neq$  positions équivalentes et indépendantes

# Programmation dynamique avec 3 états

$A(i,j)$  = score du meilleur alignement  
de  $x_1, \dots, x_i$  avec  $y_1, \dots, y_j$   
qui se termine sur un **appariement**

I	G	A	$x_i$
L	G	V	$y_j$

G	A	-	$x_i$
L	G	V	$y_j$

$I_x(i,j)$  = score du meilleur alignement  
de  $x_1, \dots, x_i$  avec  $y_1, \dots, y_j$   
qui se termine sur une **insertion**

A	I	G	$x_i$
G	V	$y_j$	-

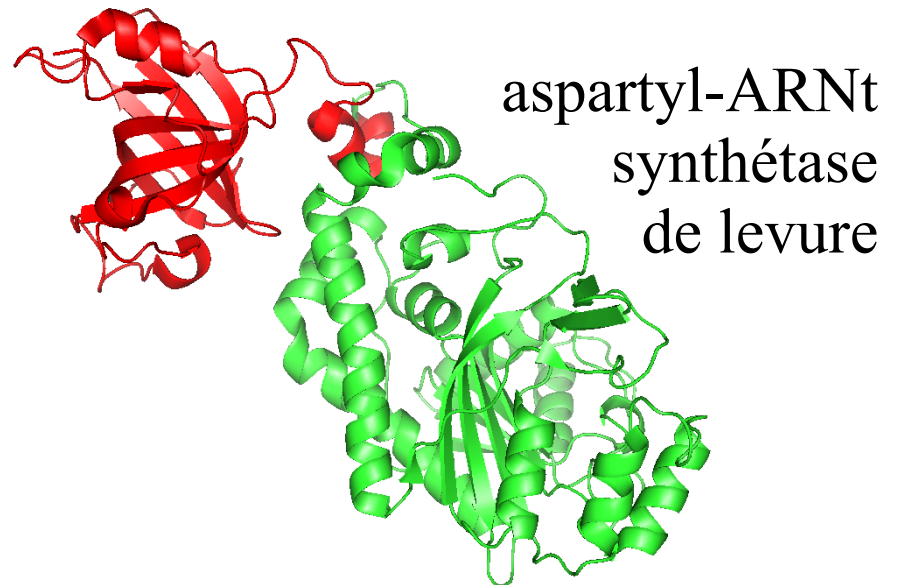
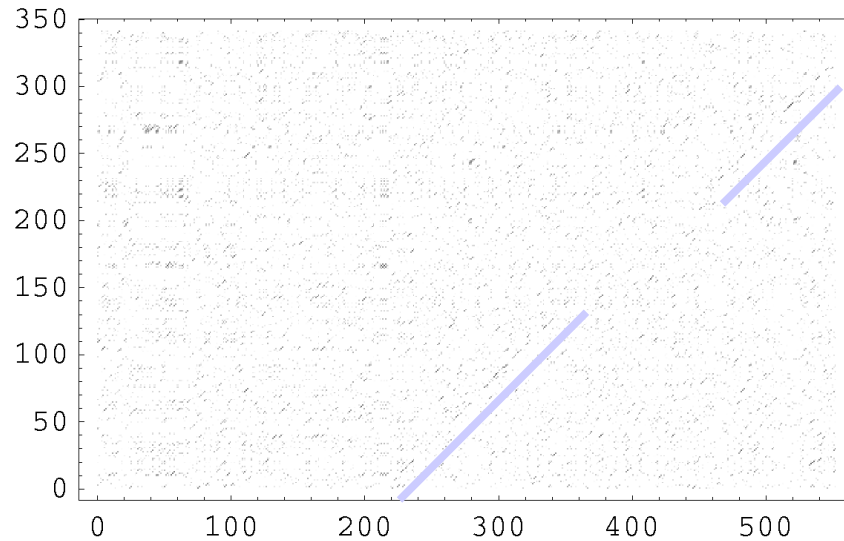
$I_y(i,j)$  = score du meilleur alignement  
de  $x_1, \dots, x_i$  avec  $y_1, \dots, y_j$   
qui se termine sur une **insertion**

G	A	$x_i$	-
S	L	G	$y_j$

**Exercice: trouver l'algorithme correspondant**

# C'est souvent une partie d'une protéine qui est conservée: il faut chercher un alignement "local"

Domaine du site actif, aspartyl-ARNt synthétase d'E coli



Dans cet exemple, chaque protéine possède un domaine absent chez l'autre (un seul est visible dans les vues 3D)

# Alignement local: méthode de Smith-Waterman

$$F(i,j) = \max \left\{ \begin{array}{l} F(i-1,j-1) + s(x_i, y_j) \\ F(i-1,j) - d \\ F(i,j-1) - d \\ 0 \end{array} \right.$$

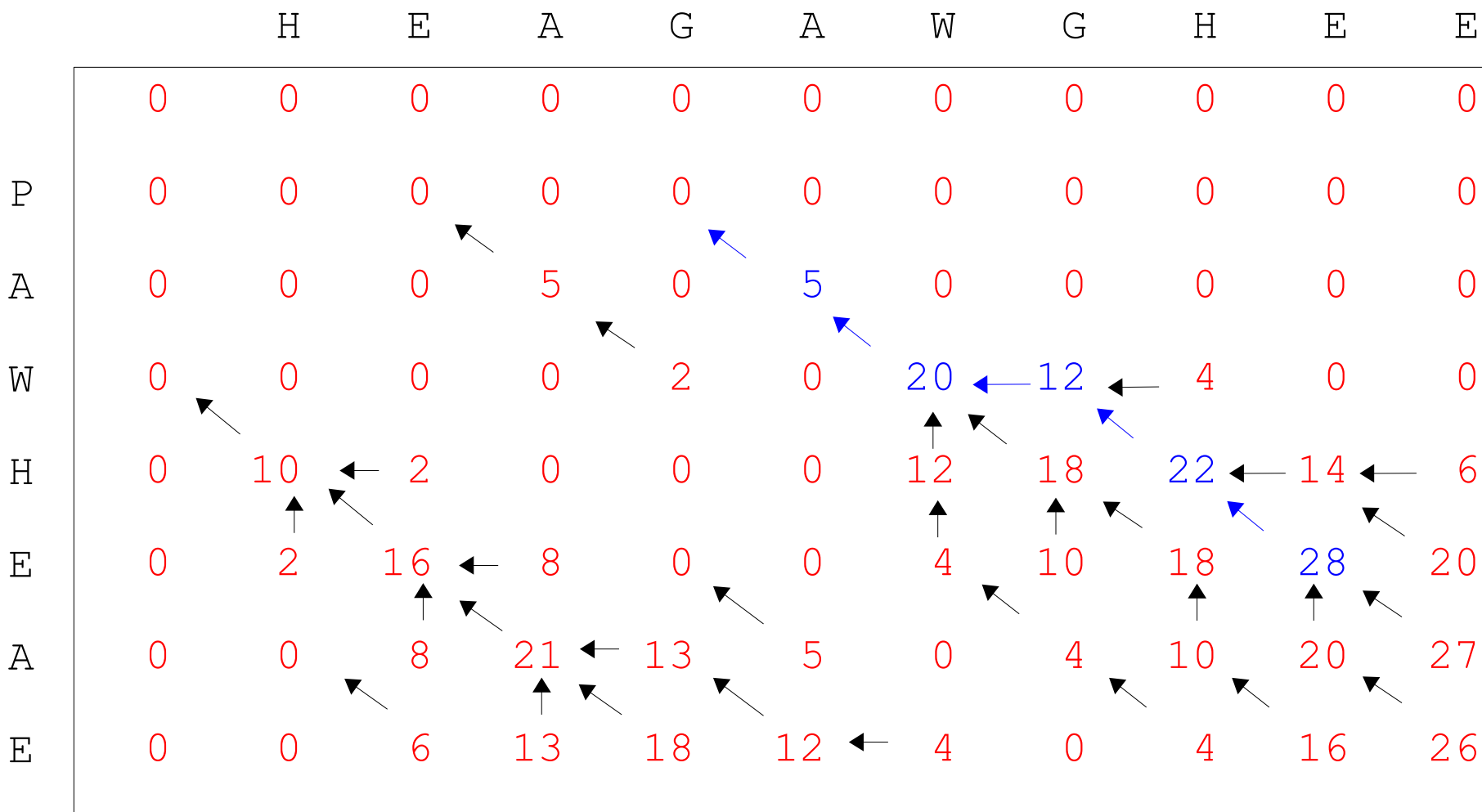
I	G	A	$x_i$
L	G	V	$y_j$

A	I	G	$x_i$
G	V	$y_j$	-

G	A	$x_i$	-	-
S	L	G	V	$y_j$

**Nouvel alignement**

On ne pénalisera pas les morceaux non-alignés en début et fin de séquence.



Alignement local optimal:

AWGHE  
AW-HE



La méthode Smith-Waterman est utilisée en stratigraphie, pour aligner les prélèvements rocheux.

## **En résumé:**

- **Comparaisons de séquences pour prédire la structure, la fonction, et l'origine évolutive des protéines**
- **L'alignement de séquences comme modèle d'un processus évolutif**
- **Alignement de séquences: algorithmes récursifs exacts et algorithmes heuristiques**

# **L'alignement de séquences: algorithmes heuristiques**

# Needleman-Wunsch: “exacte” mais coûteuse

- Atteint le maximum absolu de  $F(p,q)$

## Mais:

- Coût  $\mu$   $pq$  en mémoire et calcul
- Hypothèse de positions équivalentes
- Hypothèse de positions indépendantes:

$$P(\text{alignement}) = \prod_{i,j} P(x_i, y_j)$$

# Méthodes heuristiques rapides

## **BLAST:** *Basic Local Alignment Search Tool*

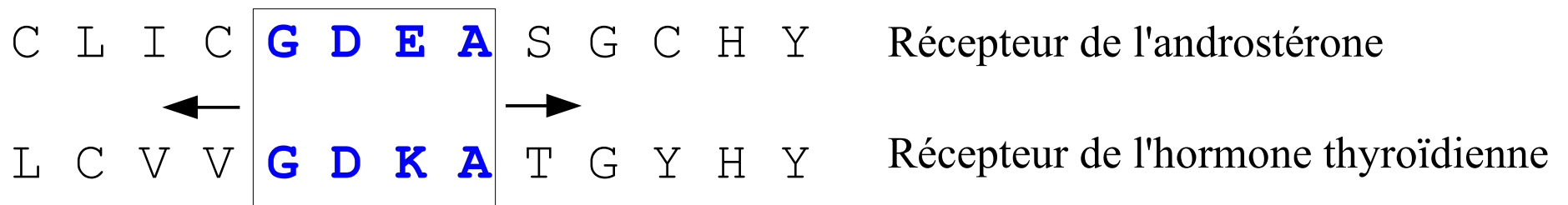
- recherche de térapeptides homologues
- extension de chaque térapeptide tant que la similarité > seuil

## **FASTA**

- recherche de dipeptides homologues
- extension tant que la similarité > seuil
- alignement exact avec gaps, limité à une bande autour de l'extension trouvée

# **BLAST:** *Basic Local Alignment Search Tool*

- recherche de térapeptides homologues
- extension de chaque térapeptide tant que la similarité > seuil



# Tétrapeptides homologues à un tétrapeptide de référence

<b>G</b>	<b>D</b>	<b>E</b>	<b>A</b>	<b>score</b>
				<b>BLOSUM62</b>
<b>G</b>	<b>D</b>	<b>E</b>	<b>A</b>	<b>21 = 6+6+5+4</b>
<b>G</b>	<b>D</b>	<b>D</b>	<b>A</b>	<b>18</b>
<b>G</b>	<b>D</b>	<b>Q</b>	<b>A</b>	<b>18</b>
<b>G</b>	<b>E</b>	<b>E</b>	<b>A</b>	<b>17</b>
<b>G</b>	<b>D</b>	<b>E</b>	<b>G</b>	<b>17</b>
<b>G</b>	<b>D</b>	<b>K</b>	<b>A</b>	<b>17</b>
<b>G</b>	<b>D</b>	<b>E</b>	<b>V</b>	<b>17</b>

Extrait de BLOSUM62

	<b>V</b>	<b>A</b>	<b>G</b>	<b>D</b>	<b>E</b>	<b>Q</b>	<b>K</b>
<b>A</b>	0	4					
<b>G</b>			6				
<b>D</b>				6	2		
<b>E</b>					5	2	1

7 homologues (avec BLOSUM62 et un seuil de 17).

# On identifie dans la séquence y un tétrapeptide homologue à un tétrapeptide de la séquence x

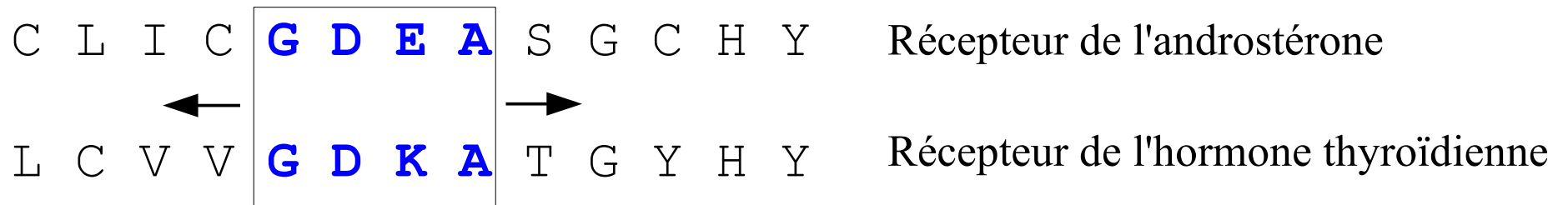
G	D	E	A
G	D	D	A
G	D	Q	A
G	E	E	A
G	D	E	G
G	D	K	A
G	D	E	V

Eg: Tétrapeptides homologues  
à la référence **GDEA**

>sp|P10827|THA\_HUMAN **Thyroid hormone receptor** alpha Homo sapiens.  
MEQKPSKVECGSDPEENSARSPDGKRKRKNGQCSLKTSMGYIPSYLDKDEQCVVCGDKA  
TGYHYRCITCEGCKGFFRRTIQKNLHPTYSCKYDSCCVIDKITRNQCQLCRFKKCIAVGM  
AMDLVLDDSKRVAKRKLIEQNRERRRKEEMIRSLQQRPEPTPEEWDLIHIATEAHRSTNA  
QGSHWKQRRKFLPDDIGQSPIVSMFDGDKVDLEAFSEFTKIITPAITRVVDFAKKLPMS  
ELPCEDQIILLKGCCMEIMSLRAAVRYDPESDTLTLSEGEMAVKREQLKNGGLGVVSDAIF  
ELGKSLSAFNLDDETEVALLQAVLLMSTDRSGLLCVDKIEKSQEAYLLAFEHYVNHHRKHNI  
PHFWPKLLMKEREVQSSILYKGAAAEGRPGGSLGVHPEGQQLGMHVVQGPQVRQLEQQL  
GEAGSLQGPVLQHQSPPKSPQQRLELLHRSGILHARAVCGEDDSSEADSPSSSEEEPEVC  
EDLAGNAASP

Séquence y

**On étend le “micro-alignement” dans  
chaque direction, tant que le score > seuil**



- Chaque térapeptide conduit à un alignement local sans gaps
- On retient les meilleurs alignements

# Homologues du récepteur de l'androstérone identifiés par recherche BLAST dans la banque SwissProt

#	ID Swissprot	Hit	Description	Score	E	% Identity	Match Length
1	P15207	ANDR_RAT	Androgen receptor.	162	1e-40	100	73
6	P19091	ANDR_MOUSE	Androgen receptor.	162	1e-40	100	73
14	Q63449	PRGR_RAT	Progesterone receptor (PR)	136	1e-32	80	72
17	P06401	PRGR_HUMAN	Progesterone receptor (PR)	136	1e-32	80	72
21	P08235	MCR_HUMAN	Mineralocorticoid receptor (MR)	136	1e-32	79	72
33	P04150	GCR_HUMAN	Glucocorticoid receptor (GR)	131	3e-31	77	72
41	Q9YH32	ESR2_ORENI	Estrogen receptor beta (ER-beta)	99	3e-21	58	72
42	Q9YH33	ESR1_ORENI	Estrogen receptor (ER-alpha)	98	4e-21	55	72
:	:	:	:	:	:	:	:
:	:	:	:	:	:	:	:
343	Q9N4Q7	NH13_CAEEL	Nuclear hormone receptor nhr-13	54	8e-08	39	66
344	Q23294	NH11_CAEEL	Nuclear hormone receptor nhr-11	54	8e-08	42	66
345	O45460	NH54_CAEEL	Nuclear hormone receptor nhr-54	54	1e-07	37	67
346	Q09565	NH20_CAEEL	Nuclear hormone receptor nhr-20	51	7e-07	34	66
347	Q09587	NH22_CAEEL	Nuclear hormone receptor nhr-22	45	5e-05	32	66
349	P17672	E75B_DROME	Ecdysone-induced protein 75B	40	0.001	37	47
351	P20659	TRX_DROME	Trithorax protein.	31	0.74	26	49
355	P98164	LRP2_HUMAN	Lipoprotein receptor.	30	1.7	27	65

# **FASTA: intermédiaire entre BLAST et Needleman-Wunsch**

- recherche de dipeptides homologues
- extension tant que la similarité  $>$  seuil
- alignement exact avec gaps, limité à une bande autour de l'extension trouvée

Cf. document écrit

## **En résumé:**

- **Comparaisons de séquences pour prédire la structure, la fonction, et l'origine évolutive des protéines**
- **L'alignement de séquences comme modèle d'un processus évolutif**
- **Alignement de séquences: algorithmes récursifs exacts et algorithmes heuristiques**