

# Fondements de l'annotation par homologie

# Comment les gènes apparaissent?

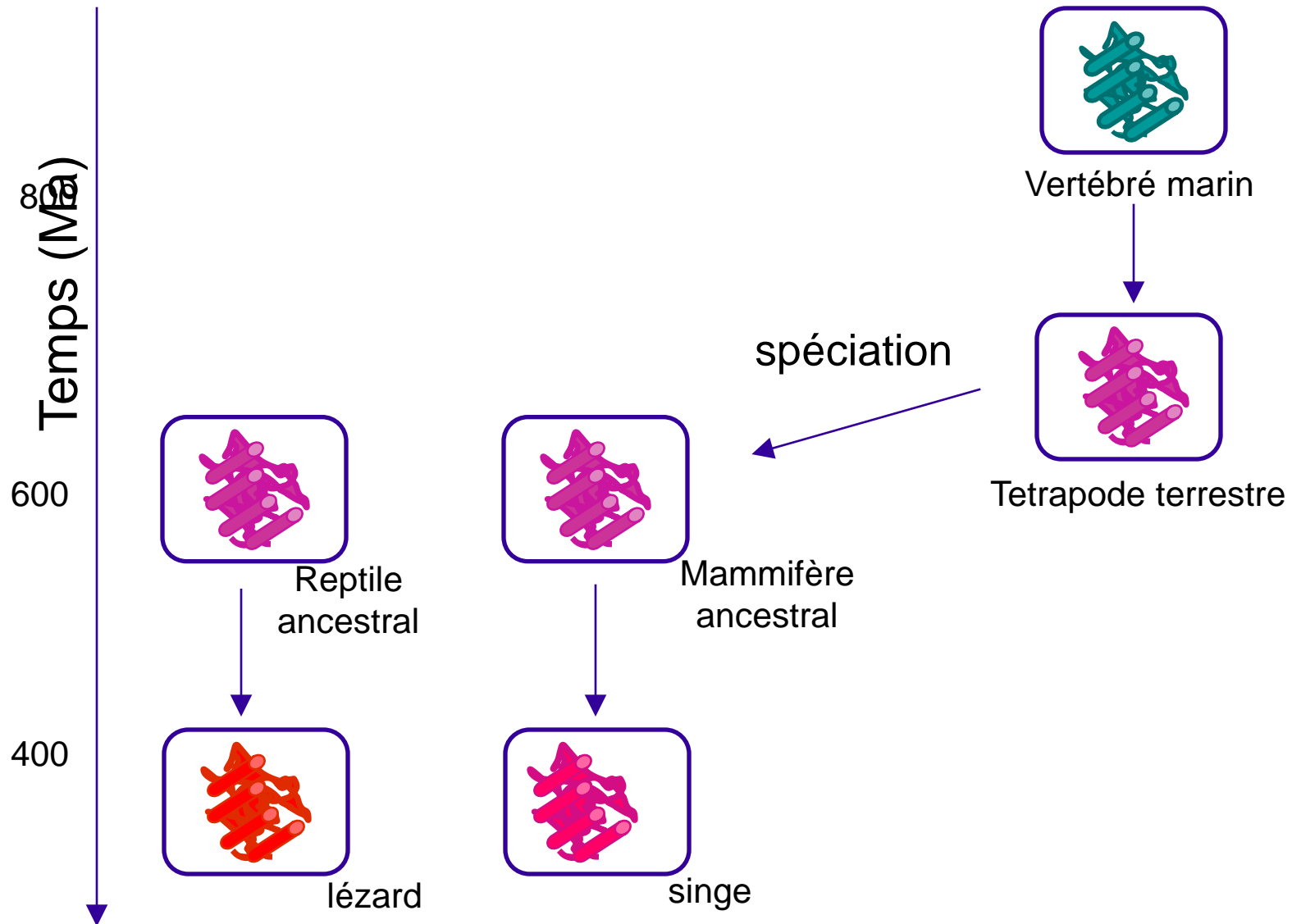
## ✦ Duplications

- ★ Duplication du génome entier ou polyploïdisation (plusieurs cas chez les eucaryotes, par exemple chez les téléostéens, principale classe de poissons)
- ★ Duplication d'un gène ou d'un groupe de gènes (fréquent)
- ★ Duplication d'un chromosome ou d'une partie (rare car délétère)
- ★ La duplication est suivie le plus fréquemment de la perte de gènes: 90% des gènes dupliqués à l'origine des vertébrés auraient été perdus depuis.

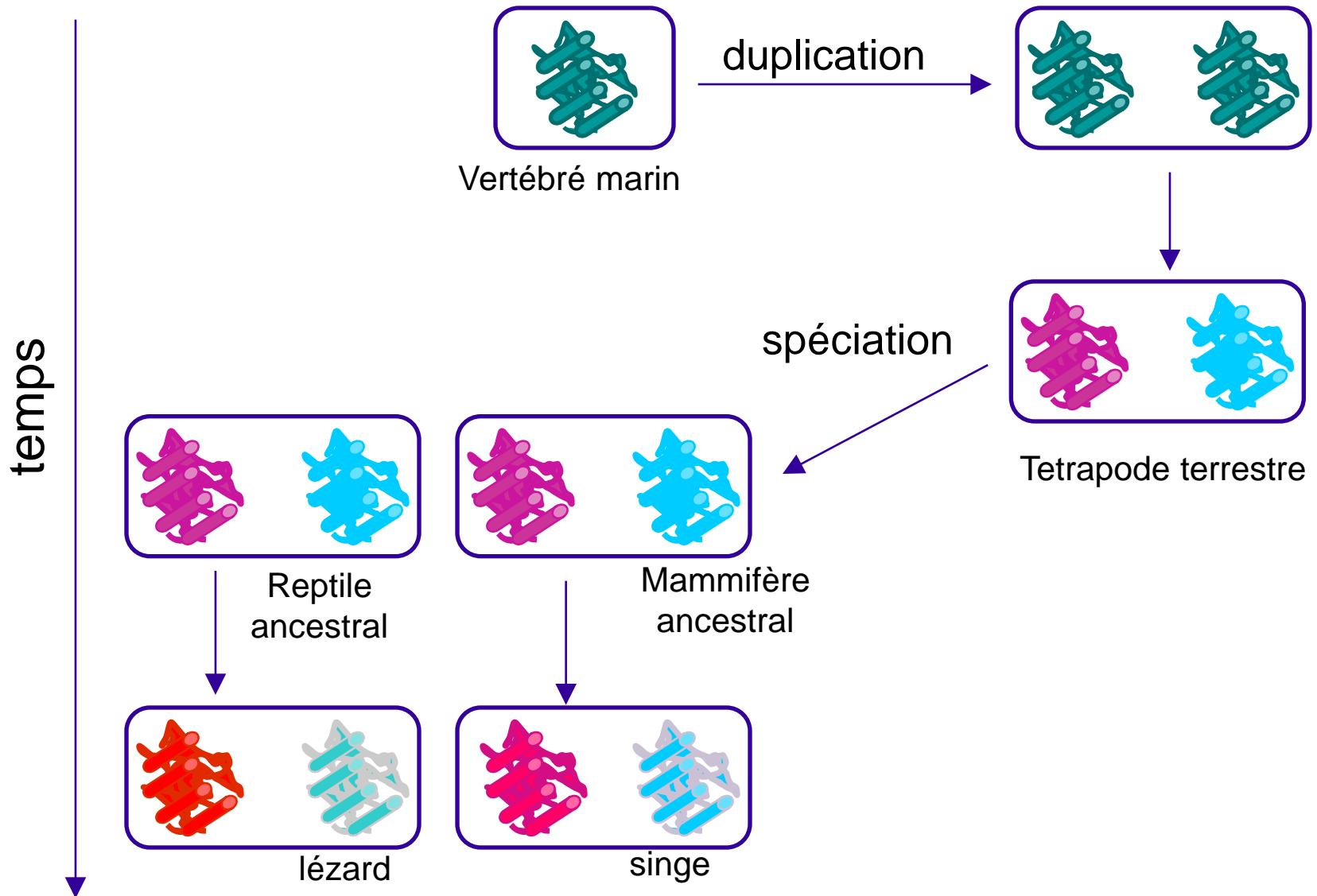
## ✦ Transferts horizontaux

- ★ Très important entre génomes procaryotes
- ★ Survient parfois de procaryote à eucaryote

# Evolution d'un gène

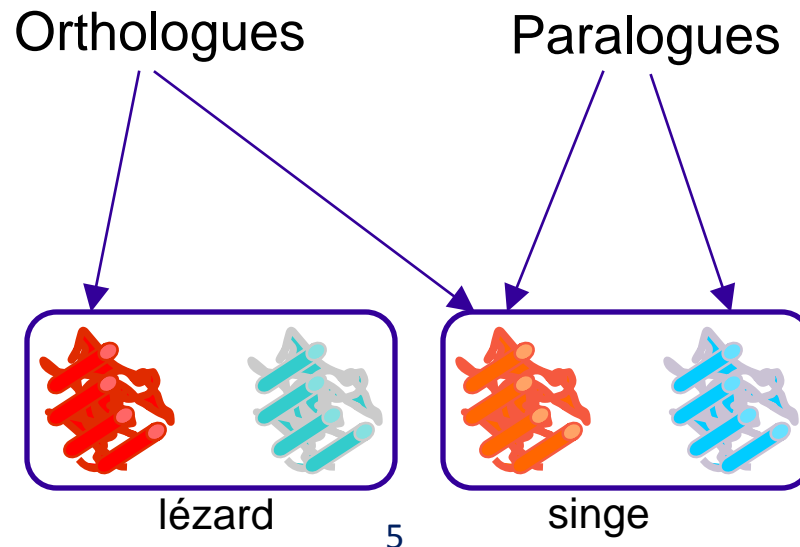


# Apparition de nouveaux gènes par duplication



# Paralogues et orthologues (Fitsch, 1970)

- ★ Homologues: gènes provenant d'un ancêtre commun
- ★ Paralogues: gènes homologues issus d'un phénomène de duplication
- ★ Orthologues: gènes homologues issus de la spéciation
- ★ Transfert horizontal: par endosymbiontes, etc. Fitch a aussi introduit "xénologue" pour évoquer ce cas.

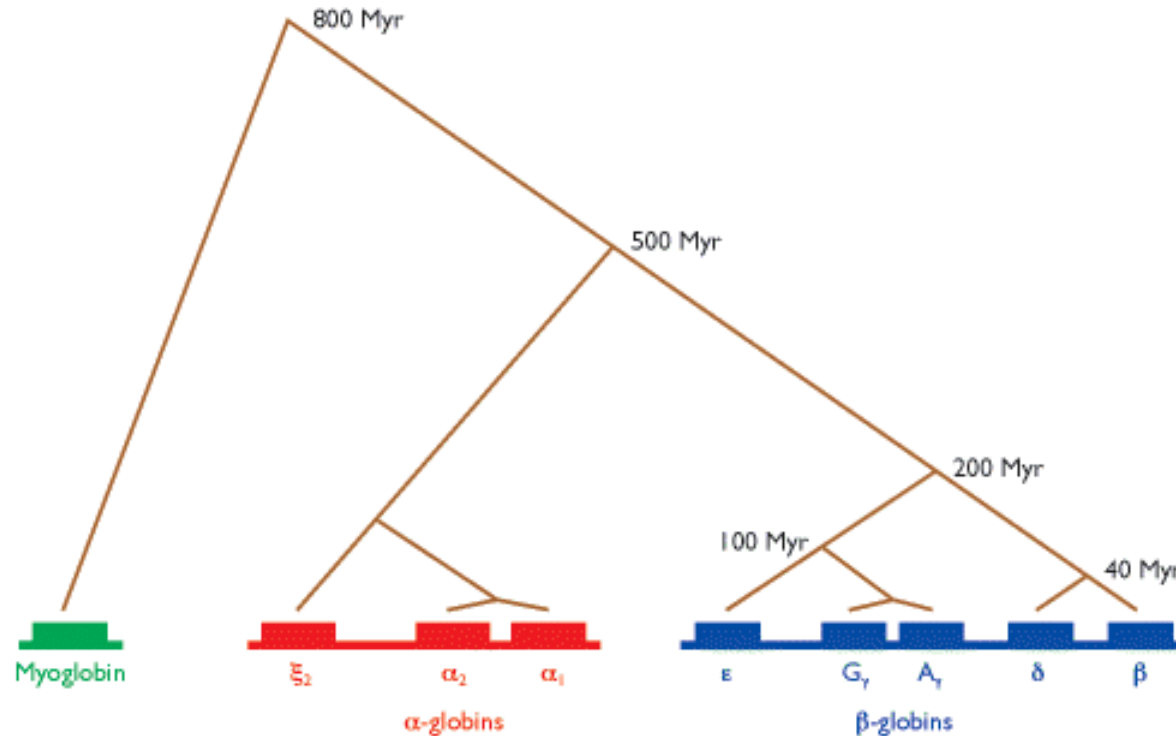


# Fonction et homologie

- ❖ Homologie n'implique pas même fonction: par exemple l'aile de l'oiseau et le bras humain n'ont pas la même fonction
- ❖ Des orthologues rapprochés (p. ex. homme/souris) ont le plus souvent la même fonction dans l'organisme.
- ❖ Des orthologues distants (p. ex. homme/mouche) ont plus rarement le même rôle *phénotypique*, mais peuvent exercer le même rôle dans une *voie* donnée.
- ❖ Les paralogues acquièrent rapidement des fonctions différentes

# Exemple: les gènes de globine humains

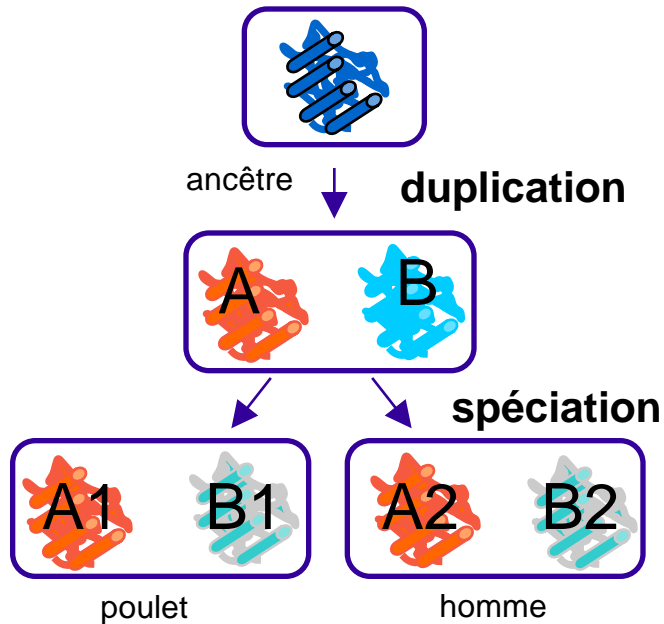
★ Tous paralogues



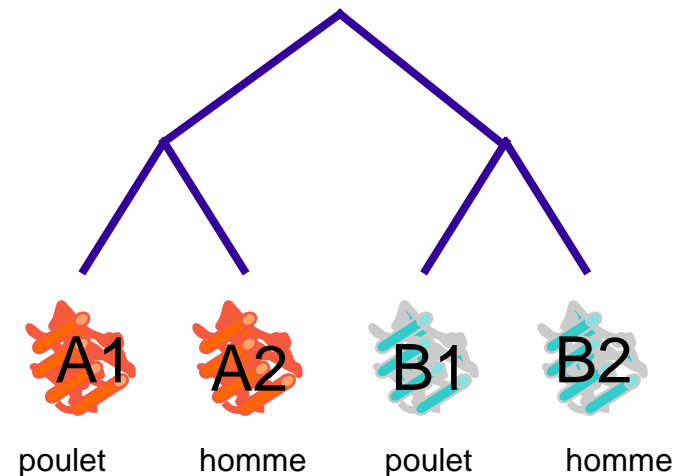
Les gènes se trouvent maintenant sur différents chromosomes: le gène de myoglobine est sur le chromosome 22, les gènes de  $\alpha$ -globines sont sur le chromosome 16 et les gènes de  $\beta$ -globine sont sur le chromosome 11.

# Arbres avec paralogues et orthologues

★ Admettons le schéma évolutif suivant (à gauche) ayant produit deux gènes paralogues présents chez tous les vertébrés.



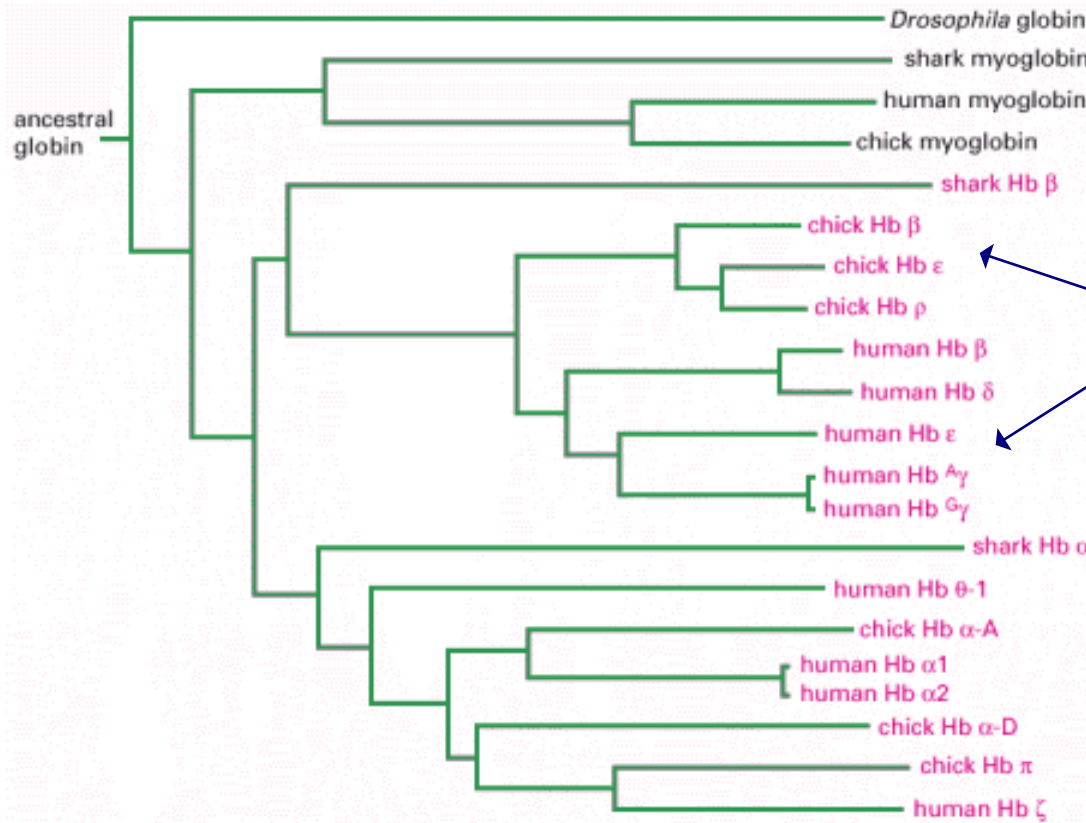
★ Etant donné que la duplication (ayant produit les paralogues) a eu lieu AVANT la spéciation (ayant produit les orthologues), les orthologues devraient être plus proches entre eux que les paralogues. L'arbre devrait donc ressembler à ceci:





# Les gènes de globine chez # espèces

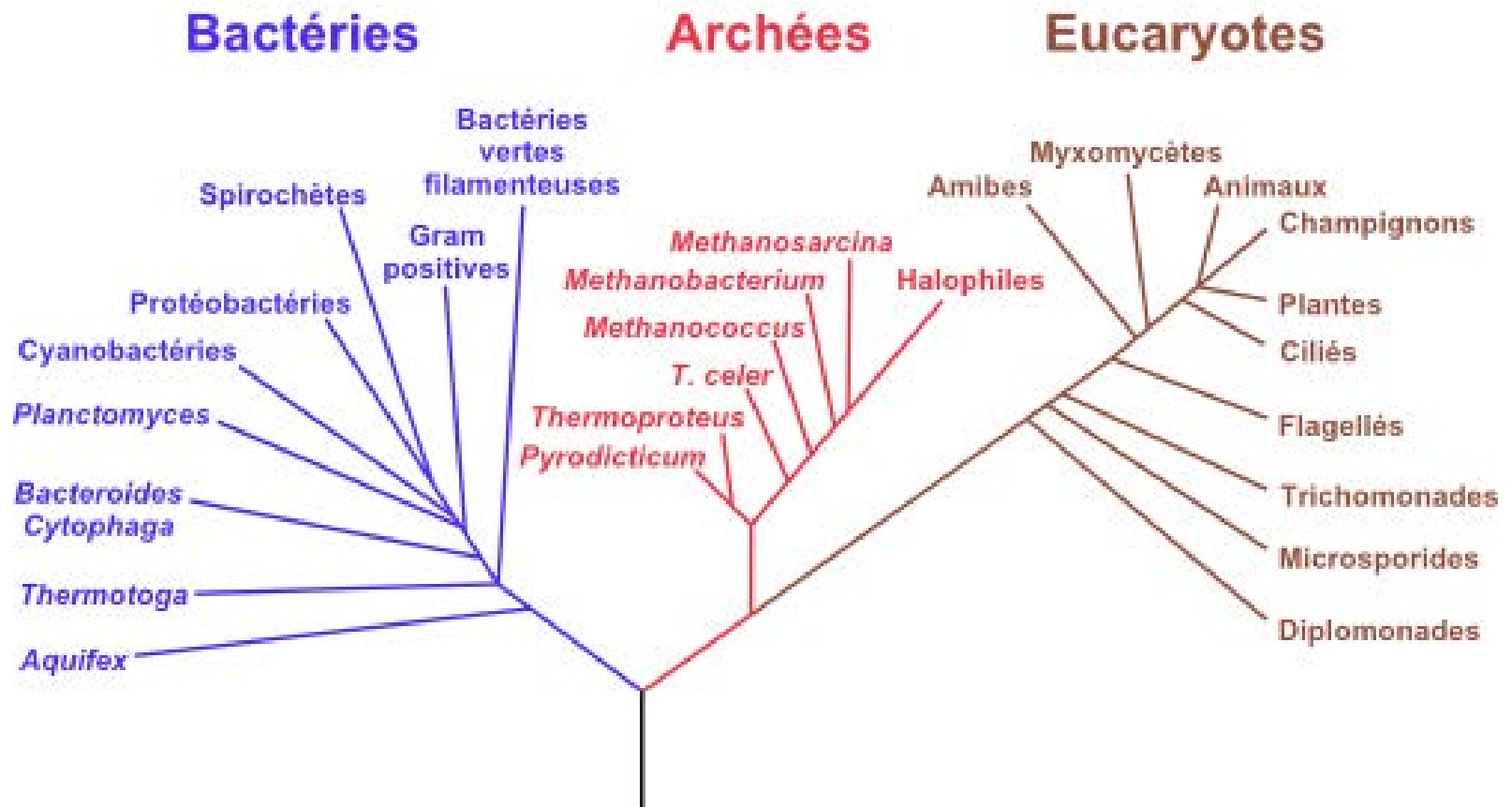
« Outgroup » = groupe extérieur. Indispensable pour placer la racine.



Myoglobines toutes orthologues.

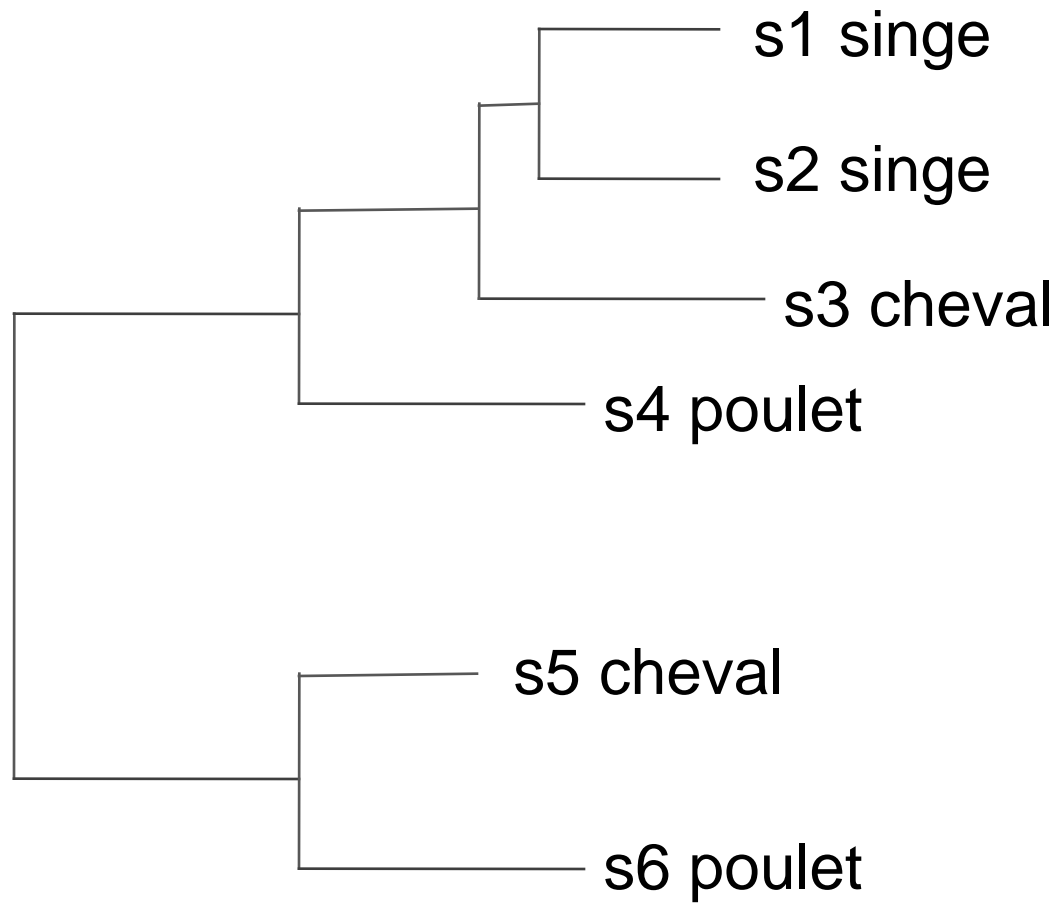
Globines epsilon homme et poulet non orthologues.  
=> Deux évènements de duplication

# Importance de la phylogénie

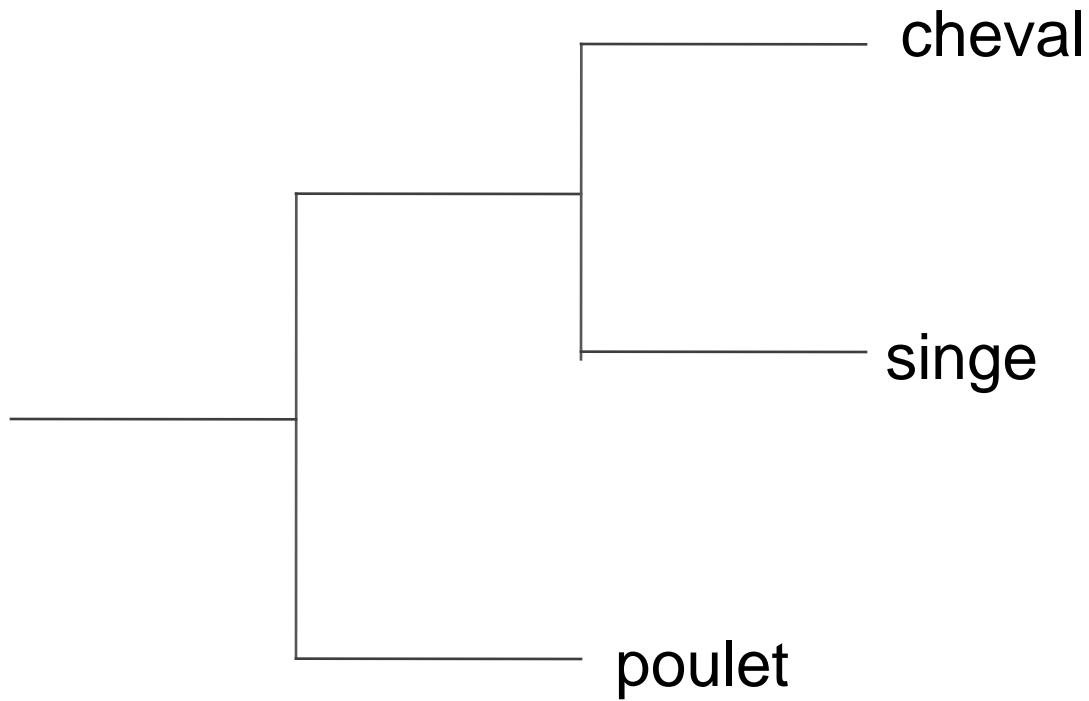




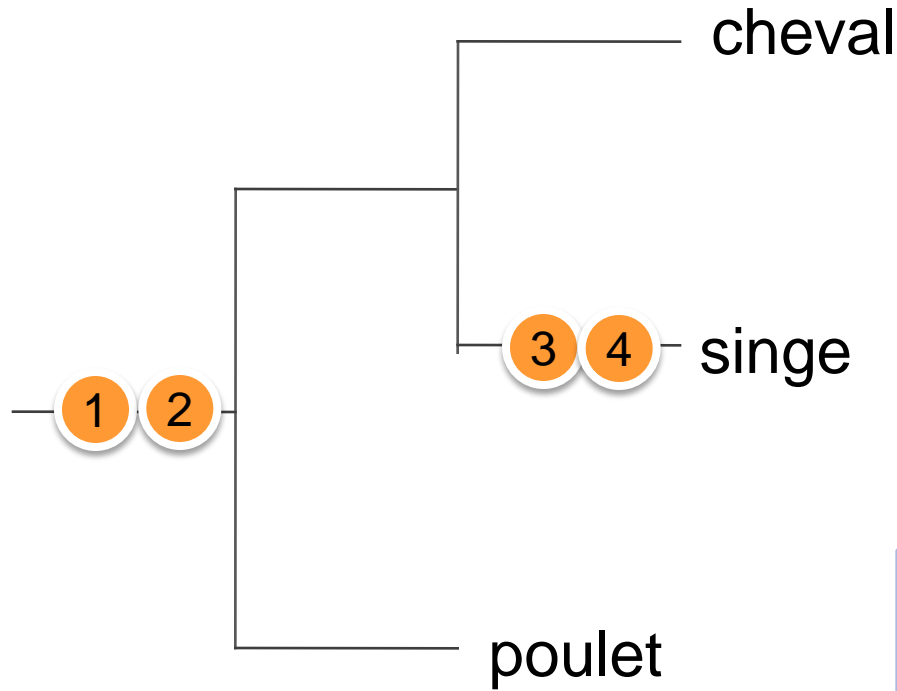
# Exercice



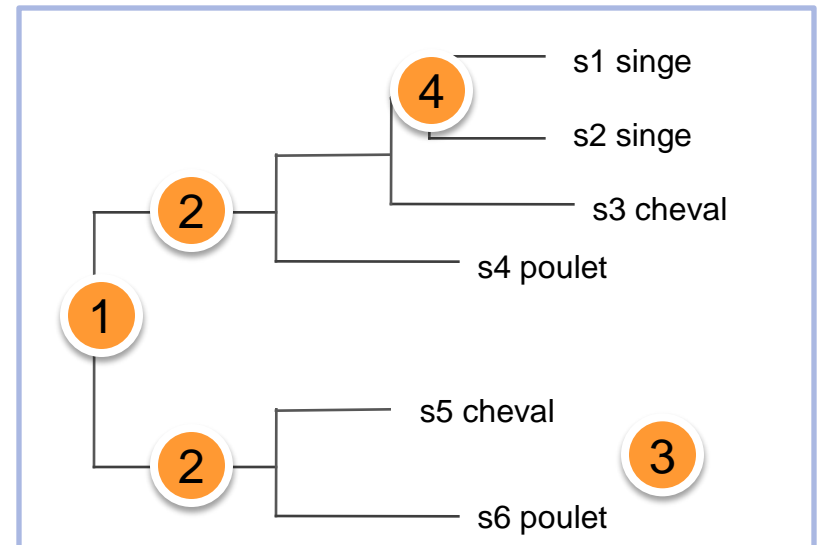
# D'abord faire l'arbre des espèces



# Puis placer les évènements



- 1 Apparition du gène ancestral P+C+S
- 2 Duplication chez ancêtre mamm+oiseaux
- 3 Perte d'une copie chez le singe
- 4 Duplication « récente » chez le singe



# Vaut-il mieux comparer les protéines ou l'ADN pour rechercher des homologues d'une séquence?

- ★ La meilleure façon de détecter des similitudes entre séquences est généralement la *comparaison au niveau protéique*.
  1. Il existe 20 aa contre 4 bases. La probabilité de trouver une "lettre" donnée par hasard est donc plus importante pour les bases.
  2. Plusieurs codons produisent le même aa. 134 / 549 substitutions de bases sont synonymes. Les séquences protéiques sont plus informatives.
  3. La raison principale est en fait l'existence d'outils de comparaison plus puissants pour les aa: utilisation des propriétés physicochimiques ou des substitutions observées dans l'évolution. Même lorsque les aa sont différents, on est capable de retrouver des similitudes. On en est tout à fait incapable au niveau des bases.
- ★ Il existe en fait des cas où la séquence d'ADN est plus conservée que la séquence protéique, ce qui enlève du poids à l'argument 1
- ★ Les comparaisons avec les séquences protéiques ne permettent de détecter que les régions codantes. Evidemment, on utilisera toujours la séquence ADN/ARN pour analyser ce qui n'est pas traduit!

# L'analyse des domaines et de l'alignement multiple



# Alignement multiple

- ✦ **Pourquoi réaliser un alignement multiple?**
- ✦ L'alignement multiple révèle des fonctions que l'on ne pouvait pas visualiser en comparant 2 séquences
  - ★ Identifier les positions et les acides aminés importants.
  - ★ Visualiser les domaines
  - ★ Distinguer paralogues et orthologues
  - ★ Etablir la phylogénie des séquences, et même parfois des organismes
  - ★ Comme une aide à la modélisation: Les algorithmes de prédiction de structures secondaires exploitent beaucoup mieux les alignements multiples. Connaître les acides aminés permis à telle ou telle position facilite l'inférence 3D.

# Exemple d'analyse d'alignement multiple: les Glutamine Aminotransferases (GAT)

Horvath & Grishin, Proteins, 2001

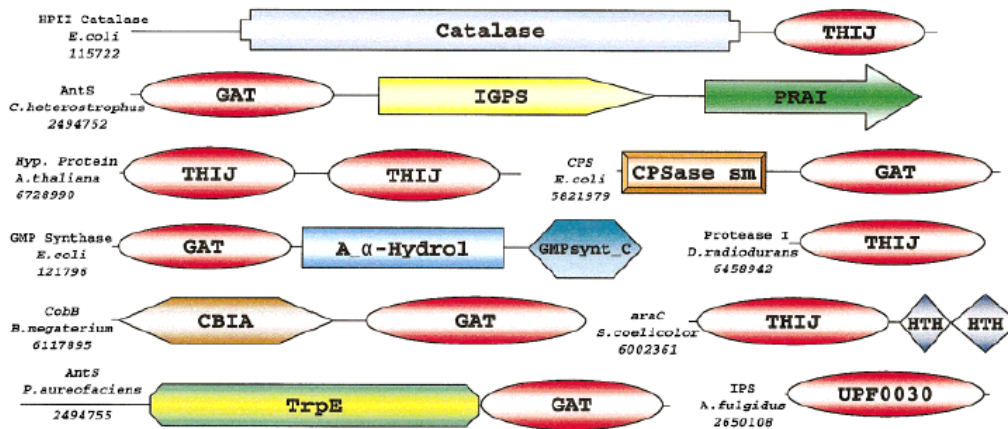
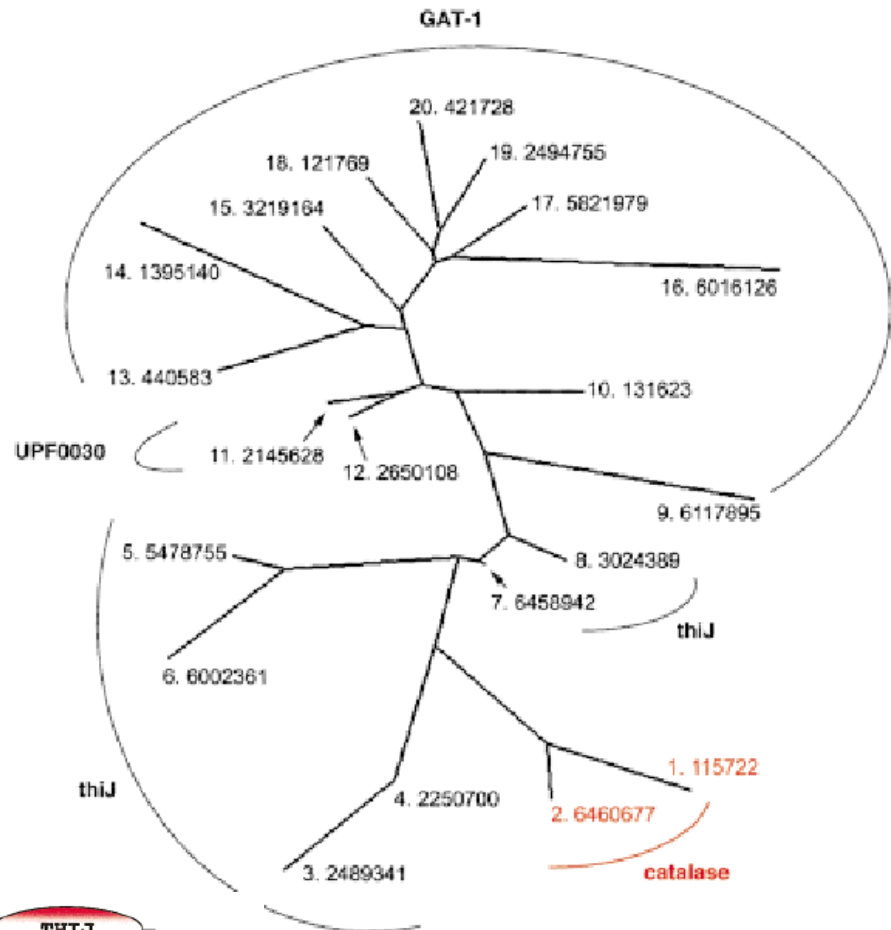
Accession	Gene	Residues	Protein	Sequence	Consensus	SS	βa	αA	βb	αB	βc																																												
1.	CAT	115722	753	599	GRVVAITLNDNDEVR	----	SADLLA	ILKAL	KAKGV	----	HAKLLYSRM	[14]ATFAGAPSLTVDAVIVPCGNIAADIADN	678																																										
2.	CAT	6460677	772	590	GRKVAVLVADGV	----	AAGVKAL	QDAL	KKADV	----	KYDIVAPHL	[10]ATLSNTDPVVYDGVVVVAGGAAAVRELA	662																																										
3.	THI	2498341	270	35	NTNIAVVVPSGCGW	[ 5]	IHEAA	YTM	YHL	SRN	GA	----	RQIPAPNQ	[37]NDLSKLDANSFDAVIFPGGHGIVKMS	141																																								
4.	THI	2250700	226	2	AARVALVLSGGCV	[ 5]	IHEASA	ILV	HLS	SRG	GA	----	EVOIFAPDV	[33]TDLANLSAANHDAATFPGGFGAAKNLS	144																																								
5.	THI	5478755	189	3	SKRALVILAKGAE	----	EMETV	IFVD	IM	RR	GI	----	KVTVAGLAG	[16]SLEBAKTQGPYDVVVVLPGGNLGAQNLS	83																																								
6.	THI	6002361	327	10	PHRVVVVGFVDFG	MK	----	LLDL	SGP	AEV	FSE	ANR	[5]RLSIVSADG	[12]ADTDARAAAAAHTLVVVVGGDALPGSPV	91																																								
7.	THI	6458942	190	9	GKKIATLLAADG	VE	----	EIEL	TS	PRA	IEA	CG	----	TTELSLEP	[19]HVVSEVQVSDYDGLLLPGGTVPNDKLR	92																																							
8.	THI	3024389	166	1	MKILFLSANEF	----	DVEL	L	YP	HRL	KEE	GH	----	EVYIASFEK	[14]LTFDEVNPEFDALVLPGGRAPEVRVL	78																																							
9.	GAT	6117895	486	287	RRRVAMASGA	AF	----	TFSYA	EHTE	LLAAA	GA	----	EVVTFDP	----	LRDEELPECTQGLVIGGGFPVEYASE	347																																							
10.	GAT	131623	227	2	KFAVIVLPGS	SNC	----	DIDM	YH	AVK	DEL	GH	----	EVEYVWH	----	EETSLDGF	DVGLIPGGFSGDYLR	58																																					
11.	UPF	2145628	219	24	FPRVGVLLAQ	G	----	DTRE	HL	TAL	REA	GA	----	DSMPVRR	----	RCELDEV	DALVLPGGESTTISHL	78																																					
12.	UPF	2650108	198	1	MKVAVVGVQ	GDV	----	EEHV	L	TKR	AL	GI	----	DGEVVT	----	RRRGV	SRSDAVILPGGESTTISKL	60																																					
13.	GAT	440583	579	303	TVKIRLVGKY	TNL	----	KDSV	L	SVI	KALE	HSS	SM	[6]	DIKWVEATD	[12]PHEAMN	MVSTADGILIPGGFCVRCGT	393																																					
14.	GAT	1395140	242	2	SKRFALLWC	SEEE	[ 1]	FDY	RE	MV	NA	KTENS	----	DWEV	ISAF	----	TDLNK	IIDNYDGFVISGSEYSVNADK	65																																				
15.	GAT	3219164	593	61	DSVVTLLDY	GAG	----	NVR	S	IR	NAL	RHL	GF	----	SIKDVCT	----	PGDIL	NADRLIFPGVGFPA	PAMD	116																																			
16.	GAT	6016126	326	32	QTGVVYSD	HPGNC	[11]	PSIA	A	S	VKL	ABS	GA	----	FVIL	IFNEP	----	GBIL	FQRL	ELVNGVILTGGWAKEGLY	107																																		
17.	GAT	5821979	382	191	PFHVVA	YDFGA	----	KRN	L	R	ML	VDR	GC	----	RLT	IVPAQ	----	TS	AE	DV	LKMN	PDGIFLSNPGPDPAFC	248																																
18.	GAT	121769	525	7	KHRILLDF	PGG	----	QY	T	L	V	RRV	REL	GV	----	YCEL	W	AWD	----	VTE	AQ	IRD	PNP	GGIILSGGPESTTE	65																														
19.	GAT	2494755	637	435	GRQVLIV	DAEDT	----	FTS	M	I	A	K	Q	AL	GL	----	VVT	V	CSF	----	SDE	Y	S	F	E	G	Y	D	L	V	M	G	P	G	P	N	P	S	E	V	Q	492													
20.	GAT	421728	195	1	MDLTLI	IDNYD	----	SFV	Y	N	I	A	Q	V	G	E	L	S	----	Y	P	I	V	I	R	N	D	----	I	S	I	K	G	I	E	R	I	D	P	D	R	L	I	S	P	G	P	C	T	P	E	K	R	E	62

Accession	Gene	Residues	Protein	Sequence	Consensus	SS	βa	αA	βb	αB	βc																																																																
1.	CAT	115722	679	----	GD	NY	Y	L	M	E	K	H	L	----	K	P	I	A	L	A	G	D	A	R	K	F	K	-	A	T	I	[ 8]	----	G	I	V	E	A	D	-	S	A	-	----	G	S	F	N	D	E	L	L	T	M	A	A	----	739																	
2.	CAT	6460677	663	----	Q	H	P	E	S	F	N	P	V	V	C	S	V	R	H	A	----	K	P	I	G	S	L	G	B	A	E	I	V	-	T	G	S	[ 8]	----	V	A	A	D	S	P	-	A	K	-	----	G	A	T	A	P	Q	N	L	S	E	V	A	G	V	R	L	A	729							
3.	THI	2498341	142	[10]	N	H	S	E	R	V	L	K	D	F	H	R	A	----	K	P	I	G	L	S	S	A	M	P	L	L	A	C	R	V	L	----	P	S	L	E	V	T	M	G	Y	E	R	D	E	S	S	R	W	G	R	W	P	N	T	M	V	Q	A	V	K	S	M	G	A	218					
4.	THI	2250700	145	[10]	V	N	K	E	V	E	R	V	L	K	E	P	H	Q	A	C	----	K	P	I	G	L	C	I	A	P	V	L	A	A	K	V	L	----	R	G	V	E	V	T	V	G	H	E	Q	E	----	G	G	K	W	P	Y	A	G	T	A	E	A	I	K	A	L	G	A	218					
5.	THI	5478755	84	----	E	S	A	L	V	K	E	I	L	K	E	Q	E	N	R	K	----	G	L	I	A	A	C	A	G	P	T	A	L	L	A	H	E	[ 31]	K	D	G	L	I	L	T	S	R	G	F	G	T	S	----	F	E	F	A	L	A	I	V	E	A	L	S	187									
6.	THI	6002361	92	----	D	P	V	L	G	A	A	K	E	L	A	E	R	A	----	G	R	V	A	S	V	C	T	G	A	F	V	L	G	A	A	-	[ 34]	K	D	C	S	T	Y	T	S	A	C	V	T	A	G	----	I	D	L	A	L	L	E	E	D	H	G	237											
7.	THI	6458942	93	----	L	E	E	G	A	M	K	F	V	R	D	M	Y	D	A	G	----	K	P	I	A	A	I	C	H	G	P	N	S	L	S	E	T	G	[ 30]	I	D	K	G	V	V	T	S	R	K	P	D	D	----	P	A	F	N	K	I	V	E	E	F	A	E	182									
8.	THI	3024389	79	----	N	E	K	A	V	E	I	A	R	K	M	F	T	E	G	----	K	P	V	A	T	I	C	H	G	P	O	I	L	I	S	A	G	[ 31]	V	D	G	N	V	S	S	R	H	P	G	D	----	Y	A	W	M	R	E	F	U	K	L	L	K	-	166										
9.	GAT	6117895	348	[ 2]	A	N	E	C	L	R	K	S	V	A	E	L	A	F	S	G	----	A	P	V	A	E	C	A	G	L	L	Y	L	C	R	E	L	[ 74]	E	R	G	V	H	A	S	Y	T	H	T	-	H	W	A	----	A	E	P	G	V	A	R	R	F	V	E	R	C	R	T	485					
10.	GAT	131623	59	[ 5]	R	F	A	N	I	M	P	A	V	K	O	A	A	E	A	G	----	K	P	V	L	G	V	C	N	G	F	O	I	L	O	E	L	G	[ 88]	K	G	N	V	L	G	M	M	P	H	P	-	E	R	A	V	-	D	E	L	L	G	S	A	D	G	L	K	L	F	O	S	I	V	K	217
11.	UPF	2145628	79	[ 1]	L	D	C	E	L	L	E	P	L	R	A	R	L	A	D	G	----	L	P	A	G	A	C	T	G	M	I	L	L	A	S	E	[ 76]	Q	G	S	M	L	A	T	A	P	H	P	-	E	M	T	----	S	D	R	R	I	H	E	L	F	V	D	I	V	N	216							
12.	UPF	2650108	61	[ 1]	F	S	D	G	I	A	E	I	L	Q	A	E	E	G	----	K	P	V	M	G	T	C	A	G	L	I	L	S	K	-	[ 75]	Q	K	N	V	L	G	L	A	P	H	P	-	E	L	T	----	D	D	T	R	I	H	E	F	F	L	K	I	G	E	196									
13.	GAT	440583	384	----	E	C	M	V	L	A	A	R	W	A	R	E	N	H	----	I	P	L	G	V	C	L	L	Q	I	A	T	I	E	F	[111]	H	P	Y	I	A	T	Q	V	H	P	-	E	Y	T	S	-	K	V	L	D	P	S	K	P	F	L	G	L	V	A	A	S	558							
14.	GAT	1395140	66	[ 1]	K	F	S	G	L	F	E	F	I	R	A	V	H	K	K	E	----	K	P	I	V	G	L	C	F	G	C	S	L	A	V	A	L	[ 69]	G	P	Y	A	N	C	I	S	G	H	P	-	E	I	S	----	K	K	T	L	E	Q	D	F	L	R	V	H	L	E	D	G	N	199			
15.	GAT	3219164	117	[ 2]	N	R	T	C	M	A	E	A	L	C	K	Y	L	E	N	D	----	R	P	L	G	I	C	L	G	L	L	P	D	S	[ 85]	R	C	N	V	H	A	V	O	D	H	P	-	E	K	S	----	C	E	V	G	L	S	V	L	R	R	F	L	H	P	K	L	P	267						
16.	GAT	6016126	108	----	F	E	I	V	K	K	I	L	N	K	V	L	E	R	N	[ 5]	F	T	A	I	C	I	G	E	L	L	T	M	E	T	[ 92]	K	Y	P	V	T	G	F	Q	W	H	P	-	E	K	N	[12]	-	E	D	A	I	Q	V	T	Q	H	A	A	N	H	L	V	278							
17.	GAT	5821979	249	----	D	T	A	N	A	T	Q	K	E	E	T	S	----	I	P	V	F	G	I	C	L	G	H	Q	L	L	A	L	A	S	[ 62]	D	K	P	A	F	S	F	Q	H	P	-	E	A	S	P	-	G	H	D	A	L	F	D	H	P	I	E	L	E	Q	376									
18.	GAT	121769	68	----	E	N	S	P	R	A	F	Q	V	F	E	A	G	----	V	P	V	F	G	V	C	G	M	Q	T	M	A	M	Q	L	[ 75]	E	K	R	F	Y	G	V	Q	H	P	-	E	V	T	----	H	T	R	Q	M	R	L	E	R	F	V	R	D	I	201										
19.	GAT	2494755	493	[ 2]	K	I	N	H	L	H	V	A	I	R	S	L	S	Q	Q	----	R	P	L	A	V	C	L	S	H	Q	V	L	S	L	C	L	[ 62]	G	P	S	F	A	S	M	Q	H	A	-	E	S	L	----	L	T	Q	E	G	P																	

# Mobilité des domaines entre protéines différentes

## Exemple des Glutamine Aminotransferases (GAT):



100 amino acids

... de l'importance de séparer les domaines pour l'analyse

# Banques de domaines

## ✦ Prosite

Expression régulière

SIGMA70\_2, PS00716, Sigma-70 factors family signature 2 (PATTERN)

Consensus pattern: [STN]-x(2)-[DENQ]-[LIVMT]-[GAS]-x(4)-[LIVMF]-[PSTG]-x(3)-[LIVMA]-x-[NQR]-[LIVMA]-[EQH]-x(3)-[LIVMFWK]-x(2)

## ✦ PFAM

There are 346 sequences with the following architecture: DUF1785, PAZ, Piwi

[AGOL\\_ARATH](#) [Arabidopsis thaliana (Mouse-ear cress)] Argonaute-like protein At2g27880 (997 residues)



[Show](#) all sequences with this architecture.

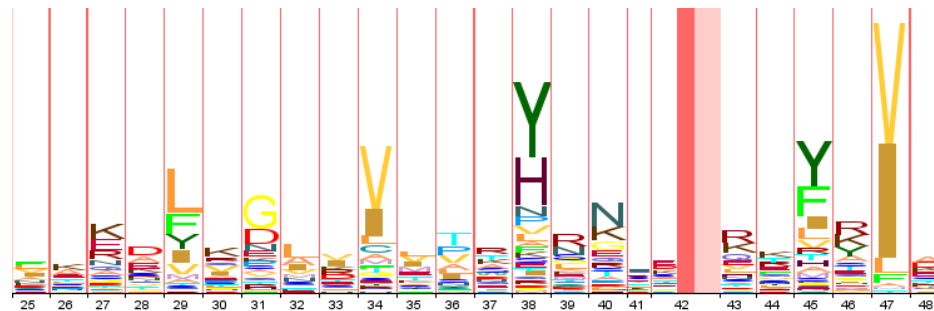
There are 232 sequences with the following architecture: PAZ, Piwi

[PIWL2\\_HUMAN](#) [Homo sapiens (Human)] Piwi-like protein 2 (973 residues)



[Show](#) all sequences with this architecture.

HMM logo



# Informations précieuses présentes dans les banques de domaines

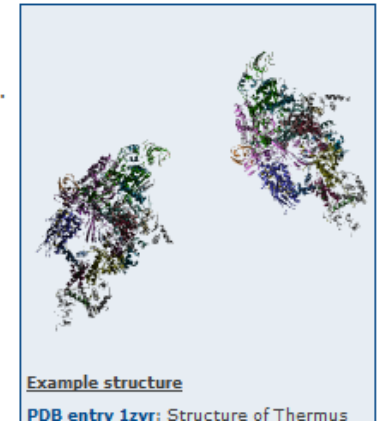
## Sigma-70 region 2

Add annotation

Region 2 of sigma-70 is the most conserved region of the entire protein. All members of this class of sigma-factor contain region 2. The high conservation is due to region 2 containing both the -10 promoter recognition helix and the primary core RNA polymerase binding determinant. The core binding helix, interacts with the clamp domain of the largest polymerase subunit, beta prime [1,2]. The aromatic residues of the recognition helix, found at the C-terminus of this domain are thought to mediate strand separation, thereby allowing transcription initiation [1,2].

### Literature references

1. Campbell EA, Muzzin O, Chlenov M, Sun JL, Olson CA, Weinman O, Trester-Zedlitz ML, Darst SA; , Mol Cell 2002;9:527-539.: Structure of the bacterial RNA polymerase promoter specificity sigma subunit. [PUBMED:11931761](https://pubmed.ncbi.nlm.nih.gov/11931761/)
2. Malhotra A, Severinova E, Darst SA; , Cell 1996;87:127-136.: Crystal structure of a sigma 70 subunit fragment from E. coli RNA polymerase. [PUBMED:8858155](https://pubmed.ncbi.nlm.nih.gov/8858155/)



Exemple d'annotation de domaine dans PFAM

# Interpro: la banque de domaines intégrée de l'EBI

EBI > Databases > InterPro

Search InterPro:



Jump to: [InterProScan](#) [Databases](#) [Documentation](#) [FTP site](#) [Help](#) [Advanced search](#)

## IPR005034 Dicer double-stranded RNA-binding fold

### Protein matches

<b>UniProtKB Matches:</b> 379 proteins	Overview:	<a href="#">sorted by AC</a> , <a href="#">sorted by name</a> , <a href="#">of known structure</a> , <a href="#">proteins with splice variants</a>
	Detailed:	<a href="#">sorted by AC</a> , <a href="#">sorted by name</a> , <a href="#">of known structure</a> , <a href="#">proteins with splice variants</a>
	Table:	<a href="#">For all matching proteins</a> , <a href="#">of known structure</a>
	<a href="#">Architectures</a> <a href="#">Accession List</a> <a href="#">Matches in BioMart</a>	
<b>Accession</b>	IPR005034 Dicer_dsRNA_binding_fold	

<b>Type</b>	Domain
-------------	--------

<b>Signatures</b>	Database	ID	Name	Proteins
	<a href="#">Pfam</a>	<a href="#">PF03368</a>	dsRNA_bind	374
	<a href="#">PROSITE profile</a>	<a href="#">PS51327</a>	DICER_DSRBF	369
<a href="#">Signatures in BioMart</a>				

### GO Term annotation

<b>Function</b>	<a href="#">GO:0016891</a> endoribonuclease activity, producing 5'-phosphomonoesters
-----------------	--

### InterPro annotation

	<a href="#">Entry Details in BioMart</a>
--	--

<b>Abstract</b>	<p>This domain is found in members of the Dicer protein family of dsRNA nucleases. This entry represents a dsRNA-binding domain. RNA interference (RNAi) is an ancient gene-silencing process that plays a fundamental role in diverse eukaryotic functions including viral defence, chromatin remodelling, genome rearrangement, developmental timing, brain morphogenesis, and stem cell maintenance. All RNAi pathways require the multidomain ribonuclease Dicer, which initiates RNAi by cleaving double-stranded RNA (dsRNA) substrates into small fragments ~25 nucleotides in length. A typical eukaryotic Dicer consists of a helicase domain (<a href="#">PDOC51192</a>), a domain of unknown function, and a PAZ domain (<a href="#">PDOC50821</a>) at the amino (N)-terminus as well as two ribonuclease III domains (<a href="#">PDOC00448</a>) and a dsRNA-binding domain (dsRBD) (<a href="#">PDOC50137</a>) at the carboxy (C)-terminus. The domain of unknown function of ~100 amino acids is predicted to adopt the canonical alpha-beta-beta-beta-alpha-fold found in all dsRBDs <a href="#">[1, 2, 3, 4]</a>.</p>
-----------------	---

<b>Database links</b>	Enzyme: <a href="#">EC:3.1.26</a> PANDIT: <a href="#">PF03368</a> Blocks: <a href="#">IPB005034</a>
-----------------------	---



# Interpro: la banque de domaines intégrée de l'EBI

## Example proteins

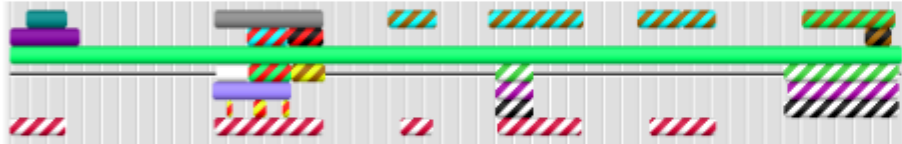
P34529 Endoribonuclease dcr-1



P84634 Dicer-like protein 4



Q12923 Tyrosine-protein phosphatase non-receptor type 13



Q8R418 Endoribonuclease Dicer



Q9VUQ5 Protein argonaute-2



SCOP: [b.34.14.1](#) 591 - 716

Chaque rectangle coloré: un domaine identifié par une des banques (par exemple: banque SCOP)