

Evolution et biodiversité des microorganismes : travaux dirigés

L'évolution des protéines universelles Sua5/TsaC : quand la simplicité est l'ultime sophistication ?

Objectifs : Savoir chercher les séquences à partir de la base de données NCBI. Savoir utiliser les logiciels d'alignement de séquences. Savoir construire les arbres phylogénétiques et les interpréter.

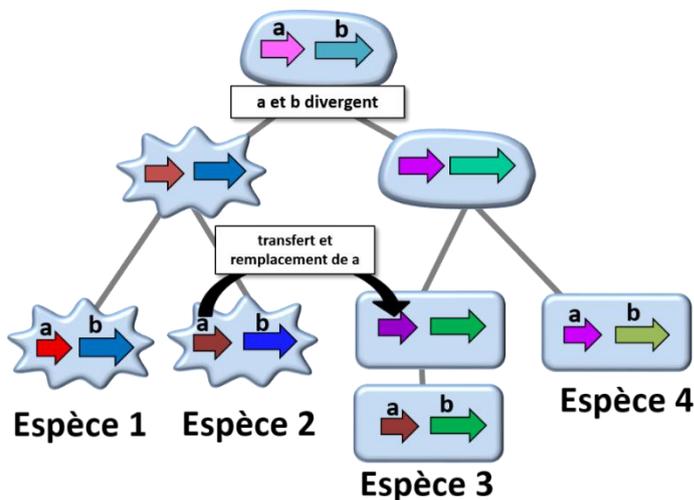
Consignes : Donnez vos réponses en couleur pour mieux les repérer, utilisez la police Courier New (taille de lettres constante) pour écrire les séquences.

PARTIE II

Le transfert horizontal de gènes (HGT) a été découvert grâce à l'étude des microorganismes à la fin des années 1940. Depuis, l'amélioration des méthodes de détection d'HGT a mis en évidence l'étendue surprenante de l'HGT dans l'évolution du contenu génomique chez les virus, les procaryotes et les eucaryotes (Soucy et al., 2015, Nature Reviews Genetics). Par exemple, de nombreuses duplications de gènes sont maintenant connus pour être le résultat de l'HGT, et non pas de la duplication des gènes autochtones.

Dans cette deuxième partie du TD nous allons tenter d'identifier l'HGT au sein de gènes codant pour la famille Sua5/TsaC. En effet, la mobilité des gènes entre espèces pourrait expliquer en partie la distribution étonnante de Sua5 et de TsaC au sein du vivant.

1. Les gènes « voyageurs » qui ont subi un HGT sont généralement éliminés lors des analyses phylogénétiques dont l'objectif est d'inférer la relation de parenté entre différents organismes. Expliquez pourquoi à l'aide du schéma ci-dessous.



Sur le schéma on observe un transfert du gène *a* de l'espèce 2 à l'espèce 3 suivi par la perte du gène *a* endogène. Cela implique que les gènes *a* des espèces 2 et 3 sont plus proches que les gènes *a* des espèces 3 et 4. En effet l'arbre fait avec les séquences des gènes *a* va regrouper les espèces 2 et 3 en groupe sœur ; cet arbre représente l'histoire évolutive du gène *a* mais ne représente pas l'histoire évolutive des espèces qui le portent.

2. La question précédente nous montre qu'un HGT peut être identifié si la topologie de l'arbre d'un gène ne correspond pas à celle des espèces porteuses de ce gène. Nous avons vu également qu'avec le temps le gène endogène peut être remplacé par le gène acquis par HGT. Cependant, quand le transfert est relativement récent les deux copies de gène (endogène et exogène) peuvent co-exister. L'analyse de distribution de *Sua5* et *TsaC* a ainsi mis en évidence la présence de gènes codant pour *Sua5* et pour *TsaC* chez les champignons du genre *Aspergillus*. Or, l'analyse de distribution de *tsaC* et *sua5* au sein du vivant montre que tous les autres champignons codent seulement pour la *Sua5*. Cela indique que le gène codant pour la *TsaC* a été introduit par HGT chez le genre *Aspergillus*. Votre mission est de tester cette hypothèse grâce à l'analyse phylogénétique.

Voici les séquences protéiques pour le *Sua5* (UniProt ID Q4X1Q3) et *TsaC* (Uniprot ID A4D9M1) d'*Aspergillus fumigatus* Z5.

```
>tr|Q4X1Q3|Q4X1Q3_ASPFU Threonylcarbamoyl-AMP synthase OS=Aspergillus
fumigatus (strain ATCC MYA-4609 / CBS 101355 / FGSC A1100 / Af293)
OX=330879 GN=AFUA_2G09160 PE=3 SV=1
MTPRETRILSVRRLNKEDPGTRSLAEWWASERSQKTPESSAIEEAAHLLRTSDIPVAFPT
ETVYGLGADATRSAAVQGIYKAKQRPSDNPLIVHIDSIEMLERLLNPASSSPTRTRTAK
NTIPAIYQPLIERFWPGPLTILLHNPSGSLLADEV TANLTTFGVRMPASPLARLLIHVAD
RPLAAPSANASTKPSPTAAEHVYHDLEGRINLILDGGPCGVGVESTVVDGLSKPPAILRP
GGVGIEELRTVPGWENVQIAYHDGNLDVKEVPRAPGMKYRHYS PKARVVLFCAGSSEEEI
AKYVYKDLEDTAIRAHMIGVVRTRQWKRG LGLVSEEDIQKTLKPIPSLVDDLVCFSVPVK
GRINNREIARAADFCHLGDNIESIARGLFSALRAMDEMEVDVIYVEGVSDSQDLAAAVM
NRLRKAAGTVLKL
```

```
>tr|A4D9M1|A4D9M1_ASPFU Threonylcarbamoyl-AMP synthase OS=Aspergillus
fumigatus (strain ATCC MYA-4609 / CBS 101355 / FGSC A1100 / Af293)
OX=330879 GN=AFUA_4G13765 PE=3 SV=1
MQTTIDIPTDAARVFSILAQGGIGIVPSSVGYGIIATEPPALQRIYTVKRRQPHKRHAI I
GSYALHREIHVLPDPKMLVRWLTVDLNLPLGVIARYRRDHPLLARLDEETRAASSMDGT
MAMLVNGGPFQEELVRVAAAAGRAVLGSSANLTGQGTKTVEAIEEEI REAADI VVDYGR
VRDGWPRASSTMVDFEAMRVVRVGACYEAIHDVVKRFAGLNWPDPSV
```

3. Vérifiez par alignement multiple de séquences que *AfSua5* et *AfTsaC* contiennent les motifs conservés ($P^{227}-G^{228}-M^{229}$ et $H^{234}-Y^{235}$ et $K^{56}-X-R^{58}/S^{143}-X-N^{145}$). Utilisez la séquence de *Sua5* de *P. abyssi* (ci-dessous) comme référence. Copiez-collez le résultat d'alignement et surlignez les motifs trouvés en jaune. Que pouvez-vous conclure ?

```
>WP_010868714.1 threonylcarbamoyl-AMP synthase [Pyrococcus abyssi]
MTIIINVRERIEEWKIRIAAGFIREGKLVAFPTETVYGLGANALDENAVKRIFEAKGRPADNPLIIHIAS
FEQLEVLAKEIPEEAEMLAKRFPWPGPLTLVLPKSEVVPRVITGGLDTVAVRMPAHEIALKLIELSERPIA
APSANISGKPSPTSAHHVAEDFYGKIECIIDGGETRIGVESTVIDLTEWPPVLLRPGGLPLEEIEKVI GE
IRIHPAVYGKSVDTAKAPGMKYRHYP SAEVIVVEGPRDKVRRKIEELIAKFKEEGKVKVIGSGSYDAD
EVFYLGDTVVEEIIARNL FKALRHMDRTGVDVILAEGVEEKGLGLAVMNRLRKASGYRI IKV
```


<input type="checkbox"/>	hypothetical protein [Streptomyces sp. PRh6]	223	223	96%	8e-70	50%	WP_037952691.1
<input type="checkbox"/>	hypothetical protein [Streptomyces canus]	220	220	96%	6e-69	50%	WP_020124679.1
<input type="checkbox"/>	hypothetical protein [Mycobacterium sp. Soil638]	215	215	96%	8e-67	50%	WP_05149493.1
<input type="checkbox"/>	hypothetical protein [Streptomyces canus]	214	214	96%	3e-66	49%	WP_020123940.1
<input type="checkbox"/>	hypothetical protein [Mycobacterium chubuense]	213	213	96%	7e-66	50%	WP_048420378.1
<input type="checkbox"/>	hypothetical protein [Mycobacterium iranicum]	212	212	96%	2e-65	49%	WP_085174403.1
<input type="checkbox"/>	hypothetical protein [Mycobacterium iranicum]	211	211	96%	4e-65	48%	WP_024446495.1
<input type="checkbox"/>	hypothetical protein [Mycobacterium smegmatis]	211	211	99%	5e-65	46%	WP_003891541.1
<input type="checkbox"/>	hypothetical protein [Candidimonas nitroreducens]	207	415	96%	1e-63	49%	WP_088605544.1
<input type="checkbox"/>	hypothetical protein [Advenella kashmirensis]	207	207	97%	1e-63	49%	WP_024004843.1

La recherche dans la NCBI Taxonomy nous apprend que les 8 premières espèces appartiennent à la classe Actinobacteria et que les 2 dernières espèces appartiennent à la classe Betaproteobacteria. Cela suggère qu'en dehors de champignons, AftSaC semble être le plus similaire à des orthologues bactériens. On peut donc émettre l'hypothèse qu'un gène ancestral tsaC a été transféré depuis les bactéries dans l'ancêtre commun de la lignée conduisant aux Aspergillus.

6. Pour commencer les analyses nous allons récupérer la séquence protéique qui correspond au meilleur résultat du BLASTp ainsi que les séquences de TsaC issues de bactéries *E. coli*, *Thermus aquaticus* et *Prochlorococcus marinus* :

```
>WP_037952691.1 hypothetical protein [Streptomyces sp. PRh5]
MGRHDINADAKRVFDTITAGGAVILPGDIGYGAGASSPEALQRLFVAKQGRAPHKRRHAMVGNIELHRELVH
LGSREQEIVDAITIDADLPLTVVAAAYRADHPVVAAVEPETLAASTVGNTLALLVNGGRLQDEVVRLCHQA
GLPFLGSSANLTGTGTGKFRVEDIQQPILDVAVLDVLDYGLRKYHHYRRSSTIIDFSTMEVVRIGTCYELIS
DIMATQFGITLTPADPGRDALPSGHLREQTQ
```

```
>tr|A0A163U4K7|A0A163U4K7_PROMR Threonylcarbamoyl-AMP synthase
OS=Prochlorococcus marinus str. MIT 1342 OX=1801627 GN=yw1C PE=3 SV=1
MAVLPSAVLDASALASRLHAGSAALLPTDPLALAAAPVHAAQLWTIKQRSADKPLILMG
ATPEELLAHVLPLEALDAWSMARRYWPGALTLVVPASGVLVEALNPGAFTLGLRVPDCGM
LRNLLKQSGPLATTSANLSGSAPTFSADAHTCFPGLPLGPLPWPTPSGLASTVVAWQG
AGRWHELRRGAVVLEL
```

```
>sp|P45748|TSAC_ECOLI Threonylcarbamoyl-AMP synthase OS=Escherichia coli
(strain K12) OX=83333 GN=tsaC PE=1 SV=2
MNNNLQRDAIAAAIDVLNEERVIAYPTEAVFGVGCDDPSETAVMRLLELQRPVVDKGLIL
IAANYEQLKPYIDDTMLTDVQRETIFSRWPGPVTFVFPAPATTPRWLTGRFDSLAVRVTD
HPLVVALCQAYGKPLVSTSANLSGLPPCRTVDEVRAQFGAAFPVVPGETGGRLNPFSEIRD
ALTGELFRQG
```

```
>tr|A0A0M9AGB3|A0A0M9AGB3_THEAQ Threonylcarbamoyl-AMP synthase OS=Thermus
aquaticus OX=271 GN=yw1C PE=3 SV=1
```

MTEVYQKELAQAAEVLKKGGLVAFPTDVTWVGLARMEDEAACRRIYALKGREEKKPLQVL
 VAGLEDALRLSDLGPLEERFLRLAEAFWPGALTIVVPPGRGIPPWISRDSVGLRMPAHEA
 LRELLRRVGGHAAATSLNRSGEPPVRTEAEARAFPVDFVFPGEATGLASSVVDLRTGEIL
 REGAIPKEALLPYLHG

Nous allons compléter ce jeu de séquences par celles issues d'eucaryotes. Nous allons utiliser les Tsac issues des Metazoa (animaux) car ce groupe est phylogénétiquement proche des champignons et il contient des organismes qui codent exclusivement pour le Tsac. Voici les séquences de Tsac d'*Homo sapiens*, *Mus musculus*, et *Drosophila melanogaster*.

```
>sp|Q86U90|YRDC_HUMAN YrdC domain-containing protein, mitochondrial OS=Homo sapiens OX=9606 GN=YRDC PE=1 SV=1
MSPARRCRGMRAAVAASVGLSEGPAGSRSGRLFRPPSPAPAAPGARLLRLPGSGAVQAAS
PERAGWTEALRAAVAELRAGAVVAVPTDTLYGLACAASCSAALRAVYRLKGRSEAKPLAV
CLGRVADVRYRCVRVPEGLLKDLLPGPVTLVMERSEELNKDLNPFPTLVGIRIPDHAFM
QDLAQMFEGPLALTSANLSSQASSLNVEEFQDLWPHLSLVIDGGQIGDQGSPECRLGSTV
VDLSVPGKFGIIRPGCALESTTAILQOKYGLLPSHASYL
```

```
>AAP37032.1 ischemia/reperfusion inducible protein [Mus musculus]
MGLSDGPASSGRGCRLLLPPEPAPALPGARLLRLPESEPVAAASPERAGWTEALRAAVAELRAGAVVAVP
TDTLYGLACSASSAALSCVYRLKGRSEAKPLAVCLGRVADVRYRCQVRVPRELLEDLFFPGPVTLVMERS
EELNKDLNPFTRLVIRIPDHAFMLDLAQMFEGPLALTSANLSSQASSLSVEEFQDLWPHLSLVIDGGPI
GDSQSPECRLGSTVVVDLSVPGKFGIIRPGCALENTTSTILQOKYGLLPSQSCS
```

```
>tr|Q8SYJ9|Q8SYJ9_DROME RE55868p OS=Drosophila melanogaster OX=7227
GN=CG10438 PE=2 SV=1
MRRQLTSLYRLLRAHHTSSRMQHQASELRTPVCAVGDEAALQLARQCLLGGQVIALPTDPT
VYGLACDANNETAIQQLYEIKGRDEHKPVAICVHNIDALRRFGQAAHLSDELLTRLLPGP
LTIVIERSTQLSNRFLNPSTSKIGIRIPDFNFMRDLCVWQEKPLALTSANRSSAPSSLQ
VSEFRSLWPLGAVFDAGRIGLTEERRLASTVIDLATPGYIEIVRAGVALKPTLSLMEEF
GIRELKMM
```

8. Aalignez maintenant toutes les séquences y compris AftTsac (donc huit séquences au total), puis construisez un arbre avec IQ-Tree. Comme précédemment, laissez le programme choisir le meilleur modèle d'évolution de séquences et demandez 1000 répliques d'analyse ultrabootstrap. Copiez-collez l'alignement et l'arbre et commentez-les.

```
tr|A0A0J5PIE0|A0A0J5PIE0_ASPFM ----- 0
WP_037952691.1 ----- 0
sp|P45748|TSAC_ECOLI ----- 0
tr|A0A0M9AGB3|A0A0M9AGB3_THEAQ ----- 0
tr|A0A163U4K7|A0A163U4K7_PROMR -----MA 2
tr|Q8SYJ9|Q8SYJ9_DROME -----MRRQLTSS-----LYRLLRAHHTSSRMQHQASELRTP----- 31
sp|Q86U90|YRDC_HUMAN MSPARRCRGMRAAVAASVGLSEGPAGSR-SGRLFRPPSPAPA-APGARLLRLPGSGAVQA 58
AAP37032.1 -----MGLSDGPASSGRGCRLLLPPEPAPA-LPGARLLRLPESEPVEA 42
```

```
tr|A0A0J5PIE0|A0A0J5PIE0_ASPFM -MQTTIDIPTDAARVFSILAQGGIGIVPSSVGYGIIA--TEPPALQRIYTVKRRQPHKRH 57
WP_037952691.1 --MGRHDINADAKRVFDITAGGAVILPGDIGYGAGA--SSPEALQRLVFAKQRAPHKRH 56
sp|P45748|TSAC_ECOLI --MNNNLQRDAIAAAIDVLEERVIAYPTEAVFVGVGCDPDESETAVMRLLELKRQRPVVKGL 58
tr|A0A0M9AGB3|A0A0M9AGB3_THEAQ ---MTEVYQKELAQAAEVLKKGGLVAFPTDVTWVGLARMEDEAACRRIYALKGREEKKPL 57
tr|A0A163U4K7|A0A163U4K7_PROMR VLPSSAVLDA---SALASRLHAGSAALLPTDTLPALAAAPV---HAAQLWTIKQRSADKPL 56
tr|Q8SYJ9|Q8SYJ9_DROME --VCAVGDEAALQLARQCLLGGQVIALPTDPTVYGLACDANNETAIQQLYEIKGRDEHKPV 89
sp|Q86U90|YRDC_HUMAN ASPERAGWTEALRAAVAELRAGAVVAVPTDTLYGLACAASCSAALRAVYRLKGRSEAKPL 118
AAP37032.1 ASPERAGWTEALRAAVAELRAGAVVAVPTDTLYGLACSASSAALSCVYRLKGRSEAKPL 102
: * . . . : * * *
```

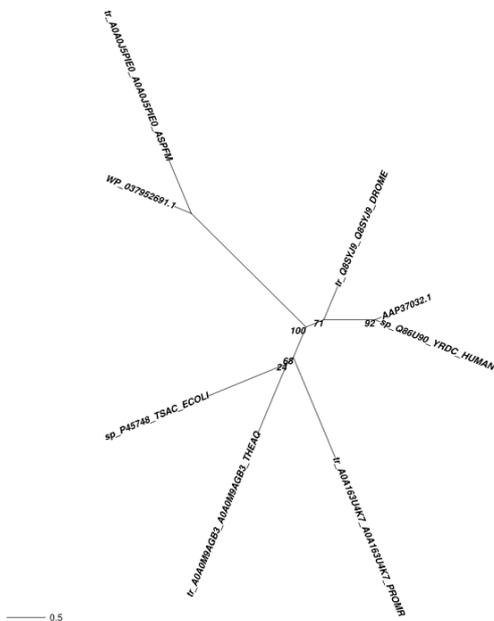
```
tr|A0A0J5PIE0|A0A0J5PIE0_ASPFM AIIIGS-YALH-REIHVLPDPKMLVRLWLVLDLNL--PLGVIARYRRDHPLLARLDEETRA 113
WP_037952691.1 AMVGN-YELH-REVHVLSGREQEIVDAITIDADL--PLTVVAAVYRADHPVVAVEPETLA 112
sp|P45748|TSAC_ECOLI ILIAANYEQLKPYIDDTML---TDVQRETIFSRWPGPVTFVFPAPA--TT-----PRW 106
```

Université Paris-Sud, Faculté des sciences, M1 Biologie-Santé, Evolution et Biodiversité des Microorganismes, 2024-2025

tr A0A0M9AGB3 A0A0M9AGB3_THEAQ	QVLVAGLEDALRLSDLGPL---EERFLRLAEAFWPGALTIIVVPRG--I-----PPW	104
tr A0A163U4K7 A0A163U4K7_PROMR	ILMGATPEELLAHVLPAL----EDAWSMARRYWPGALTLVVPASG--VL-----VEA	103
tr Q8SYJ9 Q8SYJ9_DROME	AICVHNIDALRRFGQAA-----HLSDELLTRLLPGPLTIIVERTS--QL-----S--NRF	135
sp Q86U90 YRDC_HUMAN	AVCLGRVADVRYCRV-----RVPEGLLKDLLPGPVTLVMERSE--EL-----NKD	162
AAP37032.1	AVCLGRVADVRYCQV-----RVPRELLEDLFPGPVTLVMERSE--EL-----NKD	146
	: : :	
tr A0A0J5PIE0 A0A0J5PIE0_ASPFM	ASSMDGTMAMLVNGGPFQEELVRAAAA-GRAVLGS SANLTGQGTKTVEAIEEEEIREAA	172
WP_037952691.1	ASTVGNLALLVNGGRLQDEVVRLCHQA-GLPFLGSSANLTGTGTFKFRVEDIQPILDAV	171
sp P45748 TSAC_ECOLI	LTGRFDSLAVRVTDH---PLVVALCQAY-GKPLVSTSANLSGLPPCRTVDEVRA-QFGAA	161
tr A0A0M9AGB3 A0A0M9AGB3_THEAQ	IS-RDGSVGLRMPAH---EALRELLRRV-GGHAATSLNRSGEPPVTEAEARA-F-FVD	157
tr A0A163U4K7 A0A163U4K7_PROMR	LNPGAFTLGLRVPCD---GMLRNLLK--QSGPLATTSANLSGSAPTFSADAHT-CFPG	157
tr Q8SYJ9 Q8SYJ9_DROME	LNPSTSKIGIRIPDF---NFMRDLCVWQEKPLALTSANRSSAPSSLQVSEFRS-LWPQL	191
sp Q86U90 YRDC_HUMAN	LNPFTPLVGIRIPDH---AFMQDLAQM-FEGPLALTSANLSSQASSLNVEEFQD-LWPQL	217
AAP37032.1	LNPFTRLVGIRIPDH---AFMFLDLAQM-FGGPLALTSANLSSQASSLSVEEFQD-LWPHL	201
	. : : : : : : * * .	
tr A0A0J5PIE0 A0A0J5PIE0_ASPFM	DIVVDYGRVRDVGW---PRASSTMVDFEAM---RVVRVGCYEAIHVDVVKRFAGLNWPDPS	226
WP_037952691.1	DLVIDYGLRKYHH---YRRSSTIIDFSTM---EVVRIGTCYELISDIMATQFGITLPADP	225
sp P45748 TSAC_ECOLI	FPVVPGETGGRNLN-----PSEIRDALTG---ELFRQG-----	190
tr A0A0M9AGB3 A0A0M9AGB3_THEAQ	F-VFPGEATGL-----ASSVVDLRTG---EILREGAIPKEALLPY--LHG-----	196
tr A0A163U4K7 A0A163U4K7_PROMR	PLLGP---LP--WFTPSGLASTVVAWQAGRWHELRRGAVVLEL-----	196
tr Q8SYJ9 Q8SYJ9_DROME	GAVFDAGRIGLT---EERRLASTVIDLATPGYYEIVRAGVALKPTLSLME-EFGIRELKM	248
sp Q86U90 YRDC_HUMAN	SLVIDGGQIGDQSPCECLGSTVVDLSVPGKFGIIRPGCALESTAILQQKYGLLPSHAS	277
AAP37032.1	SLVIDGGPIGDSQSPCECLGSTVVDLSVPGKFGIIRPGCALENTSILQQKYGLLPSQGS	261
	: * : * *	
tr A0A0J5PIE0 A0A0J5PIE0_ASPFM	V-----	227
WP_037952691.1	GRDALPSGHLREQTQ	240
sp P45748 TSAC_ECOLI	-----	190
tr A0A0M9AGB3 A0A0M9AGB3_THEAQ	-----	196
tr A0A163U4K7 A0A163U4K7_PROMR	-----	196
tr Q8SYJ9 Q8SYJ9_DROME	-----	248
sp Q86U90 YRDC_HUMAN	YL-----	279
AAP37032.1	CS-----	263

Sur l'alignement on peut identifier le motif $K^{56}\text{-X-R}^{58}/S^{143}\text{-X-N}^{145}$ (souligné en jaune) impliqué dans la fixation de l'ATP. Il s'agit donc probablement des TsAc fonctionnels.

Résultat d'analyse IQ-Tree



Modèle d'évolution de séquences utilisé : WAG+F+I+G4

(tr_A0A0J5PIE0_A0A0J5PIE0_ASPFM:0.7517048483,WP_037952691.1:0.2308037019, ((sp_P45748_TSAC_ECOLI:1.0986000172,tr_A0A0M9AGB3_A0A0M9AGB3_THEAQ:0.9665203323)24:0.1156960630,tr_A0A163U4K7_A0A163U4K7_PROMR:1.4094376756)68:0.4281641237,(tr_Q8SYJ9_Q8SYJ9_DROME:0.4666681366,(sp_Q86U90_YRDC_HUMAN:0.070739394

3, AAP37032.1:0.0901858430) 92:0.6619580960) 71:0.2487546783) 100:2.0983180858) ;

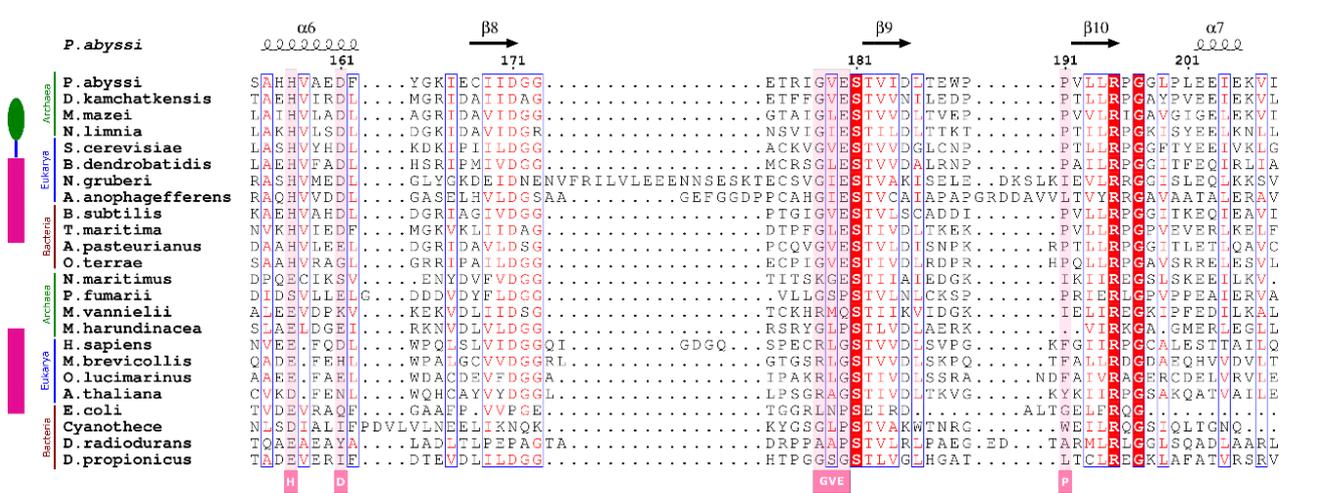
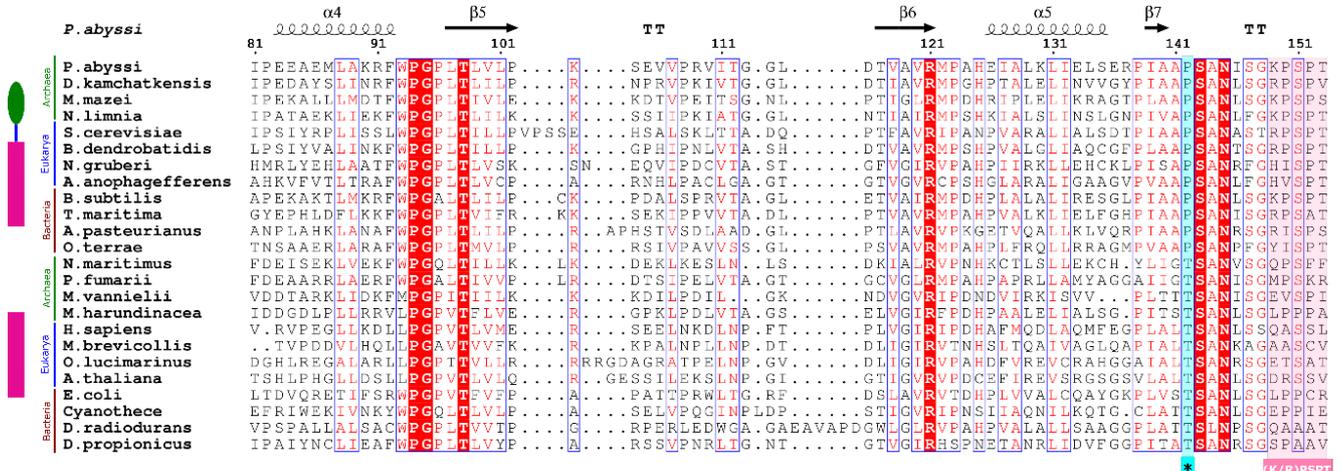
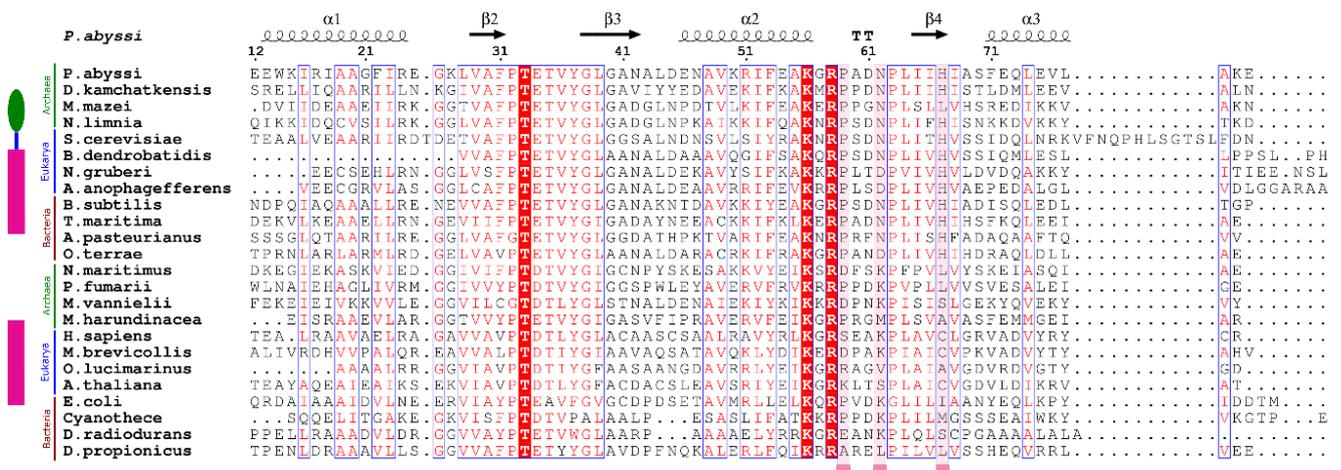
9. Editez votre arbre dans iTOL pour faire apparaître les noms d'espèces. Quelle topologie votre arbre devrait-il avoir si AfTsaC était effectivement transféré aux *Aspergillus* à partir de Bactéries? Votre arbre correspond-il à cette prédiction ? Enracinez votre arbre, interprétez-le et concluez.



Si l'hypothèse de HGT est correcte on devrait trouver un arbre composé de deux clades : l'un comportant toutes les séquences d'eucaryotes (hormis celle d'*Aspergillus*) et l'autre comportant toutes les séquences de bactéries y compris celle d'*Aspergillus*. L'arbre ci-dessus correspond à cette topologie, cependant le clade d'eucaryotes et celui de bactéries sont faiblement soutenus (74 et 63% respectivement), de plus il n'y a aucun support pour la branche reliant le clade *Aspergillus*/*Streptomyces* avec le clade des bactéries. Le positionnement du clade *Aspergillus*/*Streptomyces* et donc incertain (on ne sait pas s'il est plus proche des bactéries ou des eucaryotes). L'analyse confirme bien cependant que les TsaC issues d'*Aspergillus* et *Streptomyces* sont les séquences les plus proches parmi celles que nous avons analysées. En conclusion, l'arbre semble bien indiquer un transfert horizontal de TsaC depuis les bactéries vers *Aspergillus* cependant le signal phylogénétique n'est pas assez robuste. Il faudrait faire d'autres arbres avec plus de séquences pour chaque groupe pour essayer d'obtenir un meilleur soutien pour les branches.

PARTIE III

Dans la première partie de ces travaux dirigés vous avez constaté que les séquences de Sua5 et TsaC se ségrégent en deux groupes vraisemblablement parce qu'il existe des résidus conservés spécifiques d'un ou de l'autre variant. En effet, l'analyse des séquences de TsaC et des domaines TsaC-like par alignement sur un échantillon taxonomique large montre l'existence de ces résidus (figure ci-dessous). Ainsi, on trouve 14 résidus conservés (indiqués en rectangles roses sous l'alignement) spécifiques des protéines Sua5. Un seul résidu est **conservé et spécifique des protéines TsaC**, il s'agit de la thréonine (Thr) adjacente au motif S¹⁴³-X-N¹⁴⁵ (souligné en bleu). De manière intéressante, à la même position dans l'alignement on trouve un résidu Proline (Pro) conservé chez les protéines Sua5. Cela suggère que la présence du domaine supplémentaire chez les protéines Sua5 nécessite la présence de Pro à côté du motif S¹⁴³-X-N¹⁴⁵. L'inspection d'un grand nombre de séquences de la base de données UniProt confirme cette règle (données non-montrées). Mais, comme souvent en biologie, il y a des exceptions. Dans cette partie du TD, nous allons étudier l'une de ces exceptions : les Sua5 issues d'archées hyperthermophiles, les Archaeoglobales.



1. Pour commencer vous avez fait un BLASTp en utilisant Sua5 de *P. abyssi* et en restreignant la recherche uniquement aux Archaeoglobi. Vous avez récupéré les séquences de toutes les Sua5 d'espèces cultivées (*A. sulfaticallidus*, *G. acetivorans*, *G. ahangari*, *F. placidus*, *A. veneficus*, *A. profundus*, *A. fulgidus*). Vous allez inclure également la Sua5 de *P. abyssi* pour servir d'outgroup.

```
>WP_015590325.1 threonylcarbamoyl-AMP synthase [Archaeoglobus sulfaticallidus]
MKKTRIKVNPENPEEDRIMAGAEIKNGGLVAFPTETVYGLGANALDRKAVLGIYKAKERPADNPLIVHVCNKDMIEDI
AELNDVAEKLIKRFPPGPLTVVVKKKECIPQEVTAGLDTVAIRMPDNKVALKLIELSEVPI SAPSANLAGKPSPTKAEHV
IEDLYGRIDAILDGGPTNIGLESTVVDTTVYPVELLRPGGLSVEELEKVVDIKIPPELDEGIPRSPGMKYRHYAPSAELIV
IYGKREDVVSKI IKTAGQLIDKSNKKIGLVVSDSEPKDIDVEIVEIGRSVEEVARNLFSALRELDRRGIDLIIAEGVEE
KGLGLAVMNRLLKASNYRIFRV

>WP_048093555.1 threonylcarbamoyl-AMP synthase [Geoglobus acetivorans]
MIRISSEKFSEEELSPAELIKAGNLVAFPTETVYGLGANALEGKAVRKIFAAGKRPDNLIVHVHEHEQIYELAKPNR
VAEKLVEEFFPGPLTLVMRKKEVVPGETTGGLNTVAVRMPAHKVALKLIELSGVPIAAPSANKSGKPSPTRAEHVIEDFG
NEIDCII DAGKTRIGLESTVVDTTTYPLEILRPGAITREMLEEFFEVKVIKSDVARSPGMKYRHYSPADVIVLVGDTA
RDEMRELAERLSRDGRKVGIAAMKADEFEGIATYNLSTLREFAERLFDALRELDRCVDVIIVEGVEEKGLGLAIMNRLS
KAGRVYRV

>WP_048095299.1 threonylcarbamoyl-AMP synthase [Geoglobus ahangari]
MRVIRVAPERFRDEELEPAADIIRQGLVAFPTETVYGLGADALNERVRRIFEAKGRPSDNPLIVHVSSIEEVYRIAKP
NRVAEKLMEEFFPGPLTLVMEKREAVPEVTTGGLRVTAVRMPKHRVALKLIELSGTPIAAPSANKSGKPSPTRAEHVME
FESIECVIDAGRTEVGLLESTVVDTTVYPIEILRPGAITREMLEELFEVVRVAGKSDVARSPGMKYRHYSPADTIVIVGE
KRREEIKRLAHLRAGEGRVTGVAAAMNGEEFESVAPHVYDLGASLEEFERLFDALRTLDRFCDVIVVEGVQKIGLAIM
NRLSKAGRVYRV

>WP_012965368.1 threonylcarbamoyl-AMP synthase [Ferroglobus placidus]
MKIRTRIKVDPENPEDDKIGIAAEEIKKGNLVAFPPTETVYGLGADALNEKAVRKIFEVKGKRPADNPLIVHVSSVDQIYE
IAQNEIAKKLIDAFFPGPLTLVMKKNKVPKITTAGLDTVAVRMPDVKVALKLIELSETPIAAPSANKSGKPSPTKAEHV
LEDLGNLVDVIIDGGETKIGLESTVVDTTVYVEILRPGAITKKEELEYVDVKYAEDEFSVAKSPGMKYKHYAPDAETIVL
VGKNFLKATQIAEDFVKKGYRVGVAGLNLEEKRNVDVIYKNFGRNLEEMAKNLFKVLRELDKECDVIVQGVEEKGIGKA
IMNRLYKAGRVFRI

>WP_013683929.1 threonylcarbamoyl-AMP synthase [Archaeoglobus veneficus]
MIVIRVDPVNPERRGKIARAEEIKKGLVAFPTETVYGLGGDALNENAVRRIFEAKERPPRNPLIVHVSSVEQVYRIAEV
NEVAEKLMEEFFPGPLALVLKDKDVPDITTAGMKKAVRMPDVKVALTLIELSETPIAAPSANISGRPSPTKAEHVIED
LAGRIDAILDAGEVKIGIESTVLDVTSKPAKILRLGAITPEMLVERGIEVEVIDTRPFRHYQTKAKLYVNAENLAEFVG
SLREKGVKVGVARITAECADRIVELGKSIEDVAKNLFSAALRELDRCVVDIIVVEAVERKGLGKVMISKLEEEAGEIV

>WP_010878284.1 threonylcarbamoyl-AMP synthase [Archaeoglobus fulgidus]
MAEVLEPSREGIRKAVAVLRAGGIVAFPTETVYGMGCATNEEALRRLYEIKGRSLNKPFIIVGVWSDRYVKGIAEVDERA
EKLMFAFFPGPLTLVLKSKGVMPSSLSPKGIIVRMPAHEVPLQLMMLRKP IIVVPSANLSGRPSLMRWEHVVEELGSR
DAVVKGECKVGVSTIVDLTETPAKVLRVGAVSVESIKKHVEVVVEPRKETYSLSQVYVFGERALQRIKEFVDEAEKR
GRVVVIAREKITDETIVIGKSAEYSANLFSAVREAESRRPDIIVIEGVENEAIMDRLRRLAGERVFRV

>WP_012939514.1 threonylcarbamoyl-AMP synthase [Archaeoglobus profundus]
MTKIYKVDPRNPDESVLREVAHMIKEGKLVAYPTETVYGLGTNALDENAVKRLFVKGKRPKPKVSVVSDLDHVERIAEP
NETALKLMEKFFPGPITIIIVKVKVIPPVITAGTDKIGIRLPDYKIPKLAEFSGVPTSTSANVSGKPSPTKPEHIMVD
FMGKIDAILDAGETPLKIESTVIDTTEPPRVLRVGLPLNEIEKVVGEVEVLDRAYKPKAKVAVFVKGDVVEEVASRFE
GRVCIIEKVKDLMDVLRCKDCDVIVFECIDLSEAVRARIMSIADVIYD

>WP_010868714.1 threonylcarbamoyl-AMP synthase [Pyrococcus abyssi]
MTIINVRERIEEWKIRIAAGFIREGKLVAFPTETVYGLGANALDENAVKRIFEAKGRPADNPLIIHIAS
FEQLEVLAKEIPEEAEMAKRFWPGPLTLVLPKSEVVRVITGGLDTVAVRMPAHEIALKLIELSERPIA
APSANISGKPSPTSAHHVAEDFYGKIECIIDGGETRIGVESTVIDLIEWPPVLLRPGGLPLEEIEKVI
IRIHPAVYKGSVDTAKAPGMKYRHYAPSAEVIIVVEGPRDKVRRKIEELIAKFKEEGKVGVIIGSGSYDAD
EVFYLGDVVEEIANRNLFKALRHMDRTGVDVILAEGVEEKGLGLAVMNRLLKASGYRIIKV
```

2. Faites un alignement de ces séquences avec Clustal Omega et copiez-collez l'alignement ci-dessous. Identifiez et soulignez les motifs conservés (P^{227} -G²²⁸-M²²⁹ et H²³⁴-Y²³⁵ et K⁵⁶-X-R⁵⁸/S¹⁴³-X-N¹⁴⁵). Que remarquez-vous ?

3. Sachant que les protéines Sua5/TsaC sont essentielles pour la survie des cellules et que chez ces Archaeoglobi chaque espèce code pour un gène *sua5*, que suggère votre analyse de conservation de motifs importants pour la fonction des protéines Sua5 ?

Elle suggère que, chez 3 espèces *A. profundus*, *A. fulgidus* et *A. veneficus*, le linker (qui relie le domaine catalytique au domaine SUA5) a perdu sa fonction mais que la protéine est fonctionnelle tout de même. Autrement dit ces trois protéines n'ont plus besoin de linker pour fonctionner alors que le linker est essentiel pour les Sua5 classiques.

4. Regardez de nouveau l'alignement. Comparez la taille du linker et domaine SUA5 entre les Sua5 des trois espèces que vous avez identifiées comme porteuses de mutations dans le linker et la Sua5 d'*A. sulfaticallidus*. Pour information, chez *P. abyssi* le linker débute au résidu His²¹² de Sua5. Il comporte environ 20 résidus. Que remarquez-vous ?

On remarque que le linker et de domaine SUA5 chez *A. profundus*, *A. fulgidus* et *A. veneficus* comportent des délétions et que le nombre de résidus perdus est le plus grand chez *A. profundus*.

5. Enfin, trouvez dans l'alignement le résidu signature (adjacent au motif S¹⁴³-X-N¹⁴⁵) et soulignez-le en rouge ; ce résidu est toujours une Pro (P) pour les protéines Sua5 et toujours une Thr (T) pour les protéines TsaC. Que remarquez-vous ? À la vue de l'ensemble de vos observations quelle hypothèse pouvez-vous émettre ?

On remarque que le résidu signature est une Pro comme attendu pour six Sua5, une exception remarquable est le Sua5 de *A. profundus* qui code pour une Thr au lieu d'une Pro. De plus, c'est chez le même organisme que le linker et le domaine SUA5 portent le plus de délétions. A partir de ces données on peut proposer que le Sua5 d'*A. profundus* est une forme intermédiaire entre la variante longue Sua5 et la variante courte TsaC. Il semblerait que dans ce groupe d'organismes les Sua5 sont en train de perdre leur domaine SUA5 pour devenir la forme courte TsaC ; la Sua5 de *A. profundus* est le cas le plus avancé de ce processus.

Nous allons maintenant tenter de reconstituer le processus évolutif conduisant aux Sua5 plus courtes en utilisant l'analyse cladistique. Pour cela vous allez identifier parmi les résidus spécifiques et conservés des Sua5 ceux qui sont des synapomorphies.

6. Question préliminaire : dans l'analyse cladistique qu'est qu'un caractère plésiomorphe ? Apomorphe ? Synapomorphe ?

Selon l'approche cladistique de classification un caractère peut présenter soit un état ancestral (plésiomorphe) ou alors un état dérivé, évolué (apomorphe). Quand l'état apomorphe est partagé par deux ou plusieurs espèces cela s'appelle une synapomorphie.

7. Parmi ces trois catégories de caractères lequel peut être utilisé pour inférer la phylogénie ?

Les synapomorphies : Le partage par au moins deux taxons certifie que la dérive du caractère n'est pas simplement contingente, mais est apparue chez l'ancêtre commun des espèces qui partagent la synapomorphie.

8. Le tableau ci-dessous répertorie quelques résidus spécifiques (pas tous pour simplifier l'exercice) et conservés pour les protéines Sua5. Complétez-le en vous servant de votre alignement des Sua5 d'Archaeoglobi. Quand le résidu est présent mettez un X dans la case, quand il est absent coloriez la case en bleu, quand le résidu a été substitué coloriez la case en rouge. La numérotation des résidus correspond à celle de la Sua5 de *P. abyssi*.

	P59	N62	H67	K/R149	P152	T153	H157	D161	G178	P228	G229	M230	H234	Y235
<i>P.abyssi</i>	X	X	X	X	X	X	X	X	X	X	X	X	X	X
<i>A. sulfaticallidus</i>	X	X	X	X	X	X	X	X	X	X	X	X	X	X
<i>A. veneficus</i>	X	X	X	X	X	X	X	X	X				X	X
<i>A. fulgidus</i>														
<i>A. profundus</i>														

Correction:

	P59	N62	H67	K/R149	P152	T153	H157	D161	G178	P228	G229	M230	H234	Y235
<i>P.abyssi</i>	X	X	X	X	X	X	X	X	X	X	X	X	X	X
<i>A. sulfaticallidus</i>	X	X	X	X	X	X	X	X	X	X	X	X	X	X
<i>A. veneficus</i>	X	X	X	X	X	X	X	X	X				X	X
<i>A. fulgidus</i>				X			X	E	X					X
<i>A. profundus</i>	X			X	X	X	X	X						X

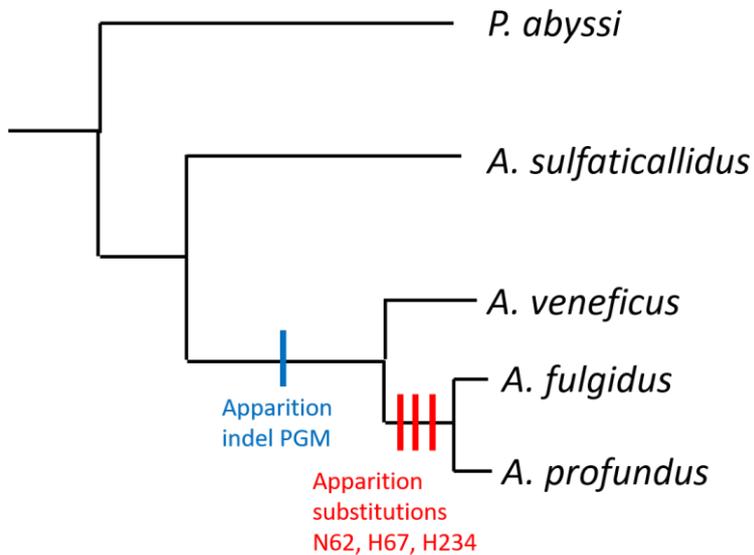
9. Faites maintenant une liste de caractères synapomorphes et des espèces chez lesquelles on les trouve. Pour cela complétez le texte ci-dessous :

Délétion PGM est présente chez : *A. veneficus*, *A. fulgidus* et *A. profundus*

Substitution N62 et H67 est présente chez : *A. fulgidus* et *A. profundus*

Substitution H234 est présente chez : *A. fulgidus* et *A. profundus*

10. A partir de cette liste dessinez maintenant un cladogramme décrivant l'évolution des protéines Sua5 issues des 5 espèces étudiées. Placez sur les branches l'apparition des synapomorphies chez l'ancêtre d'un clade. Vous pouvez faire votre dessin avec le logiciel Power Point puis le copier-coller ci-dessous.



11. A partir de l'analyse de l'ensemble des données au cours de ce TD proposez un scénario pour l'histoire évolutive de la famille Sua5/TsaC permettant d'expliquer l'existence de deux variantes et leur distribution chaotique (qui ne suit pas l'arbre des espèces) au sein du vivant.

Les données suggèrent que la forme ancestrale est Sua5 et qu'elle était donc présente chez LUCA, l'ancêtre des toutes les cellules modernes. La forme courte TsaC a pu apparaître au cours de l'évolution par perte du domaine SUA5 au moins une fois. Ce phénomène semble s'opérer actuellement chez certaines espèces d'Archaeoglobi. Les gènes *tsaC* et *sua5* semblent pouvoir se propager au sein du vivant par HGT (exemple d'*Aspergillus*). Les deux phénomènes (perte du domaine SUA5 et HGT) combinés peuvent expliquer la distribution actuelle de Sua5 et TsaC au sein du vivant.

12. Ces TDs sont intitulés « L'évolution des protéines universelles Sua5/TsaC : quand la simplicité est l'ultime sophistication ? » Pouvez-vous maintenant expliquer pourquoi les termes simplicité et sophistication ont été utilisés ?

Simplification : car les données suggèrent que les protéines Sua5 peuvent perdre un domaine, devenir donc des polypeptides plus simples, contenant un seul domaine au lieu de deux.

Sophistication : on peut spéculer sur les « raisons » pour lesquelles l'évolution conduit vers une perte du domaine SUA5 ; peut-être que les TsaC sont des versions « améliorées » plus performantes ? Cela reste à vérifier par des expériences biochimiques, structurales et *in vivo*.