

Enquête bioinformatique sur les origines de SARS-CoV-2

CM1 – Analyse des génomes de coronavirus

*Cours donné par **Jacques van Helden**, **Emese Meglécz** et **Gabriel Neve***

*Sur base d'une enquête menée par **Erwan Salard**, **José Haloy**, **Didier Casane**,
Etienne Decroly et **Jacques van Helden***

03/11	10:00	12:00	CM1	<ul style="list-style-type: none"> (1) La question des origines de SARS-CoV-2. (2) Biologie de SARS-CoV-2. (3) Événements évolutifs. (4) Bases de données biologiques. (5) Alignement de séquences par paires. (6) Recherche de séquences par similarité.
10/11	8:30	11:30	TP1 groupe 1	<ul style="list-style-type: none"> (1) Bases de données de séquences biologiques (Uniprot, NCBI). (2) Alignement par paires (needle). (3) Recherche de similarités (BLAST).
	12:30	15:30	TP1 groupe 2	cf groupe 1
	15:45	18:00	CM2	<ul style="list-style-type: none"> (1) Interprétation des résultats du TP 1. (2) Profils de pourcents de positions identiques. (3) Alignements multiples. (4) Inférence phylogénétique
17/11	8:30	11:30	TP2 groupe 1	<ul style="list-style-type: none"> (1) Alignements 1 à N et profils de PPI (PIPprofiler) (2) Alignements multiples (clustal). (3) Inférence phylogénétique (phylogeny.fr)
	12:30	15:30	TP2 groupe 2	cf groupe 1
	16:00	17:00	CM3	Interprétation, résumé et conclusion du cours

Pourquoi l'origine de SARS-CoV-2 suscite-t-elle des débats ?

Pourquoi la question des origines se pose-t-elle ?

- Dès le début de la pandémie, la question des origines du virus a suscité de fortes controverses : origine naturelle (transfert d'animal à humain) ou artificielle (produit dans un laboratoire) ?
- Les éléments de la controverse
 - ❑ Le virus a émergé dans la ville de Wuhan, où est situé le laboratoire de référence pour les recherches sur les coronavirus
 - ❑ Ce laboratoire réalise des expériences sur les coronavirus, notamment de gain de fonction
 - ❑ L'hypothèse d'une contamination dans le marché de Wuhan est fortement remise en cause
 - ❑ Certains chercheurs affirment que le génome de SARS-CoV-2 résulte de manipulations génétiques
 - ❑ Certains chercheurs pensent que le virus pourrait s'être échappé accidentellement d'un laboratoire. D'autres affirment carrément qu'il a été conçu pour servir d'arme biologique
 - ❑ Ces affirmations sont fortement médiatisées par les réseaux sociaux, mais ne sont pas publiées dans les revues scientifiques avec comité de lecture
 - ❑ Le débat scientifique est fortement perturbé par le contexte politique et géopolitique (relations Chine - USA, rôle de l'OMS, ...)
- Comment distinguer le vrai du faux ?
 - ❑ **Approche: ignorer délibérément les débats médiatiques et analyser les séquences des virus**

A la recherche de l'hôte manquant

- Les coronavirus responsables des dernières émergences épidémiques chez l'humain trouvaient leur origine dans un "réservoir naturel" constitué par les chauves-souris.
- La transmission à l'homme passe généralement par un hôte animal intermédiaire: la civette pour le SRAS (2002) et le dromadaire pour le MERS (2012).
- Pour la pandémie COVID-19 on a invoqué le pangolin comme hôte intermédiaire, en suggérant que le passage à l'homme provenait d'animaux vendus sur le marché de Wuhan.
- Cette hypothèse est cependant remise en cause pour différentes raisons.
 - Les premiers patients ne fréquentaient pas le marché
 - Les génomes des coronavirus de pangolin dont on dispose sont beaucoup plus éloignés de SARS-CoV-2 que ceux des chauves-souris

Expériences de gain de fonction

- Les coronavirus responsables des dernières émergences épidémiques chez l'humain trouvaient leur origine dans un "réservoir naturel" constitué par les chauves-souris.
- La transmission à l'homme passe généralement par un hôte animal intermédiaire: la civette pour le SRAS (2002) et le dromadaire pour le MERS (2012).
- Pour la pandémie COVID-19 on a invoqué le pangolin comme hôte intermédiaire, en suggérant que le passage à l'homme provenait d'animaux vendus sur le marché de Wuhan.
- Certains laboratoires de virologie pratiquent des expériences dites "de gain de fonction" qui consistent à modifier un virus pour le rendre plus virulent ou plus contagieux, afin d'étudier les mécanismes moléculaires de l'infection virale.
- Ceci suscite de fortes réticences au sein même de leur communauté.
- 2011: Ron Fouchier annonce qu'il a produit un virus H5N1 modifié (9 mutations ponctuelles) pour augmenter sa contagiosité chez le furet.
- Moratoire:
 - En 2014 le NIH annonce qu'il ne financera plus les expériences de gain de fonction
 - En 2017 il lève cette mesure, en invoquant l'intérêt de ces expériences pour comprendre les mécanismes de l'infection virale.
- Plusieurs incidents ont déjà été rapportés concernant des fuites accidentelles de virus de laboratoires.

- Andersen et collègues publient dès janvier 2020 un article affirmant que SARS-CoV-2 est sans aucun doute d'origine naturelle.
- Argument principal : la séquence du domaine de liaison au récepteur (RBD) est optimale pour se lier au récepteur ACE2, mais d'une façon différente de celles qu'on connaissait jusqu'alors. D'après les auteurs, si on avait conçu un virus dans le but de le rendre infectieux pour l'homme, on n'aurait pas pu concevoir cette séquence.
- Autre argument: on ne trouve pas dans ce génome de traces d'ingénierie moléculaire (par exemple des sites de restriction)

correspondence



The proximal origin of SARS-CoV-2

NATURE MEDICINE | VOL 26 | APRIL 2020 | 450-455 | www.nature.com/naturemedicine

While the analyses above suggest that SARS-CoV-2 may bind human ACE2 with high affinity, computational analyses predict that the interaction is not ideal and that the RBD sequence is different from those shown in SARS-CoV to be optimal for receptor binding. Thus, the high-affinity binding of the SARS-CoV-2 spike protein to human ACE2 is most likely the result of natural selection on a human or human-like ACE2 that permits another optimal binding solution to arise. This is strong evidence that SARS-CoV-2 is not the product of purposeful manipulation.

Andersen, K.G., Rambaut, A., Lipkin, W.I., Holmes, E.C., and Garry, R.F. (2020). The proximal origin of SARS-CoV-2. *Nature Medicine* 26, 450–452.

- Andersen et collègues publient dès janvier 2020 un article affirmant que SARS-CoV-2 est sans aucun doute d'origine naturelle.
- Argument principal : la séquence du domaine de liaison au récepteur (RBD) est optimale pour se lier au récepteur ACE2, mais d'une façon différente de celles qu'on connaissait jusqu'alors. D'après les auteurs, si on avait conçu un virus dans le but de le rendre infectieux pour l'homme, on n'aurait pas pu concevoir cette séquence.
 - **Oui mais** les auteurs n'envisagent pas une autre voie possible : en cultivant des virus sur des cellules (en labo) on peut réaliser une sélection artificielle, qui peut déboucher sur des souches adaptées sans nécessiter de connaissance a priori des séquences.
- Autre argument: on ne trouve pas dans ce génome de traces d'ingénierie moléculaire (par exemple des sites de restriction)
 - **Oui mais** les techniques de biologie synthétique permettent depuis 15 ans de générer une molécule d'ADN de novo, sans recourir à des enzymes de restriction (et donc sans trace).

correspondence



The proximal origin of SARS-CoV-2

NATURE MEDICINE | VOL 26 | APRIL 2020 | 450-455 | www.nature.com/naturemedicine

While the analyses above suggest that SARS-CoV-2 may bind human ACE2 with high affinity, computational analyses predict that the interaction is not ideal and that the RBD sequence is different from those shown in SARS-CoV to be optimal for receptor binding. Thus, the high-affinity binding of the SARS-CoV-2 spike protein to human ACE2 is most likely the result of natural selection on a human or human-like ACE2 that permits another optimal binding solution to arise. This is strong evidence that SARS-CoV-2 is not the product of purposeful manipulation.

Andersen, K.G., Rambaut, A., Lipkin, W.I., Holmes, E.C., and Garry, R.F. (2020). The proximal origin of SARS-CoV-2. *Nature Medicine* 26, 450–452.

- Sirotkin & Sirotkin discutent des faiblesses des arguments d'Andersen
- Ils soulignent également l'absence générale d'évaluation sérieuse des hypothèses alternatives concernant les possibilité d'échappement de laboratoire.
- Ils développent l'historique des accidents de laboratoire, et évaluent les scénarios qui permettraient également de comprendre la nature des données en notre possession.
- Des passages successifs d'une souche virale d'une espèce à l'autre (en culture cellulaire ou sur animaux) donneraient le même effet : une divergence accélérée par rapport aux taux de mutations en milieu naturel.

PROBLEMS & PARADIGMS

Prospects & Overviews

BioEssays

www.bioessays-journal.com

Might SARS-CoV-2 Have Arisen via Serial Passage through an Animal Host or Cell Culture?

A potential explanation for much of the novel coronavirus' distinctive genome

Karl Sirotkin* and Dan Sirotkin

Despite claims from prominent scientists that SARS-CoV-2 indubitably emerged naturally, the etiology of this novel coronavirus remains a pressing and open question: Without knowing the true nature of a disease, it is impossible for clinicians to appropriately shape their care, for policy-makers to correctly gauge the nature and extent of the threat, and for the public to appropriately modify their behavior. Unless the intermediate host necessary for completing a natural zoonotic jump is identified, the dual-use gain-of-function research practice of viral serial passage should be considered a viable route by which the novel coronavirus arose. The practice of serial passage mimics a natural zoonotic jump, and offers explanations for SARS-CoV-2's distinctive spike-protein region and its unexpectedly high affinity for angiotensin converting enzyme (ACE2), as well as the notable polybasic furin cleavage site within it. Additional molecular clues raise further questions, all of which warrant full investigation into the novel coronavirus's origins and a re-examination of the risks and rewards of dual-use gain-of-function research.

same genetic signatures behind as a natural jump but occurring in a much shorter period of time.

The genetic signatures in question includes two distinctive features possessed by SARS-CoV-2's spike-protein: the unique sequence in the receptor binding domain (RBD), a region known to be critical for SARS-CoV-2's utilization of human angiotensin converting enzyme (ACE2), which is the cell surface receptor used by both SARS-CoV and SARS-CoV-2 for fusion with target cells and subsequent cell entry. The second feature is the presence of a polybasic furin cleavage site, which is also known as a multibasic cleavage site (MBS)—a four amino acid insertion with limited sequence flexibility—within the coronavirus's novel spike-protein, that is not found in SARS-CoV or other lineage B coronaviruses. This furin cleavage site,

Un virus synthétique avec des bouts de HIV ?

Le 17 avril 2020, le Professeur Luc Montagnier, Prix Nobel de médecine pour sa contribution à la découverte du HIV (le virus responsable du SIDA), défraie la chronique en annonçant sur plusieurs médias (Pourquoi Docteur, CNEWS) que le génome du coronavirus SARS-CoV-2, agent de la pandémie COVID-19, comporte quatre fragments de séquences provenant du HIV. De plus, il affirme que la présence de ces séquences ne résulte pas d'une recombinaison naturelle (fréquente chez les virus) ou d'un accident, mais d'un vrai travail d'ingénieur, effectué intentionnellement, vraisemblablement dans le cadre de recherches visant à développer des vaccins contre le HIV.

Pour appuyer sa théorie, Luc Montagnier cite deux études :

- le travail d'un collègue mathématicien, Jean-Claude Perez, qui "a fouillé les moindres détails de la séquence",
- une analyse des séquences génomiques et protéiques des coronavirus préalablement publiée par une équipe indienne, qui a, selon lui, "été forcée de rétracter" sa publication.

Professeur Luc Montagnier : Le virus covid19 est une manipulation humaine

(<https://www.youtube.com/watch?v=qSWCLHIOiMo>).

"Je suis arrivé à la conclusion qu'il y avait eu une manipulation de ce virus. [...] Il y a un modèle qui est évidemment le virus classique, et là c'était un modèle venant de la chauve-souris, et là, à ce modèle on a par-dessus ajouté les séquences du VIH, du SIDA. ... Non, ce n'est pas naturel, c'était un travail de professionnel, de biologiste moléculaire, très minutieux, on peut dire d'horloger, au niveau des séquences. Dans quel but ce n'est pas clair. Mon travail c'est d'exposer les faits, c'est tout. Je n'accuse personne, je ne sais pas qui a fait ça et pourquoi. La possibilité c'est qu'on a voulu faire un vaccin contre le SIDA. Donc on a pris des petites séquences du virus [HIV] et on les a installées dans la séquence plus grande du coronavirus. [...] Il y a quand même une volonté d'étouffement, nous ne sommes pas les premiers. Un groupe de chercheurs indiens très renommés avaient publié la même chose, on les a forcés à rétracter. Si vous regardez leur publication vous voyez une grande bande "annulé"."

Une arme biologique ?

- Li-Meng Yan
 - chercheuse chinoise réfugiée aux Etats-Unis
 - travaillait dans le laboratoire de référence de l'OMS pour la Chine
- En septembre et octobre 2020, elle publie sur zenodo (dépôt d'oeuvres électroniques, sans révision par des experts) deux articles, où elle affirme que
 - le virus résulte de manipulations génétiques reposant sur les méthodes classiques de biologie moléculaire (recombinaison d'ADN au moyen d'enzymes de restriction) ;
 - ce virus est une arme biologique.
- Ses articles n'ont pas encore été publiés, mais ils sont téléchargés et fortement médiatisés.
- Les arguments sous-jacents font cependant l'objet de critiques par les spécialistes.

Yan, Li-Meng, Kang, Shu, Guan, Jie & Hu, Shanchang. 2020a. SARS-CoV-2 Is an Unrestricted Bioweapon: A Truth Revealed through Uncovering a Large-Scale, Organized Scientific Fraud. , doi: [10.5281/ZENODO.4073131](https://doi.org/10.5281/ZENODO.4073131). Zenodo.

Yan, Li-Meng, Kang, Shu, Guan, Jie & Hu, Shanchang. 2020b. Unusual Features of the SARS-CoV-2 Genome Suggesting Sophisticated Laboratory Modification Rather Than Natural Evolution and Delineation of Its Probable Synthetic Route. , doi: [10.5281/ZENODO.4028830](https://doi.org/10.5281/ZENODO.4028830). Zenodo.

zenodo Search Upload Communities Log in Sign up

September 14, 2020 Working paper Open Access

Unusual Features of the SARS-CoV-2 Genome Suggesting Sophisticated Laboratory Modification Rather Than Natural Evolution and Delineation of Its Probable Synthetic Route

Yan, Li-Meng; Kang, Shu; Guan, Jie; Hu, Shanchang

The COVID-19 pandemic caused by the novel coronavirus SARS-CoV-2 has led to over 910,000 deaths worldwide and unprecedented decimation of the global economy. Despite its tremendous impact, the origin of SARS-CoV-2 has remained mysterious and controversial. The natural origin theory, although widely accepted, lacks substantial support. The alternative theory that the virus may have come from a research laboratory is, however, strictly censored on peer-reviewed scientific journals. Nonetheless, SARS-CoV-2 shows biological characteristics that are inconsistent with a naturally occurring, zoonotic virus. In this report, we describe the genomic, structural, medical, and literature evidence, which, when considered together, strongly contradicts the natural origin theory. The evidence shows that SARS-CoV-2 should be a laboratory product created by using bat coronaviruses ZC45 and/or ZXC21 as a template and/or backbone. Building upon the evidence, we further postulate a synthetic route for SARS-CoV-2, demonstrating that the laboratory-creation of this coronavirus is convenient and can be accomplished in approximately six months. Our work emphasizes the need for an independent investigation into the relevant research laboratories. It also argues for a critical look into certain recently published data, which, albeit problematic, was used to support and claim a natural origin of SARS-CoV-2. From a public health perspective, these actions are necessary as knowledge of the origin of SARS-CoV-2 and of how the virus entered the human population are of pivotal importance in the fundamental control of the COVID-19 pandemic as well as in preventing similar, future pandemics.

776,634 views 597,172 downloads See more details...

Indexed in OpenAIRE

Publication date: September 14, 2020
DOI: [10.5281/zenodo.4028830](https://doi.org/10.5281/zenodo.4028830)
Communities: Coronavirus Disease Research Community - COVID-19
License (for files):

zenodo Search Upload Communities Jacques.van-Helden@univ-amu.fr

October 8, 2020 Working paper Open Access

SARS-CoV-2 Is an Unrestricted Bioweapon: A Truth Revealed through Uncovering a Large-Scale, Organized Scientific Fraud

Yan, Li-Meng; Kang, Shu; Guan, Jie; Hu, Shanchang

Two possibilities should be considered for the origin of SARS-CoV-2: natural evolution or laboratory creation. In our earlier report titled "Unusual Features of the SARS-CoV-2 Genome Suggesting Sophisticated Laboratory Modification Rather Than Natural Evolution and Delineation of Its Probable Synthetic Route", we disproved the possibility of SARS-CoV-2 arising naturally through evolution and instead proved that SARS-CoV-2 must have been a product of laboratory modification. Despite this and similar efforts, the laboratory creation theory continues to be downplayed or even diminished. This is fundamentally because the natural origin theory remains supported by several novel coronaviruses published after the start of the outbreak. These viruses (the RaTG13 bat coronavirus, a series of pangolin coronaviruses, and the RmYN02 bat coronavirus) reportedly share high sequence homology with SARS-CoV-2 and have altogether constructed a seemingly plausible pathway for the natural evolution of SARS-CoV-2. Here, however, we use in-depth analyses of the available data and literature to prove that these novel animal coronaviruses do not exist in nature and their sequences have been fabricated. In addition, we also offer our insights on the hypothesis that SARS-CoV-2 may have originated naturally from a coronavirus that infected the Mojiang miners.

Revelation of these virus fabrications renders the natural origin theory unfounded. It also strengthens our earlier assertion that SARS-CoV-2 is a product of laboratory modification, which can be created in approximately six months using a template virus owned by a laboratory of the People's Liberation Army (PLA). The fact that data fabrications were used to cover up the true origin of SARS-CoV-2 further implicates that the laboratory modification here is beyond simple gain-of-

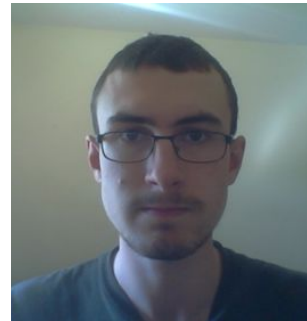
161,914 views 90,966 downloads See more details...

Indexed in OpenAIRE

Publication date: October 8, 2020
DOI: [10.5281/zenodo.4073131](https://doi.org/10.5281/zenodo.4073131)
License (for files): [Creative Commons Attribution 4.0 International](https://creativecommons.org/licenses/by/4.0/)

A l'origine de ce cours

- Ce cours repose en grande partie sur un travail d'équipe mené durant le confinement d'avril-mai 2020.
- Nous remercions les nombreux collègues qui nous ont suggéré des améliorations sur les premières versions du manuscrit.
- Afin d'assurer la traçabilité et la reproductibilité de nos analyses, les données utilisées et les logiciels développés sont en libre accès
 - https://ivanheld.github.io/SARS-CoV-2_origins/



Erwan Sallard

Étudiant à l'Institut de Biologie
Ecole Normale Supérieure, Paris



Didier Casane

Biologie évolutive
Professeur, Université de Paris



José Halloy

Modélisation des systèmes biologiques
Professeur, Université de Paris



Etienne Decroly

Virologie
Dir. de Recherches, CNRS



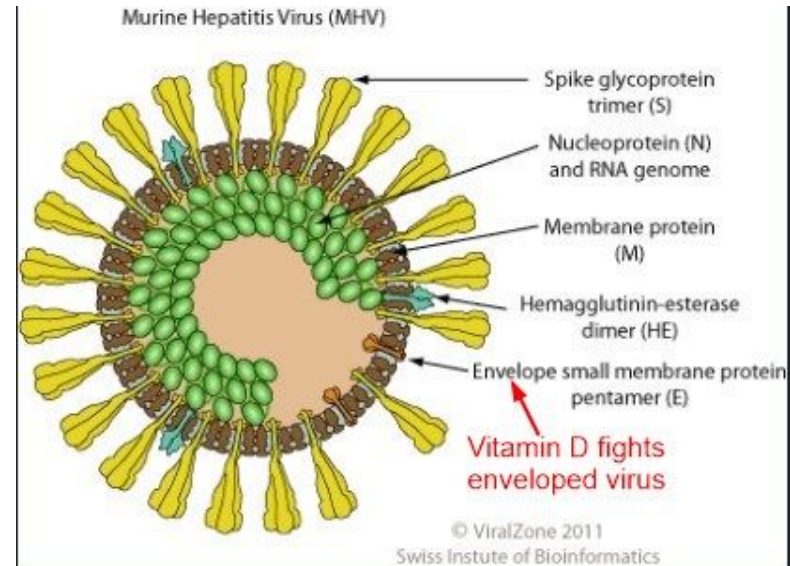
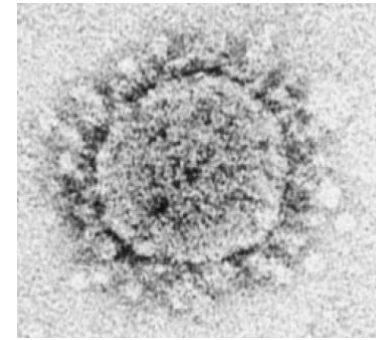
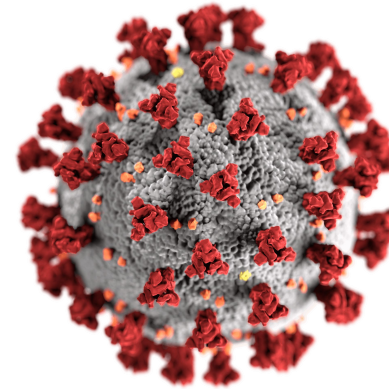
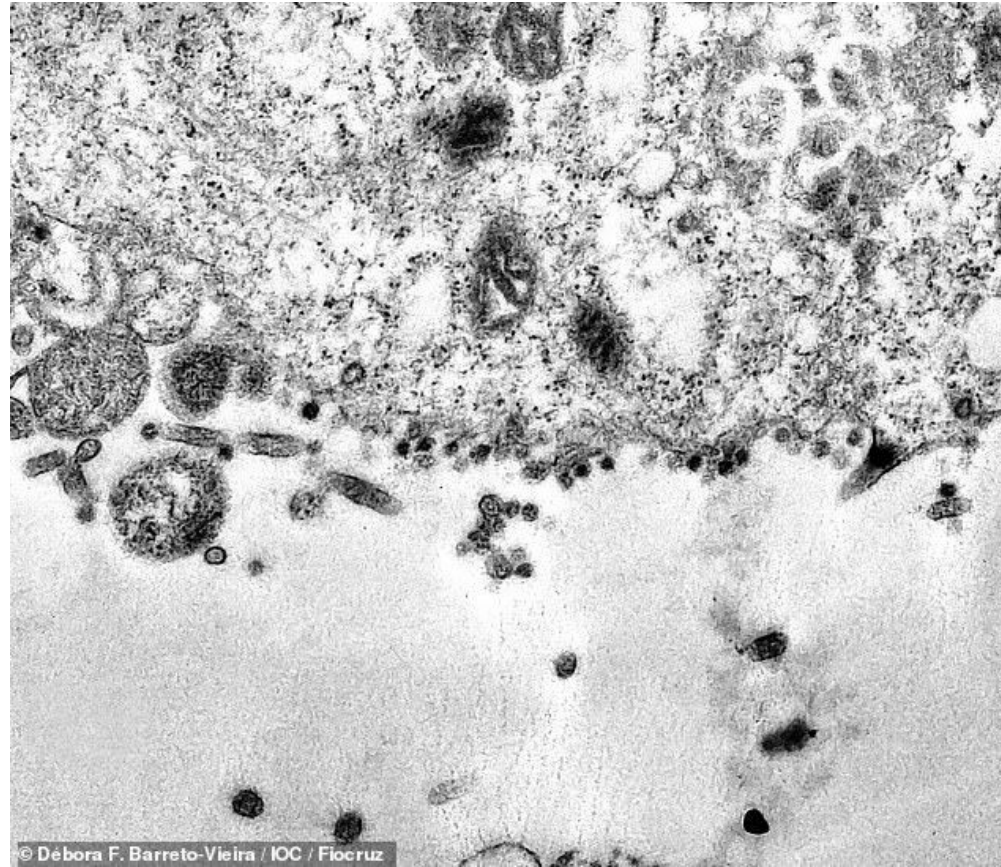
Jacques van Helden

Bioinformatique
Professeur, Aix-Marseille Université

- Sallard, E., Halloy, J., Casane, D., van Helden, J. & Decroly, É. 2020. **Retrouver les origines du SARS-CoV-2 dans les phylogénies de coronavirus.** *Med Sci (Paris)* 36: 783–796.
- English version : Erwan Sallard, José Halloy, Didier Casane, Etienne Decroly, Jacques van Helden. **Tracing the origins of SARS-CoV-2 in coronavirus phylogenies.** [hal-02891455](https://doi.org/10.1007/s12250-020-00145-5)

Biologie de SARS-CoV-2

Coronavirus



Le génome des coronavirus est constitué d'ARN

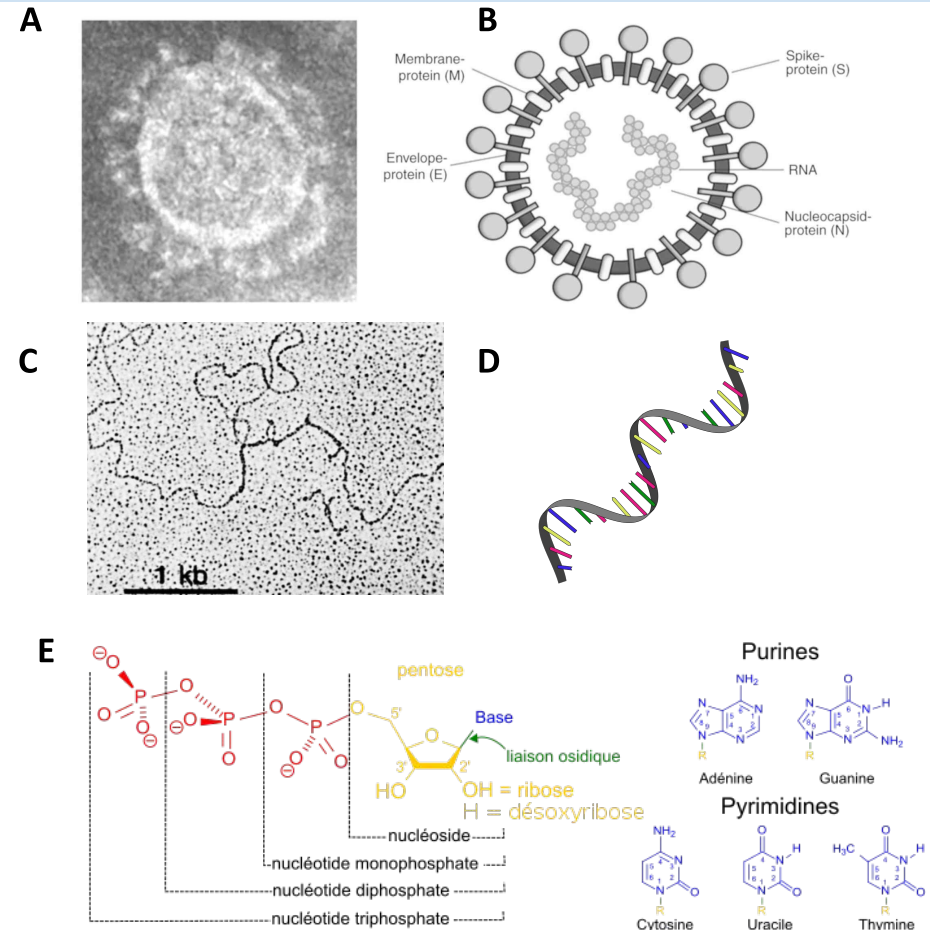
- Le génome est l'ensemble du matériel génétique d'un organisme
- Le génome des coronavirus est constitué d'acide ribonucléique (ARN).

Légende de la figure

- A. Micrographie électronique d'un virion de coronavirus.
- B. Schéma de la structure du virion. Le chapelet à l'intérieur symbolise l'ARN.
- C. Micrographie électronique de l'ARN viral extrait de son enveloppe
- D. Schéma de la structure de l'ARN. Les bâtonnets de couleur représentent les nucléotides

Adénine
Cytosine
Guanine
Uracile

- E. Structure chimique des nucléotides



Génome de SARS-CoV-2 de référence

- Début de la séquence génomique de SARS-CoV-2
- La taille totale fait 29 899 nucléotides
- Un des plus grands génomes parmi les virus à ARN
 - ❑ A Adénine
 - ❑ C Cytosine
 - ❑ G Guanine
 - ❑ T Uracile (on remplace le U par un T)

```
>MT019529.1 Severe acute respiratory syndrome coronavirus 2 isolate
BetaCoV/Wuhan/IPBCAMS-WH-01/2019, complete genome
ATTTAAAGGTTTATACCTTCCAGGTAACAAACCACTTTTCGATCTCTTGTAGATCTGTTCTCTAAA
CGAACTTTAAAATCTGTGTGGCTGTCACCTCGGCTGCATGCTTAGTGCACCTCACGCAGTATAATTAATAAC
TAATTAAGTGTGCTGTGACAGGACACGAGTAACTCGTCTATCTTCTGCAGGCTGCTTACGGTTTCGTCCTGTG
TTGCAGCCGATCATCAGCACATCTAGGTTTCGTCGGGTGTGACCGAAAGGTAAGATGGAGAGCCTTGTCTC
CCTGGTTTCAACGAGAAAAACACAGTCCAACTCAGTTTGCCTGTTTTACAGGTTTCGCGACGTGCTCGTAC
GTGGCTTTGGAGACTCCGTGGAGGAGTCTTATCAGAGGCACGTCAACATCTTAAAGATGGCACTTGTGG
CTTAGTAGAAGTTGAAAAGGCGTTTTCCTCAACTTGAACAGCCCTATGTGTTTCAACAGTTTCGGAT
GCTCGAAGTGCACCTCATGGTTCATGTTATGGTTGAGCTGGTAGCAGAACTCGAAGGCATTTCAGTACGGTC
GTAGTGGTGGAGACTTGGTGTCTTGTCCCTCATGTGGGCGAAATACAGTGGCTTACCGCAAGGTTCT
TCTTCGTAAGAACGGTAATAAAGGAGCTGGTGGCCATAGTTACGGCGCCGATCTAAAGTCAATTTGACTTA
GGCGACGAGCTTGGCACTGATCCTTATGAAGATTTTCAAGAAAAGTGGAACTAAACATAGCAGTGGTG
TTACCCGTGAACTCATGCGTGAGCTTAAACGGAGGGGCATACACTCGCTATGTCGATAACAACCTTCTGTGG
CCCTGATGGCTACCCTCTTGTAGTGCATTAAGACCTTCTAGCACGTGCTGGTAAAGCTTTCATGCACTTTG
TCCGAACAAGTGGACTTTATTGACACTAAGAGGGGTGTATACTGCTGCCGTGAACATGAGCATGAAATTTG
CTTGGTACACGGAACGTTCTGAAAAGAGCTATGAATTGCAGACACCTTTTGAATTAATTTGGCAAAGAA
ATTTGACACCTTCAATGGGGAATGTCCAAATTTTGTATTTCCCTTAAATTCATAATCAAGACTATTCAA
CCAAGGTTGAAAAGAAAAGCTTGTATGGCTTTATGGGTAGAATTCGATCTGTCTATCCAGTTGCGTCAC
CAAATGAATGCAACCAAATGTGCCTTTCAACTCTCATGAAGTGTGATCATTGTGGTGAACCTTTCATGGCA
GACGGGCGATTTTGTAAAGCCACTTGCGAATTTTGTGGCACTGAGAATTTGACTAAAGAAGGTGCCACT
ACTTGTGGTTACTTACCCCAAATGCTGTGTTTAAATTTATTGTCCAGCATGTCACAATTCAGAAGTAG
GACCTGAGCATAGTCTTGCAGAAATACATAATGAATCTGGCTTGAACCACTTCTTCGTAAGGGTGGTCCG
CACTATTGCCTTTGGAGGCTGTGTGTTCTCTTATGTTGGTTGCCATAACAAGTGTGCCATTATGGGTTCCA
CGTGCTAGCGCTAACATAGGTTGTAACCATACAGGTGTTGTTGGAGAAGGTTCCGAAGGCTTAAATGACA
ACCTTCTTGAAATACTCCAAAAGAGAAAGTCAACATCAATATGTTGGTGACTTTAACTTAAATGAAGA
GATCGCCATTTATTTGGCATCTTTTCTGCTTCCACAAGTGTCTTTGTGGAAACTGTGAAAGGTTTGGAT
TATAAAGCATTCAAAACAAATGTTGAATCCTGTGGTAATTTTAAAGTTACAAAAGGAAAAGCTAAAAAG
GTGCTTGAATATTTGGTGAACAGAAATCAATACTGAGTCTCTTATGCAATTTGCATCAGAGGCTGCTCG
```

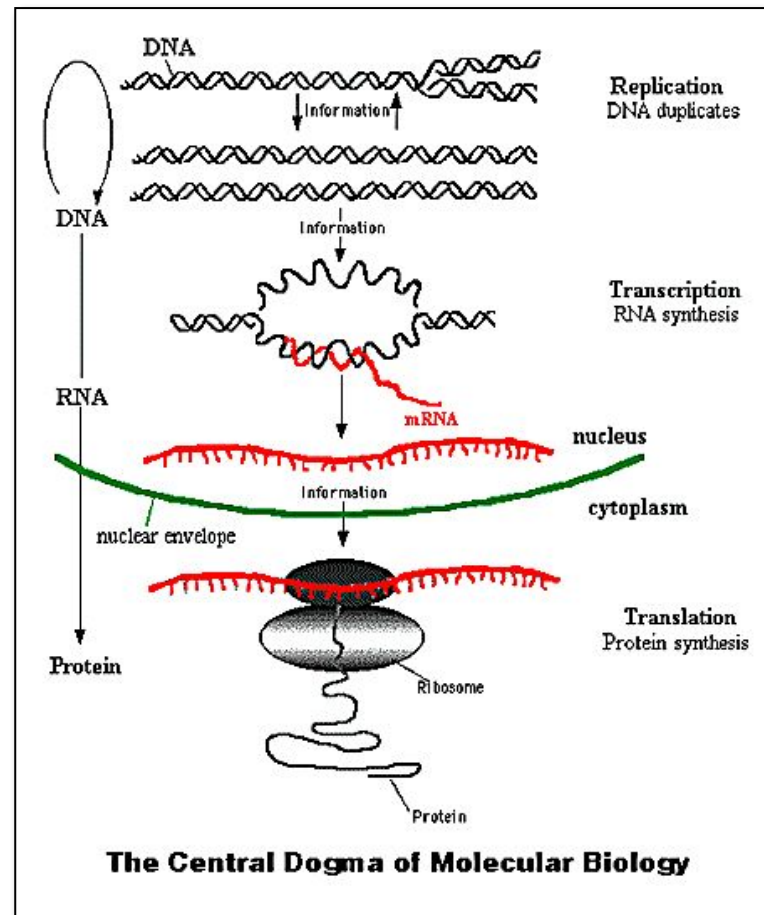
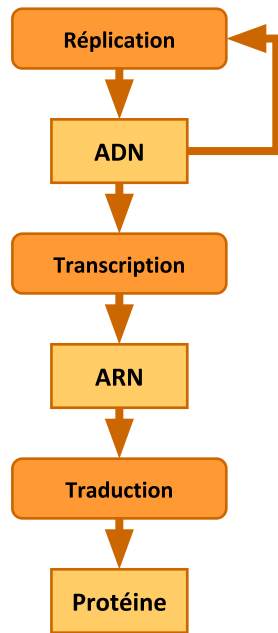
- Le séquençage d'un génome viral représente un coût modique et se fait très rapidement.
- On dispose de
 - plusieurs dizaines de génomes de coronavirus non-humains
 - plusieurs dizaines de milliers de génomes de SARS-CoV-2 (humain)
- Ces données permettent de poser des questions concernant
 - l'origine animale de SARS-CoV-2
 - sa propagation dans la population humaine

Le "dogme central" de la biologie

- Formulé en 1958 par Francis Crick
 - Crick, F. H. (1958). On protein synthesis. Symp Soc Exp Biol 12, 138-63.
 - Je recommande également de lire cette discussion ultérieure :
Crick, F. (1970). Central dogma of molecular biology. Nature 227, 561-3.
- On le résume souvent de la façon suivante
 - **DNA makes RNA makes protein**
(l'ADN fait l'ARN qui fait les protéines)
- Cependant le dogme ne dit pas exactement cela. Il énonce les transferts d'information qui sont possibles ou impossibles entre les séquences d'acides nucléiques et celles des protéines.

Le dogme central stipule que, une fois que l' « information » est passée dans la protéine elle ne peut pas en ressortir. Plus précisément, le transfert d'information serait possible d'acide nucléique à acide nucléique, ou d'acide nucléique à protéine, mais le transfert de protéine à protéine, ou de protéine à acide nucléique est impossible. Information signifie ici la détermination précise de la séquence, soit des bases dans l'acide nucléique, soit des résidus aminoacides dans la protéine.

Crick, F. H. (1958). On protein synthesis. Symp Soc Exp Biol 12, 138-63.

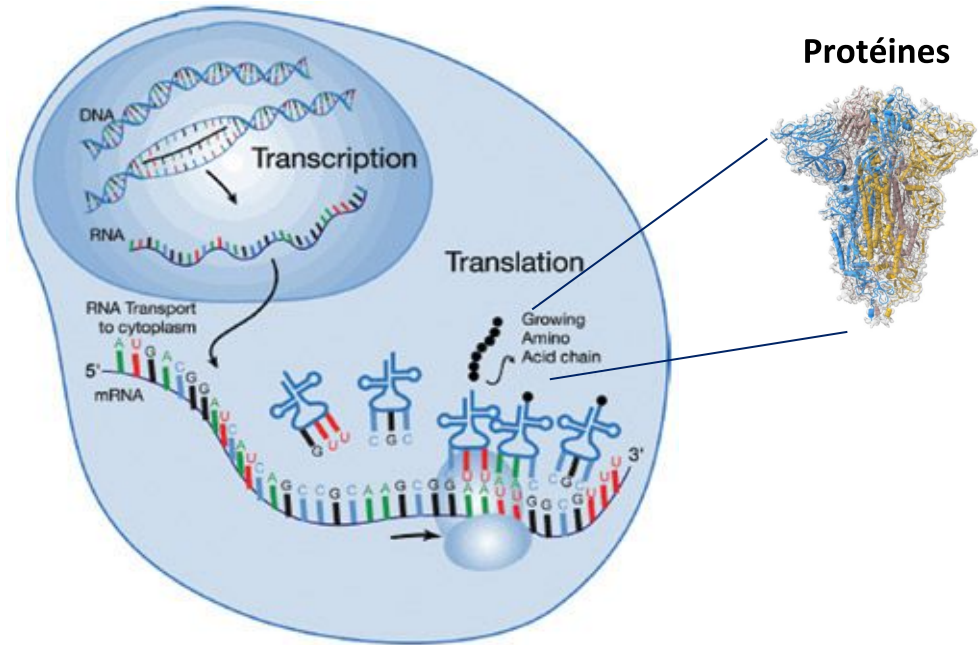


<http://www.accessexcellence.org/AB/GG/central.html>

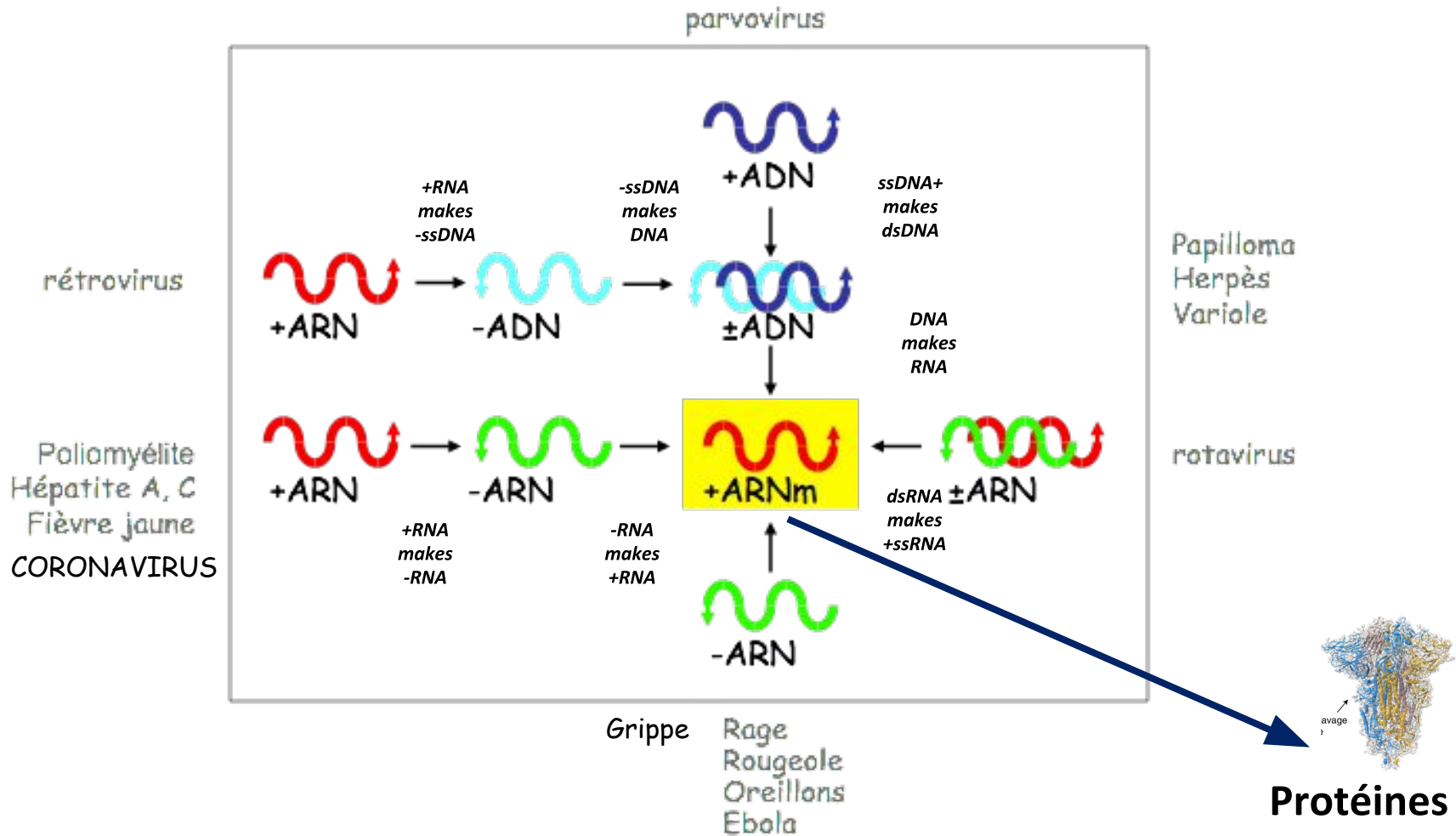
DNA makes RNA makes protein*

* Formulation compacte des flux de l'information moléculaire, par Francis Crick

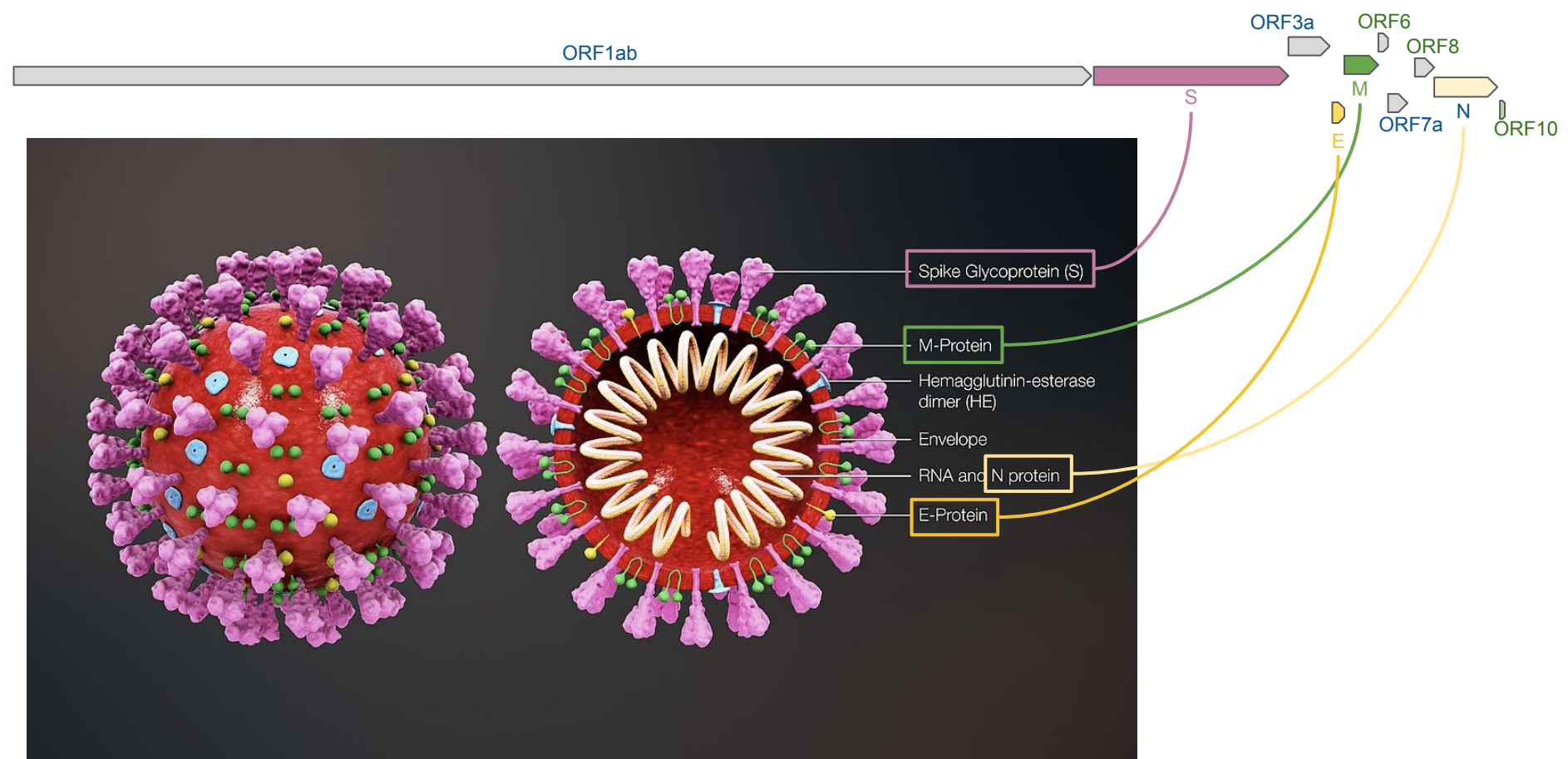
- Dans tous les organismes cellulaires
 - L'ADN sert de modèle à la synthèse d'ARN (**transcription**)
 - L'ARN sert de modèle à la synthèse des protéines (**traduction**)
- Les protéines sont les principaux acteurs moléculaires des organismes vivants
 - Enzymes
 - Transporteurs
 - Régulateurs
 - Cycle cellulaire
 - Différenciation cellulaire
 - ... un tas d'autres fonctions



Le matériel génétique des virus

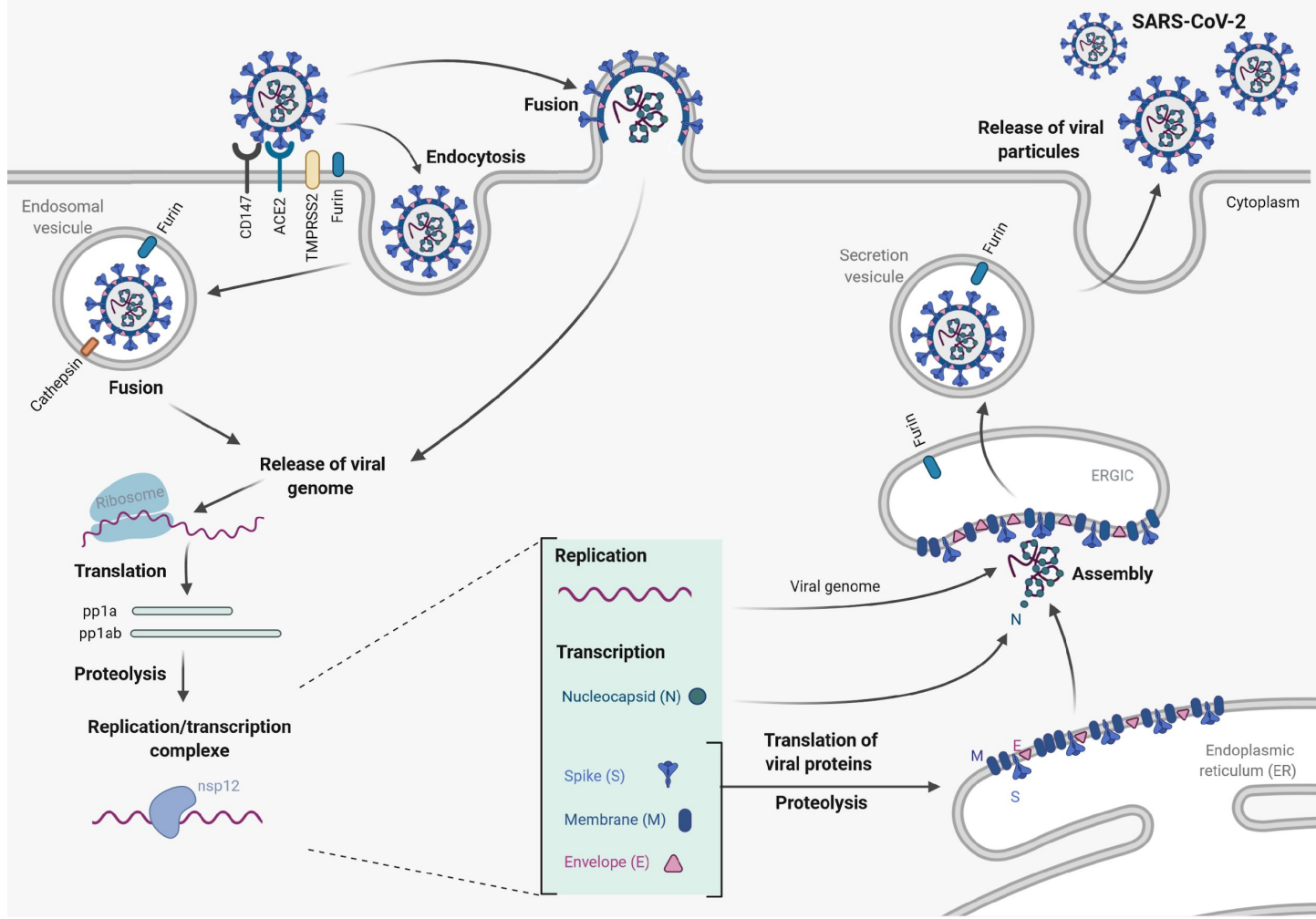


Fonction des gènes de SARS-CoV-2 – Gènes structuraux

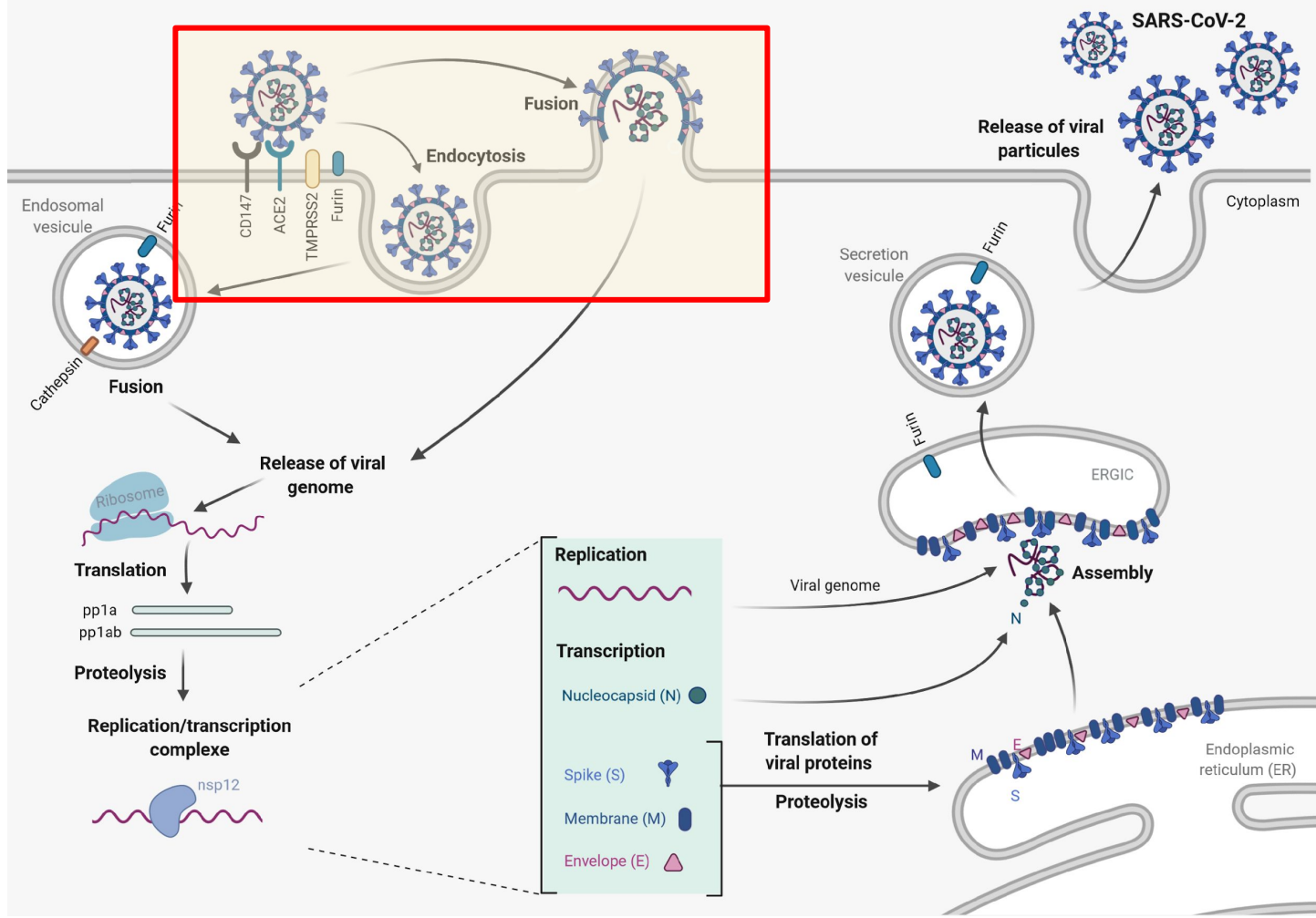


Mécanismes d'infection

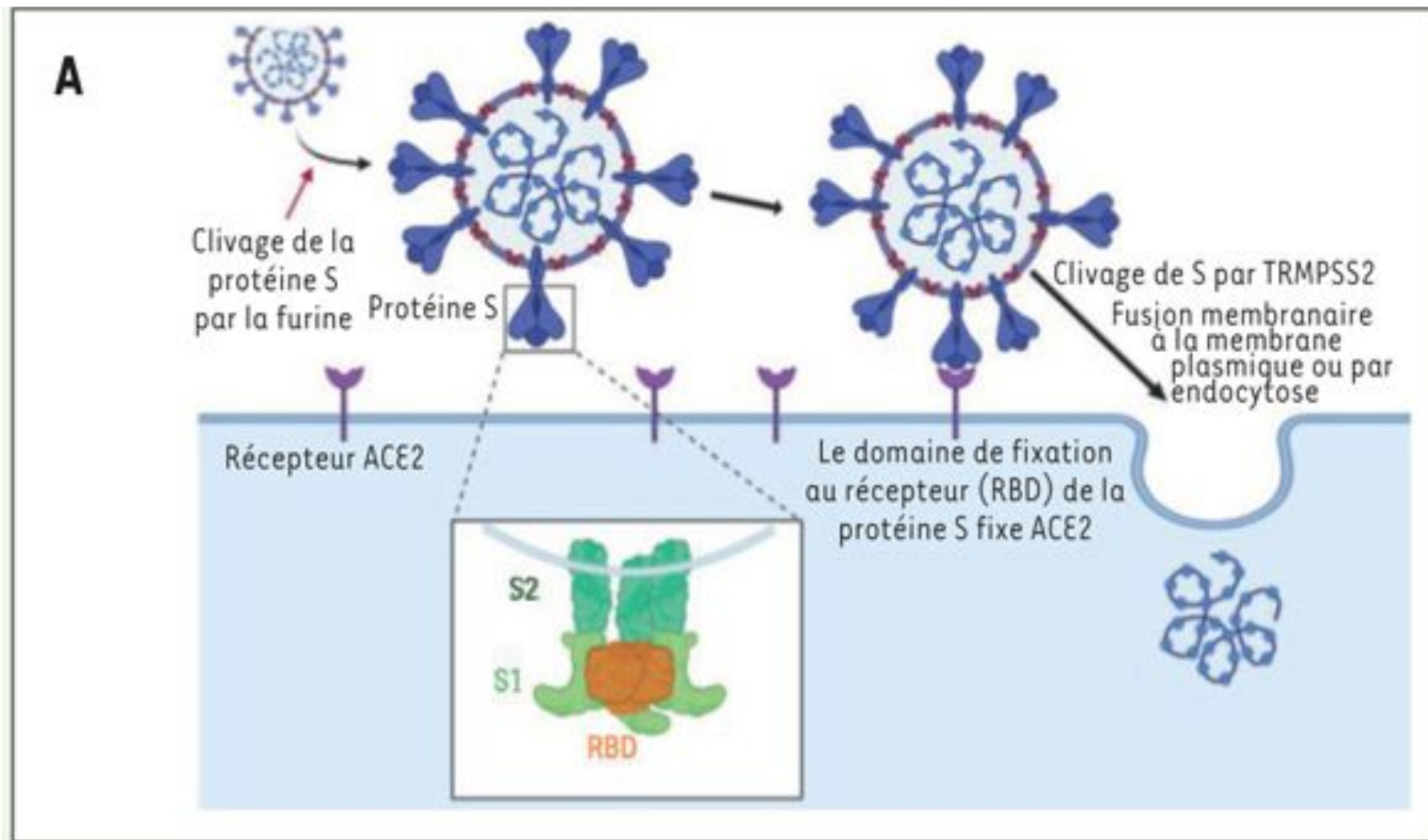
Cycle de répliation des coronavirus



Cycle de réplication des coronavirus - Entrée dans la cellule hôte

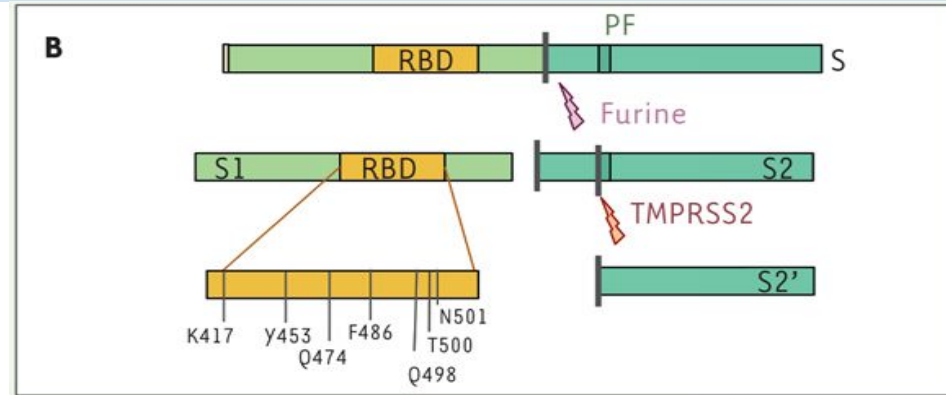


La protéine S reconnaît un récepteur cellulaire (SARS likes ACE2)

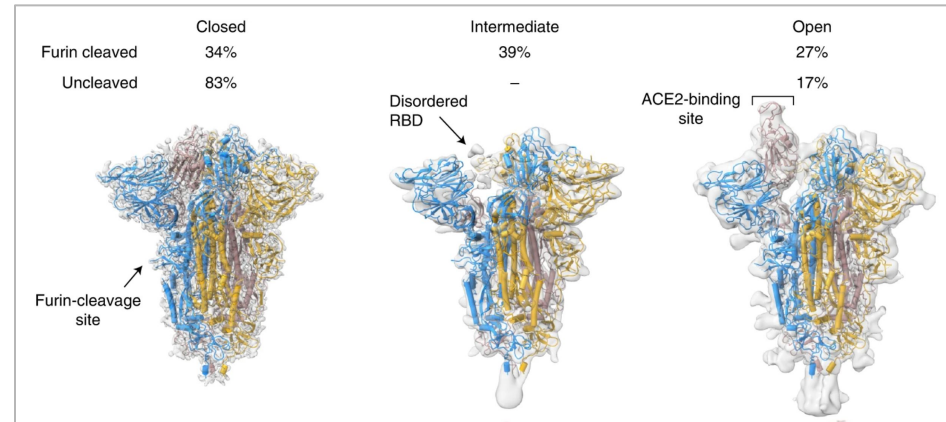


La protéine S doit être activée par un clivage protéolytique (priming)

- Juste après la traduction, le produit du gène S est une protéine inactive.
- L'activation requiert le clivage (coupure) de la protéine en deux parties (S1, S2).
- Chez SARS-CoV-2, ce clivage est réalisé par une enzyme appelée furine.
- Ceci a une grave conséquence, car cette enzyme est ubiquitaire dans les cellules humaines, ce qui explique en partie que les symptômes de la Covid-19 ne se limitent pas aux voies respiratoires.



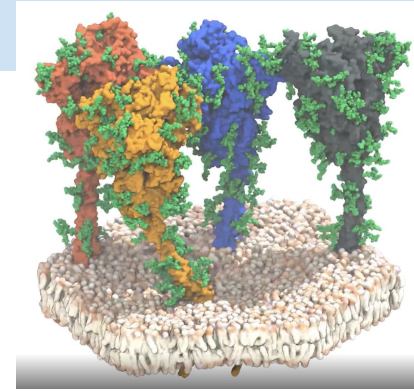
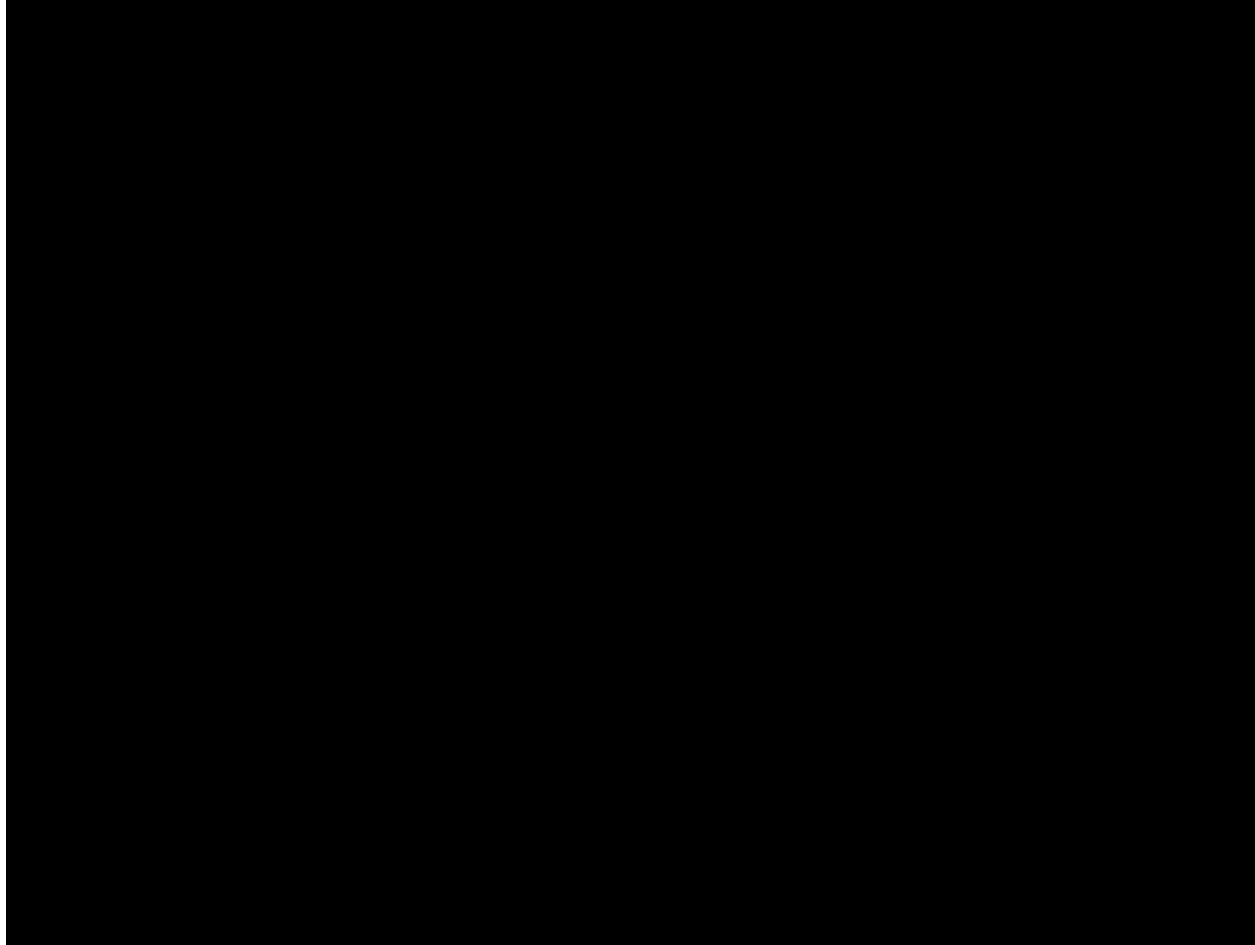
Sallard, E., Halloy, J., Casane, D., van Helden, J. & Decroly, É. 2020. Retrouver les origines du SARS-CoV-2 dans les phylogénies de coronavirus. *Med Sci (Paris)* 36: 783–796



Wrobel, A.G., Benton, D.J., Xu, P., Roustan, C., Martin, S.R., Rosenthal, P.B., et al. 2020. SARS-CoV-2 and bat RaTG13 spike glycoprotein structures inform on virus evolution and furin-cleavage effects. *Nat Struct Mol Biol* 27: 763–767.

Coutard, B., Valle, C., de Lamballerie, X., Canard, B., Seidah, N.G. & Decroly, E. 2020. The spike glycoprotein of the new coronavirus 2019-nCoV contains a furin-like cleavage site absent in CoV of the same clade. *Antiviral Res.* 176: 104742.

Structure de la protéine S & dynamique moléculaire



Simulation moléculaire de l'interaction spicule - récepteur

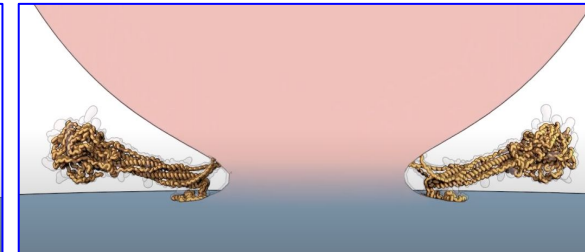
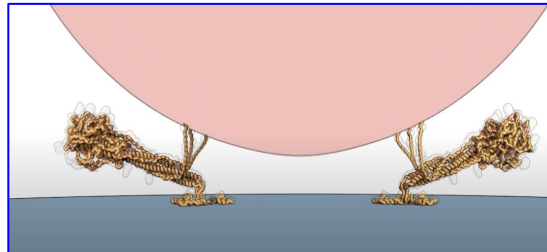
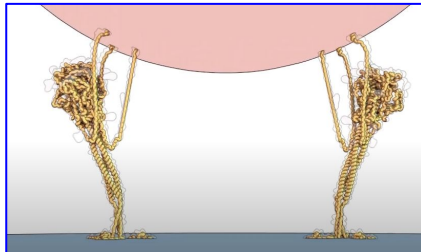
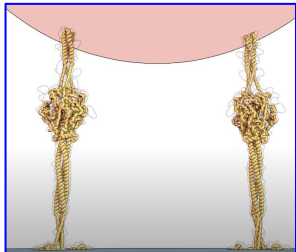
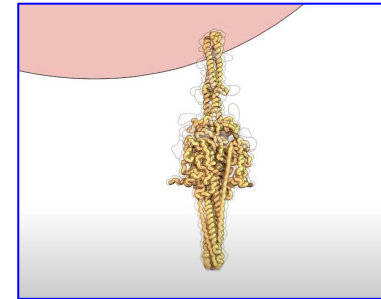
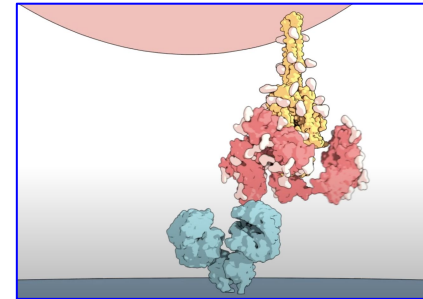
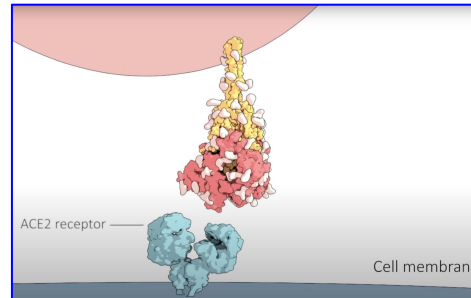
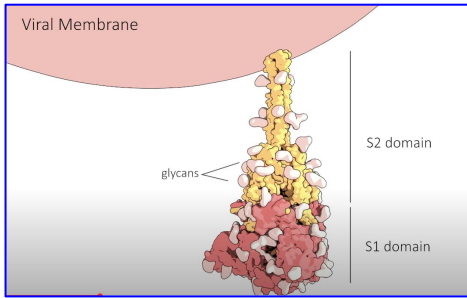
<https://youtu.be/e2Qi-hAXdJo?t=12>

Un modèle visuel de simulation dynamique illustrant la façon dont la protéine de spicule (spike protein) assure la fusion des membranes virale et cellulaire.

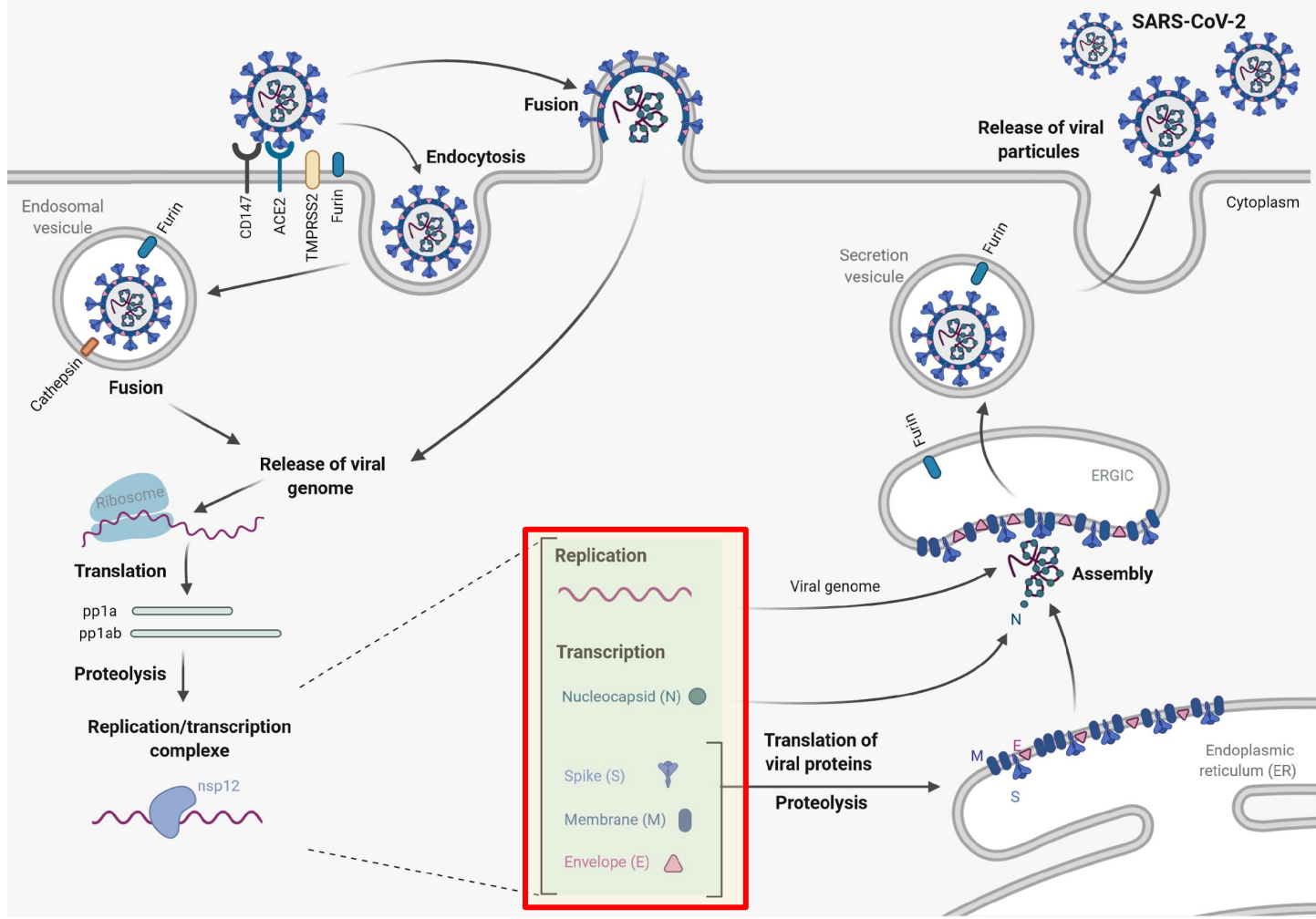
Created by Jonathan Khao, Ph.D. & Gaël McGill, Ph.D.
Digizyme Inc.
www.digizyme.com

Modeled & Simulated with Molecular Maya (Modeling & Rigging kits)
www.clarafi.com/tools/mmaya

We wish to thank Bing Chen, Ph.D. and Stephen Harrison, Ph.D.
for their guidance and sharing data prior to publication.



Cycle de réplication des coronavirus - Réplication et transcription

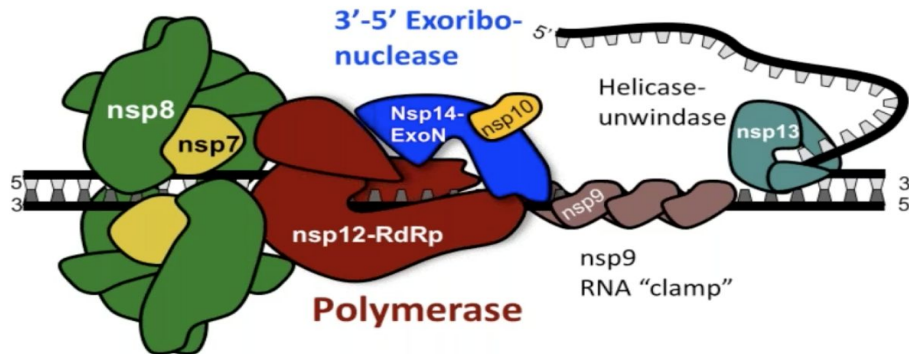
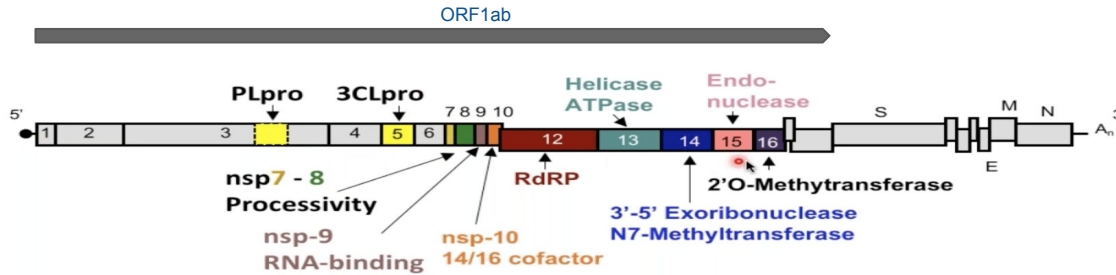


Le complexe réplication/transcription

SARS-CoV-2

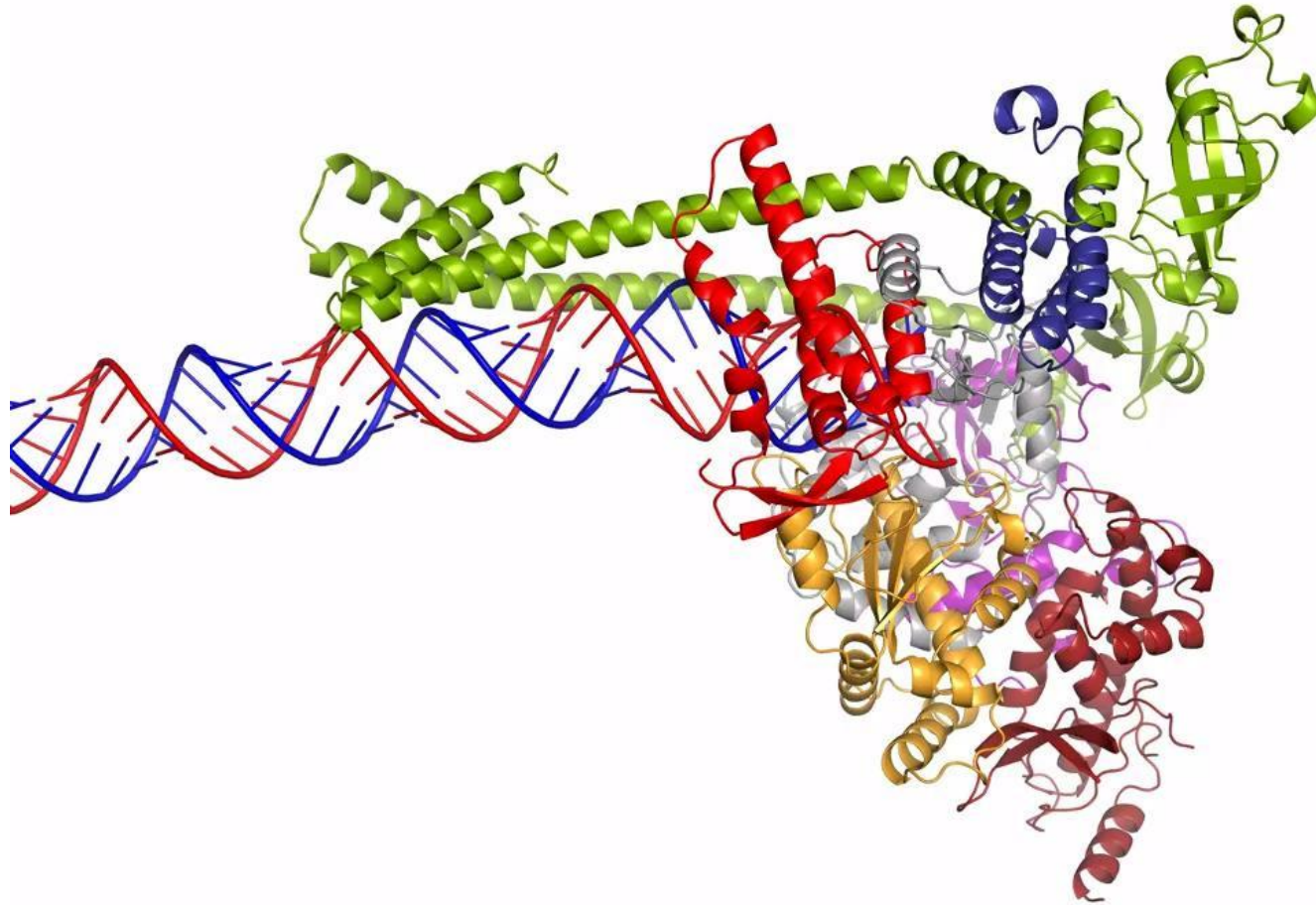


SARS-CoV

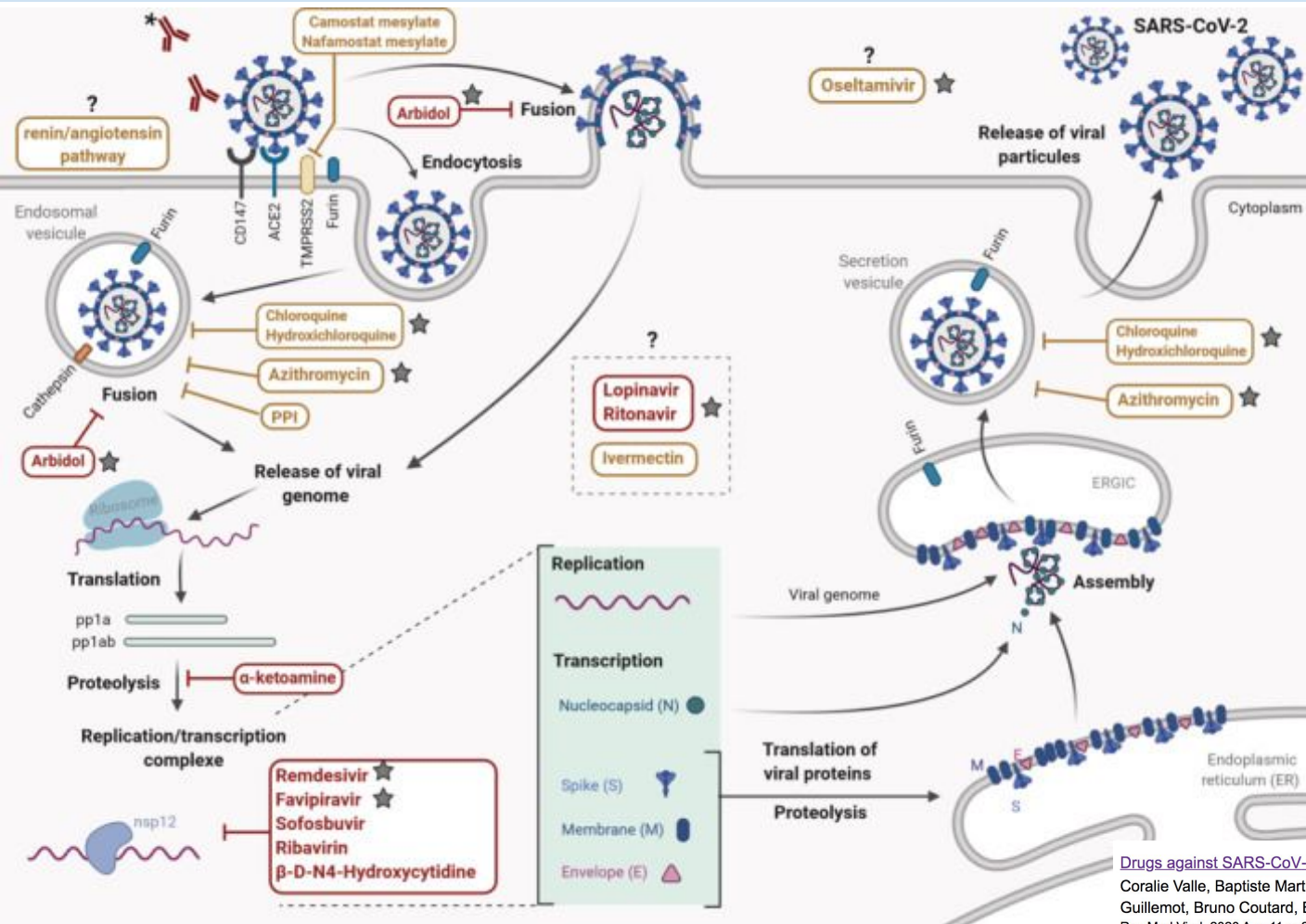


- Le gène ORF1ab code pour une “polyprotéine”, qui contient 16 protéines distinctes
- Une dizaine de ces protéines forment un complexe (figure du bas) qui assure la réplication et la transcription de l’ARN.

SARS-CoV : le complexe répllication/transcription



Mécanismes d'action des médicaments contre le SARS-CoV-2



- Des études sont menées dans un grand nombre de laboratoires pour identifier des molécules qui pourraient bloquer l'infection par SARS-CoV-2.
- Ces études partent de médicaments antiviraux connus, qui interagissent à différents niveaux du cycle infectieux des virus.

Drugs against SARS-CoV-2: What do we know about their mode of action?

Coralie Valle, Baptiste Martin, Franck Touret, Ashleigh Shannon, Bruno Canard, Jean-Claude Guillemot, Bruno Coutard, Etienne Decroly
 Rev Med Virol. 2020 Aug 11 : e2143. doi: 10.1002/rmv.2143 [Epub ahead of print]

Événements évolutifs

Typologie des mutations

- Substitution

- Remplacement d'un résidu (une lettre) par un autre

Avant répllication

ATGACCATGA



Après répllication

ATGACCA**G**GA

- Délétion

- Perte d'un fragment de la molécule

ATG**ACCAT**GA



↓
ATGGA

- Insertion

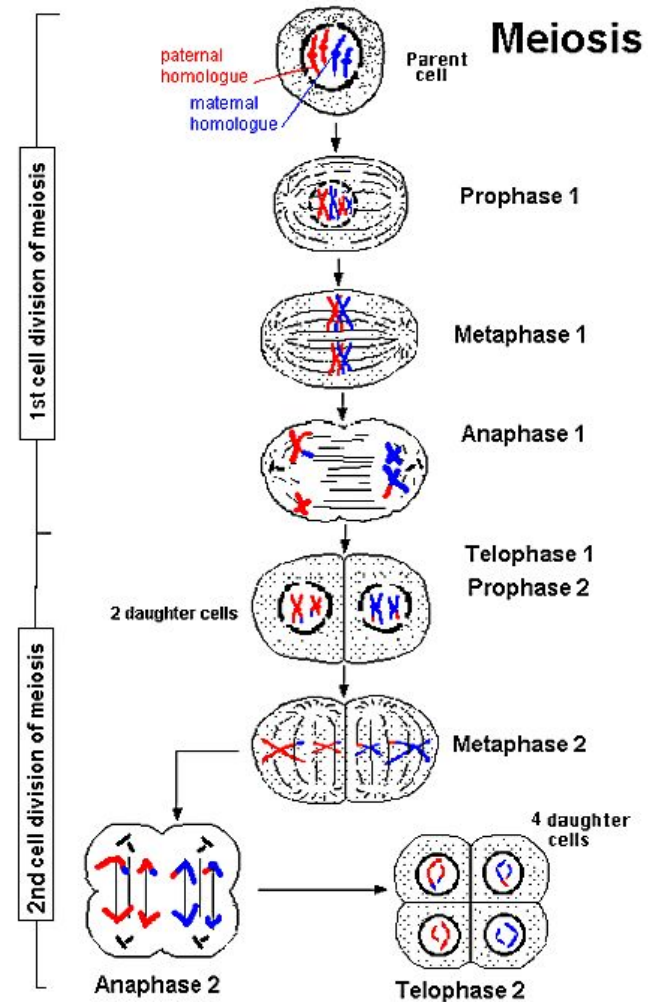
- Ajout d'un fragment de molécule

ATGACCATGA



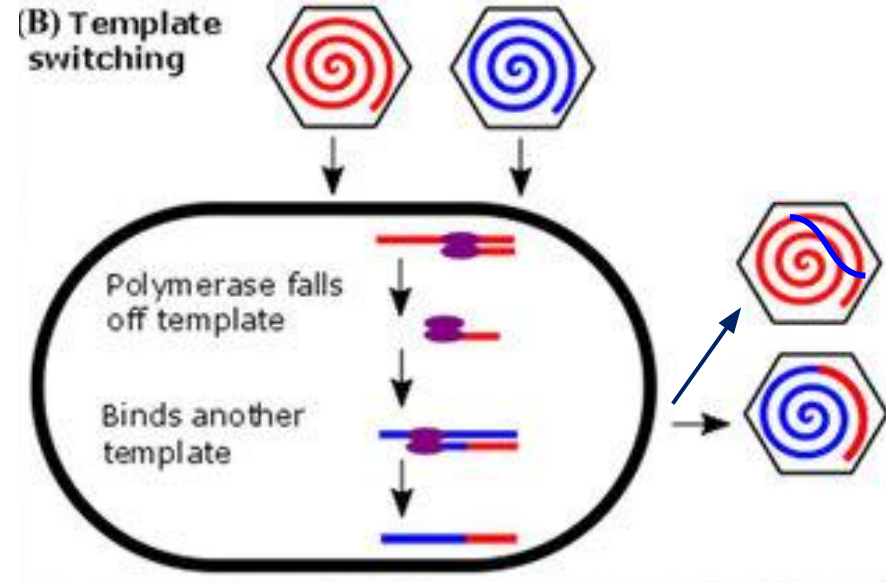
ATGA**CAAA**CATGA

- Chez les organismes cellulaires, lors de la méiose, une cellule mère diploïde forme 4 cellules-filles haploïdes.
- Les chromosomes parentaux sont distribués aléatoirement et de façon indépendante entre les 4 cellules-filles.
- Des événements de “crossing-over” (croisements) provoquent une recombinaison de fragments de chromosomes.
- La liaison génétique entre les gènes d'un même chromosome n'est pas complète.



Recombinaisons chez les coronavirus

- Une chauve-souris peut se retrouver infectée par plusieurs coronavirus en même temps.
- Pendant la réplication des coronavirus, il arrive que la polymérase de l'ARN "saute" d'un virus à l'autre.
- Ceci donne naissance à un virus "chimérique", dont le génome est composé de fragments d'origines différentes.
- Ceci complique l'analyse de la phylogénie des virus, car différents fragments génomiques résultent d'histoires évolutives différentes.

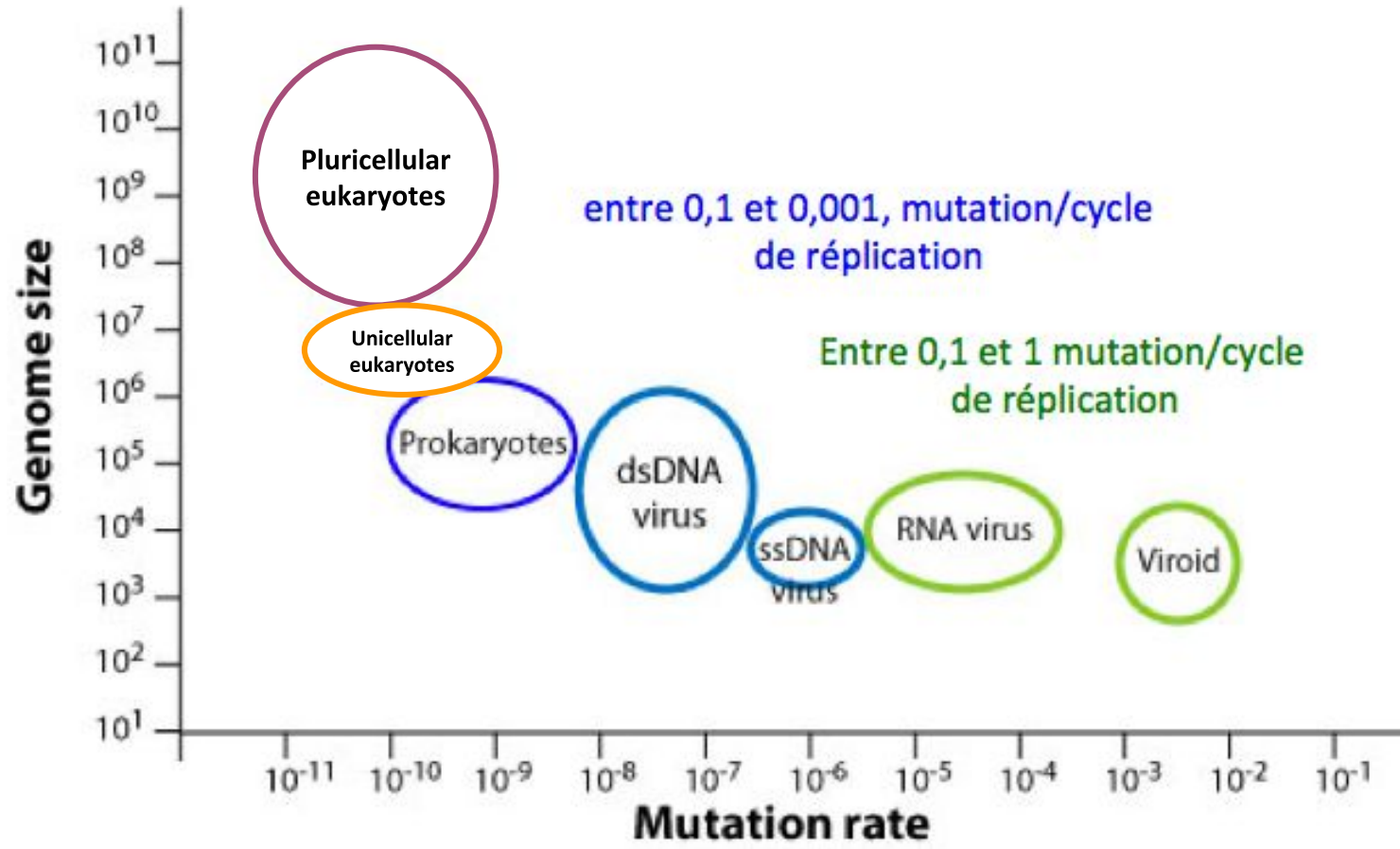


Discovery of a rich gene pool of bat SARS-related coronaviruses provides new insights into the origin of SARS coronavirus

Ben Hu ,et al Plos path : November 30, 2017

J. Dennehy, Evolutionary ecology of virus emergence: Virus emergence, 2016, Annals of the New York Academy of Science

Taux de substitutions représentatifs chez différents groupes taxonomiques



Alignement de séquences – Gènes S de SARS-CoV-2 et RaTG13

```
# Aligned_sequences: 2
# 1: Human_SARS-CoV-2_BetaCoV/Wuhan/IPBCAMS-WH-01/2019
# 2: Bat_RaTG13
#
# Length: 3822
# Identity: 3549/3822 (92.9%)
# Similarity: NA/3822 (NA%)
# Gaps: 12/3822 (0.3%)
# Score: 5435.624
```

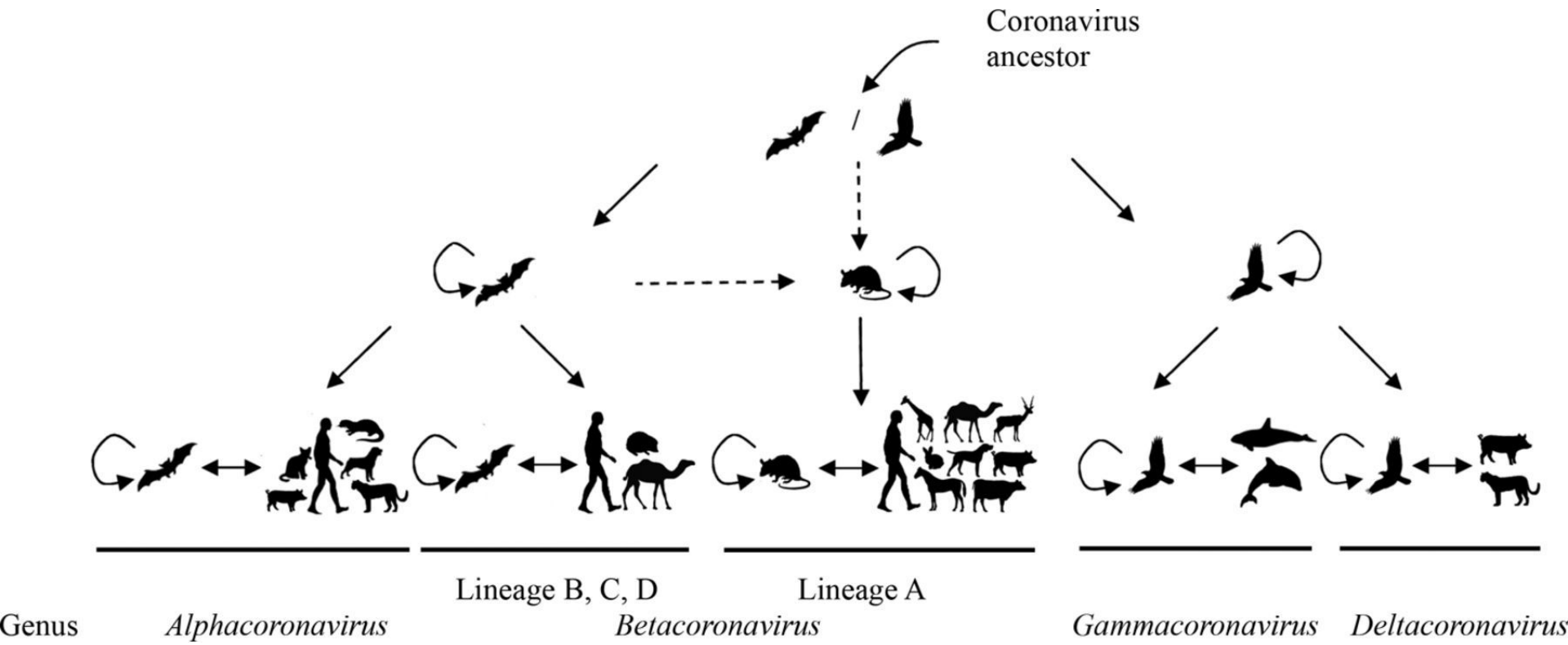
Human_SARS-CoV-2	1	ATGTTTGT	TTTTCTT	TGTTTTAT	TGCCACTAGT	CTCTAGTC	CAGTGTGT	TAA	50		
Bat_RaTG13	21545	ATGTTTGT	TTTTCTT	TGTTTTAT	TGCCACTAGT	TTCTAGTC	CAGTGTGT	TAA	21594		
...											
Human_SARS-CoV-2	2001	TGCAGGT	TATATG	CGCTAG	TATCAG	ACTCAG	ACTAAT	TCTCCT	CGGCGGG	2050	
Bat_RaTG13	23545	TGCAGGA	AATATG	CGCCAG	TATCAG	ACTCAA	ACTAAT	TCTCCT	CGGCGGG	23583	
...											
Human_SARS-CoV-2	2051	CACGTAG	TGTAG	GCTAGT	CAATCC	ATCATT	GCCTAC	ACTATG	TCACTT	TGGT	2100
Bat_RaTG13	23584	-ACGTAG	TGTGG	CCAGT	CAATCT	ATTATT	GCCTAC	ACTATG	TCACTT	TGGT	23632

Note

- “Indel” signifie “Insertion ou délétion”
- Sur base de ce résultat, la différence observée peut provenir soit d’une insertion chez un ancêtre de SARS-CoV-2, soit d’une délétion chez un ancêtre de RaTG13

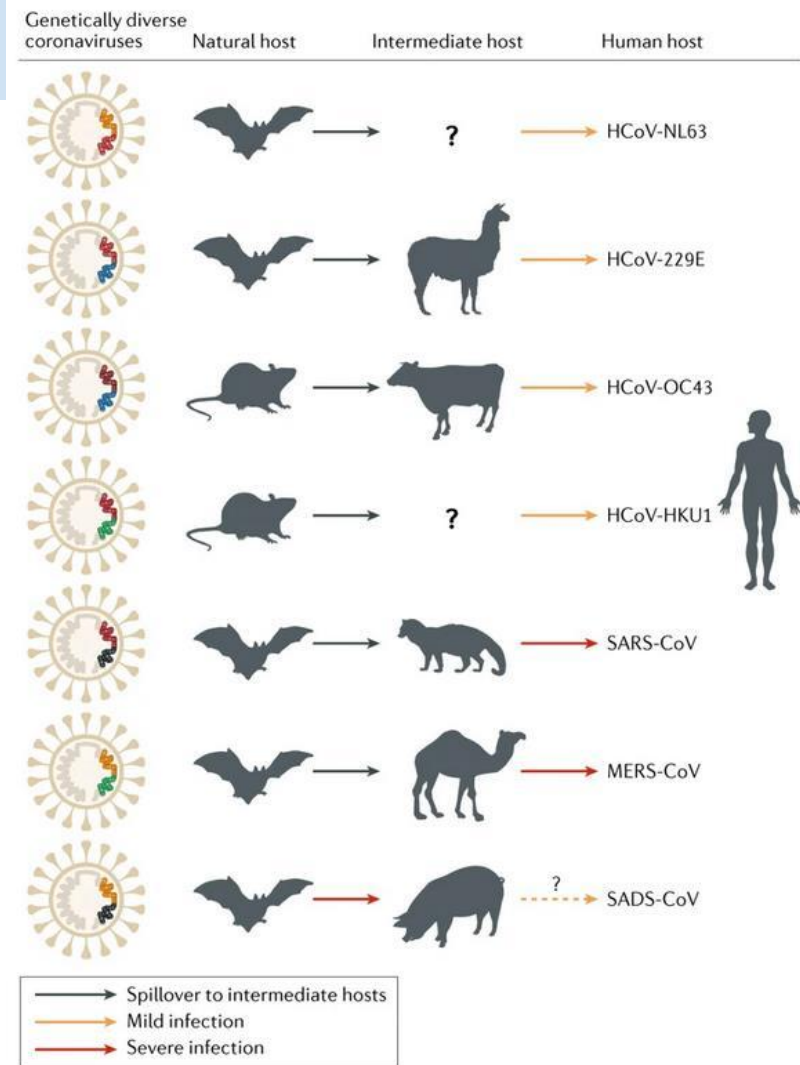
Des chauves-souris et des hommes

Cycle zoonotique des coronavirus



- Modèle du “débordement”
- Modèle de la circulation de « quasi espèces”

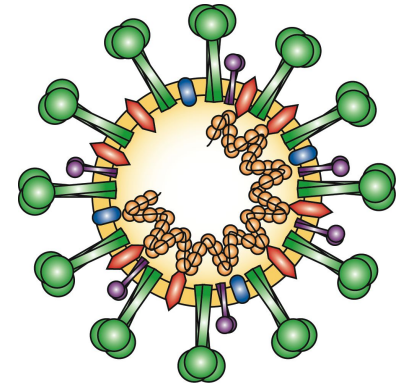
Mécanismes d'émergence des CoV humains



Origin and evolution of pathogenic coronaviruses. Jie Cui, Fang Li & Zheng-Li Shi.

Nature Reviews Microbiology volume 17, pages 181–192(2019)

Emergences de coronavirus humains



Common cold
OC43 can infect
lower respiratory
track
229E
OC43

HKU1; Pneumonia
NL63; Bronchiolitis

SARS-CoV

HKU1

NL63

MERS-CoV

2019-nCoV

1967/1970

1980

2002/2004/2005

2012

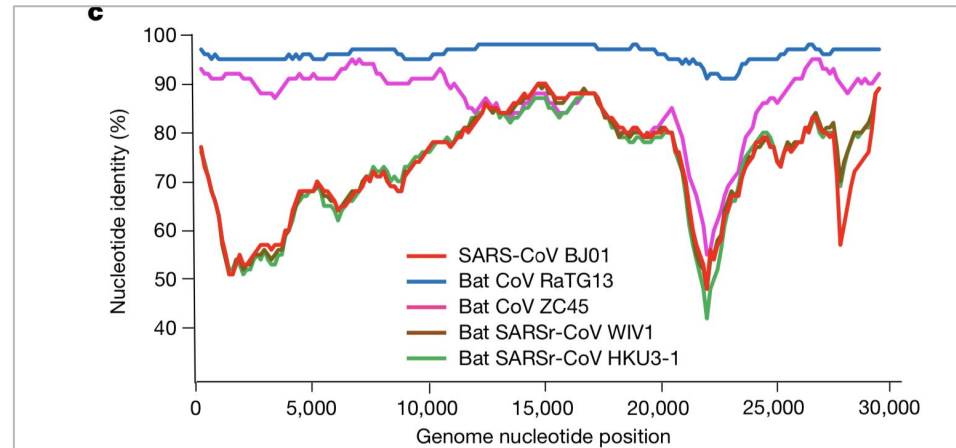
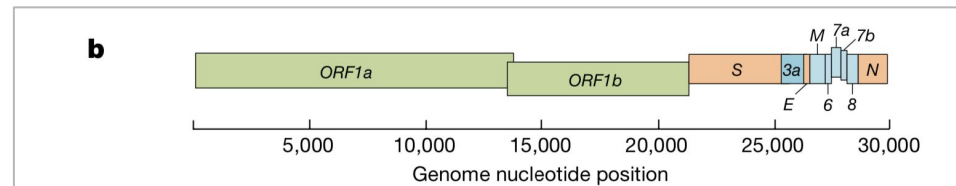
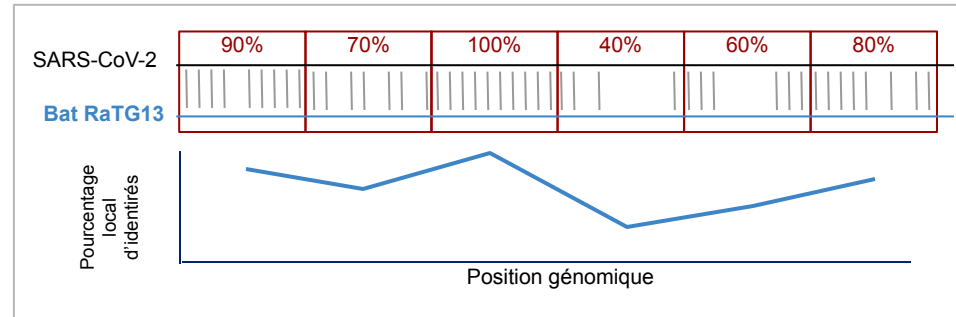
2019

OC43 genome similar to
Bovine coronavirus

SARS-CoV, MERS-CoV, 2019-nCoV
Severe respiratory disease

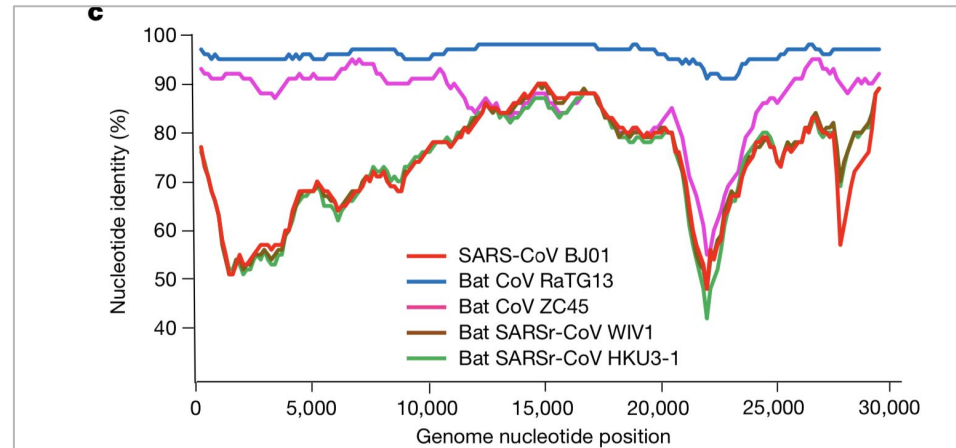
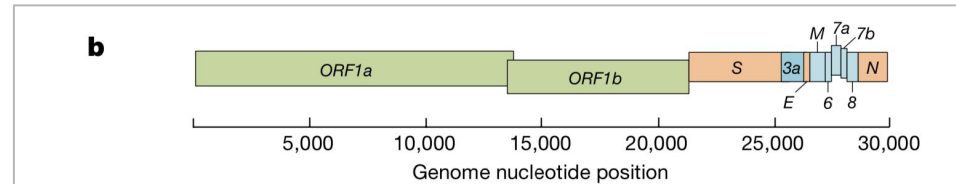
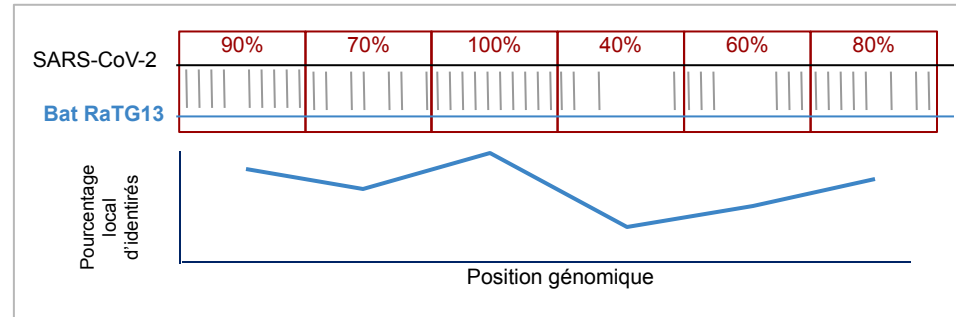
Profils de positions identiques (PPI)

- Haut: principe de calcul du PPI
 - alignement d'une paire de séquences
 - découpage de la séquence en "fenêtres"
 - calcul du pourcentage local d'identité de chaque fenêtre
 - dessin du profil de positions identiques (PPI)
- Milieu : positions des gènes de SARS-CoV-2 sur le génome
- Bas: PPI de quelques génomes de coronavirus sur celui de SARS-CoV-2
- Commentaires dans la diapo suivante



Profils de positions identiques (PPI) de génomes de coronavirus

- RaTG13 (virus de chauve-souris) est le génome le plus proche de SARS-CoV-2
- On a identifié d'autres virus de chauve-souris relativement proches de SARS-CoV-2 (Cov ZC45)
- Les virus SARS-CoV humains (pandémie 2002-2003) sont moins proches
- Pour chaque espèce, on observe des fluctuations le long du profil de PPI
- Entre 22.000 et 25.000 : chute brutale des PPI



Une origine probable: la chauve-souris

- 3 février 2020: publication du génome complet de SARS-CoV-2
- Recherche de virus similaires dans les bases de données de séquence
 - Les virus les plus proches sont des virus de chauves-souris (Bat CoV ZC45)
- Dans le même article, les auteurs décrivent un nouveau génome de chauve-souris: **RaTG13**
 - A ce jour la souche virale la plus proche de SARS-CoV-2 connue
- Figures du bas: profil de positions identiques (PPI), expliqué ci-après.

Article

A pneumonia outbreak associated with a new coronavirus of probable bat origin

<https://doi.org/10.1038/s41586-020-2012-7>

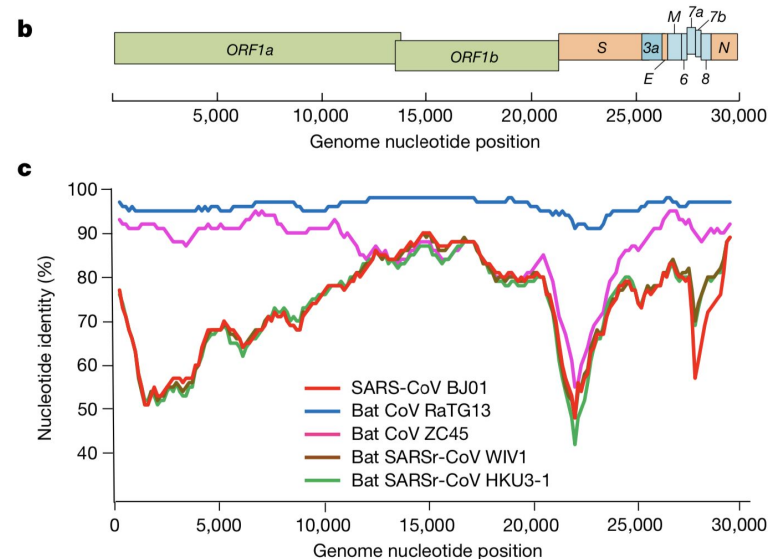
Received: 20 January 2020

Accepted: 29 January 2020

Published online: 3 February 2020

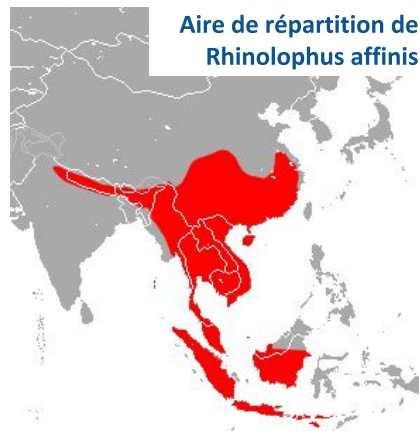
Open access

Peng Zhou^{1,5}, Xing-Lou Yang^{1,5}, Xian-Guang Wang^{2,5}, Ben Hu¹, Lei Zhang¹, Wei Zhang¹, Hao-Rui Si^{1,5}, Yan Zhu¹, Bei Li¹, Chao-Lin Huang², Hui-Dong Chen², Jing Chen^{1,3}, Yun Luo^{1,3}, Hua Guo^{1,3}, Ren-Di Jiang^{1,3}, Mei-Qin Liu^{1,5}, Ying Chen^{1,5}, Xu-Rui Shen^{1,3}, Xi Wang^{1,3}, Xiao-Shuang Zheng^{1,3}, Kai Zhao^{1,3}, Qian-Jiao Chen¹, Fei Deng¹, Lin-Lin Liu⁴, Bing Yan¹, Fa-Xian Zhan⁴, Yan-Yi Wang¹, Geng-Fu Xiao¹ & Zheng-Li Shi^{1,5*}



De Yunnan à Wuhan

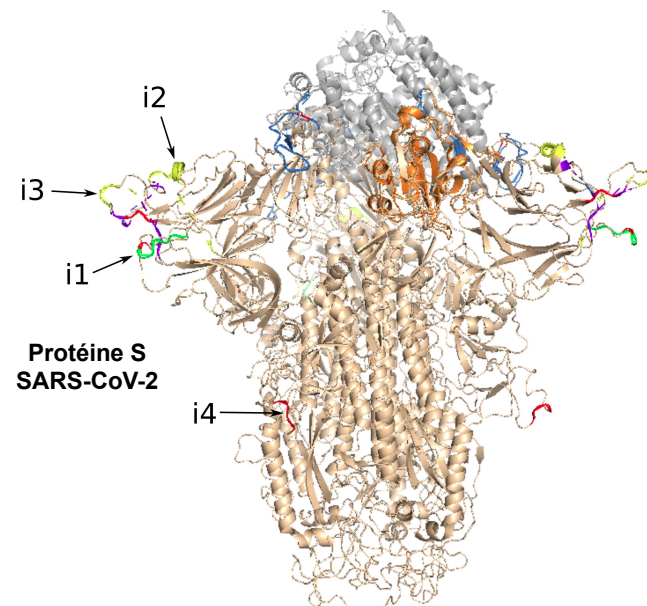
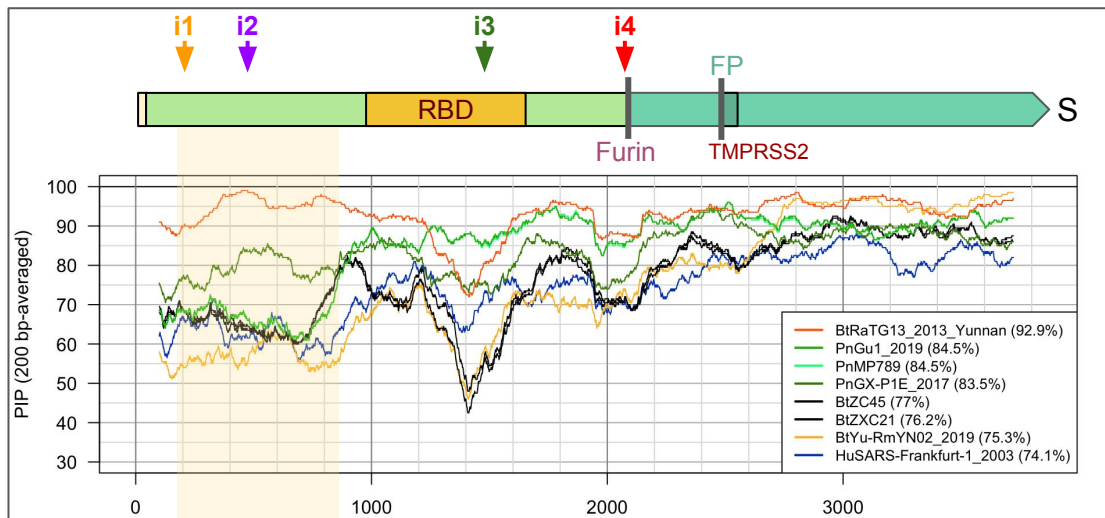
- 2013
 - ❑ Pneumonie atypique chez 6 mineurs dans la province de Yunnan, 3 décès
 - ❑ Plusieurs pistes sont évoquées (levures, virus) dont un coronavirus
 - ❑ Collecte d'échantillons de chauves-souris dans la mine
- 2016 : publication d'un fragment de séquence (360 nucléotides, 1% du génome) de virus de *Rhinolophus affinis*
- 2018 : dépôt des fragments de séquençage (reads) dans une base de données, pour ~90% du génome
- 2020 :
 - ❑ publication de la séquence complète du génome viral, sous l'identifiant "**BatCoV RaTG13**".
 - ❑ Ce génome est le plus proche connu de celui de SARS-CoV-2 (96,2% nucléotides identiques).
- Note: l'implication de coronavirus dans le décès des mineurs fait encore l'objet de débats



Insertions dans les séquences du gène S

Quatre insertions dans le gène S de SARS-CoV-2

- Les flèches indiquent la position des 4 insertions sur le gène S (gauche) et sur la protéine spicule (droite).
- Les 3 premières sont situées à l'extérieur de la protéine, dans des régions "exposées".



- Le site de clivage par la furine qu'on observe dans la protéine spicule de SARS-CoV-2 ne se trouve dans aucun autre coronavirus.
- Il résulte de l'insertion de 12 nucléotide à un endroit particulier du le gène S.

≡ EL PAÍS

CORONAVIRUS

ccu cgg cgg gca

The 12 letters that changed the world

The genome of the new coronavirus harbors a short sequence suspected of being the main culprit of its uniquely infectious and aggressive nature



MANUEL ANSEDE  | ARTUR GALOCHA | MARIANO ZAFRA 

19 MAY 2020 - 18:25 CEST

Des insertions bizarres?

- Figure from Pradhan et al (2020), initially published on bioRxiv and retracted.
- The “multiple alignment” is actually a pairwise alignment + a consensus.
- The gaps obtained from a multiple alignment overlap with these ones, but they start and end at different positions.
- It is precisely because they did not do a multiple alignment that they did not realize that 3 of these insertions were not unique to SARS-CoV-2.

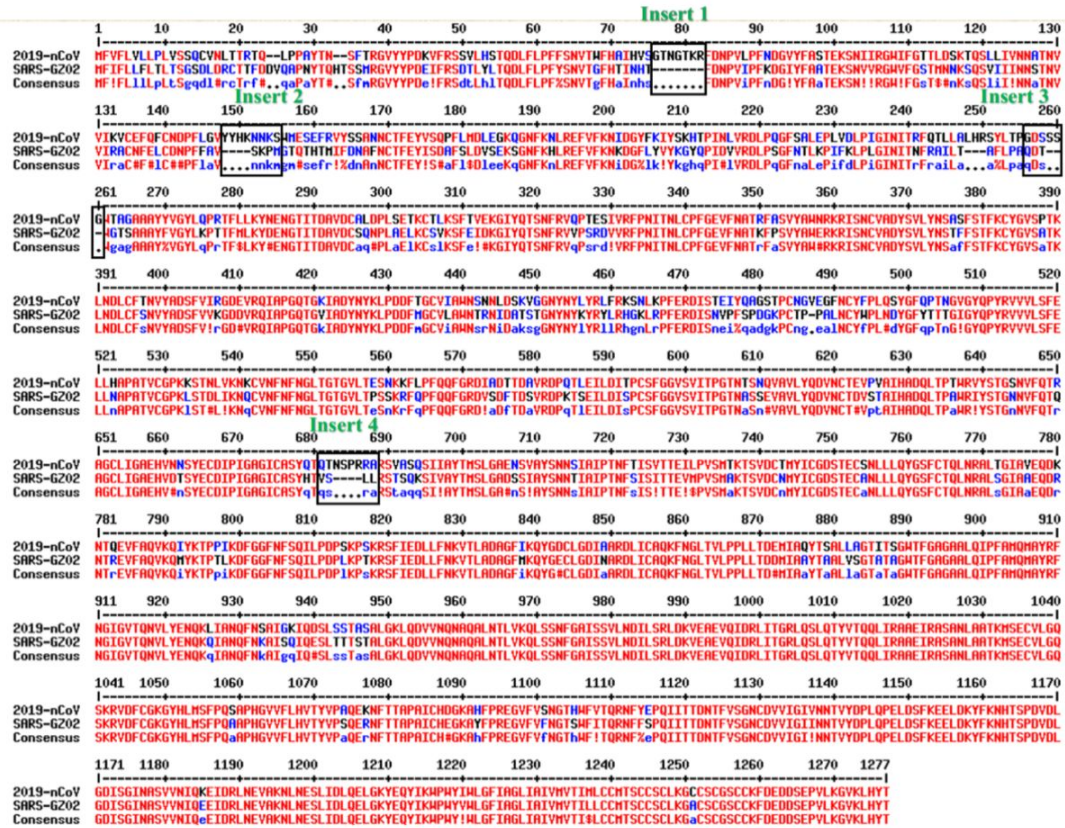


Figure 2: Multiple sequence alignment between spike proteins of 2019-nCoV and SARS. The sequences of spike proteins of 2019-nCoV (Wuhan-HU-1, Accession NC_045512) and of SARS CoV (GZ02, Accession AY390556) were aligned using MultiAlin software. The sites of difference are highlighted in boxes.

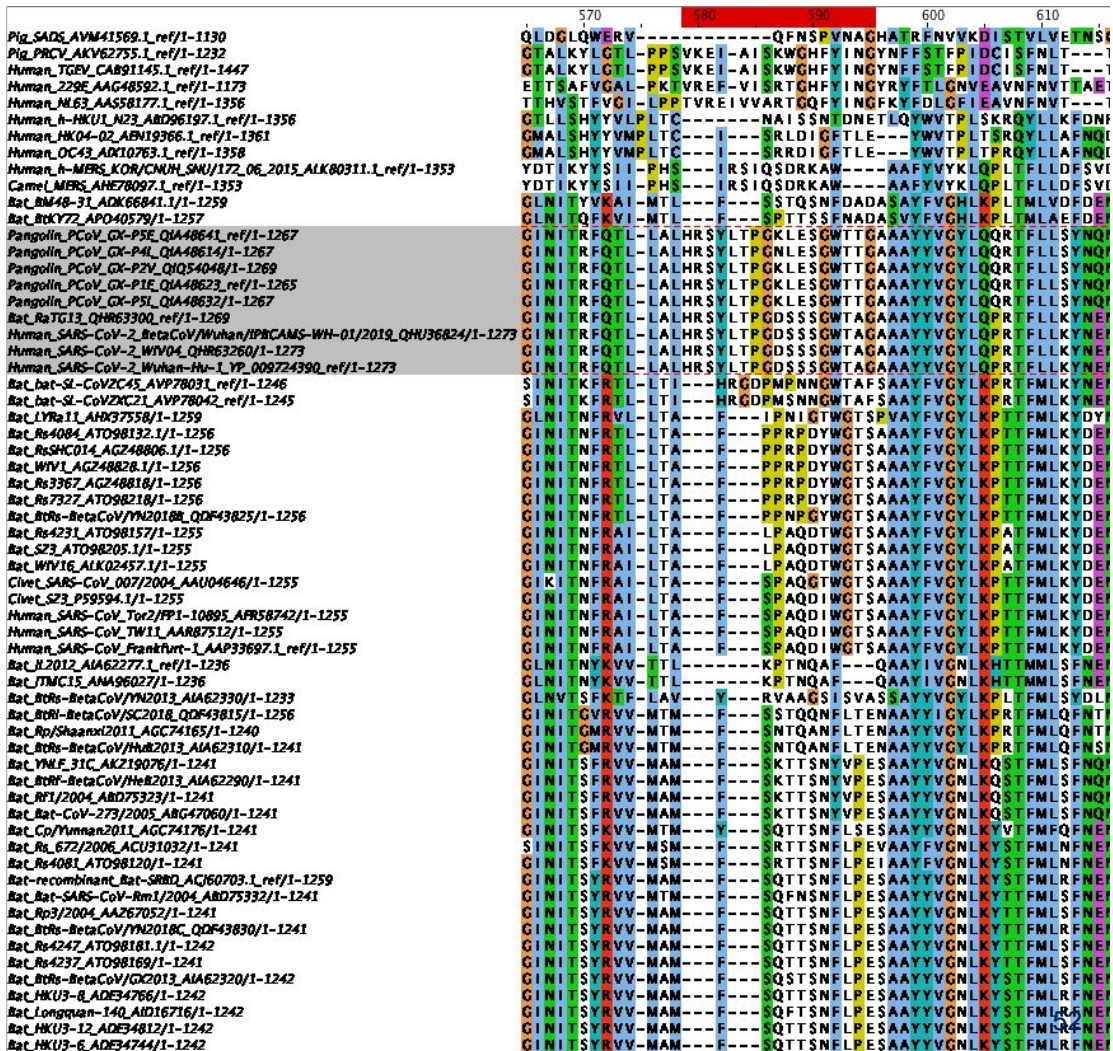
Insertion partagée entre tous les virus du groupe CoV-2

- Position: 153-158 de SARS-CoV-2
- Cette insertion se trouve chez les virus de pangolin + plusieurs chauve-souris
- Les résidus sont identiques entre SARS-CoV-2 et la souche RaTG13 de chauve-souris (la plus proche de SARS-CoV-2)
- Par contre elle présente 3 substitutions entre les souches de pangolin et SARS-CoV-2.

	420	430	440	450	460	470	480
Pig_SADS_AYM41569.1_ref1-1130	NV	VY	VY	LRCR	---	WW	---
Pig_PRCV_AKV62755.1_ref1-1232	SG	KL	VY	YKQF	---	LW	---
Human_TGVB_CAB91145.1_ref1-1447	SG	KL	VY	YKQF	---	LW	---
Human_229E_AAG48592.1_ref1-1173	SF	QPL	LL	LNCL	---	W	---
Human_ML63_AAS58177.1_ref1-1356	LY	QF	LR	LC	---	W	---
Human_h-HKU1_N23_ABD96197.1_ref1-1356	NG	VY	EI	TACQ	---	Y	---
Human_HK04-02_ABN19366.1_ref1-1361	QGL	LE	YS	VY	---	V	---
Human_OC43_ADX10763.1_ref1-1358	QGL	LE	YS	VY	---	V	---
Human_h-MERS_KOR/CMNH_SNU/172_06_2015_ALK80311.1_ref1-1353	I	V	LL	PD	---	CG	---
CameL_MERS_AHE78097.1_ref1-1353	I	V	LL	PD	---	CG	---
Rat_RM48-31_ADK66841.1/1-1259	G	T	H	I	---	D	---
Rat_RoKY22_APO40579/1-1257	G	T	H	I	---	D	---
Pangolin_PCoV_GX-PSE_QIA48641.1-1267	A	N	V	V	---	K	---
Pangolin_PCoV_GX-P4L_QIA48614/1-1267	A	N	V	V	---	K	---
Pangolin_PCoV_GX-P2V_QHQ54048/1-1269	A	N	V	V	---	K	---
Pangolin_PCoV_GX-P1E_QIA48623.1-1265	A	N	V	V	---	K	---
Pangolin_PCoV_GX-P5L_QIA48632/1-1267	A	N	V	V	---	K	---
Rat_RaTG13_QHR63300.ref1-1269	A	N	V	V	---	K	---
Human_SARS-CoV-2_BetaCoV/Huhan/IPBCAMS-WH-01/2019_QHU36824/1-1273	A	N	V	V	---	K	---
Human_SARS-CoV-2_WIV04_QHR63260/1-1273	A	N	V	V	---	K	---
Human_SARS-CoV-2_Wuhan-Hu-1_YP_009724390.ref1-1273	A	N	V	V	---	K	---
Rat_bat-SL-CoVZC45_AVIP7803.ref1-1246	A	N	V	V	---	K	---
Rat_bat-SL-CoVZMC12_AWP78042.ref1-1245	A	N	V	V	---	K	---
Rat_LYba11_AHK37558/1-1259	S	N	V	V	---	L	---
Rat_Ro4084_ATO98132.1/1-1256	S	N	V	V	---	L	---
Rat_RoShC014_AG248806.1/1-1256	S	N	V	V	---	L	---
Rat_WIV1_AG248828.1/1-1256	S	N	V	V	---	L	---
Rat_Ro3367_AG248818/1-1256	S	N	V	V	---	L	---
Rat_Ro7927_ATO98218/1-1256	S	N	V	V	---	L	---
Rat_RoRc-BetaCoV/YN2010R_QDQ43825/1-1256	S	N	V	V	---	L	---
Rat_Ro4231_ATO98157/1-1255	S	N	V	V	---	L	---
Rat_S23_ATO98205.1/1-1255	S	N	V	V	---	L	---
Rat_WIV16_AIK02457.1/1-1255	S	N	V	V	---	L	---
Civet_SARS-CoV_007/2004_AAU04646/1-1255	S	N	V	V	---	L	---
Civet_S23_P59594.1/1-1255	S	N	V	V	---	L	---
Human_SARS-CoV_Tor2/FP1-10895_APR58742/1-1255	S	N	V	V	---	L	---
Human_SARS-CoV_TW11_AAR87512/1-1255	S	N	V	V	---	L	---
Human_SARS-CoV_Frankfurt-LAAP33697.1.ref1-1255	S	N	V	V	---	L	---
Ra_T.2012_AIA62277.1.ref1-1236	G	S	A	I	---	E	---
Rat_TMC15_ANA96027/1-1236	G	S	A	I	---	E	---
Rat_RoRc-BetaCoV/YN2013_AIA62330/1-1233	S	N	V	V	---	L	---
Rat_RoRc-BetaCoV/SC201R_QDQ43815/1-1256	S	N	V	V	---	L	---
Rat_Ro/Shaanxt2011_AGC74163/1-1240	S	N	V	V	---	L	---
Rat_RoRc-BetaCoV/HuH2013_AIA62310/1-1241	S	N	V	V	---	L	---
Rat_YNF_31C_AK219076/1-1241	S	N	V	V	---	L	---
Rat_RoRc-BetaCoV/HeB2013_AIA62290/1-1241	S	N	V	V	---	L	---
Rat_RY1/2004_ABD75323/1-1241	S	N	V	V	---	L	---
Rat_Bat-CoV-273/2005_ABG47060/1-1241	S	N	V	V	---	L	---
Rat_Co/Yunnan2011_AGC74176/1-1241	S	N	V	V	---	L	---
Rat_Ro_672/2006_ACU31032/1-1241	S	N	V	V	---	L	---
Rat_Ro4081_ATO98120/1-1241	S	N	V	V	---	L	---
Rat-recombinant_Bat-SRBD_AJ60703.1_ref1-1259	S	N	V	V	---	L	---
Rat_Bat-SARS-CoV-Rm1/2004_ABD75332/1-1241	S	N	V	V	---	L	---
Rat_Rp3/2004_AA267052/1-1241	S	N	V	V	---	L	---
Rat_RoRc-BetaCoV/YN2018C_QDQ43830/1-1241	S	N	V	V	---	L	---
Rat_Ro4247_ATO98181.1/1-1242	S	N	V	V	---	L	---
Rat_Ro4237_ATO98169/1-1241	S	N	V	V	---	L	---
Rat_RoRc-BetaCoV/GZ2013_AIA62320/1-1242	S	N	V	V	---	L	---
Rat_HKU3-8_ADE34766/1-1242	S	N	V	V	---	L	---
Rat_Longquan-14Q_AID16716/1-1242	S	N	V	V	---	L	---
Rat_HKU3-12_ADE34812/1-1242	S	N	V	V	---	L	---
Rat_HKU3-6_ADE34744/1-1242	S	N	V	V	---	L	---

Insertion partagée par la majorité des virus du groupe CoV-2

- Position: 245-251 de SARS-CoV-2
- Cette insertion se trouve chez les virus de pangolin + la souche RaTG13 de chauve-souris
- Elle est cependant absente de 2 souches de chauves-souris appartenant au groupe CoV-2 : CoVZC45 et CoVZXC21
- Au-delà de l'insertion on trouve un bloc conservé (jusqu'à la position 595 de l'alignement).
- Au sein de ce bloc, une paire de résidus distingue les pangolins du groupe SARS2 + Bat RaTG13 .



Insertion i3

- Position
 - 470-486 de SARS-CoV-2
 - 855-872 sur l'alignement
- Commune au groupe pangolin + Bat_RaTG13 + SARS-CoV-2
- 2 substitutions uniques à Bat_RaTG13

```

%ig_SARS_AVN41569_1_ref/1-1130
%ig_PRCV_AKV62755_1_ref/1-1232
human_TGEV_CAB91145_1_ref/1-1447
human_229E_AAC48592_1_ref/1-1173
human_NL63_AAS58177_1_ref/1-1356
human_h-PKU1_N23_ABD96197_1_ref/1-1356
human_HK04-02_ABN19366_1_ref/1-1361
human_OC43_ADI0763_1_ref/1-1358
Human_h-MERS_KOR/CNH/SMU/172_06_2015_ALK80311.1_ref/1-1353
Camel_MERS_AHE78097_1_ref/1-1353
Bat_BM48-31_ADK66041.1/1-1259
Bat_BatKY72_APO40579/1-1257
Pangolin_PCoV_GX-PSE_QIA48641_ref/1-1267
Pangolin_PCoV_GX-P4I_QIA48614/1-1267
Pangolin_PCoV_GX-P2V_QIQ54048/1-1269
Pangolin_PCoV_GX-P1E_QIA48623_ref/1-1265
Pangolin_PCoV_GX-PSL_QIA48632/1-1267
Bat_RaTG13_QHR63300_ref/1-1269
Human_SARS-CoV-2_BetaCoV/Wuhan/HPBCAMS-WH-01/2019_QHU36824/1-1273
Human_SARS-CoV-2_WIV04_QHR63260/1-1273
Human_SARS-CoV-2_Wuhan-Hu-1_YP_009724390_ref/1-1273
Bat_bat-SL-CoVZC45_AVP78031_ref/1-1246
Bat_bat-SL-CoVZXC21_AVP78042_ref/1-1245
Bat_LYRa11_AHX97558/1-1259
Bat_Rs4084_ATO98132.1/1-1256
Bat_RsSHC014_AG248806.1/1-1256
Bat_WIV1_AG24882B.1/1-1256
Bat_Rs336Z_AG24881B/1-1256
Bat_Rs722Z_ATO98218/1-1256
Bat_BtRs-BetaCoV/YN2018B_QDF43825/1-1256
Bat_Rs423L_ATO98157/1-1255
Bat_SZ3_ATO98205.1/1-1255
Bat_WIV16_ALK02457.1/1-1255
Clivec_SARS-CoV_007/2004_AAU04646/1-1255
Clivec_SZ3_P59594.1/1-1255
Human_SARS-CoV_Tor2/FP1-10895_AFR58742/1-1255
Human_SARS-CoV_TW11_AAR87512/1-1255
Human_SARS-CoV_Frankfurt_1_AAP33697_1_ref/1-1255
Bat_H2012_AIA62277_1_ref/1-1236
Bat_JTMC15_AIA96027/1-1236
Bat_BtRs-BetaCoV/YN2013_AIA62330/1-1233
Bat_BtRI-BetaCoV/SC2018_QDF43815/1-1256
Bat_Rp/Shaanxi2011_AGC74165/1-1240
Bat_BtRs-BetaCoV/Han2013_AIA62310/1-1241
Bat_YNF_E_31C_AK219076/1-1241
Bat_BtRI-BetaCoV/Han2013_AIA62290/1-1241
Bat_Rf1/2004_ABD75323/1-1241
Bat_Bat-CoV-273/2005_ABG47060/1-1241
Bat_Cp/Myman2011_AGC74176/1-1241
Bat_Rs_672/2006_ACU31032/1-1241
Bat_Rs408L_ATO98120/1-1241
Bat-recombinant_Bat-SRB0_ACJ60703_1_ref/1-1259
Bat_Bat-SARS-CoV-Rm1/2004_ABD75332/1-1241
Bat_Rp3/2004_AAZ67052/1-1241
Bat_BtRs-BetaCoV/YN2018C_QDF43830/1-1241
Bat_Rs424Z_ATO98181.1/1-1242
Bat_Rs423Z_ATO98169/1-1241
Bat_BtRs-BetaCoV/GX2013_AIA62320/1-1242
Bat_HKU3-B_ADE34766/1-1242
Bat_Longman-140_AID16716/1-1242
Bat_HKU3-12_ADE34812/1-1242
Bat_HKU3-6_ADE34744/1-1242

```

```

--VLRVGRG---KAVNRITVTRYLPKYP-----
DAIVIKITG---CPFSFDKLNYY-----LTFNKFCFLSS--VGAN
DAVAVIKITG---CPFSFDKLNYY-----LTFNKFCFLSIPVGA--N
SASINTG---NCPFSF---GKYNFF-----VKFGSVCFKLDIPGG--
IYLLKSG---ICPFSFKLNNF-----IKFKTICFTVEVPSS--
RFGNFNFNLSHSHSVVSRVCFSVNNTGCPKAPKSFASSCKSKHPPSASCFPIGTNYSRCESTTVLDHTDWC
FVFPQTGVFNTHVSVAQHCFFKAPKNCPC---KLNIGSCPGKNGGICTCPAGTNLTCDNL--
FVFPQTGVFNTHVSVAQHCFFKAPKNCPC---SSCPGKNGGICTCPAGTNLTCDNL--
VPHNLTTITKPLKYSYINKCSRLLSDD-----RTEVPQLVLMANQYSPCVS--VTPS
VPHNLTTITKPLKYSYINKCSRLLSDD-----RTEVPQLVLMANQYSPCVS--VTPS
IINSVSDSN---EFFFYRFRHKGIKPY-----GKDLNVLFFNPSGGTCSA-EGLN
INSVSDSKGN---NFYYRFLRHGRIKPY-----ERDINSVLYNSAGGTCCSSIISGLG
SVKQDALTDGNYGGLYRLLFRKSKLKP-----ERDIDSTEIYQAGSTPFCNGQVGLN
SVKQDALTDGNYGGLYRLLFRKSKLKP-----ERDIDSTEIYQAGSTPFCNGQVGLN
SVKQDALTDGNYGGLYRLLFRKSKLKP-----ERDIDSTEIYQAGSTPFCNGQVGLN
SVKQDALTDGNYGGLYRLLFRKSKLKP-----ERDIDSTEIYQAGSTPFCNGQVGLN
SVKQDALTDGNYGGLYRLLFRKSKLKP-----ERDIDSTEIYQAGSTPFCNGQVGLN
SVKQDALTDGNYGGLYRLLFRKSKLKP-----ERDIDSTEIYQAGSTPFCNGQVGLN
SVKQDALTDGNYGGLYRLLFRKSKLKP-----ERDIDSTEIYQAGSTPFCNGQVGLN
SKHIDAKEGGFNLYRLLFRKANLKP-----ERDIDSTEIYQAGSKPTNGVQLN
SNNLDSKVGGNYNYLFRKSNLKP-----ERDIDSTEIYQAGSTPFCNGVEGFN
SNNLDSKVGGNYNYLFRKSNLKP-----ERDIDSTEIYQAGSTPFCNGVEGFN
SNNLDSKVGGNYNYLFRKSNLKP-----ERDIDSTEIYQAGSTPFCNGVEGFN
SKAQDVG---SYFYRSSRSKSLKP-----ERDLS--EE--N
TAKQDVG---HYFYRSSRSKSLKP-----ERDLS--DE--N
TRNIDATSSGNFNKYRSLRHGKLRP-----ERDINSVFP--PDGKPCFP--PAFN
INSKDSSTSGNMYLRYRWRKSKLNP-----ERDLSNDIY--PGQSCSA-VGFN
INSKDSSTSGNMYLRYRWRKSKLNP-----ERDLSNDIY--PGQSCSA-VGFN
TRNIDATQTGNMYNYRSLRHGKLRP-----ERDINSVFP--PDGKPCFP--PAFN
TRNIDATQTGNMYNYRSLRHGKLRP-----ERDINSVFP--PDGKPCFP--PAFN
TRNIDATQTGNMYNYRSLRHGKLRP-----ERDINSVFP--PDGKPCFP--PAFN
TRNIDATQTGNMYNYRSLRHGKLRP-----ERDINSVFP--PDGKPCFP--PAFN
TRNIDATSGNYNKKYRSLRHGKLRP-----ERDINSVFP--PDGKPCFP--PAFN
TRNIDATSGNYNKKYRSLRHGKLRP-----ERDINSVFP--PDGKPCFP--PAFN
TRNIDATSGNYNKKYRSLRHGKLRP-----ERDINSVFP--PDGKPCFP--PAFN
TRNIDATSGNYNKKYRSLRHGKLRP-----ERDINSVFP--PDGKPCFP--PAFN
TRNIDATSGNYNKKYRSLRHGKLRP-----ERDINSVFP--PDGKPCFP--PAFN
TRNIDATSGNYNKKYRSLRHGKLRP-----ERDINSVFP--PDGKPCFP--PAFN
TAKQDVG---SYFYRSSRSKSLKP-----ERDLS--EE--N
TAKQDVG---SYFYRSSRSKSLKP-----ERDLS--EE--N
TAKQDVG---SYFYRSSRSKSLKP-----ERDLS--DE--N
TAKQDVG---SYFYRSSRSKSLKP-----ERDLS--DDG--N
TAKQDVG---QYYRSSHRSKSLKP-----ERDLS--DE--N
TAKQDVG---YYYRSSHRSKSLKP-----ERDLS--DDG--N
TAKQDVG---SYFYRSSRSKSLKP-----ERDLS--EE--N
TAKQDVG---SYFYRSSRSKSLKP-----ERDLS--EE--N
TAKQDVG---SYFYRSSRSKSLKP-----ERDLS--EE--N
TAKQDVG---SYFYRSSRSKSLKP-----ERDLS--EE--N
TAKQDVG---SYFYRSSRSKSLKP-----ERDLS--VE--E
TAKQDVG---QYYRSSHRSKSLKP-----ERDLS--DE--N
TAKQDVG---QYYRSSHRSKSLKP-----ERDLS--DE--N
TAKQDVG---QYYRSSHRSKSLKP-----ERDLS--DE--N
TRNIDATSGNYNKKYRSLRHGKLRP-----ERDINSVFP--PDGKPCFP--PAFN
TAKQDVG---QYYRSSHRSKSLKP-----ERDLS--DE--N
TAKQDVG---QYYRSSHRSKSLKP-----ERDLS--DE--N
TAKQDVG---QYYRSSHRSKSLKP-----ERDLS--DE--N
TAKQDVG---QYYRSSHRSKSLKP-----ERDLS--DDG--N
TAKQDVG---NYYRSSHRSKSLKP-----ERDLS--DDG--N
TAKQDVG---NYYRSSHRSKSLKP-----ERDLS--DDG--N
TAKQDVG---NYYRSSHRSKSLKP-----ERDLS--DDG--N
TAKQDVG---NYYRSSHRSKSLKP-----ERDLS--DDG--N
TAKQDVG---NYYRSSHRSKSLKP-----ERDLS--DDG--N
TAKQDVG---NYYRSSHRSKSLKP-----ERDLS--DDG--N

```


Insertion d'un site Furine (i4)

- Positions : 1181-1184 de l'alignement
- On trouve chez SARS-CoV-2 un site unique SPRRAR, qui résulte d'une insertion SPRR et d'une substitution L -> A
- La séquence PRRA correspond au motif reconnu par la furine (protéase).
- Cette insertion est à l'origine du site de clivage responsable du caractère particulièrement virulent de SARS-CoV-2

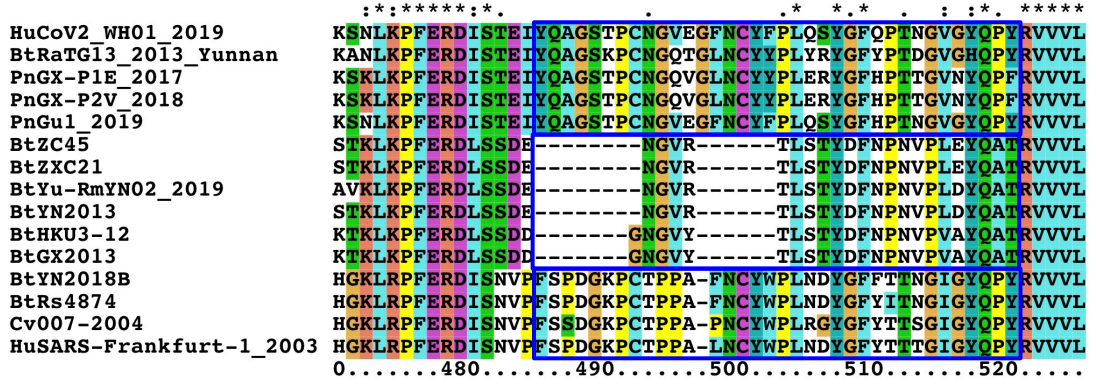
Cm_MERS_AHE78097.1_ref	SLCALP-DTPST----LTPRSVRSV	20
Hu_MERS_172-06_2015_ALK80311.1_ref	SLCALP-DTPST----LTPRSVRSV	20
Bt_BM48-31_ADK66841.1	GICAKYTNVSSST----LVRSGGHSI	21
Bt_BtKY72_APO40579	GICAKF-GSDKI-----RMGOESI	18
BtYu-RmYN02_2019_S-gene_21544-25227_1	GVCASY-NSPAA-----RVGTNSI	18
Bt_LYRa11_AHX37558	GICASY-HTASL----LRNTDQKSI	20
Bt_YN2018B_QDF43825	GICASY-HTVSS----LRSTSQKSI	20
Bt_Rs4874_ATO98205.1	GICASY-HTVSS----LRSTSQKSI	20
Cv_007-2004_AAU04646	GICASY-HTVSS----LRSTSQKSI	20
Hu_SARS-Frankfurt-1_2003_AAP33697.1_ref	GICASY-HTVSL----LRSTSQKSI	20
Bt_rec-SARS_2008_ACJ60694.1_ref	GICASY-HTVSL----LRSTSQKSI	20
Bt_ZC45_AVP78031_ref	GICASY-HTASI----LRSTSQKAI	20
Bt_ZXC21_AVP78042_ref	GICASY-HTASI----LRSTGQKAI	20
PnGu1_2019_S-gene_21541-25338_1	GICASY-QTQTN----SRSVSSQAI	20
Pn_GX-P1E_2017_QIA48623_ref	GICASY-HSMSS----LRSVNORSI	20
Pn_GX-P2V_2018_QIQ54048	GICASY-HSMSS----FRSVNORSI	20
Bt_RaTG13_2013_Yunnan_QHR63300_ref	GICASY-QTQTN----SRSVASQSI	20
Hu_CoV2_WH01_2019_QHU36824_ref	GICASY-QTQTN SPRR ARSVASQSI	24
Bt_JL2012_AIA62277.1_ref	GICASY-HTASL----LRSTGQKSI	20
Bt_YN2013_AIA62330	GICASY-HTAST----LRSIGQKSI	20
Bt_Rp-Shaanxi2011_AGC74165	GICASY-HTASV----LRSTGQKSI	20
Bt_SC2018_QDF43815	GICASY-HTAST----LRSTGQKSI	20
Bt_YNLF_31C_AKZ19076	GICASY-HTASV----LRSTGQKSI	20
Bt_Cp-Yun_2011_AGC74176	GICASY-HTASL----LRNTGQKSI	20
Bt_Rs_672-2006_ACU31032	GICASY-HTAST----LRSVGQKSI	20
Bt_Rm1/2004_ABD75332	GICASY-HTASV----LRSTGQKSI	20
Bt_YN2018C_QDF43830	GICASY-HTAST----LRSVGQKSI	20
Bt_Rp3-2004_AAZ67052	GICASY-HTAST----LRSVGQKSI	20
Bt_GX2013_AIA62320	GICASY-HTASV----LRSTGQKSI	20
Bt_HKU3-12_ADE34812_ref	GICASY-HTASV----LRSTGQKSI	20

0.....790.....800....



Un site recombinant?

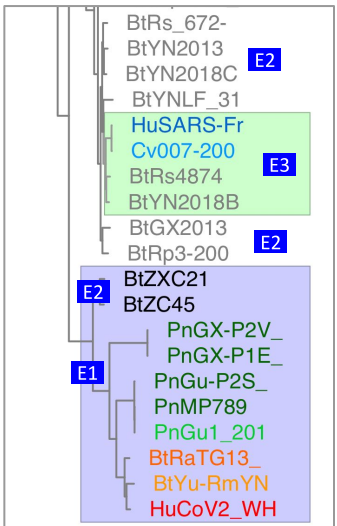
- L'interprétation de ce site est plus complexe.
- Il existe clairement trois groupes de séquences.
- Ceux-ci s'étendent au-delà des deux indels.
- L'arbre construit à partir de cette région est peu robuste, et incohérent avec celui des génomes.
- La répartition des sous-blocs de séquences est plus cohérente avec l'arbre des espèces.
- Cette région a échappé à Pradhan et al. parce qu'ils ont réalisé un alignement par paire SARS-CoV-2 vs SARS plutôt qu'un alignement multiple.



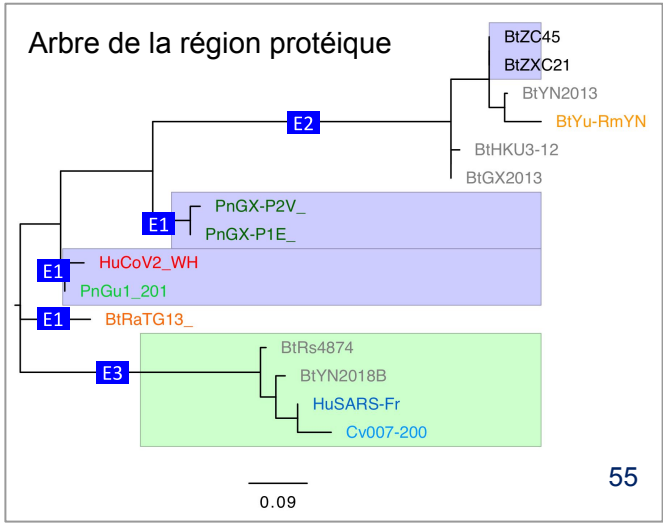
E1
E2
E3



Arbre des génomes



Arbre de la région protéique

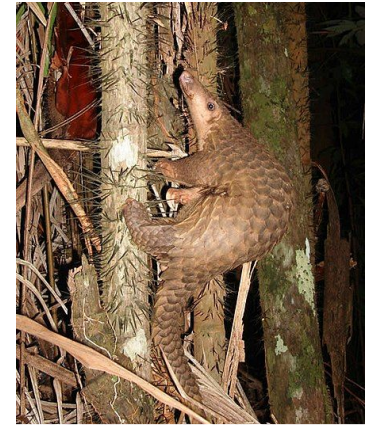


Des chauves-souris et des hommes ... et des pangolins ?

Pangolins

- Mode de vie
- Aire géographique
- Migration
- Contacts avec les chauves-souris
- Contacts avec l'humain

Répartition géographique de différentes espèces de pangolin



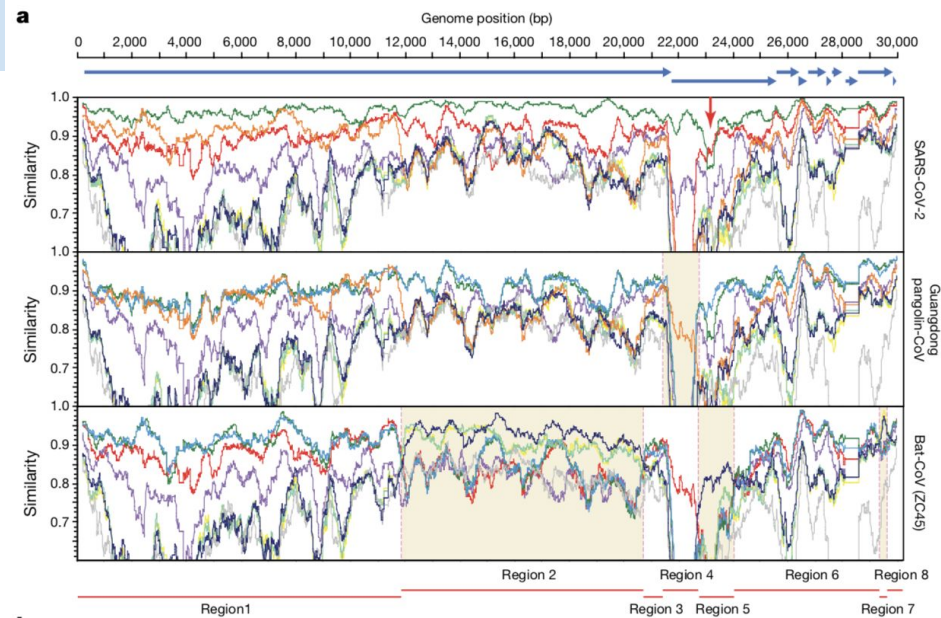
Manis javanica
(Pangolin malais)

Aire de répartition de *Manis javanica*



... et des pangolins ?

- Lam et collègues comparent le génome de SARS-CoV-2 à des génomes de virus isolés à partir de pangolins.
- Ces virus sont globalement plus éloignés de SARS-CoV-2 que ceux de chauve-souris.
- Cependant, on observe une identité plus élevée dans la région particulière où les PPI des autres coronavirus s'affaissent.
- Ceci suggère la possibilité d'une recombinaison entre des virus de chauve-souris et de pangolin.



b

- Guangdong pangolin-CoV
- Guangxi pangolin-CoV
- SARS-CoV-2
- Bat-CoV (RaTG13)
- Bat-CoV (ZXC21, ZC45)
- Bat-SL-CoV (273, Rs3367)
- Bat-SL-CoV (HKU3, Rf1, 273)
- SARS-CoV
- Bat-CoV from Kenya and Bulgaria

Et si les pangolins n'y étaient pour rien ?

L'hypothèse du pangolin est fortement remise en cause pour plusieurs raisons.

- Les génomes des virus de pangolin les plus proches dont on dispose sont plus éloignés de SARS-CoV-2 que ceux des génomes de chauves-souris.
- La région où la similarité est la plus forte correspondent au gène S, qui code pour la protéine spicule. Or la protéine spicule de ces virus de pangolin n'est pas capable d'adhérer au récepteur de cellules humaines.

franceinfo:

vidéos

radio

jt

magazines

DIRECTV

DIRECT RADIO



le billet sciences

Anne-Laure Barral

◆ / replay_radio / Le billet sciences

Covid-19 : et si le pangolin n'y était pour rien ?

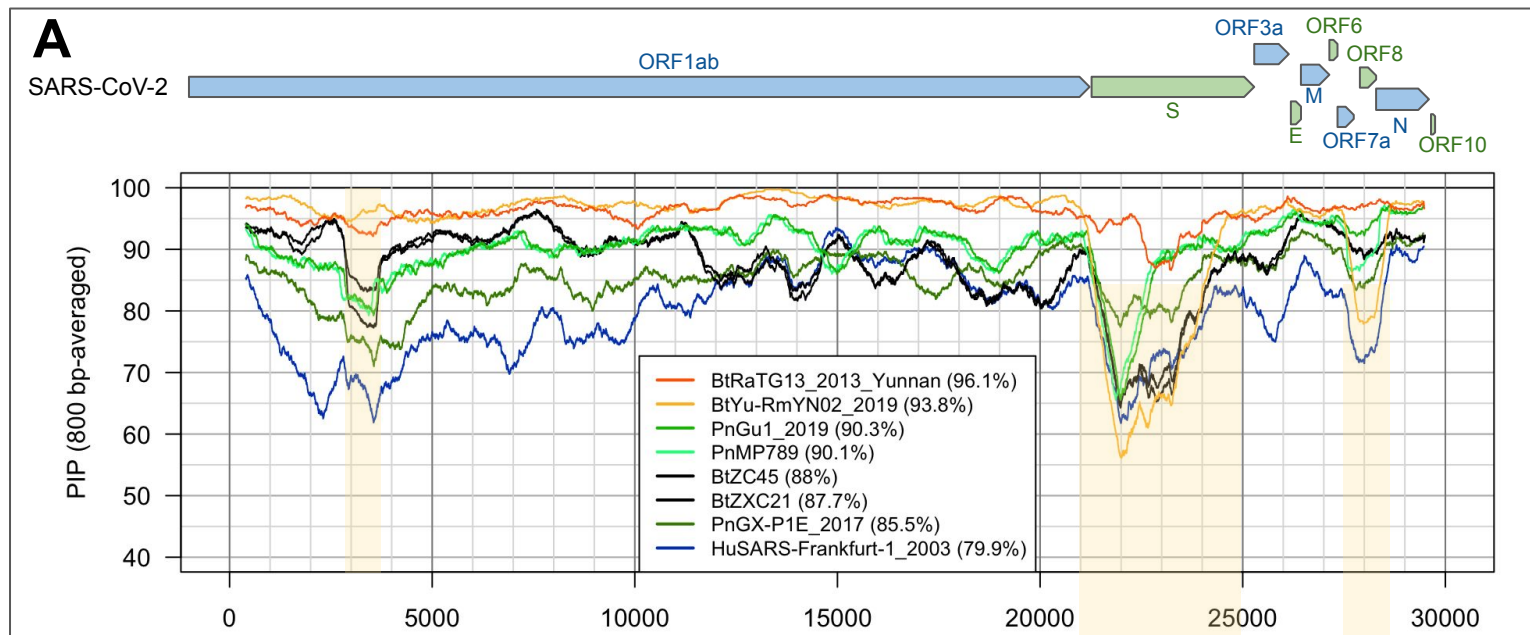
Les chercheurs ont peut-être accusé un peu trop vite le pangolin d'être à l'origine de la pandémie qui a débuté à Wuhan en Chine fin décembre 2019. Ce petit mammifère à écailles est de plus en plus disculpé par des études scientifiques.



Détecter les recombinaisons par comparaison de génomes

Détection des recombinaisons génomiques

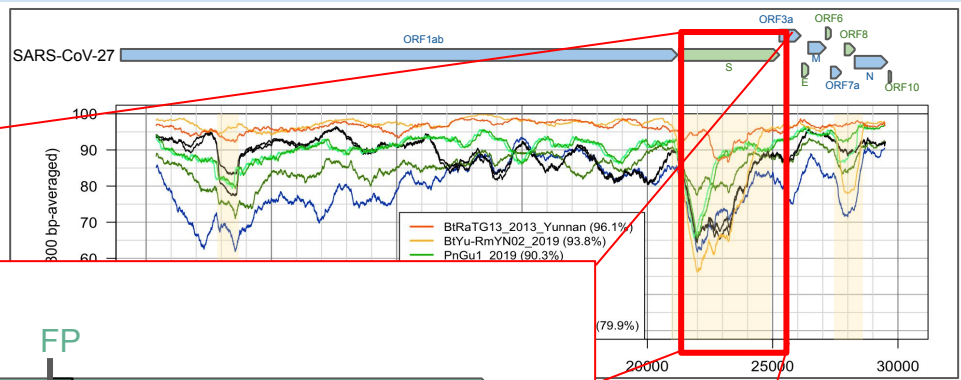
Les chutes brutales d'identité sur les profils PPI (fond jaune) dénotent des régions résultant vraisemblablement de recombinaisons.



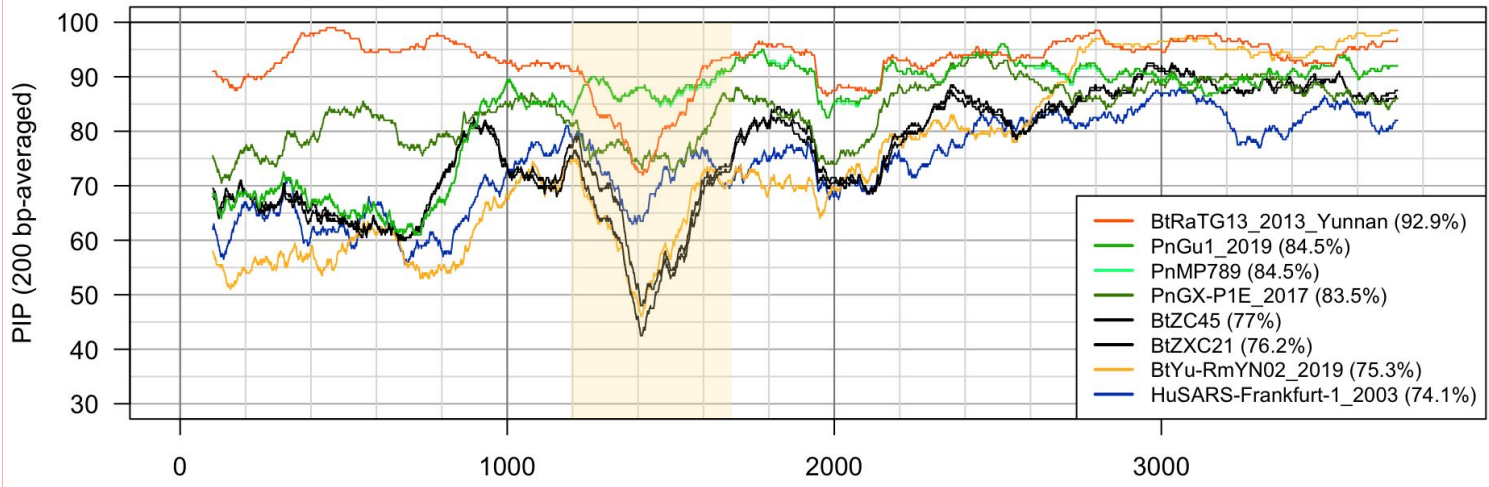
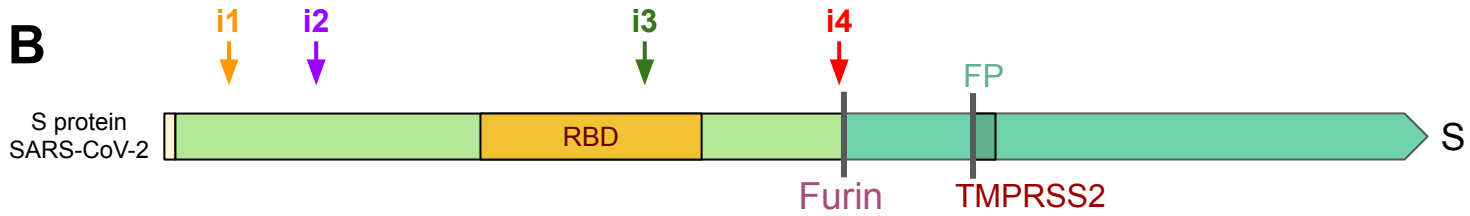
- Sallard, E., Halloy, J., Casane, D., van Helden, J. & Decroly, É. 2020. **Retrouver les origines du SARS-CoV-2 dans les phylogénies de coronavirus.** [Med Sci \(Paris\) 36: 783–796.](#)
- English version : Erwan Sallard, José Halloy, Didier Casane, Etienne Decroly, Jacques van Helden. **Tracing the origins of SARS-CoV-2 in coronavirus phylogenies.** [hal-02891455](#)

Comparaison entre coronavirus - gène S

Profils de pourcentages de positions identiques (PPI) entre régions génomiques du gène S de différents coronavirus et SARS-CoV-2 (la référence à 100%).



B



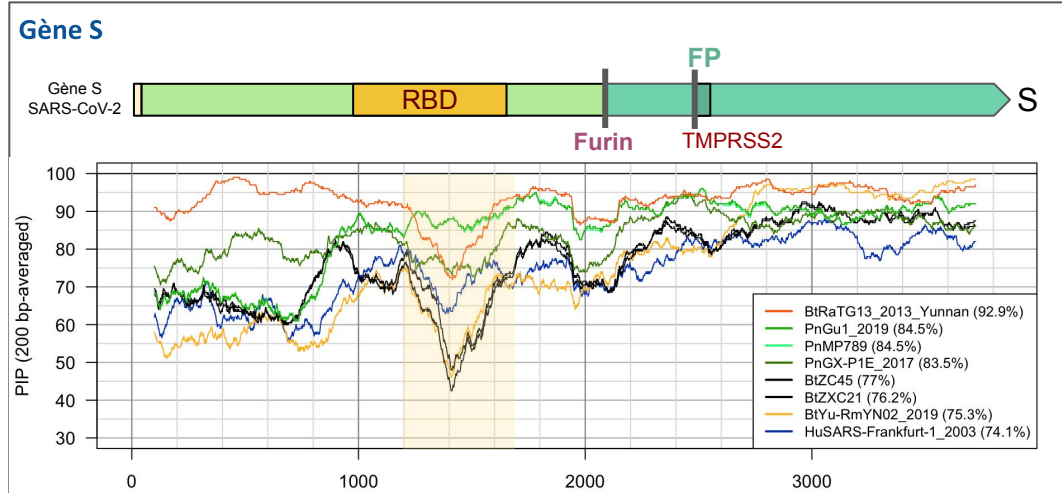
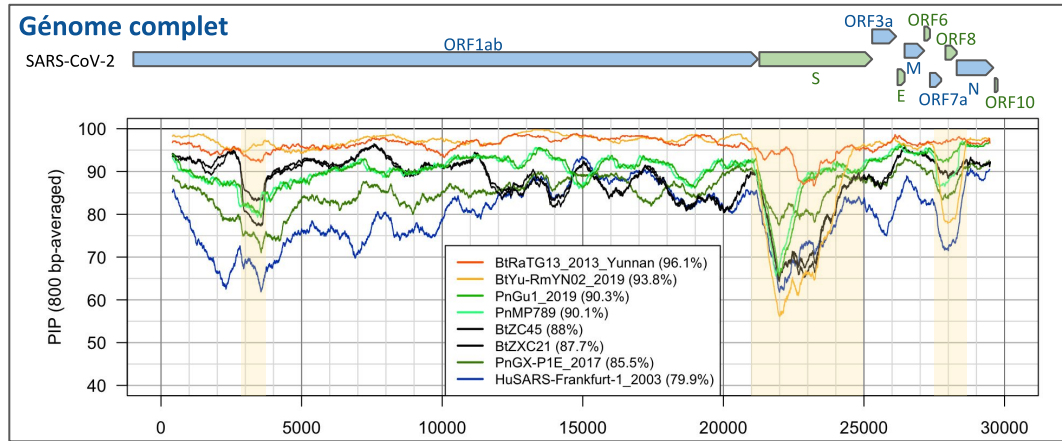
Recombinaisons génomiques dans les génomes de coronavirus

Le profil génomique

- régions ayant vraisemblablement fait l'objet de recombinaisons (fond jaune).

Profil PPI du gène spicule (S)

- **S:** spicule
- **RBD:** receptor binding domain
- Dans la région du RBD forte baisse des PPI
- Le RBD est en évolution rapide, pourquoi ?
 - Immunogène → forte pression sélective en faveur de variations qui permettent d'échapper à l'immunité
 - Spécificité d'espèce → modifications permettent de changer d'hôte



Un métagénome de virus de chauve-souris (RmYN02)

Current Biology

A Novel Bat Coronavirus Closely Related to SARS-CoV-2 Contains Natural Insertions at the S1/S2 Cleavage Site of the Spike Protein

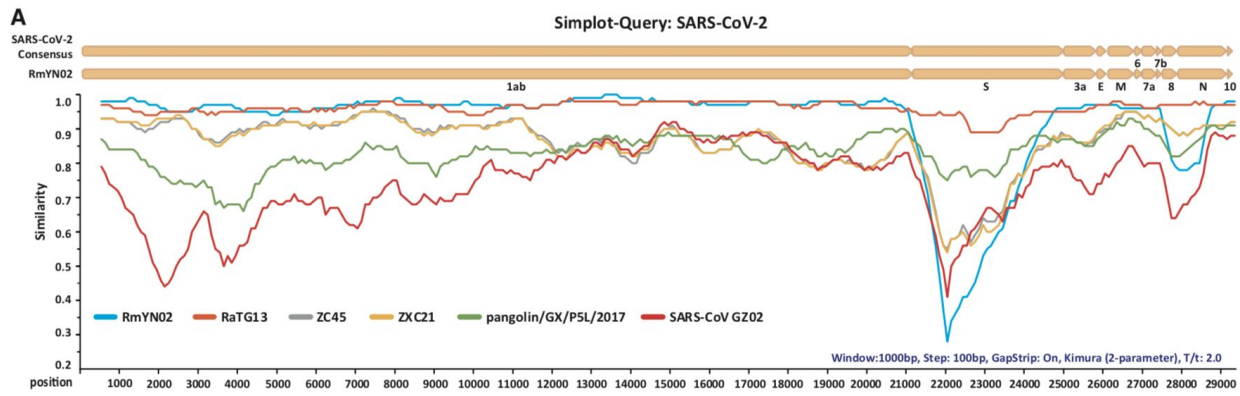
Highlights

- Metagenomic analysis identified a novel coronavirus, RmYN02, from *R. malayanus*

Authors

Hong Zhou, Xing Chen, Tao Hu, ..., Alice C. Hughes, Yuhai Bi, Weifeng Shi

- Hong Zhou et coll. (2019)
- Séquence métagénomique d'un coronavirus de chauve-souris assemblée à partir de 11 échantillons de chauves-souris de Yunnan
- Génome très proche de SARS-CoV-2 **sauf** dans la région du gène spicule (S).
 - ce virus serait donc un recombinant de RaTG13
 - oui mais** ce métagénome est reconstruit à partir de 11 échantillons différents. La recombinaison pourrait donc être un artéfact
- Les auteurs soulignent aussi la présence d'une insertion similaire au site furine de SARS-CoV-2.
 - Oui mais** ce scénario est peu vraisemblable, car il aurait nécessité quatre événements évolutifs (une insertion et trois délétions).



H

Polybasic cleavage site

SARS-CoV-2 Numbering	667	668	669	670	671	672	673	674	675	676	677	678	679	680	681	682	683	684	685	686	687	688	689	690	691	692	693
Consensus SARS-CoV-2	G	A	G	I	C	A	S	Y	Q	T	Q	T	N	S	P	R	R	A	R	S	V	A	S	Q	S	I	I
RmYN02	G	A	G	V	C	A	S	Y	-	-	-	-	N	S	P	-	A	A	R	-	V	G	T	N	S	I	I
RaTG13	G	A	G	I	C	A	S	Y	Q	T	Q	T	N	S	-	-	-	-	R	S	V	A	S	Q	S	I	I
ZC45	G	A	G	I	C	A	S	Y	H	T	A	S	I	L	-	-	-	-	R	S	T	S	Q	K	A	I	V
ZXC21	G	A	G	I	C	A	S	Y	H	T	A	S	I	L	-	-	-	-	R	S	T	G	Q	K	A	I	V
pangolin/MP789/2019	G	A	G	I	C	A	S	Y	Q	T	Q	T	N	S	-	-	-	-	R	S	V	S	S	X	A	I	I
pangolin/GX/P5L/2017	G	A	G	I	C	A	S	Y	H	S	M	S	S	F	-	-	-	-	R	S	V	N	Q	R	S	I	I
SARS-CoV GZ02	G	A	G	I	C	A	S	Y	H	T	V	S	L	L	-	-	-	-	R	S	T	S	Q	K	S	I	V
RmYN01	G	A	G	I	C	A	S	Y	H	T	A	S	L	L	-	-	-	-	R	N	T	G	Q	K	S	I	V

O-linked glycan residues

Bases de données de séquences biologiques

- KB: knowledge base
- Swiss-prot contient des séquences protéiques “annotées” par des biologistes. L’annotation consiste à associer à une séquence les connaissances résultant d’expérimentation.
- Deux limitations
 - Le nombre de publications augmente tellement qu’il n’est pas possible à l’équipe de Swiss-prot de tout annoter
 - Le nombre de séquences augmente de façon tellement rapide qu’il est impossible de toutes les caractériser expérimentalement
- TREMBL: annotation automatique des séquences traduites de la base de données EMBL.
- Uniprot = Swiss-prot + TREMBL

Le NCBI (USA) propose une série de bases de données couvrant les différents domaines de la biologie : taxonomie, séquences protéiques, séquences nucléiques, publications biologiques, ...

Nous utiliserons le site Web “Entrez” pour consulter les séquences génomiques de coronavirus et leurs annotations.

Alignement d'une paire de séquences

Exercice

- On dispose des deux séquences suivantes
 - **Seq1** TTTGCGTTAAATCGTGTAGCAATTTAA
 - **Seq2** AAGAATGGCGTTTTTAATAGCAATAT
- Questions
 1. En décalant progressivement les séquences, identifiez le(s) décalage(s) qui révèlent des régions de similarité.
 2. A chaque position de décalage, identifiez les segments parfaitement conservés (successions ininterrompue de résidus identiques).
 3. Au vu du résultat, pensez-vous que l'insertion d'un gap permettrait d'augmenter le score d'alignement?

Solution de l'exercice

- On dispose des deux séquences suivantes
 - Seq1 **TTT**GCGTTAAATCGTGTAGCAATTTAA
 - Seq2 **AAGAATGGCGTTTTTAATAGCAATAT**
- Questions
 1. En décalant progressivement les séquences, identifiez le(s) décalage(s) qui révèlent des régions de similarité.
 2. A chaque position de décalage, identifiez les segments parfaitement conservés (successions ininterrompue de résidus identiques).
 3. Au vu du résultat, pensez-vous que l'insertion d'un gap permettrait d'augmenter le score d'alignement?
- Décalage -4
 - Position -4 123456789
 - Seq1 1234**TTT**GCGTTAAATCGTGTAGCAATTTAA
 - Seq2 **AAGAATGGCGTTTTTAATAGCAATAT**
- Décalage -1
 - Seq1 **TTT**GCGTTAAATCGTGTAGCAATTTAA
 - Seq2 **AAGAATGGCGTTTTTAATAGCAATAT**

Alignement avec « gaps » (brèches)

- Les alignements sans gaps sont rarement pertinents, car les divergences entre séquences incluent souvent des insertions et délétions.
- Les gaps permettent de mettre en évidence les régions de similarités multiples.

```
----TTTGCGGTT--AAATCCGTGTAGCAATTTAA      s=substitution; |=identité
1111s|s|||||11s||22222|s|22                1=gap dans la 1ère séquence
AAGAATGGCGGTTTTAA-----TAGCAATAT--      2=gap dans la 2de séquence
```

- Gaps, insertions et délétions
 - Les “**gaps**” (**brèches**) reflètent soit une insertion dans l’une des séquences, soit une délétion dans l’autre.
 - Sur seule base de l’alignement d’une paire de séquences, on ne peut pas déterminer si un gap correspond à une délétion ou une insertion.
 - On utilise le terme **indel** pour désigner cet événement évolutif de nature indéterminée (insertion ou délétion) qui a donné lieu à un gap observé dans un alignement.

Alignements globaux (Needleman-Wunsch) versus locaux (Smith-Waterman)

■ Alignement global

- Approprié, par exemple, pour les protéines homologues qui sont conservées sur toute leur longueur.
- L'alignement final inclut obligatoirement les deux séquences complètes.

```
LQGPSKTGKGS-SRSWDN
|----|-|||---|--|
LN-ITKAGKGAIMRLGDA
```

- Algorithme: **Needleman-Wunsch** (1970).
- Outil web EMBOSS : **needle** (nucleic acids(nucleic acids or proteins)).

■ Alignement local

- Approprié, par exemple, pour les protéines qui partagent un domaine commun, restreint à un segment de chaque séquence.

```
LQGPSSKTGKGS-SSRIWDN
  | - | | |
LN-ITKKAGKGAIMRLGDA
```

- L'alignement final est restreint aux segments conservés.

```
KTGKG
| - | | |
KAGKG
```

- Algorithme: **Smith-Waterman** (1981).
- Outil Web EMBOSS : **water** (nucleic acids(nucleic acids or proteins))

■ Needleman, S. B. & Wunsch, C. D. (1970). A general method applicable to the search for similarities in the amino acid sequence of two proteins. J Mol Biol 48, 443-53.

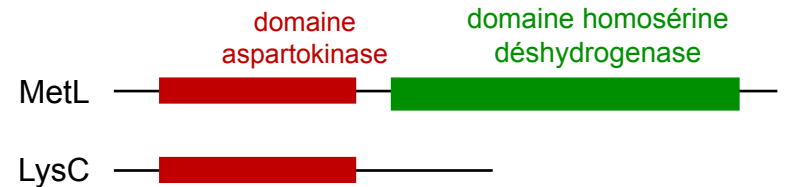
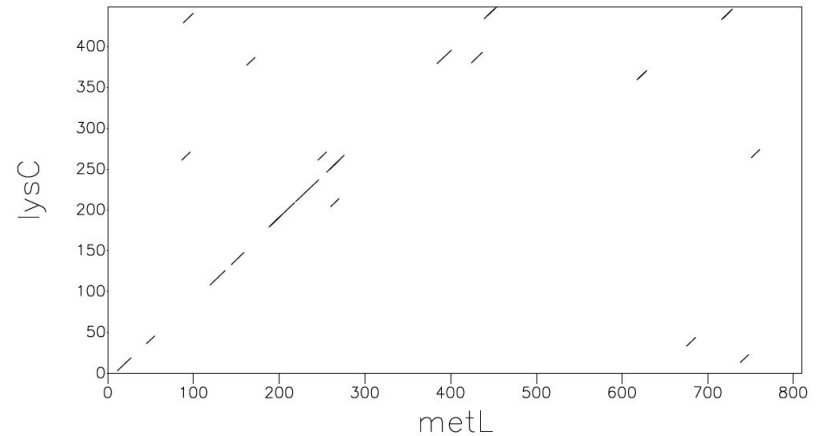
■ Smith, T. F. & Waterman, M. S. (1981). Identification of common molecular subsequences. J Mol Biol 147, 195-7.

Aspartokinases: dot plot avec matrice de substitution (BLOSUM62)

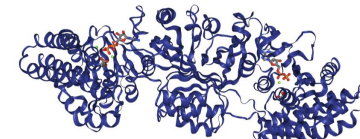
- Avec le logiciel *dotmatcher*, une matrice de substitution est utilisée pour assigner un score à chaque paire de résidus. Les segments de lignes indiquent des régions de correspondance entre séquences.
- Ceci révèle la similarité entre les **domaines aspartokinase** de LysC (l'ensemble de la séquence) et de MetL (positions 1 à ~450).
- La région de similarité ne recouvre que la partie N-terminale (gauche) de MetL, car il s'agit d'une enzyme bi-fonctionnelle. La région C-terminale de MetL contient un **domaine homosérine déshydrogenase** qui est absent de la protéine LysC.
- Sur base de ce dessin, on comprend qu'un alignement local sera plus pertinent qu'un alignement global car il révélera le domaine que ces protéines ont en commun.

Dotmatcher: metL vs lysC

(windowsize = 10, threshold = 23.00 08/09/04)



Structure 3D du domaine aspartokinase



Alignement global : exemple (Needleman-Wunsch)

```
# Matrix: EBLOSUM62
# Gap_penalty: 10.0
# Extend_penalty: 0.5
# Length: 867
# Identity:      254/867 (29.3%)
# Similarity:    423/867 (48.8%)
# Gaps:         104/867 (12.0%)
# Score: 929.0
```

```
metL      1  MSVIAQAGAKGRQLHKFGSSSLADVKCYLRVAGIMAEYSQPDDM-MVVSA      49
          . . . . | | | | : | : | : . . . . : | | | . | : . . . . : . | : | |
thrA      1  MRVLLKFGGTSVANAERFLRVADILESNAEQGVATVLSA      39

metL     50  AGSTTNQLINWLK-----LSQTDRLSAHQVQQTLLRRYQC DLISG      88
          . . . . | | . | : . . . . : | : . . | : . | : : : | : : |
thrA     40  PAKITNHLVAMIEKTISGQDALPNISDAERIFA-----ELLTG      77

metL     89  LLPAAEADSL--ISAFV-SDLERLAALLDSGIN-----DAVYAEVVGHG      129
          | . | : . . . . | . . | | . . . . . . . | . | | : | : . . | : . . . . |
thrA     78  LAAAQPGFPLAQLKTFVDQEFAQIKHVL-HGISLLGQCPDSINAALICRG      126

metL    130  EVWSARLMSAVLNQOGLPAAWLD-AREFLRAERAAQPQVD--EGLSYPLL      176
          | . . | . . : | : . | | . . . . . | . . . . | . . . . . . . | | | . . . . .
thrA    127  EKMSIAIMAGVLEARGHNVTVIDPVEKLLAVGHYLESTVDIAESTRRIAA      176

metL    177  QQLLVQHHPGKRLVVTGFISRNNAGETVLLGRNGSDYSATQIGALAGVSRV      226
          . : . . . . | . . . . . | | . . . | . . | | . | : | | | | | | | | | . . . . | . . . . .
thrA    177  SRIPADH---MVL MAGFTAGNEKGELVVLGRNGSDYSAAVLAACL RADCC      223

metL    227  TIWSDVAGVYSADPRKVKDA CLLPLRLDEASELARLAAPV LHARTLQPV      276
          . | | : | | . | | | | : | . | | . | | . | | . | | . | | . | | . | | . | | : | : | : | : |
```

- Alignement des protéines metL et thrA d'E.coli avec l'algorithme de Needleman-Wunsch.
- Barres verticales « | »
 - **Identité**: les deux résidus alignés sont identiques.
- Doubles points « : »
 - **Substitution « conservative »**
 - Les deux résidus alignés sont différents mais **similaires** (la paire de résidus a un score positif dans la matrice de substitution utilisée (ici, BLOSUM62). Voir plus loin pour comprendre ces matrices.
- Points « . »
 - **Substitution non-conservative**
 - Cette paire de résidus (distincts) a un score négatif dans la matrice de substitution.
- Espace: « »
 - **Gap**: les résidus d'une des deux séquences ne correspondent à aucun résidu sur l'autre.
 - Le gap peut provenir soit d'une délétion, soit d'une insertion, on parle donc d'**indel**, pour désigner l'événement évolutif d'où provient ce gap.

Needleman-Wunsch with partial similarities

```
# Matrix: EBLOSUM62
# Gap_penalty: 10.0
# Extend_penalty: 0.5
# Length: 854
# Identity:      136/854 (15.9%)
# Similarity:    209/854 (24.5%)
# Gaps:          449/854 (52.6%)
# Score: 351.0

metL      1 MSVIAQAGAKGRQLHKFGGSSLADVKCYLRVAGIMA EYSQPDDMMVVSAA      50
  ||.|.      :.|||||:|:|.....|.|. |:.....  .::|:|:|:
lysC      1 MSEIV-----VSKFGGT SVADFDAMNRSADIVLSDANV-RLVVL SAS      41

metL     51 GSTTNQLINWLK-LSQTDRLSAHQVQQT LRRYQC DLI SGL----LP AEEA      95
  ...||.|:....: |....|.  :.....:|.|. :.....|  :..||.
lysC     42 AGITNLLVALAEGLEPGERF---EKLDAIRNIQFAILERLRYPNVIREEI      88

metL     96 DSLISAFVSDLERLAALLDSGINDAVYA E VVGHG EVWSARLMSAVLNQQG      145
  :.:...  ::.|...|||.|.  .|:..:|:|.|||:|.|.|. :...:|:..:
lysC     89 ERLLEN-ITVLA EAAALATS---PALTDELVSHGELMSTLLFVEILRERD      134

metL    146 LPAAWLDAREFLRA-ERAAQPQVDEGLSYPL LQQLLVQH P GKRLVVT-GF      193
  :.|.|.|.:...:  :|.....:|. :.....|.....|:.....:|:|:|  ||
lysC    135 VQAQWFDVRKVMRTNDRFGRAEPDIAALAE LAALQLL PRLNEGLVITQGF      184

metL    194 ISRN NAGETVLLGRNGSDYSATQIGALAGVSRVTI WSDVAGVYSADPRKV      243
  |...|.|.|. .|||.|||:|:..:.....| |.|. |:|:|. |:|:|. |
lysC    185 IGSENKGR TTTTLGRGGS DYTAALLAEALHASRVDI WTDVPGIYTTDPRVV      234

metL    244 KDA CL LPLLR LDEASELARLAAPV LHARTLQPVSGSEIDLQLRCSYTPDQ      293
```

- Alignment of *E. coli* lysC and metL proteins with Needleman-Wunsch algorithm.
- metL contains two domains: aspartokinase and homoserine dehydrogenase.
- LysC only contains the aspartokinase domains.
- With Smith-Waterman, the %similarity is calculated over the whole length of the alignment (854aa), which gives 24.5%.
- Actually, most of the alignment length is in the terminal gap (the homoserine dehydrogenase domain of metL).
- This percentage is lower than the usual threshold for considering two proteins as homolog.

Smith-Waterman with partial similarities

```
# Matrix: EBLOSUM62
# Gap_penalty: 10.0
# Extend_penalty: 0.5
# Length: 482
# Identity:      133/482 (27.6%)
# Similarity:    205/482 (42.5%)
# Gaps:          85/482 (17.6%)
# Score: 353.5
metL      16  KFGGSSLADVKCYLRVAGIMA EYSQPDDMMVVSAAGSTTNQLINWLK-LS      64
      ||||:|:|.....|.|. |:..... .::|:|:|:..|. |:...: |.
lysC      8  KFGGTSVADFDAMNRSADIVLSDANV-RLVVLSASAGITNLLVALAEGLE      56

metL     65  QTDRLSAHQVQQTLLRRYQC DLISGL----LPAAEADSLISAFVSDLERLA    110
      .:|. :.....|..|.....| :..||:.. |:.. :.:|. .|. |
lysC     57  PGERF---EKLDAIRNIQF AILERLRYPNVIREEIERLLEN-ITVLAEAA      102

metL    111  ALLDSGINDAVYAEVVGHGEVWSARLMSAVLNQQGLPAAWLDAREFLRA-    159
      ||..| . |:.. |:|. |||:|. |..|.....|:.....|. |. |. |:..|.
lysC    103  ALATS---PALTDELVSHGELMSTLLFVEILRERDVQAQWFDVRKVMRTN    149

metL    160  ERAAQPVQDEGLSYPLLQQLLVQH PGKRLVVT-GFISRNNAGETVLLGRN    208
      :|:.....|.....|.....|:.....:| |:| |||...|. |. |. |||.
lysC    150  DRFGRAEPDIAALAEALQLLPRLNEGLVITQGFIGSENKGRTTTTLGRG    199

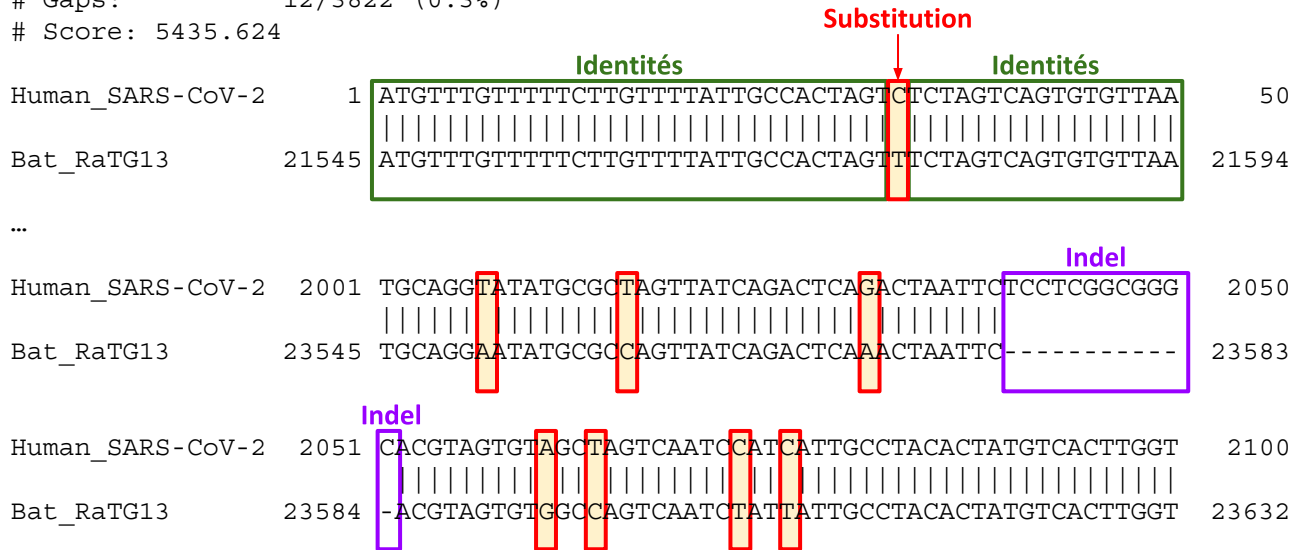
metL    209  GSDYSATQIGALAGVSRVTI WSDVAGVYSADPRKVKDA CLLPLRLDEAS    258
      ||||:|:|:.....| |||. |:|:|:|:|. |:|:|. |..|.....:..| |:
lysC    200  GSDYTAALLAEALHASRVDI WTDVPGIYTTDPRVVSAAKRIDEIAFAEAA    249

metL    259  ELARLAAPVLHARTLQP VSGSEIDLQLRCSYTPDQGSTRI-----E      299
```

- Alignment of *E.coli* lysC and metL proteins with Smith-Waterman algorithm.
- The alignment is almost identical to the one reported by Needleman-Wunsch, but the score is now considered on the aligned segments only (482 aa).
- On this region, there is 42.5% of similarity.

Alignement de séquences – Gènes S de SARS-CoV-2 et RaTG13

```
# Aligned_sequences: 2
# 1: Human_SARS-CoV-2_BetaCoV/Wuhan/IPBCAMS-WH-01/2019
# 2: Bat_RaTG13
#
# Length: 3822
# Identity: 3549/3822 (92.9%)
# Similarity: NA/3822 (NA%)
# Gaps: 12/3822 (0.3%)
# Score: 5435.624
```



Note

- “Indel” signifie “Insertion ou délétion”
- Sur base de ce résultat, la différence observée peut provenir soit d’une insertion chez un ancêtre de SARS-CoV-2, soit d’une délétion chez un ancêtre de RaTG13

Recherche de séquences par similarité

- Les diapos précédentes illustraient des alignements entre une paire de séquences choisies a priori.
- Dans certains cas on dispose de la séquence d'un gène ou d'une protéine d'intérêt, et on désire l'aligner avec **toutes** les séquences connues.
- Ceci permet par exemple de détecter des ressemblances entre une séquence inconnue (par exemple une séquence codante qu'on vient d'identifier dans un génome nouvellement séquencé) et des séquences déjà présentes dans une base de données, et dont la fonction est connue.
- La recherche par similarité est aujourd'hui l'approche principale pour assigner des fonctions aux gènes (alignement d'ADN) et aux protéines (alignement de séquences peptidiques).

- La similarité entre deux traits (organes, séquences) peut s'interpréter par deux hypothèses alternatives: homologie et analogie.
- **Homologie**
 - La similarité s'explique par le fait que les deux séquences divergent d'un ancêtre commun.
 - Les différences entre les deux caractères homologues résultent de l'accumulation de mutations à partir de l'ancêtre commun. Il s'agit donc d'une évolution par **divergence évolutive**.
- **Analogie**
 - Ressemblance entre deux traits (organes, séquence) qui ne résulte pas d'une origine ancestrale commune (par opposition à l'homologie).
 - Les traits similaires sont apparus de façon **indépendante**. Leur ressemblance peut éventuellement manifester l'effet d'une pression évolutive qui a sélectionné les mêmes propriétés.
 - Dans ce cas, on parle de **convergence évolutive**.

Matrices de substitutions

- Une **matrice de substitution** associe un score à chaque paire de résidus qu'on peut trouver dans un alignement.
 - Chaque ligne et chaque colonne représente l'un des résidus (4 nucléotides, 20 acide aminés).
 - La diagonale correspond aux identités.
 - Le triangle inférieur correspond à des substitutions.
 - Le triangle supérieur est symétrique au triangle inférieur, il n'est pas nécessaire d'indiquer les nombres.
- Les **scores négatifs** sont considérés comme des pénalités associées à certaines substitutions qu'on n'observe que rarement dans les alignements. Les algorithmes d'alignements tenteront donc d'éviter ces substitutions.
- Les **scores positifs** correspondent à des substitutions qu'on observe plus souvent que prévu, dans les alignements d'un grand nombre de séquences. Ceci suggère que ces substitutions particulières sont moins dommageable que d'autres, et on les qualifie donc de « **substitutions conservatives** » ou encore de « **mutations ponctuelles acceptées** » (*PAM*).
- Au sein d'un alignement, le terme **similarité** désigne les positions où se superposent des résidus ayant un score positif dans la matrice de substitution (identité ou substitution conservative).

Matrice de substitutions entre nucléotides

	A	C	G	T
A	2			
C	-2	2		
G	-2	-2	2	
T	-1	-2	-2	2

Matrice de substitutions entre acides aminés

Ala	A	4																			
Arg	R	-1	5																		
Asn	N	-2	0	6																	
Asp	D	-2	-2	1	6																
Cys	C	0	-3	-3	-3	9															
Gln	Q	-1	1	0	0	-3	5														
Glu	E	-1	0	0	2	-4	2	5													
Gly	G	0	-2	0	-1	-3	-2	-2	6												
His	H	-2	0	1	-1	-3	0	0	-2	8											
Ile	I	-1	-3	-3	-3	-1	-3	-3	-4	-3	4										
Leu	L	-1	-2	-3	-4	-1	-2	-3	-4	-3	2	4									
Lys	K	-1	2	0	-1	-3	1	1	-2	-1	-3	-2	5								
Met	M	-1	-1	-2	-3	-1	0	-2	-3	-2	1	2	-1	5							
Phe	F	-2	-3	-3	-3	-2	-3	-3	-3	-1	0	0	-3	0	6						
Pro	P	-1	-2	-2	-1	-3	-1	-1	-2	-3	-3	-1	-2	-4	7						
Ser	S	1	-1	1	0	-1	0	0	0	-1	-2	-2	0	-1	-2	-1	4				
Thr	T	0	-1	0	-1	-1	-1	-1	-2	-2	-1	-1	-1	-1	-2	-1	1	5			
Trp	W	-3	-3	-4	-4	-2	-2	-3	-2	-2	-3	-2	-3	-1	1	-4	-3	-2	11		
Tyr	Y	-2	-2	-2	-3	-2	-1	-2	-3	2	-1	-1	-2	-1	3	-3	-2	2	7		
Val	V	0	-3	-3	-3	-1	-2	-2	-3	-3	3	1	-2	1	-1	-2	-2	0	-3	-1	4
	A	R	N	D	C	Q	E	G	H	I	L	K	M	F	P	S	T	W	Y	V	

Utilisation d'une matrice de substitution pour calculer le score d'un alignement

$$S = \sum_{i=1}^L S_{r_{1,i}, r_{2,i}}$$

Ala	A	4																							
Arg	R	-1	5																						
Asn	N	-2	0	6																					
Asp	D	-2	-2	1	6																				
Cys	C	0	-3	-3	-3	9																			
Gln	Q	-1	1	0	0	-3	5																		
Glu	E	-1	0	0	2	-4	2	5																	
Gly	G	0	-2	0	-1	-3	-2	-2	6																
His	H	-2	0	1	-1	-3	0	0	-2	8															
Ile	I	-1	-3	-3	-3	-1	-3	-3	-4	-3	4														
Leu	L	-1	-2	-3	-4	-1	-2	-3	-4	-3	2	4													
Lys	K	-1	2	0	-1	-3	1	1	-2	-1	-3	-2	5												
Met	M	-1	-1	-2	-3	-1	0	-2	-3	-2	1	2	-1	5											
Phe	F	-2	-3	-3	-3	-2	-3	-3	-3	-1	0	0	-3	0	6										
Pro	P	-1	-2	-2	-1	-3	-1	-1	-2	-2	-3	-3	-1	-2	-4	7									
Ser	S	1	-1	1	0	-1	0	0	0	-1	-2	-2	0	-1	-2	-1	4								
Thr	T	0	0	0	-1	-1	-1	-1	-2	-2	-1	-1	-1	-1	-2	-1	1	5							
Trp	W	-3	-3	-4	-4	-2	-2	-3	-2	-2	-3	-2	-1	-1	-1	1	-4	-3	-2	11					
Tyr	Y	-2	-2	-2	-3	-2	-1	-2	-3	2	-1	-1	-2	-1	3	-3	-2	-2	2	7					
Val	V	0	-3	-3	-3	-1	-2	-2	-3	-3	3	1	-2	1	-1	-2	-2	0	-3	-1	4				

- Les matrices de substitution sont utilisées pour calculer le score d'un alignement.
- Ce score est la somme, pour toutes les positions de l'alignement (*i de 1 à L*), des scores des paires de résidus ($r_{1,i}$ et $r_{2,i}$).
- Les "gaps" sont traités par une règle spécifique reposant sur deux paramètres de pénalité:
 - Pénalité d'ouverture de gap (**go**)
 - Valeurs typiques: entre -10 et -15
 - Pénalité d'extension de gap (**ge**)
 - Valeurs typiques: entre -0.5 et -2

<i>i</i>	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	
	R	L	A	S	V	E	T	D	M	P	-	-	-	-	-	L	T	L	R	Q	H	
	.		.		:	:		.	:	.	go	ge	ge	ge	ge		
	T	L	T	S	L	Q	T	T	L	K	N	L	K	E	M	A	H	L	G	T	H	
S	-1	+4	+0	+4	+1	+2	+5	-1	+2	-1	-10	-1	-1	-1	-1	-1	-2	+4	-2	-1	+8	= 7

Résultat de BLAST – Requête peptidique vs DB de peptides

```
>gi|16127996|ref|NP_414543.1| bifunctional: aspartokinase I  
(N-terminal); homoserine dehydrogenase I (C-terminal)  
[Escherichia coli K12]  
Length = 820
```

```
Score = 344 bits (882), Expect = 2e-95
```

```
Identities = 247/821 (30%) Positives = 410/821 (49%) Gaps = 44/821 (5%)
```

```
Query: 16 KFGGSSLADVKCYLRVAGIMA EYSQ PDDMM-VVSAAGSTTNQLINWLKLSQTDRLSAHQV 74  
KFGG+S+A+ + +LRVA I+ ++ + V+SA TN L+ ++ + + + + +  
Sbjct: 5 KFGGTSVANAERFLRVADILESMAROGQVATVLSAPAKITNHLVAMIEKTISGQDALPNI 64
```

```
Query: 75 QQTLRRYQCDLISGLLPAREADS L--ISAFVSDLERLAALLDSGIN-----DAVYAEVV 126  
R + +L++GL A+ L + FV + GI+ D++ A ++  
Sbjct: 65 SDAERIF-AELLTGLAAACPGFPLAQLKTFVDQEFQAQIKHVLHGISLLGQCPSINAALI 123
```

```
Query: 127 GHGEVWSARLMSAVLNQOGLPAAWLDAREFLRAER---AAQPQVDEGLSYPLLQQLLVQH 183  
GE S +M+ VL +G +D E L A + + E ++ H  
Sbjct: 124 CRGEKMSIAIMAGVLEARGHNVTVIDPVEKLLAVGHYLESTVDIAESTRRIAASRIPADH 183
```

```
Query: 184 PGKRLVVTGFI SRN NAGETVLLGRNGSDYSATQIGALAGVSRVTI WSDVAGVYSADPRKV 243  
+++ GF + N GE V+LGRNGSDYSA + A IW+DV GVY+ DPR+V  
Sbjct: 184 ---MVL MAGFTAGNEKGELVVLGRNGSDYSAAVLAACLRADCCEIWTDVDGVYTC DPRQV 240
```

```
Query: 244 KDACLLPLLRLDEASELARLAAPVLHARTLQPVSGSEIDLQLRCSYTPDQ-----GSTRI 298  
DA LL + EA EL+ A VLH RT+ P++ +I ++ + P G++R  
Sbjct: 241 PDARLLKMSYQEAMELSYFGAKVLHPRTITPIAQFQIPCLIKNTGNPQAPGTLIGASRD 300
```

```
Query: 299 ERVLASGTGARIVTSHDDVCLIEFQVPASQDFKLAHKEIDQILKRAQVRPLAVGVHNDRO 358  
E L + + + + + + P + + + RA++ + + +  
Sbjct: 301 EDELP----VKGISNLNNMAMFSVSGPMKGMVGMGAARVFAAMSRARISVVLITOSSSEY 356
```

- La ligne entre les séquences “Query” et “Sbjct” indique les correspondances entre acides aminés.

- Identités
- Substitutions “conservatives”:
paires de résidus distincts mais dont la substitution est *généralement* moins délétère que pour d’autres paires de résidus.

- Substitutions non conservatives
- Positives: identités + substitutions conservatives.

- Gaps: lacunes insérées dans une séquence afin d’optimiser l’alignement des fragments avoisinants.

Résultat de BLAST – Requête peptidique vs DB de peptides

```
>gi|16127996|ref|NP_414543.1| bifunctional: aspartokinase I  
  (N-terminal); homoserine dehydrogenase I (C-terminal)  
  [Escherichia coli K12]  
Length = 820
```

```
Score = 344 bits (882), Expect = 2e-95
```

```
Identities = 247/821 (30%), Positives = 410/821 (49%), Gaps = 44/821 (5%)
```

```
Query: 16  KFGGSSLADVKCYLRVAGIMAEYSQPDDMM-VVSAAGSTTNQLINWLKLSQTDRLSAHQV 74  
          KFGG+S+A+ + +LRVA I+  ++  +  V+SA  TN L+  ++  +  +  +  +  +  
Sbjct: 5   KFGGTSVANAERFLRVADILESNDARQGQVATVLSAPAKITNHLVAMIEKTISGQDALPNI 64
```

```
Query: 75  QQTLRRYQCDLISGLLPAAEADSL--ISAFVSDLERLAALLDSGIN-----DAVYAEVV 126  
          R  +  +L++GL  A+  L  +  FV          +  GI+  D++  A  ++  
Sbjct: 65  SDAERIF-AELLTGLAAAQPGFPLAQLKTFVDQEFQAQIKHVLHGISLLGQCPDSINAALI 123
```

```
Query: 127 GHGEVWSARLMSAVLNQOGLPAAWLDAREFLRAER--AAQPQVDEGLSYPLLQQLLVQH 183  
          GE  S  +M+  VL  +G  +D  E  L  A  +  +  E  ++  H  
Sbjct: 124 CRGEKMSIAIMAGVLEARGHNVTVIDPVEKLLAVGHYLESTVDIAESTRRIAASRIPADH 183
```

```
Query: 184 PGKRLVVTGFI SRN NAGETVLLGRNGSDYSATQIGALAGVSRVTI WSDVAGVYSADPRKV 243  
          +++  GF  +  N  GE  V+LGRNGSDYSA  +  A  IW+DV  GVY+  DPR+V  
Sbjct: 184 ---MVL MAGFTAGNEKGELVVLGRNGSDYSAAVLAACLRADCCEIWTDVDGVYTC DPRQV 240
```

```
Query: 244 KDACLLPLLRRLDEASELARLAAPVLHARTLQPVSGSEIDLQLRCSYTPDQ-----GSTRI 298  
          DA  LL  +  EA  EL+  A  VLH  RT+  P++  +I  ++  +  P  G++R  
Sbjct: 241 PDARLLKMSYQEAMELSYFGAKVLHPRTITPIAQFQIPCLIKNTGNPQAPGTLIGASRD 300
```

```
Query: 299 ERV LASGTGARIVTSHDDVCLIEFQVPASQDFKLAHKEIDQILKRAQVRPLAVGVHNDRO 358  
          E  L  +  +++  +++  +  P  +  +  +  RA++  +  +  +  
Sbjct: 301 EDELP----VKGISNLNNMAMFSVSGPMKGMVMAARVFAAMSRARISVVLITOSSSEY 356
```

- Exemple de résultat de recherche par similarité de séquences.
 - Requête (query): metA
 - Protéine identifiée dans la base de données: (subject): thrA.
- Le premier critère d'évaluation d'un résultat de BLAST:
 - La **e-valeur (expect)** indique le nombre de faux-positifs attendus au hasard, si l'on plaçait le seuil **au niveau du score observé** (344 bits dans ce cas-ci).
 - Plus la e-valeur est faible, plus le résultat est fiable.
 - Si la e-valeur est ≥ 1 , le résultat est peu fiable (on aurait facilement obtenu un « aussi bon » alignement avec des séquences aléatoires.

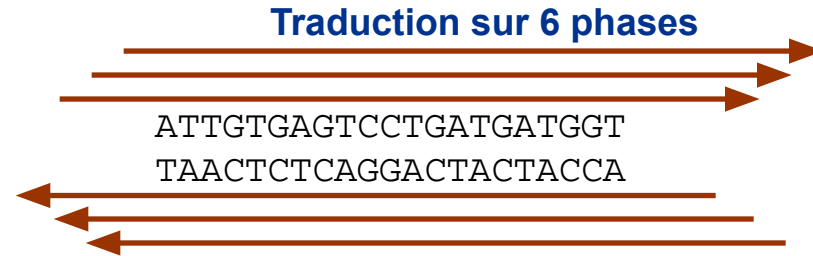
Modalités de BLAST

DNA versus protein searches

- When the query is a coding DNA sequence, it is recommended to apply searches with the translated rather than raw DNA sequences
 - This allows to introduce a substitution matrix (PAM, BLOSUM, ...), which better reflects the evolutionary changes.
 - It has been shown that some distant relationships can be detected with translated searches, but escape detection with the DNA search.
 - It is easier to filter out low complexity regions from proteins than from DNA sequences.

Traduction d'une séquence nucléique dans les 6 phases

- Si l'on dispose d'une séquence nucléique, on peut facilement déduire la séquence de la protéine qui pourrait être produite par sa traduction, sur chacun des 6 brins.
- Si cette séquence n'est pas codante, on s'attend à trouver des codons stop assez fréquemment (3 codons sur 64).
- Cependant, rien n'empêche d'aligner les 6 séquences ainsi produites avec d'autres séquences peptidiques.

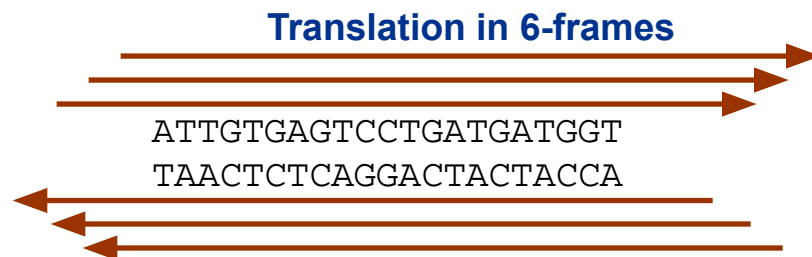


Résultat

F1	I	V	S	P	D	D	G
F2	L	*	V	L	M	M	V
F3	C	E	S	*	*	W	X
1	ATTGTGAGTCCTGATGATGGT	21					
	----:---- ----:---- -						
1	TAACTCTCAGGACTACTACCA	21					
F6	X	T	L	G	S	S	P
F5	X	Q	S	D	Q	H	H
F4	N	H	T	R	I	I	T

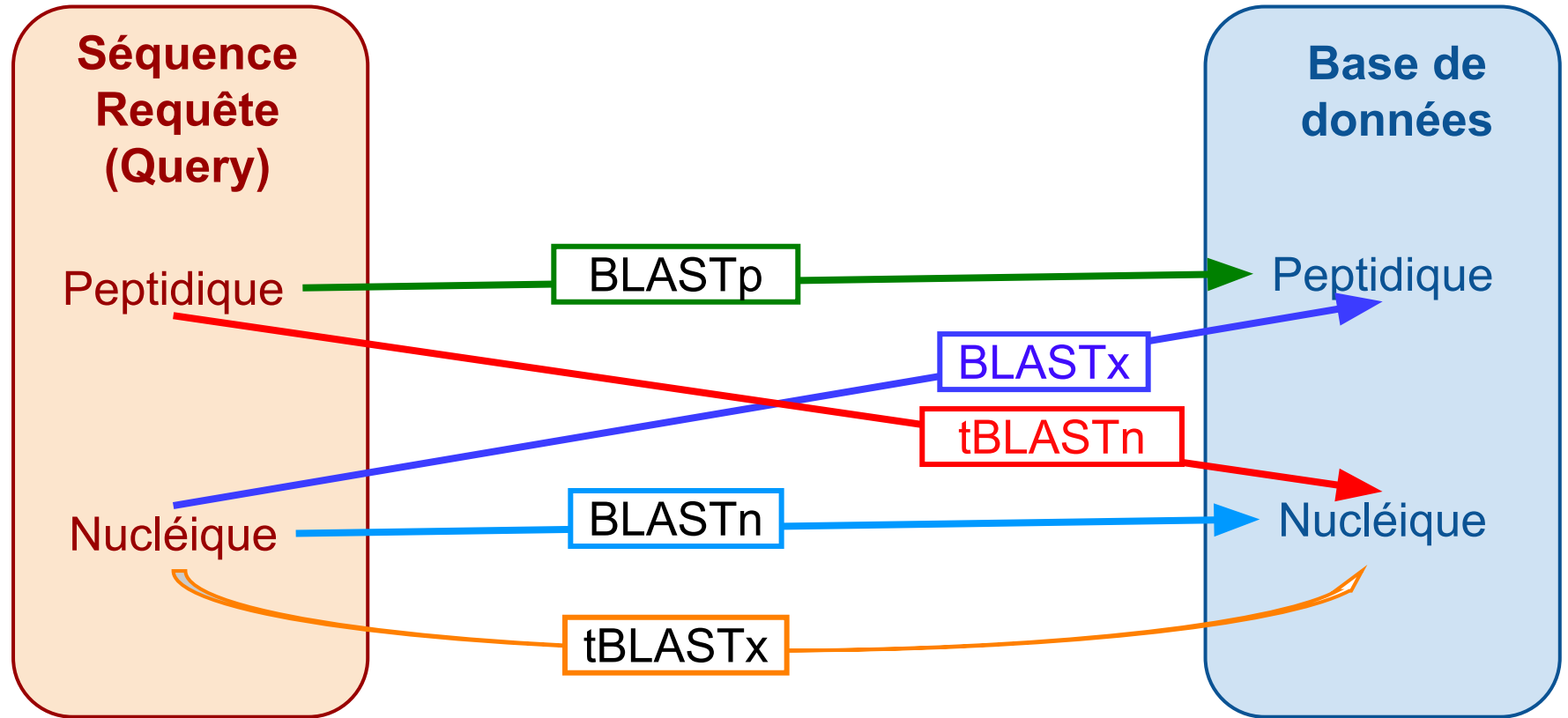
BLAST - a family of purpose-specific programs

- Different program names exist, depending on the type (protein or nucleic acid) of query and database sequences.
- For comparison between nucleic acids and proteins, the nucleic acid is translated in the 6 frames (3 frames per strand)



Query	Database	Program	Application examples	Study cases
protein	protein	blastp	Starting from a protein of known function detect putative homologs in the whole Uniprot database.	Collect sequences similar to the blue-sensitive opsin in all human proteins.
nucleic acid	nucleic acid	blastn	Match RNAi against a genome. Match mRNA (or EST) against a genome.	
nucleic acid (translated)	protein	blastx	After having sequenced a piece of DNA, search potentially coding fragments + their putative homologs without any prior knowledge of gene positions in the query sequence.	
protein	nucleic acid (translated)	tblastn	- Identify a genomic region likely to code for an homolog of a protein of interest. - Identify pseudo-genes (defective genes, with many stop codons) for a protein of interest in a genome.	Do cats see colors ? Get Human blue-sensitive opsin protein, connect UCSC genome browser, use BLAT to find similarities in Cat genome
nucleic acid (translated)	nucleic acid (translated)	tblastx		

Les modalités de BLAST



Comment interpréter la similarité entre séquences de HIV et de SARS-CoV-2 ?

Alignement de séquences de SARS-CoV-2 sur le génome du HIV

Haut: fragment le plus significatif de l'alignement de la séquence du gène S sur le génome du VIH. Noter le score Expect = 7.5. Ce score n'est significatif que s'il est nettement inférieur à 1.

Bas: fragment le plus significatif de l'alignement d'une séquence aléatoire sur le génome du VIH. Noter le score Expect = 2.1, supérieur à 1 et donc non-significatif (comme on s'y attend, puisque la séquence est aléatoire).

Conclusion: l'alignement sur lequel s'appuient Perez et Luc Montagnier correspond à ce qu'on s'attend à trouver par hasard en alignant des séquences de cette taille.

HIV-1 isolate 19828.PPH11 from Netherlands envelope glycoprotein (env) gene, partial cds				
Sequence ID: HQ644953.1		Length: 1143	Number of Matches: 1	Range 1: 967 to 994
Score	Expect	Identities	Gaps	Strand
38.3 bits(41)	7.5	25/28(89%)	0/28(0%)	Plus/Plus
Query	86	AATGGTACTAAGAGGTTTGATAACCCTG	113	
Sbjct	967	AATGGTACTAAAAGGTTAGATAACACTG	994	

HIV-1 isolate patient B clone 16.3 from Netherlands envelope glycoprotein (env) gene, complete cds				
Sequence ID: HQ386166.1		Length: 2580	Number of Matches: 1	Range 1: 2493 to 2523
Score	Expect	Identities	Gaps	Strand
39.2 bits(42)	2.1	27/31(87%)	0/31(0%)	Plus/Minus
Query	351	CCTAAAAGTTCTTTGTAATAACTGTATTATT	381	
Sbjct	2523	CCTAAAAGTTCTTTGTAATATTTCTATAATT	2493	

Travaux pratiques

Nous mobiliserons une série d'outils bioinformatiques accessibles en ligne pour analyser les séquences de coronavirus et pour tenter de trouver des éléments informatifs concernant l'origine de SARS-CoV-2.

Supports de ce cours	Diapos, tuto, données	https://ivanheld.github.io/shnc-origines-sars-cov-2/
Uniprot	Base de donnée de séquences protéiques	https://www.uniprot.org/
NCBI Entrez	Bases de données biologiques	https://www.ncbi.nlm.nih.gov/
EMBOSS needle	Alignement de paires de séquences	https://www.ebi.ac.uk/Tools/psa/emboss_needle/
NCBI BLAST	Recherche de séquences par similarité	https://blast.ncbi.nlm.nih.gov/Blast.cgi
PIPprofiler	Profils de pourcentages de positions identiques	https://pipprofiler.france-bioinformatique.fr/
Clustal	Alignement de séquences multiples	https://www.ebi.ac.uk/Tools/msa/clustalo/
phylogeny.fr	Phylogénie moléculaire	https://www.phylogeny.fr/
AMU	page AMETICE de N&C3	https://ametice.univ-amu.fr/course/view.php?id=62928