

BIN702

ALGORITHMES POUR LA

BIOINFORMATIQUE

Manuel Lafond



Qu'est-ce que la bioinformatique?

- Utilisation de méthodes informatiques pour approfondir notre compréhension des processus biologiques.

Qu'est-ce que la bioinformatique?

- Utilisation de méthodes informatiques pour approfondir notre compréhension des processus biologiques.
 - **Comparaison** de séquences biologiques
 - Découverte de **nouveaux gènes**
 - **Séquençage** et assemblage de génomes
 - Analyse d'**expression** des gènes
 - Prédiction de **structure des ARN** et des **protéines**
 - Prédiction de **fonctions** des gènes
 - Reconstruction de **l'évolution des espèces**
 - Recherche **d'interactions** entre protéines
 - ...

Qu'est-ce que la bioinformatique?

- Utilisation de méthodes informatiques pour approfondir notre compréhension des processus biologiques.
- Combine des expertises en **biologie**, en **informatique**, en **mathématiques/statistiques**.

Qu'est-ce que la bioinformatique?

- Utilisation de méthodes informatiques pour approfondir notre compréhension des processus biologiques.
- Combine des expertises en **biologie**, en **informatique**, en **mathématiques/statistiques**.
- Applications à l'agriculture, la pharmacologie, la médecine, la virologie, etc.

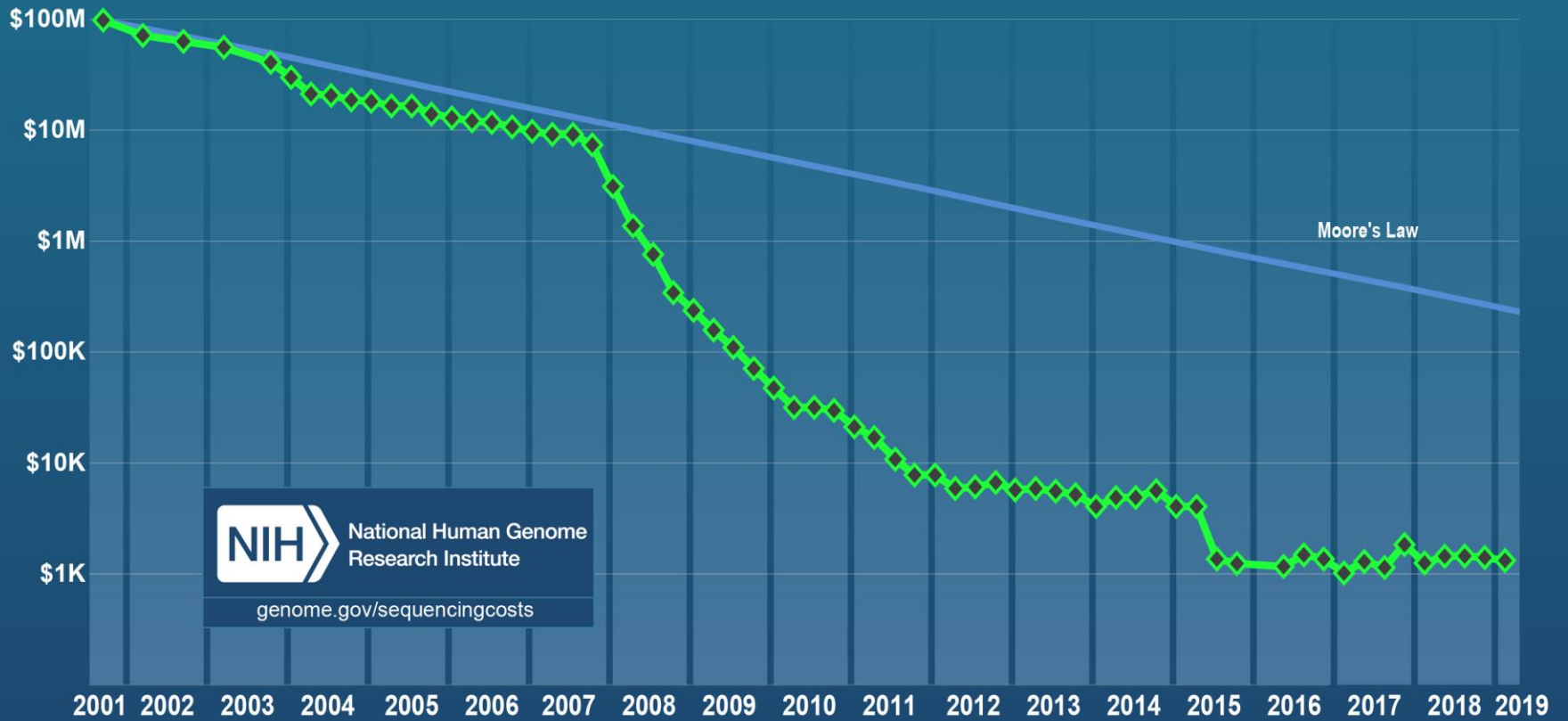
Micro-historique

- Années 1950, Frédérick Sanger détermine la séquence des acides aminés de l'insuline.
- En 1965, Margaret Dayhoff: premier atlas de séquences de protéines.
- De 1990 à 2003: *Human Genome Project* : séquençage de l'ADN de l'homme. Coûts estimés: \$2.7 milliards.

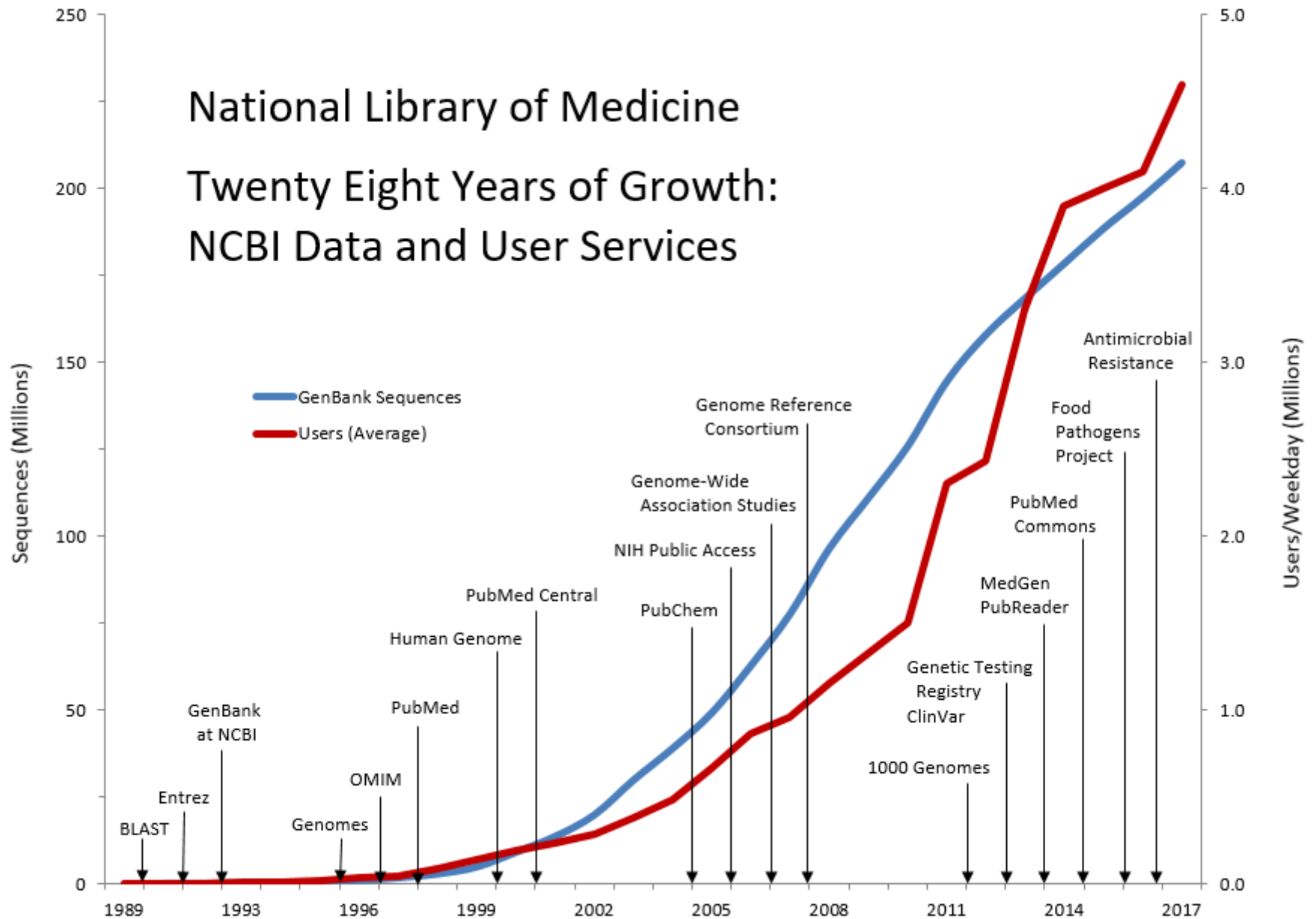
Micro-historique

- Années 1950, Frédérick Sanger détermine la séquence des acides aminés de l'insuline.
- En 1965, Margaret Dayhoff: premier atlas de séquences de protéines.
- De 1990 à 2003: *Human Genome Project* : séquençage de l'ADN de l'homme. Coûts estimés: \$2.7 milliards.
- Black Friday de 2018: Dante Labs peuvent séquencer votre génome pour \$200.

Cost per Genome



National Library of Medicine Twenty Eight Years of Growth: NCBI Data and User Services



Des données, des données

- Croissance **exponentielle** du nombre de séquences de gènes, d'ARN et de protéines.
- Cette croissance ne va probablement pas s'arrêter bientôt
 - de plus en plus d'espèces sont séquencées
 - génomes d'individus différents => multiplication des séquences
- Uniprot:
 - 167,761,270 séquences sur 1,175,367 espèces
 - 56,302,024,981 acides aminés
- GenBank:
 - 213,383,758 séquences, 329,835,282,370 paires de bases

Le spectre de la bioinformatique



Le spectre de la bioinformatique

Utilisation de logiciels
Requêtes à des
bases de données

Développement de
logiciels, implémentation
d'algorithmes connus,
pipelines, analyse de
données

Modélisation de problèmes,
création de nouveaux
algorithmes, implémentation
et essais sur données

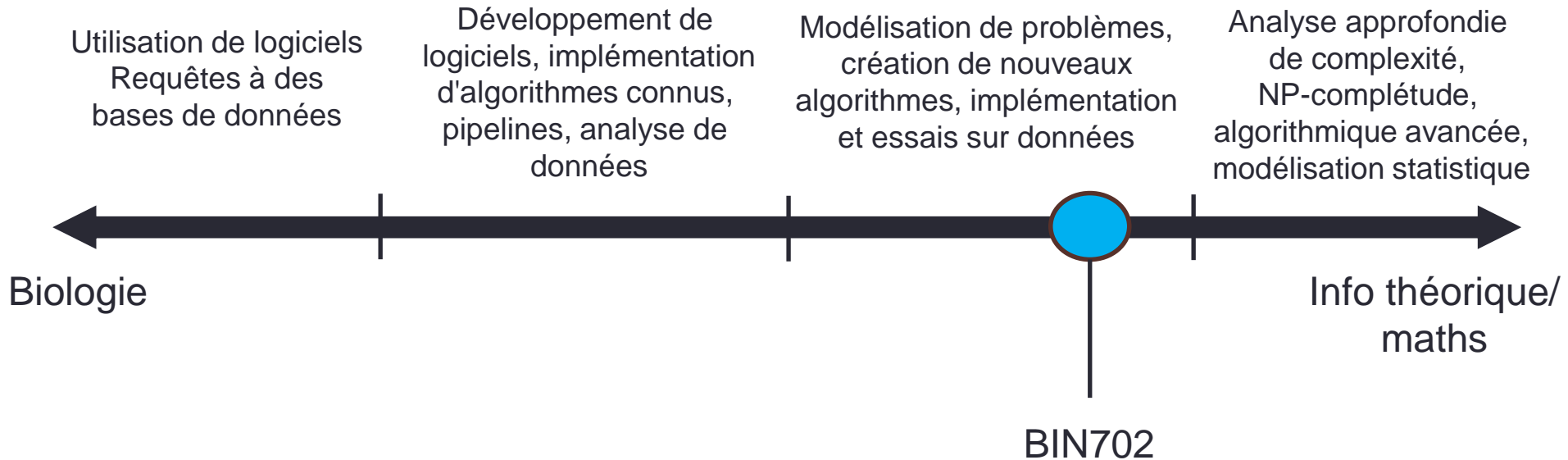
Analyse approfondie
de complexité,
NP-complétude,
algorithmique avancée,
modélisation statistique



Biologie

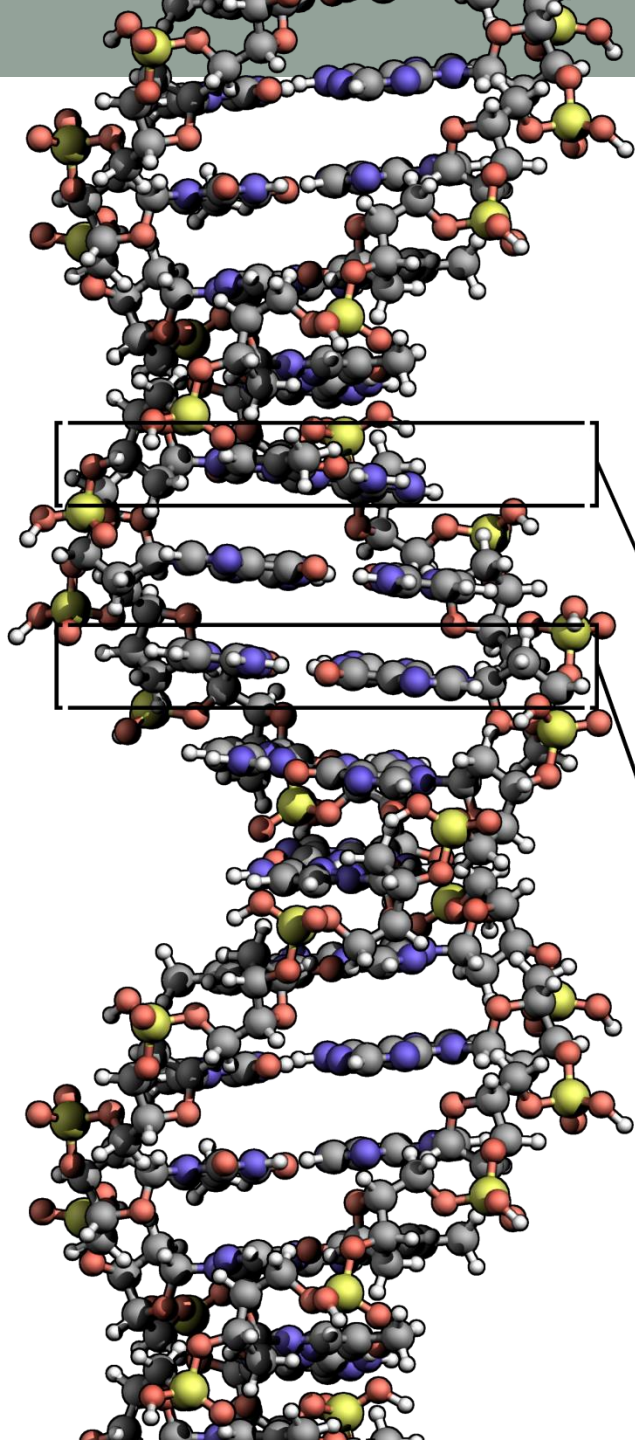
Info théorique/
maths

Le spectre de la bioinformatique

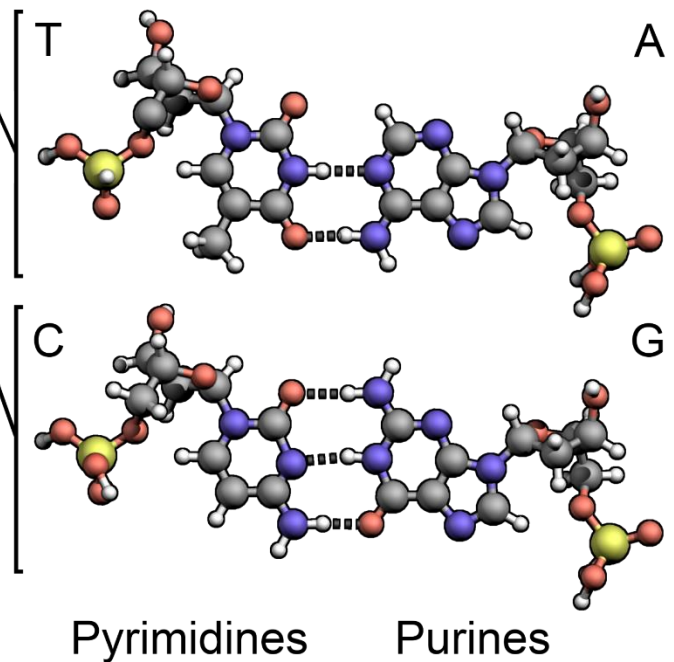


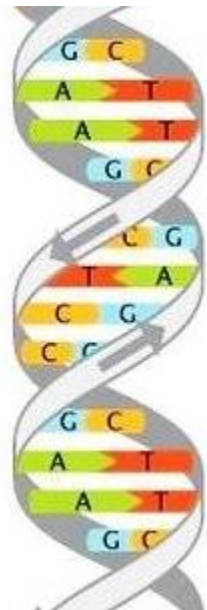
Minor groove

Major groove

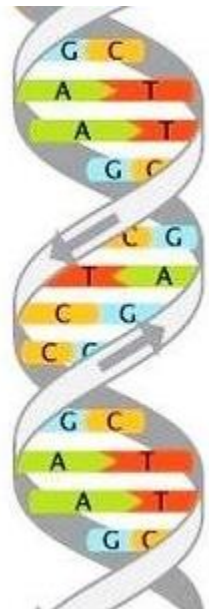


- Hydrogen
- Oxygen
- Nitrogen
- Carbon
- Phosphorus





- Pour nous informaticiens, l'ADN n'est qu'une chaîne de caractères sur un alphabet de 4 lettres.
 - (plus tard, nous verrons que les protéines sont des séquences sur un alphabet de taille 20)
- Cette conversion "entité bio" => "concept informatique" est très fréquente.



```
> seq1  
GAAGCTCCGAAG...
```


>>agp1

>S.cerevisiae

CTTATGTTAATACTATTACTACTATCATTACATTTTAATATCCTGTACATATATACTCTTTGAAGAAGTTTGGATTTTTCTTTCCCATTATAGATCACTT
TTTTTTTTACCTTATCTTTCTGTAAGTTCTTTTGTGGCTATATTTCTCTTTGGCGGCTAACTCAAGAAGCAGTAAGGGAATAATTTCGTTACACATTA
CAGCTAGCTAAAATCATCAAGAGGCAGGAACACACAGAGATATACTGATAATGTTGGAATGCTAAGCGCTGCTGGTTCTCTTTGCCGGTGGTGGT
GGGAGCGTTCTCGCAGCGGTTCAAAGCAAAACTGTGGCAGATCCTAACTTATGTCCAGGTTACAATTACAACCTGATTTCCCCTTCTCTCCAGTT
GCGTATGTACAACAAATTTACCGAGGTACTAATATTCTCTAGATTTTATTACTGACAAGATCCATGTGAAGAAACGAAATTTAAGTGAATGCGTGTCCCG
TTATTTTGATGAACAATATGCATTTTGTCCGGTCATGCGTTACTGTGAATAATGAAACA

>S.paradoxus

CCTATGTTAATACTATCATTACATTAGTATCCTGTACATATATACTCTATTGGAAGAAGTCTGAATTTTTCCCCTTATGGATCATCTTTTTCTATTTTTT
TTTTTTTTGCTAGTCTTTTTGTGGCTATATTTCTCTTTGGCGGCTAACTTAAGAAGTAGCATGGGAATAATTTGTTACATATTGCAGTTATCGAAAGT
CATCAAGAAGCTGGAACACACGAAATACACTAATAATGTTGGGATGTTTAAGCGCTCTGCTGGTTCTCTTTGGCTGGCGGGGGTAGGAGCATTCTCG
CGGCGGTTCAAAGTAAAACCGTGACAGATCCTAACTTATGCCAGGTTACAATTCACAACCTAATTTCCCCTTCTCTTTCTAGTTGCGTATGTACAATA
ACTTGCTGAGGTACTAATTTAATCTTTAAATTTATTACTGACAAGATCCACTTGAAGAAACGAGATTTGAATGAATGCATATTCATTATTTGATGAACA
ATATGCATTTTGTCCGGTCATGTGTTACTGTGAATAACGAAACA

>S.mikatae

CTTACGTTAATACTGTTATTATAATTACTATTACTACATTAGTATCTTGTACATATATACTCTACTTGAAGAAGTGAATATTTCTATATCCTGTTATGGTT
TAGAATTTTTTTCTTTGATCATTTCTTTTTGGTTAGTCATTTTTGTGGTCATATCGTCTCTTTGGCGGCTAACTTGAAATGTCGAGGGAGACACCT
GTTATATATTAGAACTACCAAAATTTATCAAGAACTGGAACACAAGAAATATACCAGCAATGTTGGGATGTTTAAGCACACTGCTGGTTCCTTAGCT
GGCGGAGGTAGTTTGATTCTCCAGCAGTTCAAAGCAAACCGTGTGAGATCCTAACTTATGTCCAGGTTACAATTCACGAAGAGTTTCTCCCTTCC
TTCCAGTTGTGTATGTACAATGAACCTACTGATATACTAGTGAATTTCTAGATTTTATTACTAACAAGAGCCACGTTTAGAAAACAAATCTGAATAAA
GC

>S.kudriavzevii

CCACAGCCATTGATTGCCCAAGCAGCTTTCTTTTTCTTCTCCATAACAGCAGCCCAACAGATAAACACCCGAAGTTATGCCCAATCGAAATGCATTAA
CGACATCAACAACGATCCTAATACGATGATCTTTCATCCTCGATTCCGTAACCATGTCTTTTCGTTGTTTTTTTCCATACGAAATCGTTTGCTAAAGAC
AGTTGATAGCATCCGTTGATGTATCTTCTAGAAAAACAATTCTGCAGAAGTGTTATTTTATCTCTTTTTTATAGCAAGCGGCACGGCCTCGGCGGCT
AACTCCTTAAGGCCTCGGCGGCTAAAAACGGGCTCTTCCATTTACACAGACTCCGTACCTCCCACGGCATTTTTTGCATCTTTTTTTTTTTTTCTTG
AGACATAATAGAAGGAAGGGGGAACGCGCGCCATTAGACGGCGCCGCTCAGAAGGCGCCGCTAAGCGGCACCGATCTGCTGCGAAACATGCAG
CTTTGCCGTTAATGGCGCTCCAGTGGCACCATCGCAGTAATTGTGGCCTGATCGACCCATAGATAAGCTGCCGCTACCACCAGGCAGTTCTCAAT
GCCCCACAGTCCACTTGATTGCTGCTCTTTGGTAGTCACACGCTTCTCGTTTTAGCTGCCGTAAGGCTCTTTTTCTTCTTTCTTCTGCTTTTTTTA
GTGCCTTCTGTCCATCTCTCCAGCTTGGCAATGGCTATCTTTTCACTAAATCGAAAATCAGTGCTGATAAGATGGCTCAGGTGCATGGGATAGGTCG
GCATACATCTGTTTTACGCTATCTTCTCCTCGTCCGCTAGCGGACTCACTTATTCTGAATAATTTCTTTTTTTTTCAAGAGTAACTTTGTATTTTACA
TACGTATATAAGTGGTTTGTCTGGATTTGCTGCTTTACAAGGGTGAGGAATAGCATTGTCATGTATTTTTTTTTTGAACGCAAGTAGTATAGACTA
AACAAAAAGCTTCGCATA

>>spo16

>S.cerevisiae

GGATCGTATATTTACCAGTATGCTACCTCCTTGAATCTTGCCCAATAACCTAATGACTCCATTGAATATACAATATAGAACATTGTGATGATTTATTTA
TTGGCAAAAAATGTGACGCGAAAAACAACATTAGATGAGTAAAAAGAAAAGTACATCTGATATAGCTTACCCTGTTAACGCGTAAAAAGTGGGCG
GCTAAAACCGAGAAAATACGAAATAGTGTAAAGGTTCAAGAGAATAAACTTCTGGAACAACATTTCTGGAATGGTTCAATAGAA

>S.paradoxus

GGAGTATATTTGACCAGTACACTACCTCCTTGAATCTTATCCGATGACCCAATTACTTACATTGAATACAACCTATAGAACATCTATGATGATTTATTT
ATTGGCAAAAAATGTGACGCGAAAAACAGAACATTAGGTGAGTAAAAAGAAAAGTACATCTGATATAGCTTACCCTGTTAACGCGTAAAAAGTGG
GCGGCTAAAACCGGAAAAATATGAAAGAGTAAAGGTTCAAGAGAATAAACTTCTGGGAGAACGTTTCTCGAATGCTTAACTGGAC

>S.mikatae

AGCTTATATCTCGACCAGAACACTATTTCTTGGAAATCTTACTTGGAGTGTTAATAACCTTGCATTGAGTATAAACTATAGCATGTATATGATGATTTATT
TATTAGCAGAGAAATGTGACGCGAAAAATCCAACACTTCAGTGAGTCAAAAATAAAACCGTCTAAATTTCTTACCCTGTTGACGCGTGAATAAATG
TGGGCGGCTAAAACCAAAAAATCTTGAACAGGATAAAGAGAAGATGAGCTACAAGTATAAAGCTACAGAAAACCGCTTGTGAATAAACCGGAATATAC

>S.kudriavzevii

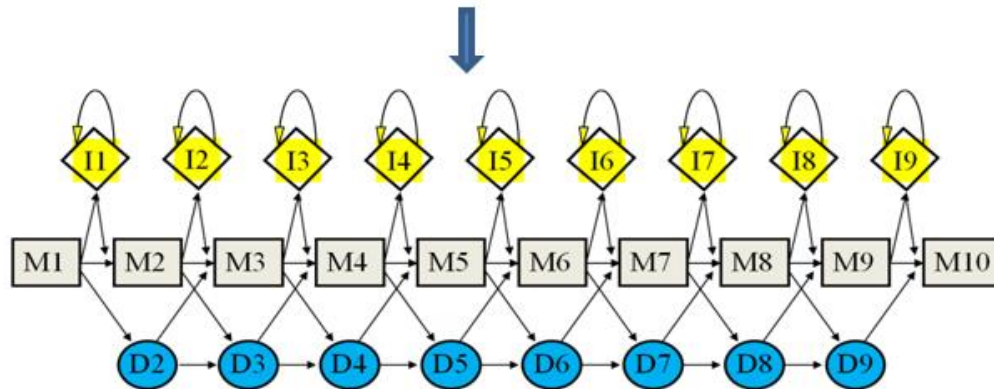
Alignement multiple de séquences

Q5E940_BOVIN	-----MPREDRATWKSNYFLKIIQLLDDYPKCFIVGADNVGSKOMQOIRMSLRGK-AVVLGMGKNTMMRKAIRGHLENN--PALE	76
RLA0_HUMAN	-----MPREDRATWKSNYFLKIIQLLDDYPKCFIVGADNVGSKOMQOIRMSLRGK-AVVLGMGKNTMMRKAIRGHLENN--PALE	76
RLA0_MOUSE	-----MPREDRATWKSNYFLKIIQLLDDYPKCFIVGADNVGSKOMQOIRMSLRGK-AVVLGMGKNTMMRKAIRGHLENN--PALE	76
RLA0_RAT	-----MPREDRATWKSNYFLKIIQLLDDYPKCFIVGADNVGSKOMQOIRMSLRGK-AVVLGMGKNTMMRKAIRGHLENN--PALE	76
RLA0_CHICK	-----MPREDRATWKSNYFMKIIQLLDDYPKCFVVGADNVGSKOMQOIRMSLRGK-AVVLGMGKNTMMRKAIRGHLENN--PALE	76
RLA0_RANSY	-----MPREDRATWKSNYFLKIIQLLDDYPKCFIVGADNVGSKOMQOIRMSLRGK-AVVLGMGKNTMMRKAIRGHLENN--SALE	76
Q7ZUG3_BRARE	-----MPREDRATWKSNYFLKIIQLLDDYPKCFIVGADNVGSKOMQOIRMSLRGK-AVVLGMGKNTMMRKAIRGHLENN--PALE	76
RLA0 ICTPU	-----MPREDRATWKSNYFLKIIQLLNDYPKCFIVGADNVGSKOMQOIRMSLRGK-AIVLMGKNTMMRKAIRGHLENN--PALE	76
RLA0_DROME	-----MVRENKAAWKAQYFIKVVLFDFEFPKCFIVGADNVGSKOMQOIRMSLRGL-AVVLGMGKNTMMRKAIRGHLENN--PQLE	76
RLA0_DICDI	-----MSGAG-SKRKKLFIEKATKLFITTYDKMIVAEADTVGSQLOKIRKSIRGI-GAVLMGKKTMRKVIIRDLDLADSK--PELD	75
Q54LP0_DICDI	-----MSGAG-SKRKNVFIKATKLFITTYDKMIVAEADTVGSQLOKIRKSIRGI-GAVLMGKKTMRKVIIRDLDLADSK--PELD	75
RLA0_PLAF8	-----MAKLSKQKKQMYIEKLSL IQQYSKILIVHVDNVGSNQMASVRKSLRGK-ATIILMGKNTIRRTALKKNLQAV--PQIE	76
RLA0_SULAC	----MIGLAVTTTKKIAKWKVDEVAELTEKPKLTKHTIIIANIEGFPADKLHEIRKKLRGK-ADIKVTKNLNFNIALKNAG----YDTK	79
RLA0_SULTO	----MRIMAVITQERKIAKWKIEEVKELEOKLREYHTIIIANIEGFPADKLHDIRKMRGM-AEIKVTKNTLFGIAAKNAG----LDVS	80
RLA0_SULSO	----MKRLALALKQRKVASWKLKEVKELELTKNSNTILIGNLEGFPADKLHEIRKKLRGK-ATIKVTKNTLFGIAAKNAG----IDIE	80
RLA0_AERPE	MSVVSIVGQMYKREKPIPEWKTLMLELEELFSKHRVVLFDLGTPTFVVQVRKKLWKK-YPMVAKKRIILRAMKAAGLE---LDDN	86
RLA0_PYRAE	-MMLAIGKRRYVTRQYPARKVKIVSEATELLQKYPYVFLFDLHGLSSRILHEYRYRLRRY-GVIKIKPTLFGIAFTKVYGG---IPAE	85
RLA0_METAC	-----MAEERHTEHIPQWKKDEIENIKELIQSHKVFGMVIEGILATKMKIRRDLDKV-AVLKVSRLNTLTERALNQLG----ETIP	78
RLA0_METMA	-----MAEERHTEHIPQWKKDEIENIKELIQSHKVFGMVRIEILATKMKIRRDLDKV-AVLKVSRLNTLTERALNQLG----ESIP	78
RLA0_ARCFU	-----MAAVRGS--PPEYKVRAVEEIKRMISKPVVAIVSFRNVPAGOMKIRREFRGK-AEIKVKNTLLEALDALG----GDYL	75
RLA0_METKA	MAVKAKGQPPSGYE PKVAEWKRREVKEKELMDEYENVGLVDLEGIPAPQLQEIIRAKLRERDTIIRMSRLNTLMRIALEEKLDER--PELE	88
RLA0_METHH	-----MAHVAEWKKKEVQELHDLIKGYEVVGIANLADIPARQLQKMRQTLRDS-ALIRMSKKTLLISLAEKAGREL--ENVD	74
RLA0_METTL	----MITAESEHKIAPWKIEEVNKLKELKNGQIVALVDMMEVPARQLQEIIRDKIR-GTMTLKMSRLNTLIERAIKEVAEETGNPEFA	82
RLA0_METVA	----MIDAKSEHKIAPWKIEEVNALKKELLSANVIALIDMMEVPARQLQEIIRDKIR-DQMTLKMSRLNTLIERAIKEVAEETGNPEFA	82
RLA0_METJA	----METKVKAHVAPWKIEEVNKLKELIKSKPVVAIVDMMDVPAPQLQEIIRDKIR-DKVKLRMSRLNTLIERAIKEVAEELNNPKLA	81
RLA0_PYRAB	-----MAHVAEWKKKEVEELANLKSYPVIALVDVSSMPAYPLSQMRLIRENGLLRVSRLNTLIELAIKKAQELGKPELE	77
RLA0_PYRHO	-----MAHVAEWKKKEVEELAKLKSYPVIALVDVSSMPAYPLSQMRLIRENGLLRVSRLNTLIELAIKKAQELGKPELE	77
RLA0_PYRFU	-----MAHVAEWKKKEVEELANLKSYPVIALVDVSSMPAYPLSQMRLIRENGLLRVSRLNTLIELAIKKAQELGKPELE	77
RLA0_PYRKO	-----MAHVAEWKKKEVEELANIKSYPVIALVDVAGVPAYPLSKMRDKLR-GKALLRVSRLNTLIELAIKRAAQELGQPELE	76
RLA0_HALMA	----MSAESERKTETIPEWKQEEVDIVEMIESYESVGVVNIAGIPSRQLQDMRRDLHGT-AELRVSRLNTLLEALDDVD----DGLE	79
RLA0_HALVO	----MSESEVRQTEVIPQWKREEVDELVDFIESYESVGVVGVAGIPSRQLQSMRRE LHGS-AAVRMSRLNTLVNRALEEVN----DGFE	79
RLA0_HALSA	----MSAEEQRTTEEVPEWKQEEVAELVDLLETYSVGVVNVGTIPSKLQDMRRGLHGO-AALRMSRLNTLLVRALEEAG----DGLD	79
RLA0_THEAC	-----MKEVSOQKKELVNEITARIKASVAIVDTAGIRTRQIQDIRGKNRGK-INLKVIKKTLFLKALENLGD----EKLS	72
RLA0_THEVO	-----MRKINPKKKEIVSELAQDITKSKAVAVDIKGVIRTRQIQDIRAKNRDK-VKIKVKKTLFLKALDSIND----EKLT	72
RLA0_PICTO	-----MTEPAQWKIDFVKNLENEINSRKVAIVS IKGLRNNFQKIRNSIRDK-ARIKVSARLLRLAIENTGK----NNIV	72
ruler	1.....10.....20.....30.....40.....50.....60.....70.....80.....90	

Modélisation par Modèle de Markov

Multiple sequence alignment

Sequence 1:	F	K	L	L	S	H	C	L	L	V
Sequence 2:	F	K	A	F	G	Q	T	M	F	Q
Sequence 3:	Y	P	I	V	G	Q	E	L	L	G
Sequence 4:	F	P	V	V	K	E	A	I	L	K
Sequence 5:	F	K	V	L	A	A	V	I	A	D
Sequence 6:	L	E	F	I	S	E	C	I	I	Q
Sequence 7:	F	K	L	L	G	N	V	L	V	C



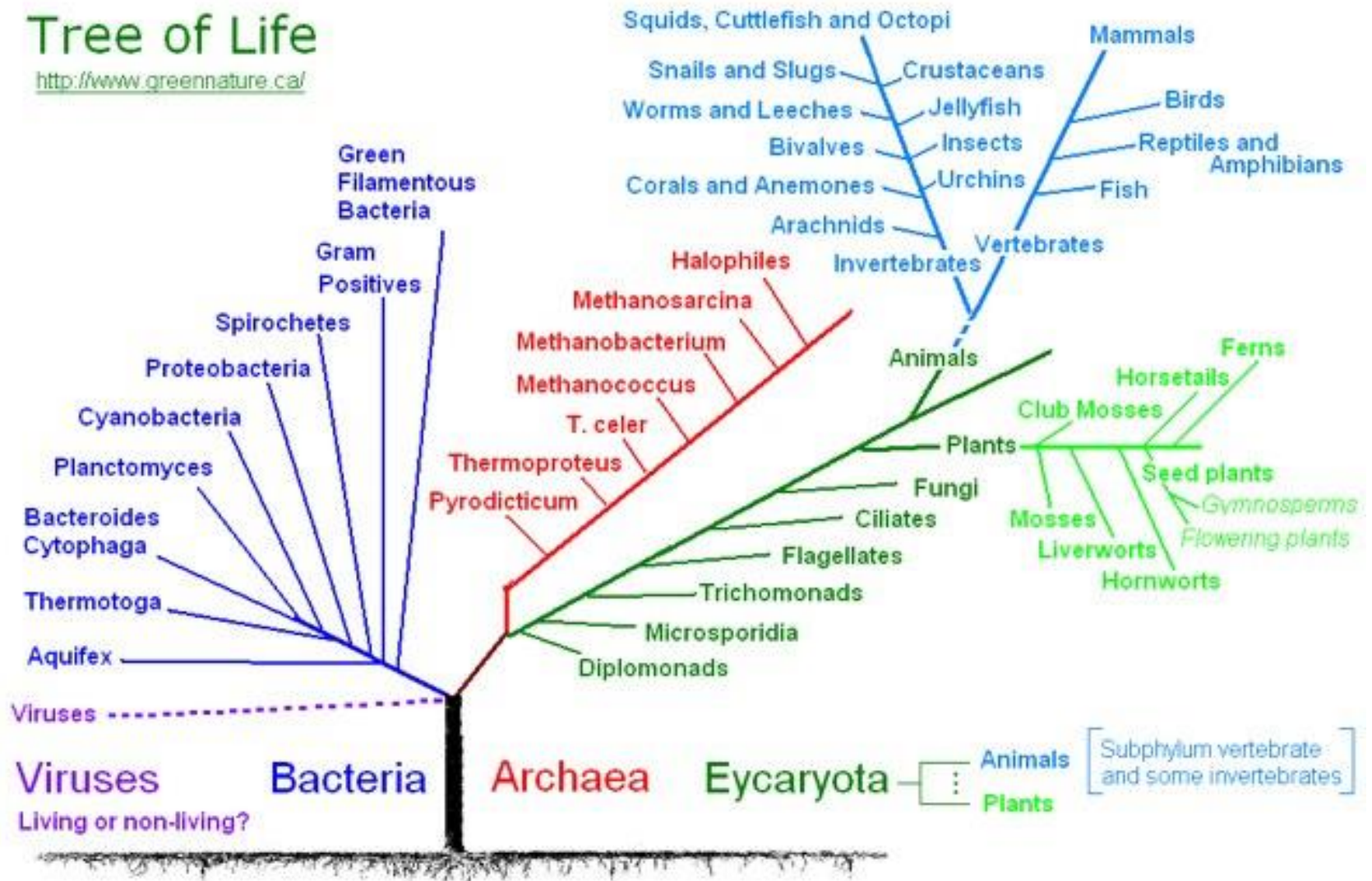
I = insert state

M = match state

D = delete state

Tree of Life

<http://www.greennature.ca/>



Qu'est-ce que BIN702?

- Nous nous intéressons aux fondements des méthodes utilisées par les biologistes et bioinformaticien(ne)s.
 - Quels algorithmes sont implémentés? [IFT436]
 - Quelle est leur complexité? [IFT339/IFT436]
 - Quelles structures de données sont utilisées? [IFT339]

Qu'est-ce que BIN702?

- Nous nous intéressons aux fondements des méthodes utilisées par les biologistes et bioinformaticien(ne)s.
 - Quels algorithmes sont implémentés? [IFT436]
 - **Préalable souhaité**: concepts de base en algorithmique (approches diviser-pour-régner, programmation dynamique, algorithmes gloutons, graphes, ...)
 - Quelle est leur complexité? [IFT339/IFT436]
 - **Préalable requis**: savoir analyser le temps d'un algorithme en notation O
 - Quelles structures de données sont utilisées? [IFT339]
 - **Préalable requis**: concepts de base en structures de données (listes, tableaux, arbres binaires de recherche, tables de hachage)

Qu'est-ce que BIN702 n'est pas?

Qu'est-ce que BIN702 n'est pas?

- Nous n'utilisons **pas d'outils concrets** - ce sont les méthodes derrière ces outils qui nous intéressent.
 - des algorithmes des années 1970 sont encore très utilisés
 - très peu de logiciels ont une durée de vie de plus de 10 ans
 - compétences utiles dans d'autres domaines
- Nous ne discutons **pas de sources de données** biologiques (ex: Uniprot, Ensembl, Cancer Genome Atlas, ...).
- Nous n'entrons pas dans des **détails d'implémentation**.
- Nous n'utilisons **pas l'IA ni l'apprentissage-machine** appliqué à la bioinformatique.
 - Quelques exceptions peuvent survenir.
 - Vous pouvez explorer le côté pratique au cours de votre **projet**.

Projet en BIN702

- Vous avez un maximum de liberté!
 - Une banque d'idée sera proposée
 - Si vous amenez votre propre idée, le correcteur sera content :)
- Projet pratique
 - Implémenter un algorithme qui n'existe qu'en pseudo-code
 - Améliorer les performances d'un algorithme connu (ou du moins, essayer)
 - Tester plusieurs logiciels accomplissant la même tâche (*benchmarking*)
 - Développer un algo + programme pour un nouveau problème
 - Développer un *pipeline* pour une tâche commune en bioinformatique

Projet en BIN702

- Vous avez un maximum de liberté!
 - Une banque d'idée sera proposée
 - Si vous amenez votre propre idée, le correcteur sera content :)
- **Projet théorique**
 - Résoudre un problème ouvert
 - Revue approfondie de littérature sur un sujet donné
 - Améliorer la complexité à un problème déjà étudié
 - Montrer qu'un problème est NP-complet

Projet en BIN702

- Rapport de projet
 - Mise en contexte - où se place votre projet par rapport à la littérature? Quelle est sa pertinence?
 - Qu'est-ce qui rend votre projet original?
 - Quel est son niveau de difficulté?
 - Quels défis devez-vous surmonter?
 - Quelles solutions avez-vous apportées?
 - Projets pratiques: comment évaluez-vous votre approche? Quels résultats attendiez-vous, et qu'avez-vous obtenu?
 - Projets théoriques: si vous n'avez pas résolu votre problème, quelles pistes avez-vous tentées? Pourquoi ont-elles échoué?