

Centre de recherche

INRIA Rennes – Bretagne Atlantique

L'INRIA

Le Centre

La bio-informatique

Le bio-informaticien

Pierre Peterlongo
Mai 2011
Lycée Descartes Rennes

INSTITUT NATIONAL
DE RECHERCHE
EN INFORMATIQUE
ET EN AUTOMATIQUE



centre de recherche
RENNES - BRETAGNE ATLANTIQUE

Les instituts de recherche en France

Etablissements publics à caractère scientifique et technologique (EPST) [EMARPEF](#) Centre national du machinisme agricole, du génie rural, des eaux et des forêts

CNRS Centre national de la recherche scientifique
INED Institut national d'études démographiques
INRA Institut national de la recherche agronomique
INRETS Institut national de recherche sur les transports et leur sécurité
INRIA Institut national de recherche en informatique et en automatique
INSERM Institut national de la santé et de la recherche médicale
INSTITUT Institut de recherche pour le développement
IFSTTAR Laboratoire central des ponts et chaussées

Etablissements publics à caractère industriel et commercial (EPIC) **ADEME** Agence de l'environnement et de la maîtrise de l'énergie

ADIT Agence pour la diffusion de l'information technologique
ANDRA Agence nationale de gestion des déchets radioactifs
BREVI Bureau de recherches géologiques et minières
CEA Commissariat à l'énergie atomique
INRAE Centre national d'études spatiales
INRS Cité des sciences et de l'industrie
IFSTTAR Centre scientifique et technique du bâtiment
IFP Institut français du pétrole
IFREMER Institut français de recherche pour l'exploitation de la mer
INERIS Institut national de l'environnement industriel et des risques
IRSN Institut de radioprotection et de sûreté nucléaire
ONERA Office national d'études et de recherches aérospatiales

OSSEO Anvar, ex Agence nationale de valorisation de la recherche)

Etablissements publics à caractère administratif (EPA)

CEA Centre d'études de l'emploi
CNRS Centre informatique national de l'enseignement supérieur
CNRS Centre national de documentation pédagogique (*Etablissement rattaché à l'Enseignement scolaire*)
CNRS Centre national d'enseignement à distance (*Etablissement rattaché à l'Enseignement scolaire*)
CNRS Centre National des Œuvres Universitaires et Scolaires
CE Groupe des Ecoles des Télécommunications (Recherche en technologie de l'information et de la communication)
INRS Institut national de recherche pédagogique
EPA **Assas** : désamiantage, mise en sécurité et rénovation du site
INRAE Institut national de recherches archéologiques préventives **Fondations**
CEPH Centre d'étude du polymorphisme humain

INRAE Institut Pasteur de Paris

INRAE Institut Pasteur de Paris

INRAE Institut Pasteur de Paris

INRAE Institut Pasteur de Paris

INRAE Institut Pasteur de Paris

INRAE Institut Pasteur de Paris

INRAE Institut Pasteur de Paris

INRAE Institut Pasteur de Paris

INRAE Institut Pasteur de Paris

INRAE Institut Pasteur de Paris

INRAE Institut Pasteur de Paris

INRAE Institut Pasteur de Paris

INRAE Institut Pasteur de Paris

INRAE Institut Pasteur de Paris

INRAE Institut Pasteur de Paris

INRAE Institut Pasteur de Paris

INRAE Institut Pasteur de Paris

INRAE Institut Pasteur de Paris

INRAE Institut Pasteur de Paris

INRAE Institut Pasteur de Paris

INRAE Institut Pasteur de Paris

INRAE Institut Pasteur de Paris

INRAE Institut Pasteur de Paris

INRAE Institut Pasteur de Paris

INRAE Institut Pasteur de Paris

INRAE Institut Pasteur de Paris

INRAE Institut Pasteur de Paris

INRAE Institut Pasteur de Paris

INRAE Institut Pasteur de Paris

INRAE Institut Pasteur de Paris

INRAE Institut Pasteur de Paris

INRAE Institut Pasteur de Paris

INRAE Institut Pasteur de Paris

INRAE Institut Pasteur de Paris

INRAE Institut Pasteur de Paris

INRAE Institut Pasteur de Paris

INRAE Institut Pasteur de Paris

INRAE Institut Pasteur de Paris

INRAE Institut Pasteur de Paris

INRAE Institut Pasteur de Paris

INRAE Institut Pasteur de Paris

INRAE Institut Pasteur de Paris

INRAE Institut Pasteur de Paris

INRAE Institut Pasteur de Paris

INRAE Institut Pasteur de Paris

INRAE Institut Pasteur de Paris

INRAE Institut Pasteur de Paris

INRAE Institut Pasteur de Paris

INRAE Institut Pasteur de Paris

INRAE Institut Pasteur de Paris

INRAE Institut Pasteur de Paris

INRAE Institut Pasteur de Paris

INRAE Institut Pasteur de Paris

<http://www2.enseignementsup-recherche.gouv.fr>

Les instituts de recherche en France

- CNRS: un peu tout
- INED: Démographie
- IRD: Développement
- CNES: Etudes spatiales
- IFREMER: Mer
- INRA: Agronomie
- INRIA: Informatique
- 82 Universités...

<http://www2.enseignementsup-recherche.gouv.fr>

Etablissements publics à caractère scientifique et technologique (EPST) IFREMER Centre national du machinisme agricole, du génie rural, des eaux et des forêts

CNRS Centre national de la recherche scientifique

INED Institut national d'études démographiques

INRA Institut national de la recherche agronomique

INRETS Institut national de recherche sur les transports et leur sécurité

INRIA Institut national de recherche en informatique et en automatique

INSERM Institut national de la santé et de la recherche médicale

INSTITUT Institut de recherche pour le développement

IFSTTAR Laboratoire central des ponts et chaussées

Etablissements publics à caractère industriel et commercial (EPIC) ADEME Agence de l'environnement et de la maîtrise de l'énergie

ADIF Agence pour la diffusion de l'information technologique

ANDRA Agence nationale de gestion des déchets radioactifs

BRETAGNE Bureau de recherches géologiques et minières

CEA Commissariat à l'énergie atomique

IRIA Centre de coopération international en recherche agronomique

IRISA Centre national d'études spatiales

INRAE Cité des sciences et de l'industrie

INRS Centre scientifique et technique du bâtiment

IFP Institut français du pétrole

IFREMER Institut français de recherche pour l'exploitation de la mer

INRS Institut national de l'environnement industriel et des risques

IRSN Institut de radioprotection et de sûreté nucléaire

ONIS Office national d'études et de recherches aérospatiales

OSSEO (OSSEO Anvar, ex Agence nationale de valorisation de la recherche)

Etablissements publics à caractère administratif (EPA)

CEA Centre d'études de l'emploi

CNRS Centre informatique national de l'enseignement supérieur

CNRS Centre national de documentation pédagogique (Etablissement rattaché à l'Enseignement scolaire)

CNRS Centre national d'enseignement à distance (Etablissement rattaché à l'Enseignement scolaire)

CNRS Centre National des Œuvres Universitaires et Scolaires

GET Groupe des Ecoles des Télécommunications (Recherche en technologie de l'information et de la communication)

INRS Institut national de recherche pédagogique

EPA : désamiantage, mise en sécurité et rénovation du site

INSTITUT Institut national de recherches archéologiques préventives **Fondations**

CEPH Centre d'étude du polymorphisme humain

INSTITUT Institut Pasteur de Paris

INSTITUT Institut Pasteur de Lille

Groupements d'intérêt public (GIP) ANRS Agence nationale de la recherche sur le sida

CNRS Consortium national de recherche en génomique (dont fait partie le Centre national de génomique) avec :

GENOSCOPE Centre national de séquençage

CNRS Réseau national des génopoles

CNRS Centre National de Recherche sur les Sites et Sols Pollués

INRAE Institut polaire français Paul-Emile Victor

INRAE GIP consacré à la recherche en génomique et au développement d'entreprises de biotechnologies

OS Observatoire des sciences et techniques

RSN Réseau national pour la technologie, l'enseignement et la recherche

Etablissement professionnel

Etablissements d'enseignement supérieur et de recherche 82 Universités (liste complète: <http://www.education.gouv.fr/sup/univ.htm>)

Instituts nationaux polytechniques

Grenoble <http://www.inpg.fr>

Lorraine Nancy <http://www.inpl-nancy.fr>

Toulouse <http://www.inpl-toulouse.fr>

4 Ecoles normales supérieures

Ecole normale supérieure <http://www.ens.fr>

Ecole normale supérieure de Cachan <http://www.ens-cachan.fr>

Ecole normale supérieure Lettres et Sciences humaines <http://www.ens-lyh.fr>

Ecole normale supérieure de Lyon <http://www.ens-lyon.fr>

5 Ecoles françaises à l'étranger

Casa de Velazquez <http://www.casafrvelazquez.org>

Ecole française de Rome <http://www.efa-rome.org>

Ecole française archéologie athènes <http://www.efa.gr>

Ecole Française d'Extrême-Orient <http://www.efeo.org>

Institut français d'archéologie orientale du Caire IFAO <http://www.ifao.org>

14 Grande Etablissements de statuts divers

CNAM (Conservatoire National des Arts et Métiers) <http://www.cnam.fr>

Collège de France <http://www.college-de-france.fr>

Ecole Centrale des Arts et Manufactures <http://www.ecp.fr>

Ecole Nationale des Chartes <http://www.chartes.fr>

(ENSAM) Ecole Nationale Supérieure d'Arts et Métiers <http://www.enscm.fr>

Ecole Nationale Supérieure des Sciences de l'Information et des Bibliothèques <http://www.ensib.fr>

(EPHE) Ecole Pratique des Hautes Etudes <http://www.ephe.fr>

(EHESP) Ecole des Hautes Etudes en Sciences Sociales <http://www.ehess.fr>

(GET) Groupe des écoles des Télécommunications <http://www.get.fr>

(IEP) Institut d'Etudes Politiques de Paris <http://www.sciences-po.fr>

(INALCO) Institut National des Langues et Civilisations Orientales <http://www.inalco.fr>

(Muséum National d'Histoire naturelle) <http://www.museum.org>

Observatoire de Paris <http://www.obspm.fr>

Palais de la Découverte <http://www.palais-decouverte.fr>

(Groupe des Ecoles des Télécommunications (GET) (établissement public placé sous la tutelle de la ministre déléguée à l'industrie: <http://www.get.fr>))

Etablissements publics à caractère scientifique, culturel et professionnel (EPSCP)

Ecole centrale de Lille <http://www.ec-lille.fr>

Ecole centrale de Lyon <http://www.ec-lyon.fr>

Ecole centrale de Nantes <http://www.ec-nantes.fr>

Ecole nationale des Ponts et Chaussées <http://www.enpc.fr>

Ecole nationale supérieure des arts et industries de Strasbourg <http://www.ensai.u-strasbg.fr>

Institut de Physique du Globe de Paris <http://www.ipg.fr>

Institut national des sciences appliquées de Lyon <http://www.inria-lyon.fr>

Institut national des sciences appliquées de Rennes <http://www.inria-rennes.fr>

Institut national des sciences appliquées de Toulouse <http://www.inria-toulouse.fr>

Institut national des sciences appliquées de Rouen <http://www.inria-rouen.fr>

Institut supérieur des matériaux et de la construction mécanique <http://www.cmi.fr>

Université de technologie de Compiègne <http://www.utc.fr>

Université de technologie de Belfort-Montbéliard <http://www.univ-belfort.fr>

Université de technologie de Troyes <http://www.univ-troyes.fr>

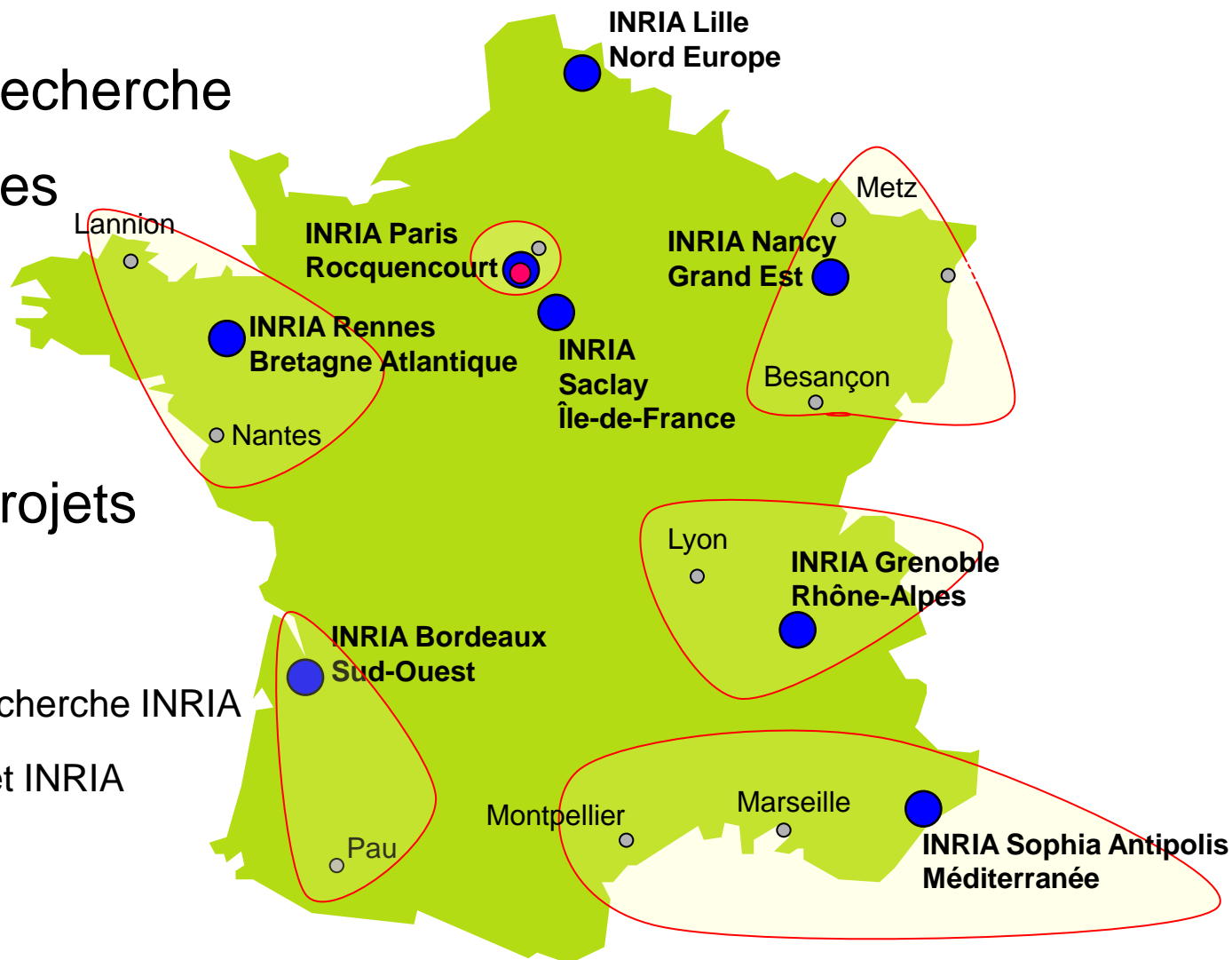
L'INRIA :

8 centres de recherche

3800 personnes

150 équipes-projets

- Siège
- Centre de recherche INRIA
- Équipe-projet INRIA hors site



Chiffres : mai 2008

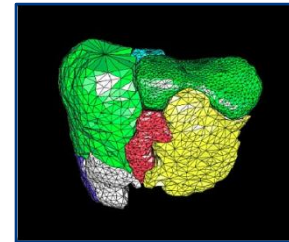
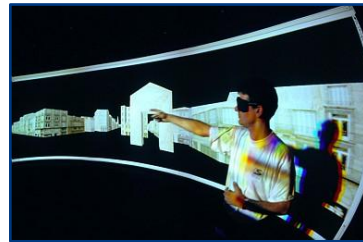
L'INRIA

Ministère de la recherche et ministère de l'industrie



chercher

créer la référence



expérimenter

partager & transmettre

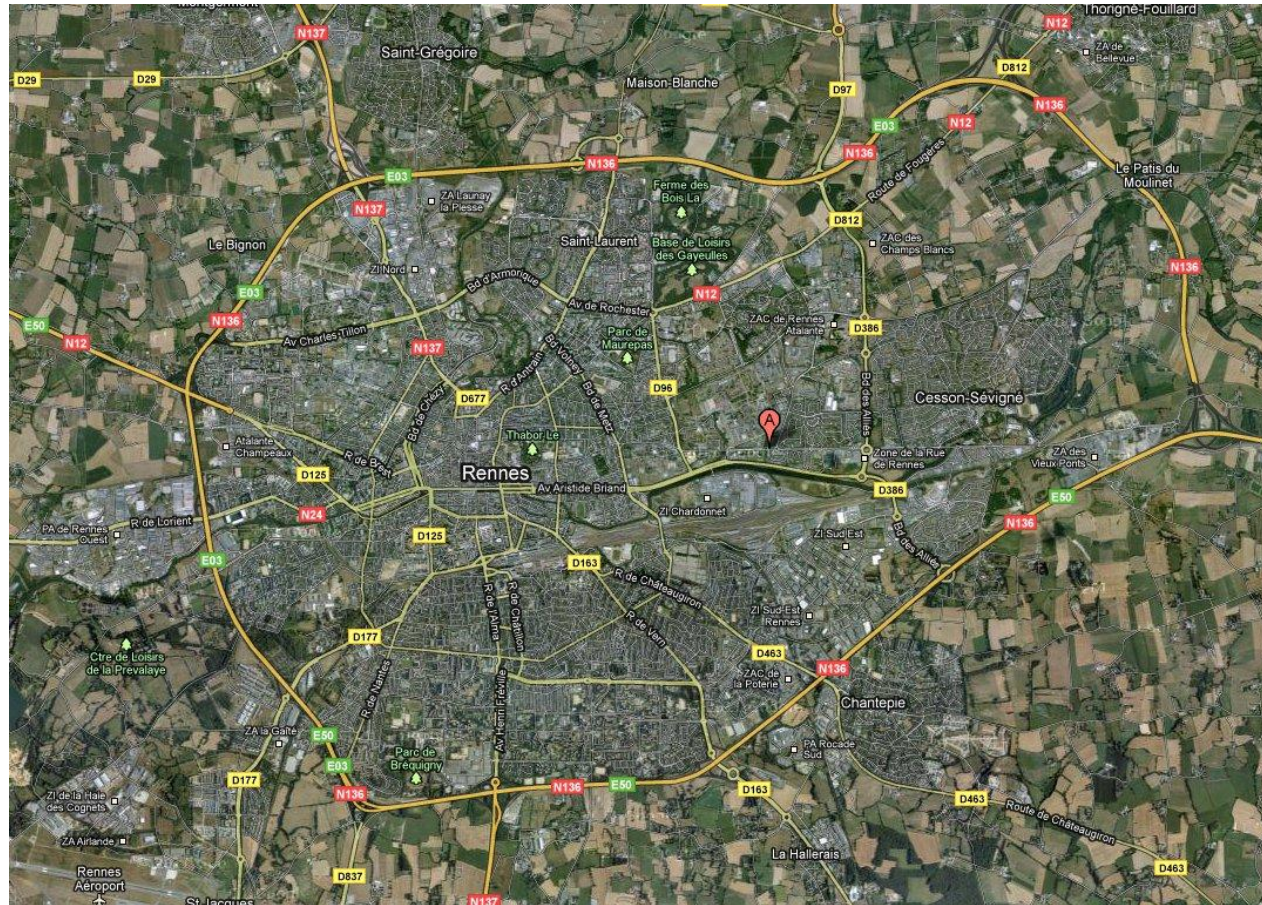


Le centre de recherche INRIA Rennes – Bretagne Atlantique

Présentation générale

Où ?

- Campus de beaulieu, Rennes.



Où ?

- Campus de beaulieu, Rennes.



Qui ?

Environ 600 personnes:

- 200 (235) Chercheurs et enseignants chercheurs
- 100 (75) Ingénieurs, techniciens et administratifs



Quoi ?

Sciences du **numérique**

32 équipes de recherche

- Réalité virtuelle
- Réseaux / internet
- Robotique
- Sciences du vivant
- ...



Les plate-formes d'expérimentation

Multimédia

- Système d'immersion virtuelle
- Plate-forme de capture, stockage et indexation de données vidéo

Robotique

- Robots avec capteurs de vision, véhicule Cycab
- Système robotisé d'échographie 3D

Informatique

- Un des 9 nœuds de la grille de calcul nationale Grid'5000
- Plate-forme bio-informatique de Ouest-génopole

Biomédical

- Plate-forme Neurinfo



© CNRS – Photothèque / Hubert Raguet

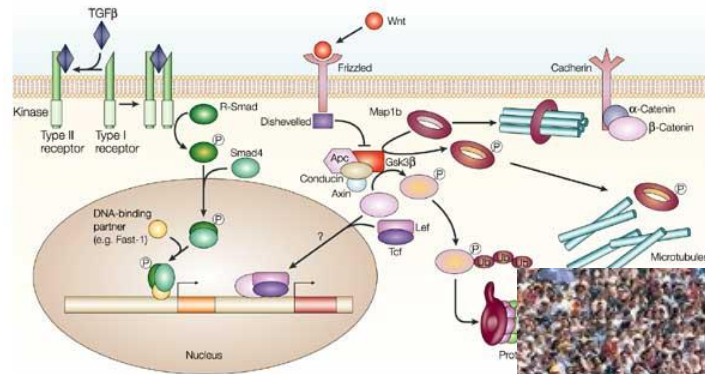


L'équipe projet Symbiose

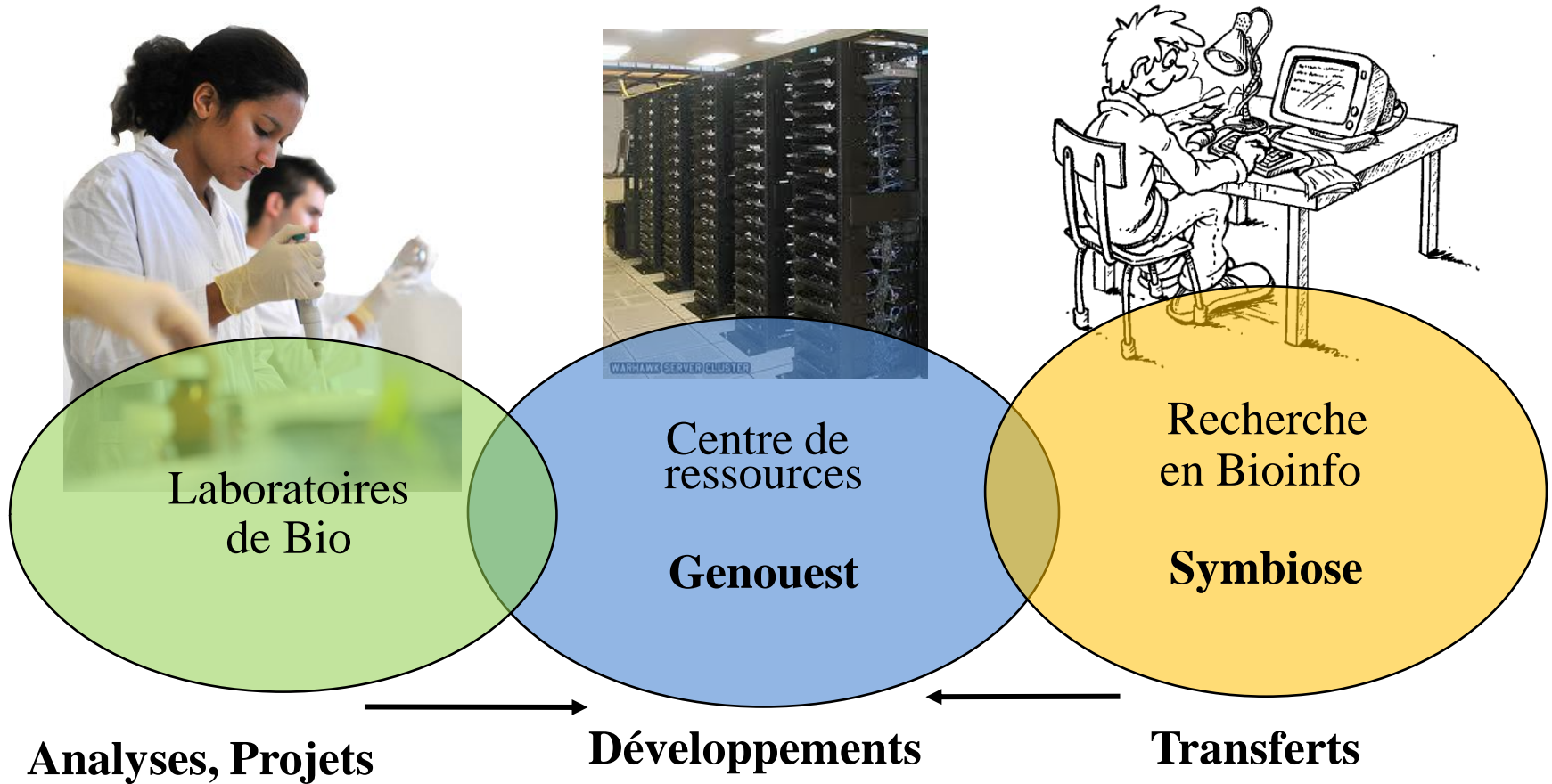


Symbiose, 3 axes de recherche

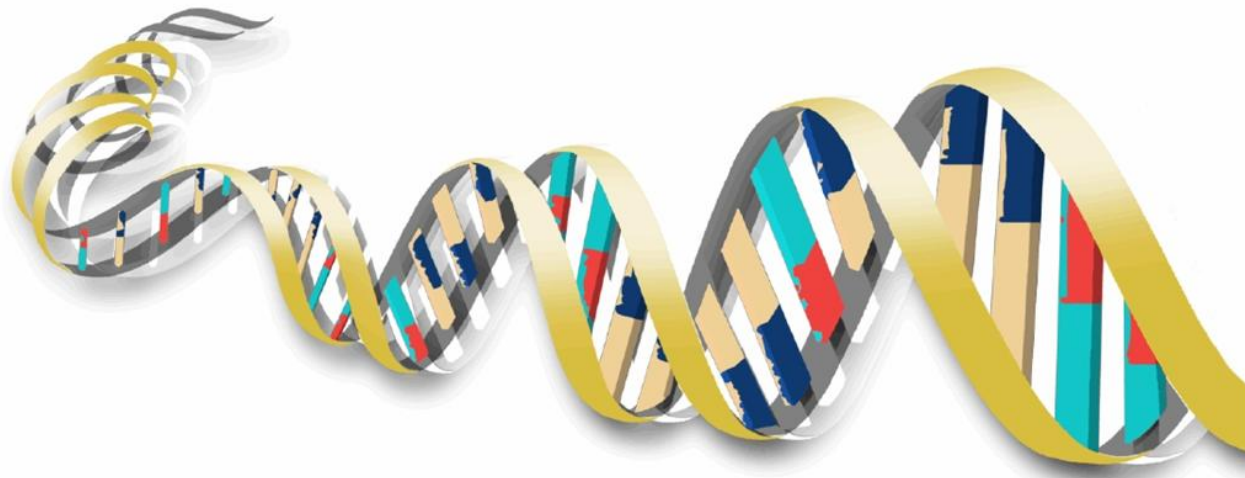
- Interactions
- Grosses masses de données
- Calculs intensifs



Symbiose: la recherche, la plateforme, les biologistes



C'est quoi la Bioinformatique ?



Quand ?

- Naissance : années 80 – 90
- Avant : expérimentations
 - macro ou micro scopique
 - Comportements

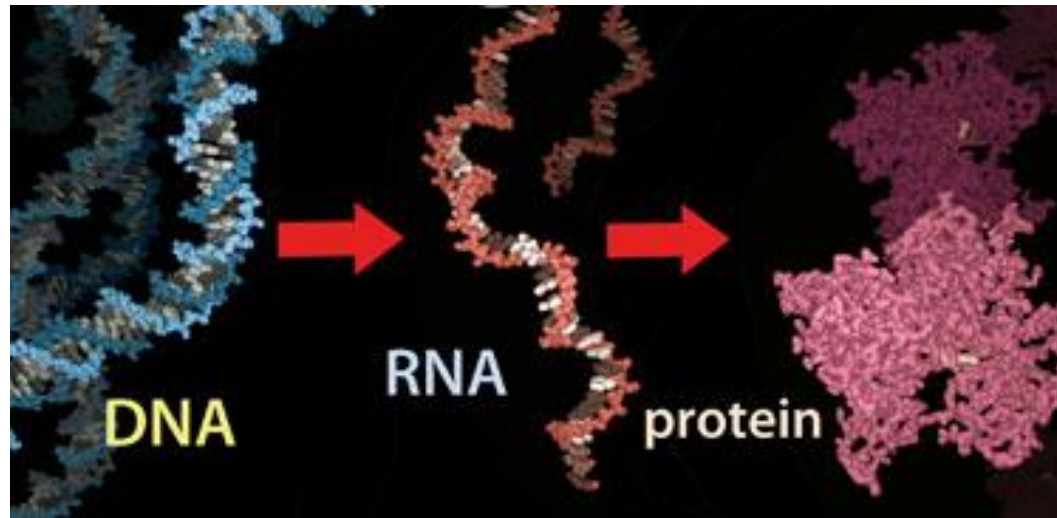


Quand ?

- Naissance : années 80 – 90
- Après : comme avant +
 - ADN, ARN, Acides Aminés...



L'ADN, et après ?

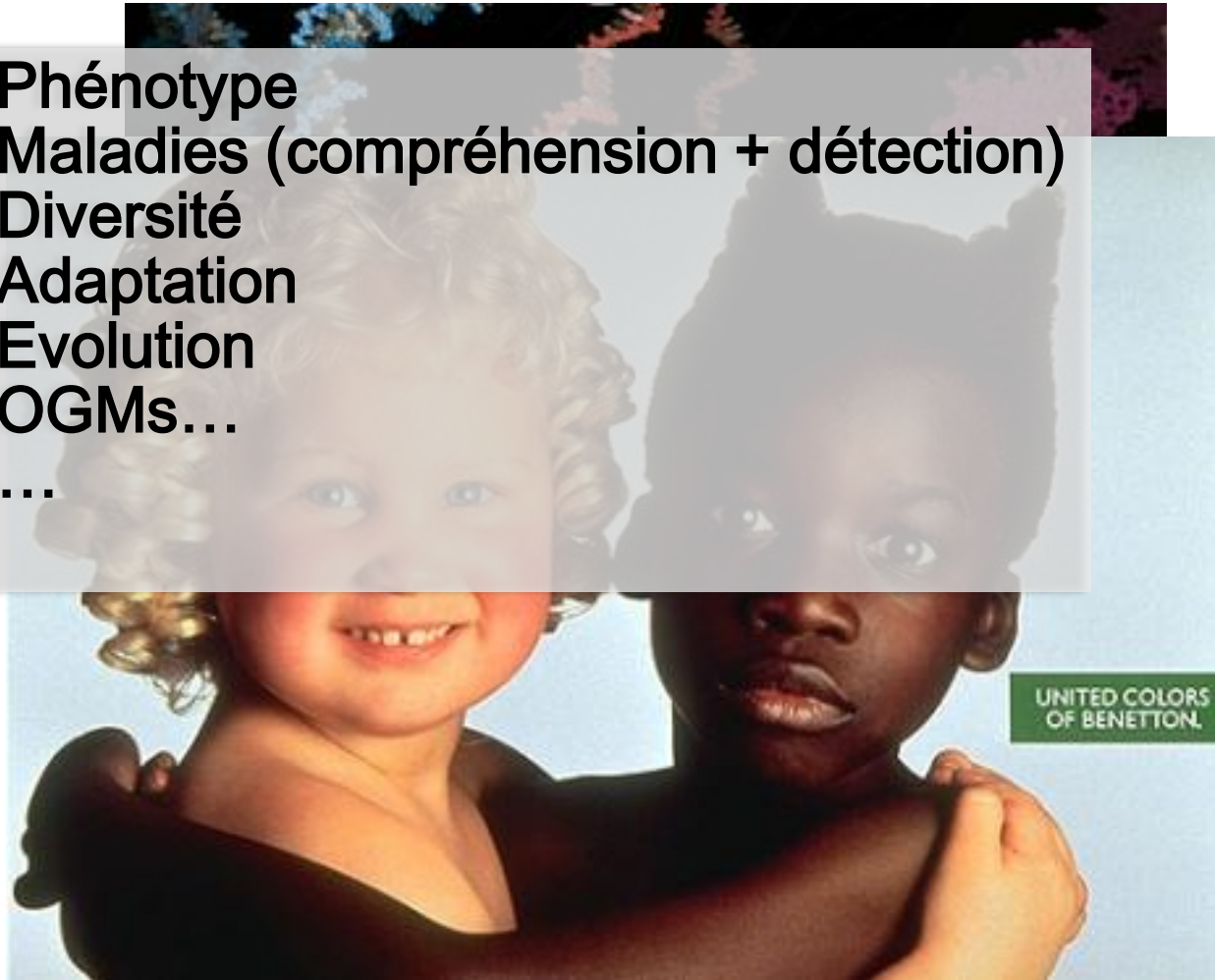


L'ADN, et après ?



L'ADN, et après ?

Phénotype
Maladies (compréhension + détection)
Diversité
Adaptation
Evolution
OGMs...
...



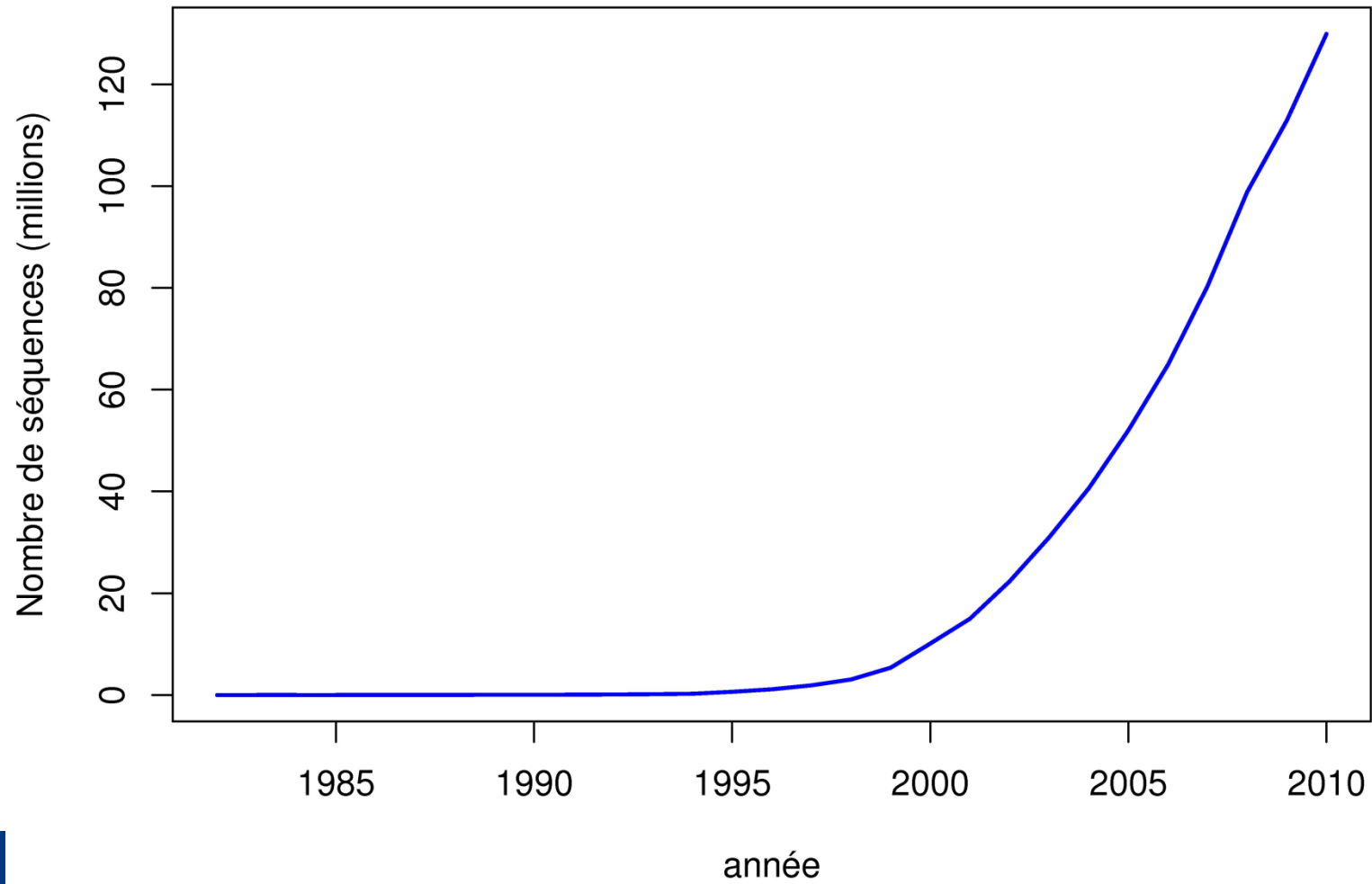
Lire l'ADN ?

- 1953 Découverte de la structure en double hélice de l'ADN
- 1977 Développement des techniques de séquençage d'ADN
- 1977 Premier génome séquencé : virus (5 386 nt)
- 1984 Génome du virus du sida (170 000 nt)
- 1990 Début du projet Génome humain
- 1996 Premier génome eucaryote (levure 12 000 000 nt)
- 2001 Premier génome de mammifère (homme 3 000 000 000 nt)
- 2011 >2600 génomes complets (200 mammifères)



Lire l'ADN ?

croissance de Genbank



Lire l'ADN = Séquencer

- Techno les plus connues:
 - 454: Genome Sequencer
 - Illumina Solexa Genome Analyser
 - Applied Biosystem: SOLiD



Lire l'ADN, pour quel prix ?

- Génome humain:
- 2000: 3 milliards \$



- Actuellement: 10000 \$



AACAAAAATGTTCTAAAATGGTCTTTTTTAAAACCTTTTTTAAATTAGAAAATATTACCATATAACCTTGTCAAAT
ACTATTAGTTTTTTTTTTTATAAACCTTCTTTTACATACACAGGAAAACAATGAACTAAGTTCTGAATGTATC
ATCTATAATACAAAAAGGAATTGTAGCTTATGGCTGGTATAATATTTTGCTGTTATTTCTCTAATTTCTACTC
CAAAATAGAGTCACAAGTTTTTTATAAAAACATCTTTGTACTCTAGGTTCCAACCTCTAACATGTAAAGAG
TTTTGAAATCATCAACCTTATCCTTACAACAGGAATAAAGCTGAAGACACTGAAAATCAACAACCTTTTCT
GGGACCCTTCGGACAACCTGCAGTCACAAGGTAAACTACCATTTTTTAAATATGGAGACAGAGGAATACAG
AGTACAACCTGACATCAGCTTACCTGGAGCAGGAACCACGGAACCTAAAAGTGGATGGATCACCTGAAT
GGTAATTCAGATAAATTGCTGTAAGTTGAGTGTGACATAGCTCAAGAATAAAAACCTCCTTGGAGTCACAG
TCTTAGAAGGGGTTCTGCATATTCAGAGGTTTTTACCTCCAAGAACCCACCAAGTTCTCATCATGAAGA
GCCAAGAAAAATCATCTCGTGTCTTTGGCAGGGGGAGGAGAAAAGTAACCCTTTGAAATATGACAGTT
TTTTTCTTGTAACAAAGTCCTATGCTCCACTGAAAAATATTTTATCAGAGCCTTATCTTATGTTGACAGAA
GGGCAGTTACCCAGTTCCAGCCACCCATAGCCTTTCTGTCTTACCTAAGTGGGGATGAAAAAAAAGGTG
AAAGCTCTTGTGAAAGTCACAGCCACATGGAAAGACCTGCTAAAAAACTGAAATTTAATCACAAAGATT
ACAGAATTCTCTTTTTCTCCATTCCTTGCTACCAAATCAATGGGCATCTAGTATAATAACAGTGGATTTGTA
GCCAAAAGATCTGCAAACCTTAGACTCTATTTAAGAATGAGCTCTTAGGGAAACCTAAAGACAACACCAG
AGAAAAAAACAAGGACATTAAGAATTGAAAGACTCTGACACATACAGCTACAGAAAACATTGGACA
CAGCCCAAATCCTAGCCAGGTTAAGATAACAACCTCACACTAAAGGCCTATGTGTCTCATTCTCTTTTCCA
AATATATGCATATAGCTTTCAATAACAACAAAAAACTCAAGGTATGCTAAAAAGTAGGAAAAACACAG
AAAATACAAAGCAAGCATTAAACCAAGACTCAGATTTGACAGAAATGTTGGACTTAGAGATGGAATGTA
AAATAACTATGGTTCATATTTAAGGGTCCTAATGGAAAAAAAGTAGACATTATGTAAAAACAGAGATAA
TGTAAGCAGAAATTTGACATTGAAAAATAGACACAAGAACCAAGTTAGGATTAATATTAGGGTCCACC
TGTTGCTAAAGCCCTAAATTGACCCTGGCTCATGTACAGATCTCCTTTAACTCCACTCAATCCTGCACCTT
GGGGCTGCATGTGTAGCAAGACTTGGGTCTATAGAATTGTACAAGATATAGTTAATTTGCAGTTGGACTA
GAGCAGCCCTGTGCTTATCTCTGGGGAACATGGAAATCTAGTGCTGACTCTTGGCCAAGATAAGGGTAG
TGGAATAAATAATCATTATAGAGTGCAATTTAACTTTGAAACTCCTATGGTGGCACTCAGGCCCAGCTT
AATTATTGATGGTTTTTACCAGCATAATTAGCTAGAAAGCCTTTGTTGTCTATATAATCAAACCTGTTAGGG
AGCTCCCCCTAAACCAAGGTTTCCCTGTACAAACCTTGTCATCTAATTTGAAGGCAAACCTATTTCTTTGTG
GGCCTAAGGATGACCCTGAAGAGCTATGTCTGGGAGAATCTAAGGGCTCTGACACTTGCTGGCATTATCT
TGCACCTATCTCTTGCACCTATCCTTTGGCTTGAATACTCTGTGGGATCTTAGAAGTCCTTTCAATGATC
TGATCAATGAACTTCATAAAAATAAATATTTTTTAGGCTAACTAATTGGAAGTTCACAGACTGTTCA



AACAAAAATGTTCTAAAATGGTCTTTTTTAAAACCTTTTTTAAATTAGAAAATATTACCATATAACCTTGTCAT
ACTATTAGTTTTTTTTTTATAAACCTTCTTTTACATACACAGGAAAACAATGAACTAAGTTCTGAATGTATC
ATCTATAATACAAAAAGGAATTGTAGCTTATGGCTGGTATAATTTTTGCTGTTATTTCTCTAATTTCTACTC
CAAATAGAGTCACAAGTTTTTTATAAAAACATCTTTGTACTCTAGGTTCCAACCTCTAACATGTAAAGAG
TTTTGAAATCATCAACCTTATCCTTACAACAGGAATAAAGCTGAAGACACTGAAAATCAACAACCTTTTCT
GGGACCCTTCGGACAACCTGCAGTCACAAGGTAAACTACCATTTTTTAAATATGGAGACAGAGGAATACAG
AGTACAACCTGACATCAGCTTACCTGGAGCAGGAACCACGGAACCTAAAAGTGGATGGATCACCTGAAT
GGTAATTCAGATATAGCTCAAGAATAAAAACCTCCTTGGAGTCACAG
TCTTAGAAGGGGTCCAAGAACCCCAAGTTCTCATCATGAAGA
GCCAAGAAAAATCACTCTCGTGTCTTTGGCAGGGGGAGGAGAAAAGTAACCACTTTGAAATATGACAGTT
TTTTTCTTGTAACAAAGTCCTATGCTCCACTGAAAAATATTTTATCAGAGCCTTATCTTATGTTGACAGAA
GGGCAGTTACCCAGTTCCAGCCACCCATAGCCTTTCTGTCTTACCTAAGTGGGGATGAAAAAAAAGGTG
AAAGATT
ACAGTGTA
GCCACAG
AGAACACA
CAGCCCA
AATAAG
AAAAGTA
AAATAACTATGGTTCATAITTTAAGGGTCCTAATGGAAAAAAGTAGACATTATGTAAAAACAGAGATAA
TGTAAGCAGAAATTTGACATTGAAAAATAGACACAAGAACCAAGTTAGGATTAATATTAGGGTCCACC
TGTTGCTAAAGCCCTAAATTGACCCTGGCTCATGTACAGATCTCCTTTAACTCCACTCAATCCTGCCTTT
GGGGCTGCATGTGTAGCAAGACTTGGGTCTATAGAATTGTACAAGATATAGTTAATTTGCAGTTGGACTA
GAGCAGCCCTGTGCTTATCTCTGGGGAACATGGAAATCTAGTGCTGACTCTTGGCCAAGATAAGGGTAG
TGGAATAAATAATCATTTATAGAGTGCAATTTAACTTTGAAACTCCTATGGTGGCACTCAGGCCAGCTT
AATTATTGATGGTTTTTACCAGCATAATTAGCTAGAAAGCCTTTGTTGTCTATATAATCAAACCTGTTAGGG
AGCTCCCCCTAAACCAAGGTTTCCCTGTACAAACCTTGTCATCTAATTTGAAGGCAAACCTATTTCTTTGTG
GGCCTAAGGATGACCCTGAAGAGCTATGTCTGGGAGAATCTAAGGGCTCTGACACTTGCTGGCATTATCT
TGCACTCTATCTCTTGCACTCTATCCTTTGGCTTGAATACTCTGTGGGATCTTAGAAGTCCTTTCAATGATC
TGATCAATGAACTTCATAAAAATAAATATTTTTTAGGCTAACTAATTGGAAGTTCACAGACTGTTCA

Ici: 2000 nucléotides

Génome humain: 3 milliards de nucléotides: X 1,5 Millions

AACAAAAATGTTCTAAAATGGTCTTTTTTAAAACCTTTTTTAAATTAGAAAATATTACCATATAACCTTGTCAA
ACTATTAGTTTTTTTTTTATAAACCTTCTTTTACATACACAGGAAAACAATGAACTAAGTTCTGAATGTATC
ATCTATAATACAAAAAGGAATTGTAGCTTATGGCTGGTATAATA
CAAAATAGAGTCACAAGTTTTTTATAAAAACATCTTTGTACT
TTTTGAAATCATCAACCTTATCCTTACAACAGGAATAAAGCT
GGGACCCTTCGGACAACCTGCAGTCACAAGGTAAACTACCAT
AGTACAACCTGACATCAGCTTACCTGGAGCAGGAACCACGG
GGTAATTCAGAT
TCTTAGAAGGGG
GCCAAGAAAAATCAICTCGTGTCTTTGGCAGGGGGGAGGAG
TTTTTCTTGTAACAAAGTCCTATGCTCCACTGAAAAATATTT
GGGCAGTTACCCAGTTCCAGCCACCCATAGCCTTTCTGTCTT
AAA
ACA
GCC
AGA
CAG
AATA
AAA
AAATAACTATGGTTCATAITTTAAGGGTCCTAATGGAAAAAA
TGTAAGCAGAAATTTGACATTGAAAAATAGACACAAGAACC
TGTTGCTAAAGCCCTAAATTGACCCTGGCTCATGTACAGATC
GGGGCTGCATGTGTAGCAAGACTTGGGTCTATAGAATTGTAC
GAGCAGCCCTGTGCTTATCTCTGGGGAACATGGAAATCTAG
TGGAATAAATAATCATTATAGAGTGCAATTTAACTTTGAA
AATTATTGATGGTTTTTACCAGCATAATTAGCTAGAAAGCCT
AGCTCCCCCTAAACCAAGGTTTCCCTGTACAAACCTTGTCAT
GGCCTAAGGATGACCCTGAAGAGCTATGTCTGGGAGAATCT
TGC ACTCTATCTCTTGC ACTCTATCCTTTGGCTTGAATACTCT
TGATCAATGAACTTCATAAAAATAAATATTTTTTAGGCTAACT

Ici: 2000 nucléotides

Génome humain: 3 milliards de nucléotides

* ~ 200 livres de 500 pages...



AACAAAAATGTTCTAAAATGGTCTTTTTTAAAACCTTTTTTAAATTAGAAAATATTACCATATAACCTTGTCAA
ACTATTAGTTTTTTTTTTTATAAACCTTCTTTTACATACACAGGAAAACAATGAACTAAGTTCTGAATGTATC
ATCTATAATACAAAAAGGAATTGTAGCTTATGGCTGGTATAATA
CAAAATAGAGTCACAAGTTTTTTATAAAAACATCTTTGTACT
TTTTGAAATCATCAACCTTATCCTTACAACAGGAATAAAGCT
GGGACCCTTCGGACAACCTGCAGTCACAAGGTAAACTACCAT
AGTACAACCTGACATCAGCTTACCTGGAGCAGGAACCACGG
GGTAATTCAGAT
TCTTAGAAGGGG
GCCAAGAAAAATCAICTCGTGTCTTTGGCAGGGGGGAGGAG
TTTTTCTTGTAACAAAGTCCTATGCTCCACTGAAAAATATTT
GGGCAGTTACCCAGTTCCAGCCACCCATAGCCTTTCTGTCT
AAA
ACAC
GCCA
AGA
CAG
AATA
AAA
AAATAACTAT
TGTAAG
TGT
GAAAAAA
SACACAAGAACC
GGCTCATGTACAGATC
TTGGGTCTATAGAATTGTAC
CTCTGGGGAACATGGAAATCTAG
TTATAGAGTGCAATTTAACTTTGAA
TTTTTCACCAGCATAATTAGCTAGAAAGCCT
CCCTAAACCAAGGTTTCCCTGTACAAACCTTGTCAT
GGCTAAGGATGACCCTGAAGAGCTATGTCTGGGAGAATCT
TGC ACTCTATCTCTTGC ACTCTATCCTTTGGCTTGAATACTCT
TGATCAATGAACTTCATAAAAATAAATATTTTTTAGGCTAACT

Ici: 2000 nucléotides

Génome humain: 3 milliards

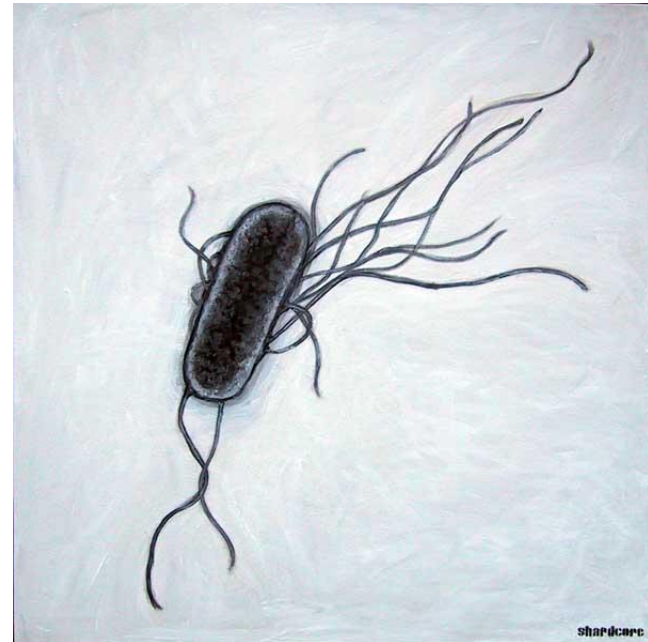
* ~ 200 livres de



Que faire de tout ça ???
Comment extraite de l'information ?

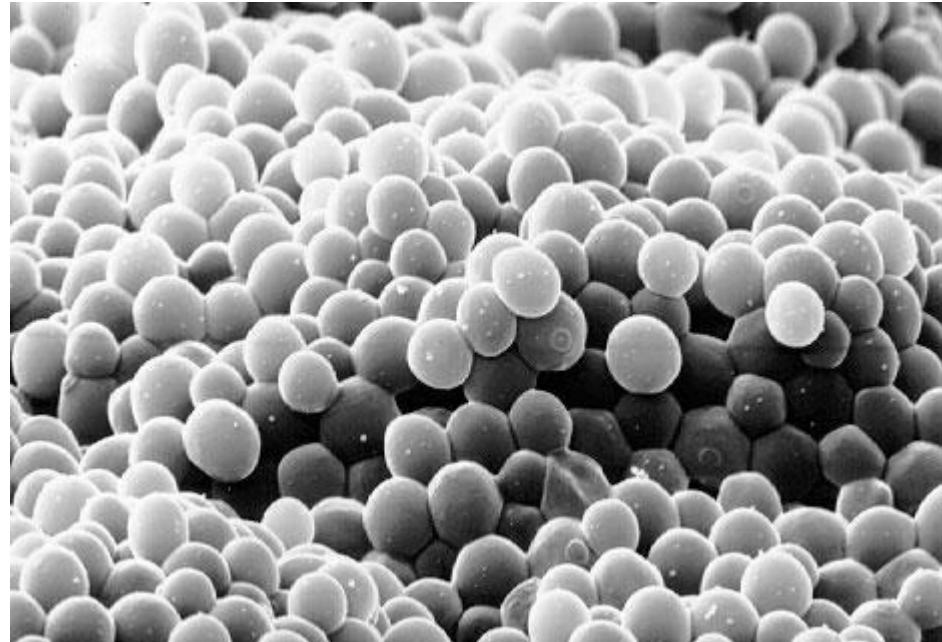
Et on n'est pas seuls...

- E. Coli: 5 Mb



Et on n'est pas seuls...

- E. Coli: 5 Mb
- Levure: 12 Mb



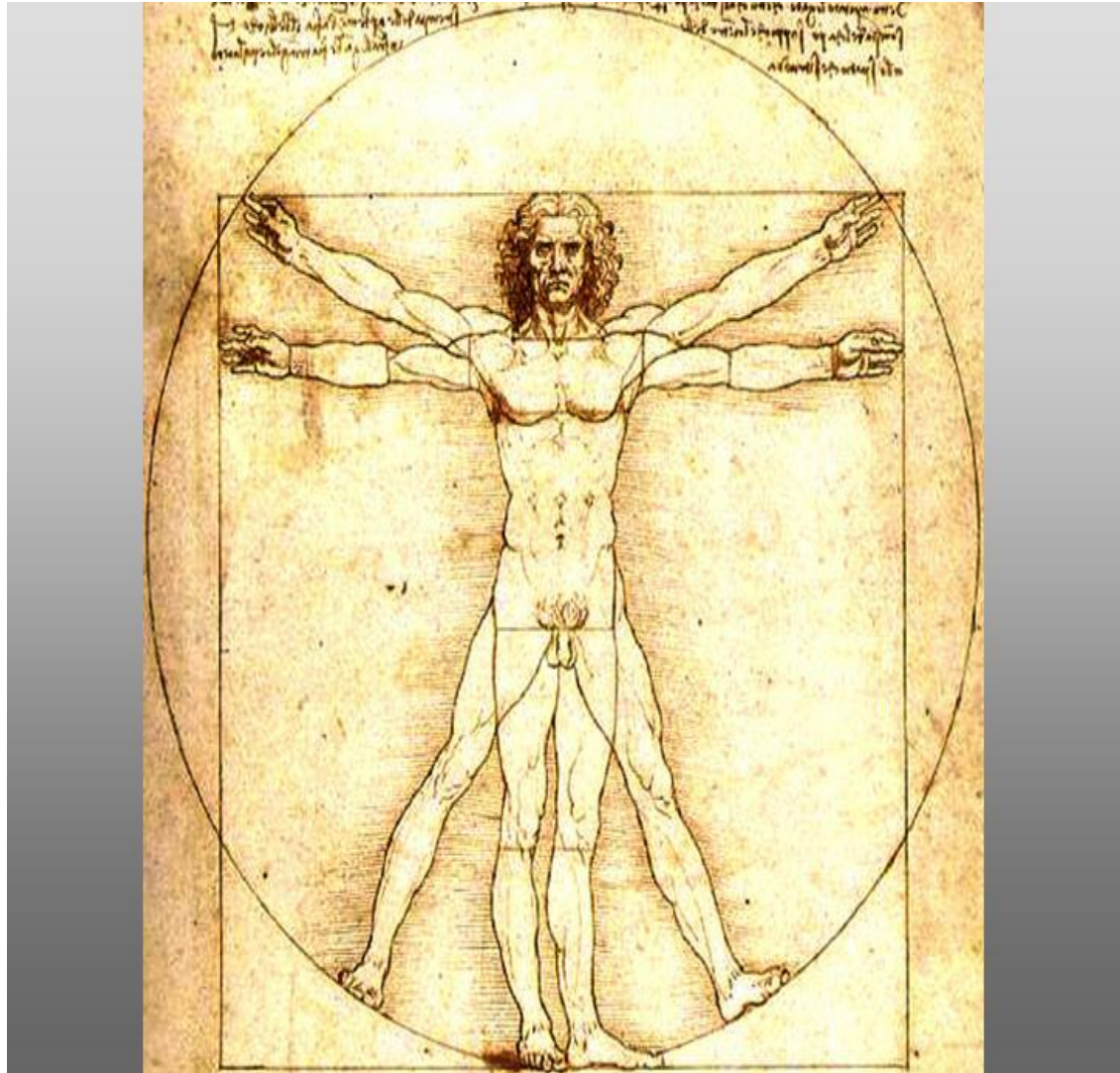
Et on n'est pas seuls...

- E. Coli: 5 Mb
- Levure: 12 Mb
- Mouche: 150 Mb



Et on n'est pas seuls...

- E. Coli: 5 Mb
- Levure: 12 Mb
- Mouche: 150 Mb
- Homme: 3000 Mb



Et on n'est pas seuls...

- E. Coli: 5 Mb
- Levure: 12 Mb
- Mouche: 150 Mb
- Homme: 3000 Mb
- Souris: 3400 Mb



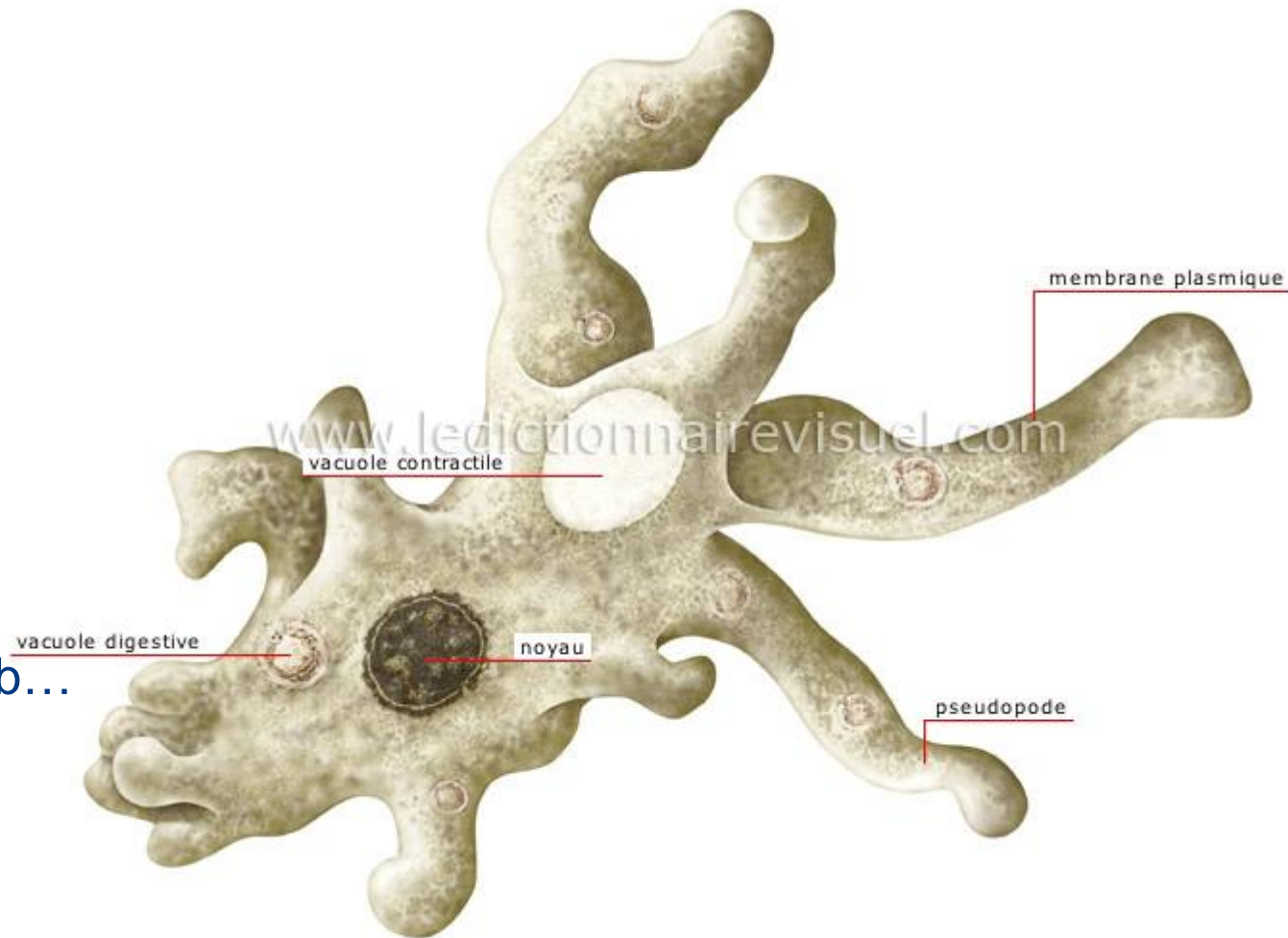
Et on n'est pas seuls...

- E. Coli: 5 Mb
- Levure: 12 Mb
- Mouche: 150 Mb
- Homme: 3000 Mb
- Souris: 3400 Mb
- Maïs: 5000 Mb



Et on n'est pas seuls...

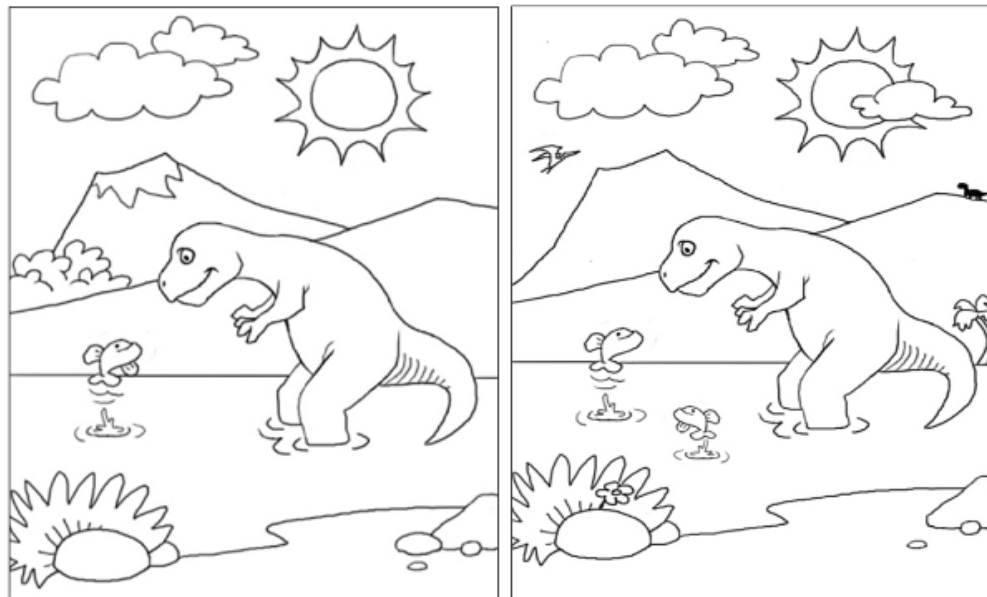
- E. Coli: 5 Mb
- Levure: 12 Mb
- Mouche: 150 Mb
- Homme: 3000 Mb
- Souris: 3400 Mb
- Maïs: 5000 Mb
- Amibe: 675000 Mb...



Que faire de ces données?

→ Les comparer !

- Apprendre en comparant les séquences:
 - Commun = conservé par l'évolution = fonctionnel
 - Différent = "inutile" ou différences en espèces



Que faire de ces données?

→ Les comparer !

- Apprendre en comparant les séquences:
 - Commun = conservé par l'évolution = fonctionnel
 - Différent = "inutile" ou différences en espèces
- Mais aussi :
 - Ne pas recommencer: inférer les connaissances
 - Relations de parenté, mécanismes d'évolution
 - ...

Similarité de séquences

.TGATGCCA = TGATGCCA

.TGATGCCA ≈ TGAAGCCA

.TGATGCCA ≈ TGAGCCA



Similarité de séquences (Level 1)

.ATGCTAAATGGGGAAAAGCTGGGTACCATACG
.ACGCATGAACGATAAATGGGGAAAATAGGAGT

- Deux fragments identiques de taille >5 ?



Similarité de séquences (Level 1)

.ATGCTTAAATGGGGAAAAGCTGGGTACCATACG

.ACGCATGAACGATTAAATGGGGAAAATAGGAGT

TAAATGGGGAAA
 | | | | | | | | | | | | | | | |
 TAAATGGGGAAA



Similarité de séquences (Level 2)

.ATGCTAAATGGGGAAAAGCTGGGTACCATACG

.ACAACAGCTGTACTACGCGAGCAGACTACTC

- Deux fragments similaires (85 % identité minimum) de taille >5 ?



Similarité de séquences (Level 2)

.ATGCTAAATGGGGAAAAGCTGGTACCATACG

.ACAACAGCTGTACTACGCGAGCAGACTACTC

```

AAAAGCTGGTAC
|| . ||||| |||
AACAGCTG - TAC
  
```



Similarité de séquences (Level 3)

.ATGCATCGCTAACGCATCAGATGGCTGACTGACTGCATTTTCTCT
 ACGATCGGCATCGATCTCAGCATCAGCATCATGCAGCATCATCTA
 TACGACCAGCATACTCAGCAGCAGCATAACGCATCATCAGCAGCAG
 CGAGCAGCAGCAGGCAGCAGCAGCGAGCAGCATACTGACTACGAC
 TCAGACTACCAACGTTAGCAGCAGCAGGACCAGGCAGGAGCGGC
 GAGACGATGCACATCAGCATCAGGCAGCATCGAACGTACGGCAC
 GAGGCATATACTACTACTCATTTCGCAGAGCGACGAAGCTA

.ACTGCGCAGTACGATCGATCGTACGTACGATCGCATCGATCGAT
 CGTAGCATCGTACGTACGATCGATCGATCGATCGATCTAGCTACG
 ATCGATCGATCGATCATGGCTGACTGACTGACTCGATCGATCGAT
 CGATCGATCGTACGTACGTACGTACGTACGTACGACTGATCGTCG
 GCGCGCGATCGTACTACGATCGTACGACGTACTCACGATCGGGCAT
 CGATCTACGATCGATCGATCGAGCTACGTTCGTACGATCGTACGGC
 GCGATATACGATCGATACTGACTCGTACGATCGATCGTACGTACG



Similarité de séquences (Level 3)

.ATGCATCGCTAACGCATCAGATGGCTGACTGACTGCATTTTCTCT
ACGATCGCATCGATCTCAGCATCAGCATCATGCAGCATCATCTATA
 CGACCAGCATACTCAGCAGCAGCATAACGCATCATCAGCAGCAGC
 GAGCAGCAGCAGGCAGCAGCAGCGAGCAGCATAACGACTACGACT
 CAGACTACCAACGTTAGCAGCAGCAGGACCAGGCAGGAGCGGCG
 AGACGATGCACATCAGCATCAGGCAGCATTCGAACGTACGGCACG
 AGGCATATACTACTACTCATTTCGCAGAGCGACGAAGCTA

.ACTGCGCAGTACGATTCGATCGTACGTACGATCGCATCGATCGAT
 CGTAGCATCGTACGTACGATCGATCGATCGATCGATCTAGCTACG
 ATCGATCGATCGATCATGGCTGACTGACTGACTCGATCGATCGAT
 CGATCGATCGTACGTACGTACGTACGTACGTACGACTGATCGTCG
 GCGCGCGATCGTACTACGATCGTACGACGTACTCACTATCGGCAT
CGACCTACGATCGATCGATCGAGCTACGTTCGTACGATCGTACGGC
 GCGATATACGATCGATACTGACTCGTACGATCGATCGTACGTACG



Similarité de séquences (Niveau 3)

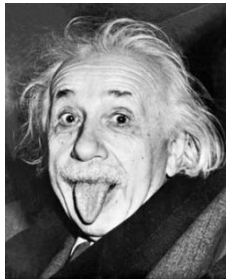
.ATGCATCGCTAACGCATCAGATGGCTGACTGACTGCATTTTCTCT
 ACGATCGGCATCGATCTCAGCATCAGCATCATGCACCATC
 TACGACCAGCATACTC

En pratique:

- Deux genomes complets
- Un génome contre une banque
- Plus de deux génomes

CGATCGATCGTACGTACGTACGTACGTACGTACGACTGATCGTTCG
 GCGCGCGATCGTACTACGATCGTACGACGTACTCACGATCGGCAT
 CGATCTACGATCGATCGATCGAGCTACGTTCGTACGATCGTACGGC
 GCGATATACGATCGATACTGACTCGTACGATCGATCGTACGTACG

Similarité : comparer 2 séquences



- Trouver les similarités entre deux génomes:
- 95 ans sur un ordinateur classique actuel
- Besoins algorithmes
- Besoins machines spécialisées



Similarité : comparer 2 séquences

- Blast: www.ncbi.nlm.nih.gov/BLAST/
- Outil de comparaison de séquences.
- Article scientifique le plus cité, toutes disciplines confondues !

BLAST Human Sequences. - Mozilla Firefox

http://www.ncbi.nlm.nih.gov/genome/seq/BlastGen/BlastGen.cgi?taxid=9606

NCBI Home > Genomic Biology > Human Genome Resources > BLAST

Search
Map Viewer

BLAST Human Sequences.

Enter an accession, gi, or a sequence in FASTA format:

Or, choose a file to upload

Set subsequence: (optional)
From: To:

Database:
genome (all assemblies) 11496 sequences

Program:
megaBLAST: Compare highly related nucleotide sequences

Optional parameters

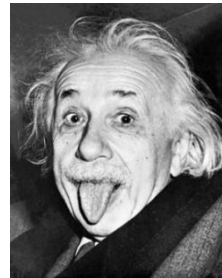
Expect	Filter	Descriptions	Alignments
0.01	default	100	100

Advanced options:

Begin Search Clear Input auto-check for results

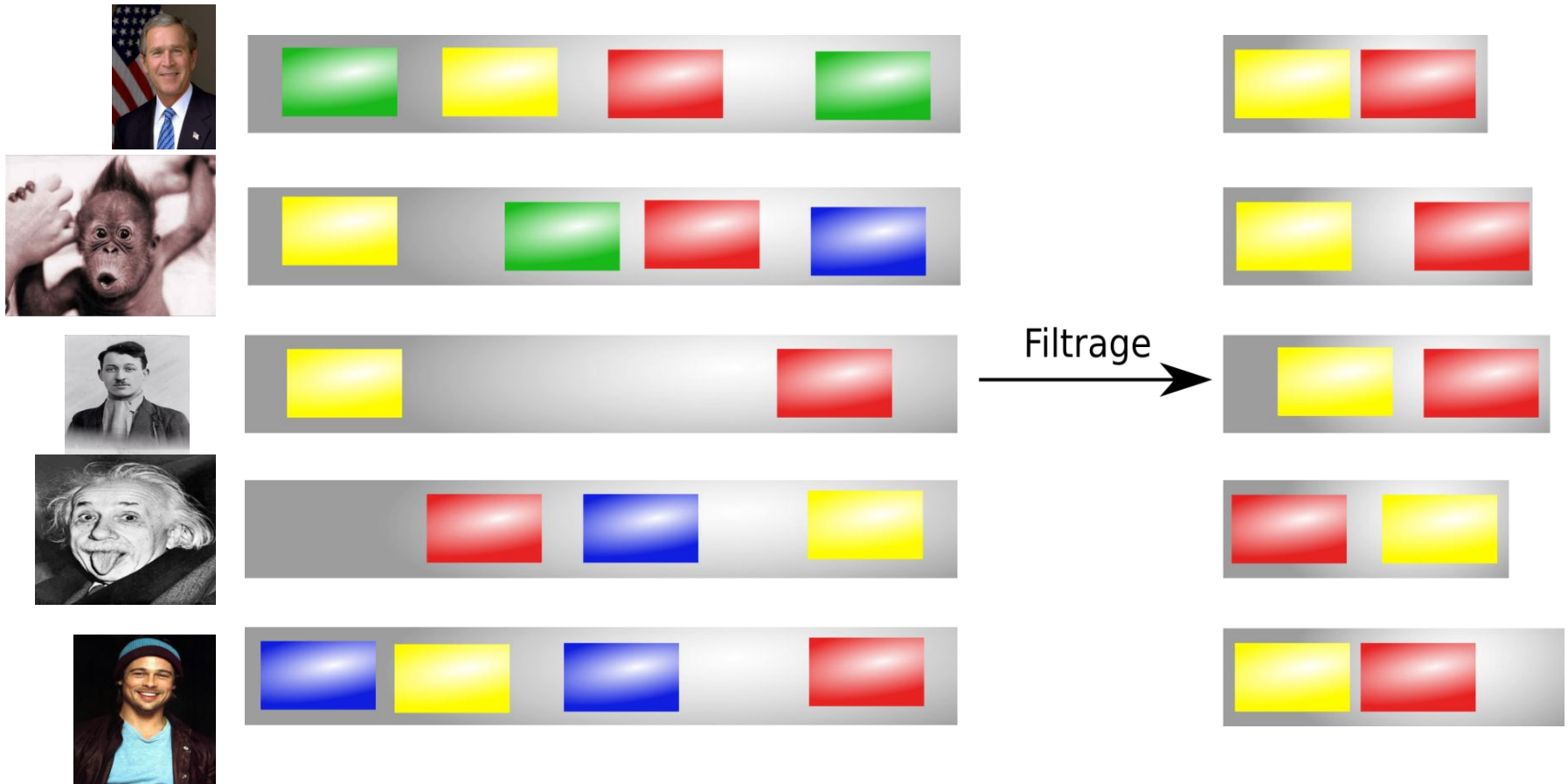
Get the URL with preset values ? Get URL

Similarité : comparer plus de 2 séquences



- Similarités de 5 génomes
- 6.10^{19} x âge univers

Similarité : comparer plus de 2 séquences



Les étapes de la construction du filtre.

- Conception modèle (Biologistes)
- Conception méthode (Condition nécessaire mais non suffisante)
- Conception algorithme (Des concepts au langage algorithmique)
- Codage (Ecriture du programme) → Filtrage →
- Tests (programme seul / résultats biologiques)
- Diffusion
 - Article
 - Site Web



Similarité : comparer plus de 2 séquences

- Résultats
- Exemple : 24 minutes au lieu de 11h30
- Env. 10 utilisations quotidienne

Lossless filter for multiple repeats with bounded edit distance

Pierre Peterlongo^{*1}, Gustavo Akio Tominaga Sacomoto², Alair Pereira do Lago³,
Nadia Pisanti⁴, Marie-France Sagot⁵

¹Équipe-projet Symbiose, IRISA / CNRS, Campus de Beaulieu, Rennes, France

²Curso Experimental de Ciências Moleculares da Universidade de São Paulo, Brazil

³Instituto de Matemática e Estatística da Universidade de São Paulo, Brazil

⁴Dipartimento di Informatica, Università di Pisa, Italy

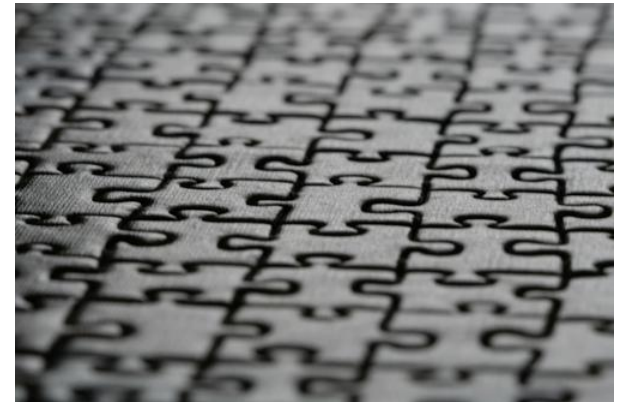
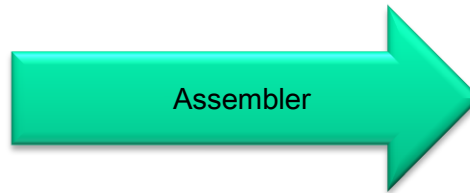
⁵Équipe BAOBAB, Laboratoire de Biométrie et Biologie Evolutive (UMR 5558); CNRS; Univ. Lyon 1, Villeurbanne Cedex, France and Équipe-Projet BAMBOO, INRIA Rhône-Alpes, France and King's College, London, UK

Email: Pierre Peterlongo* - pierre.peterlongo@irisa.fr; Gustavo Akio Tominaga Sacomoto - sacomoto@gmail.com;
Alair Pereira do Lago - alair@ime.usp.br; Nadia Pisanti - pisanti@di.unipi.it; Marie-France Sagot - Marie-France.Sagot@inria.fr;

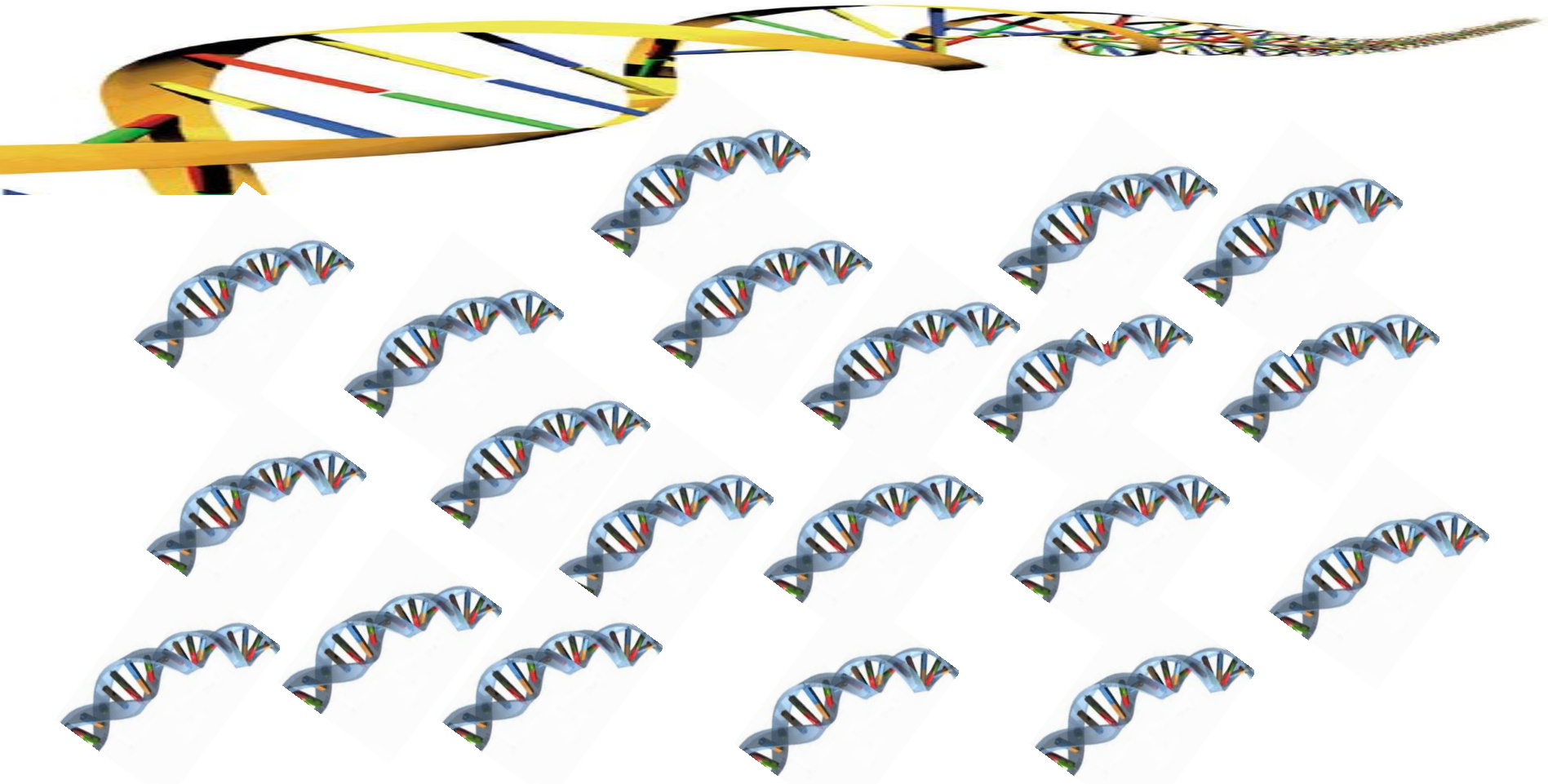
* Corresponding author



Autre exemple: Assembler des données



Les défauts des séquenceurs...





Les défauts des séquenceurs...



Que faire de 2 fragments ?



AGTATAGGAAGATAGACAGCAAGC

GAAGATAGACAGCAAGCAGAGATAG



Que faire de 2 fragments ?



AGTATAGGAAGATAGACAGCAAGC

GAAGATAGACAGCAAGCAGAGATAG



AGTATAGGAAGATAGACAGCAAGCAGAGATAG





En pratique: millions/millards de fragments

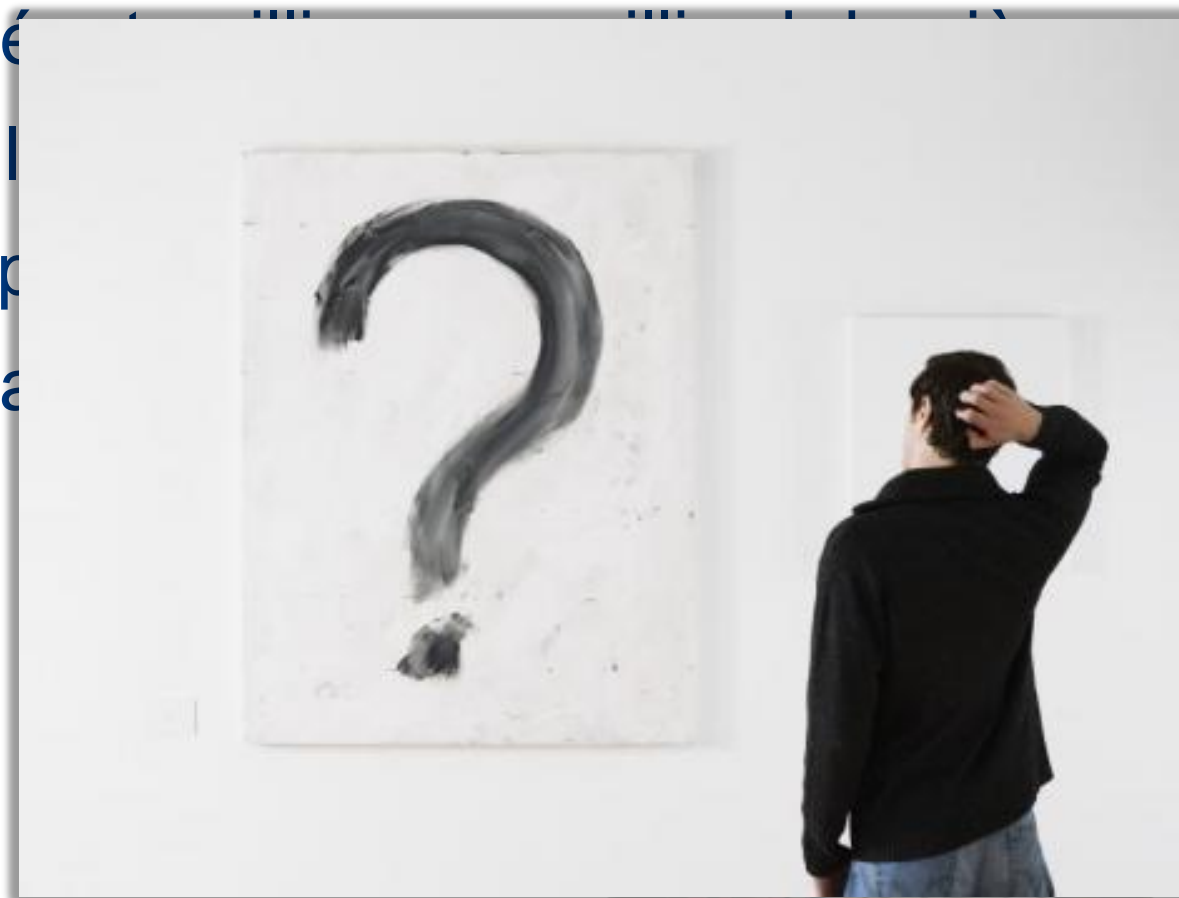
- Puzzle géant: millions ou milliard de pièces
- Manque la boîte avec l'image
- Chaque pièce: pile ou face ?
- Assemblage des pièces imparfait (erreurs)
- ...



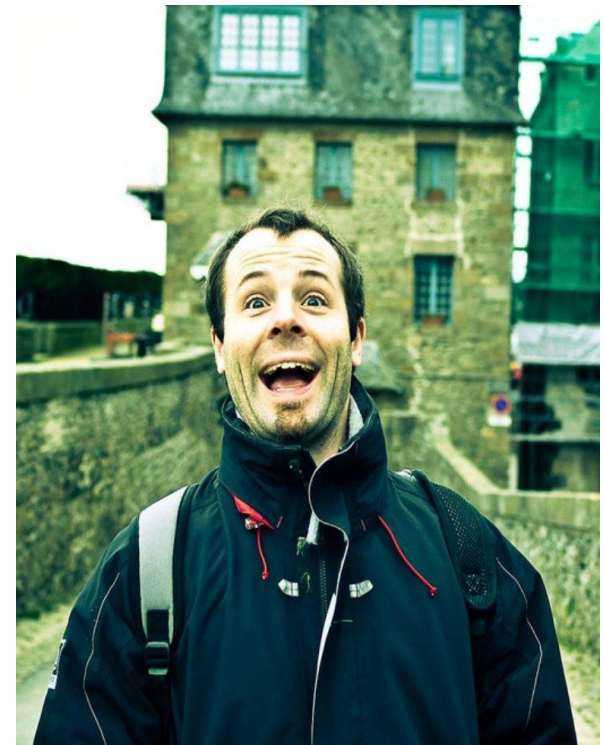


En pratique: millions/millards de fragments



- Puzzle gé
- Manque l
- Chaque p
- Assembla
- ...



C'est quoi un Bioinformaticien ?



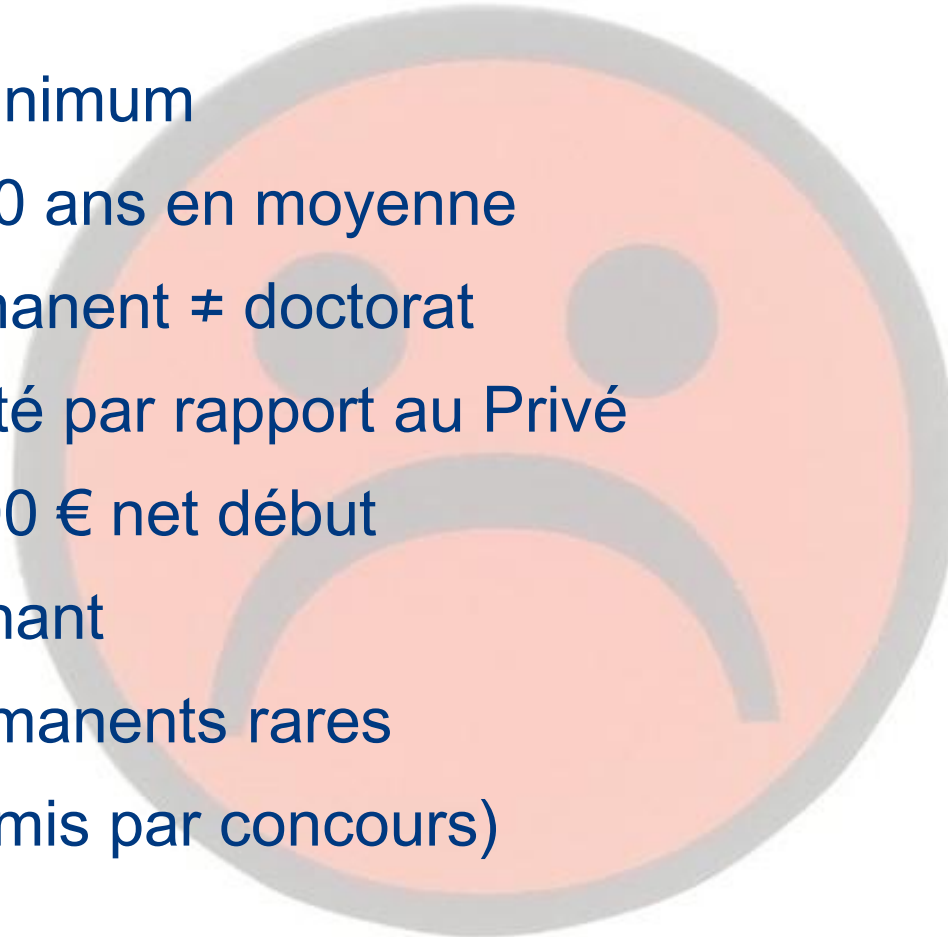
(mon) Parcours

- BAC
- Licence 3 ans (4 en fait)
- Master 2 ans
- Doctorat (thèse) 3 ans 
- Post Doc 2 ans 
- « Chargé de recherche »
 - Concours
 - Fonctionnaire



Les « moins »

- BAC + 8 minimum
- Fixé vers 30 ans en moyenne
- Poste permanent ≠ doctorat
- Salaire limité par rapport au Privé
- 1600 à 2000 € net début
- Travail prenant
- Postes permanents rares
 - (5 % admis par concours)



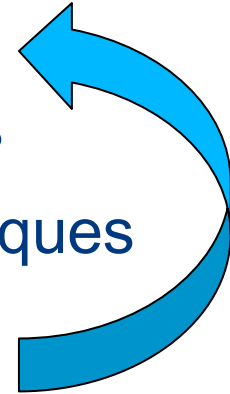
Les « plus »

- Plus de « plus » que de « moins »
- Travail passionnant
- Voyages
- Liberté (thématiques)
- Sécurité emploi
- Peu de pression par l'argent
- De moins en moins vrai
- Rencontres (internationales)



Le quotidien

- Derrière un ordinateur (pas de manipulations)
 - Question biologique
 - Inventer de nouveaux algorithmes / logiciels
 - Tester ces logiciels sur des données biologiques
 - Réponse... et d'autres questions
- Discussions avec :
 - des biologistes
 - des informaticiens
 - etc.



Pas le même langage !



Le quotidien, en plus de la « recherche »...

- Recherche dans un contexte international, et collaboratif
→ **échanges**
 - Publications et exposés dans des conférences
 - Lire / assister / évaluer
 - Écrire / présenter
 - Mais aussi : relire (évaluations), organiser, etc.
 - Bien sûr, en anglais !
- Grand public
- Chercher de l'argent
- Encadrer des étudiants / enseigner à l'université



Merci !

Contact: pierre.peterlongo@inria.fr