

CONNAITRE LE GENOME

La génomique

Structure physique

Structure génétique

Les espèces modèles

Plantes

Champignons

Prokaryotes

Animaux



Un ancêtre commun ?
LUCA

Last Universal Common Ancestor

Les objectifs

- Localiser et séquencer les gènes



- Étudier la fonction des gènes

Les applications

Construction de génotypes élites

Par :

- la fourniture d'une grande quantité de marqueurs moléculaires
- un meilleur contrôle de la régulation des gènes
- l'identification de gènes candidats pour l'analyse des QTL de caractères majeurs
- l'identification d'allèles favorables

Analyse fonctionnelle

La Biologie moléculaire

et

le développement de la
bioinformatique

La Biologie moléculaire (mi 80- début 1990)

Le clonage et le séquençage d'ADN deviennent progressivement pratiques courantes.

L'utilisation de l'outil informatique est devenu indispensable pour réaliser les expériences de Biologie Moléculaire et pour en analyser les résultats



→ Au départ, c'est la BM qui a besoin de la Bioinfo

Internet, BLAST et séquencage

Internet : L'accès aux bases de données et aux outils d'analyse devient transparent et facile pour tous

BLAST : rechercher dans les bases de données l'existence de séquences similaires à une séquence donnée c'est la fonction utilisée dans plus de 90% des analyses. BLAST (1990) permet à tout un chacun de le faire, en classant les résultats significatifs

Le séquençage génomique : il donne accès d'un seul coup à toute l'information génétique d'un organisme ... mais comment la gérer, la déchiffrer et l'utiliser ?

*explosion de nouveaux besoins pour la bioinformatique
explosion de données d'un type nouveau à exploiter*



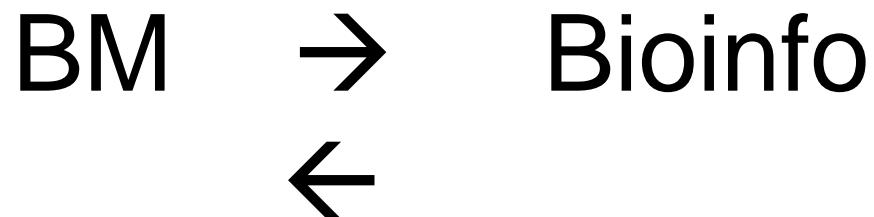
Besoins

Gérer les programmes de séquençage

Annoter les génomes : où sont les gènes? Que font-ils ?

Etablir des ontologies et les peupler : de quoi parle t'on au juste ?

Permettre la comparaison de génomes entre eux



- Si au départ, c'est la BM qui avait besoin de la Bioinfo

- aujourd'hui, la bioinfo a aussi besoin des données générées par la BM

CONNAITRE LE GENOME

La génomique

Structure physique

Structure génétique

Les espèces modèles

Quels modèles choisir ?

Plantes

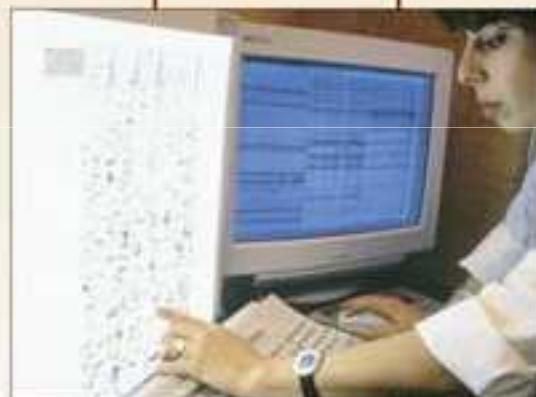
Champignons

Prokaryotes

Animaux

Les objectifs

- Localiser et séquencer les gènes



- Étudier la fonction des gènes

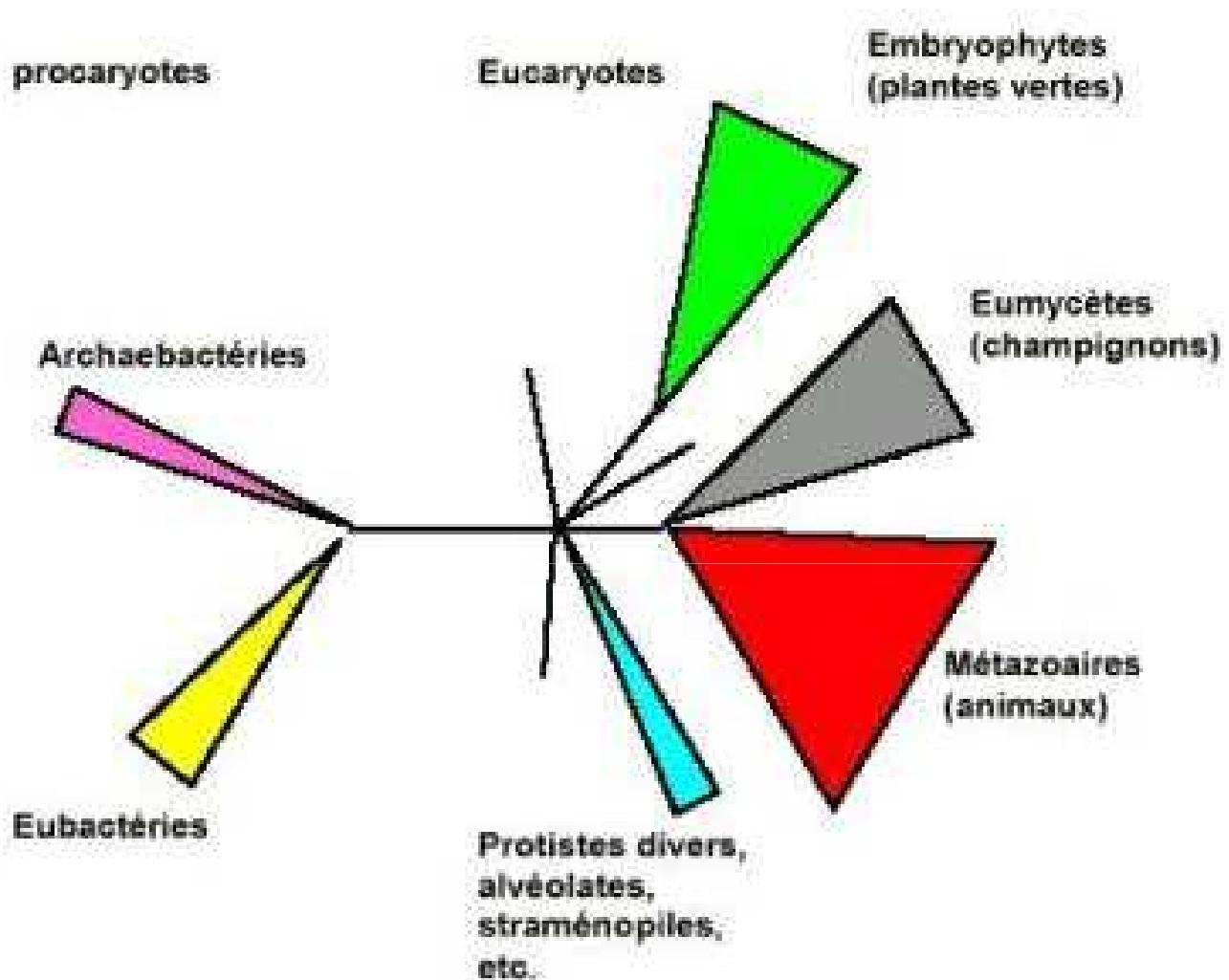
Analyse fonctionnelle

Les applications

Construction de génotypes élites

Par :

- la fourniture d'une grande quantité de marqueurs moléculaires
- un meilleur contrôle de la régulation des gènes
- l'identification de gènes candidats pour l'analyse des QTL de caractères majeurs
- l'identification d'allèles favorables



Les triangles représentent l'effectif approximatif en espèces de chaque groupe

Le choix des espèces modèles se fera sur la taille des génomes et sur leur facilité de « culture » en laboratoire

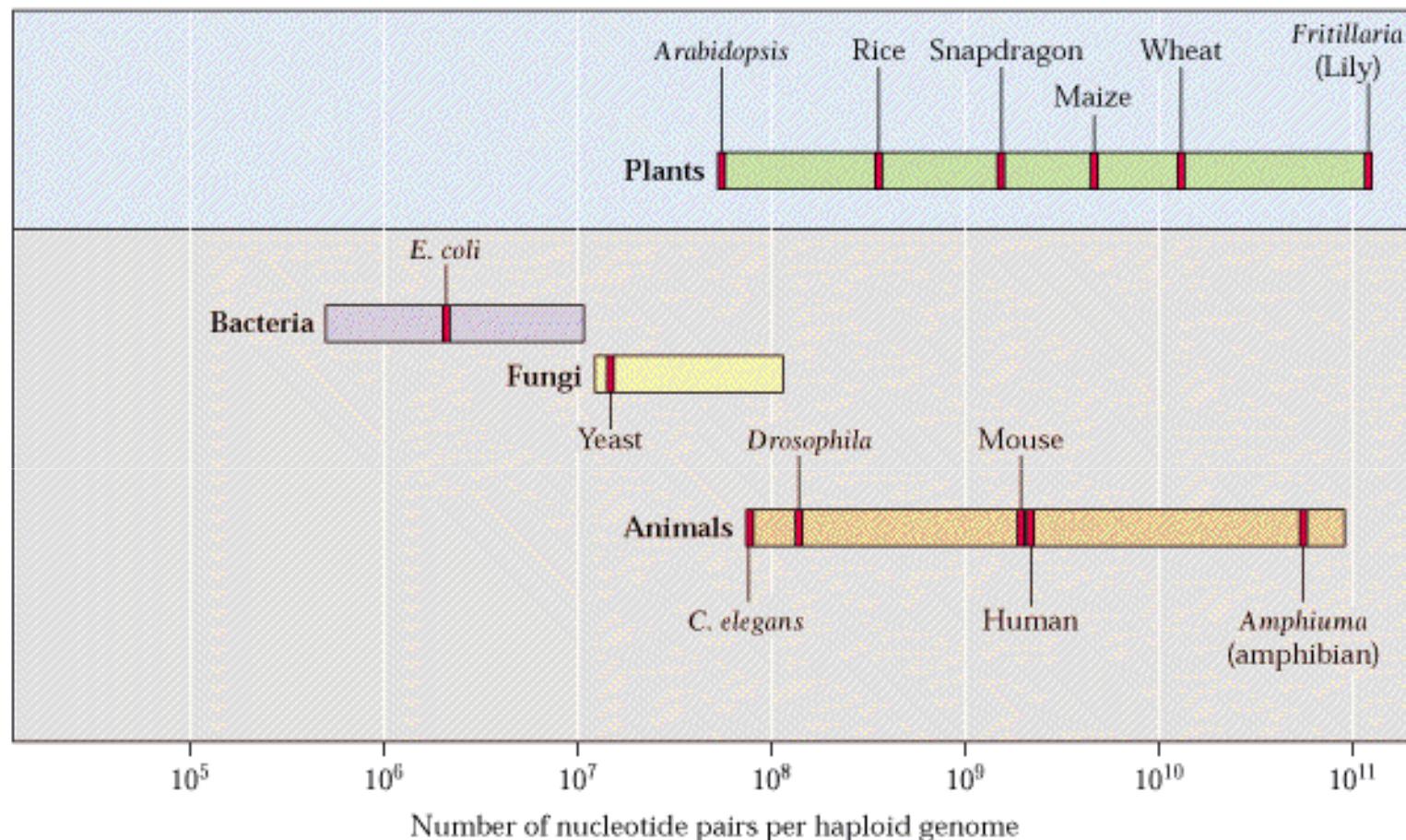
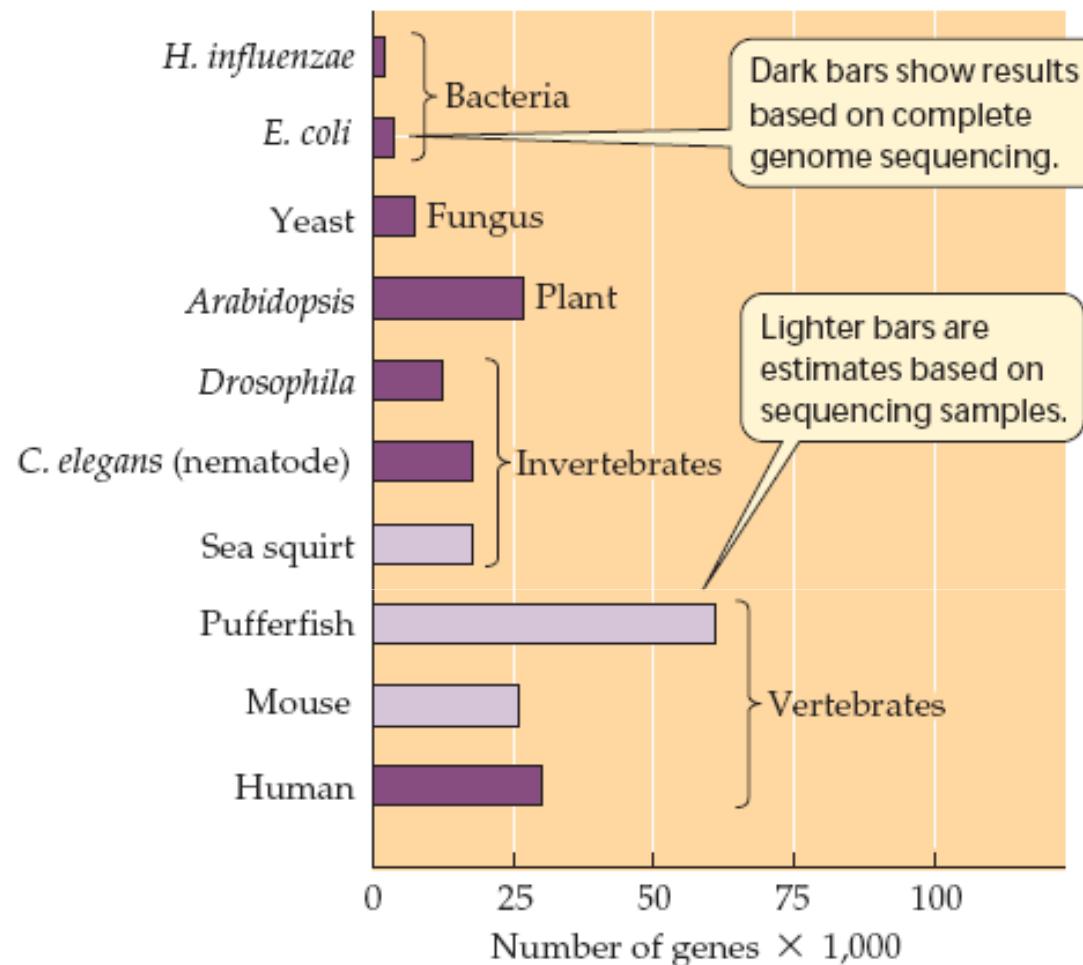


Figure 7.19
C values (haploid genome size in basepairs) from various organisms.

Combien de gènes par organismes ?



26.7 Complex Organisms Have More Genes than Simpler Organisms Genome sizes have been measured or estimated in a variety of organisms ranging from single-celled prokaryotes to vertebrates.

gènes

Percent of genome that is functional genes

100

80

60

40

20

0

0.001 0.01 0.1 1 10 100 1000

Genome size ($\times 10^9$ base pairs)*E. coli*

Yeast

*Drosophila**Arabidopsis**C. elegans*

Human

Newt

Lungfish

Lily

Taille du
génome

Quel part pour l'ADN codant et non-codant ?

26.8 A Large Proportion of DNA Is Noncoding

Most of the DNA of bacteria and yeasts encodes RNAs or proteins, but most of the DNA of more complex organisms is noncoding. Most noncoding DNA is probably nonfunctional.

Les programmes de séquençage du génome

- Les genomes center ?
- Les programmes ESTs
- Le séquençage des génomes
- L'annotation des génomes
- Les centres de ressources

Les genomes centers

Human Genome Sequencing Center / Baylor College of Medicine, Houston (Texas) ; USA

Human Genome Center / Beijing Genomics Institute, Académie chinoise des sciences, Beijing ; Chine

Lita Annenberg Hazen Genome Center / Cold Spring Harbor Laboratory, Cold Spring Harbor (N.Y.), USA

Gesellschaft fur Biotechnologische Forschung mbH, Braunschweig ; Allemagne

Genoscope, Evry ; France

GTC Sequencing Center / Genome Therapeutics Corp., Waltham (Mass.) ; USA

Department of Genome Analysis / Institute of Molecular Biotechnology, Jena ; Allemagne

Joint Genome Institute / U.S. Department of Energy, Walnut Creek (Calif.) ; USA

Département de biologie moléculaire / Ecole de médecine de l'université Keio, Tokyo ; Japon

Max Planck Institute for Molecular Genetics, Berlin ; Allemagne

Multimegaparsec Sequencing Center / The Institute for Systems Biology, Seattle (Wash.) ; USA

RIKEN Genomic Sciences Center, Yokohama ; Japon

The Wellcome Trust Sanger Institute (Sanger Center), Hinxton ; Royaume-Uni

Stanford Human Genome Center, Stanford (Calif.) ; USA

University of Oklahoma / Advanced Center for Genome Technology, Norman (Okla.), USA

Whitehead Institute / MIT Center for Genome Research, Cambridge (Mass.) ; USA

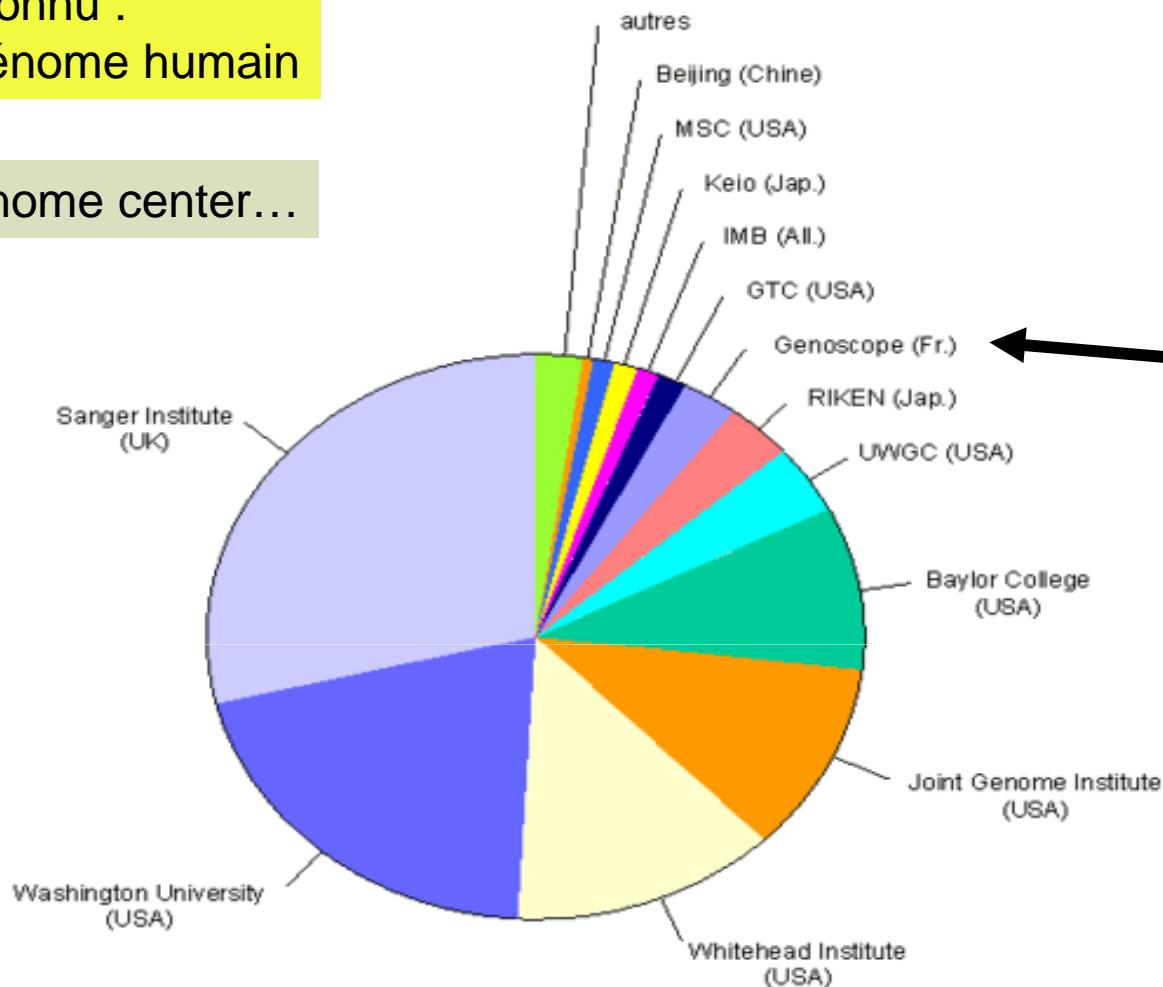
Washington University / Genome Sequencing Center, St Louis (Mo.) ; USA

USA

Les capacités de séquençage sont exponentielles : et chaque centre ou fédération de laboratoires a ses propres séquenceurs → Vers du séquençage de paillasse

L'exemple le mieux connu : Le séquençage du génome humain

Par genome center...



Sur le plan international, les contributions des 6 pays impliqués dans le projet sont les suivantes :

Par pays ...

Etats-Unis	60,8 %
Royaume-Uni	28,9 %
Japon	4,9 %
France	2,8 %
Allemagne	1,5 %
Chine	0,7 %

Le séquençage du génome humain

Le coût total du projet Génome humain est d'environ 2,7 milliards de dollars alors qu'il avait été estimé à 3 milliards de dollars au début du projet, en 1990.

Cette économie résulte de **progrès techniques considérables et de l'accélération du projet, terminé avec deux ans d'avance sur les prévisions.**

Une grande part de la somme a été dépensée pour la finition de l'ébauche génomique obtenue en 2000.

Le séquençage du chromosome 14 a coûté environ 10 millions d'euros,
+
plusieurs millions d'euros pour l'analyse et l'annotation.

ESTs

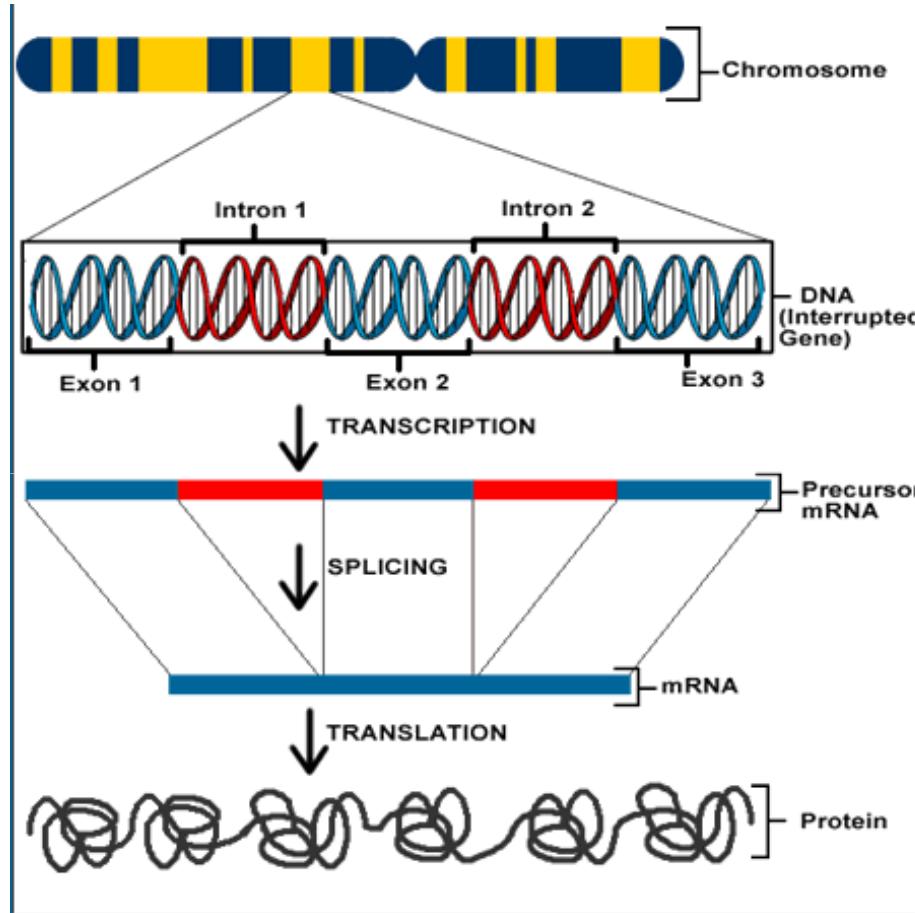


figure 1. An overview of the process of protein synthesis.

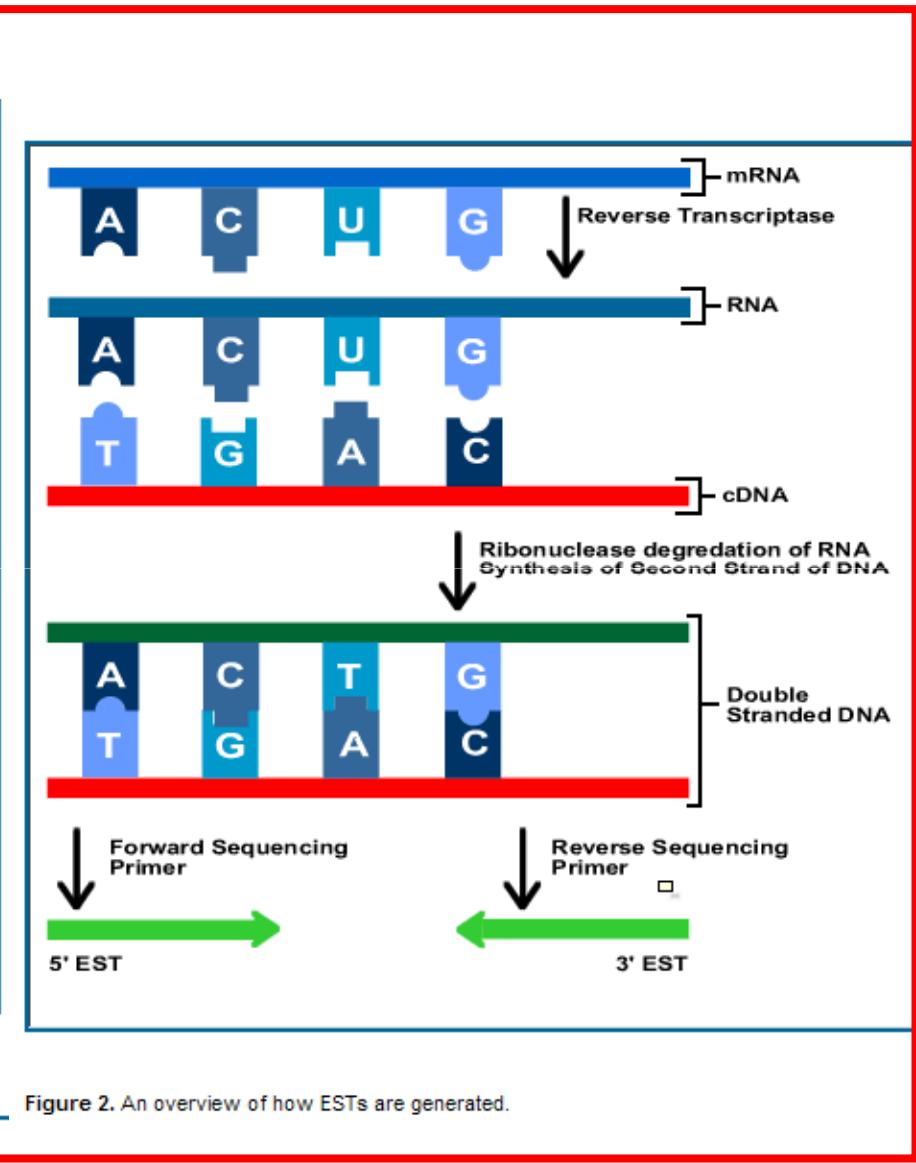


Figure 2. An overview of how ESTs are generated.

Banques ADNc (multiples conditions cultures ...)

→ Clone pris au hasard, séquencage et dépôt de la séquence dans une base de donnée et du clone dans un centre de ressource génétique

Séquençage (plutot en 5') de 300-800n en 1 seule réaction

Summary by Organism - Jan 2, 1997

Number of public entries: 775,836

Homo sapiens (human)	539,706
Mus musculus + domesticus (mouse)	120,098
Caenorhabditis elegans (nematode)	30,196
→ Arabidopsis thaliana (thale cress)	29,166 —
→ Oryza sativa (rice)	12,672 —
Toxoplasma gondii	8,165
Rattus sp. (rat)	6,615
Brugia malayi (parasitic nematode)	5,184
Saccharomyces cerevisiae (baker's yeast)	2,944
Caenorhabditis briggsae	2,424
Trypanosoma brucei rhodesiense	2,175
Plasmodium falciparum (malaria parasite)	2,011
Schistosoma mansoni (blood fluke)	1,958
→ Zea mays (maize)	1,757 —
Drosophila melanogaster (fruit fly)	1,690
Brassica napus (oilseed rape)	1,427
Sus scrofa (pig)	1,055
→ Brassica campestris (field mustard)	965 —
→ Ricinus communis (castor bean)	750 —
Pyrococcus furiosus (hyperthermophilic archaeon)	577
Leishmania major	527
→ Citrus unshiu	426 —
Onchocerca volvulus	306
Emericella nidulans	272
→ Pinus taeda (loblolly pine)	256 —
Capra hircus (goat)	245
Strongylocentrotus purpuratus (purple urchin)	242
Oryctolagus cuniculus (rabbit)	206
Gallus gallus (chicken)	191
Danio rerio (zebrafish)	128
Schistosoma japonicum (blood fluke)	109
Saccharum sp.	105
Wuchereria bancrofti	99
→ Glycine max (soybean)	97 —
→ Sorghum bicolor (sorghum)	91 —
→ Citrus sinensis	87 —
Aspergillus niger	85
Pristionchus pacificus	78
Cryptosporidium parvum	65
Bos taurus (cattle)	60
Strongyloides stercoralis	57
Trypanosoma brucei brucei	56
→ Solanum tuberosum (potato)	49 —
Filobasidiella neoformans (Cryptococcus neoformans)	44
Macropus eugenii (tammar wallaby)	42

Summary by Organism - January 21, 2000

Number of public entries: 3,483,128

Homo sapiens (human)	1,715,552
Mus musculus + domesticus (mouse)	904,535
Rattus sp. (rat)	132,680
Caenorhabditis elegans (nematode)	101,232
Drosophila melanogaster (fruit fly)	86,121
Danio rerio (zebrafish)	54,600
Lycopersicon esculentum (tomato)	50,303
Oryza sativa (rice)	47,474
Zea mays (maize)	47,231
Arabidopsis thaliana (thale cress)	45,757
Glycine max (soybean)	40,046
Brugia malayi (parasitic nematode)	20,773
Dictyostelium discoideum	19,182
Emericella nidulans	12,993
Schistosoma mansoni (blood fluke)	12,509
Chlamydomonas reinhardtii	11,690
Toxoplasma gondii	10,741
Neurospora crassa	10,229
— Pinus taeda (loblolly pine)	9,870 —
— Gossypium hirsutum (upland cotton)	9,335 —
Trypanosoma cruzi	8,796
Onchocerca volvulus	8,674
Schizosaccharomyces pombe (fission yeast)	8,118
Bombyx mori (domestic silkworm)	6,508
— Medicago truncatula (barrel medic)	6,059 —
Xenopus laevis (African clawed frog)	5,853
Sorghum bicolor (sorghum)	5,428
Sus scrofa (pig)	5,194
Pristionchus pacificus	4,989
— Populus tremula x Populus tremuloides	4,809
Trypanosoma brucei rhodesiense	4,742
— Mesembryanthemum crystallinum (common ice plant)	4,429
Gallus gallus (chicken)	4,194
Halocynthia roretzi	4,163
Eimeria tenella	3,721
Saccharomyces cerevisiae (baker's yeast)	3,041
Plasmodium falciparum (malaria parasite)	2,871
Caenorhabditis briggsae	2,424
Leishmania major	2,191
Oryctolagus cuniculus (rabbit)	2,094
Physcomitrella patens	1,909
Paralichthys olivaceus	1,787
Brassica napus (oilseed rape)	1,702
Metarhizium anisopliae	1,693
Ceratodon purpureus	1,663
Bos taurus (cattle)	1,610
Schistosoma japonicum (blood fluke)	1,395
Aedes aegypti (yellow fever mosquito)	1,327

X4.5

Summary by Organism - October 13, 2000

Number of public entries: 6,051,828

<i>Homo sapiens</i> (human)	2,519,935
<i>Mus musculus + domesticus</i> (mouse)	1,703,709
<i>Rattus</i> sp. (rat)	188,794
<i>Bos taurus</i> (cattle)	134,604
<i>Glycine max</i> (soybean)	122,843
<i>Arabidopsis thaliana</i> (thale cress)	112,467 *
<i>Caenorhabditis elegans</i> (nematode)	101,252
<i>Drosophila melanogaster</i> (fruit fly)	95,211
<i>Lycopersicon esculentum</i> (tomato)	87,680
<i>Medicago truncatula</i> (barrel medic)	81,694 *
<i>Zea mays</i> (maize)	73,965
<i>Danio rerio</i> (zebrafish)	73,703
<i>Oryza sativa</i> (rice)	62,798
<i>Chlamydomonas reinhardtii</i>	55,860
<i>Sorghum bicolor</i> (sorghum)	45,265
<i>Triticum aestivum</i> (wheat)	44,132
<i>Xenopus laevis</i> (African clawed frog)	43,935
<i>Sus scrofa</i> (pig)	33,465
<i>Lotus japonicus</i>	26,844
<i>Neurospora crassa</i>	24,635
<i>Hordeum vulgare</i> (barley)	23,519
<i>Brugia malayi</i> (parasitic nematode)	22,392
<i>Pinus taeda</i> (loblolly pine)	22,030
<i>Dictyostelium discoideum</i>	19,183
<i>Solanum tuberosum</i> (potato)	14,043
<i>Bombyx mori</i> (domestic silkworm)	14,807
<i>Gossypium arboreum</i>	13,662
<i>Emericella nidulans</i>	12,993
<i>Schistosoma mansoni</i> (blood fluke)	12,959
<i>Gallus gallus</i> (chicken)	12,840
<i>Onchocerca volvulus</i>	12,552
<i>Toxoplasma gondii</i>	12,177
<i>Mesembryanthemum crystallinum</i> (common ice plant)	11,115
<i>Ciona intestinalis</i>	10,347
<i>Porphyra yezoensis</i>	10,185
<i>Trypanosoma cruzi</i>	9,919
<i>Gossypium hirsutum</i> (upland cotton)	9,438
<i>Schizosaccharomyces pombe</i> (fission yeast)	8,118
<i>Physcomitrella patens</i>	7,496 *
<i>Strongyloides stercoralis</i>	7,343
<i>Meloidogyne incognita</i> (southern root-knot nematode)	6,626
<i>Secale cereale</i>	6,574
<i>Anopheles gambiae</i> (African malaria mosquito)	6,023
<i>Eimeria tenella</i>	5,499
<i>Pristionchus pacificus</i>	4,989
<i>Trypanosoma brucei rhodesiense</i>	4,821
<i>Populus tremula x Populus tremuloides</i>	4,809
<i>Halocynthia roretzii</i>	4,163

X 1.7

Summary by Organism - September 26, 2003

Number of public entries: 18,666,654

Homo sapiens (human)	5,425,176
Mus musculus + domesticus (mouse)	3,874,479
Rattus sp. (rat)	537,607
Triticum aestivum (wheat)	499,979
Ciona intestinalis	492,488
Gallus gallus (chicken)	451,109
Zea mays (maize)	362,549
Danio rerio (zebrafish)	362,279
Hordeum vulgare + subsp. vulgare (barley)	348,225
Glycine max (soybean)	341,573
Bos taurus (cattle)	319,805
Xenopus laevis (African clawed frog)	311,905
Drosophila melanogaster (fruit fly)	261,404
Oryza sativa (rice)	260,879
Saccharum officinarum	246,301
Caenorhabditis elegans (nematode)	215,200
Silurana tropicalis	203,883
Arabidopsis thaliana (thale cress)	188,782
Medicago truncatula (barrel medic)	187,763
Dictyostelium discoideum	155,032
Chlamydomonas reinhardtii	154,600
Sorghum bicolor (sorghum)	151,874
Lycopersicon esculentum (tomato)	150,228
Sus scrofa (pig)	150,133
Schistosoma mansoni (blood fluke)	139,064
Anopheles gambiae (African malaria mosquito)	130,731
Vitis vinifera	128,211
Oryzias latipes (Japanese medaka)	103,098
Oncorhynchus mykiss (rainbow trout)	102,218
Pinus taeda (loblolly pine)	100,541
Solanum tuberosum (potato)	94,473
Toxoplasma gondii	72,859
Lactuca sativa	68,188
Helianthus annuus	59,837
Salmo salar	58,330
Populus tremula x Populus tremuloides	56,013
Strongylocentrotus purpuratus (purple urchin)	51,744
Physcomitrella patens subsp. patens	49,583
Schistosoma japonicum (blood fluke)	45,902
Ascaris suum (pig roundworm)	39,242
Gossypium arboreum	38,915
Brassica napus (oilseed rape)	37,108
Lotus corniculatus var. japonicus	36,262
Hydra magnipapillata	35,154
Magnaporthe grisea	31,397
Bombyx mori (domestic silkworm)	28,978
Eimeria tenella	28,550
Neurospora crassa	28,089
Canis familiaris (dog)	27,010

X 3



dbEST: database of "Expressed Sequence Tags"

dbEST release 091908

Summary by Organism - September 19, 2008

X 3

Number of public entries: 55,785,287

Homo sapiens (human)	8,138,094
Mus musculus + domesticus (mouse)	4,850,258
Arabidopsis thaliana (thale cress)	1,526,133
Bos taurus (cattle)	1,517,053
Sus scrofa (pig)	1,476,546
Zea mays (maize)	1,464,859
Danio rerio (zebrafish)	1,379,829
Xenopus (Silurana) tropicalis (western clawed frog)	1,271,375
Oryza sativa (rice)	1,220,876
Ciona intestinalis	1,204,893
Triticum aestivum (wheat)	1,051,300
Rattus norvegicus + sp. (rat)	895,094
Glycine max (soybean)	839,041
Xenopus laevis (African clawed frog)	677,784
Oryzias latipes (Japanese medaka)	616,739
Gallus gallus (chicken)	599,610
Brassica napus (oilseed rape)	596,249
Drosophila melanogaster (fruit fly)	573,749
Hordeum vulgare + subsp. vulgare (barley)	478,734
Panicum virgatum (switchgrass)	436,535
Salmo salar (Atlantic salmon)	433,337
Canis lupus familiaris (dog)	365,909
Vitis vinifera (wine grape)	352,984
Caenorhabditis elegans (nematode)	346,109
Branchiostoma floridae (Florida lancelet)	334,502
Pinus taeda (loblolly pine)	328,628
Physcomitrella patens subsp. patens	305,606
Aedes aegypti (yellow fever mosquito)	301,342
Ictalurus punctatus (channel catfish)	300,678
Gasterosteus aculeatus (three spined stickleback)	276,992
...	...

dbEST release 091109

Summary by Organism - September 11, 2009

Number of public entries: 62,950,194

+ ~ 8 000 000 /année 2008

Homo sapiens (human)	8,296,272
Mus musculus + domesticus (mouse)	4,852,144
Zea mays (maize)	2,018,634
Sus scrofa (pig)	1,536,375
Arabidopsis thaliana (thale cress)	1,527,298
Bos taurus (cattle)	1,517,145
Danio rerio (zebrafish)	1,481,930
Glycine max (soybean)	1,422,497
Xenopus (Silurana) tropicalis (western clawed frog)	1,271,375
Oryza sativa (rice)	1,248,955
Ciona intestinalis	1,205,674
Triticum aestivum (wheat)	1,067,115
Rattus norvegicus + sp. (rat)	1,009,817
Drosophila melanogaster (fruit fly)	821,005
Xenopus laevis (African clawed frog)	677,806
Oryzias latipes (Japanese medaka)	665,382
Brassica napus (oilseed rape)	631,983
Gallus gallus (chicken)	600,075
Hordeum vulgare + subsp. vulgare (barley)	501,614
Salmo salar (Atlantic salmon)	494,152
Panicum virgatum (switchgrass)	436,535
Phaseolus coccineus	391,138
Canis lupus familiaris (dog)	365,909
Physcomitrella patens subsp. patens	362,131
Vitis vinifera (wine grape)	357,849
Caenorhabditis elegans (nematode)	354,744
Ictalurus punctatus (channel catfish)	354,434
Branchiostoma floridae (Florida lancelet)	334,502
Pinus taeda (loblolly pine)	328,628
Malus x domestica (apple tree)	324,308
Ovis aries (sheep)	323,866



dbEST: database of "Expressed Sequence Tags"

dbEST release 082710

Summary by Organism - August 27, 2010

Number of public entries: 66,682,428

Homo sapiens (human)	8,304,690
Mus musculus + domesticus (mouse)	4,852,147
Zea mays (maize)	2,019,105
Sus scrofa (pig)	1,621,000
Bos taurus (cattle)	1,559,485
Arabidopsis thaliana (thale cress)	1,529,555
Danio rerio (zebrafish)	1,481,936
Glycine max (soybean)	1,459,936
Xenopus (Silurana) tropicalis (western clawed frog)	1,271,375
Oryza sativa (rice)	1,249,124
Ciona intestinalis	1,205,674
Rattus norvegicus + sp. (rat)	1,162,136
Triticum aestivum (wheat)	1,071,199
Drosophila melanogaster (fruit fly)	821,005
Xenopus laevis (African clawed frog)	677,806
Oryzias latipes (Japanese medaka)	665,382
Brassica napus (oilseed rape)	643,884
Gallus gallus (chicken)	600,418
Panicum virgatum (switchgrass)	546,245
Hordeum vulgare + subsp. vulgare (barley)	501,620
Salmo salar (Atlantic salmon)	496,300
Caenorhabditis elegans (nematode)	393,714
Phaseolus coccineus	391,150
Canis lupus familiaris (dog)	382,618
Vitis vinifera (wine grape)	362,193
Physcomitrella patens subsp. patens	362,131
Ictalurus punctatus (channel catfish)	354,466
Ovis aries (sheep)	335,950
Branchiostoma floridae (Florida lancelet)	334,502
Pinus taeda (loblolly pine)	328,628
-----	-----

dbEST release 120701
Summary by Organism
- 01 July 2012

Number of public entries:
73,360,923

X 100 en 15 ans

An inventory of 1152 expressed sequence tags obtained by partial sequencing of cDNAs from *Arabidopsis thaliana*[†]

Herman Höfte*, Thierry Desprez, Joëlle Amselem, Hélène Chiapello and Michel Caboche, Annick Moisan, Marie-Françoise Jourjon and Jean-Louis Charpenteau, Pierre Berthomieu, Danièle Guerrier and Jérôme Giraudat, Françoise Quigley, Frank Thomas, De-Yao Yu and Régis Mache, Monique Raynal, Richard Cooke, Françoise Grellet and Michel Delseny, Yves Parmentier, Guy de Marcillac, Claude Gigot, Jacqueline Fleck and Gabriel Philipp, Michèle Axelos, Claude Bardet, Dominique Tremousaygue and Bernard Lescure

As part of the goal to generate a detailed transcript map for *Arabidopsis thaliana*, 1152 single run sequences (expressed sequence tags or ESTs) have been determined from cDNA clones taken at random in libraries prepared from different sources of plant material: developing siliques, etiolated seedlings, flower buds, and cultured cells. Eight hundred and ninety-five different genes could be identified, 32% of which showed significant similarity to existing sequences in *Arabidopsis* and an array of

Genes Galore: A Summary of Methods for Accessing Results from Large-Scale Partial Sequencing of Anonymous *Arabidopsis* cDNA Clones¹

Tom Newman, Frans J. de Bruijn, Pam Green, Ken Keegstra, Hans Kende, Lee McIntosh, John Ohlrogge, Natasha Raikhel, Shauna Somerville, Mike Thomashow, Ernie Retzel, and Chris Somerville*

Arabidopsis Expressed Sequence Tag Project, Department of Energy Plant Research Laboratory, Michigan State University, East Lansing, Michigan 48824 (T.N., F.J.d.B., P.G., K.K., H.K., L.M., J.O., N.R., M.T.); Computational Biology Center, Medical School, University of Minnesota, 1460 Mayo, UMHC 196, 420 Delaware Street S.E., Minneapolis, Minnesota 55455-0312 (E.R.); and Carnegie Institution of Washington, Department of Plant Biology, 290 Panama Street, Stanford, California 94305-4101 (S.S., C.S.)

En 1990, on pouvait publier des séquences ESTs !!

Volume 9 Issue 1

□ 1: [FEBS Lett.](#) 1997 Mar 24;405(2):129-32.

The *Arabidopsis thaliana* cDNA sequencing projects.

[Delseny M](#), [Cooke R](#), [Raynal M](#), [Grellet F](#).

University of Perpignan, Laboratoire de Physiologie et Biologie, Moléculaire des Plantes, UMR 5545 CNRS, France. delseny@univ-perp.fr

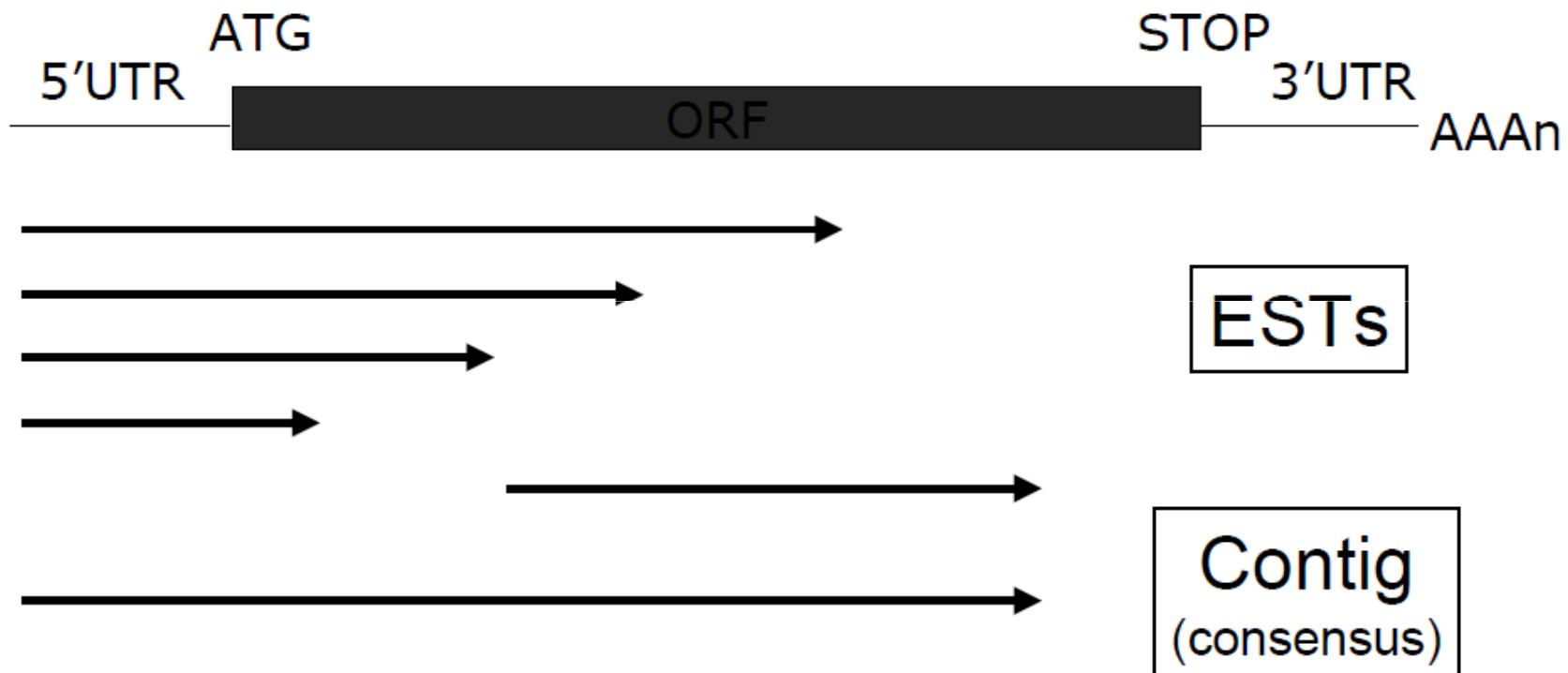
Nearly 30000 *Arabidopsis thaliana* EST (Expressed Sequence Tags) have been produced by a French and an American co. Despite redundancy these sequences tag about half the expected *Arabidopsis* genes. Approximately 40% of the non-redundant can be assigned a putative function by simple homology search. This programme allowed the identification of a large number which would have been very difficult to isolate by other classical techniques. It considerably stimulated many areas of plant biology, the rapid discovery a large number of genes, by revealing multigene families and by allowing the analysis of differential expression of the different members. Finally this programme facilitated construction of physical maps of the chromosomes and opened complete sequencing of the *Arabidopsis* genome and comparative mapping of the major plant crops.

Further progress towards a catalogue of all *Arabidopsis* genes: analysis of a set of 5000 non-redundant ESTs

Richard Cooke^{1,*}, Monique Raynal¹, Michele Laudié¹, Françoise Grellet¹, Michel Delseny¹, Peter-Christian Morris², Danièle Guerrier², Jérôme Giraudat², Françoise Quigley³, Gérard Clabault³, You-Fang Li³, Régis Mache³, Micheline Krivitzky⁴, Isabelle Jean-Jacques Gy⁴, Martin Kreis⁴, Alain Lecharny⁴, Yves Parmentier⁵, Jacqueline Marbach⁵, Jacqueline Fleck⁵, Bernadette Clément⁶, Gabriel Philipp⁶, Christine Hervé⁷, Claude Bardet⁷, Dominique Tremousaygue⁷, Bernard Lescure⁷, Christophe Lacomme⁸, Dominique Roby⁸, Marie-Françoise Jourjon⁹, Patrick Chabrier⁹, Jean-Louis Charpenteau⁹, Thierry Desprez¹⁰, Joëlle Amselem¹⁰, Helen Chiapello¹⁰ and Herman Höfte¹⁰

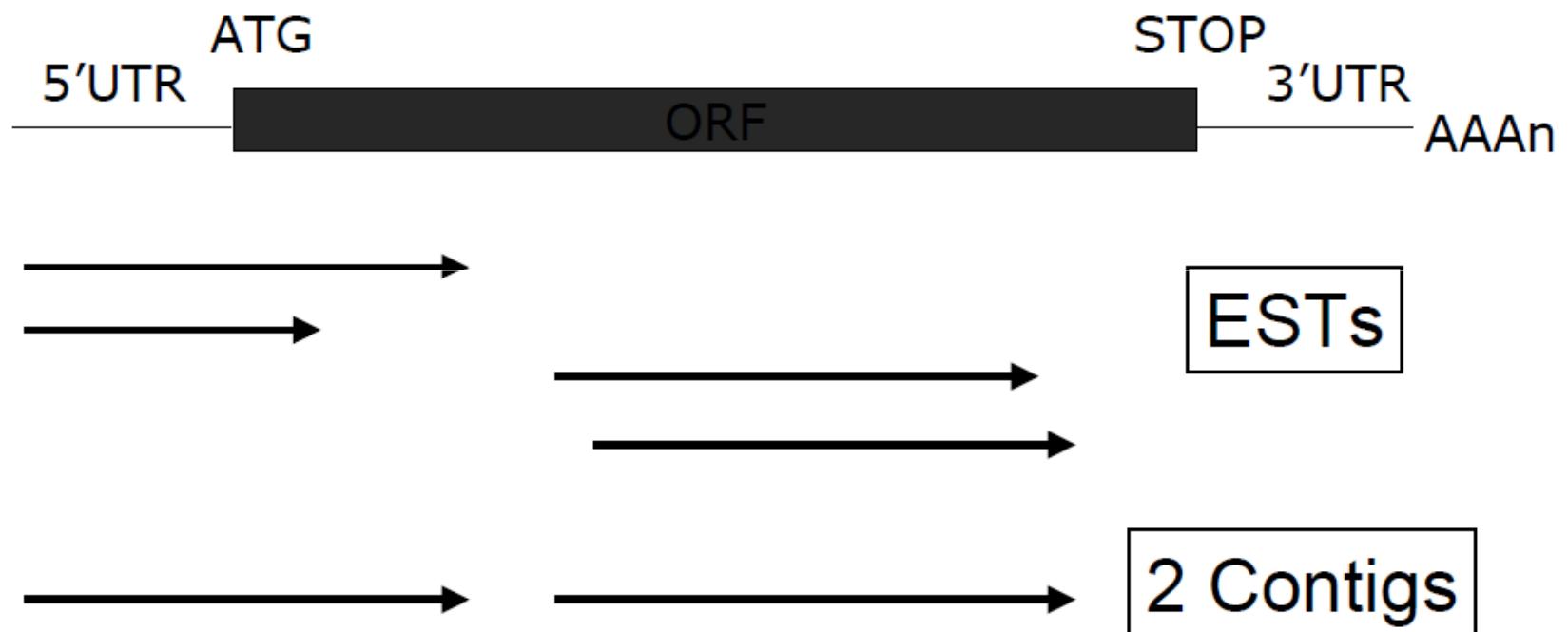
Nearly 7000 *Arabidopsis thaliana*-expressed sequence tags (ESTs) from 10 cDNA libraries have been sequenced, of which almost 5000 non-redundant tags have been submitted to the EMBL data bank. The quality of the cDNA libraries used is analysed. Similarity searches in international protein data banks have allowed the detection of significant similarities to a wide range of proteins from many organisms. Alignment with ESTs from the rice systematic sequencing project has allowed the detection of amino acid motifs which are conserved between the two organisms, thus identifying tags to genes encoding highly conserved proteins. These genes are candidates for a common framework in genome mapping projects in different plants.

Example: EST contig



Donc de plusieurs ESTs → 1 contig

Example: nonoverlapping ESTs



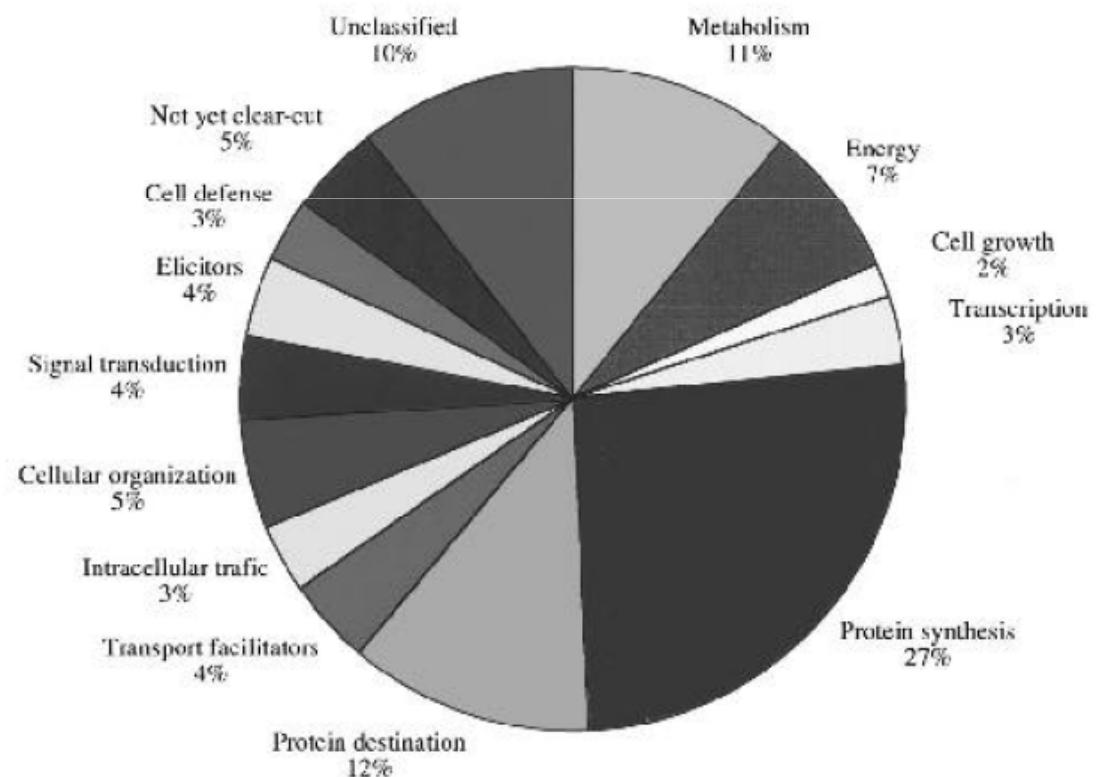
Mais aussi ... plusieurs ESTs → ... plusieurs contigs

De la séquence de l'Est à son annotation et à la base de données

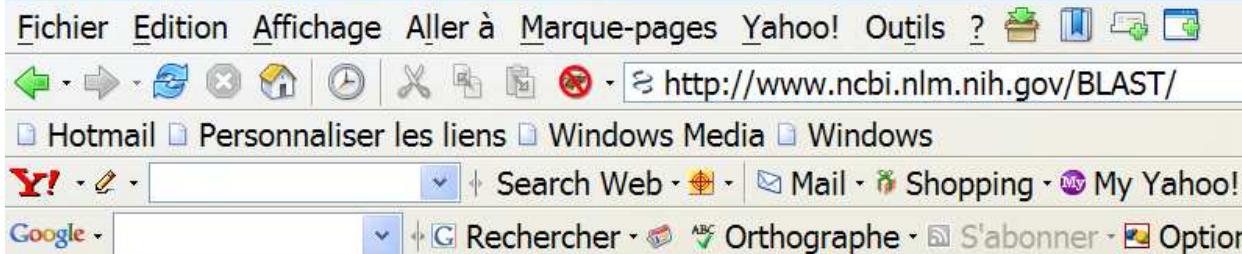
Functional categories: example

01 Metabolism

- 01.01 Amino acid metabolism
- 01.02 Nitrogen and sulfur metabolism
- 01.03 Nucleotide metabolism
- 01.04 Phosphate metabolism
- 01.05 Carbohydrate metabolism
- 01.06 Lipid and sterol metabolism
- 01.07 Biosynthesis of vitamins, coenzymes and prosthetic groups



NCBI BLAST - Mozilla Firefox



The **Basic Local Alignment Search Tool (BLAST)** finds regions of local similarity between sequences. The program compares nucleotide or protein sequences to sequence databases and calculates the statistical significance of matches. BLAST can be used to infer functional and evolutionary relationships between sequences as well as help identify members of gene families.

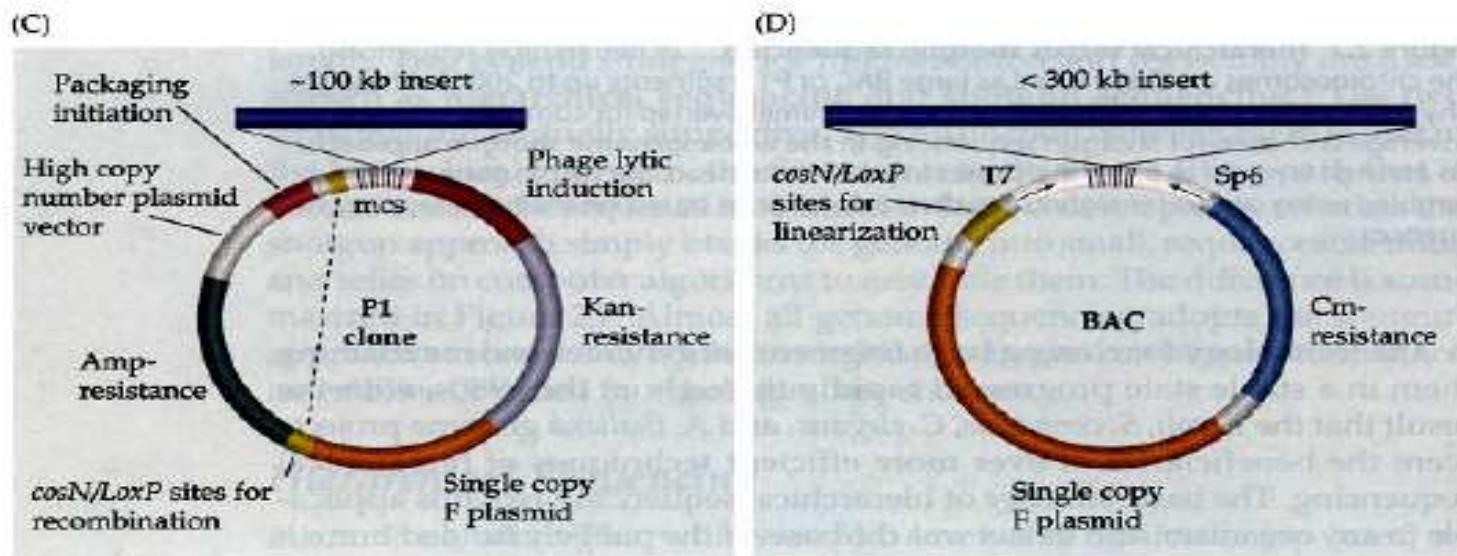
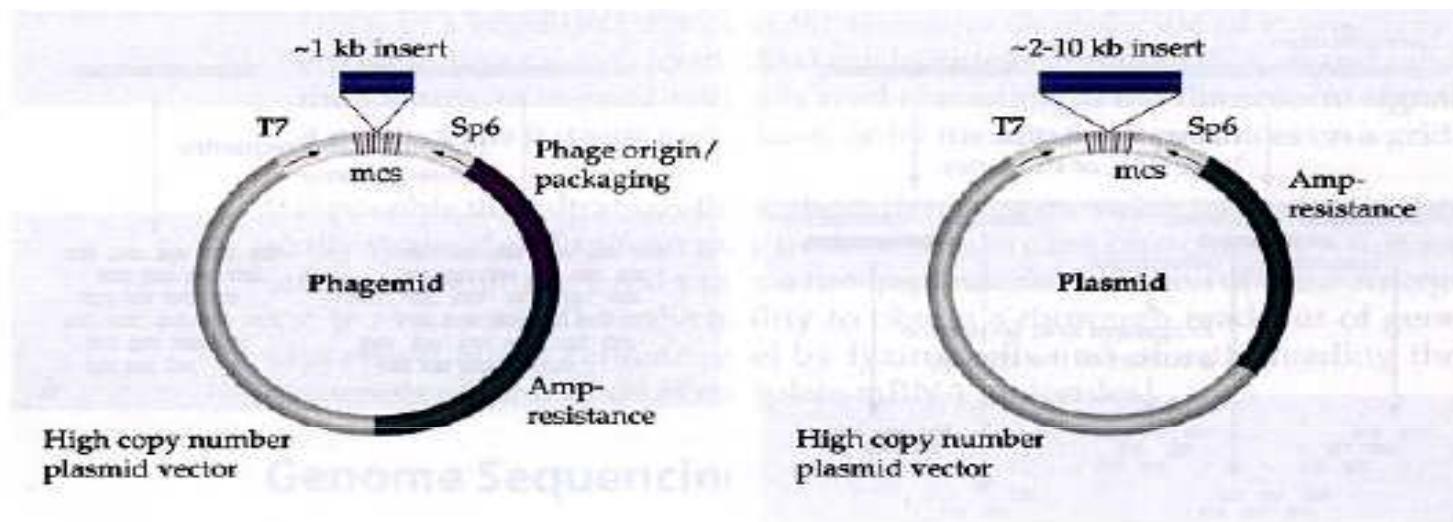
Nucleotide	Protein
<ul style="list-style-type: none">Quickly search for highly similar sequences (megablast)Quickly search for divergent sequences (megablast)Nucleotide-nucleotide BLAST (blastn)Search for short, nearly exact matchesSearch trace archives with megablast or contiguous megablast	<ul style="list-style-type: none">Protein-protein BLAST (blastp)Position-specific iterated and pattern-hit initiated BLAST (PSI- and PHI-BLAST)Search for short, nearly exact matchesSearch the conserved domain database (rpsblast)Protein homology by domain architecture (cdart)
Translated	Genomes
<ul style="list-style-type: none">Translated query vs. protein database (blastx)Protein query vs. translated database (tblastn)Translated query vs. translated database (tblastx)	<ul style="list-style-type: none">Human, mouse, rat, chimp, cow, pig, dog, sheep, catChicken, puffer fish, zebrafishFly, honey bee, other insectsMicrobes, environmental samplesPlants, nematodesFungi, protozoa, other eukaryotes
Special	Meta

Le séquençage des génomes

Major differences between prokaryote and eukaryote genome

	Genome size	Chromosome	Centromere	Telomere	Organization	Repetitive Sequence
Prokaryotes	Small	Single circular, few linear	No	No	High gene density & lack introns	Non or very low
Eukaryotes	Large	Linear	yes	yes	Low gene density & disrupted by introns	High

- Le développement de vecteurs pour le séquençage des génomes



Bacterial Artificial Chromosome (BAC)

- **Bacterial cloning system based on E. coli single copy F factor (PNAS 1992. 89: 8794-8797)**
- **Easy to manipulate**
- **Easy to screen**
- **Stable maintenance of cloned DNA**
- **Nonchimeric**
- **High transformation efficiency**

BAC : 100 à 350kb

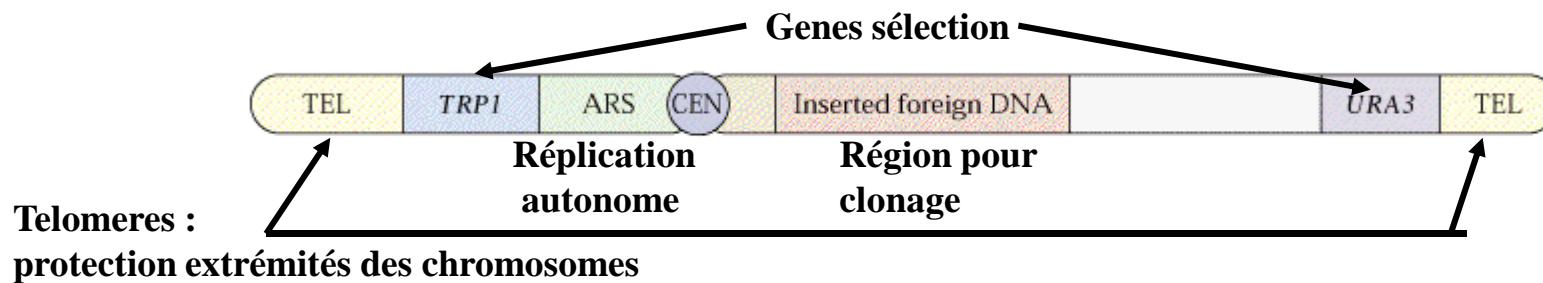
Repose sur épisome F, présent en copie unique chez certaines souches d'E Coli.

→ pb de conservation or pas de recombinaison et monocopie

YAC : de 250kb à 2900 kb !

Propagation possible d'ADN exogene dans la levure sous forme d'un chromosome artificiel

→ faible efficacité de transfo, pb de recombinaison, clonage de grands fragments



Gérome humain :

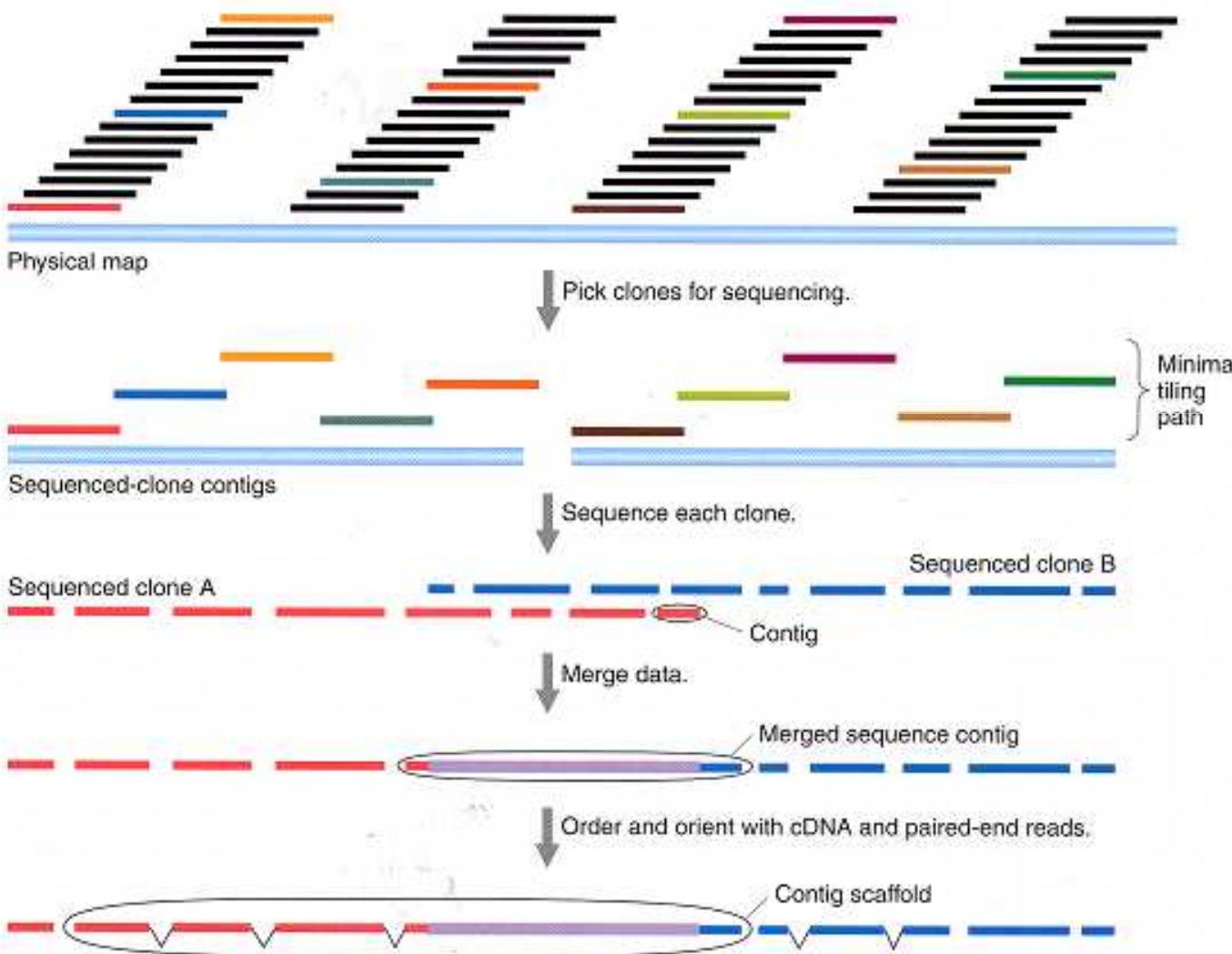
3,2 milliards de nucléotides, soit, en caractères, le contenu de 2000 livres de 500 p.

Pour couvrir une banque génomique humaine :

920 000 clones	phages lambda 15kb
307 000	cosmides 45 kb
92 000	BAC 150kb
14 000	YAC 1000kb

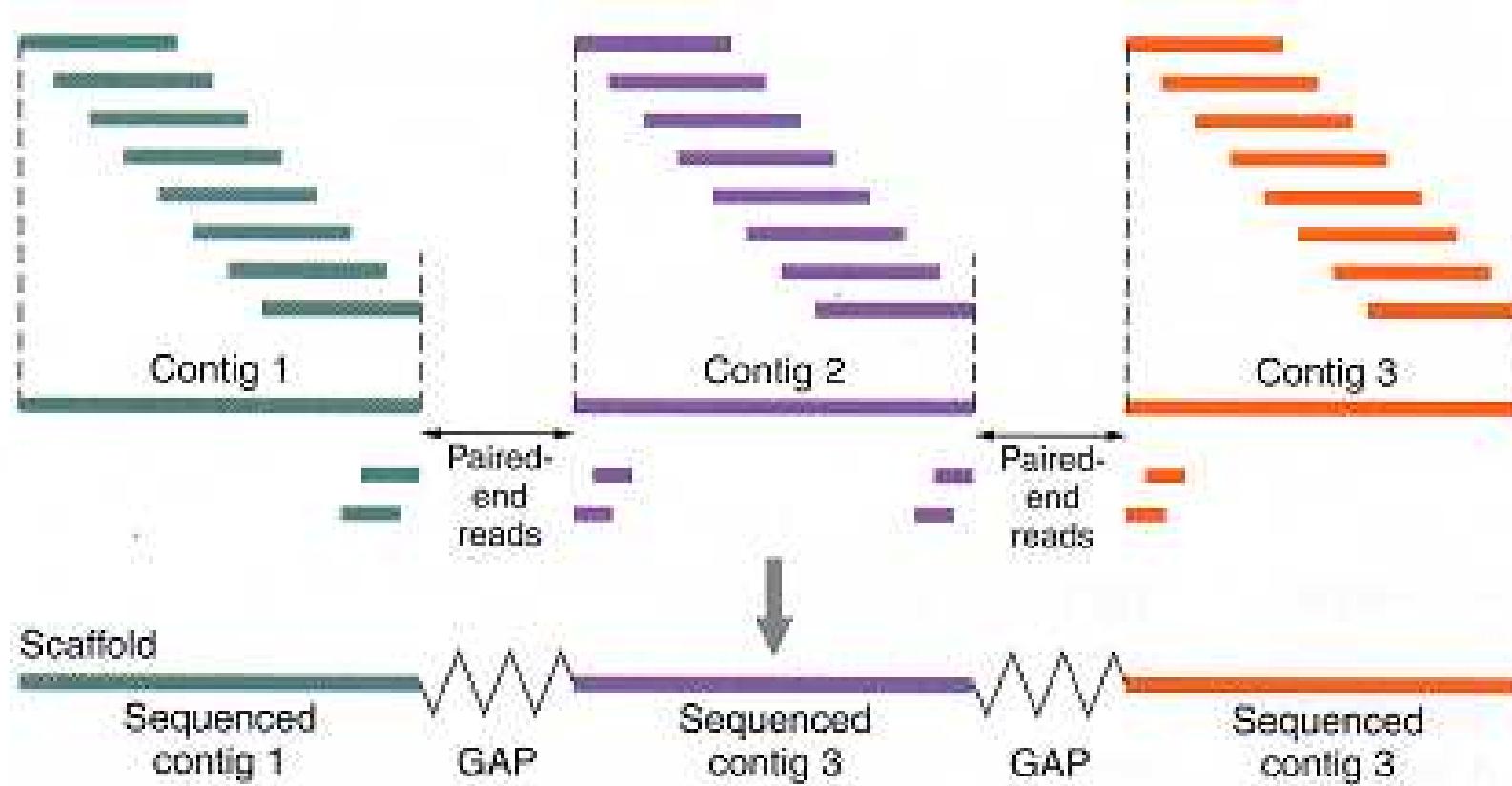
Genome Sequencing by 1st Physical Mapping (Clone à clone)

Extraction A DN, digestion, clonage dans vecteur adapté, assemblage des clones



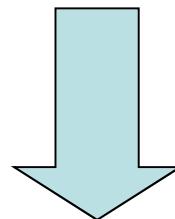
Genome Sequencing by “Shotgun” Method

Extraction ADN, digestion en fragment de 1 à 10 kb,
sequencage (6-8 fois genome), assemblage



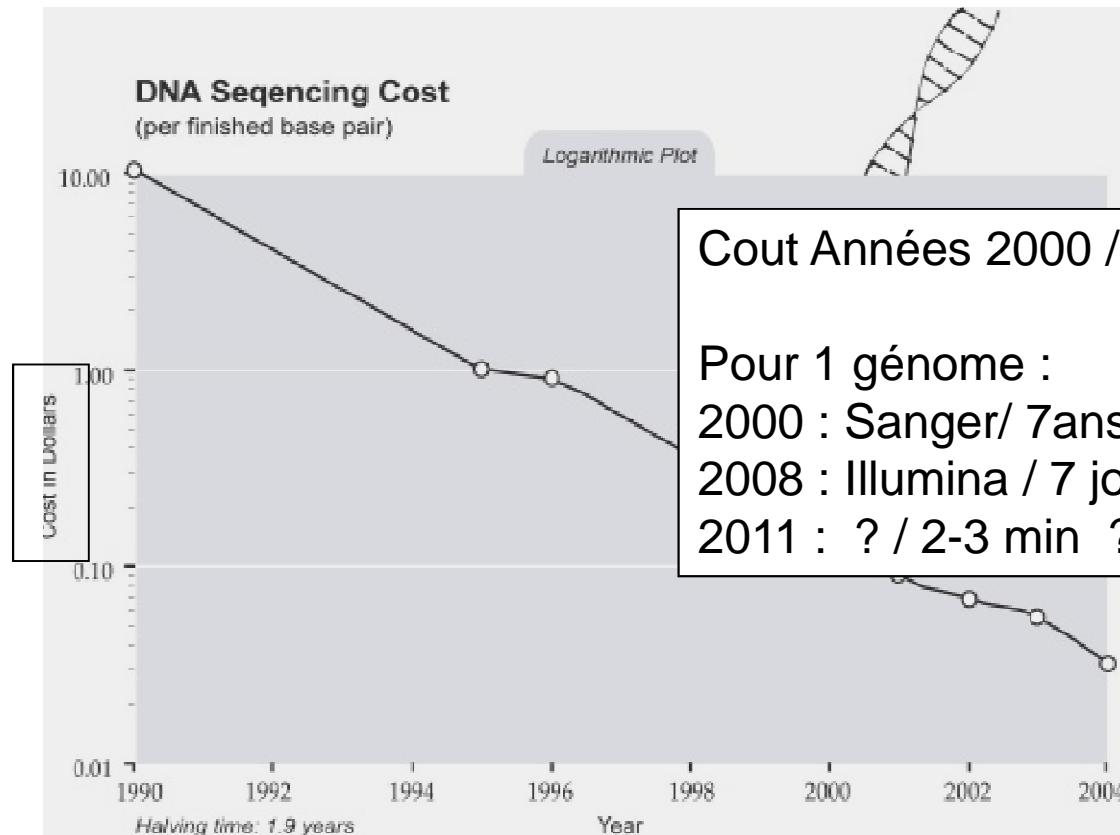
Méthodes de séquençage reposent sur la technique de SANGER

- Points - {
- Nécessite d'isoler un clone par PCR, clonage ... avant de séquencer
 - time consuming → donc cher
- Points + {
- mais génère des fragments de 700 à 1000 bp avec peu d'erreurs



UHTS

ULTRA-HIGH-THROUGHPUT SEQUENCING



Cout Années 2000 /~ 1000 depuis 2008 ...

Pour 1 génome :

2000 : Sanger / 7ans / 500 personnes / 70 000 000 dollars

2008 : Illumina / 7 jours / 2 post-docs / 7000 dollars

2011 : ? / 2-3 min ? / 70 dollars

Table 1 Second-generation DNA sequencing technologies

Feature generation	Sequencing by synthesis	Cout			Taille	
		Cost per megabase	Cost per instrument	Paired ends?	1° error modality	Read-length
454	Emulsion PCR	~\$60	\$500,000	Yes	Indel	250 bp
Solexa	Bridge PCR	~\$2	\$430,000	Yes	Subst.	36 bp
SOLID	Emulsion PCR	~\$2	\$591,000	Yes	Subst.	35 bp
Polonator	Emulsion PCR	~\$1	\$155,000	Yes	Subst.	13 bp
HeliScope	Single molecule	~\$1	\$1,350,000	Yes	Del	30 bp

The pace with which the field is moving makes it likely that estimates for costs and read-lengths will be quickly outdated. Vendors including Roche Applied Science, Illumina, and Applied Biosystems have major upgrade releases currently in progress. Estimated costs-per-megabase are approximate and inclusive only of reagents. Read-lengths are for single tags. Subst., substitutions; indel, insertions or deletions; del, deletions.

Pour quelles applications ?

Table 2 Applications of next-generation sequencing

Category	Examples of applications	Refs
Complete genome resequencing	Comprehensive polymorphism and mutation discovery in individual human genomes	44
Reduced representation sequencing	Large-scale polymorphism discovery	45
Targeted genomic resequencing	Targeted polymorphism and mutation discovery	46–52
Paired end sequencing	Discovery of inherited and acquired structural variation	53,54
Metagenomic sequencing	Discovery of infectious and commensal flora	55
Transcriptome sequencing	Quantification of gene expression and alternative splicing; transcript annotation; discovery of transcribed SNPs or somatic mutations	56–63
Small RNA sequencing	microRNA profiling	64
Sequencing of bisulfite-treated DNA	Determining patterns of cytosine methylation in genomic DNA	60,65,66
Chromatin immunoprecipitation–sequencing (ChIP-Seq)	Genome-wide mapping of protein-DNA interactions	67–70
Nuclease fragmentation and sequencing	Nucleosome positioning	69
Molecular barcoding	Multiplex sequencing of samples from multiple individuals	61,71

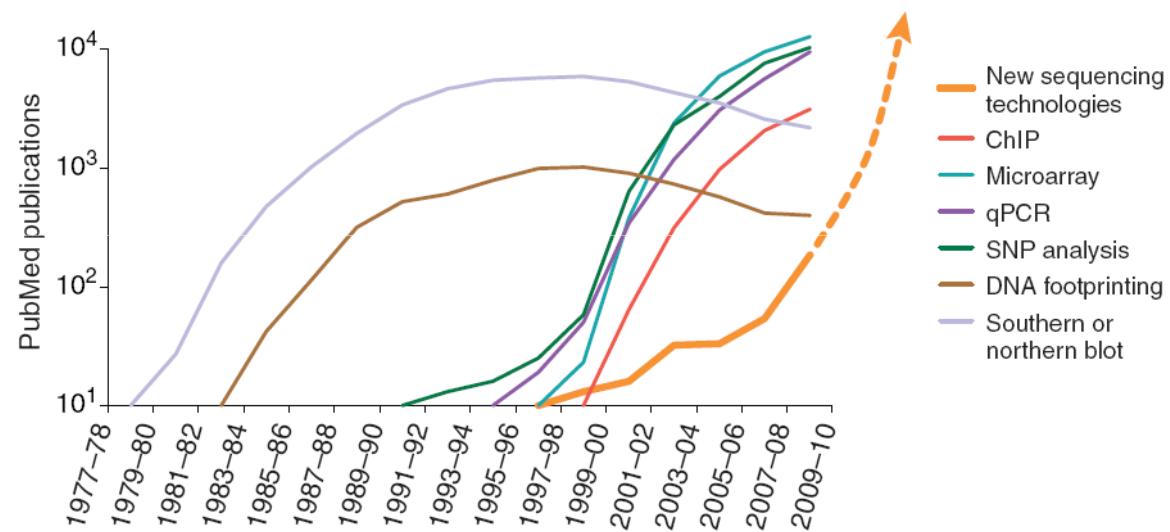
La bonne question ?

What would you do if you could sequence everything?

Avak Kahvejian¹, John Quackenbush² & John F Thompson¹

Nature biotech, 2008

Figure 1 The number of publications with keywords for nucleic acid detection and sequencing technologies. PubMed (<http://www.ncbi.nlm.nih.gov/sites/entrez>) was searched in two-year increments for key words and the number of hits plotted over time. For 2007–2008, results from January 1–March 31, 2008 were multiplied by four and added to those for 2007. Key words used were those listed in the legend except for new sequencing technologies ('next-generation sequencing' or 'high-throughput sequencing'), ChIP ('chromatin immunoprecipitation' or 'ChIP-Chip' or 'ChIP-PCR' or 'ChIP-Seq'), qPCR (TaqMan or qPCR or 'real-time PCR') and SNP analysis (SNPs or 'single-nucleotide polymorphisms' and not nitroprusside (nitroprusside is excluded because sodium nitroprusside is sometimes abbreviated as 'SNP' but is generally unrelated to genetics)).



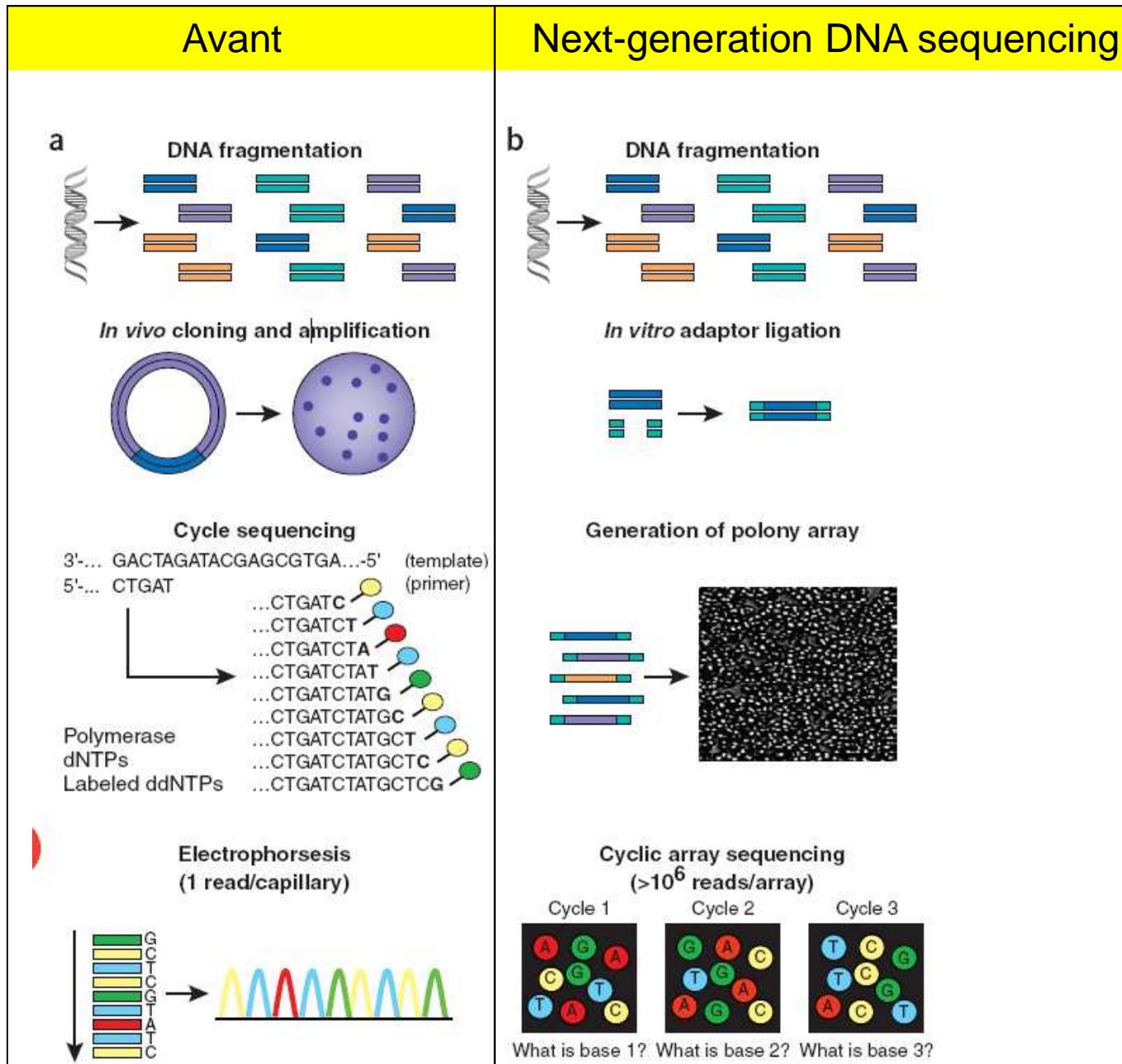
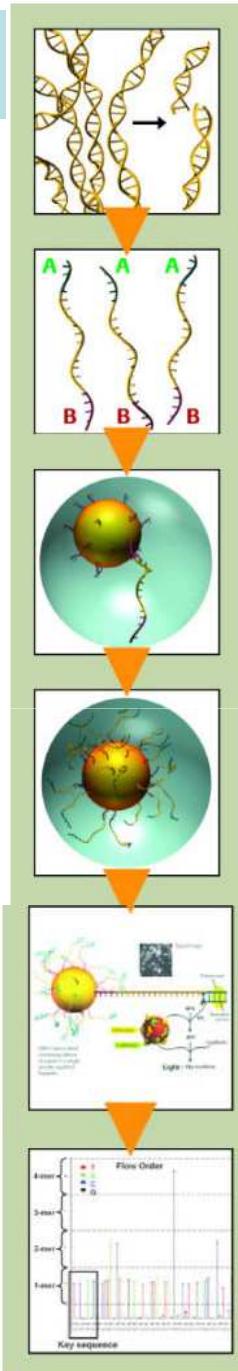


Figure 1 Work flow of conventional versus second-generation sequencing. (a) With high-throughput

454 pyrosequencing



Sample Input and Fragmentation

The Genome Sequencer FLX System supports the sequencing of samples from a wide variety of starting materials including genomic DNA, PCR products, BACs, and cDNA. Samples such as genomic DNA and BACs are fractionated into small, 300- to 800-basepair fragments. For smaller samples, such as small non-coding RNA or PCR amplicons, fragmentation is not required. Instead, short PCR products amplified using Genome Sequencer fusion primers can be used for immobilization onto DNA capture beads as shown below under "One Fragment = One Bead".

The GS FLX System supports multiple sample prep options, [click here to see the full list.](#)

Library Preparation

Using a series of standard molecular biology techniques, short adaptors (A and B) - specific for both the 3' and 5' ends - are added to each fragment. The adaptors are used for purification, amplification, and sequencing steps. Single-stranded fragments with A and B adaptors compose the sample library used for subsequent workflow steps.

One Fragment = One Bead

The single-stranded DNA library is immobilized onto specifically designed DNA Capture Beads. Each bead carries a unique single-stranded DNA library fragment. The bead-bound library is emulsified with amplification reagents in a water-in-oil mixture resulting in microreactors containing just one bead with one unique sample-library fragment.

emPCR (Emulsion PCR) Amplification

Each unique sample library fragment is amplified within its own microreactor, excluding competing or contaminating sequences. Amplification of the entire fragment collection is done in parallel; for each fragment, this results in a copy number of several million per bead. Subsequently, the emulsion PCR is broken while the amplified fragments remain bound to their specific beads.

One Bead = One Read

The clonally amplified fragments are enriched and loaded onto a PicoTiterPlate device for sequencing. The diameter of the PicoTiterPlate wells allows for only one bead per well. After addition of sequencing enzymes, the fluidics subsystem of the Genome Sequencer FLX Instrument flows individual nucleotides in a fixed order across the hundreds of thousands of wells containing one bead each. Addition of one (or more) nucleotide(s) complementary to the template strand results in a chemiluminescent signal recorded by the CCD camera of the Genome Sequencer FLX Instrument. For a detailed explanation of this reaction see [Sequencing Chemistry](#).

Data Analysis

The combination of signal intensity and positional information generated across the PicoTiterPlate device allows the software to determine the sequence of more than 1,000,000 individual reads per 10-hour instrument run simultaneously. For sequencing-data analysis, three different bioinformatics tools are available supporting the following applications: de novo assembly up to 400 megabases; resequencing genomes of any size; and amplicon variant detection by comparison with a known reference sequence.

Fragmentation ADN

Ajout d'adaptateur

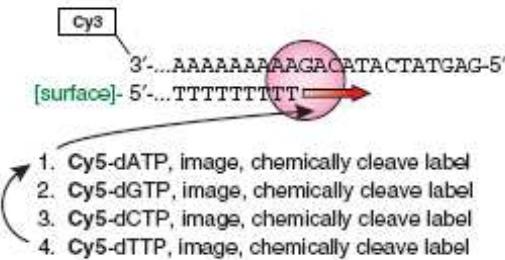
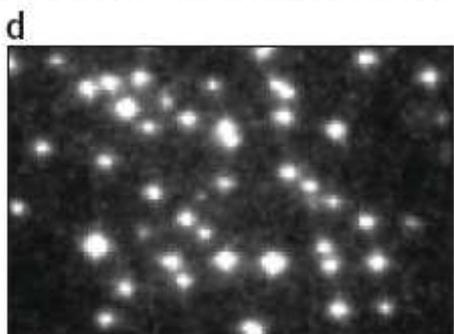
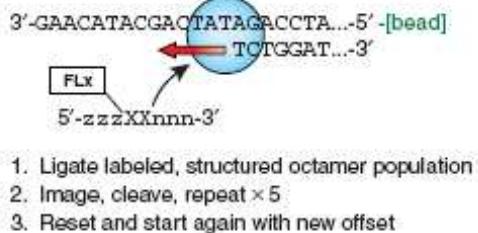
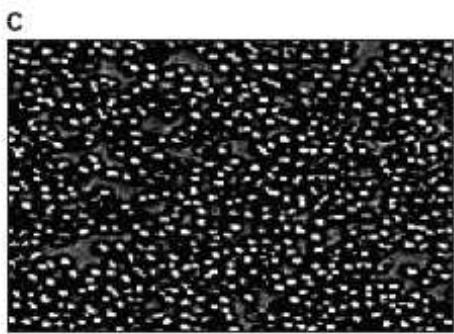
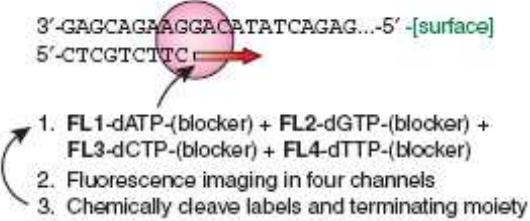
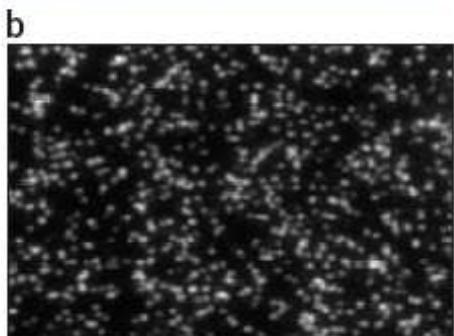
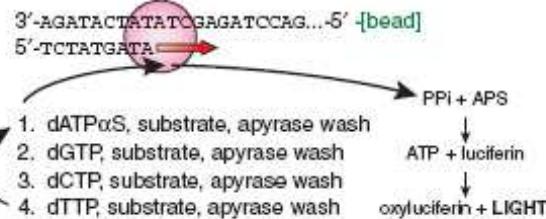
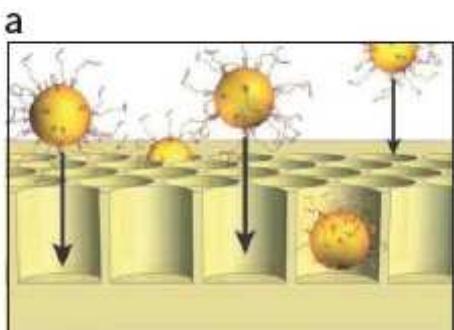
Couplage a des billes
(1 fragment = 1 bille)
Puis billes dans uplaque

PCR

Pyrosequencage : A chaque cycle on ajoute 1 dNTP (dATP, ou dGTP, ou dCTP, ou dTTP)

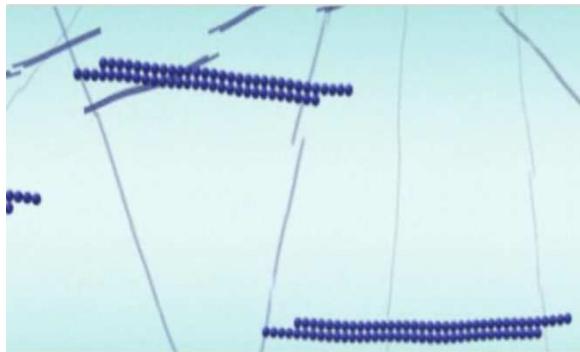
Si il est incorporé alors pyrophosphate est libéré et 1 signal est enregistré grace à un systeme luminescent

Illustration 454

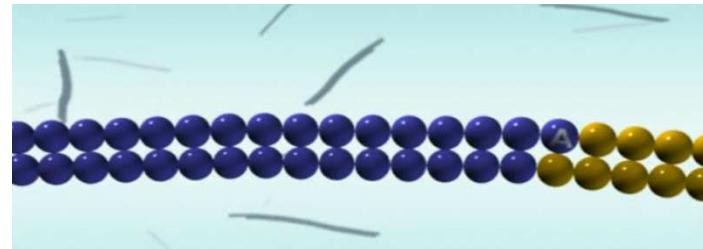


ILLUMINA

1. DNA fragmentation →



2. ajout d'adaptateur →



3. Chips avec oligos complementaires →
aux adaptateurs



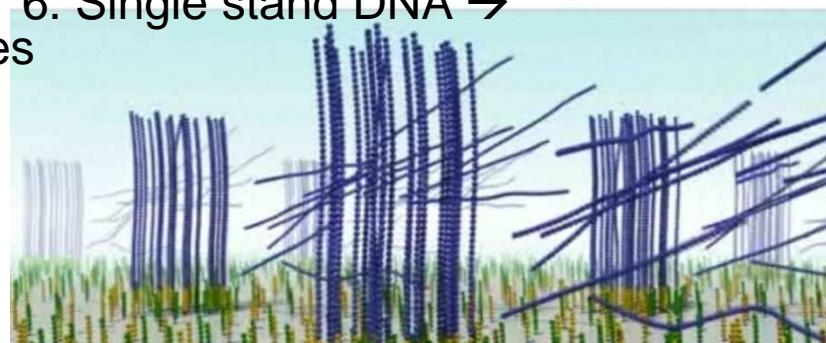
4. Fixation ADN au chip



5. PCR → molécules ADN identiques
Localisées à 1 endroit du chip



6. Single stand DNA →



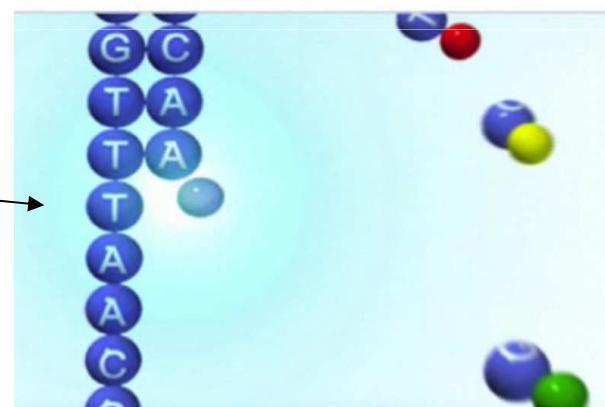
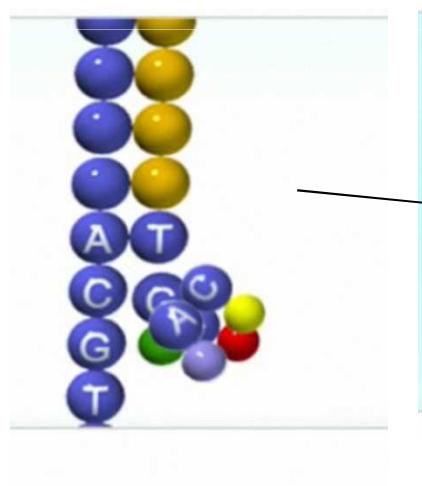
Sequencage

7. Ajout des 4 dNTP (chacun avec un fluorochrome spécifique)

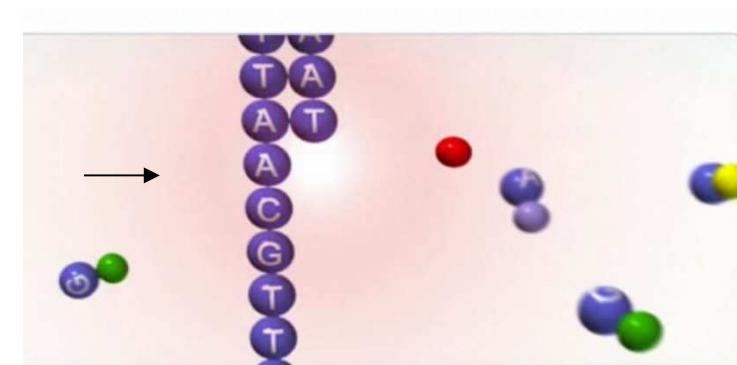
Séquencage en Synthèse

8. Fixation du dNTP libération du fluorochrome

9. Détection

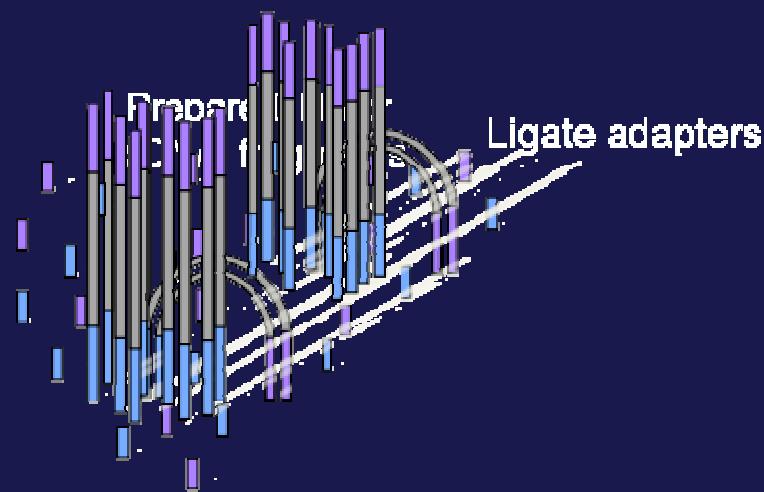


A...



AT... etc ...

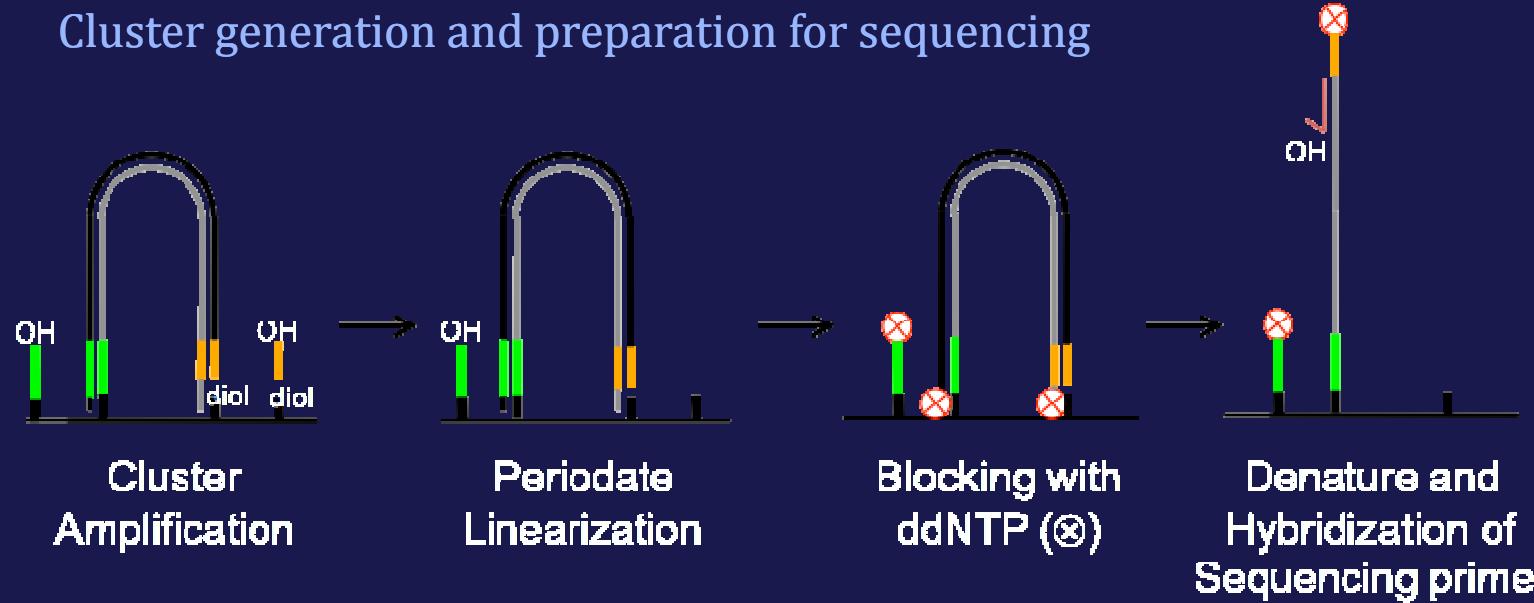
Sample generation and cluster generation



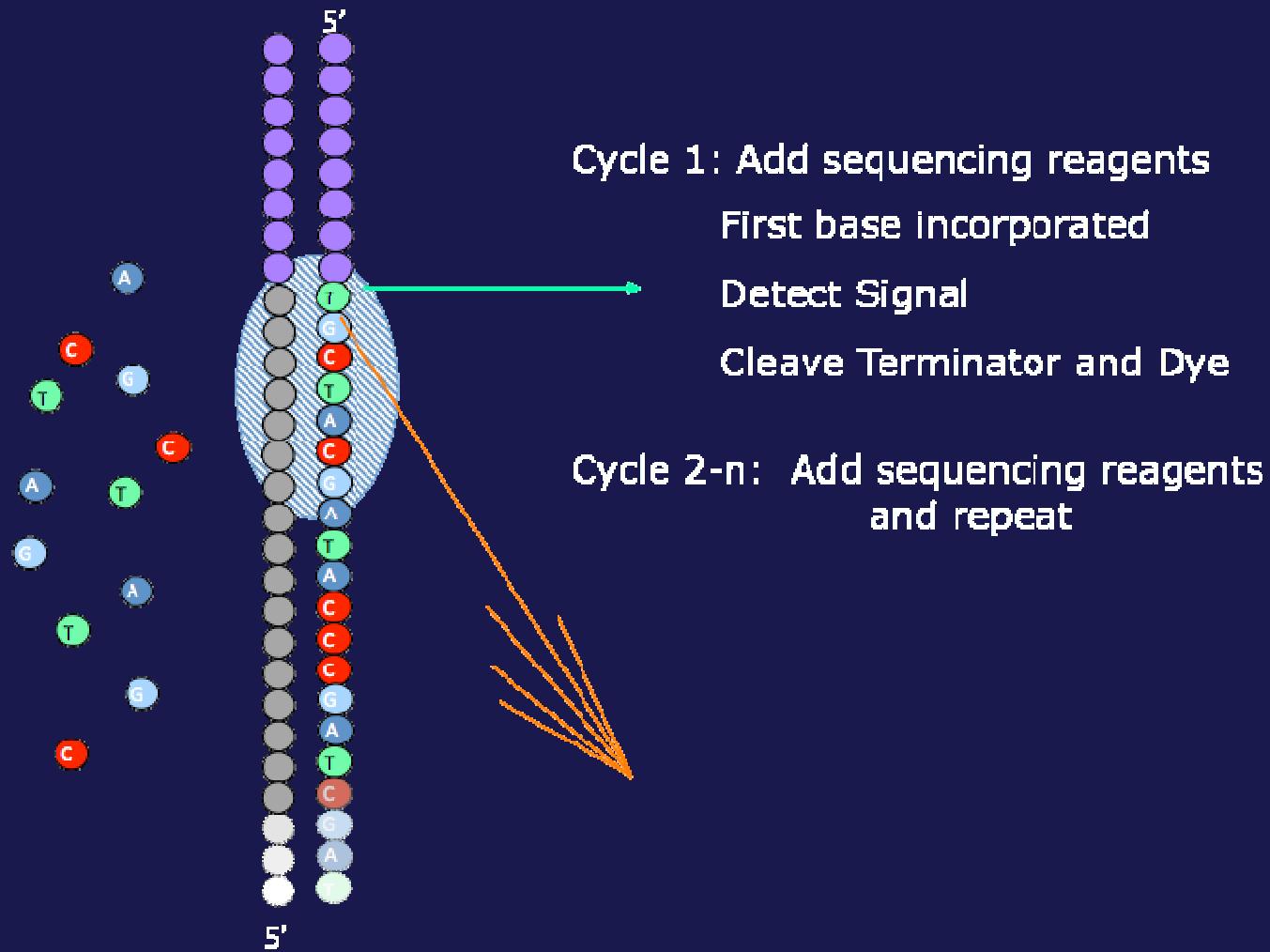
Attach single molecules to surface

Amplify to form clusters

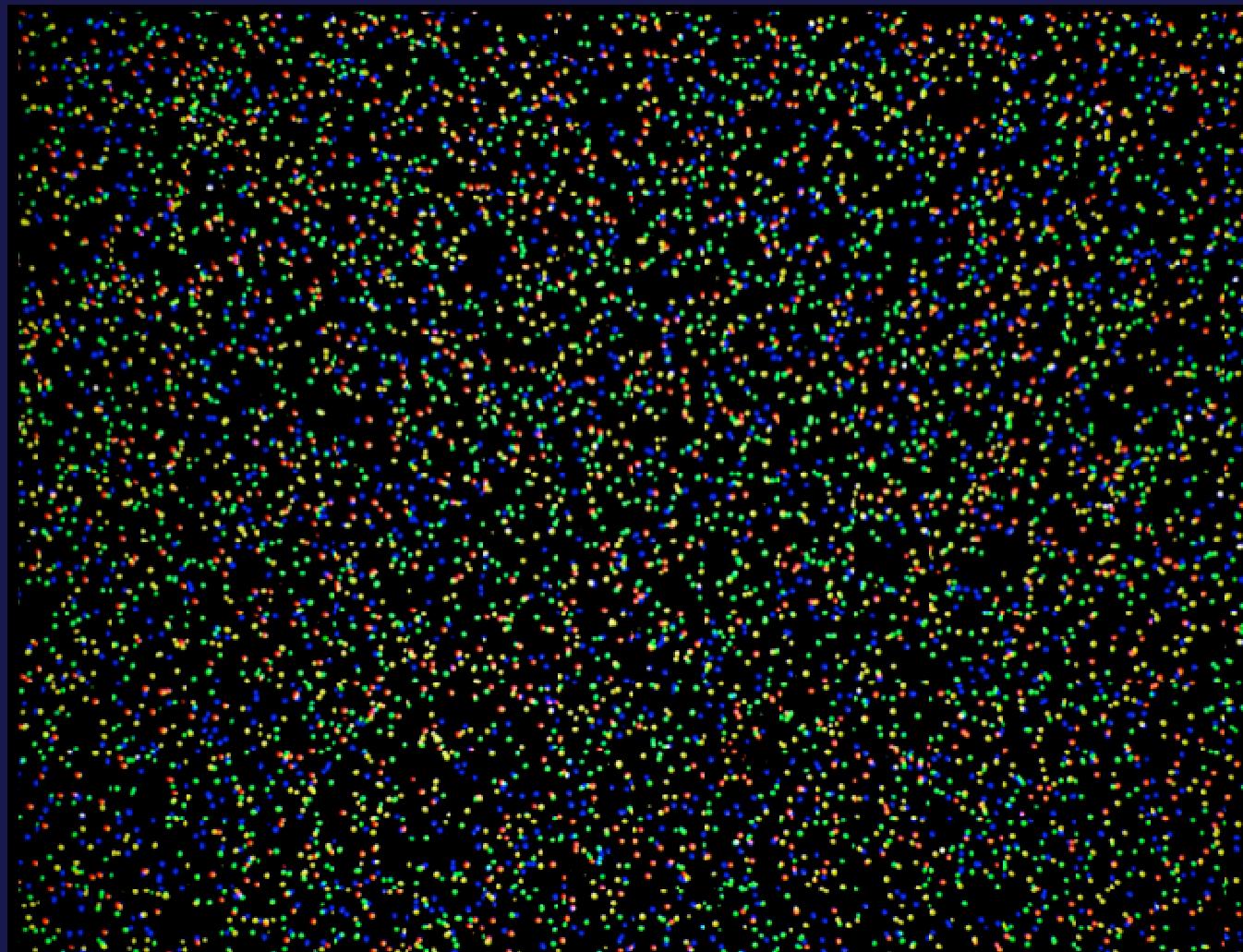
Cluster generation and preparation for sequencing



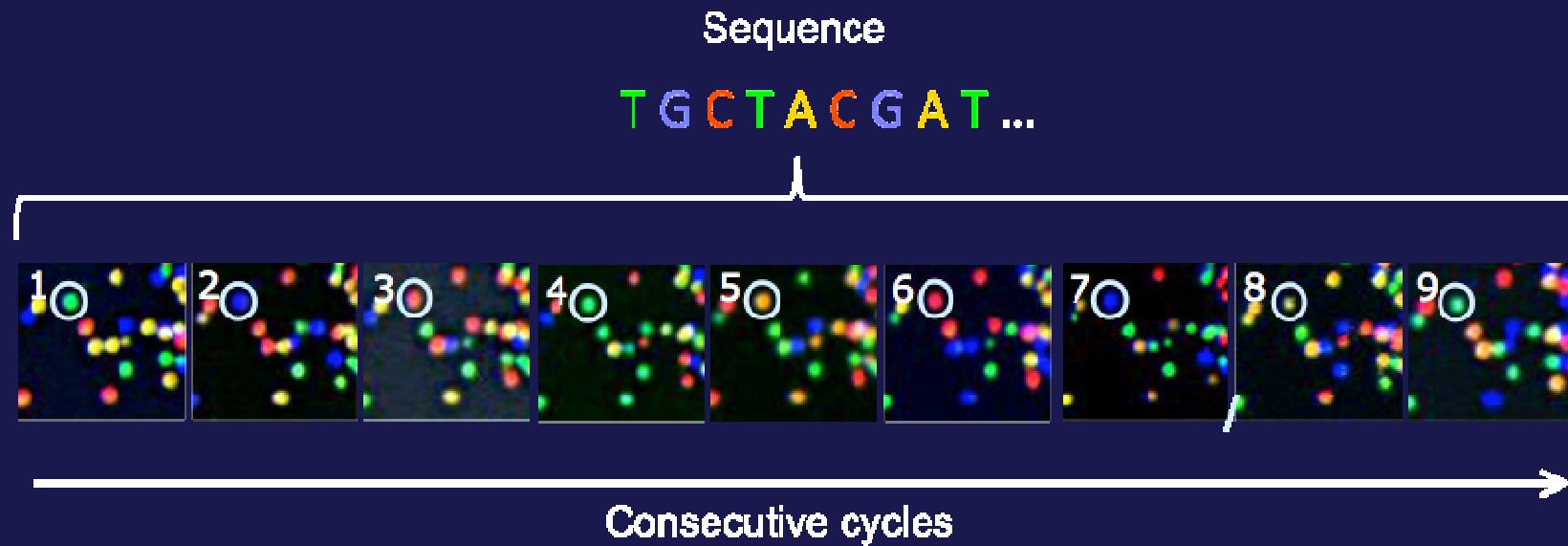
Sequencing by Synthesis (SBS)



Sequencing by Synthesis (SBS)



Base Calling



The identity of each base of a cluster is read from stacked sequential images

An integrated semiconductor device enabling non-optical genome sequencing

Jonathan M. Rothberg¹, Wolfgang Hinz¹, Todd M. Rearick¹, Jonathan Schultz¹, William Mileski¹, Mel Davey¹, John H. Leamon¹, Kim Johnson¹, Mark J. Milgrew¹, Matthew Edwards¹, Jeremy Hoon¹, Jan F. Simons¹, David Marran¹, Jason W. Myers¹, John F. Davidson¹, Annika Branting¹, John R. Nobile¹, Bernard P. Puc¹, David Light¹, Travis A. Clark¹, Martin Huber¹, Jeffrey T. Branciforte¹, Isaac B. Stoner¹, Simon E. Cawley¹, Michael Lyons¹, Yutao Fu¹, Nils Homer¹, Marina Sedova¹, Xin Miao¹, Brian Reed¹, Jeffrey Sabina¹, Erika Feierstein¹, Michelle Schorn¹, Mohammad Aljmary¹, Eileen Dimalanta¹, Devin Dressman¹, Rachel Kasinskas¹, Tanya Sokolsky¹, Jacqueline A. Fidanza¹, Eugeni Namsaraev¹, Kevin J. McKernan¹, Alan Williams¹, G. Thomas Roth¹ & James Bustillo¹

DNA sequencing technology in semiconductor able to directly perform non-optical DNA sequencing of genomes.

Sequence data are obtained by directly sensing the ions produced by template-directed DNA polymerase synthesis using all-natural nucleotides on this massively parallel semiconductor-sensing device or ion chip.

The ion chip contains ion-sensitive, field-effect transistor-based sensors in perfect register with 1.2 million wells, which provide confinement and allow parallel, simultaneous detection of independent sequencing reactions.

We show the performance of the system by sequencing three bacterial genomes, its robustness and scalability by producing ion chips with up to 10 times as many sensors and sequencing a human genome.

Nature 2011, 475, 348-352

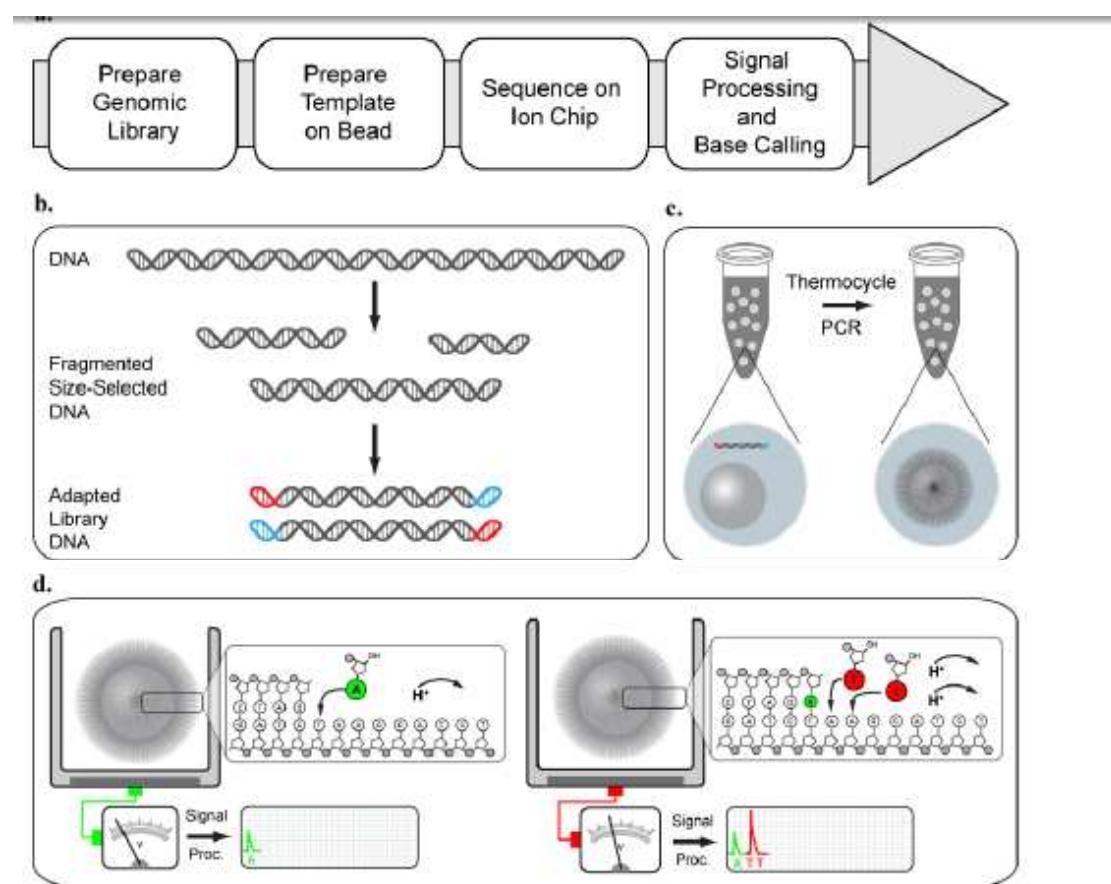
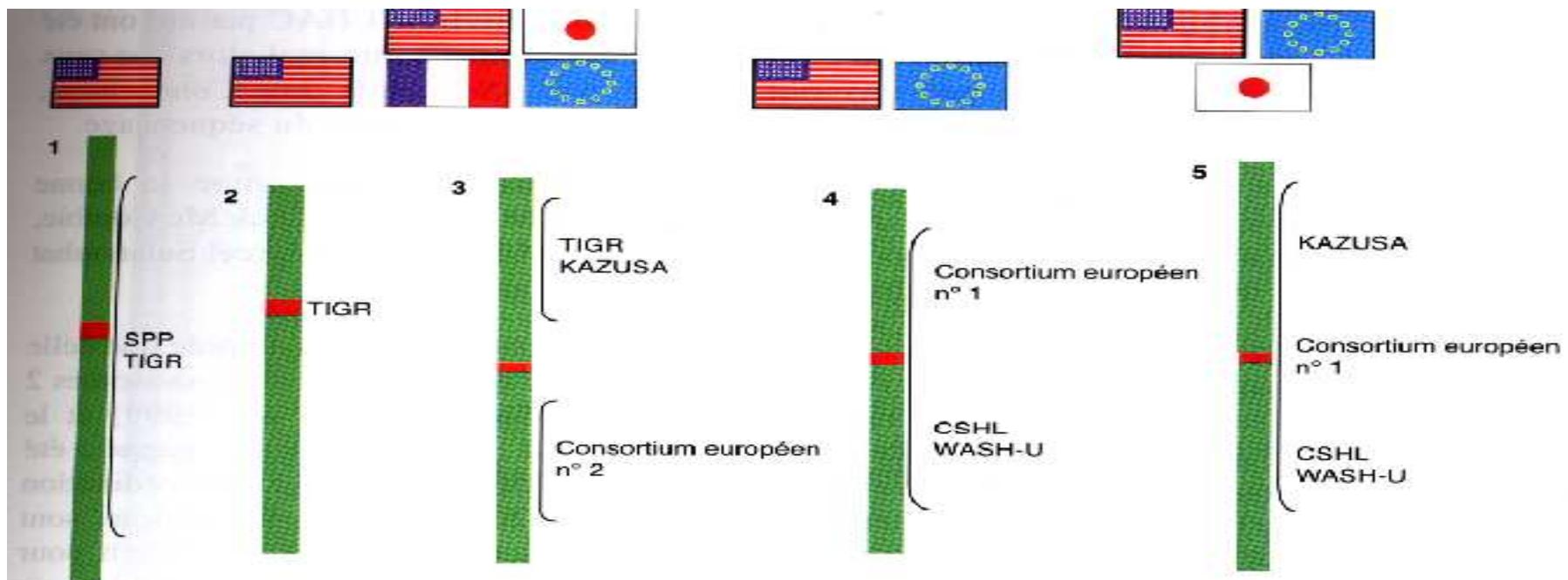


Figure S1 Process overview

a, Overview of ion sequencing work flow. b, Prepare genomic library, DNA is fragmented, sized, and forward and reverse adapters ligated. c, Amplify Template on bead, adapter-ligated libraries are clonally amplified onto beads. A magnetic bead-based enrichment process selects template-carrying beads. d, Sequence on ion chip, sequencing primers and DNA polymerase are bound to the template-carrying beads, beads are pipetted into the chip's loading port. The chip is installed in the sequencing instrument; all four nucleotides cyclically flowed in an automated 2-hour run. Signal processing, software converts the raw data into measurements of incorporation in each well for each successive nucleotide flow. After bases are called, each read is passed through a filter to exclude low-accuracy reads and per-base quality values are predicted.

Séquençage des génomes : Le Yalta des chromosomes chez *Arabidopsis*

Arabidopsis thaliana



SPP Consortium : Stanford DNA Sequencing & Technology Center, Stanford University, *Arabidopsis thaliana* Genome Center, University of Pennsylvania, Plant Gene Expression Center, USDA.

TIGR : The Institute for Genomic Research, Washington.

KAZUSA DNA Research Institute, Japan.

Consortium européen n°1 : coordinateur : John Innes Center, Norwich.

Consortium européen n°2 : coordinateur : Génoscope, Évry.

Wash-U : Washington University, St Louis.

CSHL : Cold Spring Harbor Laboratory, NY.

Figure 6. Répartition du séquençage des 5 chromosomes d'*Arabidopsis thaliana* entre les différents centres de séquençage de l'AGI.

Stratégie clone à clone

Séquençage des génomes : Le Yalta des chromosomes

Oriza sativa

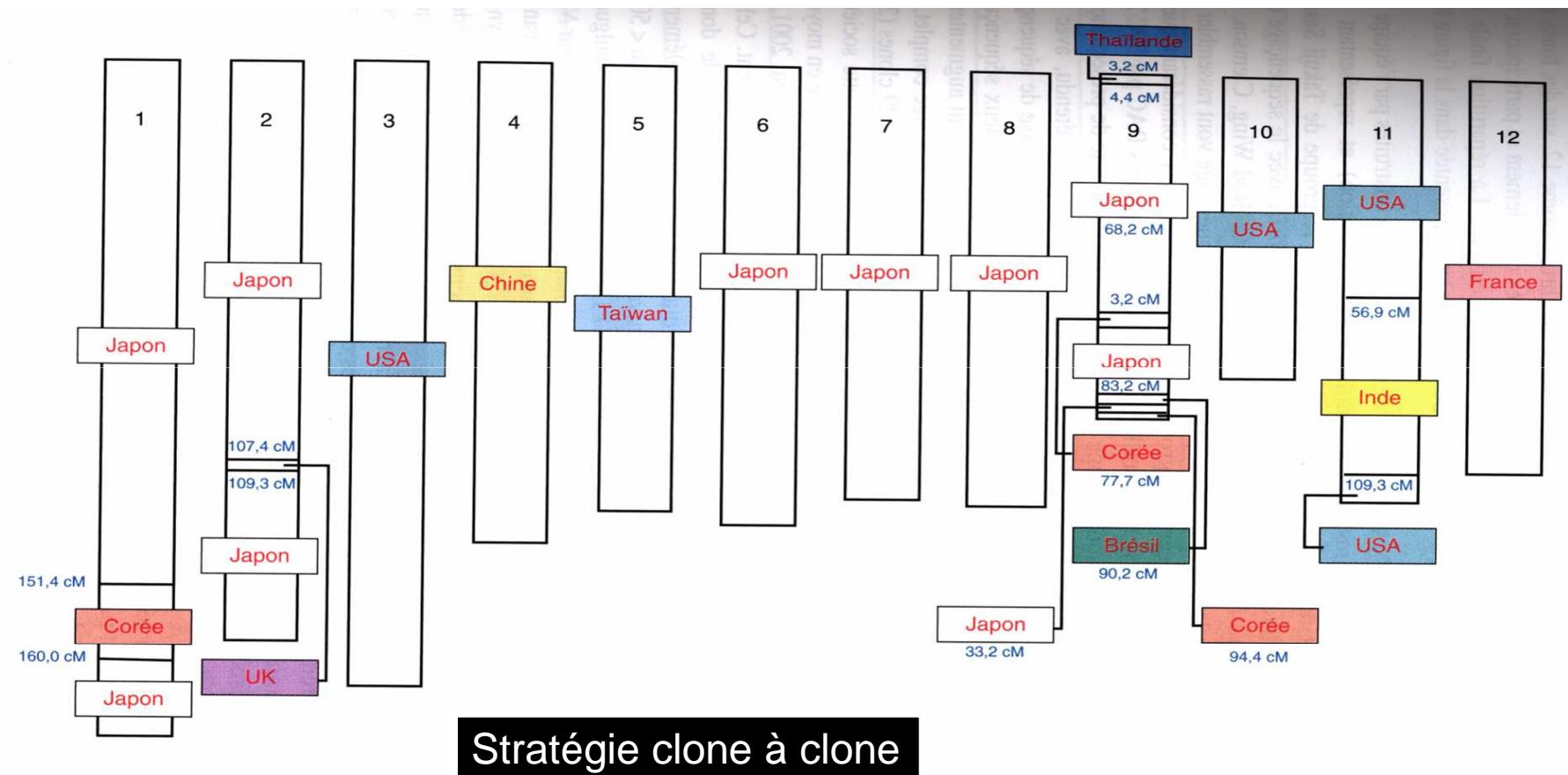


Figure 8. État en novembre 2002 de la répartition du séquençage des 12 chromosomes du riz entre les différents pays membres de l'IRGSP.

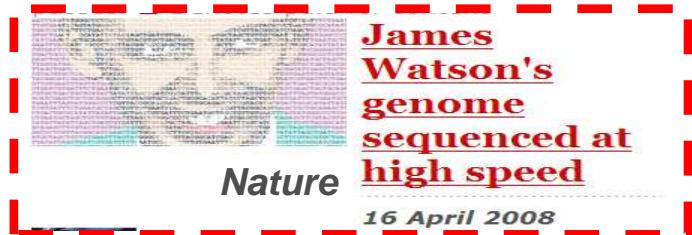
Le séquençage du Génome de Watson , n'est pas le 1er génome humain séquencé (C. Venter en 2007) mais ...

nature
Vol 452 | 17 April 2008 | doi:10.1038/nature06884

LETTERS

The complete genome of an individual by massively parallel DNA sequencing

David A. Wheeler^{1,*}, Maithreyan Srinivasan^{2,*}, Michael Egholm^{2,*}, Yufeng Shen¹, Lei Chen¹, Amy McGuire³, Wen He², Yi-Ju Chen², Vinod Makhijani², G. Thomas Roth², Xavier Gomes², Karrie Tartaro^{2†}, Faheem Niazi², Cynthia L. Turcotte², Gerard P. Iryzk², James R. Lupski^{4,5,6}, Craig Chinault⁴, Xing-zhi Song¹, Yue Liu¹, Ye Yuan¹, Lynne Nazareth¹, Xiang Qin¹, Donna M. Muzny¹, Marcel Margulies², George M. Weinstock^{1,4}, Richard A. Gibbs^{1,4} & Jonathan M. Rothberg^{2†}



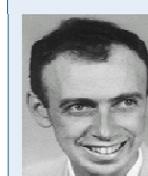
Stratégie UHTS

J. Craig Venter



James Dewey Watson

Biochemiker



* 26. April 1928 in Chicago, USA

QUICKER, SMALLER, CHEAPER

Genome sequenced (publication year)	HGP (2003)	Venter (2007)	Watson (2008)
Time taken (start to finish)	13 years	4 years	4.5 months
Number of scientists listed as authors	> 2,800	31	27
Cost of sequencing (start to finish)	\$2.7 billion	\$100 million	< \$1.5 million
Coverage	8-10 ×	7.5 ×	7.4 ×
Number of institutes involved	16	5	2
Number of countries involved	6	3	1

Human Genome Project

©2008 Nature Publishing Group

6 Nov 2008



Current Issue

Volume 455
Number 7218

This week's news

- [What Obama's win means for science](#)
- [HIV vaccine failure explained?](#)
- [Blink and you'll miss it](#)
- [More from Nature News](#)

Latest Research

nature

Vol 456 | 6 November 2008 | doi:10.1038/nature07484

ARTICLES

The diploid genome sequence of an Asian individual

Jun Wang^{1,2,3,4*}, Wei Wang^{1,3*}, Ruiqiang Li^{1,3,4*}, Yingrui Li^{1,5,6*}, Geng Tian^{1,7}, Laurie Goodman¹, Wei Fan¹, Junqing Zhang¹, Jun Li¹, Juanbin Zhang¹, Yiran Guo^{1,7}, Binxiao Feng¹, Heng Li^{1,8}, Yao Lu¹, Xiaodong Fang¹, Huiqing Liang¹, Zhenglin Du¹, Dong Li¹, Yiqing Zhao^{1,7}, Yujie Hu^{1,7}, Zhenzhen Yang¹, Hancheng Zheng¹, Ines Hellmann⁹, Michael Inouye⁸, John Pool⁹, Xin Yi^{1,7}, Jing Zhao¹, Jinjie Duan¹, Yan Zhou¹, Junjie Qin^{1,7}, Lijia Ma^{1,7}, Guoqing Li¹, Zhentao Yang¹, Guojie Zhang^{1,7}, Bin Yang¹, Chang Yu¹, Fang Liang^{1,7}, Wenjie Li¹, Shaochuan Li¹, Dawei Li¹, Peixiang Ni¹, Jue Ruan^{1,7}, Qibin Li^{1,7}, Hongmei Zhu¹, Dongyuan Liu¹, Zhike Lu¹, Ning Li^{1,7}, Guangwu Gu¹, Jianguo Zhang¹, Jia Ye¹, Lin Fang¹, Qin Hao^{1,7}, Quan Chen^{1,5}, Yu Liang^{1,7}, Yeyang Su^{1,7}, A. san^{1,7}, Cuo Ping^{1,7}, Shuang Yang¹, Fang Chen^{1,7}, Li Li¹, Ke Zhou¹, Hongkun Zheng^{1,4}, Yuanyuan Ren¹, Ling Yang¹, Yang Gao^{1,6}, Guohua Yang^{1,2}, Zhuo Li¹, Xiaoli Feng¹, Karsten Kristiansen⁴, Gane Ka-Shu Wong^{1,10}, Rasmus Nielsen⁹, Richard Durbin⁸, Lars Bolund^{1,11}, Xiuqing Zhang^{1,6}, Songgang Li^{1,2,5}, Huanming Yang^{1,2,3} & Jian Wang^{1,2,3}

Here we present the first diploid genome sequence of an Asian individual. The genome was sequenced to 36-fold average coverage using massively parallel sequencing technology. We aligned the short reads onto the NCBI human reference genome to 99.97% coverage, and guided by the reference genome, we used uniquely mapped reads to assemble a high-quality consensus sequence for 92% of the Asian individual's genome. We identified approximately 3 million single-nucleotide polymorphisms (SNPs) inside this region, of which 13.6% were not in the dbSNP database. Genotype analysis showed that SNP identification had high accuracy and consistency, indicating the high sequence quality of the assembly. We also carried out heterozygote phasing and haplotype prediction against HapMap CHB and JPT haplotypes (Chinese and Japanese, respectively), sequence comparison with the two available individual genomes (J. D. Watson and C. Venter), and structural variation identification. These variations were considered for their potential biological impact. Sequence data and analyses demonstrate the potential usefulness of next-generation sequencing technologies for personal genomics.

Vol 456 | 6 November 2008 | doi:10.1038/nature07517

nature

ARTICLES

Accurate whole human genome sequencing using reversible terminator chemistry

A list of authors and their affiliations appears at the end of the paper

Vol 456 | 6 November 2008 | doi:10.1038/nature07485

nature

ARTICLES

DNA sequencing of a cytogenetically normal acute myeloid leukaemia genome

Timothy J. Ley^{1,2,3,4*}, Elaine R. Mardis^{2,3*}, Li Ding^{2,3}, Bob Fulton³, Michael D. McLellan³, Ken Chen³, David Dooling³, Brian H. Dunford-Shore³, Sean McGrath³, Matthew Hickenbotham³, Lisa Cook³, Rachel Abbott³, David E. Larson³, Dan C. Koboldt³, Craig Pohl³, Scott Smith³, Amy Hawkins³, Scott Abbott³, Devin Locke³, LaDeana W. Hillier^{3,8}, Tracie Miner³, Lucinda Fulton³, Vincent Magrini^{2,3}, Todd Wylie³, Jarret Glasscock³, Joshua Conyers³, Nathan Sander³, Xiaoqi Shi³, John R. Osborne³, Patrick Minx³, David Gordon⁸, Asif Chinwalla³, Yu Zhao¹, Rhonda E. Ries¹, Jacqueline E. Payton⁵, Peter Westervelt^{1,4}, Michael H. Tomasson^{1,4}, Mark Watson^{3,4,5}, Jack Baty⁶, Jennifer Ivanovich^{4,7}, Sharon Heath^{1,4}, William D. Shannon^{1,4}, Rakesh Nagarajan^{4,5}, Matthew J. Walter^{1,4}, Daniel C. Link^{1,4}, Timothy A. Graubert^{1,4}, John F. DiPersio^{1,4} & Richard K. Wilson^{2,3,4}

Acute myeloid leukaemia is a highly malignant haematopoietic tumour that affects about 13,000 adults in the United States each year. The treatment of this disease has changed little in the past two decades, because most of the genetic events that initiate the disease remain undiscovered. Whole-genome sequencing is now possible at a reasonable cost and timeframe to use this approach for the unbiased discovery of tumour-specific somatic mutations that alter the protein-coding genes. Here we present the results obtained from sequencing a typical acute myeloid leukaemia genome, and its matched normal counterpart obtained from the same patient's skin. We discovered ten genes with acquired mutations; two were previously described mutations that are thought to contribute to tumour progression, and eight were new mutations present in virtually all tumour cells at presentation and relapse, the function of which is not yet known. Our study establishes whole-genome sequencing as an unbiased method for discovering cancer-initiating mutations in previously unidentified genes that may respond to targeted therapies.

NATURE | Vol 456 | 20 November 2008

DNA SEQUENCING

Mammoth genomics

Michael Hofreiter

Reconstruction of most of the genome sequence of the woolly mammoth illustrates how such investigations will pave the way for a deeper understanding of the biology and evolution of extinct species.

En décembre 2011 :

Génomes séquencés et en cours de séquençage :

8847 procaryotes

2987 virus

1518 Eucaryotes

29 groupes

15 groupes

561 animaux
449 champignons
235 protistes
259 plantes



Le projet HapMap :

- comparer les séquences génétiques de différents individus afin de relever les régions chromosomiques où des variations génétiques sont partagées
- découvrir les gènes qui jouent un rôle dans la santé, la maladie et la réponse des individus aux médicaments et aux facteurs environnementaux.
- information sera dans le domaine public,

collaboration entre
Japon, du Royaume-Uni, du Canada, de la Chine, du Nigeria et des États-Unis.

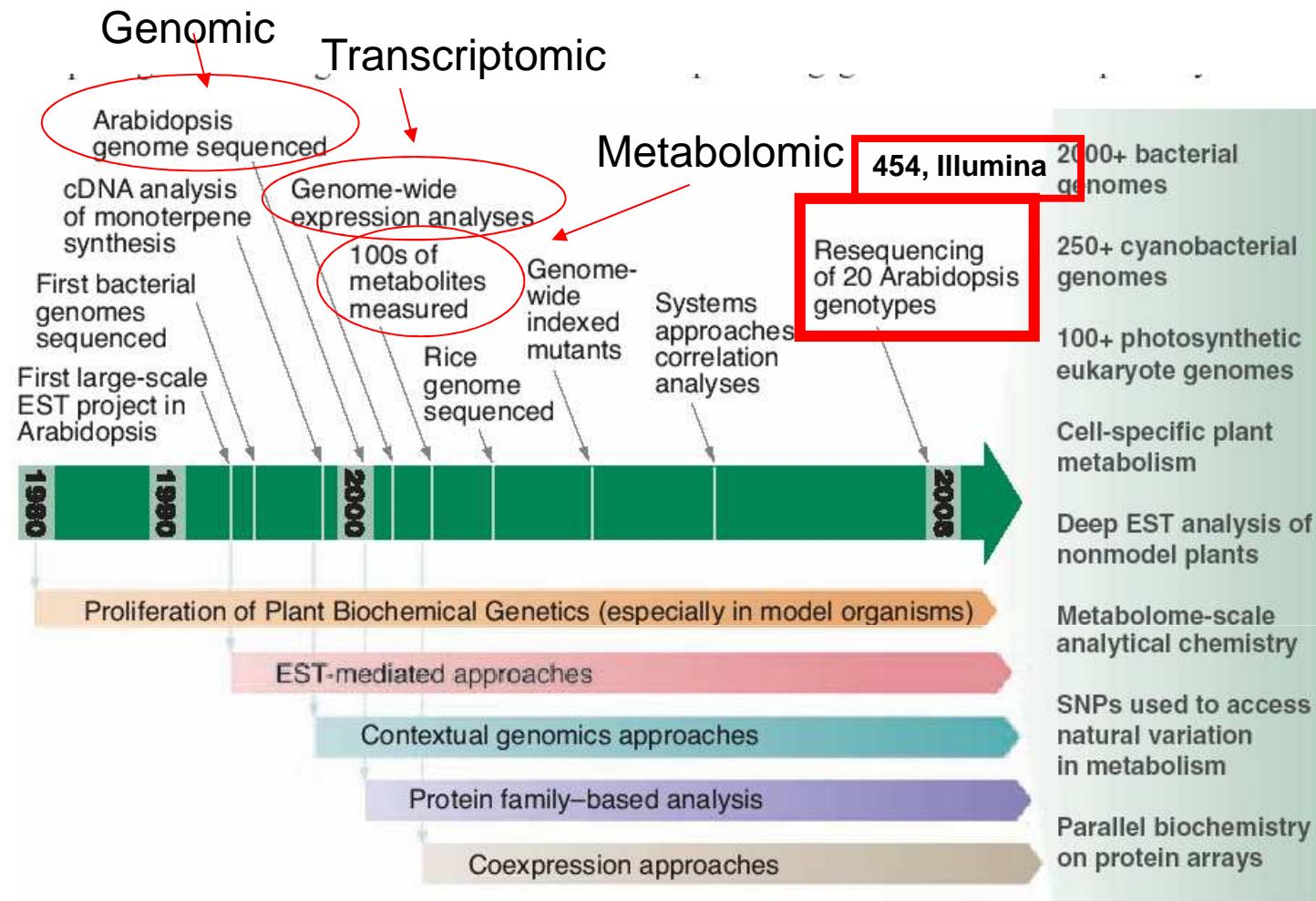
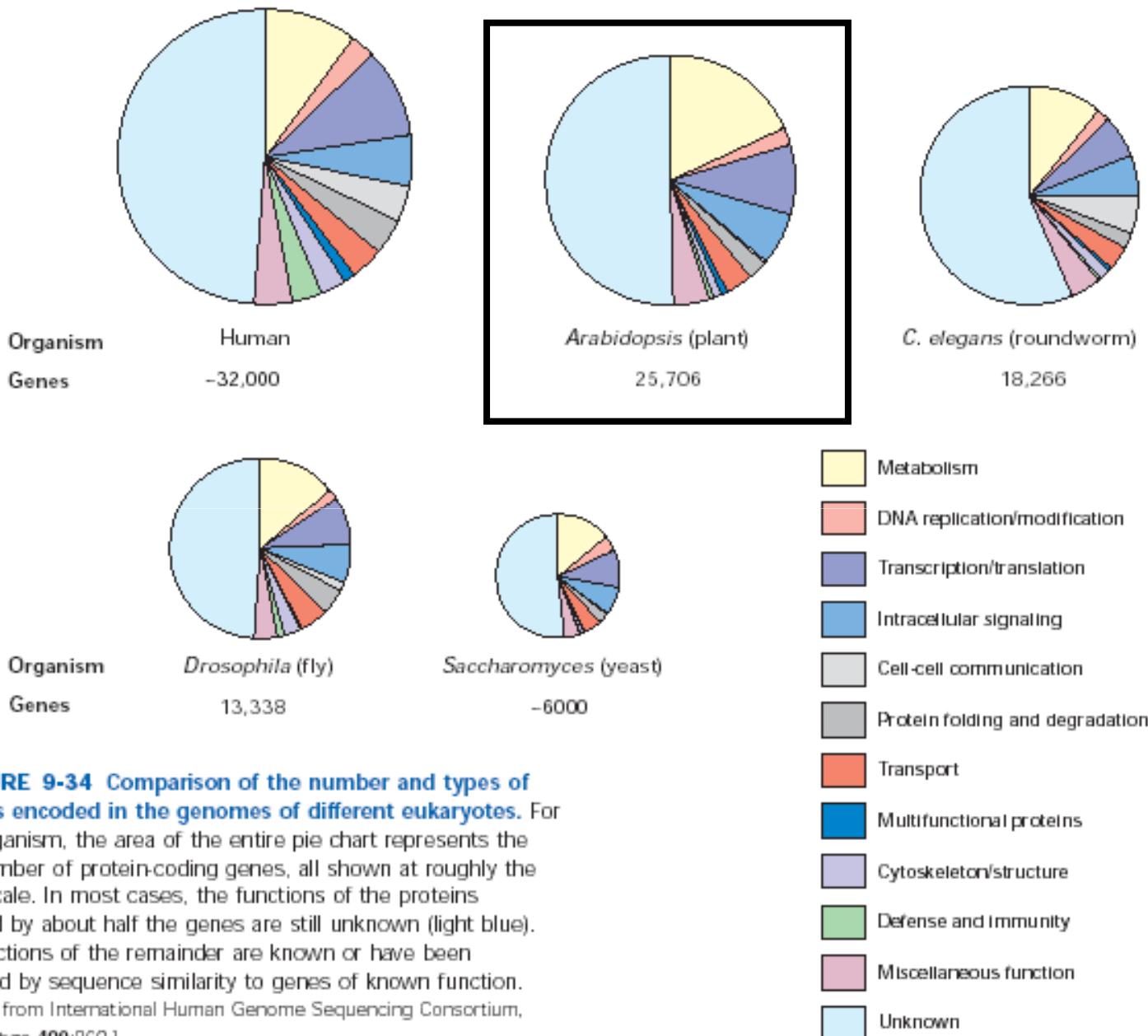


Fig. 1. Time line of genomics-enabled plant biochemistry. Selected major advances to date are indicated above and below the time line. Some approaches and tools likely to further understanding of plant biochemistry during the coming decade are indicated to the right of the time line.

Comparaison des génomes



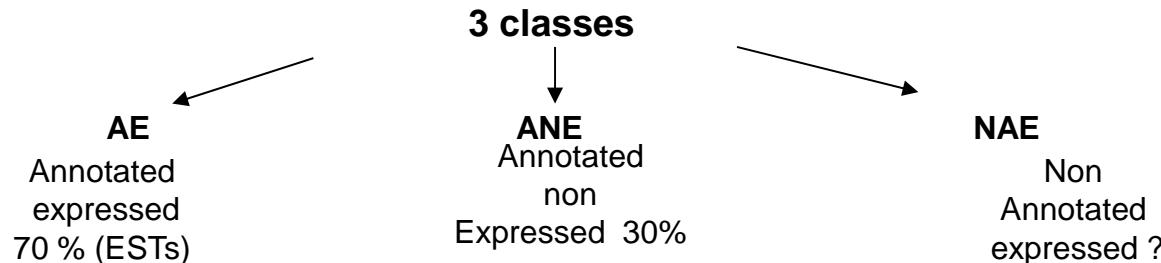
▲ **FIGURE 9-34 Comparison of the number and types of proteins encoded in the genomes of different eukaryotes.** For each organism, the area of the entire pie chart represents the total number of protein-coding genes, all shown at roughly the same scale. In most cases, the functions of the proteins encoded by about half the genes are still unknown (light blue). The functions of the remainder are known or have been predicted by sequence similarity to genes of known function. [Adapted from International Human Genome Sequencing Consortium, 2001, *Nature* 409:860.]

De la connaissance du génome...

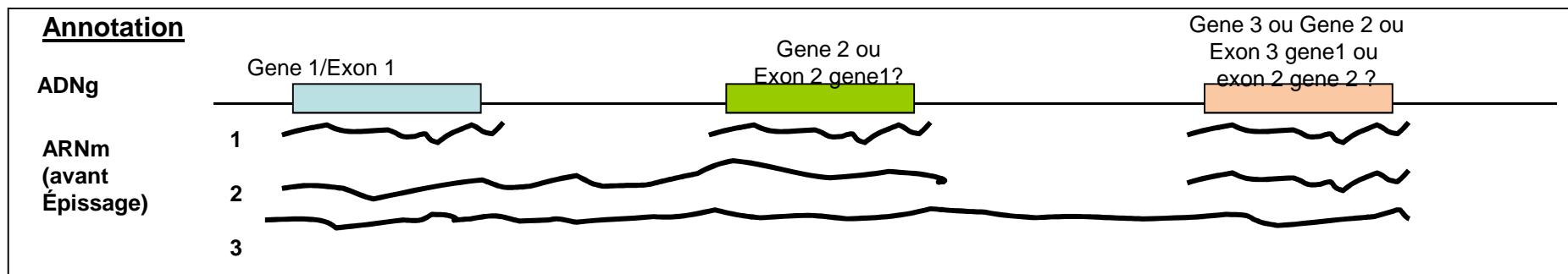
- A son annotation ...
- A son étude de l'expression ...
- Vers l'analyse fonctionnelle...

Yamada et al. Science, 302, 842-846 , 2003. Empirical Analysis of Transcriptional Activity in the *Arabidopsis* Genome

Arabidopsis : 26828 gènes prédis dont 25540 prédis comme gènes codants

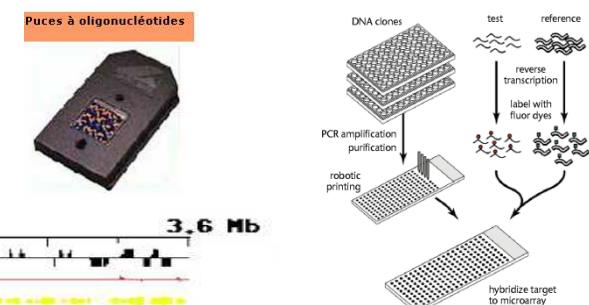
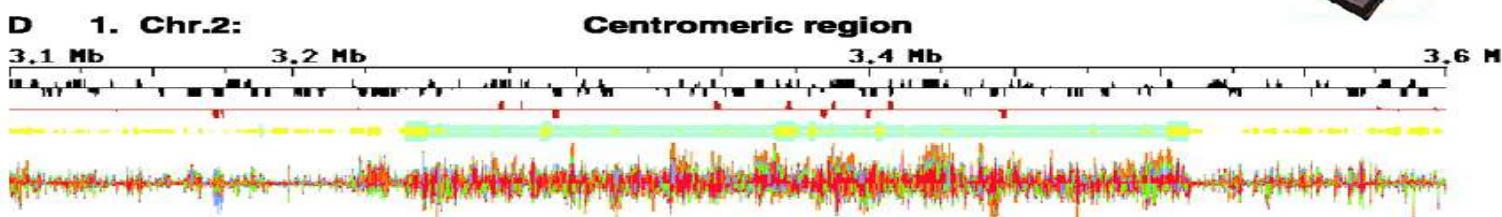


→ intérêt des collections d'EST et de cDNA pleine taille (full lenght)



Solution ? High density oligo arrays qui couvre 94 % génome Arabidopsis

Soit 12 arrays avec 834 000 oligos de 25 mers / arrays
 Hybridation avec 4 populations d'ARNm



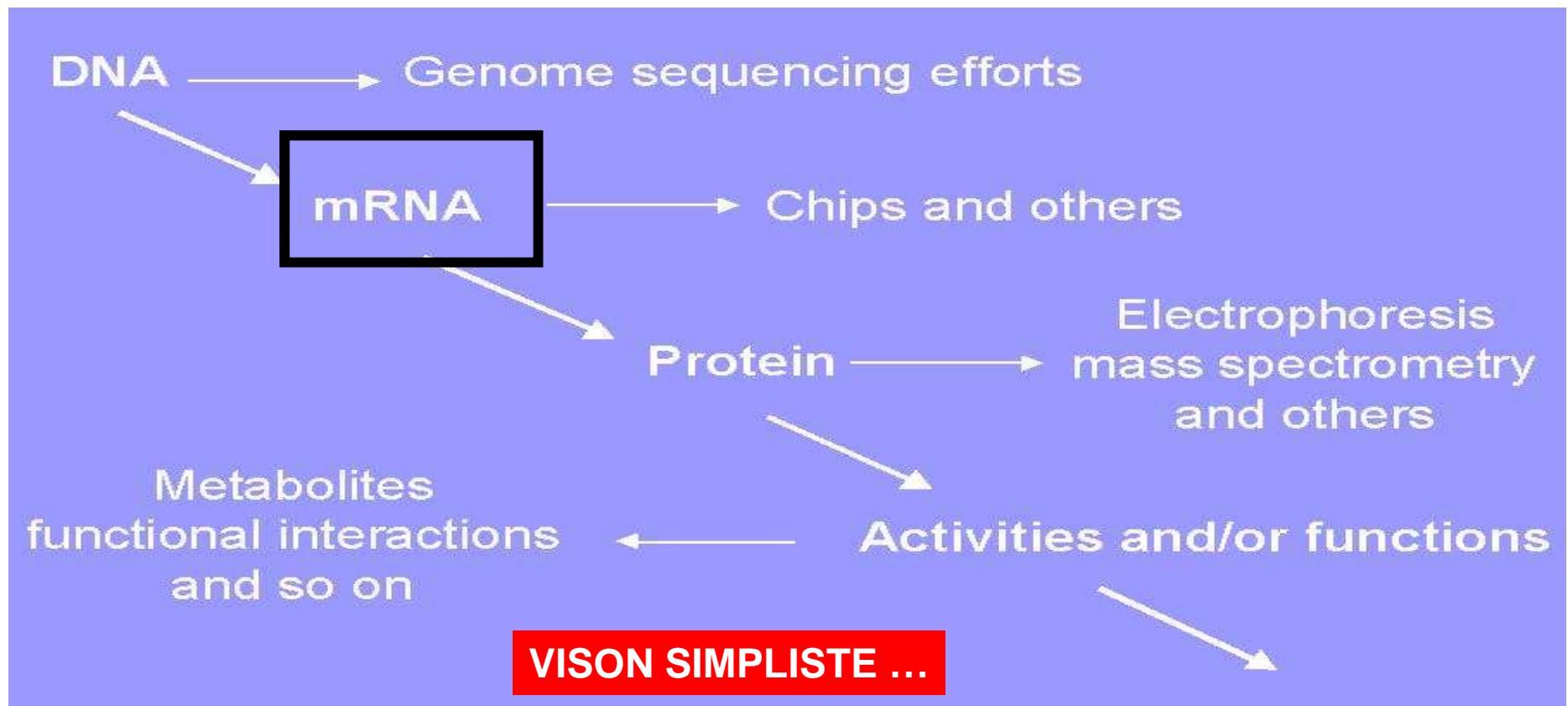
Detection d'activité transcriptionnelle dans région centromérique (40 gènes) et dans régions intergéniques (2000)

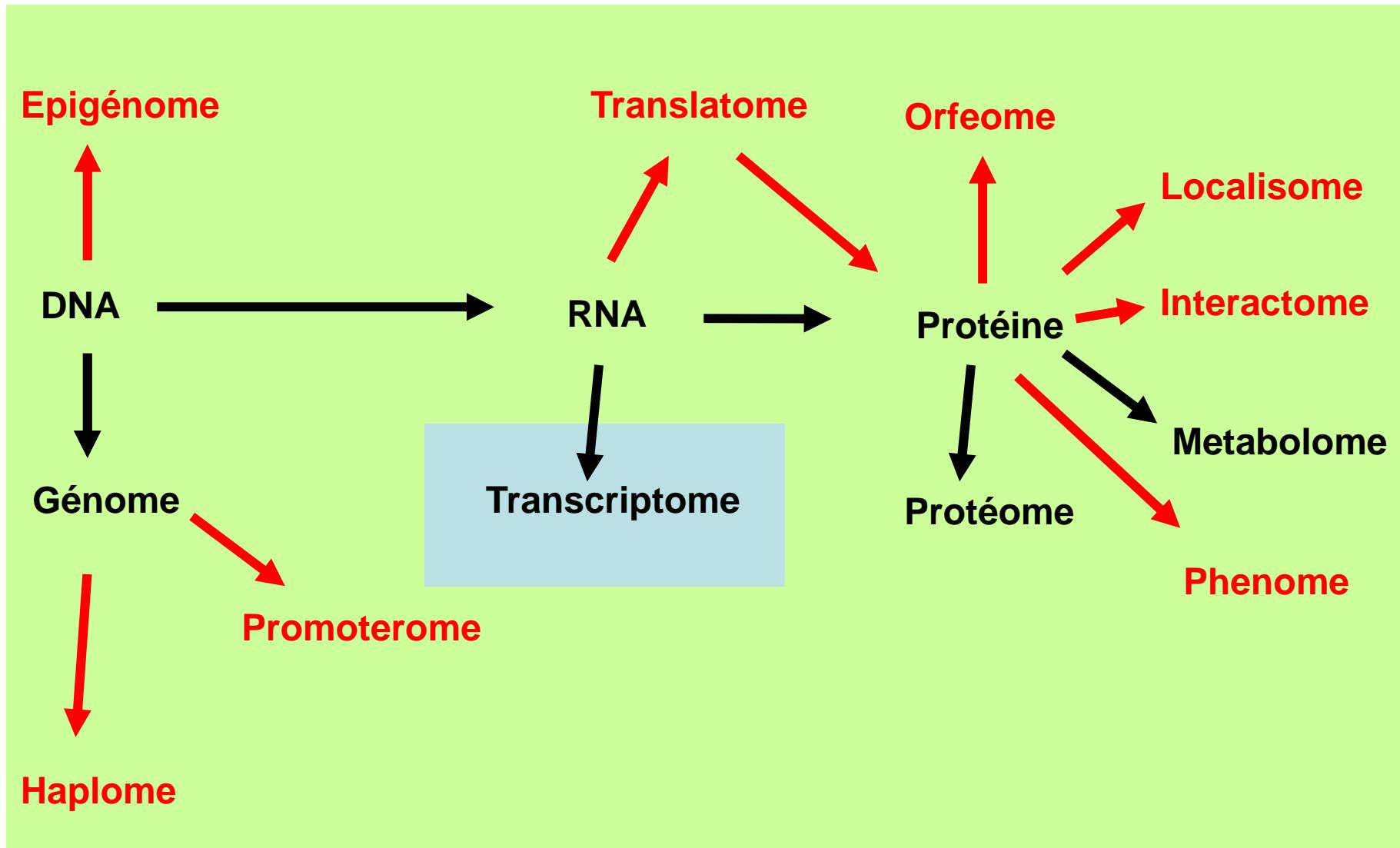
Bilan : 5817 nouvelles unités de transcription soit 30 % de plus que prédit par annotation

De la génomique à la protéomique ...

De la relation : 1 gène } 1 protéine...

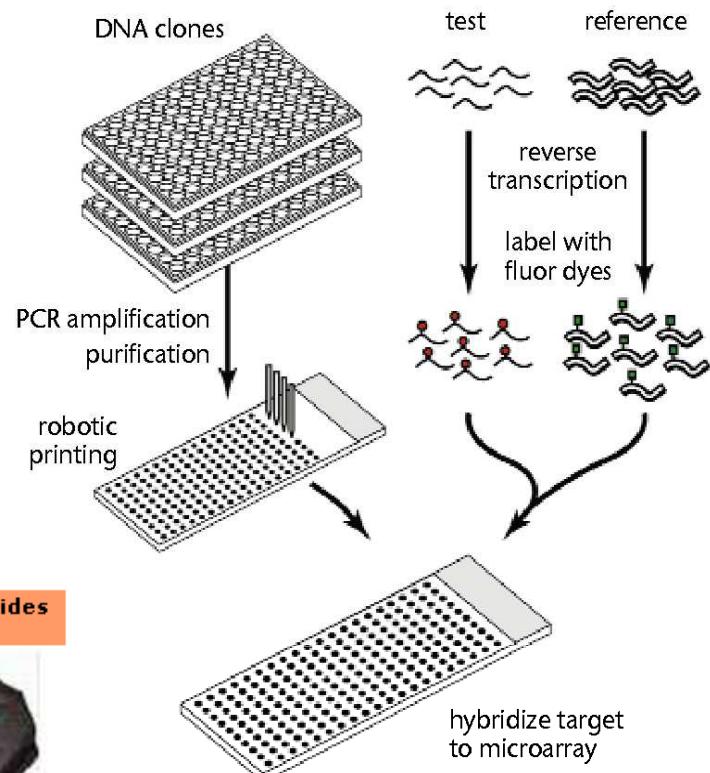
De la relation : 1 protéine } 1 fonction ...



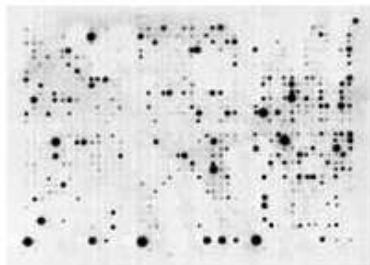


Puces à ADN :

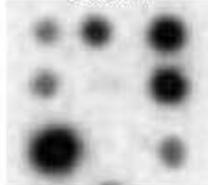
- des cDNAs à la puce ?
- Quels supports ?



Filtres haute densité
(macroarrays)



Détail :



Taille : 12cm x 8cm

- 2400 clones par membrane
- marquage radioactif
- 1 condition expérimentale par membrane

Lames de verre (microarrays)



Détail :



Taille : 5,4cm x 0,9cm

- 10000 clones par lame
- marquage fluorescent
- 2 conditions expérimentales par lame

Puces à oligonucléotides

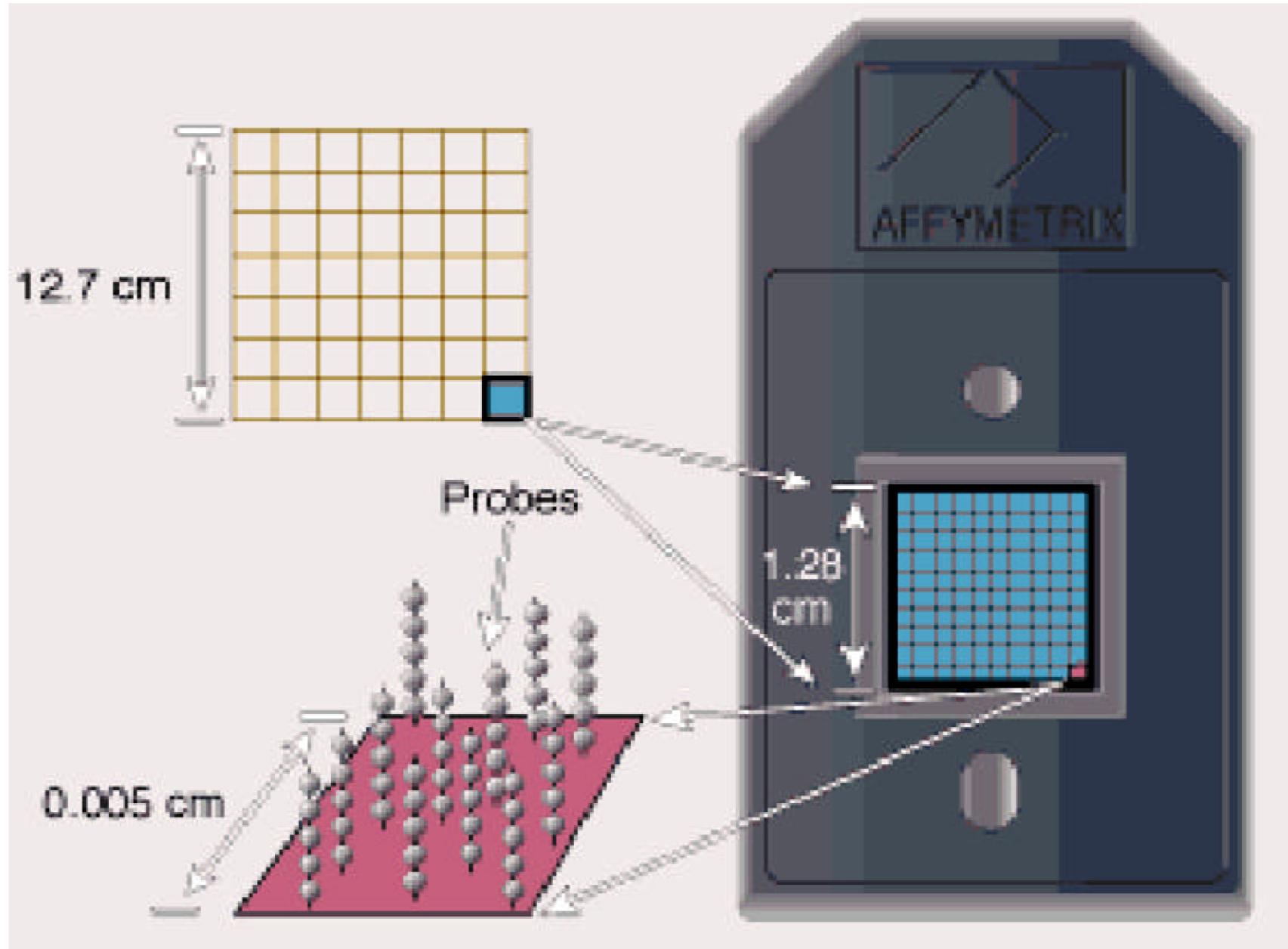


Détail :

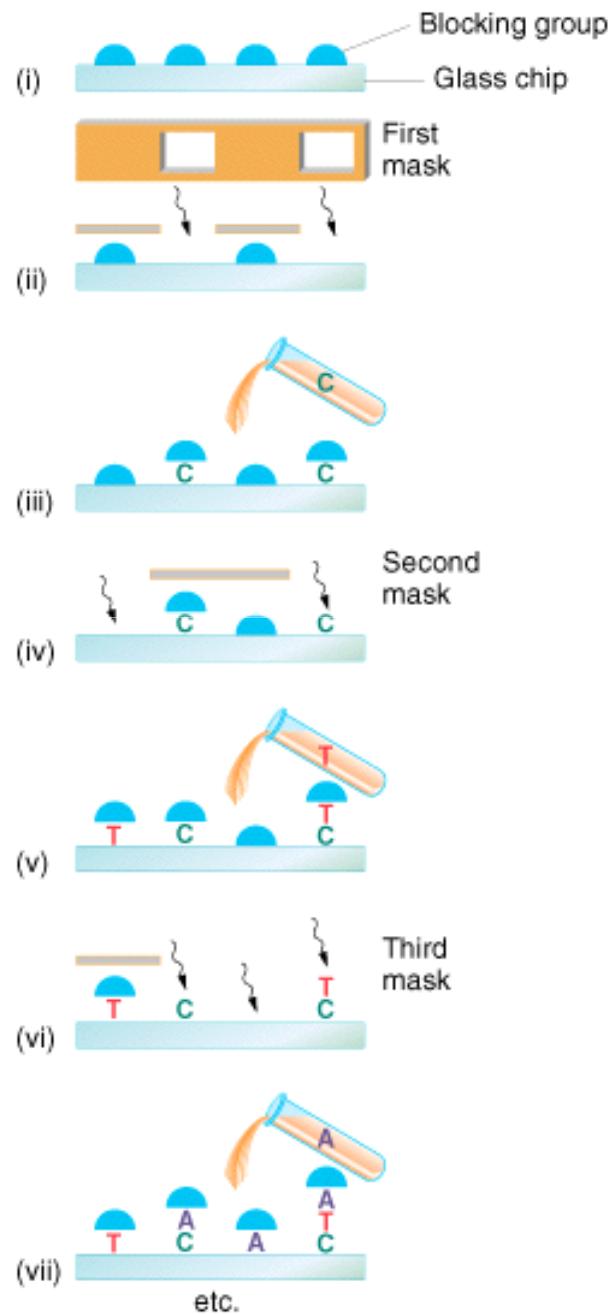


Taille : 1,28cm x 1,28cm

- 300000 oligonucléotides par lame
- marquage fluorescent
- 1 condition expérimentale par lame



Method of oligonucleotide synthesis



Synthesis of an oligonucleotide gene chips - part I

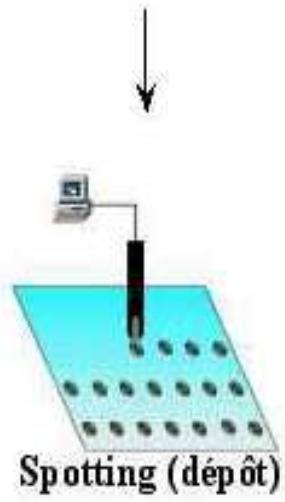
Puces à ADN : Principe

Fabrication des puces à ADN

Lames de verre
recouvertes de polylysine



+
6116 ORFs de levure
amplifiées par PCR



Hybridation

Souche 1 Souche 2

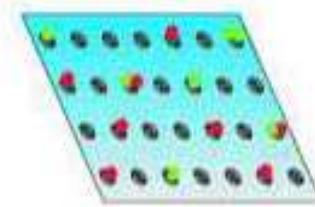


Extraction
des ARN

Cy3

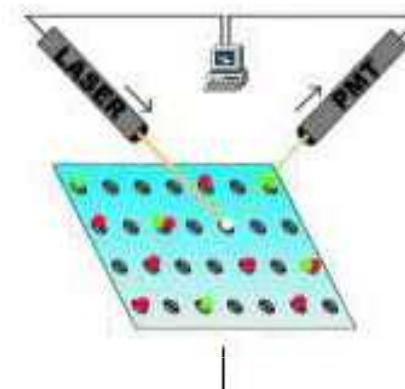
Transcription
des ARNm
en ADNc

Cy5



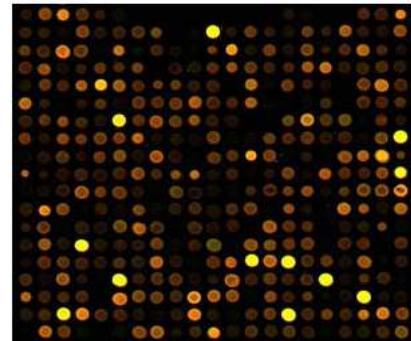
Obtention des résultats

Lecture (scanner)



Analyses des résultats

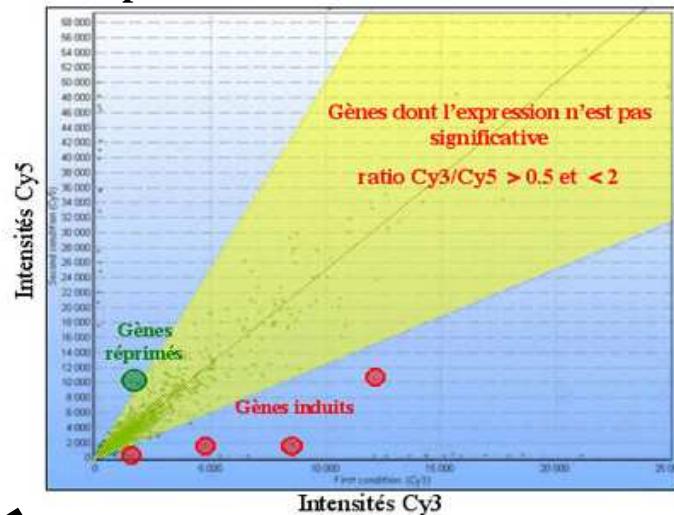
Puces à ADN: du résultat à la signification biologique ?



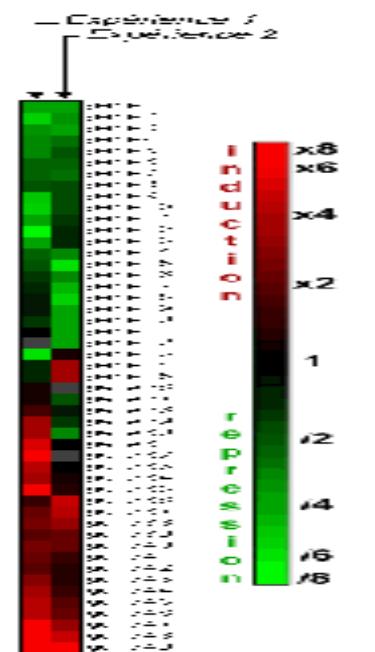
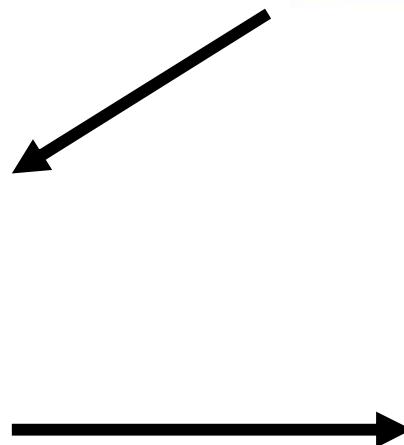
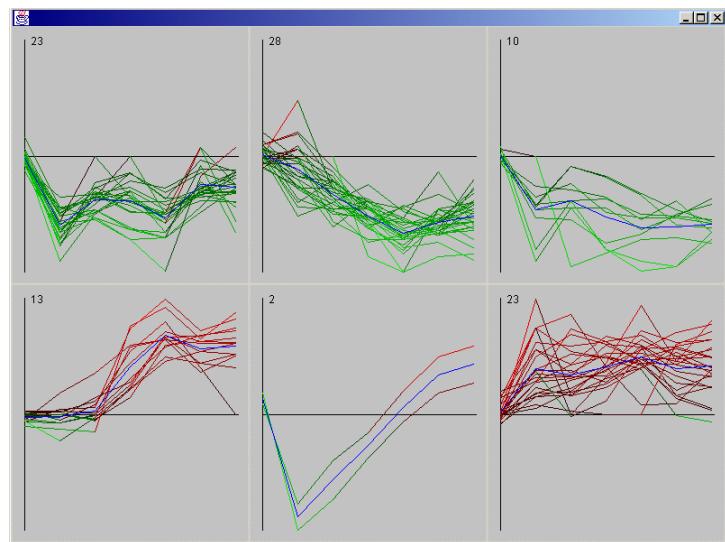
Rouge : Spé A
Vert : Spé B
Jaune : Commun

Sonde A →
Sonde B

Quels sont les gènes dont l'expression est réprimée ? Induite ? Invariable ?



Quels sont les gènes dont le profil d'expression est identique ?



Le clustering d'expression

MPSS : ? Massively Parallel Signature Sequencing

UHTS method

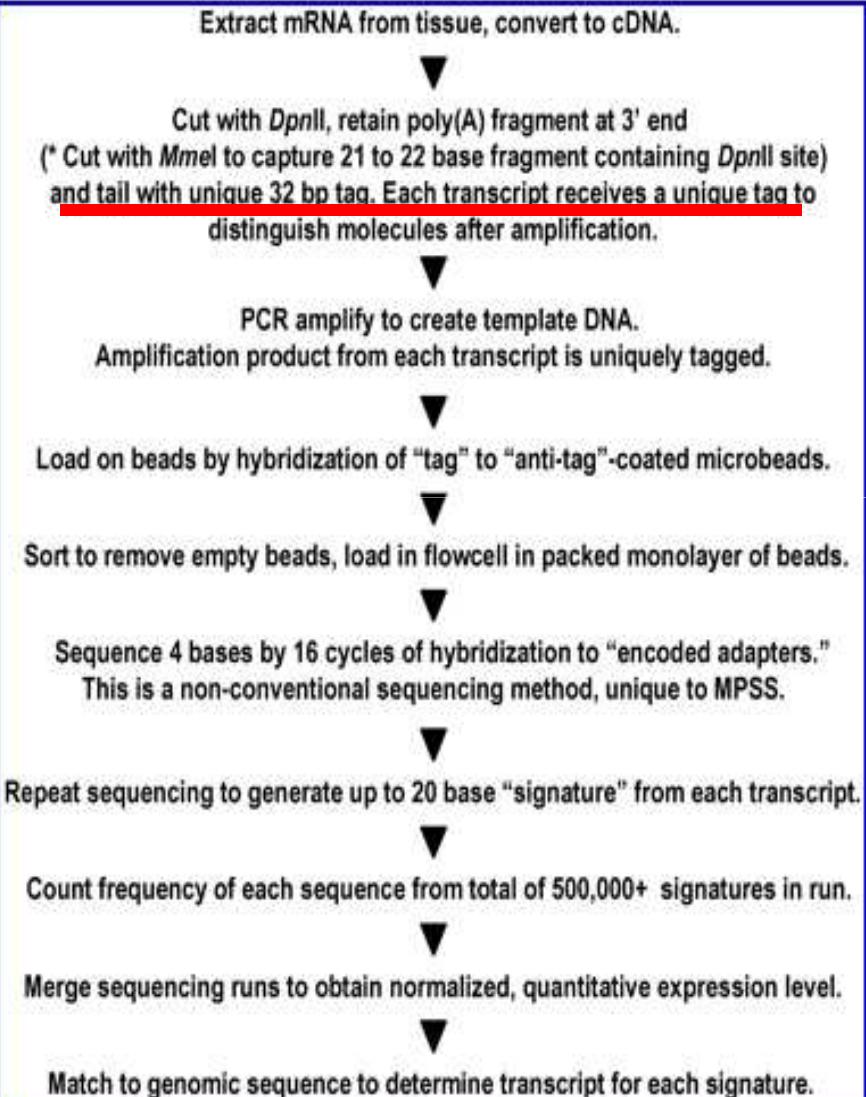
Genome Research 14:1641-1653, 2004
ISSN 1088-9051 / \$5.00

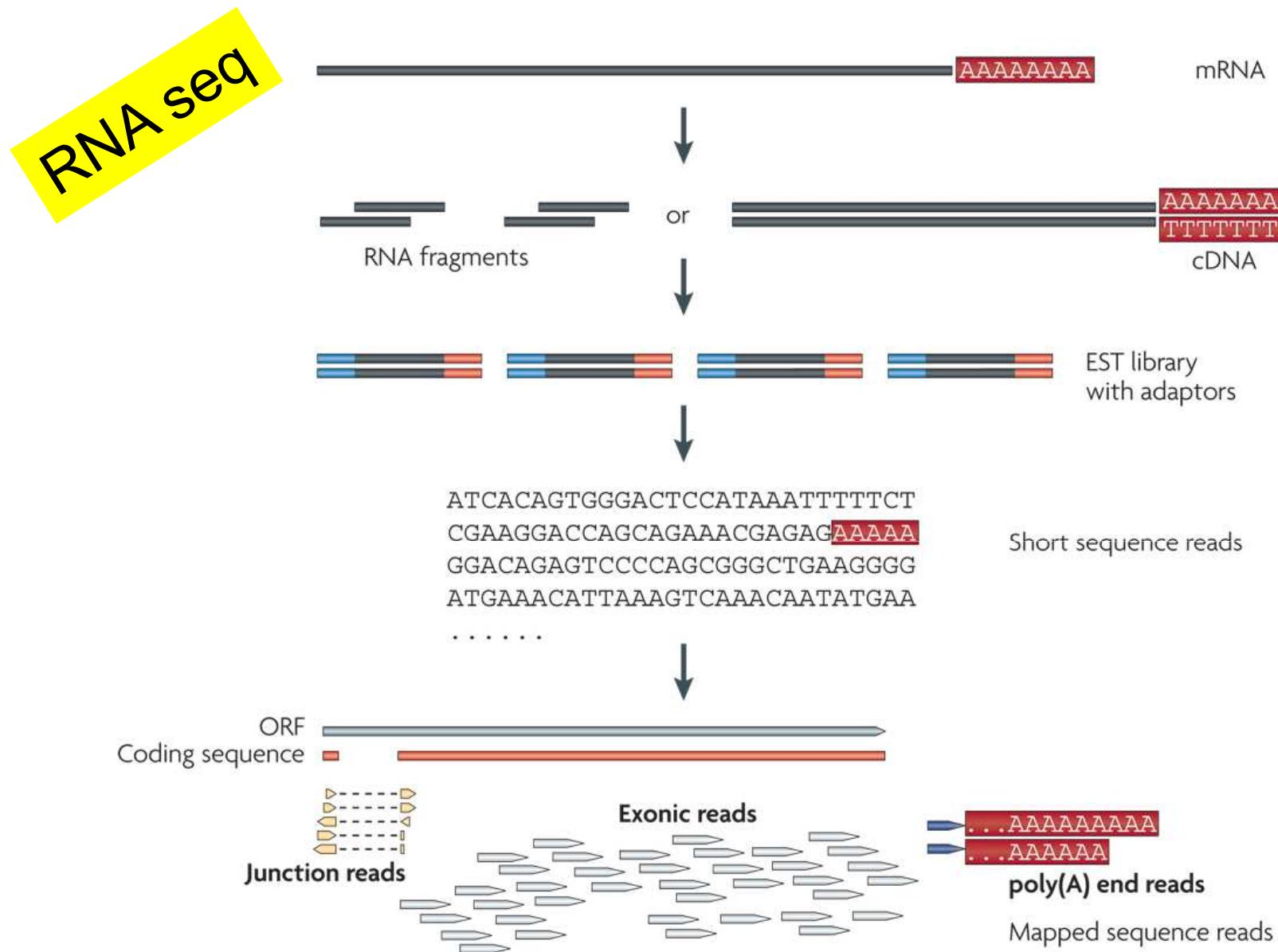
Methods

The Use of MPSS for Whole-Genome Transcriptional Analysis in *Arabidopsis*

Blake C. Meyers^{1,4}, Shivakundan Singh Tej¹, Tam H. Vu¹, Christian D. Haudenschi
Vikas Agrawal¹, Steve B. Edberg², Hassan Ghazal¹ and Shannon Decola³

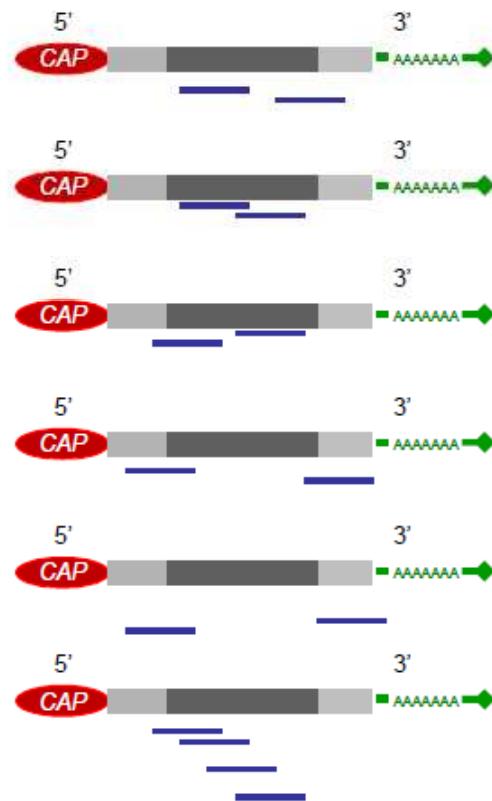
¹ Department of Plant and Soil Sciences, and Delaware Biotechnology Institute, University of Delaware, Newark, Delaware 19714, USA; ² Department of Vegetable Crops, University of California, Davis, California 95616, USA; ³ Lynx Therapeutics, Inc., Hayward, California 94545, USA



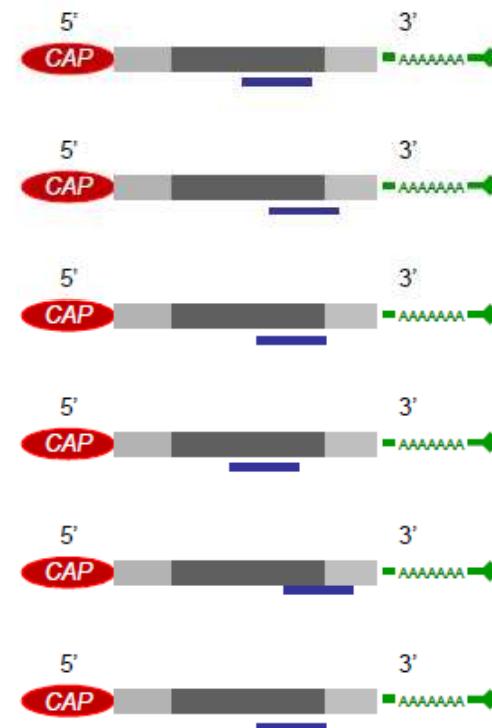


RNA seq

RNA seq



DGE

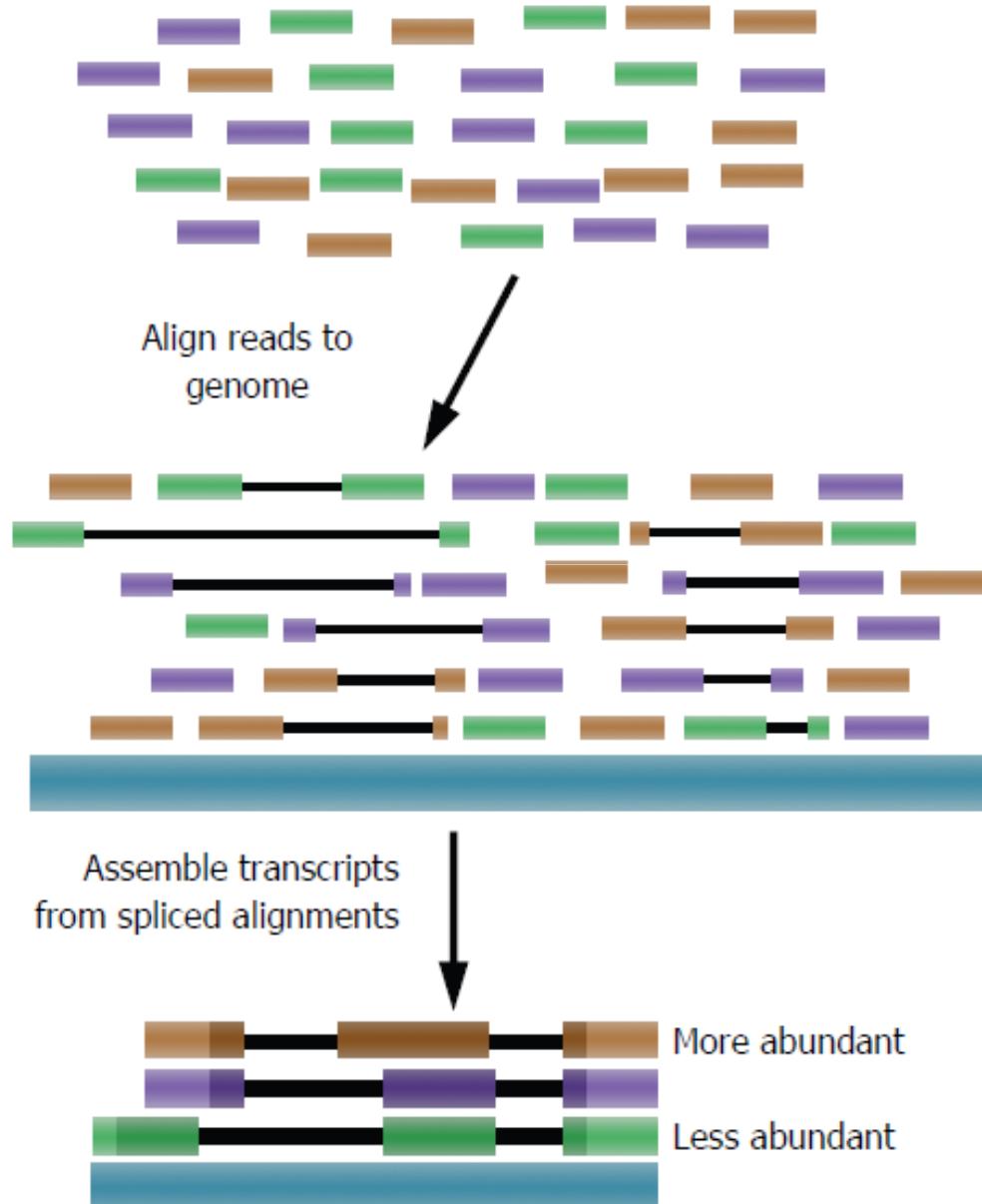


Resequencing Strategy Summary



RNA seq

RNA-Seq Reads



ChIP on chip

Chromatin ImmunoPrecipitation on Chip

ChIP and sequencing high throughput

utilisées pour repérer des sites de fixation de facteurs de transcription
localiser ces sites et d'étudier les séquences d'ADN correspondantes

- 1.Liaison covalente in vivo des protéines à l'ADN
- 2.Extraction de l'ADN
- 3.Découpage de l'ADN par sonication
- 4.Sélection des fragments grâce à un anticorps
- 5.Précipitation des complexes ADN-protéine-anticorps,
6. Séparation du complexe ADN-protéine pour garder ADN (proteinase K)
7. On obtient une collection de fragments d'ADN qui interagissent avec une protéine d'intérêt .

combinaison de la technique de Chromatin Immunoprécipitation avec la méthode des puces à ADN.

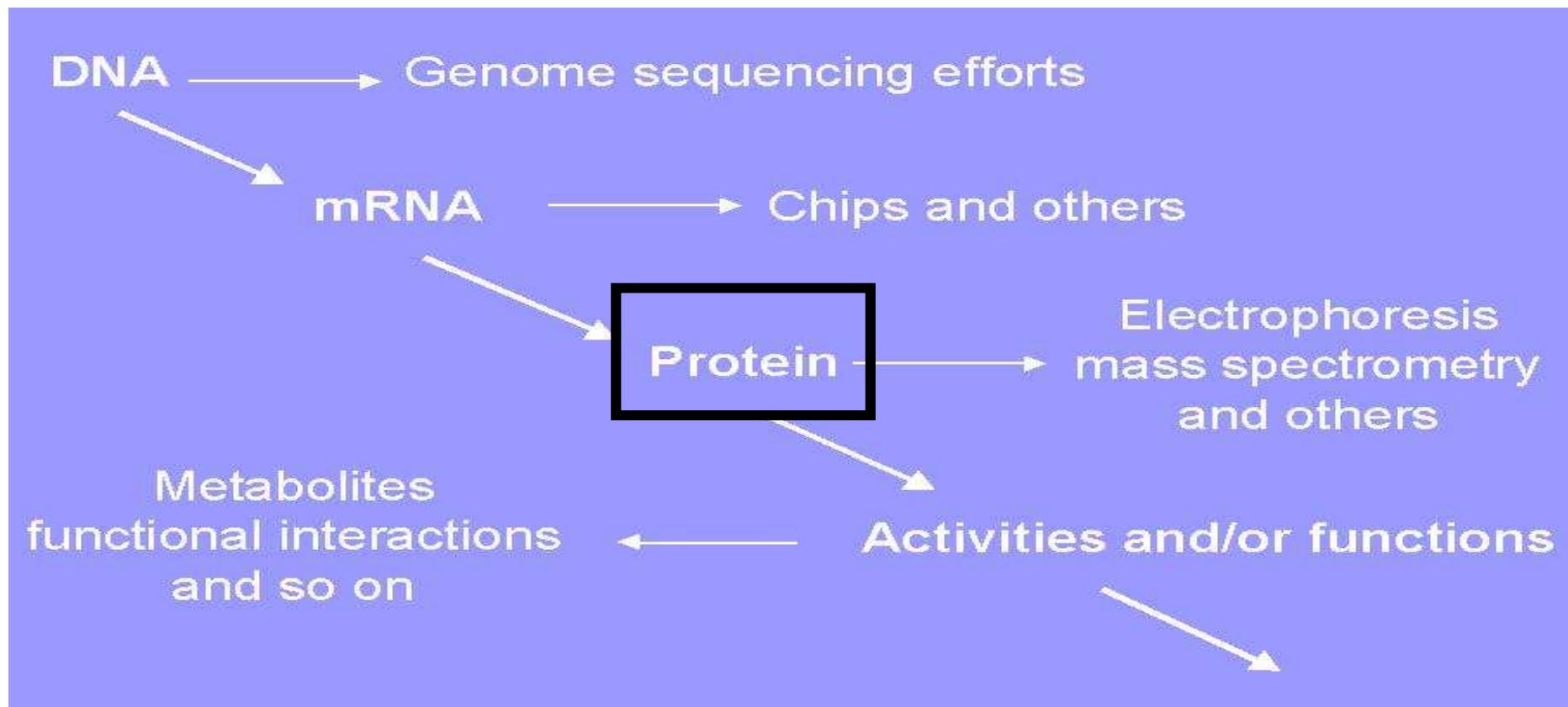
8.On utilise une puce à ADN pour identifier les fragments

8.On séquence (illumina)... pour identifier les fragments

De la génomique à la protéomique ...

De la relation : 1 gène } 1 protéine...

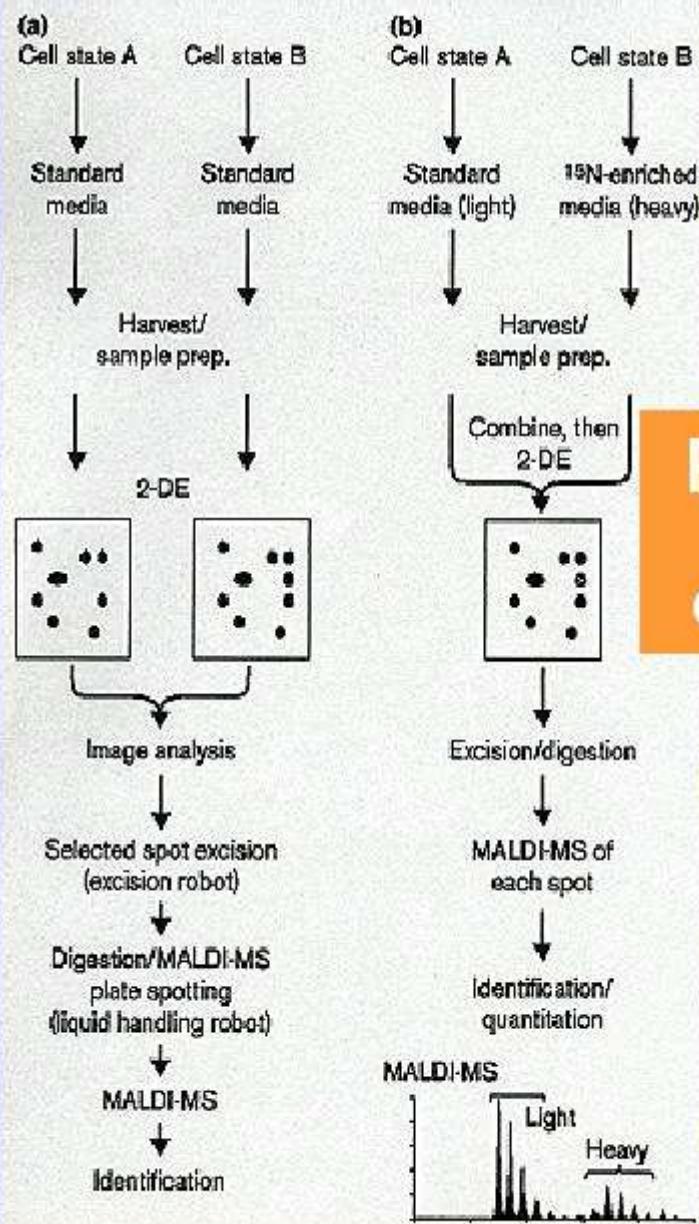
De la relation : 1 protéine } 1 fonction ...



Defining Proteomics

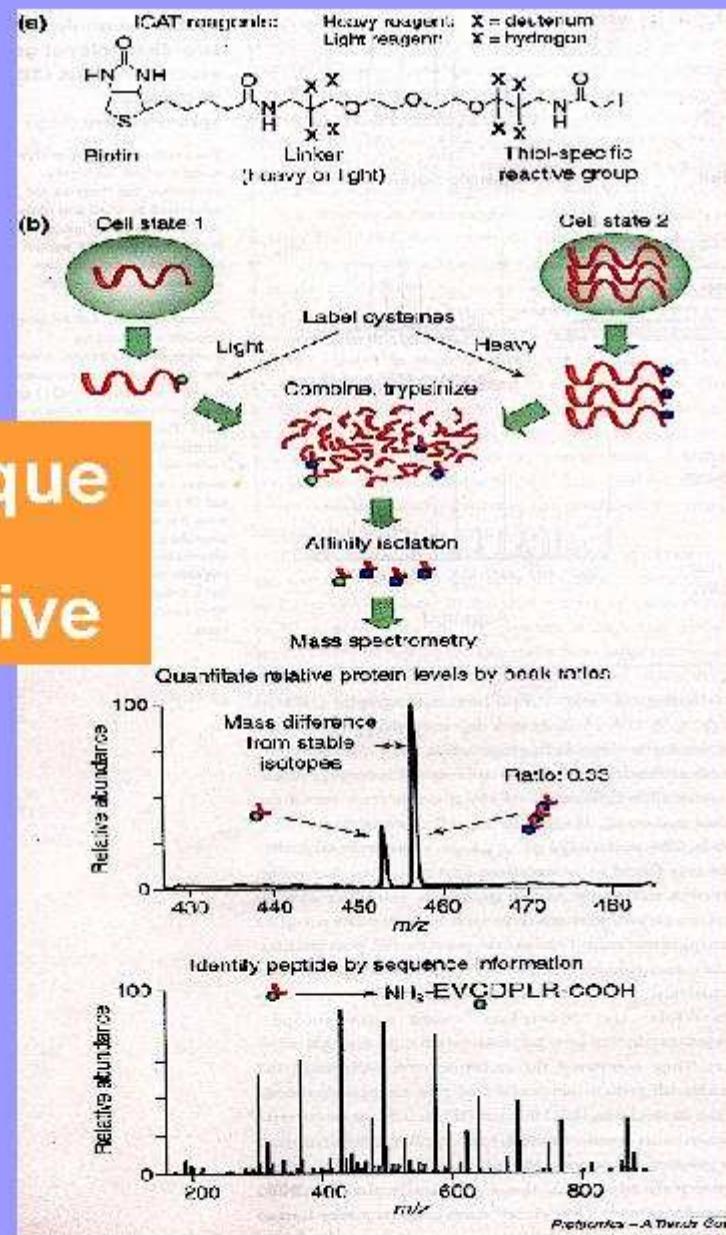
- Proteomics - in the post-genomic era
 - ◆ Differential display of proteins on 2-dimensional gels
 - ♦ Application to mutants, knockouts, different stimuli, etc.
 - ◆ Protein identification
 - ♦ Affinity purification (analysis on 1-D or 2-D gels)
 - ♦ Proteins chips (coated with proteins or antibodies)
 - ◆ Protein modification
 - ♦ Phosphorylation, glycosylation,
 - ◆ Protein-protein interaction
 - ♦ Yeast two-hybrid, phage display
 - ♦ Affinity purification of protein complexes
 - ◆ Protein identification by mass spectrometry

Méthode "standard"



Marquage différentiel

Isotope-Coded Affinity Tag (ICAT)



Protéomique comparative

NEWS

NATURE | Vol 452 | 24 April 2008

Biologists initiate plan to map human proteome

Ambitious plans to catalogue and characterize all proteins in the human body — a Human Proteome Project — are being drawn up by a small group of researchers. But with a price tag of around US\$1 billion, some question whether the organizers can raise enough money or momentum for such an undertaking.

Researchers looked into the idea in the mid-1990s as the Human Genome Project was taking shape — the human proteome seemed a natural successor. However, a coordinated effort to index human proteins never emerged. One reason is that the scale and complexity of the problem proved daunting and nebulous. Protein-coding genes in the body can make tens of different versions of a protein, and each of these can be modified by the addition of chemical groups in countless different ways. All these proteins are being manufactured at differing levels, and at different moments in time, in the 200 or so types of human cell. "It was thought to be beyond comprehension," says John Bergeron of McGill University in Montreal, Canada, former president of the Human Proteome Organisation (HUPO).

Now Bergeron and a group of leading proteomics researchers are putting together a pro-

because estimates of the number of protein-coding genes have shrunk. It was once thought that there might be around 50,000 or 100,000, but now, just 21,000 or so are thought to exist, making the scale of human proteomics more manageable. And the group plans to focus on only a single protein produced from each gene, rather than its many forms. "We got rid of all this complexity," Bergeron says. "We tried to craft a project that would be doable with easy-to-track milestones."

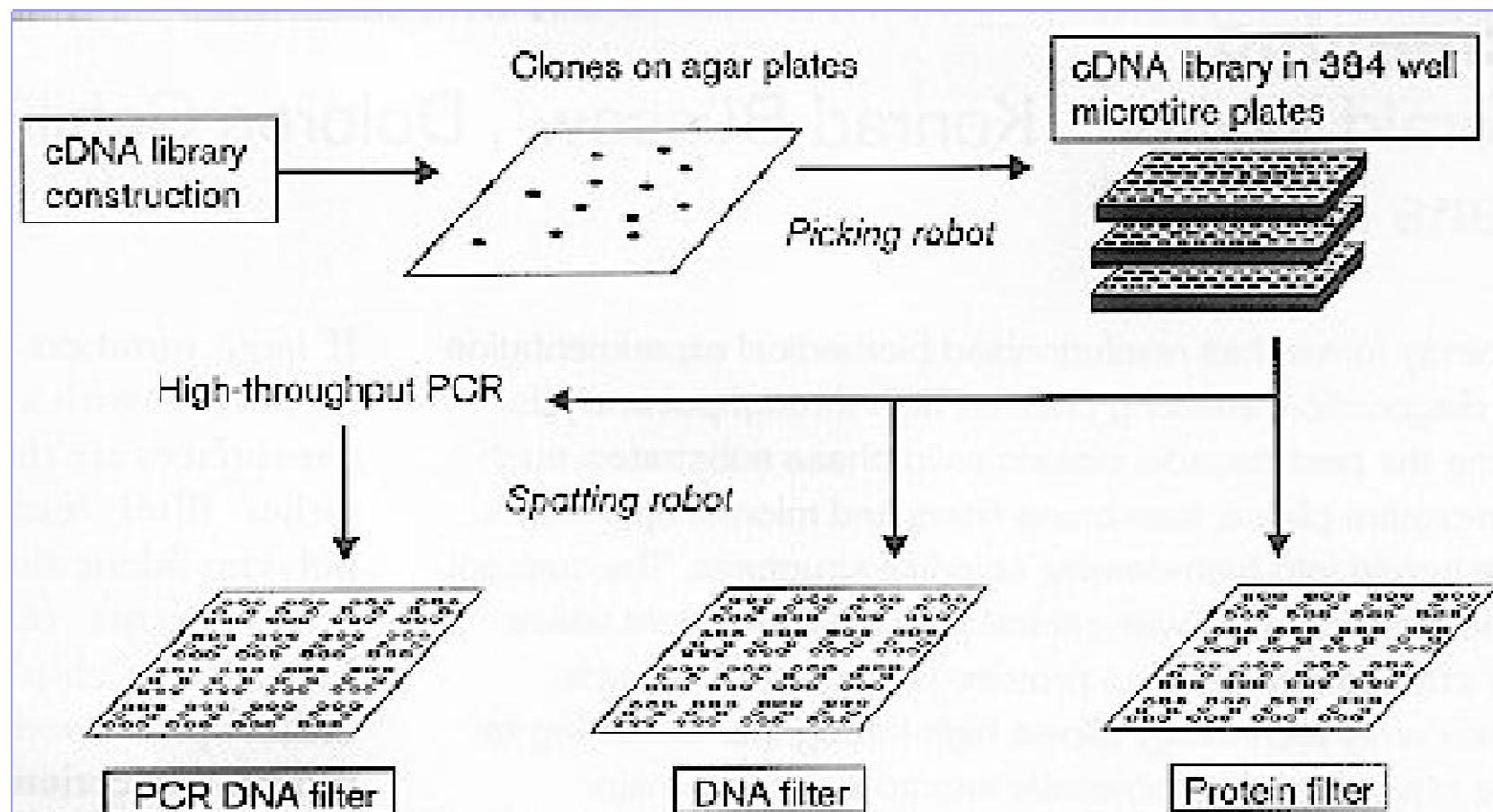
The plan is to tackle this with three different experimental approaches. One would use mass spectrometry to identify proteins and their quantities in tissue samples; another would generate antibodies to each protein and use these to show its location in tissues and cells; and the third would systematically identify, for each protein, which others it interacts with in protein complexes. The project would also involve a massive bioinformatics effort to ensure that the data could be pooled and accessed, and the production of shared reagents.

Bergeron envisages the work being divvied up between labs around the world. He says that the first stage of the project — which involves amassing existing mass spectrometry

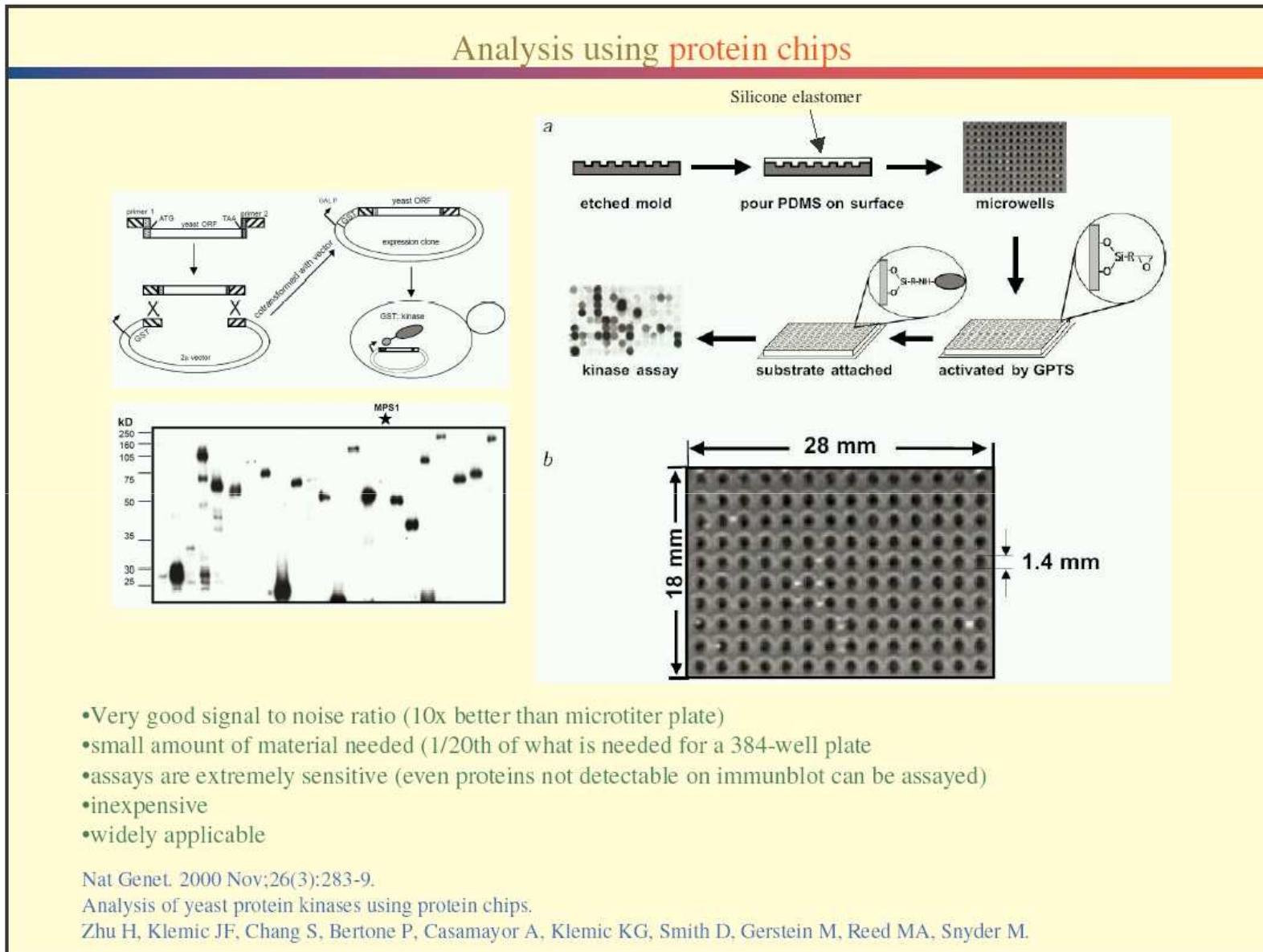


Des approches à haut débit pour la recherche de partenaires protéiques ?

Protein arrays



Ex: utilisation de protein chips pour la recherche de substrat d'une kinase



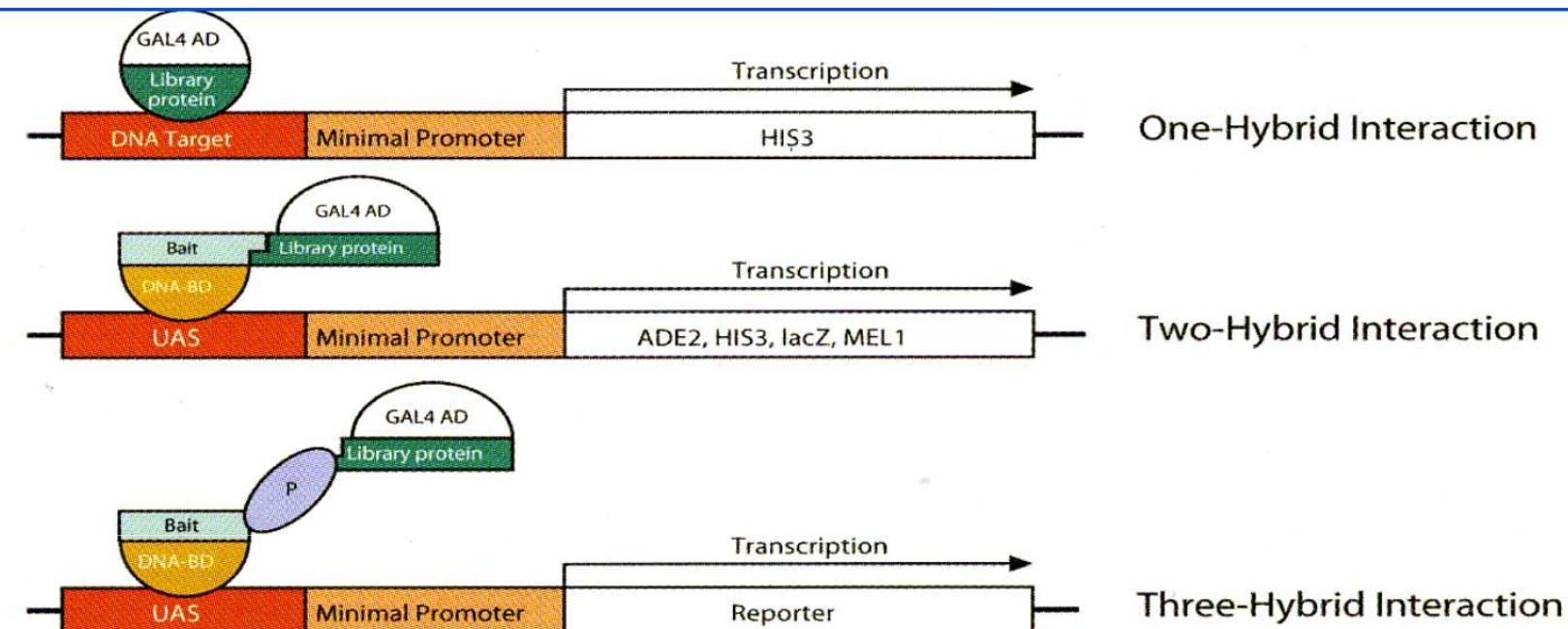
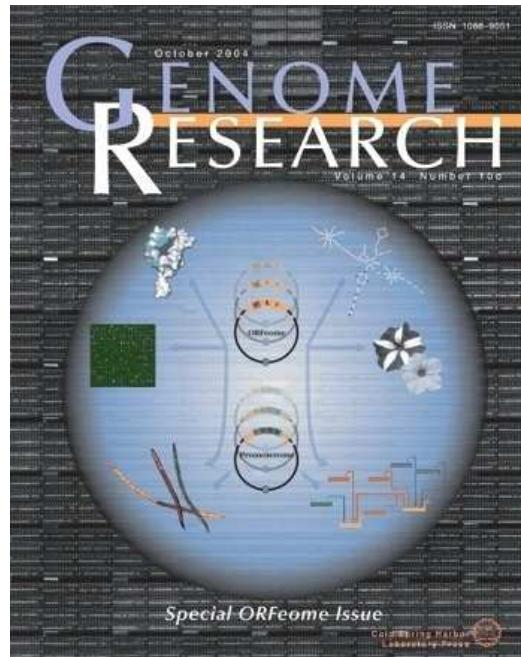


Figure 1. Yeast hybrid technology. Matchmaker systems use a sensitive transcriptional assay to detect one-hybrid, two-hybrid, and three-hybrid interactions. A yeast reporter strain transformed with both a bait and library plasmid will express the plasmids' gene inserts as fusions to either the GAL4 DNA-BD (DNA-binding domain) or AD (transcription-activating domain), depending on the plasmid (see below). If a library protein interacts with a bait protein (two-hybrid) or DNA target (one-hybrid) the host strain actively expresses the reporter gene located downstream of the promoter. Three-hybrid interactions, or protein interactions that occur via a third protein (P), can be detected using our pBridge Vector (Cat. No. 630404). UAS = GAL4-responsive upstream activating sequence.



Methods

High-Throughput Expression of *C. elegans* Proteins

Chi-Hao Luan,^{1,3} Shihong Qiu,¹ James B. Finley,¹ Mike Carson,¹ Rita J. Gray,¹ Wenyng Huang,¹ David Johnson,¹ Jun Tsao,¹ Jérôme Reboul,² Philippe Vaglio,² David E. Hill,² Marc Vidal,² Lawrence J. DeLucas,¹ and Ming Luo^{1,3}

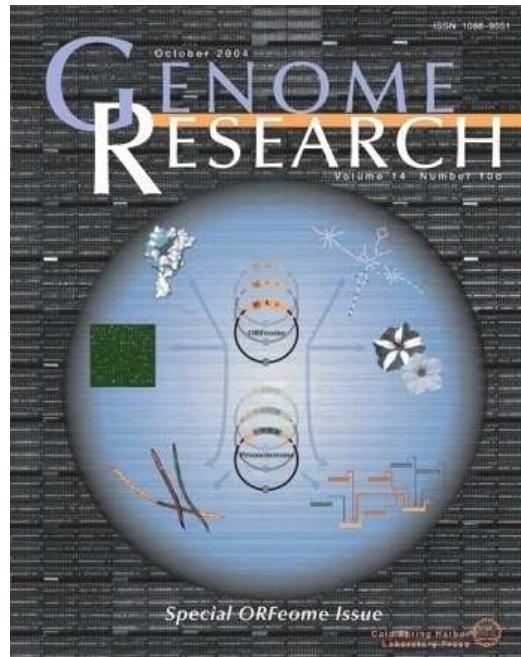
¹Center for Biophysical Sciences and Engineering, Southeast Collaboratory for Structural Genomics, University of Alabama at Birmingham, Birmingham, Alabama 35294, USA; ²Center for Cancer Systems Biology and Department of Cancer Biology, Dana-Farber Cancer Institute, and Department of Genetics, Harvard Medical School, Boston, Massachusetts 02115, USA

Proteome-scale studies of protein three-dimensional structures should provide valuable information for both investigating basic biology and developing therapeutics. Critical for these endeavors is the expression of recombinant proteins. We selected *Caenorhabditis elegans* as our model organism in a structural proteomics initiative because of the high quality of its genome sequence and the availability of its ORFeome, protein-encoding open reading frames (ORFs). In a flexible recombinational cloning format. We developed a robotic pipeline for recombinant protein expression, applying the Gateway cloning/expression technology and utilizing a stepwise automation strategy on an integrated robotic platform. Using the pipeline, we have carried out heterologous protein expression experiments on 10,167 ORFs of *C. elegans*. With one expression vector and one *Escherichia coli* strain, protein expression was observed for 4854 ORFs, and 1536 were soluble. Bioinformatics analysis of the data indicates that protein hydrophobicity is a key determining factor for an ORF to yield a soluble expression product. This protein expression effort has investigated the largest number of genes in any organism to date. The pipeline described here is applicable to high-throughput expression of recombinant proteins for other species, both prokaryotic and eukaryotic, provided that ORFeome resources become available.

2102 **Genome Research**
www.genome.org

14:2102–2110 ©2004 by Cold Spring Harbor Laboratory Press ISSN 1088-9051/04; www.genome.org

→ Pour Interactome



Resource

Generation of the *Brucella melitensis* ORFeome Version 1.1

Amélie Dricot,¹ Jean-François Rual,^{1,2} Philippe Lamesch,^{1,2} Nicolas Bertin,² Denis Dupuy,² Tong Hao,² Christophe Lambert,¹ Régis Hallez,¹ Jean-Marc Delroisse,¹ Jean Vandenhoute,¹ Ignacio Lopez-Goñi,³ Ignacio Moriyon,³ Juan M. Garcia-Lobo,⁴ Félix J. Sangari,⁴ Alastair P. MacMillan,⁵ Sally J. Cutler,⁵ Adrian M. Whatmore,⁵ Stephanie Bozak,⁶ Reynaldo Sequerra,⁶ Lynn Doucette-Stamm,⁶ Marc Vidal,² David E. Hill,² Jean-Jacques Letesson,¹ and Xavier De Bolle^{1,7}

¹Research Unit in Molecular Biology (URBM), University of Namur, 5000 Namur, Belgium; ²Center for Cancer Systems Biology and Department of Cancer Biology, Dana-Farber Cancer Institute and Department of Genetics, Harvard Medical School, Boston, Massachusetts 02115, USA; ³Departamento de Microbiología, Universidad de Navarra, Pamplona 31008, Spain; ⁴Departamento de Biología Molecular, Universidad de Cantabria, Santander 39011, Spain; ⁵Department of Statutory and Exotic Bacterial Diseases, FAO/WHO Collaborating Centre for Reference and Research on Brucellosis, Veterinary Laboratories Agency, Weybridge, Surrey KT15 3NB, United Kingdom; ⁶Agencourt Biosciences Corporation, Beverly, Massachusetts 01915, USA

→ Pour Interactome

Gateway & Plant biology

Technologie gateway & RNAi



ACADEMIC
PRESS

Available online at www.sciencedirect.com



Methods 30 (2003) 289–295

METHODS

www.elsevier.com/locate/ymeth

Constructs and methods for high-throughput gene silencing in plants

Chris Helliwell and Peter Waterhouse*

CSIRO Plant Industry, GPO Box 1600, Canberra ACT 2601, Australia

Accepted 7 February 2003

Abstract

Gene silencing can be achieved by transformation of plants with constructs that express self-complementary (termed hairpin) RNA containing sequences homologous to the target genes. The DNA sequences encoding the self-complementary regions of hairpin (hp) RNA constructs form an inverted repeat. The inverted repeat can be stabilized in bacteria through separation of the self-complementary regions by a “spacer” region. When the spacer sequence encodes an intron, the efficiency of gene silencing is very high. There are at least three ways in which hpRNA constructs can be made. The construct may be generated from standard binary plant transformation vectors in which the hairpin-encoding region is generated *de novo* for each gene. Alternatively, generic gene-silencing vectors such as the pHANNIBAL and the pHELLSGATE series can be used. They simply require the insertion of PCR products, derived from the target gene, into the vectors by conventional cloning or by using the Gateway directed recombination system. In this article, we describe and evaluate the advantages of these vectors and then provide the protocols for their efficient use.

© 2003 Published by Elsevier Science (USA).

Keywords: Gene silencing; Recombination cloning; Plant functional genomics; pHELLSGATE

DNA



mRNA



Protein



Activities and/or functions

METABOLOMICS IN SYSTEMS BIOLOGY

Wolfram Weckwerth

Max-Planck-Institut für Molekulare Pflanzenphysiologie, 14424 Potsdam, Germany;
email: weckwerth@mpimp-golm.mpg.de

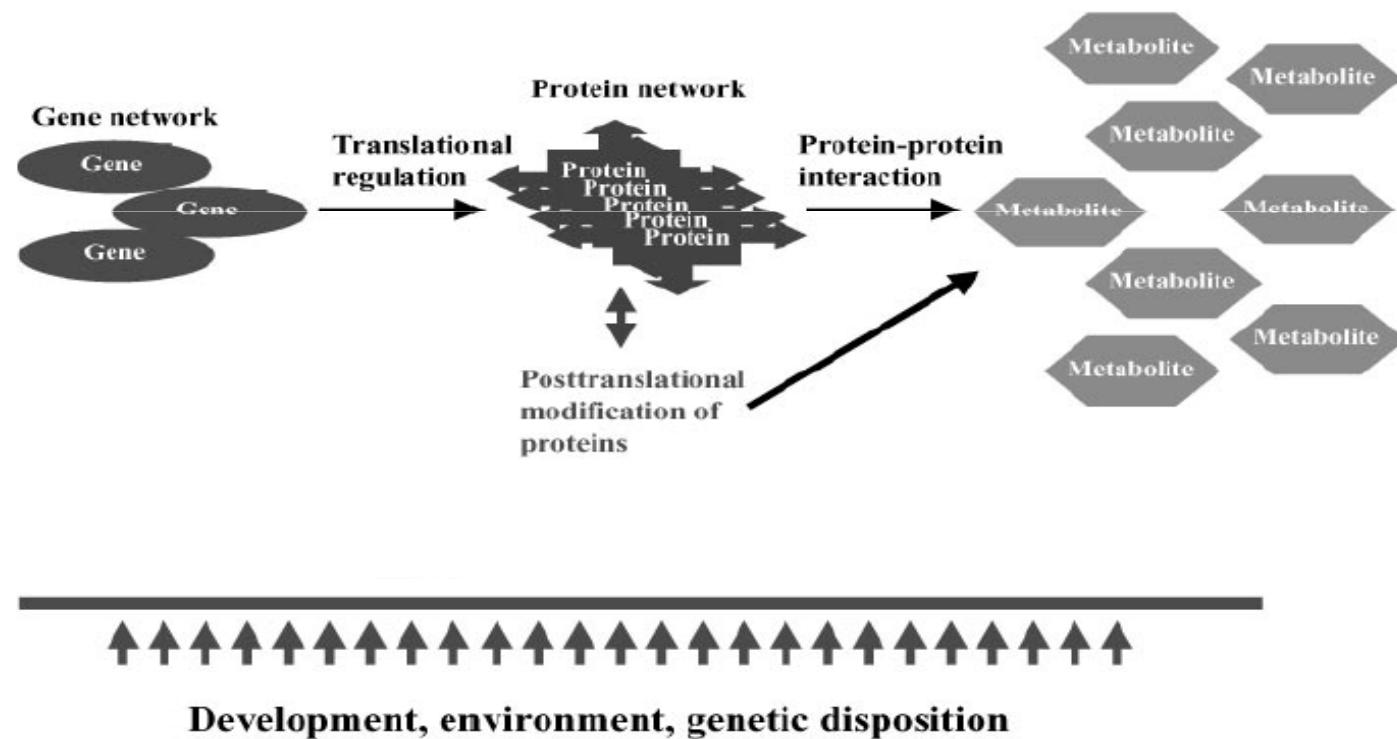


Figure 1 Amplification of a metabolic network and feedback regulation in response to developmental and environmental conditions.



Plant Molecular Biology 48: 155–171, 2002.
© 2002 Kluwer Academic Publishers. Printed in the Netherlands.

Metabolomics – the link between genotypes and phenotypes

Oliver Fiehn

Max Planck Institute of Molecular Plant Physiology, 14421 Potsdam, Germany
(e-mail fiehn@mpimp-golm.mpg.de)

Key words: functional genomics, mass spectrometry, metabolism, metabolite profiling

Growth under controlled conditions and statistically sound plots*

Harvest with rapid freezing

Tissue homogenization

Comprehensive extraction

Aliquoted

Derivatization

GC/MS

RPLC and HILIC[†]
+ UV/ECD/MS

Data deconvolution

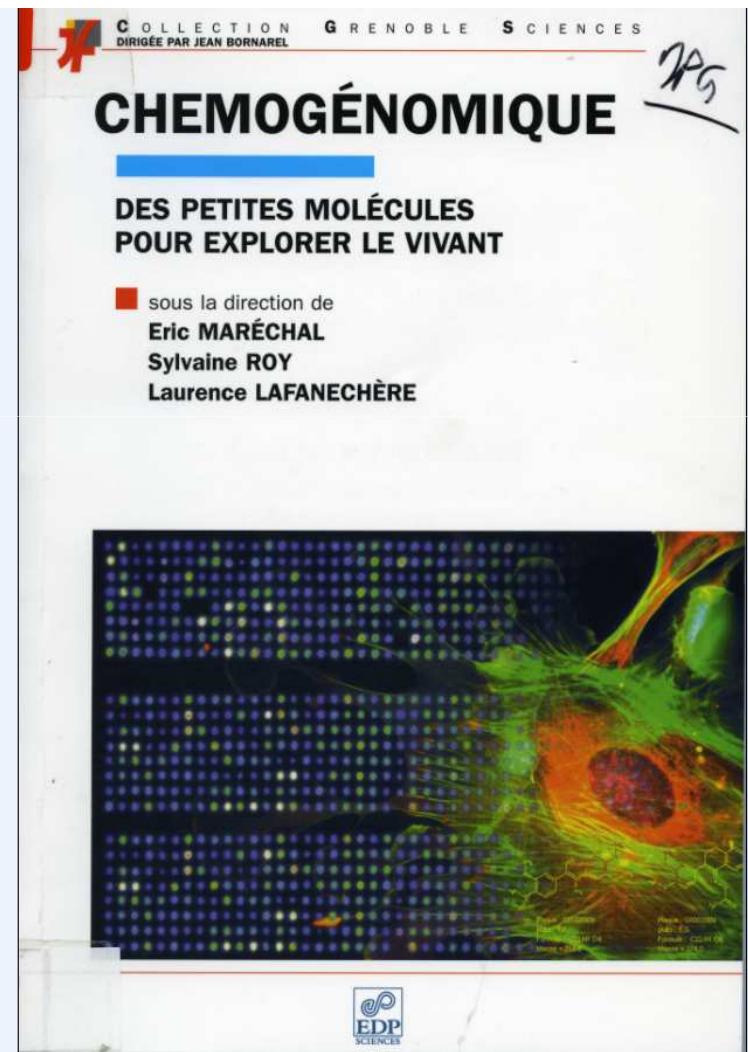
Data deconvolution

Network computation

Current Opinion in Biotechnology

Proposed scheme for comprehensive metabolomic data acquisition.

La chémogénomique



= criblage à haut débit de molécules naturelles ou de synthèses pour découvrir de nouvelles cibles thérapeutiques et de nouveaux médicaments

...

Ou de nouvelles molécules interférant avec Fonctions cellulaires ou des molécules d'intérêt en agronomie et microbiologie (pesticides, herbicides, etc ...)

- Les chimiothèques

- les criblages de molécules et la mesure de l'effet biologique

- Intérêt en post-génomique

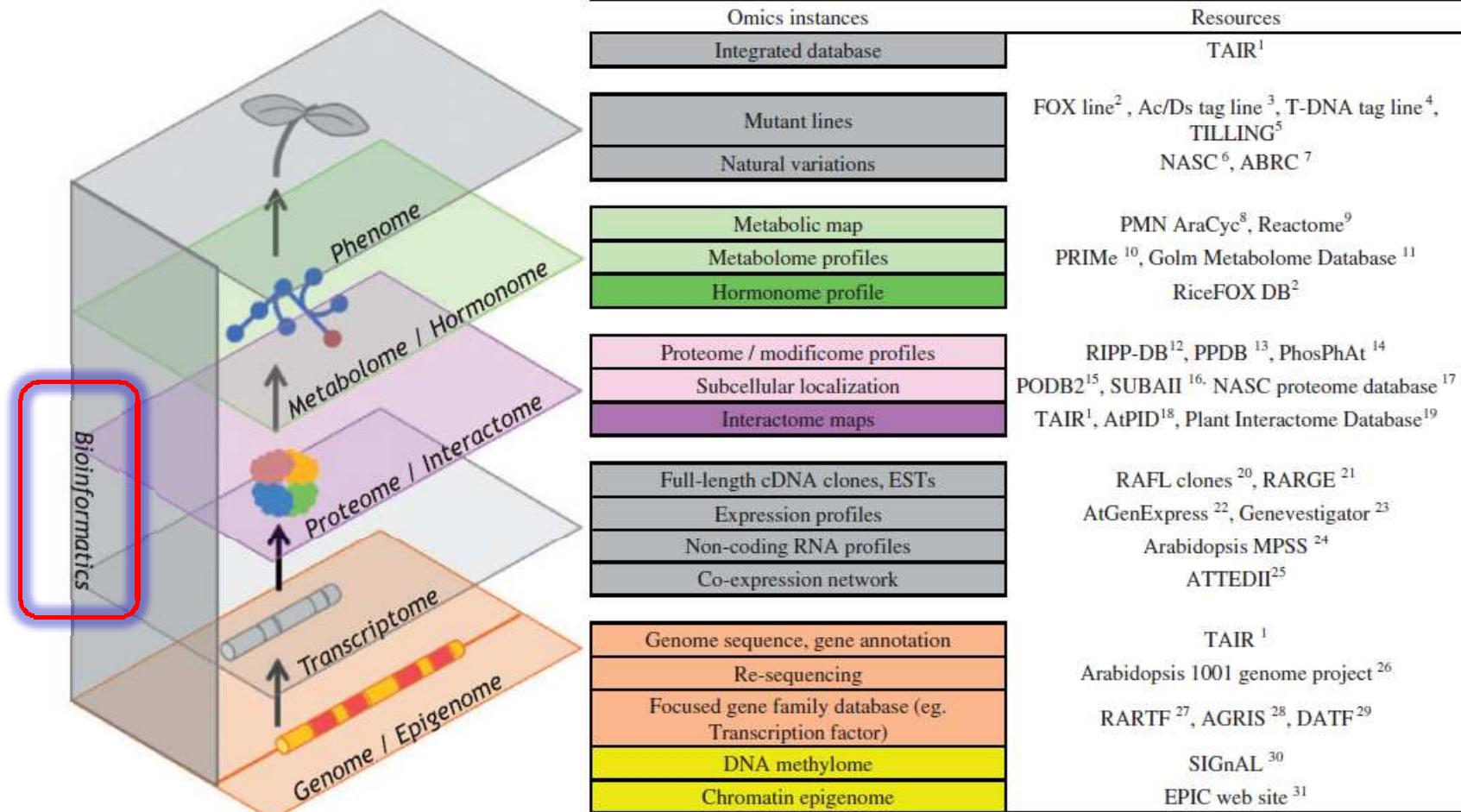


Fig. 1 An updated omic space with emerging omics layers: epigenome, interactome and hormonome added to each of the closely related layers with illustrative resources for *Arabidopsis* available on the web. ¹<http://www.arabidopsis.org/>, ²<http://ricefox.psc.riken.jp/>, ³<http://rarge.gsc.riken>