

Genómica funcional

Introducción

Valencia, Septiembre 2011
CEFIRE

Jorge Jiménez

jjimenez@cipf.es

<http://www.bioinformatico.es>

Bioinformatics and Genomics Department
Centro de Investigación Principe Felipe (CIPF)
(Valencia, Spain)

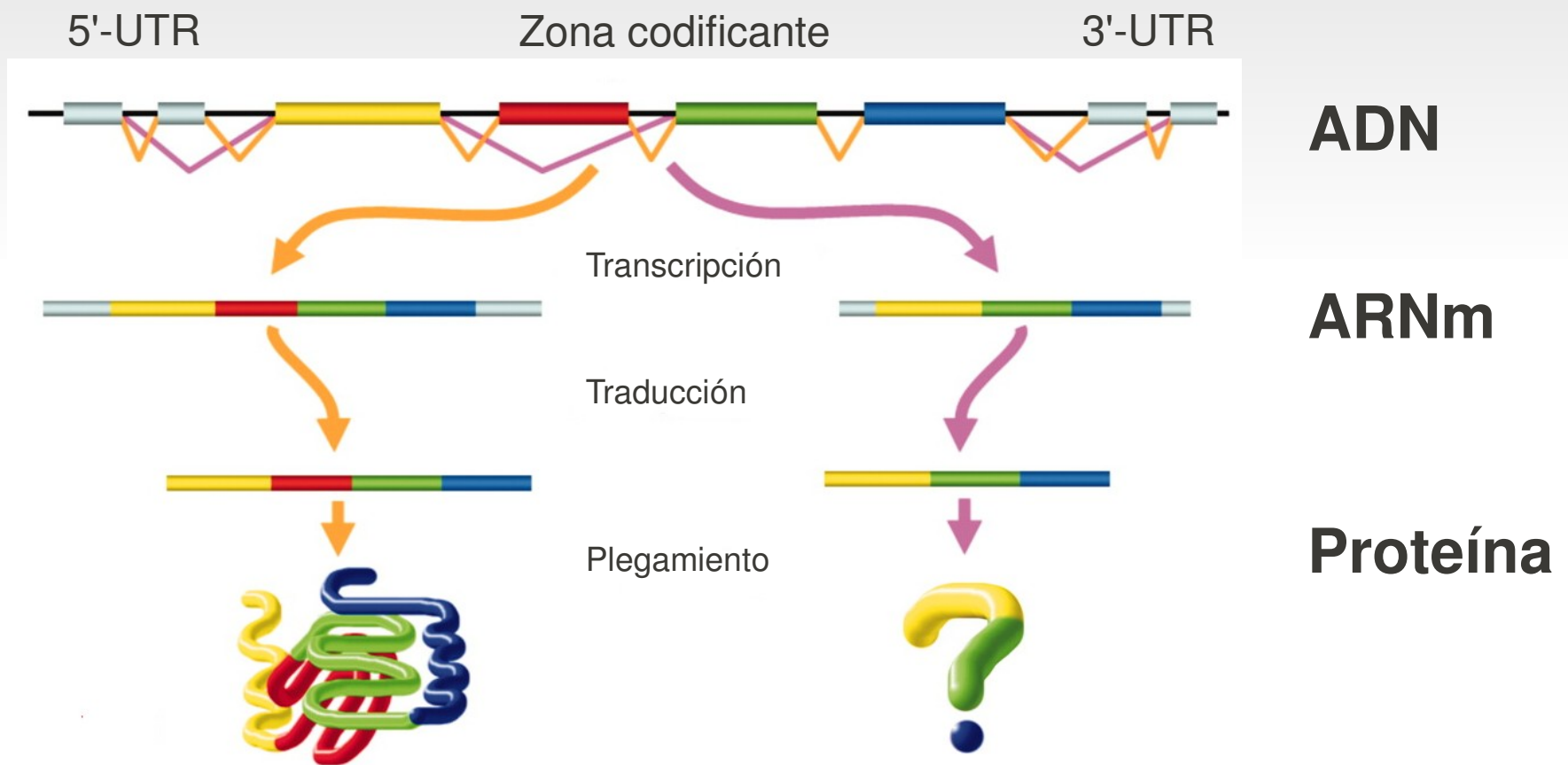


Definición y objetivos

- Entender la relación entre el genotipo y el fenotipo
- Entender la función de los genes y la de los demás componentes del genoma
- Estudio a nivel de ADN, ARN y proteína
- Utiliza datos generados por máquinas de alto rendimiento

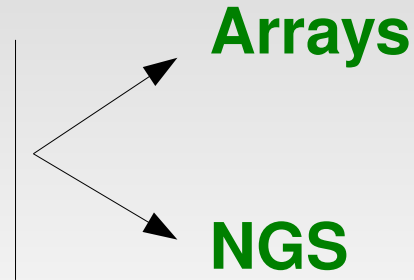
Estructura de un gen

- Exones, intrones, zonas UTR



Niveles

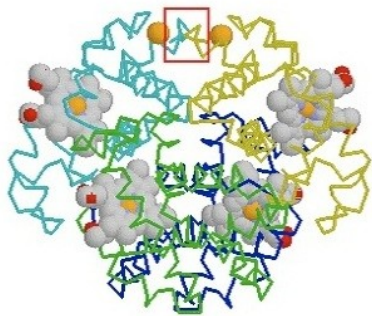
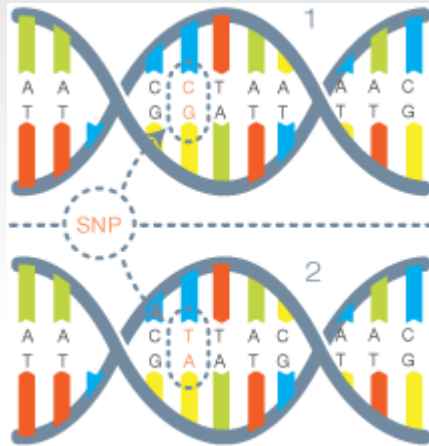
- ADN: SNPs, CNVs, etc.
- ARN: expresión, microARN, etc.
- Proteínas: interacciones proteína-proteína
- Anotación genómica
- Anotación funcional



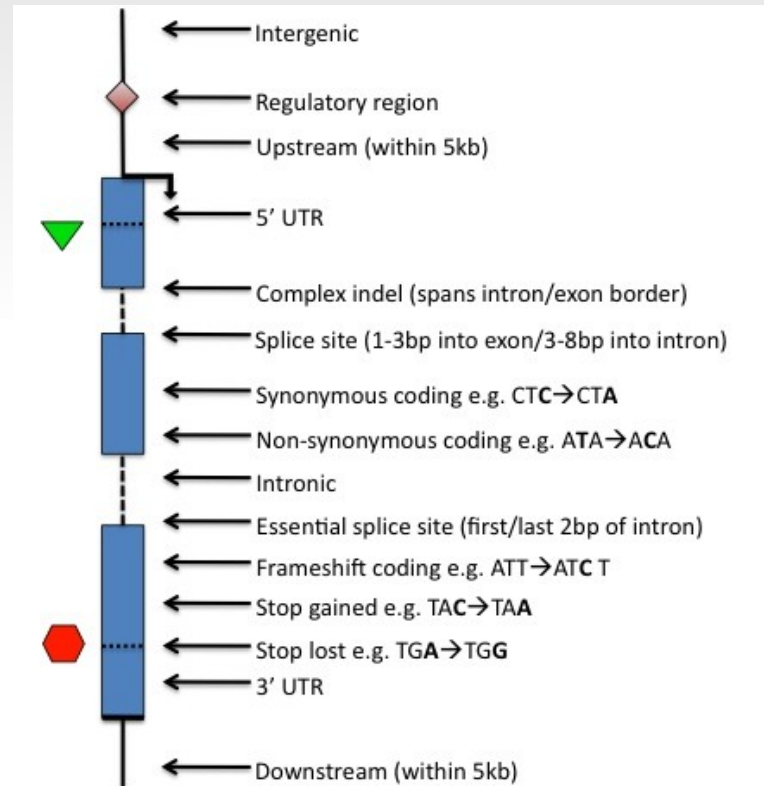
BIOINFORMÁTICA!!

ADN

SNP



Consecuencias



Others: Within non-coding gene, Within mature miRNA, NMD transcript

SNP

Técnicas de detección:

- Secuenciación: Sanger
- Espectrofotometría de masas
- single-strand conformation polymorphism (SSCP)
- Análisis electroquímico
- HPLC y electroforesis en gel
- restriction fragment length polymorphism
- Análisis por hibridación

Técnicas de alto rendimiento

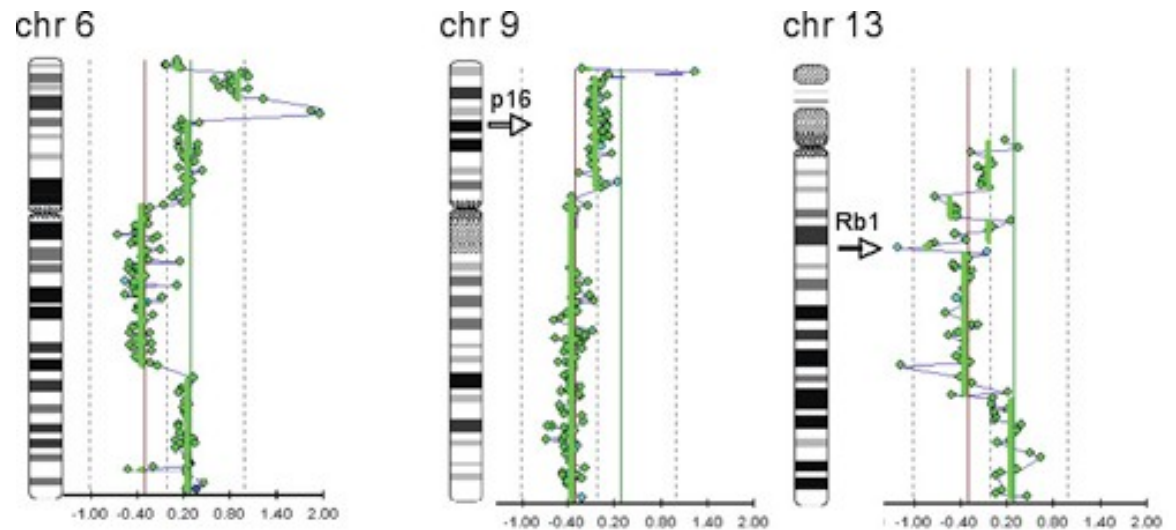
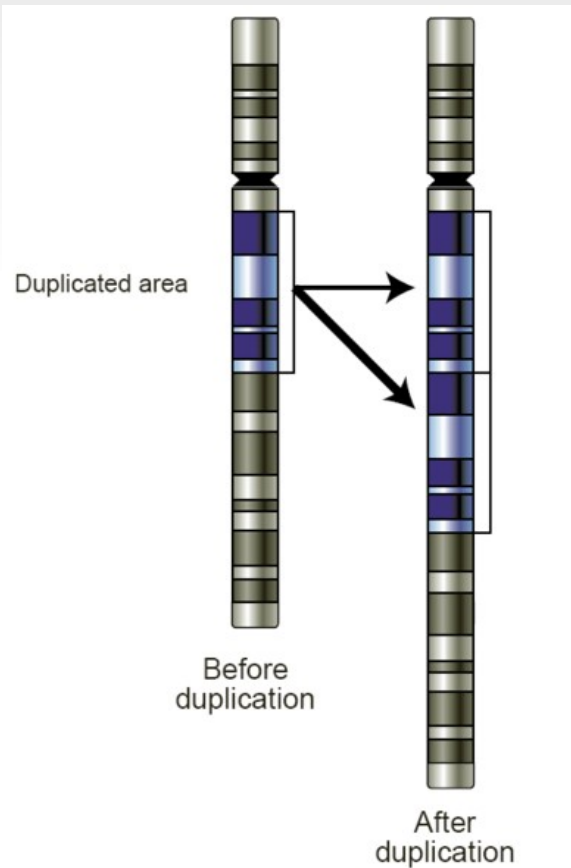
- Chips de SNPs
- NGS: análisis de exoma, secuenciación genoma completo, resecuenciación dirigida

ADN - CNVs

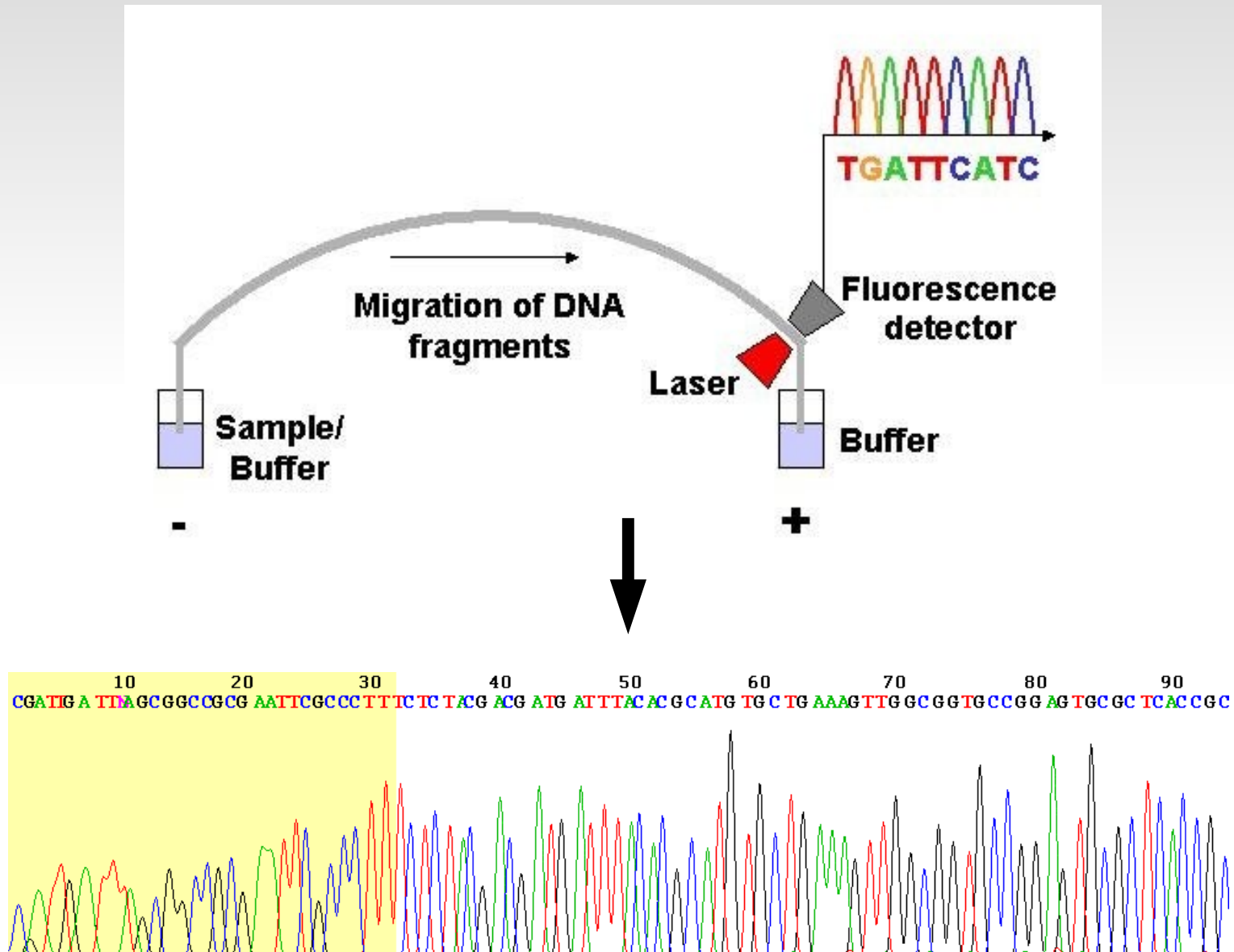
CNVs (12%)

Detección:

- CGH array: comparación entre dos muestras.
- Arrays de SNPs
- NGS



Secuenciación de primera generación



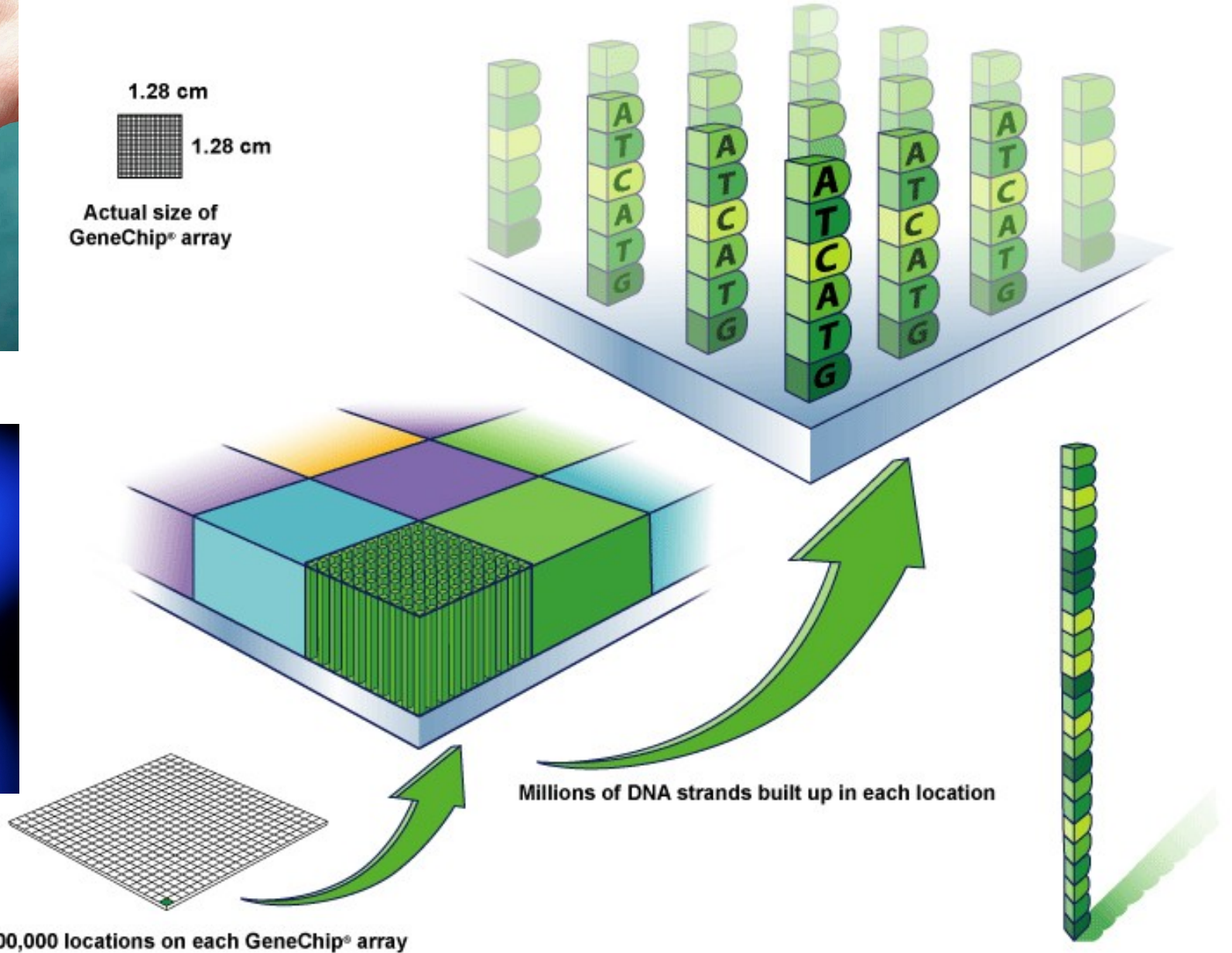
Chips de ADN

- Primer método de alto rendimiento
- Capaces de determinar genes, SNP, mRNA, CNVs,...
- Medida en diferentes condiciones
- Análisis a escala genómica

Microarrays de ADN



1.28 cm
1.28 cm
Actual size of
GeneChip® array

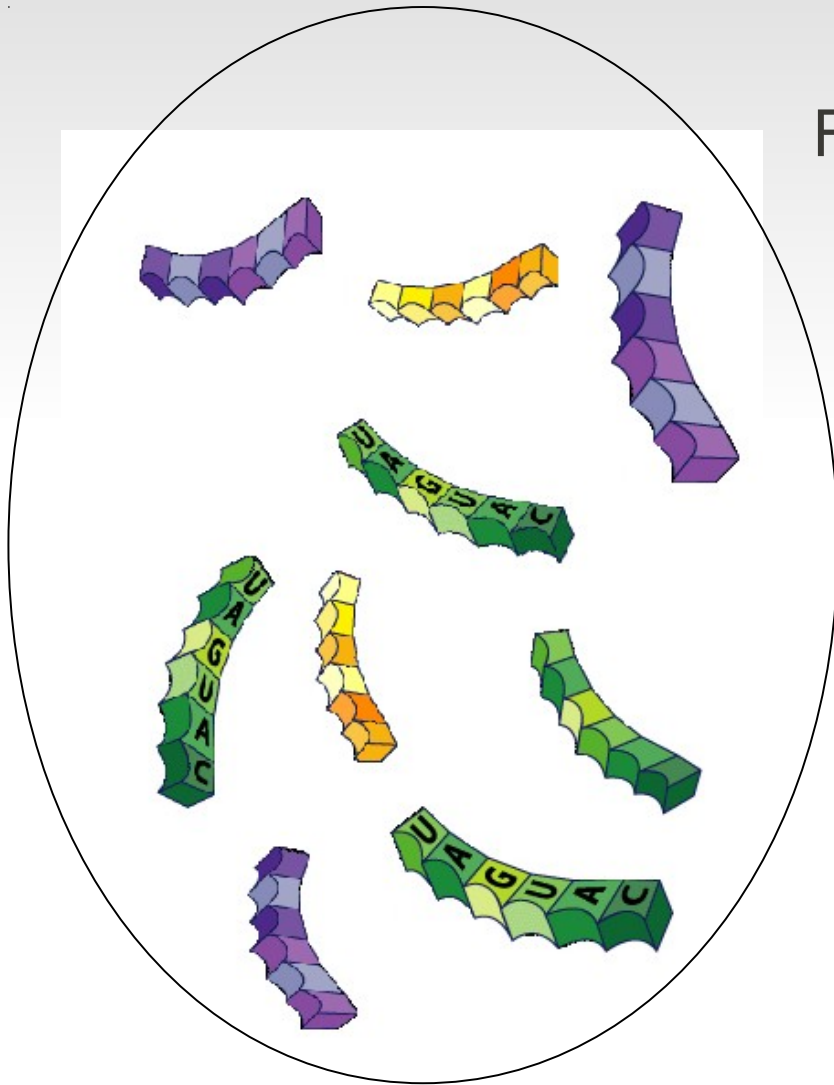


500,000 locations on each GeneChip® array

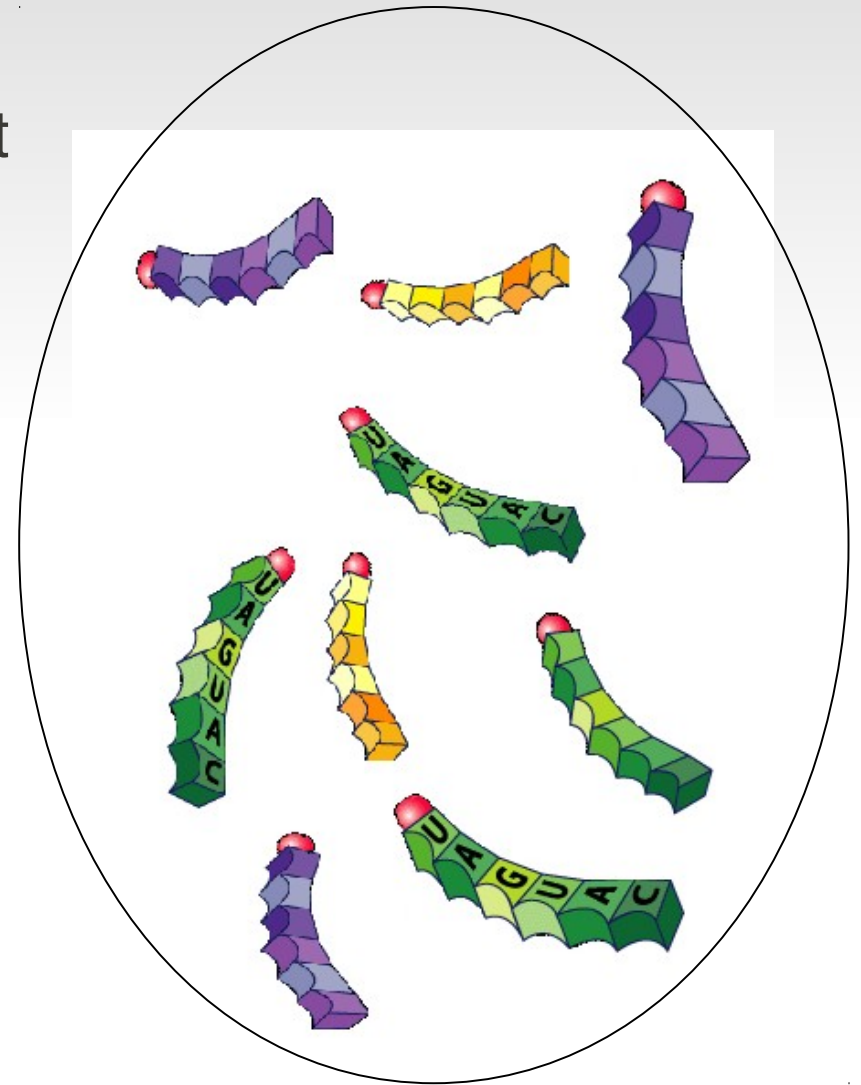
Millions of DNA strands built up in each location

Actual strand = 25 base pairs

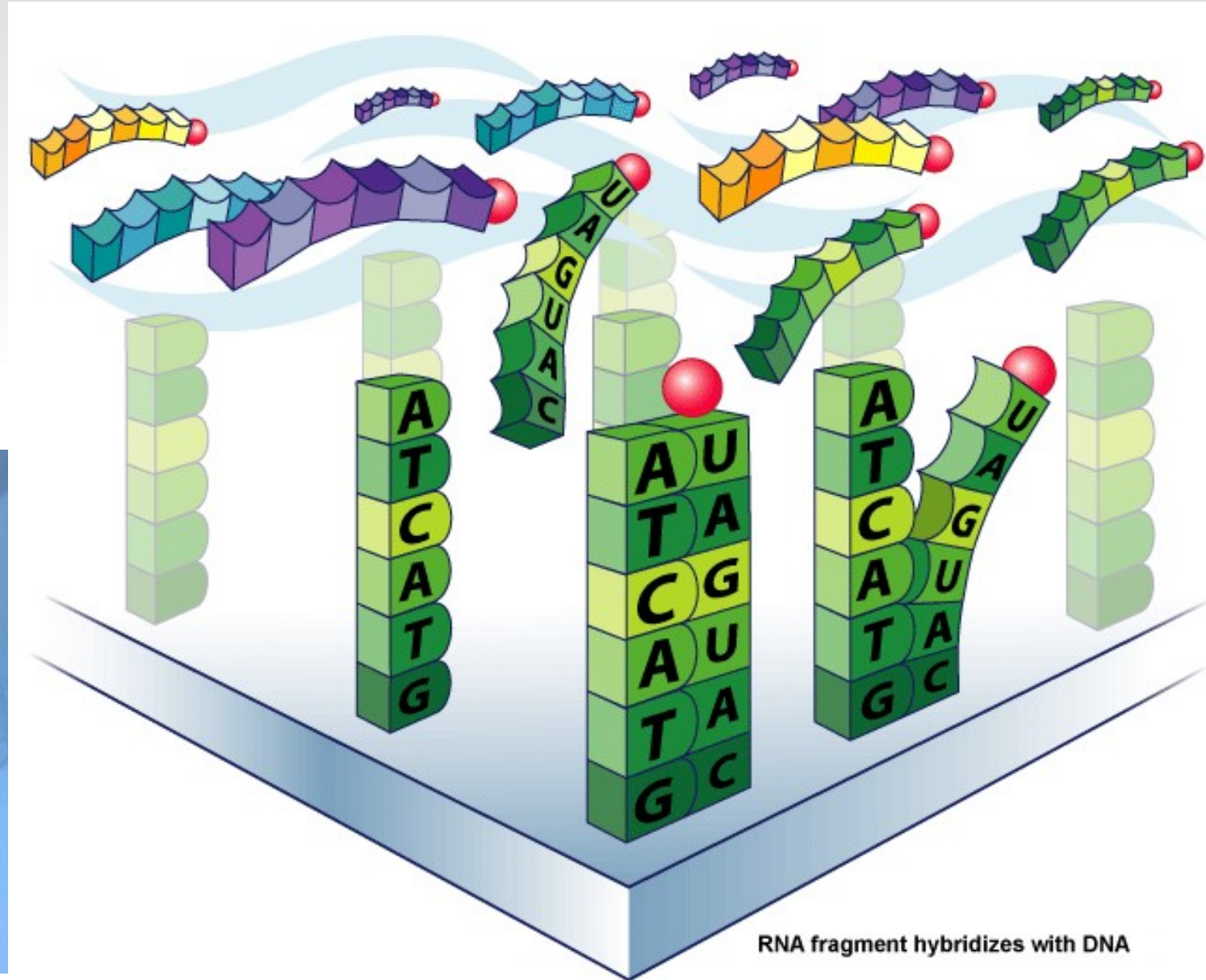
Labelling the Sample



Fluorescent
Dye



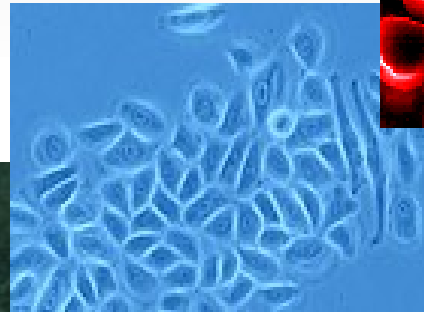
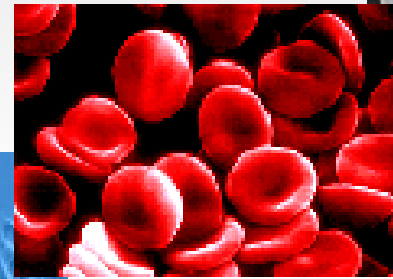
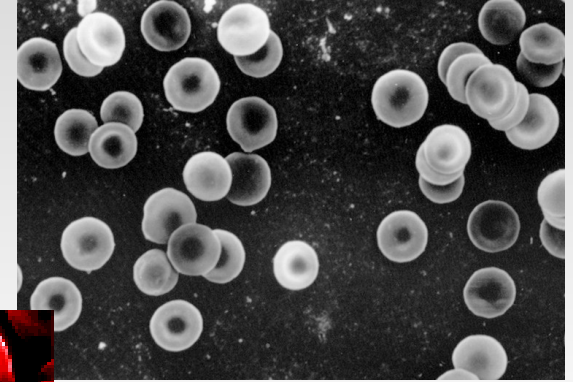
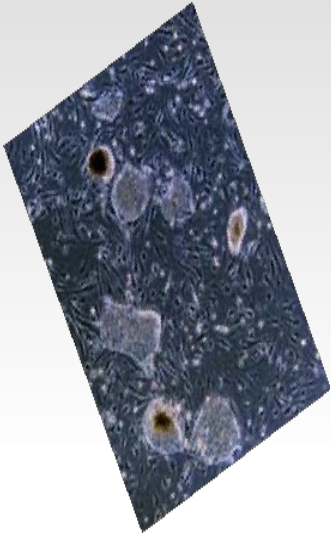
Hibridación



RNA fragment hybridizes with DNA

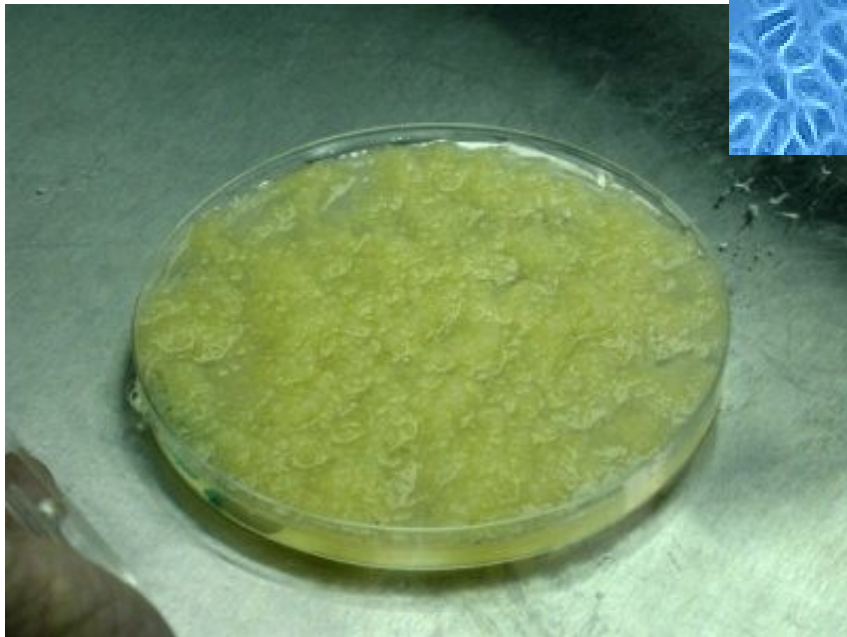
ARN

Queremos saber qué genes
están expresados en
condiciones celulares muy
concretas



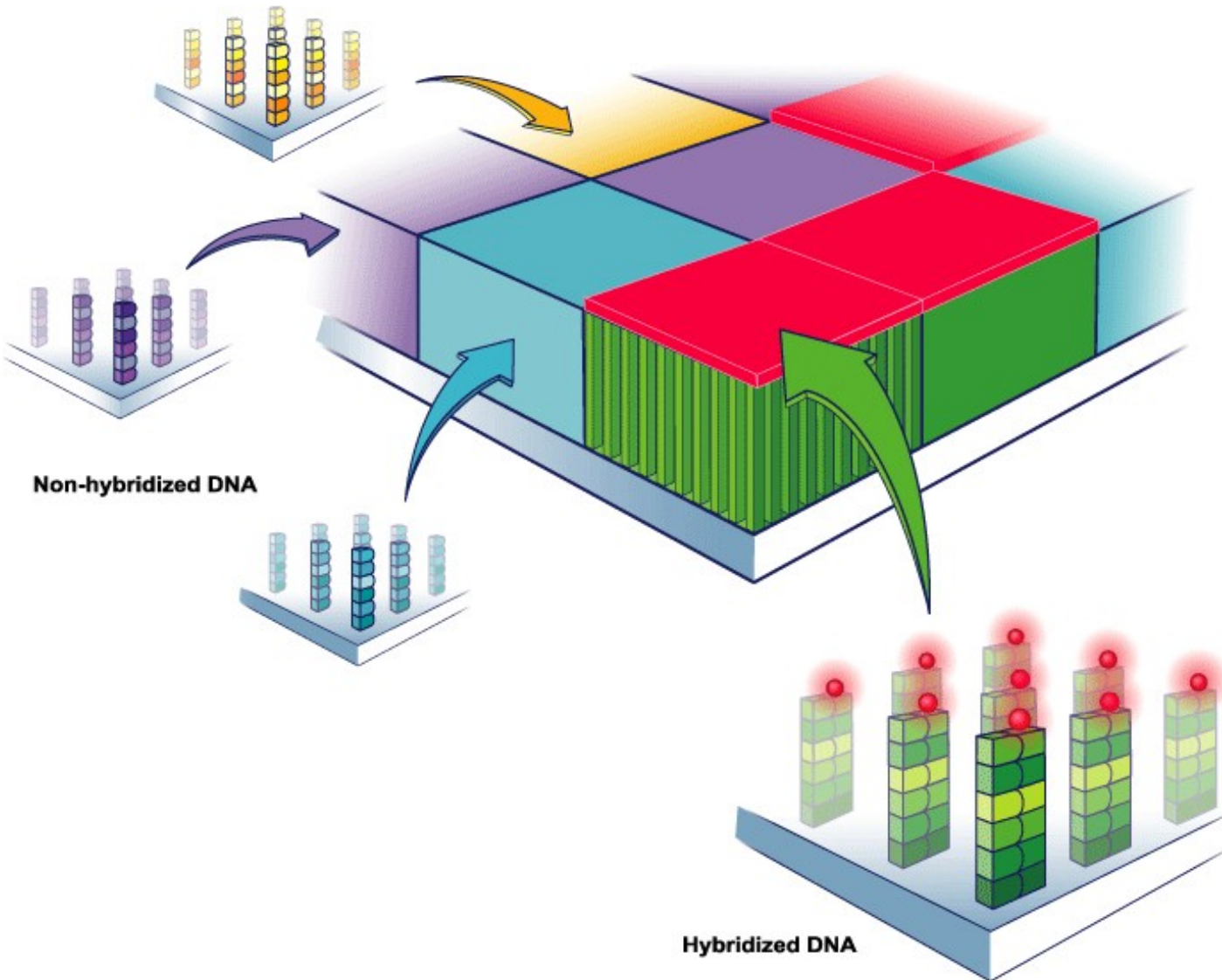
Somos capaces
de extraer todas las
moléculas de mRNA
presentes en las células

y conocer el nivel de expresión
como el indicador de la
concentración en la muestra
biológica



Expression Measurement

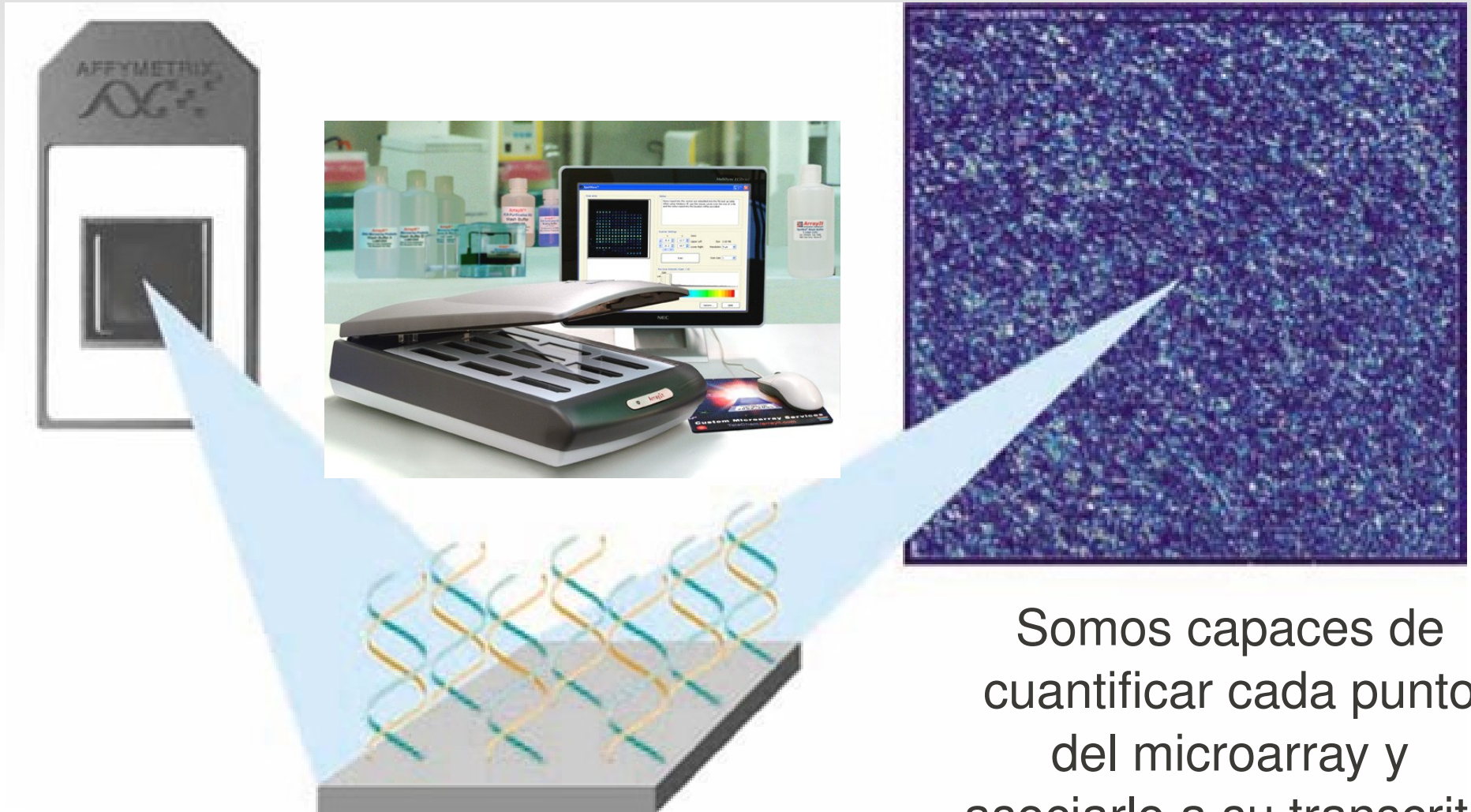
Shining a laser light at GeneChip® array causes tagged DNA fragments that hybridized to glow



Detectaremos aquellos genes que se estén expresando cuando se unan a las sondas marcadas

Cuánto mayor sea la fluorescencia, mayor será la concentración del transcrito

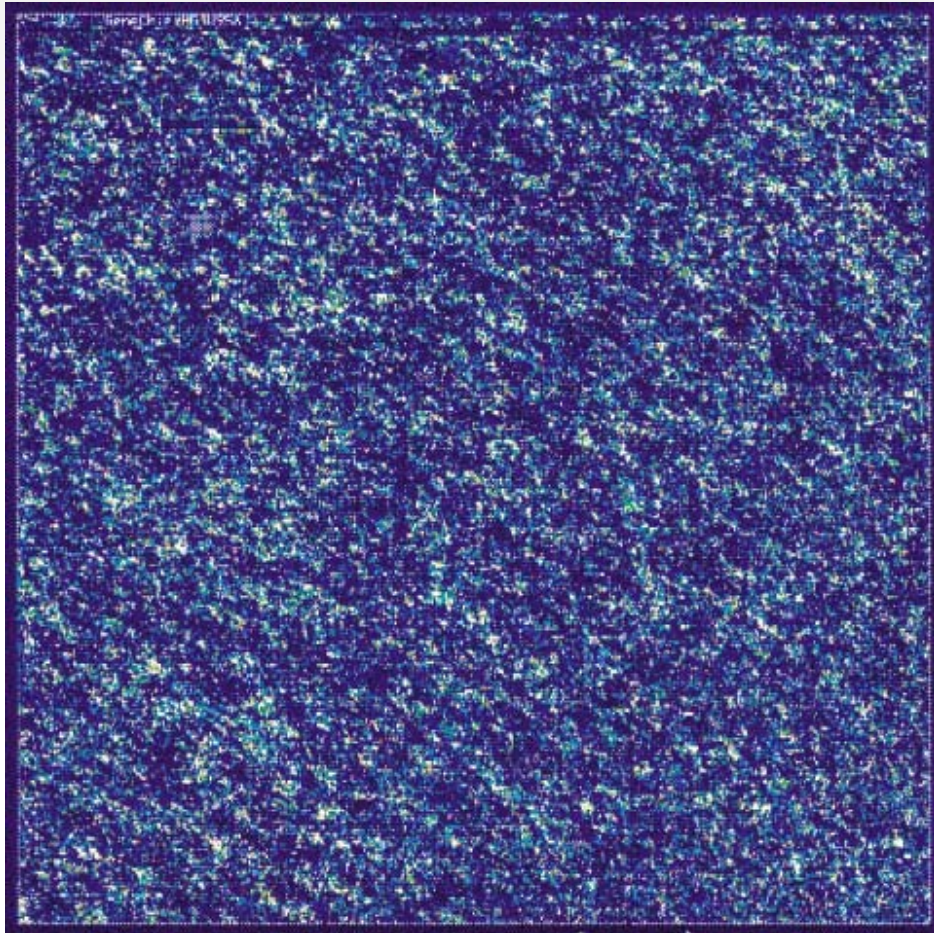
Escáner del microarray



Somos capaces de cuantificar cada punto del microarray y asociarlo a su transcrito

Los datos

Para cada muestra
biológica

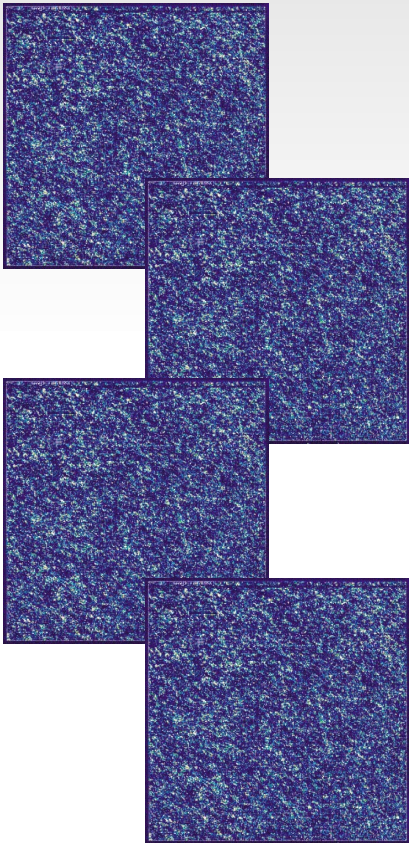


Obtenemos la
intensidad de
fluorescencia de miles
de transcritos

La medida de
intensidad es un
indicador de la
expresión génica

200000_s_at	134.4
200001_at	586.5
200002_at	1868.4
200003_s_at	1232.7
200004_at	1071.6
200005_at	312.8
200006_at	1712.6
200007_at	606.5
200008_s_at	421.9
200009_at	395.6
200010_at	1228.6
200011_s_at	132.5
200012_x_at	2606.3
200013_at	1572.9
200014_s_at	138.7
200015_s_at	124.1
200016_x_at	1058.7
200017_at	889.4
200018_at	3964.2
200019_s_at	1069.9
200020_at	212.1
200021_at	1018.1
200022_at	1254.8
200023_s_at	1202.8
200024_at	2460.6

Varios Microarrays



31307_at	5.53
31308_at	7.07
31309_r_at	6.05
31310_at	7.42
31311_at	7.77
31312_at	9.47
31313_at	8.58
31314_at	7.14
31315_at	9.62
31316_at	5.45
31317_r_at	10.27
31318_at	5.7
31319_at	9.25
31320_at	11.5
31321_at	7.79
31322_at	6.98
31323_r_at	11.18
31324_at	7.97
31325_at	9.53
31326_at	9.67
31327_at	6.48
31328_at	8.92
31329_at	6.11
31330_at	14.44
31331_at	6.3

31307_at	5.66
31308_at	7.14
31309_r_at	5.33
31310_at	7.02
31311_at	7.83
31312_at	9.43
31313_at	8.67
31314_at	7.3
31315_at	9.62
31316_at	5.53
31317_r_at	10.75
31318_at	5.53
31319_at	9.19
31320_at	11.51
31321_at	7.91
31322_at	6.93
31323_r_at	10.27
31324_at	8.12
31325_at	9.37
31326_at	10.16
31327_at	6.2
31328_at	9.11
31329_at	5.86
31330_at	14.32
31331_at	6.4

31307_at	5.52
31308_at	7.05
31309_r_at	5.35
31310_at	7.02
31311_at	7.79
31312_at	9.34
31313_at	8.52
31314_at	7.19
31315_at	9.22
31316_at	5.3
31317_r_at	10.41
31318_at	5.59
31319_at	9.24
31320_at	11.35
31321_at	7.76
31322_at	6.91
31323_r_at	10.32
31324_at	8.06
31325_at	9.33
31326_at	9.92
31327_at	6.2
31328_at	8.81
31329_at	5.81
31330_at	14.3
31331_at	6.42

Chips de SNPs

- Análisis de miles de SNPs a la vez
- Aplicaciones
 - Genética de poblaciones
 - GWAS
 - Análisis de enfermedades
 - Análisis del número de copias del SNP

NGS

- Características más importantes:
 - Paralelización del proceso
 - Generación de millones de secuencias por experimento
- Ventajas
 - Barato
 - Gran versatilidad de análisis

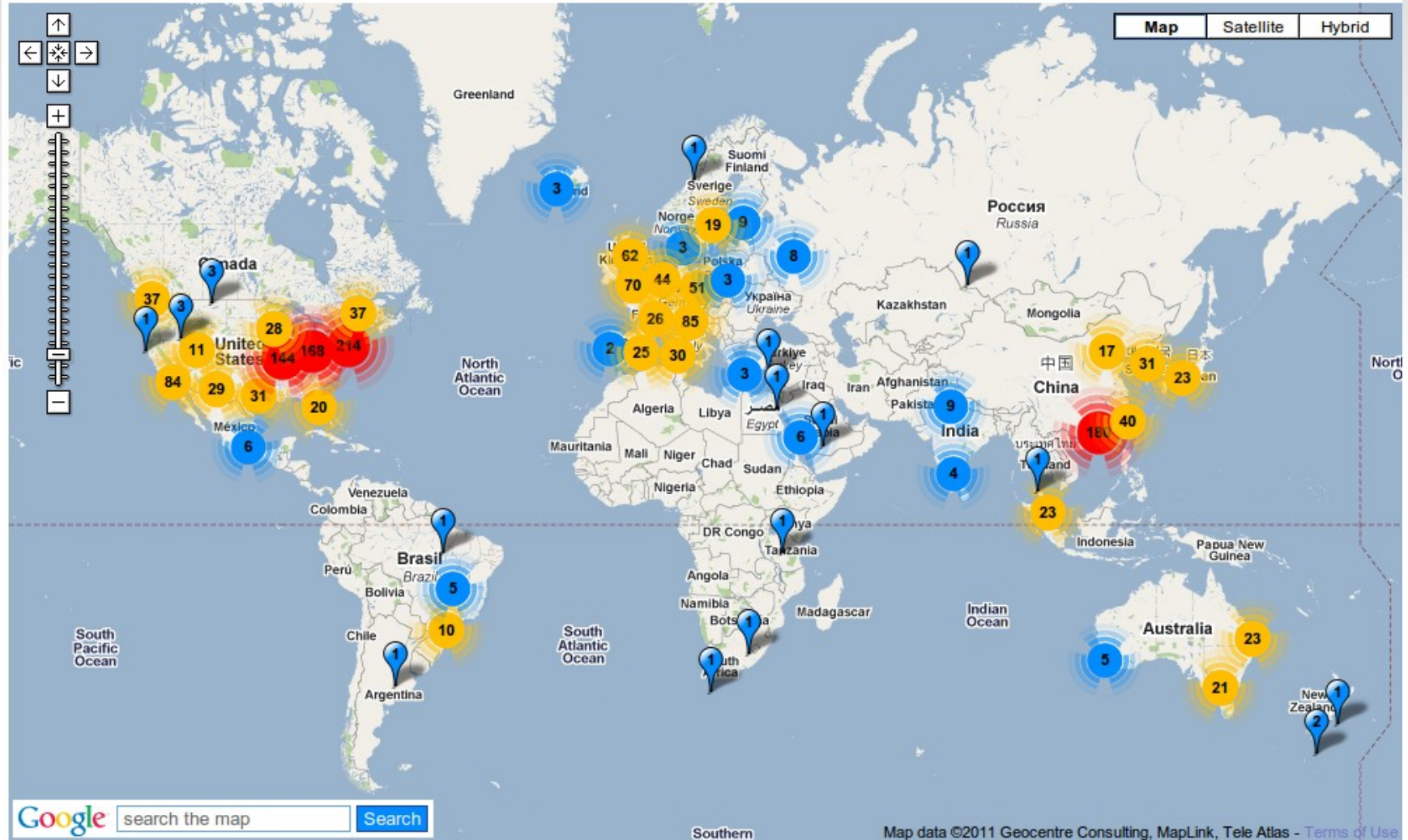
NGS

- Se utilizan para secuenciación de novo, resecuenciación dirigida de ADN (ej: exoma) y estudio de expresión (RNA-Seq)
- Secuenciadores de segunda generación
- Plataformas más importantes:
 - 454 – pirosecuenciación
 - Illumina
 - SOLiD
 - etc.

NGS

Next Generation Genomics: World Map of High-throughput Sequencers

Show all platforms Illumina GA2 Illumina HiSeq Ion Torrent PacBio Polonator Roche/454 SOLID Service Provider



Revolución tecnológica

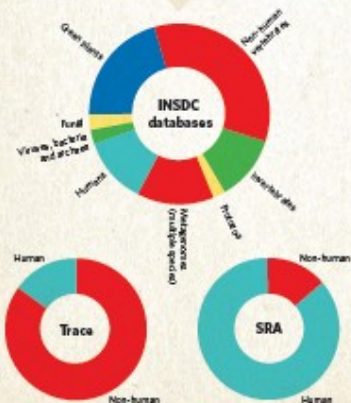


THE SEQUENCE EXPLOSION

At the time of the announcement of the first drafts of the human genome in 2000, there were 8 billion base pairs of sequence in the three main databases for 'finished' sequence: GenBank, run by the US National Center for Biotechnology Information; the DNA Data Bank of Japan; and the European Molecular Biology Laboratory (EMBL) Nucleotide Sequence Database. The databases share their data regularly as part of the International Nucleotide Sequence Database Collaboration (INSDC). In the subsequent first post-genome decade, they have added another 2.70 billion bases to the collection of finished sequence, doubling the size of the database roughly every 18 months. But this number is dwarfed by the amount of raw sequence that has been created and stored by researchers around the world in the Trace Archive and Sequence Read Archive (SRA). See Editorial, page 640, and human genome special site www.nature.com/humangenome

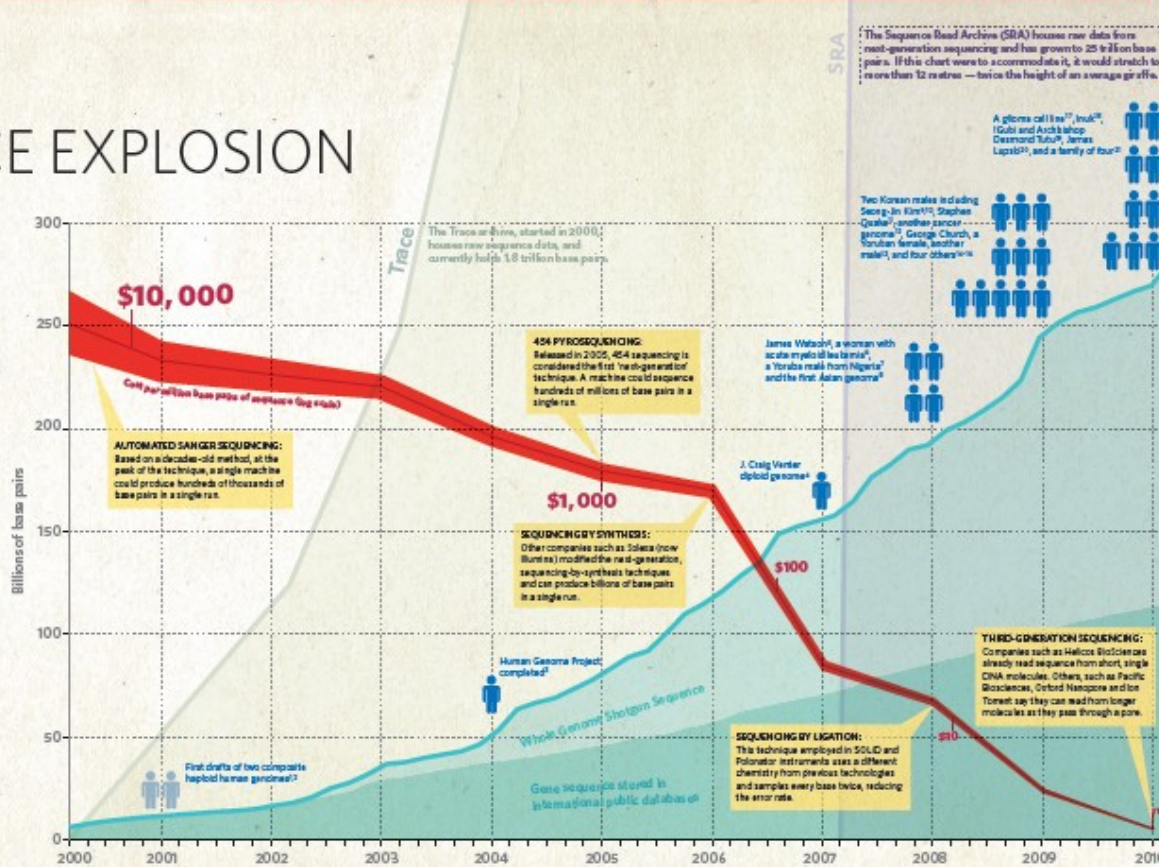
DNA SEQUENCES BY TAXONOMY

International Nucleotide Sequence Database Collaborators
The main repositories of 'finished' sequence spans a wide range of organisms, representing the many priorities of scientists worldwide.



Trace Archive Developed to house the raw output of high-throughput sequencers built in the late 1990s, the trace archive spans a wide range of taxa.

Sequence Read Archive Houses raw data from next-generation sequencing, including multiple coverage for more than 170 people.



HOW MANY HUMAN GENOMES?

The graphic shows all published, fully sequenced human genomes since 2000, including nine from the first quarter of 2010. Some are resequencing efforts on the same person and the list does not include unpublished completed genomes.

1. Venter J. C. et al. *Science* 291, 1304-1305 (2001).
2. International Human Genome Sequencing Consortium *Nature* 409, 845-921 (2006).
3. International Human Genome Sequencing Consortium *Nature* 428, 921-942 (2004).
4. Levy S. et al. *PLoS Biol* 5, e254 (2007).
5. Wheeler D. A. et al. *Nature* 421, 879-876 (2008).
6. Liu T. et al. *Nature* 456, 65-72 (2008).
7. Bentley D. R. et al. *Nature* 456, 51-59 (2008).
8. Wang J. et al. *Nature* 456, 65-68 (2008).
9. Ahe S.-M. et al. *Genome Res* 19, 1622-1629 (2009).
10. Kim J.-I. et al. *Nature* 460, 1071-1075 (2009).
11. Pathak S. et al. *Nature* 460, 1071-1075 (2009).
12. Hardy C. R. et al. *N. Engl. J. Med.* 361, 1058-1066 (2009).
13. Derracq, R. et al. *Science* 327, 78-81 (2009).
14. McKernan K. J. et al. *Genome Res* 19, 1527-1541 (2009).
15. Plessence, C. D. et al. *Nature* 462, 101-106 (2010).
16. Plessence, C. D. et al. *Nature* 462, 104-109 (2010).
17. Clark, M. J. et al. *PLoS Genet* 6, e1000823 (2010).
18. Saarnanen, M. et al. *Nature* 462, 757-763 (2010).
19. Schaefer, J. C. et al. *Nature* 462, 943-947 (2010).
20. Lagan, J. R. et al. *N. Engl. J. Med.* doi:10.1056/NEJMed090904 (2010).
21. Beach, J. C. et al. *Science* doi:10.1126/science.1186802 (2010).

Page size by comparison

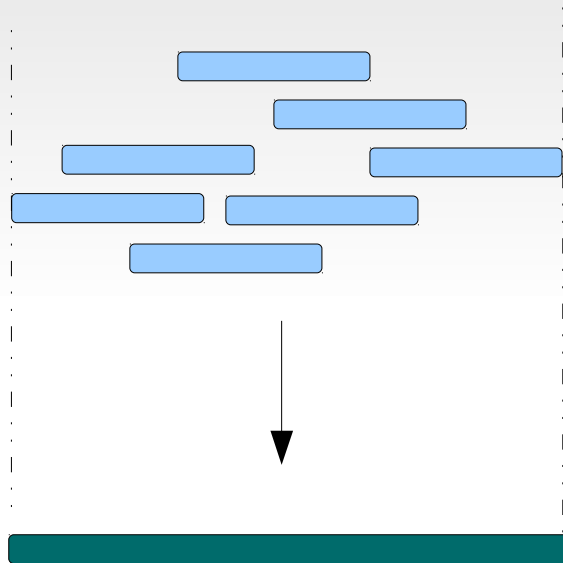
Comparativa de plataformas

Platform	Library/ template preparation	NGS chemistry	Read length (bases)	Run time (days)	Gb per run	Machine cost (US\$)	Pros	Cons	Biological applications	Refs
Roche/454's GS FLX Titanium	Frag, MP/ emPCR	PS	330*	0.35	0.45	500,000	Longer reads improve mapping in repetitive regions; fast run times	High reagent cost; high error rates in homo- polymer repeats	Bacterial and insect genome <i>de novo</i> assemblies; medium scale (<3 Mb) exome capture; 16S in metagenomics	D. Muzny, pers. comm.
Illumina/ Solexa's GA _{II}	Frag, MP/ solid-phase	RTs	75 or 100	4 [‡] , 9 [§]	18 [‡] , 35 [§]	540,000	Currently the most widely used platform in the field	Low multiplexing capability of samples	Variant discovery by whole-genome resequencing or whole-exome capture; gene discovery in metagenomics	D. Muzny, pers. comm.
Life/APG's SOLiD 3	Frag, MP/ emPCR	Cleavable probe SBL	50	7 [‡] , 14 [§]	30 [‡] , 50 [§]	595,000	Two-base encoding provides inherent error correction	Long run times	Variant discovery by whole-genome resequencing or whole-exome capture; gene discovery in metagenomics	D. Muzny, pers. comm.
Polonator G.007	MP only/ emPCR	Non- cleavable probe SBL	26	5 [§]	12 [§]	170,000	Least expensive platform; open source to adapt alternative NGS chemistries	Users are required to maintain and quality control reagents; shortest NGS read lengths	Bacterial genome resequencing for variant discovery	J. Edwards, pers. comm.
Helicos BioSciences HeliScope	Frag, MP/ single molecule	RTs	32*	8 [‡]	37 [‡]	999,000	Non-bias representation of templates for genome and seq-based applications	High error rates compared with other reversible terminator chemistries	Seq-based methods	91
Pacific Biosciences (target release: 2010)	Frag only/ single molecule	Real-time	964*	N/A	N/A	N/A	Has the greatest potential for reads exceeding 1 kb	Highest error rates compared with other NGS chemistries	Full-length transcriptome sequencing; complements other resequencing efforts in discovering large structural variants and haplotype blocks	S. Turner, pers. comm.

*Average read-lengths. [‡]Fragment run. [§]Mate-pair run. Frag, fragment; GA, Genome Analyzer; GS, Genome Sequencer; MP, mate-pair; N/A, not available; NGS, next-generation sequencing; PS, pyrosequencing; RT, reversible terminator; SBL, sequencing by ligation; SOLiD, support oligonucleotide ligation detection.

NGS

Ensamblaje

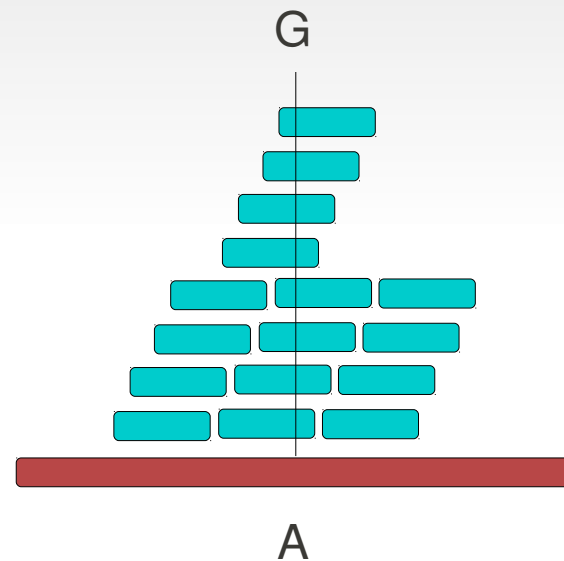


454-Pirosecuenciación

Lecturas largas

Secuenciación de genomas

Alineamiento

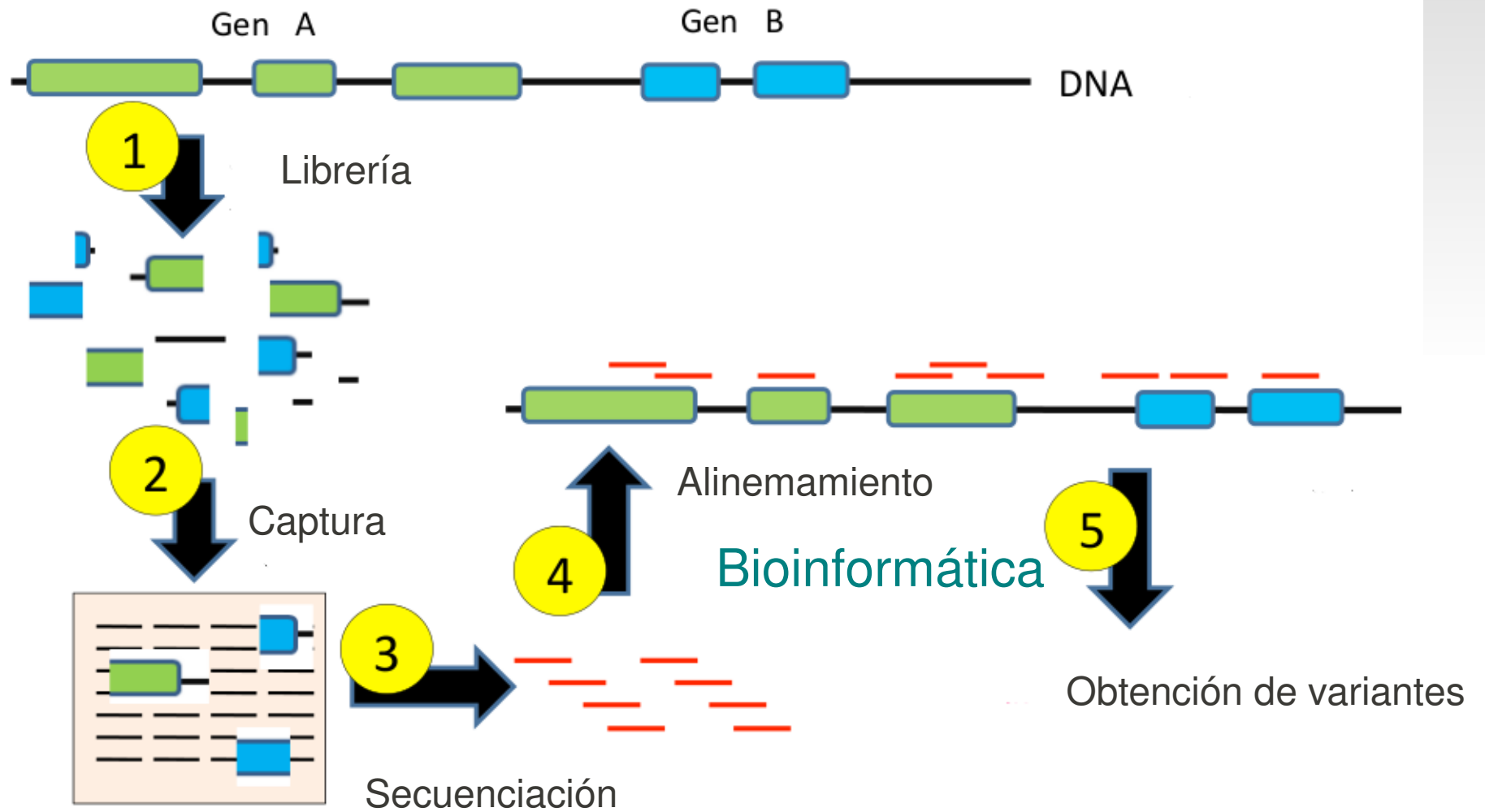


SOLiD - Illumina

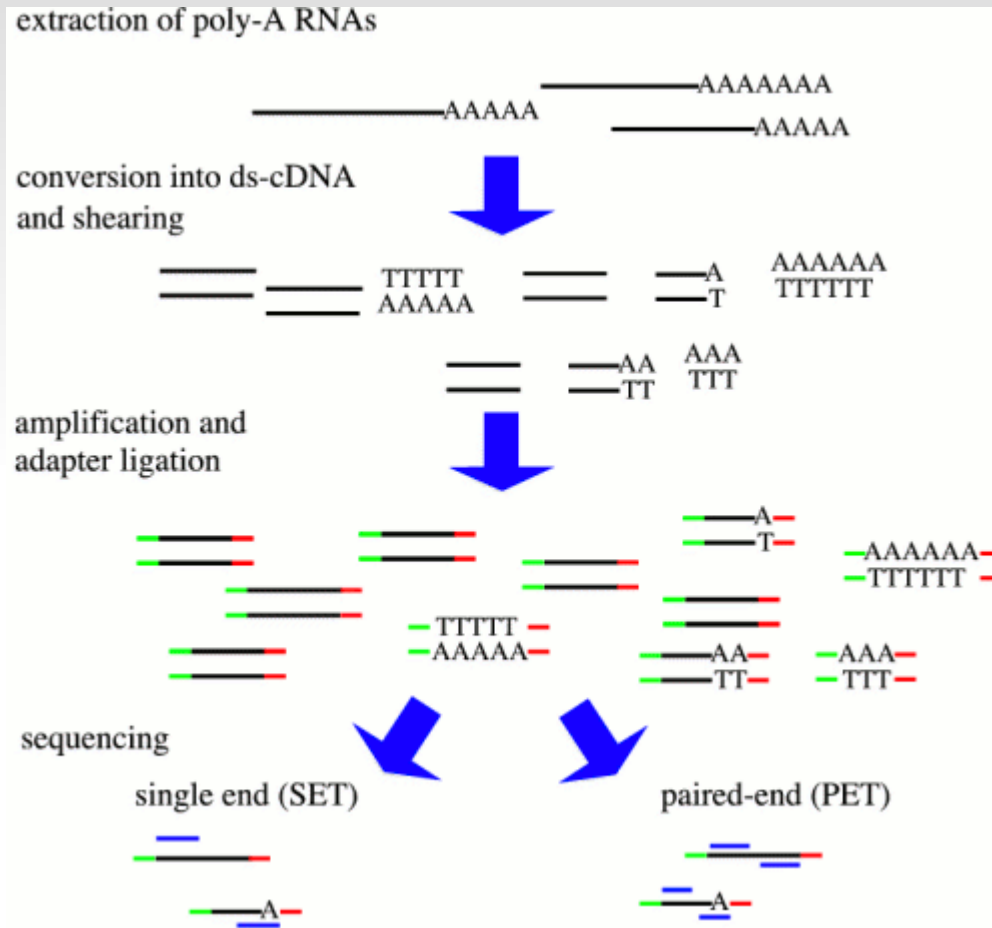
Lecturas cortas

Resecuenciación ADN

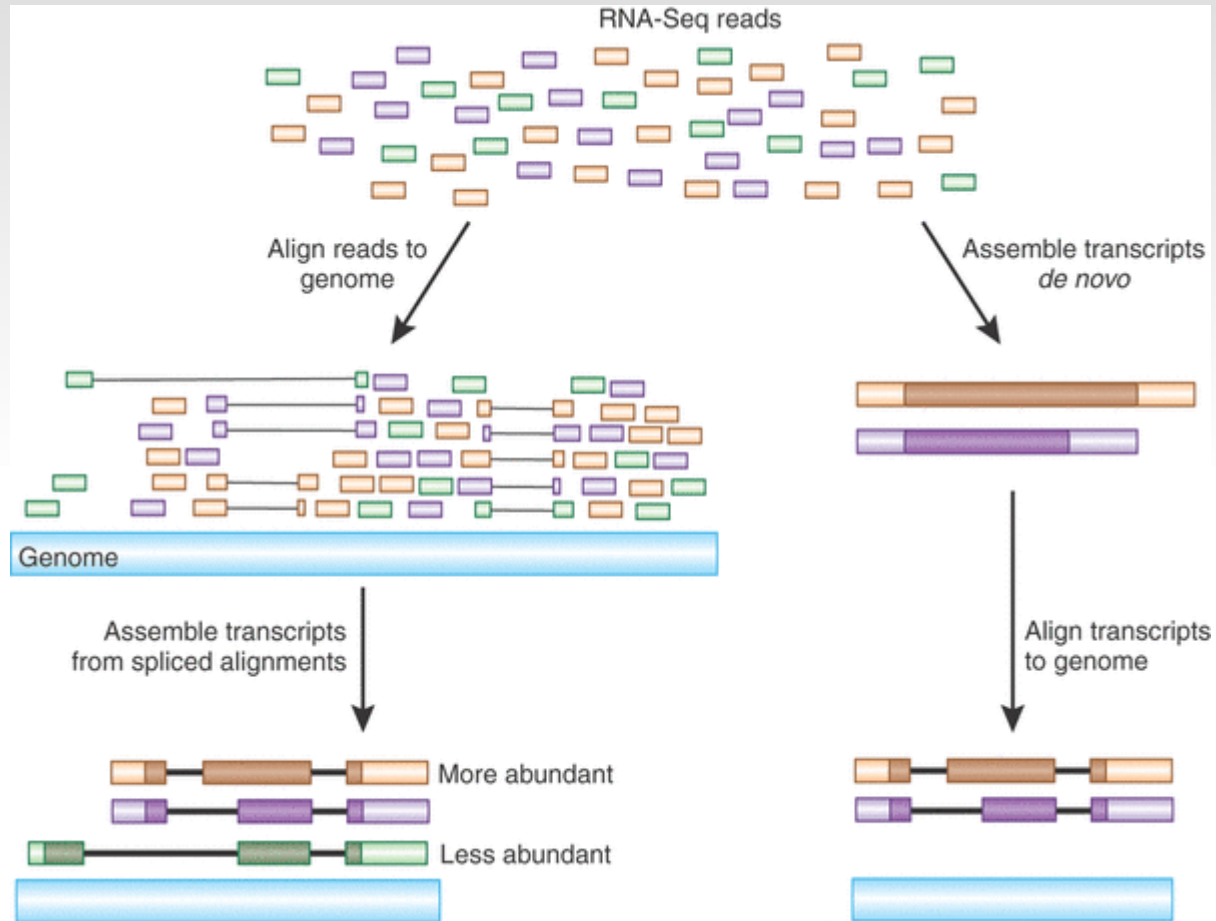
Exoma



RNA-Seq



RNA-Seq



RNA-Seq

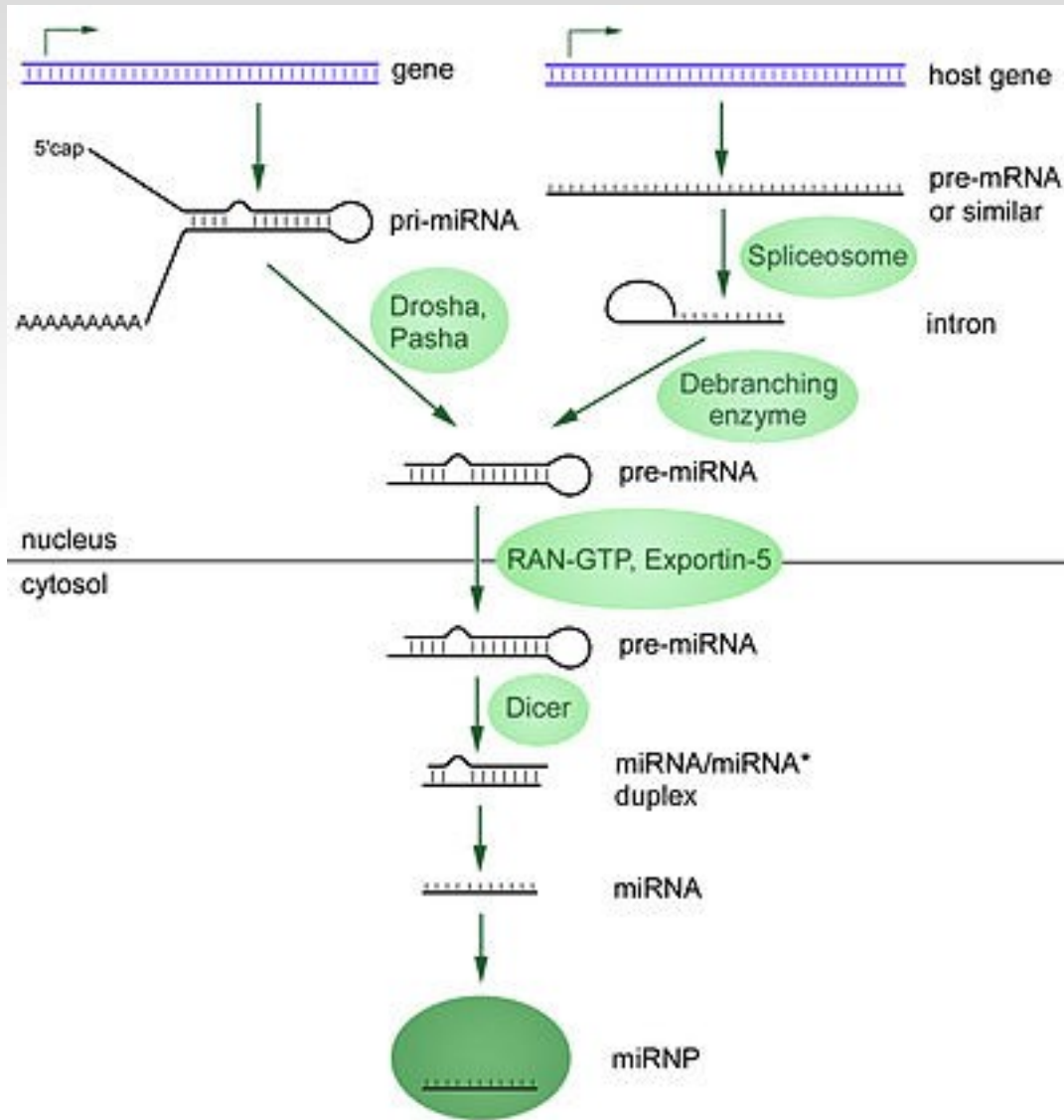
RNA-Seq

- detección de transcrito no dirigida
- no necesita genoma de referencia
- específico de cadena
- descubrimiento de splicing
- detecta variaciones en el ADN
- detecta transcritos raros

Microarrays de ARN

- restringido a las sondas en el diseño
- necesita conocer el genoma
- no es específico de cadena
- difícil trabajar con exones
- no aporta información sobre la secuencia de los transcritos
- difícil de detectar transcritos poco raros
- técnica y análisis muy optimizado

microRNAs



Mecanismos de actuación:

- degradación de la proteína durante la traducción
- inhibición de la elongación de la traducción
- terminación prematura de la traducción (disgregación de los ribosomas)
- inhibición de la iniciación de la traducción

Técnicas de detección:

- RT-PCR + qRT-PCR
- Microarrays de expresión
- NGS

Interacciones proteína-proteína

- Técnicas de detección moleculares
 - Doble híbrido
 - Espectrofotometría de masas
- Técnicas de alto rendimiento
 - Microarrays de expresión

Anotación genómica

- Identificación de elementos del genoma
- Anotación estructural: búsqueda de genes y otros elementos en el genoma
 - Búsqueda de secuencias de traducción
 - Búsqueda de dominios
 - Comparación con bases de datos: BLAST

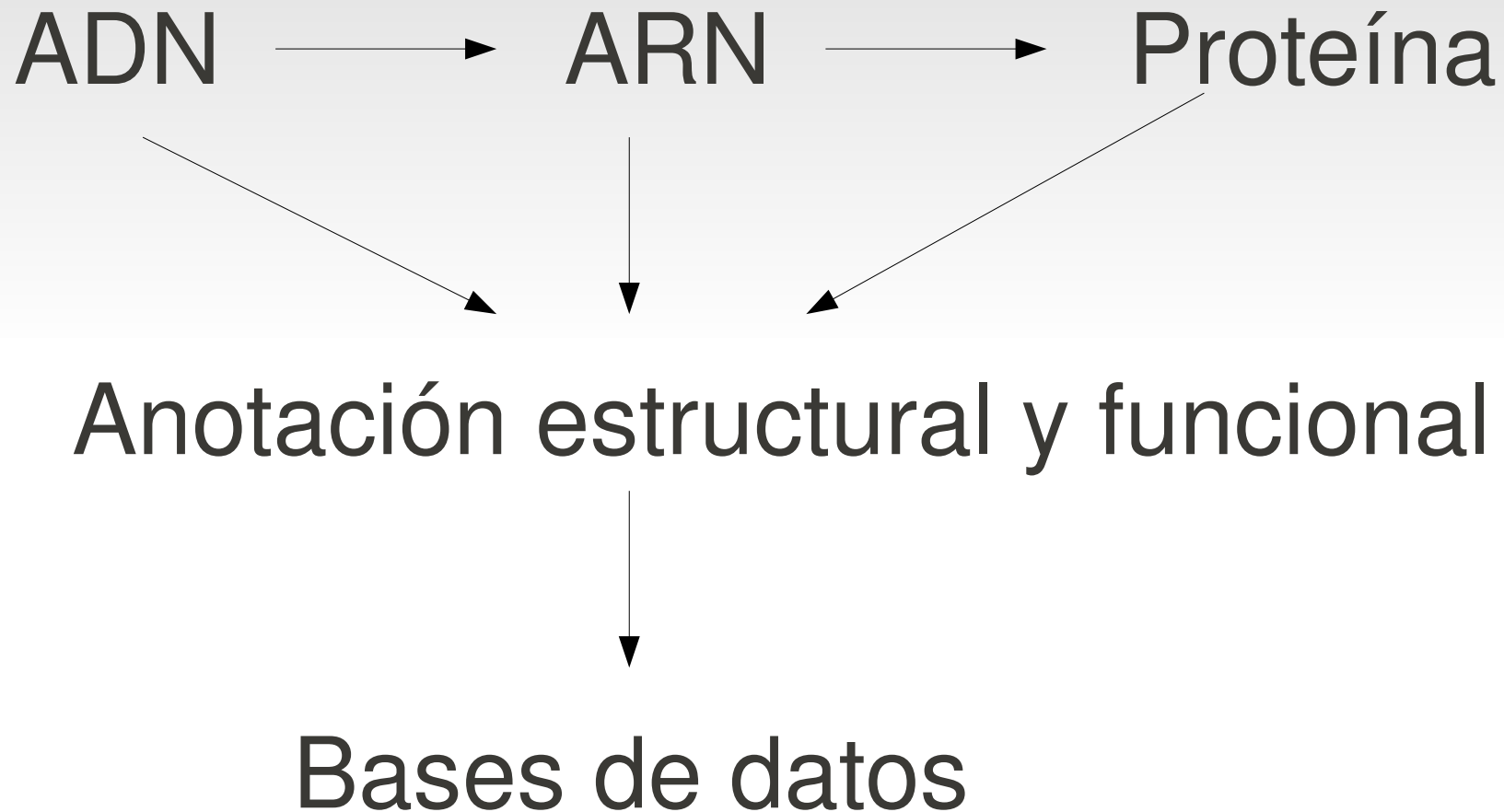
Anotación genómica

- Una vez se han encontrado los posibles elementos, se procede a la anotación funcional:
 - Definir posibles funciones
 - Comparar con bases de datos
 - Anotar con bases de datos (Gene Ontology)
- Ejemplo de grandes proyectos:
 - ENCODE (<http://www.genome.gov/10005107>)
 - ENSEMBL (<http://www.ensembl.org>)

Bases de datos

- Toda esta información se recopila en bases de datos:
 - GenBank (<http://www.ncbi.nlm.nih.gov/genbank/>)
 - dbSNP (<http://www.ncbi.nlm.nih.gov/projects/SNP/>)
 - Reactome (<http://www.reactome.org>)
 - PDB (<http://www.rcsb.org>)
 - Uniprot (<http://www.uniprot.org>)
 - Etc.

Resumen



Gracias por vuestra atención!!

¿Preguntas?