

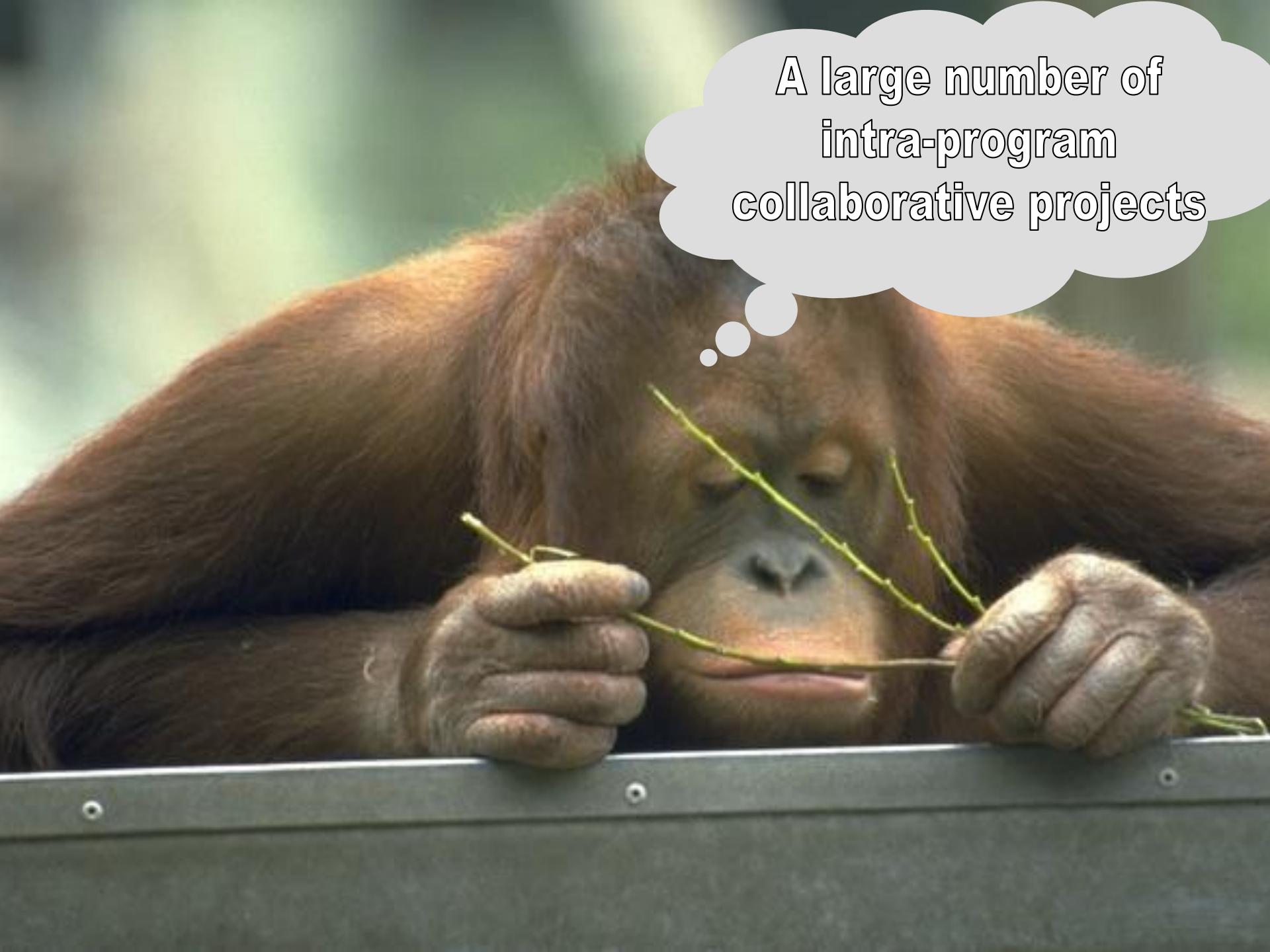
Life Sciences Workpackage

Consolider E-Science
Final meeting





More than 100 papers
20 organized meetings
>10 MEuros secure funding



A large number of
intra-program
collaborative projects

Whole-genome sequencing identifies mutations in chronic lymphocytic leukemia

Xose S. Puente¹, Magda Pinyol², Víctor Quesada¹, Laura Conde³, Pedro Jares³, Sílvia Bea³, Marcos González-Díaz⁵, Laia Bassagany³, Dolors Colomer³, José M. C. Tubío^{4,8}, Cristina López³, Alba Navarrete³, Jesús M. Hernández⁵, Diana A. Puente¹, José M. P. Freije⁴, Gloria Vicent³, Sara Guijarro³, Anna Enjuanes³, Lluís Hernández³, Jordi Yagüe⁷, Iñaki Himmelbauer¹⁰, Ester Castillo¹⁰, Juliane C. Dohm¹⁰, Silvia Jesú San Miguel⁹, Romina Royo¹³, Josep L. Gelpí¹³, David Torren¹³, Roderic Guigó¹⁵, Mónica Bayés¹⁶, Simon Heath¹⁶, Marta Gut¹⁶, Peter

nature
genetics

Exome sequencing identifies splicing factor SF3B1 gene mutations in chronic lymphocytic leukemia

Víctor Quesada¹, Laura Conde², Neus Villamor¹, Andrew J Ramsay¹, Sílvia Bea², Magda Pinyol⁴, Alba Navarro², Tycho Baumann⁵, Marta Ayme⁶, Jesús M Hernández⁶, Marcos González-Díaz⁶, José M C Tubío³, Romina Royo⁷, Josep L Gelpí⁸, Miguel Vázquez⁸, Alfonso Valencia⁸, Heinz Hünig⁹, Ivo Gut¹⁰, Xavier Estivill³, Armando López-Guerra¹¹

Here we perform whole-exome sequencing of sample 105 individuals with chronic lymphocytic leukemia (the most frequent leukemia in adults in Western coun-

OPEN  ACCESS Freely available online

Evidence for Transcript Networks Composed of Chimeric RNAs in Human Cells

Sarah Djebali^{1,9}, Julien Lagarde^{1,9}, Philipp Kapranov^{2,9,10}, Vincent Lacroix^{1,9,11}, Christelle Borel³, Jonathan M. Mudge⁴, Cédric Howald⁵, Sylvain Foissac^{1,2,12}, Catherine Ucla³, Jacqueline Chrast⁵, Paolo Ribeca¹, David Martin¹, Ryan R. Murray⁶, Xinpeng Yang⁶, Lila Ghamsari⁶, Chenwei Lin⁶, Ian Bell², Erica Dumais², Jorg Drenkow⁷, Michael L. Tress⁸, Josep Lluís Gelpí⁹, Modesto Orozco⁹, Alfonso Valencia⁸, Nynke L. van Berkum¹⁰, Bryan R. Lajoie¹⁰, Marc Vidal⁶, John Stamatoyannopoulos¹¹, Philippe Batut⁷, Alex Dobin⁷, Jennifer Harrow⁴, Tim Hubbard⁴, Job Dekker¹⁰, Adam Frankish⁴, Kourosh Salehi-Ashtiani^{6,12}, Alexandre Reymond⁵, Stylianos E. Antonarakis^{3,*}, Roderic Guigó^{1,13,*}, Thomas R. Gingeras^{2,7}

1 Bioinformatics and Genomics, Centre for Genomic Regulation and Universitat Pompeu Fabra, Barcelona, Catalonia, Spain, **2** Affymetrix Inc., Santa Clara, California, United States of America, **3** Department of Genetic Medicine and Development, University of Geneva Medical School, University Hospitals of Geneva, Geneva, Switzerland, **4** Wellcome Trust Sanger Institute, Hinxton, Cambridgeshire, United Kingdom, **5** The Center for Integrative Genomics, University of Lausanne, Lausanne, Switzerland,

LETTERS

Epigenomic analysis detects widespread gene-body DNA hypomethylation in chronic lymphocytic leukemia

na C Queirós^{4,11}, Alba Navarro¹, Guillem Clot¹, Ille Brun-Heath², Magda Pinyol⁶, Sergio Barberán-Soler⁷, Rico⁸, Simone Ecker⁸, Miriam Rubio⁸, Romina Royo⁹, Conde¹, Mónica López-Guerra¹, Dolors Colomer¹, Sílvia Bayés², Marta Gut², Josep L Gelpí⁹, Xose S Puente¹⁰, David G Pisano⁸, Alfonso Valencia⁸, Hünig¹⁰, Elías Campo¹ & José I Martín-Subero⁴

DNA methylation is a major mechanism in cell differentiation and neoplastic transformation^{4–7}. The epigenetic modifications contributing to CLL development, however, are not fully understood^{8–10}. Here

 PLOS ONE

Resource

GENCODE: The reference human genome annotation for The ENCODE Project

Jennifer Harrow,^{1,9} Adam Frankish,¹ Jose M. Gonzalez,¹ Electra Tapanari,¹ Mark Diekhans,² Felix Kokocinski,¹ Bronwen L. Aken,¹ Daniel Barrell,¹ Amonida Zadissa,¹ Stephen Searle,¹ If Barnes,¹ Alexandra Bignell,¹ Veronika Boychenko,¹ Toby Hunt,¹ Mike Kay,¹ Gaurab Mukherjee,¹ Jeena Rajan,¹ Gloria Despacio-Reyes,¹ Gary Saunders,¹ Charles Steward,¹ Rachel Harte,² Michael Lin,³ Cédric Houwald,⁴ Andre

ARTICLE

An integrated encyclopedia of DNA elements in the human genome

The ENCODE Project Consortium*

The human genome encodes the blueprint of life, but the function of the vast majority of its nearly three billion bases is unknown. The Encyclopedia of DNA Elements (ENCODE) project has systematically mapped regions of transcription, transcription factor association, chromatin structure and histone modification. These data enabled us to assign biochemical functions for 80% of the genome, in particular outside of the well-studied protein-coding regions. Many discovered candidate regulatory elements are physically associated with one another and with expressed genes, providing new insights into the mechanisms of gene regulation. The newly identified elements also show a statistical correspondence to sequence variants linked to human disease, and can thereby guide interpretation of this variation. Overall, the project provides new insights into the organization and regulation of our genes and genome, and is an expansive resource of functional annotations for biomedical research.

The human genome sequence provides the underlying code for human biology. Despite intensive study, especially in identifying protein-coding genes, our understanding of the genome is far from complete, particularly with



ENCODE
Encyclopedia of DNA Elements
nature.com/encode

Research

Chimeras taking shape: Potential functions of proteins encoded by chimeric RNA transcripts

Milana Frenkel-Morgenstern,¹ Vincent Lacroix,² Iakes Ezkurdia,¹ Yishai Levin,³ Alexandra Gabashvili,³ Jaime Prilusky,⁴ Angela del Pozo,¹ Michael Tress,¹ Rory Johnson,⁵ Roderic Guigo,⁵ and Alfonso Valencia^{1,6}

Science AAAS.org | FEEDBACK | HELP | LIBRARIANS | FOR EDITORIAL JOURNALISTS | CAREER OPPORTUNITIES | ALERTS | ACCESS RIGHTS | MY ACCOUNT | SIGN IN

Science The World's Leading Journal of Original Scientific Research, Global News, and Commentary

Science Home Current Issue Previous Issues Science Express Science Products My Science About the Journal

Home > Science Magazine > 7 September 2012 | Perissi, 337 (6099) 1158-1161

Article Views

- Summary
- Full Text
- Full Text (PDF)

NEWS & ANALYSIS

GENOMICS

ENCODE Project Writes Eulogy for Junk DNA

Elizabeth Pennisi

This week, 30 research papers, including six in *Nature* and additional papers published online by *Science*, sound the death knell for the idea that our DNA is mostly littered with useless bases. A decade-long project, the Encyclopedia of DNA Elements (ENCODE), has found that 80% of the human genome serves some purpose, biochemically speaking. Beyond defining proteins, the DNA bases highlighted by ENCODE specify landing spots for proteins that influence gene activity, strands of RNA with myriad roles, or simply places where chemical modifications serve to silence stretches of our chromosomes.

Read Full Text

Genome Research
www.genome.org

molecular
systems
biology

artners

Alfonso Valencia^{1,*}

NIO), Madrid, Spain, ² Structural Bioinformatics Group, Centre for ecomputing Center, Barcelona, Spain, ⁴ Computational systems Biology Group, National Centre for Biotechnology

ore
Research Centre (CNIO), C/Melchor Fernández Almagro 3,
ail: valencia@cnio.es

given proteome ('interactome') is the ms Biology. Separately the prediction is a well-established scientific area. action partners. We provide a proof of *as representing known interactors in*

95% of the genome lies within 8 kilobases (kb) of a DNA–protein interaction (as assayed by bound ChIP-seq motifs or DNase I footprints), and 99% is within 1.7 kb of at least one of the biochemical events measured by ENCODE.



,... and a large number of
inter-program
collaborative projects



Cell-Dock: high-performance protein–protein docking

[« Previous | Next Article »](#)
[Table of Contents](#)

Search this journal

Carles Pons^{1,2}, Daniel Jiménez-González^{3,4,*}, Cecilia C. Sánchez-Orive¹,
Harald Servat^{3,4}, Daniel Cabrera-Benítez⁴, Xavier Agulló¹,
Juan Fernández-Recio^{1,*}

Author Affiliations

*To whom correspondence should be addressed.

Published online 9 May 2011

Nucleic Acids Research

T-Coffee: a web server for alignment of protein and structural information and its applications

Paolo Di Tommaso¹, Sébastien Moretti^{2,3},
Alberto Montanyola⁴, Jia-Ming Chang¹, Jean-Pierre

¹Centre For Genomic Regulation (Pompeu Fabra University), Barcelona, Spain, ²Vital-IT, Swiss Institute of Bioinformatics, Quartier de l'Informatique, CH-1211 Geneva 4, Switzerland, ³Department of Ecology and Evolution, Biology Department, University of Edinburgh, Edinburgh, UK, ⁴Department of Computer Science and Engineering, Campus de Cappont, C. de Jaume II 69, E-25001 Lleida, Spain

BIOINFORMATICS APPLICATIONS NOTE

Vol. 26 no. 15 2010, pages 1903–1904
doi:10.1093/bioinformatics/btq304

Sequence analysis

Cloud-Coffee: implementation of a parallel consistency-based multiple alignment algorithm in the T-Coffee package and its benchmarking on the Amazon Elastic-Cloud

Paolo Di Tommaso¹, Miquel Orobioig², Fernando Guirado², Fernando Cores², Toni Espinosa³ and Cedric Notredame^{1,*}

Software

Highly accessed

Open Access

A user-friendly web portal for T-Coffee on supercomputers

Josep Rius^{1,*}, Fernando Cores¹, Francesc Solsona¹, Jano I van Hemert², Jos Koetsier² and Cedric Notredame³

* Corresponding author: Josep Rius jrius@diei.udl.cat

Author Affiliations

¹ Department of Computer Science and Industrial Engineering, University of Lleida, C/Jaume II 69, E-25001 Lleida, Spain

² UK National e-Science Centre, University of Edinburgh, 10 Crichton Street, EH8 9AB, Edinburgh, UK

³ Centre For Genomic Regulation (Pompeu Fabra University), C/Doctor Aiguader 88, 08003 Barcelona, Spain

For all author emails, please [log on](#).

Barcelona,
pont, C. de
niversitat

t me t

urnal

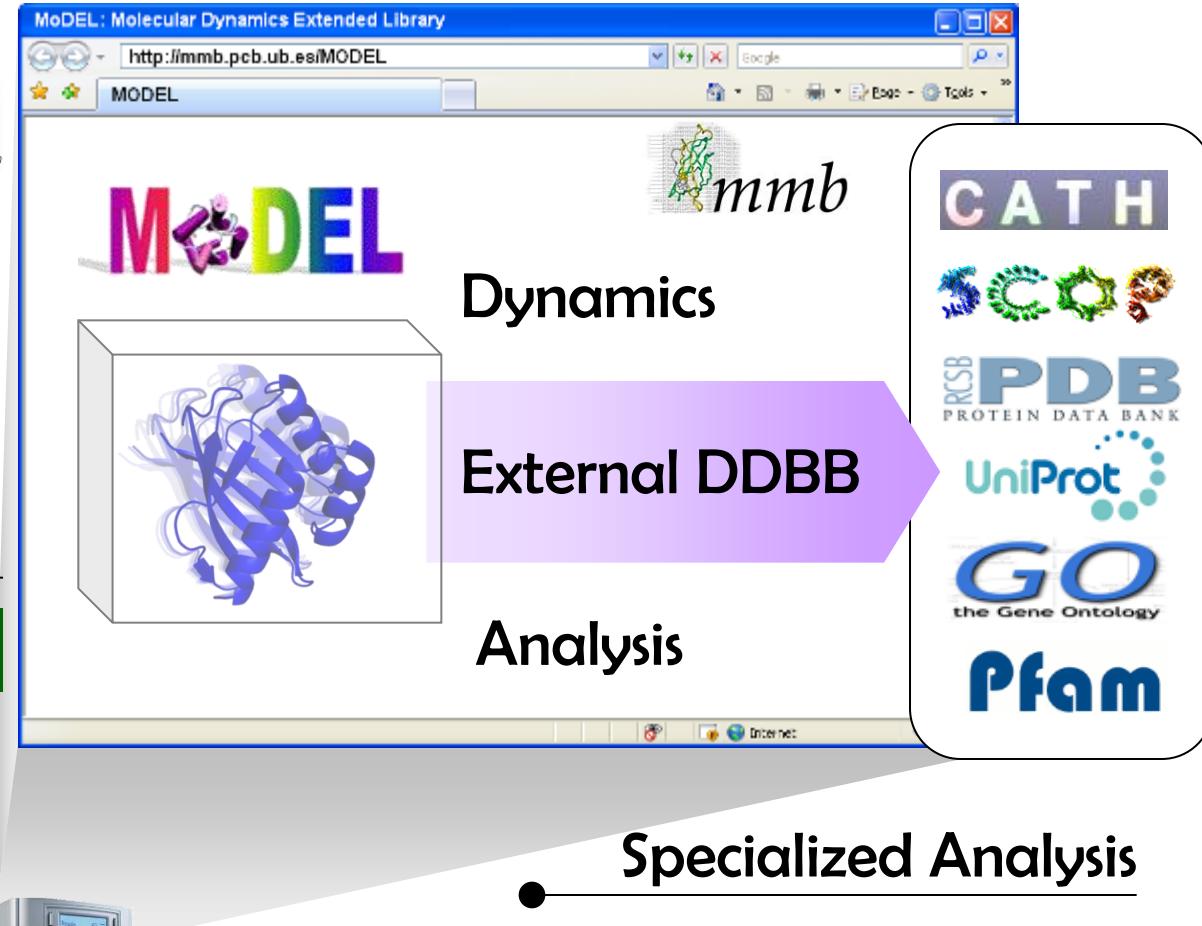
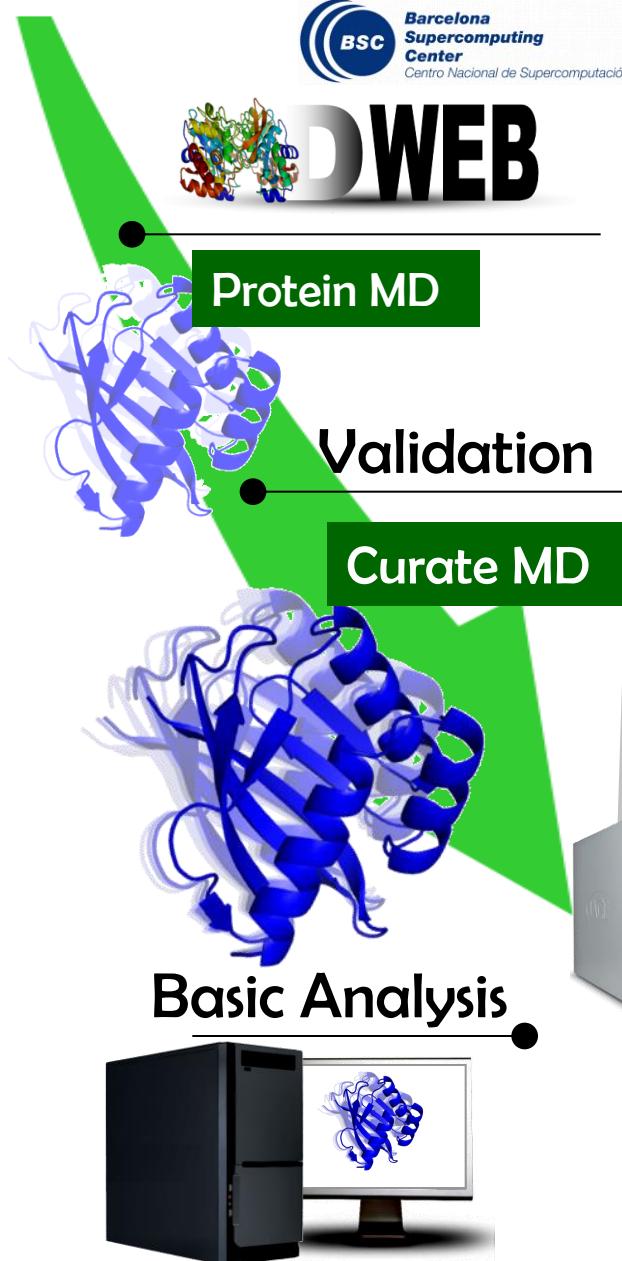
is jour
r of ev

Perm
i date
urnal is

Main research lines

- Molecular Simulations
- Alternative Splicing
- Genomics from HiSeq
- Genetic basis of complex diseases
- *Software development*

Protein Structure



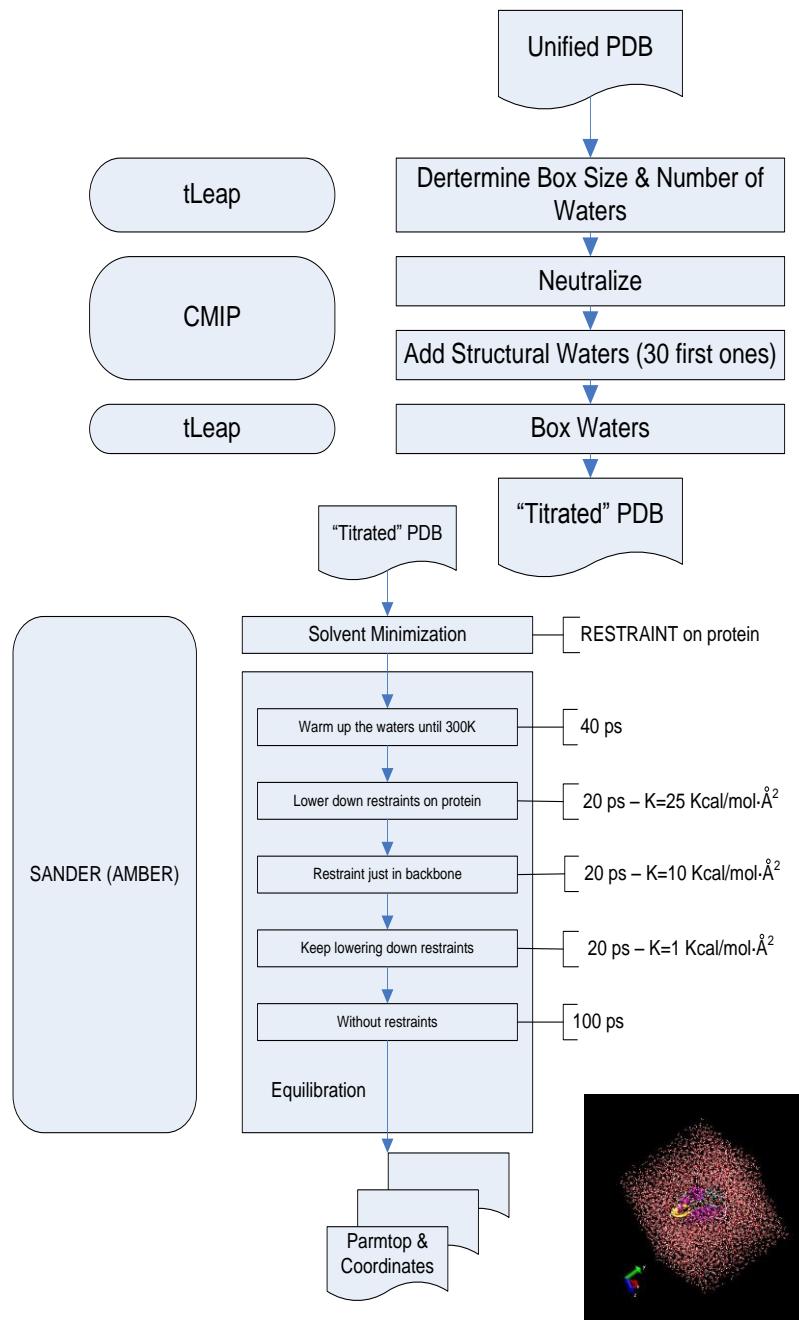
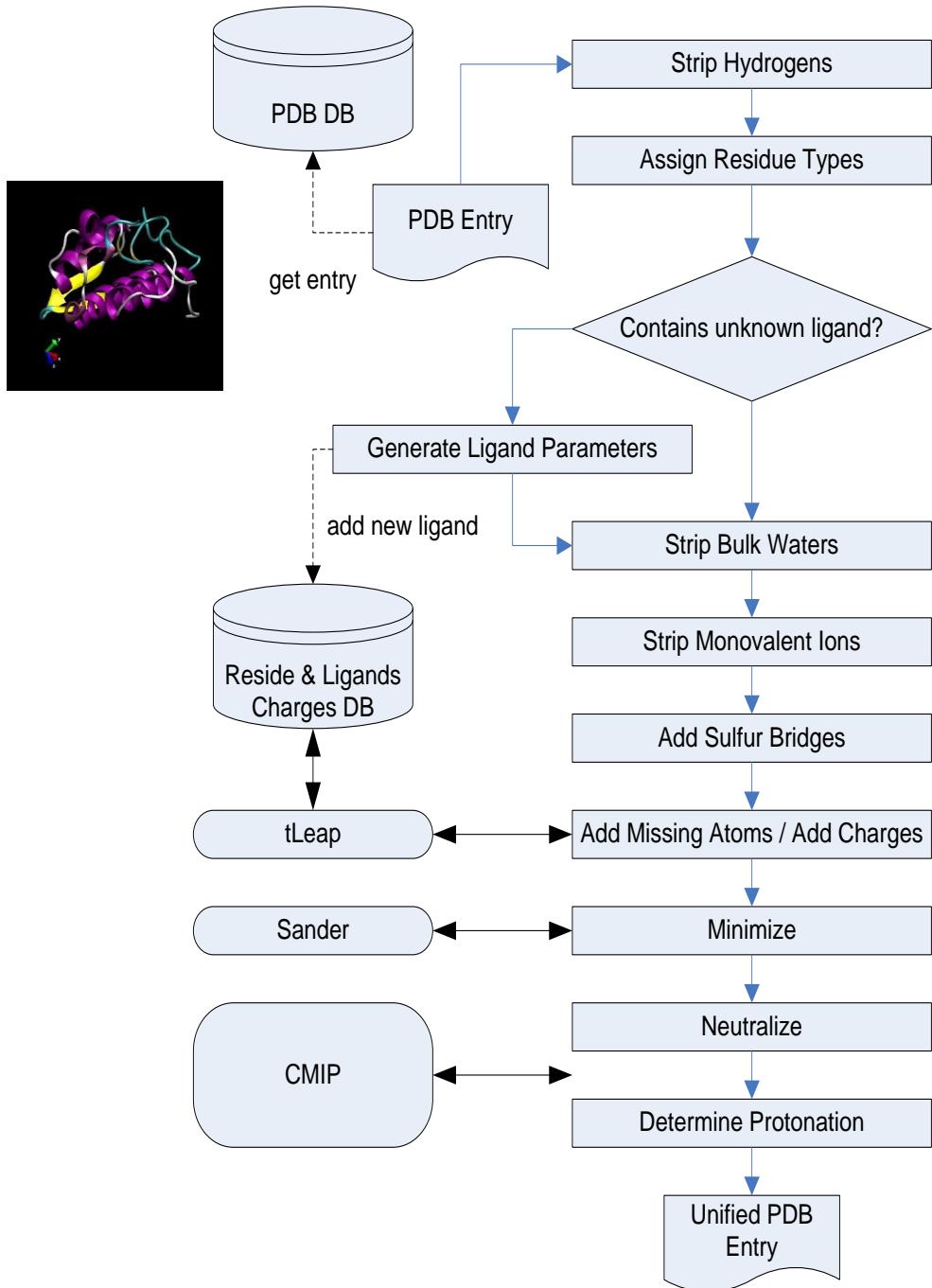
Specialized Analysis

Coarse-Grained

MySQL
Server

Flex
Serv

INB®



NAFlex: A web server for the study of Nucleic Acids Flexibility.



All-Atom
Molecular Dynamics

Powered by
NAMD
Parallel Molecular Dynamics

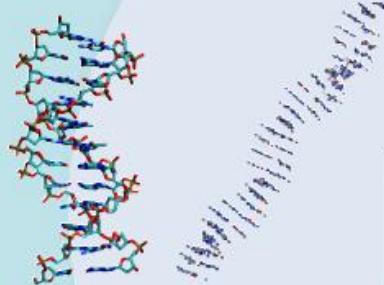
AmberTools

GROMACS FAST.
FLEXIBLE.
FREE.

Input Structure:

PDB
PROTEIN DATA BANK

Dynamics Simulations

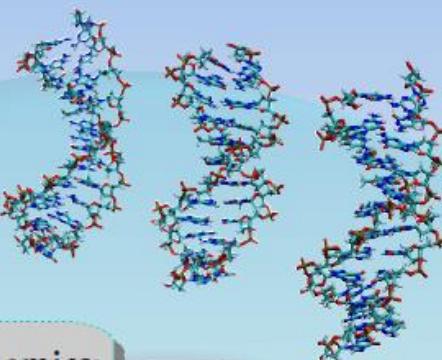


3D Structure Generation

B-DNA, A-DNA, A-RNA
User defined Nucleic Acids
NA From Helical Parameters
Coarse-Grained Models

Input Trajectory:
crd, dcd, netcdf, xtc, binpos

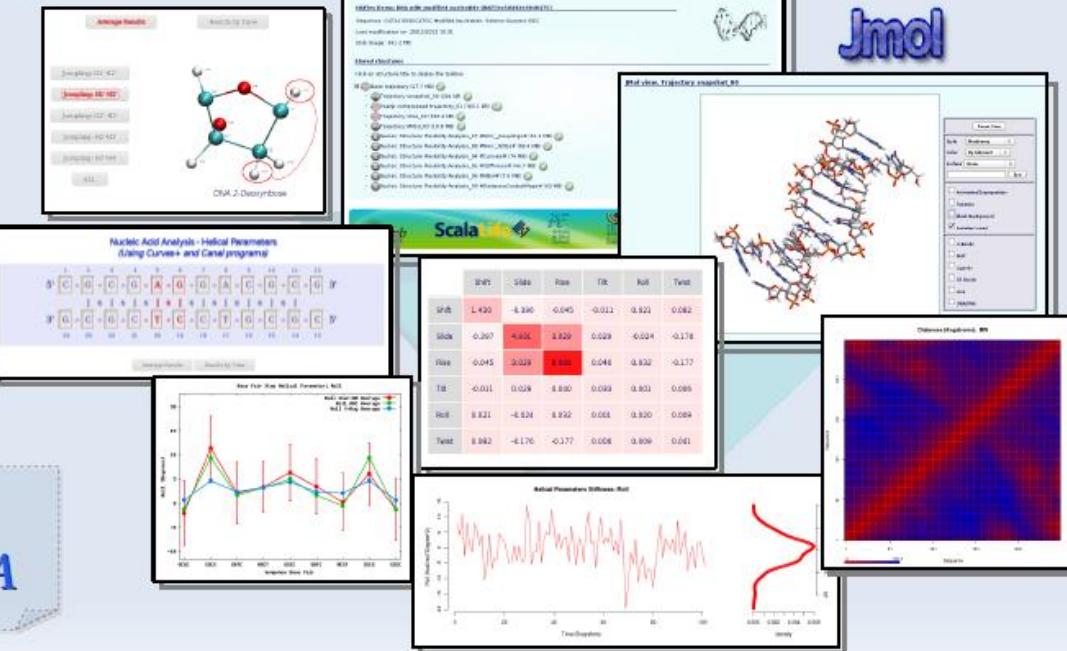
Coarse-Grained Dynamics:
Mesoscopic Elastic Model
Worm-Like Chain Model



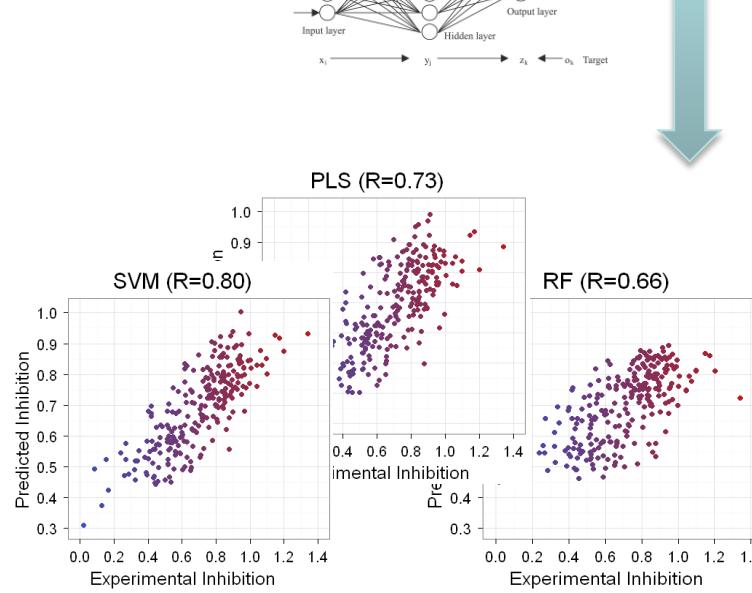
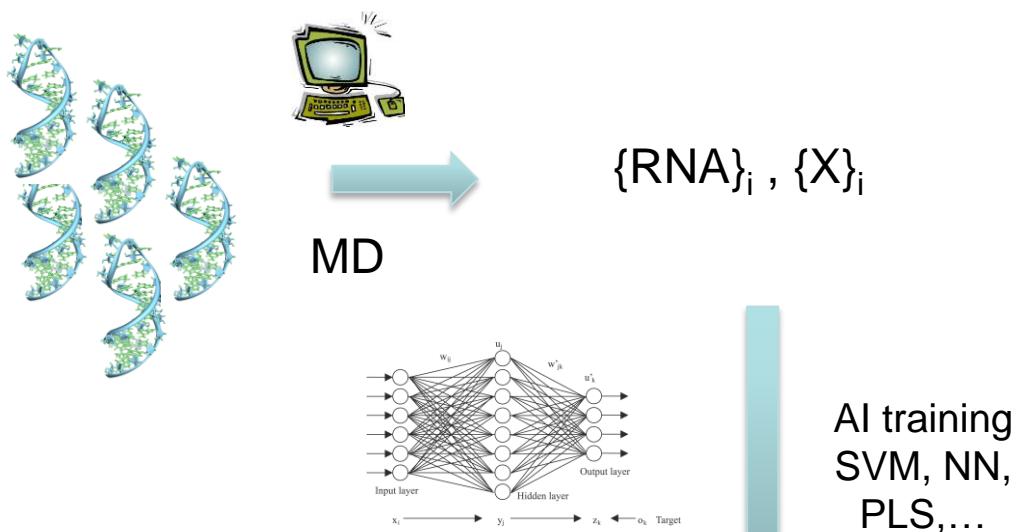
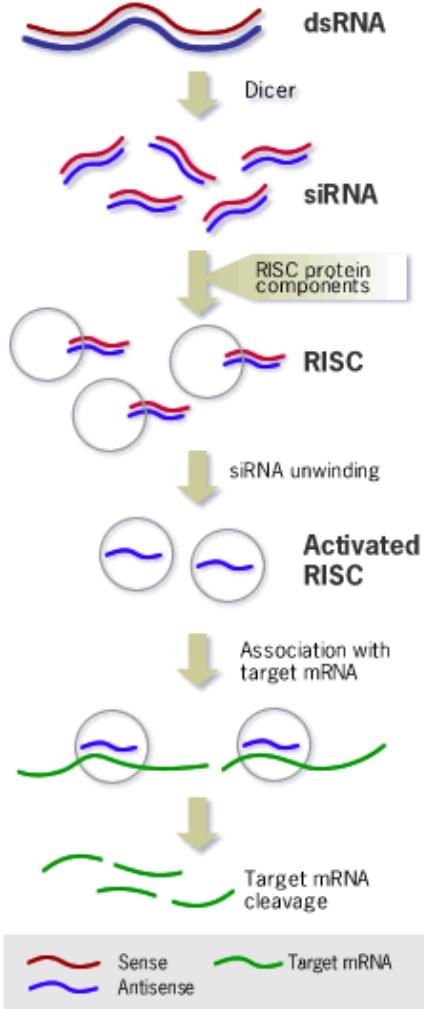
Analysis & Visualization

Cartesian Analysis
Helical Description
Principal Component Analysis
Stiffness Analysis
Energy Analysis
NMR Descriptors

Input Sequence:
GATTACCA

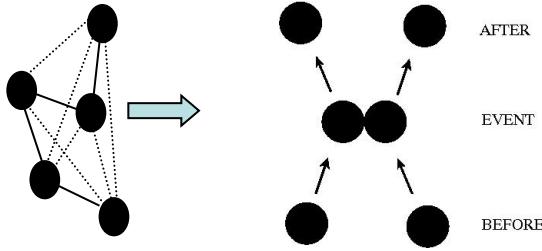


Predicting SiRNA efficiency



Validation

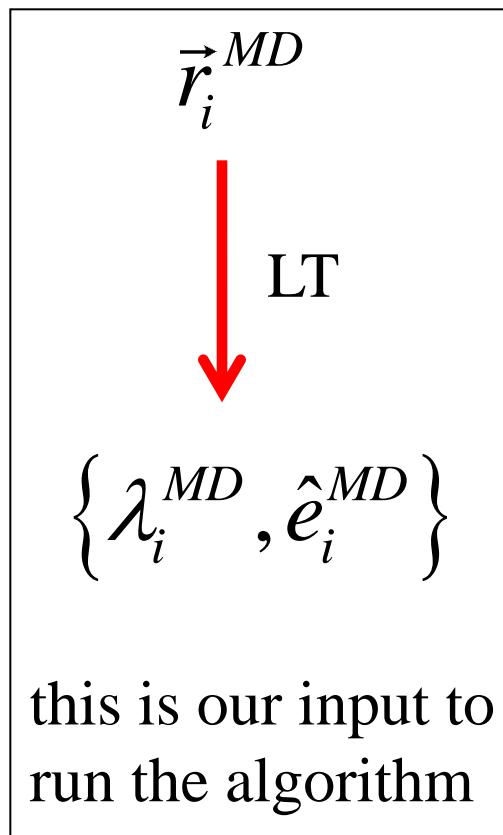




$$\vec{r}_i(t + t_c) = \vec{r}_i(t) + \vec{v}_i(t)t_c$$

Collision: $m_j \vec{v}_j + \Delta \vec{p} = m_j \vec{v}_j'$ $m_i \vec{v}_i = m_i \vec{v}_i' + \Delta \vec{p}$

$$t_{ij} = \frac{-b_{ij} \pm \sqrt{b_{ij}^2 - v_{ij}^2(r_{ij}^2 - d^2)}}{v_{ij}^2} \quad \Delta p = \frac{2m_i m_j}{m_i + m_j} (v_i - v_j)$$

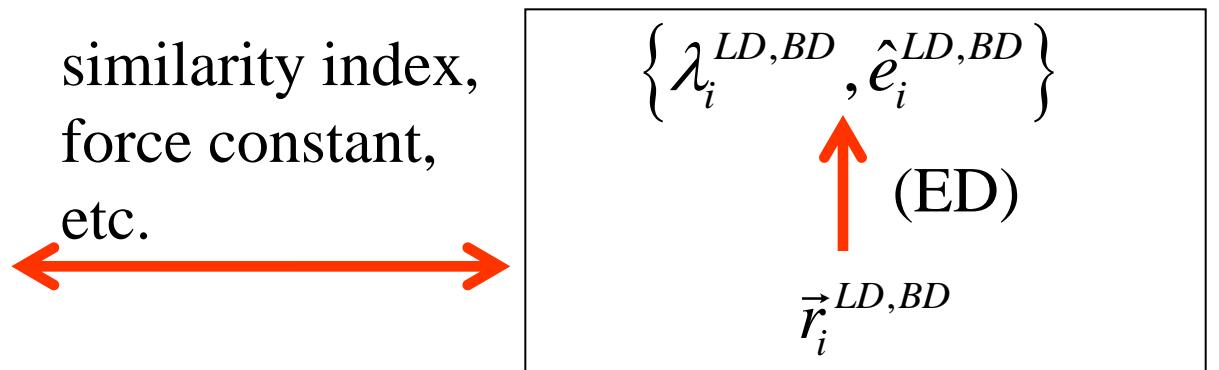


similarity index,
force constant,
etc.

initial condition: all atoms in their equilibrium positions

$$r_i = r_i^0 + \tau \left(1 - e^{-\Delta t/\tau}\right) v_i^0 + \frac{\Delta t}{\tau} \left(1 - \frac{\tau}{\Delta t} \left(1 - e^{-\Delta t/\tau}\right)\right) F_i^0 + \Delta r_i^G$$

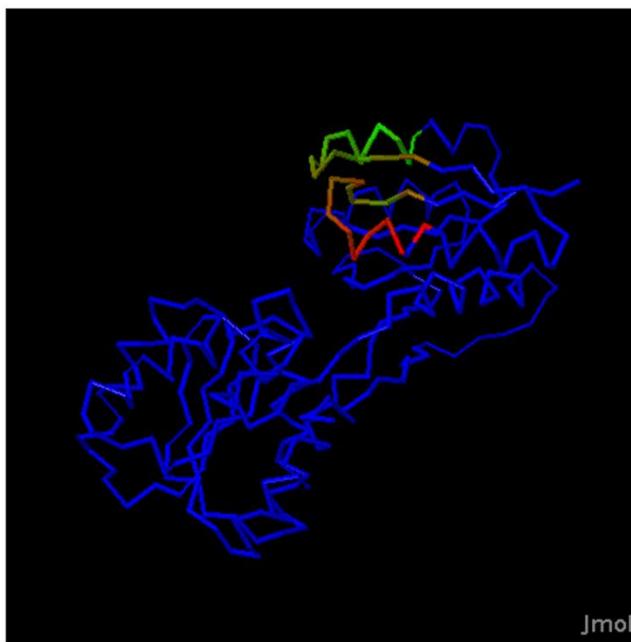
$$v_i = e^{-\Delta t/\tau} v_i^0 + \frac{1}{\gamma} \left(1 - e^{-\Delta t/\tau}\right) F_i^0 + \Delta v_i^G$$



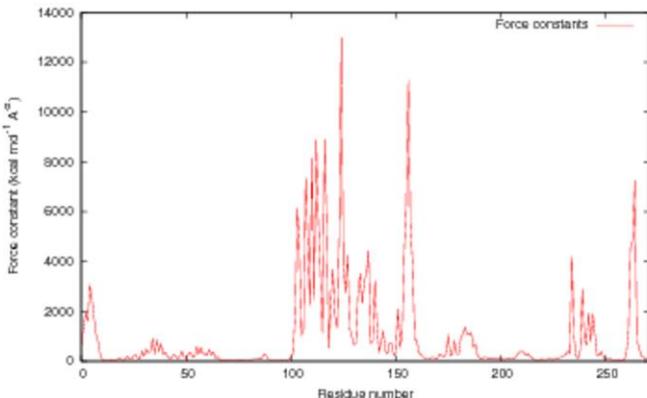
Speed-up by a factor of 1000

Chained correlations

Residue: 50 Search width: 3 Search depth: 10 Threshold: 0.97 RMS type: Gaussian



Hinge prediction

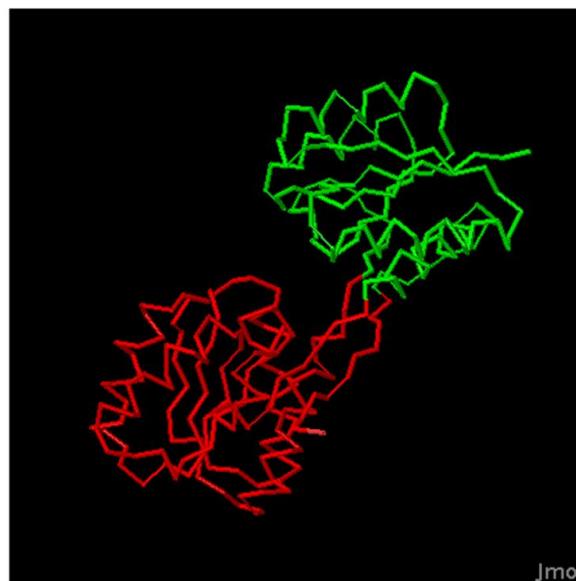


The peak constant corresponds to residue 125 with value 13002.110352

Residues in the top 20% are: 125, 157

Dynamic domain detection method

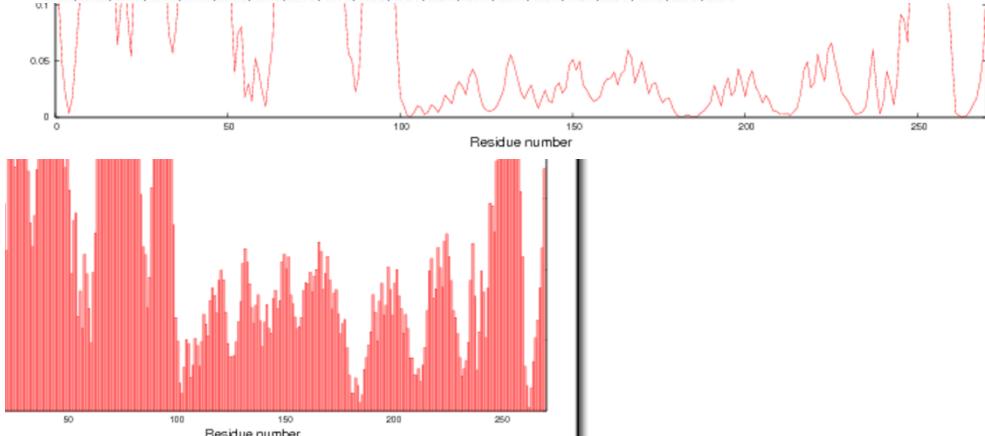
Movement by normal mode 1. Gaussian RMS applied

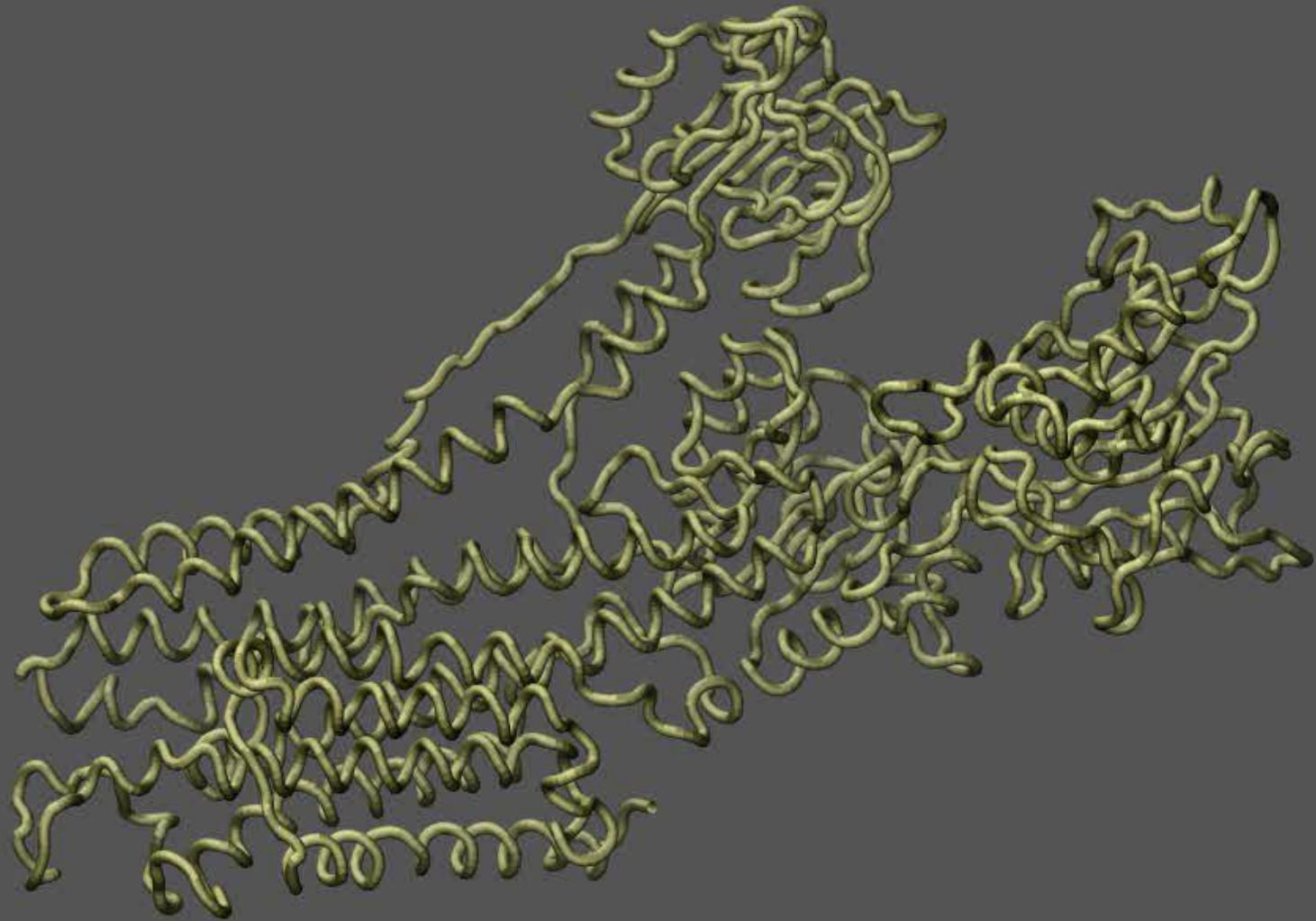


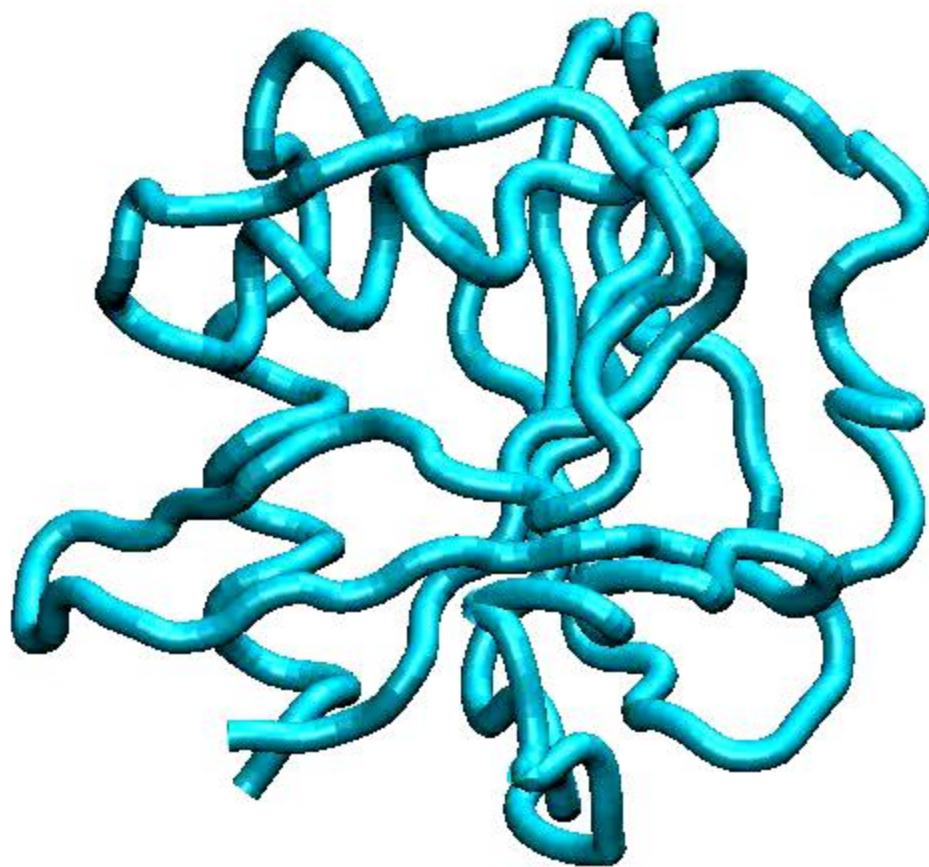
Residues in each cluster:

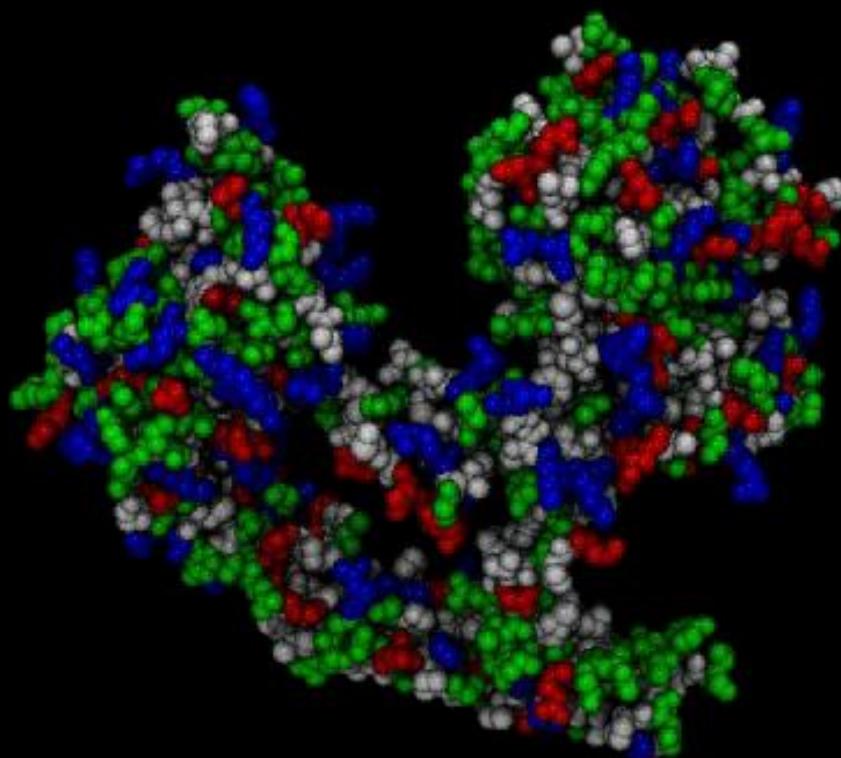
Cluster 1: 102, 103, 104, 105, 106, 107, 108, 109, 110, 111, 112, 113, 114, 115, 116, 117, 118, 119, 120, 121, 122, 123, 124, 125, 126, 127, 128, 129, 130, 131, 132, 133, 134, 135, 136, 137, 138, 139, 140, 141, 142, 143, 144, 145, 146, 147, 148, 149, 150, 151, 152, 153, 154, 155, 156, 157, 158, 159, 160, 161, 162, 163, 164, 165, 166, 167, 168, 169, 170, 171, 172, 173, 174, 175, 176, 177, 178, 179, 180, 181, 182, 183, 184, 185, 186, 187, 188, 189, 190, 191, 192, 193, 194, 195, 196, 197, 198, 199, 200, 201, 202, 203, 204, 205, 206, 207, 208, 209, 210, 211, 212, 213, 214, 215, 216, 217, 218, 219, 220, 221, 222, 223, 224, 225, 226, 227, 228, 229, 230, 231, 232, 233, 234, 235, 236, 261, 262, 263, 264, 265, 266, 267, 268, 269, 270

Cluster 2: 0, 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34, 35, 36, 37, 38, 39, 40, 41, 42, 43, 44, 45, 46, 47, 48, 49, 50, 51, 52, 53, 54, 55, 56, 57, 58, 59, 60, 61, 62, 63, 64, 65, 66, 67, 68, 69, 70, 71, 72, 73, 74, 75, 76, 77, 78, 79, 80, 81, 82, 83, 84, 85, 86, 87, 88, 89, 90, 91, 92, 93, 94, 95, 96, 97, 98, 99, 100, 101, 237, 238, 239, 240, 241, 242, 243, 244, 245, 246, 247, 248, 249, 250, 251, 252, 253, 254, 255, 256, 257, 258, 259, 260

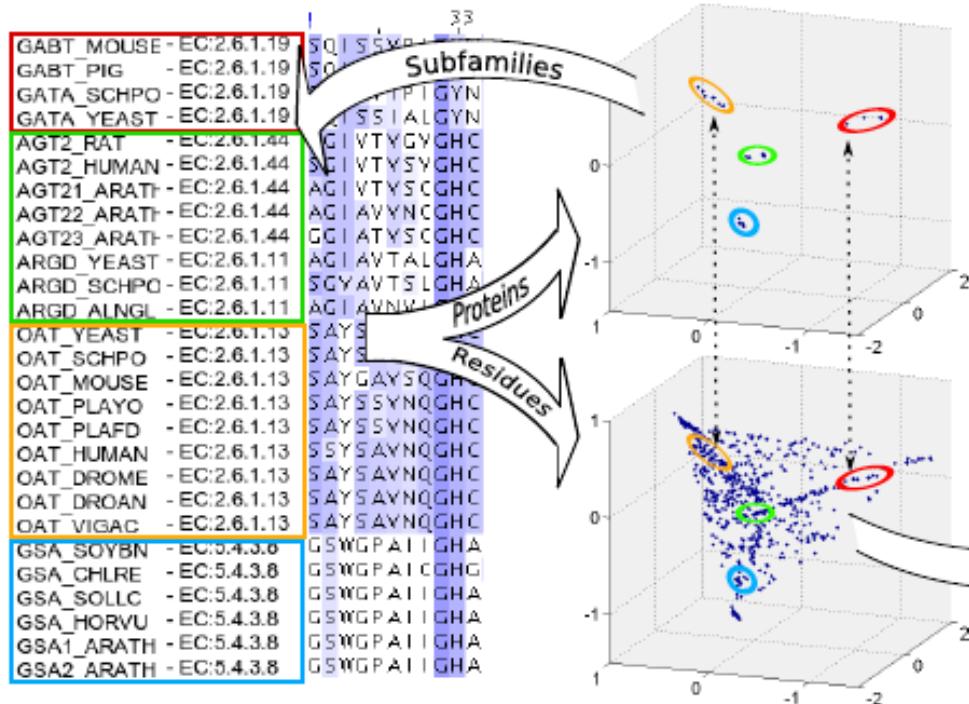




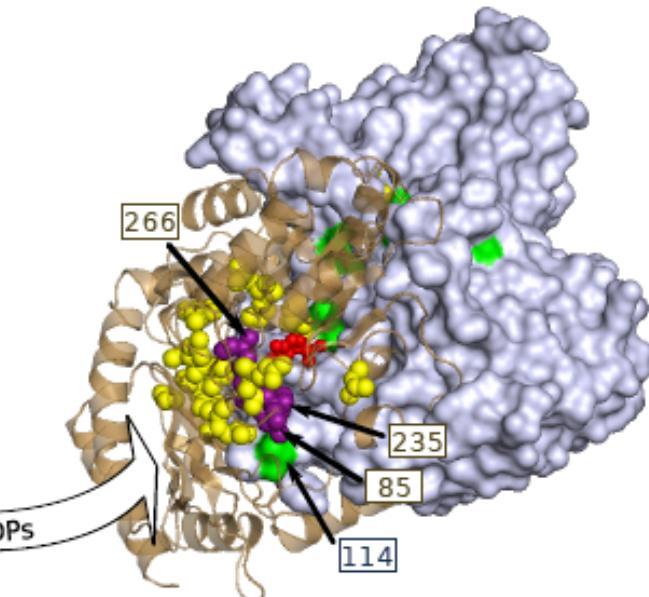




Prediction of protein interaction sites from protein family sequence analysis



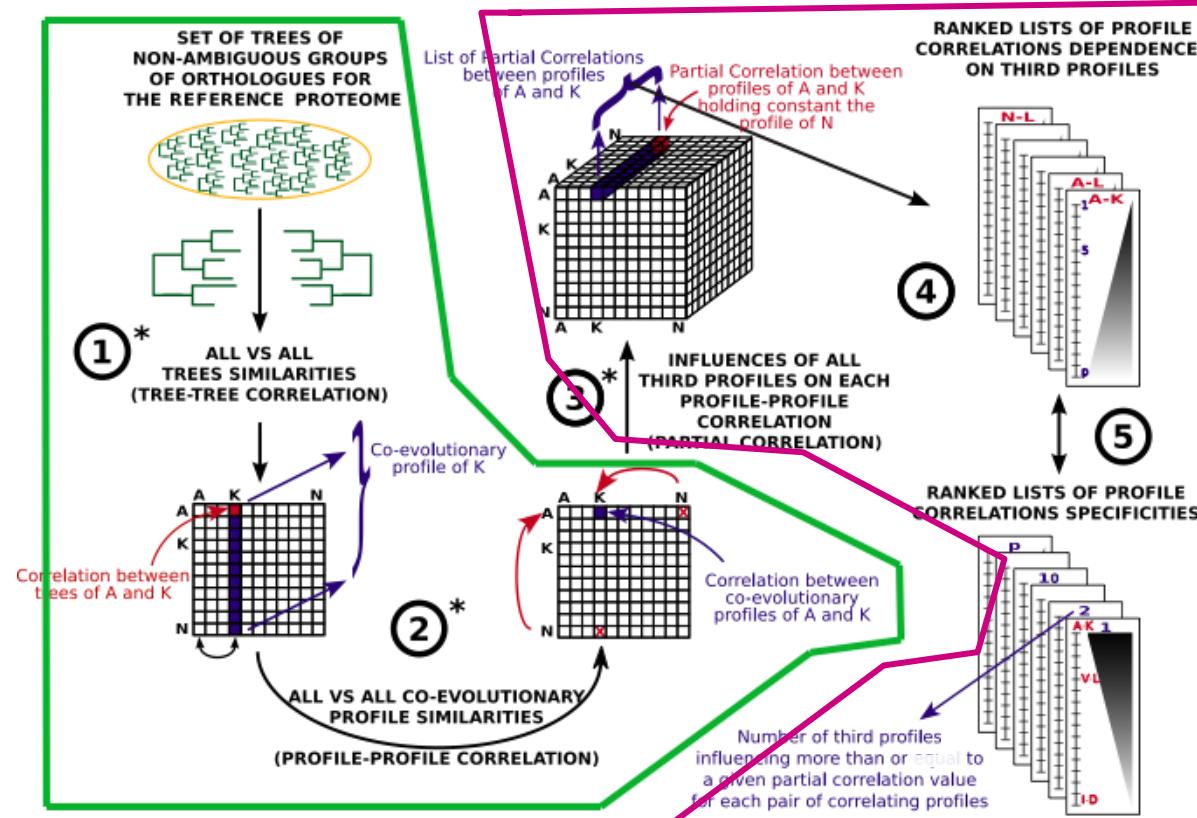
Correspondence analysis of multiple sequence alignments to identify sequence groups and associated significant residues (SDPs)



human ornithine aminotransferase (PDB 1oat) bound to Pyridoxal-5'-phosphate (red spheres). The two subunits of the complex are shown in brown cartoons and grey surface. Key predicted residues (SDPs) in yellow/violet spacefill and green surface respectively.

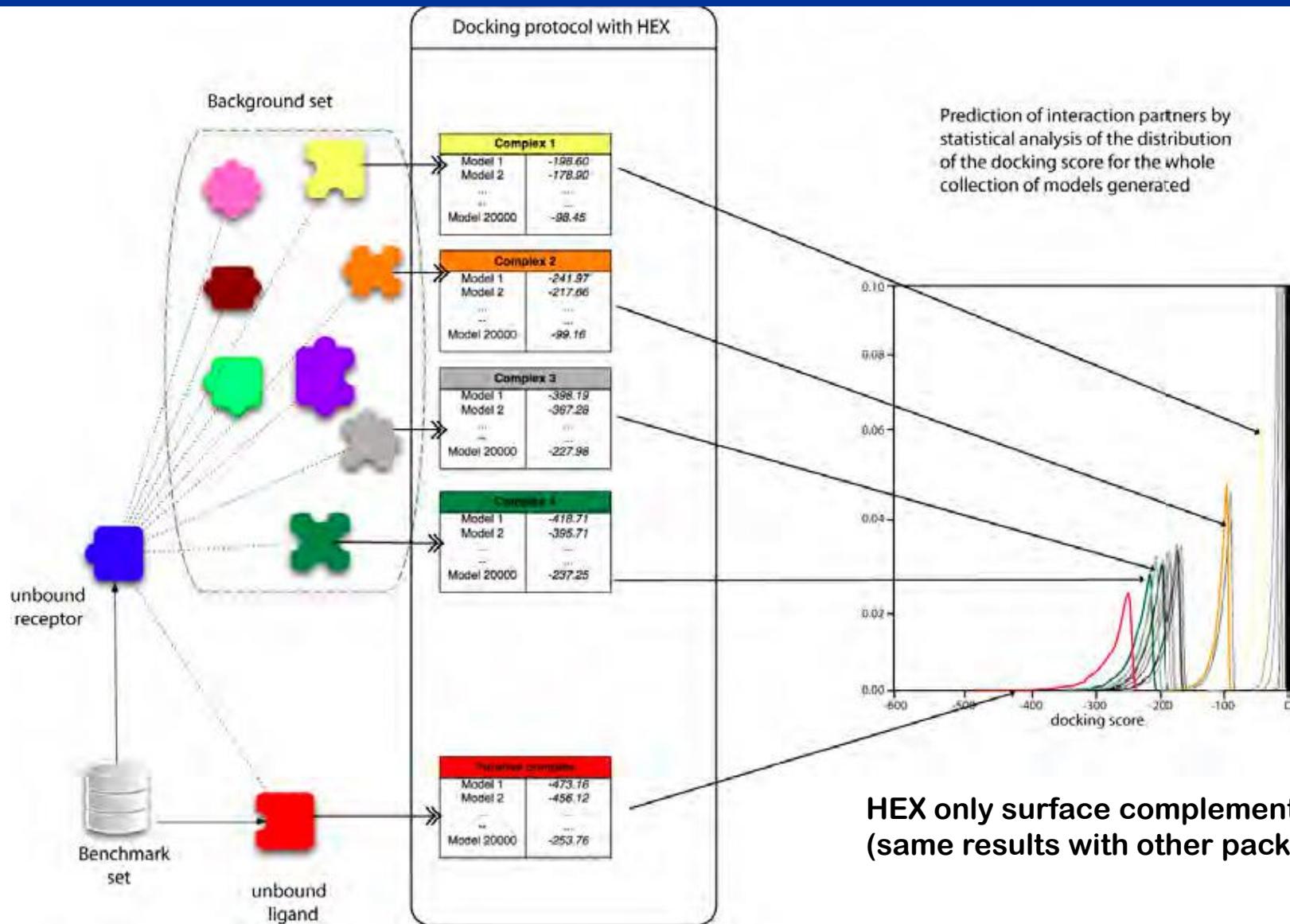
Predicting interaction networks

Detecting Evolutionary Similarities



Contrast function

HT docking: all against all protein docking. 56 unbound proteins against a benchmark of 633 proteins producing 20,000 models for each one of them (total of more than 700 M models)



Alternative Splicing

In allowing for the generation of a diverse range of mature RNAs from the same gene, alternative splicing allows has considerable theoretical potential to expand the range of cellular proteins.

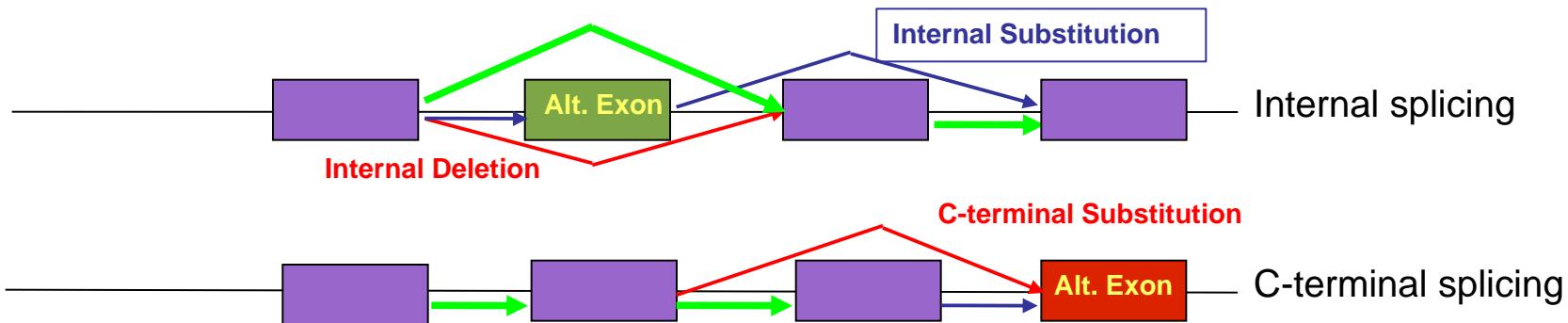
Recent studies estimate that **40–80%** of multi-exon human genes can produce differently spliced mRNAs.

There are many studies that implicate alternative transcripts in biological processes such as development.

This has lead to the hypothesis that **alternative splicing can explain the apparent lack of correlation between organism complexity and numbers of genes.**

*In the past few years, it has become clear that a phenomenon called alternative splicing is one reason human genomes can produce such complexity with so few genes.**

*Science, July 2005

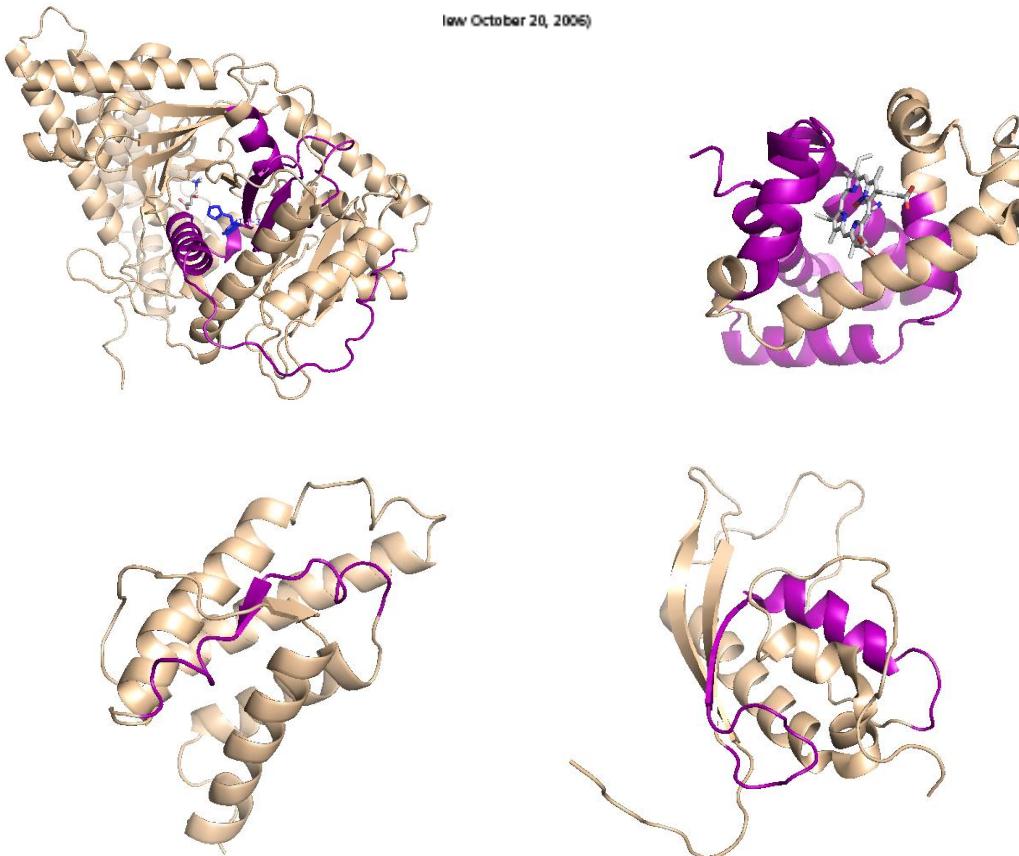
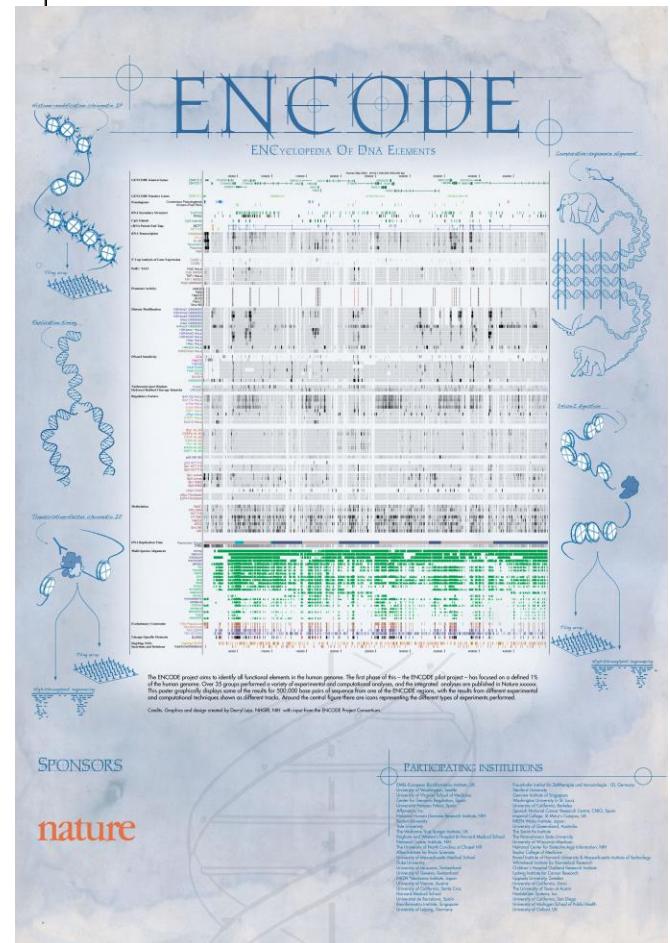


The implications of alternative splicing in the ENCODE protein complement

Michael L. Tress^{a,b}, Pier Luigi Martelli^c, Adam Frankish^d, Gabrielle A. Reeves^e, Jan Jaap Wesselink^a, Corin Yeats^f, Páll Ísólfur Olason^g, Mario Albrecht^h, Heidi Hegyiⁱ, Alejandro Giorgetti^j, Domenico Raimondi^k, Julien Lagarde^k, Roman A. Laskowski^l, Gonzalo López^m, Michael I. Sadowskiⁿ, James D. Watson^o, Piero Fariselli^p, Ivan Ross^q, Alinda Nagy^r, Wang Kal^s, Zenia Storling^t, Massimiliano Orsini^u, Yassen Assenov^v, Hagen Blankenburg^w, Carola Huthmacher^x, Fidel Ramirez^y, Andreas Schlicker^z, France Denoué^z, Phil Jones^z, Samuel Kerrien^z, Sandra Orchard^z, Stylianos E. Antonarakis^z, Alexandre Reymond^z, Ewan Birney^z, Søren Brunak^z, Rita Casadio^z, Roderic Guigó^z, Jennifer Harrow^z, Henning Hermjakob^z, David T. Jones^z, Thomas Lengauer^z, Christine A. Orengo^z, László Patthy^z, Janet M. Thornton^z, Anna Tramontano^z, and Alfonso Valencia^{a,z}

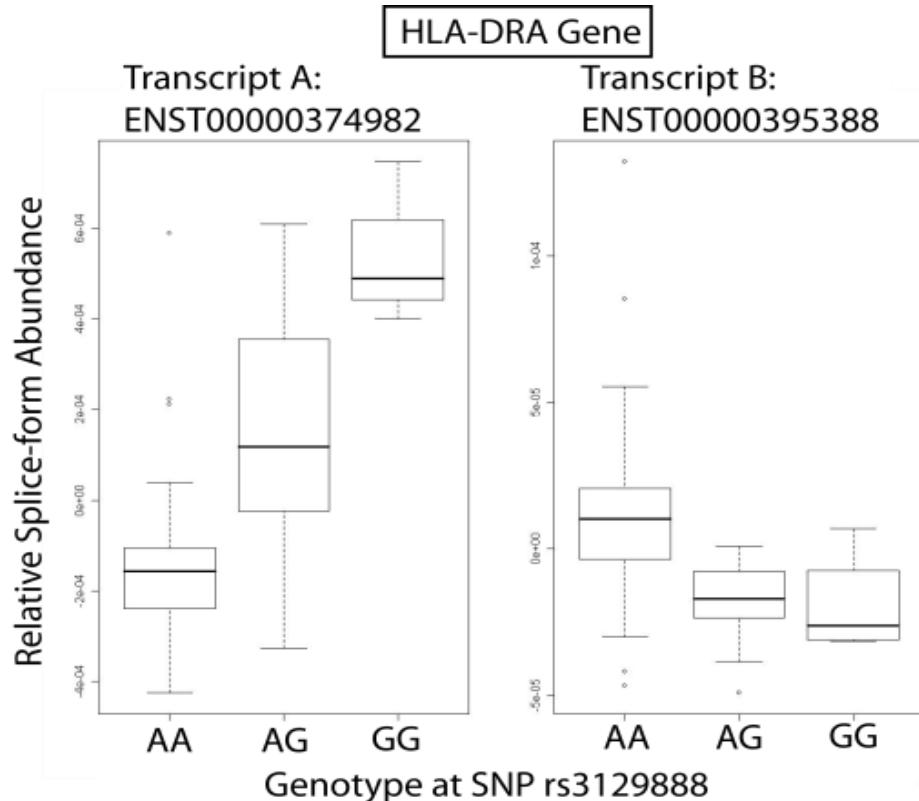
^aStructural Computational Biology Programme, Spanish National Cancer Research Centre, E-28029 Madrid, Spain; ^bDepartment of Biology, University of Bologna, 33-40126 Bologna, Italy; ^cHAVANA Group, The Wellcome Trust Sanger Institute, Hinxton, Cambridge CB10 1SA, United Kingdom; ^dEuropean Bioinformatics Institute, Hinxton, Cambridge CB10 1SD, United Kingdom; ^eDepartment of Biochemistry and Molecular Biology and Bioinformatics Unit, University College London, London WC1E 6BT, United Kingdom; ^fCenter for Biological Sequence Analysis, BioCentrum-DTU, DK-2800 Lyngby, Denmark; ^gMax Planck Institute for Informatics, 66123 Saarbrücken, Germany; ^hBiological Research Center, Hungarian Academy of Sciences, 1113 Budapest, Hungary; ⁱDepartment of Biochemical Sciences, University of Rome "La Sapienza," 00185 Rome, Italy; ^jResearch Unit on Biomedical Informatics, Institut Municipal d'Investigació Mèdica, E-08003 Barcelona, Spain; ^kCenter for Advanced Studies, Research and Development in Sardinia (CRS4), 09016 Pula, Italy; ^lDepartment of Genetic Medicine and Development, University of Geneva Medical School, 1211 Geneva, Switzerland; ^mCenter for Integrative Genomics, Genopode, ⁿCentre de Regulació Genòmica, Universitat Pompeu Fabra, E-08003 Barcelona, Spain

(Received October 20, 2006)



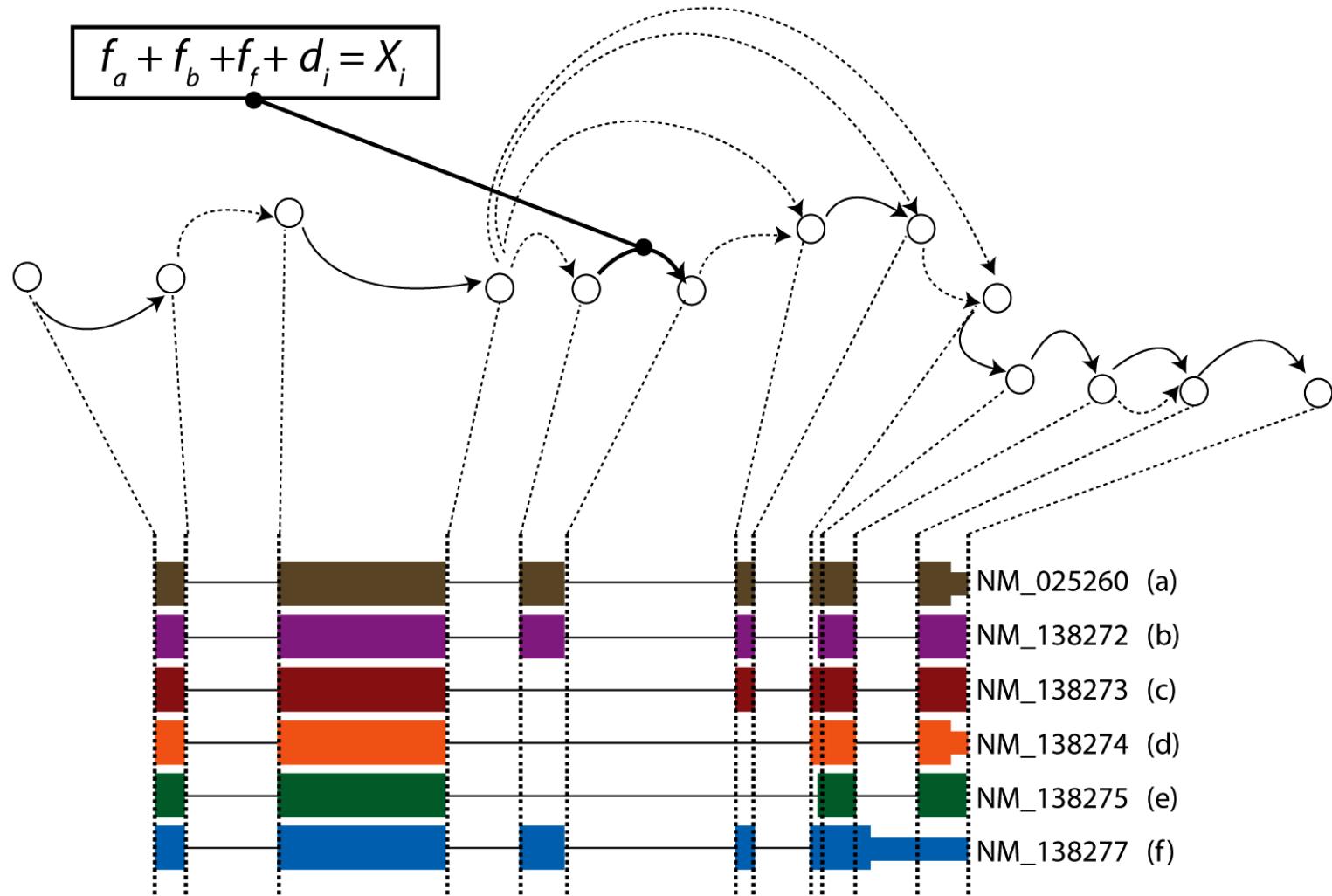
① Transcriptome genetics using second generation sequencing in a Caucasian population

Stephen B. Montgomery^{1,2*}, Micha Sammeth³, Maria Gutierrez-Arcelus¹, Radoslaw P. Lach², Catherine Ingle², James Nisbett², Roderic Guigo³ & Emmanouil T. Dermitzakis^{1,2*}



Using Flux, CRG group has been able to discover mutations with opposite effects on alternative transcript variants within the same gene

The Flux Capacitor. Compute transcript abundance from RNASeq experiments



Gene expression

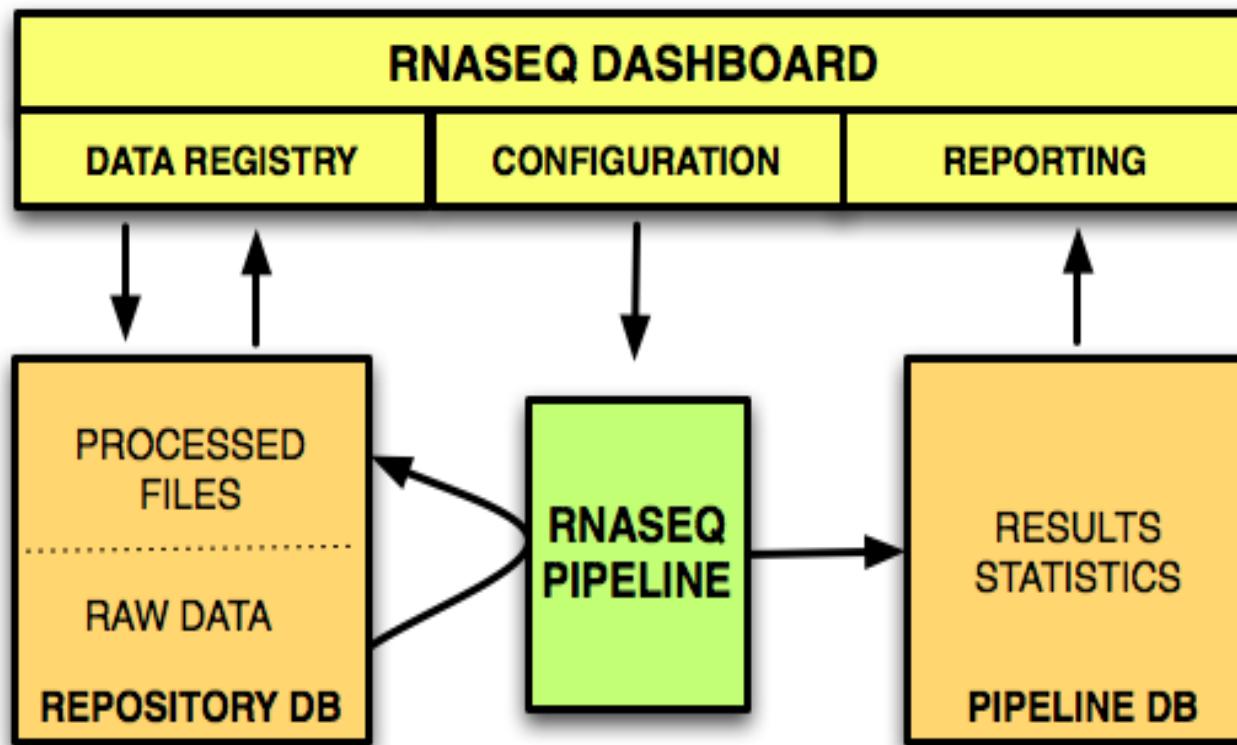
Advance Access publication January 17, 2013

Grape RNA-Seq analysis pipeline environment

David G. Knowles^{1,2,*}, Maik Röder^{1,2}, Angelika Merkel^{1,2} and Roderic Guigó^{1,2,*}

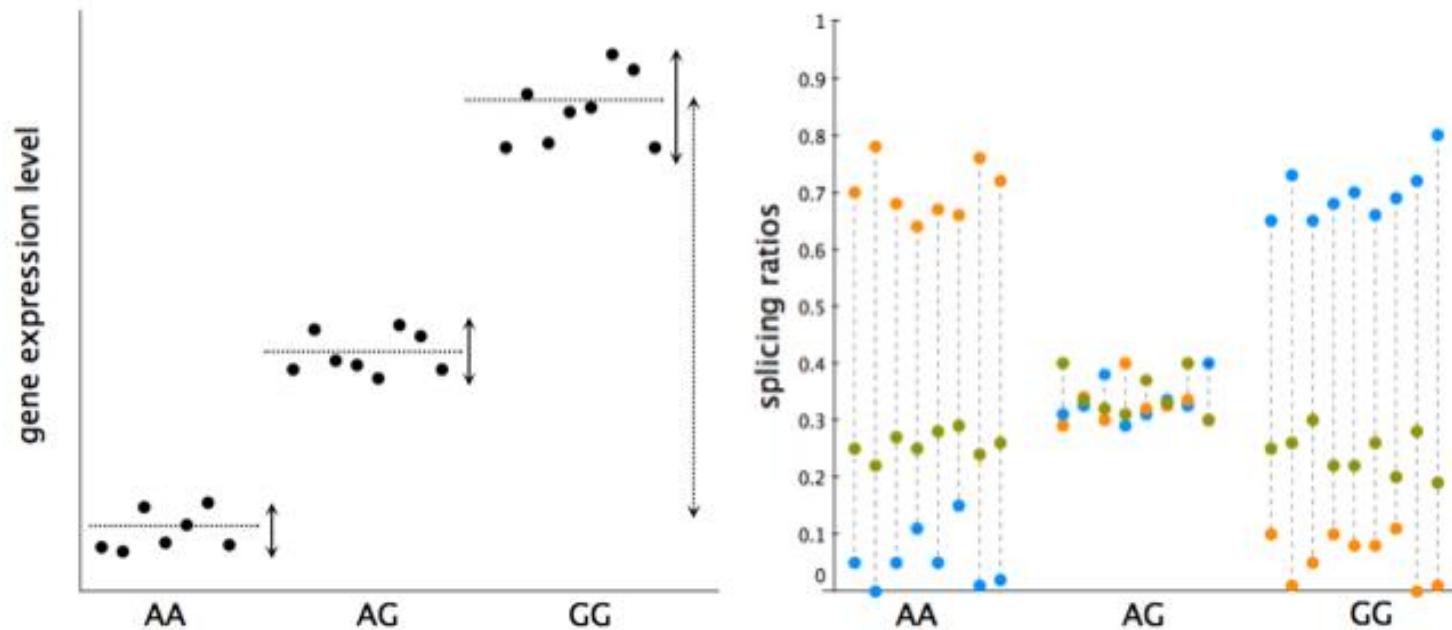
¹Bioinformatics and Genomics Group, Centre for Genomic Regulation (CRG) and ²Universitat Pompeu Fabra (UPF), Dr Aiguader 88, 08003 Barcelona, Spain

Associate Editor: Ivo Hofacker



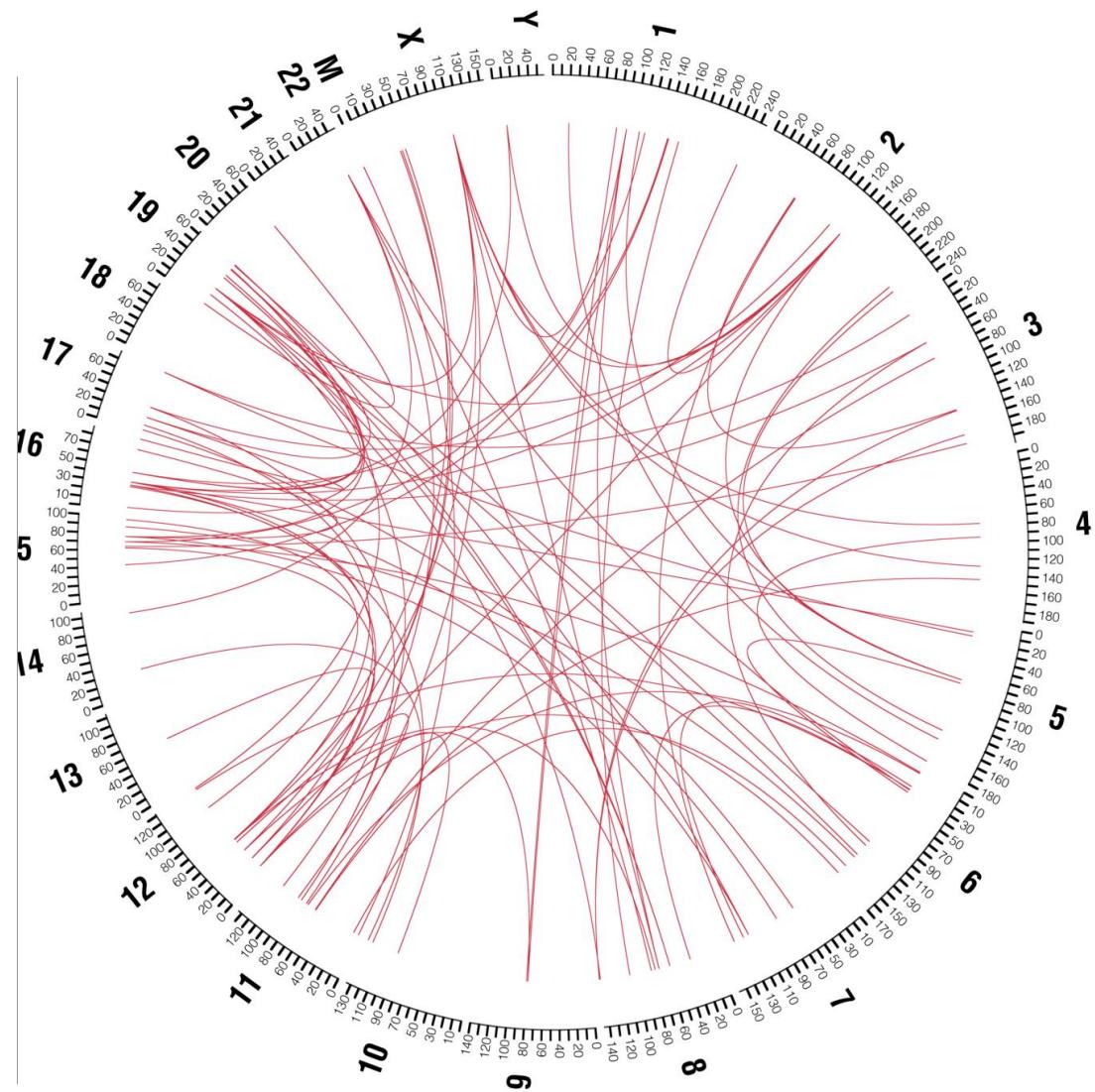
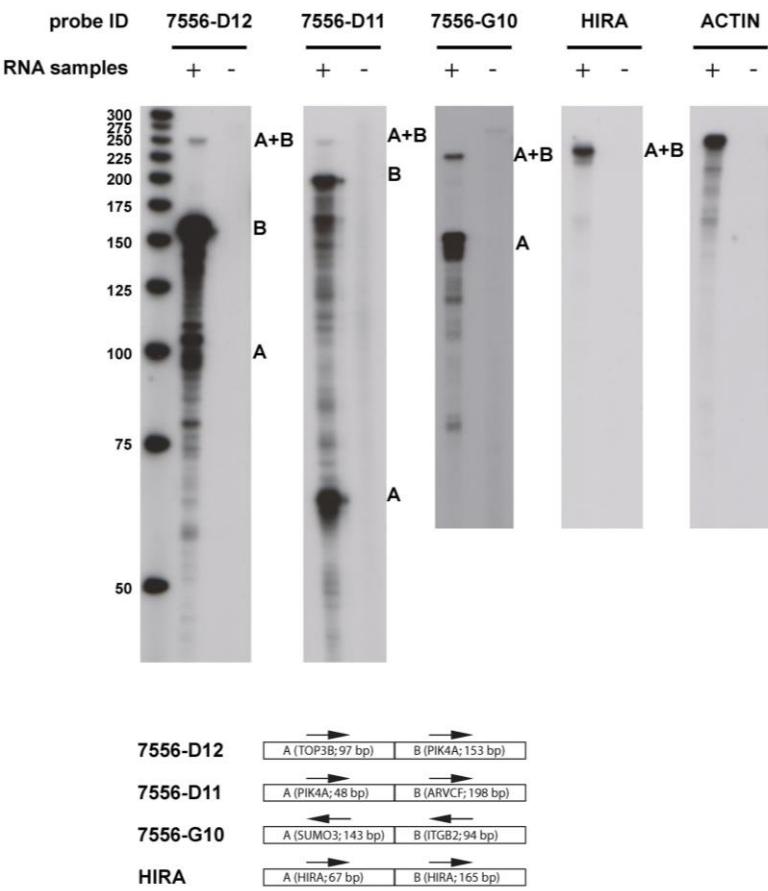
sQTLseekeR: a general method to detect genetic variants associated with alternative splicing in RNA sequencing population studies

Jean Monlong^{1,4}, Miquel Calvo³ and Roderic Guigó*^{1,2}



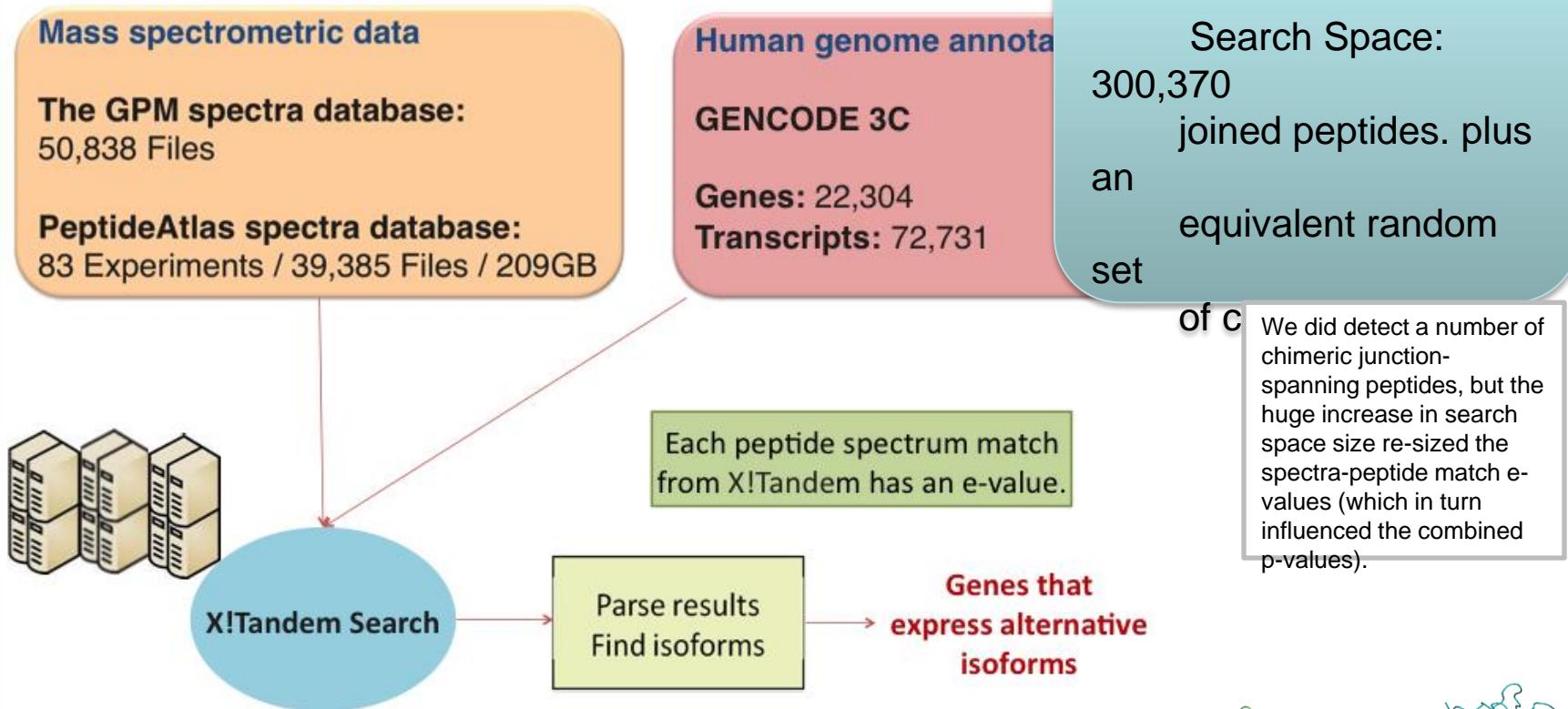
Inter-genic (trans) splicing

Experimental verification through RNAase protection assays



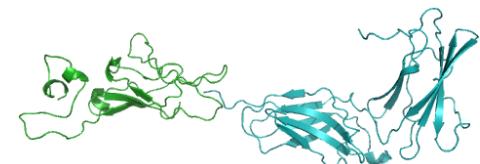
Trans-splicing and the human genome

Workflow



Human genome: **72h**

Possible exon joins: $190e^9$. A subset of 300,370 junctions (1/58 of the combinations for chromosome 21 and 22) 68h = **43000 days**

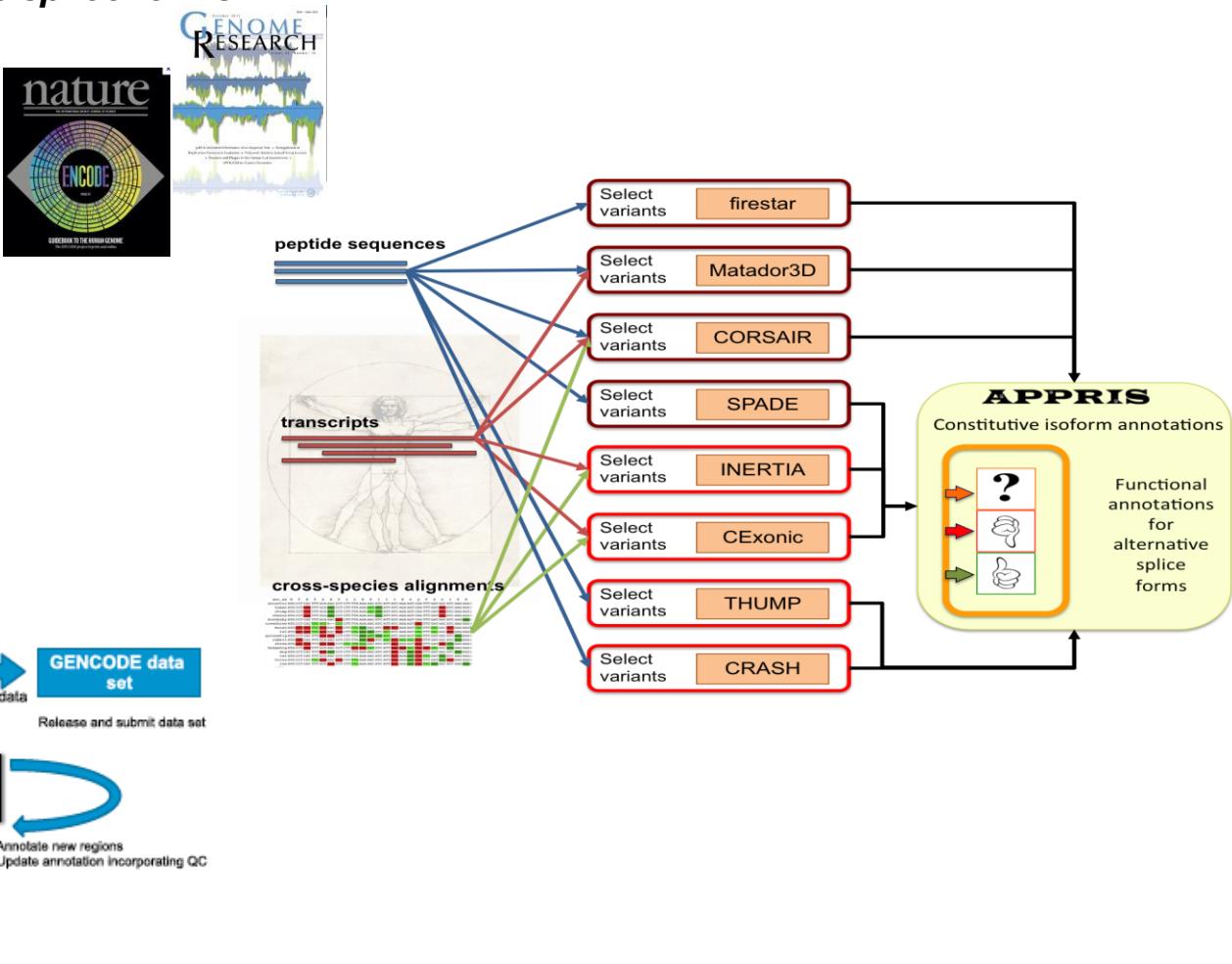


Canonical Splicing variants

APPRIS: analysis of functional annotations for alternative splice forms

Rodriguez et al., *NAR* 2012

GENCECODE/ENCODE:
Gen Res 2012, *Nature* 2012



Translicing/chimeric mRNAs

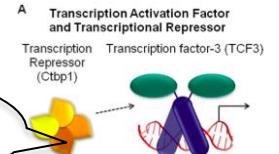
POST TRANSCRIPTIONAL REGULATION

Chimeric protein production

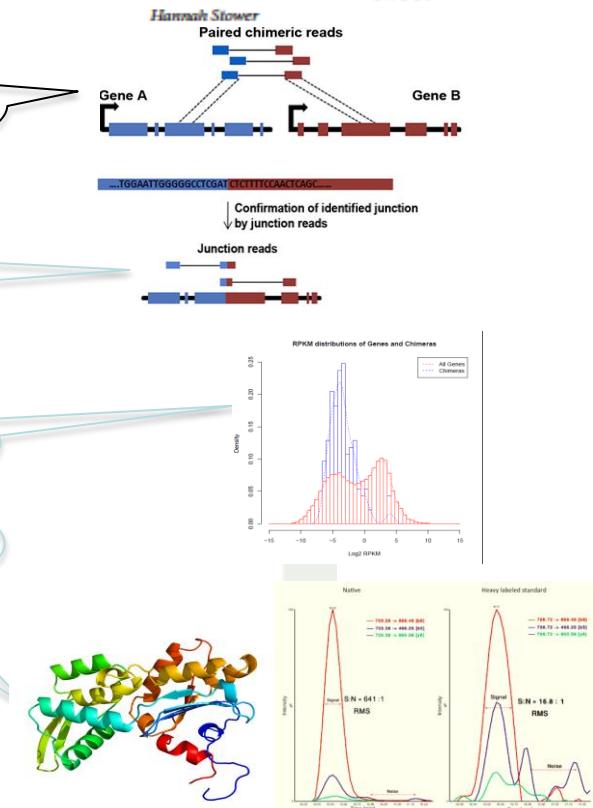
Chimeric transcripts arise from the joining of exons from two or more different genes. Reporting of such transcripts has become more widespread as data from genome-wide transcriptional analyses has increased. However, the number of known chimeric transcripts far outnumbers the reported number of chimeric proteins. Here, Frenkel-Morgenstern et al increase the number of identified translated chimeric transcripts and describe features indicative of their biological functionality.

Firstly the authors analysed previously reported human tissue RNA-seq datasets for the presence of chimeric reads, that is, those that do not align to annotated transcripts and include a chimeric exon-exon junction. This approach confirmed the expression of 175 out of 7,424 previously reported chimeric transcripts from 16 human tissues. Analysis of the level of expression of these transcripts revealed that while chimeric transcripts are themselves expressed at a low level, they incorporate transcripts that are normally expressed at a high level; they are also expressed in a highly tissue-specific manner. To confirm the translation of these transcripts, the authors both searched mass spectrometry databases and generated their own shotgun mass spectrometry datasets. They initially searched for peptides that spanned a junction site of the human chimeric transcripts and they found 12 chimeric proteins. Following this, they generated targeted mass spectrometry data and this confirmed the expression of a further three novel chimeric proteins. The 175 expressed chimeras are enriched in signal and transmembrane peptides suggesting that generating chimeric transcripts

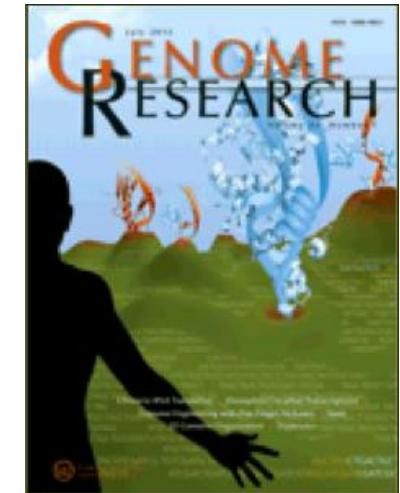
is a mechanism for altering the localisation of the translated protein. Thus the authors' data suggests a potential function for chimeric proteins. The search is on to elucidate precise biological roles.



Chimeras may produce a dominant negative effect



Two overlapping mass-spec peptides, (18 experiments, Pvalue<10⁻⁷, FDR<1%)
Dr. Levin, Weizmann Institute



Research

Chimeras taking shape: Potential functions of proteins encoded by chimeric RNA transcripts

Milana Frenkel-Morgenstern,¹ Vincent Lacroix,² Iñaki Ezkurdia,¹ Yishai Levin,³ Alexandra Gabashvili,³ Jaime Prilusky,⁴ Angela del Pozo,¹ Michael Tress,¹ Rory Johnson,⁵ Roderic Guigo,⁵ and Alfonso Valencia^{1,6}

Genome Research
www.genome.org

ChiTaRS: chimeric transcripts/proteins Database

Morgenstern et al., NAR 2012

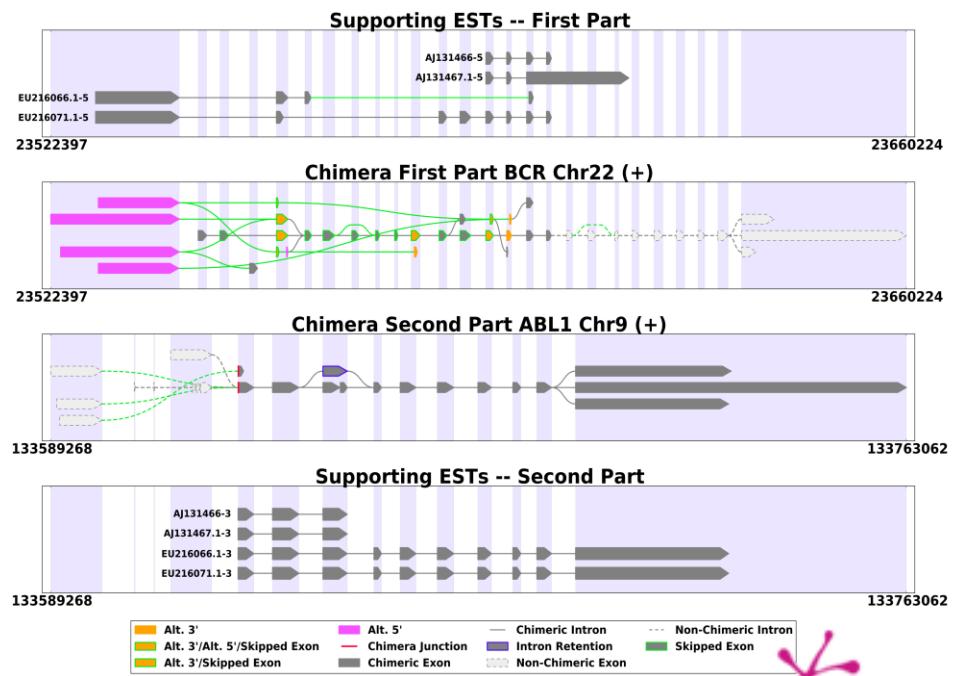
The screenshot shows the ChiTaRS homepage with a search interface for the Chimera Full Collection. Key search parameters include:

- Search by Keyword, Gene/Synonym, Tissue Name, Identity, Gene Name, Chimera Full Collection, ChimeraID.
- Dataset updates: All.
- Rank: Junction Consistency: RNAseq evidence, Breakpoints, Mass-spec Hits.
- Organisms: Homo Sapiens, Mus musculus, D. Melanogaster.

RESULT FOR THE SEARCH: CHIMERA FULL COLLECTION:

Total sequences: 16188

Organism	Graphical View	Dataset updates	First Gene (1)	Second Gene (2)	Deviation	Rank	RNAseq and/or Mass-spec evidences	Cancer Breakpoint, Pubmed Reference
Homo Sapiens	SpIGraphs	AA [1] 2012-07-16	6983 CARD10	217 99.6 CAMK2N1	211 486 100.0	0 0 0	0 NA	
Homo Sapiens	SpIGraphs	EF [2] 2012-07-16	632110 HNRNPA2B1	175 100.0 ETV1	174 417 100.0	0 0 0	0 NA	
Homo Sapiens	SpIGraphs	DA [1] 2012-07-16	134735 ZMYM2	193 99.5 RAB1A	194 551 100.0	0 0 1	Spes-1 Spes-2	
Homo Sapiens	SpIGraphs	EF [2] 2012-07-16	428111 PRKAR1A	182 100.0 RAB1A	182 417 100.0	0 0 1	Human lung total RNA, lot 0904002 causasian BestTissue = HS440 NumberOfReads = 4 NumberOfDistinctReads = 2 NumberOfTissues = 3 NumberOfReadsInBestTissue = 2 NumberOfDistinctReadsInBestTissue = 1 TissueSpecificity = 1.039720 RPKM = 0.026422	
Homo Sapiens	SpIGraphs	DA [1] 2012-07-16	092511 CHL1	272 99.7 ELAVL1	272 99.7 CHL1	0 0 1	HS440	
Homo Sapiens	SpIGraphs	BG [2] 2012-07-16	678110 GSTP1	204 95.5 PSMB1	205 459 98.9	0 0 0	from dbCRID: aberration = t(11;12)(p15;q13) location1 = p15 location2 = q13	
Homo Sapiens	SpIGraphs	AJ [1] 2012-07-16	438985 NUP98	737 100.0 HOXC13	734 992 100.0	0 0 0	NA	12619167



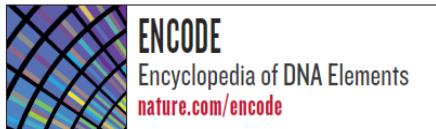
ARTICLE

An integrated encyclopedia of DNA elements in the human genome

The ENCODE Project Consortium*

The human genome encodes the blueprint of life, but the function of the vast majority of its nearly three billion bases is unknown. The Encyclopedia of DNA Elements (ENCODE) project has systematically mapped regions of transcription, transcription factor association, chromatin structure and histone modification. These data enabled us to assign biochemical functions for 80% of the genome, in particular outside of the well-studied protein-coding regions. Many discovered candidate regulatory elements are physically associated with one another and with expressed genes, providing new insights into the mechanisms of gene regulation. The newly identified elements also show a statistical correspondence to sequence variants linked to human disease, and can thereby guide interpretation of this variation. Overall, the project provides new insights into the organization and regulation of our genes and genome, and is an expansive resource of functional annotations for biomedical research.

The human genome sequence provides the underlying code for human biology. Despite intensive study, especially in identifying protein-coding genes, our understanding of the genome is far from complete, particularly with



95% of the genome lies within 8 kilobases (kb) of a DNA–protein interaction (as assayed by bound ChIP-seq motifs or DNase I footprints), and 99% is within 1.7 kb of at least one of the biochemical events measured by ENCODE.

Science AAAS.org | Feedback | Help | Librarians | AAAS SCIENTIFIC JOURNALS | EDITOR'S CHOICE | SCIENCE NOW | CAL UNIVERSITAT DE BARCELONA | ALERTS | ACCESS RIGHTS | MY ACCOUNT | SIGN IN

Science The World's Leading Journal of Original Scientific Research, Global News, and Commentary.

Science Home Current Issue Previous Issues Science Express Science Products My Science About the Journal

Home > Science Magazine > 7 September 2012 > Pennisi, 337 (6099): 1159-1161

Article Views Vol. 337 no. 6099 pp. 1159-1161 DOI: 10.1126/science.337.6099.1159

• Summary NEWS & ANALYSIS

• Full Text GENOMICS

• Full Text (PDF) ENCODE Project Writes Eulogy for Junk DNA

Article Tools Elizabeth Pennisi

• Save to My Folders

• Download Citation

• Alert Me When Article Is Cited

• Post to CiteULike

• E-mail This Page

• Rights & Permissions

• Cited by (100)

Read the Full Text

This week, 30 research papers, including six in *Nature* and additional papers published online by *Science*, sound the death knell for the idea that our DNA is mostly littered with useless bases. A decade-long project, the Encyclopedia of DNA Elements (ENCODE), has found that 80% of the human genome serves some purpose, biochemically speaking. Beyond defining proteins, the DNA bases highlighted by ENCODE specify landing spots for proteins that influence gene activity, strands of RNA with myriad roles, or simply places where chemical modifications serve to silence stretches of our chromosomes.

The GEM mapper: fast, accurate and versatile alignment by filtration

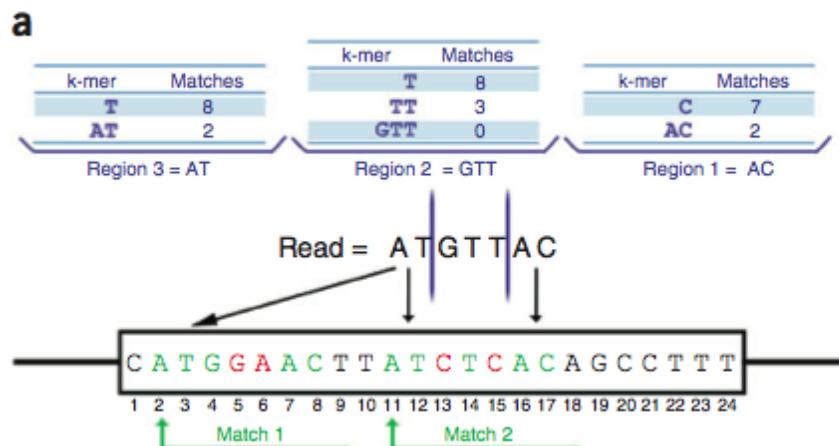
Santiago Marco-Sola¹, Michael Sammeth¹, Roderic Guigó² & Paolo Ribeca^{1,3}

Because of ever-increasing throughput requirements of sequencing data, most existing short-read aligners have been designed to focus on speed at the expense of accuracy. The Genome Multitool (GEM) mapper can leverage string matching by filtration to search the alignment space more efficiently, simultaneously delivering precision (performing fully tunable exhaustive searches that return all existing matches, including gapped ones) and speed (being several times faster than comparable state-of-the-art tools).

In recent years, the superexponential increase in worldwide sequencing capacity¹ has driven a substantial amount of research into the development of efficient algorithms for the analysis of short-read sequence data. In particular, rapid alignment of reads

alignment to the rest of the sequence; as few mismatches are usually allowed in the seed, this approach provides good performance but is also inflexible and incapable of returning all existing matches.

As the needs of the field are rapidly shifting, the assumptions described above are now constantly challenged. Some current biological problems (such as nonmodel-organism studies for which matching reference sequences might be incomplete, inaccurate or missing, or cross-species comparisons in evolutionary studies²) and new experimental protocols (data in color space, bisulfite-converted sequences or RNA sequencing³) require a higher tolerance to errors or more flexible alignment models^{4,5}. Other applications (prediction of genomic variation or



GEM: crystal-clear DNA alignment

Gregory G Faust & Ira M Hall

The Genome Multitool (GEM) mapper rapidly and accurately provides all alignments of a read within a user-defined number of mismatches.

Over the past several years, the flood of data produced by second-generation DNA sequencers has motivated the development

of a new generation of DNA alignment tools designed for mapping short reads to large reference genomes¹. These aligners have allowed

Gregory G. Faust is in the Department of Computer Science, Ira M. Hall is at the Center for Public Health Genomics and both are in the Department of Biochemistry and Molecular Genetics, University of Virginia, Charlottesville, Virginia, USA.

e-mail: irahall@virginia.edu

of heuristics can provide an excellent cost-performance trade-off and in general achieve high-quality results. Yet, it is inevitable that in certain situations, heuristic strategies will fail to identify biologically relevant alignments. The key advance of the GEM aligner is that it uses a strategy that is guaranteed to identify all alignments within a user-specified edit distance. Remarkably, GEM accomplishes this difficult task while also achieving a significant speed increase over other popular aligners.

GEM's impressive speed and accuracy derive from its use of a novel strategy for pruning the

partitioning, each match is extended into a full alignment using dynamic programming. This latter step often dominates the run time of traditional 'seed-and-extend' aligners that rely on fixed-length seeds for initial matching. In contrast, GEM adaptively partitions queries into variably sized subsequences according to their uniqueness in the reference genome. This drastically reduces the total number of reference-genome loci requiring the computationally expensive extension step, thereby dramatically improving overall speed.

1. Li, H. & Homer, N. *Brief. Bioinform.* **11**, 473–483 (2010).
2. Marco-Sola, S., Sammeth, M., Guigó, R. & Ribeca, P. *Nat. Methods* **9**, 1185–1188 (2012).
3. Li, R. et al. *Bioinformatics* **25**, 1966–1967 (2009).
4. Langmead, B., Trapnell, C., Pop, M. & Salzberg, S.L. *Genome Biol.* **10**, R25 (2009).
5. Langmead, B. & Salzberg, S.L. *Nat. Methods* **9**, 357–359 (2012).
6. Li, H. & Durbin, R. *Bioinformatics* **25**, 1754–1760 (2009).
7. Alkan, C. et al. *Nat. Genet.* **41**, 1061–1067 (2009).
8. Li, H. & Durbin, R. *Bioinformatics* **26**, 589–595 (2010).
9. Alkan, C., Coe, B.P. & Eichler, E.E. *Nat. Rev. Genet.* **12**, 363–376 (2011).

Genetic basis of complex diseases

Chronic Lymphocytic Leukemia

CLL-ICGC. Cancer Project



50 cancers
25000 cancer genomes

International Cancer Genome Consortium - Windows Internet Explorer
http://www.icgc.org/

Archivo Edición Ver FAVORITOS Herramientas Ayuda

IRB Barcelona Mail - d - m... Resultados de la Búsqueda... Internet Explorer no pue... Internat...

International Cancer Genome Consortium

Overview | Cancer Genome Projects | Committees | Policies and Guidelines

International Cancer Ger...

- Brain Cancer United States
- Breast Cancer European Union / United Kingdom
- Breast Cancer France
- Breast Cancer United Kingdom
- Chronic Lymphocytic Leukemia Spain
- Colon Cancer United States
- Gastric Cancer China
- Leukemia United States
- Liver Cancer France
- Liver Cancer Japan

and control of cancer. The ICGC will facilitate communication among members and provide a forum for coordination and efficiency among the scientists working to understand diseases.

nature International network projects. Nature 4 (2010) HTML

ICGC Public Presentation April 15, 2011
International Cancer Genome Consortium (ICG)
Guidelines HTML PDF

Members of the ICG Committed Projects

Spain

Funding Organizations

Spain: Institute of Health Carlos III

Spain: Spanish Ministry of Science and Innovation

Research Organizations

Spain: Barcelona Supercomputing Centre (BSC)

Spain: ICO (Institut Català d'Oncologia)

Spain: CIC (Centro de Investigación del Cáncer)

Spain: Center for Cancer Research, University Hospital

Spain: National Genome Analysis Centre

Spain: Hospital Clinic of Barcelona

Spain: Pompeu Fabra University

Spain: Spanish National Cancer Research Centre

Spain: University of Deusto

Spain: University of Oviedo

Spain: University of Santiago de Compostela

Kingston (F)

Inicio | International Cancer ... | Presentación | KINGSTON (F)

Chronic Lymphocytic Leukemia - CLL with r...

Spain

Funding Organizations

Spain: Institute of Health Carlos III

Spain: Spanish Ministry of Science and Innovation

Research Organizations

Spain: Barcelona Supercomputing Centre (BSC)

Spain: ICO (Institut Català d'Oncologia)

Spain: CIC (Centro de Investigación del Cáncer)

Spain: Center for Cancer Research, University Hospital

Spain: National Genome Analysis Centre

Spain: Hospital Clinic of Barcelona

Spain: Pompeu Fabra University

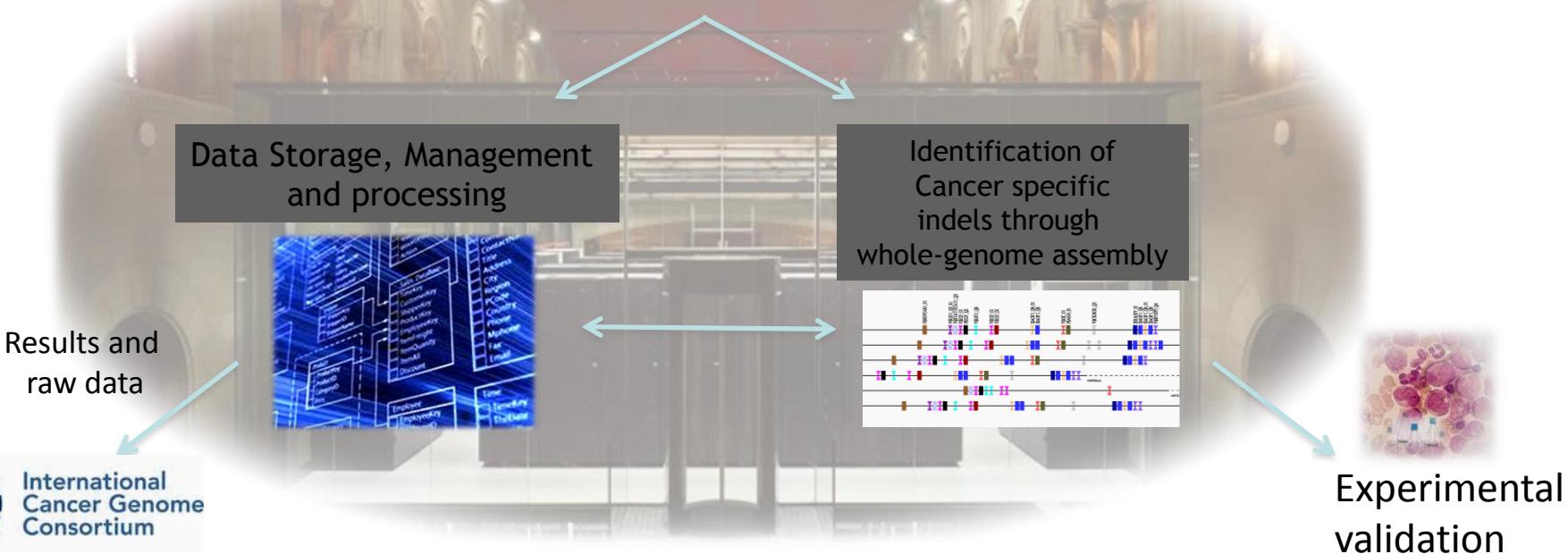
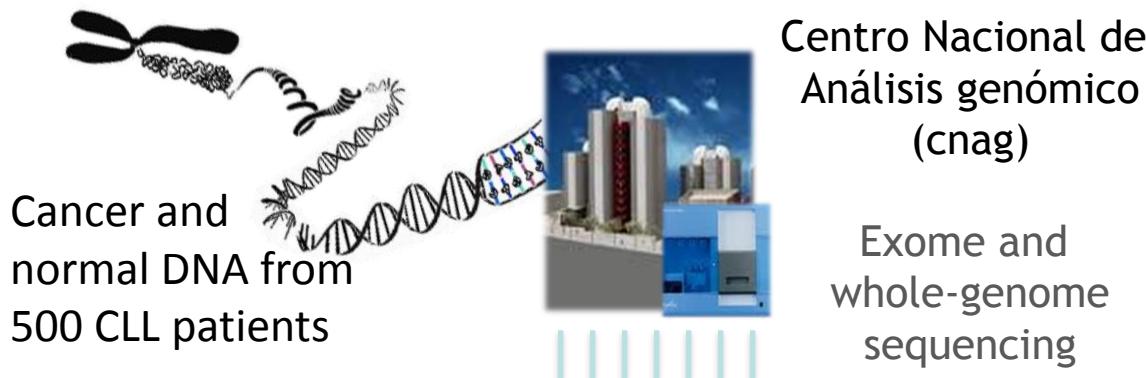
Spain: Spanish National Cancer Research Centre

Spain: University of Deusto

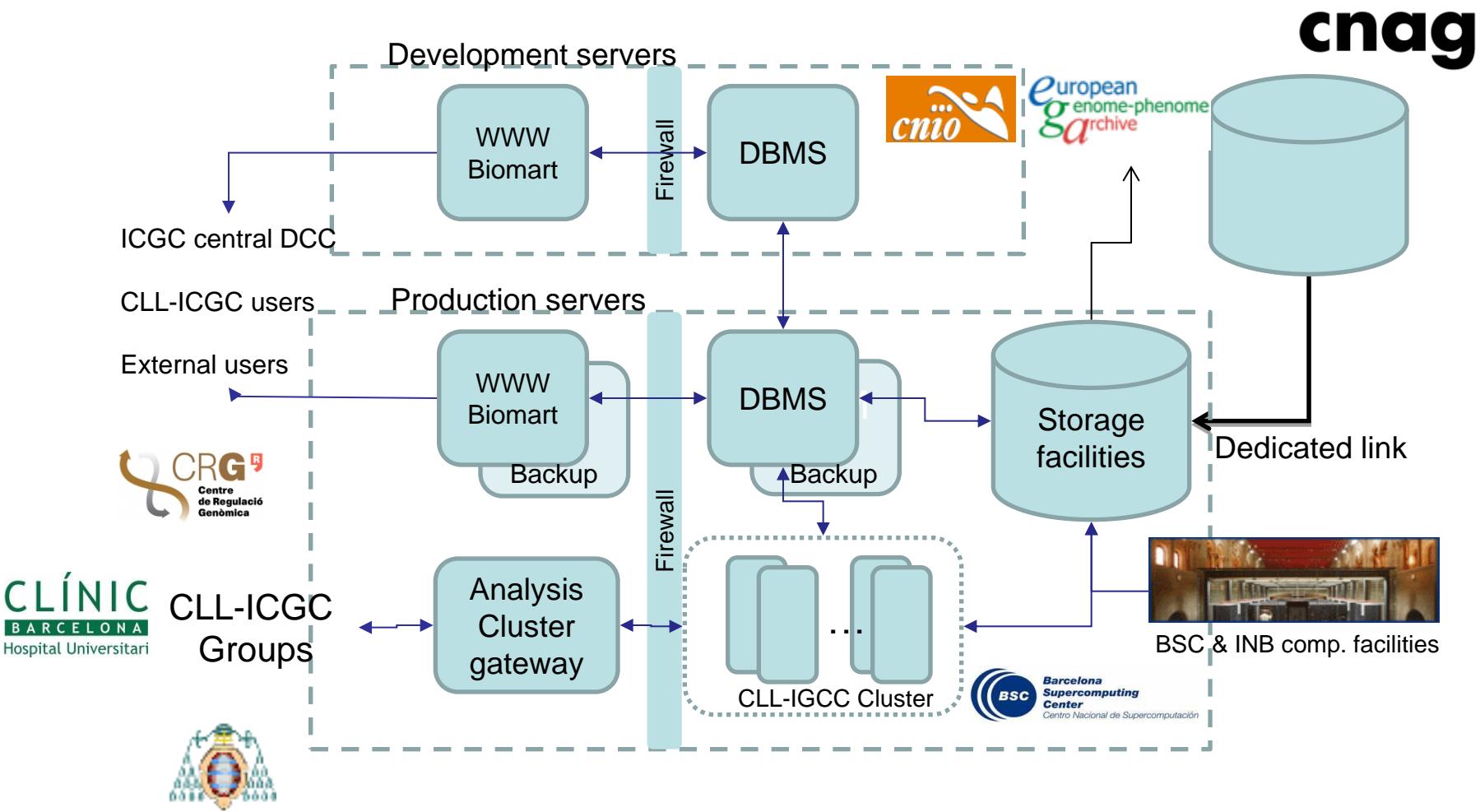
Spain: University of Oviedo

Spain: University of Santiago de Compostela

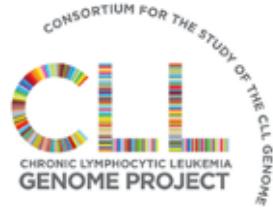
ICGC DATA PROCESSING WORKFLOW



CLL INFRASTRUCTURE SETUP



Chronic Lymphocytic Leukemia



Análisis de los datos de Hseq de **Chronic Lymphocytic Leukemia (CLL) Genome Project**: Objetivo entender las bases genéticas de la Leucemia



Secuenciación
Genoma



500 Paciente

BSC

Data Management
esperable 1.5 Pb



HPC Computing
Hasta 10% uso de MN

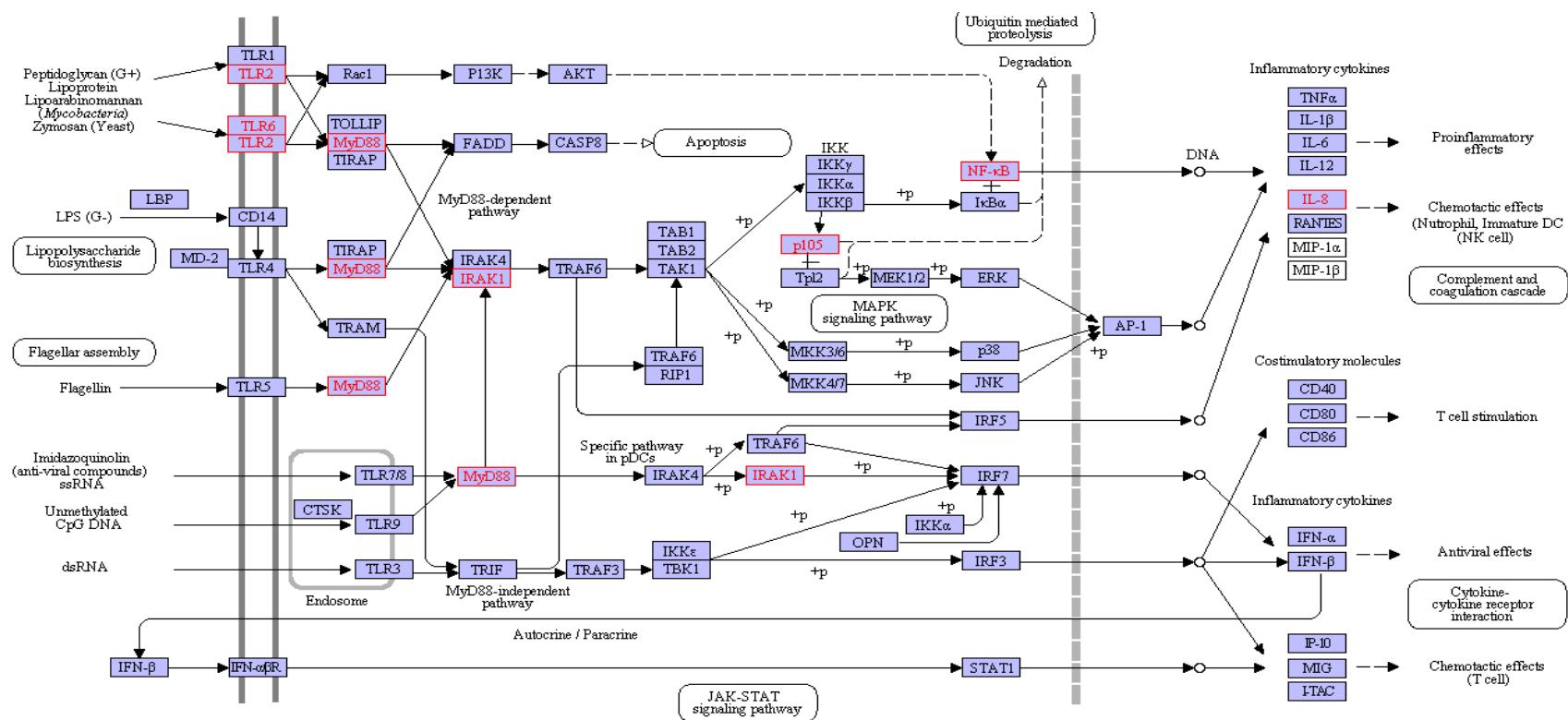
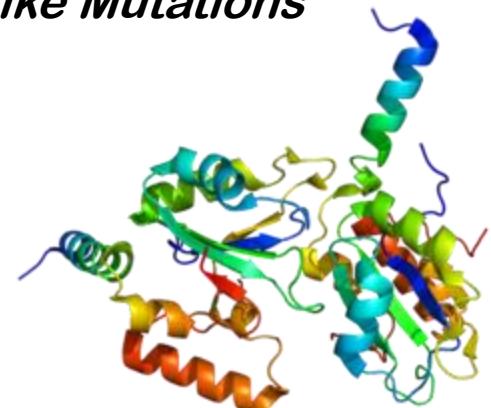
Validación
Experimental



Objetivo: Encontrar la bases genéticas de la Leucemia

Exome sequencing identifies recurrent mutations of the splicing factor *SF3B1* gene in chronic lymphocytic leukemia

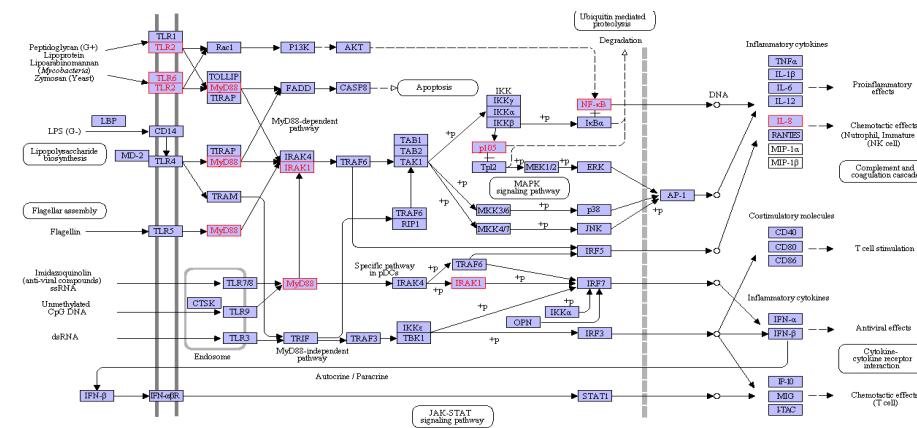
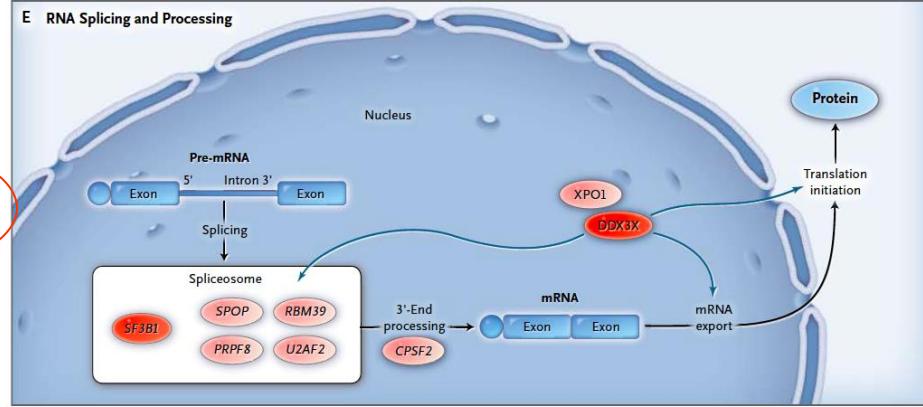
Victor Quesada¹, Laura Conde², Neus Villamor², Gonzalo R Ordóñez¹, Pedro Jares², Laia Bassaganyas³, Andrew J Ramsay¹, Silvia Bea², Magda Pinyol⁴, Alejandra Martínez-Trillo⁵, Mónica López-Guerra², Dolors Colomer², Alba Navarro², Tycho Baumann⁵, Marta Aymerich², María Rozman², Julio Delgado⁵, Eva Gimé⁵, Jesús M Hernández⁶, Marcos González-Díaz⁶, Diana A Puente¹, Gloria Velasco¹, José M P Freije¹, José M C Tubio³, Romina Royo⁷, Josep L Gelpí⁷, Modesto Orozco⁷, David G Pisano⁸, Jorge Zamora⁸, Miguel Vázquez⁸, Alfonso Valencia⁸, Heinz Himmelbauer⁹, Mónica Bayés¹⁰, Simon Heath¹⁰, Marta Gut¹⁰, Ivo Gut¹⁰, Xavier Estivill³, Armando López-Guillermo⁵, Xose S Puente¹, Elias Campo^{2,11} & Carlos López-Otín^{1,11}



Exome sequencing identifies recurrent mutations of the splicing factor *SF3B1* gene in chronic lymphocytic leukemia

Víctor Quesada¹, Laura Conde², Neus Villamor², Gonzalo R Ordóñez¹, Pedro Jares², Laia Bassaganyas³, Andrew J Ramsay¹, Silvia Beá², Magda Pinyol⁴, Alejandra Martínez-Trillo⁵, Mónica López-Guerra², Dolors Colomer², Alba Navarro², Tycho Baumann⁵, Marta Aymerich², María Rozman², Julio Delgado⁵, Eva Giné⁵, Jesús M Hernández⁶, Marcos González-Díaz⁶, Diana A Puente¹, Gloria Velasco¹, José M P Freije¹, José M C Tubio³, Romina Royo⁷, Josep L Gelpí⁷, Modesto Orozco⁷, David G Pisano⁸, Jorge Zamora⁸, Miguel Vázquez⁸, Alfonso Valencia⁸, Heinz Himmelbauer⁹, Mónica Bayés¹⁰, Simon Heath¹⁰, Marta Gut¹⁰, Ivo Gut¹⁰, Xavier Estivill³, Armando López-Guillermo⁵, Xose S Puente¹, Elias Campo^{2,11} & Carlos López-Otín^{1,11}

Here we perform whole-exome sequencing of samples from 105 individuals with chronic lymphocytic leukemia (CLL)^{1,2}, the most frequent leukemia in adults in Western countries. We found 1,246 somatic mutations potentially affecting gene function and identified 78 genes with predicted functional alterations in more than one tumor sample. Among these genes, *SF3B1*, encoding a subunit of the spliceosomal U2 small nuclear ribonucleoprotein (snRNP), is somatically mutated in 9.7% of affected individuals. Further analysis in 279 individuals with CLL showed that *SF3B1* mutations were associated with faster disease progression and poor overall survival. This work provides the first comprehensive catalog of somatic mutations in CLL with relevant clinical correlates and defines a large set of new genes that may drive the development of this common form of leukemia. The results reinforce the idea that targeting several well-known genetic pathways, including mRNA splicing, could be useful in the treatment of CLL and other malignancies.



The NEW ENGLAND JOURNAL OF MEDICINE
N ENGL J MED 365;26 NEJM.ORG DECEMBER 29, 2011

ORIGINAL ARTICLE

SF3B1 and Other Novel Cancer Genes in Chronic Lymphocytic Leukemia

Lili Wang, M.D., Ph.D., Michael S. Lawrence, Ph.D., Youzhong Wan, Ph.D., Petar Stojanov, B.A., Carrie Sougnez, B.S., Kristen Stevenson, M.S., Lillian Werner, M.S., Andrey Sivachenko, Ph.D., David S. DeLuca, Ph.D., Li Zhang, Ph.D., Wandi Zhang, M.D., Alexander R. Vartanov, B.A., Stacey M. Fernandes, B.S., Natalie R. Goldstein, B.A., Eric G. Folco, Ph.D., Kristian Cibulskis, B.S., Bethany Tesar, M.S., Quinlan L. Sievers, B.A., Erica Shefler, B.S., Stacey Gabriel, Ph.D., Nir Hacohen, Ph.D., Robin Reed, Ph.D., Matthew Meyerson, M.D., Ph.D., Todd R. Golub, M.D., Eric S. Lander, Ph.D., Donna Neuberg, Sc.D., Jennifer R. Brown, M.D., Ph.D., Gad Getz, Ph.D., and Catherine J. Wu, M.D.

