



***Workflows management:
new abilities for the biological
information overflow***

Luciano Milanesi

*Institute Biomedical Technologies CNR
luciano.milanesi@itb.cnr.it*



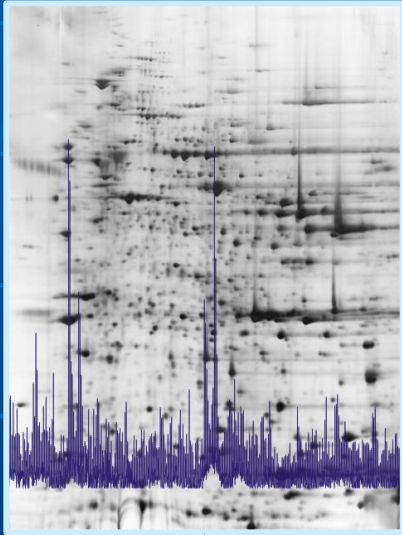
CNR - National Research Council
ITB - Institute of Biomedical Technologies

*Workflows management:
new abilities for the biological
information overflow*

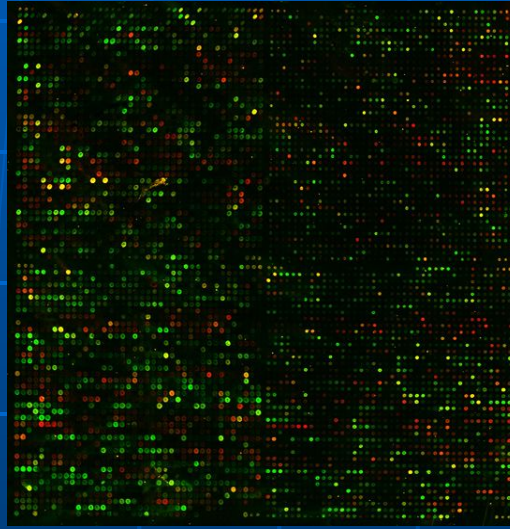
Luciano Milanesi

*Institute Biomedical Technologies CNR
luciano.milanesi@itb.cnr.it*

Complex Disease Mapping



Proteins
(Proteome)



Microarray
(Genome)



Gene & SNPs
(Genome)

HTS Data Project

Bioinformatics: Emerging Opportunities and Emerging Gaps

Paula E. Stephan and Grant Black

- A typical gene lab can produce 100 terabytes of information a year, the equivalent of 1 million encyclopedias.
- Few biologists have the computational skills needed to fully explore such an astonishing amount of data; nor do they have the skills to explore the exploding amount of data being generated from clinical trials.
- The immense amount of data that are available, and the knowledge that this is but the tip of the data iceberg.

HTS Data Project



EST

Microsatellite

MSMS

**DNA High Throughput
Sequencing**

SNP's

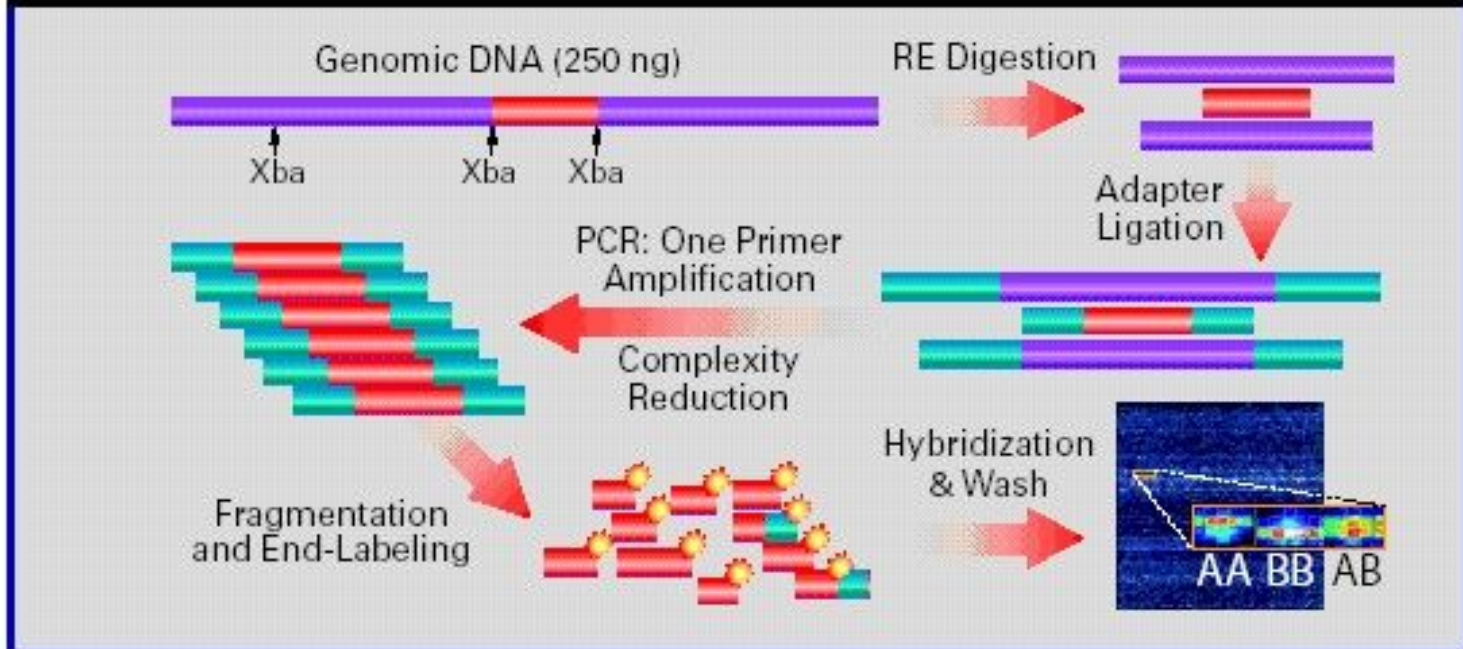
Microarray

HTS Data Project

DNA High Throughput Sequencing



Figure 1: GeneChip® Mapping Assay Overview.

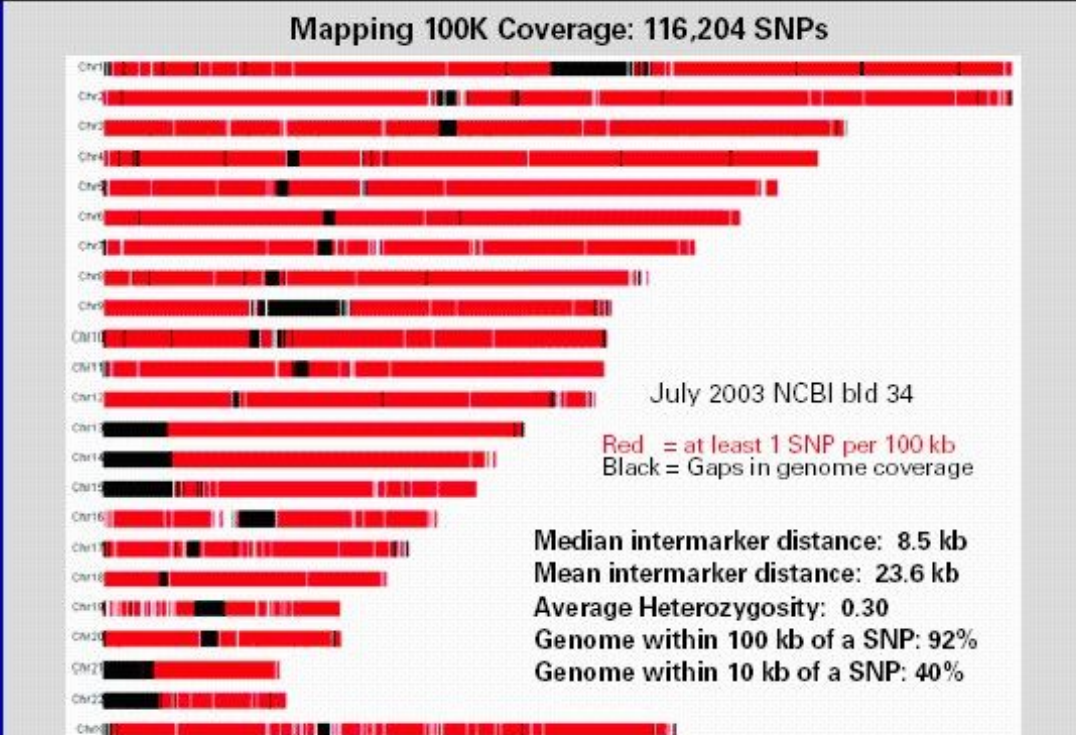


HTS Data Project

SNP's



Figure 2: Genome Coverage of Mapping 100K SNPs by chromosome. Black areas represent gaps in the human genome sequence, primarily centromeres and telomeres.



HTS Data Project

SNP's



The Sardinian Challenge



From homologies of surnames, linguistic roots and genetic markers, Sardinians have been subdivided into 31 sub-populations. Long term isolated villages within these sub-populations can be considered as true "Mendelian Breeding Units (MBUs)" just like the precious collection of the Jackson Lab (J.H.) highly inbred mouse strains whose contribution to the study of mammalian genetics (*Homo sapiens* included) has been, and still is, simply outstanding.



ANDTIASQDT PAK **VIK** **ANK LKI** **LKDYVDD** **LKTYNNTYSN** **VVTVAGEDRI** **ETAIELSSK** **YVNSDDKNAIT**
DKAVNDIVLV **GSTSIVDGLV** **ASPLASEKTA** **PLLLTSK** **DKL** **DSSVKSEIKR** **VMNLKSDTGI** **NTSKKVYLAG** **GVNSISKDVE**
NELKNMGLKV **TRLSGEDRYE** **TSLALADEIG** **LDNDKAFVVG** **GTGLADAMSI** **APVASQLK** **DG** **DATPIVVVDG**
KAKEISDDAK **SFLGTSDDVI** **IGGKNSVSKE** **IEESIDSATG** **KTPDRISGDD** **RQATNAEVLK** **EDDYFTDGEV**
VNYFVAK **DGS** **TKEDQLVDAL** **AAPIAGR** **FK** **ESPAPIILAT** **DTLSSDQNVAVSK** **AVPKDGG** **TNLVQVGKGI**
ASSVINKMKD **LLDM**

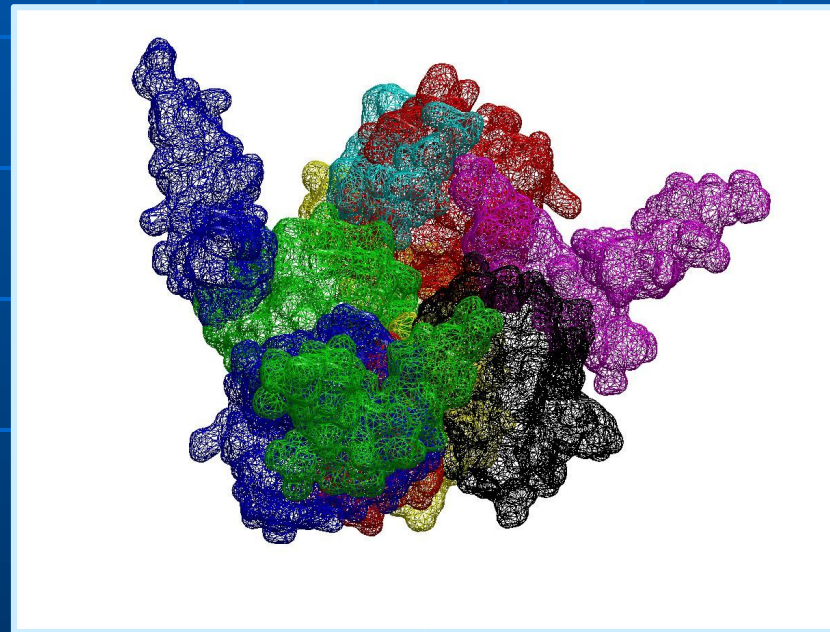
>monoisotopic mass = **39480**

position sequence (NCBI BLAST link)

```

-----
18- 22  ANKLK
245- 249 NSVSK
262- 265 TPDR
116- 119 SEIK
344- 347 AVPK
298- 302 DGSTK
266- 271 ISGDDR
163- 168 LSGEDR
116- 120 SEIKR
60- 66  YVNSDDK
224- 230 EISDDAK
155- 159 NMGLK
67- 72  NAITDK
126- 134 SDTGINTSK
160- 168 VTRLSGEDR
21- 25  LKDLK
108- 115 DKLDSSVK
135- 147 RVYLAGGVNSISK
14- 17  VVIK
1- 13   ANDTIASQDTPAK
126- 147 SDTGINTSKKVYLAGGVNSISK
14- 20  VVIKANK

```



Disease Gene

- Mutations in target sequences are usually revealed by either phenotypic selection in experimental test systems or, in case of disease-causing genes in humans, by clinical studies in which certain genes are sequenced in groups of patients and in control groups.
- Both the experimental test systems and the clinical studies rely on detectable (mutable) positions, which are sites where DNA sequence changes cause phenotypic changes.

Human genetic identity

- Genomic sequence **99.9%** identical
- **3,200,000** nucleotides different
- Single base differences in genomes between any two individuals: ca. **3 million**
- Amino acid differences in proteomes between any two individuals: ca. **100,000**

Variation types

- Macro:
 - Chromosome numbers
 - Segmental duplications, rearrangements, and deletions
- Medium:
 - Sequence Repeats
 - Transposable Elements
 - Short Deletions, Sequence and Tandem Repeats (including microsatellites)
- Micro:
 - **Single Nucleotide Polymorphisms (SNPs)**
 - **Single Nucleotide Insertions and Deletions (Indels)**

What are Single Nucleotide Polymorphisms (SNPs)?

ATGGTAA**C**CTGAG**C**TGACTTAGCGT-AT

ATGGTAA**A**CTGAG**T**TGACTTAGCGTCAT



snp



snp



indel

SNPs result from replication errors and DNA damage

Are all SNPs really SNPs?

- A SNP is found by aligning overlapping DNA sequences and identifying variable positions
- Two types of errors in SNP finding
 - Inclusion of paralogs in sequence alignment
 - Sequencing errors

GCATGCAAGCAGATA
GCATGCA^CGCAGATA
GCATGCAAGCAGATA
GCATGCAAGCAGATA
GCATGCAAGCAGATA
GCATGCAAGCAGATA
GCATGCAAGCAGATA

Application of SNP data

- Study of evolution
 - Traces evolutionary history of different populations
- DNA fingerprinting
 - criminal or parental verification
- Markers for mapping of polygenic traits
- Genotype-specific medication
- Most genes contain SNPs
 - 93% genes have one or more SNPs
 - 39% have more than 10 SNPs)

Types of SNPs

Genic, coding SNPs

- non-synonymous

 - Maintaining vs. altering protein structure/function

- synonymous

 - Maintaining vs. altering splicing

Genic, non-coding SNPs

- Regulatory SNPs

 - Maintaining vs. altering gene expression

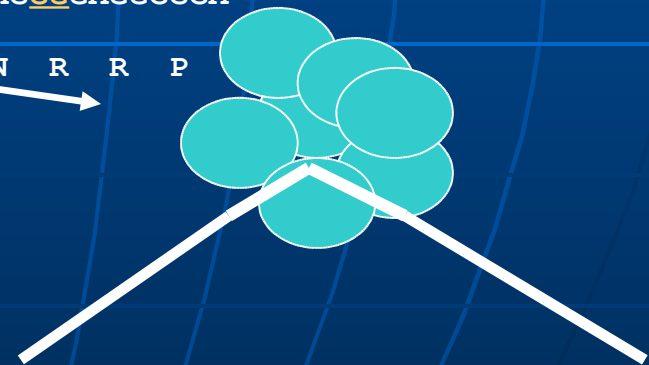
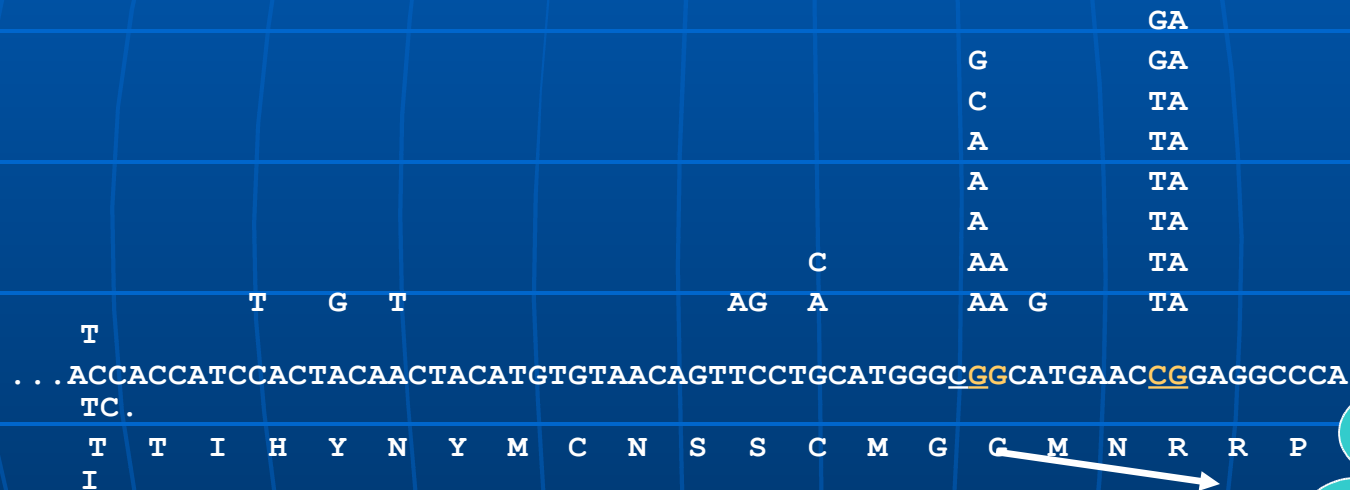
- Intronic SNPs

 - Maintaining vs. altering gene expression/splicing

Linked SNPs

- usually intergenic

Context-dependence of mutations



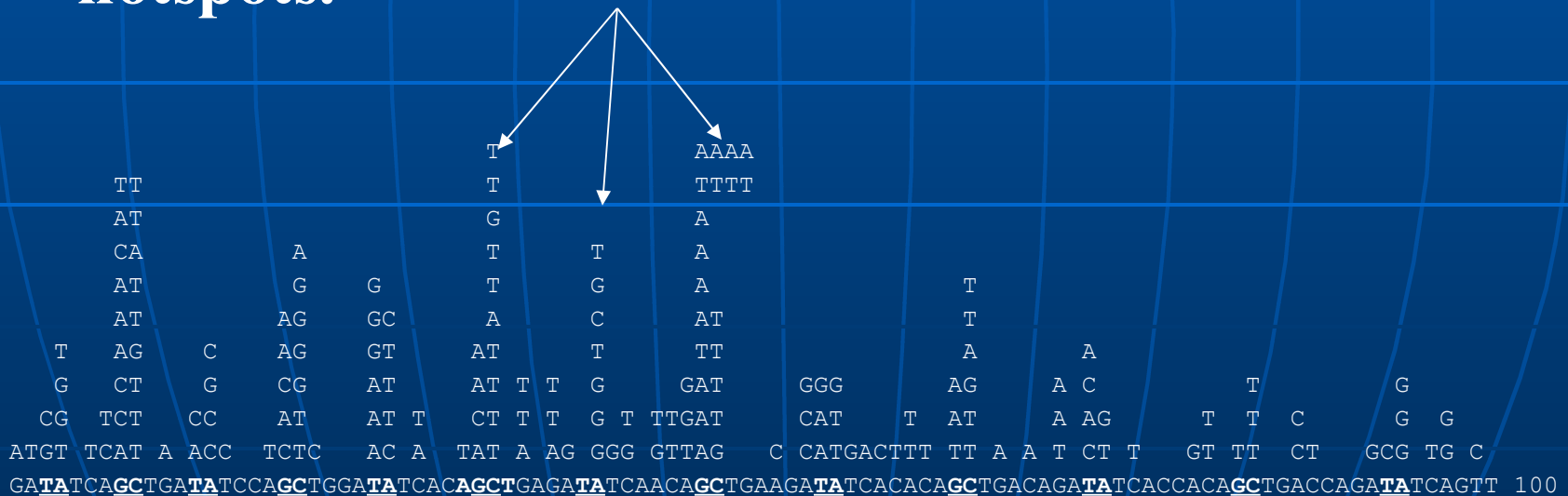
Mutational spectra

- **A mutation spectrum** is a distribution of mutation frequencies along nucleotide sequences and is compiled by the analysis of a large number of mutant target sequences.

[illegible]

Mutational spectra

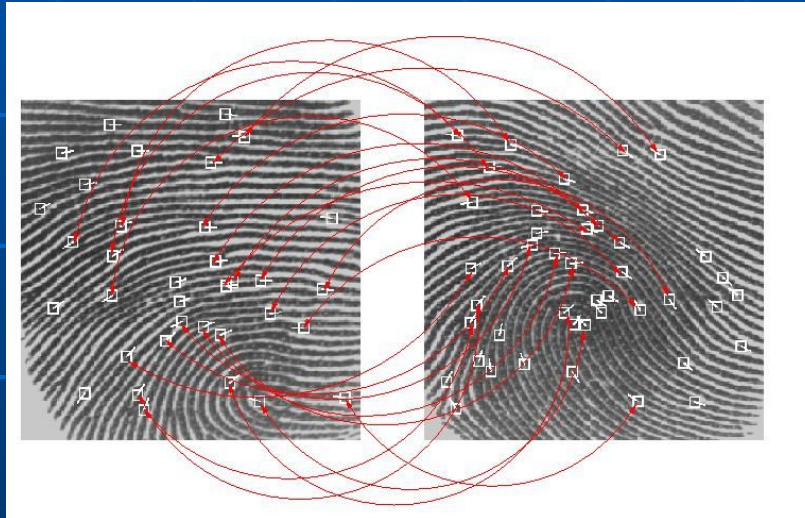
- Mutation frequencies vary significantly along nucleotide sequences such that mutations often concentrate at certain positions called **mutation hotspots**.



Hotspots in immunoglobulin V genes (Betz et al. 1993)

Mutational spectra

Mutation hotspots in DNA reflect intrinsic properties of the mutation process, such as sequence specificity, that manifests itself at the level of interaction between mutagens, DNA, and the action of the repair and replication machineries.



The nucleotide sequence context of mutation hotspots is a fingerprint of interactions between DNA and repair/replication/modification enzymes, and the analysis of hotspot context provides evidence of such interactions.

Mutational spectra

Context factors may influence mutation rates:

- Homonucleotide runs and microsatellites
- Direct and inverted repeats
- Local mononucleotide composition
- DNA conformation
- Oligonucleotide content
- higher-level features of gene
- chromatin structure

Methods

Various classification and statistical approaches are used for analysis of the nucleotide sequence context of mutations hotspots

(Rogozin, Babenko, Milanesi, Pavlov 2003 *Brief. Bioinform.* 4, 210-227).

DNA polymerase η mutation hotspots in *lacZ*

Sequence	Hotspot position	Type of changes	Number of mutations
CA <u>A</u> TT	3	A→G, T, C	15, 1, 1
TT <u>A</u> TC	14	A→G, C, T	14, 1, 1
GA <u>A</u> AT	21	A→G, T	16, 2
AT <u>A</u> GC	38	A→G, T, C	9, 2, 1
CA <u>T</u> AG	39	T→G, A, C	9, 9, 2
TC <u>A</u> TG	46	A→G, T	13, 1
GT <u>A</u> AT	50	A→G, T	16, 4
GA <u>A</u> TT	56	A→G	17
AA <u>A</u> CG	70	A→G, T	18, 3
GT <u>A</u> AA	73	A→G, T	14, 1
CG <u>A</u> CG	80	A→G, T	11, 2

W <u>A</u>	Consensus		

Methods

Correlation between nucleotide sequence features and mutation spectra. Nucleotide sequence features (mutable motifs) can be correlated with a mutation spectrum and the correlation can be tested for statistical significance.

				Motifs	
				R <u>G</u> Y <u>W</u> /W <u>R</u> C <u>Y</u> +W <u>A</u> / <u>T</u> W	
				+	-
T	AT	C			
G	AG	G	Mutations	22	7
CG	CT	CC			
ATGT	TCT	ACC	Positions	6	12
GA	TCAT	A			
GA	<u>T</u> A	<u>G</u> C			
<u>W</u> A	<u>R</u> G <u>Y</u> W	<u>W</u> A			
<u>T</u> W	<u>W</u> R <u>C</u> Y	<u>T</u> W			
				Fisher exact test P = 0.006	
				(non-random targeting of	
				mutations to mutable motifs)	

Webmutation - Microsoft Internet Explorer

Indirizzo http://www.itb.cnr.it/webmutation/wwwmsa_ex.html

Webmutation server

MUTATIONAL SPECTRUM ANALYSIS

- MSA Program
- MSA Example
- MSA Description

MSA - Mutational Spectrum Analysis - Example

Select the desired values from the list and press the "Submit" button.

Enter Spectrum Name:

Distribution: ☒ Binomial ☐ Poisson

Mail To : ☒ No ☐ Yes

Press to analyze the data, to reset form

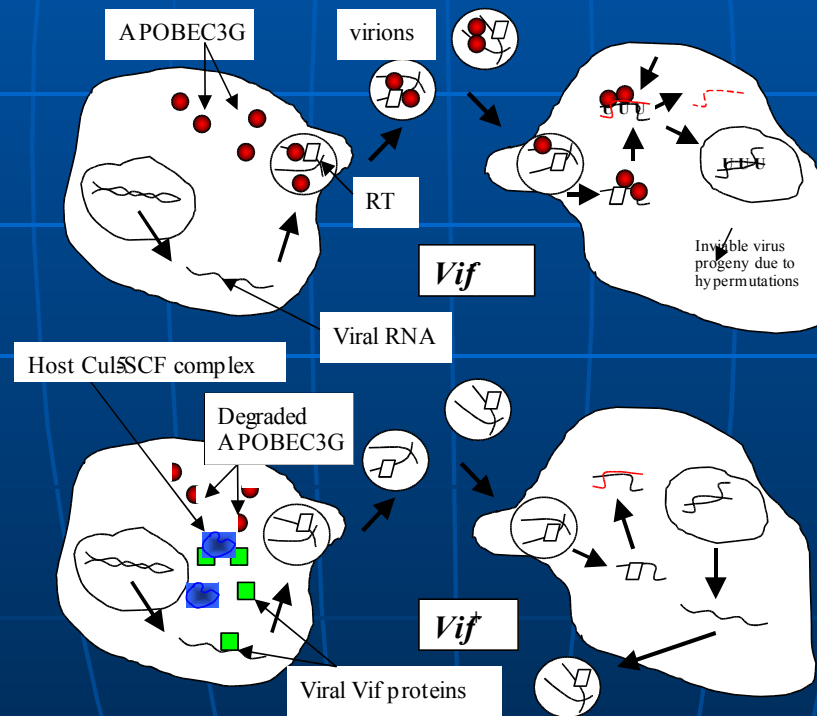
1. Number of mutation in a site (X)	<input type="text" value="0"/>	Number of sites with (X) mutations	<input type="text" value="6"/>
2. Number of mutation in a site (X)	<input type="text" value="1"/>	Number of sites with (X) mutations	<input type="text" value="3"/>
3. Number of mutation in a site (X)	<input type="text" value="2"/>	Number of sites with (X) mutations	<input type="text" value="4"/>
4. Number of mutation in a site (X)	<input type="text" value="3"/>	Number of sites with (X) mutations	<input type="text" value="1"/>
5. Number of mutation in a site (X)	<input type="text" value="4"/>	Number of sites with (X) mutations	<input type="text" value="2"/>
6. Number of mutation in a site (X)	<input type="text" value="0"/>	Number of sites with (X) mutations	<input type="text" value="0"/>
7. Number of mutation in a site (X)	<input type="text" value="0"/>	Number of sites with (X) mutations	<input type="text" value="0"/>
8. Number of mutation in a site (X)	<input type="text" value="0"/>	Number of sites with (X) mutations	<input type="text" value="0"/>
9. Number of mutation in a site (X)	<input type="text" value="8"/>	Number of sites with (X) mutations	<input type="text" value="2"/>
10. Number of mutation in a site (X)	<input type="text" value="9"/>	Number of sites with (X) mutations	<input type="text" value="1"/>
11. Number of mutation in a site (X)	<input type="text" value="10"/>	Number of sites with (X) mutations	<input type="text" value="1"/>
12. Number of mutation in a site (X)	<input type="text" value="11"/>	Number of sites with (X) mutations	<input type="text" value="1"/>
13. Number of mutation in a site (X)	<input type="text" value="12"/>	Number of sites with (X) mutations	<input type="text" value="1"/>
14. Number of mutation in a site (X)	<input type="text" value="0"/>	Number of sites with (X) mutations	<input type="text" value="0"/>
15. Number of mutation in a site (X)	<input type="text" value="0"/>	Number of sites with (X) mutations	<input type="text" value="0"/>

The Webmaster: milan@itb.cnr.it

Internet

Hypermutation in HIV-1

The cytidine deaminase APOBEC3G confer resistance to HIV. Its antiviral action could be overcome by the presence of virion infectivity factor (Vif), encoded by the viral genome.



Hypermutation in HIV-1

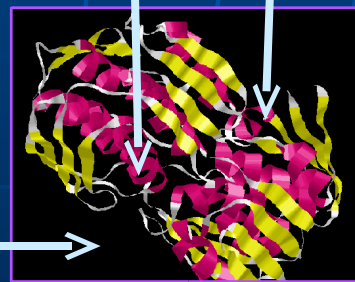
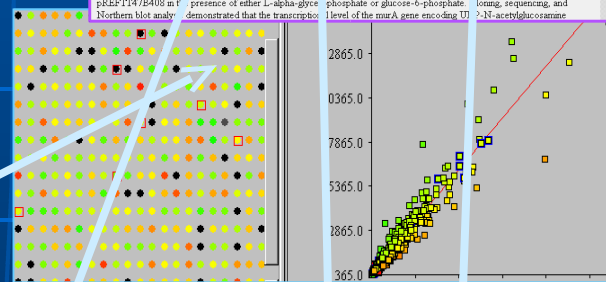
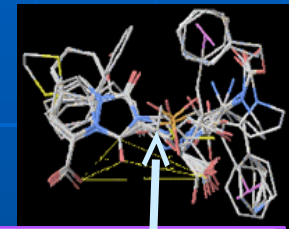
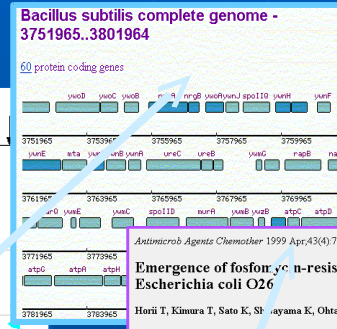
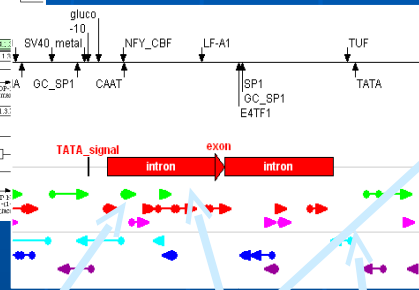
Specificity of APOBEC3G for GG sequences, which is frequently a part of TGG tryptophan codons, results in a frequent generation of TAG nonsense codons which leads to a premature termination of protein synthesis.

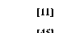
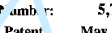
The cytidine deaminase APOBEC3G may cause mutations in HIV-1

W
... TGG ...



*
... TAG ...

[illegible]

			
United States Patent		5,753,483	
Elbein		May 19, 1998	
PERFECTED HOMOGENEOUS UDP-GLCNAC (GALNAC) PYROPHOSPHORYLASE		References Cited	
Inventor: Alan D. Elbein, Little Rock, Ark.		U.S. PATENT DOCUMENTS	
Assignee: University of Arkansas, Little Rock, Ark.		4,813,349 8/1986 Smet et al. 435/518 5,461,701 1/1994 Kiek et al. 352/8	
Appl. No.: 437,146		Primary Examiner—Blaine Lankford Attorney Agent, or Firm—Benjamin A. Orr, Dallas	
Filed: May 8, 1995		ABSTRACT	
Int. Cl.⁶ C12N 916; C12N 9/14		The present invention provides the enzyme UDP-N-acetylglucosaminide pyrophosphorylase in a purified and	
U.S. Cl. 435/196; 435/194		parted form. Also provided is a various method of using and	
Field of Search: 435/196; 233, 435/194		preparing this purified, homogeneous enzyme.	
3 Claims, 9 Drawing Sheets			

[illegible]



- Home Page
- Library details
- **Processing, Assembly & Annotation Protocol**
- Statistics on Sequences
- Sequence report
- Contig report
- SNPs
- Text Search
- Blast Search
- Links
- References
- Download
- DB contents and help

PROCESSING, ASSEMBLY AND ANNOTATION PROTOCOL

A fully automated pipeline has been prepared to process EST sequences using public software integrated by in-house developed Perl scripts.

Sequence files, produced by Dr. Pozzi's lab (libraries S3 and S4), together with the quality files produced by the base-calling program **phred** (Ewing et al, 1998) have been processed in order to identify vector contamination and low quality regions with the program **Lucy** (Chou et al, 2001) using default parameters.

Vector-free high quality sequences have been added to peach EST sequences produced by the groups of Padua (libraries Pp-S4 and S3II) and by the university of Clemson (library PP_Lea) and submitted to the program **CAP3** (Huan et al, 1999) to perform contig assembly. Stringency parameters have been modified (-p 95, -d 60) to identify paralogs.

All the input EST sequences and all the contig consensus sequences have been submitted to the **BLASTx** program for annotation (Altschul et al, 1990). BLASTx compares the six-frame conceptual translation products of a nucleotide query sequence against a protein sequence database. BLASTx is run locally against the Genbank nr protein database (nr contains all non-redundant GenBank CDS translations + RefSeq Proteins + PDB + SwissProt + PIR + PRF). Blast output has been parsed and stored in a MySQL database together with other data produced in the intermediate steps of the pipeline.

The ESTree DB Unigene dataset is defined as the collection of all the non-redundant sequences present in the DB and is derived from CAP3 output marking as Unigene all the singleton sequences and the longest sequence of each contig.

SNP detection has been performed with the program AutoSNP version 9 (Baker et al, 2003). Tgid parameters have been set to -p 95, -l 60, -v 20 for appropriate EST clustering.

A php-based web interface has been prepared to query the database. Users can view **sequence** data, **BLAST** outputs, **contig** alignments, and global **statistics** on sequences.

Single **sequence** and **contig** report pages are available, and fasta formatted sequences can be copied from here. Sequence quality files are available upon request.

A **text search** utility is available and queries can be performed on the **sequence** and **contig** report tables. BLAST **E-value** intervals can also be selected by users.

Local **BLASTn**, **BLASTp**, **BLASTx**, **tBLASTn** and **tBLASTx** programs are also available to perform BLAST and batch BLAST searches on the ESTree database.

In the web pages that present data in tabular form an EXCEL **spreadsheet** view is available, to retrieve and copy data in EXCEL format.

In the SNP report page SNP data can be accessed by sequence GI number, by sequence name or by SNP number.

Complete sequence and contig download is allowed from the Download page both in multifasta , CSV and NCBI format.

- Home Page
- Library details
- Processing, Assembly & Annotation Protocol
- Statistics on Sequences
- Sequence report
- Contig report
- SNPs
- Text Search
- Blast Search
- Links
- References
- Download
- DB contents and help

CONTIG

Contig name: Contig19

Contig length: 1311

Sequences in this contig: Pp-S4_EST1698, PP_LEa0030J08f, PP_LEa0011B19f, PP_LEa0021I03f, S410N4, S319G11, PP_LEa0028N23f, PP_S3IIsel_A2_C10, S37D16, S411L16, S311K5, PP_S3IIall_A6_C05

Blast Output

Best Blast Hit: BAD33762.1

E-value: 1e-102

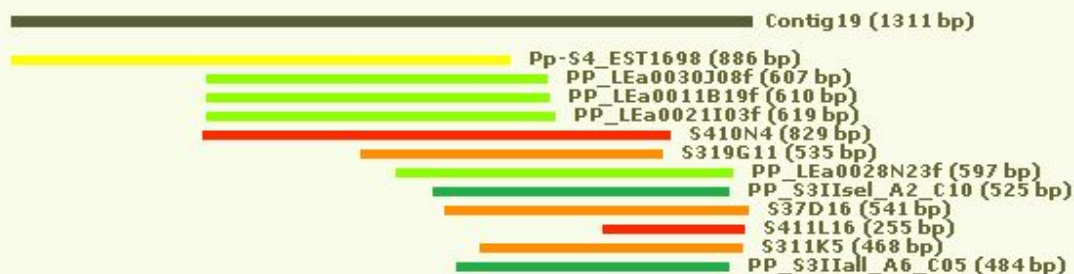
Annotation: putative 6-phosphogluconolactonase [*Oryza sativa* (japonica cultivar-group)]

Source: *Oryza sativa* (japonica cultivar-group)...

View alignment details

Libraries and developmental stages:

- S3
- S4
- PP_Lea
- Pp_S4
- S3II



>Contig19 length=1311

```
ATCGGTCGCGCTACNTCTCGCGCGCGTTGATCTCACTCCGAGTGTCACCAATATTCCAA
GCTCTCATCGATTAAATCCGTTTCGCACCAATTTCAAGAACCCAAATTTCTGCATCGGCAT
CAGCATCAGCTAGGATGGCGGGTCAGAATAAGAAGATAGAGAAGTTTGAGACGGAGGAGG
AAGTGGCGGTGCGTTTGGCCAAGTACACCGCAGATCTGTCCGCTAAGTTNGTGAAAGAGA
```



eSTuberdb

- [Home Page](#)
- [Library details](#)
- [Processing, Assembly & Annotation Protocol](#)
- [Statistics on Sequences](#)
- [Sequence report](#)
- [Contig report](#)
- [Text Search](#)
- [Blast Search](#)
- [Links](#)
- [Download](#)
- [DB contents and help](#)

Welcome to ESTuber

A data bank of Est sequences developed at the Istituto Biologia e Biotecnologia Agraria in the frame of the CNR Strategic Project TUBER: biotecnologia della micorrizzazione.

The ESTuber Consortium:

ESTuber is a Consortium of several research centers in Italy devoted to the implementation of genomics and functional genomics in truffle species.

The primary objectives of the Consortium are:

the development of an extensive EST database for Tuber species (ESTuber DB);
the analysis of several biochemical pathways based on oligonucleotide microarrays derived from ESTs collection;
the analysis of conserved genetic modules in truffle species and related filamentous fungi.

The current members of the Consortium are:

Istituto Biologia e Biotecnologia Agraria, IBBA-CNR, Milano
Istituto Genetica Vegetale, IGV-CNR, Section Perugia
Istituto per la Protezione delle Piante, IPP- CNR, Section Torino
Istituto di Biologia e Patologia Molecolari, IBPM-CNR, Roma
Dipartimento di Biochimica e Biologia Molecolare Università degli Studi di Parma
Istituto di Chimica Biologica, Centro di Biochimica delle Proteine, Università di Urbino

The ESTuber DB:

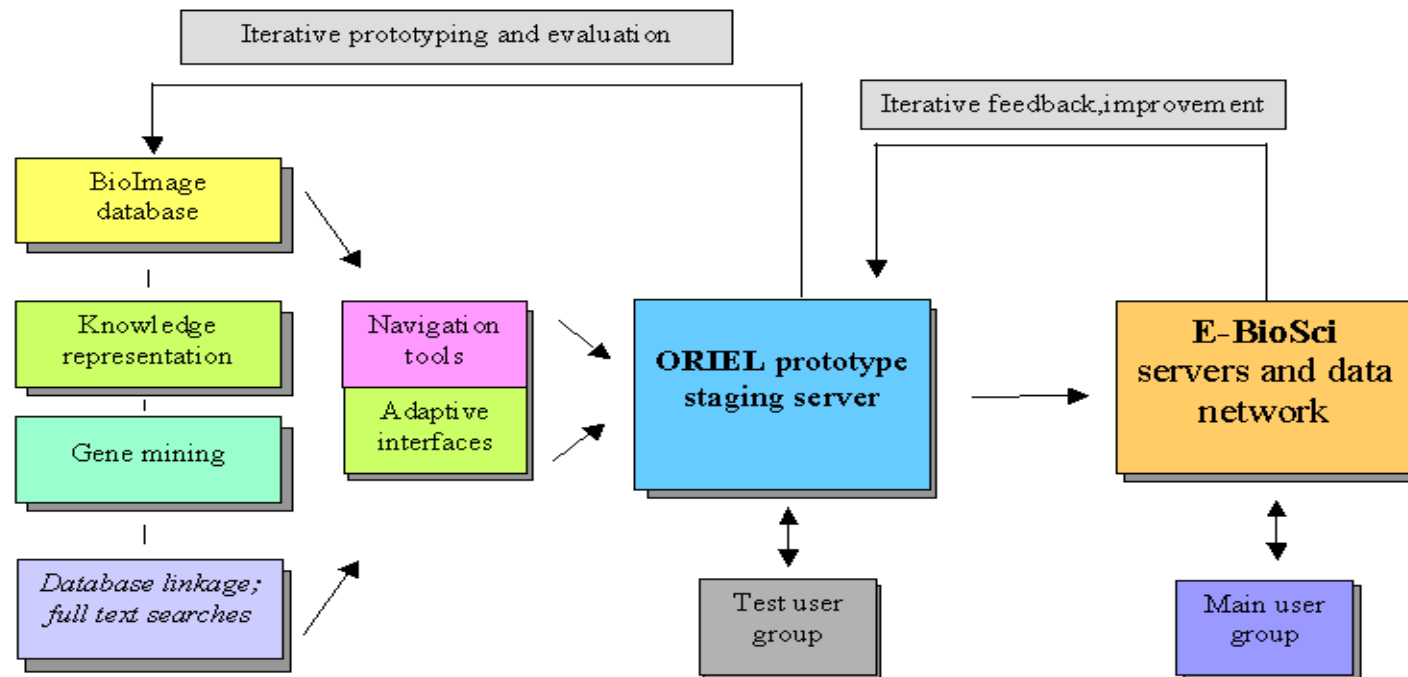
The ESTuber DB is a collection of *Tuber borchii* EST sequences.

The ESTuber DB has been developed by:

Dr. Barbara Lazzari and Dr. Andrea Caprera: DB designers and webmasters
Dr. Cristian Cosentino and Dr. Barbara Lazzari: cDNA library production and sequences management
Dr. Luciano Milanesi: Bioinformatics
Dr. Angelo Viotti: ESTuber project manager

Italy

ORIEL & E-BIOSCI



Or / e l

GRID Networked data

Laboratory for Interdisciplinary Technologies in Bioinformatics

The GRID: networked data processing centres and "middleware" software as the "glue" of resources.

Researchers perform their activities regardless geographical location, interact with colleagues, share and access data

Scientific instruments and experiments provide huge amount of data



NETTAB 2005 5-7 October, 2005 Naples
Italy

LITBIO Laboratory for Interdisciplinary Technologies in BIOinformatics



<http://www.litbio.org>

Molecular medicine has progressed to the point where the majority of human genes and proteins have been characterized and computational methods are necessary for further understanding of structure, mechanism and function.

Our aim is to develop a Laboratory for Interdisciplinary Technologies in Bioinformatics (LITBIO) applied to Genomics, Transcriptomics, Proteomics and Metabolomics.