# Introduction

# GBIO0002 Archana Bhardwaj University of Liege a.bhardwaj@uliege.be

# Overview

- **1. Introduction to Bioinformatics**
- 2. Introduction to public databases
- 3. Intro to basic R

# **Bioinformatics**

**Definition 1**: the collection, classification, storage, and analysis of biochemical and biological information using computers especially as applied to molecular genetics and genomics (*Merriam-Webster dictionary*)



: Shotgun Whole-Genome Sequencing



# **Definition 2:** a field that works on the problems involving intersection of Biology/Computer Science/Statistics

									В	I	0	L	0	G	y				
					Ρ	H	Y	S	Ĵ	0	M	I	C	S					
					C	E	L	L	0	M	I	C	S						
								B	l	0	T	E	C	Н					
	E	V	0	L	U	T	Ι	0	Ν										
							I	N	F	0	T	E	C	H					
						0	N	T	0	L	0	G	У						
								P	R	0	Ť	E	0	M	I	C	S		
M	0	L	E	C	U	L	A	R	М	0	D	E	L	I	N	G			
								M	Α	Т	H	E	M	A	T	I	C	S	
							M	E	Т	A	B	0	L	0	M	I	C	S	
		T	R	A	N	S	C	R		Ρ	Т	0	M	I	C	S			
			G	E	N	0	M	I	С	S									
									S	T	A	Т	Ι	S	Т	Ι	C	S	

# What "unit of information" do we deal within bioinformatics ?

- DNA
- RNA
- Protein



- Sequence
- Structure
- Evolution



- Pathways
- Interactions
- Mutations







# Central Dogma of Molecular Biology

#### https://www.genome.gov/human-genome-project



## Human Genome- 1990-2003

The first printout of the human genome to be presented as a series of books, displayed at the <u>Wellcome Collection</u>, London



#### **Genomic information**



#### **More information :**

DNA sequence, RNA sequence, Protein sequence



#### http://humanproteomemap.org/ (Human Proteome Map (HPM)

 $\leftarrow$   $\rightarrow$  C (i) Not secure | humanproteomemap.org



Home O

Query Download

FAQs Contact us

#### About Human Proteome Map

The Human Proteome Map (HPM) portal is an interactive resource to the scientific community by integrating the massive peptide sequencing result from the draft map of the human proteome project. The project was based on LC-MS/MS by utilizing of high resolution and high accuracy Fourier transform mass spectrometry. All mass spectrometry data including precursors and HCD-derived fragments were acquired on the Orbitrap mass analyzers in the high-high mode. Currently, the HPM contains direct evidence of translation of a number of protein products derived from over 17,000 human genes covering >84% of the annotated protein-coding genes in humans based on >290,000 non-redundant peptide identifications of multiple organs/tissues and cell types from individuals with clinically defined healthy tissues. This includes 17 adult tissues, 6 primary hematopoietic cells and 7 fetal tissues. The HPM portal provides an interactive web resource by reorganizing the label-free quantitative proteomic data set in a simple graphical view. In addition, the portal provides selected reaction monitoring (SRM) information for all peptides identified.

Statistics	
otatistics	

Organs/cell types	30
Genes identified	17,294
Proteins identified	30,057
Peptide sequences	293,700
N-terminal peptides	4,297
Splice junctional peptides	66,947
Samples	85
Adult tissues	17
Fetal tissues	7
Cell types	6



Adult tissues







# **Bioinformatics Significance**

#### **RESEARCH NEWS**

's dis-

cts 17

han 2

ponsi-

Isher-

ge 40.

olecu-

of the

id the

, and

eneral

10 and

osome

aining

re re-

ted to

182. ning so

#### **Missing Alzheimer's Gene Found**

Researchers find the gene that causes Alzheimer's disease in "Volga German" families. It shows a remarkable similarity to another recently discovered Alzheimer's gene

pinpointed as the likely site of the Alzheimer's gene. "That was like a sledgehammer to the forehead," says Schellenberg. "It went from being a ho-hum project to ... saying 'oh my God this is the gene.' "

Within a few days, the team sequenced the gene from Volga German family members, with help from David Galas and his col-

> close on the heels of the chromosome 14 gene discovery," says Alzheimer's researcher Dennis Selkoe of Harvard Medical School. "It is very important that the new gene on chromosome 1 has high homology to \$182," he adds. The similarity between the two genes may mean that the proteins they encode have similar functions. According to Selkoe, the resemblance "suggests that something about this type of ... protein is very important for the biology of Alzheimer's disease."

discovery was provocative because it provided a direct link to a characteristic feature of e, has Altheimer's pathology: APP is the source of a peptide called B-amyloid that is found in the abnormal "senile plaques" that stud Alzcovery. heimer's patients' brains. But mutant APP genes turned out to account for only 2% to 3% of familial Alzheimer's cases. orm of

About a year later, several teams, including Schellenberg's, showed that many more cases of familial Alzheimer's are caused by an unknown defective gene on chromosome 14. That gene was identified earlier this year by a team led by Peter St. George-Hyslop of the

University of Toronto; the results were reported in the 29 June issue of Nature.

Intriguing as these discoveries were, they left untouched one handful of Alzheimer's-carrying families, which had been identified by Thomas Bird at the Veterans Affairs Medical Center in Seattle: the socalled Volga Germans, who were all descended from a colony of ethnic Germans liv-

sequence tagged (EST) sequences, short DNA sequences known to come from active genes. Wasco found an EST with a sequence similar to \$182, Tanzi recalls, and said, "maybe this is the Volga German gene."

After the S182 sequence was published, Tangi and Wasco told Schellenberg about Wasco's idea. "Having seen a zillion candidates [for the Volga German gene] come and go, I wasn't excited," Schellenberg recalls. But Ephrat Levy-Lahad, in his lab group, went ahead and checked. She found that the new gene was not only on chromosome 1, but was in the very stretch of DNA that she had



Family resemblance. Mutations in the similar proteins made by the genes S182 and STM2 cluster around the membrane-spanning regions.

# Changes in the number and order of genes (A-D) create genetic diversity within and between populations.



# Why do we need DATABASES ?



# Genome sequencing generates lots of data



# DATABASES



# What are Biological Databases??

#### **Biological Database**

- It is a collection of data that is structured, searchable, updated periodically and cross-referenced.
- Stores biological data in electronic form.
- · Purpose-
- Systemization of database
- Availability of biological data
- Analysis of computed biological data

#### Features of Biological

#### Databases

- 1. Heterogeneity
- 2. High volume data
- 3. Uncertainity
- 4. Data curation
- 5. Data integration
- 6. Data sharing
- 7. Dynamics

# **DATABASE ARCHITECTURE**



# **Types of Biological Databases??**

There are many different types of database but for routine sequence analysis, the following are initially the most important.

Primary databases
 Secondary databases
 Composite databases



# **Interconnections between Databases**



# **Primary Databases**

Theses are the primary sources of data used to store nucleic acid, protein sequences and structural information of biological macromolecules.

Some primary databases-

- NCBI(The National Centre for Biotechnology Information)
  - GenBank
  - DDBJ (DNA data bank of Japan)
- SWISS-PROT(Swiss-Prot)
- PIR (Protein Information Resource)
- PDB(Protein Data Bank)

This sequence collection of this database is due to the efforts of basic research from academic industrial and sequencing lab)

# **Classification : Primary Databases**

- ✓ Sequence Information
  - ✓ DNA: EMBL, Genbank, DDBJ
  - ✓ Protein: SwissProt, TREMBL, PIR, OWL
- ✓ Genome Information
  - ✓ GDB, MGD, ACeDB
- ✓ Structure Information
   ✓ PDB, NDB, CCDB/CSD

# The National Center for Biotechnology Information





#### Created in 1988 as a part of the National Library of Medicine at NIH

- Establish public databases
- Research in computational biology
- Develop software tools for sequence analysis
- Disseminate biomedical information

# **Primary Databases - GenBank**

# Database from NCBI, includes sequences from publicly available resources

S NCBI Resources	How To 🕑		
GenBank	Nucleotide 🗸	Search	
GenBank 🔻 Submi	✓ Genomes ▼ WGS ▼ Metagenomes ▼ TPA ▼ TSA ▼ INSDC ▼ Other ▼		
GenBank Overvie	GenBank Resources		
What is GenBank?		GenBank Home	
GenBank <sup>®</sup> is the NIH ge Research, 2013 Jan;41(E	etic sequence database, an annotated collection of all publicly available DNA sequences ( <u>Nucleic Acids</u> 1):D36-42). GenBank is part of the International Nucleotide Sequence Database Collaboration, which comprises	Submission Tools	
the DNA DataBank of Ja	an (DDBJ), the European Nucleotide Archive (ENA), and GenBank at NCBI. These three organizations exchange	Search GenBank	
data on a daily basis.		Update GenBank Records	
A GenBank release occu	s every two months and is available from the <u>ftp site</u> . The <u>release notes</u> for the current version of GenBank		

provide detailed information about the release and notifications of upcoming changes to GenBank. Release notes for <u>previous GenBank</u> releases are also available. GenBank growth statistics for both the traditional GenBank divisions and the WGS division are available from each release. GenBank growth <u>statistics</u> for both the traditional GenBank divisions and the WGS division are available from each release.

An <u>annotated sample GenBank record</u> for a Saccharomyces cerevisiae gene demonstrates many of the features of the GenBank flat file format.

### ✓ Open « Gene » and Search KRAS

S NCBI Resources	🖸 How T	'o 🕑					
Gene	Gene	∽ K	RAS reate RSS Create alert A	Advanced			× 😒 Search
<b>Gene sources</b> Genomic Mitochondria		Tabular - 20 pe	er page - Sort by Relevance			Send to: 🗸	Filters: <u>Manage Filters</u>
Organelles Categories Alternatively spliced		See <u>KRAS K</u> kras in <u>Homo</u>	RAS proto-oncogene, GTF sapiens Mus musculus Ra	Pase in the Gene database attus norvegicus All 238 Gene	records		Results by taxon Top Organisms [Tree]
Annotated genes Non-coding Protein-coding Pseudogene		Search resul Items: 1 to 20 o See also 16 o	ts of 1257 discontinued or replaced iter	<< First	< Prev Page 1 of 63 Next	> Last >>	Homo sapiens (755) Mus musculus (134) Rattus norvegicus (14) Cricetulus griseus (8) Xenopus laevis (7)
Sequence content		Name/Gene ID	Description	Location	Aliases	MIM	All other taxa <i>(339)</i> More
Ensembl RefSeq RefSeqGene Status	clear	☐ <u>KRAS</u> ID: 3845	KRAS proto-oncogene, GTPase [ <i>Homo sapiens</i> (human)]	Chromosome 12, NC_000012.12 (2520478925251003, complement)	C-K-RAS, CFC2, K-RAS2A, K-RAS2B, K-RAS4A, K-RAS4B, K-Ras, KI-RAS1, KRAS2, NS, NS3, RALD, RASK2, c-Ki-ras2, KRAS	190070	Find related data Database: Select Find items
<u>Clear all</u> Show additional filters		☐ <u>Kras</u> ID: 16653	Kirsten rat sarcoma viral oncogene homolog [ <i>Mus musculus</i> (house mouse)]	Chromosome 6, NC_000072.6 (145216699145250291, complement)	Al929937, K-Ras, K-Ras 2, K-ras, Ki-ras-2, Kras2, c-K-ras, c-Ki-ras, p21B, ras, Kras		Search details

ocation: 12p12.1 con count: 6					See	KRAS in <u>Genome Dat</u>	<u>i View</u>
Annotation release	Status	Assembly	Chr	Location			
09	current	GRCh38.p12 (GCF_000001405.38)	12	NC_000012.12 (2520478925251	003, complement)		
05	previous assembly	GRCh37.p13 (GCF_000001405.25)	12	NC_000012.11 (2535818025403	870, complement)		
Genomic regions, tran	scripts, and products 00012.12 Chromosome 12 Reference	e GRCh38.p12 Primary Assembly 🖂			Go	to <u>reference</u> quenc	
Genomic regions, tran	scripts, and products	e GRCh38.p12 Primary Assembly ∨			Go Go to nucleotide: (	to <u>reference</u> quence Graphics <u>FASTA</u>	ienBa
Genomic regions, tran nomic Sequence: NC_0	d:	e GRCh38.p12 Primary Assembly ∨ ⇒ Q Q Q Q Q Q Q Q Q Q Q Q Q Q Q Q Q Q Q	25,225 K	25,220 K	Go Go to nucleotide: ( 25,210 K	to <u>reference</u> quence <u>Graphics</u> <u>FASTA</u> (C cools •   ( Tracks • ) (	ienBa

Format Homo sapiens chromosome 12, GRCI	h38.p12 Primary Assembly
NCBI Reference Sequence: NC_000012.12	
FASTA Graphics	
LOCUS NC_000012 46215 bp DNA	linear CON 26-MAR-2018
DEFINITION Homo sapiens chromosome 12, GRCh38.p12 Prim	ary Assembly.
ACCESSION <u>NC_000012</u> REGION: complement(252047892525	1003)
Accession – DBLINK BioProject: PRINA168	
Assembly: GCF_000001405.38	
Key Identifier KEYWORDS RefSeq.	
SOURCE Homo sapiens (human)	
Spocios	tebrata: Euteleostomi:
Mammalia; Eutheria; Euarchontoglires; Prima	tes; Haplorrhini;
Catarrhini; Hominidae; Homo.	
REFERENCE 1 (bases 1 to 46215)	P (noo A Ding V
Dugan-Rocha.S., Gill.R., Gunaratne.P., Harr	is.R.A., Hawes.A.C.,
Hernandez,J., Hodgson,A.V., Hume,J., Jackso	n,A., Khan,Z.M.,
Kovar-Smith,C., Lewis,L.R., Lozado,R.J., Me	tzker,M.L.,
Milosavljevic,A., Miner,G.R., Montgomery,K. Nazareth IV Scott G. Sodergren E. Song	T., Morgan,M.B.,
Lovering, R.C., Wheeler, D.A., Worley, K.C., Y	/uan,Y., Zhang,Z.,
Adams,C.Q., Ansari-Lari,M.A., Ayele,M., Bro	wn,M.J., Chen,G.,
Chen,Z., Clerc-Blankenburg,K.P., Davis,C.,	Delgado,O., Dinh,H.H.,
Draper,H., Gonzalez-Garay,M.L., Havlak,P.,	Jackson,L.R.,
Maheshwari,M., Nguyen,B.V., Okwuonu,G.O., P	asternak,S., Perez,L.M.,
Plopper,F.J., Santibanez,J., Shen,H., Tabor	,P.E., Verduzco,D.,
Waldron,L., Wang,Q., Williams,G.A., Zhang,J	., Zhou,J., Allen,C.C.,
Amin,A.G., Anyalebechi,V., Balley,M., Barba Bryant N.P., Burch P.E., Burkett C.E., Burr	rla,J.A., Bimage,K.E.,
Cardenas,V., Carter,K., Casias,K., Cavazos,	I., Cavazos,S.R.,
Ceasar,H., Chacko,J., Chan,S.N., Chavez,D.,	Christopoulos,C.,
Chu,J., Cockrell,R., Cox,C.D., Dang,M., Dat	horne,S.R., David,R.,
Eaves.K.A., Egan.A., Emery-Cohen.A.J., Esco	tto.M., Flagg.N.,
Forbes,L.D., Gabisi,A.M., Garza,M., Hamilto	n,C., Henderson,N.,
Hernandez, O., Hines, S., Hogues, M.E., Huang,	M., Idlebird,D.G.,
Johnson, R., Jolivet, A., Jones, S., Kagan, R.,	King,L.M., Leal,B.,
Lovensubewa.L.M., Louiseged.H., Lovett.D.A.	. Lucier.A
Lucier,R.L., Ma,J., Madu,R.C., Mapua,P., Ma	rtindale,A.D.,
Martinez,E., Massey,E., Mawhiney,S., Meador	,M.G., Mendez,S.,

	##Genome-Annotation-Data-END##	
FEATURES	Location/Qualifiers	
source	146215	
	/organism="Homo sapiens"	
	/mol_type="genomic DNA"	
	/db_xref="taxon: <u>9606</u> "	
	/chromosome="12"	
gene	146215	
	/gene="KRAS"	
	/gene_synonym="C-K-RAS; c-Ki-ras2; CFC2; K-Ras; K-RAS2A;	
	K-RAS2B; K-RAS4A; K-RAS4B; KI-RAS; KRAS1; KRAS2; NS; NS3;	
	RALD; RASK2"	
	/note="KRAS proto-oncogene, GTPase; Derived by automated	
	computational analysis using gene prediction method:	
	BestRefSeq,Gnomon."	
	/db_xref="GeneID: <u>3845</u> "	
	/db_xref="HGNC: <u>HGNC:6407</u> "	
	/db_xref="MIM: <u>190070</u> "	
mRNA	join(1240,56095730,2359223770,2523125390,	
	3544435567,4109341179)	
	/gene="KRAS"	
	/gene_synonym="C-K-RAS; c-Ki-ras2; CFC2; K-Ras; K-RAS2A;	
	K-RAS2B; K-RAS4A; K-RAS4B; KI-RAS; KRAS1; KRAS2; NS; NS3;	
	RALD; RASK2"	
	/product="KRAS proto-oncogene, GTPase, transcript variant	
	X1"	
	/note="Derived by automated computational analysis using	
	gene prediction method: Gnomon. Supporting evidence	
	includes similarity to: 3 mRNAs, 1 long SRA read, 13	
	Proteins, and 100% coverage of the annotated genomic	
	feature by RNAseq alignments, including 39 samples with	
	support for all annotated introns"	
	/transcript_id="XM_006/19069.4"	
	(db_xret= denerb:3645	
	(db_xpet="HGNC:6407"	
mPNIA	100_XTET= MIM: 1900/0 1010(50 240 5600 5720 22502 22770 25221 25200	
	11003 45750)	
	(general VDAS"	
	(gene skappyme"C-K-PAS; c-Ki-pas2; CEC2; K-Pas; K-PAS2A;	
	$V_{\rm PAS2R}$ , $V_{\rm PAS4A}$ , $V_{\rm PAS2R}$ , $V_{\rm$	
	RAID: RASK2"	
	/product="KRAS proto-opcogene GTPase transcript variant	
	xy"	
	/note="Derived by automated computational analysis using	
	gene prediction method: Gnomon, Supporting evidence	
	includes similarity to: 6 mRNAs, 234 ESTs, 539 long SRA	
	reads, 18 Proteins, and 97% coverage of the annotated	
	genomic feature by RNAseg alignments, including 60 samples	
	with support for all annotated introns"	
	/transcript_id="XM_011520653.3"	
	/db_xref="GeneID:3845"	
	/db_xref="HGNC:6407"	
	/db_xref="MIM:190070"	
mRNA	join(73253,56095730,2359223770,2523125390,	

#### FASTA -

#### Homo sapiens chromosome 12, GRCh38.p12 Primary Assembly

NCBI Reference Sequence: NC\_000012.12

GenBank Graphics

>NC\_000012.12:c25251003-25204789 Homo sapiens chromosome 12, GRCh38.p12 Primary Assembly Header stars with ">" sign

GGAACGCATCGATAGCTCTGCCCTCTGCGGCCCGGCCCCGAACTCATCGGTGTGCTCGGAGCTCGAT CGGCGGGGCCAGAGGCTCAGCGGCTCCCAGGTGCGGGAGAGAGGGACCGGCGGACCACCCCTCCTGGGC AGGCTTCTGGGGAGAAACTCGGGCCGGGCCGGCTGCCCCTCGGAGCGGTGGGGGTGCGGTGGAGGTTACTC CCGCGGCGCCCCGGCCTCCCCCTCTCCCCGCTCCCGCACCTCTTGCCTCCCTTTCCAGCACTCGG CTGCCTCGGTCCAGCCTTCCCTGCTGCATTTGGCATCTCTAGGACGAAGGTATAAACTTCTCCCTCGAGC GCAGGCTGGACGGATAGTGGTCCTTTTCCGTGTGTAGGGGATGTGTGAGTAAGAGGGGAGGTCACGTTTT GGAAGAGCATAGGAAAGTGCTTAGAGACCACTGTTTGAGGTTATTGTGTTTGGAAAAAAATGCATCTGCC TCCGAGTTCCTGAATGCTCCCCCCCCCCATGTATGGGCTGTGACATTGCTGTGGCCACAAAGGAGGAGGT GGAGGTAGAGATGGTGGAAGAACAGGTGGCCAACACCCTACACGTAGAGCCTGTGACCTACAGTGAAAAG GAAAAAGTTAATCCCAGATGGTCTGTTTTGCTTGGTCAAGTTAAACCCGAAGAAAACCCGCAGAGCAGAA GCAAGGCTTTTTCCTTGCTAGTTGAGTGTAGACAGCAATAGCAAAAATAGTACTTGAAGTTTAATTTACC TGTTCTTGTCCTTTCCCCTATTTCTTATGTATTACCCCTCATCCCCTCGTCTCTTTTATACTACCCCTCATT TTGCAGATGTGTTCTACATCTCAAGAGTTATTACAGTACTCCAAAACAGCACTTACATGATTTTTTAAAC TTACAGAGGAATTGTAGCAATCCACCAGCTAACCGCCTGAAATAGACTTAAACATGTGCATCTCCTTTT TTTTTTTTTTTGAGACACAGTCTCGCTCTGTTGCCCAGGCTGGAGTGCAATGGCGCGGTATCGGCTCAC TGAAACCTCCGCCTCCTGGGTTCAAGCAATTCTCCTGCCTCAGCCTCCCGAGTAGCTGGGACTAGTAGGT GCACGCCACCATGCCCAGCTAATTTTTGTATTTTTAGTAGAGACAGAGTTTCATCATGTTGGTCAGGATG CTGCATTCAAGCAATTCTCCTGCCTCAGCCTCCCGAATAACTGGGATTACAGGTGTCTGCCGCCATGCCC GGCTAATTTTTTGTATTTTTAGTAGAGAGAGGGGGTTTCACCATGTTGGTCAGGCTGGTCTAGAACTCCTG

The FASTA format is now universal for all databases and software that handles
DNA and protein sequences
Specifications:
One header line

•starts with > with a ends with [return]



Search '6Q6I' : Lysine decarboxylase A from Pseudomonas aeruginosa Classification: OXIDOREDUCTASE (type) Organism(s): Pseudomonas aeruginosa Expression System: Escherichia coli

https://www.rcsb.org/

# **OMIM** database

- Online Mendelian Inheritance in Man (OMIM)
- "information on all known mendelian disorders linked to over 12,000 genes"
- "Started at 1960s by Dr. Victor A. McKusick as a catalog of mendelian traits and disorders"
- Linked disease data
- Links disease phenotypes and causative genes
- Used by physicians and geneticists



# **OMIM-search results**

• Look for the entires that link to the genes. Apply filters if needed



# **OMIM-entries**



#### Description

Shondwloarthronathy (ShA) one of the commonest chronic rheumatic diseases includes a spectrum of related

# **OMIM Gene ID -entries**



TEXT

For background information on the major histocompatibility complex (MHC) and human leukocyte antigens
### **OMIM-Finding disease linked genes**

#### Mapping

Gu et al. (2009) conducted a genomewide scan followed by fine mapping analysis in a 4-generation Han Chinese family with ankylosing spondylitis and obtained a maximum lod score of 4.02 at D6S273 (theta = 0.0) on chromosome 6, verifying the HLA-B locus.

#### Linkage Heterogeneity

To identify major loci controlling clinical manifestations of AS, Brown et al. (2003) performed genomewide linkage analysis on 188 affected sib-pair families containing 454 affected individuals. Heritabilities of the traits studied were as follows: age at symptom onset, 0.33 (p = 0.005); disease activity assessed by the Bath Ankylosing Spondylitis Disease Activity Index (BASDAI), 0.49 (p = 0.0001); and functional impairment assessed by the Bath Ankylosing Spondylitis Functional Index (BASFI), 0.76 (p = 0.000001). No linkage was observed between the MHC and any of the traits studied. Significant linkage (lod = 4.0) was observed between a region on chromosome 18p and the BASDAI. Age at symptom onset showed suggestive linkage to chromosome 11p (lod = 3.3). Maximum linkage with the BASFI was seen at chromosome 2q (lod = 2.9; see SPDA3, new). Brown et al. (2003) concluded that these clinical manifestations are largely determined by a small number of genes not encoded within the MHC.

In a multistage study involving 12,701 SNPs and patients with autoimmune diseases, including ankylosing spondylitis, the Wellcome Trust Case Control Consortium and the Australo-Anglo-American Spondylitis Consortium (2007) identified significant association with SNPs in the ARTS1 gene (ERAP1; 606832) (combined results,  $p = 1.2 \times 10(-8)$  to  $3.4 \times 10(-10)$ ) on chromosome 5q15. Association was also found with SNPs in the IL23R gene (607562) on chromosome 1p31.3: in combined analysis, the strongest association was at rs11209032 (odds ratio, 1.3;  $p = 7.5 \times 10(-9)$ ). The association remained strong when only individuals who self-reported as not having inflammatory bowel disease (see IBD17, 612261) were considered, and was still strongest at rs11209032 ( $p = 6.9 \times 10(-7)$ ).

### **Secondary Databases**



### **Secondary Database : PROSITE**

### ✓ Open link <u>https://prosite.expasy.org/</u>



Database of protein domains, families and functional sites

PROSITE consists of documentation entries describing protein domains, families and functional sites as well as associated patterns and profiles to identify them [More... / References / Commercial users].

PROSITE is complemented by ProRule, a collection of rules based on profiles and patterns, which increases the discriminatory power of profiles and patterns by providing additional information about functionally and/or structurally critical amino acids [More...].

Release 2018\_08 of 12-Sep-2018 contains 1814 documentation entries, 1309 patterns, 1222 profiles and 1245 ProRule.

Search	Browse
e.g. PDOC00022, PS50089, SH3, zinc finger Search	<ul> <li>by documentation entry</li> <li><u>by ProRule description</u></li> <li>by taxonomic scope</li> <li>by number of positive hits</li> </ul>



### **Primary vs Secondary Databases**



### **Composite Databases**

- Collection of various primary
   Renders sequence searching databases sequences
  - highly efficient as it searches multiple resources



### **Other Databases**



### PubMed database

- <u>PubMed</u> is one of the best known database in the whole scientific community
- Most of biology related literature from all the related fields are being indexed by this database
- It has very powerful mechanism of constructing search queries
  - Many search fields Logical operators (AND, OR)
- Provides electronic links to most journals
- Example of searching by author articles published within 2012-2013

```
Search results
Items: 11
PLANET-SNP pipeline: PLants based ANnotation and Establishment of True SNP pipeline
1. Bhardwaj A, Bag SK.
    Genomics. 2019 Sep;111(5):1066-1077. doi: 10.1016/j.ygeno.2018.07.001. Epub 2018 Jul 3.
    PMID: 31533899
    Similar articles
Transcriptome analysis provides insight into prickle development and its link to defense and
2. secondary metabolism in Solanum viarum Dunal.
    Pandey S, Goel R, Bhardwaj A, Asif MH, Sawant SV, Misra P.
    Sci Rep. 2018 Nov 20;8(1):17092. doi: 10.1038/s41598-018-35304-8.
    PMID: 30459319 Free PMC Article
    Similar articles
    In Silico identification of SNP diversity in cultivated and wild tomato species: insight from molecular
3. simulations.
    Bhardwai A, Dhar YV, Asif MH, Bag SK.
    Sci Rep. 2016 Dec 8;6:38715. doi: 10.1038/srep38715.
```

### **Applications of Bioinformatics : Medical Implications**

### ✓ Pharmacogenomics

- ✓Not all drugs work on all patients, some good drugs cause death in some patients
- ✓ So by doing a gene analysis before the treatment the offensive drugs can be avoided
- ✓ Also drugs which cause death to most can be used on a minority to whose genes that drug is well suited volunteers wanted!
- ✓Customized treatment
- ✓ Gene Therapy
  - ✓ Replace or supply the defective or missing gene
  - ✓ E.g: Insulin and Factor VIII or Haemophilia

### **Applications of Bioinformatics : Diagnosis of Disease**

### ✓ Diagnosis of disease

□Identification of genes which cause the disease will help detect disease at early stage e.g. Huntington disease -

- Symptoms uncontrollable dance like movements, mental disturbance, personality changes and intellectual impairment
- ✓ Death in 10-15 years
- ✓ The gene responsible for the disease has been identified
- ✓ Contains excessively repeated sections of CAG
- $\checkmark$  So once analyzed the couple can be counseled

# **Applications of Bioinformatics : Drug Design**

- ✓ Can go up to 15yrs and \$700million
- ✓One of the goals of bioinformatics is to reduce the time and cost involved with it.
- $\checkmark$  The process
  - ✓ Discovery
    - ✓ Computational methods can improves this
  - ✓ Testing

### Introduction to



A basic tutorial

### Statistical languages GUIs

WPS - test/Scriptl.sas - WPS Workbench	Strok	_survival.sav (DataSe	et2] - IBM SPSS Statistics Data Ed	litor	The second s	9.28
Eile Edit Navigate Search Project WPS Bun Window Help	Elle Edit	View Data Transfo	rm <u>Analyze</u> Direct Barketing Grap	ns itilies Add-gris Windsw H		
] 🖿 ♥ 🗑 🛆   O ≡ 😡 😡 🛃   Q ♥   A ♥   E ♥ ↓ Y ♥ ↓ ♥ ♥ ↔ ♥ ↔ ♥	Gen I	1.0.	Reports P	AA 2011 1111 1111	A men	A ABC
EI 🖉 WPS			Descriptive Statetics			1H @ @ ~0
Project Explorer 🖾 🔍 😨 Scriptl.sas 🕸 🕡 Log 🛛 🖓 Script Compatibility R 🛛 🎧 Language Usage Report 📄 🗖 😫 Outline 🖄			Tables b		and the second second	Visible 42 of 42 Variables
😑 🎭 🔽 1 <sup>5</sup> data a: 💌 💌	000	1	Commen Massa			
🛿 🗃 test 🧄 2 string="Hadoop"; "input;		brateg	Constrail Inter Hadel	gender active	obesity	diabetes bp
R DATASETS 4 len=length(string);	mat 1	9735702127	29 Deceratories water	-54 Female Yes	No	No Hypotension
COMMAND S Call symputx('len', len);     Call symputx('len', len);     Call symputx('len', len);	2	4852351830	79 Generalized Linkar Models	-74 Male Yes	Yes	No Hypertension
B DOS 7	3	3434994256	79 Majed Hodels P	-74 Female Yes	Yes	Yes Hypertension
🐵 🎃 OUTPUTS 🚽 🗧 🗧 data aa:	4	6053971728	82 Correlate P	.74 Mala Yes	No	No Normal
PROCENTS     Set as     PROCENTS     PR	ort	9370757269	29 Begression *	Astomatic Linear Vodeling.	No	No Hypertension
- C 109001 Sas	6	1617105720	ng Laginear >	E.F. Linear	Vas	tio biomal
- T09003.sas 12 STRING_master=upcase(string);		0275365230	Neural Networks	R	No	Ver Alemat
-S T09004.sas 14 *capitalization 1: Map;		0215365329	dz Classify *	Quive Estratus	NO	Tes Pagimai
- 2 T09005, sas 15 array str[4len] \$;	8	3906563332	79 Dimension Reduction +	Partial Leagt Squares	No	No Normal
C 10900, Sas 10 do 1 =1 to len;	9	4785366661	82 Scale >	Bhary Logett.	No	No Normal
-2 T09008, sas 18 end;	10	9589919145	82 Maximum Tests	Definition of a state	No	No Hypertension
-Se T09010. sas 19 20 Incentral institution 2: Reduces	11	4598012219	79 Farmanton h	and growth and any other	Yes	No Normal
TOODIL sas 21 STRING workers are strillen;	12	3629441662	79	Crimal.	No	No Normal
	13	630781658R	70	Prest .	No	No Hundansion
-2 T09014. sas 23 arop ien 17		6367060660	9.2 million filesporter #	Real Mandanear	Var	No Normal
	14	5357063059	02 Mesing Value Analyse.	TR and all the starting	res	reg regimal
26 Output Feed annu 97	15	5132742071	29 Mattole Imputation *	has medal canades.	Yes	Yes Normal
	16	2660586207	29 Complex Samples P	hill 2-Stage Loost Squares	Yes	No Hypertension
	17	5408312498	79 Quality Control *	Qutimal Scaling (CATREG)	No	No Hypertension
	18	9069087682	29 ROC Carve.	-64 M1310 T 05	No	No Hypertension
- Listing	19	8173197592	799998 58 1	55-64 Female No	No	No Normal
lie Results in a Local i ⇒ Libraries	20	8808732689	822229 83	75+ Male Yes	No	No Hypotension
* 🗊 db	21	3400443333	822220 67	65.74 Famala Yas	Var	his historial
	21	1	01 1	0074 Tendle Tes	145	
	Concession of the local division of the loca	A DESCRIPTION OF THE OWNER OWNER OF THE OWNER OWNER OF THE OWNER OF	THE REAL PROPERTY OF THE PARTY	664		
Filerofs	Oata Viev	Variable View				STRUCTURE STOCKED A STRUCT
0° & test/Scriptl.sas	Litear.			6V	SPSS Statistics Pro	cassor is ready





### **R** GUI



#### Less fancy and no frills, but free!



✓ "R is a free software environment for statistical computing and graphics"

✓ R is considered to be one of the most widely used languages amongst statisticians, data miners, bioinformaticians and others.

✓ R is free implementation of S language

✓ Other commercial statistical packages are SPSS, SAS, MatLab

## Why to learn R?

- ✓Since it is free and open-source, R is widely used by bioinformaticians and statisticians
- $\checkmark$  It is multiplatform and free
- ✓ Has wide very wide selection of additional libraries that allow it to use in many domains including bioinformatics
- ✓ Main library repositories CRAN and BioConductor

### Install R

http://www.r-project.org/

and do the following (assuming you work on a windows computer):

- click download CRAN in the left bar
- choose a download site
- choose Windows as target operation system
- click base

 choose Download R 3.0.3 for Windows <sup>+</sup> and choose default answers for all questions

### **Install RStudio**

http://www.rstudio.org/

and do the following (assuming you work on a windows computer):

- click Download RStudio
- click Download RStudio Desktop
- click Recommended For Your System
- download the .exe file and run it (choose default answers for all questions)

### **RStudio layout**

#### The RStudio interface consists of several windows

RStudio <u>File Edit View W</u> orkspace <u>Plots T</u> ools <u>H</u> elp		
testscript.R ×	Workspace History	Import Dataset + 🖌 Clear All
11 12 13 for(i in 1:3) 14 { 15 print(i) 16 }	Values a 1 b numer c chara	ric[3] acter[3]
17 18 19 4:1 € (Top Level) ≎ Console ~/ ↔	e numer i 3L	es Help
<pre>console ~/ ~~ ~~ ~~ ~~ ~~ ~~ ~~ ~~ ~~ ~~ ~~ ~~</pre>		E Thep Export • ♥ € Clear All © Clear All © 4 6 8 10 d

- Bottom left: console window (also called command window). Here you can type simple commands after the ">" prompt and R will then execute your command. This is the most important window, because this is where R actually does stuff.
- Top left: editor window (also called script window). Collections of commands (scripts) can be edited and saved. When you don't get this window, you can open it with File → New → R script

Top right: workspace / history window. In the workspace window you can see which data and values R has in its memory. You can view and edit the values by clicking on them. The history window shows what has been typed before.

 Bottom right: files / plots / packages / help window. Here you can open files, view plots (also previous plots), install and load packages or use the help function.

### **Working directory**

 Your working directory is the folder on your computer in which you are currently working.

```
setwd("C:/Users/archana/Desktop/")
```

### **Libraries**

- R can do many statistical and data analyses.
- They are organized in so-called packages or libraries.
- With the standard installation, most common packages are installed.

### **Libraries Installation**

- If you want to install and use a package (for example, the package called "geometry") you should
- Install the package:
- click install packages in the packages window and type geometry or type install.packages("geometry") in the command window.

#### RStudio

File Edit Code View Plots Session Build Debug Profile Tools Help		
📀 🔹 😪 🖙 🗣 🔚 🔚 📄 👘 Go to file/function		🔋 Project: (None) 🔻
Untitled1* ×	Environment History	Connections
<□ □ □   🚛   小 🖓 🔍   💀 Preview 👻 🕫 🔹 👘 😨 Insert 🗸   ☆ 🖓   🖶 Run 🗸   � 🕫 🖛	🚰 🔒 🐨 Import Datas	set 🕶 🚽 📃 List 🕶 📿 🕶
	🜗 Global Environment 🗸	Q
2 title: "R Notebook" 3 output: html notebook	Data	
4	🜔 adesignMat… 2 ob	s. of 2 variables
5	🜔 adesignMat 3527	6 obs. of 2 variables 🔲
6 This is an [R Markdown]( <u>http://rmarkdown.rstudio.com</u> ) Notebook. When you execute code	adjacency num	[1:202, 1:202] 1.00 9
7	🜔 aisoRepCou 4585	6 obs. of 3 variables 🔲
8 Try executing this chunk by clicking the *Run* button within the chunk or by placing your	⊙all⊤raits 18 o	bs. of 22 variables 🔲 💌
cursor inside it and pressing *Ctrl+Shift+Enter*.	Files Plots Packages	Help Viewer
4:1 📆 R Notebook 💠 R Markdown 🛊	New Folder 🕴 Delet	re Rename 👶 More 🗸 📿
Console Terminal V John V		
	A Name	Size Mod
The tollowing object is masked from package:BlocGenerics :	workplace.rdaTm	66.1 MB See
1		13 B Sep
P I OTMA		71 6 MD
Loading required package: edgeR		
<pre>&gt; install.packages("geometry")</pre>		122.3 MB Dec
Installing package into 'C:/Users/archana/Documents/R/win-library/3.6'	.RDataTmp1	7.9 MB Jul 1
trying URL 'https://cran.rstudio.com/bin/windows/contrib/3.6/geometry_0.4.4.zip'	Rhistory	22.6 KB Sep
Content type 'application/zip' length 1529464 bytes (1.5 MB)	🔲 🗋 1-enriched	1.5 KB Mar
downloaded 1.5 MB	🗌 🗋 1.txt	249 KB Mar
package 'geometry' successfully unpacked and MD5 sums checked	🔲 🗋 1.txt.bak	249 KB Mar
The deviated biogeneration in the	🗍 🗋 10.txt	62.6 KB Mar
ine downloaded binary packages are in C:\Users\archana\AppData\Local\Temp\RtmpOw1c5g\downloaded packages	18-19.txt	957 KB Jul 3
>	18-entrezz.csv	4.3 MB Mar

Load the package: check box in front of geometry or type library("geometry") in the command window.

## Variables/Operators

• Variables store one element

Here x variable is assigned value 25

• Check value assigned to the variable x

>x

### [1] 25

- Basic mathematical operators that could be applied to variables: (+),(-),(/),(\*)
- Use parenthesis to obtain desired sequence of mathematical operations

### Arithmetic operators

• What is the value of small z here?

x <- 25y <- 15z <- (x + y) \*2Z <- z\*zz[1] 80

#### **Calculator**

R can be used as a calculator. You can just type your equation in the command window after the ">":

> 10^2 + 36

#### <u>Workspace</u>

You can also give numbers a name. By doing so, they become so-called variables which can be used later. For example, you can type in the command window:

You can also ask R what a is (just type a ENTER in the command window):

> a [1] 4

or do calculations with a:

> a \* 5 [1] 20

To remove all variables from R's memory, type

> rm(list=ls())

### Vectors

 ✓ Vectors have only 1 dimension and represent enumerated sequence of data. They can also store variables

```
v1 <- c(1, 2, 3, 4, 5)
mean(v1)
[1] 3
```

✓ The elements of a vector are specified /modified with braces (e.g. [number]) v1[1] <- 48 v1

```
[1] 48 2 3 4 5
```

### Logical operators

- ✓ These operators mostly work on vectors, matrices and other data types
- ✓ Type of data is not important, the same operators are used for numeric and character data types

Description
less than
less than or equal to
greater than
greater than or equal to
exactly equal to
not equal to

### Logical operators

✓ Can be applied to vectors in the following way. The return value is either True or False

v1
[1] 48 2 3 4 5
v1 <= 3
[1] FALSE TRUE TRUE FALSE FALSE

### R workspace

# ✓ Display all workplace objects (variables, vectors, etc.) via ls():

ls() [1] "Z" "v1" "x" "y" "z"

✓ Useful tip: to save "workplace" and restore from a file use:
✓ save.image(file = "workplace.rda")
✓ load(file = "workplace.rda")

# How to find help info?

✓ Any function in R has help information

- ✓ To invoke help use ? Sign or help():
  - ? function\_name()

```
? mean
```

```
help(mean, try.all.packages=T)
```

- ✓ To search in all packages installed in your R installation always use try.all.packages=T in help()
- ✓ To search for a key word in R documentation use help.search():

help.search("mean")

### Basic data types

- ✓ Data could be of 3 basic data types:
  - √numeric
  - ✓ character
  - ✓ logical
- ✓ Numeric variable type:

x <- 1 mode(x) [1] "numeric"

### Basic data types

✓ Logical variable type (True/False):

y <- 3<4 mode(y) [1] "logical"

✓ Character variable type:
 z <- "Hello class"</li>
 mode(z)
 [1] "character"

### Data structures

✓ The main data objects in R are:

- ✓ Matrices (single data type)
- ✓ Data frames (supports various data types)
- ✓ Lists (contain set of vectors)
- ✓ Other more complex objects

✓ Matrices are 2D objects (rows/columns)

✓ Lists contain various vectors. Each vector in the list can be accessed by double braces [[number]]
### Data Frames

Data frames are similar to matrices but can contain various data types

x <- c(1,5,10)y <- c("A", "B", "C") z <-data.frame(x,y) ХУ 1 1 A 2 5 B 3 10 C

# Input/Output

✓ To read data into R from a text file use read.table()
• read help(read.table) to learn more

```
Data_test <- read.table(header=TRUE,
text='subject sex size
1 M 7
2 F NA
3 F 9
4 M 11 ')
```

✓ To write data into R from a text file use read.table()

write.table(Data\_test, "data\_test.csv", row.names=FALSE)

### Plots generation in R

 $\checkmark R$  provides very rich set of plotting possibilities

✓ The basic command is plot()

✓ Each library has its own version of plot() function

✓ When R plots graphics it opens "graphical device" that could be either a window or a file

# **Plotting functions**

#### ✓ R offers following array of plotting functions

Function	Description
plot(x)	plot of the values of x variable on the y axis
	bi-variable plot of x and y values (both axis scaled based
plot(x,y)	on values of x and y variables)
pie(y)	circular pie-char
boxplot(x)	Plots a box plot showing variables via their quantiles
hist(x)	Plots a histogram(bar plot)



### plot : Plotting functions

 $\checkmark$  Lets work on plot, hist and pie chart x <- c(1,2,3,4) y <- c(5,6,7,8) plot(x,y) plot(x,y,col="red") pie(x) pie(y) hist(y)

### **Boxplot : Plotting functions**

#### ✓ Lets work on boxplot

```
x <- c(1,2,3,4)
y <- c(5,6,7,8)
boxplot(x)
boxplot(y)
boxplot(x)
boxplot(x,y)
boxplot(x,y,col="grey")
boxplot(x,y,col="red")
boxplot(x,y,col=c("red",blue))
```

# References

1.<u>https://media.readthedocs.org/pdf/a-little-book-of-r-for-</u> bioinformatics/latest/a-little-book-of-r-for-bioinformatics.pdf

2.https://cran.r-project.org/doc/manuals/r-release/R-intro.pdf