Introduction to Biological Big Data

Simon Andrews simon.andrews@babraham.ac.uk

v2025-05

























Course Structure

Central Dogma Data Sources

- Genomes and Annotations
- Protein Domains and Structures
- Reactions, Pathways and Interactions
- Experimental Techniques, Datatypes and Repositories
 - Sequencing and Variants
 - Proteomics and Metabolites
 - Flow and Imaging

Central Dogma Data Resources





Genomes and Annotations

- Genome Assemblies
 - Underlying sequence of the organism's chromosomes
 - Often starts as scaffolds / contigs
 - Eventually assembled into chromosomes (still with holes)
 - Only one chromosome sequence per chromosome
 - Represents an 'average' individual (unless backcrossed)
 - Variations (natural or clinical) are stored separately)
 - Assembly is refined and improved over time, new releases get new names

Genome Assembly Nomenclature

- Chromosome / Scaffold sequences
 - Originally deposited with ENA / NCBI as sequence records

- Genome Assembly
 - Given an official name by a supervising group (sometimes two!)
 - Fixed coordinates at that point

Current Human Genome

- Assembly Name: GRCh38
- Current Patch: GRCh38.p16
- Managed by: Genome Reference Consortium
 - Assembly type:
 - Chromosome:
 - Genome:

Chromosomal Chr1 = CM000663.2 = NC_000001.11 GCF_000001405.40 (Assembly Refseq) GCA_000001405.29 (Assembly Genbank)

Genome Annotation Sets

- Built on top of a specific assembly
- Combination of prediction tools and real data
- Main annotation is Genes, Transcripts, Coding Sequences
- Many other tracks often added

- Different sites will have different annotations
- Annotations updated more frequently than assemblies

Genome Annotation Details

Genome-Annotation-Data

##Genome-Annotation-Data-START##							
Annotation Provider::NCBI							
Annotation Status::Updated annotation							
Annotation Name::Homo sapiens Updated Annotation Release 109.20210226							
Annotation Version::109.20210226							
Annotation Pipeline::NCBI eukaryotic genome annotation pipeline							
Annotation Software Version::8.6							
Annotation Method::Best-placed RefSeq; propagated RefSeq model							
Features Annotated::Gene; mRNA; CDS; ncRNA							
##Genome-Annotation-Data-END##							

General stats

- pseudogenes

Total No of Genes	60649	Total No of Transcripts	237012
Protein-coding genes	19955	Protein-coding transcripts	86757
Long non-coding RNA genes	17944	- full length protein-coding	61015
Small non-coding RNA genes	7567	- partial length protein-coding	25742
Pseudogenes	14773	Nonsense mediated decay transcripts	18881
- processed pseudogenes	10667	Long non-coding RNA loci transcripts	48752
- unprocessed pseudogenes	3565		
- unitary pseudogenes	241		
- polymorphic pseudogenes	49		(00/0
- pseudogenes	15		03700
Immunoglobulin/T-cell receptor gene segments		Genes that have more than one distinct translations	13007
- protein coding segments	409		

236

Assembly	GRCh38.p14 (Genome Reference Consortium Human Build 38), INSDC Assembly <u>GCA_000001405.29</u> &, Dec 2013
Base Pairs	3,099,750,718
Golden Path Length	3,099,750,718
Assembly provider	Genome Reference Consortium &
Annotation provider	Ensembl
Annotation method	Full genebuild
Genebuild started	Jan 2014
Genebuild released	Jul 2014
Genebuild last updated/ patched	Nov 2024
Database version	114.38
Gencode version	GENCODE 48

Gene counts (Primary assembly)

Coding genes	19,871 (excl 661 readthrough)
Non coding genes	42,126
Small non coding genes	4,866
Long non coding genes	35,044 (excl 301 readthrough)
Misc non coding genes	2,216
Pseudogenes	15,198 (excl 1 readthrough)
Gene transcripts	387,954

Viewing Annotated Genomes

- Mostly web based
 - Species specific sites
 - Generic multi-species sites

- Often adds more information
 - Regulation, conservation, repeats
 - Experimental datasets
 - Upload your own

Species specific genome viewer sites



Arabidopsis

https://www.arabidopsis.org





https://flybase.org/

WormBase

Nematode worms

https://wormbase.org

Generic genome viewer sites



https://www.ensembl.org



UCSC Browser https://genome.ucsc.edu



Name 🍦	Transcript ID 💧	bp 🍦	Protein 🖕	Biotype	CCDS	UniProt Match 🖕
BRCA2-201	ENST0000380152.8	11954	<u>3418aa</u>	Protein coding	<u>CCDS9344</u> &	<u>P51587</u> 🗗
BRCA2-210	ENST0000680887.1	11880	<u>3418aa</u>	Protein coding	<u>CCDS9344</u> 🖗	-
BRCA2-206	ENST0000544455.6	11854	<u>3418aa</u>	Protein coding	<u>CCDS9344</u> @	<u>P51587</u> ൽ
BRCA2-204	ENST0000530893.6	2011	<u>481aa</u>	Protein coding	-	<u>A0A590UJI7</u> &
BRCA2-207	ENST0000614259.2	11763	<u>2649aa</u>	Nonsense mediated decay	-	-



Chromosome 13: 32,315,086-32,400,268

🌣 < 🖃 🇞							
Assembly exceptions Chr. 13	p13 p11.2	· •	018.3 014.11 014.3 02	Q24 US(S)	q31.1 q31.3	- <u>-</u>	q34
Assembly exceptions							

Region in detail @





Track Based Displays





- Large scale querying and export of genomic data
- Annotations, Sequences, Variants etc.
 - Select data type (eg genes)
 - Select genome species
 - Select genes / regions / identifiers
 - Select attributes to export
 - Generate report

Genome File Formats

- Genome Assemblies
 - Chr sequence, FastA format
 - A small header plus DNA bases
 - Also used for RNA / protein

★ ≑	Species	DNA (FASTA)	cDNA (FASTA)	CDS (FASTA)	ncRNA (FASTA)	Protein sequence (FASTA)	Annotated sequence (EMBL)	Annotated sequence (GenBank)	Gene sets
Y	<u>Human</u> Homo sapiens	<u>FASTA</u> ₽	<u>FASTA</u> ₽	<u>FASTA</u> ₽	<u>FASTA</u> ₽	<u>FASTA</u> ₽	<u>EMBL</u> ₽	<u>GenBank</u> &	<u>GTF</u> ଜ <u>GFF3</u> ଜ
Y	<u>Mouse</u> Mus musculus	<u>FASTA</u> ₽	<u>FASTA</u> ₽	<u>FASTA</u> ₽	<u>FASTA</u> ₽	<u>FASTA</u> ₽	<u>EMBL</u> &	<u>GenBank</u> &	<u>GTF</u> & <u>GFF3</u> &
Y	<u>Zebrafish</u> Danio rerio	<u>FASTA</u> ₽	<u>FASTA</u> ₽	<u>FASTA</u> ₽	<u>FASTA</u> ₽	<u>FASTA</u> ₽	<u>EMBL</u> ₽	<u>GenBank</u> ⊮	<u>GTF</u> മ <u>GFF3</u> മ

- Gene Annotations
 - GFF or GTF format (both very similar)
 - Hierarchical format linking exons to transcripts to genes

https://www.ensembl.org/info/data/ftp/index.html

FastA Format Data

>I dna:chromosome chromosome:R64-1-1:I:1:230218:1 REF CATCCTAACACTACCCTAACACAGCCCTAATCTAACCCTGGCCAACCTGTCTCTCAACTT ACCCTCCATTACCCTGCCTCCACTCGTTACCCTGTCCCATTCAACCATACCACTCCGAAC CACCATCCATCCCTCTACTTACTACCACTCACCCACCGTTACCCCTCCAATTACCCATATC >II dna:chromosome chromosome:R64-1-1:II:1:813184:1 REF AAATAGCCCTCATGTACGTCTCCTCCAAGCCCTGTTGTCTCTTACCCGGATGTTCAACCA AAAGCTACTTACTACCTTTATTTTATGTTTACTTTTTATAGGTTGTCTTTTTATCCCACT TCTTCGCACTTGTCTCTCGCTACTGCCGTGCAACAACACTAAATCAAAACAAT CTACTACATCAAAACGCATTTTCCCCTAGAAAAAATTTTCCTTACAATATACTATACTAC

Annotation Descriptions



Gene Exon (combined into transcript) Coding Exon

GTF

- Targeted at gene structure definition
- Tab delimited text file
- Multi-level structure
 - Genes > Transcripts > Exons

GTF File Fields

- 1. Chromosome
- 2. Source
- 3. Feature Type
- 4. Start
- 5. End
- 6. Score
- 7. Strand (+/-)
- 8. Frame (1,2,3)
- 9. Group/Attributes

1	havana	gene	11869	14409	. + .	gene_id	"ENSG223972";	gene_name "DD>	(11L1";			
1	havana	transcript	11869	14409	. + .	gene_id	"ensg223972";	transcript_id	"enst456328";	transcript_nam	e "DDX11L1-	-202";
1	havana	exon	11869	12227	. + .	gene_id	"ensg223972";	transcript_id	"enst456328";	exon_number "1	"; exon_id	"ense2234944";
1	havana	exon	12613	12721	. + .	gene_id	"ensg223972";	transcript_id	"enst456328";	exon_number "2	"; exon_id	"ense3582793";
1	havana	exon	13221	14409	. + .	gene_id	"ensg223972";	transcript_id	"enst456328";	exon_number "3	; exon_id	"ense2312635";

Genome Exploration Exercise



mRNA Translation into Protein

Start Codon

GACACC ATG AGC ACT GAA ... CTG TGAUTRMet Ser Thr GluArg Stp

- Most species use the same code
- Some have minor differences

https://www.ncbi.nlm.nih.gov/Taxonomy/Utils/wprintgc.cgi



Codon Usage



Genes analyzed: Nuclear genes Ribosomal proteins Mitochondrial genes

Species: Homo sapiens Taxonomy ID: 9606 Assembly: GCF_000001405.39 | GRCh38.p13 Genetic code: 1 Number of genes: 19850 Number of codons: 11577026

Amino Acid	Codon	Count	RSCU	Preferred
Ala	GCA	187108	0.921	Unpreferred
Ala	GCC	323249	1.590	Preferred
Ala	GCG	89097	0.438	Preferred
Ala	GCT	213559	1.051	Unpreferred
Arg	AGA	142934	1.303	Unpreferred
Arg	AGG	140481	1.281	Unpreferred
Arg	CGA	70319	0.641	Unpreferred
Arg	CGC	119972	1.094	Preferred
Arg	CGG	132275	1.206	Preferred
Arg	CGT	52129	0.475	Unpreferred
Asn	AAC	212987	1.037	Preferred
Asn	AAT	197831	0.963	Unpreferred
Asp	GAC	287974	1.059	Preferred
Asp	GAT	255933	0.941	Unpreferred



Often Binding or catalytic sites



- A single protein can have more than one functional unit
 - Proteins are annotated with functional 'domains'
 - A domain is normally linked with a globular folded structure
- Domain structures are re-used to provide modular functionality across multiple proteins.
 - Often linked to exon structures or splice variation
- It can be useful to know the key functional amino acids
 - Binding pockets
 - Active sites

Protein Domain Databases



http://smart.embl-heidelberg.de/



https://www.ebi.ac.uk/interpro/



Tryp_SPc domain

This is a SMART Tryp_SPc domain (full annotation).

 Position:
 437 to 675

 E-value:
 4.3565597296112e-75 (HMMER2)



SMART ACC:	SM000020								
Definition:	Trypsin-like serine protease								
Description:	Many of these are synthesised as inactive precursor zymogens that are cleaved during limited proteolysis to generate their active forms. A few, however, are active as single chain molecules, and others are inactive due to substitutions of the catalytic triad residues.								

Types of domain

- Globular
 - Forms a concerted 3D structure
 - Most catalytic and some binding domains
- Semi-ordered
 - Coiled coil
 - Many binding domains
- Transmembrane
 - Threaded through a membrane
 - Transmembrane regions, then internal and external segments
- Disordered / Low Complexity
 - Linker regions
 - Intrinsically disordered proteins







Key Residue Databases



Protein Structure Databases



Sequence of	2LUL Soluti	• Chain	Ф 1: Тул	osine-p 🍳	A 0	0	C Structure
1 MGHHHHHHSH 1 81	1 INFNTILEEI LIP 91	RSQQKKK TSP	LNYKERL FVL	TKSMLTY YEGI	RAEKKYR KGFI	DVSKIK CVEIVK	2LUL Solution
NDDG VIPCONK 161 CEKYNLFE SSI	(ypf qvvhdantl ir	Y IFAPSPQSR	D LWVKKLKEE:	I KNNNNIMIKY	(HPKFWTDGSY (QCCRQTEKLA PG	, Model Inde
H / \	Madel 4 (00					C.	Тур



Ð	🗘 Structure								
^	2LUL Solution NMR Structure of P								
¥	Model Index	•		-	1				
	Туре	Model							
	Nothin	ng Focused			\odot				
	X Measureme	nts							
	Q Structural Motif Search								
	© Component	5			21.01				
	기 Preset	+ Add			Ð				
	Polymer	Cartoon	0	D					
	lon	Ball & Stick	0	D					
	Assembly Sy	mmetry							

Export Animation



C 3D View

NMR structure of the SH3 domain from the Tec protein tyrosine kinase

Mulhern, T.D., Pursglove, S.E., Booker, G.W.

(2002) J Biol Chem 277: 755-762

Released	2001-11-28
Method	SOLUTION NMR
Organisms	Mus musculus
Macromolecule	TYROSINE-PROTEIN KINASE TEC (protein)

2LUL

Solution NMR Structure of PH Domain of Tyrosine-protein kinase Tec from Homo sapiens, Northeast Structural Genomics Consortium (NESG) Target HR3504C

Liu, G., Xiao, R., Janjua, H., Hamilton, K., Shastry, R., Kohan, E., Acton, T.B., Everett, J.K., Lee, H., Pederson, K., Huang, Y.J., Montelione, G.T., Northeast Structural Genomics Consortium (NESG)

To be published

1GL5

Released	2012-08-15
Method	SOLUTION NMR
Organisms	Homo sapiens
Macromolecule	Tyrosine-protein kinase Tec (protein)
Unique Ligands	ZN



Download File View File

Protein Structure Classification Databases



https://scop.mrc-Imb.cam.ac.uk/

https://www.cathdb.info/

Predicted Structure Database

Currently (Mar 2022), only 7914/22818 protein coding genes have an experimental 3D structure available



ARTIFICIAL INTELLIGENCE

DeepMind's protein-folding Al has solved a 50-year-old grand challenge of biology

AlphaFold can predict the shape of proteins to within the width of an atom. The breakthrough will help scientists design drugs and understand disease.

https://alphafold.ebi.ac.uk/

Q5VSL9 E. coli Help: AlphaFold DB search help

Examples: Free fatty acid receptor 2 At1g58602

Feedback on structure: Contact DeepMind

By Will Douglas Heaven

Protein Annotation Exercise





Hierarchy of Reaction Annotations

- Components (Reactants / Products)
- Proteins (Enzymes)
- Reactions
- Pathways
- Processes



Rhea

Reactions



Enzymes ⁽²⁾ 82,966 proteins (UniProtKB) EC 2.7.10.1 Receptor protein-tyrosine kinase EC 2.7.10.2 Non-specific protein-tyrosine kinase EC 2.7.12.1 Dual-specificity kinase EC 2.7.12.2 Mitogen-activated protein kinase kinase



Enzyme Databases

- Enzymes are described by an Enzyme Commission (EC) number
 - EC 2.7.1.10 is phosphoglucokinase
 - Hierarchical structure
- Main Enzyme databases
 - Expasy^a Expasy Enzyme



EC Tree

2 Transferases
2.7 Transferring phosphorus-containing groups

- 2.7.1 Phosphotransferases with an alcohol group as acceptor
 - └፼2.7.1.10 phosphoglucokinase




Chemical entities of biological interest

A database of "small" molecules with biological relevance Natural or synthetic products which intervene in the processes of living organisms

CHEBI:58392 - α-D-glucose 1,6-bisphosphate(4-)



Pathways









Disease-associated

Reactome



Summation

One hallmark of cancer is altered cellular metabolism. Malic enzymes (MEs) are a family of homotetrameric enzymes that catalyse the reversible oxidative decarboxylation of L-malate to pyruvate, with a simultaneous reduction of NAD(P)+ to NAD(P)H. As MEs generate NADPH and NADH, they may play roles in energy production and reductive biosynthesis. Humans possess three ME isoforms; ME1 is cytosolic and utilises NADP+, ME3 is mitochondrial and can utilise NADP+ and ME2 is mitochondrial and can utilise either NAD+ or NADP+ (Chang & Tong 2003, Murugan & Hung 2012).

NADP-dependent malic enzyme (ME1, aka c-NADP-ME) is a cytosolic enzyme that oxidatively decarboxylates (s)-malate (MAL) to pyruvate (PYR) and CO2 using NADP+ as cofactor (Zelewski & Swierczynski 1991). ME1 exists as a dimer of dimers (Murugan & Hung 2012, Hsieh et al. 2014) and a divalent metal such as Mg2+ is essential for catalysis (Chang & Tong 2003).

Background literature references...

KEGG databases

Category	Entry point
	KEGG PATHWAY
Systems information	KEGG BRITE
	KEGG MODULE
	KEGG RModule
	KEGG ORTHOLOGY
	KEGG Annotation
Genomic	KEGG GENES
information	KEGG SeqData
Chamical	KEGG GENOME
	KEGG Virus
	KEGG COMPOUND
	KEGG GLYCAN
information	KEGG REACTION
	KEGG Enzyme
Health information	KEGG NETWORK
	KEGG DISEASE
	KEGG DRUG

KEGG



Functional Gene Sets



Molecular-level activities performed by gene products. Molecular function terms describe activities that occur at the molecular level, such as "catalysis" or "transport". GO molecular function terms represent activities rather than the entities (molecules or complexes) that perform the actions, and do not specify where, when, or in what context the action takes place. Molecular functions generally correspond to activities that can be performed by individual gene products (*i.e.* a protein or RNA), but some activities are performed by molecular complexes composed of multiple gene products. Examples of broad functional terms are *catalytic activity* and *transporter activity*, examples of narrower functional terms are *adenylate cyclase activity* or *Toll-like receptor binding*. To avoid confusion between gene product names and their molecular functions, GO molecular functions are often appended with the word "activity" (a *protein kinase* would have the GO molecular function *protein kinase activity*).

Cellular Component

Biological

Process

Molecular

Function

The locations relative to cellular structures in which a gene product performs a function, either cellular compartments (e.g., *mitochondrion*), or stable macromolecular complexes of which they are parts (e.g., the *ribosome*). Unlike the other aspects of GO, cellular component classes refer not to processes but rather a cellular anatomy.

The larger processes, or 'biological programs' accomplished by multiple molecular activities. Examples of broad biological process terms are *DNA repair* or *signal transduction*. Examples of more specific terms are *pyrimidine nucleobase biosynthetic process* or *glucose transmembrane transport* . Note that a biological process is not equivalent to a pathway. At present, the GO does not try to represent the dynamics or dependencies that would be required to fully describe a pathway.



Synonyms hexose anabolism, hexose biosynthesis, hexose formation, hexose synthesis

Alternate IDs None

Definition The chemical reactions and pathways resulting in the formation of hexose, any monosaccharide with a chain of six carbon atoms in the molecule. Source: ISBN:0198506732

Gene/product	Gene/product name	Organism
Sds	serine dehydratase	Mus musculus
G6pc	glucose-6-phosphatase, catalytic	Mus musculus
Gnpda1	glucosamine-6-phosphate deaminase 1	Mus musculus
Nr3c1	nuclear receptor subfamily 3, group C, member 1	Mus musculus
Gpt	glutamic pyruvic transaminase, soluble	Mus musculus
Ranbp2	RAN binding protein 2	Mus musculus
Ptpn2	protein tyrosine phosphatase, non-receptor type 2	Mus musculus
Stk11	serine/threonine kinase 11	Mus musculus
Gm10768	predicted gene 10768	Mus musculus
Fbp1	fructose bisphosphatase 1	Mus musculus



Genes assigned to ontology terms

Nanog homeobox [Source:HGNC Symbol;Acc:HGNC:20857]

- Cellular Component
 - GO:0005634 nucleus
 - GO:0005654 nucleoplasm
 - GO:0005730 nucleolus

- Molecular Function
 - GO:0003677 DNA binding
 - GO:0003700 transcription factor activity, sequence-specific DNA binding
 - GO:0003714 transcription corepressor activity
 - GO:0005515 protein binding
 - GO:0043565 sequence-specific DNA binding

- Biological Process
 - GO:0001714 endodermal cell fate specification
 - GO:0006351 transcription, DNA-templated
 - GO:0006355 regulation of transcription, DNAtemplated
 - GO:0007275 multicellular organism development
 - GO:0008283 cell proliferation
 - GO:0019827 stem cell population maintenance
 - GO:0030154 cell differentiation
 - GO:0035019 somatic stem cell population maintenance
 - GO:0045595 regulation of cell differentiation
 - GO:0045944 positive regulation of transcription from RNA polymerase II promoter
 - GO:1903507 negative regulation of nucleic acid-templated transcription

Reactions and Pathways Exercise



Regulation and Interactions

- The regulation of genes is as important as their structure or function
- Several sources of useful information
 - Regulatory binding proteins, mostly transcription factors
 - Interactions with other proteins to form complexes
 - Composition of known complexes

Transcription Factor Information

3 JASPAR²⁰²⁴

Profile sur	nmary 🏾 🏋 Add	Seq	uence l	ogo								📩 Do	wnload S	SVG
Name:	ATF3			2.0										
Matrix ID:	MA0605.2			1.5		T/			١T					
Class:	Basic leucine zipper factors (bZIP)			10			1	'	1	1				
Family:	Fos-related			a 10	٨	Ш	•		•		Т			
Collection:	CORE			0.5	P		Л	U	Л	UΓ				
Taxon:	Vertebrates			a.a 🚍	1 2	3		6 3		9 10		12		
Species:	Homo sapiens													
Data Type:	HT-SELEX				_		_			_		_		
Validation:	12815047	m	equend atrix	Ly		JASPAR	≵ ⊺	RANSFAC	🛓 м	ЕМЕ	🛓 RAW PF	M		
						- Reverse	comp.							
Uniprot ID:	P18847													
Source:	P18847 28473536	Α[8505	24741	0	0	40546	0	891	0	1520	40546	0	691
Source: Comment:	P18847 28473536	A[C[8505 8220	24741 1354	0	0	40546 0	0 40546	891 0	0	1520 40546	40546 0	0 15737	691 168
Source: Comment:	P18847 28473536	A [C [G [8505 8220 16894	24741 1354 15805	0 0 1	0 0 40546	40546 0 0	0 40546 0	891 0 40546	0 0 0	1520 40546 0	40546 0 0	0 15737 1094	691 168 824
Oniprot ID: Source: Comment:	P18847 28473536	A[C[G[T[8505 8220 16894 6926	24741 1354 15805 0	0 0 1 40546	0 0 40546 1366	40546 0 0 0	0 40546 0 556	891 0 40546 0	0 0 0 40546	1520 40546 0 0	40546 0 0 0	0 15737 1094 24808	691 168 824 856

Transcription Factor Information



Genes Regulated by a Transcription Factor

- Difficult to predict lots of false positives
 - Swiss Regulon



BRCA2

Entrez	Description	Chrom.	Strand	Promoter (Start - Stop)	TSS
675	breast cancer 2, early onset	chr13	+	32884616 - 32890116	32889616

Transcription Factor Binding Sites

Download all TFBS in the BRCA2 promoter

now 10 🗸 entries						Search:	
Motif	Source 👙	Strand 🌢	Start 🍦	Stop 🍦	PValue 🔺	Match Sequence	Overlap w/ Footprints
Pax4_MA0068.1	JASPAR	+	32888990	32889019	0.0E+00	AAAAAAAAAGCAAAAGATACTACCAAGCC	30
V_GC_01_M00255	TRANSFAC	-	32889167	32889180	0.0E+00	AGTGGGCGGGGCTG	14
V_LDSPOLYA_B_M00317	TRANSFAC	-	32889437	32889452	0.0E+00	AGTGTGTGTTCTCTTC	16
V_SOX2_Q6_M01272	TRANSFAC	-	32889284	32889299	0.0E+00	AATACCTTTGTTCTGA	16
V_SP1_Q4_01_M00932	TRANSFAC	-	32889989	32890001	0.0E+00	AAGGGGCGGGGCT	13
V_STAT5A_01_M00457	TRANSFAC	+	32890101	32890115	0.0E+00	AATTTCTTGGAAACA	15
V_STAT5A_Q6_M01890	TRANSFAC	+	32890100	32890112	0.0E+00	AAATTTCTTGGAA	13
V_STAT5B_01_M00459	TRANSFAC	+	32890101	32890115	0.0E+00	AATTTCTTGGAAACA	15
V_STAT_Q6_M00777	TRANSFAC	+	32890099	32890111	0.0E+00	GAAATTTCTTGGA	13
SP1_C2H2_DBD_monomeric_11_1	SELEX	+	32889169	32889179	1.0E-05	GCCCCGCCCAC	11
nowing 1 to 10 of 196 entries						Pr	evious Next

ReMap2022

GTRD

Gene Transcription Regulation Database

Gene Interactions

- Many genes form stable or transitory interactions with others
- Knowing the genes that interact helps understand biology





Interactor Statistics	
Proteins/Genes 208	Public 1
Interactors w/ Physical (HTP) Evi	
Interactors w/ Physical (HTP) Evid	Jence (101)

Interactor	 Organism / Chemical Type 	♦ Aliases	Description	e Evidence
RAD51	H. sapiens	RECA, BRCC5, MRMV2, HRAD51, RAD51A, HsR ad51, HsT16930	RAD51 recombinase	1 74 View
PALB2	H. sapiens	PNCA3, FANCN	partner and localizer of BRCA2	34 View
BRCA1	H. sapiens	IRIS, PSCP, FANCS, RNF53, BRCC1, PNCA4, BR CAI, PPP1R53, BROVCA1	breast cancer 1, early onset	2 11 View
FANCD2	H. sapiens	FA4, FAD, FAD2, FACD, FANCD, FA-D2	Fanconi anemia, co	20K2
HMG20B	H. sapiens	SOXL, HMGX2, HMGXB2, PP7706, BRAF25, BRAF 35, pp8857, SMARCE1r	high mobility group	PCNA RPA1
PLK1	H. sapiens	PLK, STPK13	polo-like kinase 1 (RAZ) BARDI (R	AD51

Interactors w/ Genetic (LTP) Evidence (3)
 Interactors w/ More than One Evidence Type (12)

Interactors w/ Physical (LTP) Evidence (77)
Interactors w/ Genetic (HTP) Evidence (15)

Types of Interaction

• Physical

- Two proteins directly interact, either stably or transiently

- Genetic
 - One gene influences another, normally after modification
 - Co-expression
 - Knockout compensation

Complex Prediction

- Many proteins interact with several others, but at different times
- Complexes suggest that multiple proteins directly associate
 - Can't always be clearly predicted from pairwise interactions
 - Other experimental methods are required

Complex Portal



Regulatory Information Exercise



Sequence Variants



Sequence Variants

ReferenceGATCTTAGVariantGATCTTAC

- Germline variants
 - Happen in sperm or eggs
 - Completely inherited into the next generation
 - Can cause genetic disease
- Somatic variants
 - Happen in other tissues
 - Partially penetrant
 - Common cause of cancer

Types of Variant

Ref GATCTTA<mark>G</mark>CTGA Var GATCTTA<mark>C</mark>CTGA

Substitution Single Nucleotide Polymorphism SNP

Ref GATCTTAG.CTGA Var GATCTTACAACTGA Insertion InDel Ref GATCTTAGCTGA Deletion Var GATCTTAC.GA

Functional Variant Consequences

- Within Coding Region
 - Silent (codon changes, but same translation)
 - Missense (change translation from one amino acid to another)
 - Nonsense (change translation from one amino acid to STOP)
 - Frameshift (InDel changing the translation frame)
- Outside CDS
 - Breaks or adds splice junction
 - Changes functional binding site

Structural variants

- Chromosomal copy number change
 - Gain or loss of a chromosome
 - Leads to serious genetic disease

- Segmental Deletion / Duplication
 - Large parts of chromosomes deleted, duplicated, inverted, translocated
 - 1kb to 3Mbp
 - Affects many genes, can lead to gene fusions

Databases of Variants

- Common genomic variants
 - Measured across a large population
 - Shows natural variation
 - Not necessarily linked to disease
 - Used for studying populations and families
- Functional variants
 - Variants with an associated phenotype
 - Often disease related but can be any measurable phenotype

Variant Databases

- Single Variants
 - dbSNP (<u>https://www.ncbi.nlm.nih.gov/snp/</u>)
 - Full reference for any reported SNPs, mix of functional and non-functional
 - HGMD (<u>http://www.hgmd.cf.ac.uk</u>)
 - Human genetic disease focussed database
 - COSMIC (<u>https://cancer.sanger.ac.uk/cosmic</u>)
 - Mutations observed in Cancer
 - Also has details of mutations in immortalised cell lines
- Larger Regions
 - dbVar (<u>https://www.ncbi.nlm.nih.gov/dbvar/</u>)
 - Counterpart to dbSNP for larger variants
 - ClinVar (<u>https://www.ncbi.nlm.nih.gov/clinvar/</u>)
 - Larger variants with clinical relevance
 - OMIM (<u>https://www.ncbi.nlm.nih.gov/omim</u>)
 - A more wide ranging collection of the phenotypic variation linked to genes

Variant Terminology

• Minor Allele Frequency (MAF)

– How prevalent the variant is in the population

- Impact scores (SIFT / PolyPhen etc)
 - A quantitative value assessing the likely biological impact of a variant

Variant Exercise



Other Information and Data Sources







Gene Expression Information Genevisible



Gene Expression Information

Expression Atlas

Organism part

Showing 29 experiments:





Post translational Modifications

- Many proteins are modified after they have been translated
 - Phosphorylation
 - Glycosylation
 - Ubiquitination
 - Nitrosylation
 - Methylation
 - Acetylation
 - Lipidation
 - Proteolysis

iPTMnet

Both PTMs observed on a protein and proteins modified by a query gene.

¥.	Site All -	PTM Type	PTM Enzyme	Score 2 selected -
	S21	Phosphorylation		****
	S115	Phosphorylation		****
	S180	Phosphorylation	P05771 (PRKCB)	****
	Y223	Phosphorylation	Q06187 (BTK) , Q08881 (ITK) , P07948 (LYN) , P00519 (ABL1) , A0A173G4P4 (Abl fusion) , P42680 (TEC)	****
	Y551	Phosphorylation	P07948 (LYN) , Q06187 (BTK) , P12931 (SRC) , P43405 (SYK)	****



Combined Gene/Protein Centric Datasources



https://www.uniprot.org/

https://www.genecards.org



https://www.ncbi.nlm.nih.gov/gene/



https://www.wikigenes.org/



Names & Taxonomy

Subcell. location

Pathol./Biotech

PTM / Processing

Expression

Interaction

Structure

Family & Domains

Sequences (1+)

Similar proteins

Cross-references

Entry information

Miscellaneou:



Disease Relevance High Impact publication summaries Biological Context Anatomical Context Chemical Compound Associations Physical Interactions Enzymatic Interactions Regulatory relationships Analytical, diagnostic and therapeutic context References



Summary

Genomic context
Genomic regions, transcripts, and products
Expression
Bibliography
Variation
Pathways from PubChem
Interactions
General gene information Markers, Homology, Gene Ontology
General protein information
NCBI Reference Sequences (RefSeq)
Related sequences
Additional links



Research Antibodies Assays Proteins Inhib. RNA CRISPR Exp. Assays miRNA Drugs Animal M	section	Aliases Paralogs	Disorders Pathways	Domains Products	Drugs Proteins	Expression Publications	Function Sources	Genomics Summaries	Localization Transcripts	Orthologs Variants	
Deschuster and the second se	Research	Antibodies	Assays	Proteins	Inhib. RNA	CRISPR	Exp. Assays	miRNA	Drugs	Animal Mode	ls
Products Cell Lines Clones Primers Genotyping	Products	Cell Lines	Clones	Primers	Genotyping						

Experimental Data Types and Repositories

Simon Andrews, Chris Hall, Judith Webster, Eoin Fahy, Laura Biggins, Hanneke Okkenhaug, Simon Walker



Big Data Generation

- High throughput sequencing
 - Genomics, Transcriptomics, Epigenetics
- Multi-channel Flow Cytometry
 - Cell surface proteomics
- Mass Spectrometry
 - Proteomics, Metabolomics
- Biological Imaging
 - Cell / Tissue structure, Proteomics, Metabolomics

Data Repositories

- For many techniques deposition of data in a suitable repository is a condition of publication
- Repositories are more developed and complete for some techniques than others
- Still a growing area

Mandatory deposition	Suitable repositories
Protein sequences	Uniprot
DNA and RNA sequences	Genbank
	DNA DataBank of Japan (DDBJ)
	EMBL Nucleotide Sequence Database (ENA)
DNA and RNA sequencing data	NCBI Trace Archive
	NCBI Sequence Read Archive (SRA)
Genetic polymorphisms	dbSNP
	dbVar
	European Variation Archive (EVA)
Linked genotype and phenotype data	dbGAP
	The European Genome-phenome Archive (EGA)
Macromolecular structure	Worldwide Protein Data Bank (wwPDB)
	Biological Magnetic Resonance Data Bank (BMRB)
	Electron Microscopy Data Bank (EMDB)
Gene expression data (must be MIAME compliant)	Gene Expression Omnibus (GEO)
	ArrayExpress
Crystallographic data for small molecules	Cambridge Structural Database
Proteomics data	PRIDE
*Earth, space & environmental sciences	Recommended Repositories

FAIR Data Principles

- Designed to make data as useful as possible to future researchers
 - -Findable
 - Unique accession code
 - Rich metadata
 - -Accessible
 - Automated query and download API
 - Iteroperable
 - Use of open formats
 - Standard Ontologies for descriptions
 - -Reusable
 - Clear licensing
 - Annotated to common community standards
High Throughput Sequencing



Element Aviti

Data Generation Capacity

Sequencer	Read Length	Bases per run
Illumina NovaSeq	50-250bp	3000 Gbp
ONT Promethion 48	1kb - 80Mbp	48 x 20-90 Gbp
PacBio Revio	1kb - 20kb	90 Gbp

What can you measure?

- Genomics
 - Whole genome sequencing, Targeted Sequencing
- Transcriptomics
 - RNA-Sequencing
- Regulation
 - Accessible DNA (ATAC-Seq), Histone Modifications, Transcription Factor binding sites
- Epigenetics
 - DNA Methylation, Chromatin Structure

Types of Sequencing Library

- DNA-Based
 - Genome-Seq: Variants
 - Exome-Seq: Variants
 - ATAC-Seq: Accessible DNA
 - ChIP, Cut n Run:
 - DNA binding sites
 - Epigenetic Marks
 - Polymerase Attachment
 - BS-Seq, EM-Seq: DNA Methylation
 - 3C, 4C, Hi-C: Genome Structure
 - TrAEL-Seq: DNA-replication

- RNA Based
 - RNA-Seq: RNA transcription
 - Ribo-Seq: Ribosome attachment
 - CAGE-Seq: Transcription start sites
 - VDJ-Seq: Antibody repertoires
 - CLIP-Seq: RNA-binding protein sites
 - sRNA-Seq: Small RNA abundance
 - SLAM-Seq: RNA dynamics

Genome Sequencing



RNA-Sequencing



Enrichment Sequencing



Bisulphite/EM Sequencing



10X Single Cell RNA-Seq



Gel Beads in Emulsion (GEMs)

10X Single Cell RNA-Seq Adapter System



Multi-measure single cell

> Nat Commun. 2018 Feb 22;9(1):781. doi: 10.1038/s41467-018-03149-4.

scNMT-seq enables joint profiling of chromatin accessibility DNA methylation and transcription in single cells

Stephen J Clark ¹, Ricard Argelaguet ² ³, ChantrioInt-Andreas Kapourani ⁴, Thomas M Stubbs ⁵, Heather J Lee ⁵ ⁶ ⁷, Celia Alda-Catalinas ⁵, Felix Krueger ⁸, Guido Sanguinetti ⁴, Gavin Kelsey ⁵ ⁹, John C Marioni ¹⁰ ¹¹ ¹², Oliver Stegle ¹³, Wolf Reik ¹⁴ ¹⁵ ¹⁶ Chromium Single Cell Multiome ATAC + Gene Expression

Unify the Transcriptome and Epigenome in Every Cell



Spatial Transcriptomics



- 10X Visium
- Nanostring CosMX
- Vizgen Merscope

FastQ Format Data

@HWUSI-EAS611:34:6669YAAXX:1:1:5069:1159 1:N:0: TCGATAATACCGTTTTTTCCGTTTGATGTTGATACCATT +

DF=DBD<BBFGGGGGGGGGGBD@GGGGD4@CA3CGG>DDD:D,B
@HWUSI-EAS611:34:6669YAAXX:1:1:5266:1162 1:N:0:
GGAGGAAGTATCACTTCCTTGCCTGCCTCCTCGGGGCCT
+

:GBGGGGGGGGGGGGGGDEDGGGGGGGHHDHGHHGBGG:GG

Public Sequencing Databases

- GEO (NCBI)
- Array Express (EBI)
 - Databases for quantitated sequencing data.
 Provide experimental annotation and metadata and processed quantitated data







- SRA (NCBI)
- ENA (EBI)
 - Provide raw sequencing data as fastq files

Accession Codes

Transcription-induced formation of extrachromosomal DNA during yeast ageing

Ryan M. Hull¹^{1°a}, Michelle King¹, Grazia Pizza^{1°b}, Felix Krueger², Xabier Vergara^{1°c}, Jonathan Houseley¹*

Data Availability Statement: All relevant data are within the paper and its Supporting Information files. All sequencing files are available from the GEO database (accession number GSE135542).



Series GSE13554	Query DataSets for GSE135542
Status	Public on Oct 18, 2019
Title	Transcription-induced formation of extrachromosomal DNA during yeast ageing
Organism	Saccharomyces cerevisiae
Overall design	Aged cell samples analysed in pairs of -/+ Cu, for both wt and various mutants. 3 replicates of the 3xCUP1 experiment are included
Contributor(s)	Hull R, King M, Houseley J
Platforms (1)	GPL17342 Illumina HiSeq 2500 (Saccharomyces cerevisiae)
Samples (30)	GSM4015617 3xCUP1_24hr_1_REC-seq
± More	GSM4015618 3xCUP1_24hr_2_REC-seq
	GSM4015619 3xCUP1_24hr_300uM_Cu_1_REC-seq
Delettere	

Re	la	tion	15	
	_		-	

BioProject	PRJNA559191
SRA	SRP217740

Supplementary file	Size	Download	File type/resource
GSE135542_3xCUP1_processed_data_report.txt.gz	1.3 Mb	(ftp)(http)	TXT
GSE135542_cu_and_gal_processed_data_report.txt.gz	12.4 Mb	(ftp)(http)	ТХТ
GSE135542_mutants_processed_data_report.txt.gz	1.4 Mb	(ftp)(http)	TXT

SRA Run Selector 🛽

Raw data are available in SRA

Processed data are available on Series record

SRA Run Selector

elec	tor	BioSample	Bases ³	Bytes	Experiment	ہ GEO_Accession	Sample Name	source_name	° strain
1	SRR9924096	SAMN12529574	1.01 G	344.69 Mb	SRX6673092	GSM4015624	GSM4015624	Cells aged 24 hours in SD media	MEP mus81
2	SRR9924097	SAMN12529572	994.71 M	338.32 Mb	SRX6673093	GSM4015625	GSM4015625	Cells aged 24 hours in SD media	MEP mus81
3	SRR9924098	SAMN12529570	838.88 M	294.83 Mb	SRX6673094	GSM4015626	GSM4015626	Cells aged 24 hours in SD media	MEP mus81
4	SRR9924099	SAMN12529569	631.87 M	250.37 Mb	SRX6673095	GSM4015627	GSM4015627	Cells aged 48 hours in YPD media	MEP [Pgal1-3HA cup1]/CUP1
5	SRR9924100	SAMN12529567	1.11 G	407.87 Mb	SRX6673096	GSM4015628	GSM4015628	Cells aged 48 hours in YPD media	MEP [Pgal1-3HA cup1]/CUP1
6	SRR9924101	SAMN12529565	903.88 M	343.05 Mb	SRX6673097	GSM4015629	GSM4015629	Cells aged 48 hours in YPGal media	MEP [Pgal1-3HA cup1]/CUP1
7	SRR9924102	SAMN12529564	1.45 G	529.28 Mb	SRX6673098	GSM4015630	GSM4015630	Cells aged 48 hours in YPGal media	MEP [Pgal1-3HA cup1]/CUP1
8	SRR9924103	SAMN12529561	646.51 M	227.76 Mb	SRX6673099	GSM4015631	GSM4015631	Cells aged 24 hours in SD media	MEP sae2
9	SRR9924104	SAMN12529560	1.05 G	357.64 Mb	SRX6673100	GSM4015632	GSM4015632	Cells aged 24 hours in SD media	MEP sae2
10	SRR9924105	SAMN12529558	829.11 M	281.75 Mb	SRX6673101	GSM4015633	GSM4015633	Cells aged 24 hours in SD media	MEP sae2
11	SRR9924106	SAMN12529557	823.58 M	284.92 Mb	SRX6673102	GSM4015634	GSM4015634	Cells aged 24 hours in SD media	MEP sae2
12	SRR9924107	SAMN12529555	1.03 G	354.90 Mb	SRX6673103	GSM4015635	GSM4015635	Cells aged 24 hours in SD media	MEP spt3
13	SRR9924108	SAMN12529553	994.63 M	344.93 Mb	SRX6673104	GSM4015636	GSM4015636	Cells aged 24 hours in SD media	MEP spt3
14	SRR9924109	SAMN12529608	551.82 M	242.67 Mb	SRX6673105	GSM4015637	GSM4015637	Cells aged 24 hours in SD media	MEP
15	SRR9924110	SAMN12529606	961.46 M	360.10 Mb	SRX6673106	GSM4015638	GSM4015638	Cells aged 24 hours in SD media	MEP
16	SRR9924111	SAMN12529604	1.33 G	454.02 Mb	SRX6673107	GSM4015639	GSM4015639	Cells aged 24 hours in SD media	MEP
17	SRR9924112	SAMN12529602	1.06 G	395.94 Mb	SRX6673108	GSM4015640	GSM4015640	Cells aged 24 hours in SD media	MEP
18	SRR9924113	SAMN12529601	563.99 M	258.47 Mb	SRX6673109	GSM4015641	GSM4015641	Cells aged 24 hours in SD media	MEP
19	SRR9924114	SAMN12529599	886.36 M	336.01 Mb	SRX6673110	GSM4015642	GSM4015642	Cells aged 24 hours in SD media	MEP
20	SRR9924115	SAMN12529598	1.30 G	446.44 Mb	SRX6673111	GSM4015643	GSM4015643	Cells aged 24 hours in SD media	MEP
21	SRR9924116	SAMN12529596	1.33 G	483.90 Mb	SRX6673112	GSM4015644	GSM4015644	Cells aged 24 hours in SD media	MEP
22	SRR9924117	SAMN12529594	1.13 G	389.68 Mb	SRX6673113	GSM4015645	GSM4015645	Cells aged 24 hours in SD media	MEP
23	SRR9924119	SAMN12529593	921.09 M	324.51 Mb	SRX6673114	GSM4015646	GSM4015646	Cells aged 24 hours in SD media	MEP



Project: PRJNA559191

Extrachromosomal circular DNA (eccDNA) facilitates adaptive evolution by allowing rapid and extensive gene copy number variation, and is implicated in the pathology of cancer and ageing. Here, we demonstrate that yeast aged under environmental copper accumulate high levels of eccDNA containing the copper resistance gene CUP1. Transcription of CUP1 causes CUP1 eccDNA accumulation, which occurs in the absence of phenotypic selection. We have developed a sensitive and quantitative eccDNA sequencing pipeline that reveals CUP1 eccDNA accumulation on copper exposure to be exquisitely site specific, with no other detectable changes across the eccDNA complement. eccDNA forms de novo from the CUP1 locus through processing of DNA double-strand breaks (DSBs) by Sae2 / Mre11 and Mus81, and genome-wide analyses show that other protein coding eccDNA species in aged yeast share a similar biogenesis pathway. Although abundant we find that CUP1 eccDNA does not replicate efficiently, and high copy numbers in aged cells arise through frequent formation events combined with asymmetric DNA segregation. The transcriptional stimulation of CUP1 eccDNA formation shows that age-linked genetic change varies with transcription pattern, resulting in gene copy number profiles tailored by environment. Overall design: Aged cell samples analysed in pairs of -/+ Cu, for both wt and various mutants. 3 replicates of the 3xCUP1 experiment are included.

Show More

Organism:	Saccharomyces cerevisiae (baker's yeast)
Secondary Study Accession:	SRP217740
Study Title:	Transcription-induced formation of extrachromosomal DNA during yeast ageing
Center Name:	Bioinformatics, The Babraham Institute
Study Name:	Transcription-induced formation of extrachromosomal DNA during yeast ageing

Read Files								2
Show Column	Selection							~
Download report: JSON TSV Download Files as ZIP Download selected files							l files	
						Ŧ	Download All	±٥
Study Accession	Sample Accession	Experiment Accession	Run Accession	Tax Id	Scientific Name		FASTQ FTP	Sul
					Saccharomyces		R992409fastq.gz	
PRJNA559191	NA559191 SAMN12529574 SRX66		SRR9924096 4932		cerevisiae	SR	R992409fastq.gz	

SRA Explorer

This tool aims to make datasets within the Sequence Read Archive more accessible.

Search for:	SRP21774	SRP217740[All Fields]				
Max Results	100 0	Start At Record	0			
Need inspiration? Try g	SE30567 , SRP0435:	10, PRJEB8073, ERP009109 0	r human liver miRNA.			

Select relevant datasets and click add to collection. When you're finished, view all saved datasets with the button in the top right of the page, where you can copy the SRA URLs.

Showing 30 results.

•

Title	Accession	Instrument	Total Bases (Mb)	Date Created
GSM4015617: 3xCUP1_24hr_1_REC-seq; Saccharomyces cerevisiae; OTHER	SRR9924120	Illumina HiSeq 2500	8685	21 Oct 2019
GSM4015618: 3xCUP1_24hr_2_REC-seq; Saccharomyces cerevisiae; OTHER	SRR9924121	Illumina HiSeq 2500	10212	21 Oct 2019
GSM4015619: 3xCUP1_24hr_300uM_Cu_1_REC-seq; Saccharomyces cerevisiae; OTHER	SRR9924122	Illumina HiSeq 2500	9693	21 Oct 2019
GSM4015620: 3xCUP1_24hr_300uM_Cu_2_REC-seq; Saccharomyces cerevisiae; OTHER	SRR9924123	Illumina HiSeq 2500	9602	21 Oct 2019

SRA Downloader

sradownloader SRR9924120

Sequencing Data Exercise



Flow Cytometry



Flow Cytometry



Small Scale Measurement



https://flowjo.com

Using Flow for Sorting Cells



- Cell subpopulations
- CRISPR screens
- Cow sexing!

Large Scale Measurement





https://doi.org/10.3389/fimmu.2019.01194

Problems with multiple fluorescent markers



Traditionally filters measure one wavelength per fluor

Spectral Flow Cytometry measures the whole spectrum and can deconvolve overlapping emissions spectra

Allows for 40+ markers to be used simultaneously.



Public Flow Data Repository



- Deposition of FCS files
 - Instrument details
 - Raw data
 - Analysis details

• Basic description of experiment structure

ID: FR-FCM-Z2KP Prim

« Back to All Public Experiments



« Back to Start Page

Help

The following open access article describes how to upload and annotate flow cytometry data sets: Spidlen J, Breuer K and Brinkman R. Preparing a Minimum Information about a Flow Cytometry Experiment (MIFlowCyt) Compliant Manuscript Using the International Society for Advancement of Cytometry (ISAC) FCS File Repository (FlowRepository.org). <u>Current</u> <u>Protocols in Cytometry, UNIT 10.18,</u> July 2012.

We also have a <u>Quick start guide</u> and a <u>FAQ</u> section.

You may download <u>slides</u> from our Workshop at CYTO 2012: Publishing MIFlowCyt Compliant Data to ISAC's FlowRepository.org for Cytometry A and Other Journals

• Experiment O	VEIVIEW					
Repository ID:	FR-FCM-Z2KP	Experiment name:	Human COVID-19 Immune Phenotyping	MIFlowCyt score:	97.60%	
Primary researcher:	Stephanie Humblet-Baron	Stephanie Humblet-Baron PI/manager:		Uploaded by:	Oliver Burton	
Experiment dates:	2020-04-21 - 2020-04-24 Dataset uploaded:		May 2020	Last updated:	May 2020	
Keywords:	[Intracellular Cytokine Staining] [human PBMCs] [COVID-19] [SARS-CoV2] [Coronavirus] Manuscripts:					
Organizations:	VIB/KU Leuven, Leuven, Leu Babraham Institute, Babraha	iven (Belgium) am Institute, Cambridge, Ca	ambridge (United Kingdom)			
Purpose:	Analysis of cytokine producti	ion by PBMC from COVID-1	9 patients			
Conclusion:	Cytokines are produced by PBMC from SARS-CoV2-infected patients					
Comments:	None					
Funding:	VIB, KU Leuven					
Quality control:	Unstimulated controls, health	ny controls, Automated com	pensation			

Experiment variables

Conditions



export_COVID19 samples 23_04_20_ST3_COVID19_HC_001 ST3 230420_017_Live_cells.fcs · export_COVID19 samples 23_04_20_ST3_COVID19_HC_005 ST3 230420_016_Live_cells.fcs · export_COVID19 samples 23_04_20_ST3_COVID19_HC_006 ST3 230420_015_Live_cells.fcs · export_COVID19 samples 23_04_20_ST3_COVID19_HC_007 ST3 230420_014_Live_cells.fcs · export_COVID19 samples 23_04_20_ST3_COVID19_HC_008 ST3 230420_013_Live_cells.fcs · export_COVID19 samples 23_04_20_ST3_COVID19_HC_009 ST3 230420_013_Live_cells.fcs · export_COVID19 samples 23_04_20_ST3_COVID19_HC_009 ST3 230420_012_Live_cells.fcs

> export_COVID19 samples 21_04_20_ST3_COVID19_ICU_changedW_019_O ST3 210420_040_Live_cells.fcs · export_COVID19 samples 21_04_20_ST3_COVID19_ICU_changedW_027_O ST3 210420_039_Live_cells.fcs · export_COVID19 samples 21_04_20_ST3_COVID19_ICU_changedW_036_O ST3 210420_035_Live_cells.fcs · export_COVID19 samples 21_04_20_ST3_COVID19_W_033_O ST3 210420_036_Live_cells.fcs · export_COVID19

Flow Exercise



Mass Spec

- General purpose method to measure the accurate masses of small molecules
- Can be used to identify
 - Proteins (plus modifications)
 - Metabolites
 - Sugars
 - Nucleotides
 - Amino Acids
 - Lipids

Protein Mass Spec



Peptides



Too Big

Too Many

Non-specific

A peptide

Protein Mass Spec Workflow



Protein Mass Spec Results



MS features view

https://www.maxquant.org/

Protein Identification



http://www.ohri.ca/proteomics/

Post Translational Modifications

- When doing tandem mass spectrometry you can also identify modified peptides
- Phosphorylation
- Acetylation
- Methylation
- Palmitylation
- Acylation
- Ubiquitination
- etc.



High throughput proteomics

J Proteome Res. 2019 May 3;18(5):2346-2353. doi: 10.1021/acs.jproteome.9b00082.
 Epub 2019 Apr 12.

Evosep One Enables Robust Deep Proteome Coverage Using Tandem Mass Tags while Significantly Reducing Instrument Time

Jonathan R Krieger, Leanne E Wybenga-Groot, Jiefei Tong, Nicolai Bache¹, Ming S Tsao²³⁴, Michael F Moran⁵

30 samples per day Evosep workflow, >12 000 proteins were identified in 48 h of mass spectrometry time

J Proteome Res. 2021 May 7;20(5):2964-2972. doi: 10.1021/acs.jproteome.1c00168. Epub 2021 Apr 26.

TMTpro-18plex: The Expanded and Complete Set of TMTpro Reagents for Sample Multiplexing

Jiaming Li ¹, Zhenying Cai ² ³, Ryan D Bomgarden ⁴, Ian Pike ⁵, Karsten Kuhn ⁵, John C Rogers ⁴, Thomas M Roberts ² ³, Steven P Gygi ¹, Joao A Paulo ¹

> Nat Protoc. 2018 Jul;13(7):1632-1661. doi: 10.1038/s41596-018-0006-9.

Reproducible workflow for multiplexed deep-scale proteome and phosphoproteome analysis of tumor tissues by liquid chromatography-mass spectrometry

Philipp Mertins ¹ ² ³, Lauren C Tang ¹, Karsten Krug ¹, David J Clark ⁴, Marina A Gritsenko ⁵, Lijun Chen ⁴, Karl R Clauser ¹, Therese R Clauss ⁵, Punit Shah ⁴, Michael A Gillette ¹, Vladislav A Petyuk ⁵, Stefani N Thomas ⁴, D R Mani ¹, Filip Mundt ¹, Ronald J Moore ⁵, Yingwei Hu ⁴, Rui Zhao ⁵, Michael Schnaubelt ⁴, Hasmik Keshishian ¹, Matthew E Monroe ⁵, Zhen Zhang ⁴, Namrata D Udeshi ¹, Deepak Mani ¹, Sherri R Davies ⁶, R Reid Townsend ⁶, Daniel W Chan ⁴, Richard D Smith ⁵, Hui Zhang ⁴, Tao Liu ⁵, Steven A Carr ⁷

10,000 proteins per sample 37,000 phosphosites per sample

Expanding Mass Spec Technology

Ultra-high sensitivity mass spectrometry quantifies single-cell proteome changes upon perturbation

O Andreas-David Brunner, Marvin Thielert, Catherine G. Vasilopoulou, Constantin Ammar, Fabian Coscia, Andreas Mund, Ole B. Hoerning, Nicolai Bache, Amalia Apalategui, Markus Lubeck, Sabrina Richter, David S. Fischer, Oliver Raether, Melvin A. Park, Florian Meier, Fabian J. Theis, Matthias Mann

1,400 proteins measured from a single cell

Mass Spectral Imaging

Fix a sample to a surface and then do scanning Ionisation over it to get a spectrum for each point.

You can then pick any fragment and image its distribution over the original sample


Data Repositories for Proteomics Mass Spec



Data Repositories for Proteomics Mass Spec

- Varying amounts of experimental annotation
- Good description of processing and preparation
- Raw data files available
 - Mass spec still uses a lot of proprietary vendor file formats
 - Open mzML format is defined but often not used
 - Converters exist but often lose information.





Dataset Identifier	Title 🔶	Repos 🔶	Species 🔶	Instrument 🔶	Publication \$	LabHead \$	Announce Date	Keywords	\$
PXD026962	Multi-omic Profiling of Plasma Identify Biomarkers and Pathogenesis of COVID-19 in Children	iProX	Homo sapiens	QTRAP 6500+	Dataset with its publication pending	Xi Zhou	2021-06-28	Multi-omic Profiling, COVID-19, Childr	en,
PXD026928	Mycoplasma gallisepticum WhiA knockdown and overexpression	PRIDE	Mycoplasma gallisepticum S6	Q Exactive Plus	Dataset with its publication pending	Gleb Fisunov	2021-06-25	mycoplasma, transcription factor, Whi overexpression, knockdown,	٨,
PXD022361	Recombinant SWATH library for identification of low abundant human plasma proteins	MassIVE	Homo sapiens	TripleTOF 6600	Ahn et al. (2021)	Prof Mark S. Baker	2021-06-24	SWATH, Recombinant Protein, Plasma Proteome,	
PXD021581	Prognostic accuracy of Mass Spectrometric Analysis of Plasma in COVID-19	PRIDE	Homo sapiens	LTQ Orbitrap Velos	Dataset with its publication pending	Giuseppe Palmisano	2021-06-21	Sars-cov-2, Covid-19, COVID-19, SAR CoV-2, Mass spectrometry, Biomarker, Plasma, Prognosis,	.S-

	PXD026962			
PXD026962 is an original	dataset announced via ProteomeXchange.			
ataset Summary				
Title	Multi-omic Profiling of Plasma Identify Biomarkers and Pathogenesis of COVID-19 in Children			
Description	Although people of all ages are susceptible to COVID-19, children usually develop less severe disease than adults Little is known about the pathogenesis of COVID-19 in children. Herein, we conduct the plasma proteomic and metabolomic profiling of a cohort of COVID-19 children patients with mild symptoms, and uncovered that many proteins involved in immune response are significantly up-regulated in a stronger extent than in adults with COVID Interestingly, more molecules involved in protective processes of reducing inflammation are also stimulated to antagonize the deleterious effect in both proteomic and metabolomic levels. By developing a machine learning-bas pipeline, we prioritize two set of biomarker combinations, and identify 5 proteins and 5 metabolites as potentially children-specific biomarkers. Further experiments demonstrate these protective metabolites not only inhibit the expression of pro-inflammatory factors, but also suppress the viral replication. Taken together, our study not only discover the protective mechanisms in children with COVID-19, but also shed light on potential therapies targets for treating COVID-19.			
HostingRepository	iProX			
AnnounceDate	2021-06-28			
AnnouncementXML	Submission_2021-06-28_01:26:39.414.xml			
DigitalObjectIdentifier				
ReviewLevel	Peer-reviewed dataset			
DatasetOrigin	Original dataset			
Repository Support	Unsupported dataset by repository			
Primary Submitter	Yang Qiu			
SpeciesList	scientific name: Homo sapiens; NCBI TaxID: 9606;			

Project Information	\sim
Project ID	
IPX0002673000	
ProteomeXchange ID	
PXD026962	
Project Title	
Multi-omic Profiling of Plasma Identify Biomarkers and Pathogenesis of COVID-19 in Children	
XML File PX_IPX0002673000.xml	
Download All Files (7.55G)	
	Ťĸ.
IPro	

Metabolite Mass Spectrometry

- Similar concepts to protein mass spec
- Range of starting material
 - Serum
 - Urine
 - Cerebrospinal fluid
 - Saliva
- Different separations
- Up to 5000 different metabolites to find

Data Repository for Metabolomics Data

MoNA - MassBank of North America

- Reference spectra for biological molecules
 - Used for searching and quantitation



- Experimental datasets of MassSpec Studies
 - Used to answer biological questions
 - Also provides visualisations and tools

Metabolomics Workbench

Data for ST001140 Perform analyte scaling Perform sample normalization

(Analysis AN001870): Average values per metabolite and experimental factor (Units:uM)

Metabolite structure	All data	F1	F2	F3	F4
CE(16:0)	ME271966	2.67	2.94	2.03	1.56
CE(16:1)	ME271967	0.47	0.52	0.44	0.37
CE(17:0)	ME271968	0.06	0.06	0.05	0.03
CE(17:1)	ME271969	0.05	0.06	0.06	0.05
CE(18:0)	ME271970	0.52	0.66	0.46	0.38
CE(18:1)	ME271971	22.69	22.21	22.53	20.25
CE(18:2)	ME271972	85.96	85.68	95.16	74.09







Mass Spec Data Exercise



Imaging Analysis

- What can you measure with imaging?
- Cell structure and morphology

 In both live and fixed cells
- Targeted molecules (fluorescence microscopy)
 - Antibodies to proteins
 - Fluorescent fusion proteins
- Functional readouts
 - Redox state
 - pH





Types of Microscopy

- Light Microscopy
 - Sample is illuminated, some light goes to the viewer
 - Biological samples are generally clear, so hard to see
 - Can use stains (often toxic) or reflection or phase shift to see better



https://zeiss-campus.magnet.fsu.edu/articles/basics/contrast.html

Types of Microscopy

- Fluorescence Microscopy
 - Uses molecules which excite at one wavelength and emit at another
 - Allow the tagging of specific biological molecules
 - Confocal microscopes allow clear views of a single plane in the sample



DOI:10.3390/ijms20082033

Ultra-plex fluorescence imaging









Simon Walker – BI Imaging Facility

Types of Microscopy

- Electron Microscopy
 - Fixed and processed samples only (not live)
 - Very high resolution





High Content Imaging

- Microscopy traditionally operated on small numbers of individual samples
- Improved equipment and automation now allows for more ambitious studies
 - 384 well plates
 - 30 images per well
 - 5 different markers
 - Thousands of images
 - Hundreds of measured features per cell





Imaging Flow Cytometry High Content Imaging from Flowed Cells



High Content Applications

- Screening for drugs with specific phenotypic effects
- Measuring CRISPR library phenotypes
- Measuring RNAi library phenotypes

BioImage Archive

Submit Help 🔻 Browse

Release Date: 1 July 2019

Accession

S-BIAD9

About us 🔻

BIOSTUDIES / BIOIMAGES / S-BIAD9 The BioImage that are useful archiving servi including adde IDR and Tissue

Home

Recordings of locomotor behaviour in wild-type and mutant Caenorhabditis elegans

Akihiro Mori 1, Yee Lian Chew 2, Laura Grundy 1, Eviatar Yemini 3, Andr? E.X. Brown 4, William R. Schafer 1

Ø Data files

¹ Division of Neurobiology, MRC Laboratory of Molecular Biology ^{2 3 4 5} Current Address: Illawarra Health and Medical Research Institute & School of Chemistry and Molecular Bioscience, University of Wollongong, Wollongong, Australia ⁶ Current Address: Columbia University, New York, NY USA ⁷ Current Address: Imperial College London, UK

🖉 Data files		
Show 5 v entries	E Search:	
Name		

{JSON}

→PageTab↓ ④FTP

×

Description The analysis of behaviou sh	ow 5 v entries							+ Search:	
genes affecting nervous system functio	Name	🔺 Size 🔶	Section 🔶	Sample Name	Protocol REF 🝦	Assay Name	🔷 Raw Data File	Comment[Data Repository]	Comment[ventral side]
have been identified whose loss of func					Tracking of				
high-content, quantitative phenotypes	A02047 A02047 on food L 2012 02 00 10 25 26 1 1 co	129.4	Sereen A	402047	wild-type and	A02047 mutant replicat	0 A02047 A02047 on food L 2012 02 00 10 25 26 1 1 0	i Diastudias	Loft
underlying this database, consisting of	A6544, A6544, 01-1000 - 2015 05-04 - 10-52-20 - 1-1-36	MB	ScreentA	AQ2947	Caenorhabditis	AQ2947_Initiant_replicat	e Mőzadi Wőzadi Olliond r solts ős ővel solts solt solts solt solts solt solt	VI Biostudies	Len
wild-type controls. Each genotype is re					elegans Tracking of				
accessory files containing records of st	_	60.7			wild-type and				
represent a useful resource for investig	AQ2947_AQ2947_on_food_L_2012_02_0910_57_3882_se	g.avi MB	Screen A	AQ2947	mutant Caenorhabditis	AQ2947_mutant_replicat	e AQ2947_AQ2947_on_food_L_2012_02_0910_57_3882.a	vi Biostudies	Left
specific genes on locomotion.					elegans				
Study type high content screen	AQ2947_AQ2947_on_food_L_2012_02_0911_19_463_seg.avi	129.5	Screen A	AQ2947	Tracking of wild-type and mutant	AQ2947_mutant_replicat	e AQ2947_AQ2947_on_food_L_2012_02_0911_19_463.avi	Biostudies	Left
Key words C. elegans, locomotor be		MB			Caenorhabditis elegans				
Study Organism Caenorhabditis elegan					Tracking of				
	AQ2947_AQ2947_on_food_L_2012_02_0911_4834_seg.av	86.2 MB	Screen A	AQ2947	mutant Caenorhabditis elegans	AQ2947_mutant_replicat	e AQ2947_AQ2947_on_food_L_2012_02_0911_4834.avi	Biostudies	Left
					Tracking of				
E	AQ2947_AQ2947_on_food_L_2012_02_0912_13_275_se	g.avi 346.0 MB	Screen A	AQ2947	wild-type and mutant Caenorhabditis elegans	AQ2947_mutant_replicat	e AQ2947_AQ2947_on_food_L_2012_02_0912_13_2725.a	vi Biostudies	Left
c+	powing 1 to 5 of 10 105 optrios						Previous	2 3 4 5	2021 Next

Showing 1 to 5 of 10,105 entries

DI	R Studies Genes Phenotypes C	ell Lines	siRNAs	Antibodie	s Compo	unds Or	ganisms	About
ub	lie Public data 🔹							
re	Tags Shares	Index	Field#1	~				
		3	1	2	3	4	5	6
lid	r0018-neff-histonathology/experimentA 248	^ ^			- /4	1		- 1
id	r0019-sero-nfkappab/screenA 7				1117	para 1	2.	12
id	r0020-barr-chtoo/screenA 4	в		. e.	:: 7	P	-	1
id	r0021-lawo-pericentriolarmaterial/experimentA	A 10		× _	1/3	2+ J -	11, 2	47
id	r0022-koedoot-cellmigration	C.14	_	.*	C3	117	1 Barris	34
id	r0022-koedoot-cellmigration/screenA 524		-		0.1.	5. 501	3 50 50	1.
id	r0022-koedoot-cellmigration/screenB 152		*		N CAN	2.3	2	
id	r0023-szymborska-nuclearpore/experimentA	55 D	* v.		120	i geti	14 a.	100
id	r0025-stadler-proteinatlas/screenA 3			10		MA TH	bry in	1971
id	r0026-weigelin-immunotherapy/experimentA 1	8 =	*-	r	ier	2-2-2	t.	the fight
id	r0027-dickerson-chromatin/experimentA 8				1 -2	5.6	Nest 1	
id	r0028-pascualvargas-rhogtpases		- 1	· /	· · · ·	Ary Sich	うる	. 5
id	r0028-pascualvargas-rhogtpases/screenA 4	F	5 .		the a	the the	10 15	1
id	r0028-pascualvargas-rhogtpases/screenB 4	3				NV Y	18 185 ·	Ya
id	r0028-pascualvargas-rhogtpases/screenC 4	G	10	2	1 - A		Kr J.	15 15-5
8	MDA-MB-231_siGENOME_1A		-	3 1		A . Y	11 2 -	
8	MDA-MB-231_siGENOME_1B		1	10 pt	" Pro	1 12		4 4 4
8	MDA-MB-231_siGENOME_2A	н	100.	0	- , ' ab			R 32-
8	MDA-MB-231_siGENOME_2B	*		0-0	States 1	- As	10 187	1943
id	r0028-pascualvargas-rhogtpases/screenD 4	1	2 10	1.10	-7 · M	. 771	ALL CO	CA STA
id	r0030-sero-yap/screenA 10	4		-		- A - A	143 2	· 240
id	r0032-yang-meristem/experimentA 115	c			R	No at	A -may 77	Arriv
id	r0033-rohban-pathways/screenA 12	1		1. 0.	814	=	V.1 - 8	they we
id	Ir0034-kilpinen-hipsci/screenA 29			8	TAKE	- AV	-17-6	12
id	r0035-caie-drugresponse/screenA 55	к			H Di	and the	wet	-
id	r0036-gustafsdottir-cellpainting/screenA 20	4		-		A ATAL	TATE	2-
IC	r0037-vigilante-hipsci/screenA 69		E		2	54 1		N.
IC	r0038-heid-kidneylightsheet	7	1	14 14	h	1 11-	-1 1-	~. Er
10	r0038-heid-kidneylightsheet/experimentA 7		34	L13.21	-1 1	14	· ·····	59 7
10	r0038-heid-kidneylightsheet/experimentB 4	> ~ M	1 m	Notice -	in a	12	16 53 6	
	Field positions in well	XM			2 01 1	the .	: La Tal	
			.te	21 M	1 34	AN 1- 1	3-64	14.1
		o	and the		in the second	4 4 4 4 4	TRAT	2 inte
			10	The Real	- We de	W.	123 1 1	40
		P	PER E	The st	a la	*.	the state	20
				1 Kenter		X	1.2 -	14 m/4

-

Explo

Image Data Resource MDA-MB-231_siGENOME_1A [0451 1A Well G10 Field #1] Image Info Image ID: 2969393 Well: 1250727 Owner: Public data Acquisition Date: 2016-05-28 03:18:27 Import Date: 2016-09-29 17:14:15 Dimension (XY): 667 x 500 Pixels Type: uint16

Pixels Type:uint16Pixels Size (XYZ) (µm):0.65 x 0.65 x -Z-sections:1Timepoints:1Channels:Hoechst, AlexaFluor568, Phalloidin488.

AlexaFluor647

Imaging Data Exercise

