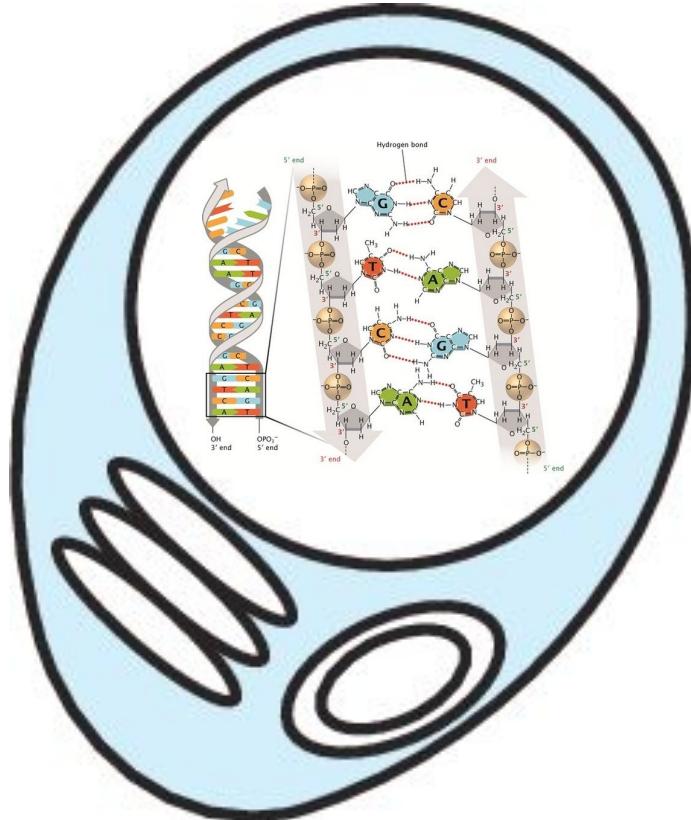


Computational approaches for integrative cancer genomics

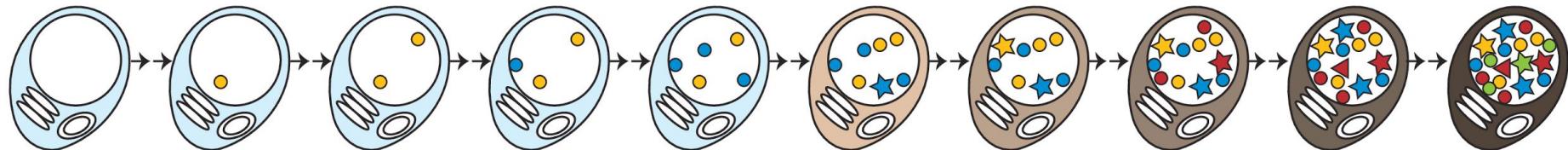
Christian Pérez Llamas

18th Dec 2015

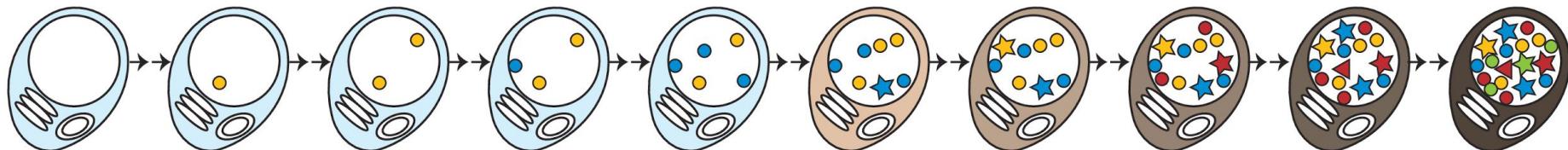
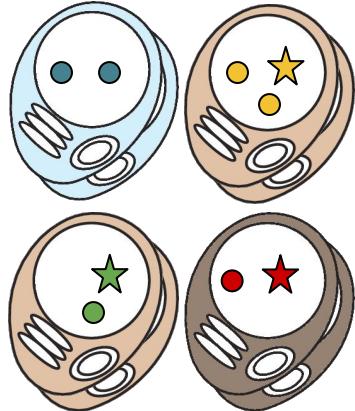
Integrative Cancer Genomics



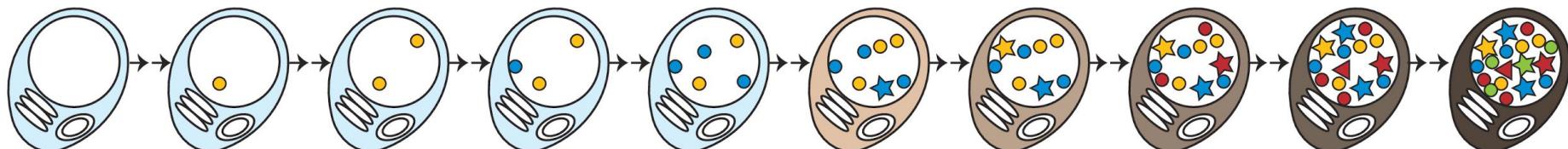
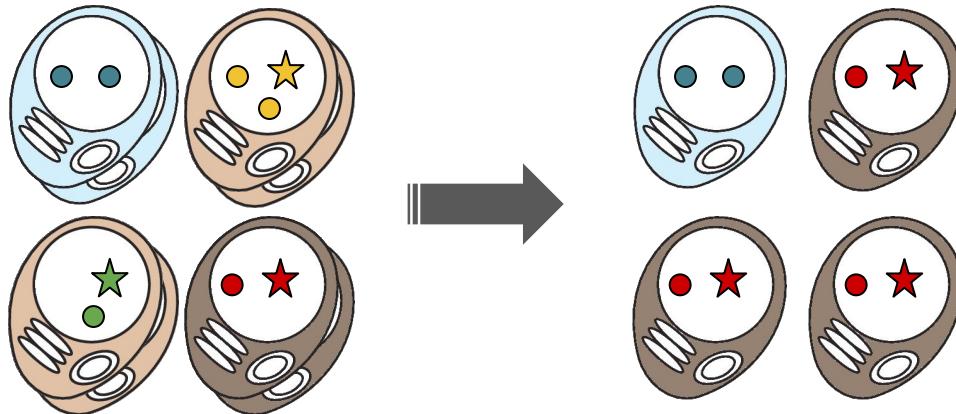
Cancer development



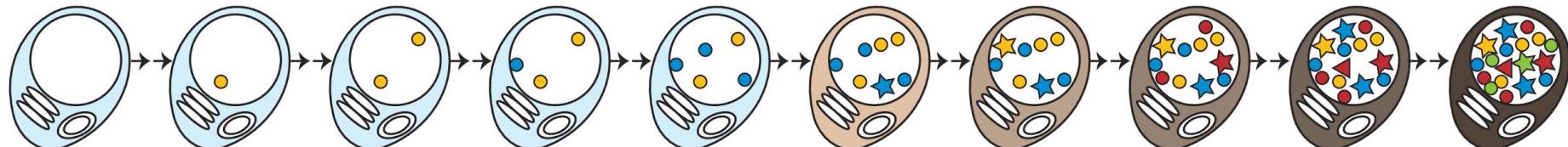
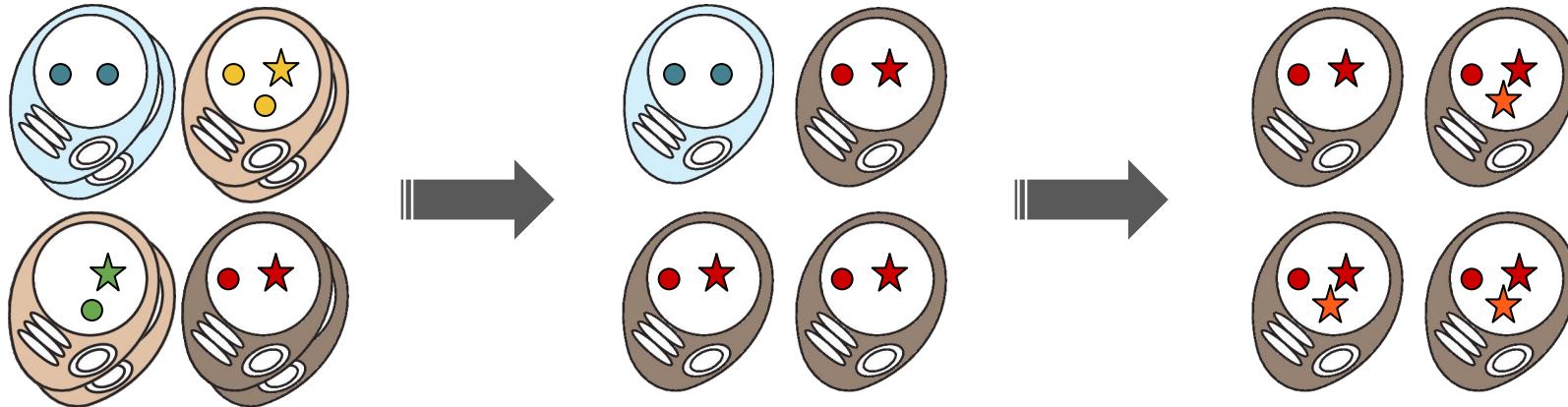
Cancer development



Cancer development



Cancer development



Genomic alterations

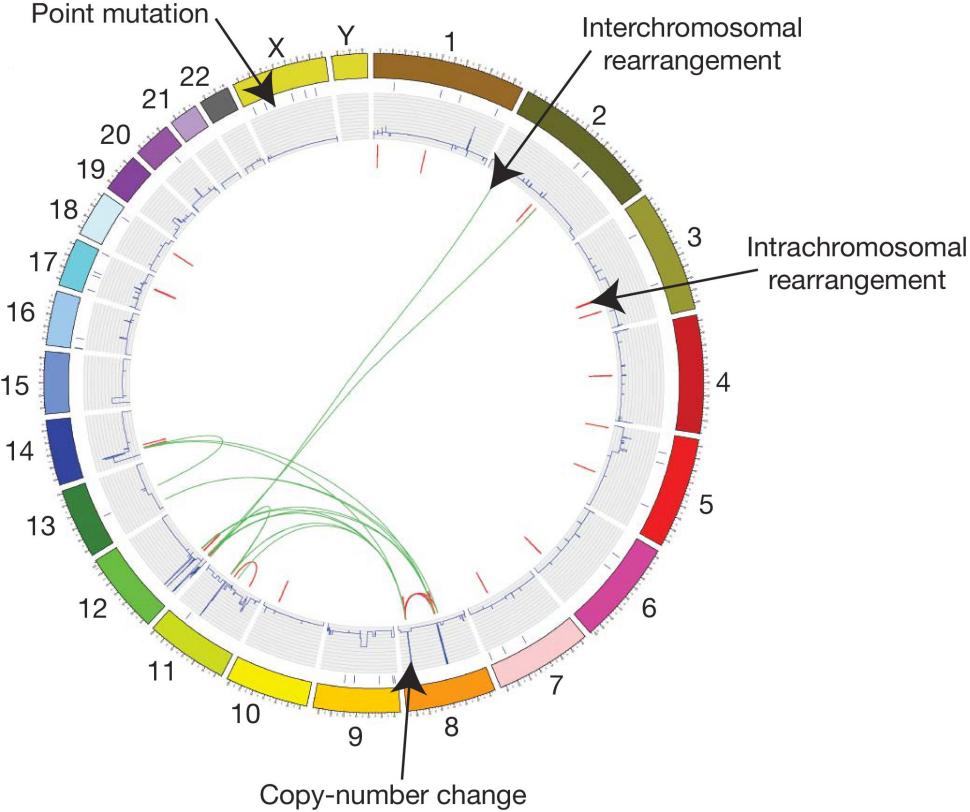


Somatic alterations

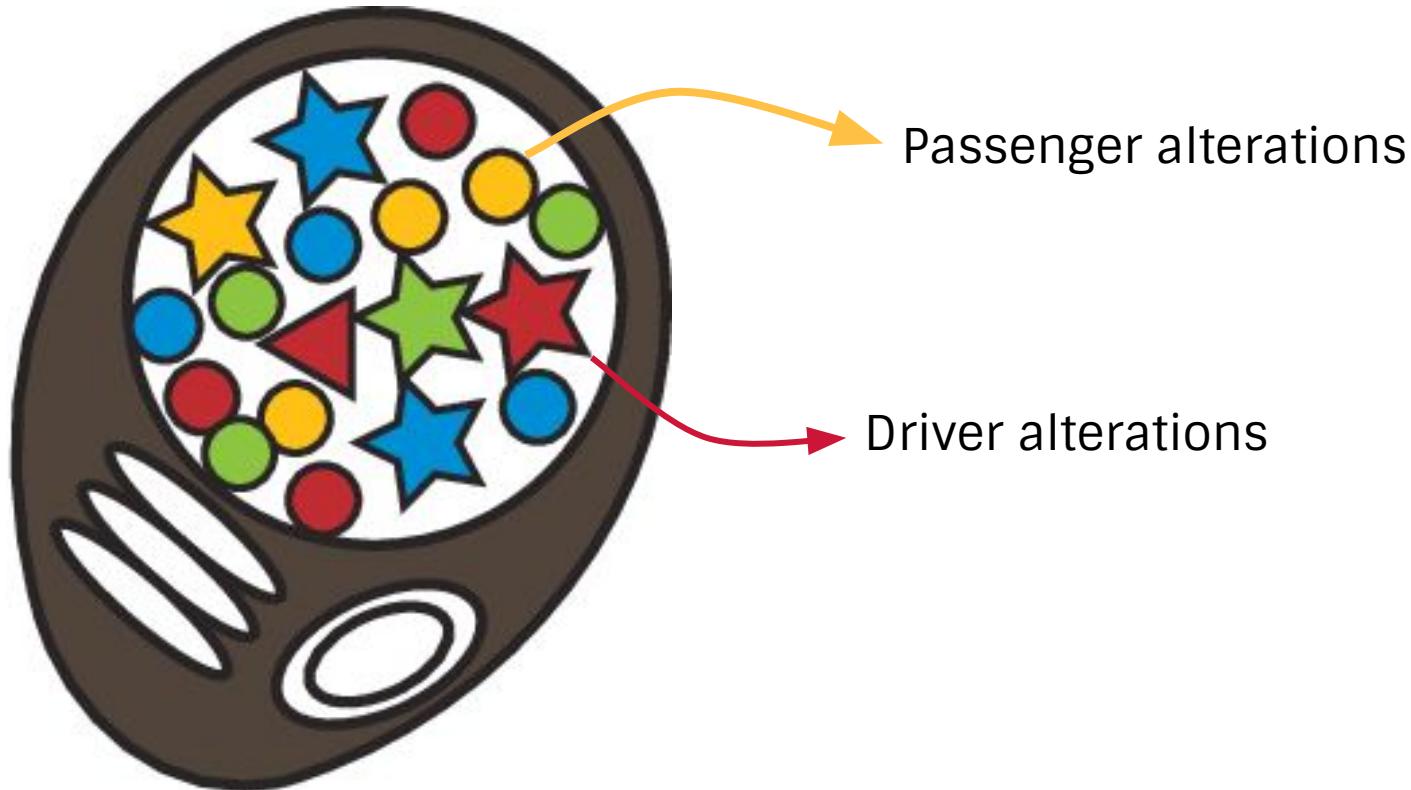
Germline alterations



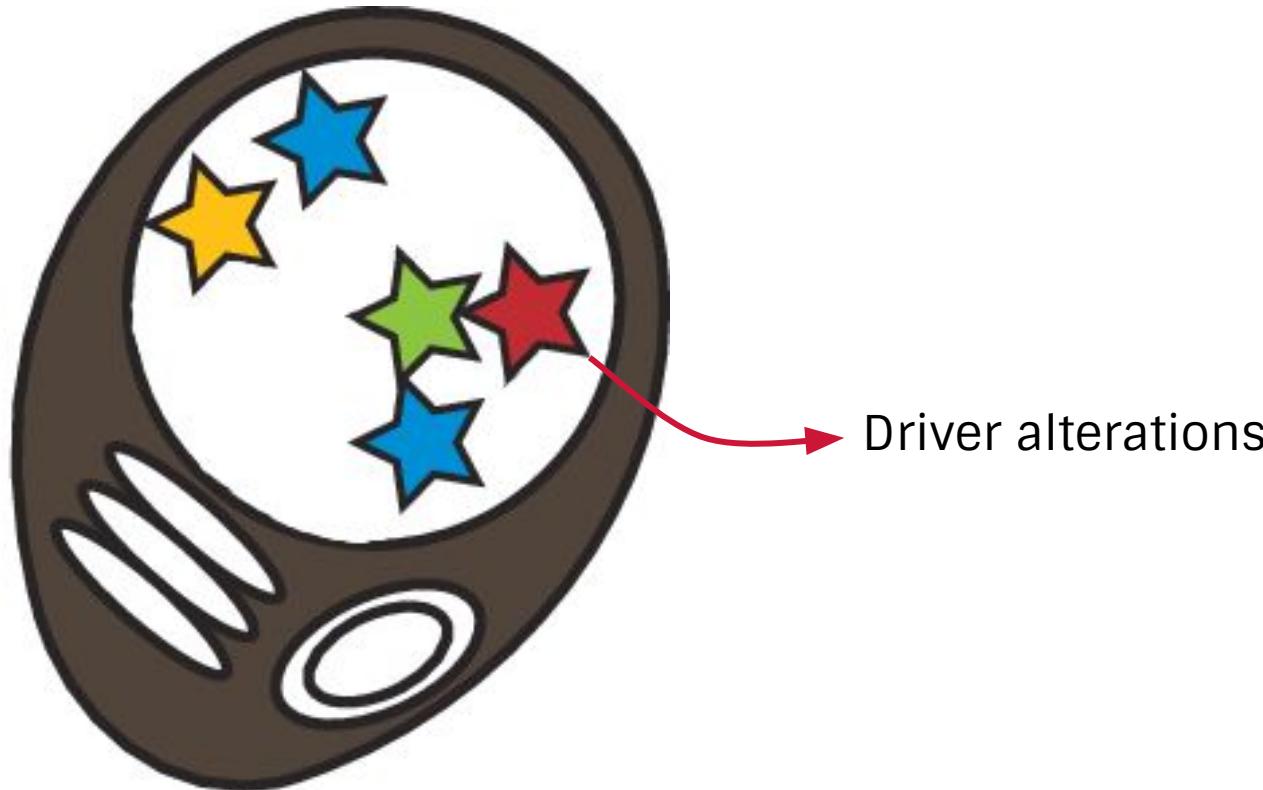
Genomic alterations



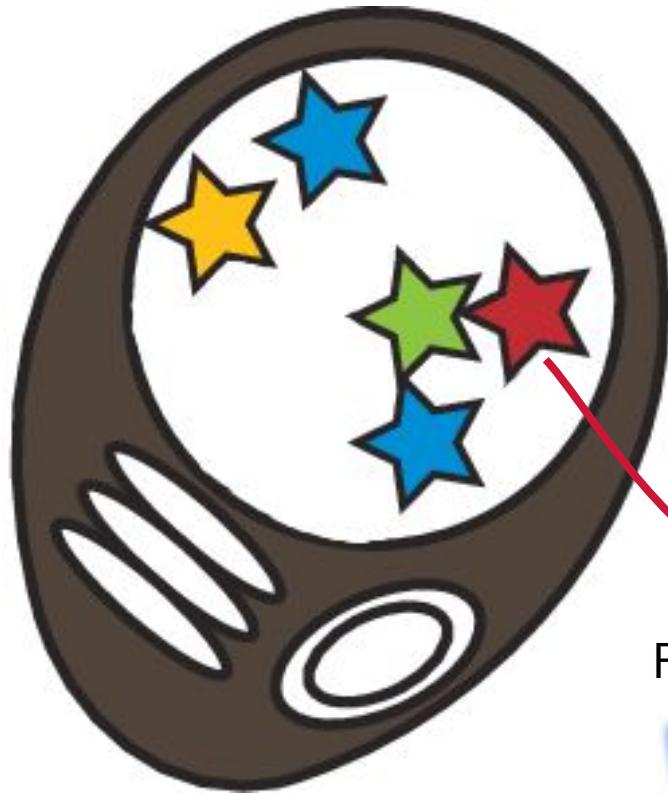
Genomic alterations



Identification of driver alterations



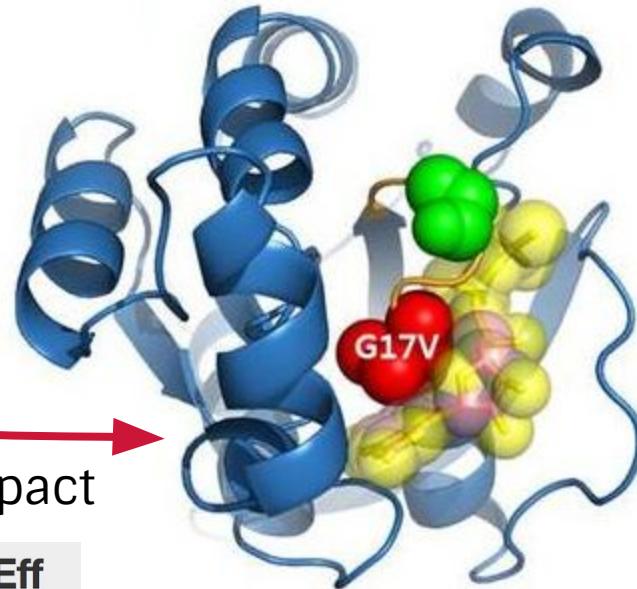
Identification of driver alterations



Functional Impact

Ve!P

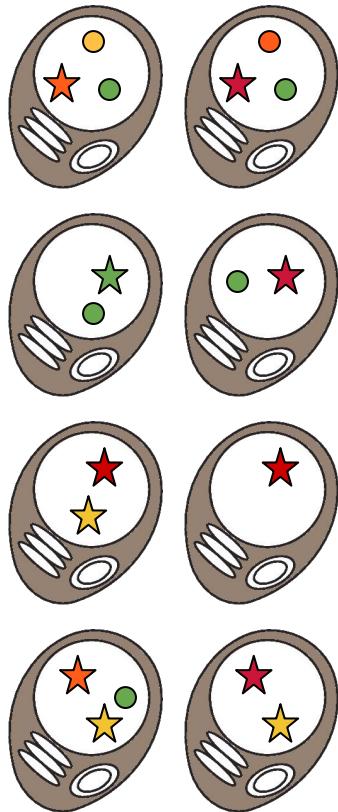
SnpEff
Genetic variant annotation
and effect prediction toolbox.



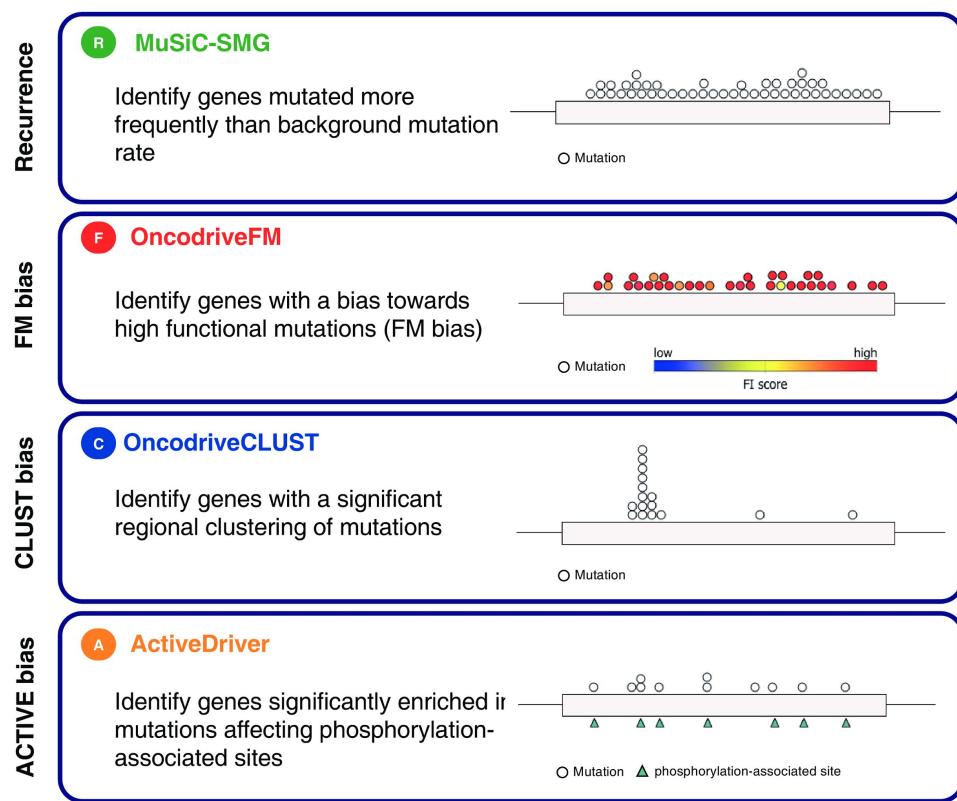
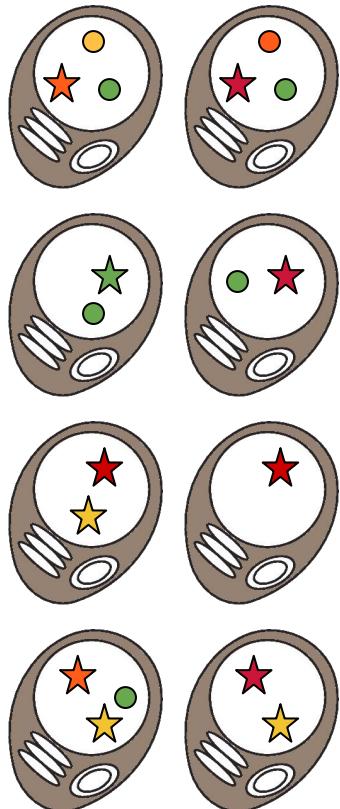
Identification of driver alterations



Identification of driver alterations



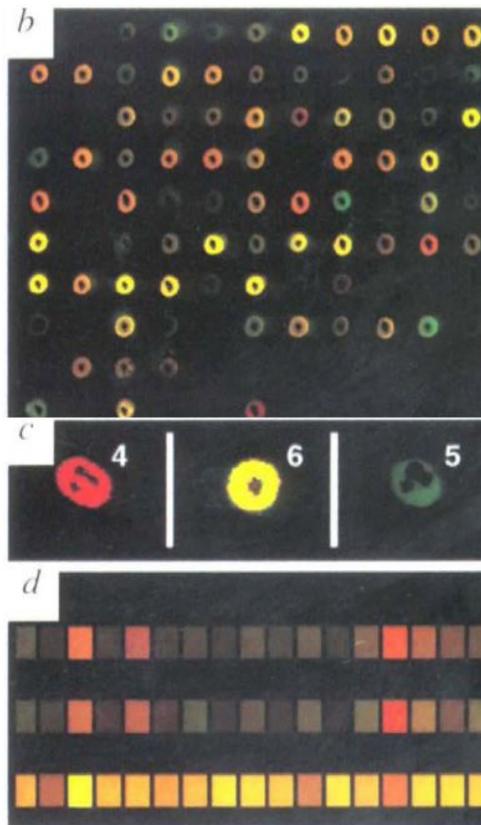
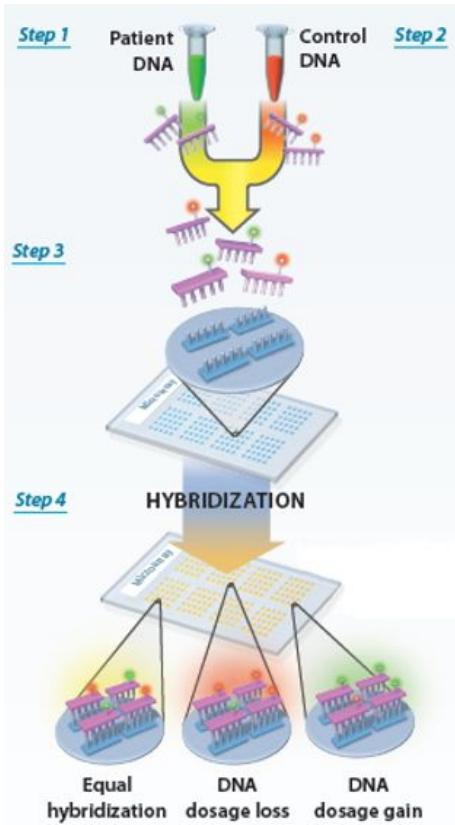
Identification of driver genes



Advances in Technology

	Begin	End	+10 years
Human Genome Project	Genome sequencing		
	Cost to Generate a Human Genome Sequence	-\$1 billion	-\$10-50 million
	Time to Generate a Human Genome Sequence	-6-8 years	-3-4 months
	Human Genome Sequences	0	1
	Genome Sequence Data		
	Total DNA Bases in GenBank	-49 Mbases	-31 Tbases
	Human Single-Nucleotide Polymorphisms	-4.4 K	-3.4 M
	Genomic Medicine		
	Drugs with Pharmacogenomics Information on Label	4	46
			106

Microarrays



Gene Expression Profiling

Comparative Genomic Hybridization

Chromatin Immunoprecipitation

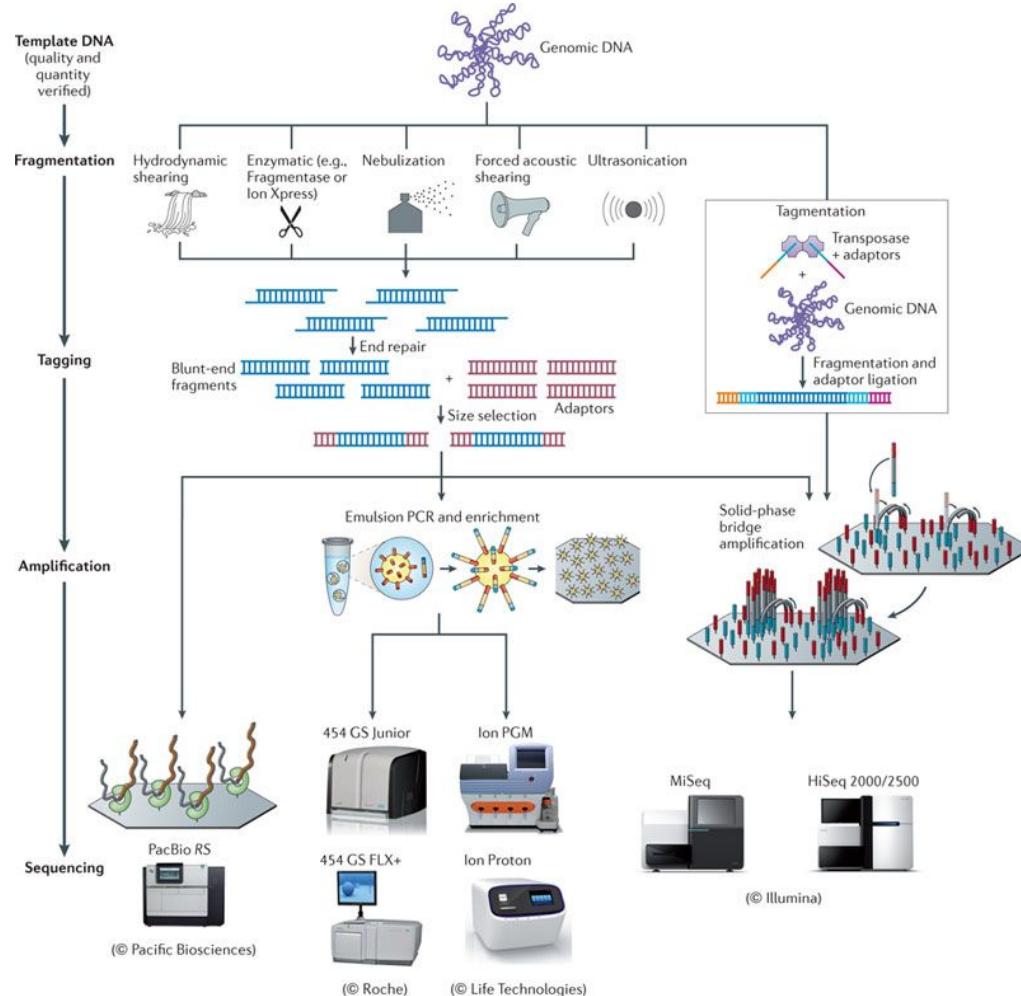
SNP detection

Next Generation Sequencing

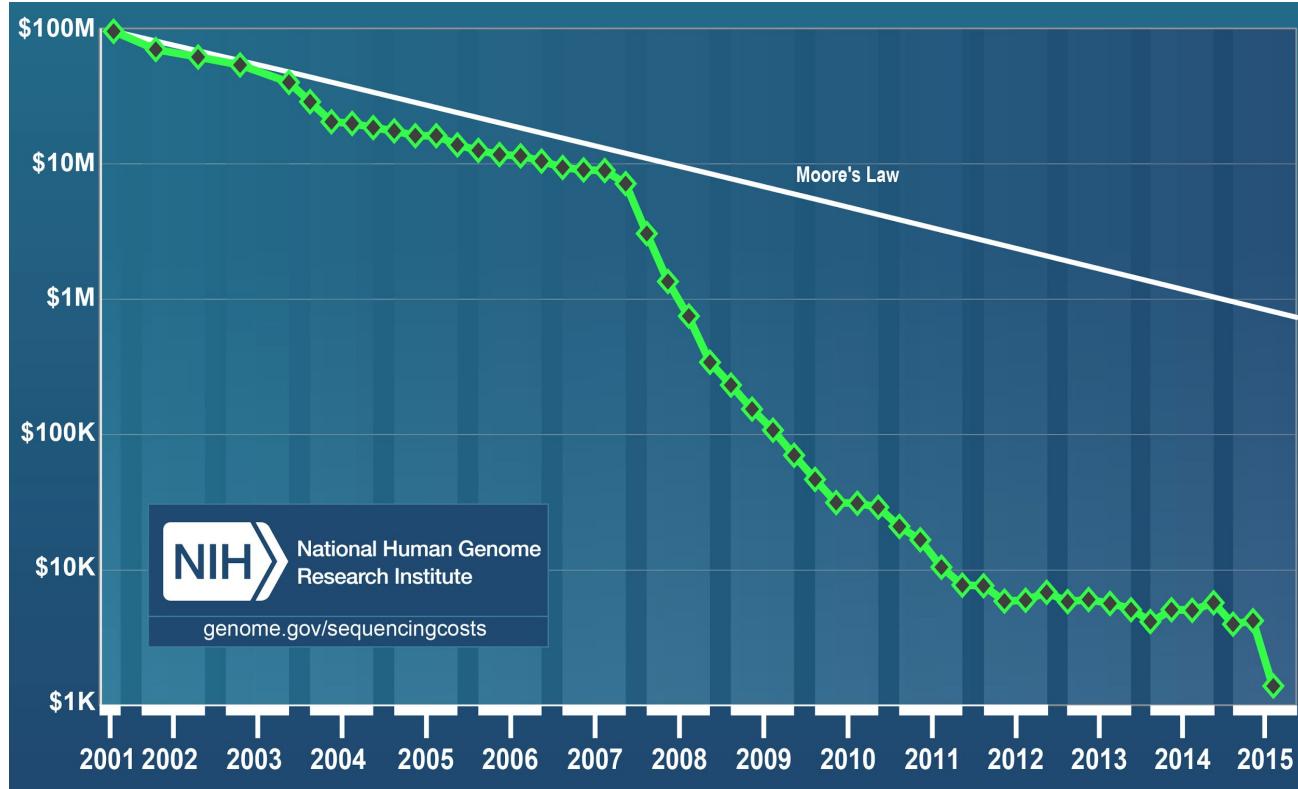
Whole Genome Sequencing

Whole Exome Sequencing

Whole Transcriptome Sequencing



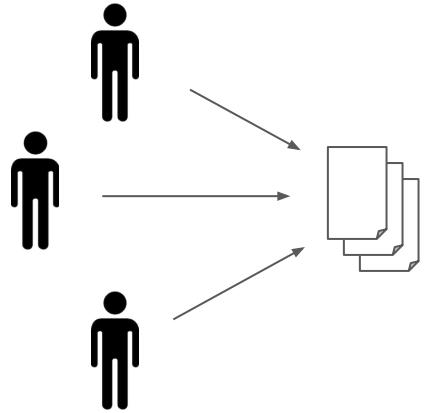
Next Generation Sequencing



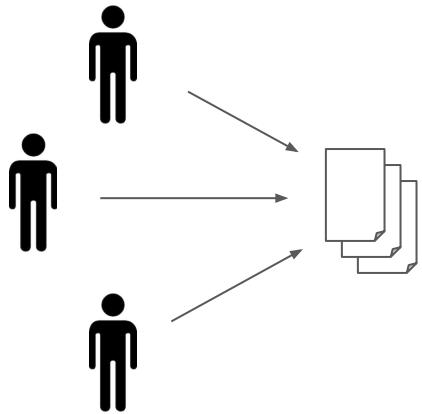
Integration of data



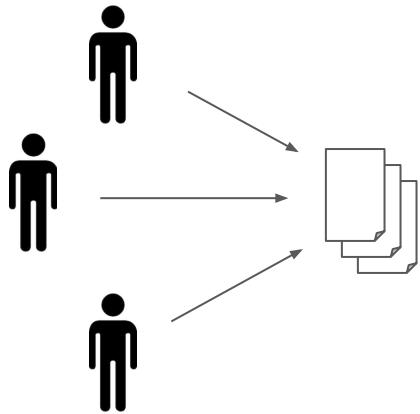
Integration of data



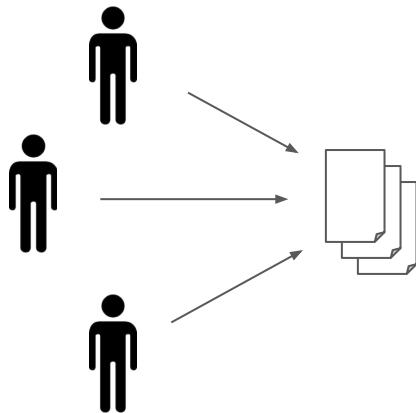
Integration of data



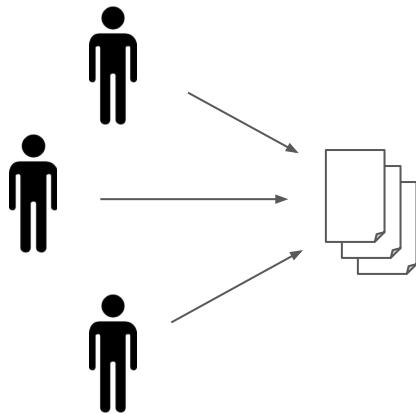
Integration of data



Integration of data



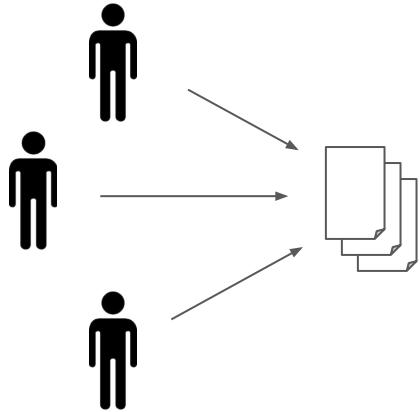
Integration of data



UCSC Genome Bioinformatics



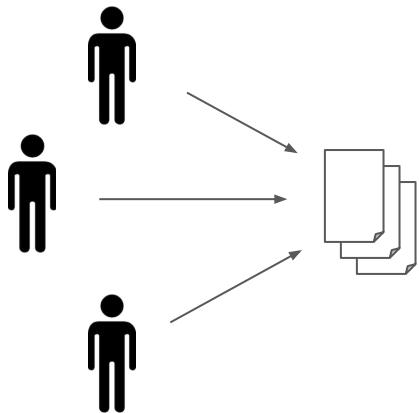
Integration of data



UCSC Genome Bioinformatics



Integration of data

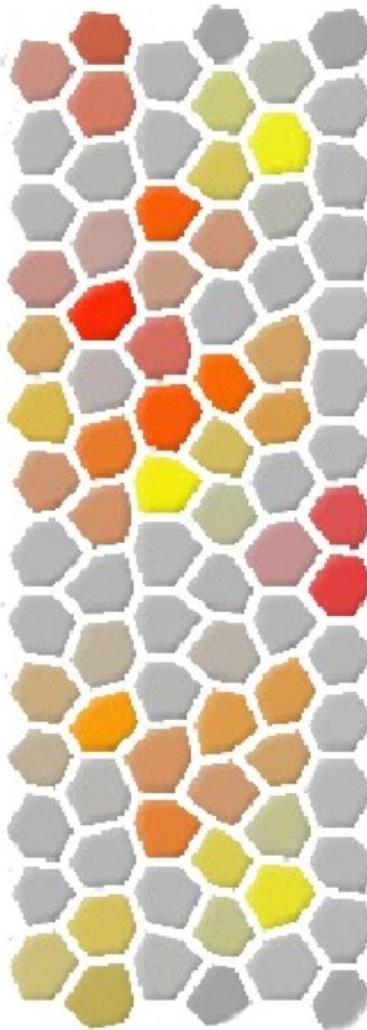


UCSC Genome Bioinformatics



e!Ensembl







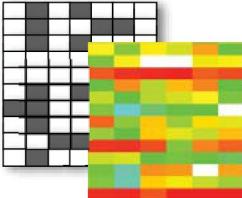
```
graph LR; A[Data] --> B[Analysis]; B --> C[Visualization]
```

Data

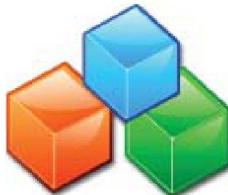
Analysis

Visualization

Matrices



Modules



Tables



Data

Analysis

Visualization

Binary Data Matrix (BDM)				
	S1	S2	S3	S4
ENSG000000000003	0	1	0	1
ENSG000000000005	-	1	0	0
ENSG00000000419	1	0	1	0
ENSG00000000457	0	1	0	0

Continuous Data Matrix (CDM)				
	S1	S2	S3	S4
ENSG000000000003	0,01	1,00	0,40	0,94
ENSG000000000005	-	1,00	1,00	1,00
ENSG00000000419	1,00	0,95	0,98	1,00
ENSG00000000457	1,00	1,00	1,00	1,00

Tab Separated Values (TSV)			
id	symbol	chromosome	band
ENSG000000000003	TSPAN6	X	q22.1
ENSG000000000005	TNMD	X	q22.1
ENSG00000000419	DPM1	20	q13.13
ENSG00000000457	SCYL3	1	q24.2

Gene Matrix (GMX)		
bound_0h	bound_27h	bound_96h
genes bound at 0h	genes bound at 27h	genes bound at 96h
ENSG000000000005	ENSG00000000971	ENSG00000001036
ENSG00000000419	ENSG00000001167	ENSG00000003096
ENSG00000000457		ENSG00000004455

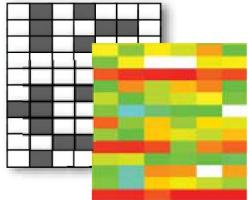
Two Columns Mapping (TCM)	
ENSG0000115204	GO:0000002
ENSG0000025708	GO:0000002
ENSG0000151729	GO:0000002
ENSG0000137074	GO:0000012
ENSG0000174405	GO:0000012
ENSG0000042088	GO:0000012
ENSG0000169621	GO:0000012
ENSG0000073050	GO:0000012
ENSG0000118245	GO:0000012

Indexed Mapping (IXM)	
This format is designed to use less space but it can not be easily edited. There is a command line converter that can generate it from other formats.	

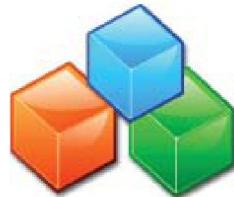
Gene Matrix Transposed (GMT)				
bound_0h	bound at 0h	ENSG000000000005	ENSG00000000419	ENSG00000000457
bound_27h	bound at 27h	ENSG00000000971	ENSG00000001167	
bound_96h	bound at 96h	ENSG00000001036	ENSG00000003096	ENSG00000004455

Data

Matrices



Modules



Tables



Analysis



BioMart Central Portal

Ensembl
Ensembl Bacteria
Ensembl Metazoa
Ensembl Protists
Ensembl Plants
Ensembl Fungi
Phytozome
Gramene
Europhenome
UniProt
InterPro
HGNC

CyanoBase
Wormbase
DroSpeGe
ArrayExpress DW
Eurexpress
HapMap
Dictybase
COSMIC
Rat Genome Database
GermOnLine
PRIDE
PepSeeker
VectorBase

HTGT
Cildb
Pancreatic Expression Database
Reactome
EU Rat Mart
Paramecium DB
International Potato Center (CIP)
Mouse Genome Informatics (MGI)

International Cancer
Genome Consortium



Integrative
Onco
Genomics

Genes significantly altered
Modules of genes significantly altered

Experiments
Tumour types

Up-regulation
Down-regulation
Gain
Loss

levels
alterations



Kyoto Encyclopedia of
Genes and Genomes

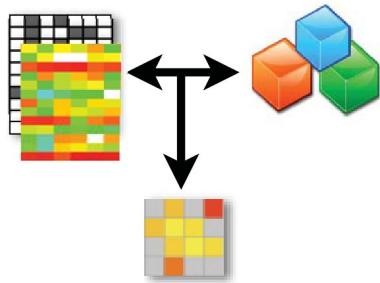
Pathways

Data

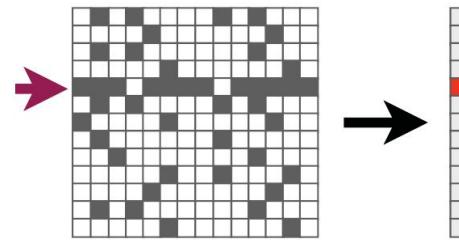
Analysis

Visualization

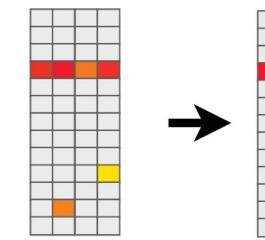
Enrichment Analysis

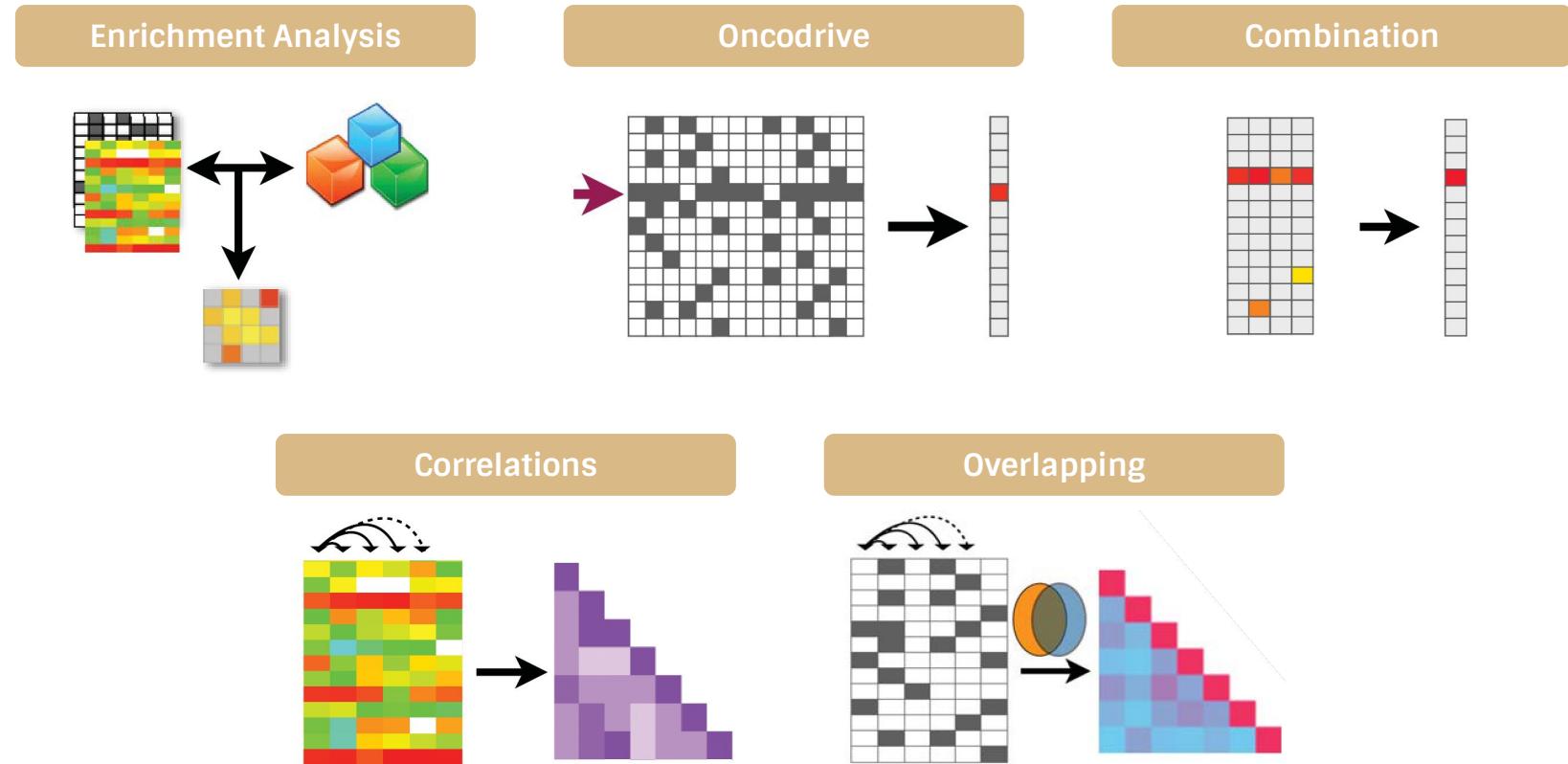


Oncodrive



Combination

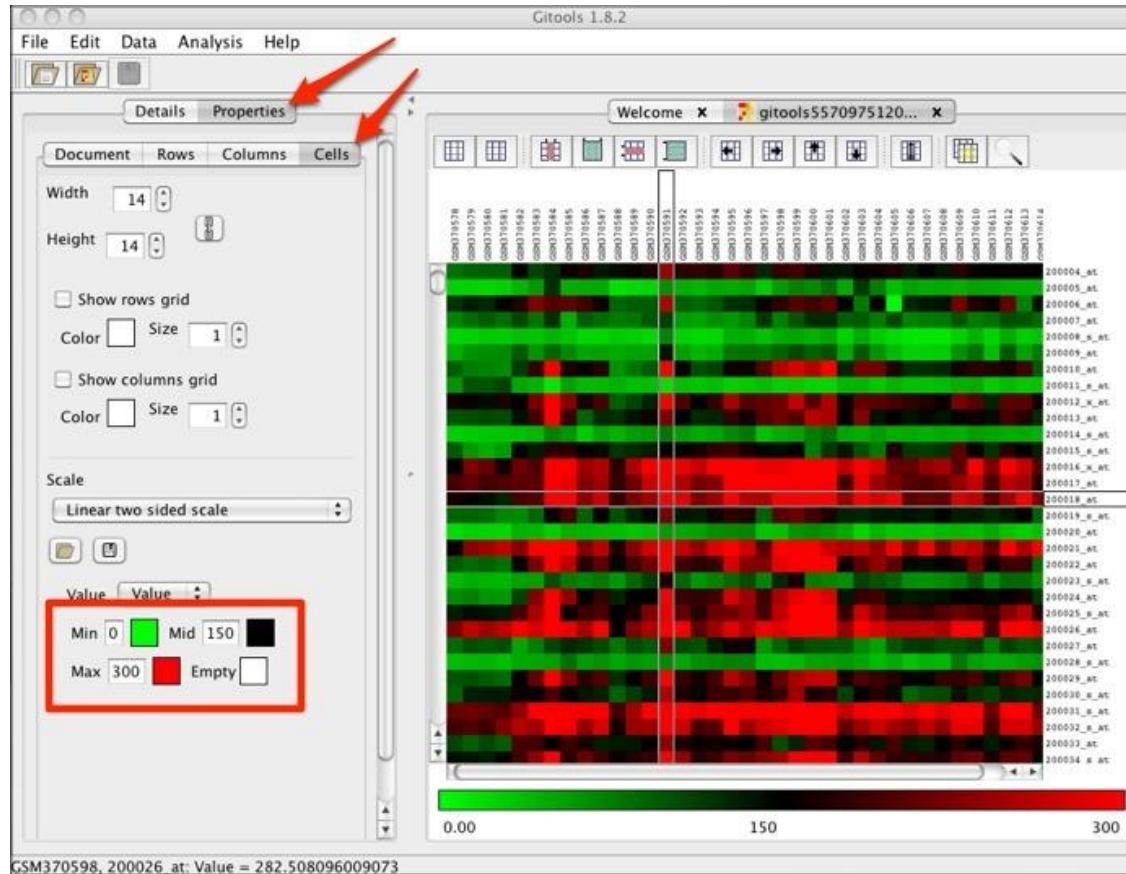




Data

Analysis

Visualization



→ Scale

→ Grid

→ Sizes

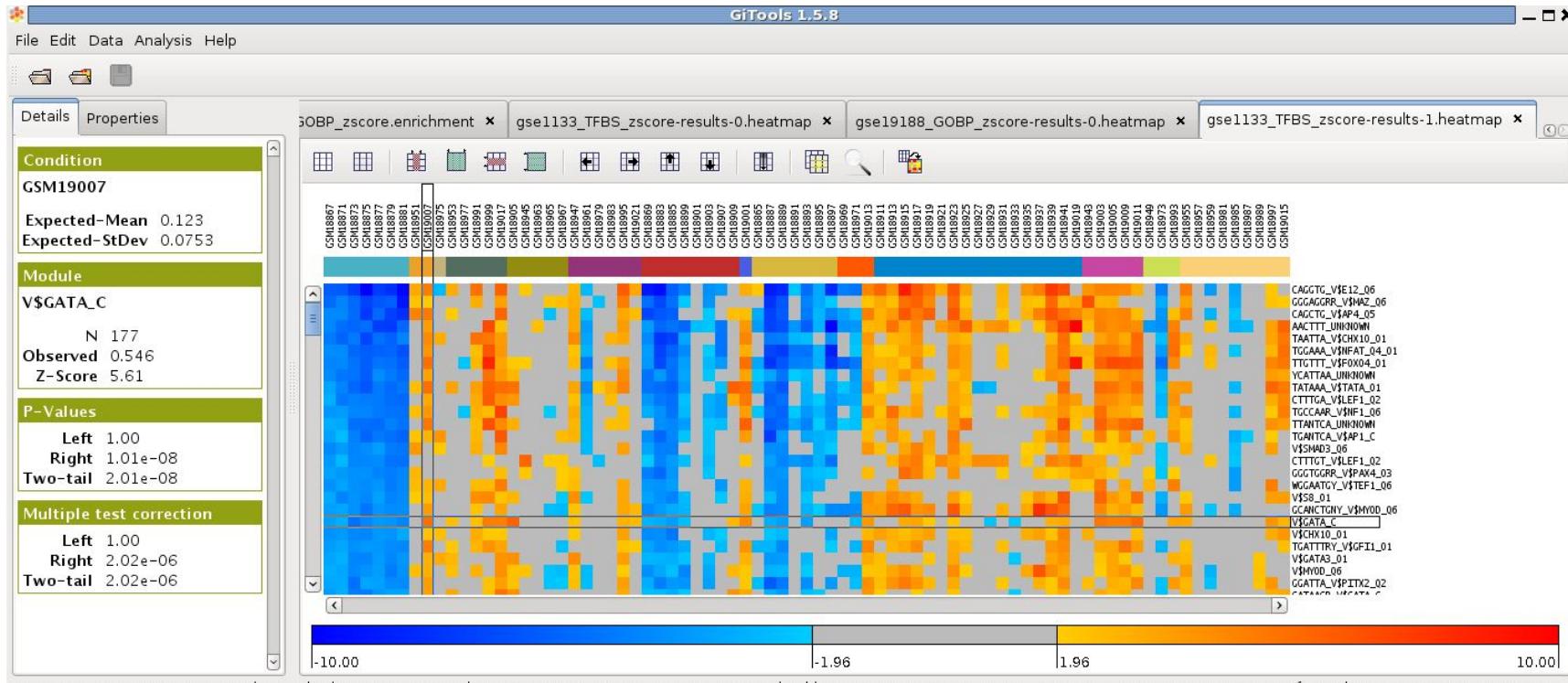
→ Font

→ Labels

Data

Analysis

Visualization



Explore the data by searching, filtering, sorting, and reordering columns and rows

Gitoools: Analysis and Visualisation of Genomic Data Using Interactive Heat-Maps.

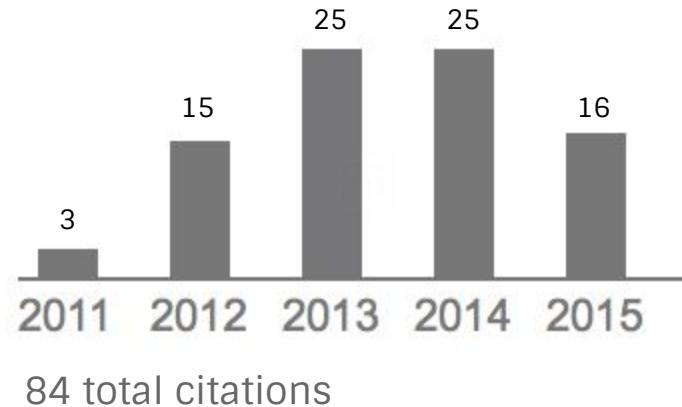
Perez-Llamas, C. & Lopez-Bigas, N. PLoS ONE 6, e19541 (2011)

DIANA miRPath v. 2.0: investigating the combinatorial effect of microRNAs in pathways.

A polycomb group protein is retained at specific sites on chromatin in mitosis.

Transcriptomic characterization of temperature stress responses in larval zebrafish.

温度刺激对斑马鱼仔鱼基因转录表达的影响.



● Users www.gitools.org



Gitoools: Analysis and Visualisation of Genomic Data Using Interactive Heat-Maps.

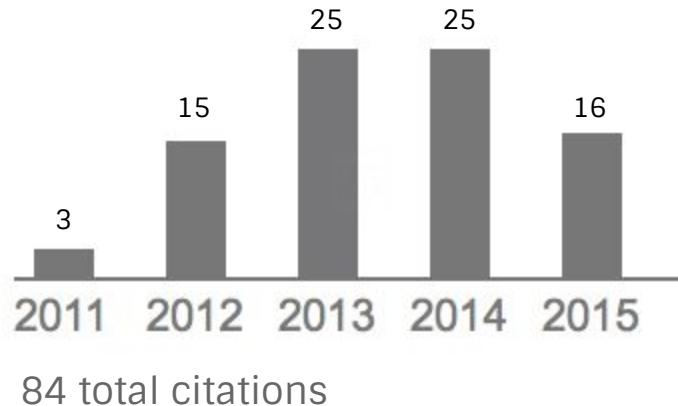
Perez-Llamas, C. & Lopez-Bigas, N. PLoS ONE 6, e19541 (2011)

DIANA miRPath v. 2.0: investigating the combinatorial effect of microRNAs in pathways.

A polycomb group protein is retained at specific sites on chromatin in mitosis.

Transcriptomic characterization of temperature stress responses in larval zebrafish.

温度刺激对斑马鱼仔鱼基因转录表达的影响.



● Users www.gitoools.org

1,500

750

Publication

January 2011

January 2012

January 2013

January 2014

January 2015





Data

A large orange rectangular block containing the word "Data". A white-outlined, orange puzzle piece shape is attached to its right edge, partially overlapping the next block.

Analysis

A light yellow rectangular block containing the word "Analysis". A white-outlined, yellow puzzle piece shape is attached to its right edge, partially overlapping the next block.

Visualization

An orange rectangular block containing the word "Visualization". A white-outlined, orange puzzle piece shape is attached to its left edge, partially overlapping the previous block.



Cancer
Data

Tumour
Alterations

Multiple
Views



Nuria López-Bigas

Abul Islam

Gunes Gundem

Alba Jene-Sanz

Jordi Deu-Pons

Abel Gonzalez-Perez

Carlota Rubio-Perez

David Tamborero

Simon J Furney

Michael P Schroeder

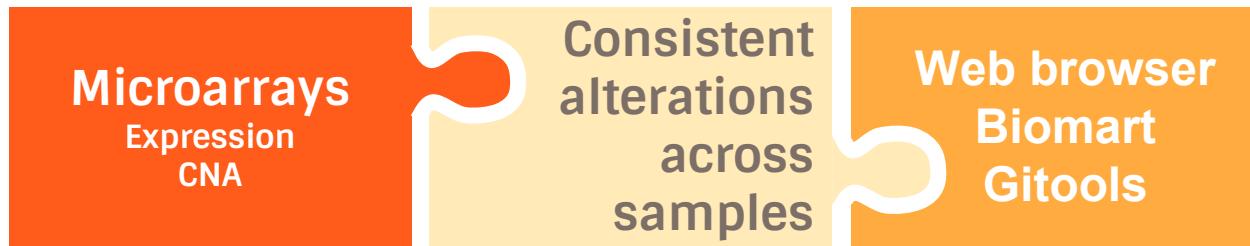
Alberto Santos

Anna Kedzierska

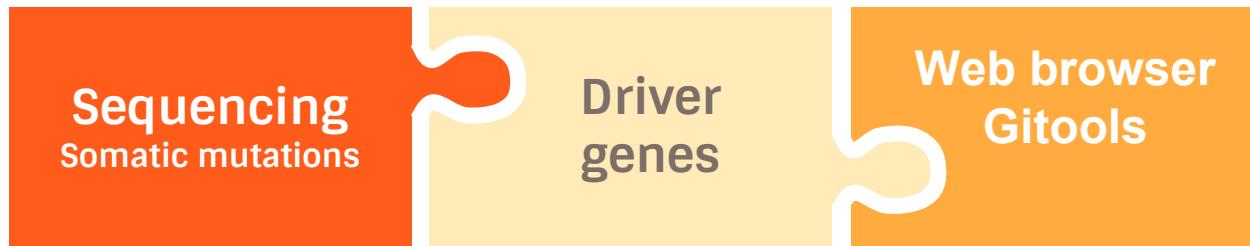
Christian Perez-Llamas



Arrays

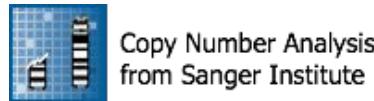


Mutations





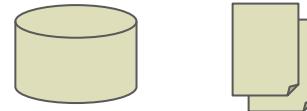
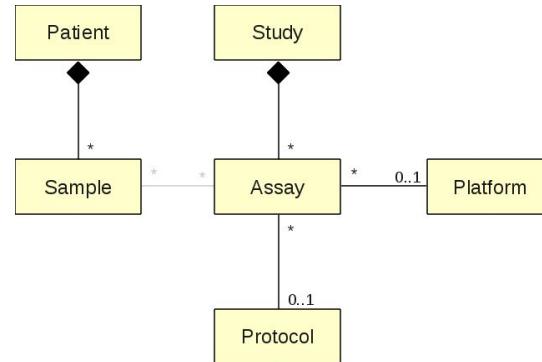
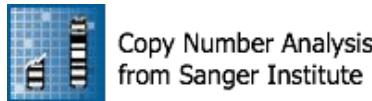
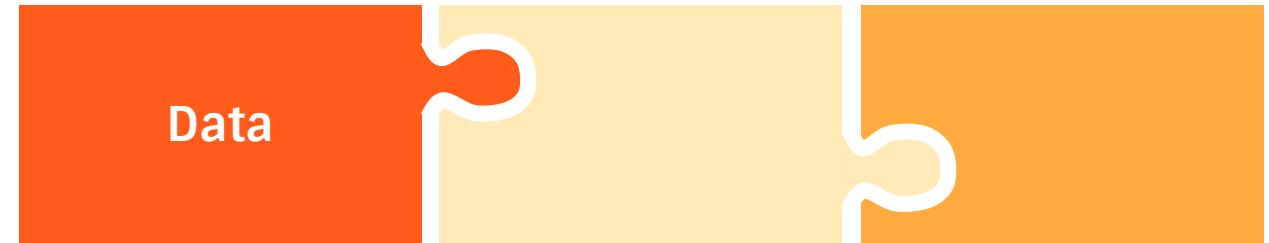
Data



Collect

Organize

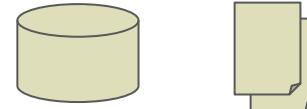
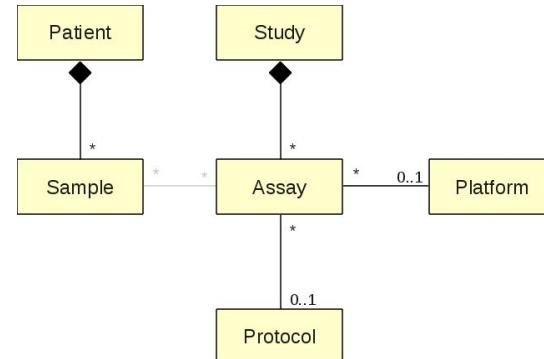
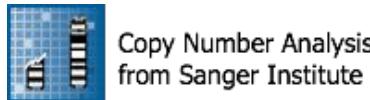
Annotate



Collect

Organize

Annotate



Collect

Organize

Annotate

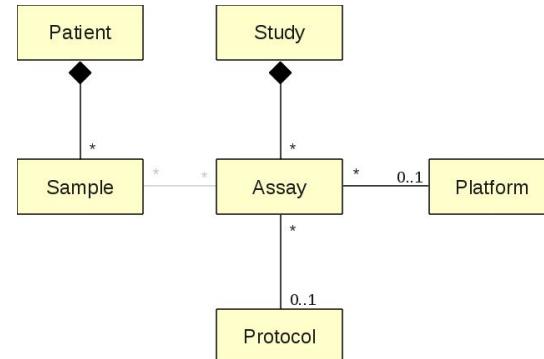


Arrays

Data



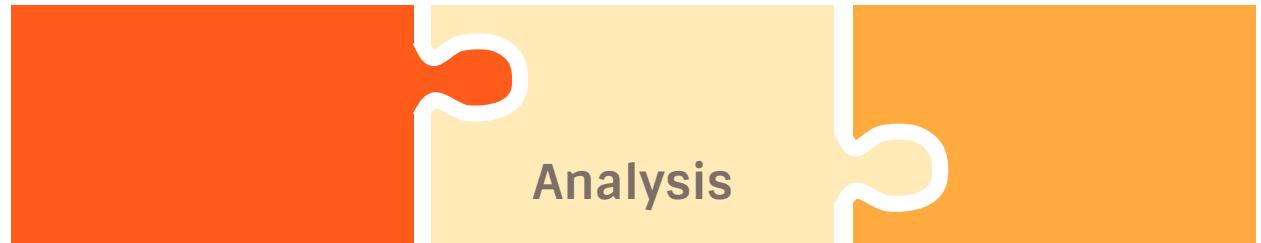
Gene Expression Omnibus



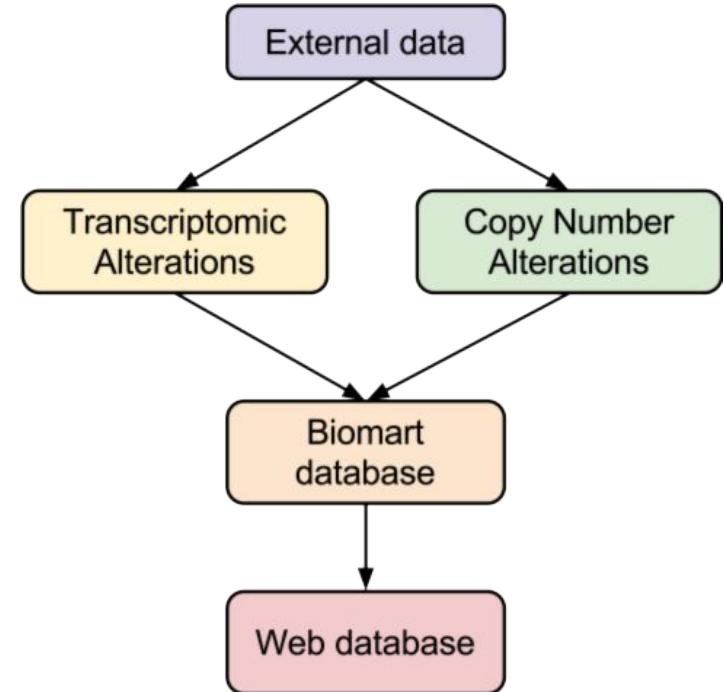
Collect

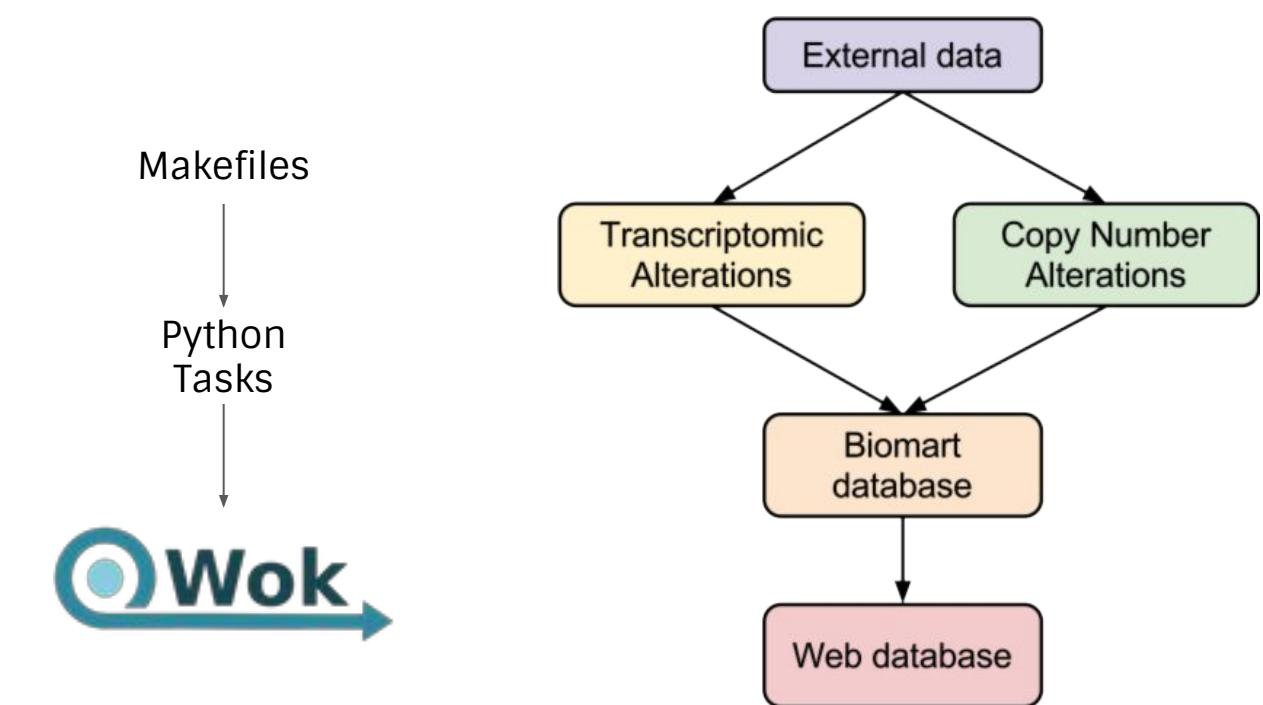
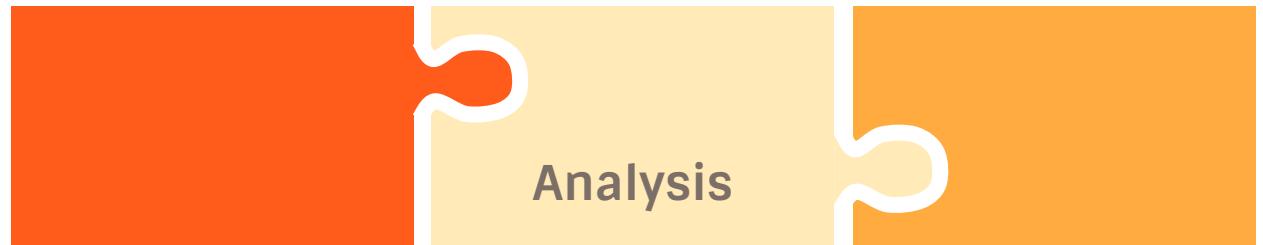
Organize

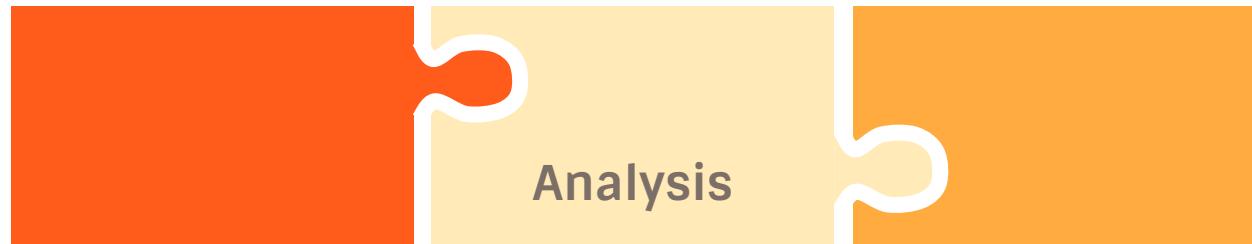
Annotate



Analysis







BPEL



Makefiles

Python
Tasks



Analysis

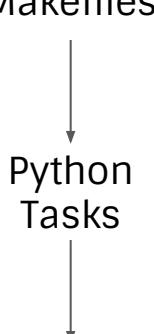
External data

Transcriptomic
Alterations

Copy Number
Alterations

Biomart
database

Web database



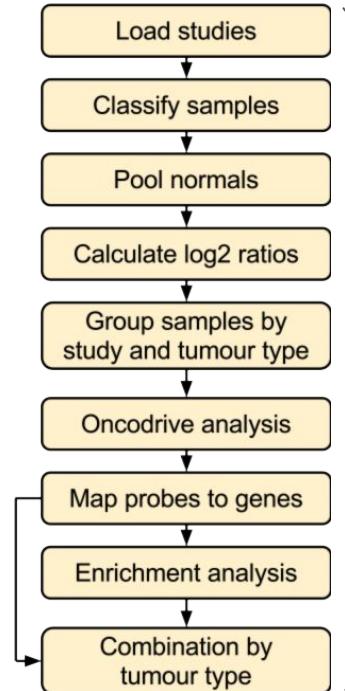
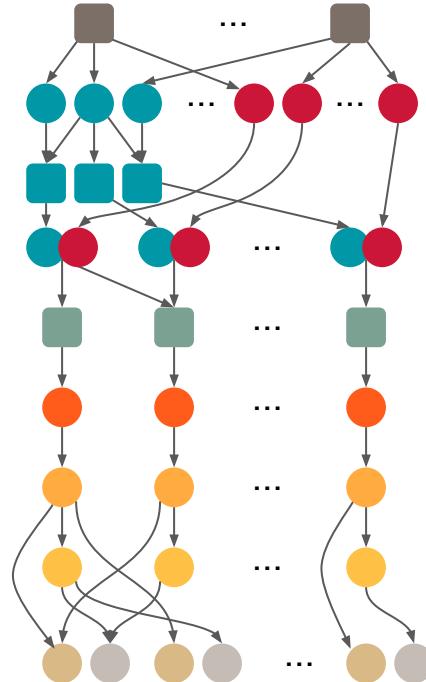
Wok

Python

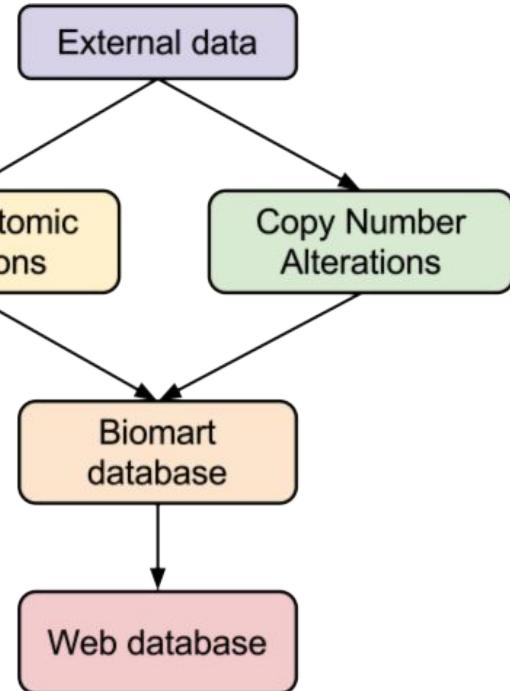
R

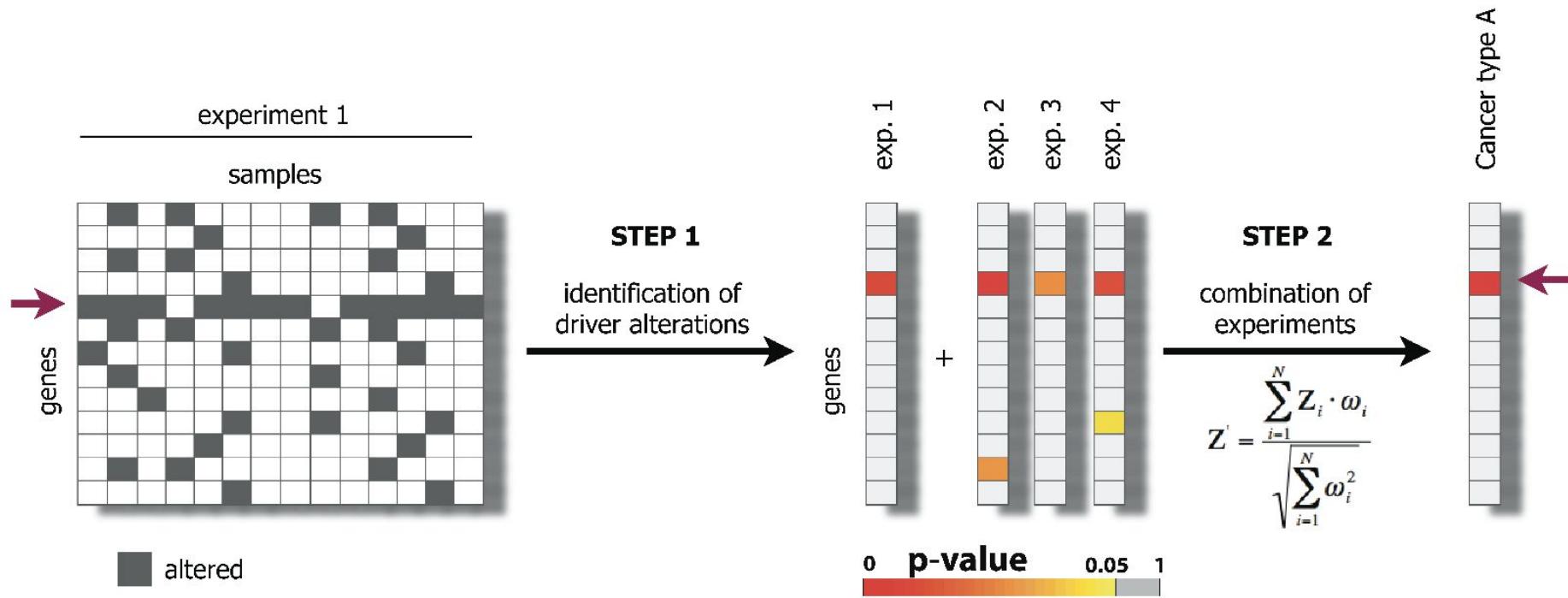
gitools

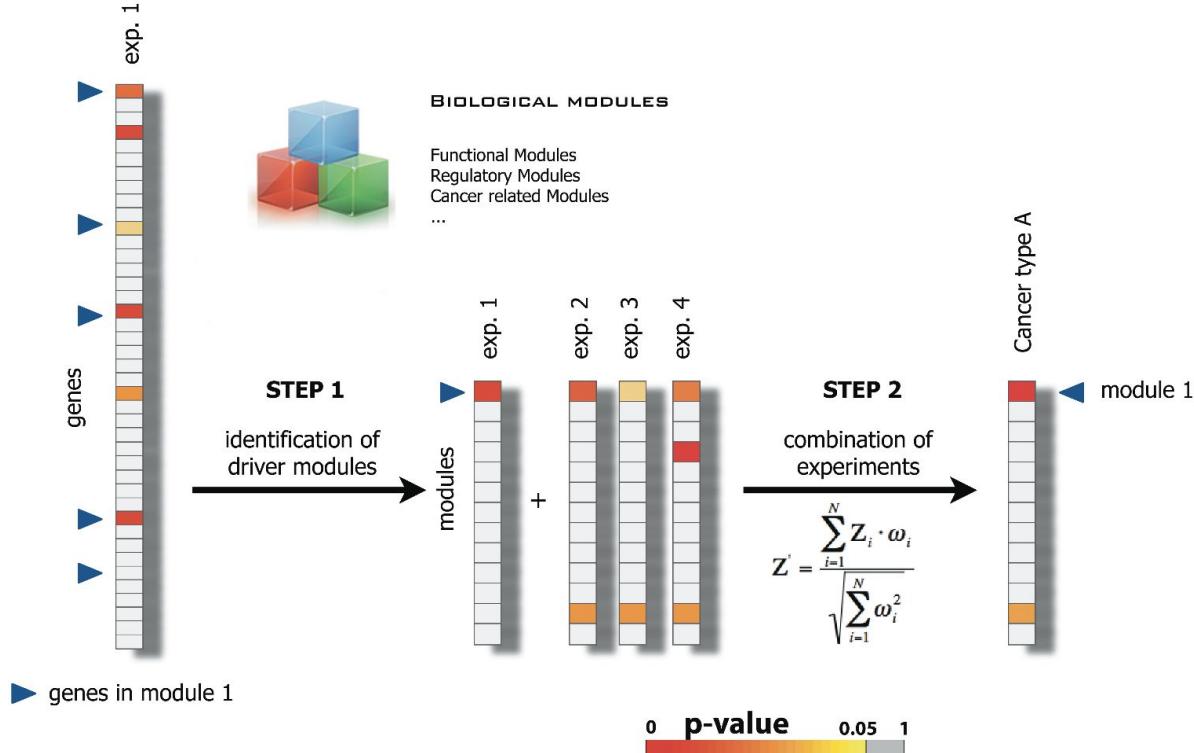
GRID ENGINE



Analysis





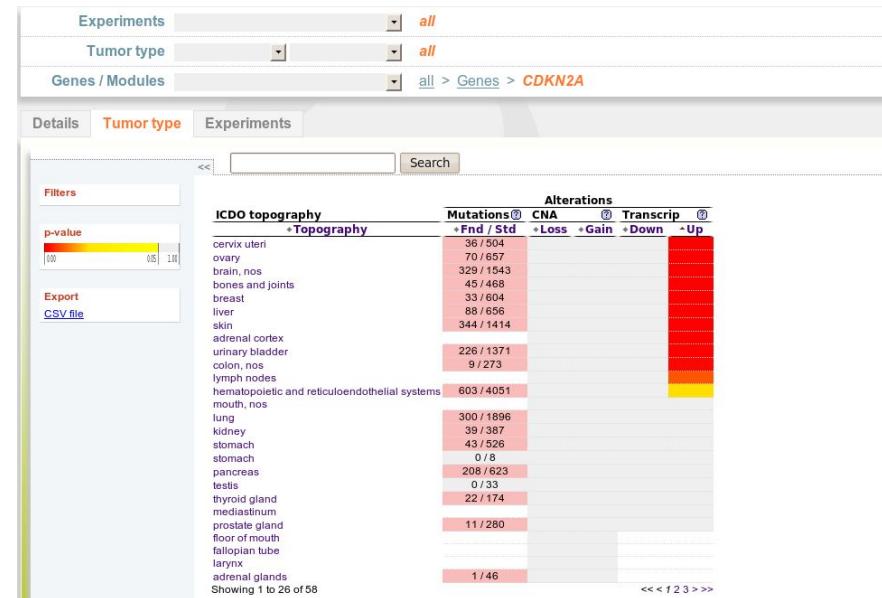
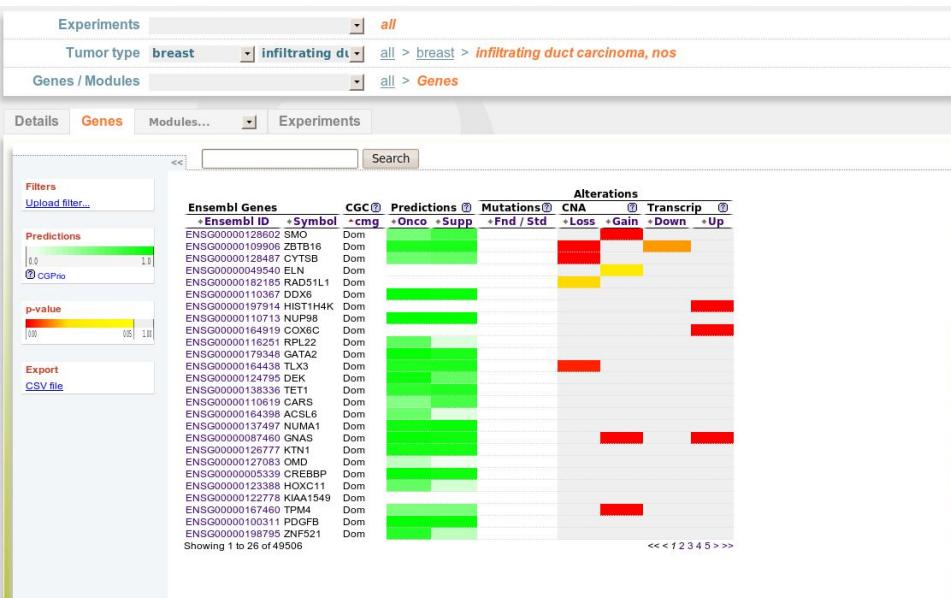




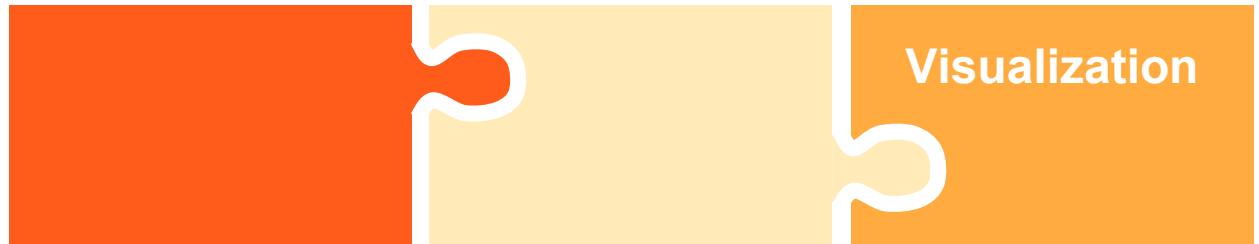
Arrays

intogen

Integrative Onco Genomics

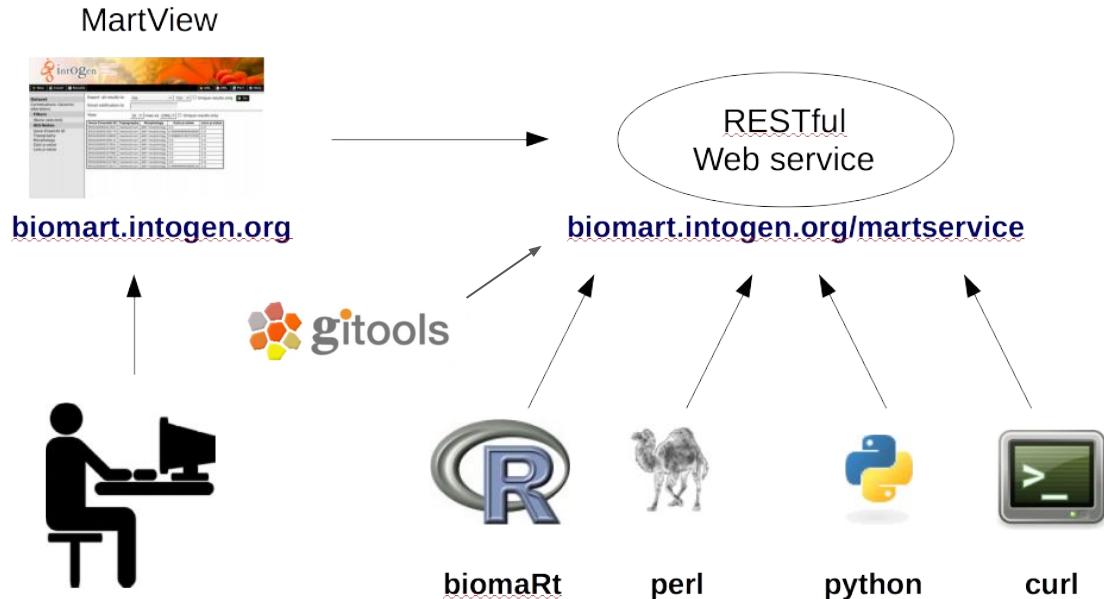


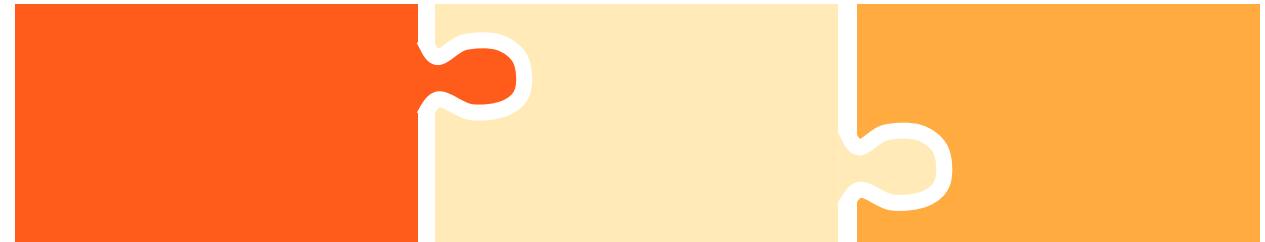
Web browser by Jordi Deu-Pons



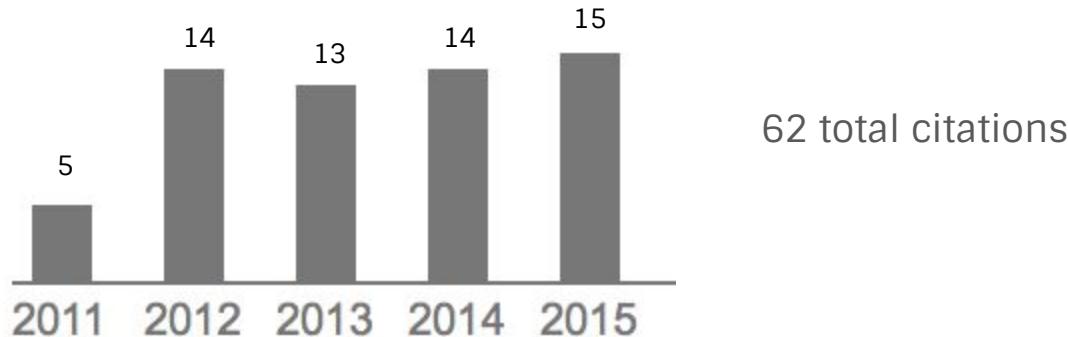
Perez-Llamas, C., Gundem, G., & Lopez-Bigas, N. (2011).
Integrative cancer genomics (IntOGen) in Biomart Database
The Journal of Biological Databases and Curation, 2011,
bar039. <http://doi.org/10.1093/database/bar039>

Guberman, J. M., Ai, J., Arnaiz, O., Baran, J., Blake, A.,
Baldock, R., ... Kasprzyk, A. (2011).
BioMart Central Portal: An open database network for the biological community.
Database, 2011.





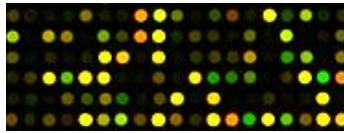
Gundem, G., Perez-Llamas, C., Jene-Sanz, A., Kedzierska, A., Islam, A., Deu-Pons, J., ... Lopez-Bigas, N. (2010).
IntOGen: integration and data mining of multidimensional oncogenomic data.
Nature Methods. <http://doi.org/10.1038/nmeth0210-92>





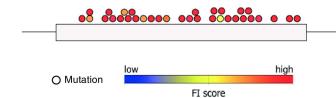
Sequencing
Somatic mutations

Driver
genes

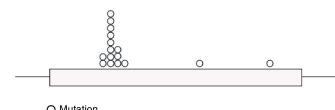


International
Cancer Genome
Consortium

Oncodrive**FM**



Oncodrive**CLUST**





Data

Independent
projects



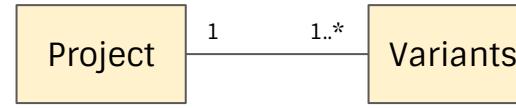
Collect

Organize



Data

Independent
projects

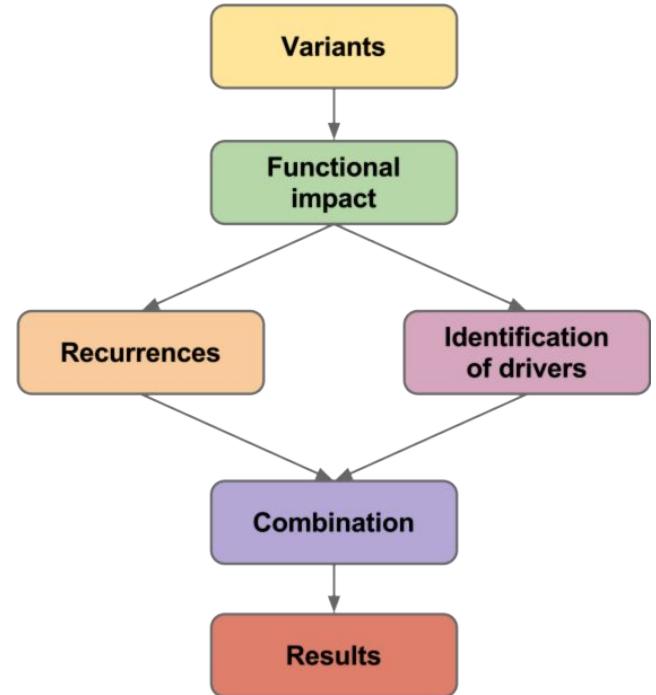
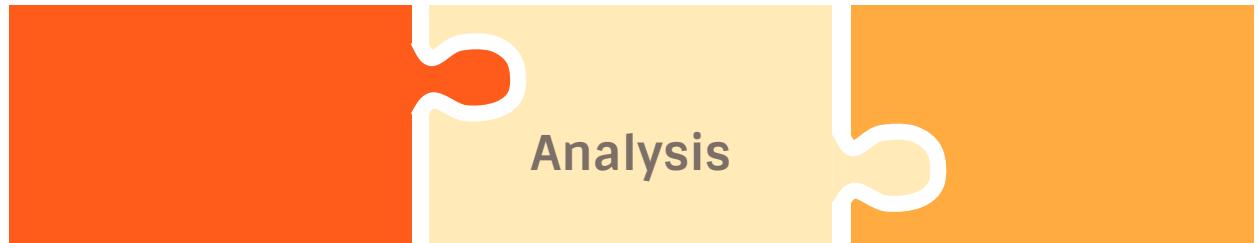


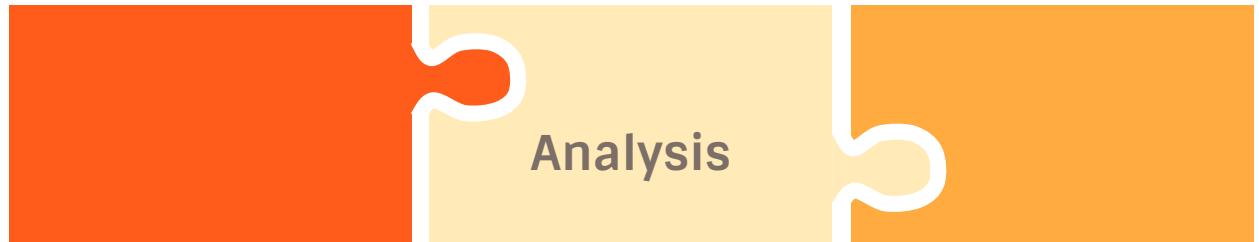
JSON

TSV
VCF
MAF

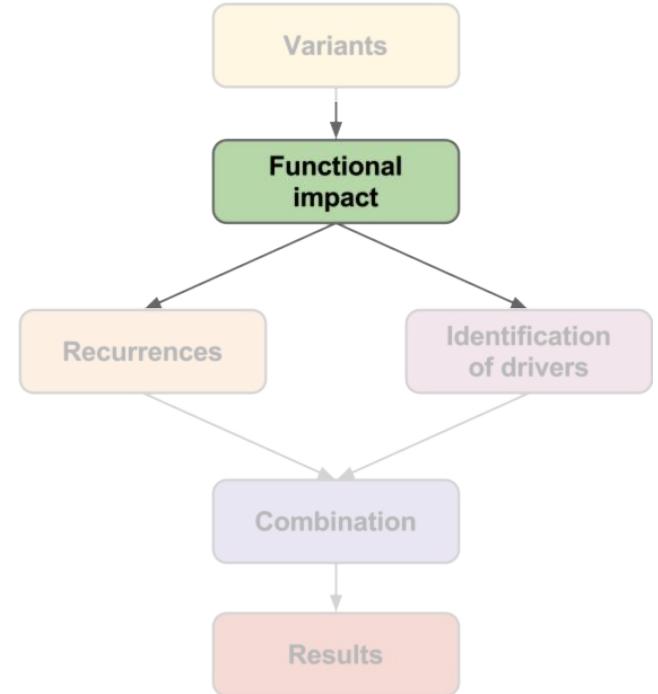
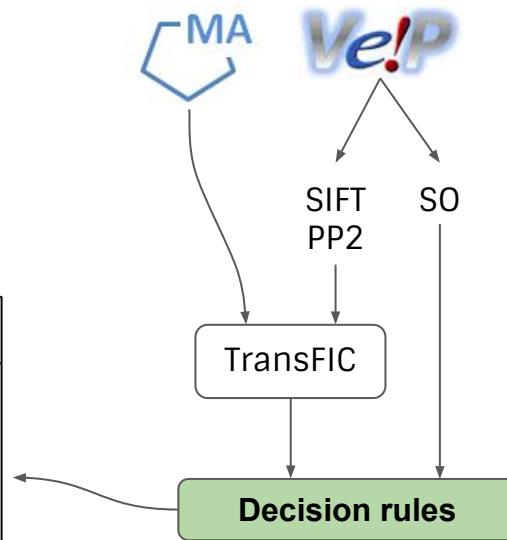
Collect

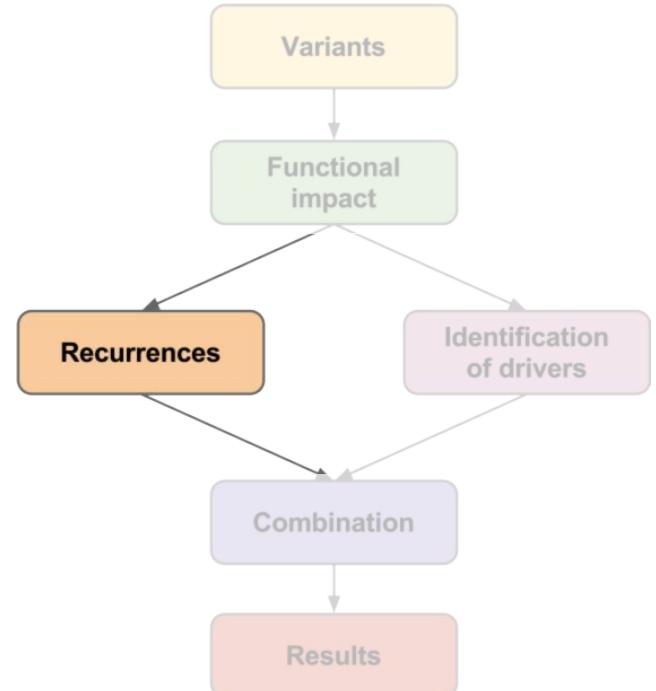
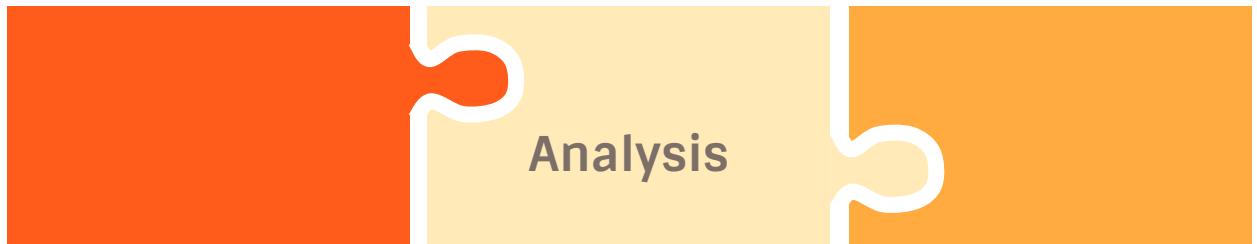
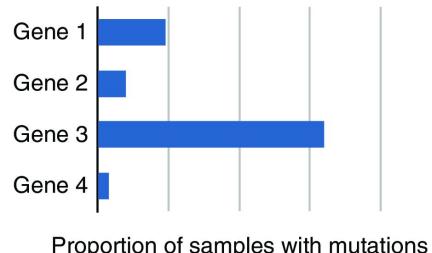
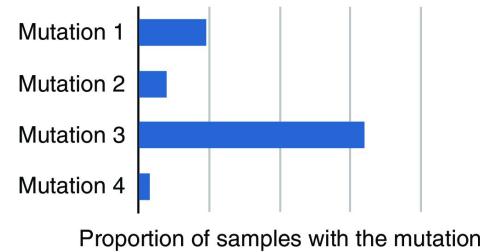
Organize

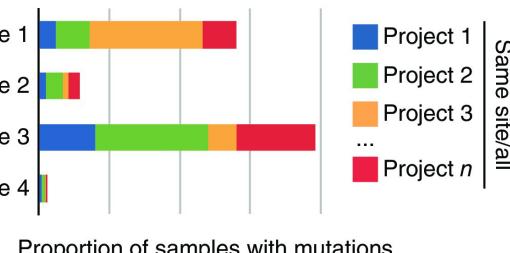
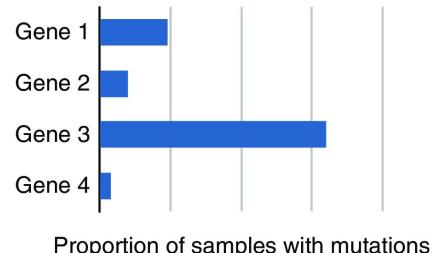
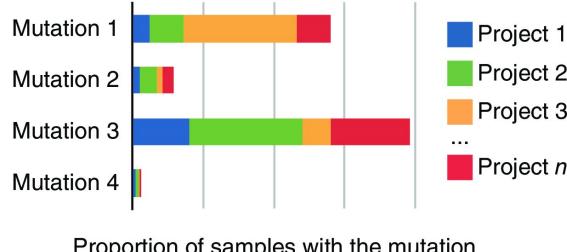
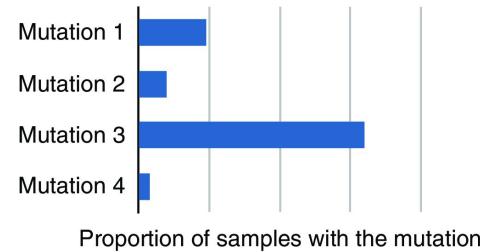




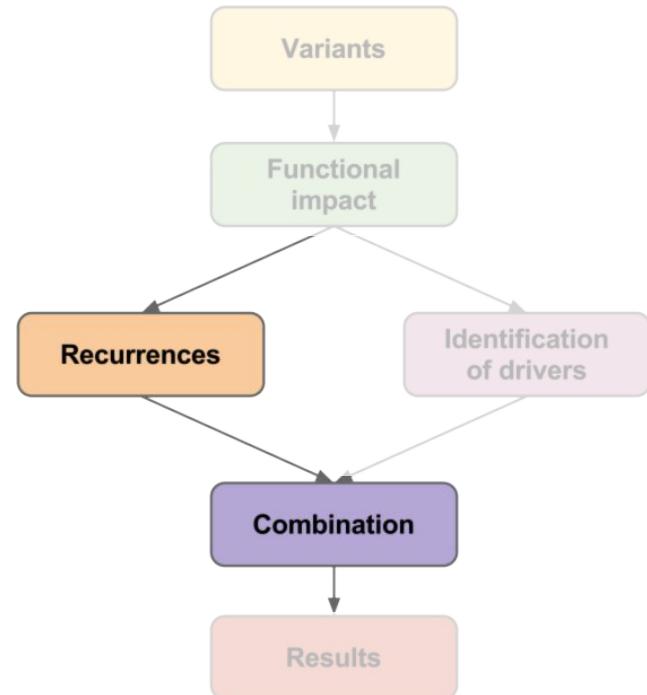
Mutation	Gene	Consequence	Functional impact
Mutation 1	Gene 1	Missense	Medium
Mutation 2	Gene 2	Synonymous	None
Mutation 3	Gene 2	Missense	Low
Mutation 4	Gene 3	STOP gain	High
Mutation 5	Gene 3	Missense	High
Mutation 6	Gene 3	Frameshift	High
Mutation 7	Gene 4	Synonymous	None
...			
Mutation n	Gene n	Missense	High

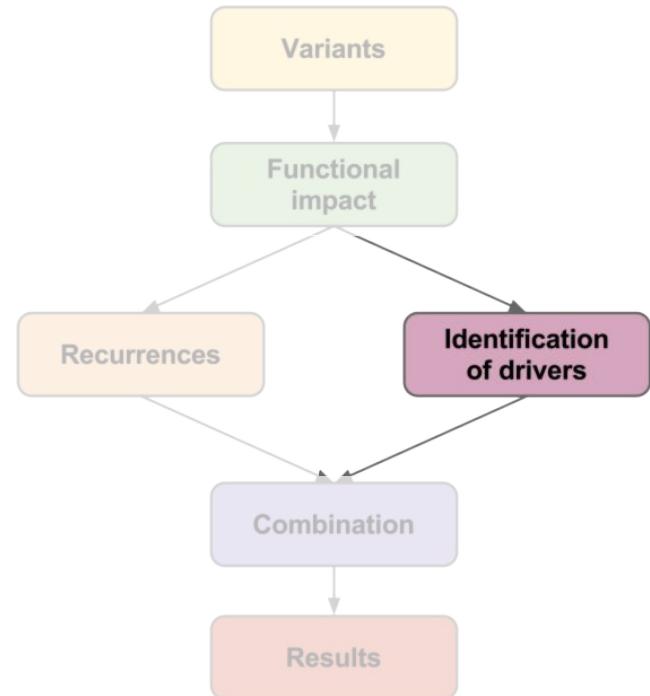
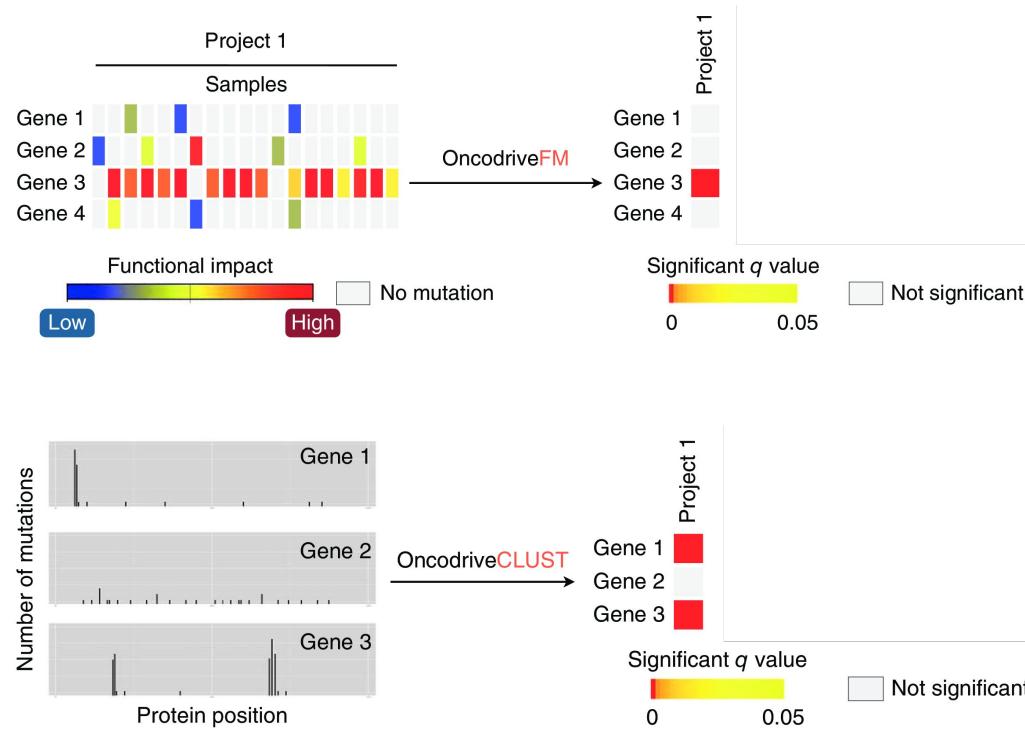
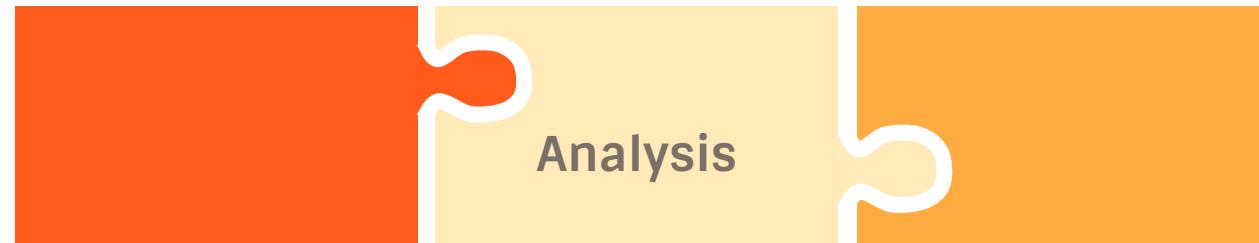


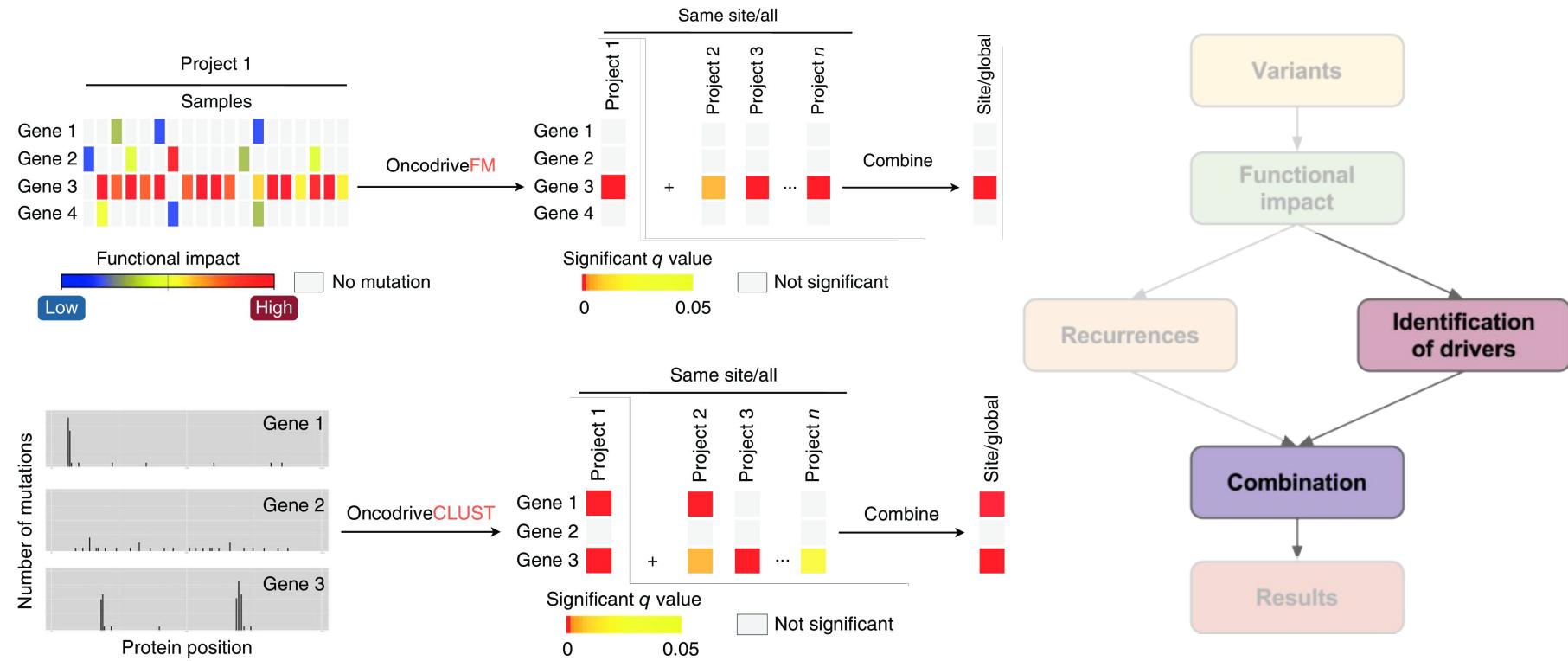


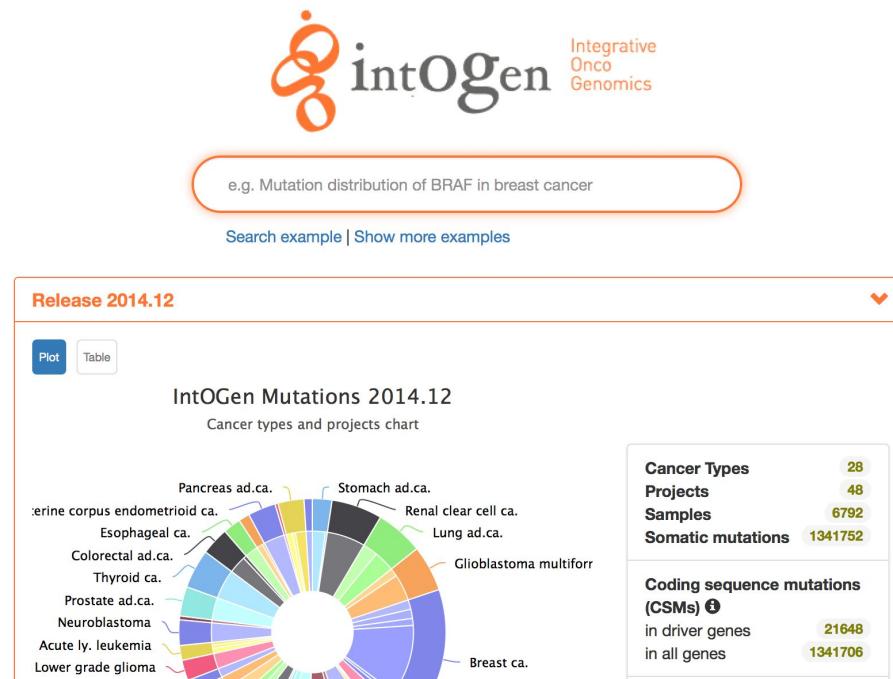


Analysis









Q Search Downloads Analysis About

Analysis

e.g. Mutation distribution of BRAF in breast cancer

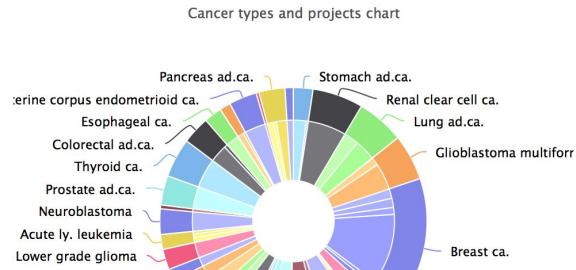
Search example | Show more examples

Release 2014.12

Plot Table

IntOGen Mutations 2014.12

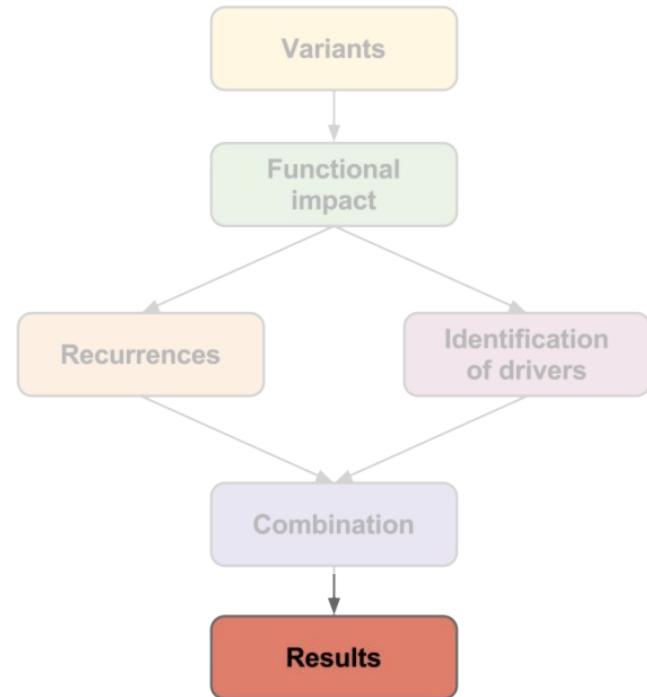
Cancer types and projects chart



Cancer Types	28
Projects	48
Samples	6792
Somatic mutations	1341752

Coding sequence mutations (CSMs) ⓘ	
in driver genes	21648
in all genes	1341706

Visualization



IntOGen Mutations Analysis

 Download

To interpret catalogs of cancer somatic mutations.

Cohort analysis



Use this if you have a list of somatic mutations for a cohort of tumors and want to identify driver mutations, genes and pathways.

 [View an example](#)

 [Analyse your data](#)

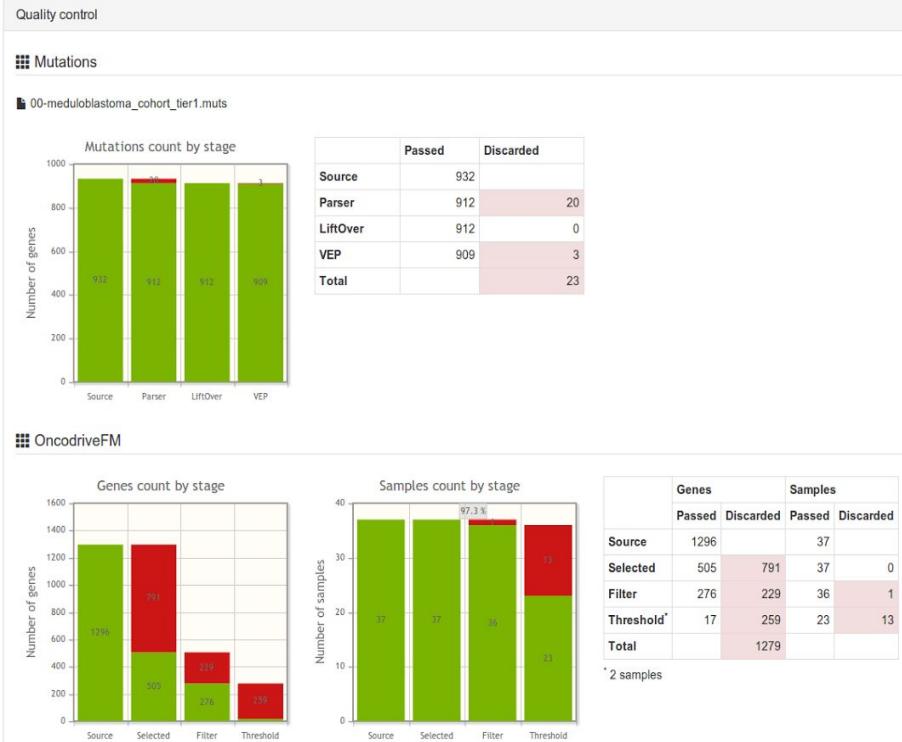
Single tumor analysis



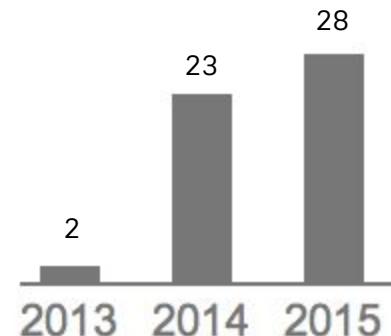
Use this if you have a list of somatic mutations for a single tumor and want to rank them based on their implication in cancer development.

 [View an example](#)

 [Analyse your data](#)

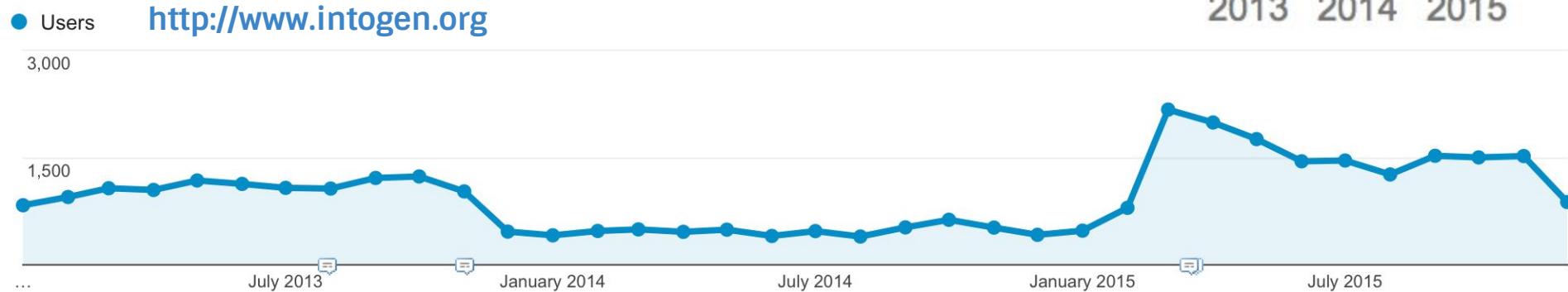


Gonzalez-Perez, A., Perez-Llamas, C., Deu-Pons, J., Tamborero, D., Schroeder, M. P., Jene-Sanz, A., ... Lopez-Bigas, N. (2013).
IntOGen-mutations identifies cancer drivers across tumor types.
Nature Methods, 10(11), 1081–2. <http://doi.org/10.1038/nmeth.2642>

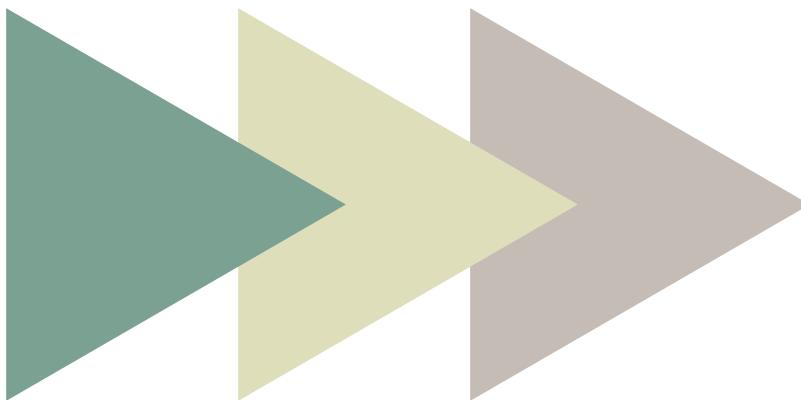


54 total citations

653 registered users (Oct 2015)



Benchmark of impact predictor tools



Prediction tools

General functional impact

SIFT, PolyPhen2, MutationTaster, FATHMM for diseases

Cancer functional impact

Mutation Assessor, FATHMM for cancer, CHASM, InCa

General conservation scores

GERP RS, PhyloP, C-Score

Combined / transformed scores

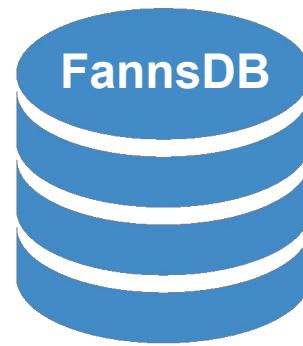
Condel, TransFIC

Precalculated scores



SIFT
PolyPhen2
MutationTaster
Mutation Assessor
GERP RS
PhyloP
FATHMM for disease

241 million possible protein variants with annotations for SwissProt, RefSeq, HGNC, OMIM, Ensembl



CHASM 3
InCa

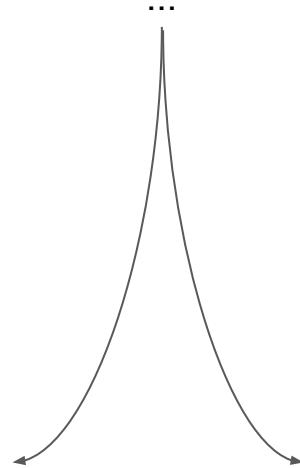


Condel

TransFIC

Proxy datasets

CHROM	POS	STRAND	REF	CHANGE	LABEL
1	10000	1	A	T	DRIVER
12	20000	0	C	G	PASSENGER



DRIVERS

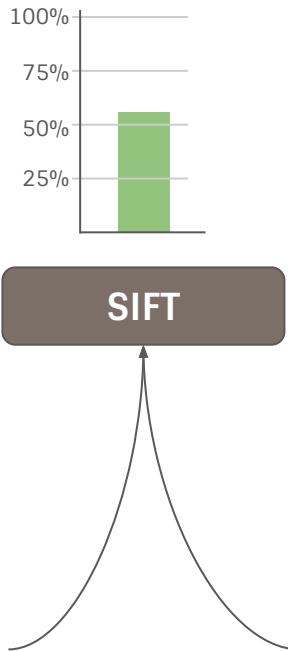


PASSENGERS



Proxy datasets

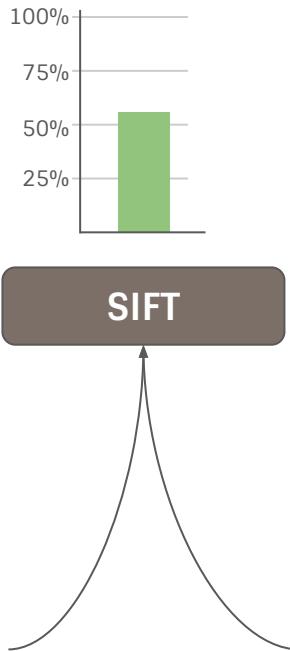
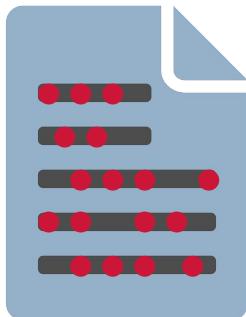
DRIVERS



PASSENGERS

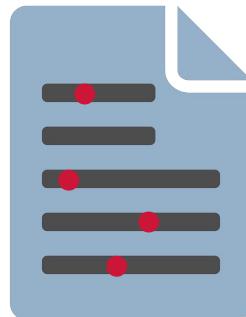
Proxy datasets

ENRICHED
FOR DRIVERS

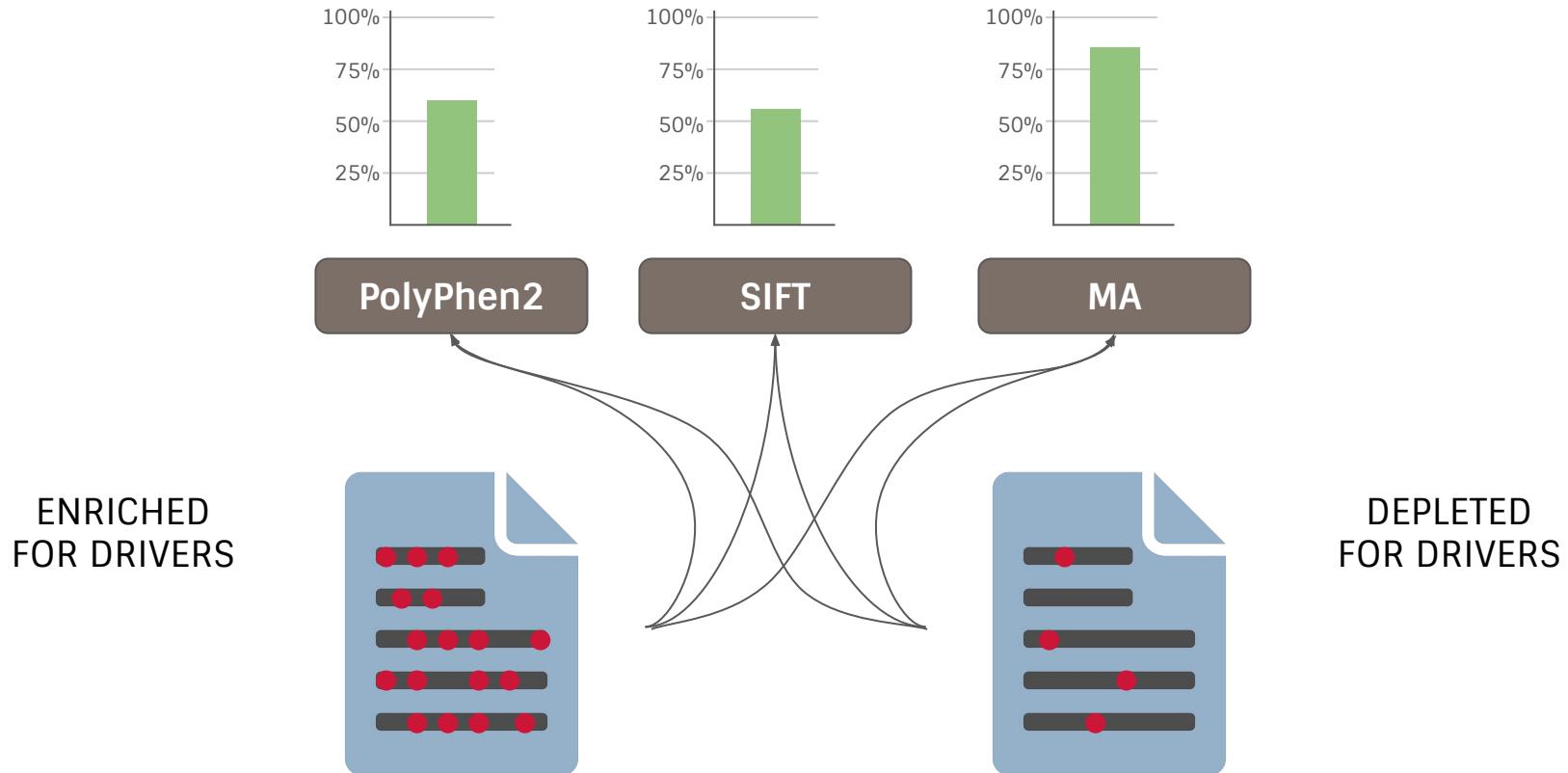


SIFT

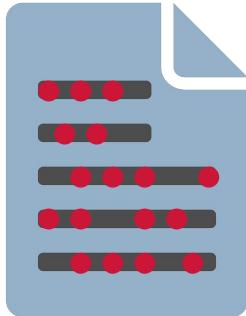
DEPLETED
FOR DRIVERS



Proxy datasets

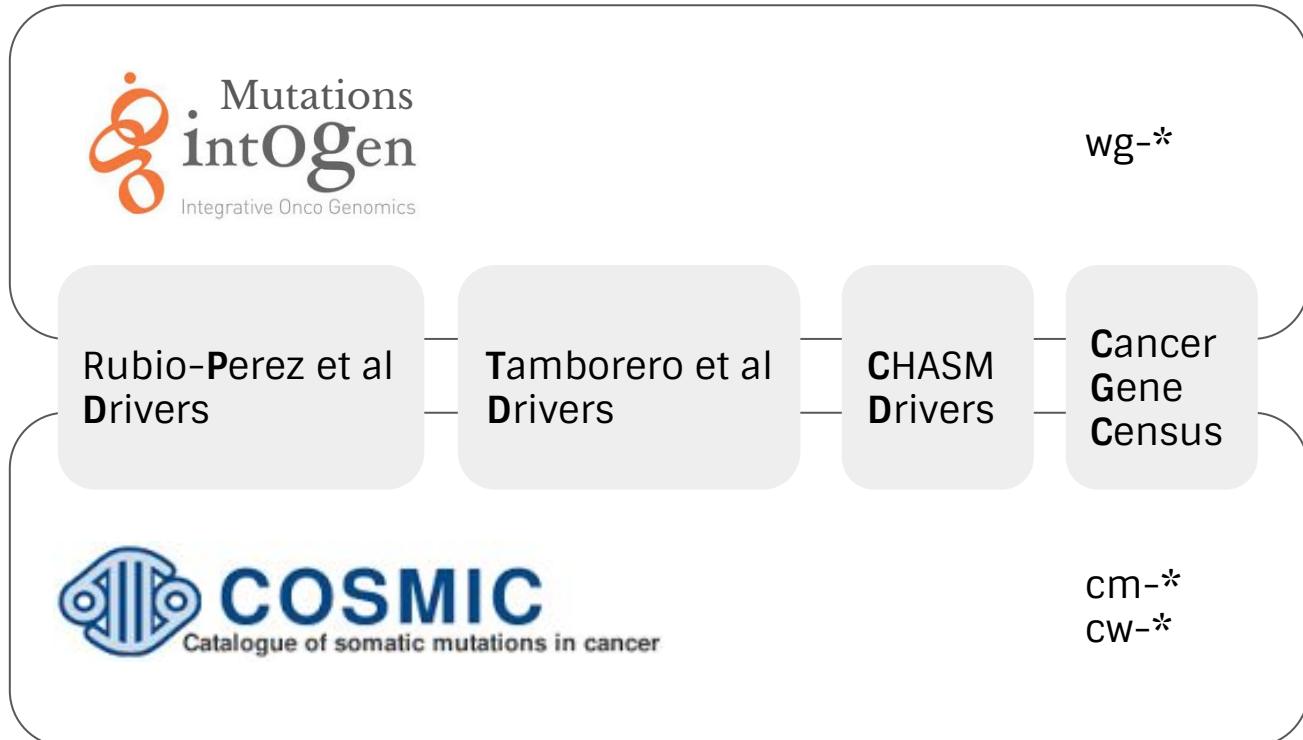


Proxy datasets



ENRICHED FOR
DRIVERS

Humvar



Evaluation of performance

		Proxy Dataset	
		Enriched	Reference
Predictor	Func Impact Driver	TP	FP
	Non Func Impact Non Driver	FN	TN

Evaluation of performance

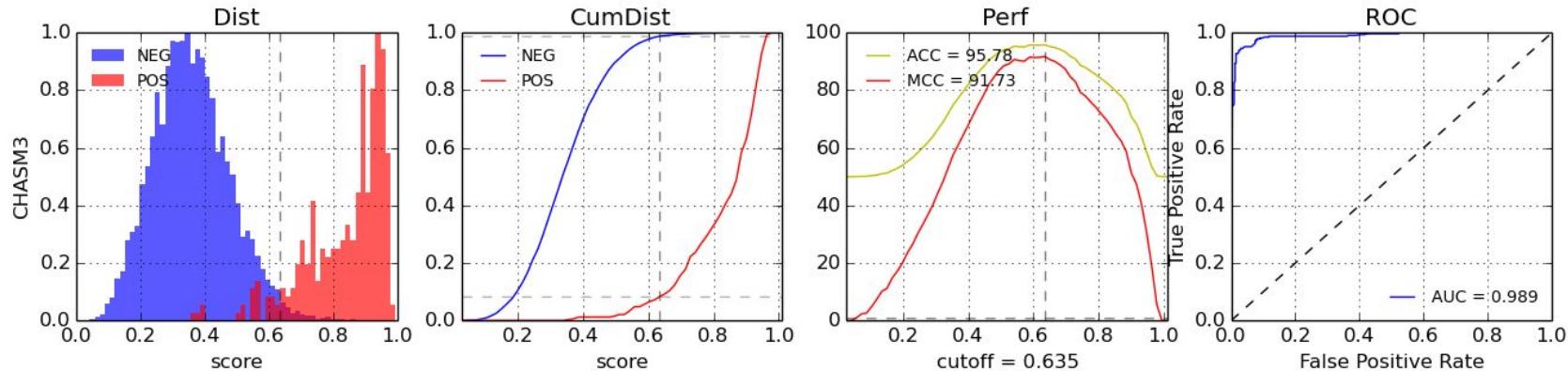
		Proxy Dataset	
		Enriched	Reference
Predictor	Func Impact Driver	TP	FP
	Non Func Impact Non Driver	FN	TN

Accuracy (ACC)	Sensitivity	Specificity
$\frac{\sum \text{TP} + \sum \text{TN}}{\text{TOTAL}}$	$\frac{\sum \text{TP}}{\sum \text{Enriched}}$	$\frac{\sum \text{TN}}{\sum \text{Enriched}}$

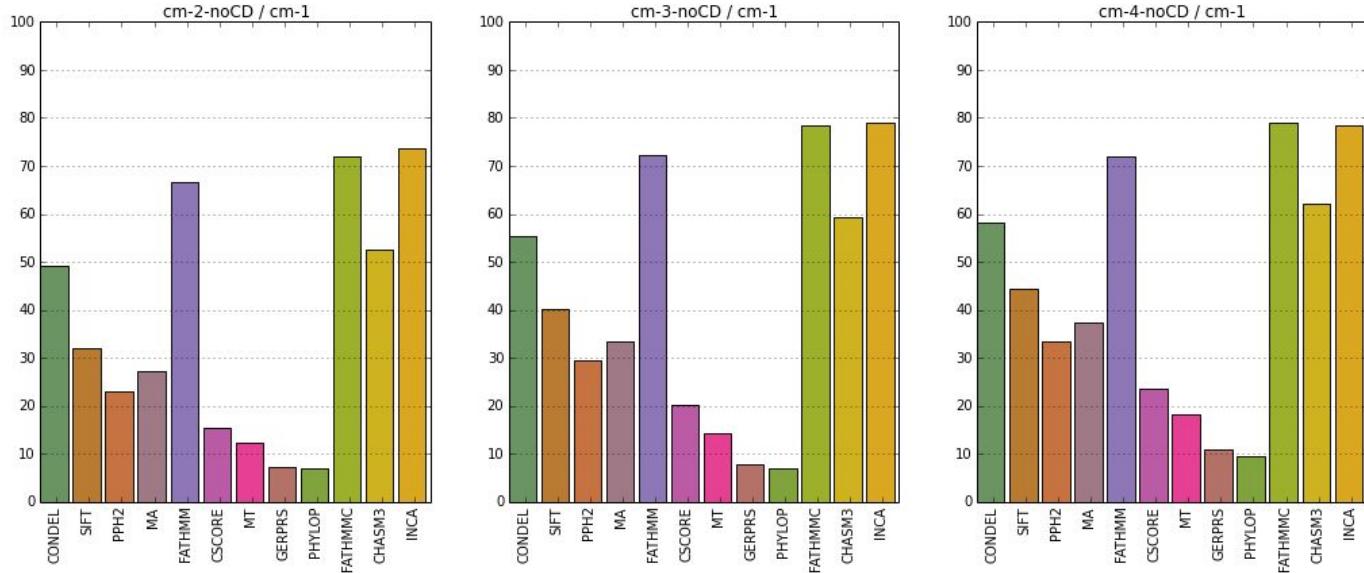
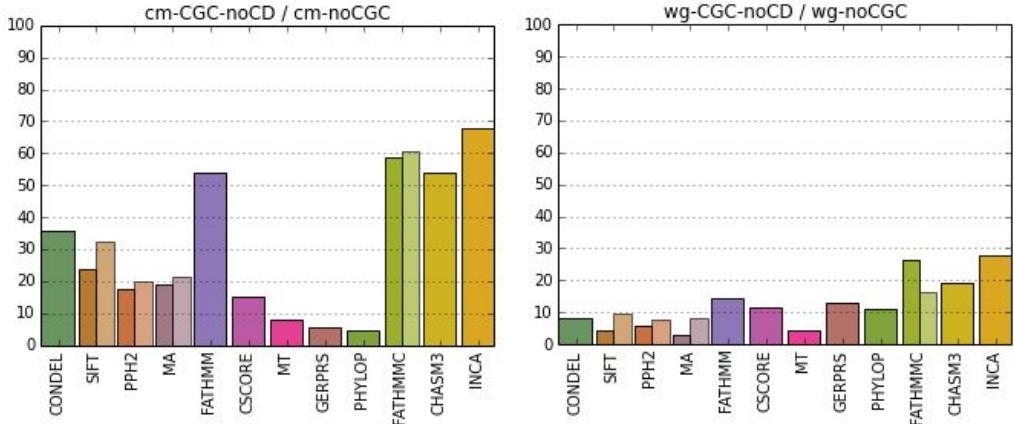
$$\text{MCC} = \frac{TP \times TN - FP \times FN}{\sqrt{(TP + FP)(TP + FN)(TN + FP)(TN + FN)}}$$

Evaluation of performance

		Proxy Dataset	
		Enriched	Reference
Predictor	Func Impact Driver	TP	FP
	Non Func Impact Non Driver	FN	TN

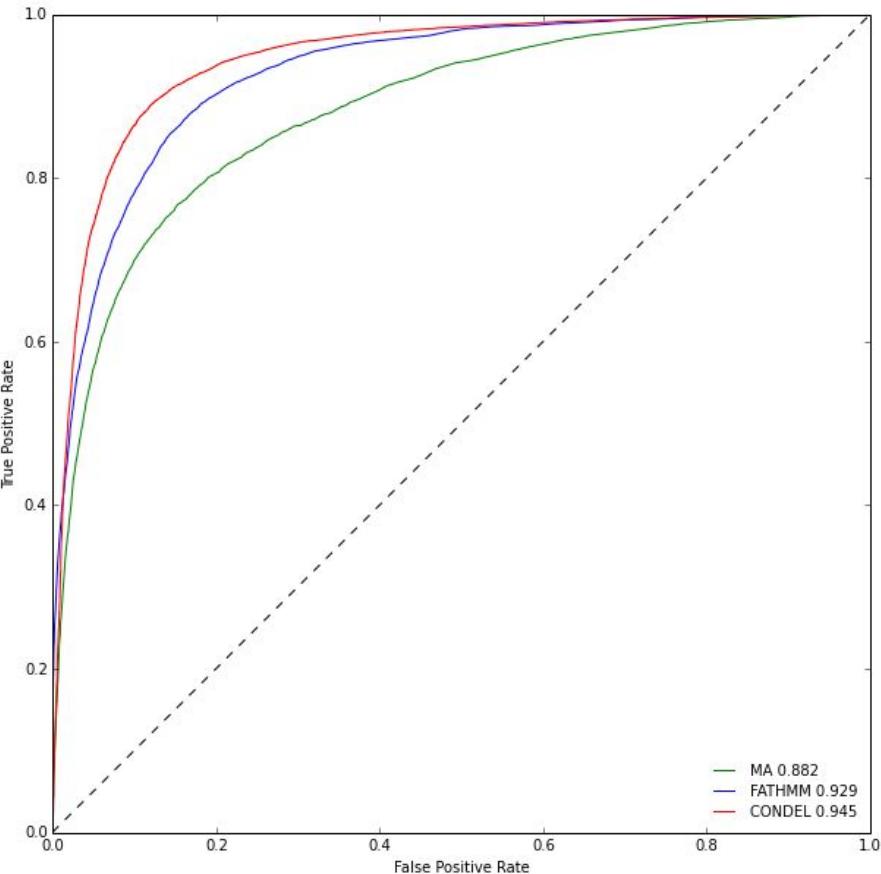


Benchmark results

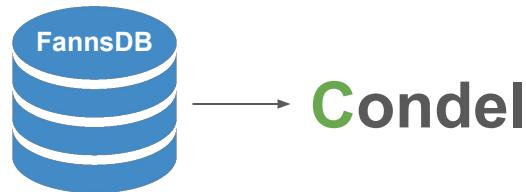


Condel

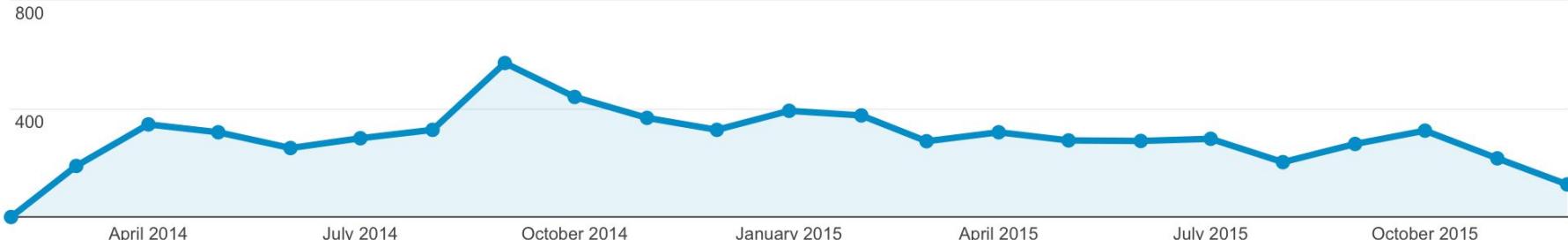
González-Pérez, A., & López-Bigas, N. (2011).
**Improving the Assessment of the Outcome of
Nonsynonymous SNVs with a Consensus
Deleteriousness Score, Condel.**
The American Journal of Human Genetics, 88(4),
440–449. <http://doi.org/10.1016/j.ajhg.2011.03.004>



241 million possible protein variants with annotations for SwissProt, RefSeq, HGNC, OMIM, Ensembl



- Users <http://bg.upf.edu/fannsdb/>



FannsDB

FannsDB is a database for Functional ANnotations for Non Synonymous SNVs which contains precalculated scores for several predictors.

[Documentation](#)[Download](#)

Condel

Condel is a method to assess the outcome of non-synonymous SNVs using a CONsensus DEleteriousness score that combines various tools ([MutationAssessor](#), [FATHMM](#)). This is the second version of Condel which includes an update of the combined tools and a new web

Conclusions

- I created **Gitools** for accessing and analysing biological data, as well as for visualizing multi-dimensional results with **interactive-heatmaps**.
- We collected and **integrated cancer data** from several repositories and large scale projects under the framework of **IntOGen**. The use of simple models and controlled vocabularies was of vital importance for the integration.
- I implemented the **IntOGen workflows** using state of the art tools and methodologies to detect consistently altered genes and identify candidate genes that drive the tumorigenesis for thousand of samples.
- **IntOGen results were easily available** to researchers in different ways. The IntOGen Mutations methods were also available for researchers to analyse their data.
- I integrated scores from several functional impact, cancer driver, and conservation predictors in a **database**, and evaluated their **performance**, and **updated Condel**. Results suggesting that **they alone are not enough for finding drivers** and need to be combined with other methodologies.

Questions

Gitools Case Study

