Developing Original Cloudbased Bioinformatics Software Applications

Daniel Clarke, Alexander Lachmann, Avi Ma'ayan



Department of Pharmacological Sciences Mount Sinai Center for Bioinformatics

Icahn School of Medicine at Mount Sinai New York, NY USA



Icahn Cen School of Bic Medicine at Mount Sinai





BD2K-LINCS DATA COORDINATION AND INTEGRATION CENTER









TRENDS in Pharmacological Sciences

Ma'ayan et al. Trends Pharmacol Sci. 2014 Sep;35(9):450-60.





https://maayanlab.cloud/Enrichr/





https://maayanlab.cloud/Enrichr/









Search for genes or proteins and their functional terms extracted and organized from over a hundred publicly available resources. Learn more.



Example searches achilles STAT3 breast cancer

Ma'ayan Laboratory of Computational Systems Biology | Contact Us | Terms



Please acknowledge the Harmonizome in your publications by citing the following reference:

Rouillard AD, Gundersen GW, Fernandez NF, Wang Z, Monteiro CD, McDermott MG, Ma'ayan A. The harmonizome: a collection of processed datasets gathered to serve and mine knowledge about genes and proteins. Database (Oxford). 2016 Jul 3;2016. pii: baw100.

Funding

BD2K-LINCS Data Coordination and Integration Center Illuminating the Druggable Genome, Knowledge Management Center

Functional Associations

STAT3 has 14,374 functional associations with biological entities spanning 8 categories (molecular profile, organism, functional term, phrase or reference, disease, phenotype or trait, chemical, structural feature, cell line, cell type or tissue, gene, protein or microRNA) extracted from 100 datasets.

Click the + buttons to view associations for STAT3 from the datasets below.

If available, associations are ranked by standardized value 🚱

	Dataset	Summary
+	Achilles Cell Line Gene Essentiality Profiles	Cell lines with fitness changed by STAT3 gene knockdown relative to other cell lines from the Achilles Cell Line Gene Essentiality Profiles dataset.
+	Allen Brain Atlas Adult Human Brain Tissue Gene Expression Profiles	Tissues with high or low expression of STAT3 gene relative to other tissues from the Allen Brain Atlas Adult Human Brain Tissue Gene Expression Profiles dataset.
+	Allen Brain Atlas Adult Mouse Brain Tissue Gene Expression Profiles	Tissues with high or low expression of STAT3 gene relative to other tissues from the Allen Brain Atlas Adult Mouse Brain Tissue Gene Expression Profiles dataset.
+	Allen Brain Atlas Developing Human Brain Tissue Gene Expression Profiles by Microarray	Tissue samples with high or low expression of STAT3 gene relative to other tissue samples from the Allen Brain Atlas Developing Human Brain Tissue Gene Expression Profiles by Microarray dataset.
+	Allen Brain Atlas Developing Human Brain Tissue Gene Expression Profiles by RNA- seq	Tissue samples with high or low expression of STAT3 gene relative to other tissue samples from the Allen Brain Atlas Developing Human Brain Tissue Gene Expression Profiles by RNA-seq dataset.
+	Allen Brain Atlas Prenatal Human Brain Tissue Gene Expression Profiles	Tissues with high or low expression of STAT3 gene relative to other tissues from the Allen Brain Atlas Prenatal Human Brain Tissue Gene Expression Profiles dataset.
+	Biocarta Pathways	Pathways involving STAT3 protein from the Biocarta Pathways dataset.
+	BioGPS Cell Line Gene Expression Profiles	Cell lines with high or low expression of STAT3 gene relative to other cell lines from the BioGPS Cell Line Gene Expression Profiles dataset.
+	BioGPS Human Cell Type and Tissue Gene Expression Profiles	Cell types and tissues with high or low expression of STAT3 gene relative to other cell types and tissues from the BioGPS Human Cell Type and Tissue Gene Expression Profiles dataset.
+	BioGPS Mouse Cell Type and Tissue Gene Expression Profiles	Cell types and tissues with high or low expression of STAT3 gene relative to other cell types and tissues from the BioGPS Mouse Cell Type and Tissue Gene Expression Profiles dataset.
+	CCLE Cell Line Gene CNV Profiles	Cell lines with high or low copy number of STAT3 gene relative to other cell lines from the CCLE Cell Line Gene CNV Profiles dataset.
+	CCLE Cell Line Gene Expression Profiles	Cell lines with high or low expression of STAT3 gene relative to other cell lines from the CCLE Cell Line Gene Expression Profiles dataset.



Impact of the Ma'ayan Lab Tools

>41 published bioinformatics tools and databases
>1.2 million unique users across all tools
>3,000 unique users per day
>30,000 unique users per month
>3,000 papers that cite the tools in 2020

Ψ	L1000CDS2 BD2K-LINCS DCIC An ultra-tast LINCS L1000 Characteristic Direction signature search engine	Drug-Pathway Bro Has LINS Interactive may of key sign transduction pathways are drug-target data	ere LUNCS BOXELINCS DCIC An integrative web platform for available of the standard signatures O	Drug Gene Budger BDX-LINCS DCIC Identifies drugs & small molecules to regulate expression of target grees	Omics Integrator Domics Integrator Omics Hold Information Discours networks Information regenomic data	GEO2Enrichr BD2KLINGS DCIC A web app and browser signatures from GEO	RTK Profile Browser HMSINCS Online tool for browing breast career cell line RTK profile
	Repurposing App LINCS Transcriptomics Tool to septore repurposing collection of ~5000 tool compounds and drugs	Drug/Cell-line Brow BD2/CLNCS DCIC DCCP provide interactive viability data	ine • • • • • • • • • • • • • • • • • • •	Harmonizome BD2KLINCS DCIC Web portal with a collection of 14 disease shortracked into gene function tables	CREEDS BD2KLINCSDCIC Collections of processed gene, dy and disease signatures from GEO	TEDIE HNSUNCS Software for computing required transcript data With the software for the	LINCS Project Mobile BD2K LINCS DCIC A mobile application to explore LINCS centers and resources
	Query App LINCS Transcriptomics Connections to user-defined signatures	LINCS Data Portal BO2K UNCS DCIC Provide unified interface secting all UNCS data scidages and entities	CLUE Platform LINCS Transcriptonics Computational environment to interface with the L1000 data	Network2Canvas BO2K-LINCS DCIC Network visualization on a canvas with enrichment analysis	SEP L1000 BD24-LNCSDCIC Web portal for searching and browsing predictive small- molecule/ADR connections	Herast Cancer Browser HISUNCS Order tool for borowing multiple datasets relevant to breast cancer biology	Clarion MEPUNCS Million of microenvironment perturbations at a click
*	L1000FWD BD2K-LNCS DCIC Large scale visualization of drug-induced transcriptomic signatures	Drug Response Bro HMS LINCS Online tool for froming to cancer cell line dug dose- response data	erset	HMS LINCS Database HMS LINCS Of the database HMS LINCS Center databases and reagent information	LINCS Canvas Browser BD24-UNCSDCIC LCD provides interaction cluster ing with enrichment of L1000 Signatures	Touchstone App LINCS Transcriptomics Callection of CMap reference signatures	UP-BCNB BD2KLINCSDCC LitiCs Joint Project-Breast Carter Network Brower is an interactive viol 2544 agrutures:
X2Kweb	eXpression2Kinases BD2K-LINCS DCIC Linking expression signatures to ugstream cell signaling networks	PAEA BD2KLINCS DCIC Briothemst analysis tool invidenmenting the principal angle method	Adaptation Browser HMSUNCS Colline tool for forwards metamona celline adaptive resistance mechanisms	Silcr BO2K-UNCS DCIC A metadata search engine that provides easy access to L1000 OE0 data	cynetworkbrowser BD2xLUNC5DCIC An interactive tool for generating and viewing protein interaction networks	SynergySeq B2AC-UNCS DOIC Underly syner gistic drug combinations	Datasets2Tools D2XLINCSOCC Aplatom for the discovery and evaluation of biomedical digital objects
	PIUMet NeuroLINCS Tool for discovering pathways of peaks from untargeted metabolomics data	Morpheus App UNCS Transcriptomics An interactive version of the ICV that fersy comanipulate and annotate an existing dataset	ICV App LINCS Transcriptomics A matrix-based interactive heatmup to optione relationships within the data	GEN3VA BD2KLINCSDCIC Agregate and analyses gene expression signitures extracted from GEO	Cell Dynamics Browser HMS LINCS Online tool for live-cell image data visualization	AchroMap NeuroLINCS A data integration tool for transcriptomic and epigenomic data	Panorama UNCS PCCSE Proteomics repository that provide access to the PCCSE ptoto and GCP data
	GF Response Browser	BioJupies	GUldock	GR Browser	MEPmosaic	slinky	PINET
	HMS LINCS Online tool for browsing growth factor-stimulated cell signaling profiles	BD2K-LINCS DCIC Automates RNA seq data analysis natebooks	BDZK-UNKS DCIC Coder package containing a computational environment to run apps with a COI	HMS LINCS BOZK-LINCS DOIC CALCULATOR	B2X-LINCS DCIC Mosaic visualization of high- cancer cell types	R package to query the L1000 metadata via the CLUE io REST API	BD2X-LINES DOCL Used to annotate, map and analyze a set of peptide moieties

https://lincsproject.org/LINCS/tools





Step 1. Upload or Fetch RNA-seq Data

Upload your raw or processed RNA-seq data
Fetch >8,000 public RNA-seq datasets published in the Gene Expression Omnibus

Step 2. Select Data Analysis Tools

- Select from multiple state-of-the-art RNA-seq data analysis tools
- Contribute your computational tool as a plugin

Step 3. Generate Your Notebook

 Access and share your results through a permanent URL
 Download, rerun and customize your notebook using Docker

• • • • • • • •

BioJupies Automatically Generates RNA-seq Data Analysis Notebooks

With BioJupies you can produce in seconds a customized, reusable, and interactive report from your own raw or processed RNA-seq data through a simple user interface

Get Started



To acknowledge Biolupies in your publications, please use the following reference: Torre, D., Lachmann, A., and Ma'ayan, A. (2018). Biolupies: Automated Generation of Interactive Notebooks for RNA-Seq Data Analysis in the Cloud. Cell Systems.

Θ



BioJupies now supports user accounts! Sign in from the top right corner of the page for access to unlimited private notebooks, RNA-seq datasets and alignment jobs.

0

FIRCHS⁴

Download Chrome Extension Help

Search





5



https://maayanlab.cloud/archs4/data.html Nat Commun. 2018 Apr 10;9(1):1366.

Cell Lines

Integumentary System Muscular System Nervous System Respiratory System Urogenital System

The L1000FWD Map - ~17K Signatures, ~5K Drugs



https://maayanlab.cloud/L1000FWD/main

Bioinformatics. 2018 Jun 15;34(12):2150-2152

What are the molecular effects of hydroxychloroquine on human cells?

OME SEARCH	SITE MAR	GEO Publications FAQ MIAME Email GEO
NCBI > GEO >	Acces	sion Display 🛛 Not logged in Login 🗹
GEO help: Mous	se over	screen elements for information.
Scope: Self	T	Format: HTML Amount: Quick GEO accession: GSE74235 GO GO
Series GSE	74235	Query DataSets for GSE74235
Status		Public on Oct 21, 2016
Title		Hydroxychloroquine inhibits responses to group A streptococcus in peripheral blood mononuclear cells
Organism		Homo sapiens
Experiment t	type	Expression profiling by high throughput sequencing
Summary		Immune responses to group A streptococcus in humans can lead to the development of acute rheumatic fever and rheumatic heart disease. Immune pathways that are activated by group A streptococcus are potential targets for inhibiting autoimmune responses to group A streptococcus. This experiment tests the impact of the drug hydroxychloroquine on immune responses to group A streptococcus in peripheral blood mononuclear cells
Overall desig	jn	Peripheral blood mononuclear cells from three healthy participants were stimulated with rheumatogenic, heat-killed group A streptococcus for 24 hours, MOI 10. The effect of hydroxychloroquine (20 μ M) (HCQ) was measured, both alone and in combination with group A streptococcus (GAS).
Contributor(s	5)	Martin W, Pacini G, Smyth GK
Citation miss	sing	Has this study been published? Please login to update or notify GEO.
Submission of	date	Oct 21, 2015
Last update (date	May 15, 2019
Contact nam	e	Gordon K Smyth
E-mail(s)		smyth@wehi.edu.au
Phone		(+61 3) 9345 2326
Fax		(+61 3) 9347 0852
URL		http://www.wehi.edu.au
Organization	name	Walter and Eliza Hall Institute of Medical Research
Department		Bioinformatics
Lab		Smyth 10 Devel Dda
Street addres	SS	To Koyai Pue
City State/provin	~~	Vie
State/provin	CC .	vic

https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE74235



0

Q



Step 1. Upload or Fetch RNA-seq Data

- · Upload your raw or processed RNA-seq data
- Fetch >8,000 public RNA-seq datasets published in the Gene Expression Omnibus



- Select from multiple state-of-the-art RNA-seq data
 analysis tools
- Contribute your computational tool as a plugin

Step 3. Generate Your Notebook

- · Access and share your results through a permanent URL
- Download, rerun and customize your notebook using Docker

BioJupies Automatically Generates RNA-seq Data Analysis Notebooks

 With BioJupies you can produce in seconds a customized, reusable, and interactive report from your own raw or processed RNA-seq data through a simple user interface

Get Started

• • • • • • • • •

https://maayanlab.cloud/biojupies/

Cell Systems 2018 Nov 28;7(5):556-561.e3.





 \mathbf{c}

Which dataset would you like to analyze?

Use the form below to search 9,145 publicly available datasets published in the Gene Expression Omnibus database and processed by ARCHS4.

hydroxychlo	proquine							Q
Displaying 1	-1 of 1 results	Organism 🛛 :	All 🗸	Sort by $oldsymbol{ heta}_{:}$	Newest	~	Samples 0:	35 70+
A	Hydroxychloroq GSE74235	uine inhibits resp 15 samples	ponses to group Published Octi	A streptococcu	ıs in peripheral bl	ood mononucl	ear cells	Analyze > More Info ~





Search GEO datasets (e.g. cancer, p53) 🛛 🔍

K Back

Continue >

More Info 🛩

Add

Which analyses would you like to perform?

Use the form below to **add or remove data analysis and visualization tools** to your notebook. These tools will analyze the selected dataset and embed interactive results in your notebook. Once you have selected the desired tools, click **Continue** to proceed.

Differential expression analysis between two groups of samples

Volcano Plot





Help

Which samples would you like to compare?

One or more of the selected tools require generating a gene expression signature Θ . To generate one, you must define two groups of samples whose gene expression you wish to compare by using the form below. Once you have defined the desired groups, click **Continue** to proceed.

Back
 Continue >

What are the names of the groups?

First, name the groups of samples you wish to compare using the text boxes below. It is recommended to use names descriptive of the groups' experimental condition: for example, *WT*, *DMSO*, *Vehicle* for the **Control** group; *Metastatic*, *p53 KO*, *Dasatinib-6h* for the **Perturbation** group.



Which groups do the samples belong to?

Second, select the samples you wish to assign to the groups by using the dropdown menus below. To improve the statistical validity of the gene expression signature, each group must have at least three samples.

Group	*	Accession	Sample Title 👙	Length Of Culture	Subject Id#	Treated With	¢
Untreated -		GSM1915090	Untreated 0h Patient 2	0 hours	2	none (untreated)	
Perturbation 👻]	GSM1915091	HCQ 24h Patient 2	24 hours	2	20 µM hydroxychloroquine (HCQ)	
Untreated -		GSM1915092	Untreated 24h Patient 1	24 hours	1	none (untreated)	
Perturbation 👻]	GSM1915093	HCQ 24h Patient 1	24 hours	1	20 µM hydroxychloroquine (HCQ)	
Untreated -		GSM1915094	Untreated Oh Patient	0 hours	1	none (untreated)	



Search GEO datasets (e.g. cancer, p53) Q



Help

Results

New Notebook Open Notebook

Open	
Notebook	

Untreated vs Perturbation Analysis Notebook | BioJupies

Dataset GSE74235

Signature: Untreated vs Perturbation

Analysis Tools:	PCA Clustergrammer Library S	ize Analysis	Differential Expression Table
Enrichr Links	Gene Ontology Enrichment Analysis	Pathway En	richment Analysis
L1000FWD Qu	ery		



Effects of Hydroxychloroquine on Uninfected PBMCs

Untreated vs Perturbation | Gene Ontology Biological Process (2018 version)

Up-regulated in Perturbation

*cholesterol biosynthetic process (GO:0006695)
*regulation of alcohol biosynthetic process (GO:1902930)
*sterol biosynthetic process (GO:0016126)
*secondary alcohol biosynthetic process (GO:1902653)
*regulation of cholesterol biosynthetic process (GO:0045540)
*regulation of cholesterol metabolic process (GO:0090181)
*regulation of steroid biosynthetic process (GO:0050810)
*cholesterol metabolic process (GO:0008203)
*lipid biosynthetic process (GO:0008610)
*fatty acid derivative biosynthetic process (GO:1901570)
*sterol metabolic process (GO:0016125)
*organic hydroxy compound biosynthetic process (GO:1901617)
*fatty-acyl-CoA biosynthetic process (GO:0046949)
*acetyl-CoA metabolic process (GO:0006084)
*acyl-CoA biosynthetic process (GO:0071616)
5 10 15 20
-log10P

Down-regulated in Perturbation

	lipid homeostasis (GO:0055088)
ľ	anion homeostasis (GO:0055081)
Ì	cholesterol homeostasis (GO:0042632)
Ì	sterol homeostasis (GO:0055092)
Ì	negative regulation of cholesterol storage (GO:0010887)
Ì	spliceosomal complex assembly (GO:0000245)
Ì	reverse cholesterol transport (GO:0043691)
Ì	phospholipid homeostasis (GO:0055091)
Ì	positive regulation of cell activation (GO:0050867)
Ì	interleukin-1 secretion (GO:0050701)
Ì	regulation of cholesterol transport (GO:0032374)
Ì	mRNA splice site selection (GO:0006376)
Ì	response to lipid (GO:0033993)
Ì	regulation of cholesterol storage (GO:0010885)
	phospholipid efflux (GO:0033700)
0	1 2 3 4

-log10P

https://amp.pharm.mssm.edu/biojupies/notebook/ZxbvpHAn5

The L1000FWD Map - ~17K Signatures, ~5K Drugs



https://maayanlab.cloud/L1000FWD/main

Bioinformatics. 2018 Jun 15;34(12):2150-2152

Projecting the Hydroxychloroquine Signatures onto the L1000FWD Map



Projecting the SARS-CoV-2 Signatures onto the L1000FWD Map



Drugs Hitting the Same Region in Expression Space



The L1000FWD Map for the A549 Cell-Line



https://maayanlab.cloud/L1000FWD/graph_page/A549-tSNE_layout.csv

Bioinformatics. 2018 Jun 15;34(12):2150-2152



TenOever Lab (Mount Sinai)



Daisy A. Hoagland, Daniel J.B. Clarke, Rasmus Møller, Yuling Han, Liuliu Yang, Megan L. Wojciechowicz, Alexander Lachmann, Kasopefoluwa Y. Oguntuyo, Christian Stevens, Benhur Lee, Shuibing Chen, Avi Ma'ayan, Benjamin R tenOever. **Modulating the transcriptional landscape of SARS-CoV-2 as an effective method for developing antiviral compounds.** bioRxiv 2020.07.12.199687

L The COVID-19 Drug and Gene Set Library

A collection of drug and gene sets related to COVID-19 research contributed by the community



Drug sets	Gene sets				
A Experim	nental drug sets	🖬 Computational drug sets	♥ Twitter drug sets	🌐 All drug sets	
Draw a Ve	enn diagram		Search in description	n, metadata or drugs: e.g	. Remdesivir
Descri	iption			Drugs	DrugEnrichr link
0 15 SA	RS-CoV-2 inhibito	ors in Vero E6 cells from Jan et a	il. 2021	15 drugs	ď
24 FD/	24 FDA-approved drugs inhibiting SARS-CoV-2 from Xiao et al.			24 drugs	ď
□ <u>25 cor</u>	mpounds against	SARS-CoV-2 via AlphaLISA RBE	-ACE2 assay from Hans	on et 25 drugs	C.



https://maayanlab.cloud/covid19/

Patterns 2020 Sep 11;1(6):100090.

Drugs from Published In-Vitro Screens are Hitting the Same Region in Expression Space



Enriched Upregulated GO Biological Processes

https://appyters.maayanlab.cloud/#/Drugmonizome_Consensus_Terms



Search appyters... Drugmonizome RNA-seq scRNA-seq Enrichr Machine Learning Harmonizome L1000 Compare Sets microRNAs Kinome Aging Drugmonizome Extract, Transform, Load Drugmonizome **Drug Set Consensus** Drugmonizome-ML Drugmonizome ETL: Drugmonizome **Consensus Terms** DrugRepurposingHub Views: 377 Starts: 357 Runs: 169 Retrievals: 131 Views: 47 Starts: 26 Views: 45 Starts: 49 Machine learning pipeline to predict novel Runs: 9 Retrievals: 9 Runs: 29 Retrievals: 82 drug indications from small molecule An appyter to process drug-target and An appyter that queries an input collection attributes drug-mechanism of action associations of drug sets against libraries in from the Drug Repurposing Hub. Drugmonizome and returns the top v0.0.10 Apache-2.0 Machine Learning enriched consensus terms Drugmonizome v0.0.2 MIT Drug Repurposing Hub v0.0.3 MIT Drugmonizome ETL Script Drugmonizome Select Drug Set Enrichment Select Select

https://appyters.maayanlab.cloud/#/?tags=Drugmonizome

PLK1 Inhibitor as a Potential Drug to Treat Diabetic Kidney Disease



https://appyters.maayanlab.cloud/#/L1000FWD_Consensus_Drugs





Drug

Diabetes. 2020 Oct 16;db200580.





https://appyters.maayanlab.cloud/#/?g=consensus



https://lincsproject.org/LINCS/tools

Please cite this article in press as: Clarke et al., Appyters: turning Jupyter Notebooks into data-driven web apps, Patterns (2021), https://doi.org/ 10.1016/j.patter.2021.100213

Patterns



Article

Appyters: turning Jupyter Notebooks into data-driven web apps

Daniel J.B. Clarke,¹ Minji Jeon,¹ Daniel J. Stein,¹ Nicole Moiseyev,¹ Eryk Kropiwnicki,¹ Charles Dai,¹ Zhuorui Xie,¹ Megan L. Wojciechowicz,¹ Skylar Litz,¹ Jason Hom,¹ John Erol Evangelista,¹ Lucas Goldman,¹ Serena Zhang,¹ Christine Yoon,¹ Tahmid Ahamed,¹ Samantha Bhuiyan,¹ Minxuan Cheng,¹ Julie Karam,¹ Kathleen M. Jagodnik,¹ Ingrid Shu,¹ Alexander Lachmann,¹ Sam Ayling,² Sherry L. Jenkins,¹ and Avi Ma'ayan^{1,3,*} ¹Department of Pharmacological Sciences, Mount Sinai Center for Bioinformatics, Icahn School of Medicine at Mount Sinai, One Gustave L. Levy Place, Box 1603, New York, NY 10029, USA ²Pencil Worx Design, 345 West 88th Street, New York, NY 10024, USA ³Lead contact *Correspondence: avi.maayan@mssm.edu https://doi.org/10.1016/j.patter.2021.100213

So What is an Appyter and How You Can Create One?



Developing an Appyter: Start a jupyter-notebook

(venv) my-first-ap	ppyter % pip install appyter	
(venv) my-first-ap	ppyter % jupyter notebook .	
[I 16:46:10.066 No [I 16:46:10.067 No	otebookApp] The Jupyter Note otebookApp] http://localhost	book is running at: :8888/?token=797ac2fae005b3a1754cbeb43c9cc8277964a7ed9e3e1564
Ċ Jupyter	Quit Logout	
Files Running Clusters		Cjupyter Appyter Tutorial (unsaved changes)
elect items to perform actions on them.	Upload New - C	File Edit View Insert Cell Kernel Widgets Help Trusted my-first-appyter O Image: State of the state of
🗌 0 💌 🖿 I	Python 3	
🗋 🗅 venv	Python 3.7.6 64-bit	
requirements.txt	Python 3.7.6 64-bit ('conda-env': conda) 6 B my-first-appyter	<pre>In [1]: #%%appyter init from appyter import magic magic.init(lambda _=globals: _())</pre>
	Other:	In []:
	Folder	
	Terminal	

Developing an Appyter: Start the Appyter Server

(venv) my-first-appyter % appyter --profile=biojupies Appyter\ Tutorial.ipynb

* Running on http://127.0.0.1:5000/ (Press CTRL+C to quit)

. . .

Appyter Tutorial-Jupyter Not 🗙 🕂	_ 🗆 ×	Appyter	× +				- 0	×
← → C ③ localhost:8888/notebooks/Appyter%20Tutorial.ipynb	🖈 🔜 📕 🖉 🚺 E	← → Ç ③ 127.0.0.1:5	000			☆ 🔜	i # D	:
🛗 Apps 📀 Signature Comm 🥝 Signature Comm 🥝 Signature Comm	omm 🕲 Swagger UI 🌹 Rancher 🛛 »	Apps 🕥 Signature Comm	. 🕲 Signature Comm 🕻	Signature Comm	Signature Comm	🚯 Swagger UI 🕴	Rancher	
Jupyter Appyter Tutorial (autosaved)	Cogout Logout	*		Cubmit				
File Edit View Insert Cell Kernel Widgets Help	Trusted my-first-appyter O			Submit				
(5) + 3< 2 K + ↓ NRun ■ C ⇒ Code ↓								
<pre>In [1]: #Mappyter init from appyter import magic magic.init(lambda _=globals: _()) In []:</pre>								
/home/u8sand/Programs/work/my-first-appyter - Konsole		/home/u8sand/Programs/work/my-first-	appyter – Konsole					
File Edit View Bookmarks Settings Help Most Likely you need to configure your SUID sandbox correctly [T 16:52:53.20 NotebookApp] 404 GET /nbextensions/jupyter-widget-from-react/ex .43ms referer-http://localhost:8888/notebook/Untitled.ipynb?kernel_name=my-fir [T 16:52:53.55] NotebookApp] 404 GET /nbextensions/widgets/notebook/js/ettensio eferer-http://localhost:8888/notebooks/Untitled.ipynb?kernel_name=my-firSi27d [W 16:52:55.258.NotebookApp] 404 GET /nbextensions/widgets/notebook/js/ettensio eferer-http://localhost:8888/notebooks/Untitled.ipynb?kernel_name=my-firSt-appy [T 16:53:56.238 NotebookApp] Saving file at /Appyter Tutorial.ipynb [T 16:53:59.976 NotebookApp] Saving file at /Appyter Tutorial.ipynb [T 16:54:00.886 NotebookApp] Kernel restarted: f69d740b-88b2-43e8-a083-c7a7e32 [T 16:54:01.467 NotebookApp] Kernel restarted: f69d740b-88b2-43e8-a083-c7a7e32	File Edit View Bookmarks (venv) UBsand@daniel-pc.my- Server initialized for thre WebSocket transport not ava serving Flask app "appyt Environment: production mANING: This is a used Use a production WSGI se Debug mode: on Running on http://127.0. Restarting with stat Server initialized for thre WebSocket transport not ava	Settings Help first-appyter % appyter anding. nilable. Install eventle cer.render.flask_app.app opment server. Do not i rever instead. 0.1:5000/ (Press CTRL+(eading. nilable. Install eventle	rprofile=biojupi et or gevent and ge p" (lazy loading) use it in a product C to quit) et or gevent and ge	es <u>Appyter\ Tutoria</u> event-websocket for ion deployment. event-websocket for	<u>l.ipynb</u> improved perfo improved perfo	[prmance, ormance,	(0)	

Use Appyter Magic for jinja2 Meta-Programming

Minor changes can turn an existing jupyter notebook into a full-blown web application allowing others to use your data processing pipeline.



Single Source of Truth: Jupyter Notebook (ipynb)



Platform Agnostic Multi-Appyter Orchestration



The Appyters Catalog Currently Contains 77 Appyters



Clarke et al., Appyters: turning Jupyter Notebooks into data-driven web apps, Patterns (2021), https://doi.org/10.1016/j.patter.2021.100213

https://appyters.maayanlab.cloud/

ARCHS⁴ Resource



•Samples can be searched by:

- Meta data, text annotation from SRA
- Data driven sample query (highly expressed and low expressed genes)
- Sample query through functional enrichment
- Manual selection

nature communications

nature > nature communications > articles > article

Article | Open Access | Published: 10 April 2018

Massive mining of publicly available RNA-seq data from human and mouse

Alexander Lachmann, Denis Torre, Alexandra B. Keenan, Kathleen M. Jagodnik, Hoyjin J. Lee, Lily Wang, Moshe C. Silverstein & Avi Ma'ayan 🖂

Nature Communications 9, Article number: 1366 (2018) | Cite this article 26k Accesses | 112 Citations | 69 Altmetric | Metrics



Gene **Expression** Atlas for **Tissues and Cell Lines**





IMR90 A549 MRG5 H1299 NHBE NGH450 BFAS28

JURIZAT MITA BIAB HELGEIIIS HELGO RAII

MCF7 MCF10 MDAME2E1 SXBR3 MDAME433 MDAME433

BJ KB HIF10300 NHDF NHDF

LNCAP FCB DULAS CA2

NALME DAOY JEGB BEWO

11191 119335 (CEM 1194314)

SKOVB 2003 OVCAR3

AABI

LHCNM2

REH RPMI8226

HCTU16
 HT29

SHSY5

K562 1208

ERWO

PANGI

- 0937

HERC2

- 23

C401

• 124

KGLGEIKS

20517 HEX293 20517 FURN/TREX293

Lung

Prediction of Gene Function using ARCHS4 Data



GO | ChEA | Mouse Phenotype | Human Phenotype | KEA | KEGG Tissue Expression | Cell Line Expression

Search

Appliprotein (spo) A-V gene contains 3 exons separated by two introns. A sequence polymorphism has been identified in the 3'UTR of the third exon. The primary translation product is a 398-residue preprotein which after proteosity for processing is accented as primary site of synthesis, the intestine, in association with chylomicron particles. Although its precise function is not known, app A-V is a potent advator of leath-in-direction objectivates in vitro.

Functional Annotation Prediction

Predicted biological processes (GO)

Rank		GO term	Z-score					APOA4	104	
1	exogenous drug catabolic process (GO:0042738)			9.67360960	2	AU	C = 0.898			_
2	digestion (GO:0007586)			8.63908848	Ι.		1			
3	* regulation of triglyceride catabolic process (GO:0010896)			8.26290027	3.	1,	r			ł
4	regulation of guanylate cyclase activity (GO:0031282)			7.94133341	, s.	1				
5	glucocorticoid biosynthetic process (GO:0006704)			7.56956866	÷.,					
6	drug catabolic process (GO:0042737)			7.20406798		11				
7	positive regulation of guany	riate cyclase activity (GO:0	0031284)	6.59351313	8	1				
8	* triglyceride-rich lipoprot	ein particle remodeling (GO:0034370)	6.49610724		10				
9	* regulation of cholesterol	esterification (GO:0010)	372)	6.23647715	1.	0.0	0.2	0.4	0.6	,
10	* phospholipid efflux (GO:	:0033700)		6.03232892				Specifi	iky	

GO term Z-score APOA4 26 BP1_19119308_ChIP-ChIP_Hs578T_Human 9.94549333 AUC = 0.688 ESR1 17901129 ChIP-ChIP LIVER Mouse 5.60571253 4.60205684 CDX2 20551321 ChIP-Seg CACO-2 Human TRIM28 21343339 ChIP-Seq HEK293 Human 4.57369677 P63_17297297_ChIP-ChIP_HaCaT_Human 3.46234525 CTNNB1_24651522_ChIP-Seq_LGR5+_INTESTINAL_STEM_Human 3.29785656 EGR1_23403033_ChIP-Seq_LIVER_Mouse 3.29326701 RXR_22158963_ChIP-Seq_LIVER_Mouse 3.21230272 LXR_22158963_ChIP-Seq_LIVER_Mouse 2.90574163 HNE4A 19761587 ChiP-ChiP CACO-2 Huma 2 861996

Predicted mouse phenotypes (MGI)

Rank		GO term	Z-scor	e	AP0A4 11			
1	MP0005360_urolthiasis			7.36699685	2 AUC = 0.792			
2	MP0002139_abnormal_hepatobilary_system			6.89648342				
3	MP0001666_abnormal_nutrient_absorption			6.71050281				
4	MP0005085_abnormal_gallbladder_physiolo			5.62082136	2 8-			
5	MP0003283_abnormal_digestive_organ			4.66159451	poor a			
6	* MP0009840_abnormal_fe	oam_cell		4.35705070				
7	MP0005365_abnormal_bile	salt		3.64401650	8-			
8	MP0003868_abnormal_fece	is_composition		3.31594048	2			
9	* MP0010329_abnormal_li	poprotein_level		3.09092425	0.0 0.2 0.4 0.8 0.8 10			
10	MP0001986_abnormal_tast	te_sensitivity		3.09004587	Specificity			

Predicted upstream transcription factors (ChEA)

Prediction of Gene Function using ARCHS4 Data



Mining Massive Gene Expression Repositories

- ARCHS4 was used to align 16 trillion reads (16*10^12) to human and mouse reference genomes
- > petabyte of data (1 million GB)
- 900,809 samples from 23,739 experiments



Cost challenges



- CPU
 - Hardware purchases are high upfront and require maintenance
 - CPU rentals are high if used continuously



- Memory
 - Aligning reads to a reference genome requires efficient in memory index structures
 - Depending on alignment algorithm this can vary from 4GB (kallisto/De Bruijn Graph) to 16GB (STAR/suffix tree)

- Storage
 - Raw read files need to be stored for processing
 - Required storage varies significantly between samples
- Network bandwidth
 - Fast networking is required to retrieve raw read data
 - Slow bandwidth will increase overall processing time



Hybrid RNA-seq Pipeline Design

- Processing pipelines are dockerized for easy deployment to resource pool
- Control server hosting job commands serves job descriptions upon request of worker nodes
- Job descriptions are JSON objects
- Results are stored at an online location accessible via URL



Hybrid RNA-seq Pipeline Design

- Parallelization through deployment of dockerized workflows
- Instances are precisely chosen for required resources
 - memory
 - bandwidth
 - storage

Cost comparison of RNA-seq resequencing projects

- The average cost of processing an RNA-seq sample is below \$0.01
- Minimal resource allocation and low-cost high network bandwidth of cloud resources is a key advantage of the ARCHS4 pipeline
- Choice of efficient alignment algorithm (kallisto) results in low memory footprint

RNA-seq resource	ARCHS4	Recount	Toil Recompute
Human samples	84,863	61,350	19,931
Mouse samples	103,083	0	0
Total samples	187,946	61,350	19,931
Cost per sample	< \$0.01	\$0.73	\$1.30



Mount Sinai Center for Bioinformatics Com

Ma'ayan Laboratory of Computational Systems Biology

2021 Summer Research Training Program in Biomedical Big Data Science



Icahn Center for School of Bioinformatics Medicine at Mount Sinai

Research intensive 10-week training program for undergraduate and master's students



Program Dates: June 7 - August 13, 2021

Students selected for summer session 2021 of our **research training program in the Ma'ayan Laboratory at the Icahn School of Medicine at Mount Sinai** will conduct faculty-mentored independent research projects in the following areas:

Data Harmonization

Cloud Computing

Application Deadline: February 1, 2021 at 5 PM Eastern Time

Who Should Apply:

Students majoring in Computer Science, Informatics, Mathematics, Statistics, Physics, Engineering, Chemistry/Chemical Sciences or Biological Sciences and have an interest in Biomedical Big Data Science.

Contact: Sherry Jenkins, MS Program Manager E-mail: <u>sherry.jenkins@mssm.edu</u>

Machine Learning

Dynamic Data Visualization

Faculty Mentor and Principal Investigator: Avi Ma'ayan PhD, Professor and Director Mount Sinai Center Bioinformatics Icahn School of Medicine at Mount Sinai New York, New York

Trainee Salary:

\$8,000 salary for the 10-week training period

Eligibility:

To be considered for this program, applicants must be:

- U.S. citizen or U.S. permanent resident
- Undergraduate or master's student in good academic standing

APPLICATION DETAILS: http://labs.icahn.mssm.edu/maayanlab/summer-research-program/

Summary

- We systematically convert publicly available omics datasets into an abstract format centered on genes and drugs.
- The resources we developed have made a big impact on the research community.
- Appyters can enable the rapid development of bioinformatics applications not just for the Mount Sinai Center for Bioinformatics but also for others at Mount Sinai.
- Machine Learning to impute knowledge about gene and drug functions.
- SignatureCommons is a template system to quickly bring up a data portal to serve data and metadata.

Acknowledgements

Ma'ayan Lab

Sherry Jenkins, MS - Project Manager Daniel Clarke, MS - Data Science Analyst Alexander Lachmann, PhD - Assistant Professor Kathleen Jagodnik, PhD - Postdoctoral Fellow Minji Jeon, PhD - Postdoctoral Fellow Megan Wojciechowicz, MS - PhD Student John Erol Evangelista, MS - Bioinformatician Sherry Xie, BS - Bioinformatician Eryk Kropiwnicki, MS - Bioinformatician Maxim Kuleshov, MS - Bioinformatician Ingrid Shu, BS - Bioinformatician Allison Bailey, MPH - Associate Researcher



BD2K-LINCS DATA COORDINATION AND INTEGRATION CENTER









Icahn Cen School of Bio Medicine at **Mount** Sinai

Center for Bioinformatics