# Canadian Bioinformatics Workshops

www.bioinformatics.ca

Slide Concept by Cameron Neylon, who has waived all copyright and related or neighbouring rights. This slide only ccZero.
Social Media Icons adapted with permission from originals by Christopher Ross. Original images are available under GPL at;
http://www.thisismyurl.com/free-downloads/15-free-speech-bubble-icons-for-popular-websites



Module 6 part 1
Galaxy

BF Francis Ouellette
BF Francis Ouellette
Informatics on High Throughput Sequencing
June 9-10, 2014

E-mail | francis@oicr.on.ca

@bffo

#IHTSD14

#CBW2014

#usegalaxy

Module 6

**bio**informatics.ca

## Disclaimer

- I do not (and will not) profit in any way, shape or form, from any of the brands, products or companies I may mention.
- I am on the Galaxy Scientific Advisory Board (Galaxy's NIH grant), but I do that for free.
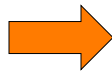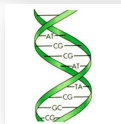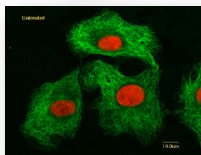
Module 6

**bio**informatics.ca

# Outline

- Workflows & an examples on using Galaxy platform for DNA sequence manipulations.
- Reproducible Science
- Galaxy Public server; Galaxy @home;
  Galaxy in the cloud
- Putting and getting data in and out of Galaxy
- Processing Data in Galaxy
- Example of a Galaxy pipeline on RNA-Seq
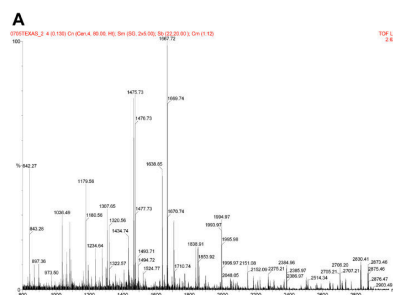- Lab

Module 6

**bio**informatics.ca

# What biologist do:



- Make observations
- Make hypothesis
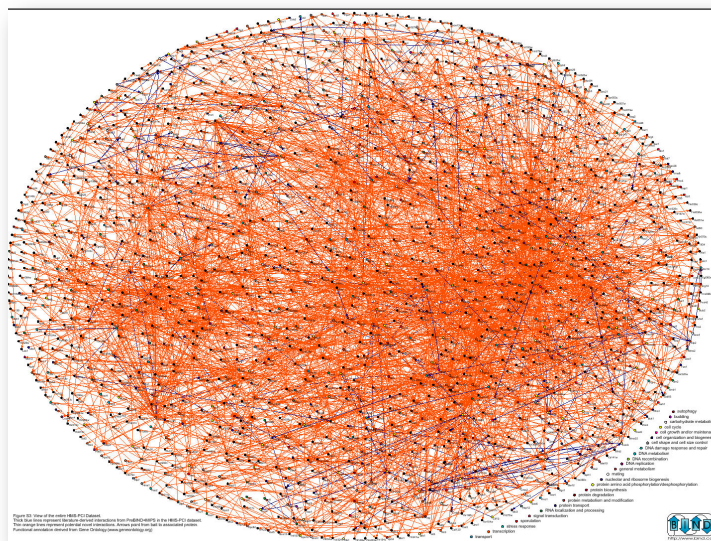- Test them
- Challenge them
- Conclude things
- Write papers

http://goo.gl/7sCUI

Module 6

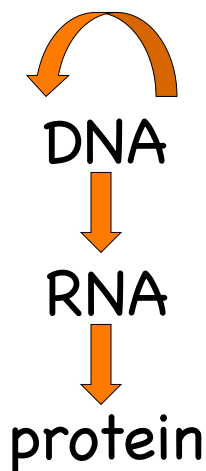**bio**informatics.ca

# Central Dogma

DNA

RNA

protein

Module 6

**bio**informatics.ca

# Central Dogma

DNA

RNA

protein

Then you
write a
paper
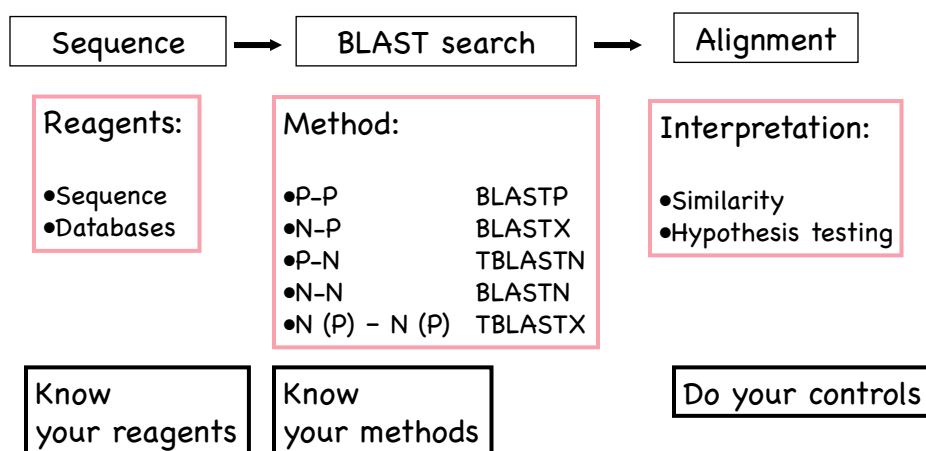about it

Module 6

**bio**informatics.ca

# Some of the things we do when we try and understand the cell ...

- We do experiments
- Some of these are bioinformatics experiments
- We all want these to be reproducible
- We want people to find our data
- We want people to find our methods
- ... and we want them to be able to rerun our experiments, validate our work, move the science forward.

Module 6

**bio**informatics.ca

---

# Bioinformatics experiments:

| Sequence | → | BLAST search | → | Alignment |

**Reagents:**

- Sequence
- Databases

**Method:**

- P–P          BLASTP
- N–P          BLASTX
- P–N          TBLASTN
- N–N          BLASTN
- N (P) – N (P)   TBLASTX

**Interpretation:**

- Similarity
- Hypothesis testing

Know your reagents

Know your methods

Do your controls

Module 6      14      **bio**informatics.ca

# Doing and redoing experiments

- If you do something once, you usually don't need a script. Do it hundreds or thousands of times, you will want something to help you.
- Want to share what you did, providing a script is usually a good way.
- Sometimes though, scripts are too complicated, and don't capture all that is need to do an experiment. For example: the version of a tool you used!
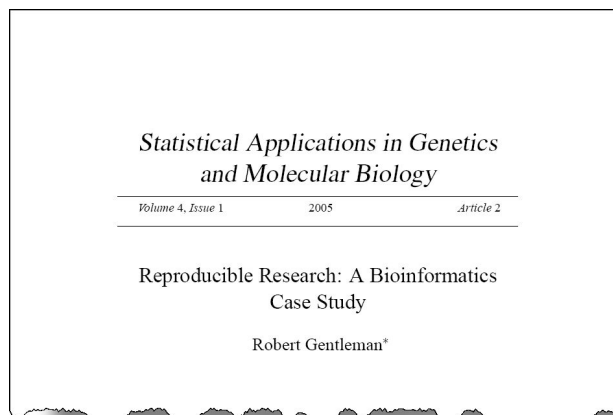
Module 6

**bio**informatics.ca

# Some requirements:

- Open Source
- Solution should be useful to large community
- Well supported (by community and funding agency)
- Flexible
- Expandable
- Scalable
- Cloud-aware
- User friendly?

Open Source

Module 6

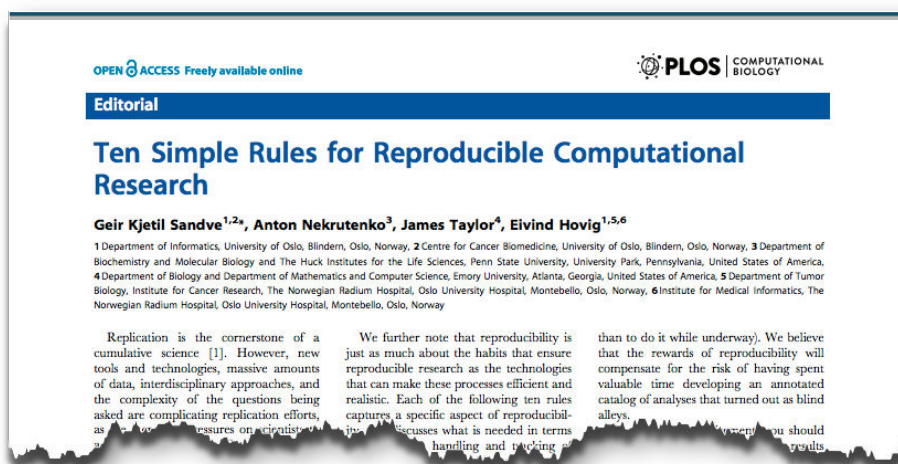**bio**informatics.ca

# Some solutions

1. R and bioconductor (#rstat)

*Statistical Applications in Genetics
and Molecular Biology*

| Volume 4, Issue 1 | 2005 | Article 2 |

Reproducible Research: A Bioinformatics
Case Study

Robert Gentleman*

http://www.ncbi.nlm.nih.gov/pubmed/16646837

Module 6

**bio**informatics.ca

---

OPEN ACCESS Freely available online

PLOS | COMPUTATIONAL BIOLOGY

**Editorial**

# Ten Simple Rules for Reproducible Computational Research

Geir Kjetil Sandve[1,2*], Anton Nekrutenko[3], James Taylor[4], Eivind Hovig[1,5,6]

1 Department of Informatics, University of Oslo, Blindern, Oslo, Norway, 2 Centre for Cancer Biomedicine, University of Oslo, Blindern, Oslo, Norway, 3 Department of Biochemistry and Molecular Biology and The Huck Institutes for the Life Sciences, Penn State University, University Park, Pennsylvania, United States of America, 4 Department of Biology and Department of Mathematics and Computer Science, Emory University, Atlanta, Georgia, United States of America, 5 Department of Tumor Biology, Institute for Cancer Research, The Norwegian Radium Hospital, Oslo University Hospital, Montebello, Oslo, Norway, 6 Institute for Medical Informatics, The Norwegian Radium Hospital, Oslo University Hospital, Montebello, Oslo, Norway

Replication is the cornerstone of a cumulative science [1]. However, new tools and technologies, massive amounts of data, interdisciplinary approaches, and the complexity of the questions being asked are complicating replication efforts, as are the various pressures on scientists

We further note that reproducibility is just as much about the habits that ensure reproducible research as the technologies that can make these processes efficient and realistic. Each of the following ten rules captures a specific aspect of reproducibil- ity, and discusses what is needed in terms of handling and tracking

than to do it while underway). We believe that the rewards of reproducibility will compensate for the risk of having spent valuable time developing an annotated catalog of analyses that turned out as blind alleys. We argue that you should results

http://goo.gl/j8kCgd

Module 6

**bio**informatics.ca

**Ten Simple Rules for Reproducible Computational Research**

Rule 1: For Every Result, Keep Track of How It Was Produced
Rule 2: Avoid Manual Data Manipulation Steps
Rule 3: Archive the Exact Versions of All External Programs Used
Rule 4: Version Control All Custom Scripts
Rule 5: Record All Intermediate Results, When Possible in
       Standardized Formats
Rule 6: For Analyses That Include Randomness, Note Underlying
       Random Seeds
Rule 7: Always Store Raw Data behind Plots
Rule 8: Generate Hierarchical Analysis Output, Allowing Layers
       of Increasing Detail to Be Inspected
Rule 9: Connect Textual Statements to Underlying Results
Rule 10: Provide Public Access to Scripts, Runs, and Results

Module 6

**bio**informatics.ca

# Some solutions (2)

- SeqWare : http://seqware.github.io/



Module 6

**bio**informatics.ca

# Some solutions (3)

- Galaxy



Module 6

**bio**informatics.ca



Goecks et al. Genome Biology 2010, **11**:R86
http://genomebiology.com/2010/11/8/R86

Genome **Biology**

SOFTWARE — Open Access

Galaxy: a comprehensive approach for supporting accessible, reproducible, and transparent computational research in the life sciences

Jeremy Goecks[1], Anton Nekrutenko[2*], James Taylor[1*], The Galaxy Team

http://genomebiology.com/2010/11/8/R86

Module 6

**bio**informatics.ca

Using Cloud Computing Infrastructure with CloudBioLinux, CloudMan, and Galaxy

**UNIT 11.9**

Enis Afgan,[1,5] Brad Chapman,[2] Margita Jadan,[3] Vedran Franke,[4] and James Taylor[5]

[1]Center for Informatics and Computing, Ruđer Bošković Institute (RBI), Zagreb, Croatia
[2]Harvard School of Public Health, Boston, Massachusetts
[3]Division of Materials Chemistry, Laboratory for Ichthyopathology–Biological Materials, Ruđer Bošković Institute (RBI), Zagreb, Croatia
[4]Department of Biology, University of Zagreb, Zagreb, Croatia
[5]Department of Biology and Department of Mathematics and Computer Science, Emory University, Atlanta, Georgia

ABSTRACT

Cloud computing has revolutionized availability and access to computing and storage ...

http://onlinelibrary.wiley.com/doi/10.1002/0471250953.bi1109s38/pdf

Module 6

**bio**informatics.ca

---

# Which Galaxy?

- galaxyproject.org: Galaxy home page
- usegalaxy.org: main Galaxy public server
- getgalaxy.org: source for installing local Galaxy
- usegalaxy.org/cloud: use galaxy in the cloud
- http://goo.gl/mlyOC : Other public Galaxy servers

| | Main | Local | Cloud | Other |
|---|---|---|---|---|
| Your data sets are moderately sized | Yes | Yes | Yes | ? |
| Your computational requirements are moderate | Yes | Yes | Yes | ? |
| You want to share your Galaxy objects with others | Yes | Yes | Yes | ? |
| All needed Tools are installed on Main. | Yes | ? | Yes | ? |
| Your data sets are very large | No | ? | Yes | ? |
| Your computational requirements are very large | No | ? | Yes | ? |
| You have absolute data security requirements | No | Yes | Yes | ? |

http://goo.gl/x3DXm

Module 6

**bio**informatics.ca

Module 6

bioinformatics.ca



Module 6

bioinformatics.ca

# getgalaxy.org



Module 6

bioinformatics.ca

# usegalaxy.org/cloud



Module 6

bioinformatics.ca

# http://wiki.galaxyproject.org/PublicGalaxyServers



Module 6

bioinformatics.ca

---



- Galaxy integrates input data sources
- Galaxy allows you to use many tools that you don't need to install and maintain.
- Galaxy allows you to maintain workflows, reuse them, and share them.
- Galaxy lets you "publish" experiments.
- Galaxy has fully entered the "next-gen" space.
- Galaxy works in the cloud.

Module 6

bioinformatics.ca

**Galaxy = collaboration and reproducibility**

Best of all, Galaxy's history system provides a complete analyses record that can be shared. Every history is an analysis workflow, which can be used to reproduce the entire experiment...

- **History is an analysis record** | Every step of your analyses is recorded in Galaxy's history system. You can have any number of histories saved. This way you can go back to your analyses anytime.
- **Share your analyses** | Alice works at Penn State, while Bob suffers from the terrible San Diego climate. Alice wants Bob to see her analyses. Alice clicks the "share" link and enters Bob's e-mail address. Now Alice's history is visible to Bob (see "Sharing history" screencast).
- **Now your results are reproducible!** | When publishing results, replace "the data were analyzed using a collection of in-house scripts" with a URL pointing to Galaxy's history. Your reviewers will have no further questions. That's reproducible genomics!

- Galaxy strongly believes on reproducibility!
- Galaxy is very good at keeping a history of what you did, and allow you do it again when you need to, or allow somebody else to do it again.
- Galaxy makes it very easy to work with collaborators down the hall, or across the globe.

Module 6                                                    **bio**informatics.ca

---

**Designed for biologists and developers**

Yep, sometimes you can mix water with oil...

- **Biologists** | Use our public site to access popular sources of data like the UCSC Table Browser. Run analyses right on the spot using a variety of integrated tools. Your results are never deleted and can be easily shared with others.
- **Developers** | Galaxy is an easy-to-use, open-source, scalable framework for tool and data integration. Stop wasting time writing interfaces and get your tools used by biologists! Galaxy includes everything you need to get started, so download and start integrating!

- Galaxy is designed with biologists in mind, and basically thinks like we do (most of the time!)
- Galaxy has a healthy developer community, and is very present in forums of other Open Source initiatives.

Module 6                                                    **bio**informatics.ca

**Why did we do it?**

You are an experimental biologist. You keep watching databases fill with more and more data. You keep thinking: *even if I knew how to use Excel as a pro, it would probably not load 12,435,654 SNPs.* So how do you perform analyses without calling somebody on the Computer Science side of campus? Suppose you want to find human promoters with the highest SNP density. There is no straightforward way of doing it without learning programming first. And this is why...

• **Databases are not analyses tools** | Databases are where you get the data. Browsers are where you visualize the results. For a bench biologist there is not much in between besides spreadsheets or Perl scripting.
• **No tools for new datatypes** | Some datatypes generated by high throughput genomics are so new that there are no tools to analyze them. For example, how do you extract sequences of coding exons from the latest 28-way alignments of vertebrate genomes or analyze quality scores from 454/Solexa/SOLiD? With Galaxy.
• **Genomics is not really reproducible** | The Methods section of too many papers sound like *the data were analyzed using a collection of in-house scripts.* How do you repeat such a study? Galaxy saves every step of your analysis and allows you to share these workflows with others.
• **Too many tools** | *Bioinformatics* publishes hundreds of application notes per year. How does one know which tool to use? Galaxy integrates a multitude of different tools by giving them the same "look and feel" and linking them to data warehouses.

- To help biologists deal with tools and data.
- Funding: NIH, NSF, & Penn State University.
- Development: Emory University and Penn State
- http://wiki.galaxyproject.org/
- http://wiki.galaxyproject.org/Learn

Module 6

**bio**informatics.ca

# Challenge with multiple sites/model

- Not all galaxy are created the same
- Galaxy team moving to an "empty" shell, and cafeteria model: take only what you need.
- Adding tools and updating tools causes problems sometimes, but Galaxy team is working to make this easier
- The Toolshed is a great solution for this!

Module 6

**bio**informatics.ca

**Galaxy Toolshed:**
**http://toolshed.g2.bx.psu.edu/**

Module 6



# Galaxy Toolshed: SAM

http://toolshed.g2.bx.psu.edu/

Module 6

# General workflow for Galaxy



Module 6

**bio**informatics.ca

# Time for sponsor announcement!



http://www.pmgenomics.ca/

Zhibin Lu

Module 6

**bio**informatics.ca

Click on the link

Module 6

**bio**informatics.ca



Enter the user name and password you used on the first form. i.e cbw#

Module 6

**bio**informatics.ca

Keep default
Press in Choose platform type

Module 6

bioinformatics.ca



Module 6

bioinformatics.ca

Press on Access Galaxy

Module 6

**bio**informatics.ca



Module 6

**bio**informatics.ca

# What is next?

- I'm going to tell you about getting data in and out of Galaxy
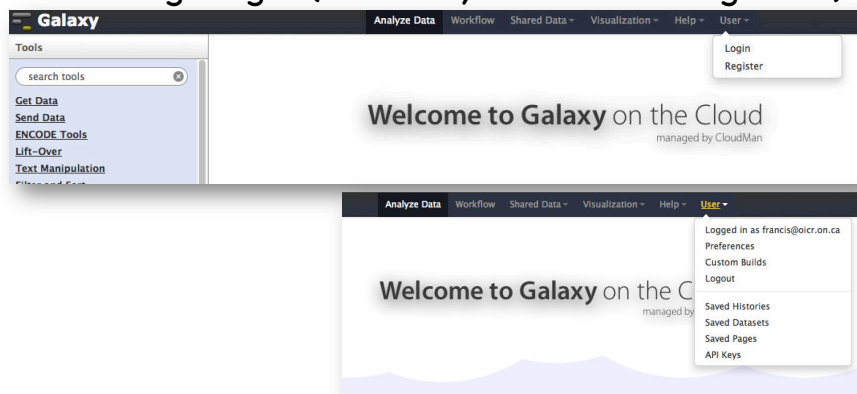- Doing operations in Galaxy
- Understanding the user interface.
- Linking multiple steps into "pipelines"
- Do an RNASeq mapping experiment
- Sharing pipeline with colleagues, and making them public.
- How to learn more ...

Module 6                                    **bio**informatics.ca
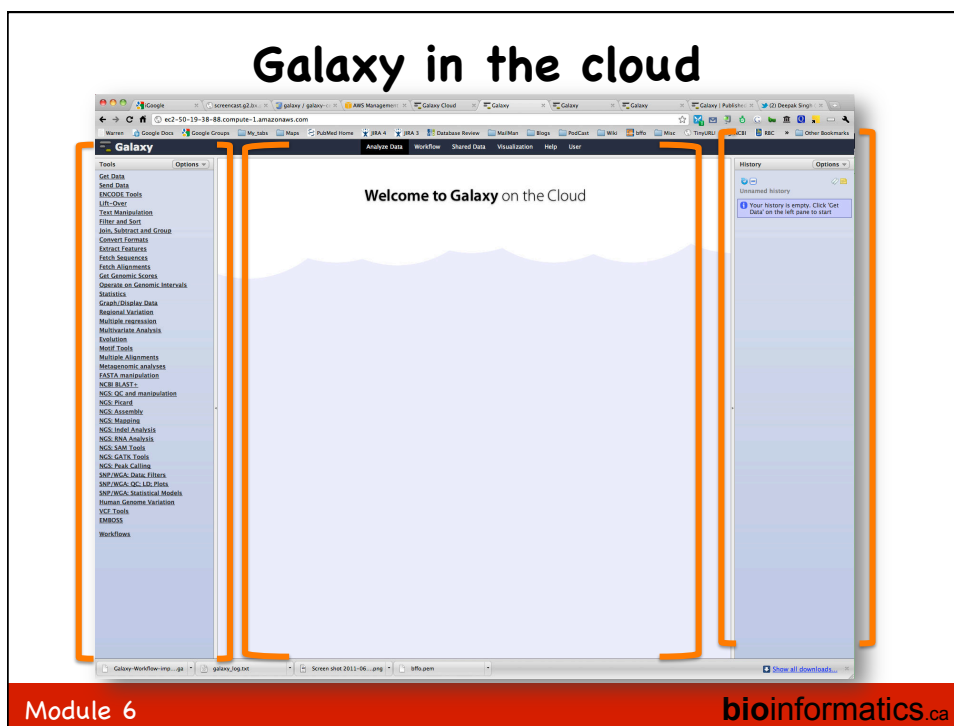
# 1st thing to do before we start:

- This is important, irrespective of which cloud you are using: Login (1st time you need to "register")



Module 6                                    **bio**informatics.ca

24

# Galaxy in the cloud



Module 6

**bio**informatics.ca

---

- Get Data
- Send Data
- ENCODE Tools
- Lift–Over
- Text Manipulation
- Filter and Sort
- Join, Subtract and Group
- Convert Formats
- Extract Features
- Fetch Sequences
- Fetch Alignments
- Get Genomic Scores
- Operate on Genomic Intervals
- Statistics
- Graph/Display Data
- Regional Variation
- Multiple regression

- Multivariate Analysis
- Evolution
- Motif Tools
- Multiple Alignments
- Metagenomic analyses
- FASTA manipulation
- NGS: QC and manipulation
- NGS: Assembly
- NGS: Mapping
- NGS: Indel Analysis
- NGS: RNA Analysis
- NGS: SAM Tools
- NGS: GATK Tools
- NGS: Peak Calling
- SNP/WGA: Data; Filters
- SNP/WGA: QC; LD; Plots
- SNP/WGA: Statistical Models
- Human Genome Variation
- VCF Tools

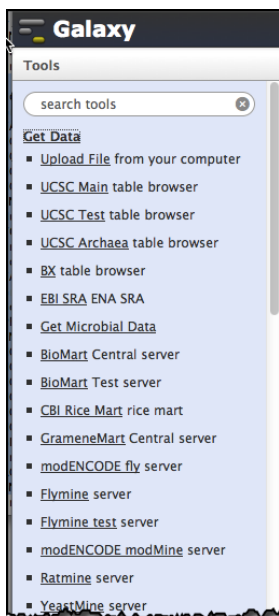Module 6

**bio**informatics.ca

## Galaxy cloud    usegalaxy.org

- < NGS: Assembly
- < NGS: GATK Tools
- < SNP/WGA: Statistical Models
- < Human Genome Variation
- < VCF Tools

- > Genome Diversity
- > Phenotype Association
- > EMBOSS
- > NGS Toolbox Beta
- > NGS: GATK Tools (beta)
- > NGS: Variant Detection
- > NGS: Picard (beta)
- > BEDTools
- > snpEff
- > RGENETICS
- > SNP/WGA: Statistical Models

Module 6

**bio**informatics.ca

---

**Galaxy**

Tools

search tools

Get Data
- Upload File from your computer
- UCSC Main table browser
- UCSC Test table browser
- UCSC Archaea table browser
- BX table browser
- EBI SRA ENA SRA
- Get Microbial Data
- BioMart Central server
- BioMart Test server
- CBI Rice Mart rice mart
- GrameneMart Central server
- modENCODE fly server
- Flymine server
- Flymine test server
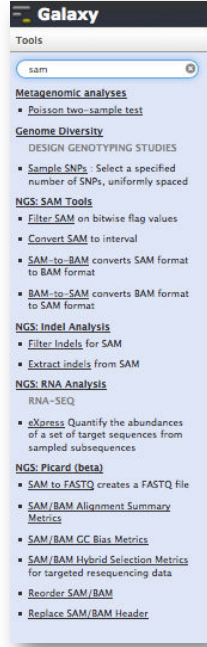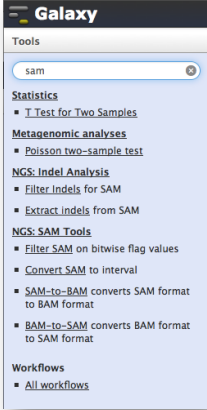- modENCODE modMine server
- Ratmine server
- YeastMine server

- ... and each item, when you click on it expands to lots more choices!
- What I find most useful when I know the name of the tool I'm looking for is to simply using the search tool.
- E.g. look for "sam"

Module 6

**bio**informatics.ca

usegalaxy.org

Galaxy cloud

Module 6

bioinformatics.ca

---

# UCSC Genome Browser: source of data for Galaxy

- Browse many Eukaryotic genomes (yeast to human)
- Most annotations are there
- Important evolutionary and variation data representation.
- Very flexible and configurable views
- Graphical and table views (Galaxy uses this)
- Upload your data into custom tracks and share with colleagues
- Client/server application with it's issues, but a great app!

Module 6

bioinformatics.ca

Module 6

bioinformatics.ca

# Other Examples of Data Format outputs from UCSC:

- Tab-separated
- Sequence (FASTA)
- Browser Extensible Data format (BED)
- General Feature Format (GFF)
- Gene Transfer Format (GTF)

Module 6

bioinformatics.ca

## Examples of Data Formats for UCSC:

- Sequence (FASTA):

```
>gi|89058412|ref|NT_028395.3| Homo sapiens chromosome 22 genomic contig, GRCh37.p5
Primary Assembly
GATCTGATAAGTCCCAGGACTTCAGAAGAGCTGTGAGACCTTGGCCAAGTCACTTCCTCCTTCAGGAACA
TTGCAGTGGGCCTAAGTGCCTCCTCTCGGGACTGGTATGGGGACGGTCATGCAATCTGGACAACATTCAC
CTTTAAAAGTTTATTGATCTTTTGTGACATGCACGTGGGTTCCCAGTAGCAAGAAACTAAAGGGTCGCAG
GCCGGTTTCTGCTAATTTCTTTAATTCCAAGACAGTCTCAAATATTTTCTTATTAACTTCCTGGAGGGAG
GCTTATCATTCTCTCTTTTGGATGATTCTAAGTACCAGCTAAAATACAGCTATCATTCATTTTCCTTGAT
TTGGGAGCCTAATTTCTTTAATTTAGTATGCAAGAAAACCAATTTGGAAATATCAACTGTTTTGGAAACC
TTAGACCTAGGTCATCCTTAGTAAGATCTTCCCATTTATATAAATACTTGCAAGTAGTAGTGCCATAATT
ACCAAACATAAAGCCAACTGAGATGCCCAAAGGGGGCCACTCTCCTTGCTTTTCCTCCTTTTTAGAGGAT
TTATTTCCCATTTTTCTTAAAAAGGAAGAACAAACTGTGCCCTAGGGTTTACTGTGTCAGAACAGAGTGT
GCCGATTGTGGTCAGGACTCCATAGCATTTCACCATTGAGTTATTTCCGCCCCCTTACGTGTCTCTCTTC
AGCGGTCTATTATCTCCAAGAGGGCATAAAACACTGAGTAAACAGCTCTTTTATATGTGTTTCCTGGATG
AGCCTTCTTTTAATTAATTTTGTTAAGGGGATTTCCTCTAGGGCCACTGCACGTCATGGGGAGTCACCCC
AGACACTCCCAATTGGCCCCTTGTCACCCAGGGGCACATTTCAGCTATTTGTAAAACCTGAAATCACTAG
AAAGGAATGTCTAGTGACTTGTGGGGGCCAAGGCCCTTGTTATGGGGATGAAGGCTCTTAGGTGGTAGCC
CTCCAAGAGAATAGATGGTGAATGTCTCTTTTCAGACATTAAAGGTGTCAGACTCTCAGTTAATCTCTCC
TAGATCCAGGAAAGGCCTAGAAAAGGAAGGCCTGACTGCATTAATGGAGATTCTCTCCATGTGCAAAATT
TCCTCCACAAAAGAAATCCTTGCAGGGCCATTTTAATGTGTTGGCCCTGTGACAGCCATTTCAAAATATG
TCAAAAAATATATTTTGGAGTAAAAATACTTTCATTTTCCTTCAGAGTCTGCTGTCGTATGATGCCATACC
AGAGTCAGGTTGGAAGTAAGCCACATTATACAGCGTTAACCTAAAAAAACAAAAAACTGTCTAACAAGA
TTTTATGGTTTATAGAGCATGATTCCCCGGACACATTAGATAGAAATCTGGGCAAGAGAAGAAAAAAAGG
```

## Browser Extensible Data format (BED)

```
track name=pairedReads description="Clone Paired Reads" useScore=1
chr22 1000 5000 cloneA 960 + 1000 5000 0 2 567,488, 0,3512
chr22 2000 6000 cloneB 900 - 2000 6000 0 2 433,399, 0,3601
```

The first three required BED fields are:

1. **chrom** - The name of the chromosome (e.g. chr3, chrY, chr2_random) or scaffold (e.g. scaffold10671).
2. **chromStart** - The starting position of the feature in the chromosome or scaffold. The first base in a chr
3. **chromEnd** - The ending position of the feature in the chromosome or scaffold. The *chromEnd* base is *chromStart=0, chromEnd=100*, and span the bases numbered 0-99.

The 9 additional optional BED fields are:

4. **name** - Defines the name of the BED line. This label is displayed to the left of the BED line in the Gen mode.
5. **score** - A score between 0 and 1000. If the track line *useScore* attribute is set to 1 for this annotation da darker gray). This table shows the Genome Browser's translation of BED score values into shades of gr

| shade | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| score in range | ≤ 166 | 167-277 | 278-388 | 389-499 | 500-611 | 612-722 | 723-833 | 834-944 | ≥ 945 |

http://goo.gl/agfWu

# General Feature Format (GFF)

**GFF format**

GFF (General Feature Format) lines are based on the GFF standard file format. GFF lines have nine required fields that *must* be tab-separated. If the fields are separated by spaces instead of tabs, the more information on GFF format, refer to http://www.sanger.ac.uk/resources/software/gff/.

If you would like to obtain browser data in GFF (GTF) format, please refer to Genes in gtf or gff format on the Wiki.

Here is a brief description of the GFF fields:

1. **seqname** - The name of the sequence. Must be a chromosome or scaffold.
2. **source** - The program that generated this feature.
3. **feature** - The name of this type of feature. Some examples of standard feature types are "CDS", "start_codon", "stop_codon", and "exon".
4. **start** - The starting position of the feature in the sequence. The first base is numbered 1.
5. **end** - The ending position of the feature (inclusive).
6. **score** - A score between 0 and 1000. If the track line *useScore* attribute is set to 1 for this annotation data set, the *score* value will determine the level of gray in which this feature is displayed there is no score value, enter ".".
7. **strand** - Valid entries include '+', '-', or '.' (for don't know/don't care).
8. **frame** - If the feature is a coding exon, *frame* should be a number between 0-2 that represents the reading frame of the first base. If the feature is not a coding exon, the value should be '.'.
9. **group** - All lines with the same group are linked together into a single item.

*Example:*
Here's an example of a GFF-based track. This example can be pasted into the browser without editing. NOTE: Paste operations on some operating systems will replace tabs with spaces, which will track is uploaded. You can circumvent this problem by pasting the URL of the above example (http://genome.ucsc.edu/goldenPath/help/regulatory.txt) instead of the text itself into the custom anno encounter an error when loading a GFF track, check that the data lines contain tabs rather than spaces.

```
browser position chr22:10000000-10025000
browser hide all
track name=regulatory description="TeleGene(tm) Regulatory Regions"
visibility=2
chr22  TeleGene enhancer  10000000  10001000  500 + .  touch1
chr22  TeleGene promoter  10010000  10010100  900 + .  touch1
chr22  TeleGene promoter  10020000  10025000  800 - .  touch2
```

Click here to display this track in the Genome Browser.

http://goo.gl/agfWu

Module 6                                                                 **bio**informatics.ca

---

# Gene Transfer Format (GTF)

- Like GFF, but specific to exon and CDS features, and has one extra field:

The attribute list must begin with the two mandatory attributes:

- **gene_id** *value* - A globally unique identifier for the genomic source of the sequence.
- **transcript_id** *value* - A globally unique identifier for the predicted transcript.

gene_id "Em:U62317.C22.6.mRNA"; transcript_id "Em:U62317.C22.6.mRNA"; exon_number 1

Module 6                                                                 **bio**informatics.ca

**General workflow for Galaxy**

login | Get data | Upload your data | Manipulate you data | Save your output | Save your workflow | Publish your page

Module 6 — bioinformatics.ca



**Pages in Galaxy**

- https://usegalaxy.org/page/list_published

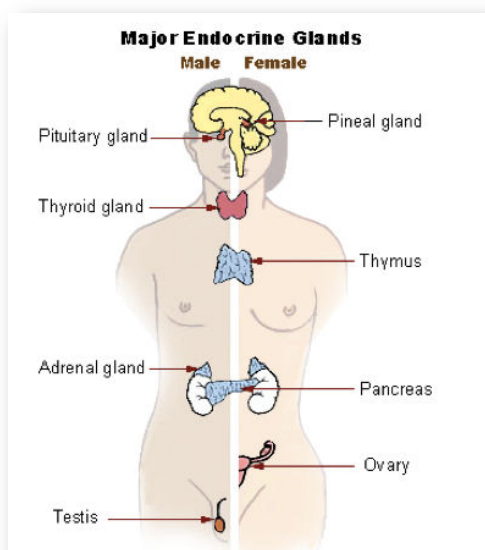Module 6 — bioinformatics.ca

**Module 6**

bioinformatics.ca



**Module 6**

bioinformatics.ca

# RNA-Seq Analysis Exercise

- Human BodyMap 2.0 data from Illumina.
- **paired-end** 50bp reads from **adrenal** and **brain** tissues. The sampled reads map mostly to a 500Kb region of chromosome 19, positions 3–3.5 million (chr19:3000000–3500000).



Module 6                                                                 **bio**informatics.ca



http://en.wikipedia.org/wiki/Adrenal_gland

Module 6                                                                 **bio**informatics.ca

# Getting data

- Most of time, you will get from a file on your computer, or from a URL.



Module 6

**bio**informatics.ca

# Get 4 files

- adrenal_1
  https://usegalaxy.org/dataset/display?
  dataset_id=d44d2a324474d1aa&to_ext=fastqsanger
- adrenal_2
  https://usegalaxy.org/dataset/display?
  dataset_id=d08360a1c0ffdc62&to_ext=fastqsanger
- brain_1
  https://usegalaxy.org/dataset/display?
  dataset_id=f187acb8015d6c7f&to_ext=fastqsanger
- brain_2
  https://usegalaxy.org/dataset/display?
  dataset_id=08c45996966d7ded&to_ext=fastqsanger

Module 6

**bio**informatics.ca

Load file(s) to Galaxy

Module 6



Module 6

"Poke the eye"

"Edit attribute"

"Delete"

"Numbers may vary with usage"

Module 6

**bio**informatics.ca

# "poke the eye"



Module 6

**bio**informatics.ca

## "Edit attributes"

Attributes   Convert Format   Datatype   Permissions

**Edit Attributes**

**Name:**

brain_1.fastqsanger

**Info:**

https://usegalaxy.org/dataset
/display?dataset_id=f187acb8015d6c

**Annotation / Notes:**

Add an annotation or notes to a dataset; annotations are available when a h

**Database/Build:**

Human Feb. 2009 (GRCh37/hg19) (hg19)

Save

Auto-detect
This will inspect the dataset and attempt to correct the above column values

Module 6

**bio**informatics.ca

## General workflow for Galaxy

login

Get data

Upload your data

Manipulate you data

Save your output

Save your workflow

Publish your page

Module 6

**bio**informatics.ca

37

# Need to remove bad bases in reads?

- Assume a median quality score of below 20 to be unusable.
- Given this criterion, is trimming needed for the datasets?
- If so, which base pairs should be trimmed?
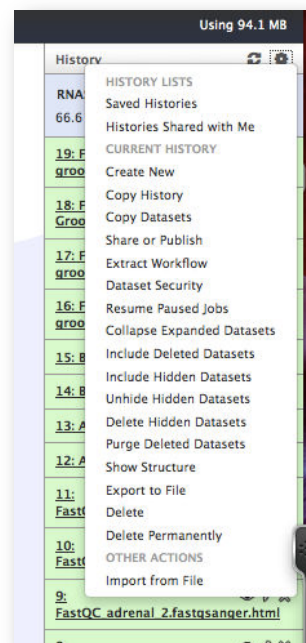- [NGS: QC and manipulation >] FASTQ Trimmer

Module 6

**bio**informatics.ca



"Numbers may vary with usage"

Module 6

**bio**informatics.ca

# [NGS: RNA Analysis >] Tophat tool

- Step 1
- Use the [NGS: RNA Analysis >] Tophat tool
- To map RNA-seq reads to the hg19 Canonical Female build.
- Because the reads are paired, you'll need to set mean inner distance between pairs; this is the average distance in base pairs between reads, not the total insert/fragment size.
- Use a mean inner distance of 110 for BodyMap data.

Module 6     **bio**informatics.ca



Module 6     **bio**informatics.ca

Module 6

bioinformatics.ca



Module 6

bioinformatics.ca

Module 6

bioinformatics.ca



sharing

Module 6

bioinformatics.ca

- Share history with colleagues
- Extract workflow

Module 6

bioinformatics.ca



Module 6

bioinformatics.ca

Module 6

**bio**informatics.ca

# Remember, lots of tutorials, videos, mailing list, twitter etc ...

• https://vimeo.com/galaxyproject



Module 6

**bio**informatics.ca

Module 6

bioinformatics.ca



Module 6

bioinformatics.ca

Module 6 — bioinformatics.ca



http://genomespace.org/

Module 6 — bioinformatics.ca

```
ArrayExpress:   http://www.ebi.ac.uk/arrayexpress/
Cystrome:       http://www.cistrome.org
Cytoscape:      http://www.cytoscape.org/
Galaxy:         http://usegalaxy.org
GenePattern:    http://www.broadinstitute.org/cancer/software/genepattern/
Genomica:       http://genomica.weizmann.ac.il/
geWorkbench:    http://www.geworkbench.org
Gitools:        http://www.gitools.org/
IGV:            http://www.broadinstitute.org/igv/
InSilico DB:    https://insilico.ulb.ac.be/
ISACreator      http://isatab.sourceforge.net/tools.html
MSigDB:         http://www.broadinstitute.org/gsea/msigdb/
UCSC GB:        http://genome.ucsc.edu/
```

Module 6  **bio**informatics.ca

# Useful Resources

- Galaxy
  - usegalaxy.org and usegalaxy.org/cloud
  - Twitter: @galaxyproject #usegalaxy
  - User's mailing list:
    http://lists.bx.psu.edu/listinfo/galaxy-user

- BioStaR
  - biostars.org
  - Twitter: @biostarquestion

Module 6  **bio**informatics.ca

## Useful Resources

- OpenHelix
  - http://www.openhelix.com/
  - Twitter: @openhelix
  - Blog: http://blog.openhelix.com/

- UCSC
  - http://genome.ucsc.edu/
  - Twitter: @GenomeBrowser
  - More tutorials: http://genome.ucsc.edu/training.html

- SEQanswers
  - Forum for NGS technologies
    http://seqanswers.com/

Module 6

**bio**informatics.ca

---

# Papers of interest:

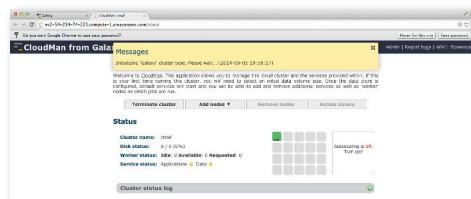- Robert Gentleman, 2005, Reproducible research: a bioinformatics case Source, Stat Appl Genet Mol Biol. 2005;4:Article2.
  http://www.ncbi.nlm.nih.gov/pubmed/?term=16646837
- Goecks J, Nekrutenko A, Taylor J; Galaxy Team. (2010) Galaxy: a comprehensive approach for supporting accessible, reproducible, and transparent computational research in the life sciences. Genome Biology 2010, 11:R86
  http://www.ncbi.nlm.nih.gov/pubmed/?term=20738864
- Afgan E, Chapman B, Jadan M, Franke V, Taylor J. (2012) Using cloud computing infrastructure with CloudBioLinux, CloudMan, and Galaxy. Curr Protoc Bioinformatics. Chapter 11:Unit11.9. doi: 10.1002/0471250953.bi1109s38.
  http://www.ncbi.nlm.nih.gov/pubmed/22700313
- Goecks J1, Eberhard C, Too T; Galaxy Team, Nekrutenko A, Taylor J. (2013) Web-based visual analysis for high-throughput genomics.
  BMC Genomics. 2013 Jun 13;14:397
  http://www.ncbi.nlm.nih.gov/pubmed/23758618

Module 6

**bio**informatics.ca

# Before Coffee Break

- Go to page 20 (or there about) and do:
  **Log onto Galaxy**
- Login info will be on wiki (at that time)
- Once you have this image, you can go on break:



Module 6      **bio**informatics.ca

---

# After Break we will be doing lab

- Want to acknowledge Florence Cavalli and Zhibin Lu for great work they have done to help me with the cloud, some of the slides and with the accuracy of the slides.
- That said, all errors, mistakes, old URLs etc are my fault, entirely!

  @bffo

Module 6      **bio**informatics.ca

# We are on a Coffee Break & Networking Session

- For those of you not here, watching video, maybe you want to register for workshop?
- More details at http://bioinformatics.ca

Module 6

**bio**informatics.ca