# Jenkins as a Scientific Data and Image Processing Platform

Ioannis K. Moutsatsos, Ph.D., M.SE.
Novartis Institutes for Biomedical Research
www.novartis.com

June 18, 2014

#jenkinsconf

# Life Sciences
# are Computational Sciences

- Modern life sciences (biomedical research, systems biology) are heavily dependent on
  - Data Management
  - Computational Analysis
  - Computational Modeling
- Modern laboratory technologies and instrumentation generate data that are
  - Big
  - Heterogeneous
  - Complex

# Computational Challenges & Opportunities

## Scientists

- Face daily challenges by continuing increases in computational complexity
- Focused on the biology and not the compute problem
- Have varying and rapidly changing requirements

## Life Sciences Research

- Benefits from computational systems that are
  - Easy to use
  - Fast to implement
  - Flexible
  - Support
    - Collaboration
    - Transparency
    - Automation
    - Reproducible Research
    - Open standards

# Talk Outline

- A life sciences computational challenge
  - High Content Image Analysis
    - What is it?
- Jenkins-CI as a scientific data/image processing platform
  - Functionality with standard plugins
  - How Jenkins-CI provided a HP image analysis platform for lab scientists
- Jenkins as a data analytics platform
  - Domain specific analysis and visualization plugins
  - The Jenkins pros and cons
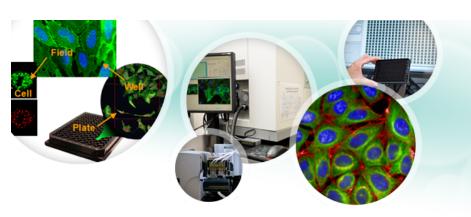    - What are we missing?
- Where do we want to take Jenkins?

# High Content Screening: HCS

*High throughput automated fluorescent microscopy for drug discovery*

- Wet Lab Workflow
  - Cells grown on high density arrays
  - Cells treated with large number of chemical or biological factors
  - Cells stained with fluorescent antibodies

- Data Acquisition
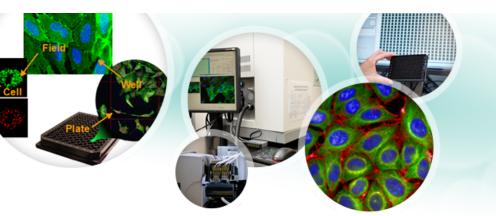  - Stained cells are imaged in high throughput mode using a computerized microscope

- Computational Workflow
  - Cell images processed to extract phenotypic measurements
  - Measurements analyzed to understand factor effects
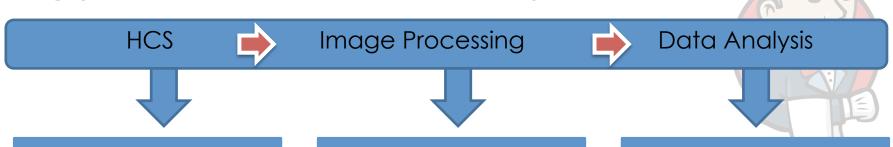
# High Content Screening

- Novartis
  - High Throughput Biology (my group)
    - Data from 2010-2013

- Captured
  - 83 Terabytes of high content image data
    - 17.5 million wells
    - 27 million images
    - ~540 days of imaging time
    - ~1.5 years of computing time

# HCS: Workflow and Data Stream

| HCS | ➡ | Image Processing | ➡ | Data Analysis |
|---|---|---|---|---|

**Raw Data**

- Images
  - channels
  - fields
- Metadata
  - Acquisition
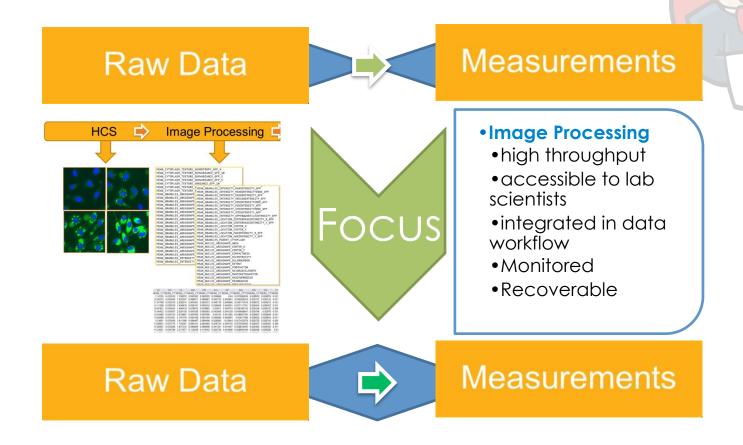  - Experiment

**Measurements**

- Raw (>500 parameters)
  - Aggregated or cell by cell
  - filtered
- Metadata
  - Image Processing

**Results**

- Assay QC
- Hit Identification
- Multi-parametric Statistics
- Correlations
- Machine Learning etc.

# HCS-High Performance Image Analysis

*Intial Focus: Remove Image Processing Bottleneck*



- **Image Processing**
  - high throughput
  - accessible to lab scientists
  - integrated in data workflow
  - Monitored
  - Recoverable

# HCS: Image Measurements and Analytics
*Easily Accessible, High Performance Image Analytics*

- Vision
  - Image and data analysis using high performance (HP) image processing tools
    - Accessible, scalable, affordable, flexible and well-supported
- Strategy
  - Evaluate and adopt open-source, community supported tools
    - CellProfiler, ImageJ, Jenkins-CI
  - Utilize NIBR-IT systems and resources
    - Linux Compute Engine (cluster) / Network Attached Storage
    - Development expertise (UI, data management and web-services)
  - Increase usability of NIBR-IT systems and resources
  - Engage and provide timely and practical functionality to both expert and casual imaging platform users
- Tactics
  - Develop functional prototypes (Jenkins-CellProfiler, Test Mosaic, R-Analytics)
  - Collaborate to develop new image/data analysis systems
  - Explore imaging tools and data space. Define HP image processing requirements.
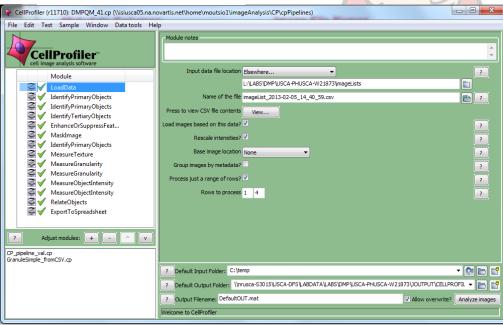  - Provide training, support and engage in community building

# CellProfiler



- Open Source Image Processing

- Platform independent

- Desktop client for defining an arbitrarily complex image processing pipeline

- Pipeline can be used by the command line CellProfiler executable

  – Suitable for high throughput analysis

  – Suitable for deployment on a Linux grid engine

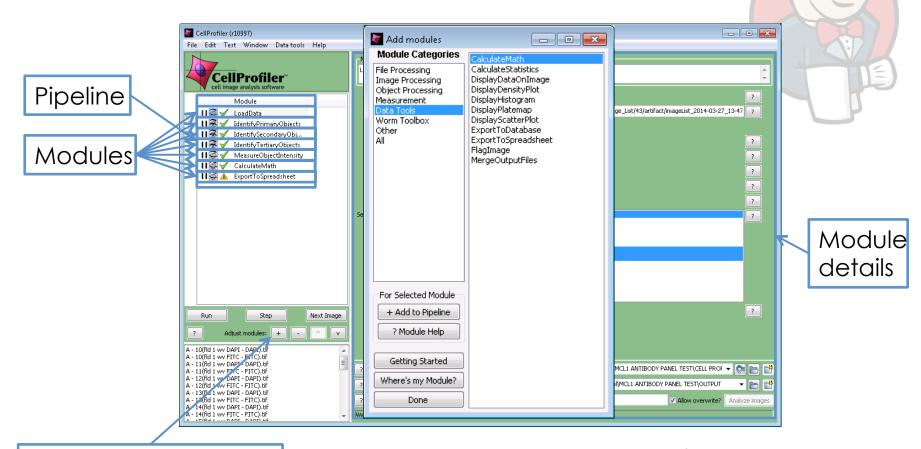  – Can process large image sets (300K + images)

  – Developed and supported by the Broad Institute and a sizable scientific user community

  – Supports additional imaging tools (ImageJ)

# CellProfiler – general anatomy
*Nuclear Translocation Assay*



Pipeline

Modules

Module details

Add/Subtract Modules

# HCS Image/Data Processing
*Programming and Prototyping Functional Requirements*

- Scripting
  - Pros
    - Quick prototyping
    - Flexibility
    - Platform independence
  - Cons
    - Unsuitable for end users
      - Requires installation of scripting tools
      - Command Line Interface

- Scripting for end users
  - Requires a user interface
    - Most UI prototypes are either
      - » hard
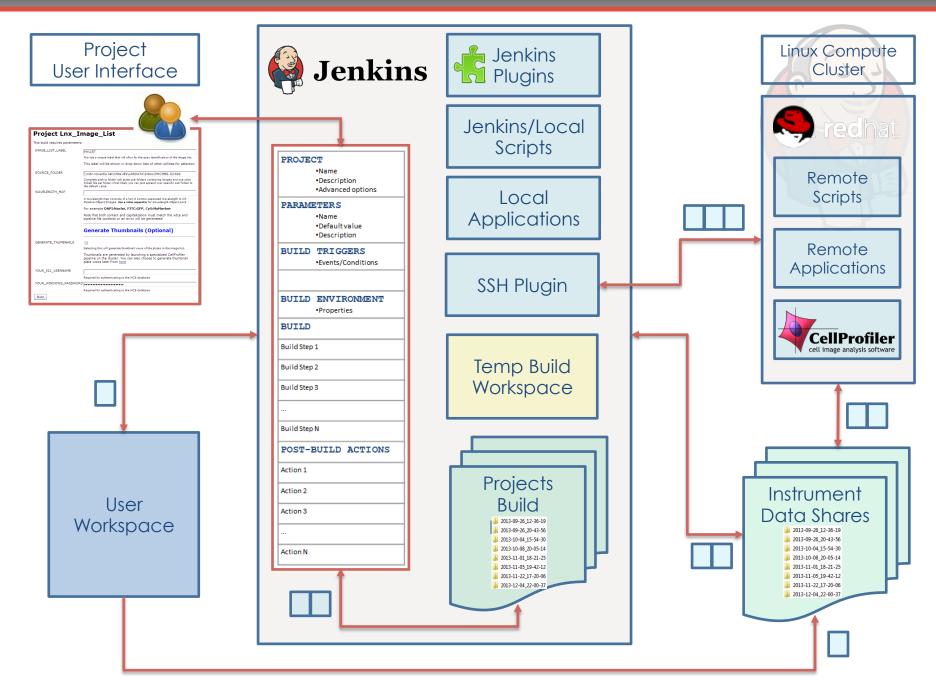      - » Pretty but not functional
      - » Or...

An extendable open source continuous integration server

- ...not very pretty
  - But quite functional

# Why choose Jenkins-CI?

- Why Jenkins-CI?
  - Jenkins allows us to rapidly wrap any command line script or program in a web interface
    - Excellent support for Groovy a Java based, dynamic, modern scripting language
    - Straight forward integration with other languages, tools, OS, frameworks
  - Jenkins has broad community support that provides access to over 800 plugins
    - Plugins allow easy customization of Jenkins for a variety of tasks
  - Jenkins provides basic workflow and web server functionality
    - Which works well in combination with CellProfiler
  - Jenkins is used extensively by the NIBR-IT group to build all kinds of internal software
    - Many software developers know a lot about Jenkins
  - Jenkins is now emerging as a useful Bioinformatics tool
    - The BioUno project

# Jenkins-CellProfiler
*HP Image Processing Workflow: Outline*



| Contribute Image Processing Pipeline | → | Assemble Images, Metadata | → | CellProfiler HP Image Processing | → | Retrieve & Use Data |
|---|---|---|---|---|---|---|

**Upload your CellProfiler pipeline**

Upload a CellProfiler image processing pipeline from y...

The pipeline will be available for use on the Jenkins-H...

Contribute Pipeline

**Generate a CellProfiler image list**

The source folder must contain one or more subfolders

Generate Linux Image List Now
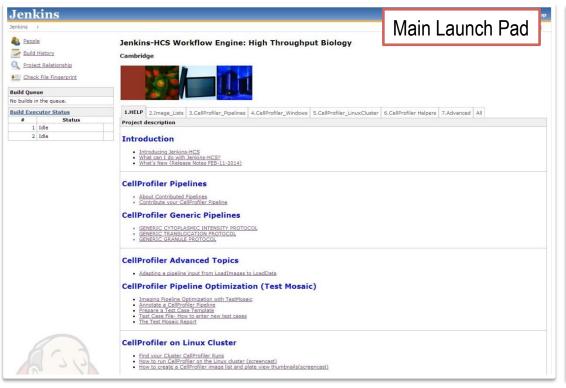
**Process images using CellProfiler**

The CellProfiler pipeline and the image list are selected

Processing can be **restricted** to a subset of the plates

Help Screencast

CellProfiler
cell image analysis software

Execute CellProfile...

# Jenkins-HCS Workflow Engine
*High Level UI Components*



Main Launch Pad

Project Launch Pad

Data Pipeline Visualization

# Typical Workflow
*Step 1: Contribute a pipeline*

- Project: Contribute_Pipeline
  - Upload and annotate a standard CellProfiler image analysis pipeline. Uploaded pipelines are usable in other projects
  - Assumptions
    - The pipeline has been designed and successfully tested on the CellProfiler desktop client
  - Outcome
    - The CellProfiler pipeline file will be uploaded and stored on Jenkins
    - Additional annotation will be extracted and attached to the pipeline

## Project Contribute_Pipeline

This build requires parameters:

Like ☐

### Pipeline Information

PIPELINE_LABEL    MY CONTRIBUTED PIPELINE

**[Required]**
A label that can be used as a pipeline identifier.

PIPELINE_DESCRIPTION
A short description of what the pipeline does or how it is used.

AUTHOR    Aaron ▼
Select the name of the contributing pipeline author

AUTHOR_COMMENTS

Provide any additional comments that may be useful to others in using this pipeline. H

### Pipeline File

user_pipeline.cp    Choose File  No file chosen

**[Required]**
Choose a working CellProfiler image analysis pipeline to share

example_image_list.csv    Choose File  No file chosen
An example image list for this pipeline.

This is only required if the uploaded pipeline image list is not accessible fro

An image list is not required if the pipeline uses an image list from a URL.

Build

# Build report from a contributed pipeline
*Uses: Summary Display Plugin*



- Contributed pipelines are annotated by a combination of user provided and auto-extracted metadata

  - Presented as a tab panel

  - Pipeline can be downloaded and further modified

*Use PIPELINE FILE tab to download or quickly browse the pipeline*

# A CellProfiler Pipeline from Jenkins Server
*Additional Usage*

- # CellProfiler pipelines on the Jenkins server can be used as follows:

  - For inspection

  - For re-use
    - On CellProfiler desktop client
    - On Jenkins-CellProfiler

  - For further experimentation
    - Load in desktop client and further customize



**Build GENERIC GRANULE PROTOCOL**

| Parameters | |
|---|---|
| user_pipeline.cp | GENERIC GRANULE PROTOCOL.cp view |
| PIPELINE_LABEL | GENERIC GRANULE PROTOCOL |

Download

Inspect

Copy URL

CellProfiler (r11710)

File   Edit   Test   Window   Data tools   Help

Load Pipeline...                    ctrl+O

Load Pipeline from URL

Save Pipeline              ctrl+shift+S

Save Pipeline as...

Customize

Desktop

2) Select a processing pipeline:

PIPELINE   CellProfiler_Pipeline GENERIC GRANULE PROTOCOL

CellProfiler_Pipeline GENERIC GRANULE PROTOCOL
CellProfiler_Pipeline GENERIC TRANSLOCATION PROTOCOL
CellProfiler_Pipeline GENERIC CELL INTENSITY PROTOCOL
CellProfiler_Pipeline YAP Test
CellProfiler_Pipeline pAKT Cell Panel
CellProfiler_Pipeline HUH-1 pS6
CellProfiler_Pipeline SNU878 P62 Endogenous
CellProfiler_Pipeline HUH-1 P62 Endogenous
CellProfiler_Pipeline KDM4_ImageStats
CellProfiler_Pipeline DMPQM_33
CellProfiler_Pipeline DMPQM_14
CellProfiler_Pipeline DMPQM_25
CellProfiler_Pipeline DMPQM_40

GENERATED_IMAGE_LIST

QUERY_KEY

Groups are composed from any combination of the following Keys: **Barcode, RowNum**

- Key names are case sensitive

Jenkins-CP

# Typical Workflow
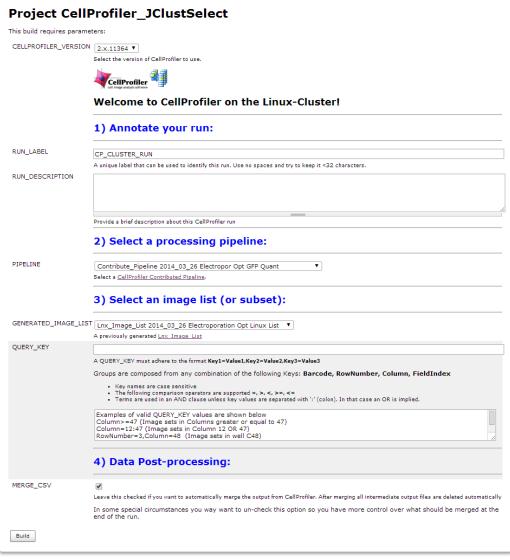## *Step 3: Execute CellProfiler on the Linux Cluster*

- Project: CellProfiler_JClustSelect
  - Executes a series of image processing steps using the Jenkins-CI CellProfiler
  - Uses the SSH Plugin
  - Typical Assumptions
    - CellProfiler pipeline and a CP formatted image list are stored on the Jenkins server
      - Jenkins build artifacts
  - Outcome
    - Summary report
    - A file containing combined measurements from all the images processed.
      - Results file is in CSV format

**Project CellProfiler_JClustSelect**

This build requires parameters:

CELLPROFILER_VERSION   [2.x.11364 ▼]
Select the version of CellProfiler to use.

**CellProfiler**
cell image analysis software

**Welcome to CellProfiler on the Linux-Cluster!**

**1) Annotate your run:**

RUN_LABEL   [CP_CLUSTER_RUN]
A unique label that can be used to identify this run. Use no spaces and try to keep it <32 characters.

RUN_DESCRIPTION
Provide a brief description about this CellProfiler run

**2) Select a processing pipeline:**

PIPELINE   [Contribute_Pipeline 2014_03_26 Electropor Opt GFP Quant ▼]
Select a CellProfiler Contributed Pipeline

**3) Select an image list (or subset):**

GENERATED_IMAGE_LIST   [Lnx_Image_List 2014_03_26 Electroporation Opt Linux List ▼]
A previously generated Lnx_Image_List

QUERY_KEY
A QUERY_KEY must adhere to the format **Key1=Value1,Key2=Value2,Key3=Value3**

Groups are composed from any combination of the following Keys: **Barcode, RowNumber, Column, FieldIndex**

- Key names are case sensitive
- The following comparison operators are supported =, >, <, >=, <=
- Terms are used in an AND clause unless key values are separated with ':' (colon). In that case an OR is implied.

Examples of valid QUERY_KEY values are shown below
Column>=47 (Image sets in Columns greater or equal to 47)
Column=12:47 (Image sets in Column 12 OR 47)
RowNumber=3,Column=48 (Image sets in well C48)

**4) Data Post-processing:**

MERGE_CSV   [✓]
Leave this checked if you want to automatically merge the output from CellProfiler. After merging all intermediate output files are deleted automatically

In some special circumstances you way want to un-check this option so you have more control over what should be merged at the end of the run.

[Build]

# Monitoring CellProfiler runs on the cluster
*Uses: [Build Pipeline Plugin](#)*

> Users switch to the graphical review of the workflow!
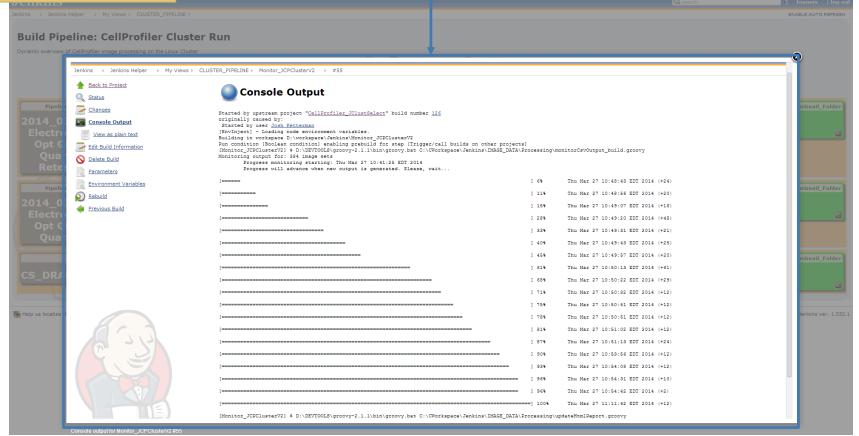


Schedule, run and monitor CellProfiler jobs executed on the Cambridge cluster

| 1.HELP | 2.Image_Lists | 3.CellProfiler_Pipelines | 4.CellProfiler_Windows | **5.CellProfiler_LinuxCluster** |

**Name** ↓ | **Project description**

_5_TabGuide

## CellProfiler on Linux Cluster: Helpful Links

### Multi-Stage Pipelines

_5_TabGuide

- Graphical Review of Cluster Pipeline
- Graphical Review of Image List Generation

- Graphical Review Yokogawa Cluster Pipeline
- Graphical Review Yokogawa Image List Generation

Jenkins › Jenkins Helper › My Views › CLUSTER_PIPELINE ›                                                 ENABLE AUTO REFRESH

**Build Pipeline: CellProfiler Cluster Run**

Dynamic overview of CellProfiler image processing through the cluster

# Monitoring CellProfiler runs on the cluster

*Uses: Build Pipeline Plugin and the Console*

# Run Report & Measurement Retrieval
*Uses: Associated Files and HTML Publisher plugins*

# Jenkins –CI: CellProfiler Image Processing
*Uses:* [HTML Publisher]() *plugin*

# Advanced/Experimental Functionality
*Exploring the parameter space (a.k.a. Test Mosaic)*

- Optimization of imaging module parameters

  – A typical pipeline development requirement

- Test Mosaic

  – Allows systematic and documented exploration of the parameter space

  – Evaluation is based on visual and quantitative interpretation of the results

# HCS-Multi-Parametric Data Analysis

*Current Focus:* *Prototype powerful and easy to use analytics*



Contributed by
Stanley Lazic

# Statistics, Visualization, Reporting

*My current Jenkins toolkit*

- ## Jenkins R-Plugin

  - Supplies build step for executing R scripts

    - This plug-in was created by the BioUno project (sponsored by TupiLabs), and released to Jenkins as well.

- ## Image Gallery Plugin

  - This plug-in reads a job workspace and collects images to produce an image gallery

  - Useful for visualizing various statistical plots and graphs

    - This plug-in was created by the BioUno project (sponsored by TupiLabs), and released to Jenkins as well.

- Reporting Plugins

  - HTML Publisher, Summary Display

# Jenkins for Interactive Analytics
*Using R in a Jenkins pipeline interactively*

- Opportunities

  - Quickly prototype functional analysis for multi-parametric data

    - Improve analysis requirements
    - Experiment with required data management and analysis workflows

  - Provide lab scientists with an easy to use, yet sophisticated, standardized and validated platform for MP data analysis tools

# Jenkins for Interactive Analytics
*Using R in a Jenkins pipeline interactively*

- Challenges

  – Limitations of the Jenkins user interface

    • Limited interaction between UI controls

  – Large and varied HC measurement metadata

    • A challenge for creating HC data schemata as well

- Strategies

    • Open source collaboration with BioUno project

      – Uno-choice UI control greatly facilitates dynamic updating of the UI

    • Initial design supports flexible (but still controlled) data schema

      – Low tech, cumulative, shared key-value Java properties

# Analytical Builds

*A build may create a new transform of the data or simply add metadata*



**CSV-01**

Jenkins output
or external

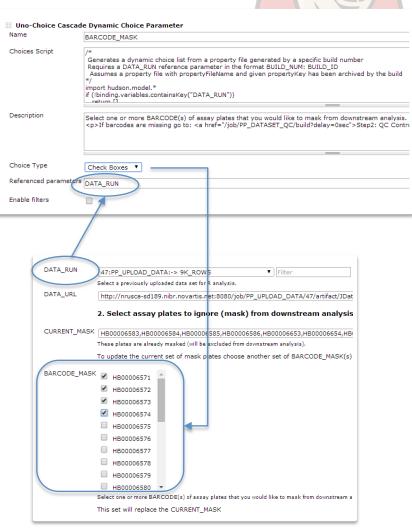| Name | Project description ↓ |
|---|---|
| PP_UPLOAD _DATA | **Step1: Upload Data:-: Step2: QC Control W** This project will upload a data set for analysis and w |
| PP_DATASET _QC | **Step2: QC Control Wells Step3: Review res** This project illustrates the spread and differences of QC Dataset Now |
| PP_DATASET _QCHEATMAP | **Step3.0: Review response across plates St** This project generates comparative heatmaps for se |
| PP_DATASET _HISTOGRAM | **Step3.1: Review response across plates St** This project generates histogram distributions for se |
| PP_MASK _BARCODES | **Step4: Remove bad plates Step5: Review** This project illustrates the spread and differences of Mask Asay Plates Now |
| PP_FEATURE _ANALYSIS | **Step5: Analyze data set features Step5: N** This project analyzes the predictive characteristics o |
| PP_LOESS _NORMALIZATION | **Step6: Normalize plate data** This project applies LOESS correction to normalize th |

**Data-01**

**Metadata-01**

**Metadata-02**

**Metadata-03**

.......

**Metadata-N**

.......

**Result-06**

**Metadata-01**

.......

## *Workflow Requirement*

*Ability to select ad hoc Artifacts (data, metadata, results) from previous project builds*
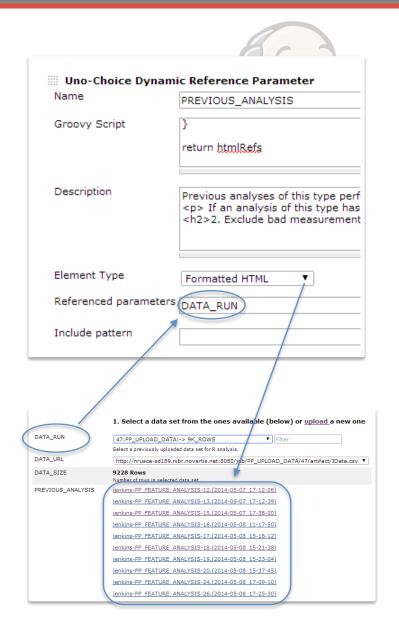
# The Uno-Choice plugin

- Provides a list of dynamically generated options
  - Driven by a Groovy script
  - Single/Multi-select (Check Boxes, Radio Buttons)
  - References one or more other UI parameters
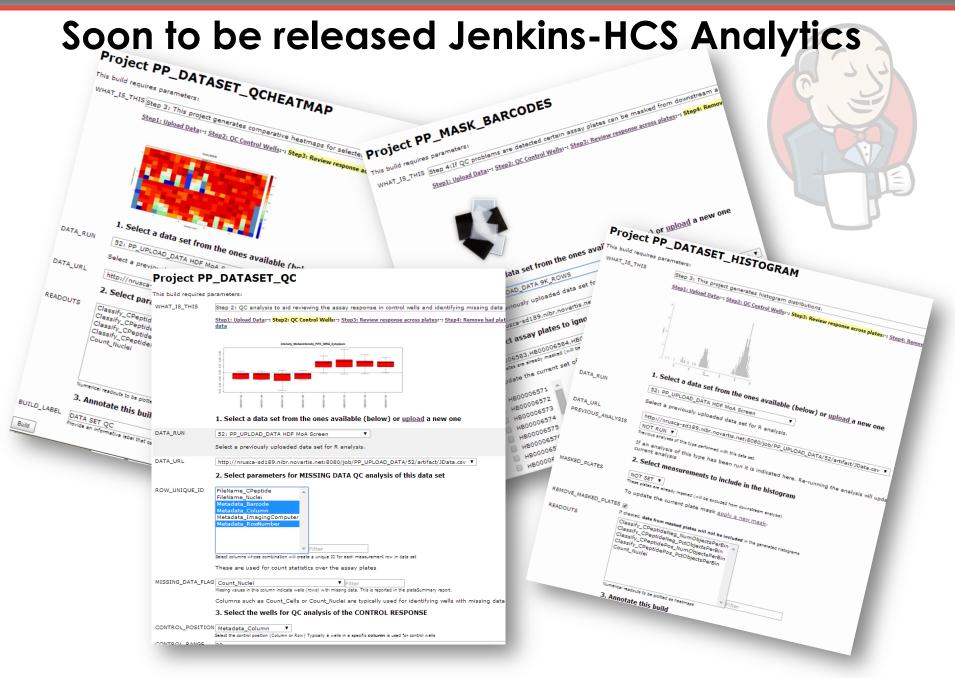  - Dynamically refreshes when referenced UI parameters change

# The Uno-Choice plugin

- Provides reference parameters
  - Dynamically rendered in the UI but not used in the build
  - Rendered as lists, 'free-form' HTML, or an image gallery

# Soon to be released Jenkins-HCS Analytics

# Introducing Jenkins to Life Sciences!

*Let's start by explaining away 'artifacts'!*

**ar·ti·fact** 🔊 [**ahr**-t*uh*-fakt] ?   Show IPA

*noun*

1. any object made by human beings, especially with a view to subsequent use.
2. a handmade object, as a tool, or the remains of one, as a shard of pottery, characteristic of an earlier time or cultural stage, especially such an object found at an archaeological excavation.
3. any mass-produced, usually inexpensive object reflecting contemporary society or popular culture: *artifacts of the pop rock generation.*
4. a substance or structure not naturally present in the matter being observed but formed by artificial means, as during preparation of a microscope slide.
5. a spurious observation or result arising from preparatory or investigative procedures.

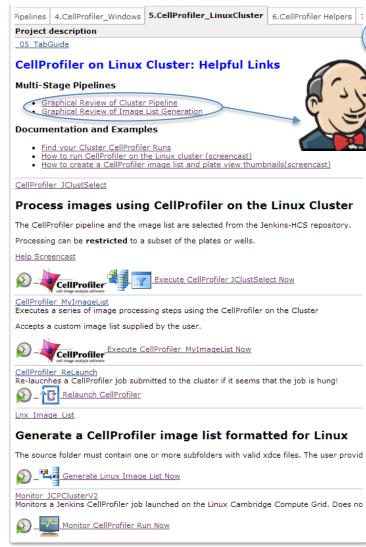*http://dictionary.reference.com/browse/artifact*

**Developer**

**Impedance Mismatch !**

**Scientist**

# Introducing Jenkins to Life Sciences

*Let's improve the User Interface/Experience*

- Let's start by improving the default Jenkins UI
  - Layout
  - Navigation
  - Refreshing
  - Interactivity
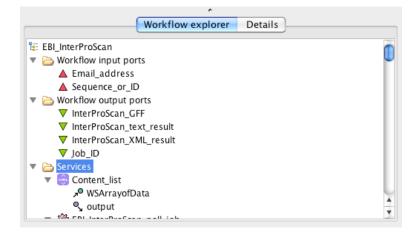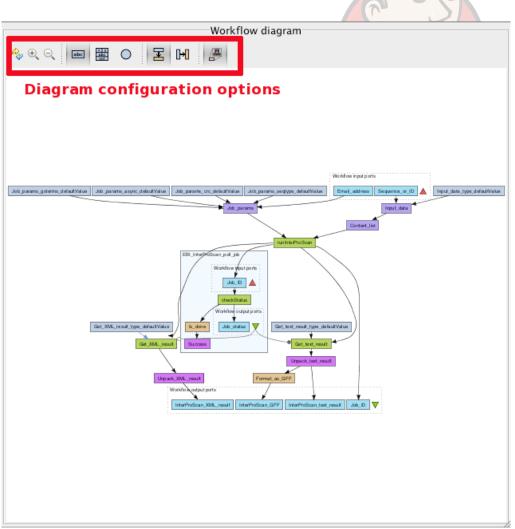- This is an active Jenkins community discussion

# What We are Missing
*Configuration Explorer*

- Structured

- Graphical

- Dynamic

# What We are Missing
## *Bi-Directional Build Interaction*

- Build C produces an intermediate report that will get updated once Build D is finished successfully.

■ Build A **uses** artifacts of Build B

- Limited support by **Run Type** parameter
  - Missing flexible and dynamic filtering

■ Build D **modifies** build/ publisher artifacts of Build C

- Sometimes not do-able

- Sometimes requires a reload

Back to CellProfiler_MyImageList | cpReport

**CellProfiler**
cell image analysis software

**2014_05_21_PGC1A_Set1Reimg**

CellProfiler-on Cluster Report: 2014-05-23_13-31-11 (build 27)

| Build Parameters | Review |
|---|---|
| CellProfiler Pipeline | Review |
| Total Source Images | 4608 |
| Image Source | D:\DEVTOOLS\Jenkins\workspace\CellProfiler_MyImageList |
| Measurements Folder | /labdata/incell/cluster_runs/CPJENKINS/JCP_2014-05-23_13-31-11 |
| Merged Data Folder | \\nibr.novartis.net\usca-dfs\LABDATA\LABS\incell\cluster_runs\JOUTPUT\CELLPROFILER\2014-05-23-13-31-11\ALL |
| Progress Monitor | Progress Monitor |

- Build D monitors output of long running job and updates report of Build C

| Progress Monitor | Progress Monitor |
|---|---|

'Progress Monitor' Link and cell color are updated

# What We are Missing

*A good, deep search and metadata framework*

- ## Supported

  - ### View Searches

  - ### Build Browsing

    - By timeline
    - By view
    - By user

- ## Missing

  - ### Build Searching

    - Parameter Search

    - Metadata Search

      – Metadata plugin (currently limited to adding metadata at project level)

    - Artifact Search

  - ### Tagging

  - ### Dynamic Metadata

# What We are Missing
*Life-Sciences Domain Plugins (Bio/Chem Informatic)*

| Phylogenetics | |
|---|---|
| MrBayes Plug-in | Integrates MrBayes and Jenkins. |
| FigTree Plug-in | Integrates FigTree and Jenkins. |
| **Genetic Analysis** | |
| Structure Plug-in | Integrates Structure and Jenkins. |
| Structure Harvester Plug-in | Integrates Structure Harvester an |
| CLUMPP Plug-in | Integrates CLUMPP and Jenkins |
| Distruct Plug-in | Integrates Distruct and Jenkins. |

- The BioUno project is filling the gap

- Interested in plugins that

  - Integrate bio-informatic, statistical and visualization tools

  - Connect to life-science data repositories

  - Generate artifacts and reports in LifeSci formats

| UI | |
|---|---|
| Uno-Choice Plug-in | A proposal for a new Jenkins UI plugin for selecting one or multiple parameters. Attempting to fill the gaps left by current plugin options. |
| Image Gallery Plug-in | This plug-in reads a job workspace and collects images to produce an image gallery using colorbox lightbox Javascript library. |
| **Misc** | |
| R Plug-in | A simple plug-in to invoke R interpreter and execute an R script. |

# In Summary

- We have demonstrated that Jenkins-CI can be used for life-science applications
  - Using standard functionality
  - Using domain specific plugins
  - In demanding environments of big data and high performance
- We have observed that scientist are able and willing to use the platform despite it's 'domain impedance mismatch'
- There is some fundamental interest in the larger Jenkins-CI community to expand the boundaries of the framework beyond continuous integration

# Where do we want to take Jenkins-CI?

- Discussion
  - No changes?
  - Gradual improvements?
    - User interface
    - API
    - New life-science plugins
  - Fundamental changes?
  - Integration framework for orchestrating more granular pipelines?
    - CellProfiler
    - Galaxy
    - Knime
    - Others?

# Acknowledgments

- Novartis
  - Fred Harbinski
  - Christian Parker
  - Stanley Lazic
  - Imtiaz Hossain
  - Josh Snyder
  - Erik Sassaman

- BioUno
  - Bruno Kinoshita

- The Jenkins Community

# Thank You To Our Sponsors