

How do I automate my job pipeline to ensure reproducibility?

Ilkay Altintas

altintas@sdsc.edu

Welcome to Day 2!

- What have you learned so far?
 - Launching and managing jobs
 - Managing data on the file system
- What if you wanted to combine these steps?



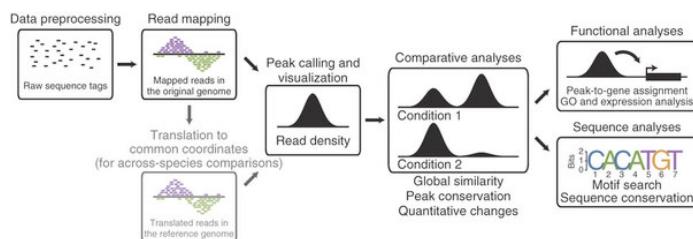
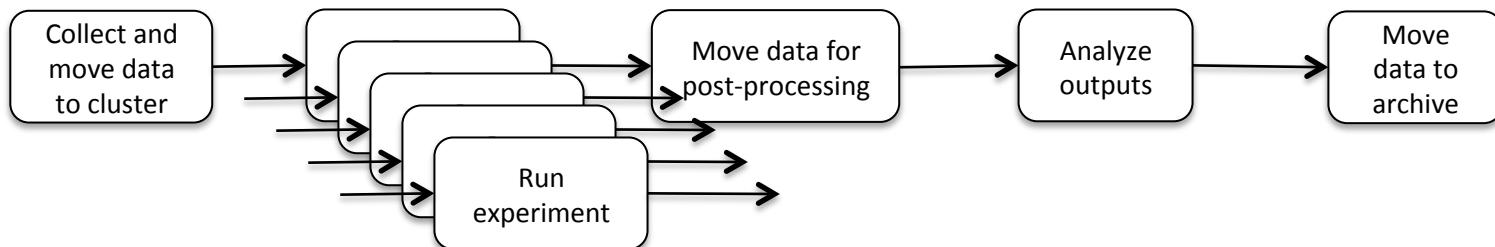


What is a Pipeline?

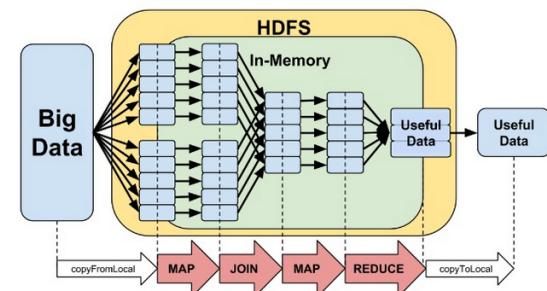
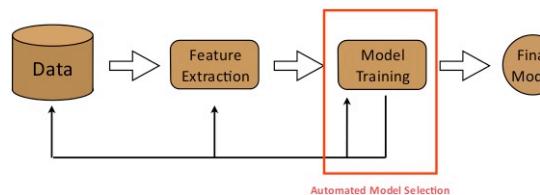
cat  sort

A UNIX pipe provides one-way communication
between two processes on the same computer

Configure
software

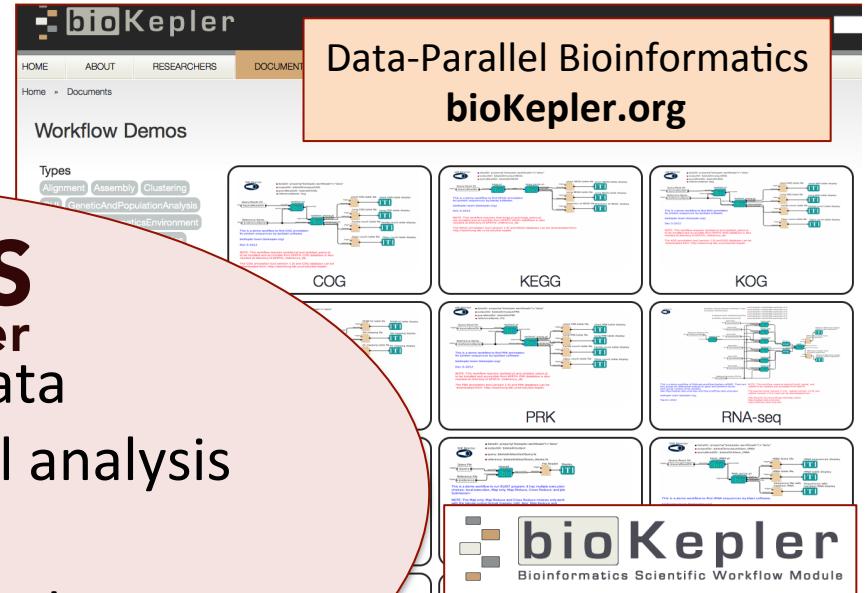
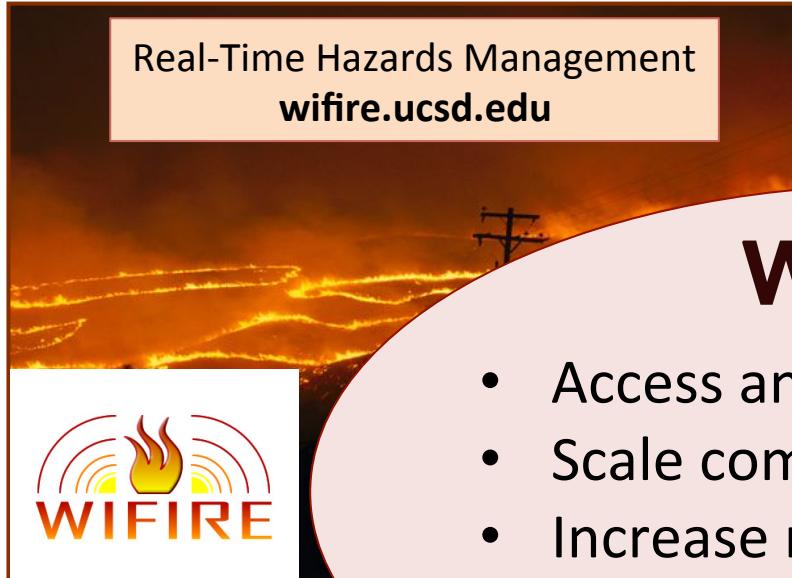


A Standard ML Pipeline



Data Science Workflows

- Programmable and Reproducible Scalability -

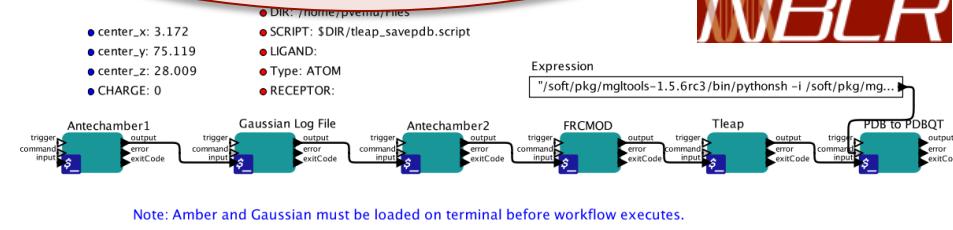


- Access and query data
- Scale computational analysis
- Increase reuse
- Save time, energy and money
- Formalize and standardize

kepler-project.org



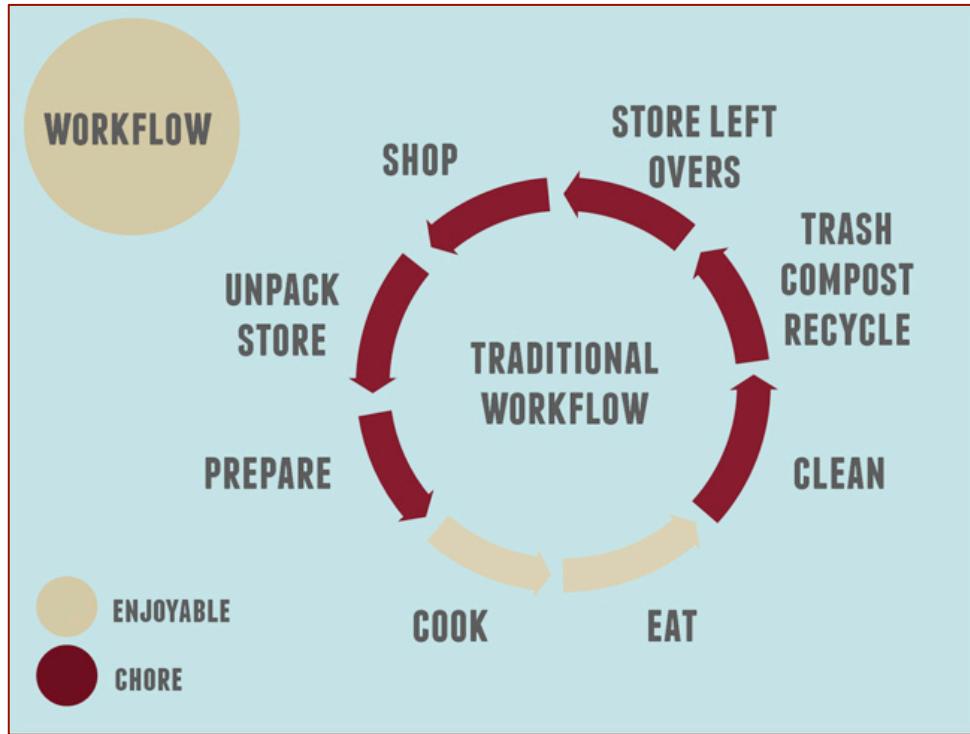
Scalable Automated Molecular Dynamics and Drug Discovery
nbcr.ucsd.edu



WorDS.sdsc.edu



So, what is a workflow?



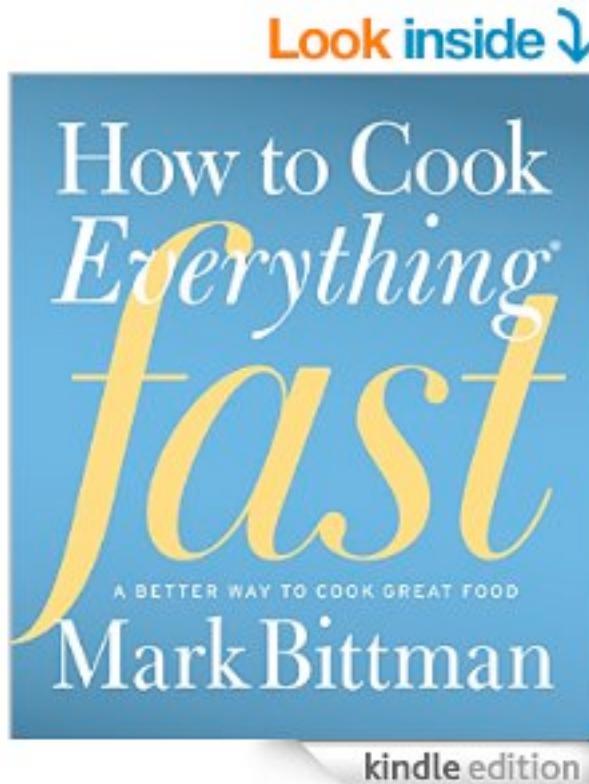
Source:
<http://www.fastcodesign.com/1663557/how-a-kitchen-design-could-make-it-easier-to-bond-with-neighbors>



Let's make pasta this evening!



How to Cook Everything Fast



Look inside ↓

“How to Cook Everything Fast is a book of kitchen innovations. **Time management**—the essential principle of fast cooking—is woven into revolutionary recipes that do the thinking for you. You’ll learn how **to take advantage of downtime to prepare vegetables while a soup simmers or toast croutons while whisking a dressing**. Just cook as you read—and let the recipes guide you quickly and easily toward a delicious result.”

Image and quote source: [amazon.com](https://www.amazon.com)

What if you have more than one cooks?





MAP

- **Input:** veggies
- **User defined function(UDF):** chop
- **Output:** Chopped groups of each kind of veggie



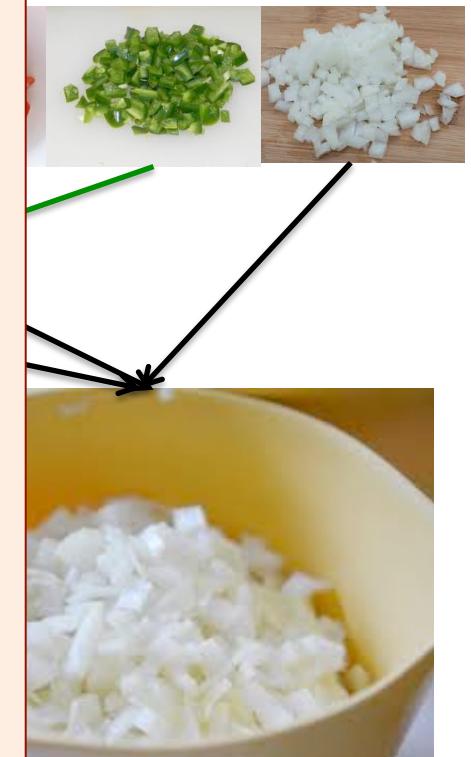
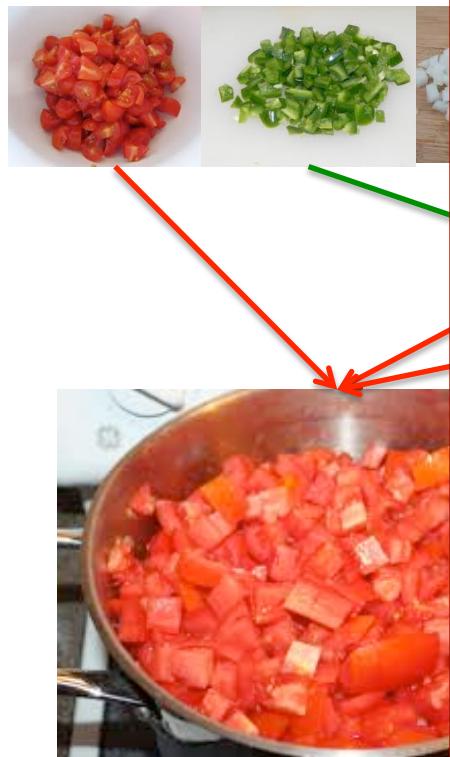
...





REDUCE

- **Input:** chopped batches for each veggie type
- **User defined function(UDF):** combine based on veggie type as key
- **Output:** a bowl of veggies per veggie kind



Thanksgiving dinner preparation: more planning and tasks?



PHOTOS COURTESY OF LOL FOODIE, KING ARTHUR FLOUR, SAVORY SWEET LIFE, MY RECIPES, FOOD NETWORK, WHAT WE'RE EATING AND FOODIE TOTS

Menu Item	Preparation Time	Cooking Time	Cooling Time
Turkey	30 minutes	4 hours	15 minutes
Veggies	30 minutes	45 minutes	None
Cranberry Sauce	5 minutes	30 minutes	2 hours
Soup	20 minutes	30 minutes	None
Pie	30 minutes	5 minutes	1 day



- When do you start cooking?
- What order do you cook?
- Can you cook some menu items in parallel?
- Who cooks what?
- ...

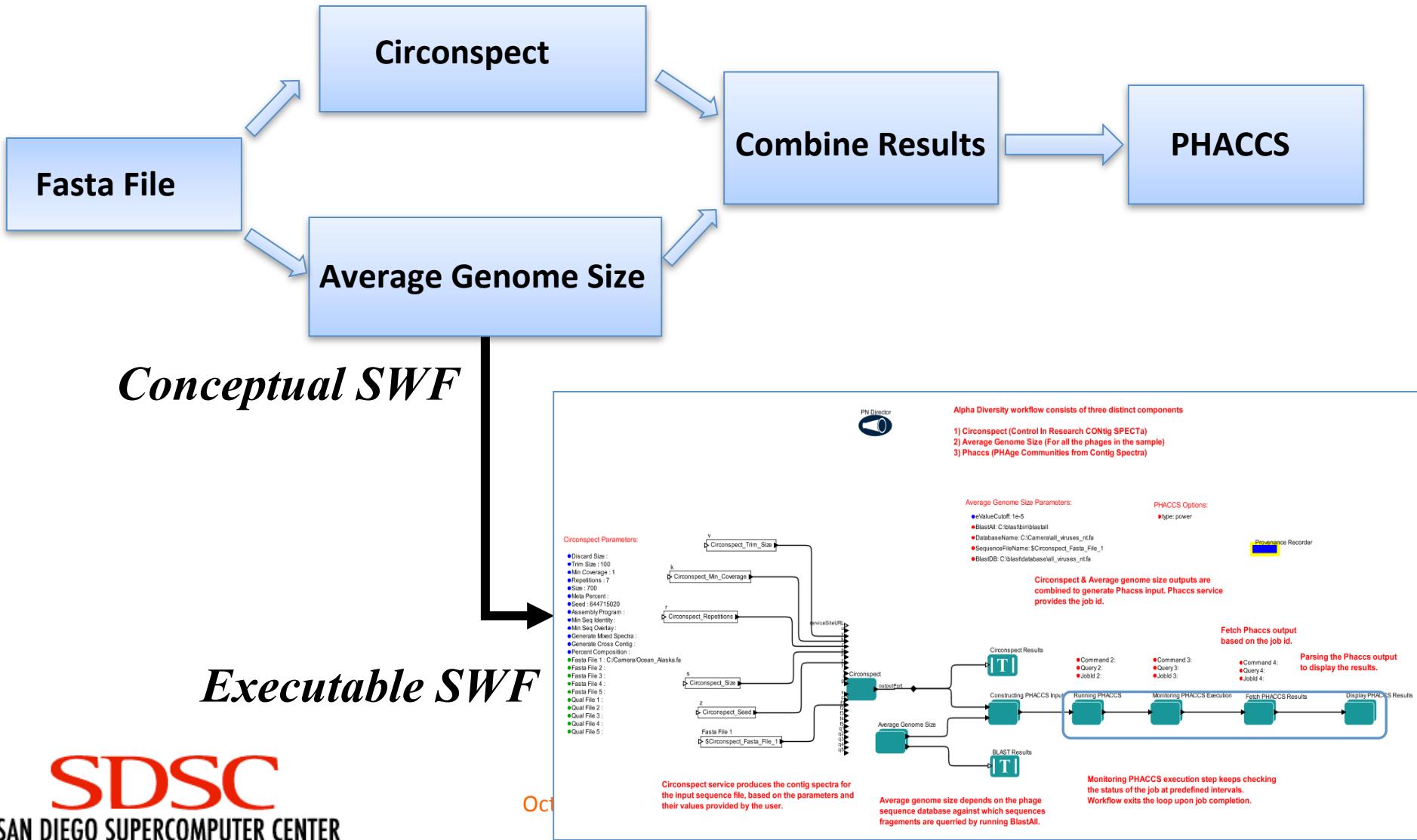
What is a scientific workflow?

Scientific workflows emerged as an answer to the need to **combine multiple data and computing components** in automated process networks.

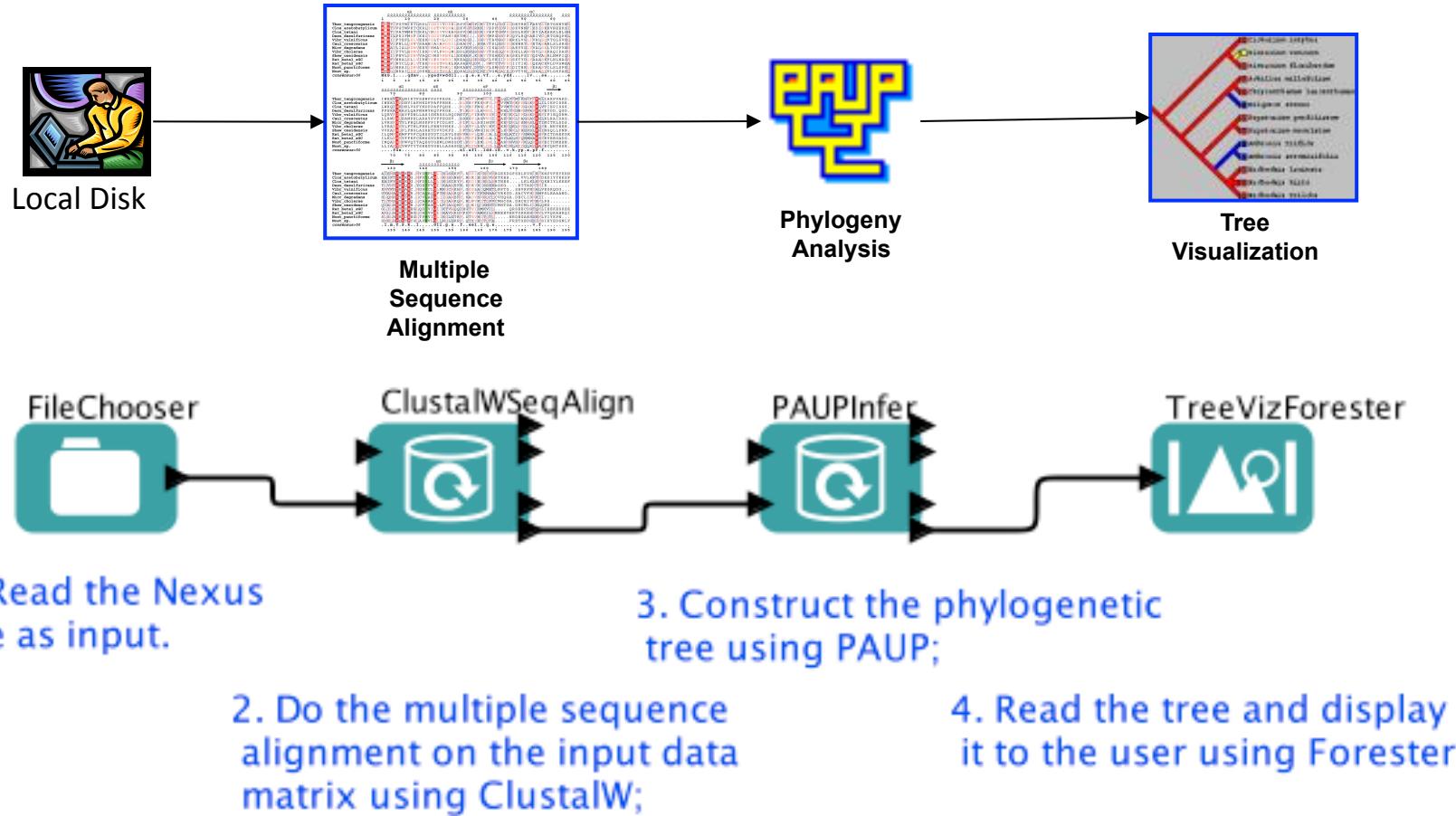
A scientific workflow is a **series of computational steps** that scientists use to generate results. That may involve accessing multiple applications and databases, and processing the data using computationally intensive jobs on high-performance clusters.

The Big Picture is Supporting the Scientist

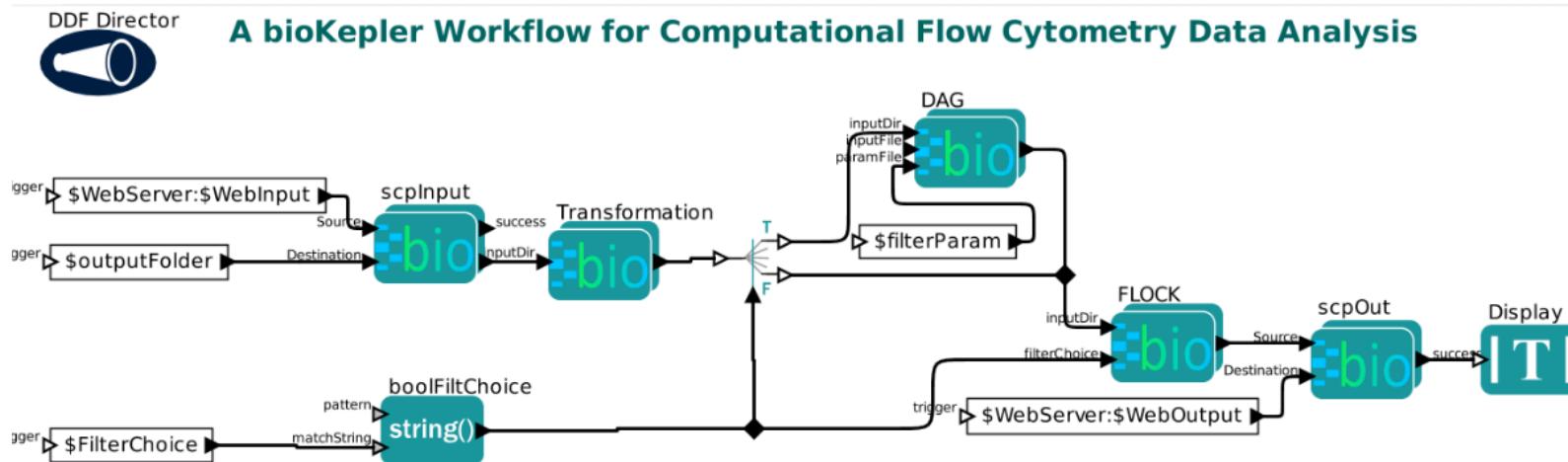
From “*Napkin Drawings*” to *Executable Workflows*



They can be simple pipelines.



They can integrate components to analyze local data.



This workflow implements a computational procedure for automated identification of cell populations from multidimensional flow cytometry data with the FLOCK algorithm. It contains three major steps including Transformation (FCSTrans), Filtering (DAG), and Clustering (FLOCK). The workflow transfers user input data from front-end web portal to the backend workflow cloud virtual machine or to the cluster, carries computational analysis in background for multiple input files on Cloud VM or the SDSC Gordon Supercomputer, and uploads result back to the web portal. The workflow configures or by-pass filter stage depending user filter choice. The Dynamic Dataflow (DDF) director manages workflow execution.

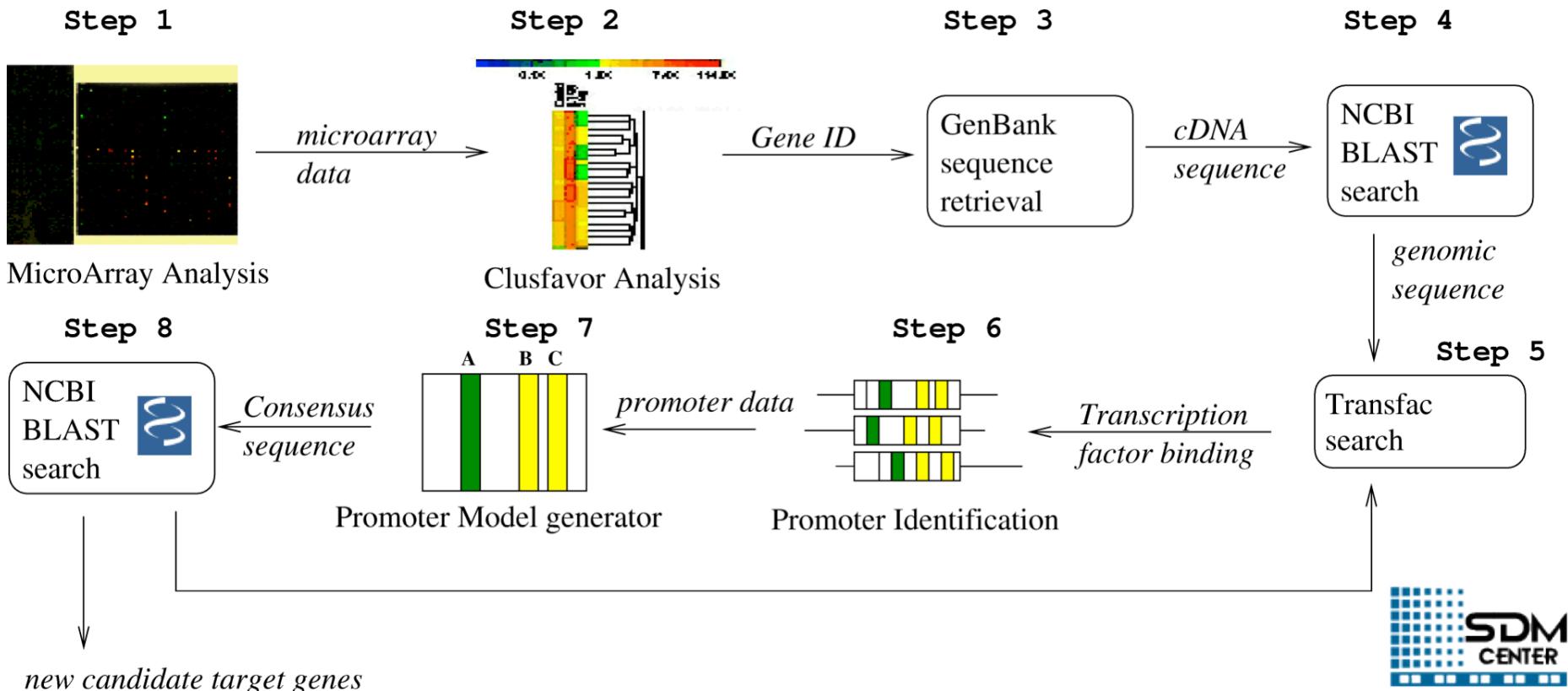
- FilterChoice: True
- filterParam: \$toolDir/b_filter_popDef.csv
- outputFolder: /home/spurawat/JCVI/CWDflowER
- toolDir: /home/spurawat/JCVI/CDUCModExec
- WebInput: /export/CDUCWebInput/User1_InputTest.z
- WebOutput: /export/CDUCWebOutput
- WebServer: ubuntu@132.249.230.30

See <http://flowgate.jcvi.org>

Authors: Shweta Purawat, Jianwu Wang, Ilkay Altintas, Robert Sinkovits @ SDSC
Yu Qian, Rick Stanton, Hyunsoo Kim, Richard Scheuermann @ JCVI

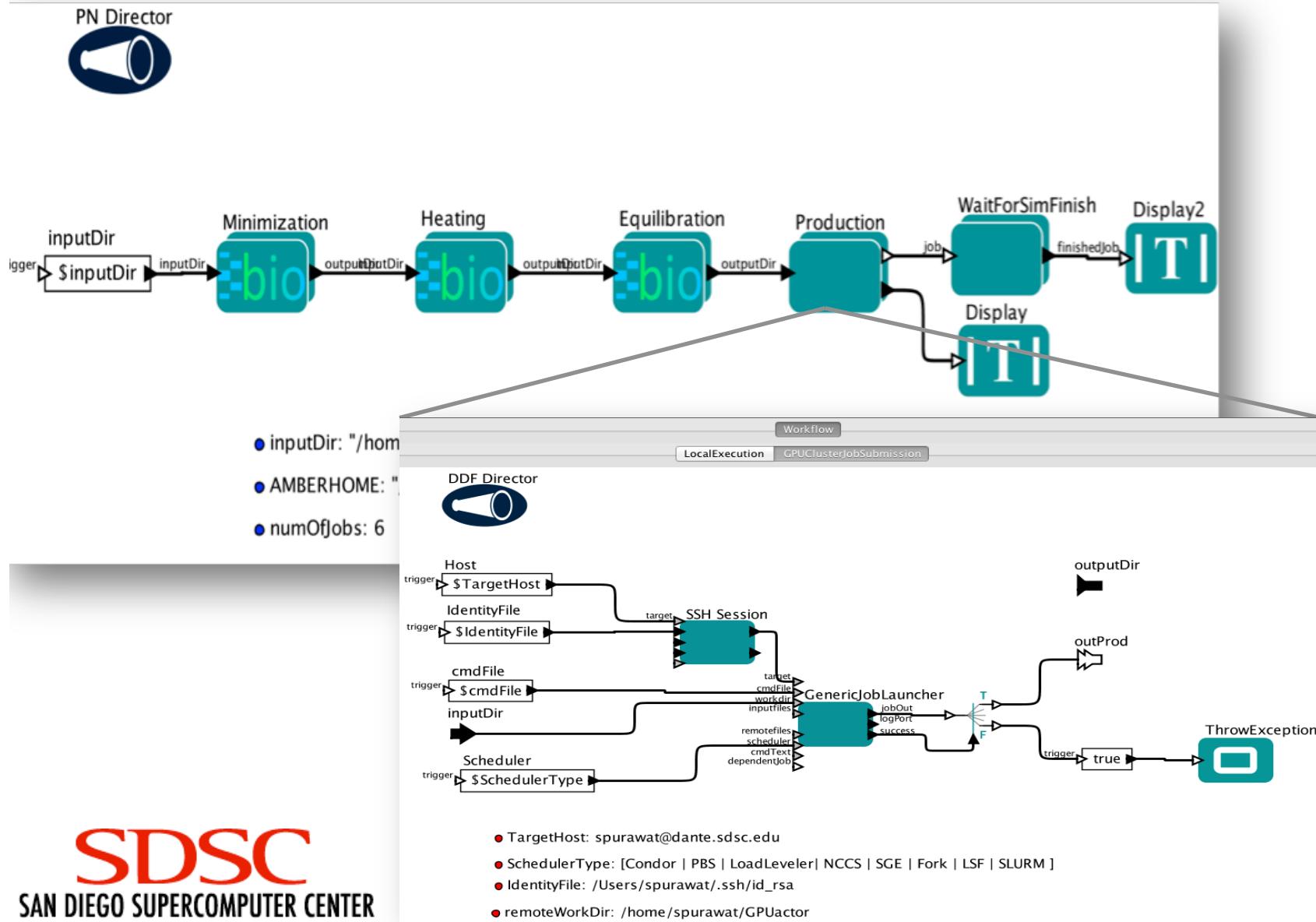
April-25-2014

They can integrate online tools and databases using complex programming structures.

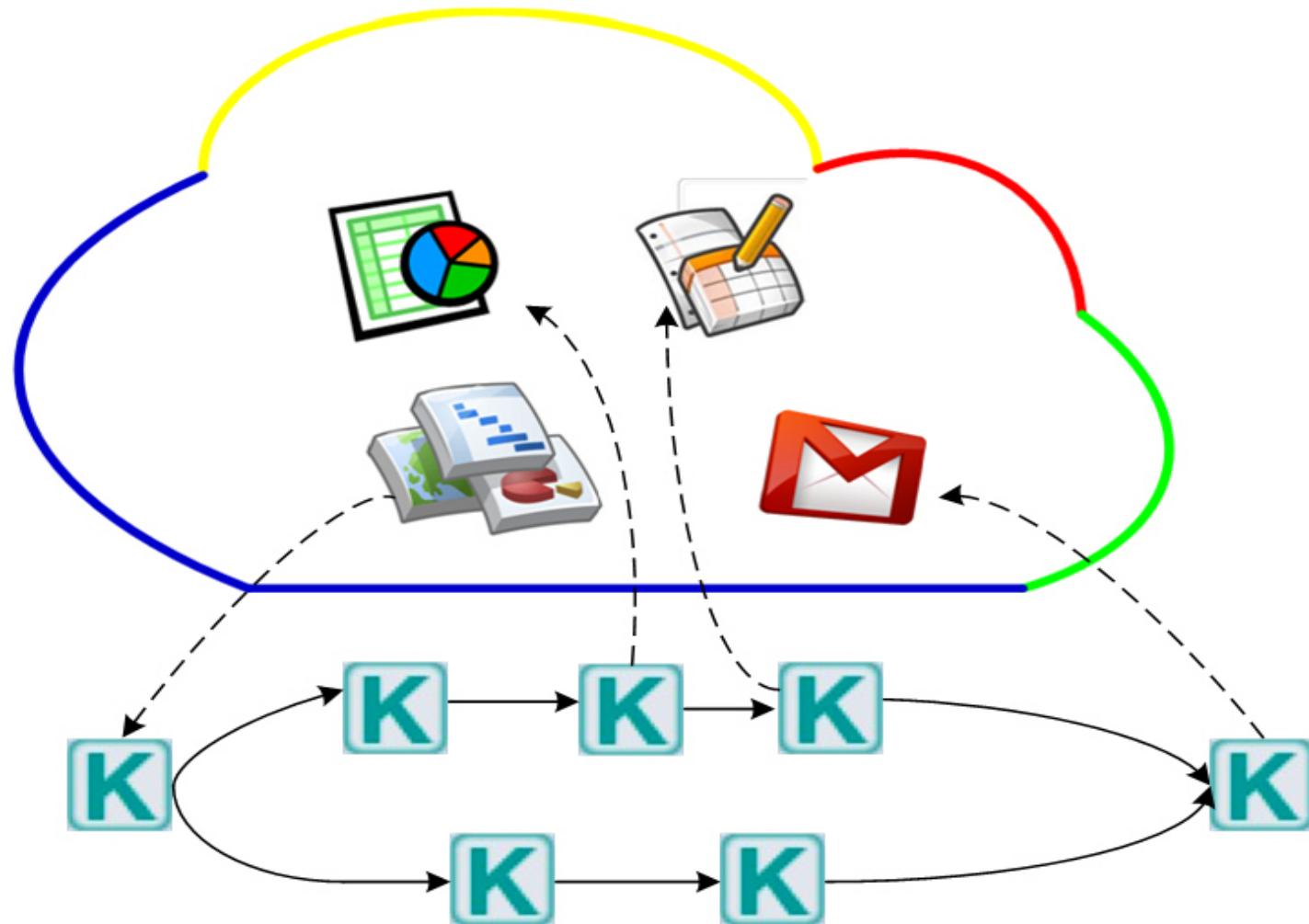


Source: Matt Coleman (LLNL)

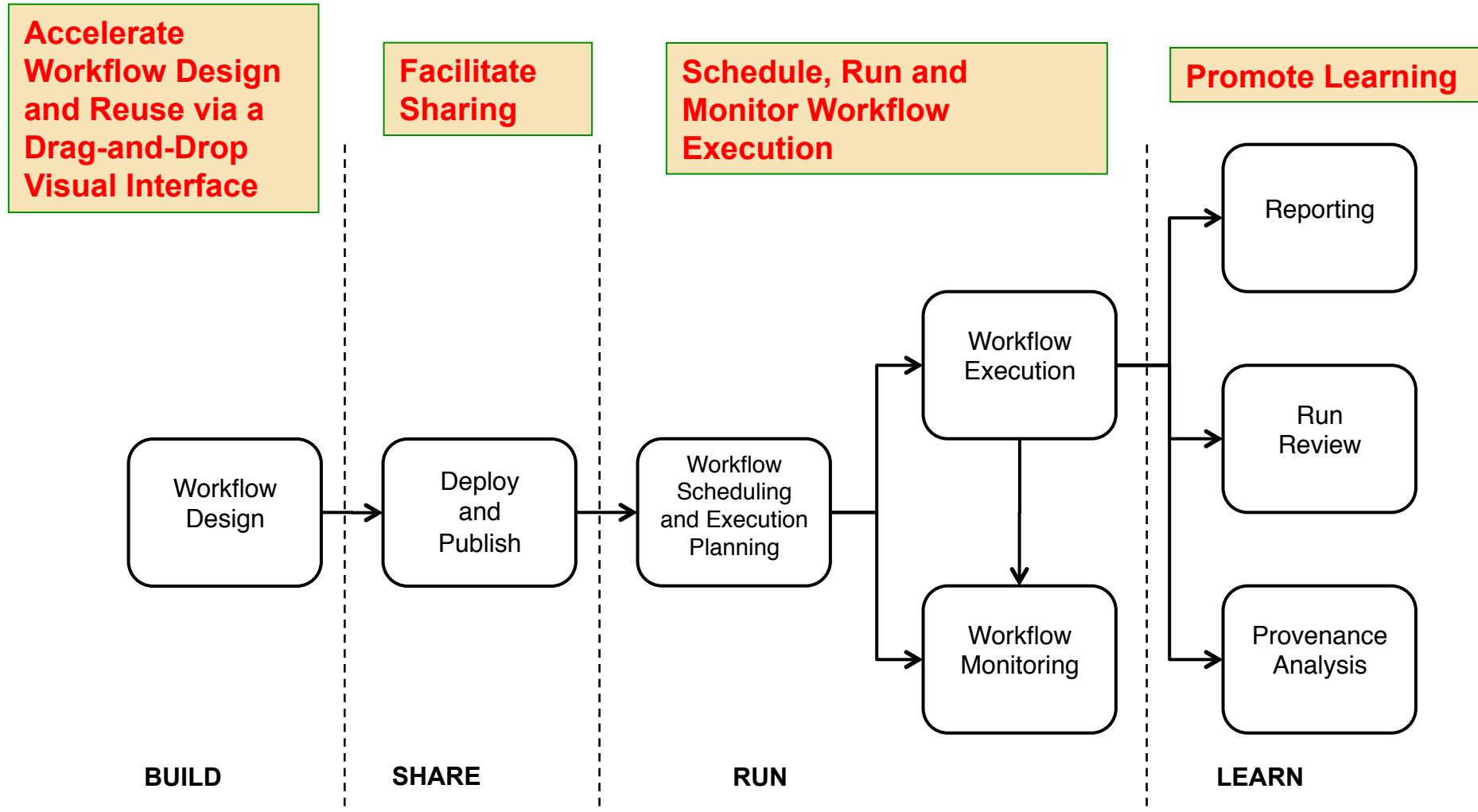
They can have distributed concurrent analysis and data management steps.



They can integrate Cloud resources.



Workflows are a Part of Cyberinfrastructure



Support for end-to-end computational scientific process

Requirements are similar for many domains

-- with slight variations --

Facilitating and Accelerating XXX-Info or Comp-XXX Research using Scientific Workflows

- Important Attributes



Assemble complex processing easily



Access transparently to diverse resources



Incorporate multiple software tools



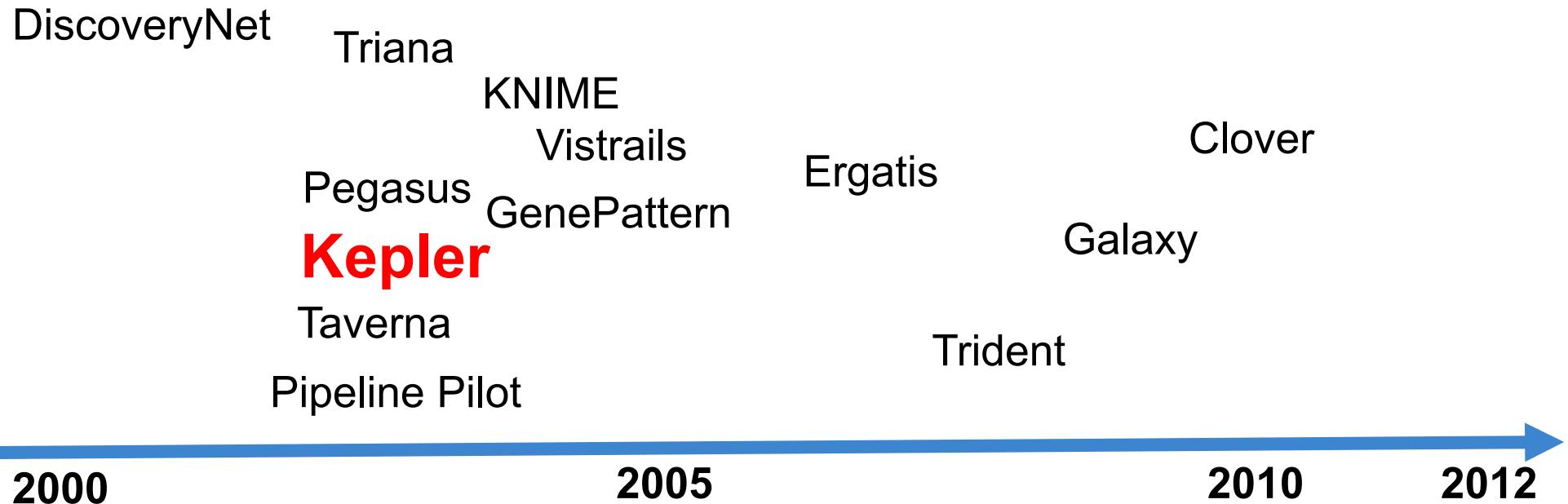
Assure reproducibility



Build around community development model



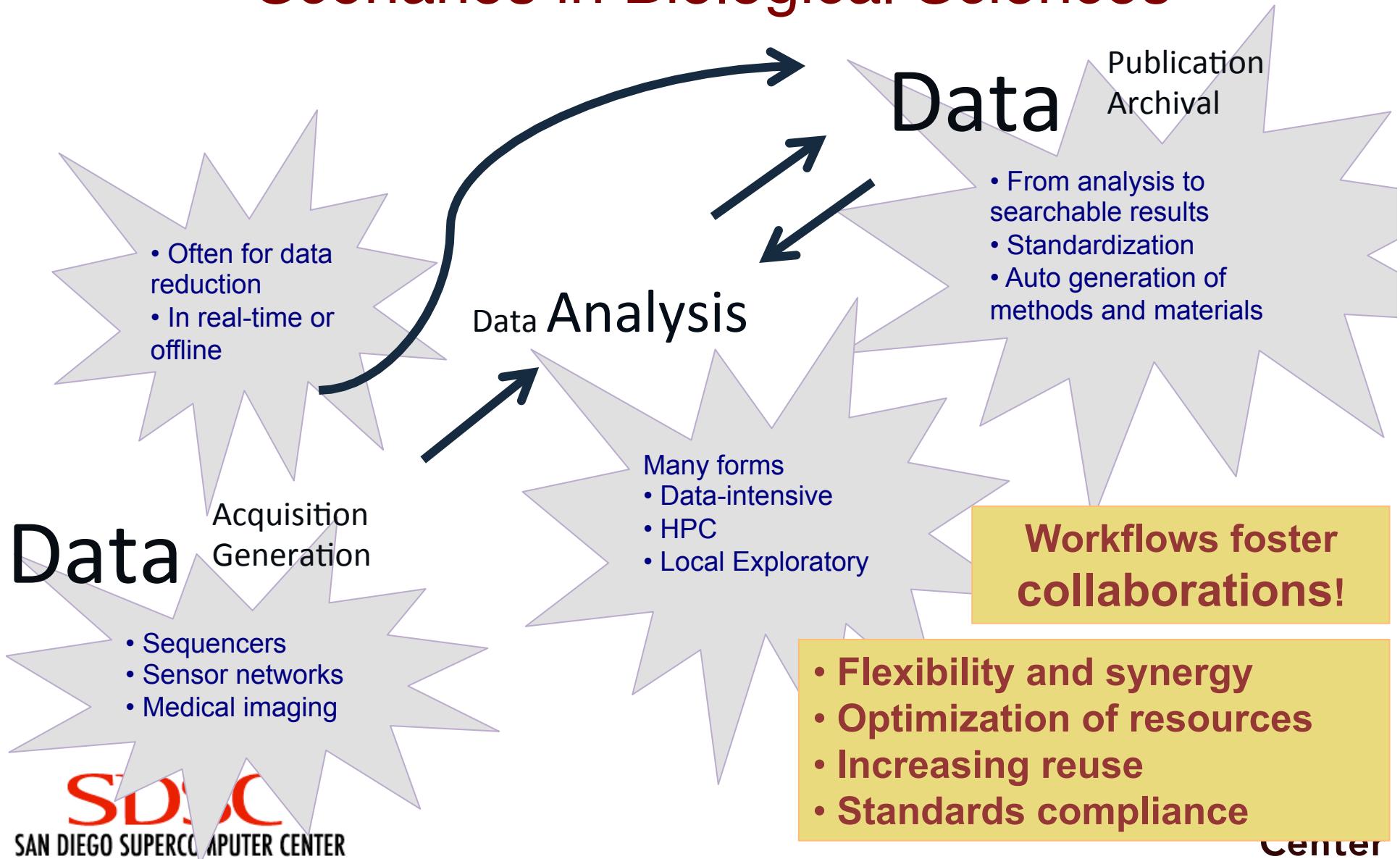
Many Scientific Workflow Systems



Kepler

- A diverse library of scientific components and usecases
- Transparent support for multiple workflow engines
- Used by many communities, specialized gateways and individuals

Workflows are Used in These Diverse Scenarios in Biological Sciences



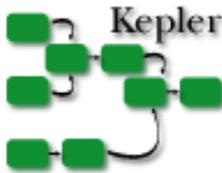
A Toolbox with Many Tools



- Data
 - Search, database access, IO operations, streaming data in real-time...
- Compute
 - Data-parallel patterns, external execution, ...
- Network operations
- Provenance and fault tolerance

Need expertise to identify which tool to use when and how!
Require computation models to schedule and optimize execution!

Kepler is a Scientific Workflow System



www.kepler-project.org

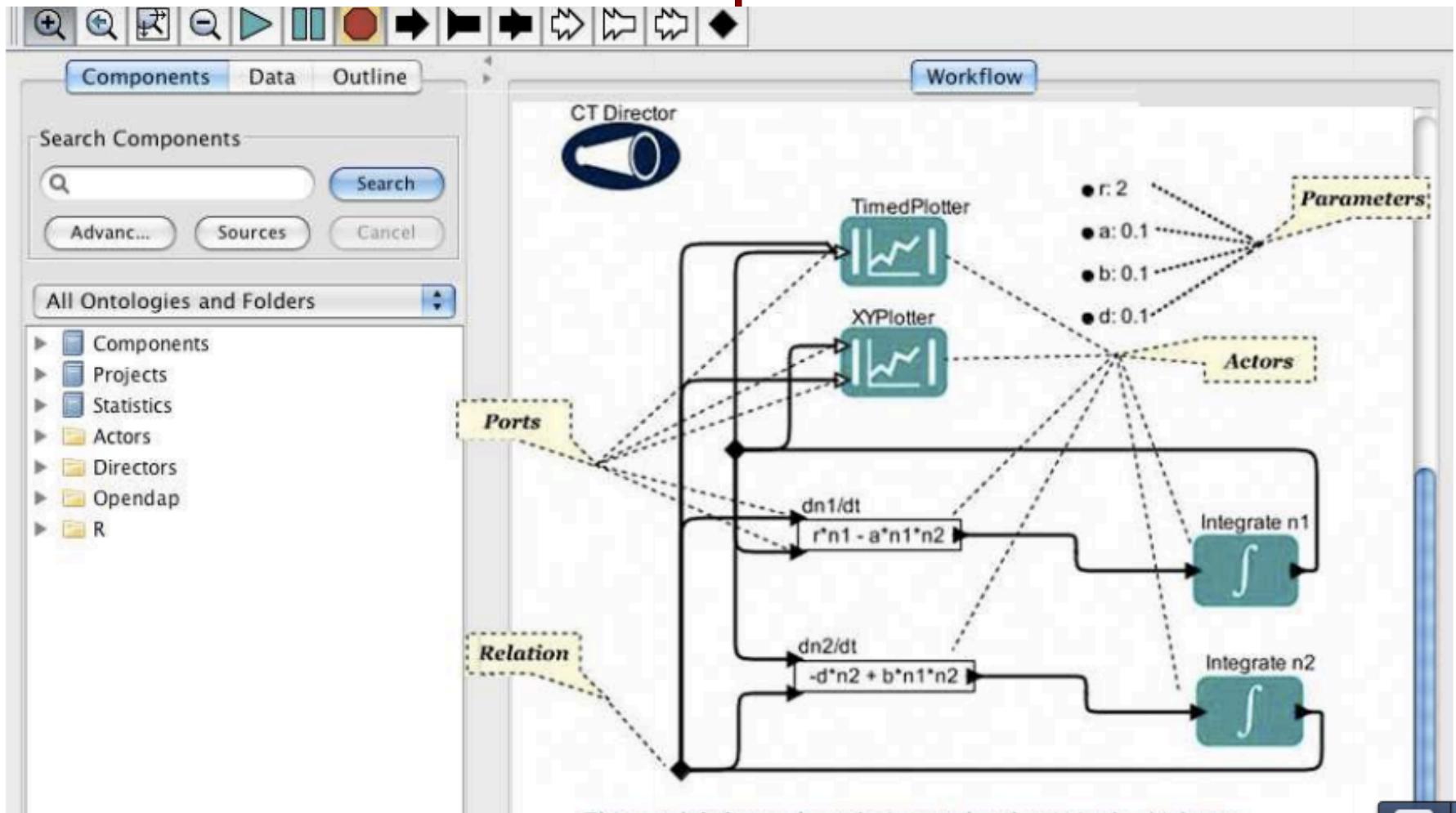
- A cross-project collaboration
... initiated August 2003
- 2.4 released 04/2013
- Builds upon the open-source Ptolemy II framework

Ptolemy II: A laboratory for investigating design

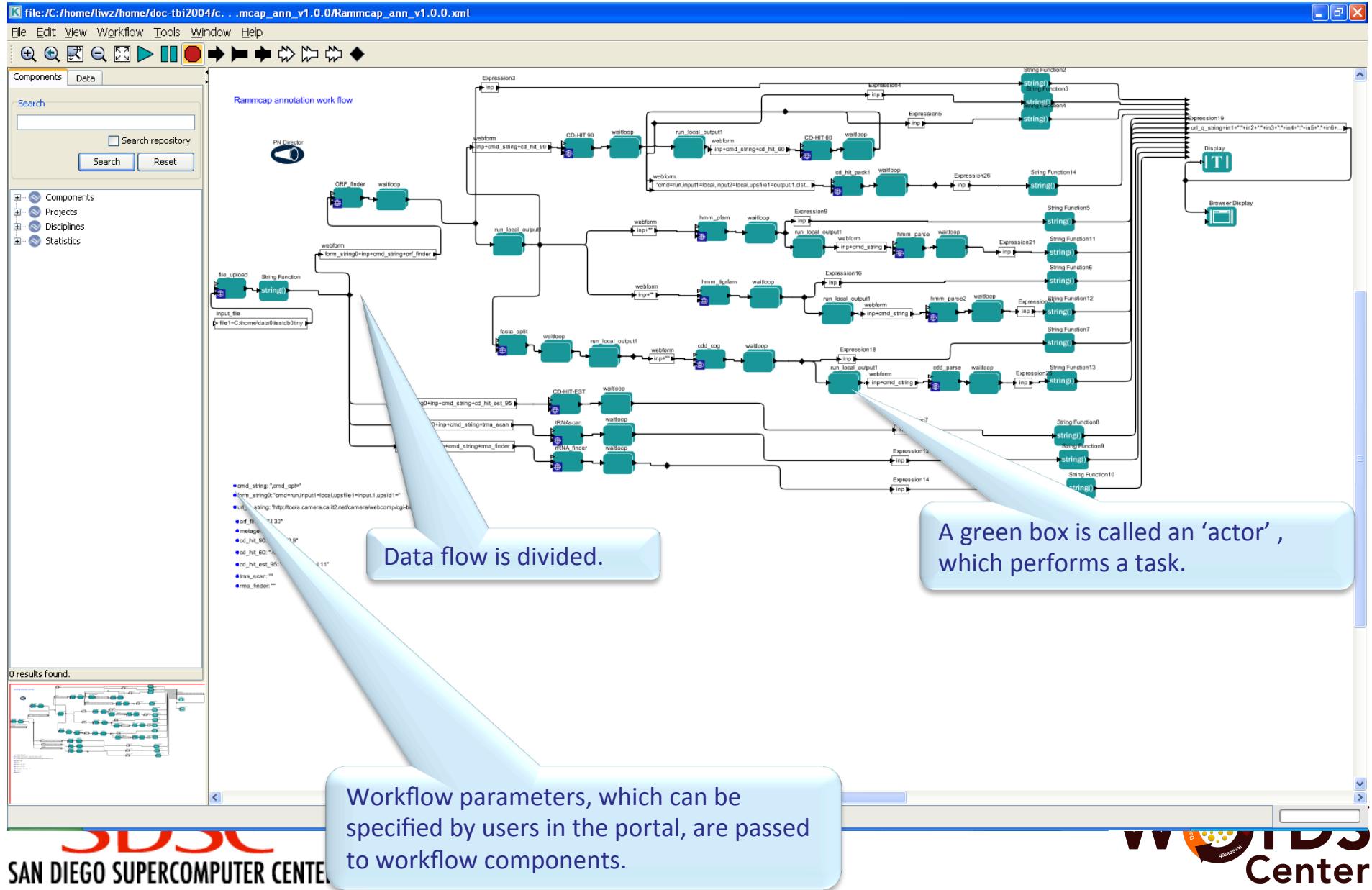
KEPLER: A problem-solving environment for Scientific Workflow

KEPLER = “Ptolemy II + X” for Scientific Workflows

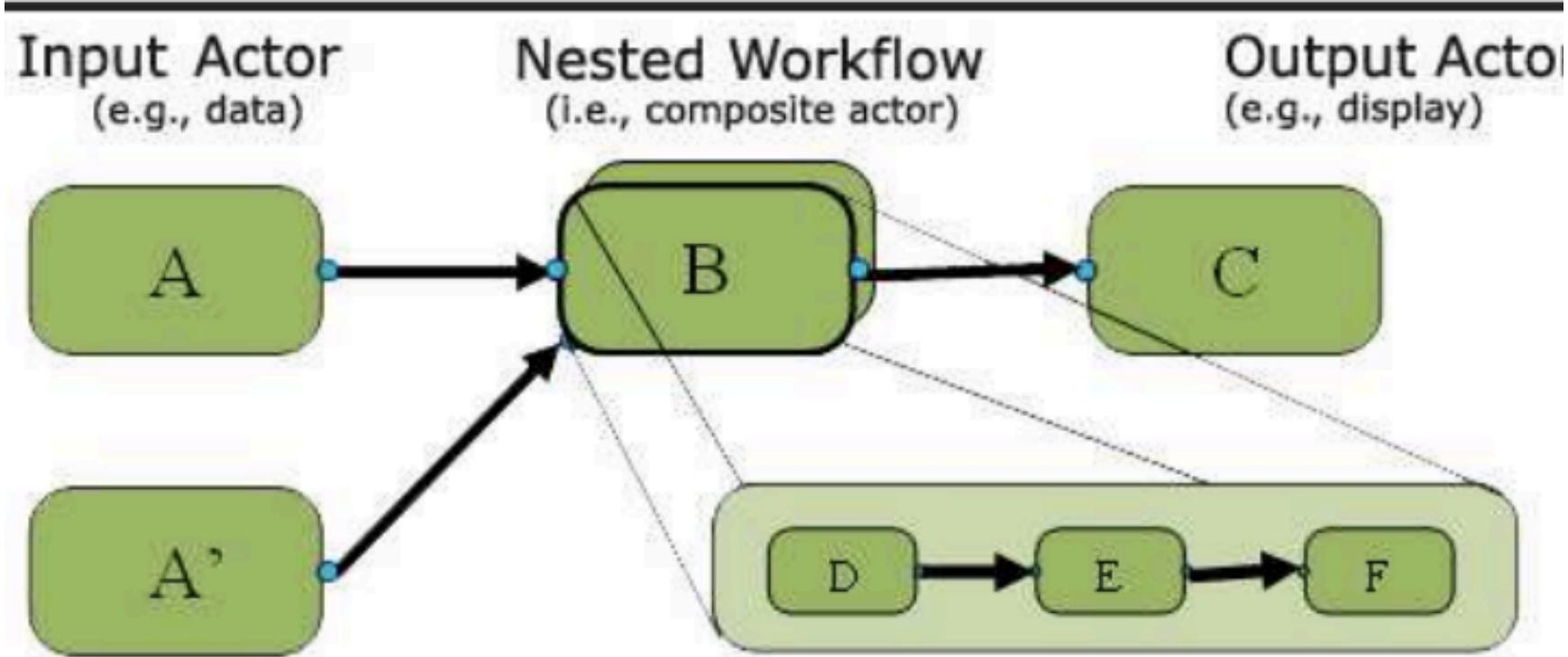
Basic components and terminologies in Kepler



A Typical Kepler Workflow



Kepler Actor Types



Some actors in place for...

- Command Line wrapper tools (**local execution, ssh, scp, ftp, etc.**)
- Generic Web Service Clients for **SOAP** and **REST**
- A suite of **cloud computing** actors for VM instantiation and management
- Job management actors for **HPC, GPU, SGE** and other commodity clusters
- Customizable **RDBMS** query and update
- Distributed data parallel patterns, e.g., Map, Reduce, Cross
- **Hadoop**, Stratosphere, and **Spark** integration
- iRODS support
- Native **R** and **Matlab** support
- Communication with external workflow engines, e.g., **KNIME**
- Communication with sensor data loggers through actors and services
- Imaging, Gridding, Vis Support
- Textual and Graphical Output
- Integration with **Jython, JavaScript, Java, JRuby**
- ...more generic and domain-oriented actors...

Workflow Execution across Multiple Environments

- Execution Choice Actor: Multiple types of executions within one workflow
 - Local execution
 - Hadoop execution
 - EC2 execution
 - Remote job execution
- Useful for **heterogeneous execution requirements**



Running on Heterogeneous Computing Resources

- Execution of models on where they run most efficiently -

Different models have different computing architecture needs!

e.g., memory-intensive, compute-intensive, I/O-intensive

**Local Cluster
Resources**



**NSF/DOE: TeraScale
Resources (XSEDE)**



(Gordon)



(Comet)

SDSC



(Lonestar)

TACC



(Stampede)

**Private Cluster:
User Owned
Resources**



Google Cloud Platform





More on Heterogeneous
Computing at the
Workflow Management
Session Thursday
afternoon

To summarize...

Workflows are used for

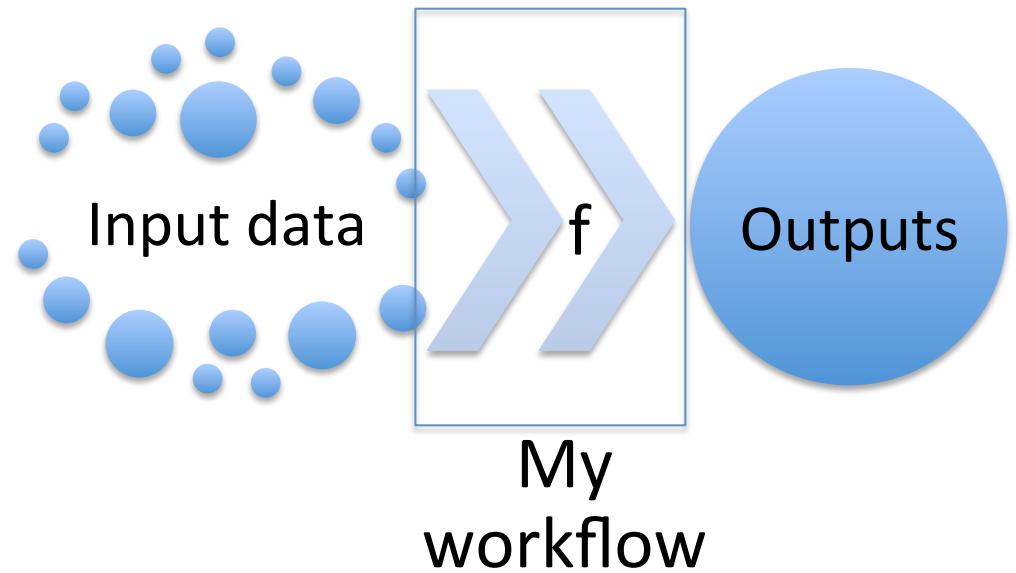
- **documentation** of the analysis
- **visual representation** of analytical steps
- ability to work **across multiple software and platforms**
- **distributed data-parallel** scheduling and execution
- **reproducibility** of a given project with little effort
- **reuse** of part or all of a workflow in a different project



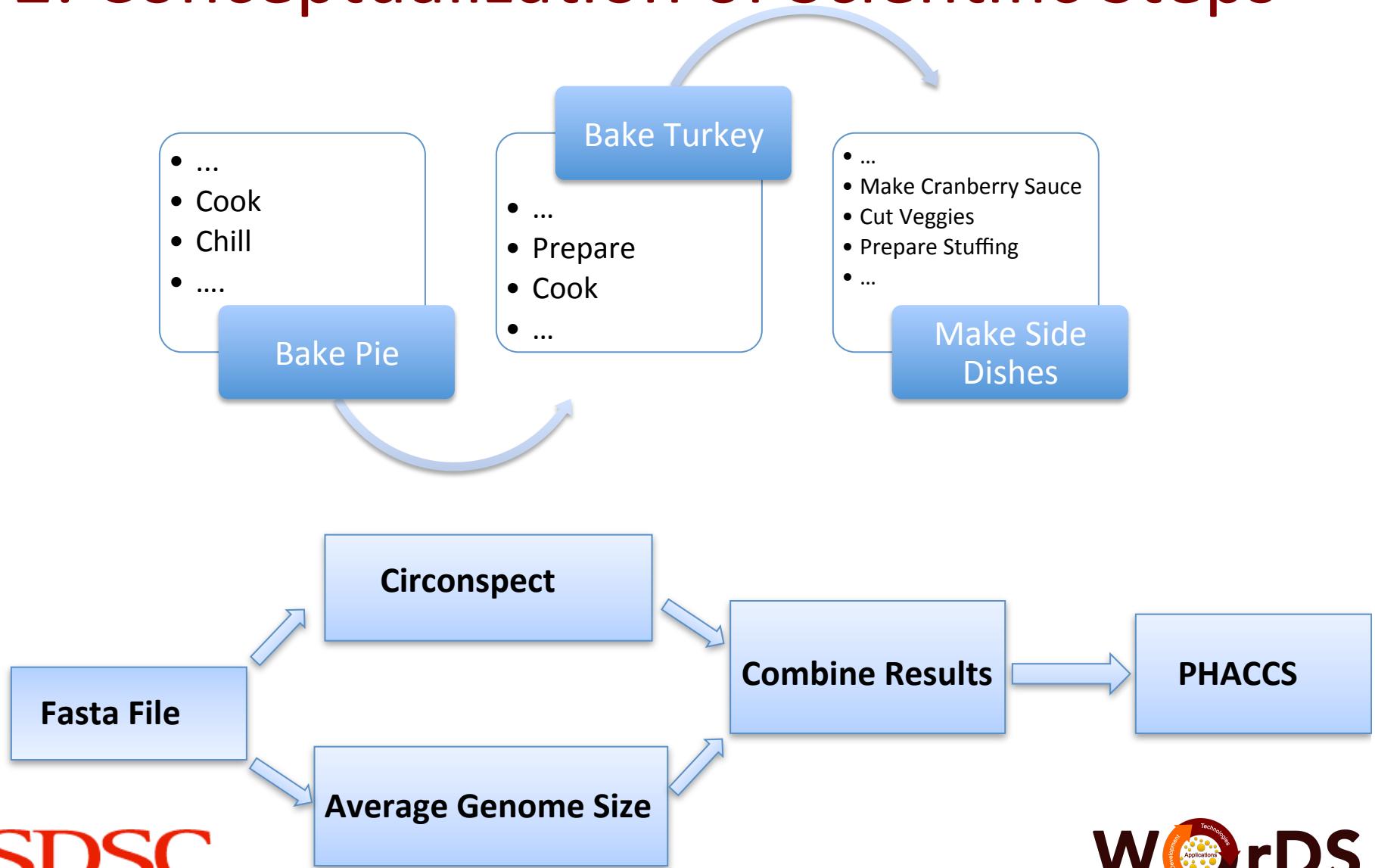
How Do I Start Conceptualizing a Computational Workflow?

1: Start with the Workflow As a Blackbox

- Treat the whole workflow as a blackbox
 - What is the usecase/application?
 - What is the science question this workflow is solving?
 - What is the input data?
 - What are the expected outcomes?
- Give the workflow a title based on initial assessment!

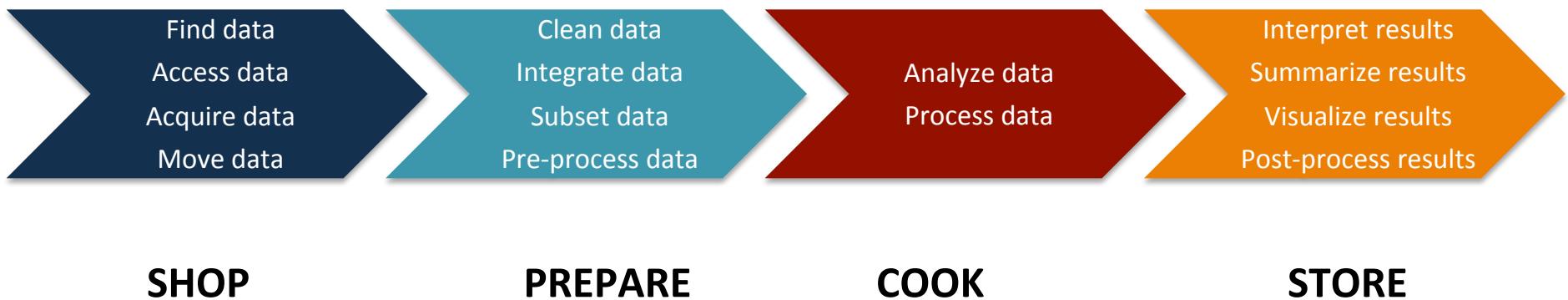


2: Conceptualization of Scientific Steps



3: Treat Each Step Like a Workflow

- until you reach an atomic functional step -



Some questions to ask:

- Where and how do I get the data?
- What format is the data in?
- How do I integrate or subset datasets?
- How do I analyze the data and what is the analysis function?
- What are the parameters to customize each step?
- What are the computing needs to schedule and run each step?
- How do I make sure the results are useful for the next step or as scientific products?

Searching for Data

- You need to know:
 - What you are looking for
 - Some data formats: FASTA, FASTQ, PHYLIP, PDB, PFAM, SAM, BAM,...
 - Meaning of data: how it scientifically/conceptually relates to the analysis you are implementing.
 - Where to find what you are looking for
 - Searchable data archives/repositories
 - A database system
 - Sensor networks or other scientific instruments



Moving Data



- **Data access technologies:**
 - ftp, scp, web services, ...
- **Data querying technologies:**
 - rbdms, opendap, ...

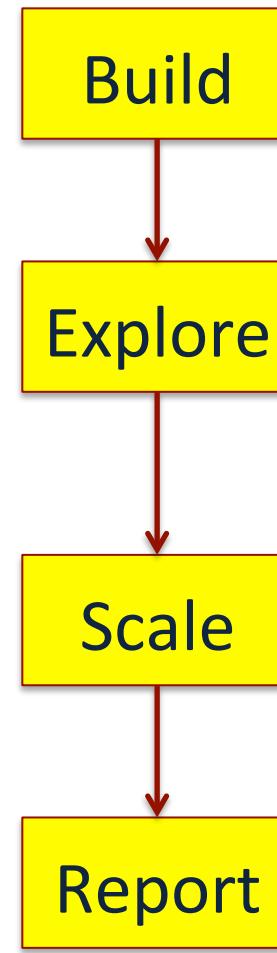
Pre-processing data

- **Clean:** Reformat, QA/QC, duplicate removal, etc.
- **Divide/chop/slice-n-dice:** query, subset, filter, ...
- **Transform:** fasta->fastq; .fcs->.txt
- **Reduce:** statistical operations, and
- **Integrate:** combining data from previous steps

Step 4: Start Building Each Step Including the Alternatives

-- Integration of Many Tools to Serve a Purpose --

- Alternative tools
- Multiple modes of scalability
- Support for each step of the development and production process
- Different reporting needs for exploration and production stages



Need toolboxes with many tools for:

- data access,
- analysis,
- **scalable execution**,
- fault tolerance,
- provenance tracking,
- reporting
- ...

Build Once, Run Many Times...

- The same workflow should support **experimental work** and **dynamic scalability** on many platforms
- Scalability based on:
 - data volume and velocity
 - dynamic modeling needs based on various optimization criteria
 - changes in network, storage and computing availability

Step 5: Save and Share Reports and Final Products with your Team

- Computational data science is in the middle bridging the gap between data and insight
 - Computational Data Scientists tackle challenging questions and define the steps to achieve the results as a workflow
- Developers/computer scientists use their favorite tools to implement the methods in the workflow
- The process is kept sharable, standardized, scalable and accountable

Why Accountable Science?

- Scientific experiments involve many:
 - Data
 - Which data came from which source?
 - Which version of the data?
 - Processes
 - What processes ran in which order?
 - Which libraries were used?
 - Collaborators
 - Who produced what?

Provenance Helps with Accountability and Reproducibility

- *Chronology of ownership, custody, or location of historical object* – Wikipedia
- Data and Process Provenance
 - Inputs, outputs, intermediate results
 - Workflow: actors, links, parameters, etc.

Why is Provenance Useful?

- Keep the association of results to processes
- Make it easier to validate/regenerate results and processes
- Enable comparison between different workflow versions
- Smart re-runs
- Failure recovery

The Need for Provenance Interchange

- Applications and Tools:
 - Workflow systems
 - Databases
 - Programming Languages
 - File Systems
- Vocabularies and Ontologies:
 - Open Provenance Model (OPM)
 - Dublin Core



PROV

- Set of W3C recommendations
- Finalized September 2013
- Core concepts
 - Identifying objects, processes, etc.
- Conceptual data model & serializations:
 - OWL2 ontology
 - XML schema
 - Human-readable notation



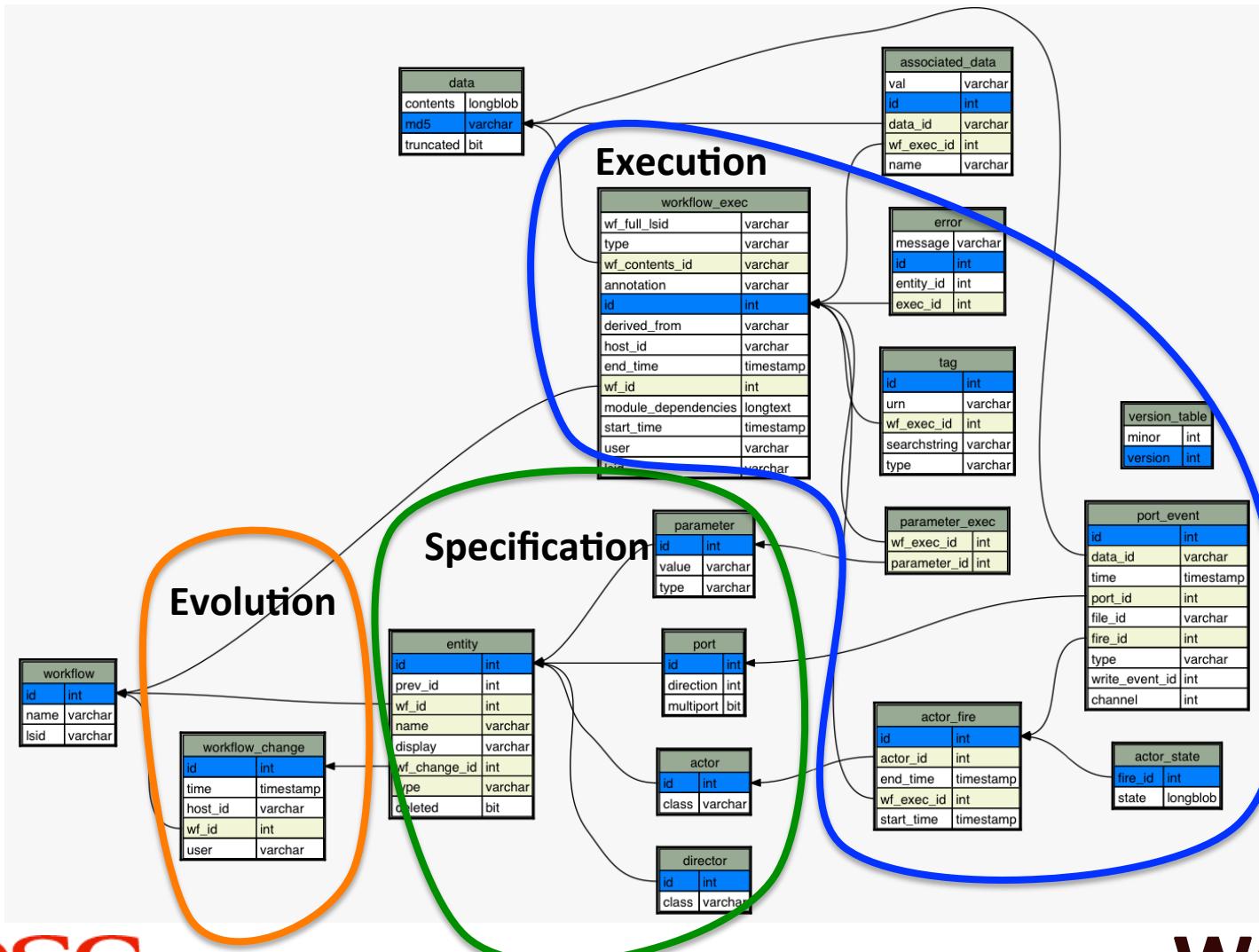
Kepler Provenance Framework

- **What** provenance is recorded:
 - Workflow Specification: actors, ports, connections, parameters, etc.
 - Workflow Evolution: parameter values that change over time, addition/removal of actors, ports, etc.
 - Workflow Execution:
 - Start/stop of workflow, individual actor executions
 - Data exchanged between actors

Kepler Provenance Framework

- **Where** provenance recorded:
 - Modular interface supports saving to different output types.
 - SQL Database: MySQL, Postgres, Oracle, HSQL
 - Text
 - XML

Provenance SQL Schema



Example Queries

- How long did workflow run n take?
- What parameter values were used in run n ?
- What actors does workflow m have?
- How long does a given actor usually take to run?
- What workflows did a specific user run on a given date?

Sharing Process and Data Products

- Are you a dog person or parrot person? -



There are many reasons for sharing workflows

- Collaborative scientific research
 - Inform collaborators of your progress
 - Formalize and standardize application development
 - Increase reuse and repurposing of existing workflows
 - Run collaborators' workflows on similar infrastructure
- Reproducible publications
 - Rerun/replay workflow executions that generated the published products
 - Validate outputs

Sharing Workflow Results

- Version Control Repositories
 - e.g., Git
- Cloud File Hosting Services
 - e.g., Dropbox, Google Drive, Microsoft OneDrive
 - File locker
 - Synchronization across devices
 - Share files and directories
- SeedMe
 - Web platform to share results in near real-time
 - Specialized viewers for plots, sequences, videos, etc.
- Create an iPython notebook to share the execution history, input parameters and results

HAPPY ENDING?

There are no happy endings.
Endings are the saddest part,
So just give me a happy middle
And a very happy start.

İlkay Altintas

altintas@sdsc.edu

