# *Distributed Workflow-Driven Analysis of Large-Scale Biological Data using bioKepler*

**Ilkay ALTINTAS, Ph.D.**

*Deputy Coordinator for Research, San Diego Supercomputer Center, UCSD*

*Lab Director, Scientific Workflow Automation Technologies*

*altintas@sdsc.edu*

SDSC    UC San Diego    NSF

*bioKepler.org*

# *Welcome to SDSC!*

– Workshop website

http://www.biokepler.org/workshops/2012-sep

–Logistics for the next two days

SDSC

UC San Diego

NSF

*bioKepler.org*
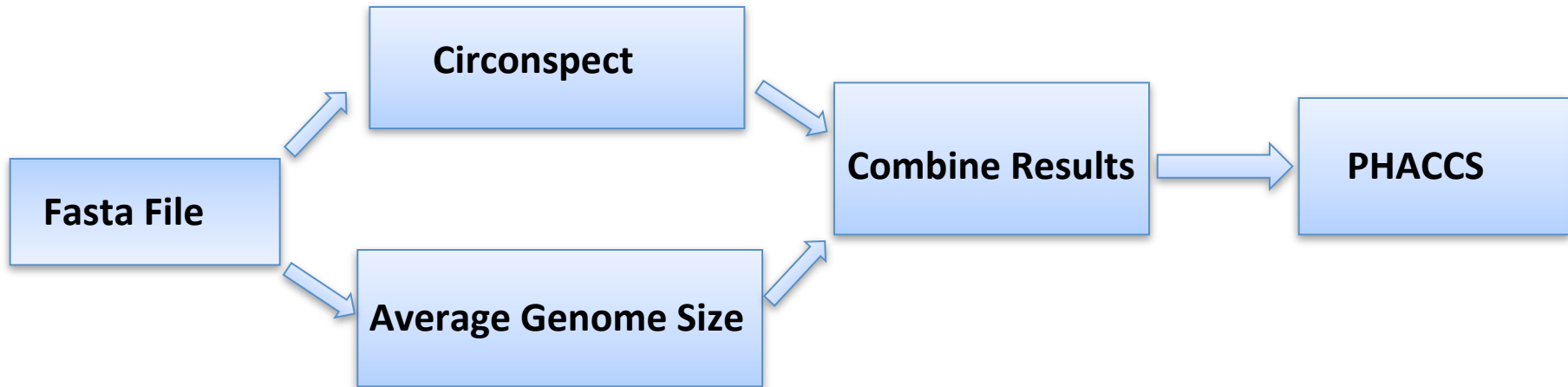
# *So, what is a scientific workflow?*

Scientific workflows emerged as an answer to the need to **combine** multiple Cyberinfrastructure components in **automated process networks**.
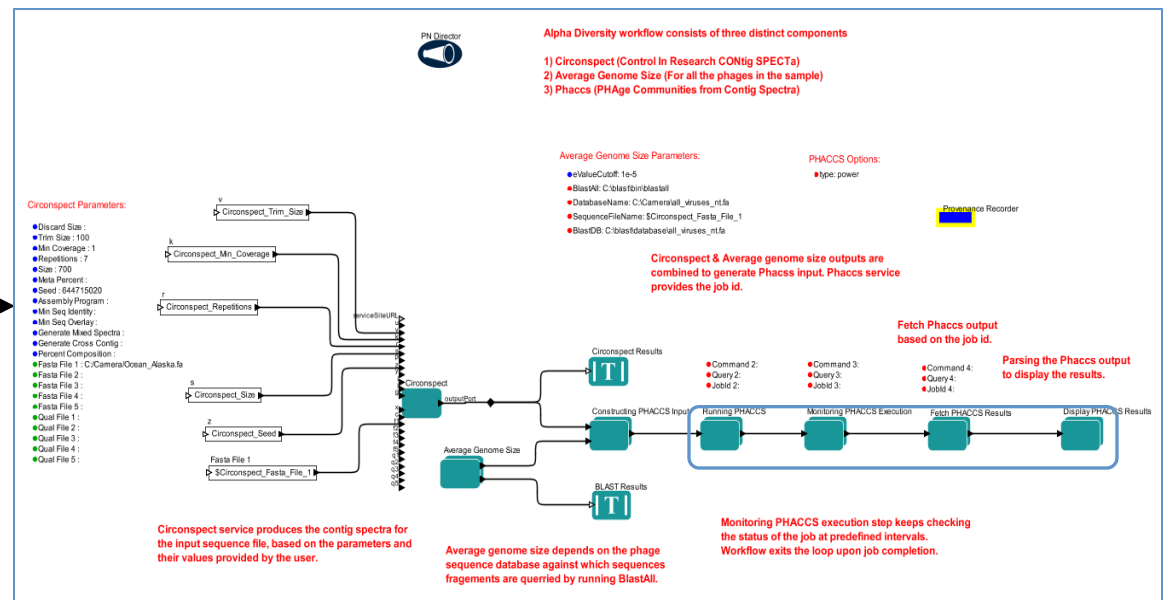
*bioKepler.org*

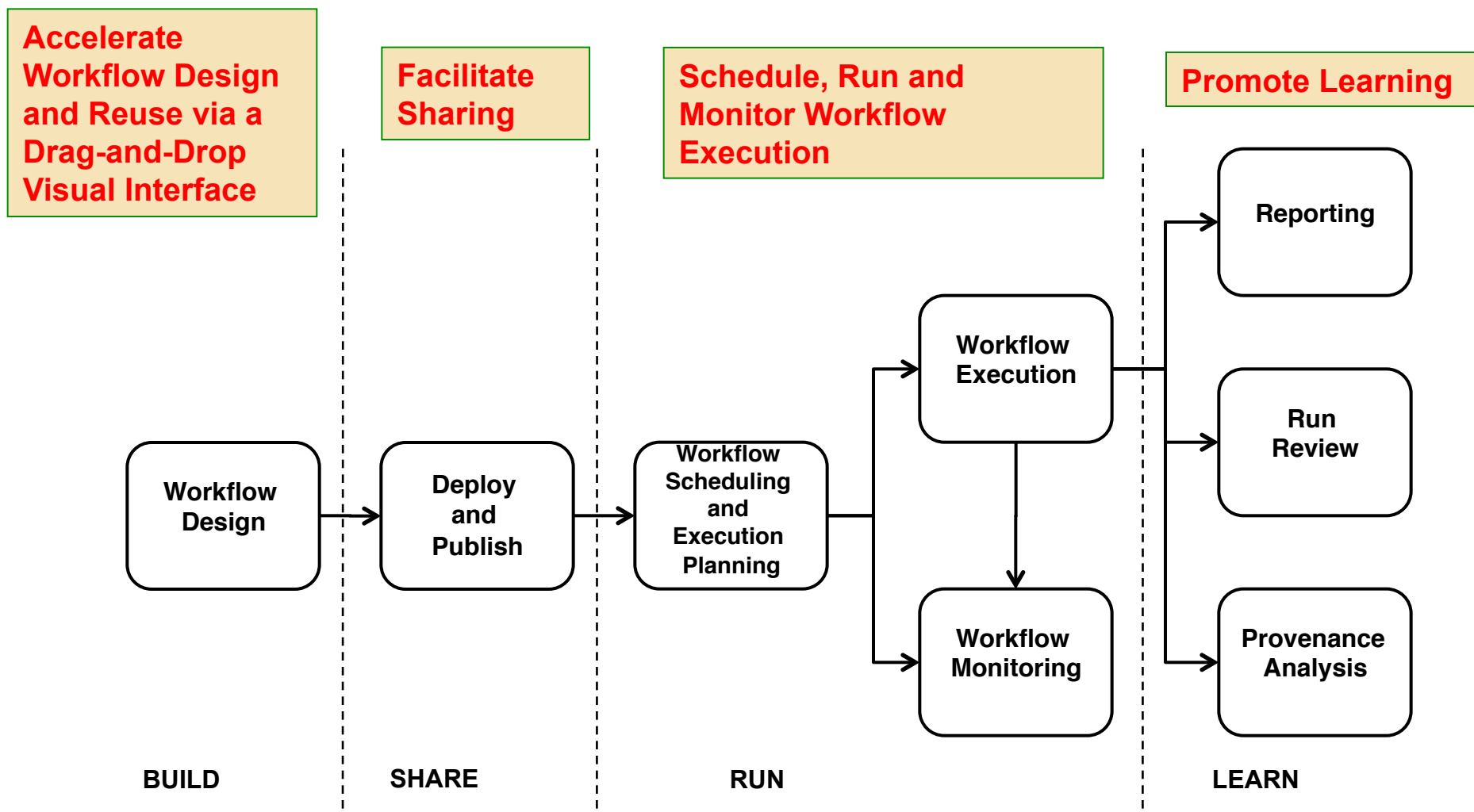# The Big Picture is Supporting the Scientist

## From "Napkin Drawings" to Executable Workflows



Fasta File → Circonspect → Combine Results → PHACCS

Fasta File → Average Genome Size → Combine Results

*Conceptual SWF*

*Executable SWF*

SDSC · UC San Diego · NSF

# *Workflows are a Part of Cyberinfrastructure*

**Accelerate Workflow Design and Reuse via a Drag-and-Drop Visual Interface**

**Facilitate Sharing**

**Schedule, Run and Monitor Workflow Execution**

**Promote Learning**

```
Workflow Design  →  Deploy and Publish  →  Workflow Scheduling and Execution Planning  →  Workflow Execution  →  Reporting
                                                                                            ↓                      Run Review
                                                                                       Workflow Monitoring        Provenance Analysis
```

**BUILD**          **SHARE**          **RUN**          **LEARN**

**Support for end-to-end computational scientific process**

# *Kepler is a Scientific Workflow System*

**Kepler**

www.kepler-project.org

- A cross-project collaboration
    … initiated August 2003
- 2.3 release released 01/2012
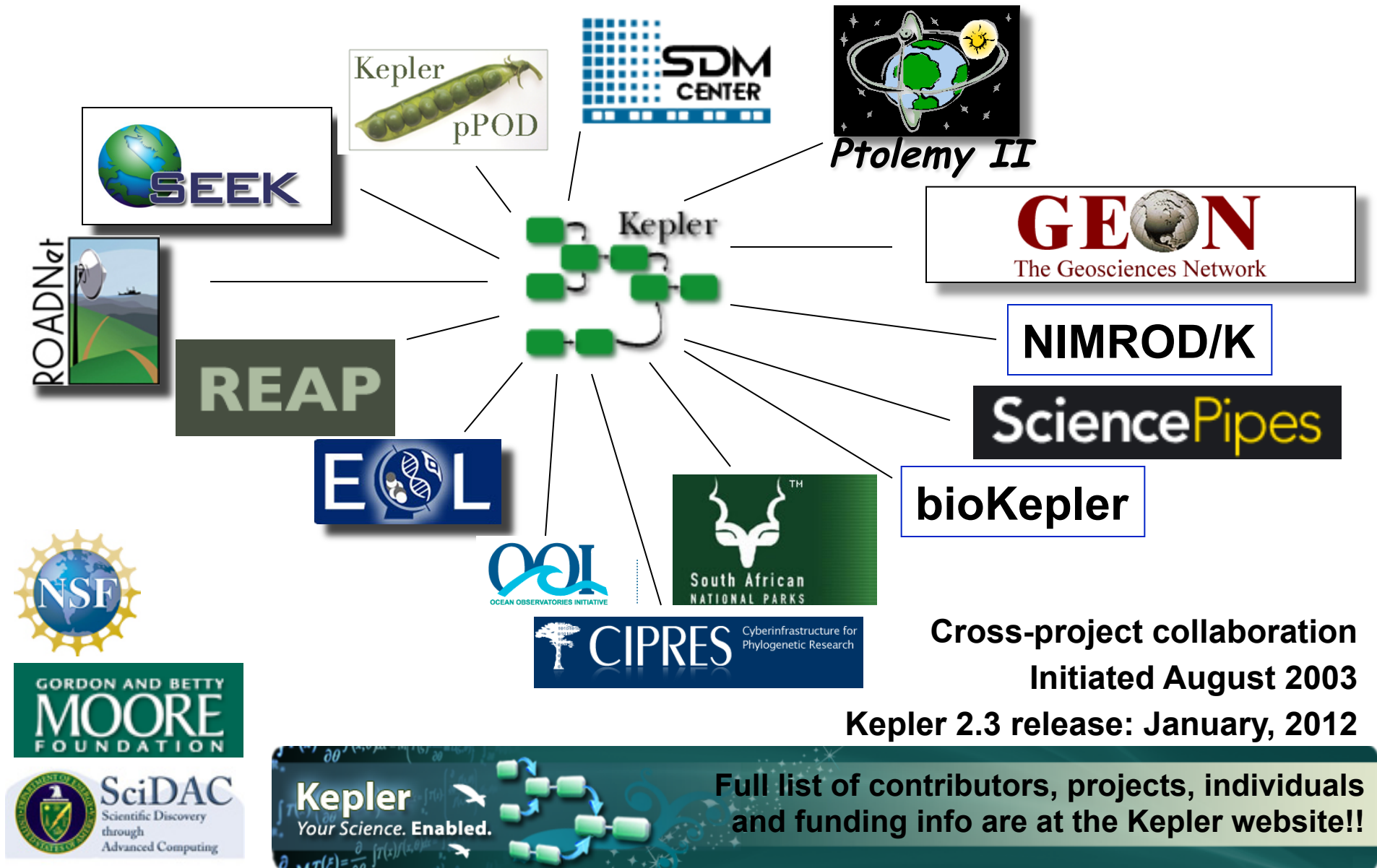
- Builds upon the open-source Ptolemy II framework

Ptolemy II: A laboratory for investigating design

KEPLER: A problem-solving environment for Scientific Workflow

KEPLER = "Ptolemy II + X" for Scientific Workflows

SDSC

UC San Diego

NSF

*bioKepler.org*

# A Typical Kepler Workflow



**Data flow is divided.**

**A green box is called an 'actor', which performs a task.**

**This special actor represents an annotation component, such as BLAST search.**

**Workflow parameters, which can be specified by users in the portal, are passed to workflow components.**

bioKepler.org

# *Kepler is a Team Effort and Modular*



SEEK

Kepler pPOD

SDM CENTER

Ptolemy II

ROADNet

GEON The Geosciences Network

Kepler

REAP

NIMROD/K

SciencePipes

E⊙L

bioKepler

OOI OCEAN OBSERVATORIES INITIATIVE

South African NATIONAL PARKS

CIPRES Cyberinfrastructure for Phylogenetic Research

NSF

GORDON AND BETTY MOORE FOUNDATION

SciDAC Scientific Discovery through Advanced Computing

Kepler Your Science. Enabled.

**Cross-project collaboration**
**Initiated August 2003**
**Kepler 2.3 release: January, 2012**

**Full list of contributors, projects, individuals and funding info are at the Kepler website!!**

# Requirements are similar for many domains

## -- with slight variations --

*bioKepler.org*

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

- Important Attributes

  <span style="color:red">Assemble</span> complex processing easily

  <span style="color:red">Access transparently</span> to diverse resources

  <span style="color:red">Incorporate</span> multiple software <span style="color:red">tools</span>

  <span style="color:red">Assure reproducibility</span>

  Build around <span style="color:red">community development</span> model

SDSC    UC San Diego    NSF

*bioKepler.org*

# *Many Bioinformatics Workflow Systems*

DiscoveryNet

Triana

Vistrails

Clover

Pegasus

Ergatis

**Kepler**

Galaxy

Taverna

Trident

Pipeline Pilot

**2000**　　　　　　　　**2005**　　　　　　　**2010**　　**2012**

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

SDSC　　UC San Diego　　NSF

*bioKepler.org*

# *Workflows are Used in These Diverse Scenarios in Biological Sciences*

**Data** Publication Archival

- From analysis to searchable results
- Standardization
- Auto generation of methods and materials

**Data** Analysis

- Often for data reduction
- In real-time or offline

Many forms
- Data-intensive
- HPC
- Local Exploratory

**Data** Acquisition Generation

- Sequencers
- Sensor networks
- Medical imaging

**Workflows foster collaborations!**

- Flexibility and synergy
- Optimization of resources
- Increasing reuse
- Standards compliance

SDSC

UCSanDiego

NSF

bioKepler - September, 2012

12

# *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!**

# *CAMERA Example:*

# *Using Scientific Workflows and Related Provenance for Collaborative Metagenomics Research*

**Community Cyberinfrastructure for Advanced Microbial Ecology Research and Analysis**

**(CAMERA)**

http://camera.calit2.net

SDSC

UC San Diego

NSF

camera
Marine Microbial Ecology

# CAMERA is a Collaborative Environment

**Data Discovery**
GIS and Advanced query options

**Data Cart**
Multiple Available Mixed collections of CAMERA Data (e.g. projects, samples)

**User Workspace**
Single workspace with access to all data and results (private and shared)

**Data Analysis**
Workflow based analysis

**Group Workspace**
Share specified User Workspace data with

bioKepler - September, 2012

# *Workflows are a Central Part of CAMERA*

- **CAMERA-supported**
  - 28 existing workflows
- **Workflows under development**
  - Fragment Recruitment Viewer
  - Next Generation Sequencing
  - VIROME Pipeline
  - Standalone b tools
  - National Cent Research
  - Joint Genome Institute
- **User built**
  - Currently running in a sandbox
  - Will be ported to a virtual cloud environment

**More than 1500 workflow submissions monthly!**

QC filter

Taxonomy Binning

Duplicate filtering

Metagenomic Annotation and Clustering

Assembly

Comparison, Statistical analysis, and more workflows

- **Inputs:** from local or CAMERA file systems; user-supplied parameters
- **Outputs:** sharable with a group of users and links to the semantic database

All can be reached through the CAMERA portal at:http://portal.camera.calit2.net

SDSC

UC San Diego

NSF

bioKepler - September, 2012

16

# *CAMERA Portal - Workflows*

| Home | Browse Data | Data Analysis | Sharing | | Submit Data to CAMERA |
|---|---|---|---|---|---|

**Main**  Workflows  Blast Results

Data Analysis ⟩ Main          Quick Navigation ▾    Search CAMERA Data 🔍    **Help?**

## DATA ANALYSIS

CAMERA utilizes workflows to launch data analysis tools. Workflows are configurable analysis packages that can be applied to data within the CAMERA workspace or to data uploaded from the local system. CAMERA workflows include:

- Metagenomic data annotation and clustering (RAMMCAP)
- BLAST tools (Click here for a complete list of CAMERA Reference Datasets ⌐)
- RNA and Orf Prediction
- Click here for a complete list of CAMERA workflows ⌐
- Job Submission Policy ⌐
- FASTA Validation Guidelines ⌐

**Queue Status**  (Last Update: 3/20 20:30:03 PDT)

| | |
|---|---|
| **CAMERA Cluster Utilization:** | **95%** |
| Workflows / Jobs* Running: | 4 / 281 |
| Workflows / Jobs* Queued: | 4 / 2231 |
| | |
| **CAMERA Compute Cloud:** | **Enabled** |
| Workflows / Jobs* Running: | 1 / 320 |

\* Each workflow is composed of numerous jobs.

### To Launch a Workflow:

1. Select an analysis from the workflow menu and click the 'Start >' button at the bottom of the page.

2. Fill out the analysis parameters and click the 'Submit Workflow!' button at the bottom of the page. For additional information regarding the parameters, mouse-over the circular 'i' button.

3. Use the Results and Status page to view and share results. For BLAST, use the dedicated Blast Results viewer.

### Upload User Workflows (Beta):

CAMERA provides a collaborative environment for analysis and data. As part of this environment, users can upload and share their own workflows/analysis with their colleagues or with the greater scientific community. Please note this new feature is in a BETA state and may have problems. Initially, this area is for those who already have an understanding of the Kepler workflow system.

Click here to upload a workflow. Note that you must create a group to associate with the user workflow.

To get started with workflow development, please go here ⌐
By submitting or running a workflow, you are agreeing to the Terms and Conditions ⌐

17

# CAMERA Workflows

Launch CAMERA
Supported Workflows

**Execute Workflow: Metagenomic data annotation and clustering**

- Quick start 📕
- Scientific example 📕
- Full documentation 📕

**Parameter Help** — Close

```
Parameters for ORF call by six reading frame translation, default value is
"-l 30 -L 30 -t 11"

-l 30 -L 30 means the cutoff length of ORFs is 30 amino acid
-t 11 means translation table 11

Options:
  -l minimal length of orf, default 20
  -L minimal length of orf between 2 stop codons, default 40
  -t translation table, default 1
  -b ORF begin option: default 2
     1: start at the begining of DNA sequence or after pervious stop codon
     2: start with the first ATG if there is a stop codon upstream
     We don't know which ATG is the real start, but for prokaryotic DNA,
     a fragment between a stop codon and the first ATG can not be part of real g
     Therefore, -b 2 is recommanded for prokaryotic
  -e ORF end option: default 1
     1: end at the end of DNA sequence or at a stop codon
     2: must end at a stop codon
```

This is the full RAMMCAP pipeline for analysis of metagenomic sequences. It accepts a FASTA file of raw reads. The pipeline identifies the tRNA, rRNA, and ORFs from the reads. It then performs clustering analysis on the reads and the ORFs. The ORFs are

ORF clustering first run ⓘ — -d 0 -n 5 -p 1 -g 1 -G 0 -c 0.90 -aS 0.8

ORF clustering second run ⓘ — -d 0 -n 4 -p 1 -g 1 -G 0 -c 0.60 -aS 0.8

ORF Finder ⓘ — -l 30 -L 30 -t 11

Read clustering ⓘ — -d 0 -n 10 -l 11 -r 1 -p 1 -g 1 -G 0 -c 0.95 -aS 0.8

E-value cutoff for Pfam ⓘ — 0.001

E-value cutoff for Tigrfam ⓘ — 0.001

E-value cutoff for COG ⓘ — 0.001

Submit Workflow!

# CAMERA Job Status

# *CAMERA Workflow Results*

# *Pushing the boundaries of existing infrastructure and workflow system capabilities*

SDSC

UC San Diego

NSF

*bioKepler.org*

# *Requirements from the User Community*

- Increase reuse
  - best development practices by the scientific community
  - other bio packages

- Increase programmability by end users
  - users with various skill levels
  - to formulate actual domain specific workflows

- Increase resource utilization
  - optimize execution across available computing resources
  - in an efficient, transparent and intuitive manner

- Make analysis a part of the end-to-end scientific model from data generation to publication

**SDSC**  UC San Diego  **NSF**

*bioKepler.org*

# *bioKepler responds to these requirements!*

www.bioKepler.org

CAMERA and other user environments

**Kepler and Provenance Framework**

**bioKepler**

| BioLinux | Galaxy | ... | Clovr | Stratosphere |

**CLOUD and OTHER COMPUTING RESOURCES**
e.g., SGE, Amazon, FutureGrid, XSEDE

**A coordinated ecosystem of biological and technological packages for microbiology!**

# *Reuse, Programmability, Execution*

www.bioKepler.org

CAMERA and other user environments

## Kepler and Provenance Framework

### bioKepler

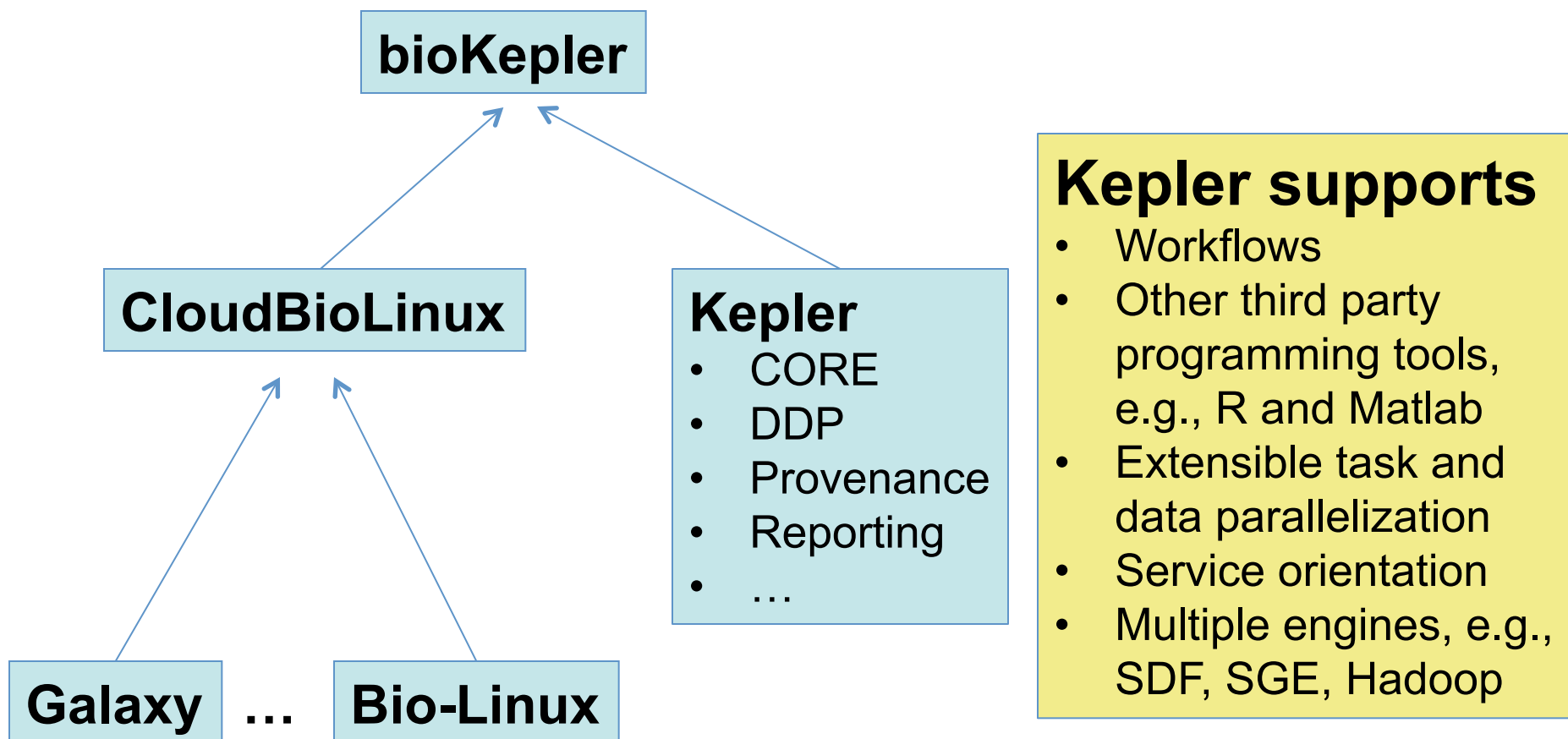| BioLinux | Galaxy | ... | Clovr | Stratosphere |
|----------|--------|-----|-------|--------------|
| bio-linux | Galaxy | | CloVR | StratoSphere *Above the Clouds* |

- Funded by NSF ABI & CI Reuse programs ($1.4M through 2015)
  - Ilkay Altintas (PI) and Weizong Li (Co-PI)
- Development of a comprehensive bioinformatics scientific workflow module for distributed analysis of large-scale biological data

**Will be a huge improvement on usability and programmability by end users!**

SDSC

UC San Diego

NSF

*bioKepler.org*

# *bioKepler and Other Related Systems*



**bioKepler**

**CloudBioLinux**

**Galaxy** ... **Bio-Linux**

**Kepler**
- CORE
- DDP
- Provenance
- Reporting
- …

**Kepler supports**
- Workflows
- Other third party programming tools, e.g., R and Matlab
- Extensible task and data parallelization
- Service orientation
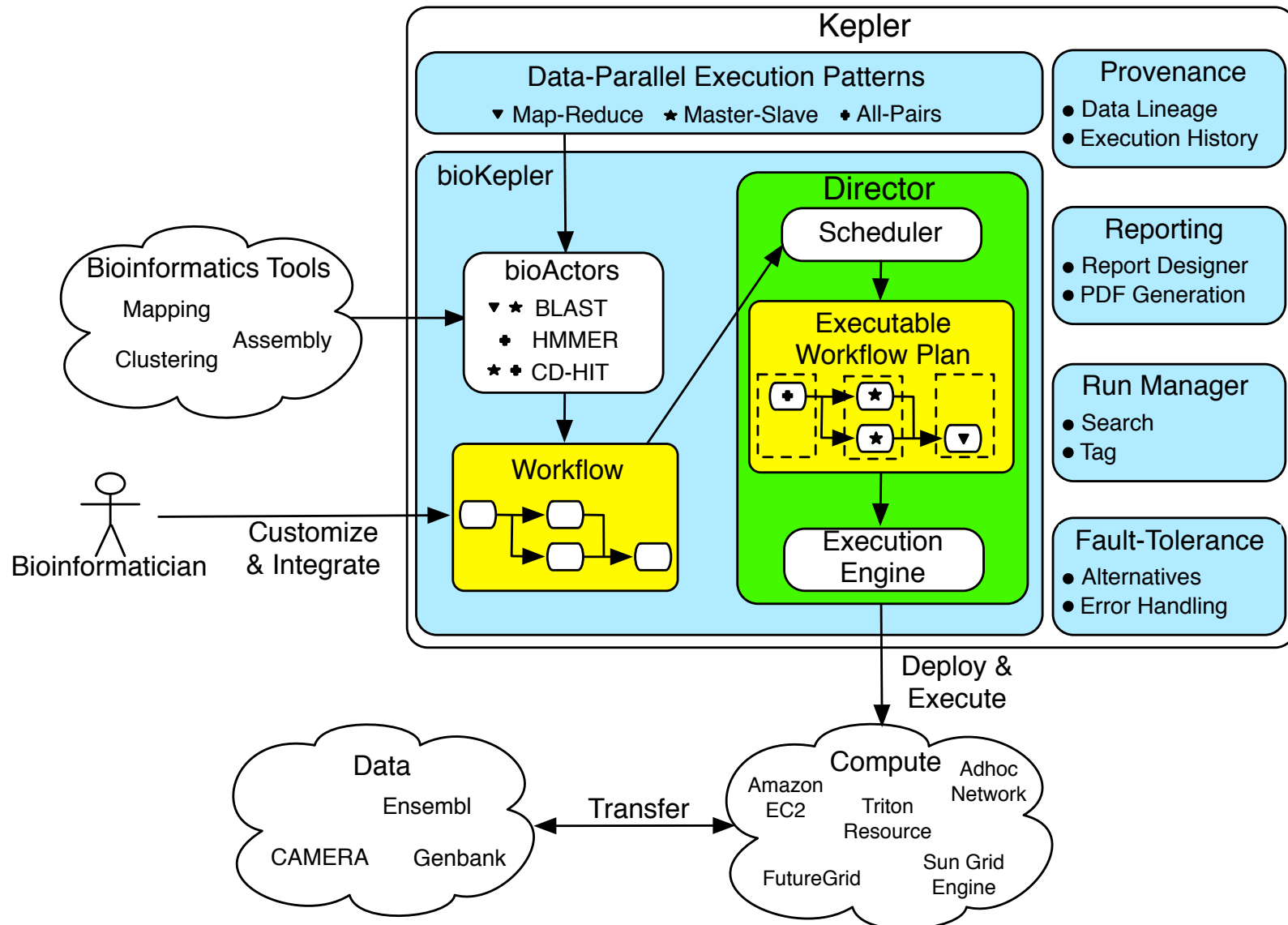- Multiple engines, e.g., SDF, SGE, Hadoop
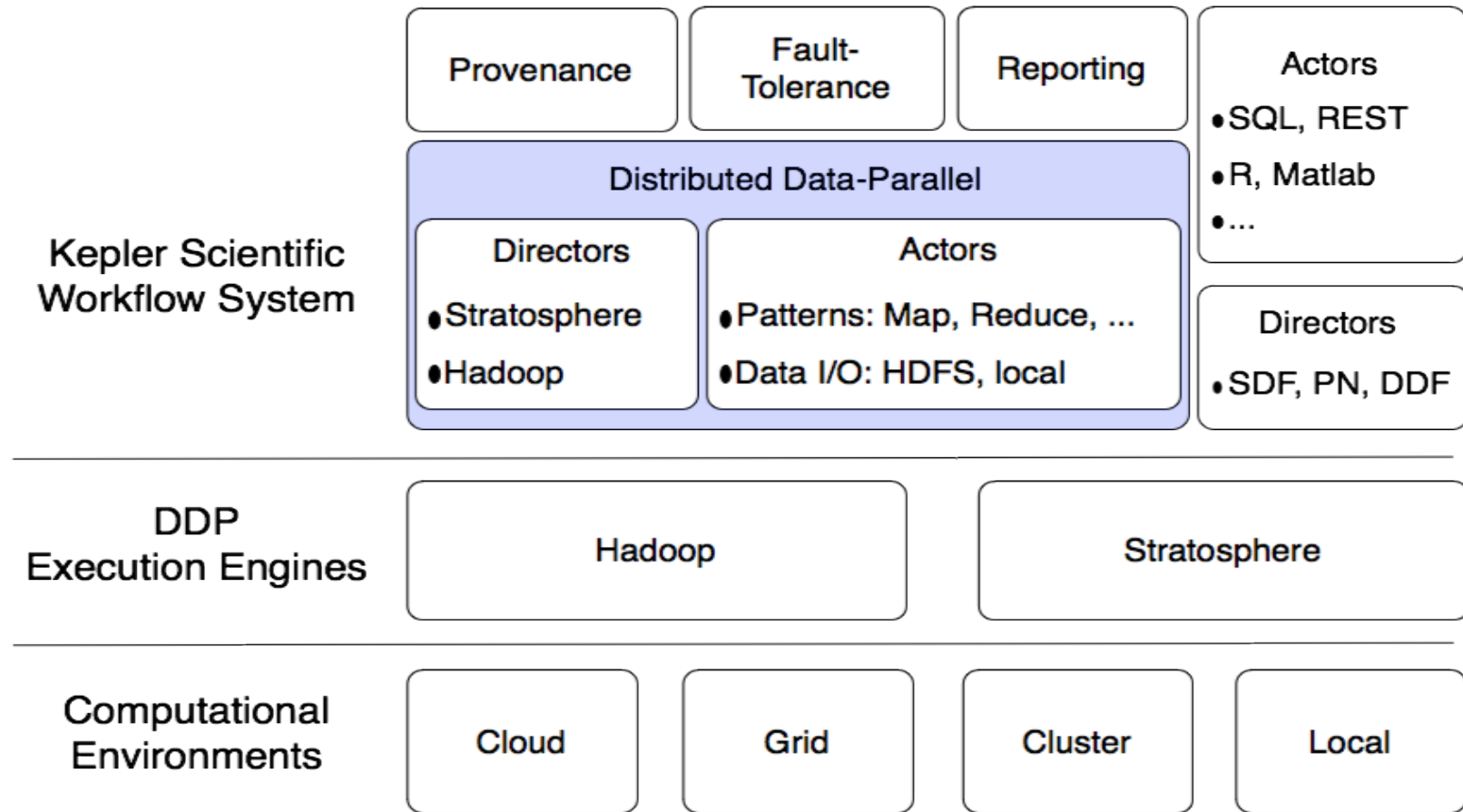
SDSC

UC San Diego

NSF

*bioKepler.org*

# The bioKepler Approach

- Parallel Computation Framework
  - Use Distributed Data-Parallel (DDP) frameworks, e.g., MapReduce, and other parallelization methods to execute subworkflows

- bioActors
  - Configurable and reusable higher-order components for bioinformatics and computational biology

- Transparent support for different execution engines and computational environments

- Deployment on diverse environments

SDSC

UC San Diego    NSF

*bioKepler.org*

# *bioKepler's Conceptual Framework*

UC San Diego

# *bioKepler's Software Architecture*

SDSC

UC San Diego

NSF

*bioKepler.org*

# *bioActors*

- Set of steps to execute a bioinformatics tool locally or in an external environment
  - Locally executable
  - Parallelized external execution
- Customizable by the user based on external packages
  - Tools imported from CloudBioLinux
- Tools are evaluated on their computational requirements

SDSC

UCSanDiego

NSF

# *Example bioActors*

Disciplines
- Biology
  - BlastTabularResultMerge
  - hmm rRNA
  - Acd
  - Alignment
    - FastTree
    - dialign
    - Alignment Consensus
    - Alignment Differences
    - Alignment Dot Plots
    - Alignment Editing
    - Alignment Global
    - Alignment Graphical
    - Alignment Local
    - Alignment Multiple
    - Alignment Profiles
    - Alignment Statistics
  - Clustering
    - Clustering Graph
    - Clustering Sequences
  - Databases
  - Display
  - Edit
  - Enzyme Kinetics
  - Feature Tables
  - Hmm
    - hmmalign
    - hmmbuild
    - hmmcalibrate
    - hmmcalibrate-pvm
    - hmmconvert
    - hmmemit
    - hmmfetch
    - hmmindex
    - hmmpfam
    - hmmpfam-pvm

- **Alignment:** BLAST, BLAT

- **Profile-Sequence Alignment:** PSI-BLAST

- **Hidden Markov Model:** HMMER

- **Mapping:** Bowtie, BWA, Samtools

- **Multiple Alignment:** ClustalW, Muscle

- **Clustering:** CD-HIT, Blastclust

- **Gene Prediction:** Glimmer, Genescan, Fraggenescan

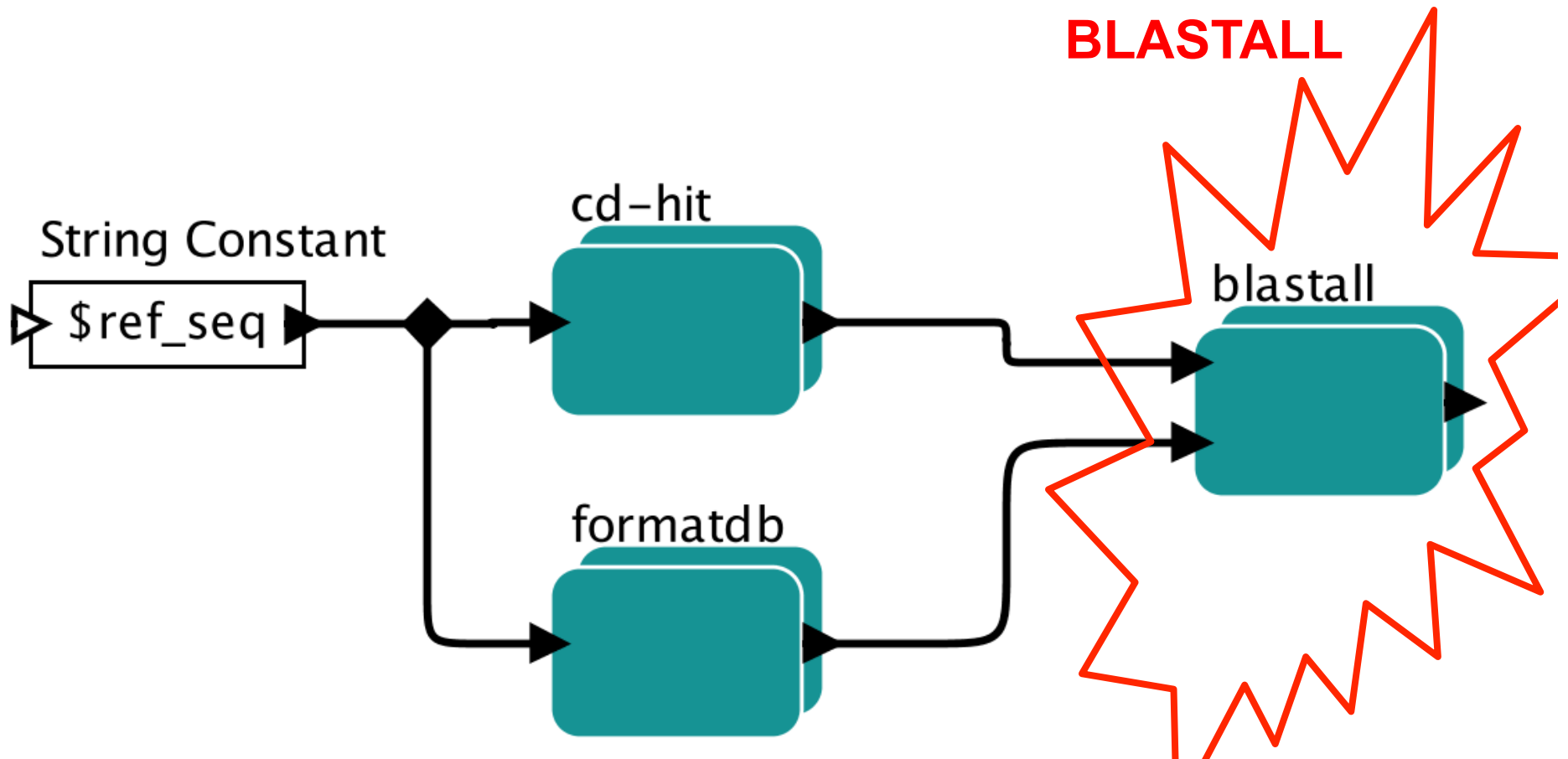- **tRNA prediction:** tRNA-scan, Meta-RNA

- **Phylogeny:** FastTree, RAxML

SDSC

UC San Diego

NSF

*bioKepler.org*

# A Workflow with Three bioActors

SDF Director

ref_seq: small.faa

BLASTALL

String Constant

$ref_seq

cd-hit

formatdb

blastall

# *Current Progress and Release*

- A bioKepler VM executable on Amazon EC2 and FutureGrid
  - Builds upon CloudBioLinux including Bio-Linux and Galaxy
- A bioActor template that can be customized for different execution choices
  - e.g., local vs. Map/Reduce on a specific environment
- Example usecases

SDSC

UC San Diego        NSF

*bioKepler.org*

bioKepler - September, 2012

34

# 1st Workshop on bioKepler Tools and Its Applications

## September 5-6, 2012
## SDSC/UCSD La Jolla, CA

[http://www.biokepler.org/workshops/2012-sep](http://www.biokepler.org/workshops/2012-sep)

### Introductions

**SDSC**   UC San Diego   NSF

**bioKepler.org**

# NEXT:

# *Introduction to bioActors*

*Weizhong Li*

**1st Workshop on bioKepler Tools and Its Applications**

SDSC

UC San Diego

NSF

*bioKepler.org*