

Kepler Scientific Workflow System

Introduction to Workflow-Driven Scientific Computing in Kepler

Ilkay Altintas, Ph.D.
altintas@sdsc.edu

Shweta Purawat
shpurawat@sdsc.edu

San Diego Supercomputer Center, UC San Diego

Slides for Blue Waters Webinar Series, 04/26/2017

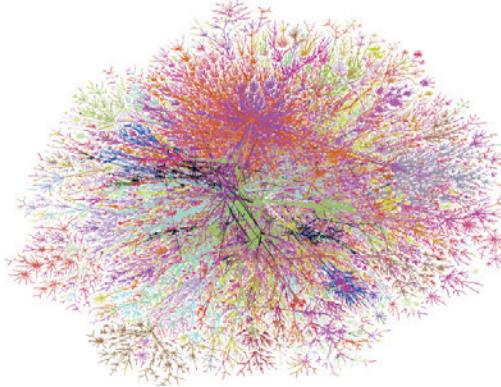
Computing Today has Many Shapes and Sizes



 **hadoop**



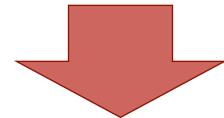
**COMPUTING AT
SCALE**



BIG DATA

Enables dynamic data-driven applications

Requires software
for dynamic coordination
and resource optimization



Workflow Systems



Smart Manufacturing

Computer-Aided Drug Discovery



Personalized Precision Medicine

Smart Cities

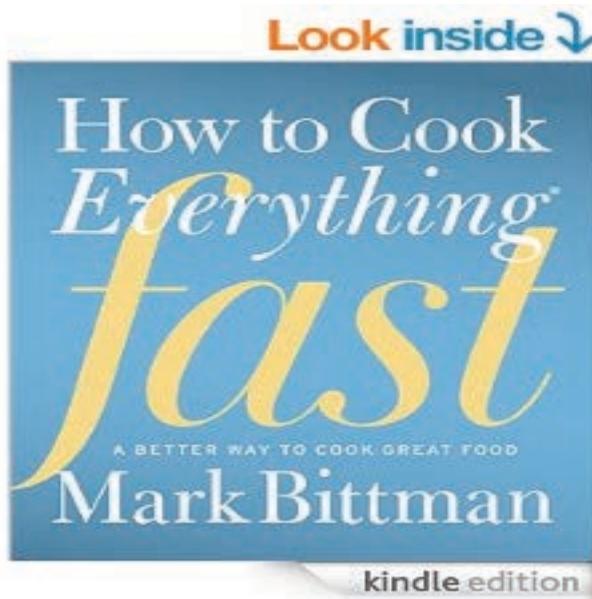


Smart Grid and Energy Management

Let's make pasta this evening!



How to Cook Everything Fast



"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](#)

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?*
- ...

"Big" Data Engineering

Find data
Access data
Acquire data
Move data

ACCESS

Clean data
Integrate data
Subset data
Pre-process data

MANAGE

Computational "Big" Data Science

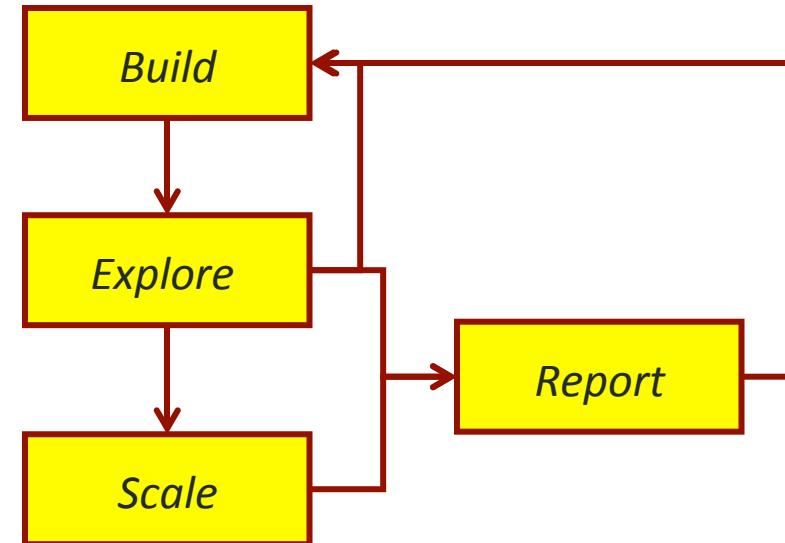
Analyze data
Process data

ANALYZE

Interpret results
Summarize results
Visualize results
Post-process results

REPORT

Many ways to look at the process for HPC and Big Data... not every step is automatable!



Facilitating and Accelerating XXX-Data Science or Comp-XXX Research using Scientific Workflows

Important Attributes



Assemble complex processing easily



Access transparently to diverse resources



Incorporate multiple software tools

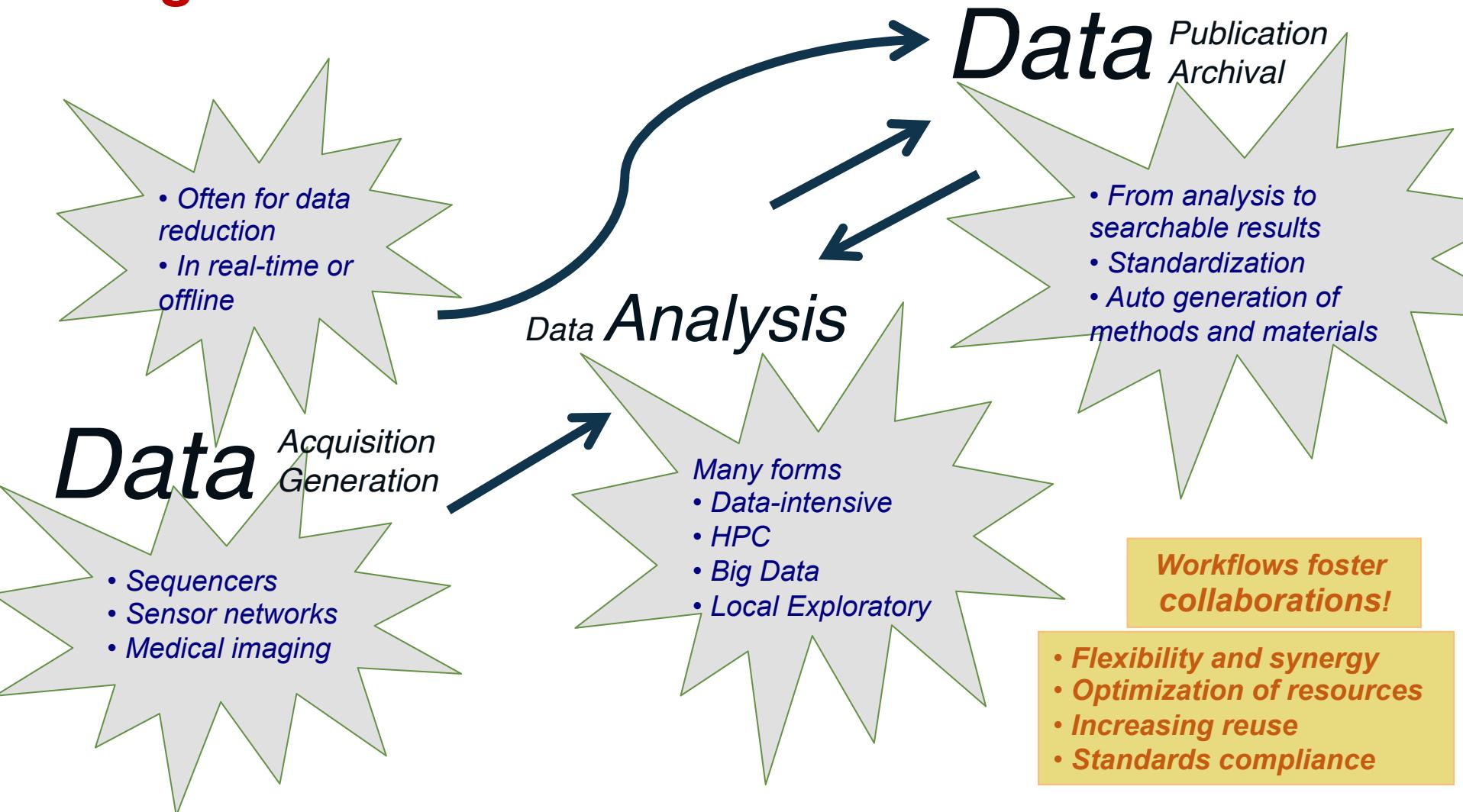


Assure reproducibility and robustness

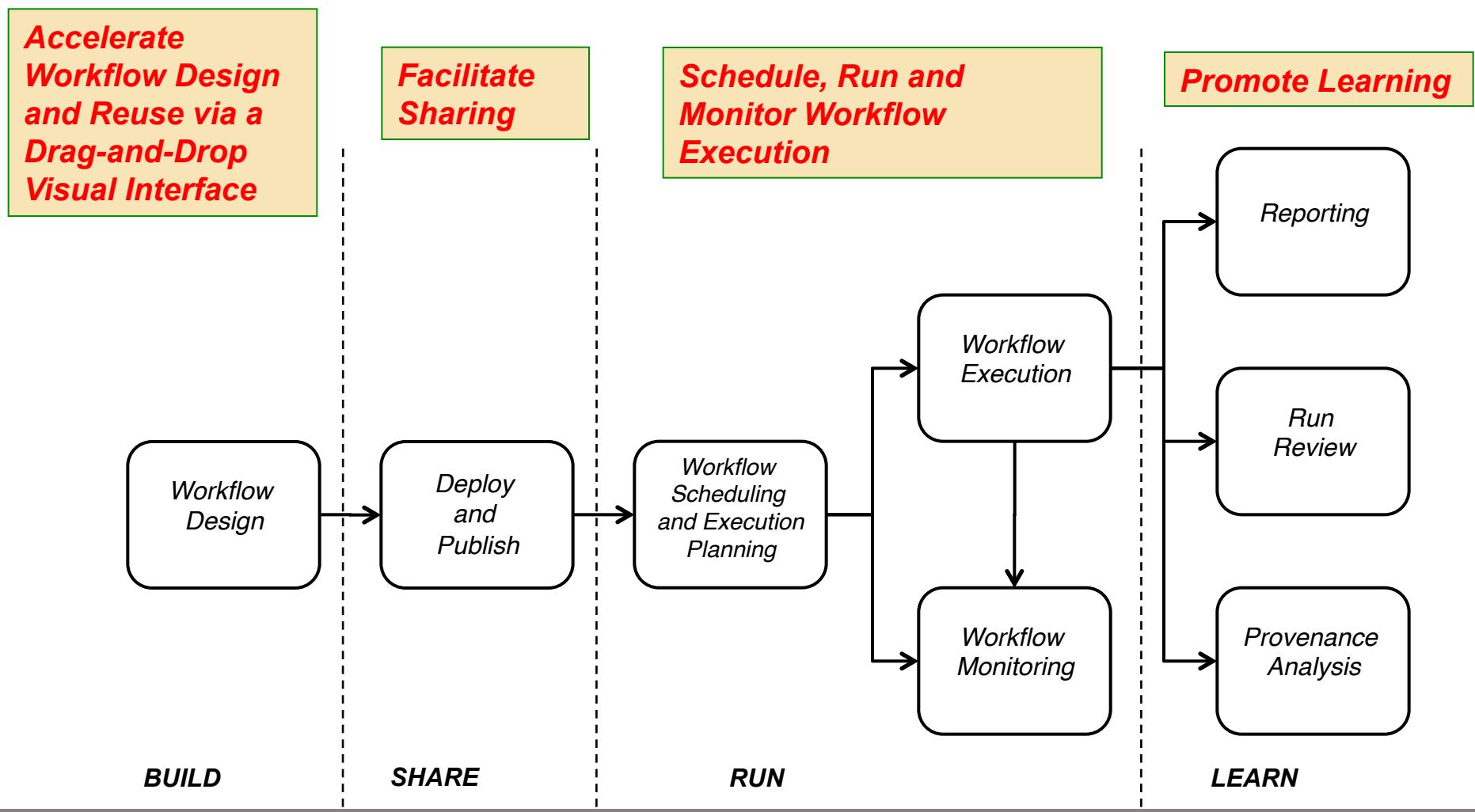


Build around community development model

Workflows are Used in These Diverse Scenarios in Biological Sciences



Workflows are a Part of Cyberinfrastructure



Support for end-to-end computational scientific process

Graphical Workflow Systems

-Toolboxes 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 in 2015
 - Frequent module release updates
- *Builds upon the open-source Ptolemy II framework*

Ptolemy II: A laboratory for investigating design

KEPLER: A problem-solving environment for Scientific Workflow

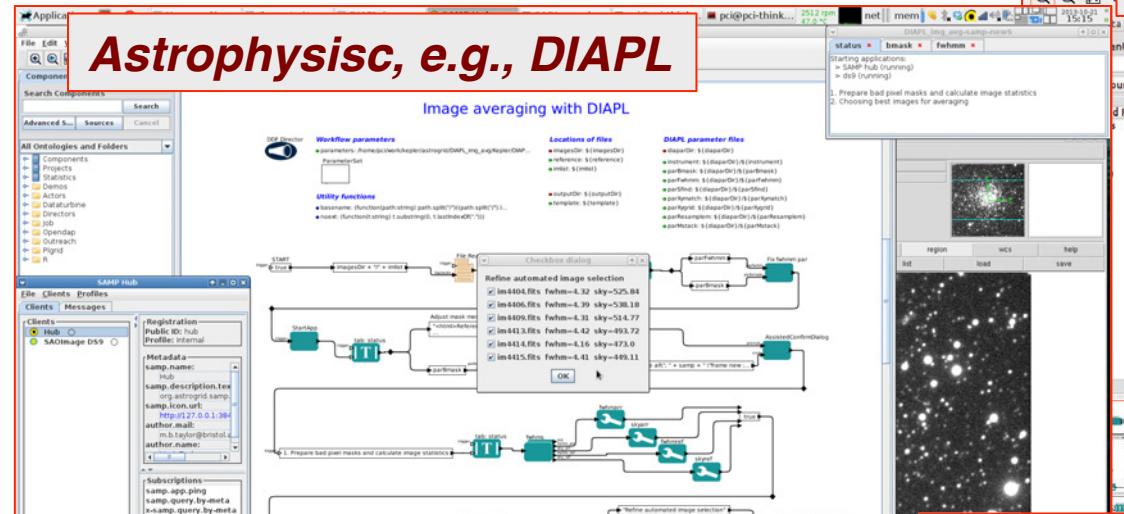
KEPLER = “Ptolemy II + X” for Scientific Workflows

Kepler can be applied to problems in different scientific disciplines: some here and many more...

Noanotechnology, e.g., ANELLI

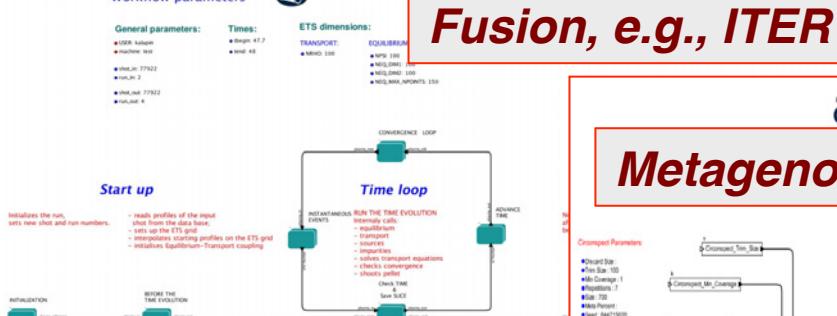
Astrophysics, e.g., DIAPL

Image averaging with DIAPL



European Transport Simulator

Workflow parameters



Fusion, e.g., ITER

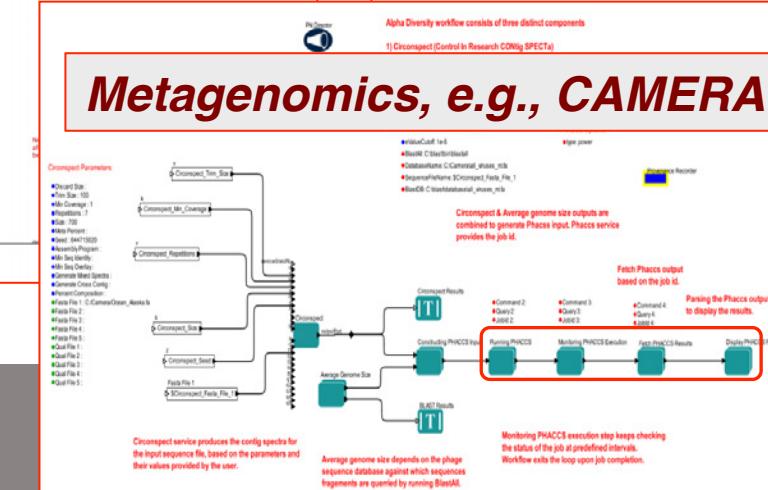
EP Model Parameters

- ConductRatioLV: '1 25 50 75 100'
- ConductBulk: '0.0001; 0.000075; 0.00005; 0.000025; 0.00001'
- ConductRatioRV: '1 25 50 75 100'
- ConductScar: '0.5 0.1 0.005 0.001'

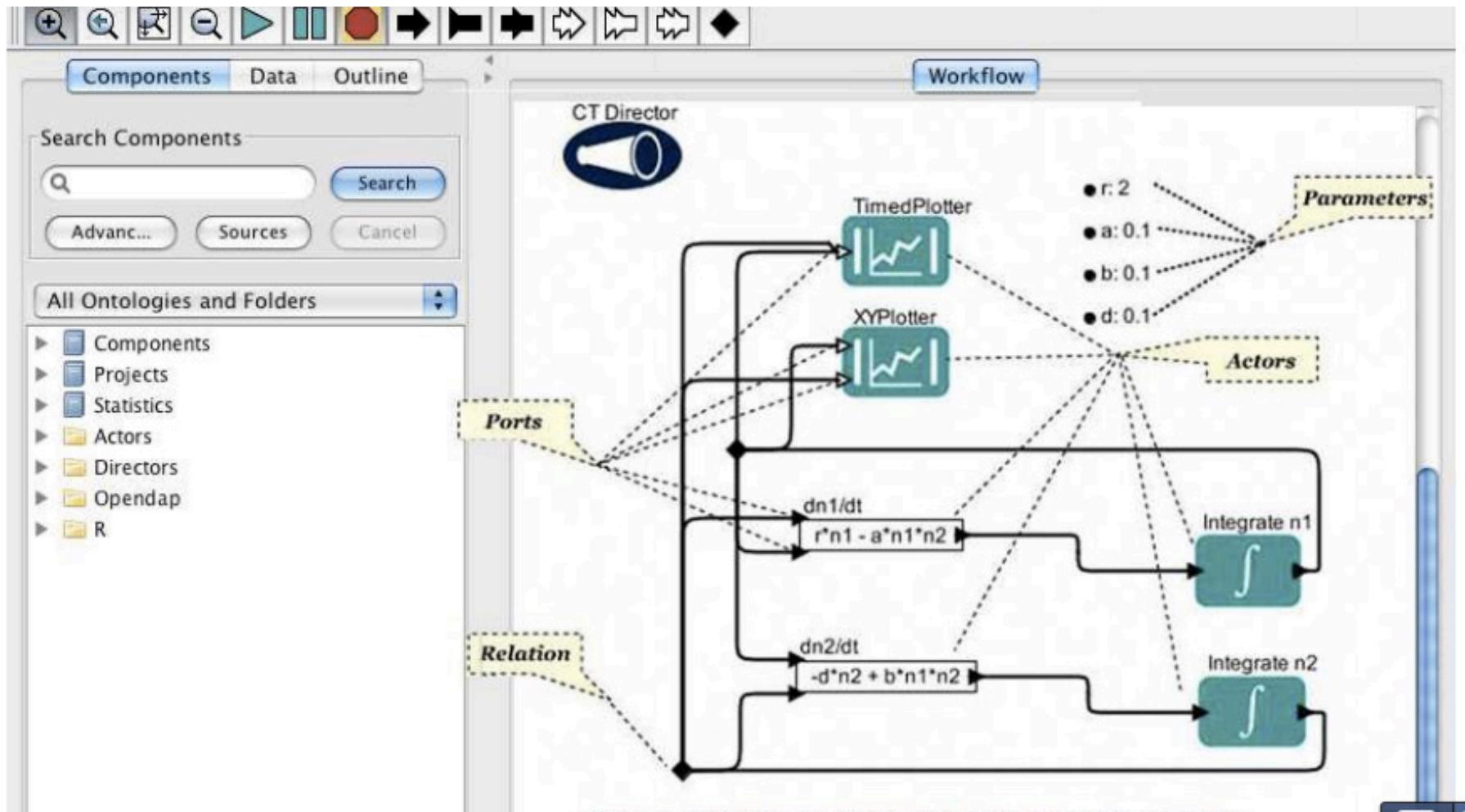
Optimizer Parameters

- GridSpacing: '10 6 8'
- Delta: 1/4
- StoppingCriteria: 1/8
- MinParamValues: '0.5 0.0 -1.3'

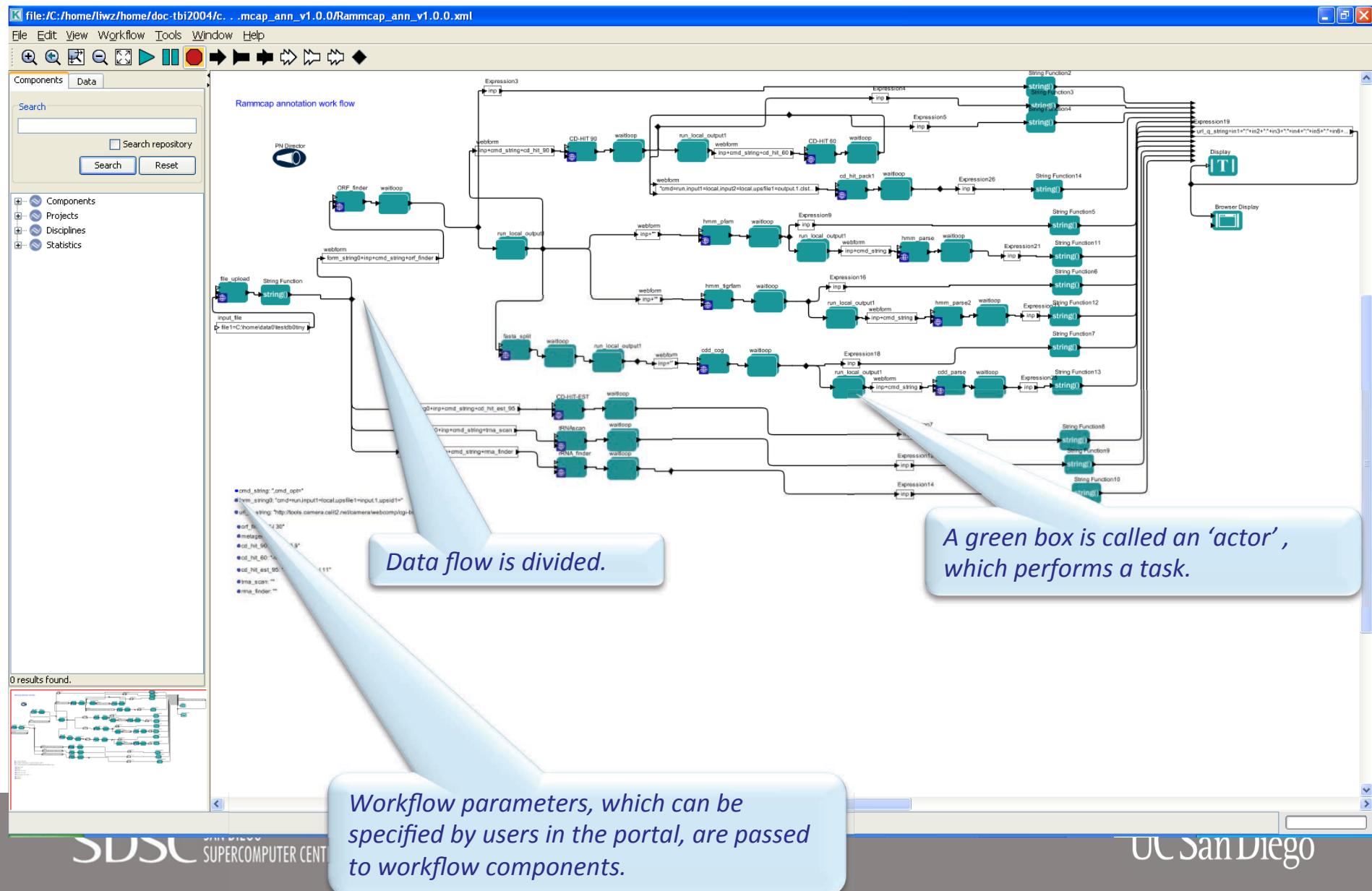
Metagenomics, e.g., CAMERA



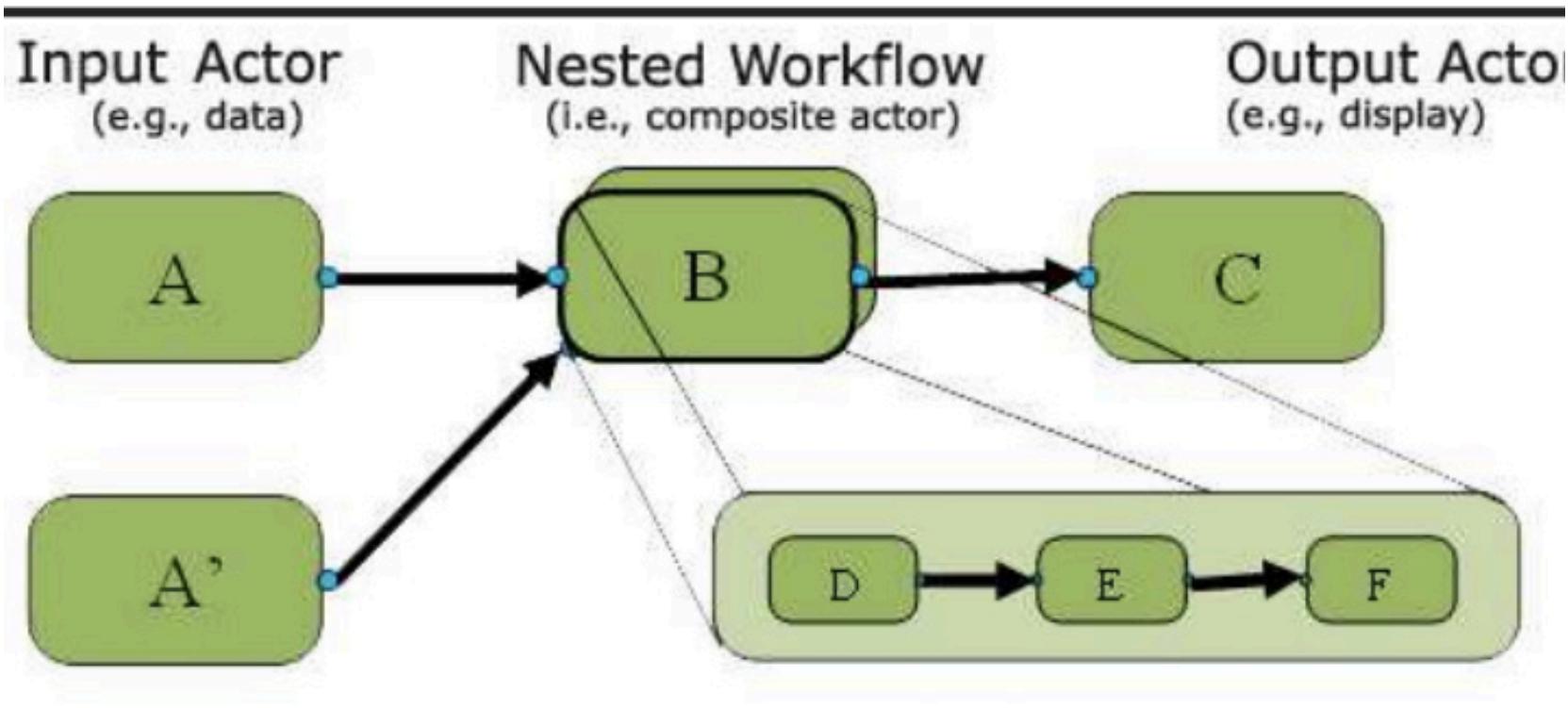
Basic components and terminologies in Kepler



A Typical Kepler Workflow



Kepler Actor Types



Some actors in place for...

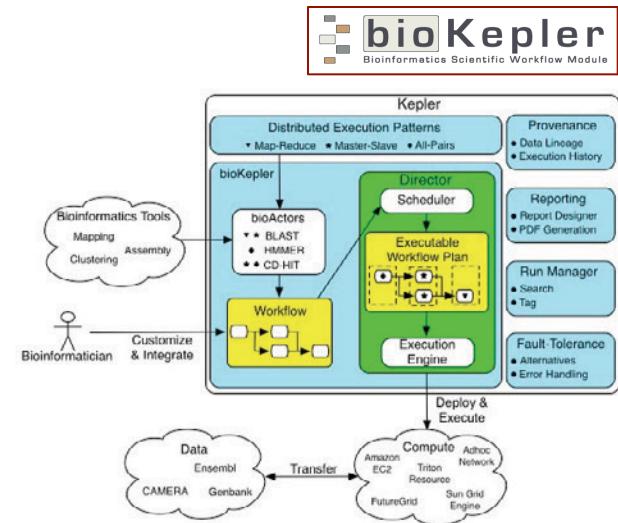
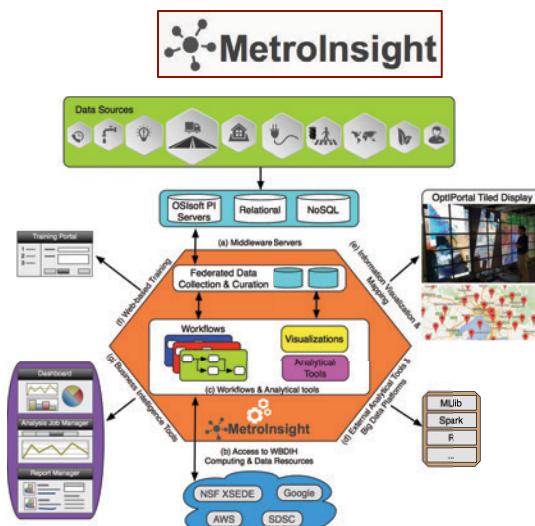
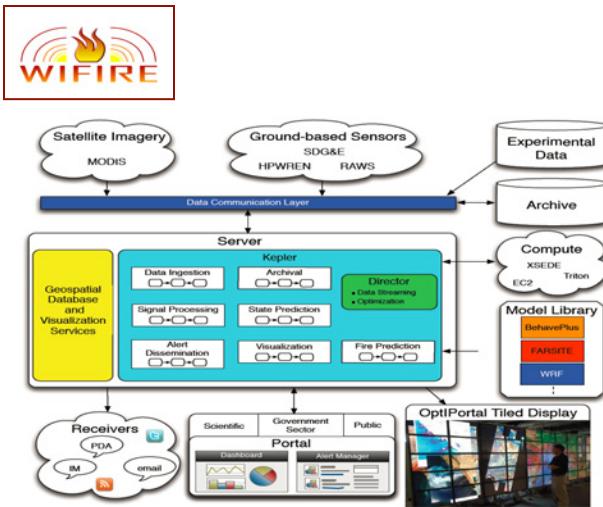
- *Command Line wrapper tools (**local execution, ssh, scp, ftp, etc.**)*
- *Generic Web Service Clients for **SOAP** and **REST***
- *A suite of **cloud computing** actors for VM instantiation and management*
- *Job management actors for **HPC**, **GPU**, **SGE** and other commodity clusters*
- *Customizable **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*
- *Open **GIS** tools*
- *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**

Some Recent Kepler Workflow Examples

Examples: Use of Workflows as an Application Integration Tool for “Big” Data and Computational Science



Using Workflows and Cyberinfrastructure for Wildfire Resilience

- A Scalable Data-Driven Monitoring and Dynamic Prediction Approach -

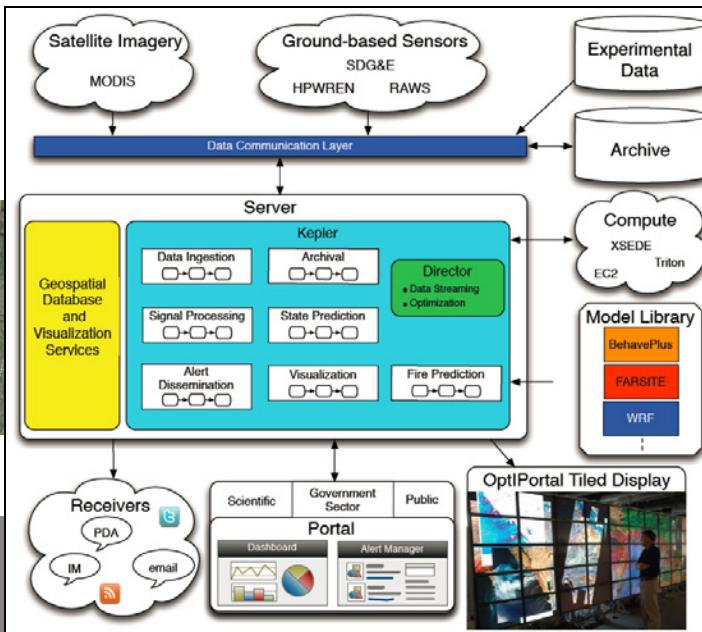




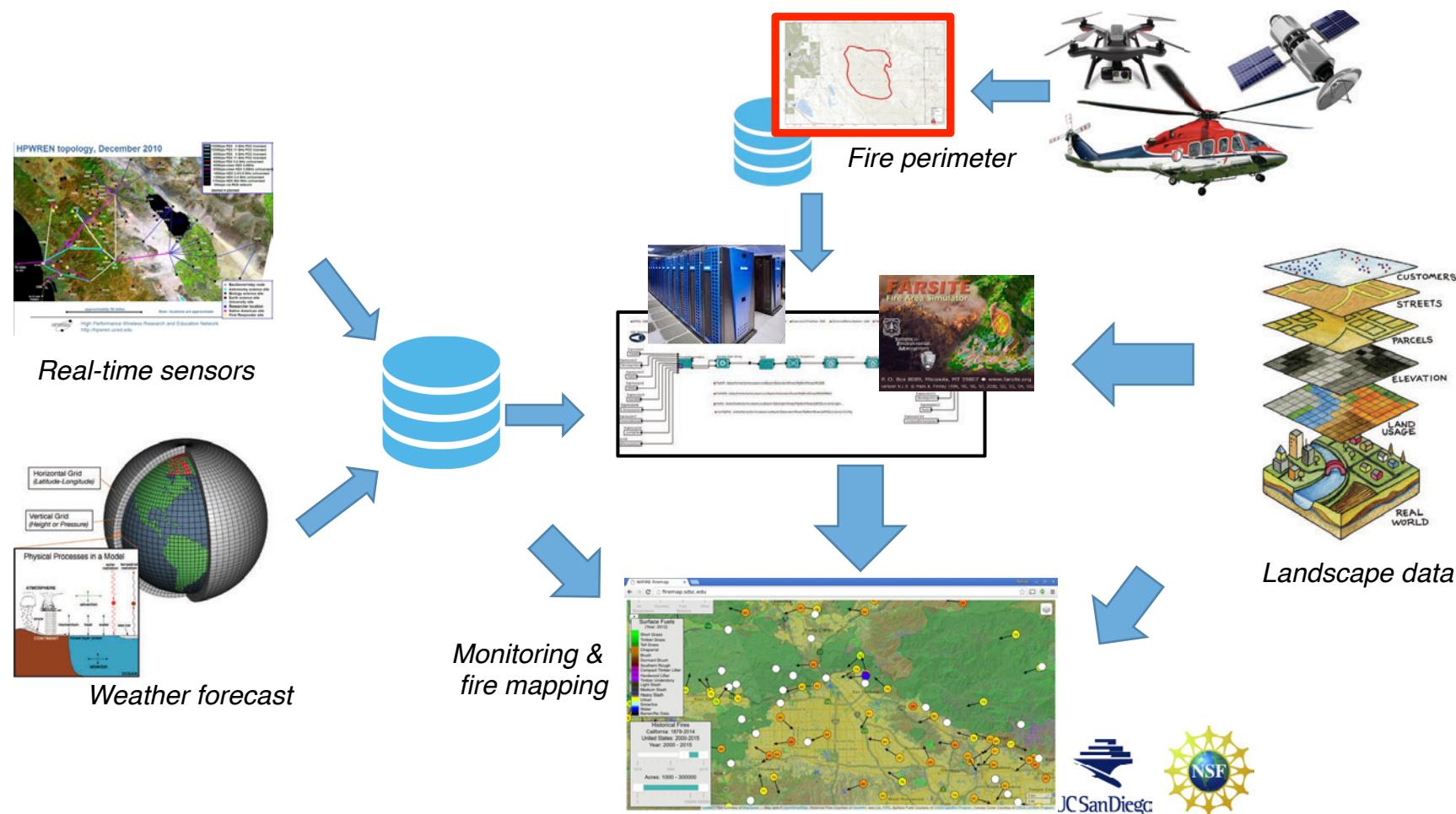
*Big
Data*



Monitoring Visualization Fire Modeling

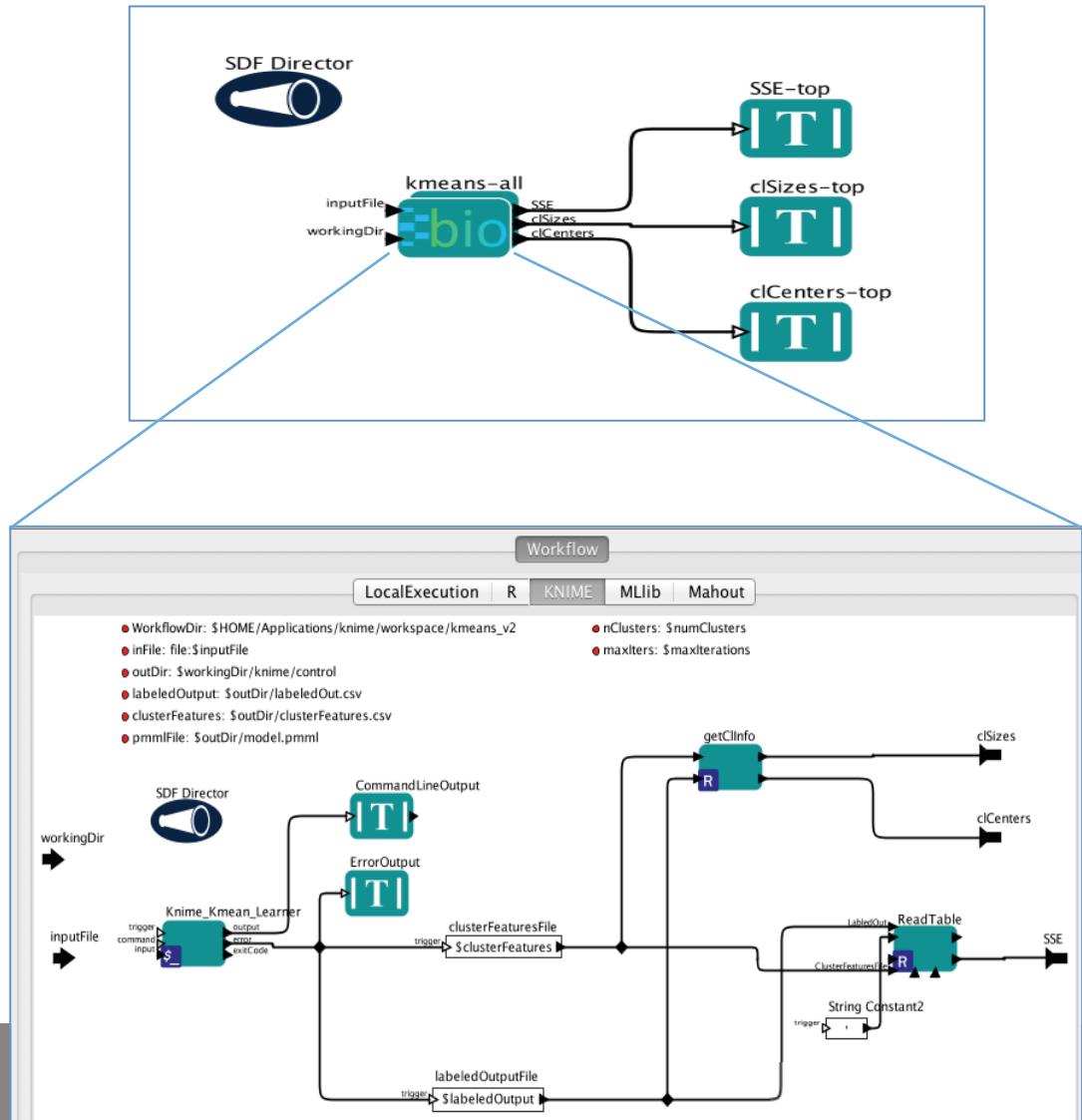


Fire Modeling Workflows in WIFIRE



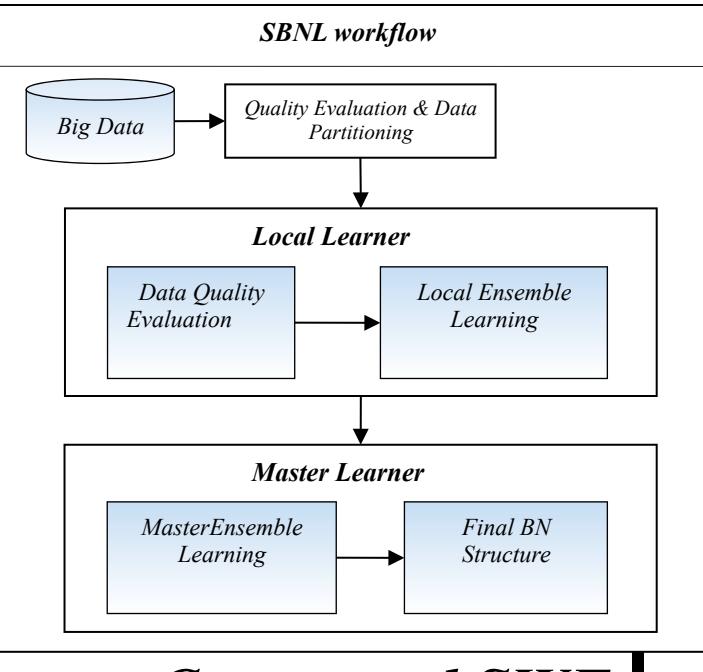
Flexible programming of K-means

- **R:** Programming language and software environment for statistical computing and graphics.
- **KNIME:** Platform for data analytics.
- **MLlib:** Scalable machine learning library running on Spark cluster computing framework
- **Mahout:** Scalable machine learning library based on MapReduce.



Scalable Bayesian Network Learning

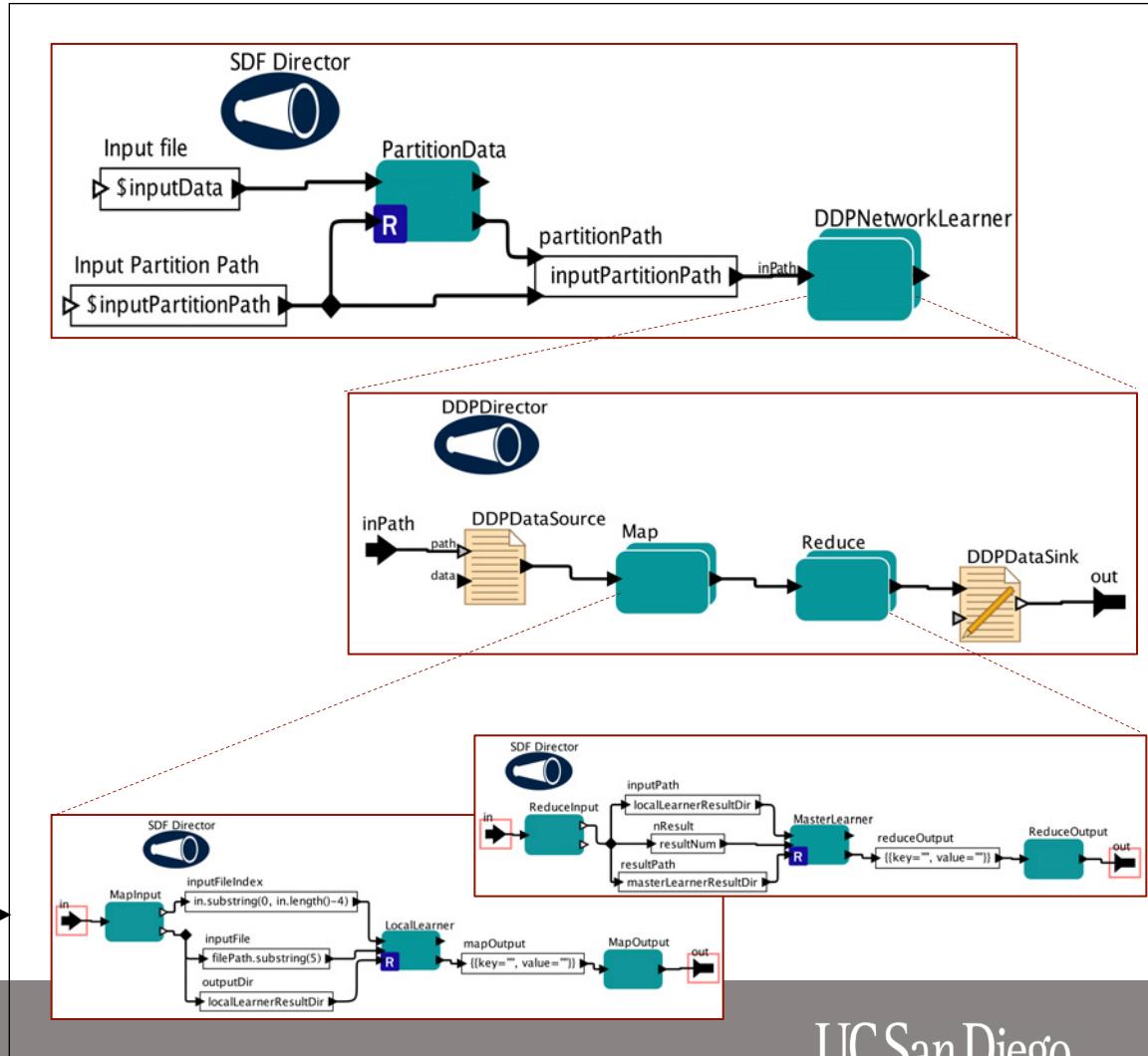
From “Napkin Drawings” to Executable Workflows...

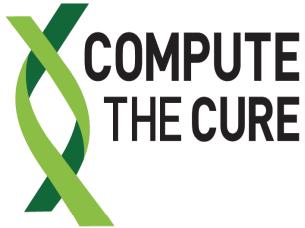


Conceptual SWF

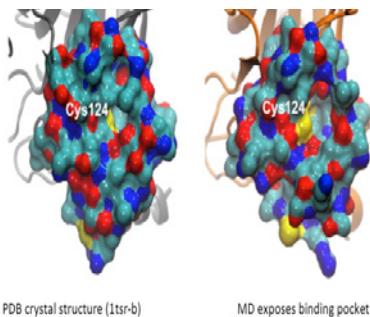
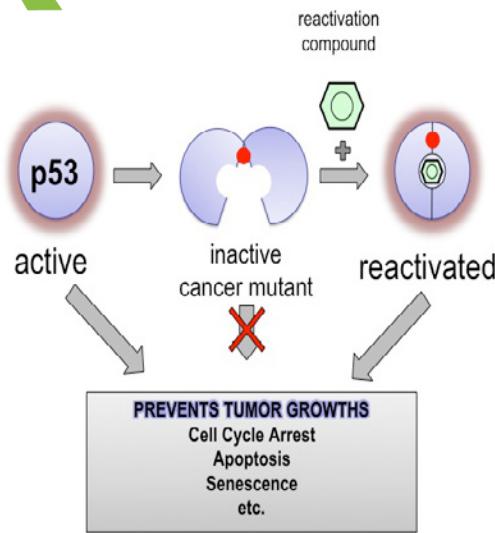
**Insurance and Traffic
Data Analytics using
Big Data Bayesian
Network Learning**

Executable SWF



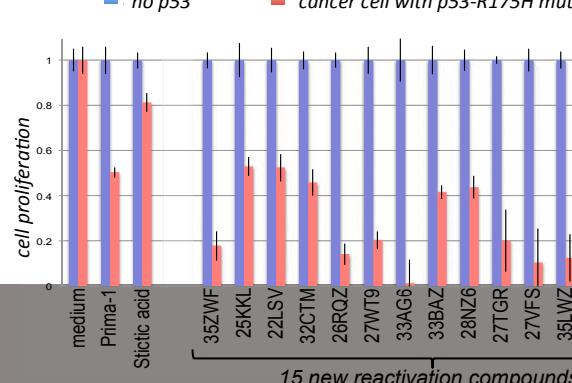
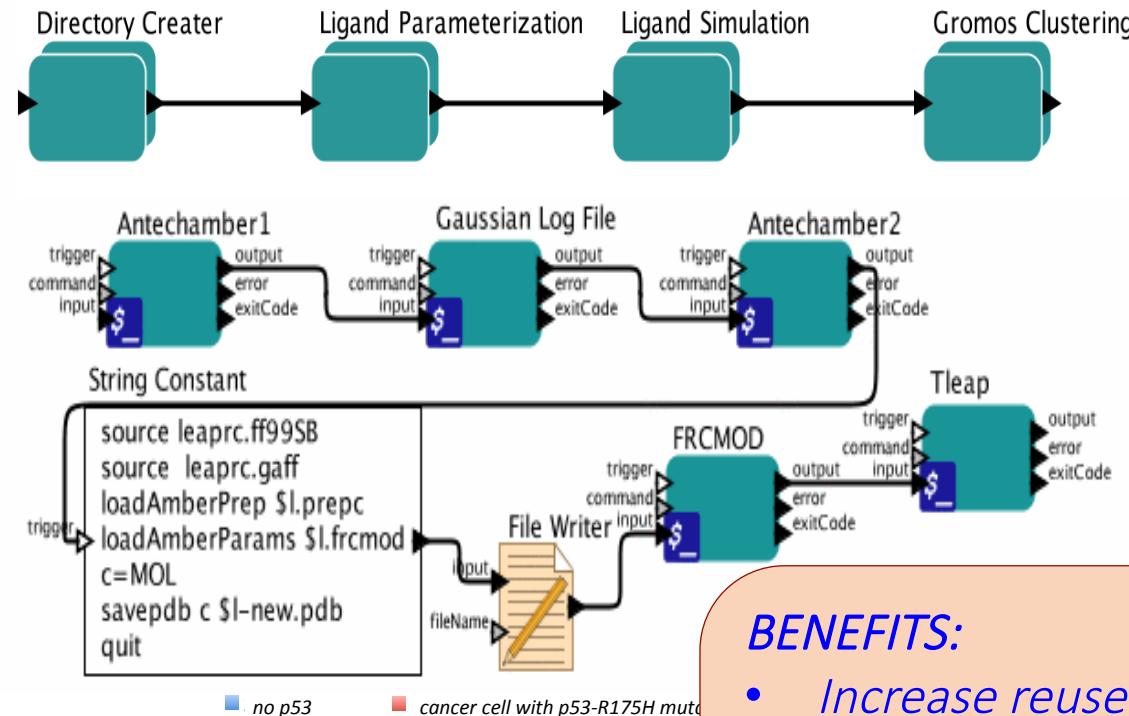


Scalable Drug Discovery



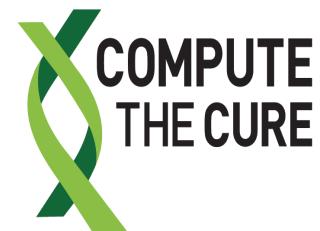
Closed in PDB crystal structure (1tsr-b)

MD exposes binding pocket

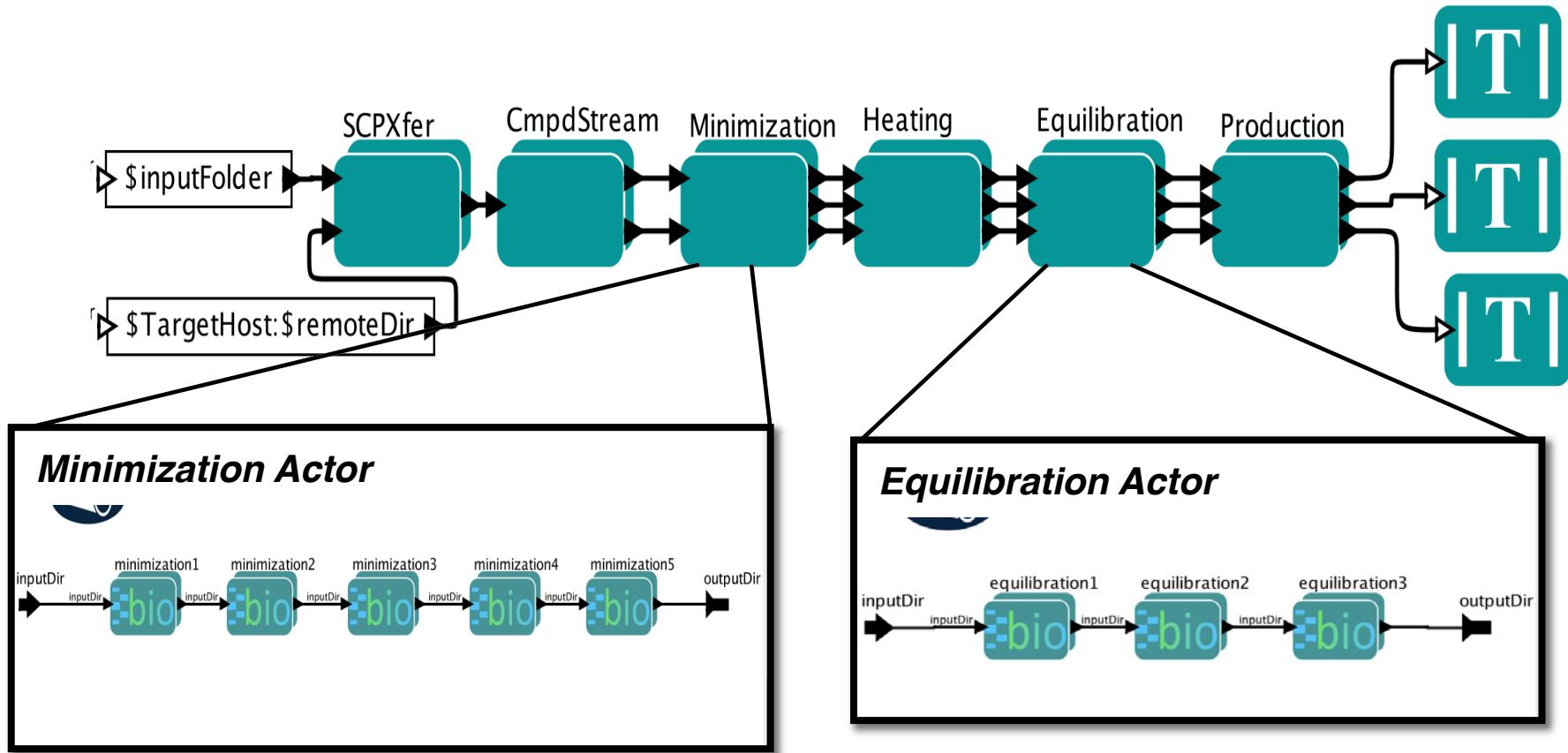


BENEFITS:

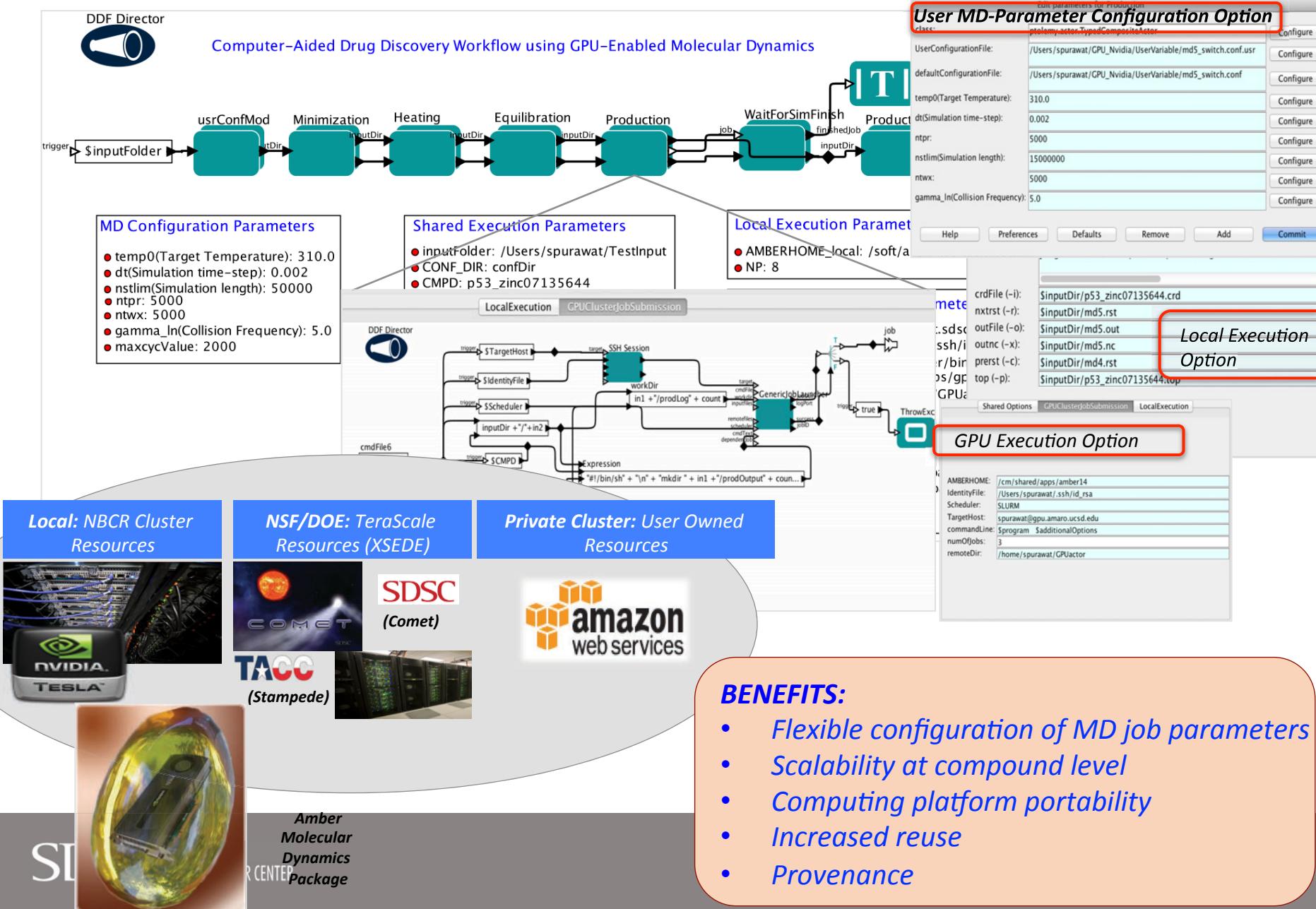
- Increase reuse
- Reproducibility
- Scale execution, problem & solution
- Compare methods
- Training



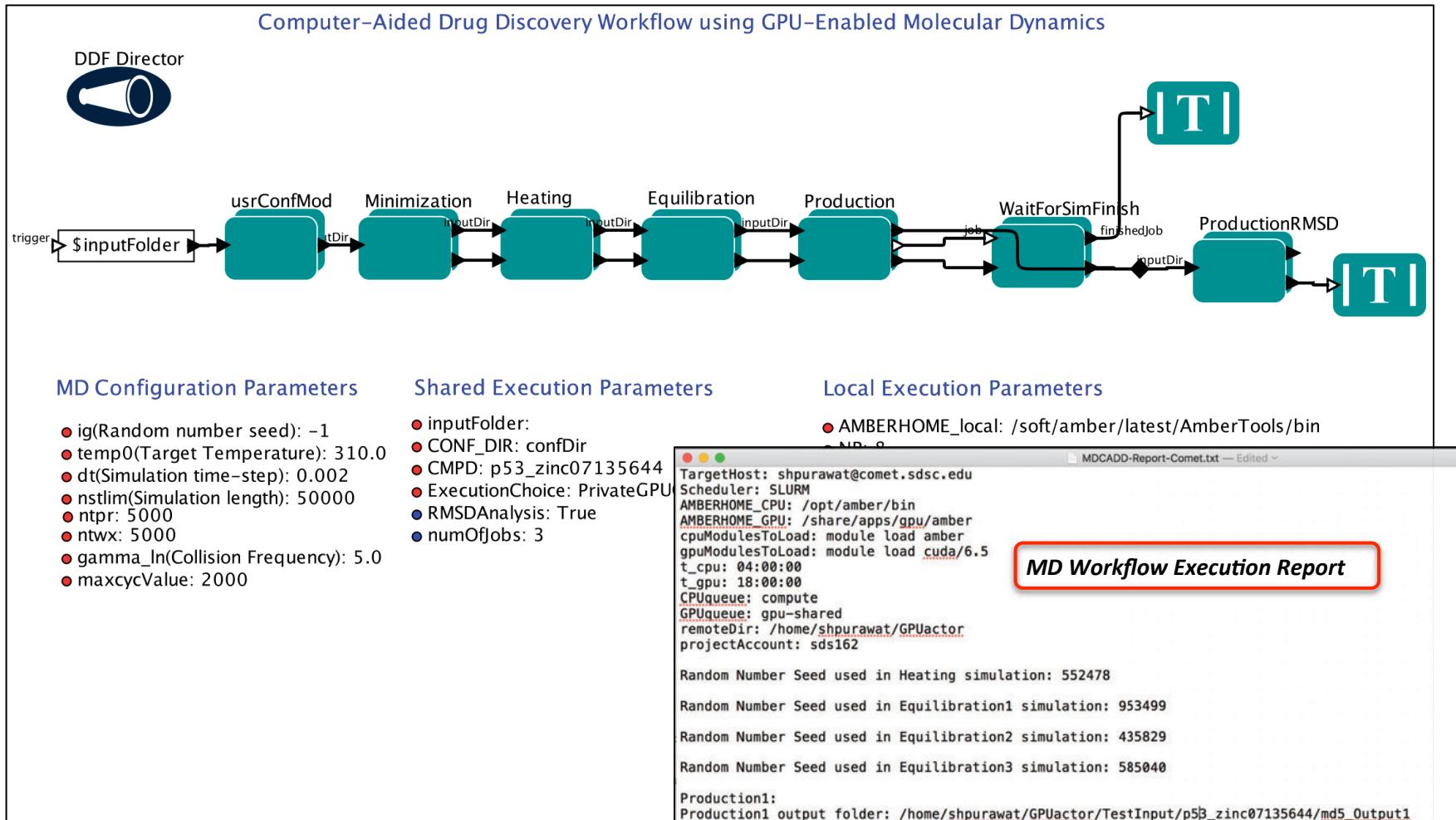
AMBER GPU Molecular Dynamics Workbench



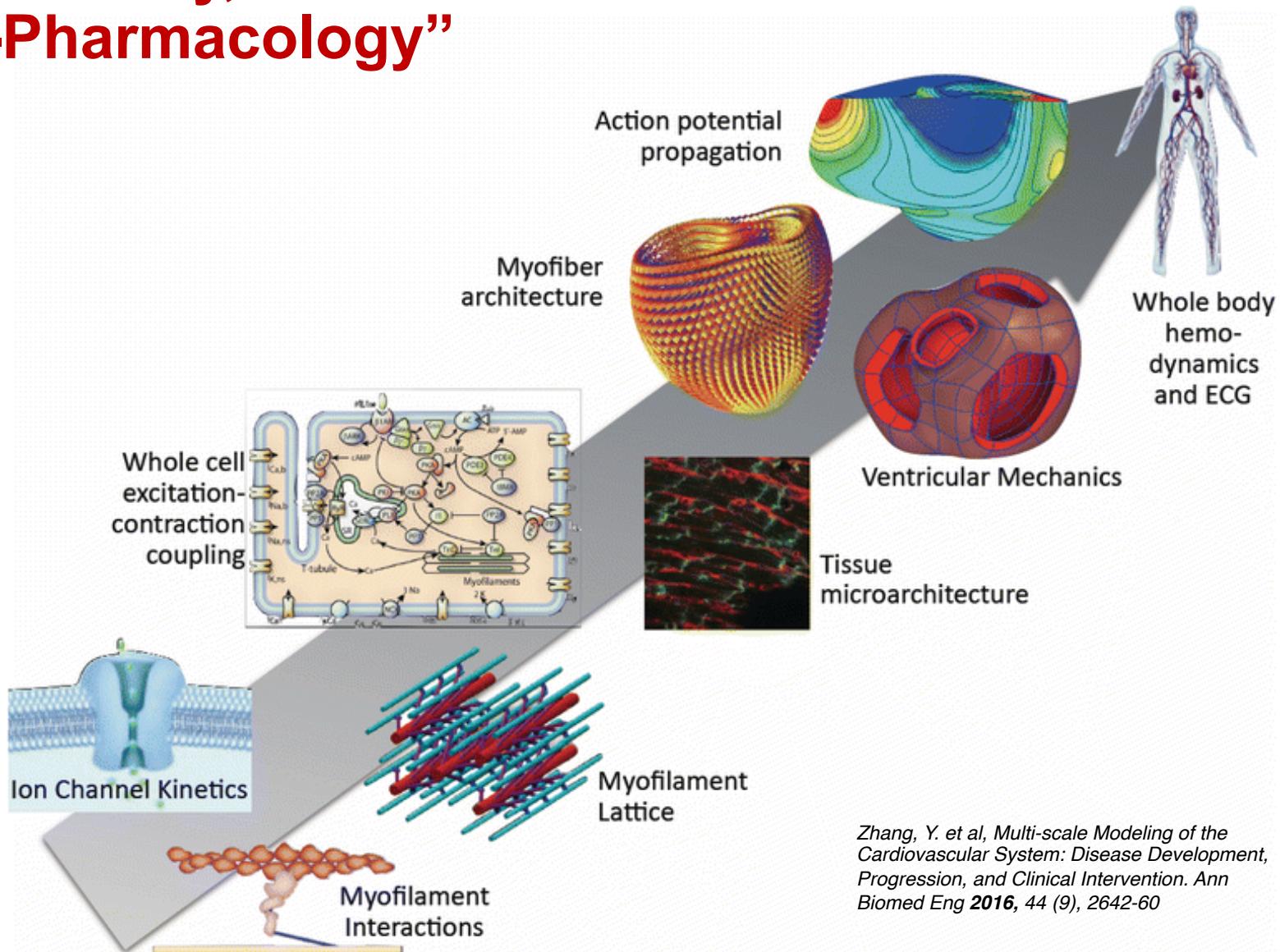
Parametric execution of each step



MDCADD WF – Workflow Execution Report



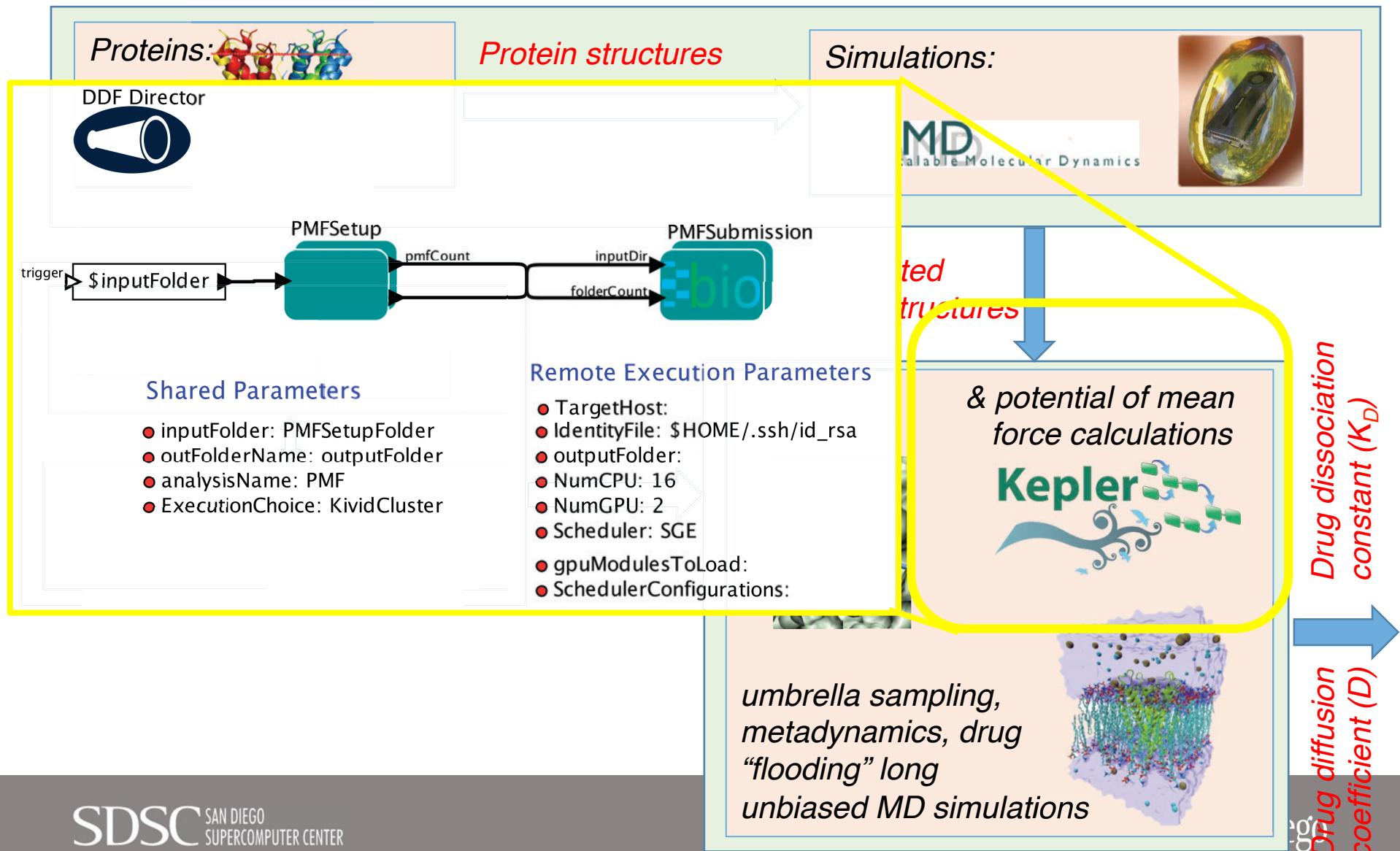
Colleen Clancy, “Predictive Multiscale in silico Cardio-Pharmacology”



Zhang, Y. et al, *Multi-scale Modeling of the Cardiovascular System: Disease Development, Progression, and Clinical Intervention*. *Ann Biomed Eng* 2016, 44 (9), 2642-60

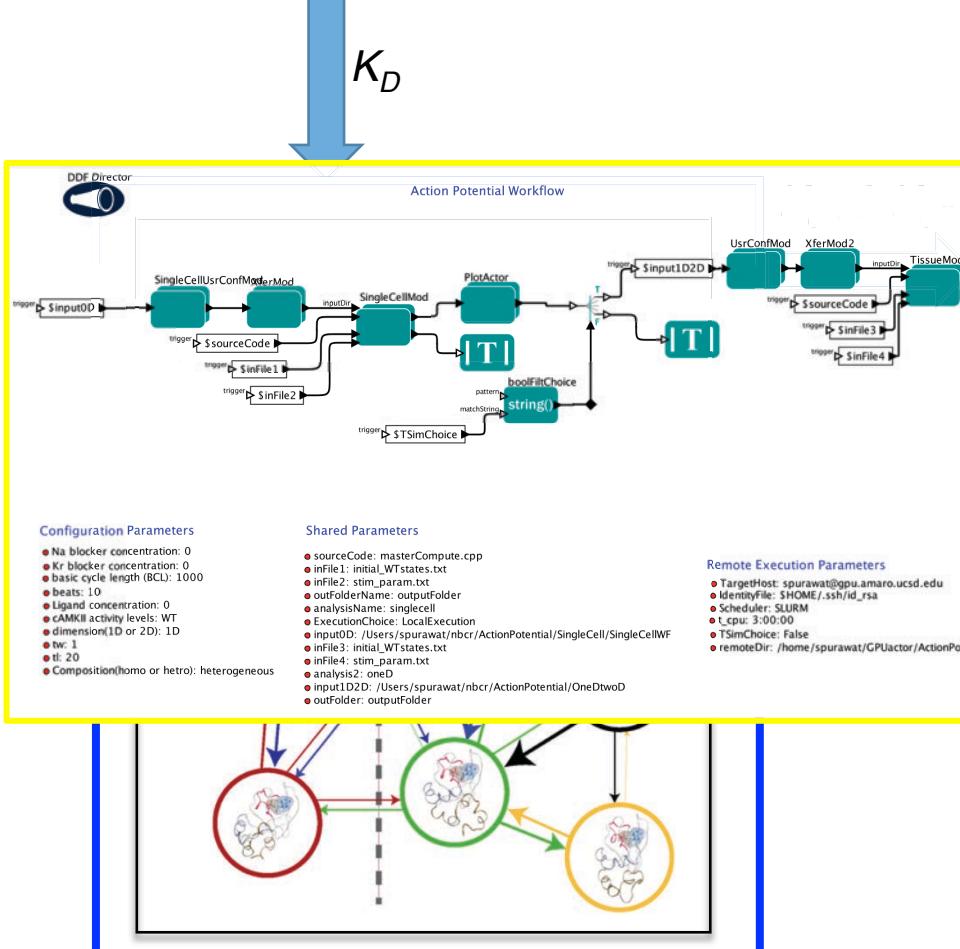
A flow-chart of proposed multi-scale simulations.

Part 1: Atomistic level simulations.



Atomistic level simulations

Part 2: Functional level simulations



Cell simulations

Markov model into cell
• validations
Generates
“populations” cells)



Tissue simulations Kepler

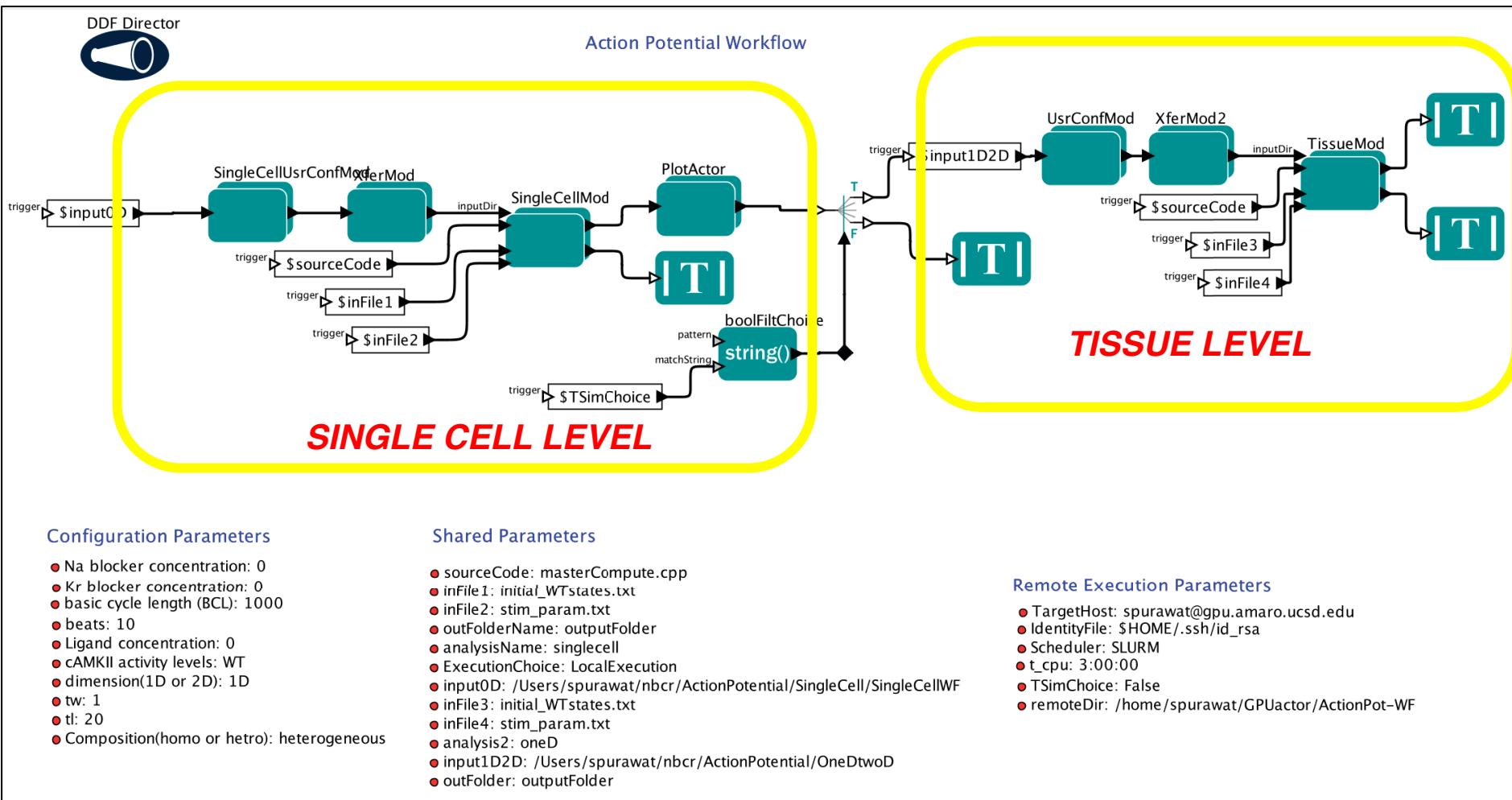
Cable/tissue maps
Compare to experiments



Pseudo ECG

Compare to clinical
ECG

Cell and Tissue Simulation Workflow



Integration with Scientific Process and Communication

Aim: Integrate the workflow-driven, scalable and reproducible science with notebooks and exploratory nature of the scientific process.

jupyter Notebooks as Workflow Execution and Reporting Interfaces

- ✓ Jupyter Kepler Magics:

https://github.com/words-sdsc/Jupyter_Kepler_Integration/blob/master/README.md

NBCR Demo Kepler from Jupyter Last Checkpoint: a minute ago (autosaved)

View Insert Cell Kernel Widgets Help

CellToolbar

RUNNING MDCADD WORKFLOW IN JUPYTER NOTEBOOKS

Kepler-MDCADD Workflow from Jupyter

Run Kepler without leaving the notebook:

```
import KeplerMagicFunction

%KpConf /Users/spurawat/Kepler_Repository/bioKepler-trunk-Dec2015/kepler.modules/kepler.sh

%WpConf /Users/spurawat/nbcr/rocce-NBCR-Product/MDCADD.xml

%Kepler
```

SDSC SUPERCOMPUTER CENTER

UC San Diego

jupyter Notebooks as Workflow Execution and Reporting Interfaces

jupyter NBCR Demo Kepler from Jupyter Last Checkpoint: 3 minutes ago (autosaved) Logout

File Edit View Insert Cell Kernel Widgets Help Kernel

Plot out graphs from Kepler execution:

```
In [12]: graph = '/Users/spurawat/NBCR_Demo_Feb17/IPython-Kepler-Magic-Function/p53_zinc07135644/prod2_total_energy.png'
from IPython.display import Image
Image(filename=graph)
```

Out[12]:

Total energy of Production2

VIEWING WORKFLOW RESULTS IN JUPYTER NOTEBOOKS

jupyter Notebooks as Workflow Execution and Reporting Interfaces

jupyter NBCR Demo Kepler from Jupyter Last Checkpoint: 3 minutes ago (autosaved) [Logout](#)

File Edit View Insert Cell Kernel Widgets Help

Kernel O

In [6]: `output = open('MDCADD-Report.txt','r')
print(output.read())`

TargetHost: shpurawat@comet.sdsc.edu
Scheduler: SLURM
AMBERHOME_CPU: /opt/amber/bin
AMBERHOME_GPU: /share/apps/gpu/amber
cpuModulesToLoad: module load amber
gpuModulesToLoad: module load cuda/6.5
t_cpu: 04:00:00
t_gpu: 18:00:00
CPUqueue: compute
GPUqueue: gpu-shared
remoteDir: /home/shpurawat/GPUactor
projectAccount: sds162

Random Number Seed used in Heating simulation: 552478
Random Number Seed used in Equilibration1 simulation: 953499
Random Number Seed used in Equilibration2 simulation: 435829
Random Number Seed used in Equilibration3 simulation: 585040

Production1:
Production1 output folder: /home/shpurawat/GPUactor/TestInput/p53_zinc07135644/md5_Output1
Random Number Seed used in Production1 simulation: 241619

Production3:
Production3 output folder: /home/shpurawat/GPUactor/TestInput/p53_zinc07135644/md5_Output3
Random Number Seed used in Production3 simulation: 740887

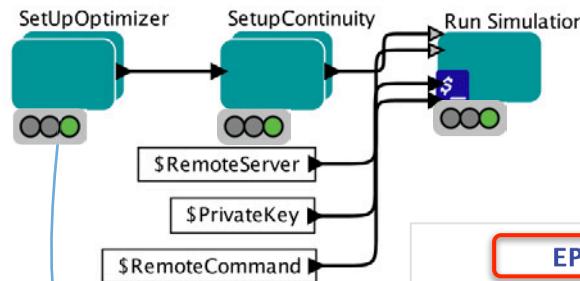
Amber Version:
pmemd.cuda: Version 14.0

nvcc: NVIDIA (R) Cuda compiler driver
Copyright (c) 2005-2014 NVIDIA Corporation
Built on Thu Jul 17 21:41:27 CDT 2014

GENERATING WORKFLOW EXECUTION REPORTS IN JUPYTER NOTEBOOKS

Run Parameters

- RunId: "Cond-1"
- LocalWorkingDirectory: "/Users/jeffvandorn/work/myStuff/psm_workflow/"
- StimulusLocationsFile: "stimLocations-RVendo-38"
- RemoteServer: "rocce.ucsd.edu"
- RemoteWorkingDir: "/home/jvandorn/psm_workflow/"
- RemoteCommand: "python runEPbatch.py"
- PrivateKey: "/Users/jeffvandorn/.ssh/id_rsa"



*Details of
SetUpOptimizer*

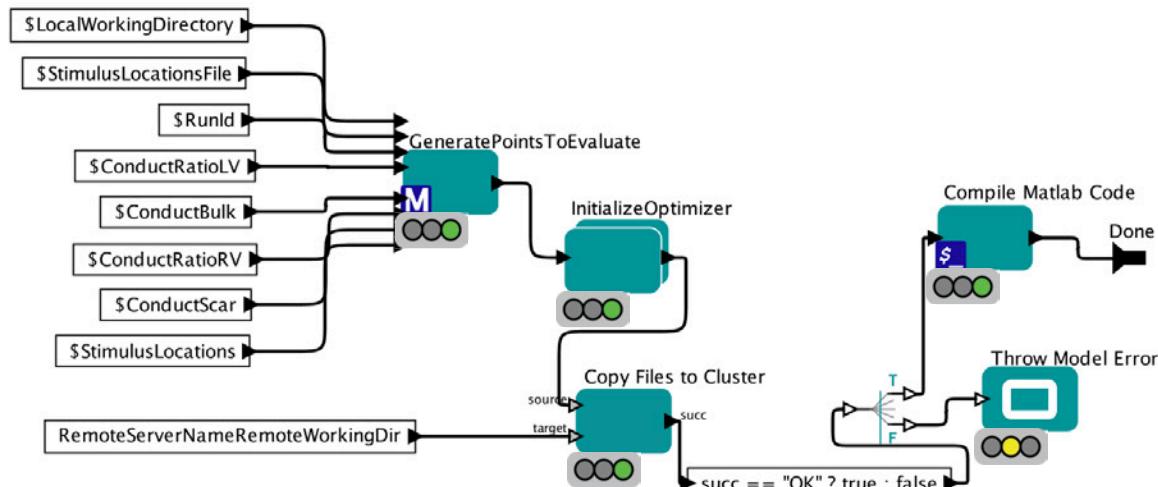


EP Model Parameters

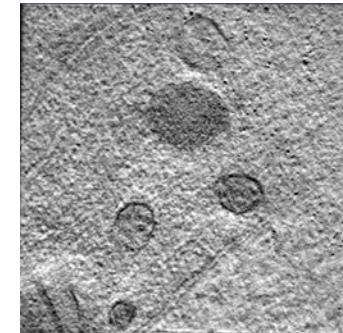
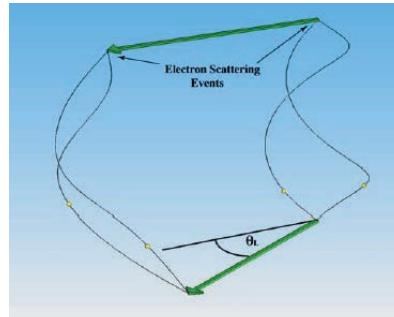
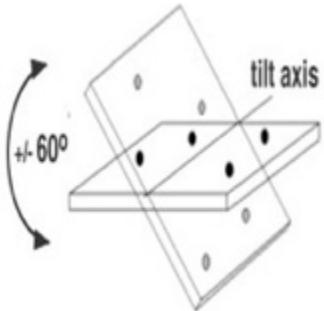
- ConductRatioLV: '1 25 50 75 100'
- ConductBulk: '0.0001; 0.000075; 0.00005; 0.000025; 0.00001'
- ConductRatioRV: '1 25 50 75 100'
- ConductScar: '0.5 0.1 0.005 0.001'
- StimulusLocations: 50

Optimizer Parameters

- GridSpacing: '10 6 8'
- Delta: 1/4
- StoppingCriteria: 1/8
- MinParamValues: '0.5 0.0 -1.3'
- MaxParamValues: '1.25 1.5 1.3'

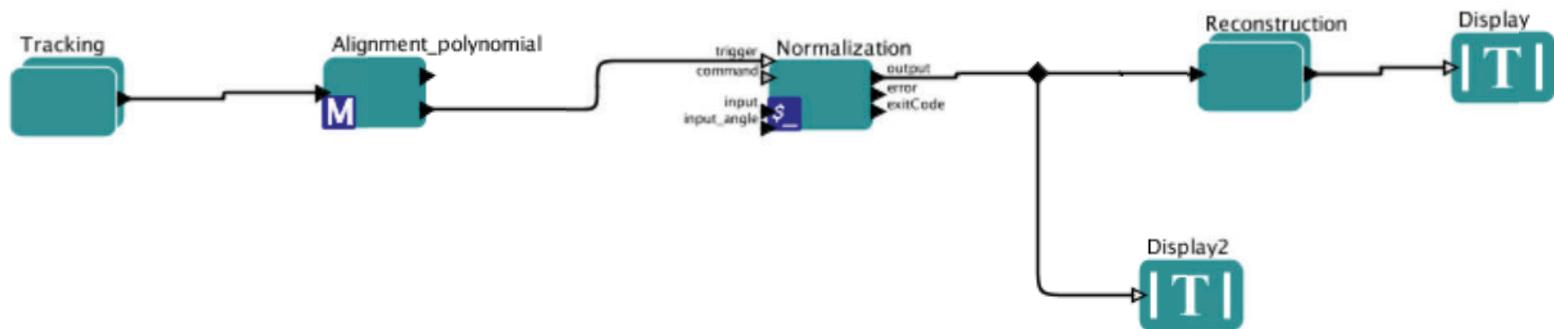


Example: Integration of Software Tools for Iterative Reconstruction



*From “Conceptual Steps”
To Integrated, Executable and Portable
Workflows*

Electron Microscope Processing in Kepler (EPiK) Workflow for TxBR Streamlines Iterative Reconstruction



*In Partnership with
Sun Yat Sen Foundation*

BENEFITS:

- *Improved programmability*
- *Ability to use and compare multiple methods in one integrated workflow*
- *Reduction of complexity*

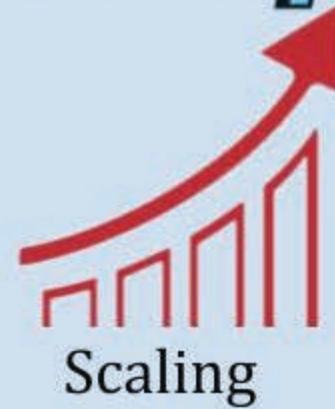
Kepler workflows are used for:

- **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

Kepler Demo!



PBS



Apache Flink

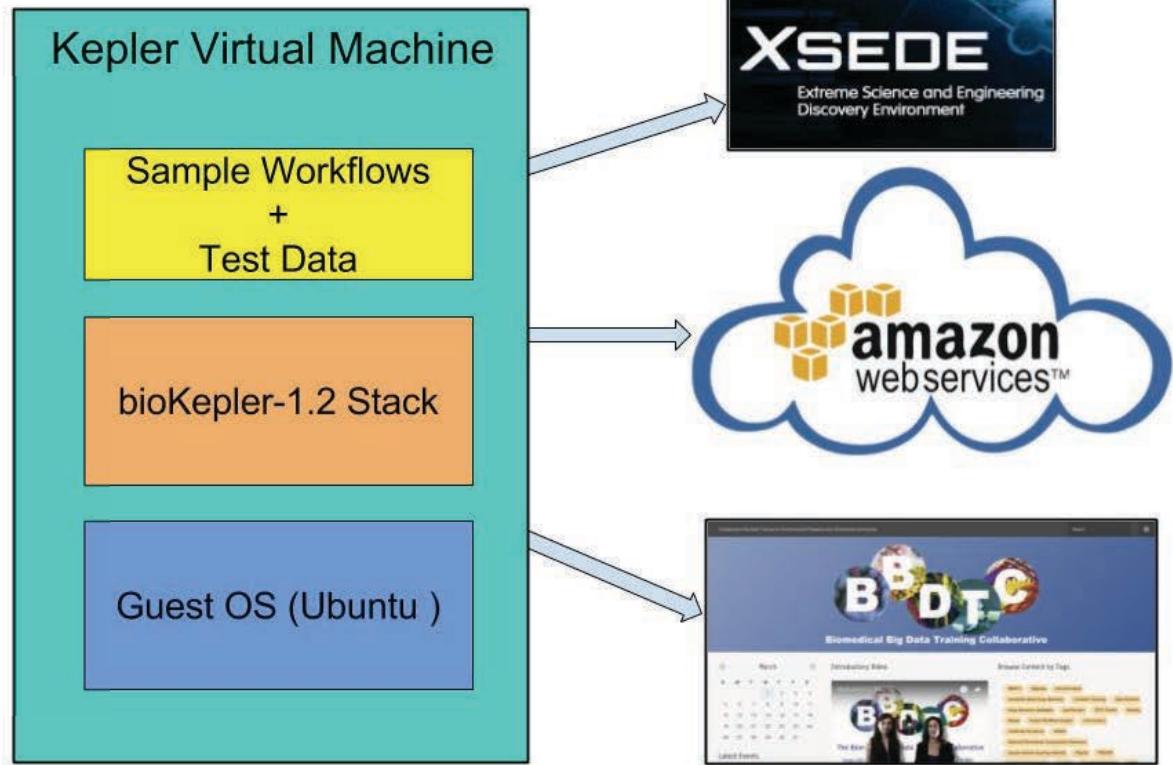


SGE

???



Kepler – XSEDE Architecture



<https://biobigdata.ucsd.edu/>

Downloading and Installing Kepler

- Java 1.8 or later is a prerequisite to run Kepler 2.5
- Download Kepler installer for Windows, Mac, or Linux-based platforms.

<https://kepler-project.org/users/downloads>

- Add-on Modules can be added to Kepler to enhance its functionality. Example: bioKepler, Provenance, Reporting.

Useful Links

- <https://kepler-project.org/users/downloads>
- <https://kepler-project.org/users/documentation>
- <http://www.biokepler.org/userguide>
- <https://biobigdata.ucsd.edu/>
- <http://words.sdsc.edu>
- [https://github.com/words-sdsc/Jupyter Kepler Integration](https://github.com/words-sdsc/Jupyter_Kepler_Integration)
- <https://words.sdsc.edu/publications>

Please use the following reference:

- To cite Kepler:

Kepler: an extensible system for design and execution of scientific workflows.
Altintas, I. and Berkley, C. and Jaeger, E. and Jones, M. and Ludascher, B. and
Mock, S. (2004) Scientific and Statistical Database Management, 2004.
Proceedings. 16th International Conference on. Pages: 423--424.

- To cite bioKepler:

Altintas, J. Wang, D. Crawl, W. Li, "Challenges and approaches for distributed workflow-driven analysis of large-scale biological data", in: Proceedings of the Workshop on Data analytics in the Cloud at EDBT/ICDT 2012 Conference, DanaC2012, 2012, pp 73-78.

- To cite Biomedical Big Data Training Collaborative (BBDTC):

S. Purawat, C. Cowart, R. E. Amaro, and I. Altintas, "Biomedical Big Data Training Collaborative (BBDTC): An effort to bridge the talent gap in biomedical science and research", Journal of Computational Science, March 2017. <http://dx.doi.org/10.1016/j.jocs.2017.03.010>

Please use the following reference:

- To cite Molecular Dynamics Computer Aided Drug-Discovery Workflow:

S. Purawat, P. leong, R. Malmstrom, G. Chan, R. Walker, I. Altintas, and R. Amaro,
“A Kepler Workflow Tool for Reproducible Molecular Dynamics”, Biophysical
Journal, 2017

- To cite WIFIRE:

Altintas I., Block J., de Callafon R., Crawl D., Cowart C., Gupta A., Nguyen M., Braun H.W., Schulze J., Gollner M., Trouve A., Smarr L., Towards an Integrated Cyberinfrastructure for Scalable Data-Driven Monitoring, Dynamic Prediction and Resilience of Wildfires. In Proceedings of the Workshop on Dynamic Data-Driven Application Systems (DDDAS) at the 15th International Conference on Computational Science (ICCS 2015). doi:10.1016/j.procs.2015.05.296

Crawl, D., Block, J., Lin, K., Altintas, I., Firemap: A Dynamic Data-Driven Predictive Wildfire Modeling and Visualization Environment, In Proceedings of the Workshop on Urgent Computing (UC) at the 17th International Conference on Computational Sciences (ICCS 2017), 2017.

Please use the following reference:

- To cite Machine Learning Integration in the Kepler

Nguyen, M., Crawl, D., Masoumi, T., Altintas, T., Integrated Machine Learning in the Kepler Scientific Workflow System, In proceedings of the Third International Workshop on Advances in the Kepler Scientific Workflow System and Its Applications at the International Conference on Computational Science (ICCS 2016). doi:10.1016/j.procs.2016.05.545

- To cite Scalable Bayesian Network Learning

Wang J., Tang Y., Nguyen M., Altintas I., A Scalable Data Science Workflow Approach for Big Data Bayesian Network Learning. Accepted by the International Symposium on Big Data Computing (BDC 2014)

- To cite EPiK - Electron Tomography Kepler Workflow:

Chen R., Wan X., Lawrence A., Wang J., Crawl D., Phan S., Altintas. I, Ellisman M., EPiK - a Workflow for Electron Tomography in Kepler. In Proceedings of the Second International Workshop on Advances in the Kepler Scientific Workflow System and Its Applications at the 14th International Conference on Computational Science (ICCS 2014).

Workflows for Data Science Center of Excellence

- Exploring Programmable, Reusable and Reproducible Scalability using Workflows-

Real-Time Hazards Management
wifire.ucsd.edu

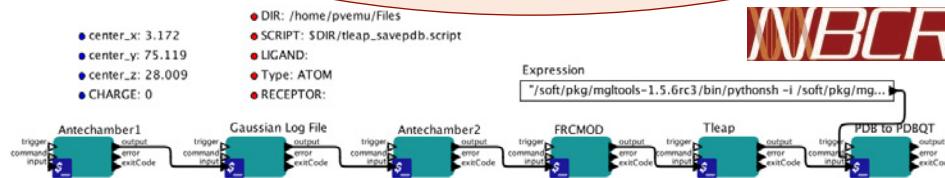


Data-Parallel Bioinformatics
biokepler.org

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



kepler-project.org



WorDS.sdsc.edu

SDSC SAN DIEGO SUPERCOMPUTING CENTER

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

UC San Diego