# Workflows as an Operational Tool Scientific Computing using Data Sci

İlkay ALTINTAŞ, Ph.D.
*Chief Data Science Officer, San Diego Supercomputer Center*
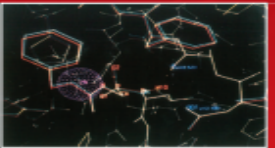*Founder and Director, Workflows for Data Science Center of E*

SDSC SAN DIEGO SUPERCOMPUTER CENTER

UC San Di

# SC is 31 Years Young!
## viding Cyberinfrastructure for Research and Educati

- tablished as a national percomputer resource center in 1985 NSF

- world leader in HPC, data-intensive mputing, and scientific data anagement

- rrent strategic focus on "Big Data" d "HPC Cloud" : versatile computing



SAN DIEGO SUPERCOMPUTER CENTER at UC SAN D

SDSC
30years
of Turning Data to Disco

I ♥ BIG DATA
N DIEGO SUPERCOMPUTER CENTER at UC SAN DIEGO

In pioneering efforts in drug design, Paul Bash, et. al., using SDSC supercomputers, determine free energies of solvation for proteins and nucleic acids, and relative free energies for binding, published in *Science*.
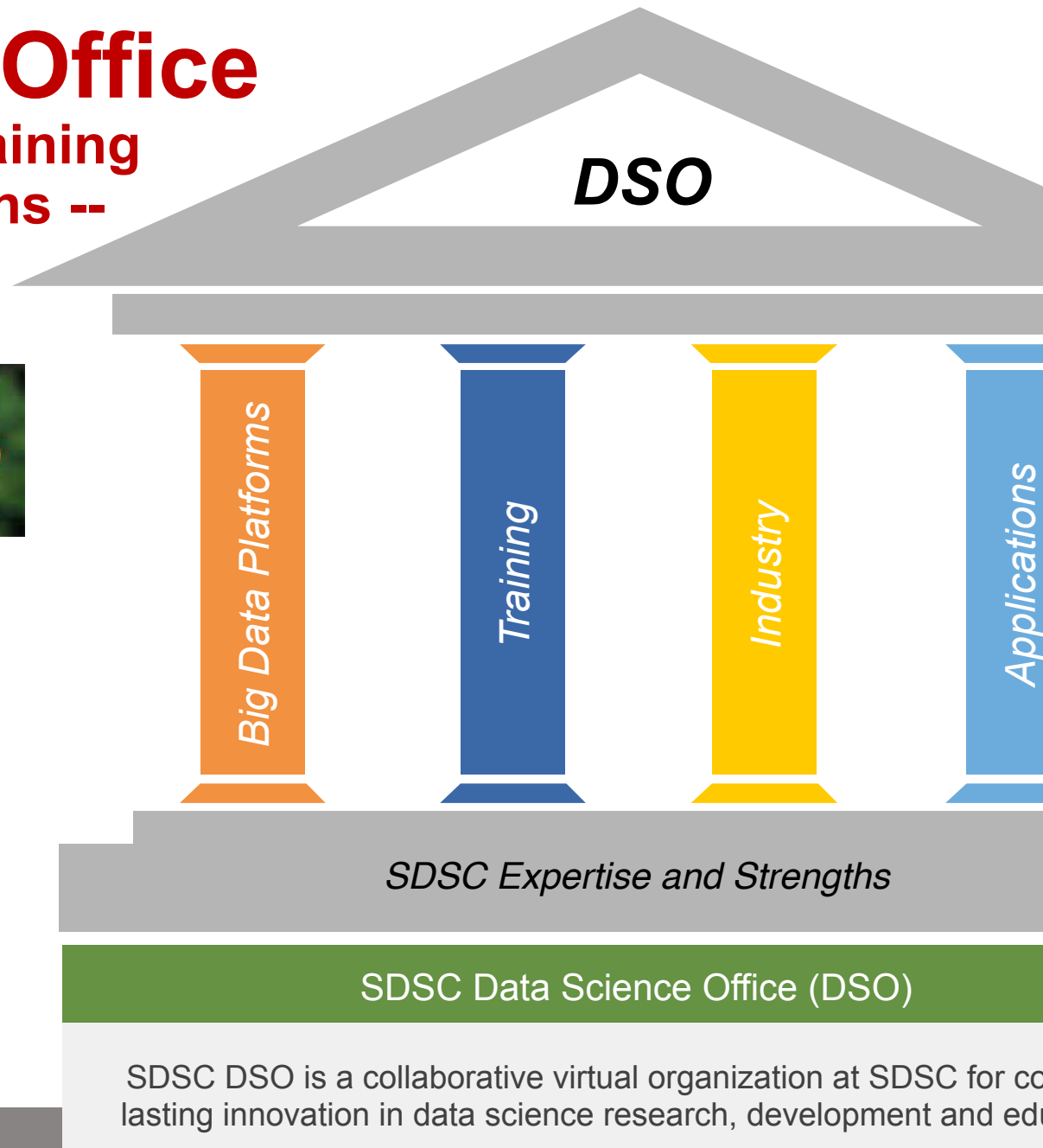
1987

1991
With large-scale computer simula SDSC, researchers led by J. Andrev at UCSD show how one of the fas in the world, acetylcholinesterase results are published in the *Proce National Academy of Sciences.*

**Two discoveries in design from 1987 an**
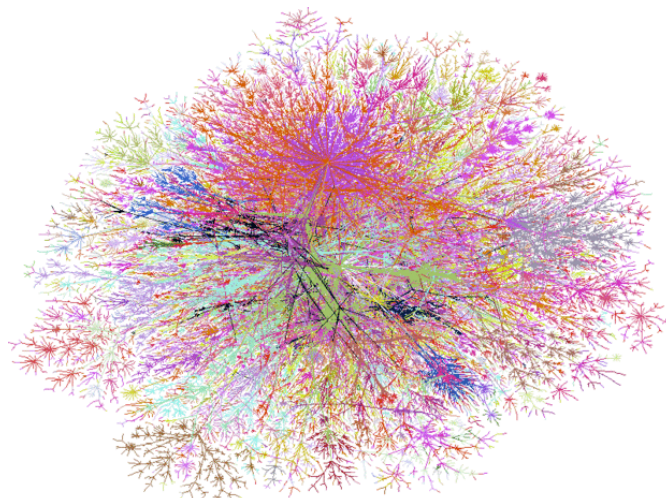
## ...e-Step Molecular Dynamics through Hydrogen Mass ...oning

4 fs

Time Step

*...alker Group*

*...nanosecond level snapshots of financial markets*

...oi[2], David O'Neal[3], Mao Ye[1] and Robert S. Sinkovits[2,*,†]

| ...original code | Wall time (s) modified code | Speedup |
|---|---|---|
| ...400 | 128 | 66x |
| ...200 | 437 | 126x |
| ...914 | 1145 | 113x |

**1 Comment** 💬 149

**CORPORATE RESPONSIBILITY**

## San Diego Supercomputer Center's Quake Research Wins $150,000 Global Impact Award

By Tonie Hansen on March 16, 2015

## Multi-GPU Implementation of a 3D Finite Difference Time Domain Earthquake Code on Heterogeneous Supercomputers

Jun Zhou[a,b,*], Yifeng Cui[a], Efecan Poyraz[a,b], Dong Ju Choi[a], Clark C. Guest[b]

In-Order Communication
(first west/east, then south/north)

MPI:
MPI:
MPI:
MPI:
CPU:
GPU:

- 🟥 Data copy from GPU to CPU
- 🟥 MPI Communication between CPUs
- 🟩 Data copy from CPU to GPU
- 🟦 Computing inside GPU

**Gordon**: First ...Flash-based Supercomput... for Data-intensive Apps

...riant calling leveraging next-gene... ...g for large-scale whole-genome se...

...[2], Tristan M. Carland[2], Glenn K. Lockwood[3], Wayne... ...Mahidhar Tatineni[?], C Chris Huang[4], Sarah Lamberth[4], Yauheniya Cherkas... ...Brodmerkel[4], Ed Jaeger[45], Lance Smith[45], Gunaretnam Rajagopal[45], Mark... ...Nicholas J. Schork[2,*]

| Step | Tool | Memory per command (GB) | Cores per command |
|---|---|---|---|
| Map | BWA | 32 | 8 |
| Bam | Samtools | 4 | 1 |
| Merge | Samtools | 4 | 1 |
| Sort | Samtools | 4 | 1 |
| MarkDuplicates | PicardTools | 7 | 2 |
| TargetCreator | GATK | 7 | 2 |
| IndelRealigner* | GATK | 12 | 3 |
| BaseRecalibrator | GATK | 30 | 8 |
| PrintReads* | GATK | 30 | 8 |
| HaplotypeCaller | GATK | 60 | 16 |

*Smaller memory allocation and more samples per node may prove ...computationally efficient

## ...net: *Serving the Long ...of Science*

...andard racks ...= 1944 nodes ...= 46,656 cores ...= 249 TB DRAM ...= 622 TB SSD ...flop/s

- 36 GPU nodes
- 4 Large Memory nodes
- 7 PB Lustre storage
- High performance virtualization

...SC SUPERCOMPUTER CENTER

UC San Di...

# SDSC Data Science Office

## -- Expertise, Systems and Training for Data Science Applications --



**DSO**

Big Data Platforms | Training | Industry | Applications

SDSC Expertise and Strengths

SDSC Data Science Office (DSO)

SDSC DSO is a collaborative virtual organization at SDSC for co[...] lasting innovation in data science research, development and edu[...]

SAN DIEGO SUPERCOMPUTER CENTER

UC San Di[...]

# mputing Today has Many Shapes and Siz

*COMPUTING AT SCALE*

**+**

*BIG DATA*

*Requires:*

- *Data manageme*
- *Data-driven met*
- *Scalable tools fo*
  *dynamic coordi*
  *and stateful reso*
  *optimization*
- *Skilled interdisci*
  *workforce*

*Enables dynamic data-driven applications*

*New era*
*data scie*

*Computer-Aided Drug Discovery*    *Smart Cities*    *Disaster Resilience and Response*

*anufacturing*    *Personalized Precision Medicine*    *Smart Grid and Energy Management*

# Needs and Trends
# for Scientific Computing under the Influence of Big Data and Cloud Systems

*New era of data science!*

- More data-driven
- More dynamic
- More process-driven
- More collaborative
- More accountable
- More reproducible
- More interactive
- More heterogeneous

BIG DATA

Size

Complexity

Volume

Speed

Variety

Velocity

Application-Specific
**Value**

Veracity

Valence

Quality

Connected

SAN DIEGO
SUPERCOMPUTER CENTER

# ta Management and Processing in the Big Data has Unique Challenges!

| Volume | → | Scalable batch processing |
| Velocity | → | Stream processing |
| Variety | → | Extensible data stora access and integrati |

# ...ese challenges come with new tools to tackle the...

Higher levels:
Expression and interactivity

Hive

Pig

Giraph

Storm

Spark

Flink

HBase

Cassandra

MongoDB

MapReduce

YARN

hadoop

HDFS

hadoop HDFS

Lower levels:
Storage and scheduling

**How do we use these new tools and combine them with existing solutions in scientific computing and data science?**

**COORDINATION AND WORKFLOW MANAGEMENT**

**DATA INTEGRATION AND PROCESSING**

**DATA MANAGEMENT AND STORAGE**

# ample Big Data Processing Pipelin



ta Processing Pipeline

Application level code

Logstore

Datastore

Backends + MapReduce

BigQuery

SQL    API

The big data pipeline

ML

PREDICT

Data source

CSV

Text

JSON

Data analytics pipeline

Data cleaning

Python

MR

Spark

Data preprocessing

MR_v4

Bash    MR

Spark

MR_v1

Data analysis

Python

MR

Spark

HDFS

clean

create graph    analyze graph

S3

clean    transform    mahout    result

HDFS

load    clean    transform

APACHE HBASE

# COORDINATION AND WORKFLOW MANAGEMENT

ACQUIRE → PREPARE → ANALYZE → REPORT → ACT

OOZiE

Apache Zookeeper

Kepler

http://kepler-project.org

# Research Challenges

w to easily program a workflow using the Big Data Patterns?

w to parallelize legacy tools for Big Data?

nich pattern(s) to use under which Big Data engine to use, e.g., as Hado
nk or Spark?

d-to-end performance prediction for Big Data applications/workflows (
g to run)

- Knowledge based: Analyze performance using profiling techniques and dependency analysis
- Data driven: Predict performance based on execution history (provenance) using machine learning
  techniques

demand resource provisioning and scheduling for Big Data application
here and how to run)

- Find the best resource allocation based on execution objectives and performance predictions
- Find the best workflow and task configuration on the allocated resources

# ng Big Data Patterns in Kepler Workflo

 define a separate DDP

stributed Data-Parallel) task/actor
each pattern

ese DDP actors partition input
a and process each partition
parately

r-defined functions are described
sub-workflows of DDP actors

P director: executes DDP
rkflows on top of Big Data
gines

- *Visual programming*
- *Parallel execution of th sub-workflows*
- *Existing actors can eas reused for new tools*



(a) Top-level Workflow

(b) Sub-workflow for tRNAscan-SE

(c) Sub-workflow f

*...rkflow is a combination of modules running in places and intera... with each other via data or message passing via a connection ...*

**Individual Executable**
-- Often defined as a part of module on a workflow --

**+**

**Computational Resource**
-- Often HPC, commodity or Cloud based systems --

**Module**

**Individual Executable**
-- Often defined as a part of module on a workflow --

**+**

**Computational Resource**
-- Often HPC, commodity or Cloud based systems --

**Module**

**Individual Executable**
-- Often defined as a part of module on a...

**+**

**Computational Resour...**
-- Often HPC, commodity or Cloud based...

**Module**

**Individual Executable**
-- Often defined as a part of module on a workflow --

**+**

**Computational Resource**
-- Often HPC, commodity or Cloud based systems --

**Module**

**Individual Executable**
-- Often defined as a part of module on a workflow --

**+**

**Computational Resource**
-- Often HPC, commodity or Cloud based systems --

**Module**

*...kflow Performance == Composed Module Performance on an Infrastructure Inst...*

# RAMMCAP

# nization of Heterogeneous Resource Utilization using bioKep

# more traditional HPC and HTC workloads to th

*Dynamic data-driven coordination & resource optimization*

**Requires:**

*Ability to explore and scale on multiple platforms*

**? Are workflows increasingly becoming the dynamic operations research tool for scien**

**hallenge:** Make workflows more aware of stributed system and application state!

# me steps to get there…

Analyze each task in a workflow as an individual odule based on all past executions of that executable sk.

Model workflow performance as an aggregate of edictions of individual tasks to form prediction for tire workflow.

Include system level analytics at the workflow level t ake sure scheduling can use system level information o account in a dynamic data-driven way.

**1. Profiling Framework**

**3. Module Performance Prediction Workflow Composition (f)**

Uses existing tools and computing systems!

Computing is just one part of big data workflows…

… new methods needed!

Feature Selection and Training

3-a: Module Prediction: Single Predictor For Two Independent Software Tools

# Workflows for Data Science Center of Excellence at SDSC

*WorDS.sdsc.edu*

**Technologies**

**Development**

*Focus on the question, not the technology!*

*Real-Time Hazards Management*
*wifire.ucsd.edu*

*Data-Parallel Bioinformatics*
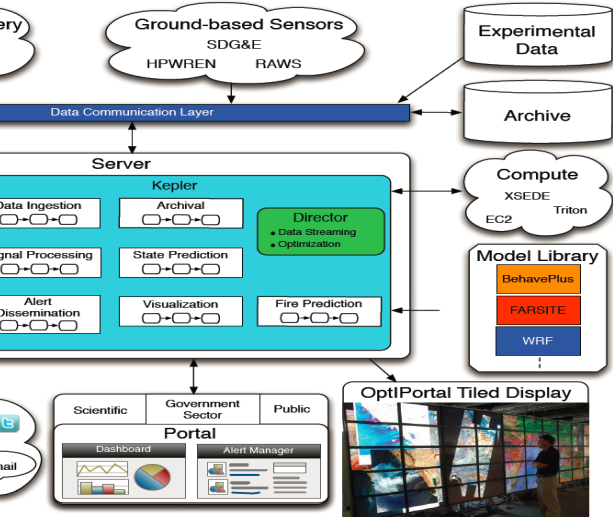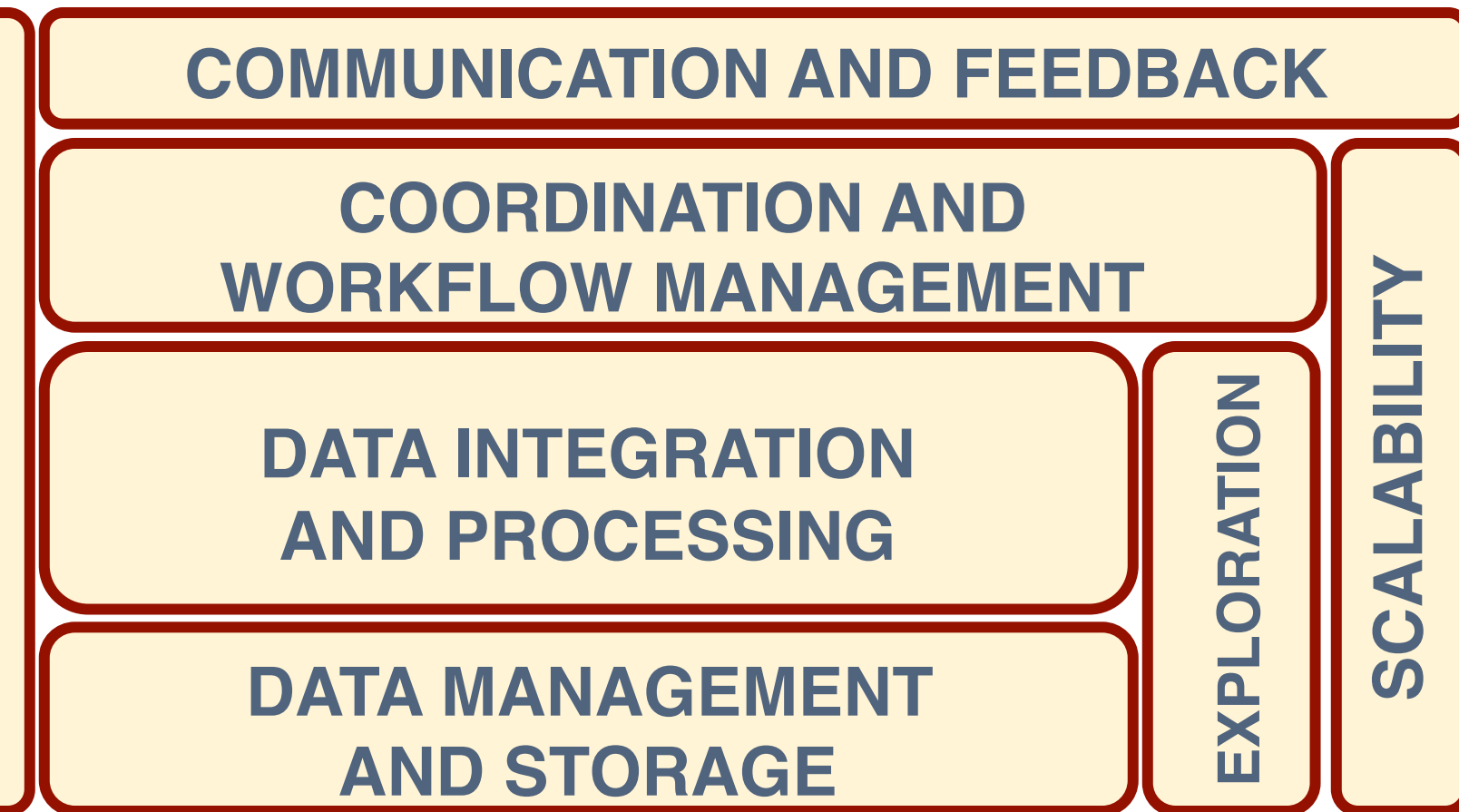*bioKepler.org*

## WorDS Center

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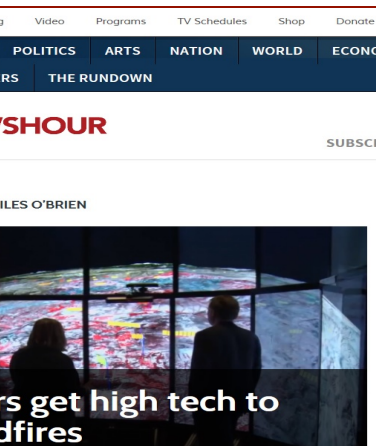- *Save time, energy and money*
- *Formalize and standardize*

**Goal:** Methodology and tool development to build automated operational workflow-driven solution architectures on big data and HPC platforms.

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

- center_x: 3.172
- center_y: 75.119
- center_z: 28.009
- CHARGE: 0

- DIR: /
- SCRIPT: $DIR/tleap_savepdb.script
- LIGAND:
- Type: ATOM
- RECEPTOR:

Expression
"/soft/pkg/mgltools-1.5.6rc3/bin/pythonsh –i /soft/pkg/mg...

Antechamber1  Gaussian Log File  Antechamber2  FRCMOD  Tleap  PDB to PDBQT

Note: Amber and Gaussian must be loaded on terminal before workflow executes.

SAN DIEGO SUPERCOMPUTER CENTER

# Examples: Use of Workflows as an Application Integr Tool for "Big" Data and Computational Science

# Towards an Integrated Cyberinfrastructure
## for Scalable Data-Driven Monitoring,
## Dynamic Prediction and Resilience of Wildfires

Ilkay Altintas[1], Jessica Block[2], Raymond de Callafon[3], Daniel Crawl[1], Charles Cowart[1], Amarnath Gupta[1], Mai Nguyen[1], Hans-Werner Braun[1], Jurgen Schulze[2], Michael Gollner[4], Arnaud Trouve[4] and Larry Smarr[2]

[1]San Diego Supercomputer Center, University of California San Diego, U.S.A.
[2]Qualcomm Institute, University of California San Diego, U.S.A.
[3]Dept. of Mechanical and Aerospace Engineering, University of California San Diego, U.S.A.
[4]Fire Protection Engineering Dept., University of Maryland, U.S.A.

*wifire.ucsd.edu*

# Hybrid Data Processing Architecture

**SPEED LAYER**

- Stream processing
- Real-time data interfaces

**BATCH LAYER**

- Batch processing on all data
- Batch data collection generation

**SERVING LAYER**
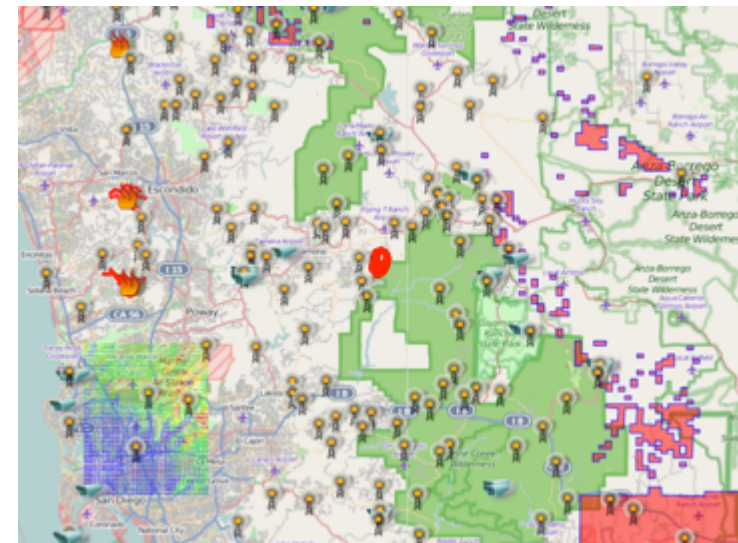
- Querying

a sources formally described

a merged from multiple sources into a single, unified model

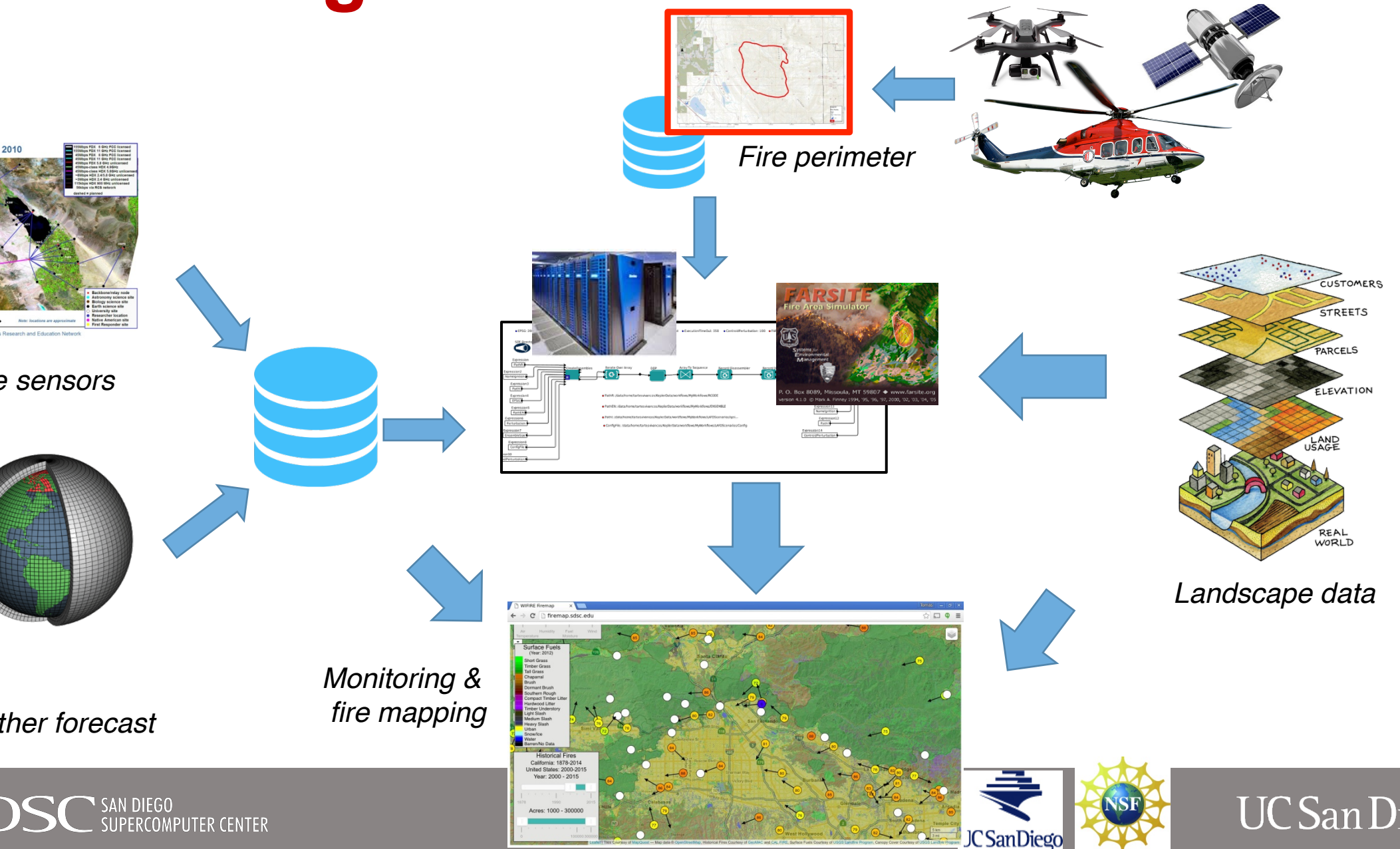Measurements from weather stations and cameras

Fire perimeters, e.g., InciWeb , GeoMac, SANDAG

Model output, e.g., FARSITE, Firefly, etc.

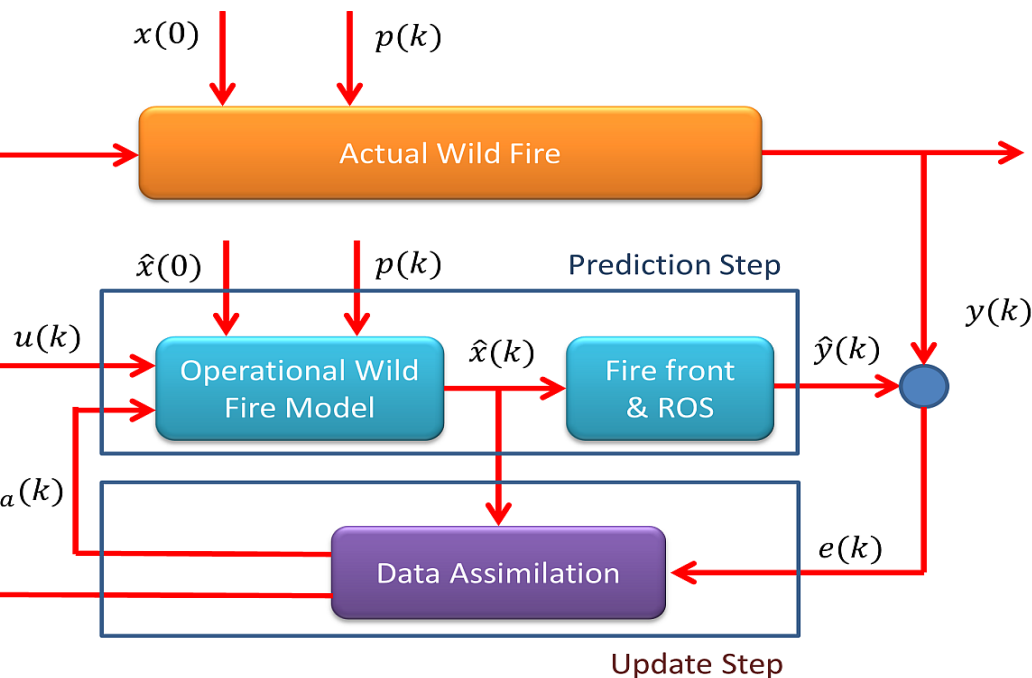nified REST interface to access data multiple formats

# Modeling Workflows in WIFIRE



Fire perimeter

Fire sensors

Weather forecast

Monitoring & fire mapping

Landscape data

CUSTOMERS

STREETS

PARCELS

ELEVATION

LAND USAGE

REAL WORLD

# Closing the Loop using Big Data

## -- Wildfire Behavior Modeling and Data Assimilation --



$x(0)$    $p(k)$

**Actual Wild Fire**

$\hat{x}(0)$    $p(k)$    **Prediction Step**

$u(k)$

**Operational Wild Fire Model**   $\hat{x}(k)$   **Fire front & ROS**   $\hat{y}(k)$

$y(k)$

$a(k)$

**Data Assimilation**   $e(k)$

**Update Step**

*Conceptual Data Assimilation Workflow with Prediction and Update Steps using Sensor Data*

- Computational costs for existin models too high for real-time analysis

- *a priori -> a posteriori*
  - Parameter estimation to make adjustments to the (input) param
  - State estimation to adjust the simulated fire front location with posteriori update/measurement actual fire front location

# Summary: Three questions about converge workflow applications! (Out of many…)

| | | |
|---|---|---|
| y exploratory an in the loop ponents: | Needs to run different parts of the workflow on changing distributed platforms: | Accountability an reporting needed each step: |
| **ow can we scale the roducts of xploratory steps in roduction mode?** | **Is workflow scheduling a closed control loop problem?** | **What does provenance and reproducibility r in dynamic applications?** |

**Questions?**

Work funded by NSF, DOE, NIH, UC San Diego and industry partners.

*WorDS Director:  Ilkay Altintas, Ph.D.*
*Email: altintas@sdsc.edu*