# Processing Streaming Data In High Energy Physics Workflows

Nathan Tallent, Darren Kerbyson, Mahantesh Halappanavar
Malachi Schram, Kevin Barker, Luis de la Torre, Ryan Friese, Jian Yin
Eric Stephan, Kerstin Kleese van Dam
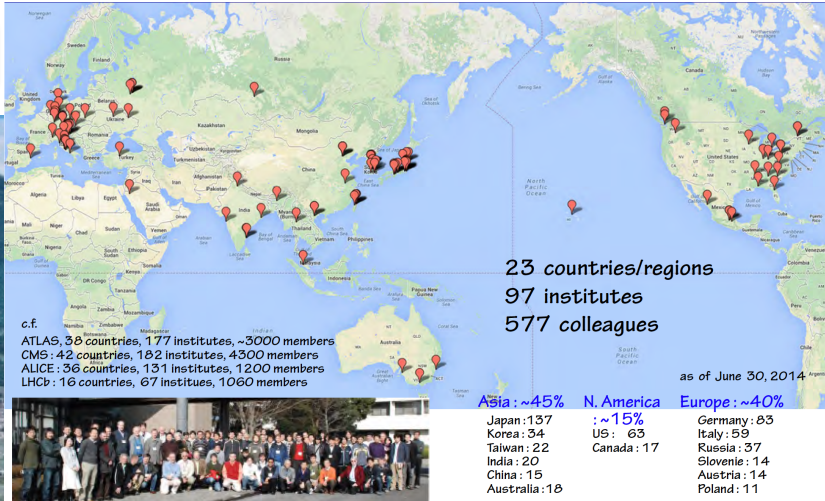
Pacific Northwest National Lab

*STREAM '16 Workshop*

March 22-23

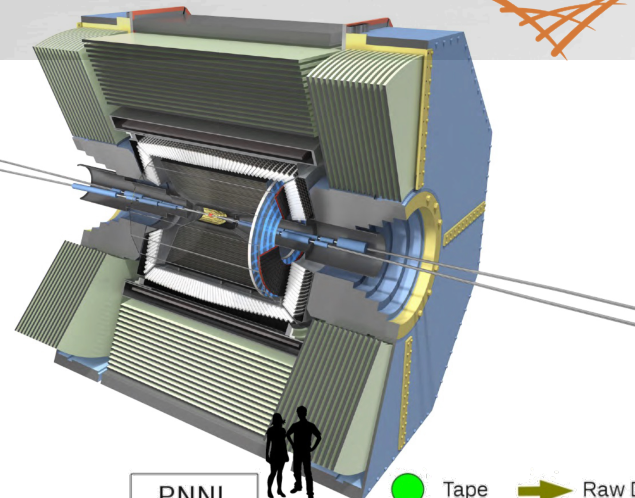# High Energy Physics: Belle II Analysis Workflow
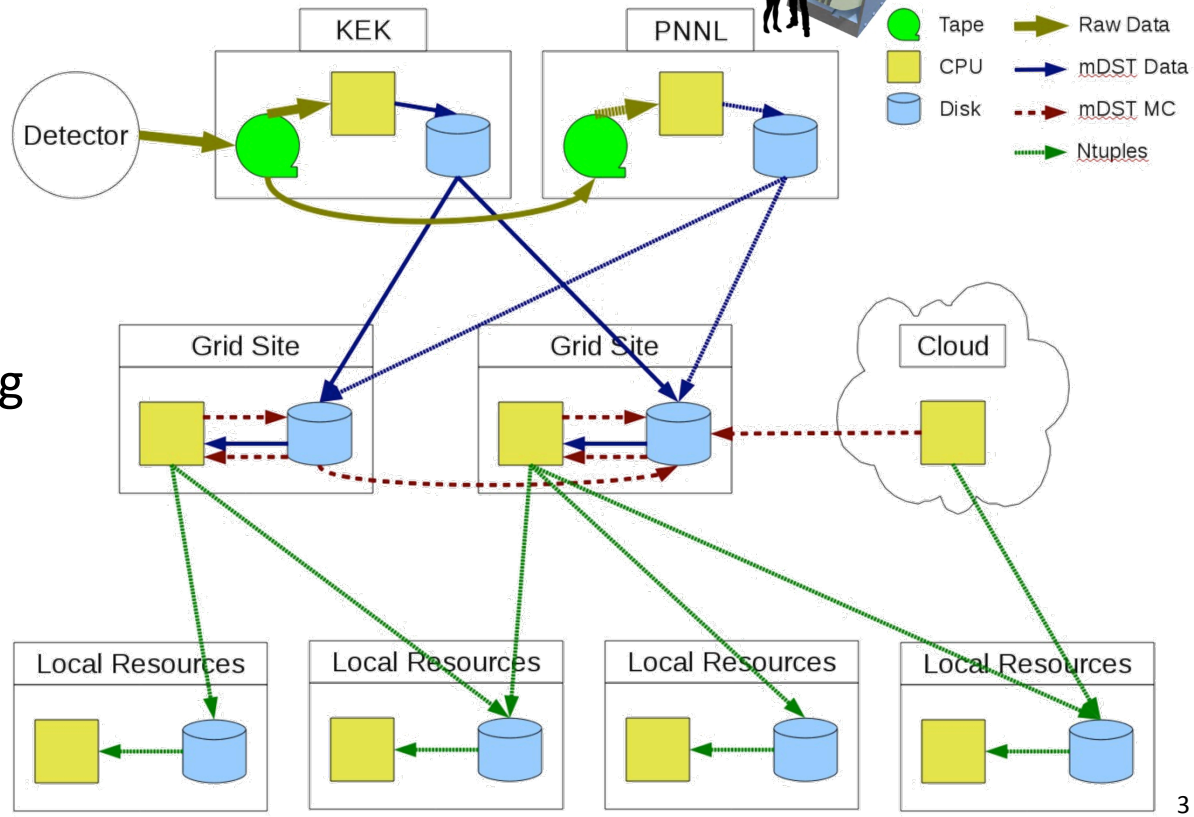
**International effort to advance particle physics**

23 countries/regions
97 institutes
577 colleagues

c.f.
ATLAS, 38 countries, 177 institutes, ~3000 members
CMS : 42 countries, 182 institutes, 4300 members
ALICE : 36 countries, 131 institutes, 1200 members
LHCb : 16 countries, 67 institues, 1060 members

as of June 30, 2014

| Asia : ~45% | N. America | Europe : ~40% |
|---|---|---|
| Japan : 137 | : ~15% | Germany : 83 |
| Korea : 34 | US : 63 | Italy : 59 |
| Taiwan : 22 | Canada : 17 | Russia : 37 |
| India : 20 | | Slovenia : 14 |
| China : 15 | | Austria : 14 |
| Australia : 18 | | Poland : 11 |

Credit: Malachi Schram

**KEK**
High Energy Accelerator Reseach Organization

# Belle II: Geographically Distributed Analytics

▶ Belle II Workflow: Extensive data analysis

▶ Data! 25 PB/year of raw data

- Stored data expected to reach 350 PB

▶ Many analysis pipelines run concurrently

- Normalize raw data

- Physics analysis
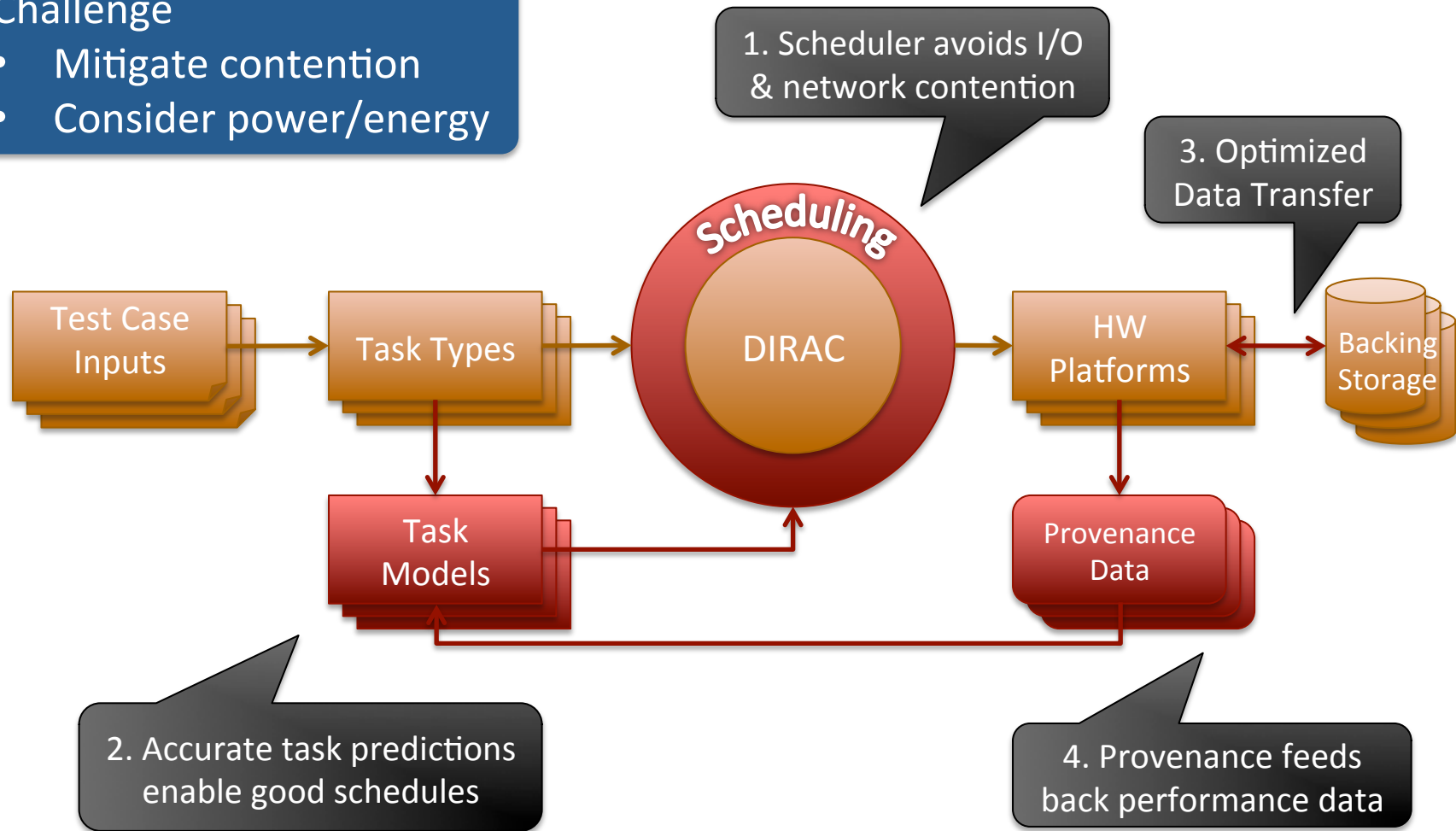
- Monte Carlo simulations

- Data storage/archiving

Contention! Many independent data accesses in small window.

# IPPD's 'Enhanced' Belle II Workflow Execution

**Challenge**
- Mitigate contention
- Consider power/energy

1. Scheduler avoids I/O & network contention

3. Optimized Data Transfer

**Scheduling**

DIRAC

Test Case Inputs → Task Types → DIRAC → HW Platforms → Backing Storage

Task Models

Provenance Data

2. Accurate task predictions enable good schedules

4. Provenance feeds back performance data

IPPD: Integrated End-to-End Performance Prediction and Diagnosis

# Hierarchical Scheduler Avoids I/O Contention

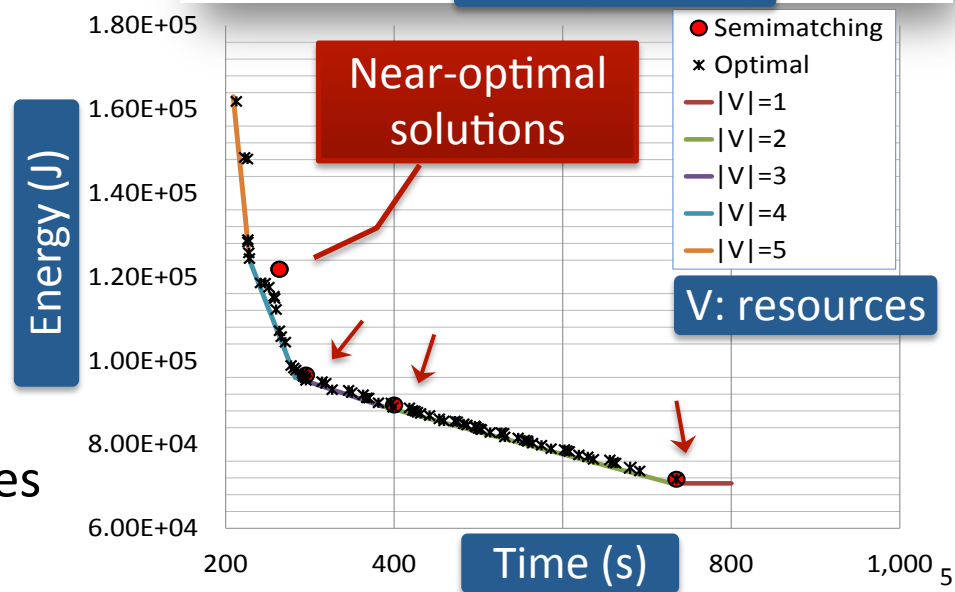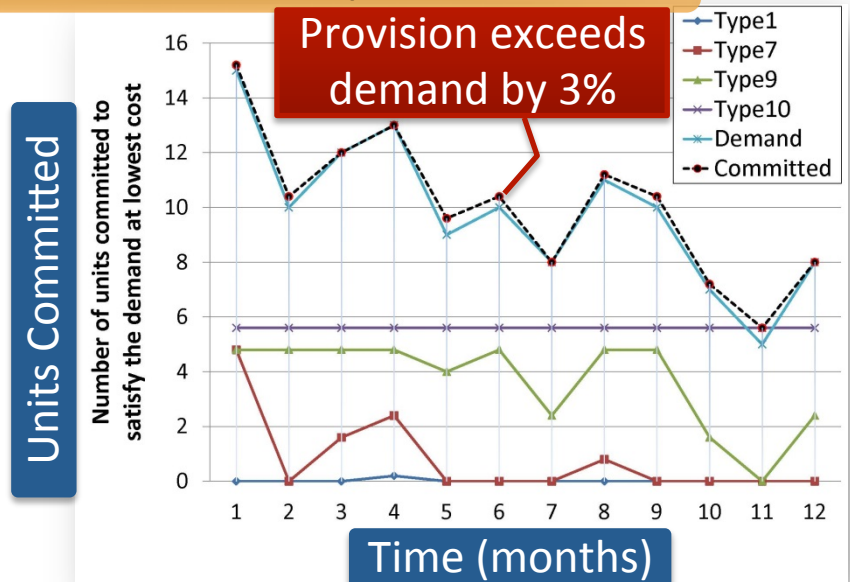## Approximate, two-level solution for NP-hard problem

**Challenge**
- Demand & supply vary considerably
- Hard to estimate task times
- Congestion dilates execution time

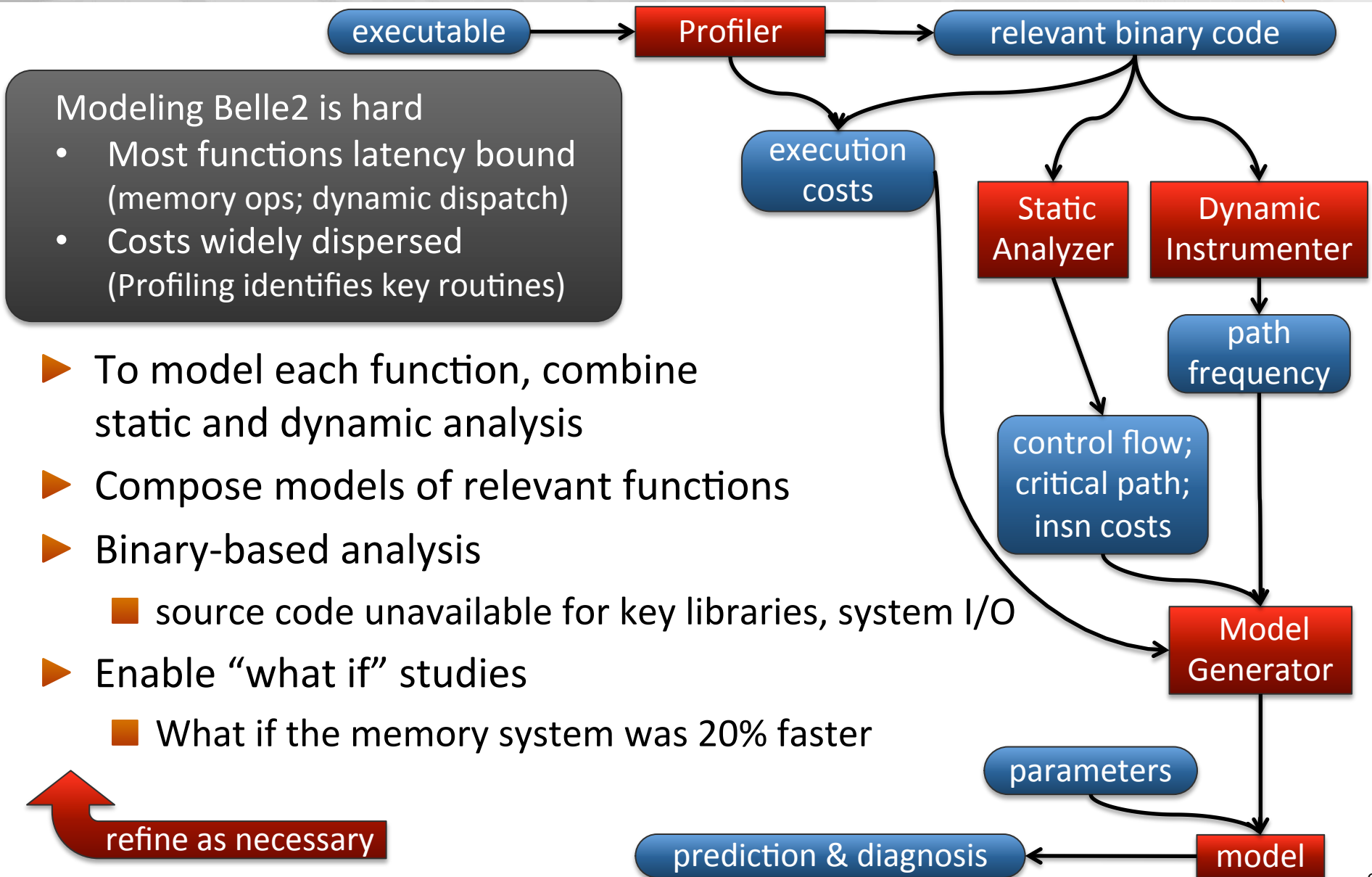**1** Most *cost-efficient subset* of compute resources that meets the tasks' demand

- unit commitment (power grid)
- mixed int/linear programming

**2** Best assignment of tasks to compute resources

- bi-objective: energy & time
- semi-matching: tasks ⇸ resources



Provision exceeds demand by 3%

Units Committed

Number of units committed to satisfy the demand at lowest cost

Legend: Type1, Type7, Type9, Type10, Demand, Committed

Time (months)



Near-optimal solutions

Energy (J)

Legend: Semimatching, Optimal, |V|=1, |V|=2, |V|=3, |V|=4, |V|=5

V: resources

Time (s)

Time in seconds

# Analytical Modeling Predicts Task Execution Time

executable → Profiler → relevant binary code

**Modeling Belle2 is hard**
- Most functions latency bound (memory ops; dynamic dispatch)
- Costs widely dispersed (Profiling identifies key routines)

execution costs

Static Analyzer

Dynamic Instrumenter

path frequency

control flow; critical path; insn costs

▶ To model each function, combine static and dynamic analysis

▶ Compose models of relevant functions

▶ Binary-based analysis
  ■ source code unavailable for key libraries, system I/O

Model Generator

▶ Enable "what if" studies
  ■ What if the memory system was 20% faster

parameters

refine as necessary

prediction & diagnosis ← model

# Optimize Data Transfer via 'Paced' Prefetching
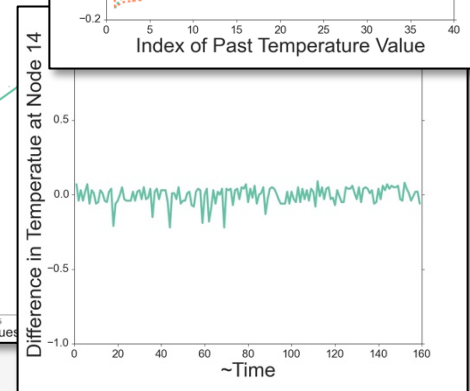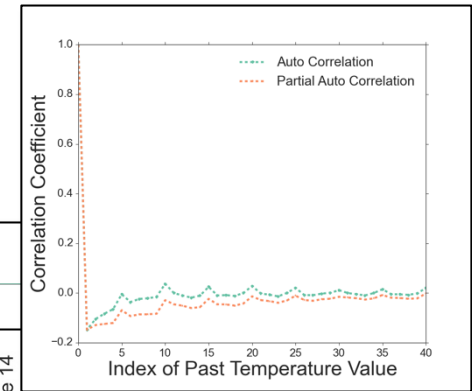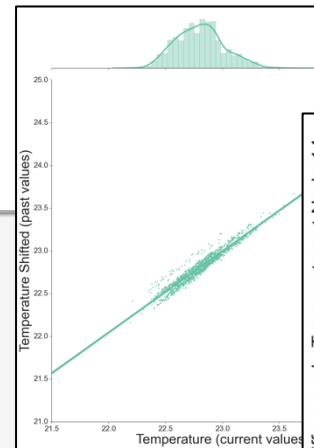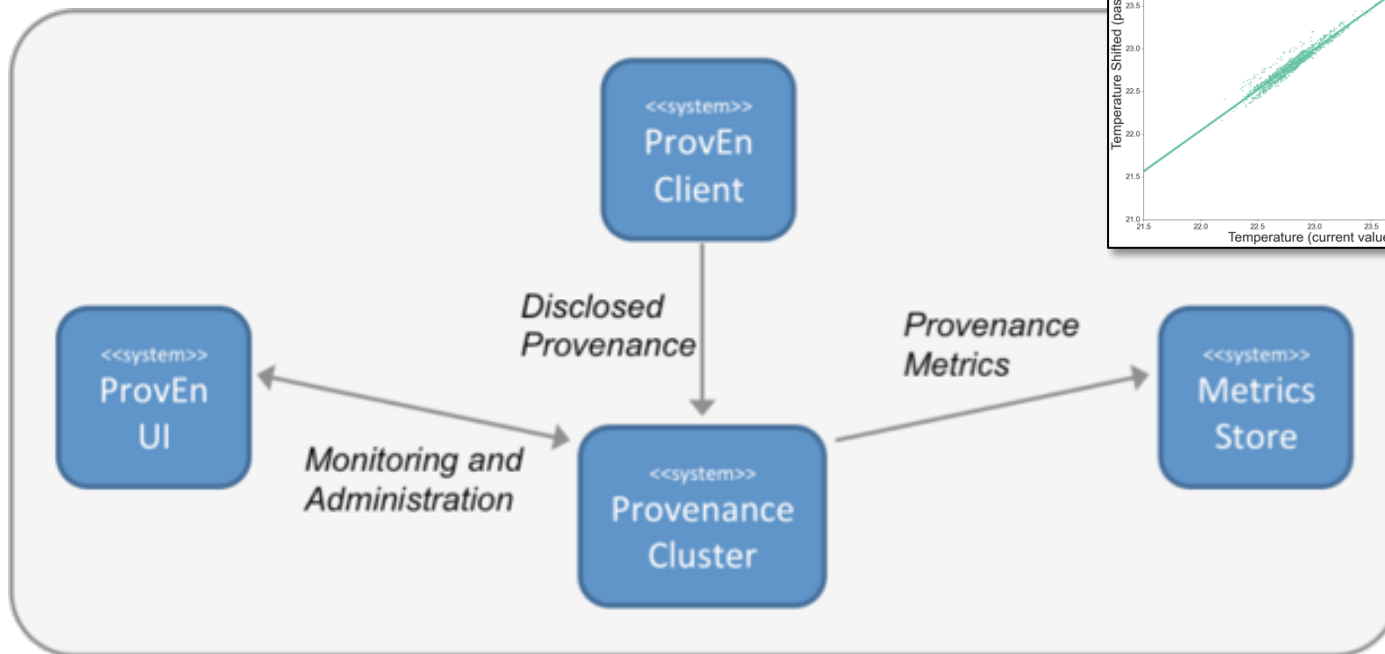
## Challenge: I/O Requests Create Blocking Time

▶ Prefetch data to reduce I/O blocking time

- Overlap remote data transfer and computation

- Retrieve only the needed part of a file

- Split data transfer across multiple internet connections

- Dynamically adjust given load on each connection

- Pace I/O request to improve end-to-end performance

# Provenance Feeds Back to Scheduler/Modeling

► Provenance delivers execution statistics to scheduler & modeling

► ProvEn (Provenance Environment) collects:

■ Time series-based information for system/host
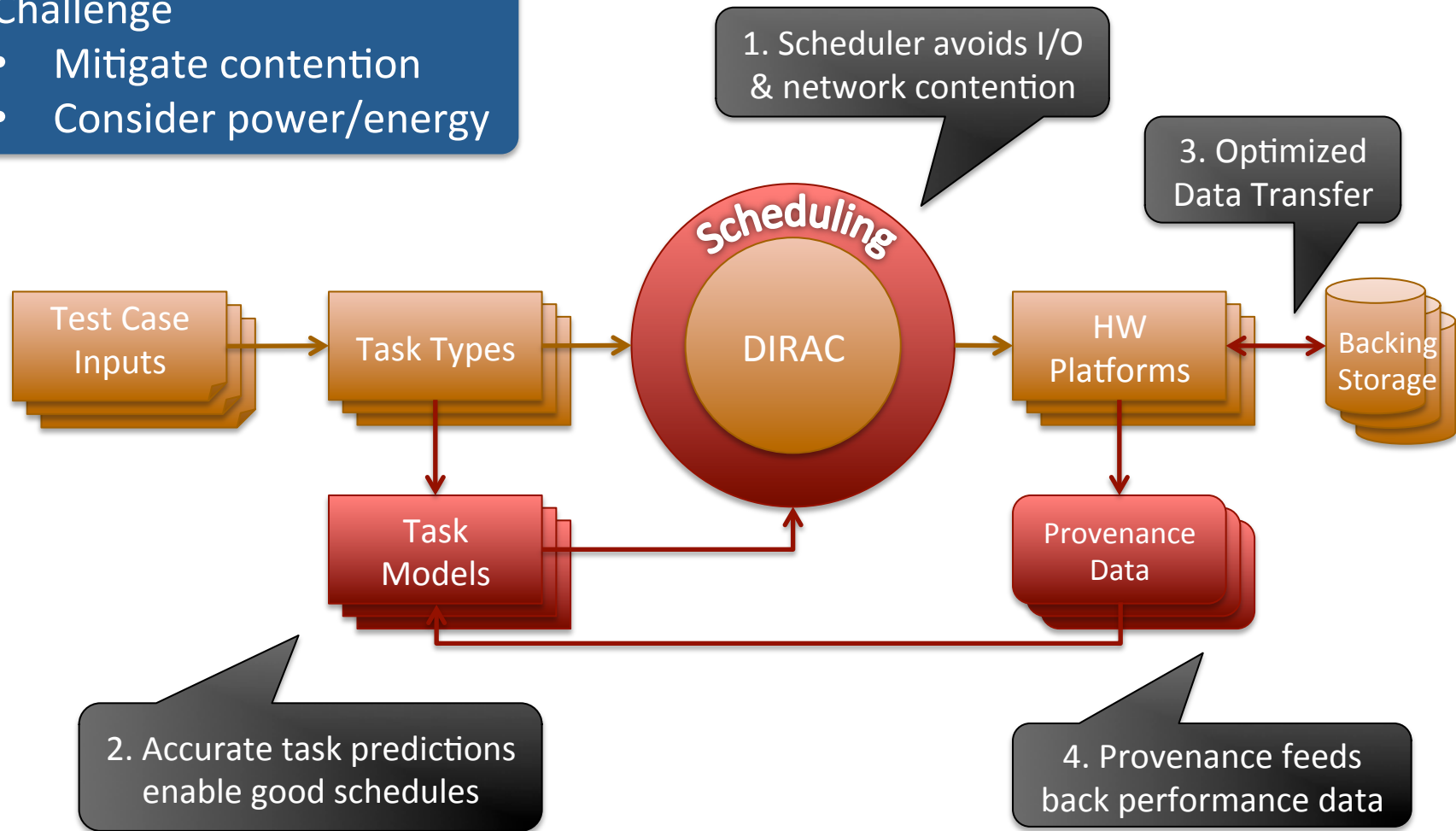
■ Performance metrics for application/workflow



Predictive
Analytics