# Collaborative Modeling, Simulation, and Analytics with Splash

Nicole Barberis, **Peter J. Haas**, Cheryl Kieliszewski, Yinan Li, Paul Maglio, Piyaphol Phoungphol, Pat Selinger, Yannis Sismanis, Wang-Chiew Tan, Ignacio Terrizzano, Haidong Xue, SJSU CAMCOS

**IBM Research – Almaden**

**Splash**
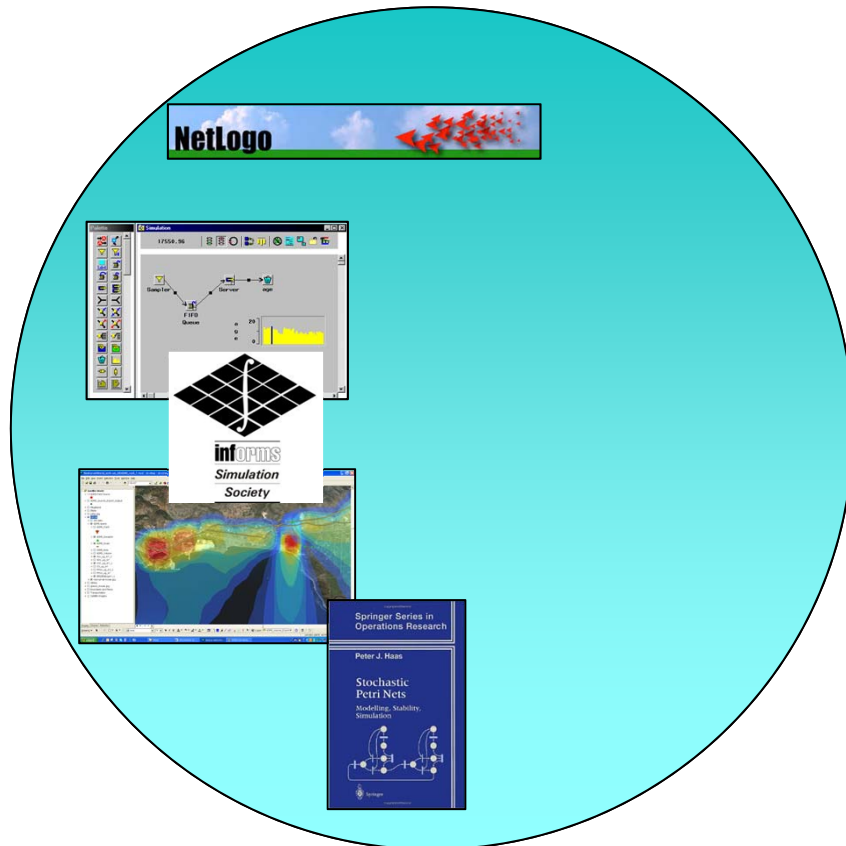Smarter Planet Platform
for Analysis and
Simulation of Health

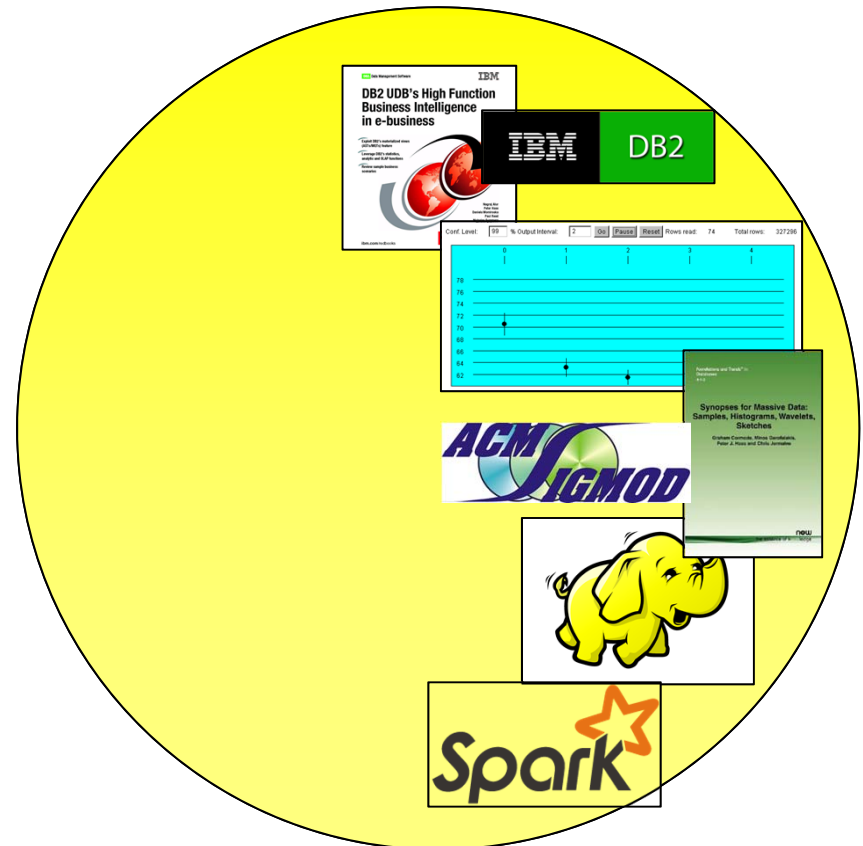http://researcher.watson.ibm.com/researcher/view_project.php?id=3931

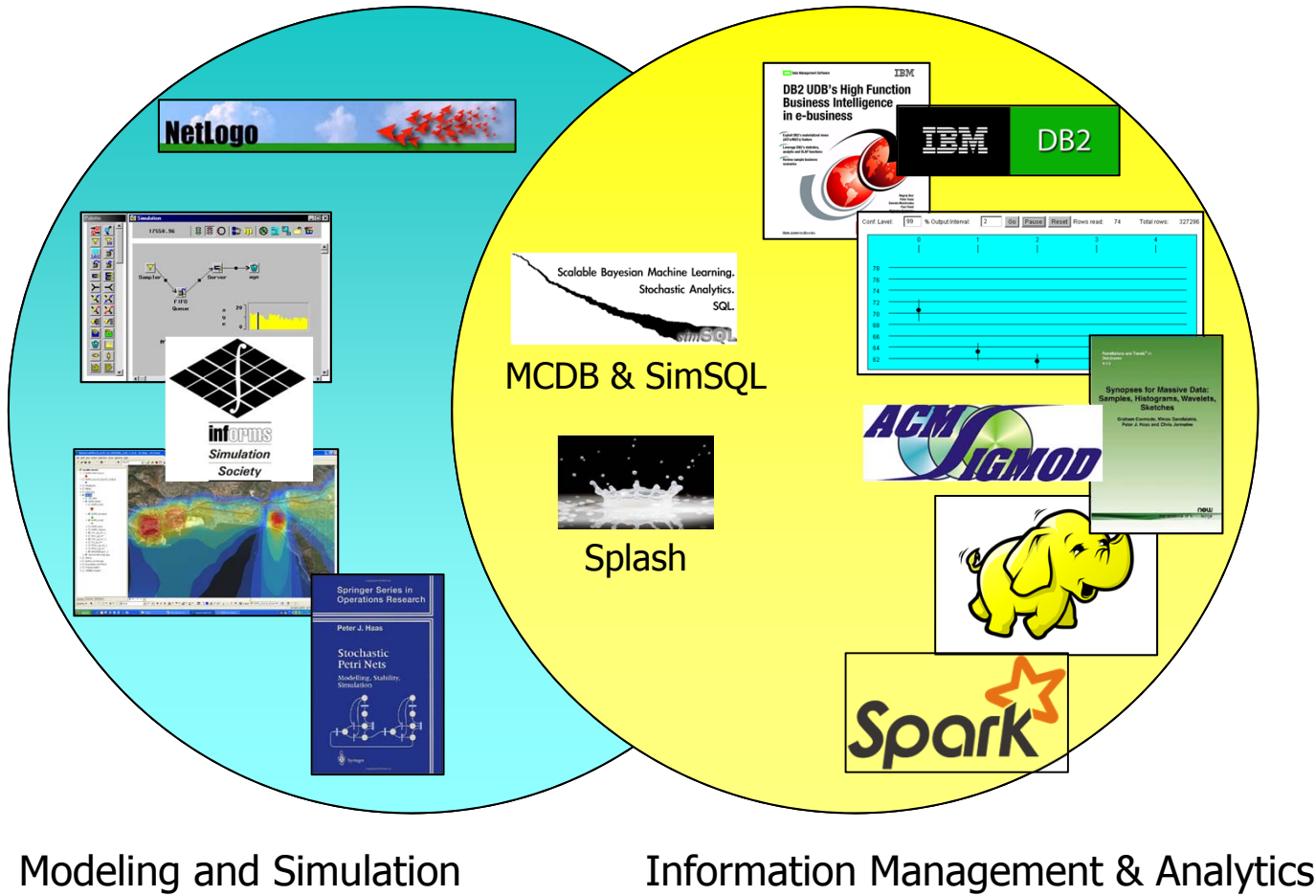# Some Context: Model-Data Ecosystems

# My Two Communities



Modeling and Simulation

Information Management & Analytics

# Opportunities for Innovation at the Intersection



Modeling and Simulation
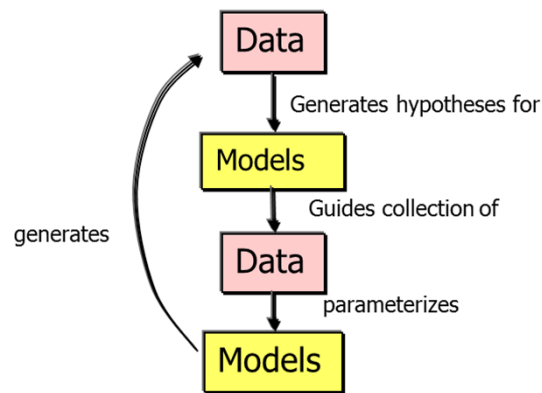
Information Management & Analytics

# Some Further Thoughts and Examples [PODS 2014 Tutorial]

(In addition to large-scale scientific environments)

- **Data-intensive simulation**
  - Simulations within databases
  - Databases within simulations
  - Data harmonization at scale

- **Information integration**
  - Simulation as an information-integration tool
  - Combining real and simulated data

- **And more!**



Ecosystem of Data and Models

# Motivation for Splash

# The Setting: Analytics for Decision Support

**informs online**

**Analytics Section**

"Analytics is…a complete [enterprise] problem solving and decision making process"

**Descriptive Analytics**: Finding patterns and relationships in historical and existing data

**Predictive analytics**: predict future probabilities and trends to allow what-if analysis

**Prescriptive analytics**: deterministic and stochastic optimization to support better decision making

Splash

# Shallow Versus Deep Predictive Analytics



Extrapolation

**United States House Prices**

Actual prices

Extrapolation of 1970-2006
median U.S. housing prices



NCAR Community
Atmosphere Model (CAM)

**3.3  Eulerian Dynamical Core**

$$\frac{\partial \zeta}{\partial t} = \boldsymbol{k} \cdot \nabla \times (\boldsymbol{n}/\cos\phi) + F_{\zeta_H},$$

$$\frac{\partial \delta}{\partial t} = \nabla \cdot (\boldsymbol{n}/\cos\phi) - \nabla^2 (E + \Phi) + F_{\delta_H},$$

$$\frac{\partial T}{\partial t} = \frac{-1}{a\cos^2\phi} \left[ \frac{\partial}{\partial\lambda}(UT) + \cos\phi\frac{\partial}{\partial\phi}(VT) \right] + T\delta - \dot{\eta}\frac{\partial T}{\partial\eta} + \frac{R}{c_p^*}T_v\frac{\omega}{p}$$
$$+ Q + F_{T_H} + F_{F_H},$$

$$\frac{\partial q}{\partial t} = \frac{-1}{a\cos^2\phi} \left[ \frac{\partial}{\partial\lambda}(Uq) + \cos\phi\frac{\partial}{\partial\phi}(Vq) \right] + q\delta - \dot{\eta}\frac{\partial q}{\partial\eta} + S,$$

$$\frac{\partial \pi}{\partial t} = \int_1^{\eta_s} \boldsymbol{\nabla} \cdot \left( \frac{\partial p}{\partial\eta}\boldsymbol{V} \right) d\eta.$$

# Big, Difficult, Important Problems Span Many Disciplines

## Need collaborative cross-disciplinary modeling and simulation



IBM analysis based on OECD data.

### GLOBAL FOOD SUPPLY

# Linking Policy on Climate and Food

H. C. J. Godfray, [1] J. Pretty, [2] S. M. Thomas, [3]* E. J. Warham, [3] J. R. Beddington[3]

At the United Nations (UN) climate negotiations in Cancún, Mexico, in December 2010, the parties agreed to a global target of no more than 2°C warming above preindustrial levels. In an important new step, both developed and developing countries agreed to take urgent action to reduce greenhouse gas (GHG) emissions to meet this long-term goal. They also set important milestones on reducing deforestation and providing funds to help developing countries adapt to climate change.



sions as delegates prepare for the next UN negotiations in December 2011 in South Africa. We need to rethink the way we use land to produce food, and to bring the challenges of sustainability and reducing emissions to the fore. This has been a central theme of the UK Government's Foresight Programme on the Future of Food and Farming to which we, along with experts from 35 countries, have been contributors. The study took a broad approach to the food system, including its impact on the environment and especially climate change, as well as the special needs of the world's poorest. It demonstrates both the importance of incorporating agriculture into climate change discussions, and the urgency for action (3).

**Agriculture and Climate Change**

Agriculture is a major source of $CO_2$ emissions and contributes a disproportionate amount of other GHGs with high impact on warming [about 47% and 58% of total $CH_4$ and $N_2O$

emissions by 20% by 2020 (8), whereas the UK has set the legally binding target of reducing emissions by 34% by 2020 and at least 80% by 2050 (9). Ambitious goals such as these cannot be achieved without involving the food system. Policies for mitigating climate change will have a substantial effect on production. If applied inappropriately, these could have a detrimental effect on food availability, especially for the 925 million (3) who already experience chronic hunger and for the additional billion or so who suffer nutrient and vitamin deficiencies.

**Land Use**

The Cancún meeting made notable progress in an area with important ramifications for the food system. Pressure from expanding agriculture has led to much recent tropical deforestation, especially in South America and Southeast Asia. Land conversion releases large amounts of GHGs and is one of the most serious, although indirect, ways that pressure from the food system contributes to global warming. The UN initiative on Reducing Emissions from Deforestation and Forest Degradation (REDD) offers financial

**POLICY**FORUM

## Linking Policy on

H. C. J. Godfray,[1] J. Pretty,[2] S. M. Thomas,[3]* E. J.

At the United Nations (UN) climate negotiations in Cancún, Mexico, in December 2010, the parties agreed to a global target of no more than 2°C warming above preindustrial levels. In an important new step, both developed and developing countries agreed to take urgent action to reduce greenhouse gas (GHG) emissions to meet this long-term goal. They also set important milestones on reducing deforestation and providing funds to help developing countries adapt to climate change.

The food system is complex, and interventions often have unintended and deleterious effects on food security, or have major consequences that affect GHS emissions. Agricultural, economic, and climate modelers must compare their models more systematically, share results, and integrate their work to meet the needs of policy-makers.

on warming [about 47% and 58% of total CH, and N.O.

Reducing Emissions from Deforestation and Forest Degradation (REDD) offers financial

**World Health Organization**

Health is a state of complete physical, mental, and social well-being and not merely the absence of disease or infirmity.

# Example: Unintended Outcomes in Healthcare Optimization



Simulation model of
Calgary Lab Services

Re-design

Avg. Patient Delay

Time (months)

0      6      12      18

T. R. Rohleder & D. P. Bischak & L. B. Baskin (2007). Modeling patient service centers with simulation and system dynamics. Health Care Manage. Sci., 10:1–12.

# Example: Unintended Outcomes in Healthcare Optimization



Simulation model of
Calgary Lab Services

Re-design

Avg. Patient Delay

Time (months)

System-dynamics social model of lab use

T. R. Rohleder & D. P. Bischak & L. B. Baskin (2007). Modeling patient service centers with simulation and system dynamics. Health Care Manage. Sci., 10:1–12.

# Example: Unintended Outcomes in Healthcare Optimization

Avg. Patient Delay



**Moral:**

**Combine models across disciplines for more robust decision making**

System-dynamics social model of lab use

T. R. Rohleder & D. P. Bischak & L. B. Baskin (2007). Modeling patient service centers with simulation and system dynamics. Health Care Manage. Sci., 10:1–12.

# Combining Models Across Disciplines is HARD

- Domain experts have different worldviews

- Use different vocabularies

- Sit in different organizations

- Develop models on different platforms

- Don't want to rewrite existing models!



Huang, T. T, Drewnowski, A., Kumanyika, S. K., & Glass, T. A., 2009,
"A Systems-Oriented Multilevel Framework for Addressing Obesity in the 21st Century,"
Preventing Chronic Disease, 6(3)

# Prior approaches to Combining Models

**Monolithic models**

- Create a monolithic model that encompasses all relevant domains

**Integrated models**

- Create modules that can be compiled into one
    - SpatioTemporal Epidemiological Modeler (STEM)
    - Community Atmospheric Model (CAM)

**Tightly-coupled models**

- Create modules that understand standard interfaces
    - DOD High Level Architecture (HLA)
    - Discrete-Event System Specification (DEVS)
    - Open Modeling Interface (OpenMI).

# Splash: An Alternative Approach

**Loosely** couple models and data via **data exchange**



Simulation model

Statistical model

Decision/optimization model

Dataset

Data transformation

**Splash = data integration + workflow management + simulation**

Re-use heterogeneous models and heterogeneous data that are curated by different domain experts

# Some Benefits of Loose Coupling

Facilitates cross-disciplinary modeling, analytics, and simulation for robust decision making under uncertainty

Enables re-use of models and datasets

Encourages comprehensive documentation and curation of models via metadata

Allows model flexibility:
– Upgrading to state-of-the-art
– Customizing for different users

| Weather Model V3.0 | → | Traffic Model | → | NYC Emergency-Services Model |
| Weather Model V4.0 | → | Traffic Model | → | NYC Emergency-Services Model |
| Weather Model V4.0 | → | Traffic Model | → | San Jose Emergency-Services Model |

# Splash

A prototype platform and service for integrating existing data, models, and simulations to gain insight needed for complex decision making related to policy, planning, and investment.



Splash Platform

Models
Data

SADL

Model and Data Curation

Model and Data Discovery

Model Composition

Composite-Model Execution

Experiment Management

Analysis
Visualization

DBMS, Hadoop, Visualization Tools, Information-Integration Tools, Stats Packages

# Model and Data Curation



Splash Platform

Models
Data

SADL

Model and Data Curation

Model and Data Discovery

Model Composition

Composite-Model Execution

Experiment Management

Analysis
Visualization

DBMS, Hadoop, Visualization Tools, Information-Integration Tools, Stats Packages

# Splash Actor Description Language (SADL)

- SADL provides "schemas and constraints" for models, transformations, and data, enabling interoperability

- SADL file for data (can exploit XSD)
  - **Attribute names, semantics, units**
  - **Constraints**
  - **How to access**
  - **Security**
  - **Experiment-management info**

- SADL file for a model:
  - **Inputs and outputs** (pointers to SADL files for data sources and sinks)
  - **How to execute** (info needed to synthesize command line)
  - **Semantics and assumptions**
  - **Provenance** (papers, ratings, ownership, security, change history, ...)
  - **RNG info**

```
<Actor name="BMI Model" type = "model" model_type = "simulation"
       sim_type = "continuous-deterministic" owner="Jane Modeler">
  <Description>
  Predict weight change over time based on an individual's energy    and food
  intake. Implemented in C. Reference: http://csel/asu.edu/?q=Weight
  </Description>
  <Environment>
    <Variable name="EXEC_DIR" default="/Splash" description="executable
  directory path"/>
    <Variable name="SADL_DIR" default="/Splash/SADL" description="schema
  directory path"/>
  </Environment>
  <Execution>
    <Command>$EXEC_DIR/Models/BMIcalc.out</Command>
    <Title>Run BMI model</Title>
  </Execution>
'  <Arguments>
    <Input name="demographics" sadl="$SADL_DIR/BMIInput.sadl"
          description="demographics data"/>
    <Output name="people" sadl="$SADL_DIR/BMIOutput.sadl"
          description="people's daily calculated BMI"/>
  </Arguments>
</Actor>
```
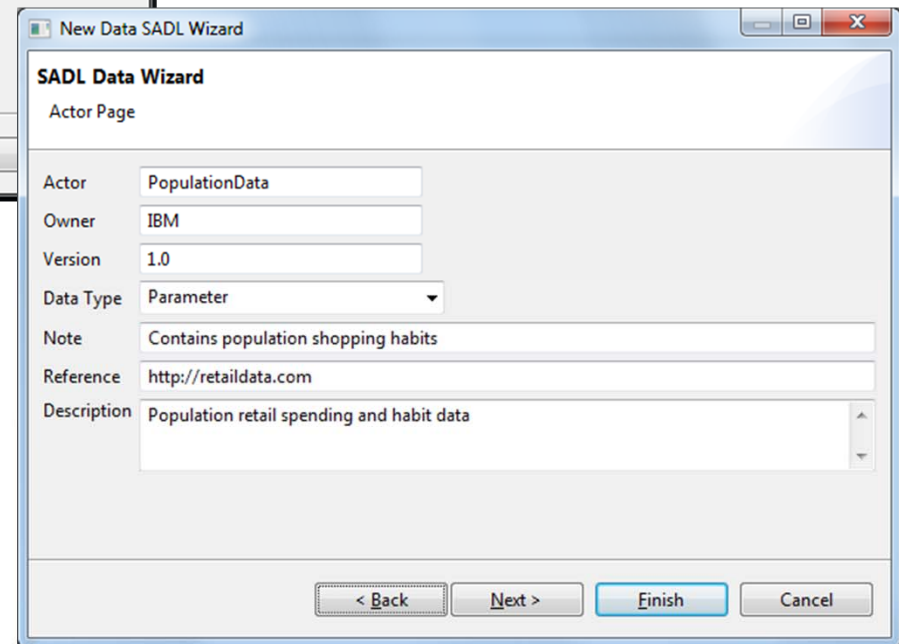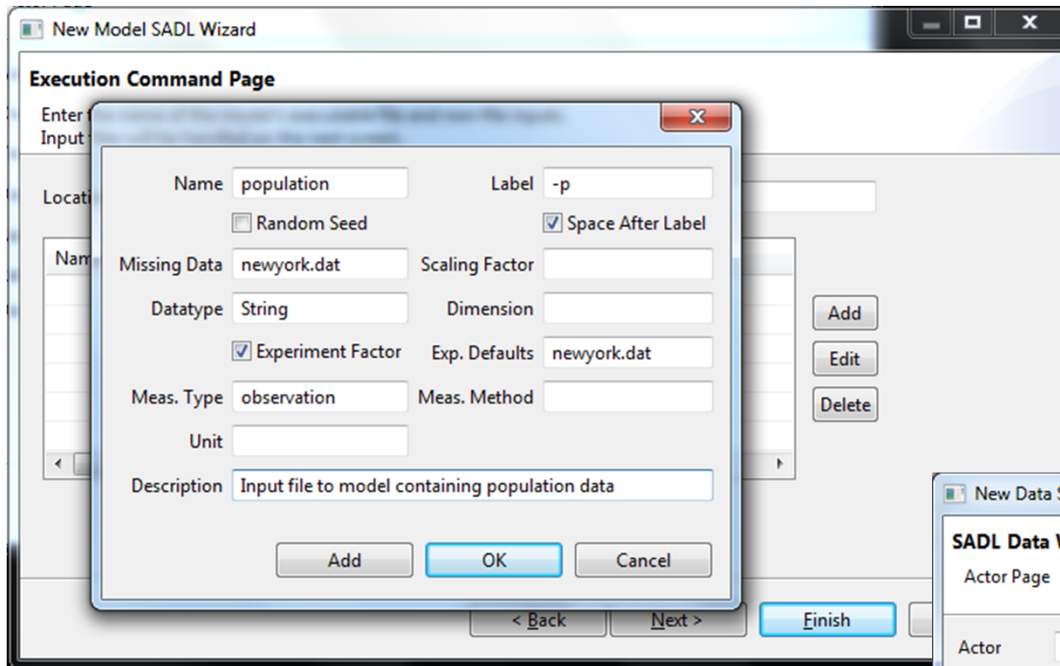
# Registration: Use Wizards to Create Model and Data SADL Files



Model Wizard offers step by step guidance to generate the Model's SADL, and the command line for invocation

Data Wizard generates SADL for model input and output files

# Model Composition

# Obesity Example

| Data source | Dataflow | Simulation model | Dataflow | Data Transformation |
|---|---|---|---|---|

# Sample Results

## If we open a new "healthy" food store in a "bad" neighborhood…



Without traffic model                    Including traffic model

# Implemented Obesity Example



- **Data actors:** input and output files, databases, web services, etc.
- **Model actors:** simulation, optimization, statistical models
- **Mapping actors:** data transformations, time and space alignment
- **Visualization actors:** graphs, reports, etc.

# Implemented Obesity Example



- **Data actors:** input and output files, databases, web services, etc.

- **Model actors:** simulation, optimization, statistical models

- **Mapping actors:** data transformations, time and space alignment

- **Visualization actors:** graphs, reports, etc.

# Implemented Obesity Example



- *Data actors:* input and output files, databases, web services, etc.
- *Model actors:* simulation, optimization, statistical models
- *Mapping actors:* data transformations, time and space alignment
- *Visualization actors:* graphs, reports, etc.
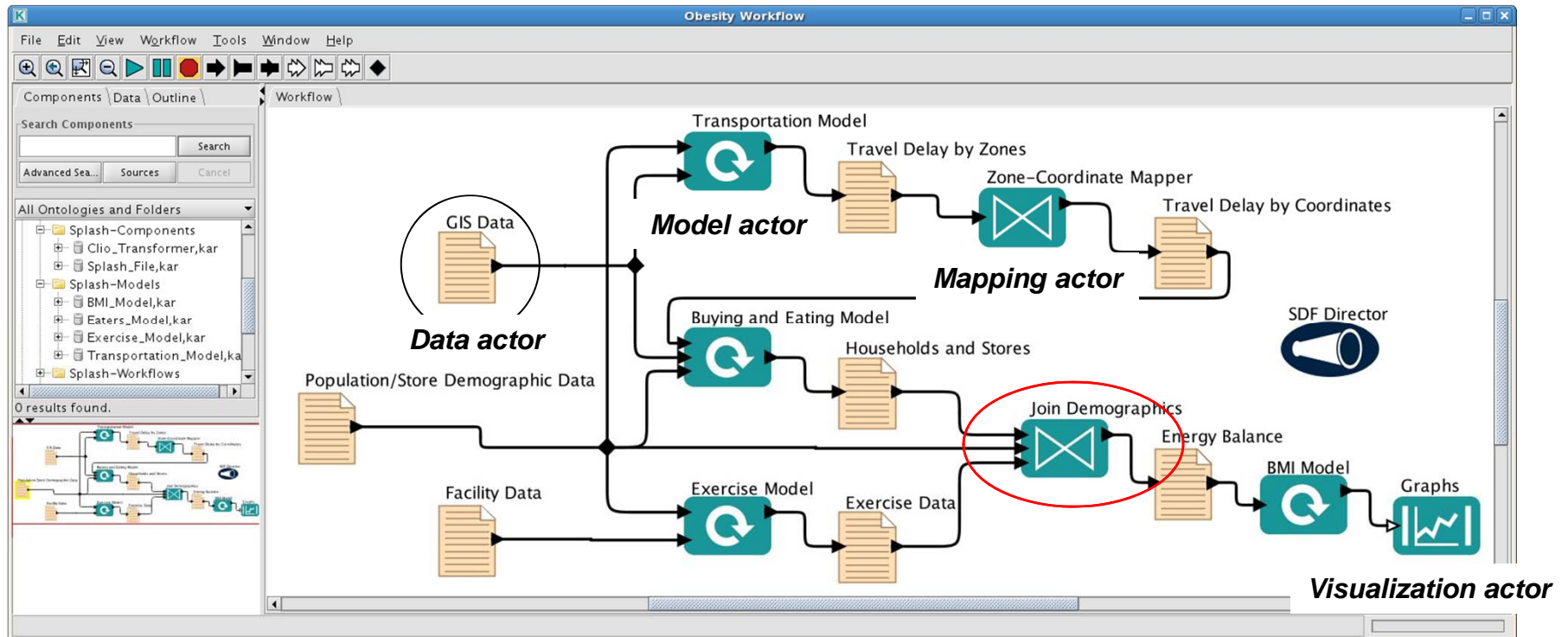
# Data Transformations Between Models

- Transformation design tools for structural (schema) and time alignments
- SADL metadata used to automatically detect mismatches
- Splash generates code for massive-scale transformation on Hadoop at simulation time



Clio++: Schema mapping
& unit corrections

Time Aligner: Time-series
harmonization

# Composite-Model Execution

# Executing a Composite Model: The Need for Runtime Efficiency

**A huge parameter space to explore (many model runs)**

- Ex: 3 models + 10 params/model + 2 vals/param = over 1 billion model runs

- Even worse for stochastic models (multiple Monte Carlo replications)

- Experimental design can help

**Each model run can be extremely time consuming**

- Large-scale, high resolution models produce and consume massive amounts of time-series and other data

- CPU-intensive computations

- Composing models (with data transformations) intensifies the problem




T-cell biology model


Regional traffic model


NCAR Community Atmosphere Model (CAM)


Agent-based social model

# Time alignment with MapReduce

$s_0$

Irregular source time series

$t_0$

Regular target time series to be calculated.

Interpolation, nearest neighbor, aggregation (since-last, since-start)

# Time alignment with MapReduce

$s_0$

Irregular source time series

Sliding window by size (n=4)

$t_0$

Regular target time series to be calculated.

Interpolation, nearest neighbor, aggregation (since-last, since-start)

# Time alignment with MapReduce

$s_0$

Irregular source time series

Sliding window by size (n=4)

$t_0$

Regular target time series to be calculated.

Interpolation, nearest neighbor, aggregation (since-last, since-start)

# Time alignment with MapReduce

$s_0$

Irregular source time series

Sliding window by size (n=4)

$t_0$

Regular target time series to be calculated.

Interpolation, nearest neighbor, aggregation (since-last, since-start)

# Time alignment with MapReduce

$s_0$

Irregular source time series

Sliding window by size (n=4)

$t_0$

Regular target time series to be calculated.

Interpolation, nearest neighbor, aggregation (since-last, since-start)

# Time alignment with MapReduce

$s_0$

Irregular source time series

Sliding window by size (n=4)

$t_0$

Regular target time series to be calculated.

Interpolation, nearest neighbor, aggregation (since-last, since-start)

# Time alignment with MapReduce

$s_0$

Irregular source time series

Sliding window by size (n=4)

$t_0$

Regular target time series to be calculated.

Interpolation, nearest neighbor, aggregation (since-last, since-start)

# Time alignment with MapReduce

$s_0$

Irregular source time series

Sliding window by size (n=4)

$t_0$

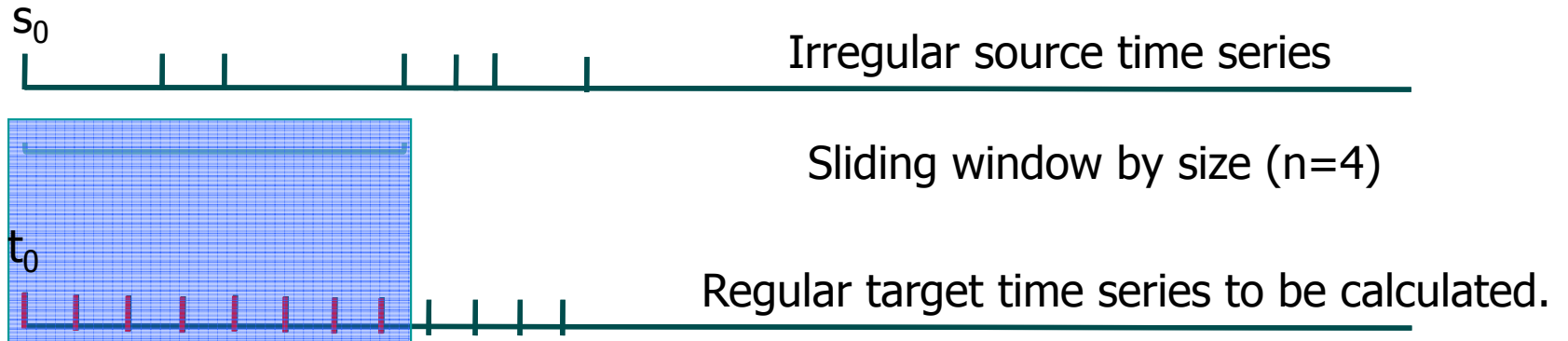Regular target time series to be calculated.



Interpolation, nearest neighbor, aggregation (since-last, since-start)

# Time alignment with MapReduce

$s_0$

Irregular source time series

Sliding window by size (n=4)
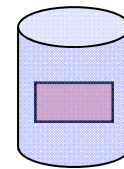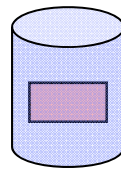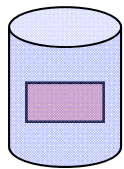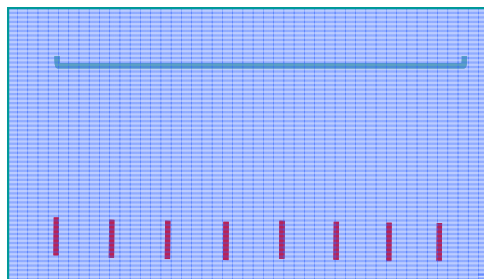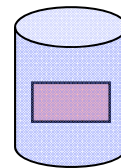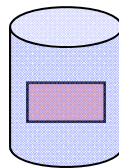
$t_0$

Regular target time series to be calculated.



Interpolation, nearest neighbor, aggregation (since-last, since-start)

# Cubic-Spline Interpolation in MapReduce

- **Recall:** Source outputs 1 tick per two days; target needs one tick per day

- **(Natural) cubic spline** widely used
  - Uniformly approximates $f$ and $f'$
  - Error of $O(h^4)$ as knot spacing $h \to 0$
  - Default method in SAS

Cubic-spline interpolation

- Given **source** and **target** time series:

$$S = \langle (s_0, d_0), (s_1, d_1), \ldots, (s_m, d_m) \rangle \ \text{ and } \ T = \langle (t_0, \tilde{d}_0), (t_1, \tilde{d}_1), \ldots, (t_n, \tilde{d}_n) \rangle$$

- Given **window** $W_i$ for $t_i$: $W_i = \langle (s_j, d_j, \sigma_j), (s_{j+1}, d_{j+1}, \sigma_{j+1}) \rangle$ where $[s_j, s_{j+1})$ contains $t_i$

$$\tilde{d}_i = f(W_i) = \frac{\sigma_j}{6 h_j} (s_{j+1} - t_i)^3 + \frac{\sigma_{j+1}}{6 h_j} (t_i - s_j)^3 + \left( \frac{d_{j+1}}{h_j} - \frac{\sigma_{j+1} h_j}{6} \right)(t_i - s_j) + \left( \frac{d_j}{h_j} - \frac{\sigma_j h_j}{6} \right)(s_{j+1} - t_i)$$

$$h_j = s_{j+1} - s_j$$

# Question: How to Compute Spline Constants?

- Must **solve** $Ax = b$ ($m$-1 rows and columns):

$$A = \begin{pmatrix} \frac{h_0 + h_1}{3} & \frac{h_1}{6} & 0 & \cdots & 0 & 0 & 0 \\ \frac{h_1}{6} & \frac{h_1 + h_2}{3} & \frac{h_2}{6} & \cdots & 0 & 0 & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots & \vdots \\ 0 & 0 & 0 & \cdots & \frac{h_{m-3}}{6} & \frac{h_{m-3} + h_{m-2}}{3} & \frac{h_{m-2}}{6} \\ 0 & 0 & 0 & \cdots & 0 & \frac{h_{m-2}}{6} & \frac{h_{m-2} + h_{m-1}}{3} \end{pmatrix} \quad b = \begin{pmatrix} \frac{d_2 - d_1}{h_1} - \frac{d_1 - d_0}{h_0} \\ \frac{d_3 - d_2}{h_2} - \frac{d_2 - d_1}{h_1} \\ \vdots \\ \frac{d_m - d_{m-1}}{h_{m-1}} - \frac{d_{m-1} - d_{m-2}}{h_{m-2}} \end{pmatrix}$$

- **Prior work**
  - Some solutions require **evenly spaced** source points
  - Some solutions require **precomputation** (somehow) of $A^{-1}$
  - Other solutions for **vector machines, MPI architectures, GPUs**
    - Require a lot of **data shuffling** (reduce steps) in Hadoop adaptation
    - Example: **Parallel Cyclic Reduction (PCR)** uses $\log_2 m$ map-reduce jobs

- **Our approach:  minimize** $L(x) = \left\| Ax - b \right\|_2^2 = \sum_i (A_i. x - b_i)^2 = \sum_i L_i(x)$

# Our Solution: Distributed Stochastic Gradient Descent (DSGD)

- Originally for **matrix completion**, e.g., Netflix ratings problem [GHS KDD11]



- Uses **stochastic gradient descent (SGD)** to minimize $L$
  - **Deterministic** gradient descent (**DGD**): $$x^{(n+1)} = x^{(n)} - \varepsilon_n L'(x^{(n)})$$

  where $L'(x^{(n)}) = \sum_{i=1}^{m-1} L_i'(x^{(n)})$

  - **Stochastic** gradient descent: $$x^{(n+1)} = x^{(n)} - \varepsilon_n \hat{L}'(x^{(n)})$$

  where $\hat{L}'(x^{(n)}) = (m-1)L_I'(x^{(n)})$

  and $I$ is randomly chosen from $[1..m-1]$

  - Avoids getting stuck at local minima
  - **Problem**: SGD is not a parallel algorithm



- **Idea**: run SGD on subsets (strata) of rows, randomly switch strata; choose "sparse" strata that allow parallel execution of SGD
  - Converges to overall solution with probability 1 under mild conditions

# Choosing Strata

**Goal:** Permit parallel execution of SGD within each stratum

**Key observation:** $L'_i(x) = \begin{pmatrix} 0 & \ldots & 0 & u_{i,i-1} & u_{i,i} & u_{i,i+1} & 0 & \ldots & 0 \end{pmatrix}$

Updating $x_i$ only affects (and is affected by) $x_{i-1}$ and $x_{i+1}$

where $u_{i,j} = 2a_{i,j}(a_{i,i-1}x_{i-1} + a_{i,i}x_i + a_{i,i+1}x_{i+1})$

**Stratum choice:**

- Can implement as map-only Hadoop job (almost no data shuffling)
- Exploit discrepancy between logical splits and physical blocks

**Empirical study:**

- 2x-3x faster than best-of-breed PCR alg.
- 10 scans vs $\log m$ for PCR
- PCR requires extra sort
- PCR requires massive data shuffling (network bottleneck)

node 1
node 2
node 3

# Speeding up Composite Simulations: Result Caching

**Motivating example: Two models in series, 100 reps**

Model 1 → Model 2

Deterministic       Stochastic

- Naïve approach: execute composite model (i.e., Models 1 & 2) 100 times

- A better approach:

Model 1 → Cache → Model 2

- Execute Model 1 once and cache result
- Read from cache when executing Model 2

**Question: Can result-caching idea be generalized?**

Identical

# General Method for Two Stochastic Models in Series



**Goal: Estimate** $\theta = E[Y_2]$ **based on n replications**

**Result-caching approach:**

1. Set $m_n = \lceil \alpha n \rceil$ for some $\alpha \in (0,1]$ (the re-use factor)    Ex: n=10, $m_n$ = 4

2. Generate $m_n$ outputs from Model 1 and cache them

3. To execute Model 2, cycle through Model 1 outputs

4. Estimate $\theta$ by $\theta_n = \sum_{i=1}^{n} Y_{2;i} / n$

# Optimizing the Re-Use Factor for Maximum Efficiency

**Q: How to trade off cost and precision?**

- Assume a (large) fixed computational budget $c$
- Random cost model: correlated pair $(\tau_i, Y_i)$
  - $\tau_i =$ (random) cost of producing an observation $Y_i$
  - $N(c) =$ # of observations of $Y_2$ generated under $c$
  - $\hat{\theta}(c) = \sum_{j=1}^{N(c)} Y_{2;j} / N(c)$

- Approx. distribution of $\hat{\theta}(c)$:

variance $= g(\alpha) / c$

$\theta$

$$g(\alpha) = (\alpha E[\tau_1] + E[\tau_2]) \left\{ \mathrm{Var}[Y_2] + (2r_\alpha - \alpha r_\alpha (r_\alpha + 1)) \mathrm{Cov}[Y_2, \tilde{Y}_2] \right\} \qquad r_\alpha = \lfloor 1/\alpha \rfloor$$

(cost per obs.) x (contributed variance per obs.)

# The Optimal Re-Use Factor

**Optimal solution**

- Assume that $\mathrm{Cov}[Y_2, \tilde{Y}_2] \geq 0$

- Optimal value of $\alpha$:

$$\alpha^* \approx \left( \frac{E[\tau_2] / E[\tau_1]}{\left( \mathrm{Var}[Y_2] / \mathrm{Cov}[Y_2, \tilde{Y}_2] \right) - 1} \right)^{1/2}$$

(truncate at 1/n or 1)

**Observations**

- If E[Model 1 cost] >> E[Model 2 cost], then high re-use of output

- If Model 2 insensitive to Model 1 (Cov << Var), then high re-use

- If Model 1 is deterministic (Cov = 0), then total re-use

# Experiment Management (and Optimization)



**Splash Platform**

Models
Data

SADL

- Model and Data Curation
- Model and Data Discovery
- Model Composition
- Composite-Model Execution
- Experiment Management

Analysis
Visualization

DBMS, Hadoop, Visualization Tools, Information-Integration Tools, Stats Packages

# Experiment Design and Efficiency

Trades off execution cost versus level of detail that can be estimated

Coarse resolution is OK for sensitivity analysis etc.

Resolution III design

| Run | Factors | | | | | | |
|---|---|---|---|---|---|---|---|
| | $A$ | $B$ | $C$ | $D=AB$ | $E=AC$ | $F=BC$ | $G=ABC$ |
| def | -1 | -1 | -1 | 1 | 1 | 1 | -1 |
| afg | 1 | -1 | -1 | -1 | -1 | 1 | 1 |
| beg | -1 | 1 | -1 | -1 | 1 | -1 | 1 |
| abd | 1 | 1 | -1 | 1 | -1 | -1 | -1 |
| cdg | -1 | -1 | 1 | 1 | -1 | -1 | 1 |
| ace | 1 | -1 | 1 | -1 | 1 | -1 | -1 |
| bcf | -1 | 1 | 1 | -1 | -1 | 1 | -1 |
| abcdefg | 1 | 1 | 1 | 1 | 1 | 1 | 1 |

**Example: 1st-order polynomial metamodel for scaled data (7 factors)**

$$Y = \beta_0 + \beta_1 x_1 + \cdots \beta_7 x_7$$
$$+ \beta_{1;2} x_1 x_2 + \cdots + \beta_{6;7} x_6 x_7 + \beta_{1;2;3} x_1 x_2 x_3 + \cdots + \text{noise}$$
$$x_1, \ldots, x_7 \in \{-1, 1\} \quad \text{(full factorial = 128 runs)}$$

Fractional-factorial experimental designs

| To estimate | If you can ignore | Resolution | # runs |
|---|---|---|---|
| Main effects | All high-order effects | III | 8 |
| Main effects | 3rd-order and higher | IV | 16 |
| Main effects + 2-way interactions | 3rd-order and higher | V | 64 |

# Running experiments in Splash

**Goal**

- Provide a facility that gives the illusion of executing **one** coherent simulation model



**Main Challenges**

- Automate the coordination between experiment conditions and inputs to different submodels.

- Automate the combination of different replications of different submodels.

# Example: Healthcare Payer Model

**Composition of two models**

- Emory/Georgia Tech Predictive Health Institute model [Park et al. 2012]
    - Simple agent-based model of prevention and wellness program
    - For investigation of payment systems (capitated vs outcome-based)
- Simple logarithmic random walk model of interest & inflation rates

# Experiment Manager (Specifying Experimental Factors)

SADL

```
<attribute name="paymentModel"
  measurement_type="numerical"
  missing_data="0"
  experiment_default_values=""
  experiment_factor="true"
  datatype="double"
  random_seed="false" />
```

**Splash Experiment Manager**

**Experiment Factors**

Select experiment factors

| PHI_Model | Financial_Rate_Model |

▽ PHI_Model.CommandLine(1)

SADL 🔍

☐ population — Value ∨ /default_dir/populationdata.csv

▽ PHI_Model.parameters(12)

SADL 🔍

☑ paymentModel — Value ∨ 0.5 1.0 1.5

☐ capitationPerParticipant — Value ∨ 500

☐ costModel — Value ∨ 1

☐ terminalAge — Value ∨ 65

☐ diabetesRiskThreshhold — Value ∨ 0.25

☐ diabetesRiskReduction — Value ∨ 0.55

[ < Back ] [ Next > ] [ Finish ]

**GUI collects simulation parameters from all component models experiment_factor = TRUE in SADL file**

**User selects values for each experiment factor**

**User selects subset of parameters as experiment factors**

# Experiment Design in Splash



**Design Persistence**

EML

<model name= PHI>…

<factor name="Tage">

<values>"65"</values>
<values>"85"</values>
</factor>…
<rep n="10">…
</experiment>

**Editable design**

(Factor values and
# of Monte Carlo reps
for each condition)

**Execution Engine**

# Experiment Manager (Running an Experiment)

**Technical challenges include:**

- **Routing parameter values to models**
  - Different sources: command line args, parameter files, stdin, ...
  - Synthesizing the parameter files that a model expects (templating)

- **Managing PRNG seeds**
  - Avoiding cycle overlaps
  - PRNG info in SADL file
  - Diagnostics (future work)

> **Experiment Manager invokes Splash execution engine to run experiments**

Pro...

! Execute Experiments

running experiment 3, replicaton 9

```
writing stdout to /Users/pmac/Desktop/exec-dir/experiments/Financial Rate Model-06
writing stdout to /Users/pmac/Desktop/exec-dir/experiments/Financial Rate Model-06
writing stdout to /Users/pmac/Desktop/exec-dir/experiments/Financial Rate Model-06
writing stdout to /Users/pmac/Desktop/exec-dir/experiments/Financial Rate Model-06
writing stdout to /Users/pmac/Desktop/exec-dir/experiments/Financial Rate Model-06
writing stdout to /Users/pmac/Desktop/exec-dir/experiments/Financial Rate Model-06
writing stdout to /Users/pmac/Desktop/exec-dir/experiments/Financial Rate Model-06
writing stdout to /Users/pmac/Desktop/exec-dir/experiments/Financial Rate Model-06
```

Cancel    OK

> **Intermediate and final outputs can be saved in a file tree for**
> - **Provenance tracking**
> - **Traceability**
> -  **Drill down**

# Template-Based Data File Generation Process

...
<attributes>
<attribute name=temperature
Datatype=numeric…/>
<attribute name=pressure
Datatype=numeric…/>
...

Input data for city of Detroit
Temperature=$$temperature$$&&%0.1&&
Pressure=$$pressure$$&&%0.1&&
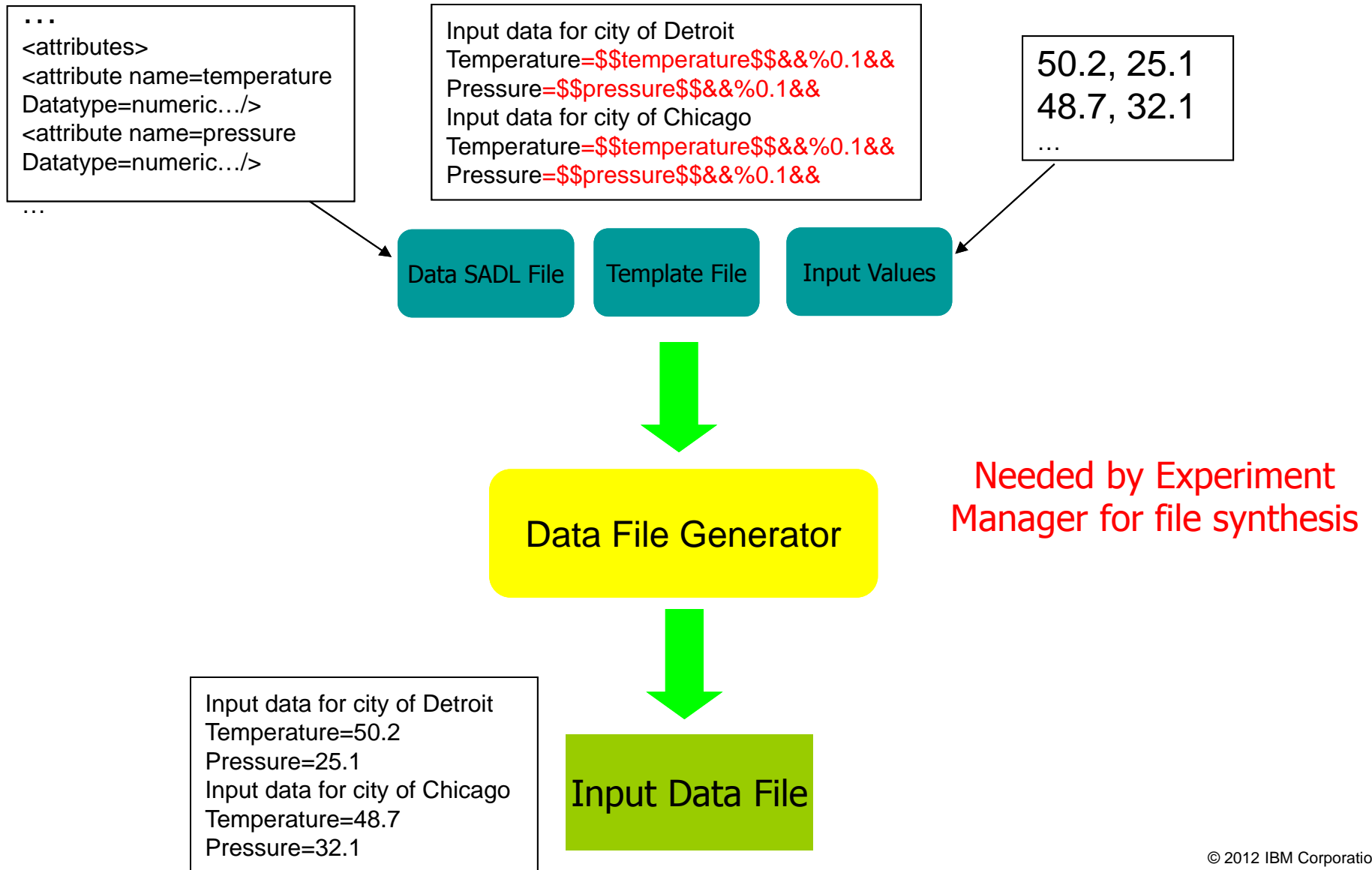Input data for city of Chicago
Temperature=$$temperature$$&&%0.1&&
Pressure=$$pressure$$&&%0.1&&

50.2, 25.1
48.7, 32.1
...

**Data SADL File**   **Template File**   **Input Values**

## Data File Generator

**Needed by Experiment
Manager for file synthesis**

Input data for city of Detroit
Temperature=50.2
Pressure=25.1
Input data for city of Chicago
Temperature=48.7
Pressure=32.1

## Input Data File

# Template-Based Data Extraction Process

Input data for city of Detroit
Temperature=$$temperature$$&&%0.1&&
Pressure=$$pressure$$&&%0.1&&
Input data for city of Chicago
Temperature=$$temperature$$&&%0.1&&
Pressure=$$pressure$$&&%0.1&&

Input data for city of Detroit
Temperature=50.2
Pressure=25.1
Input data for city of Chicago
Temperature=48.7
Pressure=32.1

**Template File**

**Unstructured Data File**

**Data Extractor**

Needed to extract performance measures of interest for optimization, visualization, etc.

Extracted Values

50.2, 25.1
48.7, 32.1
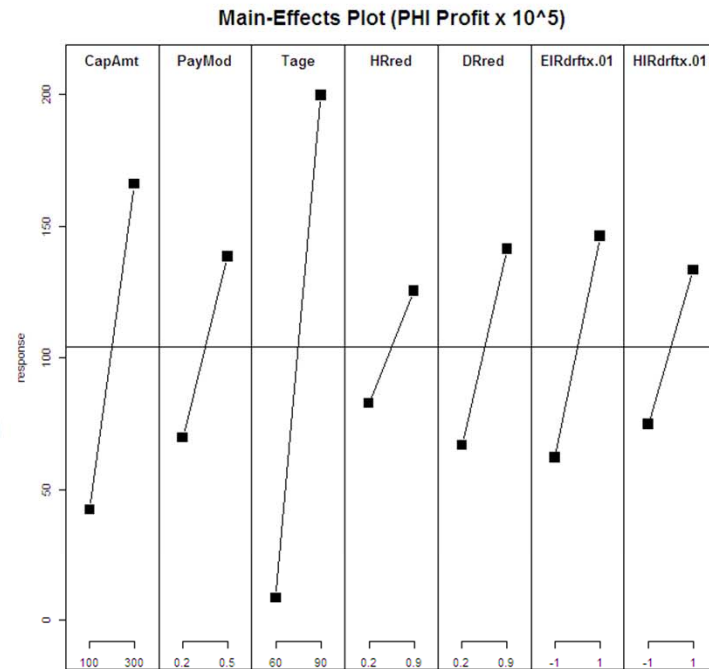…

# Efficient Sensitivity Analysis

- Main-effects plots:
  - High/low values
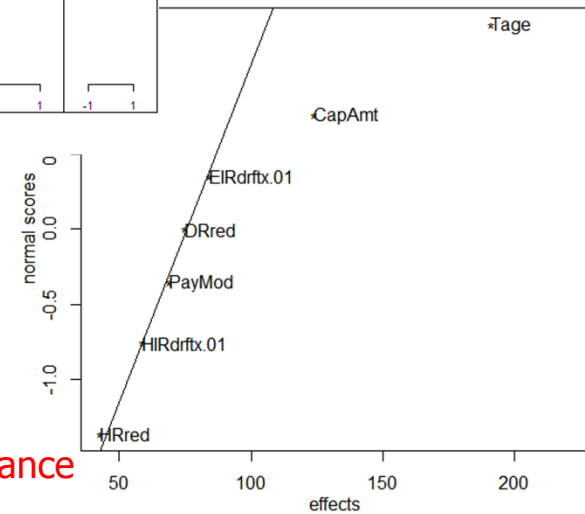  - Orthogonal fractional factorial experiment design (160 vs 2560 runs)

PHI healthcare payer model + interest-rate model

(Park et al., *Service Science*, 2012)



**Main-Effects Plot (PHI Profit x 10^5)**

Identify the most important profit drivers
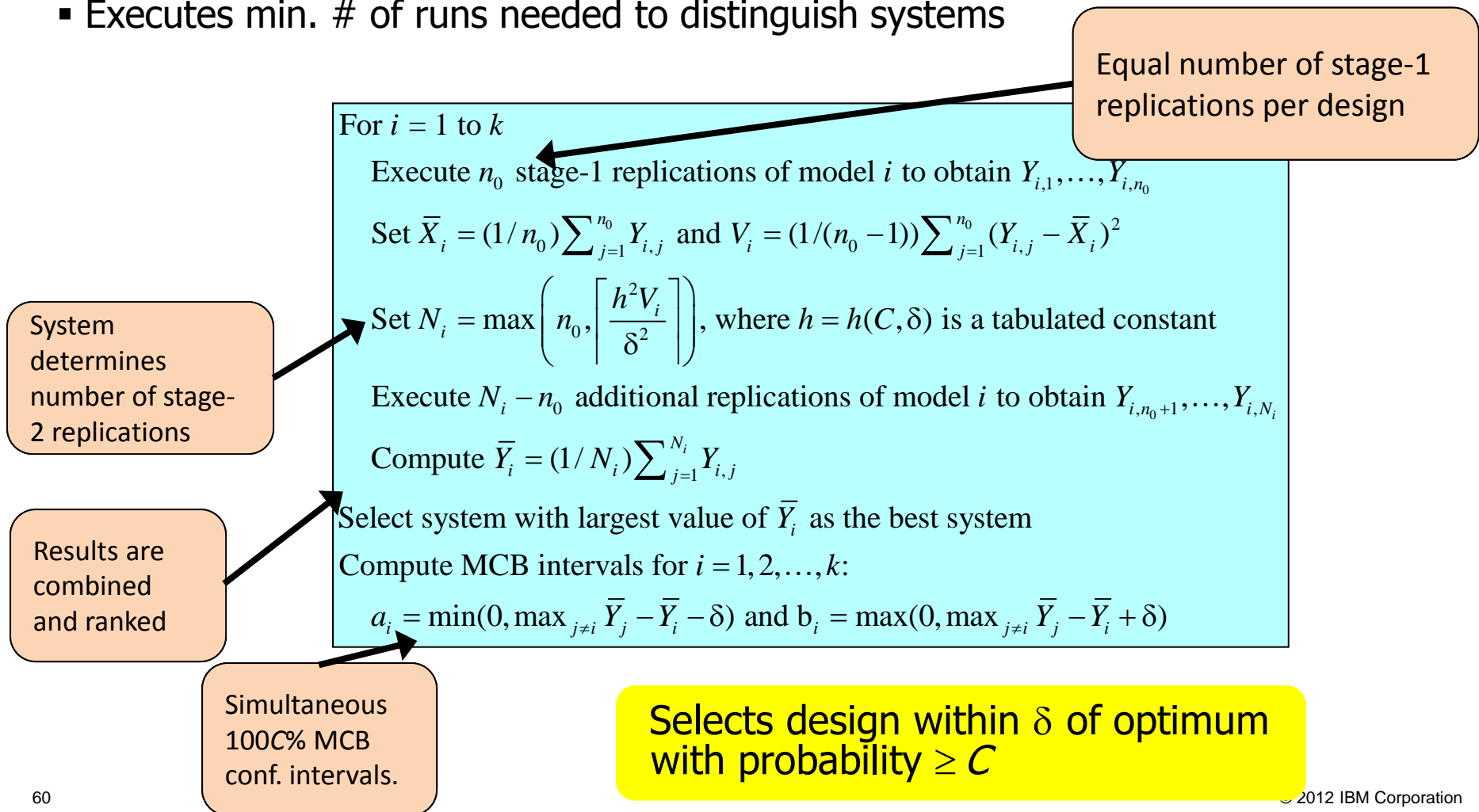(CapAmt & Tage)

**Normal Effects Plot for PHI Profit**

Check statistical significance of graphical results
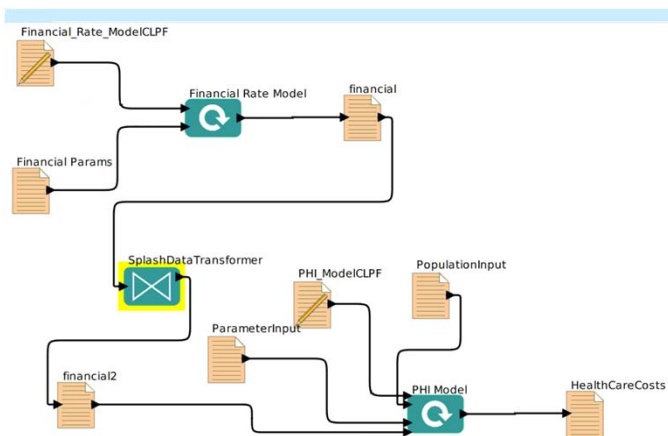
# Optimization Functionality: Ranking and Selection

- Rinott procedure for finding best among small number of designs
- Executes min. # of runs needed to distinguish systems

Equal number of stage-1 replications per design

For $i = 1$ to $k$

Execute $n_0$ stage-1 replications of model $i$ to obtain $Y_{i,1}, \ldots, Y_{i,n_0}$

Set $\bar{X}_i = (1/n_0)\sum_{j=1}^{n_0} Y_{i,j}$ and $V_i = (1/(n_0-1))\sum_{j=1}^{n_0}(Y_{i,j} - \bar{X}_i)^2$

Set $N_i = \max\left(n_0, \left\lceil \dfrac{h^2 V_i}{\delta^2} \right\rceil\right)$, where $h = h(C, \delta)$ is a tabulated constant

Execute $N_i - n_0$ additional replications of model $i$ to obtain $Y_{i,n_0+1}, \ldots, Y_{i,N_i}$

Compute $\bar{Y}_i = (1/N_i)\sum_{j=1}^{N_i} Y_{i,j}$

Select system with largest value of $\bar{Y}_i$ as the best system

Compute MCB intervals for $i = 1, 2, \ldots, k$:

$a_i = \min(0, \max_{j \neq i} \bar{Y}_j - \bar{Y}_i - \delta)$ and $b_i = \max(0, \max_{j \neq i} \bar{Y}_j - \bar{Y}_i + \delta)$

System determines number of stage-2 replications

Results are combined and ranked

Simultaneous 100$C$% MCB conf. intervals.

Selects design within $\delta$ of optimum with probability $\geq C$
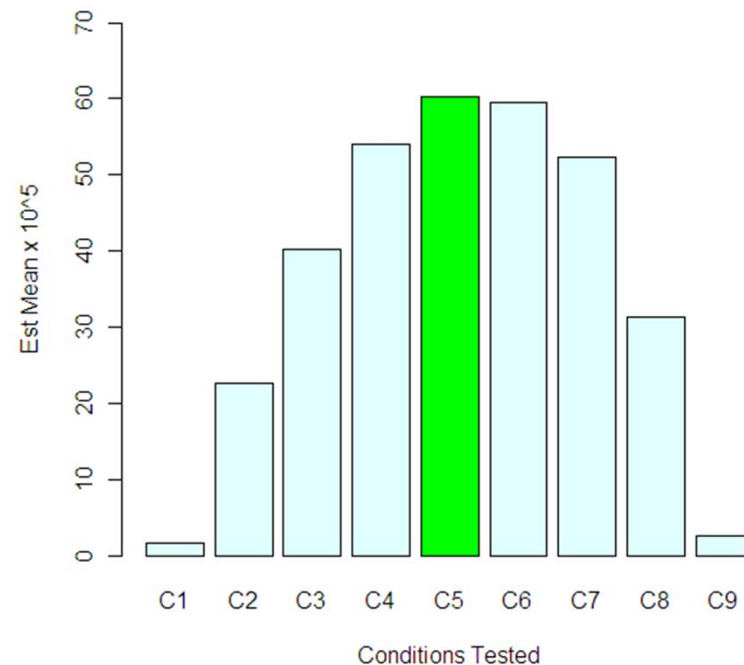
# Results for PHI Profitability: Estimated Best System

"Conditions" = payment schemes for wellness program
(0 = full capitation, 1 = pay-for-outcome)

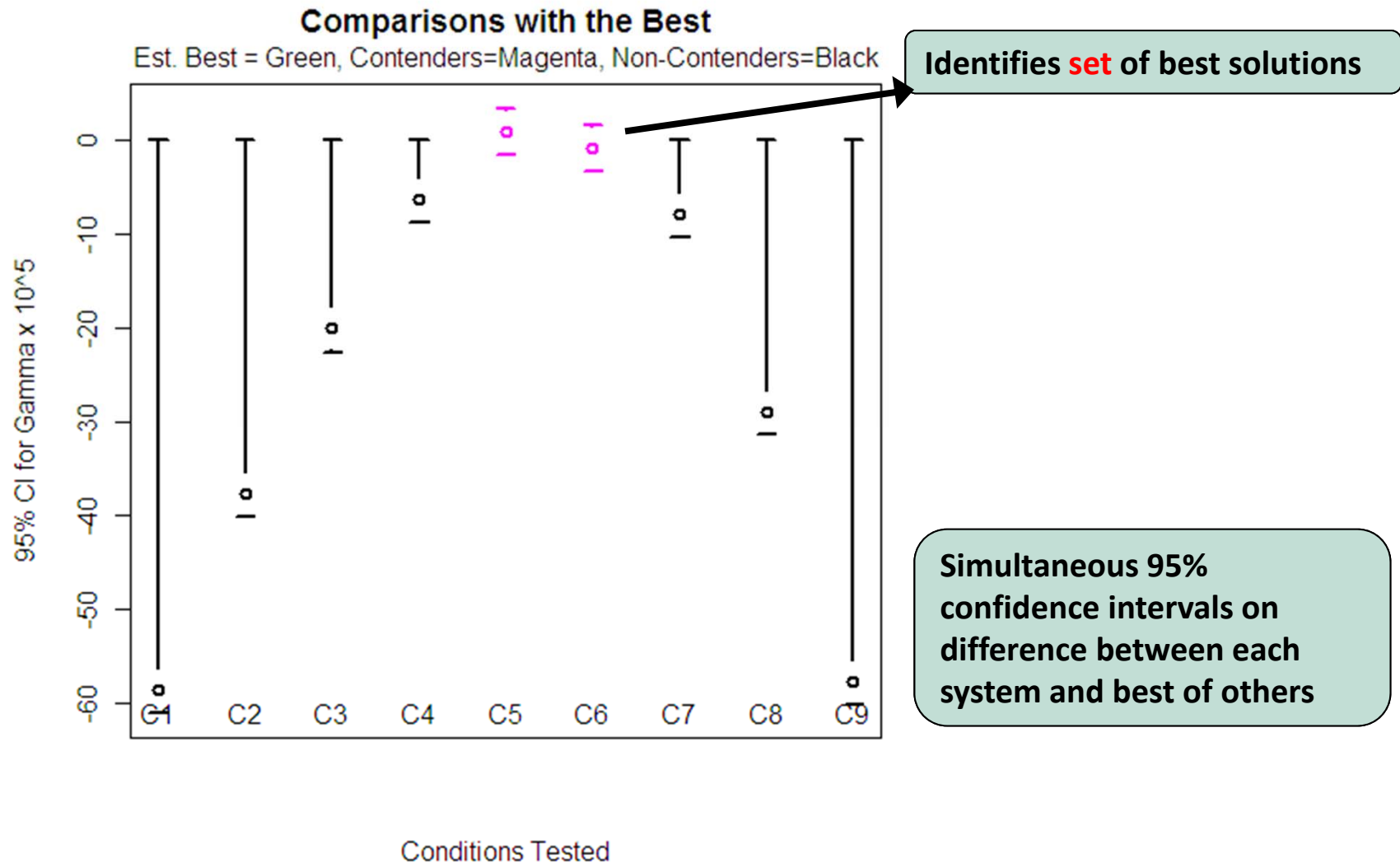Look at weighted schemes: 0.1, 0.2, … , 0.9



PHI healthcare payer model +
interest-rate model
(Park et al., *Service Science*, 2012)

**Est Best System (in green)**



With prob = 95%, C5 = 0.5 is the "best system"
(within indifference zone = $250K)

# Results Continued: Multiple Comparisons with the Best



Identifies **set** of best solutions

Simultaneous 95% confidence intervals on difference between each system and best of others

# Simulation Metamodeling (Joint Work with SJSU CAMCOS)

**"Simulation on demand"**

1. Run simulations in advance to get values at multiple "design points"

2. Fit a (stochastic) response surface

3. Decision maker can explore surface in real time

4. Can apply stochastic optimization techniques to find peaks and valleys

5. Can use for factor screening

**Technique: Stochastic Kriging**
(Ankenman et al., *Oper. Res.*, 2010)

- Robust, global fit

- Gives approximate model response + uncertainty estimates (MSE)

- Efficient allocation to of runs to minimize integrated mean-square error (IMSE)

- Metamodel added to Splash repository



Image: SJSU CAMCOS

Models uncertainty due to both interpolation and simulation variability

# Assessment of PHI metamodel



Metamodel gives good approximation to real results (1.6% error in this example)

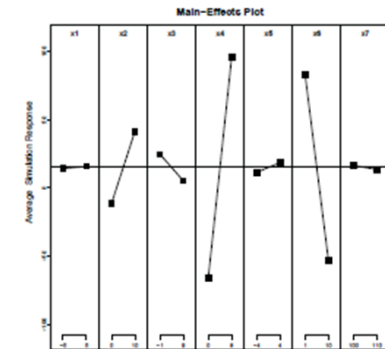Faster by over two orders of magnitude

# Factor screening (Joint with SJSU CAMCOS)

## Goal: identify most important subset of drivers
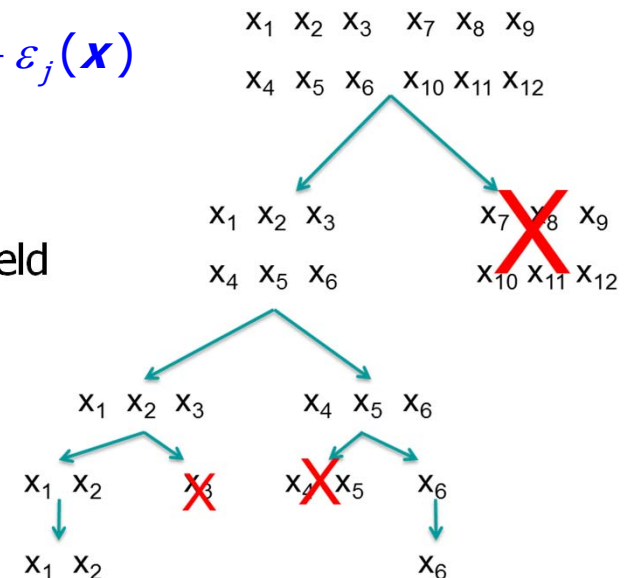
- Drivers captured in metamodel parameters

## Ex: Linear models $Y(x) = \beta_0 + \beta_1 x_1 + \cdots + \beta_7 x_7 + \varepsilon$

- Main effects used for screening

- For Gaussian noise, positive effects: <span style="color:red">sequential bifurcation</span>

## Ex: Gaussian process models $Y_j(x) = \beta_0 + M(x) + \varepsilon_j(x)$

- Special case of stochastic kriging

- $\varepsilon_j(x)$ = simulation noise
- $M(x)$ = interpolation uncertainty, modeled as Gaussian field
  - For any $x_1, x_2, \ldots, x_r$ vector $V = (M(x_1), \ldots, M(x_r))$ is multivariate normal
  - $\mathrm{Cov}[M(x_i), M(x_j)] = \tau^2 \prod_{k=1}^{n} \exp(-\theta_k (x_{i,k} - x_{j,k})^2)$
- Small $\theta_k \Rightarrow$ small effect of $k^{\text{th}}$ factor

- Bayesian "posterior quantiles" method for screening
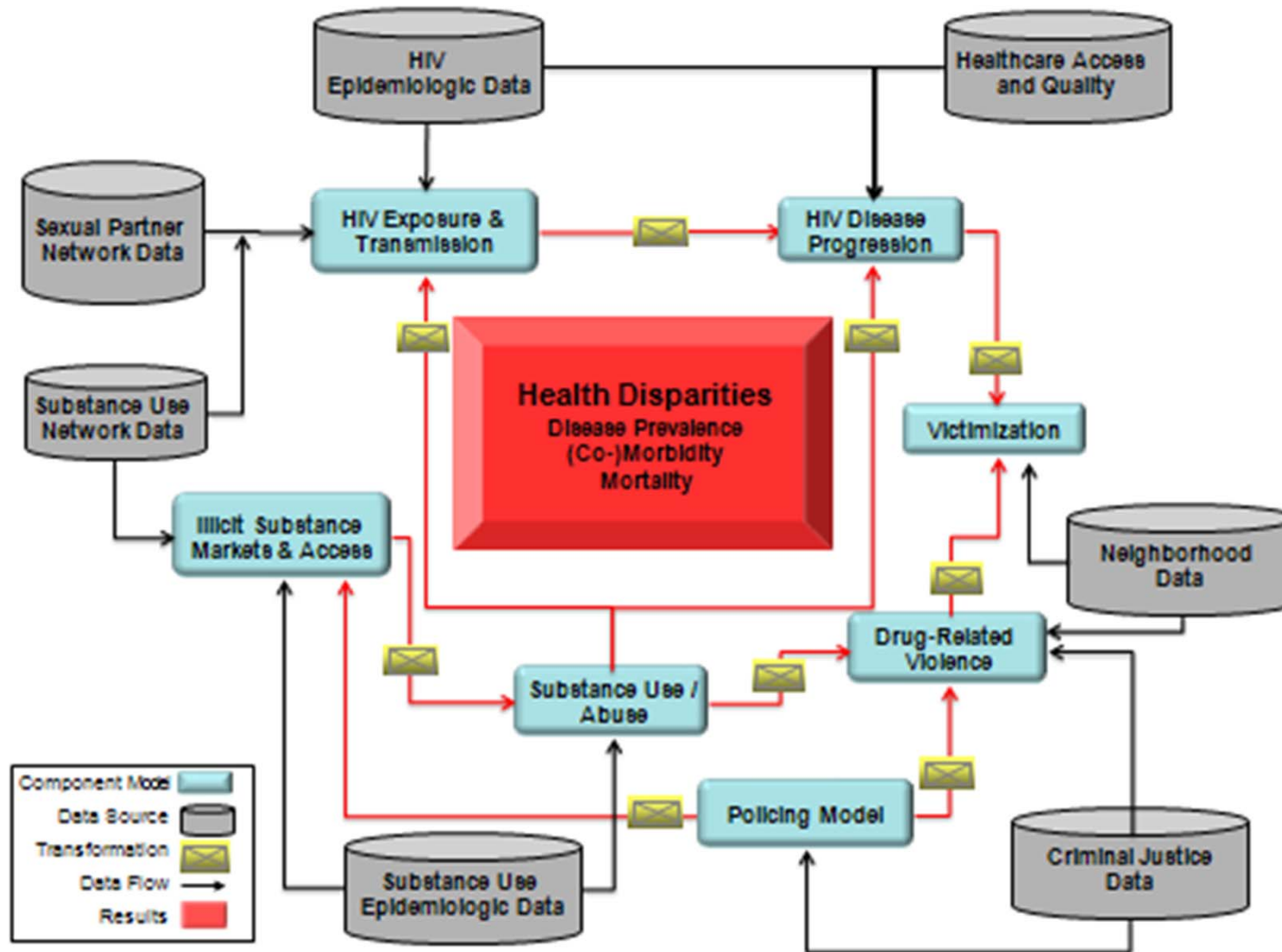
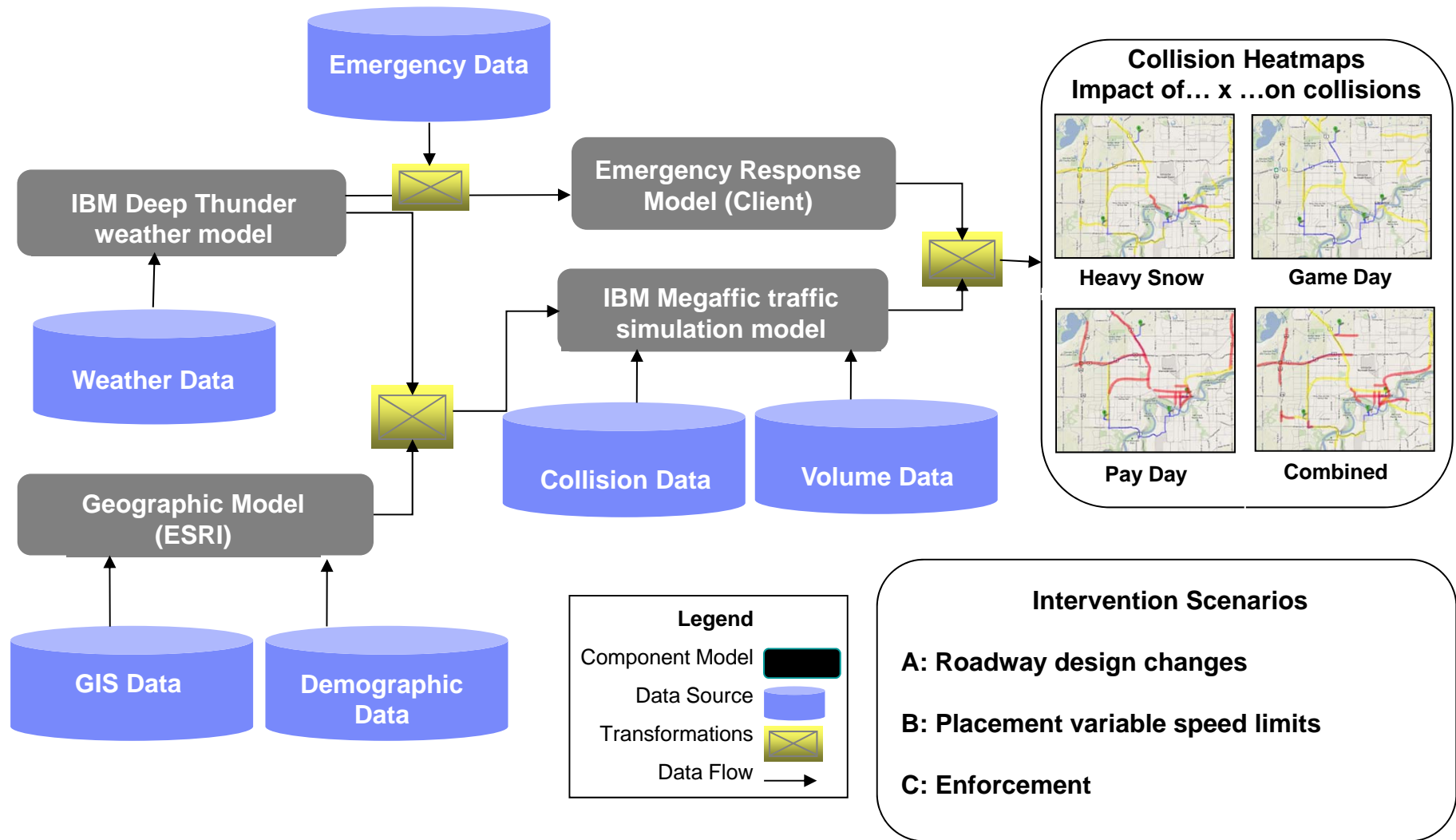# Some Potential Splash Applications

# Multi-level, End-to-End Modeling



Business Models

Socio-Economic Models

**Healthcare Ecosystem**
(Society)

**System Structure**
(Organizations)

**Delivery Operations**
(Processes)

**Clinical Practices**
(People)

Careflow Models
(*Flow of Patients, Money, Information*)

Disease Progression Models

Lever1   Lever2   Lever3

Policy "Flight Simulator"

Personalized Medicine
(*Targeted interventions*)

Rouse, W. B. & Cortese, D. A. (2010). Introduction, in W. B. Rouse & D. A. Cortese (Eds.), *Engineering the System of Healthcare Delivery*. IOS Press.

# Cross-domain, Syndemic Modeling

# Composite model for traffic safety

**Emergency Data**

**IBM Deep Thunder weather model**

**Emergency Response Model (Client)**

**Weather Data**

**IBM Megaffic traffic simulation model**

**Geographic Model (ESRI)**

**Collision Data**

**Volume Data**

**GIS Data**

**Demographic Data**

**Collision Heatmaps
Impact of… x …on collisions**

**Heavy Snow**    **Game Day**

**Pay Day**    **Combined**

## Legend

| Component Model | |
| Data Source | |
| Transformations | |
| Data Flow | → |

**Intervention Scenarios**

**A: Roadway design changes**

**B: Placement variable speed limits**

**C: Enforcement**

# Open Research Questions

# How to Determine User Requirements?

## Common to Analysts and Scientists

- Examine schemas (data) and variables (models) prior to selection
- Compare output of simulation results to examine trade-offs and simulation selection
- Dashboard with summary of models and data sources used to run a simulation



## Specific to Analysts

- Guidance and recommendations
- Pre-defined templates for simulation set-up and analyzing simulation output
- Recommendations for what template to use and the steps to run a simulation
- Recommended output visualization – suggest one chart style would be better than another style to explain relationships in data



## Specific to Scientists

- Feature to assess the veracity and provenance of model and data sources
- Ability to upload their own sources to supplement the existing sources
- High levels of interaction with the models & data when previewing search results prior to running the simulation

# Database Research++

- **Data search → model-and-data search**
  - Find compatible models, data, and mappings (using metadata)
  - Involves semantic search technologies, repository management, privacy and security

- **Data integration → model integration**
  - Simulation-oriented data mapping
  - Geospatial alignment [e.g., Howe & Maier 2005]
  - Hierarchical models with different resolutions
  - Complex data transformations (e.g., raw simulation output to histogram)

- **Query optimization → simulation-experiment optimization**
  - Optimally configure workflow among distributed data and models
  - Factoring common operations across different mappings in the workflow
  - Avoiding redundant computations across experiments (e.g., result caching)
  - Statistical issues: managing pseudorandom numbers and Monte Carlo replications

# Some Deep Problems

- **Causality approximation**
  - Fixed-point + perturbation approaches
  - System support
  - Theoretical support

- **Deep collaborative analytics**
  - Visualizing and mining the results
  - Understanding and explaining results:
    - Provenance [e.g., J. Friere et al.]
    - Root-cause analysis
  - Trusting results
    - Model validation
    - ManyEyes++, Swivel++

Transportation Model → Buying & Eating Model

$$\dot{f}_n(t) = \Lambda_1\left(f_n(t), g_{n-1}(t)\right)$$
$$\dot{g}_n(t) = \Lambda_2\left(f_{n-1}(t), g_n(t)\right)$$

$$\left.\begin{array}{l}\dot{f}(t) = \Lambda_1\left(f(t), g(n\Delta t)\right)\\ \dot{g}(t) = \Lambda_2\left(f(n\Delta t), g(t)\right)\end{array}\right\} \text{ for } t \in \left[n\Delta t, (n+1)\Delta t\right)$$

# Conclusion



- **Splash:**
  - composition of heterogeneous models and data
    to support cross-disciplinary decision making in complex systems
  - Loose coupling of models through data exchange
  - Combines data-integration, simulation, and workflow technologies

- **Key features**
  - SADL metadata language for curation and functionality
  - Automated detection of data mismatches
  - Semi-automated design of scalable data transformations (schema and time alignment)
  - Runtime accelerators
    - MapReduce framework for scalable data transformations
    - Map-only Hadoop method for cubic-spline interpolation
    - Result-caching to minimize # of model executions
  - Experiment-manager allows sensitivity analysis, factor screening and optimization
  - Simulation metamodeling for real-time model exploration

- **Many open research questions!**

# Questions?



Splash project page:
http://researcher.watson.ibm.com/researcher/view_project.php?id=3931

# Backup Slides

# Splash Technology for Loose Coupling via Data Exchange



SADL metadata language

Kepler adapted for model composition

Design-time components

Run-time components:

- Kepler adapted for model execution
- Experiment Manager
  (sensitivity analysis, metamodeling, optimization)

Data transformation tools:
  - Clio++
  - Time Aligner (MapReduce algorithms)
  - Templating mechanism

# Distributed SGD, Continued

- Divide the $m$-1 rows into three **strata**: $U^1$, $U^2$, $U^3$

- **Decompose** loss function:

  $$L(x) = \tfrac{1}{3} L^1(x) + \tfrac{1}{3} L^2(x) + \tfrac{1}{3} L^3(x)$$

  where $L^s(x) = 3\sum_{i \in U^s} L_i(x)$



- Define (random) **stratum sequence** $\gamma_1, \gamma_2, \ldots$

- Execute SGD **w.r.t.** $L^{\gamma_k}$ at $k^{th}$ step in parallel

- **Theorem:** Suppose that $x^* = A^{-1}b$ exists and

  - $\varepsilon_n = O(n^{-\alpha})$ for some $\alpha \in (0.5, 1)$
  - $(\varepsilon_n - \varepsilon_{n+1}) / \varepsilon_n = O(\varepsilon_n)$
  - $\{\gamma_n : n \geq 0\}$ is regenerative

    with $E[\tau_1^{1/\alpha}] < \infty$ and $E[X_1(s)] = 0$

  - Stratum sequence occasionally restarts probabilistically
  - Time $\tau$ between restarts has finite $1/\alpha$ moment
  - Sequence spends $\approx 1/3$ of its time on each stratum

  Then $x^{(n)} \to x^*$ with probability 1

- **Proof:** [GHS11] + Liapunov-function argument

# Hadoop Implementation

- Physical **blocks** and logical **splits**
  - InputFormat operator creates splits (one split per mapper)
  - A split is mostly on one block
  - Splits are usually disjoint
  - Map job: each mapper first obtains all split data (small amount of data movement)
  - Reduce job: massive shuffling of data over network

- We allow splits to **overlap** by two rows

- DSGD is implemented as a **map-only** job (no data shuffling!)

| | | | | |
|---|---|---|---|---|
| 1 | $a_{1,1}$ | $a_{1,2}$ | $b_1$ | $x_1$ |
| 2 | $a_{2,1}$ | $a_{2,3}$ | $b_2$ | $x_2$ |
| 3 | $a_{3,2}$ | $a_{3,4}$ | $b_3$ | $x_3$ |
| 4 | $a_{4,3}$ | $a_{4,5}$ | $b_4$ | $x_4$ |
| 5 | $a_{5,4}$ | $a_{5,6}$ | $b_5$ | $x_5$ |
| 6 | $a_{6,5}$ | $a_{6,7}$ | $b_6$ | $x_6$ |
| 7 | $a_{7,6}$ | $a_{7,8}$ | $b_7$ | $x_7$ |
| 8 | $a_{8,7}$ | $a_{8,9}$ | $b_8$ | $x_8$ |
| 9 | $a_{9,8}$ | $a_{9,10}$ | $b_9$ | $x_9$ |
| 10 | $a_{10,9}$ | $a_{10,11}$ | $b_{10}$ | $x_{10}$ |
| 11 | $a_{11,10}$ | $a_{11,12}$ | $b_{11}$ | $x_{11}$ |
| 12 | $a_{12,11}$ | $a_{12,13}$ | $b_{12}$ | $x_{12}$ |
| 13 | $a_{13,12}$ | $a_{13,14}$ | $b_{13}$ | $x_{13}$ |

split 1

split 2

stratum $s = 1$

(mapper 2 modifies $x_7$)

# Hadoop Implementation

- Physical **blocks** and logical **splits**
  - InputFormat operator creates splits (one split per mapper)
  - A split is mostly on one block
  - Splits are usually disjoint
  - Map job: each mapper first obtains all split data (small amount of data movement)
  - Reduce job: massive shuffling of data over network

- We allow splits to **overlap** by two rows

- DSGD is implemented as a **map-only** job (no data shuffling!)

| | | | | |
|---|---|---|---|---|
| 1 | $a_{1,1}$ | $a_{1,2}$ | $b_1$ | $x_1$ |
| 2 | $a_{2,1}$ | $a_{2,3}$ | $b_2$ | $x_2$ |
| 3 | $a_{3,2}$ | $a_{3,4}$ | $b_3$ | $x_3$ |
| 4 | $a_{4,3}$ | $a_{4,5}$ | $b_4$ | $x_4$ |
| 5 | $a_{5,4}$ | $a_{5,6}$ | $b_5$ | $x_5$ |
| 6 | $a_{6,5}$ | $a_{6,7}$ | $b_6$ | $x_6$ |
| 7 | $a_{7,6}$ | $a_{7,8}$ | $b_7$ | $x_7$ |
| 8 | $a_{8,7}$ | $a_{8,9}$ | $b_8$ | $x_8$ |
| 9 | $a_{9,8}$ | $a_{9,10}$ | $b_9$ | $x_9$ |
| 10 | $a_{10,9}$ | $a_{10,11}$ | $b_{10}$ | $x_{10}$ |
| 11 | $a_{11,10}$ | $a_{11,12}$ | $b_{11}$ | $x_{11}$ |
| 12 | $a_{12,11}$ | $a_{12,13}$ | $b_{12}$ | $x_{12}$ |
| 13 | $a_{13,12}$ | $a_{13,14}$ | $b_{13}$ | $x_{13}$ |

split 1

split 2

stratum $s = 2$

(mapper 2 modifies $x_7$)

# Hadoop Implementation

- Physical **blocks** and logical **splits**
  - InputFormat operator creates splits (one split per mapper)
  - A split is mostly on one block
  - Splits are usually disjoint
  - Map job: each mapper first obtains all split data (small amount of data movement)
  - Reduce job: massive shuffling of data over network

- We allow splits to **overlap** by two rows

- DSGD is implemented as a **map-only** job (no data shuffling!)

| 1 | $a_{1,1}$ | $a_{1,2}$ | $b_1$ | $x_1$ |
| 2 | $a_{2,1}$ | $a_{2,3}$ | $b_2$ | $x_2$ |
| 3 | $a_{3,2}$ | $a_{3,4}$ | $b_3$ | $x_3$ |
| 4 | $a_{4,3}$ | $a_{4,5}$ | $b_4$ | $x_4$ |
| 5 | $a_{5,4}$ | $a_{5,6}$ | $b_5$ | $x_5$ |
| 6 | $a_{6,5}$ | $a_{6,7}$ | $b_6$ | $x_6$ |
| 7 | $a_{7,6}$ | $a_{7,8}$ | $b_7$ | $x_7$ |
| 8 | $a_{8,7}$ | $a_{8,9}$ | $b_8$ | $x_8$ |
| 9 | $a_{9,8}$ | $a_{9,10}$ | $b_9$ | $x_9$ |
| 10 | $a_{10,9}$ | $a_{10,11}$ | $b_{10}$ | $x_{10}$ |
| 11 | $a_{11,10}$ | $a_{11,12}$ | $b_{11}$ | $x_{11}$ |
| 12 | $a_{12,11}$ | $a_{12,13}$ | $b_{12}$ | $x_{12}$ |
| 13 | $a_{13,12}$ | $a_{13,14}$ | $b_{13}$ | $x_{13}$ |

split 1

split 2

stratum $s = 3$

($x_7$ affects mapper 1)

# Other Implementation Details

- **Initial guess**
  - Ignore off-diagonal elements
  - Works well due to "diagonal dominance"

- **Stratum sequence** as in [GHS11]
  - Meander in a stratum for a while, then jump to next stratum
  - Tension between thorough exploration of stratum and randomness
  - Visit all $k$ rows in stratum: at each "sub-epoch" select one of $k$! orders at random
  - Similar strategy for jumping between strata
  - Convergence Theorem still applies

- **Step-size sequence**
  - Constant during sub-epoch
  - "Bold driver" heuristic
  - Experiment with initial step size
    (in parallel on small subsequences)

# Optimizing the Re-Use Factor for Maximum Efficiency

**To define (asymptotic) efficiency, consider budget-constrained setting [Fox & Glynn 1990; Glynn & Whitt 1992]**

- Cost of producing n outputs from Model 2:

$$C_n = \sum_{j=1}^{m_n} \tau_{1;j} + \sum_{j=1}^{n} \tau_{2;j}$$

$\tau_{i;j} =$ (random) cost of producing $j^{th}$ observation of $Y_i$

- Under (large) fixed computational budget c

  - Number of Model 2 outputs produced:

$$N(c) = \max\{n \geq 0 : C_n \leq c\}$$

  - Estimator:

$$U(c) = \theta_{N(c)} = N(c)^{-1} \sum_{j=1}^{N(c)} Y_{2;j}$$

# Optimizing the Re-Use Factor II

**The key limit theorem as budget increases to infinity**

Suppose that $E[\tau_1 + \tau_2 + Y_2^2] < \infty$. Then $U(c)$ is asymptotically $N(\theta, g(\alpha)/c)$.

where $r_\alpha = \lfloor 1/\alpha \rfloor$ and

$$g(\alpha) = (\alpha E[\tau_1] + E[\tau_2])\left\{ \mathrm{Var}[Y_2] + \left(2r_\alpha - \alpha r_\alpha(r_\alpha + 1)\right)\mathrm{Cov}[Y_2, \tilde{Y}_2]\right\}$$
(cost per obs.)  x  (contributed variance per obs.)

$\mathrm{Cov}[Y_2, \tilde{Y}_2] =$ covariance of two Model 2 outputs that share a Model 1 input

- Thus, minimize $g(\alpha)$ [or maximize asymptotic efficiency $= 1/g(\alpha)$ ]

# Proof Outline

- Set $W_{n,j} = \sum_{i=1}^{n} Y_{2;i} I[\text{input for ith run of Model 2 is } Y_{1;j}]$

- Thus $\theta_n = \left(\dfrac{m_n}{n}\right) m_n^{-1} \sum_{j=1}^{m_n} W_{n;j} \approx \alpha \cdot m_n^{-1} \sum_{j=1}^{m_n} W_{n;j}$

- By Theorem 1 in [Glynn & Whitt 1992], it suffices to show that

  - $C_n / n \xrightarrow{\text{a.s.}} \alpha c_1 + c_2$  (straightforward to show)

  - $W_{n,1}, W_{n,2}, \ldots, W_{n,m_n}$ obeys a "Lindeberg-Feller" FCLT

- Can establish standard "Lindeberg condition" which suffices for FCLT (Billingsley 1999)

- Some additional fussy details due to the cycling through Model 1 outputs

$$W_{n,1}$$
$$W_{n,2}$$
$$W_{n,3}$$
$$W_{n,4}$$

$W_{n,j}$ and $W_{n,j'}$ are independent for $j \neq j'$

# Point and Interval Estimates

**Typical scenarios:**

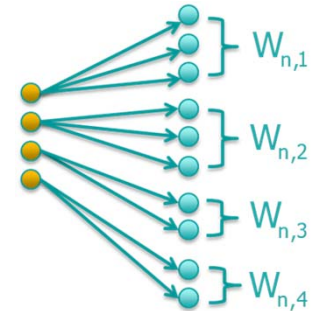- Compute $100(1-\delta)\%$ confidence interval for $\theta$ under fixed budget c
- Estimate $\theta$ to within $\pm 100\varepsilon\%$ with probability $100(1-\delta)\%$

$W_{n,1}$
$W_{n,2}$
$W_{n,3}$
$W_{n,4}$

**Issue: n is unknown a priori (so can't compute $m_n$)**

- Solution: estimate n from $n_0$ pilot (or prior) runs

$W_{n,j}^{(c)}$ is "centered" version of $W_{n,j}$

- Can show: $\dfrac{\sqrt{n}(\theta_n - \theta)}{\sqrt{h_n(\alpha)}} \Rightarrow N(0,1)$ where $h_n(\alpha) = n^{-1}\sum_{j=1}^{m_n}\left(W_{n,j}^{(c)}\right)^2$
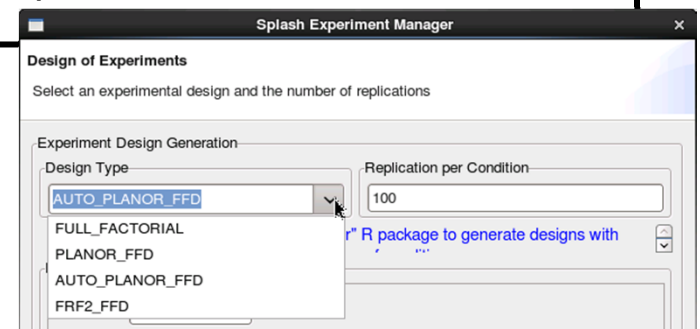
so that CI from n runs is $\left[\theta_n - z_\delta\left(h_n(\alpha)/n\right)^{1/2}, \theta_n + z_\delta\left(h_n(\alpha)/n\right)^{1/2}\right]$

where $z_\delta$ is $(1+\delta)/2$ normal quantile

- Can set
  - $n \approx c/(\alpha c_1 + c_2)$      for fixed budget
  - $n \approx h_{n_0}(\alpha)\left(z_\delta/\varepsilon\theta_{n_0}\right)^2$   for fixed precision
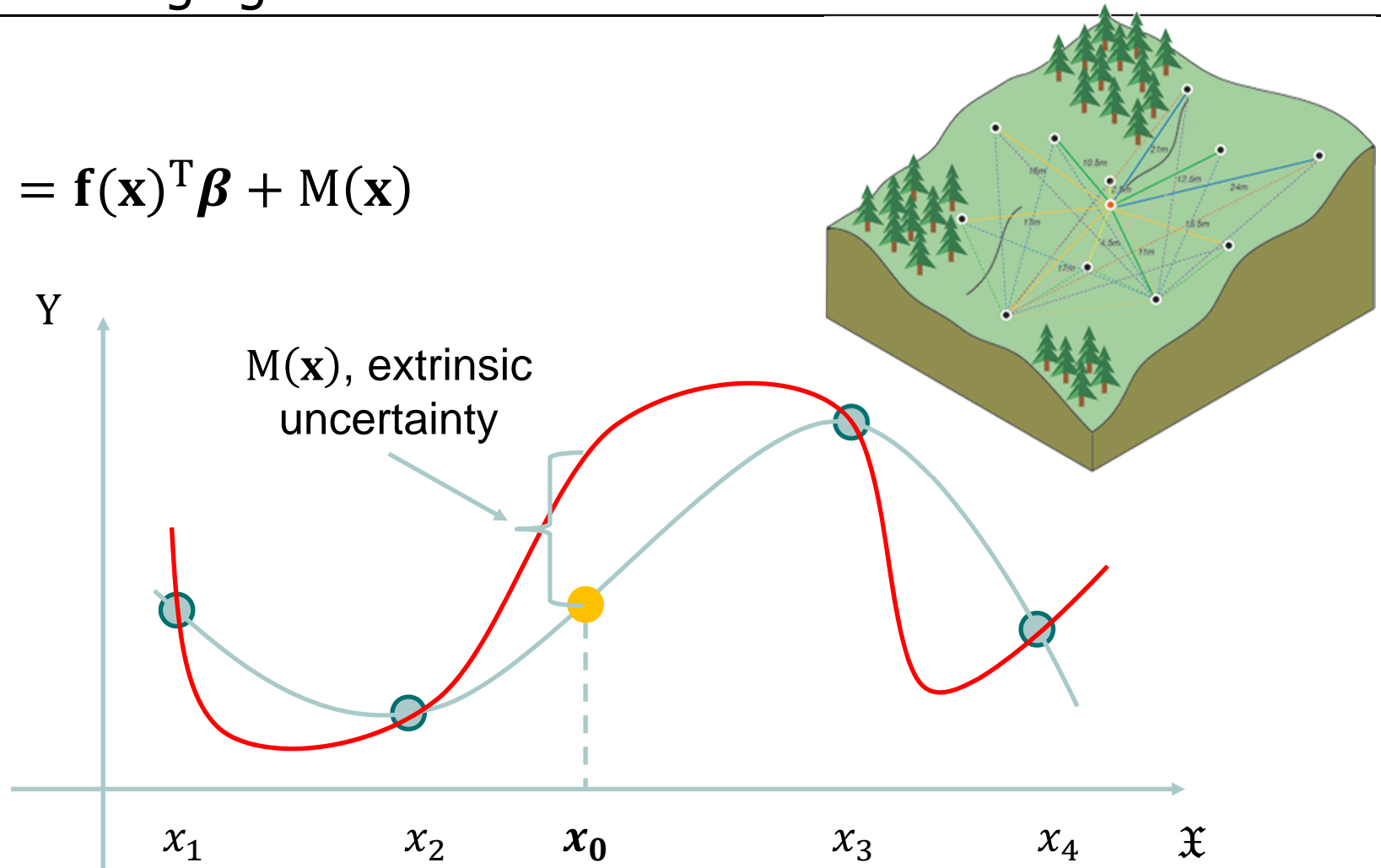
# Interface to R system for experimental design

| Method | Provider | Notes |
|---|---|---|
| Full Factorial Design | Experiment Manager | ▪Simple, fast design generation<br>▪Exhaustive factor combinations -> slow execution |
| Planor Fractional Factorial Design | R – planor package<br>http://cran.r-project.org/web/packages/planor/vignettes/PlanorInRmanual.pdf | ▪Supports arbitrary factor levels<br>▪Leverages R design generation<br>▪Checks statistical feasibility of user's proposed design<br>▪Slow design generation, fast experiment execution |
| Auto Planor Fractional Factorial Design | R – planor package<br>http://cran.r-project.org/web/packages/planor/vignettes/planorVignette.pdf | ▪Supports arbitrary factor levels<br>▪Leverages R design generation<br>▪Automatically finds smallest feasible experiment<br>▪Slower design generation, fast experiment execution |
| FRF2 Fractional Factorial Design | R – FrF2 package<br>http://cran.r-project.org/web/packages/FrF2/FrF2.pdf | ▪Only supports 2-level factors<br>▪Fast generation<br>▪Fast execution |
| Custom | User Specified | ▪Any design above may be used as basis |

As new designs are introduced in R, the

interface is in place to take advantage of these.

**Splash Experiment Manager**

**Design of Experiments**

Select an experimental design and the number of replications

Experiment Design Generation

Design Type

AUTO_PLANOR_FFD

Replication per Condition

100

" R package to generate designs with

FULL_FACTORIAL
PLANOR_FFD
AUTO_PLANOR_FFD
FRF2_FFD

# Standard Kriging

$$\mathcal{Y}(\mathbf{x}) = \mathbf{f}(\mathbf{x})^{\mathrm{T}}\boldsymbol{\beta} + \mathrm{M}(\mathbf{x})$$



Y

$\mathrm{M}(\mathbf{x})$, extrinsic uncertainty

$x_1$    $x_2$    $\boldsymbol{x_0}$    $x_3$    $x_4$    $\mathfrak{x}$

Images: SJSU

# Stochastic Kriging

$$\mathcal{Y}_j(\mathbf{x}) = \mathbf{f}(\mathbf{x})^{\mathrm{T}}\boldsymbol{\beta} + \mathrm{M}(\mathbf{x}) + \boldsymbol{\varepsilon}_j(\mathbf{x})$$



$\boldsymbol{\varepsilon}_j(\mathbf{x})$, intrinsic uncertainty

Y

$x_1 \qquad x_2 \qquad \boldsymbol{x_0} \qquad x_3 \qquad x_4$

MLE estimate:

$$\widehat{\widehat{\mathrm{Y}}}(\mathbf{x}_0) = \beta_0 + \Sigma_{\mathrm{M}}(\mathbf{x}_0, \cdot)^{\mathrm{T}}\left[\Sigma_{\mathrm{M}} + \hat{\Sigma}_{\varepsilon}\right]^{-1}(\bar{\mathcal{Y}} - \beta_0 \mathbf{1}_{\mathrm{k}})$$

Images: SJSU

# Optimization Process Flow



- Optimizer is R code,
- Orchestration via Python scripts

● = template-based data extraction