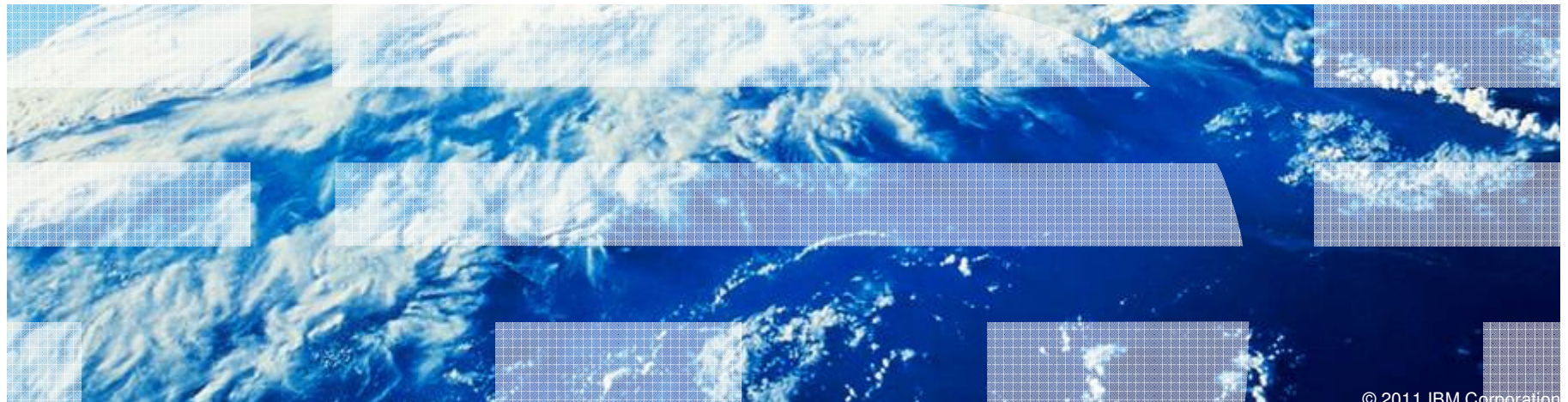




California Water Quality Monitoring Council – Data  
Management Workgroup:

## Data Integration Approaches and Technologies

*Peter Williams, CTO, “Big Green Innovations”, IBM  
Tuesday, November 22<sup>nd</sup> 2011*



# Begin with what we know: "My Water Quality" portal

The California Water Quality Web Portal - hosted by the Surface Water Ambient Monitoring Program - Mozilla Firefox: IBM Edition

File Edit View History Bookmarks Tools Help

http://www.waterboards.ca.gov/mywaterquality/

California water quality monitoring council

CA.GOV State of California ENVIRONMENTAL PROTECTION AGENCY RESOURCES AGENCY CALIFORNIA WATER QUALITY MONITORING COUNCIL

Skip to: [Content](#) | [Footer](#) | [Accessibility](#)

Search  GO

California  This Site

Home Safe to Drink Safe to Swim Safe to Eat Fish Ecosystem Health Stressors & Processes Contact Us

My Water Quality - hosted by the Surface Water Ambient Monitoring Program (SWAMP)

GOVERNOR SCHWARZENEGGER  
Visit his Website

- » Cal/EPA
- » The Resources Agency
- » About the California Water Quality Monitoring Council
- » State & Regional Water Boards
  - » Performance Report
- » Web Portal Partners
- » Monitoring & Assessment Programs, Data Sources & Reports
- » Water Quality Standards, Plans and Policies
- » Regulatory Activities
- » Enforcement Actions
- » Research
- » About SWAMP
- » SWAMP Tools

**SWAMP**  
Surface Water Ambient Monitoring Program

## Welcome to My Water Quality

This web portal, supported by a wide variety of public and private organizations, presents California water quality monitoring data and assessment information that may be viewed across space and time. Initial web portal development concentrates on four theme areas, with web portals to be released one at a time. Click the [Contact Us](#) tab for more information.

The Monitoring Council seeks to provide multiple perspectives on water quality information and to highlight existing data gaps and inconsistencies in data collection and interpretation, thereby identifying areas for needed improvement in order to better address the public's questions. Questions and comments should be addressed through the [Contact Us](#) tab.

### IS OUR WATER SAFE TO DRINK?

Safe drinking water depends on a variety of chemical and biological factors regulated by a number of local, state, and federal agencies. [More>>](#)

### IS IT SAFE TO SWIM IN OUR WATERS?

Swimming safety of our waters is linked to the levels of pathogens that have the potential to cause disease. [More >>](#)

### IS IT SAFE TO EAT FISH AND SHELLFISH FROM OUR WATERS?

Aquatic organisms are able to accumulate certain pollutants from the water in which they live, sometimes reaching levels that could harm consumers. [More>>](#)

### ARE OUR AQUATIC ECOSYSTEMS HEALTHY?

The health of fish and other aquatic organisms and communities depends on the chemical, physical, and biological quality of the waters in which they live. [More>>](#)

### WHAT STRESSORS AND PROCESSES AFFECT OUR WATER QUALITY?

Beneficial uses of our waters are affected by emerging contaminants, invasive species, trash, global warming, acidification, pollutant loads, and flow. [More>>](#)

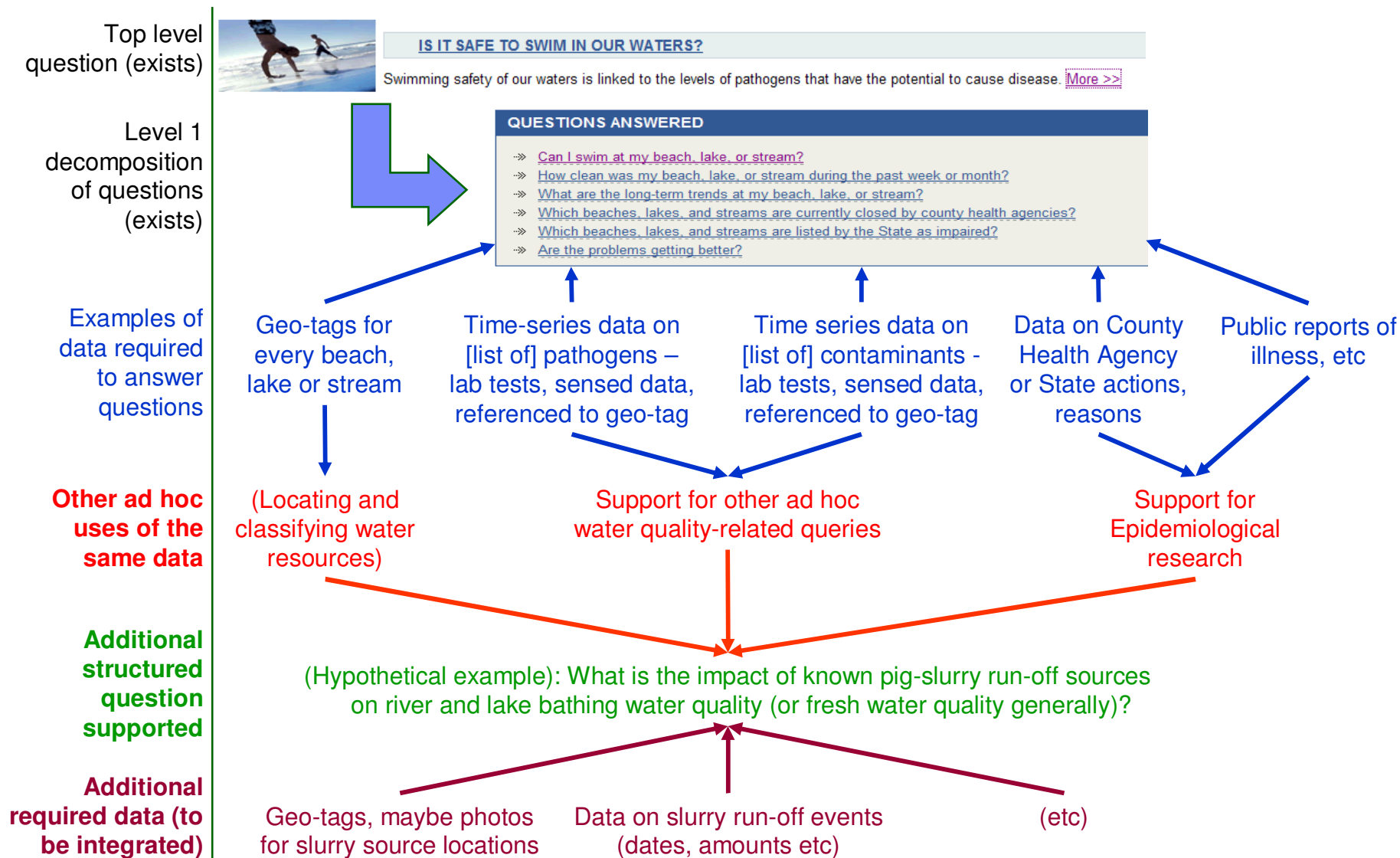
(Updated 3/16/10)

Back to Top Help Contact Us Site Map

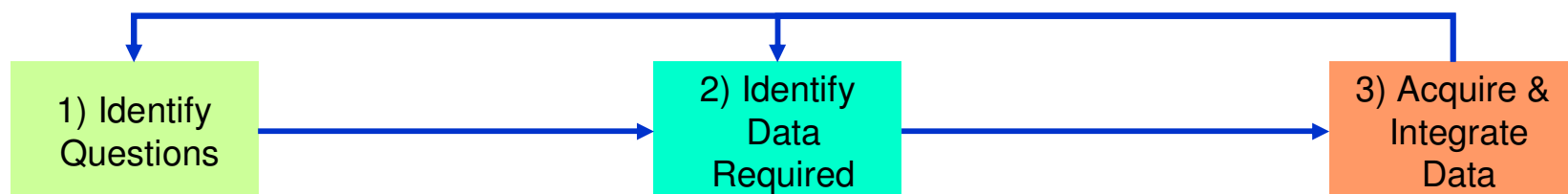
## Why “My Water Quality” gets it right...

- |  |  |
|--|--|
| <ul style="list-style-type: none"><li>■ Uses <b>questions</b> as the expression of business need for data:<ul style="list-style-type: none"><li>- Captures that need from the user’s perspective</li><li>- Can be broken down into subsidiary questions</li></ul></li></ul>  | <ul style="list-style-type: none"><li>■ Uses <b>questions</b> as “integration points” for multiple data sources:<ul style="list-style-type: none"><li>- Drives data architecture</li><li>- Forces de facto standards around semantics and scale</li><li>- (Can incorporate other standards)</li><li>- Data aggregated for each question should then support additional related, but less structured queries too</li><li>- Provides a focus to collaborate around</li><li>- Allows value to be derived faster</li></ul></li></ul> |
| <ul style="list-style-type: none"><li>■ Uses questions to identify data that matters and prioritize data sets to work with:<ul style="list-style-type: none"><li>- Can be made to use results of previous questions - builds over time to a wider level of integration</li><li>- If you can’t frame the question, you don’t know the need (or the value)</li></ul></li></ul> |  |
| <ul style="list-style-type: none"><li>■ Avoids GADWITS - the “Great Amorphous Data Warehouse In The Sky”!!</li><li>■ Enables “purpose-driven” data federation</li></ul>  |  |

# Example: “Is it safe to swim in our waters?”



## Structured questions as the core of a data integration process?



- Identify key questions to be answered
  - Identify decomposition questions
- Test with users
  - Identify value of answers to questions (who benefits; how; \$\$\$ if available)
- Prioritize key questions and decompositions
- Add questions to catalog

- Identify data required to answer key and decomposition questions
- Identify consuming applications, models, portals & needs arising
- Identify available data-sets & gaps
  - Search datasets already integrated
  - Identify proxy data if needed
- Confirm scale and semantic match with need
- Identify other possible ad-hoc uses and other questions from catalog that could be answered with this dataset
- Confirm cost and feasibility of acquiring and using each data-set
- Confirm question priorities in light of technical feasibility

- Acquire data
- Create adaptors etc necessary for semantic matches
- “Fill in” or “widen out” to address scale mismatches
- Test with consuming applications, models, portals –
- Confirm that question is “answered”/user need is met
- Create pub/sub or other interfaces to enable continuous acquisition
- Populate data warehouse/ data-set catalog as applicable (depends on degree of data federation)
- Repeat the above for additional uses

## Core principles

- “Question-driven”: data only integrated as needed to answer an unambiguously articulated questions – never “just because”
  - Can be adapted for historical perspective: for example, “how has water cleanliness in lakes and rivers changed since 1950?”
  - Can be adapted for causation: for example, “what has been the correlation between levels of contaminant x and high rainfall events?”
- Data warehouse, if used, is populated as data and datasets are assembled to answer questions – not before
- Each dataset is catalogued by ref to content, format, scale etc and also questions it can be used to answer
- Value of any given data set and integration activity =  
$$\frac{\text{(Formal questions + ad hoc uses supported)}}{\text{(Technical risk x cost of integration)}}$$
- Assembling data to answer each question builds on data sets integrated – cumulative coverage of the field, driven by value of integration activity

- So what technological trends and developments do we need to take into account as we frame integration goals?

## Data – in the past, often the elephant in the room...?

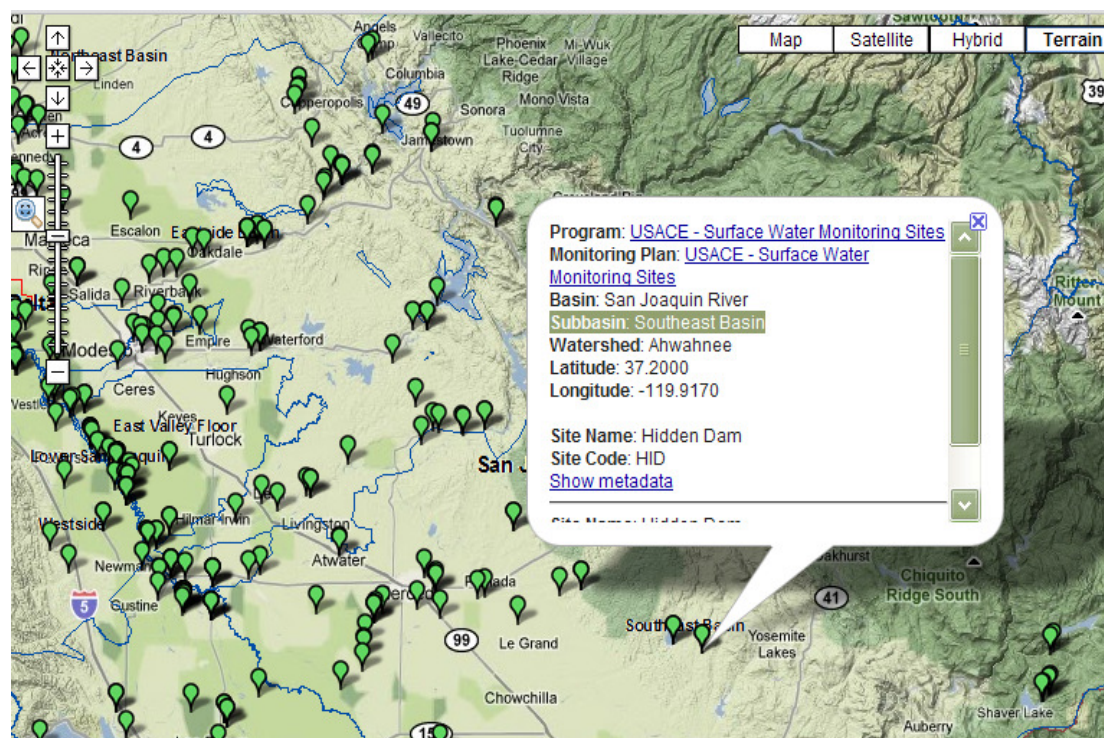
1. No data (rarer than one might think!)
2. Data is in the wrong scale (spatial or temporal) for the decision – too slow, late or infrequent, too scattered.
3. Data is fragmented between different stakeholders: different formats, scales, frequencies, standards; re-capture many times.
4. Too much data to use or analyze.
5. Incompatible or incomplete models mean that data is not leveraged, information is incomplete, or that solutions are partial or ineffective.
6. Poor visualization of information impedes effective decision-making: “so what’s this telling us?” syndrome.
7. Lack of awareness that 1-6 may be a problem. Possibly the most insidious problem of all!





## But data is becoming more available – even democratic

- Data is getting cheaper (Moore's law, more interest => more data and data sources)
- More and more data in the public domain – the open data movement.
- More and more ability to “mash it up” to create applications with it – often in the public domain.
- Socrata offers one of the most advance sets of open data/publishing tools around – starting to attract major attention.
  - San Francisco, Washington, Oregon (see over), Oklahoma, Chicago and others
- Especially useful in publishing data for consumption by other apps – eg smart phone apps



<http://www.centralvalleymonitoring.org/>

# Socrata example: Oregon's open data portal

**OREGON.gov**

Home Help Geo Data Metrics Suggest A Dataset APIs Oregon.gov

**AISP One & Two Year Tyvek Tag Dealers**  
Based on AISP One and Two Year Tyvek Tag Dealers  
Where to get the Aquatic Invasive Species Prevention Permit. All boats 10 feet long and longer (canoes, kayaks, drift boats, inflatable boats, etc.), out-of-state motor boats, and sailboats are required to purchase

Visualize

Conditional Formatting

Map

Views with locations can be displayed as points on a map

Map Setup

- Map Type: Google Maps
- Plot Style: Point Map

Location

- Location: Location

Details

- Title: Name
- Flyout Details: Location
- Flyout Details: Phone
- Flyout Details: County
- Flyout Details: Type of Tyvek Permi...
- Flyout Details: Hours
- Flyout Details: Website

+ Add Flyout Details

- Example shows location of invasive species tag dealers.
- Note user configurability – panel on the right.

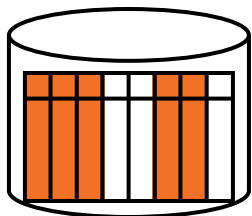
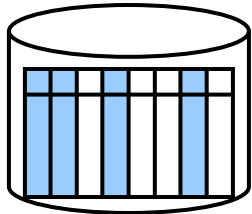
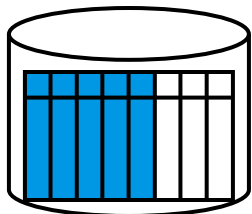
## Data integration – a choice of approaches

- Materialization: create integrated data set via extract, transform, load (ETL) from multiple data sources or replication
- Enterprise Application Integration (EAI): writing a special application linked to workflow
- Federation: virtual representation of integrated data set, materializing only what is needed, when needed
- Indexing: single index, enables data or docs to be fetched dynamically at user request

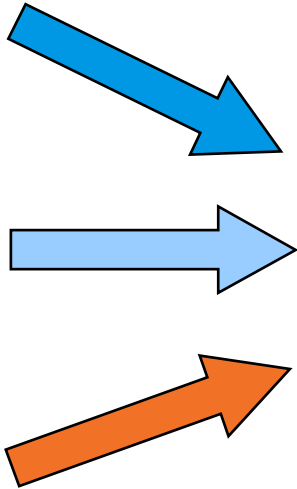
# Data warehousing is the oldest form of information integration

*Materialization, Data Exchange, ETL (Extract, Transform, Load), Integration*

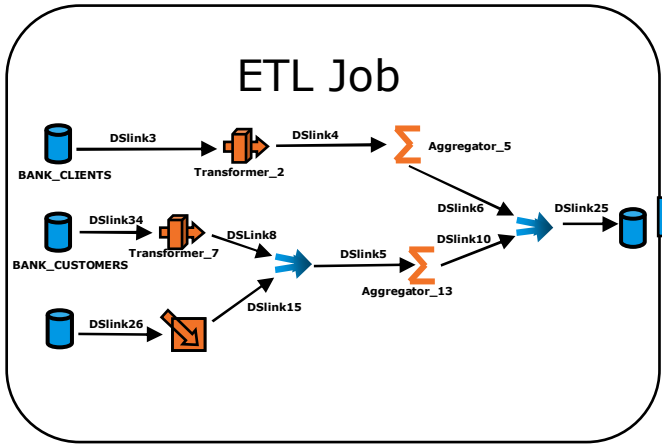
Data Sources  
(DBs)



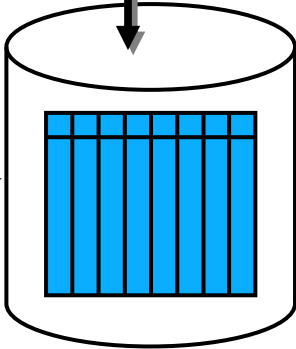
some data  
is relevant



*Integrated or  
Global Schema*



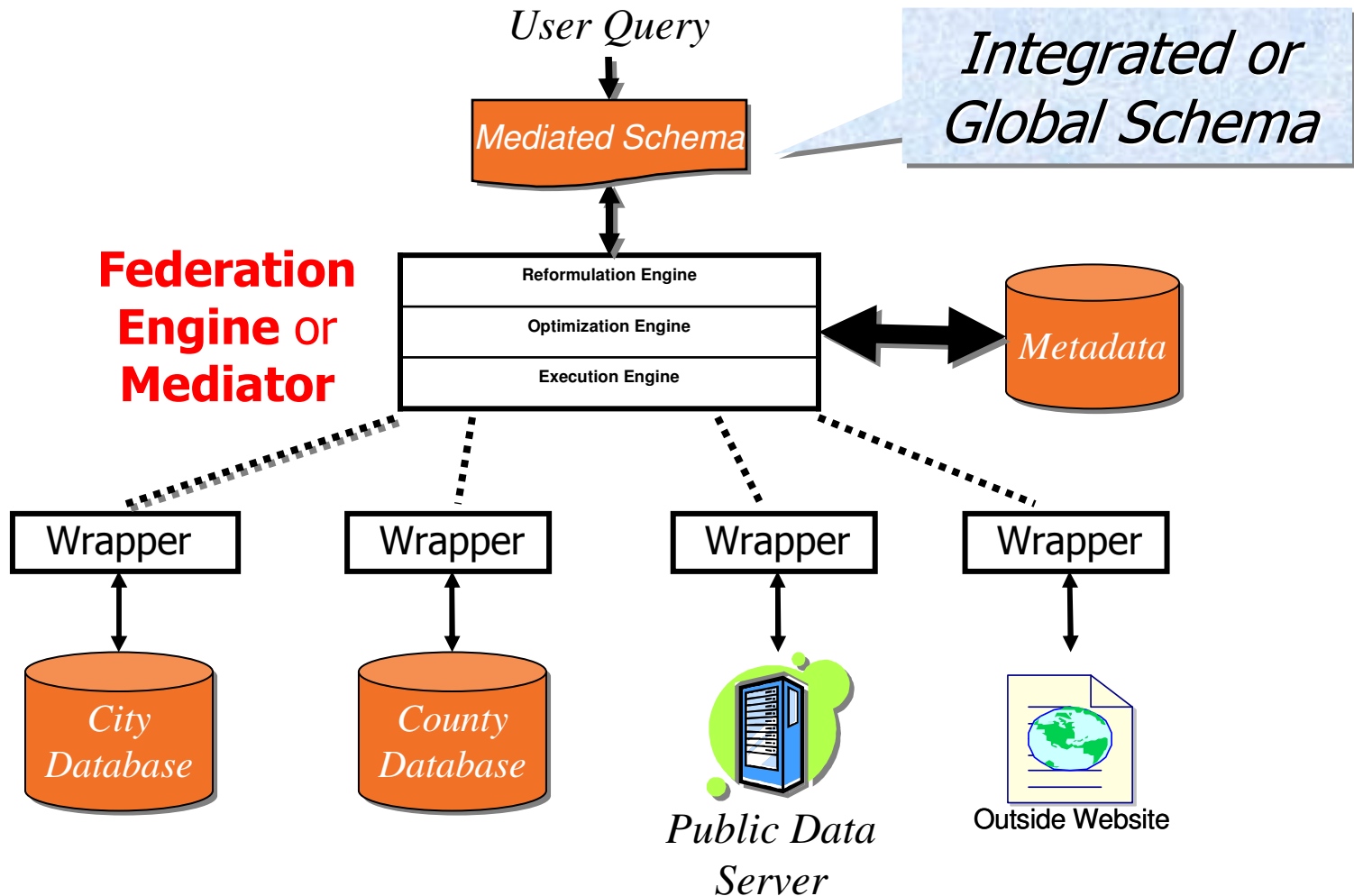
User Query  
Warehouse Schema



All relevant  
data collected

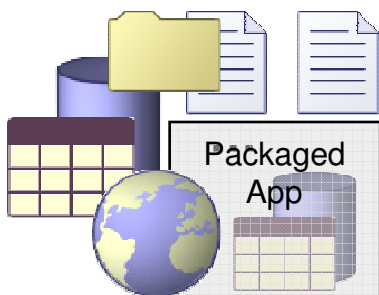
# Data federation and indexing is the other extreme

*Mediation, Virtual Integration, Data Integration, Lazy Integration*

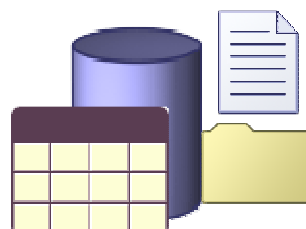


# Pros and cons of integration approaches

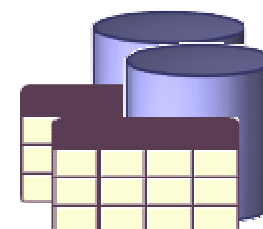
## Federation Services



...



...



Federation or index	Dimension	Consolidation (eg ETL)*
Latest, most current data – some network delay can be tolerated	<b>Key need</b>	Response time critical - rapid queries with no network delays
Small to medium result sets	<b>Data volumes</b>	Medium to very large result sets
Heterogenous - any	<b>Data types</b>	Largely relational
“Best endeavors”	<b>Availability</b>	High availability
Ad hoc	<b>Usage</b>	Predictable
Publishers’/originators’ systems	<b>Commonest failure mode</b>	Internal errors
Medium – tend to be faster	<b>Implementation complexity</b>	High – long complex projects
Stays with the originator	<b>Data control/politics</b>	Moves to the aggregator

## Complex integration targets - MIDAS

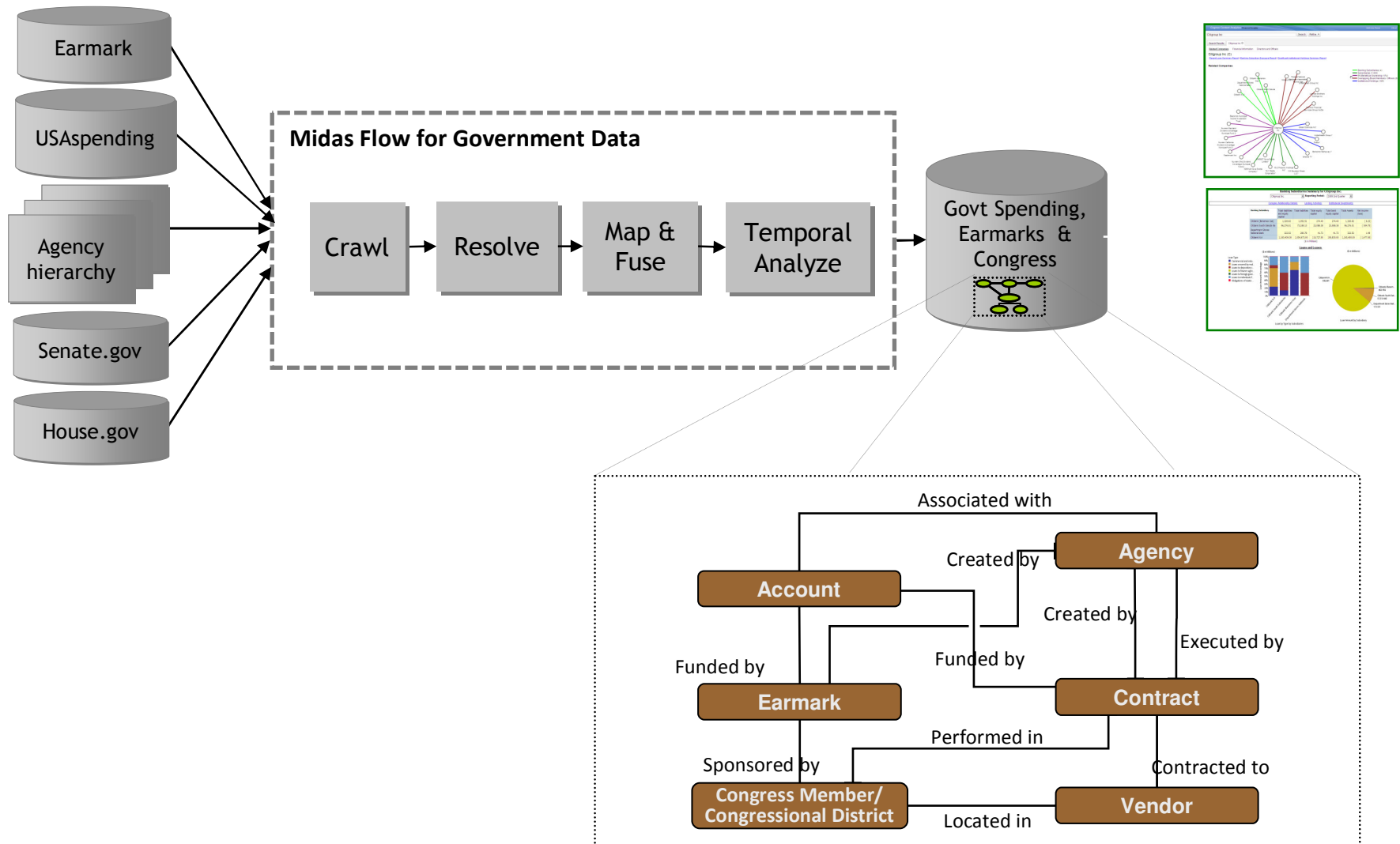
- MIDAS – IBM approach to enable queries over multiple, non-rationalized, data sources: integrates multiple partial or overlapping references into single complex entities.
  - May need to integrate multiple data types – text sensor data, numerical and so on
  - These data types may be both structured (tables) and unstructured (text, images, etc etc)
- Tested on work with financial services and government data, to answer queries like:
  - “How many and what value of earmarks in 2008 were solely sponsored by Republican or Democrat congress members?”
  - “Which public companies currently share one or more board member?”
- Significant open-sourced component – based on Apache Hadoop (data-intensive distributed application processing)

# Technical challenges of complex entity creation

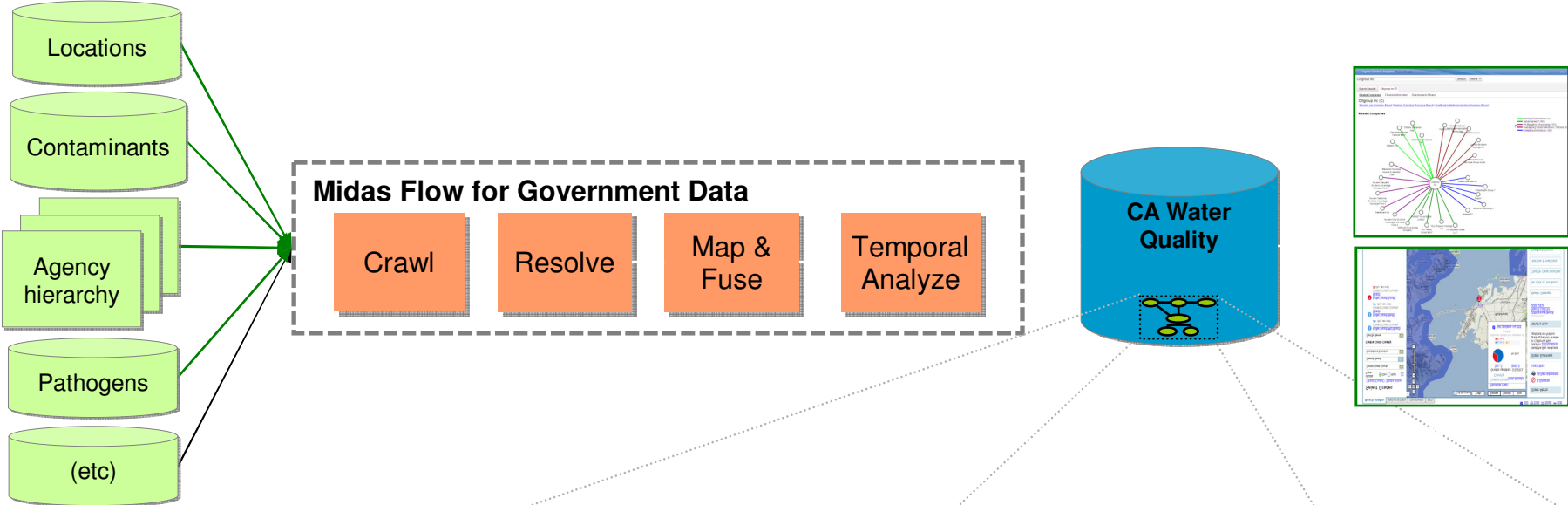
- Entity Extraction
  - Extracting structured data from unstructured texts about key entities, and relationships between key entities of the same type or different types
- Entity Resolution
  - Identifying that two instances (i.e., references) of the same entity type refer to the same real-world entity instance
  - Example: “Tahoe” vs “Lake Tahoe”, or whether “Tahoe” is a reference to the lake, the geological basin, or the region
- Data Mapping and Fusion
  - Map each extracted record (from its own format) to the target format
  - Fusion: aggregate multiple records of an entity into one (complex) object
- Temporal Analysis
  - Creation and management of timeline for relationships
- Scalable architecture
  - Complex analysis over millions of documents in a scalable manner



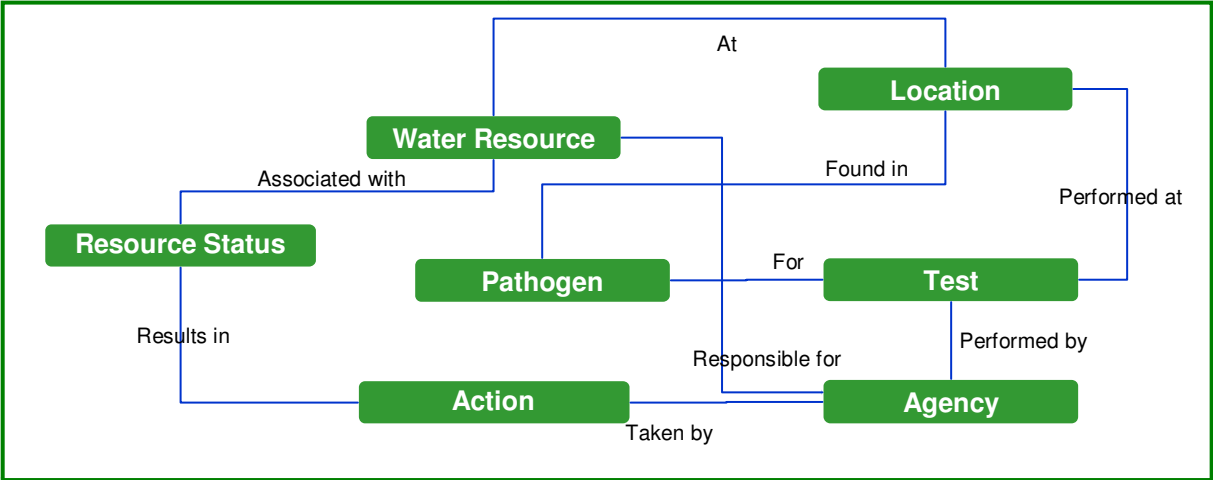
# Example: MIDAS flow for government spending data



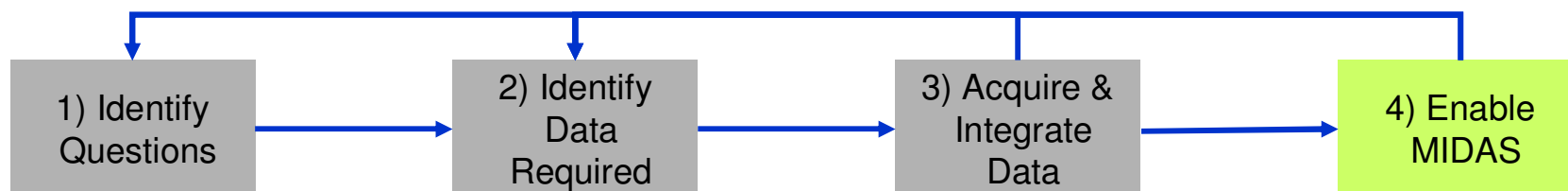
# MIDAS for water management...?



Q: "How clean was my beach, lake or stream during the past week or month"?



# Incorporation of MIDAS into the process...



- Identify key questions to be answered
  - Identify decomposition questions
- Test with users
  - Identify value of answers to questions (who benefits; how; \$\$\$ if available)
- Prioritize key questions and decompositions
- Add questions to catalog

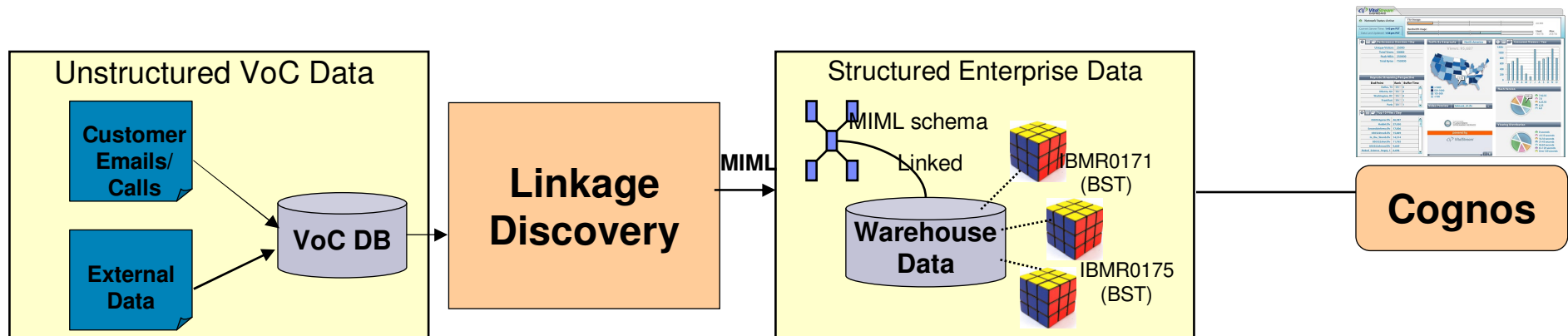
- Identify data required to answer key and decomposition questions
- Identify consuming applications, models, portals
- Identify available data-sets and gaps
  - Search datasets already integrated
  - If no direct data, identify proxies, if any
- Confirm scale and semantic match with need
- Identify other possible ad-hoc uses and other questions from catalog that could be answered with this dataset
- Confirm cost and feasibility of acquiring and using each data-set
- Confirm question priorities in light of technical feasibility

- Acquire data
- Create adaptors etc necessary for semantic matches
- “Fill in” or “widen out” as needed to address scale mismatches
- Test with consuming applications, models, portals –
- Confirm that question is “answered”/user need is met
- Create pub/sub or other interfaces to enable continuous acquisition (**see right**)
- Populate data warehouse/ data-set catalog as applicable (depends on degree of data federation)
- Repeat the above for additional uses

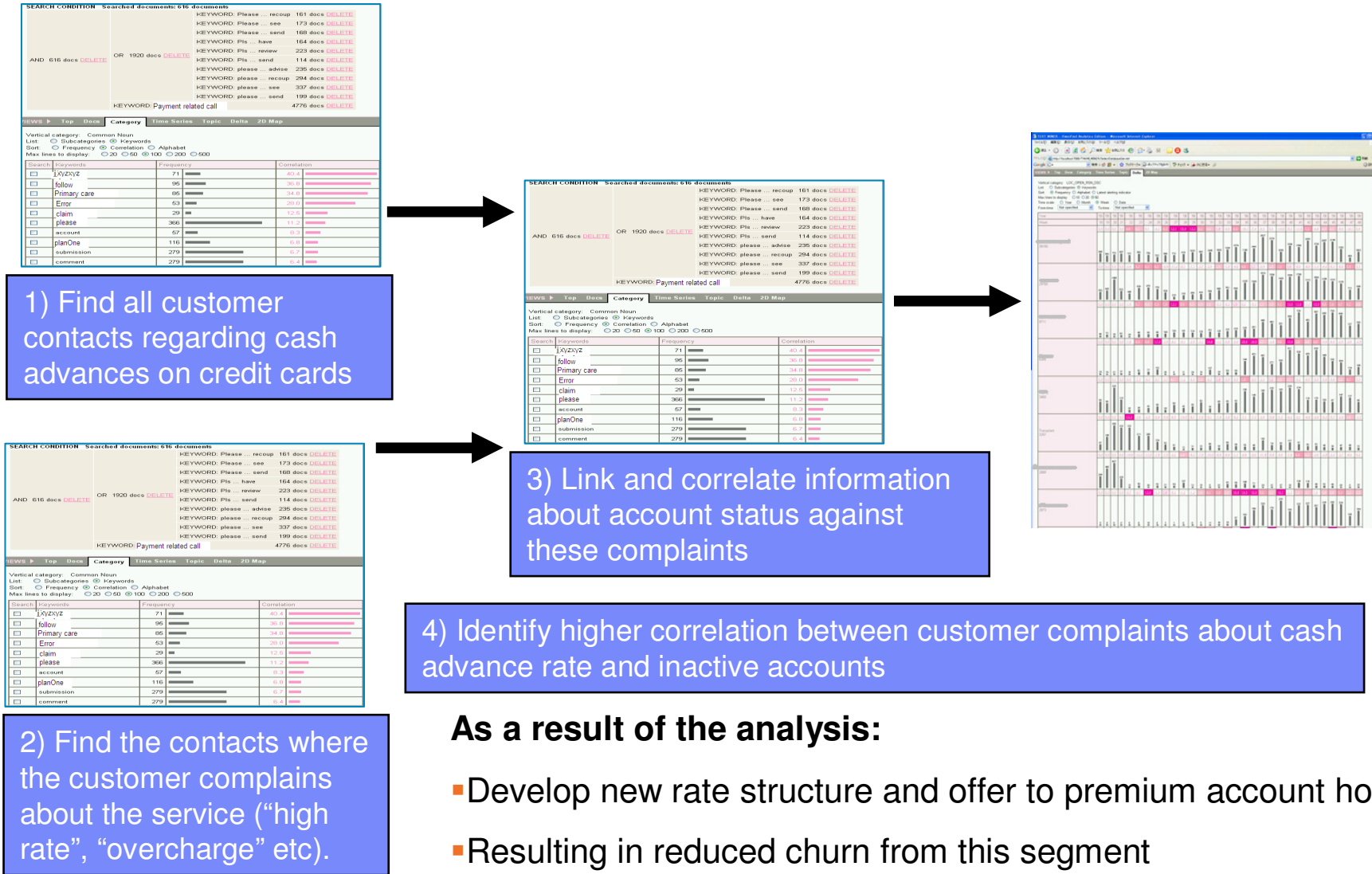
- Identify entities and sub entities
- Create nested entity relationships
- Aggregate records for each entity into complex objects and index these
  - By question supported
  - By other ad hoc use type supported
- Create search tool that enables queries
- If required, export to RDB to enable conventional BI tools to operate

## Enhanced reporting using unstructured data, in Cognos

<i>Example Reports generated from Structured Data</i>	<i>Example Reports generated from Unstructured Data</i>
<b>Location Profitability</b> - To identify the contribution to profit of geographic areas served by the financial institution.	<b>Analyze complaint categories and their distribution with region and location.</b>
<b>Channel Profitability</b> - To identify the contribution to profit of the Financial Institution's channel including branch networks, agencies, correspondents, and electronic channels	<b>Distribution of complaints across service channels, user sentiments vs. service channels</b>
<b>Customer Attrition Analysis</b> - To understand the reason and impact of customers ceasing to use the Financial Institution's products and services	<b>Reporting on key complaints, opinion features (specific feature of the product/ product that customers have opinion on )</b>
<b>Customer Complaint Analysis</b> - To understand the pattern of complaints and the effectiveness of the resolution process.	<b>How are complaints distributed across customer segments? Is there a pattern to the type of problems faced by customers in a specific segment</b>



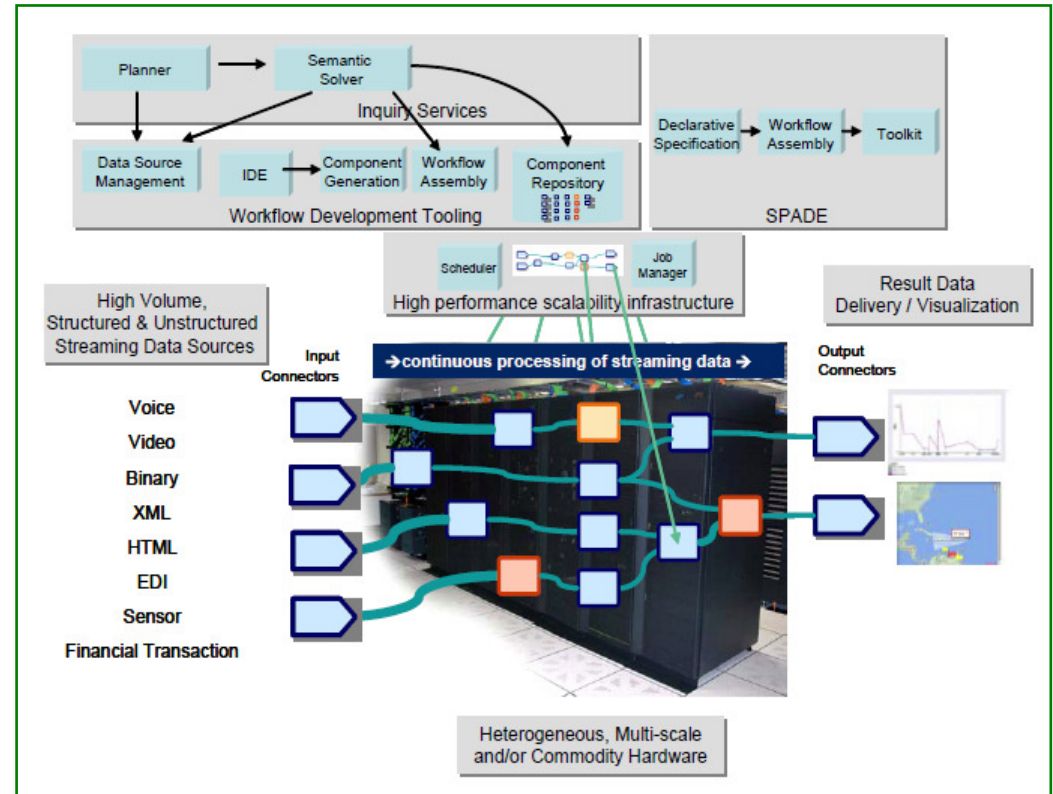
# Example analysis scenario in a typical financial services company. What's the equivalent for water?



- From integrating data to integrating models?

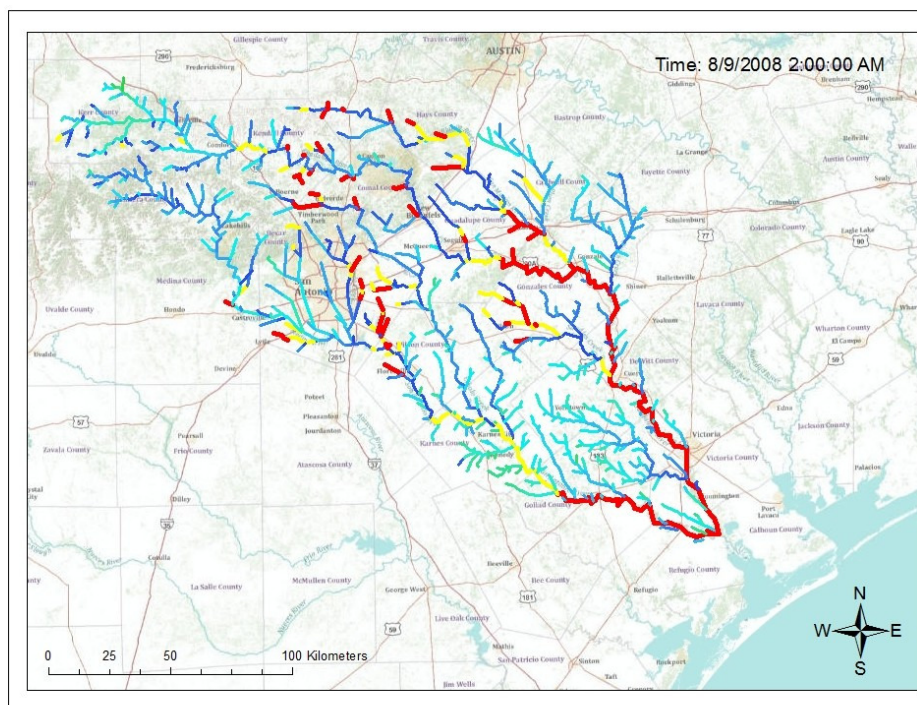
# The role of models that consume data is changing

- Advent of the model as a management tool
  - Used to be an “off-line”, back room scenario generation or exploration tool
  - Now more likely to be “on-line”, part of the decision stream
- What’s driving this?
  - Moore’s law: models that took hours to run can now run in minutes (combining computing power with model tuning).
  - Extreme case – can run models in real time as part of event processing system
  - Those models increasingly have sufficiently granular data from sensors and other systems to produce “operational-grade” conclusions
  - Commercial modeling and visualization software increasingly focused on supporting operational decision-making. Examples: IBM/ILOG, Optimatics, GL, Bentley and others.

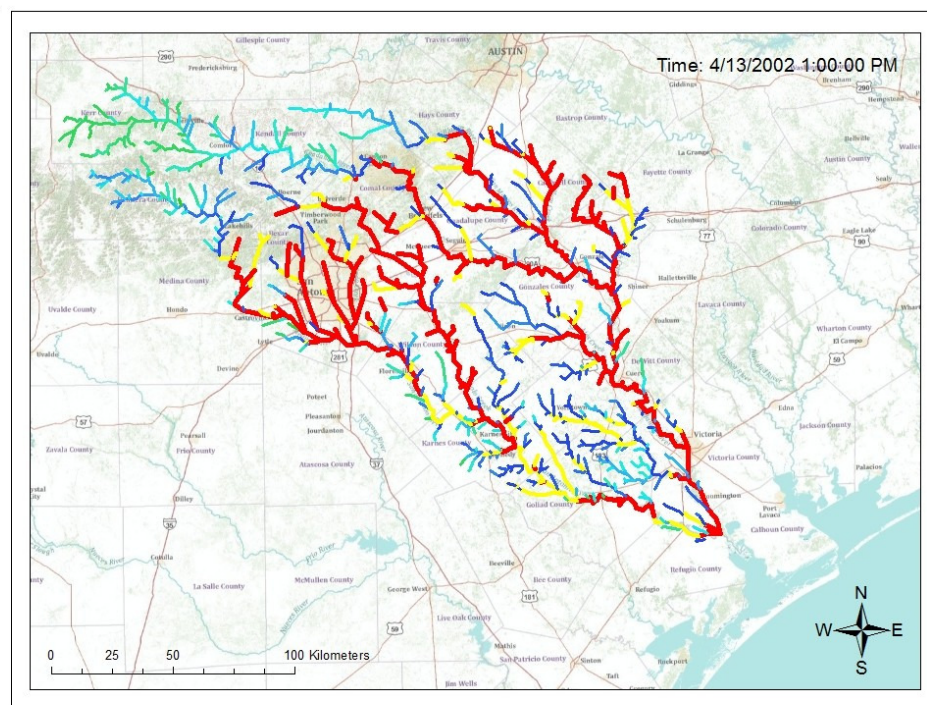


# Model tuning example

- Guadalupe-San Antonio river basin in central Texas
- Approximately 3,500 reaches, 15,000 km total length
- Modeled by 110K segments
- Fully dynamic computation for a seven day event simulated in one hour on a one-node x86 computer.
- Adapts “netlist” techniques from semiconductor industry to allocate processing power to segments of the model highly efficiently



Snapshot of event 1



Snapshot of event 2

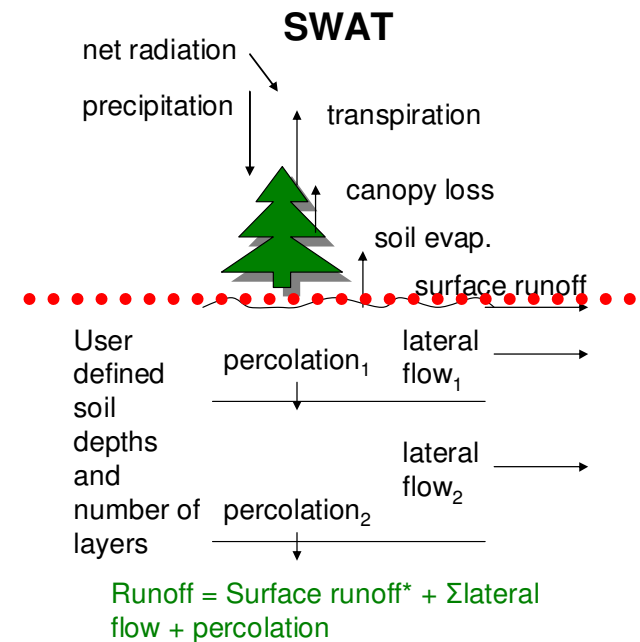
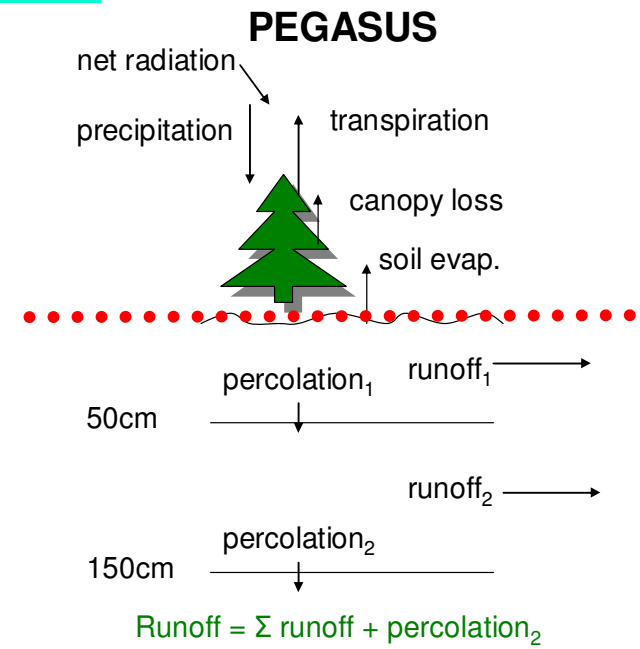
Mpeg files available



\*See final slide

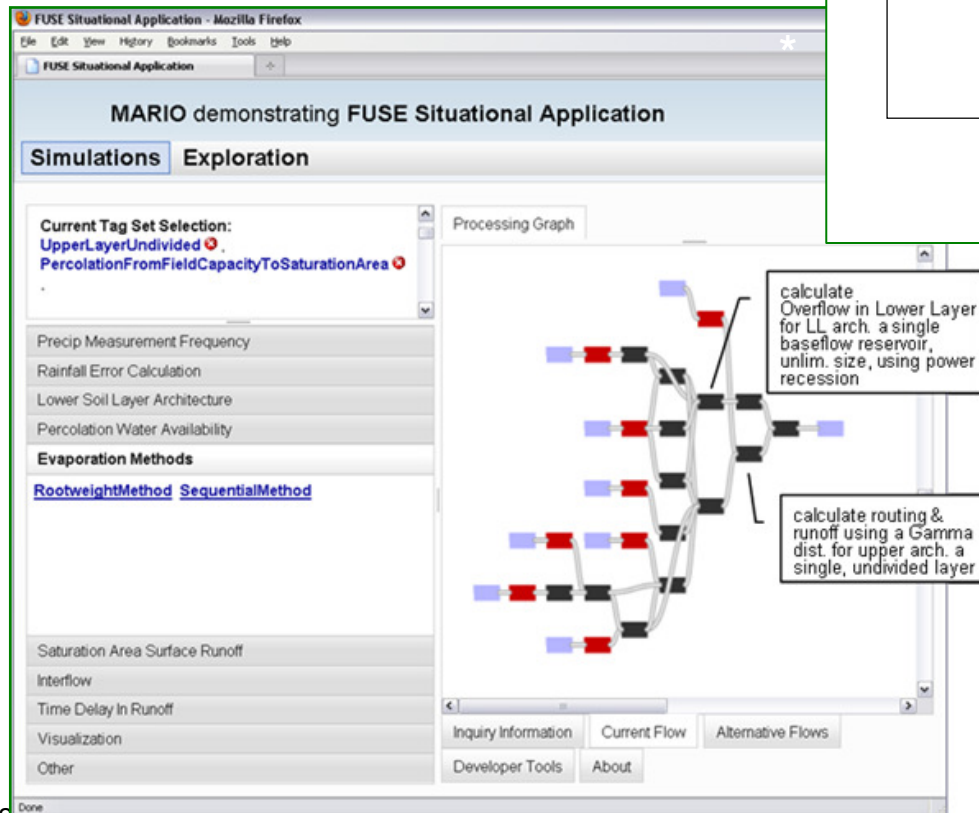
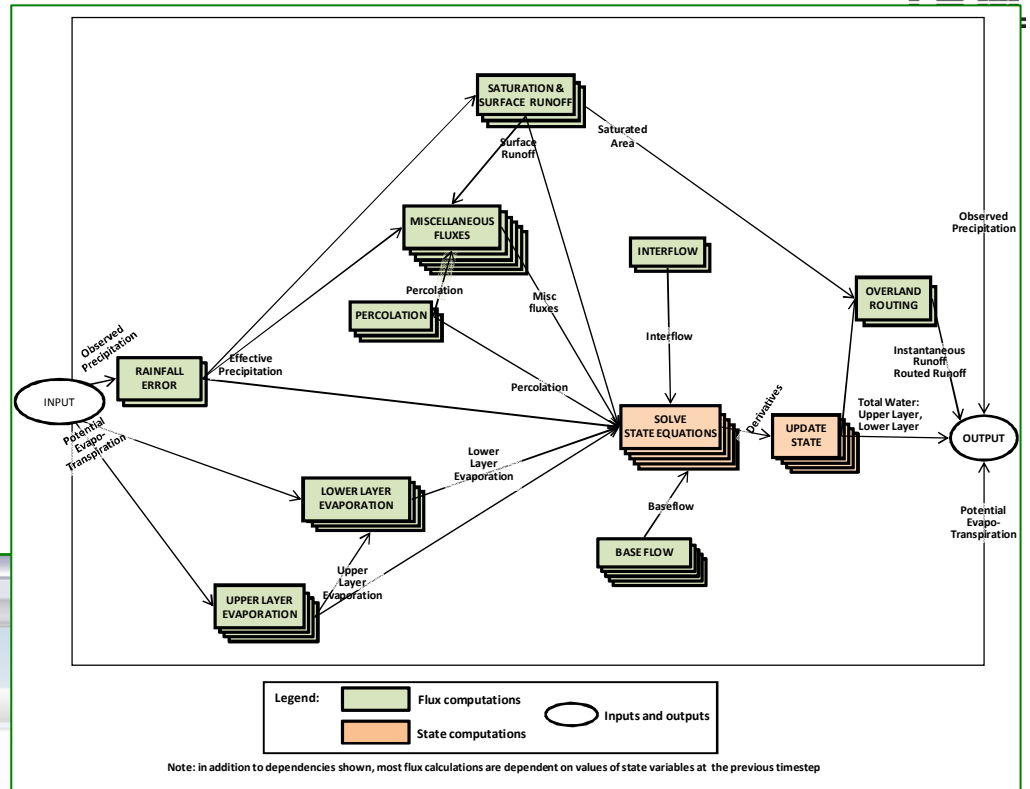
# Integrative modeling frameworks\*

- Model codes are typically monolithic: hard to integrate
- A more component based approach would allow components to be integrated into chains or compositions specific to the task at hand.
  - See right: the above-ground components of these two run-off models are identical, but the below-ground components are different
- Needed:
  - A way to break models into components, and compile and manage component sets and (larger) libraries
  - A way to allow only scientifically valid combinations of components - but include all such valid combinations
  - A consistent interface set and set of services to enable modules to interoperate
  - A concept of workflow to enable simulations using the chain of components
  - A rich semantic framework to express the above – captures context, not just checking inputs and outputs
  - Automatic co-calibration of model components
  - Ideally, ease of use for the non-scientist



# Mashup Automation with Runtime Invocation & Orchestration (MARIO)\* \*See final slide

- The models on the previous page were integrated manually.
- Now demonstrated as an automated process using FUSE\* framework for estimating stream discharge given observed precipitation:



- 250,000 theoretical module combinations, of which just 316 are scientifically valid (for example, tying flux equations to appropriate soil layer architectures)
- MARIO enables selection of processing goals via a set of tags, then selects and ranks possible chains of model components that satisfy those goals according to a set of rules

- Visualization – the afterthought

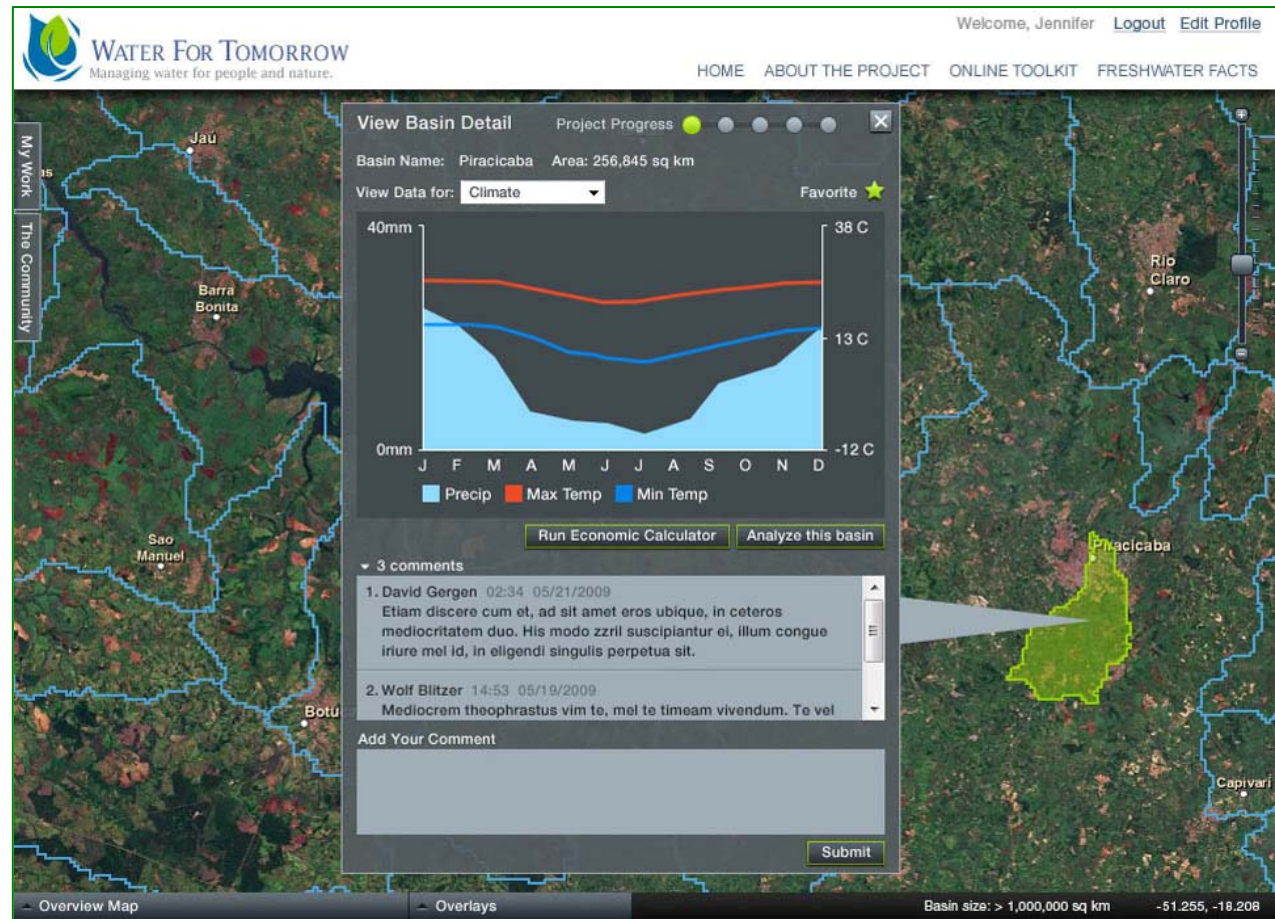
# Visualization

- **“A method of computing in which the enormous bandwidth and processing power of the human visual (eye-brain) system becomes an integral part of extracting knowledge from complex data”**,\* for example in:
  - Feature, trend or event identification \*See final slide
  - Comparison or fusion of data from multiple sources (visual fusion)
  - Decision support
  - Communication
- Visualization is too often an afterthought. It should be considered from the start when designing any application of environmental informatics – as distinct from “ok, here’s the data, now how do we display it?”

# Visualization as more than an afterthought: Water for Tomorrow\*

\*See final slide

- Helps planners and scientists analyze river basins and visualize the effects of different management scenarios on basin health
  - In so doing, initiates collaboration to develop sustainable water management policies
- Marries state-of-the-art technology with collaboration tools to support conservation of freshwater ecosystems
  - Combines rich graphics and dynamic mapping capabilities

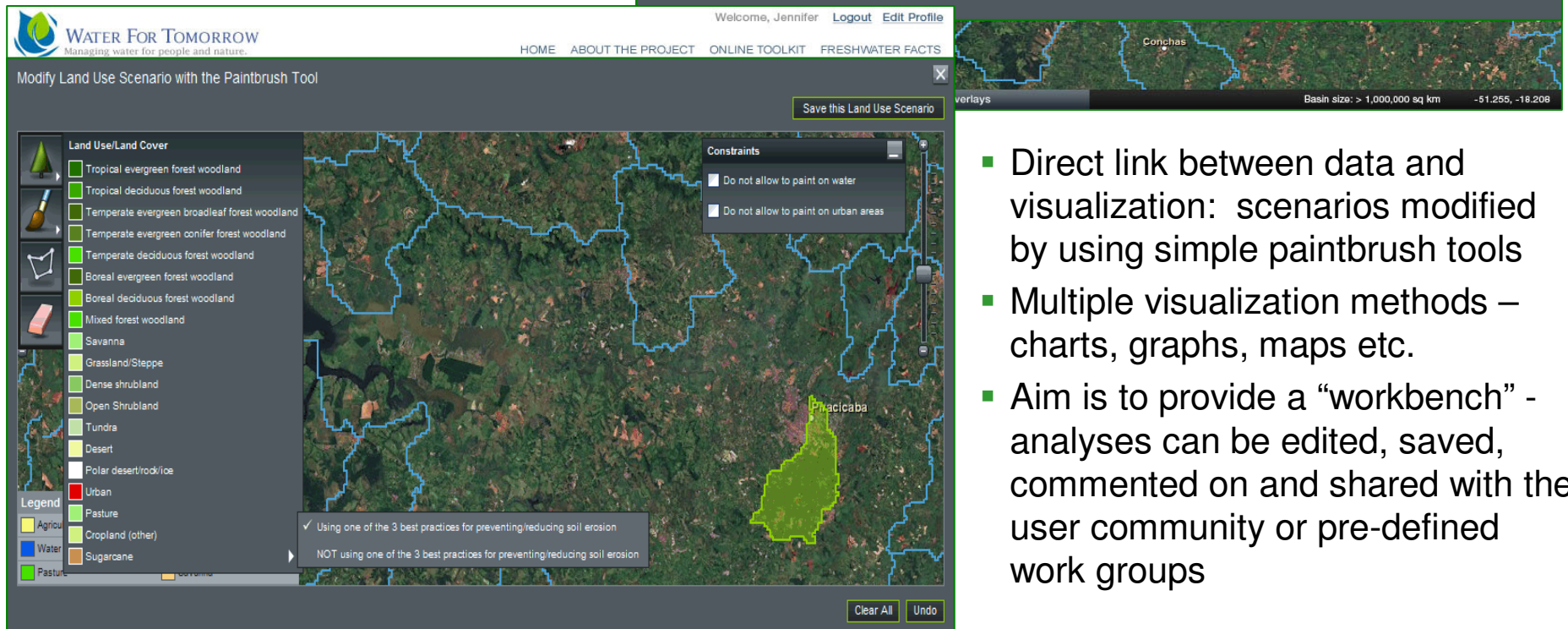
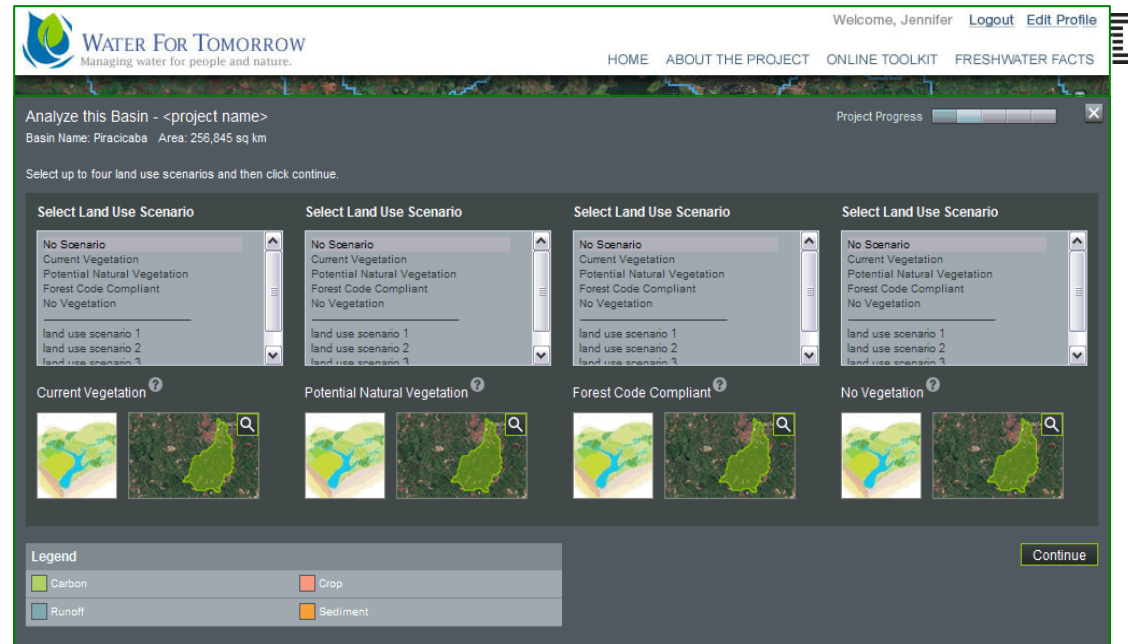


- Based on the Pegasus and SWAT hydrological modules reviewed earlier
- Basins selected by delineating areas on the map, or uploading a basin file
- This basin is part of the Piracicaba, Capivari and Jundiaí rivers (PCJ) pilot site in Brazil



# Scenario Selection and modification

- A land use scenario is a set of land use practices. Scenarios are used by the WFT hydrology model to compute crop production, water quality, water balance and other variables.
- Pre-canned scenarios are provided in the tool, or custom scenarios can be created or imported.



- Direct link between data and visualization: scenarios modified by using simple paintbrush tools
- Multiple visualization methods – charts, graphs, maps etc.
- Aim is to provide a “workbench” - analyses can be edited, saved, commented on and shared with the user community or pre-defined work groups

# Interactive Visualizations

- Contains an ROI calculator that enables the economic value/impact of land use scenarios, and cost benefit of current policies vs ecosystem protection vs remediation to be calculated and weighed into the discussion.

Water For Tomorrow  
Managing water for people and nature.

Welcome, Jennifer [Logout](#) [Edit Profile](#)

HOME ABOUT THE PROJECT ONLINE TOOLKIT FRESHWATER FACTS

Analyze this Basin - <project name>  
Basin Name: Piracicaba Area: 256,845 sq km

Project Progress

Land Use Scenarios  
[Edit Scenario Selections](#)

Current Vegetation  
Simulated Forest Code Compliant  
No Vegetation

<Chart Title> [Change Variables](#) [View Other Charts](#)

Carbon Description of this variable Units  
Runoff Description of this variable Units  
Crops Description of this variable Units

Basin Map  
[Explore Maps](#)

Legend  
Carbon Runoff Crop Sediment

Basin size: > 1,000,000 sq km -51,255, -18,208

Water For Tomorrow  
Managing water for people and nature.

Welcome, Jennifer [Logout](#) [Edit Profile](#)

HOME ABOUT THE PROJECT ONLINE TOOLKIT FRESHWATER FACTS

Economic Calculator  
Project Progress

Etiam discere cum et, ad sit amet eros ubique, in ceteros mediocritatem duo. His modo zrril suscipiantur ei, illum congue iriure mel id, in eligendi singulis perpetua sit. Mediocrem theophrastus vim te, mel te timeam vivendum. Te vel amet probo sanctus, ne movet tritani fastidii per, mel nulla minimum te. Partiendo consequuntur has at, idque labores eum ad, ne pro affert sapientem.

Basin Name: Piracicaba Area: 256,845 sq km What is your budget? \$1,640,000 [Recalculate](#)

Land Use	Units	Cost per Unit	Results	Notes
<input type="checkbox"/> Tropical evergreen forest woodland				
<input type="checkbox"/> Tropical deciduous forest woodland				
<input type="checkbox"/> Temperate evergreen broadleaf forest woodland				
<input type="checkbox"/> Temperate evergreen conifer forest woodland				
<input type="checkbox"/> Boreal evergreen forest woodland				
<input checked="" type="checkbox"/> Mixed forest woodland	365	\$1,525.00	Lorem	Lorem ipsum dolor sit amet consectetur adipiscing sed diam ...
<input type="checkbox"/> Savanna				
<input type="checkbox"/> Grassland / Steppe				
<input type="checkbox"/> Dense Shrubland				
<input type="checkbox"/> Open Shrubland				
<input type="checkbox"/> Tundra				
<input type="checkbox"/> Desert				
<input type="checkbox"/> Polar desert / rock / ice				
<input type="checkbox"/> Urban				
<input type="checkbox"/> Pasture				
<input type="checkbox"/> Cropland (other)				
<input type="checkbox"/> Sugarcane				

[Save & Share These Results](#) [Select Land Use Scenario](#)

- Visualization as the integration of data, models and decision processes
- Visualization drives the definition of the system, function offered and model integration – by implication, also used to structure underlying data

# Scorecards and dashboards

- Commonly used metaphors in business reporting, and for infrastructures. Excellent for depicting environmental performance of a business.
- But, they are under-investigated uses for environmental informatics:
  - They rely on key performance indicators: these could be defined for ecosystems, say, as easily as for a business.
  - Factors affecting performance against those indicators can be derived by systems dynamics modeling, or other methods.
- Can include considerable logic, up to and including models
  - Can operate in real time, with streaming data if needed – for process control, say
- Dashboards and scorecards are a beguilingly easy concept to grasp – in fact, unless all data is coming from an ERP system, they can require VERY significant integration work with systems generating core data
  - Direct linkages with sensors would in many cases be easier!!





# References

- Slide 25:
  - “Toward an integrative software infrastructure for water management in the Smarter Planet” B Eckman, M Feblowitz, A Mayer, and A Riabov, IBM Systems Journal 2009.
  - Deconstruction of Pegasus and SWAT produced by Eckman and others in course of IBM/Nature Conservancy “Water For Tomorrow” project (see below)
  
- Slide 26:
  - Eckman, et al. The authors demonstrated the application of MARIO to automating the chaining of modules from the hydrological models incorporated in the Framework for Understanding Structural Errors (FUSE) (Clark et al, referenced in Eckman et al). In the graphic, stacks of green boxes are options for flux calculations; stacks of brown boxes are options for state equations.
  - Other model integration frameworks are starting to appear in hydrology and water management:
    - OpenMI: standard software component interface that allows water models to exchange data on a time-step-controlled basis – allows modeling of process interactions, irrespective of temporal and spatial resolutions of component models. See <http://www.openmi.org/reloaded/about/what-is-openmi.php>
    - The Consortium of Universities for the Advancement of Hydrologic Science, Inc. (CUAHSI) has launched the a Community Hydrologic Modeling Platform (CHyMP). See <http://www.cuahsi.org/chymp-20090331.html>
  
- Slide 28:
  - “Visual Data Fusion for Applications of High Resolution Numerical Weather Prediction”, L Treinish, IBM, date unknown. Available from: <http://www.research.ibm.com/weather/vis/df.htm>
  
- Slide 29:
  - Water for Tomorrow: <http://www.nature.org/ourinitiatives/habitats/riverslakes/howwework/The-Partnership.xml> and [http://www.ibm.com/ibm/ideasfromibm/us/environment/100807/images/IFI\\_10082007.pdf](http://www.ibm.com/ibm/ideasfromibm/us/environment/100807/images/IFI_10082007.pdf)