# Mind Your Metadata

Exploiting Semantics for Configuration, Adaptation, and Provenance in Scientific Workflows

Yolanda Gil
Pedro Szekely
Craig Knoblock
Varun Ratnakar
Shubham Gupta
Maria Muslea
Fabio Silva

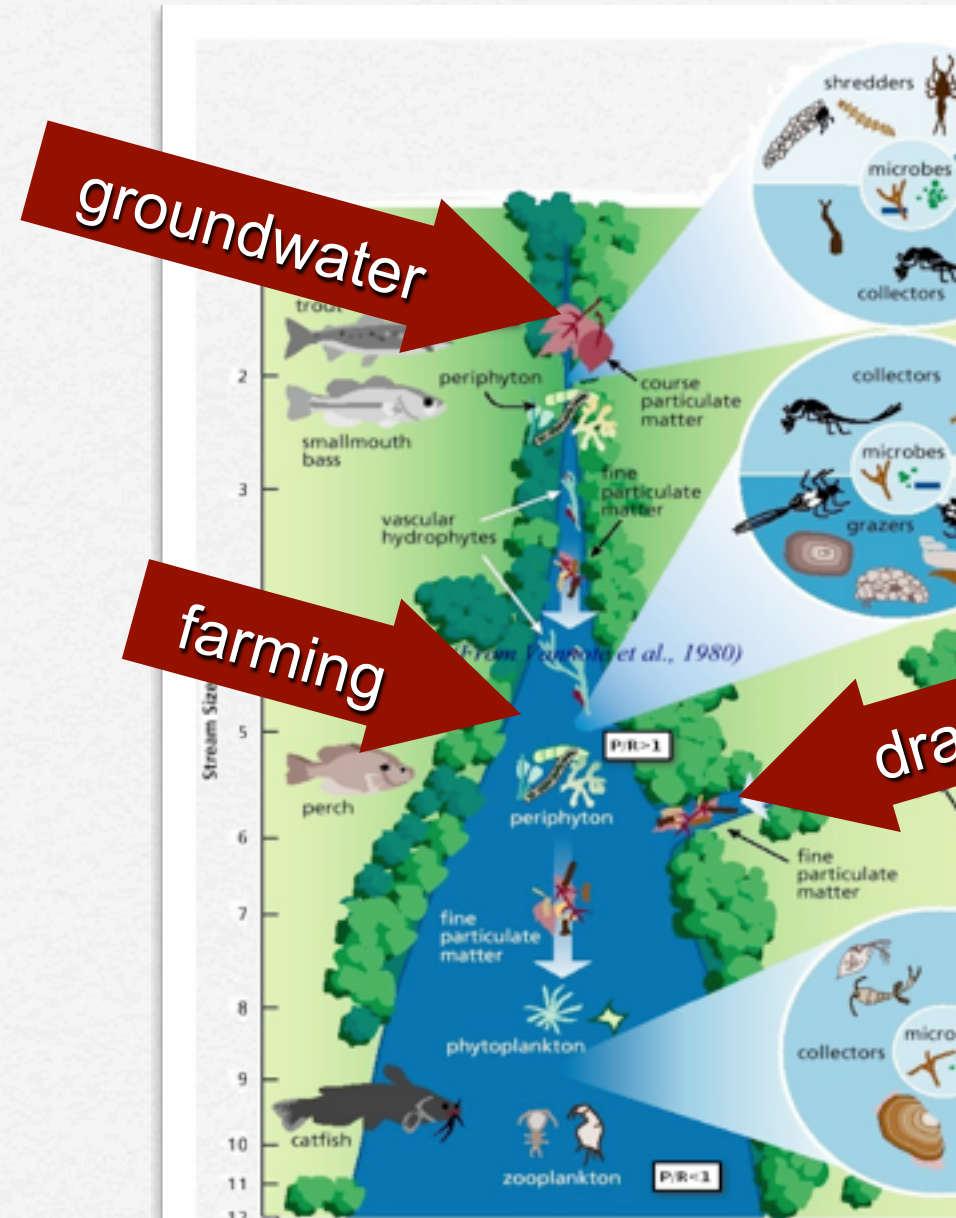Tom Harmon
Sandra Villamizar

UC Merced

# River Continuum vs Human Activities

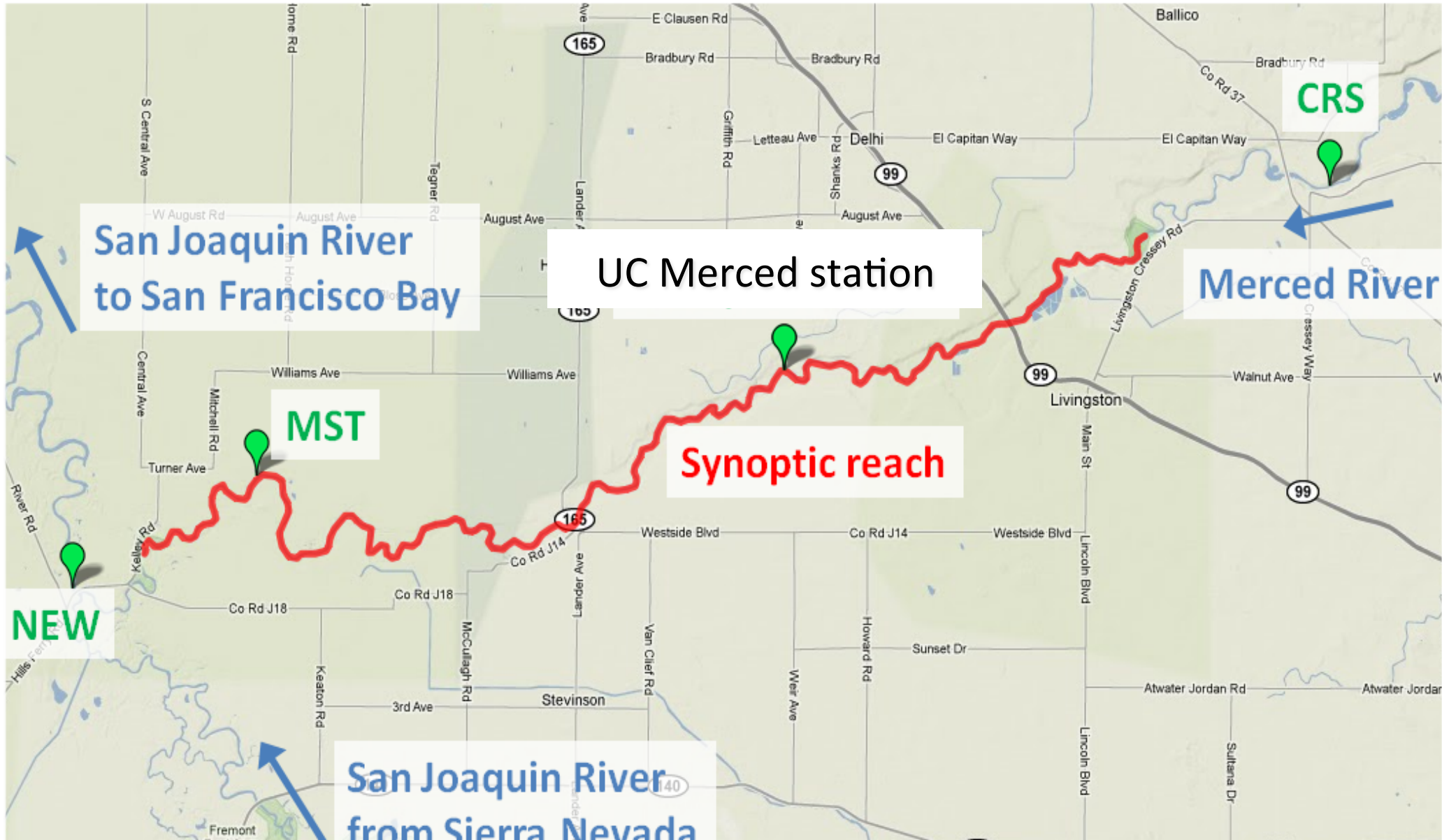River continuum: natural inputs, reactive transport

Human intervention: Agricultural, industrial, municipal

What management practices help/hurt?

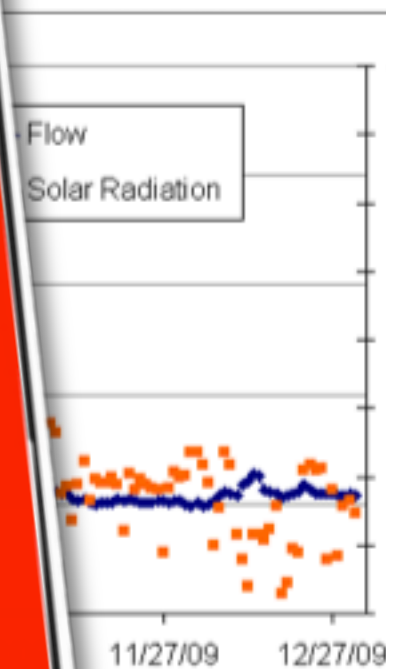Can we restore natural behavior?

# Case Study

# Stream Metabolism Response to Human Disturbances



... but how does this affect the ecology of the river?

... how about the effect of farmers?

se releases in the spring and fall to help the salmon r
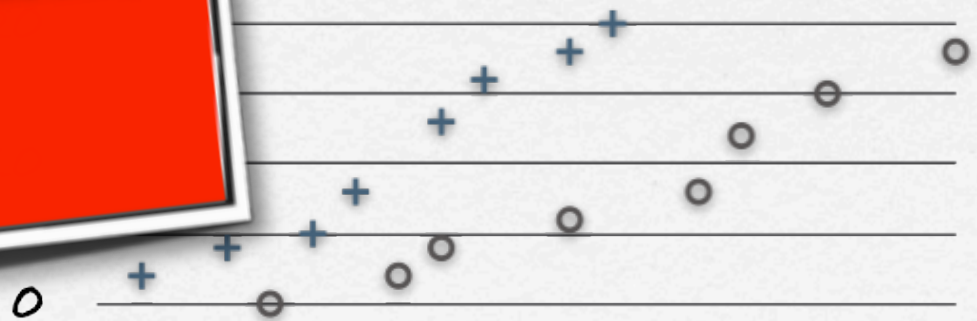
# Aquatic Photosynthesis

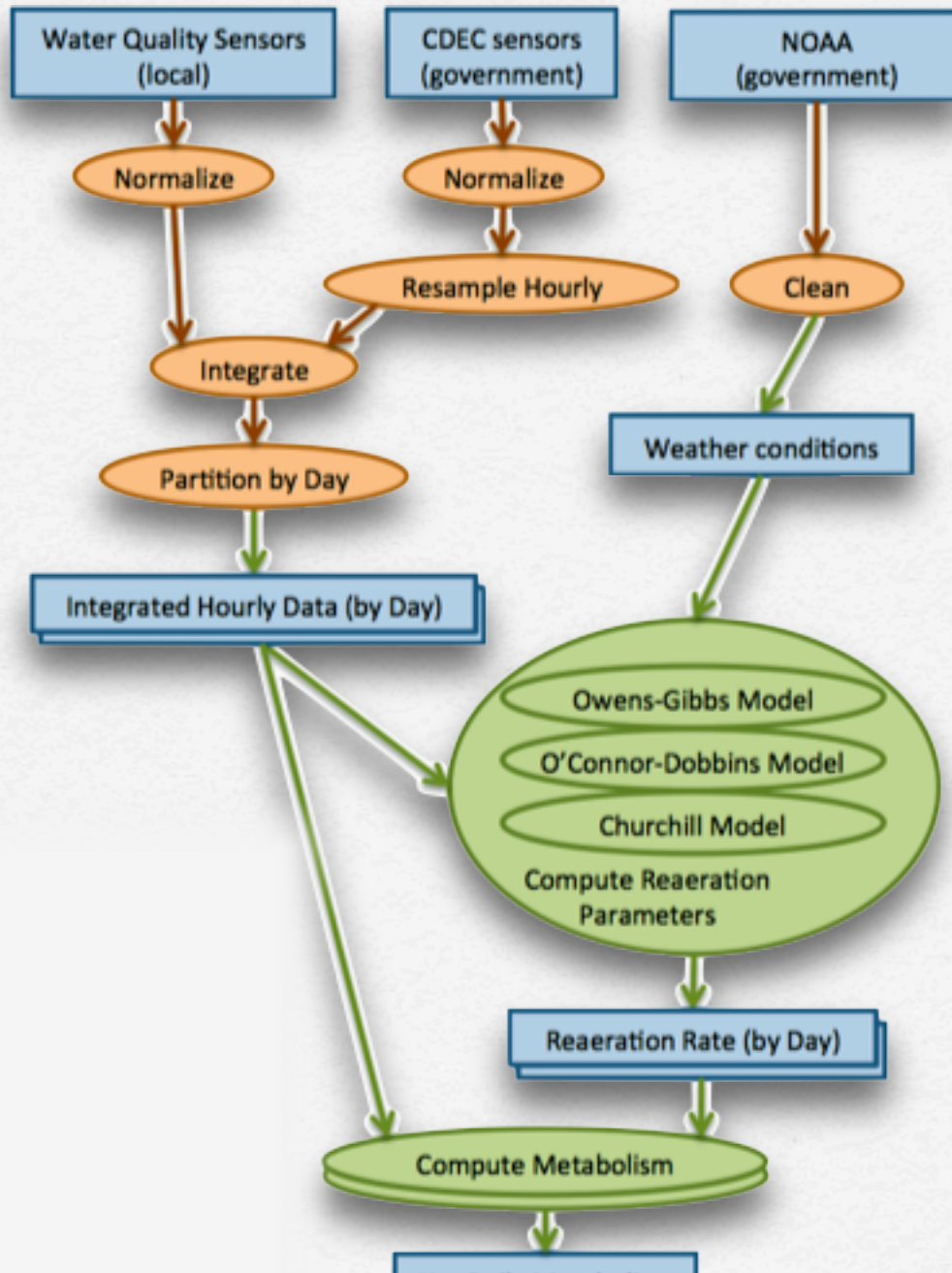## Models of gross primary production (GPP), community respiration (CR24)

### Sensors

### Analysis

*alyses must be fast to produce actionable information*

# Workflow



Water Quality Sensors (local) → Normalize → Integrate

CDEC sensors (government) → Normalize → Resample Hourly → Integrate

Integrate → Partition by Day → Integrated Hourly Data (by Day)

NOAA (government) → Clean → Weather conditions

Owens-Gibbs Model
O'Connor-Dobbins Model
Churchill Model
Compute Reaeration Parameters

Reaeration Rate (by Day)

Compute Metabolism

Tom Harmon
environmental system

# Vision: Automated & Fast

# Reality: Difficult & Time Consuming

?

# Current Method

Water Quality Sensors (local)

CDEC sensors (government)

NOAA (government)

Manua

Normalize

Normalize

Normalize

Resample Hourly

Clean

Integrate
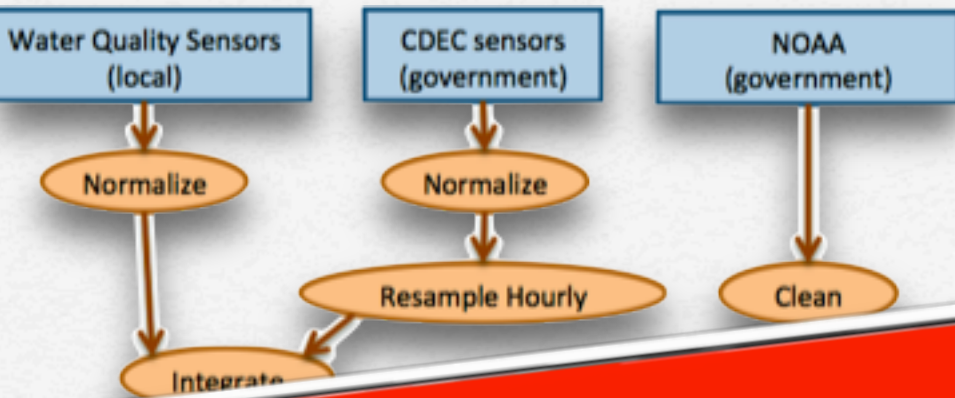
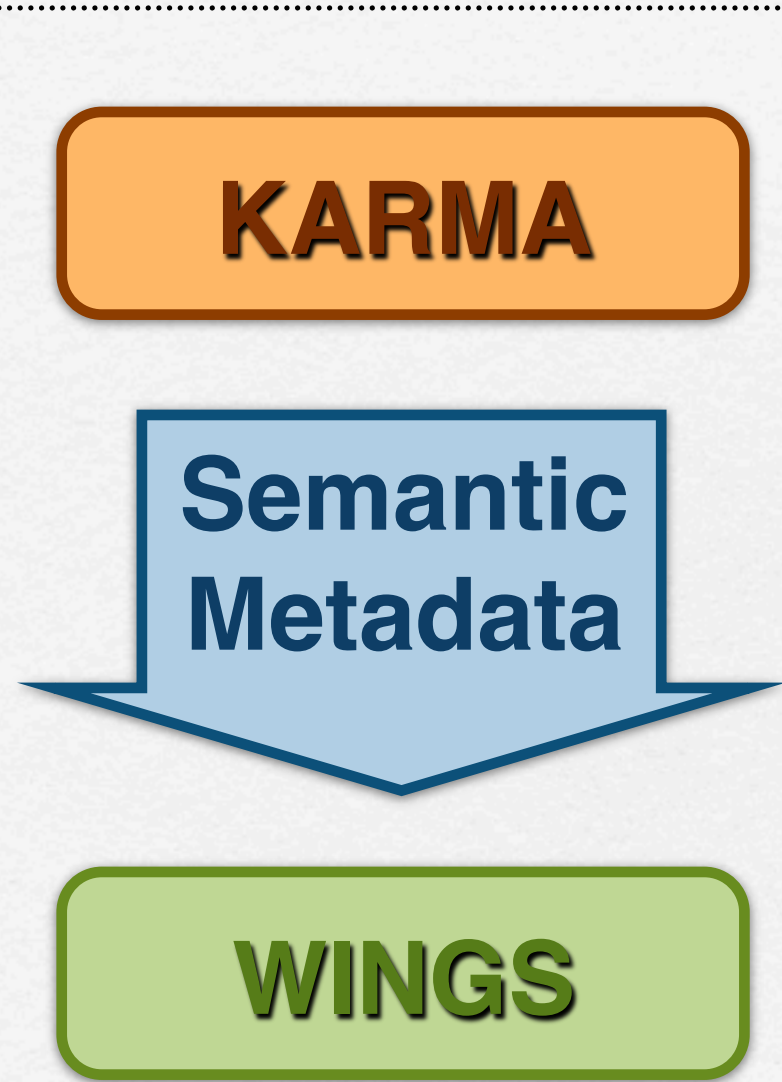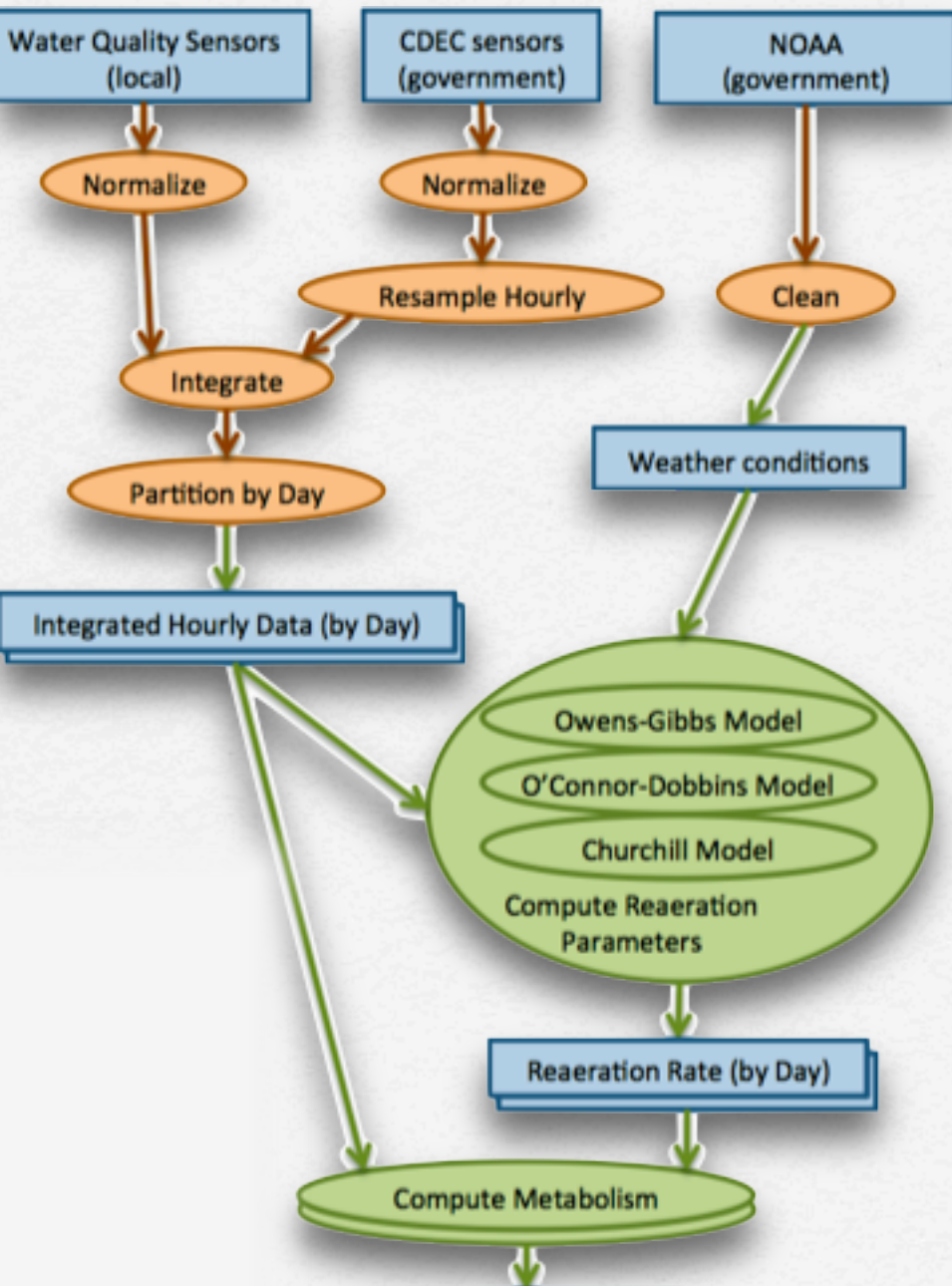Multiple, separate tools

High learning costs

Ad hoc, by-hand movement of data & tool invocation

Data does not "flow" across tools

Scripts

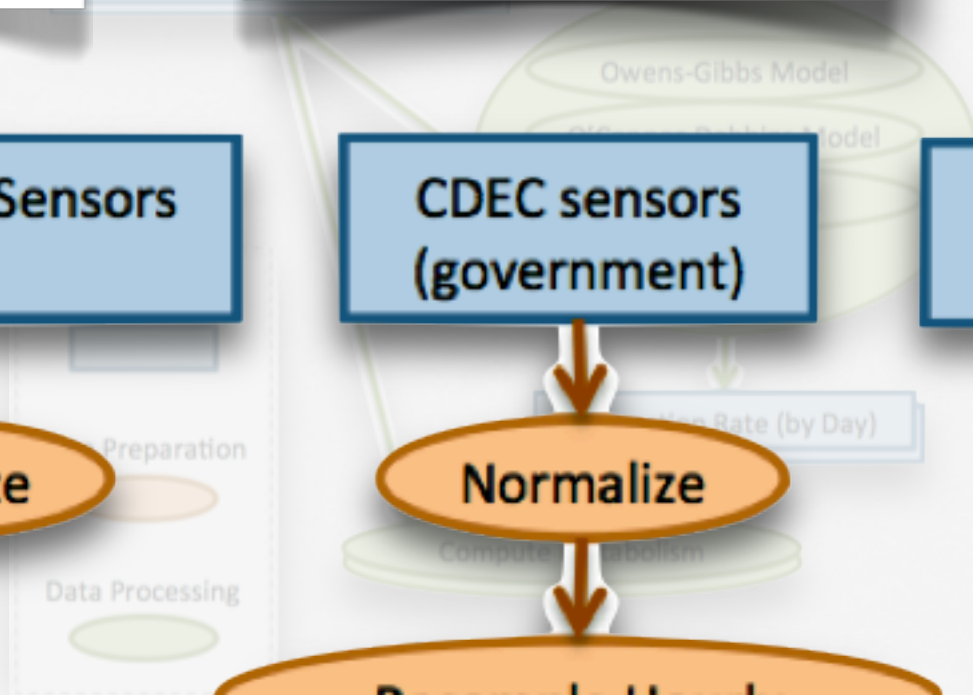Compute Metabolism

# Our Approach

# Data Sources
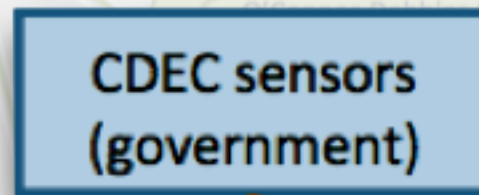
Water Quality Sensors
(local)

CDEC sensors
(government)

NOAA
(government)

Normalize

Normalize

Normalize

# KARMA

# Data Import

# Need to Clean Data

| | Time | | GE (feet) |
|---|---|---|---|
| 0309 | 2300 | | |
| 0309 | | | |

CI

HYDRO

| | | Temp | Cond |
|---|---|---|---|
| 010 | 23:00 | 13.4 | 1181.00 |
| 010 | 23:15 | 13.4 | 1179.00 |

60 Files for 1 month!

Requ

Fo

# Data Cleaning with KARMA

**Karma_v0.4**

Table   Script   Alignment   Column

| CDEC - Event Data0 | Source2 | Source3 | Cleaning Table |

| String | Data Type | Data Type |
|---|---|---|
| Date | User Defined Values | Final Values |
| 20100309 | 03/09/2010 | |
| 20100309 | | |
| 20100309 | | |
| 20100309 | | |
| 20100310 | | |
| 20100310 | | |
| 20100310 | | |
| 20100310 | | |
| 20100310 | | |
| 20100310 | | |
| 20100310 | | |
| 20100310 | | |

Import   Clean   Integrate   Publish

Data Cleaning for:  Column 2

Final result:

◯ Use original extracted values

# Data Cleaning with KARMA

# Integrated Dataset

pt   **Alignment**   **Column**

| | CDEC - Event Data1 | CDEC - Event Data2 | CDEC - Event Data3 | CDEC - Event Data4 | weatherstationdata.csv | CDEC |
|---|---|---|---|---|---|---|
| | forDate | String | String | String | String | |
| D | **Start Date** | **Date** | **Time** | **Temp** | **Cond** | |
| | 03/10/2010 | 03/09/2010 | 23:00 | 13.4 | 1181.00 | |
| | 03/10/2010 | 03/09/2010 | 23:15 | 13.4 | 1179.00 | |
| | 03/10/2010 | 03/09/2010 | 23:30 | 13.4 | 1184.00 | |
| | 03/10/2010 | 03/09/2010 | 23:45 | 13.3 | 1185.00 | |
| | 03/10/2010 | 03/10/2010 | 00:00 | 13.3 | 1185.00 | |
| | 03/10/2010 | 03/10/2010 | 00:15 | 13.3 | 1183.00 | |

**Import**   **Clean**   **Integrate**   **Publish**

| SMN | 03/10/2010 | 03/10/2010 | 00:30 | 13.2 |
| SMN | 03/10/2010 | 03/10/2010 | 00:45 | 13.2 |
| SMN | 03/10/2010 | 03/10/2010 | 01:00 | 13.2 |
| SMN | 03/10/2010 | 03/10/2010 | 01:15 | 13.2 |
| SMN | 03/10/2010 | 03/10/2010 | 01:30 | 13.1 |

| mport | Clean | Integrate | Publish |

| HTML | KML | XML | CSV Text File | Database | RDF | WebService |

## Web Services

### WebService Name
WINGS Portal

## Inputs

File Name

WEATHER_2010_03_10

File Content

CDEC - Event Data0

# Semantic Metadata for Input Files

```xml
<?xml version="1.0" encoding="UTF-8" ?>

...
...base="http://www.isi.edu/dc/Water/library...
...s:rdf="http://www.w3.org/1999/02/22-rdf-s...
...s:rdfs="http://www.w3.org/2000/01/rdf-sch...
...s:owl="http://www.w3.org/2002/07/owl#"
...s:xsd="http://www.w3.org/2001/XMLSchema#"
...s:dc="http://www.isi.edu/dc/ontology.owl#"
...s:dcdom="http://www.isi.edu/dc/Water/ontol...
...s="http://www.isi.edu/dc/Water/library.owl#...

...Daily_Sensor_Data rdf:ID="FILENAME">
...om:forDate rdf:datatype="http://www.w3.org/...
...om:forSite rdf:datatype="http://www.w3.org/...
...om:siteLatitude rdf:datatype="http://www.w3...LATITUDE</dcdom:siteLatit...
...om:siteLongitude rdf:datatype="http://www.w3...2001/XMLSchema#float">LONGITUDE</dcdom:siteLon...
...om:slope rdf:datatype="http://www.w3.org/2001/XMLSchema#float">SLOPE</dcdom:slope>
...om:velocity rdf:datatype="http://www.w3.org/2001/XMLSchema#float">VELOCITY</dcdom:velocity>
...om:depth rdf:datatype="http://www.w3.org/2001/XMLSchema#float">DEPTH</dcdom:depth>
...om:flow rdf:datatype="http://www.w3.org/2001/XMLSchema#float">FLOW</dcdom:flow>
...om:barpress rdf:datatype="http://www.w3.org/2001/XMLSchema#float">760</dcdom:barpress>
...Daily_Sensor_Data>
```

Automatically Generated by KARMA

# Workflows with WINGS



Conceptual

WINGS

Workflow

Metadata automatically associated with each input file

Data

DataBrowser | CDEC_WEATHER ×

View File | Delete File

Metadata for CDEC_WEATHER_2010_03_10

Daily_Data
  Daily_Parameters
    NTM_Parameters
  Daily_Sensor_Data
    CDEC_WEATHER_2010_03_02
    CDEC_WEATHER_2010_03_03
    CDEC_WEATHER_2010_03_04
    CDEC_WEATHER_2010_03_05
    CDEC_WEATHER_2010_03_06
    CDEC_WEATHER_2010_03_07
    CDEC_WEATHER_2010_03_08
    CDEC_WEATHER_2010_03_09

| Name ▲ | Value |
| --- | --- |
| barpress | 760 |
| depth | 1.6564940150390628 |
| flow | 1213.7113 |
| forDate | 2010-03-10 |
| forSite | SMN |
| siteLatitu... | 37.347214 |
| siteLongi... | -120.976181 |
| slope | 0.0001 |
| usedAlg... | |

# Workflow



Integrated Hourly Data (by Day)

Owens-Gibbs Model

O'Connor-Dobbins Model

Churchill Model

Compute Reaeration Parameters

Reaeration Rate (by Day)

Compute Metabolism

Meta



```
<dcdom:Hydrolab_Sensor_Data rdf:ID="Hydrolab-CDEC-04272011">
    <dcdom:siteLong rdf:datatype="float">-120.931</dcdom:siteLongitu
    <dcdom:siteLatitude rdf:datatype="float">37.371</dcdom:siteLatit
    <dcdom:dateStart rdf:datatype="date">2011-04-27</dcdom:dateStart
    <dcdom:forSite rdf:datatype="string">MST</dcdom:forSite>
    <dcdom:numberOfDayNights rdf:datatype="int">1</dcdom:numberOfDay
    <dcdom:avgDepth rdf:datatype="float">4.523957</dcdom:avgDepth>
    <dcdom:avgFlow rdf:datatype="float">2399</dcdom:avgFlow>
</dcdom:Hydrolab_Sensor_Data>
```

Meta



```
<dcdom:Hydrolab_Sensor_Data_rdf:ID="Hydrolab-CDEC-04272011">
    <dcdom:siteLong rdf:datatype="float">-120.931</dcdom:siteLongitu
    <dcdom:siteLatitude rdf:datatype="float">37.371</dcdom:siteLatit
    <dcdom:dateStart rdf:datatype="date">2011-04-27</dcdom:dateStart
    <dcdom:forSite rdf:datatype="string">MST</dcdom:forSite>
    <dcdom:numberOfDayNights rdf:datatype="int">1</dcdom:numberOfDay
    <dcdom:avgDepth rdf:datatype="float">4.523957</dcdom:avgDepth>
    <dcdom:avgFlow rdf:datatype="float">2399</dcdom:avgFlow>
</dcdom:Hydrolab_Sensor_Data>
```

etting
meters

Meta

FormattedData

FilterTimestampsAndData

FilteredData

CalculateHourlyAverages
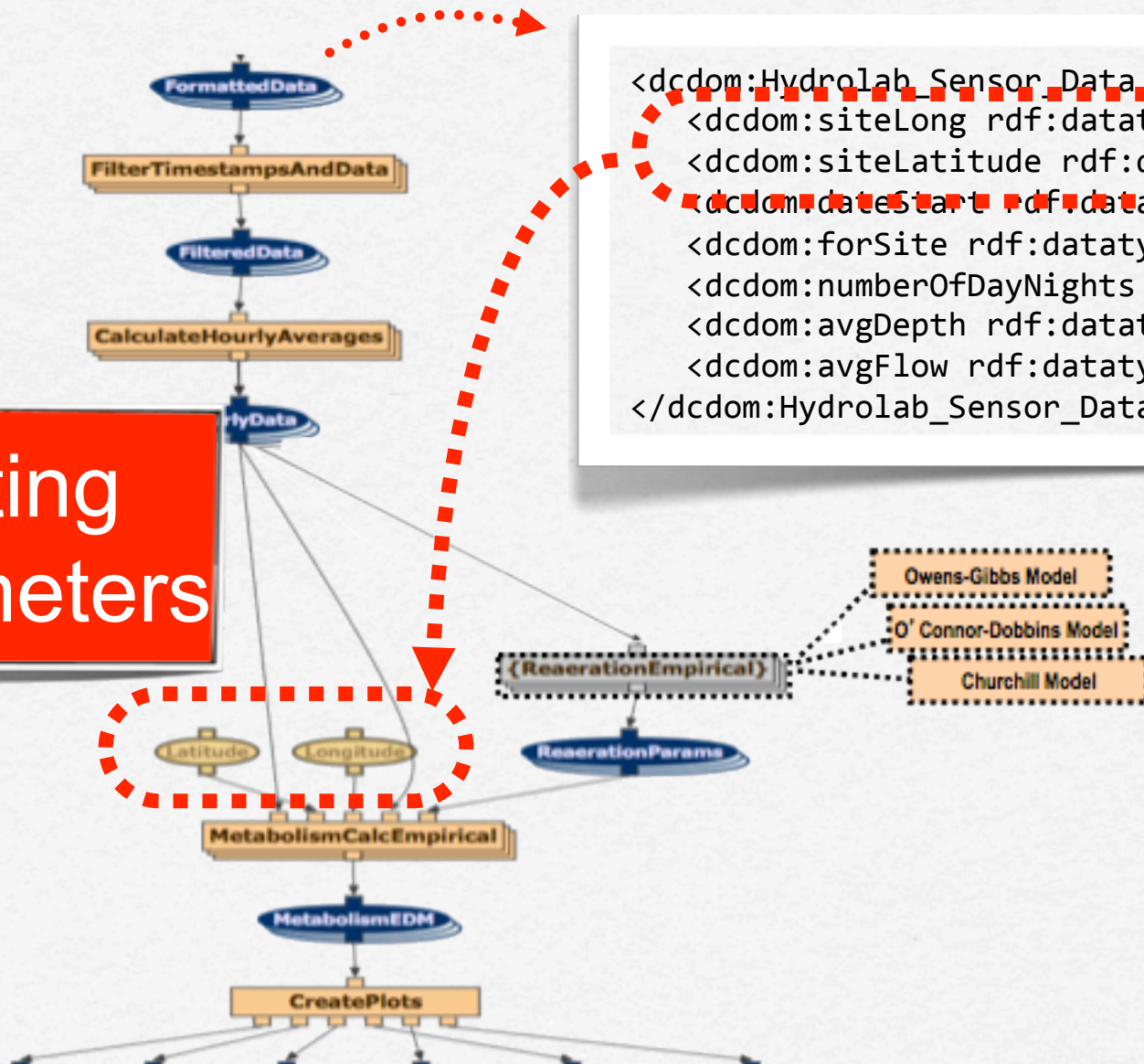
rlyData

```
<dcdom:Hydrolab_Sensor_Data_rdf:ID="Hydrolab-CDEC-04272011">
    <dcdom:siteLong rdf:datatype="float">-120.931</dcdom:siteLongitu
    <dcdom:siteLatitude rdf:datatype="float">37.371</dcdom:siteLatit
    <dcdom:dateStart rdf:datatype="date">2011-04-27</dcdom:dateStart
    <dcdom:forSite rdf:datatype="string">MST</dcdom:forSite>
    <dcdom:numberOfDayNights rdf:datatype="int">1</dcdom:numberOfDay
    <dcdom:avgDepth rdf:datatype="float">4.523957</dcdom:avgDepth>
    <dcdom:avgFlow rdf:datatype="float">2399</dcdom:avgFlow>
</dcdom:Hydrolab_Sensor_Data>
```

etting
meters

Owens-Gibbs Model

O' Connor-Dobbins Model

{ReaerationEmpirical}

Churchill Model

Latitude    Longitude

ReaerationParams

Choosing
Models

MetabolismCalcEmpirical

MetabolismEDM

CreatePlots

# Workflow Results

DO_MST_2011-01-01_0 *(3 KB, 🖫Save)*, **DO_MST_2011-01-01_1** *(3 KB, 🖫Save)*, **DO_MST_2011-01-01_2** *(3 KB, 🖫Save)*,
MST_2011-01-01_3 *(3 KB, 🖫Save)*, **DO_MST_2011-01-01_4** *(3 KB, 🖫Save)*, **DO_MST_2011-01-01_5** *(3 KB, 🖫Save)*,
MST_2011-01-01_6 *(3 KB, 🖫Save)*, **DO_MST_2011-01-01_7** *(3 KB, 🖫Save)*, **DO_MST_2011-01-01_8** *(3 KB, 🖫Save)*,
MST_2011-01-01_9 *(2 KB, 🖫Save)*, **DO_MST_2011-01-01_10** *(2 KB, 🖫Save)*, **DO_MST_2011-01-01_11** *(2 KB, 🖫Save)*,
MST_2011-01-01_12 *(2 KB, 🖫Save)*, **DO_MST_2011-01-01_13** *(3 KB, 🖫Save)*, **DO_MST_2011-01-01_14** *(3 KB, 🖫Save)*,
MST_2011-01-01_15 *(3 KB, 🖫Save)*, **DO_MST_2011-01-01_16** *(3 KB, 🖫Save)*, **DO_MST_2011-01-01_17** *(3 KB, 🖫Save)*,
MST_2011-01-01_18 *(3 KB, 🖫Save)*, **DO_MST_2011-01-01_19** *(3 KB, 🖫Save)*, **DO_MST_2011-01-01_20** *(3 KB, 🖫Save)*,
MST_2011-01-01_21 *(3 KB, 🖫Save)*, **DO_MST_2011-01-01_22** *(3 KB, 🖫Save)*, **DO_MST_2011-01-01_23** *(3 KB, 🖫Save)*,
MST_2011-01-01_24 *(3 KB, 🖫Save)*, **DO_MST_2011-01-01_25** *(3 KB, 🖫Save)*, **DO_MST_2011-01-01_26** *(3 KB, 🖫Save)*,
MST_2011-01-01_27 *(3 KB, 🖫Save)*, **DO_MST_2011-01-01_28** *(3 KB, 🖫Save)*, **DO_MST_2011-01-01_29** *(3 KB, 🖫Save)* }

ams_MST_2011-01-01_0 *(59 B, 🖫Save)*, **Params_MST_2011-01-01_1** *(68 B, 🖫Save)*, **Params_MST_2011-01-01_2** *(66 B, 🖫Sa*
ams_MST_2011-01-01_3 *(68 B, 🖫Save)*, **Params_MST_2011-01-01_4** *(66 B, 🖫Save)*, **Params_MST_2011-01-01_5** *(66 B, 🖫Save*
ams_MST_2011-01-01_6 *(66 B, 🖫Save)*, **Params_MST_2011-01-01_7** *(66 B, 🖫Save)*, **Params_MST_2011-01-01_8** *(66 B, 🖫Save*
ams_MST_2011-01-01_9 *(56 B, 🖫Save)*, **Params_MST_2011-01-01_10** *(46 B, 🖫Save)*, **Params_MST_2011-01-01_11** *(56 B, 🖫Sa*
ams_MST_2011-01-01_12 *(66 B, 🖫Save)*, **Params_MST_2011-01-01_13** *(57 B, 🖫Save)*, **Params_MST_2011-01-01_14** *(68 B, 🖫*
ams_MST_2011-01-01_15 *(66 B, 🖫Save)*, **Params_MST_2011-01-01_16** *(66 B, 🖫Save)*, **Params_MST_2011-01-01_17** *(68 B, 🖫*
ams_MST_2011-01-01_18 *(68 B, 🖫Save)*, **Params_MST_2011-01-01_19** *(68 B, 🖫Save)*, **Params_MST_2011-01-01_20** *(68 B, 🖫*
ams_MST_2011-01-01_21 *(68 B, 🖫Save)*, **Params_MST_2011-01-01_22** *(66 B, 🖫Save)*, **Params_MST_2011-01-01_23** *(68 B, 🖫*
ams_MST_2011-01-01_24 *(66 B, 🖫Save)*, **Params_MST_2011-01-01_25** *(66 B, 🖫Save)*, **Params_MST_2011-01-01_26** *(57 B, 🖫*
ams_MST_2011-01-01_27 *(68 B, 🖫Save)*, **Params_MST_2011-01-01_28** *(66 B, 🖫Save)*, **Params_MST_2011-01-01_29** *(66 B, 🖫*

Params_MST_2011-01-01_6 *(66 B, 🖫Save)*, Params_MST_2011-01-01_7 *(66 B, 🖫Save)*, Params_MST_2011-01-01_8 *(66 B, 🖫Save)*,
Params_MST_2011-01-01_9 *(56 B, 🖫Save)*, Params_MST_2011-01-01_10 *(46 B, 🖫Save)*, Params_MST_2011-01-01_11 *(56 B, 🖫Save)*,
Params_MST_2011-01-01_12 *(66 B, 🖫Save)*, Params_MST_2011-01-01_13 *(57 B, 🖫Save)*, Params_MST_2011-01-01_14 *(68 B, 🖫Save)*,
Params_MST_2011-01-01_15 *(66 B, 🖫Save)*, Params_MST_2011-01-01_16 *(66 B, 🖫Save)*, Params_MST_2011-01-01_17 *(68 B, 🖫Save)*,
Params_MST_2011-01-01_18 *(68 B, 🖫Save)*, Params_MST_2011-01-01_19 *(68 B, 🖫Save)*, Params_MST_2011-01-01_20 *(68 B, 🖫Save)*,
Params_MST_2011-01-01_21 *(68 B, 🖫Save)*, Params_MST_2011-01-01_22 *(66 B, 🖫Save)*, Params_MST_2011-01-01_23 *(68 B, 🖫Save)*,
Params_MST_2011-01-01_24 *(66 B, 🖫Save)*, Params_MST_2011-01-01_25 *(66 B, 🖫Save)*, Params_MST_2011-01-01_26 *(57 B, 🖫Save)*,
Params_MST_2011-01-01_27 *(68 B, 🖫Save)*, Params_MST_2011-01-01_28 *(66 B, 🖫Save)*, Params_MST_2011-01-01_29 *(66 B, 🖫Save)* }

7   OutputHourlyAvgedData   { AvgHourly_MST_2011-01-01_0 *(882 B, 🖫Save)*, **AvgHourly_MST_2011-01-01_1** *(871 B, 🖫Save)*, **AvgHourly_MST_2011-01-01_2** *(887 B,*

WINGS automatically generates metadata for each output file

```
<dcdom:Metabolism_Results  rdf:ID="Metabolism_Results-CDEC-04272011">
    <dcdom:siteLong rdf:datatype="float">-120.931</dcdom:siteLongitude>
    <dcdom:siteLatitude rdf:datatype="float">37.371</dcdom:siteLatitude>
    <dcdom:dateStart rdf:datatype="date">2011-04-27</dcdom:dateStart>
    <dcdom:forSite rdf:datatype="string">MST</dcdom:forSite>
    <dcdom:numberOfDayNights rdf:datatype="int">1</dcdom:numberOfDayNights>
    <dcdom:avgDepth rdf:datatype="float">4.523957</dcdom:avgDepth>
    <dcdom:avgFlow rdf:datatype="float">2399</dcdom:avgFlow>
</dcdom: Metabolism_Results>
```
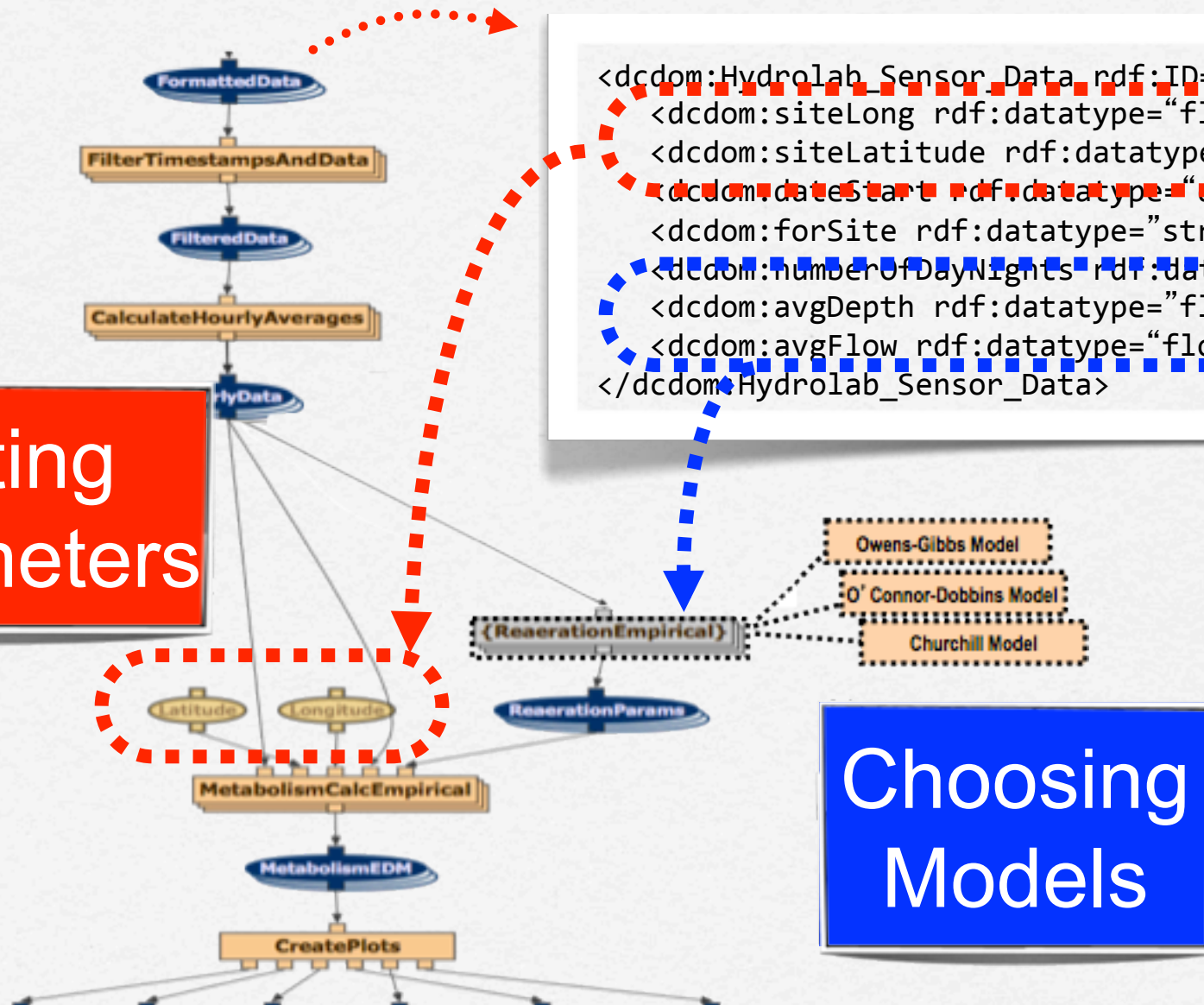
```
LECT ?url WHERE {
ata dcdom:usedAlgorithm dcdom:ODM .
ata rdf:type dcdom:Metabolism_Estimates .
ata wflow:hasLocation ?url
```
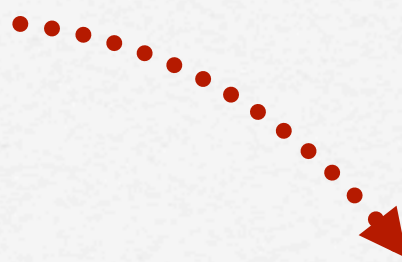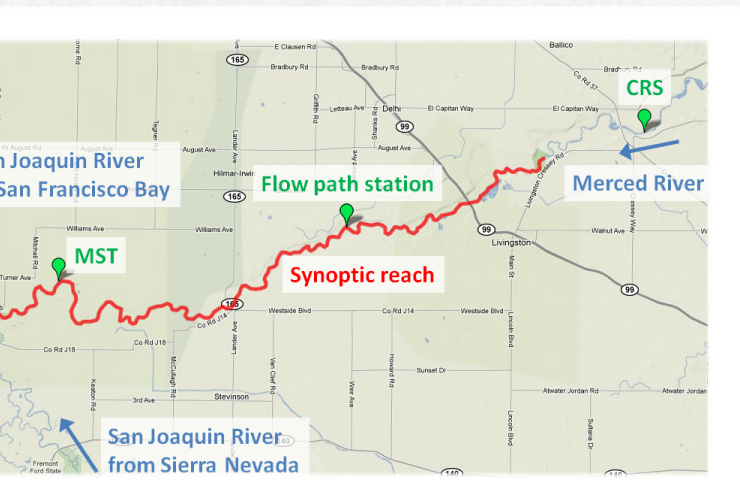
**Metadata for Metabolism_SMN_2010_03_03Z_0**

Save Metadata

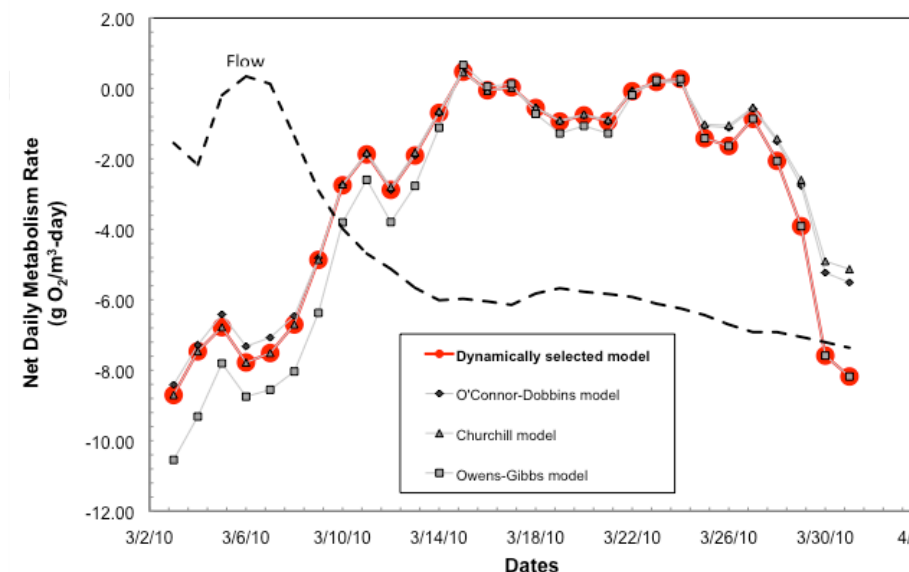| Name | Value |
|---|---|
| velocity | 0.66163415 |
| usedAlgorithm | dcdom:ODM |
| slope | 1.0E-4 |
| siteLongitude | -120.97618 |
| siteLatitude | 37.347214 |
| forSite | SMN |
| forDate | 2010-03-03Z |
| flow | 1581.6842 |
| depth | 1.0403947 |

# Aquatic Photosynthesis

## Models of gross primary production (GPP), community respiration (CR24)
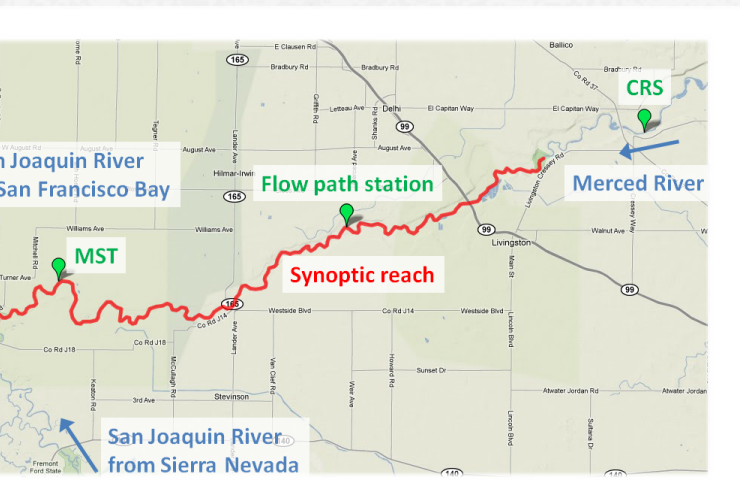
Sensors



Workflow Results

# Aquatic Photosynthesis

## Models of gross primary production (GPP), community respiration (CR24)

Sensors



Workflow Results

# Summary


Metadata Inside

- Tools for end-users

- End to end support

- Data import, cleaning, integration

- Automated workflow execution

- Captures metadata provenance

# Related Work

- Data integration:
  Data Wrangler [Kandel et al 2011]
  Google Refine [Huynh et al]

- Workflow systems:
  VisTrails [Howe et al 2008],
  Kepler [Barseghian et al 2010]

- Many tools generate provenance metadata, often in RDF

  - None generate other kinds of metadata