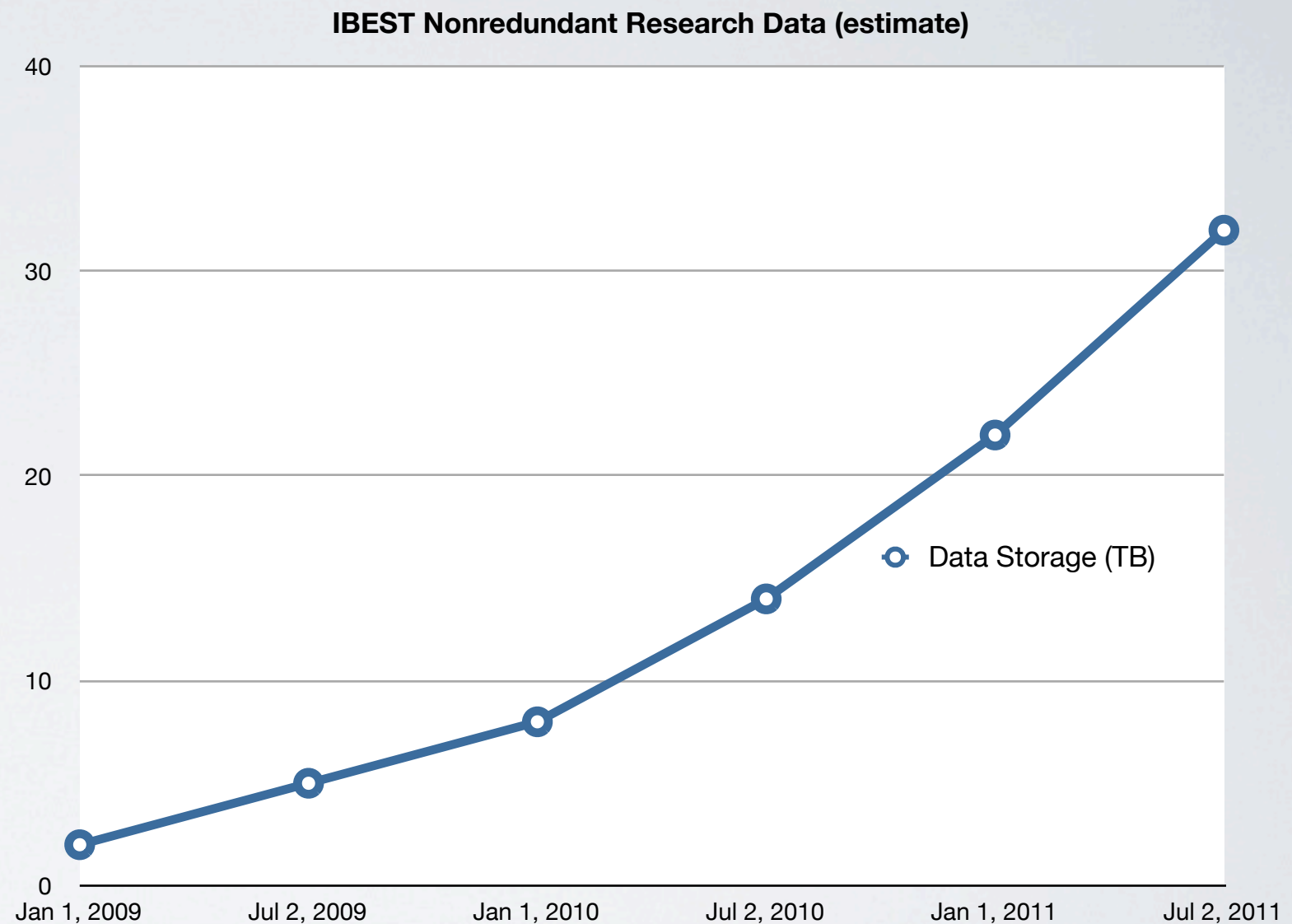# HATCH

USGS-CRRP

**?**

Why?

# Introduction

- Exponential scientific data growth

- CPU computer power cannot keep up

- Huge emerging gap between researchers and results

- Big problems will remain unanswered

- Less investment in computing equipment, more in staff

[1] Szalay, Alex; Gray, Jim. **"2020 Computing: Science in an exponential world"**. Nature 440, 413-414 (23 March 2006)

[2] Foster, James. **Visualizing Human Microbiome Ecosystems**. University of Idaho: Computer Science Colloquium, December 7th 2010. Seminar.

# Exponential data growth

**IBEST Nonredundant Research Data (estimate)**



- Doubling every 2 years [1]

- U of I IBEST data is ~33 TB, growing

- <10% data used in final research [2].

[1] Szalay, Alex; Gray, Jim. **"2020 Computing: Science in an exponential world"**. Nature 440, 413-414 (23 March 2006)

[2] Foster, James. **Visualizing Human Microbiome Ecosystems**. University of Idaho: Computer Science Colloquium, December 7th 2010. Seminar.
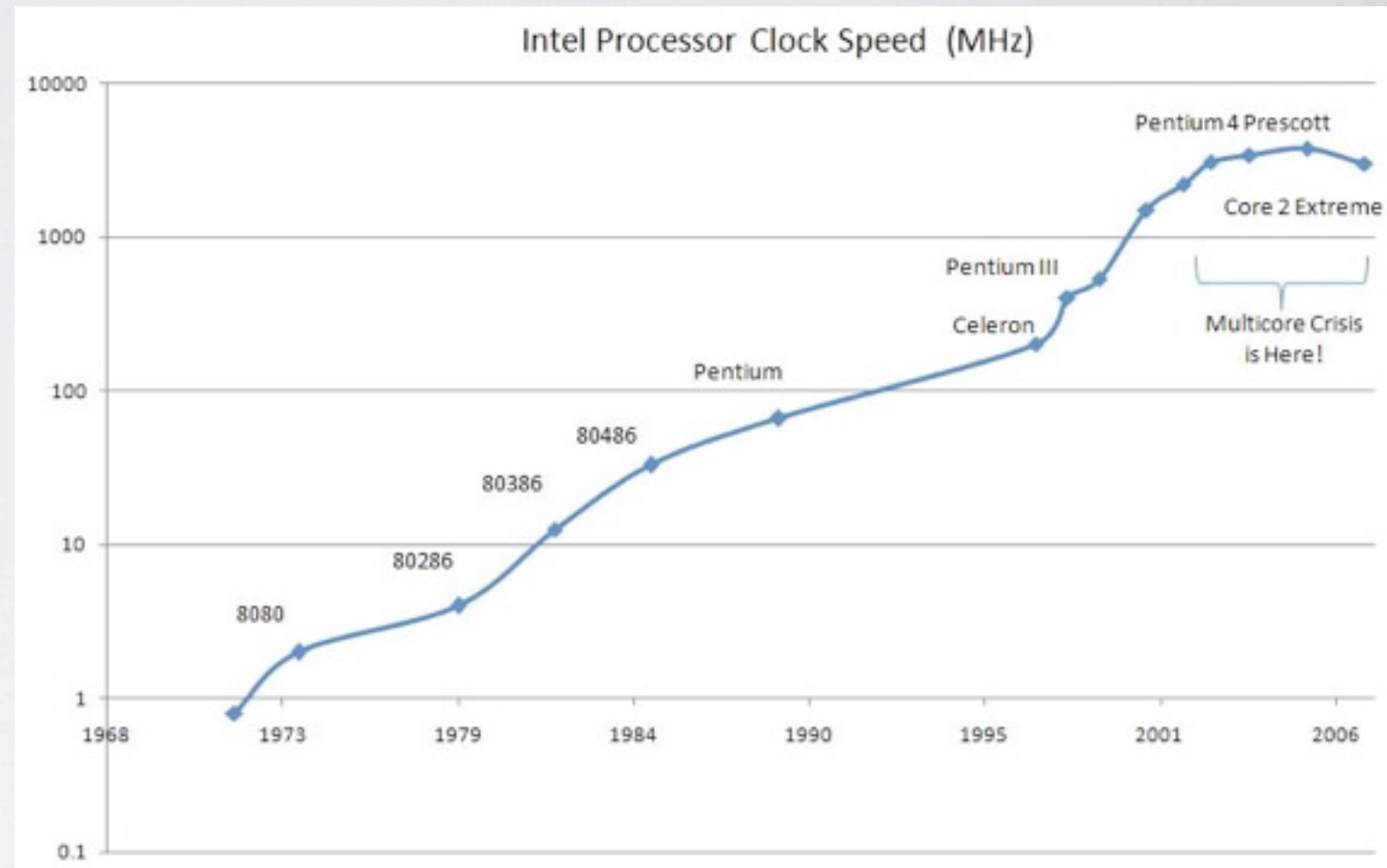
[1] Szalay, Alex; Gray, Jim. **"2020 Computing: Science in an exponential world"**. Nature 440, 413-414 (23 March 2006)

[2] Foster, James. **Visualizing Human Microbiome Ecosystems**. University of Idaho: Computer Science Colloquium, December 7th 2010. Seminar.

[3] Manek Dubash (2005-04-13). "Moore's Law is dead, says Gordon Moore". Techworld. Retrieved 2006-06-24

# The death of Moore's Law

- Back in 1965, Moore's Law declared [3]

- Transistors density doubled every 2 years

- CPU power roughly the same

- Growth continued for 40 years

- Transistors approach to size of atom

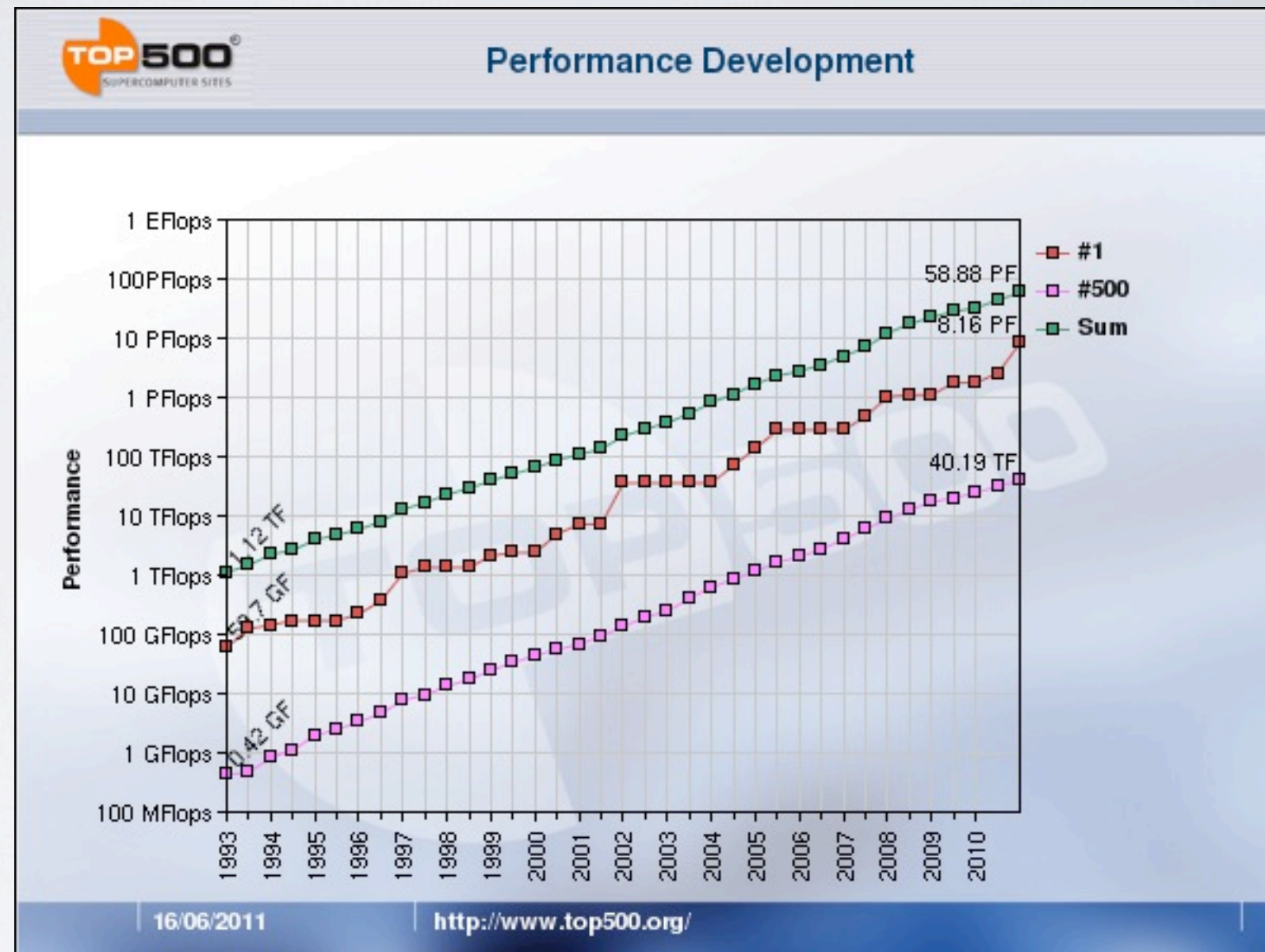- Who boldly declared the law dead?

Intel Processor Clock Speed (MHz)

[3] Manek Dubash (2005-04-13). "Moore's Law is dead, says Gordon Moore". Techworld. Retrieved 2006-06-24

[3] Manek Dubash (2005-04-13). "Moore's Law is dead, says Gordon Moore". Techworld. Retrieved 2006-06-24

# CPU Chips cannot keep up

- High Performance Compute Clusters

- Groups of computers, one huge CPU

- Amazingly powerful and complicated

Today's scientific computation is
too complicated

- the gap between data and analysis is broad
- data assets become management liabilities
- unnecessary computation and data redundancy is common
- the research workspace is more fractured than it is unified

# An interactive solution is needed

- requires no new training
- maximizes use of existing standards and libraries
- takes advantage of the existing use cases
- answers the current administrative questions
- doesn't create new ones

Hatch

?

# What is Hatch?

# Hatch is a simplified:

- •data visualization and filtration system
- •entry for local data to a research cloud
- •minimization of redundancy in data sharing
- •linkage between data acquisition and analysis

A graphical way to manage:

- data sets
- computational resources
- cloud network topologies
- processes

What is Hatch?

Networking Protocol

Web GUI

Analysis (future)

Metadata Server

Data IO

Visualization

Database

Import/Export

Graphing UI

- decentralized cloud
- seamless data management
- communication encouragement
- preservation of group autonomy

# Completed so far

Import Data
•Currently, EcoData supports data that is formatted in well-structured Comma-Separated Values (CSV) format.
•A CSV data file is uploaded via the web interface (pictured above) and stored in the database.

Manipulate Data

•Currently all data must be manipulated before being imported into EcoData. This limitation will be overcome as the database backend is finalized.

Visualize Data
•A visualization framework has been put in place and is nearly ready to accept user data.  Once the database structure is finalized, data will be graphed as seen in "Visualizations in the Web Browser."

# Data IO, Search, Filter

data to the cloud and back

# Sample data

```
# just draw a test/example chart
sample0 = [0, 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23#
sample1 = [34, 891, 274, 569, 724, 967, 599, 777, 896, 481, 844, 888, 691, 1045, 860, 596, 1211#
sample2 = [907, 414, 387, 32, 612, 943, 430, 844, 472, 938, 768, 319, 611, 930, 511, 813, 979, #
@data = sample0.zip(sample1)
@hc = LazyHighCharts::HighChart.new('visualization') do |f|
    f.options[:chart][:defaultSeriesType] = 'spline'
    f.series(:name=>'test series', :data=>@data)
    f.series(:name=>'test series 2', :data=>sample0.zip(sample2))
    f.options[:title] = {:text=>'test chart'}
    f.options[:xAxis][:title] = {:text=>'x axis'}
    f.options[:yAxis][:title] = {:text=>'y axis'}
end
```

# Sample graph

# Sample graphs

Alex Fremier, Associate Professor at the University of Idaho College of Natural Resources

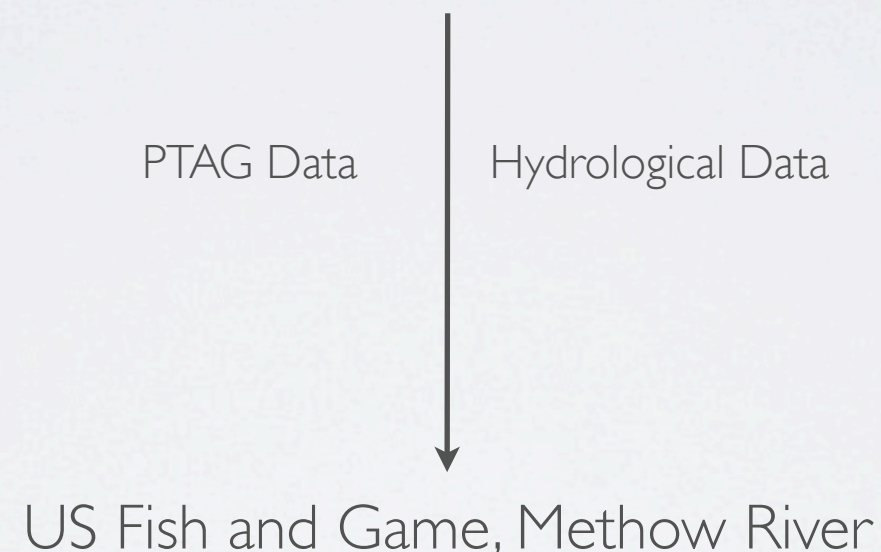PTAG Data | Hydrological Data

US Fish and Game, Methow River

Customers

Who else?

# Trial and error

- coupling design too closely to the data
- tool dependencies
- coupling data too closely to the tools
- tool module coupling

# Long term goals

While the simplest use of EcoData is a simple process involving importing and storage of data followed by manipulations and visualizations, each stage of the process can and will be elaborated upon.

- Import data
  o Support other data types
- Manipulate Data
  o Filtering
  o Joins between tables of data
- Visualize data
  o Add more types of visualizations

Other objectives include:
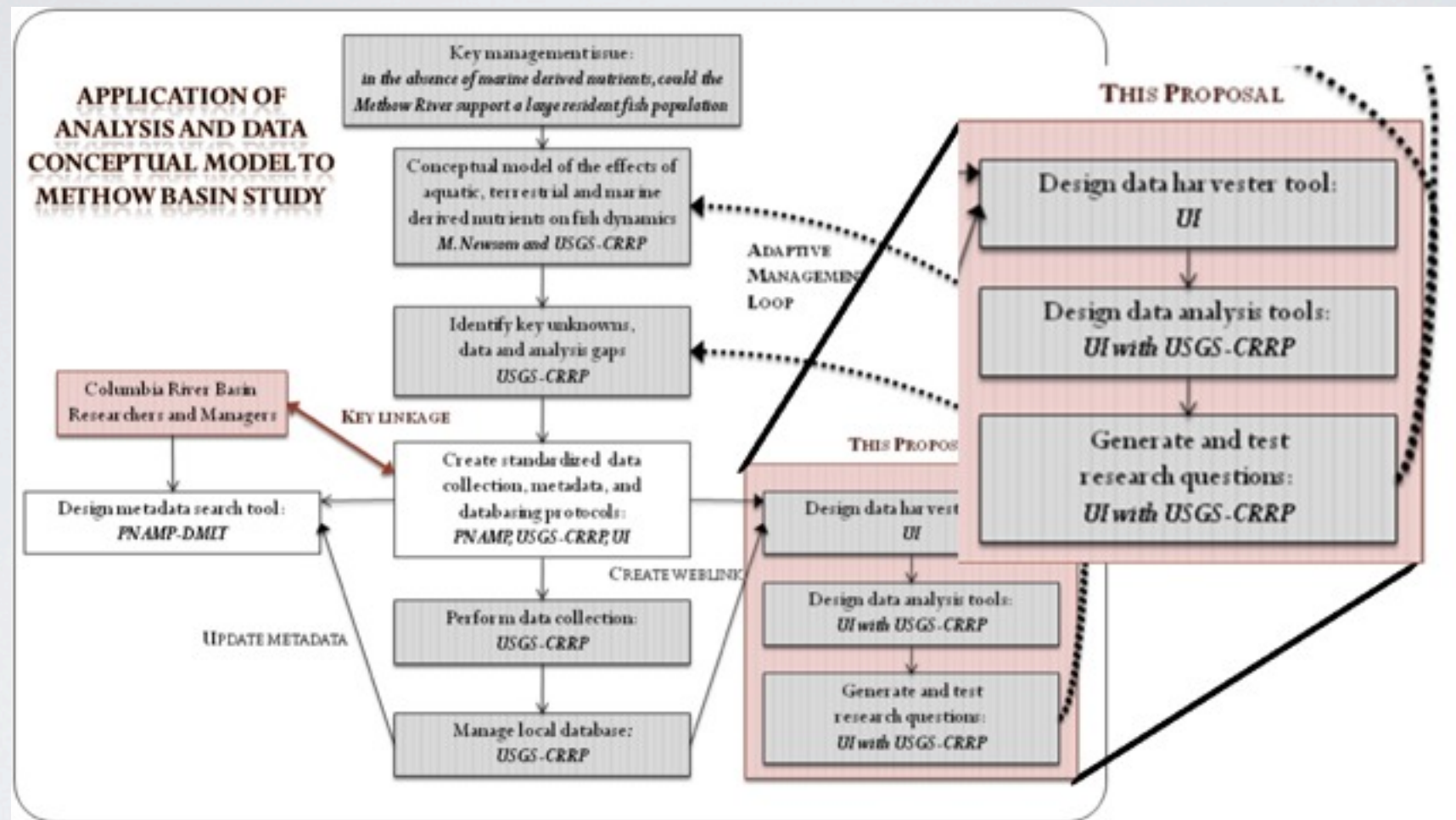- Authentication and security for data
- Availability of data between servers
- Cross-server data availability
  o Fast transfers of data between remote servers

# Future of Portal



- Portal is the first piece in a larger data management and analysis system.
- It provides the tools necessary to convert and store data in a structured format and to produce basic visualizations of that data.
- Eventually, it will be used as one piece in a more complex data analysis pipeline, allowing for more sophisticated data manipulations and visualizations.

# Conclusion

- Dire need of more processing power

- Data collection eclipsed by the inability to analyze it

- Big answers lie in the balance

- Data deluge will hold back advances

- People disciplined in computers are needed

- NSF and NIH grants favor funding data collection

- Actually analyzing the data

- Research will become less and less meaningful