# Workflow and Provenance

(anything profound, and the cool slides, is from Bertram Ludäscher. Everything else is from Renear

# What is data workflow?

Much of our work with data, especially in scientific applications, consists in *transforming one data set into another*

# Data curation and data workflow

Data curation is concerned with transformations in *two* ways:

managing and documenting transformations involved in data analytics

performing transformation to realize data curation objectives.
(preservation, integration, format conversion, etc.)

# Kinds of data transformations

**Transformations where input and output datasets are** <span style="color:red">identical</span> **in propositional content**

      transformation to a different data description language       (or new version of a language)

      transformation to a different serialization            (or new version of a serialization)

**Transformations where the input dataset** <span style="color:red">mathematically contains</span> **the output dataset**

      transformation to a subset matching specific conditions
          e.g. simple queries

      transformation to a logically or mathematically entailed data of the same kind
          e.g., summaries, statistics, visualizations

**Transformations where the input dataset** <span style="color:red">scientifically contains</span> **in the output dataset**

      transformation to scientifically entailed data of the same kind
          here the resulting data set typically contains information different in kind
              e.g., a data set about air pressure is transformed to a dataset about altitudes.

# Example Bioinformatics Workflow:

## *Motif-Catcher*

Marc Facciotti *et al.*
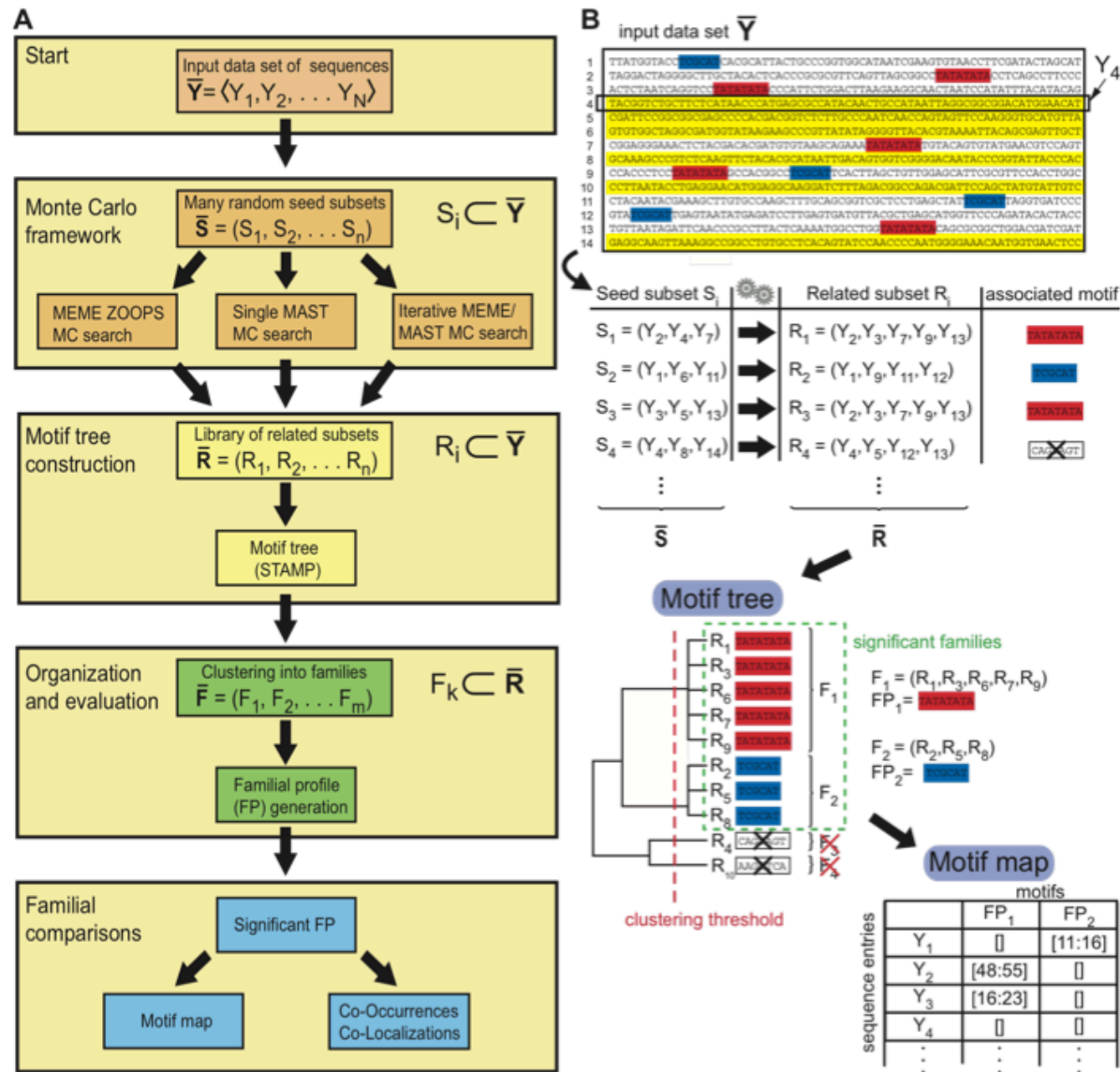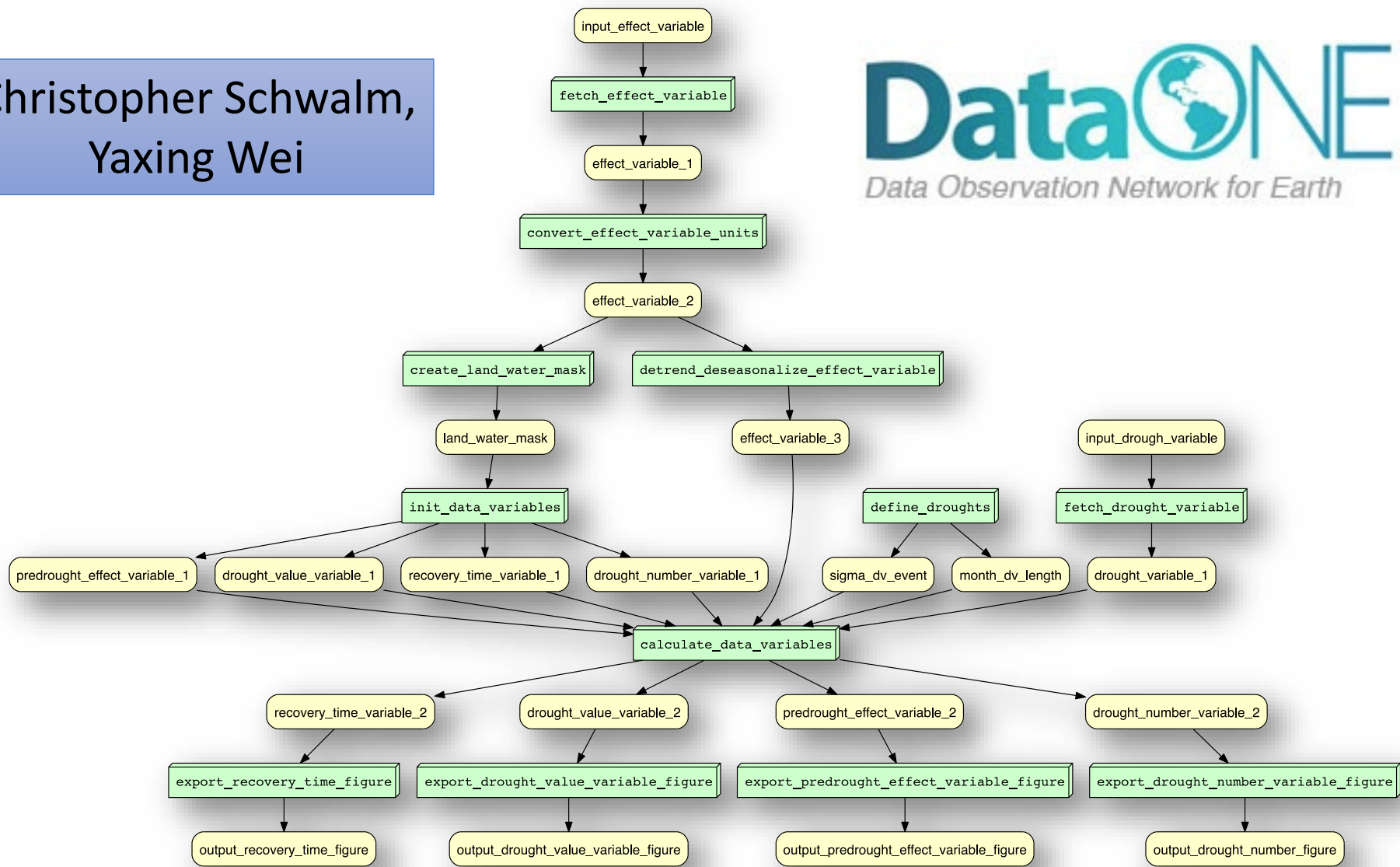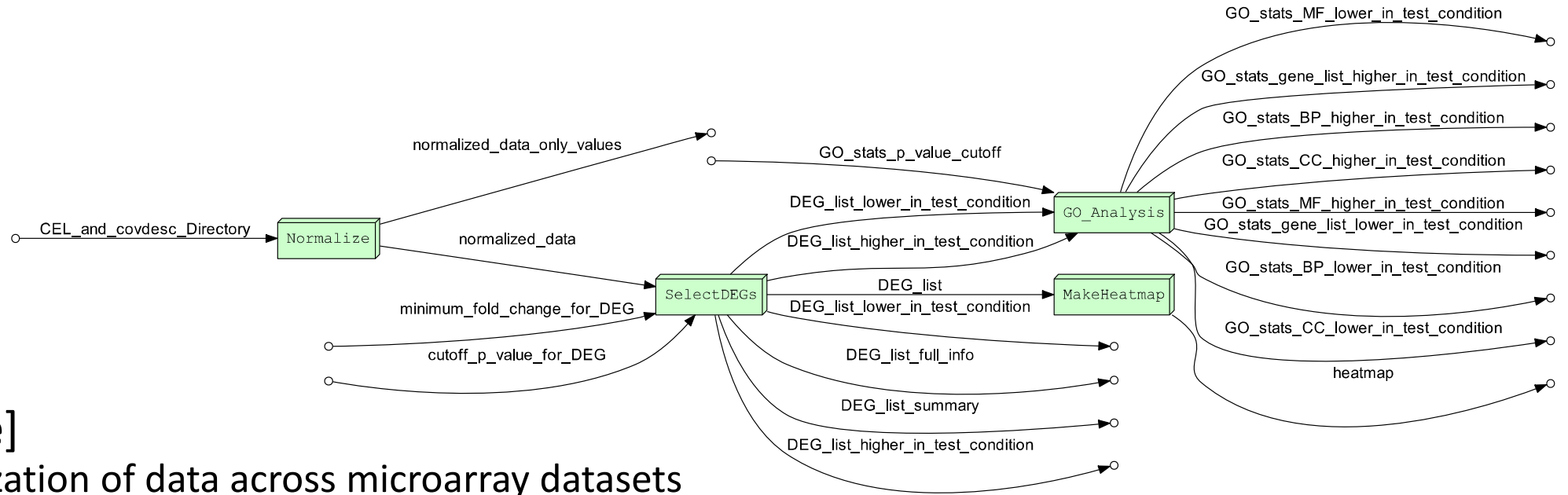UC Davis Genome Center



Figure 1: Concept of Monte-Carlo based detection and interpretation of motifs.
A) Abstract description of MotifCatcher process. B) Examples illustrating the
process with sample data.

# Multi-Scale Synthesis and Terrestrial Model Intercomparison Project (MsTMIP)

# Gene Expression Microarray Data Analysis



- **[Normalize]**
  - Normalization of data across microarray datasets
- **[SelectDEGs]**
  - Selection of differentially expressed genes between conditions
- **[GO Analysis]**
  - determination of gene ontology statistics for the resulting datasets
- **[MakeHeatmap]**
  - creation of a heatmap of the differentially expressed genes.

Tyler Kolisnik, Mark Bieda

# Why is workflow important

Thoughtfully designed organized workflows support:

Efficiency

Reliability

Modifiability

Reuse

Reproducibility

# Computational provenance

The heart of computational provenance:

*What data was used?*

*What calculations were performed?*

and also: "What in the world exactly happened just now?!"

# Why is provenance important?

Access to provenance information supports

>       understanding
>
>       reliability
>
>       reproducibility
>
>       trust
>
>       attribution and credit
>
>       discovery and reuse of data, tools, and algorithms

# Prospective vs Retrospective provenance (Ludaescher)

**Prospective**

a specification the workflow scenario

**Retrospective**

generated data on the execution of the workflow scenario

# Scripts are (or can be) workflow!!
# [try.yesworkflow.org]

# At least remember this

don't just sit there typing at the command line,
write a script, and document it

[for crying out loud]

# Communication

# Scientific communication is how data gets noticed



[After all, if the results of analyzing data are not communicated, then what's the point of it all?]

Scientific and technical communication is a critical part of the data lifecycle, with effects flowing both ways:
　　　— from the research process,
　　　— and back into the research process.

image from "Science 2.0 Repositories: Time for a Change in Scholarly Communication" Massimiliano Assante et al., *D-Lib Magazine* 2015.

# . . . the crisis

*But scientific and technical publishing is in crisis*

a problem caused by data
and that can be addressed with data

as we'll see in the next video

# Lisa's problem



from John Wilbanks, adapting a KEGG diagram

# *Are you kidding me???*



[Axis is x10$^{-2}$ for total Medline abstracts]

Adapted from Jensen, Saric, & Bork; *Nature* (2006).

# Faster, faster, faster, more more more



Tenopir et al. 1977-2005

# Responses to the problem

**One response:**        *text mining [*instead of L2R/T2B *reading]*

information extraction

"undiscovered public knowledge"
and hypothesis generation
(Swanson and Smalheiser

**Another response:**        *tools for **<span style="color:red">strategic reading</span>***

# Necessary data standards are now, finally, in place to support reading tools

**Character** encoding interoperability
       Unicode/UTF-xx                               [Adoption: nearly total]

**Data structure serialization** interoperability
       XML, JSON                                      [Adoption: nearly total]

**Syntactic** interoperability
       i.e. RDF(S), OWL                              [Adoption: underway]

**Semantic** interoperability
       RDF/OWL ontologies; linked data.               [Adoption: substantial]

**Document markup meta-languages**
       XML                                         [Adoption: nearly total]

**Document markup languages**
       e.g,   NLM/DTD, XHTML, TEI, DocBook, DITA    [Adoption: widely adopted]

**Metaphysical** interoperability
       "upper" ontologies                             [Adoption: (hard to say)]

**Domain ontologies and terminologies**
       hundreds                                    [Adoption: steady improvements]

Hoffmann, R; Valencia, A (Jul 2004).
"A gene network for navigating the literature.".
*Nature Genetics*. **36** (7): 664.

Muller HM, Kenny EE, Sternberg PW "Textpresso: an ontology-based information retrieval and extraction system for biological literature" *PLoS Biol*. 2004 Nov;2(11)..

# Data Mining?

*We wouldn't have to mine the data if we didn't bury in the first place.*

Barend Mons, "Which gene did you mean?" *BMC Bioinfomatics* (2005)

And finally. . .

Automate like you are going to live forever

Document like you are going to  die tomorrow

*— Michael Sperberg-McQueen*

# and for fun. . .

- https://www.youtube.com/watch?v=66oNv_DJuPc