# Collection-Oriented Scientific Workflows

## for

## Integrating and Analyzing Biological Data

**Shawn Bowers[1]**

**Timothy McPhillips[1]**

**Bertram Ludaescher[1,2]**

[1] UC Davis Genome Center

[2] Dept. of Computer Science

University of California, Davis

UCDAVIS genome center

**UC DAVIS** Department of Computer Science

**DILS06**

# Cutting to the chase: What are Scientific Workflows?

- **Probably the single-most important concept you hear about @DILS'06**

- **Attempto-Plain-English Definition:**
  - *SWFs := System designs and/or executable programs/scripts*
  - *… aiming to solve complex scientific data integration, analysis, management, visualization tasks*
    - in plainer English: doing hard and/or messy stuff …
  - *… while doing it in a scientist-friendly way*
    - that is: making it **look** easy
  - *… with the ultimate goal to*
  - *… do new, more, and better (e-)Science,*
  - *… faster!*

- **In short: SWFs are nothing less than MIRACLE-IT® to make scientists (biologists, physicists, ... ) happy.**

  Bio to CS/IT guy: "Please MIRACLE-IT® !" and "MIRACLE-IT®-Happen-**NOW**!!"

- **Attempto-DL Definition:**

  HappyScientist ≡ Scientist ⊓ ∀hasWF.HappyWF
  HappyWF ≡ Workflow ⊓ Reusable ⊓ Configurable ⊓ Efficient ⊓ . . .

UCDAVIS
UNIVERSITY OF CALIFORNIA
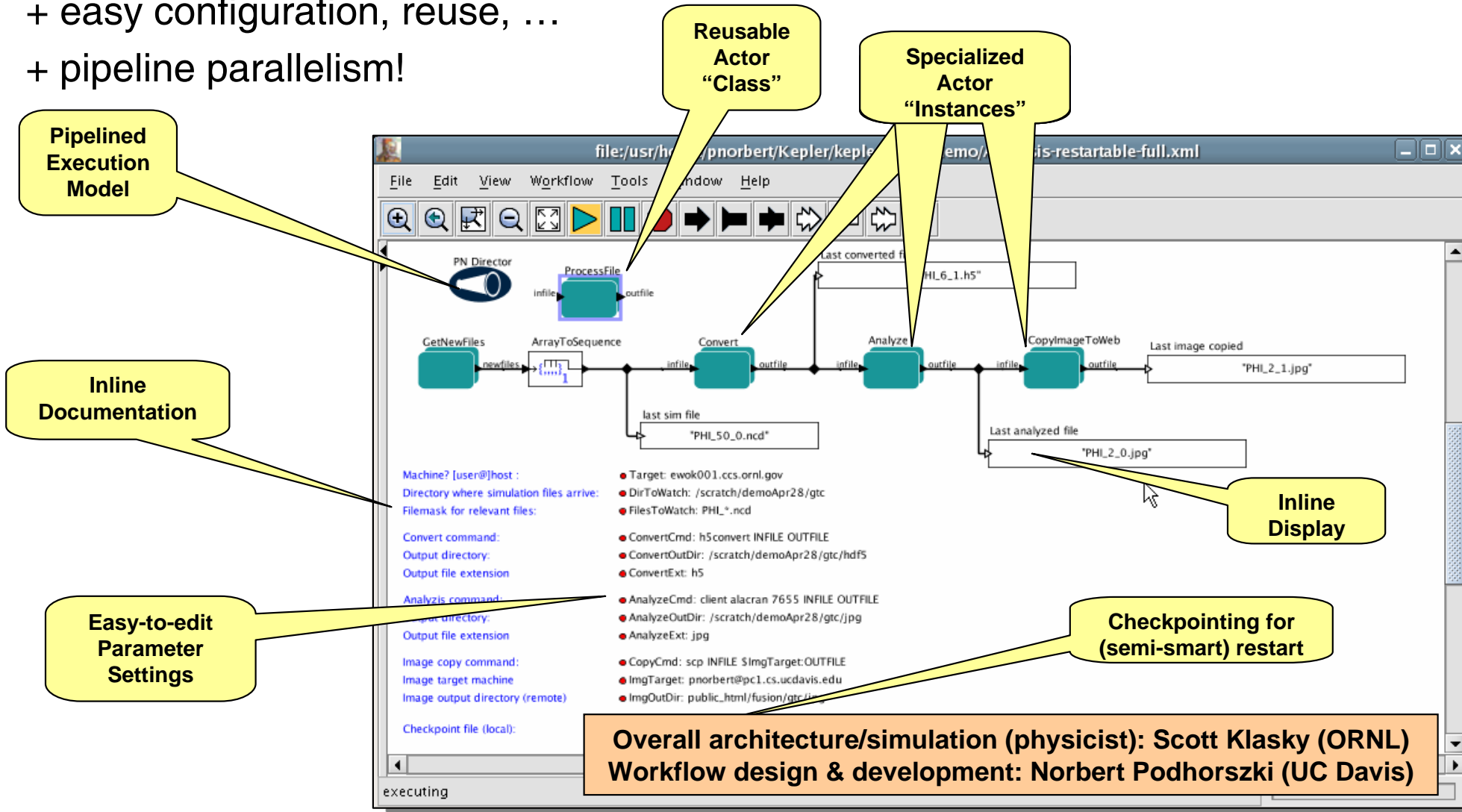
# Ex: A Happy Fusion Simulation Workflow
## Subspecies/Variety: *"Plumbing WF" (flux-laboris plumbiensis)*

- **Implements concurrent analysis pipeline (@2$^{ndary}$ cluster):**

  Tasks: *convert*; *analyze*; *copy-to-Web-portal*   (makes scientists **really** happy!)
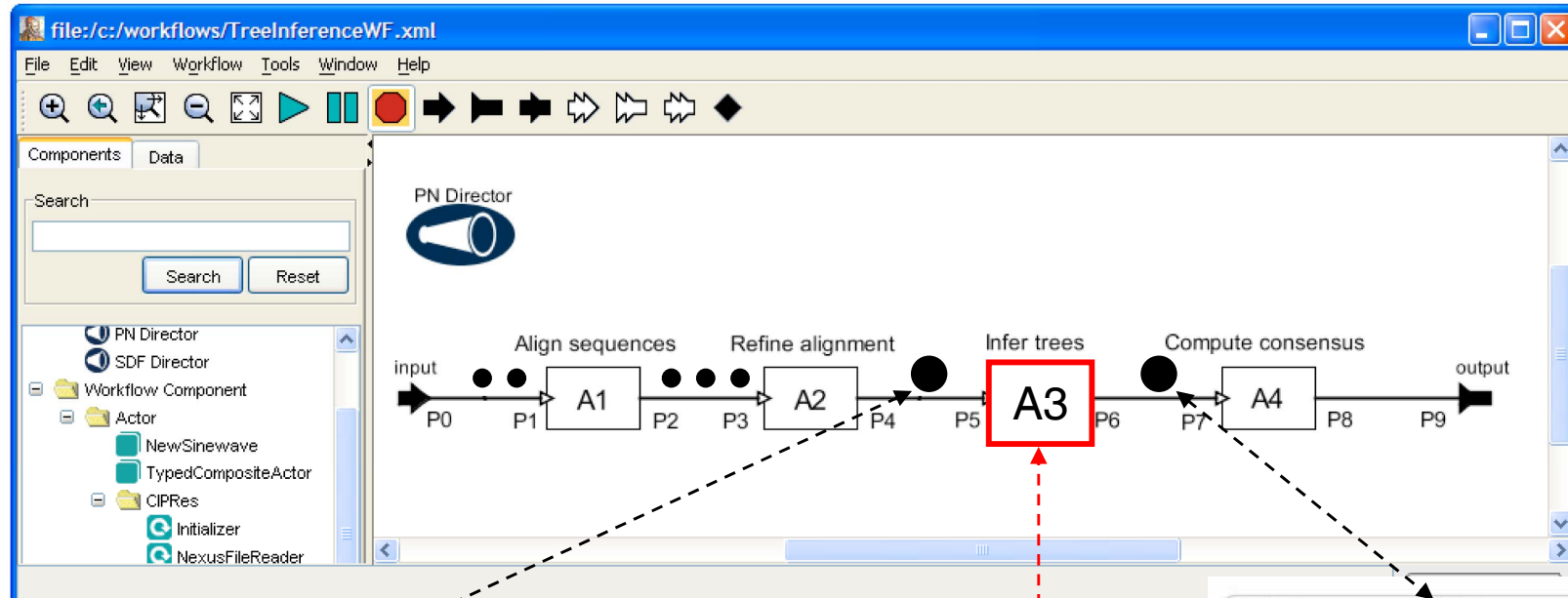
  + easy configuration, reuse, …
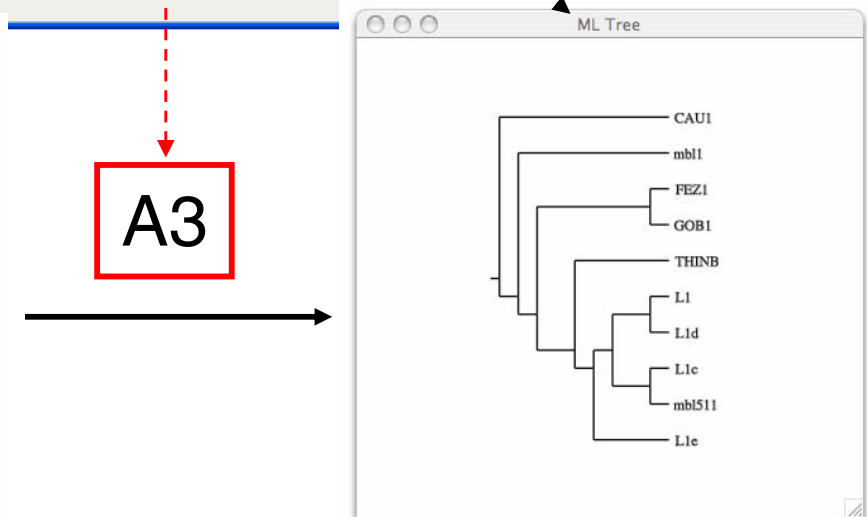
  + pipeline parallelism!



Reusable Actor "Class"

Specialized Actor "Instances"

Pipelined Execution Model

Inline Documentation

Inline Display

Easy-to-edit Parameter Settings

Checkpointing for (semi-smart) restart

**Overall architecture/simulation (physicist): Scott Klasky (ORNL)**
**Workflow design & development: Norbert Podhorszki (UC Davis)**

# Ex: Pipelined workflow for inferring phylogenetic trees



**Aligned sequences "token"**

**Phylogenetic tree "token"**

**Collection-Oriented Scientific Workflows,** Bowers, McPhillips, Ludäscher

UCDAVIS
UNIVERSITY OF CALIFORNIA
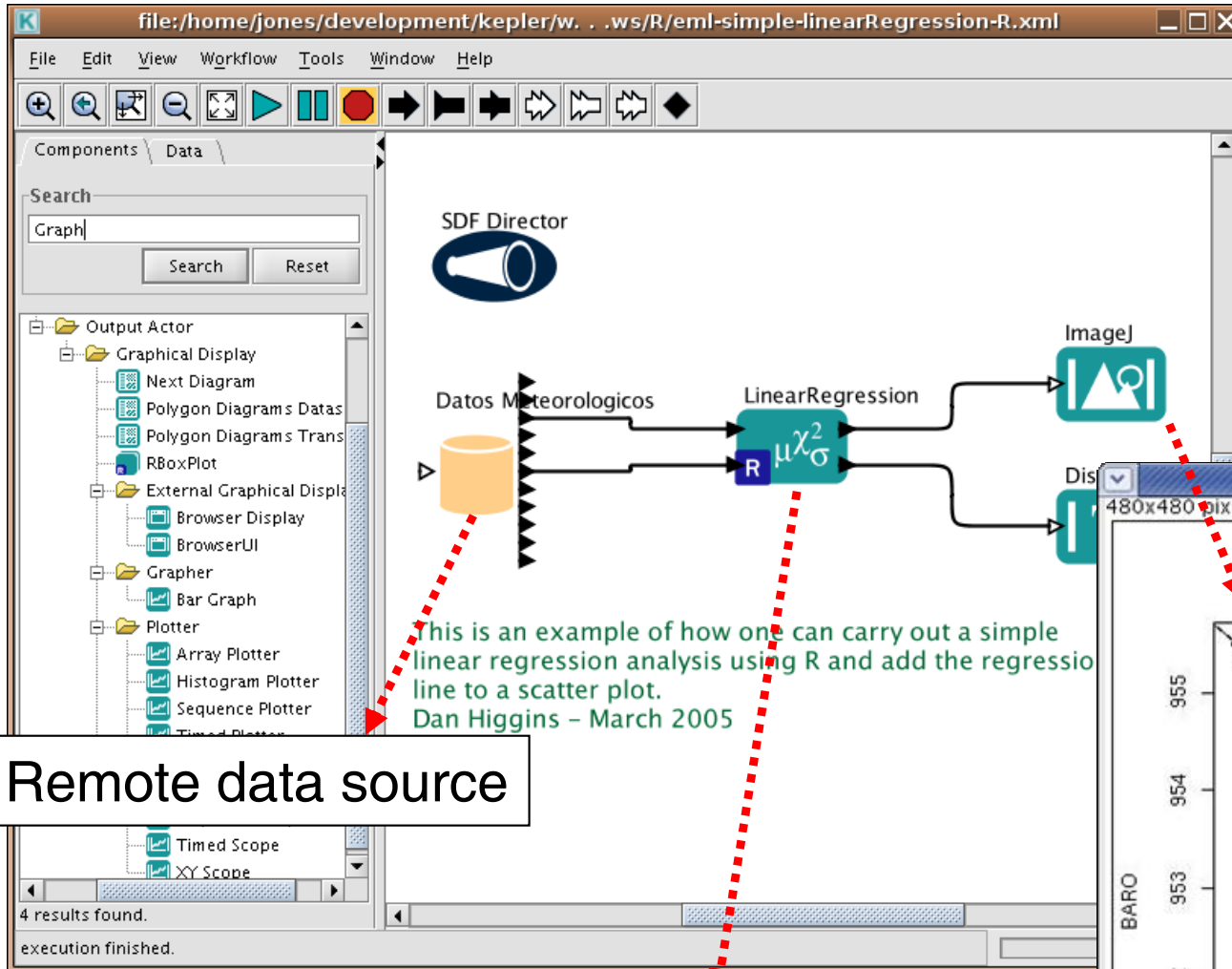
# What about scripts instead of SWFs?

- **Why not just use …**
  - MIRACLE-Perl®
  - or MIRACLE-Python®
  - or MIRACLE-BPEL4WS® ???

- **Perl/Python in the hand of a gifted** (and in the case of Perl: masochistic) **programmer are hard to beat …**

- **… but (MIRACLE-) Scientific Workflows offer some new features:**
  - parameter configuration, parameter studies
  - actor (component)-oriented workflow design  [Bowers-Ludaescher-ER'05]
    - component and workflow reuse & repurposing
  - semantic extensions (smart search/link/…)    [Bowers-Ludaescher-QLQP'06]
  - data (and workflow) provenance support   ([Altintas-*et-al*], [McPhillips-*et-al*]) @ IPAW'06
    - explain data dependencies/lineage, debug "strange" results, smart rerun,…
  - data-, task-, **pipeline**-parallelism
  - comprehensibility, documentation

- **e.g. check out demos by Kepler, Taverna, …**

# Kepler SWF using remote datasets, 3rd-party software ...
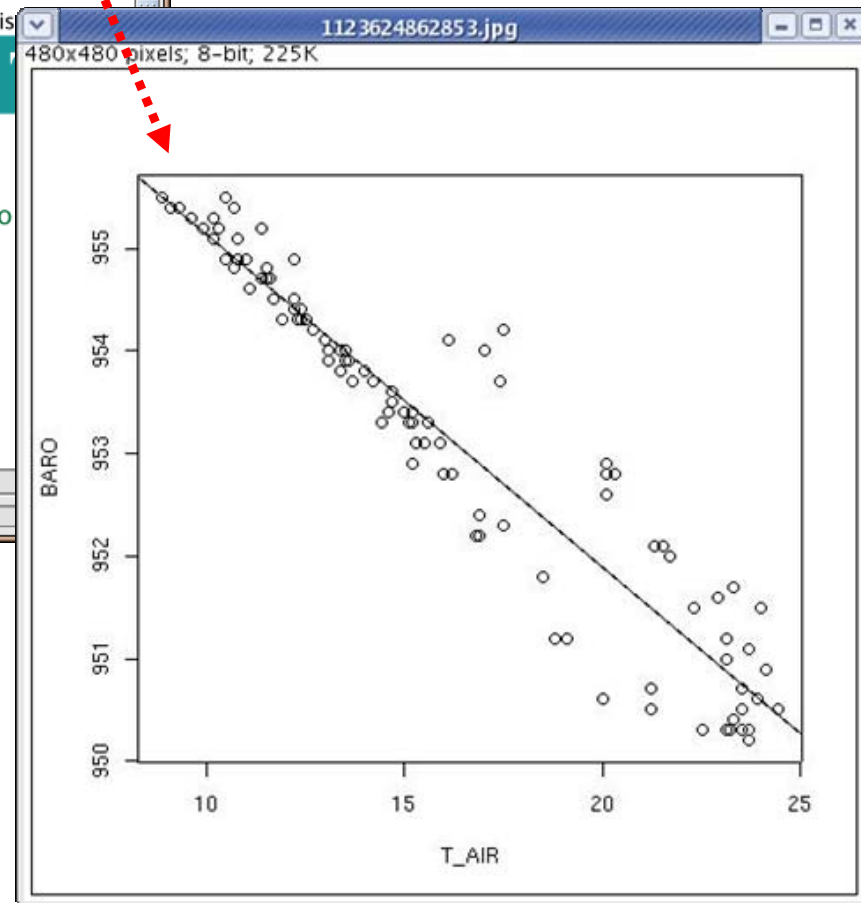


**Remote data source**

This is an example of how one can carry out a simple linear regression analysis using R and add the regression line to a scatter plot.
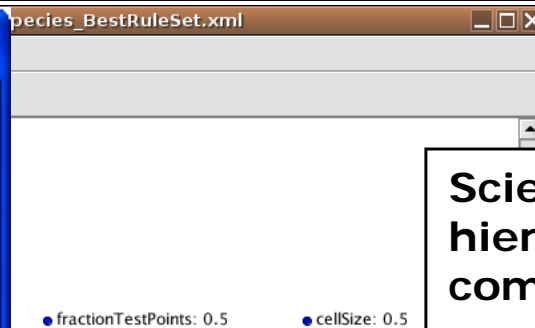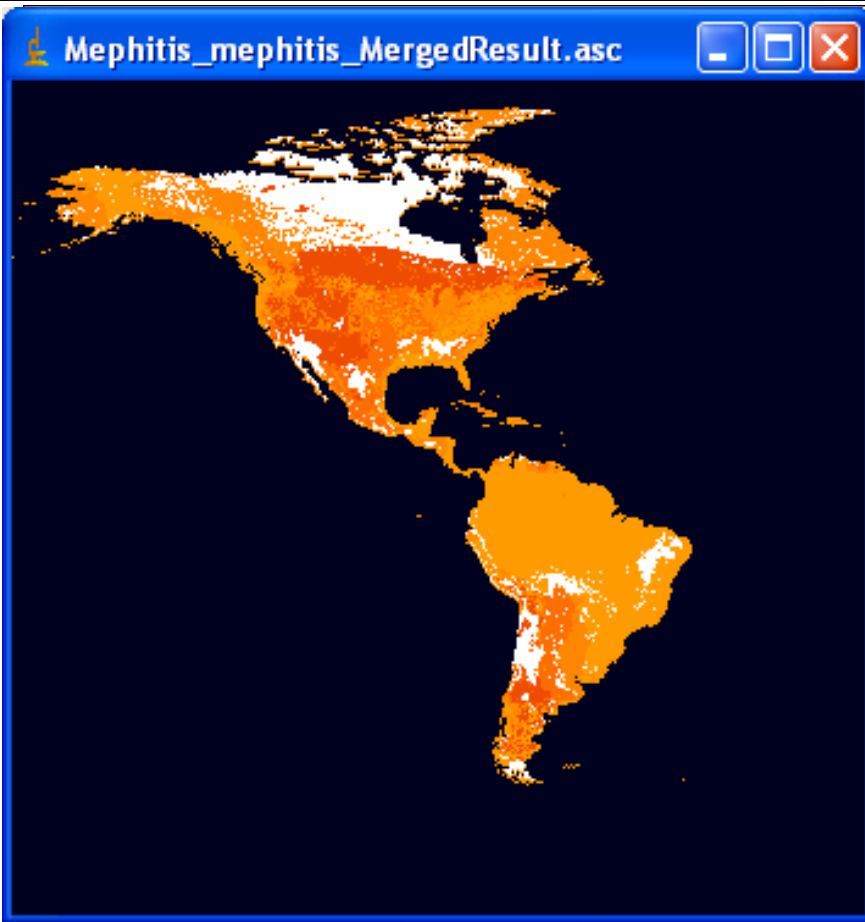Dan Higgins – March 2005

**R** processing script

```
res ← lm(BARO ~ T_AIR)
res
plot(T_AIR, BARO)
abline(res)
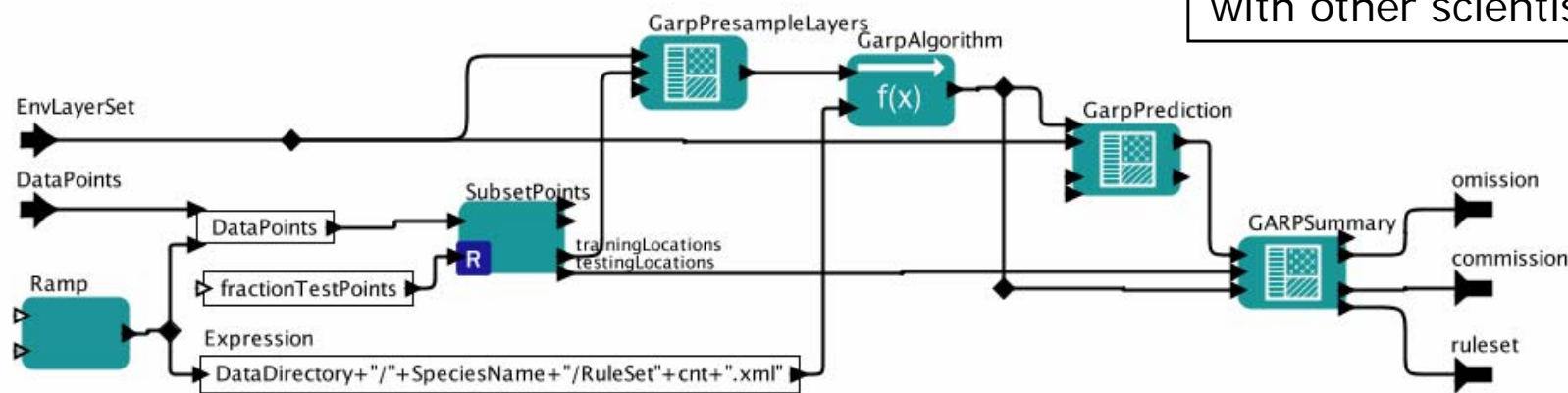```

# Managing complexity: Actor-oriented Modeling & Design

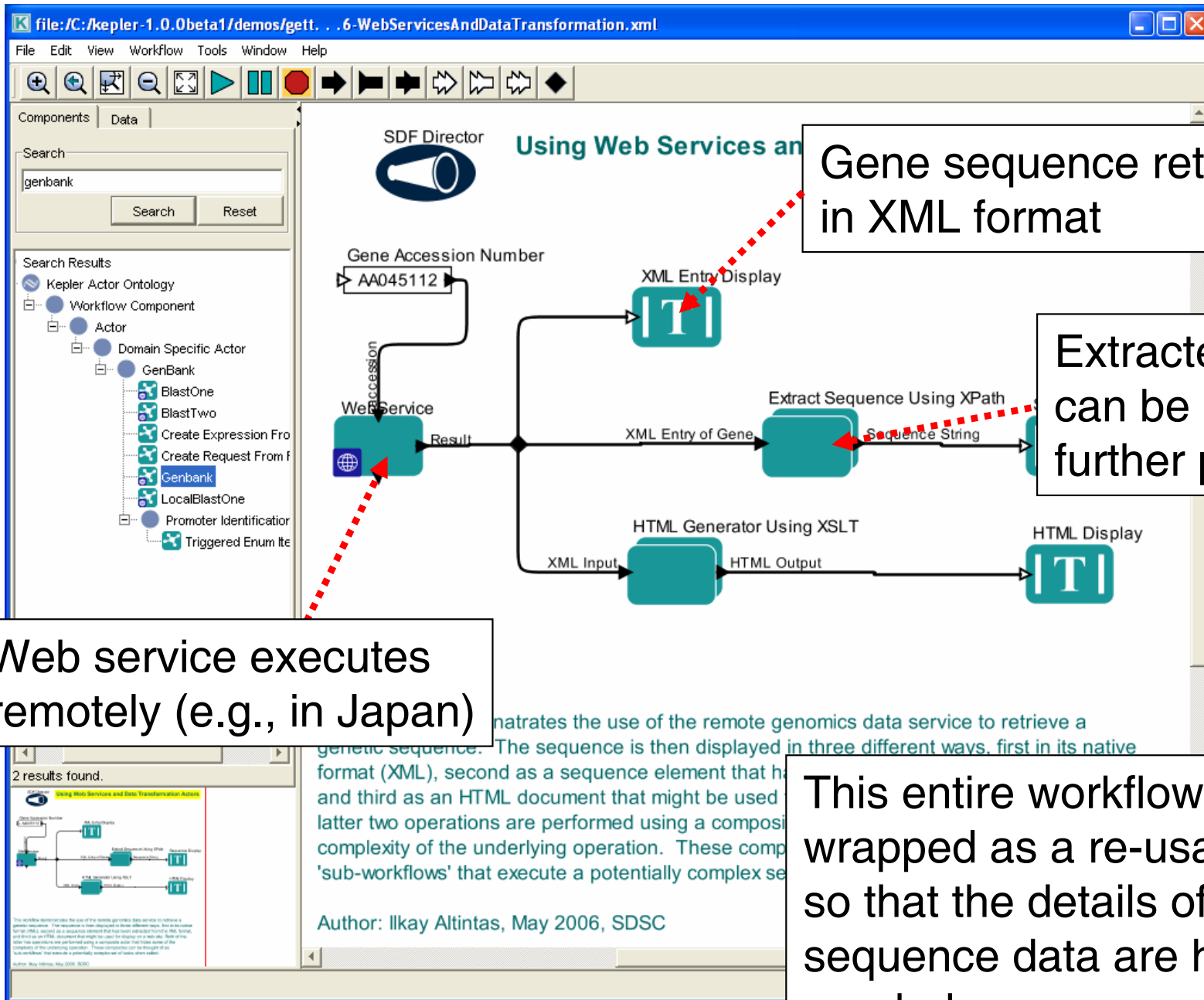**Scientific workflows use hierarchy to hide complexity:**

...rkflows can be a ...resentation of ...ocess that is easy ...d at a glance

... into sub- ...eals increasing

...models using hierarchy promotes the development of re-usable components that can be shared with other scientists

fractionTestPoints: 0.5          cellSize: 0.5

Mephitis_mephitis_MergedResult.asc

...pecies_BestRuleSet.xml

EnvLayerSet

DataPoints

GarpPresampleLayers

GarpAlgorithm

f(x)

GarpPrediction

SubsetPoints

DataPoints

R

trainingLocations
testingLocations

GARPSummary

omission

commission

Ramp

fractionTestPoints

Expression

DataDirectory+"/"+SpeciesName+"/RuleSet"+cnt+".xml"

ruleset

# Retrieving gene sequences via web services



Gene sequence returned in XML format

Extracted sequence can be returned for further processing

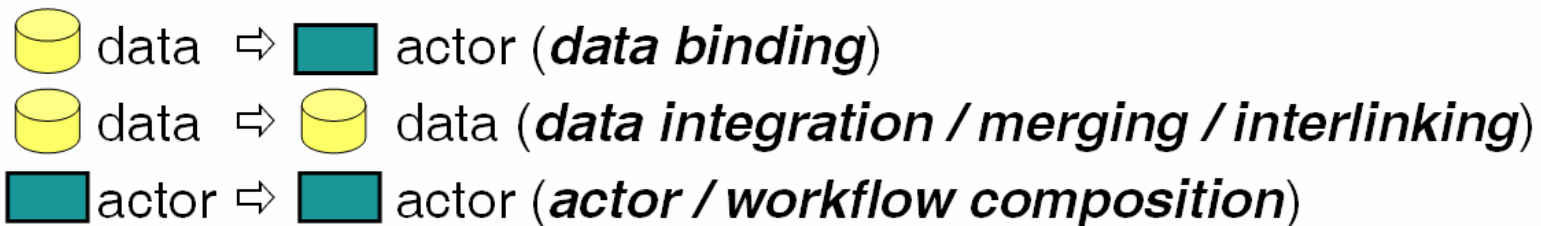Web service executes remotely (e.g., in Japan)

This entire workflow can be wrapped as a re-usable component so that the details of extracting sequence data are hidden unless needed.

# Kepler/SEEK Semantic Extensions

- **Employ semantic extensions (ontologies) for ..**



| | | |
|---|---|---|
| data ⇨ | actor (*data binding*) | |
| data ⇨ | data (*data integration / merging / interlinking*) | |
| actor ⇨ | actor (*actor / workflow composition*) | |

  – Smart **S**earch (➜ *Resource Discovery*)

  – Smart **A**ttach (➜ *Data Binding*)

  – Smart **I**ntegration (➜ *Merge Actor*)

  – Smart **L**inks (➜ *Actor Composition*)



*S*earch-*A*ttach-*I*ntegrate-*L*ink

# "Hybrid" Types ... Semantic + Structural Typing
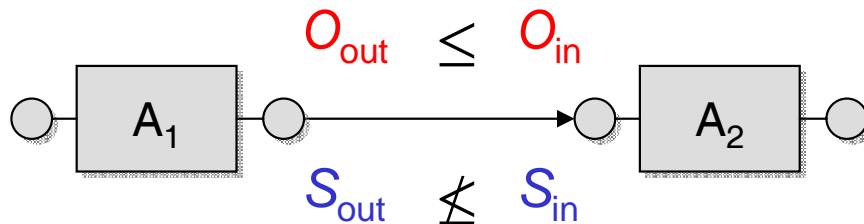
$O$

$S$

$O$

$S$

$O_{out}$

$S_{out}$

$O$ : Observation $\sqcap$
$\quad$ $\forall$obsProperty.SpeciesOccurrence
$S$ : SpeciesData(site, day, spp, occ)

**Structural Types: Given a structural type language** $S$

– Datasets, inputs, and outputs can be assigned structural types $S \in \mathbf{S}$

**Semantic Types: Given an ontology language** $O$ **(e.g., OWL-DL)**

– Datasets, inputs, and outputs can be assigned ontology types $O \in \mathbf{O}$

$O_{out} \leq O_{in}$

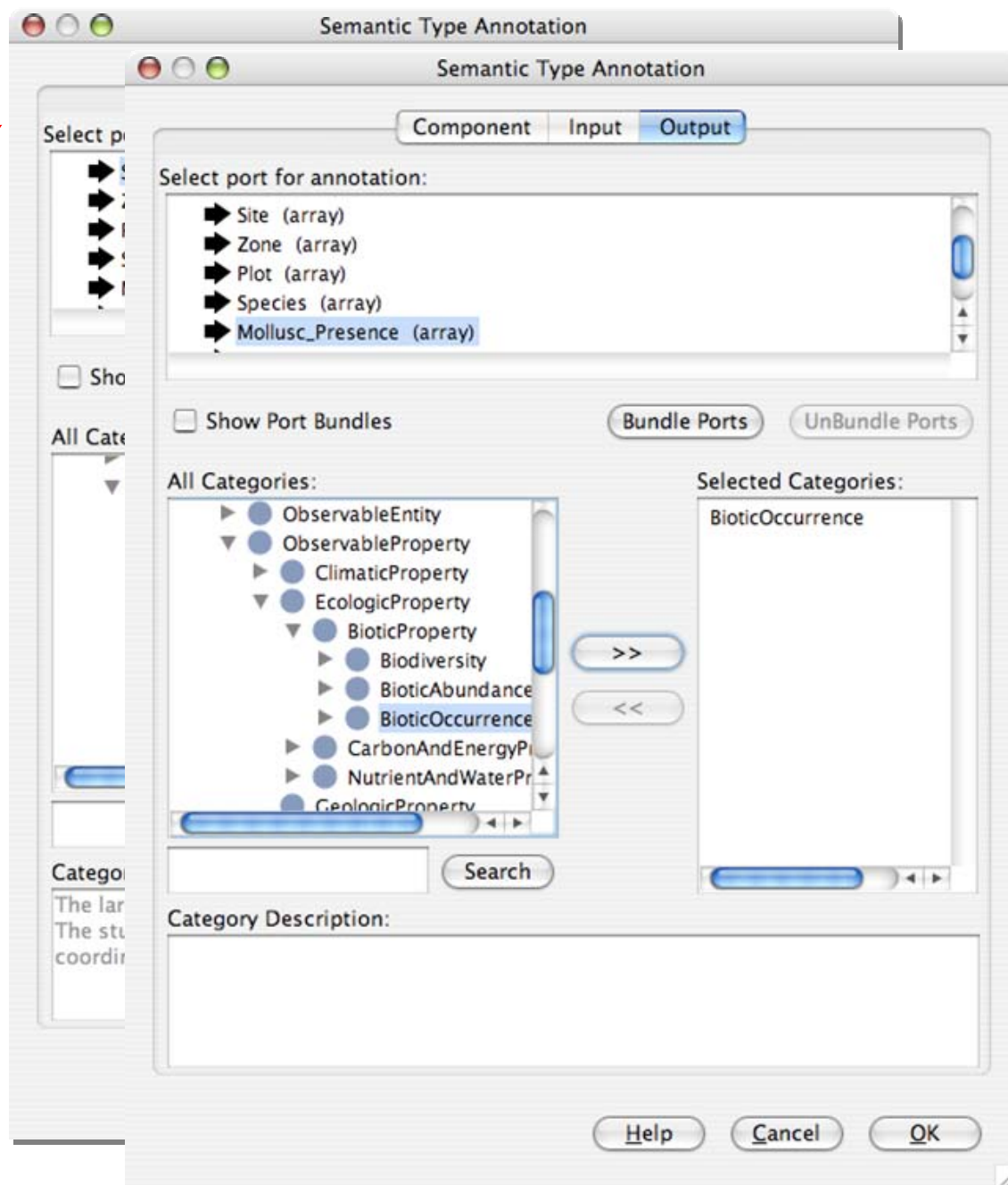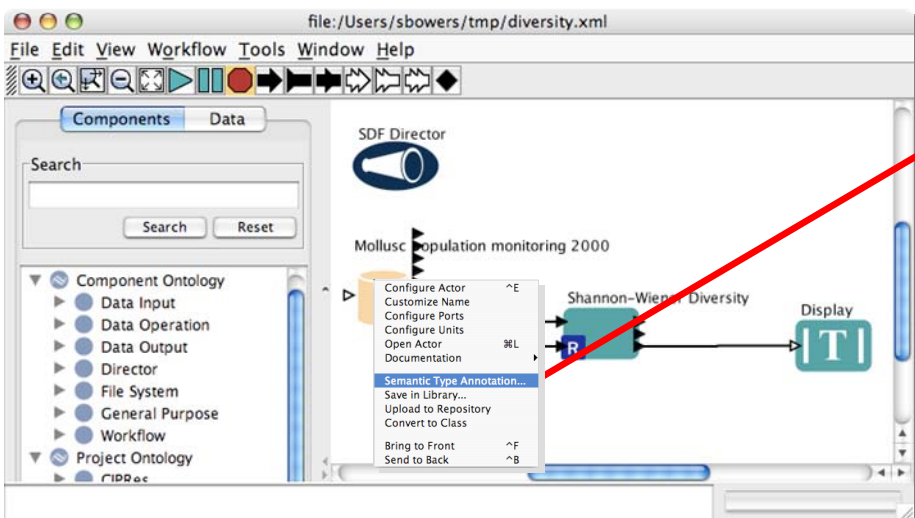$A_1$

$A_2$

$S_{out} \not\leq S_{in}$

Semantically compatible
but structurally incompatible

Semantic & structural types can be **combined using logic constraints**

$$\alpha := (\forall site, day, sp, occ) \ \text{SpeciesData}(site,\ day,\ sp,\ occ) \rightarrow$$
$$(\exists y) \ \text{Observation}(y),\ \text{obsProp}(y,\ occ),\ \text{SpeciesOccurrence}(occ)$$

# Semantic Type Annotation in Kepler



- **Component input and output port annotation**
  - Each port can be annotated with multiple classes from multiple ontologies
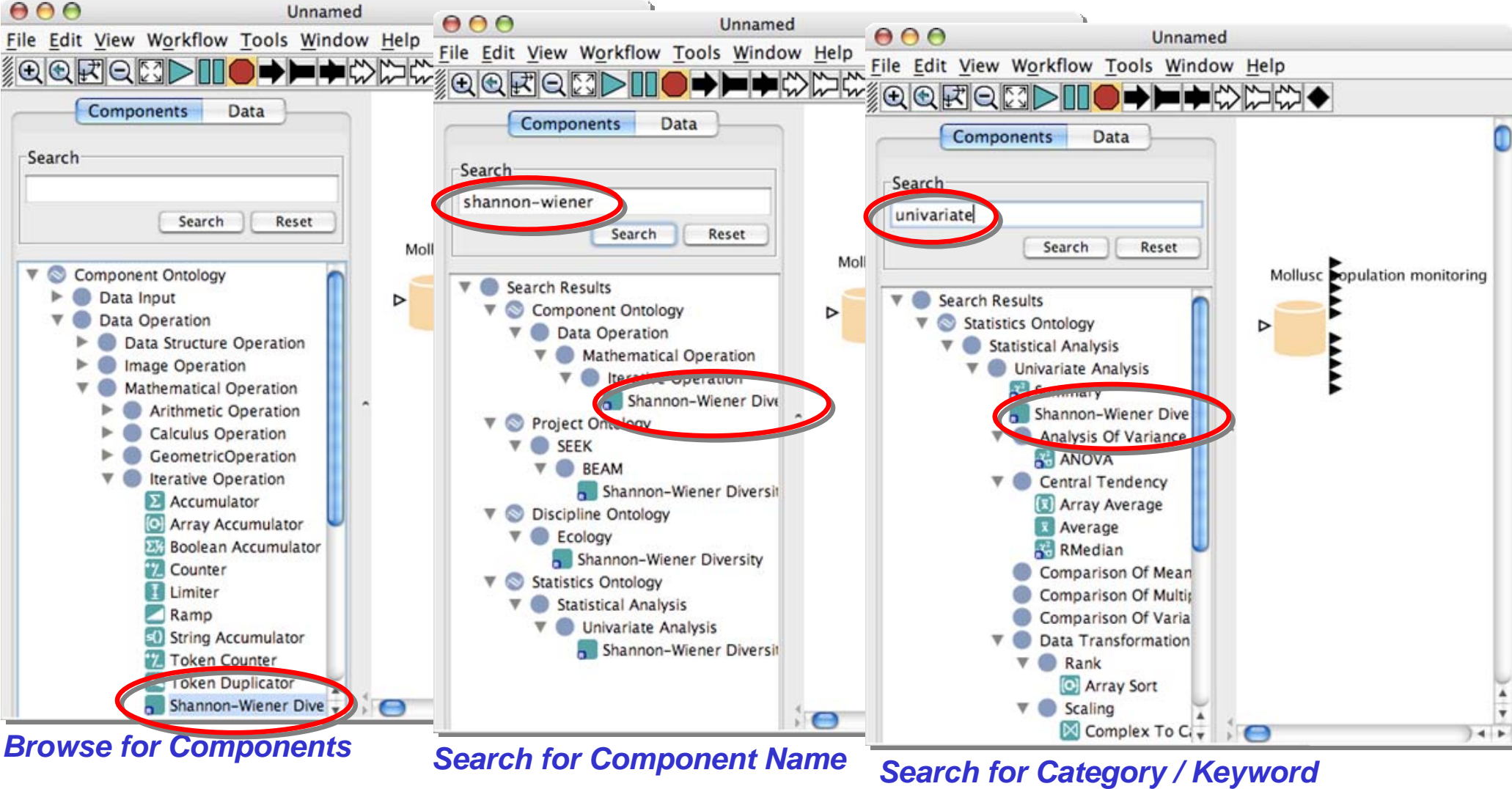  - Annotations are stored within the component metadata

# Component Annotation and Indexing



- **Component Annotations**
  - New components can be annotated and indexed into the component library (e.g., specializing generic actors)
  - Existing components can also be revised, annotated, and indexed (hiding previous versions)
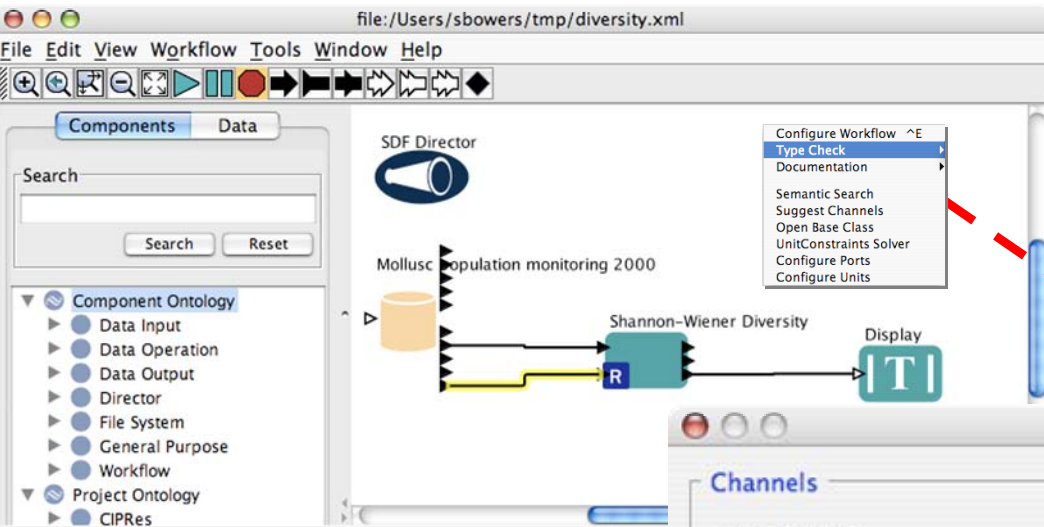
# Smart Search



**Browse for Components**

**Search for Component Name**

**Search for Category / Keyword**

Find a component (here: an actor) **in different locations ("categories")**

- … based on the **semantic annotation** of the component (or its ports)
- … registered to one or more ontologies (controlled vocabularies)

DILS06 **Collection-Oriented Scientific Workflows,** Bowers, McPhillips, Ludäscher

UCDAVIS
UNIVERSITY OF CALIFORNIA

# Smart Linking (Workflow Design)



- **Statically perform semantic and structural type checking**

- **Navigate errors and warnings within the workflow**

  - Search for and insert "adapters" to fix (structural and semantic) errors …

# Some KEPLER Actors ... (oh, the good old days ... )

UCDAVIS
UNIVERSITY OF CALIFORNIA

# So ...

- **… a question: If scientific workflows are so great, why haven't they taken over the world??**
  - A1: just wait …
  - A2: they already have …
  - A3: The problem of creating flexible, reusable, comprehensible, efficient, … workflows
  - … is akin to the problem of creating modular, reusable, maintainable, … software!
  - … it's complex systems engineering (as in: **difficult**)
    - … and using UML, XML, WS-foo, SOA-bar, and BPEL-baz are no substitute for solving your design problem!

# Complexity in Scientific Workflow Design

## The use of "control-flow" primitives

– Managing complex, **nested** data structures (select/filter/transform)

– Fault-tolerance and exception handling

– …



*Custom actors, hand-crafted control flow limited to sequential execution (SSDBM'03)*

*Fault-tolerance control-flow **"wired-in"**, e.g. via Boolean switches, **complex branching and looping***

**UCDAVIS**
UNIVERSITY OF CALIFORNIA

# Modeling Control-Flow Constructs in Dataflow

- **Dataflow** in Kepler
  - Based on **dataflow process networks** (Kahn et al, Lee et al)
  - Supports **pipeline parallelism** (streaming data)
  - Natural paradigm for **data-driven** workflows
  - Efficient analysis and scheduling
  - **Intuitive model** for workflow designers

  
  Dataflow&Ctrl-flow Marriage!

- **Control-Flow** in Kepler
  - Branching via if-then-else and switch-case statements
  - Iteration with multiple entry and exit points
  - Low-level actors for manipulating structure (e.g., record-to-array)

- **Problems modeling Control-Flow directly using Dataflow**
  - Overly complicated workflows; hard to understand (low-level programming), maintain, debug, extend ➔ limited reusability; complex re-configuration



*Note to self: Check out the new Taverna control-flow layered architecture !*

[**Bowers**-*et-al*, **SciFlow'06**]

*New Design Primitives (e.g., templates & frames)*

UC DAVIS
UNIVERSITY OF CALIFORNIA

# Scientific Workflow Design: Challenges



"And that's why our scientific workflows are much **easier to develop, understand and maintain!**"

# A Simple Motivating Example



- **Take the services (actors, components) in (a)**
- **… and chain them together in a scientist friendly form a la (b)**
- **… considering the following signatures (cf. Haskell, ML, … )**
  - (c) BLAST :: DNA→ [DNA]
  - (d) MotifSearch :: DNA → [Motif]
  - (e) BLAST o  MotifSearch   = \x. BLAST(MotifSearch)(x)
- **… oops: (e) is <u>not type correct</u>: note the signatures of (c) and (d)!**
- **a neat solution: implicit or explicit iteration / map(f)[x1,…,xn]**
  - cf. Kepler and Taverna, Kepler solutions

# Extended Example: Workflow Evolution



Conventional workflows

- **(a) => (b): replace *A::a→b* with *A'::a→[b]***
  - need to call B iteratively i.e. wrap B inside a component or add control-flow
- **(b) => (c): upstream produces *[a], [a], …* instead of *a, a, …***
- **(d): need to "bypass" data components since *B* can't handle *d*s**
- **This gets messy quickly …**

# A Realistic Example (ChIP-chip workflow)

# So how to get from messy to clean & reusable designs?



*before…*

*AFTER!*

UCDAVIS
UNIVERSITY OF CALIFORNIA

# Answer: Collection-Oriented Modeling & Design

- ## Collection-Oriented Modeling & Design (COMAD)
  - starting point: dataflow / actor-oriented modeling & design
  - embrace the **assembly line metaphor fully**
    - ➔ Flow-based Programming (J. Morrison)
  - data = **tagged nested collections**
    - e.g. represented as flattened, **pipelined** (XML) **token streams:**



➔ Multi-level Pipeline Parallelism!

UCDAVIS
UNIVERSITY OF CALIFORNIA

# Multi-level Pipeline Parallelism



U.S. Department of the Interior, National Park Service, Edison National Historic Site

Metadata for collection **b**

Actors 4 and 5 are processing contents of collection **a concurrently**
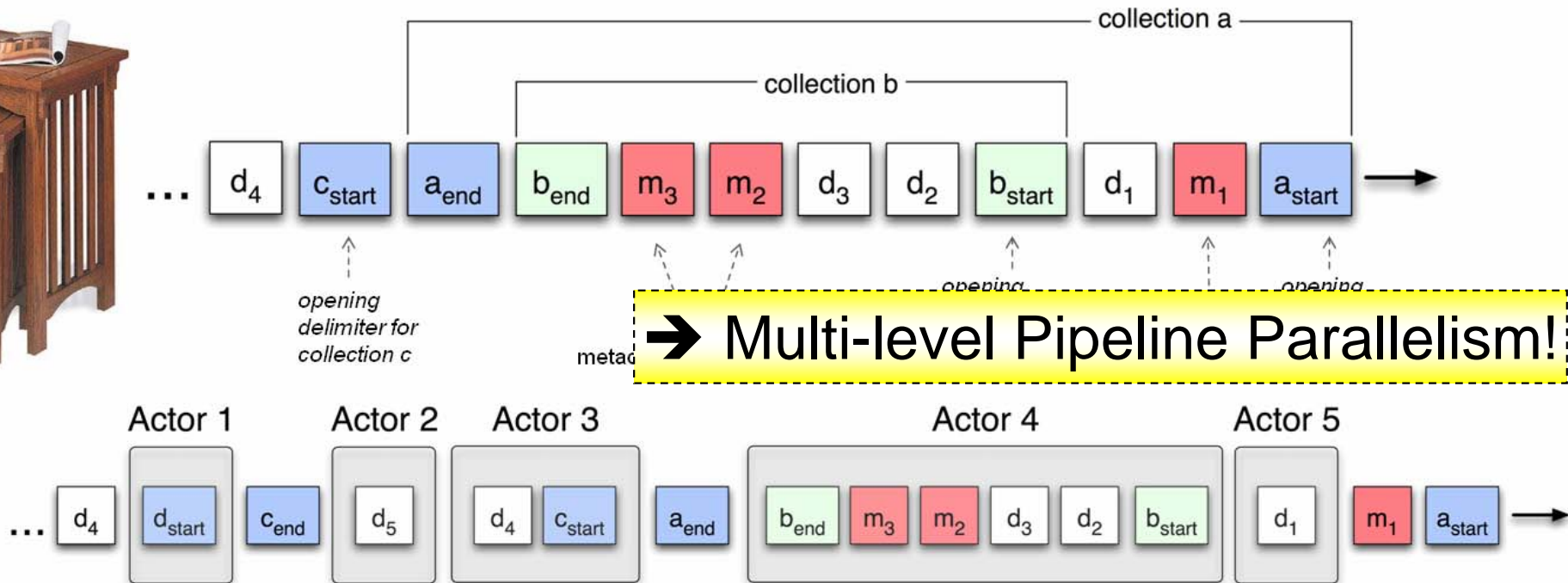
Actor 1 — $d_{start}$ | $c_{end}$

Actor 2 — $d_5$

Actor 3 — $d_4$ | $c_{start}$

$a_{end}$

Actor 4 — $b_{end}$ | $m_3$ | $m_2$ | $d_3$ | $d_2$ | $b_{start}$

Actor 5 — $d_1$

$m_1$ | $a_{start}$

... $d_4$

*Closing delimiter for top-level collection* **c**

*Opening delimiter for top-level collection* **c**

*Actor 4 processes* **entire collections** *(of a particular type) at one time.*

*Actor 5 processes* **one data token** *at a time.*

McPhillips, T., Bowers, S., An Approach for Pipelining Nested Collections in Scientific Workflows. *SIGMOD Record*, 34(3), 2005.

UCDAVIS
UNIVERSITY OF CALIFORNIA

# How does COMAD work?

- **Some COMAD principles:**
  - data = tagged, flattened, nested collections (token streams)
    - data tokens
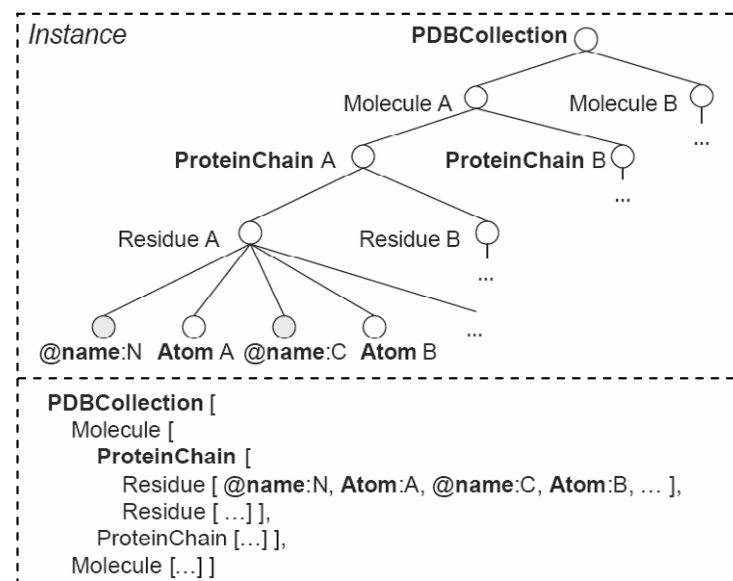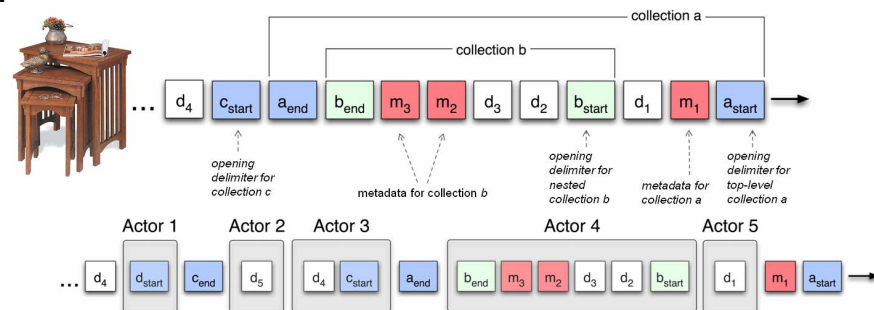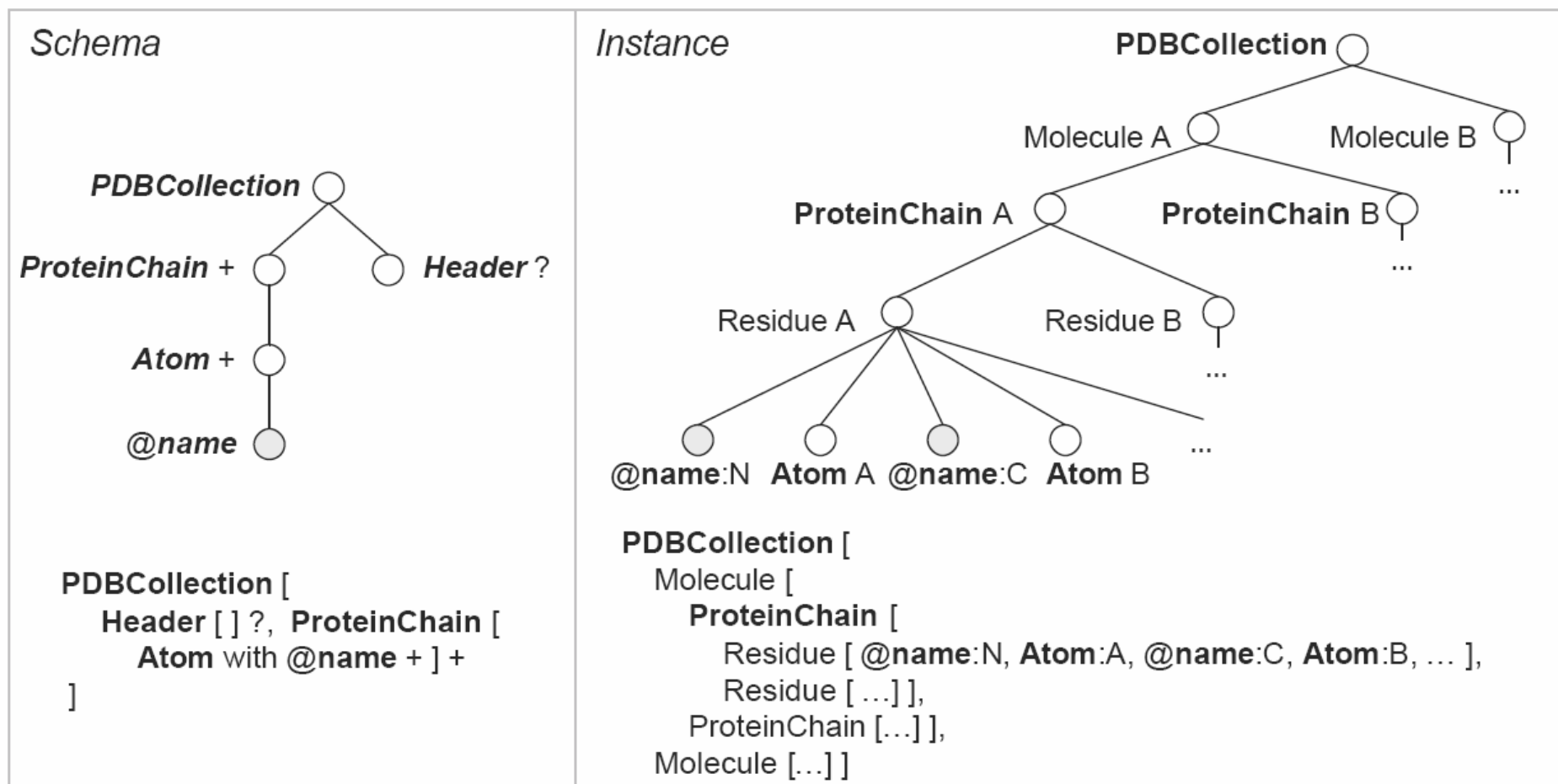    - metadata tokens
      - inherited downwards into (sub)collections
  - define an actor's **read scope** via an (X)Path-like expression:
    - default actor behavior:
      - not mine?
        - ➔ don't do anything: **just pass the buck!**
      - stuff within my scope? ➔
        - » add-only to it (default)
        - » consume scope; write-out result
        - (but remember the bypass!)
  - **iteration scope** is a query involving group-by and further refines the granularity/subtrees that constitute the tokens consumed by an actor firing
  - has aspects of implicit iteration (a la Taverna)
    - default iteration level to fix signature mismatches
  - but also:
    - granularity/grouping is definable
    - works on "anything" (assuming scope is matched correctly)
    - add-only and replace modes

# Read Scoping with Collection Schemas (Types)



Collection schema
(defining a query pattern)

… and a sample matching instance

# Iteration Scoping via Queries

The ComposeNexus actor converts an optional character matrix, weight vector, and a list of zero-or-more trees into a Nexus file. Note here that the actor is invoked exactly once for each input Nexus collection, unlike the TreeReporter actor, which is invoked once per tree. This invocation pattern is described by the following iteration scope.

ComposeNexus iteration-scope ($c, $v, collect($t in $n)) :=
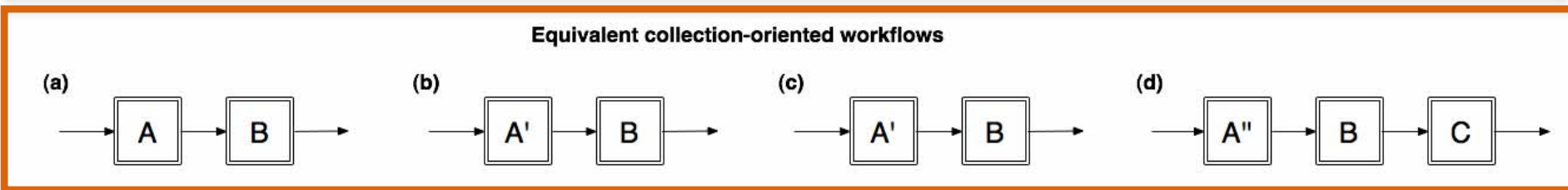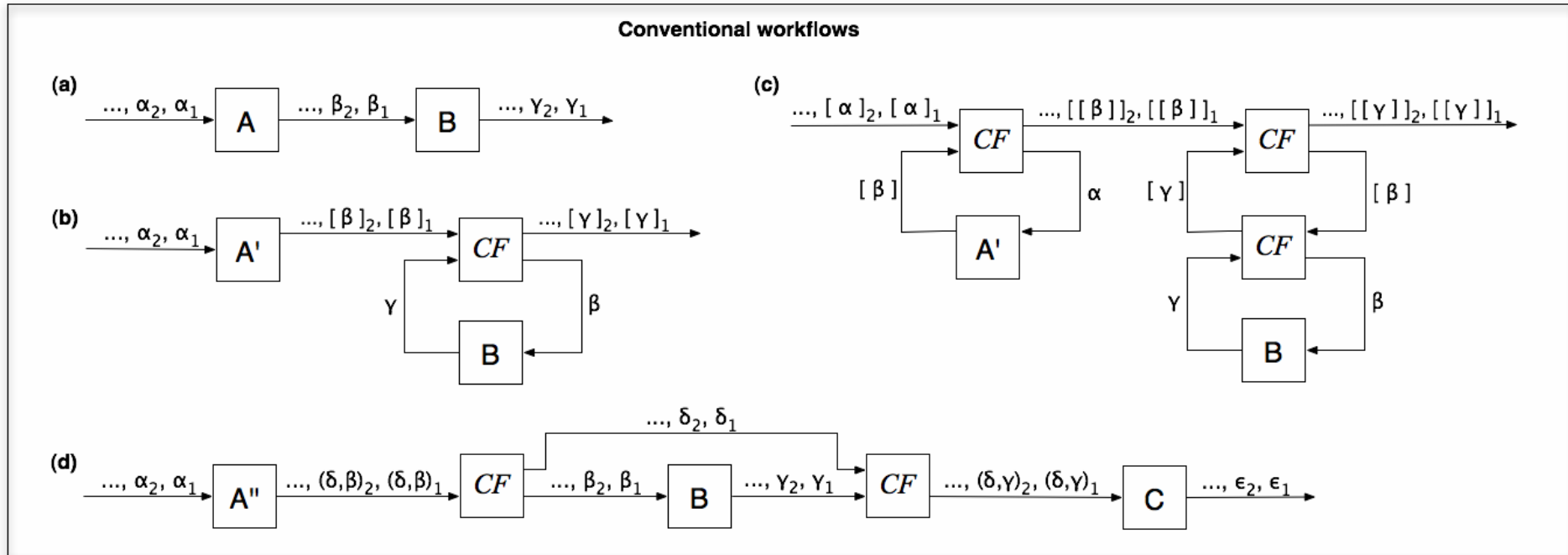    Nexus{$n}[ CharacterMatrix {$c},  WeightVector {$v}, Tree {$t} ]

The collect expression constructs a list of trees, where each tree is contained in the given Nexus collection. Every collect expression in an iteration scope consists of a data or metadata variable (in this case $t) combined with a collection variable (in this case $n).

- **ingredients / related:**
  - tree pattern queries / XML selection queries
  - … with group-by (here "group-by $c, $v")
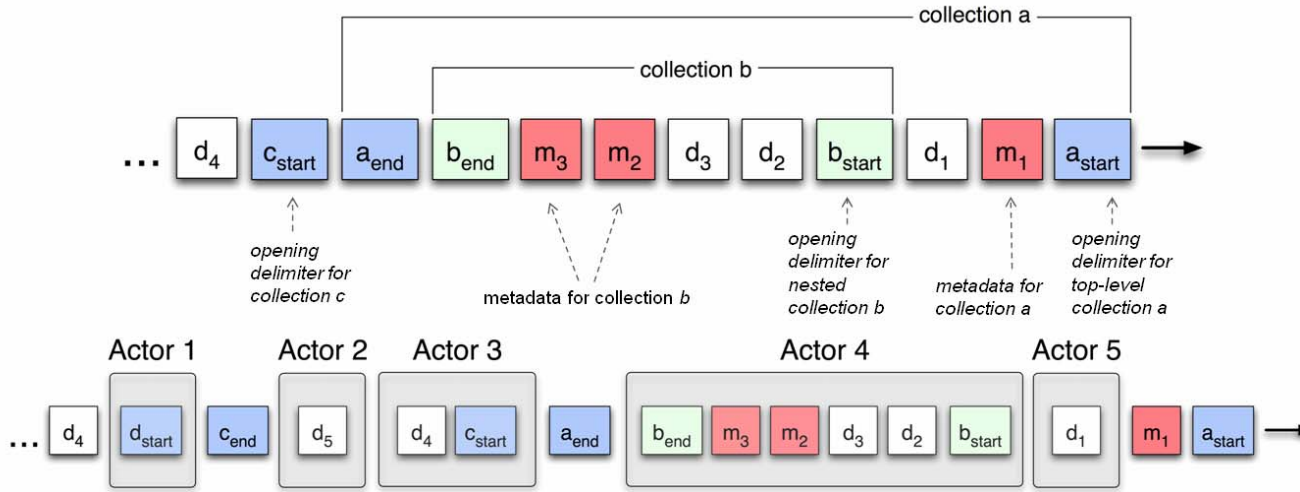  - see also: list comprehensions as queries (cf. CPL/Biokleisli)

# COMAD: What we gained
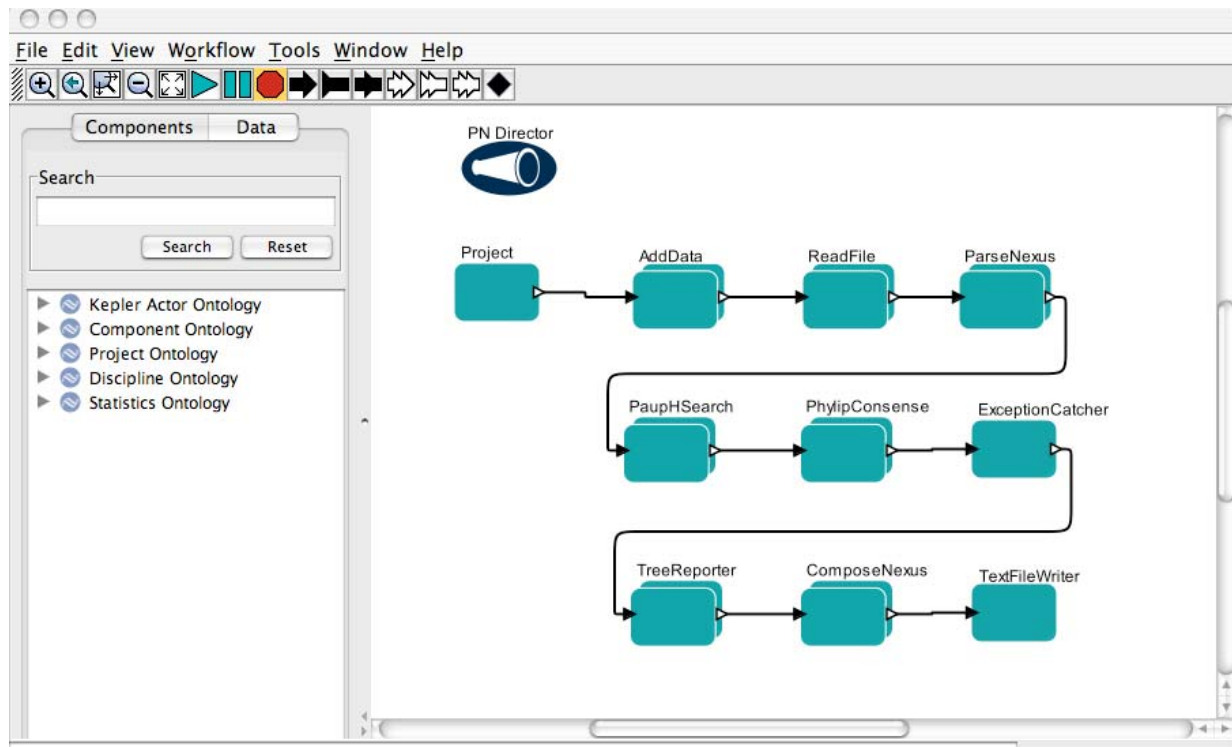


Conventional workflows

Equivalent collection-oriented workflows

- **from fragile, messy workflow designs**
- **… to more reusable actors**
  - just change the scopes
  - sometimes not even that is needed
- **… and cleaner workflow design**
- **Crux: keep the nesting structure of data (pass through, add-only)**
- **… and let it drive the (semi-)implicit iteration**
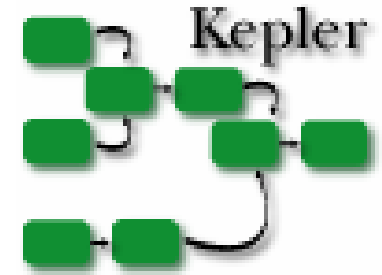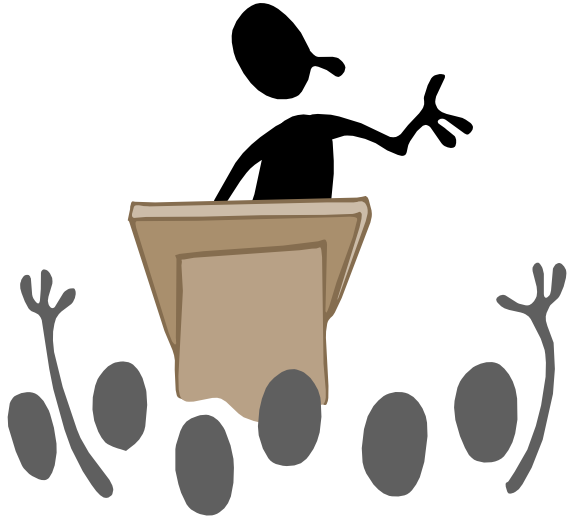
# Summary: Collection-Oriented WF Modeling & Design



- **Assembly line metaphor: Dataflow + XML …**
  - Streams are nested collections ($\approx$ XML)
  - Less "messy" WFs (more linear, less branching)
  - Pipelined parallelism (stream lists)

**Collection-Oriented Scientific Workflows,** Bowers, McPhillips, Ludäscher

UCDAVIS
UNIVERSITY OF CALIFORNIA

# Acknowledgements and Q&A ...

- **NSF/ITR Science Environment for Ecological Knowledge (SEEK)**
- **NSF/ITR Geosciences Network (GEON)**
- **DOE/SciDAC Scientific Data Management Center (SDM)**
- **U.S. Dept. of Energy, LLNL**