

X-CSR: Dataflow Optimization for Distributed XML Process Pipelines

Daniel Zinn,
Shawn Bowers, Timothy McPhillips, Bertram Ludäscher

ICDE 2009

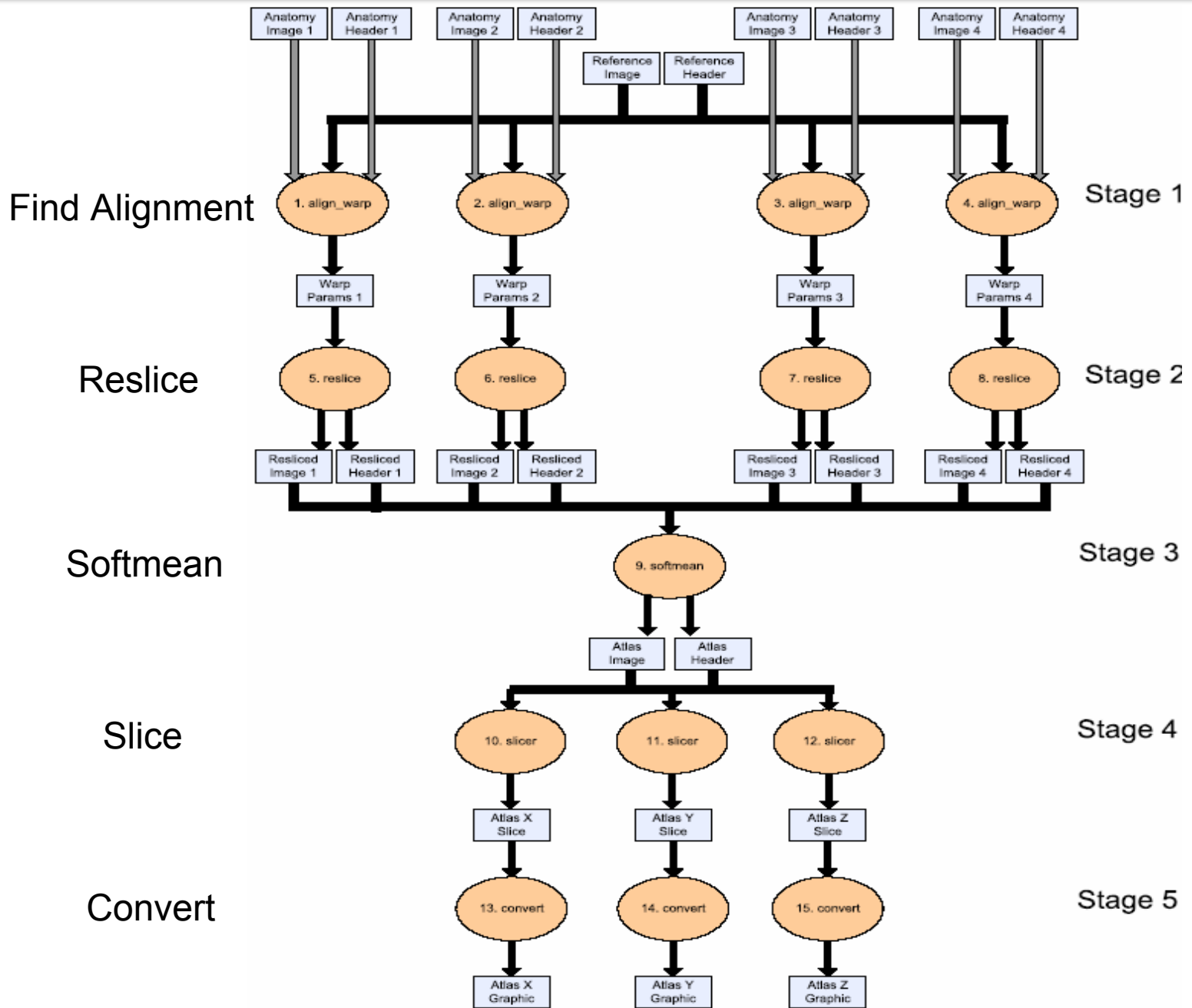
Department of Computer Science
University of California at Davis

April 1st, 2009

- **XML Processing Pipelines (Δ -XML)**
- **X-CSR: Dataflow-Optimization for Δ -XML**

Motivation: Scientific Workflows

- Phylogenetics
 - Gather DNA sequence data
 - Align sequences
 - Reconstruct phylogeny using multiple alignments – often obtaining a large number of trees
 - Compute consensus (or otherwise estimate the reliable components of the evolutionary history)
 - Perform post-tree analyses
- **Functional Magnetic Resonance Imaging (fMRI)**
 - Brain scans while performing activities
 - Voxel-picture post-processing (command-line tools)



Motivation for XML Processing Pipelines

Combining existing software to build complex scientific analysis systems while saving machine cycles and human cycles

- To save developer's time, systems should be...
 - Easy to build
 - Easy to evolve and maintain
 - Easy to share and re-use
 - Self-documenting
- To save execution time,
 - Systems should efficiently execute

Setting Characteristics

Software components are heterogeneous:

- Languages: Matlab, C, Java, Pascal, R, ...
- Invocation: Portals, Web-Services, local executables, libraries (BioPerl)

Process is exploratory:

- Refining steps and their parameter after new insights

Workflows are data- and/or compute-intensive:

- Phylogenetics (NP-hard problems, exponential convergence, ...)
- Astronomy (LSST project: nightly data of 15TB)

Current Approaches

Manual

With
Scripting
Languages

With
Scientific
Workflow
System

Δ -XML

XML Processing Pipelines

- Design, evolution, validation
- Optimization
- Automatic provenance

Approach: Scientific Workflow Systems

The image shows a screenshot of a scientific workflow system interface. The main window displays a workflow diagram with several steps: ReadFastaFile, CipresClustal, Gblocks, ComposeMatrix, CipresRAxML, PhylipDrawgram, WriteNexusFile, and SaveTrace. A yellow callout bubble labeled "Run analysis" points to the workflow execution controls at the top. Another yellow callout bubble labeled "Algorithms and data sources" points to the workflow steps. A third yellow callout bubble labeled "Analysis pipeline" points to the CipresRAxML step. A fourth yellow callout bubble labeled "Provenance information" points to the "Traces" panel on the left, which lists various workflow traces like "Crazy-loop.trace" and "Clustal_Gblocks_RAxML". The interface also shows a "Workflows" panel on the left with a tree view of workflow files, and a "Design" tab at the top. The title bar indicates the file path: "file:/Users/sbowers/cvs/daks/dev/kepl...leWorkflows/Clustal_Gblocks_RAxML.moml".

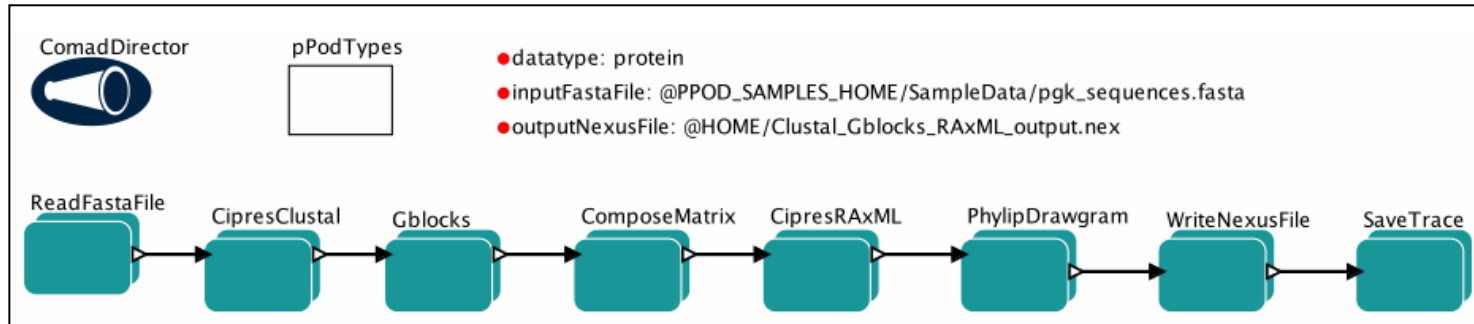
Run analysis

Algorithms and data sources

Analysis pipeline

Provenance information

Workflow 101



Actors

- Components with ports

Channels with Types

- Explicit communication between actors

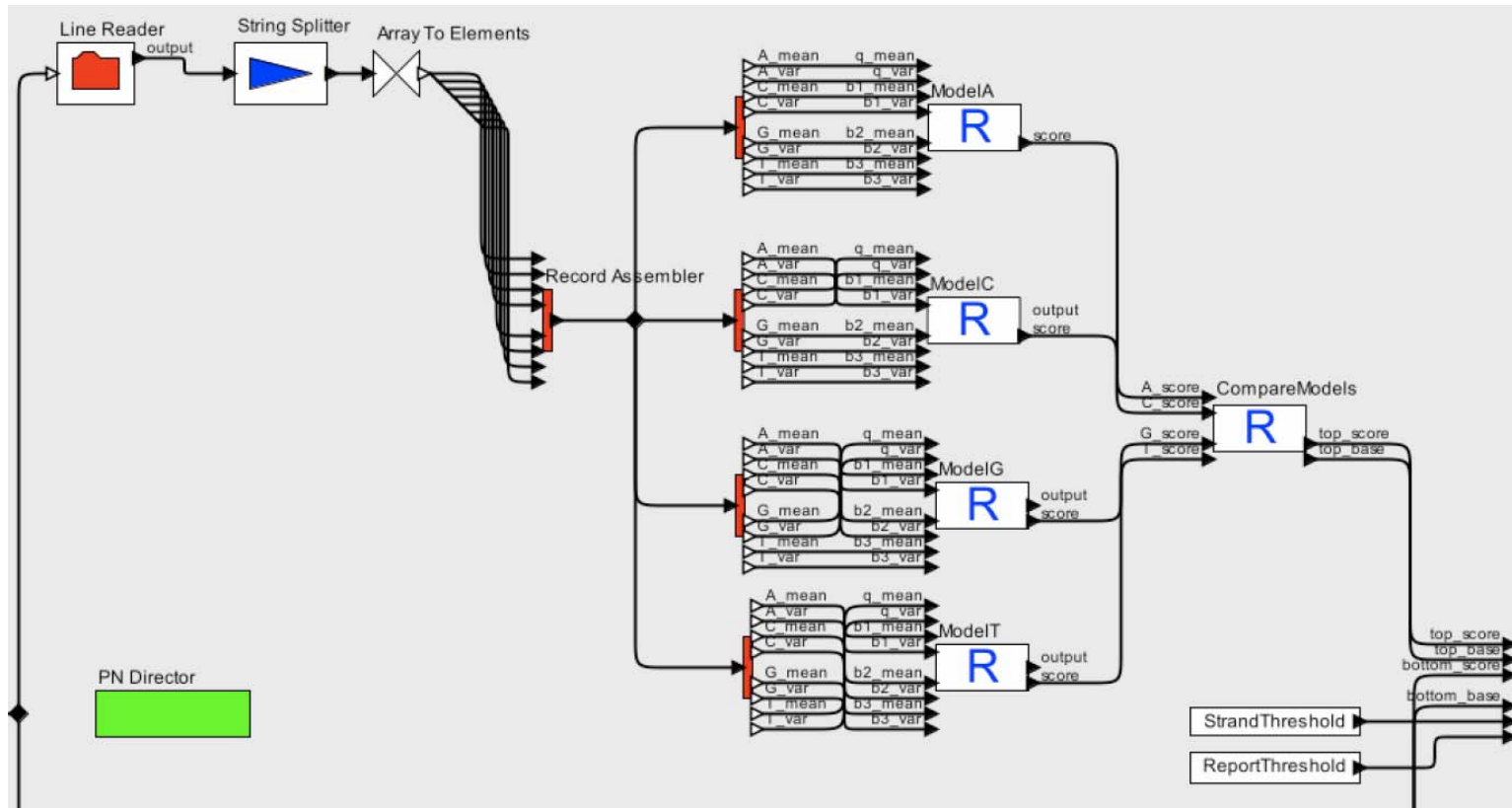
Tokens

- Data entities sent on channels

Model of Computation (MoC)

- Determines how workflow graph is executed

But...Problems with Nested Data



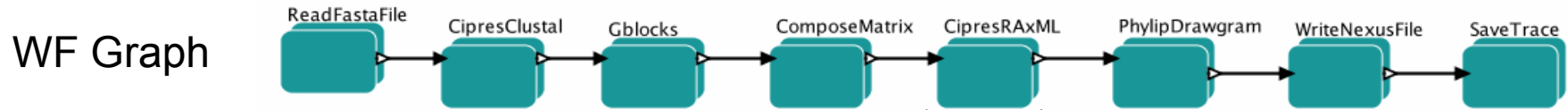
Problems

- Problems:
 - a) **Data associations**: workflow graph exposes details of low-level data “massaging”
 - b) **Complex-control flow**: workflow wiring can be overly complex to deal with exceptions, loops, etc.
- → Overly complex workflow graphs
 - Tedious to develop
 - Difficult to maintain and evolve
 - Data massaging encoded into workflow graph
 - Impossible to feed with structurally different input data

Solutions

- Ideally, switch to a paradigm / MoC that ...
 - a) automatically handles complex data
 - b) simplifies control-flow specification
- Solutions:
 - a) Nested collections (e.g. a la XML)
 - b) Virtual Assembly-Lines with actor configurations
 - Mostly linear sequence of actors
 - Every actor sees everything but only works on selected data
 - ➔ Exception handling e.g. via “tagging” faulty parts that will be removed or replaced, handled downstream (or simply ignored)
 - ➔ Increased flexibility for adding, removing, swapping actors

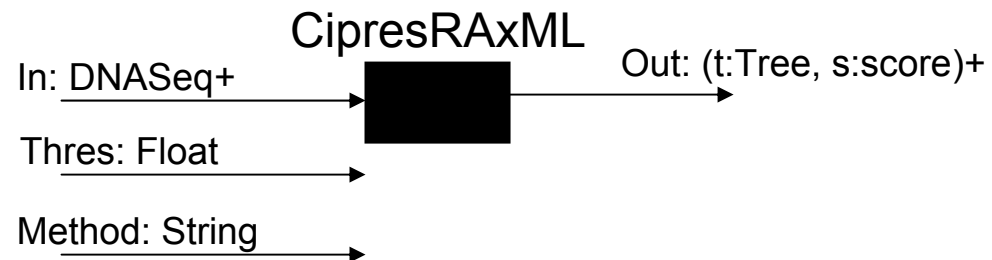
3-Layer Architecture of Δ -XML Pipelines



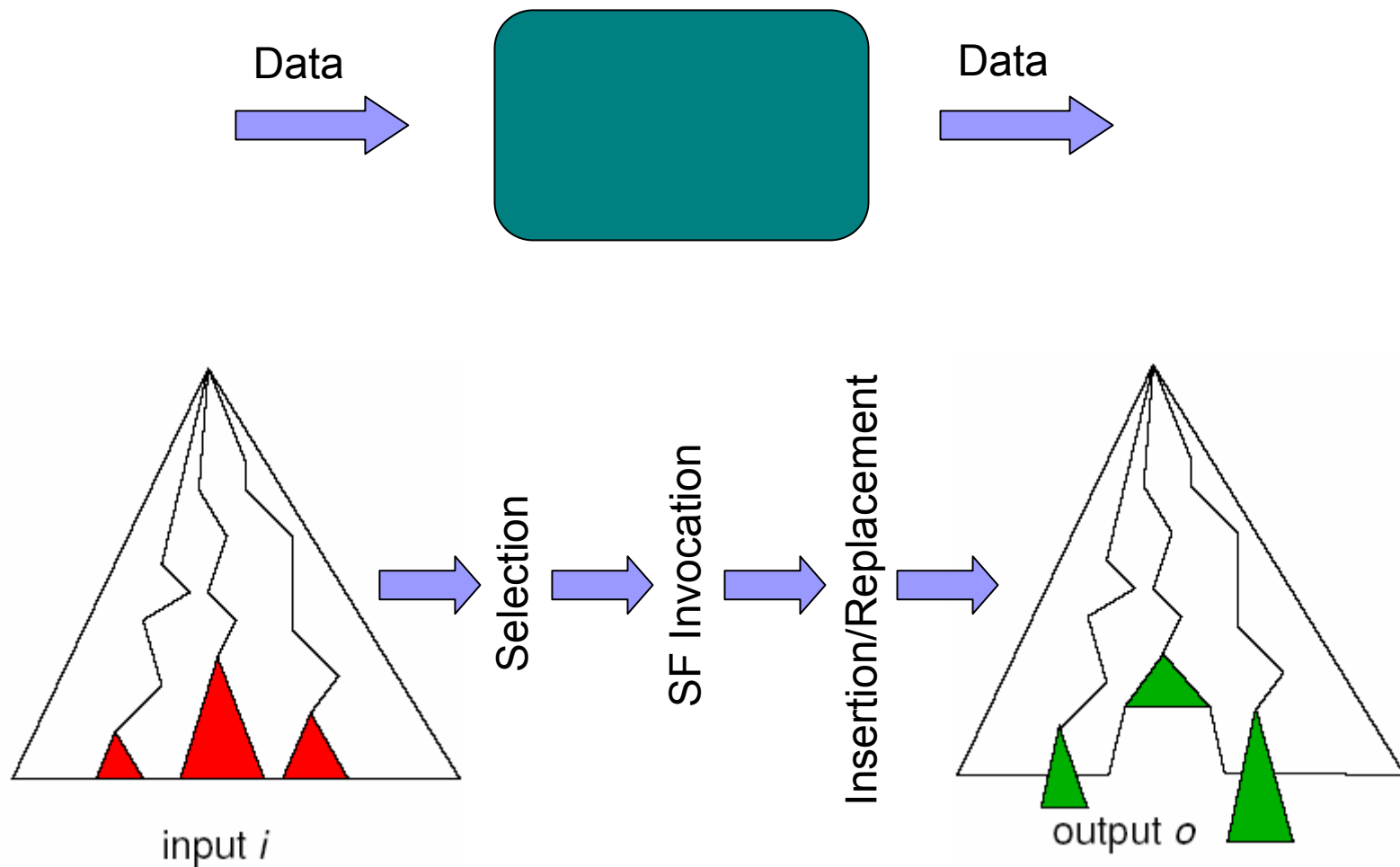
Configurations
(white-box)

- Access data in XML stream
- Call Scientific Functions (Services)
- Put results back into stream

Scientific Functions
(black-boxes)



Δ -XML Actor



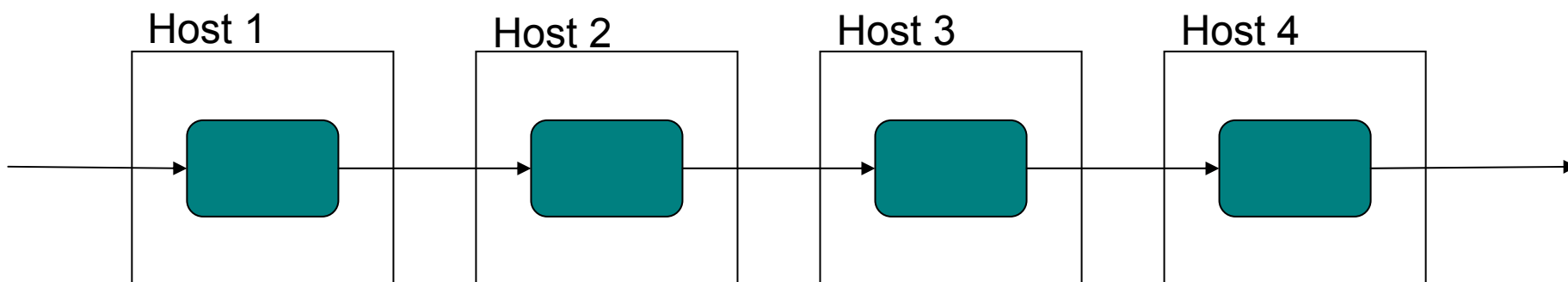
XML Processing Pipelines

- Δ -XML = XML processing with
 - Sequence of updates instead of queries
 - Calls to external functions
 - Data-centric XML data

- **Δ -XML: XML Processing Pipelines**
- **X-CSR: Dataflow-Optimization for Δ -XML**

Dataflow Optimization X-CSR

- Δ -XML execution scenario



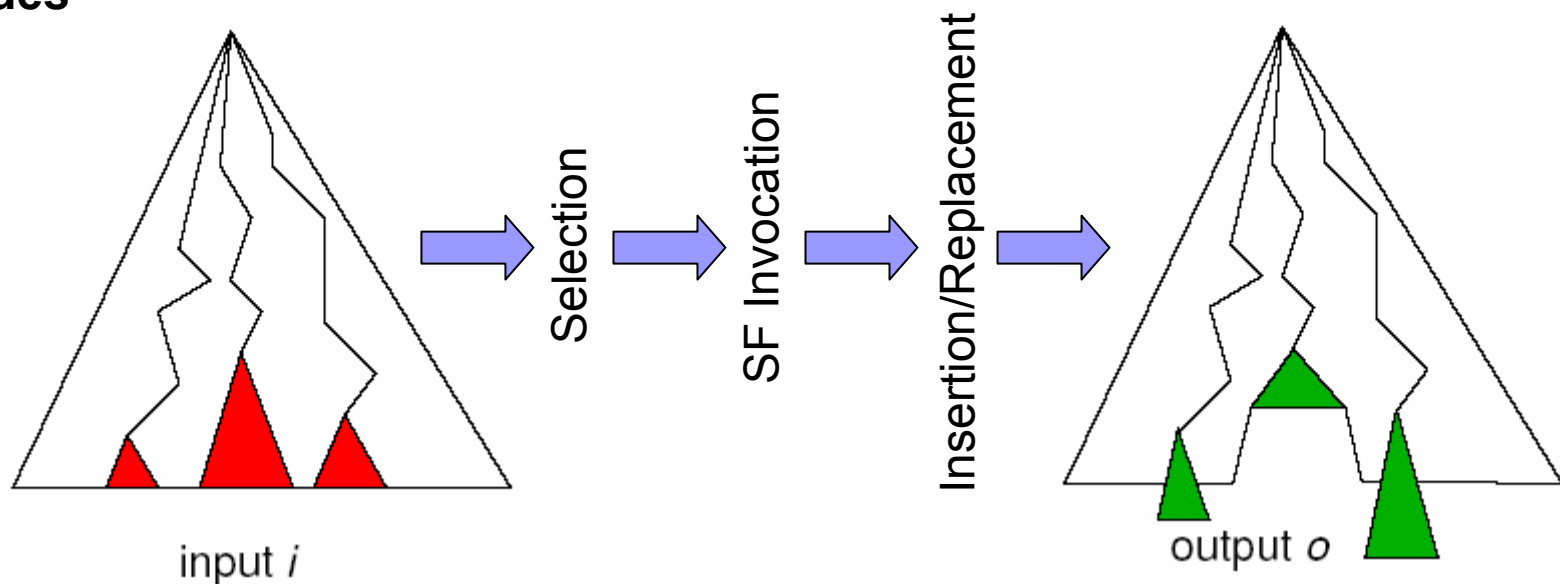
- Δ -XML advantages come with overhead:
 - Unnecessary data shipments (selective reads)
- X-CSR: Only ship ***selected*** data

Abstracting Actor Behavior

- Q1: What is *selected* by A?
- Q2: How does A change the stream?
- → Signature of A

Schemas

Values

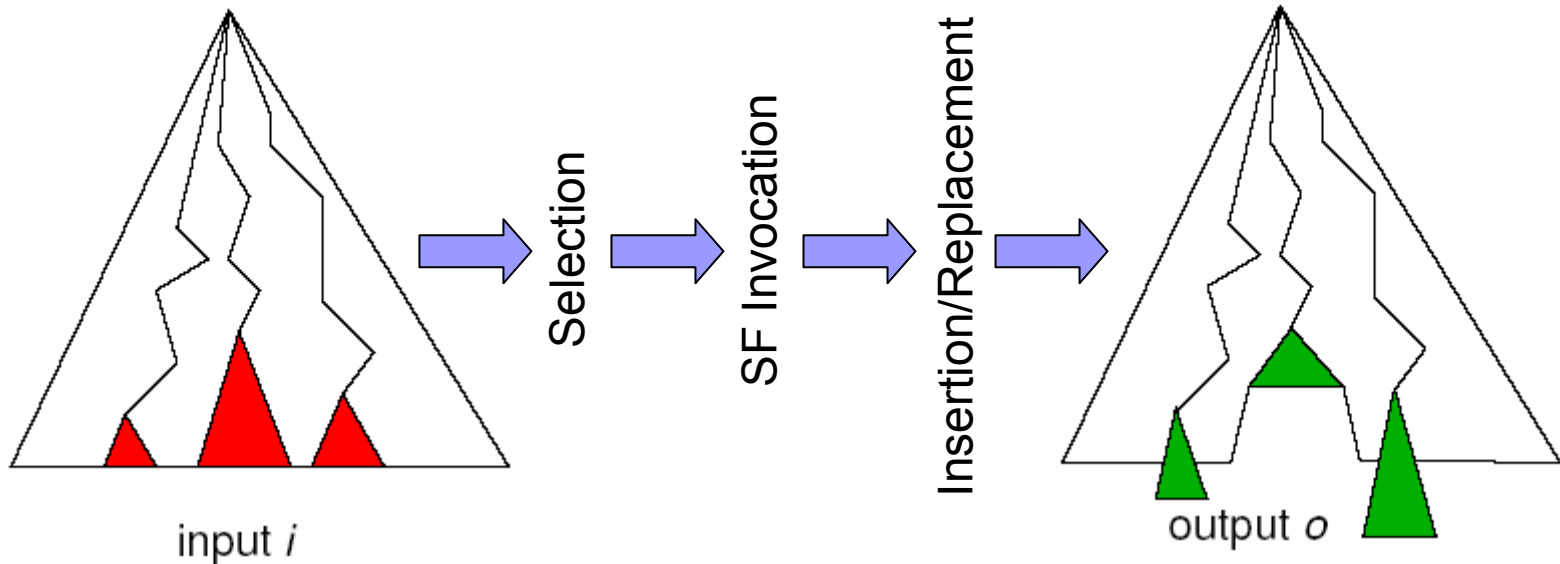


Abstracting Actor Behavior



Schemas

Values



Tree Grammars

- Characterize sets of labeled trees / valid input instances for an actor
- Similar to DTDs, but can capture upper-context information
- Roughly equivalent to XMLSchema

Tree Grammars

Formal Definition

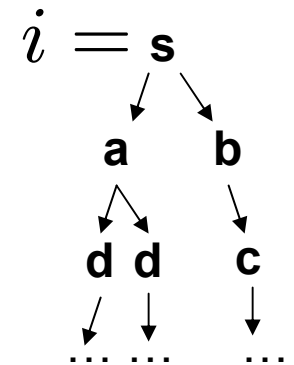
- Set of *type names* \mathbf{T}
- Set of *labels* \mathbf{L}
- Set of *type declarations* of the form

$$A ::= \langle a \rangle \mathcal{R}(\mathbf{T}) \quad , \text{ with } A \in \mathbf{T}, \langle a \rangle \in \mathbf{L}$$

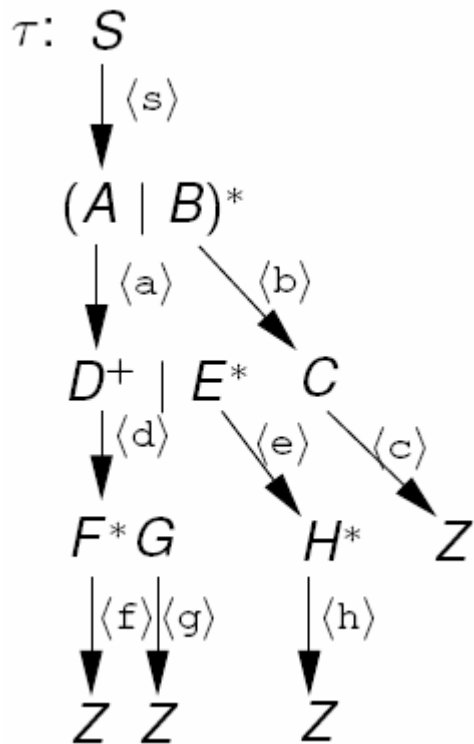
Example

$$\tau = \left\{ \begin{array}{l} S ::= \langle s \rangle (A \mid B)^* , \\ A ::= \langle a \rangle D^+ \mid E^* , \\ B ::= \langle b \rangle C , \\ C ::= \langle c \rangle Z , \\ D ::= \langle d \rangle F^* G , \dots \end{array} \right\}$$

$i = \langle s \rangle$
 $\langle a \rangle$
 $\langle d \rangle \dots \langle /d \rangle$
 $\langle d \rangle \dots \langle /d \rangle$
 $\langle /a \rangle$
 $\langle b \rangle$
 $\langle c \rangle \dots \langle /c \rangle$
 $\langle /b \rangle$
 $\langle /s \rangle$



Example Schema



- Use type Z to denote base data (holds the actual data)
- Non-ambiguous
- Non-recursive

$$\mathbf{T}_\tau = \{S, A, \dots, Z\}$$

$$\mathbf{L}_\tau = \{\langle s \rangle, \langle a \rangle, \dots, \langle h \rangle\}$$

Actor Signatures

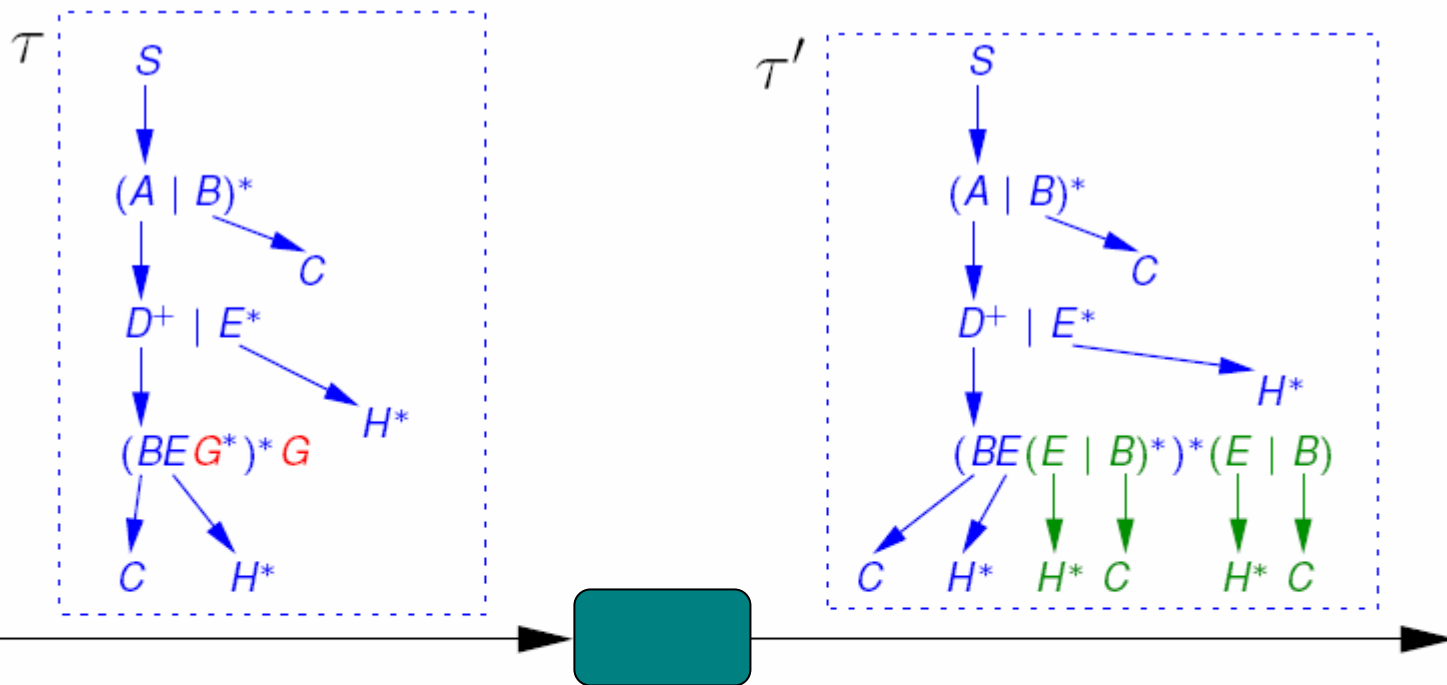
Signature captures ...

- What the actor reads, and
- How stream is modified

Formal:

- $\Delta_A = \langle \sigma, \tau_\alpha, \tau_\omega \rangle$
- τ_α – input selection schema
- τ_ω – new output schema parts
- σ – set of match rules, each of the form $X \rightarrow \mathcal{R}$ with
 - $X \in \tau_\alpha$, and
 - \mathcal{R} regexp. over types in $\tau_\alpha \cup \tau_\omega$

Example



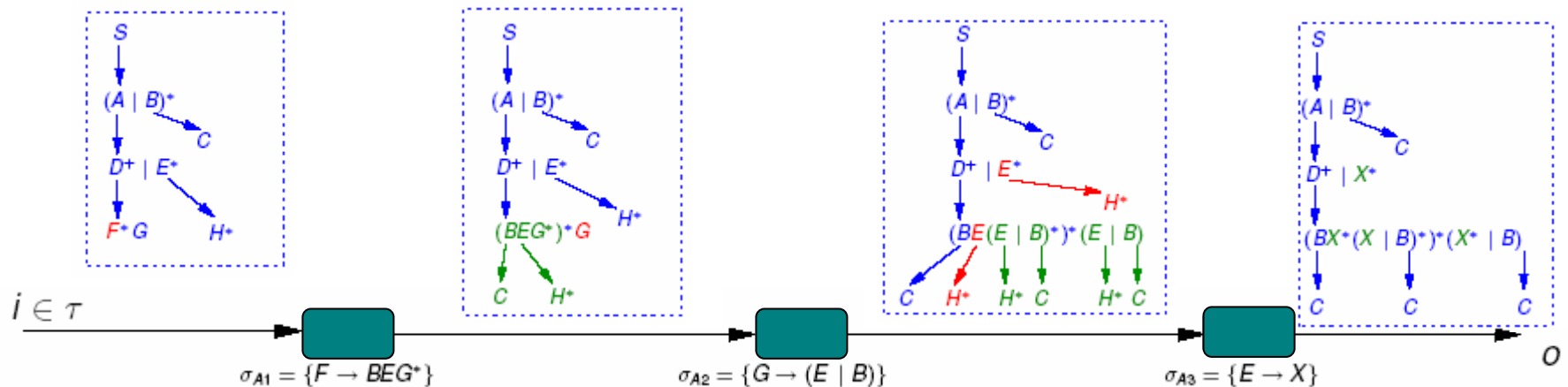
$\Delta_A = \langle \sigma, \tau_\alpha, \tau_\omega \rangle$ with

$\sigma : \{G \rightarrow E | B\},$

$\tau_\alpha : \{G ::= \langle g \rangle Z\},$ and

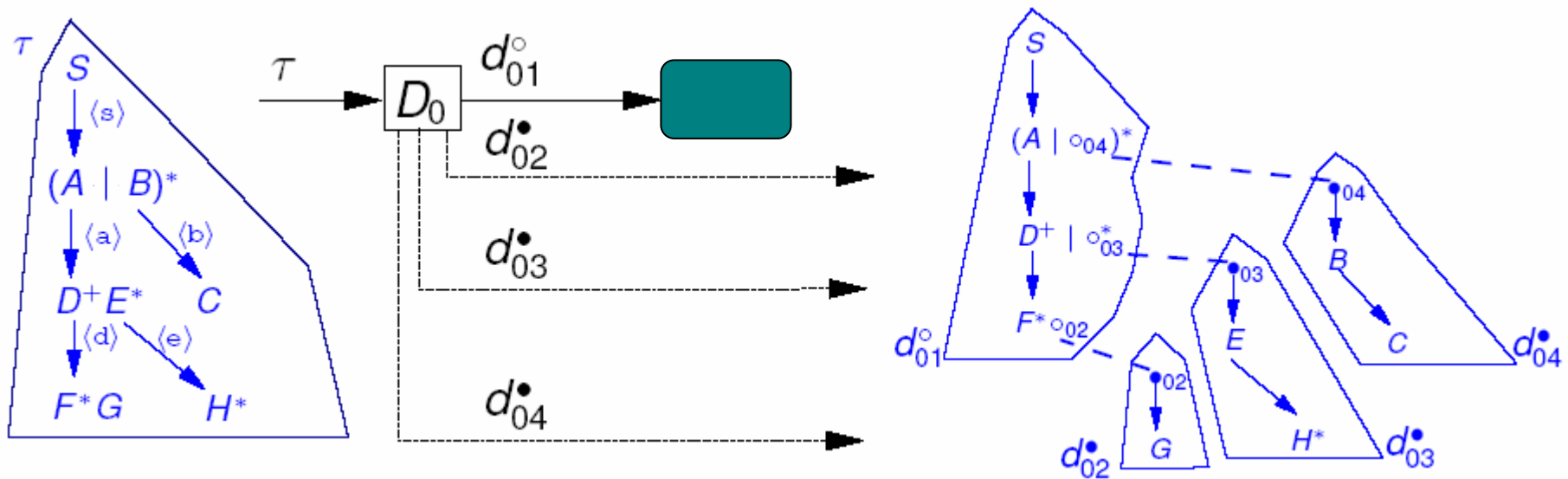
$\tau_\omega : \{E ::= \langle c \rangle H^*, B ::= \langle b \rangle C, H ::= \langle h \rangle Z, C ::= \langle c \rangle Z\}.$

Modeled Workflow



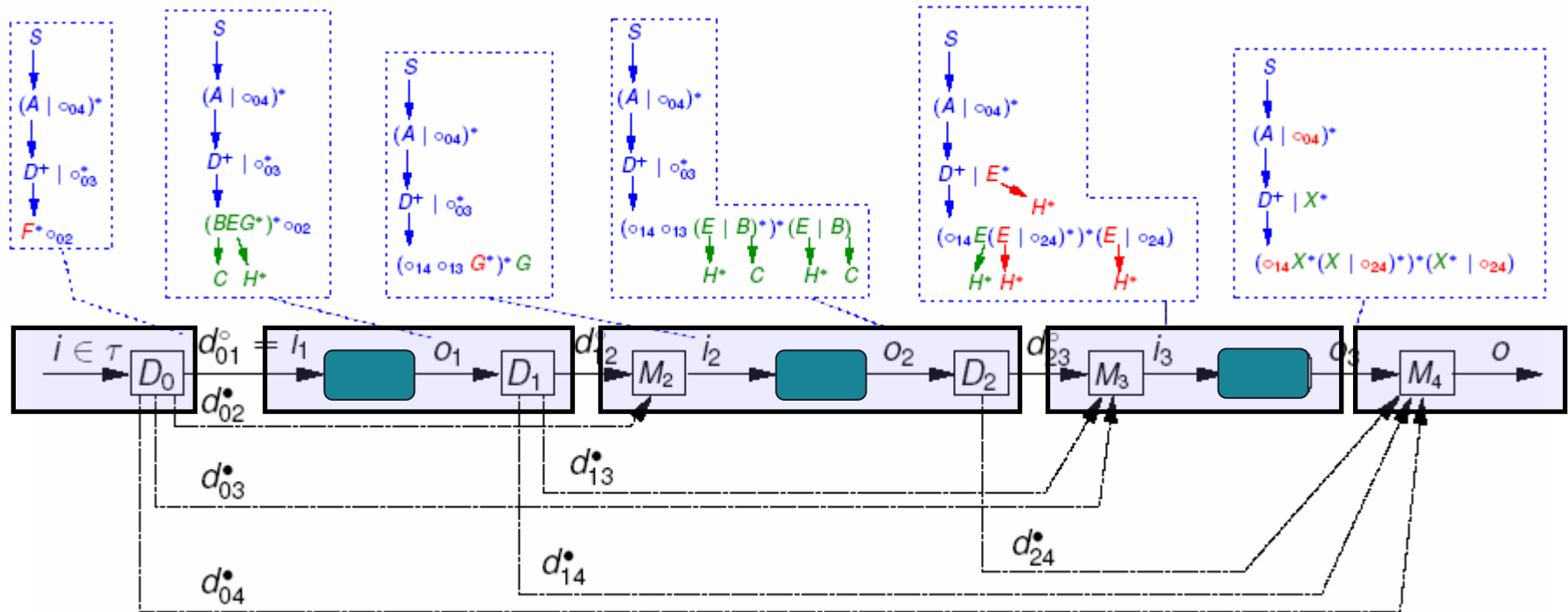
- Type propagation
- Data dependency analysis

Shipping Idea



- Bypassing approach
- Keep associations via holes and fillers

Auto-generated Shipping-optimized Workflow



- Distributors and mergers
- Direct connections between all hosts

Compile-time

- Routing specification computed by dataflow-analysis algorithm
- Input schema, annotated with destination for each type

Runtime

- Glushkov automata to parse input stream
- Data is sent to different output channels and holes and fillers are inserted

Holes Need Two Indexes

i -- From which distributor

- If not, merger would not know from where to read the data

j -- To which merger

- If not, merger would not know whether it should fill the hole (or pass on)

Mergers

Merger j :

- 1. Read from main line until i,j -hole is found (forward all other holes and data)
- 2. Read from bypass line i until closing i,j -filler
- 3. Go back to Step 1

Correctness and Optimality

Correctness

- If data fragment X is bypassed around A , then X is not read by A

Optimality

- If data fragment X is merged into A , then X is read by A

$O(\text{shipping savings})$?

- Linear in size of bypassed data
- Linear in number of bypassed actors

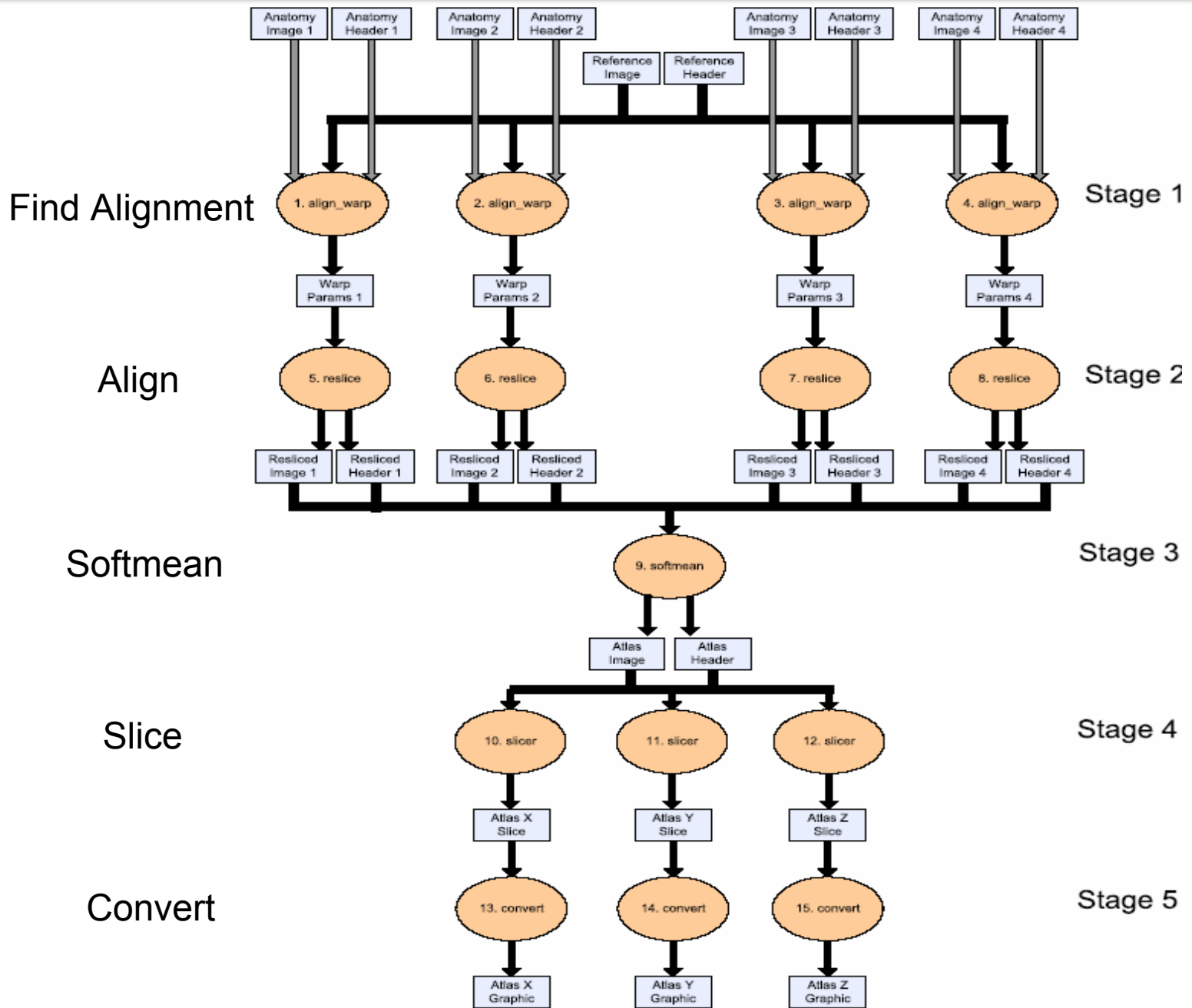
Experimental Evaluation

Environment

- Cluster of 40 Linux nodes
- 100 MBit/s switched network

System Implementation

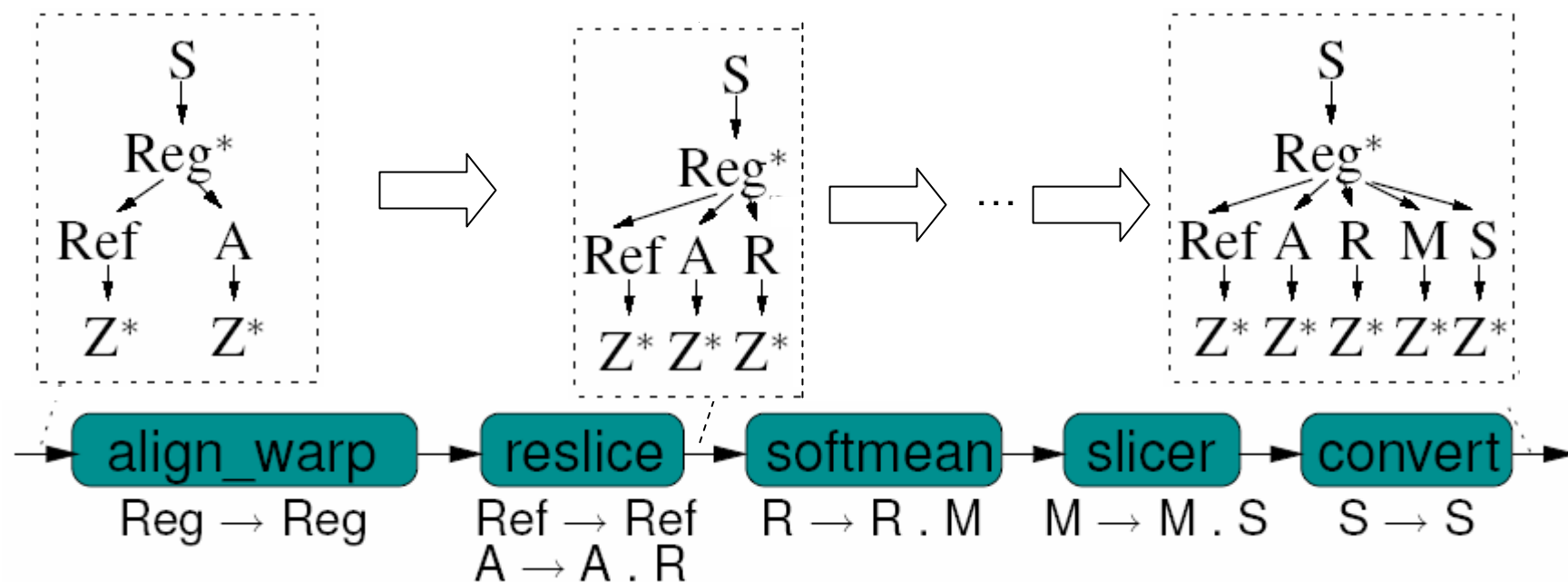
- Each actor as Unix process
- Written in C++/Perl on top of PVM
- PVM for passing Collection Tokens
- Base data as local files, scp to move



fMRI Workflow

Functional Magnetic Resonance Imaging

- Post-processing of voxel images of the human brain
- Data intensive
- Computationally expensive



Experimentation Results for fMRI Workflow

Input/Output data size:

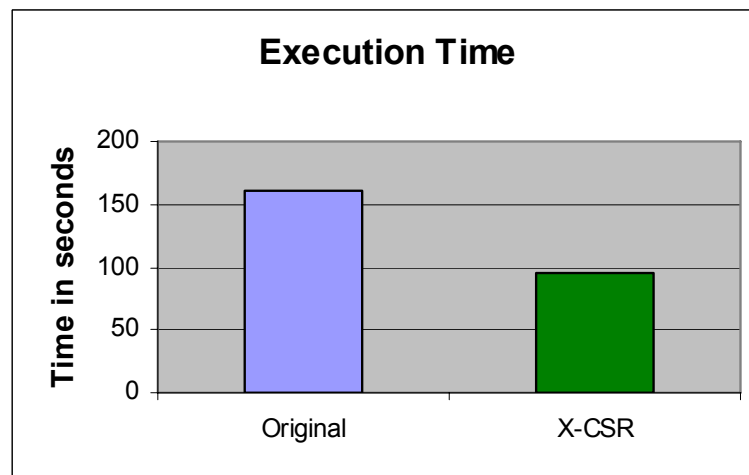
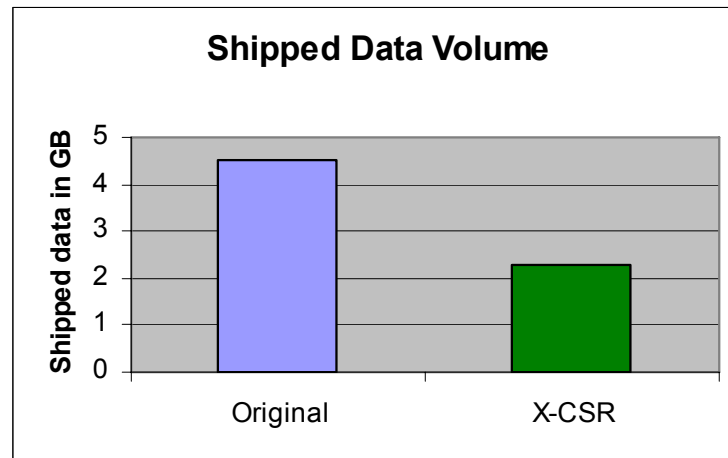
400MB/900MB

Data shippings:

4.5GB vs. 2.3GB (-49%)

Execution time:

160s vs. 95s (-40%)



Summary

Δ -XML

- Flexible paradigm for processing tree structured data
- Currently used in scientific workflow system Kepler

Static Analysis

- Signatures approximate actor behavior
- More than type-in + type-out
- Track data dependencies

Shipping Optimization

- Automatic optimization informed by static knowledge
- 40% faster for scientific example

Thank you.



Questions?