

Some Visualization Models applied to the Analysis of Parallel Applications

Lucas Mello Schnorr

Advisors:

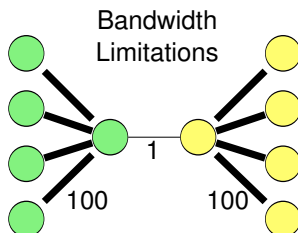
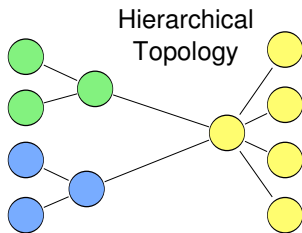
Philippe O. A. Navaux & Denis Trystram & Guillaume Huard

Federal University of Rio Grande do Sul, Brazil
Grenoble Institute of Technology, France



Introduction - Context

- Distributed Systems → Grids
- Grid Interconnection and Scalability
 - Topology and Connectivity
 - Performance: bandwidth and latency
 - New resources can be added very easily



- Influence in the application execution
- Visualization – Performance Analysis

Introduction - Existing Tools/Techniques

■ Statistical Techniques

- ParaGraph (1990) – bar charts, utilization Count
- Pablo (1993) – bar charts + 3D scatter plot

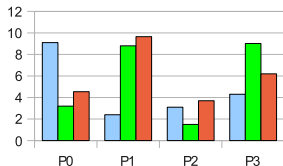
■ Behavioral Techniques

- Vampir (1996) – time-line system view
- Jumpshot (1999), Pajé (2000) – space-time

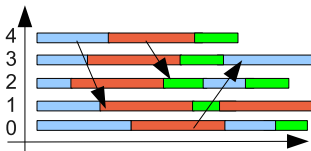
■ Structural Techniques

- ParaGraph (1990) – network display / hypercube

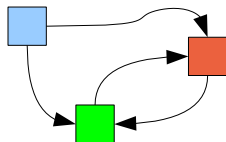
Statistical



Behavioral



Structural



Introduction - Problem Identification

- Lack of a network-aware analysis
 - Difficult to analyze using space-time views
 - Structural techniques undeveloped
- Problems of visualization scalability
 - Visualization techniques limitations reached
 - Analysis are limited to hundreds of entities

Desirable Characteristics for Application Analysis

→ The Objectives

- Consider network properties
- Visualization scalability in the analysis

Introduction - The Thesis Approach

- Explore techniques from Information Visualization
- Context of parallel application analysis
 - Grid resources
 - Thread/Process parallel applications

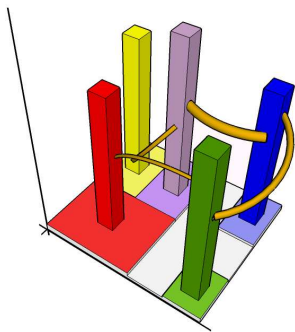
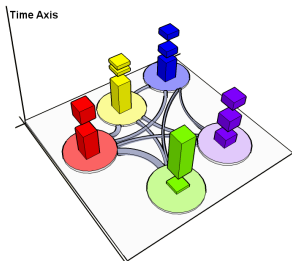
Proposed Visualization Models

- Behavioral and Structural/Statistical (3D)
 - Communication Pattern
 - Network topology + Communication Pattern
 - Logical representation
- Visual Aggregation
 - Large-scale traces
 - Local and Global summaries

Outline

3D Model - Visual Conception

- Resources represented in 2D
 - Structural (e.g. a graph)
 - Statistical
- Vertical dimension is time
 - Objects' Behavior Evolution
 - States and Links
- Interaction Techniques
 - Notion of a Camera
 - Rotation
 - Translation
 - Objects Animation
 - Replay step-by-step



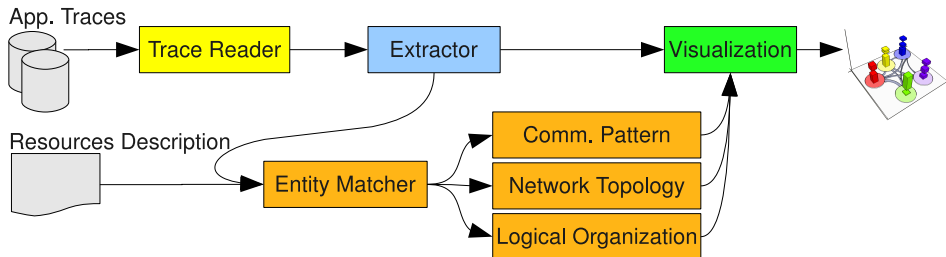
3D Model - Differences from existing tools

- 3D Statistical Representation
 - Pablo → 3D Scatter Plot
 - Paradyn → 3D Terrain
 - ParaProf → Triang Mesh, 3D Bar and 3D Scatter Plot
- 3D Behavioral Representation
 - ParaProf → 2 metrics and time
 - Virtue → the time-tunnel view

Our Approach

- Presence of a timeline to show objects' evolution
- Multiple Configurations in the visualization base

3D Model - Abstract Component Model

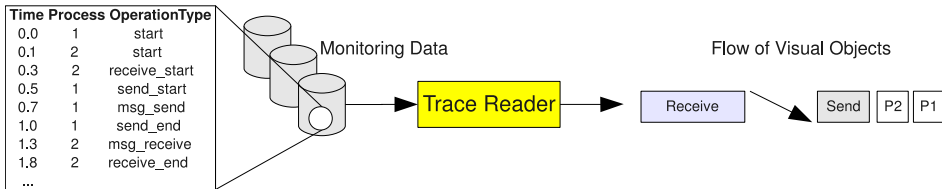


Input Data

- Application Traces
 - Timestamp-based events
 - Behavior registered
- Resources Description
 - Network topology: graph
 - Logical resource organization: tree

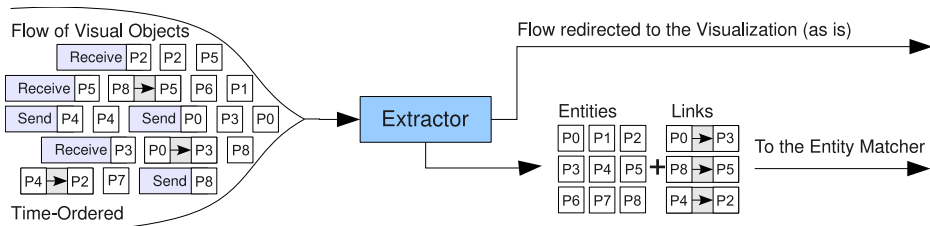
3D Model - The Trace Reader

- Deals directly with application traces and events
- Only trace-dependent part of the model
- Transform events into high-level visual objects
 - Container → Entities
 - State/Variable/Event → Evolution
 - Link → Communications
- No semantics → Visualization is generic



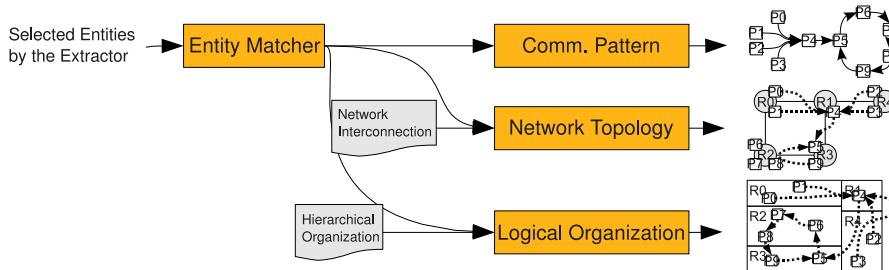
3D Model - The Extractor

- Supply entity matcher needs: links and entities
- Attribute entities with location data
 - where a process is executed
 - which process a thread belongs to
- Input is also redirected to the Visualization module



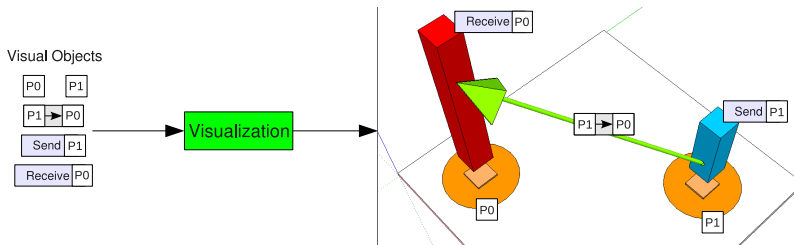
3D Model - The Entity Matcher

- Responsible for the Visualization Base layout
- Three possibilities of configuration are proposed
 - Communication Pattern (deadlocks, ...)
 - Network Topology (network utilization, routes, ..)
 - Logical Organization (load balancing, ...)



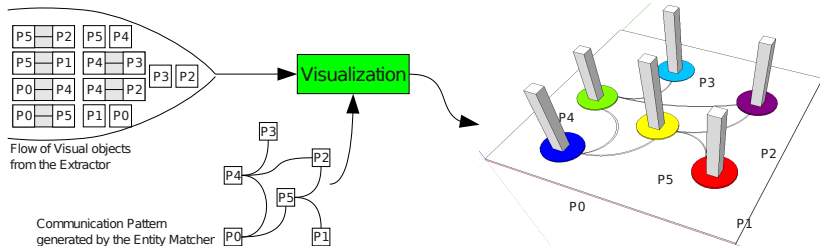
3D Model - Visualization

- How the visual objects are represented in 3D



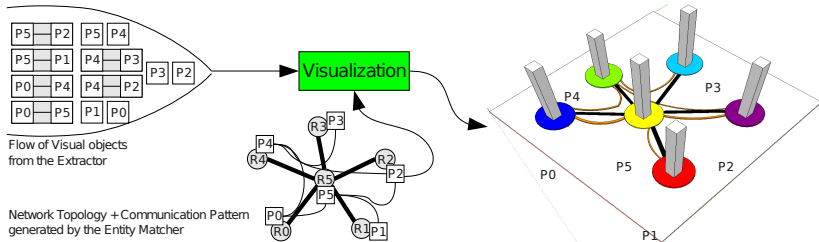
3D Model - Visualization

- How the visual objects are represented in 3D
- Rendering the visualization base
 - Application Communication Pattern



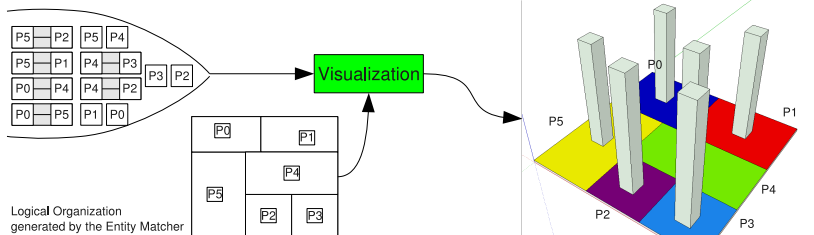
3D Model - Visualization

- How the visual objects are represented in 3D
- Rendering the visualization base
 - Application Communication Pattern
 - Network Topology + App. Communication Pattern



3D Model - Visualization

- How the visual objects are represented in 3D
- Rendering the visualization base
 - Application Communication Pattern
 - Network Topology + App. Communication Pattern
 - Logical Organization of Resources



Outline

Aggregation Model - Overview

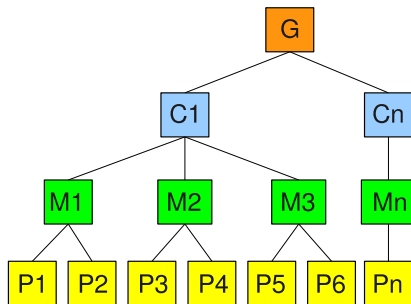
- Enable large-scale trace analysis
- Visually compare entities behavior
- Detect global and local characteristics

Steps of the Model

- 1 Hierarchical Monitoring Data
 - 2 Time-Slice algorithm (temporal integration)
 - 3 Aggregation model (spatial integration)
 - 4 Treemap representation
- Visualization differences from existing tools
 - PlanetLab's CoVisualize → resources
 - Treemap for Workload Visualization [Stephen 2003]
 - Lack of configurable time intervals, aggregated data

Hierarchical Monitoring Data

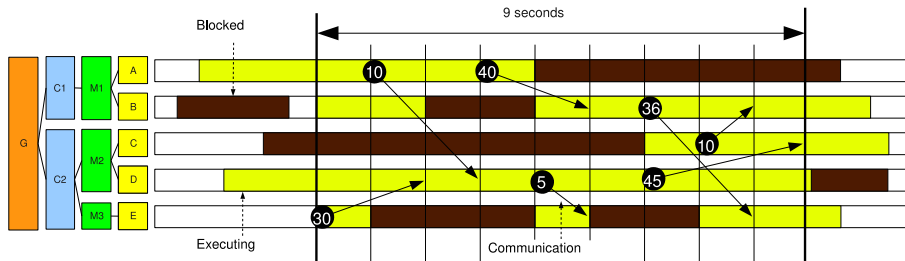
- Monitoring systems register entities behavior
- Entities can be processes and threads
- They can be organized as a hierarchy
 - Logical hierarchy
 - Geographical Location hierarchy
 - Other possibilities: libraries, components
- Grid'5000 example



Time-Slice Algorithm - Basics

Objective: annotate leaf nodes of the hierarchy

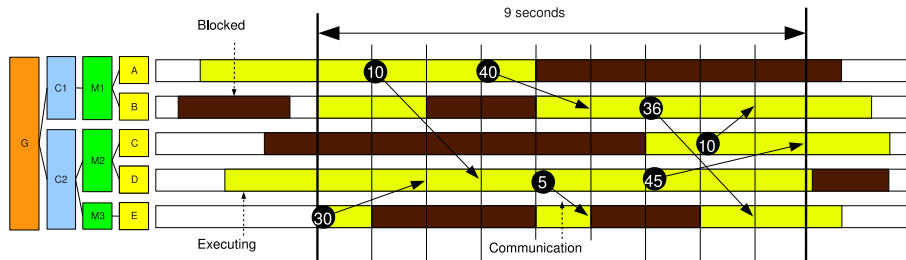
- Time-slice definition
- Summary of trace events on the interval
 - States, Variables, Links, Events, ...



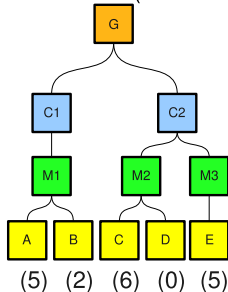
Output of the Algorithm

- Hierarchy of input + computed values on leaves

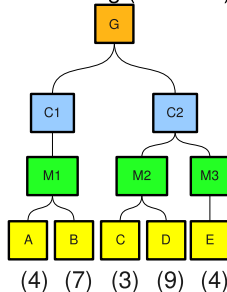
Time-Slice Algorithm - Example



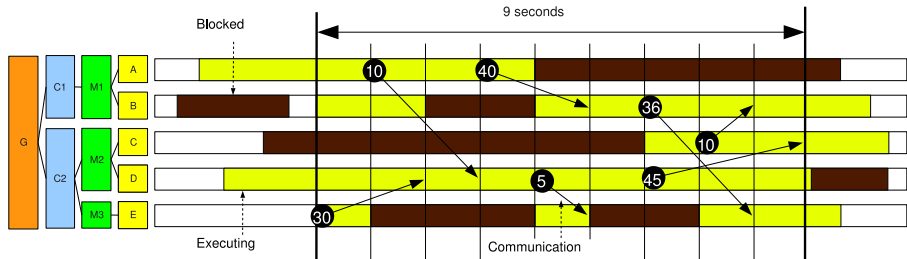
Blocked (seconds)



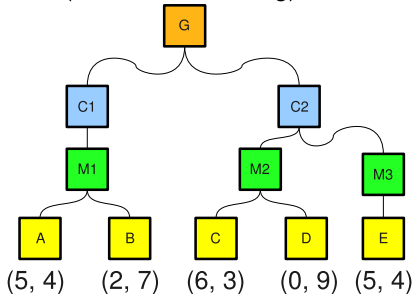
Executing (seconds)



Time-Slice Algorithm - Example



(Blocked, Executing)

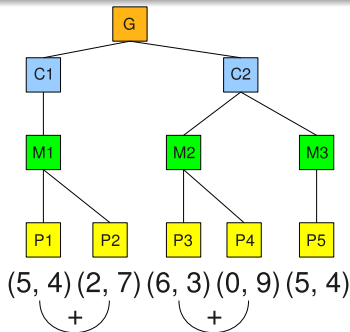


Aggregation Model

- Objective: **aggregated values** at intermediary levels

Aggregation Functions

- add, subtract, multiply, divide, max, min, median, ...
- Depends on
 - what type of value the leaves have
 - the desired statistical result

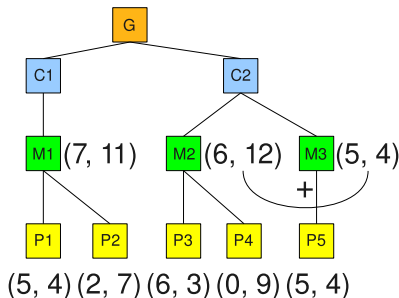


Aggregation Model

- Objective: **aggregated values** at intermediary levels

Aggregation Functions

- add, subtract, multiply, divide, max, min, median, ...
- Depends on
 - what type of value the leaves have
 - the desired statistical result

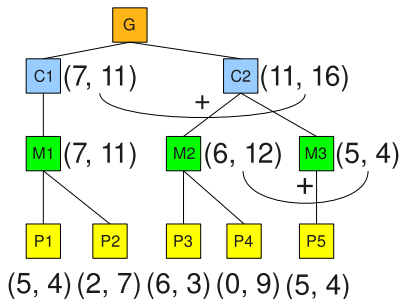


Aggregation Model

- Objective: **aggregated values** at intermediary levels

Aggregation Functions

- add, subtract, multiply, divide, max, min, median, ...
- Depends on
 - what type of value the leaves have
 - the desired statistical result

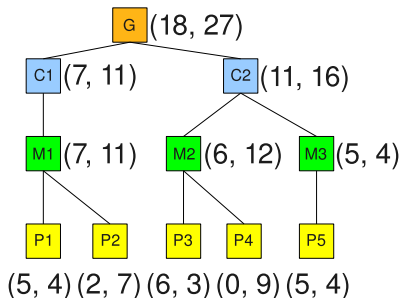


Aggregation Model

- Objective: **aggregated values** at intermediary levels

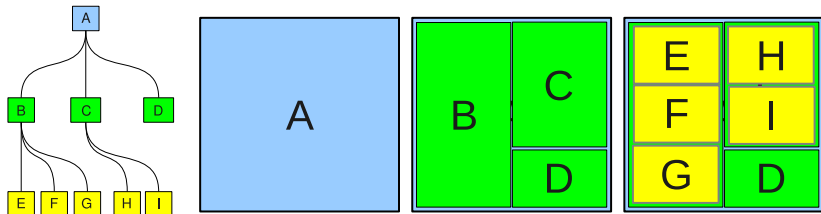
Aggregation Functions

- add, subtract, multiply, divide, max, min, median, ...
- Depends on
 - what type of value the leaves have
 - the desired statistical result



Visualization of the Approach - Treemaps

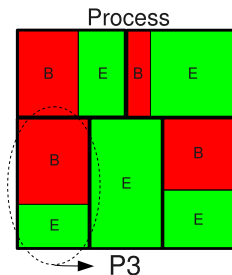
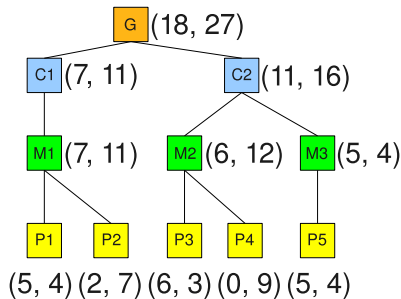
- Technique created in 1991
- Scalable hierarchical representation
- Algorithm
 - Top-down drawing
 - For a given node, split screen space among children



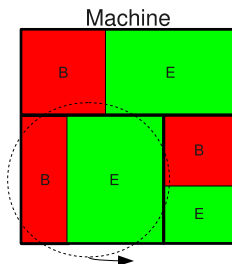
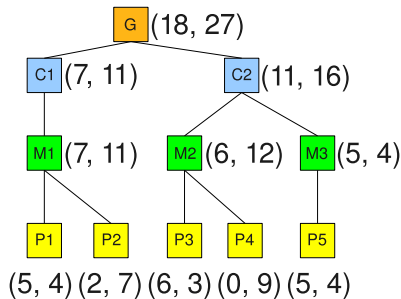
Original algorithm has several evolutions

- Squarified treemap is used here
 - Keeps rectangles as close to squares as possible

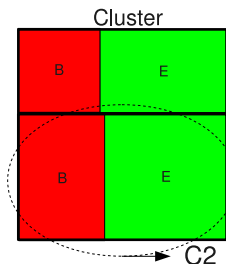
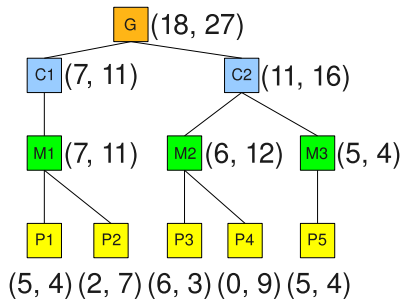
Treemap to view the Aggregated Hierarchy



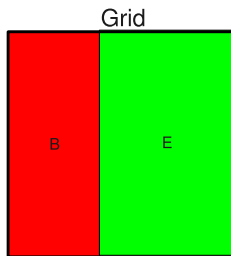
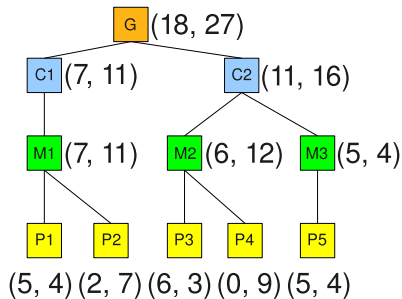
Treemap to view the Aggregated Hierarchy



Treemap to view the Aggregated Hierarchy



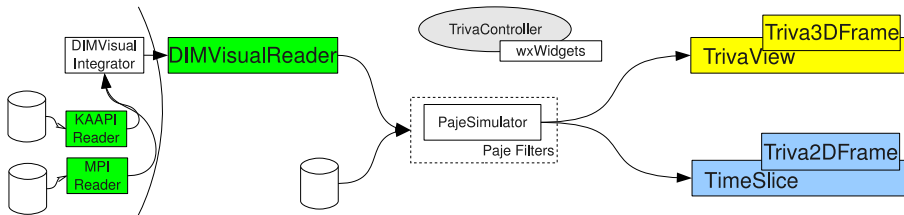
Treemap to view the Aggregated Hierarchy



Outline

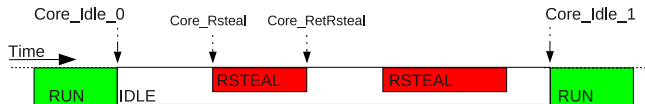
Triva Prototype Implementation

- Developed in Objective-C and C++
- Combine several existing tools
 - DIMVisual library
 - Pajé Components (the Simulator)
 - Graphviz, Ogre3D, wxWidgets
- Performance evaluation of Pajé
 - Able to handle large-scale traces
 - Small response-time
 - Memory limitations

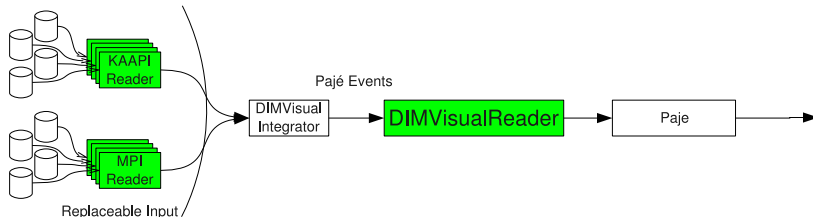
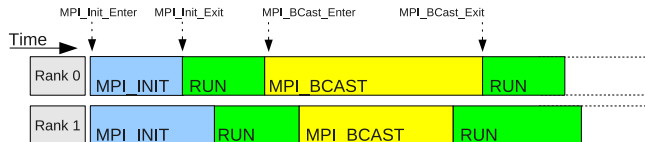


DIMVisualReader - Trace Reader

■ Built-in instrumentation of KAAPI library

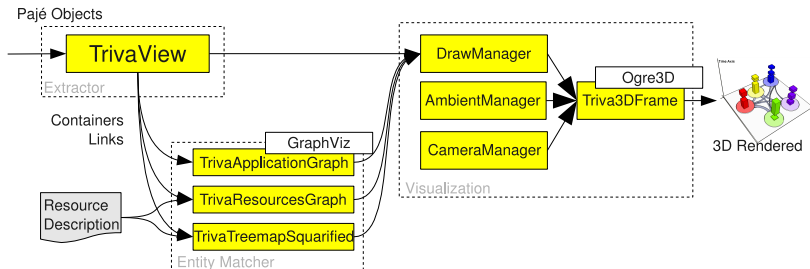


■ MPIRastro wrapper for MPI applications



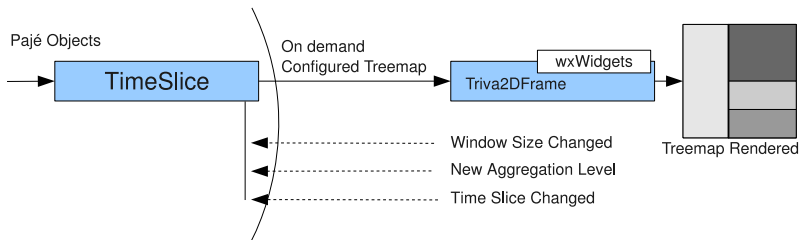
TrivaView - The 3D Approach

- Model: Extractor, Entity Matcher & Visualization
- Interaction Techniques (Ambient, CameraManager)
- Base configuration
 - Application Comm. Pattern created with GraphViz
 - Network Topology description (dot format)
 - Logical Organization (plist format)
- Placement on the Visualization Base
- Rendering the 3D Timestamped Pajé Objects



TimeSliceView - The Aggregation Model

- Only two components
 - TimeSlice Filter
 - Triva2DFrame
- Time-Slice Algorithm and Aggregation Model
- Implementation of the Squarified Treemap Algorithm
- Drawing the rectangles with the wxWidgets



Outline

Results

- Different application traces are used as input
- Results are composed of screenshots of the prototype

Objective

- Check if 3D visualizations enable a better understanding of traces with the network topology
 - Check if large-scale analysis are possible with the aggregation model
-
- Traces Description
 - 3D Visualization
 - Treemap Visualization

Results - Trace Description

- Synthetic traces

- Large-scale hierarchies (up to 100 thousand)
- Typical Communication Patterns



- Real traces

- KAAPI Traces



- MPI Traces

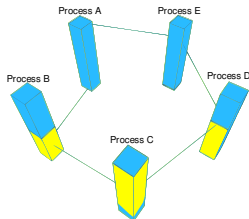
- Grid'5000 platform in France

- Xiru Cluster at Porto Alegre

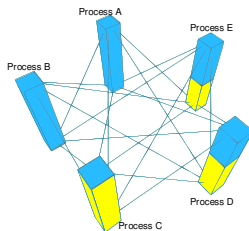
3D Visualization - Communication Patterns

- Differences from the space-time diagram
 - 3D enables Graph-like representations
 - with time evolution

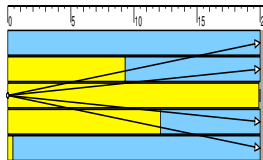
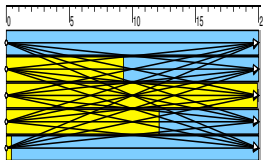
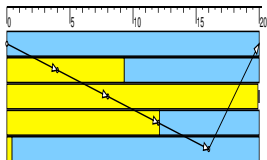
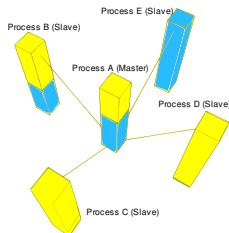
Ring Communication Pattern



Fully-Connected

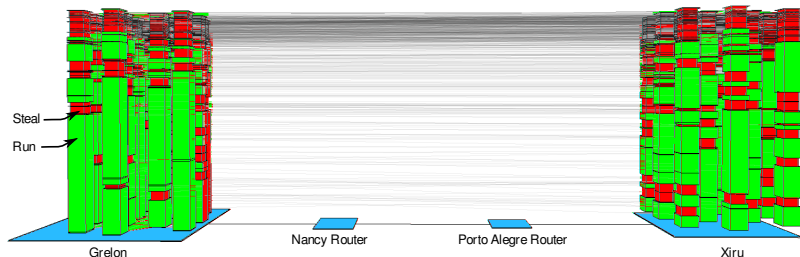


Star



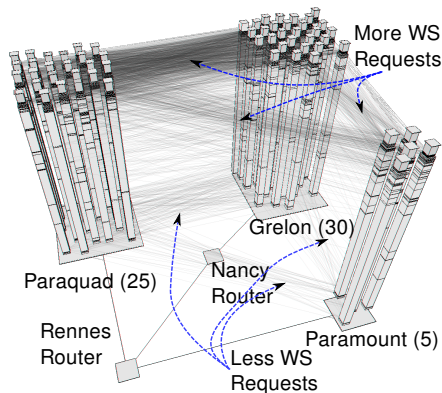
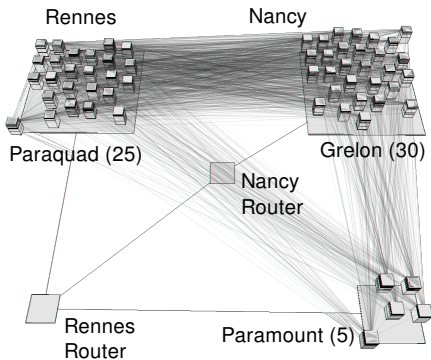
3D Visualization - KAAPI Trace

- Fibonacci Application
- 26 processes, two sites, two clusters
- Lines represent steal requests
- Different number of communication between clusters
 - beginning → big tasks, less communication
 - end → smaller tasks, more communication



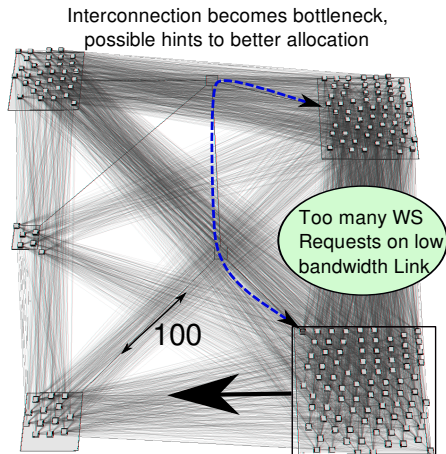
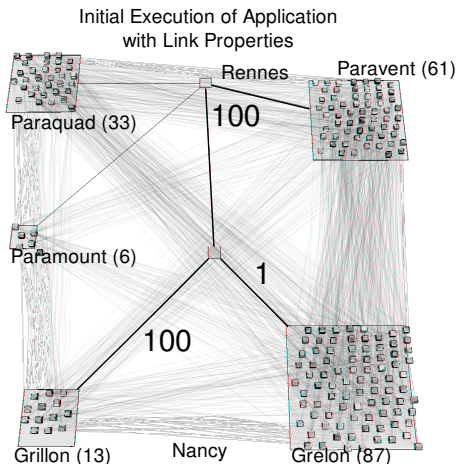
3D Visualization - KAAPI Trace

- 60 processes, two sites, three clusters
- Total execution time of a KAAPI fibonacci application
- Observe number of requests in time



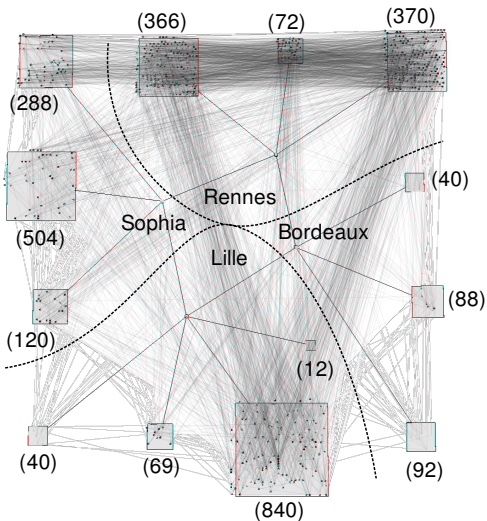
3D Visualization - KAAPI Trace

- 200 processes, 200 machines, two sites, five clusters
- Annotated manually with bandwidth limitations

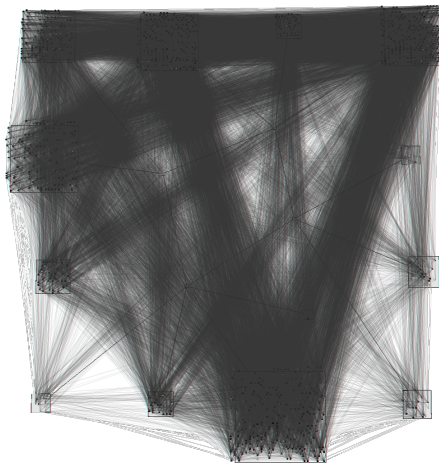


3D Visualization - KAAPI Trace

- 2900 processes, four sites, thirteen clusters



End of Execution

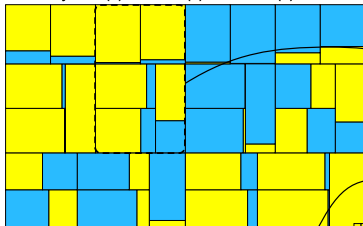


Treemap Visualization - Description

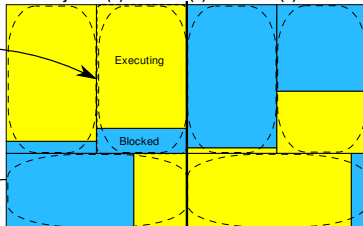
Time-Slice and Aggregated Hierarchies

- Interaction Techniques: mouse wheel, mouse over
- Detailed information is available in the status bar

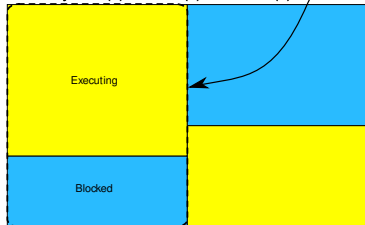
Hierarchy: Site (2) - Cluster (3) - **Machine** (5)



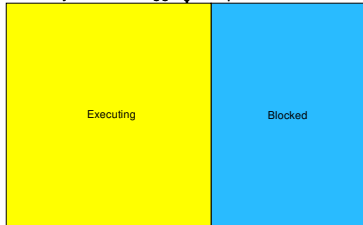
Hierarchy: Site (2) - **Cluster** (3) - Machine (5)



Hierarchy: **Site** (2) - Cluster(3) - Machine (5)

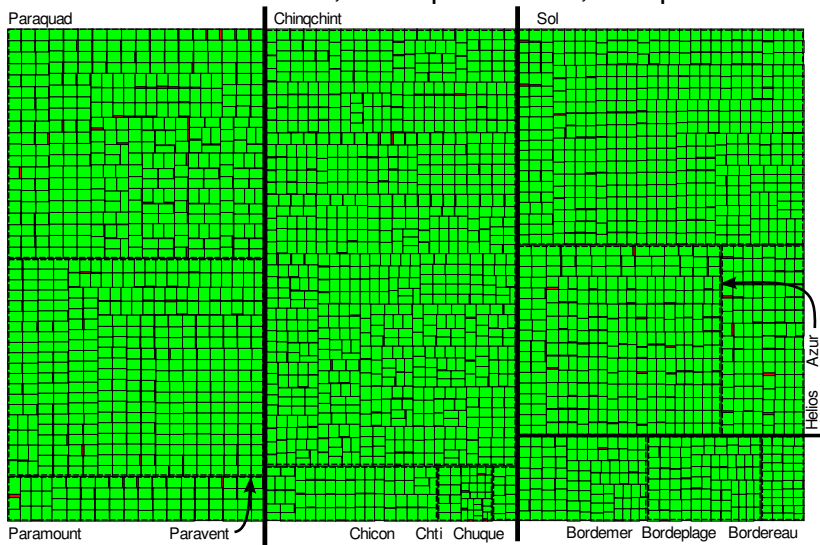


Hierarchy: maximum aggregation possible



Treemap Visualization - KAAPI Trace

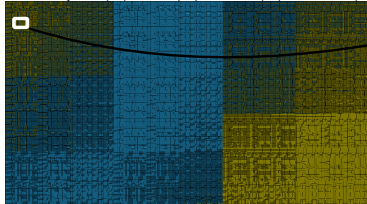
Run and **RSteal** states, 2900 processes, 310 processors



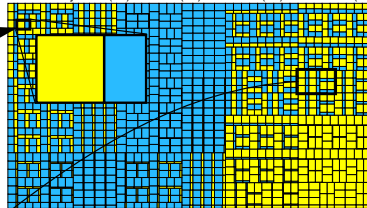
Treemap Visualization - Large-Scale

- Synthetic trace with 100 thousand processes
- Two states, four-level hierarchy

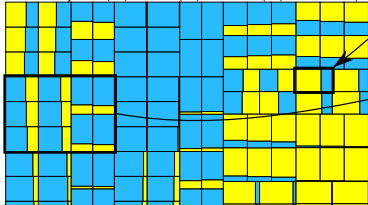
A Hierarchy: Site (10) - Cluster(10) - Machine (10) - **Processor**(100)



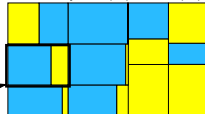
B Hierarchy: Site (10) - Cluster(10) - **Machine** (10) - Processor (100)



C Hierarchy: Site (10) - **Cluster**(10) - Machine (10) - Processor (100)



D Hierarchy: **Site** (10) - Cluster(10) - Machine (10) - Processor (100)



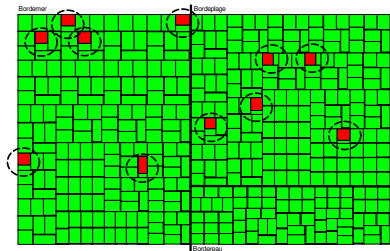
E Maximum Aggregation



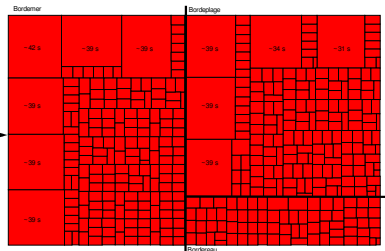
Treemap Visualization - KAAPI Trace

- 400 processes, 50 machines, one site
- 8 processes per machine
 - Overload of some machines with 2 CPUs
 - Unusual amount of time in Steal state
- Machines with 4 CPUs show normal behavior

A Larger **RSteal** states, for each K-Processor



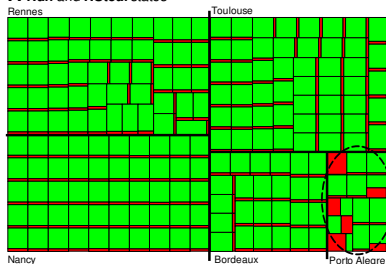
B Showing only **RSteal** state, for each K-Processor



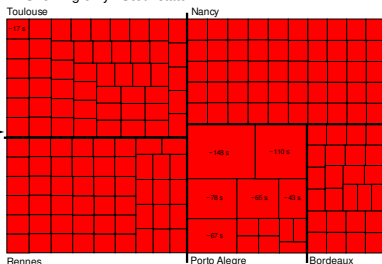
Treemap Visualization - KAAPI Trace

- 188 processes, 188 machines, five sites
- Different behavior at Porto Alegre
- Probably due to the interconnection
 - Latency for Grid'5000 in France: ~ 10 ms
 - Latency between Porto Alegre and France: ~ 300 ms
- More time spent in work stealing functions

A Run and RSteal states



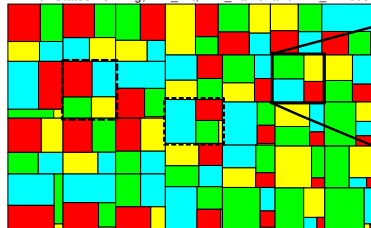
B Showing only RSteal state



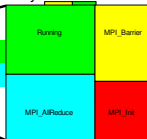
Treemap Visualization - MPI Trace

- Traces from the EP application – NAS Benchmark
- 32 processes – time spent in each MPI operation
- Init/Barrier: might indicate a linear implementation

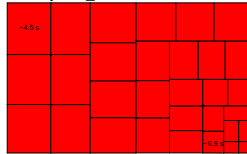
A With States Running, MPI_Init, MPI_Barrier and MPI_AllReduce



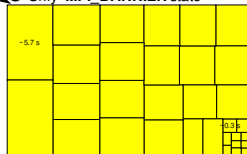
Only Process Rank 21



B Only MPI_INIT state



C Only MPI_BARRIER state



Maximum Aggregation



Conclusions

The problem identified in the Thesis

- Lack of structural visualization analysis
- Visualization scalability

Main Achievements

- Behavioral with Structural/Statistical Model (3D)
 - Analysis considering network structure
 - Experiments using Grid'5000 platform
 - Identification of behavior in KAAPI work stealing
- Time-Slice Technique & Aggregation Model
 - Validated with real-scenario with 2900 processes
 - Tested with synthetic traces up to 100K processes
 - Load-balance efficiency / global and local summaries

Perspectives and Implications

■ Perspectives

- Show aggregated objects in the 3D visualization
- Other types of information for the time-slice technique
- Use of other aggregation functions
- Aggregation model to merge communication patterns
- ...

■ Implications

- Better understanding of parallel applications
 - consider execution environment details
 - large-scale visual analysis
- Re-thinking behavioral visualization
 - Do we need a timeline in representations?
 - Aggregated data
- Use of information visualization techniques

Publications

2	PARALLEL APPLICATION VISUALIZATION	23
2.1	Historical Evolution	23
2.2	Examples of Performance Visualization Tools	26
2.2.1	ParaGraph	27
2.2.2	TraceView	29
2.2.3	Pablo	29
2.2.4	Paradyne	30
2.2.5	Vampir	31
2.2.6	Virtue	32
2.2.7	Jumpshot	32
2.2.8	ParaProf	33
2.2.9	Pajé	35
2.3	Summary of Visualization Techniques	35
2.3.1	Behavioral	36
2.3.2	Structural	38
2.3.3	Statistical	39
2.4	Summary	41
3	THE THREE-DIMENSIONAL MODEL	43
3.1	Visual Conception	44
3.2	Model Overview	46
3.3	The Trace Reader	47
3.4	The Extractor	48
3.5	The Entity Matcher	49
3.5.1	Case 1: Parallel Application's	50
3.5.2	Case 2: Network Topology of	51
3.5.3	Case 3: Logical Organization	52
3.6	The Visualization	53
3.6.1	Rendering the Visualization	54
3.6.2	Interaction Mechanisms	56
3.7	Summary	57

Future Generation
Computer Systems
Journal

Sbac 2009

Grid 2008

4	VISUAL AGGREGATION MODEL	59
4.1	Hierarchical Organization of Monitoring Data	60
4.2	The Time-Slice Algorithm	62
4.2.1	States	63
4.2.2	Variables	64
4.2.3	Links	65
4.2.4	Events	66
4.2.5	More statistics	66
4.2.6	Example	67
4.3	The Aggregation Model	68
4.3.1	Aggregation Functions	69
4.4	Visualization of the Approach	70
4.4.1	Treemaps Basic Concepts	70
4.4.2	The Scalability Issue	71
4.4.3	Using Treemap in the Example	72
4.5	Summary	74
5	TRIVA PROTOTYPE IMPLEMENTATION	75
6	RESULTS AND EVALUATION	95
6.1	Traces Description	95
6.1.1	Synthetic Traces	96
6.1.2	KAAP Traces	99
6.1.3	MPI Traces	100
6.2	3D Visualizations	101
6.2.1	Description of the Visualization	101
6.2.2	Communication Patterns Analysis	102
6.2.3	KAAP and the Grid'5000 Topology	105
6.3	Treemap Visualizations	110
6.3.1	Description of the Visualization	111
6.3.2	Large Scale Visualizations	112
6.3.3	KAAP Work Stealing Analysis	114
6.3.4	MPI Operations Analysis	119

CCGrid 2009

Submit to
Journal of Grid
Computing

Acknowledgements

- Thesis financed with scholarships by
 - CAPES and CNPq
 - CAPES/Cofecub – Project 4602/06-4
- Thanks to
 - Advisors: Navaux, Guillaume and Denis
 - Family and friends