

XII International Workshop on Advanced Computing and
Analysis Techniques in Physics Research

Job Centric Monitoring for ATLAS jobs in the LHC Computing Grid

Prof. Erich Ehses², Dr. Torsten Harenberg¹, Prof. Peter Mättig¹, Markus Mechtel¹,
Tim München^{1,3,*}, Martin Rau^{1,2}, Prof. Peer Ueberholz⁴, Prof. Nikolaus Wulff³

¹ Bergische Universität Wuppertal, ² Fachhochschule Köln, ³ Fachhochschule Münster, ⁴ Hochschule Niederrhein, * Speaker



Agenda

- Introduction and Motivation
- JEM overview
- JEM architecture
- Usage of JEM
- Current development and outlook



Motivation

- ATLAS data analysis
 - High number of relatively short user jobs
 - Athena: C++ analysis code with python launcher scripts + shellscripts for environment setup
 - Many users, many WNs, high submission frequency
 - Ganga: ATLAS distributed analysis frontend
 - Presented today by Mr Shank
 - Still quite high failure rate during RUNNING state

Motivation

- How to help the user finding the failure reason?
 - Done (non-zero exitcode) jobs
 - Error may be hidden in tons of stderr output
 - Stacktraces may be cluttered with ROOT, Athena, Ganga wrappers
 - Or error may not be to determine at all
 - How to handle „hanging“ Jobs
 - Abort it / Failure due to expiry of time limit or certificate
 - In both cases: sandbox lost

JEM Overview

- Functionality of JEM:
 - Supervised line-by-line execution of user scripts
 - Monitoring of vital system data
 - Transmission of these data to the grid UI
 - In nearly real-time
 - Ready to view as HTML/Graphs
 - As XML for further processing (e.g. view in Ganga)
 - Stored in the job's output sandbox
 - Allowing post-mortem analysis



JEM Overview: Data provided

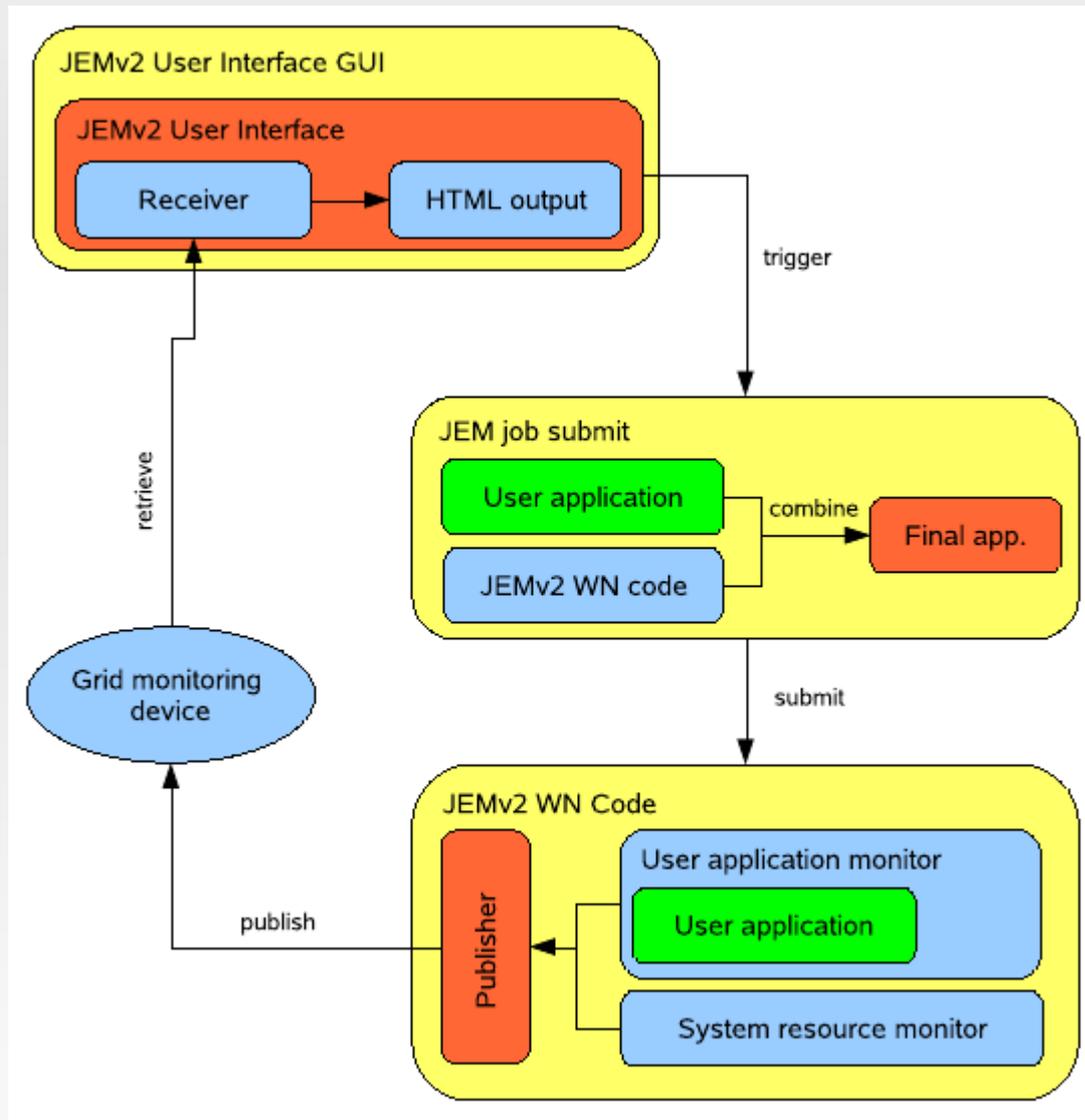
- From Bash-Scripts
 - Executed commands and built-ins (args, exitcode)
- From Python Scripts
 - Calls (args, location, caller), returns
 - Exceptions (type, reason, location, local vars)
- From Binaries (C/C++) (when prepared)
 - Calls (args, location, caller, local vars), returns

JEM Architecture

- Python application run in user space
- Little dependencies
 - Python 2.3.x
 - optional: (py-) R-GMA, (py-) OpenSSL, rrdtool (on the UI)
- Two main parts
 - JEM_UI, run on the UI machine
 - JEM_WN, submitted with the user job



JEM Architecture



JEM Architecture: JEM_UI

- Submission of decorated user job
- Reception and storage of monitoring data
- Flavours
 - Command line tools (Submission only),
 - Interactive shell (Nearly real-time data reception),
 - Ganga-Integrated (Also nearly real-time data reception; preferred way)

JEM Architecture: JEM_WN

- Part of JEM that gets submitted with the user job to the CE
- Runs the user job line-by-line
- Gathers execution and exception data
- Periodically gathers environment data
 - CPU / RAM / HDD usage, Network traffic, ...
- Transmits the data to the UI Part



JEM Architecture: JEM_WN

- *.py monitor
 - Wrapper written in Python
 - Uses python's built-in trace functionality
 - A bit overhead, but very rich information (e.g. all local variable values in case of an exception)
- *.sh monitor
 - Enhanced bash
 - Direct analysis of shell-commands
 - Little overhead

JEM Architecture: JEM_WN

- Binary tracer
 - Shared library written in C
 - Exploits special functionality provided by GCC
 - Must be dynamically linked to the traced object
 - Traced object must contain debug symbols
 - Instruments calls/returns, gathers arguments and local data and feeds these data into JEM
 - Little overhead, big information gain

JEM Architecture: Data transfer

- Monitoring data is transferred from the WNs to the UI:
 - in (almost) real-time
 - R-GMA (deprecated, does not scale well)
 - MonALISA
 - Direct TCP socket connection (debug purposes)
 - HTTP(S) PUT (preferred: encrypted, lightweight)
 - and / or via the output-sandbox

JEM Usage

- Since mid-2008: Tightly integrated job submission with JEM via Ganga
- Handles creation of listener and reception of monitoring data automatically in the background
- Allows for (nearly) real-time access to monitoring data from within Ganga
- If the user quits Ganga during job runtime, data is nevertheless stored and available later

Current development / Outlook

- Scalability of the different data transfer alternatives still gets evaluated (r-gma, HTTP PUT, ...)
- Per-Site aggregation and filtering of data before it is sent to the UI (e.g. on the MON-BOX)
 - Chronological ordering of the data
 - Elimination of duplicate (or related) entries
 - e.g. when 100s of splitjobs fail at the same code location
 - Combination of related data (e.g. Exception stack frames, call/return data only if exception occurred...)

Current development / Outlook

- The use of JEMs monitoring data as (one) input for an expert system (GridXP)
 - Another project in development at the University of Wuppertal
 - Error classification / severity ranking
 - More human-readable output and generated solution suggestions



Current development / Outlook

- „Post-mortem“ analysis of Job runs
 - Export of CSV data
 - For analysis in the usual spreadsheet apps
 - JEM-Provided specialized tool („JEMpole“)
 - Arbitrary filters, Full-text search, Code vicinity viewer
- Improvements to the data transfer
 - e.g. compression

Current development / Outlook

JEMpole - The JEM post-mortem log explorer

File Edit Settings Help

Filter mode

- Hide matching
- Show only matching
- Disable filter

File (only filename) equals simple_cpp.cpp
Type is Exception
<doubleclick to add rule>

Timestamp	File	Frame	Code / Command	Arguments / Return Value / Exc	Local
15:28:38.874	simple_cpp.cpp:20	int MyClass::square	{		
15:28:38.884	simple_cpp.cpp:56	main	a.doStuffWithString(name, x);	{'s':('char*','0x08048b68 [= "... {'a':('	
15:28:38.894	simple_cpp.cpp:28	void MyClass::doStuffWithString	{		{'bar'
15:28:38.904	simple_cpp.cpp:54	main	if (a.square(x) < 10)	{'arg':('int','4')}	{'a':('
15:28:38.914	simple_cpp.cpp:22	MyClass::square	foo(arg*arg, &arg);	{'arg':('int','16'),'ptr':('int*','0x...	
15:28:38.924	simple_cpp.cpp:13	void MyClass::foo	{		{'loca
15:28:38.934	simple_cpp.cpp:20	int MyClass::square	{		
15:28:38.944	simple_cpp.cpp:54	main	if (a.square(x) < 10)	{'arg':('int','5')}	{'a':('
15:28:38.954	simple_cpp.cpp:22	MyClass::square	foo(arg*arg, &arg);	{'arg':('int','25'),'ptr':('int*','0x...	
15:28:38.964	simple_cpp.cpp:13	void MyClass::foo	{		{'loca
15:28:38.974	simple_cpp.cpp:20	int MyClass::square	{		
15:28:38.984	simple_cpp.cpp:62	main	b->doStuffWithString("Poin'te...	{'s':('char*','0x08048b75 [= "... {'a':('	
15:28:38.994	simple_cpp.cpp:28	void MyClass::doStuffWithString	{		{'bar'

Local Variables

Type	Name	Value
char*	name	0x08048b68 [= "Heinrich XII"
int	x	5
MyClass	a	<MyClass instance @ 0xffffc71
MyClass*	b	0x00000001 [= <MyClass inst

Arguments

Type	Name	Value
int	arg	5

Code Vicinity

```
000049 int x = 0;
000050 char* name = "Heinrich XII";
000051
000052 for (x = 0; x < 6; ++x)
000053 {
000054     if (a.square(x) < 10)
000055     {
000056         a.doStuffWithString(name,
000057     }
000058 }
```

Discussion

Thank you for your attention!

More information:

<http://www.grid.uni-wuppertal.de/grid/jms/index.html>

