



TEXAS ADVANCED COMPUTING CENTER

WWW.TACC.UTEXAS.EDU



TEXAS

The University of Texas at Austin

# XALT: Understanding HPC Usage via Job Level Collection

Robert McLay

January 22, 2018

# XALT: What runs on the system

- ▶ A U.S. NSF Funded project: PI: Mark Fahey and Robert McLay
- ▶ A Census of what programs and libraries are run
- ▶ Running at TACC, NICS, U. Florida, KAUST, ...
- ▶ Integrates with TACC-Stats.
- ▶ Has commercial support from Ellexus

# History of XALT

- ▶ Mark Fahey (was NICS, now ANL): ALT-D (MPI only)
- ▶ Robert McLay (TACC) Lariat (MPI only)
- ▶ Reuben Budiardja (was NICS now ORL)
- ▶ ALT-D + Lariat  $\Rightarrow$  XALT 1: (MPI only)
- ▶ XALT 2: All programs

# Design Goals

- ▶ Be extremely light-weight
- ▶ Provide provenance data: How?
- ▶ How many use a library or application?
- ▶ Collect Data into a Database for analysis.

# Understanding what your users are doing

- ▶ What programs, libraries are your users using?
- ▶ What are the top programs by node-hours? by counts?
- ▶ Are they building their own programs or using someone else's?
- ▶ Are Executables implemented in C/C++/Fortran?
- ▶ Track MPI: tasks? nodes?
- ▶ Track Threading via \$OMP\_NUMTHREADS

# Design: Linker

- ▶ XALT wraps the linker to enable tracking of exec's
- ▶ The linker (ld) wrapper intercepts the user link line.
- ▶ Generate assembly code: key-value pairs
- ▶ Capture tracemap output from ld
- ▶ Transmit collected data in \*.json format
- ▶ Adds codes that executes before main() and after main() completes

# Design: Transmission to DB

- ▶ File: collect nightly
- ▶ Syslog: Use Syslog filtering (or ELK)
- ▶ Direct to DB. (Not in XALT 2)
- ▶ Future: RabbitMQ

# Lmod to XALT connection

- ▶ Lmod spider walks entire module tree.
- ▶ Can build a reverse map from paths to modules
- ▶ Can map program & libraries to modules.
- ▶ `/opt/apps/i15/mv2_2_1/phdf5/1.8.14/lib/libhdf5.so.9` ⇒ `phdf5/1.8.14(intel/15.02:mvapich2/2.1)`
- ▶ Also helps with function tracking.
- ▶ Tmod Sites can still use Lmod to build the reverse map.



# Protecting XALT: Python to C++

- ▶ XALT 1 used python scripts
- ▶ It was difficult to protect Python from users in everytime
- ▶ Solution: LD\_LIBRARY\_PATH="@ld\_lib\_path@"  
PATH=/usr/bin:/bin C++-exec ...
- ▶ Everything that depends on PATH must be hard coded

# Using XALT Data

- ▶ Targetted Outreach: Who will be affected
- ▶ Largemem Queue Overuse
- ▶ XALT and TACC-Stats

# Tracking Non-mpi jobs (I)

- ▶ Originally we tracked only MPI Jobs
- ▶ By hijacking mpirun etc.
- ▶ Now we can use ELF binary format to track jobs

# ELF Binary Format Trick

```
void myinit(int argc, char **argv)
{
    /* ... */
}
void myfini()
{
    /* ... */
}
static __attribute__((section(".init_array")))
    typeof(myinit) *__init = myinit;
static __attribute__((section(".fini_array")))
    typeof(myfini) *__fini = myfini;
```

# Using the ELF Binary Format Trick

- ▶ This C code is compiled and linked in through the hijacked linker
- ▶ It can also be used with `LD_PRELOAD`
- ▶ We are using both...

# Challenges (I)

- ▶ Do not want to track mv, cp, etc
- ▶ Only want to track some executables on compute nodes
- ▶ Do not want to get overwhelmed by the data.

# Answers

- ▶ XALT Tracking only when told to
- ▶ Compute node only by host name filtering
- ▶ Executable Filter based on Path
- ▶ Protection against closing stderr before fini.
- ▶ Site configurable!

# Path Filtering

- ▶ Uses FLEX to compile in patterns
- ▶ Use regex expression to control what to keep and ignore.
- ▶ Three files containing regex patterns, converted to code.
- ▶ Accept List Tests: Track `/usr/bin/ddt`, `/bin/tar`
- ▶ Ignore List Tests: `/usr/bin`, `/bin`, `/sbin`, ...



# TACC\_config.py

```
hostname_patterns = [  
    ['KEEP', '^c[0-9][0-9][0-9]-[0-9][0-9][0-9]:* ']  
]  
path_patterns = [  
    ['SPSR', r'*/python[0-9][^/][^/]* '],  
    ['SPSR', r'*/R'],  
    ['KEEP', r'^/usr/bin/ddt'],  
    ['SKIP', r'^/usr/. *'],  
    ['SKIP', r'^/bin/. *'],  
]  
env_patterns = [  
    ['SKIP', r'^MKLROOT=. * '],  
    ['SKIP', r'^MKL_DIR=. * '],  
    ['KEEP', r'^I_MPI_INFO_NUMA_NODE_NUM=. * '],  
]
```

# Protecting XALT (I): UTF8 Characters

- ▶ Linux supports UTF8 Characters in file names, env. vars.
- ▶ Python supports UTF8 if you know what you are doing.
- ▶ Switch XALT to use prepared statements
- ▶ Where query="INSERT INTO table VALUE(?, ?)"
- ▶ This prevent SQL injection: "johnny drop tables;"
- ▶ Also supports UTF8 characters.

# Protecting XALT (II): Python to C++

- ▶ Difficult to protect Python from users in every case
- ▶ Solution: `LD_LIBRARY_PATH="@ld_lib_path@"`  
`PATH=/usr/bin:/bin C++-exec ...`
- ▶ Everything that depends on `PATH` must be hard coded
- ▶ `basename`  $\Rightarrow$  `/bin/basename`
- ▶ Unique install for each operating system.
- ▶ Certain programs aren't in the same place: `basename`

# Speeding up XALT 2

- ▶ Minimal impact on jobs ( $> 0.09$  secs)
- ▶ Non-mpi jobs only produce end record
- ▶ Filter by job\_id after transport.
- ▶ No more than 100 executions per job\_id. (Changable)
- ▶ A launcher job 5 million executions in two hours.

# Tracking R packages

- ▶ XALT 2 can now track R package usage
- ▶ James McComb & Michael Scott from IU developed the R part
- ▶ They do this by intercepting the “imports”
- ▶ Still in preliminary stages
- ▶ Plan to support Python later.

# New program: xalt\_extract\_record

- ▶ This program reads the watermark.
- ▶ Find out who built this program on what machine
- ▶ Find out what modules were used.

# Example of xalt\_extract\_record output

\*\*\*\*\*

XALT Watermark: hello

\*\*\*\*\*

```
Build_Epoch                1510257139.4624
Build_LMFILES              /opt/apps/modulefiles/intel/17.0
Build_LOADEDMODULES       intel/17.0.4:impi/17.0.3:pytho
Build_OS                   Linux 3.10.0-
514.26.2.el7.x86_64
Build_Syshost              stamped2
Build_UUID                 586d5943-67eb-480b-a2fe-35e87a1f2
Build_User                 mclay
Build_compiler             icc
Build_date                 Thu Nov 09 13:52:19 2017
Build_host                 c455-011.stamped2.tacc.utexas.ec
XALT_Version               1.7.7-devel
```

# Conclusion

- ▶ Lmod:
  - ▶ Source: [github.com/TACC/lmod.git](https://github.com/TACC/lmod.git), [lmod.sf.net](http://lmod.sf.net)
  - ▶ Documentation: [lmod.readthedocs.org](http://lmod.readthedocs.org)
- ▶ XALT:
  - ▶ Source: [github.com/Fahey-McLay/xalt.git](https://github.com/Fahey-McLay/xalt.git), [xalt.sf.net](http://xalt.sf.net)
  - ▶ Documentation: [doc/\\* .pdf](http://doc/* .pdf), [xalt.readthedocs.org](http://xalt.readthedocs.org)