

Collectl Intro and Demo

Abhinav Thota
SciAPT
Indiana University



**RESEARCH
TECHNOLOGIES**

INDIANA UNIVERSITY
University Information Technology Services



**PERVASIVE TECHNOLOGY
INSTITUTE**

INDIANA UNIVERSITY

What is Collectl?

- Collectl is a systems monitoring tool
- It can monitor a broad set of subsystems including:
 - buddyinfo, cpu, disk, inodes, infiniband, lustre, memory, network, nfs, processes, quadrics, slabs, sockets and tcp.
- The collected data can be stored in compressed or uncompressed data files
 - which themselves can be in either raw format or in a space-delimited format that enables plotting using gnuplot or Microsoft Excel.
- A KB article is available: <https://kb.iu.edu/d/bedc>



**RESEARCH
TECHNOLOGIES**

INDIANA UNIVERSITY
University Information Technology Services



**PERVASIVE TECHNOLOGY
INSTITUTE**

INDIANA UNIVERSITY

Indirect monitoring of applications

- This sounds like a monitoring tool for system admins
- But, we are using this to monitor some of the subsystems to indirectly monitor/evaluate our application (s) running on the system
 - Instead of profiling or instrumenting the code
- By looking at memory consumption on a compute node, we can tell what the application is consuming
- By looking at CPU utilization, we can again tell what the application is doing
- This can be useful, even if multiple users are using a node
 - We can filter processes by who owns them, for example



**RESEARCH
TECHNOLOGIES**

INDIANA UNIVERSITY
University Information Technology Services



**PERVASIVE TECHNOLOGY
INSTITUTE**

INDIANA UNIVERSITY

On Karst -

- module load gnuplot
 - This is so that we can generate plots at the end of a collectl run
- module load collectl
- Feel free to run “collectl” to get a feel for how this works



**RESEARCH
TECHNOLOGIES**

INDIANA UNIVERSITY
University Information Technology Services



**PERVASIVE TECHNOLOGY
INSTITUTE**

INDIANA UNIVERSITY



Commonly used options

- -Z would be the most obvious option
 - Includes CPU, memory and IO data

Option	Description																				
<code>-F</code>	Flushes the output buffers at the specified interval (in seconds); <code>-F0</code> (zero) causes a flush at each sample interval																				
<code>-i</code>	Indicates how frequently a sample occurs (in seconds) The value preceding the colon (<code>:</code>) is the primary sampling interval. The value following the colon is the rate at which subsystem, subprocess (<code>-sZ</code>), and slab (<code>-sY</code>) data is collected, and must be a multiple of the value before the colon.																				
<code>-s</code>	Determines what subsystem data (summary or detail) is collected; the default is <code>cdn</code> , which stands for CPU (<code>c</code>), disk (<code>d</code>), and network (<code>n</code>) summary data; lowercase characters indicate summary data collection, uppercase characters indicate detail data collection; to get both, include lowercase and uppercase letters (e.g., <code>-s Zl</code>). Options for summary and detailed subsystem data include: <table border="1"> <thead> <tr> <th>Summary subsystems</th> <th>Detail subsystems</th> </tr> </thead> <tbody> <tr> <td><code>b</code> - memory fragmentation</td> <td><code>C</code> - CPU</td> </tr> <tr> <td><code>c</code> - CPU</td> <td><code>D</code> - disk</td> </tr> <tr> <td><code>d</code> - disk</td> <td><code>J</code> - interrupts</td> </tr> <tr> <td><code>j</code> - interrupts</td> <td><code>L</code> - Lustre</td> </tr> <tr> <td><code>l</code> - Lustre</td> <td><code>M</code> - memory node data (or numa data)</td> </tr> <tr> <td><code>m</code> - memory</td> <td><code>N</code> - networks</td> </tr> <tr> <td><code>n</code> - networks</td> <td><code>Y</code> - slabs (system object caches)</td> </tr> <tr> <td><code>s</code> - sockets</td> <td><code>Z</code> - processes</td> </tr> <tr> <td><code>y</code> - slabs (system object caches)</td> <td></td> </tr> </tbody> </table> Note: In the above example script, <code>-s Zl</code> directs <code>collectl</code> to collect detail data for processes (<code>Z</code>) and summary data for Lustre I/O (<code>l</code>).	Summary subsystems	Detail subsystems	<code>b</code> - memory fragmentation	<code>C</code> - CPU	<code>c</code> - CPU	<code>D</code> - disk	<code>d</code> - disk	<code>J</code> - interrupts	<code>j</code> - interrupts	<code>L</code> - Lustre	<code>l</code> - Lustre	<code>M</code> - memory node data (or numa data)	<code>m</code> - memory	<code>N</code> - networks	<code>n</code> - networks	<code>Y</code> - slabs (system object caches)	<code>s</code> - sockets	<code>Z</code> - processes	<code>y</code> - slabs (system object caches)	
Summary subsystems	Detail subsystems																				
<code>b</code> - memory fragmentation	<code>C</code> - CPU																				
<code>c</code> - CPU	<code>D</code> - disk																				
<code>d</code> - disk	<code>J</code> - interrupts																				
<code>j</code> - interrupts	<code>L</code> - Lustre																				
<code>l</code> - Lustre	<code>M</code> - memory node data (or numa data)																				
<code>m</code> - memory	<code>N</code> - networks																				
<code>n</code> - networks	<code>Y</code> - slabs (system object caches)																				
<code>s</code> - sockets	<code>Z</code> - processes																				
<code>y</code> - slabs (system object caches)																					
<code>--procfilt</code>	Tells <code>collectl</code> to get data only for the processes that are specified by the filter parameters In the above example, the <code>--procfilt</code> option indicates that only processes for the <code>\$UID</code> user ID (<code>u\$UID</code>) should be monitored. Other filters and options for specifying what data to monitor and how to record them include <code>--diskfilt</code> , <code>--memopts</code> , <code>--lustropts</code> , and <code>--procopts</code> .																				
<code>-f</code>	Sets the directory for <code>collectl</code> output For jobs on Big Red II, Karst, or Mason, set <code>COLLECTLDIRECTORY</code> to a temporary directory in your Data Capacitor II scratch space (e.g., <code>/N/dc2/scratch/username/temp/</code> ; replace <code>username</code> with your IU Network ID username).																				



As a background process in a batch job

```
module load collectl
cd /N/dc2/scratch/username/temp
SAMPLEINTERVAL=10
COLLECTLDIRECTORY=/N/dc2/scratch/username/temp/
collectl -F1 -i$SAMPLEINTERVAL:$SAMPLEINTERVAL -sZ --procfilt
u$UID -f $COLLECTLDIRECTORY &
./my_binary
collectl_stop.sh
```

- You should adjust SAMPLEINTERVAL according to the expected runtime for your application. If your job will run for less than a day, [UITS](#) recommends setting SAMPLEINTERVAL=10; if your job will run for multiple days, UITS recommends setting SAMPLEINTERVAL=30 or SAMPLEINTERVAL=60.



**RESEARCH
TECHNOLOGIES**

INDIANA UNIVERSITY
University Information Technology Services



**PERVASIVE TECHNOLOGY
INSTITUTE**

INDIANA UNIVERSITY

CAFÉ

- CAFÉ is the application we will use today to demonstrate Collecti and other tools
- Used to analyze changes in gene family size and provide a statistical foundation for evolutionary inferences.
 - <https://hahnlab.github.io/CAFE/>
- Please copy this directory to your \$HOME:
 - /N/dc2/scratch/athota/incoming/cafe
 - Unzip v4.0.zip and run "make" to build the binary
 - The binary will be built in the CAFE-4.0/release directory
 - From the top level directory, to run cafe:
 - CAFE-4.0/release/cafe caferun.sh



**RESEARCH
TECHNOLOGIES**

INDIANA UNIVERSITY
University Information Technology Services



**PERVASIVE TECHNOLOGY
INSTITUTE**

INDIANA UNIVERSITY

Karst Desktop Beta

- Linux remote desktop service based on VNC
- We will be using this service for the hands-on part of the presentation
- The KB article: <https://kb.iu.edu/d/bfwp>
- TL;DL:
 - <https://www.cendio.com/thinlinc/download>
 - Server: test-desktop.karst.uits.iu.edu
 - Password: your IU CAS password
 - Jobs can be submitted to Karst from the desktop nodes



**RESEARCH
TECHNOLOGIES**

INDIANA UNIVERSITY
University Information Technology Services



**PERVASIVE TECHNOLOGY
INSTITUTE**

INDIANA UNIVERSITY

Running CAFÉ with Collectl

```
#!/bin/bash
cd /N/u/athota/Karst/perf-tuning/cafe
collectl -F1 -i5:5 -sZ --procfilt u1224111 -f collectl_log/ &
CAFE-4.0/release/cafe caferun.sh
killall collectl
sleep 10
killall perl
sleep 10
```

- This writes a log file to the collectl_log directory in this format:
 - hostname-YYYYMMDD-HHMMSS.raw.gz



RESEARCH
TECHNOLOGIES

INDIANA UNIVERSITY
University Information Technology Services



PERVASIVE TECHNOLOGY
INSTITUTE

INDIANA UNIVERSITY



Plotting the data

- Can be done with the `collectl_plot.sh` script that is included on IU systems
- From the `collectl_plot` directory, run:
 - `collectl_plot.sh hostname-YYYYMMDD-HHMMSS.raw.gz .`
 - Don't forget the "dot" at the end :)
- There is a lot more data that we are not plotting in the log file
 - Can adjust gnuplot scripts to change the plots



**RESEARCH
TECHNOLOGIES**

INDIANA UNIVERSITY
University Information Technology Services



**PERVASIVE TECHNOLOGY
INSTITUTE**

INDIANA UNIVERSITY

