

SANDIA REPORT

2006-xxxx

Unlimited Release

Printed November 2006

Monitoring Computational Clusters with OVIS

J. M. Brandt, A. C. Gentile, P. P. Pébay, and M. H. Wong

Prepared by
Sandia National Laboratories
Albuquerque, New Mexico 87185 and Livermore, California 94550

Sandia is a multiprogram laboratory operated by Sandia Corporation,
a Lockheed Martin Company, for the United States Department of Energy's
National Nuclear Security Administration under Contract DE-AC04-94-AL85000.

Approved for public release; further dissemination unlimited.



Sandia National Laboratories

Issued by Sandia National Laboratories, operated for the United States Department of Energy by Sandia Corporation.

NOTICE: This report was prepared as an account of work sponsored by an agency of the United States Government. Neither the United States Government, nor any agency thereof, nor any of their employees, nor any of their contractors, subcontractors, or their employees, make any warranty, express or implied, or assume any legal liability or responsibility for the accuracy, completeness, or usefulness of any information, apparatus, product, or process disclosed, or represent that its use would not infringe privately owned rights. Reference herein to any specific commercial product, process, or service by trade name, trademark, manufacturer, or otherwise, does not necessarily constitute or imply its endorsement, recommendation, or favoring by the United States Government, any agency thereof, or any of their contractors or subcontractors. The views and opinions expressed herein do not necessarily state or reflect those of the United States Government, any agency thereof, or any of their contractors.

Printed in the United States of America. This report has been reproduced directly from the best available copy.

Available to DOE and DOE contractors from
U.S. Department of Energy
Office of Scientific and Technical Information
P.O. Box 62
Oak Ridge, TN 37831

Telephone: (865) 576-8401
Facsimile: (865) 576-5728
E-Mail: reports@adonis.osti.gov
Online ordering: <http://www.doe.gov/bridge>

Available to the public from
U.S. Department of Commerce
National Technical Information Service
5285 Port Royal Rd
Springfield, VA 22161

Telephone: (800) 553-6847
Facsimile: (703) 605-6900
E-Mail: orders@ntis.fedworld.gov
Online ordering: <http://www.ntis.gov/ordering.htm>



Monitoring Computational Clusters with OVIS

J. M. Brandt

Sandia National Laboratories
M.S. 9152, P.O. Box 969
Livermore, CA 94550, U.S.A.
brandt@sandia.gov

A. C. Gentile

Sandia National Laboratories
M.S. 9152, P.O. Box 969
Livermore, CA 94550, U.S.A.
gentile@sandia.gov

P. P. Pébay

Sandia National Laboratories
M.S. 9051, P.O. Box 969
Livermore, CA 94550, U.S.A.
pppebay@sandia.gov

M. H. Wong

Sandia National Laboratories
M.S. 9152, P.O. Box 969
Livermore, CA 94550, U.S.A.
mhwong@sandia.gov

Abstract

Traditional cluster monitoring approaches consider nodes in singleton, using manufacturer-specified extreme limits as thresholds for failure “prediction”. We have developed a tool, OVIS, for monitoring and analysis of large computational platforms which, instead, uses a statistical approach to characterize single device behaviors from those of a large number of statistically similar devices.

Baseline capabilities of OVIS include the visual display of deterministic information about state variables (*e.g.*, temperature, CPU utilization, fan speed) and their aggregate statistics. Visual consideration of the cluster as a comparative ensemble, rather than as singleton nodes, is an easy and useful method for tuning cluster configuration and determining effects of real-time changes.

Contents

1	Introduction	7
2	Installation	7
2.1	Instructions for compiling Qt 4.x.y for Linux	7
2.2	Setting the environment	8
2.3	Instructions for compiling OVIS 1.1 on Linux	9
2.4	OVIS ancillary files	9
3	Data Collection	9
3.1	Data File Format	10
3.2	Data Collection	11
4	Operation Guide for OVIS 1.1	14
4.1	Description of menus and controls	14
4.2	Cluster Definition and Representation	15
4.3	Getting Data Into the Display	16
4.4	The OVIS Display	16
4.5	Customizing Color Scales	17
4.6	Use Of Data Interpolation	18
5	Tools	18
5.1	Statistics	18
5.2	Correlation	19
6	Calibration Code	19
7	Cluster Analysis	20
7.1	Viewing Raw Data	20
7.2	Viewing Statistical Outliers	21
7.3	Looking for Correlations	21
8	OVIS Tutorial	21

This page left intentionally blank

Monitoring Computational Clusters with OVIS

1 Introduction

This document describes OVIS: A tool for intelligent analysis and monitoring of large computational clusters. Through an intuitive user-friendly GUI, OVIS provides 2D visualization capabilities of cluster health parameters as well as a number of statistical tools for analysis of real-time (monitoring) or stored (diagnostic) data. Moreover, owing to the statistical similarity and the large number of nodes in a cluster, the user can examine the health of each individual node using the statistics tools provided with OVIS together with the ability to bin variable values by color. This allows at-a-glance detection of abnormalities in a system or node depending on how the user has tuned the interface. Furthermore, the correlation engine provided with the distribution allows the user to analyze statistical relationships between different parameters of the system. Knowledge of these relationships can then greatly reduce the parameter space that the user needs to actively monitor in order to detect abnormal node behavior.

This is in contrast to current cluster monitoring methodologies in which data is obtained from each node and a pre-defined rule set is applied, on a per-node basis, to any node whose value(s) cross a pre-determined threshold(s). Though this methodology is well suited to single nodes and small clusters, OVIS uses the statistical properties of these large collections of statistically similar devices to add a great deal of intelligence to the process of monitoring and analysis as well as being able to determine many problems sooner than is possible using static thresholds.

2 Installation

2.1 Instructions for compiling Qt 4.x.y for Linux

Download the latest version of the Qt 4.x.y source from:

<http://www.trolltech.com/developer/downloads/qt/x11>

Select **qt-x11-opensource-src-4.x.y.tar.gz** under "Download from ftp.trolltech.com"

Move the tar file to an appropriately named directory e.g. ~/qt:

```
mkdir ~/qt
```

```
mv < path_to_file >/qt-x11-opensource-src-4.x.y.tar.gz ~/qt
```

Untar the qt-x11-opensource-src-4.x.y.tar.gz file:

```
cd ~/qt
tar xvfz qt-x11-opensource-src-4.x.y.tar.gz
```

Generate a Makefile:

Move to the source directory

```
cd qt-x11-opensource-src-4.x.y
```

If you are building this in your home directory

```
./configure -prefix ~/qt/qt-x11-opensource-src-4.x.y
```

This will put bin, lib, doc, etc. in ~/qt/qt-x11-opensource-src-4.x.y

If you are building this as a system resource

```
./configure (This will put bin, lib, doc, etc. in /usr/local/Trolltech/Qt-4.x.y)
```

Compile and install Qt 4.x.y according to the instructions given when **configure** completes. **Note: You must have root privileges if this is a system build**

2.2 Setting the environment

Set \$QTDIR to the directory to which you have installed Qt. The default Qt 4.x.y installation directory is /usr/local/Trolltech/Qt-4.x.y.

In the default case:

If you are using tcsh, do the following:

```
setenv QTDIR /usr/local/Trolltech/Qt-4.x.y
```

If you are using bash, do the following:

```
export QTDIR=/usr/local/Trolltech/Qt-4.x.y
```

In the case you installed into ~/qt/qt-x11-opensource-src-4.x.y:

If you are using tcsh, do the following:

```
setenv QTDIR ~/qt/qt-x11-opensource-src-4.x.y
```

If you are using bash, do the following:

```
export QTDIR=~/qt/qt-x11-opensource-src-4.x.y
```

Add \$QTDIR/bin to your system path.

If you are using tcsh, do the following:

```
set PATH=($QTDIR/bin $PATH)
```

If you are using bash, do the following:

```
export PATH=$QTDIR/bin:$PATH
```


Note: Make sure that \$QTDIR/bin comes first, to make sure that you are not using the built-in Qt executables, which may be the incorrect version of Qt.

2.3 Instructions for compiling OVIS 1.1 on Linux

If you havent already done so, compile Qt 4.x.y for Linux. See **”Instructions for compiling Qt 4.x.y for Linux”**

Untar the OVIS 1.1.x source:

```
tar xvfz ovis-1.1.x-src.tgz
```

Move to the ”src” directory:

```
cd ovis-1.1.x/src
```

Generate the Makefile using qmake:

```
qmake
```

Compile OVIS 1.1:

```
make
```

Launch OVIS

```
./Ovis
```

2.4 OVIS ancillary files

The ancillary files needed for meaningful display by OVIS are:

1) An XML description of the geometry of the machine being analysed

This should reside in the ovis-1.1.x/geometries directory to be seen using the default configuration but can reside anywhere and be defined on the General tab in the **Edit**—>**Preferences** pull down menu.

2) Data from the machine in OVIS format for display and analysis

The directory in which this data resides needs to be specified on the General tab in the **Edit**—>**Preferences** pull down menu.

3 Data Collection

OVIS 1.1 reads its data from files in the OVIS data format. Data collection mechanisms are not part of OVIS. For convenience, however, we do include in the OVIS release a few example scripts of how data can be obtained from typical sources and written into the OVIS

data format. These scripts are stand-alone codes and do not interact with OVIS itself.

3.1 Data File Format

OVIS 1.1 reads its data from files in the OVIS data format. This requires that a data directory exist containing data files, one for each compute node of the cluster. While the data directory can take any name, the data files must be named `proc.<hostname>.out`, for example, `proc.cn105.out` for `hostname = cn105`. Each data file consists of chronological lines, one for each timestep, containing the data values for each quantity being monitored. The quantities must be in the same order each time and the timestamp must be the first value in the line. The identities of the quantities are specified in a single header file called `header` that resides in the data directory. The quantity specifying time must be called `Time` and its units are in seconds since the epoch; all other labels (and units) can be arbitrary and are read in by OVIS when the data files are loaded.

For example:

```
> more header
Time CPU1_Temp CPU2_Temp FAN1 FAN2 RX TX
```

```
> more proc.cn105.out
1112340000 45 40 12 12 17 31
1112340002 45 45 12 12 13 22
1112340004 45 50 12 12 20 100
1112340005 45 54 12 20 100 100
1112340007 45 52 12 20 100 100
```

```
> more proc.cn106.out
1112340000 42 25 12 12 25 1000
1112340003 43 25 12 12 33 1000
1112340004 42 25 12 12 100 1000
1112340006 42 25 12 12 300 2000
1112340008 42 25 12 12 200 1000
```

The times do not have to be equally spaced, nor do times in different data files have to have a one-to-one correspondence with each other.

Notes:

- OVIS will be changing to a binary format, rather than ascii. A converter will be supplied to convert between the two formats.
- Currently, we keep sequential directories holding the data files, e.g., `main_dir/0`, `main_dir/1`, etc. each having a set of data files called `proc.<hostname>.out`. We are debating options for allowing for sequential (in time) files all in the same dir, with a number extension on the file name, e.g., in `main_dir` is `proc.<hostname>.out.0`, `proc.<hostname>.out.1` etc. In either case, there is currently no mechanism for OVIS to automatically roll over from one

file to the next one in the sequence, however, we intended to support this in a future release.

- The header file cannot have extraneous whitespace and must have single space delimiters. Tabs are not recognized. This will be fixed in a future release.
- The data files must use single space whitespace as a delimiter, and extra single space whitespace is allowed. Tabs are not recognized. This will be fixed in a future release.

3.2 Data Collection

It is up to the user to arrange for data collection and conversion to the OVIS data format. Most systems have onboard sensors and counters with a well known interface from which the information can be obtained. For convenience we provide example scripts for some of these formats that can be customized to fit your scenario. Currently we address reading from /proc or hplog and reading from RRDs, such as those populated by Ganglia. Supermon and IPMI readers are in the works.

3.2.1 Reading from /proc, hplog etc

Summary: The gatherdata.c code reads and writes data for one node. It can be launched to all nodes using a utility like pdsh or it can be started via cron. However it is launched, a check should periodically be made to ensure that if it dies it is restarted.

Details: The gatherdata.c code gathers data for a single node on our Shasta Cluster. It is customized for calls to 1) hplog to get CPU temperatures for the two processors and fan speeds of the processor zone and the system board, 2) /proc to get user, system, and idle utilizations for the two CPUs, and 3) ifconfig eth to get received and transmitted bytes. The calls and the output of the calls are in a known format from which the data is extracted. (Example output from the calls is given at the top of the files, so that one can understand the parsing of the results). The script is run on each compute node. The collection frequency is specified in the script.

The user should customize this code for his own data collecting interface and metrics of interest. This code should probably have as small a footprint as possible so as not to compete with the user processes on the nodes. Thus it is reasonable to take advantage of the known output formats of the data being collected.

3.2.2 Reading from Ganglia populated RRDs

Summary: ovis_RRD_reader reads and writes data for one node; ovis_RRD_launcher launches ovis_RRD_readers, one for each node. ovis_RRD_killer can be used to kill ovis_RRD_readers.

Details: The ovis_RRD_XXX scripts work together to read from RRDs, such as those populated by Ganglia, via the RRDtool perl interface (reference RRDtool). They do not

interact with Ganglia itself. The assumption is that there exists a main directory containing subdirectories, one for each compute node. Within those subdirectories are the RRDs for that node, one for each different variable to be monitored. Each RRD could have multiple timescales being considered (e.g., for data taken on 5 second intervals, the RRD may contain: 1 single datapoint AVERAGE, 2 minute AVERAGE, 10 minute AVERAGE, etc.), but it must correspond to one quantity only (e.g., `cpu1_temp`). The rrd names are assumed to be `<variable>.rrd`, e.g., `cpu1_temp.rrd`.

The `ovis_RRD_reader` will read from the RRDS for a single node. It will extract info from the RRDS for that node, and determine for each RRD, what the minimum measurement interval is (e.g., in the example above it will discern that the 1 datapoint average is the shortest interval (taken on 5 second intervals). It will then compare across all the RRDS and only collect data from those whose shortest interval corresponds to the global shortest interval. For example, if 5 seconds is the shortest interval across all RRDs, then data will only be collected from RRDs who have a quantity whose shortest interval is 5 seconds. In theory, this could result in totally different things being collected (e.g., the 5 second AVERAGE of `cpu1_temp` with data taken on one second intervals and the 5 second XXX of `fan_speed` with data being taken on 5 second intervals), but in practice one usually sets up all the Ganglia RRDS with the same time granularity. Also there is no reason that any two things have to be related. The only issue is if you might want a few different metrics for the same quantity (e.g., the 1 minute AVERAGE and the 10 minute AVERAGE of `fan_speed`), however, there are so many special cases of this that the user should just alter the scripts to fit his special circumstances.

These readers produce the data files, and optionally the header file. If the header file is being produced by the `ovis_RRD_reader` then it will check all possible RRDS for inclusion in the set. If an existing headerfile is specified on the command line, then only those quantities will be considered for inclusion in the set.

Command line syntax for the `ovis_RRD_reader`:

Usage: `./ovis_RRD_reader.pl <parentcndir> <cndir> <outputdir> [options]`

Required arguments:

`<parentcndir>` directory holding the directories of all the cn.

`<cndir>` the directory of the particular cn (relative to the `parentcndir`)

`<outputdir>` output directory

Options:

`-t <starttime> <duration>` starttime can be integer or relative to now i.e, now-300 duration = 0 runs infinitely

`-s <outputfilesize>` outputfilesize, if no flag, defaults to infinite

`-useheader` if the `useheader` flag exists, it will use the headerfile called `header` in the `outputdir`

`-useparameters` if the `useparameter` flag exists, it will use the `queryparameterfile` called `parameters.in` in the `outputdir`

The parameter file contains optional values for parameters related to the data collection algorithm. These provide constraints on how many timesteps the code will wait for new

data before moving on and how old data can be before it is considered stale. These are in the getParams function, their defaults and roles are:

agewindow_steps = 1 Values will be considered stale if the obtained time is older than targettime - agewindow_steps.

maxwait_steps = 4 Longest amount of time after the target time we are willing to wait for valid data. Otherwise, assume we will get nothing valid and move on....

leadtime_steps = 2 How much ahead of the target time we will ask for, in order to see if subsequent data has come in – in order to discover that a point may never come in

backlook_steps = 3*agewindow_steps; How much behind a target time we will ask for to use as potential data

badval = -1; What will be written to the file in case no valid data for a rrd at a given time is found

The ovis_RRD_launcher launches readers, one for each subdirectory of the main directory - again, these correspond to compute nodes. One can specify subdirectories to include or exclude in a file specified on the command line. See the code Usage for this. The ovis_RRD_launcher coordinates the start times of all the ovis_RRD_readers, if a start time is not specified. It can use a predefined header file, as mentioned, with respect to the ovis_RRD_reader or it will create a header that is the union of all the subdirs rrd names.

If different compute nodes have different intervals, then the output files won't have corresponding times. In practice, this is usually not an issue as one typically coordinates the setup of all the RRDS and the times are regular within the RRD (see RRD documentation on how time is handled within an RRD).

Usage: **./ovis_RRD_launcher.pl** <parentdir> <outputdir> [options]

Required arguments:

<parentdir> Directory holding the directories of all the cn's.

<outputdir> Directory where the output will be written Options:

-t <starttime> <duration> Start time. starttime can be integer or relative to current time i.e, <current time> - 300 duration = 0 runs infinitely

-s <outputfilesize> Output file size. If no flag, defaults to infinite

-useheader If the useheader flag exists, it will use the header file called header in the output directory

-usesubdirs|nousesubdirs Specifies whether or not to use the nodes listed in the file subdirs.in in the output directory

The ovis_RRD_launcher starts the ovis_RRD_readers each as a separate processes and then exits itself. The ovis_RRD_killer script will kill all processes with "ovis_RRD_reader" in the name for ease of killing everything off.

NOTES: As mentioend earlier, the algorithmic parameters of the ovis_RRD_reader (e.g., agewindow_steps) have defaults or can be specified in a parameter file. In the future these values will be determined more intelligently by considering the usual frequency of updates, etc. The exact methodology for doing this has yet to be determined.

4 Operation Guide for OVIS 1.1

4.1 Description of menus and controls

This subsection contains a comprehensive list of the menus and controls available via the OVIS GUI.

File -- > **Open Dataset** – Specify the **directory** of a particular data set.
Exit – Exit OVIS

Edit -- > **Preferences [General tab]** – This is where the user specifies defaults for geometry, data results and system types.

Preferences [Variable Scales tab] – This is where the user sets up the value to color bindings.

View -- > **Use Gradient Colormap** – Interpolate between colors (bin if not checked). (See "Customizing Color Scales" section for more information.)

Interpolate Missing Data – Options for displaying data in a given time range about a time point if no data exists for the exact time point. If no data available in the selected range, no data (gray) is shown for this time step. (See "Use Of Data Interpolation" section for more information.)

Tools -- > **Statistics and Correlation** – Allows the selection of various statistical analysis packages

Time -- > **Reset Time Boundaries** – Reset time boundaries to beginning and end specified in the data files. Use this if user has re-defined boundaries and wants to go back to whole file.

Help -- > **About Ovis** – Gives version

System: – Select system to display. The name(s) displayed in this combo box are the <descriptive name> part of the <descriptive name>.geom.xml in the geometries directory.

Variable: – Variable name that is currently being displayed on the GUI. This comes from the header file located in the data directory in which the variable names correspond positionally with data in the data files.

Speed: – This dial allows the user to adjust the rate at which stored data are being read from data files. It's extremes are zero delay (limit on the rate at which the display changes is only limited by disk I/O and processing capacity) and approximately 10 minutes between time steps. The dial is log based with faster being in the clockwise direction and the default speed being one second between time steps. This control allows the user to scroll quickly through data or slowly examine the changes. The stop button can be used if one wants to spend even more time on a particular time step. If the "Sliding End" box is checked and the end of the data file is reached, this dial will automatically be adjusted to match the variable

data update frequency.

Time: – Displays the current time in the files. Since there are typically many files this represents the latest time for a given time step.

Go button – Used to position the display at a user defined time within the file if it is in bounds.

From: – Is either the earliest data point across all data files or a user defined beginning of an analysis interval.

To: – Is either the latest end data point across all data files or a user defined end of an analysis interval.

Sliding End check box – Tells OVIS that data is still being written to the file and to keep checking even though it has reached the end.

Sampling Period: – Eliminates the processing/IO bottleneck in scrolling through data by allowing the user to sample on the user defined time interval (in seconds?) displayed in the spinbox. In this case the speed dial defines the time between successive samples.

Play button – Continues to step through the data, from the current time, in a manner consistent with all other controls.

Stop button – Stops stepping through the data but leaves the display set according to the data at the current time as defined in the Time box in the upper right hand corner.

To beginning button – Takes the display directly to the data points defined by the time displayed in the From box in the lower left of the display.

To end button – Takes the display directly to the data points defined by the time displayed in the To box in the lower left of the display.

4.2 Cluster Definition and Representation

The first thing OVIS needs is a spatially meaningful description of the cluster it is to represent. This must be defined in an xml file whose name is <descriptive name>.geom.xml. A template of this file is located in the ovis-1.1.x/src/templates directory. Once this has been written OVIS must be made aware of it by selecting the directory containing it and possibly other geometry definitions via the Edit-->Preferences-->[General Tab]Default Geometry Location. All names for which there exists an xml file in the Default Geometry Location will be available in the System combo box in the upper left of the display. Choose the system for which you will be viewing data by choosing its descriptive name from the System selector in the upper left of the display. Note that currently it is up to the user to ensure that the data being loaded corresponds to the system selected as there is no checking to make sure there is a match between the number of data files and the number of nodes

defined in a geometry file.

4.3 Getting Data Into the Display

The next thing needed is a data source to apply to the chosen cluster description. If there isn't any data saved or being saved for the system in the OVIS format (see section on Data Collection) you must first perform this step. To access a dataset currently being written (active) or an old dataset that has been closed, choose File-->Open Dataset from the pull down menu in the upper left of the display to launch a file browser with which you can locate and specify the appropriate dataset. Note that you are selecting a directory containing data files, not the files themselves. If you are choosing an active data set and want to watch real time updates you will need to check the Sliding End box located in the lower left of the display after To. Once a dataset has been chosen it will be automatically loaded. If a directory is chosen that contains no data files a pop up dialog box will appear informing user of the error and his choices. The beginning and end times of the file will be displayed as From and To in the lower left section of the display. These fields are editable so that the user can choose to view/analyze a subset of the dataset. Note that when the Sliding End box is checked the end time is grayed out and cannot be edited. If the user wishes to change the end of an analysis interval within an active file the Sliding End box will need to be unchecked.

4.4 The OVIS Display

The OVIS display depicts spatially correct scaling with respect to the vertical dimension. In the horizontal dimension the width is currently a function of how many racks (whether populated with compute equipment or not) are to be displayed in a plane (we recommend ≤ 20). The color displayed represents the value of the variable chosen in the Variable combo box (upper left of the display), the value-to-color mapping chosen for that variable in the Variable Scales tab in the Edit-->Preferences pull down menu (upper left of the display) and the status of the View-->Use Gradient Colormap checkbox. The default numeric values for all variables correspond to currently reasonable Celsius temperature values for a CPU and may or may not be reasonable for any particular variable on your system so should be reviewed in light of actual data. In order to quickly discover reasonable value limits for a given variable, use the statistics tool (Tools-->Statistics). After making the variable in question active by selecting it in the Variable combo box, set a time interval and node, rack, or system and click OK. When data is returned you will have the minimum and maximum values for the data (over the interval). Note that this can take quite some time if the time interval is long and you are looking at the whole system. Though there is no guarantee that these are the absolute bounds, it is a good place to start. See the Customizing Color Scales subsection for details on adjusting the color map to suit your needs and Statistics tool subsection for more detail on proper operation of this tool.

This organizational and display technique allows for intuitive visualization of patterns across a large number of nodes simultaneously. It is expected that given a homogeneous environment and workload the nodes will exhibit similar behavior with respect to most variables. Hence, given a similar workload across a set of nodes, the expectation is a largely homogeneous (in color) display across those nodes. Deviation from this and especially patterns indicate lack of such homogeneity and possible problems. Note that one can force the display to look homogeneous even in the presence of grossly varying data by improper setting of the color maps so care should be taken to insure that these are set correctly for your data.

4.5 Customizing Color Scales

In order to customize the value-to-color mappings choose the Variable Scales tab from the Edit—>Preferences pull down menu (upper left). In this tab there are five value fields associated with five color fields. This menu allows the user to choose value to color mappings that make sense for him in the context of the variable being viewed. There are five value spinboxes associated with five color boxes. Left clicking on a color box gives access to a color palette from which to choose a color to associate with the corresponding value spinbox. Linear interpolation is used to generate colors for values between two consecutive spinboxes. The default colors are Blue and Red for the extremes with magenta, green and yellow as intermediates. The default numeric values for all variables correspond to currently reasonable Celsius temperature values for a CPU and may or may not be reasonable for any other variable so should be reviewed in light of actual data. There is a persistent popup, that stays in the foreground, of the color map for the current working variable. This value-to-color map applies to the working variable for all nodes within the cluster so that a visual comparison can be made based on color alone. These fields can be used in several ways which are described in order of simplest to most complex:

- 1) You can choose a smooth gradient between two extremes of a variable by setting a lower bound and an associated color in one of the end fields and an upper bound and associated color in the adjacent fields and checking the gradient box in the main display upper center.

- 2) If the desire is to give finer granularity to the values in a certain region you can define up to four such regions. The color gradient will always be from a left value/color pair to the right adjacent value/color pair. The leftmost color refers to the leftmost value field and values to the left of it (lower numerical values) and similarly the rightmost color refers to the rightmost value field and values further to the right (higher numerical values). Colors in this case scale smoothly in hue from one color field to the adjacent color field as the value of the variable moves between the associated value fields.

- 3) Another option is to bin values in order to better visualize where values lie on a much coarser scale. Currently OVIS allows definitions of five color bins. Each bin contains the values starting from and including that in a variable box and up to (but not including) the value in the variable box to the right. A bin's color is associated with the range that begins

with the value of the variable field above it. There are two special cases: 1) The leftmost bin additionally contains all values to the left of (less than) it. 2) The rightmost bin begins with the value of its associated variable field and contains all values to the right of (greater than) it.

4.6 Use Of Data Interpolation

Options for data interpolation are implemented in OVIS because data returned for a given request over a large number of nodes will have some time skew (order of a few seconds) and hence will not be displayed simultaneously even though they represent the same region in time. There are three types of options: 1) only data for the current time is displayed, if available, and no data is displayed if there is no data available for this time 2) data for the nearest point in time within a given time range is displayed if available, and no data is displayed if there is no data available for any of the times within the range 3) data for the current time is displayed if available and, if no data is available or a given time, the last valid value is displayed. In case 2, the range is centered about the given time, that is, a choice of $t \pm X$ seconds is made, while in case 3, the range includes only previous times. Also in case 2, a few options are given for convenience or an option to enter in a custom time in seconds is available.

Though the use of interpolation makes for a nice continuous display and is valid for some types of data over some time periods, it must be used cautiously with understanding of the implications. For instance due to the thermal mass, the temperature of a CPU will not change appreciably over a few second interval. Thus if the data collection rate is one sample per second it may make sense to carry a missing point forward a time step or two whereas for a data collection rate of one sample per minute it wouldn't. For variables such as network transmits and receives there is no inertia but due to the nature of data gathering and applications sending in bursts it may still make sense if the interval is short enough. **Use with caution.**

5 Tools

The tools currently available in OVIS are described below.

5.1 Statistics

This tool allows the user to choose a single node, a rack of nodes or all nodes in the system and a time period over which to analyse a specified variable with respect to the following basic statistics: sample size, minimum, maximum, mean, median, mode, unbiased variance, sample variance, sample skewness and sample kurtosis. This tool gives the user

insight not only into a particular variable's behavior with respect to a given node but can also yeild insight about phenomena affecting a variable with respect to spatial and temporal variation.

Note that node and rack IDs must be entered as they appear in the XML geometry file (and therefore also in the User Interface). For example, if a cluster has nodes labeled from `cn-1` through `cn-100`, then, in order for node 10, the user must enter `cn-10` (and not only 10). Blank spaces must be entered as well, when present in the ID. Note also that a common mistake is related to the case (*e.g.*, `rack` instead of `Rack`).

The results as calculated by the Statistics tool can be saved into a text file. The prefix of the filename is set, by default, to `OvisStatistics-`. A counter is incremented after each save operation, and the current value is appended to the filename prefix. The dialog allows the user to modify the filename prefix.

5.2 Correlation

This tool allows the user to choose a single node, a rack of nodes or all nodes in the system and a time period over which to analyse two specified variables for linear dependencies. The following information is calculated and displayed: Variables being evaluated, sample size, means of the two variables, sample variance of each variable, covariance, variable 2 on 1 regression slope, variable 2 on 1 regression intersect, variable 1 on 2 regression slope, variable 1 on 2 regression intersect and the linear correlation coefficient. This tool can give the user insight into the linear dependencies between sets of variables which in turn may allow inference about cause and effect relationships within a cluster.

The same syntactic rules as for the Statistics tool apply. Note that the filename prefix is set by default to `OvisCorrelation-`.

6 Calibration Code

A calibration code is supplied as a convenience for the user but is not necessary to the operation of OVIS. The `dialutil.c` code is used to impart a known CPU Utilization to the nodes. This value is specified by the user on the command line (See the code Usage). This code can then be used to calibrate things such as the Temperature-change vs change-in-utilization values. For instance, on our cluster Shasta, we find this is generally a constant, allowing one to factor out the CPU Utilization when considering the temperature profile, thus allowing comparison across the cluster, even if the cluster is running different jobs.

As an example, one can simultaneously run:

```
./util/dial/dialutil 30 120 temp           (starts the code, writes to dir ./temp)
./util/gatherdata temp                   (gathers data)
```

tail -f ./temp/proc*.out

(you can see the CPU util in the output).

7 Cluster Analysis

7.1 Viewing Raw Data

Viewing raw data in gradient mode over many nodes layed out in a spatially correct fashion can yeild insight about the environment where "environment" here refers to any external influence(s) on the nodes. It can also assist in the detection of gross outliers but these are covered in the next section. In this section we will address the detection of thermal, radiation, and network environmental factors though there are many more.

7.1.1 Thermal Influence

This is perhaps the most obvious. Color gradients on a display of the raw temperature values for any node parameter while nodes within the display area are performing similar tasks are likely due to spatial temperature gradients in the external environment. The strength of such gradients will vary depending on the parameter being viewed and the relative influence of external temperature on that parameter. Thus one would expect a linear relationship between intake temperature and any other temperature variable within the node. The relationship between intake temperature and memory error rates may be non-linear and may not even be apparent over the normal spatial temperature distribution.

7.1.2 Radiation Influence

Finding color gradients accross nodes when viewing memory error related parameters can be indicative of thermal gradients as mentioned above but also can be, among other things, due to cosmic radiation gradients. This would be expected in high altitude installations such as Los Almos and Albuquerque where a cluster might have large enough extent that some nodes sit in areas of higher shielding than others. Of course to see these effects the display must simultaneously display all of the nodes over which the gradient occurs.

7.1.3 Network Influence

Depending on the technologies of the networking equipment being employed in a particular cluster, color gradients may occur over adjacent nodes network parameter values when the connections to the networking gear are made to adjacent network gear ports or bundles of such connections take different paths to the network gear. These gradients could be due to

problems in a particular group of interfaces or thermal gradients across the blades of such gear.

7.2 Viewing Statistical Outliers

The Statistics tool not only allows intelligent choice of bounds for a given parameter but also gives information about the distribution of values for a particular parameter over a set of nodes. This information allows the user to create color bins defined by the relative probability of the values within that bin occurring. For instance: The statistics tool for a specified interval of time over a statistically meaningful number of nodes running with the same operational parameters says that the values for a given parameter are normally distributed with a mean of μ and a standard deviation of σ . Binning colors such that values greater than $(\mu + 2 * \sigma)$ are red and values less than $(\mu - 2 * \sigma)$ are blue allows the user to detect at a glance, nodes for which the value of this parameter has less than five percent probability. The spatial distribution of such nodes gives information as to the underlying cause. A sparse seemingly random distribution would probably say that the nodes should be looked at while a clustered distribution probably points to environmental factors as described above.

7.3 Looking for Correlations

The Correlation tool allows to assess the linear correlation between two system variables, sampled over the same time interval and across the same spatial layout (same node, same rack, or entire system). The Correlation tool calculates the linear regression of each of the 2 variables on the other; this calculation thus yields 2 regression slopes and 2 regression intercepts. In addition to these 2 regressions that provide information as to how one variable can be linearly guessed from the knowledge of the other variable, the tool also calculates the linear correlation coefficient. This coefficient, ranging in $[-1, 1]$, measures to goodness of fit of the linear regression. This tool can thus provide valuable insight as to correlations within the system.

8 OVIS Tutorial

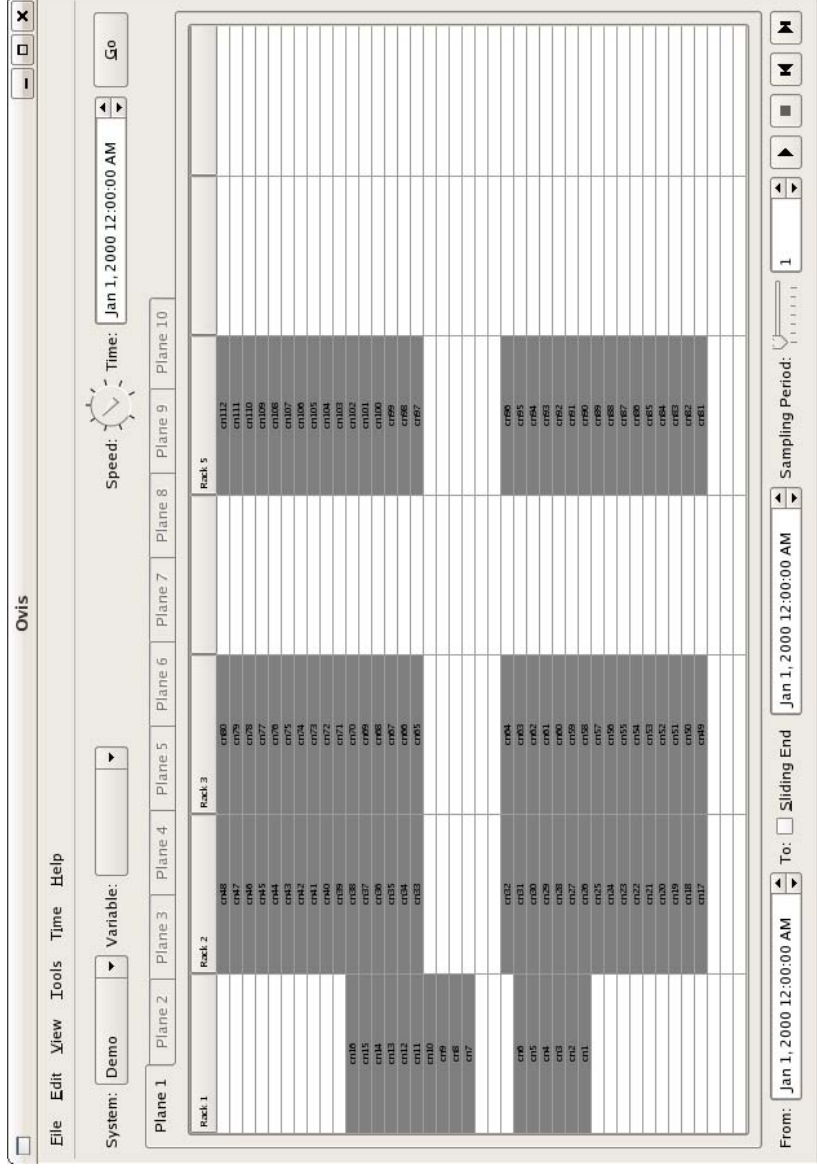
The following slides are a tutorial guide to setting up OVIS. The demo files included with the source tarball are used.

OVIS Tutorial

- **Note:** This tutorial assumes you have properly installed QT4 and ovis-0.1 according to the installation instructions in the Users Guide.

Starting OVIS

- `cd <install_path>/ovis-0.1/src`
- `./Ovis -->` Starts the OVIS GUI shown below
 - The configuration shown is for demonstration



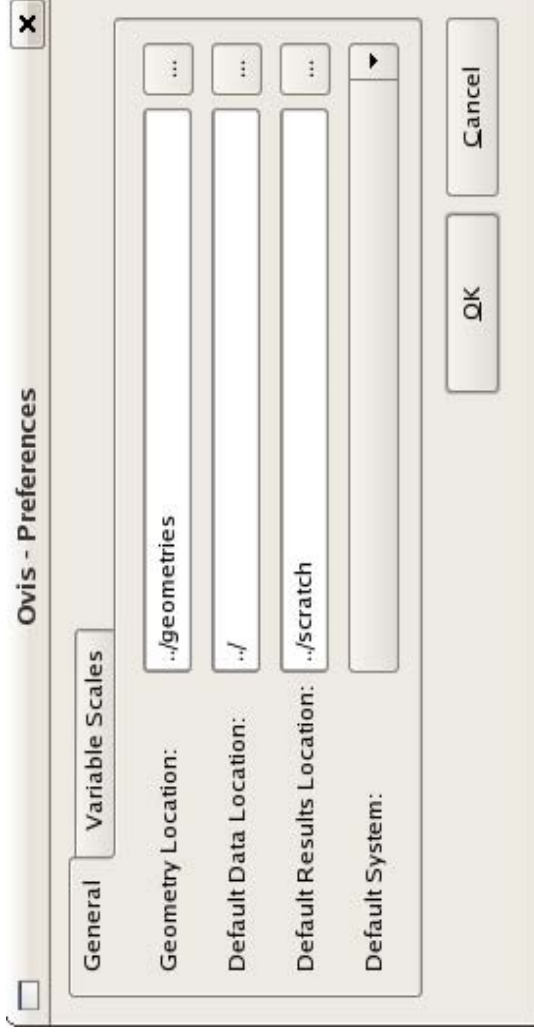
OVIS Configuration

Will be covered over the next 11 slides

- Before using OVIS to visualize data you must:
 - Specify a geometry file which describes the physical layout of the cluster
 - Specify a dataset, that you want to view, corresponding to the geometry selected
- To make the visualization meaningful you must:
 - Define color maps for the variables in accordance with the range of values in the dataset

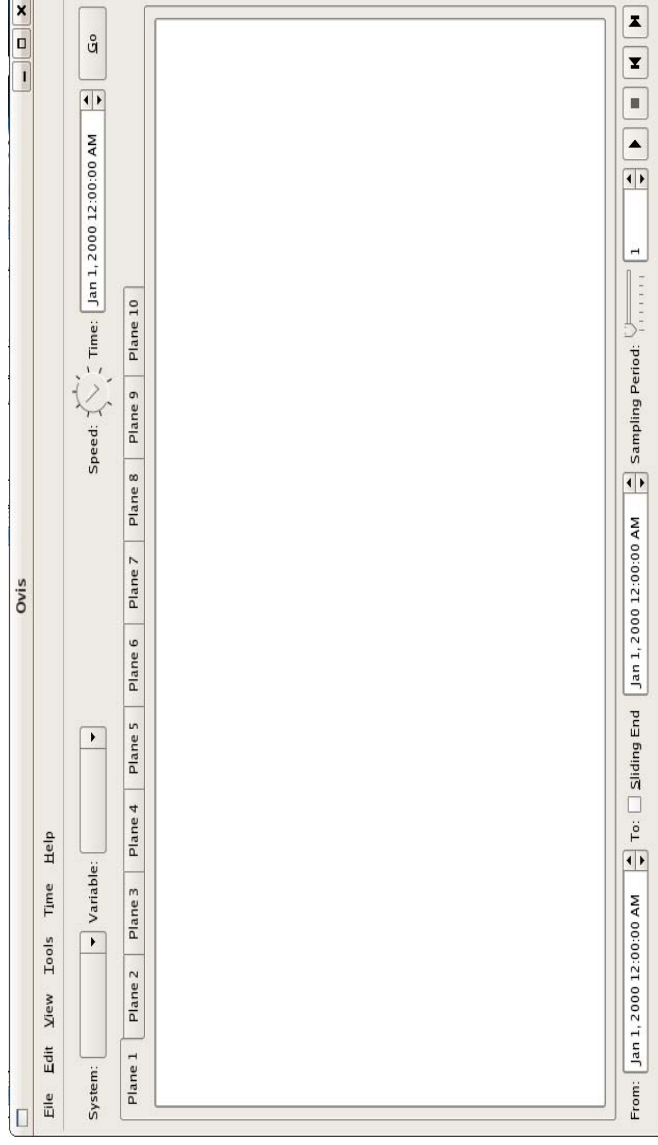
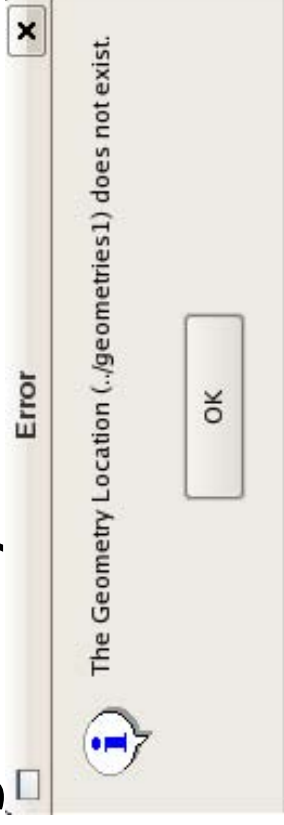
Specifying Geometry File

- From the upper left pull down menu select:
 - Edit-->Preferences
 - Type the path in the “Default Geometry Location” dialog box or click on ... to browse the filesystem
 - Click “Open”

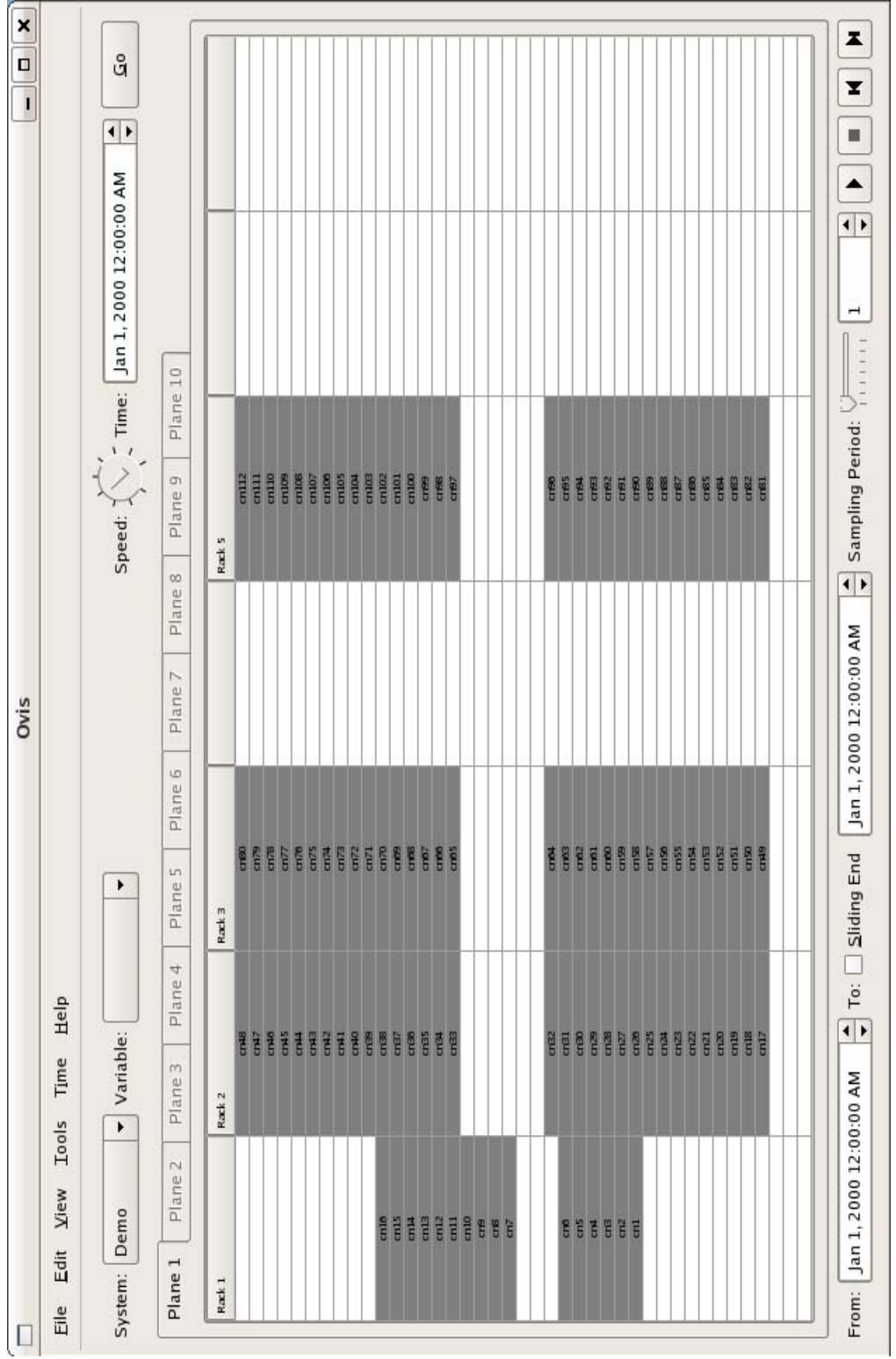


Specifying Geometry File Cont.

- If no valid geometry file exists in the location given
- And the resulting display looks like this

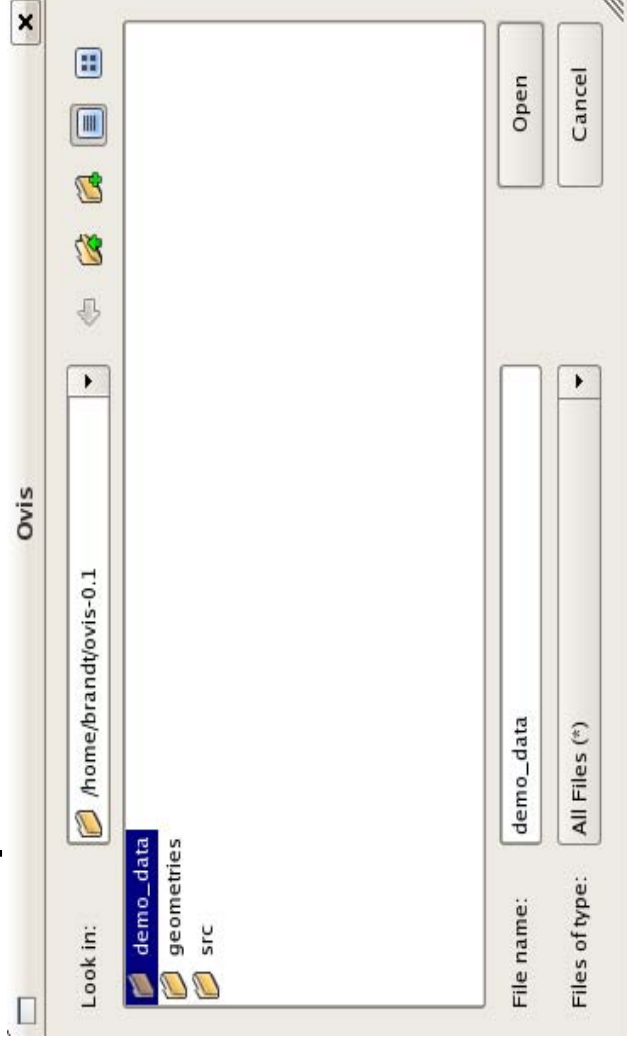


System Display With Geometry File But Without Data



Select Dataset

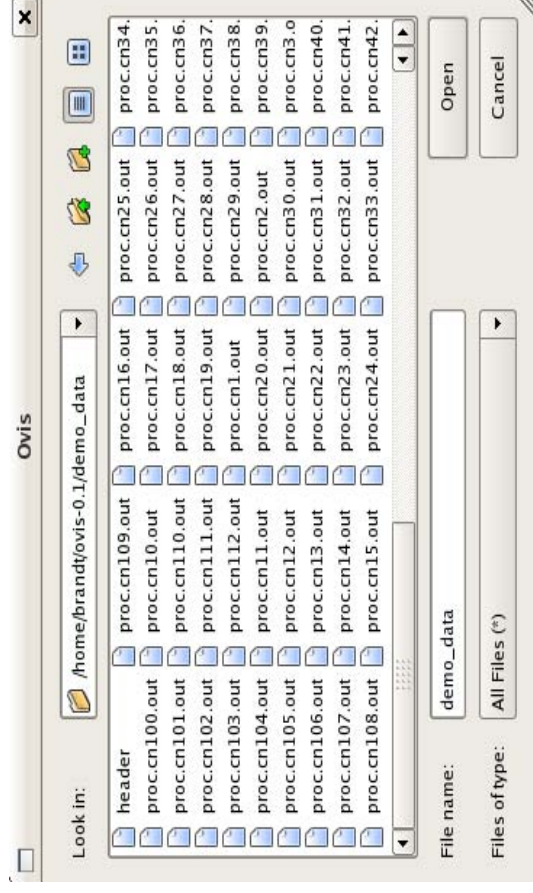
- File-->Open Dataset



- Browse for data directory and click on Open
 - Note: Only click “Open” when you have selected the directory containing the data files. **Do NOT** use “Open” as a navigation tool.

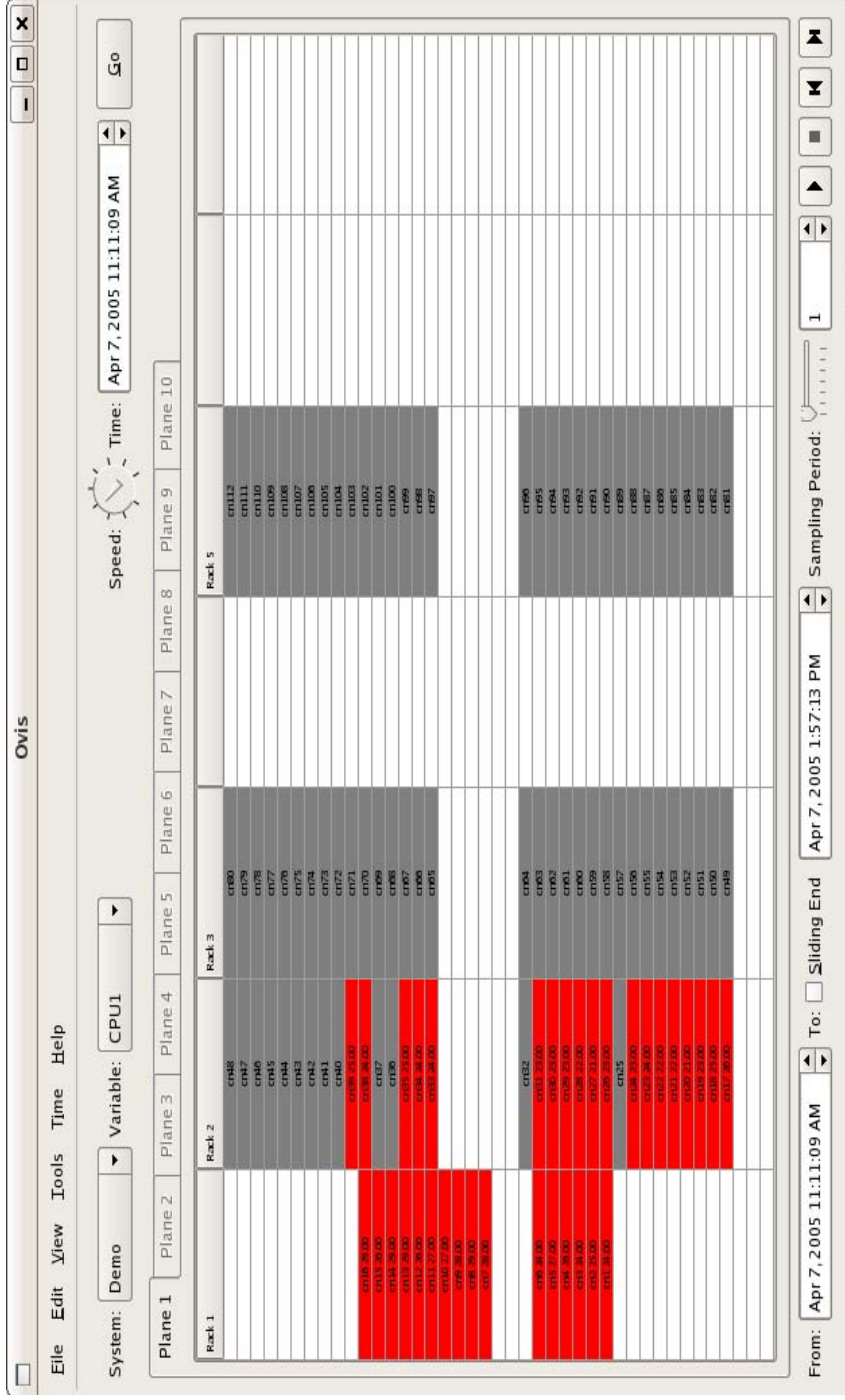
Select Dataset Cont.

- You are opening a **DIRECTORY** not a file. If preferences specifies the directory containing the actual data files or you double click on this directory when specifying the dataset, “Open” will do nothing and you will get:



- In this case go up one level then select the **directory** containing these files and click “Open”

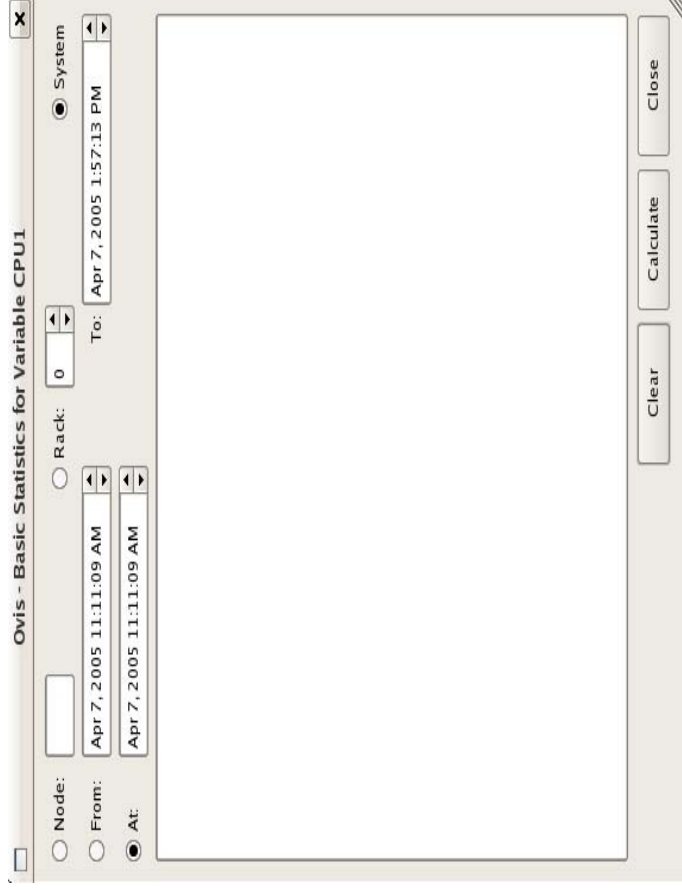
OVIS with demo Geometry file and demo_data Dataset loaded



- Note that not all nodes have data on the first timestep as displayed by gray boxes

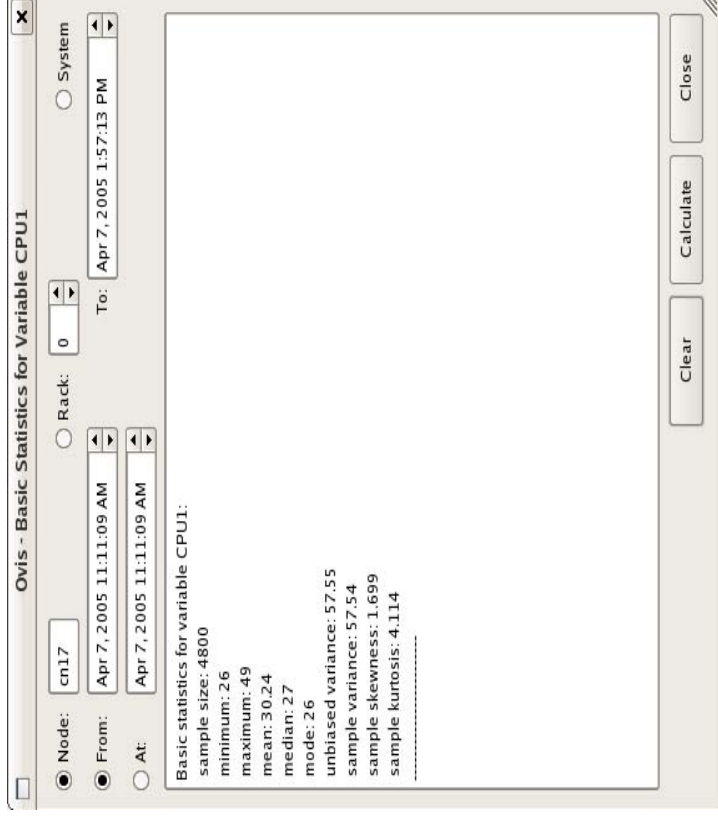
Define meaningful color maps

- Determine reasonable data value limits
 - First select a variable from the “Variable” combo box in the upper left of the GUI
 - Next start the statistics tool: Tools-->Statistics



Configure statistics tool

- The fastest analysis will be on a single node and if data has been taken over a period sufficiently long for a node to have gone through periods of idle and high CPU utilization it will get you in the ballpark.

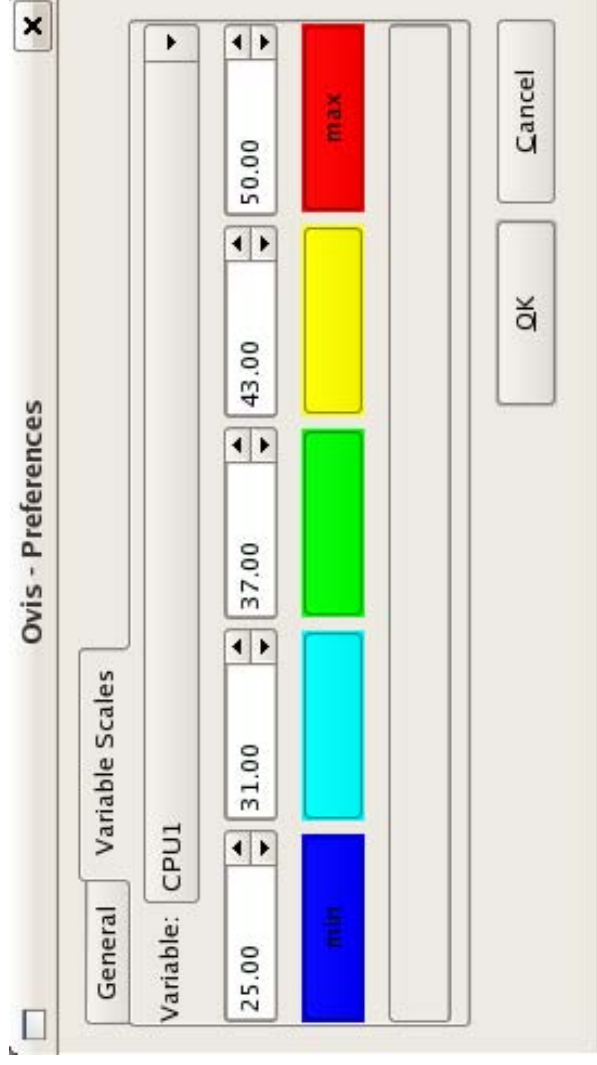


Configure Color Maps

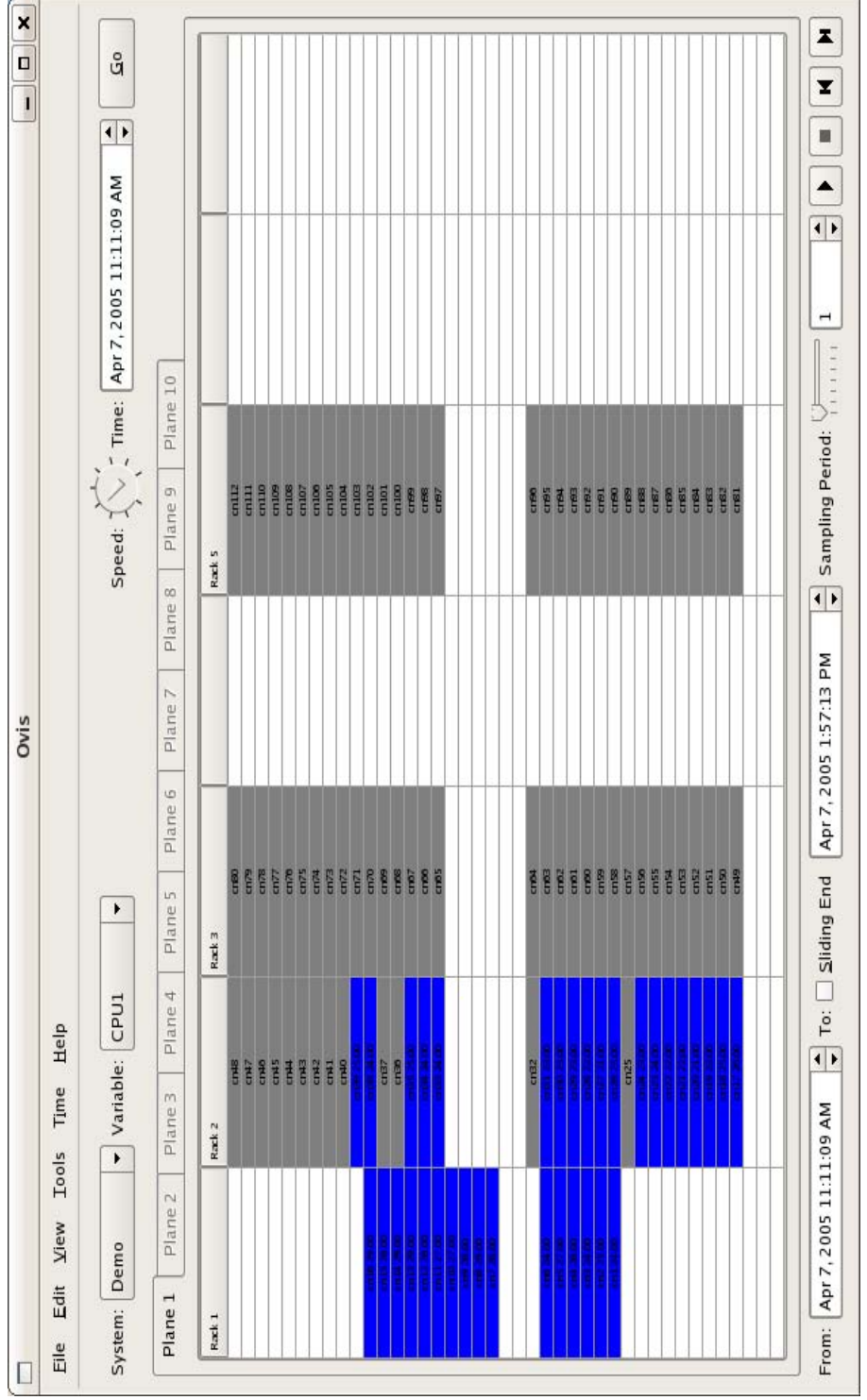
- After gathering the minimum and maximum values for the variables of interest select:
 - Edit-->Preferences (this will pop up a dialog box)
- Click on the “Variable Scales” tab
- From here you can select any variable available in the dataset and distribute the spectrum over the value range of the variable however you see fit. See next slide for an example using the statistics data from the previous slide.

Color Map for variable CPU1

- Note that this scale has been set from 25 to 50 to reflect the range for on node's statistics over some period. This may need to be adjusted slightly if a lot of the nodes spend significant time colored on either extreme.



OVIS with demo Geometry file, demo_data Dataset loaded and color map defined for variable "CPU1"



At this point OVIS is configured for meaningful display of this dataset on this cluster specification

- The following slides will explain the controls functionality and use as well as how to use this tool to:
 - Evaluate the cluster environment
 - Evaluate cluster configuration
 - Use the statistics tool to do meaningful analysis of system parameters

OVIS GUI Controls

- There are four basic control types
 - DVD like “playing” controls allow you to step through a dataset like you would a movie
 - Time Location controls allow fine grained control on region and position
 - Sampling and Speed controls allow definition of how often, in time, to sample and at what rate to display these samples
 - View controls affect how and what data are displayed

OVIS GUI Controls Continued

- Playing controls allow you to “play” through a dataset like you would a movie
 - “Play” -- Right pointing arrow starts playing from “Time” towards “To” in accordance with “Speed” and “Sampling Period” rules
 - “Stop” -- Square stops playing but remains at current “Time”
 - “Rewind” -- Left pointing arrow with verticle line rewinds to “From” time
 - “Forward” -- Right pointing arrow with verticle line forwards to “To” time

OVIS GUI Controls Continued

- Time location
 - “From” bounds time on the left and is located in the lower left of the display
 - Set by default to earliest time in dataset
 - Can be set by the user anywhere between the default and the time in the “To” box
 - This allows the user to effectively zoom in on a time region
 - “To” bounds time on the right and is located to the right of the From box
 - Set by default to the latest time in the dataset
 - Can be set by the user anywhere between the default and the time in the “From” box
 - This allows the user to effectively zoom in on a time region

OVIS GUI Controls Continued

- Time location continued
- Time points to the time of the data currently being displayed
 - Set by default to earliest time but steps through the dataset when “Playing” according to rules established by the “Speed” and “Sampling Period” controls defined below
 - Can be set by the user to any time between and including “From” and “To”
 - When user defined, “Go” causes the seek
 - Moves forward in accordance with “Speed” and “Sampling Period” rules when “Playing”
 - Remains at current time when “Stop” is depressed
 - Moves to “From” time when “Rewind” is depressed
 - Moves to “To” when “Forward” is depressed

OVIS GUI Controls Continued

- Sampling and Speed controls
 - “Sampling Period” is a spinbox located in the lower right of the display. The number displayed in this box defines the subsampling interval, in seconds, for display of data from a dataset. Setting an interval shorter than the interval between successive datapoints yields no data for the timesteps between even though the display will update for every such timestep it would be empty for the times between actual data unless interpolation was turned on
 - “Speed” is a radio dial located in the upper right of the display. This dial affects the time between successive updates of the display according to the sampling period chosen. Clockwise is faster and actual maximum will be determined by processor and computational load.

OVIS GUI Controls Continued

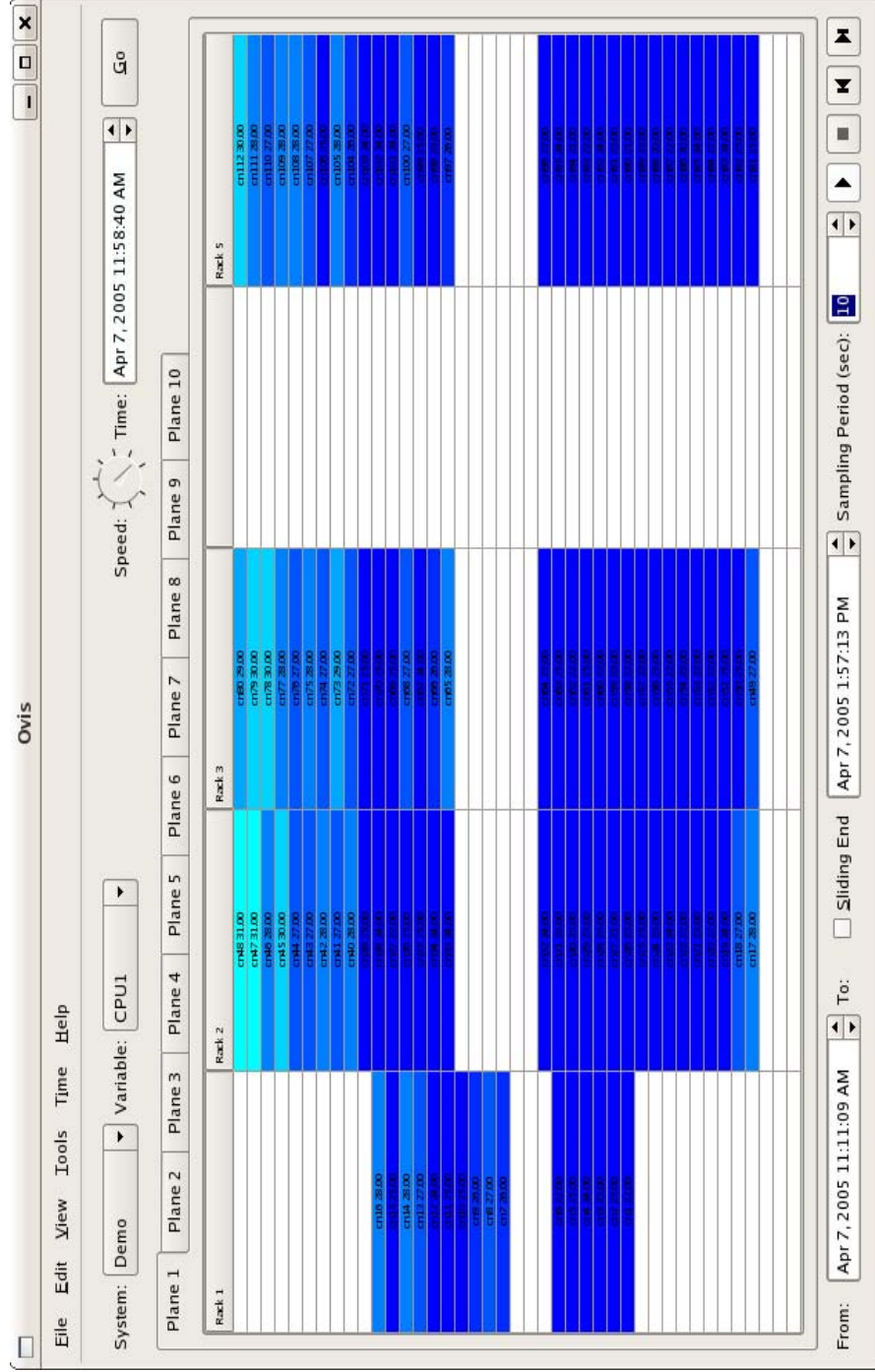
- View controls affect how and what data are displayed
 - Bin or gradient according to colormap
 - If the “Use Gradient Colormap” box is checked (default) the color displayed is calculated by interpolating between the colors the variable value falls between
 - See the “Customizing Color Scales” section of the Operation Guide for more comprehensive description of these options
 - Show missing data or use interpolation to hide the gaps
 - Options under “Interpolate Missing Data” allow data in a range about the current time to be displayed if there is no data for a node for the current timestep. These include: No interpolation, a finite range centered about the current time, and carryover of the previous valid value regardless of age of the data point.
 - See the “Use of Interpolation” section of the Operation Guide for discussion on where and when this should be used.

Evaluation of Cluster Environment

- Find a place in time where the operating conditions are stable and similar across the system.
 - Since the demo dataset was taken using our calibration code running on all nodes in the system we step through the data and note times when the system looks stable at idle and at high CPU utilization (apparent from reddish colored nodes in this example)
 - The current version of OVIS does not include a mechanism for normalizing one variable to some function of others.

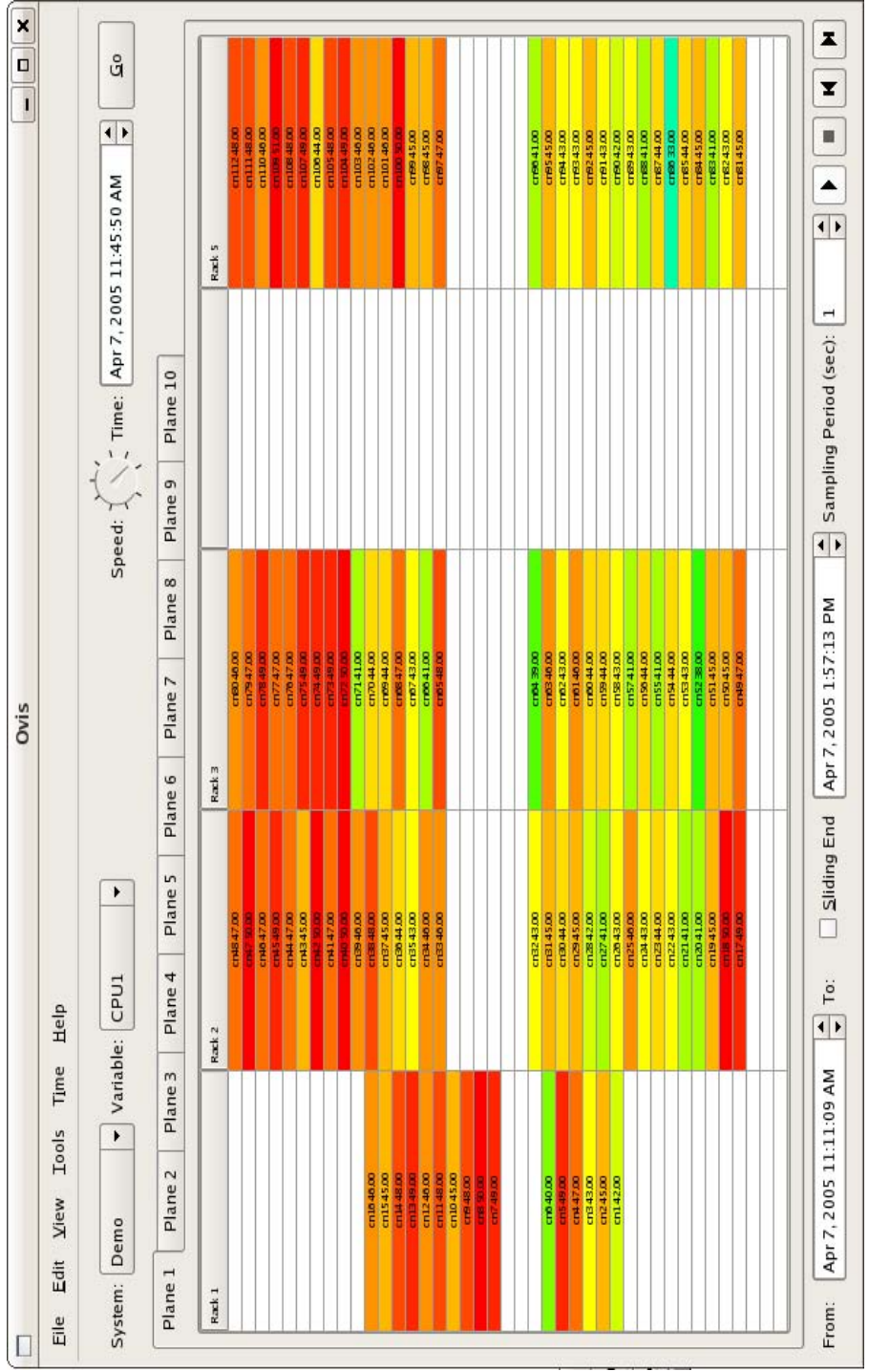
All Nodes Idle

- Note the thermal gradient with vertical dependency
- Looks here like this cluster needs more airflow



All Nodes 100% Utilization

- Note the thermal gradient with vertical dependency
 - Cluster needs more air volume but lower air velocity as the current high velocity cooling air is causing hot exhaust to be sucked under racks and into bottom nodes



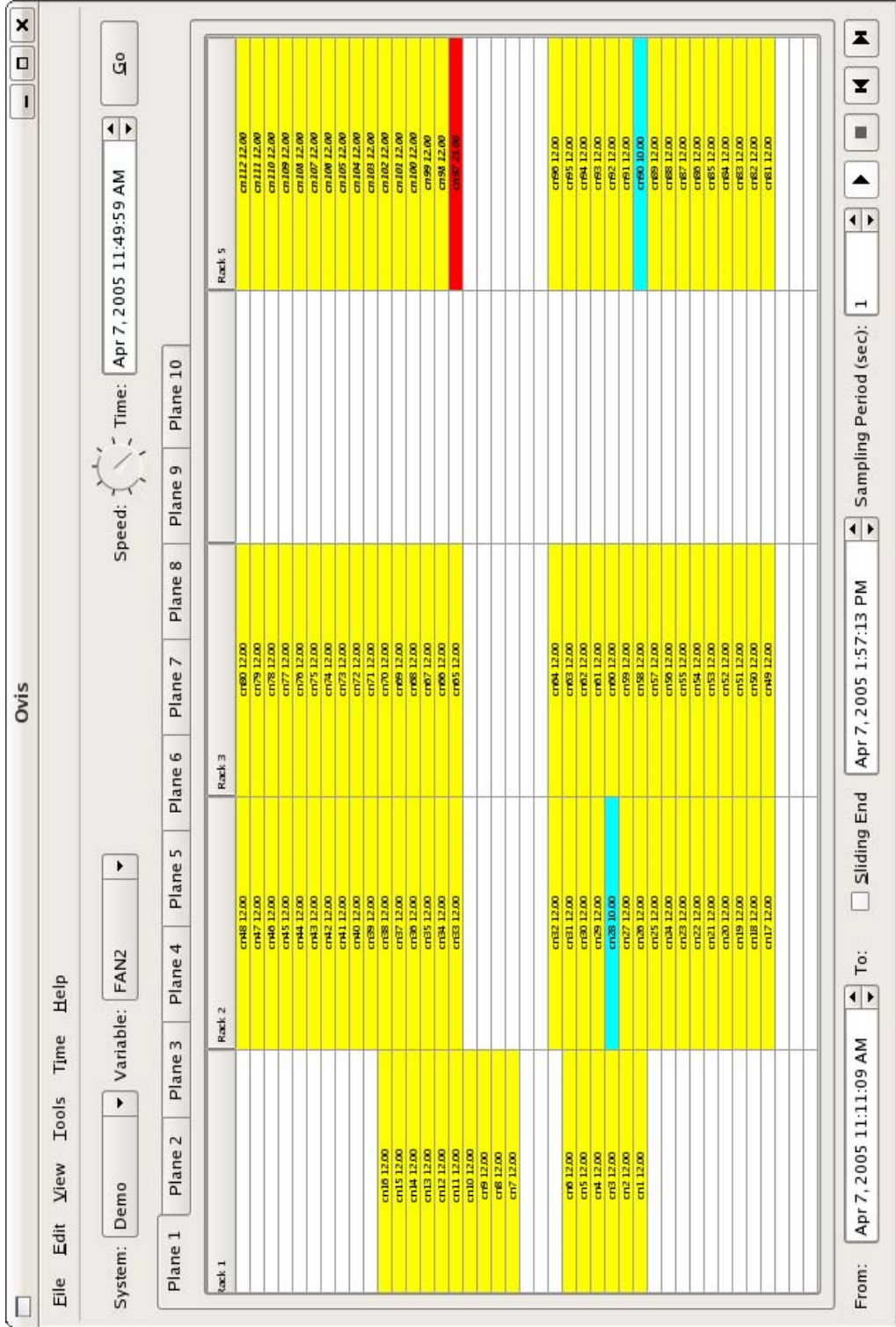
Evaluation of Cluster Configuration

- Using the same method to find regions of operational stability we look for variation not explained by the relatively smooth gradients we would expect from environmental phenomena.
 - In the display showing the cluster at idle there is nothing remarkable.
 - In the display showing the cluster at 100% utilization the nodes just above the empty space in the middle run a little hotter than the ones before the space and the ones above them. Since this phenomenon spans all of the racks it is worth investigating what role these spaces might play.

Cluster Configuration Continued

- In the case of this particular cluster it turned out that these spaces had no panels in them and hot exhaust air was being recirculated. This particular data was taken after blocking this space in racks 1, 2 and 3 with cardboard which still allowed some leakage. In the case of rack 5, though the temperature profile looks the same a look at the fan speed profile (next slide) shows that fan speed of the node above the space is markedly higher.
 - Note: The color maps for the fan speeds was set up in the same way as those for the CPU temperatures and ran from 9 to 13 as the nodes selected showed a high and low of 12.

Fan Speed Profile



DISTRIBUTION:

- | | |
|------------------------------------|---|
| 1 MS 9051
Dawn Manley, 8351 | 1 MS 9158
David A. Evensky, 8963 |
| 1 MS 9152
Jerry Friesen, 8963 | 2 MS 9152
David C. Thompson, 8963 |
| 1 MS 0817
James Ang, 9224 | 2 MS 9152
Jim M. Brandt, 8963 |
| 1 MS 0823
Josh England, 4324 | 2 MS 9152
Ann C. Gentile, 8963 |
| 1 MS 1319
Jon Stearley, 1422 | 2 MS 9152
Matthew Wong, 8963 |
| 1 MS 0807
Bob Ballance, 4328 | 1 MS 9051
Youssef Marzouk, 8351 |
| 1 MS 9011
Ed Talbot, 8965 | 2 MS 9051
Philippe P. Pébay, 8351 |
| 1 MS 9151
Howard Hirano, 8960 | 3 MS 9018
Central Technical Files, 8940-1 |
| 1 MS 9158
Curtiss Jansen, 8961 | 1 MS 0899
Technical Library, 9616 |
| 1 MS 9159
Heidi Ammerlahn, 8962 | 1 MS 9021
Classification Office, 8511, for
Technical Library, MS 0899, 9616
DOE/OSTI via URL |
| 1 MS 9159
Mike Hardwick, 8964 | |