

Runtime Improvements

Chapel Team, Cray Inc.

Chapel version 1.9 summary

April 17th, 2014 (released) / May 2014 (documented)



COMPUTE | STORE | ANALYZE

Safe Harbor Statement

The Cray logo is located in the top right corner of the slide. It consists of the word "CRAY" in a blue, sans-serif font, followed by a stylized graphic of a network or cluster of nodes and lines.

This presentation may contain forward-looking statements that are based on our current expectations. Forward looking statements may include statements about our financial guidance and expected operating results, our opportunities and future potential, our product development and new product introduction plans, our ability to expand and penetrate our addressable markets and other statements that are not historical facts. These statements are only predictions and actual results may materially vary from those projected. Please refer to Cray's documents filed with the SEC from time to time concerning factors that could affect the Company and these forward-looking statements.



COMPUTE | STORE | ANALYZE

Copyright 2014 Cray Inc.

2

Runtime Improvements



- **Accurate running task counts**
- **Symmetry in multilocale execution**
- **Tasking-layer-independent default stack size**
- **Hide the threading layer**
- **Remove 'none' tasking layer**
- **Qthreads status**
 - Bug fixes and other changes
 - Progress towards making qthreads the default tasking layer



COMPUTE | STORE | ANALYZE

Copyright 2014 Cray Inc.

3

Note that though most of this work was motivated by preparing Qthreads to become the default tasking layer, the solutions were about making behavior consistent across tasking layers, regardless of which one was the default.

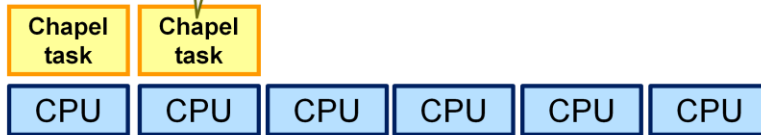
Accurate Running Task Counts

CRAY

```
forall a in A {  
  ...  
}
```

If we want to fully utilize the hardware,
how many tasks should we create?

Answer: #CPUs - #active tasks



COMPUTE | STORE | ANALYZE

Copyright 2014 Cray Inc.

4

Accurate Running Task Counts

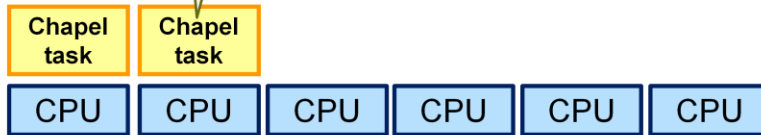
CRAY

```
forall a in A {  
  ...  
}
```

If we want to fully utilize the hardware,
how many tasks should we create?

Answer: #CPUs - #active tasks

We were getting this wrong.



COMPUTE | STORE | ANALYZE

Copyright 2014 Cray Inc.

5

Accurate Running Task Counts

The Cray logo is located in the top right corner of the slide. It consists of the word "CRAY" in a blue, sans-serif font, followed by a stylized graphic of a network or cluster of nodes and connections.

Background:

- New default in 1.9: dataParIgnoreRunningTasks=false
- When deciding how many tasks to create for a forall-stmt, take into account how many are already running
- Our running-task counting was inaccurate

Why Is This a Problem?

- Not enough tasks: under-perform due to not using available CPUs
- Too many tasks: under-perform due to increased overhead



COMPUTE | STORE | ANALYZE

Copyright 2014 Cray Inc.

6

Accurate Running Task Counts

CRAY

What a mess!

fifo	Decrement on task termination late and not synchronized with end of parallel statement; sometimes not done before next parallel statement encountered. Intermittently over-counted tasks, under-utilized CPUs.
massivethreads	Did not maintain running task count at all, always said 0. Under-counted tasks, over-utilized CPUs.
muxed	Decrement on task termination late and not synchronized with end of parallel statement; sometimes not done before next parallel statement encountered. Intermittently over-counted tasks, under-utilized CPUs.
qthreads	Only counted running tasks on current shepherd. But often we have more than one shepherd. Under-counted tasks, over-utilized CPUs.



COMPUTE | STORE | ANALYZE

Copyright 2014 Cray Inc.

7

Accurate Running Task Counts



Solution: Moved running task count into the modules

- Reduces duplicate code
- Gets the right answer!
(for all tasking implementations)

Next Step:

- Remove little bit of remnant task counting code from runtime



COMPUTE | STORE | ANALYZE

Copyright 2014 Cray Inc.

8

Symmetry in Multilocale Execution



Background: Asymmetric execution scheme across locales

- 0 and non-0 locales differed in what they used the locale process for, and where they did Active Message handling
- Asymmetry led to more complicated start up code in the runtime
- Behavior on locale 0 differed from that of other locales



COMPUTE | STORE | ANALYZE

Copyright 2014 Cray Inc.

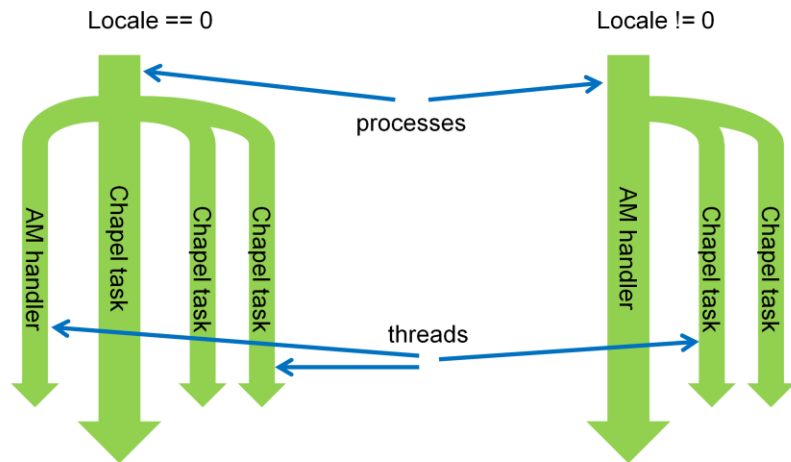
9

“Behavior” (last bullet) in the sense of what happened there and why. For example, if you noticed that on locale 0 you were seeing interference between Active Message handling and Chapel tasks, you could not deduce that that was also happening on other locales.

Symmetry in Multilocale Execution

CRAY

v1.8 and before



COMPUTE | STORE | ANALYZE

Copyright 2014 Cray Inc.

10

- Cannot configure tasking layers to behave the same on all locales
- Hard to reason about behavior on non-0 locales based on measurements taken on locale 0 and vice-versa

Symmetry in Multilocale Execution



This Effort: Make execution scheme symmetric across locales

- All locales have the initial process waiting for completion
- Start the AM handler as a separate thread/task on all locales



COMPUTE | STORE | ANALYZE

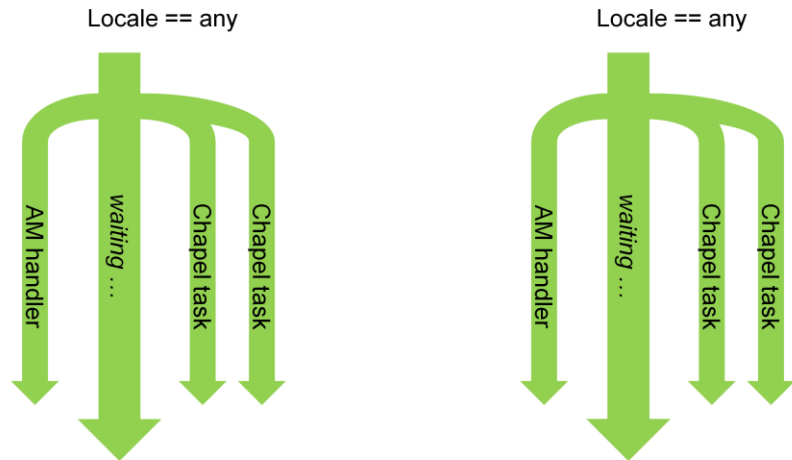
Copyright 2014 Cray Inc.

11

Symmetry in Multilocale Execution

CRAY

v1.9



- Tasking layers behave the same on all locales
- Knowledge gained about behavior on one locale applies to other locales also



Copyright 2014 Cray Inc.

12

Symmetry in Multilocale Execution

The Cray logo is located in the top right corner of the slide. It consists of the word "CRAY" in a blue, sans-serif font, followed by a stylized graphic of a network or cluster of nodes and connections.

Impact:

- Hides knowledge of comm layer AM handling inside comm layers
- Still have locale 0/non-0 asymmetry
 - Okay, because inherent to Chapel's execution model
- Improves ability to reason about overall behavior from observations or measurements done on one locale

Next Steps:

- Further simplification of runtime start up code



COMPUTE | STORE | ANALYZE

Copyright 2014 Cray Inc.

13

Tasking Layer Independent Default Stack Size

CRAY

Problem:

- Tasking implementations set default call stack sizes themselves
- Duplicated effort without duplicating results
- Led to surprises when switching from one implementation to another

Solution:

- Moved default call stack size selection into common code

Benefits:

- Less code to maintain
- Fewer surprises; more dependability



COMPUTE | STORE | ANALYZE

Copyright 2014 Cray Inc.

14

Hide the Threading Layer



Background:

- For some time, we've had both tasks and threads as 1st-class entities:

```
[gbt@host]$ printchplenv
```

```
...
```

```
CHPL_COMM: none
```

```
CHPL_TASKS: fifo
```

```
CHPL_THREADS: pthreads
```

```
CHPL_LAUNCHER: none
```

```
...
```
- But this is unnecessarily complicated



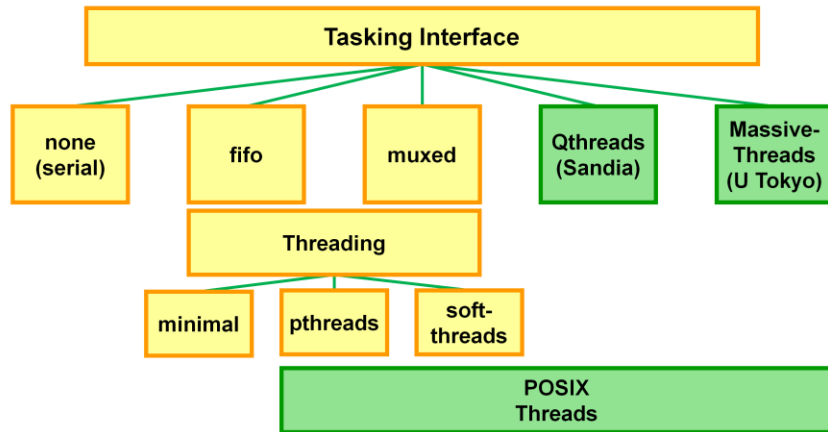
COMPUTE | STORE | ANALYZE

Copyright 2014 Cray Inc.

15

Hide the Threading Layer

CRAY



COMPUTE | STORE | ANALYZE

Copyright 2014 Cray Inc.

16

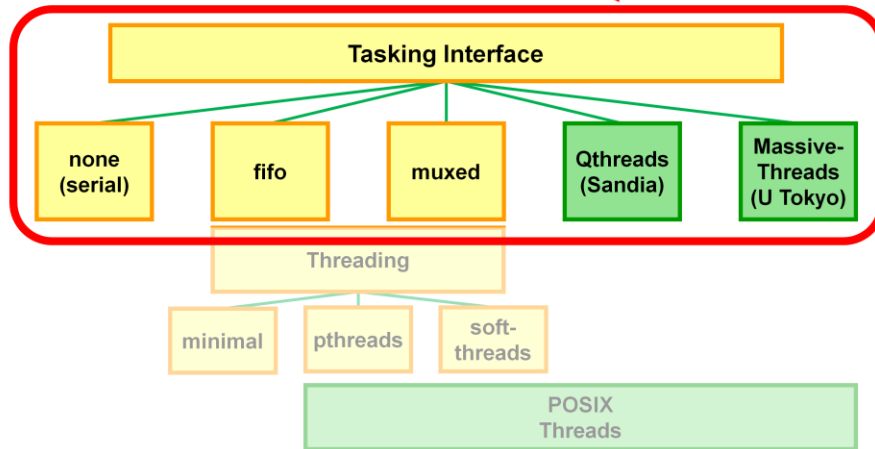
Here's a picture of how tasking is implemented.

From top to bottom, how far down do we want users to be thinking about this?

Hide the Threading Layer

CRAY

want users to think about this, at most

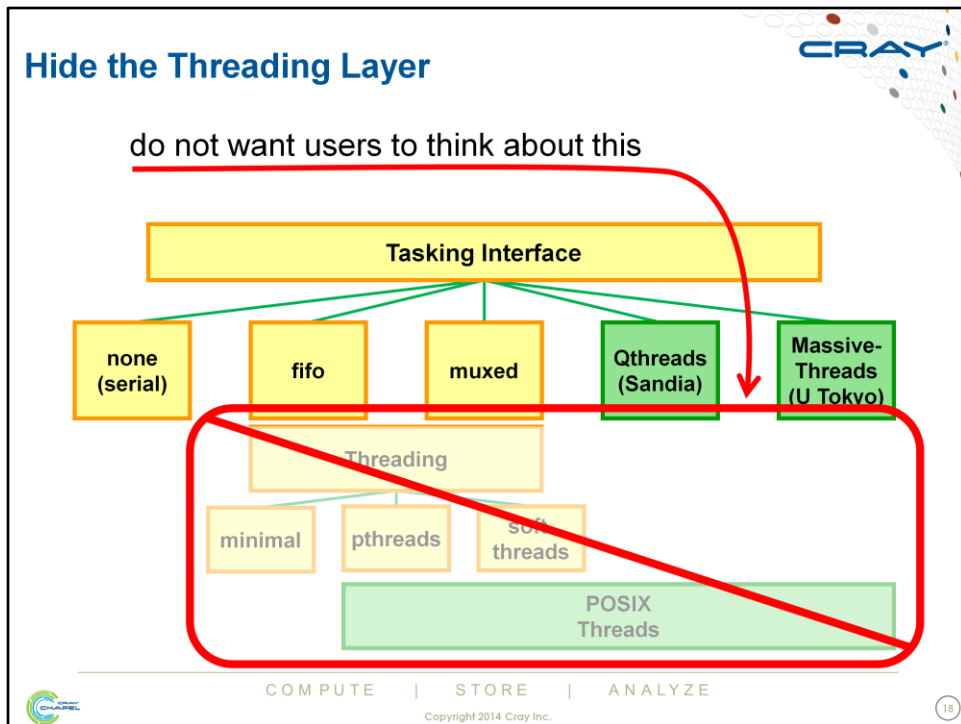


COMPUTE | STORE | ANALYZE

Copyright 2014 Cray Inc.

17

- **Tasks are the Chapel abstraction of execution**
 - Behave as described in the Chapel Language Specification
 - Parallelism and synchronization is expressed in terms of tasks
 - Want programmers to reason on this level



- **Threads are an underlying software abstraction by which tasking layers make use of hardware processors**
 - Defined and used differently by each tasking layer
 - For fifo, a thread is a Linux (UNIX) pthread and a Chapel task is bound to a single thread throughout its existence
 - For qthreads, a thread is a worker qthread and a Chapel task may shared its thread with other threads and/or change host threads during its life
 - Etc.
 - Don't really have anything to do with Chapel programming
 - Do not want programmers to be burdened with this level of detail

Hide the Threading Layer



Result:

- printchplenv doesn't show CHPL_THREADS any more:

```
[gbt@host]$ printchplenv
```

```
...
```

```
CHPL_COMM: none
```

```
CHPL_TASKS: fifo
```

```
CHPL_LAUNCHER: none
```

```
...
```
- If CHPL_THREADS is set, an error message results
- The notion of threading still exists, but only deep in the implementation (and will likely fade away over time)



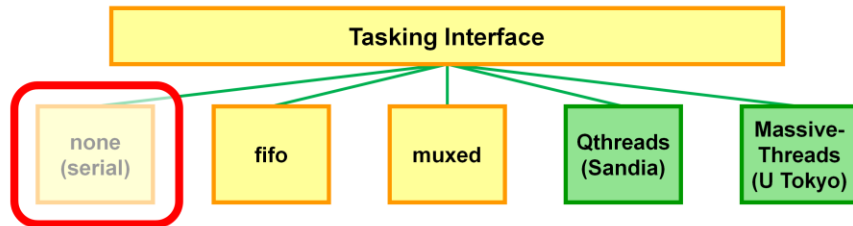
COMPUTE | STORE | ANALYZE

Copyright 2014 Cray Inc.

19

Remove 'none' Tasking Layer

CRAY



- Minimalist tasking implementation
- Not regularly tested
- Not used
- Had a small but non-zero maintenance cost
- Got rid of it



COMPUTE | STORE | ANALYZE

Copyright 2014 Cray Inc.

20

Qthreads Status



- **Bug fixes and other changes**
 - Better multilocale integration
 - Other optimization and tuning
- **Progress towards making qthreads the default tasking layer**



COMPUTE | STORE | ANALYZE

Copyright 2014 Cray Inc.

21

Better Multilocale Integration for qthreads



Background: Multilocale qthreads program were hanging intermittently

- Tasks were being starved



COMPUTE | STORE | ANALYZE

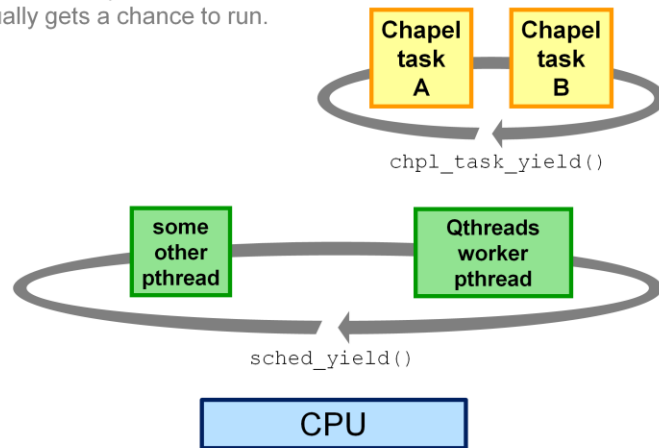
Copyright 2014 Cray Inc.

22

Better Multilocale Integration for qthreads

CRAY

This works well: the nested yields ensure that everything eventually gets a chance to run.



- We're time-sharing execution vehicles at two levels here: Chapel tasks on Qthreads worker pthreads, and pthreads of various kinds on the CPU.



COMPUTE | STORE | ANALYZE

Copyright 2014 Cray Inc.

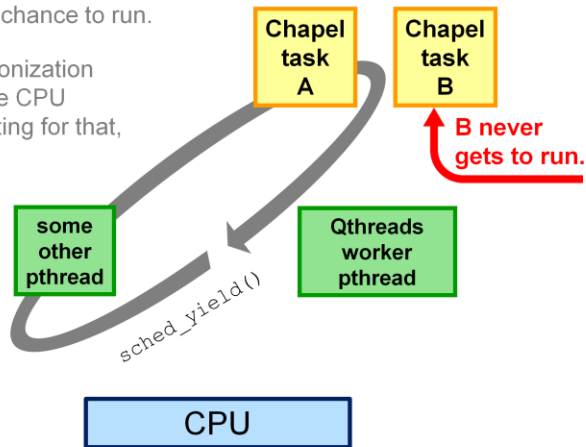
23

Better Multiscale Integration for qthreads

CRAY

This works well: the nested yields ensure that everything eventually gets a chance to run.

But if A and B have a synchronization dependence and we yield the CPU instead of the task while waiting for that, we no longer run both.



- Since we're yielding the pthread on the CPU instead of the task on the worker, we never go back into the Qthreads code to do a task switch to task B.
- Note that adding resources (more worker pthreads or more CPUs) doesn't solve the basic problem, it just changes how many things you need going on for it to happen.



COMPUTE | STORE | ANALYZE

Copyright 2014 Cray Inc.

24

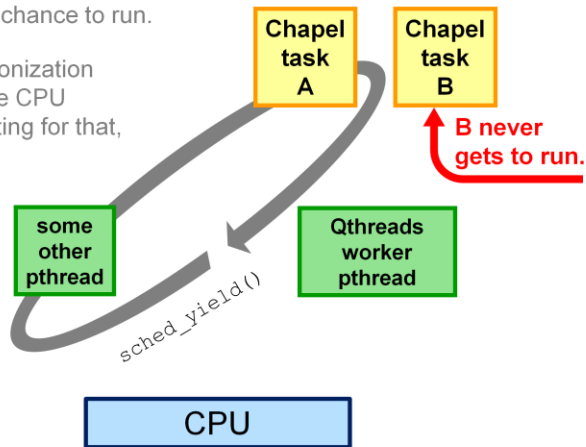
Better Multiscale Integration for qthreads

CRAY

This works well: the nested yields ensure that everything eventually gets a chance to run.

But if A and B have a synchronization dependence and we yield the CPU instead of the task while waiting for that, we no longer run both.

We did this in the GASNet comm layer, while waiting for remote task completion. Result: hang.
Solution: yield the task while waiting in the comm layer instead.



COMPUTE | STORE | ANALYZE

Copyright 2014 Cray Inc.

25

Qthreads Optimization and Tuning

The Cray logo is located in the top right corner of the slide. It consists of the word "CRAY" in a blue, sans-serif font, followed by a stylized graphic of a network or cluster of nodes and connections.

Done In v1.9: Optimize

- Inline several small, frequently used utility functions
- Build with oversubscription enabled to support multilocale testing
- Configure and enable guard pages by default for functional testing, but publicize how to disable for performance testing

Done After v1.9: Tune the defaults

- Assign worker pthreads to cores
(was assigning workers to hyperthreads, when those were present)
- Default number of workers = number of cores
(was = number of hyperthreads, due to a bug)



COMPUTE | STORE | ANALYZE

Copyright 2014 Cray Inc.

26

Qthreads as the Default Tasking Layer



Background

- Current default FIFO tasking has flaws:
 - No support for NUMA or other node hierarchy
 - Sync variable synchronization slow (requires kernel intervention)
- 3 candidate tasking layers to replace FIFO:
 - Muxed: lacks NUMA, not open source (Cray specific)
 - Massivethreads: lacks NUMA, immature
 - Qthreads: has NUMA, fairly mature

Qthreads benefits

- NUMA support (with hwloc)
- Fast sync variable synchronization (done at user level)
- Open source and fairly mature
- Most mature as a Chapel tasking layer



COMPUTE | STORE | ANALYZE

Copyright 2014 Cray Inc.

27

- “Fairly mature”: qthreads itself is quite mature, but the qthreads-based Chapel tasking layer is less so. It was at least being minimally tested in nightly testing, though, which was/is not the case for either muxed or massivethreads.

Qthreads as the Default Tasking Layer



Status:

- With comm=none
 - Passes same nightly tests as fifo tasking, and in 7.5 hr vs. 8.5 hr
- With comm=gasnet
 - Passes multilocale tests in nightly testing

Next Steps:

- Expand to full nightly testing
- Characterize performance
- Fix any problems
- Switch!
- Longer term:
 - Tie Chapel sync vars more directly to qthreads sync vars



COMPUTE | STORE | ANALYZE

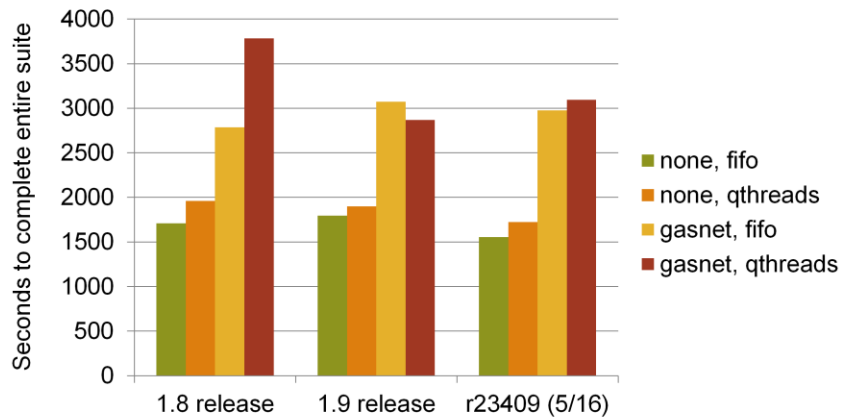
Copyright 2014 Cray Inc.

28

Qthreads as Default Tasking

CRAY

Performance on the examples suite



COMPUTE | STORE | ANALYZE

Copyright 2014 Cray Inc.

29

Note that the bulk of the time here is actually compilation, but it still gives an indication that we're in the ballpark.

Legal Disclaimer

The Cray logo is located in the top right corner of the page. It consists of the word "CRAY" in a blue, sans-serif font, followed by a stylized graphic of a network or cluster of nodes and connections.

Information in this document is provided in connection with Cray Inc. products. No license, express or implied, to any intellectual property rights is granted by this document.

Cray Inc. may make changes to specifications and product descriptions at any time, without notice.

All products, dates and figures specified are preliminary based on current expectations, and are subject to change without notice.

Cray hardware and software products may contain design defects or errors known as errata, which may cause the product to deviate from published specifications. Current characterized errata are available on request.

Cray uses codenames internally to identify products that are in development and not yet publically announced for release. Customers and other third parties are not authorized by Cray Inc. to use codenames in advertising, promotion or marketing and any use of Cray Inc. internal codenames is at the sole risk of the user.

Performance tests and ratings are measured using specific systems and/or components and reflect the approximate performance of Cray Inc. products as measured by those tests. Any difference in system hardware or software design or configuration may affect actual performance.

The following are trademarks of Cray Inc. and are registered in the United States and other countries: CRAY and design, SONEXION, URIKA, and YARCDATA. The following are trademarks of Cray Inc.: ACE, APPRENTICE2, CHAPEL, CLUSTER CONNECT, CRAYPAT, CRAYPORT, ECOPHLEX, LIBSCI, NODEKARE, THREADSTORM. The following system family marks, and associated model number marks, are trademarks of Cray Inc.: CS, CX, XC, XE, XK, XMT, and XT. The registered trademark LINUX is used pursuant to a sublicense from LMI, the exclusive licensee of Linus Torvalds, owner of the mark on a worldwide basis. Other trademarks used in this document are the property of their respective owners.

Copyright 2014 Cray Inc.



COMPUTE | STORE | ANALYZE

Copyright 2014 Cray Inc.

30



CRAY
THE SUPERCOMPUTER COMPANY

<http://chapel.cray.com>

chapel_info@cray.com

<http://sourceforge.net/projects/chapel/>