

# Beowulf cluster computing for all: The HiPerCiC project

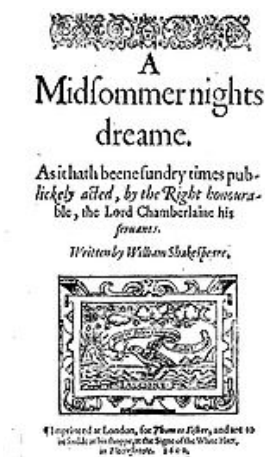
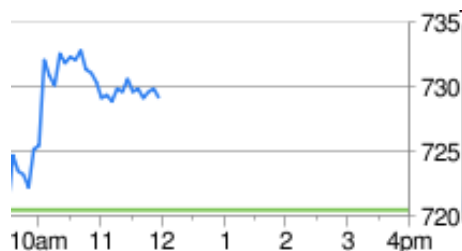
Stephanie Tanner '10, Todd Frederick '09, Jeremy Gustafson '08  
Richard Brown, Project Director

*Note: This document is adapted from the contents of a poster presented in Fall, 2009. The poster may be viewed in the Regents Hall "Link" near RNS 203. -- RAB*

## Abstract

The HiPerCiC (or **High Performance Computing in the Classroom**) project brings the results of Beowulf cluster computing to students and faculty who may have little or no knowledge of how to operate a computing cluster directly. In HiPerCiC, a large-scale *problem* or computing goal is identified in a target field, which may be in any academic discipline. Undergraduate research students develop programs for St. Olaf's Beowulf cluster computers that address the problem, then develop a web-based *user interface* that enables students and faculty in that target discipline to use the programs and explore the results conveniently, yielding a HiPerCiC *application*. Example HiPerCiC applications for problems in Environmental Science and in Political Science are presented.

## Imagining the possibilities



Potential large-scale computing problems can be identified in virtually any academic discipline or combination of disciplines. The availability of powerful Beowulf cluster computing on campus at St. Olaf makes it feasible to consider ambitious investigations that might automatically have been ruled out as impractical only a few years ago. Here are some examples:

- Examine all combinations of two or three consecutive words appearing in all Shakespeare plays. Include context information such as the line or foot in which the words appear.

Figure 2: Biological model of nitrogen flow in a riparian plant

## HiPerCiC/Riparian

The riparian-plant application described above has become the first HiPerCiC application. A prototype version of this application was created by Todd Frederick '09 and Jeremy Gustafson '08 in the Fall 2007 offering of CS 390, *Senior Capstone Seminar*. Stephanie Tanner '10 rewrote the user interface and produced a complete application as a summer research project in 2009, and continues to work with Prof. Schade to refine and extend the application.

First, a professor or advanced student produces a *data set* to examine, created by Beowulf computing through an automated procedure controlled by that user (Figure 3). Then, that user and optionally other users can *explore* that data set. In the case of the Riparian application, a data set is generated by providing parameters for Prof. Schade's model, and exploration includes graphing different combinations of the parameters and result values (Figure 4).

The screenshot shows the 'Riparian' application interface. On the left is a sidebar with the title 'High-Performance Computing in the Classroom' and contact information for Computer Science. The main area has tabs for 'Home', 'New', 'Explore', and 'Edit'. The 'New Data Set' section contains a form to create a new data set. It includes a text input for 'Name of data set:' and a table of parameters with 'Start', 'End', and 'Count' values, all set to 0.0, 1.0, and 1.0 respectively. The parameters are: Nitrogen concentration in water (p max), Half saturation coefficient (m r), %n coefficient (k p), Root:shoot coefficient (k d), Root turnover rate (k n), Root decomposition rate (k r s), and Plant productivity (n c). A 'Create' button is at the bottom right.

	Start	End	Count
Nitrogen concentration in water (p max)	0.0	1.0	1.0
Half saturation coefficient (m r)	0.0	1.0	1.0
%n coefficient (k p)	0.0	1.0	1.0
Root:shoot coefficient (k d)	0.0	1.0	1.0
Root turnover rate (k n)	0.0	1.0	1.0
Root decomposition rate (k r s)	0.0	1.0	1.0
Plant productivity (n c)	0.0	1.0	1.0

Figure 3: Creating a data set in HiPerCiC/Riparian

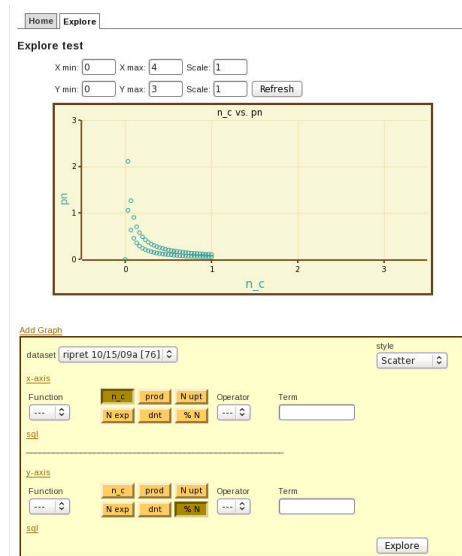


Figure 4: Exploring data in HiPerCiC/Riparian

## Political Blogs



Blogs have become a formidable factor in political discourse, and have many interesting features from a political science viewpoint (e.g., anyone has access to post, no guarantee of fact checking or editorial review, rapid and wide distribution). Yet few political-science studies have been conducted to date, at least in part because it is difficult or impossible to use traditional computing methods to process the thousands of potentially significant blog pages produced each day.

Therefore, Megan Goebel '11 (co-director Christopher Chapp) is using map-reduce programming strategies (see below) on a St. Olaf's Beowulf cluster to perform political-science analysis of numerous political blogs over time, beginning with a study of approximately 400 editions of some 60 prominent liberal and conservative blogs during the 2008 election year.

## HiPerCiC/Political Blogs, and beyond

Mike Holm '11 and Mary Scaramuzza '12 are developing a HiPerCiC application that will enable students and faculty in Political Science to perform their own analyses of the blog data. In this case, data sets will be generated using providing *dictionaries* or word lists that indicate some political science issue, e.g., a “horserace” dictionary indicating competitive language or a dictionary of “health care” terms. The Beowulf programming tabulates appearances of dictionary entries among the blogs, producing a data set. Students and faculty will explore those results in HiPerCiC, and will be able to download those results for further analysis with a statistics package or in a spreadsheet.

The primary technique used for our Beowulf computations on political blogs is called *map-reduce*. Google Corp. developed this strategy for processing vast quantities of data in a reasonable amount of time using computing clusters, using undergraduate-level programming. Google employs map-reduce for analyzing everything from web-page contents for its search engine to graphical images together with business information for map and GPS services. We see rich and exciting possibilities for applying powerful cluster-computing techniques creatively to other disciplines across campus, in collaboration with student researchers.

## References

- Beowulf Overview: Frequently Asked Questions. Downloaded Nov. 1 2009 from <http://www.beowulf.org/overview/faq.html#17>
- Jeffrey Dean and Sanjay Ghemawat, MapReduce: Simplified Data Processing on Large Clusters. Retrieved August 22, 2008, from <http://labs.google.com/papers/mapreduce.html>
- Schade JD, Lewis DB. 2006. Plasticity in resource allocation and nitrogen-use efficiency in riparian vegetation: Implications for nitrogen retention. *Ecosystems* 9:740-755

## Acknowledgments

Funding sources Kay Winger-Blair, St. Olaf College, HHMI; faculty John Schade (Environmental Science), Chris Chapp (Political Science), Shilad Sen and Libby Shoop (Macalester Computer Science); students Mike Gesme '10, Megan Goebel '11, Tony Waldschmidt '08, and Summer 2009 CS undergraduate researchers.