

# Lawrence Berkeley National Laboratory

## LBL Publications

### Title

Efficient Graph Based Assembly of Short-Read Sequences on a Hybrid Core Architecture

### Permalink

<https://escholarship.org/uc/item/2ft9w97r>

### Authors

Sczyrba, Alex  
Pratap, Abhishek  
Canon, Shane  
et al.

### Publication Date

2011-03-22

# Efficient Graph Based Assembly of Short-Read Sequences on a Hybrid Core Architecture

**Alex Sczyrba**\*<sup>1,2</sup>, Abhishek Pratap<sup>1,2</sup>, Shane Canon<sup>2,3</sup>, James Han<sup>1,4</sup>, Alex Copeland<sup>1,2</sup>, Zhong Wang<sup>1,2</sup>, Tony Brewer<sup>5</sup>, David Soper<sup>5</sup>, Mike D'Jamoos<sup>5</sup>, Kirby Collins<sup>5</sup>, George Vacek<sup>5</sup>

<sup>1</sup>DOE Joint Genome Institute, Walnut Creek, CA, USA

<sup>2</sup>Lawrence Berkeley National Laboratory, Berkeley, CA, USA

<sup>3</sup>National Energy Research Scientific Computing Center (NERSC), Oakland, CA, USA

<sup>4</sup>Lawrence Livermore National Laboratory, Livermore, CA, USA

<sup>5</sup>Convey Computer Corporation, Richardson, TX, USA

March 2011

The work conducted by the U.S. Department of Energy Joint Genome Institute is supported by the Office of Science of the U.S. Department of Energy under Contract No. DE-AC02-05CH11231

## DISCLAIMER

This document was prepared as an account of work sponsored by the United States Government. While this document is believed to contain correct information, neither the United States Government nor any agency thereof, nor The Regents of the University of California, nor any of their employees, makes any warranty, express or implied, or assumes any legal responsibility for the accuracy, completeness, or usefulness of any information, apparatus, product, or process disclosed, or represents that its use would not infringe privately owned rights. Reference herein to any specific commercial product, process, or service by its trade name, trademark, manufacturer, or otherwise, does not necessarily constitute or imply its endorsement, recommendation, or favoring by the United States Government or any agency thereof, or The Regents of the University of California. The views and opinions of authors expressed herein do not necessarily state or reflect those of the United States Government or any agency thereof or The Regents of the University of California.

## Efficient Graph Based Assembly of Short-Read Sequences on a Hybrid Core Architecture

Alex Sczyrba\*<sup>1,2</sup>, Abhishek Pratap<sup>1,2</sup>, Shane Canon<sup>2,3</sup>, James Han<sup>1,4</sup>, Alex Copeland<sup>1,2</sup>, Zhong Wang<sup>1,2</sup>, Tony Brewer<sup>5</sup>, David Soper<sup>5</sup>, Mike D’Jamoos<sup>5</sup>, Kirby Collins<sup>5</sup>, George Vacek<sup>5</sup>

<sup>1</sup>DOE Joint Genome Institute, Walnut Creek, CA, USA

<sup>2</sup>Lawrence Berkeley National Laboratory, Berkeley, CA, USA

<sup>3</sup>National Energy Research Scientific Computing Center (NERSC), Oakland, CA, USA

<sup>4</sup>Lawrence Livermore National Laboratory, Livermore, CA, USA

<sup>5</sup>Convey Computer Corporation, Richardson, TX, USA

Advanced architectures can deliver dramatically increased throughput for genomics and proteomics applications, reducing time-to-completion in some cases from days to minutes. One such architecture, hybrid-core computing, marries a traditional x86 environment with a reconfigurable coprocessor, based on field programmable gate array (FPGA) technology. In addition to higher throughput, increased performance can fundamentally improve research quality by allowing more accurate, previously impractical approaches.

We will discuss the approach used by Convey’s de Bruijn graph constructor for short-read, *de-novo* assembly. Bioinformatics applications that have random access patterns to large memory spaces, such as graph-based algorithms, experience memory performance limitations on cache-based x86 servers. Convey’s highly parallel memory subsystem allows application-specific logic to simultaneously access 8192 individual words in memory, significantly increasing effective memory bandwidth over cache-based memory systems. Many algorithms, such as Velvet and other de Bruijn graph based, short-read, *de-novo* assemblers, can greatly benefit from this type of memory architecture. Furthermore, small data type operations (four nucleotides can be represented in two bits) make more efficient use of logic gates than the data types dictated by conventional programming models.

JGI is comparing the performance of Convey’s graph constructor and Velvet on both synthetic and real data. We will present preliminary results on memory usage and run time metrics for various data sets with different sizes, from small microbial and fungal genomes to very large cow rumen metagenome. For genomes with references we will also present assembly quality comparisons between the two assemblers.

The work conducted by the U.S. Department of Energy Joint Genome Institute is supported by the Office of Science of the U.S. Department of Energy under Contract No. DE-AC02-05CH11231