

Lawrence Berkeley National Laboratory

Recent Work

Title

DUK - A Fast and Efficient Kmer Matching Tool

Permalink

<https://escholarship.org/uc/item/4ff4315f>

Authors

Li, Mingkun
Copeland, Alex
Han, James

Publication Date

2011-03-22

DUK – A Fast and Efficient Kmer Matching Tool

Mingkun Li¹, Alex Copeland¹, James Han²

¹ Lawrence Berkeley National Laboratory

² Lawrence Livermore National Laboratory

February 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.

DUK – A Fast and Efficient Kmer Matching Tool

***Mingkun Li**¹, Alex Copeland¹ and James Han²

¹Lawrence Berkeley National Laboratory

²Lawrence Livermore National Laboratory

US Department of Energy Joint Genome Institute, 2800 Mitchell Drive, Walnut Creek, California 94598, USA

Email: mli@gmail.com

Abstract

A new tool, DUK, is developed to perform matching task. Matching is to find whether a query sequence partially or totally matches given reference sequences or not. Matching is similar to alignment. Indeed many traditional analysis tasks like contaminant removal use alignment tools. But for matching, there is no need to know which bases of a query sequence matches which position of a reference sequence, it only need know whether there exists a match or not. This subtle difference can make matching task much faster than alignment. DUK is accurate, versatile, fast, and has efficient memory usage. It uses Kmer hashing method to index reference sequences and Poisson model to calculate p-value. DUK is carefully implemented in C++ in object oriented design. The resulted classes can also be used to develop other tools quickly. DUK have been widely used in JGI for a wide range of applications such as contaminant removal, organelle genome separation, and assembly refinement. Many real applications and simulated dataset demonstrate its power.

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

LLNL-ABS-475246