

















High-Performance Computing for Assembly and Analysis of Big Genomics Data

René Warren, Benjamin Vandervalk, Anthony Raymond, Shaun Jackman, Hamid Mohamadi, Daniel Paulino, Justin Chu, Ewan Gibb and İnanç Birol

Abstract

DNA Sequencing technology is developing at an unprecedented pace, surpassing the rate of advances in computer hardware development. Limited compute resources for storing, processing and analyzing omics data have spurred the improvements of file compression formats, low memory footprint data structures, algorithms that use communication protocols for parallel programming, and sature approaches for handling large data on commodity hardware. Our research team oversees the development of such bioinformatics technologies. Past accomplishments include: enabling the first assembly with millions of very short sequence reads (Warren et al. 2006), assembly of the human genome from short reads with the first parallel assembler (Simpson et al. 2009) and last year, assembly of the then largest genome, that of the 20 Gbp white spruce (Birol et al. 2013). We discuss key enabling algorithms, specifically introducing data structures, processes, compression schemes within ABPSS (Simpson et al. 2009).

BBT (Chu et al. 2014), DIDA (Mohamadi et al. submitted), Konnector (Vandervalk et al. 2014) and TASR (Warren et al. 2011) that are tailored to the needs of today's big sequence data reality. DNA Sequencing technology is developing at an unprecedented pace, surpassing the

Funding

BC Cancer Foundation



References

Reirol I, Raymond A, Jackman SD, Pleasance S, Coope R, Taylor GA, Yuen MM, Keeling CI, Brand D, Vandervalk BP et al. 2013. Assembling the 20 Gb white spruce (Picea glauca) genome from whole-genome shotgun sequencing data. Bioinformatics. 29:1492

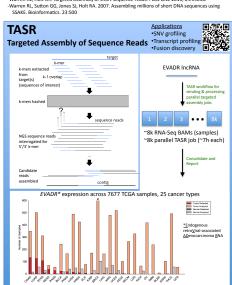
-Bloom BH. 1970. Space/time trade-offs in hash coding with allowable errors. Commun ACM. 13:422

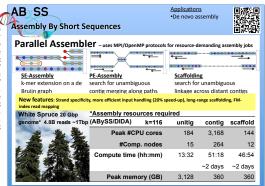
-Chu J. Sadeghi S, Raymond A, Jackman SD, Wi KM, Mar R, Mohamadi H, Butterfield YS, Robertson A. Brid 1. 2014. Bolloom took: fast, accurate and memory-efficient host Species sequence screening Birol 1. 2014. BioBloom took: fast, accurate and memory-efficient host species sequence screen using bloom filters, Bioinformatics, 30-3002
Simpson JT, Wong K, Jackman SD, Schein JF, Jones SJ, Birol I. 2009. ABySS: a parallel assembler for short read sequence data. Genome Res. 19-1117
Vandervalk B, Jackman SD, Raymond A, Mohamadi H, Yang C, Attail D A, Chu J, Warren RJ, Birol I. 2014. Konnector: connecting paired end reads using a Bloom filter de Bruijn graph. in Bioinform and Biomedicine (BIBM 2014), Belfast UK, 2014.

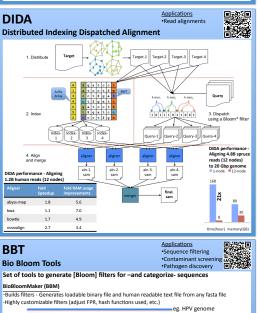
Amounte (Note 2014), Bettles UN, 2014.

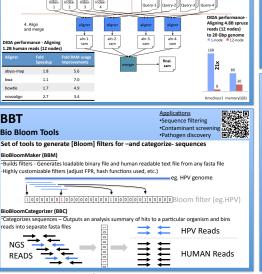
-Warren RL, Holt RA: Targeted assembly of short sequence reads. PLoS ONE 2011, 6:e19816.

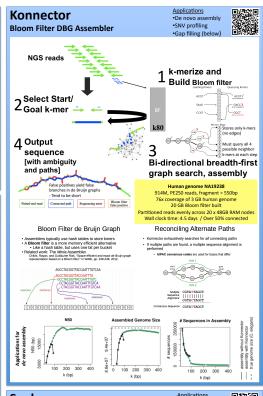
-Warren RL, Sutton GG, Jones SJ, Holt RA: 2007. Assembling millions of short DNA sequences using SSAKE. Bioinformatics. 23:500

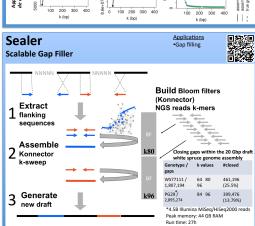












Genome Sciences Centre • British Columbia Cancer Agency • 100 - W 7th Ave Vancouver BC V5Z 4S6 Canada • tel. 604-707-5900 • fax. 604-876-3561 • www.bcgsc.ca