

# Developing an Arvados BWA-GATK pipeline

Pjotr Prins<sup>1,2</sup>, Joep de Lig<sup>3</sup>, Isaac Nijman<sup>1</sup>, Ward Vandewege<sup>4</sup>, Bryan Cosca<sup>4</sup> and Peter Amstutz<sup>4</sup>

<sup>1</sup> University Medical Center Utrecht, Utrecht, The Netherlands

<sup>2</sup> University of Tennessee Health Science Center, Memphis, USA

<sup>3</sup> Hubrecht Institute, Utrecht, The Netherlands

<sup>4</sup> Curoverse, Boston, USA<sup>†</sup>

Contact: [pjotr.public05@thebird.nl](mailto:pjotr.public05@thebird.nl)

Project website: <http://arvados.org/>

Source code: <https://arvados.org/projects/arvados/repository>

License: AGPLv3 for core, Apache 2.0 for the SDK

Arvados is a free and open source platform for bioinformatics and big data science. In this (lightning) talk, we present the porting of an existing HiSeq BWA-GATK based variant calling pipeline from Sun Grid Engine (SGE) to Arvados, resulting in a pipeline that is faster, scalable, simpler and more robust.

The main reason the pipeline runs faster is that the Arvados file system (named Keep) is decentralized. In most bioinformatics cluster environments, the central file storage is the bottleneck because of resource contention, i.e. many cluster nodes are hitting the storage for data requests. When uploading a file into Keep the file is chunked and distributed across nodes. When a node requires data, it can fetch it from different locations in the network.

Because of Keep, the new pipeline is scalable and runs in flat time, i.e. processing one genome takes the same amount of time as running ten or a hundred. Running pipelines in parallel is now only a function of adding new nodes. Keep is distributed and scales accordingly.

The original SGE pipeline consisted of hundreds of lines of Perl and bash code, which mostly dealt with the plumbing of submitting jobs, checking conditions and job completion. In contrast, the Arvados version is mostly single lines for command line invocation as Arvados takes care of the plumbing.

We also needed to speed up GATK using GATK/Queue. Queue chunks BAM files and fans work out to multiple nodes. For this, a special Queue adapter had to be written for Arvados similar to the one that exists for SGE. The Arvados implementation is an improvement over the SGE implementation because it avoids using mounted NFS for sharing and collating Queue results, thereby making GATK/Queue more robust.

---

<sup>†</sup>Disclosure: Curoverse is a major contributor to the Arvados open source project and a sponsor of BOSC 2015.