

Nextflow and Singularity enabling heterogeneous CPU/GPU bioinformatics containers in workflows

Kevin Sayers
September, 2017



Overview

- Technologies
 - Nextflow
 - Containers
 - GPU
 - AWS
 - Machine learning
- SRAGPU-nf workflow
- Results
- Conclusions

nextflow



Nextflow

- Domain specific workflow language for bioinformatics
- Deploy workflows locally, in an HPC cluster, or on Amazon web services (AWS)
- Supports Docker and Singularity containers
- GitHub integration
- Automatic parallelization

Nextflow script

```
1 process setup{
2     container = "docker://sayerskt/samtools"
3     publishDir './', mode: 'copy', overwrite: true
4
5     output:
6     file "Homo_sapiens.GRCh38.cdna.all.fa" into reference
7     file "Homo_sapiens.GRCh38.cdna.all.fa.fai" into refindex
8
9     """
10    wget ftp://ftp.ensembl.org/pub/release-88/fasta/homo_sapiens/cdna/Homo_sapiens.GRCh38.cdna.all.fa.gz
11    gunzip Homo_sapiens.GRCh38.cdna.all.fa.gz
12    samtools faidx Homo_sapiens.GRCh38.cdna.all.fa
13    """
14
15 }
16 process barracudaIndex{
17     container = 'shub://KevinSayers/BarraCUDA_Singularity'
18
19     storeDir 'index/'
20     input:
21     file ref from reference
22
23     output:
24     file "${reference.baseName}.*" into indexOut, indexFiles
25
26     """
27     barracuda index -p ${reference.baseName} ${ref}
28
29     """
30
31 }
```

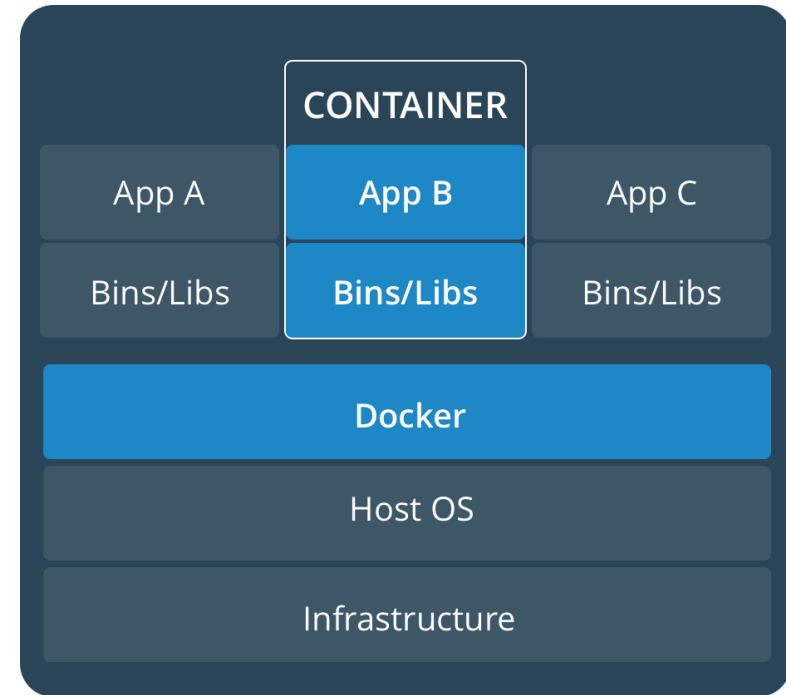
```
1 singularity{
2     enabled = true
3     runOptions = '--nv'
4 }
5
6 cloud {
7     imageId = 'ami-f379948a'
8     instanceType = 'p2.xlarge'
9     spotPrice = 0.40
10    bootStorageSize = '800 GB'
11 }
```

Containers

- Software portability
 - Packaged together
 - No conflicting dependencies
- Reproducibility
 - Created from an image
 - Control over software versions
- Bioinformatics respositories



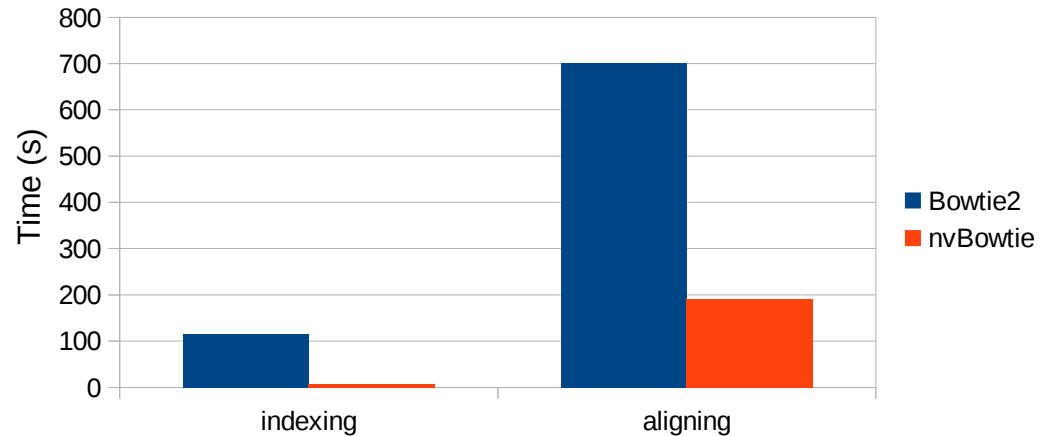
Biocontainers



<https://www.docker.com/what-container>

GPUs in Bioinformatics

- Thousands of parallel cores
- Can substantially reduce processing time
- Bioinformatics tools
 - Short read aligners (BarraCUDA, nvBowtie, SOAP3)
 - GPU-BLAST
 - Molecular modeling
- Machine learning



GPU containers

- Portability of difficult to compile GPU programs
- Recent work
- Implementations
 - Nvidia-docker
 - Shifter
 - Singularity
- Containerized single tools



Portable, high-performance containers for HPC

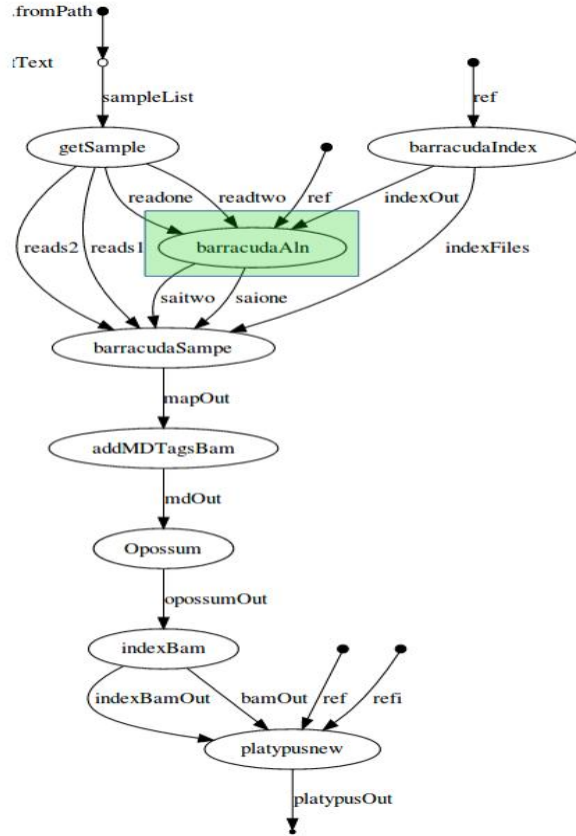
Lucas Benedicic*, Felipe A. Cruz, Alberto Madonna, Kean Mariotti
Systems Integration Group
CSCS, Swiss National Supercomputing Centre
Lugano, Switzerland
Email: *lucas.benedicic@cscs.ch

Experimental data

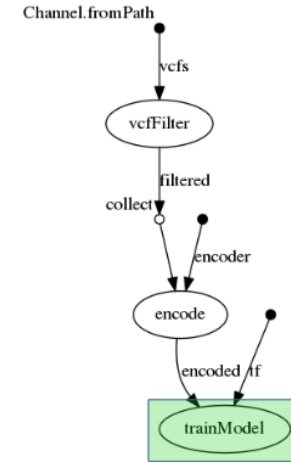
- Single cell RNA-seq
- 4 classes
 - Estrogen receptor positive (ER+)
 - Human epidermal growth factor receptor 2 (HER2+)
 - Triple-negative breast cancer (TNBC)
 - ER2+ and HER2+
- 381 samples
- Reference transcriptome
- BarraCUDA GPU enabled short read aligner
- TensorFlow machine learning classifier

SRAGPU-nf

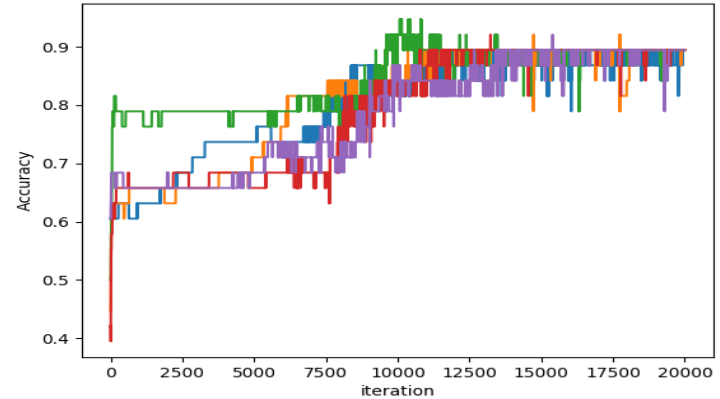
Sample processing



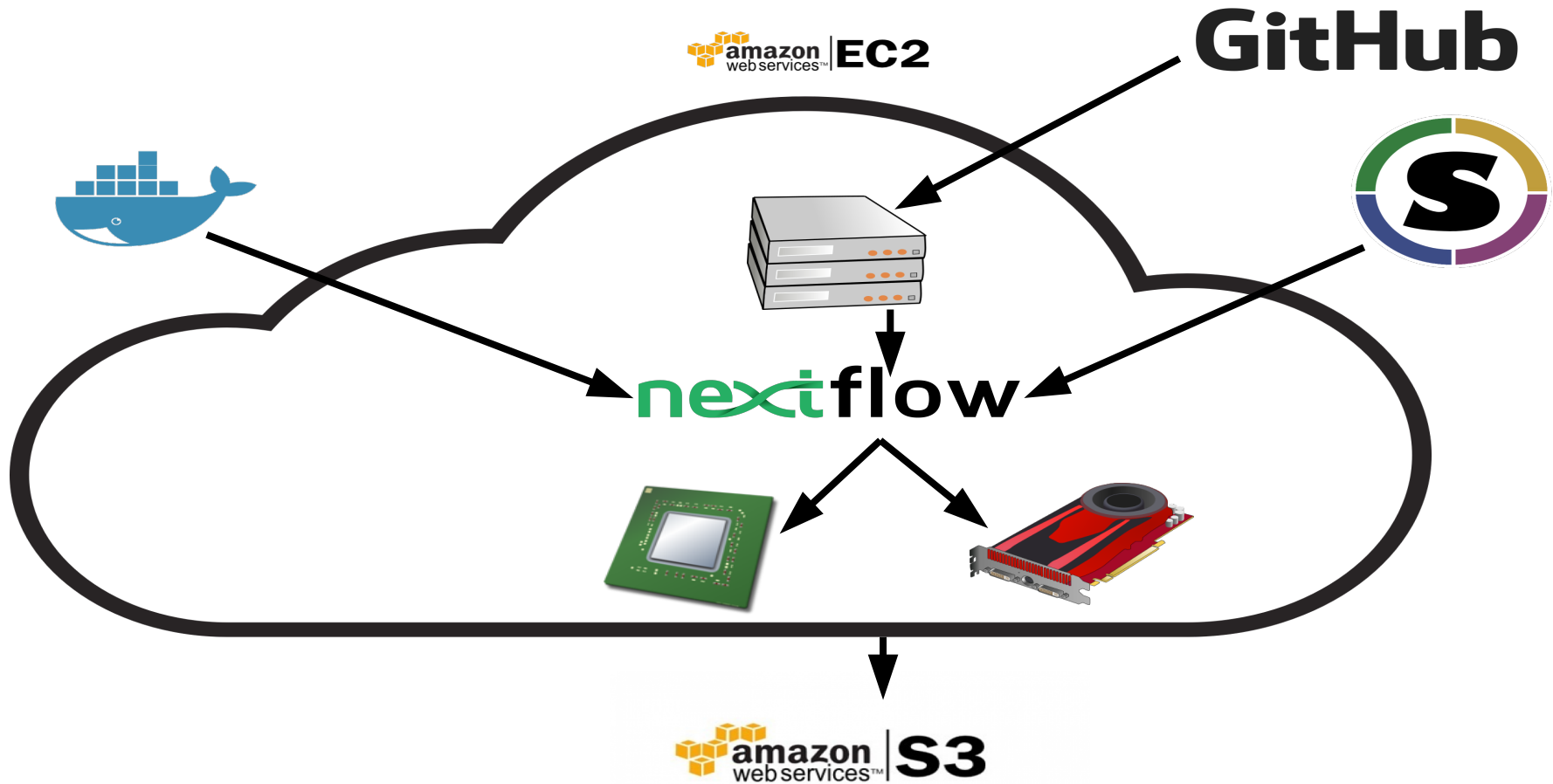
Classification



Testing accuracy



Cloud configuration



Conclusions

- GPU bioinformatics tools can decrease processing times
- Nextflow and Singularity can be used to deploy heterogeneous container workflows
- Publicly published nvBowtie and BarraCUDA GPU containers
- Machine learning steps can be incorporated into Nextflow using GPU containers
- Opening up new workflow possibilities that rely heavily on GPU based tools

Future works

- Add support for each process being a different instance type
- Singularity still actively developing GPU features
- Improve the parameters for the SRAGPU-nf workflow
- Improve the machine learning models
- Validate with other RNA-seq data

Thank you for listening!

1. Di Tommaso, P., Chatzou, M., Floden, E. W., Barja, P. P., Palumbo, E., & Notredame, C. (2017). Nextflow enables reproducible computational workflows. *Nature Biotechnology*, 35(4), 316-319.
2. Benedicic, L., Cruz, F. A., Madonna, A., & Mariotti, K. (2017). Portable, high-performance containers for HPC. *arXiv preprint arXiv:1704.03383*.
3. Kurtzer, G. M., Sochat, V., & Bauer, M. W. (2017). Singularity: Scientific containers for mobility of compute. *PloS one*, 12(5), e0177459.
4. Oikkonen, L., & Lise, S. (2017). Making the most of RNA-seq: Pre-processing sequencing data with Opossum for reliable SNP variant detection. *Wellcome open research*, 2.
5. Abadi, M., Agarwal, A., Barham, P., Brevdo, E., Chen, Z., Citro, C., ... & Ghemawat, S. (2016). Tensorflow: Large-scale machine learning on heterogeneous distributed systems. *arXiv preprint arXiv:1603.04467*.

Questions?

- BarraCUDA: [shub://KevinSayers/BarraCUDA_Singularity](https://github.com/KevinSayers/BarraCUDA_Singularity)
- nvBowtie: [shub://KevinSayers/nvBowtie_Singularity](https://github.com/KevinSayers/nvBowtie_Singularity)
- SRAGPU-nf: <https://github.com/KevinSayers/SRAGPU-nf>
- OneHotVCF: <https://github.com/KevinSayers/OneHotVCF>
- Paper in the works!
- I presented portions of this work at the 2017 Bioinformatics Open Source Conference (BOSC) in Prague as a lightning talk.