

Harnessing cloud computing for high capacity analysis of neuroimaging data from NDAR

Daniel Clark¹, Christian Haselgrove², David Kennedy², Zhizhong Liu³,
Michael Milham¹, Petros Petrosyan⁴, Carinna Torgerson³, John Van Horn³, Cameron Craddock^{1,5}

¹Child Mind Institute, New York, NY, ² University of Massachusettes Medical School, Worcester, MA, ³University of Southern California, Los Angeles, CA, ⁴UCLA, Los Angeles, CA, ⁵Nathan S. Kline Institute for Psychiatric Research, Orangeburg, NY



Introduction

- ▶ The National Database for Autism Research (NDAR)¹ hosts a vast collection of neuroimaging datasets that can be processed and utilized to yield significant scientific discoveries.
- ▶ This amount of resources necessitates a high-performance computing (HPC) infrastructure, which is not always readily available for researchers in-house.
- ▶ Amazon Web Services (AWS) Elastic Compute Cloud (EC2) computing service offers a “pay as you go” model that allows researchers to utilize HPC performance without the up-front capital costs and maintenance of an in-house solution.
- ▶ The developers of the Laboratory of Neuro Imaging (LONI) Pipeline, the Neuroimaging Informatics Tools and Resources Clearinghouse (NITRC) Computational Environment (CE) and the Configurable Pipeline for the Analysis of Connectomes (C-PAC) have implemented pipelines in EC2 that interface with NDAR.

Methods

LONI Pipeline

- ▶ The LONI Pipeline software was extended to include new pipeline modules to access data from the NDAR database, transfer input data out of Amazon S3 (Simple Storage Service), and to load results back into S3².
- ▶ A pipeline was constructed to extract cortical thickness and subcortical region volume data from structural MRI images in the NDAR database, which included:
 1. Reorient images to standard orientation using FSL's reorient2std module
 2. Extract cortical thickness using FreeSurfer recon-all
 3. Calculate volumes of subcortical regions using FSL's BET and FIRST all
- ▶ The resulting pipeline was used to process 780 T1-weighted structural images and return the results to NDAR.

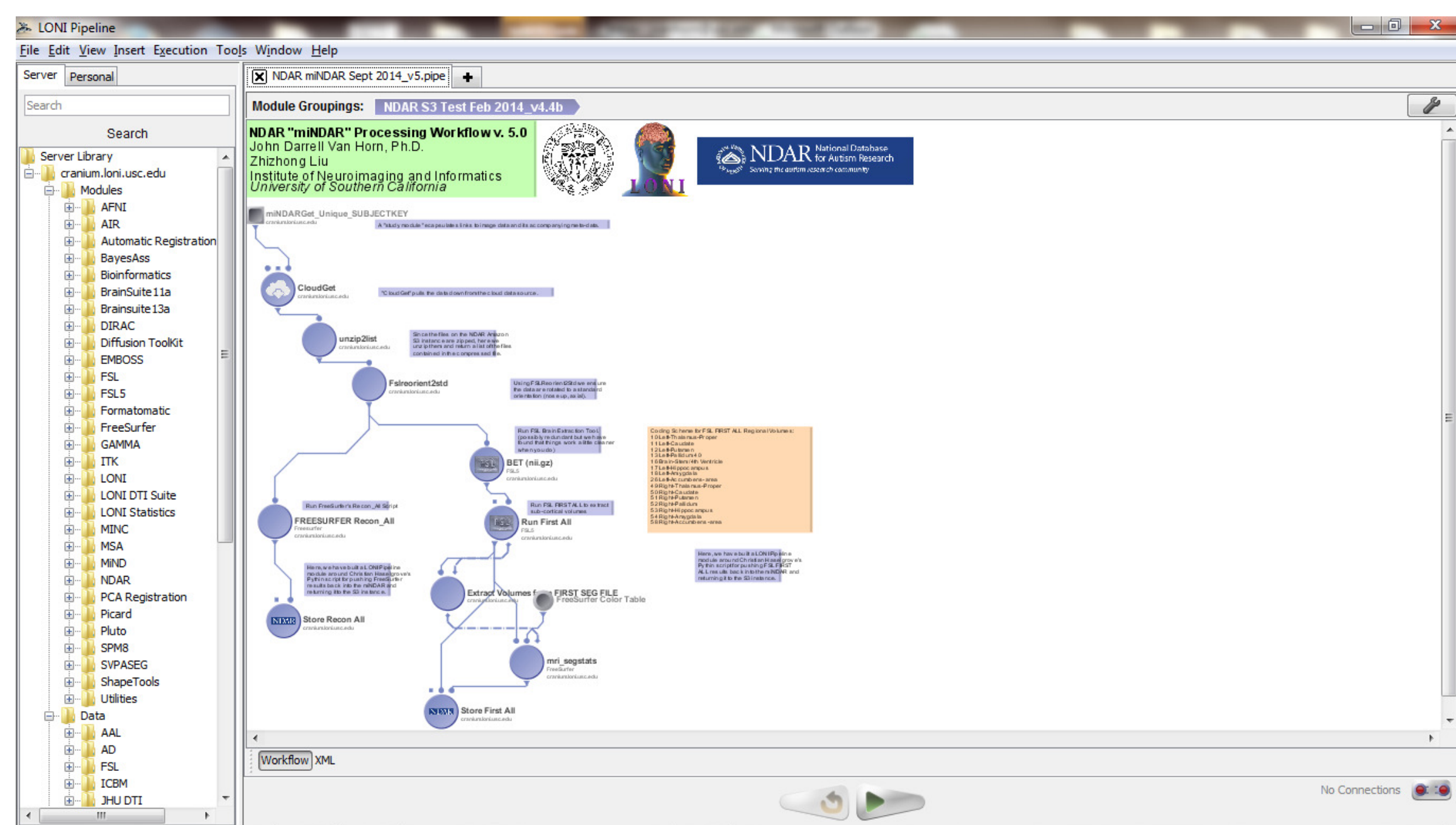


Figure 1 : Graphical layout of the constructed pipeline

Figure 2 : miNDAR database

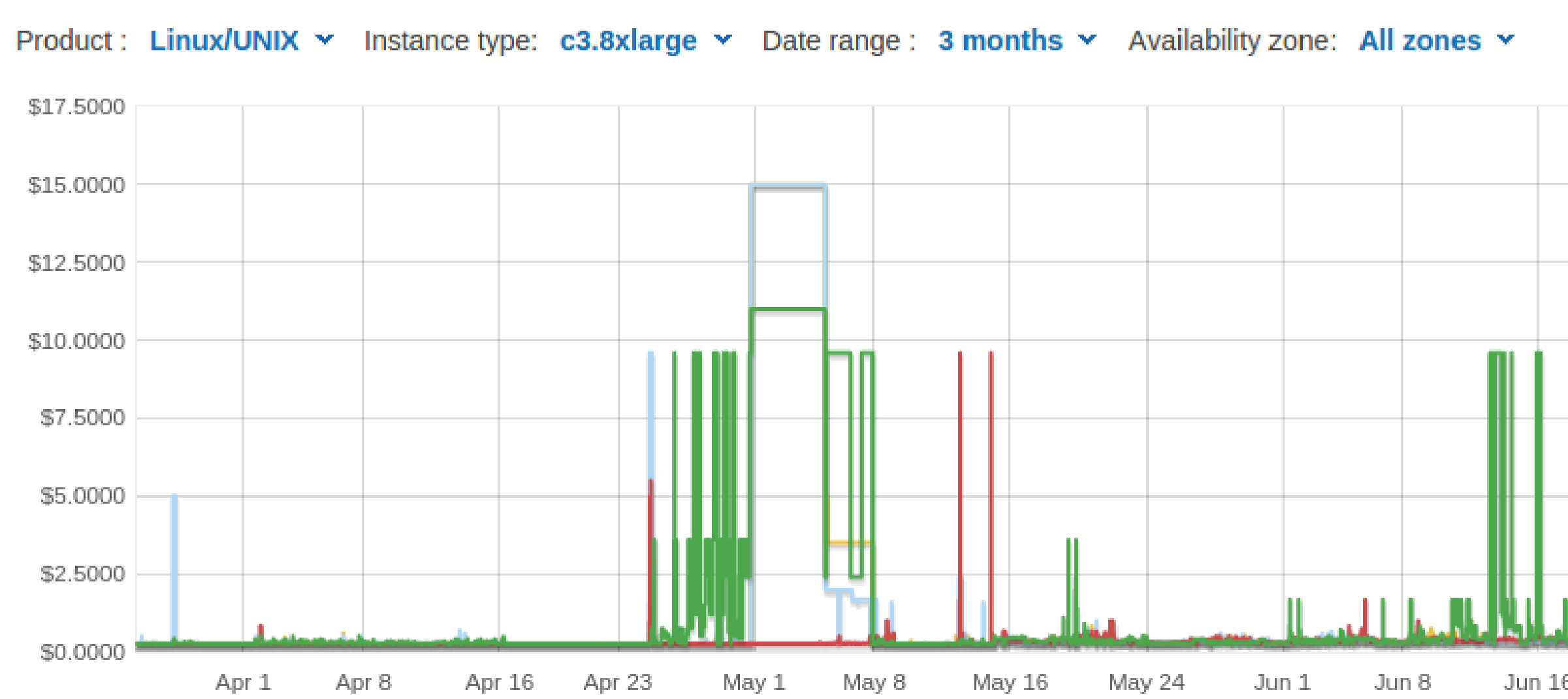


Figure 3 : Spot price history for the c3.8xlarge instance in the us-east-1 region on AWS across the four different availability zones for the past three months

Interacting with the NDAR database

- ▶ An AWS-hosted miNDAR database was created by querying the NDAR website for the data of interest (e.g. from a particular study).
- ▶ The developed modules were then used to build a subject list by querying the database for subjects of interest to pass to the pipeline.
- ▶ The launch and configuration of an AWS EC2 HPC cluster was handled using Starcluster (<http://star.mit.edu/cluster/>).
- ▶ A process run was submitted from the cluster head node as a Sun Grid Engine job using our pipeline software and the subject list.
- ▶ The pipeline software then processed the data, stored raw outputs in an S3 bucket, and inserted the S3 filepaths and output measures into the miNDAR database.

EC2 Spot Pricing

- ▶ In addition to the “pay as you go” compute pricing model, AWS offers users to bid on instances that are not being used.
- ▶ As opposed to paying an “on-demand” fixed price, the hourly charge for the instances will vary over time based on what people in different regions are willing to pay.
- ▶ If the spot market price goes above what the user bids, their instance will be terminated by Amazon and all of their processes will be stopped and their data lost.
- ▶ Typically, one can expect to pay one-eighth of the hourly price of a standard on-demand instance when using spot-pricing.
- ▶ AWS offers cloud solutions based on geographic region, where the price for services vary.
- ▶ Some regions fluctuate on spot price more than others, but there are ways to minimize process interruptions and lost data in how one configures their cluster and pipelines.

Results

Table 1 : Processing completed as a part of the effort. Nodes corresponds to the number of hosts used in the calculation. PF is parallelization factor and corresponds to the number of jobs ran in parallel on each node. On demand instances were used for the master node and spot instances were used for all computation nodes in the cluster. CPU Time is the total amount of time required to perform the computation and Wall Time is the amount of time that passed. # DS: Number of datasets. CPD: Cost Per Dataset. C-PAC: Configurable Pipeline for the Analysis of Connectomes. NITRC-CE: NITRC Computational Environment

Processing	# DS	Platform	Nodes	PF	CPU Time	Wall Time	Cost	CPD
ANTS Cortical Thickness	3197	C-PAC	20	8	23,018	147	\$760.24	\$0.24
Resting state fMRI processing w/ 4 strategies	1112		20	3	834	22	\$80.54	\$0.07
Quality Assessment Protocol	1112		20	4	380	14	\$19.02	\$0.02
FreeSurfer recon-all	986	NITRC-CE	4	32	23,664	193	\$211.44	\$0.21
FSL FIRST	1247		4	32	208	3	\$2.19	> \$0.01
Temporal QA	1349		4	32	450	13	\$4.69	> \$0.01
FreeSurfer recon-all and FSL FIRST	780	LONI Pipeline	20	32	18,720	49	\$252.36	\$0.32

Benefits of the cloud

- ▶ The scalability and payment model for processing this data on a cloud platform allowed for an efficient processing of the NDAR datasets, both in time and cost.

Results cont.

Costs and run time models

- ▶ Models assume that users upload their input data to the cloud, run their pipelines and (optionally) download their outputs as they become available.
- ▶ Outputs are stored on an AWS Elastic Block Store (EBS) hard drive mounted to the master node and is NFS-shared across all of the slave nodes.
- ▶ Users can also upload their processed outputs to a cloud storage solution, like AWS S3, directly from the cluster - this was the approach taken in processing data for NDAR.
- ▶ Storing results in S3 avoids costly download time and provides for a viable solution for backing up data; however this does incur additional storage costs.

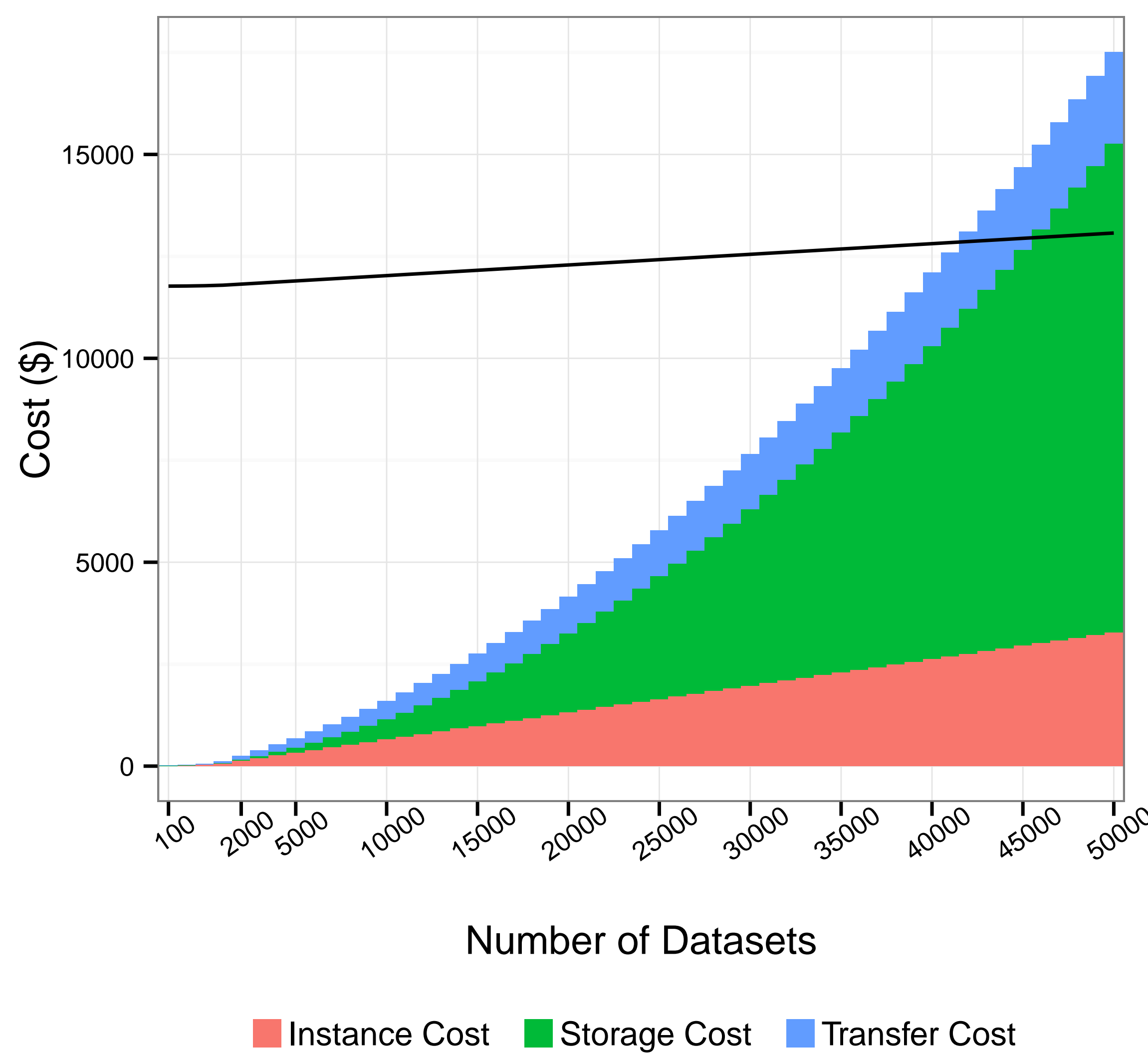


Figure 4 : AWS EC2 costs, grouped by cost type, for a typical C-PAC pipeline for different sized datasets versus owning and maintaining own server; costs are based on using the cluster configuration shown in Table 1 for the cloud, and a single compute node locally

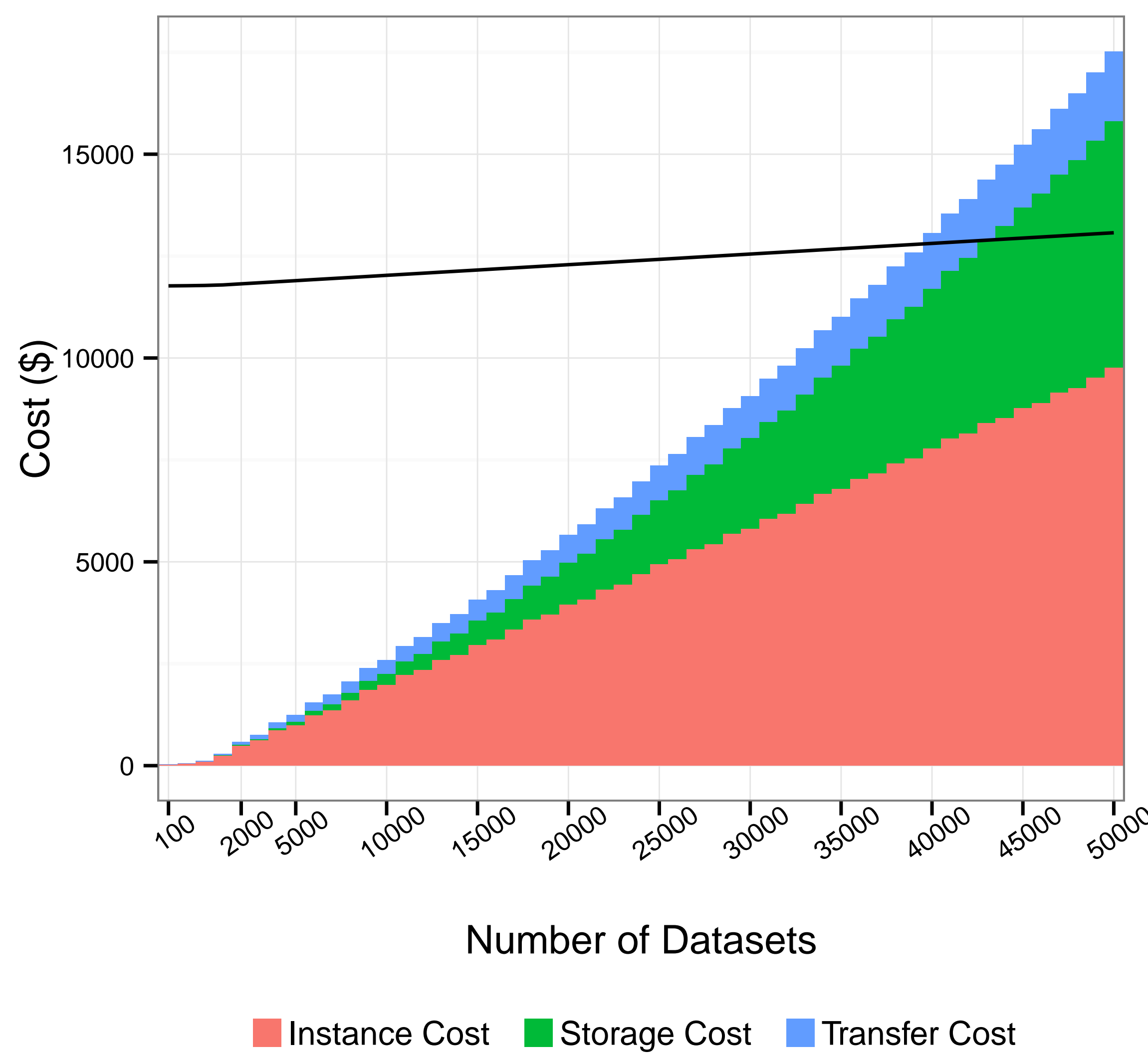


Figure 5 : AWS EC2 costs, grouped by cost type, for a typical FreeSurfer pipeline for different sized datasets versus owning and maintaining own server; costs are based on using the cluster configuration shown in Table 1 for the cloud, and a single compute node locally

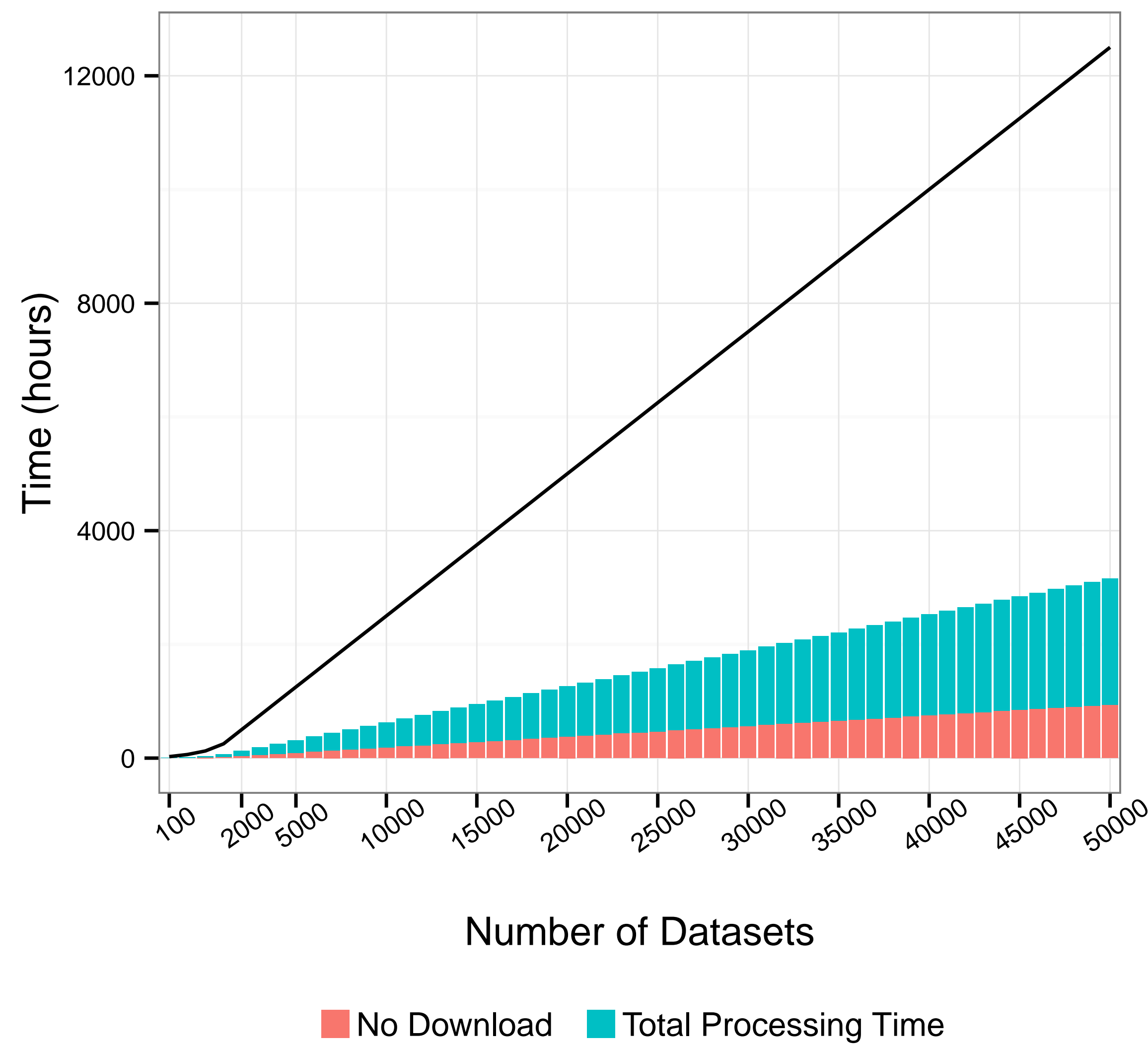


Figure 6 : Run times using the C-PAC pipeline, grouped by downloading output data or not, for different sized datasets versus running locally; times are based on using the cluster configuration shown in Table 1 for the cloud, and a single compute node locally

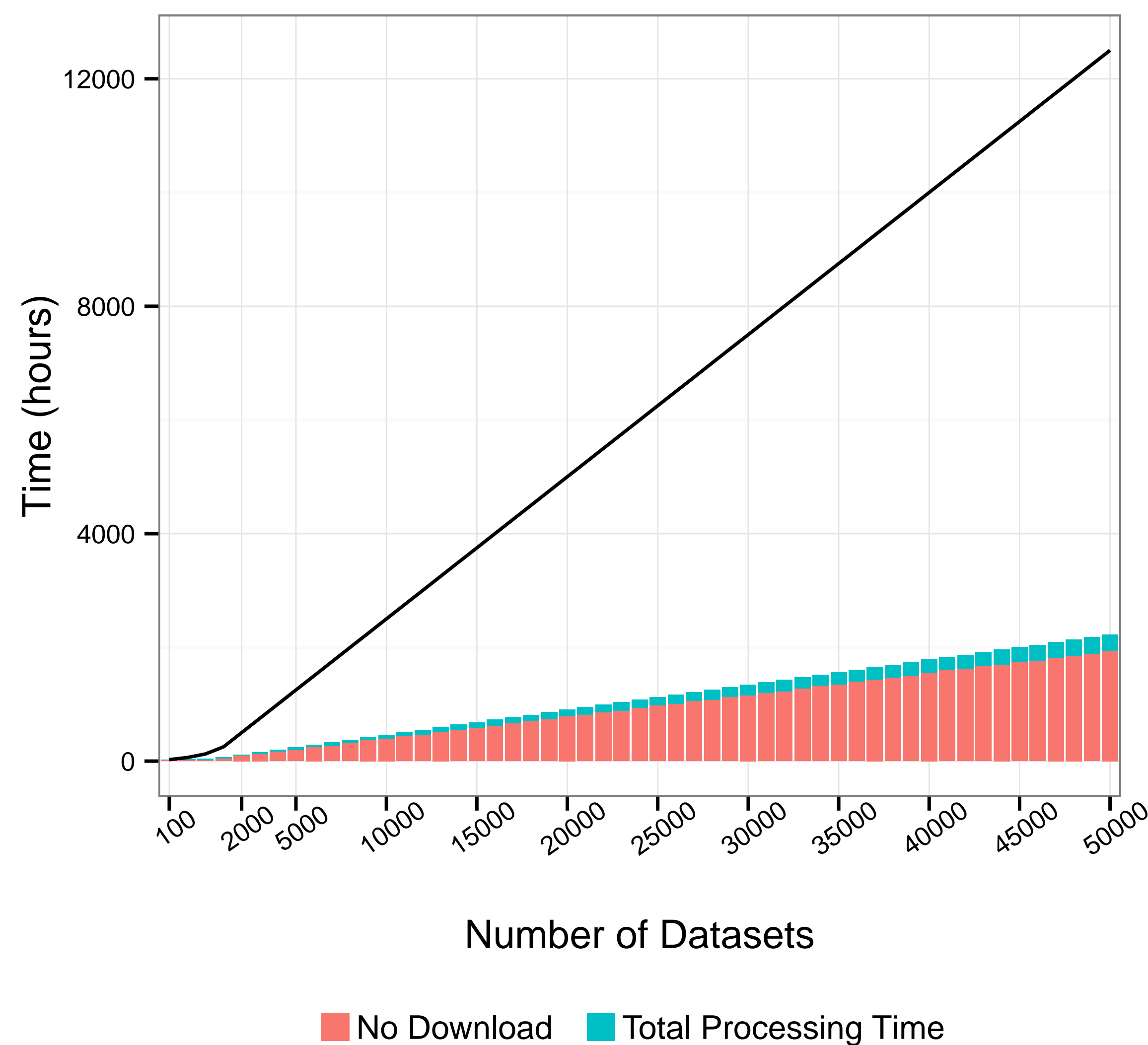


Figure 7 : Run times using the FreeSurfer pipeline, grouped by downloading output data or not, for different sized datasets versus running locally

Conclusion

- ▶ NDAR provides a wealth of heterogenous data for studying the effects of autism and an excellent framework for the querying and storing of scientific results.
- ▶ Cloud computing services, like AWS EC2, are viable solutions for the processing and analysis of large amounts of data.
- ▶ A time and cost analysis argues for the benefits of using cloud computing over purchasing and maintaining one's own computing resource.

References and Acknowledgements

1. Hall, D. et al. (2012), Neuroinformatics. 10(4):331-9
2. Torgerson, C.M. et al. (2015), Brain Imaging and Behavior 9:89-103
3. Craddock, R.C. et al. (2013), Frontiers in Neuroinformatics 42
4. Das, S.R. et al. (2009), Neuroimage, 45(3): 867-79
5. Shehzad, Z. et al. (2014), Resting State Conference, Boston
6. Haselgrove, C. et al. (2014), Front. Neuroinform. 8:52

Data processing and salary support was provided by NDAR.