

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

Daniel Clark<sup>1</sup>, Christian Haselgrove<sup>2</sup>, David Kennedy<sup>2</sup>, Zhizhong Liu<sup>3</sup>, Michael Milham<sup>1</sup>, Petros Petrosyan<sup>4</sup>, Carinna Torgerson<sup>3</sup>, John Van Horn<sup>3</sup>, Cameron Craddock<sup>1</sup> <sup>1</sup>Child Mind Institute, New York, NY, <sup>2</sup> University of Massachuttes Medical School, Worcester, MA, <sup>3</sup>University of Southern California, Los Angeles, CA, <sup>4</sup>UCLA, Los Angeles, CA, <sup>5</sup>Nathan S. Kline Institute for Psychiatric Research, Orangeburg, NY





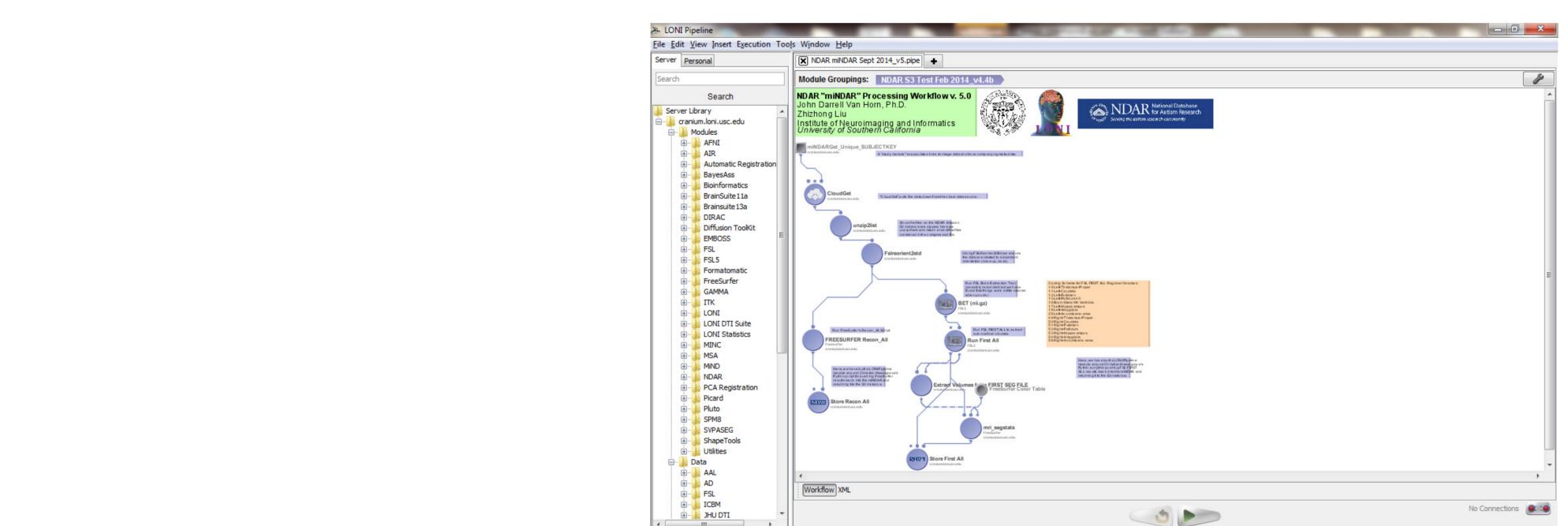
Introduction

- $\triangleright$  The National Database for Autism Research (NDAR)<sup>1</sup> 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 captial
- ▶ The developers of the Laboratory of Neuro Imaging (LONI) Pipeline, the Neuroimaging Informatics Tools and Resources Clearinghouse (NITRC) Computational Environment (CE) and the

costs and maintenance of an in-house solution.

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 S3<sup>2</sup>
- ► A pipeline was constructed to extract cortical thickness and subcortical region volume data from structural MRI images in the NDAR database, which included:
- .. Reorient images to standard orientation using FSL's reorient2std module 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

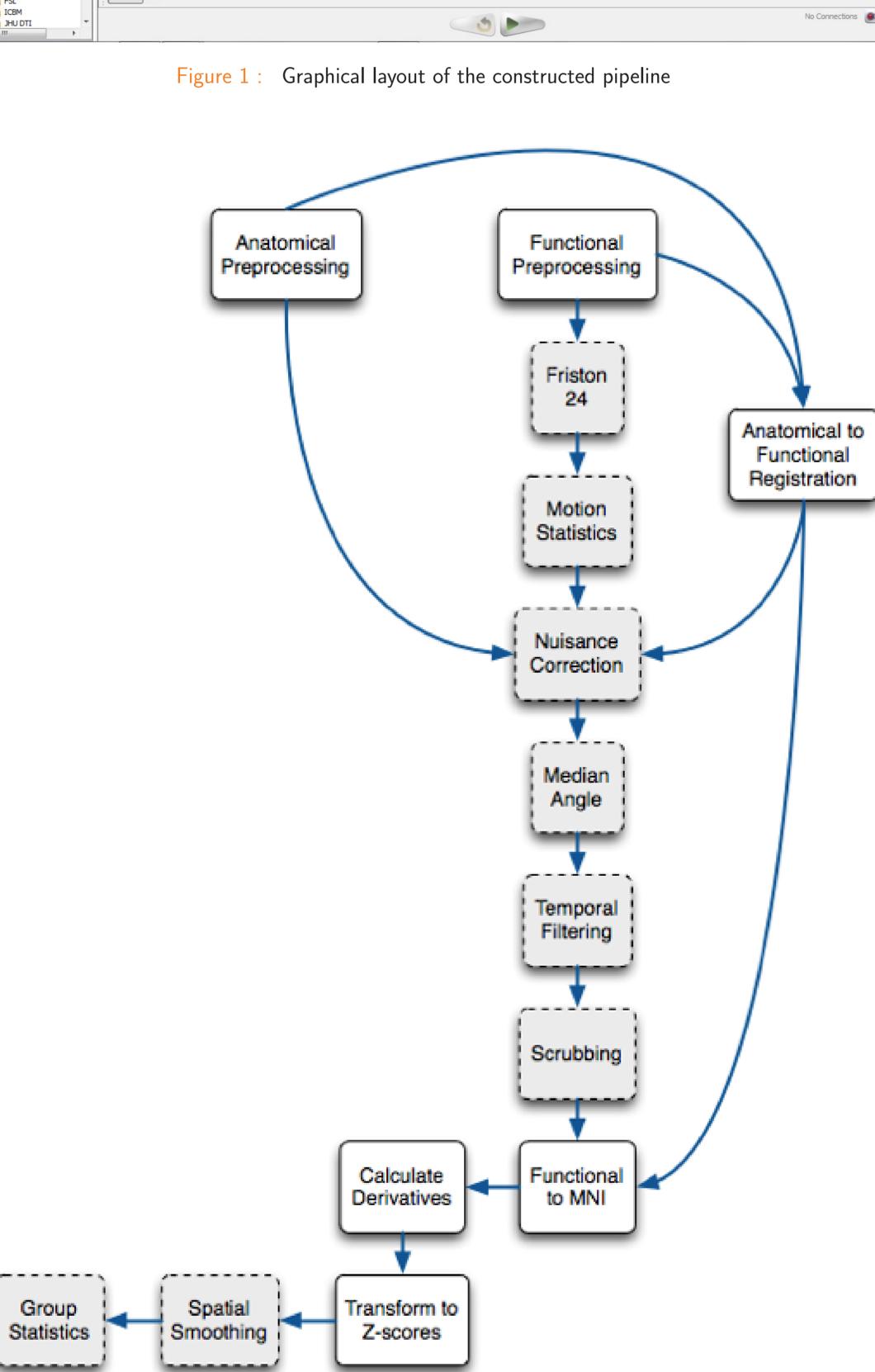


Figure 2: C-PAC preprocessing stages

## Configurable Pipeline For the Analysis of Connectomes (C-PAC)<sup>3</sup>

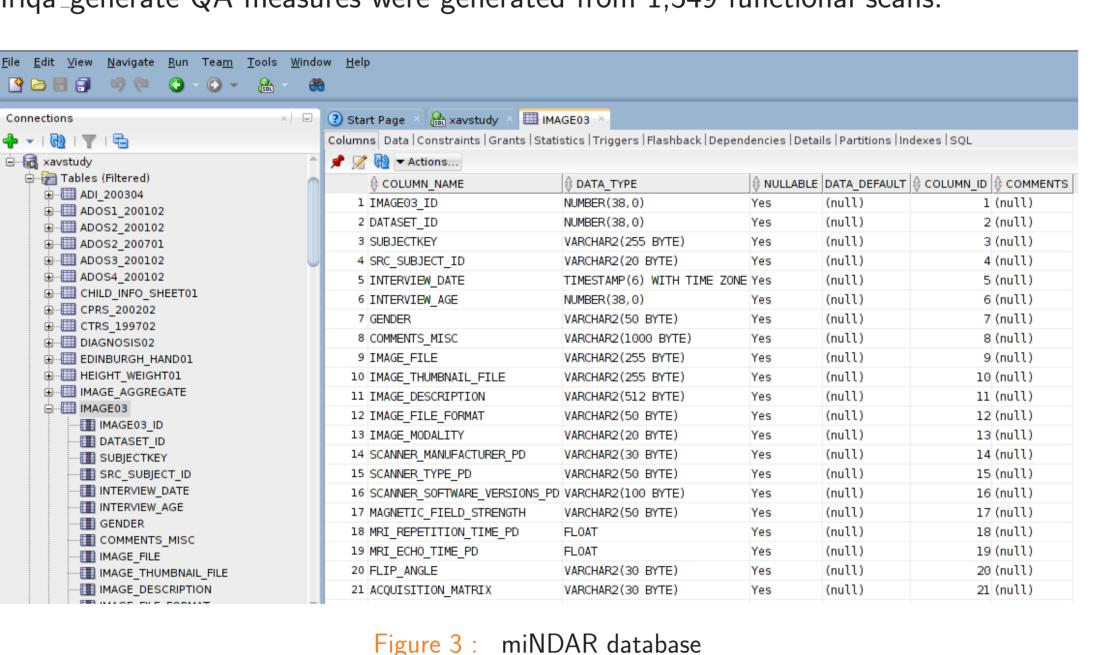
- ► C-PAC modules were written in Python to build input data lists by querying NDAR, read input data from S3, write processed results to S3 and write values back to the NDAR database ▶ New pipelines were created to perform the ANTS cortical thickness⁴ procedure and the Preprocessed Connectomes Project's Quality Assessment Protocol
- $(http://preprocessed-connectomes-project.github.io/quality-assessment-pipeline)^5$ ▶ The resulting pipelines and modules were used to process several datasets and return the results to NDAR
- 1. Cortical extraction from 3,197 T1-weighted structural images
- 2. Structural and functional processing for 1,112 datasets from ABIDE

### 3. Automated quality asssessment of 1,112 datasets from ABIDE Neuroimaging Informatics Tools and Resources Clearinghouse Computational Environment (NITRC-CE)<sup>6</sup>

- ► The NITRC pipeline processed data using three primary utilities
- .. Extract anatomical and surface-base measures with Freesurfer recon-all 2. Segment subcortical structures using FSL's FIRST to produce volumetric and mesh outputs
- 3. Time series QA measures using the fmriqa\_generate.pl utility from the BXH/XCEDE Tools suite, including mean intensty, center of mass, per-slice variation, and others > Python modules were created to query and download data from NDAR as well as to store results back to their database
- ► The recon-all and FIRST tools processed 986 and 1,247 T1-weighted anatomical scans, respectively; the fmriqa\_generate QA measures were generated from 1,349 functional scans.

## Interacting with the NDAR database

- ► Launched an AWS-hosted miNDAR database by querying NDAR website for the data of interest (e.g. from a particular study)
- ▶ Built a subject list by querying the database for subjects of interest to pass to our pipeline ► Launch an AWS EC2 HPC cluster using Starcluster<sup>7</sup>
- ► Log into the cluster and submit a Sun Grid Engine job using our pipeline software and the subject list ▶ The pipeline software will process the data, store raw outputs in an AWS S3 bucket and
- insert S3 filepaths and output measures into 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 fluxuate 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



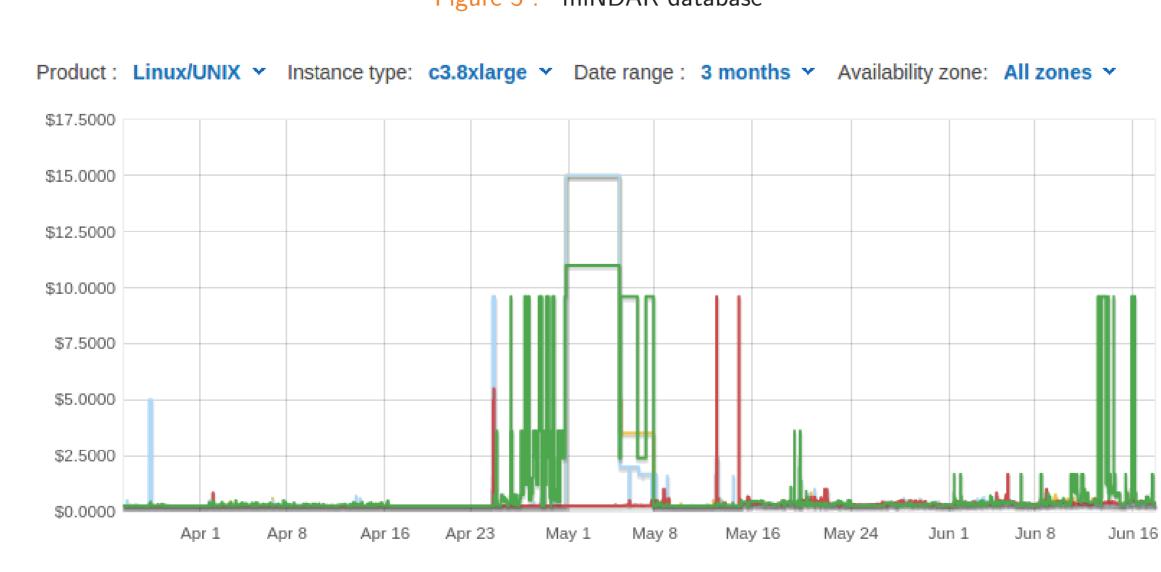


Figure 4: Spot price history for the c3.8xlarge instance in the us-east-1 region on AWS across the four

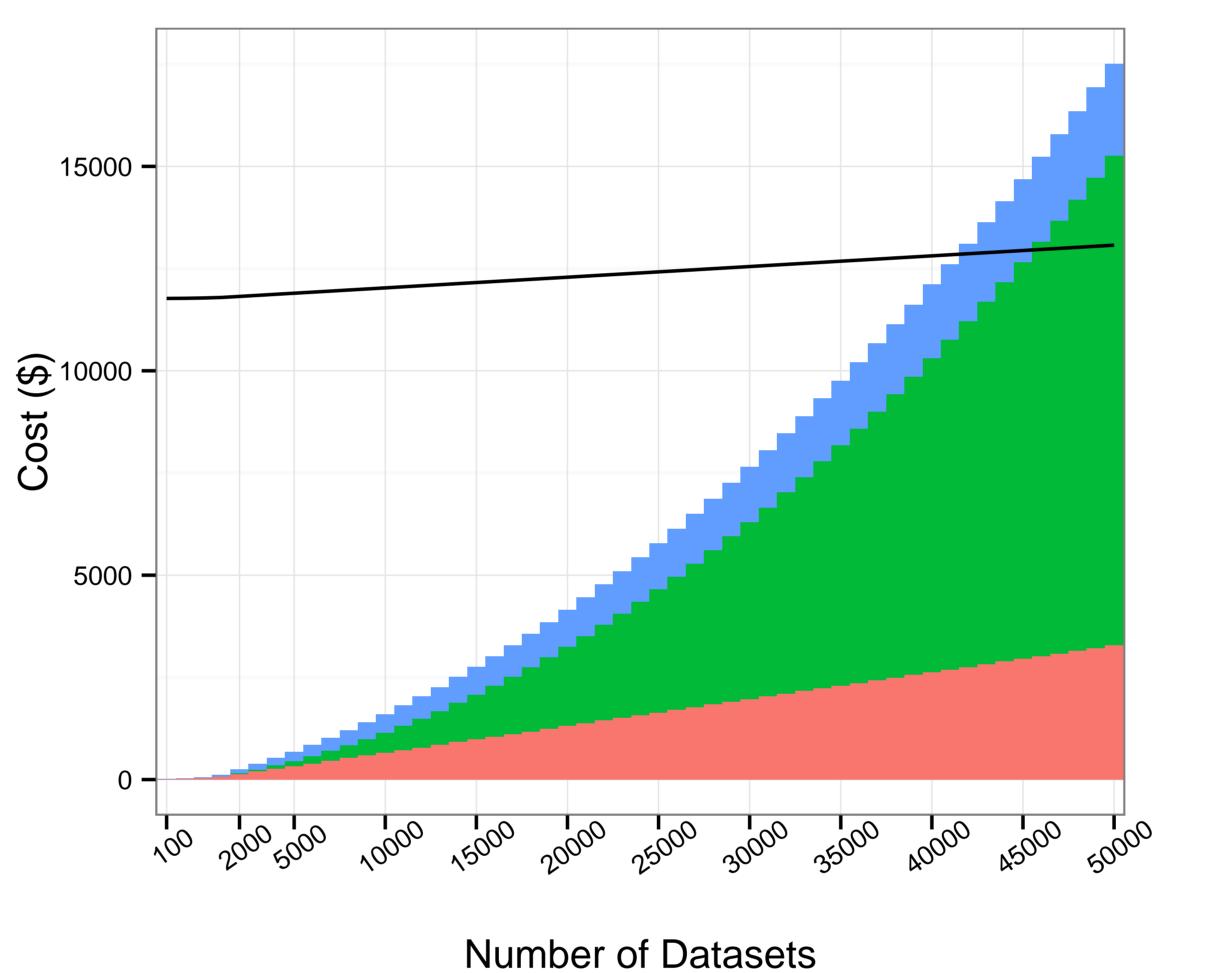
different availability zones for the past three months

## Results

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

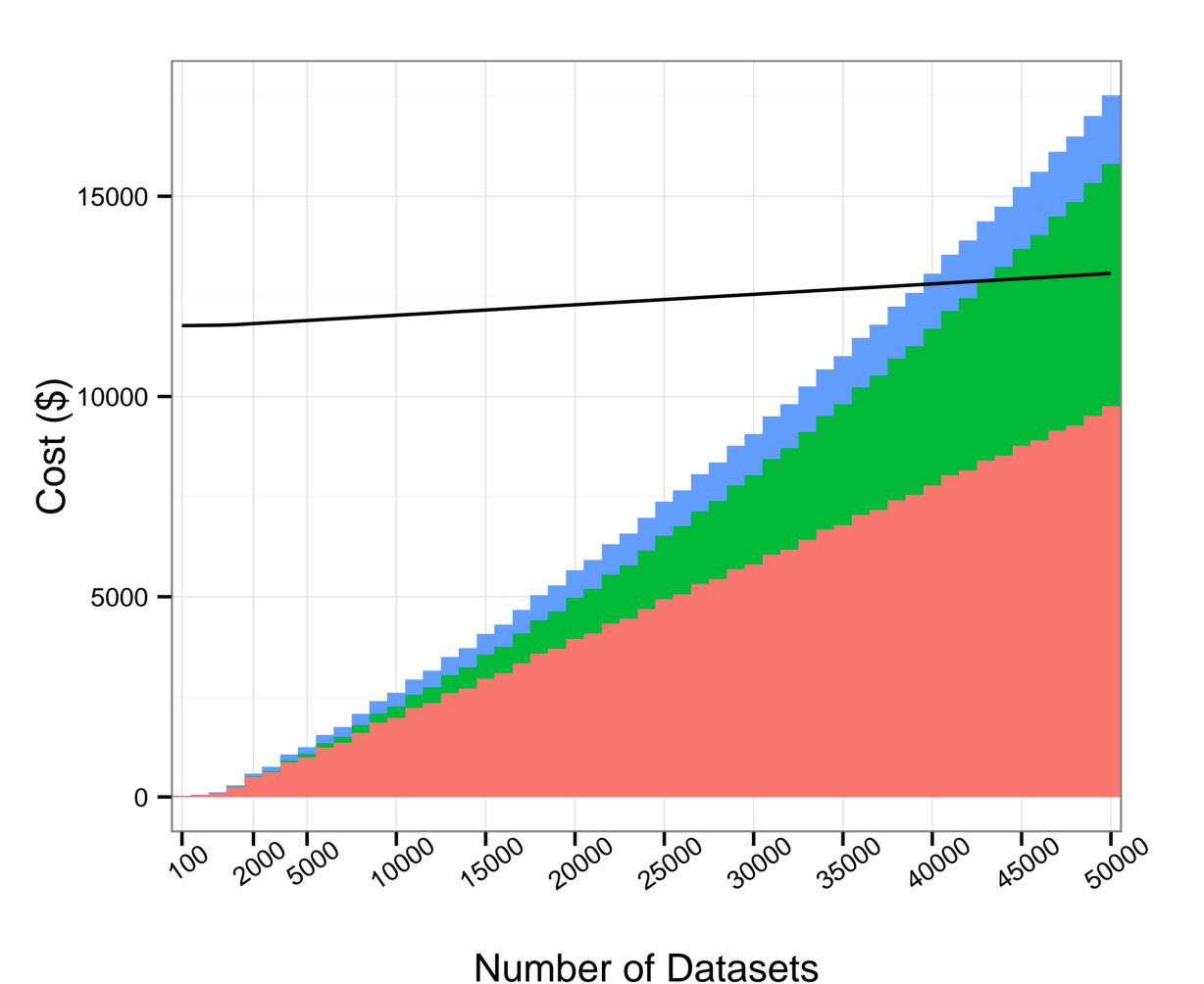
passed. # DS: Number of datasets. CPD: Cost Per Dataset. C-PAC: Configurable Pipeline for the Analysis of Connectomes. NITRC-CE: NITRC Computational Environment **Processing** Nodes PF CPU Time Wall Time Cost CPD # DS ANTS Cortical Thickness 147 \$760.24 \$0.24 C-PAC 20 3 22 \$80.54 \$0.07 Resting state fMRI processing w/ 4 strategies Quality Assessment Protocol 20 4 14 \$19.02 \$0.02 Freesurfer recon-all 23,664 193 \$211.44 \$0.21  $3 \quad $2.19 > $0.01$ NITRC-CE **FSL FIRST** 13 \$4.69 > \$0.01 Temporal QA Freesurfer recon-all and FSL FIRST 20 32 49 \$252.36 \$0.32 LONI Pipeline Benefits of the cloud

- an efficient processing of the NDAR datasets, both in time and cost ► Additionally, there was no need to install, configure, and maintain machines locally, saving
- on overhead ► Competition for computing resources was controlled as every cluster launch was done by
- the user doing the processing; no fellow researchers needed to login and use additional computing power ► Cost and time models were produced to demonstrate the benefits of using cloud
- computing for increasingly large datasets



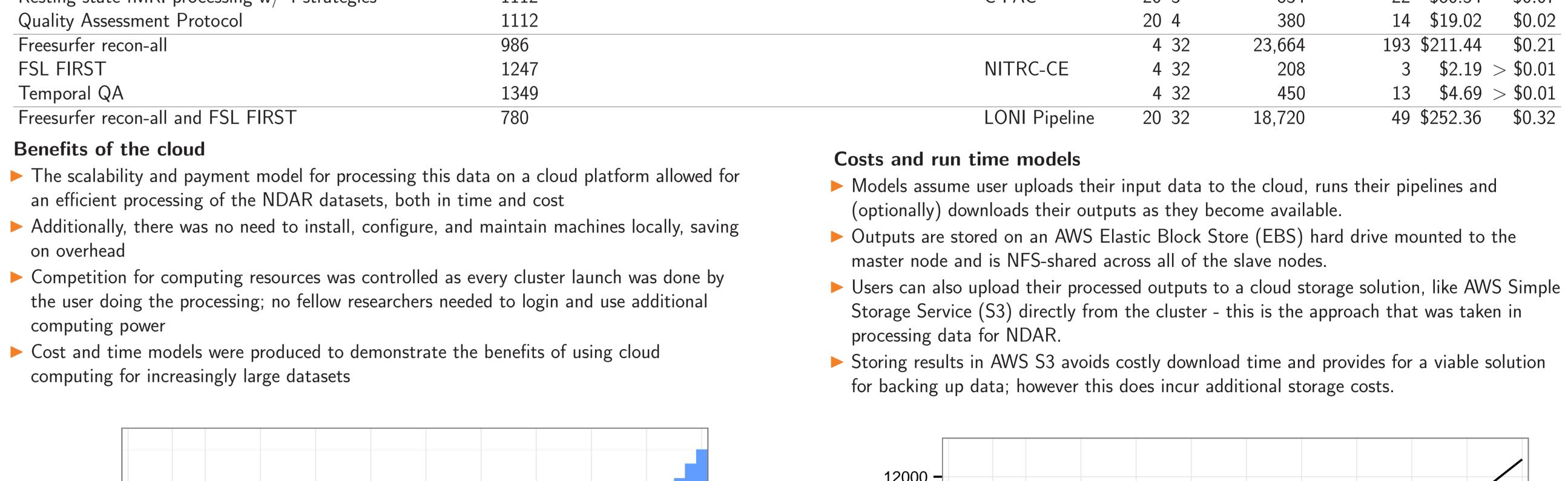
Instance Cost Storage Cost Transfer Cost

Figure 5: 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

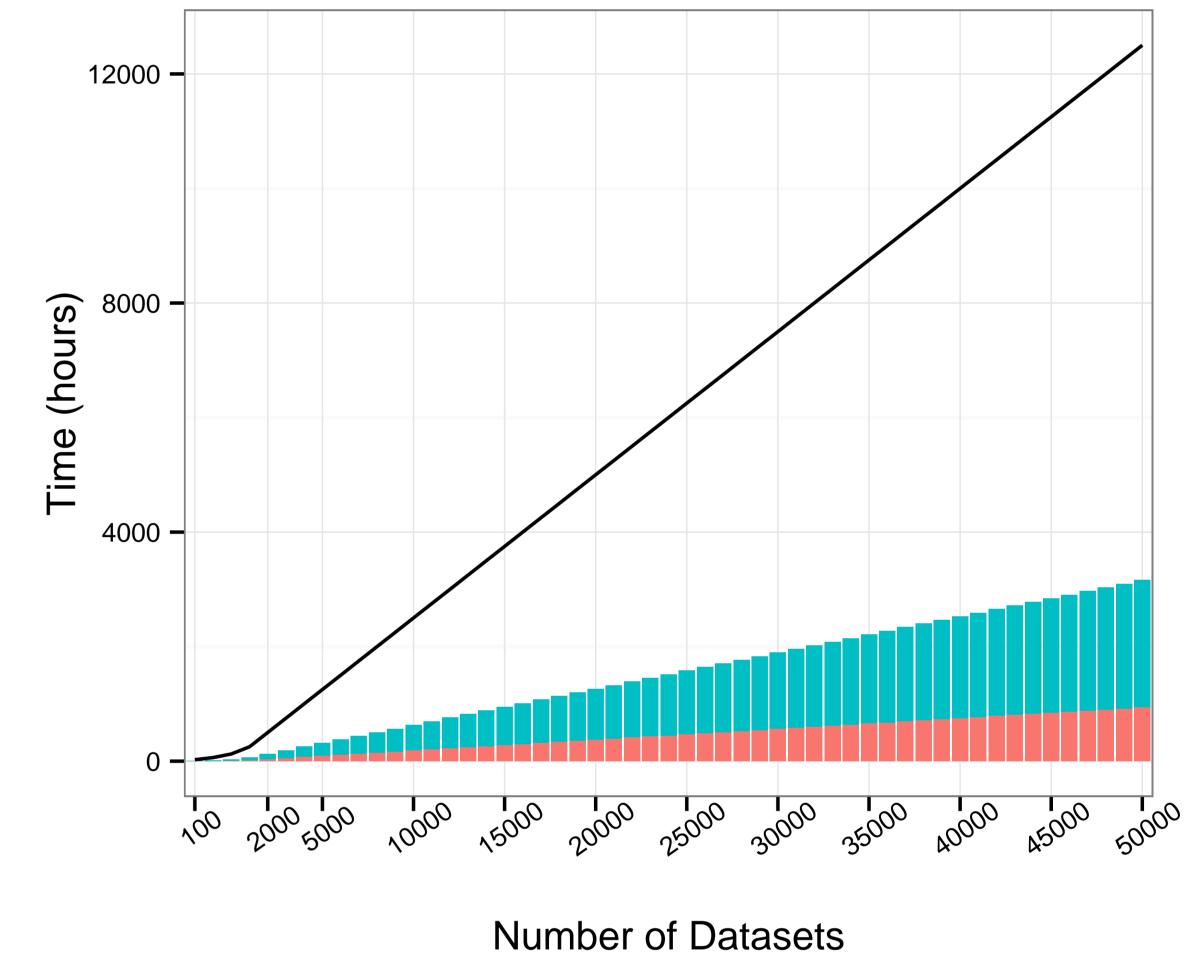


Instance Cost Storage Cost Transfer Cost

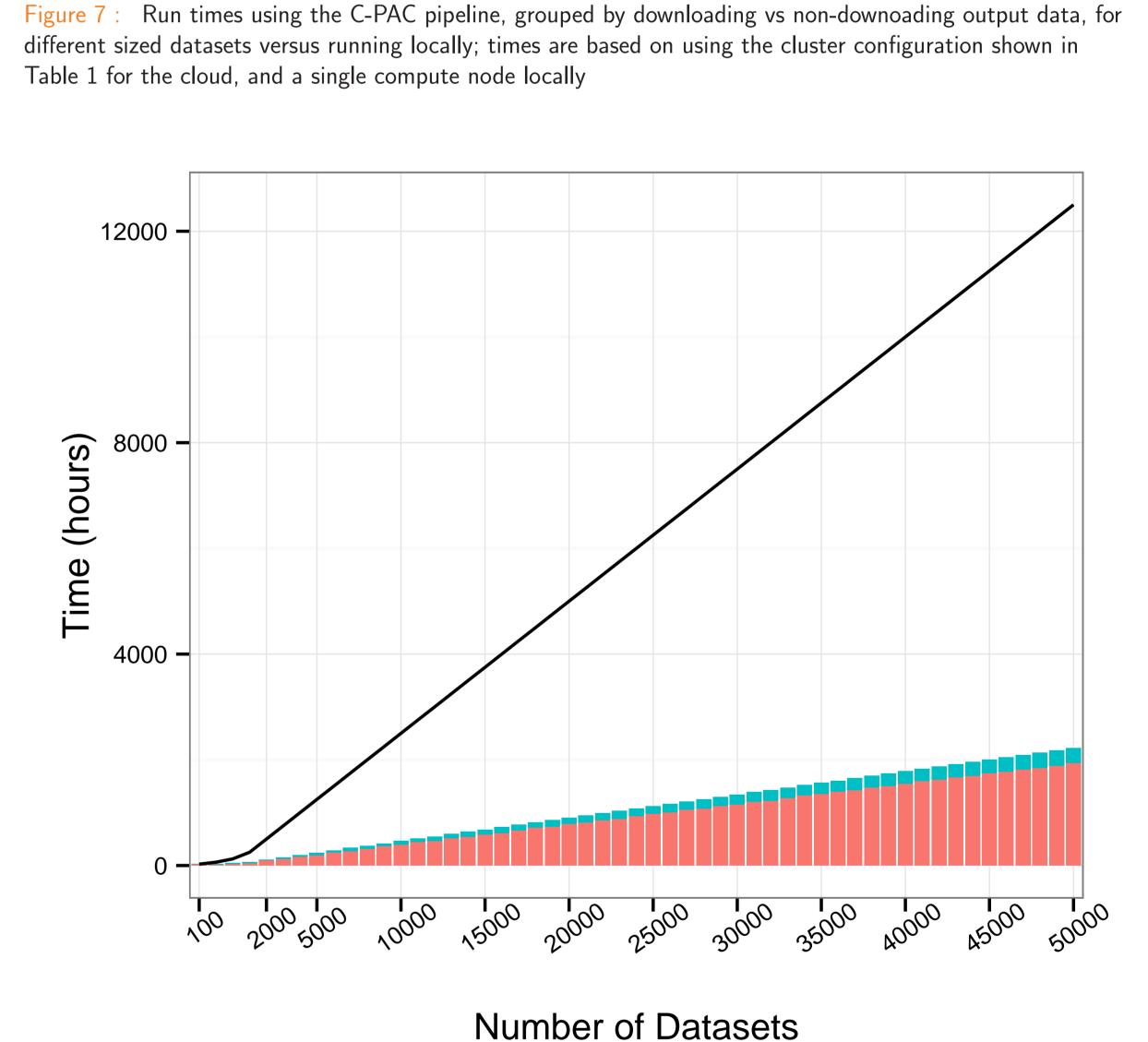
Figure 6: AWS EC2 costs, grouped by cost type, for a typical Freesurfer pipeline for different sized datasets versus owning and maintaining own server; costs are based on using the cluster configuration shown in Table 1 for the cloud, and a single compute node locally



► Storing results in AWS S3 avoids costly download time and provides for a viable solution



No Download Total Processing Time



No Download Total Processing Time

Figure 8: Run times using the Freesurfer pipeline, grouped by downloading vs non-downoading output data,

for different sized datasets versus running locally

## Conclusion

▶ Cloud computing services like AWS EC2 are viable solutions for the processing and analysis of large amounts of data

- . Hall, D. et al. (2012), Neuroinformatics. 10(4):331-9 2. Torgerson, C.M. et al. (2015), Brain Imaging and Behavior 9:89-103
- 4. Das, S.R. et al. (2009), Neuroimage, 45(3): 867-79
- 7. Starcluster ref
- Data collection and salary support was provided by a NARSAD Young Investigator Award to RCC.

- 3. C-PAC ref
- 5. QAP ref 6. Haselgrove, C. et al. (2014), Front. Neuroinform. 8:52