THE EFFECT OF COMPUTATIONAL ENVIRONMENTS ON BIG DATA PROCESSING PIPELINES IN NEUROIMAGING

Mohammad Ali Salari

A THESIS IN THE DEPARTMENT OF COMPUTER SCIENCE AND SOFTWARE ENGINEERING

Presented in Partial Fulfillment of the Requirements For the Degree of Doctor of Philosophy Concordia University Montréal, Québec, Canada

> April 2022 © Mohammad Ali Salari, 2022

CONCORDIA UNIVERSITY School of Graduate Studies

This is to certify that the thesis prepared

 By:
 Mohammad Ali Salari

 Entitled:
 The effect of Computational Environments on Big Data

 Processing Pipelines in Neuroimaging

and submitted in partial fulfillment of the requirements for the degree of

Doctor of Philosophy (Computer Science)

complies with the regulations of this University and meets the accepted standards with respect to originality and quality.

Signed by the final examining committee:

	Chair
Dr. Sebastien Le Beux	
	External Examiner
Dr. Etienne Roesch	
	Examiner
Dr. Yann-Gael Gueheneuc	
	Examiner
Dr. Gregory Butler	
	Examiner
Dr. Christophe Grova	
	Supervisor
Dr. Tristan Glatard	
Approved by	
Dr. Leila Kosseim, Graduate Program Director	
February 16, 2022	
Dr. Mourad Debbabi, Dean	

Gina Cody School of Engineering and Computer Science

Abstract

The effect of Computational Environments on Big Data Processing Pipelines in Neuroimaging

Mohammad Ali Salari, Ph.D. Concordia University, 2022

Variations in computational infrastructures, including operating systems, software versions, and hardware architectures, introduce variability in neuroimaging analyses that could affect the reproducibility of scientific conclusions. These variations are due to the creation, propagation, and amplification of numerical instabilities in analysis pipelines. It is critical to identify numerical instabilities to make experiments computationally reproducible. In this thesis, we characterize the numerical stability of commonly-used complex pipelines in the context of neuroimaging analysis across operating systems and provide accessible tools for developers and researchers to evaluate their pipelines and findings. First, we present the Spot tool that identifies the processes from which differences originate and the path along which they propagate in a pipeline. In the next step, to study the numerical instabilities more comprehensively, we introduce controlled numerical perturbations to the floating-point computations using the Monte-Carlo arithmetic method. For this purpose, we propose an interposition technique to model the effect of operating system updates on analysis pipelines using the Monte-Carlo arithmetic. Finally, leveraging the interposition technique, we compare numerical variability with tool variability in an fMRI analysis. We show that the results of analyses are sensitive to computational environment changes originating from numerical errors. All the methods implemented in this thesis are publicly available and can be used to facilitate further investigations toward stabilizing pipelines.

Acknowledgments

First and foremost, I would like to thank my research supervisor Tristan Glatard for all the invaluable support, understanding, guidance, and encouragement he has given me over the past years. Without his assistance and dedicated involvement in every step throughout the process, the success of this thesis would not be possible. I would also like to thank my research team members, Gregory Kiar and Yohan Chatelain, for their generous advice and ongoing contributions to my work. Special thanks to all the /bin lab members for a cherished time spent together in the lab and social settings, with a special mention to Valerie and Martin; you guys have been wonderful labmates. Getting through my dissertation required more than academic support, and I have many people to thank for listening to and, at times, having to tolerate me over the past years. My appreciation also goes out to my family and friends for their unconditional love, encouragement, and support throughout my studies.

Contribution of Authors

I was responsible for software development, data processing, and analysis of all findings, and drafting manuscripts. Tristan Glatard was responsible for supervising and supporting all of my contributions. The contributions of authors to each publication are described below.

C.I – File-based localization of numerical perturbations in data analysis pipelines [107]

I was responsible for the tool development, data processing, analysis, drafting the manuscript, and designing the figures. Lindsay B. Lewis and Alan C. Evans provided input on the dataset and pipelines, reviewed the results, and approved the final version of the manuscript. Gregory Kiar and Tristan Glatard supported development processes and data visualization. Tristan Glatard edited the manuscript, contributed to the interpretation of results, and supervised the findings of this work.

C.II – Accurate simulation of operating system updates in neuroimaging using Monte-Carlo arithmetic [106]

I was responsible for the software implementation, data processing, analysis, drafting the manuscript, and designing the figures. All authors contributed to the editing of the manuscript, experimental design and discussed the results. Yohan Chatelain helped with Monte-Carlo arithmetic simulations and software testing. Gregory Kiar and Tristan Glatard provided software development support. Tristan Glatard supervised the findings of this work.

C.III – Software variability in fMRI analysis: comparing between-tool and numerical errors

I was responsible for reproducing the experiments, data processing, drafting the manuscript, and designing the figures. Alexander Bowring supported the implementation of fMRI analyses and provided valuable feedback. Gregory Kiar, Yohan Chatelain, Alexander Bowring, Camille Maumet, and Tristan Glatard contributed to the experimental design, statistical analysis, and interpretation of results. Tristan Glatard edited the manuscript and supervised the findings of this work.

Contents

Li	st of	Figure	es	ix
Li	st of	Table	S	xii
1	Intr	oducti	ion	1
	1.1	Repro	ducibility Definitions	1
	1.2	Repro	ducibility Crisis	3
	1.3	Main	Causes of Irreproducibility	4
	1.4	Analy	zing Neuroimaging Data	5
	1.5	Thesis	Outline	6
2	Lite	rature	e Review	8
	2.1	Comp	utational Reproducibility	8
		2.1.1	Effect of Hardware Resources	9
		2.1.2	Effect of Parallelization	10
		2.1.3	Effect of Operating Systems	11
		2.1.4	Effect of Analysis Software	13
		2.1.5	Effect of Small Data Perturbations	14
	2.2	Techn	iques to Improve Reproducibility	15
		2.2.1	Code and Data Sharing	16
		2.2.2	Portability	17
		2.2.3	Numerical Instability	18
	2.3	Prove	nance Capture	21
		2.3.1	System-level Provenance Management Tools	21
		2.3.2	Provenance Formats	22
		2.3.3	Neuroimaging-specific Workflow Engines	23

3	File	e-based localization of numerical perturbations in data analysis pipelines \mathbb{C}	25			
	3.1	Introduction	27			
	3.2	Tool description	28			
		3.2.1 Recording provenance graphs	29			
		3.2.2 Capturing transient files	30			
		3.2.3 Labeling processes	31			
		3.2.4 Implementation	31			
	3.3	Experiments	32			
		3.3.1 HCP pipelines and dataset	32			
		3.3.2 Data processing	33			
	3.4	Results	33			
	3.5	Discussion	37			
		3.5.1 Key findings	37			
		3.5.2 Spot evaluation \ldots \ldots \ldots \ldots \ldots \ldots \ldots \ldots \ldots	40			
	3.6	Conclusion	41			
	3.7	Availability of Source Code and Requirements	42			
4	Acc	curate simulation of operating system undates in neuroimaging using				
Monte-Carlo arithmetic		onte-Carlo arithmetic	43			
	4 1	Introduction	45			
	4.2	Simulating OS updates with Monte-Carlo arithmetic	46			
	4.3	HCP Pipelines & Dataset 47				
	4.4	4 Results				
		4.4.1 Fuzzy librath accurately simulates the effect of OS updates	49			
		4.4.2 Fuzzy librath preserves between-subjects image similarity	50^{-0}			
		4.4.3 Results are stable across virtual precision	50			
	4.5	Conclusion & Discussion	52			
5	Cor	mparing software variability across and within fMRI analysis packages 5	54			
	5.1	Introduction	56			
	5.2	Materials and Methods	57			
		5.2.1 fMRI analysis & Dataset	57			
		5.2.2 Within-tool software variability simulation with Fuzzy Libmath	58			
		5.2.3 Data processing $\ldots \ldots \ldots$	60			
	5.3	Results	61			

	5.3.1	Validation of replication	61
	5.3.2 In the group analysis, BT was larger than WT		61
	5.3.3	In subject analyses, WT approached BT for some subjects $\ . \ . \ .$.	63
	5.3.4	Previous results were confirmed in thresholded group maps	63
	5.3.5	Brain masking instability was triggered by WT and BT	64
5.4	Discussion		65
S1	Reproduced results		69
S2	Maps of t-statistics for subject with highest WT variability		69
Б.			
Dise	cussior	1	71
6.1	The In	npact of Numerical Perturbations	71
6.2	The Importance of Numerical Instability		72
6.3	Recon	nmendations for Future Research	73
6.4	Conclu	usion	75

6

List of Figures

1	Surface maps of four metrics, standard-deviation and mean absolute differ-	
	ences, t-statistic and RFT significance values, indicate the inter-OS differences	
	for the cortical thickness extracted with CIVET over 146 subjects [47]. \ldots	12
2	Comparison of the thresholded statistic maps of two different analyses within	
	AFNI, FSL and SPM. Each row shows the results of each reanalyses, and	
	the last column shows the main figure from the original publication. Total	
	16 subjects and 21 subjects are participated in the first study (first row) and	
	second study (second row) respectively [15]	14
3	Provenance graphs created from the example pipeline in Listing 1. Processes	
	are represented with circles, files with rectangles, and read/write accesses with	
	plain edges. For convenience, the process tree is also shown, with gray dashed	
	edges. Processes forked by bet were captured by ReproZip while they did	
	not appear in Listing 1. Processed associated with executables located in	
	/usr/bin/ or /bin/ are not shown	29
4	Heatmap of non-reproducible processes across PreFreeSurfer pipeline steps.	
	Each cell represents the occurrence of a particular command line in a pipeline	
	step among Anatomical Average (AAve), Anterior/Posterior Commissure Align-	
	ment (ACPC-A), Brain Extraction (BExt), Bias Field Correction (BFC), or	
	Atlas-Registration (AR). Cell labels indicate the fraction of subjects for which	
	the corresponding process wasn't reproducible. For example, the flirt tool	
	was invoked 6 times in step DC for each of the 20 subjects: 2 instances weren't	
	reproducible in 19 subjects, 3 instances were always reproducible, and 1 in-	
	stance wasn't reproducible in 17 subjects. Grey cells indicate that the process	
	did not occur in the corresponding pipeline step	35

5	A complete provenance graph from the PreFreesurfer pipeline. Node labels	
	use the same abbreviations as in Figure 4. For better visualization, processes	
	associated with commands in /bin or /usr/bin were omitted, as well as	
	imtest, imcp, remove_ext, fslval, avscale, and fslhd	36
6	Differences between T2 fnirt results in PreFreeSurfer's Brain Extraction	
	(CentOS6 vs CentOS7). The colored squares indicate results obtained with	
	CentOS6 (in purple) and CentOS7 (in green). The red boxes highlight regions	
	with significant differences between the two OSes. An animated version of the	
	comparison is available here for better visualization	37
7	Sum of binarized differences between whole-brain FreeSurfer segmentations	
	obtained from PreFreeSurfer processings in CentOS6 vs CentOS7 (N=20).	
	Segmentations were resampled and overlaid to the MNI152 volume template.	
	Each voxel shows the number of subjects for which different results were ob-	
	served between CentOS 6 and CentOS 7. An animated comparison of segmen-	
	tations obtained for a particular subject is available here for better visualization.	38
8	Dice coefficients between regions segmented by FreeSurfer in CentOS6 vs Cen-	
	tOS7 (N=20), ordered by increasing median values. Each point represents the	
	Dice coefficient between segmentations of a particular region obtained in Cen-	
	tOS 6 vs CentOS 7 for a given subject. Boxes brightness is proportional to	
	the logarithm of the corresponding brain region size. \ldots \ldots \ldots \ldots	39
9	PreFreeSurfer pipeline steps.	48
10	Comparison of OS and FL effects on the precision of PreFreeSurfer results	
	for $n=20$ subjects. FL samples were obtained at the global nearest virtual	
	precision of t=37 bits. \ldots	49
11	RMSE-based hierarchical clustering of OS (left) and FL (right) samples. Col-	
	ors identify different subjects, showing that similarities between subjects are	
	preserved by the numerical perturbations. Horizontal gray lines represent	
	average RMSEs between (top line) and within (bottom line) subject clusters.	50
12	Comparison of RMSE values computed between OS and FL results for differ-	
	ent virtual precisions	51
13	${\bf A}$ and ${\bf B}:$ Bland-Altman plots comparing group-level differences computed	
	between tools (\mathbf{A}) and within tools at machine error (\mathbf{B})	62
14	Voxel-wise comparison of group-level differences in BT and WT	63

15	For subject with highest WT variability, unthresholded subject-level variabil-		
	ity computed between tools (A), and within tools at machine error (B)	64	
16	$\mathbf{A}, \mathbf{B}, \mathbf{C}$: Thresholded group-level t-statistic within tools at machine error for		
	FSL (A), SPM (B) and AFNI (C). Arrows point to activation clusters im-		
	pacted by both within- and between-tool variability. \mathbf{D} : Confusion matrices		
	of activation instability in BT and WT among the 360 regions of the HCP- $$		
	MMP1.0 parcellation.(with masking)	67	
17	${f A}$ and ${f B}$: Bland-Altman plots comparing group-level differences computed		
	over the intersections between tools (\mathbf{A}) and within tools at machine error (\mathbf{B}) .	68	
S1	Differences between reproduced and original results obtained in [15] of un-		
	thresholded group-level t-statistic for SPM (left) and AFNI (right). The		
	highest areas of difference in AFNI seem to be due to differences in brain		
	masks	69	
S2	For subject with highest WT variability, unthresholded subject-level t-statistic $\$		
	within tools at machine error for FSL (A), SPM (B), AFNI (C)	70	

List of Tables

1	Overview of definitions	3
2	Execution statistics of the pipelines per subject.	34
3	Types of provenance graphs in PreFreeSurfer.	34
4	Software processing steps (adapted from [15])	58
5	Voxel-wise mean and standard deviation of BT and WT variability in t-	
	statistic maps.	62
6	Voxel-wise mean and standard deviation of BT and WT variability in inter-	
	sected t-statistic maps	65

Chapter 1

Introduction

Reproducibility is regarded as a fundamental concept in the scientific community. Research findings are expected to be reproducible so that their authenticity and reliability can be ensured. The goal of our research is to investigate the reproducibility of computational analyses in general across different computing environments. In particular, we are mostly interested in neuroimaging as a case study. We present techniques to evaluate the numerical instability of analysis across different computing environments instead of masking the reproducibility problem by fixing parameters.

In this chapter, we summarize the main definitions and principles relevant to reproducibility. We describe the context of the current "reproducibility crisis" acknowledged in several scientific disciplines. Multiple studies have shown that some research findings could not be reproduced by independent researchers, or even by the original researchers themselves. We discuss the main causes for this lack of reproducibility, focusing on the computational aspects. Additionally, we describe different kind of analyses of neuroimaging data and their implemented software which are used throughout this thesis.

1.1 Reproducibility Definitions

There are different definitions for the terms reproducibility, repeatability, and replicability, which leads to confusion because the same words are used for different concepts. We present different terminologies found in the literature and summarized by Plesser [102] (see Table 1).

According to Peng's definition [99], reproducibility is defined as the ability to regenerate the same results as the original findings when the experiment is reanalyzed given exactly the same analytic methods and data. Reproducibility ensures that independent scientists can reproduce the same results using the same data and procedure as published in the original publication. Replicability is defined as the ability to obtain similar results as published in the original study when the experiment is reimplemented using independent data and analytic methods. Replicability confirms scientific claims and ensures that independent investigators can produce consistent results, using new data and methods. Peng introduced the idea of reproducibility spectrum based on his definition of reproducibility, which defines a minimum standard to evaluate the authenticity of scientific claims. In this spectrum, according to what data and sources are available, a full replication or no replication of a study can be achieved. The same definitions of reproducibility and replicability are also used by Schwab et al. [109].

Donoho et al. [33] defined reproducible computational research as a process by which "all details of computations such as code and data are made conveniently available to others". The authors associate reproducible research with open science, including open code and data. They observe that reproducibility can be achieved by publishing the experimental resources over the Internet, which facilitates versioning, testing, discovery, and access to the research materials.

In addition, Goodman et al. [48] renamed Peng's reproducibility and replicability as methods reproducibility and results reproducibility respectively, and adopted a new terminology called inferential reproducibility. From Goodman's terminology, exactly the same data and procedure are reanalyzed in methods reproducibility. Result reproducibility is equivalent to Peng's replicability terminology, which is defined as getting almost the same results compared to the original study from an independent replication of a study. Inferential reproducibility is defined as getting the same conclusions from either a reanalysis of the original study or an independent replication of a study with different data and procedures.

Furthermore, the Association for Computing Machinery (ACM) [3] proposes three different categories of repeatability, replicability, and reproducibility. Repeatability is defined as repeating computation on the same experimental setup including operator team, operating conditions, location, and measuring system. Similar to repeatability, replicability uses identical experimental conditions except performer team, which means that an independent group can achieve the same results through the same experimental parameters. Reproducibility is also defined as performing computation on different experimental setups via different teams independently. Reproducibility and replicability are used inversely compared to Peng's definitions.

Schwab et al. ₍₂₀₀₀₎	Donoho et al. $_{(2009)}$	Peng ₍₂₀₁₁₎
Reproducibility Replicability	Open code and data	Reproducibility spectrum
Revol et al. ₍₂₀₁₃₎	$ACM_{(2016)}$	Goodman et al. ₍₂₀₁₆₎
Numerical reproducibility	Repeatability Replicability Reproducibility	Method reproducibility Results reproducibility Inferential reproducibility

Table 1: Overview of definitions.

In addition, numerical reproducibility is defined as the ability to regenerate bit for bit identical results from multiple runs [104]. Two files will be considered numerically identical if they have identical binary contents. Binary comparison is calculated by comparing checksums. A computation might be reproducible based on Peng's definition, but not be numerically reproducible. For example, small numerical errors created during the pipeline execution may hamper numerical reproducibility, but be negligible in the final results.

Reproducibility, as the cornerstone of scientific research, guarantees the reliability, a level of accuracy accepted by the user, of results. A reproducible study provides a context in which one can get results consistent with the original work. In addition, it allows researchers to perform similar analyses more quickly by sharing resources rather than spend months figuring out current solutions. This enables others to use and modify existing works as a part of their experiments [102]. In our work, we follow Peng's definition of reproducibility unless we directly refer to numerical reproducibility. We seek to identify why reproducibility may not be ensured, focusing particularly on computational aspects.

1.2 Reproducibility Crisis

Recently, scientists began to realize that the results of many scientific experiments were neither replicable nor reproducible. This realization led to the so-called the reproducibility crisis. We provide an overview of evidence for the reproducibility crisis, which has raised important concerns in the scientific community.

Ioannidis [60] introduces an important framework to demonstrate the probability that research findings are false, and the propagation of valid findings in a given research field. He defined biased research as "the combination of various design, data, analysis, and presentation factors that tend to produce research findings when they should not be produced". Consequently, biased research, focused on an individual discovery rather than on broader evidences, decreases the chance of true findings. He concluded that "most of the research claims are less likely to be true than false for most fields and research designs". The author argued that the probability of true findings is highly dependent on the number of similar studies in a scientific field, the number of researchers/teams involved in the study, and the flexibility of analytic models, definitions, and outcomes. For example, the smaller the studies conducted in a scientific field, the less likely the research findings are to be true.

To highlight the importance of scientific reproducibility, a survey [8] collected data from 1,500 scientists among different disciplines mostly from biology, medicine, and engineering. This survey found that 70% of the scientists could not replicate another scientist's findings, and even 50% failed to reproduce their own results. This survey listed some of the main reasons that lead to irreproducibility of analysis such as poor statistics, the pressure to publish and selective analysis. With this, over 50% of the scientists believed that there was a significant crisis.

Furthermore, some studies underlined the reproducibility issues of current analysis methods in neuroimaging [68, 92, 35]. For example, to evaluate the reproducibility of a group of functional MRI (fMRI) analyses, a study [35] collected resting-state fMRI data from 499 healthy controls. Using this dataset, it reports that the most common software packages for fMRI analysis (SPM, FSL, AFNI) can result in a high degree of false positives, up to 70% compared with the expected 5%. These results question the validity of some 40,000 fMRI studies and may have a large impact on the interpretation of neuroimaging results.

The impact of Alzheimer disease and semantic dementia on Grey Matter volume changes in [80] show similar changes between volumes of specific structures compared to the discrepancies caused by computation environment variability in [52]. In addition, differences in cortical thickness caused by various operating systems, software versions and workstation types were roughly of the same order of magnitude than findings [78] from patients who suffered from schizophrenia. All these evidences show a significant crisis in reproducibility of experiments that should be taken into consideration in scientific communities.

1.3 Main Causes of Irreproducibility

There are a number of reasons that limit the reproducibility of most research, which can be simplified into two primary categories: human errors like an incomplete specification of processing detail and the computational causes of irreproducibility such as software/hardware changes.

The main barrier to reproducibility in many cases is that the source code, data, and analytic method description are no longer available. Addressing this problem requires the development of a culture of reproducibility in the scientific community to make sure that data can later be used appropriately, which would enable any third party to reproduce the same experiment [99, 111].

From the computational point of view, reasons such as the lack of details of the computational environments can contribute to the irreproducibility of research results. Analyses must provide sufficient information on code, software, hardware, and implementation details to be computationally reproducible. However, capturing such information is complicated, particularly in domains where results rely on a sequence of complex analyses such as neuroimaging pipelines. To overcome this complexity, a mechanism called provenance capture is designed to encompass all dependency information of the computational analysis such as input/output data, processing steps, and detail of computing environments.

Furthermore, the variety of computational infrastructures, including workstation types, parallelization methods, operating systems, and analysis packages, are known to influence reproducibility because of the creation of small numerical errors [52, 32, 47]. For example, we explain in the next chapter that different order of summation operation of floating-point numbers can lead to creation of small numerical differences. The propagation and amplification of these tiny differences by analysis pipelines may cause reproducibility issues. We will discuss in more detail the effect of computational environments on big data processing pipelines in Chapter 2.

1.4 Analyzing Neuroimaging Data

There are many different kinds of imaging techniques to acquire brain image data. The most common techniques are structural magnetic resonance imaging (sMRI), functional magnetic resonance imaging (fMRI) and diffusion magnetic resonance imaging (dMRI).

Structural neuroimaging deals with the anatomical structure of the brain and helps diagnose brain injury and certain diseases, such as tumor and stroke. The main software packages used for sMRI are CIVET [4], FreeSurfer [41], and FSL (FMRIB Software Library) [65]. Functional imaging is used to measure brain function based on specific tasks completed by subjects such as listening to sounds, reading, or small movements. Functional imaging identifies the areas of the brain that are involved in these tasks. The main software packages that implement fMRI processing are SPM (Statistical Parametric Mapping) [2], FSL (FMRIB Software Library) [65], and AFNI (Analysis of Functional NeuroImages) [23]. Diffusion imaging is another kind of MRI analysis that measures the anatomical connectivity between regions, and its main toolboxes are DIPY (Diffusion Imaging in Python) [42], MRtrix [115], and FSL.

Depending on the MRI analysis, several steps can be involved in a neuroimaging study. Generally, the analysis procedure can be divided into pre-processing and statistical steps. Pre-processing steps prepare data for the statistical analyses and are common between all MRI analyses, including brain extraction to separate the brain tissues from the other parts, or brain alignment which aligns a brain image with a reference image such as one produced by MNI (Montreal Neuroimaging Institute) [37]. After the pre-processing steps, depending on the modality of analyses (e.g., sMRI, fMRI, and dMRI), statistical analyses are applied to understand the nature of the data and obtain relevant results that can be used and interpreted by neuroscientists.

The various pre-processing and analysis steps involved in a neuroimaging experiment are often combined in workflows or pipelines. Pipelines are used to automate data analysis and accelerate the processing of complicated analyses.

1.5 Thesis Outline

The objective of this thesis is to understand the role that numerical instability plays in the reproducibility of results by focusing on the effect of operating system variability. For this purpose, we leverage system call interception techniques including the ReproZip tool [103], a tool that tracks the operating system calls. We also use perturbation models such as Monte-Carlo arithmetic (MCA) [97] as an extension of standard floating-point arithmetic that exploits randomness in basic floating-point operations to simulate the numerical errors. This thesis is manuscript-based, meaning that each of the three chapters is an exact copy of either a published or to be submitted manuscript.

C.I – File-based localization of numerical perturbations in data analysis pipelines [107] (Chapter 3)

C.II – Accurate simulation of operating system updates in neuroimaging using Monte-Carlo arithmetic [106] (Chapter 4)

C.III – Comparing software variability across and within fMRI analysis packages (Chapter 5) In Chapter 2, we review the background material related to this thesis in general. Chapter 3 introduce Spot, a tool to detect the source of numerical differences in complex pipelines executed on different operating systems. This chapter is completed and published in the GigaScience journal. Chapter 4 then study whether the MCA method is a good perturbation model for evaluating pipeline stability across operating systems. This chapter is published in the MICCAI workshop on Uncertainty for Safe Utilization of Machine Learning in Medical Imaging (UNSURE). Chapter 5 present a comparison of numerical and software variability through MCA. This chapter will be submitted to the Human Brain Mapping (HBM) journal by the end of Winter 2022. The thesis then provide a discussion and a conclusion in Chapter 6.

Chapter 2

Literature Review

In this chapter, we present previous works that investigated the effect of computational environments on scientific results: showing the magnitude of the effect of computing environment changes such as hardware and software implementations. Next, we review techniques and tools to enhance the reproducibility of experiments including code and data sharing methods using version control systems, and virtualization techniques to encapsulate computational variability of the analysis. Finally, we describe provenance management tools to collect and represent the analysis dependencies.

2.1 Computational Reproducibility

Many works have investigated the reproducibility of computational pipelines in the past few years. In general, analysis results are not reproducible because of the numerical errors caused by computing environment changes, including hardware configuration or operating system.

Changes in the computational conditions may introduce small numerical errors, subsequently propagated and amplified by pipelines. The analysis pipelines are said to be numerically unstable. Numerical instability is a characteristic of the pipelines which amplify small numerical errors and then hamper the reproducibility of the analyses depending on the complexity of the pipeline and magnitude of the errors. In many cases, numerical instability is an important issue for reproducibility.

The following sections discuss the effect of influential elements on reproducibility, in particular workstation type, parallelization techniques, operating system changes, analysis software variety, and perturbations applied in input data.

2.1.1 Effect of Hardware Resources

The hardware configuration of computers is an influential source of irreproducibility [59]. The numerical errors are particularly noticeable across computing processors such as CPUs (Central Processing Units), GPUs (Graphics Processing Units) and APUs (Accelerated Processing Units), mainly due to incompatibilities of floating-point units (FPU) with the IEEE-754 standard when arithmetic precision of the floating-point values are not specified uniformly.

Even using the same arithmetic precision, it is difficult to achieve bitwise identical results across different hardware resources. Recent studies show that hardware developments to improve computational performance sacrifice numerical reproducibility [34, 27]. For instance, code optimization techniques embedded inside CPUs, known as out-of-order execution (dynamic scheduling), impede reproducibility because processors might execute instructions out of the original order in which they appear based on the availability of input data and execution units to use resources efficiently [121]. Therefore, they might compute floating-point operations in different order, which often leads to different results due to different rounding of the intermediate floating-point arithmetics. Also, this has been shown in several papers [34, 27] that some operations in particular sum and division are not associative.

A study [67] implemented acoustic wave equation to see the effect of processor architecture on results. The authors illustrate irreproducible results across different processors including AMD CPU, NVIDIA GPU, and AMD APU, even using the IEEE-754 standard. The results numerically vary from one architecture to another, the maximal relative difference between results in the range [0.1, 1] and its mean value is 10^{-5} . Such differences often occur due to rounding errors generated by different orders in the sequence of arithmetic operations.

In neuroimaging, it is important to evaluate the consistency of results when they are executed on heterogeneous computing systems that use more than one kind of processor or core. A number of tests were conducted [52] to gain insight into the variability of results from neuroimaging packages based on different data processing conditions like different workstation types. Two different types of workstations were compared: an HP (Hewlett Packard) one using Centos 5.3 and 8 CPU cores, and a Mac one using OSX 10.5.8 and 2 CPU cores. This study showed significant absolute differences among the volumes of anatomical structures obtained on the two different workstations. These differences were on average $8.8\pm6.6\%$ (range 1.3-64.0%) (volume) and $2.8\pm1.3\%$ (1.1-7.7%) (cortical thickness).

2.1.2 Effect of Parallelization

Developers leverage parallelization techniques to accelerate the execution performance at different levels, from multi-threaded programming to high-performance computing (HPC). With parallelization techniques, contrary to sequential implementations, the execution order of the processes may change in different runs. Consequently, several runs of the parallelized code may produce different results, even on the same computer.

To show the existence of such issues, the impact of the number of processors on numerical reproducibility was studied [32]. This study simulated the process of deformation of metal sheets in the packaging industry to measure local change of the sheet thickness using different number of processors. Results obtained significant differences in maximal and minimal values of sheet thickness changes originating from the nondeterministic behavior of a program code linked to a component of the Intel Math Kernel Library. Their findings showed the amplification of rounding errors in summations after running the same simulation on the same computers with a different number of processors. This also proved that the summation operation is not associative because of different rounding of the intermediate floating-point results, even using the standard IEEE double-precision arithmetics. Therefore, final result of the summation depends on the order in which values are processed, which changed by the number of processors.

Another statistical simulation showed reproducibility failures in multi-core processing performed on GP-GPUs and multi-core CPUs [113]. Multi-core architectures enable multithreaded environments for running numerical intensive applications at high speeds. This study showed that the stability of molecular dynamics simulation results is not guaranteed in multi-core processors due to different orders of floating-point operations (e.g., division and square root operations) leading to different rounding and truncation.

In addition, parallel programming may lead to vulnerabilities like race conditions that further impede reproducibility. A race condition is a situation in concurrent programming where two concurrent threads or processes have access to the same resources and attempt to change it at the same time. When one thread is performing read on a particular data element, another thread is allowed to modify or delete this element. The resulting final state depends on the order of the operations, which is not correctly programmed by the application developers. In addition to race condition, some other problems have been listed as the main sources of numerical differences in many parallelized experiments such as outof-order execution, and message buffering non-blocking communication operations [104].

Message buffering is a type of communication using send/receive functions in parallel

programming, which can be blocking and non-blocking. Non-blocking communication means that computing and transferring data can happen at the same time for a single process. This allows communication to overlap, which generally can speed up the process but also can lead to different computing orders and irreproducible results for different runs.

Furthermore, some experiments [52] determined the effect of parallelization on neuroimaging pipelines, most precisely in different versions of FreeSurfer. They showed that concurrent running would not make statistically significant differences based on the comparison of voxel volume of specific brain structures for the same conditions. This is an example where Peng's reproducibility is achieved while numerical reproducibility is not. However, further experiments are needed to investigate the effect of parallelization on neuroimaging.

2.1.3 Effect of Operating Systems

We summarize results [47] that quantified the reproducibility of computational analyses across operating systems. In particular, the authors determined the reproducibility of three neuroimaging workflow packages, FSL, FreeSurfer, and CIVET between CentOS 5.10 and Fedora 20.

Using FSL, cortical and subcortical tissue classifications resulted in Dice values as low as 0.59 between the classified tissues on CentOS and Fedora operating systems. These differences mainly correspond to the mathematical functions implemented in different operating system libraries.

The results of <u>RS-fMRI</u> analysis revealed significant inter-OS differences in the second experiment, which showed that each pre-processing step could introduce small numerical variations and their accumulation creates important differences. These numerical differences are caused by changes in the implementation of mathematical functions like $\underline{sinf()}$ between operating systems.

Using FreeSurfer and CIVET, cortical thickness extractions had important differences in some specific brain regions across operating systems. Figure 1 shows localized regions of these differences for CIVET, which are quantified by mean absolute difference, standard deviation of absolute difference, t-statistic and random field theory (RFT). Areas in shades of blue on the RFT map are significant at the cluster level.

Additionally, inter-build differences are measured in this study. A static build of a pipeline refers to its compiled version where libraries are statically linked. The authors used the static builds of FreeSurfer CentOS 4 and CentOS 6 to measure their reproducibility. Results of FreeSurfer show that building static programs improves reproducibility across OSes, but small differences still remain (the t-statistic values of 2). The main cause of such differences is dynamic libraries that are loaded by the static executable at run-time.



Figure 1: Surface maps of four metrics, standard-deviation and mean absolute differences, t-statistic and RFT significance values, indicate the inter-OS differences for the cortical thickness extracted with CIVET over 146 subjects [47].

The work [47] detected that most of neuroimaging pipelines are sensitive to operating systems. The effect size of the variations is changed based on the complexity of the analysis pipeline. For instance, shorter analyses like brain extraction have much less significant disagreement compared to longer ones like subcortical tissue classification and RSfMRI analysis because of the accumulation of numerical errors in the complex analyses.

Furthermore, in future works, the authors expect similar reproducibility issues for the other Linux distributions including Debian and Ubuntu as long as they are based on glibc, the GNU C library, which includes mathematical libraries. Other studies [52, 77] reported similar issues for non-Linux operating systems.

2.1.4 Effect of Analysis Software

Reproducibility of computations also depends on the executed analysis software, even using the same operating system and hardware resources. Different version of an analysis software used in a computation may produce different results. Also, the re-implementation of the same experiment through different software packages can introduce discrepancies between their results. We summarize the impact of software variability including different software versions and a wider range of software packages on reproducibility of results.

Effect of software versions

In addition to comparing hardware and operating system variability, [52] studied the impact of using different pipeline versions. Significant volume differences are quantified across the FreeSurfer versions for both anatomical brain structures and cortical thickness measures.

The same study [52] showed that the effect sizes of different operating systems or software versions are close to the ones measured in neuropsychiatric diseases. For example, the impact of Alzheimer disease and semantic dementia on grey matter volume changes show similar changes between volumes of specific structures compared to the discrepancies caused by computation environment variability [52]. In addition, differences in cortical thickness caused by various operating systems, software versions and workstation types were roughly of the same order of magnitude than findings [78] from patients who suffered from schizophrenia. There are many other proofs in different domains that show the influence of software updates on results [110, 119].

Effect of software packages

In all aforementioned analyses, the choice of the software package remained fixed for carrying out the analyses in each study. To understand out the impact of analysis software variations on task fMRI results, several tests were conducted [15]. The authors investigated differences produced across three of the most popular neuroimaging software packages, AFNI, FSL, and SPM. They replicated specific analyses, a number of image processing steps, as closely matched to the original study as possible.

The statistical comparisons show a substantial disagreement between software package results, producing different location of activation regions. Figure 2 shows the substantial variation between each main activation area found in the original study and the reanalyses. Results indicate that the precise location of the significantly activated regions is highly dependent on the choice of software package and inference method.



Figure 2: Comparison of the thresholded statistic maps of two different analyses within AFNI, FSL and SPM. Each row shows the results of each reanalyses, and the last column shows the main figure from the original publication. Total 16 subjects and 21 subjects are participated in the first study (first row) and second study (second row) respectively [15].

Analyses [15] found that the size of datasets can contribute to the variation of results. For instance, results obtained from analyses that use smaller sample size are less likely to be reproducible than analyses in which more subjects participated. This is likely explained by the fact that group analyses benefit from regularization of numerical noise, which is expected to increase with sample size. Therefore, variation in the outcome of an fMRI analysis depends not only on the choice of software package used, but also on the dataset being analyzed.

2.1.5 Effect of Small Data Perturbations

Neuroimaging pipelines are sensitive to changes in the computing environment. Studies were conducted to show the instability of some specific steps of MRI analysis through the simulation of minor perturbations in input data. For instance, reproducibility of the cortical surface reconstruction analysis in the presence of small perturbations is measured in [82]. The authors investigated results of two pipelines, CIVET and FreeSurfer, after applying 1% intensity modification on one voxel located in a non-cortical region. Contrary to expectations, widespread surface changes were observed across the cortex.

Similarly, another study [45] observed substantial variability of motion correction algorithms in fMRI analyses by applying one-voxel perturbation. Results demonstrate significant differences for Niak and FSL. These variations may result in wrong activation maps and increase the prevalence of false activations on the subsequent steps of fMRI processing. Recently, the processing of high-resolution images has been made possible through a new version of pipelines. [83] quantifies the variability of analysis results across different image resolutions. The authors investigated the partial volume effects¹ of various image resolutions on the automated cortical surface extraction through CIVET and FreeSurfer pipelines. They shows significant variability in results for the same analysis using images with different resolutions. For both pipelines, mean absolute error, signed error, and standard deviation are mostly reduced as a function of increasing resolution. Also, comparison of projected distance error maps between histological ground truth surfaces and MRI-derived surfaces confirms that the accuracy of analysis is increased in higher resolutions. Further research is needed to minimize partial volume effects along with magnifying the resolution to get more accurate results.

2.2 Techniques to Improve Reproducibility

Reproducibility is mainly ensured through three properties: source sharing for both code and data, research portability, and pipeline stability. Source sharing and research portability can be related to the FAIR principles [123] according to which scientific sources have to be findable, accessible, interoperable, and re-usable.

The first step in reproducibility is finding and accessing research products related to Findable and Accessible principles in FAIR. Code and data must be publicly available in a machine-readable structure. It enables the verification of scientific results by independent investigators.

Analysis pipelines must integrate with other execution environments. This is termed research portability and can be achieved using virtual machines and containerization technologies. Portability enables researchers to re-run analyses in a variety of execution conditions. This can be matched with the Interoperability and Re-usability of FAIR principles.

Analysis pipelines must be numerically stable across computing environments to be reproducible. Although many solutions currently exist to address analysis sharing and portability, the effect of numerical instability remains largely unexplained. We discuss a number of techniques and tools used to enhance sharing, portability, and stability of the analyses.

 $^{^{1}}$ Defined as the loss of contrast between two adjacent tissues in an image because of the limited resolution of the imaging system

2.2.1 Code and Data Sharing

To successfully reproduce a computational experiment, analysis sources must be accessible in a machine-readable structure [111, 56]. The importance of a proper structure is clear, specifically when we aim to share code with others or contribute to a wider group.

One foundation of code sharing is modern software engineering, which includes practices like version control systems (VCS). Version control ensures that the history of the code is available and archived. Git [114], as one of the most popular VCS frameworks, provides a distributed means to manage project files. Git facilitates the collaboration of developers on the same project using GitHub. GitHub is a Web-based service for Git, which hosts Git repositories.

Developers may share programs instead of source code because of commercial reasons, simplifying its usage, reproducibility improvement, etc. Several sharing tools exist to maintain a set of packages. PyPI (Python Package Index) is a software repository for the Python programming language. PyPI helps to share python packages and allows users to search for packages by keywords. There are more specific sharing tools for neuroimaging programs including Boutiques [46], a system to publish and integrate command-line applications automatically using a JSON descriptor, or NITRC-CE [71], which provides a number of pre-installed neuroimaging tools such as AFNI, FSL, and FreeSurfer into a standardized computational environment.

However, there are some challenges associated with data sharing including concerns about the privacy of personal information, lack of incentive for researchers, and technical issues associated with sharing of large datasets.

Git is a very efficient tool for managing textual information such as code, text, and configuration, but it is inefficient for storing large data. Therefore, extensions of Git, like git-annex [57] and Git-LFS (Large File Storage) [6] were developed to address the problem of sharing and versioning large data collections. git-annex uses Git to store and index files without committing large files into the Git repository. Similarly, Git-LFS reduces the impact of file size in the repository by replacing large files with lightweight pointer files, which refer to the actual file location.

Both Git and git-annex allow collaboration on a single repository, but sharing code and data between multiple projects can be an issue. Also, they lack advanced meta-data search capabilities. For instance, they cannot crawl through domain-specific repositories. Datalad [54] is an efficient tool for data sharing and versioning multiple datasets. This tool is built on top of git-annex and provides a unified access to data regardless of its origin. It guarantees that the content of a same version of a file would be the same across all clones of a dataset, regardless of where the content was obtained. Datalad supports multiple redundant data providers for each file in a dataset and transparently attempt to obtain data from an alternative location if a particular data provider is not available. Furthermore, a provenance record is provided by Datalad with all necessary information about input data to help reproduce the analysis results.

Higher-level platforms were designed to help scientists share neuroimaging data and make them public on the Web, such as openNeuro [49], LORIS [26], and XNAT [86]. These tools usually have a Web-based user interface and can integrate with processing platforms. Most of these tools use the Brain Imaging Data Structure (BIDS) [51] to describe and organize neuroimaging data. BIDS provides a standard for organizing and representing MRI data that reduces the effort of data sharing.

There exists a number of projects that promote open data-sharing initiatives in the field of neuroimaging including the International Neuroimaging Data-Sharing Initiative (INDI) [88, 90], the Alzheimer's Disease Neuroimaging Initiative (ADNI) [62], and the Human Connectome Project (HCP) [117]. Researchers who once struggled to access restricted datasets can now explore thousands of subjects using data published by these projects. Public access to this amount of brain imaging data is invaluable for specialists to test a variety of scientific hypotheses and evaluate novel image processing algorithms.

2.2.2 Portability

Although source code and data used in original experiments are available, re-executing a computational analysis is still not straightforward. To reproduce computational analyses, information about the computing environment is needed, in particular the operating system configuration, the hardware system architecture, and the specific versions of tools. Therefore, virtual machines and containerization techniques are suggested to ensure that specific computing parameters are completely preserved.

Virtual machines (VMs) encapsulate the entire context of computations, which provides an exact replica of the computational environment where analyses took place. The most popular implementations of VMs include VMware [122], VirtualBox [120], and KVM [76]. VMs may produce large images because they hold a copy of all the operating system files including the kernel, system libraries, and system configuration files. VMs bring an extra performance overhead such as I/O, CPU, and memory. Containerization tools like Docker [13] and Singularity [79] reduce the performance overhead and image size of VMs by sharing the kernel of the host system across the containers. These tools have emerged to build lightweight and portable images. Container images can be version controlled as the analysis code so that the exact same computing environment can be used to re-execute an analysis. However, similar to VMs, users are faced with the burden of ensuring that all necessary dependencies are collected inside the containers. Some workflow management systems exist to record the computational dependencies: we describe them in Section 2.3.

As an example of a container-based tool, Nextflow [29] is implemented to ensure workflow reproducibility. Nextflow uses Docker to containerize pipeline dependencies including data, code, and the computing environment. Nextflow can be integrated with public repositories in GitHub and cloud computing infrastructures to provide a rapid computation and effective scaling.

Containers are excellent lightweight technologies to solve portability issues. However, they are not perfect solutions to address the reproducibility of analyses across computing environments because they mask differences instead of fixing them so that differences can emerge in another source of variabilities, as explained hereafter.

2.2.3 Numerical Instability

Containers and VMs are good to mask the effect of variations in the hardware, parallelization, and operating system. However, this effect is due to numerical instabilities in the data analysis pipelines. These effects are the combined results of 1) the creation of numerical errors between conditions and 2) the amplification of these numerical errors throughout the pipelines.

Creation of Numerical Errors

The main causes of numerically irreproducibility stem from floating-point operations, in particular, using a finite precision in their arithmetic operations like summation [59, 113]. In this section, we summarise several solutions proposed to improve numerical reproducibility related to the floating-point operation and discuss their limitations.

Floating-point numbers are composed of a mantissa (significand) as the significant digits of the number, a base and an exponent that specifies a finite precision representation, and a sign. Floating-point numbers are an approximation of real numbers on computers [59]. Each computing system provides standardized math libraries necessary the floating-point computations with a finite precision. Finite precision computations create numerical errors mostly due to truncation and round-off errors.

Computers represent floating-point numbers with limited precision; they must round numerical results to the closest number that they can represent [38]. Truncation errors are due to the difference between actual (analytical) and truncated (approximated) values of computation [75]. When truncating a number to a limited number of decimal places, say x, the first x digits of the mantissa are reserved, simply chopping off the remainder. When rounding a number, the computer chooses the closest number that it can represent. Although rounding error is in the order of magnitude of $e > 10^{-7}$ for IEEE-754 single-precision and $e > 10^{-16}$ for IEEE-754 double-precision, their accumulation can be significant.

Due to the non-associativity of floating-point addition, rounding errors can lead to different results depending on the order in which operations are performed. For instance, assuming a computer with 4 decimal digits of precision, the following summation in different orders leads to different results.

 $(4.127 \oplus 100.2) \oplus -104.2 = 104.3 \oplus -104.2 = 0.100$

$$4.127 \oplus (100.2 \oplus -104.2) = 4.127 \oplus -4.000 = 0.127$$

In the first summation, a rounding error is introduced in the truncation of 104.327 to 104.3.

In the following, we explain several approaches [93] to solve these numerical errors, such as using higher precision, deterministic order of operations, arbitrary-precision operations, and fixed-point arithmetic.

Higher precision. Using high-precision numbers, for instance, calculations using doubleprecision produce more accurate results than single-precision. However, rounding errors still exist for higher precision.

Deterministic order of operations. It is possible to make computations deterministic in the order in which floating-point operations are performed. This can lead to more numerically stable results across runs, but this solution add memory overhead and affect the performance of executions[9].

Fixed-point arithmetic. Fixed-point numbers reserve a fixed number of digits after the decimal point, assuming that numbers are integers multiple of some common denominators

with similar orders of magnitude. It is common to use fixed-point arithmetic to represent large fractional numbers. Fixed-point operations are often faster than floating-point ones because they do not depend on the availability of an FPU. Although they helps reduce rounding errors, they limits the range of values, and overflows can occur if the result of an operation is larger or smaller than the numbers in that range.

Arbitrary-precision operations. We can use high-precision or even arbitrary-precision operations to push the precision limitation of floating-point arithmetic. The precision of numbers is limited only by the memory of the host system. This requires many hardware instructions for each arithmetic operation and the difficulty of handling the variable-width storage.

Amplification of Numerical Errors

There is evidence showing that analyses are not stable to small numerical errors because of the propagation and amplification of these errors. For instance, propagation of rounding errors from the initial value in numerical computations were studied by performing different experiments in [38]. In this paper [38], several computational experiments are presented to demonstrate the rapid growth of rounding errors in iterative computations like iterative addition. The accumulation of rounding error from summation operation indicates that analyses may produce different results. Due to similar reasons, the propagation of rounding error when simulating the metal sheet thickness changes in [32] turned to different results.

Another study [45] evaluates the stability of different neuroimaging pipelines in presence of one voxel perturbations. This study showed that iterative initialization schemes in motion correction algorithms lead to the propagation and amplification of numerical errors along the time series.

To address the numerical instability of pipelines, we can use the bootstrap technique. In [45], the authors explained that bootstrapping is an efficient technique to improve the robustness of motion estimation. The bootstrap version of the pipelines computes the median transformation results from the 30 samples from the medians of the parameters of the 30 transformations. It is, however, a compute-intensive technique that should be used only when no other solution to the instability is available.

In addition to bootstrapping, the bagging technique can reduce the effect of perturbations [16, 17]. Bagging, also called bootstrap aggregating, is a simple and powerful ensemble method. It helps reduce both bias and variance in the results, but it adds computational overheads. So, we can possibly stabilize pipelines and improve their accuracy using aggregates of results obtained with data perturbations.

2.3 Provenance Capture

Portability requires comprehensive information about the computational analysis in a machineexecutable form. This information can be obtained by provenance capturing tools.

Provenance is defined as the collected information about objects and processes involved in workflow results. This information can be used to verify the reliability and reproducibility of executions [91]. Provenance information can contain metadata that displays what data processing is undergone [94], for example, which parameters are used for the analysis, what form of image is used, how the image was registered/aligned to a standard space, how noise was eliminated, how a specific feature has been recognized. In this section, we discuss different aspects of provenance capturing such as system-level provenance capturing and workflow specifications. Finally, we give examples of some specific workflow engines that provide these features.

2.3.1 System-level Provenance Management Tools

Automated provenance capturing of computational analyses that contain a complicated sequence of dependencies is challenging. Packaging tools can automate the configuration capturing of an experiment by tracing the executed process using system call interceptions, such as ReproZip [22], CDE (Code, Data, and Experiment) [53], and CARE (Comprehensive Archiver for Reproducible Execution) [63]. These tools support reproducibility of research projects in a system-level provenance capturing.

ReproZip provides a lightweight solution that simplifies the process of making experiments reproducible. ReproZip creates a self-contained package for experiments by tracking processes and identifying all system dependencies automatically.

ReproZip packs all the necessary information of the experiment in a single package including input/output data, executable programs and steps, and computing environments. Using this provenance information, readers/reviewers can then extract the packages and reproduce the analysis. In addition, ReproZip generates a workflow specification for experiments that models the processes involved in the workflow. Users can explore the reproducibility of experiments or test other configurations.

ReproZip suffers from limitations as it cannot deal with packing experiments in other

operating systems than Linux-based OS. Also, packages may not be re-executed if they use absolute path hard-coded in the underlying experiment missing in the target environment.

In addition, ReproZip is unable to capture values processed in-memory and not written to disk, and temporary files that are removed during the execution. Therefore, full replication of the experiments may be impossible because these files and variables are not available in the provenance template. Furthermore, ReproZip cannot identify the execution order of files that are written by multiple processes concurrently. So, there is no guarantee to reproduce analysis in this condition as well.

Similar to ReproZip, CARE is a packaging tool that enables users to reproduce Linuxbased experiments by making a compressed archive of all the software dependencies. CARE is a portable tool that does not need any installation process, neither administrative privileges [63]. With the same purpose, CDE relies on system call interception to capture and make an independent package of computing environments [53]. In contrast to CDE, which is able to capture dependencies of simple analyses, CARE is more practical for complex analyses because of tracking the history of processes.

2.3.2 Provenance Formats

All aforementioned provenance capturing tools tracking, bundling, and sharing all the necessary dependencies of a project automatically and systematically. A few works have been conducted recently to introduce an integrated and standard provenance specification.

It is important to define a standard data model for representing and exchanging provenance information produced by workflow engines. Therefore, the World Wide Web Consortium (W3C) designed the PROV data model based on the history of three captured elements: entities, activities, and agents. The PROV model contains a set of documents to define various aspects of provenance information in heterogeneous environments such as the Web. For instance, PROV can make a relational model of provenance elements as an XML format. Also, PROV is not tailored to any specific application domains [21, 91].

Using PROV, we can check the reproducibility of scientific workflows by comparing results in different conditions. Also, this specification can provide information about the processes that lead to execution failures [91]. Similarly, a number of projects specific to neuroimaging proposed to support reproducibility. We discuss two popular neuroimaging provenance specifications: NIDM-Results and BIDS-Derivatives.

NIDM-Results, as a part of the Neuroimaging Data Model (NIDM) project [1], is a domain-specific extension of PROV based on semantic Web technologies. NIDM-Results provides a machine-readable representation of neuroimaging results. This specification encode the provenance results of some specific neuroimaging software such as SPM and FSL. It is also suitable for different neuroimaging modalities including functional MRI, structural MRI, and diffusion MRI.

NIDM-Results uses the same elements in PROV to provide an interpretable data provenance across heterogeneous neuroimaging workflows.

BIDS-Derivatives provides a standard data provenance compatible with BIDS [51] raw data format. BIDS-Derivatives simplifies both provenance capturing and representing. The specification is created as a JSON file based on a simple file format and folder structure. Researchers can easily share derived data, statistical models, and computational results automatically.

2.3.3 Neuroimaging-specific Workflow Engines

Capturing and documenting provenance information in neuroimaging pipelines is a challenging issue for reproducibility. Therefore, workflow engines were developed to address these issues using the specification models and capturing techniques introduced in the previous sections. These engines facilitate workflow composition and document them in a machinereadable form, which significantly enhance reproducibility. Some of the existing workflow engines in neuroimaging are explained in this section.

Nipype [50] is a Python package that introduces a framework to 1) make uniform access to neuroimaging analysis software and usage information, which allows mixing components from other packages developed in different programming languages through interfaces provided by Nipype; 2) simplify the design of workflows and facilitate the interaction between workflow modules; 3) reduce the training time of how use the packages. Nipype represents the provenance information using the W3C-Prov specification. VisTrails [18] and Taverna [95] perform similar methodology but are not specific to neuroimaging, and they are different in the way of data representation and the type of information they capture.

LONI [105] is a provenance framework for documenting data flows in computational environments without user intervention [85]. The relationship between processes is captured in an XML file format and re-executed later similarly. LONI is a Java-based program that facilitates the process of provenance capturing using a graphical user interface. The XML extension provided by LONI can be interpreted across different environments. Additionally, LONI supports parallel executions and provides a simple mechanism for researchers, particularly in the neuroimaging field, to disseminate their experiments. ReproNim [70] is an integration of tools to ensure reproducibility at different stages of the analysis including data acquisition, annotation, processing, publication. ReproNim helps researchers to comprehensively describe data and analysis workflows in machine-readable form (with ReproIn and Brainverse), manage the computational environments (with NICEMAN), find and share data in a FAIR fashion (with NeuroBlast). This framework facilitates the implementation of analysis in a reproducible fashion.
Chapter 3

File-based localization of numerical perturbations in data analysis pipelines

Ali Salari¹, Gregory Kiar²³, Lindsay Lewis⁴, Alan C. Evans⁴², Tristan Glatard¹

Published in: GigaScience journal https://doi.org/10.1093/gigascience/giaa106

 $^{^1 \}rm Department$ of Computer Science and Software Engineering, Concordia University, Montreal, Canada $^2 \rm McGill$ University, Montreal, Canada

³Montreal Neurological Institute, Montreal, Canada

Abstract

Data analysis pipelines are known to be impacted by computational conditions, presumably due to the creation and propagation of numerical errors. While this process could play a major role in the current reproducibility crisis, the precise causes of such instabilities and the path along which they propagate in pipelines are unclear. We present Spot, a tool to identify which processes in a pipeline create numerical differences when executed in different computational conditions. Spot leverages system-call interception through ReproZip to reconstruct and compare provenance graphs without pipeline instrumentation. By applying Spot to the structural pre-processing pipelines of the Human Connectome Project, we found that linear and non-linear registration are the cause of most numerical instabilities in these pipelines, which confirms previous findings.

3.1 Introduction

Numerical perturbations resulting from variations in computational environments impact data analyses in various fields, but identifying the origin of these perturbations in complex pipelines remains challenging. In some cases, small perturbations resulting from changes in operating system versions [47], hardware [67], or parallelization parameters [31], result in substantially different analysis outcomes, due to the propagation and amplification of floating-point errors. While the existence of such numerical errors is well known [112], their impact on scientific computations has multiplied with the rise of the Big Data era, due to the sustained growth of data sets, the increasing complexity of analysis pipelines, and the diversification of computing infrastructures. To better understand and correct these effects, efficient tools are needed to assist pipeline developers in the comparison of results obtained across different conditions.

In neuroimaging, our primary application field, data analyses often consist of hundreds of computational processes – often coming from multiple toolboxes – that are aggregated to perform a specific function. For instance, the fMRIprep pipeline [36] assembles software blocks from FSL [65], AFNI [24], FreeSurfer [41] and ANTs [7] to provide a state-of-the art functional MRI processing tool with minimal user input. Another example are the pipelines of the Human Connectome Project [44] that combine tools from FSL and FreeSurfer to pre-process structural, functional and diffusion data from their uniquely high-fidelity open dataset. In both cases, pipelines leverage toolboxes that are widely trusted in the community, yet, at the same time substantial variations in results have been observed in these toolboxes resulting from minor data or infrastructure perturbations [52, 47, 82, 70], suggesting that further investigation of their numerical conditioning is required. For such complex pipelines, a lightweight solution has to be found to perform such evaluations with limited code instrumentation.

Numerical evaluations are traditionally performed using techniques such as interval arithmetics [58] that require complete code re-writes and are therefore barely applicable to complex pipelines. Recently, Monte-Carlo Arithmetic [97, 28] provided a practical way to evaluate the uncertainty of numerical results without the need to rewrite the application in a different paradigm. By perturbating floating-point computations, it introduces a controllable amount of noise in the pipelines, effectively sampling results from a random distribution. While this technique is very appealing, it suffers from two main issues that make it impractical at the scale of a complete pipeline. First, it requires that all software components be recompiled for MCA instrumentation, which is not always feasible. Second, it multiplies the execution time by a factor of 10 to 100, which is impractical when executions already take a few hours to complete.

We present Spot, a tool to identify the source of numerical differences in complex pipelines without instrumentation. Using system-call interception through the ReproZip tool [103], Spot traverses graphs of processes and intermediary files to pinpoint the pipeline components that are unstable across execution conditions. When differences start accumulating, effectively masking any further instability, it restores clean data copies through a set of wrapper scripts. Wrapper scripts are also used to restore temporary data that might have been deleted during the execution, and to disambiguate files that have been written by multiple processes. The remainder of this paper presents the design of Spot, and its application to pre-processing pipelines of the HCP project.

3.2 Tool description

Spot identifies the components in a pipeline, at the resolution level of a system process, that produce different results in different execution conditions. First, a directed bipartite provenance graph is recorded for each pipeline execution, where nodes represent application processes and files, and edges represent read and write file accesses (Figure 3a). Second, transient files, i.e., files that are either deleted during pipeline execution or modified by multiple processes, are identified and disambiguated, resulting in a provenance DAG (Directed Acyclic Graph) in which file nodes have a single parent (in-degree of 1) (Figure 3b). DAGs produced in different conditions are then compared, in a step-by-step execution that prevents the propagation of differences in the pipeline (Figure 3c). The resulting labeled graph identifies the non-reproducible processes in the pipeline.

To ensure that a file can be unambiguously associated with the process that created it, we assume that the pipeline can be transformed such that:

- 1. Processes don't run concurrently;
- 2. Each process sequentially reads, computes, and writes.

In practice, pipeline processes may still run concurrently provided that they don't write concurrently to the same files. A process may also interleave file writes with computing, for instance when different file blocks are processed sequentially. However, only a single version of the file must eventually be made available to the other processes. In particular, in case a process deletes a file that it had created itself, this file must not be used by any other







(a) Raw provenance graph (ReproZip output), with transient files shown in gray boxes.

(b) Provenance DAG, with disambiguated transient files.

(c) Labeled DAG comparing 2 execution conditions, showing 1 non-reproducible process.

Figure 3: Provenance graphs created from the example pipeline in Listing 1. Processes are represented with circles, files with rectangles, and read/write accesses with plain edges. For convenience, the process tree is also shown, with gray dashed edges. Processes forked by bet were captured by ReproZip while they did not appear in Listing 1. Processed associated with executables located in /usr/bin/ or /bin/ are not shown.

process. Finally, we also require that processes are associated to a command line (executable and arguments), to facilitate process instrumentation.

3.2.1 Recording provenance graphs

We use ReproZip [103] to capture: (1) the set of processes created by the pipeline, and (2) the set of files read and written by each process, including temporary files. ReproZip collects this information through the ptrace() system call, with no required instrumentation of the pipeline. Using the ReproZip trace, Spot reconstructs a provenance graph by creating process and file nodes and by adding directed edges corresponding to file reads and writes (Figure 3a). We assume that provenance graphs are identical for the ReproZip traces obtained from the same subjects in different operating systems.

Provenance graphs are often data-dependent, due to variations in input data that may trigger differing branching or looping patterns across executions, for example. Some of these differences can be neglected: for instance, when a data decompression step is present at the beginning of the execution for some subjects only. Other differences cannot: for instance, when entirely different processing paths are used for different datasets. Spot includes helpers to identify different instances of provenance graphs, such as supporting the clustering of process trees, where nodes are processes and edges are fork() or clone() system calls,

Listing 1 Example pipeline that computes the volume of the brain from a T1 image.

```
#!/usr/bin/env bash
if [ $# != 1 ]
then
    echo "usage: $0 <input_image.nii.gz>"
    exit 1
fi
# Parse argument, set output file names
input_image=$1
# Run FSL bet, put result in £{bet_output}
bet ${input_image} output.nii.gz
# Create binary mask
fslmaths output.nii.gz -bin output.nii.gz
echo "Voxels / volume in binarized brain mask:"
fslstats output.nii.gz -V > voxels.txt
# Remove temporary file
\rm output.nii.gz
```

using the tree edit distance [128] implemented in Python's zss package.

3.2.2 Capturing transient files

We capture temporary files by replacing every process P by a wrapper that first calls P and then saves the produced temporary files to a read-only directory. This process replacement is done by pre-pending to the PATH environment variable a directory that contains a wrapper script named after the executable called by P.

Files written by multiple processes are disambiguated using a similar technique. For a file F written by the processes in $\mathbf{P} = \{P_1, \ldots, P_n\}$, we first check that processes in \mathbf{P} do not write concurrently to F, which would violate our assumptions. Then, we replace every process P_i by a PATH-based wrapper that first calls P_i and then saves F to a read-only directory. In this way, successive versions of F are preserved for comparison. We finally update the provenance graph accordingly, so that all files in the graph have an in-degree of 1 (Figure 3b). This operation also makes the provenance graph acyclic, since we assumed that a process could only release a single version of a file.

3.2.3 Labeling processes

After capturing transient files in the first condition (i.e. operating system, library version, etc.), we re-run the pipeline step by step in the second one to label processes. The output files created by a process in both conditions are compared: if no differences are found, the process is marked as reproducible; otherwise, the process is marked as non-reproducible, and the output files produced in the first condition are copied to the second one, to ensure that differences do not propagate further in the pipeline. Processes are instrumented transparently through a modification of the PATH variable similar to the one described previously. By default, differences in output files are identified by comparing file checksums. Other comparison functions can also be defined for specific file types, for instance to ignore file headers or file sections containing timestamps. Spot finally creates a labeled provenance graph highlighting non-reproducible processes.

Figure 3c illustrates a hypothetical incremental labeling of the example in Listing 1. Process **bet2** is labeled as non-reproducible (red) as it produces files with differences. To prevent the propagation of these differences, the files produced by **bet2** in Condition 2 are replaced with the files produced by **bet2** in Condition 1. Processes **fslmaths** and **fslstats** are then executed and labeled as reproducible (green) as they produce files without differences.

The labeled graph can differ depending on the order of executions in which condition we capture transient files or execute the pipeline to pinpoint the propagation of differences. Therefore, we run the comparison in both condition orders, and we label a process as nonreproducible (red) if it creates different results in at least one condition order.

3.2.4 Implementation

Spot is implemented in Python (i=3.6). In this work we used Spot version 0.2 and the following version of the Python package dependencies: NumPy v1.19.0 [96] and Pandas v1.0.5 [87], for data manipulations, SciPy v1.5.1 [118] and Scikit-learn v.0.23.1 [98] for the clustering of provenance graphs, Zss v1.2.0 [128] for tree distances, ReproZip v1.0.11 for the capture of provenance traces, Docker v17.05 [89] for the edition of container images, and Boutiques v0.5.25 [46] for uniform pipeline executions.

Software users will mostly have to interact with the Boutiques and ReproZip packages. Boutiques is a flexible description framework for containerized pipelines, required by the pipelines analyzed in Spot. It provides a JSON schema to describe inputs, outputs and their dependencies. Examples, tutorials and usage documentation are available at http: //boutiques.github.io. ReproZip intercepts system calls to identify the files and processes involved in a pipeline execution. Before using Spot, users have to collect ReproZip traces of their pipeline executions. Examples in the Spot documentation include ReproZip provenance capture. More documentation on ReproZip is available at https://www.reprozip.org.

3.3 Experiments

We applied Spot to the minimal pre-processing pipelines released by the Human Connectome Project (HCP), a leading initiative in neuroimaging.

3.3.1 HCP pipelines and dataset

The HCP developed a set of pre-processing pipelines to process structural, functional, and diffusion MRI data acquired in the project. We focus on HCP pre-processing pipelines for structural data, and particularly on PreFreeSurfer and FreeSurfer. A detailed description of the analyses done by these pipelines is available in [44]. In summary, the PreFreeSurfer pipeline consists of the following steps:

- Gradient Distortion Correction (DC),
- Alignment and Anatomical Average (AAve), T1w(s), T2w(s),
- Anterior/Posterior Commissure Alignment (ACPC-A),
- Brain Extraction (BExt),
- Bias Field Correction (BFC),
- Atlas-Registration (AR).

And the FreeSurfer pipeline consists of the following:

- Image downsampling,
- T1w image registration,
- T1w image segmentation,
- Surface placement,
- Surface registration.

We randomly selected 20 unprocessed subjects from the HCP data release S500 available in the ConnectomDB repository as a subset of the 1200 Subject Release (see Supplementary Table S1). For each subject, available data consisted of 1 or 2 T1-weighted images and 1 or 2 T2-weighted images, with $256 \times 320 \times 320$ voxels of size $0.7 \times 0.7 \times 0.7$ mm. Acquisition protocols and parameters are detailed in [116].

3.3.2 Data processing

We built Docker images for the HCP pre-processing pipelines v3.19.0 (PreFreeSurfer and FreeSurfer) in CentOS 6.9 (Final) and CentOS 7.4 (Core), available on DockerHub. Container images contain the HCP software dependencies, including FSL (version 5.0.6), FreeSurfer (version 5.3.0-HCP, CentOS4 build), and Connectome Workbench (version 1.0).

We processed the 20 subjects with PreFreeSurfer and FreeSurfer, using the 2 CentOS versions. The PreFreesurfer results obtained in CentOS6 were used as the input of FreeSurfer in both conditions. We also used the ReproZip trace file captured in CentOS6 for labeling the processes in both pipelines. Each subject was processed twice on the same operating system to detect within-OS variability coming from pseudo-random operations. We compared pipeline results using FreeSurfer tools mri_diff, mris_diff, and lta_diff, to ignore execution-specific information such as file path or timestamps. To compare segmentations X and Y, we used the Dice coefficient defined as follows:

$$DICE = \frac{2|X \cap Y|}{|X| + |Y|}$$

The Dice coefficient [30] is a commonly used metric to validate medical image segmentation. Dice values range from 0 to 1, with 1 indicating a perfect overlap between two segmentation results and 0 indicating no overlap. Alternatively, the Jaccard coefficient [61] could be used; there is a direct correspondence between both metrics.

3.4 Results

All experiments were run on a machine with a 3.4GHz, 8-core Intel Core i7 processor, 32GB of RAM, CentOS 7.3.1611, and Linux kernel version 3.10. The processing time, output file size, number of file accesses and number of processes observed in PreFreeSurfer and FreeSurfer are shown in Table 5. The scripts and analyses used to create the figures in this section are available at https://github.com/big-data-lab-team/HCP-reproducibility-paper.

	PreFreeSurfer		FreeSurfer	
	Mean	Standard error	Mean	Standard error
Processing time (mins)	106.67	2.68	650.25	8.88
Output file size (GB)	2.8	0.10	4.15	0.15
Number of file accesses	$94,\!089$	$2,\!645$	62,729	984
Number of processes	8,731	198	4,031	47

Table 2: Execution statistics of the pipelines per subject.

Within-OS differences

We did not observe any within-OS difference in PreFreeSurfer. In FreeSurfer, we identified 2 processes leading to within-OS differences due to the use of pseudo-random numbers: image registration with mri_segreg, and cortical surface curvature estimations with mris_curvature. Fixing the random seed used in FreeSurfer removed these differences.

Between-OS differences in PreFreeSurfer

We identified four types of subjects with different PreFreeSurfer provenance graphs (Table 3). Differences between subject types came from different numbers of T1 and T2 images in the raw data. We verified that the provenance graphs were identical for all subjects of the same type, for both versions of CentOS.

Figure 4 shows the frequency of non-reproducible pipeline processes in PreFreeSurfer. The processes identified as non-reproducible were observed in linear registration with FSL flirt (in ACPC-Alignment, Brain Extraction, Distortion Correction, and Atlas Registration), in non-linear registration with FSL fnirt (in Brain Extraction and Atlas Registration), and in image warping with FSL new_invwarp (in Brain Extraction and Atlas Registration). Differences were also observed in image mean computations with FSL maths (in Anatomical Average). Figure 5 shows a complete PreFreeSurfer labeled DAG, localizing the observed differences in the entire pipeline, for a given subject.

Type	Number of Subjects	Number of T1w images	Number of T2w images
1	9	2	2
2	8	1	1
3	1	1	2
4	2	2	1

Table 3: Types of provenance graphs in PreFreeSurfer.



Figure 4: Heatmap of non-reproducible processes across PreFreeSurfer pipeline steps. Each cell represents the occurrence of a particular command line in a pipeline step among Anatomical Average (AAve), Anterior/Posterior Commissure Alignment (ACPC-A), Brain Extraction (BExt), Bias Field Correction (BFC), or Atlas-Registration (AR). Cell labels indicate the fraction of subjects for which the corresponding process wasn't reproducible. For example, the flirt tool was invoked 6 times in step DC for each of the 20 subjects: 2 instances weren't reproducible in 19 subjects, 3 instances were always reproducible, and 1 instance wasn't reproducible in 17 subjects. Grey cells indicate that the process did not occur in the corresponding pipeline step.

Figure 6 compares fnirt results in Brain Extraction for a particular subject using the checkerboard pattern, a common method to illustrate the magnitude of the differences in registration results. Differences appear to be visually important, in particular in the areas framed in red, to the point that most experimenters would likely reject such a registration following visual quality control.

Between-OS differences in FreeSurfer

The only non-reproducible process identified by Spot in FreeSurfer was mris_make_surfaces (cortical and white matter surfaces generation), a dynamically-linked executable that produced different results for 10 out of 20 subjects.

However, FreeSurfer results still differ between conditions, due to the propagation of differences created in PreFreeSurfer. We observed the effect of this propagation in FreeSurfer results, as shown in Figure 7 for whole-brain segmentations. The Dice coefficients associated with the 44 regions segmented by FreeSurfer are shown in Figure 8, showing that Dice coefficients below 0.9 are observed in most regions, and particularly in the smallest ones. However, no significant correlation between the Dice values and the region sizes was found



Figure 5: A complete provenance graph from the PreFreesurfer pipeline. Node labels use the same abbreviations as in Figure 4. For better visualization, processes associated with commands in /bin or /usr/bin were omitted, as well as imtest, imcp, remove_ext, fslval, avscale, and fslhd.



Figure 6: Differences between T2 fnirt results in PreFreeSurfer's Brain Extraction (CentOS6 vs CentOS7). The colored squares indicate results obtained with CentOS6 (in purple) and CentOS7 (in green). The red boxes highlight regions with significant differences between the two OSes. An animated version of the comparison is available here for better visualization.

(Pearson's coefficient = 0.12, p-value = 0.43).

3.5 Discussion

Our results provide insights on the reproducibility of neuroimaging pipelines, and on the relevance of the approach implemented in Spot for reproducibility studies.

3.5.1 Key findings

Linear and non-linear registration with FSL were found to frequently lead to differences between results obtained with different operating systems. This does not come as a surprise given the instabilities associated with these processes. It also corroborates our previous findings in [47], where fMRI pre-processing with FSL was found to vary across operating systems starting from the motion correction step, a step that uses FSL's flirt tool internally. It would be relevant to investigate if the observed instability of registration processes generalizes to other toolkits, or if it remains specific to FSL. In view of the effect of small data perturbations in a variety of toolboxes and processes, such as cortical surface extraction using FreeSurfer and CIVET [82] or connectome estimation using Dipy [74], it is probable that this observation generalizes widely across toolboxes and requires a deeper investigation of the stability of linear and non-linear registration.

While only a handful or processes were found non-reproducible across the tested operating systems, the effect of such instabilities were found to propagate widely in the pipelines, and



Figure 7: Sum of binarized differences between whole-brain FreeSurfer segmentations obtained from PreFreeSurfer processings in CentOS6 vs CentOS7 (N=20). Segmentations were resampled and overlaid to the MNI152 volume template. Each voxel shows the number of subjects for which different results were observed between CentOS 6 and CentOS 7. An animated comparison of segmentations obtained for a particular subject is available here for better visualization.

to substantially impact the segmentations created by FreeSurfer. This illustrates the need to conduct reproducibility studies on entire pipelines rather than isolated processes. It also highlights the need for a deeper stability analysis of pipeline processes.

As is shown in Figure 4, the reproducibility of a given tool may vary across subjects and across processing parameters. For instance, linear registration with flirt seems to be fully reproducible in the Anatomical Average sub-pipeline, while it is highly non-reproducible in ACPC Alignment. In Brain Extraction, the same tool was found reproducible for some subjects only. Therefore, reproducibility studies need to be performed on several subjects. While this is common practice to some extent in neuroimaging, software tests are often executed only on a single dataset to reduce the associated computational load. Our results show that pipeline tests should encompass enough subjects to cover execution paths adequately.

Our results illustrate the type of variability that can be introduced in neuroimaging results due to operating system updates. The numerical noise introduced by operating system updates is realistic, as such updates are likely to occur throughout the time span of a neuroscience study, but it is also uncontrolled, as it originates in updates of low-level libraries by third-party developers. A possible method to study this problem more comprehensively would be to introduce controlled numerical perturbations in pipelines, which could be done by introducing noise either in the data, or in floating-point computations through Monte-Carlo arithmetic [97]. The work in [74] discusses and compares these two techniques.



Figure 8: Dice coefficients between regions segmented by FreeSurfer in CentOS6 vs CentOS7 (N=20), ordered by increasing median values. Each point represents the Dice coefficient between segmentations of a particular region obtained in CentOS 6 vs CentOS 7 for a given subject. Boxes brightness is proportional to the logarithm of the corresponding brain region size.

3.5.2 Spot evaluation

The processes identified by Spot as non-reproducible were all associated with dynamicallylinked executables. This makes complete sense as statically-linked executables are not impacted by library updates. Moreover, the hypothetical effects of hardware or Linux kernel updates were not measured, as the different operating systems were deployed in Docker containers on the same host, that is, using the same kernel and hardware.

To evaluate the reproducibility of a pipeline, Spot needs to execute it 5 times in order to (1) record a first ReproZip trace, (2) save transient files in the first condition, (3) compare results in the second condition, and repeat steps (2) and (3) for the other order of execution. It might be possible to further reduce this overhead by executing at step (2) only the processes depending on transient files, and capturing the transient files for the second condition simultaneously at step (3).

The target users of the Spot tool are primarily pipeline developers and users who have technical skills for creating Docker containers and Boutiques JSON files. We demonstrated the applicability of our approach by evaluating two of the arguably most complex pipelines in neuroimaging. Technically, these pipelines consist of a mix of tools assembled from different toolboxes through a variety of scripts written in different languages. Our file-based approach, notably enabled by ReproZip, was able to analyze these pipelines without requiring their instrumentation, which saved a very substantial technical effort. The assumptions made on the pipeline structure, related to the absence of concurrent writes, were not violated in our analysis, and are likely to not impede Spot's applicability to the most common neuroimaging pipelines.

Spot only tests pipeline reproducibility in the scope of a particular dataset. However, it is very plausible for pipeline processes to exhibit different reproducibility behaviors when executed on different datasets. Therefore, only the lack of reproducibility of a pipeline process could be guaranteed from an analysis with Spot, since proving reproducibility would require testing the pipeline on all possible datasets, in all possible environments, which is not feasible. Two elements could be considered in future work to address this issue. First, similar to conventional software testing, a code coverage metric could be developed to assess the fraction of the pipeline code involved in the tested dataset and parameters. This would quantify the representativity of the dataset and pipeline parameters used in the evaluation. Second, statistical risk models could be used to estimate the probability for a process to be reproducible, given a set of observations with no numerical differences. For instance, models described in [10] could be leveraged for this purpose. File-based analyses also have limitations related to the granularity at which they operate. Indeed, differences can only be identified at the level of an entire operating-system process, which can correspond to arbitrary amounts of code. Narrowing down the analysis to particular libraries, functions, or even code sections would require another approach. Similarly, Spot would not be able to detect differences in data not saved in files but instead passed to subsequent processes in memory. A common scenario in neuroimaging pipelines is that tools return results in their standard output, which is parsed by the calling process and passed to subsequent ones through variables.

Computational environments are only one of many factors contributing to the on-going reproducibility crisis. In fact, sample size selection, publication bias, or methodological flexibility in the analysis are likely to have a stronger effect than numerical perturbations, although to our knowledge no evidence of this is available. We refer to the studies in [14, 15, 12, 70 for deeper analyses of the associated effects on neuroimaging analyses. It should also be noted that the effects of computational environments and these other factors manifest at different levels: referring to the terminology used in [99], computational environments are associated with reproducibility, the minimal standard by which identical results should be obtainable from identical data and parameters, while the other aforementioned factors belong to replicability, the ultimate standard by which independent experimenters should be able to draw similar conclusions from similar experiments. In practice, variability resulting from computational environments manifests during software testing (test results depend on execution platform), deployment on HPC systems (results obtained on local vs HPC systems differ), or software version updates (results obtained before vs after the update differ), while factors related to replicability impact the community more broadly. Ultimately, both reproducibility and replicability should be understood and improved.

3.6 Conclusion

We presented Spot, a tool to detect the source of numerical differences in complex pipelines executed in different computational conditions. Spot leverages system-call interception through the ReproZip tool, and therefore can be applied to the most complex pipelines with-out requiring their instrumentation. It is available at https://github.com/big-data-lab-team/spot under MIT license.

By applying Spot to the pre-processing pipelines of the Human Connectome Project, compared in different operating systems, we showed that between-OS differences are mostly originating in linear and non-linear image registration tools. Moreover, differences introduced during image registration propagate widely in the pipelines, leading to important variability in whole-brain segmentations.

Future work will investigate in more details the numerical stability of registration algorithms. Additionally, we plan on using Monte-Carlo arithmetic to inject controlled amounts of noise in pipelines and monitor uncertainty propagation and amplification in their results.

3.7 Availability of Source Code and Requirements

- Project name: Spot
- Project home page: https://github.com/big-data-lab-team/spot
- Operating system: Linux
- Programming language: Python (3.6 or higher)
- Main dependencies: ReproZip, Docker, and Boutiques
- Other dependencies: see setup.py
- License: MIT License
- Biotools identifier: spottool
- SciCrunch ID: RRID:SCR_018915
- DOI: 10.5281/zenodo.3873219

Chapter 4

Accurate simulation of operating system updates in neuroimaging using Monte-Carlo arithmetic

Ali Salari¹, Yohan Chatelain¹, Gregory Kiar², Tristan Glatard¹

Published in:

MICCAI workshop on Uncertainty for Safe Utilization of Machine Learning in Medical Imaging (UNSURE)

https://doi.org/10.1007/978-3-030-87735-4_2

¹Department of Computer Science and Software Engineering, Concordia University, Montreal, Canada ²Center for the Developing Brain, Child Mind Institute, New York, NY, USA

Abstract

Operating system (OS) updates introduce numerical perturbations that impact the reproducibility of computational pipelines. In neuroimaging, this has important practical implications on the validity of computational results, particularly when obtained in systems such as high-performance computing clusters where the experimenter does not control software updates. We present a framework to reproduce the variability induced by OS updates in controlled conditions. We hypothesize that OS updates impact computational pipelines mainly through numerical perturbations originating in mathematical libraries, which we simulate using Monte-Carlo arithmetic in a framework called "fuzzy libmath" (FL). We applied this methodology to pre-processing pipelines of the Human Connectome Project, a flagship open-data project in neuroimaging. We found that FL-perturbed pipelines accurately reproduce the variability induced by OS updates and that this similarity is only mildly dependent on simulation parameters. Importantly, we also found between-subject differences were preserved in both cases, though the between-run variability was of comparable magnitude for both FL and OS perturbations. We found the numerical precision in the HCP pre-processed images to be relatively low, with less than 8 significant bits among the 24 available, which motivates further investigation of the numerical stability of components in the tested pipeline. Overall, our results establish that FL accurately simulates results variability due to OS updates, and is a practical framework to quantify numerical uncertainty in neuroimaging.

4.1 Introduction

Numerical round-off and cancellation errors are ubiquitous in floating-point computations. In neuroimaging, they contribute to results uncertainty along with other sources of variability, including population selection, scanning devices, sequence parameters, acquisition noise, and methodological flexibility [15, 14]. Numerical errors manifest particularly through variations in elementary mathematical libraries resulting from operating system (OS) updates. Indeed, due to implementation differences, mathematical functions available in different OS versions provide slightly different results. The impact of such epsilonesque differences on image analysis depends on the conditioning of the problem and the pipeline's numerical implementation. In neuroimaging, established image processing pipelines have been shown to be substantially impacted: for instance, differences in cortical thicknesses measured by the same Freesurfer version in different execution platforms were shown to reach statistical significance in some brain regions [52], and Dice coefficients as low as 0.6 were observed between FSL or Freesurfer segmentations obtained in different platforms [47, 107]. Such observations threaten the validity of neuroimaging results by revealing systematic instabilities.

Despite its possible implications on results validity, the effect of OS updates remains seldom studied due to (1) the lack of closed-form expressions of condition numbers for complex pipelines and non-differentiable non-linear analyses, (2) the technical challenge associated with experimental studies involving multiple OS distributions and versions, (3) the uncontrolled nature of OS updates. As a result, the effect of OS updates on neuroimaging analyses is generally neglected or handled through the use of software containers (Docker or Singularity), static executable builds, or similar approaches. While such techniques improve experiment portability, they only mask numerical instabilities and do not tackle them. Numerical perturbations are bound to reappear due to security updates [69], obsoleting software [101], or parallelization. Therefore, the mechanisms through which numerical instabilities propagate need to be investigated and eventually addressed.

This paper presents "fuzzy libmath" (FL), a framework to simulate OS updates in controlled conditions, allowing software developers to evaluate the robustness of their tools with respect to likely-to-occur numerical perturbations. As we hypothesize that numerical perturbations resulting from OS updates primarily come from implementation differences in elementary mathematical libraries, we leverage Monte-Carlo arithmetic (MCA) [97] to introduce controlled amounts of noise in these libraries. FL enables MCA in mathematical functions used by existing pipelines without the need to modify or recompile them. To demonstrate the approach, we study the effect of common OS updates on the numerical precision of structural MRI pre-processing pipelines of the Human Connectome Project [116], a major neuroimaging initiative.

4.2 Simulating OS updates with Monte-Carlo arithmetic

MCA models floating-point roundoff and cancellations errors through random perturbations, allowing for the estimation of error distributions from independent random result samples. MCA simulates computations at a given virtual precision using the following perturbation:

$$inexact(x) = x + 2^{e_x - t}\xi \tag{1}$$

where e_x is the exponent in the floating-point representation of x, t is the virtual precision and ξ is a random uniform variable of $\left(-\frac{1}{2}, \frac{1}{2}\right)$.

MCA allows for three perturbation modes: Random Rounding (RR) introduces the perturbation in function outputs, simulating roundoff errors; Precision Bounding (PB) introduces the perturbation in function operands, allowing for the detection of catastrophic cancellations; and, Full MCA combines RR and PB, resulting in the following perturbation:

$$mca_mode(x \circ y) = inexact_{RR}(inexact_{PB}(x) \circ inexact_{PB}(y))$$
 (2)

To simulate OS updates, we introduce random perturbations in the GNU mathematical library, the main mathematical library in GNU/Linux systems. Instrumenting mathematical libraries with MCA raises a number of issues as many functions assume deterministic arithmetic. For instance, applying random perturbations around a discontinuity or within piecewise approximations results in large variations and a total loss of significance that are not relevant in our context. Therefore, we have applied MCA to proxy mathematical functions wrapping those in the original library, such that only the outputs of the original functions were perturbed but not their inputs or the implementations themselves. This technique allows us to control the magnitude of the perturbation as perceived by the application.

We instrumented the GNU mathematical library with MCA using Verificarlo [28], a tool that (1) uses the Clang compiler to generate an LLVM (http://llvm.org) Intermediate Representation (IR) of the source code, (2) replaces floating-point operations in the IR by a call to the Verificarlo API, and (3) compiles the modified IR to an executable using LLVM.

The perturbation applied by the Verificarlo API can be configured at runtime, for instance to change the virtual precision applied to single- and double-precision floating-point values.

The resulting MCA-instrumented mathematical library, "fuzzy libmath" (FL), is loaded in the pipeline using LD_PRELOAD, a Linux mechanism to force-load a shared library into an executable. As a result, functions defined in fuzzy libmath transparently overload the original ones without the need to modify or recompile the analysis pipeline. Fuzzy libmath functions call the original functions through dlsym, a function that returns the memory address of a symbol. To trigger MCA instrumentation, a floating-point zero is added to the output of the original function and the result of this sum is perturbed and returned.

Finally, we measure results precision as the number of significant bits among result samples, as defined in [97]:

$$s = -\log_2 \left| \frac{\sigma}{\mu} \right| \tag{3}$$

where σ and μ are the observed cross-sample standard deviation and average.

4.3 HCP Pipelines & Dataset

We apply the methodology described above to the minimal structural pre-processing pipeline associated with the Human Connectome Project (HCP) dataset [44], entitled "PreFreeSurfer". This pipeline consists of many independent components, including: spatial distortion correction, brain extraction, cross-modal registration, and alignment to standard space. Each high-level component of this pipeline (Fig. 9) consists of several function calls using FSL, the FMRIB Software Library [65]. The pipeline requires T1w and T2w images for each subject. A full description of the pipeline is available at [44].

It should be noted that the PreFreeSurfer pipeline uses both single and double precision functions from the GNU mathematical library. Among the pre-processing steps in the pipeline, it has been shown that linear and non-linear registrations implemented in FSL FLIRT [66, 64] and FNIRT [5] are the most sensitive to numerical instabilities [107].

We selected 20 unprocessed subjects from the HCP data release S500 available in the ConnectomDB repository. We selected these subjects from different subject types to cover execution paths sufficiently. For each, the available data consisted of 1 or 2 T1w and T2w images each, with spatial dimensions of $256 \times 320 \times 320$ and voxel resolution of 0.7 mm. Acquisition protocols and parameters are detailed in [116]. Two distinct experimental configurations were tested:



Figure 9: PreFreeSurfer pipeline steps.

- **Operating Systems (OS):** subjects were processed on three different Linux operating systems inside Docker images: CentOS7 (glibc v.2.17), CentOS8 (glibc v.2.28), and Ubuntu20 (glibc v.2.31).
- Fuzzy libmath (FL): the dataset was processed on an Ubuntu20 system using fuzzy libmath. The virtual precision (t) for the perturbations was swept from 53 bits (the full mantissa for double-precision data) down to 1 bit by steps of 2. For $t \ge 24$ bits, only double-precision was altered and single-precision was set to 24 bits, and for t < 24 bits, both double- and single-precision simultaneously were changed. Three FL-perturbed samples were generated for each subject and virtual precision, to match the number of OS samples.

After conducting both experiments, we selected the virtual precision that most closely simulated the variability observed across OSes via the root-mean-square error (RMSE) between the number of significant bits per voxel in all subjects and conditions. This precision is referred to as the global nearest virtual precision and was used to compare results obtained in both the FL and OS versions.

4.4 Results

The fuzzy libmath source code, Docker image specifications, and analysis code to reproduce the results are available at https://github.com/big-data-lab-team/MCA-libmath-paper. All experiments were conducted on the Béluga HPC computing cluster made available by Compute Canada through Calcul Québec. Béluga is a general-purpose cluster with 872 available nodes. All nodes contain 2× Intel Gold 6148 Skylake @ 2.4 GHz (40 cores/node) CPU, and node memory can range between 92 to 752 GB. The average processing time of



(a) Distribution of significant bits (b) Significance map (subject average)

Figure 10: Comparison of OS and FL effects on the precision of PreFreeSurfer results for n=20 subjects. FL samples were obtained at the global nearest virtual precision of t=37 bits.

the pipeline without FL instrumentation was 69 minutes (average of 3 executions). The FL perturbation increased it to 93 minutes.

We ensured that the pipeline does not use pseudo-random numbers by processing each subject twice on the same operating system. To validate that FL was correctly instrumented with Verificarlo, we used Veritracer [20], a tool for tracing the numerical quality of variables over time. For one subject, the traces showed that the number of significant bits in the function outputs varied over time, confirming the instrumentation with MCA. Throughout the pipeline execution, Veritracer reported approximately 4 billion calls to FL, with the following ratio of calls: 47.12% log, 40.96% exp, 6.92% expf, 3.39% logf, 1.55% sincosf, and 0.06% of cumulated calls to atan2f, pow, sqrt, exp2f, powf, log10f, log10, cos, and asin. We also checked that long double types were not used.

4.4.1 Fuzzy librath accurately simulates the effect of OS updates

Fuzzy libmath accurately reproduced the effect of OS updates, both globally (Fig. 10a) and locally (Fig. 10b). The distributions of significant bits in the atlas registered T1w images were nearly identical (p > 0.05, KS test) on the average and individual subject distributions for 15/20 subjects, after correcting for multiple comparisons. Locally, the spatial distribution of significant digits also appeared to be preserved. Losses in significance were observed mainly at the brain-skull interface and between brain lobes, indicating spatial dependency of numerical properties.

The average number of significant bits in either the FL or OS conditions were 7.76 out of 24 available, which corresponds to 2.32 significant (base 10) digits. This relatively low



Figure 11: RMSE-based hierarchical clustering of OS (left) and FL (right) samples. Colors identify different subjects, showing that similarities between subjects are preserved by the numerical perturbations. Horizontal gray lines represent average RMSEs between (top line) and within (bottom line) subject clusters.

precision motivates future investigations of the stability of pipeline components, in particular for image registration.

4.4.2 Fuzzy librath preserves between-subjects image similarity

Numerically-perturbed samples remained primarily clustered by individual subjects (Fig. 11), indicating that neither FL nor OS perturbations were impactful enough to blur the differences between subjects. Notably, the similarity between subjects was also preserved by the numerical perturbation, leading to the same subject ordering in the dendrograms. However, the average RMSE within samples of a given subject was approximately $13 \times$ lower than the average RMSE between different subjects. The fact that between-subject variabilities were nearly on the same order of magnitude as OS and FL variability demonstrates the potential severity of these instabilities.

4.4.3 Results are stable across virtual precision

The FL results presented previously were obtained at the global nearest virtual precision of t=37 bits, determined as the precision which minimized the RMSE between FL and OS average maps of significant bits. We varied the virtual precision in steps of 2 between t=1 and t=53 bits (Fig. 12). On average, no noticeable RMSE change was observed between the FL and OS variability for precisions ranging from t=21 to t=53 bits, which shows that FL can robustly approximate OS updates.



Figure 12: Comparison of RMSE values computed between OS and FL results for different virtual precisions.

The observed plateau suggests the existence of an "intrinsic precision" for the pipeline, above which no improvement in results precision is expected. For the tested pipeline, this intrinsic precision was observed at t=21 bits, which indicates that the pipeline could be implemented exclusively with single-precision floating-point representations (24 bits of mantissa) without loss of results precision. This would substantially decrease the pipeline memory footprint and computational time, as approximately 88% of operations used in this pipeline made use of double-precision data. In addition, the presence of such a plateau suggests that numerical perturbations introduced by OS updates might be in the range of machine error (t=53 bits), although it is also possible that the extent of the plateau results from the numerical conditioning of the tested pipeline. It is possible in contrast that the absence of such a plateau would suggest an unstable pipeline that would benefit either from correction or larger datatypes. The ability to capture stability across a range of precisions importantly demonstrates a key advantage of using FL to simulate OS variability.

The relationship between RMSE of individual subjects was generally consistent with the average line, with the notable exception of subject 18. The observed discrepancies between this subject and potential others might be leveraged for quality control checks and, as a result, inform tool development.

The pipeline failed to complete for at least one subject below the virtual precision of t=13 bits, also referred to as the tolerance of the pipeline. Specifically, 51% of pipeline executions crashed among all subjects for precisions ranging from 1–11 bits, and there was

no relationship between tolerance-level and precision. The error raised was in the Readout Distortion Correction portion of the pipeline, and appears to stem from the FSL FAST tissue segmentation. The specific source of the error within this component is presently unknown, but is an open question for further exploration.

4.5 Conclusion & Discussion

We demonstrated fuzzy librath as an accurate method to simulate variability in neuroimaging results due to OS updates. Alongside this evaluation, fuzzy librath can be used by pipeline developers or consumers to evaluate the numerical uncertainty of tools and results. Such evaluations may also help decrease pipeline memory usage and computational time through the controlled use of reduced numerical precision. Fuzzy librath does not require any modification of the pipeline as it operates on the level of shared libraries. The accuracy of the simulations were shown to be robust across a wide range of virtual precisions, which reinforces the applicability of the method.

The proposed technique is directly applicable to MATLAB code executed with GNU Octave, to Python programs executed on Linux, and to C programs that depend on GNU libmath. Numerical noise can be introduced in other libraries, such as OpenBLAS or NumPy, using our https://github.com/verificarlo/fuzzy environment.

A commonly used approach to address instabilities resulting from OS version updates in practice is to sweep the issue under the rug of software containers or static linking. While such solutions are undoubtedly helpful to improve code portability or strict re-executability, a more honest position is to consider computational results as realizations of random variables depending on numerical error. The presented technique enables estimating result distributions, a first step toward making analyses reproducible across heterogeneous execution environments. While this work did not investigate the precise cause of numerical instabilities by tracing the system function calls, this is a topic for future work.

The tested OS versions span a timeframe of 7 years (2012–2020) and focused on GNU/Linux, a widely-used platform in neuroimaging [55]. Given that our experiments focused on numerical perturbations applied to mathematical functions, which are implemented similarly across OSes, our findings are likely to generalize to OS/X or MS Windows, although future work would be needed to confirm that. The tested pipeline is the official solution of the HCP project to pre-process data, and is considered the state-of-the-art. This pipeline assembles software components from the FSL toolbox consistent with common practice in neuroimaging, such as in fMRIPrep [36] or the FSL feat workflow [65], to which fuzzy libmath can be directly applied. Efforts are on-going to use fuzzy libmath in fMRIPrep software tests, to guarantee that bug fixes do not perturb results beyond numerical uncertainty.

The fact that the induced numerical variability preserves image similarity between subjects is reassuring and, in fact, exciting. OS updates provide a convenient, practical target to define a virtual precision leading to a detectable but still reasonable numerical perturbation. However, it is also of importance that OS- and FL-induced variability were on a similar order of magnitude as subject-level effects. This suggests that the preservation of relative betweensubject differences may not hold in all pipelines, and such a comparison could be used to evaluate the robustness of a pipeline to OS instabilities. The fact that the results observed across OS versions and FL perturbations arise from equally-valid numerical operations also suggests that the observed variability may contain meaningful signal. In particular, signal measured from these perturbations might be leveraged to enhance biomarkers, as suggested in [73] where augmenting a diffusion MRI dataset with numerically-perturbed samples was shown to improve age classification.

Chapter 5

Comparing software variability across and within fMRI analysis packages

Ali Salari¹, Yohan Chatelain¹, Alexander Bowring², Camille Maumet³, Gregory Kiar⁴, Tristan Glatard¹

¹Department of Computer-Science and Software Engineering, Concordia University, Montreal, Canada ²Li Ka Shing Centre for Health Information and Discovery, Nuffield Department of Population Health, Big Data Institute, University of Oxford, Oxford, UK

 $^{^3}$ Inria, Univ Rennes, CNRS, Inserm, IRISA UMR 6074, Empenn ERL U 1228, Rennes, France 4 Center for the Developing Brain, Child Mind Institute, New York, NY, USA

Abstract

Variability has been broadly observed in functional MRI analyses as a result of software differences both between and within analytic tools. However, the relationship between within-tool and between-tool software variabilities for a given analysis is unclear. Withintool variability has multiple origins including variations in algorithms, parametrizations, and low-level software dependencies. We focus on the effects resulting from software dependencies — in particular due to operating system and other infrastructural updates as they often remain unnoticed while algorithms and parametrizations are more controllable. We extended a previous comparison of fMRI analysis software libraries (namely FSL, AFNI, and SPM) and related the observed differences to within-tool software variability simulated using Monte-Carlo arithmetic. In group analyses, we found that between-tool software variability was consistently larger than within-tool software variability. In subject analyses, between-tool software variability also dominated within-tool variability overall, however, within-tool software variability approached between-tool software variability in some regions for some subjects. Interestingly, the brain masking instability was triggered by both withinand between-tool software variabilities. Our findings motivate the continued investigation of within-tool software variability in neuroimaging.

5.1 Introduction

Recent explorations of the analytical flexibility in brain imaging across tools, platforms, or teams, have demonstrated unexpected variability in results, even when analyzing identical data [14]. Possible explanations for such discrepancies include methodological flexibility [19] and software variability, the focus of this paper. A recommendation to address this variability is to adopt a "multiverse" approach that analyzes the same dataset multiple times in different software environments, and ultimately conclude from the resulting set of outcomes. However, the range of analytical conditions to be included in such multiverse analyses remains poorly defined, both because of the boundless set of tools and configurations, and that the precise causes of software variability remain unclear. This lack of clarity is in part because fMRI analyses depend on complex software stacks that leverage low-level libraries provided by the operating system (e.g., mathematical functions), general scientific computing methods (e.g., optimization toolboxes), and specific fMRI analysis tools (e.g., spatial image normalization methods). At each level, conceptual and implementation differences across experiments may each create substantial variability in the analysis outcomes.

As a result of these factors, the variability resulting from the use of different fMRI analysis tools implementing similar analytical approaches (between-tool software variability), or different versions of the same tool, can reach worrying magnitudes. For instance, the study in [15] compared the results produced by the three main fMRI analysis toolboxes, namely SPM [100], AFNI [23] and FSL [65], using similar pipelines. It reported limited similarity between the activation clusters produced by these tools, measured by Dice coefficients ranging from 0.0 to 0.769, where 0 indicates non-overlapping clusters and 1 means identical clusters. More recently, the work in [84] also showed a low similarity between the results produced by five different tools (ABCD [39], CCS [124], CPAC [25], DPARSF [125], fMRIPrep [36]) and identified the main factors contributing to these differences. The magnitude of the highlighted differences suggest that between-tool software variability may be playing a critical role in the reproducibility of an fMRI analysis overall.

The variability resulting from differences in lower-level software libraries (within-tool software variability) has also been quantified in fMRI. The study in [47] mentions a low similarity between the activation clusters produced by FSL using different versions of the GNU/Linux system, measured by Dice coefficients ranging from 0.0 to 1.0, covering the full spectrum of possible similarities. This variability resulted from updates in the GNU mathematical library and can be properly simulated by introducing small numerical perturbations on the results returned by mathematical functions [106]. Here again, these observations

could partly explain the differences reported in [14].

The relationship between within- and between-tool software variability are poorly understood, but could play an important role in the construction of multiverse study environments. Importantly, if associations exist between these two types of variability, they may both originate to some degree from the inherent instability of brain activity estimation from BOLD signal variations, and shed light on the numerical confidence of results. Under this hypothesis, small perturbations introduced by low-level software updates could trigger effects correlated with those created by tool variations.

This paper investigates the relationship between numerical stability — used as a proxy for within-tool software variability — and between-tool software variability through two main questions:

- 1. what is the relative magnitude of within- and between-tool software variability, and
- 2. is there an association between within- and between-tool software variability.

We address these two questions by reproducing the study in [15] and extending it with the addition of numerical perturbations of controlled magnitude.

5.2 Materials and Methods

5.2.1 fMRI analysis & Dataset

We replicated the analysis described as study 'ds000001' in [108], relying on the data publicly available in OpenNeuro at https://openneuro.org/datasets/ds000001 and using three widely-used software packages for fMRI data processing, namely FMRIB Software Library (FSL) [65], Analysis of Functional NeuroImages (AFNI) [23], and Statistical Parametric Mapping (SPM) [100]. We selected this dataset because comparable analysis pipelines implemented in FSL, AFNI and SPM were already publicly available and extensively described in [15]. Furthermore, the work in [15] already evaluated the effect of tool variability for this dataset, which we intended to extend with the present quantification of within-tool variability.

In the selected study, 16 healthy adult subjects participated in the balloon analog risk task [81] to measure risk-taking behavior over three scanning sessions [108]. We reused the preprocessing, first-level, and second-level analyses implemented by [15] consistently across all three software packages. Table 4, adapted from [15], summarizes the analytical steps in each pipeline.

		FSL	AFNI	SPM
Preprocessing	Motion Correction		\checkmark	\checkmark
	Segmentation			\checkmark
	Brain Extraction (Anatomical)	\checkmark	\checkmark	\checkmark
	Brain Extraction (Functional)		\checkmark	
	Intra-subject Coregistration	\checkmark	\checkmark	\checkmark
	Inter-subject Registration	\checkmark	\checkmark	\checkmark
	Smoothing	\checkmark	\checkmark	\checkmark
First-level	Model Specification	\checkmark	\checkmark	\checkmark
	Inclusion of 6 Motion Parameters	\checkmark	\checkmark	\checkmark
	Model Estimation	\checkmark	\checkmark	\checkmark
	Contrasts	\checkmark	\checkmark	\checkmark
Second-level	Model Specification	\checkmark	\checkmark	\checkmark
	Model Estimation	\checkmark	\checkmark	\checkmark
	Contrasts	\checkmark	\checkmark	\checkmark
	Second-level Inference	\checkmark	\checkmark	\checkmark

Table 4: Software processing steps (adapted from [15]).

5.2.2 Within-tool software variability simulation with Fuzzy Libmath

We simulated within-tool software variability by introducing numerical noise in the analyses using Fuzzy Libmath [106], a version of the GNU mathematical library (libmath) instrumented with Monte-Carlo arithmetic. Monte-Carlo arithmetic simulates numerical errors by introducing a controlled amount of noise in floating-point operations through the following perturbation [97]:

$$inexact(x) = x + 2^{e_x - t}\xi,\tag{4}$$

where e_x is the exponent in the floating-point representation of x, t is the virtual precision (the number of unperturbed bits in the mantissa of x), and ξ is a random uniform variable of $\left(-\frac{1}{2}, \frac{1}{2}\right)$. We introduced the perturbation using Verificarlo [28], an LLVM compiler supporting Monte-Carlo arithmetic and other types of numerical instrumentations.

We loaded the instrumented librath functions in the pipeline using LD_PRELOAD, a Linux mechanism to force-load a shared library into an executable. This mechanism allows functions defined in Fuzzy Librath to transparently overload the original ones without the need to modify or recompile the analysis pipeline. Fuzzy Libmath introduces numerical perturbations in the values returned by mathematical functions but not in their input values or within their implementation. This is done by wrapping the original functions and applying function inexact to their returned values. Listing 5.1 shows an example of this wrapping for the log function in single and double precision. In this wrapper, the original function is called through dlsym, a function that returns the memory address of a symbol — in our case RTLD_NEXT, the address of the next occurrence of the function in memory. Compiling function wrappers with Verificarlo instruments the result of the addition between the original function output and the floating-point zero.

Listing 5.1: Sample wrapper function (C code)

```
#include <dlfcn.h>
#include <math.h>
static double (*real_log)(double dbl);
static float (*real_logf)(float dbl);
// Override
double log(double dbl);
{
        real_log = dlsym(RTLD_NEXT, "log");
        return real_log(dbl) + 0.0;
}
float logf(float dbl);
{
        real_logf = dlsym(RTLD_NEXT, "logf");
        return real_logf(dbl) + 0.0f;
}
```

In [106], Fuzzy Libmath was shown to accurately simulate the effect of Linux operating system updates in structural pre-processing pipelines of the Human Connectome Project which are largely based on FSL. To validate our pipeline instrumentations for the present study, we first verified that non-instrumented executions of the same pipeline on the same dataset led to identical results. We also listed the pipeline library dependencies using the 1dd Linux utility and verified that (1) the tested pipelines were dynamically linked to the GNU libmath library, and (2) there was no alternative implementation of elementary mathematical functions in the pipeline dependencies. Finally, we verified that the use of Fuzzy Libmath affected computational results.

5.2.3 Data processing

We measured between-tool software variability (BT) by running the pipelines described in [15] with FSL version 5.0.10, AFNI version 18.1.09, and SPM12 version r7771 executed with GNU/Octave version 5.2. These software versions were identical to the ones used in [15] except for SPM for which we had to use a more recent version in order to be able to use it with GNU/Octave instead of MATLAB to enable mathematical function instrumentation using Fuzzy Libmath. Indeed, MATLAB uses its own built-in mathematical functions, which prevents the use of Fuzzy Libmath. In AFNI, we set the number of threads to 7 even though AFNI executions in [15] were single-threaded. This was meant to reduce the time overhead resulting from Fuzzy Libmath instrumentation. All the analyses were conducted on the CentOS 7.3 operating system. The computations were performed on Compute Canada's Béluga cluster nodes, each with $2 \times$ Intel Gold 6148 Skylake @ 2.4 GHz (40 cores/node) CPU and 8 GB of RAM per core. To facilitate portability and reproducibility, we encapsulated the above-mentioned software packages in Docker container images based on CentOS 7.3 which we converted to Singularity images to enable running on cluster nodes.

In the IEEE-754 format, a floating-point number N is represented by a sign bit (s), an exponent (e), and a pseudo-mantissa (m) as:

$$N = (-1)^s \times 2^E \times (1.m),\tag{5}$$

where E = e - 127 for single precision and E = e - 1023 for double precision. In single precision, the pseudo-mantissa has 23 bits and the exponent has 8 bits, whereas in double precision the pseudo-mantissa has 52 bits and the exponent has 11 bits. Machine error, the smallest absolute difference between two floating-point numbers that can be represented, is in general $2^{-24} \times e$ in single precision and $2^{-53} \times e$ in double precision. We simulated machine error by introducing the random perturbation defined in Equation 4 at the virtual precision of t = 24 bits for single-precision values and t = 53 bits for double-precision ones. The perturbed pseudo-mantissa is obtained by extending the size of the original pseudo-mantissa, applying the perturbation at the desired virtual precision, and rounding the perturbed result back to the original precision. We simulated within-tool software variability (WT) by running the same analyses 10 times using Fuzzy Libmath with a virtual precision of t = 53 bits for double-precision values and t = 24 bits for single-precision values. Perturbations were applied at the same magnitude of floating-point precisions during the pipeline execution regardless of what it corresponds to. These values were chosen such that the numerical perturbation affects all the pipeline steps and simulates machine error originating from software, hardware,
and operating system changes. The resulting samples are equally plausible estimates of the true numerical result at the precision used by the pipelines.

We evaluated BT for thresholded as well as unthresholded group-level and subject-level t-statistic maps by computing the differences of t-statistic maps across tools. For WT, we computed the average difference across the 10 Fuzzy Libmath samples.

Further, from the thresholded maps, we determined regional instability between activation clusters in the 360 regions in the Human Connectome Project Multi-Modal Parcellation atlas version 1.0 (HCP-MMP1.0) [43]. For BT, we considered a region unstable for a pair of tools if it contained activated voxels for a tool but not for the other one. For WT, we considered a region unstable for a pair of tools (A, B) if for tool A or tool B it contained activated voxels only for some Fuzzy Libmath samples.

5.3 Results

The scripts, Docker images, and derived datasets that were used to reproduce the results are available in our GitHub repository at https://github.com/big-data-lab-team/ fuzzy-neurotools.

5.3.1 Validation of replication

We verified the correctness of our analyses by comparing our unperturbed t-statistic group maps with the ones obtained in [15]. For FSL, we found strictly identical results with the same MD5 checksums. For SPM, the files (as measured by MD5 checksums) were different but differences were visually unnoticeable. For AFNI, the overall patterns of activations were preserved, however, differences were noticeable visually. The observed differences remained small (see Supplemental Material S1), and might be due to the use of GNU/Octave vs MATLAB in SPM, and of multithreading in AFNI. We performed visual quality control of the AFNI and SPM results for each individual subject and confirmed that T1-weighted images were correctly skull-stripped and registered to the MNI template.

5.3.2 In the group analysis, BT was larger than WT

Table 5 presents summary statistics for variability in the group-level t-statistic. Mean of differences in t-statistic maps was centered around zero for BT and WT in both thresholded and unthresholded maps (t-test $p < 10^{-5}$). However, for each tool pair (A, B), BT variability

was significantly larger than WT in tool A or B (Wilcoxon signed-rank test and t-test $p < 10^{-5}$). These global differences were confirmed by Bland-Altman plots showing a clear dominance of BT over WT (Figure 13-A,B). In addition, results show that BT and WT were not correlated across voxels (Pearson's $r \approx 0$, $p < 10^{-5}$, Figure 14).

		Thresh	Grou nolded	o map Unthresholded		Subject maps Unthresholded	
		μ	σ	μ	σ	$\mid \mu$	σ
Between Tools	FSL vs. SPM	-0.845	2.639	-0.052	1.122	-0.059	0.911
(BT)	FSL vs. AFNI	-0.871	3.216	0.048	1.405	0.027	1.097
	AFNI vs. SPM	0.426	3.214	-0.113	1.543	-0.104	1.208
Within Tool	FSL	-0.397	1.218	-0.020	0.189	0.017	0.207
(WT)	SPM	0.142	1.035	0.000	0.139	0.000	0.126
	AFNI	-0.045	1.282	0.000	0.373	0.000	0.355

Table 5: Voxel-wise mean and standard deviation of BT and WT variability in t-statistic maps.



Figure 13: **A** and **B**: Bland-Altman plots comparing group-level differences computed between tools (\mathbf{A}) and within tools at machine error (\mathbf{B}) .



Figure 14: Voxel-wise comparison of group-level differences in BT and WT.

5.3.3 In subject analyses, WT approached BT for some subjects

Table 5 also presents summary statistics for subject-level unthresholded t-statistic maps (mean over subjects). As for group-level maps, the mean of differences in t-statistic maps was approximately zero centered in both BT and WT (t-test $p < 10^{-4}$ for 13 subjects), and BT was consistently larger than WT (Wilcoxon signed-rank test and t-test $p < 10^{-4}$ for 14 subjects). However, for some subjects and for AFNI, WT approached BT in some regions (Figure 15, and see Supplemental Material S2).

5.3.4 Previous results were confirmed in thresholded group maps

We also compared BT and WT in thresholded maps since these maps are commonly used instead of unthresholded ones to conclude on the activation of specific regions. Thresholding is an unstable operation that introduced variability at the edges of active regions for both BT and WT. Except at the edges, BT remained consistently larger than WT (Figure 16-A,B,C and Table 5).

Moreover, to compare the effects of BT and WT on the detection of activated regions, we measured WT instability and BT instability in each region of the HCP-MMP1.0 parcellation. The confusion matrices in Figure 16-**D** report these instabilities for the 360 tested regions. The average ratio of unstable regions was 26.2% for BT and 20.6% for WT, which confirmed that BT was larger than WT even in thresholded maps.

Finally, we measured agreement between BT and WT from the confusion matrices since



Figure 15: For subject with highest WT variability, unthresholded subject-level variability computed between tools (\mathbf{A}) , and within tools at machine error (\mathbf{B}) .

interpreting the correlation of thresholded maps is difficult due to the discontinuity introduced by thresholding. The average Cohen's kappa score⁵ computed from the confusion matrices between WT instability and BT instability was $\kappa = 0.2$, indicating a moderate agreement between WT instability and BT instability consistent with the correlations observed in unthresholded maps.

5.3.5 Brain masking instability was triggered by WT and BT

While the previous results were obtained over the union of t-statistic maps between the pair of tools, we also computed results in the intersection of maps and found that the brain masking instability was triggered by both BT and WT (Table 6). BT and WT variabilities were generally smaller than the ones obtained previously (Table 5) in both thresholded and unthresholded maps. The effect of masking appeared in Bland-Altman plots by showing that blob lines were removed after fixing masks (Figure 17- \mathbf{A} , \mathbf{B}).

 $^{{}^5\}kappa \leq 0$ denotes chance agreement, $-1 \leq \kappa \leq 1$

		Group map				Subject maps	
		Inresnolaed		Unthresholded		Unthresholded	
		μ	σ	μ	σ	μ	σ
Between Tools	FSL vs. SPM	-0.668	1.056	-0.033	1.121	-0.068	0.937
(BT)	FSL vs. AFNI	-0.823	1.857	0.048	1.439	0.034	1.146
	AFNI vs. SPM	0.645	1.428	0.064	1.489	-0.065	1.206
Within Tool	FSL	0.012	0.127	0.000	0.165	0.017	0.199
(WT)	SPM	0.002	0.154	0.000	0.121	0.000	0.110
	AFNI	-0.003	0.3548	0.003	0.285	0.000	0.328

Table 6: Voxel-wise mean and standard deviation of BT and WT variability in intersected t-statistic maps.

5.4 Discussion

In fMRI group analyses, within-tool software variability remains an order of magnitude smaller than between-tool variability. This is likely explained by the fact that group analyses benefit from regularization of numerical noise, which is expected to increase with sample size. This finding is consistent with observations made in [72] from diffusion MRI data where connectome graph statistics were found to be substantially unstable at the subject level while group distributions remained consistent. Therefore, for fMRI studies with large sample sizes, within-tool software variability may be neglected with respect to between-tool variability. In particular, multiverse analyses aggregating the outcome of multiple analysis tools are likely to successfully correct for machine error in such group studies. Nevertheless, within-tool software variability remains substantial in group analyses that are based on a single tool, as is commonly the case in current fMRI studies. In particular, in our study, the inherent instability of thresholding was triggered by within-tool software variability in 20% of 360 brain regions, which indicates that it might have impacted neuroscientific conclusions related to these regions.

In subject-level analyses, within- and between-tool software variabilities can become of comparable magnitude for some subjects in some regions. This observation is particularly relevant to the development of fMRI-based biomarkers aiming at individualized phenotype predictions. Machine error may play a non-negligible role in such analyses, even when predictions combine results produced by multiple tools.

For both group- and subject-level analyses, between- and within-tool software variabilities similarly trigger the brain masking instability. Even though both types of variability are different in nature, this result suggests that in some cases they may have a common cause that might be related to the conditioning of fMRI analysis in a specific dataset. In such cases, numerical stability may be a suitable proxy to study between-tool variability, which could be of practical value.

Our results are limited by the type of numerical noise introduced in the analyses. Indeed, we only perturbed the outputs of elementary mathematical functions while numerical noise could creep in any floating-point operation. Therefore, our estimation of within-tool software variability should be considered a lower bound. Likewise, our estimation of tool varability is likely to be underestimated, having tested only 3 analytical pipelines among the thousands available [19].

In conclusion, between-tool variability clearly dominates within-tool variability overall. However, within-tool variability still reaches magnitudes that may affect scientific conclusions derived from data analyses — in particular in subject analyses and in group analyses that rely on thresholded maps. Moreover, sources of variability such as brain masking affect both within-tool and between-tool variability: in these specific cases only, within-tool variability may be used as a proxy to between-tool variability. These findings motivate further investigations of within-tool variability in fMRI analyses.



Figure 16: **A**,**B**,**C**: Thresholded group-level t-statistic within tools at machine error for FSL (A), SPM (B) and AFNI (C). Arrows point to activation clusters impacted by both withinand between-tool variability. **D**: Confusion matrices of activation instability in BT and WT among the 360 regions of the HCP-MMP1.0 parcellation.(with masking)



Figure 17: **A** and **B**: Bland-Altman plots comparing group-level differences computed over the intersections between tools (\mathbf{A}) and within tools at machine error (\mathbf{B}) .

Supplemental Materials

S1 Reproduced results

Figure S1 shows the difference in SPM and AFNI group analyses maps between the results in [15] and our replication. Numerical perturbations of 1 ulp are likely to have been introduced by our replication due to the use of GNU/Octave vs MATLAB for SPM and multithreading for AFNI. However, the group maps remained very similar overall, which led us to conclude that our results correctly reproduced the ones in [15].



Figure S1: Differences between reproduced and original results obtained in [15] of unthresholded group-level t-statistic for SPM (left) and AFNI (right). The highest areas of difference in AFNI seem to be due to differences in brain masks.

S2 Maps of t-statistics for subject with highest WT variability

Figure S2 plots the maps of t-statistic for each run of tools for the subject with highest WT variability.



Figure S2: For subject with highest WT variability, unthresholded subject-level t-statistic within tools at machine error for FSL (\mathbf{A}) , SPM (\mathbf{B}) , AFNI (\mathbf{C}) .

Chapter 6

Discussion

The aim of this thesis was to study the numerical stability of neuroimaging pipelines focusing on the effect of operating system variations. We leveraged system call interception techniques, the ReproZip tool and perturbation models such as Monte-Carlo arithmetic (MCA) [97]. We showed that This chapter first discusses the findings and contributions. It then discuss the implication of the results and limitations of the current methodology. This chapter concludes with recommendations for future studies.

6.1 The Impact of Numerical Perturbations

We showed the role of numerical errors in different computational environments in neuroimaging pipelines. In Chapter 3, we introduced Spot, a tool to detect the source of numerical differences in pipelines executed on different operating systems. This work localized the origin of processes that hamper bitwise reproducibility due to OS updates. We observed registration processes as the highly contributing source of instability in the FSL tool. We found that differences can propagate during the pipeline execution and substantially impact the results of the entire pipeline (e.g., segmentations created by FreeSurfer). We also obtained different results between subjects, showing the sensitivity of the analyses to dataset selection. The irreproducible processes identified by Spot were all associated with dynamically linked executables.

The numerical differences were caused by truncation and round-off errors due to the precision limitations of the floating-point arithmetic. This is uncontrolled error that originated from the updates of the system libraries. The observed instability across OSes depends on which operating systems are used. This limits the evaluation of the underlying stability of a particular tool.

In Chapter 4, we presented a framework to model the numerical uncertainty induced by OS updates in a control condition using Monte-Carlo arithmetic. We obtained an accurate simulation of OS variations using the fuzzy libmath library, a version of the GNU mathematical library (libmath) instrumented with Monte-Carlo arithmetic. We found that the pipeline could be implemented exclusively with lower precision of floating-point representations (single-precision) without loss of results precision. This would substantially decrease the pipeline memory footprint and computational time. However, the finding showed a very low number of significant digits, 5 out of 15 available digits. This motivates further investigation of the numerical stability of the main components of the pipelines, such as linear and non-linear registrations. Also, it is notable that OS- and FL-induced variability were on a similar order of magnitude as subject-level effects, showing that we applied a reasonable amount of perturbations.

This framework works on shared libraries, so there is no need for recompilation or modification of the pipeline or any other sources. We believe that pipelines that depend on different third-party mathematical libraries could be studied similarly by building fuzzy versions of these dependencies. This motivated us to investigate numerical quality in more applications using the proposed MCA-based method.

Finally, in Chapter 5, we investigated the changes to pipeline results due to variations between software packages. We compared numerical variability with tool variability across three of the most popular software packages in neuroimaging. In fMRI group analyses, numerical variability was found to be an order of magnitude smaller than tool variability. We also found that numerical instability in individual analyses was attenuated in group analyses. We obtained more uncertainty on thresholded results than unthresholded maps, probably due to different thresholds used in different tools.

6.2 The Importance of Numerical Instability

This study was a contribution to uncertainty quantification in medical imaging. Numerical errors are often neglected or only partially studied due to the associated engineering challenges among the various sources of uncertainty involved in medical imaging results, including population selection, scanning devices and sequence parameters, acquisition noise, image reconstruction algorithms, and methodological flexibility. Our work proposed methodologies and implementations to address this issue.

The current approach to address numerical instability resulting from OS updates is mainly to ignore the issue and hide it using Docker containers or other types of virtualization. Although building static program and containerization techniques improve reproducibility across OSes, but small differences remained. It remains that computational results should be understood as realizations of a random variable resulting from floating-point arithmetic. The presented techniques in this thesis enable estimating result distributions, the initial step toward making analyses reproducible across execution environments, including HPC systems, GPU accelerators, or merely different workstations.

Results showed the significance of numerical instabilities in neuroimaging pipelines and demonstrated numerical analysis techniques such as MCA as valuable methods for evaluating the associated variability. Moreover, a related work [72] suggested that capturing this variability may improve the robustness of scientific findings. This finding highlights how numerical variability may be a feature that should be taken into considerations by the pipeline developers.

We demonstrated the numerical noise sensitivity of two of the most complex pipelines in neuroimaging, HCP preprocessing pipelines, PreFreeSurfer and FreeSurfer. These pipelines consist of a mix of tools assembled from different toolboxes through a variety of scripts written in different languages. In addition to the preprocessing steps, we evaluated the numerical stability of a complete fMRI analysis using three of the most popular tools in neuroimaging, FSL, SPM, and AFNI. Thus, the results presented in this thesis would apply to a wide range of other pipelines. However, the current methodology is limited to Linux operating systems. Our findings are likely to generalize to OS/X or MS Windows, although future work would be needed to confirm that.

Moreover, the interposition technique is only applicable to intercept system calls in dynamically linked programs. Further investigations are needed to evaluate the stability of pipelines with statically linked executables.

6.3 Recommendations for Future Research

Among the processes that contribute to pipeline instability, some of them substantially have a higher effect on results. For instance, iterative computations accumulate rounding errors and significantly amplify these errors. Considering this, the next steps in this area could be evaluating the stability of pipeline components as well as the entire pipeline. This would help to identify processes that propagate and amplify errors and then substitute them for more stable tools that perform a similar function, ultimately improving the quality of the pipeline. It is important to propose numerical debugging tools to identify such numerical bugs during the execution of pipelines. Numreical errors are a well-studied problem in the floating-point error characterization research field, which has resulted in the development of debugging tools such as VeriTracer [20], FpDebug [11], Verrou [40]. These tools automatically replace all the floating-point operations with their MCA counterparts or other stochastic arithmetics to evaluate the numerical quality of the computation and pinpoint the parts of the source code. They instrument pipelines at compile time or runtime using Valgrind-based tools, which have limitations, including the need to access the source code, pre-knowledge of the functions or variables in the pipeline, and computation time overhead.

Using a combination of techniques developed in this thesis, in future work, we can create a debugging tool at the level of system call processes without recompiling the source code to identify the processes with the highest impact on results in the pipeline due to the numerical errors. This functionality could leverage the tool developed in Chapter 3 and the interposition technique to inject MCA perturbations described in Chapter 4. Moreover, we could use the principle of minimization by delta-debugging to search through a large number of processes. The delta-debugging algorithm is a general technique that can automatically narrow downthe failure caused by changing circumstances that are critical to producing a bug, such as program input, program code, environmental configurations, etc. [126, 127]. Instead of working on the source code, we can apply delta-debugging to the program history by comparing various versions until the faulty change is found. Finding these processes could narrow down further investigations toward stabilizing the pipelines.

Improving stability needs actions different from evaluation and localization of instability. There is not a general approach to stabilize the entire pipeline; we should look at the relevant code section to find an appropriate solution. For example, we found that linear and non-linear registration processes introduce errors in both FreeSurfer and PreFreeSurfer pipelines. We know that the registration procedure is sensitive to the initialization of the optimization method used [45]. A possible method to address such instabilities is using the bootstrapping technique. In [45], the authors explained that bootstrapping is an efficient technique to improve the robustness of motion estimation. The bootstrap version of the pipelines computed the median transformation results from the 30 samples. In addition to bootstrapping, it is shown that the bagging technique can reduce the effect of the medians of the parameters of the 30 transformations. Bagging, also called bootstrap aggregating, is a simple and powerful ensemble method to reduce both bias and variance in the results. So, we can possibly stabilize pipelines and improve their accuracy using aggregates of results obtained with data perturbations. However, it is a compute-intensive technique that adds a significant computation time overhead, so that should be used only when no other solution to stability is available. Further study must be conducted into the processes involved in analysis instability.

An exciting research topic for future work could be finding that the variability arising from numerical perturbations may contain meaningful signals. The perturbation model is not only helpful as a method for measuring the stability of analyses but that it can improve their quality. We can also apply other types of perturbations. Given that FL only perturbs basic mathematical functions, we expect more numerical variability by perturbing operations that rely on the linear algebra libraries BLAS and LAPACK. In future work, this can be evaluated using MCA-instrumented versions of BLAS and LAPACK along with other libraries available in the Fuzzy project in Verificarlo's GitHub repository at github.com/verificarlo/fuzzy.

6.4 Conclusion

The numerical stability of computational analyses plays an important role in the reproducibility of scientific findings. These analyses are sensitive to the computing environment changes such as operating system and analysis toolboxes, particularly in computationally intensive domains where results rely on a series of complex computations. Throughout this thesis, we demonstrated the impact of numerical instabilities on neuroimaging results. We implemented Spot, a tool that localizes irreproducible processes between operating system variations. Moreover, we presented an MCA-based method to apply numerical perturbations in floating-point operations, simulating OS variability and studying the numerical instabilities more comprehensively. We used the Linux interception utility LD_PRELOAD, which transparently interposes system calls to an instrumented counterpart. We expanded our findings by capturing numerical variability among different software packages and comparing software variability across and within fMRI analysis packages. As the successful completion of this thesis, we built freely available tools that enable software developers and researchers to evaluate the stability of their pipelines and results when dynamically linked mathematical libraries are changed. The findings of this thesis could be used to narrow down further investigations toward stabilizing pipelines.

Bibliography

- [1] Neuroimagin Data Model. NIDM-Results. http://nidm.nidash.org/specs/ nidm-results_130.html. [Online; accessed 6-October-2021].
- [2] SPM-Statistical Parametric Mapping. https://www.fil.ion.ucl.ac.uk/spm/. [Online; accessed 6-October-2021].
- [3] Association for Computing Machinery. Artifact Review and Badging., 2021. [Online; accessed 6-October-2021].
- [4] Yasser Ad-Dab'bagh, O. Lyttelton, JS Muehlboeck, C. Lepage, D. Einarson, K. Mok, O. Ivanov, RD Vincent, J. Lerch, E. Fombonne, et al. The CIVET Image-Processing Environment: A Fully Automated Comprehensive Pipeline for Anatomical Neuroimaging Research. In *Proceedings of the 12th Annual Meeting of the Organization for Human Brain Mapping*, page 2266. Florence, Italy, 2006.
- [5] Jesper LR Andersson, Mark Jenkinson, and Stephen Smith. Non-Linear Registration, aka Spatial Normalisation FMRIB. Technical Report TR07JA2, FMRIB Analysis Group of the University of Oxford, 2007.
- [6] Atlassian. Git LFS Large File Storage Atlassian Git Tutorial, 2021. [Online; accessed 6-October-2021].
- Brian B. Avants, Nick Tustison, and Gang Song. Advanced Normalization Tools (ANTS). Insight j, 2:1–35, 2009.
- [8] Monya Baker. 1,500 Scientists Lift the Lid on Reproducibility. Nature News, 533(7604):452, 2016.
- Cagri Balkesen. In-memory Parallel Join Processing on Multi-Core Processors. PhD thesis, 2014.

- [10] Bernard Beauzamy. Méthodes Probabilistes pour l'étude des Phénomènes Réels. Société de Calcul Mathématiques, SA" Algorithmes et Optimisation", 2004.
- [11] Florian Benz, Andreas Hildebrandt, and Sebastian Hack. A Dynamic Program Analysis to Find Floating-Point Accuracy Problems. ACM SIGPLAN Notices, 47(6):453–462, 2012.
- [12] Nikhil Bhagwat, Amadou Barry, Erin W. Dickie, Shawn T. Brown, Gabriel A. Devenyi, Koji Hatano, Elizabeth DuPre, Alain Dagher, M. Mallar Chakravarty, Celia MT Greenwood, et al. Understanding the Impact of Preprocessing Pipelines on Neuroimaging Cortical Surface Analyses. *bioRxiv*, 2020.
- [13] Carl Boettiger. An Introduction to Docker for Reproducible Research. ACM SIGOPS Operating Systems Review, 49(1):71–79, 2015.
- [14] Rotem Botvinik-Nezer, Felix Holzmeister, Colin F. Camerer, Anna Dreber, Juergen Huber, Magnus Johannesson, Michael Kirchler, Roni Iwanir, Jeanette A. Mumford, R. Alison Adcock, et al. Variability in the Analysis of a Single Neuroimaging Dataset by Many Teams. *Nature*, 582(7810):84–88, 2020.
- [15] Alexander Bowring, Camille Maumet, and Thomas E. Nichols. Exploring the Impact of Analysis Software on Task fMRI Results. *Human Brain Mapping*, 40:1–23, 2019.
- [16] Leo Breiman. Bagging Predictors. Machine Learning, 24(2):123–140, 1996.
- [17] Leo Breiman. Heuristics of Instability and Stabilization in Model Selection. The Annals of Statistics, 24(6):2350–2383, 1996.
- [18] Steven P. Callahan, Juliana Freire, Emanuele Santos, Carlos E. Scheidegger, Cláudio T. Silva, and Huy T. Vo. VisTrails: Visualization Meets Data Management. In Proceedings of the 2006 ACM SIGMOD International Conference on Management of Data, pages 745–747. ACM, 2006.
- [19] Joshua Carp. On the Plurality of (Methodological) Worlds: Estimating the Analytic Flexibility of fMRI Experiments. *Frontiers in Neuroscience*, 6:149, 2012.
- [20] Yohan Chatelain, Pablo de Oliveira Castro, Eric Petit, David Defour, Jordan Bieder, and Marc Torrent. VeriTracer: Context-Enriched Tracer for Floating-Point Arithmetic Analysis. In 2018 IEEE 25th Symposium on Computer Arithmetic (ARITH), pages 61–68. IEEE, 2018.

- [21] James Cheney, Anthony Finkelstein, Bertram Ludäscher, and Stijn Vansummeren. Principles of Provenance (Dagstuhl Seminar 12091). In *Dagstuhl Reports*, volume 2. Schloss Dagstuhl-Leibniz-Zentrum fuer Informatik, 2012.
- [22] Fernando Chirigati, Rémi Rampin, Dennis Shasha, and Juliana Freire. Reprozip: Computational reproducibility with Ease. In *Proceedings of the 2016 International Conference on Management of Data*, pages 2085–2088. ACM, 2016.
- [23] Robert W. Cox. AFNI: Software for Analysis and Visualization of Functional Magnetic Resonance Neuroimages. Computers and Biomedical Research, 29(3):162–173, 1996.
- [24] Robert W. Cox. AFNI: What a Long Strange Trip It's Been. Neuroimage, 62(2):743– 747, 2012.
- [25] Cameron Craddock, Sharad Sikka, Brian Cheung, Ranjeet Khanuja, Satrajit S Ghosh, Chaogan Yan, Qingyang Li, Daniel Lurie, Joshua Vogelstein, Randal Burns, et al. Towards Automated Analysis of Connectomes: The Configurable Pipeline for the Analysis of Connectomes (c-pac). Front Neuroinform, 42:10–3389, 2013.
- [26] Samir Das, Alex P. Zijdenbos, Dario Vins, Jonathan Harlap, and Alan C. Evans. LORIS: A Web-Based Data Management System for Multi-Center Studies. *Frontiers in Neuroinformatics*, 5:37, 2012.
- [27] James Demmel and Hong Diep Nguyen. Numerical Reproducibility and Accuracy at Exascale. In 2013 IEEE 21st Symposium on Computer Arithmetic, pages 235–237. IEEE, 2013.
- [28] Christophe Denis, Pablo de Oliveira Castro, and Eric Petit. Verificarlo: Checking Floating Point Accuracy through Monte Carlo Arithmetic. In 2016 IEEE 23nd Symposium on Computer Arithmetic (ARITH), pages 55–62, 2016.
- [29] Paolo Di Tommaso, Maria Chatzou, Evan W. Floden, Pablo Prieto Barja, Emilio Palumbo, and Cedric Notredame. Nextflow Enables Reproducible Computational Workflows. *Nature Biotechnology*, 35(4):316, 2017.
- [30] Lee R. Dice. Measures of the Amount of Ecologic Association Between Species. *Ecology*, 26(3):297–302, 1945.
- [31] Kai Diethelm. The Limits of Reproducibility in Numerical Simulation. Computing in Science & Engineering, 14(1):64–72, 2011.

- [32] Kai Diethelm. The Limits of Reproducibility in Numerical Simulation. Computing in Science & Engineering, 14(1):64–72, 2012.
- [33] David L. Donoho, Arian Maleki, Inam Ur Rahman, Morteza Shahram, and Victoria Stodden. Reproducible Research in Computational Harmonic Analysis. *Computing in Science & Engineering*, 11(1), 2009.
- [34] Peter D. Düben, Hugh McNamara, and Tim N. Palmer. The Use of Imprecise Processing to Improve Accuracy in Weather & Climate Prediction. *Journal of Computational Physics*, 271:2–18, 2014.
- [35] Anders Eklund, Thomas E. Nichols, and Hans Knutsson. Cluster Failure: Why fMRI Inferences for Spatial Extent Have Inflated False-Positive Rates. Proceedings of the National Academy of Sciences, page 201602413, 2016.
- [36] Oscar Esteban, Christopher J. Markiewicz, Ross W. Blair, Craig A. Moodie, A. Ilkay Isik, Asier Erramuzpe, James D. Kent, Mathias Goncalves, Elizabeth DuPre, Madeleine Snyder, et al. fMRIPrep: A Robust Preprocessing Pipeline for Functional MRI. Nature Methods, 16(1):111–116, 2019.
- [37] Alan C. Evans, Sean Marrett, Peter Neelin, Louis Collins, Keith Worsley, Weiqian Dai, Sylvain Milot, Ernst Meyer, and Daniel Bub. Anatomical Mapping of Functional Activation in Stereotactic Coordinate Space. *Neuroimage*, 1(1):43–53, 1992.
- [38] Suvarna Fadnavis. Some Numerical Experiments on Round-Off Error Growth in Finite Precision Numerical Computation. arXiv preprint physics/9807003, 1998.
- [39] Eric Feczko, Greg Conan, Scott Marek, Brenden Tervo-Clemmens, Michaela Cordova, Olivia Doyle, Eric Earl, Anders Perrone, Darrick Sturgeon, Rachel Klein, et al. Adolescent Brain Cognitive Development (ABCD) Community MRI Collection and Utilities. *bioRxiv*, 2021.
- [40] François Févotte and Bruno Lathuilière. Debugging and Optimization of Hpc Programs with the Verrou Tool. In 2019 IEEE/ACM 3rd International Workshop on Software Correctness for HPC Applications (Correctness), pages 1–10. IEEE, 2019.
- [41] Bruce Fischl. FreeSurfer. *Neuroimage*, 62(2):774–781, 2012.

- [42] Eleftherios Garyfallidis, Matthew Brett, Bagrat Amirbekian, Ariel Rokem, Stefan Van Der Walt, Maxime Descoteaux, and Ian Nimmo-Smith. Dipy, a Library for the Analysis of Diffusion MRI Data. *Frontiers in Neuroinformatics*, 8:8, 2014.
- [43] Matthew F. Glasser, Timothy S. Coalson, Emma C. Robinson, Carl D. Hacker, John Harwell, Essa Yacoub, Kamil Ugurbil, Jesper Andersson, Christian F. Beckmann, Mark Jenkinson, et al. A Multi-Modal Parcellation of Human Cerebral Cortex. *Nature*, 536(7615):171–178, 2016.
- [44] Matthew F. Glasser, Stamatios N. Sotiropoulos, J. Anthony Wilson, Timothy S. Coalson, Bruce Fischl, Jesper L. Andersson, Junqian Xu, Saad Jbabdi, Matthew Webster, Jonathan R. Polimeni, et al. The Minimal Preprocessing Pipelines for the Human Connectome Project. *Neuroimage*, 80:105–124, 2013.
- [45] Tristan Glatard and Pierre Bellec. Numerical Stability of Motion Estimation in fMRI Time Series. In Annual Meeting of the Organization for Human Brain Mapping, 2018.
- [46] Tristan Glatard, Gregory Kiar, Tristan Aumentado-Armstrong, Natacha Beck, Pierre Bellec, Rémi Bernard, Axel Bonnet, Shawn T. Brown, Sorina Camarasu-Pop, Frédéric Cervenansky, et al. Boutiques: A Flexible Framework to Integrate Command-Line Applications in Computing Platforms. *GigaScience*, 7(5):giy016, 2018.
- [47] Tristan Glatard, Lindsay B. Lewis, Rafael Ferreira da Silva, Reza Adalat, Natacha Beck, Claude Lepage, Pierre Rioux, Marc-Etienne Rousseau, Tarek Sherif, Ewa Deelman, et al. Reproducibility of Neuroimaging Analyses Across Operating Systems. *Frontiers in Neuroinformatics*, 9:12, 2015.
- [48] Steven N. Goodman, Daniele Fanelli, and John PA Ioannidis. What Does Research Reproducibility Mean? Science Translational Medicine, 8(341):341ps12–341ps12, 2016.
- [49] K. Gorgolewski, Oscar Esteban, Gunnar Schaefer, B. Wandell, and R. Poldrack. Open-Neuro—a Free Online Platform for Sharing and Analysis of Neuroimaging Data. Organization for Human Brain Mapping. Vancouver, Canada, 1677, 2017.
- [50] Krzysztof Gorgolewski, Christopher D. Burns, Cindee Madison, Dav Clark, Yaroslav O. Halchenko, Michael L. Waskom, and Satrajit S. Ghosh. Nipype: A Flexible, Lightweight and Extensible Neuroimaging Data Processing Framework in Python. *Frontiers in neuroinformatics*, 5:13, 2011.

- [51] Krzysztof J. Gorgolewski, Tibor Auer, Vince D. Calhoun, R. Cameron Craddock, Samir Das, Eugene P. Duff, Guillaume Flandin, Satrajit S. Ghosh, Tristan Glatard, Yaroslav O. Halchenko, et al. The Brain Imaging Data Structure, a Format for Organizing and Describing Outputs of Neuroimaging Experiments. *Scientific Data*, 3:160044, 2016.
- [52] Ed H. B. M. Gronenschild, Petra Habets, Heidi I. L. Jacobs, Ron Mengelers, Nico Rozendaal, Jim van Os, and Machteld Marcelis. The Effects of FreeSurfer Version, Workstation Type, and Macintosh Operating System Version on Anatomical Volume and Cortical Thickness Measurements. *PloS One*, 7(6):e38234, January 2012.
- [53] Philip Guo. CDE: A Tool for Creating Portable Experimental Software Packages. Computing in Science & Engineering, 14(4):32–35, 2012.
- [54] Yaroslav O. Halchenko, Kyle Meyer, Benjamin Poldrack, Debanjum Singh Solanky, Adina S. Wagner, Jason Gors, Dave MacFarlane, Dorian Pustina, Vanessa Sochat, Satrajit S. Ghosh, et al. DataLad: Distributed System for Joint Management of Code, Data, and Their Relationship. *Journal of Open Source Software*, 6(63):3262, 2021.
- [55] Michael Hanke and Yaroslav O. Halchenko. Neuroscience Runs on GNU/Linux. Frontiers in Neuroinformatics, 5:8, 2011.
- [56] Khawar Hasham, Kamran Munir, and Richard McClatchey. Cloud Infrastructure Provenance Collection and Management to Reproduce Scientific Workflows Execution. *Future Generation Computer Systems*, 86:799–820, 2018.
- [57] Hess, Joey. git-annex: A Distributed File Synchronization System Written in Haskell., 2021. [Online; accessed 6-October-2021].
- [58] Timothy Hickey, Qun Ju, and Maarten H. Van Emden. Interval Arithmetic: From Principles to Implementation. Journal of the ACM (JACM), 48(5):1038–1068, 2001.
- [59] David RC Hill. Numerical Reproducibility of Parallel and Distributed Stochastic Simulation Using High-Performance Computing. In *Computational Frameworks*, pages 95–109. Elsevier, 2017.
- [60] John PA Ioannidis. Why Most Published Research Findings Are False. PLoS Medicine, 2(8):e124, 2005.

- [61] Paul Jaccard. The Distribution of the Flora in the Alpine Zone. 1. New phytologist, 11(2):37–50, 1912.
- [62] Clifford R. Jack Jr, Matt A. Bernstein, Nick C. Fox, Paul Thompson, Gene Alexander, Danielle Harvey, Bret Borowski, Paula J. Britson, Jennifer L. Whitwell, Chadwick Ward, et al. The Alzheimer's Disease Neuroimaging Initiative (ADNI): MRI Methods. Journal of Magnetic Resonance Imaging: An Official Journal of the International Society for Magnetic Resonance in Medicine, 27(4):685–691, 2008.
- [63] Yves Janin, Cédric Vincent, and Rémi Duraffort. CARE, the Comprehensive Archiver for Reproducible Execution. In Proceedings of the 1st ACM SIGPLAN Workshop on Reproducible Research Methodologies and New Publication Models in Computer Engineering, page 1. ACM, 2014.
- [64] Mark Jenkinson, Peter Bannister, Michael Brady, and Stephen Smith. Improved Optimization for the Robust and Accurate Linear Registration and Motion Correction of Brain Images. *Neuroimage*, 17(2):825–841, 2002.
- [65] Mark Jenkinson, Christian F. Beckmann, Timothy EJ Behrens, Mark W. Woolrich, and Stephen M. Smith. FSL. *Neuroimage*, 62(2):782–790, 2012.
- [66] Mark Jenkinson and Stephen Smith. A Global Optimisation Method for Robust Affine Registration of Brain Images. *Medical Image Analysis*, 5(2):143–156, 2001.
- [67] Fabienne Jézéquel, Jean-Luc Lamotte, and Issam Saïd. Estimation of Numerical Reproducibility on CPU and GPU. In 2015 Federated Conference on Computer Science and Information Systems (FedCSIS), pages 675–680. IEEE, 2015.
- [68] Jorge Jovicich, Silvester Czanner, Xiao Han, David Salat, Andre van der Kouwe, Brian Quinn, Jenni Pacheco, Marilyn Albert, Ronald Killiany, Deborah Blacker, et al. MRI-Derived Measurements of Human Subcortical, Ventricular and Intracranial Brain Volumes: Reliability Effects of Scan Sessions, Acquisition Sequences, Data Analyses, Scanner Upgrade, Scanner Vendors and Field Strengths. *Neuroimage*, 46(1):177–192, 2009.
- [69] Bhupinder Kaur, Mathieu Dugré, Aiman Hanna, and Tristan Glatard. An Analysis of Security Vulnerabilities in Container Images for Scientific Data Analysis. *GigaScience*, 10(6):giab025, 2021.

- [70] David N. Kennedy, Sanu A. Abraham, Julianna F. Bates, Albert Crowley, Satrajit Ghosh, Tom Gillespie, Mathias Goncalves, Jeffrey S. Grethe, Yaroslav O. Halchenko, Michael Hanke, et al. Everything Matters: The ReproNim Perspective on Reproducible Neuroimaging. *Frontiers in Neuroinformatics*, 13:1, 2019.
- [71] David N. Kennedy, Christian Haselgrove, Jon Riehl, Nina Preuss, and Robert Buccigrossi. The NITRC Image Repository. *Neuroimage*, 124:1069–1073, 2016.
- [72] Gregory Kiar, Yohan Chatelain, Pablo de Oliveira Castro, Eric Petit, Ariel Rokem, Gael Varoquaux, Bratislav Misic, Alan C. Evans, and Tristan Glatard. Numerical Uncertainty in Analytical Pipelines Lead to Impactful Variability in Brain Networks. *PloS One*, 16(11), 2021.
- [73] Gregory Kiar, Yohan Chatelain, Ali Salari, Alan C. Evans, and Tristan Glatard. Data Augmentation through Monte Carlo Arithmetic Leads to More Generalizable Classification in Connectomics. *Neurons, Behavior, Data analysis, and Theory*, 2021.
- [74] Gregory Kiar, Pablo de Oliveira Castro, Pierre Rioux, Eric Petit, Shawn T. Brown, Alan C. Evans, and Tristan Glatard. Comparing Perturbation Models for Evaluating Stability of Neuroimaging Pipelines. *The International Journal of High Performance Computing Applications*, 34(5):491–501, 2020.
- [75] Jaan Kiusalaas. Numerical Methods in Engineering with Python 3. Cambridge University Press, 2013.
- [76] Avi Kivity, Yaniv Kamay, Dor Laor, Uri Lublin, and Anthony Liguori. kvm: The Linux Virtual Machine Monitor. In *Proceedings of the Linux symposium*, volume 1, pages 225–230. Dttawa, Dntorio, Canada, 2007.
- [77] Dagmar Krefting, Michael Scheel, Alina Freing, Svenja Specovius, Friedemann Paul, and Alexander Brandt. Reliability of Quantitative Neuroimage Analysis Using Freesurfer in Distributed Environments. In MICCAI Workshop on High-Performance and Distributed Computing for Medical Imaging. (Toronto, ON), 2011.
- [78] Gina R. Kuperberg, Matthew R. Broome, Philip K. McGuire, Anthony S. David, Marianna Eddy, Fujiro Ozawa, Donald Goff, W. Caroline West, Steven CR Williams, Andre JW van der Kouwe, et al. Regionally Localized Thinning of the Cerebral Cortex in Schizophrenia. Archives of General Psychiatry, 60(9):878–888, 2003.

- [79] Gregory M. Kurtzer, Vanessa Sochat, and Michael W. Bauer. Singularity: Scientific Containers for Mobility of Compute. *PloS One*, 12(5):e0177459, 2017.
- [80] Manja Lehmann, Abdel Douiri, Lois G. Kim, Marc Modat, Dennis Chan, Sebastien Ourselin, Josephine Barnes, and Nick C. Fox. Atrophy Patterns in Alzheimer's Disease and Semantic Dementia: A Comparison of FreeSurfer and Manual Volumetric Measurements. *Neuroimage*, 49(3):2264–2274, 2010.
- [81] Carl W. Lejuez, Jennifer P. Read, Christopher W. Kahler, Jerry B. Richards, Susan E. Ramsey, Gregory L. Stuart, David R. Strong, and Richard A. Brown. Evaluation of a Behavioral Measure of Risk Taking: The Balloon Analogue Risk Task (BART). Journal of Experimental Psychology: Applied, 8(2):75, 2002.
- [82] Lindsay Lewis, Claude Lepage, Najmeh Khalili-Mahani, Mona Omidyeganeh, Seun Jeon, Patrick Bermudez, Alex Zijdenbos, Robert Vincent, Reza Adalat, and Alan Evans. Robustness and Reliability of Cortical Surface Reconstruction in CIVET and FreeSurfer. In Annual Meeting of the Organization for Human Brain Mapping, 2017.
- [83] Lindsay B. Lewis, Claude Y. Lepage, and Alan C. Evans. Utilizing the BigBrain as Ground Truth for Evaluation of CIVET & Freesurfer Structural MRI Pipelines. In Annual Meeting of the Organization for Human Brain Mapping, 2018.
- [84] Xinhui Li, Lei Ai, Steve Giavasis, Hecheng Jin, Eric Feczko, Ting Xu, Jon Clucas, Alexandre Franco, Anibal Sólon Heinsfeld, Azeez Adebimpe, , et al. Moving Beyond Processing and Analysis-Related Variation in Neuroscience. *bioRxiv*, 2021.
- [85] Allan J. MacKenzie-Graham, Arash Payan, Ivo D. Dinov, John D. Van Horn, and Arthur W. Toga. Neuroimaging Data Provenance Using the LONI Pipeline Workflow Environment. In *International Provenance and Annotation Workshop*, pages 208–220. Springer, 2008.
- [86] Daniel S. Marcus, Timothy R. Olsen, Mohana Ramaratnam, and Randy L. Buckner. The Extensible Neuroimaging Archive Toolkit. *Neuroinformatics*, 5(1):11–33, 2007.
- [87] Wes McKinney. pandas: A Foundational Python Library for Data Analysis and Statistics. Python for High Performance and Scientific Computing, 14(9), 2011.
- [88] Maarten Mennes, Bharat B. Biswal, F. Xavier Castellanos, and Michael P. Milham. Making Data Sharing Work: The FCP/INDI Experience. *Neuroimage*, 82:683–691, 2013.

- [89] Dirk Merkel. Docker: Lightweight Linux Containers for Consistent Development and Deployment. *Linux journal*, 2014(239):2, 2014.
- [90] Michael Peter Milham. Open Neuroscience Solutions for the Connectome-Wide Association Era. Neuron, 73(2):214–218, 2012.
- [91] Paolo Missier, Khalid Belhajjame, and James Cheney. The W3C PROV Family of Specifications for Modelling Provenance Metadata. In Proceedings of the 16th International Conference on Extending Database Technology, pages 773–776. ACM, 2013.
- [92] Veronika I. M"uller, Edna C. Cieslik, Ilinca Serbanescu, Angela R. Laird, Peter T. Fox, and Simon B. Eickhoff. Altered Brain Activity in Unipolar Depression Revisited: Meta- Analyses of Neuroimaging Studies. JAMA Psychiatry, 74(1):47–55, 2017.
- [93] Ingo Müller, Andrea Arteaga, Torsten Hoefler, and Gustavo Alonso. Reproducible Floating-Point Aggregation in RDBMSs. arXiv Preprint arXiv:1802.09883, 2018.
- [94] Thomas E. Nichols, Samir Das, Simon B. Eickhoff, Alan C. Evans, Tristan Glatard, Michael Hanke, Nikolaus Kriegeskorte, Michael P. Milham, Russell A. Poldrack, Jean-Baptiste Poline, et al. Best Practices in Data Analysis and Sharing in Neuroimaging Using MRI. *Nature Neuroscience*, 20(3):299, 2017.
- [95] Tom Oinn, Matthew Addis, Justin Ferris, Darren Marvin, Martin Senger, Mark Greenwood, Tim Carver, Kevin Glover, Matthew R. Pocock, Anil Wipat, et al. Taverna: A Tool for the Composition and Enactment of Bioinformatics Workflows. *Bioinformatics*, 20(17):3045–3054, 2004.
- [96] Travis E. Oliphant. A Guide to NumPy, volume 1. Trelgol Publishing USA, 2006.
- [97] Douglass Stott Parker. Monte Carlo Arithmetic: Exploiting Randomness in Floating-Point Arithmetic. University of California (Los Angeles). Computer Science Department, 1997.
- [98] Fabian Pedregosa, Gaël Varoquaux, Alexandre Gramfort, Vincent Michel, Bertrand Thirion, Olivier Grisel, Mathieu Blondel, Peter Prettenhofer, Ron Weiss, Vincent Dubourg, et al. Scikit-learn: Machine learning in Python. The Journal of Machine Learning Research, 12:2825–2830, 2011.
- [99] Roger D. Peng. Reproducible Research in Computational Science. Science, 334(6060):1226–1227, 2011.

- [100] William D. Penny, Karl J. Friston, John T. Ashburner, Stefan J. Kiebel, and Thomas E. Nichols. Statistical Parametric Mapping: The Analysis of Functional Brain Images. Elsevier, 2011.
- [101] Jeffrey M. Perkel. Challenge to Scientists: Does Your Ten-Year-Old Code Still Run? Nature, 584(7822):656–658, 2020.
- [102] Hans E. Plesser. Reproducibility vs. Replicability: A Brief History of a Confused Terminology. Frontiers in Neuroinformatics, 11:76, 2018.
- [103] Rémi Rampin, Fernando Chirigati, Dennis Shasha, Juliana Freire, and Vicky Steeves. ReproZip: The Reproducibility Packer. Journal of Open Source Software, 1(8):107, 2016.
- [104] Nathalie Revol and Philippe Théveny. Numerical Reproducibility and Parallel Computations: Issues for Interval Algorithms. *IEEE Transactions on Computers*, 63(8):1915– 1924, 2014.
- [105] David E. Rex, Jeffrey Q. Ma, and Arthur W. Toga. The LONI Pipeline Processing Environment. *Neuroimage*, 19(3):1033–1048, 2003.
- [106] Ali Salari, Yohan Chatelain, Gregory Kiar, and Tristan Glatard. Accurate Simulation of Operating System Updates in Neuroimaging Using Monte-Carlo arithmetic. In Uncertainty for Safe Utilization of Machine Learning in Medical Imaging, and Perinatal Imaging, Placental and Preterm Image Analysis, pages 14–23. Springer, 2021.
- [107] Ali Salari, Gregory Kiar, Lindsay Lewis, Alan C. Evans, and Tristan Glatard. File-Based Localization of Numerical Perturbations in Data Analysis Pipelines. *Giga-Science*, 9(12), 12 2020.
- [108] Tom Schonberg, Craig R. Fox, Jeanette A. Mumford, Eliza Congdon, Christopher Trepel, and Russell A. Poldrack. Decreasing Ventromedial Prefrontal Cortex Activity During Sequential Risk-Taking: An fMRI Investigation of the Balloon Analog Risk Task. Frontiers in Neuroscience, 6:80, 2012.
- [109] Matthias Schwab, N. Karrenbach, and Jon Claerbout. Making Scientific Computations Reproducible. Computing in Science & Engineering, 2(6):61–67, 2000.
- [110] Ji Suk Shim, Jin Sook Lee, Jeong Yol Lee, Yeon Jo Choi, Sang Wan Shin, and Jae Jun Ryu. Effect of Software Version and Parameter Settings on the Marginal and Internal

Adaptation of Crowns Fabricated with the CAD/CAM System. *Journal of Applied Oral Science*, 23(5):515–522, 2015.

- [111] Victoria Stodden, Marcia McNutt, David H. Bailey, Ewa Deelman, Yolanda Gil, Brooks Hanson, Michael A. Heroux, John PA Ioannidis, and Michela Taufer. Enhancing Reproducibility for Computational Methods. *Science*, 354(6317):1240–1241, 2016.
- [112] Josef Stoer and Roland Bulirsch. Introduction to Numerical Analysis, volume 12. Springer Science & Business Media, 2013.
- [113] Michela Taufer, Omar Padron, Philip Saponaro, and Sandeep Patel. Improving Numerical Reproducibility and Stability in Large-Scale Numerical Simulations on GPUs. In Parallel & Distributed Processing (IPDPS), 2010 IEEE International Symposium on, pages 1–9. IEEE, 2010.
- [114] Torvalds, Linus and Hamano, Junio. Git: A Free and Open Source Distributed Version Control System, 2021. [Online; accessed 6-October-2021].
- [115] J-Donald Tournier, Fernando Calamante, and Alan Connelly. MRtrix: Diffusion Tractography in Crossing Fiber Regions. International Journal of Imaging Systems and Technology, 22(1):53-66, 2012.
- [116] David C. Van Essen, Stephen M. Smith, Deanna M. Barch, Timothy EJ Behrens, Essa Yacoub, Kamil Ugurbil, and Wu-Minn HCP Consortium. The WU-Minn Human Connectome Project: An Overview. *Neuroimage*, 80:62–79, 2013.
- [117] David C. Van Essen, Kamil Ugurbil, Edward Auerbach, Deanna Barch, Timothy EJ Behrens, Richard Bucholz, Acer Chang, Liyong Chen, Maurizio Corbetta, Sandra W. Curtiss, et al. The Human Connectome Project: A Data Acquisition Perspective. *Neuroimage*, 62(4):2222–2231, 2012.
- [118] Pauli Virtanen, Ralf Gommers, Travis E. Oliphant, Matt Haberland, Tyler Reddy, David Cournapeau, Evgeni Burovski, Pearu Peterson, Warren Weckesser, Jonathan Bright, et al. SciPy 1.0: Fundamental Algorithms for Scientific Computing in Python. Nature Methods, 17:261–272, 2020.
- [119] Lina Wadi, Mona Meyer, Joel Weiser, Lincoln D. Stein, and Jüri Reimand. Impact of Outdated Gene Annotations on Pathway Enrichment Analysis. *Nature Methods*, 13(9):705, 2016.

- [120] Jon Watson. Virtualbox: Bits and Bytes Masquerading as Machines. Linux Journal, 2008(166):1, 2008.
- [121] Wikipedia contributors. Out-of-Order Execution Wikipedia, The Free Encyclopedia, 2021. [Online; accessed 6-October-2021].
- [122] Wikipedia contributors. VMware Workstation Wikipedia, The Free Encyclopedia, 2021. [Online; accessed 6-October-2021].
- [123] Mark D. Wilkinson, Michel Dumontier, IJsbrand Jan Aalbersberg, Gabrielle Appleton, Myles Axton, Arie Baak, Niklas Blomberg, Jan-Willem Boiten, Luiz Bonino da Silva Santos, Philip E. Bourne, et al. The FAIR Guiding Principles for Scientific Data Management and Stewardship. *Scientific Data*, 3, 2016.
- [124] Ting Xu, Zhi Yang, Lili Jiang, Xiu-Xia Xing, and Xi-Nian Zuo. A Connectome Computation System for Discovery Science of Brain. *Science Bulletin*, 60(1):86–95, 2015.
- [125] Chaogan Yan and Yufeng Zang. DPARSF: A Matlab Toolbox for" Pipeline" Data Analysis of Resting-State fMRI. Frontiers in Systems Neuroscience, 4:13, 2010.
- [126] Andreas Zeller. Yesterday, My Program Worked. Today, It Does Not. Why? ACM SIGSOFT Software Engineering Notes, 24(6):253–267, 1999.
- [127] Andreas Zeller. Isolating Cause-Effect Chains from Computer Programs. ACM SIG-SOFT Software Engineering Notes, 27(6):1–10, 2002.
- [128] Kaizhong Zhang and Dennis Shasha. Simple Fast Algorithms for the Editing Distance between Trees and Related Problems. SIAM Journal on Computing, 18(6):1245–1262, 1989.