

Assessing the Impact of MDE and Code Generation on the Sustainability of Scientific Computing Software: A Research Proposal

Spencer Smith

September 25, 2019

Abstract

...

Contents

1	Introduction	3
2	Background and Literature Review	6
2.1	Software Qualities of Interest	7
2.2	Desirable Qualities of Documentation	9
2.3	Literature Review on Current State of the Practice for SCS Development	10
2.4	Documentation for SCS	12
2.5	MDE and Code Generation	12
2.6	Empirical Methods for Software Engineering	14
3	Research Design and Methodologies	14
3.1	State of the Practice	14
3.1.1	Medical Imaging	18
3.1.2	Climate Modelling Software	18
3.1.3	Computational Medicine	18
3.1.4	Finite Element Analysis	18
3.1.5	Psychology Software	18

3.1.6	Chemical Engineering Software	18
3.1.7	Geographic Information Systems (GIS)	18
3.1.8	Physical Metallurgy Software	19
3.1.9	Quantum Chemistry Software	19
3.2	Impact of MDE on the Sustainability of SCS	19
3.2.1	Fitting MDE into the Scientific Software Development Process	19
3.2.2	Measuring the Qualities of MDE Developed Artifacts .	21
3.2.3	Case Study	21
3.2.4	Development of Small Physics Examples	22
3.3	Impact of Assurance Cases for Building Confidence	22
3.4	Impact of Computational Variability Testing	23
4	Implications and Contributions to Knowledge	24
5	Research Schedule	25
5.1	State of Practice Exercise	25
5.2	Impact of MDE on the Sustainability Timeline	25

1 Introduction

Without dramatic intervention, our collective confidence in Scientific Computing Software (SCS) is due for a catastrophic collapse. We are increasingly trusting ever more complex and ambitious computations, but our software foundations are built on sand. Areas of concern include nuclear safety analysis and computational medicine. Successfully building such software requires communication between software developers and experts from multiple domains. Collaboration is difficult at the best of times, and is made worse because developers favour handcrafted solutions over adapting software engineering processes, methods and tools (Faulk et al., 2009). Handcrafted solutions do not account for the inevitable changes in requirements, design and implementation. The twin challenges of changing requirements and inadequate documentation conspire to make computational results notoriously difficult to reproduce, especially in the important endeavour of one researcher independently replicating the results of another (Smith, 2018). Inadequate traceability within documentation means recertification of previously certified engineering software is as expensive and time consuming as the original certification exercise, because of a lack of explicit, effort reducing, traceability information. The existing challenges for modification, maintenance and extension point to problems with the sustainability of SCS.

Documentation can improve software sustainability. For software in general, documentation, written before and during development, provides many benefits Parnas (2010): easier reuse of old designs, better communication about requirements, more useful design reviews, easier integration of separately written modules, more effective code inspection, more effective testing, and more efficient corrections and improvements. For SCS in particular, documentation provides significant benefits. For instance, Smith et al. (2015) shows that the quality of statistical software for psychology is generally improved when developed using the structured CRAN (Comprehensive R Archive Network) process and tools, versus an ad hoc approach. Smith and Koothoor (2016) highlights the value of proper documentation by redeveloping nuclear safety analysis software. Twenty seven (27) worrisome documentation problems were found, including problems with incompleteness, ambiguity, inconsistency, verifiability, modifiability, traceability and a lack of abstraction. A redevelopment experiment with five existing projects (Smith et al., 2016a) enabled the code owners (as ascertained through interviews) to clearly see the value of documentation. However, mirroring other studies

(Carver et al., 2007), the code owners felt documentation takes too much time. Due to the time commitment and resource limitation, documentation is rarely emphasized in SCS.

One promising approach to gain the sustainability benefits of documentation, but significantly reduce the time commitment, is to use a Model Driven Engineering (MDE) process. MDE is an engineering approach that exploits structured representations (metamodels for abstract syntax, and constraints to capture static semantics) to automate repetitive, tedious and error prone tasks. In this research proposal, we will use MDE as a short-hand to include knowledge models, model transformation, Domain Specific Languages (DSLs) and code generation. Several current papers point to DSLs and code/document generation as a transformative technology for documentation, design and verification of SCS (Johanson and Hasselbring, 2018; Smith, 2018). MDE has the potential to dramatically reduce the cost of developing reliable, reproducible and re-certifiable software.

One example of applying MDE to SCS is Drasil, which is introduced in Szymczak et al. (2016) and available at <https://github.com/JacquesCarette/Drasil>. Drasil is implemented as an open source set of Domain Specific Languages (DSLs). The Drasil approach involves: i) creating infrastructure for a knowledge base of scientific and computing models; and, ii) writing explicit “recipes” that weave together this knowledge to generate theoretical models, design documents, code, test cases and build scripts. The generator can render to multiple languages, such as html, and LaTeX for documentation and Python, Java, C++ and Lua, for code. One source of knowledge, with rules for transformation, means completeness, consistency and traceability can be achieved by construction. Moreover, these qualities can be maintained as requirements are modified, design decisions are changed, documentation standards are varied, and software is recertified. Although creating the knowledge base is time consuming, the knowledge can be built up incrementally. New projects will reuse existing knowledge and expand as necessary.

As described above, sustainability is a concern for SCS; MDE holds promise for addressing this concern. However, a transformative change in the development of SCS requires more than promise. Empirical evidence is necessary to demonstrate that an MDE process can improve the sustainability of SCS. This leads to the long-term objective for this research proposal: *assess and measure the impact of MDE on the sustainability of SCS*.

The scope of SCS software for this proposal includes software intended for a long life (10+ years), as implied by the sustainability objective. The

effort involved in MDE does not provide enough of a payoff for short-term projects. To gain the full benefits of generating documentation, the scope also includes software that interests multiple stakeholders, with different interests and backgrounds. Finally, the scope of this proposal emphasizes safety related software, such as software for nuclear safety analysis, medical imaging and computational medicine. Safety related software is prioritized because documentation is often expected as part of a certification exercise. Although this study is not specifically restricted to open-source software, for practical reasons (especially the lack of access to the code), commercial software will not be strongly emphasized.

Implicit in the scope is that we are interested in software for domains that are well-understood. This is necessary for knowledge capture and the associated construction of models. We can refine the notion of well-understood by using the fact that MDE is simply a technology based realization of the older concept of program families. Success with an MDE approach requires satisfaction of the same three hypotheses as success for program family development (Weiss, 1997):

Redevelopment Hypothesis – most software development involved in producing an individual family members should be redevelopment because the family members have so much in common.

Oracle Hypothesis – the changes that are likely to occur during the software’s lifetime are predictable.

Organizational Hypothesis – designers can organize the software and the development effort so that predicted changes can be made independently.

SCS can be shown to satisfy these three hypotheses (Smith and Chen, 2004; Smith et al., 2007); therefore, SCS is suitable for a program family, or MDE software development approach. Although our focus is on SCS, an MDE approach could be applied to any domain that satisfies the above three hypotheses.

Given the large number of potential MDE processes, techniques and technologies, together with the wide variety in purpose, scope and context for SCS, the long-term objective of measuring the impact of MDE is decomposed into four short-term objectives. For each objective, a reference is given to the corresponding section where details are provided.

1. Assess the current state of the software development practice for SCS. How is the software developed in different sub-domains of SCS? Are there existing development methods that lead to higher software sustainability? (Section 3.1)
2. Assess the impact of an MDE process on end user developed SCS with respect to developer productivity and software sustainability for long-lived lived software. (Section 3.2)
3. Assess the impact on safety when the MDE generated documentation targets a safety assurance case. An assurance case is an explicit structured argument pertaining to specific properties, such as trustworthiness. The evidence for the argument will come from the traceability between artifacts generated by model transformations, expert reviews, test cases, etc. (Section 3.3)
4. Assess the impact on software quality when using MDE to facilitate Computational Variability Testing (CVT). CVT uses code generation to build confidence in the generated code in an analogous way to the use of grid refinement studies. Grid refinement looks at how the solution changes by varying the run-time parameter of grid density and comparing the results to theoretical expectations. CVT, on the other hand, can generate code to “refine” build time parameters, such as order of interpolation, or degree of implicitness. These parameters can be systematically varied and the results compared against the expected trend. (Section 3.4)

Before presenting the details of the short-term objectives (Section 3), the necessary background information is provided, along with a literature review (Section 2). The literature review covers software quality (Sections 2.1 and 2.2), the literature on the current state of the practice for SCS development (Section 2.3), potential target documentation artifacts for SCS (Section 2.4), MDE and code generation (Section 2.5) and empirical methods for software engineering (Section 2.6). This document ends with a summary of the potential contributions to knowledge from this study (Section 4) and a research schedule (Section 5).

2 Background and Literature Review

To assess the impact of MDE on SCS quality, we need a clear definition of what we mean by quality. Section 2.1 shows how the concept of quality is decomposed into a set of separate qualities. This set of qualities can be applied to the software artifacts (documentation, test cases, etc) and to the software development process itself. Several of the qualities from the list in Section 2.1 cannot be measured directly, such as maintainability. Therefore, Section 2.2 introduces measurable documentation qualities that are believed to contribute to the overall software qualities. In some cases it may be necessary to use the qualities in Section 2.2 to indirectly measure qualities listed in Section 2.1.

The other subsection in this background section provide an overview of the literature on the current approaches to developing SCS software 2.3, recommended document artifacts for SCS 2.4, an overview of MDE and Code Generation in an SCS context 2.5, and a summary of existing literature on empirical methods for software engineering 2.6.

2.1 Software Qualities of Interest

Our analysis is centred around a set of software qualities. Quality is not considered as a single measure, but a collection of different qualities, often called “ilities.” These qualities highlight the desirable nonfunctional properties for software artifacts, which include both documentation and code. Some qualities, such as visibility and productivity, apply to the process used for developing the software. The following list of qualities is based on Ghezzi et al. (2003). To the list from Ghezzi et al. (2003), we have added three qualities important for SC: installability, reproducibility and sustainability.

Installability A measure of the ease of installation.

Correctness Software is correct if it matches its specification.

Verifiability involves “solving the equations right” (Roache, 1998, p. 23); it benefits from rational documentation that systematically shows, with explicit traceability, how the governing equations are transformed into code.

Validatability means “solving the right equations” (Roache, 1998, p. 23). Validatability is improved by a rational process via clear documentation

of the theory and assumptions, along with an explicit statement of the systematic steps required for experimental validation.

Reliability is a critical quality for scientific software, since the results of computations are meaningless, if they are not dependable. Reliability is closely tied to verifiability, since the key quality to verify is reliability, while the act of verification itself improves reliability.

Performance considerations can make certification challenging, since QA becomes more difficult for more complex code. However, as Roache ([Roache, 1998](#), p. 355) points out, using simpler algorithms and reducing the number of options in general purpose code, is not always a practical option.

Usability can be a problem. Different users, solving the same physical problem, using the same software, can come up with different answers, due to differences in parameter selection ([Roache, 1998](#), p. 370). To reduce misuse, a rational process must state expected user characteristics, modelling assumptions, definitions and the range of applicability of the code.

Maintainability is necessary in scientific software, since change, through iteration, experimentation and exploration, is inevitable. Models of physical phenomena and numerical techniques necessarily evolve over time [Carver et al. \(2007\)](#); [Segal and Morris \(2008\)](#). Proper documentation, designed with change in mind, can greatly assist with change management.

Reusability provides support for the quality of reliability, since reliability is improved by reusing trusted components [Dubois \(2005\)](#). (Care must still be taken with reusing trusted components, since blind reuse in a new context can lead to errors, as dramatically shown in the Ariane 5 disaster ([Oliveira and Stewart, 2006](#), p. 37–38).) The odds of reuse are improved when it is considered right from the start.

Understandability is necessary, since reviewers can only certify something they understand. Scientific software developers have the view “that the science a developer embeds in the code must be apparent to another scientist, even ten years later” [Kelly \(2013\)](#). Understandability applies to the documentation and code, while usability refers to the executable

software. Documentation that follows a rational process is the easiest to follow.

Reproducibility is a required component of the scientific method [Davison \(2012\)](#). Although QA has, “a bad name among creative scientists and engineers” ([Roache, 1998](#), p. 352), the community need to recognize that participating in QA management also improves reproducibility. Reproducibility, like QA, benefits from a consistent and repeatable computing environment, version control and separating code from configuration/parameters [Davison \(2012\)](#).

Productivity [This needs to be filled in. Productivity is an important quality for this proposal, since part of what MDE promises is increased productivity. —SS]

Sustainability [This needs to be filled in. A search needs to be done to find a good definition. A starting point might be the recent paper that defines sustainability as a combination of other qualities. —SS]

2.2 Desirable Qualities of Documentation

To achieve the qualities listed in [Section 2.1](#), the documentation should achieve the qualities listed in this section. All but the final quality listed (abstraction), are adapted from the IEEE recommended practise for producing good software requirements [IEEE \(1998\)](#). Abstraction means only revealing relevant details, which in a requirements document means stating what is to be achieved, but remaining silent on how it is to be achieved. Abstraction is an important software development principle for dealing with complexity ([Ghezzi et al., 2003](#), p. 40). [Smith and Koothoor \(2016\)](#) present further details on the qualities of documentation for SCS.

Completeness Documentation is said to be complete when all the requirements of the software are detailed. That is, each goal, functionality, attribute, design constraint, value, data, model, symbol, term (with its unit of measurement if applicable), abbreviation, acronym, assumption and performance requirement of the software is defined. The software’s response to all classes of inputs, both valid and invalid and for both desired and undesired events, also needs to be specified.

Consistency Documentation is said to be consistent when no subset of individual statements are in conflict with each other. That is, a specification of an item made at one place in the document should not contradict the specification of the same item at another location.

Modifiability The documentation should be developed in such a way that it is easily modifiable so that likely future changes do not destroy the structure of the document. Also it should be easy to reflect the change, wherever needed in the document to maintain consistency, traceability and completeness. For documentation to be modifiable, its format must be structured in a way that repetition is avoided and cross-referencing is employed.

Traceability Documentation should be traceable, as this facilitates maintenance and review. If a change is made to the design or code of the software, then all the documentation relating to those segments have to be modified. This property is also important for recertification.

Unambiguity Documentation is said to be unambiguous only when every requirement's specification has a unique interpretation. The documentation should be unambiguous to all audiences, including developers, users and reviewers.

Correctness There is no direct tool or method for measuring correctness. One way of building confidence in correctness is by reviewing to ensure that each requirement stated is one that the stakeholders and experts desire. By maintaining traceability, consistency and unambiguity, we can reduce the occurrence of errors and make the goal of reviewing for correctness easier.

Verifiability Every requirement in the documentation must be the one fulfilled by the implemented software. Therefore all the requirements should be clear, unambiguous and testable, so that a person or a machine can verify whether the software product meets the requirements.

Abstract Documented requirements are said to be abstract if they state what the software must do and the properties it must possess, but do not speak about how these are to be achieved. For example, a requirement can specify that an Ordinary Differential Equation (ODE) must be solved, but it should not mention that Euler's method should

be used to solve the ODE. How to accomplish the requirement is a design decision, which is documented during the design phase.

2.3 Literature Review on Current State of the Practice for SCS Development

The SCS community is finally realizing that current practices are not sustainable. [Faulk et al. \(2009\)](#) observe, “growing concern about the reliability of scientific results based on ... software.” Embarrassing failures have occurred, like a retraction of derived molecular protein structures ([Miller, 2006](#)), false reproduction of sonoluminescent fusion ([Post and Votta, 2005](#)), and fixing and then reintroducing the same error in a large code base three times in 20 year ([Milewicz and Raybourn, 2018](#)). A recent report on directions for SCS research and education states: “While the volume and complexity of [SCS] have grown substantially in recent decades, [SCS] traditionally has not received the focused attention it so desperately needs ... to fulfill this key role as a cornerstone of long-term collaboration and scientific progress” ([Rüde et al., 2018](#)). Estimates suggest that the number of released faults per thousand executable lines of code during a given program’s life cycle is at best 0.1, and more likely 10 to 100 times worse ([Hatton, 2007](#)).

Although reproducibility is the cornerstone of the scientific method, until recently it has not been treated seriously in software ([Benureau and Rougier, 2017](#)). Fortunately, in recent years multiple conferences, workshops and individuals are calling for dramatic change ([Bailey et al., 2016](#)). The need for action is highlighted by a study of 402 computer systems papers - only 48.3% of the code was both available and compilable ([Collberg et al., 2015](#)). (Drasil addresses this problem because as programming languages evolve the code renderers in Drasil can be updated.) Reproducibility problems are even more extreme when the goal is replicability. A third party should be able to repeat a study using only the description of the methodology from a published article, with no access to the original code or computing environment ([Benureau and Rougier, 2017](#)). However, replicability is rarely achieved, as shown for microarray gene expression ([Ioannidis et al., 2009](#)) and for economics modelling ([Ionescu and Jansson, 2012](#)). Drasil addresses completeness and ambiguity problems, since it emphasizes capturing and documenting all of the required knowledge, including derivation of equations and rationales. [Crick et al. \(2014\)](#) point out potential roadblocks for

reproducibility, including page length constraints and differing detail needs depending on the audience. Again Drasil addresses these concerns because the recipes used to generate the documentation can be tailored to the level of detail required.

Although scientists recognize the seriousness of their SCS problems, their corrective steps are too incremental. For instance, a recent proposal for the future of High Energy Physics (HEP) software [Stewart et al. \(2017\)](#) uses words like sustainability, maintainability and reproducibility, but is almost completely silent on how these qualities are to be achieved. The proposal mentions developing new and improved algorithms (including parallel computing and machine learning), programming tools, recruitment and training, but there is little on documentation, design or verification techniques (other than unit testing). Similarly, a recent proposal on future directions for SCS research and education ([Rüde et al., 2018](#)) recognizes the desperate need for change, but then only suggests training on project management tools, open sharing and ethics. *Incremental change is not adequate; we need transformative change.*

Thankfully SCS leaders recognize that an interdisciplinary approach provides the path forward. They believe that the solution to SCS quality problems is applying, adapting and developing SE methods, tools and techniques. However, typical software processes are a barrier to progress. “To break the gridlock, we must establish a degree of cooperation and collaboration with the [SE] community that does not yet exist” ([Faulk et al., 2009](#)). “There is a need to improve the transfer of existing practices and tools from ... [SE] to scientific programming. In addition, ... there is a need for research to specifically develop methods and tools that are tailored to the domain” ([Storer, 2017](#)). This tailored research will require teams with a multidisciplinary background.

2.4 Documentation for SCS

Table 1 shows the recommended documentation for a scientific software project. The documents are typical of what is suggested for scientific software certification, where certification consists of official recognition by an authority, or regulatory body, that the software is fit for its intended use. For instance, the Canadian Standards Association (CSA) requires a similar set of documents for quality assurance of scientific programs for nuclear power plants [CSA \(1999\)](#).

Table 1: Recommended Documentation	
<i>Problem Statement</i>	Description of problem to be solved
<i>Development Plan</i>	Overview of development process/infrastructure
<i>Requirements</i>	Desired functions and qualities of the software
<i>V&V Plan</i>	Verification that all documentation artifacts, including the code, are internally correct. Validation, from an external viewpoint, that the right problem, or model, is being solved.
<i>Design Specification</i>	Documentation of how the requirements are to be realized, through both a software architecture and detailed design of modules and their interfaces
<i>Code</i>	Implementation of the design in code
<i>V&V Report</i>	Summary of the V&V efforts, including testing
<i>User Manual</i>	Instructions on installation, usage; worked examples

[Add figure showing the V-model of the documentation. —SS]

2.5 MDE and Code Generation

DSLs and code/document generation provide a transformative technology for documentation, design and verification (Johanson and Hasselbring, 2018; Smith, 2018). Scientists prefer focusing on science not software (Kelly, 2007). DSLs allow them to do this. A generative approach removes the maintenance nightmare of documentation duplicates and near duplicates (Luciv et al., 2018), since knowledge is only captured once and automatically transformed as needed. Code generation has previously been applied to improve SCS (Whaley et al., 2001; Veldhuizen, 1998; Püschel et al., 2001). For instance, ATLAS (Automatically Tuned Linear Algebra Software) (Whaley et al., 2001) and Blitz++ (Veldhuizen, 1998) produces efficient and portable linear algebra software. Spiral (Püschel et al., 2001) uses software/hardware generation for digital signal processing. Carette and Kiselyov (2011) shows how to generate a family of efficient, type-safe

Gaussian elimination algorithms. FEniCS (Finite Element and Computational Software) (Logg et al., 2012) uses code generation when solving differential equations. Unlike previous work on SCS code generation, Drasil (<https://github.com/JacquesCarette/Drasil>) focus on generating all software artifacts (requirements, design etc.), not just code.

Drasil, as presented in Szymczak et al. (2016) removes excuses for avoiding documentation by providing transformative SCS development technology. Drasil provides an infrastructure for knowledge capture and document/code generation. Drasil is implemented via Domain Specific Languages (DSLs) embedded in Haskell. Drasil generates requirements documentation and code for several case studies, including simulating the temperature of a solar water heating system, glass breakage and two-dimensional game physics. The full documentation (requirements, design etc), code and test cases have been created manually for each case study. These manual case studies (many from Smith et al. (2016a)) provide the ‘gold standard’ against which Drasil is tested. All of these artifacts are available publicly at <https://jacquescarette.github.io/Drasil/>. The generation techniques for Drasil began to take shape during work on generating geometric data structures and functions (Carette et al., 2011).

2.6 Empirical Methods for Software Engineering

The empirical methods in this study will follow the guidelines given in Kitchenham et al. (2002). Since case studies will be part of the methodology, the guidelines for conducting and reporting case study research in software engineering (Runeson and Höst, 2009) will be followed.

3 Research Design and Methodologies

This section is decomposed into 4 subsections – one for each of the different approaches for assessing the impact of MDE on the sustainability of SCS, as outlined in the Introduction (Section 1). The goal of the first subsection is to establish a “baseline” on the current state of the practice for SCS. This is done by identifying domains within SCS and then systematically assessing the quality of the work in those domains. The second subsection looks specifically at the impact of MDE on SCS through a series of case studies. An “ideal” MDE process is proposed and then followed with different stakeholders and

qualitative data is collected. The last two sub-sections discuss assessing the impact on SCS of assurance cases and computational variability testing, respectively.

3.1 State of the Practice

To understand the current quality of SCS and to assist in developing methods for measuring qualities, we will identify several domains within SCS and endeavour to summarize the current state of the practice within those domains. This “state of the practice” exercise will build off of prior work on measuring/assessing the state of software development practice in several SCS domains. We will update the work that was done previously for domains such as Geographic Information Systems ([Smith et al., 2018a](#)), Mesh Generators ([Smith et al., 2016b](#)), Seismology software ([Smith et al., 2018d](#)), and Statistical software for psychology ([Smith et al., 2018c](#)). We may return to these domains and/or introduce new domains. Potential domains/collaborators are listed in the subsections below.

MEng students will be asked to assess the state of practice in each domain. During the previous state of the practice exercise each student measured two domains, but with the higher demands from the new measurement protocol, we will likely switch to one domain per student.

In the previous project, we measured 30 software projects for each domain. With the increased scrutiny required in this re-boot, we will not likely be able to measure the details on this high a number of software projects. We may still start with 30 software examples in each domain, but then use some criteria to create a shorter sub-list for detailed measurement. The details of how to proceed here would be determined and documented in the measurement procedures. We still would like to keep the number of measured software projects above 10.

In the previous state of the practice assessing exercise ([Smith et al., 2018a](#)), a series of questions and simple metrics was devised to measure the quality of the documentation and adherence to best software development practices. The previous set of questions is provided [here](#). The new project will critically assess the previous set of questions and revise as necessary. In addition, the following data will also be considered for collection for each domain:

- Characterization of the functionality provided by the software in the

domain. What services does each member of the domain provide? Ideally this information would be summarized in a commonality analysis document. The document will summarize the domain software via its commonalities, variabilities and parameters of variation.

- Usability testing of each software in the domain. Specific metrics for assessing usability still need to be researched, but measures that involve the time to complete a task and the users overall impression will be considered.
- Collect empirical software engineering related data, such as the number of files, number of lines of code, cyclomatic complexity, number of open issues, number of closed issues, etc. As a starting point for tool support, HubListener (<https://github.com/pjmc-oliveira/HubListener>) could be used.
- Collect at least one empirical measure of the quality of the documentation - the number of lines of documentation. We should investigate to see if this is a metric that has previously been collected. So many of the existing metrics are completely focused on the code, or possibly on issue creation and management; it seems that little work has been done to assess the documentation quality.
- One of the guidelines for successfully on-boarding new developers is to “Keep knowledge up to date and findable” (Sholler et al., 2019). The selected measurements should be able to assist with judging whether a project succeeds on this front.

To focus the revised measurement procedure, we will go through each of the quality criteria listed in Section 2.1. For each criteria we will look at how it was previously measured by us, but also at the literature and current best practices for its measurement.

The above data will be combined to rank the software in the domain. The specifics are yet to be determined, but in the past the Analytic Hierarchy Process (AHP) (Saaty, 1980) proved helpful in this. In the past the AHP results were calculated by a grade out of 10 from each student evaluator for each quality for each measured product. A better approach may be to have a group of people familiar with the software do a pairwise comparison between each product to come up with a number. The number would be based on the

data collected, but not in a numerical way. The data collected would serve the same role as a decision support system. It would provide insight, but human beings would do the final pairwise comparisons, and thus the final ranking.

Key to the success of this exercise will be to involve and engage a partner for each state of practice project. In the previous effort we did not involve domain experts with the rationale of excluding potential bias, but this advantage is not worth not being able to evaluate the functionality/usability of the software. Moreover, not having an expert makes publication more difficult, since there is no one to advise on how best to approach journals and publishers. The domain expert will be asked to help in the following ways:

- Review the protocol for assessing/measuring quality. This same protocol will be used in each domain.
- Provide their expertise on potential publication. Specifically, they will recommend a suitable journal and act as the corresponding author for any paper submissions.
- Provide an authoritative list of domain software, possibly augmented by existing on-line lists.
- Provide some assistance from their own team of supervised students to facilitate software testing and domain characterization. Assessing the functionality and measuring usability will require an individual with domain knowledge. Multiple measurements will be necessary to have confidence. The measurements will likely take a few hours over a few days of time from each student volunteer.

There is no budget for this project, but the student volunteers will be considered co-authors in the resulting paper. Having them as co-authors also means that ethics approval should not be necessary.

GitHub will be used to coordinate the work of the large team of people that will be involved in this project. In addition to the project record left on GitHub, the final data will be exported to Mendeley (as was done for the previous quality measurement exercise, for instance the grades for GIS software are summarized at <https://data.mendeley.com/datasets/6kprpv7r7/1>.)

As a next step for this work, the procedure for measuring software qualities in SCS will be brainstormed, discussed and written up. The procedure

will be applied to each domain within SCS that is assessed. The write-up on the procedure will be published. Ideally at a relevant conference, but if that is not a feasible option, it will be posted on the arxiv.

Potential domains for future measurement are listed below. Determination of these domains are good candidates will require evaluation of the following criteria:

- Availability of a domain expert that is willing to assist, and who has at least a few graduate students that can assist with the measurement tasks.
- A domain where end-user developed software is common.
- A domain where open-source software is available. Commercial products are not as effective for measurement, since we cannot generally gain access to the source code, developer documentation, or issue tracking system.
- A domain where at least 10 software projects are available for measurement.
- A domain where scientific computing is used. The selected software should fit in the larger domain of SCS.

The data for each of the different subdomains within SCS will be standardized to facilitate a meta-analysis that combines all of the data.

3.1.1 Medical Imaging

Medical imaging and analysis software. This work could build off of the previous project on this topic completed by Vasudha Kapil. (Vasudha's report is located in the `mmsc` svn repo) .

3.1.2 Climate Modelling Software

Climate modelling software. This could be software used for multi-model ensemble aggregating ([Samouly et al., 2018](#)) or possibly with software that uses physical models to predict climate change.

3.1.3 Computational Medicine

Potential sub-domains with computational medicine could include: extraction of 3D geometries from medical images, or the use of Lattice Boltzmann Solvers.

3.1.4 Finite Element Analysis

Finite element solvers and/or mesh generators. Potential software includes FeNiCS, MOOSE, FreeFEM++, Dolfin, etc. Install the packages and use them to do something non-trivial.

3.1.5 Psychology Software

Statistics software for psychology. A project in this domain would allow us to update [Smith et al. \(2018c\)](#).

3.1.6 Chemical Engineering Software

Software is used for many purposes in chemical engineering. We would want to find an application where non-commercial software is available.

3.1.7 Geographic Information Systems (GIS)

We studied GIS software in the past ([Smith et al., 2018a](#)), so this would be an interesting area to update. The key for the updated effort will be to identify a domain expert available for collaboration.

3.1.8 Physical Metallurgy Software

Solidification and casting software and other software used for physical metallurgy. An example of software written for analyzing the solidification of a molten alloy is given in [Smith et al. \(2019\)](#).

3.1.9 Quantum Chemistry Software

Quantum chemistry software. One example is software for searching for conformers, like [Kaplan](#).

3.2 Impact of MDE on the Sustainability of SCS

Once we have details on the current state of the practice in SCS, we can begin to apply MDE to SCS projects and assess the ramifications of this change. This section starts by laying out the details on how MDE can be applied to SCS (Section 3.2.1). We plan to start with an ideal process and then approximate it through different experiments to measure and assess the impact of MDE. The second subsection (Section 3.2.2) outlines how the quality of MDE artifacts can be assessed. These measurements will be used in the various experiments to assess the impact of MDE. The third and fourth subsections describe two experiments for collecting data: i) conducting case studies with real SCS projects (Section 3.2.3); and, ii) performing an experiment where multiple developers create a small SCS project in one day (Section 3.2.4).

3.2.1 Fitting MDE into the Scientific Software Development Process

Many researchers have stated that a document driven process is not used by, nor suitable for, scientific software. These researchers argue that scientific developers naturally use an agile philosophy [Ackroyd et al. \(2008\)](#); [Carver et al. \(2007\)](#); [Easterbrook and Johns \(2009\)](#); [Segal \(2005\)](#), or an amethodical process [Kelly \(2013\)](#), or a knowledge acquisition driven process [Kelly \(2015\)](#). However, as discussed in the introduction to this proposal (Section 1, documentation provides many benefits. An MDE approach holds promise to provide these benefits, while removing the drudgery from the process of documentation and facilitating inevitable changes in the documentation and the software.

The general idea that MDE will be helpful is appealing, but to study the impact of MDE, we need to know how it will actually be used. How exactly would MDE be fit into the scientific software development process? The new ideal process should encourage an early intensity of thought and effort, but at the same time, it should not impose waterfall development, since waterfall is an unrealistic expectation for a problem that is not yet well understood. MDE supports frequent change and incremental development, but where should developers start and how should the work proceed? How does one create an almost empty project in MDE and then borrow from existing knowledge and create new knowledge? To create a specific MDE

process for SCS, the Drasil tools will be utilized. A good starting project could be [Projectile](#), which models simple projectile motion.

Compared to existing SCS development techniques, the new Drasil process should emphasize requirements and domain analysis earlier in the process. The quality of the resulting software depends on the quality of the requirements. As Boehm’s data shows [Boehm \(1981\)](#), early identification of errors implies a significant return on investment, since the cost of fixing errors increases dramatically for later development stages. An SRS is invaluable for verification, since the quality of software cannot be assessed without a standard against which to judge it. Better design decisions are facilitated by the information captured in an SRS. Moreover, the SRS aids the software lifecycle by facilitating incremental development. This last advantage is important for SCS software, since changes frequently occur as developers explore the problem domain. Unfortunately, developers do not generally write requirements. Quite possibly they do not like the intensity of thought that is necessary to get the requirements right. Inspection techniques, such as task based inspection ([Kelly and Shepard, 2000, 2004](#)) hold promise for getting the scientist to think deeply at an earlier stage in the development process.

The ideal process using the Drasil tool will focus on the ideal version of Drasil, since Drasil is still a work on progress. The design of the ideal process for the ideal Drasil, will help to guide future development of Drasil.

Since Drasil is not the only tool that can be used for MDE, we will also consider what an MDE approach would look like using other tools. As a starting point, we will look at [Epsilon](#). Epsilon is a family of languages and tools for code generation, model-to-model transformation, model validation, comparison, migration and refactoring that work out of the box with EMF, UML, Simulink, XML and other types of models.

3.2.2 Measuring the Qualities of MDE Developed Artifacts

Some of the quality measures developed in Section [3.1](#) for the “state of the practice” exercise will also be used to measure the quality of the generated documentation and code. Section [3.1](#) mentions using code quality metrics, like cyclomatic complexity, to attempt to quantify code quality. Given the importance of documentation, it seems reasonable to also assess its quality. Measuring documentation quality is not something that is often promoted. ([Sanga \(2003\)](#)) does discuss quantifying the connectivity within an SRS doc-

ument by analyzing the traceability graph.) Part of the reason for the lack of effort on quantifying the quality of documentation is that is not represented formally. With Drasil much of the documentation is formal. This should facilitate automatic measures of documentation quality, or at least metrics that might be correlated with quality.

To quantify the quality of the generated documentation, measures of information compression will be used. The measure will be based on the length of the (family of) generated documentation, code and test cases, compared to the length of the generator itself. Our measure will be based on the Minimum Description Length (MDL) Principle, which states that “the more we are able to compress the data, the more we have learned about the data” ([Grünwald, 2004](#)). Another important measure will be whether the generated documentation achieves reproducibility. The proposed measure is reproducibility depth ([Soergel, 2014](#)). A depth of 1 means that an independent researcher can start from the generated code and calculate the same results. Achieving deeper reproducibility will mean the same results can be calculated for an independent researcher starting with the design documentation. Each higher measure gets further from the source code, until the highest depth (reproducibility), which starts from the requirements and leads to identical (within acceptable error) results.

3.2.3 Case Study

For assessing the impact of MDE, experiments where a given project is developed in parallel by two independent teams are infeasible due to expense and challenges for having a truly controlled experiment. Instead, we will focus on the more practical, but still very useful, case study approach. Guidelines for conducting and reporting case study research in software engineering are given in [Runeson and Höst \(2009\)](#). As explained there, a case study should be built around a research question. The research question should be specific and measurable. For the current case study exercise, we propose the question: “What tools and methods from MDE are perceived by end user developers to improve their productivity and the sustainability of their software in the long term?”

We will identify multiple case studies, with potential projects coming from the domains listed for the “state of the practice” exercise (Section 3.1), like computational medicine, medical imaging or quantum chemistry software. We will first identify the partners and look for feasible projects that can be

developed using an MDE approach. We will use the ideal process identified in Section 3.2.1. The respective advantages and disadvantages of Drasil versus Epsilon will be compared.

To assess the impact of MDE, we will use techniques like pre and post interviews, the measurements of quality used for the “state of the practice” exercise, ethnographic study, etc.

[More information is needed in this section. —SS]

3.2.4 Development of Small Physics Examples

To learn more about MDE for SCS, and to assess the current and future state of Drasil, a workshop could be conducted where the participants are all expected to build a small physics application, together with full documentation and test cases, in a one day, or possibly half day workshop. The participants would be surveyed before and after the development exercise, and observations would be made throughout the exercise.

[More detail is needed. —SS]

3.3 Impact of Assurance Cases for Building Confidence

Sketching out an Assurance Case in advance of developing the software could be just what is needed to manufacture the up-front intensity that is needed for successful projects.

Experimentally Determine Appropriate Documentation Template: The best approach for documenting requirements and design is an open question. Drasil could be used to attempt to address this question. Different recipes for documentation could be developed and compared. Some comparison could be automated between the documents, such as empirical measures of their quality. Other measures would require some kind of usability experiments.

Some SCS needs to be verifiably safe. We will achieve this by building a safety assurance case. An assurance case is an explicit structured argument pertaining to specific properties, such as trustworthiness

Assurance cases have potential to be used in SC to improve our collective confidence in the developed software [Smith et al. \(2018b\)](#). Further work is necessary to demonstrate the value of assurance cases, but they certainly hold promise. Drasil could go a long way to supporting development of assurance cases. Some potential Drasil related support ideas include the following:

- Automatic generation of visual depiction of the assurance case from the goals, claims, evidence etc.
- Automated reporting of missing evidence, incomplete arguments, etc.
- Coordination of requirements, design decisions, etc. Assurance cases use the same information as recommended in “document driven” design, but the traceability and connection between the data is even more important. A tool like Drasil could go a long way in facilitating the presentation and navigation of an assurance case.

Assurance cases in Drasil could potentially go beyond the scientific context. For instance, they could also incorporate safety cases for control systems. This could potentially be of interest to McSCert.

3.4 Impact of Computational Variability Testing

This activity supports objective S3. The applicant’s students will apply LSS and code generation to Computational Variability Testing (CVT), using the example of a family of finite element programs. As proposed by the applicant, CVT uses code generation to build confidence in the generated code in an analogous way to the use of grid refinement studies. Grid refinement looks at how the solution changes by varying the run-time parameter of grid density and comparing the results to theoretical expectations. CVT, on the other hand, can generate code to “refine” build time parameters, such as order of interpolation, or degree of implicitness. These parameters can be systematically varied and the results compared against the expected trend. Varying computational variabilities is not a new idea, but it is only recently that the option has become available of coupling this with code generation, so that the validity of an entire program family can be assessed at the same time. The experimentation with CVT will initially use the FEniCS package ([Logg et al., 2012](#)) for elliptical partial differential equations to compare elements with different orders of interpolation to verify that the convergence to the solution agrees with the theoretical expectations.

Design and Documentation For Family of Data Fitting Algorithms: The fitting used in GlassBR and in SFS (Software for Solidification) could be made much more generic. We could have a family of fitting algorithms that could be used in any situation where fitting is required. A proper commonality analysis of this domain could potentially show the potential design

decisions that bridge between the requirements and the design. In the SFS example many different fitting routines were tried. If the experiments could have been done easily via a declarative specification, considerable time would have been saved. If the experiments are combined with automated testing and “properties of a correct solution” the human involvement could be reduced, so that we have partially automated algorithm selection.

4 Implications and Contributions to Knowledge

We are currently putting too much trust in SCS. Examples such as nuclear safety analysis and computational medicine show the significance of SCS for health and safety. Although SCS developers do excellent work, we do not currently have enough checks and balances in place for full confidence. More should be expected for documentation, design and testing. Fortunately, SE offers techniques, like DSLs and code generation, that can be employed to not only improve SCS quality, but also reduce development costs. Now is the time for change. SCS practitioners have recognized problems with the status quo and multidisciplinary researchers are needed to bridge the gap between SE and SCS.

Codifying medical, biomechanical and computational knowledge is challenging, but success will completely transform the development of safety-related software in medicine, science and engineering. MDE will remove the drudgery of generating and maintaining documentation, code, test cases and build environments, enabling scientists to focus on science, engineers to focus on engineering and medical professionals to focus on medicine. MDE will detect inconsistencies between models via inter-model consistency constraints. As a consequence, mistakes like inconsistent units or assumptions cannot occur. MDE will raise the bar so that we can expect an explicit argument for safety, not just code that we are expected to blindly trust. With the right up-front investment of knowledge capture, we can have software that is long lived because the stable knowledge is separated from the rapidly changing modelling assumptions and design decisions.

5 Research Schedule

[Schedules are incomplete —SS]

5.1 State of Practice Exercise

- Fall 2019
 - Review previous measurement template
 - Update measurement process
 - Meet with stakeholders
 - Identify list of software to be measured
- Winter 2020
 - Complete measurements
 - Use AHP to rank software products for each quality, and for overall quality
- Summer 2020
 - Write up results, MEng reports, papers

5.2 Impact of MDE on the Sustainability Timeline

A high-level timeline of the research plan for assessing the impact of MDE on sustainability (see Section 3.2) follows:

- Fall 2019
 - Write research protocol
- Winter 2020
 -
- Summer 2020
 -
- Fall 2020

- Winter 2021
- Summer 2021
- Fall 2021

References

- Karen S. Ackroyd, Steve H. Kinder, Geoff R. Mant, Mike C. Miller, Christine A. Ramsdale, and Paul C. Stephenson. Scientific software development at a research facility. *IEEE Software*, 25(4):44–51, July/August 2008.
- David H. Bailey, Jonathan M. Borwein, and Victoria Stodden. *Reproducibility: Principles, Problems, Practices*, chapter Facilitating reproducibility in scientific computing: principles and practice, pages 205–232. John Wiley and Sons, New York, 2016.
- F. Benureau and N. Rougier. Re-run, Repeat, Reproduce, Reuse, Replicate: Transforming Code into Scientific Contributions. *ArXiv e-prints*, August 2017.
- Barry Boehm. *Software Engineering Economics*. Prentice Hall, Englewood Cliffs, New Jersey, 1981.
- Jacques Carette and Oleg Kiselyov. Multi-stage programming with Functors and Monads: Eliminating abstraction overhead from generic code. *Sci. Comput. Program.*, 76(5):349–375, 2011.
- Jacques Carette, Mustafa ElSheikh, and W. Spencer Smith. A generative geometric kernel. In *ACM SIGPLAN 2011 Workshop on Partial Evaluation and Program Manipulation (PEPM’11)*, pages 53–62, January 2011.
- Jeffrey C. Carver, Richard P. Kendall, Susan E. Squires, and Douglass E. Post. Software development environments for scientific and engineering software: A series of case studies. In *ICSE ’07: Proceedings of the 29th International Conference on Software Engineering*, pages 550–559, Washington, DC, USA, 2007. IEEE Computer Society. ISBN 0-7695-2828-7. doi: <http://dx.doi.org/10.1109/ICSE.2007.77>.
- Christian Collberg, Todd Proebsting, and Alex M Warren. Repeatability and benefaction in computer systems research. Technical Report TR 14-04, Department of Computer Science, University of Arizona, Tucson, AZ, 2015. URL <http://repeatability.cs.arizona.edu/v2/RepeatabilityTR.pdf>.

- Tom Crick, Benjamin A. Hall, and Samin Ishtiaq. “Can I implement your algorithm?”: A model for reproducible research software. *CoRR*, abs/1407.5981, 2014. URL <http://arxiv.org/abs/1407.5981>.
- CSA. Quality assurance of analytical, scientific, and design computer programs for nuclear power plants. Technical Report N286.7-99, Canadian Standards Association, 178 Rexdale Blvd. Etobicoke, Ontario, Canada M9W 1R3, 1999.
- A. P. Davison. Automated capture of experiment context for easier reproducibility in computational research. *Computing in Science & Engineering*, 14(4):48–56, July-Aug 2012.
- P. F. Dubois. Maintaining correctness in scientific programs. *Computing in Science & Engineering*, 7(3):80–85, May-June 2005. doi: 10.1109/MCSE.2005.54.
- Steve M. Easterbrook and Timothy C. Johns. Engineering the software for understanding climate change. *Computing in Science & Engineering*, 11(6): 65–74, November/December 2009. ISSN 0740-7475. doi: <http://dx.doi.org/10.1109/MCSE.2009.193>.
- S. Faulk, E. Loh, M. L. V. D. Vanter, S. Squires, and L. G. Votta. Scientific computing’s productivity gridlock: How software engineering can help. *Computing in Science Engineering*, 11(6):30–39, Nov 2009. ISSN 1521-9615. doi: 10.1109/MCSE.2009.205.
- Carlo Ghezzi, Mehdi Jazayeri, and Dino Mandrioli. *Fundamentals of Software Engineering*. Prentice Hall, Upper Saddle River, NJ, USA, 2nd edition, 2003.
- Peter Grünwald. A tutorial introduction to the minimum description length principle. *ArXiv e-prints*, June 2004.
- Les Hatton. The chimera of software quality. *Computer*, 40(8), August 2007.
- IEEE. Recommended practice for software requirements specifications. *IEEE Std 830-1998*, pages 1–40, October 1998. doi: 10.1109/IEEESTD.1998.88286.

- John PA Ioannidis, David B Allison, Catherine A Ball, Issa Coulibaly, Xiangqin Cui, Aedín C Culhane, Mario Falchi, Cesare Furlanello, Laurence Game, Giuseppe Jurman, et al. Repeatability of published microarray gene expression analyses. *Nature genetics*, 41(2):149–155, 2009.
- Cezar Ionescu and Patrik Jansson. Dependently-Typed Programming in Scientific Computing — Examples from Economic Modelling. In *Revised Selected Papers of the 24th International Symposium on Implementation and Application of Functional Languages*, volume 8241 of *Lecture Notes in Computer Science*, pages 140–156. Springer International Publishing, 2012. doi: 10.1007/978-3-642-41582-1_9.
- Arne N. Johanson and Wilhelm Hasselbring. Software engineering for computational science: Past, present, future. *Computing in Science & Engineering*, Accepted:1–31, 2018.
- Diane Kelly. Industrial scientific software: A set of interviews on software development. In *Proceedings of the 2013 Conference of the Center for Advanced Studies on Collaborative Research*, CASCON '13, pages 299–310, Riverton, NJ, USA, 2013. IBM Corp. URL <http://dl.acm.org/citation.cfm?id=2555523.2555555>.
- Diane Kelly. Scientific software development viewed as knowledge acquisition: Towards understanding the development of risk-averse scientific software. *Journal of Systems and Software*, 109:50–61, 2015. doi: 10.1016/j.jss.2015.07.027. URL <http://dx.doi.org/10.1016/j.jss.2015.07.027>.
- Diane Kelly and Terry Shepard. Task-directed software inspection technique: an experiment and case study. In *CASCON '00: Proceedings of the 2000 conference of the Centre for Advanced Studies on Collaborative research*, page 6. IBM Press, 2000. URL <http://portal.acm.org/citation.cfm?id=782040#>.
- Diane Kelly and Terry Shepard. Eight maxims for software inspectors. *Software Testing, Verification and Reliability*, 14(4):243–256, 2004.
- Diane F. Kelly. A software chasm: Software engineering and scientific computing. *IEEE Software*, 24(6):120–119, 2007. ISSN 0740-7459. doi: <http://dx.doi.org/10.1109/MS.2007.155>.

- Barbara A. Kitchenham, Shari Lawrence Pfleeger, Lesley M. Pickard, Peter W. Jones, David C. Hoaglin, Khaled El Emam, and Jarrett Rosenberg. Preliminary guidelines for empirical research in software engineering. *IEEE Trans. Softw. Eng.*, 28(8):721–734, August 2002. ISSN 0098-5589. doi: 10.1109/TSE.2002.1027796. URL <https://doi.org/10.1109/TSE.2002.1027796>.
- A. Logg, K.-A. Mardal, and G. N. Wells, editors. *Automated Solution of Differential Equations by the Finite Element Method*, volume 84 of *Lecture Notes in Computational Science and Engineering*. Springer, 2012. doi: 10.1007/978-3-642-23099-8. URL <http://dx.doi.org/10.1007/978-3-642-23099-8>.
- D. V. Luciv, D. V. Koznov, G. A. Chernishev, A. N. Terekhov, K. Yu. Romanovsky, and D. A. Grigoriev. Detecting near duplicates in software documentation. *Programming and Computer Software*, 44(5):335–343, Sep 2018. ISSN 1608-3261. doi: 10.1134/S0361768818050079. URL <https://doi.org/10.1134/S0361768818050079>.
- Reed Milewicz and Elaine M. Raybourn. Talk to me: A case study on coordinating expertise in large-scale scientific software projects. *ArXiv e-prints*, September 2018.
- Greg Miller. SCIENTIFIC PUBLISHING: A Scientist’s Nightmare: Software Problem Leads to Five Retractions. *Science*, 314(5807):1856–1857, 2006. doi: 10.1126/science.314.5807.1856. URL <http://www.sciencemag.org>.
- Suely Oliveira and David E. Stewart. *Writing Scientific Software: A Guide to Good Style*. Cambridge University Press, New York, NY, USA, 2006. ISBN 0521858968.
- David Lorge Parnas. Precise documentation: The key to better software. In *The Future of Software Engineering*, pages 125–148, 2010. doi: 10.1007/978-3-642-15187-3_8. URL http://dx.doi.org/10.1007/978-3-642-15187-3_8.
- D. E. Post and L. G. Votta. Computational Science Demands a New Paradigm. *Physics Today*, 58(1):35–41, January 2005. doi: 10.1063/1.1881898.

- Markus Püschel, Bryan Singer, Manuela Veloso, and José M. F. Moura. Fast automatic generation of DSP algorithms. In *International Conference on Computational Science (ICCS)*, volume 2073 of *Lecture Notes In Computer Science*, pages 97–106. Springer, 2001.
- Patrick J. Roache. *Verification and Validation in Computational Science and Engineering*. Hermosa Publishers, Albuquerque, New Mexico, 1998.
- U. Rüde, K. Willcox, L. McInnes, and H. Sterck. Research and education in computational science and engineering. *SIAM Review*, 60(3):707–754, 2018. doi: 10.1137/16M1096840. URL <https://doi.org/10.1137/16M1096840>.
- Per Runeson and Martin Höst. Guidelines for conducting and reporting case study research in software engineering. *Empirical Software Engineering*, 14(2):131–164, Dec 2009. ISSN 1573-7616. doi: 10.1007/s10664-008-9102-8. URL <https://doi.org/10.1007/s10664-008-9102-8>.
- T. L. Saaty. *The Analytic Hierarchy Process: Planning, Priority Setting, Resource Allocation*. McGraw-Hill Publishing Company, New York, New York, 1980.
- Aly Al Samouly, Chanh Nien Luong, Zhong Li, Spencer Smith, Brian Baetz, and Maysara Ghaith. Performance of a multi-model ensemble for the simulation of temperature variability over Ontario, Canada. *Environmental Earth Sciences*, 77(524):1–12, 2018. doi: <https://doi.org/10.1007/s12665-018-7701-2>.
- B. Sanga. Assessing and improving the quality of software requirements specification documents (SRSDs). Thesis for M.Sc. program, Computing and Software Department, McMaster University, Hamilton, ON, August 2003.
- Judith Segal. When software engineers met research scientists: A case study. *Empirical Software Engineering*, 10(4):517–536, October 2005. ISSN 1382-3256. doi: 10.1007/s10664-005-3865-y. URL <http://dx.doi.org/10.1007/s10664-005-3865-y>.
- Judith Segal and Chris Morris. Developing scientific software. *IEEE Software*, 25(4):18–20, July/August 2008.

- Dan Sholler, Igor Steinmacher, Denae Ford, Mara Averick, Mike Hoye, and Greg Wilson. Ten simple rules for helping newcomers become contributors to open projects. *PLOS Computational Biology*, 15(9):1–10, 09 2019. doi: 10.1371/journal.pcbi.1007296. URL <https://doi.org/10.1371/journal.pcbi.1007296>.
- Spencer Smith, Malavika Srinivasan, and Sumanth Shankar. Debunking the myth that upfront requirements are infeasible for scientific computing software debunking the myth that upfront requirements are infeasible for scientific computing software. <https://arxiv.org/abs/1906.07812>, June 2019.
- W. Spencer Smith. Beyond software carpentry. In *2018 International Workshop on Software Engineering for Science (held in conjunction with ICSE’18)*, pages 1–8, 2018.
- W. Spencer Smith and Chien-Hsien Chen. Commonality and requirements analysis for mesh generating software. In F. Maurer and G. Ruhe, editors, *Proceedings of the Sixteenth International Conference on Software Engineering and Knowledge Engineering (SEKE 2004)*, pages 384–387, Banff, Alberta, 2004.
- W. Spencer Smith and Nirmitha Koothoor. A document-driven method for certifying scientific computing software for use in nuclear safety analysis. *Nuclear Engineering and Technology*, 48(2):404–418, April 2016. ISSN 1738-5733. doi: <http://dx.doi.org/10.1016/j.net.2015.11.008>. URL <http://www.sciencedirect.com/science/article/pii/S1738573315002582>.
- W. Spencer Smith, John McCutchan, and Fang Cao. Program families in scientific computing. In Jonathan Sprinkle, Jeff Gray, Matti Rossi, and Juha-Pekka Tolvanen, editors, *7th OOPSLA Workshop on Domain Specific Modelling (DSM’07)*, pages 39–47, Montréal, Québec, October 2007.
- W. Spencer Smith, Yue Sun, and Jacques Carette. Comparing psychometrics software development between CRAN and other communities. Technical Report CAS-15-01-SS, McMaster University, January 2015. 43 pp.
- W. Spencer Smith, Thulasi Jegatheesan, and Diane F. Kelly. Advantages, disadvantages and misunderstandings about document driven design for

- scientific software. In *Proceedings of the Fourth International Workshop on Software Engineering for High Performance Computing in Computational Science and Engineering (SE-HPCCE)*, November 2016a. 8 pp.
- W. Spencer Smith, Adam Lazzarato, and Jacques Carette. State of practice for mesh generation software. *Advances in Engineering Software*, 100:53–71, October 2016b.
- W. Spencer Smith, Adam Lazzarato, and Jacques Carette. State of the practice for GIS software. <https://arxiv.org/abs/1802.03422>, February 2018a.
- W. Spencer Smith, Mojdeh Sayari Nejad, and Alan Wassyn. Assurance cases for scientific computing software (poster). In *ICSE 2018 Proceedings of the 40th International Conference on Software Engineering*, May 2018b. 2 pp.
- W. Spencer Smith, Yue Sun, and Jacques Carette. Statistical software for psychology: Comparing development practices between CRAN and other communities. <https://arxiv.org/abs/1802.07362>, 2018c. 33 pp.
- W. Spencer Smith, Zheng Zeng, and Jacques Carette. Seismology software: State of the practice. *Journal of Seismology*, 22(3):755–788, May 2018d.
- David Soergel. Confirmation depth as a measure of reproducible scientific research. <http://davidsoergel.com/posts/confirmation-depth-as-a-measure-of-reproducible-scientific-research>, October 2014.
- Graeme Stewart et al. A Roadmap for HEP Software and Computing R&D for the 2020s. *arXiv*, 2017.
- Tim Storer. Bridging the chasm: A survey of software engineering practice in scientific programming. *ACM Comput. Surv.*, 50(4):47:1–47:32, August 2017. ISSN 0360-0300. doi: 10.1145/3084225. URL <http://doi.acm.org/10.1145/3084225>.
- Daniel Szymczak, W. Spencer Smith, and Jacques Carette. Position paper: A knowledge-based approach to scientific software development. In *Proceedings of SE4Science’16 in conjunction with the International Conference on*

- Software Engineering (ICSE)*, Austin, Texas, United States, May 2016. In conjunction with ICSE 2016. 4 pp.
- Todd. L. Veldhuizen. Arrays in Blitz++. In *Proceedings of the 2nd International Scientific Computing in Object-Oriented Parallel Environments (ISCOPE'98), Lecture Notes in Computer Science*. Springer-Verlag, 1998.
- David M. Weiss. Defining families: The commonality analysis. *Submitted to IEEE Transactions on Software Engineering*, 1997. URL <http://www.research.avayalabs.com/user/weiss/Publications.html>.
- R. C. Whaley, A. Petitet, and J. J. Dongarra. Automated empirical optimization of software and the ATLAS project. *Parallel Computing*, 27(1–2): 3–35, 2001.