

When Capturing Knowledge Improves Productivity

Jacques Carette
Computing and Software
McMaster University
Hamilton, Ontario, Canada
carrette@mcmaster.ca

Spencer Smith
Computing and Software
McMaster University
Hamilton, Ontario, Canada
smiths@mcmaster.ca

Jason Balaci
Computing and Software
McMaster University
Hamilton, Ontario, Canada
balacij@mcmaster.ca

ABSTRACT

Current software development is often quite code-centric and aimed at short-term deliverables, due to various contextual forces. We're interested in contexts where different forces are at play. **Well understood domains** and **long-lived software** provide such an opportunity. By applying generative techniques, aggressive knowledge capture has the real potential to greatly increase long-term productivity.

Key is to recognize that currently hand-written software artifacts contain considerable knowledge duplication. With proper tooling and appropriate codification of domain knowledge, greatly increasing productivity is feasible. We present an example of what this looks like, and just some of the benefits (reuse, traceability, change management) thus gained.

CCS CONCEPTS

• **Software and its engineering** → *Application specific development environments; Requirements analysis; Specification languages; Automatic programming.*

KEYWORDS

code generation, document generation, knowledge capture, software engineering

ACM Reference Format:

Jacques Carette, Spencer Smith, and Jason Balaci. 2022. When Capturing Knowledge Improves Productivity. In *Pittsburgh '22: NIER - New Ideas and Emerging Results (ICSE 2022), May 21–29, 2022, Pittsburgh, PA, USA*. ACM, New York, NY, USA, 5 pages. <https://doi.org/10.1145/nnnnnnn.nnnnnnn>

1 THE CONTEXT

1.1 “Well understood” software?

DEFINITION 1. A software domain is well understood if

- (1) its Domain Knowledge (DK) is codified,
- (2) the computational interpretation of the DK is clear, and
- (3) writing code to perform said computations is well understood.

By *codified*, we mean that the knowledge exists in standard form in a variety of textbooks. For example, many engineering domains

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

Pittsburgh '22, May 21–29, 2022, Pittsburgh, PA, USA

© 2022 Association for Computing Machinery.

ACM ISBN 978-x-xxxx-xxxx-x/YY/MM...\$15.00

<https://doi.org/10.1145/nnnnnnn.nnnnnnn>

use ordinary differential equations as models, the quantities of interest are known, given standard names and standard units. In other words, standard vocabulary has been established over time and the body of knowledge is uncontroversial.

We can refine these high level ideas, using the same numbering, although the refinement should be understood more holistically.

- (1) Models in the DK *can be* written formally.
- (2) Models in the DK *can be* turned into functional relations by existing mathematical steps.
- (3) Turning these functional relations into code is an understood transformation.

Most importantly, the last two parts deeply involve *choices*: What quantities are considered inputs, outputs and parameters to make the model functional? What programming language? What software architecture data-structures, algorithms, etc.?

In other words, *well understood* does not imply *choice free*. Writing a small script to move files could just as easily be done in Bash, Python or Haskell. In all cases, assuming fluency, the author's job is straightforward because the domain is well understood.

1.2 Long-lived software?

For us, long-lived software is software that is expected to be in continuous use and evolution for 20 or more years. The main characteristic of such software is the *expected turnover* of key staff. This means that all tacit knowledge about the software will be lost over time if it is not captured.

1.3 Productivity?

We adapt the standard definition of productivity [4], where inputs are labour, but adjust the outputs to be knowledge and user satisfaction, where user satisfaction acts as a proxy for effective quality.

This explicit emphasis on all knowledge produced, rather than just the operationalizable knowledge (aka code) implies that human-reusable knowledge, i.e. documentation, is crucial. This is why the *long-lived* context is important.

1.4 Documentation

Our definition of well understood also applies to **documentation** aimed at humans. Explicitly:

- (1) The meaning of the models is understood at a human-pedagogical level, i.e. it is explainable.
- (2) Combining models is explainable. Thus *transformers* simultaneously operate on mathematical representations and on explanations. This requires that English descriptions also be captured in the same manner as the formal-mathematical knowledge.

- (3) Similarly, the *transformers* that arise from making software oriented decisions should be captured with a similar mechanism, and also include English explanations.

We dub these *triform theories*, as a nod to *biform theories* [6]. We couple (1) an axiomatic description, (2) a computational description, and (3) an English description of a concept.

1.5 Software artifacts

Software currently consists of a whole host of artifacts: requirements, specifications, user manual, unit tests, system tests, usability tests, build scripts, READMEs, license documents, process documents, as well as code.

Whenever appropriate, we use standards and templates for each of the generated artifacts. For requirements, we use a variant [10] of the IEEE [7] and Volere templates [9].

1.6 Examples of context

When are these conditions fulfilled? One example is *research software* in science and engineering. While the results of running various simulations is entirely new, the underlying models and how to simulate them are indeed well-known. One particularly long-lived example is embedded software for space probes (like Pioneer 10).

2 A NEW DEVELOPMENT PROCESS

Given appropriate infrastructure, what would be an *idealized process* (akin to Parnas' ideas of faking a rational design process [8])?

- (1) Have a task to achieve where *software* can play a central part in the solution.
- (2) The underlying problem domain is *well understood*.
- (3) Describe the problem:
 - (a) Find the base knowledge (theory) in the pre-existing library or, failing that, write it if it does not yet exist, for instance the naturally occurring known quantities and associated constraints.
 - (b) Assemble the ingredients into a coherent narrative,
 - (c) Describe the characteristics of a good solution,
 - (d) Come up with basic examples (to test correctness, intuitions, etc).
- (4) Describe, by successive refinement transformations, how the above can be turned into a deterministic¹ input-output process.
 - (a) Some refinements will involve *specialization* (eg. from n -dimensional to 2-dimensional, assuming no friction, etc). These *choices* and their *rationale* need to be documented, as a crucial part of the solution. Whether these choices are (un)likely to change in the future should be recorded.
 - (b) Choices tend to be dependent, and thus (partially) ordered. *Decisions* frequently enable or reveal downstream choices.
- (5) Describe how the process from step 4 can be turned into code. The same kinds of choice can occur here.
- (6) Turn the steps (i.e. from items 4 and 5) into a *recipe*, aka program, that weaves together all the information into a

¹A current meta-design choice.

variety of artifacts (documentation, code, build scripts, test cases, etc). These can be read, or executed, or ... as appropriate.

While this last step might appear somewhat magical, it isn't. The whole point of defining *well understood* is to enable it! A *suitable* knowledge encoding is key. This is usually tacit knowledge that entirely resides in developers' heads.

What is missing is an explicit *information architecture* of each of the necessary artifact. In other words, what information is necessary to enable the mechanized generation of each artifact? It turns out that many of them are quite straightforward.

Often steps 1 and 3 are skipped; this is part of the **tacit knowledge** of a lot of software. Our process requires that this knowledge be made explicit, a fundamental step in *Knowledge Management* [5].

3 AN EXAMPLE

We have built the needed infrastructure. It consists of 60KLoc of Haskell implementing a series of interacting Domain Specific Languages (DSLs) for knowledge encodings, mathematical expressions, theories, English fragments, code generation and document generation.² A full description would take too much space. Instead, we provide an illustrative example.

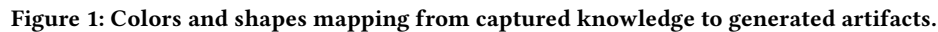
We will focus on information capture and the artifacts we can generate. For concreteness, we'll use a single example from our suite: GlassBR, used to assess the risk for glass facades subject to blast loading. The requirements are based on an American Standard Test Method (ASTM) standard [1–3]. GlassBR was originally a Visual Basic code/spreadsheet created by colleagues in a civil engineering research lab. We added their domain knowledge to our framework, along with recipes to generate relevant artifacts. Not only can we generate code for the necessary calculations (in C++, C#, Java, Python and Swift), we added documentation that was absent in the original (Software Requirements Specification, doxygen, README.md and a Makefile). Moreover, our implementation is actually a family of implementations, since some design decisions are explicitly exposed as changeable variabilities, as described below.

The transformation of captured knowledge is illustrated in Figure 1. This is read starting from the upper right box. Each piece of information in this figure has its own shape and colour (orange-cloud, pink lozenge, etc). It should be immediately clear that all pieces of information reappear in multiple places in the generated artifacts. For example, the name of the software (GlassBR) ends up appearing more than 80 times in the generated artifacts (in the folder structure, requirements, README, Makefile and source code). Changing this name would traditionally be extremely difficult; we can achieve this by modifying a single place, and regenerating.

The first box shows the directory structure of the currently generated artifacts; continuing clockwise, we see examples of Makefiles for the Java and Python versions, parts of the fully documented, generated code for the main computation in those languages, user instructions for running the code, and the processed \LaTeX for the requirements.

The name GlassBR is probably the simplest example of what we mean by *knowledge*: here, the concept “program name” is internally

²We will provide a link if the paper is accepted.



Step 2: Understood The details are extensively documented in [1–3].

Step 3a: Base Knowledge A recurring idea is the different types of Glass:

Concept	Term (Name)	Abbrev.	Domain
fullyT	Fully Tempered	FT	[Glass]
heatS	Heat Strengthened	HS	[Glass]
iGlass	Insulating Glass	IG	[Glass]
lGlass	Laminated Glass	LG	[Glass]
glassTypeFac	Glass Type Factor	GTF	[Glass]

The “Risk of Failure” is definable, using the following template for a *data definition*:

Label	Risk of Failure
Equation	$B = \frac{k}{(ab)^{m-1}} (Eh^2)^m LDFe^J$
Description	<p>B is the Risk of Failure (Unitless)</p> <p>k is the surface flaw parameter ($\frac{m^{12}}{N^7}$)</p> <p>a & b are the plate length & width (m)</p> <p>...</p>
Source	[1], [3]

Step 3b: Coherent narrative The descriptions in GlassBR are produced using an experimental language using specialized markup for describing relations between knowledge. For example, the goal of GlassBR (“Predict-Glass-Withstands-Explosion”) is to “Analyze and predict whether the *glass slab* under consideration will be able to withstand the **explosion** of a certain **degree** which is calculated based on *user input*”, where italicized names are “named ideas”, and bold-faced names are “concept chunks” (named ideas with a domain of related ideas). We call this goal a “concept instance” (a concept chunk applied in some way). This language lets us perform various static analyses on our artifacts.

Step 3c: Characteristics of a good solution One of our outputs is a probability, which should be checked to be between 0 and 1 (code not shown).

Step 4: Specialization of theories In the GlassBR example, the main specialization is that the thickness parameter is not free to vary, but must take one of a specific set of values. The rationale for this specialization is that manufacturers have standardized the glass thickness they will provide.

The phenomenological nature of the GlassBR example is atypical in this way; most research software examples involve more significant specialization, such as partial differential equations to ordinary, elimination of variables, use of closed-forms instead of implicit equations, and so on.

Step 5: Code-level choices We can choose the output programming language, how “modular” the generated code is, whether we want programs or libraries, the level of logging and comments, etc. Here we show the actual code we use for this, as it is reasonably direct.

```
code :: CodeSpec
code = codeSpec fullSI choices allMods

choices :: Choices
choices = defaultChoices {
  lang = [Python, Cpp, CSharp, Java, Swift],
  modularity = Modular Separated,
  impType = Program, logFile = "log.txt",
  logging = [LogVar, LogFunc],
  comments = [CommentFunc, CommentClass, CommentMod],
  doxVerbosity = Quiet,
  dates = Hide,
```

```
onSfwrConstraint = Exception, onPhysConstraint = Exception,
inputStructure = Bundled,
constStructure = Inline, constRepr = Const
}
```

Step 6: Recipe to generate artifacts

The various pieces of knowledge specified in the previous steps are assembled into *recipes* for generating the artifacts, for example specification, code, dependency diagrams and log of all choices used. Each one of these has its own DSL that encodes the parts of a specification, the requirements for code³, etc.

The venerable Makefile language is a good example of a DSL for specifying build scripts, but where our ideas are more properly compared with the kind of thinking behind *Rattle* [11].

And finally, we can put all of these things together, and generate **everything**:

```
main :: IO ()
main = do
  setLocaleEncoding utf8
  gen (DocSpec (docChoices SRS [HTML, TeX]) "GlassBR_SRS") srs
  printSetting
  genCode choices code
  genDot fullSI
  genLog fullSI printSetting
```

4 CONCLUDING REMARKS

Opportunity : well-understood + long-lived
Tools : knowledge capture + DSLs + generators
Result : long-term productivity gain

We get productivity gains because there is in fact an incredible amount of *knowledge duplication* present in software, especially in well understood domains. For these domains building software really should be a matter of assembly-line style engineering. Furthermore, there is also inherent knowledge duplication between the artifacts of a project because *they are about the same topic*.

This means that if we spend some up-front time capturing the fundamental knowledge of specific domains (such as mechanics of rigid body motion, dynamics of soil, trains, etc), most later development time can be spent on the *specifics* of the given requirements.

As a side-effect, we obtain traceability and consistency, by construction. Traceability is illustrated by Figure 1, which shows we can track where each of the concepts is used. An example of consistency is found in the table of symbols that appears in our generated documentation. The list of symbols is not entered by the user; it is automatically constructed from the base knowledge (Step 3a).

There are further ideas that co-exist smoothly with our framework, most notably software families and correctness. Our process will remove human errors from generating and maintaining documentation, code, test cases and build environments, since these mundane details are handled by the generator. Furthermore, our work makes it easy to experiment with “what if” scenarios, which make it easy to explicitly track the ramifications proposed changes.

With the right up-front investment, we can have sustainable software development because stable knowledge is separated from local context, which is where most of the rapidly changing assumptions and requirements reside.

³which has been already published and will be cited later

REFERENCES

- [1] ASTM. 2009. Standard Practice for Determining Load Resistance of Glass in Buildings.
- [2] ASTM. 2015. Standard Practice for Specifying an Equivalent 3-Second Duration Design Loading for Blast Resistant Glazing Fabricated with Laminated Glass.
- [3] W. Lynn Beason, Terry L. Kohutek, and Joseph M. Bracci. 1998. Basis for ASTM E 1300 Annealed Glass Thickness Selection Charts. *Journal of Structural Engineering* 124, 2 (February 1998), 215–221.
- [4] Barry W. Boehm. 1987. Improving Software Productivity. *Computer* 20, 09 (1987), 43–47.
- [5] Kimiz Dalkir. 2011. *Knowledge Management in Theory and Practice* (2nd ed.). The MIT Press, Cambridge, Massachusetts, USA.
- [6] William M. Farmer. 2007. Biform Theories in Chiron. In *Towards Mechanized Mathematical Assistants*, Manuel Kauers, Manfred Kerber, Robert Miner, and Wolfgang Windsteiger (Eds.). Springer Berlin Heidelberg, Berlin, Heidelberg, 66–79.
- [7] IEEE. 1998. Recommended Practice for Software Requirements Specifications. *IEEE Std 830-1998* 830 (Oct. 1998), 1–40. <https://doi.org/10.1109/IEEESTD.1998.88286>
- [8] David Lorge Parnas and Paul C Clements. 1986. A rational design process: How and why to fake it. *IEEE transactions on software engineering* 2 (1986), 251–257.
- [9] Suzanne Robertson and James Robertson. 1999. *Mastering the Requirements Process*. ACM Press/Addison-Wesley Publishing Co, New York, NY, USA, Chapter Volere Requirements Specification Template, 353–391.
- [10] W. Spencer Smith, Lei Lai, and Ridha Khedri. 2007. Requirements Analysis for Engineering Computation: A Systematic Approach for Improving Software Reliability. *Reliable Computing, Special Issue on Reliable Engineering Computation* 13, 1 (Feb. 2007), 83–107. <https://doi.org/10.1007/s11155-006-9020-7>
- [11] Sarah Spall, Neil Mitchell, and Sam Tobin-Hochstadt. 2020. Build Scripts with Perfect Dependencies. *Proc. ACM Program. Lang.* 4, OOPSLA, Article 169 (18 November 2020), 28 pages. <https://doi.org/10.1145/3428237>