

An Edit Calculus for Probabilistic Program Repair

Breandan Considine



School of Computer Science
McGill University
Montreal, Quebec, Canada

October 28, 2024

A thesis proposal submitted to McGill University in partial
fulfillment of the requirements of the degree of
Doctor of Philosophy

©Breandan Considine, 2024

Abstract

Official McGill Guidelines: If the language of the thesis is neither English nor French (only allowed for specific language Units) then a third abstract in the language of the thesis is required.

Abstracts in English and French are mandatory and must be text only, i.e. no images, special characters (apart from the West European character set excluding the “Œ” and “œ”), chemical or mathematical formulae, or special formatting (e.g. lists, tables). Abstracts have a maximum limit of 4000 characters.

Abrégé

Official McGill Guidelines: La même chose en français.

Contribution

Official McGill Guidelines: A doctoral thesis must clearly state the elements of the thesis that are considered original scholarship and distinct contributions to knowledge.

- Contributions of the student to each chapter must be explicitly stated.
- Contributions of any co-authors to each chapter must be explicitly stated.

Acknowledgements

Official McGill Guidelines: Among other acknowledgements, the student is required to declare the extent to which assistance (paid or unpaid) has been given by members of staff, fellow students, research assistants, technicians, or others in the collection of materials and data, the design and construction of apparatus, the performance of experiments, the analysis of data, and the preparation of the thesis (including editorial help).

- In addition, it is appropriate to recognize the supervision and advice given by the thesis supervisor(s) and advisors.

Contents

1	Formal Language Theory	3
2	Deterministic Program Repair	4
3	Probabilistic Program Repair	5
4	Discussion	6
5	Conclusions and Future Work	7

List of Figures

List of Tables

Introduction

Pray, Mr. Babbage, if you put
into the machine wrong figures,
will the right answers come out?

—Charles Babbage (1791–1871)

Computer programs are instructions for performing a chore that humans would rather avoid doing ourselves. In order to persuade the computer to do them for us, we must communicate our intention in a way that is plain and unambiguous. Programming languages are protocols for this dialogue. Every program is written in a language that was designed to facilitate the exchange of information between programmers and computers.

Programs are seldom written from left-to-right in one go. During the course of writing a program, the programmer often revisits and revises code as they write, sharing additional information and receiving feedback. Often, during this process, the programmer makes a mistake, causing the program to behave in an unexpected manner. This mistake can be as simple as a typographic or syntactic error, or a more subtle runtime error.

To intercept these errors, programming language designers have adopted a convention for specifying valid programs, called a grammar, which serves two essential purposes. The first is to reject ill-formed programs, and the second is to transform the sequence into a treelike intermediate representation that can be handled by a compiler. We will focus on the first case.

When a program enters an invalid state, a series of unfortunate events occur. The compiler halts, and an error message is raised. To rectify this situation, the programmer must pause their work, plan a fix, apply it, then try to remember what they were doing beforehand. The cognitive overhead of this simple but repetitive chore can be tiresome. To make matters worse, the error message may be unhelpful or challenging to diagnose.

Program repair attempts to address this problem by inferring the author’s intent from an approximate solution. We can think of this as playing a kind of language game. Given an invalid piece of source code for some programming language, the objective of this game is to modify the code to satisfy the language specification. The game is won when the proposed solution is both valid and the author is satisfied with the result. We want to play this game as efficiently as possible, with as little human feedback as possible.

Our goal in this thesis is to introduce a theory of program repair. Broadly, our approach is to repair faulty programs by combining probabilistic language models with exact combinatorial methods. We do so by reformulating the problem of program repair in the parlance of formal language theory. In addition to being a natural fit for syntax repair, this also allows us to encode and compose static analyses as grammatical specifications.

Program repair is a highly underdetermined problem, meaning that the validity constraints do not uniquely determine a solution. A proper theory of program repair must be able to resolve this ambiguity to infer the user’s intent from an incomplete specification, and incrementally refine its guess as more information becomes available from the user.

This theory we propose has a number of desirable properties. It is highly compositional, meaning that users can manipulate constraints on programs while retaining the algebraic closure properties, such as union, intersection, and differentiation. It is well-suited for probabilistic reasoning, meaning we can use any probabilistic model of language to guide the repair process. It is also amenable to incremental repair, meaning that we can repair programs in a streaming fashion, while the user is typing.

1

Formal Language Theory

In computer science, it is common to conflate two distinct notions for a set. The first is a collection of distinct objects sitting on some storage device. The second is a lazy construction: not an explicit collection of objects, but a representation that allows us to efficiently determine membership on demand. This lets us represent infinite sets without requiring an infinite amount of memory. Inclusion then, instead of being simply a lookup query, becomes a decision procedure. This is the basis of formal language theory.

The representation we are chiefly interested in are grammars, which are a common metanotation for specifying the syntactic constraints on programs shared by nearly every programming language. Programming language grammars are overapproximations to the true language of interest, providing a fast procedure for rejecting invalid programs and parsing valid ones.

Like sets, it is possible to abstractly combine languages by manipulating their grammars, mirroring the setwise operations of union, intersection, and difference over languages. These operations are convenient for combining, for example, syntactic and semantic constraints on programs. For example, we might have two grammars representing two properties that are both necessary for a program to be considered valid. We can treat valid programs as a subset of the intersection between the two languages.

2

Deterministic Program Repair

Parsimony is a guiding principle in program repair that comes from the 14th century Fransiscan friar named William of Ockham. In keeping with the Fransiscan minimalist lifestyle, Ockham’s principle basically says that when you have multiple hypotheses, the simpler one is the better. It is not precisely clear what “simpler” should mean in the context of program repair, but a first-order approximation is to strive for the smallest number of changes required to transform an invalid program into a valid one.

Levenshtein distance is one such metric for measuring the number of edits between two strings. First proposed by the Soviet scientist Vladimir Iosifovich Levenshtein, it quantifies how many insertions, deletions, and substitutions are required to transform one string into another. As it turns out, there is an automaton, called the Levenshtein automaton [?], that recognizes all strings within a certain Levenshtein distance of a given string. We can use this automaton to find the most likely repair consistent with the observed program and the grammar.

Given the source code for a computer program $\hat{\sigma}$ and a grammar G , our goal is to find the most likely valid string σ consistent with the grammar G and the observed program $\hat{\sigma}$. We can formalize all possible repairs as a language intersection problem.

3

Probabilistic Program Repair

As we have seen, the problem of program repair is highly underdetermined. To resolve this ambiguity, we will use a probabilistic model to induce a distribution over the language of valid programs. This distribution will guide the repair process by assigning a likelihood to each possible repair. Then, taking the maximum over all possible repairs, we can find the most likely repair consistent with the constraints and the observed program.

Specifically, we will define an ordering over strings by their likelihood under the probabilistic model. We then define a repair as the most likely string consistent with the observed program and the grammar. We factorize the probability of a string as the product of the probability of each token in the string, conditioned on the previous tokens. This allows us to compute the likelihood of a string in a left-to-right fashion.

This probabilistic model will generally admit programs that are locally probable, but globally inconsistent with the grammar. To enforce syntactic validity, we will use the probabilistic language model to “steer” a generative sampler through the automaton representing the repair language. This has two advantages: first, it allows us to sample from the repair language incrementally, and second, it ensures that subsequences with high probability are retrieved first, and all trajectories are syntactically valid.

4

Discussion

Official McGill Guidelines: The discussion of findings must be in line with disciplinary expectations. A comprehensive discussion is expected to be a minimum of 10 pages, double-spaced for doctoral students and a minimum of 5 pages, double-spaced for Master's students (including figures, images, and tables). It pertains to the entirety of a thesis. The discussion of findings should provide an final, overarching summary of study themes, limitations, and future directions.

In the case of a manuscript-based thesis, the comprehensive discussion should encompass all of the chapters of the thesis and should not be a repetition of the individual chapters. This section can be used to address issues not sufficiently covered in the preceding chapters or papers (e.g., critiques raised by reviewers that could not be incorporated into published works, or reintroducing discussion arguments removed from published papers upon reviewer request). This section can also be used to elaborate on the practical/applied aspects of published findings in a manner that is more accessible to less expert readers.

5

Conclusions and Future Work

Official McGill Guidelines: Clearly state how the objectives of the research were met and discuss implications of findings.

Publications

This part is optional, but it gives a nice touch to list all the publications (official venues down to poster sessions) throughout your PhD.

Keep this in the same style as publications in your academic CV: Conference / Year - Title - Authors. And here comes a sample ref for the bibliography: [Con23]

Acronyms

This part is likewise optional. But it does not hurt to provide a list of all acronyms, e.g.:

- **REST**: **R**epresentational **S**tate **T**ransfer

Bibliography

- [Con23] Breandan Considine. A pragmatic approach to syntax repair. In *Companion Proceedings of the 2023 ACM SIGPLAN International Conference on Systems, Programming, Languages, and Applications: Software for Humanity*, pages 19–21, 2023.