



UNIWERSYTET
IM. ADAMA MICKIEWICZA
W POZNANIU

Wydział Matematyki i Informatyki

Bohdan Bondar, Marcin Jałowski, Aleksander Mendoza-Drosik
Numer albumu: sS432778,s434701,s434749

Solomonoff - kompilator transduktorów skończenie stanowych

Solomonoff - Finite state transducer compiler

Praca inżynierska na kierunku **informatyka**
napisana pod opieką
dr Bartłomieja Przybylskiego

Poznań, luty 2021

Poznań, 19 stycznia 2021 r.

Oświadczenie

Ja, niżej podpisana **Bohdan Bondar, Marcin Jałowski, Aleksander Mendoza-Drosik**, studentka Wydziału Matematyki i Informatyki Uniwersytetu im. Adama Mickiewicza w Poznaniu oświadczam, że przedkładaną pracę dyplomową pt. *Solomonoff - kompilator transduktorów skończenie stanowych* napisałam samodzielnie. Oznacza to, że przy pisaniu pracy, poza niezbędnymi konsultacjami, nie korzystałam z pomocy innych osób, a w szczególności nie zlecałam opracowania rozprawy lub jej części innym osobom, ani nie odpisywałam tej rozprawy lub jej części od innych osób. Oświadczam również, że egzemplarz pracy dyplomowej w wersji drukowanej jest całkowicie zgodny z egzemplarzem pracy dyplomowej w wersji elektronicznej. Jednocześnie przyjmuję do wiadomości, że przypisanie sobie, w pracy dyplomowej, autorstwa istotnego fragmentu lub innych elementów cudzego utworu lub ustalenia naukowego stanowi podstawę stwierdzenia nieważności postępowania w sprawie nadania tytułu zawodowego.

[TAK/TAK/TAK] wyrażam zgodę na udostępnianie mojej pracy w czytelni Archiwum UAM

[TAK/TAK/TAK] wyrażam zgodę na udostępnianie mojej pracy w zakresie koniecznym do ochrony mojego prawa do autorstwa lub praw osób trzecich

Streszczenie

Słowa kluczowe: klasa

Abstract

Nam dui ligula, fringilla a, euismod sodales, sollicitudin vel, wisi. Morbi auctor lorem non justo. Nam lacus libero, pretium at, lobortis vitae, ultricies et, tellus. Donec aliquet, tortor sed accumsan bibendum, erat ligula aliquet magna, vitae ornare odio metus a mi. Morbi ac orci et nisl hendrerit mollis. Suspendisse ut massa. Cras nec ante. Pellentesque a nulla. Cum sociis natoque penatibus et magnis dis parturient montes, nascetur ridiculus mus. Aliquam tincidunt urna. Nulla ullamcorper vestibulum turpis. Pellentesque cursus luctus mauris.

Keywords: klasa

Tu możesz umieścić swoją dedykację.

Spis treści

Rozdział 1. Introduction	7
Rozdział 2. Transducers	12
2.1. Mealy automata and transducers	12
Rozdział 3. Build system	17
Rozdział 4. Web technologies	21

ROZDZIAŁ 1

Introduction

This project focuses on research in the field of automata theory and inductive inference. The main product of our work is the "Solomonoff" regular expression compiler for finite state transducers. Plenty of research has gone into development of the theory behind this system. As a result the transducers contain several features not known before.

The most innovative achievements is the lexicographic arctic semiring of weights, specialized adaptation of Glushkov's construction for subsequential transducers and the most significant flagship feature - built-in support for inductive inference and machine learning of transducers. Thanks to the cooperation with LearnLib and Dortmund University, Solomonoff supports learning algorithms such as RPNI and several of its derivatives. Solomonoff contributes its own more specialized for transducers inference. We implemented OSTIA for efficient learning of deterministic transducers. For nondeterministic ones we developed our own OFTIA algorithm, that was not known before.

All those features together make Solomonoff a unique library that stands out from all the alternatives. We support most of the features of UNIX regexes, including look-aheads and look-behinds, which is unusual for automata-based regex engine. The key that allows Solomonoff for doing that is the possibility of emulating look-aheads with careful placement of transducer outputs. As a result Solomonoff can compete and do much more than existing projects such as RE2 developed by Google, BRICKS automata developed at Aarhus University or even to certain extent with Pearl/Java based regular expression engines. Another, much stronger competitor for Solomonoff is the OpenFST project developed by Google. Their Thrax grammars are capable of doing most of the things that Solomonoff can and they also support probabilistic automata. OpenFST is much older and more established in the scientific community. They support a lot more features that were developed by scientists, by the course of many years. Solomonoff cannot compete with this level of sophistication, but perhaps, what might seem like a limitation, is in fact our strongest advantage. Solomonoff

focuses on functional transducers and enforces this property at compilation time. Any arising nondeterministic output is automatically rejected as an error. This allows Solomonoff to perform a lot more optimisations, the automata are smaller and their behaviour is more predictable. Moreover, lexicographic weights allow for precise disambiguation of nondeterministic paths whenever necessary and their most important advantage is that lexicographic semiring is not commutative and hence it does not "propagate" throughout entire automaton. This only increases Solomonoff's robustness and allows for (exponentially) smaller automata, without sacrificing predictability of the regular expression. On the contrary, probabilistic weights in Thrax, make the whole system, more heavyweight, unpredictable and difficult to maintain. The results are especially palpable when comparing our benchmarks. Solomonoff was written in Java and Thrax uses C++, but despite this our compiler is several magnitudes more efficient. We performed efficiency tests on a large corpus of linguistic data (dictionary with 6000 records). Solomonoff compilation times were around 2 seconds, whereas Thrax took 19 minutes. Solomonoff's automata were also much smaller, as thanks to Glushkov's construction, there is a 1:1 relationship between size of regular expression and size of transducer. As a result Thrax's transducer takes of 6336K of RAM, whereas Solomonoff only takes 738K. Execution time for such large corpora was about 10 milliseconds in Solomonoff (which is roughly the same performance as using Java's HashMap), while Thrax took about 250 milliseconds. One might argue that, such great differences are achievable, only because Thrax supports a lot more features. OpenFST uses epsilon transitions, while Solomonoff does not implement them and instead all automata are always and directly produced in epsilon-free form. OpenFST performs operations such as sorting of edges, determinization, epsilon-removal and minimization. Solomonoff always has all of its edges sorted and it doesn't need a special routine for it (which makes for additional performance gains), it has no epsilons to remove and it does not need determinization procedure, because it can pseudo-minimise nondeterministic transducers, using heuristics inspired by Kameda-Weiner's NFA minimization. One could, half-jokingly, summarize that the difference between Thrax and Solomonoff is like that between "Android and Apple". Thrax wants to support "all of the features at all cost", whereas Solomonoff carefully chooses the right features to support. We believe this approach will be our strongest asset, that will make our compiler a serious alternative to the older and more established OpenFST. An additional strength that favours Solomonoff over OpenFST is that we support inductive inference out-of-the-box, require no programming (regular expressions are the primary user interface and Java API is minimal and optional) and we provide automatic

conversion from Thrax to Solomonoff, so that existing codebases can be easily migrated.

The characteristic feature of Solomonoff, is that its development focuses on the compiler and regular expressions instead of library API. Everything can be done without writing any Java code. It also allows the developers for much greater flexibility, because the internal implementation can be drastically changed at any time, without breaking existing regular expressions. There is very little public API that needs to be maintained. As a result backwards-compatibility is rarely an issue.

The primary philosophy used in implementing this library is the top-down approach. Features are added conservatively in a well thought-through manner. No features will be added ad-hoc. Everything is meant to fit well together and follow some greater design strategy. For comparison, consider the difference between OpenFst and Solomonoff.

- OpenFst has Matcher that was meant to compactify ranges. In Solomonoff all transitions are ranged and follow the theory of (S,k) -automata. They are well integrated with regular expressions and Glushkov's construction. They allow for more efficient squaring and subset construction. Instead of being an ad-hoc feature, they are well integrated everywhere.
- OpenFst had no built-in support for regular expression and it was added only later in form of Thrax grammars, that aren't much more than another API for calling library functions. In Solomonoff the regular expressions are the primary and only interface. Instead of having separate procedures for union, concatenation and Kleene closure, there is only one procedure that takes arbitrary regular expression and compiles it in batches. This way everything works much faster, doesn't lead to introduction of any ϵ -transitions and allows for more optimisation strategies by AST manipulations. This leads to significant differences in performance.

While, most of other automata libraries (RE2/BRICKS) were not designed for large codebases and have no build system, Solomonoff comes with its very own build system out-of-the-box. Thrax is the only alternative that does have a "build system" but it's very primitive and relies on generating Makefiles. Solomonoff has a well integrated tool that assembles large projects, detects cyclic dependencies, allows for parallel compilation and performs additional code optimisations. It also ships with syntax highlighting and tools that assist code refactoring. Our project strives to make automata as easy to use and accessible to mass audiences as possible. For this reason we developed a website with online playground where visitors could test Solomonoff and experiment without

any setup required. The backend technology we used is Spring Boot, because it allowed for convenient integration with existing API in Java.

While at the beginning our attempts focused on implementing the library in C, switching to Java turned out to be a major advantage. While, it's true that C allows for more manual optimisations, having a garbage collector proved to be a strong asset. The compilation relies on building automata in form of directed graphs. Each vertex itself is a separate Java object. The automaton is always kept trim because all the unreachable vertices are lost and free to be garbage collected at any time. If we tried to implement such data structure in C we would need to reinvent a simple (and most likely, less efficient) garbage collector ourselves. Manipulation of linked lists, linked graphs and other recursive data structures is not as convenient in C as it is in Java. Another advantage that Java gave us over C is that the plenty of existing infrastructure was already implemented in Java. Hence we became more compatible with Samsung's systems where Solomonoff could be deployed. We could also easily integrate our compiler with LearnLib.

For environments where Java is not an option, Solomonoff will allow for generating automata in form of C header files that can be easily included and deployed in production code. It will not allow for manipulation of automata from C level, but our compiler was never intended to be used programmatically. There is, in fact, not much benefit from doing so. If user needs to create their custom automata manually, it's much easier to generate them in AT&T format and then have them read by the compiler. Later they can be manipulated with regular expressions like any other automaton. Generating AT&T has also additional benefits, as such custom script for generating AT&T can be attached to Solomonoff build system as an external routine and then compilation can be optimised more efficiently. An example scenario would be having two independent routines that the build system could decide to run in parallel and then it could cache the results for subsequent rebuilds. For comparison, if anybody decided to use OpenFST manually from C++, they would need to first learn the API (which is less user friendly than learning AT&T format) and then they would also need to implement caching and parallelization themselves. Moreover, the users are also tied to C++, whereas, AT&T format can be generated in any language. This shows, why API for programmatic manipulation offers little to no benefit, while imposes much more limitations. It should also be mentioned that Solomonoff allows user for implementing their own Java functions that could then be incorporated into the regular expressions as "native" calls.

To guarantee the better utility of Solomonoff, it has to provide ability of interactive evaluation in order to enable efficient and convenient work with the compiler. It can be also vastly useful while learning. The mentioned feature is

implemented with a REPL. The biggest challenge here is adjusting the entire stack of compilation processes to work interactively with a sensible manner. The other condition of usefulness is a fast build automation. To make our system meet the need it has to support both parallel compilation and caching of previously compiled units. Also the ability to resolve dependencies may be important in some more complex projects. In this case the main challenge is to find a satisfying solution of ordering in a parallel compilation. This problem can be broken down to a topological sort in a directed acyclic graph. To make the use of Solomonoff even more practical it provides an export option for transducers in a way that allows to include them into a C code during a recompilation phase.

ROZDZIAŁ 2

Transducers

The theory behind implementation and design of Solomonoff has been studied and developed over the course of 2 years. The early versions and our initial ideas looked very different to the final results we've achieved.

At the beginning it was meant to be a simple tools that focused only on deterministic Mealy automata. Such a model was very limited and it quickly became apparent that certain extensions would need to be made. In the end we implemented a compiler that supports nondeterministic functional weighted symbolic transducers.

2.1. MEALY AUTOMATA AND TRANSDUCERS

The standard definition of automaton found in introductory courses states that finite state automaton is a tuple $(Q, I, \delta, F, \Sigma)$ where Q is the set of states, $I \subset Q$ are the initial states, $F \subset Q$ are the final (or accepting) states, $\delta : Q \times \Sigma \rightarrow Q$ is the transition function and Σ is some alphabet.

A Mealy machine extends the above definition with output $(Q, I, \delta, F, \Sigma, \Gamma)$ where Γ is some output alphabet and transition function has the form of $\delta : Q \times \Sigma \rightarrow Q \times \Gamma$.

Such a model is frequently used in the field of formal methods. Many complex state-based systems can be modelled and simplified using Mealy machines. As an example consider a black-box computer program whose logs can be observed. The current snapshot of the program's memory determines its state. Depending on subsequent user input, we might observe different log traces. There are many existing machine learning and inference algorithms that can build an equivalent model of Mealy machine only by interacting with such a black-box system and reading its logs.

It can be proved that the expressive power of deterministic automata with output is strictly less than that of their nondeterministic counterparts. It is

known as the prefix-preserving property. If a deterministic automaton read input string $\sigma_1\sigma_2\sigma_1$ and printed $\gamma_1\gamma_1\gamma_2$, then at the next step it is only allowed to append one more symbol to the previously generated output. For instance we could not suddenly change the output to $\gamma_2\gamma_2\gamma_2\gamma_1$ after reading one more input symbol $\sigma_1\sigma_2\sigma_1\sigma_2$. The prefix $\gamma_1\gamma_1\gamma_2$ must be preserved.

The problems that we wanted to tackle with Solomonoff revolved around building sequence-to-sequence models. For instance we might want to translate from numbers written as English words into digits. A sentence like "one apple" should become "1 apple". The prefix preserving property would be too limiting because it often happens that the suffix of string has decisive impact on the translation. The phrase "once again" also starts with prefix "one" but it should not be translated as "1ce again"!

We intended to find the right balance between expressive power of nondeterministic machines and strong formal guarantees of Mealy automata. To achieve this, we initially decided to use multitape automata. The idea was to write all possible continuations of given output and store each in a separate output tape. Then upon reaching the end of string, the state would decide which tape to use as output.

Over the course of research and development we have discovered that the power of multitape Mealy machines was still too limiting for our purposes. In particular we could define congruence classes on pairs of strings similar to those in Myhill-Nerode theorem. Then it can be easily noticed that as long as the number of output tapes is finite, the number of congruence classes must be finite as well. A very simple counter-example would be the language that for every string a, aa, aaa, \dots respectively prints output c, bc, bbc, \dots . It could not be expressed using only a finite number of output tapes, because there are infinitely many ways to continue the output and none of them is a prefix of the other.

Yet another limitation of Mealy machines is that their δ function enforces outputs of the exact same length as inputs. In the field of natural language processing such an assumption is too strict. For instance, we might want to build a machine that translates sentences from one language to another. A word "fish" in English might have 4 letters but Spanish "pescado" is much longer. Our first idea was to use sequential transducers instead of the plain Mealy automata. Their definition allows the transition function to print output strings of arbitrary length $\delta : Q \times \Sigma \rightarrow Q^*$.

The nondeterministic transducers don't suffer from any of the above problems. Their expressive power exactly corresponds to that of regular transductions. It's a very strong and expressive class. The only way to obtain a stronger model would be by introducing context-free grammars and pushdown

automata. The reason why we were hesitant to use this approach was because of the possible ambiguity of output. Nondeterministic transducers may contain epsilon transitions, which could lead to infinite number of outputs for a any given input. Without epsilons, the number of outputs is finite but it's still possible to return more than one ambiguous output.

The model of automata that was finally implemented in Solomonoff turned out to be the functional nondeterministic transducers. Functionality means that for any input, there may be at most one output. Such a model proved to provide the perfect balance of power with many strong formal guarantees. While epsilon transitions strictly increase power of transducers, when restricted only to functional automata, the erasure of epsilons becomes possible. Using advance-and-delay algorithm one can test functionality of any automaton in quadratic time. There exists a special version of powerset construction that can take any functional transducers and produce an equivalent unambiguous automaton. One can take advantage of unambiguity to build an inference algorithm for learning functional automata from sample data. The automata are closed under union, concatenation, Kleene closure and composition. Unlike nonfunctional transducers, they are not closed under inversion but we developed a special algebra that uniquely defines an invertible bijective transduction for any automaton. Glushkov's construction can be augmented to produce functional transducers. Lexicographic semiring of weights can be used to make functional automata more compact.

Once we decided to use the power of functional transducers, the next problem we had to solve was their optimisation. One of the most important operations in natural language processing is the context-dependent rewrite. The standard way of implementing it is by building a large transducer that handles all possible cases. In Solomonoff we've developed our own approach that produces much smaller automata. It is done with lexicographic weights.

Another important optimisation is the state minimisation. Many existing libraries implement separate functions for all regular operations and additional one for minimisation. A regular expression like $A(B + C)^*$ would be the translated to a series of function calls like

$$\text{minimise}(\text{concatenate}(A, \text{kleene_closure}(\text{union}(B, C))))$$

In Solomonoff we don't have separate implementation for those operations. Instead everything is integrated in form of one monolithic procedure that implements Glushkov's construction. This way the produced automata are very small even without the need for minimisation. If the regular expression consists of n input symbols, then the resulting transducers has $n + 1$ states. Because there is one-to-one correspondence between regular expression symbols and automata

states, we are able to retain plenty of useful meta information. In particular each state can tell us exactly which source file it comes from and the precise text line and column in that file. This way, whenever compilation fails, user can see meaningful error messages.

The standard minimisation procedure used by most libraries works by finding the unique smallest deterministic automaton. Nondeterministic automata do not admit unique smallest representative and might be even exponentially smaller than their equivalent minimal deterministic counterparts. Finding the smallest possible nondeterministic automaton is a hard problem. For this reason Solomonoff implements a heuristic pseudo-minimisation algorithm that attempts to compress nondeterministic transducers and does not attempt to determinise them. Glushkov's construction already produces very small automata, hence any attempt at minimising them by determinisation would result in larger automata than the initial ones.

Lexicographic weights allow to make the transducers even more compact. We proved that there exist weighted automata exponentially smaller than even the smallest unweighted nondeterministic ones.

In the tasks of natural language processing it's common to work with large alphabets. User input might contain unexpected sequences like math symbols, foreign words or even emojis and other UNICODE entities. The regular expressions should handle such cases gracefully, especially when using wildcards such as `.*` or `\p{Lu}`. Representation of large character classes is a challenging task for finite state automata. In order to use the dot wildcard in UNICODE, the automaton would require millions of transitions, one for each individual symbol. In order to optimise this, Solomonoff employs symbolic transitions. While classic automata have edges labelled with individual symbols, our transducers have edges that span entire ranges. A range is easy to encode in computer's memory. It's enough to store the first and last symbol.

Compilation of transducers is the core part of our library but in order to make the automata useful there must be a way to execute them. Deterministic transducers and Mealy machines could be evaluated in linear time. Nondeterministic automata are more expensive. The computation might branch and automaton could be in multiple states simultaneously. Using dynamic programming it's possible to build a table with rows representing symbol on input string and each column keeping track of one state. At each step of execution, one input symbol is read and one row in the table is filled based on the contents of previous row. By the end of evaluating such a table, it's enough to scan the last row to find a column of some accepting state and then the table can be backtracked from that

state backwards. The backtracking step is necessary to correct output produced from each transition.

This algorithm has been made even more efficient by using techniques from graph theory. Every automaton is a directed graph that could be represented as adjacency matrix. Sparse graphs can be optimised and instead of using matrix it's possible to only store the list of adjacent vertices. Automata very often are a perfect example of sparse graphs. Hence instead of using a table, it's better to use a list of states nondeterministically reached at any given step of evaluation. The backtracking can be made efficient by storing pointer to source state taken at any transition.

Glushkov's construction relies on building three sets: the set of initial symbols, final symbols and 2-factor strings. The prototype version of Solomonoff would represent sets as bitsets, where each bit specifies whether element belongs to the set or not. This representation simplified implementation of many algorithms but was highly inefficient. The current implementation uses hash sets instead. While hash maps have constant insertions and deletions, they ensure it at the cost of larger memory consumption. During benchmarks on real-life datasets, Solomonoff would often run out of memory and crash. Further optimisation was obtained by representing sparse sets as arrays of elements. This approach proved to be the best choice, because Glushkov's construction inherently ensures uniqueness of all inserted elements (hence using hashes to search for duplicates before insertion was not necessary) and it does not use deletions. As a result any array was guaranteed to behave like a set.

The standard definition of Glushkov's construction produces set of 2-factors as its output. Then a separate procedure would be necessary to collect all such strings and convert them into a graph of automaton. We found a way to make it more efficient and build the graph directly. The step of building 2-factors was entirely bypassed. This provided even further optimisation.

One of the core features of Solomonoff is the algorithm for detecting ambiguous nondeterminism. Advance-and-delay procedure can check functionality of automaton in quadratic time. While the compiler does provide implementation of this procedure, it is not used for checking functionality. Instead we use a simpler algorithm for checking strong functionality of lexicographic weights. While the performance difference between the two is negligible, the advantage of our approach comes from better error messages in case of nondeterminism. If some lexicographic weights are in conflict with each other, compiler can point user to the exact line and column of text where the conflict arises, whereas advance-and-delay might miss the origin of problem and only fail further down the line.

ROZDZIAŁ 3

Build system

The primary goal of this project was to develop the compiler backend. Initially little attention was paid to structuring larger projects composed of multiple source files and dependencies. While smaller applications can work fine with a single script of regular expressions, in the tasks on natural language processing it's common to build large codebases composed of millions of linguistic rules. Being able to manage and organise them efficiently becomes a major issue.

The compiler builds automata using our special version of Glushkov's construction, which is capable of compiling subexpressions independently of each other. The results are returned in form of a singly linked graph. This format allows one to parallelise compilation and split work across multiple files. For instance, a user could create two files `file1.mealy` and `file2.mealy`, define some variable in the first one and use it in the second.

```
1 | file1.mealy:
2 |     var1 = 'abc'
3 | file2.mealy:
4 |     var2 = ('prefix' var1 'suffix')*
```

Such a feature is not native to the compiler backend itself (it only parses continuous string of input) and has been delegated to the build system as an independent application instead.

The easiest approach to implementing a build system would be by concatenating all source files into one large stream of input and then feed it into the compiler. For instance

```
1 | concatenatedFiles.mealy:
2 |     //from file1.mealy
3 |     var1 = 'abc'
4 |     //from file2.mealy
5 |     var2 = ('prefix' var1 'suffix')*
```

It could be easily done with a simple Bash script or a couple of Makefiles.

Despite being very straightforward, such a solution has multiple flaws. The order of concatenation matters a lot. The compiler requires that every variable is defined before being used. Should the order get mistakenly swapped, the compilation would fail. For instance

```
1 concatenatedFiles.mealy:
2     //from file2.mealy
3     var2 = ('prefix' var1 'suffix')*
4     // var1 used before definition!
5     //from file1.mealy
6     var1 = 'abc'
```

As a result, it would be user's obligation to ensure correct order of concatenation, which might quickly become unmaintainable in large projects. Second problem stems from linear types. The compiler follows the semantics of linear logic and a variable once consumed should not be reused. An explicit copy is always necessary. Hence the order of definition and usage of all variables is even more important than in other general-purpose languages that do not have linear types. For instance this will compile

```
1 X = 'a'
2 Y = !!X 'b'
3 Z = X 'c'
```

but this one will fail

```
1 X = 'a'
2 Z = X 'c'
3 Y = !!X 'b'
```

Third problem is that managing dependencies and building packages would become impossible. For example we might imagine code, which includes some library

```
1 include libX
2 X = f // f is defined in libX
```

If in the future a new version of the library is released, it might happen that variable X is added and suddenly our code becomes invalid because we're trying to redefine X.

Our initial thought was to introduce a special function `import!('filepath')`, which would load any variable from a precompiled binary file. Then the user would write code as follows

```
1 X = import!('lib/libX/f')
```

This approach suffers from one logistic problem. If user decided to split their code into several files and then import the necessary automata at will, the build system would need to detect the correct compilation order of all files. For instance if there were two files like these

```
1 | file X.mealy:
2 |     a = 'a'
3 | file Y.mealy:
4 |     b = import!('X/a')
```

then the build system would need to first compile `X.mealy` and produce binary file `X` and only then compilation of `Y.mealy` would become possible. The major problem appears when cyclic dependencies between files arise

```
1 | file X.mealy:
2 |     X = import!('Y/Y')
3 | file Y.mealy:
4 |     Y = import!('X/X')
```

While it's possible to create a built tool capable of detecting such problems and notifying the user, its main disadvantage is that in some cases, there might be cyclic dependency between files, despite not introducing cyclic dependencies between automata themselves. As a result, logically valid regular expressions would be prematurely rejected by build system. For example

```
1 | file X.mealy:
2 |     X1 = import!('Y/Y1')
3 |     X2 = 'a'
4 | file Y.mealy:
5 |     Y1 = 'b'
6 |     Y2 = import!('X/X2')
```

The final solution we settled for was to split compilation into three phases. First the build system scans all files (in parallel) and builds abstract syntax trees of all variable definitions. In those trees there would be references to other variables that could be defined anywhere else in the project. Then in the second phase the variables would be scanned and a dependency graph would be built. We used a specialised data structure for efficiently working with directed acyclic graphs. If at any point a cycle was introduced, the build system would throw an error. Finally in the last phase every node in the graph would be compiled in parallel. For optimal performance, multiple threads would take vertices in topological order. If there is a directed edge from vertex v_1 to v_2 (meaning that variable v_2

depends on v_1), then v_1 will be added to the queue before v_2 . Threads process queue entries in the FIFO order.

In the first phase, the parser that scans source files has been reused from the compiler backend. More precisely, the grammar of syntax is taken from the library but the build system has its own implementation of `SolomonoffGrammarListener`. Parsing is then performed using the standard ANTLR functions

```
1  SolomonoffGrammarLexer lexer =
2      new SolomonoffGrammarLexer(sourceFile);
3  SolomonoffGrammarParser parser =
4      new SolomonoffGrammarParser(new CommonTokenStream(lexer));
5
6  parser.addErrorListener(new BaseErrorListener() {
7      public void syntaxError(Recognizer, ? recognizer,
8          Object offendingSymbol,
9          int line,
10         int charPositionInLine,
11         String msg,
12         RecognitionException e) {
13         System.err.println("line " + line
14             + ":" + charPositionInLine
15             + " " + msg + " " + e);
16     }
17 });
18
19 final SolomonoffWeightedParser listener =
20     new SolomonoffWeightedParser(collector, sourceFile, compiler);
21 ParseTreeWalker.DEFAULT.walk(listener, parser.start());
```

ROZDZIAŁ 4

Web technologies

Creation of user-friendly interface was an important part of the project. The greatest challenge lied in finding the most intuitive way of presenting a complicated and highly advanced system. The main component was the language of regular expressions itself. User should be able to edit its code with ease. The second key feature was the ability to execute the code. In many Turing-complete languages, every expression can be evaluated into some value, which could be printed back to the user. Here the problem was not so trivial. The regular expressions could in principle be evaluated down to formal languages. For example

```
1 | 'a' ('b' | 'c' | 'ef' ) 'd'
```

would return a language consisting of strings

```
1 | 'abd', 'acd', 'aefd'
```

The issue with such approach is that not all languages are finite. The expression

```
1 | 'a'*
```

would be evaluated as infinite set

```
1 | '', 'a', 'aa', 'aaa', ...
```

Some regexes, might be finite but of exponential size. For instance

```
1 | ('0' | '1') ('0' | '1') ('0' | '1') ('0' | '1')
```

yields set of all bit-strings of length 4. Presenting user with the result in form of formal languages would be often impractical or impossible.

As a result, our REPL does not evaluate expressions. The results of compilation are not printed in any form. Instead the interface is meant to be silent when compilation is successful. Only errors are printed.

There are many different approaches to implement user interface for REPL. One of them would be having a single editor window with all the code in

it and the REPL output printed on the margins next to each respective line. This provides a very immersive user experience for Turing-complete languages. For regular expressions it's not as spectacular. Instead we decided to use two windows - one for code editor and the other for REPL console. All interaction with regular expressions is performed via special commands built into the console. Those commands could not be used inside the code editor, as they are not part of the language itself. In order to evaluate a transducer user would type the following line into REPL input

```
1 | :eval NAME 'input string'
```

Sometimes, user might want to see all the strings that belong to a given language. While there might be infinitely many of them, it's possible to ask user how large sample to generate. To achieve this the following line can be used

```
1 | :rand_sample NAME of_size NUMBER
```

Automata can also be interpreted as directed graphs. This property makes was used to further enhance user interface. Automaton's graph will be shown after typing this command

```
1 | :vis NAME
```

Those and many other functionalities have been implemented in the browser-based version of REPL.

The implementation is not trivial. One of the ways to achieve such results would be by implementing a parser that could halt mid-parsing. For example user could first type

```
1 | x = ('x' |
```

and hit return button. The parser should notice that the expression is not finished and it has to wait for the next line of input. Then as the user types the next line

```
1 | 'y' )
```

a full and valid expression could be recognised and parser could return. This approach is used by some programming languages. It's difficult to implement and requires the grammar to be appropriately structured. We later abandoned this idea due to the problematic nature of Solomonoff's grammar. In particular, it does not use semicolons to separate statements. For example

```
1 | x = 'a'  
2 | y = 'b'
```

could be written in a single line

```
1 | x = 'a' y = 'b'
```

The equality sign determines start of new statement. By its very nature, parsing this, requires a lookahead of one token into the future. When the input is read in fragments, line by line, such a lookahead is not possible to obtain. User could first type

```
1 | x = 'a' y
```

which would be recognized by parser as concatenation of string 'a' with variable y. If the user then types

```
1 | = 'b'
```

in the upcoming line, then the previous results of parsing would have to be discarded and the entire input reparsed again. Hence we decided to simplify the REPL and assume that every line of input fully defines the entirety of expression. As a result it's not possible to split input into multiple lines when using console. This is not a serious limitation, because multiline expressions could still be written in the editor window instead of console.

The division of user interface into editor and console has one more advantage. It closely mimics the layout of command-line interface, where the typical workflow is to edit source code in local files using any text editor of user's choice and the REPL is kept open all the time alongside the editor. Many existing modes for Emacs follow similar convention.

The REPL is implemented on the server-side as a REST API endpoint.

```
1 | @PostMapping("/repl")
2 | public ReplResponse repl(HttpSession httpSession,
3 |     @RequestBody String line)
```

Every user has their own instance of REPL

```
1 | Repl repl = (Repl) httpSession.getAttribute("repl");
```

which holds a reference to the compiler and a set of built-in commands

```
1 | public static class Repl {
2 |     private static class CmdMeta<Result> {
3 |         final ReplCommand<Result> cmd;
4 |         final String help;
5 |         final String template;
6 |
7 |         private CmdMeta(ReplCommand<Result> cmd,
```

```

8         String help, String template) {
9             this.cmd = cmd;
10            this.help = help;
11            this.template = template;
12        }
13    }
14
15    HashMap<String, Repl.CmdMeta<String>> commands;
16    OptimisedHashLexTransducer compiler;
17 }

```

whenever user types some command on the REPL console

```
1 | :cmd arg1 arg2 arg3
```

it gets parsed as

```

1 | String firstWord = "cmd";
2 | String remaining = "arg1 arg2 arg3";

```

and then the appropriate command implementation is looked up in the map

```

1 | final Repl.CmdMeta<String> cmd = commands.get(firstWord);
2 | return cmd.cmd.run(httpSession, compiler, log, debug, remaining);

```

The rest controller contains implementations of many such commands

```

1 | public static final ReplCommand<String> REPL_LIST = ...
2 | public static final ReplCommand<String> REPL_EVAL = ...
3 | public static final ReplCommand<String> REPL_RUN = ...
4 | public static final ReplCommand<String> REPL_EXPORT = ..
5 | public static final ReplCommand<String> REPL_IS_DETERMINISTIC = ...
6 | public static final ReplCommand<String> REPL_LIST_PIPES = ...
7 | public static final ReplCommand<String> REPL_EQUAL = ...
8 | public static final ReplCommand<String> REPL_RAND_SAMPLE = ...
9 | public static final ReplCommand<String> REPL_CLEAR = ...
10 | public static final ReplCommand<String> REPL_UNSET = ...
11 | public static final ReplCommand<String> REPL_RESET = ...
12 | public static final ReplCommand<String> REPL_LOAD = ...
13 | public static final ReplCommand<String> REPL_VIS = ...

```

All of those definitions above are lambda expressions that use library functions of the compiler. The parameters taken by those lambda expressions are as follows

```

1 | public interface ReplCommand<Result> {
2 |     Result run(
3 |         HttpSession httpSession,

```



```

4      OptimisedHashLexTransducer compiler,
5      Consumer<String> log,
6      Consumer<String> debug,
7      String args) throws Exception;
8  }

```

As an example, consider the command

```
1  :eval f 'abc'
```

which evaluates transducer `f` for input string `'abc'`. On the frontend JavaScript will perform REST query

```

1  const response = await fetch('repl', {
2      method: 'POST',
3      body: ":eval f 'abc'"
4  })

```

which will be received by server

```

1  @PostMapping("/repl")
2  public ReplResponse repl(HttpSession httpSession,
3      @RequestBody String line) {
4      Repl repl = (Repl) httpSession.getAttribute("repl");
5      final String result = repl.run(
6          httpSession,
7          line, // ":eval f 'abc'"
8          s -> out.append(s).append('\n'), // console output
9          s -> { } // debug logs are not displayed
10     );
11     ...
12 }

```

and the `repl.run` method will query the appropriate implementation to call for the `:eval` command.

```

1  final String firstWord = "eval";
2  final String remaining = "f 'abc'";
3  final Repl.CmdMeta<String> cmd = commands.get(firstWord);
4  return cmd.cmd.run(httpSession, compiler, log, debug, remaining);

```

This in turn will trigger the following lambda function

```

1  ReplCommand<String> REPL_EVAL =
2      (httpSession, compiler, logs, debug, args) -> {
3      String[] parts = args.split("\\s+", 2); // f 'abc'
4      String name = parts[0]; // f

```

```
5   String input = parts[1]; // 'abc'
6   G transducer = compiler.getTransducer(name);
7   String output = compiler.specs.evaluate(transducer, input);
8   return output == null ? "No match!" : output;
9 };
```

The output is sent back to JavaScript in form of JSON

```
1 | const replResult = JSON.parse(await response.text())
```

Remaining commands are implemented in a similar way.

Several of the functions may require access to HTTP session. In particular it's worth analysing the `:load` command. It's purpose is to emulate the process of loading source code from file. Whenever user types some code in the editor, it needs to be transported to server, then parsed and compiled. All the defined transducers results need to be saved for later use. The simplest way of achieving this, would be by extending the REST API as

```
1 | class ReplInput{
2 |     String command;
3 |     String editorContent;
4 | }
5 | public ReplResponse repl(
6 |     HttpSession httpSession,
7 |     @RequestBody ReplInput)
```

and query it using

```
1 | const response = await fetch('repl', {
2 |     method: 'POST',
3 |     body: {
4 |         command: replCommand,
5 |         editorContent: editor.getValue()
6 |     }
7 | })
```

The downside of such solution is that the editor content could become large and sending it would require more internet bandwidth and time. Often user only wants to execute simple short REPL commands that do not require sending the entire code. Sometimes the code might be required but resending it might be omitted as long as user has not modified it. Hence the process of uploading code to the server has been delegated to a separate REST call.

```
1 | const response = await fetch('upload_code', {
2 |     method: 'POST',
```

```
3         body: code
4     })
```

The code is then stored in HTTP session, so the REST endpoint has a very simple implementation

```
1 @PostMapping("/upload_code")
2 public void uploadCode(HttpSession httpSession,
3     @RequestBody String text) {
4     httpSession.setAttribute("code", text);
5 }
```

Aside from the editor and REPL there is one more window on the webpage. It is dedicated for tutorial and short documentation. While it does not enhance the functionality of the website per se, it plays an important role. The Solomonoff compiler is a very niche and specialised tool. There are no similar tools and any user coming to the website is not expected to be familiar with its usage. The primary purpose of the website is not to be a replacement for user's IDE and terminal. Instead it serves as an all-in-one introductory tutorial, interactive playground and a marketing campaign. We want to make the learning materials easily accessible and abundant. Building a strong community is the back-bone of every open-source project.