



OSTBAYERISCHE
TECHNISCHE HOCHSCHULE
REGENSBURG

Building Domain Specific Languages and Polyglot Applications with GraalVM

Presented to the Faculty of Computer Science and Mathematics
University of Applied Sciences Regensburg
Study Programme:
Master Computer Science

Master Thesis

In Partial Fulfillment of the Requirements for the Degree of
Master of Science (M.Sc.)

Presented by: Christian Paling
Student Number: 3213285

Primary Supervising Professor: Prof. Dr. Michael Bulenda
Secondary Supervising Professor: Prof. Dr. Carsten Kern

Submission Date: January 15, 2021

THESIS DECLARATION

ABSTRACT

Table of Contents

| | | |
|----------|--|-----------|
| 1 | Introduction | 1 |
| 2 | Domain Specific Languages | 2 |
| 2.1 | Definition of Domain Specific Languages | 2 |
| 2.2 | Benefits and Problems of Domain Specific Languages | 4 |
| 2.3 | Implementation of Domain Specific Languages | 6 |
| 2.3.1 | Internal Domain Specific Languages | 6 |
| 2.3.2 | External Domain Specific Languages | 9 |
| 3 | Overview of GraalVM | 20 |
| 3.1 | Motivation | 20 |
| 3.2 | Features | 20 |
| 3.2.1 | GraalVM Compiler | 20 |
| 3.2.2 | Native Images | 20 |
| 3.2.3 | Truffle Framework | 21 |
| 3.2.4 | Polyglot Applications | 21 |
| 4 | Domain Specific Languages in GraalVM | 22 |
| 4.1 | Technical Overview | 22 |
| 4.2 | <INSERT NAME OF DSL> | 22 |
| 4.3 | Implementation of <INSERT NAME OF DSL> | 22 |
| 4.4 | Evaluation | 22 |
| 5 | Integration of Domain Specific Languages | 23 |
| 5.1 | Technical Overview | 23 |
| 5.2 | Integration of <INSERT NAME OF DSL> | 23 |
| 5.3 | Evaluation | 23 |
| 6 | Conclusion | 24 |
| A | Completion of Code Listings | 26 |
| A.1 | Internal Timer DSL | 26 |
| A.2 | External Timer DSL | 28 |

1 Introduction

Business as usual

2 Domain Specific Languages

Before diving into the technical details and the implementation of DSLs on top of GraalVM, some background information is necessary to lay a foundation for the upcoming chapters of this thesis. First, the term *domain specific language* is properly defined and a distinction between different types of DSLs is made. Afterwards, benefits as well as problems of DSLs and the usage thereof are discussed. Lastly, basic techniques and examples of how to implement the different types of DSLs are introduced. The contents of this chapter are based heavily on *Domain-Specific Languages* by Martin Fowler [1] and *Crafting Interpreters* by Bob Nystrom [2] which can be consulted for further information about the creation of DSLs or custom languages in general.

2.1 Definition of Domain Specific Languages

To establish boundaries to a term with a generally vague meaning, Martin Fowler defines DSLs as follows [1]:

Domain specific language: a computer programming language of limited expressiveness focused on a particular domain.

A DSL is therefore characterized by first being *a computer programming language*. Its primary usage is to allow humans to instruct a computer to perform a certain action. Contrary to a *general-purpose language* like Java or Ruby, however, a DSL only has a *limited expressiveness* and is specialized on a *particular domain*. In other words, a DSL only supports a small amount of features and syntax which are tailored to the domain where it should be employed.

Fowler furthermore distinguishes DSLs into three categories [1]:

- **External DSLs** are separate from the main language of the application and usually have a custom syntax. They therefore have to be parsed by the host application in order to execute them.
- **Internal DSLs** use capabilities of the general-purpose language of the application to try to offer the feeling of a custom language. The code of the DSL is valid code in its general-purpose language as well, so no additional parsing is necessary.
- **Language Workbenches** offer environments for defining and building DSLs as well as writing scripts for the DSLs. Since language workbenches do not play any role in this thesis, they will not be given further attention.

For all these types of DSLs, the boundary which determines whether something is or is not a DSL is quite blurry. For internal DSLs, the distinction has to be made between a normal *application programming interface* (API) and an actual internal DSL. For Fowler [1] the difference lies in the nature of a DSL to define a new language in form of a grammar. The documentation of an API can offer a good indication whether the module or library exposes a normal API or a DSL. In the case of APIs, methods usually can be documented by themselves and therefore have a self-sufficient meaning. In a DSL, however, methods usually do not hold any meaning by themselves but can only be interpreted in context of a larger expression.

Listing 1 depicts a testing library offered by the Spring framework to check the behaviour of a RESTful backend. The testing library offers a variety of static methods combined with

```
apiTestClient.perform(get("/users"))
    .andDo(print())
    .andExpect(status().isOk());
```

Listing 1: The Spring framework offers internal DSLs for testing purposes.

elegant method chaining to fluently define a test. For instance, the *andDo* method expects an object that implements a *ResultHandler* interface. The static method *print* constructs such an instance and passes it to the *andDo* method. It is therefore apparent that a standalone executing of the *print* method would not result in anything meaningful. The *print* method as well as the *andDo* method can only be reasonably evaluated when they are both combined with each other.

```
public class PersonBuilder {
    private String name;
    private Integer age;
    private String placeOfBirth;

    public static PersonBuilder newPerson() {
        return new PersonBuilder();
    }

    public PersonBuilder name(String name) {
        this.name = name;
        return this;
    }

    public PersonBuilder age(Integer age) {
        this.age = age;
        return this;
    }

    public PersonBuilder placeOfBirth(String placeOfBirth) {
        this.placeOfBirth = placeOfBirth;
        return this;
    }

    public Person build() {
        return new Person(this.name, this.age, this.placeOfBirth);
    }
}

// Usage
PersonBuilder.newPerson()
    .name("John Doe")
    .age(21)
    .build();
```

Listing 2: It is arguable whether *PersonBuilder* can be considered to be a DSL.

On the other hand, listing 2 shows the definition and usage of a *builder pattern* to create

instances of a hypothetical *Person* class. In this case it is arguable whether *PersonBuilder* exposes an internal DSL. Each method of the builder such as *name* or *age* can be independently described by setting an attribute of the resulting person, i.e. each method has a self-sufficient meaning by itself. Additionally, except having to call *newPerson* at the beginning and *build* at the end, the creation of a new person is not dependent on any grammar which an actual language should be composed of.

For external DSLs one has to differentiate between a DSL and a general-purpose language, though the boundary is not as blurry as with internal DSLs. A good example presented by Martin Fowler [1] is the *R language*¹, which is a programming language for statistical computing. While focusing on a particular domain, the R language is not limited in its expressiveness and can be employed for purposes it was not initially intended for. Therefore, though it partly complies with the definition of a DSL, it should be categorized as a general-purpose language. A popular and widely spread example for an external DSL is *regular expressions*. It is specialized on matching text and offers only the amount of features and syntax to excel for this purpose. As a general rule, external DSLs are not *Turing-complete*. They usually do not offer mechanisms for control flow such as loops or conditions combined with the possibility to define variables and functions.

2.2 Benefits and Problems of Domain Specific Languages

After defining and categorizing DSLs, the question arises why developers of software systems should actually build and use DSLs. What are potential benefits as well as problems of DSLs? By weighing each of the advantages and downsides of DSLs, software professionals will be able to decide whether or not a DSL could potentially help to solve a certain problem.

According to Martin Fowler, DSLs offer the following advantages [1]:

- **Improving Development Productivity:** Since DSLs are specialized to express a certain aspect of a system, the code of the DSL will be more easy to write, read, and understand. This leads to an improvement of productivity by both making less mistakes as well as fixing defects more quickly. In Fowler's words: "The limited expressiveness of DSLs makes it harder to say wrong things and easier to see when you've made an error."
- **Communication with Domain Experts:** Good communication in software projects is, according to research, a very important critical success factor for projects to succeed [3]. Since software professionals develop systems for a wide variety of industries, they have often to be in contact with experts of the particular industry, so called *domain experts*. Due to their specialized syntax, DSLs offer the possibility for domain experts to read and correct source code and therefore highly improve the communication between tech and non-tech project stakeholders.
- **Change in Execution Context:** A common reason for external configuration files written in XML and similar formats, is the ability to read and evaluate them at runtime. This way the system does not have to be recompiled for every change of its configuration. DSLs offer a resembling advantage: they can also shift changes of logic of a system from compile time to the execution of that system.

¹<https://www.r-project.org/>

- **Alternative Computational Model:** Most general-purpose languages follow the *imperative style* of computation: the computer is told what to do in a certain sequence with features such as control flow and variables. For some problems, however, different approaches are more suitable and easier to utilize. Build automation is one of these problems: build tools such as *Apache Maven*² generally offer a *declarative style* to describe the build of a software system. Instead of focusing on *how* something should be done, the declarative style of programming concentrates on *what* should happen, leaving the *how* to a different layer of the system. When creating and using DSLs, it is also possible to employ a different computational model than the main language of the application with which it is easier to express or define certain aspects of the respective domain.

Contrary to these advantages, the usage of DSLs also comes with some problems and threats. Among them are the following [1]:

- **Language Cacophony:** This term was coined by Martin Fowler and states that learning new languages is generally hard. Therefore, it is apparent that combining multiple language for a project complicates the development compared to only using a single language. It is therefore necessary to determine whether or not learning a DSL is less costly opposed to understanding and working on the problems at hand without a DSL.
- **Cost of Building:** The most obvious problem of creating a DSL is the initial cost of building it. However, not only the initial costs of implementing the DSL has to be taken into account. Throughout time the DSL will have to be maintained and extended as well. Moreover, according to Fowler, it is not common for developers to know the techniques which are necessary to build DSLs which further aggravates the cost of implementing one.
- **Ghetto Language:** With the *ghetto language problem*, Martin Fowler refers to an issue which contrasts with the language cacophony problem. The term describes a language, built in-house, which is being utilized in more and more systems of the company as well as being continually extended with features and therefore slowly evolving into a general-purpose language. In the long run, this will lead the company to be inflexible regarding technological innovations and shifts in the industry as well as making it harder to hire staff. As a consequence, companies should clearly define the purpose and boundaries of their DSL and refrain from breaching these decisions.
- **Blinkered Abstraction:** The last problem Fowler highlights is the situation where developers are too confident about their DSL and try to fit the world to work with their language, instead of changing the language in accordance to the world. Thus, software professionals must view their DSL to be constantly under development, instead of regarding it as being finished.

As a conclusion, there are two possible reasons not to use a DSL. First, in case none of the benefits of a DSL applies to the problem at hand it is naturally not a fitting tool to solve that problem. Secondly, if the costs and risks of building a DSL outweigh its potential benefits. Otherwise it can be worthwhile to consider building or using a DSL to benefit from the potential prospects as set out in this section.

²<https://maven.apache.org/>

2.3 Implementation of Domain Specific Languages

In order to compare and evaluate the implementation of DSLs with the frameworks offered by GraalVM, an overview of how DSLs can be built without additional technologies is necessary. The following section explains how internal and external DSLs can be implemented. For both types, a language for the same and rather simple problem will be built. The Java SDK ships with a powerful timer facility to schedule tasks for future and recurring execution. A *TimerTask* defines such a task which can be run once or repeatedly in the future. Listing 3 displays how a *TimerTask* can be created and scheduled. In this example, the string *Hello World* will be printed periodically every 1000 milliseconds with a delay of 5000 milliseconds. If the last parameter is omitted, *Hello World* would be only printed once after 5000 milliseconds have elapsed.

```
var timer = new Timer();

timer.schedule(new TimerTask() {
    @Override
    public void run() {
        System.out.println("Hello World");
    }
}, 5000, 1000);
```

Listing 3: After five seconds print *Hello World* every second by using a *TimerTask*.

The internal and external DSLs which will be presented in the further course of this section will serve as a layer on top of this API and will enable developers to schedule tasks in a more fluent manner. The primary objective of both upcoming DSLs, however, is to illustrate prevalent approaches to implement both types of DSLs.

When creating DSLs, a common strategy described by Fowler [1] is to first write some code to exemplify how the DSL should look like. Using these examples, the developer can iteratively verify whether the abilities of the DSL already fulfil all requirements or whether the DSL has to be modified and adapted. After the design of the language is set, it can be implemented incrementally, feature after feature, to its intended form.

2.3.1 Internal Domain Specific Languages

Internal DSLs are generally more approachable than external DSLs due to the fact that external DSLs require more techniques such as grammars and parsers in order to build them. On the flip side, internal DSLs are largely constrained by their host language. There are general-purpose languages such as Ruby or Lisp which are very flexible regarding their syntax or offer specialized functionalities, such as macros in Lisp, to create custom languages. Other programming languages like Java or C++ have more restrictive syntactic rules in comparison which affects the look and feel of internal DSLs.

To build and structure internal DSLs different approaches exist and are employed. However, since this thesis covers GraalVM, a technology based on Java, a common way to build internal DSLs using *object-oriented programming* (OOP) will be illustrated. To create internal DSLs using an OOP host language, Martin Fowler argues [1] that the DSL itself and the actual objects which the DSL utilizes should be separate from each other. Internal DSLs should be built in

form of so called *expression builders* which should not define any domain logic but only offer constructs to build expressions of the DSL. The actual logic should be located in another layer hidden behind the expression builder which the builder utilizes once the DSL expression should be executed. This approach enables separate testing of the domain logic and the expression builder as well as the possibility to replace the expression builder with an external DSL if necessary. In the context of the timer scheduling DSL, the Java timer API represents the layer of the domain logic while a separate layer of expression builders has to be implemented.

As previously mentioned, the first step of building a DSL is to write some example code. Listing 4 depicts how the internal timer scheduling DSL should look like. The timer itself is configured using an API similar to a builder pattern while static methods act as descriptive parameters, like setting what the timer should execute or the delay of the timer.

```
timer()
    .execute(print("Hello World repeatedly!"))
    .repeatedly()
    .every(minutes(1))
    .after(seconds(30))
    .setup();

timer()
    .execute(print("Hello World once!"))
    .once()
    .after(seconds(10))
    .setup();

timer()
    .execute(print("Hello World once now!"))
    .once()
    .rightNow()
    .setup();
```

Listing 4: Some expressions to schedule future and potentially periodic tasks.

Because the static methods for the different units of time and for the timer tasks provide the more simpler functionalities of the DSL, they will be attended to first. Listing 5 and 6 depict two classes which are structured in a similar fashion. Both classes are final and therefore cannot and should not be extended. Furthermore, both have private constructors to prohibit the creation of instances of both classes. The implementation of the *Duration* class is self-explanatory and converts different units of time to milliseconds, since the Java SDK expects milliseconds for the scheduling of timers. Static methods of the *Tasks* class should create instances of the *TimerTask* class offered by the Java SDK which will be scheduled and executed after the configuration of the timer has completed. In this example only a simple *print* task exists, though more complex tasks like syncing databases or sending emails would be possible.

The method chaining with which the timer is constructed is built using separate classes. Each class offers the developer one or more possibilities to configure the timer and returns an instance of a new class which defines the next step of configuration. Each step therefore acquires a part of the configuration and passes it on to the next step. In the final step and class, all the obtained information is used to configure and schedule an actual timer using the Java API. The first class in this hierarchy is shown in listing 7. It offers the static *timer* method

```
public final class Duration {
    private Duration() {}

    public static long seconds(long n) {
        return n * 1000;
    }

    public static long minutes(long n) {
        return seconds(60 * n);
    }

    public static long hours(long n) {
        return minutes(60 * n);
    }
}
```

Listing 5: *Duration* offers static methods for different units of time.

```
public final class Tasks {
    private Tasks() {}

    public static TimerTask print(String message) {
        return new TimerTask() {
            @Override
            public void run() {
                System.out.println(message);
            }
        };
    }
}
```

Listing 6: *Tasks* offers static methods for different timer tasks, here only a print task.

which was the initial method with which each DSL expression has to start according to the language design of listing 4. This method creates the actual instance of the builder class which only possesses one instance method called *execute*. Since *execute* expects an instance of type *TimerTask*, it fits perfectly to the static methods of the *Tasks* class from listing 6 which should return predefined objects of type *TimerTask*.

The *execute* method creates an instance of another class called *TimerExpressionBuilderWithTask* which is displayed in listing 8 and defines the next possible steps of the timer configuration. The developer can choose between either calling *repeatedly* or *once* which both create different subsequent objects to differentiate between a timer task that should be executed only once and one that should be run multiple times.

Since each step of the DSL is in a separate class, the type system makes it impossible to create invalid DSL expressions. If all methods would be defined in a single class, a developer could potentially call the methods *once* and *repeatedly* after each other which would result in ambiguous code. Furthermore, considering that code completion is offered by nearly every *integrated development environment* nowadays, the developer is piloted through the creation of the expression, since the code completion will only offer the next methods according to the

```
public final class TimerExpressionBuilder {  
    private TimerExpressionBuilder() {}  
  
    public static TimerExpressionBuilder timer() {  
        return new TimerExpressionBuilder();  
    }  
  
    public TimerExpressionBuilderWithTask execute(TimerTask task) {  
        return new TimerExpressionBuilderWithTask(task);  
    }  
}
```

Listing 7: *TimerExpressionBuilder* defines the starting point of the DSL.

```
public final class TimerExpressionBuilderWithTask {  
    private final TimerTask task;  
  
    public TimerExpressionBuilderWithTask(TimerTask task) {  
        this.task = task;  
    }  
  
    public RepeatableTimerExpressionBuilder repeatedly() {  
        return new RepeatableTimerExpressionBuilder(this.task);  
    }  
  
    public SingleTimerExpressionBuilder once() {  
        return new SingleTimerExpressionBuilder(this.task);  
    }  
}
```

Listing 8: *TimerExpressionBuilderWithTask* marks the next step of configuration of the timer.

hierarchy of the expression builder classes.

All remaining steps and expression builder classes follow a similar structure and can be viewed in listing 26 and 27 of the appendix.

2.3.2 External Domain Specific Languages

External DSLs compared to internal ones come with a much greater syntactic freedom. This liberality concerning the syntax, however, goes along with a more complex implementation. The basic principles with which external DSLs are build are very similar to the ones of general-purpose languages, though developers of DSLs do not have to know the techniques as in depth as general-purpose language developers. Interestingly, according to Bob Nystrom [2], the techniques with which languages are build have not really changed since the early days of computing.

Before explaining the approach with which the external DSL for scheduling timers is implemented, the structure and syntax of the intended language will be presented first. Listing 9 presents some example code of the external DSL. It is apparent that the syntax of the DSL does not follow the syntactic rules of Java anymore. Timers are grouped in *timer* and *end* pairs and allow the same configurable features as with the internal DSL.

```
timer
  print "Hello World"
  repeatedly
  every 30 seconds
  after 2 minutes
end

timer
  print "Hello World once!"
  once
  after 10 seconds
end

timer
  print "Hello World now!"
  once
  right now
end
```

Listing 9: Some external DSL expressions to schedule future and potentially periodic tasks.

To build this DSL, a process based on Bob Nystrom’s online book *Crafting Interpreters* [2] was employed. The book uses widespread techniques to build languages which are also highlighted in Fowler’s work about DSLs [1]. This process divides the evaluation of language expressions into at least three steps.

The first step is called *lexing*. A *lexer* takes the code of the language and splits it into individual tokens. A token is a data structure which is associated to a certain type and might contain a value. Listing 10 lists all types of tokens of the DSL as an enum. Every keyword is a different token type, in addition to the two datatypes which the DSL supports: strings and numbers. Lastly, an *EOF* token type marks the end of the source code.

```
public enum TokenType {
    TIMER, REPEATEDLY, ONCE, RIGHT, NOW,
    PRINT, AFTER, EVERY, STRING, NUMBER,
    SECONDS, MINUTES, HOURS, END, EOF
}
```

Listing 10: All types of tokens of the DSL.

The token itself is a simple class with, as previously mentioned, attributes for the type of the token and the value. It is presented in listing 11. Note that the value will be *null* for most types of tokens except strings and numbers since keywords do not hold any literal values.

The lexer moves character by character through the source code, tries to identify tokens and stores them in a list, and in the end returns that list of tokens. Listing 12 depicts the basic structure of such a lexer. The attributes include the start position of the current read as well as the end position, the source code itself, and the list of tokens which will be returned in the end.

As long as the lexer has not reached the end of the source code, i.e. the start position is

```
public class Token {
    private final TokenType type;
    private final Object value;

    public Token(TokenType type, Object value) {
        this.type = type;
        this.value = value;
    }

    public TokenType getType() {
        return type;
    }

    public Object getValue() {
        return value;
    }
}
```

Listing 11: The *Token* class for the lexer.

```
public class Lexer {
    private int startOfToken = 0;
    private int endOfToken = 0;
    private final String code;
    private final List<Token> tokens = new ArrayList<>();

    public Lexer(String code) {
        this.code = code;
    }

    public List<Token> getTokens() throws TimerDSLException {
        while (!isAtEnd()) {
            readNextToken();
            this.startOfToken = this.endOfToken + 1;
            this.endOfToken = this.startOfToken;
        }

        tokens.add(new Token(EOF, null));
        return tokens;
    }
}
```

Listing 12: Basic structure of the *Lexer* class.

greater than the length of the source code, the lexer tries to read the next token. Listing 13 illustrates how the lexer identifies the next token. By comparing the character of the current position, the lexer can judge what it will expect as a next token. If for example the current character is a double quote, the lexer can assume that the next token should be a string.

After the decision has been made regarding the expectation of the next token, the lexer tries to find the end of this token. Listing 14 shows how this is accomplished for strings.

With the help of the peek method which returns the character of the current end position,

```
private void readNextToken() throws TimerDSLEXception {
    var nextChar = code.charAt(this.startOfToken);

    if (List.of(' ', '\r', '\t', '\n').contains(nextChar)) {
        // Ignore whitespaces
    } else if ( '"' == nextChar) {
        string();
    } else if (isDigit(nextChar)) {
        number();
    } else if (isAlpha(nextChar)) {
        keyword();
    } else {
        throw new TimerDSLEXception("Unexpected character");
    }
}
```

Listing 13: The lexer identifies the next token by checking the first character of the next token.

```
private void string() throws TimerDSLEXception {
    endOfToken++;
    while (peek() != '"' && !isAtEnd()) endOfToken++;

    if (isAtEnd()) throw new TimerDSLEXception("Unterminated string");

    endOfToken++;
    var value = code.substring(startOfToken + 1, endOfToken - 1);
    tokens.add(new Token(STRING, value));
}
```

Listing 14: The lexer tries to find the end of the string to then get the value between the start and end position.

the lexer is able to find the end of the string by searching for the second double quote. In case it reaches the end of the source code before finding the second double quote, the lexer throws an exception, otherwise the value of the string is extracted from the source code and saved as a string token in the list of tokens.

The approach for identifying numbers or keywords is using a very similar approach and can be viewed in the complete definition of the lexer class in listing 28 and 29 of the appendix.

In the second step of the whole evaluation, a *parser* takes this list of tokens to generate an *abstract syntax tree* (AST) according to the grammatical rules of the language. The grammar is generally a *context-free grammar* (CFG) which is often notated in a flavour of the *Backus-Naur form* (BNF). Listing 15 illustrates how a grammar could be defined using a version of the BNF which Bob Nystrom uses in his work [2].

A CFG has *terminals* and *nonterminals*. A terminal is like a literal value of the grammar, for example *mozzarella cheese* or *mushrooms*. Terminals mark end points and cannot be replaced with more symbols. Nonterminals on the other hand are references to other rules which allow the construction of more complex expressions. The *pizza* nonterminal is the starting point of the grammar with a *crust* nonterminal at the beginning. The *crust* nonterminal offers two

```

pizza  → crust "with" cheese "and" (topping "and" | topping)+
crust  → "thin crust" | "thick crust"
cheese → "mozzarella cheese" | "provolone cheese"
topping → "mushrooms" | "extra cheese" | "salami" | "ham"

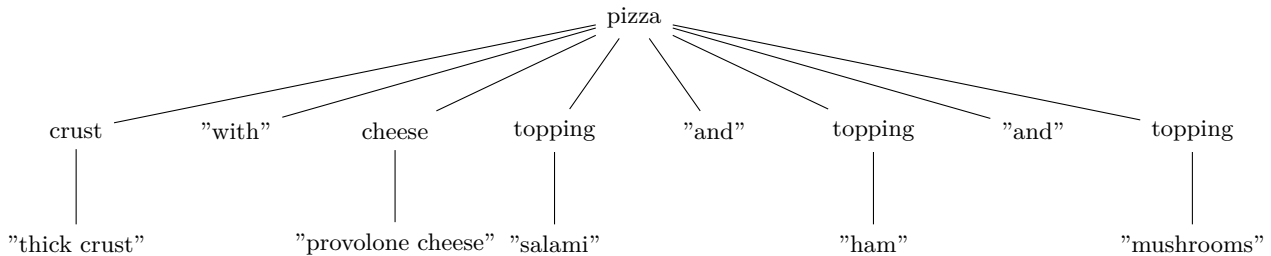
```

Listing 15: A simple grammar for configuring pizzas

possible terminals (specified by the `|` sign): either a *thin crust* or a *thick crust*. At the end of the pizza nonterminal there are again two possibilities. It is either allowed to choose a topping combined with an *and* terminal (to be able to have multiple toppings) or just a single topping. The `+` sign specifies the same as with regular expressions. It marks that a certain rule can occur once or more times while a `*` would indicate that a rule can be utilized zero or more times. The parentheses group these possibilities regarding the toppings together and signify that the `+` sign can only be applied to the toppings. This way an arbitrary amount of toppings is possible. The following sentences would be valid according to the grammar:

- thin crust with mozzarella cheese and mushrooms
- thick crust with provolone cheese and salami and ham and mushrooms

With the help of a grammar, it is also possible to represent an expression in form of a tree, the AST. Figure 1 visualizes the second sentence from above in the form of an AST which conforms to the defined grammar.

Figure 1: *thick crust with provolone cheese and salami and ham and mushrooms* represented as an AST

The main task of the parser in the process presented by Nystrom [2] is to build an AST representation of the tokens for easier future processing. To understand how such a parser can be build, the implementation of a parser for the timer scheduling DSL will be subsequently illustrated. Listing 16 depicts a possible grammar for the DSL (as seen in listing 9) in BNF.

A program written in the DSL consists of one or more *timer statements*. Each *timer statement* has to start with the terminal *timer* and has to end with the terminal *end*. Between *timer* and *end*, the first expected nonterminal is the command. Currently only the *print* command is supported which expects a string. After the command, two different possibilities exist to configure the timer: a *once timer* and a *repeated timer*. The *once timer* only expects a configuration for the delay of the command while the *repeated timer* expects the configuration of the period of the command in addition.

| | |
|---------------------|---|
| program | → timer_stmt+ |
| timer_stmt | → "timer" command (once_timer repeated_timer) "end" |
| command | → "print" STRING |
| once_timer | → "once" after_configuration |
| repeated_timer | → "repeatedly" "every" NUMBER time_unit after_configuration |
| after_configuration | → "right" "now" "after" NUMBER time_unit |
| time_unit | → "seconds" "minutes" "hours" |

Listing 16: The grammar of the external timer DSL in BNF.

The implementation is surprisingly simple, once well understood. The first step is to define the AST datastructure. Listing 17 shows the root element of the tree: a *timer statement*. The class has two attributes which resemble the children of the root: a *command* and the configuration of the timer.

```
public class TimerStmt {
    private final Command command;
    private final TimerConfiguration configuration;

    public TimerStmt(Command command, TimerConfiguration configuration) {
        this.command = command;
        this.configuration = configuration;
    }

    public Command getCommand() {
        return command;
    }

    public TimerConfiguration getConfiguration() {
        return configuration;
    }
}
```

Listing 17: The root element of the AST.

Since for the purposes of this example only a *print* command is supported, the command class is rather simple, although it is laid out to be extended at will. Listing 18 depicts the *Command* class which is abstract and which includes the *print* command as a static nested class. Naturally, the *print* command only has one "child" which is the message it should print.

The *timer configuration* has a resembling structure and is presented in listing 30 in the appendix for the sake of completion. The *TimerConfiguration* class itself is again abstract but has different subclasses. Analogous to the grammar, a timer configuration is either a *once timer* or a *repeated timer*. The *once timer* has only a time setting for the delay, while the *repeated timer* has an additional time setting for the period.

It is noticable that the composition of the AST is very similar to the composition of the grammar. This is due to the AST being a representation of the syntactic structure of the code, as previously mentioned. The question now arises, however, how the AST of some concrete code can actually be constructed. To address this problem, Bob Nystrom presents a popular

```
public abstract class Command {  
    public static class PrintCommand extends Command {  
        private final String message;  
  
        public PrintCommand(String message) {  
            this.message = message;  
        }  
  
        public String getMessage() {  
            return message;  
        }  
    }  
}
```

Listing 18: All commands are subclasses of the *Command* class.

technique in his work [2] which is called *recursive descent*. In simple words, recursive descent parsing is a translation of the grammar into programming language code. Many of today's programming language implementations are based on the recursive descent parsing technique, such as the GCC or the Roslyn C# compiler [2].

As presented in listing 19, the parser for the timer DSL has only two attributes: the list of tokens and the current position of the parser in this aforementioned list.

```
public class Parser {  
    private final List<Token> tokens;  
    private int current;  
  
    public Parser(List<Token> tokens) {  
        this.tokens = tokens;  
        this.current = 0;  
    }  
}
```

Listing 19: The basic structure of the timer DSL parser.

As was mentioned, the recursive descent technique is a translation of the grammar into code. The first rule of the grammar specifies that a program consists of one or more timer statements. Therefore the method with which the parser will be called has to reflect this rule, as shown in listing 20.

First, a new list of timer statements, which are the root nodes of the AST, is created. The same list will be returned at the end of the method. Afterwards, since the rule expects at least one timer statement, the code also adds at least one element to that list. Subsequently, additional timer statements are added to the list until the end of the list is reached, i.e. an *EOF* token is encountered.

The *timerStmt* method corresponds to the next rule in the grammar and is outlined in listing 21.

Throughout the parser there are two helpful methods: *consume* and *match*. The *consume* method expects a token of a certain type at the current position. In case the type of the current

```
// program → timer_stmt+
public List<TimerStmt> parse() throws TimerDSLEXception {
    var timerStatements = new ArrayList<TimerStmt>();
    timerStatements.add(timerStmt());

    while (!isAtEnd()) {
        timerStatements.add(timerStmt());
    }

    return timerStatements;
}
```

Listing 20: The initial method with which the parser will be called and which reflects the first rule of the grammar.

```
// timer_stmt → "timer" command (once_timer | repeated_timer) "end"
private TimerStmt timerStmt() throws TimerDSLEXception {
    consume(TIMER, "Expected 'timer' at the beginning of definition.");

    var command = command();

    TimerConfiguration config = null;

    if (match(ONCE)) config = onceTimer();
    else {
        consume(REPEATEDLY, "Expected 'once' or 'repeatedly' after command.");
        config = repeatedTimer();
    }

    consume(END, "Expected 'end' at the end of definition.");

    return new TimerStmt(command, config);
}
```

Listing 21: The *timerStmt* method which corresponds to the *timer_stmt* rule of the grammar.

token corresponds to this expected type, the token is returned and the *current* attribute of the parser incremented, if not then an exception with a given message is thrown. The *match* method, however, only returns a boolean which is true if the given type is equal to the type of the current token. It does not change the position of the parser inside the list of tokens. Therefore, in the first line of the *timerStmt* method a *TIMER* token is expected, since every timer statement has to start with the *timer* keyword. In case no *TIMER* token exists at that position, an exception is thrown with the message *Expected 'timer' at the beginning of definition*. Since the *command* nonterminal follows the *timer* keyword, the method calls a *command* method in the next step which handles the *command* rule. After the command, there are two possibilities: either a *once timer* or a *repeated timer* configuration. Since the once timer has to start with the *once* keyword, the method checks whether the current token is of type *ONCE*. If yes it calls the *onceTimer* method, otherwise it expects a *REPEATEDLY* token and calls the corresponding method. At the end of the timer statement, the *END* token

must be consumed and the whole timer statement is returned.

As a final example, listing 22 depicts the *command* method.

```
// command → "print" STRING
private Command command() throws TimerDSLEXception {
    consume(PRINT, "Expected 'print' command.");
    var message = consume(STRING, "Expected 'string' after 'print'.");
    return new PrintCommand((String) message.getValue());
}
```

Listing 22: The *command* method currently only has the print command as a possibility.

The method first expects a *PRINT* token, followed by a *STRING* token. In the case of the string, the returned token of the *consume* method is actually saved in a variable, to pass it to the *PrintCommand* AST element. The remaining methods of the recursive descent parser (see 31 in the appendix) for this DSL work very similar to the examples that were presented above. All methods of the parser correspond to one rule of the grammar. The parser then utilizes these methods to descent recursively according to the grammar to construct an AST in the end.

The final step of the processing of the DSL is to walk through the AST returned by the parser and interpret it. Since this DSL is rather simple, the *interpreter* is implemented using a very naive approach. The basic structure of the interpreter is visible in listing 23.

```
public class Interpreter {
    private List<TimerStmt> timerStatements;

    public Interpreter(List<TimerStmt> timerStatements) {
        this.timerStatements = timerStatements;
    }

    public void interpret() throws TimerDSLEXception {
        for(var stmt: timerStatements) {
            evaluate(stmt);
        }
    }
}
```

Listing 23: The basic structure of the interpreter.

The interpreter receives the list of statements, which was built by the parser, through its constructor. It has a public *interpret* method which iterates over each statement and evaluates it. Listing 24 shows the *evaluate* method which accepts a single timer statement and performs the actual evaluation.

The interpreter first "walks" to the *command* node of the statement to build an instance of the *TimerTask* class provided by the JDK. It then checks whether the configuration is a *once timer* or a *repeated timer* and schedules the timer using the remaining nodes of the AST. Listing 25 presents the remaining methods of the interpreter which are used by the *evaluate* method.

```
private void evaluate(TimerStmt stmt) throws TimerDSLException {
    var timer = new Timer();
    var timerTask = buildTask(stmt.getCommand());

    if (stmt.getConfiguration() instanceof TimerConfiguration.OnceTimer) {
        var onceTimer = (TimerConfiguration.OnceTimer) stmt.getConfiguration();
        timer.schedule(
            timerTask,
            getMillis(
                onceTimer.getAfterSetting().getNumber(),
                onceTimer.getAfterSetting().getUnit()
            )
        );
    } else {
        var repeatedTimer = (TimerConfiguration.RepeatedTimer) stmt.getConfiguration();
        timer.schedule(
            timerTask,
            getMillis(
                repeatedTimer.getAfterSetting().getNumber(),
                repeatedTimer.getAfterSetting().getUnit()
            ),
            getMillis(
                repeatedTimer.getEverySetting().getNumber(),
                repeatedTimer.getEverySetting().getUnit()
            )
        );
    }
}
```

Listing 24: The *evaluate* method with which a statement is evaluated by the timer.

It's noticable that the *instanceof* checks could make the code of the interpreter quite obscure if the DSL is much more complex. For this reason, Bob Nystrom presents the *visitor pattern* in his work [2] as a possibility to cleanly structure the interpreter without having to resort to *instanceof* checks when walking through the AST. However, in the case of this simple DSL, the visitor pattern would have been overkill as a solution.

Naturally, walking through the AST is not the only possibility to interpret the code of the DSL. Although no other techniques to write interpreters for external DSLs will be introduced in this section, Martin Fowler [1] highlights more possibilities which can be consulted for further information regarding this topic.

```
private TimerTask buildTask(Command command) throws TimerDSException {
    if (command instanceof Command.PrintCommand) {
        var message = ((Command.PrintCommand) command).getMessage();
        return new TimerTask() {
            @Override
            public void run() {
                System.out.println(message);
            }
        };
    } else throw new TimerDSException("Unknown command type");
}

private long getMillis(long number, TimerConfiguration.TimeUnit unit) {
    if (unit == TimerConfiguration.TimeUnit.SECONDS) {
        return number * 1000;
    } else if (unit == TimerConfiguration.TimeUnit.MINUTES) {
        return number * 1000 * 60;
    } else {
        return number * 1000 * 60 * 60;
    }
}
```

Listing 25: Remaining methods of the interpreter.

3 Overview of GraalVM

3.1 Motivation

Why do we need Graal?

Write more of Java in Java itself.

3.2 Features

3.2.1 GraalVM Compiler

Explanation and some benchmarks

source: <https://www.youtube.com/watch?v=sFf15TvSXZ0>

1. What is a JIT compiler

When compile Java Source with `javac` -> Java Bytecode At Runtime -> Bytecode is compiled in Machine Code Machine Code delivers usually much better perf

2. Why write JIT compiler in Java

C2 is the JIT compiler written in C++ Developers of JIT think C2 is too old now, too hard to maintain Developing in Java tends to be easier and more productive than in C++

3. JVM compiler interface

Allows to plugin a custom JIT compiler for the JVM written in Java In thesis interface can shown here: <https://github.com/openjdk/jdk/blob/master/src/jdk.internal.vm.ci/share/classes/jdk.vm.ci.runtime/src/jdk/vm/ci/runtime/JVMCICompiler.java> Takes bytecode and returns new bytecode

4. Graal JIT compiler process

The compiler first represents code in graphs Every node will then be transformed into machine code

5. Optimisations

- Basically changes the nodes
- Canonicalisation: e.g. `-x` -> `x`
- Global value numbering: remove redundant code: `(a + b) * (a + b)` -> only `(a + b)` once calculated
- Lock coarsening: two synchronized locks immediately after each other -> change to only once

3.2.2 Native Images

Explanation and some benchmarks yet again

3.2.3 Truffle Framework

Basic explanations: why is there a truffle framework and what is achievable

3.2.4 Polyglot Applications

Basic explanations: what's possible here

4 Domain Specific Languages in GraalVM

4.1 Technical Overview

How to build DSLs with GraalVM?

4.2 <INSERT NAME OF DSL>

Introduce the DSL of this thesis here 1 2 3 4 5

$$A = \{ x \mid x \in (A \cap B) \} \tag{1}$$

4.3 Implementation of <INSERT NAME OF DSL>

Highlight some key aspects of implementation

4.4 Evaluation

Evaluate the DSL and GraalVM, highlight pain points etc.

5 Integration of Domain Specific Languages

5.1 Technical Overview

How do polyglot applications technically work?

5.2 Integration of <INSERT NAME OF DSL>

Showcase how it's done using the thesis DSL

5.3 Evaluation

Evaluation how good this actually works

6 Conclusion

Business as usual

References

- [1] Martin Fowler. *Domain-Specific Languages*. Pearson Education, 2010.
- [2] Bob Nystrom. Crafting interpreters. <http://craftinginterpreters.com/>. Accessed: January 15, 2021.
- [3] Goparaju Purna Sudhakar. A model of critical success factors for software projects. *Journal of Enterprise Information Management*, 2012.

A Completion of Code Listings

A.1 Internal Timer DSL

```
public final class RepeatablTimerExpressionBuilder {
    private final TimerTask task;

    public RepeatablTimerExpressionBuilder(TimerTask task) {
        this.task = task;
    }

    public PeriodicRepeatablTimerExpressionBuilder every(long period) {
        return new PeriodicRepeatablTimerExpressionBuilder(this.task, period);
    }
}

public final class PeriodicRepeatablTimerExpressionBuilder {
    private final TimerTask task;
    private final long period;

    public PeriodicRepeatablTimerExpressionBuilder(TimerTask task, long period) {
        this.task = task;
        this.period = period;
    }

    public FinalizedRepeatablTimerExpressionBuilder rightNow() {
        return after(0);
    }

    public FinalizedRepeatablTimerExpressionBuilder after(long delay) {
        return new FinalizedRepeatablTimerExpressionBuilder(this.task, this.period, delay);
    }
}

public final class FinalizedRepeatablTimerExpressionBuilder {
    private final TimerTask task;
    private final long period;
    private final long delay;

    public FinalizedRepeatablTimerExpressionBuilder(TimerTask task, long period, long delay) {
        this.task = task;
        this.period = period;
        this.delay = delay;
    }

    public void setup() {
        var timer = new Timer();
        timer.schedule(this.task, this.delay, this.period);
    }
}
```

Listing 26: All remaining classes to define a periodic timer task.

```
public static class SingleTimerExpressionBuilder {
    private final TimerTask task;

    public SingleTimerExpressionBuilder(TimerTask task) {
        this.task = task;
    }

    public FinalizedSingleTimerExpressionBuilder rightNow() {
        return after(0);
    }

    public FinalizedSingleTimerExpressionBuilder after(long delay) {
        return new FinalizedSingleTimerExpressionBuilder(this.task, delay);
    }
}

public static class FinalizedSingleTimerExpressionBuilder {
    private final TimerTask task;
    private final long delay;

    public FinalizedSingleTimerExpressionBuilder(TimerTask task, long delay) {
        this.task = task;
        this.delay = delay;
    }

    public void setup() {
        var timer = new Timer();
        timer.schedule(this.task, delay);
    }
}
```

Listing 27: All remaining classes to define a single timer task.

A.2 External Timer DSL

```
public class Lexer {
    private static final Map<String, TokenType> KEYWORDS = new HashMap<>();
    private int startOfToken = 0;
    private int endOfToken = 0;
    private final String code;
    private final List<Token> tokens = new ArrayList<>();

    static {
        KEYWORDS.putAll(Map.of(
            "timer", TIMER, "print", PRINT, "repeatedly", REPEATEDLY, "once", ONCE,
            "every", EVERY, "after", AFTER, "seconds", SECONDS, "minutes", MINUTES,
            "hours", HOURS, "right", RIGHT
        ));
        KEYWORDS.putAll(Map.of(
            "now", NOW, "end", END
        ));
    }

    public Lexer(String code) {
        this.code = code;
    }

    public List<Token> getTokens() throws TimerDSLException {
        while (!isAtEnd()) {
            readNextToken();
            this.startOfToken = this.endOfToken + 1;
            this.endOfToken = this.startOfToken;
        }

        tokens.add(new Token EOF, null);
        return tokens;
    }

    private void readNextToken() throws TimerDSLException {
        var nextChar = code.charAt(this.startOfToken);

        if (List.of(' ', '\r', '\t', '\n').contains(nextChar)) {
            // do nothing
        } else if ("'" == nextChar) {
            string();
        } else if (isDigit(nextChar)) {
            number();
        } else if (isAlpha(nextChar)) {
            keyword();
        } else {
            throw new TimerDSLException("Unexpected character");
        }
    }

    // Continues on the next page
}
```

Listing 28: The whole lexer class of the external timer scheduling DSL.


```
private void string() throws TimerDSLException {
    endOfToken++;
    while (peek() != '"' && !isAtEnd()) endOfToken++;

    if (isAtEnd()) throw new TimerDSLException("Unterminated string");

    endOfToken++;
    var value = code.substring(startOfToken + 1, endOfToken - 1);
    tokens.add(new Token(String, value));
}

private void number() {
    while (isDigit(peek())) endOfToken++;
    tokens.add(
        new Token(NUMBER, Integer.parseInt(code.substring(startOfToken, endOfToken)))
    );
}

private void keyword() throws TimerDSLException {
    while (isAlpha(peek())) endOfToken++;
    var text = code.substring(startOfToken, endOfToken);

    if (KEYWORDS.containsKey(text))
        tokens.add(new Token(KEYWORDS.get(text), null));
    else
        throw new TimerDSLException("Unexpected keyword.");
}

private char peek() {
    return code.charAt(endOfToken);
}

private boolean isDigit(char c) {
    return c >= '0' && c <= '9';
}

private boolean isAlpha(char c) {
    return c >= 'a' && c <= 'z';
}

private boolean isAtEnd() {
    return startOfToken >= code.length();
}
}
```

Listing 29: The whole lexer class of the external timer scheduling DSL (continuation).

```
public abstract class TimerConfiguration {
    public enum TimeUnit {
        SECONDS, MINUTES, HOURS
    }

    public static class OnceTimer extends TimerConfiguration {
        private final TimeSetting afterSetting;

        public OnceTimer(TimeSetting afterSetting) {
            this.afterSetting = afterSetting;
        }

        public TimeSetting getAfterSetting() {
            return afterSetting;
        }
    }

    public static class RepeatedTimer extends TimerConfiguration {
        private final TimeSetting everySetting;
        private final TimeSetting afterSetting;

        public RepeatedTimer(TimeSetting everySetting, TimeSetting afterSetting) {
            this.everySetting = everySetting;
            this.afterSetting = afterSetting;
        }

        public TimeSetting getEverySetting() {
            return everySetting;
        }

        public TimeSetting getAfterSetting() {
            return afterSetting;
        }
    }

    public static class TimeSetting {
        private final long number;
        private final TimeUnit unit;

        public TimeSetting(long number, TimeUnit unit) {
            this.number = number;
            this.unit = unit;
        }

        public long getNumber() {
            return number;
        }

        public TimeUnit getUnit() {
            return unit;
        }
    }
}
```

Listing 30: The timer configuration classes of the AST.

```
public class Parser {
    private final List<Token> tokens;
    private int current;

    public Parser(List<Token> tokens) {
        this.tokens = tokens;
        this.current = 0;
    }

    // program → timer_stmt+
    public List<TimerStmt> parse() throws TimerDSLEXception {
        var timerStatements = new ArrayList<TimerStmt>();
        timerStatements.add(timerStmt());

        while (!isAtEnd()) {
            timerStatements.add(timerStmt());
        }

        return timerStatements;
    }

    // timer_stmt → "timer" command (once_timer | repeated_timer) "end"
    private TimerStmt timerStmt() throws TimerDSLEXception {
        consume(TIMER, "Expected 'timer' at the beginning of definition.");

        var command = command();

        TimerConfiguration config = null;

        if (match(ONCE)) config = onceTimer();
        else {
            consume(REPEATEDLY, "Expected 'once' or 'repeatedly' after command.");
            config = repeatedTimer();
        }

        consume(END, "Expected 'end' at the end of definition.");

        return new TimerStmt(command, config);
    }

    // command → "print" STRING
    private Command command() throws TimerDSLEXception {
        consume(PRINT, "Expected 'print' command.");
        var message = consume(STRING, "Expected 'string' after 'print'.");
        return new PrintCommand((String) message.getValue());
    }

    // once_timer → "once" after_configuration
    private OnceTimer onceTimer() throws TimerDSLEXception {
        current++;
        return new OnceTimer(afterConfig());
    }

    // Continues on the next page
}
```

Listing 31: The complete recursive descent parser.

```
// repeated_timer → "repeatedly" "every" NUMBER time_unit after_configuration
private RepeatedTimer repeatedTimer() throws TimerDSLEException {
    consume(EVERY, "Expected 'every' after 'repeatedly'.");
    var number = consume(NUMBER, "Expected 'number' after 'every'.");
    return new RepeatedTimer(
        new TimeSetting(Long.valueOf((Integer) number.getValue()), timeUnit()),
        afterConfig()
    );
}

// after_configuration → "right" "now" | "after" NUMBER time_unit
private TimeSetting afterConfig() throws TimerDSLEException {
    if (match(RIGHT)) {
        current++;
        consume(NOW, "Expected 'now' after 'right'.");
        return new TimeSetting(0, TimeUnit.SECONDS);
    } else {
        consume(AFTER, "Expected 'right now' or 'after' as a time setting.");
        var number = consume(NUMBER, "Expected 'number' after 'after'.");
        return new TimeSetting(Long.valueOf((Integer) number.getValue()), timeUnit());
    }
}

// time_unit → "seconds" | "minutes" | "hours"
private TimeUnit timeUnit() throws TimerDSLEException {
    if (match(SECONDS)) {
        current++;
        return TimeUnit.SECONDS;
    } else if (match(MINUTES)) {
        current++;
        return TimeUnit.MINUTES;
    } else {
        consume(HOURS, "Expected 'minutes', 'seconds', or 'hours' as time unit.");
        return TimeUnit.HOURS;
    }
}

private Token consume(TokenType type, String message) throws TimerDSLEException {
    if (match(type)) {
        current++;
        return tokens.get(current-1);
    }

    throw new TimerDSLEException(message);
}

private boolean match(TokenType type) {
    return tokens.get(current).getType() == type;
}

private boolean isAtEnd() {
    return tokens.get(current).getType() == EOF;
}
}
```

Listing 32: The complete recursive descent parser (Continuation).