

Grammar Zoo: A Corpus of Experimental Grammarware

Vadim Zaytsev

*Software Analysis & Transformation Team (SWAT),
Centrum Wiskunde & Informatica (CWI), The Netherlands;
Universiteit van Amsterdam, The Netherlands*

Abstract

In this paper we describe composition of a corpus of grammars in a broad sense in order to enable reuse of knowledge accumulated in the field of grammarware engineering. The Grammar Zoo displays the results of grammar hunting for big grammars of mainstream languages, as well as collecting grammars of smaller DSLs and extracting grammatical knowledge from other places. It is already operational and publicly supplies its users with grammars that have been recovered from different sources of grammar knowledge, varying from official language standards to community-created wiki pages.

We summarise recent achievements in the discipline of grammarware engineering, that made the creation of such a corpus possible. We also describe in detail the technology that is used to build and extend such a corpus. The current contents of the Grammar Zoo are listed, as well as some possible future uses for them.

Keywords: grammarware engineering, grammar recovery, experimental infrastructure, curated corpus

1. Introduction

This paper contains a description of a method to compose a corpus of grammars in a broad sense. Having such a corpus could be profitable for mining new properties and patterns from the existing body of grammatical knowledge, for comparing grammar-based techniques and developing new ones. Formal grammars are inherently complex software artefacts, and until recently it was technically unfeasible to create such a large scale corpus, so in existing literature most case studies involve one, two or no more than a handful of grammars, and many statements about software language design remain statistically unchecked and empirically unvalidated or even unprovable.

The main contributions of this paper are:

Email address: vadim@grammarware.net (Vadim Zaytsev)

- The Grammar Zoo as the big and still growing corpus of hundreds of grammars in a broad sense.
- An open source toolkit for supporting the creation and expansion of the Grammar Zoo.
- A Grammar Zoo entry metadata model for enabling efficient sampling and reuse.

The paper is organised as follows: §2 explains the problem in detail, motivates the need for its solution, sets goals, defines context and envisions possible problems. In §3 we revisit those grammarware engineering challenges that have already been addressed in prior work, and have made the current development possible¹. In §4, the metadata model for the corpus is presented and the tools available for grammar extraction, recovery and evolution are highlighted². §5 lists the current contents of the Grammar Zoo and sketches directions for future work. §6 concludes the paper.

2. A repository of grammars

In [3], Klint et al. have defined the field of “grammarware” and identified its set of problems, promises, principles and challenges. The foundation of their work was formed by the vast existing body of knowledge about formal grammars, compiler construction, metaprogramming, source code analysis, term rewriting, parsing techniques, generative programming, attribute grammars, graph transformation and other adjacent fields. In the years after that, there have appeared many publications contributing directly to this domain, and the “engineering discipline for grammarware” from the ideal long term goal has turned into a technically achievable and partially even achieved objective. However, comparison of different grammar-based methods is still hindered by the relative lack of stability in grammar metrics and their sensitivity to many factors ranging from grammar development style (e.g., horizontal or vertical style of writing production rules has a substantial impact on the number of rules) to the choice of grammar-based technology (some syntactic notations are more expressive than others; some technologies implicitly expect grammars to be written with a specific kind of recursion, etc).

In contemporary software engineering, especially in empirical studies thereof, a similar problem has been addressed by introducing a curated collection of code artefacts [4]. In model-driven engineering, many metamodels — artefacts commonly compared to grammars in the literature — have been collected in one place to form a corpus available in many formats [5]. Such reference corpora can be used as an input of various newly proposed analysis and transformation

¹Parts of the section (e.g., §3.1) have previously appeared in a workshop publication [1].

²Parts of the section (e.g., §4.3) have previously appeared in a workshop publication [2].

techniques, allowing their output to be measured and reported in a systematic manner.

In [6], van Wijngaarden states that a powerful and elegant language should not contain many concepts and should be explainable in few words. In [7], Wirth concludes that language simplicity should be achieved through modularity and not through generalisation. In [8], Hoare claims that in programming language design, simplicity is different from and more important than modularity. In [9], Mernik et al. argue that language modularity is positively influenced by the presence of a textual notation. In [10], Völter et al. tie the multitude and diversity of general purpose programming languages to their domain-specific optimisations. In [11], Chomsky elaborates that semantic and statistical considerations should be of no consequence to the grammatical structure of the language. In [12], Erwig and Walkingshaw maintain that language design should be semantics driven. In [13], Tratt establishes that the evolution of a domain specific language mostly involves adding functionality found in general purpose programming languages. In [14], Herrmannsdörfer et al. observe that modelling languages evolution is bound to requirements creep and technological progress. In [15], Hutchinson et al. claim that effectiveness of a domain specific language and narrowness of its domain are in inverse proportion. Many more claims like these can be found in academic and engineering publications about grammarware and related topics — most are backed by expert opinions and case studies of manageable size. However, we still lack the luxury of (re)formulating them as research hypotheses and subsequently validating against (a chosen part of) the corpus of grammars and languages. We construct the Grammar Zoo in order to enable such activities in the future.

It has been previously noted by Do et al. that obtaining the right kind of infrastructure for setting up experiments is nontrivial and labour intensive, and its usefulness has huge impact on future experiments [16]. According to Do et al., the users of such infrastructure mostly face the following challenges:

Supporting replicability across experiments. Homogeneity or well-documented heterogeneity of the collected artefacts and completeness of metadata are the key factors for the creators of the infrastructure, to help addressing this challenge [16, 17, 18].

Supporting aggregation of findings. Systematic capture of the experimental context is required to complement high replicability, in order to guarantee correct aggregation of findings from different experiments [16, 19].

Reducing the cost of controlled experiments. In order to facilitate painless artefact reuse, artefact organisation needs to be standardised, they need to be complete in some sense (preferably by conforming to a well-defined completeness level) and require as little manual handling as possible [16].

Obtaining sample representativeness. The main problems foreseen by [16] in allowing the users to acquire representative samples, are small sample sizes and sampling bias. These are to be addressed by including many artefacts obtained from different heterogeneous sources.

Isolating the effects of individual factors. Isolating software language design concerns and decoupling conceptual modules within one software language have always been challenging problems, and still pose great difficulty. Since this is an open research question, we will not be able to prevent all problems that it leads to.

2.1. Illustration: Grammar Zoo utilisation

Suppose that we have assembled the Grammar Zoo as a collection of various grammars. What kind of research questions we can answer with it and what kind of problems can we address? We provide some scenarios below.

Interoperability testing. Suppose that we have identified multiple grammars of the same intended language that correspond to its different frontends. To test their interoperability, one can do code reviews or develop a test suite, but a better, more systematic, way is to generate such a suite and compare or converge those grammars automatically. An approach for that has been proposed in [20] and evaluated by two case studies involving 4 Java grammars and 33 TESCOL grammars, which were extracted from parser specifications and became one of the early fragments of the Grammar Zoo.

Grammar recovery heuristics. There have been many successful attempts of reusing grammatical knowledge embedded in various software artefacts like parser specifications, data format descriptions or metamodels. Some focused on idiosyncratic properties of the source notation, others tried to generalise the relaxed ways of treating the baseline artefact with heuristic rules for splitting/combining names, matching parentheses, etc [1]. The more grammars can be recovered with such heuristics, the better validated and motivated they become.

Empirical grammar analysis. Grammar metrics is a mature field of research, but more elaborate characterisations such as “top” or “bottom” nonterminals are common in grammar-based papers. Given a large enough repository of various grammars, the micropattern mining methodology [21] can be applied to infer characterisations by mining the repository [22]. They can in turn be used for clustering grammars based purely on statistical data about sets of indicators.

Grammar components. There is ongoing work on identifying semantic components of software languages that correspond to concepts like loops, variables, exception throwing, etc. [23]. By using a combination of program slicing and clone detection techniques on a large enough corpus, we can identify syntactic components of software languages and investigate whether there is any correspondence with semantic components.

2.2. Software Language Processing Suite and GrammarLab

The creation of the Grammar Zoo was facilitated by the tools accumulated within two projects: the Software Language Processing Suite, also abbreviated

as SLPS [24], and the GrammarLab [25]. SLPS started in 2008³ on the Sourceforge platform and migrated to GitHub in 2012. GrammarLab is a grammar manipulation library for Rascal [26] developed within a project with the same name⁴, conceptually based on the same techniques and also hosted at GitHub.

In particular, they contain the following groups of tooling:

- *Formats* for storing grammars [27], grammar transformations [28], grammar mutations [29], metalanguage notations [30], grammar documentation [31], etc. These formats have evolved by being used in methods and tools working with grammars in a broad sense and thus abstract from many technical idiosyncrasies. The formats are designed to describe essential grammatical knowledge and are therefore compatible with many different platforms and technologies.
- *Transformers* for grammars [28, 29], metasyntactic specifications [32], language documents [31], etc. These can be used either as first class entities to encapsulate grammar evolution or correction steps, or just as a technical aid for changing grammar-related artefacts.
- *Megamodels* for modelling a linguistic architecture of a system (a megamodel is model of a software system with some elements denoting very complex notions like languages, grammars, technologies and stakeholders) when nontrivial differences between similar technologies need to be spotted, modelled and resolved [33]. SLPS contains a Rascal megamodeling library that is compatible with the most recently developed general purpose megamodeling language MegaL [34, 35, 36].
- *Extractors* for locating and obtaining fragments of grammatical knowledge from various software artefacts such as parser specifications, data type definitions, grammars of various kinds, data schemata, etc. Since extractors are essential for the process of creation of the Grammar Zoo, they will be described in more detail in §4.7.
- *Recovery tools* are advanced extractors that use heuristics or context conditions to infer corrections to the source during or immediately after the extraction of a grammar. In the past, grammar recovery has been done manually [37] or semi-automatically [38]. The state of the art tools allow fully automatic grammar recovery to work based on a specification of the expected metalanguage [1].
- *Analysers* for investigating grammar properties, calculating grammar metrics [39, 40], expressing and collecting micropatterns [22], to aid in other grammar-based activities.

³Ralf Lämmel, *Software Language Processing Samples (SLPS)*. MSDN blog post in *Grammarware, Haskellware, XMLware*, May 2008, <http://blogs.msdn.com/b/ralfflammel/archive/2008/05/23/software-language-processing-samples-slps.aspx>.

⁴Paul Klint, Jurgen Vinju, Tijs van der Storm, Vadim Zaytsev, *Foundations for a Grammar Laboratory*, NWO 612.001.007, 2010–2013, <http://grammarware.github.io/lab>.

- *Exporters* take care of grammar visualisation and producing grammars in a form recognisable by other language workbenches and metaprogramming environments.
- *Test generators* can automatically support a test suite exercising all features of a software language described by a grammar and possibly limited by coverage criteria [41, 20].
- *Documentation* support is present as well: in its foundation we have a unified format for language documents [31] that was constructed by analysing and comparing hundreds of software language specifications, descriptions, standards, studybooks and alike.
- *Frameworks* for working with the above are present in languages like Prolog, Python and Rascal.

Both the Software Language Processing Suite and the GrammarLab are noncommercial efforts mainly aimed at developing tool prototypes and proofs of concepts [24, 25]. They are publicly available on the internet for distribution through a free and open software license (CC-BY: Attribution⁵).

2.3. Main objectives

We define the **main goals** of the Grammar Zoo as follows:

- Collecting grammars in the broad sense — structural definitions of software languages.
- Annotating each grammar with information about its source, original format and authors.
- Complementing each grammar with details about how it was fetched, extracted, recovered, adapted, etc.
- Documenting usages of each grammar — its derivatives, tools, documents and other grammars.
- Making all grammars publicly available in a variety of formats.

What the Grammar Zoo is **not** about:

- It is not about collecting parsers. Not all grammars in a broad sense are meant to be used for syntactic analysis of textual data, and collecting a large number of them systematically would mean committing to a specific parsing technology or even a specific grammar manipulation framework (metagrammarware).

⁵CC Attribution 3.0 Unported, <http://creativecommons.org/licenses/by/3.0/>.

- It is not about unifying syntactic notation for grammars. Numerous attempts to unify textual representations of context-free grammars have failed in the past. Instead of fighting the diversity of notations, we embrace it and develop tools that can deal with it. We ultimately aim at storing the pure grammatical knowledge and exporting it on demand according to the users' needs.
- It is not about enforcing the quality and level of grammars. Heterogeneous content is also inherent to the field of grammarware engineering in its current state: different tasks expect different properties from grammars. Instead, we aim at properly documenting such differences so that the corpus can be used to obtain representative sets of grammars that are similar in some required sense.

2.4. Summary

Recent advances in the field of grammarware engineering make it possible to design and engineer a corpus of grammars in a broad sense. Having such a repository of grammars, each annotated with metadata about its source, means of extraction and recovery, its evolution and linked tooling, would allow us to mine it for similarities and singularities, as well as to use it as a common testing ground for grammar-based methods. Based on the experience gained from several such experiments in the past, in the next sections we will compose the Grammar Zoo, which will collect as many grammars as we can secure.

3. Previously addressed challenges

The Grammar Zoo is a relatively new initiative which was infeasible until recently. The following subsections are dedicated to the challenges which were solved during the last years by various researchers, and helped us to found this initiative on top of their methodologies. We credit the following research directions:

Grammar extraction. Automated recovery of grammars from existing artefacts is required, otherwise adding each grammar to the corpus will always be a separate project with its own specifics.

Grammar evolution. When grammar recovery goes decidedly beyond extraction, it involves bringing systematic changes to the grammar. Proper documentation of such adjustments relies on an advanced transformational infrastructure.

Metalinguistic evolution. With many notations for syntactic definitions being used in various grammar-based toolkits, it is crucial to be able to export each grammar in a variety of notations, or perhaps even in a new notation defined on-the-fly.

Generating browsable documentation. One of the most common uses for a grammar, beside generating grammarware code, is facilitating its inspection.

3.1. Grammar extraction and automated recovery

Basically, there are three main challenges in developing grammar extraction: the unnecessary diversity of notation [54, 30]; the error-prone manual process of grammar creation and typesetting; and the semantic gap between the metalanguages. The first (notational) challenge is addressed by the notation-parametric approach [1] that requires a formal specification of a notation and can base the extraction steps on it. The second (lexical) challenge is solved by heuristics expressing relaxed conformance to the intended notation. The third (semantic) challenge is avoided by using the “extraction through abstraction” [27] approach of extracting pure grammatical knowledge from the available artefacts while abstracting from the technology-specific details. All three challenges are addressed by automated tools from SLPS/GrammarLab, a more detailed overview of them will follow in §4.7.

The progress of the grammar extraction methods with automated error recovery, also summed up in Table 1, was as follows⁶ [1]:

Message Sequence Charts was a DSL described in a Word document, which was converted to an ASCII file in 1996, processed by a Perl script and produced BNF rules, which were in turn manually edited with all 14 changes claimed to be documented. Another script was used to generate a hypertext form of a grammar suitable for browsing [44].

COBOL grammar capable of handling a range of language dialects was recovered in 1997 by converting 1100 production rules of the ANSI COBOL 85 standard to SDF [56]. A long and sophisticated process of forced coupling followed, leading to (disciplined) changes brought both to the codebase and to the grammar, and resulting in capability of the adjusted grammar to parse the adjusted source code [37].

Switching System Language was a proprietary DSL documented in a set of HTML files containing its grammar in an BNF dialect they called SBNF. The recovery endeavour was reported in 2000 and is a remarkable milestone in a way that it was an attempt to use precise parsing on an unreliable source. A range of (as we are now aware) typical issues arose such as naming convention violations and non-matching brackets, and significant amount of interactive grammar adjustments was needed. The project succeeded also due to development support of the ASF+SDF Meta-Environment, resulting again in the situation where an adjusted SBNF grammar was used to parse adjusted syntax rules [44].

⁶NB: we limit the mentioned initiatives to those which were extracting grammatical knowledge from software artefacts that were already known to contain it and to *define structures*: context-free grammars, compiler sources, parser definitions, XML schemata, etc. For a comprehensive overview of methods and initiatives of extracting/inferring grammatical knowledge from a collection of *language instances*, we redirect the reader to [55].

Language	Source	Extracted	Connected	Adopted	Browsable	Result	Ref
ANSI COBOL III	spec in "general format"	manually	manually	—	—	no	[42]
IBM COBOL	SC26-9046-03 spec	Prolog	—	—	Prolog	[43]	[43]
IBM OS PL/I V2R3	SC26-4308-02 spec	Prolog	—	—	Prolog	[43]	[43]
Ericsson MSC	ITU Z.120 (MS-Word document)	Perl	manually	—	Perl	no	[44]
ANSI COBOL 85	hard copy	reused	reused	ASF	—	no	[37]
Ada 95	ISO/IEC 8652:1995 spec	Prolog	—	—	Prolog	[43]	[44]
Ericsson SSL	hypertext documentation	bnf2sdf	manually	ASF	Box	no	[44]
Ericsson PLEX	BNF rules in compiler sources	6×SDF	interactive	ASF	—	no	[45]
IBM VS COBOL II	language spec [46]	Prolog	interactive	ASF	Prolog	[43]	[38]
C#	ECMA-334 spec [47]	manually	FST	FST	ASF	[43]	[48]
FL	various grammars in a broad sense	automated	XBGF	—	XSLT	[24]	[27]
Java	Sun language specs [49, 50, 51]	Python	XBGF	—	BGF	[24]	[28]
MediaWiki	community-created grammar	5×EDD	XBGF	—	XSLT	[24]	[52]
SLPS Zoo	42 sources	varies	varies	—	XSLT	[24]	[53, ...]
SLPS Tank	53 sources	varies	varies	—	XSLT	[24]	[53, ...]
Grammar Zoo	1710 sources	varies	XBGF or GLUE	—	XSLT	[24]	here

Table 1: The summary of grammar recovery activities that took place before the Grammar Zoo. The early SLPS Zoo and SLPS Tank, as well as the current Grammar Zoo, are added for comparison to the bottom of the table. In the "extracted" column we denote a method of obtaining an *extracted* (level 1) grammar, a freshly extracted CFG corrected from misspellings. The next column is a method of getting a maximally connected (level 2 or 3) grammar, while "adapted" usually refers to a level 4 or 5 grammar strongly linked to a tool that has been exercised on large codebase. §4.3 elaborates on the notion of a grammar level.

Programming Language for EXchanges was a complex DSL consisting of 20 sublanguages (“sectors”) and having over 60 Mb of grammarware source code. The mining process delivered fragments of BNF found in the comments, which with the help of six parsers were transformed to pure aggregated BNF and subsequently to SDF, which was combined with a lexer. The project took only two weeks and resulted in parsing 8 MLOC of unmodified PLEX, as reported in 2001 [45].

IBM VS COBOL II recovery project is one of the most complicated among those reported in academic sources. A raw grammar was extracted from the language documentation, which was not trivial since it used “railroad track” kind of syntax diagrams instead of purely textual (E)BNF. The recovery went in stages: static error fixing, adding lexical syntax, test-driven correction/completion, beautification, modularisation, disambiguation, adaptation, etc. The recovery was reported in 2001 [38], and its outcome was made freely available for reuse from the authors’ website [43].

C# recovery targeted a non-legacy language in 2005, but witnessed similar problems. In order to parse *C#* code, the project involved manual transition from the ECMA-produced PDF to LLL and intensive grammar transformation with FST and GDK, as reported in [48] and [57, §3].

FL was an artificial toy functional language used to demonstrate the principles of grammar convergence in [27]. The project assumed sources to be reliably correct and focused on the lightweight extraction process from concrete syntax definitions in SDF [58], parser specifications in ANTLR [59], definite clause grammars in Prolog [60], grammars in TXL [61], object models in Java [51], document schemata of XML [62].

Java grammar recovery required tolerance to overcome layout inconsistencies and other lexical deviations of the source grammars, which was expressed as a list of heuristics and described in detail [28]. The same technology was used later for grammars of C, C++ and C# found in other ISO standards written in the same syntactic notation [57].

Notation-parametric recovery method [1] relies on encapsulating commonly varying details of the syntactic notation in a notation specification [30] and binding the recovery heuristics to those variation points. This approach allows to extract a grammar in a never-seen-before notation in a matter of several minutes required to compose such a specification.

3.2. Grammar evolution support

The ability to express grammar evolution steps as first class executable entities, was identified as one of the crucial components of the engineering discipline

for grammarware [3]. Below we try to cover the existing spectrum of the grammar evolution support⁷:

Attribution as a claim that one grammar is “derived from” or “written from” some other source like another extracted grammar or language documentation, is a common way to represent grammar evolution in many existing grammar repositories [67, 68, 5]. However, these derivation steps are rarely documented, so this level of detail is not informative enough for any automated verification of such claims or even for their consistency management.

Documented patches as lists of changes that were applied to the grammar in order to get from the original to the final version, are much more useful, even if they are not entirely formal. They have been used in early grammar engineering projects [44], and are also not uncommon in grammar-related bug reports [69]. What is often missing in such lists, is justification for the proposed solution: for example, compare [69] with [28].

Grammar transformation operators are a functional way of representing patches. Each change step is expressed as a function application, with a function being one of the predefined operators from a grammar transformation operator suite. First such operator suites that were published, are FST [70] and the one used for COBOL grammar recovery [71, 38]. The differences between them are insignificant for the current paper — a discussion about them can be found in [28]. An ideologically similar approach was demonstrated by TXL [61], a framework where grammar specialisation for each task is an essential part of the grammarware engineering paradigm.

Higher-level grammar transformation operators like “fold a nonterminal” or “perform a safe refactoring” were shown to be more useful and better maintainable than the low-level ones like “remove any part of a grammar” or “replace any expression by another expression everywhere” [72]. This approach was used in Grammar Deployment Kit, experimental metagrammarware that was used successfully in a number of projects [73]. A similarly advanced operator suite for the modelware technological space, has also been developed and published [74]: even though always taking coupled evolution of models into account when dealing with metamodel changes, has its challenges [75], they seem to be addressed in recent research [76].

Recovery domain-specific operators were introduced for Grammar Recovery Kit as a demonstration of how a successful grammar recovery project can be undertaken and documented with a minimal number of them [77].

General purpose grammar evolution was the opposite attempt to cover all possible use cases for grammar evolution, recovery, convergence, adaptation,

⁷NB: we limit our overview on *programmable* grammar evolution, without considering adjacent initiatives such as incremental grammar refinement [63, 64, 65, 66].

etc. The resulting grammar manipulation language is called XBGF, for Transformations of BNF-like Grammar Format [57, §7], and was also used extensively throughout the Software Language Processing Suite, for many tasks in many projects [24], some of which involved operating on grammars of industrial size (Java, C#, C++, etc).

Bidirectional grammar evolution operator suite Ξ BGF is a bidirectional variation of XBGF, that has shown its usefulness for metalinguistic evolution and derivation of transformation steps [32].

3.3. Metalinguistic evolution

Even though many metagrammarware tools and language documents claim to use EBNF [54], it does not mean that they agree on a metalanguage. EBNF has grown to become a family of textual notations for defining context-free grammars with possible extensions. We speak of “metalinguistic evolution” as a specific case of grammar evolution, when the language defined by the grammar, does not change, but the metalanguage in which it was written, does. This is a known problem at least since [44] (translation between SBNF and SDF) and [48] (translation from EBNF used by ISO standards to LLL), which received a relatively straightforward solution with the introduction of the notation specification [30].

Given a grammar $G_N(L)$ written in a metalanguage N , we can express metalinguistic evolution as a metalanguage transformation σ that transforms specification $S(N)$ into a specification $S(N')$. A new metalanguage N' , defined by the transformed specification; its own grammar $G(N')$ and an updated grammar $G_{N'}(L)$ can all be automatically derived from this σ , as shown in [32].

3.4. Generating browsable artefacts

In [44], the generation of browsable artefacts was claimed to be possible with the Box functionality of the ASF+SDF Meta-Environment [78]. In [77], the author went to great length to inject the changes in the grammar back to its documentation. As a part of research on language documentation in [31], a case study was completed that involved extracting all the available information from a language manual for the purpose of regenerating it after necessary manipulations.

In general, we have expected and achieved the following qualities in order to claim useful browsability:

Metasyntax highlighting means that different entities are displayed differently — for example, terminal and nonterminal symbols use different colours.

Interactiveness means that if an element can be observed, it can also be interacted with, if such interaction makes sense — for example, one can get to a definition of a nonterminal from its occurrence.

Metrics as simple as top and bottom nonterminals proposed by [44, 79] or as complex as grammatical level depth and others listed in [39, 40], aid comprehension of a grammar just as much as traditional software metrics help estimating code quality and detecting code smells.

Full automation is expected to take care of the abovementioned qualities without asking the end user to set up any complex infrastructure or configuration.

All four objectives can be achieved in at least two different ways: hypertext and IDEs. Hypertextual rendering of grammars is a pretty straightforward mapping from the language of their internal representation to XHTML supported by a predefined CSS. There is hardly anything scientifically challenging in this mapping, but we have of course engineered and automated it within the SLPS/GrammarLab to support the Grammar Zoo [24]. The other way to achieve browsability is relying on a sufficiently advanced IDE framework such as Rascal language workbench [26]. Using the framework functionality, one can quickly prototype a powerful grammarware engineering environment, which would be highly domain-specific and yet extensible and programmable.

4. Creating a grammar repository

Having set up the objectives in §2 and revisited previously addressed challenges that enable the creation of the Grammar Zoo in §3, we move on to explain the repository creation process. In particular, we will define the model of the metadata that is supplied for each grammar in the repository (§4.1), explain different kinds of sources for grammar extraction and the levels of grammars extracted from them (§4.5), list the currently available automated extractors (§4.7) and exporters (§4.8) and explain the recovery process (§4.10). Several intervening sections will present illustrative examples to demonstrate our approach.

4.1. The model of metadata

We present the data model behind the Grammar Zoo frontend on Figure 1. (An intuitively readable dialect of EBNF is used with ? denoting “zero or one”, * denoting “zero or more”, + denoting “one or more” and using | for choice; x:y means referencing a nonterminal y with the name x). Details follow:

repository — the root element;

entry — the first class entity, an entry in a repository; can be one language or language family (e.g., “Java”, “C#”, “(E)BNF”), or one version of a language with all its grammars and metadata (e.g., “Java 5”), in both cases containing subdirectories with other entries; can also be a particular identifiable software language (e.g., “a Java 5 grammar from §18 of the Java Language Specification”);

```

// A repository is hierarchically structured
repository ::= entry+;
// Each directory has a name for its hierarchy level
// (which corresponds to a language, a dialect or any group)
// and lists resources and possibly grammars.
entry ::= name source:resource (item:resource)* grammar*;
// Each (re)source refers to its origin and may state title, subtitle,
// publication venue, edition number, etc.
resource ::= origin title? subtitle? venue? version? edition?
           date specific? link+;
// Resource origin is a list of authors, a reference number of a standard
// or an organisation name that produced it.
origin ::= author+ | standard | organisation;
// Each grammar refers back to the descriptive name of the language,
// its corresponding status, a method of obtaining this grammar,
// as well as tools and files used.
grammar ::= of? status level method toolused* fileused*;
// There are many methods of obtaining a grammar:
method ::= copy|download|automatic|semi-automatic|git|...;
// Each link can be named and is generally a URI with shorthand notations.
link ::= name? (uri|doi|tool|xbgf|wiki|readme|slps|...);
// Each related grammarware tool has a name (e.g., parser),
// a technology it is based on (e.g., ANTLR) and a list of links.
tool ::= name technology link+;
// There is also a repository of tools known to the system
toolused ::= tool | toolref;

```

Figure 1: The data model of Grammar Zoo displayed as an annotated grammar.

grammar — one grammar usually with files used in the process of its extraction, metadata annotations, execution scripts, grammar manipulation scripts, etc;

name — a string naming something; important for identification of a language, a language version or a grammar within the repository, although does not have to be unique: it is possible and useful to have several grammars of the same intended language in a repository;

source — the primary resource used for grammar extraction (several resources are allowed if the grammar fragments needed to be collected); discussed in detail in §4.5;

item — a related resource of secondary importance;

resource — a publication related to the grammar or using it extensively, a website dedicated to its recovery process, etc;

origin — any entity responsible for creating a resource;

title — the title of a resource;

subtitle — a string necessary to identify the resource, but not a part of the title itself;

venue — the name of a conference, a workshop or a journal;

edition/version — language specifications often have editions or versions;

date — when the resource was created, usually the year of publication: timestamps can be important (among other things) to identify the version of the intended language, if it cannot be derived through other means;

specific — specific coordinates for the extraction source within the generally identified resource (i.e., a chapter of a book, page numbers, an important note);

link — essentially a customary named URI;

author — the name of each of the authors stored as a string: possibly can be matched with DBLP, ORCID or similar framework and appropriately linked;

standard — a reference number for a language standard (e.g., for the ISO standard of EBNF it is “ISO/IEC 14977:1996(E)” [80]);

organisation — the name of a company or a standardisation body responsible for the creation of the resource (and possibly the holder of the copyright);

of — an explicitly verbose human readable identification of a software language intended to be defined by the grammar (used in the visualisation);

status/level — the maturity status of the recovered grammar; discussed in detail in §4.3; the status also serves as the name of the subdirectory where the grammar is located together with its attachments;

method — a particular way of obtaining the grammar of this status and level: manual copy-paste, automated download, automatic use of a toolkit, semi-automatically programming the necessary steps, digging it out of the repository history, etc;

toolused — usually a grammar refers to at least one set of extraction tools, but there can possibly be more (recovery toolkits, disambiguation tools, grammar correction scripts, etc); discussed in detail in §4.7 for extractors and §4.10 for recovery tools;

fileused — the name of a particular file used for extraction: especially crucial for multipart sources such as modularised grammars or complicated modularised adaptation scenarios;

- `uri` — a uniform resource identifier, a link to a webpage;
- `doi` — a digital object identifier, the easiest way to refer to most academic publications: can be easily resolved to an official publisher’s page with <http://dx.doi.org>;
- `tool` — a grammarware tool coupled with the grammar or related to it; for example, if this grammar was extracted from a parser specification, this will link to the executable parser; it can be a validator for data models, a refactoring tool, a migration tool, etc.;
- `xbgf` — a shorthand notation for referencing grammar transformation scripts, which are then also automatically rendered as hypertext; such link does not need a name since it can be inherited from the transformation script;
- `wiki` — a shorthand notation for referencing a page on the SLPS wiki at <http://github.com/grammarware/slps/wiki>, usually such a page is automatically generated and contains links to multiple files in the repository that comprise one tool;
- `readme` — a shorthand notation for referencing a verbose `README.txt` file supplied with more details described in natural language;
- `slps` — a shorthand notation for referencing other files in the SLPS repository on GitHub [24];

The specifications of the Grammar Zoo conforming to the data model described above, can be found as `zoo.xml` files in SLPS in directories with each extracted grammar. The root `zoo.xml` document⁸ in the GitHub project [slps.github.com](http://github.com) collects links to all those. The table with a brief overview of the current contents of the Grammar Zoo will be presented later as ??.

4.2. Illustration: notation-parametric recovery of a Java grammar

Consider Figure 2 as an example of one of the 1710 entities comprising the Grammar Zoo. In this example, dots represent omissions. Additional resources refer to [28] and other related publications. The (post-)extraction grammar transformation scripts `recover.xbgf` and `correct.xbgf` are inherited from the Java grammar convergence case study in [28] and pretty-printed in hypertext as a part of the Grammar Zoo. Both represent grammar evolution steps that lift this Java grammar from the *extracted* (freshly recovered) to *connected* (checked and tested) — we will list all supported statuses later in §4.5. The main difference between them is that `correct.xbgf` fixes the mistakes made by the creators of the grammar source, while `recover.xbgf` fixes the mistakes that were not automatically handled by recovery heuristics. For instance, the non-terminal `FormalParameter` was erroneously left undefined in the specification

⁸http://github.com/slps/slps.github.com/blob/master/_dev/zoo.xml

```

<entry>
  <name>Implementable</name>
  <source>
    <author>James Gosling</author>
    <author>Bill Joy</author>
    <author>Guy Steele</author>
    <author>Gilad Bracha</author>
    <title>Java Language Specification</title>
    <edition>3</edition>
    <date>2004</date>
    <specific>Ch. 18: Syntax, pages 585-596</specific>
    <link>
      <uri>http://java.sun.com/.../syntax.html</uri>
      <name>HTML</name>
    </link>
  </source>
  <readme/>
  <grammar>
    <dir>fetched</dir> <level>0</level>
    <method>download</method>
  </grammar>
  <grammar>
    <of>Java (J2SE 5.0 "Tiger")</of>
    <dir>extracted</dir> <level>1</level>
    <method>automatic</method>
    <toolused>html2bgf</toolused>
  </grammar>
  <grammar>
    <of>Java (J2SE 5.0 "Tiger")</of>
    <dir>connected</dir> <level>2</level>
    <method>semi-automated</method>
    <fileused>recover.xbgf</fileused>
    <fileused>correct.xbgf</fileused> <toolused>xbgf</toolused>
  </grammar>
</entry>

```

Figure 2: A typical entry in the Grammar Zoo, corresponding to the “more implementable” grammar of Java 5 extracted from the Java Language Specification [51] for the purpose of reverse engineering its relationships with other Java grammars [28].

— hence, `correct.xbgf` contains a call to a **define** operator with appropriate arguments. On the other hand, on several occasions curly bracket terminal symbols were mistaken by the extractor to be repetition metasymbols (i.e., in the JLS notation “`{x}`” means “zero or more `x`”) — hence, these mistakes were fixed in the `recover.xbgf` script.

4.3. Grammar maturity status

Lämmel and Verhoef proposed the notion of a *grammar level*⁹ to specify a quality level or a recovery status of a grammar. For the Grammar Zoo, we

⁹These “grammar levels” are essentially CMM-like levels applied to grammars, unrelated to well-known “grammatical levels” used for a range of grammar metrics [81, 39].

have essentially inherited that hierarchy and extended it to accommodate more important details, yielding a maturity model for grammars [2]. We distinguish among the following grammar levels:

- A grammar is **fetch**ed if it can be put in a file which we claim to contain grammatical knowledge. A fetched grammar is usually written in an (E)BNF-like notation, but it can also be an XML Schema schema, an Ecore model, a parser specification, etc. A grammar from an undisclosed ISO standard or a grammar built in a proprietary tool is *not* fetched, since we have no possible way to extract the knowledge from it. A compiler is therefore not a fetched grammar since the grammatical knowledge is ingrained too deep in it and requires special techniques to be fetched. The sources of a compiler, however, can be considered fetched, since further extraction can be semi-automated, and the result will depend mostly on the source and not on the extraction algorithm. Hence, a corpus of programs in a given language is also *not* a fetched grammar, but an APTA (Augmented Prefix Tree Acceptor) [82] or a DFA (Deterministic Finite Automaton) constructed from it by a grammatical inference algorithm, is a fetched grammar. A fetched grammar can contain unreadable symbols, incorrect indentation, parts written in an unknown notation or a natural language, or even be present in a form of an image or a manuscript.
- A grammar is **extract**ed, if it was fetched and then successfully processed by a (hopefully automatic) grammar extractor, possibly also corrected of typographical, text recognition and similar errors and converted into a context-free grammar or, more broadly speaking, to a Boolean grammar [83]. (Not venturing beyond context freedom is simply a consequence of the current lack of theoretical foundations for linking classic context-dependent grammars to generalised types — in general, aligning the Chomsky hierarchy [84] with Barendregt λ -cube [85]). An extracted grammar is suitable for automated processing: it can be pretty-printed in a range of different ways and transformed by general means, without writing a tool specific for its peculiar notation or format: syntax diagrams, Relax NG schemata, algebraic data types, parser specifications all lose their notational differences when they are being extracted, but they retain all structural peculiarities such as using a particular style of recursion (syntax diagrams are incapable of expression left recursion, and some parsing algorithms tend to avoid it as well), the lack or presence of terminal symbols (anything that defines an abstract syntax, has no terminals) or nonterminal symbols (classic regular expressions [86] have no notion of named subcomponents), etc. An extracted grammar corresponds to a level 1 grammar from [38].
- A grammar is **connect**ed, if it was extracted and then processed to not contain unwanted top sorts (defined but never used) and bottom sorts (used but not defined). These two quality indicators were proposed in

[79, 44] and discussed in more detail in [71] before being formalised as micropatterns [22]. Connecting a grammar can be done automatically with a mutation [32, 29], semi-automatically with a sequence of transformation steps, or by editing it inline. Any connected grammar corresponds to level 2 from [38]. Connecting is a simple procedure that allows to start making some claims about the grammar, since it enables its formalisation (the classic $\langle \mathcal{N}, \mathcal{T}, \mathcal{P}, s \rangle$ model of a grammar requires it to have one known starting symbol) and possible application of grammar-based algorithms (in particular grammar-based test data generation expects the grammar to be connected because otherwise it is futile to use any coverage criteria). In a broader sense, a connected grammar always relies on some underlying mechanism of testing or validation which ensures its general quality — as opposed to the extracted grammar which can be an output of an automated extractor and never checked nor inspected further.

- A grammar is **adapted**, if it is connected and then transformed towards satisfying some constraints: it could be complemented with a lexical part, or its naming convention can be adjusted, or certain metaconstructs can be introduced to or removed from its syntax. The adaptation has a clear intent: adding a lexical part can lead to automated generation of a parser or at least a recogniser; conforming to a naming convention can enable the use of the grammar in specific language workbenches, etc. An adapted grammar corresponds to level 3, or to level 4 if it has been tested on a large scale [38].
- A grammar is **exported**, if it was adapted and then a piece of grammarware was generated from it. An exported grammar bidirectionally and possibly nontrivially corresponds to a real piece of grammarware such as a compiler or a code analysis or transformation tool. This is a level 5 grammar [38] which either demonstrates the absence of manual steps in grammar deployment, or documents them by its existence.

Each Grammar Zoo entry has one *fetch*d grammar: ones with less than one are “non-entries” that can perhaps be referred to, but can under no circumstances be machine processed; having more than one *fetch*d grammar can happen for cases such as multiple websites mirroring one another, but then a simple check is required to assert them to be equal. If several extractors are available (e.g., one straightforward one and one heuristic-based error-correcting one), there can be multiple *extract*d grammars per entry. Similarly, there can be several grammars of level *connect*ed and up per entry, varying per their extracted source and methods of processing.

At this point in the history of the Grammar Zoo we have not yet experienced the need to explicitly distinguish the reason for adaptation of each grammar: some are massaged for better readability, some adjusted with parsing in mind, some are disambiguated [87], some adapted for testing purposes [88, 20], etc. We intentionally leave the hierarchy as general as it is, and leave its extension to future work.

4.4. Illustration: a grammar life cycle

Suppose we would like to have a piece of grammarware to parse and analyse programs in a particular software language — say, COBOL or PHP. Being constrained in time, we usually start by looking for existing grammars: once we find one that seems reasonably suitable for our needs, we can declare it *fetches*. If a fetched grammar of our intended language is already in the Grammar Zoo, it can save us the search, the frustration from websites having been taken down, as well as the ambiguity about the true source of the grammar.

Once the grammar is fetched, it is usually necessary for us to extract it. In the simplest cases, grammar extraction methods and tools can be applied with reasonable success. There is quite a collection of them readily available within Software Language Processing Suite and GrammarLab (see §4.7), and it is fairly straightforward to use notation-parametric grammar extraction [1], if the input notation is anything like BNF or EBNF. If all available methods fail, we can attempt to apply grammar recovery tools (see §4.10), which have heuristics known to overcome frequent erroneous patterns. Once some reasonable kind of non-empty grammar is obtained or if it was in the Grammar Zoo to begin with, the grammar can be considered to be *extracted*.

An extracted grammar is a full-fledged tangible software artefact that can be processed further, analysed, transformed, exported, imported, visualised etc — there are many tools in the GrammarLab and SLPS that can do it directly, and they can also help to export it to a format readily consumable by other metagrammarware. However, it does not mean that this grammar would “work” there, whatever that might mean. There are some sensible metrics, constraints and grammar analyses established in state of the art grammar recovery [38, 44, 57], that are almost universally useful in improving the quality of a grammar. For instance, we would like to identify the starting symbol of a grammar, establish it being unique. Furthermore, for each parts unreachable from it, we would like to make a decision and either remove them or connect to the rest of the grammar. This is usually done by programming the corresponding steps in XBGF [27], SLEIR [29] or any other grammar manipulation language. This usually implies manual examination of a grammar and its metrics by an expert, making the appropriate decisions and then documenting the changes. Once this is completed, we speak of having a *connected* grammar.

The next step is grammar adaptation: a goal-specific continuation of grammar transformation activities. For example, if we have decided to parse and analyse code in COBOL or PHP, this is our goal, and in case of Rascal it will mean having a complete concrete syntax specification, and a suitable algebraic data type. Both can be obtained from a connected grammar, but the adaptation strategies are different. For a syntax specification, we need to add the lexical part, specify layout, increasingly disambiguate the grammar, etc. For a data type, we should think of its suitability for specifying our analyses later, and we can easily eliminate all terminal symbols and massage the remaining abstract grammar to enable more concise and readable patterns. These streaks of activity end up with an *adapted* grammar each.

Finally, our two grammars (or a syntax spec and a data type, or a grammar and a schema — terminology may vary) are ready to be *exported* — we do this with out of the box renderers (see §4.8), possibly followed by manual polishing such as adding documenting annotations and inserting copyright notices. It is not unusual for an exported grammar to be linked to a specific tool which it forms a part of.

4.5. Grammar sources

So far, we have encountered the following kinds of grammar sources:

Language standard is a language document that was developed under supervision or received acknowledgement from a standardisation body (ANSI, ECMA, IEEE-SA, ISO, IEC, ITU, IETF, OASIS, OMG, WSA, W3C, etc). There are two additional factors that play important roles:

Centralised or distributed? Grammar knowledge can be concentrated in an appendix or a specific section of the language standards, but it can also be distributed all over the document (e.g., when it is used for explaining language constructs one at a time). In the second case, the extraction process is prone to missing grammar fragments due to incorrect markings and other reasons.

Open or closed? When a standardisation body commits to public disclosure of a language standard, it goes through a certain process which usually comprises sanitising the contents at least to some extent: clean up, mark up, linking and similar activities improve the quality of the grammar source. If the standard is a close publication, it can be unavailable for inspection for a broad audience (require payment or special subscription), and there is an additional step of reentering the data from its printed copy back into a computer. Both manual retyping and automated text recognition processes are error-prone. A grammar from an uncompromisingly closed publication cannot be fetched at all.

Industrial specification is in many aspects the same as a language standard, but it is developed inside a commercial company (Ericsson, Google, IBM, Microsoft, Oracle, etc). The same additional factors from language standards apply, with conditions for disclosure usually being even more strict.

Browsable documentation is often found in many corners of the Internet. Many people spend their own time on extracting grammar knowledge from the artefacts they were able to obtain, sanitising it and reformatting the resulting grammar as hypertext. Some of such endeavours that we mentioned before, are well-documented and linked to a published scientific report [43], others contain conformance and validity claims that require thorough verification.

Parser specification is an executable grammar that contains many annotations that often take it beyond the context-free class. Grammars specified in ANTLR [59, 67], Bison [89], JavaCC [90], Kiama [91], Rascal [26, 92], SDF [56, 58, 93], TXL [61, 68], YACC [94] and many other metagrammarware frameworks can be located in their corresponding repositories or just anywhere close to end users of these products. Such specifications can be stripped from excessive information and extracted in the form expected by the representation central for the repository.

Metamodel is a grammar in a broad sense used in the modelware technological space. Just like a parser specification, it can contain details that transcend structural definitions: constraints, certain relations, etc. However, that information can be abstracted from, and the grammars can be extracted. AtlanMod already started an initiative of accumulating metamodels from model-driven open source projects: we have referenced the EMF XMI part of it as [5], but the same repository also contain metamodels in KM3 [95], MSchema [96], Clojure [97], SBVR [98], UML 2.1 [99], GraphML [100], OWL [101], MOF [102], etc.

Wiki pages can also be a source for grammar extraction, if the grammar was developed by a community. So far we have encountered only one such initiative, with several reverse engineered grammars of MediaWiki, and reported it in detail in [52].

Scientific papers often contain small grammars or grammar fragments. We have not yet attempted a big scale mining process of recovering all possible grammar fragments published in a certain set of venues. However, at least once [32] it was useful to compare the grammar published in a workshop paper [72] with its updated version published electronically in a programmer's manual [73].

4.6. Illustration: fetching grammar knowledge

Apart from grammar sources occasionally found in various places, we have systematically inherited and extracted grammars from the existing collections, listed in Table 2.

ANTLR Grammar Lists are collections of context-free and lexical definitions intended to be used in a parser generator called ANTLR [59]. The current version has a separate grammar repository, and thus is growing and exposing all versions for each grammar, which can eventually be useful for mining their evolution. For the previous version (ANTLR 3), there is an operational website where these grammars can still be downloaded. A corresponding website for ANTLR 2 has already been taken down, so we are left with several grammars from it that have been previously used in research [20].

Collection	Fetches	Format(s)
<i>ANTLR 2 Grammar List</i> ^a	7 (offline)	ANTLR2
<i>ANTLR 3 Grammar List</i> ^b	154	ANTLR3
<i>ANTLR 4 Grammar List</i> ^c	24 (growing)	ANTLR4
<i>TXL World</i> ^d	21	TXL [61]
<i>SDF Library</i> ^e	21	SDF2 [58]
<i>JavaCC grammars</i> ^f	27	JavaCC [90]
<i>SableCC Grammars</i> ^g	18	SableCC [103]
<i>Rascal Language Library</i> ^h	13	Rascal syntax def. ⁱ
	8	Rascal ADT ^j
	1	EBNF
<i>Atlantic Metamodel Zoo</i> ^k	303	Ecore [104]
<i>ReMoDD metamodels</i> ^l	15	various
<i>TCS Zoo</i> ^m	24	Ecore
	24	TCS [105]
<i>Concrete Syntax Zoo</i> ⁿ	195	EMFText [106], Ecore
<i>RelaxNG Schemas</i> ^o	90	RELAX NG [107]
<i>ISO/IEC JTC1/SC22</i> ^p	40	EBNF
	WIP (C++ ^q)	EBNF
<i>OMG Specifications</i> ^r	614	KM3 [95], XSD [62], MOF [102]
<i>SLPS</i> ^s	12 (FL [27, 108])	various
	33 (TESCOL [20])	ANTLR3 [109]
	35 (LDF [31])	BGF [27]
	23 (other)	various

^a<http://web.archive.org/web/20130115233436/http://antlr.org/grammar/list.html>

^b<http://www.antlr3.org/grammar/list.html>

^c<https://github.com/antlr/grammars-v4>

^d<http://www.txl.ca/nresources.html>

^e<https://github.com/cwi-swat/meta-environment/tree/master/sdf-library/library/languages>

^f<https://java.net/projects/javacc/downloads/directory/contrib/grammars>

^g<http://sablecc.sourceforge.net/grammars.html>

^h<https://github.com/cwi-swat/rascal/tree/master/src/org/rascalmpl/library/lang>

ⁱ<http://tutor.rascal-mpl.org/Rascal/Declarations/SyntaxDefinition/SyntaxDefinition.html>

^j<http://tutor.rascal-mpl.org/Rascal/Declarations/AlgebraicDataType/AlgebraicDataType.html>

^k<http://www.emn.fr/z-info/atlanmod/index.php/Ecore>

^l<http://www.cs.colostate.edu/remodd/v1/category/artifact-types/metamodel>

^m<http://wiki.eclipse.org/TCS/Zoo>

ⁿhttp://emftext.org/index.php/EMFText_Concrete_Syntax_Zoo

^o<http://relaxng.org/schemas>

^p<http://open-std.org/jtc1/sc22/>

^q<https://github.com/cplusplus/draft/commits/master/source/grammar.tex>

^r<http://www.omg.org/spec/index.htm>

^s<http://github.com/grammarware/slps>

Table 2: Existing collections of grammars in a broad sense that have been incorporated (at least *fetches*, and at least *extracted* if an extractor is present) in the Grammar Zoo.

TXL World is similarly accessed through the internet, but it deploys grammars in tarballs of varying structure, so the process involved less automation: each source must have been explored individually to locate the files relevant for grammar extraction. TXL grammars typically are of the adapted level, and can be considered exported if shipped together with a corresponding tool [61].

SDF Library is a component of the *Meta-Environment* [70] distribution, bearing similar structure and properties as the ANTLR 4 repository.

JavaCC and *SableCC* are other compiler compilers, popular among in communities different from ANTLR but sharing a Java ecosystem [90, 103].

Rascal is a metaprogramming system with many components, containing at least two forms of grammar-related artefacts: concrete syntax definitions (in fact, concrete grammars with identifiable context-free and lexical parts) and algebraic data type declarations (similar in notation and intention to Haskell ADTs) [26]. Both forms are suitable for grammar extraction, and both end up immediately as both *fetched* and *extracted* grammars.

Atlantic Metamodel Zoo is web-based, so obtaining a fetched grammar is easy with semi-standard widespread tools like `wget` or `curl`. All the metadata about the grammars (i.e., authorship, creation dates, etc) was recovered by parsing the front webpage. Most of the metamodels yielded connected grammars, and those that did not, were deliberately designed to consist of several non-connected components.

ReMoDD is a repository for Model-Driven Development, a recent initiative to collect any kinds of models used in papers on MDD [110]. We only fetched models from the artefact category “metamodel” since they have the most straightforward mapping to grammars.

TCS Zoo is a collection of grammars in Textual Concrete Syntax [105] and a part of the Eclipse/GMT library. A TCS model is basically a metamodel annotated with syntactic information.

EMFText Concrete Syntax Zoo is a similar collection of grammars in EMF-Text [106], and it also contains for most (but not all) grammars corresponding Ecore and GenModel metamodels generated from it — if provided, the triple can be converged [27] for external validation.

RelaxNG Schemas is a collection of links: some of them became dysfunctional over time and needed to be located via The Wayback Machine¹⁰, otherwise processing it was similar to TXL, ANTLR and SDF grammars.

¹⁰The Wayback Machine, <http://archive.org/web/web.php>.

ISO and ECMA standards contain grammars in textual form, written in a variety of informally described (E)BNF dialects (far from ISO EBNF), and are often compromised by typesetting artefacts, misspellings, etc. Still, advanced notation-parametric grammar recovery [1] combined with some basic cleanup actions specific for each grammar, was enough to obtain a number of such extracted grammars.

OMG is a standardisation body similar to the two mentioned above. However, it exposes much more of its standards freely online, with minimalistic annotations about their role and format.

SLPS is again an open source project with its repositories also freely exposed to the public, which allows us not only to extract grammars from them in a very robust way, but also collect and analyse the whole history of commits concerning each of them. This was successfully performed as an experiment for some SLPS grammars, but not yet for other sources.

4.7. Grammar extractors

As of now, we have the following grammar extractors available in the Software Language Processing Suite:

ADT to BGF. Without loss of generality, one can assume that abstract data types in Rascal are conceptually the same as in Haskell or any other advanced functional language. In this extractor, types are mapped to non-terminals and constructors are mapped to alternative right hand sides.

Example	Bert Lisser, Dot, lang::dot::Dot, 2012
<pre>public data DotGraph = graph(Id id, Stms stmts) digraph(Id id, Stms stmts);</pre>	
<pre>DotGraph ::= [graph]::([id]::Id [stmts]::Stms) [digraph]::([id]::Id [stmts]::Stms);</pre>	
<p>(All types are treated as nonterminals, constructor names and argument names as labels).</p>	

ANTLR to BGF. In order to be able to extract grammars from ANTLR parser definitions, we reused the standard ANTLR grammar for ANTLR grammars by attaching appropriate semantic actions to it. The semantic actions were programmed for using XML API to serialise the parse tree as a BGF grammar and abstract from the parsed semantic actions.

Example	Oliver Kellogg, Ada, <code>ada.g</code> , 2003
<pre> use_clause : u:USE^ (TYPE! subtype_mark (COMMA! subtype_mark) * { Set(#u, USE_TYPE_CLAUSE); } c_name_list { Set(#u, USE_CLAUSE); }) SEMI! ; </pre>	
<pre> use_clause ::= [u]:USE ((TYPE subtype_mark (COMMA subtype_mark)*) c_name_list) SEMI; </pre>	
<p>(Both context-free and lexical nonterminals become nonterminals, all AST-building annotations are abstracted from, as are semantic predicates and actions).</p>	

DCG to BGF. Definite clause grammars are a way of specifying a parser in Prolog [60]. Their clauses are mapped straightforwardly to production rules by an extractor written also in Prolog.

Example	Ralf Lämmel, FL, <code>Parser.pro</code> , 2008
<pre> function((N,Ns,E)) --> name(N), +(name,Ns), @("="), expr(E), +(newline). </pre>	
<pre> function ::= name name+ "=" expr newline+; </pre>	
<p>(Left hand sides of definite clauses become left hand sides of a grammar, all predicates on the right hand side that are not recognised as metasyntax, are treated as nonterminals, except for “@”, “reserved” and “keyword” that are treated as specifying terminal symbols. All information on variables is abstracted from).</p>	

Ecore to BGF. Since Ecore models are by default serialised as XMI, we only needed to express the mapping between Ecore and BGF, which was done in XSLT.

Example	Hugo Brunelière, Cobol, <code>COBOL.ecore</code> , 2005
<pre> <eClassifiers xsi:type="ecore:EClass" name="COBOL88Element"> <eStructuralFeatures xsi:type="ecore:EAttribute" name="name" ordered="false" unique="false" lowerBound="1" eType="/1/String"/> <eStructuralFeatures xsi:type="ecore:EReference" name="has" ordered="false" lowerBound="1" upperBound="-1" eType="/0/COBOL88ElementValue" containment="true"/> </eClassifiers> </pre>	
<pre> COBOL88Element ::= [name]::String <has>:COBOL88ElementValue+; </pre>	
<p>(The extractor works with the XMI serialisation of an Ecore model. All non-abstract <code>eClassifiers</code> tags of <code>EClass</code> type become nonterminals, with each of <code>eStructuralFeatures</code> referring to a named element in the sequence on the right hand side in the grammar production rule).</p>	

Java to BGF. The object model of a Java program is extracted from a Java source by the use of reflection. Classes are treated as nonterminals, and their visible interfaces (public members and getters/setters) serve as the right hand side. This mapping helped to trivially converge the structure

defined by the Java source generated by a data binding framework (JAXB) with the structure defined by the original schema (XSD) in [27].

Example	Ralf Lämmel, FL, types/Binary.java, 2008
<pre>public class Binary extends Expr { public Ops ops; public Expr left, right; public Binary(Ops o, Expr left, Expr right) { this.ops = o; this.left = left; this.right = right; } public void accept(Visitor v) { v.visit(this); } }</pre>	
<pre>Expr ::= Binary; Binary ::= [ops]::Ops [left]::Expr [right]::Expr;</pre>	
<p>(Each class is mapped to a nonterminal. Inheritance is treated as a chain production rule. Every public element of a class becomes a sequentially composed element on the right hand side of the grammar production rule, with types becoming nonterminals and names becoming selectors).</p>	

LDF to BGF. Since we assume that any language document does contain grammar knowledge explicitly, i.e., in BGF, we use a special extractor to take out the BGF bits and compose a grammar from them. In the past, the LDF to BGF extractor was mostly used for testing purposes. We leave it without an example since extraction from LDF is basically selective copying.

LLL to BGF. The first extractor from LLL was developed as a means of importing grammars manipulated by GDK [72]. Later it has been retired in favour of an EBNF Dialect Definition (EDD) of LLL that serves as a parameter for Grammar Hunter (see below).

Example	Vadim Zaytsev, C#, Final.111 extracted from ECMA-344, 2005
<pre>goto-statement : "goto" lex-csharp/identifier ";" "goto" "case" expression ";" "goto" "default" ";" ;</pre>	
<pre>goto-statement ::= "goto" lex-csharp/identifier ";" "goto" "case" expression ";" "goto" "default" ";";</pre>	
<p>(The notations are very close: BGF covers strictly more metasyntactic features than LLL).</p>	

Python to BGF. A Python library called PyParsing allows to define a PEG inside Python code. This extractor, written in Rascal, relies on the structure expected by PyParsing, in order to recover grammar knowledge from a Python program.

Example	Ruwen Hahn, FL, <code>parser.py</code> , 2008
<pre>ifThenElse = (_IF + expr + _THEN + expr + _ELSE + expr).setParseAction(lambda tok: t.IfThenElse(*tok))</pre>	
<pre>ifThenElse ::= _IF expr _THEN expr _ELSE expr;</pre>	
<p>(The extractor abstracts from all parse actions of the PyParsing library and treats all components as nonterminals, while also recognising combinators for sequential composition (“~” and “+”), optionality (<code>Optional</code>), negation (<code>NotAny</code>), Kleene star (<code>ZeroOrMore</code>), etc).</p>	

Rascal grammar to BGF. By reusing Rascal grammar for Rascal and internal interfaces for accessing it, this extractor delivers a platonic grammar extracted from a Rascal [26] concrete syntax definition (which is essentially an annotated parser specification).

Example	Jurgen Vinju, Pico, <code>lang::pico::syntax::Main</code> , 2011–2012
<pre>syntax Expression = id: Id name strcon: String string natcon: Natural natcon bracket "(" Expression e ")" > left concat: Expression lhs " " Expression rhs > left (add: Expression lhs "+" Expression rhs min: Expression lhs "-" Expression rhs) ;</pre>	
<pre>Expression ::= [id>::([name)::Id) [strcon>::([string)::String) [natcon>::([natcon)::Natural) "(" e:Expression ")" <left>:[concat>::([lhs)::Expression " " [rhs)::Expression) <left>:[add>::([lhs)::Expression "+" [rhs)::Expression) [min>::([lhs)::Expression "-" [rhs)::Expression]);</pre>	
<p>(The extractor abstracts from associativity rules and priority specifications and treats names of alternatives as labels. Any annotations are also neglected during extraction).</p>	

RelaxNG schema to BGF. RELAX NG is a schema language for XML, alternative to a much more popular XML Schema [107]. It is in general believed to be simpler and thus more understandable and making the users less prone to mistakes, but all commonly used features are present in both languages, making the choice between them usually platform-dependent or library-dependent. The extractor was developed incrementally by covering the functionality found in the schemata referenced at the RELAX NG home page [111]; it was not a scientifically challenging process.

Example	Norman Walsh, DocBook, dbhier.rng, 2002
<pre> <define name="tocfront"> <element name="tocfront"> <ref name="tocfront.attlist"/> <zeroOrMore> <ref name="para.char.mix"/> </zeroOrMore> </element> </define> </pre>	
<pre> tocfront ::= [tocfront]::(tocfront.attlist para.char.mix*); </pre>	
<p>(All <code>define</code> clauses with a nested <code>element</code> are mapped to labelled production rules. References to other definitions are treated as nonterminals, repetition kinds are also mapped straightforwardly. The <code>interleave</code> construct is abstracted from and mapped to a sequence, which is still a reasonably good representation in the case of an abstract grammar. Mixed content is mapped to a Kleene star over the choice of the actual content and a raw string).</p>	

SDF to BGF. We encoded the necessary traversal functions for crawling the parse trees of SDF grammars and producing BGF and reused the SDF module and the XML module from the standard package of the Meta-Environment. Separate command line tools are used to make a parse table, to compile ASF formulæ, to parse the input grammar, to rewrite the parse tree and to serialise the transformed parse tree into a file. They are bundled together and wrapped in a black box extraction script.

Example	Jurgen Vinju, Taeke Kooiker, Mark van den Brand, C, ansi-c/syntax/Declarations.sdf, 2006–2008
<pre> context-free syntax Specifier+ {InitDeclarator ","}+ ";" -> Declaration Specifier+ ";" -> Declaration {avoid} </pre>	
<pre> Declaration ::= Specifier+ {InitDeclarator ","}+ ";" ; Declaration ::= Specifier+ ";" ; </pre>	
<p>(The notations are very similar, except the right hand side and the left hand sides are flipped to a more common position, and priority specifications are neglected. Constructor names, if present, becomes labels for production rules).</p>	

W3C Specification to BGF. As a part of the initiative to create a unified data model for language documentation, a case study was completed to map the W3C specification of XPath to it [31]. This extractor looks for all `<scrap>` elements inside a W3C standard specification and maps its `<prod>` elements to production rules.

Example	W3C, XPath, REC-xpath-19991116.xml, 1999
<pre> <prod id="NT-LocationPath"> <lhs>LocationPath</lhs> <rhs> <nt def="NT-RelativeLocationPath">RelativeLocationPath</nt> </rhs> <rhs> <nt def="NT-AbsoluteLocationPath">AbsoluteLocationPath</nt> </rhs> </prod> </pre>	
<pre> LocationPath ::= [NT-LocationPath]::(RelativeLocationPath AbsoluteLocationPath); </pre>	
<p>(Production rules are formed from <code><prod></code> tags, with their IDs becoming production labels and each reference to a nonterminal being mapped to a nonterminal).</p>	

TXL to BGF. We reused the TXL grammar for TXL grammars and made use of the TXL engine's option of returning the parse tree in an XML form. The mapping between TXL XML and our XML (i.e., BGF) was straightforwardly encoded in XSLT.

Example	William Waite, James Cordy, Fortran, <code>fortran.grm</code> , 2009
<pre> define BlockDataStmt [LblDef] 'blockdata [opt BlockDataName] [EOS] end define </pre>	
<pre> BlockDataStmt ::= LblDef "blockdata" BlockDataName? EOS; </pre>	
<p>(Each <code>define</code> and <code>redefine</code> clause is mapped to a production rule. All combinators of the latest version of TXL (<code>opt</code>, <code>repeat</code>, <code>list</code>, etc) are supported. All rule clauses are ignored).</p>	

XML Schema to BGF. Not all elements of the XML Schema can be mapped to grammar concepts efficiently, but the most used ones easily find their counterparts. The mapping is thus partial and bidirectional at best (e.g., XML elements and XSD complex types are both mapped to nonterminal symbols).

Example	Vadim Zaytsev, LDF, <code>ldf.xsd</code> , 2010
<pre> <xsd:complexType name="named-link"> <xsd:sequence> <xsd:element name="title" type="xsd:string"/> <xsd:choice minOccurs="0"> <xsd:element name="version" type="xsd:string"/> <xsd:element name="edition" type="xsd:string"/> </xsd:choice> <xsd:element name="uri" type="xsd:anyURI" minOccurs="0"/> </xsd:sequence> </xsd:complexType> </pre>	
<pre> named-link ::= [title]::str ([version]::str [edition]::str)? [uri]::str?; </pre>	
<p>(All elements, groups, complex and simple types are mapped to nonterminals. Combinations of <code>minOccurs</code> and <code>maxOccurs</code> constraints are mapped to the closest possible variant available in BGF: “+”, “*” or “?”).</p>	

4.8. Grammar renderers

Technically, a grammar renderer (pretty-printer, exporter) is a part of a bidirectional transformation [112, 113] between the source grammar specified in a notation of a particular framework or methodology, and the platonic piece of grammatical knowledge stored in the repository. When the grammar is recovered, corrected and otherwise changed, the two become desynchronised and can be brought back in sync by a renderer. A more simplistic view would be to assume that the source grammar is lost or empty, and that the synchronisation process only entails generating the corresponding definition entirely. To illustrate this, suppose that we have a syntax definition D in SDF [56]: it contains a context-free part, a lexical part, possibly variable definitions, priority descriptions, AST constructor annotations and similar details, which are together enough to parse textual data unambiguously to form a tree correctly representing its structure in the described language. A Grammar Zoo entry G extracted from that source, would have less information, because it has undergone through the “extraction through abstraction” process [27] and lost its idiosyncratic details. Still, if something changes — for instance, a production rule is removed from the extracted grammar G , we can automatically detect which rule needs to be removed from D if we wish them synchronised. If we represent the bidirectional transformation classically as a relation and two update functions, then the pair of D and G will be an element of this relation and two functions would be the extractor and the renderer.

For the sake of simplicity we adopt the simplistic view and consider that a rendered grammar is generated from scratch from a repository entity and does not need to be matched with its older version and updated. SLPS contains such simple renderers for exporting grammars in a readable EBNF dialect like the one displayed in this paper; in the BNF dialect used in the DMS Toolkit [114]; in Graphviz dot format for visualisation purposes; in Rascal as a syntax definition [26]; in SDF [56]; in TXL [61]; in LaTeX for publishing purposes, etc. Conceptually each one of those is a reverse of an extractor — that is, a pretty-printer. We will only consider hypertext rendering in detail.

Browsable grammars are not a new idea: in particular, they have been motivated in [38] and [115], and even the bare observation of the abundance of such grammars being displayed all over the internet, with metasyntax highlighting and hyperlinked nonterminals, demonstrates practical need for them.

The hypertext rendering of a grammar in the Grammar Zoo, besides the main part, contains the following additional computed information:

- Number of production rules — the actual number of different definitions for nonterminals physically present in the grammar;
- Number of top alternatives — the PROD metric from [39], calculated as the number of top level alternatives in order to account for “horizontal” definitions (i.e., “ $X ::= Y \mid Z$,” instead of “ $X ::= Y$; $X ::= Y$ ”);
- Number of defined nonterminal symbols — the number of unique left hand sides of all production rules;

- Root nonterminal symbols — nonterminals explicitly marked as starting symbols of a grammar;
- Other top nonterminal symbols — defined nonterminals not used anywhere in a grammar and not marked as roots; in [38] it is suggested that top nonterminals should be examined during grammar recovery and either elevated to the root status or removed from the grammar;
- Bottom nonterminal symbols — nonterminals forming the difference between the abovementioned number of defined nonterminal symbols and the VAR metric from [39]; undefined nonterminals encountered exclusively within the right hand sides of production rules;
- Number of used terminal symbols — the TERM metric from [39].

The question of how to render a grammar perfectly, is an open one, but these properties have been quoted as helpful for grammar comprehension: they help a human engineer to estimate the grammar’s complexity and perform some superficial analyses like examination and elimination of multiple top nonterminals.

4.9. Illustration: convergence of JLS grammars

The method of grammar convergence is used to reverse engineer true relationships between grammars by transforming them toward equality and examining the properties of these transformation steps (in particular, preservation of the assumed semantics). For example, if the only changes between two grammars concern renaming nonterminals, we conclude that the grammars are equivalent in the sense of generating or accepting the same string language, but not equivalent in the sense of XML document types. When the method of grammar convergence was first proposed in [27], it was demonstrated on a small case study with six grammars of such a small size that all of them would fit on one page. Obviously, the next step in demonstrating the viability of the method was to perform a full scale case study with industrial size grammars of a mainstream language: Java was chosen as one of the languages having three published editions of the official language specification, each containing two formally unrelated alternative grammars for the same language version.

In the Java convergence paper [28], one quarter of it, 11 full pages out of 42 (not counting the bibliography), is dedicated to the extraction and recovery process. However, the recovery of the six grammars from three editions of the Java Language Specification is not among the main contributions of the paper, so the amount of attention it received is due to the sheer complexity of the process, which entailed:

- tag elimination — source L0 grammars were typed in manually in ill-formed HTML;
- indentation processing — the input notation relied on indentation to separate top level choices in production rules;

- robust parsing — in order to deal with incorrect markup and unexpected artefacts;
- matching parentheses — to disambiguate brackets-as-terminals from brackets-as-metasympols;
- adjusting symbol roles — incorrect markup and improper context could lead to metasympols being treated as terminal symbols or vice versa, or terminals as nonterminals;
- composing sibling symbols and decomposing compound symbols — to address the issue of erroneous markup being interjected in the middle of a token;
- removing duplicates — also those with slightly different typesetting.

If anyone is to replicate the case study and, for instance, to converge several C# grammars taken from different sources (ECMA standards, ISO standards and Microsoft language specifications), the whole effort would need to be reinvested both in the engineering side of the case study (the actual programming of an appropriate extractor) and in the description of it (even if only for the sake of honest report). By saving this kind of information together with the extracted grammar in a publicly accessible place, we save the space in any future replication and allow such papers to fully focus on their main contributions.

4.10. Grammar recovery and evolution

All of the tools listed so far in §4.7, did not go beyond simple extraction of a level 1 grammar: that is, in the presence of an error in the grammar source used for extraction, they give up after reporting it to the user, which then has to go back to the source and figure out a way to solve it. Unlike them, the tools listed below are capable to identify and even resolve some of the commonly encountered issues.

BNF to LLL. After we have noticed that many ISO standards of programming languages (C, C++, C#) share the same metalanguage, this was the first tool to be developed. It normalised some lexical singularities and translated the EBNF dialect used by ISO grammar developers to LLL used in GDK [72] (for which we have another extractor ready). It should be noted here that the EBNF used in ISO standards is not the same as ISO EBNF defined by [80].

PDF to BGF. A grammar copy-pasted from a PDF of an ISO standard of C, C++, C# or any other that uses the same metalanguage, can be extracted with some tolerance regarding lexical imperfections. This extractor is essentially a composition of BNF to LLL and LLL to BGF.

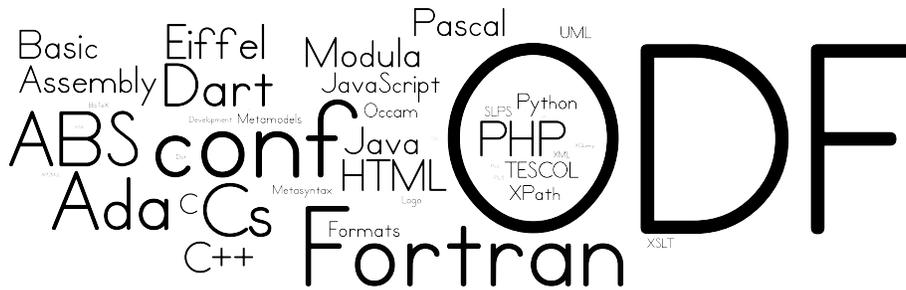


Figure 3: A tag cloud showing language and language group names according to the average sizes of their grammars (made with [Wordle](#)). “Cs” denotes C#, “conf” is a collection of metamodels for the domain of conference organisation. “ODF” is Open Document Format.

HTML to BGF. This advanced extractor had to work on a manually and loosely hypertext source. It comprised a set of generalised heuristics in a pattern form that it tried to apply for automated recovery. In the Java Language Specification case study [28] its use was prolific, the extractor fixed 669 errors before we started to program the main body of grammar transformations.

EDD to Rascal. This tool aids semi-automated interactive grammar recovery [1]. It requires a specification of the input metalanguage in a form of EDD, an EBNF Dialect Definition [30], from which it generates a Rascal plugin that enables manipulation of grammars written in the specified EBNF dialect with standard means of Eclipse. This method does not automatically solve any problems, but it helps identifying them and facilitates a grammar engineer in fixing them.

Grammar Hunter. This tool used for notation-parametric grammar recovery [1], requires a specification of the input metalanguage in a form of EDD, an EBNF Dialect Definition [30]. Then it consumes the input text, treating it as a grammar text written in the specified EBNF dialect, applies all appropriate heuristics and delivers a recovered grammar automatically.

XBGF scripts are more or less a standard way of programmable grammar manipulations in SLPS: they can be generated or manually developed, and always can be re-executed or pretty-printed for inspection. The XBGF language is described in detail in [28], [57, §7] and [24, XBGF Manual]. Essentially it allows to use predefined operators for (un)folding nonterminals, factoring choices, renaming symbols, etc.

5. Using the Zoo

[Figure 3](#) visualises the software languages currently represented in the Grammar Zoo. In order to avoid multiple diagrams, we have used a cumulative size

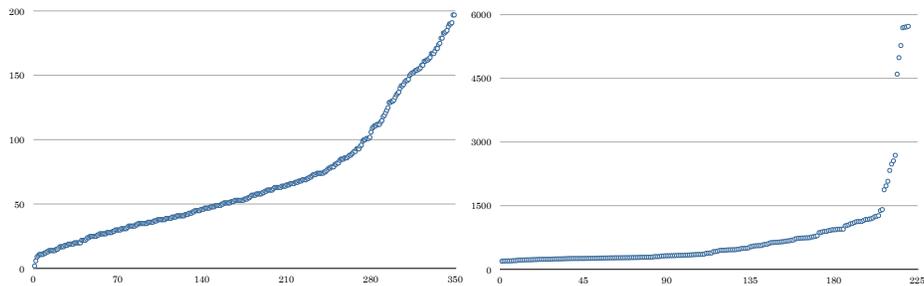


Figure 4: Sizes of individual extracted grammars in Grammar Zoo, in ascending order (up to 200 on the left, from 200 up on the right).

metric conceptually akin to Halstead vocabulary [116] and calculated as the sum of the number of unique nonterminal symbols, the number of unique terminal symbols, the number of unique labels for production rules and top alternatives, and the number of selector names for grammatical subexpressions. As one can see from Figure 4, the Grammar Zoo contains grammars of all sizes, so clustering based on this cumulative metric is unfeasible. The outrageously large outliers on both figures are Open Document Format grammars, all containing around 1000 nonterminals, around 700 terminals and around 2000 production rules.

A more detailed summary of the minimum, average and maximum sizes of grammars per language group, can be found on Figure 5: the smallest grammars of ODF are those of its Manifest Schema and its Digital Signature Schema, so a “language” in this context should be understood as a domain, not as a string language in the formal grammar sense.

Unfortunately, the entire contents of the Grammar Zoo are too long and too actively changing in order to present them in a paper. They are, however, always available online at <http://slps.github.io/zoo>.

Future possible uses for the Grammar Zoo and any other grammar repository built on the same principles we have proposed in this paper, are diverse and include, but are not limited to, the following:

- Improving the overall accessibility of the repository by providing grammars in more notations, with better visualisation and filtering strategies.
- Assimilating other grammar collections — for instance, the Grammar Zoo contains several ANTLR grammars from the ANTLR Grammar List [67], but not yet all of them.
- Curating the corpus and raising the maturity level [2] of grammars there for the sake of enabling bulk processing of only connected or only adapted grammars, in the manner that was performed in [20] with differential testing of parsers.
- Adding a time dimension: for grammars which source is a repository, all versions can be extracted and somehow presented in the Grammar Zoo.

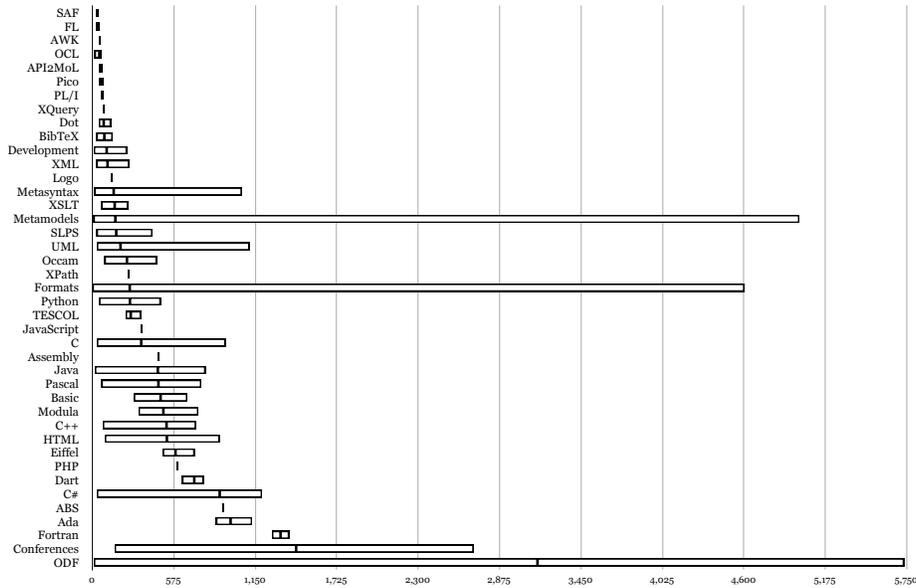


Figure 5: Grammar sizes: minimum, average and maximum — per language group. By “grammar size” we understand a cumulative metric summing the number of nonterminals, number of terminals and number of expression labels in a grammar. Such a remarkably small size for PL/I can be explained by the fact that at the moment we only have one simplified metamodel of the language.

This has been done as an experiment for some of the DSL grammars used within SLPS: BGF, a BNF-like grammar format, and LDF, a unified format for language documentation.

- Mining the corpus for various properties — for example, defining a collection of micropatterns based on the current content of the Grammar Zoo has been done in [22].
- Performing comparative studies on grammar-based techniques.
- Empirical studies of grammar improvement: user studies or genetic algorithms.
- Developing methods of inferring notation specifications.
- Adding interactive on demand grammar export instead of providing a static collection of predefined notations.
- Bridging technological spaces by investigating grammars from them, megamodelling the corpus itself and grammars in it.
- Improving interoperability of metagrammarware.

6. Conclusion

Many claims about design of software languages and their descriptions (grammars in a broad sense, per [3, 57]) can be found in existing literature [6, 7, 8, 9, 10, 11, 12, 13, 14, 15, et al.], all backed up by case studies and expert opinions. In order to support, challenge or generalise them, one could collect relevant statistical evidence, based on a sufficiently large corpus of diverse language definitions. Such a corpus, referred to as the **Grammar Zoo**, was proposed in this paper. This corpus will support replicability across grammar-based experiments (by providing open access to its grammars to the public), support aggregation of findings (by annotating the grammars extensively), reduce the cost of controlled experiments (by shifting the focus of future research from obtaining the grammars for case studies to newly proposed techniques), aid to obtain sample representativeness (by making the grammarbase large and versatile) and help to isolate the effects of individual factors (by metadata-based filtering). The Grammar Zoo initiative is modelled after successful projects like Qualitas Corpus [4] in source code analysis and Atlantic Metamodel Zoo [5] in modelware.

By relying on previous experiences in grammar extraction and grammar recovery [44, 37, 45, 38, 43, 48, 27, 28, 1, 30], we provide methods and tools for relatively semi-automatic easy extraction of grammars in a broad sense from various software artefacts: syntax specifications, type definitions, data schemata, etc. These methods and tools rely on systematic manipulation of syntactic notations, on reproducible specification of grammar evolution steps, on advanced IDE support, as well as on some other less recently developed technologies. A unified data model for systematic accumulation of grammar knowledge has been designed, presented and exemplified.

The Grammar Zoo, publicly available as <http://slps.github.io/zoo>, is a collection of big grammars of mainstream languages such as Ada, C, C++, C#, Dart, Modula, Fortran, as well as small grammars used for various purposes, mostly for demonstrating certain software language engineering techniques on a proof-of-concept scale. Its name stems from the activity known as “*grammar hunting*” [1] or “*grammar stealing*” [45] and hints at the fact that the result of a hunt is not cooked, eaten and gone, but rather carefully put on display. The Grammar Zoo at the time of submission has 1710 fetched grammars and 500+ grammars of extracted level and higher, and it continues growing.

References

- [1] V. Zaytsev, Notation-Parametric Grammar Recovery, in: A. Sloane, S. Andova (Eds.), Post-proceedings of the 12th International Workshop on Language Descriptions, Tools, and Applications (LDTA 2012), ACM Digital Library, 2012, pp. 9:1–9:8. doi:10.1145/2427048.2427057.
- [2] V. Zaytsev, Grammar Maturity Model, in: A. Pierantonio, D. Tamzalit, B. Schätz (Eds.), Ninth Workshop on Models and Evolution (ME’14), CEUR, 2014.
- [3] P. Klint, R. Lämmel, C. Verhoef, Toward an Engineering Discipline for Grammarware, ACM Transactions on Software Engineering Methodology (TOSEM) 14 (3) (2005) 331–380.

- [4] E. Tempero, C. Anslow, J. Dietrich, T. Han, J. Li, M. Lumpe, H. Melton, J. Noble, Qualitas Corpus: A Curated Collection of Java Code for Empirical Studies, in: Asia Pacific Software Engineering Conference (APSEC 2010), 2010, pp. 336–345.
- [5] J. Cabot, M. Tisi, H. Brunelière, et al., AtlantEcore Metamodel Zoo, <http://www.emn.fr/z-info/atlanmod/index.php/Ecore> (2003).
- [6] A. van Wijngaarden, Generalized ALGOL, in: R. Goodman (Ed.), Annual Review in Automatic Programming 3, Pergamon Press, 1963, pp. 17–26.
- [7] N. Wirth, On the Design of Programming Languages, in: IFIP Congress, 1974, pp. 386–393.
- [8] C. A. R. Hoare, Hints on Programming Language Design, Tech. rep., Stanford University, Stanford, CA, USA (1973).
- [9] M. Mernik, J. Heering, A. M. Sloane, When and How to Develop Domain-Specific Languages, ACM Computing Surveys 37 (4) (2005) 316–344. doi:10.1145/1118890.1118892.
- [10] M. Völter, S. Benz, C. Dietrich, B. Engelmam, M. Helander, L. C. L. Kats, E. Visser, G. Wachsmuth, DSL Engineering: Designing, Implementing and Using Domain-Specific Languages, dslbook.org, 2013.
- [11] N. Chomsky, Syntactic Structures, Mouton, 1957.
- [12] M. Erwig, E. Walkingshaw, Semantics First!, in: Proceedings of the Fourth International Conference on Software Language Engineering, SLE’11, Springer, 2012, pp. 243–262. doi:10.1007/978-3-642-28830-2_14.
- [13] L. Tratt, Evolving a DSL Implementation, in: Generative and Transformational Techniques in Software Engineering, Vol. 5235 of LNCS, Springer, 2008, pp. 425–441. doi:10.1007/978-3-540-88643-3_11.
- [14] M. Herrmannsdörfer, S. D. Vermolen, G. Wachsmuth, An Extensive Catalog of Operators for the Coupled Evolution of Metamodels and Models, in: Proceedings of the Third International Conference on Software Language Engineering, SLE’10, Springer, 2011, pp. 163–182.
- [15] J. Hutchinson, J. Whittle, M. Rouncefield, S. Kristoffersen, Empirical Assessment of MDE in Industry, in: Proceedings of the 33rd International Conference on Software Engineering, ICSE ’11, ACM, 2011, pp. 471–480. doi:10.1145/1985793.1985858.
- [16] H. Do, S. Elbaum, G. Rothermel, Supporting Controlled Experimentation with Testing Techniques: An Infrastructure and its Potential Impact, Journal of Empirical Software Engineering 10 (4) (2005) 405–435.
- [17] C. Wohlin, P. Runeson, M. Höst, M. C. Ohlsson, B. Regnell, A. Wesslén, Experimentation in Software Engineering: an Introduction, Kluwer Academic Publishers, Norwell, MA, USA, 2000.
- [18] V. R. Basili, F. Shull, F. Lanubile, Building Knowledge through Families of Experiments, IEEE Transactions on Software Engineering 25 (4) (1999) 456–473.
- [19] L. M. Pickard, B. A. Kitchenham, P. W. Jones, Combining Empirical Results in Software Engineering, Journal of Information and Software Technology 40 (14) (1998) 811–821.
- [20] B. Fischer, R. Lämmel, V. Zaytsev, Comparison of Context-free Grammars Based on Parsing Generated Test Data, in: U. Aßmann, A. Sloane (Eds.), Post-proceedings of the Fourth International Conference on Software Language Engineering (SLE 2011), Vol. 6940 of LNCS, Springer, 2012, pp. 324–343. doi:10.1007/978-3-642-28830-2_18.

- [21] J. Gil, I. Maman, Micro Patterns in Java Code, in: Proceedings of OOPSLA'05, ACM, 2005, pp. 97–116.
- [22] V. Zaytsev, Micropatterns in Grammars, in: M. Erwig, R. F. Paige, E. V. Wyk (Eds.), Proceedings of the Sixth International Conference on Software Language Engineering (SLE 2013), Vol. 8225 of LNCS, Springer, 2013, pp. 117–136. doi:10.1007/978-3-319-02654-1_7.
- [23] A. Johnstone, P. D. Mosses, E. Scott, An Agile Approach to Language Modelling and Development, Innovations in Systems and Software Engineering 6 (1-2) (2010) 145–153. doi:10.1007/s11334-009-0111-6.
- [24] V. Zaytsev, R. Lämmel, T. van der Storm, L. Renggli, R. Hahn, G. Wachsmuth, Software Language Processing Suite¹¹, <http://slps.github.io>. Contains, among other works: Grammar Zoo (V. Zaytsev, 2009–2014), <http://slps.github.io/zoo>. (2008–2014).
- [25] V. Zaytsev, GrammarLab, <http://grammarware.github.io/lab> (2013–2014).
- [26] P. Klint, T. van der Storm, J. Vinju, EASY Meta-programming with Rascal, in: J. M. Fernandes, R. Lämmel, J. Visser, J. Saraiva (Eds.), Post-proceedings of the Third International Summer School on Generative and Transformational Techniques in Software Engineering (GTTSE 2009), Vol. 6491 of LNCS, Springer, 2011, pp. 222–289.
- [27] R. Lämmel, V. Zaytsev, An Introduction to Grammar Convergence, in: M. Leuschel, H. Wehrheim (Eds.), Proceedings of the Seventh International Conference on Integrated Formal Methods (iFM 2009), Vol. 5423 of LNCS, Springer, 2009, pp. 246–260. doi:10.1007/978-3-642-00255-7_17.
- [28] R. Lämmel, V. Zaytsev, Recovering Grammar Relationships for the Java Language Specification, Software Quality Journal (SQJ) 19 (2) (2011) 333–378. doi:10.1007/s11219-010-9116-5.
- [29] V. Zaytsev, [Software Language Engineering by Intentional Rewriting](#), Electronic Communications of the European Association of Software Science and Technology (EC-EASST); Software Quality and Maintainability 65. URL <http://journal.ub.tu-berlin.de/eceasst/article/view/903>
- [30] V. Zaytsev, BNF WAS HERE: What Have We Done About the Unnecessary Diversity of Notation for Syntactic Definitions, in: S. Ossowski, P. Lecca (Eds.), Programming Languages Track, Volume II of the Proceedings of the 27th ACM Symposium on Applied Computing (SAC 2012), ACM, Riva del Garda, Trento, Italy, 2012, pp. 1910–1915. doi:10.1145/2245276.2232090.
- [31] V. Zaytsev, R. Lämmel, A Unified Format for Language Documents, in: B. A. Malloy, S. Staab, M. G. J. van den Brand (Eds.), Post-proceedings of the Third International Conference on Software Language Engineering (SLE 2010), Vol. 6563 of LNCS, Springer, 2011, pp. 206–225. doi:10.1007/978-3-642-19440-5_13.
- [32] V. Zaytsev, [Language Evolution, Metasyntactically](#), Electronic Communications of the European Association of Software Science and Technology (EC-EASST) 49. URL <http://journal.ub.tu-berlin.de/eceasst/article/view/708>
- [33] J. Bézivin, F. Jouault, P. Valduriez, On the Need for Megamodels, OOPSLA & GPCE, Workshop on best MDSO practices, 2004.

¹¹The authors are given according to the list of contributors at <http://github.com/grammarware/slps/graphs/contributors>.

- [34] J.-M. Favre, R. Lämmel, A. Varanovich, Modeling the Linguistic Architecture of Software Products, in: Proceedings of the 15th international conference on Model Driven Engineering Languages and Systems (MoDELS), LNCS, Springer, 2012, pp. 151–167.
- [35] V. Zaytsev, Renarrating Linguistic Architecture: A Case Study, in: C. Hardebolle, E. Syriani, J. Sprinkle, T. Mészáros (Eds.), Post-proceedings of the Sixth International Workshop on Multi-Paradigm Modeling (MPM 2012), ACM Digital Library, 2012, pp. 61–66. doi:10.1145/2508443.2508454.
- [36] R. Lämmel, V. Zaytsev, [Language Support for Megamodel Renarration](#), in: J. De Lara, D. Di Ruscio, A. Pierantonio (Eds.), Post-proceedings of the Second Workshop on Extreme Modeling (XM 2013), Vol. 1089 of CEUR Workshop Proceedings, CEUR-WS.org, 2013, pp. 36–45.
URL <http://ceur-ws.org/Vol-1089/5.pdf>
- [37] M. G. J. van den Brand, M. P. A. Sellink, C. Verhoef, Obtaining a COBOL Grammar from Legacy Code for Reengineering Purposes, in: M. P. A. Sellink (Ed.), Proceedings of the Second International Workshop on the Theory and Practice of Algebraic Specifications, Springer, 1997, pp. 6–17.
- [38] R. Lämmel, C. Verhoef, Semi-automatic Grammar Recovery, *Software—Practice & Experience* 31 (15) (2001) 1395–1438.
- [39] J. F. Power, B. A. Malloy, A Metrics Suite for Grammar-based Software, *Journal of Software Maintenance and Evolution: Research and Practice* 16 (2004) 405–426.
- [40] C. Julien, M. Črepinšek, R. Forax, T. Kosar, M. Mernik, G. Roussel, On Defining Quality Based Grammar Metrics, in: Proceedings of the International Multiconference on Computer Science and Information Technology, IMCSIT 2009, 2009, pp. 651–658.
- [41] R. Lämmel, W. Schulte, Controllable Combinatorial Coverage in Grammar-Based Testing, in: Proceedings of the 18th IFIP TC6/WG6.1 International Conference on Testing of Communicating Systems (TestCom’06), Vol. 3964 of LNCS, Springer, 2006, pp. 19–38.
- [42] P. I. Manuel, ANSI Cobol III in SDF + an ASF Definition of a Y2K Tool, Master’s thesis, Universiteit van Amsterdam, The Netherlands (Nov. 1996).
- [43] R. Lämmel, C. Verhoef, Browsable grammars, <http://www.cs.vu.nl/grammarware/browsable/>, contains: VS COBOL II grammar Version 1.0.4 (Lämmel, Verhoef, 1999–2003), <http://www.cs.vu.nl/grammarware/browsable/vs-cobol-ii>; COBOL grammar Version 0.1.1 (Lämmel, Verhoef, 1999), <http://www.cs.vu.nl/grammarware/browsable/cobol>; OS PL/I V2R3 grammar Version 0.1 (Lämmel, Verhoef, 1999), <http://www.cs.vu.nl/grammarware/browsable/os-pli-v2r3>; Browsable Ada 95 Grammar (Lämmel, Verhoef, 2000), <http://www.cs.vu.nl/grammarware/browsable/ada>; C# Grammar Recovered (Zaytsev, 2005), <http://www.cs.vu.nl/grammarware/browsable/CSharp> (1999).
- [44] M. P. A. Sellink, C. Verhoef, Development, Assessment, and Reengineering of Language Descriptions, in: J. Ebert, C. Verhoef (Eds.), Proceedings of the Fourth European Conference on Software Maintenance and Reengineering (CSMR 2000), IEEE Computer Society, 2000, pp. 151–160.
- [45] R. Lämmel, C. Verhoef, Cracking the 500-Language Problem, *IEEE Software* (2001) 78–88.
- [46] IBM Library, [SX26-3721-05: VS COBOL II Application Programming Reference Summary, Release 4](#) (1987).
URL http://publibz.boulder.ibm.com/cgi-bin/bookmgr_OS390/BOOKS/IGYR1101

- [47] Standard ECMA-334, *C# Language Specification*, 4th Edition, available at <http://ecma-international.org/publications/standards/Ecma-334.htm> (June 2006).
- [48] V. Zaytsev, Correct C# Grammar too Sharp for ISO, in: Participants Workshop, Part II of the Pre-proceedings of the International Summer School on Generative and Transformational Techniques in Software Engineering (GTTSE 2005), Technical Report, TR-CCTC/DI-36, Universidade do Minho, Braga, Portugal, 2005, pp. 154–155, extended abstract.
- [49] J. Gosling, B. Joy, G. L. Steele, *The Java Language Specification*, Addison-Wesley, 1996, available at <http://java.sun.com/docs/books/jls>.
- [50] J. Gosling, B. Joy, G. L. Steele, G. Bracha, *The Java Language Specification*, 2nd Edition, Addison-Wesley, 2000, available at <http://java.sun.com/docs/books/jls>.
- [51] J. Gosling, B. Joy, G. L. Steele, G. Bracha, *The Java Language Specification*, 3rd Edition, Addison-Wesley, 2005, available at <http://java.sun.com/docs/books/jls>.
- [52] V. Zaytsev, *MediaWiki Grammar Recovery*, Computing Research Repository (CoRR) 1107.4661 (2011) 1–47.
- [53] V. Zaytsev, *The Grammar Hammer of 2012*, Computing Research Repository (CoRR) 1212.4446 (2012) 1–32.
- [54] N. Wirth, What Can We Do about the Unnecessary Diversity of Notation for Syntactic Definitions?, *Communications of the ACM* 20 (11) (1977) 822–823.
- [55] A. Stevenson, J. R. Cordy, *A Survey of Grammatical Inference in Software Engineering*, Science of Computer Programming In print.
- [56] J. Heering, P. R. H. Hendriks, P. Klint, J. Rekers, *The Syntax Definition Formalism SDF—Reference Manual*, ACM SIGPLAN Notices 24 (11) (1989) 43–75.
- [57] V. Zaytsev, *Recovery, Convergence and Documentation of Languages*, Ph.D. thesis, Vrije Universiteit (Oct. 2010).
- [58] E. Visser, *Syntax Definition for Language Prototyping*, Ph.D. thesis, University of Amsterdam (Sep. 1997).
- [59] T. Parr, *The Definitive ANTLR Reference: Building Domain-Specific Languages*, 1st Edition, Pragmatic Programmers, Pragmatic Bookshelf, 2007.
- [60] F. Pereira, D. Warren, *Definite Clause Grammars for Language Analysis*, in: B. J. Grosz, K. Sparck-Jones, B. L. Webber (Eds.), *Readings in Natural Language Processing*, Morgan Kaufmann Publishers Inc., San Francisco, CA, USA, 1986, pp. 101–124.
- [61] T. R. Dean, J. R. Cordy, A. J. Malton, K. A. Schneider, *Grammar Programming in TXL*, in: *Proceedings of the Second IEEE International Conference on Source Code Analysis and Manipulation (SCAM 2002)*, IEEE, 2002, pp. 93–102.
- [62] S. Gao, C. M. Sperberg-McQueen, H. S. Thompson, N. Mendelsohn, D. Beech, M. Maloney, *W3C XML Schema Definition Language (XSD) 1.1 Part 1: Structures*, W3C Candidate Recommendation. URL <http://www.w3.org/TR/2009/CR-xmlschema11-1-20090430>
- [63] O. Nierstrasz, M. Kobel, T. Girba, M. Lanza, H. Bunke, *Example-Driven Reconstruction of Software Models*, in: R. Krikhaar, C. Verhoef, G. D. Lucca (Eds.), *Proceedings of the 11th European Conference on Software Maintenance and Reengineering (CSMR 2007)*, IEEE Computer Society, 2007, pp. 275–284.

- [64] M. Herrmannsdörfer, D. Ratiu, M. Kögel, Metamodel Usage Analysis for Identifying Metamodel Improvements, in: B. A. Malloy, S. Staab, M. G. J. van den Brand (Eds.), Post-proceedings of the Third International Conference on Software Language Engineering (SLE 2010), Vol. 6563 of LNCS, Springer, 2011, pp. 62–81.
- [65] M. Lungu, M. Lanza, O. Nierstrasz, Evolutionary and Collaborative Software Architecture Recovery with SoftwareNaut, *Science of Computer Programming; WASDeTT 2010 Special Issue (EST 4)* 79 (2014) 204–223. doi:<http://dx.doi.org/10.1016/j.scico.2012.04.007>.
- [66] M. Črepinšek, M. Mernik, F. Javed, B. R. Bryant, A. Sprague, Extracting Grammar from Programs: Evolutionary Approach, *SIGPLAN Notices* 40 (4) (2005) 39–46.
- [67] T. Parr, et al., ANTLR Grammar List, <http://www.antlr.org/grammar/list> (2003).
- [68] J. R. Cordy, et al., TXL World: Grammars, <http://www.txl.ca/nresources.html> (2003).
- [69] R. Bosworth, Syntax specification (section 18) is inconsistent with other sections, http://bugs.sun.com/bugdatabase/view_bug.do?bug_id=6442525 (2006).
- [70] R. Lämmel, G. Wachsmuth, Transformation of SDF Syntax Definitions in the ASF+SDF Meta-Environment, *Electronic Notes in Theoretical Computer Science* 44 (2) (2001) 9–33.
- [71] R. Lämmel, Grammar Adaptation, in: *Proceedings of the International Symposium of Formal Methods Europe on Formal Methods for Increasing Software Productivity*, Vol. 2021 of LNCS, Springer, 2001, pp. 550–570.
- [72] J. Kort, R. Lämmel, C. Verhoef, The Grammar Deployment Kit. System Demonstration, *Electronic Notes in Theoretical Computer Science* 65 (3) (2002) 117–123, Second Workshop on Language Descriptions, Tools and Applications (LDTA 2002).
- [73] J. Kort, *Grammar Deployment Kit: Reference Manual*, Universiteit Amsterdam (May 2003).
URL <http://gdk.sourceforge.net/gdkref.pdf>
- [74] G. Wachsmuth, Metamodel Adaptation and Model Co-adaptation, in: E. Ernst (Ed.), *Proceedings of the 21st European Conference on Object-Oriented Programming (ECOOP 2007)*, Vol. 4609 of LNCS, Springer, 2007, pp. 600–624.
- [75] A. Cicchetti, F. Ciccozzi, T. Leveque, A. Pierantonio, On the concurrent versioning of metamodels and models: challenges and possible solutions, in: *Proceedings of the Second International Workshop on Model Comparison in Practice (IWMCP 2011)*, ACM, New York, NY, USA, 2011, pp. 16–25.
- [76] M. Herrmannsdörfer, S. Vermolen, G. Wachsmuth, An Extensive Catalog of Operators for the Coupled Evolution of Metamodels and Models, in: B. A. Malloy, S. Staab, M. G. J. van den Brand (Eds.), Post-proceedings of the Third International Conference on Software Language Engineering (SLE 2010), Vol. 6563 of LNCS, Springer, 2011, pp. 163–182.
- [77] R. Lämmel, The Amsterdam Toolkit for Language Archaeology, *Electronic Notes in Theoretical Computer Science* 137 (3) (2005) 43–55.
- [78] M. van der Graaf, A Specification of Box to HTML in ASF+SDF, Master’s thesis, Universiteit van Amsterdam, The Netherlands (Aug. 1997).
- [79] A. Sellink, C. Verhoef, Generation of Software Renovation Factories from Compilers, in: *Proceedings of 15th International Conference on Software Maintenance (ICSM 1999)*, 1999, pp. 245–255.

- [80] ISO/IEC 14977:1996(E), Information Technology—Syntactic Metalanguage—Extended BNF, available at <http://www.cl.cam.ac.uk/~mgk25/iso-14977.pdf>.
- [81] A. Kelemenová, Grammatical Levels of the Position Restricted Grammars, in: Proceedings on Mathematical Foundations of Computer Science, Springer, 1981, pp. 347–359.
- [82] F. Coste, J. Nicolas, Regular Inference as a Graph Coloring Problem, in: Workshop on Grammar Inference, Automata Induction, and Language Acquisition, 1997, pp. 9–17. doi:10.1.1.34.4048.
- [83] A. Okhotin, *Boolean Grammars*, Information and Computation 194 (1) (2004) 19–48. doi:10.1016/j.ic.2004.03.006. URL http://users.utu.fi/aleokh/papers/boolean_grammars_ic.pdf
- [84] N. Chomsky, On Certain Formal Properties of Grammars, Information and Control 2 (2) (1959) 137–167.
- [85] H. P. Barendregt, Introduction to Generalized Type Systems, Journal of Functional Programming 1 (2) (1991) 125–154.
- [86] S. C. Kleene, Representation of Events in Nerve Nets and Finite Automata, Automata Studies (1956) 3–42.
- [87] M. G. J. van den Brand, J. Scheerder, J. J. Vinju, E. Visser, Disambiguation Filters for Scannerless Generalized LR Parsers, in: N. Horspool (Ed.), Compiler Construction 2002 (CC 2002), 2002, pp. 143–158.
- [88] P. Purdom, A Sentence Generator for Testing Parsers, BIT 12 (3) (1972) 366–375.
- [89] J. Levine, flex & bison, O’Reilly Media, 2009.
- [90] T. Copeland, Generating Parsers with JavaCC: An Easy to Use Guide for Programmers, Centennial Books, 2007.
- [91] A. M. Sloane, L. C. Kats, E. Visser, A Pure Embedding of Attribute Grammars, Science of Computer Programming (2011) 18 pages. In press. Available online since 29 November 2011.
- [92] P. Klint, J. J. Vinju, A. Lankamp, A. Izmaylova, D. Landman, T. van der Storm, A. van der Ploeg, M. Hills, B. Lissner, E. Balland, A. Shahi, A. H. Bagge, M. Steindorfer, W. Venema, J. van den Bos, V. Zaytsev, J. Timmer, B. Basten, J. Peeters, C. Berghuizen, D. Meers, P. I. Valdera, R. van Rozen, J. Stoel, M. Bierlee, T. Hooper, K. van der Vlist, J. van der Woning, P. Hijma, Rascal Language Library¹², <http://github.com/cwi-swat/rascal>, the `src/org/rascalmpl/library/lang` directory (2010–2014).
- [93] J. Vinju, M. van den Brand, T. van der Storm, M. Bravenboer, SDF Library¹³, <http://github.com/cwi-swat/meta-environment>, the `sdf-library` directory (1996–2012).
- [94] S. C. Johnson, YACC—Yet Another Compiler Compiler, Computer Science Technical Report 32, AT&T Bell Laboratories, Murray Hill, New Jersey, 1975.
- [95] F. Jouault, J. Bézivin, KM3: A DSL for Metamodel Specification, in: R. Gorrieri, H. Wehrheim (Eds.), Proceedings of the International Conference on Formal Methods for Open Object-Based Distributed Systems (FMOODS 2006), Vol. 4037 of LNCS, Springer, 2006, pp. 171–185.

¹²The authors are given according to the list of contributors at <http://github.com/cwi-swat/rascal/graphs/contributors>.

¹³The authors are given according to the AUTHORS file provided in the distribution.

- [96] Microsoft Corporation, The Microsoft code name "M" Modeling Language Specification, <http://msdn.microsoft.com/en-us/library/dd285271.aspx>, obsolete.
- [97] M. Fogus, C. Houser, The Joy of Clojure, Manning, 2010.
- [98] Object Management Group, [Semantics Of Business Vocabulary And Business Rules](http://www.omg.org/spec/SBVR/1.0), v1.0 Edition (2008).
URL <http://www.omg.org/spec/SBVR/1.0>
- [99] Object Management Group, [Unified Modeling Language](http://schema.omg.org/spec/UML/2.1.1), 2nd Edition (2007).
URL <http://schema.omg.org/spec/UML/2.1.1>
- [100] U. Brandes, M. Eiglsperger, I. Herman, M. Himsolt, M. S. Marshall, Graphml progress report: Structural layer proposal, in: P. Mutzel, M. Jnger, S. Leipert (Eds.), Revised Papers of the 9th International Symposium on Graph Drawing (GD 2001), Vol. 2265 of LNCS, Springer, 2001, pp. 501–512.
- [101] D. L. McGuinness, F. van Harmelen, [OWL: Web Ontology Language. Overview](http://www.w3.org/TR/owl-features), W3C Recommendation. .
URL <http://www.w3.org/TR/owl-features>
- [102] Object Management Group, [Meta-Object Facility \(MOF™\) Core Specification](http://www.omg.org/spec/MOF/2.0), 2nd Edition (Jan. 2006).
URL <http://www.omg.org/spec/MOF/2.0>
- [103] É. M. Gagnon, L. J. Hendren, SableCC, an Object-Oriented Compiler Framework, in: TOOLS 1998: 26th International Conference on Technology of Object-Oriented Languages and Systems, 3-7 August 1998, Santa Barbara, CA, USA, IEEE Computer Society, 1998, pp. 140–154. doi:10.1109/TOOLS.1998.711009.
- [104] Eclipse, Eclipse Modeling Framework Project (EMF 2.4), <http://www.eclipse.org/modeling/emf/> (2008).
- [105] F. Jouault, J. Bézivin, I. Kurtev, TCS: a DSL for the Specification of Textual Concrete Syntaxes in Model Engineering, in: S. Jarzabek, D. C. Schmidt, T. L. Veldhuizen (Eds.), Proceedings of the Fifth International Conference on Generative Programming and Component Engineering (GPCE), ACM, 2006, pp. 249–254.
- [106] F. Heidenreich, J. Johannes, S. Karol, M. Seifert, C. Wende, Derivation and Refinement of Textual Syntax for Models, in: R. F. Paige, A. Hartman, A. Rensink (Eds.), Model Driven Architecture — Foundations and Applications, Vol. 5562 of LNCS, Springer, 2009, pp. 114–129. doi:10.1007/978-3-642-02674-4_9.
- [107] J. Clark, M. Murata, RELAX NG Specification, OASIS Committee Specification Available at <http://relaxng.org/spec-20011203.html>, also being standardised by ISO as multiple parts of ISO/IEC 19757.
- [108] V. Zaytsev, Guided Grammar Convergence, JOT. In print (2014).
- [109] T. Parr, ANTLR v3 — ANother Tool for Language Recognition, <http://antlr3.org> (2008).
- [110] R. B. France, J. M. Bieman, S. P. Mandalaparty, B. H. C. Cheng, A. C. Jensen, Repository for model driven development (remodd), in: M. Glinz, G. C. Murphy, M. Pezzè (Eds.), 34th International Conference on Software Engineering (ICSE), IEEE, 2012, pp. 1471–1472.
- [111] M. Murata, RELAX NG Schemas, <http://relaxng.org/#schemas> (2013).

- [112] K. Czarnecki, J. Foster, Z. Hu, R. Lämmel, A. Schürr, J. Terwilliger, Bidirectional Transformations: A Cross-Discipline Perspective, in: R. Paige (Ed.), Theory and Practice of Model Transformations, Vol. 5563 of LNCS, Springer, 2009, pp. 260–283.
- [113] J. N. Foster, M. B. Greenwald, J. T. Moore, B. C. Pierce, A. Schmitt, Combinators for Bidirectional Tree Transformations: A Linguistic Approach to the View-Update Problem, ACM Transactions on Programming Languages and Systems (TOPLAS) 29.
- [114] I. D. Baxter, C. Pidgeon, M. Mehlich, Dms: Program transformations for practical scalable software evolution, in: Proceedings of the 26th International Conference on Software Engineering, ICSE '04, IEEE Computer Society, Washington, DC, USA, 2004, pp. 625–634.
- [115] S. Klusener, V. Zaytsev, [Language Standardization Needs Grammarware](#), JTC1/SC22 Document N3977, ISO/IEC (2005).
URL <http://www.open-std.org/jtc1/sc22/open/n3977.pdf>
- [116] M. H. Halstead, Elements of Software Science, Elsevier Science Inc., New York, NY, USA, 1977.