

Efficient Pattern Matching in Python

Manuel Krebber, Henrik Barthels, Paolo Bientinesi

Aachen Institute for Advanced Study in Computational Engineering Science

High-Performance and Automatic Computing Group

RWTH Aachen University

Abstract—Pattern matching is a powerful tool for symbolic computations. Applications include term rewriting systems, as well as the manipulation of symbolic expressions, abstract syntax trees, and XML and JSON data. It also allows for an intuitive description of algorithms in the form of rewrite rules. We present the open source Python module MatchPy, which offers functionality and expressiveness similar to the pattern matching in Mathematica. In particular, it includes syntactic pattern matching, as well as matching for commutative and/or associative functions, sequence variables, and matching with constraints. MatchPy uses new and improved algorithms to efficiently find matches for large pattern sets by exploiting similarities between patterns. The performance of MatchPy is investigated on several real-world problems.

1. Introduction

Pattern matching is a powerful tool which is part of many functional programming languages as well as computer algebra systems such as Mathematica. It is useful for many applications including symbolic computation, term simplification, term rewriting systems, automated theorem proving, code generation, and model checking. In this paper, we present the Python pattern matching library MatchPy [1] and its underlying algorithms.

The applications of pattern matching are similar to those of regular expressions, but on symbolic tree structures instead of strings. Moreover, unlike regular expressions, pattern matching can handle nested expressions up to arbitrary depth. The goal of pattern matching is to find a match substitution given a subject term and a pattern, which is a term with variables [2]. The substitution maps variables in the pattern to replacement terms. A match is a substitution that, when applied to the pattern, yields the original subject. As an example consider the subject $f(a)$ and the pattern $f(x)$ where f is a function symbol, a is a constant and x is a variable. Then the substitution $\sigma = \{x \mapsto a\}$ is a match because $\sigma(f(x)) = f(a)$. The most basic form of pattern matching where functions are neither associative nor commutative is called syntactic matching.

In addition to regular variables which can match a single term, MatchPy supports sequence variables that can match a sequence of terms. For example, if x^+ is such a sequence variable, then $f(a, b)$ matches the pattern $f(x^+)$ with the

substitution $\{x \mapsto (a, b)\}$, whereas $f(x)$ does not match. MatchPy enables further control over what variables can match by supporting symbol variables, i.e. variables that only match a specific user-defined class of symbols instead of any term. Furthermore, variables can have a default value, as for instance in the pattern $x:1 \cdot y$. Here, $x:1$ denotes that if x does not match anything else, then it has a default value of 1. This pattern matches both a and $a \cdot b$. For a , the match is $\{x \mapsto 1, y \mapsto a\}$, whereas for the subject $a \cdot b$, it is $\{x \mapsto a, y \mapsto b\}$.

In practice, functions commonly have properties such as commutativity and associativity. These properties affect the pattern matching, e.g. $f(a, b)$ and $f(b, a)$ match iff f is commutative. To our knowledge, no existing work covers pattern matching with function symbols which are either commutative or associative but not both at the same time. However, several common functions possess such properties, e.g. matrix multiplication and arithmetic mean.

Among the existing systems, Mathematica [3] arguably offers the most powerful pattern matching. Patterns are used widely in Mathematica, e.g. in function definitions or for manipulating terms. Mathematica offers support for associativity, commutativity, sequence variables, optional variables, symbol variables, constraints, alternative patterns, and arbitrarily repeated patterns. MatchPy aims to replicate most of this functionality.

Pattern matching forms the basis of term rewriting systems (TRS), where it is necessary to determine which rewrite rules can be applied. Since a TRS maps input expressions to output expressions, they can be seen as algorithms or programs operation on symbolic expressions. As a sufficiently powerful TRS can be Turing complete, they can be used as a programming language [4].

In many applications, a fixed pattern set is matched repeatedly against different subjects. The simultaneous matching of multiple patterns is called many-to-one matching, as opposed to one-to-one matching. Many-to-one matching can provide a significant speedup over one-to-one matching by exploiting similarities between patterns. While there has been research on many-to-one matching for some of MatchPy's features, previously no many-to-one algorithms existed for pattern matching as expressive as the one offered by MatchPy (and Mathematica). We generalize existing algorithms to support associativity, commutativity, sequence/optional/symbol variables, and constraints.

2. Related Work

While some basic forms of pattern matching are common in functional programming languages, support for pattern matching in popular imperative programming languages is fairly rare. Most of the existing pattern matching libraries for Python only support syntactic patterns. Among them are several that allow for functional-style syntactic pattern matching for native data structures, e.g. MacroPy [5], patterns [6] or PyPatt [7]. The pattern matching in SymPy [8] can work with associative and commutative functions, but it is limited to finding a single match—to enable further processing, it is often useful to find all possible matches for a pattern—and there are no algorithms for many-to-one matching. Furthermore, SymPy does not support sequence variables and is limited to a predefined set of mathematical operations.

Previous research on many-to-one matching either focused on syntactic matching (functions with fixed arity and no special properties) [9], [10], [11] or AC matching (variadic, associative and commutative functions) [4], [12], [13], [14], [15], [16]. This research did not consider function symbols which only have some of those properties, nor sequence variables. The work of Kutsia does include sequence variables, but the focus is on theoretical aspects of one-to-one matching [17], [18].

Even though Mathematica has powerful pattern matching features, it has some drawbacks. The possibilities to access Mathematica’s features from other programming languages is rather limited; writing large programs in Mathematica can be cumbersome and slow; furthermore, Mathematica is a commercial and proprietary product. Instead, it is desirable to have a free and open source pattern matching implementation that also enables other researchers to use and extend it.

Mathics [19] is an open source computer algebra system written in Python that aims to replicate the syntax and functionality of Mathematica. Unfortunately, currently there is little active development in this project.

3. Preliminaries

The notation and definitions used are based on what is used in term rewriting systems literature [2], [20], [21]. Pattern matching works on terms that consist of function symbols \mathcal{F} and variables \mathcal{X} . The function symbol set is composed of function symbols with different arities, i.e. $\mathcal{F} = \bigcup_{n \geq 0} \mathcal{F}_n$ where \mathcal{F}_n contains symbols with arity n . The function symbols can either have a fixed arity (i.e. they only occur in one \mathcal{F}_n) or be *variadic* (i.e. occur in all \mathcal{F}_n for $n \geq n_0$ and some fixed n_0). Specifically, \mathcal{F}_0 contains all constant symbols. The set of all terms $\mathcal{T}(\mathcal{F}, \mathcal{X})$ is the smallest set such that

- $\mathcal{X} \subseteq \mathcal{T}(\mathcal{F}, \mathcal{X})$ and
- for all $n \geq 0$, all $f \in \mathcal{F}_n$, and all $t_1, \dots, t_n \in \mathcal{T}(\mathcal{F}, \mathcal{X})$ we have $f(t_1, \dots, t_n) \in \mathcal{T}(\mathcal{F}, \mathcal{X})$.

Whenever the actual symbols and variables are not important, we use \mathcal{T} instead of $\mathcal{T}(\mathcal{F}, \mathcal{X})$. Terms in $\mathcal{X} \cup \mathcal{F}_0$ are called *atomic terms*, all others are *compound terms*. The set $\mathcal{T}(\mathcal{F}, \emptyset) \subseteq \mathcal{T}(\mathcal{F}, \mathcal{X})$ is the set of *ground terms*. We usually use $\mathcal{G}(\mathcal{T})$ or simply \mathcal{G} to denote it. Since ground terms are variable-free, they are also called *constant terms*. The set of variables occurring in a term t is denoted by $\mathcal{V}ar(t)$. A pattern t is called *linear* if every variable in $\mathcal{V}ar(t)$ occurs at most once.

In the following, we usually use f, g , and h as function symbols and a, b , and c as constant symbols. Common mathematical functions such as $+$ and \times are usually written in infix notation, and we omit the parenthesis if unnecessary, i.e. we write $a + b$ instead of $+(a, b)$.

A substitution is a function $\sigma: \mathcal{X} \rightarrow \mathcal{T}$. We extend the substitution to a function $\sigma: \mathcal{T} \rightarrow \mathcal{T}$ over all terms by using $\sigma(f(t_1, \dots, t_n)) = f(\sigma(t_1), \dots, \sigma(t_n))$ for every $n \geq 0$, all $f \in \mathcal{F}_n$ and all $t_1, \dots, t_n \in \mathcal{T}$. We often write substitutions as $\sigma = \{x_1 \mapsto t_1, \dots, x_n \mapsto t_n\}$ where $\{x_1, \dots, x_n\} = \text{Dom}(\sigma)$ are the variables *instantiated* by the substitution. The set of all substitutions is called $\text{Sub}(\mathcal{T}(\mathcal{F}, \mathcal{X}))$ or simply Sub .

A pattern term $t \in \mathcal{T}$ *matches* a subject term $s \in \mathcal{G}$, iff there exists a substitution σ such that $\sigma(t) = s$. Such a substitution is also called a *match*. The objective of pattern matching is to find such match substitutions.

We extend this basic definition of syntactic pattern matching to support associativity, commutativity, and one-identity. If a function symbol f is associative, i.e. $f(t_1, f(t_2, t_3)) = f(f(t_1, t_2), t_3)$ holds for all $t_1, t_2, t_3 \in \mathcal{T}$, we always use the equivalent flattened form $f(t_1, t_2, t_3)$. If a function symbol f is commutative, i.e. $f(t_1, t_2) = f(t_2, t_1)$ holds for all $t_1, t_2 \in \mathcal{T}$, we always use the equivalent form $f(t_1, t_2)$ with lexicographically sorted arguments (assuming $t_1 \leq t_2$). Finally, if a function symbol f is equivalent to the identity function with a single argument, i.e. $f(t_1) = t_1$ holds for all $t_1 \in \mathcal{T}$, we always write t_1 instead. We use $\mathcal{F}_a, \mathcal{F}_c$, and \mathcal{F}_I to denote the set of all associative functions, the set of all commutative functions, and the set of identity functions, respectively. Finally, we use $=_{ACI}$ to denote equivalence of terms under associativity, commutativity, and one-identity and consider σ to be a match if there exist $t' \in \mathcal{T}$ and $s' \in \mathcal{G}$ such that $\sigma(t') = s'$ and $t =_{ACI} t'$ and $s =_{ACI} s'$.

In contrast to regular variables, sequence variables match a sequence of terms instead of a single term. Sequence variables are denoted analogously to the notation used in regular expressions: x^* matches any number of terms including the empty sequence, x^+ requires at least one term to match.¹ The definition of a substitution is extended to allow sequences of terms as well: $\sigma: \mathcal{X} \rightarrow \mathcal{T}^*$. When applying a substitution, the replacement of a sequence variable is integrated into the sequence of function symbol arguments: $\sigma(f(a, x^*, b)) = f(a, c, d, b)$ for $\sigma = \{x^* \mapsto (c, d)\}$. Patterns with sequence variables may yield multiple matches.

1. Corresponding to `x_____` and `x__` in Mathematica, respectively.

Matching sequence variables is equivalent to matching regular variables within associative functions. As an example, consider $f(a, x^+)$ and $f_A(a, x)$ where f_A is associative. For the subjects $f(a, b, c)$ and $f_A(a, b, c)$, substitutions are $\{x \mapsto (b, c)\}$ and $\{x \mapsto f_A(b, c)\}$, respectively.

MatchPy also supports constraint predicates to further filter matches. A constraint is a function $\varphi: \text{Sub} \mapsto \{0, 1\}$. We say a match σ is *valid* for a pattern with constraint φ iff $\varphi(\sigma) = 1$. Symbol variables can be used to match only specific symbol classes, i.e. a specific subset $C \subseteq \mathcal{F}_0$. Furthermore, matching a variable can be made optional by providing a default value for when it does not match. This is useful to reduce the number of patterns. As an example, $a \cdot x + b$ can cover a linear equation with a defaulting to 1 and b to 0. Then the subject x can be matched with $\{a \mapsto 1, b \mapsto 0, x \mapsto x\}$.

4. Methods/Optimizations

4.1. One-to-one

4.1.1. Associativity/Sequence variables. Associativity enables arbitrary grouping of arguments for matching: For example, $1 + a + b$ matches $1 + x$ with $\{x \mapsto a + b\}$ because we can group the arguments as $1 + (a + b)$. When regular variables are arguments of an associative function, they behave like sequence variables. Both can result in multiple distinct matches for a single pattern. To enumerate all matches, we need to backtrack while matching for every choice that we make for such variables. To accomplish this, we use Python generators to return a value and be able to resume at the same state when backtracking. Associative matching is NP-complete [22].

4.1.2. Commutativity. Matching commutative terms is difficult because matches need to be found independently of the argument order. Commutative matching has been shown to be NP-complete [22]. It is possible to find all matches by matching all permutations of the subjects arguments against all permutations of the pattern arguments. If n is the number of subject arguments, and m is the number of pattern arguments, this naive approach leads to testing a total of $n!m!$ combinations, and it is likely that most of them either do not match or yield redundant matches. To address this challenge, we interpret the arguments as a multiset, i.e. an orderless collection that allows repetition of elements.

Additionally, we use the following sequence of steps for matching the subterms of a commutative term:

- 1) Constant arguments
- 2) Matched variables, i.e. variables that already have a value assigned in the current substitution
- 3) Non-variable arguments
- 4) Go back to 2 if there are new matched variables.
- 5) Regular variables
- 6) Sequence variables

The idea behind this sequence of steps is to reduce the size of the search space as quickly as possible, starting

with the simplest steps. As an example, matching constant arguments is equivalent to testing multiset membership, which can be done very efficiently. Each step reduces the search space for successive steps. If one step finds no match, the remaining steps do not have to be performed. Note that steps 3, 5 and 6 can yield multiple matches and backtracking is employed to check every combination. Since step 6 is the most involved, it is described in more detail in the next section.

4.1.3. Sequence Variables in Commutative Functions.

The distribution of n subjects subterms onto m sequence variables within a commutative function symbol can yield up to m^n distinct solutions. Enumerating all of the solutions is accomplished by generating and solving several linear Diophantine equations. As an example, lets assume we want to match $f(a, b, b, b)$ with $f(x^*, y^+, y^+)$ where f is commutative. This means that the possible distributions are given by the non-negative integer solutions of these equations:

$$\begin{aligned} 1 &= x_a + 2y_a \\ 3 &= x_b + 2y_b \end{aligned}$$

x_a determines how many times a is included in the substitution for x . Because y requires at least one term, we have the additional constraint $y_a + y_b \geq 1$. The only possible solution $x_a = x_b = y_b = 1 \wedge y_a = 0$ corresponds to the match substitution $\{x \mapsto (a, b), y \mapsto (b)\}$.

Extensive research has been done on solving linear Diophantine equations and linear Diophantine equation systems [23], [24], [25], [26], [27]. In our case, the equations are actually independent except for the additional constraints for plus variables. Furthermore, only the non-negative solutions are of interest in this case, and they can be found more easily. We use an adaptation of the algorithm used in SymPy which recursively reduces a linear Diophantine equation to a set of equations of the form $ax + by = d$, which can be solved efficiently with the Extended Euclidian algorithm [28]. The solutions are then combined into a solution for the original equation.

Since the coefficients in the equations correspond to the multiplicity of sequence variables, they are likely very small. Similarly, the number of variables in the equations is usually small as they map to sequence variables. The constant is the multiplicity of a subject term and hence also usually small. Overall, the number of distinct equations that are solved is small. Thus, by caching solutions, the impact of matching sequence variables on the overall run time can be reduced.

4.2. Many-to-one

MatchPy includes several many-to-one matching approaches: A deterministic discrimination net that only works for syntactic patterns, a generic many-to-one matcher and a code generator based on the many-to-one matcher. A discrimination net is a data structure similar to a decision tree or a finite automaton [9], [10], [11]. Generally, the many-to-one matchers enable matching multiple patterns against

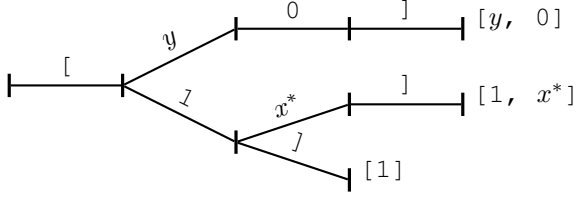


Figure 1. Example Discrimination Net.

a single subject much faster than matching each pattern individually using the one-to-one matching. The generic matcher utilizes a generalized form of non-deterministic discrimination nets that support sequence variables and associative function symbols. Furthermore, as described in the next section, it can also match commutative terms.

In Figure 1, an example for a non-deterministic discrimination net is shown. It contains three patterns that match Python lists: One matches the list that consists of a single 1, the second one matches a list with exactly two elements where the last element is 0, and the third pattern matches any list where the first element is 1. Note that these patterns can also match nested lists, e.g. the second pattern would also match $[[2, 1], 0]$.

Matching starts at the root and proceeds along the transitions. Simultaneously, the subject is traversed in preorder and each symbol is checked against the transitions. Only transitions matching the current subterm can be used. Once a final state is reached, its label gives a list of matching patterns. For non-deterministic discrimination nets, all possibilities need to be explored via backtracking. The discrimination net allows to reduce the matching costs, because common parts of different patterns only need to be matched once. For non-matching transitions, their whole subtree is pruned and all the patterns are excluded at once, further reducing the match cost.

In Figure 1, for the subject $[1, 0]$, there are two paths and therefore two matching patterns: $[y, 0]$ matches with $\{y \mapsto 1\}$ and $[1, x^*]$ matches with $\{x \mapsto 0\}$. Both the y -transition and the 1 -transition can be used in the second state to match a 1.

Compared to existing discrimination net variants, we added transitions for the end of a compound term to support variadic functions. Furthermore, we added support for both associative function symbols and sequence variables. Finally, our discrimination net supports transitions restricted to symbol classes, i.e. symbol variables. We decided to use a non-deterministic discrimination net instead of a deterministic one, since the number of states of the latter grows exponentially with the number of patterns. While the deterministic discrimination net also has support for sequence variables, in practice the net became too large to use with just a dozen patterns.

4.2.1. Commutative Many-to-one Matching. Many-to-one matching for commutative terms is more involved. We use a nested commutative matcher which in turn uses another generic many-to-one matcher to match the subterms.

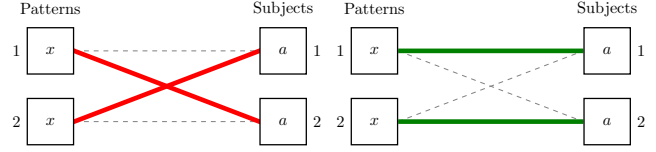


Figure 2. Example for Order in Bipartite Graph.

Our approach is similar to the one used by Bachmair and Kirchner in their respective works [4], [15]. We match all the subterms of the commutative function in the subject with a many-to-one matcher constructed from the subpatterns of the commutative function in the pattern (except for sequence variables, which are handled separately). The resulting matches form a bipartite graph, where one set of nodes consists of the subject subterms and the other contains all the pattern subterms. Two nodes are connected by an edge iff the pattern matches the subject. Such an edge is also labeled with the match substitution(s). Finding an overall match is then accomplished by finding a maximum matching in this graph. However, for the matching to be valid, all the substitutions on its edges must be compatible, i.e. they cannot have contradicting replacements for the same variable. We use the Hopcroft-Karp algorithm [29] to find an initial maximum matching. However, since we are also interested in all matches and the initial matching might have incompatible substitutions, we use the algorithm described by Uno, Fukuda and Matsui [30], [31] to enumerate all maximum matchings.

To avoid yielding redundant matches, we extended the bipartite graph by introducing a total order over its two node sets. This enables determining whether the edges of a matching maintain the order induced by the subjects or whether some of the edges "cross". Formally, for all edge pairs $(p, s), (p', s') \in M$ we require $(s \equiv s' \wedge p > p') \implies s > s'$ to hold where M is the matching, s, s' are subjects, and p, p' are patterns. An example of this is given in Figure 2. The order of the nodes is indicated by the numbers next to them. The only two maximum matchings for this particular match graph are displayed. In the left matching, the edges with the same subject cross and hence this matching is discarded. The other matching is used because it maintains the order. This ensures that only unique matches are yielded. Once a matching for the subpatterns is obtained, the remaining subject arguments are distributed to sequence variables in the same way as for one-to-one matching.

4.3. Code Generation

Both the one-to-one and many-to-one matching algorithm use some in-memory representation of the patterns at run time. This allows for the dynamic construction of the pattern set at run time. For fixed pattern sets, however, many-to-one matching can be sped up even further by generating Python code, similar to how parser generators generate code that parses a given grammar. This is done by converting the

Table 1. LINEAR ALGEBRA OPERATIONS

Operation	Symbol	Arity	Properties
Multiplication	\times	variadic	associative
Addition	$+$	variadic	associative, commutative
Transposition	T	unary	
Inversion	-1	unary	
Inversion and Transposition	$-T$	unary	

discrimination net structure of the many-to-one matcher to code. While this code generation is expensive and yields large code files, the resulting code offers a significant speedup over regular many-to-one matching.

5. Experimental Results

We perform multiple experiments to evaluate the performance of different pattern matching strategies. The applications comprise linear algebra expressions, Python source code transformation, converting logic formulas to algebraic normal form (ANF), and symbolic integration. The evaluated pattern matching methods are one-to-one matching, many-to-one matching, code generation, and parallelized one-to-one matching. In every experiment, a single subject is matched against a fixed pattern set. Since all methods can be parallelized over multiple subjects by using standard methods (e.g. `multiprocessing` or `dask`), the experiments do not cover this. We do not consider the singular setup time for many-to-one matching or code generation here, because previous experiments have shown that the setup time is easily amortized [32], [33].

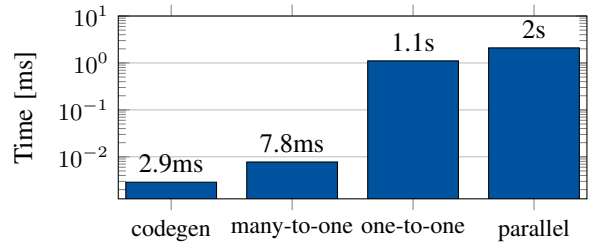
All times are measured on a server machine with 2 Intel Xeon E5-2670 v2 2.5GHz CPUs with 10 cores each and 64GB of RAM running Cent OS 6.8 and Python 3.6.2. The parallelized matching uses all 20 cores using the `multiprocessing` module. Because of the global interpreter lock (GIL) that prevents multithreading within Python code itself, this uses process based parallelism. We use a pool of processes and a queue for every subject-pattern pair. Since the pattern set is fixed, the patterns are available directly in the code. The subjects are considered to be dynamic and dependent on the use of the application. Hence, they are serialized to be sent to the worker processes.

5.1. Linear Algebra

BLAS is a collection of optimized routines that can compute specific linear algebra operations efficiently [34], [35], [36]. As an example, assume we want to match all subexpressions of a linear algebra expression which can be computed by the `?TRMM` BLAS routine. These have the form $\alpha \times op(A) \times B$ or $\alpha \times B \times op(A)$ where $op(A)$ is either the identity function or transposition, and A is a triangular matrix. For this example, we leave out all variants where $\alpha \neq 1$.

In order to model the linear algebra expressions, we use the operations shown in Table 1. In addition, we have special symbol subclasses for scalars, vectors, and matrices and use symbol variables to match them. Matrices also have a set of properties, e.g. they can be triangular, symmetric, square, etc. We use constraints to check the properties of the matched symbols. Finally, sequence variables are used to capture the remaining operands of multiplications or additions.

In total, we have about 200 patterns with about 60 being patterns for sums and 140 for products. A lot of these patterns only differ in terms of constraints, e.g. there are ten distinct patterns matching $A \times B$ with different constraints on the two matrices. As subjects, we use about 140 terms representing different linear algebra problems.

Figure 3. Total times for `LinAlg`

The results for those `LinAlg` problems are shown in Figure 3. Both many-to-one matching and the generated code are significantly faster than sequential and parallel one-to-one matching. Both are about two orders of magnitude faster than one-to-one matching. The generated code is about 2.7 times faster than the many-to-one matcher. The parallelized one-to-one matching is slower than the regular one-to-one matching. This is caused by the overhead of the process-based parallelism which requires the subjects to be serialized to be passed between processes.

5.2. Abstract Syntax Trees

Python includes a tool to convert code from Python 2 to Python 3. It is part of the standard library package `lib2to3` which has a collection of "fixers" that each convert one of the incompatible cases. To find matching parts of the code, those fixers use pattern matching on the abstract syntax tree (AST). Such an AST can be represented in the `MatchPy` data structures. We converted some of the patterns used by `lib2to3` both to demonstrate the generality of `MatchPy` and to evaluate the performance of many-to-one matching. Because the fixers are applied one after another and can modify the AST after each match, it would be difficult to use many-to-one matching for `lib2to3` in practice.

The following is an example of such a pattern:

```
power< 'isinstance'
  trailer< '(' arglist< any ','
    atom< '('
      args=testlist_gexp< any+ >
    ')' > > ')' > >
```

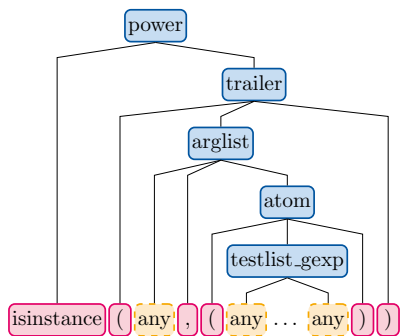


Figure 4. AST of the `isinstance` pattern.

It matches an `isinstance` expression with a tuple as second argument. Its tree structure is illustrated in Figure 4. The corresponding fixer cleans up duplications generated by previous fixers. For example, `instance(x, (int, long))` would be converted by another fixer into `instance(x, (int, int))` which in turn is then simplified to `instance(x, int)` by this fixer.

Out of the original 46 patterns, 36 could be converted to MatchPy patterns. Some patterns could not be converted because they contain features that MatchPy does not support yet. Those features include negated subpatterns (e.g. `not atom<'(' [any] ')>`) and subpatterns that allow an arbitrary number of repetitions (e.g. `any (',' any)+`).

Furthermore, some of the AST patterns contain alternative or optional subpatterns, e.g. `power<'input' args=trailer<'(' [any] ')>`. While these features are not directly supported by MatchPy either, they can be replicated by using multiple patterns. For those `lib2to3` patterns, all combinations of the alternatives were generated and added as individual patterns. This resulted in about 1200 MatchPy patterns to cover the original 36 `lib2to3` patterns. As shown by our experiments, many-to-one matching can mostly compensate for this substantial increase in the number of patterns. Note that the `lib2to3` patterns do not use some features of MatchPy, e.g. associativity or commutativity.

For the experiments, we use 613 examples from the unittests of `lib2to3` with a total of about 900 non-empty lines of Python code. The original `lib2to3` matcher using the set of 36 patterns is compared with MatchPy’s matchers using the 1200 patterns. A total of about 560 matches are found.

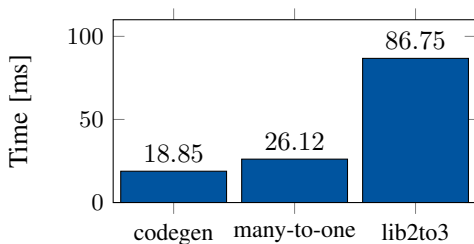


Figure 5. Total times for AST

The total times for matching all subjects are shown in Figure 7. One-to-one matching takes a total of more than five seconds to match and when parallelized the time goes up to 26s. Therefore we leave these times out of the plot as they are worse by several magnitudes and we want to focus on the comparison of the other three cases. The fastest method is the generated code with a speedup of about 4.7 over the `lib2to3` matcher. Many-to-one matching is also about 3.4 times faster than the original implementation. Overall, the many-to-one methods are about 200 times faster than the one-to-one matching. This is largely caused by the significant overlap of most of the patterns which were generated from the alternatives in the original patterns. Furthermore, less than 1% of the patterns match any given subject.

5.3. Rubi

Rubi is a rule-based symbolic integrator that uses more than 6000 replacement rules [37]. It significantly outperforms Mathematica and Maple in terms of the quality of the integration. The rules are available as Mathematica code. However, as part of a project to integrate those rules into SymPy, they are currently being ported to MatchPy.² We use a subset of these rules for our experiments. A very basic example of such a rule is $x^m \rightarrow \frac{1}{m+1}x^{m+1}$ with the additional constraints that x does not appear in m and $m+1 \neq 0$. We use about 100 of the problems from Section 1.2 of the Rubi integration test suite as subjects.

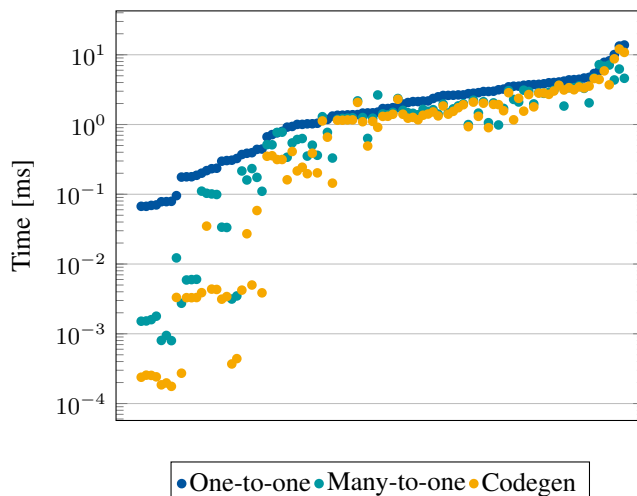


Figure 6. Times for Rubi

The resulting times are displayed in Figure 6 sorted by the time that the one-to-one matching takes. Here the speedups are less uniform than for the other applications. The longer the overall matching takes, the smaller is the difference between the different methods. The reason is

2. We would like to thank Arihant Parsoya, Abdullah Javed Nesar and Francesco Bonazzi, who have been working on integrating the Rubi rules into SymPy as part of the Google Summer of Code 2017.

that there is a significant overhead not related to pattern matching itself. There are some expected correlations: The more steps (i.e. rule applications) the integration takes, the longer the overall time is. The larger the term resulting from the integration is, the longer the integration takes on average. Overall, on average the many-to-one methods are about 50% faster than one-to-one matching. For the fastest third of subjects (with regards to one-to-one matching time) the speedup is significantly higher with about 5.6 for code generation and 2.2 for many-to-one matching on average.

5.4. Logic

We use the `Prop` example from [4] and add more examples to benchmark the performance of pattern matching on a small pattern set. The example uses a TRS with ten rules to convert boolean formulas to algebraic normal form (ANF), i.e. formulas that only contain \wedge and \oplus . Starting from either \top or \perp , we generate more complicated versions of these tautologies and contradictions by applying the simplifications backwards. For example, instead of t we can write the equivalent $\neg\neg t$ or $t \wedge \perp$ for any formula t . One characteristic of this problem is that formulas can become very large in the process of simplification.

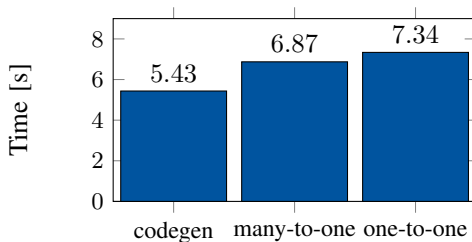


Figure 7. Average times for `Prop`

The resulting times for the logic benchmark are displayed in Figure 7. We only look at averages because the speedup is very uniform over all subjects. Because there is not a lot of overlap between the patterns and the number of patterns is small, the speedup of many-to-one matching and code generation is relatively small with 7% and 35%, respectively. In this case, it is also easy to write Python code by hand that performs the conversion to ANF. Our hand-coded solution is more than three orders of magnitude faster than the one-to-one matching approach. Therefore, we conclude that in cases with few patterns and little overlap of patterns, pattern matching is likely not a good choice in terms of performance. It might still be interesting in terms of productivity, though, because the TRS is significantly shorter and more readable than the hand-coded solution.

6. Productivity

Pattern matching and term rewriting systems provide a way to solve many problems in an intuitive and expressive way. In this section, we give an introduction of how MatchPy can be used.

MatchPy can be installed using `pip` and all necessary classes can be imported from the top-level module `matchpy`. Expressions in MatchPy consist of constant symbols, operations, and wildcards (i.e. variables). We use Mathematica’s notation³ for wildcards, i.e. we append underscores to wildcard names to distinguish them from symbols.

MatchPy can be used with native Python types such as `list` and `int`. The following is an example of how the subject `[0, 1]` can be matched against the pattern `[x_, 1]`. The expected match here is the replacement `0` for `x_`. Because some patterns can have multiple distinct matches, `match` is a generator that yields any match found. Hence, we use `next` because we only want to use the first (and in this case only) match of the pattern:

```
>>> x_ = Wildcard.dot('x')
>>> next(match([0, 1], Pattern([x_, 1])))
{'x': 0}
```

To illustrate how powerful pattern matching is, we show how a single rewrite rule with sequence variables and constraints can be used to implement a sorting algorithm for lists. The `CustomConstraint` class can be used to create a constraint that checks whether `a` is smaller than `b`:

```
a_lt_b = CustomConstraint(lambda a, b: a < b)
```

The lambda function gets called with the variable substitutions based on their name. The order of arguments is not important and it is possible to only use a subset of the variables in the pattern. With this constraint, we can define a replacement rule that basically describes bubble sort:

```
>>> pattern = Pattern([h_, b_, a_, t_], a_lt_b)
>>> rule = ReplacementRule(pattern,
...                          lambda a, b, h, t: [*h, a, b, *t])
```

This rule finds adjacent elements in the list which are in the wrong order and swaps them. The replacement lambda gets called with all matched variables as keyword arguments and needs to return the replacement. This replacement rule sorts a list when applied repeatedly with `replace_all`:

```
>>> replace_all([1, 4, 3, 2], [rule])
[1, 2, 3, 4]
```

Sequence variables can also be used to match subsequences that match a constraint. For example, we can use the `this` feature to find all subsequences of integers that sum up to 5. In the following example, we use anonymous wildcards which have no name and are hence not part of the match substitution:

```
>>> x_sums_to_5 = CustomConstraint(
...     lambda x: sum(x) == 5)
>>> pattern = Pattern([_, x_, _], x_sums_to_5)
>>> list(match([1, 2, 3, 1, 1, 2], pattern))
[{'x': (2, 3)}, {'x': (3, 1, 1)}]
```

More examples can be found in MatchPy’s documentation [38].

7. Conclusion & Future Work

7.1. Conclusions

We present MatchPy, a pattern matching library for Python with support for sequence variables and associa-

3. See <https://reference.wolfram.com/language/guide/Patterns.html>

tive/commutative functions. This library provides tools for both one-to-one and many-to-one matching. Because non-syntactic pattern matching is NP-hard, in the worst case the pattern matching times grows exponentially with the length of the pattern. Nonetheless, our experiments on real world examples indicate that many-to-one matching can give a significant speedup over one-to-one matching. However, the employed discrimination nets come with a one-time construction cost which needs to be amortized to benefit from their speedup. In previous experiments [32], [33], when considering typical numbers of subjects for the respective application, many-to-one matching was always faster overall. Therefore, many-to-one matching is likely to result in a compelling speedup in practice. The efficiency of using many-to-one matching heavily depends on the actual pattern set, i.e. the degree of similarity and overlap between the patterns.

We also use the discrimination nets to generate Python code to further improve the performance. As expected, the generated code is faster at matching, but comes with additional time cost for generating the code in the first place. For fixed pattern sets, code generation is the best option in terms of matching performance. However, the generated code can get large and complicated for large pattern sets to the point where it is not viable anymore.

In our experiments, the parallelization of one-to-one matching over the patterns is not beneficial. The parallel version is slower than the serial version in almost all cases. This is due to the GIL in Python that only allows process based parallelism within pure Python code. As a result, exchanging data between processes requires serializing, which introduces a significant overhead. Independent of the choice of the matching algorithm, matching can be parallelized over multiple subjects. For example, if we have multiple logic formulas which need to be converted to ANF, each one can be normalized independently and each replacement can run in parallel. There is still serialization overhead, but if each process takes enough time, there is parallel speedup to be gained. Each individual subject can then also be matched with many-to-one matching. However, how exactly this parallelization is best implemented highly depends on the application, so we do not investigate this approach in this paper.

7.2. Future Work

We plan on extending MatchPy with more powerful pattern matching features to make it useful for an even wider range of applications. The greatest challenge with additional features is likely to implement them for many-to-one matching. In the following, we discuss some possibilities for extending the library.

Additional Pattern Features. In the future, we plan to implement similar functionality to the *Repeated*, *Sequence*, and *Alternatives* functions from *Mathematica*. These provide another level of expressiveness which cannot be fully replicated with MatchPy’s current feature

set. Another useful feature are context variables as described by Kutsia [17]. They allow matching subterms at arbitrary depths which is especially useful for structures like XML. With context variables, MatchPy’s pattern matching would be as powerful as XPath [39] or CSS selectors [40] for such structures. Similarly, function variables which can match a function symbol would also be useful for those applications.

C Implementation. If pattern matching is a major part of an application, its runtime can significantly impact the overall speed. Reimplementing parts of MatchPy as a C module would likely result in a substantial speedup. Alternatively, adapting part of the code to Cython could be another option to increase the speed [41], [42]. This would also open up the possibility to benefit from parallelization by circumventing the GIL.

Code Generation. Furthermore, the code generation has room for improvement. For commutative matching, the many-to-one matching code is still mostly used directly with information about the patterns in the form of dictionaries. Instead, Python code could be generated for this as well to improve performance and remove the MatchPy dependency from the generated code. Another option is to generate Cython or C code.

References

- [1] M. Krebber, “MatchPy,” 2017. [Online]. Available: <https://github.com/HPAC/matchpy>
- [2] F. Baader and T. Nipkow, *Term Rewriting and All That*. New York, NY, USA: Cambridge University Press, 1998.
- [3] I. Wolfram Research, “Mathematica,” 2016, URL: <https://www.wolfram.com>.
- [4] H. Kirchner and P.-E. Moreau, “Promoting rewriting to a programming language: A compiler for non-deterministic rewrite programs in associative-commutative theories,” *International Journal of Foundations of Computer Science*, vol. 11, no. 2, pp. 207–251, Mar. 2001.
- [5] L. Haoyi, “Macropy,” <https://github.com/lihaoyi/macropy>, 2014. [Online]. Available: <https://github.com/lihaoyi/macropy>
- [6] A. Schepanovski, “patterns,” <https://github.com/Suor/patterns>, 2014. [Online]. Available: <https://github.com/Suor/patterns>
- [7] G. Jenks, “PyPatt,” <https://pypi.python.org/pypi/pypatt>, 2015. [Online]. Available: <https://pypi.python.org/pypi/pypatt>
- [8] A. Meurer, C. P. Smith, M. Paprocki, O. Čertík, S. B. Kirpichev, M. Rocklin, A. Kumar, S. Ivanov, J. K. Moore, S. Singh, T. Rathnayake, S. Vig, B. E. Granger, R. P. Muller, F. Bonazzi, H. Gupta, S. Vats, F. Johansson, F. Pedregosa, M. J. Curry, A. R. Terrel, v. Roučka, A. Saboo, I. Fernando, S. Kulal, R. Cimrman, and A. Scopatz, “SymPy: Symbolic computing in python,” *PeerJ Computer Science*, vol. 3, p. e103, Jan. 2017.
- [9] J. Christian, “Flatterms, discrimination nets, and fast term rewriting,” *Journal of Automated Reasoning*, vol. 10, no. 1, pp. 95–113, 1993.
- [10] A. Gräf, “Left-to-right tree pattern matching,” in *Proceedings of the 4th International Conference on Rewriting Techniques and Applications*, ser. RTA-91. New York, NY, USA: Springer-Verlag New York, Inc., 1991, pp. 323–334.
- [11] N. Nedjah, C. D. Walter, and S. E. Eldridge, *Optimal Left-to-right Pattern-matching Automata*, 1997, pp. 273–286.

- [12] E. Kounalis and D. Lugiez, "Compilation of pattern matching with associative-commutative functions," in *TAPSOFT '91*. Springer Nature, 1991, pp. 57–73.
- [13] L. Bachmair, T. Chen, and I. V. Ramakrishnan, *Associative-commutative Discrimination Nets*, 1993, pp. 61–74.
- [14] D. Lugiez and J. L. Moysset, "Tree automata help one to solve equational formulae in AC-theories," *Journal of Symbolic Computation*, vol. 18, no. 4, pp. 297–318, Oct. 1994.
- [15] L. Bachmair, T. Chen, I. V. Ramakrishnan, S. Anantharaman, and J. Chabin, "Experiments with associative-commutative discrimination nets," in *Proceedings of the 14th International Joint Conference on Artificial Intelligence - Volume 1*, ser. IJCAI'95. Morgan Kaufmann, 1995, pp. 348–354.
- [16] S. M. Eker, "Associative-commutative matching via bipartite graph matching," *The Computer Journal*, vol. 38, no. 5, pp. 381–399, May 1995.
- [17] T. Kutsia, "Context sequence matching for XML," *Electronic Notes in Theoretical Computer Science*, vol. 157, no. 2, pp. 47–65, May 2006.
- [18] —, "Solving equations with sequence variables and sequence functions," *Journal of Symbolic Computation*, vol. 42, no. 3, pp. 352–388, Mar. 2007.
- [19] J. Pöschko, "Mathics," <http://mathics.github.io/>, Oct. 2016. [Online]. Available: <http://mathics.github.io/>
- [20] N. Dershowitz and J.-P. Jouannaud, "Rewrite systems," in *Handbook of Theoretical Computer Science (vol. B)*, J. van Leeuwen, Ed. Cambridge, MA, USA: MIT Press, 1990, ch. Rewrite Systems, pp. 243–320.
- [21] J. W. Klop, R. de Vrijer, and M. Bezem, *Term Rewriting Systems*. New York, NY, USA: Cambridge University Press, 2001.
- [22] D. Benanav, D. Kapur, and P. Narendran, "Complexity of matching problems," *Journal of Symbolic Computation*, vol. 3, no. 1, pp. 203–216, Feb. 1987.
- [23] R. Weinstock, "Greatest common divisor of several integers and an associated linear diophantine equation," *The American Mathematical Monthly*, vol. 67, no. 7, p. 664, Aug. 1960.
- [24] J. Bond, "Calculating the general solution of a linear diophantine equation," *The American Mathematical Monthly*, vol. 74, no. 8, p. 955, Oct. 1967.
- [25] J. L. Lambert, "Finding a partial solution to a linear system of equations in positive integers," *Computers & Mathematics with Applications*, vol. 15, no. 3, pp. 209–212, 1988.
- [26] M. Clausen and A. Fortenbacher, "Efficient solution of linear diophantine equations," *Journal of Symbolic Computation*, vol. 8, no. 1–2, pp. 201–216, Jul. 1989.
- [27] K. Aardal, C. A. J. Hurkens, and A. K. Lenstra, "Solving a system of linear diophantine equations with lower and upper bounds on the variables," *Mathematics of Operations Research*, vol. 25, no. 3, pp. 427–442, Aug. 2000.
- [28] A. J. Menezes, S. A. Vanstone, and P. C. V. Oorschot, *Handbook of Applied Cryptography*, 1st ed. Boca Raton, FL, USA: CRC Press, Inc., 1996.
- [29] J. E. Hopcroft and R. M. Karp, "An $n^{5/2}$ algorithm for maximum matchings in bipartite graphs," *SIAM Journal on Computing*, vol. 2, no. 4, pp. 225–231, Dec. 1973.
- [30] K. Fukuda and T. Matsui, "Finding all the perfect matchings in bipartite graphs," *Applied Mathematics Letters*, vol. 7, no. 1, pp. 15–18, Jan. 1994.
- [31] T. Uno, "Algorithms for enumerating all perfect, maximum and maximal matchings in bipartite graphs," in *Algorithms and Computation*. Springer Nature, 1997, pp. 92–101.
- [32] Manuel Krebber, Henrik Barthels, and Paolo Bientinesi, "MatchPy: A Pattern Matching Library," in *Proceedings of the 15th Python in Science Conference*, Katy Huff, David Lippa, Dillon Niederhut, and M. Pacer, Eds., 2017, pp. 73–80.
- [33] M. Krebber, "Non-linear Associative-Commutative Many-to-One Pattern Matching with Sequence Variables," Master's thesis, RWTH Aachen University, 2017.
- [34] C. L. Lawson, R. J. Hanson, D. R. Kincaid, and F. T. Krogh, "Basic linear algebra subprograms for fortran usage," *ACM Transactions on Mathematical Software*, vol. 5, no. 3, pp. 308–323, Sep. 1979.
- [35] J. J. Dongarra, J. D. Croz, S. Hammarling, and R. J. Hanson, "An extended set of FORTRAN basic linear algebra subprograms," *ACM Transactions on Mathematical Software*, vol. 14, no. 1, pp. 1–17, Mar. 1988.
- [36] J. J. Dongarra, J. D. Croz, S. Hammarling, and I. S. Duff, "A set of level 3 basic linear algebra subprograms," *ACM Transactions on Mathematical Software*, vol. 16, no. 1, pp. 1–17, Mar. 1990.
- [37] A. D. Rich and D. J. Jeffrey, *A Knowledge Repository for Indefinite Integration Based on Transformation Rules*. Berlin, Heidelberg: Springer Berlin Heidelberg, 2009, pp. 480–485. [Online]. Available: https://doi.org/10.1007/978-3-642-02614-0_39
- [38] M. Krebber, "Matchpy documentation," <http://matchpy.readthedocs.io/>, 06 2017. [Online]. Available: <http://matchpy.readthedocs.io/>
- [39] J. Robie, M. Dyck, and J. Spiegel, "XML path language (XPath) 3.1," URL: <https://www.w3.org/TR/2017/REC-xpath-31-20170321/>, W3C, W3C Recommendation, Mar. 2017.
- [40] F. Rivoal, T. A. Jr., and E. Etamad, "CSS snapshot 2017," URL: <https://www.w3.org/TR/2017/NOTE-css-2017-20170131/>, W3C, W3C Note, Jan. 2017.
- [41] S. Behnel, R. W. Bradshaw, and D. S. Seljebotn, "Cython tutorial," in *Proceedings of the 8th Python in Science Conference*, G. Varoquaux, S. van der Walt, and J. Millman, Eds., Pasadena, CA USA, 2009, pp. 4–14.
- [42] I. M. Wilbers, H. P. Langtangen, and Å. Ødegård, "Using cython to speed up numerical python programs," in *Proceedings of Mekt'09*, 2009, pp. 495–512.