
Lark Documentation

Erez Shinan

Jan 10, 2024

OVERVIEW

1	Philosophy	1
1.1	Design Principles	1
1.2	Design Choices	1
2	Features	3
2.1	Main Features	3
2.2	Extra features	4
3	Parsers	5
3.1	Earley	5
3.2	LALR(1)	6
3.3	CYK Parser	7
4	JSON parser - Tutorial	9
4.1	Part 1 - The Grammar	9
4.2	Part 2 - Creating the Parser	11
4.3	Part 3 - Shaping the Tree	12
4.4	Part 4 - Evaluating the tree	13
4.5	Part 5 - Optimizing	14
4.6	Afterword	17
5	How To Use Lark - Guide	19
5.1	Work process	19
5.2	Getting started	19
5.3	Debug	19
5.4	Tools	21
6	How to develop Lark - Guide	23
6.1	Code Style	23
6.2	Unit Tests	23
7	Recipes	25
7.1	Use a transformer to parse integer tokens	25
7.2	Collect all comments with lexer_callbacks	25
7.3	CollapseAmbiguities	26
7.4	Keeping track of parents when visiting	27
7.5	Unwinding VisitError after a transformer/visitor exception	27
7.6	Adding a Progress Bar to Parsing with tqdm	28
8	Examples for Lark	29
8.1	Beginner Examples	29

8.2	Advanced Examples	37
9	Grammar Composition	39
10	Example Grammars	41
11	Standalone example	43
11.1	Advanced Examples	43
11.2	Grammar Composition	69
11.3	Example Grammars	72
11.4	Standalone example	72
12	Grammar Reference	75
12.1	Definitions	75
12.2	Terminals	76
12.3	Rules	78
12.4	Directives	79
13	Tree Construction Reference	81
13.1	Terminals	81
13.2	Shaping the tree	82
14	API Reference	85
14.1	Lark	85
14.2	Tree	89
14.3	Token	90
14.4	Transformer, Visitor & Interpreter	91
14.5	ForestVisitor, ForestTransformer, & TreeForestTransformer	91
14.6	UnexpectedInput	91
14.7	InteractiveParser	92
14.8	ast_utils	93
15	Transformers & Visitors	95
15.1	Visitor	95
15.2	Interpreter	96
15.3	Transformer	96
15.4	v_args	98
15.5	merge_transformers	99
15.6	Discard	100
15.7	VisitError	100
16	Working with the SPPF	101
16.1	SymbolNode	101
16.2	PackedNode	102
16.3	ForestVisitor	102
16.4	ForestTransformer	103
16.5	TreeForestTransformer	104
16.6	handles_ambiguity	105
17	Tools (Stand-alone, Nearley)	107
17.1	Stand-alone parser	107
17.2	Importing grammars from Nearley.js	107
18	Install Lark	109
19	Syntax Highlighting	111

20 Resources	113
Index	115

PHILOSOPHY

Parsers are innately complicated and confusing. They're difficult to understand, difficult to write, and difficult to use. Even experts on the subject can become baffled by the nuances of these complicated state-machines.

Lark's mission is to make the process of writing them as simple and abstract as possible, by following these design principles:

1.1 Design Principles

1. Readability matters
2. Keep the grammar clean and simple
3. Don't force the user to decide on things that the parser can figure out on its own
4. Usability is more important than performance
5. Performance is still very important
6. Follow the Zen of Python, whenever possible and applicable

In accordance with these principles, I arrived at the following design choices:

1.2 Design Choices

1.2.1 1. Separation of code and grammar

Grammars are the de-facto reference for your language, and for the structure of your parse-tree. For any non-trivial language, the conflation of code and grammar always turns out convoluted and difficult to read.

The grammars in Lark are EBNF-inspired, so they are especially easy to read & work with.

1.2.2 2. Always build a parse-tree (unless told not to)

Trees are always simpler to work with than state-machines.

1. Trees allow you to see the “state-machine” visually
2. Trees allow your computation to be aware of previous and future states
3. Trees allow you to process the parse in steps, instead of forcing you to do it all at once.

And anyway, every parse-tree can be replayed as a state-machine, so there is no loss of information.

See this answer in more detail [here](#).

To improve performance, you can skip building the tree for LALR(1), by providing Lark with a transformer (see the [JSON example](#)).

1.2.3 3. Earley is the default

The Earley algorithm can accept *any* context-free grammar you throw at it (i.e. any grammar you can write in EBNF, it can parse). That makes it extremely friendly to beginners, who are not aware of the strange and arbitrary restrictions that LALR(1) places on its grammars.

As the users grow to understand the structure of their grammar, the scope of their target language, and their performance requirements, they may choose to switch over to LALR(1) to gain a huge performance boost, possibly at the cost of some language features.

Both Earley and LALR(1) can use the same grammar, as long as all constraints are satisfied.

In short, “Premature optimization is the root of all evil.”

1.2.4 Other design features

- Automatically resolve terminal collisions whenever possible
- Automatically keep track of line & column numbers

FEATURES

2.1 Main Features

- Earley parser, capable of parsing any context-free grammar
 - Implements SPPF, for efficient parsing and storing of ambiguous grammars.
- LALR(1) parser, limited in power of expression, but very efficient in space and performance ($O(n)$).
 - Implements a parse-aware lexer that provides a better power of expression than traditional LALR implementations (such as ply).
- EBNF-inspired grammar, with extra features (See: [Grammar Reference](#))
- Builds a parse-tree (AST) automatically based on the grammar
- Stand-alone parser generator - create a small independent parser to embed in your project. ([read more](#))
- Flexible error handling by using an interactive parser interface (LALR only)
- Automatic line & column tracking (for both tokens and matched rules)
- Automatic terminal collision resolution
- Warns on regex collisions using the optional `interegular` library. ([read more](#))
- Grammar composition - Import terminals and rules from other grammars (see [example](#)).
- Standard library of terminals (strings, numbers, names, etc.)
- Unicode fully supported
- Extensive test suite
- Type annotations (MyPy support)
- Pure-Python implementation

[Read more about the parsers](#)

2.2 Extra features

- Support for external regex module ([see here](#))
- Import grammars from Nearley.js ([read more](#))
- CYK parser
- Visualize your parse trees as dot or png files ([see_example](#))
- Automatic reconstruction of input from parse-tree (see [example](#) and [another example](#))
- Use Lark grammars in [Julia](#) and [Javascript](#).

PARSERS

Lark implements the following parsing algorithms: Earley, LALR(1), and CYK

3.1 Earley

An [Earley Parser](#) is a chart parser capable of parsing any context-free grammar at $O(n^3)$, and $O(n^2)$ when the grammar is unambiguous. It can parse most LR grammars at $O(n)$. Most programming languages are LR, and can be parsed at a linear time.

Lark's Earley implementation runs on top of a skipping chart parser, which allows it to use regular expressions, instead of matching characters one-by-one. This is a huge improvement to Earley that is unique to Lark. This feature is used by default, but can also be requested explicitly using `lexer='dynamic'`.

It's possible to bypass the dynamic lexing, and use the regular Earley parser with a basic lexer, that tokenizes as an independent first step. Doing so will provide a speed benefit, but will tokenize without using Earley's ambiguity-resolution ability. So choose this only if you know why! Activate with `lexer='basic'`

SPPF & Ambiguity resolution

Lark implements the Shared Packed Parse Forest data-structure for the Earley parser, in order to reduce the space and computation required to handle ambiguous grammars.

You can read more about SPPF [here](#)

As a result, Lark can efficiently parse and store every ambiguity in the grammar, when using Earley.

Lark provides the following options to combat ambiguity:

1. Lark will choose the best derivation for you (default). Users can choose between different disambiguation strategies, and can prioritize (or demote) individual rules over others, using the rule-priority syntax.
2. Users may choose to receive the set of all possible parse-trees (using `ambiguity='explicit'`), and choose the best derivation themselves. While simple and flexible, it comes at the cost of space and performance, and so it isn't recommended for highly ambiguous grammars, or very long inputs.
3. As an advanced feature, users may use specialized visitors to iterate the SPPF themselves.

`lexer="dynamic_complete"`

Earley's "dynamic" lexer uses regular expressions in order to tokenize the text. It tries every possible combination of terminals, but it matches each terminal exactly once, returning the longest possible match.

That means, for example, that when `lexer="dynamic"` (which is the default), the terminal `/a+/,` when given the text `"aa"`, will return one result, `aa`, even though `a` would also be correct.

This behavior was chosen because it is much faster, and it is usually what you would expect.

Setting `lexer="dynamic_complete"` instructs the lexer to consider every possible regexp match. This ensures that the parser will consider and resolve every ambiguity, even inside the terminals themselves. This lexer provides the same capabilities as scannerless Earley, but with different performance tradeoffs.

Warning: This lexer can be much slower, especially for open-ended terminals such as `/.*/`

3.2 LALR(1)

LALR(1) is a very efficient, true-and-tested parsing algorithm. It's incredibly fast and requires very little memory. It can parse most programming languages (For example: Python and Java).

LALR(1) stands for:

- Left-to-right parsing order
- Rightmost derivation, bottom-up
- Lookahead of 1 token

Lark comes with an efficient implementation that outperforms every other parsing library for Python (including PLY)

Lark extends the traditional YACC-based architecture with a *contextual lexer*, which processes feedback from the parser, making the LALR(1) algorithm stronger than ever.

The contextual lexer communicates with the parser, and uses the parser's lookahead prediction to narrow its choice of terminals. So at each point, the lexer only matches the subgroup of terminals that are legal at that parser state, instead of all of the terminals. It's surprisingly effective at resolving common terminal collisions, and allows one to parse languages that LALR(1) was previously incapable of parsing.

(If you're familiar with YACC, you can think of it as automatic lexer-states)

This is an improvement to LALR(1) that is unique to Lark.

3.2.1 Grammar constraints in LALR(1)

Due to having only a lookahead of one token, LALR is limited in its ability to choose between rules, when they both match the input.

Tips for writing a conforming grammar:

- Try to avoid writing different rules that can match the same sequence of characters.
- For the best performance, prefer left-recursion over right-recursion.
- Consider setting terminal priority only as a last resort.

For a better understanding of these constraints, it's recommended to learn how a SLR parser works. SLR is very similar to LALR but much simpler.

3.3 CYK Parser

A [CYK parser](#) can parse any context-free grammar at $O(n^3|G|)$.

Its too slow to be practical for simple grammars, but it offers good performance for highly ambiguous grammars.

JSON PARSER - TUTORIAL

Lark is a parser - a program that accepts a grammar and text, and produces a structured tree that represents that text. In this tutorial we will write a JSON parser in Lark, and explore Lark's various features in the process.

It has 5 parts.

1. Writing the grammar
2. Creating the parser
3. Shaping the tree
4. Evaluating the tree
5. Optimizing

Knowledge assumed:

- Using Python
- A basic understanding of how to use regular expressions

4.1 Part 1 - The Grammar

Lark accepts its grammars in a format called **EBNF**. It basically looks like this:

```
rule_name : list of rules and TERMINALS to match
          | another possible list of items
          | etc.
```

```
TERMINAL: "some text to match"
```

(a terminal is a string or a regular expression)

The parser will try to match each rule (left-part) by matching its items (right-part) sequentially, trying each alternative (In practice, the parser is predictive so we don't have to try every alternative).

How to structure those rules is beyond the scope of this tutorial, but often it's enough to follow one's intuition.

In the case of JSON, the structure is simple: A json document is either a list, or a dictionary, or a string/number/etc.

The dictionaries and lists are recursive, and contain other json documents (or "values").

Let's write this structure in EBNF form:

```
value: dict
      | list
      | STRING
      | NUMBER
      | "true" | "false" | "null"

list : "[" [value ("," value)*] "]"

dict : "{" [pair ("," pair)*] "}"
pair : STRING ":" value
```

A quick explanation of the syntax:

- Parenthesis let us group rules together.
- `rule*` means *any amount*. That means, zero or more instances of that rule.
- `[rule]` means *optional*. That means zero or one instance of that rule.

Lark also supports the `rule+` operator, meaning one or more instances. It also supports the `rule?` operator which is another way to say *optional*.

Of course, we still haven't defined "STRING" and "NUMBER". Luckily, both these literals are already defined in Lark's common library:

```
%import common.ESCAPED_STRING -> STRING
%import common.SIGNED_NUMBER -> NUMBER
```

The arrow (`->`) renames the terminals. But that only adds obscurity in this case, so going forward we'll just use their original names.

We'll also take care of the white-space, which is part of the text, by simply matching and then throwing it away.

```
%import common.WS
%ignore WS
```

We tell our parser to ignore whitespace. Otherwise, we'd have to fill our grammar with WS terminals.

By the way, if you're curious what these terminals signify, they are roughly equivalent to this:

```
NUMBER : /-?\d+(\.\d+)?([eE][+-]?\d+)?/
STRING : /" .*?(?<!\\""/
%ignore /[ \t\n\f\r]+/
```

Lark will accept this way of writing too, if you really want to complicate your life :)

You can find the original definitions in [common.lark](#). They don't strictly adhere to [json.org](#) - but our purpose here is to accept json, not validate it.

Notice that terminals are written in UPPER-CASE, while rules are written in lower-case. I'll touch more on the differences between rules and terminals later.

4.2 Part 2 - Creating the Parser

Once we have our grammar, creating the parser is very simple.

We simply instantiate Lark, and tell it to accept a “value”:

```
from lark import Lark
json_parser = Lark(r"""
    value: dict
        | list
        | ESCAPED_STRING
        | SIGNED_NUMBER
        | "true" | "false" | "null"

    list : "[" [value ("," value)*] "]"

    dict : "{" [pair ("," pair)*] "}"
    pair : ESCAPED_STRING ":" value

    %import common.ESCAPED_STRING
    %import common.SIGNED_NUMBER
    %import common.WS
    %ignore WS

    """, start='value')
```

It’s that simple! Let’s test it out:

```
>>> text = '{"key": ["item0", "item1", 3.14]}'
>>> json_parser.parse(text)
Tree(value, [Tree(dict, [Tree(pair, [Token(String, "key"), Tree(value, [Tree(list,
↳ [Tree(value, [Token(String, "item0"))], Tree(value, [Token(String, "item1"))],
↳ Tree(value, [Token(Number, 3.14)]))]))]))]))))
>>> print( _.pretty() )
value
dict
pair
  "key"
  value
    list
      value      "item0"
      value      "item1"
      value      3.14
```

As promised, Lark automatically creates a tree that represents the parsed text.

But something is suspiciously missing from the tree. Where are the curly braces, the commas and all the other punctuation literals?

Lark automatically filters out literals from the tree, based on the following criteria:

- Filter out string literals without a name, or with a name that starts with an underscore.
- Keep regexps, even unnamed ones, unless their name starts with an underscore.

Unfortunately, this means that it will also filter out literals like “true” and “false”, and we will lose that information. The next section, “Shaping the tree” deals with this issue, and others.

4.3 Part 3 - Shaping the Tree

We now have a parser that can create a parse tree (or: AST), but the tree has some issues:

1. “true”, “false” and “null” are filtered out (test it out yourself!)
2. It has useless branches, like *value*, that clutter-up our view.

I’ll present the solution, and then explain it:

```
?value: dict
    | list
    | string
    | SIGNED_NUMBER      -> number
    | "true"             -> true
    | "false"            -> false
    | "null"             -> null
    ...

string : ESCAPED_STRING
```

1. Those little arrows signify *aliases*. An alias is a name for a specific part of the rule. In this case, we will name the *true/false/null* matches, and this way we won’t lose the information. We also alias *SIGNED_NUMBER* to mark it for later processing.
2. The question-mark prefixing *value* (“?value”) tells the tree-builder to inline this branch if it has only one member. In this case, *value* will always have only one member, and will always be inlined.
3. We turned the *ESCAPED_STRING* terminal into a rule. This way it will appear in the tree as a branch. This is equivalent to aliasing (like we did for the number), but now *string* can also be used elsewhere in the grammar (namely, in the *pair* rule).

Here is the new grammar:

```
from lark import Lark
json_parser = Lark(r"""
?value: dict
    | list
    | string
    | SIGNED_NUMBER      -> number
    | "true"             -> true
    | "false"            -> false
    | "null"             -> null
    ...

list : "[" [value ("," value)*] "]"

dict : "{" [pair ("," pair)*] "}"
pair : string ":" value

string : ESCAPED_STRING

%import common.ESCAPED_STRING
%import common.SIGNED_NUMBER
%import common.WS
%ignore WS
```

(continues on next page)

(continued from previous page)

```
""", start='value')
```

And let's test it out:

```
>>> text = '{"key": ["item0", "item1", 3.14, true]}'
>>> print( json_parser.parse(text).pretty() )
dict
  pair
    string      "key"
    list
      string      "item0"
      string      "item1"
      number      3.14
      true
```

Ah! That is much much nicer.

4.4 Part 4 - Evaluating the tree

It's nice to have a tree, but what we really want is a JSON object.

The way to do it is to evaluate the tree, using a Transformer.

A transformer is a class with methods corresponding to branch names. For each branch, the appropriate method will be called with the children of the branch as its argument, and its return value will replace the branch in the tree.

So let's write a partial transformer, that handles lists and dictionaries:

```
from lark import Transformer

class MyTransformer(Transformer):
    def list(self, items):
        return list(items)
    def pair(self, key_value):
        k, v = key_value
        return k, v
    def dict(self, items):
        return dict(items)
```

And when we run it, we get this:

```
>>> tree = json_parser.parse(text)
>>> MyTransformer().transform(tree)
{Tree(string, [Token(ANONRE_1, "key")]): Tree(string, [Token(ANONRE_1, "item0")]),
↳Tree(string, [Token(ANONRE_1, "item1")]), Tree(number, [Token(ANONRE_0, 3.14)]),
↳Tree(true, [])}
```

This is pretty close. Let's write a full transformer that can handle the terminals too.

Also, our definitions of list and dict are a bit verbose. We can do better:

```
from lark import Transformer

class TreeToJson(Transformer):
    def string(self, s):
        (s,) = s
        return s[1:-1]
    def number(self, n):
        (n,) = n
        return float(n)

    list = list
    pair = tuple
    dict = dict

    null = lambda self, _: None
    true = lambda self, _: True
    false = lambda self, _: False
```

And when we run it:

```
>>> tree = json_parser.parse(text)
>>> TreeToJson().transform(tree)
{'key': [u'item0', u'item1', 3.14, True]}
```

Magic!

4.5 Part 5 - Optimizing

4.5.1 Step 1 - Benchmark

By now, we have a fully working JSON parser, that can accept a string of JSON, and return its Pythonic representation.

But how fast is it?

Now, of course there are JSON libraries for Python written in C, and we can never compete with them. But since this is applicable to any parser you would write in Lark, let's see how far we can take this.

The first step for optimizing is to have a benchmark. For this benchmark I'm going to take data from json-generator.com/. I took their default suggestion and changed it to 5000 objects. The result is a 6.6MB sparse JSON file.

Our first program is going to be just a concatenation of everything we've done so far:

```
import sys
from lark import Lark, Transformer

json_grammar = r"""
?value: dict
    | list
    | string
    | SIGNED_NUMBER      -> number
    | "true"             -> true
    | "false"            -> false
```

(continues on next page)

(continued from previous page)

```

        | "null"          -> null

list : "[" [value ("," value)*] "]"

dict : "{" [pair ("," pair)*] "}"
pair : string ":" value

string : ESCAPED_STRING

%import common.ESCAPED_STRING
%import common.SIGNED_NUMBER
%import common.WS
%ignore WS
"""

class TreeToJson(Transformer):
    def string(self, s):
        (s,) = s
        return s[1:-1]
    def number(self, n):
        (n,) = n
        return float(n)

    list = list
    pair = tuple
    dict = dict

    null = lambda self, _: None
    true = lambda self, _: True
    false = lambda self, _: False

json_parser = Lark(json_grammar, start='value', lexer='basic')

if __name__ == '__main__':
    with open(sys.argv[1]) as f:
        tree = json_parser.parse(f.read())
        print(TreeToJson().transform(tree))

```

We run it and get this:

```

$ time python tutorial_json.py json_data > /dev/null

real      0m36.257s
user      0m34.735s
sys       0m1.361s

```

That's unsatisfactory time for a 6MB file. Maybe if we were parsing configuration or a small DSL, but we're trying to handle large amount of data here.

Well, turns out there's quite a bit we can do about it!

4.5.2 Step 2 - LALR(1)

So far we've been using the Earley algorithm, which is the default in Lark. Earley is powerful but slow. But it just so happens that our grammar is LR-compatible, and specifically LALR(1) compatible.

So let's switch to LALR(1) and see what happens:

```
json_parser = Lark(json_grammar, start='value', parser='lalr')
```

```
$ time python tutorial_json.py json_data > /dev/null
```

```
real      0m7.554s
user      0m7.352s
sys       0m0.148s
```

Ah, that's much better. The resulting JSON is of course exactly the same. You can run it for yourself and see.

It's important to note that not all grammars are LR-compatible, and so you can't always switch to LALR(1). But there's no harm in trying! If Lark lets you build the grammar, it means you're good to go.

4.5.3 Step 3 - Tree-less LALR(1)

So far, we've built a full parse tree for our JSON, and then transformed it. It's a convenient method, but it's not the most efficient in terms of speed and memory. Luckily, Lark lets us avoid building the tree when parsing with LALR(1).

Here's the way to do it:

```
json_parser = Lark(json_grammar, start='value', parser='lalr', transformer=TreeToJson())

if __name__ == '__main__':
    with open(sys.argv[1]) as f:
        print( json_parser.parse(f.read()) )
```

We've used the transformer we've already written, but this time we plug it straight into the parser. Now it can avoid building the parse tree, and just send the data straight into our transformer. The `parse()` method now returns the transformed JSON, instead of a tree.

Let's benchmark it:

```
real      0m4.866s
user      0m4.722s
sys       0m0.121s
```

That's a measurable improvement! Also, this way is more memory efficient. Check out the benchmark table at the end to see just how much.

As a general practice, it's recommended to work with parse trees, and only skip the tree-builder when your transformer is already working.

4.5.4 Step 4 - PyPy

PyPy is a JIT engine for running Python, and it's designed to be a drop-in replacement.

Lark is written purely in Python, which makes it very suitable for PyPy.

Let's get some free performance:

```
$ time pypy tutorial_json.py json_data > /dev/null

real      0m1.397s
user      0m1.296s
sys       0m0.083s
```

PyPy is awesome!

4.5.5 Conclusion

We've brought the run-time down from 36 seconds to 1.1 seconds, in a series of small and simple steps.

Now let's compare the benchmarks in a nicely organized table.

I measured memory consumption using a little script called [memusg](#)

I added a few other parsers for comparison. [PyParsing](#) and [funcparselib](#) fair pretty well in their memory usage (they don't build a tree), but they can't compete with the run-time speed of LALR(1).

These benchmarks are for Lark's alpha version. I already have several optimizations planned that will significantly improve run-time speed.

Once again, shout-out to PyPy for being so effective.

4.6 Afterword

This is the end of the tutorial. I hoped you liked it and learned a little about Lark.

To see what else you can do with Lark, check out the [examples](#).

Read the documentation here: <https://lark-parser.readthedocs.io/en/latest/>

HOW TO USE LARK - GUIDE

5.1 Work process

This is the recommended process for working with Lark:

1. Collect or create input samples, that demonstrate key features or behaviors in the language you're trying to parse.
2. Write a grammar. Try to aim for a structure that is intuitive, and in a way that imitates how you would explain your language to a fellow human.
3. Try your grammar in Lark against each input sample. Make sure the resulting parse-trees make sense.
4. Use Lark's grammar features to *shape the tree*: Get rid of superfluous rules by inlining them, and use aliases when specific cases need clarification.

You can perform steps 1-4 repeatedly, gradually growing your grammar to include more sentences.

5. Create a transformer to evaluate the parse-tree into a structure you'll be comfortable to work with. This may include evaluating literals, merging branches, or even converting the entire tree into your own set of AST classes.

Of course, some specific use-cases may deviate from this process. Feel free to suggest these cases, and I'll add them to this page.

5.2 Getting started

Browse the [Examples](#) to find a template that suits your purposes.

Read the tutorials to get a better understanding of how everything works. (links in the [main page](#))

Use the [Cheatsheet \(PDF\)](#) for quick reference.

Use the reference pages for more in-depth explanations. (links in the [main page](#))

5.3 Debug

Grammars may contain non-obvious bugs, usually caused by rules or terminals interfering with each other in subtle ways.

When trying to debug a misbehaving grammar, the following methodology is recommended:

1. Create a copy of the grammar, so you can change the parser/grammar without any worries
2. Find the minimal input that creates the error
3. Slowly remove rules from the grammar, while making sure the error still occurs.

Usually, by the time you get to a minimal grammar, the problem becomes clear.

But if it doesn't, feel free to ask us on [gitter](#), or even open an issue. Post a reproducing code, with the minimal grammar and input, and we'll do our best to help.

5.3.1 Regex collisions

A likely source of bugs occurs when two regexes in a grammar can match the same input. If both terminals have the same priority, most lexers would arbitrarily choose the first one that matches, which isn't always the desired one. (a notable exception is the `dynamic_complete` lexer, which always tries all variations. But its users pay for that with performance.)

These collisions can be hard to notice, and their effects can be difficult to debug, as they are subtle and sometimes hard to reproduce.

To help with these situations, Lark can utilize a new external library called `interegular`. If it is installed, Lark uses it to check for collisions, and warn about any conflicts that it can find:

```
import logging
from lark import Lark, logger

logger.setLevel(logging.WARN)

collision_grammar = '''
start: A | B
A: /a+/
B: /[ab]+/
'''

p = Lark(collision_grammar, parser='lalr')

# Output:
# Collision between Terminals B and A. The lexer will choose between them arbitrarily
# Example Collision: a
```

You can install `interegular` for Lark using `pip install 'lark[interegular]'`.

Note 1: `Interegular` currently only runs when the lexer is `basic` or `contextual`.

Note 2: Some advanced regex features, such as `lookahead` and `lookbehind`, may prevent `interegular` from detecting existing collisions.

5.3.2 Shift/Reduce collisions

By default Lark automatically resolves Shift/Reduce conflicts as Shift. It produces notifications as debug messages. when users pass `debug=True`, those notifications are written as warnings.

Either way, to get the messages printed you have to configure the `logger` beforehand. For example:

```
import logging
from lark import Lark, logger

logger.setLevel(logging.DEBUG)

collision_grammar = '''
```

(continues on next page)

(continued from previous page)

```

start: as as
as: a*
a: "a"
'''

p = Lark(collision_grammar, parser='lalr', debug=True)
# Shift/Reduce conflict for terminal A: (resolving as shift)
# * <as : >
# Shift/Reduce conflict for terminal A: (resolving as shift)
# * <as : __as_star_0>

```

5.3.3 Strict-Mode

Lark, by default, accepts grammars with unresolved Shift/Reduce collisions (which it always resolves to shift), and regex collisions.

Strict-mode allows users to validate that their grammars don't contain these collisions.

When Lark is initialized with `strict=True`, it raises an exception on any Shift/Reduce or regex collision.

If `interegular` isn't installed, an exception is thrown.

When using strict-mode, users will be expected to resolve their collisions manually:

- To resolve Shift/Reduce collisions, adjust the priority weights of the rules involved, until there are no more collisions.
- To resolve regex collisions, change the involved regexes so that they can no longer both match the same input (Lark provides an example).

Strict-mode only applies to LALR for now.

```

from lark import Lark

collision_grammar = '''
start: as as
as: a*
a: "a"
'''

p = Lark(collision_grammar, parser='lalr', strict=True)

# Traceback (most recent call last):
# ...
# lark.exceptions.GrammarError: Shift/Reduce conflict for terminal A. [strict-mode]

```

5.4 Tools

5.4.1 Stand-alone parser

Lark can generate a stand-alone LALR(1) parser from a grammar.

The resulting module provides the same interface as Lark, but with a fixed grammar, and reduced functionality.

Run using:

```
python -m lark.tools.standalone
```

For a play-by-play, read the [tutorial](#)

5.4.2 Import Nearley.js grammars

It is possible to import Nearley grammars into Lark. The Javascript code is translated using Js2Py.

See the [tools page](#) for more information.

HOW TO DEVELOP LARK - GUIDE

There are many ways you can help the project:

- Help solve issues
- Improve the documentation
- Write new grammars for Lark's library
- Write a blog post introducing Lark to your audience
- Port Lark to another language
- Help with code development

If you're interested in taking one of these on, contact us on [Gitter](#) or [Github Discussion](#), and we will provide more details and assist you in the process.

6.1 Code Style

Lark does not follow a predefined code style. We accept any code style that makes sense, as long as it's Pythonic and easy to read.

6.2 Unit Tests

Lark comes with an extensive set of tests. Many of the tests will run several times, once for each parser configuration.

To run the tests, just go to the lark project root, and run the command:

```
python -m tests
```

or

```
pypy -m tests
```

For a list of supported interpreters, you can consult the `tox.ini` file.

You can also run a single unittest using its class and method name, for example:

```
## test_package test_class_name.test_function_name
python -m tests TestLalrBasic.test_keep_all_tokens
```

6.2.1 tox

To run all Unit Tests with tox, install tox and Python 2.7 up to the latest python interpreter supported (consult the file `tox.ini`). Then, run the command `tox` on the root of this project (where the main `setup.py` file is on).

And, for example, if you would like to only run the Unit Tests for Python version 2.7, you can run the command `tox -e py27`

6.2.2 pytest

You can also run the tests using pytest:

```
pytest tests
```

6.2.3 Using setup.py

Another way to run the tests is using `setup.py`:

```
python setup.py test
```

RECIPES

A collection of recipes to use Lark and its various features

7.1 Use a transformer to parse integer tokens

Transformers are the common interface for processing matched rules and tokens.

They can be used during parsing for better performance.

```
from lark import Lark, Transformer

class T(Transformer):
    def INT(self, tok):
        "Convert the value of `tok` from string to int, while maintaining line number &
        ↪column."
        return tok.update(value=int(tok))

parser = Lark("""
start: INT*
%import common.INT
%ignore " "
""", parser="lalr", transformer=T())

print(parser.parse('3 14 159'))
```

Prints out:

```
Tree(start, [Token(INT, 3), Token(INT, 14), Token(INT, 159)])
```

7.2 Collect all comments with lexer_callbacks

lexer_callbacks can be used to interface with the lexer as it generates tokens.

It accepts a dictionary of the form

```
{TOKEN_TYPE: callback}
```

Where callback is of type `f(Token) -> Token`

It only works with the basic and contextual lexers.

This has the same effect of using a transformer, but can also process ignored tokens.

```
from lark import Lark

comments = []

parser = Lark("""
    start: INT*

    COMMENT: /#.* /

    %import common (INT, WS)
    %ignore COMMENT
    %ignore WS
""", parser="lalr", lexer_callbacks={'COMMENT': comments.append})

parser.parse("""
1 2 3 # hello
# world
4 5 6
""")

print(comments)
```

Prints out:

```
[Token(COMMENT, '# hello'), Token(COMMENT, '# world')]
```

Note: We don't have to return a token, because comments are ignored

7.3 CollapseAmbiguities

Parsing ambiguous texts with `earley` and `ambiguity='explicit'` produces a single tree with `_ambig` nodes to mark where the ambiguity occurred.

However, it's sometimes more convenient instead to work with a list of all possible unambiguous trees.

Lark provides a utility transformer for that purpose:

```
from lark import Lark, Tree, Transformer
from lark.visitors import CollapseAmbiguities

grammar = """
!start: x y

!x: "a" "b"
    | "ab"
    | "abc"

!y: "c" "d"
    | "cd"
    | "d"

```

(continues on next page)

(continued from previous page)

```

"""
parser = Lark(grammar, ambiguity='explicit')

t = parser.parse('abcd')
for x in CollapseAmbiguities().transform(t):
    print(x.pretty())

```

This prints out:

```

start
x
  a
  b
y
  c
  d

start
x      ab
y      cd

start
x      abc
y      d

```

While convenient, this should be used carefully, as highly ambiguous trees will soon create an exponential explosion of such unambiguous derivations.

7.4 Keeping track of parents when visiting

The following visitor assigns a `parent` attribute for every node in the tree.

If your tree nodes aren't unique (if there is a shared `Tree` instance), the assert will fail.

```

class Parent(Visitor):
    def __default__(self, tree):
        for subtree in tree.children:
            if isinstance(subtree, Tree):
                assert not hasattr(subtree, 'parent')
                subtree.parent = proxy(tree)

```

7.5 Unwinding VisitError after a transformer/visitor exception

Errors that happen inside visitors and transformers get wrapped inside a `VisitError` exception.

This can often be inconvenient, if you wish the actual error to propagate upwards, or if you want to catch it.

But, it's easy to unwrap it at the point of calling the transformer, by catching it and raising the `VisitError.orig_exc` attribute.

For example:

```
from lark import Lark, Transformer
from lark.visitors import VisitError

tree = Lark('start: "a"').parse('a')

class T(Transformer):
    def start(self, x):
        raise KeyError("Original Exception")

t = T()
try:
    print( t.transform(tree))
except VisitError as e:
    raise e.orig_exc
```

7.6 Adding a Progress Bar to Parsing with tqdm

Parsing large files can take a long time, even with the `parser='lalr'` option. To make this process more user-friendly, it's useful to add a progress bar. One way to achieve this is to use the `InteractiveParser` to display each token as it is processed. In this example, we use `tqdm`, but a similar approach should work with GUIs.

```
from tqdm import tqdm

def parse_with_progress(parser: Lark, text: str, start=None):
    last = 0
    progress = tqdm(total=len(text))
    pi = parser.parse_interactive(text, start=start)
    for token in pi.iter_parse():
        if token.end_pos is not None:
            progress.update(token.end_pos - last)
            last = token.end_pos
    return pi.result
```

Note that we don't simply wrap the iterable because `tqdm` would not be able to determine the total. Additionally, keep in mind that this implementation relies on the `InteractiveParser` and, therefore, only works with the `LALR(1)` parser, not `earley`.

EXAMPLES FOR LARK

How to run the examples:

After cloning the repo, open the terminal into the root directory of the project, and run the following:

```
[lark]$ python -m examples.<name_of_example>
```

For example, the following will parse all the Python files in the standard library of your local installation:

```
[lark]$ python -m examples.advanced.python_parser
```

8.1 Beginner Examples

8.1.1 Parsing Indentation

A demonstration of parsing indentation (“whitespace significant” language) and the usage of the Indenter class.

Since indentation is context-sensitive, a postlex stage is introduced to manufacture INDENT/DEDENT tokens.

It is crucial for the indenter that the NL_type matches the spaces (and tabs) after the newline.

```
from lark import Lark
from lark.indenter import Indenter

tree_grammar = r"""
    ?start: _NL* tree

    tree: NAME _NL [_INDENT tree+ _DEDENT]

    %import common.CNAME -> NAME
    %import common.WS_INLINE
    %declare _INDENT _DEDENT
    %ignore WS_INLINE

    _NL: /(\r?\n[\t ]*)+/
"""

class TreeIndenter(Indenter):
    NL_type = '_NL'
    OPEN_PAREN_types = []
    CLOSE_PAREN_types = []
```

(continues on next page)

(continued from previous page)

```

    INDENT_type = '_INDENT'
    DEDENT_type = '_DEDENT'
    tab_len = 8

parser = Lark(tree_grammar, parser='lalr', postlex=TreeIndenter())

test_tree = """
a
  b
  c
    d
    e
  f
    g
"""

def test():
    print(parser.parse(test_tree).pretty())

if __name__ == '__main__':
    test()

```

Total running time of the script: (0 minutes 0.000 seconds)

8.1.2 Lark Grammar

A reference implementation of the Lark grammar (using LALR(1))

```

import lark
from pathlib import Path

examples_path = Path(__file__).parent
lark_path = Path(lark.__file__).parent

parser = lark.Lark.open(lark_path / 'grammars/lark.lark', rel_to=__file__, parser="lalr")

grammar_files = [
    examples_path / 'advanced/python2.lark',
    examples_path / 'relative-imports/multiples.lark',
    examples_path / 'relative-imports/multiple2.lark',
    examples_path / 'relative-imports/multiple3.lark',
    examples_path / 'tests/no_newline_at_end.lark',
    examples_path / 'tests/negative_priority.lark',
    examples_path / 'standalone/json.lark',
    lark_path / 'grammars/common.lark',
    lark_path / 'grammars/lark.lark',
    lark_path / 'grammars/unicode.lark',
    lark_path / 'grammars/python.lark',
]

```

(continues on next page)

(continued from previous page)

```
def test():
    for grammar_file in grammar_files:
        tree = parser.parse(open(grammar_file).read())
        print("All grammars parsed successfully")

if __name__ == '__main__':
    test()
```

Total running time of the script: (0 minutes 0.000 seconds)

8.1.3 Handling Ambiguity

A demonstration of ambiguity

This example shows how to use get explicit ambiguity from Lark's Earley parser.

```
import sys
from lark import Lark, tree

grammar = """
    sentence: noun verb noun      -> simple
             | noun verb "like" noun -> comparative

    noun: adj? NOUN
    verb: VERB
    adj: ADJ

    NOUN: "flies" | "bananas" | "fruit"
    VERB: "like" | "flies"
    ADJ: "fruit"

    %import common.WS
    %ignore WS
"""

parser = Lark(grammar, start='sentence', ambiguity='explicit')

sentence = 'fruit flies like bananas'

def make_png(filename):
    tree.pydot__tree_to_png( parser.parse(sentence), filename)

def make_dot(filename):
    tree.pydot__tree_to_dot( parser.parse(sentence), filename)

if __name__ == '__main__':
    print(parser.parse(sentence).pretty())
    # make_png(sys.argv[1])
    # make_dot(sys.argv[1])

# Output:
#
```

(continues on next page)

(continued from previous page)

```
# _ambig
#   comparative
#     noun fruit
#     verb flies
#     noun bananas
#   simple
#     noun
#       fruit
#       flies
#     verb like
#     noun bananas
#
# (or view a nicer version at "./fruitflies.png")
```

Total running time of the script: (0 minutes 0.000 seconds)

8.1.4 Basic calculator

A simple example of a REPL calculator

This example shows how to write a basic calculator with variables.

```
from lark import Lark, Transformer, v_args

try:
    input = raw_input    # For Python2 compatibility
except NameError:
    pass

calc_grammar = """
?start: sum
      | NAME "=" sum    -> assign_var

?sum: product
     | sum "+" product  -> add
     | sum "-" product  -> sub

?product: atom
         | product "*" atom -> mul
         | product "/" atom -> div

?atom: NUMBER          -> number
      | "-" atom        -> neg
      | NAME            -> var
      | "(" sum ")"

%import common.CNAME -> NAME
%import common.NUMBER
%import common.WS_INLINE
```

(continues on next page)

(continued from previous page)

```

    %ignore WS_INLINE
    """

    @v_args(inline=True)    # Affects the signatures of the methods
    class CalculateTree(Transformer):
        from operator import add, sub, mul, truediv as div, neg
        number = float

        def __init__(self):
            self.vars = {}

        def assign_var(self, name, value):
            self.vars[name] = value
            return value

        def var(self, name):
            try:
                return self.vars[name]
            except KeyError:
                raise Exception("Variable not found: %s" % name)

    calc_parser = Lark(calc_grammar, parser='lalr', transformer=CalculateTree())
    calc = calc_parser.parse

    def main():
        while True:
            try:
                s = input('> ')
            except EOFError:
                break
            print(calc(s))

    def test():
        print(calc("a = 1+2"))
        print(calc("1+a*-3"))

    if __name__ == '__main__':
        # test()
        main()

```

Total running time of the script: (0 minutes 0.000 seconds)

8.1.5 Turtle DSL

Implements a LOGO-like toy language for Python's turtle, with interpreter.

```
try:
    input = raw_input    # For Python2 compatibility
except NameError:
    pass

import turtle

from lark import Lark

turtle_grammar = """
    start: instruction+

    instruction: MOVEMENT NUMBER          -> movement
                | "c" COLOR [COLOR]      -> change_color
                | "fill" code_block       -> fill
                | "repeat" NUMBER code_block -> repeat

    code_block: "{" instruction+ "}"

    MOVEMENT: "f"|"b"|"l"|"r"
    COLOR: LETTER+

    %import common.LETTER
    %import common.INT -> NUMBER
    %import common.WS
    %ignore WS
"""

parser = Lark(turtle_grammar)

def run_instruction(t):
    if t.data == 'change_color':
        turtle.color(*t.children)    # We just pass the color names as-is

    elif t.data == 'movement':
        name, number = t.children
        { 'f': turtle.fd,
          'b': turtle.bk,
          'l': turtle.lt,
          'r': turtle.rt, }[name](int(number))

    elif t.data == 'repeat':
        count, block = t.children
        for i in range(int(count)):
            run_instruction(block)

    elif t.data == 'fill':
        turtle.begin_fill()
        run_instruction(t.children[0])
```

(continues on next page)

(continued from previous page)

```

        turtle.end_fill()

    elif t.data == 'code_block':
        for cmd in t.children:
            run_instruction(cmd)
    else:
        raise SyntaxError('Unknown instruction: %s' % t.data)

def run_turtle(program):
    parse_tree = parser.parse(program)
    for inst in parse_tree.children:
        run_instruction(inst)

def main():
    while True:
        code = input('> ')
        try:
            run_turtle(code)
        except Exception as e:
            print(e)

def test():
    text = """
        c red yellow
        fill { repeat 36 {
            f200 1170
        }}
    """
    run_turtle(text)

if __name__ == '__main__':
    # test()
    main()

```

Total running time of the script: (0 minutes 0.000 seconds)

8.1.6 Simple JSON Parser

The code is short and clear, and outperforms every other parser (that's written in Python). For an explanation, check out the JSON parser tutorial at /docs/json_tutorial.md

```

import sys

from lark import Lark, Transformer, v_args

json_grammar = r"""
    ?start: value

    ?value: object
          | array

```

(continues on next page)

(continued from previous page)

```

    | string
    | SIGNED_NUMBER      -> number
    | "true"             -> true
    | "false"            -> false
    | "null"             -> null

array : "[" [value ("," value)*] "]"
object : "{" [pair ("," pair)*] "}"
pair : string ":" value

string : ESCAPED_STRING

%import common.ESCAPED_STRING
%import common.SIGNED_NUMBER
%import common.WS

%ignore WS
"""

class TreeToJson(Transformer):
    @v_args(inline=True)
    def string(self, s):
        return s[1:-1].replace('\\"', '"')

    array = list
    pair = tuple
    object = dict
    number = v_args(inline=True)(float)

    null = lambda self, _: None
    true = lambda self, _: True
    false = lambda self, _: False

### Create the JSON parser with Lark, using the Earley algorithm
# json_parser = Lark(json_grammar, parser='earley', lexer='basic')
# def parse(x):
#     return TreeToJson().transform(json_parser.parse(x))

### Create the JSON parser with Lark, using the LALR algorithm
json_parser = Lark(json_grammar, parser='lalr',
    # Using the basic lexer isn't required, and isn't usually recommended.
    # But, it's good enough for JSON, and it's slightly faster.
    lexer='basic',
    # Disabling propagate_positions and placeholders slightly improves
    ↪ speed
    propagate_positions=False,
    maybe_placeholders=False,
    # Using an internal transformer is faster and more memory efficient
    transformer=TreeToJson())
parse = json_parser.parse

```

(continues on next page)

(continued from previous page)

```
def test():
    test_json = '''
        {
            "empty_object" : {},
            "empty_array"  : [],
            "booleans"     : { "YES" : true, "NO" : false },
            "numbers"      : [ 0, 1, -2, 3.3, 4.4e5, 6.6e-7 ],
            "strings"      : [ "This", [ "And" , "That", "And a \\\"b\" ] ],
            "nothing"      : null
        }
    '''

    j = parse(test_json)
    print(j)
    import json
    assert j == json.loads(test_json)

if __name__ == '__main__':
    # test()
    with open(sys.argv[1]) as f:
        print(parse(f.read()))
```

Total running time of the script: (0 minutes 0.000 seconds)

8.2 Advanced Examples

GRAMMAR COMPOSITION

This example shows how to do grammar composition in Lark, by creating a new file format that allows both CSV and JSON to co-exist.

We show how, by using namespaces, Lark grammars and their transformers can be fully reused - they don't need to care if their grammar is used directly, or being imported, or who is doing the importing.

See [main.py](#) for more details.

EXAMPLE GRAMMARS

This directory is a collection of lark grammars, taken from real world projects.

- Verilog - Taken from <https://github.com/circuitgraph/circuitgraph/blob/main/circuitgraph/parsing/verilog.lark>

STANDALONE EXAMPLE

To initialize, cd to this folder, and run:

```
./create_standalone.sh
```

Or:

```
python -m lark.tools.standalone json.lark > json_parser.py
```

Then run using:

```
python json_parser_main.py <path-to.json>
```

11.1 Advanced Examples

11.1.1 LALR's contextual lexer

This example demonstrates the power of LALR's contextual lexer, by parsing a toy configuration language.

The terminals *NAME* and *VALUE* overlap. They can match the same input. A basic lexer would arbitrarily choose one over the other, based on priority, which would lead to a (confusing) parse error. However, due to the unambiguous structure of the grammar, Lark's LALR(1) algorithm knows which one of them to expect at each point during the parse. The lexer then only matches the tokens that the parser expects. The result is a correct parse, something that is impossible with a regular lexer.

Another approach is to use the Earley algorithm. It will handle more cases than the contextual lexer, but at the cost of performance. See `examples/conf_earley.py` for an example of that approach.

```
from lark import Lark

parser = Lark(r"""
start: _NL? section+
section: "[" NAME "]" _NL item+
item: NAME "=" VALUE? _NL

NAME: /\w/+
VALUE: /\./+

%import common.NEWLINE -> _NL
%import common.WS_INLINE
%ignore WS_INLINE
```

(continues on next page)

(continued from previous page)

```
        """ , parser="lalr")

sample_conf = """
[bla]
a=Hello
this="that",4
empty=
"""

print(parser.parse(sample_conf).pretty())
```

Total running time of the script: (0 minutes 0.000 seconds)

11.1.2 Templates

This example shows how to use Lark's templates to achieve cleaner grammars

```
from lark import Lark

grammar = r"""
start: list | dict

list: "[" _seperated{atom, ","} "]"
dict: "{" _seperated{key_value, ","} "}"
key_value: atom ":" atom

_seperated{x, sep}: x (sep x)* // Define a sequence of 'x sep x sep x ...'

atom: NUMBER | ESCAPED_STRING

%import common (NUMBER, ESCAPED_STRING, WS)
%ignore WS
"""

parser = Lark(grammar)

print(parser.parse('[1, "a", 2]'))
print(parser.parse('{ "a": 2, "b": 6 }'))
```

Total running time of the script: (0 minutes 0.000 seconds)

11.1.3 Earley's dynamic lexer

Demonstrates the power of Earley's dynamic lexer on a toy configuration language

Using a lexer for configuration files is tricky, because values don't have to be surrounded by delimiters. Using a basic lexer for this just won't work.

In this example we use a dynamic lexer and let the Earley parser resolve the ambiguity.

Another approach is to use the contextual lexer with LALR. It is less powerful than Earley, but it can handle some ambiguity when lexing and it's much faster. See `examples/conf_lalr.py` for an example of that approach.

```
from lark import Lark

parser = Lark(r"""
    start: _NL? section+
    section: "[" NAME "]" _NL item+
    item: NAME "=" VALUE? _NL

    NAME: /\w/+
    VALUE: /.+/+

    %import common.NEWLINE -> _NL
    %import common.WS_INLINE
    %ignore WS_INLINE
""", parser="earley")

def test():
    sample_conf = """
[bla]

a=Hello
this="that",4
empty=
"""

    r = parser.parse(sample_conf)
    print (r.pretty())

if __name__ == '__main__':
    test()
```

Total running time of the script: (0 minutes 0.000 seconds)

11.1.4 Error handling using an interactive parser

This example demonstrates error handling using an interactive parser in LALR

When the parser encounters an `UnexpectedToken` exception, it creates a an interactive parser with the current parse-state, and lets you control how to proceed step-by-step. When you've achieved the correct parse-state, you can resume the run by returning `True`.

```
from lark import Token

from _json_parser import json_parser
```

(continues on next page)

(continued from previous page)

```

def ignore_errors(e):
    if e.token.type == 'COMMA':
        # Skip comma
        return True
    elif e.token.type == 'SIGNED_NUMBER':
        # Try to feed a comma and retry the number
        e.interactive_parser.feed_token(Token('COMMA', ','))
        e.interactive_parser.feed_token(e.token)
        return True

    # Unhandled error. Will stop parse and raise exception
    return False

def main():
    s = "[0 1, 2,, 3,,, 4, 5 6 ]"
    res = json_parser.parse(s, on_error=ignore_errors)
    print(res)      # prints [0.0, 1.0, 2.0, 3.0, 4.0, 5.0, 6.0]

main()

```

Total running time of the script: (0 minutes 0.000 seconds)

11.1.5 Reconstruct a JSON

Demonstrates the experimental text-reconstruction feature

The Reconstructor takes a parse tree (already filtered from punctuation, of course), and reconstructs it into correct text, that can be parsed correctly. It can be useful for creating “hooks” to alter data before handing it to other parsers. You can also use it to generate samples from scratch.

```

import json

from lark import Lark
from lark.reconstruct import Reconstructor

from _json_parser import json_grammar

test_json = '''
{
    "empty_object" : {},
    "empty_array"  : [],
    "booleans"     : { "YES" : true, "NO" : false },
    "numbers"      : [ 0, 1, -2, 3.3, 4.4e5, 6.6e-7 ],
    "strings"      : [ "This", [ "And" , "That", "And a \\\"b\" ] ],
    "nothing"      : null
}
'''

def test_earley():

```

(continues on next page)

(continued from previous page)

```

json_parser = Lark(json_grammar, maybe_placeholders=False)
tree = json_parser.parse(test_json)

new_json = Reconstructor(json_parser).reconstruct(tree)
print (new_json)
print (json.loads(new_json) == json.loads(test_json))

def test_lalr():

    json_parser = Lark(json_grammar, parser='lalr', maybe_placeholders=False)
    tree = json_parser.parse(test_json)

    new_json = Reconstructor(json_parser).reconstruct(tree)
    print (new_json)
    print (json.loads(new_json) == json.loads(test_json))

test_earley()
test_lalr()

```

Total running time of the script: (0 minutes 0.000 seconds)

11.1.6 Custom lexer

Demonstrates using a custom lexer to parse a non-textual stream of data

You can use a custom lexer to tokenize text when the lexers offered by Lark are too slow, or not flexible enough.

You can also use it (as shown in this example) to tokenize streams of objects.

```

from lark import Lark, Transformer, v_args
from lark.lexer import Lexer, Token

class TypeLexer(Lexer):
    def __init__(self, lexer_conf):
        pass

    def lex(self, data):
        for obj in data:
            if isinstance(obj, int):
                yield Token('INT', obj)
            elif isinstance(obj, (type(''), type(u''))):
                yield Token('STR', obj)
            else:
                raise TypeError(obj)

parser = Lark("""
    start: data_item+
    data_item: STR INT*

    %declare STR INT
    """, parser='lalr', lexer=TypeLexer)

```

(continues on next page)

(continued from previous page)

```

class ParseToDict(Transformer):
    @v_args(inline=True)
    def data_item(self, name, *numbers):
        return name.value, [n.value for n in numbers]

    start = dict

def test():
    data = ['alice', 1, 27, 3, 'bob', 4, 'carrie', 'dan', 8, 6]

    print(data)

    tree = parser.parse(data)
    res = ParseToDict().transform(tree)

    print('-->')
    print(res) # prints {'alice': [1, 27, 3], 'bob': [4], 'carrie': [], 'dan': [8, 6]}

if __name__ == '__main__':
    test()

```

Total running time of the script: (0 minutes 0.000 seconds)

11.1.7 Transform a Forest

This example demonstrates how to subclass `TreeForestTransformer` to directly transform a SPPF.

```

from lark import Lark
from lark.parsers.earley_forest import TreeForestTransformer, handles_ambiguity, Discard

class CustomTransformer(TreeForestTransformer):

    @handles_ambiguity
    def sentence(self, trees):
        return next(tree for tree in trees if tree.data == 'simple')

    def simple(self, children):
        children.append('.')
        return self.tree_class('simple', children)

    def adj(self, children):
        return Discard

    def __default_token__(self, token):
        return token.capitalize()

grammar = """

```

(continues on next page)

(continued from previous page)

```

sentence: noun verb noun      -> simple
        | noun verb "like" noun -> comparative

noun: adj? NOUN
verb: VERB
adj: ADJ

NOUN: "flies" | "bananas" | "fruit"
VERB: "like" | "flies"
ADJ: "fruit"

%import common.WS
%ignore WS
"""

parser = Lark(grammar, start='sentence', ambiguity='forest')
sentence = 'fruit flies like bananas'
forest = parser.parse(sentence)

tree = CustomTransformer(resolve_ambiguity=False).transform(forest)
print(tree.pretty())

# Output:
#
# simple
#   noun  Flies
#   verb  Like
#   noun  Bananas
#   .
#
#

```

Total running time of the script: (0 minutes 0.000 seconds)

11.1.8 Simple JSON Parser

The code is short and clear, and outperforms every other parser (that's written in Python). For an explanation, check out the JSON parser tutorial at /docs/json_tutorial.md

(this is here for use by the other examples)

```

from lark import Lark, Transformer, v_args

json_grammar = r"""
?start: value

?value: object
      | array
      | string
      | SIGNED_NUMBER      -> number
      | "true"             -> true
      | "false"            -> false
      | "null"             -> null

```

(continues on next page)

(continued from previous page)

```

array : "[" [value ("," value)*] "]"
object : "{" [pair ("," pair)*] "}"
pair : string ":" value

string : ESCAPED_STRING

%import common.ESCAPED_STRING
%import common.SIGNED_NUMBER
%import common.WS

%ignore WS
"""

class TreeToJson(Transformer):
    @v_args(inline=True)
    def string(self, s):
        return s[1:-1].replace('\\"', '"')

    array = list
    pair = tuple
    object = dict
    number = v_args(inline=True)(float)

    null = lambda self, _: None
    true = lambda self, _: True
    false = lambda self, _: False

### Create the JSON parser with Lark, using the LALR algorithm
json_parser = Lark(json_grammar, parser='lalr',
    # Using the basic lexer isn't required, and isn't usually recommended.
    # But, it's good enough for JSON, and it's slightly faster.
    lexer='basic',
    # Disabling propagate_positions and placeholders slightly improves
    ↪ speed
    propagate_positions=False,
    maybe_placeholders=False,
    # Using an internal transformer is faster and more memory efficient
    transformer=TreeToJson())

```

Total running time of the script: (0 minutes 0.000 seconds)

11.1.9 Custom SPPF Prioritizer

This example demonstrates how to subclass `ForestVisitor` to make a custom SPPF node prioritizer to be used in conjunction with `TreeForestTransformer`.

Our prioritizer will count the number of descendants of a node that are tokens. By negating this count, our prioritizer will prefer nodes with fewer token descendants. Thus, we choose the more specific parse.

```
from lark import Lark
from lark.parsers.earley_forest import ForestVisitor, TreeForestTransformer

class TokenPrioritizer(ForestVisitor):

    def visit_symbol_node_in(self, node):
        # visit the entire forest by returning node.children
        return node.children

    def visit_packed_node_in(self, node):
        return node.children

    def visit_symbol_node_out(self, node):
        priority = 0
        for child in node.children:
            # Tokens do not have a priority attribute
            # count them as -1
            priority += getattr(child, 'priority', -1)
        node.priority = priority

    def visit_packed_node_out(self, node):
        priority = 0
        for child in node.children:
            priority += getattr(child, 'priority', -1)
        node.priority = priority

    def on_cycle(self, node, path):
        raise Exception("Oops, we encountered a cycle.")

grammar = """
start: hello " " world | hello_world
hello: "Hello"
world: "World"
hello_world: "Hello World"
"""

parser = Lark(grammar, parser='earley', ambiguity='forest')
forest = parser.parse("Hello World")

print("Default prioritizer:")
tree = TreeForestTransformer(resolve_ambiguity=True).transform(forest)
print(tree.pretty())

forest = parser.parse("Hello World")

print("Custom prioritizer:")
```

(continues on next page)

(continued from previous page)

```

tree = TreeForestTransformer(resolve_ambiguity=True, prioritizer=TokenPrioritizer()).
↳ transform(forest)
print(tree.pretty())

# Output:
#
# Default prioritizer:
# start
#   hello Hello
#
#   world World
#
# Custom prioritizer:
# start
#   hello_world   Hello World

```

Total running time of the script: (0 minutes 0.000 seconds)

11.1.10 Python 3 to Python 2 converter (tree templates)

This example demonstrates how to translate between two trees using tree templates. It parses Python 3, translates it to a Python 2 AST, and then outputs the result as Python 2 code.

Uses `reconstruct_python.py` for generating the final Python 2 code.

```

from lark import Lark
from lark.tree_templates import TemplateConf, TemplateTranslator

from lark.indenter import PythonIndenter
from reconstruct_python import PythonReconstructor

#
# 1. Define a Python parser that also accepts template vars in the code (in the form of
↳ $var)
#
TEMPLATED_PYTHON = r"""
%import python (single_input, file_input, eval_input, atom, var, stmt, expr, testlist_
↳ star_expr, _NEWLINE, _INDENT, _DEDENT, COMMENT, NAME)

%extend atom: TEMPLATE_NAME -> var

TEMPLATE_NAME: "$" NAME

?template_start: (stmt | testlist_star_expr _NEWLINE)

%ignore /\t \f +/           // WS
%ignore /\[\t \f *\r?\n/    // LINE_CONT
%ignore COMMENT
"""

parser = Lark(TEMPLATED_PYTHON, parser='lalr', start=['single_input', 'file_input',

```

(continues on next page)

(continued from previous page)

```

↪ 'eval_input', 'template_start'], postlex=PythonIndenter(), maybe_placeholders=False)

def parse_template(s):
    return parser.parse(s + '\n', start='template_start')

def parse_code(s):
    return parser.parse(s + '\n', start='file_input')

#
# 2. Define translations using templates (each template code is parsed to a template_
↪ tree)
#

pytemplate = TemplateConf(parse=parse_template)

translations_3to2 = {
    'yield from $a':
        'for _tmp in $a: yield _tmp',

    'raise $e from $x':
        'raise $e',

    '$a / $b':
        'float($a) / $b',
}
translations_3to2 = {pytemplate(k): pytemplate(v) for k, v in translations_3to2.items()}

#
# 3. Translate and reconstruct Python 3 code into valid Python 2 code
#

python_reconstruct = PythonReconstructor(parser)

def translate_py3to2(code):
    tree = parse_code(code)
    tree = TemplateTranslator(translations_3to2).translate(tree)
    return python_reconstruct.reconstruct(tree)

#
# Test Code
#

_TEST_CODE = '''
if a / 2 > 1:
    yield from [1,2,3]
else:
    raise ValueError(a) from e
'''

```

(continues on next page)

(continued from previous page)

```
def test():
    print(_TEST_CODE)
    print('    ----->    ')
    print(translate_py3to2(_TEST_CODE))

if __name__ == '__main__':
    test()
```

Total running time of the script: (0 minutes 0.000 seconds)

11.1.11 Grammar-complete Python Parser

A fully-working Python 2 & 3 parser (but not production ready yet!)

This example demonstrates usage of the included Python grammars

```
import sys
import os, os.path
from io import open
import glob, time

from lark import Lark
from lark.indenter import PythonIndenter

kwargs = dict(postlex=PythonIndenter(), start='file_input')

# Official Python grammar by Lark
python_parser3 = Lark.open_from_package('lark', 'python.lark', ['grammars'], parser='lalr',
↳ ', **kwargs)

# Local Python2 grammar
python_parser2 = Lark.open('python2.lark', rel_to=__file__, parser='lalr', **kwargs)
python_parser2_earley = Lark.open('python2.lark', rel_to=__file__, parser='earley',
↳ lexer='basic', **kwargs)

try:
    xrange
except NameError:
    chosen_parser = python_parser3
else:
    chosen_parser = python_parser2

def _read(fn, *args):
    kwargs = {'encoding': 'iso-8859-1'}
    with open(fn, *args, **kwargs) as f:
        return f.read()

def _get_lib_path():
    if os.name == 'nt':
```

(continues on next page)

(continued from previous page)

```

    if 'PyPy' in sys.version:
        return os.path.join(sys.base_prefix, 'lib-python', sys.winver)
    else:
        return os.path.join(sys.base_prefix, 'Lib')
else:
    return [x for x in sys.path if x.endswith('%s.%s' % sys.version_info[:2])][0]

def test_python_lib():
    path = _get_lib_path()

    start = time.time()
    files = glob.glob(path+'/*.py')
    total_kb = 0
    for f in files:
        r = _read(os.path.join(path, f))
        kb = len(r) / 1024
        print( '%s -\t%.1f kb' % (f, kb))
        chosen_parser.parse(r + '\n')
        total_kb += kb

    end = time.time()
    print( "test_python_lib (%d files, %.1f kb), time: %.2f secs"%(len(files), total_kb,
→ end-start) )

def test_earley_equals_lalr():
    path = _get_lib_path()

    files = glob.glob(path+'/*.py')
    for f in files:
        print( f )
        tree1 = python_parser2.parse(_read(os.path.join(path, f)) + '\n')
        tree2 = python_parser2_earley.parse(_read(os.path.join(path, f)) + '\n')
        assert tree1 == tree2

if __name__ == '__main__':
    test_python_lib()
    # test_earley_equals_lalr()
    # python_parser3.parse(_read(sys.argv[1]) + '\n')

```

Total running time of the script: (0 minutes 0.000 seconds)

11.1.12 Creating an AST from the parse tree

This example demonstrates how to transform a parse-tree into an AST using *lark.ast_utils*.

`create_transformer()` collects every subclass of *Ast* subclass from the module, and creates a Lark transformer that builds the AST with no extra code.

This example only works with Python 3.

```
import sys
from typing import List
from dataclasses import dataclass

from lark import Lark, ast_utils, Transformer, v_args
from lark.tree import Meta

this_module = sys.modules[__name__]

#
#   Define AST
#
class _Ast(ast_utils.Ast):
    # This will be skipped by create_transformer(), because it starts with an underscore
    pass

class _Statement(_Ast):
    # This will be skipped by create_transformer(), because it starts with an underscore
    pass

@dataclass
class Value(_Ast, ast_utils.WithMeta):
    "Uses WithMeta to include line-number metadata in the meta attribute"
    meta: Meta
    value: object

@dataclass
class Name(_Ast):
    name: str

@dataclass
class CodeBlock(_Ast, ast_utils.AsList):
    # Corresponds to code_block in the grammar
    statements: List[_Statement]

@dataclass
class If(_Statement):
    cond: Value
    then: CodeBlock

@dataclass
class SetVar(_Statement):
    # Corresponds to set_var in the grammar
    name: str
    value: Value

@dataclass
class Print(_Statement):
    value: Value

class ToAst(Transformer):
```

(continues on next page)

(continued from previous page)

```

    # Define extra transformation functions, for rules that don't correspond to an AST.
    ↪class.

    def STRING(self, s):
        # Remove quotation marks
        return s[1:-1]

    def DEC_NUMBER(self, n):
        return int(n)

    @v_args(inline=True)
    def start(self, x):
        return x

#
# Define Parser
#

parser = Lark("""
    start: code_block

    code_block: statement+

    ?statement: if | set_var | print

    if: "if" value "{" code_block "}"
    set_var: NAME "=" value ";"
    print: "print" value ";"

    value: name | STRING | DEC_NUMBER
    name: NAME

    %import python (NAME, STRING, DEC_NUMBER)
    %import common.WS
    %ignore WS
    """,
    parser="lalr",
)

transformer = ast_utils.create_transformer(this_module, ToAst())

def parse(text):
    tree = parser.parse(text)
    return transformer.transform(tree)

#
# Test
#

if __name__ == '__main__':
    print(parse("""
        a = 1;

```

(continues on next page)

(continued from previous page)

```

        if a {
            print "a is 1";
            a = 2;
        }
    """)

```

Total running time of the script: (0 minutes 0.000 seconds)

11.1.13 Example-Driven Error Reporting

A demonstration of example-driven error reporting with the Earley parser (See also: `error_reporting_lalr.py`)

```

from lark import Lark, UnexpectedInput

from _json_parser import json_grammar  # Using the grammar from the json_parser example

json_parser = Lark(json_grammar)

class JsonSyntaxError(SyntaxError):
    def __str__(self):
        context, line, column = self.args
        return '%s at line %s, column %s.\n\n%s' % (self.label, line, column, context)

class JsonMissingValue(JsonSyntaxError):
    label = 'Missing Value'

class JsonMissingOpening(JsonSyntaxError):
    label = 'Missing Opening'

class JsonMissingClosing(JsonSyntaxError):
    label = 'Missing Closing'

class JsonMissingComma(JsonSyntaxError):
    label = 'Missing Comma'

class JsonTrailingComma(JsonSyntaxError):
    label = 'Trailing Comma'

def parse(json_text):
    try:
        j = json_parser.parse(json_text)
    except UnexpectedInput as u:
        exc_class = u.match_examples(json_parser.parse, {
            JsonMissingOpening: ['{"foo": ]}',
                                '{"foor": } }',
                                '{"foo": } }'],
            JsonMissingClosing: ['{"foo": [}',
                                '{',
                                '{"a": 1',
                                '[1'],

```

(continues on next page)

(continued from previous page)

```

        JsonMissingComma: ['[1 2]',
                           '[false 1]',
                           '["b" 1]',
                           '{"a":true 1:4}',
                           '{"a":1 1:4}',
                           '{"a":"b" 1:4}'],

        JsonTrailingComma: ['[,]',
                             '[1,]',
                             '[1,2,]',
                             '{"foo":1,}',
                             '{"foo":false,"bar":true,}']

    }, use_accepts=True)
    if not exc_class:
        raise
    raise exc_class(u.get_context(json_text), u.line, u.column)

def test():
    try:
        parse('{"example1": "value"')
    except JsonMissingClosing as e:
        print(e)

    try:
        parse('{"example2": ] ')
    except JsonMissingOpening as e:
        print(e)

if __name__ == '__main__':
    test()

```

Total running time of the script: (0 minutes 0.000 seconds)

11.1.14 Example-Driven Error Reporting

A demonstration of example-driven error reporting with the LALR parser (See also: `error_reporting_earley.py`)

```

from lark import Lark, UnexpectedInput

from _json_parser import json_grammar  # Using the grammar from the json_parser example

json_parser = Lark(json_grammar, parser='lalr')

class JsonSyntaxError(SyntaxError):
    def __str__(self):
        context, line, column = self.args
        return '%s at line %s, column %s.\n\n%s' % (self.label, line, column, context)

class JsonMissingValue(JsonSyntaxError):
    label = 'Missing Value'

```

(continues on next page)

(continued from previous page)

```

class JsonMissingOpening(JsonSyntaxError):
    label = 'Missing Opening'

class JsonMissingClosing(JsonSyntaxError):
    label = 'Missing Closing'

class JsonMissingComma(JsonSyntaxError):
    label = 'Missing Comma'

class JsonTrailingComma(JsonSyntaxError):
    label = 'Trailing Comma'

def parse(json_text):
    try:
        j = json_parser.parse(json_text)
    except UnexpectedInput as u:
        exc_class = u.match_examples(json_parser.parse, {
            JsonMissingOpening: ['{"foo": ]}',
                                '{"foor": } }',
                                '{"foo": } }'],
            JsonMissingClosing: ['{"foo": [}',
                                '{',
                                '{"a": 1',
                                '[1]',
            JsonMissingComma: ['[1 2]',
                              '[false 1]',
                              '["b" 1]',
                              '{"a":true 1:4}',
                              '{"a":1 1:4}',
                              '{"a":"b" 1:4}'],
            JsonTrailingComma: ['[,]',
                               '[1,]',
                               '[1,2,]',
                               '{"foo":1,}',
                               '{"foo":false,"bar":true,}'],
        }, use_accepts=True)
        if not exc_class:
            raise
        raise exc_class(u.get_context(json_text), u.line, u.column)

def test():
    try:
        parse('{"example1": "value"')
    except JsonMissingClosing as e:
        print(e)

    try:
        parse('{"example2": ] ')
    except JsonMissingOpening as e:

```

(continues on next page)

(continued from previous page)

```

print(e)

if __name__ == '__main__':
    test()

```

Total running time of the script: (0 minutes 0.000 seconds)

11.1.15 Reconstruct Python

Demonstrates how Lark's experimental text-reconstruction feature can recreate functional Python code from its parse-tree, using just the correct grammar and a small formatter.

```

from lark import Token, Lark
from lark.reconstruct import Reconstructor
from lark.indenter import PythonIndenter

# Official Python grammar by Lark
python_parser3 = Lark.open_from_package('lark', 'python.lark', ['grammars'],
                                         parser='lalr', postlex=PythonIndenter(), start=
↳ 'file_input',
                                         maybe_placeholders=False # Necessary for
↳ reconstructor
                                         )

SPACE_AFTER = set('+-*/~@<>="|:~')
SPACE_BEFORE = (SPACE_AFTER - set(':', '~')) | set('\n')

def special(sym):
    return Token('SPECIAL', sym.name)

def postproc(items):
    stack = ['\n']
    actions = []
    last_was_whitespace = True
    for item in items:
        if isinstance(item, Token) and item.type == 'SPECIAL':
            actions.append(item.value)
        else:
            if actions:
                assert actions[0] == '_NEWLINE' and '_NEWLINE' not in actions[1:],
↳ actions

                for a in actions[1:]:
                    if a == '_INDENT':
                        stack.append(stack[-1] + ' ' * 4)
                    else:
                        assert a == '_DEDENT'
                        stack.pop()
                actions.clear()

```

(continues on next page)

(continued from previous page)

```

        yield stack[-1]
        last_was_whitespace = True
    if not last_was_whitespace:
        if item[0] in SPACE_BEFORE:
            yield ' '
    yield item
    last_was_whitespace = item[-1].isspace()
    if not last_was_whitespace:
        if item[-1] in SPACE_AFTER:
            yield ' '
        last_was_whitespace = True
yield "\n"

class PythonReconstructor:
    def __init__(self, parser):
        self._recons = Reconstructor(parser, {'_NEWLINE': special, '_DEDENT': special, '_
↳ INDENT': special})

    def reconstruct(self, tree):
        return self._recons.reconstruct(tree, postproc)

def test():
    python_reconstructor = PythonReconstructor(python_parser3)

    self_contents = open(__file__).read()

    tree = python_parser3.parse(self_contents+'\n')
    output = python_reconstructor.reconstruct(tree)

    tree_new = python_parser3.parse(output)
    print(tree.pretty())
    print(tree_new.pretty())
    # assert tree.pretty() == tree_new.pretty()
    assert tree == tree_new

    print(output)

if __name__ == '__main__':
    test()

```

Total running time of the script: (0 minutes 0.000 seconds)

11.1.16 Using lexer `dynamic_complete`

Demonstrates how to use `lexer='dynamic_complete'` and `ambiguity='explicit'`

Sometimes you have data that is highly ambiguous or ‘broken’ in some sense. When using `parser='earley'` and `lexer='dynamic_complete'`, Lark will be able parse just about anything as long as there is a valid way to generate it from the Grammar, including looking ‘into’ the Regexes.

This examples shows how to parse a json input where the quotes have been replaced by underscores: `{_foo_: {}, _bar_: [], _baz_: __}` Notice that underscores might still appear inside strings, so a potentially valid reading of the above is: `{"foo_:{}, _bar": [], "baz": ""}`

```
from pprint import pprint

from lark import Lark, Tree, Transformer, v_args
from lark.visitors import Transformer_InPlace

GRAMMAR = r"""
%import common.SIGNED_NUMBER
%import common.WS_INLINE
%import common.NEWLINE
%ignore WS_INLINE

?start: value

?value: object
      | array
      | string
      | SIGNED_NUMBER      -> number
      | "true"             -> true
      | "false"            -> false
      | "null"             -> null

array  : "[" (value ("," value)*)? "]"
object : "{" (pair ("," pair)*)? "}"
pair   : string ":" value

string: STRING
STRING : ESCAPED_STRING

ESCAPED_STRING: QUOTE_CHAR _STRING_ESC_INNER QUOTE_CHAR
QUOTE_CHAR: "\""

_STRING_INNER: /.*/
_STRING_ESC_INNER: _STRING_INNER /(?!\\)(\\\\\\\\)*?/

"""

def score(tree: Tree):
    """
    Scores an option by how many children (and grand-children, and
    grand-grand-children, ...) it has.
    This means that the option with fewer large terminals gets selected
    """
```

(continues on next page)

(continued from previous page)

```

Between
    object
        pair
            string _foo_
            object
        pair
            string _bar_: [], _baz_
            string __

and

    object
        pair
            string _foo_
            object
        pair
            string _bar_
            array
        pair
            string _baz_
            string __

this will give the second a higher score. (9 vs 13)
"""
return sum(len(t.children) for t in tree.iter_subtrees())

class RemoveAmbiguities(Transformer_InPlace):
    """
    Selects an option to resolve an ambiguity using the score function above.
    Scores each option and selects the one with the higher score, e.g. the one
    with more nodes.

    If there is a performance problem with the Tree having too many _ambig and
    being slow and too large, this can instead be written as a ForestVisitor.
    Look at the 'Custom SPPF Prioritizer' example.
    """
    def _ambig(self, options):
        return max(options, key=score)

class TreeToJson(Transformer):
    """
    This is the same Transformer as the json_parser example.
    """
    @v_args(inline=True)
    def string(self, s):
        return s[1:-1].replace('\\"', '"')

array = list
pair = tuple

```

(continues on next page)

(continued from previous page)

```

object = dict
number = v_args(inline=True)(float)

null = lambda self, _: None
true = lambda self, _: True
false = lambda self, _: False

parser = Lark(GRAMMAR, parser='earley', ambiguity="explicit", lexer='dynamic_complete')

EXAMPLES = [
    r'[_array_: [1,2,3]]',

    r'[_abc_: _array must be of the following format [_1_, _2_, _3_]_]',

    r'[_foo_: {}, _bar_: [], _baz_: __]_',

    r'[_error_: _invalid_client_, _error_description_: AADSTS7000215: Invalid '
    r'client secret is provided.\r\nTrace ID: '
    r'a0a0aaaa-a0a0-0a00-000a-00a00aaa0a00\r\nCorrelation ID: '
    r'aa0aaa00-0aaa-0000-00a0-00000aaaa0aa\r\nTimestamp: 1997-10-10 00:00:00Z_, '
    r'_error_codes_: [7000215], _timestamp_: 1997-10-10 00:00:00Z_, '
    r'_trace_id_: a0a0aaaa-a0a0-0a00-000a-00a00aaa0a00_, '
    r'_correlation_id_: aa0aaa00-0aaa-0000-00a0-00000aaaa0aa_, '
    r'_error_uri_: https://example.com_]_',

]
for example in EXAMPLES:
    tree = parser.parse(example)
    tree = RemoveAmbiguities().transform(tree)
    result = TreeToJson().transform(tree)
    pprint(result)

```

Total running time of the script: (0 minutes 0.000 seconds)

11.1.17 Syntax Highlighting

This example shows how to write a syntax-highlighted editor with Qt and Lark

Requirements:

PyQt5==5.15.8 QScintilla==2.13.4

```

import sys
import textwrap

from PyQt5.QtWidgets import QApplication
from PyQt5.QtGui import QColor, QFont, QFontMetrics

from PyQt5.Qsci import QsciScintilla
from PyQt5.Qsci import QsciLexerCustom

```

(continues on next page)

(continued from previous page)

```
from lark import Lark

class LexerJson(QsciLexerCustom):

    def __init__(self, parent=None):
        super().__init__(parent)
        self.create_parser()
        self.create_styles()

    def create_styles(self):
        deeppink = QColor(249, 38, 114)
        khaki = QColor(230, 219, 116)
        mediumpurple = QColor(174, 129, 255)
        medianturquoise = QColor(81, 217, 205)
        yellowgreen = QColor(166, 226, 46)
        lightcyan = QColor(213, 248, 232)
        darkslategrey = QColor(39, 40, 34)

        styles = {
            0: medianturquoise,
            1: mediumpurple,
            2: yellowgreen,
            3: deeppink,
            4: khaki,
            5: lightcyan
        }

        for style, color in styles.items():
            self.setColor(color, style)
            self.setPaper(darkslategrey, style)
            self.setFont(self.parent().font(), style)

        self.token_styles = {
            "COLON": 5,
            "COMMA": 5,
            "LBRACE": 5,
            "LSQB": 5,
            "RBRACE": 5,
            "RSQB": 5,
            "FALSE": 0,
            "NULL": 0,
            "TRUE": 0,
            "STRING": 4,
            "NUMBER": 1,
        }

    def create_parser(self):
        grammar = '''
        anons: ":" "{" "}" " ", " "[" "]"
        TRUE: "true"
        FALSE: "false"'''
```

(continues on next page)

(continued from previous page)

```

        NULL: "NULL"
        %import common.ESCAPED_STRING -> STRING
        %import common.SIGNED_NUMBER -> NUMBER
        %import common.WS
        %ignore WS
    ...

    self.lark = Lark(grammar, parser=None, lexer='basic')
    # All tokens: print([t.name for t in self.lark.parser.lexer.tokens])

def defaultPaper(self, style):
    return QColor(39, 40, 34)

def language(self):
    return "Json"

def description(self, style):
    return {v: k for k, v in self.token_styles.items()}.get(style, "")

def styleText(self, start, end):
    self.startStyling(start)
    text = self.parent().text()[start:end]
    last_pos = 0

    try:
        for token in self.lark.lex(text):
            ws_len = token.start_pos - last_pos
            if ws_len:
                self.setStyling(ws_len, 0)    # whitespace

            token_len = len(bytearray(token, "utf-8"))
            self.setStyling(
                token_len, self.token_styles.get(token.type, 0))

            last_pos = token.start_pos + token_len
    except Exception as e:
        print(e)

class EditorAll(QsciScintilla):

    def __init__(self, parent=None):
        super().__init__(parent)

        # Set font defaults
        font = QFont()
        font.setFamily('Consolas')
        font.setFixedPitch(True)
        font.setPointSize(8)
        font.setBold(True)
        self.setFont(font)

```

(continues on next page)

(continued from previous page)

```

# Set margin defaults
fontmetrics = QFontMetrics(font)
self.setMarginsFont(font)
self.setMarginWidth(0, fontmetrics.width("000") + 6)
self.setMarginLineNumbers(0, True)
self.setMarginsForegroundColor(QColor(128, 128, 128))
self.setMarginsBackgroundColor(QColor(39, 40, 34))
self.setMarginType(1, self.SymbolMargin)
self.setMarginWidth(1, 12)

# Set indentation defaults
self.setIndentationsUseTabs(False)
self.setIndentationWidth(4)
self.setBackspaceUnindents(True)
self.setIndentationGuides(True)

# self.setFolding(QsciScintilla.CircledFoldStyle)

# Set caret defaults
self.setCaretForegroundColor(QColor(247, 247, 241))
self.setCaretWidth(2)

# Set selection color defaults
self.setSelectionBackgroundColor(QColor(61, 61, 52))
self.resetSelectionForegroundColor()

# Set multiselection defaults
self.SendScintilla(QsciScintilla.SCI_SETMULTIPLESELECTION, True)
self.SendScintilla(QsciScintilla.SCI_SETMULTIPASTE, 1)
self.SendScintilla(
    QsciScintilla.SCI_SETADDITIONALSELECTIONTYPING, True)

lexer = LexerJson(self)
self.setLexer(lexer)

```

```

EXAMPLE_TEXT = textwrap.dedent("""\
{
    "_id": "5b05ffcbcf8e597939b3f5ca",
    "about": "Excepteur consequat commodo esse voluptate aute aliquip ad sint.
↳deserunt commodo eiusmod irure. Sint aliquip sit magna duis eu est culpa aliqua.
↳excepteur ut tempor nulla. Aliqua ex pariatur id labore sit. Quis sit ex aliqua veniam.
↳exercitation laboris anim adipisicing. Lorem nisi reprehenderit ullamco labore qui sit.
↳ut aliqua tempor consequat pariatur proident.",
    "address": "665 Malbone Street, Thornport, Louisiana, 243",
    "age": 23,
    "balance": "$3,216.91",
    "company": "BULLJUICE",
    "email": "elisekelley@bulljuice.com",
    "eyeColor": "brown",
    "gender": "female",
    "guid": "d3a6d865-0f64-4042-8a78-4f53de9b0707",

```

(continues on next page)

(continued from previous page)

```

        "index": 0,
        "isActive": false,
        "isActive2": true,
        "latitude": -18.660714,
        "longitude": -85.378048,
        "name": "Elise Kelley",
        "phone": "+1 (808) 543-3966",
        "picture": "http://placeholder.it/32x32",
        "registered": "2017-09-30T03:47:40 -02:00",
        "tags": [
            "et",
            "nostrud",
            "in",
            "fugiat",
            "incidunt",
            "labore",
            "nostrud"
        ]
    } \
"""

def main():
    app = QApplication(sys.argv)
    ex = EditorAll()
    ex.setWindowTitle(__file__)
    ex.setText(EXAMPLE_TEXT)
    ex.resize(800, 600)
    ex.show()
    sys.exit(app.exec_())

if __name__ == "__main__":
    main()

```

Total running time of the script: (0 minutes 0.000 seconds)

11.2 Grammar Composition

This example shows how to do grammar composition in Lark, by creating a new file format that allows both CSV and JSON to co-exist.

We show how, by using namespaces, Lark grammars and their transformers can be fully reused - they don't need to care if their grammar is used directly, or being imported, or who is doing the importing.

See [main.py](#) for more details.

Transformer for evaluating json.lark

```

from lark import Transformer, v_args

class JsonTreeToJson(Transformer):
    @v_args(inline=True)

```

(continues on next page)

(continued from previous page)

```
def string(self, s):
    return s[1:-1].replace('\\"', '')

array = list
pair = tuple
object = dict
number = v_args(inline=True)(float)

null = lambda self, _: None
true = lambda self, _: True
false = lambda self, _: False
```

Total running time of the script: (0 minutes 0.000 seconds)

Transformer for evaluating csv.lark

```
from lark import Transformer

class CsvTreeToPandasDict(Transformer):
    INT = int
    FLOAT = float
    SIGNED_FLOAT = float
    WORD = str
    NON_SEPARATOR_STRING = str

    def row(self, children):
        return children

    def start(self, children):
        data = {}

        header = children[0].children
        for heading in header:
            data[heading] = []

        for row in children[1:]:
            for i, element in enumerate(row):
                data[header[i]].append(element)

        return data
```

Total running time of the script: (0 minutes 0.000 seconds)

11.2.1 Grammar Composition

This example shows how to do grammar composition in Lark, by creating a new file format that allows both CSV and JSON to co-exist.

- 1) We define `storage.lark`, which imports both `csv.lark` and `json.lark`,

and allows them to be used one after the other.

In the generated tree, each imported rule/terminal is automatically prefixed (with `json__` or `csv__`), which creates an implicit namespace and allows them to coexist without collisions.

- 2) We merge their respective transformers (unaware of each other) into a new base transformer. The resulting transformer can evaluate both JSON and CSV in the parse tree.

The methods of each transformer are renamed into their appropriate namespace, using the given prefix. This approach allows full re-use: the transformers don't need to care if their grammar is used directly, or being imported, or who is doing the importing.

```
from pathlib import Path
from lark import Lark
from json import dumps
from lark.visitors import Transformer, merge_transformers

from eval_csv import CsvTreeToPandasDict
from eval_json import JsonTreeToJson

__dir__ = Path(__file__).parent

class Storage(Transformer):
    def start(self, children):
        return children

storage_transformer = merge_transformers(Storage(), csv=CsvTreeToPandasDict(),
↪ json=JsonTreeToJson())

parser = Lark.open("storage.lark", rel_to=__file__)

def main():
    json_tree = parser.parse(dumps({"test": "a", "dict": { "list": [1, 1.2] }}))
    res = storage_transformer.transform(json_tree)
    print("Just JSON: ", res)

    csv_json_tree = parser.parse(open(__dir__ / 'combined_csv_and_json.txt').read())
    res = storage_transformer.transform(csv_json_tree)
    print("JSON + CSV: ", dumps(res, indent=2))

if __name__ == "__main__":
    main()
```

Total running time of the script: (0 minutes 0.000 seconds)

11.3 Example Grammars

This directory is a collection of lark grammars, taken from real world projects.

- Verilog - Taken from <https://github.com/circuitgraph/circuitgraph/blob/main/circuitgraph/parsing/verilog.lark>

11.4 Standalone example

To initialize, cd to this folder, and run:

```
./create_standalone.sh
```

Or:

```
python -m lark.tools.standalone json.lark > json_parser.py
```

Then run using:

```
python json_parser_main.py <path-to.json>
```

11.4.1 Standalone Parser

This example demonstrates how to generate and use the standalone parser, using the JSON example.

See README.md for more details.

```
import sys

from json_parser import Lark_StandAlone, Transformer, v_args

inline_args = v_args(inline=True)

class TreeToJson(Transformer):
    @inline_args
    def string(self, s):
        return s[1:-1].replace('\\"', '"')

    array = list
    pair = tuple
    object = dict
    number = inline_args(float)

    null = lambda self, _: None
    true = lambda self, _: True
    false = lambda self, _: False

parser = Lark_StandAlone(transformer=TreeToJson())

if __name__ == '__main__':
    with open(sys.argv[1]) as f:
        print(parser.parse(f.read()))
```

Total running time of the script: (0 minutes 0.000 seconds)

GRAMMAR REFERENCE

12.1 Definitions

A **grammar** is a list of rules and terminals, that together define a language.

Terminals define the alphabet of the language, while rules define its structure.

In Lark, a terminal may be a string, a regular expression, or a concatenation of these and other terminals.

Each rule is a list of terminals and rules, whose location and nesting define the structure of the resulting parse-tree.

A **parsing algorithm** is an algorithm that takes a grammar definition and a sequence of symbols (members of the alphabet), and matches the entirety of the sequence by searching for a structure that is allowed by the grammar.

12.1.1 General Syntax and notes

Grammars in Lark are based on [EBNF](#) syntax, with several enhancements.

EBNF is basically a short-hand for common BNF patterns.

Optionals are expanded:

<code>a b? c</code> <code>-></code> <code>(a c a b c)</code>

Repetition is extracted into a recursion:

<code>a: b*</code> <code>-></code> <code>a: _b_tag</code> <code>_b_tag: (_b_tag b)?</code>

And so on.

Lark grammars are composed of a list of definitions and directives, each on its own line. A definition is either a named rule, or a named terminal, with the following syntax, respectively:

<code>rule: <EBNF EXPRESSION></code> <code> etc.</code>
<code>TERM: <EBNF EXPRESSION> <i>// Rules aren't allowed</i></code>

Comments start with either `//` (C++ style) or `#` (Python style, since version 1.1.6) and last to the end of the line.

Lark begins the parse with the rule 'start', unless specified otherwise in the options.

Names of rules are always in lowercase, while names of terminals are always in uppercase. This distinction has practical effects, for the shape of the generated parse-tree, and the automatic construction of the lexer (aka tokenizer, or scanner).

12.2 Terminals

Terminals are used to match text into symbols. They can be defined as a combination of literals and other terminals.

Syntax:

```
<NAME> [. <priority>] : <literals-and-or-terminals>
```

Terminal names must be uppercase.

Literals can be one of:

- "string"
- /regular expression+/
 - "case-insensitive string"i
 - /re with flags/imulx
- Literal range: "a".."z", "1".."9", etc.

Terminals also support grammar operators, such as |, +, * and ?.

Terminals are a linear construct, and therefore may not contain themselves (recursion isn't allowed).

12.2.1 Templates

Templates are expanded when preprocessing the grammar.

Definition syntax:

```
my_template{param1, param2, ...}: <EBNF EXPRESSION>
```

Use syntax:

```
some_rule: my_template{arg1, arg2, ...}
```

Example:

```
_separated{x, sep}: x (sep x)* // Define a sequence of 'x sep x sep x ...'  
num_list: "[" _separated{NUMBER, ","} "]" // Will match "[1, 2, 3]" etc.
```

12.2.2 Priority

Terminals can be assigned a priority to influence lexing. Terminal priorities are signed integers with a default value of 0.

When using a lexer, the highest priority terminals are always matched first.

When using Earley's dynamic lexing, terminal priorities are used to prefer certain lexings and resolve ambiguity.

12.2.3 Regexp Flags

You can use flags on regexps and strings. For example:

```
SELECT: "select"i      /// Will ignore case, and match SELECT or Select, etc.
MULTILINE_TEXT: /.+/s
SIGNED_INTEGER: /
    [+]? # the sign
    (0|[1-9][0-9]*) # the digits
/x
```

Supported flags are one of: `imslux`. See Python's regex documentation for more details on each one.

Regexps/strings of different flags can only be concatenated in Python 3.6+

Notes for when using a lexer:

When using a lexer (basic or contextual), it is the grammar-author's responsibility to make sure the literals don't collide, or that if they do, they are matched in the desired order. Literals are matched according to the following precedence:

1. Highest priority first (priority is specified as: `TERM.number: ...`)
2. Length of match (for regexps, the longest theoretical match is used)
3. Length of literal / pattern definition
4. Name

Examples:

```
IF: "if"
INTEGER : /[0-9]+/
INTEGER2 : ("0".. "9")+      /// Same as INTEGER
DECIMAL.2: INTEGER? "." INTEGER /// Will be matched before INTEGER
WHITESPACE: (" " | /\t/ )+
SQL_SELECT: "select"i
```

12.2.4 Regular expressions & Ambiguity

Each terminal is eventually compiled to a regular expression. All the operators and references inside it are mapped to their respective expressions.

For example, in the following grammar, A1 and A2, are equivalent:

```
A1: "a" | "b"
A2: /a|b/
```

This means that inside terminals, Lark cannot detect or resolve ambiguity, even when using Earley.

For example, for this grammar:

```
start      : (A | B)+
A          : "a" | "ab"
B          : "b"
```

We get only one possible derivation, instead of two:

```
>>> p = Lark(g, ambiguity="explicit")
>>> p.parse("ab")
Tree('start', [Token('A', 'ab')])
```

This is happening because Python's regex engine always returns the best matching option. There is no way to access the alternatives.

If you find yourself in this situation, the recommended solution is to use rules instead.

Example:

```
>>> p = Lark("""start: (a | b)+
...           !a: "a" | "ab"
...           !b: "b"
...           """, ambiguity="explicit")
>>> print(p.parse("ab").pretty())
_ambig
  start
    a  ab
  start
    a  a
    b  b
```

12.3 Rules

Syntax:

```
<name> : <items-to-match> [-> <alias> ]
      | ...
```

Names of rules and aliases are always in lowercase.

Rule definitions can be extended to the next line by using the OR operator (signified by a pipe: |).

An alias is a name for the specific rule alternative. It affects tree construction.

Each item is one of:

- rule
- TERMINAL
- "string literal" or /regexp literal/
- (item item ..) - Group items
- [item item ..] - Maybe. Same as (item item ..)?, but when maybe_placeholders=True, generates None if there is no match.
- item? - Zero or one instances of item ("maybe")
- item* - Zero or more instances of item
- item+ - One or more instances of item
- item ~ n - Exactly *n* instances of item
- item ~ n..m - Between *n* to *m* instances of item (not recommended for wide ranges, due to performance issues)

Examples:

```
hello_world: "hello" "world"
mul: (mul "*" )? number    ///  
Left-recursion is allowed and encouraged!
expr: expr operator expr
      | value               ///  
Multi-line, belongs to expr
four_words: word ~ 4
```

12.3.1 Priority

Like terminals, rules can be assigned a priority. Rule priorities are signed integers with a default value of 0.

When using LALR, the highest priority rules are used to resolve collision errors.

When using Earley, rule priorities are used to resolve ambiguity.

12.4 Directives

12.4.1 %ignore

All occurrences of the terminal will be ignored, and won't be part of the parse.

Using the %ignore directive results in a cleaner grammar.

It's especially important for the LALR(1) algorithm, because adding whitespace (or comments, or other extraneous elements) explicitly in the grammar, harms its predictive abilities, which are based on a lookahead of 1.

Syntax:

```
%ignore <TERMINAL>
```

Examples:

```
%ignore " "
COMMENT: "#" /[^\n]/*
%ignore COMMENT
```

12.4.2 %import

Allows one to import terminals and rules from lark grammars.

When importing rules, all their dependencies will be imported into a namespace, to avoid collisions. It's not possible to override their dependencies (e.g. like you would when inheriting a class).

Syntax:

```
%import <module>.<TERMINAL>
%import <module>.<rule>
%import <module>.<TERMINAL> -> <NEWTERMINAL>
%import <module>.<rule> -> <newrule>
%import <module> (<TERM1>, <TERM2>, <rule1>, <rule2>)
```

If the module path is absolute, Lark will attempt to load it from the built-in directory (which currently contains `common.lark`, `python.lark`, and `unicode.lark`).

If the module path is relative, such as `.path.to.file`, Lark will attempt to load it from the current working directory. Grammars must have the `.lark` extension.

The rule or terminal can be imported under another name with the `->` syntax.

Example:

```
%import common.NUMBER

%import .terminals_file (A, B, C)

%import .rules_file.rulea -> ruleb
```

Note that `%ignore` directives cannot be imported. Imported rules will abide by the `%ignore` directives declared in the main grammar.

12.4.3 %declare

Declare a terminal without defining it. Useful for plugins.

12.4.4 %override

Override a rule or terminals, affecting all references to it, even in imported grammars.

Useful for implementing an inheritance pattern when importing grammars.

Example:

```
%import my_grammar (start, number, NUMBER)

// Add hex support to my_grammar
%override number: NUMBER | /0x\w+/
```

12.4.5 %extend

Extend the definition of a rule or terminal, e.g. add a new option on what it can match, like when separated with `|`.

Useful for splitting up a definition of a complex rule with many different options over multiple files.

Can also be used to implement a plugin system where a core grammar is extended by others.

Example:

```
%import my_grammar (start, NUMBER)

// Add hex support to my_grammar
%extend NUMBER: /0x\w+/
```

For both `%extend` and `%override`, there is not requirement for a rule/terminal to come from another file, but that is probably the most common usecase

TREE CONSTRUCTION REFERENCE

Lark builds a tree automatically based on the structure of the grammar, where each rule that is matched becomes a branch (node) in the tree, and its children are its matches, in the order of matching.

For example, the rule `node: child1 child2` will create a tree node with two children. If it is matched as part of another rule (i.e. if it isn't the root), the new rule's tree node will become its parent.

Using `item+` or `item*` will result in a list of items, equivalent to writing `item item item ...`

Using `item?` will return the item if it matched, or nothing.

If `maybe_placeholders=True` (the default), then using `[item]` will return the item if it matched, or the value `None`, if it didn't.

If `maybe_placeholders=False`, then `[]` behaves like `()?`.

13.1 Terminals

Terminals are always values in the tree, never branches.

Lark filters out certain types of terminals by default, considering them punctuation:

- Terminals that won't appear in the tree are:
 - Unnamed literals (like `"keyword"` or `"+"`)
 - Terminals whose name starts with an underscore (like `_DIGIT`)
- Terminals that *will* appear in the tree are:
 - Unnamed regular expressions (like `/[0-9]/`)
 - Named terminals whose name starts with a letter (like `DIGIT`)

Note: Terminals composed of literals and other terminals always include the entire match without filtering any part.

Example:

```
start:  PNAME pname

PNAME: "(" NAME ")"
pname: "(" NAME ")"

NAME:  /\w+/
%ignore /\s+/
```

Lark will parse `"(Hello) (World)"` as:

```
start
  (Hello)
  pname World
```

Rules prefixed with ! will retain all their literals regardless.

Example:

```
expr: "(" expr ")"
      | NAME+

NAME: /\w+/

%ignore " "
```

Lark will parse “((hello world))” as:

```
expr
  expr
    expr
      "hello"
      "world"
```

The brackets do not appear in the tree by design. The words appear because they are matched by a named terminal.

13.2 Shaping the tree

Users can alter the automatic construction of the tree using a collection of grammar features.

- Rules whose name begins with an underscore will be inlined into their containing rule.

Example:

```
start: "(" _greet ")"
_greet: /\w+/ /\w+/
```

Lark will parse “(hello world)” as:

```
start
  "hello"
  "world"
```

- Rules that receive a question mark (?) at the beginning of their definition, will be inlined if they have a single child, after filtering.

Example:

```
start: greet greet
?greet: "(" /\w+/ ")"
        | /\w+/ /\w+/
```

Lark will parse “hello world (planet)” as:


```
start
  greet
    "hello"
    "world"
  "planet"
```

- Rules that begin with an exclamation mark will keep all their terminals (they won't get filtered).

```
!expr: "(" expr ")"
      | NAME+
NAME: /\w+/
%ignore " "
```

Will parse “((hello world))” as:

```
expr
(
  expr
  (
    expr
    hello
    world
  )
)
```

Using the ! prefix is usually a “code smell”, and may point to a flaw in your grammar design.

- Aliases - options in a rule can receive an alias. It will be then used as the branch name for the option, instead of the rule name.

Example:

```
start: greet greet
greet: "hello"
      | "world" -> planet
```

Lark will parse “hello world” as:

```
start
  greet
  planet
```


API REFERENCE

14.1 Lark

class `lark.Lark(grammar: Union[Grammar, str, IO[str]], **options)`

Main interface for the library.

It's mostly a thin wrapper for the many different parsers, and for the tree constructor.

Parameters

- **grammar** – a string or file-object containing the grammar spec (using Lark's ebnf syntax)
- **options** – a dictionary controlling various aspects of Lark.

Example

```
>>> Lark(r'''start: "foo" ''')
Lark(...)
```

=== General Options ===

start

The start symbol. Either a string, or a list of strings for multiple possible starts (Default: “start”)

debug

Display debug information and extra warnings. Use only when debugging (Default: `False`) When used with Earley, it generates a forest graph as “sppf.png”, if ‘dot’ is installed.

strict

Throw an exception on any potential ambiguity, including shift/reduce conflicts, and regex collisions.

transformer

Applies the transformer to every parse tree (equivalent to applying it after the parse, but faster)

propagate_positions

Propagates positional attributes into the ‘meta’ attribute of all tree branches. Sets attributes: (line, column, end_line, end_column, start_pos, end_pos,

container_line, container_column, container_end_line, container_end_column)

Accepts `False`, `True`, or a callable, which will filter which nodes to ignore when propagating.

maybe_placeholders

When `True`, the `[]` operator returns `None` when not matched. When `False`, `[]` behaves like the `?` operator, and returns no value at all. (default= `True`)

cache

Cache the results of the Lark grammar analysis, for x2 to x3 faster loading. LALR only for now.

- When `False`, does nothing (default)
- When `True`, caches to a temporary file in the local directory
- When given a string, caches to the path pointed by the string

regex

When `True`, uses the `regex` module instead of the stdlib `re`.

g_regex_flags

Flags that are applied to all terminals (both `regex` and strings)

keep_all_tokens

Prevent the tree builder from automagically removing “punctuation” tokens (Default: `False`)

tree_class

Lark will produce trees comprised of instances of this class instead of the default `lark.Tree`.

=== Algorithm Options ===**parser**

Decides which parser engine to use. Accepts “`earley`” or “`lalr`”. (Default: “`earley`”). (there is also a “`cyk`” option for legacy)

lexer

Decides whether or not to use a lexer stage

- “`auto`” (default): Choose for me based on the parser
- “`basic`”: Use a basic lexer
- “`contextual`”: Stronger lexer (only works with `parser=“lalr”`)
- “`dynamic`”: Flexible and powerful (only with `parser=“earley”`)
- “`dynamic_complete`”: Same as `dynamic`, but tries *every* variation of tokenizing possible.

ambiguity

Decides how to handle ambiguity in the parse. Only relevant if `parser=“earley”`

- “`resolve`”: The parser will automatically choose the simplest derivation (it chooses consistently: greedy for tokens, non-greedy for rules)
- “`explicit`”: The parser will return all derivations wrapped in “`_ambig`” tree nodes (i.e. a forest).
- “`forest`”: The parser will return the root of the shared packed parse forest.

=== Misc. / Domain Specific Options ===**postlex**

Lexer post-processing (Default: `None`) Only works with the basic and contextual lexers.

priority

How priorities should be evaluated - “`auto`”, `None`, “`normal`”, “`invert`” (Default: “`auto`”)

lexer_callbacks

Dictionary of callbacks for the lexer. May alter tokens during lexing. Use with caution.

use_bytes

Accept an input of type `bytes` instead of `str`.

ordered_sets

Should Earley use ordered-sets to achieve stable output (~10% slower than regular sets. Default: `True`)

edit_terminals

A callback for editing the terminals before parse.

import_paths

A List of either paths or loader functions to specify from where grammars are imported

source_path

Override the source of from where the grammar was loaded. Useful for relative imports and unconventional grammar loading

=== End of Options ===

save(*f*, *exclude_options*: *Collection[str]* = ()) → None

Saves the instance into the given file object

Useful for caching and multiprocessing.

classmethod load(*f*) → *_T*

Loads an instance from the given file object

Useful for caching and multiprocessing.

classmethod open(*grammar_filename*: *str*, *rel_to*: *Optional[str]* = None, ***options*) → *_T*

Create an instance of Lark with the grammar given by its filename

If *rel_to* is provided, the function will find the grammar filename in relation to it.

Example

```
>>> Lark.open("grammar_file.lark", rel_to=__file__, parser="lalr")
Lark(...)
```

classmethod open_from_package(*package*: *str*, *grammar_path*: *str*, *search_paths*: *Sequence[str]* = [], ***options*) → *_T*

Create an instance of Lark with the grammar loaded from within the package *package*. This allows grammar loading from zipapps.

Imports in the grammar will use the *package* and *search_paths* provided, through *FromPackageLoader*

Example

```
Lark.open_from_package(__name__, "example.lark", ("grammars",), parser=...)
```

lex(*text*: *str*, *dont_ignore*: *bool* = False) → *Iterator[Token]*

Only lex (and postlex) the text, without parsing it. Only relevant when *lexer*='basic'

When *dont_ignore*=True, the lexer will return all tokens, even those marked for %ignore.

Raises

UnexpectedCharacters – In case the lexer cannot find a suitable match.

get_terminal(*name*: *str*) → *TerminalDef*

Get information about a terminal

parse_interactive(*text*: *Optional[str]* = None, *start*: *Optional[str]* = None) → *InteractiveParser*

Start an interactive parsing session.

Parameters

- **text** (*str*, *optional*) – Text to be parsed. Required for `resume_parse()`.
- **start** (*str*, *optional*) – Start symbol

Returns

A new `InteractiveParser` instance.

See Also: `Lark.parse()`

parse(*text: str*, *start: Optional[str] = None*, *on_error: Optional[Callable[[UnexpectedInput], bool]] = None*) → `ParseTree`

Parse the given text, according to the options provided.

Parameters

- **text** (*str*) – Text to be parsed.
- **start** (*str*, *optional*) – Required if Lark was given multiple possible start symbols (using the start option).
- **on_error** (*function*, *optional*) – if provided, will be called on `UnexpectedToken` error. Return true to resume parsing. LALR only. See `examples/advanced/error_handling.py` for an example of how to use `on_error`.

Returns

If a transformer is supplied to `__init__`, returns whatever is the result of the transformation. Otherwise, returns a `Tree` instance.

Raises

UnexpectedInput – On a parse error, one of these sub-exceptions will rise: `UnexpectedCharacters`, `UnexpectedToken`, or `UnexpectedEOF`. For convenience, these sub-exceptions also inherit from `ParserError` and `LexerError`.

14.1.1 Using Unicode character classes with regex

Python's builtin `re` module has a few persistent known bugs and also won't parse advanced regex features such as character classes. With `pip install lark[regex]`, the `regex` module will be installed alongside `lark` and can act as a drop-in replacement to `re`.

Any instance of `Lark` instantiated with `regex=True` will use the `regex` module instead of `re`.

For example, we can use character classes to match PEP-3131 compliant Python identifiers:

```
from lark import Lark
>>> g = Lark(r"""
    ?start: NAME
    NAME: ID_START ID_CONTINUE*
    ID_START: /\p{Lu}\p{Ll}\p{Lt}\p{Lm}\p{Lo}\p{Nl}_]/+
    ID_CONTINUE: ID_START | /\p{Mn}\p{Mc}\p{Nd}\p{Pc}\p{C}\p{Z}]/+
    """, regex=True)

>>> g.parse('')
```

14.2 Tree

class `lark.Tree`(*data: str, children: List[Union[_Leaf_T, Tree[_Leaf_T]]], meta: Optional[Meta] = None*)

The main tree class.

Creates a new tree, and stores “data” and “children” in attributes of the same name. Trees can be hashed and compared.

Parameters

- **data** – The name of the rule or alias
- **children** – List of matched sub-rules and terminals
- **meta** – Line & Column numbers (if `propagate_positions` is enabled). meta attributes: (line, column, end_line, end_column, start_pos, end_pos, container_line, container_column, container_end_line, container_end_column)
`container_*` attributes consider all symbols, including those that have been inlined in the tree. For example, in the rule ‘a: _A B _C’, the regular attributes will mark the start and end of B, but the `container_*` attributes will also include _A and _C in the range. However, rules that contain ‘a’ will consider it in full, including _A and _C for all attributes.

pretty(*indent_str: str = ' '*) → str

Returns an indented string representation of the tree.

Great for debugging.

__rich__(*parent: Optional[rich.tree.Tree] = None*) → rich.tree.Tree

Returns a tree widget for the ‘rich’ library.

Example

```
::
from rich import print from lark import Tree
tree = Tree('root', ['node1', 'node2']) print(tree)
```

iter_subtrees() → Iterator[Tree[_Leaf_T]]

Depth-first iteration.

Iterates over all the subtrees, never returning to the same node twice (Lark’s parse-tree is actually a DAG).

iter_subtrees_topdown()

Breadth-first iteration.

Iterates over all the subtrees, return nodes in order like `pretty()` does.

find_pred(*pred: Callable[[Tree[_Leaf_T]], bool]*) → Iterator[Tree[_Leaf_T]]

Returns all nodes of the tree that evaluate `pred(node)` as true.

find_data(*data: str*) → Iterator[Tree[_Leaf_T]]

Returns all nodes of the tree whose data equals the given data.

scan_values(*pred: Callable[[Union[_Leaf_T, Tree[_Leaf_T]]], bool]*) → Iterator[_Leaf_T]

Return all values in the tree that evaluate `pred(value)` as true.

This can be used to find all the tokens in the tree.

Example

```
>>> all_tokens = tree.scan_values(lambda v: isinstance(v, Token))
```

14.3 Token

```
class lark.Token(type: str, value: Any, start_pos: Optional[int] = None, line: Optional[int] = None, column:
    Optional[int] = None, end_line: Optional[int] = None, end_column: Optional[int] = None,
    end_pos: Optional[int] = None)
```

```
class lark.Token(type_: str, value: Any, start_pos: Optional[int] = None, line: Optional[int] = None, column:
    Optional[int] = None, end_line: Optional[int] = None, end_column: Optional[int] = None,
    end_pos: Optional[int] = None)
```

A string with meta-information, that is produced by the lexer.

When parsing text, the resulting chunks of the input that haven't been discarded, will end up in the tree as Token instances. The Token class inherits from Python's `str`, so normal string comparisons and operations will work as expected.

type

Name of the token (as specified in grammar)

Type

`str`

value

Value of the token (redundant, as `token.value == token` will always be true)

Type

`Any`

start_pos

The index of the token in the text

Type

`Optional[int]`

line

The line of the token in the text (starting with 1)

Type

`Optional[int]`

column

The column of the token in the text (starting with 1)

Type

`Optional[int]`

end_line

The line where the token ends

Type

`Optional[int]`

end_column

The next column after the end of the token. For example, if the token is a single character with a column value of 4, end_column will be 5.

Type

Optional[int]

end_pos

the index where the token ends (basically start_pos + len(token))

Type

Optional[int]

14.4 Transformer, Visitor & Interpreter

See *Transformers & Visitors*.

14.5 ForestVisitor, ForestTransformer, & TreeForestTransformer

See *Working with the SPPF*.

14.6 UnexpectedInput

class lark.exceptions.UnexpectedInput

UnexpectedInput Error.

Used as a base class for the following exceptions:

- **UnexpectedCharacters**: The lexer encountered an unexpected string
- **UnexpectedToken**: The parser received an unexpected token
- **UnexpectedEOF**: The parser expected a token, but the input ended

After catching one of these exceptions, you may call the following helper methods to create a nicer error message.

get_context(text: str, span: int = 40) → str

Returns a pretty string pinpointing the error in the text, with span amount of context characters around it.

Note: The parser doesn't hold a copy of the text it has to parse, so you have to provide it again

match_examples(parse_fn: Callable[[str], Tree], examples: Union[Mapping[T, Iterable[str]], Iterable[Tuple[T, Iterable[str]]]], token_type_match_fallback: bool = False, use_accepts: bool = True) → Optional[T]

Allows you to detect what's wrong in the input text by matching against example errors.

Given a parser instance and a dictionary mapping some label with some malformed syntax examples, it'll return the label for the example that bests matches the current error. The function will iterate the dictionary until it finds a matching error, and return the corresponding value.

For an example usage, see *examples/error_reporting_lalr.py*

Parameters

- **parse_fn** – parse function (usually `lark_instance.parse`)
- **examples** – dictionary of `{'example_string': value}`.
- **use_accepts** – Recommended to keep this as `use_accepts=True`.

```
class lark.exceptions.UnexpectedToken(token, expected, considered_rules=None, state=None,
                                     interactive_parser=None, terminals_by_name=None,
                                     token_history=None)
```

An exception that is raised by the parser, when the token it received doesn't match any valid step forward.

Parameters

- **token** – The mismatched token
- **expected** – The set of expected tokens
- **considered_rules** – Which rules were considered, to deduce the expected tokens
- **state** – A value representing the parser state. Do not rely on its value or type.
- **interactive_parser** – An instance of `InteractiveParser`, that is initialized to the point of failure, and can be used for debugging and error handling.

Note: These parameters are available as attributes of the instance.

```
class lark.exceptions.UnexpectedCharacters(seq, lex_pos, line, column, allowed=None,
                                         considered_tokens=None, state=None, token_history=None,
                                         terminals_by_name=None, considered_rules=None)
```

An exception that is raised by the lexer, when it cannot match the next string of characters to any of its terminals.

```
class lark.exceptions.UnexpectedEOF(expected, state=None, terminals_by_name=None)
```

An exception that is raised by the parser, when the input ends while it still expects a token.

14.7 InteractiveParser

```
class lark.parsers.lalr_interactive_parser.InteractiveParser(parser, parser_state, lexer_thread:
                                                             LexerThread)
```

`InteractiveParser` gives you advanced control over parsing and error handling when parsing with LALR.

For a simpler interface, see the `on_error` argument to `Lark.parse()`.

feed_token(token: [Token](#))

Feed the parser with a token, and advance it to the next state, as if it received it from the lexer.

Note that `token` has to be an instance of `Token`.

exhaust_lexer() → List[[Token](#)]

Try to feed the rest of the lexer state into the interactive parser.

Note that this modifies the instance in place and does not feed an '\$END' Token

as_immutable()

Convert to an `ImmutableInteractiveParser`.

pretty()

Print the output of `choices()` in a way that's easier to read.

choices()

Returns a dictionary of token types, matched to their action in the parser.

Only returns token types that are accepted by the current state.

Updated by `feed_token()`.

accepts()

Returns the set of possible tokens that will advance the parser into a new valid state.

resume_parse()

Resume automated parsing from the current state.

class `lark.parsers.lalr_interactive_parser.ImmutableInteractiveParser`(*parser, parser_state,*
lexer_thread:
LexerThread)

Same as `InteractiveParser`, but operations create a new instance instead of changing it in-place.

feed_token(token)

Feed the parser with a token, and advance it to the next state, as if it received it from the lexer.

Note that `token` has to be an instance of `Token`.

exhaust_lexer()

Try to feed the rest of the lexer state into the parser.

Note that this returns a new `ImmutableInteractiveParser` and does not feed an '\$END' Token

as_mutable()

Convert to an `InteractiveParser`.

choices()

Returns a dictionary of token types, matched to their action in the parser.

Only returns token types that are accepted by the current state.

Updated by `feed_token()`.

pretty()

Print the output of `choices()` in a way that's easier to read.

resume_parse()

Resume automated parsing from the current state.

accepts()

Returns the set of possible tokens that will advance the parser into a new valid state.

14.8 ast_utils

For an example of using `ast_utils`, see /examples/advanced/create_ast.py

class `lark.ast_utils.Ast`

Abstract class

Subclasses will be collected by `create_transformer()`

class `lark.ast_utils.AsList`

Abstract class

Subclasses will be instantiated with the parse results as a single list, instead of as arguments.

`lark.ast_utils.create_transformer`(*ast_module: module, transformer:*
~typing.Optional[~lark.visitors.Transformer] = None,
decorator_factory: ~typing.Callable = <function v_args>) →
Transformer

Collects *Ast* subclasses from the given module, and creates a Lark transformer that builds the AST.

For each class, we create a corresponding rule in the transformer, with a matching name. CamelCase names will be converted into snake_case. Example: “CodeBlock” -> “code_block”.

Classes starting with an underscore (_) will be skipped.

Parameters

- **ast_module** – A Python module containing all the subclasses of `ast_utils.Ast`
- **transformer** (*Optional[Transformer]*) – An initial transformer. Its attributes may be overwritten.
- **decorator_factory** (*Callable*) – An optional callable accepting two booleans, inline, and meta, and returning a decorator for the methods of `transformer`. (default: `v_args`).

TRANSFORMERS & VISITORS

Transformers & Visitors provide a convenient interface to process the parse-trees that Lark returns.

They are used by inheriting from the correct class (visitor or transformer), and implementing methods corresponding to the rule you wish to process. Each method accepts the children as an argument. That can be modified using the `v_args` decorator, which allows one to inline the arguments (akin to `*args`), or add the tree `meta` property as an argument.

See: [visitors.py](#)

15.1 Visitor

Visitors visit each node of the tree, and run the appropriate method on it according to the node's data.

They work bottom-up, starting with the leaves and ending at the root of the tree.

There are two classes that implement the visitor interface:

- **Visitor**: Visit every node (without recursion)
- **Visitor_Recursive**: Visit every node using recursion. Slightly faster.

Example:

```
class IncreaseAllNumbers(Visitor):
    def number(self, tree):
        assert tree.data == "number"
        tree.children[0] += 1

IncreaseAllNumbers().visit(parse_tree)
```

class `lark.visitors.Visitor(*args, **kws)`

Tree visitor, non-recursive (can handle huge trees).

Visiting a node calls its methods (provided by the user via inheritance) according to `tree.data`

visit(*tree*: `Tree[_Leaf_T]`) → `Tree[_Leaf_T]`

Visits the tree, starting with the leaves and finally the root (bottom-up)

visit_topdown(*tree*: `Tree[_Leaf_T]`) → `Tree[_Leaf_T]`

Visit the tree, starting at the root, and ending at the leaves (top-down)

__default__(*tree*)

Default function that is called if there is no attribute matching `tree.data`

Can be overridden. Defaults to doing nothing.

class `lark.visitors.Visitor_Recursive(*args, **kws)`

Bottom-up visitor, recursive.

Visiting a node calls its methods (provided by the user via inheritance) according to `tree.data`

Slightly faster than the non-recursive version.

visit(*tree*: `Tree[_Leaf_T]`) \rightarrow `Tree[_Leaf_T]`

Visits the tree, starting with the leaves and finally the root (bottom-up)

visit_topdown(*tree*: `Tree[_Leaf_T]`) \rightarrow `Tree[_Leaf_T]`

Visit the tree, starting at the root, and ending at the leaves (top-down)

__default__(*tree*)

Default function that is called if there is no attribute matching `tree.data`

Can be overridden. Defaults to doing nothing.

15.2 Interpreter

class `lark.visitors.Interpreter(*args, **kws)`

Interpreter walks the tree starting at the root.

Visits the tree, starting with the root and finally the leaves (top-down)

For each tree node, it calls its methods (provided by user via inheritance) according to `tree.data`.

Unlike `Transformer` and `Visitor`, the `Interpreter` doesn't automatically visit its sub-branches. The user has to explicitly call `visit`, `visit_children`, or use the `@visit_children_decor`. This allows the user to implement branching and loops.

Example:

```
class IncreaseSomeOfTheNumbers(Interpreter):
    def number(self, tree):
        tree.children[0] += 1

    def skip(self, tree):
        # skip this subtree. don't change any number node inside it.
        pass

IncreaseSomeOfTheNumbers().visit(parse_tree)
```

15.3 Transformer

class `lark.visitors.Transformer(visit_tokens: bool = True)`

Transformers work bottom-up (or depth-first), starting with visiting the leaves and working their way up until ending at the root of the tree.

For each node visited, the transformer will call the appropriate method (callbacks), according to the node's data, and use the returned value to replace the node, thereby creating a new tree structure.

Transformers can be used to implement map & reduce patterns. Because nodes are reduced from leaf to root, at any point the callbacks may assume the children have already been transformed (if applicable).

If the transformer cannot find a method with the right name, it will instead call `__default__`, which by default creates a copy of the node.

To discard a node, return `Discard` (`lark.visitors.Discard`).

`Transformer` can do anything `Visitor` can do, but because it reconstructs the tree, it is slightly less efficient.

A transformer without methods essentially performs a non-memoized partial deepcopy.

All these classes implement the transformer interface:

- `Transformer` - Recursively transforms the tree. This is the one you probably want.
- `Transformer_InPlace` - Non-recursive. Changes the tree in-place instead of returning new instances
- `Transformer_InPlaceRecursive` - Recursive. Changes the tree in-place instead of returning new instances

Parameters

visit_tokens (*bool, optional*) – Should the transformer visit tokens in addition to rules. Setting this to `False` is slightly faster. Defaults to `True`. (For processing ignored tokens, use the `lexer_callbacks` options)

transform(*tree: Tree[_Leaf_T]*) → *_Return_T*

Transform the given tree, and return the final result

__mul__(*other: Union[Transformer, TransformerChain[_Leaf_U, _Return_V]]*) → *TransformerChain[_Leaf_T, _Return_V]*

Chain two transformers together, returning a new transformer.

__default__(*data, children, meta*)

Default function that is called if there is no attribute matching data

Can be overridden. Defaults to creating a new copy of the tree node (i.e. `return Tree(data, children, meta)`)

__default_token__(*token*)

Default function that is called if there is no attribute matching `token.type`

Can be overridden. Defaults to returning the token as-is.

Example:

```
from lark import Tree, Transformer

class EvalExpressions(Transformer):
    def expr(self, args):
        return eval(args[0])

t = Tree('a', [Tree('expr', ['1+2'])])
print(EvalExpressions().transform( t ))

# Prints: Tree(a, [3])
```

Example:

```
class T(Transformer):
    INT = int
    NUMBER = float
```

(continues on next page)

(continued from previous page)

```
def NAME(self, name):  
    return lookup_dict.get(name, name)  
  
T(visit_tokens=True).transform(tree)
```

class lark.visitors.**Transformer_NonRecursive**(*visit_tokens: bool = True*)

Same as Transformer but non-recursive.

Like Transformer, it doesn't change the original tree.

Useful for huge trees.

class lark.visitors.**Transformer_InPlace**(*visit_tokens: bool = True*)

Same as Transformer, but non-recursive, and changes the tree in-place instead of returning new instances

Useful for huge trees. Conservative in memory.

class lark.visitors.**Transformer_InPlaceRecursive**(*visit_tokens: bool = True*)

Same as Transformer, recursive, but changes the tree in-place instead of returning new instances

15.4 v_args

lark.visitors.v_args(*inline: bool = False, meta: bool = False, tree: bool = False, wrapper:*
Optional[Callable] = None) → Callable[[Union[Callable[[...], _Return_T], type]],
Union[Callable[[...], _Return_T], type]]

A convenience decorator factory for modifying the behavior of user-supplied visitor methods.

By default, callback methods of transformers/visitors accept one argument - a list of the node's children.

v_args can modify this behavior. When used on a transformer/visitor class definition, it applies to all the callback methods inside it.

v_args can be applied to a single method, or to an entire class. When applied to both, the options given to the method take precedence.

Parameters

- **inline** (*bool, optional*) – Children are provided as **args* instead of a list argument (not recommended for very long lists).
- **meta** (*bool, optional*) – Provides two arguments: **meta** and **children** (instead of just the latter)
- **tree** (*bool, optional*) – Provides the entire tree as the argument, instead of the children.
- **wrapper** (*function, optional*) – Provide a function to decorate all methods.

Example

```

@v_args(inline=True)
class SolveArith(Transformer):
    def add(self, left, right):
        return left + right

    @v_args(meta=True)
    def mul(self, meta, children):
        logger.info(f'mul at line {meta.line}')
        left, right = children
        return left * right

class ReverseNotation(Transformer_InPlace):
    @v_args(tree=True)
    def tree_node(self, tree):
        tree.children = tree.children[::-1]

```

15.5 merge_transformers

`lark.visitors.merge_transformers(base_transformer=None, **transformers_to_merge)`

Merge a collection of transformers into the `base_transformer`, each into its own ‘namespace’.

When called, it will collect the methods from each transformer, and assign them to `base_transformer`, with their name prefixed with the given keyword, as `prefix__methodname`.

This function is especially useful for processing grammars that import other grammars, thereby creating some of their rules in a ‘namespace’. (i.e with a consistent name prefix). In this case, the key for the transformer should match the name of the imported grammar.

Parameters

- **base_transformer** (`Transformer`, *optional*) – The transformer that all other transformers will be added to.
- ****transformers_to_merge** – Keyword arguments, in the form of `name_prefix = transformer`.

Raises

AttributeError – In case of a name collision in the merged methods

Example

```

class TBase(Transformer):
    def start(self, children):
        return children[0] + 'bar'

class TImportedGrammar(Transformer):
    def foo(self, children):
        return "foo"

```

(continues on next page)

(continued from previous page)

```
composed_transformer = merge_transformers(TBase(), imported=TImportedGrammar())

t = Tree('start', [ Tree('imported__foo', []) ])

assert composed_transformer.transform(t) == 'foobar'
```

15.6 Discard

Discard is the singleton instance of `_DiscardType`.

class `lark.visitors._DiscardType`

When the Discard value is returned from a transformer callback, that node is discarded and won't appear in the parent.

Note: This feature is disabled when the transformer is provided to Lark using the `transformer` keyword (aka Tree-less LALR mode).

Example

```
class T(Transformer):
    def ignore_tree(self, children):
        return Discard

    def IGNORE_TOKEN(self, token):
        return Discard
```

15.7 VisitError

class `lark.exceptions.VisitError(rule, obj, orig_exc)`

VisitError is raised when visitors are interrupted by an exception

It provides the following attributes for inspection:

Parameters

- **rule** – the name of the visit rule that failed
- **obj** – the tree-node or token that was being processed
- **orig_exc** – the exception that cause it to fail

Note: These parameters are available as attributes

WORKING WITH THE SPPF

When parsing with Earley, Lark provides the `ambiguity='forest'` option to obtain the shared packed parse forest (SPPF) produced by the parser as an alternative to it being automatically converted to a tree.

Lark provides a few tools to facilitate working with the SPPF. Here are some things to consider when deciding whether or not to use the SPPF.

Pros

- Efficient storage of highly ambiguous parses
- Precise handling of ambiguities
- Custom rule prioritizers
- Ability to handle infinite ambiguities
- Directly transform forest -> object instead of forest -> tree -> object

Cons

- More complex than working with a tree
- SPPF may contain nodes corresponding to rules generated internally
- Loss of Lark grammar features:
 - Rules starting with ‘_’ are not inlined in the SPPF
 - Rules starting with ‘?’ are never inlined in the SPPF
 - All tokens will appear in the SPPF

16.1 SymbolNode

class `lark.parsers.earley_forest.SymbolNode(s, start, end)`

A Symbol Node represents a symbol (or Intermediate LR0).

Symbol nodes are keyed by the symbol (`s`). For intermediate nodes `s` will be an LR0, stored as a tuple of (rule, ptr). For completed symbol nodes, `s` will be a string representing the non-terminal origin (i.e. the left hand side of the rule).

The children of a Symbol or Intermediate Node will always be Packed Nodes; with each Packed Node child representing a single derivation of a production.

Hence a Symbol Node with a single child is unambiguous.

Parameters

- **s** – A Symbol, or a tuple of (rule, ptr) for an intermediate node.
- **start** – The index of the start of the substring matched by this symbol (inclusive).
- **end** – The index of the end of the substring matched by this symbol (exclusive).

Properties:

is_intermediate: True if this node is an intermediate node. **priority:** The priority of the node's symbol.

property is_ambiguous

Returns True if this node is ambiguous.

property children

Returns a list of this node's children sorted from greatest to least priority.

16.2 PackedNode

class `lark.parsers.earley_forest.PackedNode`(*parent, s, rule, start, left, right*)

A Packed Node represents a single derivation in a symbol node.

Parameters

- **rule** – The rule associated with this node.
- **parent** – The parent of this node.
- **left** – The left child of this node. `None` if one does not exist.
- **right** – The right child of this node. `None` if one does not exist.
- **priority** – The priority of this node.

property children

Returns a list of this node's children.

16.3 ForestVisitor

class `lark.parsers.earley_forest.ForestVisitor`(*single_visit=False*)

An abstract base class for building forest visitors.

This class performs a controllable depth-first walk of an SPPF. The visitor will not enter cycles and will backtrack if one is encountered. Subclasses are notified of cycles through the `on_cycle` method.

Behavior for visit events is defined by overriding the `visit*node*` functions.

The walk is controlled by the return values of the `visit*node_in` methods. Returning a node(s) will schedule them to be visited. The visitor will begin to backtrack if no nodes are returned.

Parameters

single_visit – If True, non-Token nodes will only be visited once.

visit_token_node(*node*)

Called when a Token is visited. Token nodes are always leaves.

visit_symbol_node_in(*node*)

Called when a symbol node is visited. Nodes that are returned will be scheduled to be visited. If `visit_intermediate_node_in` is not implemented, this function will be called for intermediate nodes as well.

visit_symbol_node_out(*node*)

Called after all nodes returned from a corresponding `visit_symbol_node_in` call have been visited. If `visit_intermediate_node_out` is not implemented, this function will be called for intermediate nodes as well.

visit_packed_node_in(*node*)

Called when a packed node is visited. Nodes that are returned will be scheduled to be visited.

visit_packed_node_out(*node*)

Called after all nodes returned from a corresponding `visit_packed_node_in` call have been visited.

on_cycle(*node*, *path*)

Called when a cycle is encountered.

Parameters

- **node** – The node that causes a cycle.
- **path** – The list of nodes being visited: nodes that have been entered but not exited. The first element is the root in a forest visit, and the last element is the node visited most recently. path should be treated as read-only.

get_cycle_in_path(*node*, *path*)

A utility function for use in `on_cycle` to obtain a slice of `path` that only contains the nodes that make up the cycle.

16.4 ForestTransformer

class lark.parsers.earley_forest.ForestTransformer

The base class for a bottom-up forest transformation. Most users will want to use `TreeForestTransformer` instead as it has a friendlier interface and covers most use cases.

Transformations are applied via inheritance and overriding of the `transform*node` methods.

`transform_token_node` receives a `Token` as an argument. All other methods receive the node that is being transformed and a list of the results of the transformations of that node's children. The return value of these methods are the resulting transformations.

If `Discard` is raised in a node's transformation, no data from that node will be passed to its parent's transformation.

transform(*root*)

Perform a transformation on an SPPF.

transform_symbol_node(*node*, *data*)

Transform a symbol node.

transform_intermediate_node(*node*, *data*)

Transform an intermediate node.

transform_packed_node(*node*, *data*)

Transform a packed node.

transform_token_node(*node*)

Transform a Token.

16.5 TreeForestTransformer

```
class lark.parsers.earley_forest.TreeForestTransformer(tree_class=<class 'lark.tree.Tree'>, prioritizer=<lark.parsers.earley_forest.ForestSumVisitor object>, resolve_ambiguity=True, use_cache=False)
```

A ForestTransformer with a tree Transformer-like interface. By default, it will construct a tree.

Methods provided via inheritance are called based on the rule/symbol names of nodes in the forest.

Methods that act on rules will receive a list of the results of the transformations of the rule's children. By default, trees and tokens.

Methods that act on tokens will receive a token.

Alternatively, methods that act on rules may be annotated with `handles_ambiguity`. In this case, the function will receive a list of all the transformations of all the derivations of the rule. By default, a list of trees where each `tree.data` is equal to the rule name or one of its aliases.

Non-tree transformations are made possible by override of `__default__`, `__default_token__`, and `__default_ambig__`.

Note: Tree shaping features such as inlined rules and token filtering are not built into the transformation. Positions are also not propagated.

Parameters

- **tree_class** – The tree class to use for construction
- **prioritizer** – A ForestVisitor that manipulates the priorities of nodes in the SPPF.
- **resolve_ambiguity** – If True, ambiguities will be resolved based on priorities.
- **use_cache** (*bool*) – If True, caches the results of some transformations, potentially improving performance when `resolve_ambiguity==False`. Only use if you know what you are doing: i.e. All transformation functions are pure and referentially transparent.

__default__(*name, data*)

Default operation on tree (for override).

Returns a tree with name with data as children.

__default_ambig__(*name, data*)

Default operation on ambiguous rule (for override).

Wraps data in an `'_ambig_'` node if it contains more than one element.

__default_token__(*node*)

Default operation on Token (for override).

Returns node.

16.6 handles_ambiguity

`lark.parsers.earley_forest.handles_ambiguity(func)`

Decorator for methods of subclasses of `TreeForestTransformer`. Denotes that the method should receive a list of transformed derivations.

TOOLS (STAND-ALONE, NEARLEY)

17.1 Stand-alone parser

Lark can generate a stand-alone LALR(1) parser from a grammar.

The resulting module provides the same interface as Lark, but with a fixed grammar, and reduced functionality.

Run using:

```
python -m lark.tools.standalone
```

For a play-by-play, read the [tutorial](#)

17.2 Importing grammars from Nearley.js

Lark comes with a tool to convert grammars from [Nearley](#), a popular Earley library for Javascript. It uses [Js2Py](#) to convert and run the Javascript postprocessing code segments.

17.2.1 Requirements

1. Install Lark with the `nearley` component:

```
pip install lark[nearley]
```

1. Acquire a copy of the Nearley codebase. This can be done using:

```
git clone https://github.com/Hardmath123/nearley
```

17.2.2 Usage

The tool can be run using:

```
python -m lark.tools.nearley <grammar.ne> <start_rule> <path_to_nearley_repo>
```

Here's an example of how to import nearley's calculator example into Lark:

```
git clone https://github.com/Hardmath123/nearley
python -m lark.tools.nearley nearley/examples/calculator/arithmetic.ne main ./nearley > ncalc.py
```

You can use the output as a regular python module:

```
>>> import ncalc
>>> ncalc.parse('sin(pi/4) ^ e')
0.38981434460254655
```

The Nearley converter also supports an experimental converter for newer JavaScript (ES6+), using the `--es6` flag:

```
git clone https://github.com/Hardmath123/nearley
python -m lark.tools.nearley nearley/examples/calculator/arithmetic.ne main nearley --
↳ es6 > ncalc.py
```

17.2.3 Notes

- Lark currently cannot import templates from Nearley
- Lark currently cannot export grammars to Nearley

These might get added in the future, if enough users ask for them.

Lark is a modern parsing library for Python. Lark can parse any context-free grammar.

Lark provides:

- Advanced grammar language, based on EBNF
- Three parsing algorithms to choose from: Earley, LALR(1) and CYK
- Automatic tree construction, inferred from your grammar
- Fast unicode lexer with regexp support, and automatic line-counting

INSTALL LARK

```
$ pip install lark
```


SYNTAX HIGHLIGHTING

- Sublime Text & TextMate
- Visual Studio Code (Or install through the vscode plugin system)
- IntelliJ & PyCharm
- Vim
- Atom

RESOURCES

- *Philosophy*
- *Features*
- *Examples*
- *Third-party examples*
- *Online IDE*
- *Tutorials*
 - *How to write a DSL* - Implements a toy LOGO-like language with an interpreter
 - *JSON parser - Tutorial* - Teaches you how to use Lark
 - *Unofficial*
 - * *Program Synthesis is Possible* - Creates a DSL for Z3
 - * *Using Lark to Parse Text* - Robin Reynolds-Haertle (PyCascades 2023) (video presentation)
- *Guides*
 - *How To Use Lark - Guide*
 - *How to develop Lark - Guide*
- *Reference*
 - *Grammar Reference*
 - *Tree Construction Reference*
 - *Transformers & Visitors*
 - *Working with the SPPF*
 - *API Reference*
 - *Tools (Stand-alone, Nearley)*
 - *Cheatsheet (PDF)*
- *Discussion*
 - *Gitter*
 - *Forum (Google Groups)*

Symbols

`_DiscardType` (class in `lark.visitors`), 100
`__default__()` (`lark.parsers.earley_forest.TreeForestTransformer` method), 104
`__default__()` (`lark.visitors.Transformer` method), 97
`__default__()` (`lark.visitors.Visitor` method), 95
`__default__()` (`lark.visitors.Visitor_Recursive` method), 96
`__default_ambig__()` (`lark.parsers.earley_forest.TreeForestTransformer` method), 104
`__default_token__()` (`lark.parsers.earley_forest.TreeForestTransformer` method), 104
`__default_token__()` (`lark.visitors.Transformer` method), 97
`__mul__()` (`lark.visitors.Transformer` method), 97
`__rich__()` (`lark.Tree` method), 89

A

`accepts()` (`lark.parsers.lalr_interactive_parser.ImmutableInteractiveParser` method), 93
`accepts()` (`lark.parsers.lalr_interactive_parser.InteractiveParser` method), 93
`as_immutable()` (`lark.parsers.lalr_interactive_parser.InteractiveParser` method), 92
`as_mutable()` (`lark.parsers.lalr_interactive_parser.ImmutableInteractiveParser` method), 93
`AsList` (class in `lark.ast_utils`), 93
`Ast` (class in `lark.ast_utils`), 93

C

`children` (`lark.parsers.earley_forest.PackedNode` property), 102
`children` (`lark.parsers.earley_forest.SymbolNode` property), 102
`choices()` (`lark.parsers.lalr_interactive_parser.ImmutableInteractiveParser` method), 93
`choices()` (`lark.parsers.lalr_interactive_parser.InteractiveParser` method), 92
`column` (`lark.Token` attribute), 90
`create_transformer()` (in module `lark.ast_utils`), 94

E

`end_column` (`lark.Token` attribute), 90
`end_line` (`lark.Token` attribute), 90
`end_pos` (`lark.Token` attribute), 91
`exhaust_lexer()` (`lark.parsers.lalr_interactive_parser.ImmutableInteractiveParser` method), 93
`exhaust_lexer()` (`lark.parsers.lalr_interactive_parser.InteractiveParser` method), 92

F

`feed_token()` (`lark.parsers.lalr_interactive_parser.ImmutableInteractiveParser` method), 93
`feed_token()` (`lark.parsers.lalr_interactive_parser.InteractiveParser` method), 92
`find_data()` (`lark.Tree` method), 89
`find_pred()` (`lark.Tree` method), 89
`ForestTransformer` (class in `lark.parsers.earley_forest`), 103
`ForestVisitor` (class in `lark.parsers.earley_forest`), 102

G

`get_parser_in_path()` (`lark.parsers.earley_forest.ForestVisitor` method), 103
`get_terminal()` (`lark.Lark` method), 87

H

`handles_ambiguity()` (in module `lark.parsers.earley_forest`), 105

I

`ImmutableInteractiveParser` (class in `lark.parsers.lalr_interactive_parser`), 93
`InteractiveParser` (class in `lark.parsers.lalr_interactive_parser`), 92
`Interpreter` (class in `lark.visitors`), 96
`is_ambiguous` (`lark.parsers.earley_forest.SymbolNode` property), 102
`iter_subtrees()` (`lark.Tree` method), 89

`iter_subtrees_topdown()` (*lark.Tree* method), 89

L

`Lark` (class in *lark*), 85

`lex()` (*lark.Lark* method), 87

`line` (*lark.Token* attribute), 90

`load()` (*lark.Lark* class method), 87

M

`match_examples()` (*lark.exceptions.UnexpectedInput* method), 91

`merge_transformers()` (in module *lark.visitors*), 99

O

`on_cycle()` (*lark.parsers.earley_forest.ForestVisitor* method), 103

`open()` (*lark.Lark* class method), 87

`open_from_package()` (*lark.Lark* class method), 87

P

`PackedNode` (class in *lark.parsers.earley_forest*), 102

`parse()` (*lark.Lark* method), 88

`parse_interactive()` (*lark.Lark* method), 87

`pretty()` (*lark.parsers.lalr_interactive_parser.ImmutableInteractiveParser* method), 93

`pretty()` (*lark.parsers.lalr_interactive_parser.InteractiveParser* method), 92

`pretty()` (*lark.Tree* method), 89

R

`resume_parse()` (*lark.parsers.lalr_interactive_parser.ImmutableInteractiveParser* method), 93

`resume_parse()` (*lark.parsers.lalr_interactive_parser.InteractiveParser* method), 93

S

`save()` (*lark.Lark* method), 87

`scan_values()` (*lark.Tree* method), 89

`start_pos` (*lark.Token* attribute), 90

`SymbolNode` (class in *lark.parsers.earley_forest*), 101

T

`Token` (class in *lark*), 90

`transform()` (*lark.parsers.earley_forest.ForestTransformer* method), 103

`transform()` (*lark.visitors.Transformer* method), 97

`transform_intermediate_node()`
(*lark.parsers.earley_forest.ForestTransformer* method), 103

`transform_packed_node()`
(*lark.parsers.earley_forest.ForestTransformer* method), 103

`transform_symbol_node()`
(*lark.parsers.earley_forest.ForestTransformer* method), 103

`transform_token_node()`
(*lark.parsers.earley_forest.ForestTransformer* method), 103

`Transformer` (class in *lark.visitors*), 96

`Transformer_InPlace` (class in *lark.visitors*), 98

`Transformer_InPlaceRecursive` (class in *lark.visitors*), 98

`Transformer_NonRecursive` (class in *lark.visitors*), 98

`Tree` (class in *lark*), 89

`TreeForestTransformer` (class in *lark.parsers.earley_forest*), 104

`type` (*lark.Token* attribute), 90

U

`UnexpectedCharacters` (class in *lark.exceptions*), 92

`UnexpectedEOF` (class in *lark.exceptions*), 92

`UnexpectedInput` (class in *lark.exceptions*), 91

`UnexpectedToken` (class in *lark.exceptions*), 92

V

`v_args()` (in module *lark.visitors*), 98

`value` (*lark.Token* attribute), 90

`visit()` (*lark.visitors.Visitor* method), 95

`visit()` (*lark.visitors.Visitor_Recursive* method), 96

`visit_packed_node_in()`
(*lark.parsers.earley_forest.ForestVisitor* method), 103

`visit_packed_node_out()`
(*lark.parsers.earley_forest.ForestVisitor* method), 103

`visit_symbol_node_in()`
(*lark.parsers.earley_forest.ForestVisitor* method), 102

`visit_symbol_node_out()`
(*lark.parsers.earley_forest.ForestVisitor* method), 103

`visit_token_node()` (*lark.parsers.earley_forest.ForestVisitor* method), 102

`visit_topdown()` (*lark.visitors.Visitor* method), 95

`visit_topdown()` (*lark.visitors.Visitor_Recursive* method), 96

`VisitError` (class in *lark.exceptions*), 100

`Visitor` (class in *lark.visitors*), 95

`Visitor_Recursive` (class in *lark.visitors*), 95