

```
{match: /"""\[^]*?"""/,  
lineBreaks: true, value: x =>  
x.slice(3, -3)},  
{match: /(?:\\["\\rn]|[^"\\\n])*?"/, lineBreaks: true,  
value: x => x.slice(1, -1)},  
{match: /'(?:\\['\\rn]|[^'\\\n])*?'/, lineBreaks: true,  
value: x => x.slice(1, -1)},  
],  
// ...  
})
```

Contributing

Do check the [FAQ](#).

Before submitting an issue, [remember...](#)

Literation

```
// operators; we got 'em.  
for (let here of lexer) {  
  // here = { type: 'number', value:  
  //           '123', ... }  
  let tokens = Array.from(lexer);  
  Use it's iteration tools with Moo.  
  for (let [here, next] of  
    itt(lexer).lookahead()) { //  
    pass a number if you need more  
    tokens  
    // enjoy!
```

Moo is a highly-optimised tokeniser/lexer generator. Use it to tokenize your strings, before parsing them with a parser like [nearley](#) or whatever else you're into.

MOON

Transform

Yup! Flyings-cows-and-simged-steak fast.
Moo is the fastest JS tokenizer around. It's ~2-10x faster than
most other tokenizers; it's a couple orders of magnitude faster
than some of the slower ones.

Is it fast?

```
moo.reset('invalid')
moo.next() // -> { type: 'myError',
  value: 'invalid', text:
  'invalid', offset: 0,
  lineBreaks: 0, line: 1, col: 1 }
moo.next() // -> undefined
```

You can have a token type that both matches tokens *and* contains error values.

```
moo.compile({
  // ...
  myError: {match: /[\$?`]/, error:
    true},
})
```

Formatting errors

If you want to throw an error from your parser, you might find `formatError` helpful. Call it with the offending token:

```
throw new
  Error(lexer.formatError(token,
  "invalid syntax"))
```

It returns a string with a pretty error message.

```
Error: invalid syntax at line 2 col 15:
totally valid `syntax`  
^
```

Define your tokens **using regular expressions**. Moo will compile 'em down to a **single RegExp for performance**. It uses the new ES6 **sticky flag** where possible to make things faster; otherwise it falls back to an almost-as-efficient workaround. (For more than you ever wanted to know about this, read [adventures in the land of substrings and RegExps](#).)

You *might* be able to go faster still by writing your lexer by hand rather than using RegExps, but that's icky.

Oh, and it [avoids parsing RegExps by itself](#). Because that would be horrible.

Usage

First, you need to do the needful: `$ npm install moo`, or whatever will ship this code to your computer. Alternatively, grab the `moo.js` file by itself and slap it into your web page via a `<script>` tag; moo is completely standalone.

Then you can start roasting your very own lexer/tokenizer:

```
const moo = require('moo')

let lexer = moo.compile({
  WS:      /[ \t]+/,
  comment: /\//\//.*?$/,
  number: /0|[1-9][0-9]*/,
  string: /"(?:\\["\\]|[^\\n\\])*)"/,
  lparen: '(',
```

```
If none of your rules match, Moo will throw an Error; since it  
doesn't know what else to do.  
If you prefer, you can have Moo return an error token instead  
of throwing an exception. The error token will contain the whole  
of the rest of the buffer.  
myError: moo.error,  
// ...  
moo.compile({})
```

Errors

The racing rule is annotated with `Pop`, so it moves from the main state into either `Left` or `Right`, depending on the stack.

```

    rparen: ')',
    keyword: ['while', 'if', 'else'],
    NL: { match: /\n/, 'cows': [', moo', ','],
    linebreaks: true },
    NL: { match: /\n/, 'cows': [', moo', ','],
    rparen: ')',
    And now throw some text at it:
    lexer.reset('while (10) cows\nmoo')
    lexer.next() // -> { type: 'keyword', value: 'while' },
    lexer.next() // -> { type: 'value', value: '10' },
    lexer.next() // -> { type: 'rparen', value: ')' },
    ...
more data when that happens.
When you reach the end of Moo's internal buffer, next() will return undefined. You can always reset() it and feed it

```

States

Moo allows you to define multiple lexer **states**. Each state defines its own separate set of token rules. Your lexer will start off in the first state given to `moo.states({})`.

Rules can be annotated with `next`, `push`, and `pop`, to change the current state after that token is matched. A “stack” of past states is kept, which is used by `push` and `pop`.

- `next`: 'bar' moves to the state named bar. (The stack is not changed.)
- `push`: 'bar' moves to the state named bar, and pushes the old state onto the stack.
- `pop`: 1 removes one state from the top of the stack, and moves to that state. (Only 1 is supported.)

Only rules from the current state can be matched. You need to copy your rule into all the states you want it to be matched in.

For example, to tokenize JS-style string interpolation such as `a${{c: d}}e`, you might use:

```
let lexer = moo.states({
  main: {
    strstart: {match: ``', push:
      'lit'},
    ident:   /\w+/, 
    lbrace:  {match: '{', push:
      'main'},
    rbrace:  {match: '}', pop: 1},
    colon:   ':',
    space:   {match: /\s+/,
    lineBreaks: true},
```

On Regular Expressions

RegExps are nifty for making tokenizers, but they can be a bit of a pain. Here are some things to be aware of:

- You often want to use **non-greedy quantifiers**: e.g. `*?` instead of `*`. Otherwise your tokens will be longer than you expect:

```
let lexer = moo.compile({
  string: /".*"/,    // greedy quantifier *
  // ...
})

lexer.reset('foo" "bar')
lexer.next() // -> { type: 'string', value:
  'foo" "bar' }
```

Better:

```
let lexer = moo.compile({
  string: /".*?"/,  // non-greedy
  // ...
})

lexer.reset('foo" "bar')
lexer.next() // -> { type: 'string', value:
  'foo' }
```

Keyword Types

Keywords can also have **individual types**.

```
let lexer = moo.compile({
  name: {match: /[a-zA-Z]+/, type: preceedence},
  KW-if: 'if',
  KW-def: 'def',
  KW-class: 'class',
  moo.keywords()
})
```

```
lexer.reset('def foo')
lexer.next() // ...
{}}
// ...
{}}
,
'if',
'def',
'class',
moo.keywords()
```

```
lexer.reset('def')
lexer.next() // ...
{}}
// ...
{}}
,
'if',
'def',
'class',
moo.keywords()
```

```
lexer.reset('def')
lexer.next() // ...
{}}
// ...
{}}
,
'if',
'def',
'class',
moo.keywords()
```

```
lexer.reset('def')
lexer.next() // ...
{}}
// ...
{}}
,
'if',
'def',
'space',
name: {type: 'name',
value: 'foo'}}

lexer.next() // -> {type: 'name',
value: 'bar'}
```

```
You can use Object.fromEntries to easily construct keyword objects:
```

```
object.fromEntries([
  ['if', {map(k => [`kw-` + k, k])}]]
```

- Since an excluding character ranges like `/[^]/` (which matches anything but a space) will include newlines, you have to be careful not to include them by accident! In particular, the whitespace metacharacter `\s` includes newlines.
- If you want to match newlines too, if you want to match newlines too.

example, the `dot / .` **doesn't include newlines**. Use `[^]` instead of `dot / .`

```
moo.compile({
  number: /[0-9]+/,
  identifier: /[a-zA-Z]+/,
  reset('42').next() // -> {type: 'number', value: '42'},
  identifier: /[a-zA-Z]+/,
  reset('42').next() // -> {type: 'number', value: '42'}
```

```
moo.compile({
  identifier: /[0-9]+/,
  number: /[a-zA-Z]+/,
  reset('42').next() // -> {type: 'identifier', value: '42'},
  identifier: /[a-zA-Z]+/
```

- The **order of your rules** matters. Earlier ones will take precedence.

```
lexer.next() // -> {type: 'space', value: ' '}
{
  'string', value: 'bar'
}
lexer.next() // -> {type: 'string', value: 'bar'}
```

```

moo.compile({
  IDEN: {match: /[a-zA-Z]+/, type:
    moo.keywords({
      KW: ['while', 'if', 'else',
        'moo', 'cows'],
    })},
  SPACE: {match: /\s+/, lineBreaks:
    true},
})

```

Why?

You need to do this to ensure the **longest match** principle applies, even in edge cases.

Imagine trying to parse the input `className` with the following rules:

```

keyword: ['class'],
identifier: /[a-zA-Z]+/,

```

You'll get *two* tokens — `['class', 'Name']` — which is *not* what you want! If you swap the order of the rules, you'll fix this example; but now you'll lex `class` wrong (as an `identifier`).

The keywords helper checks matches against the list of keywords; if any of them match, it uses the type '`keyword`' instead of '`identifier`' (for this example).

Line Numbers

Moo tracks detailed information about the input for you.

It will track line numbers, as long as you **apply the `lineBreaks: true` option to any rules which might contain newlines**. Moo will try to warn you if you forget to do this.

Note that this is `false` by default, for performance reasons: counting the number of lines in a matched token has a small cost. For optimal performance, only match newlines inside a dedicated token:

```

newline: {match: '\n', lineBreaks:
  true},

```

Token Info

Token objects (returned from `next()`) have the following attributes:

- `type`: the name of the group, as passed to `compile`.
- `text`: the string that was matched.
- `value`: the string that was matched, transformed by your `value` function (if any).
- `offset`: the number of bytes from the start of the buffer where the match starts.
- `lineBreaks`: the number of line breaks found in the match. (Always zero if this rule has `lineBreaks: false`.)

It will automatically compile them into regular expressions, capturing them where necessary. **Keywords** should be written using the **keywords** transform.

Keywords

```
lexer.reset() // -> { line: 10 }
let info = lexer.save() // ->
lexer.reset('some line\n') // -> { line: 10 }
let info = lexer.next() // -> { line: 10 }
lexer.reset('a different line\n',
            ...
// 
```

If you don't want this, you can save() the state, and later pass it as the second argument to reset() to explicitly control the internal state of the Lexer.

and set the `lime`, `column`, and `offset` counts back to their initial value.

Reset

Value vs. Text

- Line: the line number of the beginning of the match, starting from 1.
 - Col: the column where the match begins, starting from 1.