

Lukas Renggli ([http://www.lukas-renggli.ch/?\\_s=r5IS9pbUk273Oern&\\_k=pN4u1qIO&\\_n&4](http://www.lukas-renggli.ch/?_s=r5IS9pbUk273Oern&_k=pN4u1qIO&_n&4))

Blog ([http://www.lukas-renggli.ch/blog?\\_s=r5IS9pbUk273Oern&\\_k=pN4u1qIO&\\_n&13](http://www.lukas-renggli.ch/blog?_s=r5IS9pbUk273Oern&_k=pN4u1qIO&_n&13))

## Writing Parsers with PetitParser

Consulting ([http://www.lukas-renggli.ch/consulting?\\_s=r5IS9pbUk273Oern&\\_k=pN4u1qIO&\\_n&14](http://www.lukas-renggli.ch/consulting?_s=r5IS9pbUk273Oern&_k=pN4u1qIO&_n&14))

Contact ([http://www.lukas-renggli.ch/contact?\\_s=r5IS9pbUk273Oern&\\_k=pN4u1qIO&\\_n&15](http://www.lukas-renggli.ch/contact?_s=r5IS9pbUk273Oern&_k=pN4u1qIO&_n&15))

Smalltalk ([http://www.lukas-renggli.ch/smalltalk?\\_s=r5IS9pbUk273Oern&\\_k=pN4u1qIO&\\_n&16](http://www.lukas-renggli.ch/smalltalk?_s=r5IS9pbUk273Oern&_k=pN4u1qIO&_n&16))

iPhone ([http://www.lukas-renggli.ch/iphone?\\_s=r5IS9pbUk273Oern&\\_k=pN4u1qIO&\\_n&17](http://www.lukas-renggli.ch/iphone?_s=r5IS9pbUk273Oern&_k=pN4u1qIO&_n&17))

(<http://www.iam.unibe.ch/pipermail/moose-dev/2010-February/003956.html>) in the Moose mailing list and after [various](http://damiencassou.seasidehosting.st/) (<http://damiencassou.seasidehosting.st/>) [people](http://stephane.ducasse.free.fr/) (<http://stephane.ducasse.free.fr/>) [have](http://www.tudorgirba.com/) (<http://www.tudorgirba.com/>) asked me to provide some introduction to PetitParser I decided to write short tutorial.

Originally I have written PetitParser as part of my work on the [Helvetia](http://scg.unibe.ch/research/helvetia) (<http://scg.unibe.ch/research/helvetia>) system. PetitParser is a parsing framework different to many other popular parser generators. For example, it is not table based such as [SmaCC](http://www.refactory.com/Software/SmaCC/index.html) (<http://www.refactory.com/Software/SmaCC/index.html>) or [ANTLR](http://www.antlr.org/) (<http://www.antlr.org/>). Instead it uses a unique combination of four alternative parser methodologies: scannerless parsers, parser combinators, parsing expression grammars and packrat parsers. As such PetitParser is more powerful in what it can parse and it arguably fits better the dynamic nature of Smalltalk. Let's have a quick look at these four parser methodologies:

1. *Scannerless Parsers* combine what is usually done by two independent tools (scanner and parser) into one. This makes writing a grammar much simpler and avoids common problems when grammars are composed.
2. *Parser Combinators* are building blocks for parsers modeled as a graph of composable objects; they are modular and maintainable, and can be changed, recomposed, transformed and reflected upon.
3. *Parsing Expression Grammars* (PEGs) provide ordered choice. Unlike in parser combinators, the ordered choice of PEGs always follows the first matching alternative and ignores other alternatives. Valid input always results in exactly one parse-tree, the result of a parse is never ambiguous.
4. *Packrat Parsers* give linear parse time guarantees and avoid common problems with left-recursion in PEGs.

## Loading PetitParser

Enough theory, let's get started. PetitParser is developed in [Pharo](http://www.pharo-project.com) (<http://www.pharo-project.com>), but is also available in [GNU Smalltalk](http://smalltalk.gnu.org/) (<http://smalltalk.gnu.org/>). To load PetitParser evaluate the following Gofer expression:

```
Gofer new
  renggli: 'petit';
  package: 'PetitParser';
  load.
```

There are other packages in the same repository that provide additional features, for example PetitSmalltalk is a Smalltalk grammar, PetitXml is an XML grammar, PetitAnalyzer can perform analysis on grammars, and PetitGui is a Glamour IDE for writing complex grammars. We are not going to use any of these packages for now.

## Writing a Simple Grammar

Writing grammars with PetitParser is simple as writing Smalltalk code. For example to write a grammar that can parse identifiers that start with a letter followed by zero or more letter or digits is defined as follows. In a smalltalk session evaluate:

Search lukas-renggli.ch

After the  
[announcement](#)

## Archive

### February 2010

([http://www.lukas-renggli.ch/blog/petitparser-1?\\_s=r5IS9pbUk273Oern&\\_k=pN4u1qIO&18](http://www.lukas-renggli.ch/blog/petitparser-1?_s=r5IS9pbUk273Oern&_k=pN4u1qIO&18))  
(2)

### January 2010

([http://www.lukas-renggli.ch/blog/petitparser-1?\\_s=r5IS9pbUk273Oern&\\_k=pN4u1qIO&19](http://www.lukas-renggli.ch/blog/petitparser-1?_s=r5IS9pbUk273Oern&_k=pN4u1qIO&19))  
(1)

### November 2009

([http://www.lukas-renggli.ch/blog/petitparser-1?\\_s=r5IS9pbUk273Oern&\\_k=pN4u1qIO&20](http://www.lukas-renggli.ch/blog/petitparser-1?_s=r5IS9pbUk273Oern&_k=pN4u1qIO&20))  
(2)

### October 2009

([http://www.lukas-renggli.ch/blog/petitparser-1?\\_s=r5IS9pbUk273Oern&\\_k=pN4u1qIO&21](http://www.lukas-renggli.ch/blog/petitparser-1?_s=r5IS9pbUk273Oern&_k=pN4u1qIO&21))  
(2)

### September 2009

([http://www.lukas-renggli.ch/blog/petitparser-1?\\_s=r5IS9pbUk273Oern&\\_k=pN4u1qIO&22](http://www.lukas-renggli.ch/blog/petitparser-1?_s=r5IS9pbUk273Oern&_k=pN4u1qIO&22))  
(3)

### August 2009

([http://www.lukas-renggli.ch/blog/petitparser-1?\\_s=r5IS9pbUk273Oern&\\_k=pN4u1qIO&23](http://www.lukas-renggli.ch/blog/petitparser-1?_s=r5IS9pbUk273Oern&_k=pN4u1qIO&23))  
(1)

### July 2009

([http://www.lukas-renggli.ch/blog/petitparser-1?\\_s=r5IS9pbUk273Oern&\\_k=pN4u1qIO&24](http://www.lukas-renggli.ch/blog/petitparser-1?_s=r5IS9pbUk273Oern&_k=pN4u1qIO&24))  
(1)

### June 2009

(<http://www.lukas-renggli.ch/blog/petitparser-1>)

more letter or digits is defined as follows. In a workspace we evaluate:

```
identifier := #letter asParser , #word asParser star.
```

If you inspect the object `identifier` you'll notice that it is an instance of a `PPSequenceParser`. This is because the `#`, operator created a sequence of a *letter* and a *zero or more word character* parser. If you dive further into the object you notice the following simple composition of different parser objects:

```
PPSequenceParser (this parser accepts a sequence of parsers)
  PPPredicateParser (this parser accepts a single letter)
  PPRepeatingParser (this parser accepts zero or more instances of another parser)
    PPPredicateParser (this parser accepts a single word character)
```

## Parsing Some Input

To actually parse a string (or stream) we can use the method `#parse::`

```
identifier parse: 'yeah'.      " --> #($y #($e $a $h))
identifier parse: 'f12'.      " --> #($f #($1 $2)) "
```

While it seems odd to get these nested arrays with characters as a return value, this is the default decomposition of the input into a parse tree. We'll see in a while how that can be customized.

If we try to parse something invalid we get an instance of `PPFailure` as an answer:

```
identifier parse: '123'.      " --> letter expected at (9)0 "
```

Instances of `PPFailure` are the only objects in the system that answer with `true` when you send the message `#isPetitFailure`. Alternatively you can also use `#parse: onError:` to throw an exception in case of an error:

```
identifier
  parse: '123'
  onError: [ :msg :pos | self error: msg ].
```

If you are only interested if a given string (or stream) matches or not you can use the following constructs:

```
identifier matches: 'foo'.      " --> true "
identifier matches: '123'.      " --> false "
```

Furthermore to find all matches in a given input string (or stream) you can use:

```
identifier matchesIn: 'foo 123 bar12'.
```

## Different Kinds of Parsers

PetitParser provide a large set of ready-made parser that you can compose to consume and transform arbitrarily complex languages. The terminal parsers are the most simple ones. We've already seen a few of those:

Terminal Parsers	Description
<code>\$a asParser</code>	Parses the character <code>\$a</code> .
<code>'abc' asParser</code>	Parses the string <code>'abc'</code> .
<code>#any asParser</code>	Parses any character.
<code>#digit asParser</code>	Parses the digits 0..9.
<code>#letter asParser</code>	Parses the letters a..z and A..Z.

```
1?
s=r5IS9pbUk273Oern& k=pN4u1qIO&25)
(1)
May 2009
(http://www.lukas-
renggli.ch/blog/petitparser-
1?
s=r5IS9pbUk273Oern& k=pN4u1qIO&26)
(1)
March 2009
(http://www.lukas-
renggli.ch/blog/petitparser-
1?
s=r5IS9pbUk273Oern& k=pN4u1qIO&27)
(4)

Tags

esug
(http://www.lukas-
renggli.ch/blog/petitparser-
1?
s=r5IS9pbUk273Oern& k=pN4u1qIO&28)
(9)
magritte
(http://www.lukas-
renggli.ch/blog/petitparser-
1?
s=r5IS9pbUk273Oern& k=pN4u1qIO&29)
(9)0 "
omnibrowser
(http://www.lukas-
renggli.ch/blog/petitparser-
1?
s=r5IS9pbUk273Oern& k=pN4u1qIO&30)
(7)
pier
(http://www.lukas-
renggli.ch/blog/petitparser-
1?
s=r5IS9pbUk273Oern& k=pN4u1qIO&31)
(15)
refactoring
(http://www.lukas-
renggli.ch/blog/petitparser-
1?
s=r5IS9pbUk273Oern& k=pN4u1qIO&32)
(7)
seaside
(http://www.lukas-
renggli.ch/blog/petitparser-
1?
s=r5IS9pbUk273Oern& k=pN4u1qIO&33)
(58)
smalltalk
(http://www.lukas-
renggli.ch/blog/petitparser-
1?
s=r5IS9pbUk273Oern& k=pN4u1qIO&34)
(22)
squeak
(http://www.lukas-
renggli.ch/blog/petitparser-
1?
s=r5IS9pbUk273Oern& k=pN4u1qIO&35)
(40)
```

The class side of `PPPredicateParser` provides a lot of other factory methods that can be used to build more complex terminal parsers.

The next set of parsers are used to combine other parsers together:

Parser	Description
<b>Combinators</b>	
<code>p1 , p2</code>	Parses <code>p1</code> followed by <code>p2</code> (sequence).
<code>p1 / p2</code>	Parses <code>p1</code> , if that doesn't work parses <code>p2</code> (ordered choice).
<code>p star</code>	Parses zero or more <code>p</code> .
<code>p plus</code>	Parses one or more <code>p</code> .
<code>p optional</code>	Parses <code>p</code> if possible.
<code>p and</code>	Parses <code>p</code> but does not consume its input.
<code>p not</code>	Parses <code>p</code> and succeed when <code>p</code> fails, but does not consume its input.
<code>p end</code>	Parses <code>p</code> and succeed at the end of the input.

So instead of using the `#word` predicated we could have written our identifier parser like this:

```
identifier := #letter asParser , (#letter asParser / #digit asParser) star.
```

To attach an action or transformation to a parser we can use the following methods:

Action Parsers	Description
<code>p ==&gt; aBlock</code>	Performs the transformation given in <code>aBlock</code> .
<code>p flatten</code>	Creates a string from the result of <code>p</code> .
<code>p token</code>	Creates a token from the result of <code>p</code> and consume any spaces.

To return a string of the parsed identifier, we can modify our parser like this:

```
identifier := (#letter asParser , (#letter asParser / #digit asParser) star) flatten.
```

These are the basic elements to build parsers. There are a few more well documented and tested factory methods in the `operations` protocol of `PPParser`. If you want browse that protocol.

## Writing a More Complicated Grammar

Now we are able to write a more complicated grammar for evaluating simple arithmetic expressions. Within a workspace we start with the grammar for a number (actually an integer):

```
number := #digit asParser plus token ==> [ :token | token value asNumber ].
```

Then we define the productions for addition and multiplication in order of precedence. Note that we instantiate the productions as `PPUnresolvedParser` upfront, because they recursively refer to each other. The method `#def:` resolves this recursion using the reflective facilities of the host language:

```
term := PPUnresolvedParser new.
prod := PPUnresolvedParser new.
prim := PPUnresolvedParser new.

term def: (prod , $+ asParser token , term ==> [ :nodes | nodes first + nodes last ])
/ prod.
prod def: (prim , $* asParser token , prod ==> [ :nodes | nodes first * nodes last ])
/ prim.
prim def: ($ ( asParser token , term , $ ) asParser token ==> [ :nodes | nodes second ])
/ number.
```

To make sure that our parser consumes all input we wrap it with the `end` parser into the start production:

```
start := term end.
```

(10)  
[talk](http://www.lukas-renggli.ch/blog/petitparser-1?_s=r5IS9pbUk273Oern&_k=pN4u1qIO&36)  
[tutorial](http://www.lukas-renggli.ch/blog/petitparser-1?_s=r5IS9pbUk273Oern&_k=pN4u1qIO&36)  
 (12)  
[tutorial](http://www.lukas-renggli.ch/blog/petitparser-1?_s=r5IS9pbUk273Oern&_k=pN4u1qIO&37)  
 (7)

That's it, now we can test our parser:

```
start parse: '1 + 2 * 3'.      " --> 7 "
```

```
start parse: '(1 + 2) * 3'.    " --> 9 "
```

As an exercise we could extend the parser to also accept negative numbers and floating point numbers, not only integers. Furthermore it would be useful to add support subtraction and division as well. All these features can be added with a few lines of PetitParser code.

Posted by Lukas Renggli at 25 February 2010, 5:34 pm with tags [tutorial \(http://www.lukas-renggli.ch/blog/petitparser-1?\\_s=r5IS9pbUk273Oern&\\_k=pN4u1qIO&38\)](http://www.lukas-renggli.ch/blog/petitparser-1?_s=r5IS9pbUk273Oern&_k=pN4u1qIO&38), [petitparser \(http://www.lukas-renggli.ch/blog/petitparser-1?\\_s=r5IS9pbUk273Oern&\\_k=pN4u1qIO&39\)](http://www.lukas-renggli.ch/blog/petitparser-1?_s=r5IS9pbUk273Oern&_k=pN4u1qIO&39), [smalltalk \(http://www.lukas-renggli.ch/blog/petitparser-1?\\_s=r5IS9pbUk273Oern&\\_k=pN4u1qIO&40\)](http://www.lukas-renggli.ch/blog/petitparser-1?_s=r5IS9pbUk273Oern&_k=pN4u1qIO&40), [pharo \(http://www.lukas-renggli.ch/blog/petitparser-1?\\_s=r5IS9pbUk273Oern&\\_k=pN4u1qIO&41\)](http://www.lukas-renggli.ch/blog/petitparser-1?_s=r5IS9pbUk273Oern&_k=pN4u1qIO&41), [link \(http://www.lukas-renggli.ch/blog/petitparser-1\)](http://www.lukas-renggli.ch/blog/petitparser-1)

## Leave your comment

Author:  \*

Website:

Contents:  \*

Confirmation: ☐ Confirm submission by clicking the checkbox

[Entries \(RSS\) \(http://www.lukas-renggli.ch/blog?view=PBEEntriesRssView\)](http://www.lukas-renggli.ch/blog?view=PBEEntriesRssView) and [Comments \(RSS\) \(http://www.lukas-renggli.ch/blog?view=PBCommentsRssView\)](http://www.lukas-renggli.ch/blog?view=PBCommentsRssView)

Copyright © 26 February 2010 Lukas Renggli ([http://www.lukas-renggli.ch/contact?\\_s=r5IS9pbUk273Oern&\\_k=pN4u1qIO&\\_n&50](http://www.lukas-renggli.ch/contact?_s=r5IS9pbUk273Oern&_k=pN4u1qIO&_n&50)) . All rights reserved.  
Powered by [Seaside \(http://www.seaside.st/\)](http://www.seaside.st/) , [Magritte \(http://www.lukas-renggli.ch/smalltalk/magritte?\\_s=r5IS9pbUk273Oern&\\_k=pN4u1qIO&\\_n&51\)](http://www.lukas-renggli.ch/smalltalk/magritte?_s=r5IS9pbUk273Oern&_k=pN4u1qIO&_n&51) , and [Pier \(http://www.piercms.com/\)](http://www.piercms.com/)