

Benchmarking Parsers

Nikolaos Bezirgiannis

March 22, 2012

Outline

- ▶ Intro
- ▶ Parser Combinator Libraries
- ▶ Benchmark Results
- ▶ Outro

State of Parsing

- ▶ Parser Generators

State of Parsing

- ▶ Parser Generators
 - ▶ Yacc / Bison
 - ▶ LALR(1) algorithm
 - ▶ Accompanied by a lexer (e.g. lex / flex)
 - ▶ Generates code for C / C++ / Java

State of Parsing

- ▶ Parser Generators

- ▶ Yacc / Bison

- ▶ LALR(1) algorithm

- ▶ Accompanied by a lexer (e.g. lex / flex)

- ▶ Generates code for C / C++ / Java

- ▶ ANTLR

- ▶ LL(*) algorithm

- ▶ Generates code for many languages

- ▶ GUI tools: editor, interpreter, debugger, visualization utilities

State of Parsing

- ▶ Parser Generators
 - ▶ Yacc / Bison
 - ▶ LALR(1) algorithm
 - ▶ Accompanied by a lexer (e.g. lex / flex)
 - ▶ Generates code for C / C++ / Java
 - ▶ ANTLR
 - ▶ LL(*) algorithm
 - ▶ Generates code for many languages
 - ▶ GUI tools: editor, interpreter, debugger, visualization utilities
 - ▶ Happy
 - ▶ Written in Haskell
 - ▶ Generates Haskell
 - ▶ Accompanied by Alex (also in Haskell)
 - ▶ Pretty solid & stable

State of Parsing

- ▶ Parser Generators
 - ▶ Yacc / Bison
 - ▶ LALR(1) algorithm
 - ▶ Accompanied by a lexer (e.g. lex / flex)
 - ▶ Generates code for C / C++ / Java
 - ▶ ANTLR
 - ▶ LL(*) algorithm
 - ▶ Generates code for many languages
 - ▶ GUI tools: editor, interpreter, debugger, visualization utilities
 - ▶ Happy
 - ▶ Written in Haskell
 - ▶ Generates Haskell
 - ▶ Accompanied by Alex (also in Haskell)
 - ▶ Pretty solid & stable
- ▶ Parser Combinators

Combinatory Parsing

- ▶ A **parser** is a function that consumes input and returns some structure.
- ▶ A **parser combinator** is a higher-order function that takes parsers to form a new larger parser.
- ▶ The **parsing program** is the application of some basic parsers and parsing combinators.

Combinatory Parsing (cont.)

- ▶ Advantages
 - ▶ Readable
 - ▶ Written in the same programming language
 - ▶ Can be extended by the user with new parsers and combinators
 - ▶ Maintainable
- ▶ Disadvantages
 - ▶ No grammar analysis
 - ▶ No left-recursion
 - ▶ Some cases have high time and space usage

History of Combinatory Parsing

- ▶ Wadler, 85: *List of Successes* method
- ▶ Wadler, 92: Monadic Parsers
- ▶ Hutton, 92: basic parser combinators (in Miranda)
- ▶ Røjemo, 95: Applicative Parsers
- ▶ Hutton,Meijer, 96: Standardizing parser combinators, a Hugs parsing library
- ▶ Leijen,Meijer, 01: Parsec, fast parsing library with good error reporting
- ▶ Swierstra, 01: Applicative parsers with error correction
- ▶ Swierstra, 04: Permutation parsing
- ▶ Swierstra, 06: Online results parsing

Parsec2

- ▶ Written by Daan Leijen
- ▶ Introduced in 2001
- ▶ Became de-facto
- ▶ Used to be distributed with GHC
- ▶ Now distributed with Haskell Platform

Parsec2 Features

- ▶ Parses context-sensitive grammars
- ▶ Predictive LL(1) by default
- ▶ LL(*) by explicit backtracking
- ▶ Monadic Interface
- ▶ Permutation parsing
- ▶ Good error messages (Position, Unexpected Input, Expected production)
- ▶ Haskell98-compatible

Parsec2 Drawbacks

- ▶ Lacks an applicative interface
- ▶ Accepts a token-list as input
- ▶ Cannot be used with ByteString or Text

Parsec3

- ▶ In most cases, API-compatible to Parsec2
- ▶ Adds an applicative interface
- ▶ Parameterized over the input type
- ▶ Provides the parser as a monad transformer
- ▶ Breaks Haskell98-compatibility
- ▶ Trying to be flexible, it pays in terms of speed

UULib General

- ▶ Written by UU
- ▶ Introduced < 2005
- ▶ Stable
- ▶ The library bundles also a lexer and a prettyprinter

UULib Features

- ▶ Applicative Interface
- ▶ Good Error reporting

UULib Drawbacks

- ▶ Lacks a monadic interface \Rightarrow difficult to construct context-sensitive grammars
- ▶ Does not have `pSatisfy`. It makes difficult to construct some basic parsers.
- ▶ Not documented
- ▶ Restricted only to `String` input

UU-parsinglib General

- ▶ Written by Doaitse Swierstra
- ▶ Introduced in 2008
- ▶ Somewhat stable

UU-parsinglib Features

- ▶ LL(*)
- ▶ Applicative and Monadic Interface
- ▶ API-compatible with uulib
- ▶ Good at dealing with ambiguity
 - ▶ Guesses out when the grammar is ambiguous and reports it
 - ▶ `amb :: Parser a -> Parser [a]`
- ▶ Compared to the other libraries, it seems to be the best abstracted away
- ▶ Great Error messages
- ▶ Error Correction
- ▶ Online parsing results
- ▶ Visits the alternatives in parallel (BFS), so it does not retain the whole input in memory
- ▶ Input can also be ByteString, Text, etc
- ▶ Extremely easy to switch between input types (thanks to ListLike)

UU-parsinglib Drawbacks

- ▶ Error Correction cannot be used with the monadic interface
- ▶ Underdocumented
- ▶ The abstractions make it harder for beginners to understand the lib

Attoparsec General

- ▶ Written by Bryan O'Sullivan
- ▶ Introduced in 2007
- ▶ Has recently become very popular
- ▶ Starting to mature

Attoparsec Features

- ▶ LL(*) , implicit backtracking
- ▶ Applicative and Monadic Interface
- ▶ Supports Text, ByteString
- ▶ Supports Incremental Input
- ▶ Nearly translatable/compatible with Parsec
- ▶ Easy to use, by having distinct modules

Attoparsec Drawbacks

- ▶ Trying to stay simple, it lacks some high-level combinators
- ▶ Not good in error reporting
- ▶ Does not accepts String input

Polyparse General

- ▶ Written by Malcolm Wallace
- ▶ Introduced in 2002
- ▶ Is actually a collection of different parsing libraries

Polyparse Features

- ▶ $LL(*)$ with implicit backtracking
- ▶ Applicative and Monadic Interface
- ▶ Partial parsing results
- ▶ Supports String, ByteString and Text for input

Polyparse Drawbacks

- ▶ The library's schema is confusing
- ▶ Lacks some high-level combinators for Text

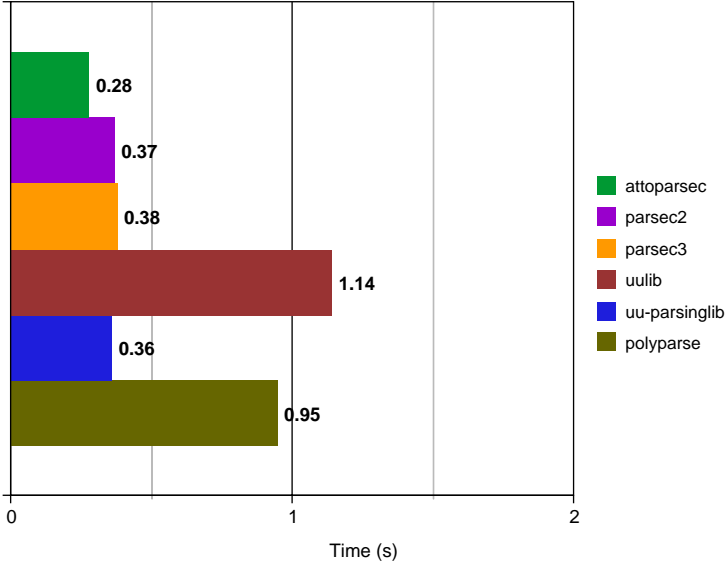
CSV - Example

- ▶ Comma-separated values
- ▶ Has a standard but many variations exist
- ▶ Our parsers support the standard with unicode
- ▶ Datasets taken from the net

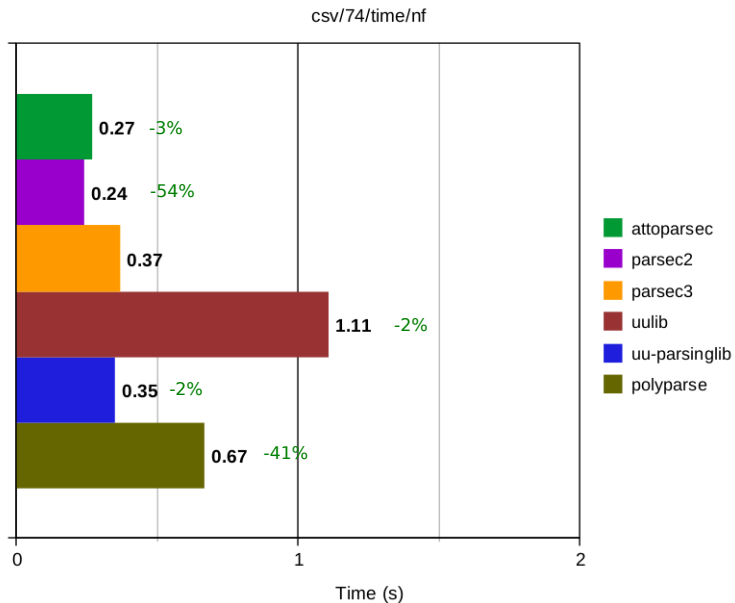
```
dma code,region,state
500,Portland-Auburn,ME
501,New York,NY
502,Binghamton,NY
503,Macon,GA
504,Philadelphia,PA
505,Detroit,MI
506,Boston,MA
507,Savannah,GA
508,Pittsburgh,PA
```

CSV - Time

csv/70/time/nf

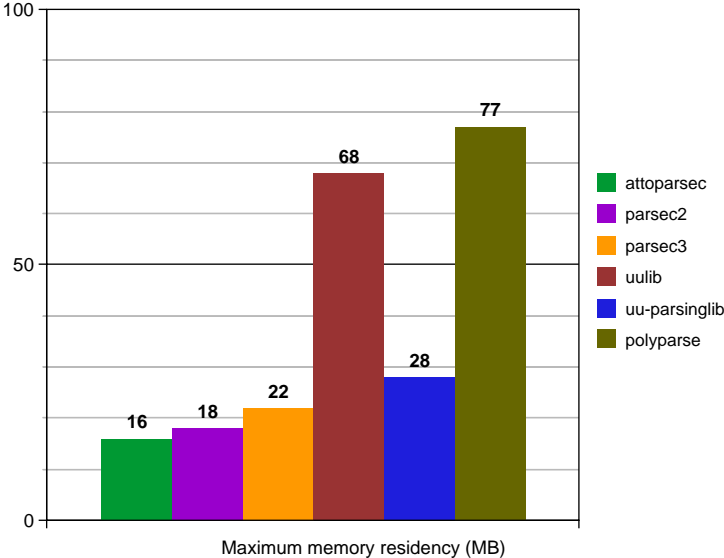


CSV - Time - 74



CSV - Space

csv/70/hp/nf

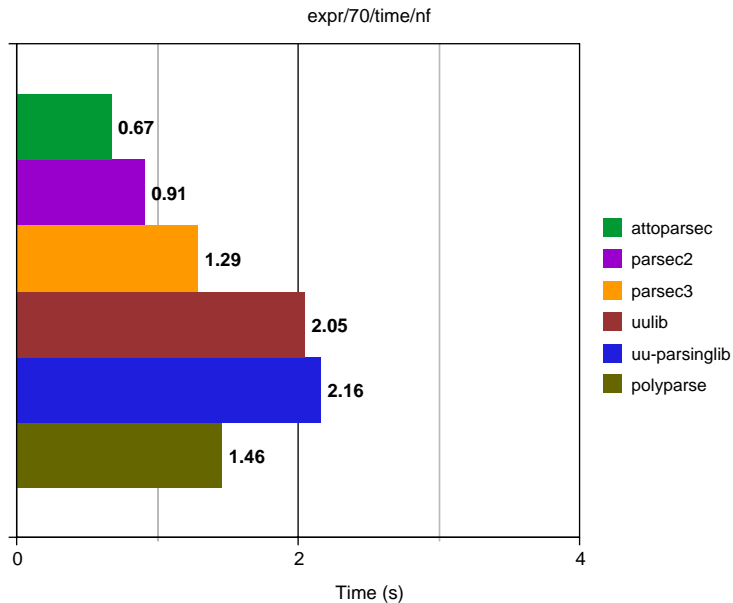


EXPR - Example

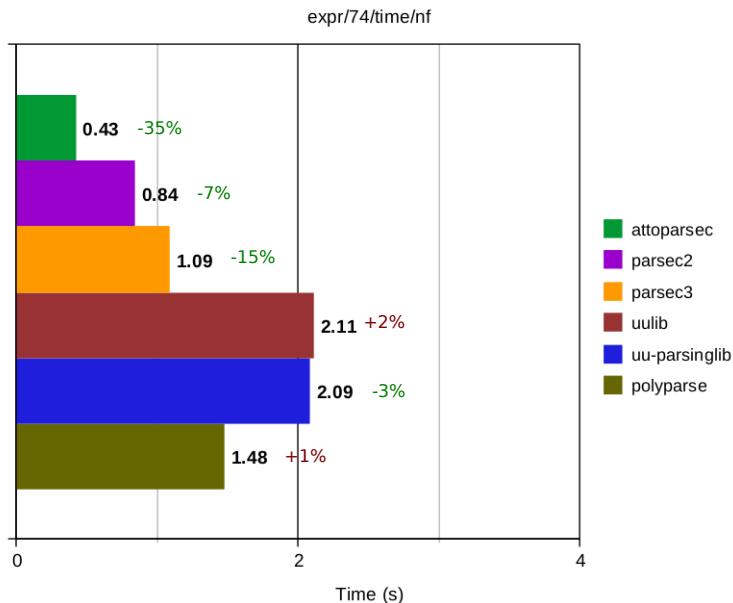
- ▶ a small custom functional language
- ▶ Datasets generated with QuickCheck's Arbitrary

```
((2 - let nngxpu = (kpjht + (lms + (medpc - 1))) in
(let sk = yui in o + ((qatzrb * 2) - (8 - 2)))) -
(((let idwxvj = let ywkldp = 5 in c in (5 * az)
* (8 * let rh = yo in h)) + (((awpcu - 2) - (tcwhem *
nyt)) * let pm = (5 * lttgyt) in (cl + 5))) -
((((bf - 3) * (7 * 7)) - ((7 * vpienh) + (dy + 6)))
+ (((8 + 6) - (7 * 7)) + let zvs = poks in (2 * sfwo))))))
```

EXPR - Time

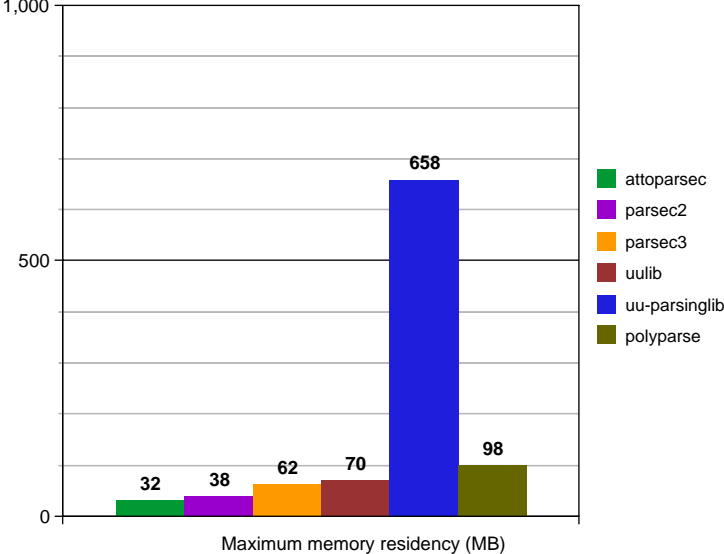


EXPR - Time - 74



EXPR - Space

expr/70/hp/nf



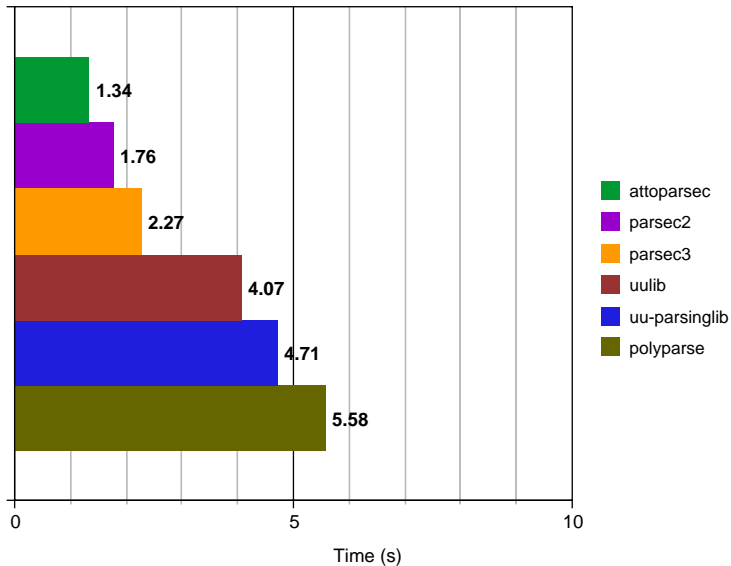
URL - Example

- ▶ URL query string
- ▶ Datasets generated from previous CSV

```
Aachen&Aalborg=aah&Aalesund=aahed&Aalst=aahing&Aalto=aahs  
&Aarau=aal&Aargau&Aarhus&Aaron=aals&Aarons+rod=aardvark  
&Aaronic=aardwolf&Ab=aargh&Abadan&Abaddon=aarrghh&Abba=  
aas&Abbasid=aasvogel&Abbevillian=aba&Abbey+Theatre=abaca&
```

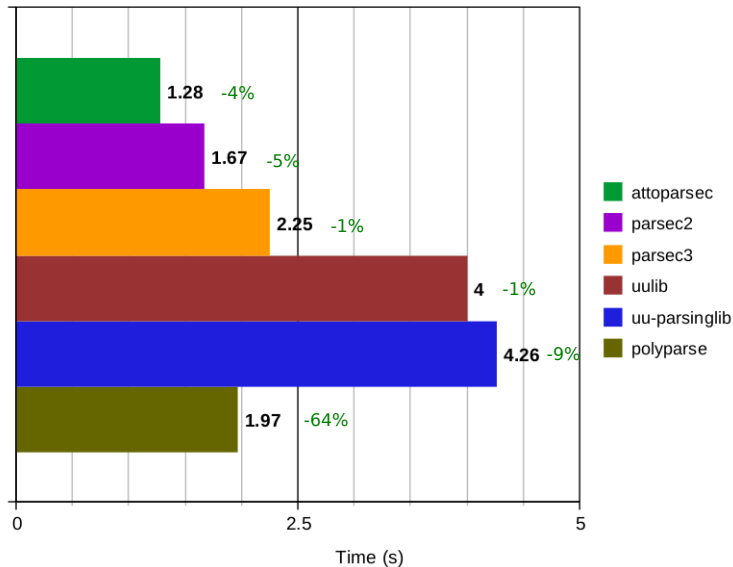
URL - Time

url/70/time/nf

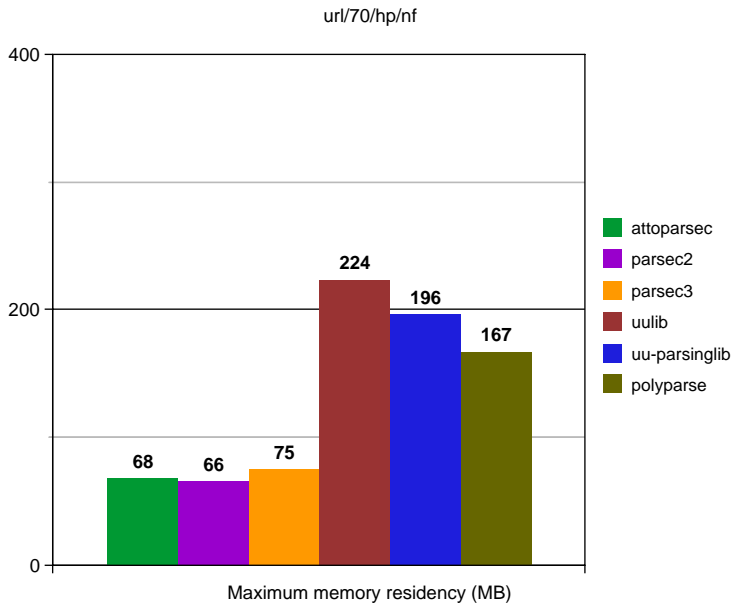


URL - Time - 74

url/74/time/nf



URL - Space



HTTP - Example

- ▶ HTTP GET and POST requests
- ▶ Datasets generated from previous CSV

```
POST http://www.w3.org/#sec3.6.1 HTTP/1.1
```

```
Aachen: aa
```

```
Aalborg: aah
```

```
Aalesund: aahed
```

```
Aalst: aahing
```

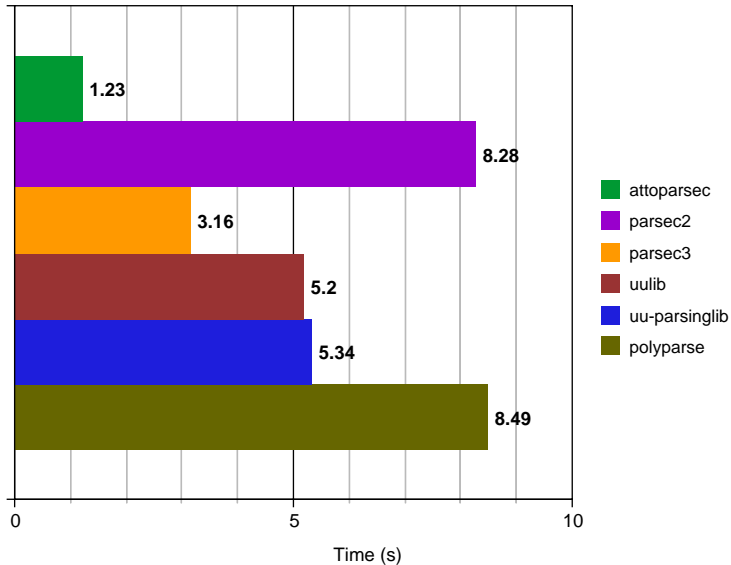
```
Aalto: aahs
```

```
Aarau: aal
```

```
Aargau: aalii
```

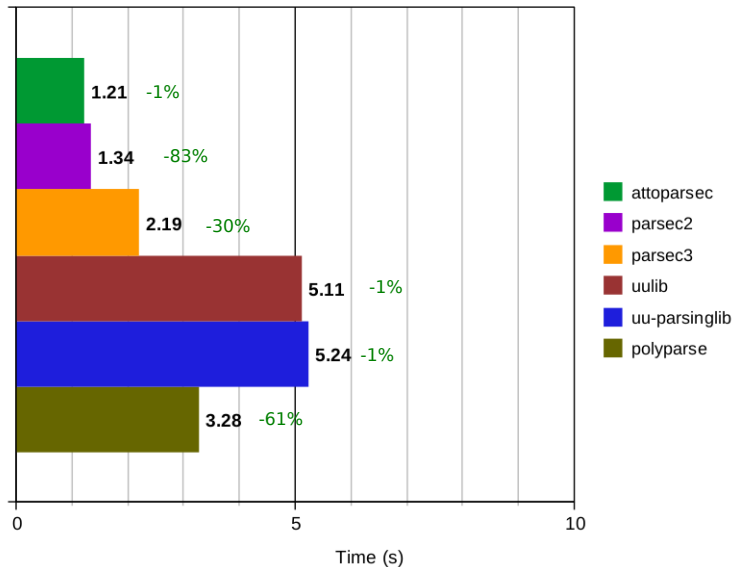
HTTP - Time

http/70/time/nf

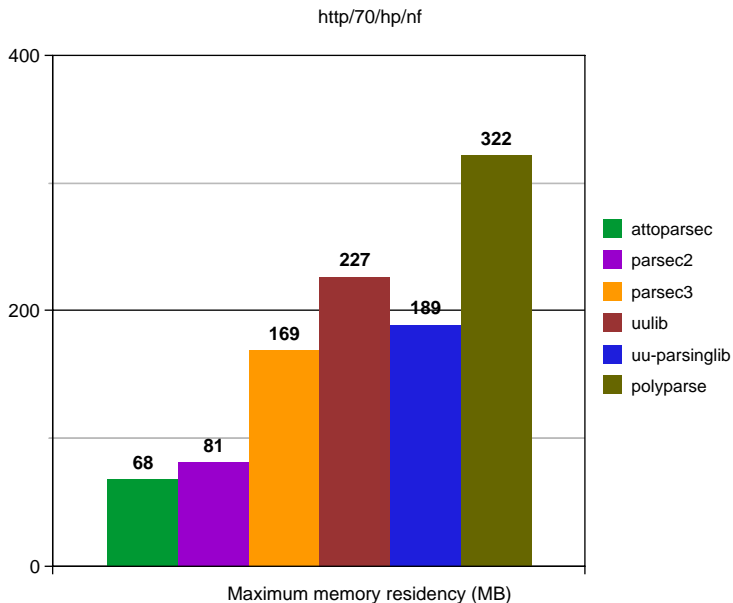


HTTP - Time - 74

http/74/time/nf



HTTP - Space

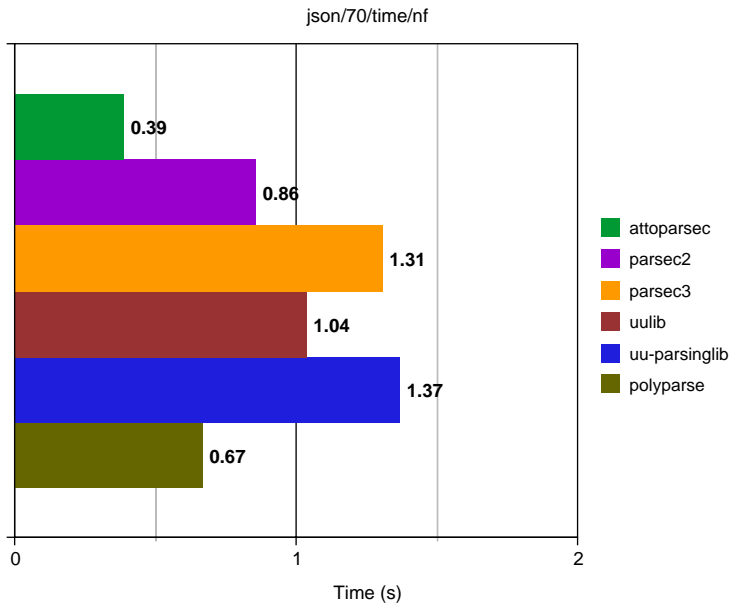


JSON - Example

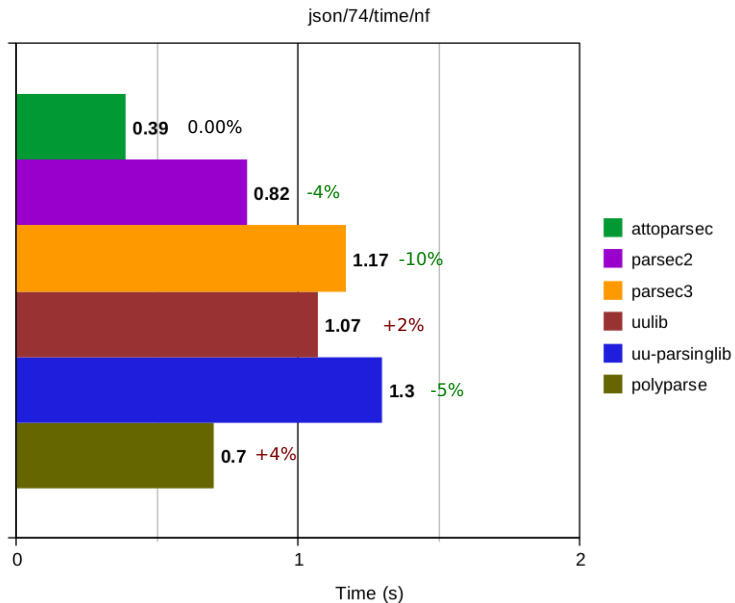
- ▶ Javascript Object Notation
- ▶ Datasets collected from GitHub

```
{"030200b_.mid": [[{"timesig": null, "keysig": 2, "st": 0, "pitch": 50, "dur": 12.0, "fermata": 0}, {"timesig": null, "keysig": 2, "st": 12.0, "pitch": 49, "dur": 2.0, "fermata": 0}]]}
```

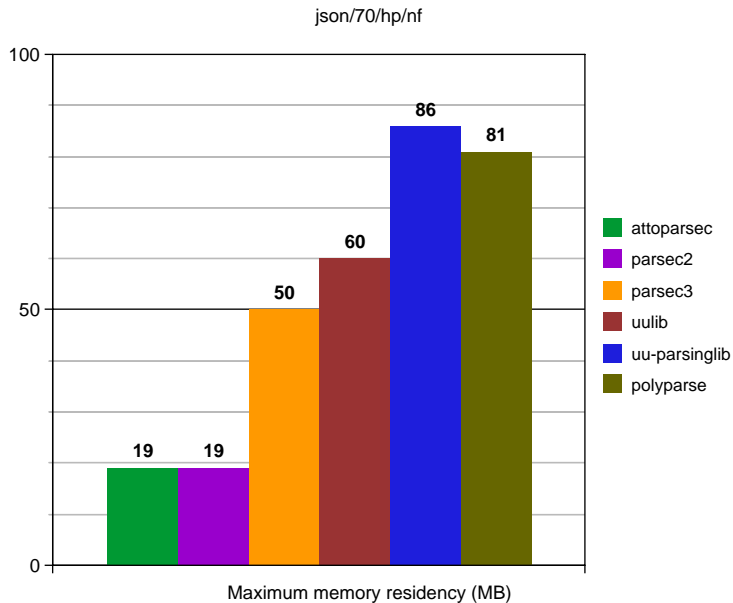
JSON - Time



JSON - Time - 74



JSON - Space

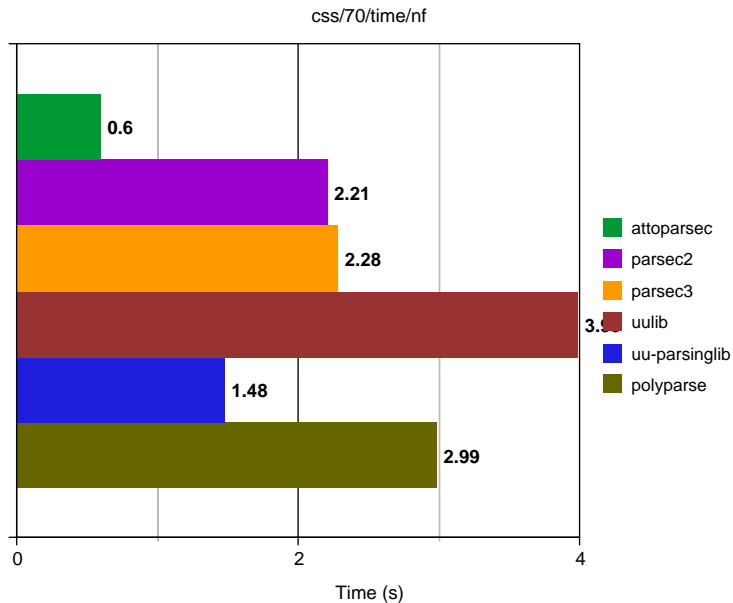


CSS - Example

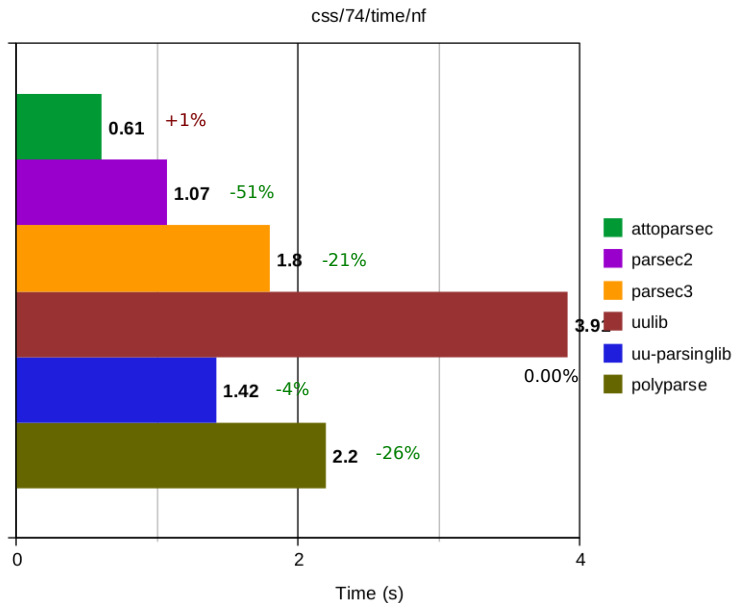
- ▶ Cascading Style Sheets
- ▶ Datasets collected from GitHub

```
body{
  background-color: #222;
  color: white;
  font-family: "Bitstream Vera Sans",sans-serif;
  margin: 0;
  padding: 0;
}
#container{
  margin: 0 10%;
  background-color: #161616;
  padding: 16px;
}
```

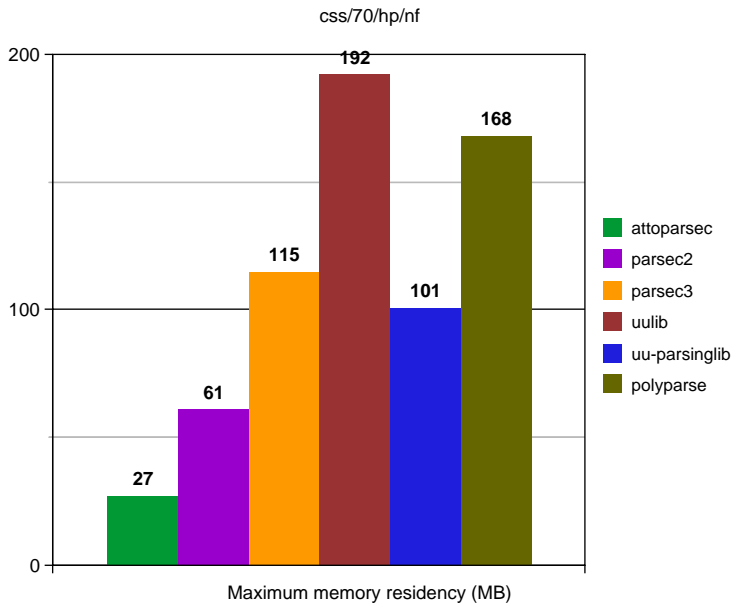
CSS - Time



CSS - Time - 74

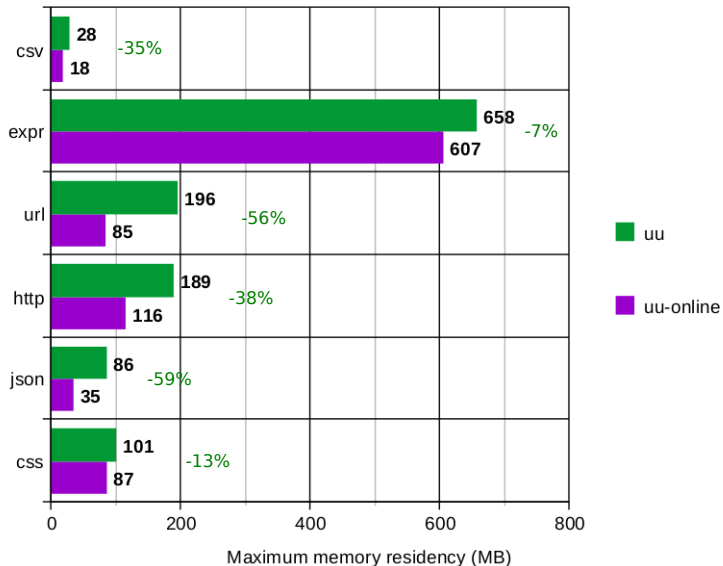


CSS - Space

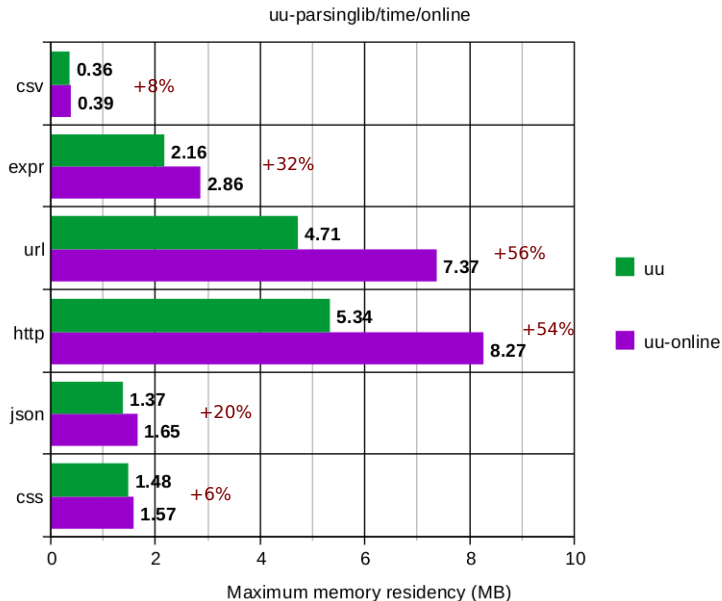


UU-parsinglib Online - Space

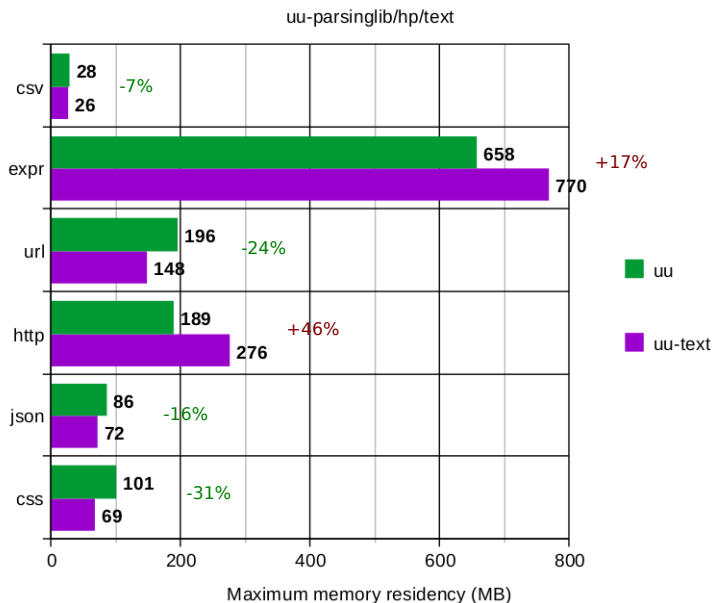
uu-parsinglib/hp/online



UU-parsinglib Online - Time



UU-parsinglib Data.Text - Space



Overall results

- ▶ GHC 7.4 brings a 15% performance increase
- ▶ WHNF 17% faster than NF on the test
- ▶ WHNF uses 16% less memory than NF

How the results were produced

- ▶ My slow laptop
- ▶ Archlinux x86_64
- ▶ GHC 7.0.3 and GHC 7.4.1
- ▶ Latest versions of parsing libraries
- ▶ Criterion
- ▶ Virtualenv
- ▶ Custom profiling scripts

Conclusion & Impressions

- ▶ attoparsec is the winner on almost all benchmarks
- ▶ parsec2 is faster than parsec 3 on all benchmarks (GHC 7.4)
- ▶ uu-parsinglib is the best abstracted away
- ▶ uu-parsinglib online feature is beneficial albeit slower
- ▶ was expecting a larger boost with Data.Text
- ▶ WHNF is slightly faster than NF for all contenstants, but does not affect the results ordering

Suggestions & Future Work

- ▶ Revise the parsers
- ▶ Try to squeeze even the last drop of speed
- ▶ Add more example cases
- ▶ It would be nice if there was a similar tool to Criterion but for profiling comparisons