

Fastest Protein Databank parser is written in Haskell

Michał J. Gajda

- Fast parsing intro
 - Tips and tricks
 - Benchmarks
- Library interface

Parser are used in many applications

- Compilers
- Code analyzers
- Text analyzers
- Mark-up processors (browsers, DocBook formatters)
- Low-level networking
- Data input/output.

Estimated 60-90% of code are parsers!

Protein DataBank

- Column and line-based file format.

```
HEADER      LIGASE                               01-JAN-01   1HTQ
TITLE      MULTICOPY CRYSTALLOGRAPHIC STRUCTURE OF A RELAXED GLUTAMINE
TITLE      2 SYNTHETASE FROM MYCOBACTERIUM TUBERCULOSIS
MODEL      1
ATOM       1  N   THR A 601      105.054  51.739 138.889  0.10 51.66      N
ATOM       2  CA  THR A 601      106.152  52.289 139.747  0.10 55.33      C
ATOM       3  C   THR A 601      107.533  51.719 139.344  0.10 77.58      C
```

- Deposition of protein and nucleic acid structures
- $1\text{\AA}=10^{-10}\text{ m}=0.1\text{nm}$ scale
- Over 10GB database.
- Parsed in under 15mins on quad core Ivy Bridge using hPDB

PugiXML – Parsing XML at speed of light

- mmap
- Avoid copying identifiers.
- Non-well-formed documents may sometimes be parsed.
- Normalizations and transformations are performed on-the-fly (entities, char. refs, newlines, attr. normaliz.).
- Char stream instead of token stream
- In-place strings, single-gap strings
- 256-byte tables + comparisons for ranges (UTF16, UTF32)
- Vectorized checks (SIMD up to 16-chars)
- Loop instead of recursion, with DOM node cursor stack
- Cold code shifting
- Null-terminated chunks
- Linked-list representation of DOM

GNU vs BSD grep

- **mmap**
- Avoid looking at all chars.
- For those that are looked up – 2-3 *i86* instructions.
- Own unbuffered output to avoid copying.
- Avoid in-kernel copying from realignment: page-aligned buffers, page-sized read chunks

Common motives

- Buffering (mmap)
- Zero-copy operations
- Number of decisions per character.
- Number of instructions per character (and vectorization).

PDB parsing

- Line is a record
- Fixed column-based fields
- Field columns depend on record type
- A lot of numbers to parse

Parser spends 60-80% of time in floating point conversions.

Practical tricks

- Zero-copy, mmap
- Instructions per byte
- Lookups/decisions per byte
- Loop check/loop unrolling
- Abstract syntax tree
- Cache and sequential lookahead

Parser structure

- `mmap` input, if can't the read the whole file as `ByteString`
- *decompress input*
- For each line:
 - Case on record type
 - Split record into fields
 - Convert types
 - Generate event
- Build `Bio.PDB.Structure` from events.
- Join fragments of the structure built in parallel.

Parallel parser structure

- `mmap` input, if can't the read the whole file as `ByteString`
- *decompress input*
- **For each input part:**
 - For each line:
 - Case on record type
 - Split record into fields
 - Convert types
 - Generate event
 - Build `Bio.PDB.Structure` from events.
- **Join fragments of the structure built in parallel.**

Tricks used

- Zero-copy input: **mmap**
- Preallocating residue's atom arrays
- Minimize lookups/decisions per byte
 - Two checks for most characters:
 - Newline check
 - Record type check
- `ByteStrings` point to the same memory
- Cache and sequential lookahead
- Using `double-conversion` library written for Google Chrome – 60-80% of total runtime

hPDB – Haskell faster than...

Table 1 - Total allocated memory in megabytes.

PDB entry	Input size	hPDB par.	hPDB seq.	BioRuby	BioJava	BioPython
1CRN	49 kB	3	1	8	240	206
3JYV	5	41	35	85	302	324
1HTQ	76	609	547	1350	1180	2409

Table 2 - Total CPU time in seconds.

PDB entry	hPDB par.	hPDB seq.	BioJava ¹	BioRuby	BioPython	PyMol	RasMol	Jmol ¹
1CRN	≤ 0.01	≤ 0.01	0.38	0.03	0.31	0.06	0.06	1.96
3JYV	0.27	0.26	1.31	0.89	1.26	0.28	0.28	3.52
1HTQ	5.08	4.63	6.66	16.52	23.41	3.94	4.90	25.82

¹ Jmol and BioJava use multiple threads, thus completion time is closer to half the CPU time than to the sum of CPU time and I/O time (as indicated in table 3).

Table 3 - Completion time after parsing in seconds.

PDB entry	hPDB par.	hPDB seq.	BioJava	BioRuby	BioPython	PyMol ²	RasMol ²	Jmol ²
1CRN	≤0.01	≤0.01	0.23	0.04	0.32	0.14	0.77	2.26
3JYV	0.09	0.28	0.71	0.94	1.43	0.38	0.86	2.81
1HTQ	1.39	4.79	3.24	17.14	24.01	4.22	5.73	12.86

² Includes the time needed for startup and closing the window.

hPDB reference

hPDB - Haskell library for processing atomic biomolecular structures in Protein Data Bank format
BMC Research Notes 2013, 6:483 DOI:[10.1186/1756-0500-6-483](https://doi.org/10.1186/1756-0500-6-483)

Library interface

- class Iterable a b where

- itmap :: (b → b) → a → a

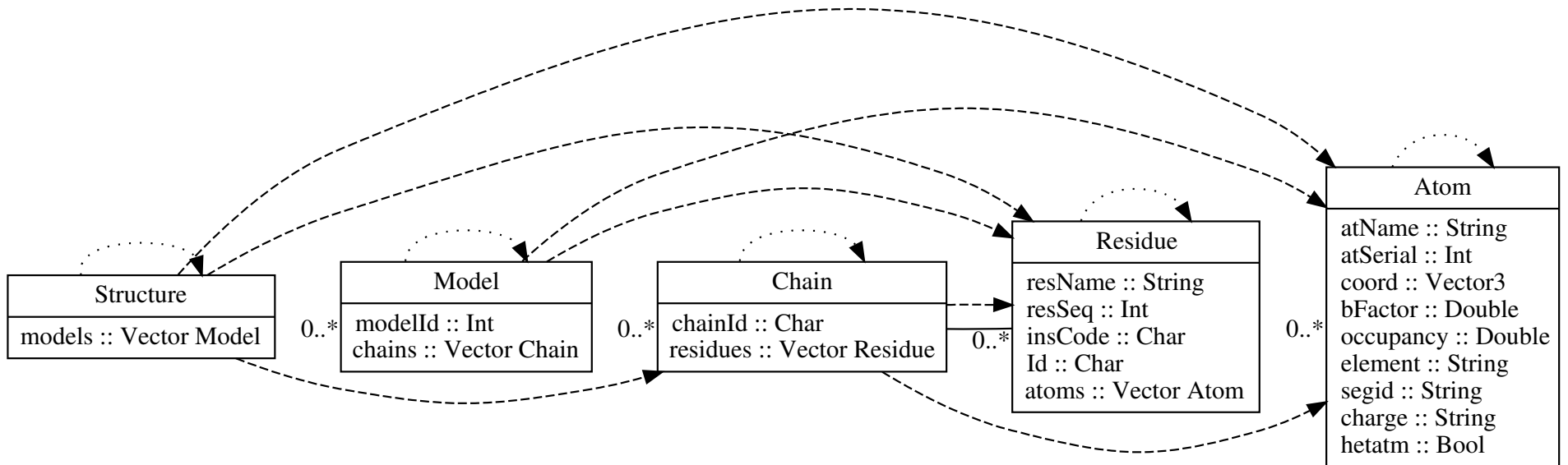
- itfoldr :: (b → c → c) → c → a → c

- itfoldl :: (c → b → c) → c → a → c

- ...

- itlength :: b → a → Int

Instance hierarchy



Example: centering

```
center :: Structure → Vector3
center s = avgVec
  where
    sumv = itfoldl' addCoord 0 (s :: Structure)
    n = fromIntegral $ numAtoms s
    addCoord v (Atom { coord = c }) = v+c
    avgVec = (1/realToFrac n) *| sumv
```

```
shift :: Double → Structure → Structure
shift v = itmap (\at →
                 at { coord = coord at - v })
```

Example – steric clash check

```
import qualified Data.Octree           as Oct
import           Bio.PDB               as PDB
import qualified Bio.PDB.Structure.Elements as PDB(vanDerWaalsRadius)

clashCheck s1 s2 = filter (/= []) . Prelude.map clashes $ itfoldr (:) [] s2
  where
    clashes (at :: PDB.Atom) = Oct.withinRange ot (radius + maxRadius) (PDB.coord
at)
      where
        radius :: Double = realToFrac . PDB.vanDerWaalsRadius . PDB.element $ at
        ot :: Oct.Octree (Int, Double)
        ot = makeOctree s1

extract :: PDB.Atom -> (Oct.Vector3, (Int, Double))
extract (PDB.Atom { coord = cvec, atSerial = ser , element = elt }) =
  (cvec, (ser, realToFrac $ PDB.vanDerWaalsRadius elt))

makeOctree structure = Oct.fromList . Prelude.map extract . itfoldr (:) []
  $ structure

main = do [input1, input2] <- Env.getArgs
  Just structure1 <- PDB.parse input1
  Just structure2 <- PDB.parse input
  print $ clashCheck structure1 structure2
```


Join

<http://www.biohaskell.org!>

- Open-source bioinformatic library for Haskell
 - Sequence, alignment parsing
 - RNA secondary structure
 - PDB, BMRB for 3D processing
- Fast!!!

✉ Mail us to request features!

biohaskell@biohaskell.org

hPDB reference

hPDB - Haskell library for processing atomic biomolecular structures in Protein Data Bank format

BMC Research Notes 2013, 6:483 DOI:[10.1186/1756-0500-6-483](https://doi.org/10.1186/1756-0500-6-483)