

# Design and Implementation of a MongoDB Driver for Prolog

Sebastian Lundström

April 11, 2011

## 1 Diary

- March 29. Wrote a decoder that decodes `{hello: world}` and `{hello: 32}`. Learned some DCG and rewrote the parser with it.
- March 30. Researched how SWI interacts with C. Wrote a `bytes8_to_double` function in C and managed to integrate it with the rest of the system. Separated bson into separate files, encoder and decoder. Should think about researching PIUnit by now. Thought about and tried PIUnit. Works fine. But how to organize tests? Implemented all bit-hacking in C for the time being. Works. Simple. Fun. Started looking at PIDoc. Generating actual docs seems like a hassle, probably won't do it. But the important thing is unit tests, and some comments inside the source. Started fleshing out the report slightly, with some ideas and stubs.
- March 31. Restructured project folder. Better including of modules now, I think. Starting to write nice unit tests. Continued with the decoder a little, the complex BSON example works. Wrote nice TextMate snippets for test boilerplating. Rewrote Makefile, yet again. Added make test, for instance.
- April 1. When trying to fix proper decoding of strings I ran into issues with Unicode. Finally found SWI-Prolog's memory files, which can be written to and then read back using a different encoding. Five line fix. Implemented a variant that might be more efficient. Might. (And like five times longer, but clean though.) Cleaned up the tests. Implementing the

rest of the decoder should be straight-forward now. Next week. It. Will. Be. Awesome.

- April 2. Looked into the Prolog module system again and started adding prefixes (“memory\_file:”, “builtin:”) to SWI predicate calls. “builtin” isn’t an actual module, but since it works anyway (it enters the module builtin and calls the predicate from there) I decided to use it. Documented predicates in bson\_bits. Fixed a bug with bytes\_to\_integer/5 because it wasn’t using a 32-bit integer. Now it uses int32\_t. Started using maplist/2. Found a library(apply\_macros) that claims to speed up maplist etc. simply by loading it.
- April 3. All use\_module directives now use an empty list in order to ensure that the module is loaded, but nothing is imported. This forces me to add the appropriate module prefixes to calls. Benchmarked the unicode converter and the old shorter code outperformed the new longer code with more than a factor of three. I guess this is due to more stuff being done by the C library and not Prolog. Anyway, the old code is faster and much shorter – win-win. Also replaced a crappy to-codes-and-then-to-atom conversion with a straight-to-utf8-atom library routine, which made it even faster. A keeper.
- April 4. Refactored decoder slightly. Added relation term\_bson/2 to bson, which calls necessary subpredicate. Moved tests into separate files: they were getting too long, and separating them makes it easier to write them simultaneously. As long as tests and implementation are close to each other, I am okay. (Impl in .pl and tests in .plt next to each other.) The decoder should now be feature complete. It does minimal error checking, but it should be able to parse all element types. It might still need some cleaning though. Added version predicates to bson. Made comments more in line with PIDoc.
- April 5. Major decoder cleanup. Took a slow day.
- April 6. Started on the encoder. Unicode issues again. Solved it after an hour or so, using even more library predicates. Factored out Unicode handling to separate module, used by both encoder and decoder. Worked a lot on the encoder, including some C hacking again.

- April 7. Lots more encoding of elements and encoder refactoring. More to come!
- April 8. Clean up a lot of bit-hacking in Prolog. Made entry-predicates into true relations and made the actual conversions into more general list predicates. Changed “builtin:” prefix to “inbuilt:”, which looks nicer and is more of a single word. Real cleanup in all the bit-hacking. It sure took most of the day, but things should be pretty solid and clean now, and the strategy is flexible but straight-forward: all integer conversions back-and-forth between signed big/little 32/64-bit and bytes is done in C. If those routines cannot be used, it fails if a negative integer is being encoded, otherwise it is treated as an unbounded unsigned, making for instance ObjectID simple. Everything goes through either the relation `FLOAT_BYTES(FLOAT, BYTES)` which is just mappings to C, or the more intricate relation `INTEGER_BYTES(INTEGER, NUMBYTES, ENDIAN, BYTES)` which can be called like: `(+,+,+,?)` or `(?,+,+,+)`. E.g. `(1, 4, little, Bytes)` gives `Bytes = [1,0,0,0]`, or `(Integer, 4, big, [0,0,0,1])` gives `Integer = 1`. The first example can be read like “the integer 1 in 4 little(-endian) bytes.”
- April 9. Factored out the unbounded unsigned handling into a separate predicate `UNSIGNED_BYTES(UNSIGNED, NUMBYTES, ENDIAN, BYTES)`. Now things are starting to look really good. `float_bytes/2` handles doubles, `integer_bytes/4` handles all signed 32/64-bit integers and `unsigned_bytes/4` handles all unbounded unsigneds. Simple. Also documented the predicates.
- April 10. Made all Unicode handling atom-only. The code became much less confusing, and using code lists for text would be bad anyway because it can be ambiguous (is `[97,98,99]` a list of numbers or `'abc'?`). Yet another implementation of `utf_to_bytes/2`. Uses less weird predicates (no `open_chars_stream` any more) and seems a tiny bit faster. The code is a bit longer though (but arguably more straight-forward).
- April 11. Encoder is now feature complete. Tomorrow I will start thinking about the report and maybe start investigating sockets.

## 2 Background

From spec:

[MongoDB is a young document-oriented database system that has started to gain much attention recently. Document-orientation involves removing rigid database schemas and advanced transactions, in favor of flexibility. Document-orientation also promotes a certain degree of denormalization which allows embedding documents into each other, leading to potentially much better performance by avoiding the need for expensive join operations.

Prolog, being an untyped language, agrees with the document-oriented approach of relaxing manifests in order to create more dynamic and flexible systems. Embedding terms in other terms is natural in Prolog, and embedding documents in other documents is natural in MongoDB.

Many drivers exist, both official and unofficial, that enable the use of MongoDB from various programming languages. At the time of writing, no such driver for Prolog seems to exist.]

## 3 Method

Research other drivers, docs.

Test-driven development.

## 4 Result

- BSON encoder/decoder
  - Some parts written in C. Why? (Basically didn't know how to easily handle bytes-to-float in Prolog.)
  - Discuss data structures, term [key:value] maps to bytelist [4,1,7,9,3,...] etc.
  - Design choices: text as “” and symbol as atom
- Network communication
  - How does the communication work? Out of scope?
- MongoDB API

- Thoroughly discuss design of wrapper API, how lists and structures are represented etc.

## 5 Discussion

- Portability
  - Tested on Mac, SWI, GCC, etc.
- Efficiency
- Future
  - BSON implementation in C? Already exists, I think.
  - More advanced features

## 6 Conclusion

Difficulties? Did it work? Is it usable?

## 7 Literature

Chodorow, K. & Dirolf, M. (2010) *MongoDB: The Definitive Guide*. Sebastopol, United States of America: O'Reilly Media, Inc.