# Comparison of serialization formats

Created by 김동구 부장(실장) on May 03, 2022

## Pros & Cons per serialization format

| name | Description | Pros | Cons | type |
|---|---|---|---|---|
| Apache Avro | • Row-oriented remote procedure call and data serialization framework developed within Apache's Hadoop project.<br>• It uses JSON for defining data types and protocols, and serializes data in a compact binary format.<br>• uses a schema to structure the data that is being encoded<br>• two different types of schema language: Avro IDL(human editing), machine-readable based on JSON<br>• be similar to Thrift and Protocol Buffers | • a linguistic-neutral serialization of data<br>• stores the schema in a file header, so the data is self-describing<br>• Easy and fast data serialization and deserialization, which can provide very good ingestion performance<br>• As with the Sequence files, the Avro files also contain synchronization markers to separate blocks. This makes it highly splittable<br>• Files formatted in Avro are splittable and compressible and are therefore a good candidate for data storage in the Hadoop ecosystem<br>• The schema used to read Avro files does not necessarily have to be the same as the one used to write the files. | • Makes you think about the schema and data types<br>• Its data is not human-readable<br>• Not integrated into every programming language | Schema-driven<br><br>ex) cpx.json<br><br>cpx is a C++ representation of the Avro schema.<br><br>{<br>  "type": "record",<br>  "name": "cpx",<br>  "fields" : [<br>    {"name": "re", "type": "double"},<br>    {"name": "im", "type" : "double"}<br>  ]<br>} |
| Apache Parquet | | • columnar format. Only the required columns will be retrieved/read, this reduces disk I/O. The concept is called projection pushdown<br>• The scheme travels with the data, so the data is self-describing<br>• Although it is designed for HDFS, data can be stored on other file systems such as GlusterFs or NFS<br>• just files, which means it's easy to work, move, backup and replicate them<br>• provides very good compression up to 75% when using even compression formats like snappy<br>• As practice shows, this format is the fastest for read-heavy processes compared to other file formats<br>• well suited for data storage solutions where aggregation on a particular column over a huge set of data is required<br>• can be read and written using the Avro API and Avro Schema<br>• provides predicate pushdown, thus reducing the further cost of transferring data from storage to the processing engine for filtering | • The column-based design makes you think about the schema and data types<br>• Parquet does not always have built-in support in tools other than Spark<br>• not support data modification (Parquet files are immutable) and scheme evolution | Schema-driven<br><br>ex)<br>enum PhoneType {<br>  HOME,<br>  WORK,<br>  MOBILE,<br>  OTHER<br>}<br>struct Phone {<br>  1: i32 id,<br>  2: string number,<br>  3: PhoneType type<br>}<br>service PhoneService {<br>  Phone findById(1: i32 id),<br>  list<Phone> findAll()<br>} |
| FlatBuffers | • Free software library implementing a serialization format similar to Protocol Buffers, Thrift, Apache Avro, SBE, Cap'n Proto | • Supports "zero-copy" deserialization | • handling of FlatBuffers requires usually more code, and some operations are not possible (like some mutation operations). | Schema-driven |
| JSON | • open standard file format and data interchange format<br>• language-independent data format | • supports hierarchical structures, simplifying the storage of related data in a single document and presenting complex relationships<br>• Most languages provide simplified JSON serialization libraries or built-in support for JSON serialization/ deserialization<br>• supports lists of objects, helping to avoid chaotic list conversion to a relational data model<br>• widely used file format for NoSQL databases such as MongoDB, Couchbase and | • consumes more memory due to repeatable column names<br>• Poor support for special characters<br>• not very splittable<br>• lacks indexing<br>• less compact as compared to over binary formats | Schema-less |

| name | Description | Pros | Cons | type |
|---|---|---|---|---|
| | | Azure Cosmos DB<br>• Built-in support in most modern tools | | |
| OpenDDL<br>(OPen Data Description Language) | • general-purpose, human-readable, and strongly-typed data language for information exchange<br>• can be used for a wide range of applications that include everything from simple configuration files to the exchange of complex information among many programs in an editable format.<br>• It was originally developed as the basis for the Open Game Engine Exchange format. | | • library license is GPL 3.0 | |
| Protocol Buffers(protobuf) | • Free and open-source cross-platform data format used to serialize structured data | • Data is fully typed.<br>• Protobuf supports schema evolution and batch/streaming processing.<br>• Mainly used to serialize data, like AVRO.<br>• Having a predefined and larger set of data types, messages serialized on Protobuf can be automatically validated by the code that is responsible to exchange them. | • not human-readable and human-editable<br>• not good for the purposes of storing something e.g. a text document, or a database dump<br>• Not splittable and not compressible.<br>• No Map Reduce support.<br>• Require a reference file with the schema.<br>• Not designed to handle large messages. Since it doesn't support random access, you'll have to read the whole file, even if you only want to access a specific item. | Schema-driven |
| Apache Thrift | • Interface definition language and binary communication protocol used for defining and creating services for numerous programming languages | • Cross-language serialization with lower overhead than alternatives such as SOAP due to use of binary format.<br>• No XML configuration files.<br>• The language bindings feel natural. For example, Java uses ArrayList<String>. C++ uses std::vector<std::string>.<br>• The application-level wire format and the serialization-level wire format are cleanly separated. They can be modified independently.<br>• The predefined serialization styles include: binary, HTTP-friendly and compact binary.<br>• Doubles as cross-language file serialization.<br>• Soft versioning[clarify] of the protocol. Thrift does not require a centralized and explicit mechanism like major-version/minor-version. Loosely coupled teams can freely evolve RPC calls.<br>• No build dependencies or non-standard software. No mix of incompatible software licenses. | • Long coding time<br>• Large volume after coding<br>• no documentation<br>• not support server push | Schema-driven |
| Smile | • Computer data interchange format based on JSON<br>• Format was specified in 2010 by Jackson JSON processor development team. | • More compact and more efficient to read and write | | Schema-less |
| XML | | • Supports batch/streaming processing.<br>• Stores meta data along the data and supports the schema evolution.<br>• Being a verbose format, It provides a good ratio of compression compare to JSON file or other text file format. | • Not splittable since XML has an opening tag at the beginning and a closing tag at the end. You cannot start processing at any point in the middle of those tags.<br>• The redundant nature of the syntax causes higher storage and transportation cost when the volume of data is large.<br>• XML namespaces are problematic to use. The support of namespace can be difficult to correctly implement in an XML parser.<br>• Difficulties to parse requiring the selection and the usage of an appropriate DOM or SAX parser. | |
| MessagePack | • Computer data interchange format<br>• binary form for representing simple | • more compact than JSON<br>• allows binary data and non-UTF-8 encoded strings<br>• any type can be a map key, | • imposes limitations on array and integer sizes<br><br>    a value of an Integer object is | Schema-less<br><br>ex) |

| name | Description | Pros | Cons | type |
|------|-------------|------|------|------|
| | data structures like arras and associative arrays. | <ul><li>including types like maps and arrays, and, like YAML, numbers</li><li>more space-efficient than BSON</li><li>MessagePack is designed for efficient transmission over the wire</li></ul> | limited from -(2^63) upto (2^64)-1<br>maximum length of a Binary object is (2^32)-1<br>maximum byte size of a String object is (2^32)-1<br>String objects may contain invalid byte sequence and the behavior of a deserializer depends on the actual implementation when it received invalid byte sequence Deserializers should provide functionality to get the original byte array so that applications can decide how to handle the object<br>maximum number of elements of an Array object is (2^32)-1<br>maximum number of key-value associations of a Map object is (2^32)-1<br><ul><li>MessagePack returns only a dynamically typed data structure and provides no automatic structure checks.</li></ul> | ```cpp\n#include <msgpack.hpp>\n#include <string>\n#include <iostream>\n#include <sstream>\n\nint main()\n{\n    msgpack::type::tuple<int, bool, std::string> src(1, true, "example");\n\n    // serialize the object into the buffer.\n    // any classes that implements write(const char*,size_t) can be a buffer.\n    std::stringstream buffer;\n    msgpack::pack(buffer, src);\n\n    // send the buffer ...\n    buffer.seekg(0);\n\n    // deserialize the buffer into msgpack::object instance.\n    std::string str(buffer.str());\n\n    msgpack::object_handle oh = msgpack::unpack(str.data(), str.size());\n\n    // deserialized object is valid during the msgpack::object_handle instance is alive.\n    msgpack::object deserialized = oh.get();\n\n    // msgpack::object supports ostream.\n    std::cout << deserialized << std::endl;\n\n    // convert msgpack::object instance into the original type.\n    // if the type is mismatched, it throws msgpack::type_error exception.\n    msgpack::type::tuple<int, bool, std::string> dst;\n    deserialized.convert(dst);\n\n    // or create the new instance\n    msgpack::type::tuple<int, bool, std::string> dst2 = deserialized.as<msgpack::type::tuple<int, bool, std::string> >();\n\n    return 0;\n}\n``` |
| Cap'n Proto | <ul><li>Data serialization format and Remote Procedure Call(RPC) framework for exchanging data between computer programs</li><li>The Cap'n Proto interface schema uses a C-like syntax and supports common primitives data types (booleans, integers, floats, etc.), compound types (structs, lists, enums), as well as generics and dynamic types.</li><li>supports Object Oriented features such as multiple inheritance, which has been criticized for its complexity.</li></ul> | <ul><li>Incremental reads: It is easy to start processing a Cap'n Proto message before you have received all of it since outer objects appear entirely before inner objects (as opposed to most encodings, where outer objects encompass inner objects).</li><li>Random access: You can read just one field of a message without parsing the whole thing.</li><li>mmap: Read a large Cap'n Proto file by memory-mapping it. The OS won't even read in the parts that you don't access.</li><li>Inter-language communication: Calling C++ code from, say, Java or Python tends to be painful or slow. With Cap'n Proto, the two languages can easily operate on the same in-memory data structure.</li><li>Inter-process communication: Multiple processes running on the same machine can share a Cap'n Proto message via shared memory. No need to pipe data through the kernel. Calling another process can be just as fast and easy as calling another thread.</li><li>Arena allocation: Manipulating Protobuf objects tends to be bogged down by memory allocation, unless you are very careful about object reuse. Cap'n Proto objects are always allocated in an "arena" or "region" style, which is faster and promotes cache locality.</li><li>Tiny generated code: Protobuf generates dedicated parsing and serialization code for every message type, and this code tends to be enormous. Cap'n Proto generated code is smaller by an order of magnitude or more. In fact, usually it's no more than some inline accessor methods!</li><li>Tiny runtime library: Due to the simplicity of the Cap'n Proto format, the runtime library can be much smaller.</li><li>Time-traveling RPC: Cap'n</li></ul> | | Schema-driven<br><br>ex)<br>@0xdbb9ad1f14bf0b36; # unique file ID, generated by `capnp id`<br><br>```capnp\nstruct Person {\n    name @0 :Text;\n    birthdate @3 :Date;\n    email @1 :Text;\n    phones @2 :List(PhoneNumber);\n\n    struct PhoneNumber {\n        number @0 :Text;\n        type @1 :Type;\n        enum Type {\n            mobile @0;\n            home @1;\n            work @2;\n        }\n    }\n}\n\nstruct Date {\n    year @0 :Int16;\n    month @1 :UInt8;\n    day @2 :UInt8;\n}\n``` |

| name | Description | Pros | Cons | type |
|---|---|---|---|---|
| | | Proto features an RPC system that implements time travel such that call results are returned to the client before the request even arrives at the server! | | |
| CSV(Comma-Separated Values) | • With only commas separating values, CSV is very concise for a human-readable format. | • human-readable and easy to edit manually<br>• provides a simple scheme<br>• can be processed by almost all existing applications<br>• easy to implement and parse<br>• compact | • allows you to work with flat data. Complex data structures have to be processed separately from the format<br>• No support for column types. No difference between text and numeric columns<br>• no standard way to present binary data<br>• Problems with CSV import (for example, no difference between NULL and quotes<br>• Poor support for special characters<br>• Lack of a universal standard | |
| YAS<br>(Yet Another Serialization) | • is created as a replacement of boost.serialization<br>• require C++ 11 support | • fast<br>• header only library<br>• YAS binary archives is endian independent. | • To do:<br>   protobuf/messagepack support<br>   limits<br>   objects versioning<br><br>• weak documentation | ex)<br>https://github.com/niXman/yas/blob/master/tests/base/include/serialize.hpp |
| cereal | • a header-only C++ 11 serialization library | • fast and compact<br>• comes with full support for c++11<br>• smart pointers (things like std::shared_ptr and std::unique_ptr) are supported.<br>• unit tested<br>• supports binary serialization, XML serialization, and JSON serialization | • a consequence of this raw pointers and references are not supported. | Schema-less<br><br>ex)<br><br>`#include <cereal/types/unordered_map.hpp>`<br>`#include <cereal/types/memory.hpp>`<br>`#include <cereal/archives/binary.hpp>`<br>`#include <fstream>`<br><br>`struct MyRecord`<br>`{`<br>`  uint8_t x, y;`<br>`  float z;`<br>`  template <class Archive>`<br>`  void serialize( Archive & ar )`<br>`  {`<br>`    ar( x, y, z );`<br>`  }`<br>`};`<br><br>`struct SomeData`<br>`{`<br>`  int32_t id;`<br>`  std::shared_ptr<std::unordered_map<uint32_t, MyRecord>> data;`<br>`  template <class Archive>`<br>`  void save( Archive & ar ) const`<br>`  {`<br>`    ar( data );`<br>`  }`<br>`  template <class Archive>`<br>`  void load( Archive & ar )`<br>`  {`<br>`    static int32_t idGen = 0;`<br>`    id = idGen++;`<br>`    ar( data );`<br>`  }`<br>`};`<br><br>`int main()`<br>`{`<br>`  std::ofstream os("out.cereal", std::ios::binary);`<br>`  cereal::BinaryOutputArchive archive( os );`<br>`  SomeData myData;`<br>`  archive( myData );`<br>`  return 0;`<br>`}` |

| name | Based on | Specification | Binary | Human-readable | Standard APIs | Supports zero-copy operations | license | Requirements |
|---|---|---|---|---|---|---|---|---|
| Apache Avro | N/A | Apache Avro Specification | Yes | Partial | C, C#, C++, Java, PHP, Python, Ruby | N/A | Apache License 2.0 | • c++ compiler and runtime libraries<br>• Boost library version 1.38 or later<br>• CMake build tool version 2.6 or later<br>• Python |
| Apache Parquet | N/A | Apache Parquet | Yes | No | Java, Python, C++ | No | Apache License 2.0 | |
| FlatBuffers | N/A | FlatBuffers GitHub | Yes | Apache Arrow | C++, Java, C#, Go, Python, Rust, JavaScript, PHP, C, Dart, Lua, | Yes | Apache License 2.0 | |

| name | Based on | Specification | Binary | Human-readable | Standard APIs | Supports zero-copy operations | license | Requirements |
|---|---|---|---|---|---|---|---|---|
| | | | | | TypeScript | | | |
| JSON | JavaScript syntax | STD90 / RFC 8259(ancillary: RFC 6901, 6902), ECMA-404, ISO/IEC 21778:2017 | No, but see BSON, Smile, UBJSON | Yes | Partial(Clarinet, JSONQuery/RQL, JSONPath), JSON-LD | No | JSON license | |
| OpenDDL | C, PHP | OpenDDL.org | No | Yes | OpenDDL library(implemented c++) | N/A | GPL 3.0 | |
| Protocol Buffers(protobuf) | N/A | Developer Guide: Encoding | Yes | Partial | C++, Java, C#, Pytho, Go, Ruby, Objective-C, C, Dart, Perl, PHP, R, Rust, Scala, Swift, Julia, Erlang, D, Haskell, Action Script, Delphi, Elixir, Elm, Erlang, GopherJS, Haskell, Haxe, JavaScript, Kotlin, Lua, Matlab, Mercurt, OCaml, Prolog, Solidity, Typescript, Vala, Visual Basic | No | BSD | |
| Apache Thrift | N/A | Original whitepaper | Yes | Partial | C++. Java, Python, PHP, Ruby, Erlang, Perl, Haskell, C#, Cocoa, JavaScript, Node.js, Smalltalk, OCaml, Delphi and other languages | N/A | Apache License 2.0 | |
| Smile | JSON | Smile Format Specification | Yes | No | Partial(via JSON APIs implemented with Smile backend, on Jackson, Python) Java, C, Go, Javascript, Python, Rust | N/A | BSD 2-Clause | |
| XML | SGML | W3C Recommendations: 1.0(Fifth Edition) 1.1(Second Edition) | Partial (Efficient XML Interchange, Binary XML, Fast Infoset, XSD base64 data) | Yes | DOM, SAX, XQuery, XPath | N/A | | |
| UBJSON | JSON, BSON | Ubjson.org | Yes | No | No | N/A | Apache 2.0 | |
| MessagePack | JSON(loosely) | MessagePack format specification (github: https://github.com/msgpack/msgpack-c/tree/cpp_master ) | Yes | No | No | Yes | Apache | • requires boost library. |
| Cap'n Proto | | Capnproto.org | | | | | MIT | |
| CSV (Comma-serperated values) | N/A | RFC 4180 (among others) | No | Yes | No | No | | |
| Cista++ | | https://github.com/felixguendling/cista | | | | | MIT | • C++17 compatible compiler<br>• CMake |
| YAS (Yet Another Serialization) | | | | | | | Boost | C++ 11 support<br>• Supported types of archives:<br>  binary<br>  text<br>  json(not fully comply)<br>• Supported compilers:<br>  GCC 4.8.5 or higher(32/64 bit)<br>  MinGW: 4.8.5 or higher(32/64 bit)<br>  Clang: 3.5 or higher(32/64 bit)<br>  Intel: (untested)<br>  MSVC : 2017(in c++14 mode), or higher(32/64 bit)<br>  Emscripten: 1.38 (clang version 6.0.1) |
| cereal | | https://uscilab.github.io/cereal/index.html | Yes | Yes (RapidJSON, RapidXML) | C++ | | BSD | • officially supports g++ 4.7.3, clang++ 3.3, and MSVC 2013 (or newer).<br>• cmake_minimum_required(VERSION 3.6...3.15) |

| name | Based on | Specification | Binary | Human-readable | Standard APIs | Supports zero-copy operations | license | Requirements |
|---|---|---|---|---|---|---|---|---|
| | | | | | | | | **Narnia current use CMake(version 3.16.2)** |

## Speed and Size comparison

(from https://github.com/thekvs/cpp-serializers)

Following results were obtained running 1000000 serialize-deserialize operations 50 times and then averaging results on a typical desktop computer with Intel Core i7 processor running Ubuntu 16.04.

For capnproto and flatbuffers since they already store data in a "serialized" form and serialization basically means getting pointer to the internal storage, we measure full build/serialize/deserialize cycle. In the case of other libraries we measure serialize/deserialize cycle of the already built data structure.

| serializer name | version | object's size(bytes) | average total time(ms) | Rank of Object size per serializer | Rank of time per serializer | Human-readable |
|---|---|---|---|---|---|---|
| thrift-binary | 0.12.0 | 17017 | 1190.22 | 6 | 5 | No |
| thrift-compact | 0.12.0 | 13378 | 3474.32 | 2 | 10 | Partial |
| protobuf | 3.7.0 | 16116 | 2312.78 | 4 | 8 | Partial (If you use text_format, you can read.) |
| boost | 1.69.0 | 17470 | 1195.04 | 9 | 6 | No |
| MessagePack | 3.1.1 | 13402 | 2560.6 | 3 | 9 | No |
| cereal | 1.2.2 | 17416 | 1052.46 | 7 | 4 | Yes (If you use JSON Archive type, you can read) |
| avro | 1.8.2 | 16384 | 4488.18 | 5 | 11 | Partial |
| capnproto | 0.7.0 | 17768 | 400.98 | 11 | 2 | No |
| flatbuffers | 1.10.0 | 17632 | 491.5 | 10 | 3 | Apache Arrow |
| yas | 7.0.2 | 17416 | 302.7 | 7 | 1 | Partial. (JSON format is supported but not fully comply) |
| yas-compact | 7.0.2 | 13321 | 2063.34 | 1 | 7 | Partial. (JSON format is supported but not fully comply) |

serializer name

Size   : yas-compact < thrift-compact < MessagePack < protobuf < avro < thrift-binary < yas , cereal < boost < flatbuffers < capn proto

Time : yas < capn proto < flatbuffers < cereal < thrift-binary < boost < yas-compact < protobuf < MessagePack < thrift-compact < avro

## JSON Performance

- CArchive is using JsonCpp library.(https://github.com/open-source-parsers/jsoncpp)

Libraries

43 libraries are successfully benchmarked. They are listed in alphabetic order: (Libraries section of https://github.com/miloyip/nativejson-benchmark)

## JsonCpp vs. RapidJSON

The followings are some snapshots from the results of MacBook Pro (Retina, 15-inch, Mid 2015, Corei7-4980HQ@2.80GHz) with clang 7.0 64-bit.

Overall

More detail information: https://rawgit.com/miloyip/nativejson-benchmark/master/sample/performance_Corei7-4980HQ@2.80GHz_mac64_clang7.0.html#1.%20Parse

**0. Overall**



| Library | Result |
|---|---|
| ArduinoJson (C++) | 69% |
| C++ REST SDK (C++11) | 87% |
| CAJUN (C++) | 79% |
| ccan/json (C) | 89% |
| cJSON (C) | 67% |
| Configuru (C++11) | 99% |
| dropbox/json11 (C++11) | 87% |
| Folly (C++11) | 86% |
| gason (C++11) | 55% |
| hjiang/JSON++ (C++) | 53% |
| Jansson (C) | 87% |
| JeayeSON (C++14) | 52% |
| jsmn (C) | 58% |
| JSON Spirit (C++) | 69% |
| JSON Voorhees (C++) | 87% |
| json-c (C) | 88% |
| JsonBox (C++) | 62% |
| jsoncons (C++) | 95% |
| JsonCpp (C++) | 85% |
| JVar (C++) | 88% |
| Jzon (C++) | 71% |
| mikeando/FastJson (C++) | 92% |
| nbsdx_SimpleJSON (C++11) | 34% |
| Nlohmann (C++11) | 96% |
| Parson (C) | 83% |
| PicoJSON (C++) | 87% |
| POCO (C++) | 94% |
| Qt (C++) | 84% |
| RapidJSON (C++) | 93% |
| RapidJSON_AutoUTF (C++) | 93% |
| RapidJSON_FullPrec (C++) | 100% |
| RapidJSON_Insitu (C++) | 93% |
| sajson (C++) | 86% |
| Scheredom json.h (C) | 74% |
| SimpleJSON (C++) | 67% |
| taocpp/json (C++11) | 100% |
| tunnuz/JSON++ (C++) | 74% |
| udp/json-parser (C) | 79% |
| ujson (C++) | 94% |
| ujson4c (C) | 63% |
| V8 (C++) | 94% |
| Vinenthz/libjson (C) | 88% |
| YAJL (C) | 87% |

*Result*

Parsing Time

**1. Parse**



Parsing Memory

**1. Parse**

| Library | Memory (byte) |
|---|---|
| ArduinoJson (C++) | 17,991,808 |
| C++ REST SDK (C++11) | 10,936,608 |
| CAJUN (C++) | 248,399,248 |
| ccan/json (C) | 15,837,328 |
| cJSON (C) | 14,985,504 |
| Configuru (C++11) | 20,040,848 |
| dropbox/json11 (C++11) | 18,022,768 |
| gason (C++11) | 10,137,696 |
| hjiang/JSON++ (C++) | 14,550,512 |
| Jansson (C) | 18,319,600 |
| jsmn (C) | 11,767,904 |
| json-c (C) | 49,886,288 |
| JsonBox (C++) | 11,782,928 |
| jsoncons (C++) | 7,564,816 |
| JsonCpp (C++) | 24,560,400 |
| JVar (C++) | 7,563,072 |
| Jzon (C++) | 23,624,624 |
| mikeando/FastJson (C++) | 13,504,128 |
| nbsdx_SimpleJSON (C++11) | 246,772,128 |
| Nlohmann (C++11) | 9,897,872 |
| Parson (C) | 9,144,496 |
| PicoJSON (C++) | 9,739,632 |
| POCO (C++) | 16,504,080 |
| Qt (C++) | 144 |
| RapidJSON (C++) | 4,833,344 |
| RapidJSON_AutoUTF (C++) | 4,870,208 |
| RapidJSON_FullPrec (C++) | 4,833,344 |
| RapidJSON_Insitu (C++) | 10,822,720 |
| sajson (C++) | 52,322,640 |
| Scheredom json.h (C) | 24,711,216 |
| SimpleJSON (C++) | 12,679,568 |
| strdup (C) | 6,602,848 |
| taocpp/json (C++11) | 11,506,176 |
| tunnuz/JSON++ (C++) | 28,651,376 |
| udp/json-parser (C) | 17,195,056 |
| ujson (C++) | 53,401,952 |
| ujson4c (C) | 22,052,960 |
| V8 (C++) | 15,971,176 |
| Vinenthz/libjson (C) | 8,594,288 |
| YAJL (C) | 17,383,568 |

*Memory (byte)*

Stringify Time

**2. Stringify**

| Library | Time (ms) |
|---|---|
| ArduinoJson (C++) | 58 |
| C++ REST SDK (C++11) | 78 |
| CAJUN (C++) | 138 |
| ccan/json (C) | 66 |
| cJSON (C) | 71 |
| Configuru (C++11) | 155 |
| dropbox/json11 (C++11) | 73 |
| gason (C++11) | 45 |
| hjiang/JSON++ (C++) | 336 |
| Jansson (C) | 70 |
| json-c (C) | 31 |
| JsonBox (C++) | 215 |
| jsoncons (C++) | 86 |
| JsonCpp (C++) | 94 |
| JVar (C++) | 50 |
| Jzon (C++) | 28 |
| mikeando/FastJson (C++) | 67 |
| nbsdx_SimpleJSON (C++11) | 79 |
| Nlohmann (C++11) | 88 |
| Parson (C) | 111 |
| PicoJSON (C++) | 80 |
| POCO (C++) | 50 |
| Qt (C++) | 152 |
| RapidJSON (C++) | 11 |
| RapidJSON_AutoUTF (C++) | 20 |
| RapidJSON_FullPrec (C++) | 11 |
| RapidJSON_Insitu (C++) | 11 |
| Scheredom json.h (C) | 34 |
| SimpleJSON (C++) | 206 |
| strdup (C) | 0 |
| taocpp/json (C++11) | 31 |
| tunnuz/JSON++ (C++) | 248 |
| udp/json-parser (C) | 45 |
| ujson (C++) | 25 |
| V8 (C++) | 34 |
| Vinenthz/libjson (C) | 69 |
| YAJL (C) | 79 |

*Time (ms)*

Code size
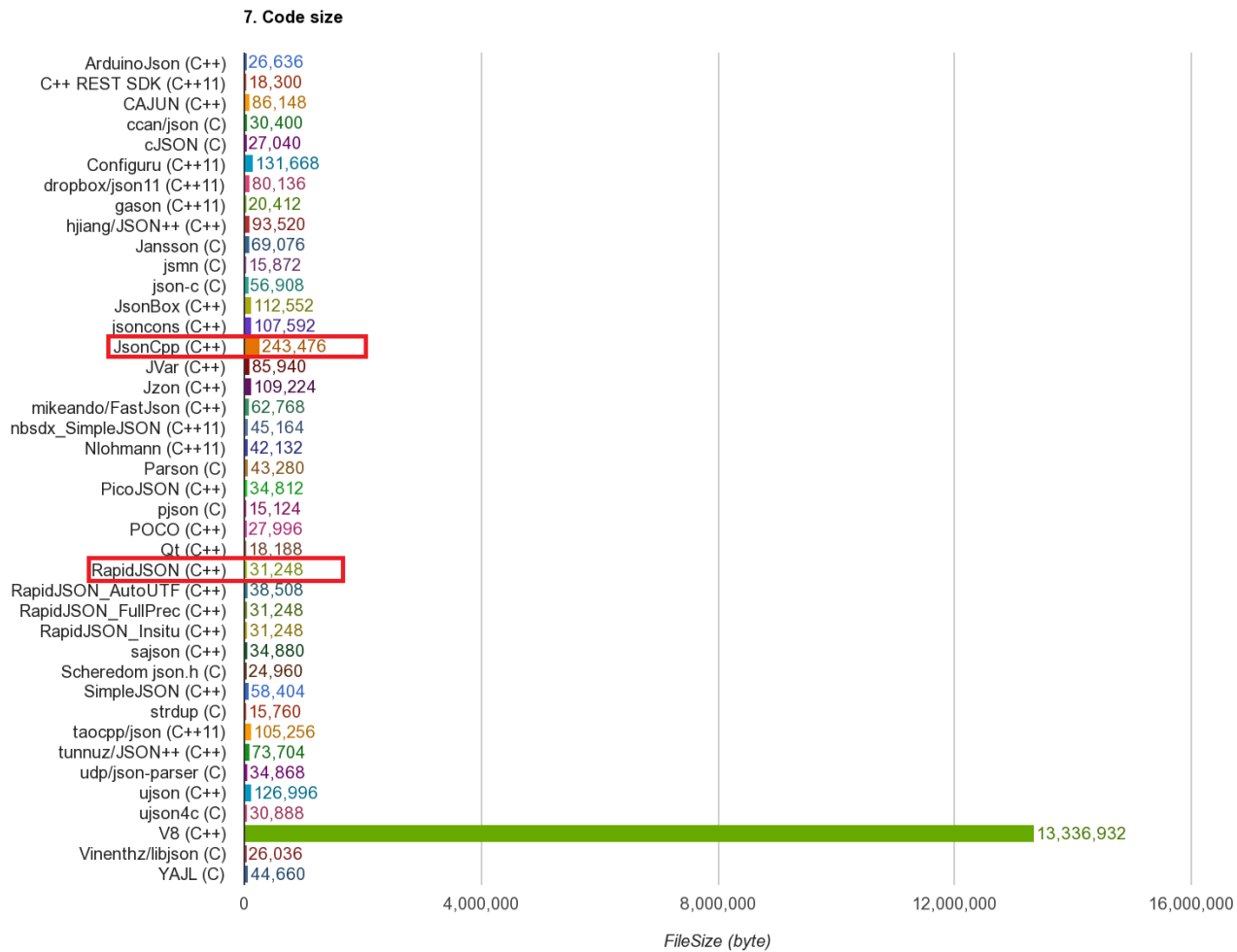
**7. Code size**



## Conclusion

According to serializer library benchmark, I found that we could use cereal library. Cereal serialization library supports JSON format. So it is human-readable.

And I searched the Json library performance benchmark result. According to the result, RapidJson library is faster and smaller than JsonCpp library.

If CArchive struct is used continuously, we can try to change native JSON library(JsonCpp library) to RapidJson.

## Reference sites

https://blog.mbedded.ninja/programming/serialization-formats/a-comparison-of-serialization-formats/

https://en.wikipedia.org/wiki/Comparison_of_data-serialization_formats

https://github.com/thekvs/cpp-serializers

https://github.com/felixguendling/cpp-serialization-benchmark

https://uscilab.github.io/cereal/

https://github.com/miloyip/nativejson-benchmark

No labels