

seriot.ch

PARSING JSON IS A MINEFIELD

[2016-10-26] First version of the article

[2016-10-28] Presentation at Soft-Shake Conference, Geneva ([slides](#))

[2016-11-01] Article and comments in [The Register](#)

Feel free to comment on [Hacker News](#) or [reddit](#).

Session Description

JSON is the de facto standard when it comes to (un)serialising and exchanging data in web and mobile programming. But how well do you really know JSON? We'll read the specifications and write test cases together. We'll test common JSON libraries against our test cases. I'll show that JSON is not the easy, idealised format as many do believe. Indeed, I did not find two libraries that exhibit the very same behaviour. Moreover, I found that edge cases and maliciously crafted payloads can cause bugs, crashes and denial of services, mainly because JSON libraries rely on specifications that have evolved over time and that left many details loosely specified or not specified at all.

Table of Contents

1. [JSON Specifications](#)
2. [Parsing Tests](#)
 - 2.1 [Structure](#)
 - 2.2 [Numbers](#)
 - 2.3 [Arrays](#)
 - 2.4 [Objects](#)
 - 2.5 [Strings](#)
 - 2.6 [RFC 7159 Ambiguities](#)
3. [Testing Architecture](#)
4. [Parsing Tests Results](#)
 - 4.1 [Full Results](#)
 - 4.2 [C Parsers](#)
 - 4.3 [Objective-C Parsers](#)
 - 4.4 [Apple \(NS\)JSONSerialization](#)
 - 4.5 [Freddy \(Swift\)](#)
 - 4.6 [Bash JSON.sh](#)
 - 4.7 [Other Parsers](#)
 - 4.8 [JSON Checker](#)
 - 4.9 [Regex](#)
5. [Parsing Contents](#)
6. [STJSON](#)
7. [Conclusion](#)
8. [Appendix](#)

1. JSON Specifications

JSON is the de facto serialization standard when it comes to sending data over HTTP, the *lingua franca* used to exchange data between heterogeneous software, both in modern web sites and mobile applications.

"Discovered" in 2001 [Douglas Crockford](#), JSON specification is so short and simple that Crockford created business cards with the whole JSON grammar on their back.

2/13

Note that since several parsers don't allow scalars at the top level ("test"), I embed strings into arrays (["test"]).

You'll find more than 300 tests in the [JSONTestSuite GitHub repository](#).

Finally, I also generated JSON files with the fuzzing software [American Fuzzy Lop](#). I then removed redundant tests that produced the same set of results, and then reduced the remaining ones to the keep the least number of characters that triggered these results (see [section 3](#)).

2.1 Structure

y_structure_lonely_string.json	"asd"
--------------------------------	-------

n_object_trailing_comma.json	{"id":0,}
n_object_several_trailing_commas.json	{"id":0,,,,,}

y_string_comments.json	["a/*b*/c/*d//e"]
n_object_trailing_comment.json	{"a": "b"}/**/
n_structure_object_with_comment.json	{"a": /*comment*/ "b"}

n_structure_object_unclosed_no_value.json	{"":
n_structure_object_followed_by_closing_object.json	{}}

```
$ python -c "print(' '*100000)" > ~/x.json
$ ./Xcode ~/x.json
Segmentation fault: 11
```

n_structure_whitespace_formfeed.json	[QC]
n_structure_whitespace_U+2060_word_joiner.json	[E281A0]
n_structure_no_data.json	

2.2 Numbers

n_number_NaN.json	[NaN]
n_number_minus_infinity.json	[-Infinity]

n_number_hex_2_digits.json	[0x42]
----------------------------	--------

i_number_very_big_negative_int.json	[-237462374673276894279832(...)]
-------------------------------------	----------------------------------

3/13

Exponential Notation - Parsing exponential notation can be surprisingly hard (see the results section). Here are some valid contents `[0E0]`, `[0e+1]` and invalid ones `[1.0e+]`, `[0E]` and `[1eE2]`.

n_number_0_capital_E+.json	[0E+]
n_number_.2e-3.json	[.2e-3]
y_number_double_huge_neg_exp.json	[123.456e-789]

2.3 Arrays

Most edge cases regarding arrays are opening/closing issues and nesting limit. These cases were discussed in section [2.1 Structure](#). Passing tests will include `[[]]`, `[[]]`, while failing tests will be like `]` or `[[]]`.

n_array_comma_and_number.json	[,1]
n_array_colon_instead_of_comma.json	["": 1]
n_array_unclosed_with_new_lines.json	[1,0A10A,1]

2.4 Objects

Duplicated Keys - [RFC 7159 section 4](#) says that "The names within an object should be unique.". It does not prevent parsing objects where the same key does appear several times `{"a":1,"a":2}`, but lets parsers decide what to do in this case. The same section 4 even mentions that "(some) implementations report an error or fail to parse the object", without telling clearly if failing to parse such objects is compliant or not with the RFC and especially [section 9](#): "A JSON parser MUST accept all texts that conform to the JSON grammar."

Variants of this special case include same key - same value `{"a":1,"a":1}`, and similar keys or values, where the similarity depends on how you compare strings. For example, the keys may be binary different but equivalent according to Unicode NFC normalization, such as in `{"C3A9":"NFC","65CC81":"NFD"}` where both keys encode "é". Tests will also include `{"a":0,"a":-0}`.

y_object_empty_key.json	{"":0}
y_object_duplicated_key_and_value.json	{"a":"b","a":"b"}
n_object_double_colon.json	{"x"::"b"}
n_object_key_with_single_quotes.json	{key: 'value'}
n_object_missing_key.json	:{ "b" }
n_object_non_string_key.json	{1:1}

2.5 Strings

File Encoding - "JSON text SHALL be encoded in UTF-8, UTF-16, or UTF-32. The default encoding is UTF-8" [section 8.1](#).

Passing tests should include text encoded in these three encodings. UTF-16 and UTF-32 texts should also include both their big-endian and little-endian variants.

The parsing of invalid UTF-8 will be implementation defined.

y_string_utf16.json	FFFE[00"00E900"00]00
i_string_iso_latin_1.json	["E9"]

[Update 2016-11-04] The first version of this article considered invalid UTF-8 as n_ tests. This classification was [challenged](#) and I eventually changed these tests into i_ tests.

Byte Order Mark - While [section 8.1](#) states "Implementations MUST NOT add a byte order mark to the beginning of a JSON text", "implementations (...) MAY ignore the presence of a byte order mark rather than treating it as an error".

Failing tests will include a plain UTF-8 BOM with no other content. Tests with implementation defined results should include a UTF-8 BOM with a UTF-8 string, but also a UTF-8 BOM with a UTF-16 string, and a UTF-16 BOM with a UTF-8 string.

n_structure_UTF8_BOM_no_data.json	EFBBBF
n_structure_incomplete_UTF8_BOM.json	EFBB{}
i_structure_UTF-8_BOM_empty_object.json	EFBBBF{}

Control Characters - Control characters must be escaped, and are defined as U+0000 through U+001F ([section 7](#)). This range does not include 0x7F DEL, which may be part of other definitions of control characters (see [section 4.6 Bash JSON.sh](#)). That is why passing tests include `["7F"]`.

n_string_unescaped_ctrl_char.json	["a\09a"]
y_string_unescaped_char_delete.json	["7F"]
n_string_escape_x.json	["\x00"]

Escape - "All characters may be escaped" ([section 7](#)) like `\uXXXX`, but some MUST be escaped: quotation mark, reverse solidus and control characters. Failing tests should include the escape character without the escaped value, or with an incomplete escaped value. Examples: `["\""]`, `["\"`, `["\`.

y_string_allowed_escapes.json	["\"\\/\b\f\n\r\t"]
n_structure_bad_escape.json	["\"

The escape character can be used to represent codepoints in the Basic Multilingual Plane (`\u005C`). Passing tests will include the zero character `\u0000`, which may cause issues in C-based parsers. Failing tests will include capital U `\U005C`, non-hexadecimal escaped values `\u123Z` and incomplete escaped values `\u123`.

y_string_backslash_and_u_escaped_zero.json	["\u0000"]
n_string_invalid_unicode_escape.json	["\uqqqq"]
n_string_incomplete_escaped_character.json	["\u00A"]

Escaped Invalid Characters

Codepoints outside of the BMP are represented by their escaped UTF-16 surrogates: U+1D11E becomes `\uD834\uDD1E`. Passing tests will include single surrogates, since they are valid JSON according to the grammar. RFC 7159 [errata 3984](#) raised the issue of grammatically correct escaped codepoints that don't encode Unicode characters.

The ABNF cannot at the same time allow non conformant Unicode codepoints (section 7) and states conformance to Unicode (section 1).

The editors considered that the grammar should not be restricted, and that warning users about the fact that parsers behaviour was "unpredictable" ([RFC 7159 section 8.2](#)) was enough. In other words, parsers MUST parse u-escaped invalid codepoints, but the result is undefined, hence the `i_` (implementation defined) prefix in the file name. According to the Unicode standard, invalid codepoints should be replaced by U+FFFD REPLACEMENT CHARACTER. People familiar with [Unicode complexity](#) won't be surprised that this replacement is not mandatory, and can be done in several ways (see [Unicode PR #121: Recommended Practice for Replacement Characters](#)). So several parsers use replacement characters, while other keep the escaped form or produce a non-Unicode character (see [Section 5 - Parsing Contents](#)).

[Update 2016-11-03] In the first version of this article, I treated non-characters such as U+FDD0 to U+10FFFF the same as invalid codepoints (`i_` tests). This classification was [challenged](#) and I eventually changed the non-characters tests into `y_` tests.

y_string_accepted_surrogate_pair.json	["\uD801\uDC37"]
n_string_incomplete_escaped_character.json	["\u00A"]
i_string_incomplete_surrogates_escape_valid.json	["\uD800\uD800\n"]
i_string_lone_second_surrogate.json	["\uDFAA"]
i_string_1st_valid_surrogate_2nd_invalid.json	["\uD888\u1234"]
i_string_inverted_surrogates_U+1D11E.json	["\uDd1e\uD834"]

Raw non-Unicode Characters

The previous section discussed non-Unicode codepoints that appear in strings, such as `"\uDEAD"`, which is valid Unicode in its u-escaped form, but doesn't decode into a Unicode character.

Parsers also have to handle raw bytes that don't encode Unicode characters. For instance, the byte `FF` does not represent a Unicode character in UTF-8. As a consequence, a string containing `FF` is not an UTF-8 string. In this case, parsers should simply refuse to parse the string, because "A string is a sequence of zero or more Unicode characters" [RFC 7159 section 1](#) and "JSON text SHALL be encoded in Unicode" [RFC 7159 section 8.1](#).

y_string_utf8.json	["€"]
n_string_invalid_utf-8.json	["FF"]
n_array_invalid_utf8.json	[FF]

RFC 7159 Ambiguities

Beyond the specific cases we just went through, finding out if a parser is RFC 7159 compliant or not is next to impossible because of [section 9 "Parsers"](#):

A JSON parser MUST accept all texts that conform to the JSON grammar. A JSON parser MAY accept non-JSON forms or extensions.

To this point, I perfectly understand the RFC. All grammatically correct inputs MUST be parsed, and parsers are free to accept other contents as well.

An implementation may set limits on the size of texts that it accepts. An implementation may set limits on the maximum depth of nesting. An implementation may set limits on the range and precision of numbers. An implementation may set limits on the length and character contents of strings.

All these limitations sound reasonable (except maybe the one about "character contents"), but contradict the "MUST" from the previous sentence. [RFC 2119](#) is crystal-clear about the meaning of "MUST":

MUST - This word, or the terms "REQUIRED" or "SHALL", mean that the definition is an absolute requirement of the specification.

RFC 7159 allows restrictions, but does not set minimal requirements, so technically speaking, a parser that cannot parse strings longer than 3 characters is still compliant with RFC 7159.

Also, RFC 7159 section 9 should require the parsers to document the restrictions clearly, and/or allow configuration by the user. These configurations would still cause interoperability issues, that's why minimal requirements should be preferred.

This lack of precision regarding allowed restrictions makes it almost impossible to say if a parser is actually RFC 7159 compliant. Indeed, parsing contents that don't match the grammar is not wrong (it's an "extension") and rejecting contents that does match the grammar is allowed (it's a parser "limit").

3. Testing Architecture

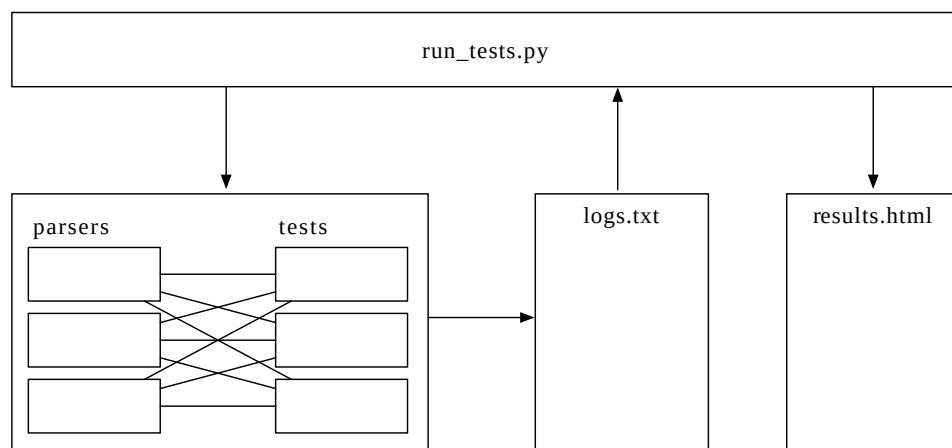
Independently from how parsers should behave, I wanted to observe how they actually behave, so I picked several JSON parsers and set up things so that I could feed them with my test files.

As I'm a Cocoa developer, I included mostly Swift and Objective-C parsers, but also C, Python, Ruby, R, Lua, Perl, Bash and Rust parsers, chosen pretty arbitrarily. I mainly tried to achieve diversity in age, popularity and languages.

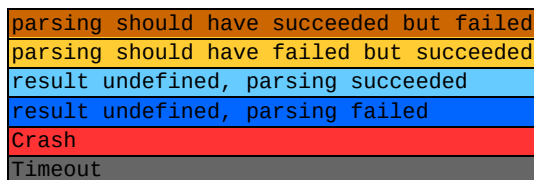
Several parsers have options to increase or decrease strictness, tweak Unicode support or allow specific extensions. I strived to always configure the parsers so that they behave as close as possible to the most strict interpretation of RFC 7159.

A Python script `run_tests.py` runs each parser with each test file (or a single test when the file is passed as an argument). The parsers are generally wrapped so that the process returns 0 in case of success, 1 in case of parsing error, yet another status in case of crash, a 5-second delay being considered as a timeout. Basically, I turned JSON parsers into JSON validators.

Python 2.7.10 SHOULD HAVE FAILED n number infinity.json



The results show one row per file, one column per parser, and one color per unexpected result. They also show detailed results by parser.



Tests are sorted by results equality, making easy to spot sets of similar results and remove redundant tests.

[illegible]

4. Results and Comments

- 4.1 [Full Results](#)
- 4.2 [C Parsers](#)
- 4.3 [Obj-C Parsers](#)
- 4.4 [Apple \(NS\)JSONSerialization](#)
- 4.5 [Swift Freddy](#)
- 4.6 [Bash](#)
- 4.7 [Other Parsers](#)
- 4.8 [JSON Checker](#)
- 4.9 [Regex](#)

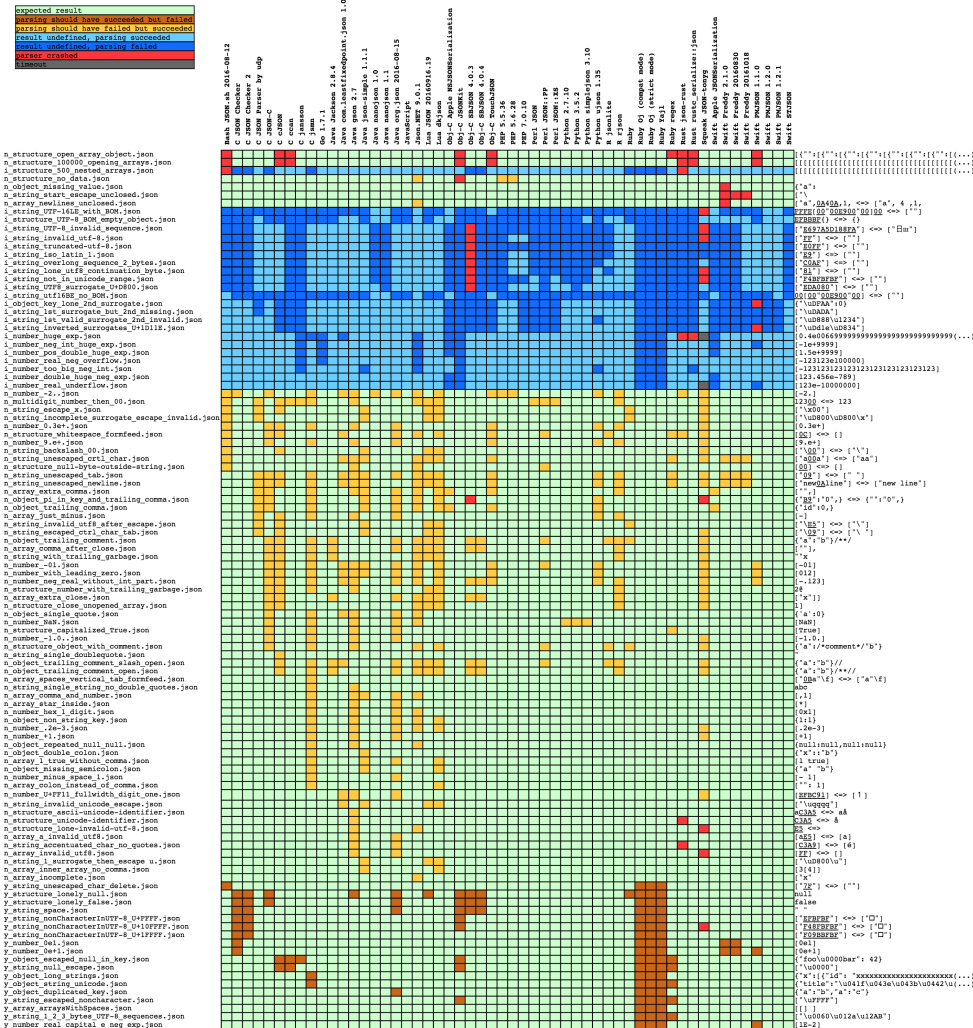
4.1 Full Results

Full results are presented in <http://seriot.ch/json/parsing.html>. The tests are vertically sorted by similar results, so it is easy to prune similar tests. An option in `run_tests.py` will produce "pruned results": when a set of tests yields the same results, only the first one is kept. Pruned results HTML file is available here: http://www.seriot.ch/json/parsing_pruned.html.

The most serious issues are crashes (in red), since parsing an uncontrolled input may put the whole process at risk. The "should have passed" tests (in brown) are also very dangerous, because an uncontrolled input may prevent the parser to parse a whole document. The "should have failed" tests (in yellow) are more benign. They indicate a JSON "extension" that can be parsed. Everything will work fine, until the parser is replaced with another parser which does not parse the same "extensions"...

JSON Parsing Tests, Pruned

Appendix to: seriot.ch [Parsing JSON is a Minefield](http://www.seriot.ch/parsing_json.php) http://www.seriot.ch/parsing_json.php
2016-11-05 00:04:08



This section highlights and comments some noticeable results.

4.2 C Parsers

Here are the five C parsers considered:

- <https://github.com/zserge/jsmn>
- <https://github.com/akheron/jansson>
- <https://github.com/rustyrussell/ccan/>
- <https://github.com/DaveGamble/cJSON>
- <https://github.com/udp/json-parser>

And here is a quick comparison between them:

	jsmn	jansson	ccan	cJSON	json-parser
Parses ["\u0000"]	YES	YES	NO	NO	YES
Too liberal	YES	NO	NO	YES	YES
Crash on nested structs.	NO	NO	YES	YES	NO
Rejects big numbers	YES	YES	NO	NO	NO

You can refer to the full results for more details.

4.3 Objective-C Parsers

Here are a couple of Objective-C parsers that used to be very popular in the early days of iOS development, especially because Apple waited until iOS 5 to release NSJSONSerialization. They are still interesting to test, since they are used in production in many applications. Let's consider:

- <https://github.com/ohnhezang/JSONKit>
- <https://github.com/TouchCode/TouchJSON>
- <https://github.com/stig/json-framework> aka SBJSON

And here is a quick comparison between them:

JSONKit TouchJSON SBJSON

Crash on nested structs.	YES	NO	YES
Crash on invalid UTF-8	NO	NO	YES
Parses trailing garbage []x	NO	NO	YES
Rejects big numbers	NO	YES	NO
Parses bad numbers [0.e1]	NO	YES	NO
Treats 0x0C FORM FEED as white space	NO	YES	NO
Parses non-char. ["\uFFFF"]	NO	YES	YES

SBJSON survived after the arrival of NSJSONSerialization, is still maintained and is available through CocoaPods, so I reported the crash when parsing non UTF-8 strings such as ["FF"] in [issue #219](#).

```
*** Assertion failure in -[SBJson4Parser parserFound:isValue:], SBJson4Parser.m:150
*** Terminating app due to uncaught exception 'NSInternalInconsistencyException', reason: 'Invalid param
*** First throw call stack:
(
  0  CoreFoundation          0x00007fff95f4b4f2 __exceptionPreprocess + 178
  1  libobjc.A.dylib         0x00007fff9783bf7e objc_exception_throw + 48
  2  CoreFoundation          0x00007fff95f501ca +[NSException raise:format:arguments:] +
  3  Foundation              0x00007fff9ce86856 -[NSAssertionHandler handleFailureInMethod:format:arguments:] + 128
  4  test_SBJSON             0x00000001000067e5 -[SBJson4Parser parserFound:isValue:] + 3
  5  test_SBJSON             0x00000001000073f3 -[SBJson4Parser parserFoundString:] + 67
  6  test_SBJSON             0x0000000100004289 -[SBJson4StreamParser parse:] + 2377
  7  test_SBJSON             0x0000000100007989 -[SBJson4Parser parse:] + 73
  8  test_SBJSON             0x0000000100005d0d main + 221
  9  libdyld.dylib           0x00007fff929ea5ad start + 1
)
libc++abi.dylib: terminating with uncaught exception of type NSException
```

4.4 Apple (NS)JSONSerialization

<https://developer.apple.com/reference/foundation/nsjsonserialization>

NSJSONSerialization was introduced with iOS 5 and is the standard JSON parser on OS X and iOS since then. It is available in Objective-C, and was rewritten in Swift: [NSJSONSerialization.swift](#). The NS prefix was then [dropped](#) in Swift 3.

Restrictions and Extensions

JSONSerialization has the following, undocumented restrictions:

- it won't parse big numbers: [123123e100000]
- it won't parse u-escaped invalid codepoints: ["\ud800"]

JSONSerialization has the following, undocumented extension:

- it does parse trailing commas: [1,] and {"a":0,}

I find the restriction about invalid codepoints to be especially problematic, especially in such a high-profile parser, because trying to parse uncontrolled contents may result in a parsing failure.

Crash on Serialization

This article is more about JSON parsing than JSON producing, yet I wanted to mention this crash that I found in JSONSerialization when writing Double.nan. Remember that NaN does not conform to JSON grammar, so in this case, JSONSerialization should throw an error, but not crash the whole process.

```
do {
    let a = [Double.nan]
    let data = try JSONSerialization.data(withJSONObject: a, options: [])
} catch let e {
}

SIGABRT

*** Terminating app due to uncaught exception 'NSInvalidArgumentException', reason: 'Invalid number value'
```

[Update 2016-10-27] The original version of this article erroneously said that JSONSerialization.isValidJSONObject(["x":"x", "x":"x"]) would crash because of a bug in the method. [@H2CO3 iOS found](#) that the crash is not related to JSONSerialization but to Swift dictionaries, that just cannot be build with the same key appearing twice.

4.5 Freddy (Swift)

Freddy (<https://github.com/bignerdranch/Freddy>) is a real JSON Parser written in Swift 3. I say real because several GitHub projects claim to be Swift JSON parsers, but actually use Apple JSONSerialization and just map JSON contents with model objects.

Freddy is interesting because it is written by a famous organization of Cocoa developers, and does leverage Swift type safety by using a Swift enum to represent the different kind of JSON nodes (Array, Dictionary, Double, Int, String, Bool and Null).

But, being [released in January 2016](#), Freddy is still young, and buggy. My test suite showed that unclosed structures such as [1, and {"a": used to crash the parser, as well as a string with a single space " ", so I opened [issue #199](#) that was fixed within 1 day!

Additionally, I found that "0e1" was incorrectly rejected by the parser, so I opened [issue #198](#), which was also fixed within 1 day.

However, Freddy does still crash on 2016-10-18 when parsing `["\.` I reported the bug in ([issue #206](#)).

The following table does summarize the evolution of Freddy's behaviour:

[illegible]

4.6 Bash JSON.sh

I tested <https://github.com/dominictarr/JSON.sh/> from 2016-08-12.

This Bash parser relies on a regex to find the control characters, which **MUST** be backslash-escaped according to RFC 7159. But Bash and JSON don't share the same definition of control characters.

The regex uses the `\c` syntax to match control characters, which is a shorthand for `[\x00-\x1F\x7F]`. But according to JSON grammar, `\x7F` DEL is not part of control characters, and may appear unescaped.

00	nul	01	soh	02	stx	03	etx	04	eot	05	enq	06	ack	07	bel
08	bs	09	ht	0a	nl	0b	vt	0c	np	0d	cr	0e	so	0f	si
10	dle	11	dc1	12	dc2	13	dc3	14	dc4	15	nak	16	syn	17	etb
18	can	19	em	1a	sub	1b	esc	1c	fs	1d	gs	1e	rs	1f	us
20	sp	21	!	22	"	23	#	24	\$	25	%	26	&	27	'
28	(29)	2a	*	2b	+	2c	,	2d	-	2e	.	2f	/
30	0	31	1	32	2	33	3	34	4	35	5	36	6	37	7
38	8	39	9	3a	:	3b	;	3c	<	3d	=	3e	>	3f	?
40	@	41	A	42	B	43	C	44	D	45	E	46	F	47	G
48	H	49	I	4a	J	4b	K	4c	L	4d	M	4e	N	4f	O
50	P	51	Q	52	R	53	S	54	T	55	U	56	V	57	W
58	X	59	Y	5a	Z	5b	[5c	\	5d]	5e	^	5f	_
60	`	61	a	62	b	63	c	64	d	65	e	66	f	67	g
68	h	69	i	6a	j	6b	k	6c	l	6d	m	6e	n	6f	o
70	p	71	q	72	r	73	s	74	t	75	u	76	v	77	w
78	x	79	y	7a	z	7b	{	7c		7d	}	7e	~	7f	del

As a consequence, JSON.sh cannot parse `["7F"]`. I reported this bug in [issue #46](#).

Additionally, JSON.sh does not limit the nesting level, and will crash when parsing 10000 times the open array character [. I reported this bug in [issue #47](#).

```
$ python -c "print(['*1000000') | ./JSON.sh
./JSON.sh: line 206: 40694 Done          tokenize
          40695 Segmentation fault: 11 | parse
```

4.7 Other Parsers

Besides C / Objective-C and Swift, I also tested parsers from other environments. Here is a synthetic review of their extensions and restrictions, with a subset of the [full tests results](#). The goal of this table is to demonstrate that there are no two parsers that agree on what is wrong and what is right.

	Go 1.7.1	JavaScript	Lua JSON 20160916.19	Lua dkjson	PHP 5.5.36	Perl JSON	Perl JSON::XS	Python 2.7.10	R jsonlite	R rjson	Ruby	Rust json-rust	Rust rustc_serialize::json
n_string_UTF8_surrogate_U+D800.json													["EPA080"] <=> [""]
n_string_invalid_utf-8.json													["FF"] <=> [""]
n_number_2.e+3.json													[2.e+3]
n_array_number_and_comma.json													[1,]
n_structure_object_with_comment.json													{"a":/*comment*/"b"}
n_structure_array_trailing_garbage.json													[1]x
n_string_invalid_backslash_esc.json													["\a"]
n_string_unicode_CapitalU.json													"\UA66D"
n_number_minus_infinity.json													[-Infinity]
y_number_real_pos_overflow.json													[123123e100000]
y_number_too_big_pos_int.json													[1000000000000000000000]

Here are the references for the tested parsers:

- Lua JSON 20160728.17 <http://regex.info/blog/lua/json> (enjoy quality of comments in source code)
- Lua dkjson 2.5.1 <https://github.com/LuaDist/dkjson>
- Go 1.7.1, json module <https://golang.org/pkg/encoding/json/>
- Python 2.7.10, json module <https://docs.python.org/2.7/library/json.html>
- JavaScript, macOS 10.12
- Perl JSON <https://metacpan.org/pod/JSON>
- Perl JSON::XS <https://metacpan.org/pod/JSON::XS>
- PHP 5.6.24, macOS 10.12
- R rjson <https://cran.r-project.org/web/packages/rjson/index.html>
- R jsonlite <https://github.com/jeroenooms/jsonlite>
- Rust json-rust <https://github.com/maciejhirsz/json-rust>
- Rust rustc_serialize::json https://doc.rust-lang.org/rustc-serialize/rustc_serialize/json/

Upon popular request, I also added the following Java parsers, which are not shown on this image but that appear in the full results:

- Java Gson 2.7 <https://github.com/google/gson>
- Java Jackson 2.8.4 <https://github.com/FasterXML/jackson>
- Java Simple JSON 1.1.1 <https://code.google.com/archive/p/json-simple/>

The Python JSON module will parse NaN or -Infinity as numbers. While this behaviour can be fixed by setting the `parse_constant` options to a function that will raise an Exception as shown below, it's such an uncommon practice that I didn't use it in the tests, and let the parser erroneously parse these number constants.

```
def f_parse_constant(o):
    raise ValueError

o = json.loads(data, parse_constant=f_parse_constant)
```

4.8 JSON Checker

A JSON parser transforms a JSON document into another representation. If the input is invalid JSON, the parser returns an error.

Some programs don't transform their input, but just tell if the JSON is valid or not. These programs are JSON validators.

json.org has a such a program, written in C, called JSON_Checker http://www.json.org/JSON_checker/, that even comes with a (small) test suite:

JSON_Checker is a Pushdown Automaton that very quickly determines if a JSON text is syntactically correct. It could be used to filter inputs to a system, or to verify that the outputs of a system are syntactically correct. It could be adapted to produce a very fast JSON parser.

Even if JSON_Checker is not a formal reference implementation, one could expect JSON_Checker to clarify JSON specifications or at least implement them correctly.

Unfortunately, JSON_Checker violates the specifications defined on same web site. Indeed, JSON_Checker will parse the following inputs: [1.], [0.e1], which do not match JSON grammar.

- u-escaped codepoints, including valid ones: `["\u002c"]`
- backslash-escaped backslash: `["\\a"]`

- a capitalized True: [True]
- an unescaped control character: ["09"]

RFC 7159 Section 9 says:

All of the above testing architecture will only tell if a parser would parse a JSON document or not, but doesn't say anything about the representation of the resulting contents.

Similarly, extreme numbers such as `0.000000000000000000000001` or `-0` can be parsed, but what should the result be? RFC 7159 doesn't make a distinction between integers and doubles, or zero and -zero. It doesn't even say if numbers can be converted into strings or not.

And what about objects with the same keys (`{"a":1,"a":2}`)? Or same keys and same values (`{"a":1,"a":1}`)? And how should a parser compare object keys?? Should it use binary comparison or a Unicode normal form such as NFC? Here again, RFC is silent

In all these cases, parsers are free to output whatever they want, leading to interoperability issues (think of what could go wrong when you decide to change your usual JSON parser with another one).

With that in mind, let's create tests for which the representation after parsing is not clearly defined. These tests serve only to document how parsers output may differ.

Contrary to the parsing tests, these tests are hard to automate. Instead, the results shown here were observed via log statements and/or debuggers.

Below is an in exhaustive list of some striking differences between resulting representations after parsing. All results are shown in appendix ["Parsing Contents"](#).

- 1.000000000000000005 is generally converted into the float 1.0, but Rust 1.12.0 / json 0.10.2 will keep the original precision and use the number 1.000000000000000005
- 1E-999 is generally converted into float or double 0.0, but Swift Freddy yields the string "1E-999". Swift Apple JSONSerialization and Obj-C JSONKit will simply refuse to parse it and return an error.
- 100000000000000000999 may be converted into a double (Swift Apple JSONSerialization), an unsigned long long (Objective-C JSONKit) or a string (Swift Freddy). It is to be noted that C cJSON will parse this number as a double, but loses precision in the process, resulting in a new number 100000000000000002048 (note the last four digits).

- In `{"C3A9": "NFC", "65CC81": "NFD"}` keys are NFC and NFD representations of "é". Most parsers will yield the two keys, except Swift parsers Apple JSONSerialization and Freddy, where dictionaries first normalize keys before testing them for equality.
- `{"a": 1, "a": 2}` does generally result in `{"a": 2}` (Freddy, SBJSON, Go, Python, JavaScript, Ruby, Rust, Lua dksjon), but may also result in `{"a": 1}` (Obj-C Apple NSJSONSerialization, Swift Apple JSONSerialization, Swift Freddy), or `{"a": 1, "a": 2}` (cJSON, R, Lua JSON).
- `{"a": 1, "a": 1}` does generally result in `{"a": 1}`, but is parsed as `{"a": 1, "a": 1}` in cJSON, R and Lua JSON.
- `{"a": 0, "a": -0}` is generally parsed as `{"a": 0}`, but can also be parsed as `{"a": -0}` (Obj-C JSONKit, Go, JavaScript, Lua) or even `{"a": 0, "a": 0}` (cJSON, R).

- ["\u0000"] contains the u-escaped form of the 0x00 NUL character, which is likely to cause problems in C-based JSON parsers. Most parsers handle this payload gracefully, but JSONKit and cJSON won't parse it. Interestingly, Freddy yields only ["A"] (the string stop after unescaping byte 0x00).
- ["\uD800"] is the u-escaped form of U+D800, an invalid lone UTF-16 surrogate. Many parsers will fail and return an error, despite the string being perfectly valid according to JSON grammar. Python leaves the string untouched and yields ["\uD800"]. Go and JavaScript replace the offending character with "❖" U+FFFD REPLACEMENT CHARACTER ["EFBFD"], R rjson and Lua dkjson simply translate the codepoint into its UTF-8 representation ["EDA080"]. R rjsonlite and Lua JSON 20160728.17 replace the offending codepoint with a question mark ["?"].

- `["EDA080"]` is the non-escaped, UTF-8 form or `U+D800`, the invalid lone UTF-16 surrogate discussed in previous point. This string is not valid UTF-8 and should be rejected (see [section 2.5 Strings - Raw non-Unicode Characters](#)). In practice however, several parsers leave the string untouched `["EDA080"]` such as cJSON, R rjson and jsonlite, Lua JSON, Lua dkjson and Ruby. Go and Javascript yield `["EFBFBDEFBFBDEFBFBFD"]` that is three replacement characters (one per byte). Interestingly, Python 2 converts the sequence into its unicode-escaped form `["\ud800"]`, while Python 3 throws a `UnicodeDecodeError` exception.
- `["\ud800\ud800"]` makes some parsers go nuts. R jsonlite yields `["\U00010000"]`, while Ruby parser yields `["F0908080"]`. I still don't get where this value comes from.

6. STJSON

STJSON is a Swift 3, 600+ lines JSON parser I wrote to see what it took to consider all pitfalls and pass all tests.

<https://github.com/nst/STJSON>

STJSON API is very simple:

```
var p = STJSONParser(data: data)

do {
    let o = try p.parse()
    print(o)
} catch let e {
    print(e)
}
```

STJSON can be instantiated with additional parameters:

```
var p = STJSON(data:data,
               maxParserDepth:1024,
               options:[.useUnicodeReplacementCharacter])
```

In fact, there is only one test where STJSON fails: `y_string_utf16.json`. This is because, as in nearly all other parsers, STJSON does not support non UTF-8 encodings, even though it should not be very difficult to add, and I may do so in the future if needed. At least, STJSON does raise appropriate errors when a file starts with a UTF-16 or UTF-32 byte order mark.

7. Conclusion

In conclusion, JSON is not a data format you can rely on blindly. I've demonstrated this by showing that the standard definition is spread out over at least six different documents ([section 1](#)), that the latest and most complete document, RFC-7159, is imprecise and contradictory ([section 2](#)), and by crafting test files that out of over 30 parsers, no two parsers parsed the same set of documents the same way ([section 4](#)).

In the process of inspecting parser results, I also discovered that `json_checker.c` from `json.org` did reject valid JSON `[0e1]` ([section 4.24](#)), which certainly doesn't help users to know what's right or wrong. In a general way, many parsers authors like to brag about how right is their parsers (including myself), but there's no way to assess their quality since references are debatable and existing test suites are generally poor.

So, I wrote yet another JSON parser ([section 6](#)) which will parse or reject JSON document according to my understanding of RFC 7159. Feel free to comment, open issues and pull requests.

This work may be continued by:

- Documenting the behaviour of **more parsers**, especially parsers that run in non-Apple environments.
- Investigating **JSON generation**. I extensively assessed what parsers do or do not parse. ([section 4](#)). I briefly assessed the contents that parsers yield when the parsing is successful ([section 5](#)). I'm pretty sure that several parsers do generate grammatically invalid JSON or even crash in some circumstances (see [Section 4.2.1](#)).
- Investigating differences in the way **JSON mappers** maps JSON contents to model objects.
- **Finding exploits** in existing software stacks (check out my [Unicode Hacks](#) presentation)
- Investigating potential interoperability issues in **other serialization formats** such as YAML, [BSON](#) or [ProtoBuf](#), which may be a potential successor to JSON. Indeed, Apple already has a Swift implementation <https://github.com/apple/swift-protobuf-plugin>.

As a final word, I keep on wondering why "fragile" formats such as HTML, CSS and JSON, or "dangerous" languages such as PHP or JavaScript became so immensely popular. This is probably because they are easy to start with by tweaking contents in a text editor, because of too liberal parsers or interpreters, and seemingly simple specifications. But sometimes, simple specifications just mean hidden complexity.

8. Appendix

1. Parsing Results <http://seriot.ch/json/parsing.html>, generated automatically for [section 4](#)
2. Transform Results <http://seriot.ch/json/transform.html>, created manually for [section 6](#)
3. JSONTestSuite <https://github.com/nst/JSONTestSuite>, contains all tests and code
4. STJSON <https://github.com/nst/STJSON>, contains my Swift 3 JSON parser

Acknowledgments

Many thanks to @Reversity, GEndignoux, @ccorsano, @BalestraPatrick and @iPlop.