

The Usage of Negation in Real-World JSON Schema Documents

Mohamed-Amine Baazizi¹, Dario Colazzo², Giorgio Ghelli³, Carlo Sartiani⁴ and Stefanie Scherzinger⁵

¹Sorbonne Université, LIP6 UMR 7606, France

²Université Paris-Dauphine, PSL Research University, France

³Dipartimento di Informatica, Università di Pisa, Italy

⁴DIMIE, Università della Basilicata, Italy

⁵Universität Passau, Passau, Germany

Abstract

Many software tools, but also formal frameworks for working with JSON Schema, do not fully support negation. This motivates us to study whether negation is actually used in practice, for which aims, and whether it could, in principle, be replaced by simpler operators. We have collected a large corpus of 80k open source JSON Schema documents. We perform a systematic analysis, quantify usage patterns of negation, and also qualitatively analyze schemas. We show that negation is indeed used, albeit infrequently, following a stable set of patterns.

Keywords

Empirical Study, Conceptual Modeling, JSON Schema

1. Introduction

JSON has become one of the most popular formats for data exchange. While many schema languages for JSON have been proposed [1], JSON Schema [2] is receiving considerable attention. In this language, a schema is a logical combination of assertions, describing classes of constraints on objects, arrays, and base values. JSON Schema is constantly evolving and new drafts always introduce new features. The language is increasingly used for defining *domain-specific* data exchange formats [3] and as a meta-language for defining other languages; a subset of JSON Schema serves as the schema language inside MongoDB [4]. As a consequence, an active and quite broad development community is releasing JSON Schema tools (validators [5], in particular).


JSON Schema is powerful but complex, and its semantics is based on an intricate interplay among logical assertions. A distinctive feature is the not operator, whereby negation can be applied to any assertion. Negation is quite rare in type and schema languages, as it poses severe challenges.

SEBD 2022: The 30th Italian Symposium on Advanced Database Systems, June 19-22, 2022, Tirrenia (PI), Italy

✉ baazizi@ia.lip6.fr (M. Baazizi); dario.colazzo@dauphine.fr (D. Colazzo); ghelli@di.unipi.it (G. Ghelli); carlo.sartiani@unibas.it (C. Sartiani); stefanie.scherzinger@uni-passau.de (S. Scherzinger)



© 2022 Copyright for this paper by its authors. Use permitted under Creative Commons License Attribution 4.0 International (CC BY 4.0).

 CEUR Workshop Proceedings (CEUR-WS.org)

<pre> 1 { "not": 2 { "required": ["DisplaceModules"] } 3 } </pre>	<pre> 1 { "title" : "Object w/ required foo.", 2 "type": "object", 3 "properties": { 4 "foo": { "type": "integer" }, 5 "bar": { "type": "string" } }, 6 "patternProperties": { 7 "f.*o": { "type": "integer" } }, 8 "required": ["foo"] 9 } </pre>
<pre> 1 { "description": "...", 2 "@errorMessages": 3 { "not": "Invalid target: ..." }, 4 "not": { "pattern": "..." } ... } </pre>	

Figure 1: Snippets of JSON Schema documents.

Example 1. One usage of `not` that startles novices (as discussed on [StackOverflow \[6\]](#)) is in combination with the keyword `required`, as shown in Figure 1(a). While “not required” may sound like “optional”, it enforces that the object must violate the assertion, so member “DisplaceModules” must be absent.

Indeed, the `not`-operator is often not fully supported, whether in academic prototype tools [7], commercial tools (e.g., [4]), or even formal frameworks [8], mostly because of the inherent complexity of handling negation. This inspired us to investigate the usage of this operator in real-world schemas, in a principled analysis of 80k JSON Schema documents crawled from GitHub. We formulate these research questions: (1) *how frequent* is negation in practice, (2) *how* is negation used, and (3) *what* are common usage patterns?

Contributions. The contribution of this systematic empirical study is threefold. We first established a method for the collecting and preparing JSON Schema documents. Next, we measured the frequency of use of JSON Schema operators and of paths that include `not`, and quantify main patterns of use. Finally, we identified well-supported *jargons*, i.e., common uses of `not` that have the potential to mature into JSON Schema *design patterns*. An extended version of this study can be found here [9].

2. Preliminaries

JSON data model. The grammar below captures the syntax of JSON values, which are basic values, objects, or arrays. Basic values B include the null value, booleans, numbers n , and strings s . Objects O represent sets of members, each member being a name-value pair, and arrays A represent sequences of values.

$J ::= B \mid O \mid A$		JSON expressions
$B ::= \text{null} \mid \text{true} \mid \text{false} \mid n \mid s$	$n \in \text{Num}, s \in \text{Str}$	Basic values
$O ::= \{l_1 : J_1, \dots, l_n : J_n\}$	$n \geq 0, i \neq j \Rightarrow l_i \neq l_j$	Objects
$A ::= [J_1, \dots, J_n]$	$n \geq 0$	Arrays

JSON Schema. JSON Schema is a language for defining constraints and requirements on the content of JSON documents. We discuss here the main keywords, and continue with two illustrative examples:

Assertions include `required`, `enum`, `const`, `pattern` and `type`, and indicate a test that is performed on the corresponding instance.

Applicators include the boolean operators `anyOf`, `allOf`, `oneOf`, `not`, the object operators `properties`, `patternProperties`, `additionalProperties`, the array operator `items`, and the reference operators `$ref`. Applicators indicate a request to apply a different operator to the same instance or to a component of the current instance.

Annotations include `title`, `description`, and `$comment`, they do not affect validation, but they indicate an annotation that should be associated with the instance. Since we are mostly interested in validation, and since, moreover, annotations are removed by the `not` operator, we will ignore them.

Example 2. *In the schema in Figure 1(c), inspired from [5], line 1 carries an annotation. In defining an object (line 2), applicators define constraints on properties (lines 3), and the type of the properties matching a pattern (see line 6). Using an assertion, it is possible to indicate required properties (line 8).*

Example 3. *JSON Schema is an open standard: in Figure 1(b), `@errorMessages` is a user-defined keyword whose value is an object that describes the error, and not a JSON Schema assertion. Hence, `not` in line 3 is just a member name, whereas negation does occur in line 4. The same string token has different semantics, depending on its context, which complicates parsing.*

2.1. Pattern Queries

To study which keywords occur below an instance of the `not` operator, we introduce a simple path language. A path such as `**.not.required` matches any path that ends with an object field named `required` found inside an object field whose name is `not`. Paths are expressed using the following language. Path matching is defined as in JSONPath [10].

$$p ::= step \mid step \ p \quad step ::= .key \mid .* \mid [*] \mid .**$$

The step `.*` retrieves all member values of an object, `[*]` retrieves all items of an array, and `**` is the reflexive and transitive closure of the union of `.*` and `[*]`, navigating to all nodes of the JSON tree to which it is applied.

Complex sub-schemas. We say that `not` has a *complex* sub-schema, when its object argument contains more than one keyword. In this case, we say these keywords *co-occur* in the negated schema; otherwise, a sub-schema is *simple*. As an example, consider the schema of Figure 3(b): the argument of `not` is complex, and we match the paths `.not.enum` and `.not.type`.

3. Methodology

Context. We explored GitHub for open source JSON Schema documents. We identified 91,6k URLs in July 2020, of which 85,6k could be retrieved (using `wget`). Discarding files with invalid syntax yields 82k files.

For each retrieved file, we analyzed the `$schema` declarations to identify the version of JSON Schema. Draft 2019-09 is still quite new, and not really represented. Draft-04 is declared in the

vast majority of the files (79%), while Draft-07, Draft-06, and the old Draft-03 are each below 5%. An analysis of the file contents showed that the actual version that a schema follows is often different from the version declared.

Data Preparation. As a first step, we renamed all references (`$ref`) by a new keyword `$eref`, with the target of the reference as its child, but we did not expand references recursively. We expanded references to external documents, provided that we were able to locate the referenced document (e.g., either contained within our corpus, or by downloading the document). References were renamed to `$fref` when expansion failed. We observed that by expanding references we lose the conceptual information encoded in the reference path itself. Thus, `$ref` is often more than just a syntactic macro.

The schema corpus contains a large share of near-duplicate schemas, with small variations in syntax. We performed duplicate elimination by comparing compact *schema signatures*, defined as a function that maps each keyword to the number of its occurrences in the schema (encoded as a vector of keyword counts); we assumed that two schemas with the same *signature* are, with high probability, versions of the same schema, and we retained just one. After duplicate elimination our corpus shrunk to 11,500 distinct schemas.

As illustrated in Example 3, correctly recognizing keywords can be a challenge. For this reason, we renamed all property names to avoid confusion when searching for patterns that involve the keyword `not`. As schema authors can define their own keywords, we have no way to know whether their value should be interpreted as an assertion. We experimented with two approaches: a “strict” approach in which we renamed everything that was inside a user-defined keyword, hence making it inaccessible by the analysis, and a “lax” approach in which we kept the content of any user-defined keyword, so that all instances of `not` in Figure 1(b) would be counted as keywords. With the strict approach, some interesting usage patterns are lost, and keyword usage is under-estimated. With the lax approach, we risk “false positives”, and hence over-estimation. We decided that the over-estimation of the lax approach was preferable.

Analysis Process. The bulk of our effort is actually invested in data preparation. After experimenting with different data analysis platforms, we resorted to a relational encoding of the JSON Schema documents in PostgreSQL. This setup met our performance expectations, and allowed us to write queries in plain SQL.

4. Results of the Study

4.1. RQ1: How frequent is negation in practice?

We study the frequency of JSON Schema keywords within our corpus, and the Boolean operators (among them, negation). The reported absolute values are mainly interesting as indicators as to the relative occurrences of operators. Figure 2 visualizes the results. From left-to-right, we sort keywords by their number of occurrence (note the log-scaled vertical axes). We also show the number of files in which keywords occur, as a further indicator of keyword relevance.

The operator `not` appears in approx. 3% of all schemas, and occupies the 30th position, out of 46 keywords analyzed. Thus, it is a comparatively rare operator. The most common

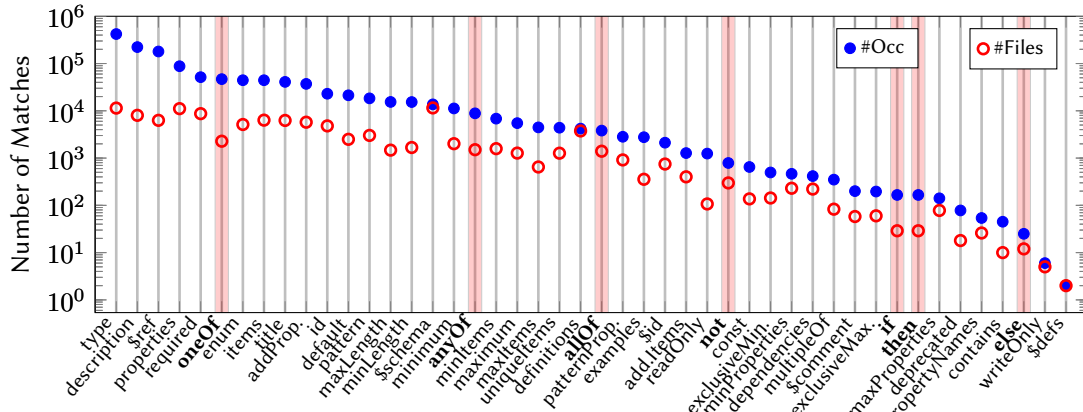


Figure 2: Number of total occurrences (#Occ), and number of files (#Files), where a JSON Schema keyword appears. Boolean operators are highlighted.

Boolean operator is `oneOf`, more frequent than `anyOf`. `allOf` is even less common. The Boolean operator `if-then-else` is even less common than `not`, but was only been introduced in Draft-07.

We found the dissemination of `oneOf` surprising, since the exclusive-disjunctive semantics of `oneOf` is more complicated than the purely disjunctive `anyOf`: `oneOf` takes as argument a collection of subschemas S_1, \dots, S_n , and a value J satisfies `oneOf` only if it matches exactly one subschema; `anyOf` is satisfied by any value J that matches at least one of the subschemas. Our hypothesis is that the description of a class as a `oneOf`-combination of a set of “subclasses” is familiar from the exclusive-subclassing mechanism of object-oriented languages.

The operator `not` appears 787 times in 298 different files out of 11,500. While not very frequent, its usage nevertheless merits a systematic study.

4.2. RQ2: How is negation used in practice?

We evaluated pattern queries to identify keywords below `not`. Table 1 summarizes the results. Consider the left half. We match the path `**.not.*` 840 times (#Occ) in 289 files (#Files). Below the top summary row, we list the individual keywords, breaking down shares of matches in percent (visualized by progress bars). The right half of the table provides statistics for subschemas that are negated and referenced, and therefore reachable via a path `**.not.$eref.*`.

In the following, we will omit the prefix `**.not.` from path queries, assuming the context is clear to our readers. We sorted the table on the total number of `not.k+not.$eref.k` occurrences, and it is interesting to compare the weight of different keywords in both parts.

A `not` may not correspond to any `not.*` pattern, when followed by `{ }`. We found 16 such occurrences, expressing the schema `false`, which is not satisfied by any instance. This use of `not` is a consequence of the fact that `false` has only been introduced with Draft-06.

Table 1 indicates a total of 840 occurrences of `not.*`, Figure 2 reported 787 occurrences of `not`. The values differ since the negated sub-schema can be complex. Most instances of `not` have a simple sub-schema. Most negated complex schemas have two keywords, but some have three or four. The situation is very different with `$eref`, i.e., references expanded in pre-processing.

Table 1
Occurrences of *not.k* paths (overall #Occ, and counting #Files).

Path	#Occ	#Files	Path	#Occ	#Files
not.*	840	289	not.\$ref.*	338	28
required	28.6 %	29.1 %	required	10.7 %	53.6 %
items	15.0 %	9.3 %	items	0.0 %	0.0 %
type	7.4 %	17.7 %	type	15.1 %	71.4 %
properties	8.5 %	16.3 %	properties	11.8 %	64.3 %
\$ref	11.1 %	9.7 %	\$ref	0.0 %	0.0 %
enum	7.3 %	18.0 %	enum	3.6 %	28.6 %
allOf	2.7 %	8.0 %	allOf	11.2 %	17.9 %
pattern	5.6 %	9.7 %	pattern	0.0 %	0.0 %
anyOf	5.4 %	12.5 %	anyOf	0.6 %	7.1 %
description	0.5 %	1.4 %	description	12.1 %	25.0 %
title	0.2 %	0.7 %	title	11.5 %	25.0 %
\$schema	0.0 %	0.0 %	\$schema	12.1 %	32.1 %
\$fref	3.2 %	4.8 %	\$fref	0.0 %	0.0 %
oneOf	0.7 %	1.4 %	oneOf	5.3 %	10.7 %
additionalProperties	1.3 %	3.8 %	additionalProperties	2.7 %	25.0 %
patternProperties	1.8 %	5.2 %	patternProperties	0.0 %	0.0 %
const	0.7 %	0.4 %	const	0.0 %	0.0 %
definitions	0.0 %	0.0 %	definitions	0.9 %	10.7 %
id	0.0 %	0.0 %	id	0.6 %	7.1 %
dependencies	0.0 %	0.0 %	dependencies	0.6 %	7.1 %
not	0.0 %	0.0 %	not	0.6 %	7.1 %
\$ref	0.0 %	0.0 %	\$ref	0.6 %	7.1 %
\$comment	0.1 %	0.4 %	\$comment	0.0 %	0.0 %

Here, 93 occurrences of *not.\$ref* correspond to 338 occurrences of *not.\$ref.**. Thanks to the mediation of *\$ref*, the schema designer implicitly applies negation to a complex argument, with an average of 3-4 members.

The most common argument of negation is *required*. The pattern *not.items* is second-most common, followed by *not.type* and *not.properties*.

While *not.required* dominates the *not.** case, the two most common cases of the *not.\$ref* group are *not.\$ref.type*, whose value is *object* in 80% of the cases, and *not.\$ref.properties*, which indicates that *not.\$ref* is mostly used to negate complex object definitions. This explains the much higher occurrence of descriptive keywords inside the referenced argument.

4.3. RQ3: What are common real-world usage patterns?

Field and value exclusion. Field exclusion via *not.required* is the most frequent path.

Paths *not.enum* and *not.const* are used to exclude values. Snippets of example schemas

(a)	<pre>"not": { "enum": ["markdown", "code", "raw"] }</pre>	(b)	<pre>{ "type": "object", "oneOf": [{ "properties": { "when": {"enum": ["delayed"]}}, "required": ["when", "start_in"] }, { "properties": { "when": { "not": {"enum": ["delayed"]}} }]] }</pre>
(b)	<pre>"not": { "enum": ["generic-linux"], "type": "string" }</pre>	(c)	<pre>{ "type": "object", "if": { "required": ["when"], "properties": { "when": {"enum": ["delayed"]}} }, "then": { "properties": { "when": {"enum": ["delayed"]}} }, "required": ["when", "start_in"] }</pre>
(c)	<pre>"not": { "items": { "not": { "type": "string", "enum": ["Dataset", "Image", "Video", "Sound", "Text"] } } }</pre>	(d)	<pre>{ "type": "object", "if": { "required": ["when"], "properties": { "when": {"enum": ["delayed"]}} }, "then": { "properties": { "when": {"enum": ["delayed"]}} }, "required": ["when", "start_in"] }</pre>

Figure 3: JSON Schema snippets exemplifying real-world usage patterns.

are shown in Figures 3(a) and (b). Such schemas have an obvious interpretation: the instance may have any type and must be different from the string or strings listed. In the majority of cases, the sub-schema is simple, as in Figure 3(a). In the complex cases, `enum` is always paired with a `"type": "string"` assertion, as in Figure 3(b). This assertion is redundant, since all values listed by `enum` are strings. This co-occurrence is not specific to negation, since also in positive schemas, `enum` is paired with a type assertion in the vast majority of cases.

Paraphrasing contains. The pattern `not.items` is among the most common `not`-paths. All such schemas have either the structure `not.items.not` (as in Figure 3(c)) or `not.items.enum`.

The `items` assertion is verified by any instance that is not an array, or that is an empty array, or that is an array where every element satisfies the schema associated with `items`. Hence, it is only violated by instances that are arrays, and which contain at least one element that violates the schema. While `items` specifies a universally quantified property, `not.items` can be used to specify an existentially quantified property, as does the `contains` keyword. The jargon `not.items.enum` specifies that the array must contain at least one value that is not listed in the argument of `enum`. The jargon `not.items.not` specifies that the instance is an array that contains at least one value that satisfies S , according to the following equivalence:

$$\text{"not": {"items": {"not": S}}} \Leftrightarrow \{\text{"type": "array", "contains": S}\}$$

These two cases cover, with minimal variations, all occurrences of `not.items`.

To sum up, `not.items` can be used to express `contains`. This is an instance of a pattern that may be replaced by a single (and thus simpler) operator.

Paraphrasing Discriminated Unions. The schema snippet in Figure 3(d) allows interesting observations about the use of `oneOf`. JSON Schema specifications do not prescribe that the branches of `oneOf` are mutually exclusive, but they state that a value must match a single branch only. However, the two branches of `oneOf` happen to be mutually exclusive: if `"when"` is absent, then only the second branch holds. If it is present, then it is associated to complementary types in the two branches, so here, `oneOf` is actually `anyOf`. Applying equivalent rewritings (from $\neg a \vee b$ to $a \Rightarrow b$, and pushing down negation), the schema can be rewritten as shown in Figure 3(e).

Now the specification is clearer: if "when" has the value "delayed", then "start_in" is required.

This suggests that oneOf is used to express a form of *discriminated unions*.

References

- [1] M. A. Baazizi, D. Colazzo, G. Ghelli, C. Sartiani, Schemas and types for JSON data: From theory to practice, in: Proc. SIGMOD 2019, 2019, pp. 2060–2063.
- [2] json-schema org, JSON Schema, 2021. Available at <https://json-schema.org>.
- [3] B. Maiwald, B. Riedle, S. Scherzinger, What Are Real JSON Schemas Like? – An Empirical Analysis of Structural Properties, in: Proc. EmpER 2019, 2019, pp. 95–105.
- [4] MongoDB, Inc., MongoDB Manual: \$jsonSchema (Version 4.4), 2021.
- [5] JSON Schema Test Suite, Available at: <https://github.com/json-schema-org/JSON-Schema-Test-Suite>, version of commit hash #09fd353., 2021.
- [6] StackOverflow, JSON Schema – valid if object does *not* contain a particular property, Available at: <https://stackoverflow.com/questions/30515253/json-schema-valid-if-object-does-not-contain-a-particular-property>, 2015.
- [7] M. Fruth, M. A. Baazizi, D. Colazzo, G. Ghelli, C. Sartiani, S. Scherzinger, Challenges in Checking JSON Schema Containment over Evolving Real-World Schemas, in: Proc. EmpER 2020, 2020, pp. 220–230.
- [8] A. Habib, A. Shinnar, M. Hirzel, M. Pradel, Finding data compatibility bugs with JSON subschema checking, in: Proc. ISSTA 2021, 2021, pp. 620–632.
- [9] M. A. Baazizi, D. Colazzo, G. Ghelli, C. Sartiani, S. Scherzinger, An empirical study on the “usage of not” in real-world JSON schema documents, in: Proceedings of ER 2021, October 18-21, 2021, 2021, pp. 102–112.
- [10] J. Friesen, Java XML and JSON: Document Processing for Java SE, Apress, 2019, pp. 299–322.