

JSON-TAB Format

Presentation

18/06/2023

Environmental Sensing



TABLE OF CONTENTS

| | |
|---|-----------|
| 1 Introduction | 2 |
| 1.1 Conventions used | 2 |
| 1.2 Terminology | 3 |
| 1.3 Rules | 3 |
| 2 Objectives | 3 |
| 2.1 Why a new textual standard for tabular data ? | 3 |
| 2.2 Key design features | 3 |
| 3 Tabular data representation | 4 |
| 3.1 Tabular data | 4 |
| 3.2 Representation of fields | 5 |
| 3.3 Level of representation | 7 |
| 4 JSON-TAB representation | 9 |
| 4.1 Structure | 9 |
| 4.2 JSON format | 10 |
| 5 Examples | 10 |
| 5.1 Field example | 10 |
| 5.2 Tabular object examples | 12 |
| 6 Parsing a JSON-value | 13 |
| Appendix : reserved values | 14 |
| Appendix : CBOR format | 14 |

1 INTRODUCTION

The JSON-TAB format is applicable to any tabular, indexed or multi-dimensional data.

This format makes it possible to:

- integrate data of a high semantic level,
- interoperability with all JSON parsers,
- avoid data duplication,
- reduce the size of data,
- integrate meta-data

This format is an application of the JSON-NTV format.

A binary version is also defined (Appendix) with CBOR format (RFC 8949)

1.1 CONVENTIONS USED

The key words "MUST", "MUST NOT", "REQUIRED", "SHALL", "SHALL NOT", "SHOULD", "SHOULD NOT", "RECOMMENDED", "NOT RECOMMENDED", "MAY", and "OPTIONAL" in this document are to be interpreted as described in BCP14 [RFC2119] [RFC8174] when, and only when, they appear in all capitals, as shown here.

The grammatical rules in this document are to be interpreted as described in [RFC5234].

1.2 TERMINOLOGY

The terms Json-Text, Json-Value (Value), Object, Member, Element, Array, Number, String, False, Null, True are defined in the JSON grammar.

The terms Geometry-type, Point, MultiPoint, LineString, MultiLineString, Polygon, MultiPolygon, GeometryCollection, GeoJSON-Types are defined in GeoJSON grammar.

The terms Table, Row, Column, Cell, Datatype are defined in W3C Recommendation 17 December 2015 : "Model for Tabular Data and Metadata on the Web".

The term Field is equivalent to Column.

Timestamp is defined in Date and Time format.

1.3 RULES

Values in Array are ordered and independent from the other Values.

Members in Objects are not ordered.

2 OBJECTIVES

2.1 WHY A NEW TEXTUAL STANDARD FOR TABULAR DATA ?

The main operational standard used to exchange textual tabular data is CSV format (RFC 4180).

Unfortunately CSV format is obsolete (last revision in 2005) and current CSV tools do not comply with the standard.

2.2 KEY DESIGN FEATURES

The format's focus is on simplicity, lightness and web usage.

The key features of this format are the following:

- JSON as the base format
 - JSON is simple and readable as simple text
 - JSON supports rich structure including nesting and basic types
 - JSON is web-native and very widely used and supported
 - JSON format has binary representation (i.e. CBOR format)
- JSON-NTV as a grammar
 - data enriched with name and type
- column-oriented structure (instead of row-oriented)
 - columns carry the semantics of the data
 - columns can be optimized based on their inter-relationships
- several representations available
 - four levels, from the simplest to the most optimized, are available
 - meta-data (header or schema) can be integrate
- high semantic level of data
 - all JSON representations can be included following JSON-NTV format

3 TABULAR DATA REPRESENTATION

3.1 TABULAR DATA

Tabular data is data that is structured into rows, each of which contains information about some things. Each row contains the same number of cells (although some of these cells may be empty), which provide values of properties of the thing described by the row. In tabular data, cells within the same column provide values for the same property of the things described by each row. This is what differentiates tabular data from other line-oriented formats. (Model for Tabular Data and Metadata on the Web - W3C - Recommendation 17 december 2015)

In tabular data, column and row are not equivalent, the columns (or fields) represent the 'semantics' of the data and the rows represent the objects arranged according to the structure defined by the columns.

If we now observe how tabular data are used, we can identify four main uses:

- **association**: this consists of coupling each value of a field to a value of another field ("coupled" relationship between two fields),
- **classification**: This involves grouping the data by category in order to be able to make a statistical use of it, for example ("derived" relationship between two fields),
- **crossing**: This consists of representing all the combinations between several fields, such as in matrix representations ("crossed" relationship between several fields),
- **characterization**: It corresponds to the documentation of defined properties (no specific relationship).

Example :

| id | produit | aliment | contenant | quantité | prix | validité | disponibilité |
|----|---------|---------|-----------|----------|------|---------------------------|---------------|
| 11 | pomme | fruit | sachet | 1 kg | 1 | du 1/7/2022 au 31/12/2022 | oui |
| 12 | pomme | fruit | carton | 10 kg | 9 | du 1/7/2022 au 31/12/2022 | oui |
| 13 | orange | fruit | sachet | 1 kg | 2 | du 1/7/2022 au 31/12/2022 | fin 2022 |
| 14 | orange | fruit | carton | 10 kg | 18 | du 1/7/2022 au 31/12/2022 | fin 2022 |
| 15 | piment | légume | sachet | 1 kg | 1.5 | du 1/7/2022 au 31/12/2022 | fin 2022 |
| 16 | piment | légume | carton | 10 kg | 13 | du 1/7/2022 au 31/12/2022 | fin 2022 |
| 17 | banane | fruit | sachet | 1 kg | 0.5 | du 1/7/2022 au 31/12/2022 | oui |
| 18 | banane | fruit | carton | 10 kg | 4 | du 1/7/2022 au 31/12/2022 | oui |

This is a price list of different foods based on packaging for the year 2022.

We find here:

- *association: between "contenant" and "quantité",*
- *classification: between "produit" and "aliment",*
- *crossing: between "produit" and "quantité",*
- *characterization: between "produit" and "disponibilité"*

3.2 REPRESENTATION OF FIELDS

A tabular object can be represented by a simple list of fields where a field represents a column. Fields can be defined by a name, a type and a list of values.

Each field has the same number of values (the length of the tabular object).

The values of a field can be represented by several formats:

Full format :

The "Full format" is the usual representation of a list of values.

Example (field "prix") :

[1, 9, 2, 18, 1.5, 13, 0.5, 4]

Complete format :

The "Full format" representation has the disadvantage of being bulky when data is duplicated in a field.

The second format is to represent a list of values by two lists:

- *codec: different values,*
- *keys: indexes of values.*

The values are reconstituted by replacing the integers in the "keys" list with the corresponding values from the "codec" list.

Example ("produit" field):

[["orange" , "piment" , "pomme" , "banane"], [2, 2, 0, 0, 1, 1, 3, 3]

Unique format :

This representation corresponds to a field composed of a single duplicate value.

The "unique format" representation consists in representing only this unique value.

It is therefore a "complete format" representation in which the "keys" list is implicit.

Example ("validité" field):

"du 1/7/2022 au 31/12/2022" (implicit Keys [0, 0, 0, 0, 0, 0, 0, 0])

Note:

This format also makes it possible to represent tabular metadata

Implicit format :

This representation is associated with "coupled" fields. These fields have a one-to-one correspondence.

The representation consists of one list and one single value:

- codec: different values,
- parent: reference (row or name) to the associated ("parent") field.

Example ("quantité" field is associated with "contenant" field) :

[["1 kg", "10 kg"], 3] (implicit Keys: "keys" list of the "contenant" field)

Relative format :

This representation is associated with "derived" fields.

The values of a "derived" field are inferred from the values of the "parent" field.

The representation consists of two lists and one single value:

- codec: different values,
- parent: reference (row or name) to the associated ("parent") field,
- keys: indexes of values

Example ("aliment" field is associated with "produit" field) :

Codec : ["fruit", "légume"]

Parent : 1 (field n° 1 : "produit")

Keys : [0, 1, 0, 0] (the absolute "keys" list is obtained by replacing the values 0, 1, 2, 3 of the "keys" list of the "product" field by 0, 1, 0, 0 i.e.: [0, 0, 0, 0, 1, 1, 0, 0])

Primary format :

This representation is associated with "crossed" fields. The values of a "crossed" field are calculated from the "codec" list and a "repetition coefficient".

The representation consists two lists:

- codec: different values,
- coefficient: list with a single integer

Example "contenant" (this field is associated to "produit" field) :

Codec : ["carton", "sachet"]

Coefficient: 1 ("contenant" is the 2nd field of type "crossed")

Keys : implicit ([0, 1, 0, 1, 0, 1, 0, 1])

Example "produit" (this field is associated to "contenant" field) :

Codec : ["pomme" , "orange" , "piment" , "banane"]

Coefficient: 2 ("produit" is the 1st field of type "crossed")

Keys : implicit ([0, 0, 1, 1, 2, 2, 3, 3])

3.3 LEVEL OF REPRESENTATION

Three levels, from the simplest to the most optimized, are available to convert tabular data in JSON structure.

- Level 1 : "full"

The fields are converted into "full format" or "unique format". A tabular object has a single full representation.

- Level 2 : "default"

The fields are converted into "full format", "complete format" or "unique format"

Nota : some fields can be converted into "full format", others into "complete format"

- Level 3 : "optimize"

This level requires an analysis of the relationships between fields.

"primary" fields ("crossed" fields that form a partition) are converted into "primary format". Other fields are converted into "full format", "unique format", "complete format", "implicit format" or "relative format" according to their position in the tree.

Level 1 is the usual representation of tabular data.

Level 2 avoids duplication of information by adding simple encoding.

Level 3 avoids duplication of information and minimizes encoding.

Several representations are available for a tabular object at level 2 or 3.

Fields values structure :

| level | | Structure | | Codec | | parent | | keys / coef | | |
|-------|----------|--------------|----------|--------------|--------------|-----------------|-----------------|------------------|------------------|---------------|
| n° | mode | Type index | format | = len parent | < len parent | implicit parent | explicit parent | Absolute (= len) | Relative (< len) | implicit |
| all | | unique | unique | | x (1) | x (root) | | | | x (list of 0) |
| 1 | full | Not unique | full | x | | x (root) | | | | x (range) |
| 2 | default | Full codec | full | x | | x (root) | | | | x (range) |
| | | Reduce codec | complete | | x | x (root) | | x | | |
| 3 | optimize | Root coupled | full | x | | x (root) | | | | x (range) |
| | | Root derived | complete | | x | x (root) | | x | | |
| | | derived | relative | | x | | x | | x | |
| | | coupled | implicit | x | | | x | | | x (parent) |
| | | primary | primary | | x | x (root) | | | x (coef) | |

4 JSON-TAB REPRESENTATION

4.1 STRUCTURE

A tabular object is defined by a name (optional) and a list of Fields with the same length.

Fields are defined by a name (optional), a type and values.

Values have one mandatory member (codec) and two optional (parent and keys/coef) as defined in chapter 3. Values can have a specific name and type.

At level 1 ('full'), Keys and Parent data is not used.

At level 2 ('default') Parent is not used.

The structure of Fields depends on the representation defined in chapter 3.2.

4.2 JSON FORMAT

With the NTV format, the objects are:

- tabular object: NV-list of Fields
- Field: NTV-list or NTV-single (if unique)
- Codec: TV-list of Values or one Value
- Value: NTV-entity
- Parent : V-single number or string
- Keys/coef: V-list of integer data

NTV fields values structure :

| Structure | | parent | | keys | | |
|-----------|--------------------|----------|----------|---------------------|---------------------|---------------|
| format | NTV | implicit | explicit | Absolute (= len) | Relative (< len) | implicit |
| full | NTV-list of values | x (root) | | | | x (range) |
| unique | NTV-single | x (root) | | | | x (list of 0) |
| complete | NV-list len=2 | x (root) | | x | | |
| relative | NV-list len=3 | | x | | x | |
| implicit | NV-list len=2 | | x | | | x (parent) |
| primary | NV-list len=2 | x (root) | | | x (coef) | |

Note:

If Codec is an NTV-single or an empty TV-list, Parent and Keys are not present.

If Parent and Keys are not present, Codec and Iindex are merged.

The JSON format of a tabular object is the JSON-NTV format.

For better readability, the JSON-value for Iindex MAY be separated by a line Separator '\n'.

5 EXAMPLES

5.1 FIELD EXAMPLE

The examples in chapter 3.2 have the following JSON format:

| Format | Representation |
|----------|--|
| Full | <pre>[1, 9, 2, 18, 1.5, 13, 0.5, 4]</pre> <pre>{ "prix": [1, 9, 2, 18, 1.5, 13, 0.5, 4] }</pre> <pre>{ "prix::float": [1, 9, 2, 18, 1.5, 13, 0.5, 4] }</pre> |
| Complete | <pre>[["orange", "piment", "pomme", "banane"], [2, 2, 0, 0, 1, 1, 3, 3]]</pre> <pre>{ "produit": [["orange", "piment", "pomme", "banane"], [2, 2, 0, 0, 1, 1, 3, 3]] }</pre> <pre>{ "produit": [{ "::string": ["orange", "piment", "pomme", "banane"] }, [2, 2, 0, 0, 1, 1, 3, 3]] }</pre> |
| Unique | <pre>"du 1/7/2022 au 31/12/2022"</pre> <pre>{ "validité": "du 1/7/2022 au 31/12/2022" }</pre> |
| Implicit | <pre>[["1 kg", "10 kg"], 3]</pre> <pre>{ "quantité": [{ "::string": ["1 kg", "10 kg"] }, "contenant"] }</pre> |
| Relative | <pre>[["fruit", "légume"], 1, [0, 1, 0, 0]]</pre> <pre>{ "aliment": [{ "::string": ["fruit", "légume"] }, "produit", [0, 1, 0, 0]] }</pre> |
| Primary | <pre>[["carton", "sachet"], [1]]</pre> <pre>{ "contenant": [{ "::string": ["carton", "sachet"] }, [1]] }</pre> <pre>[["pomme", "orange", "piment", "banane"], [2]]</pre> <pre>{ "produit": [{ "::string": ["pomme", "orange", "piment", "banane"] }, [2]] }</pre> |

```
{ "age": 25 }
```

Field with name and codec

```
{ "measure": [ 2.4, 48.9 ] }
```

Field with name and codec

```
[ [ 2.4, 48.9 ], [ 0, 0, 1, 1 ] ]
```

Field with codec and Keys

```
{ "measure": [ [ 2.4, 48.9 ], [ 0, 0, 1, 1 ] ] }
```

Field with name, codec and Keys

```
{ "::geopoint": [ [ 2.4, 48.9 ], [ 4.2, 84.9 ] ] }
```

Field with type and codec

```
{ "city:geopoint": [ 2.4, 48.9 ] }
```

Field with type, name and codec

```
{ "city": [ { "::geopoint": [ [ 2.4, 48.9 ], [ 4.2, 84.9 ] ] }, 0, [ 0, 0, 1, 1 ] ] }
```

Field with all data

5.2 TABULAR OBJECT EXAMPLES

The examples below illustrate the JSON-NTV format with 'full level' and 'optimize level'.

| Data | full level | optimize level |
|----------------------------|---|---|
| Matrix | [[['a', 'a', 'b', 'b', 'c', 'c'], [10, 20, 10, 20, 10, 20], [1, 2, 3, 4, 5, 6]] | [[['a', 'b', 'c'], [2]], [[10, 20], [1]], [1, 2, 3, 4, 5, 6]] |
| Single | [[1, 2, 3, 4, 5, 6], ['a', 'a', 'a', 'a', 'a', 'a']] | [[1, 2, 3, 4, 5, 6], ['a']] |
| Complete | [1, 2, 3, 3, 5, 5] | [[1, 2, 3, 5], [0, 1, 2, 2, 3, 3]] |
| Coupled | [[1, 2, 3, 3, 5, 5], ['a', 'b', 'c', 'c', 'e', 'e']] | [[[1, 2, 3, 5], [0, 1, 2, 2, 3, 3]], [['a', 'b', 'c', 'e'], 0]] |
| Derived | [[1, 2, 3, 4, 5, 6], ['a', 'a', 'b', 'b', 'c', 'c'], [10, 10, 10, 10, 20, 20]] | [[1, 2, 3, 4, 5, 6], [['a', 'b', 'c'], [0, 0, 1, 1, 2, 2]], [[10, 20], 1, [0, 0, 1]]] |
| Matrix + coupled | [[['a', 'a', 'b', 'b', 'c', 'c', 'd', 'd'], [10, 20, 10, 20, 10, 20, 10, 20], ['t1', 't1', 't2', 't2', 't3', 't3', 't4', 't4'], [1, 2, 3, 4, 5, 6, 7, 8]] | [[[['a', 'b', 'c', 'd'], [2]], [[10, 20], [1]], [['t1', 't2', 't3', 't4'], 0], [1, 2, 3, 4, 5, 6, 7, 8]] |
| Matrix + coupled + derived | [[['a', 'a', 'b', 'b', 'c', 'c', 'd', 'd'], [10, 20, 10, 20, 10, 20, 10, 20], ['t1', 't1', 't2', 't2', 't3', 't3', 't4', 't4'], [100, 100, 100, 100, 200, 200, 200, 200], [1, 2, 3, 4, 5, 6, 7, 8]] | [[[['a', 'b', 'c', 'd'], [2]], [[10, 20], [1]], [['t1', 't2', 't3', 't4'], 0], [[100, 200], 0], [1, 2, 3, 4, 5, 6, 7, 8]] |

The examples below show how to represent tabular data with a single value.

[]

Empty tabular data

[25] or [[25]]

data with 1 Field with 1 codec value

[2, 1] or [[2], [1]] or [2, [1]]

data with 2 Fields with 1 codec value

[[2, 1]]

data with 1 Field with 2 codec values

[[2, 1], [4,3]]

data with 2 Fields with 2 codec values

6 PARSING A JSON-VALUE

A NTV parser generates an NTV entity from a JSON-value.

A TAB decoder generates a tabular data object from an NTV entity.

Several dataset can generate inconsistent data:

- [listdata, integer, listinteger]
- [listdata, string, listinteger]
- [listdata, integer]
- [listdata, string]
- [listdata, listinteger]

To avoid this inconsistency, the order of data can be changed (e.g. [integer, listdata] or a name can be added (e.g. { 'name': [listdata, listinteger] }) or a null value can be added (e.g. [listdata, integer, listinteger, "null"]).

APPENDIX : RESERVED VALUES

to complete

APPENDIX : CBOR FORMAT

The Concise Binary Object Representation (CBOR – RFC8949) is a data format whose design goals include the possibility of extremely small code size, small message size, and extensibility without the need for version negotiation.

CBOR is based on the JSON data model: numbers, strings, arrays, maps (called objects in JSON), and a few values such as false, true, and null.

The CBOR format can be used with different options to minimize length:

- The precision of float values is adjustable from half precision (two bytes) to double precision (eight bytes),
- The datetime can be described by a standard text string (RFC3339) or by a numerical value (Epoch-based: six bytes).
- The TypeValue can be represented with a code value instead of string value
- The coordinates value can be described with integer instead of float ($\text{val_int} = \text{round}(\text{val_float}) \times 10^{**7}$: four bytes).

Example (Json format):

```
{ "type": "observation",
  "datation": ["2021-01-04T10:00:00", ["2021-01-05T08:00:00", "2021-01-5T12:00:00"]],
  "location": [[2.4123456, 48.9123456], [[[2.4123456, 48.9123456], [4.8123456, 45.8123456], [5.4123456, 43.3123456], [2.4123456, 48.9123456]]]],
  "property": [{"prp": "PM10"}, {"prp": "Temp"}],
  "result": [51.348, {"low": 2.457}, 20.88, "high"],
  "coupled": {"datation": "location"}}
```

Example optimized (Cbor format):

```
{0 : [0,1,2],
 1 : [[dt(2021, 1, 4, 10),[dt(2021,1,5,8), dt(2021, 1, 5, 12)]],
      [[2.4123456, 48.9123456], [[[2.4123456, 48.9123456], [4.8123456, 45.8123456], [5.4123456, 43.3123456], [2.4123456, 48.9123456]]]],
      [{"prp": "PM10"}, {"prp": "Temp"}] ],
 2 : [51.34375, {"low": 2.45703125}, 20.875, "high"],
 3 : {0: 1} }
```

With :

- Observation key codification:

- o 0: "order"*
 - o 1: "features"*
 - o 2: "result"*
 - o 3: "coupled"*
- *Order and coupled value codification:*
 - o 0: "datation"*
 - o 1: "location"*
 - o 2: "property"*
- *Datation value: timestamp format*
- *Location value: integer representation (four bytes)*
- *Result value: half precision (two bytes)*

Length (bytes):

- JSON: 388
- CBOR: 298
- CBOR optimized: 133