

# Draft OAGi Interoperable Mapping Specification

OAGi Members

2022/01/28

## 1 Introduction

This document describes a data mapping language designed to serve as an *interoperable exchange form* for expressing the intent of many mapping and data restructuring needs. As an interoperable exchange form, it is intended that the language can be translated (by humans and machine) into mapping specification in other languages. For example, it should be possible to translate statements in the exchange form into mapping specification used by commercial mapping tools.

The document is a draft and as of this writing (2022/01/28) is likely to be updated often. OAGi members can find the most recent version of this document in the *OAGi Mapping Specification Working Group* Confluence pages (under "Mapping Doc").

The mapping language described borrows predominantly from the JSONata mapping language [1], but also includes provisions for mapping to/from forms other than JSON. These forms include tables (e.g. Excel) XML, and networks of data. To support networks of data, the mapping language borrows ideas from the Object Management Group's Queries, Views and Transformation relational (QVT-r) [2] mapping language, and Datalog [3].

## 2 Quick start: example mapping tasks

In order to give you a sense of things, the next subsections describe the basics and how some common query, data restructuring, and mapping challenges are addressed by the language.

### 2.1 Simple operations

### 2.2 Vectors and Iteration

#### 2.2.1 Mapping over a vector

In the language, iteration is typically performed implicitly in the sense that when the argument to an operation is an array, the operation is applied to each element of the collection. The result returned is an array of the result of applying the operator to each element. Thus in Figure 1, `.zipcode` is applied to the array `$ADDRS` returning an array `["20898", "07010-35445", "10878"]`.

Figure 1: Map over an array of entities, collecting ZIP code.

```
(
  $ADDRS :=
    [{"name"      : "Peter",
      "street"    : "123 Mockingbird Lane",
      "zipcode"   : "20898"},

     {"name"      : "Bill",
      "street"    : "23 Main Street",
      "zipcode"   : "07010-3545"},

     {"name"      : "Lisa",
      "street"    : "903 Forest Road",
      "zipcode"   : "10878"}];

  $ADDRS.zipcode
)
```

In the case that the operator returns null for elements of the array, no value is collected. For example, if “Bill” in the above didn’t have a zipcode, `$ADDRS.zipcode` would have returned `["20898", "10878"]`.

## 2.2.2 Filtering an array

Suppose, for some odd reason, we didn’t like 9-digit Zip codes (what the US Postal Service calls “Zip+4”) and that we simply want to exclude them from the collection. There are a few simple ways to do this. One of those is to append a test expression to the end of the previous expression. This “filtering” test expression needs to be enclosed in square brackets to indicate that it is filtering. `$ADDRS.zipcode[$match(/^\d{5}$/)]`. Regular expression syntax is arcane but what is shown is typical of JavaScript, Java and many other regex libraries. The above matches on strings that are exactly 5 number characters, thus excluding the Zip+4 Zip Codes.

## 2.2.3 Index variables

Index variables are variables used to iterate through arrays. They are common in procedural code, for example the variable `i` in the pseudocode `for i in [1..9] {...}` is an index variable.. A goal of interoperable exchange is to allow the expression of individual mapping statements that can be defined independent of larger program context. This kind of independence is useful in conventional *mapping tables*, spreadsheet used to describe to programmers the intent of proposed mapping. The problem with index variables in this regard is that they would add meandering procedural expression where concise, independent, declarative expression is desired. But what if you needed to enumerate the things you were working with, for example, to put an index number in front of them? In cases like this, you can avoid the need for an index variable by using the language’s `$map` higher-order function; it is like JavaScript’s `map`. See Figure 2 below.

Figure 2: Using an index variable.

```
$map($ADDRS.zipcode, function($v, $i) {'zip ' & ($i+1) & ' ' & $v})
returns

[ 'zip 1 20898', 'zip 2 07010-3545', 'zip 3 10878' ].
```

In the above, `function` introduces a function, and the following code in brackets is the body of the function, a single expression is returned. In `$map` (and `$filter`, and `$reduce` described later), the function is called once each with each member of the first argument (`$ADDRS.zipcode` here) in sequence

as the value bound to `$v`, and `$i` is bound in turn to successive integers in the sequence `[1, 2, 3]`. This effectively eliminates the need for user-defined index variables.

## 2.3 Working with code lists

A mapping tasks commonly needed is translating a code list, or subsetting one to a set of code values that actually is used in your business processes. For example, there is a very large set of codes for transportation events in the code list **[don't recall the spec]**. Suppose you are receiving messages with lots of these values but your order management application only wants to know whether the shipment is likely to be late. There is no getting around the fact that you need to explicitly describe a functional relationship to do this; the functional relationship can be implemented as a lookup table. The language defines a switch statement to do this kind of thing. An example is shown in Figure 3.

Figure 3: Mapping code values.

---

```

1  $Codes2OurCodes := function($theirs)
2    {switch $theirs {
3      #{2, 5, 52, 66} : 'late';
4      #{3, 6, 77, 85} : 'on time';
5      * : 'unknown'}
6  };

```

---

The code in Lines 3 and 4 of Figure 3 implement a table. The `#{...}` syntax defines sets. On Line 3 the source defines set of code values (2, 5, 52, and 66) and its use indicates that elements of this set are mapped to the target value “late”.

Were the table to be larger, writing it could get tedious. For this reason, it might be reasonable to provide an alternative, whereby a lookup table provided elsewhere is referenced.<sup>1</sup>

## 2.4 Working with tabular data

Being able to read information from a spreadsheet is a very handy capability in mapping. Of course, Excel-like spreadsheets can have multiple sheets and the content can be non-uniform, including merged cells and formula. The language does not provide means to deal with these complexities. However simple tables with no surprises and a header row naming columns (and common-separated value (CSV) files that conform to these requirements) can easily be viewed as an array of map structures. For example, Table 1 can be viewed as the structure shown in Figure 4.

Table 1: A simple table oriented such that columns name properties.

| ShipDate | Item       | Qty | UnitPrice |
|----------|------------|-----|-----------|
| 6/15/21  | Widget 123 | 1   | \$10.50   |
| 6/15/21  | Gadget 234 | 2   | \$12.80   |
| 6/15/21  | Foobar 344 | 1   | \$100.00  |

Figure 4: Table 1 viewed as an array of map structures.

```

[{"ShipDate": "2021-06-15", "Item": "Widget 123", "Qty": 1.0, "UnitPrice": 10.5},
 {"ShipDate": "2021-06-15", "Item": "Gadget 234", "Qty": 2.0, "UnitPrice": 12.8},
 {"ShipDate": "2021-06-15", "Item": "Foobar 344", "Qty": 1.0, "UnitPrice": 100.0}]

```

---

<sup>1</sup>The code in Figure 3 also shows how names are associated with functions. Functions in the language do not have names. If you need to use a function in multiple places, assign it to a variable as shown in the figure, (use of `$Codes2OurCodes`).

The built-in function `$readSpreadsheet` reads spreadsheets. An example usage is:  
`$readSpreadsheet("data/spreadsheets/ExampleInvoiceInfo.xlsx", "Sales Info")`

where the first argument names an Excel file and the second names a sheet in that spreadsheet. If the table is transposed (so that all the properties are in its first column, and each row concerns a different object), a third argument with value `true` can be specified to access the data in the more useful orientation.

## 2.5 Working with data having complex interrelations

Some situations call for more complex query and restructuring functionality than what has been illustrated above. Everything to this point in the discussion could be performed easily with JSONata [1], a language for querying and restructuring JSON-like data. JSONata is used, for example, in the Open Integration Hub [4], an open software platform that received its initial funding in 2017 from the German Federal Government, Federal Ministry for Economic Affairs and Energy. JSONata is also used in NodeRed [5].

The JSONata viewpoint could be described as one where a tree (or equivalently, a JSON object) is navigated from the root. Decades of experience, dating back to the early days of EDI, has shown that tree-based organization such as JSON objects is a quite reasonable choice where the communication of document-like information (e.g. “messaging”) is needed. The point of an interoperable exchange for mapping such as RADmapper<sup>2</sup>, however, includes additional requirements. Among these is the ability to describe relationships to and from data possessing complex interrelationship, for example, data described by a relational schema or knowledge graph. To argue that APIs to the back-end system already do this work is to miss the point: we are not trying to replace a back-end system function, but to describe the relationships in ways that help business analysts, programmers, and machine agents.<sup>3</sup>

To illustrate the challenge of using a JSONata-like engine for data having complex interrelations, suppose, for example that you have a body of data describing a large number of products produced by a small number of vendors. With this data, you... [Later]

### 2.5.1 Data organized as triples

Relational and graph-based data can be described by triples  $[x, rel, y]$  where  $x$  is an entity identifier,  $y$  is data (string, number, object, array, another entity identifier, etc.)<sup>4 5</sup> and  $rel$  is a relationship (predicate) holding between  $x$  and  $y$ .

$rel$  might be implemented as a string. The translation of structured data such as JSON objects into such triples is relatively straightforward:

1. A unique entity identifier is associated with each object; this identifier serves as the first element of triples about that object. In JSON, for example, an object is a collection of name-value pairs delimited by curly brackets.
2. Each object attribute is represented by the string naming it<sup>6</sup>; these supply the second element of the triple, the  $rel$ .

---

<sup>2</sup>RADmapper is my shot at a name. RAD=rapid application development. Suggest a name!

<sup>3</sup>Likewise we are going to try to address situations where there are multiple sources. There are tools to do this already, but again, we are trying to describe the relationships, not replace existing software.

<sup>4</sup>What additional atomic data types shall we provide? Dates? URIs? UUIDs? etc.

<sup>5</sup>Triples provide 5th-normal form relational data. Recomposing a concept from 5-th normal form typically requires as many joins as there are attributes to recombine. The upside of using triples are that (1) substantial amounts of data are already in triples (e.g. RDF etc.), (2) some query engines helpfully index triples in multiple ways entity-attribute-value, attribute-entity-value, and attribute-value-entity, and (3) mapping engines can use triples effectively. The challenge in using triples with this language is in preserving the abstraction of a functional pipeline owed to JSONata-like language features.

<sup>6</sup>I suppose we'll only support strings, though namespaced symbols would be nice!

3. The third element of the triple provides an atomic datum (string, number, entity identifier, etc.) associated with the subject attribute of the subject object.

For example, the first object in Figure 4, `{"ShipDate": "2021-06-15", "Item": "Widget 123", "Qty": 1.0, "UnitPrice": 10.5}` can be encoded as the following four triples, where the integer 11 is the unique entity identifier for this object.

```
[11, "ShipDate", "2021-06-15"]
[11, "Item", "Widget 123"]
[11, "Qty", 1.0,]
[11, "UnitPrice", 10.5]
```

Of course, an object can have multiple values for an attribute. For example, a person can multiple phone numbers. The way you'd typically handle cardinality greater than 1 in JSON and similar schemes is to simply provide the attribute with an array value. For example,

```
{"name" : "Bob", "phoneNumbers" : ["123-456-1111", "123-456-2222"]}
```

In these cases, there would be a datum (as these triples are sometimes called in Datalog-like databases) for each of phone numbers. Thus, the above might be normalized to triples (datoms) as follows:

```
[12, "name" "Bob"]
[12, "phoneNumbers" "123-456-1111"]
[12, "phoneNumbers" "123-456-2222"]
```

In the following, an array of objects is used for phone numbers. Consequently, the `PhoneNumberObjs` datoms on entity 13 reference additional entities, not data.

```
{"name" : "Bob", "phoneNumberObjs" : [{"cell" : "123-456-1111"}, {"work" : "123-456-2222"}]}
```

Figure 5: The object for Bob as triples

```
1 [13, "name" "Bob"]
2 [13, "phoneNumberObjs" 14]
3 [13, "phoneNumberObjs" 15]
4 [14, "cell" "123-456-1111"]
5 [15, "work" "123-456-2222"].
```

It is worth taking the time to comfortable with this idea of encoding information in triples; it is how we work with complex data in RADmapper.

### 2.5.2 Queries on triples

You can query triples using Datalog-like notation.<sup>7</sup> Queries are specified using the square-bracket-form triples illustrated in the data above. For example, you can imagine `[?e "name" "Bob"]` being part of a query to get the entity identifier for the `Bob` object (where `?e` is a *query variable*). But why would you want the entity identifier; it isn't domain data? The typical answer is that you will need it to perform a sort of join on triples. For example, if we want to get the `cell` number for `Bob`, we could use the following interrelation of triples.

```
[?e "name" "Bob"]
[?e "phoneNumberObjs" ?pn]
[?pn "cell" ?cellNum]
```

If the datoms<sup>8</sup> queried are just the ones of Figure 5, then the datum on Line 1 matches the query pattern `[?e "name" "Bob"]`, binding `?e` to the entity id 13. The second query pattern, `[?e "phoneNumberObjs" ?pn]`, matches two datoms, those on Lines 2 and 3, binding `?pn` to either entity 14 or 15. Owing to specifying `"cell"` as the attribute, however, the third datum pattern, `[?pn "cell"`

<sup>7</sup>There are many implementations of Datalog, includes ones for Java and JavaScript, so what is being proposed here should not be hard to implement.

<sup>8</sup>Triples used in this context are sometimes called *datoms*. I am glossing over some detail here.

`?cellNum]`, only matches the datom on Line 4. The result of this query can be thought of as an object that looks like this: `{?e: 13, ?pn: 14, ?cellNum: "123-456-1111"}`. Such an object is called a *binding object*, a collection of them that is the result of a query is a *binding set*.

There are a few things to note about working with Datalog-like languages like this one. First, except for performance concerns, it does not matter what order you list the triple forms. Second, you can put query variables in any (or all) of the three positions. For example, `[?e "cell" "123-456-1111"]` binds `?e` to 14 in the above. Third, it is possible that there is in the data more than one match; for example, `[?e ?a ?v]` returns an array of binding objects (a binding set) matching all data that is in context. It is possible, using query parameters not yet discussed, to limit what is returned to the first match, etc. Finally, you can imagine that using attribute names like "name" and "qty" could get confusing. (Name of what? Quantity of what?) For this reason, where possible, best practice in Datalog-like languages uses namespaced attributes. We could use, for example, strings with a slash between the namespace name and the property name. Thus perhaps, "Person/name" and "Phone/cell". Where it helps reduce mistakes and the cognitive load on programmers, getting data in that form can be done with JSONata-like pre-processing.

### 2.5.3 The query function

The `query` function is a higher-order function, it returns a function that can be used to execute a query. `query` takes two kinds of information: a collection of triple forms, and (optional) parameters. In the example below there are three triples and one parameter, `$name`.

```
$queryCellByName := query($name) ([?e "name" $name]
                                   [?e "phoneNumberObjs" ?pn]
                                   [?pn "cell" ?cellNum])
```

If the context data is such as in Figure 5, then `$. $queryCellByName("Bob")` would return an array containing one binding set `[{?e: 13, ?pn: 14, ?cellNum: "123-456-1111"}]`.<sup>9</sup> There is a similar construct `queryOne` that can be used to get a binding object for the first match found.

Primarily for use in development and debugging, there is a shortcut to getting binding sets: the function `query!`. `query!` takes a collection of triple forms (and no parameters) and returns a binding object for each match to the query (a binding set). For example:

```
$.query!([?e "name" "Bob"]
         [?e "phoneNumberObjs" ?pn]
         [?pn "cell" ?cellNum])
```

returns `[{?e: 13, ?pn: 14, ?cellNum: "123-456-1111"}]` when `$` is the data of Figure 5. `query!` is very useful towards getting comfortable with triples.

## 2.6 Mapping Networked Data

Binding sets, in themselves, are not too useful; you might be wondering why we even discuss them. The answer is that they are crucial to mapping networked data, and doing in-place updates in RADmapper.

**networked data** a collection of data that contains pointers to other parts of the same data collection.

For example, we could have a collection of data that includes information about Bob. Instead of repeating everything we know about Bob each time he is referenced in the data, networked data can use references to that same data. Resolving the reference (which might be implemented as a UUID, for example) brings you to the information about Bob.

**in-place update** the idea that the mapping task involves updating an existing collection of data, rather than creating the target from scratch.

<sup>9</sup>That example uses the JSONata convention that `$` is the *context reference* always bound to whatever is the data at the time of reference.

However, before we can talk about mapping networked data, and in-place updating, it is important to recognize that these activities require some additional knowledge about the (target) data. For example, (a) how do we know that an action is intended to update some data rather than add new additional information about the object, and (b) how do we know whether to use a pointer to reference some data rather than just put the data there?

Thus to perform these more complex mapping tasks we have to know a few things: (1) the cardinality of each attribute, (2) the type of each attribute (or at least the distinction between references and data types), and (3) attributes that provide keys (that is, uniquely identify an object of a given type).

Notice that condition (3) mentions the idea of object type. It is not the case that objects need to have any inherent notion of type to use RADmapper mapping. It is enough that the programmer recognizes the type by the attributes it possesses. (This is the so called *duck typing* — if it walks like a duck and quacks like a duck it is a duck.) Thus, an object that has an email address and a cell phone might be recognized as a customer. That said, programmers are likely to have an easier time of it if they explicitly hint at the types of things using namespaced attributes. By *namespaced attributes* we mean the idea that the name of the attribute hints at the object type it appears in. Thus the attribute for a customers email might be the name-value pair "customer/email" : "bob@example.com".

The following example illustrates mapping of Resource Definition Framework (RDF) and Web Ontology Language (OWL) data to a “sort of” relational form. (We say “sort of” because the description elides discussion of some details that would need to be addressed to do a thorough job of mapping RDF and OWL data.)

The target data we’ll be creating has three kinds of things in it, schema, tables, and columns. See Figure 6. The object types define attributes (db/attrs) and keys (db/key). Each attribute has a type (db/type) and a cardinality (db/cardinality). The db/type can be any primitive type defined (e.g. strings, numbers, UUID, etc.) or “object”, meaning that the attribute is populated by either references to other objects or objects in-lined in the structure. In relational parlance the references are foreign keys. There is nothing quite like in-lining in relational technology, but we demonstrate its use in the example below nonetheless.

Figure 6: A schema for the target data

```

1 [{"schema" : {"db/attrs" : [{"schema/name" : {"db/type" : "string", "db/cardinality" : "one"}},
2   "db/key" : [{"schema/name"}]}],
3
4   {"table" : {"db/attrs" : [{"table/name" : {"db/type" : "string", "db/cardinality" : "one" },
5     "table/schema" : {"db/type" : "object", "db/cardinality" : "one",
6       "db/in-line?" : true},
7     "table/columns" : {"db/type" : "object", "db/cardinality" : "many"}]},
8     "db/key" : [{"table/schema", "table/name"}]}],
9
10  {"column" : {"db/attrs" : [{"column/name" : {"db/type" : "string", "db/cardinality" : "one"}},
11    {"column/type" : {"db/type" : "string", "db/cardinality" : "one"}},
12    {"column/table" : {"db/type" : "object", "db/cardinality" : "one"}},
13    "db/key" : [{"column/table", "column/name"}]}]}
```

The db/cardinality of all the attributes is “one” except in the case of table/columns which is “many”. (See Line 7 of the figure.) Cardinality of one means that actions on this attribute of an object either create the datum (triple) or update it. There will never be two triples about this attribute where the first two elements of the triple are identical. Cardinality of many means that there could be multiple datums for the attribute. Keep in mind that the elements of a triple are always atomic. To represent  $n$  values there are  $n$  triples having different values for the third element of the triple. Also, the underlying datalog tool preserves the order of elements in arrays.

Line 7 in the figure demonstrates the use of the optional, db/in-line? schema feature. On Line 7 db/in-line? is set to true (the default is false). db/in-line? true directs the mapper to create structures rather than references. In the example, schemas are described by just their name; entire objects are serialized like this: {schema/name ``some schema``}. It seems reasonable in these circumstances to

in-line, much less so in the case of the back-pointer `table/columns`.

The OWL and RDF data that is the source of mapping consists of objects of two types, one are OWL classes; these have the value `"owl/Class"` in their `"rdf/type"` attribute. The other are OWL properties these have either the value `"owl/ObjectProperty"` or `"owl/DatatypeProperty"` as their `"rdf/type"` attribute. An OWL class is depicted in Figure 7.

Figure 7: An object in the source population

```

1  { :resource/iri :dol/endurant,
2    :resource/name "endurant",
3    :resource/namespace "dol",
4    :owl/disjointWith [ :dol/abstract :dol/quality :dol/perdurant ],
5    :rdf/type :owl/Class,
6    :rdfs/comment
7    ["The main characteristic of endurants is that all of them are independent essential wholes..."],
8    :rdfs/subClassOf
9    [ :dol/spatio-temporal-particular
10     { :owl/onProperty :dol/participant-in, :owl/someValuesFrom [ :dol/perdurant ], :rdf/type :owl/Restriction }
11     { :owl/allValuesFrom [ :dol/endurant ], :owl/onProperty :dol/specific-constant-constituent,
12       :rdf/type :owl/Restriction }
13     { :owl/allValuesFrom [ :dol/endurant ], :owl/onProperty :dol/part, :rdf/type :owl/Restriction } ] ]

```

Figure 8 depicts the complete specification of a transformation of the source. `transform` is a side-effecting function that takes three arguments, a data context, a binding set, and an enforce function; it returns a connection to resulting data. (Returning the data itself might not be reasonable in the case that it is very large and managed by a database.) When the binding set argument is a literal call to `query!` it is possible for the parser to do syntax checking between the `query!` and `enforce`.

The `query!` produces a binding set for the source data. The `enforce` defines how values from the binding set are used in the target population.

Figure 8: Mapping OWL to the relational database schema

```

1  $.transform(
2    query!([?class rdf/type      "owl/Class"]
3           [?class resource/iri  ?class-iri]
4           [?class resource/namespace ?ns]
5           [?class resource/name  ?class-name]
6           [?rel  rdfs/domain     ?class-iri]
7           [?rel  rdf/type        ("owl/ObjectProperty" or "owl/DataProperty")]
8           [?rel  rdfs/range      ?rel-range]
9           [?rel  resource/name    ?rel-name])
10  enforce( { "table/name"       : ?class-name,
11             "table/schema"     : { "schema/name" : ?ns },
12             "table/columns"    : { "column/name"  : ?rel-name,
13                                   "column/type"   : ?rel-range,
14                                   "column/table"  : ?table-ent } } as ?table-ent))

```

We'll walk through the example line by line. On Line 1 `transform` receives data from the context variable. [Other methods will be defined.] On Line 2 we constrain the search to just those entities that have `rdf/type` containing the string `"owl/Class"`. On Lines 3, 4, and 5 we bind variables to attribute values of that entity. If instead of mapping every namespace in the source data we wanted to limit mapping to a single namespace, we could have specified a value or expression for `resource/namespace` on Line 4 rather than the variable `?ns`. For example, we could have specified `"dol"` here, that happens to refer to the DOLCE ontology namespace. (See the example data in Figure 7.)

Lines 6, 7, 8 and 9 match and bind a second entity. Line 6 uses the variable `?class-iri` which was bound on Line 3; thus a relational join is performed. According to Line 7 this entity must have `rdf/type` of either `"owl/ObjectProperty"` or `"owl/DataProperty"`. Lines 8 and 9 bind the `rdfs/range` and `resource/name` respectively.

The example dataset from DOLCE is fairly large. Let's suppose it contains 50 classes, each involved in 10 relations. That suggests a binding set of 500 binding objects. That does not necessarily entail



500 structures like Lines 10 through 14, however. In contrast to the JSONata-like operations, which maps physical structures to physical structures, the mapping engine here is mapping between logical structures. Specifically, a triple represents a fact and the database of triples need only represent a fact once. It is the knowledge of keys and references provided by the schema that allows the mapping engine to construct physical structure from the logical relationship defined by the mapping specification.

### 2.6.1 In-place updates

[**This part needs to be updated.**] Used alone in code, `query` does not fit the pattern of a JSONata-like capability. JSONata is effective and concise because it allow one to thread data through a pipeline of primitive transformations. If you place a `query` call anywhere in that pipeline, you lose all the data except what is specifically captured by `query`. There are certainly instances where that is useful, but is `query` any easier than pure JSONata in these circumstances?<sup>10</sup> Perhaps the unique value of `query` is as an argument to higher-level functions to do tasks such as in-place and bi-directional updating. Therefore, discussed next is (1) a new construct called `enforce` and (2) how `query` and `enforce` can be used as arguments to a higher-order function `transform` to do in-place updating.

## 2.7 Working with multiple sources

## 2.8 Working with knowledge graphs

# 3 Use cases

## 3.1 Interoperable mapping tables

## 3.2 Integration flows

## 3.3 Data catalogs and knowledge graphs

## 3.4 Data validation

## 3.5 Joint cognitive work / lo-code tools

# 4 Specification

[This section will include a feature-by-feature description of the language similar to a language reference manual. Right now I've just listed the functions we think are needed.]

## 4.1 Language Syntax

This section is apt to be rewritten several times as the features of the language emerge.

### 4.1.1 JSONata subset

[I haven't found a formal definition (e.g. BNF) for JSONata, but the following *seems to work*. Not yet complete.]

```
<content> ::= <code-block> | <exp>

<code-block> ::= ' (' (<exp-or-assignment> ';'*) ' ) '
```

---

<sup>10</sup>Stakeholders, do you have examples where you think `query` might be better? Nothing comes to me, at the moment. There is also the matter of getting back an object keyed by binding variables. Thus far we haven't defined ways to manipulate that.

```
<exp-or-assignment> ::= <exp> | <assignment>

<exp> ::= ( <filter-exp> | <binary-exp> | (<builtin-un-op> <exp>) | <paren-delimited-exp> |
          <square-delimited-exp> | <fn-call> | <literal> | <field> | <id> )
          <conditional-tail>?

<binary-exp> ::= <operand-exp> <binary-op> <exp>

<operand-exp> ::= <field> | <id> | literal | <fn-call> | <delimited> | <unary-op-exp>

<filter-exp> ::= <operand-exp> '[' <exp> ']'

<conditional-tail> ::= '?' <exp> ':' <exp>

<range-exp> ::= '[' <exp> '..' <exp> ']'

<fn-call> ::= <id> '(' (<exp> (',' <exp>)*)? ')'

<fn-def> ::= 'function' '(' (<id> (',' <id>)*)? ')' '{' <exp> '}'

<js-map> ::= "{" (<map-pair> (',' <map-pair>)*)? "}"

<map-pair> ::= <string> ":" <json-data>

<json-data> ::= STRING | NUMBER | <json-array>

<json-array> ::= '[' (<json-data> (',' <json-data>)*)? ']'

<literal> ::= STRING | NUMBER | 'true' | 'false' | <regex> | <js-map>

<regex> (conforms to JavaScript's implementation)

<id> (a string starting with a \$.... More on this later.)
```

## 4.2 Table of elementary language features

| Task                              | Example                                      |
|-----------------------------------|--|
| Access (or navigate)              | Address.PhoneNumber                          |
| Select from array (simple)        | [2]  |
| Select from array (from end)      | [-2]   |
| Select all from array             | (done implicitly)                            |
| Select a subset of array          | [[0..2]]                                     |
| Entire input document             | \$   |
| Simple query                      | Phone[type='mobile']                         |
| Select values of all fields       | Address.*                                    |
| Select value from any child       | *.Postcode                                   |
| Select arbitrarily deep           | **.Postcode                                  |
| Concat                            | 'a' & 'b'                                    |
| Usual numeric operators           | +, -, *, /                                   |
| Array constructors                |  |
| Object Constructors               | Phone{type: number}                          |
| Object Constructors (force array) | Phone{type: number[]}                        |
| Code Block                        | (exp1; exp2; exp3)                           |
| Value as key in result            | Account.Order.Product{'Product Name': Price} |
| Map                               | seq.exp                                      |
| Filter                            | seq[exp]                                     |
| Reduce (group and aggregate)      | seq{exp:exp, exp:exp...}                     |
| Sort                              | seq ^ (exp)                                  |
| Index                             | seq # \$var                                  |
| Join                              | seq @ \$var                                  |
| Functions (in-line)               | function(\$x) : {(...)};                     |
| Naming functions                  | <i>myfunc</i> = <i>function</i> (x)...       |
| Context variable                  | \$   |
| Root context variable             |  |
| Conditional expression            | pred ? exp : exp                             |
| Variable binding                  | \$my_var := "value"                          |
| Function chaining                 | value ~> f2 -> f3                            |
| Regular expressions (like JS)     | /ab123/                                      |
| Time                              | <i>now</i> (), <i>millis</i> ()              |
| Transform                         | head ~> — location — update [,delete] —      |
| Example variable binding          | [needs investigation]                        |
| Type coercion                     | asType(obj, typespec)                        |
| Type checking                     | isTypeOf(obj, typespec)                      |
| Coersion to integer               | toInteger(x)                                 |
| Coersion to real                  | toReal(x)                                    |
| Returning a substring             | substring()                                  |

### 4.3 Table of language features for use with collections

| Task                                | Example                                  |
|-------------------------------------|--|
| Check for inclusion                 | <code>\$includes(col, x)</code>          |
| Check for exclusion                 | <code>\$excludes(col, x)</code>          |
| Size of a collection                | <code>\$size(x)</code>                   |
| Count elements                      | <code>\$count(col, obj)</code>           |
| sum of numbers                      | <code>\$sum(x)</code>                    |
| Cartesian product                   | <code>\$product(x, y)</code>             |
| union of collections                | <code>\$union(x, y)</code>               |
| symmetric difference of sets        | <code>\$symmetricDifference(x, y)</code> |
| coerce to ordered set               | <code>\$asOrderedSet(col)</code>         |
| first element of ordered collection | <code>\$first(col)</code>                |
| last element of ordered collection  | <code>\$last(col)</code>                 |
| intersection of collections         | <code>\$intersection(x,y)</code>         |
| find one                            | <code>\$any(col, predicate)</code>       |
| check for presence (return T/F)     | <code>\$one(col, predicate)</code>       |
| filter collection                   | <code>\$filter(col, predicate)</code>    |
| sort collection                     | <code>\$sortedBy(predicate)</code>       |
| check every (return T/F)            | <code>\$forAll(predicate)</code>         |
| add value to collection             | <code>\$including(col, val)</code>       |
| flatten collection of collections   | <code>\$flatten(x)</code>                |

## References

- [1] Jsonata.org. JSONata: query and transformation language, 2021.
- [2] Object Management Group. MOF Query/View/Transformation. Technical report, Object Management Group, Needham, MA, 2016.
- [3] Serge Abiteboul, Richard Hull, Victor Vianu, and Amsterdam Bonn Sydney Singapore Tokyo Madrid San Juan Milan Paris. *Foundations of Databases*. Addison-Wesley, Reading, MA, 1995.
- [4] Open Integration Hub.org. The Open Integration Hub, 2021.
- [5] Node-Red.org. Node-Red: Low-code programming for event-drive applications, 2021.