

# OAGi Interoperable Mapping Specification

Peter Denno & OAGi Members

2023-07-22 Draft

## 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 agents) into mapping specification in other languages. For example, it should be possible to translate statements in the exchange form into mapping specifications used by commercial mapping tools.

RADmapper, the mapping language, integrates JSONata expression language [1] with Datalog [2] and large language model (LLM) capabilities. It currently supports mapping to/from JSON, XML, tables (e.g. Excel), and knowledge graphs (RDF). Use cases for the RADmapper language including in-place updating of target data sources, mapping from multiple sources, and LLM-based matching and extraction tasks.

This document describes the RADmapper language and its strategy for serving as an interoperable exchange form mapping specifications. The reference implementation of the RADmapper language can be found in the Github repository <https://github.com/pdenno/RADmapper>. An web-based “exerciser” tool to explore the language is available as a Docker image <https://hub.docker.com/r/pdenno/rm-exerciser>. The document is a draft and, as of this writing (2023-07-22), is likely to be updated often, as will the language implementation and exerciser.

## 2 Quick Start: Example Mapping Tasks

This section uses examples to describe the basic features of the RADmapper language. RADmapper provides the complete expression language and all the built-in functions of JSONata. Like JSONata (and Javascript) RADmapper is a functional language: functions can be passed to functions as values; the value returned can be a newly created function. Examples of JSONata can be found in the JSONata specification. This section only goes into detail about the capabilities of RADmapper not found in JSONata. The principal concepts to discuss are

- an end-to-end mapping task using Datalog, large language models, and the interoperable exchange form,
- the relational and graph forms of data,
- RADmapper’s `query` declaration, used to query relational data, and
- RADmapper’s `express` declaration, used to reorganize (“map”) data from its original form to the form in which it is needed.

### 2.1 An End-to-end Example

Without going into detailed discussion about how things work, this section describes a mapping task end-to-end. The example is typical of “mapping” problems: we have data encoded in one structure (the “source” form) and we’d like to express that data as another structure (the “target” form). The example uses many of the capabilities of the RADmapper language and server capabilities of the exerciser<sup>1</sup>, including

- fetching data,
- fetching stored functions,
- calling an LLM for matching and extraction,

---

<sup>1</sup>The exerciser can be viewed as a reference implementation for cases where you need executable to interact with a server.

- uses of datalog,
- general asynchronous execution, and
- human validation of AI-generated mapping functions.

The last bullet above, about validation, is in keeping with basic goals of RADmapper: we aim to produce results that expresses the essence of a mapping task in terms that aren't so procedural looking that business-oriented analysts can't read them.<sup>2</sup> Our guess is that, emerging AI capabilities notwithstanding, business-oriented validation is going to continue to be a requirement for the foreseeable future. Thus we aim for a result that is readable and easily serialized (as, say, JSON) so that it might be translated into code for use with a commercial tool.

Let's suppose that you want to interact with a business partner using data the partner can easily provide. You'd like your interface to accept the partner's form and map it to a form that your system can process directly. Suppose the business partner's data are structured as suggested by the example below, and that you'd like that data in a form shown beneath it.

Figure 1: Example source and target data

```

1 // Example data you receive from the partner:
2
3 {'Invoice':
4   {'ApplicationArea':
5     {'CreationDateTime': '2023-07-10'},
6     'DataArea':
7       {'Invoice':
8         {'InvoiceHeader': {'PurchaseOrderReference': {'ID': 'PO-1234'}},
9         'InvoiceLine': {'BuyerParty':
10                        {'Location':
11                          {'Address':
12                            {'AddressLine':
13                              '123 Mockingbird Lane, Gaithersburg MD, 20878'}},
14                            'TaxIDSet': {'ID': 'tax-id-999'}},
15                            'Item': {'ManufacturingParty': {'Name': 'Acme Widget'}}}},
16                            'Process': 'Text description here, maybe.'}}}}
17
18 // The form your system will accept:
19
20 {'Invoice':
21   {'DataArea':
22     {'ApplicationArea':
23       {'CreationDateTime': '2023-07-10'},
24       'Invoice':
25         {'InvoiceLine':
26           {'BuyerParty':
27             {'Location':
28               {'Address': {'BuildingNumber': '123',
29                           'CityName': 'Gaithersburg',
30                           'PostalCode': '20878',
31                           'StreetName': 'Mockingbird Lane'}},
32               'TaxIDSet': {'ID': 'tax-id-999'}},
33               'Item': {'ManufacturingParty': {'Name': 'Acme Widget'}},
34               'PurchaseOrderReference': {'ID': 'PO-1234'}},
35               'Process': 'Text description here, maybe.'}}}}

```

Two differences between the customer's form and yours are that:

1. the customer's 'AddressLine' (Lines 12 and 12) is decomposed into its constituent details, BuildingNumber, City-Name, (Lines 28–31) etc. in your target form and,
2. the customer's structure places the purchase order reference in a structure called InvoiceHeader (Line 8), whereas in your form a purchase order is associated with each InvoiceLine, as shown by the nesting of Line 34 in your target structure.

To get from these requirements to working code we will (1) provide the abstract "shape" of the two structures above to an LLM function, `$llmMatch`, to reconcile differences in the structures, (2) have `$llmMatch` return a "mapping function" in RADmapper language that we can study for correctness to our business requirements, and (3) document and store the verified function for use whenever it is needed.

<sup>2</sup>Or if we fall a little short of that goal, at least provide data for GUI tools that could illustrate what relationships are being asserted.

We can demonstrate use of the stored function with the exerciser web app, in which case we'd use the RADmapper function `$get` to get both data and the function. In your tooling, you might access the data and function much differently. In the exerciser, it looks like the following.<sup>3</sup>

```
1 (
2   $data := $get(['library_fn', 'bie-1-data'], ['fn_exe']).fn_exe;
3   $mappingFn := $get(['library_fn', 'invoice-match-1->2-fn'], ['fn_exe']).fn_exe;
4   $mappingFn($data)
5 )
```

Line 2 of this example calls `$get` to get the example data. `$get` works something like GraphQL, in this case just specifying one property of the argument object, `['library_fn', 'bie-1-data']`, to retrieve. Running from the exerciser, `$get` is an async call to the server serving the app. Line 3 similarly gets the mapping function. Note that both lines (and indeed the whole example and most RADmapper syntax) conforms to JSONata syntax. `$get` returns an object with the attributes listed in its second argument. Both Line 2 and Line 3, use JSONata syntax `.fn_exe` to get the one property retrieved from the object, its executable. Of course, "executable" means different things in different contexts; if the code is retrieve under Javascript, like with the web exerciser, the code is JS-executable; if `$get` is called from the server, Java-executable code is provided.<sup>4</sup> Line 4 applies the `$mappingFn` to the `$data` resulting in an object like shown in Lines of 20–35 of Figure 1. None of the above explains what `$mappingFn` does nor how it was created. We will now work backwards from this ending point (the data mapped) to see how the mapping code is generated.

In the exerciser, if you execute the expression  
`$get(['library_fn', 'invoice-match-1->2-fn'], ['fn_src']).fn_src`  
 (rather than `fn_exe` as used above) returned would be the following string value:

```
1  '// Here we assume that we validated the $llmMatch result, and stored it as this.
2
3  function($d){
4    {'Invoice':
5      {'DataArea':
6        {'ApplicationArea':
7          {'CreationDateTime': $d.Invoice.ApplicationArea.CreationDateTime},
8          'Invoice':
9            {'InvoiceLine':
10              {'BuyerParty':
11                {'Location':
12                  {'Address':
13                    {'BuildingNumber':
14                      $llmExtract(
15                        $d.Invoice.DataArea.Invoice.InvoiceLine.BuyerParty.Location.Address.AddressLine,
16                        'BuildingNumber')},
17                    'CityName':
18                      $llmExtract(
19                        $d.Invoice.DataArea.Invoice.InvoiceLine.BuyerParty.Location.Address.AddressLine,
20                        'CityName')},
21                    'PostalCode':
22                      $llmExtract(
23                        $d.Invoice.DataArea.Invoice.InvoiceLine.BuyerParty.Location.Address.AddressLine,
24                        'PostalCode')},
25                    'StreetName':
26                      $llmExtract(
27                        $d.Invoice.DataArea.Invoice.InvoiceLine.BuyerParty.Location.Address.AddressLine,
28                        'StreetName')}}},
29                    'TaxIDSet': {'ID': $d.Invoice.DataArea.Invoice.InvoiceLine.BuyerParty.TaxIDSet.ID}},
30                    'Item': {'ManufacturingParty':
31                      {'Name': $d.Invoice.DataArea.Invoice.InvoiceLine.Item.ManufacturingParty.Name}},
32                    'PurchaseOrderReference':
33                      {'ID': $d.Invoice.DataArea.Invoice.InvoiceHeader.PurchaseOrderReference.ID}},
34                    'Process': $d.Invoice.DataArea.Process}}}}
35  }'
```

This string is the value returned from a call to `$llmMatch` with the "shapes" of the source and target data. You could get the corresponding executable by either `$get`-ing the `fn_exe` or calling the function `$eval` as you would in JSONata. Note the following about this string:

<sup>3</sup>This example is currently executable from the exerciser <https://hub.docker.com/r/podenno/rm-exerciser> as the default example, if you'd like to try it yourself. Note that the exerciser has an OpenAPI interface. In the exerciser docker image, it is documented at <http://localhost:3000>; the app itself is <http://localhost:3000/app>.

<sup>4</sup>You are probably wondering what "getting the executable" means with respect to data. These `'library_fn'` objects have three attributes, `fn_exe`, `fn_src`, and `fn_doc`. `fn_src` is a string. In the case of data, `fn_exe` returns the structure represented by that string, just as you might do with a string representing JSON. Note also that there is a `$put` function in the exerciser, so that you can play with these kinds of tasks.

- The string defines a function of one argument, `$d` which should be bound to the source data, the data in Lines 3–16 of Figure 1.
- Lines 14, 18, 22, and 26 call the function `$llmExtract`, which is a LLM call to extract the information describe by the second argument string from the first argument string.
- The long expressions beginning with `$d. . .` are navigations into the source data, and are associated with target data properties.

## 2.2 Data Organized as Triples and the query Construct

`query` is the principal construct of RADmapper providing Datalog-like functionality to the language. `query` declarations are used like JSONata or Javascript function declaration in the sense that the value of the declaration (a function) can be assigned to variables and used directly. For example,

```
$addOne := function(x){x + 1}
```

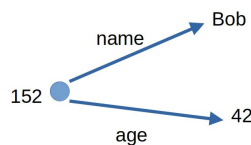
defines a function and assigns it to the variable `$addOne`. `$addOne(3)` is a call to the function with the argument 3. Similarly,

```
$myQuery := query(){[$DB1 ?person :age 42]}
```

defines a query that can be used like a function. `$myQuery($myDB)` is a call to the function with whatever “database” is assigned to `$myDB`. However, unlike ordinary functions, the body of a `query` consists of one or more *Datalog patterns* such as the pattern `[$DB1 ?person :age 42]` shown. We’ll start by talking about the database argument to the query, for example, the value assigned to `$MyDB` in the expression `$myQuery($MyDB)` then get back to talking about the patterns.

Relational (table-based) and graph-based data can be described by triples  $[x, rel, y]$  where  $x$  is an entity reference,  $y$  is data (string, number, entity reference, etc.) and  $rel$  is a relationship (predicate) holding between  $x$  and  $y$ . For example, in syntax similar to JSON we could describe the fact that Bob’s age is 42 with an object `{‘name’ : ‘Bob’, ‘age’ : 42}`. As triples, we represent Bob being 42 years old with two triples, for example, `[152 ‘name’ ‘Bob’]` and `[152 ‘age’ 42]`. As a graph, this would look like the following:

Figure 2: Information about Bob in graph form



You might wonder where the 152 in the triples came from, or for that matter, why the fact that Bob being age 42 couldn’t simply be represented by the single triple `[‘Bob’ ‘age’ 42]`. The answer is that you could represent this fact with that one triple, but in doing so you are using ‘Bob’ as a key which might not be that good if your database has more than one Bob in it. So instead, to prepare for the more general case, Datalog-like graph databases use integers to refer to entities. 152 here is like a primary key in relational DB, or an IRI for a particular entity defined in RDF. A complete RADmapper program for querying the Bob database for people age 42 is as follows:

Figure 3: A complete query

```

1 ( $myDB := [{‘name’ : ‘Bob’, ‘age’ : 42}];
2   $ageQuery := query(){[?e :age 42]};
3   $result := $ageQuery($myDB); )

```

The above needs some explanation. On Line 1 we put a JSON-like object in an Javascript-like array by wrapping it in square brackets; this is the literal form of a very small database. On Line 2 we defined the query function; here note the use of a *query variable* `?e`. Also note that whereas we used the string ‘age’ for the relation, we used `:age` — syntax we call an *attribute* — to represent that same relation in the pattern. On Line 3 we apply the query function `$ageQuery` to the database `$myDB` defined on Line 1. Of course, `$myDB` isn’t much of a database; it contains only

the data shown in Figure 2. Alternatively, you could have defined data through other means, such as reference to a file or a GraphQL query. The result of this query, assigned to `$result` is a *binding set*, a set of objects each describing a consistent binding of the search pattern's variables to values from the database. In this example, with the data in Figure 2, the binding set consists of just one binding element and it binds one variable; the binding set is `[{?e : 152}]`. This indicates that the entity indexed at 152 has an attribute `:age` with value 42. If there were more entities in the database possessing an `:age` attribute, the query would have returned one such `{?e : <whatever>}` binding object for each of them.

Amazing as that query might seem, running against a database with one entity in it and all, it didn't tell us *what person* is age 42. All we got was a binding set for every entity that has an age attribute equal to 42. Entity IDs are internal to the database and not of much value to RADmapper users.<sup>5</sup> In order to make any use of this information, we need to join the `[?e :age 42]` pattern with a pattern that grabs the name attribute in the data, `[?e :name ?name]`. This is shown in Figure 4 below.

Figure 4: A more useful query, getting the name of the 42-year old person

```
1 ( $myDB := [{ 'name' : 'Bob', 'age' : 42 }];
2   $ageQuery := query() { [?e :age 42]
3                 [?e :name ?name] }
4   $result := $ageQuery($myDB); )
```

The pattern on Line 3, by virtue of its reuse of `?e`, imposes an additional constraint on the graph match: the entity bound to `?e` must have a `:name` attribute. The value of the name attribute is bound to `?name`. Thus each element in the binding set will bind two variable, `?e` and `?name`. The binding set for the database in the graph depicted in Figure 2 is `[{?e : 152, ?name 'Bob'}]`.<sup>6</sup>

In this example, we used one variable to match on the entity and another to capture a value, however each position (entity, attribute, and value) can take a variable. Further, you can use variables in more than one of those positions. For example, `query() { [?entity ?attr ?val] }` is a query that matches on every entity, attribute, and value of the database; it represents every edge of the database's graph. `query() { [ _ ?attr _ ] }` would return the names of all the attributes in the database (without duplicates). This provides a kind of introspection that is typically more difficult to obtain in other database technology. When matching the value position you are doing a relational join, for example, we might match the social security number (SSN) in some data against a same-valued (but possibly differently named) value in other data. For example,

```
$relJoinQuery := query() { [ $DB1 ?e1 :ssn ?id ]
                          [ $DB2 ?e2 :id ?id ] }
```

You may have noticed that the query above, and one used earlier have four elements in their pattern whereas most of the examples have only three (entity, attribute, and value). If four elements are provided, the first is the database to which the pattern is applied. If there is only one database being queried, as we've been doing earlier, you don't need to use four-place patterns. Let's look at a complete `query` example that uses two databases.

Figure 5: A query that looks into two databases

```
1 ( $DBa := [{ 'email' : 'bob@example.com', 'aAttr' : 'Bob-A-data', 'name' : 'Bob' },
2           { 'email' : 'alice@alice.org', 'aAttr' : 'Alice-A-data', 'name' : 'Alice' }];
3   $DBb := [{ 'id' : 'bob@example.com', 'bAttr' : 'Bob-B-data' },
4           { 'id' : 'alice@alice.org', 'bAttr' : 'Alice-B-data' }];
5
6   $qFn := query() { [ $DBa ?e1 :email ?id ]
7                     [ $DBb ?e2 :id ?id ]
8                     [ $DBa ?e1 :name ?name ]
9                     [ $DBa ?e1 :aAttr ?aData ]
10                    [ $DBb ?e2 :bAttr ?bData ] };
11
12   $bSet := $qFn($DBa, $DBb); )
```

<sup>5</sup>In fact, I sort of told a fib for ease of exposition; by default the binding of variables in the entity position, entity IDs, aren't provided in a binding object. Using default settings, what this example returns is `[{}]` meaning "one match was found." The binding sets depicted in most examples won't include bindings for entity IDs.

<sup>6</sup>If you read the previous footnote you know it actually returns just `[{?name 'Bob'}]`.

In Figure 5 Lines 1–4 we define two small databases and assign them to variables `$DBa` and `$DBb` respectively. On Lines 6–10 we define the query. Since both databases use email addresses for customer identification, we can use the `:email` and `:id` attributes to join together information about a customer from the two databases. That is the purpose of the patterns on Lines 6 and 7; the two patterns use different variables for the entities, `?e1` and `?e2`, because the information is coming from different databases and we don't control entity IDs, but both use `?id` to force matches on email address. The remainder of the patterns in the query, Lines 8–10, pick up various information from the two databases. Line 12 calls the query function bound to `$bSet` to get the binding sets against the two databases. Unlike our previous calls to the query function, this one takes two databases as arguments. The order of the arguments in the call must be the same as the order in which the databases appear in the query statement; `$DBa` appears first on Line 6, `$DBb` appears first on Line 7, so `$DBa` is the first argument to the call.

The binding set that is produced, the value of `$bSet`, consists of two binding objects:

```
[{?id : "bob@example.com", ?name : "Bob", ?aData : "Bob-A-data", ?bData : "Bob-B-data" },
 {?id : "alice@alice.org", ?name : "Alice", ?aData : "Alice-A-data", ?bData : "Alice-B-data"}]
```

One small but very significant point before moving on to discuss `express`: writing queries like `query() {[?e :age 42] [?e :name ?name]}` could become rather tedious in the case that you might want to get data about some other age value. For this reason, the `query` declaration can serve to produce a *higher-order function*, a function that returns (query) functions as values. Figure 6 demonstrates the idea.

Figure 6: Get the names of people ages 42 and 33.

```
1 ( $myDB := [{?name : 'Bob', ?age : 42},
2           {?name : 'Alice', ?age : 33}]
3
4   $ageQueryT := query($age) {[?e :age $age]
5                             [?e :age ?age]
6                             [?e :name ?name]}
7
8   $ageQ42    := $ageQueryT(42);
9   $ageQ33    := $ageQueryT(33);
10
11  $append($ageQ42($myDB) , $ageQ33($myDB)) )
```

The key difference is that on Line 4, the `query` construct, defines a parameter, `$age`, using ordinary JSONata-like syntax. You can define as many parameters as you'd like and they can be used to substitute into any of the pattern positions, entity, attribute, or value. By convention, if we are to assign a higher-order `query` function to a variable, we end the variable name with a `T` such as shown on Line 4, `$ageQueryT`. The `T` suggests that the variable denotes a *query template*. Lines 8 and 9 define query functions for querying ages 42 and 33, respectively. Line 11 uses the JSONata-like builtin `$append` to combine the two binding sets. Note that the pattern `[?e :age ?age]` on Line 5 is used get `?age` into the binding sets. The result of running this example is the binding set

```
[{?name : 'Bob', ?age : 42},
 {?name : 'Alice', ?age : 33}].
```

## 2.3 Constructing Target Data with `express`

Binding sets produced by `query` provide ordinary JSONata-like objects<sup>7</sup> that could be used with JSONata built-in functions and operators to produce target structures. However, the `express` construct provides capabilities beyond those of the JSONata expression language. This section describes some of those features.

Let's continue with our two-database example. Lines 1–12 of Figure 7 below are as they were in Figure 5. The binding set assigned to `$bSet` is

```
[{?id : "bob@example.com", ?name : "Bob", ?aData : "Bob-A-data", ?bData : "Bob-B-data" },
 {?id : "alice@alice.org", ?name : "Alice", ?aData : "Alice-A-data", ?bData : "Alice-B-data"}].
```

Assuming that we'd just like to present the target data as nested objects indexed by the customer's email address, for example:

<sup>7</sup>Two differences: (1) the keys of binding sets are query variables, and (2) binding sets are flat objects; the values at the keys are not objects themselves though they may be entity IDs (integers). These differences notwithstanding, you can use binding sets as though they are ordinary JSONata-like objects.

```

{"alice@alice.org" { "name" : "Alice",
                    "aData" : "Alice-A-data",
                    "bData" : "Alice-B-data" },

"bob@example.com" { "name" : "Bob",
                    "aData" : "Bob-A-data",
                    "bData" : "Bob-B-data" }}

```

we could use the `express` declaration that begins on Line 14 of Figure 7 to do this.

Figure 7: Using `express` with a query that looks into two databases

```

1 ( $DBa := [{ 'email' : 'bob@example.com', 'aAttr' : 'Bob-A-data', 'name' : 'Bob' },
2   { 'email' : 'alice@alice.org', 'aAttr' : 'Alice-A-data', 'name' : 'Alice' }];
3   $DBb := [{ 'id' : 'bob@example.com', 'bAttr' : 'Bob-B-data' },
4     { 'id' : 'alice@alice.org', 'bAttr' : 'Alice-B-data' }];
5
6   $qFn := query() { [$DBa ?e1 :email ?id]
7     [$DBb ?e2 :id ?id]
8     [$DBa ?e1 :name ?name]
9     [$DBa ?e1 :aAttr ?aData]
10    [$DBb ?e2 :bAttr ?bData] };
11
12   $bSet := $qFn($DBa, $DBb);
13
14   $eFn := express() { { ?id : { 'name' : ?name,
15     'aData' : ?aData,
16     'bData' : ?bData } } };
17
18   $reduce($bSet, $eFn )

```

`express` is a function-defining construct that, like `query`, is capable of returning `express` functions when supplied with parameters. In that sense it is capable of being a higher-order function. But here on Line 14 we aren't using parameter; the declaration is simply `express() { ... }`; no parameters imply it isn't a template and can be used directly. The body of the `express` declaration, the text inside the outer curly brackets, defines the pattern of JSONata-like object structure. The target data is produced by iterating over the `express` function bound to `$eFn` using the elements of the binding set. Two of the most common ways to iterate over a function in functional languages such as JSONata and RADmapper are `map` and `reduce`. `map` and `reduce` differ in how they process the collection of arguments: `map` applies a given function to each element independently; `reduce` "summarizes" results by allowing each element to affect a summary outcome.<sup>8</sup> Whether you `$map` or `$reduce` the `express` function over the binding set can have great influence on the outcome, as will be demonstrated in subsequent examples. As written above, the call to `$reduce` on Line 18 creates a nested structure for the first binding element and then inserts a second structure into it at the second `?id` key, as shown above. Were Line 18 replaced with `$map($bSet, $eFn)`, the result would be two independent nested maps, one for each binding set in the call to `$map`:

```

[{"bob@example.com" : { "name" : "Bob", "aData" : "Bob-A-data", "bData" : "Bob-B-data" }},
 {"alice@alice.org" : { "name" : "Alice", "aData" : "Alice-A-data", "bData" : "Alice-B-data" }}].

```

## 2.4 Predicate patterns and an example from a different perspective

This example illustrates (1) simple restructuring of a data structure, (2) use of predicates as query patterns, (3) extensive joining to navigate deeply nested structures, and (4) use of query variable in the attribute position of patterns. The example is based on a discussion on the JSONata Slack channel. The goal of the example is to swap the nesting of `owners` and `systems` as shown in Figure 8.

<sup>8</sup>Examples to help you recall how these work:

`$map([1,2,3,4], function($x){$x * 2})` returns [2,4,6,8]; it multiplies each argument by 2.

`$reduce([1,2,3,4], function($x, $y){$x + $y})` returns 10; it applies + to 1 and 2, then to 3 and that sum, then to 4 and that sum.

Figure 8: The goal is to swap the nesting of 'owners' and 'systems' in the data on Lines 1–9 so that it looks Lines 13–21.

```

1 { "systems":
2   { "system1": { "owners": { "owner1": { "device1": { "id": 100, "status": "Ok" },
3                                     "device2": { "id": 200, "status": "Ok" }},
4                                     "owner2": { "device3": { "id": 300, "status": "Ok" },
5                                               "device4": { "id": 400, "status": "Ok" }}}},
6   "system2": { "owners": { "owner1": { "device5": { "id": 500, "status": "Ok" },
7                                     "device6": { "id": 600, "status": "Ok" }},
8                                     "owner2": { "device7": { "id": 700, "status": "Ok" },
9                                               "device8": { "id": 800, "status": "Ok" }}}}}}
10
11 /* ...so that it looks like: */
12
13 { "owners":
14   { "owner1": { "systems": { "system1": { "device1": { "id": 100, "status": "Ok" },
15                                     "device2": { "id": 200, "status": "Ok" }},
16                                     "system2": { "device5": { "id": 500, "status": "Ok" },
17                                               "device6": { "id": 600, "status": "Ok" }}}},
18   "owner2": { "systems": { "system1": { "device3": { "id": 300, "status": "Ok" },
19                                     "device4": { "id": 400, "status": "Ok" }},
20                                     "system2": { "device7": { "id": 700, "status": "Ok" },
21                                               "device8": { "id": 800, "status": "Ok" }}}}}}

```

Typical of a restructuring task, the goal data on Lines 13–21 does not contradict any of the facts evident in the original data on Lines 1–9. For example, there is an owner called `owner1`, and `owner1` still has the same devices associated. The idea that mapping is about how facts are viewed is a key concept in RADmapper. What has changed in restructuring the example data is not a fact but the kind of forward navigation that is possible in the two structures. In the original structure, it is possible to navigate forward from a system to an owner (such as on Line 1, from `system1` to `owner1`). In the restructured data, it is possible only to navigate forward from owner to system (such as `owner1` to `system1` on Line 14). The restructuring task can be achieved using only JSONata functions and operators as depicted in Figure 9.

Figure 9: Using JSONata to restructure the data.

```

1 {
2   "owners": $distinct(systems.*.owners.$each(function($d, $ownerName) {$ownerName}))@$o.{
3     $o: $each($$.systems, function($sys, $sysName) {{$sysName: $lookup($$.systems, $sysName).owners ~> $lookup(
4       $o)
5     }} ~> $merge()
6   } ~> $merge()
7 }

```

Though the RADmapper solution to this problems involves more lines of code (see Figure 10 below), it is arguably easier to understand. The RADmapper solution is also arguably a better candidate for *interoperable* exchange of mappings, as will be discussed later.



Figure 10: RADmapper query and express used to restructure data. Note that on Line 15 we use `express` in-lined rather than assigning that function to a variable and passing the variable into `$reduce`. The choice to in-line has no effect on the result.

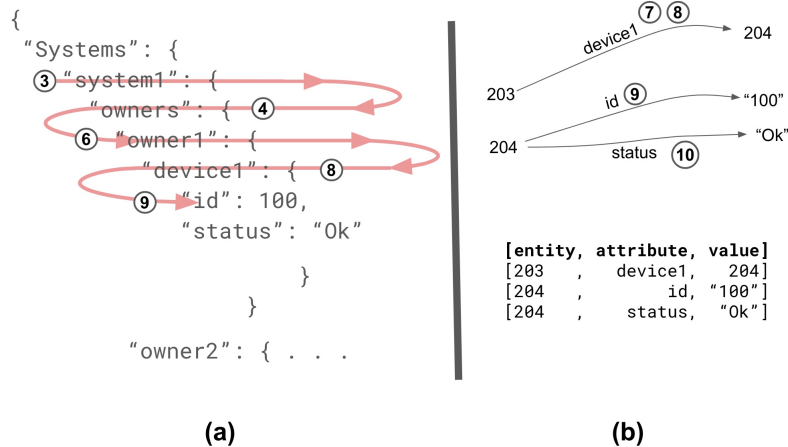
```

1  ( $data := $read('data/testing/jsonata/sTPDRs--6.json');
2    $q := query() { [?s ?systemName ?x]
3      [($match(?systemName, /system\d/))]
4      [?x :owners ?y]
5      [?y ?ownerName ?z]
6      [($match(?ownerName, /owner\d/))]
7      [?z ?deviceName ?d]
8      [($match(?deviceName, /device\d/))]
9      [?d :id ?id]
10     [?d :status ?status] };
11
12  $bsets := $q($data);
13
14  $reduce($bsets,
15    express() { { 'owners':
16      {?ownerName:
17        { 'systems':
18          {?systemName:
19            {?deviceName : { 'id' : ?id,
20                          'status' : ?status}}}}}
21      }
22    }
23  )

```

Where this example differs from earlier ones is the use of `$match` on Lines 3, 6 and 8. Instead of the usual 3- or 4-place pattern, we have a single *predicate* that is being applied using query variables from other patterns in the query. `$match` is a JSONata built-in Boolean function that returns true if the first argument, a string, matches the second argument, a regular expression. This example is also the first to use extensive joins to navigate into deeply nested objects. A “hairpin” pattern of joined variables is depicted in Figure 11 and reflected in Lines 2–10 of Figure 10.

Figure 11: Part (a): navigation of the source structure (Lines 1–9 of the data in Figure 8). Part (b): graph and triples representing device1. The circled numbers refer to lines in the code of Figure 10.



A Datalog database can be viewed as several index tables about the triples (such as those depicted in Figure 11 (b).) One such index table will index primarily by entity, another primarily by attribute, etc. Some of Lines 2–10 in Figure 10, for example, Line 2, `[?s ?systemName ?x]`, have the form of these triples, where respectively `?s`, `?systemName` and `?x` match an entity, attribute, and value. Line 4, `[?x :owners ?y]`, is similar but the attribute position is occupied by `:owners`. This triple will match any fact that has the string "owners" in its attribute position. There are two such facts in the data depicted in Figure 8, one on Line 2, and one on Line 6. In both cases the value position is occupied by another entity reference. Figure 11(b) similarly depicts a triple whose value position references another entity, `[203, device1, 204]`.

Lines 3, 6 and 8 of the query do not use the 3-place triple pattern; they consist of a single JSONata expression wrapped in parentheses. The expression should be a Boolean (should return true or false). Note that the expression uses

variable bindings (e.g. `?systemName`, `?ownerName`, and `?deviceName`) used in triples of the query. `/system\d/` is a JavaScript regular expression matching a string containing "system" followed by a single digit (the `\d` part). The following paragraph provides a line-by-line summary of how the query in Lines 2–10 works.

**Line 2** [`?s ?systemName ?x`]: All positions of this triple pattern are occupied by variables, so by itself it would match every triple in the database.

**Line 3** [`($match(?systemName, /system\d/))`]: This uses the query variable `?systemName` bound to the entity position on Line 2. Therefore, between this pattern and the one in Line 2, the only triples matching both of these patterns have "system1" or "system2" in the attribute position. (See the data to verify this.) `?s` and `?x` from Line 2 are thus bound to the entity and value positions of triples having either "system1" or "system2" in their attribute positions.

**Line 4** [`?x :owners ?y`]: Here we see `?x`, which was introduced in Line 2, reused in the entity position. We know that `?x` is bound to an entity by looking at the data. (See Line 2 in the data for example, ```owners``: {...`. The `:` `{` means the thing keyed by "owners" here is a JSON object (an "entity" in the database). This triple pattern ensures that `?x` refers to entities for which the "owners" occupies the attribute position.

**Line 5** [`?y ?ownerName ?z`]: Like Line 2, this is used to introduce a new variable, `?ownerName` in this case, which will be used in the `$match` predicate similar to Line 3. One difference here, however, is that the value of `?y` is constrained by the triple pattern on Line 4.

**Line 6** [`($match(?ownerName, /owner\d/))`]: This serves a similar purpose Line 3, but for owner keys.

**Line 7** [`?z ?deviceName ?d`]: This is like Lines 1 and 5, introducing a new variable for use in the `$match`. Like Line 5, the value of the variable in entity position is constrained by another pattern.

**Line 8** [`($match(?deviceName, /device\d/))`]: This is similar in purpose to the other two uses of `$match`.

**Line 9** [`?d :id ?id`]: Here we are finally ready to pick up some user data, the value of a device's "id" attribute.

**Line 10** [`?d :status ?id`]: Like Line 9, but to pick up the value of a device's "status" attribute.

It was suggested earlier that the encoding of the source data was a bit unusual. For example, object key values such as `owner1`, `system1` and `device1` are being used, apparently, both to indicate a unique entity and to identify the kind of entity being described, presumably an owner, system and device respectively. We might seek target data that explicitly declares the type and identifier separated, for example using attributes `type` and `id` respectively. The target data sought, for example, might be as depicted in Figure 12.

Figure 12: An alternative organization of the target data from the example.

```

1 [{"type" : "OWNER",
2   "id" : "owner1",
3   "systems": [{"type" : "SYSTEM",
4                 "id" : "system1",
5                 "devices": [{"type" : "DEVICE",
6                               "id" : "100",
7                               "status" : "Ok"},
8                               {"type" : "DEVICE",
9                                "id" : "200",
10                               "status" : "Ok"}]}],
11  {"type" : "SYSTEM",
12   "id" : "system2",
13   "devices": [{"type" : "DEVICE",
14                 "id" : "500",
15                 "status" : "Ok"},
16                 {"type" : "DEVICE",
17                  "id" : "600",
18                  "status" : "Ok"}]}]}],
19 [{"type" : "OWNER",
20   "id" : "owner2",
21   "systems": [{"type" : "SYSTEM",
22                 "id" : "system1",
23                 "devices": [{"type" : "DEVICE",
24                               "id" : "300",
25                               "status" : "Ok"},
26                               {"type" : "DEVICE",
27                                "id" : "400",
28                                "status" : "Ok"}]}],
29  {"type" : "SYSTEM",
30   "id" : "system2",
31   "devices": [{"type" : "DEVICE",
32                 "id" : "700",
33                 "status" : "Ok"},
34                 {"type" : "DEVICE",
35                  "id" : "800",
36                  "status" : "Ok"}]}]}]}

```

As the figure depicts, the output is more verbose. A guess at an `express` structure to produce this output is as follows.

Figure 13: Draft `express` structure for target data as depicted in Figure 12

```

1 $ex1 := express{{'owners' : {'type' : 'OWNER',
2                               'id' : ?ownerName,
3                               'systems': [{'type' : 'SYSTEM',
4                                              'id' : ?systemName,
5                                              'devices': [{'type' : 'DEVICE',
6                                                            'id' : ?deviceName,
7                                                            'status': ?status}]}]}]}
8 }

```

If we were to map over that `express` body, that is `$map($bsets, $ex1)`, the result would be (logically correct, of course but) even more verbose and also repetitive because there will be one full structure like the `express` body for each of the 8 binding sets. A better alternative might be to `$reduce` over the `express`, that is, `$reduce($bsets, $ex1)`; `$reduce` can have the effect of “summarizing” data, in this case by squeezing out the repetition. However, there is a problem with reducing over the `express` body as written; it would produce the following:

Figure 14: Reducing over the `express` body here results in iteratively overwriting data; the result is the same as applying `express` body to just the last binding set.

```

1 { "type"      : "OWNER",
2   "id"       : "owner1",
3   "systems"  : { "type"    : "SYSTEM",
4                 "id"      : "system1",
5                 "devices" : { "type"    : "DEVICE",
6                             "id"      : "device2",
7                             "status"  : "Ok" }}}
8
9 /* This result indicates that the last binding set processed was the following */
10 {?ownerName : "owner1", ?systemName : "system1", ?deviceName : "device2", ?status : "Ok", ?id : 200}

```

What happened to the rest of the data? Why were data not lost similarly when we reduced over the `express` body on Lines 14–20 of Figure 10? The answer is that the `express` body of Figure 10 uses unique values for the object keys. The effect of reducing over the `express` body in that example is to “drop in” new pathways for each binding set. Specifically, if we consider the Lines 13–21 of Figure 8, the result of applying the `express` body on Lines 14–20 of Figure 10 to the key values in each corresponding binding set, it is apparent that:

- the Line 14 object was created by the keys "owners", "owners1", "system1", "device1",
- the Line 15 object was created by the keys "owners", "owners1", "system1", "device2",
- the Line 16 object was created by the keys "owners", "owners1", "system2", "device5", and
- et cetera, to Line 21, created by the keys "owners", "owners2", "system2", "device8".

There are three solutions to the problem just highlighted: (1) use unique keys to “drop in” new pathways like in the original example, (2) use `$map` instead of `$reduce`, which entails accepting the fact that the result will be verbose with lots of repetition, or (3) wrap keys in the `key` construct as described below.

### 2.4.1 The key construct of `express`

We’ll continue with the running example. In order to prevent the overwriting that was depicted in Figure 14 we will adapt the code from Figure 13 by adding the `key` construct as shown.

Figure 15: The `express` structure from Figure 13 revised to identify keys

```

1 $ex1 := express({ 'owners' : { 'type'    : 'OWNER',
2                               'id'     : key(?ownerName),
3                               'systems' : [{ 'type'    : 'SYSTEM',
4                                              'id'      : key(?systemName),
5                                              'devices' : [{ 'type'    : 'DEVICE',
6                                                            'id'      : key(?deviceName),
7                                                            'status'  : ?status} ] } ] }
8
9
10 }

```

Here we modified Lines 2, 4, and 6 of Figure 13, wrapping the binding variables `?ownerName`, `?systemName`, and `?deviceName` in the `key` construct. The `key` construct declares identity conditions for the corresponding objects. With knowledge of identity conditions, reduce processing on the `express` body can distinguish between inserting a new object (one where the key hasn’t yet been seen) and updating (or mistakenly overwriting) the existing object. Thus the entire source data can be pushed into a single object structured as depicted in the figure. Incidentally, note that Lines 3 and 5 use square brackets to signify that there are possibly multiple systems and devices respectively in the target form.

## 2.5 Summary so far

`query` and `express` are two principal constructs of the RADmapper language that together provide powerful, flexible, means to restructure data. They enable mapping from multiple sources, in-place updating, and effective methods of abstraction and composition such as templating and higher-order functions. It may look like a lot to learn, and in fact there is more to the language yet to be discussed. That notwithstanding, we believe that many tasks can be programmed

with RADmapper more easily than with pure JSONata. For example, consider the running example of restructuring we discussed above. The pure JSONata solution uses lots of syntax and juggling operations (see Figure 9) to do something that is conceptually very simple. The RADmapper solution, on the other hand, reflects those simple observations: with `query` you pull out threads of data that conceptually hang together; with `express` you either `$map` or `$reduce` the individual threads into a result structure. Separating the collection of source information from the expression of target information makes the task easier to perform. Further, even these two mostly-independent steps can be approached in smaller sub-steps. For example, one could start with a `query` that binds just one or two query variables along a path. When it is clear that that works (verified in the RADmapper exerciser, for example) one could progressively add more path steps until the complete, coherent thread of domain relationships is captured.

Likewise, you can create binding sets by hand and approach programming the `express` structure in small steps. You can evaluate `express` with a single binding set to see the result, combine binding sets (even from calls to different `query` definitions) and experiment with `$reduce-ing` on the `express`.

The motivation for RADmapper, however, is more encompassing than what discussion thus far might suggest. Though you can use RADmapper's web-based exerciser to define and validate mappings, and you can include RADmapper in Java and JavaScript programs<sup>9</sup> a principal goal of RADmapper is to support *interoperable* mapping. The next section discusses interoperability; you don't need to read it if your primary goal is to learn how to use RADmapper in Java or JavaScript software.

### 3 Interoperable Specification of Mapping Requirements

RADmapper seeks to describe mapping requirements in ways that are useful to the various stakeholder and tools used in integration efforts. Describing mapping requirements is but one of the tasks of an integration effort. By “integration effort” we mean work spent to make separate things work jointly towards some goal. The goals of integration efforts are various including, for example, automating the regular purchase of items under a contract, or implementing new sensing capability in automated production. The stakeholders in integration efforts typically include at least (1) domain experts, the people who can describe what the entities (e.g. business entities, machines) need to do in the joint work, (2) back-end system administrators, that manage the databases involved in the kinds of transactions of interest, and (3) API programmers, who implement the code that achieves the goals of the integration. In small enterprises, of course, many of these roles are played by a single person.

#### 3.1 Mapping Specifications as Data

In this section, we illustrate the RADmapper approach to interoperable mapping with the discussion of three approaches to a simple mapping task. The mapping task is simply to iterate through a collection of objects in source data and, where necessary, map the source `name` key of the object to `customer` in the target data. We illustrate interoperability in the task by showing how various RADmapper implementation of the task “map to” a Mulesoft Dataweave implementation of the task. Note that in the last two sentences we used two notions of mapping: (1) mapping `name` in the source data to `customer` in target data, and (2) mapping a RADmapper implementation of the task to a Dataweave implementation. The first of these “mappings” concerns the domain requirements of a task, the second concerns describing those requirements from different viewpoints. These two notions of mapping correspond to two uses of RADmapper, one of which does the usual work of mapping as we have been doing throughout this document, the second uses the RADmapper code of the first as data. We will walk through the example now to make the above clearer.

To begin, we look at the task of mapping `name` to `customer` in a single object (say, representing orders). In RADmapper we can do this with `$reduceKV`, which is like `$reduce` but is called with a function that accepts key/value pairs of the argument object. In Figure 16, we call `$reduceKV` on Line 9 with the argument function `$name2CustomerFn`, which is defined on Line 6. In Line 7, the body of the function, a conditional expression checks whether the key is equal to “name”, and if so adds a key/value pair “customer”/`$v` to the object being reduced, `$res`, otherwise the key/value pair `$k/$v` is added. In each of these expression `$k` is a key of the argument object and `$v` the corresponding value. Note on Line 9 that `$reduceKV` is called with an empty object as its second argument.

<sup>9</sup>RADmapper is available for use in Java and JavaScript programs as a .jar file and NPM library respectively. A Docker version of the exerciser is available [URL] and hosted at [URL].

Figure 16: One implementation of changing “name” to “customer” in RADmapper.

```

1 ( $order := { 'name'           : 'Example Customer',
2               'shippingAddress' : '123 Mockingbird Lane...',
3               'item part no.'   : 'p12345',
4               'qty'              : { 'amt' : 4, 'uom' : 'unit' } };
5
6 $name2CustomerFn := function($res, $k, $v)
7                     { ($k = 'name') ? $assoc($res, 'customer', $v) : $assoc($res, $k, $v) };
8
9 $reduceKV($name2CustomerFn, {}, $order)
10 )

```

Corresponding Dataweave to Figure 16 is shown in Figure 17. The key work of mapping, which corresponds to `$reduceKV` in the above RADmapper code is Dataweave's `mapObject` on Line 6 of Figure 17.

Figure 17: One implementation of changing “name” to “customer” in Dataweave.

```

1 %dw 2.0
2 output application/json
3 ---
4 items: payload.item ->
5     {
6         item: mapObject (value, key) -> { if (key == "name")
7             { "customer" : value }
8             else
9                 { key : value }
10            }
11     }

```

The relation between the RADmapper and Dataweave viewpoints on this task is achieved by serializing the RADmapper source, Lines 1–10, of Figure 16 to a syntax tree (see Figure 18) and then mapping that data as source, to corresponding Dataweave operators and structure.

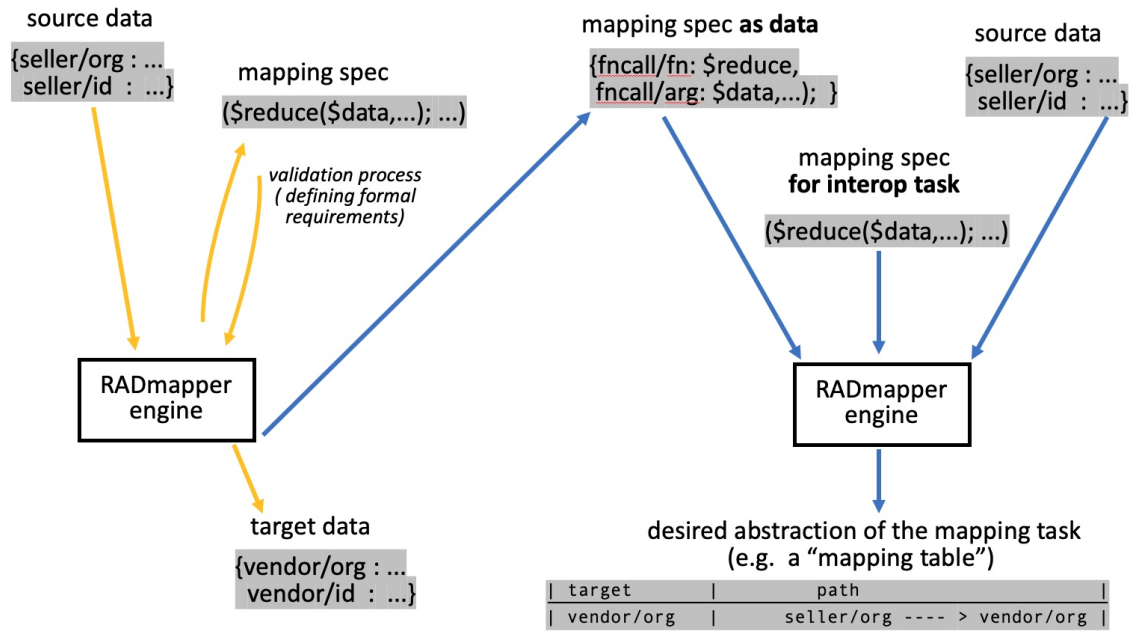
Figure 18: The code of Figure 16 as a syntax tree serialized as JSON.

```

1  {"Block":
2    [{"VarDecl":
3      {"VarName": "$order",
4        "VarValue": {"Obj": {"name": "Example Customer",
5                              "shippingAddress": "123 Mockingbird Lane...",
6                              "item part no.": "pl2345",
7                              "qty": {"Obj": {"amt": "4",
8                                              "uom": "unit"}}}}}},
9      {"VarDecl":
10       {"VarName": "$name2CustomerFn",
11         "VarValue":
12           {"FnDef":
13             {"Params": ["$res", "$k", "$v"]},
14             "Body": {"IfExp":
15               {"Predicate":
16                 {"Block": [{"BinaryExpression": ["$k", "op/eq", "name"]}}},
17               "Then":
18                 {"FnCall": {"Args": ["$res", "customer", "$v"], "FnName": "$assoc"}},
19               "Else":
20                 {"FnCall": {"Args": ["$res", "$k", "$v"], "FnName": "$assoc"}}}}}}},
21       {"FnCall":
22         {"Args": ["$name2CustomerFn", {"Obj": {}}, "$order"],
23           "FnName": "$reduceKV"}}}]

```

Figure 19: The strategy for interoperation involves a second use of RADmapper. Whereas the first use produces a mapping specification (orange arrows) and domain-specific target data, the second (blue arrows) uses that specification as data to generate an abstraction of the original mapping task useful to other mapping tools.



[Discuss simple examples of using ASTs and mapping to "mapping tables". 90 percent of use cases don't need this.]

### 3.2 Don't bother to read past here.

### 3.3 Section Summary and Additional Thoughts

RADmapper's `query` is based on Datalog, a language which was studied extensively in the 1980s [2] and has influenced languages such as SPARQL and the Shapes Constraint Language (SHACL) [?]. `query` has the power of SQL's `query` but does operates fact-at-a-time versus SQL's entity-at-a-time. That distinction entails, for example concerning the previous example, that in Datalog the facts (1) "device1 id is 100." and (2) "device1 status is Ok." are independent assertions that happen to be about the same entity, `device1`. In contrast, SQL's entity-at-a-time design entails that there is a table about devices and the table row at primary key `device1` has information about both `id` and `status`. There are advantages and disadvantages to both Datalog and SQL. For example, it should be apparent that pulling together all the information about an entity with Datalog queries requires relational joins, one each for each fact. On the other hand, RADmapper can learn Datalog-like schema by studying the data, it can create databases in milliseconds, and adding new attributes (columns in SQL) does not require data migration in Datalog.

The order of the patterns in a `query` statement does not matter at all; it does not matter with respect to the execution speed nor the result produced. However, in order to illustrate the chaining of joins being performed, it is useful to the readers of your code to keep joined things together. This is illustrated, for example, by the curvy red line in Figure 11(a).

### 3.4 Working with Tabular Data

Being able to read and write spreadsheet information 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 currently 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 20.

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 20: Table 1 viewed as an array of map structures.

```

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

```

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.

## 4 Datalog features of RADmapper

RADmapper provides a superset of the functionality of JSONata, but more importantly, it supports a second paradigm for navigating data and some very different usage scenarios. 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 form such as RADmapper, 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. RADmapper, and particularly its *AST strategy* described in Section ?? is targeted towards new integration scenarios that involve higher levels of automation and joint (human/AI) cognitive work.

### 4.1 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. Once you have the binding sets, you are halfway there; what remains is to use the binding sets to produce target data. This section describes that process, beginning with definitions of some terms just used:

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

For example, we could have a data in triples 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) connects to the information about Bob.

**in-place update** the idea that the mapping task might involve updating an existing collection of data, rather than defining new data about it.

For example, on Bob’s 30th birthday, we don’t just add the fact `{ 'person/name' : 'Bob', 'person/age' : 30 }`, we retract the fact that Bob is 29 and assert the new fact.



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, how do we know (a) that an action is intended to update some referenced data rather than add new additional information about it, and (b) whether to use a pointer to reference some data rather than just put the data there? We need more information about the data. Specifically, to perform these more complex mapping tasks we have to know: (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 phone number might be recognized as a customer.

#### 4.1.1 Complex mapping task example

The following example illustrates mapping of a network of Web Ontology Language (OWL) data to a relational form. The source OWL data used is simplified from actual OWL data to allow easier discussion. The data consists of objects of two types, one is OWL classes; these have the value `owl/Class` in their `rdf/type` attribute. The other is OWL properties (`owl/ObjectProperty`). The simplifications include using single values for `rdfs/domain`, `rdfs/range`, and `rdfs/subClassOf`. An OWL class and property is depicted in Figure 21.

Figure 21: An object (`owl/Class`) and a relation (`owl/ObjectProperty`) in the source population. These are somewhat simplified from realistic OWL data.

```

1 {'resource/iri'      : 'dol/endurant',
2  'resource/name'     : 'endurant',
3  'resource/namespace': 'dol',
4  'rdf/type'          : 'owl/Class',
5  'rdfs/comment'       : ['The main characteristic of endurants is...'],
6  'rdfs/subClassOf'   : :dol/spatio-temporal-particular,
7  'owl/disjointWith'  : ['dol/abstract', 'dol/quality', 'dol/perdurant']}
8
9 {'resource/iri'      : 'dol/participant',
10 'resource/name'     : 'participant',
11 'resource/namespace': 'dol',
12 'rdf/type'          : 'owl/ObjectProperty',
13 'rdfs/comment'       : ['The immediate relation holding between endurants and perdurants...'],
14 'owl/inverseOf'     : 'dol/participant-in',
15 'rdfs/domain'       : 'dol/perdurant',
16 'rdfs/range'        : 'dol/endurant'}
```

The target data we'll be creating consists of three kinds of things: schema, tables, and columns. Even with this simple data, there are a few options for designing the relational schema. The following are design choices that define the form of the mapping target:

1. Both `owl/Class` and `owl/ObjectProperty` can have `rdfs/comment` and the relationship is one-to-many. A single table with keys consisting of the resource IRI and comment text will suffice.
2. `rdf/type` is one-to-one with the class. Though the possible values are limited to just a few such as `owl/Class` and `owl/ObjectProperty`, we will represent it with a string naming the type.
3. `resource/name` and `resource/namespace` are also one-to-one and are just the two parts of `resource/iri` and could be computed, but we will store these in the class table too.
4. We will assume `rdfs/domain`, `rdfs/range`, and `rdfs/subClassOf` are single-valued. Typically they are not in a real OWL ontology.
5. We will assume that we want to support storage of individuals of the types defined by OWL classes. Though our approach here is not at all reflective of description logic, where class subsumption is the primary kind of inference, mapping will produce a two-column table where one column is the individual's IRI and the other is the foreign key of a class to which it belongs.
6. We will assume that all relations are conceptually binary. Thus storing individuals means that for each `owl/ObjectProperty` mapping will produce a two-column table to represent both a relation and its inverse ("inverse pairs") where an

inverse is defined. Such a table works in both directions (the relation and its inverse), so we will have to prevent creating a table for one member of each inverse pair.

With the above considerations in mind, it becomes apparent that some of the work involves nothing more than storing class and property metadata into tables we can define ahead of time. The `owl/ObjectProperty` definitions, however, entails one new table for each relation (or relation pair if an inverse is defined). The rows of these tables would be populated by instances of the classes specified by `rdfs/domain` and `rdfs/range`. This information is not in the ontology, but Figure 22 depicts the static tables in typical relational DDL. These are populated by the ontology `owl/Class` content of the Figure 23 depicts tables that would be created.

Figure 22: Static DDL for storing class and object relation metadata. The mapping will generate information equivalent to DML to populate this from the source data.

```

1  CREATE SCHEMA typicalOWL;
2
3  CREATE TABLE ObjectDefinition
4      (resourceIRI      VARCHAR(300) primary key,
5       resourceLabel    VARCHAR(300) not null,
6       resourceNamespace VARCHAR(300) not null);
7
8  CREATE TABLE ClassDefinition
9      (resourceIRI VARCHAR(300) primary key,
10     subClassOf   VARCHAR(300) references ClassDefinition);
11
12 CREATE TABLE ObjectClass
13     (resourceIRI VARCHAR(300) primary key,
14     class        VARCHAR(300) references ClassDefinition);
15
16 CREATE TABLE DisjointClass
17     (disjointID  INT primary key,
18     disjoint1   VARCHAR(300) not null references ObjectDefinition,
19     disjoint2   VARCHAR(300) not null references ObjectDefinition);
20
21 CREATE TABLE ResourceComment
22     (commentID   INT primary key,
23     resourceIRI  VARCHAR(300) not null references ObjectDefinition,
24     commentText  VARCHAR(900) not null);
25
26 CREATE TABLE PropertyDefinition
27     (resourceIRI  VARCHAR(300) primary key,
28     relationDomain VARCHAR(300) references ObjectDefinition,
29     relationRange  VARCHAR(300) references ObjectDefinition);

```

Figure 23 depicts the result of mapping data from Figure 21 using a mapping specification that will be described below.

Figure 23: Result of mapping the data depicted in Figure 21. This specifies content equivalent to (1) DML to capture metadata for owl/Class and owl/ObjectProperty objects in the static tables defined above, and (2) DDL to create tables for owl/ObjectProperty objects.

```

1  {'instance-of' : 'insert-row',
2  'table'       : 'ObjectDefinition',
3  'content'     : [{ 'resourceIRI'   : 'dol/endurant'},
4                   { 'resourceLabel' : 'endurant'},
5                   { 'resourceNamespace' : 'dol' }]}
6
7  {'instance-of' : 'insert-row',
8  'table'       : 'ClassDefinition',
9  'content'     : [{ 'resourceIRI'   : 'dol/endurant'},
10                  { 'subClassOf'    : 'dol/spatio-temporal-particular' }]}
11
12 {'instance-of' : 'insert-row',
13 'table'       : 'DisjointClass',
14 'content'     : [{ 'disjointID'    : 1},
15                  { 'disjoint1'     : 'dol/endurant'},
16                  { 'disjoint2'     : 'dol/abstract' }]} /* ... (Two more disjoints elided.) */
17
18 {'instance-of' : 'insert-row',
19 'table'       : 'ResourceComment',
20 'content'     : [{ 'commentID'     : 1},
21                  { 'resourceIRI'   : 'dol/endurant'},
22                  { 'commentText'   : 'The main characteristic of endurants is...' }]}
23
24 /* Similar content for the ObjectProperty dol/participant is elided. */
25
26 {'instance-of' : 'insert-row',
27 'table'       : 'PropertyDefinition',
28 'content'     : [{ 'resourceIRI'   : 'dol/participant'},
29                  { 'relationDomain' : 'dol/perdurant'},
30                  { 'relationRange'  : 'dol/endurant' }]}
31
32 /* The DDL for the participant table: */
33
34 {'instance-of' : 'create-table',
35 'table'       : 'DOLparticipant',
36 'columns'     : [{ 'colName'      : 'propertyID',
37                   'dtype'        : { 'type' : 'varchar', 'size' : 300, 'key' : 'primary' },
38                   { 'colName'    : 'role1',
39                     'dtype'      : { 'type' : 'varchar', 'size' : 300, 'ref' : 'ObjectDefinition' },
40                   { 'colName'    : 'role2',
41                     'dtype'      : { 'type' : 'varchar', 'size' : 300, 'ref' : 'ObjectDefinition' } }]}

```

Figure 24 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 express 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 `express`.

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

Figure 24: Mapping the example OWL to the relational database schema

```

1  ( $data := $read('data/testing/owl-example.edn');
2
3  $qtype := query($rdfType, $extraTrips)
4      { [{?class :rdf/type      $rdfType}
5         [{?class :resource/iri  ?class-iri}
6         [{?class :resource/namespace ?class-ns}
7         [{?class :resource/name    ?class-name}
8         /* ToDo: $extraTrips */
9         ]; /* Defines a higher-order function, a template of sorts. */
10
11  $settype := enforce($tableType)
12      { { 'instance-of' : 'insert-row',
13          'table'       : $tableType,
14          'content'     : { 'resourceIRI'       : ?class-iri,
15                          'resourceNamespace' : ?class-ns,
16                          'resourceLabel'      : ?class-name}}
17      }; /* Likewise, for an enforce template. */
18      /* The target tables for objects and relations a very similar. */
19
20  $quClass := $qtype('owl/Class'); /* Use the template, here and the next three assignments. */
21
22  /* This one doesn't just specify a value for $rdfType, but for $extraTrips. */
23  $quProp := $qtype('owl/ObjectProperty'); /* ToDo: ,queryTriples([?class :rdfs/domain ?domain] [{?class :rdfs
24      /range ?range]}); */
25  $enClassTable := $settype('ClassDefinition');
26  $enPropTable := $settype('PropertyDefinition');
27
28  /* Run the class query; return a collection of binding sets about classes. */
29  $clasBsets := $quClass($data);
30
31  /* We start enforcing with no data, thus the third argument is []. */
32  $star_data := $reduce($clasBsets, $enClassTable, []);
33
34  /* Get bindings sets for the ObjectProperties and make similar tables. */
35  $propBsets := $quProp($data);
36
37  /* We pass in the target data created so far. */
38  $reduce($propBsets, $enPropTable, $star_data) /* The code block returns the target data. */

```

The example dataset from DOLCE is fairly large. Let's suppose it contains 50 classes, each involved in 10 relations. That suggests a collection of 500 binding sets. 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.

## References

- [1] jsonata.org. JSONata: Query and transformation language. <https://jsonata.org>, 2021.
- [2] Serge Abiteboul, Richard Hull, and Victor Vianu. *Foundations of Databases*. Addison-Wesley, Reading, MA, 1995.