# PartiQL Tutorial

ION TEAM

# 1 Getting Started

PartiQL provides an interactive shell, or Read Eval Print Loop (REPL), that allows users to write and evaluate PartiQL queries.

## 1.1 Prerequisites

PartiQL requires the Java Runtime (JVM) to be installed on your machine. You can obtain the *latest* version of the Java Runtime from either

1. OpenJDK, and OpenJDK for Windows

2. Oracle

## 1.2 Download the PartiQL REPL

Each release of PartiQL comes with an archive that contains the PartiQL REPL as a zip file.

1. Download the latest `partiql-cli` zip archive to your machine.
2. Expand (unzip) the archive on your machine. Expanding the archive yields the following folder structure

```
 1   ├──
 2   partiql-cli
 3     ├── bin
 4     │   ├── partiql
 5     │   └── partiql.bat
 6     ├── lib
 7     │   └── ...
 8     ├── README.md
 9     └── Tutorial
10       ├── code
11       │   └── ...
12       ├── tutorial.html
13       └── tutorial.pdf
```

We have used ellipsis ... to elide files/directories.

The root folder `partiql-cli` contains a `README.md` file and 3 subfolders

1. The folder `bin` contains startup scripts `partiql` for OSX (Mac) and Unix systems and `partiql.bat` for Windows systems. Execute these files to start the REPL

2. The folder `lib` contains all the necessary java libraries needed to run PartiQL.

3. The folder `Tutorial` contains the tutorial in `pdf` and `html` form. The subfolder `code` contains 3 types of files

   1. Data files with the extension `.env`. These files contains PartiQL data that we can query

   2. PartiQL query files with the extension `.sql`. These files contain the PartiQL queries used in the tutorial.

   3. Sample query output files with the extension `.out`. These files contain sample output from running the tutorial queries on the appropriate data.

## 1.3 Running the PartiQL REPL

### 1.3.1 Windows

Run (double click) on `particl.bat`. This should open a command line prompt and start the PartiQL REPL. The PartiQL REPL prompt should look like this

```
1   Welcome to the PartiQL REPL!
2   PartiQL>
```

### 1.3.2 OSX (Mac) and Unix

1. Open a terminal and navigate to the `partiql-cli` folder we created when we extracted `partiql-cli.zip`.

2. Run the executable `partiql` file, by typing `./partiql` and hit enter. This should start the PartiQL REPL and should look like this

```
1   Welcome to the PartiQL REPL!
2   PartiQL>
```

## 1.4 Testing the PartiQL REPL

Let's write a simple query to verify that our PartiQL REPL is working. At the `PartiQL>` prompt type

```
1   PartiQL> SELECT * FROM [1,2,3]
```

and press ENTER *twice*. The output should look similar to

```
1   PartiQL> SELECT * FROM [1,2,3]
2       |
3   ==='
```

```
 4  <<
 5    {
 6      '_1': 1
 7    },
 8    {
 9      '_1': 2
10    },
11    {
12      '_1': 3
13    }
14  >>
15  ---
16  OK! (86 ms)
17  PartiQL>
```

**INFO**

An easy way to load the necessary data into the REPL while reading our Tutorial is to pass one of the

`.env` files to the REPL on invocation, e.g. on Unix or OSX while on inside `partiql-cli` folder,

```
 1   ./bin/partiql  -e code/q1.env
```

You can then see what is loaded in the REPL's global environment using the **special** REPL command

`!global_env`, e.g.,

```
 1  Welcome to the PartiQL REPL!
 2  PartiQL> !global_env
 3     |
 4  ==='
 5  {
 6    'hr': {
 7      'employees': <<
 8        {
 9          'id': 3,
10          'name': 'Bob Smith',
11          'title': NULL
12        },
13        {
14          'id': 4,
15          'name': 'Susan Smith',
16          'title': 'Dev Mgr'
17        },
18        {
```

```
19          'id': 6,
20          'name': 'Jane Smith',
21          'title': 'Software Eng 2'
22        }
23      >>
24    }
25  }
26  ---
27  OK! (6 ms)
```

Congratulations! You succesfuly installed and run the PartiQL REPL. The PartiQL REPL is now waiting for more input.

To exit the PartiQL REPL press

- `Control`+D in OSX or Unix
- `Control`+Z on Windows

or close the terminal/command prompt window.

## 2 Introduction

PartiQL provides SQL-compatible unified query access across multiple data stores containing structured, semi-structured and nested data that is supported by SQL. PartiQL separates the syntax and semantics of a query from the underlying data source and/or data format of the data. It enables users to interact with data with[1] or without regular schema.

This tutorial aims to teach SQL users the PartiQL extensions to SQL. The tutorial is primarily driven by "how to" examples.

For the reader who is interested in the full detail and formal specification of PartiQL, we recommend the 2-tiered PartiQL formal specification: The formal specification first describes the *PartiQL core*, which is a short and concise functional programming language. Then the specification layers SQL compatibility through syntactic sugar that shows how SQL features can be translated to semantically equivalent core PartiQL expressions. These translations presented as syntactic sugar enable SQL compatibility.

---

[1]The implementation currently only supports data without schema. Schema support is forthcoming.

# 3 PartiQL Queries are SQL compatible

PartiQL is backwards compatible with SQL-92[2]. We will see what compatibility means when it is used to query data found in data formats and data stores that are not SQL.

For starters, given the table `hr.employees`

| Id | name | title |
|----|------|-------|
| 3 | Bob Smith | null |
| 4 | Susan Smith | Dev Mgr |
| 6 | Jane Smith | Software Eng 2 |

the following SQL query

```
1  SELECT e.id,
2         e.name AS employeeName,
3         e.title AS title
4  FROM hr.employees e
5  WHERE e.title = 'Dev Mgr'
```

is also a valid PartiQL query. As we know from SQL, when this query operates on the table `hr.employees` it will return the result

| Id | employeeName | title |
|----|-------------|-------|
| 4 | Susan Smith | Dev Mgr |

## 3.1 PartiQL data model: Abstraction of many underlying data storage formats

PartiQL implementations operate not just on SQL tables but also on data that may have nesting, union types, different attributes across different tuples and many other features that we often find in today's nested and/or semi-structured formats, like JSON, Parquet, etc.

To capture this generality, PartiQL is based on a logical type system: the *PartiQL data model*. Each PartiQL implementation maps data formats, like JSON, Parquet etc., into a PartiQL data set that follows the PartiQL data model. PartiQL queries work on the PartiQL data set abstraction.

---

[2]SQL-92

For example, the table `hr.employees` is denoted in the PartiQL data model as this dataset

```
1  {
2      'hr': {
3          'employees': <<
4              { 'id': 3, 'name': 'Bob Smith',   'title': null }, -- a tuple is denoted by {
                     ... } in the PartiQL data model
5              { 'id': 4, 'name': 'Susan Smith', 'title': 'Dev Mgr' },
6              { 'id': 6, 'name': 'Jane Smith',  'title': 'Software Eng 2'}
7          >>
8      }
9  }
```

Notice that the `employees` is nested within `hr`.

The delimiters << ... >> denote that the data is an *unordered collection* (also known as *bag*), as is the case with SQL's tables. That is, there is no order between the three tuples. Single line comments start with `--` and end at the end of the line.

A very different kind of data source may lead to the same PartiQL dataset. For example, a set of JSON files that contain the following JSON objects

```
1  {
2      "hr" : {
3          "employees": [
4              { "id": 3, "name": "Bob Smith",   "title": null },
5              { "id": 4, "name": "Susan Smith", "title": "Dev Mgr" },
6              { "id": 6, "name": "Jane Smith",  "title": "Software Eng 2"}
7          ]
8      }
9  }
```

will likely[3] be abstracted by a PartiQL-supporting implementation into the identical PartiQL abstraction with the `hr.employees` table.

**Remark:** You will keep noticing the similarity of the PartiQL notation with the JSON notation. Notice also the subtle differences: In the interest of SQL compatibility, a PartiQL literal is quoted, while JSON literals are double-quoted.

**Remark:** You may conceptually think that a deserializer inputs JSON and outputs the PartiQL data set. But do not assume that the query processing of a PartiQL implementation will have to actually parse and

---

[3]The JSON value attached to `employee` is an *ordered* list. PartiQL implementations may provide their own mappings from popular data formats, e.g., CSV, TSV, JSON, Ion etc., to the PartiQL data model and/or allow clients to implements their own mappings.

abstract into PartiQL each and every bit of the underlying data storage. For example, AWS services like Redshift Spectrum and QLDB do much smarter things in order to evaluate your PartiQL queries efficiently.

Back to our PartiQL query

```
1   SELECT e.id,
2          e.name AS employeeName,
3          e.title AS title
4   FROM hr.employees e
5   WHERE e.title = 'Dev Mgr'
```

evaluates in PartiQL and returns

```
1   <<
2     {
3       'id': 4,
4       'employeeName': 'Susan Smith',
5       'title': 'Dev Mgr'
6     }
7   >>
8   ---
9   OK! (16 ms)
```

the result remains the same, no matter whether the `hr.employees` were the SQL table or the JSON file. All that is needed is that a catalog associates the *name* `hr.employees` with the PartiQL abstraction of the JSON data.

In the same spirit, the same PartiQL abstraction may come from a CVS file or a Parquet file, a format that has gained big traction, thanks to the efficient way in which it stores data. Again, the same query makes perfect sense, regardless of what exactly was the storage format behind `hr.employees`.

### 3.1.1 Learn more

- **PartiQL data sets look very much like JSON.**

  What are the differences? Indeed, PartiQL adopts the tuple/object and array notation of JSON. However, the PartiQL string literals are denoted by single quotes. Importantly, the scalar types of PartiQL are the ones of SQL and not just strings, numbers and booleans, as in JSON.

- **Do implementations need to have a catalog?**

  If queries refer to names, a catalog logically validates whether the name exists or not. However, we will also see PartiQL queries that refer to no names.

# 4 Querying Nested Data

SQL-92 only has tables that have tuples that contain scalar values. A key feature of many modern formats is nested data. That is, attributes whose values may be themselves tables (i.e., collections of tuples), or may be arrays of scalars, or arrays of arrays and many other combinations. We next present PartiQL's features (SQL extensions) that allow us to work with nested data.

We also include sections titled "Use Case". Such "Use Case" sections do not introduce additional features. They merely show how to combine the few novel PartiQL features with standard SQL features in order to solve a large number of problems.

## 4.1 Nested Collections

Let's now add the nested attribute `projects` into the data set.

```
 1  {
 2    'hr': {
 3        'employeesNest': <<
 4          {
 5            'id': 3,
 6            'name': 'Bob Smith',
 7            'title': null,
 8            'projects': [ { 'name': 'AWS Redshift Spectrum querying' },
 9                          { 'name': 'AWS Redshift security' },
10                          { 'name': 'AWS Aurora security' }
11                        ]
12          },
13          {
14            'id': 4,
15            'name': 'Susan Smith',
16            'title': 'Dev Mgr',
17            'projects': []
18          },
19          {
20            'id': 6,
21            'name': 'Jane Smith',
22            'title': 'Software Eng 2',
23            'projects': [ { 'name': 'AWS Redshift security' } ]
24          }
25        >>
26      }
27  }
```

Notice that the value of `'projects'` is an array. Arrays are denoted by `[ ... ]` with array elements separated by `,`. In our example the array happens to be an array of tuples. We will see that arrays may be arrays of anything, not just arrays of tuples.

### 4.1.1 Unnesting a Nested Collection

The query finds the names of employees who work on projects that contain the string `'security'` and outputs them along with the name of the `'security'` project. Notice that the query has just one extension over standard SQL – the `e.projects AS p` part.

```
1  SELECT e.name AS employeeName,
2         p.name AS projectName
3  FROM hr.employeesNest AS e,
4       e.projects AS p
5  WHERE p.name LIKE '%security%'
```

The output of our query is

```
1  <<
2    {
3      'employeeName': 'Bob Smith',
4      'projectName': 'AWS Redshift security'
5    },
6    {
7      'employeeName': 'Bob Smith',
8      'projectName': 'AWS Aurora security'
9    },
10   {
11     'employeeName': 'Jane Smith',
12     'projectName': 'AWS Redshift security'
13   }
14 >>
15 ---
16 OK! (51 ms)
```

The extension over SQL is the `FROM` clause item `e.projects AS p`. Standard SQL would attempt to find a schema named `e` with a table `projects` and since in our example there isn't an `e.projects` table, the query would fail. In contrast, PartiQL will dereference `e.projects` to the attribute `projects` of `e`.

Once we allow this extension, the semantics are alike SQL's. The alias (also called *variable* in PartiQL) `e` gets bound to each employee, in turn. For each employee, the variable `p` gets bound to each project of the employee, in turn. Thus the query's meaning, alike SQL, is

foreach employee tuple `e` from `hr.employeesNest`
    foreach project tuple `p` from `e.projects`
        if `p.name LIKE '%security%'`
        output `e.name AS employeeName`, `p.name AS projectName`

Notice that our query involved variables that were ranging over nested collections (`p` in the example), along with variables that were ranging over tables (`e` in the example), as standard SQL aliases do. All variables, no matter what they range over, can be used wherever in the `FROM`, `WHERE`, `SELECT` clauses as we will see in the examples that follow.

### 4.1.2 Learn more

- **Can I only unnest arrays of tuples?**

  No, anything can be unnested. For example, arrays of scalars, etc.

- **Does `e.projects AS p` have to appear in the same `FROM` clause that defines `e`?**

  No. For example, see below the use cases that involve subqueries. There, the `e` and `p` are defined in separate `FROM` clauses.

- **How could I force `e.projects` to refer to the nested attribute `projects` even if there were a schema named `e` with a table `projects`?**

  Use the syntax `@e.projects`. Recall, in the absence of the `@`, in the interest of SQL compatibility, PartiQL will first attempt to dereference the `e.projects` against the catalog.

- **SQL allows me to avoid writing an explicit alias `e` when I write, say, `e.name`. Can I avoid writing the `e` in PartiQL as well?**

  SQL allows us to avoid writing aliases (variables) when the schema of the tables allows correct dereferencing. PartiQL does the same. However, recall, a schema is not necessary for a PartiQL data set. Indeed, our example has not assumed a schema. Then , in the absence of a schema, you cannot omit the aliases (variables). For example, if you write just `name` and there is no schema, PartiQL cannot tell whether you mean employee name or project name. Thus you need to explicitly write the alias (variable).

  There is one exception to this rule: If your query has a single item in its `FROM` clause, you can omit the alias (variable). Eg, you can write

  ```
  1   SELECT name FROM hr.employeesNest
  ```

  In this case it is apparent that `name` may only be an employee name and thus PartiQL allows you to not provide an alias (variable).

Nevertheless, for clarity we recommend that you always use aliases (variables) and this is what this tutorial does.

- **If there is a schema, can I avoid writing the alias `p`?**

  No. The `p` has to be written in order to denote the iteration over the projects.

### 4.1.3 Unnesting Nested Collections Using `JOIN`

In this section, we simply present an alternate way to express and think about unnesting collections.

One may think that the `FROM` clause of the example executes, in a sense, a `JOIN` between employees and projects. Except that unlike a conventional SQL join that would require an **ON condition**, the employees-projects join condition is implicit in the nesting of the projects data into the employee data. If it helps you to think in terms of `JOIN`, you may replace the comma with `JOIN`. That is, the following two queries are equivalent.

```
1   SELECT e.name AS employeeName,
2          p.name AS projectName
3   FROM hr.employeesNest AS e,
4       e.projects AS p
5   WHERE p.name LIKE '%security%'
```

```
1   SELECT e.name AS employeeName,
2          p.name AS projectName
3   FROM hr.employeesNest AS e JOIN
4       e.projects AS p
5   WHERE p.name LIKE '%security%'
```

### 4.1.4 Unnesting data with LEFT JOIN always preserves parent information

Assume that we want to write a query that returns as a bag of tuples the entire employee and project information from `hr.employeesNest`. The query result we want is this bag of tuples with attributes `id`, `employeeName`, `title` and `projectName`:

```
1   <<
2     {
3       'id': 3,
4       'employeeName': 'Bob Smith',
5       'title': NULL,
6       'projectName': 'AWS Redshift Spectrum querying'
7     },
8     {
9       'id': 3,
10      'employeeName': 'Bob Smith',
```

```
11        'title': NULL,
12        'projectName': 'AWS Redshift security'
13     },
14     {
15        'id': 3,
16        'employeeName': 'Bob Smith',
17        'title': NULL,
18        'projectName': 'AWS Aurora security'
19     },
20     {
21        'id': 4,
22        'employeeName': 'Susan Smith',
23        'title': 'Dev Mgr'
24     },
25     {
26        'id': 6,
27        'employeeName': 'Jane Smith',
28        'title': 'Software Eng 2',
29        'projectName': 'AWS Redshift security'
30     }
31  >>
32  ---
33  OK! (6 ms)
```

Notice that there is a `'Susan Smith'` tuple in the result, despite the fact that Susan has no project. Susan's `projectName` is **null**. We can obtain this result by combining employees and projects using the `LEFT JOIN` operator, as follows:

```
1  SELECT e.id AS id,
2         e.name AS employeeName,
3         e.title AS title,
4         p.name AS projectName
5  FROM hr.employeesNest AS e LEFT JOIN e.projects AS p
```

The semantics of this query can be thought of as

foreach employee tuple `e` from `hr.employeesNest`
   if the `e.projects` is an empty collection then // *this part is special about LEFT JOINs*
     output `e.id AS id`, `e.name AS employeeName`, `e.title AS title`
     and output a **null** `AS projectName`
   else // *the following part is identical to plain (inner) JOINs*
     foreach project tuple `p` from `e.projects`

output `e.id AS id`, `e.name AS employeeName`, `e.title AS title`

and output a `null AS projectName`

### 4.1.5 Use Case: Checking whether a nested collection satisfies a condition

The following use cases employ the unnesting features, which we have already discussed, in new use cases. A lesson that emerges is that we can use variables (SQL aliases) that range over nested data as if they were standard SQL aliases. This realization gives us the power to solve a great number of use cases just be combining the unnesting features with features we already know from standard SQL.

In our first use case we want a query that returns the names of the employees that are involved in a project that contains the word `'security`. The solution employs SQL's "`EXISTS` (subquery)" feature, along with unnesting:

```
1  SELECT e.name AS employeeName
2  FROM hr.employeesNest AS e
3  WHERE EXISTS ( SELECT *
4                 FROM e.projects AS p
5                 WHERE p.name LIKE '%security%')
```

returns

```
 1  <<
 2    {
 3      'employeeName': 'Bob Smith'
 4    },
 5    {
 6      'employeeName': 'Jane Smith'
 7    }
 8  >>
 9  ---
10  OK! (14 ms)
```

In the second use case we want a query that outputs the names of the employees that have more than one security projects and we are aware of a key for employees (e.g., an attribute that is guaranteed to have a unique value for each employee). We can find the requested employees by utilizing a combination of `GROUP BY` and `HAVING`.[4] In our example, let's assume that the `id` attribute is a primary key for the employees. Then we could find the employees with more than one security project with this query:

```
1  SELECT e.name AS employeeName
```

---

[4]We could also have used the > operator with the subquery's result, but a current issue with the implementation currently prevents us from doing so.

```
2  FROM hr.employeesNest e,
3       e.projects AS p
4  WHERE p.name LIKE '%security%'
5  GROUP BY e.id, e.name
6  HAVING COUNT(*) > 1
```

which returns

```
1  <<
2    {
3      'employeeName': 'Bob Smith'
4    }
5  >>
6  ---
7  OK! (28 ms)
```

### 4.1.6 Use Case: Subqueries that aggregate over nested collections

Next, let's find how many querying projects (that is, projects whose name contains the word "querying") each employee has.[5]

Making the same asssumption as before, that `id` is a key for employees, we can solve the problem with the query

```
1  SELECT e.name AS employeeName,
2         COUNT(p.name) AS queryProjectsNum
3  FROM hr.employeesNest e LEFT JOIN e.projects AS p ON p.name LIKE '%querying%'
4  GROUP BY e.id, e.name
```

that returns

```
1  <<
2    {
3      'employeeName': 'Bob Smith',
4      'queryProjectsNum': 1
5    },
6    {
7      'employeeName': 'Susan Smith',
8      'queryProjectsNum': 0
9    },
10   {
```

---

[5]We could also have used the > operator with the subquery's result, but a current issue with the implementation currently prevents us from doing so.

```
11        'employeeName': 'Jane Smith',
12        'queryProjectsNum': 0
13    }
14  >>
15  ---
16  OK! (22 ms)
```

Notice, this query's result includes Susan Smith and Jane Smith, who have no querying projects.

## 4.2 Nested Tuple Values and Multi-Step Paths

A value may also be a tuple – also called object and struct in many models and formats. For example, the project value in the following tuples is always a tuple with project name and project org.

```
 1  {
 2      'hr': {
 3          'employeesWithTuples': <<
 4              {
 5                  'id': 3,
 6                  'name': 'Bob Smith',
 7                  'title': null,
 8                  'project': {
 9                      'name': 'AWS Redshift Spectrum querying',
10                      'org': 'AWS'
11                  }
12              },
13              {
14                  'id': 6,
15                  'name': 'Jane Smith',
16                  'title': 'Software Eng 2',
17                  'project': {
18                      'name': 'AWS Redshift security',
19                      'org': 'AWS'
20                  }
21              }
22          >>
23      }
24  }
```

PartiQL's multistep paths enable navigating within tuples. For example, the following query finds AWS projects and outputs the project name and employee name.

```
1  SELECT e.name AS employeeName,
2         e.project.name AS projectName
3  FROM hr.employeesWithTuples e
4  WHERE e.project.org = 'AWS'
```

The result is

```
1  <<
2    {
3      'employeeName': 'Bob Smith',
4      'projectName': 'AWS Redshift Spectrum querying'
5    },
6    {
7      'employeeName': 'Jane Smith',
8      'projectName': 'AWS Redshift security'
9    }
10 >>
11 ---
12 OK! (25 ms)
```

## 4.3 Unnesting Arbitrary Forms of Nested Collections

The previous examples have shown nested attributes that were arrays of tuples. It need not be the case that the nested attributes are collections of tuples. They may just as well be arrays of scalars, arrays of arrays, and more. In general any combination of data that one can create by composing scalars, tuples and arrays. You need not learn a different set of query language features for each case. The unnesting features, which we have already seen, are sufficient.

### 4.3.1 Use Case: Unnesting Arrays of Scalars

The list of projects associated with each employee in `hr.employeesNest` could have been simply a list of project name strings. Replacing the nested tuples with plain strings gives us

```
1  {
2      'hr': {
3          'employeesNestScalars': <<
4              {
5                  'id': 3,
6                  'name': 'Bob Smith',
7                  'title': null,
8                  'projects': [
```

```
 9                    'AWS Redshift Spectrum querying',
10                    'AWS Redshift security',
11                    'AWS Aurora security'
12                ]
13            },
14            {
15                'id': 4,
16                'name': 'Susan Smith',
17                'title': 'Dev Mgr',
18                'projects': []
19            },
20            {
21                'id': 6,
22                'name': 'Jane Smith',
23                'title': 'Software Eng 2',
24                'projects': [ 'AWS Redshift security' ]
25            }
26        >>
27    }
28 }
```

Let us repeat the previous use cases on the revised employee data.

The following query finds the names of employees who work on projects that contain the string `'security'` and outputs them along with the name of the "security" project.

```
1 SELECT e.name AS employeeName,
2        p AS projectName
3 FROM hr.employeesNestScalars AS e,
4      e.projects AS p
5 WHERE p LIKE '%security%'
```

The preceding query returns

```
 1 <<
 2   {
 3     'employeeName': 'Bob Smith',
 4     'projectName': 'AWS Redshift security'
 5   },
 6   {
 7     'employeeName': 'Bob Smith',
 8     'projectName': 'AWS Aurora security'
 9   },
10   {
```

```
11       'employeeName': 'Jane Smith',
12       'projectName': 'AWS Redshift security'
13     }
14  >>
15  ---
16  OK! (28 ms)
```

The variable `p` ranges (again) over the content of `e.projects`. In this case, since `e.projects` has strings (as opposed to tuples), the variable `p` binds each time to a project name string. Thus, this query can be thought of as executing the following snippet.

foreach employee tuple `e` from `hr.employeesNestScalars`

    foreach project `p` from `e.projects`

        if the string `p` matches `'%security%'`

        output `e.name AS employeeName` and the string `p AS projectName`

### 4.3.2 Use Case: Unnesting Arrays of Arrays

Arrays may also contain arrays, directly, without intervening tuples, as in the `matrices` data set.

```
 1  {
 2      'matrices': <<
 3          {
 4              'id': 3,
 5              'matrix': [
 6                  [2, 4, 6],
 7                  [1, 3, 5, 7],
 8                  [9, 0]
 9              ]
10          },
11          {
12              'id': 4,
13              'matrix': [
14                  [5, 8],
15                  [ ]
16              ]
17
18          }
19      >>
20  }
```

The following query finds every even number and outputs the even number and the `id` of the tuple where it was found.

```
1  SELECT t.id AS id,
2         x AS even
3  FROM matrices AS t,
4       t.matrix AS y,
5       y AS x
6  WHERE x % 2 = 0
```

The preceding query returns

```
1  <<
2    {
3      'id': 3,
4      'even': 2
5    },
6    {
7      'id': 3,
8      'even': 4
9    },
10   {
11     'id': 3,
12     'even': 6
13   },
14   {
15     'id': 3,
16     'even': 0
17   },
18   {
19     'id': 4,
20     'even': 8
21   }
22 >>
23 ---
24 OK! (59 ms)
```

Informally the query's evaluation can be thought of as

foreach tuple `t` from `matrices`
   foreach array `y` from `t.matrix`
      foreach number `x` from `y`
         if `x` is even then

output `t.id AS id` and `x AS even`

# 5  Literals

Literals of the PartiQL query language correspond to the types in the PartiQL data model:

- scalars, including **null** which follow the SQL syntax when applicable. For example:

    - 5

    - `'foo'`

- tuples, denoted by {...} with tuple elements separated by , (also known as structs and/or objects in many formats and other data models)

    - { `'id'`: 3, `'arr'`: [1, 2] }

- arrays, denoted by [...] with array elements separated by ,

    - [ 1, `'foo'`]

- bags, denoted by << ... >> with bag elements separated by a ,

    - << 1, `'foo'`>>

Notice that in the spirit of the PartiQL data model, literals compose freely and any kind of literal may appear within any tuple, array and bag literal, eg.,

```
1  {
2      'id': 3,
3      'matrix': [
4          [2, 4, 6],
5          'NA'
6      ]
7  }
```

# 6  Querying Heterogeneous and Schemaless Data

Many formats do not require a schema that describes the data – that is *schemaless data*. In such cases it is possible to have various "heterogeneities" in the data:

- One tuple may have an attribute x while another tuple may not have this attribute

- In one tuple of a collection an attribute x may be of type, e.g., string, while in another tuple of the same collection the same attribute x may be of a different type – e.g, array.

- The elements of a collection (be it a bag or array) can be heterogeneous (not have the same type). For example, the first element may be a string, the second one may be an integer and the third one an array.

- Generally, any composition is possible as we can bundle heterogeneous elements in arrays and bags.

Heterogeneities are not particular to schemaless. Schemas may allow for heterogeneity in the types of the data. For example, one of the Hive data types is the union type,[6] which allows a value to belong to any one of a list of types. For example, in the following schema the `projects` attribute may be either a string or an array of strings

```
1  CREATE TABLE employeesMixed(
2          id: INT,
3          name: STRING,
4          title: STRING,
5          projects: UNIONTYPE<STRING, ARRAY<STRING>>
6  );
```

A collection of PartiQL tuples that follows this schema could be

```
1  {
2      'hr': {
3          'employeesMixed1': <<
4              {
5                  'id': 3,
6                  'name': 'Bob Smith',
7                  'title': null,
8                  'projects': [
9                      'AWS Redshift Spectrum querying',
10                  'AWS Redshift security',
11                  'AWS Aurora security'
12                  ]
13          },
14          {
15              'id': 4,
16              'name': 'Susan Smith',
17              'title': 'Dev Mgr',
18              'projects': []
19          },
20          {
21              'id': 6,
```

---

[6]Hive Union Type

```
22                    'name': 'Jane Smith',
23                    'title': 'Software Eng 2',
24                    'projects': 'AWS Redshift security'
25               }
26          >>
27      }
28  }
```

Thus we see that data may have heterogeneities – regardless of whether they are described by a schema or not. PartiQL tackles heterogeneous data, in ways that we will see in the next use cases and feature presentations.

## 6.1  Tuples with Missing Attributes

Let's go back to the `hr.employees` table (that is, bag of tuples). Bob Smith has no title and, as is typical in SQL, the lack of title is modeled with the **null** value.

```
1  {
2      'hr': {
3          'employees': <<
4              { 'id': 3, 'name': 'Bob Smith',   'title': null }
5              { 'id': 4, 'name': 'Susan Smith', 'title': 'Dev Mgr' }
6              { 'id': 6, 'name': 'Jane Smith',  'title': 'Software Eng 2'}
7          >>
8      }
9  }
```

Nowadays, many semi-structured formats allow users to represent "missing" information in two ways.

1.  The first way is by use of **null**.
2.  The second kind is the plain absence of the attribute from the tuple.

That is, we can represent the fact that Bob Smith has no title by simply having no `title` attribute in the `'Bob Smith'` tuple:

```
1  {
2      'hr': {
3          'employeesWithMissing': <<
4              { 'id': 3, 'name': 'Bob Smith' }, -- no title in this tuple
5              { 'id': 4, 'name': 'Susan Smith', 'title': 'Dev Mgr' },
6              { 'id': 6, 'name': 'Jane Smith', 'title': 'Software Eng 2'}
7          >>
8      }
```

```
9    }
```

PartiQL does not argue about when to use **null**s and when to use "missing". Myriads of datasets already use one of the two or both. However, PartiQL enables queries to distinguish when they access a **null** Vs when they access a missing attribute. PartiQL also enables queries to create results that have both **null**s and missing attributes. Indeed, it makes it very easy to propagate source data **null**s as query result **null**s and source data missing attributes into result missing attributes.

## 6.2  Accessing and Processing Missing Attributes:  The MISSING Value

Consider again this PartiQL query, which happens to also be an SQL query.

```
1  SELECT e.id,
2           e.name AS employeeName,
3           e.title AS title
4  FROM hr.employeesWithMissing AS e
5  WHERE e.title = 'Dev Mgr'
```

What will happen when the query goes over the Bob Smith tuple, which has no `title`?

The first step to answering this question is understanding the result of the path `e.title` when the alias (variable) `e` binds to the tuple { `'id'`: 3, `'name'`: `'Bob Smith'`}. In more basic terms, what is the result of the expression { `'id'`: 3, `'name'`: `'Bob Smith'`}.`path` ? PartiQL says that it is the special value `MISSING`. `MISSING` behaves very similar to **null**.

### 6.2.1  Evaluating Functions and Conditions with MISSING

If a function (including infix functions like =) inputs a `MISSING` the function's result is also `MISSING`. In the case of the example, this means that the `WHERE` clause `e.title='Dev Mgr'` will evaluate to `MISSING` when `e` binds to { `'id'`: 3, `'name'`: `'Bob Smith'`} and, as usual in SQL, the `WHERE` clause fails when it does not evaluate to **true**. Thus the output will be

```
1  <<
2    {
3      'id': 4,
4      'employeeName': 'Susan Smith',
5      'title': 'Dev Mgr'
6    }
7  >>
8  ---
9  OK! (17 ms)
```

### 6.2.2 Propagating MISSING in Result Tuples

What would happen if a missing attribute or, more generally, an expression returning MISSING appears in the SELECT?

```
1   SELECT e.id,
2          e.name AS employeeName,
3          e.title AS outputTitle
4   FROM hr.employeesWithMissing AS e
```

The query will output one tuple for each employee. When it outputs the Bob Smith tuple, the e.title will evaluate to MISSING and then the output tuple will not even have an outputTitle attribute.

```
 1   <<
 2     {
 3       'id': 3,
 4       'employeeName': 'Bob Smith'
 5     },
 6     {
 7       'id': 4,
 8       'employeeName': 'Susan Smith',
 9       'outputTitle': 'Dev Mgr'
10     },
11     {
12       'id': 6,
13       'employeeName': 'Jane Smith',
14       'outputTitle': 'Software Eng 2'
15     }
16   >>
17   ---
18   OK! (10 ms)
```

The same treatment of MISSING would happen if, say, we had this query that converts titles to capital letters:

```
1   SELECT e.id,
2          e.name AS employeeName,
3          UPPER(e.title) AS outputTitle
4   FROM hr.employeesWithMissing AS e
```

Again, the e.title will evaluate to MISSING for 'Bob Smith', the UPPER(e.title) is then UPPER(MISSING) and also evaluates to MISSING. Thus the result will be

```
 1  <<
 2    {
 3      'id': 3,
 4      'employeeName': 'Bob Smith',
 5      'outputTitle': NULL
 6    },
 7    {
 8      'id': 4,
 9      'employeeName': 'Susan Smith',
10      'outputTitle': 'DEV MGR'
11    },
12    {
13      'id': 6,
14      'employeeName': 'Jane Smith',
15      'outputTitle': 'SOFTWARE ENG 2'
16    }
17  >>
18  ---
19  OK! (20 ms)
```

## 6.3 Variables can range over Data with Different Types

A PartiQL variable (called *alias* in SQL) can bind to data of different types during a query's evaluation. This is unlike SQL where the variables always bind to tuples. It is even different from what happened in Use Case: Unnesting Arrays of Scalars and what happened in Use Case: Unnesting Arrays of Arrays.

In the first use case, the PartiQL variable p happened to always bind to a string (given the particular sample data of the example). In the second use case, the PartiQL variable y was always bound to an array (again, given the particular sample data of the example).

To make the case for variables that bind to different types, consider the following twist in the employeesNest data set. Some of the elements of the projects array are plain strings and some are tuples. Even the employee tuples do not always have the same attributes.

```
1  {
2      'hr': {
3          'employeesMixed2': <<
4              {
5                  'id': 3,
6                  'name': 'Bob Smith',
7                  'title': null,
8                  'projects': [
```

```
 9                         { 'name': 'AWS Redshift Spectrum querying' },
10                         'AWS Redshift security',
11                         { 'name': 'AWS Aurora security' }
12                     ]
13             },
14             {
15                 'id': 4,
16                 'name': 'Susan Smith',
17                 'title': 'Dev Mgr',
18                 'projects': []
19             },
20             {
21                 'id': 6,
22                 'name': 'Jane Smith',
23                 'projects': [ 'AWS Redshift security']
24             }
25         >>
26     }
27 }
```

This query on `hr.employeesMixed2` produces employee name – employee project pairs.

```
1 SELECT e.name AS employeeName,
2        CASE WHEN (p IS TUPLE) THEN p.name
3        ELSE p END AS projectName
4 FROM hr.employeesMixed2 AS e,
5      e.projects AS p
```

Notice the sub-expression (`p IS TUPLE`). The `IS` operator can be used to check a value against it's type at evaluation time. Notice also that the variable `p` binds to different types.

In general, the `FROM` clause of a query binds its variables (aliases) to data. The variables need not bind to data that have the same types. Each binding is fed to the `SELECT` clause, which evaluates its expressions.

This table shows each variables' binding produced by the `FROM` clause and the corresponding tuple output by the `SELECT` clause.

## Variable e

```
 1  { 'id': 3,
 2
 3  'name': 'Bob Smith',
 4
 5  'title': null,
 6
 7  'projects':  [ {
 8  'name': 'AWS Redshift
 9  Spectrum querying' },
10
11  'AWS Redshift
12  security',
13
14  { 'name': 'AWS Aurora
15  security' }
16
17   ]
18
19  }
```

## Variable p

```
 1  { 'name': 'AWS
 2  Redshift Spectrum
 3  querying' }
```

## Result tuple

```
 1  {
 2
 3  'employeeName': 'Bob
 4  Smith',
 5
 6  'projectName': 'AWS
 7  Redshift Spectrum
 8  querying'
 9
10  }
```

## Variable e

```
 1  { 'id': 3,
 2
 3  'name': 'Bob Smith',
 4
 5  'title': null,
 6
 7  'projects':  [ {
 8  'name': 'AWS Redshift
 9  Spectrum querying' },
10
11  'AWS Redshift
12  security',
13
14  { 'name': 'AWS Aurora
15  security' }
16
17    ]
18
19  }
```

## Variable p

```
 1  'AWS Redshift
 2  security'
```

## Result tuple

```
 1  {
 2
 3  'employeeName': 'Bob
 4  Smith',
 5
 6  'projectName': 'AWS
 7  Redshift security'
 8
 9  }
```

| Variable e | Variable p | Result tuple |
|---|---|---|

```
1   { 'id': 3,
2
3   'name': 'Bob Smith',
4
5   'title': null,
6
7   'projects': \[ {
8   'name': 'AWS Redshift
9   Spectrum querying' },
10
11  'AWS Redshift
12  security',
13
14  { 'name': 'AWS Aurora
15  security' }
16
17  \]
18
19  }
```

```
1   { 'name': 'AWS Aurora
2   security' }
```

```
1   {
2
3   'employeeName': 'Bob
4   Smith',
5
6   'projectName': 'AWS
7   Aurora security'
8
9   }
```

```
1   { 'id': 6,
2
3   'name': 'Jane Smith',
4
5   'projects': \[ 'AWS
6   Redshift security' \]
7
8   }
```

```
1       'AWS Redshift
2   security'
```

```
1   {
2
3     'employeeName':
4   'Jane Smith',
5
6     'projectName': 'AWS
7   Redshift security'
8
9   }
```

# 7 Accessing Array Elements by Order

SQL allows us to order the output of a query using the `ORDER BY` clause. However, the SQL data model does not recognize order in the input data. In contrast, many of the new data formats feature arrays; the

array's elements have an order. We may want to find an array element according to its order, or, we may want to find the positions of certain elements in their arrays.

## 7.1 `<Array> [<number>]`

Let's consider again the dataset `hr.employeesNest`.

```
 1  {
 2    'hr': {
 3        'employeesNest': <<
 4          {
 5            'id': 3,
 6            'name': 'Bob Smith',
 7            'title': null,
 8            'projects': [ { 'name': 'AWS Redshift Spectrum querying' },
 9                          { 'name': 'AWS Redshift security' },
10                          { 'name': 'AWS Aurora security' }
11                        ]
12          },
13          {
14              'id': 4,
15              'name': 'Susan Smith',
16              'title': 'Dev Mgr',
17              'projects': []
18          },
19          {
20              'id': 6,
21              'name': 'Jane Smith',
22              'title': 'Software Eng 2',
23              'projects': [ { 'name': 'AWS Redshift security' } ]
24          }
25        >>
26    }
27  }
```

The `projects` attribute is an array of tuples; that is, each tuple has an ordinal associated with it. The following query returns each employee name, along with the first project of the employee.

```
 1  SELECT e.name AS employeeName,
 2         e.projects[0].name AS firstProjectName
 3  FROM hr.employeesNest AS e
```

The query returns

```
 1  <<
 2    {
 3      'employeeName': 'Bob Smith',
 4      'firstProjectName': 'AWS Redshift Spectrum querying'
 5    },
 6    {
 7      'employeeName': 'Susan Smith'
 8    },
 9    {
10      'employeeName': 'Jane Smith',
11      'firstProjectName': 'AWS Redshift security'
12    }
13  >>
14  ---
15  OK! (51 ms)
```

## 7.2 Multistep Paths

Technically, the structure [`<number>`] is a kind of path step. For example, notice the 4-step path `e.projects`
`[0].name`. When `e` is bound to the first tuple of `hr.employeesNest`, then the path `e.projects` results into the
array

```
 1  [
 2      { 'name': 'AWS Redshift Spectrum querying' },
 3      { 'name': 'AWS Redshift security' },
 4      { 'name': 'AWS Aurora security' }
 5  ]
```

Consequently applying the `[0]` step on `e.projects` (that is, evaluating `e.projects[0]`) leads to {`'name':`
`'AWS Redshift Spectrum querying'`}. Finally, evaluating the `.name` step on `e.projects[0]` (that is, evaluating `e.projects[0].name`) leads to `'AWS Redshift Spectrum querying'`.

## 7.3 Finding the Order of Each Element in an Array

Let's assume that the order of each employee's projects in the `projects` attribute of `hr.employeesNest`
matters. The first project is the employee's highest priority project, followed by the second and so on. The
following query finds the names of each employee involved in a security project, the security project and
its order in the `projects` array.

```
 1  SELECT e.name AS employeeName,
```

```
2           p.name AS projectName,
3           o AS projectPriority
4    FROM hr.employeesNest AS e,
5         e.projects AS p AT o
6    WHERE p.name LIKE '%security%'
```

Notice the new feature: `AT o`. While `p` ranges over the elements of the array `e.projects`, the variable `o` takes as value the ordinal number of the element in the array. The query returns

```
1    <<
2      {
3        'employeeName': 'Bob Smith',
4        'projectName': 'AWS Redshift security',
5        'projectPriority': 1
6      },
7      {
8        'employeeName': 'Bob Smith',
9        'projectName': 'AWS Aurora security',
10       'projectPriority': 2
11     },
12     {
13       'employeeName': 'Jane Smith',
14       'projectName': 'AWS Redshift security',
15       'projectPriority': 0
16     }
17   >>
18   ---
19   OK! (12 ms)
```

# 8 Pivoting & Unpivoting

Many queries need to range over and collect the attribute name/value pairs of tuples or the key/value pairs of maps.

## 8.1 Unpivoting Tuples

Consider this dataset that provides the closing prices of multiple ticker symbols.

```
1    {
2        'closingPrices': <<
3            { 'date': '4/1/2019', 'amzn': 1900, 'goog': 1120, 'fb': 180 },
```

```
4              { 'date': '4/2/2019', 'amzn': 1902, 'goog': 1119, 'fb': 183 }
5      >>
6  }
```

The following query unpivots the stock ticker/price pairs.

```
1  SELECT c."date" AS "date",
2         sym AS "symbol",
3         price AS price
4  FROM closingPrices AS c,
5       UNPIVOT c AS price AT sym
6  WHERE NOT sym = 'date'
```

Notice the use of " in this query. The double quotes allow us to disambiquate from date the keyword and "date" the identifier. Also double quote specify case sensitivity for attribute lookups.

The query returns

```
1  <<
2    {
3      'date': '4/1/2019',
4      'symbol': 'amzn',
5      'price': 1900
6    },
7    {
8      'date': '4/1/2019',
9      'symbol': 'goog',
10     'price': 1120
11   },
12   {
13     'date': '4/1/2019',
14     'symbol': 'fb',
15     'price': 180
16   },
17   {
18     'date': '4/2/2019',
19     'symbol': 'amzn',
20     'price': 1902
21   },
22   {
23     'date': '4/2/2019',
24     'symbol': 'goog',
25     'price': 1119
26   },
```

```
27    {
28      'date': '4/2/2019',
29      'symbol': 'fb',
30      'price': 183
31    }
32  >>
33  ---
34  OK! (18 ms)
```

Unpivoting tuples enables the use of attribute names as if they were data. For example, it becomes easy to compute the average price for each symbol as

```
1  SELECT sym AS "symbol",
2         AVG(price) AS avgPrice
3  FROM closingPrices c,
4       UNPIVOT c AS price AT sym
5  WHERE NOT sym = 'date'
6  GROUP BY sym
```

which returns

```
1   <<
2     {
3       'symbol': 'amzn',
4       'avgPrice': 1901
5     },
6     {
7       'symbol': 'fb',
8       'avgPrice': 181.5
9     },
10    {
11      'symbol': 'goog',
12      'avgPrice': 1119.5
13    }
14  >>
15  ---
16  OK! (31 ms)
```

## 8.2  Pivoting into Tuples

Pivoting turns a collection into a tuple. For example, consider the collection

```
1  {
2      'todaysStockPrices': <<
3          { 'symbol': 'amzn', 'price': 1900},
4          { 'symbol': 'goog', 'price': 1120},
5          { 'symbol': 'fb', 'price': 180 }
6      >>
7  }
```

Then the following PIVOT query

```
1  PIVOT sp.price AT sp."symbol"
2  FROM todaysStockPrices sp
```

produces the tuple

```
1  {
2    'amzn': 1900,
3    'goog': 1120,
4    'fb': 180
5  }
6  ---
7  OK! (43 ms)
```

Notice that the PIVOT query looks like a SELECT-FROM-WHERE-... query except that instead of a SELECT clause it has a PIVOT <value expression> AT <attribute expression>. Note also that the PIVOT query does not return a singleton collection of tuples: Rather it literally returns a tuple value.

## 8.3 Use Case: Pivoting Subqueries

(This example also uses the grouping features of PartiQL, Creating Nested Results with GROUP BY ... GROUP AS.)

Let us generalize the previous case of pivoting. We have a table of stock prices

```
1  {
2      'stockPrices':<<
3          { 'date': '4/1/2019', 'symbol': 'amzn', 'price': 1900},
4          { 'date': '4/1/2019', 'symbol': 'goog', 'price': 1120},
5          { 'date': '4/1/2019', 'symbol': 'fb',    'price': 180 },
6          { 'date': '4/2/2019', 'symbol': 'amzn', 'price': 1902},
7          { 'date': '4/2/2019', 'symbol': 'goog', 'price': 1119},
8          { 'date': '4/2/2019', 'symbol': 'fb',    'price': 183 }
```

```
 9        >>
10  }
```

and we want to pivot it into a collection of tuples, where each tuple has all the `symbol`:`price` pairs for a date, as follows

```
 1  <<
 2  {
 3      'date': date(4/1/2019),
 4      'prices': {'amzn': 1900, 'goog': 1120, 'fb': 180}
 5  },
 6  {
 7      'date': date(4/2/2019),
 8      'prices': {'amzn': 1902, 'goog': 1119, 'fb': 183}
 9  }
10  >>
```

The following query first creates one group datesPrices for each date. Then the `PIVOT` subquery pivots the group into the tuple prices.

```
 1  SELECT sp."date" AS "date",
 2         (PIVOT dp.sp.price AT dp.sp."symbol"
 3          FROM datesPrices as dp ) AS prices
 4  FROM StockPrices AS sp GROUP BY sp."date" GROUP AS datesPrices
```

For example, the `datesPrices` collection, returned from `GROUP AS` for `sp.date = date(4/1/2019)` is

```
 1      'datesPrices': <<
 2       {
 3         'sp': {
 4           'date': '4/1/2019',
 5           'symbol': 'amzn',
 6           'price': 1900
 7         }
 8       },
 9       {
10         'sp': {
11           'date': '4/1/2019',
12           'symbol': 'goog',
13           'price': 1120
14         }
15       },
16       {
```

```
17          'sp': {
18            'date': '4/1/2019',
19            'symbol': 'fb',
20            'price': 180
21          }
22        }
23    >>
```

# 9 Creating Nested and Non-SQL Results

PartiQL allows queries that create nested results as well as queries that create heterogeneous results.

## 9.1 Creating Nested Results with `SELECT VALUE` Queries

Let's consider again the dataset `hr.employeesNestScalars`:

```
 1  {
 2      'hr': {
 3          'employeesNestScalars': <<
 4              {
 5                  'id': 3,
 6                  'name': 'Bob Smith',
 7                  'title': null,
 8                  'projects': [
 9                      'AWS Redshift Spectrum querying',
10                      'AWS Redshift security',
11                      'AWS Aurora security'
12                  ]
13              },
14              {
15                  'id': 4,
16                  'name': 'Susan Smith',
17                  'title': 'Dev Mgr',
18                  'projects': []
19              },
20              {
21                  'id': 6,
22                  'name': 'Jane Smith',
23                  'title': 'Software Eng 2',
24                  'projects': [ 'AWS Redshift security' ]
25              }
```

```
26           >>
27       }
28 }
```

The following query outputs each tuple of hr.employeesNestScalars, except that instead of all projects each tuple has only the security projects of the employee. The important new feature here is the SELECT VALUE <expression>.

```
1 SELECT e.id AS id,
2        e.name AS name,
3        e.title AS title,
4        ( SELECT VALUE p
5          FROM e.projects AS p
6          WHERE p LIKE '%security%'
7        ) AS securityProjects
8 FROM hr.employeesNestScalars AS e
```

The result is

```
 1 <<
 2   {
 3     'id': 3,
 4     'name': 'Bob Smith',
 5     'title': NULL,
 6     'securityProjects': <<
 7       'AWS Redshift security',
 8       'AWS Aurora security'
 9     >>
10   },
11   {
12     'id': 4,
13     'name': 'Susan Smith',
14     'title': 'Dev Mgr',
15     'securityProjects': <<>>
16   },
17   {
18     'id': 6,
19     'name': 'Jane Smith',
20     'title': 'Software Eng 2',
21     'securityProjects': <<
22       'AWS Redshift security'
23     >>
24   }
```

```
25  >>
26  ---
27  OK! (35 ms)
```

A `SELECT VALUE <expression>` query (or subquery, as in this example) returns a collection of whatever the `<expression>` evaluates to.

Notice the difference from SQL's `SELECT`, which always produces tuples. If a SQL `SELECT` appears as a subquery, then the context of the subquery designates whether the subquery's result should be coerced into a scalar (e.g., when `5 = <subquery>`), coerced into a collection of scalars (e.g., when `5 IN <subquery>`), etc. None of this applies to `SELECT VALUE`, which produces a collection and this collection is not coerced.

## 9.2 Creating Nested Results with `GROUP BY ... GROUP AS`

Another pattern of creating nested results in PartiQL is via the `GROUP AS` extension to SQL's `GROUP BY`. This pattern is more efficient and more intuitive than the use of nested `SELECT VALUE` queries when the required nesting is not following the nesting of the input. (The example in Creating Nested Results with `SELECT VALUE` Queries is one where the nesting in the output follows the nesting of the input and, thus, an intuitive solution does not involve `GROUP BY`.)

The following query outputs each security project found in `hr.employeesNestScalars` along with the list of employee names that work on the project.

```
1  SELECT p AS projectName,
2        ( SELECT VALUE v.e.name
3          FROM perProjectGroup AS v ) AS employees
4  FROM hr.employeesNestScalars AS e JOIN e.projects AS p ON p LIKE '%security%'
5  GROUP BY p GROUP AS perProjectGroup
```

The result is

```
1  <<
2    {
3      'projectName': 'AWS Aurora security',
4      'employees': <<
5        'Bob Smith'
6      >>
7    },
8    {
9      'projectName': 'AWS Redshift security',
10     'employees': <<
11       'Bob Smith',
```

```
12        'Jane Smith'
13      >>
14    }
15  >>
16  ---
17  OK! (24 ms)
```

The `GROUP AS` generalizes SQL's `GROUP BY` by making the formulated groups available in their entirety to the query's `SELECT` and `HAVING` clauses. Contrast with SQL's `GROUP BY`, where the `SELECT` and `HAVING` clauses can have aggregate functions over grouped columns but they cannot get access to the individual values of the grouped columns.

To better understand the workings of `GROUP BY ... GROUP AS` it is best to think of PartiQL queries as a pipeline of clauses, starting with the `FROM`, continuing with the `GROUP BY` and finishing with the `SELECT`. Each clause is a function that inputs data and outputs data. In that sense, the `GROUP BY ... GROUP AS` is a function that inputs the result of the `FROM` and outputs its result to the `SELECT`.

The following query (conceptually) produces the output of the `FROM` clause.

```
1  SELECT e AS e, p AS p
2  FROM hr.employeesNestScalars AS e JOIN e.projects AS p ON p LIKE '%security%'
```

We see that the `FROM` delivers the collection of tuples consisting of an employee `e` and a project `p` that were output by the `FROM` clause, i.e., the `LEFT JOIN`. This is alike SQL's `FROM` semantics.

Variable e                                          Variable p

```
1   { 'id': 3,
2
3    'name': 'Bob Smith',
4
5    'title': null,
6
7    'projects':  [ 'AWS Redshift
8    Spectrum querying',
9
10   'AWS Redshift security',
11
12   'AWS Aurora security'
13
14     ]
15
16    }
```

```
1       'AWS Redshift security'
```

```
1   { 'id': 3,
2
3    'name': 'Bob Smith',
4
5    'title': null,
6
7    'projects':  [ 'AWS Redshift
8    Spectrum querying',
9
10   'AWS Redshift security',
11
12   'AWS Aurora security'
13
14     ]
15
16    }
```

```
1       'AWS Aurora security'
```

| Variable e | Variable p |
|---|---|

```
 1   { 'id': 6,
 2
 3   'name': 'Jane Smith',
 4
 5   'title': 'Software Eng 2',
 6
 7   'projects':  [ 'AWS Redshift
 8   security'  ]
 9
10   }
```

```
 1   'AWS Redshift security'
```

Then the `GROUP BY ... GROUP AS ...` can be thought of as outputting a table that has one column for each group-by expression (i.e., each security project p) and a last column `perProjectGroup` whose value (conceptually) is the collection of employee/project e/p tuples that correspond to the group-by expression p. Thus the `GROUP BY ... GROUP AS ...` output is the table

| p | perProjectGroup |
|---|---|

```
 1   'AWS Redshift security'
```

```
 1   <<
 2
 3   { e: { 'id': 3, 'name': 'Bob
 4   Smith', ... }, p: 'AWS Redshift
 5   security' },
 6
 7   { e: { 'id': 6, 'name': 'Jane
 8   Smith', ... }, p: 'AWS Redshift
 9   security' }
10
11   >>
```

| p | perProjectGroup |
| --- | --- |
| 1  'AWS Aurora security' | ```
1  <<
2
3  { e: { 'id': 3, 'name': 'Bob
4  Smith', ...}, p: 'AWS Aurora
5  security' },
6
7  >>
``` |

Finally the SELECT clause inputs the above and outputs the query result.