# 1 Introduction

This document will provide a walk-through of the consolidate script with dummy data, aimed to provide clarity and understanding to how the script works.

At a high-level, this script reads an input csv file and constructs a compressed csv (and json) representation for it. In the compressed csv, each cell can represent a set of values rather than a singular value.

# 2 Data

Assume the input data that we want to consolidate is as follows:

| GuardiumVersion | DatabaseName | DatabaseVersion | OSName | OSVersion | Feature |
|---|---|---|---|---|---|
| 11.0 | Cassandra | Cassandra 3.11.10 | Ubuntu | Ubuntu 21.04 | A |
| 11.0 | Cassandra | Cassandra 3.11.10 | CentOS | CentOS 7 | A |
| 11.0 | Cassandra | Cassandra 3.11.10 | CentOS | CentOS 8 | A |
| 11.0 | Cassandra | Cassandra 3.11.10 | CentOS | CentOS 8 | A |
| 11.0 | Cassandra | Cassandra 3.11.10 | Windows | Windows 10 | A |
| 11.0 | Cassandra | Cassandra 4.0 | Windows | Windows 10 | B |
| 11.1 | DB2 | DB2 11.5.7 | SUSE | 15 | C |
| 11.1 | MariaDB | MariaDB 10.5 | macOS | macOS 11.0 | C |
| 11.1 | MariaDB | MariaDB 10.6 | Ubuntu | Ubuntu 22.04 | C |
| 11.1 | MongoDB | MongoDB 4.0 | CentOS | CentOS 6 | B |
| 11.1 | MongoDB | MongoDB 4.0 | CentOS | CentOS 7 | B |
| 11.1 | MongoDB | MongoDB 4.0 | CentOS | CentOS 8 | B |
| 11.2 | MongoDB | MongoDB 4.2 | CentOS | CentOS 6 | B |
| 11.2 | MongoDB | MongoDB 4.2 | CentOS | CentOS 7 | B |
| 11.2 | MongoDB | MongoDB 4.2 | CentOS | CentOS 8 | B |
| 11.2 | MongoDB | MongoDB 4.4 | CentOS | CentOS 6 | B |
| 11.2 | MongoDB | MongoDB 4.4 | CentOS | CentOS 7 | B |
| 11.2 | MongoDB | MongoDB 4.4 | CentOS | CentOS 8 | B |
| 11.3 | MongoDB | MongoDB 4.6 | CentOS | CentOS 8 | B |
| 11.4 | MongoDB | MongoDB 4.6 | Windows | Windows 10 | B |
| 11.3 | Redis | Redis 7.2 | Debian | Debian 11 | A |

Table 1: Input Data From CSV

This script will compress that data to find this representation

| GuardiumVersion | DatabaseName | DatabaseVersion | OSName | OSVersion | Feature |
|---|---|---|---|---|---|
| 11.0 | Cassandra | Cassandra 3.11.10 | CentOS | CentOS 7, CentOS 8 | A |
| 11.0 | Cassandra | Cassandra 3.11.10 | Ubuntu | Ubuntu 21.04 | A |
| 11.0 | Cassandra | Cassandra 3.11.10 | Windows | Windows 10 | A |
| 11.0 | Cassandra | Cassandra 4.0 | Windows | Windows 10 | B |
| 11.1 | DB2 | DB2 11.5.7 | SUSE | 15 | C |
| 11.1 | MariaDB | MariaDB 10.5 | macOS | macOS 11.0 | C |
| 11.1 | MariaDB | MariaDB 10.6 | Ubuntu | Ubuntu 22.04 | C |
| 11.1 | MongoDB | MongoDB 4.0 | CentOS | CentOS 6, CentOS 7, CentOS 8 | B |
| 11.2 | MongoDB | MongoDB 4.2, MongoDB 4.4 | CentOS | CentOS 6, CentOS 7, CentOS 8 | B |
| 11.3 | MongoDB | MongoDB 4.6 | CentOS | CentOS 8 | B |
| 11.3 | Redis | Redis 7.2 | Debian | Debian 11 | A |
| 11.4 | MongoDB | MongoDB 4.6 | Windows | Windows 10 | |

Table 2: Compression version of Table 1

# WalkThru Python Code

We will start our walk-through when we call the consolidate function in main() in runner.py

```
consolidate(output_json_path, output_csv_path, input_csv_path,
                version_headers, full_key, feature_headers,
                partition_header_number, logger)
# full_key = ["Guardium Version","DatabaseName","DatabaseVersion","OSName","OSVersion","Feature"]
# version_headers = ["Guardium Version","DatabaseName","DatabaseVersion","OSName","OSVersion"]
# feature_headers = ["Feature"]
# partition_header_number = 1
```

Here, the arguments *output_json_path, output_csv_path, input_csv_path* are just paths to where we're reading the input data and writing the compressed data.

The argument *full_key* is just all the header names. Whereas, *version_headers* represent those header names that represent columns whose values are related to the version of something. On the other hand, *feature_headers* are the header names related to a feature.

*partition_header_number* is the index of the column based on which you want to partition the data. Essentially, during the compression, you will split your input data rows based on the value in column at index *partition_header_number*. This will make more sense further on.

```python
# ------------------- Inside consolidate function -------------------
    # Instantiate output variables
    output_csv = []
    output_csv.append(full_key) # Add headers to first line of csv output
    output_json = {}

    # Get unique values in a column (eg.unique database names in database column)
    # (This will be used to partition the tabular data
    # before compression for efficiency)
    uniq_vals = read_csv_get_unique_vals_in_column(input_csv_path,partition_header_number)
    logger.debug("Stored all unique vals of column #:%s ", partition_header_number)

    # uniq_vals = ["MongoDB","Cassandra","DB2","MariaDB","Redis"]
```

*uniq_vals* is a set of all the unique values in the column at index = *partition_header_number*

```python
    # Loop thru each uniq value
    for uniq_val in uniq_vals:
        output_json[uniq_val] = [] # Will store consolidated data for specific uniq_val

        # Partition data (Get subsection of all rows with uniq_val in the partition_header_number)
        partitioned_data = read_csv_for_uniq_val(input_csv_path,uniq_val,partition_header_number)
        partitioned_data = remove_duplicates_2d(partitioned_data)
        logger.debug("Partitioned Data for %s",uniq_val)

        # uniq_val = "MongoDB"
        # partitioned_data =
        # [['11.1', 'MongoDB', 'MongoDB 4.0', 'CentOS', 'CentOS 6', 'B'], ['11.1', 'MongoDB', 'MongoDB 4.0', '
            CentOS', 'CentOS 7', 'B'], ['11.1', 'MongoDB', 'MongoDB 4.0', 'CentOS', 'CentOS 8', 'B'], ['11.2', '
            MongoDB', 'MongoDB 4.2', 'CentOS', 'CentOS 6', 'B'], ['11.2', 'MongoDB', 'MongoDB 4.2', 'CentOS', '
            CentOS 7', 'B'], ['11.2', 'MongoDB', 'MongoDB 4.2', 'CentOS', 'CentOS 8', 'B'], ['11.2', 'MongoDB', '
            MongoDB 4.4', 'CentOS', 'CentOS 6', 'B'], ['11.2', 'MongoDB', 'MongoDB 4.4', 'CentOS', 'CentOS 7', 'B
            '], ['11.2', 'MongoDB', 'MongoDB 4.4', 'CentOS', 'CentOS 8', 'B'], ['11.3', 'MongoDB', 'MongoDB 4.6',
             'CentOS', 'CentOS 8', 'B'], ['11.4', 'MongoDB', 'MongoDB 4.6', 'Windows', 'Windows 10', 'B']]
```

We now use each *uniq_val* to get *partitioned_data*, which is the subsection of all rows with *uniq_val* in the *partition_header_number* column index. We continue in this for-loop with *uniq_val*="MongoDB"

```python
        # Group rows of data based on their feature values,
        # Returns Dict[str, List[List[str]]]
        # A dict with a string key (representing one rows' feature values) and,
        # Lists of lists of strings (representing version values from multiple rows) as values
        # Therefore, each list of strings
        # concatenated with the corresponding key, recreates the original row
        # CONSTRAINT: version values must come first
        grouped_partitioned_data = group_data_by_feature_set(partitioned_data
                            ,get_features=lambda x: "|+|".join(x[len(version_key):])
                            # Combine feature values into a unique string identifier
                            ,get_versions=lambda x:x[0:len(version_key)]
                            ,logger=logger) # Extract header key values
        # The 'get_features' function combines the feature  values of each row
        # into a single string, acting as a unique identifier for a grouping.
        # The 'get_versions' function extracts the version values from each row.

        # grouped_partitioned_data = { 'B':
        # [['11.1', 'MongoDB', 'MongoDB 4.0', 'CentOS', 'CentOS 6'],
        # ['11.1', 'MongoDB', 'MongoDB 4.0', 'CentOS', 'CentOS 7'],
        # ['11.1', 'MongoDB', 'MongoDB 4.0', 'CentOS', 'CentOS 8'],
        # ['11.2', 'MongoDB', 'MongoDB 4.2', 'CentOS', 'CentOS 6'],
        # ['11.2', 'MongoDB', 'MongoDB 4.2', 'CentOS', 'CentOS 7'],
        # ['11.2', 'MongoDB', 'MongoDB 4.2', 'CentOS', 'CentOS 8'],
        # ['11.2', 'MongoDB', 'MongoDB 4.4', 'CentOS', 'CentOS 6'],
        # ['11.2', 'MongoDB', 'MongoDB 4.4', 'CentOS', 'CentOS 7'],
        # ['11.2', 'MongoDB', 'MongoDB 4.4', 'CentOS', 'CentOS 8'],
        # ['11.3', 'MongoDB', 'MongoDB 4.6', 'CentOS', 'CentOS 8'],
        # ['11.4', 'MongoDB', 'MongoDB 4.6', 'Windows', 'Windows 10']]
        # }
```

We will group *partitioned_data* based on values in the *feature_headers* columns. Since, the rows in *partitioned_data* have the same value for column under "Feature", *grouped_partitioned_data* will only have one key.value pair.

```python
            # Loop thru grouped_partitioned_data dictionary
            # Each key represents one rows' feature value(s)
            # Each value represents version values from multiple rows
            for set_of_feature,version_data in grouped_partitioned_data.items():

                # Returns a consolidate/compressed representation
                # of the version values from multiple rows
                consolidated_version_rows = cartesian_decomposition(
                                                    version_data,
                                                    version_key,logger,
                                                    input_csv_path, partitioned_data)

            # version_data =
            # [['11.1', 'MongoDB', 'MongoDB 4.0', 'CentOS', 'CentOS 6'],
            # ['11.1', 'MongoDB', 'MongoDB 4.0', 'CentOS', 'CentOS 7'],
            # ['11.1', 'MongoDB', 'MongoDB 4.0', 'CentOS', 'CentOS 8'],
            # ['11.2', 'MongoDB', 'MongoDB 4.2', 'CentOS', 'CentOS 6'],
            # ['11.2', 'MongoDB', 'MongoDB 4.2', 'CentOS', 'CentOS 7'],
            # ['11.2', 'MongoDB', 'MongoDB 4.2', 'CentOS', 'CentOS 8'],
            # ['11.2', 'MongoDB', 'MongoDB 4.4', 'CentOS', 'CentOS 6'],
            # ['11.2', 'MongoDB', 'MongoDB 4.4', 'CentOS', 'CentOS 7'],
            # ['11.2', 'MongoDB', 'MongoDB 4.4', 'CentOS', 'CentOS 8'],
            # ['11.3', 'MongoDB', 'MongoDB 4.6', 'CentOS', 'CentOS 8'],
            # ['11.4', 'MongoDB', 'MongoDB 4.6', 'Windows', 'Windows 10']]
```

```python
# -------------------- Inside cartesian_decomposition function --------------------
    # Getting all uniq vals per each column from partitioned data
    uniq_column_vals_part_data = get_uniq_vals_for_each_column(key_,partitioned_data)
    logger.debug("Stored all uniq vals in each column from part. data")
    # Returns a dict, where key is the column name,
    # and value is a list of unique vals in that column

    # Getting all uniq vals per each column from version data
    uniq_column_vals_version_data = get_uniq_vals_for_each_column(key_,version_data)
    logger.debug("Stored all uniq vals in each column from version data")

    # uniq_column_vals_part_data = {
    # 'GuardiumVersion': ['11.1', '11.2', '11.3', '11.4'],
    # 'DatabaseName': ['MongoDB'],
    # 'DatabaseVersion': ['MongoDB 4.0', 'MongoDB 4.2', 'MongoDB 4.4', 'MongoDB 4.6'],
    # 'OSName': ['CentOS', 'Windows'],
    # 'OSVersion': ['CentOS 6', 'CentOS 7', 'CentOS 8', 'Windows 10']}
```

Get all unique values in each column in a dictionary.

```python
    # Generate all possible relevant ranges (ordered subsets)
    # per each column using uniq vals (Can blow up computationally)
    # eg. find_ranges([1,2,3]) = [[1],[2],[3],[1,2],[2,3],[1,2,3]]

    all_ranges = []
    for x in key_:
        ranges = find_relevant_ranges(list(uniq_column_vals_version_data[x])
                                ,list(uniq_column_vals_part_data[x]))
        all_ranges.append(ranges)
    logger.debug("Generated all possible ordered set of uniq vals of each column")


    # all_ranges = [
    # [['11.1'], ['11.1', '11.2'], ['11.1', '11.2', '11.3'], ['11.1', '11.2', '11.3', '11.4'], ['11.2'], ['11.2',
        '11.3'], ['11.2', '11.3', '11.4'], ['11.3'], ['11.3', '11.4'], ['11.4']],
    # [['MongoDB']],
    # [['CentOS'], ['CentOS', 'Windows'], ['Windows']],
    # [['CentOS 6'], ['CentOS 6', 'CentOS 7'], ['CentOS 6', 'CentOS 7', 'CentOS 8'],
    # ['CentOS 6', 'CentOS 7', 'CentOS 8', 'Windows 10'], ['CentOS 7'], ['CentOS
    # 7', 'CentOS 8'], ['CentOS 7', 'CentOS 8', 'Windows 10'], ['CentOS 8'],
    # ['CentOS 8', 'Windows 10'], ['Windows 10']]
    # ]
```

Gets each possible range for each column.

```python
    # Generate all possible combinations using ranges from each column
    # Take cartesian product of each list of ranges
    # combinations([ [[1],[2],[3],[1,2],[2,3],[1,2,3]] , [["a"],["b"],["a","b"]] ]) =
    #     [
    #         [[1], ['a']],        [[1], ['b']],        [[1], ['a', 'b']],
    #         [[2], ['a']],        [[2], ['b']],        [[2], ['a', 'b']],
```

```
7    #                [[3], ['a']],            [[3], ['b']],            [[3], ['a', 'b']],
8    #                [[1, 2], ['a']],      [[1, 2], ['b']],      [[1, 2], ['a', 'b']],
9    #                [[2, 3], ['a']],      [[2, 3], ['b']],      [[2, 3], ['a', 'b']],
10   #                [[1, 2, 3], ['a']], [[1, 2, 3], ['b']], [[1, 2, 3], ['a', 'b']]
11   #                                                                                ]
12   # (Can blow up computationally)
13   # Save each combination as a Combination Object
14   combinations_list = []
15   y = combinations(all_ranges)
16   for _,x in enumerate(y):
17       combinations_list.append(Combination(x))
18   logger.debug("Generated all possible combinations of uniq val ordered sets of each column")
19   # y = [
20   # [['11.1'], ['MongoDB'], ['MongoDB 4.0'], ['CentOS'], ['CentOS 6']],
21   # [['11.1'], ['MongoDB'], ['MongoDB 4.0'], ['CentOS'], ['CentOS 6', 'CentOS 7']],
22   # [['11.1'], ['MongoDB'], ['MongoDB 4.0'], ['CentOS'], ['CentOS 6', 'CentOS 7', 'CentOS 8', 'Windows 10']],
23   # ....,
24   # [['11.4'], ['MongoDB'], ['MongoDB 4.6'], ['Windows'], ['Windows 10']]
25   # ]
26   # len(y) = 3000
```

Using each possible range from each column, generate all possible combinations.

```
1    # Check if each row of version_data is compatible with any possible combination
2    for _,row in enumerate(version_data):
3        for _,combo in enumerate(combinations_list):
4            # Check if row is compatible with combo (Cartesian product of combo includes row)
5            # i.e row [1,a] is compatible with combo([[1], ['a', 'b']])
6            # but row [1,a] is not compatible with combo([[2, 3], ['a', 'b']])
7            if combo.combo_allows_row(row):
8                combo.add_row(row) # Store all compatible rows for each combo
9    logger.debug("Referenced all data rows with all possible combinations")
10
11   # Sort combination from higher to lower capacity
12   combinations_list = sorted(combinations_list, key=lambda x: -x.capacity)
13   # combinations_list = [
14   # Combination(length=5, key=[['11.1', '11.2', '11.3', '11.4'], ['MongoDB'],
15   # ['MongoDB 4.0', 'MongoDB 4.2', 'MongoDB 4.4', 'MongoDB 4.6'], ['CentOS',
16   # 'Windows'], ['CentOS 6', 'CentOS 7', 'CentOS 8', 'Windows 10']],
17   # capacity=128, flow=11, rows=[....]),
18   # Combination(length=5, key=[['11.1', '11.2', '11.3'], ['MongoDB'], ['MongoDB # 4.0', 'MongoDB 4.2', 'MongoDB
          4.4', 'MongoDB 4.6'], ['CentOS', 'Windows'], # ['CentOS 6', 'CentOS 7', 'CentOS 8', 'Windows 10']],
          capacity=96, flow=10, # rows=[...]),
19   ...,
20   # Combination(length=5, key=[['11.4'], ['MongoDB'], ['MongoDB 4.6'], ['Windows'], ['Windows 10']], capacity
          =1, flow=1, rows=[...])
21   ]
```

Check each *row* in *version_data* against each of the 3000 *combo* in *combinations_list*. For example, the *row* = ['11.1', 'MongoDB', 'MongoDB 4.0', 'CentOS', 'CentOS 6'] is compatible with *Combination(key=[['11.1'], ['MongoDB'], ['MongoDB 4.0'], ['CentOS'], ['CentOS 6']])* and *Combination(key=[['11.1', '11.2', '11.3', '11.4'], ['MongoDB'], ['MongoDB 4.0', 'MongoDB 4.2', 'MongoDB 4.4', 'MongoDB 4.6'], ['CentOS', 'Windows'], ['CentOS 6', 'CentOS 7', 'CentOS 8', 'Windows 10']])*. But, is not compatible with *Combination(key=[['11.1', '11.2', '11.3', '11.4'], ['MongoDB'], ['MongoDB 4.2', 'MongoDB 4.4', 'MongoDB 4.6'], ['CentOS', 'Windows'], ['CentOS 6', 'CentOS 7', 'CentOS 8', 'Windows 10']])*. We will track these compatible rows in the *row* property of the *Combination* object.

```
1    # Make copy of version_data
2    version_data_copy = version_data
3
4    # This will represent those combinations which represent exactly all rows of data
5    final_combos = []
6
7    # Loop thru each full capped combination and remove those lines from version_data_copy,
8    # that can be represented by that combination
9    # Continue till you have a list of combination that can
10   # represent all lines/rows in version_data
11   while len(version_data_copy) != 0:
12       for _,combo in enumerate(combinations_list):
13           if combo.is_full_cap(): # A combo is full_cap if len(combo.rows) == combo.capacity
14
15               # if all combo.rows are present in version_data_copy, remove these rows
16               # if not all combo.rows are present in version_data_copy, keep unchanged
17               version_data_copy_filtered = remove_if_all_present(version_data_copy, combo.rows)
18
19               if len(version_data_copy_filtered) != len(version_data_copy):
20                   # if version_data_copy_filtered was changed
21                   version_data_copy = version_data_copy_filtered
```

```
22                      combinations_list.remove(combo)
23                      #Check for subcombos
24                      final_combos.append(combo)
25
26      #
27      return [final_combo.key for final_combo in final_combos]
28
29 # returns [
30 # [['11.2'], ['MongoDB'], ['MongoDB 4.2', 'MongoDB 4.4'], ['CentOS'], ['CentOS 6', 'CentOS 7', 'CentOS 8']]
31 # [['11.1'], ['MongoDB'], ['MongoDB 4.0'], ['CentOS'], ['CentOS 6', 'CentOS 7', 'CentOS 8']]
32 # [['11.3'], ['MongoDB'], ['MongoDB 4.6'], ['CentOS'], ['CentOS 8']]
33 # [['11.4'], ['MongoDB'], ['MongoDB 4.6'], ['Windows'], ['Windows 10']]
34 # ]
35 # version_data =
36 # [['11.1', 'MongoDB', 'MongoDB 4.0', 'CentOS', 'CentOS 6'],
37 # ['11.1', 'MongoDB', 'MongoDB 4.0', 'CentOS', 'CentOS 7'],
38 # ['11.1', 'MongoDB', 'MongoDB 4.0', 'CentOS', 'CentOS 8'],
39 # ['11.2', 'MongoDB', 'MongoDB 4.2', 'CentOS', 'CentOS 6'],
40 # ['11.2', 'MongoDB', 'MongoDB 4.2', 'CentOS', 'CentOS 7'],
41 # ['11.2', 'MongoDB', 'MongoDB 4.2', 'CentOS', 'CentOS 8'],
42 # ['11.2', 'MongoDB', 'MongoDB 4.4', 'CentOS', 'CentOS 6'],
43 # ['11.2', 'MongoDB', 'MongoDB 4.4', 'CentOS', 'CentOS 7'],
44 # ['11.2', 'MongoDB', 'MongoDB 4.4', 'CentOS', 'CentOS 8'],
45 # ['11.3', 'MongoDB', 'MongoDB 4.6', 'CentOS', 'CentOS 8'],
46 # ['11.4', 'MongoDB', 'MongoDB 4.6', 'Windows', 'Windows 10']]
```

Here we loop through each *combo* and a create a set of them, such that they represent exactly all rows of *version_data*

We continue this for all *uniq_vals* to create Table 2.