# SkiRaff an ETL Testing Framework for pygrametl

June 16, 2016

Alexander Branborg `abran13@student.aau.dk`
Arash Michael Aami Kjær `ams13@student.aau.dk`
Mathias Claus Jensen `mcje13@student.aau.dk`
Mikael Vind Mikkelsen `mvmi12@student.aau.dk`

Department of Computer Science
Aalborg University
Denmark

**AALBORG UNIVERSITY**
DENMARK

Predicates
Why are they useful?
Usage/Implementation
Alternative Implementation

- ▶ Why are they useful?
- ▶ Usage/Implementation
- ▶ Alternative Forms of Implementation

- ▶ Systems level testing
  - ▶ Data loss

► Systems level testing
  ► Data loss
► Source to target test



►

- ▶ Systems level testing
  - ▶ Data loss
- ▶ Source to target test

```
┌────────┐
│Source 1│──┐
└────────┘  │
┌────────┐  │   ╱‾‾‾‾╲   ┌────┐   ╱‾‾‾‾╲   ┌──────┐
│Source 2│──┼─▶( ETL )─▶│ DW │─▶( Test )─▶│Report│
└────────┘  │   ╲____╱   └────┘   ╲____╱   └──────┘
┌────────┐  │
│Source n│──┘
└────────┘
```

- ▶
- ▶ Regression testing
- ▶ Business Rules

## Predicates available in SKiRaff

- ▶ RowCountPredicate
- ▶ ColumnNotNullPredicate
- ▶ ReferentialIntegrityPredicate
- ▶ FunctionalDependencyPredicate
- ▶ SCDVersionPredicate
- ▶ CompareTablePredicate
- ▶ RuleRowPredicate
- ▶ RuleColumnPredicate

## Predicates available in SKiRaff

- ► RowCountPredicate
- ► ColumnNotNullPredicate
- ► **ReferentialIntegrityPredicate**
- ► **FunctionalDependencyPredicate**
- ► SCDVersionPredicate
- ► CompareTablePredicate
- ► **RuleRowPredicate**
- ► RuleColumnPredicate

**6** Functional Dependency - Why is it useful?

► A, B –> C

## Functional Dependency - Why is it useful?

- ► A, B –> C
- ► DW holds certain hierarchical properties

## Setup:

```
1  FunctionalDependencyPredicate(table_name=['CountryDim','
       AuthorDim'], alpha='city', beta='country')
```

## SQL querie:

```
1  SELECT DISTINCT t1.country, t2.city
2  FROM countrydim NATURAL JOIN authordim AS t1, countrydim
       NATURAL JOIN authordim AS t2
3  WHERE t1.city = t2.city
4  AND t1.country <> t2.country
```

SkiRaff an ETL Testing
Framework for
pygrametl

Alexander Branborg,
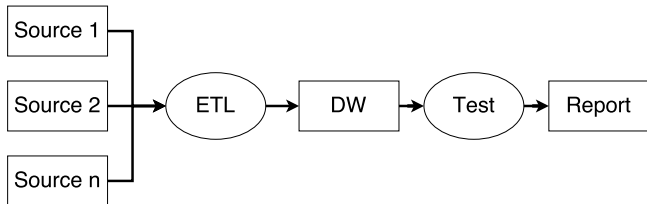Arash Michael Aami
Kjær,
Mathias Claus Jensen,
Mikael Vind Mikkelsen

Predicates
Why are they useful?
Usage/Implementation
Alternative Implementation

8

```python
1  # Creates part of select statement to get keys
2  select_alpha = ["t1." + str(a) for a in self.alpha]
3  select_beta = ["t2." + str(b) for b in self.beta]
4  select_sql = select_alpha + select_beta
5
6  # SQL setup for the left side of the dependency in WHERE—
       clause
7  alpha_sql_generator = (" t1.{} = t2.{} ".format(a, a)
8                          for a in self.alpha)
9  and_alpha = ' AND '.join(alpha_sql_generator)
10
11  # SQL setup for the right side of the dependency in WHERE—
       clause
12  beta_sql_generator = (" (t1.{} <> t2.{}) ".format(b, b)
13                         for b in self.beta)
14  or_beta = ' OR '.join(beta_sql_generator)
```

Department of Computer
Science
Aalborg University
Denmark

22

SkiRaff an ETL Testing
Framework for
pygrametl

Alexander Branborg,
Arash Michael Aami
Kjær,
Mathias Claus Jensen,
Mikael Vind Mikkelsen

Predicates
Why are they useful?
Usage/Implementation (9)
Alternative Implementation

```python
# Final setup of the entire SQL command
lookup_sql = "SELECT DISTINCT " + ','.join(select_sql) + \
             " FROM " + \
             " (" + " NATURAL JOIN ".join(self.table_name
                 ) + ") " + \
             " AS t1 ," + \
             " (" + " NATURAL JOIN ".join(self.table_name
                 ) + ") " + \
             " AS t2 " + \
             " WHERE " + and_alpha + " AND " + or_beta
```

## SQL querie:

```sql
SELECT DISTINCT t1.country, t2.city
FROM countrydim NATURAL JOIN authordim AS t1, countrydim
    NATURAL JOIN authordim AS t2
WHERE t1.city = t2.city
AND t1.country <> t2.country
```

SkiRaff an ETL Testing
Framework for
pygrametl

Alexander Branborg,
Arash Michael Aami
Kjær,
Mathias Claus Jensen,
Mikael Vind Mikkelsen

Predicates
Why are they useful?
Usage/Implementation
Alternative Implementation

```
1  cursor = dw_rep.connection.cursor()
2  cursor.execute(lookup_sql)
3  query_result = cursor.fetchall()
4  cursor.close()
5
6  # Create dict, so that attributes have names
7  names = [t[0] for t in cursor.description]
8  dict_result = []
9  for row in query_result:
10     dict_result.append(dict(zip(names, row)))
11
12 # If any rows were fetched. Assertion fails
13 if not dict_result:
14     self.__result__ = True
```

Department of Computer
Science
Aalborg University
Denmark

SkiRaff an ETL Testing
Framework for
pygrametl

Alexander Branborg,
Arash Michael Aami
Kjær,
Mathias Claus Jensen,
Mikael Vind Mikkelsen

Predicates

Why are they useful?

Usage/Implementation

Alternative Implementation

## Referential Integrity - Why is it useful?

- ▶ Most DBMS's have various referential integrity rules

Department of Computer
Science
Aalborg University
Denmark

(11) Referential Integrity - Why is it useful?

- ▶ Most DBMS's have various referential integrity rules
- ▶ Not removing the correct data from all tables

## Setup:

```
1  ReferentialIntegrityPredicate (
2      refs ={ 'FactTable': ('BookDim', 'AuthorDim'),
3              'AuthorDim': ('CountryDim')},
4      points_to_all = True ,
5      all_pointed_to = True
6  )
```

## SQL querie:

```
1  SELECT *
2  FROM facttable
3  WHERE NOT EXISTS (
4      SELECT NULL FROM author_dim
5      WHERE facttable.aid = author_dim.aid
6      )
```

SkiRaff an ETL Testing
Framework for
pygrametl

Alexander Branborg,
Arash Michael Aami
Kjær,
Mathias Claus Jensen,
Mikael Vind Mikkelsen

Predicates

Why are they useful?

Usage/Implementation (13)

Alternative Implementation

```python
1   missing_keys = []
2
3       # Maps table names to table_representations
4       refs = {}
5       for alpha, beta in self.refs.items():
6           b = []
7           if isinstance(alpha, str):
8                   a = dw_rep.get_data_representation(alpha)
9           else:
10              raise ValueError('Expected string in refs, got
                 : ' +
11                                  str(type(x)))
12          if isinstance(beta, str):
13              b.append(dw_rep.get_data_representation(beta))
```

SkiRaff an ETL Testing
Framework for
pygrametl

Alexander Branborg,
Arash Michael Aami
Kjær,
Mathias Claus Jensen,
Mikael Vind Mikkelsen

Predicates
Why are they useful?
Usage/Implementation (14)
Alternative Implementation

```
1              else:
2                  for x in beta:
3                      if isinstance(x, str):
4                          b.append(dw_rep.
5                                      get_data_representation(x
                                        ))
6                      else:
7                          raise ValueError('Expected string' + '
                                in refs, got:' + str(type(x)))
8          refs[a] = tuple(b)
9      self.refs = refs
```

SkiRaff an ETL Testing
Framework for
pygrametl

Alexander Branborg,
Arash Michael Aami
Kjær,
Mathias Claus Jensen,
Mikael Vind Mikkelsen

Predicates
Why are they useful?
Usage/Implementation
Alternative Implementation

15

```python
1  # If references not given. We check refs between all
       tables.
2  if not self.refs:
3      self.refs = dw_rep.refs
4
5  # Performs check for each pair of main table and foreign
       key table.
6  for table, dims in self.refs.items():
7      for dim in dims:
8          key = dim.key
9
10         # Check that each entry in main table has match
11         if self.points_to_all:
12             query_result = referential_check(table, dim,
                   key, dw_rep)
13
14             if query_result:
15                 for row in query_result:
16                     msg = '{}: {} in {} not found in {}' \
17                         .format(key, row[0], table.name,
                             dim.name)
18                     missing_keys.append(msg)
```

Department of Computer
Science
Aalborg University
Denmark

22

SkiRaff an ETL Testing
Framework for
pygrametl

Alexander Branborg,
Arash Michael Aami
Kjær,
Mathias Claus Jensen,
Mikael Vind Mikkelsen

Predicates

Why are they useful?

Usage/Implementation (16)

Alternative Implementation

```python
1           # Check that each entry in foreign key table has
                match
2           if self.all_pointed_to:
3               query_result = referential_check(dim, table,
                    key, dw_rep)
4
5               if query_result:
6                   for row in query_result:
7                       msg = '{}:_{}_in_{}_not_found_in_{}' \
8                           .format(key, row[0], dim.name,
                                table.name)
9                       missing_keys.append(msg)
10
11  if not missing_keys:
12      self.__result__ = True
```

SkiRaff an ETL Testing
Framework for
pygrametl

Alexander Branborg,
Arash Michael Aami
Kjær,
Mathias Claus Jensen,
Mikael Vind Mikkelsen

Predicates
Why are they useful?
Usage/Implementation
Alternative Implementation

## RuleRowPredicate - Why is it useful?

▶ Gives the user freedom to check for things our other predicate can't

▶ But with an easy setup

SkiRaff an ETL Testing
Framework for
pygrametl

Alexander Branborg,
Arash Michael Aami
Kjær,
Mathias Claus Jensen,
Mikael Vind Mikkelsen

Predicates

Why are they useful?
Usage/Implementation (17)
Alternative Implementation

## RuleRowPredicate - Why is it useful?

- ▶ Gives the user freedom to check for things our other predicate can't
- ▶ But with an easy setup
- ▶ However slower than others due to the lack of SQL implementation

SkiRaff an ETL Testing
Framework for
pygrametl

Alexander Branborg,
Arash Michael Aami
Kjær,
Mathias Claus Jensen,
Mikael Vind Mikkelsen

Predicates

Why are they useful?

Usage/Implementation

Alternative Implementation

## Setup:

```
1  def no_autobios(name, title):
2      return not name == title
3
4  RuleRowPredicate(table_name=['AuthorDim','FactTable','
       BookDim']
5                   constraint_function=no_autobios,
6                   column_names=['name', 'title'],
7                   constraint_args=[],
8                   column_names_exclude=False)
```

SkiRaff an ETL Testing
Framework for
pygrametl

Alexander Branborg,
Arash Michael Aami
Kjær,
Mathias Claus Jensen,
Mikael Vind Mikkelsen

Predicates
Why are they useful?
Usage/Implementation      (19)
Alternative Implementation

```python
1  # Gets the attribute names for columns needed for test
2  column_arg_names = self.setup_columns(dw_rep, self.
       table_name, self.column_names, self.
       column_names_exclude)
3
4  func_args = inspect.getargspec(self.constraint_function).
       args
5  if len(func_args) != len(column_arg_names) + len(self.
       constraint_args):
6      raise ValueError("""Number of columns and number of
           arguments do not match""")
```

SkiRaff an ETL Testing
Framework for
pygrametl

Alexander Branborg,
Arash Michael Aami
Kjær,
Mathias Claus Jensen,
Mikael Vind Mikkelsen

Predicates
Why are they useful?
Usage/Implementation
Alternative Implementation

20

```python
1   # Iterates over each row, calling the constraint function
        upon it
2   for row in dw_rep.iter_join(self.table_name):
3
4       # Finds parameters. First attributes then additional
            params.
5       arguments = []
6       for name in column_arg_names:
7           arguments.append(row[name])
8
9       if self.constraint_args:
10          arguments.append(*self.constraint_args)
11
12      # Runs function on parameters
13      if not self.constraint_function(*arguments):
14          wrong_rows.append(row)
15
16  if not wrong_rows:
17      self.__result__ = True
```

SkiRaff an ETL Testing
Framework for
pygrametl

Alexander Branborg,
Arash Michael Aami
Kjær,
Mathias Claus Jensen,
Mikael Vind Mikkelsen

Predicates
Why are they useful?
Usage/Implementation
Alternative Implementation (21)

## Now: SQL queries

```
25      def run(self, dw_rep):
26          pred_sql = \
27              " SELECT COUNT(*) " + \
28              " FROM " + "NATURAL JOIN".join(self.
                    table_name)
29
30          cursor = dw_rep.connection.cursor()
31          cursor.execute(pred_sql)
32          query_result = cursor.fetchall()
33          cursor.close()
34
35          if query_result[0] == self.number_of_rows:
36              self.__result__ = True
```

SkiRaff an ETL Testing
Framework for
pygrametl

Alexander Branborg,
Arash Michael Aami
Kjær,
Mathias Claus Jensen,
Mikael Vind Mikkelsen

Predicates

Why are they useful?

Usage/Implementation

Alternative Implementation 22

## Alternative: Representation objects in python

```python
21      def run ( self , dw_rep ) :
22          self . row_number = 0
23          self . table = []
24
25          for row in dw_rep . get_data_representation ( self .
                 table_name ) :
26              self . table . append ( row )
27              self . row_number += 1
28
29          if len ( self . table ) == self . number_of_rows :
30              self . __result__ = True
31          else :
32              self . __result__ = False
```

Thank you for using this theme!

AALBORG UNIVERSITY
DENMARK