

Regression Testing

of database applications under an incremental software development settings

By: RAÚL H. ROSERO^{1,2}, OMAR S. GÓMEZ², AND GLEN RODRÍGUEZ

Presenter: Aryan Ebrahimpour

Professor: Saeed Zahedi



What is regression testing?

- 🔬 verifies previous features on a software product when it is modified or new features are added to it
- 🔬 a costly process

Why not retest all?

- ⚠ Executing the total set of test cases is an expensive and laborious endeavor
- ⚠ Therefore determining a subset of test cases to be executed in a software regression test is an active topic of research

Different approaches

- ♟ **Minimization**: using greedy algorithms
- ♟ **Selection**: using graph based algorithms
- ♟ **Prioritization**: using capacity based fault detection
- ♟ Recently **Soft computing**

Why soft computing?

- ♣ the execution of the complete suite of test cases is an ***unacceptable practice*** for iterative and incremental development environments
- ♣ This raised the study of soft computing approaches, in conjunction to **agile approaches**
- ♣ Identify test cases that **Guarantee an acceptable level** of verification of software product

A more effective alternative

- ♣ Cluster test cases according to a pattern, criterion and characteristic

- ♣ Cluster: Test cases that detect the same fault

- ♣ How to cluster? Based on:

- ♣ Execution Profile

- ♣ Histories of test runs

- ♣ Function calls

- ♣ Partitioning programs

- ♣ Filtering partitions

- ♣ Access to database

- ♣ Finally: Sample the most representative test cases for each cluster

Challenges

- ⚠ Determine the quality and size of clusters to balance the cost and effectiveness of approach
- ⚠ Additional information to data such as numbers and sizes of clusters to be formed makes process expensive
 - ⚠ Alternative: use unsupervised clustering
 - ⚠ Probabilistic approaches determines number of clusters

Challenges

- ♟ Conventionally RT has been applied from **Code** perspective
 - ♟ Software is an FSM with constant Data
- ♟ Software with access to Database has two **Code** and **Data** perspective
 - ♟ **Code State** is a set of labeled memory
 - ♟ **Data State** organized by different by different data models, *such as relational*
 - ♟ Large amount of **Data** can be affected by a single line of **code**

Paper specs

♣ What are we presenting in this paper?

♣ RT for software with DB access under incremental software dev context

♣ Our approach:

♣ Clustering of DB access code along with DB Schema (by analyzing of code that accesses the DB)

♣ Unsupervised

♣ Random values that selects a set of test cases

Related Work

♣ Some of other RT approaches

Minimalization

let **p** = software product

let **t** = test cases associated with **p**

task → obtain a subset of test cases **t'** where
 t' does not include redundant or obsolete
 test cases to verify **p**

Selection

let p' = modifications to software product

task \rightarrow determine a subset of test cases t' where
 p' are tested with t'

Prioritization



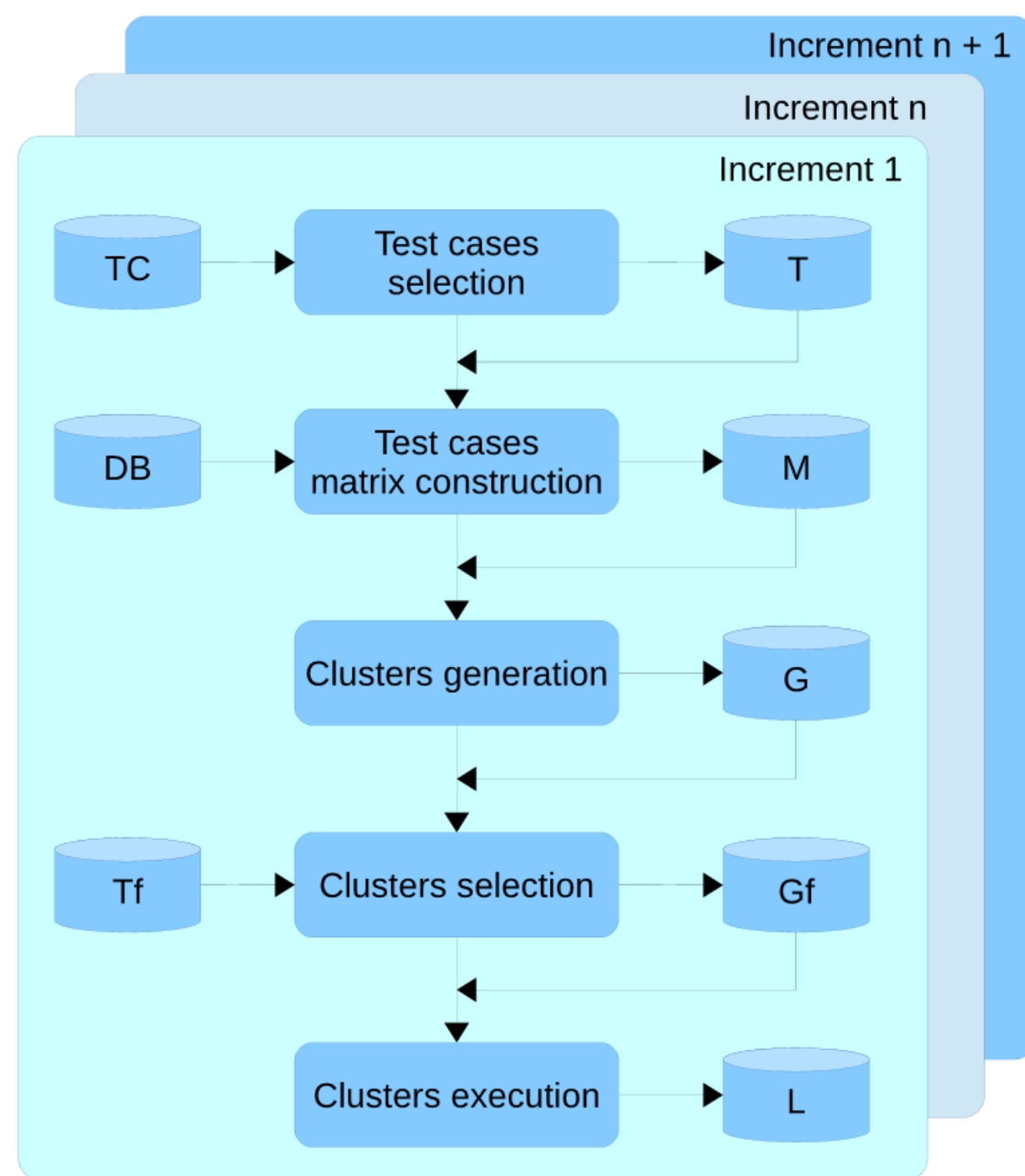
`task` -> determine an ideal permutation of sequences of test cases
to improve the RT technique

Proposed RT Approach

The background of the slide features a silhouette of a mountain range against a sunset sky. The sky transitions from a deep blue at the top to a warm orange near the horizon, where the sun is setting behind the mountains. The mountains are dark and jagged, creating a stark contrast with the bright sky.

Terminology

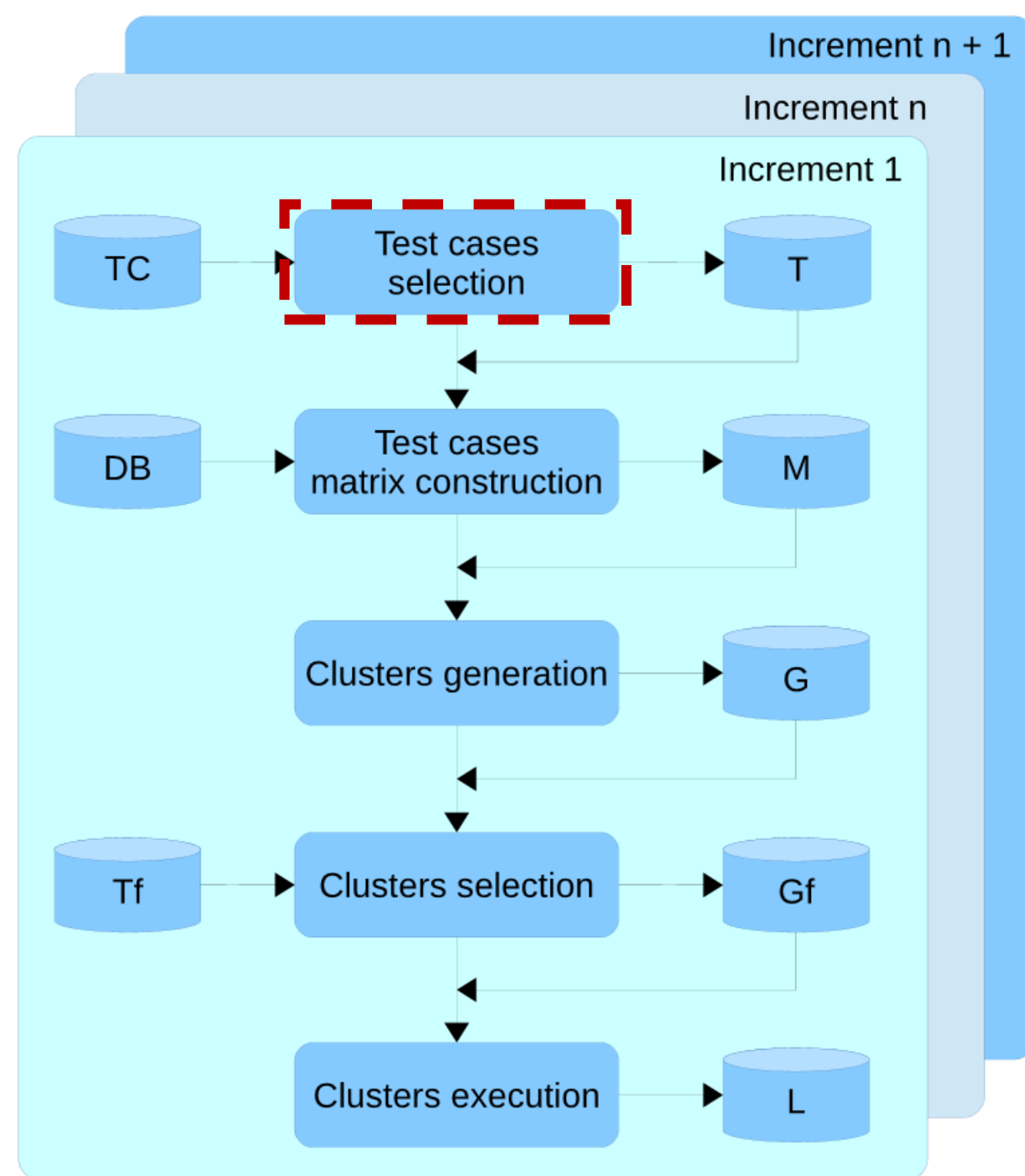
TC: Test cases of a given increment
T: Set of test cases with access to the database
DB: Database schema
M: Matrix of similarities
G: Clusters of test cases
Tf: Set of test cases with faults
Gf: Clusters of test cases with faults
L: List of faults of the regression testing



Selection

```
let new = new features and modifications
let T = from new where item access fields
      of given table of DB
```

TC: Test cases of a given increment
T: Set of test cases with access to the database
DB: Database schema
M: Matrix of similarities
G: Clusters of test cases
Tf: Set of test cases with faults
Gf: Clusters of test cases with faults
L: List of faults of the regression testing



Similarity matrix construction

task \rightarrow obtain info of tables and columns of each test case T

yield binary matrix $M[m, n]$ where

0 \rightarrow no DB access

1 \rightarrow DB access operation

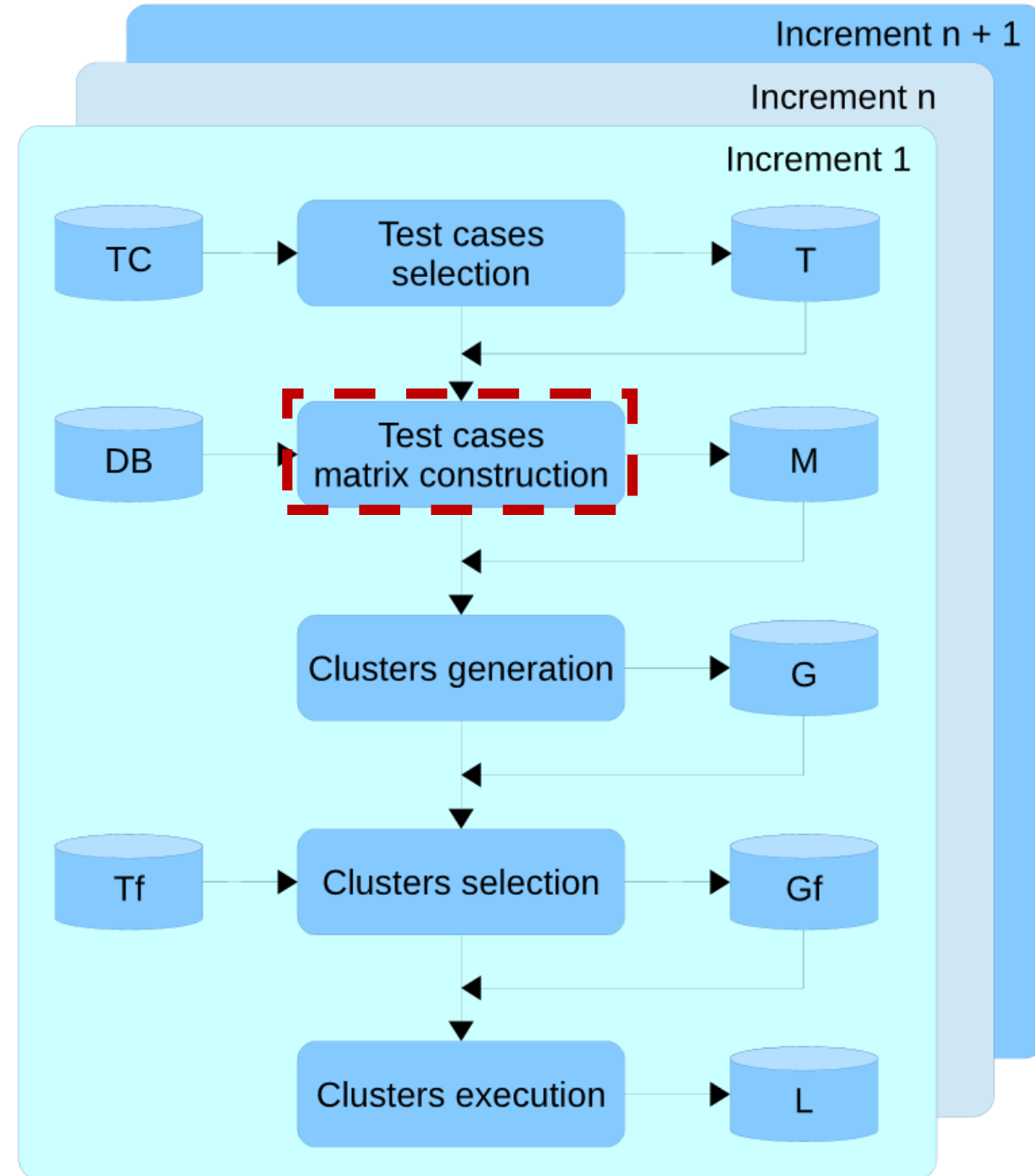
rows(m) \rightarrow test cases

cols(n) \rightarrow field of DB tables

TABLE 2. Matrix of test cases accessing the DB.

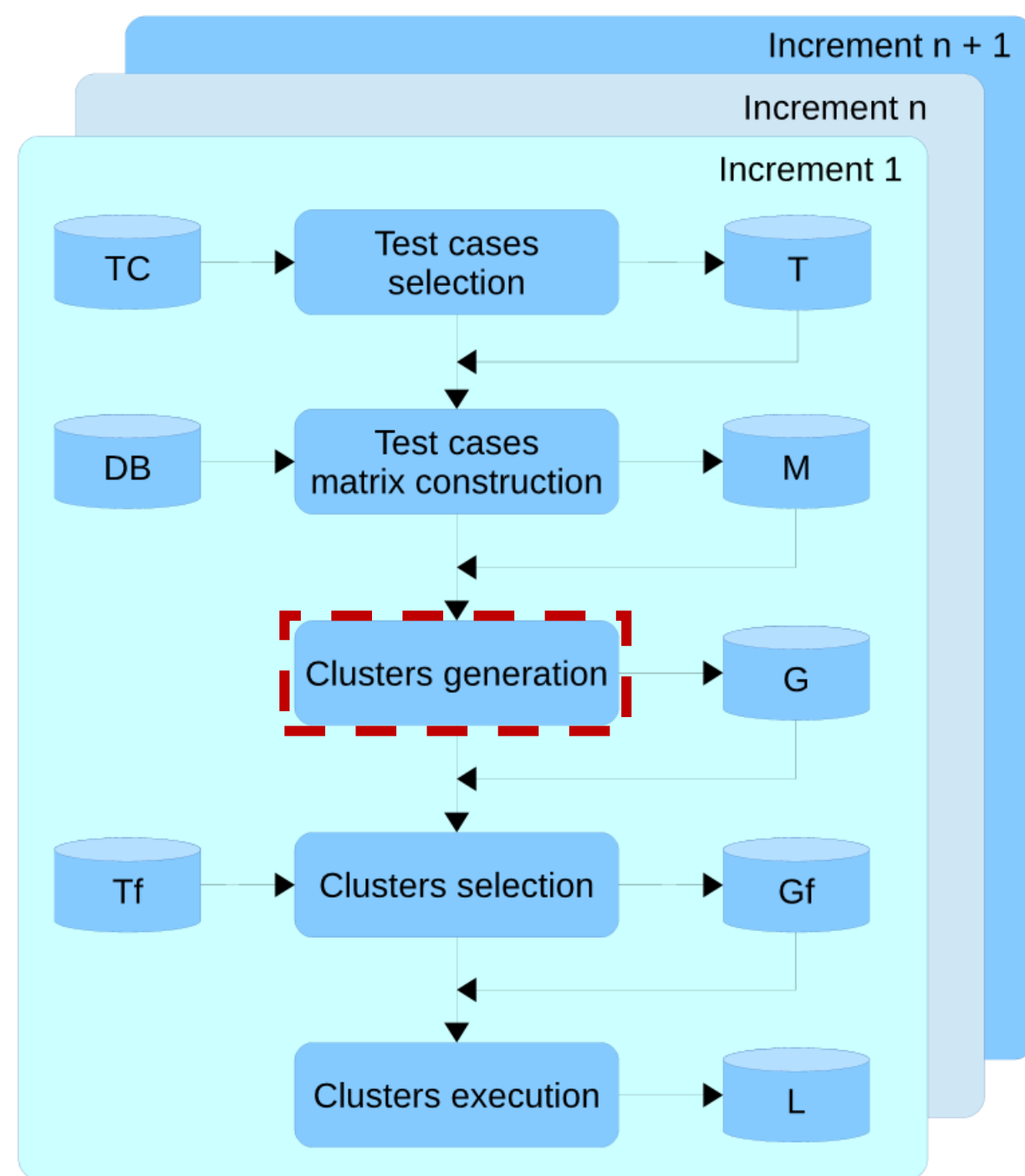
	Tables						
	Table ₁			...	Table _n		
	Col ₁	Col ₂	Col _i		Col ₁	Col ₂	Col _i
TC ₁	1	0	1	...	0	1	0

TC _n



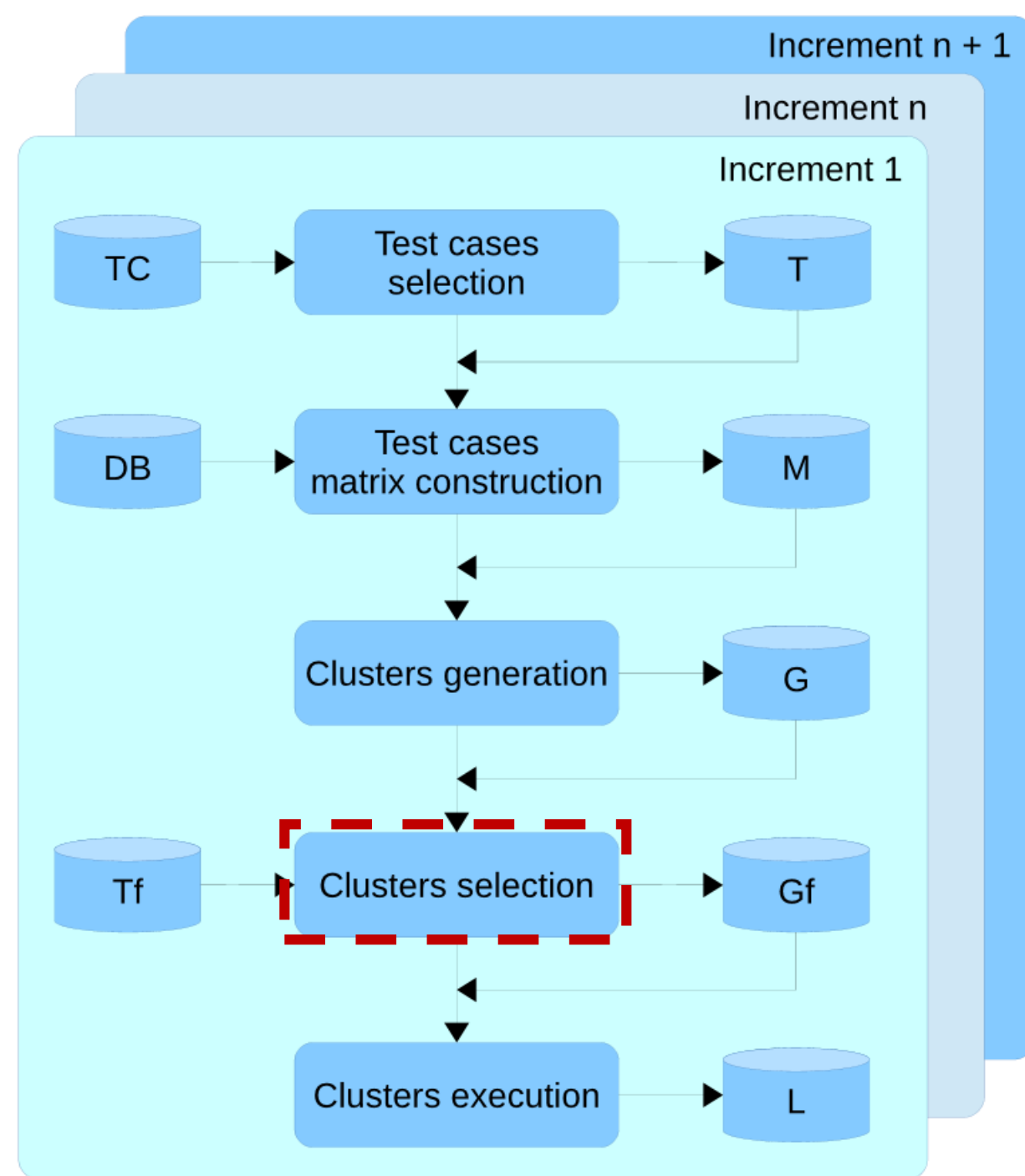
Clusters Generation

- ✓ Unsupervised clustering
- ✓ Expectation maximization
- ✓ We determine a **probabilistic distribution** function to assign a membership of a **test case** to a **cluster** based on **similarity matrix**
- ✓ Finite gaussian mixture model which all attributes are random
- ✓ **Result is one or more clusters**



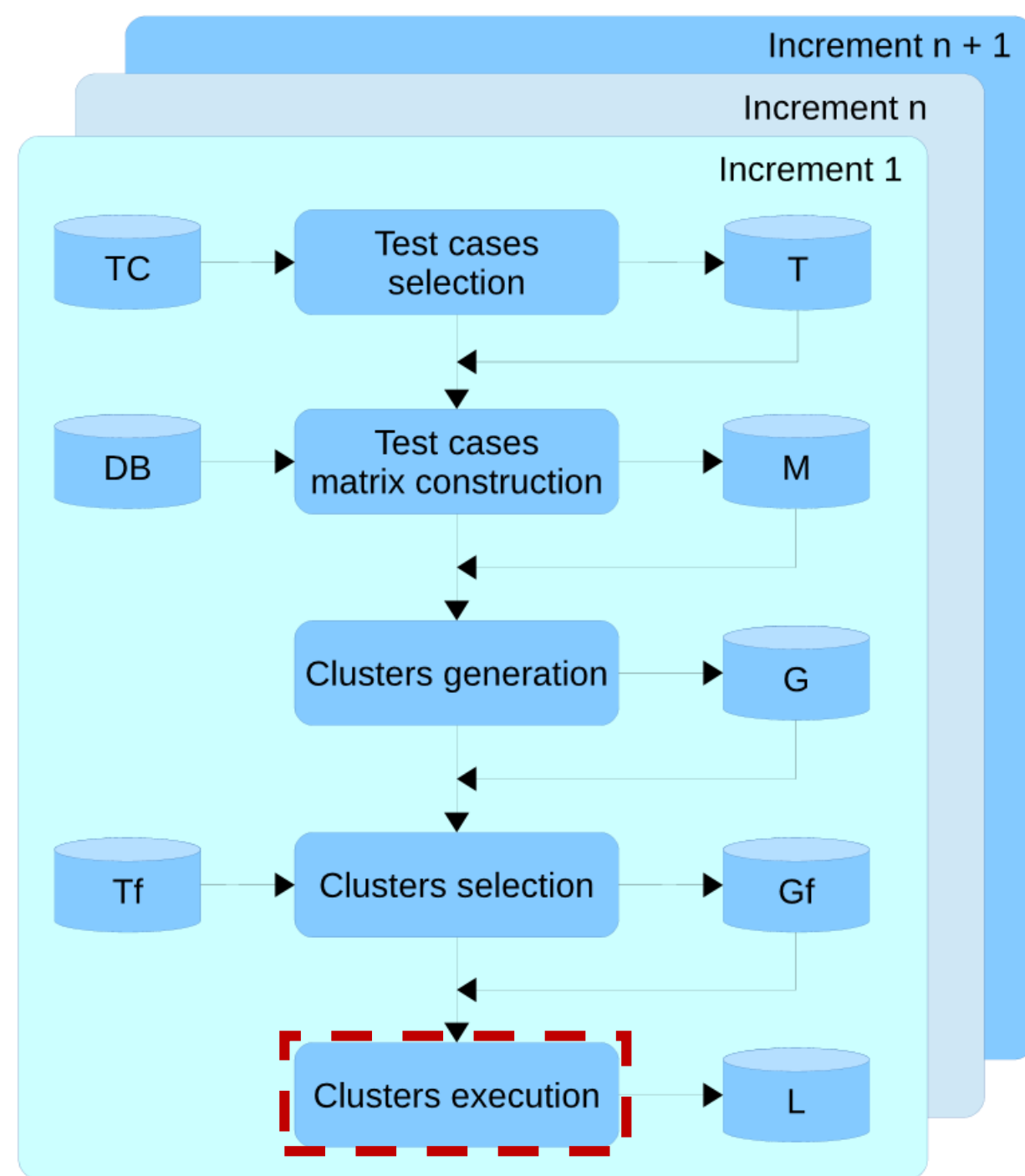
Clusters Selection

- ✓ Result: What clusters are going to be executed for RT
- ✓ Which clusters (New or modified test cases = Gf) are located
- ✓ Considering Tf: Test cases with faults



Clusters execution

- ✓ In this level, Gf are executed
- ✓ Generates a list L of faults related to test cases



Algorithm used in proposal

Algorithm Identifiers

Variable	Description
DB	Database
T	Set of test cases
F	Set of features
F_i	i-th feature, $f_i \in F$
TC_i	Set of test cases from i-th feature $f_i \in F$
tc_i	i-th Test case of a feature
$tc_{i,j}$	i-th Test case $tc_i \in TC_i$ from j-th feature
$TestDB$	Set of test cases with access to DB
$TestDB_{failed}$	Set of test cases failed with access to DB
$MatTestDB$	Test cases matrix with access DB
$MatClusterTestDB$	Test cases cluster matrix with access DB
$MatClusterRegressionTestDB$	Failed matrix of clustered test cases with access to DB
$FaultsList$	Set of faults for the regression test
$NumTestDB$	Number of test cases with access to the DB

Main algorithm

1. For each feature (f_i), a selection of test cases that access DB is performed
2. Similarity matrix construction
3. Clustering analysis
4. Execute clusters containing failed clusters, and their faults examined

Algorithm 1 clusterSelectionRegression Test

Input: f_i , TC_i

Output: Faults List

```
1: clusterSelectionRegressionTest process
2: repeat
3:   for each  $f_i$  do //  $f_i \in F$ 
4:     TestDB  $\leftarrow$  selectTestDB( $f_i$ ,  $TC_i$ )
5:     MatTestBD  $\leftarrow$  buildMatrixTestDB(testDB)
6:     MatClusterTestDB  $\leftarrow$  clusterizationEM(MatTestBD)
7:     MatClusterRegressionTestDB  $\leftarrow$  (MatCluster
        TestDB  $\cap$  TestDBfailed)
8:     executeClusterRegressionTestDB
        (MatClusterRegressionTestDB)
9:   end for
10:  until  $\nexists f_i$ 
11: end clusterSelectionRegressionTest process
```

SelectTestDb Algorithm

1. For implementing this, authors created a tool that analyzes DB Schema (yeah, it's magic)

Algorithm 2 selectedTestDB

Input: f_i , TC_i

Output: TestDB, NumTestDB

```
1: selectedTestDB process
2: numTestDB  $\leftarrow$  0
3: repeat
4:   if  $TC_i$  access DB then
5:     TestDB  $\leftarrow$  TestDB  $\cup$   $TC_i$ 
6:     NumTestDB++
7:   end if
8: until  $\nexists TC_i$ 
9: end selectedTestDB process
```

Similarity matrix construction Algorithm

1. For implementing this, authors created a tool that analyzes DB Schema (yeah, it's magic)

Algorithm 3 buildMatrixTestDB

Input: TestDB, numTestDB

Output: MatTestDB

```
1: buildMatrixTestDB process
2:  i ← 1
3:  while i < numTestBD do
4:    MatTestDB ← (MatTestDB ∪ (TestDBi))
5:  end while
6: end buildMatrixTestDB process
```

Cluster generation and selection Algorithm

1. For implementing this, authors used Weka data mining tool for practical reasons
2. Also used probabilistic algorithm EM (Expectation maximization) which has the advantage of **determining a k number of clusters** based on the information of the test cases of the similarity matrix

Algorithm 4 clusterizationEM

Input: MatTestDB

Output: MatClusterTestDB

- 1: clusterizationEM process
- 2: Θ Vector desconocido de parámetros
- 3: $\Theta^0, \Theta^1, \dots, \Theta^T$, T criterio de convergencia
- 4: $T \leftarrow 0$
- 5: $\Theta^0 \leftarrow 0$
- 6: Repetir
- 7: $Q(\Theta, \Theta^t) \leftarrow E[\log p(x^g, x^m | \Theta) | x^g, \Theta^t]$
- 8: $\Theta^{t+1} \leftarrow \arg \max_{\Theta} Q(\Theta, \Theta^t)$
- 9: $t \leftarrow t+1$
- 10: until convergence criterion
- 11: end clusterizationEM process

$$\begin{aligned}\boldsymbol{\tau}^{(t+1)} &= \arg \max_{\boldsymbol{\tau}} Q(\boldsymbol{\theta}|\boldsymbol{\theta}^{(t)}) \\ &= \arg \max_{\boldsymbol{\tau}} \left\{ \left[\sum_{i=1}^n T_{1,i}^{(t)} \right] \log \tau_1 + \left[\sum_{i=1}^n T_{2,i}^{(t)} \right] \log \tau_2 \right\}\end{aligned}$$

$$\tau_j^{(t+1)} = \frac{\sum_{i=1}^n T_{j,i}^{(t)}}{\sum_{i=1}^n (T_{1,i}^{(t)} + T_{2,i}^{(t)})} = \frac{1}{n} \sum_{i=1}^n T_{j,i}^{(t)}$$

$$\begin{aligned}(\boldsymbol{\mu}_1^{(t+1)}, \Sigma_1^{(t+1)}) &= \arg \max_{\boldsymbol{\mu}_1, \Sigma_1} Q(\boldsymbol{\theta}|\boldsymbol{\theta}^{(t)}) \\ &= \arg \max_{\boldsymbol{\mu}_1, \Sigma_1} \sum_{i=1}^n T_{1,i}^{(t)} \left\{ -\frac{1}{2} \log |\Sigma_1| - \frac{1}{2} (\mathbf{x}_i - \boldsymbol{\mu}_1)^\top \Sigma_1^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_1) \right\}\end{aligned}$$

$$\begin{aligned}{}_9\boldsymbol{\mu}_2^{(t+1)} &= \frac{\sum_{i=1}^n T_{2,i}^{(t)} \mathbf{x}_i}{\sum_{i=1}^n T_{2,i}^{(t)}} {}_9\Sigma_1^{(t+1)} = \frac{\sum_{i=1}^n T_{1,i}^{(t)} (\mathbf{x}_i - \boldsymbol{\mu}_1^{(t+1)}) (\mathbf{x}_i - \boldsymbol{\mu}_1^{(t+1)})^\top}{\sum_{i=1}^n T_{1,i}^{(t)}} {}_9\boldsymbol{\mu}_1^{(t+1)} = \frac{\sum_{i=1}^n T_{1,i}^{(t)} \mathbf{x}_i}{\sum_{i=1}^n T_{1,i}^{(t)}} \\ \Sigma_2^{(t+1)} &= \frac{\sum_{i=1}^n T_{2,i}^{(t)} (\mathbf{x}_i - \boldsymbol{\mu}_2^{(t+1)}) (\mathbf{x}_i - \boldsymbol{\mu}_2^{(t+1)})^\top}{\sum_{i=1}^n T_{2,i}^{(t)}}\end{aligned}$$

$$E_{Z|\boldsymbol{\theta}^{(t)}, \mathbf{x}}[\log L(\boldsymbol{\theta}^{(t)}; \mathbf{x}, \mathbf{Z})] \leq E_{Z|\boldsymbol{\theta}^{(t-1)}, \mathbf{x}}[\log L(\boldsymbol{\theta}^{(t-1)}; \mathbf{x}, \mathbf{Z})] + \epsilon$$



Cluster generation and selection Algorithm

1. For implementing this, authors used Weka data mining tool for practical reasons
2. Also used probabilistic algorithm EM (Expectation maximization) which has the advantage of **determining a k number of clusters** based on the information of the test cases of the similarity matrix

Algorithm 4 clusterizationEM

Input: MatTestDB

Output: MatClusterTestDB

- 1: clusterizationEM process
- 2: Θ Vector desconocido de parámetros
- 3: $\Theta^0, \Theta^1, \dots, \Theta^T$, T criterio de convergencia
- 4: $T \leftarrow 0$
- 5: $\Theta^0 \leftarrow 0$
- 6: Repetir
- 7: $Q(\Theta, \Theta^t) \leftarrow E[\log p(x^g, x^m | \Theta) | x^g, \Theta^t]$
- 8: $\Theta^{t+1} \leftarrow \arg \max_{\Theta} Q(\Theta, \Theta^t)$
- 9: $t \leftarrow t+1$
- 10: until convergence criterion
- 11: end clusterizationEM process

Thank you!

