

Storage and Data Recovery

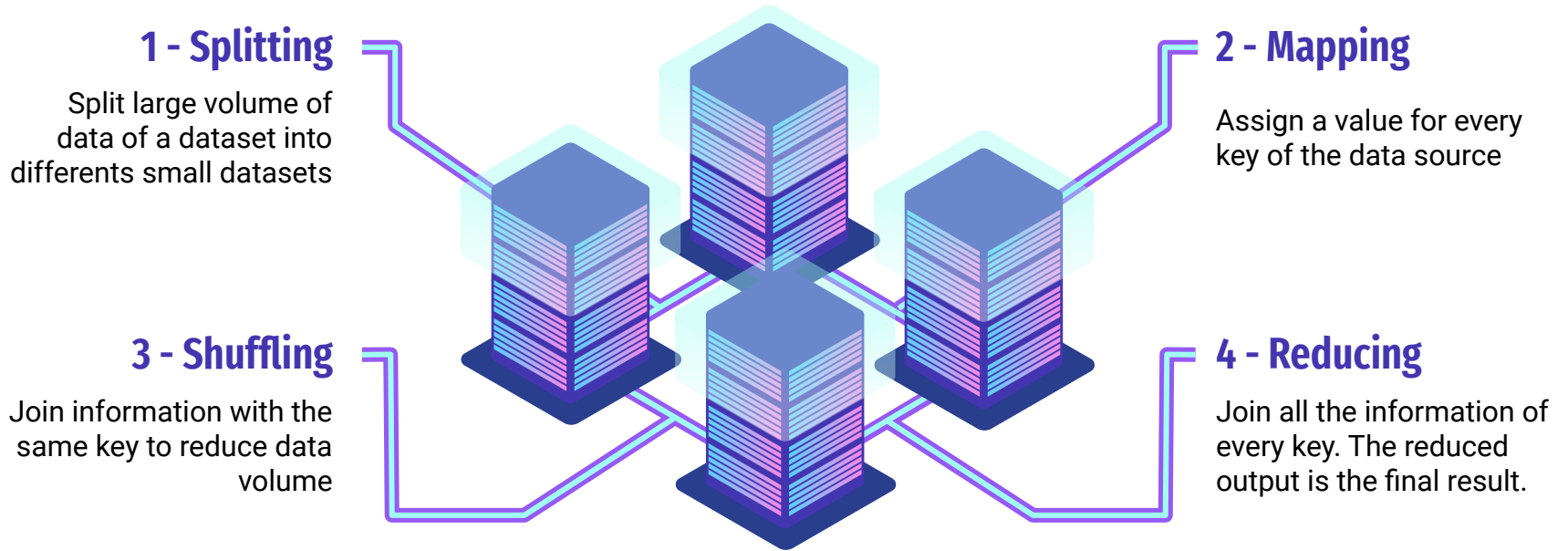
Exercise 4: Map reduce

Master's Degree in Intelligent Systems
Iván López Muñoz



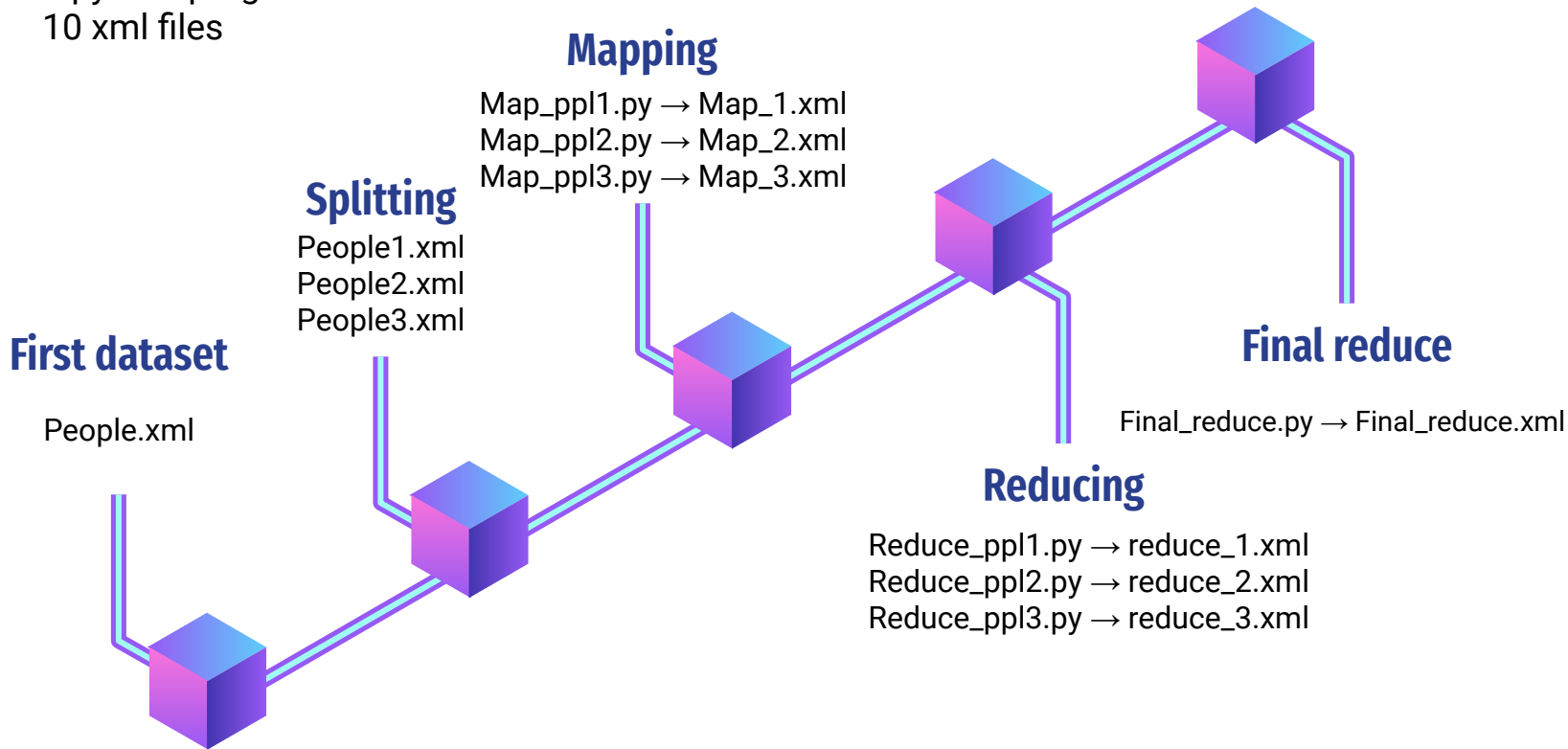
What is map reduce?

Map-reduce is a data processing paradigm for condensing large volumes of data into useful aggregated results



File summary

- 7 python programs
- 10 xml files



Splitting

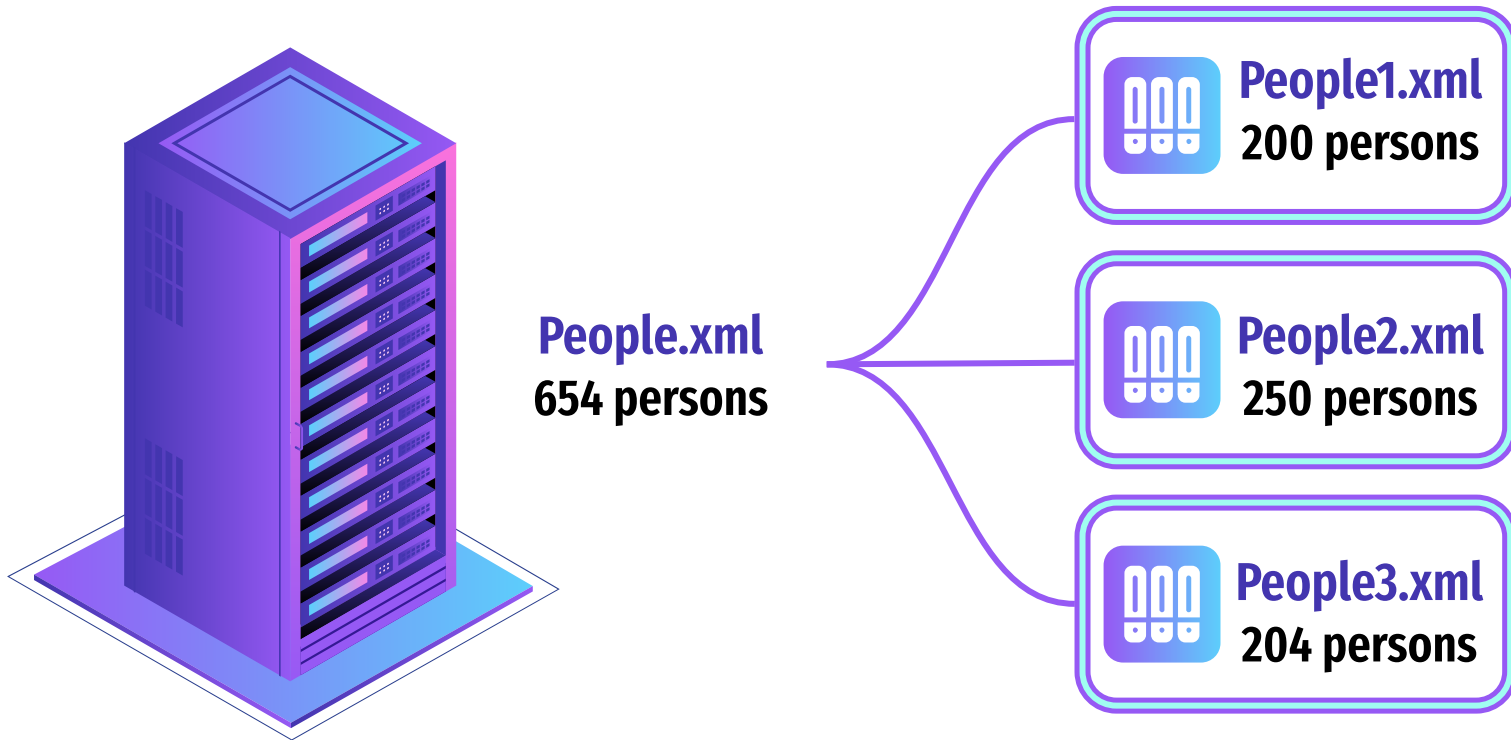
First, we need to split our dataset into smaller datasets. Here we have our People.xml dataset with 654 persons

```
<person><id>1</id><sex>H</sex><name>Juan</name><surname1>Prats</surname1><surname2>Ramis</surname2><birthdate>1816-1-1</birthdate></person>
<person><id>2</id><sex>D</sex><name>Antonia</name><surname1>Ferregut</surname1><surname2>Jaume</surname2><birthdate>1816-1-1</birthdate></person>
<person><id>3</id><sex>D</sex><name>Juana</name><surname1>Perello</surname1><surname2>Figueroa</surname2><birthdate>1816-1-6</birthdate></person>
<person><id>4</id><sex>H</sex><name>Jaime</name><surname1>Pujades</surname1><surname2>Tortella</surname2><birthdate>1816-1-7</birthdate></person>
<person><id>5</id><sex>D</sex><name>Francisca</name><surname1>Garriga</surname1><surname2>Bauza</surname2><birthdate>1816-1-8</birthdate></person>
<person><id>6</id><sex>D</sex><name>Geronima</name><surname1>Coll</surname1><surname2>Pujades</surname2><birthdate>1816-1-9</birthdate></person>
<person><id>7</id><sex>H</sex><name>Gabriel</name><surname1>Estrafy</surname1><surname2>Garau</surname2><birthdate>1816-1-9</birthdate></person>
<person><id>8</id><sex>H</sex><name>Arnaldo</name><surname1>Mulet</surname1><surname2>Perello</surname2><birthdate>1816-1-9</birthdate></person>
<person><id>9</id><sex>H</sex><name>Gabriel</name><surname1>Llinas</surname1><surname2>Alomar</surname2><birthdate>1816-1-10</birthdate></person>
<person><id>10</id><sex>D</sex><name>Margharita</name><surname1>Domenech</surname1><surname2>Pujades</surname2><birthdate>1816-1-10</birthdate></person>
<person><id>11</id><sex>H</sex><name>Pedro</name><surname1>Ferregut</surname1><surname2>Rayo</surname2><birthdate>1816-1-11</birthdate></person>
<person><id>12</id><sex>D</sex><name>Margarita</name><surname1>Genestra</surname1><surname2>Estrafy</surname2><birthdate>1816-1-11</birthdate></person>
<person><id>13</id><sex>D</sex><name>Antonia</name><surname1>Pujades</surname1><surname2>Gual</surname2><birthdate>1816-1-12</birthdate></person>
<person><id>14</id><sex>D</sex><name>Catalina</name><surname1>Llompard</surname1><surname2>Campins</surname2><birthdate>1816-1-12</birthdate></person>
<person><id>15</id><sex>D</sex><name>Juana</name><surname1>Munar</surname1><surname2>Siquier</surname2><birthdate>1816-1-12</birthdate></person>
```

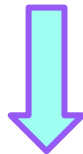
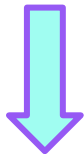
•
•
•

```
<person><id>639</id><sex>H</sex><name>Pere</name><surname1>Mulet</surname1><surname2>Perello</surname2><birthdate>1820-11-24</birthdate></person>
<person><id>640</id><sex>D</sex><name>Margalida</name><surname1>Bertran</surname1><surname2>Company</surname2><birthdate>1820-11-30</birthdate></person>
<person><id>641</id><sex>D</sex><name>Margarita</name><surname1>Gual</surname1><surname2>Vicens</surname2><birthdate>1820-12-2</birthdate></person>
<person><id>642</id><sex>H</sex><name>Miguel</name><surname1>Ferregut</surname1><surname2>Rayo</surname2><birthdate>1820-12-3</birthdate></person>
<person><id>643</id><sex>D</sex><name>Margarita</name><surname1>Pujades</surname1><surname2>Ferregut</surname2><birthdate>1820-12-4</birthdate></person>
<person><id>644</id><sex>D</sex><name>Francisca</name><surname1>Martorell</surname1><surname2>Martorell</surname2><birthdate>1820-12-4</birthdate></person>
<person><id>645</id><sex>D</sex><name>Leocadia</name><surname1>Desconeget</surname1><surname2>Desconeget</surname2><birthdate>1820-12-9</birthdate></person>
<person><id>646</id><sex>H</sex><name>Andres</name><surname1>Aguilo</surname1><surname2>Porteza</surname2><birthdate>1820-12-16</birthdate></person>
<person><id>647</id><sex>D</sex><name>Francisca</name><surname1>Bestard</surname1><surname2>Rossello</surname2><birthdate>1820-12-18</birthdate></person>
<person><id>648</id><sex>H</sex><name>Onofre</name><surname1>Llompard</surname1><surname2>Rubi</surname2><birthdate>1820-12-19</birthdate></person>
<person><id>649</id><sex>H</sex><name>Gabriel</name><surname1>Coll</surname1><surname2>Fiol</surname2><birthdate>1820-12-20</birthdate></person>
<person><id>650</id><sex>D</sex><name>Juana</name><surname1>Ferrer</surname1><surname2>Tortella</surname2><birthdate>1820-12-21</birthdate></person>
<person><id>651</id><sex>H</sex><name>Martin</name><surname1>Coll</surname1><surname2>Ferrer</surname2><birthdate>1820-12-28</birthdate></person>
<person><id>652</id><sex>H</sex><name>Gabriel</name><surname1>Matheu</surname1><surname2>Figueroa</surname2><birthdate>1820-12-28</birthdate></person>
<person><id>653</id><sex>D</sex><name>Juana</name><surname1>Vallespir</surname1><surname2>Garau</surname2><birthdate>1820-12-29</birthdate></person>
<person><id>654</id><sex>H</sex><name>Antonio</name><surname1>Martorell</surname1><surname2>Llobera</surname2><birthdate>1820-12-31</birthdate></person>
```

Splitting



Mapping



How to Map with Python

Map_pp1.py

```
1  # -*- coding: utf-8 -*-
2  """
3  Created on Wed Oct 18 11:57:25 2023
4
5  @author: Ivan
6  """
7
8  import xml.etree.ElementTree as ET
9
10 # Create an empty map element
11 map_element = ET.Element("map")
12
13 # Specify the input XML file
14 input_file = 'people1.xml'
15
16 # Parse the input XML file
17 tree = ET.parse(input_file)
18 root = tree.getroot()
19
20 # Iterate over the 'name' elements and create item elements
21 for person in root:
22     name = person.find('name').text
23
24     # Create an item element
25     item_element = ET.SubElement(map_element, "item")
26
27     # Add the 'name' and 'value' elements to the item
28     name_element = ET.SubElement(item_element, "name")
29     name_element.text = name
30
31     value_element = ET.SubElement(item_element, "value")
32     value_element.text = "1"
33
34 # Create an ElementTree object with the map element
35 result_tree = ET.ElementTree(map_element)
36
37 # Write the results to an XML file
38 result_tree.write("Map_1.xml", encoding="utf-8")
```

- Read each smaller dataset splitted previously
- Assign a value for every key of the data source, in this case for every name in the dataset
- From each 'people' dataset we obtain a list of names with a value ready to shuffle

Map_1.xml

```
<map>
  <item>
    <name>Juan</name>
    <value>1</value>
  </item>
  <item>
    <name>Antonia</name>
    <value>1</value>
  </item>
  <item>
    <name>Juana</name>
    <value>1</value>
  </item>
  <item>
    <name>Jaime</name>
    <value>1</value>
  </item>
</map>
```

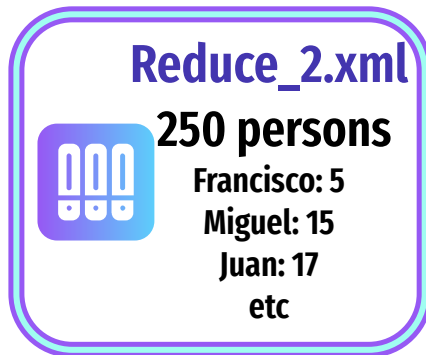
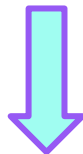
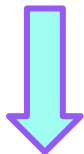
Map_2.xml

```
<map>
  <item>
    <name>Sebastian</name>
    <value>1</value>
  </item>
  <item>
    <name>Margarita</name>
    <value>1</value>
  </item>
  <item>
    <name>Cathalina</name>
    <value>1</value>
  </item>
  <item>
    <name>Ana</name>
    <value>1</value>
  </item>
</map>
```

Map_3.xml

```
<map>
  <item>
    <name>Francisco</name>
    <value>1</value>
  </item>
  <item>
    <name>Miguel</name>
    <value>1</value>
  </item>
  <item>
    <name>Juan</name>
    <value>1</value>
  </item>
  <item>
    <name>Guillermo</name>
    <value>1</value>
  </item>
</map>
```

Reducing



How to Reduce with Python

Reduce_ppl1.py

```
2  """
3  Created on Wed Oct 18 12:04:38 2023
4
5  @author: Ivan
6  """
7
8  import xml.etree.ElementTree as ET
9  from collections import defaultdict
10
11 # Create a defaultdict to store name counts
12 name_count = defaultdict(int)
13
14 # Parse the input XML file
15 tree = ET.parse("Map_1.xml")
16 root = tree.getroot()
17
18 # Iterate over the 'item' elements and
19 # extract names and values
20 for item in root.findall("item"):
21     name = item.find("name").text
22     value = int(item.find("value").text)
23
24     # Increment the count for each name
25     name_count[name] += value
26
27 # Create a new map element for the results
28 result_map = ET.Element("reduce")
29
30 # Iterate over the unique names and their counts
31 for name, count in name_count.items():
32     # Create an item element
33     item_element = ET.Element("item")
34
35     # Add the 'name' and 'value' elements to the item
36     name_element = ET.Element("name")
37     name_element.text = name
38     item_element.append(name_element)
39
40     value_element = ET.Element("value")
41     value_element.text = str(count)
42     item_element.append(value_element)
43
44     # Add the item to the result map
45     result_map.append(item_element)
46
47 # Create an ElementTree object with the result map
48 result_tree = ET.ElementTree(result_map)
49
50 # Write the results to a new XML file
51 result_tree.write("reduce_1.xml", encoding="utf-8")
```

- Read each smaller dataset splitted and mapped previously
- Iterate on each person's name and add the values for each one
- From each 'mapped' dataset we obtain a list of names with the number of times they appear in each dataset

reduce_1.xml

```
<reduce>
  <item>
    <name>Juan</name>
    <value>7</value>
  </item>
  <item>
    <name>Antonia</name>
    <value>11</value>
  </item>
  <item>
    <name>Juana</name>
    <value>13</value>
  </item>
  <item>
    <name>Jaime</name>
    <value>1</value>
  </item>
```

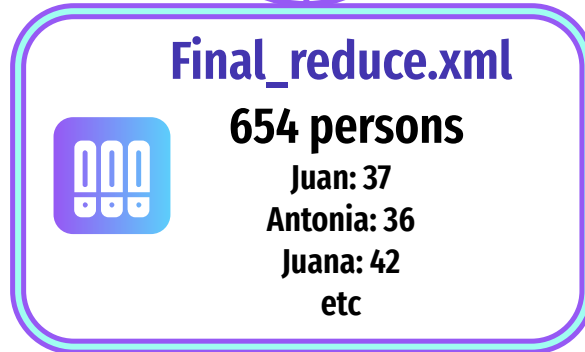
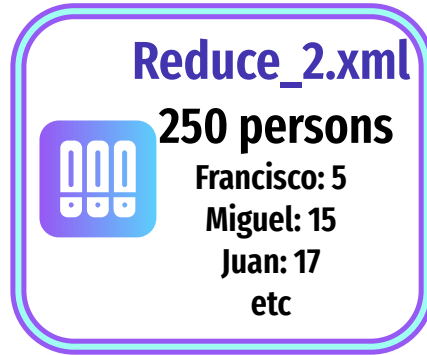
reduce_2.xml

```
<reduce>
  <item>
    <name>Francisco</name>
    <value>5</value>
  </item>
  <item>
    <name>Miguel</name>
    <value>15</value>
  </item>
  <item>
    <name>Juan</name>
    <value>17</value>
  </item>
  <item>
    <name>Guillermo</name>
    <value>5</value>
  </item>
```

reduce_3.xml

```
<reduce>
  <item>
    <name>Sebastian</name>
    <value>1</value>
  </item>
  <item>
    <name>Margarita</name>
    <value>8</value>
  </item>
  <item>
    <name>Cathalina</name>
    <value>8</value>
  </item>
  <item>
    <name>Ana</name>
    <value>5</value>
  </item>
```

Final reduce



Final reduce

Final_reduce.py

```
2  """
3  Created on Wed Oct 18 12:07:06 2023
4
5  @author: Ivan
6  """
7
8  import xml.etree.ElementTree as ET
9  from collections import defaultdict
10
11  # Create a defaultdict to store name counts
12  name_count = defaultdict(int)
13
14  # Iterate over the XML files
15  for filename in ['reduce_1.xml', 'reduce_2.xml', 'reduce_3.xml']:
16      tree = ET.parse(filename)
17      root = tree.getroot()
18
19      # Iterate over the 'name' elements and increment the count
20      for person in root.findall("item"):
21          name = person.find('name').text
22          value = int(person.find("value").text)
23          name_count[name] += value
24
25  # Create the XML result document
26  result_root = ET.Element("reduce")
27
28  # Add name count elements to the result document
29  for name, count in name_count.items():
30      # Create an item element
31      item_element = ET.SubElement(result_root, "item")
32
33      # Add the 'name' and 'value' elements to the item
34      name_element = ET.SubElement(item_element, "name")
35      name_element.text = name
36
37      value_element = ET.SubElement(item_element, "value")
38      value_element.text = str(count)
39
40  # Create an ElementTree object and write it to an XML file
41  result_tree = ET.ElementTree(result_root)
42  result_tree.write("Final_reduce.xml", encoding="utf-8")
```

- Read all reduced datasets
- Iterate on each person's name in every reduced file and add the values for each one
- Finally, we obtain a list with the number of people who have the same name

```
<reduce>
  <item>
    <name>Juan</name>
    <value>37</value>
  </item>
  <item>
    <name>Antonia</name>
    <value>36</value>
  </item>
  <item>
    <name>Juana</name>
    <value>42</value>
  </item>
  <item>
    <name>Jaime</name>
    <value>2</value>
  </item>
```

Final_reduce.xml

Conclusions about map reducing

- **Scalability:** Map Reduce efficiently processes substantial volumes of data, making it suitable for handling the vast amount of information within these XML files.
- **Data Flexibility:** Map Reduce can work with structured, semi-structured, and unstructured data, making it a versatile choice for the diverse content found in the XML files.
- **Distributed Data Processing:** The ability to process data in a distributed manner helps minimize communication overhead and maximizes computational efficiency.
- **MapReduce Phases:** The Map Reduce process involves distinct phases like splitting, mapping, shuffling, and reducing, providing a structured approach to data processing that simplifies complex tasks.

→ By modifying the final program by counting the times it writes the names of people in the final xml file, from the data set of 654 people in People.xml, we can find 86 different names

→ That's why the shuffling part has not been done, since we would have 86 xml files, each of them with a name and the number of apparitions in the dataset