

# Can Stack Overflow Posts **Capture** Library Features?

The 19th Belgium-Netherlands Software Evolution Workshop

**Camilo Velázquez-Rodríguez, Eleni Constantinou and Coen De Roover**

December 4th, 2020



# Introduction

Software products that evolve together in the same environment [1]



# Introduction

Software products that evolve together in the same environment [1]



# Introduction

Software products that evolve together in the same environment [1]



- I/O Utilities (17)**
- HTTP Clients (81)**
- Dependency Injection (51)**
- Mocking (41)**
- XML Processing (34)**
- JSON Libraries (74)**
- Collections (72)**
- Reflection Libraries (47)**
- Command Line Parsers (45)**
- ⋮

[1] - M. Lungu, "Towards reverse engineering software ecosystems," in *2008 ICSM*. IEEE, 2008, pp. 428–431.

# Introduction

Software products that evolve together in the same environment [1]



- I/O Utilities (17)**
- HTTP Clients (81)**      **Dependency Injection (51)**
- Mocking (41)**      **XML Processing (34)**
- JSON Libraries (74)**      **Collections (72)**
- Reflection Libraries (47)**
- Command Line Parsers (45)**

⋮

## JCommander Library

Version	Repository	Usages	Date
1.78	Central	20	Aug, 2019
1.77	Central	0	Aug, 2019
1.75	Central	3	Aug, 2019
1.72	Central	193	Jun, 2017
1.69	Central	18	Apr, 2017
1.64	Central	23	Mar, 2017
1.60	Central	30	Jan, 2017
1.58	Central	21	Sep, 2016
1.48	Central	259	Apr, 2015
1.47	Central	24	Dec, 2014
1.35	Central	123	Feb, 2014
1.32	Central	55	Sep, 2013
1.30	Central	94	Oct, 2012

[1] - M. Lungu, "Towards reverse engineering software ecosystems," in 2008 ICSM. IEEE, 2008, pp. 428–431.

# Introduction

Software products that evolve together in the same environment [1]



Library Selection Issues

- I/O Utilities (17)**
- HTTP Clients (81)**      **Dependency Injection (51)**
- Mocking (41)**      **XML Processing (34)**
- JSON Libraries (74)**      **Collections (72)**
- Reflection Libraries (47)**
- Command Line Parsers (45)**
- ⋮

## JCommander Library

Version	Repository	Usages	Date
1.78	Central	20	Aug, 2019
1.77	Central	0	Aug, 2019
1.75	Central	3	Aug, 2019
1.72	Central	193	Jun, 2017
1.69	Central	18	Apr, 2017
1.64	Central	23	Mar, 2017
1.60	Central	30	Jan, 2017
1.58	Central	21	Sep, 2016
1.48	Central	259	Apr, 2015
1.47	Central	24	Dec, 2014
1.35	Central	123	Feb, 2014
1.32	Central	55	Sep, 2013
1.30	Central	94	Oct, 2012

[1] - M. Lungu, "Towards reverse engineering software ecosystems," in 2008 ICSM. IEEE, 2008, pp. 428–431.

# Introduction

## Guava

But trawling through Javadoc isn't always the most effective way to learn how to make best use of a library. Here, we try to provide readable and pleasant explanations of **some of the most popular and most powerful features of Guava.**

# Introduction

## Guava

But trawling through Javadoc isn't always the most effective way to learn how to make best use of a library. Here, we try to provide readable and pleasant explanations of **some of the most popular and most powerful features of Guava**.

- Collections: Guava's extensions to the JDK collections ecosystem. These are some of the most mature and popular parts of Guava.
  - **Immutable collections**, for defensive programming, constant collections, and improved efficiency.
  - **New collection types**, for use cases that the JDK collections don't address as well as they could: multisets, multimaps, tables, bidirectional maps, and more.
  - **Powerful collection utilities**, for common operations not provided in `java.util.Collections`.
  - **Extension utilities**: writing a `Collection` decorator? Implementing `Iterator`? We can make that easier.



# Introduction

## Guava

But trawling through Javadoc isn't always the most effective way to learn how to make best use of a library. Here, we try to provide readable and pleasant explanations of **some of the most popular and most powerful features of Guava**.

- Collections: Guava's extensions to the JDK collections ecosystem. These are some of the most mature and popular
- **Graphs**: a library for modeling **graph**-structured data, that is, entities and the relationships between them. Key features include:
  - **Graph**: a graph whose edges are anonymous entities with no identity or information of their own.
  - **ValueGraph**: a graph whose edges have associated non-unique values.
  - **Network**: a graph whose edges are unique objects.
  - Support for graphs that are mutable and immutable, directed and undirected, and several other properties.

# Introduction

## Guava

But trawling through Javadoc isn't always the most effective way to learn how to make best use of a library. Here, we try to provide readable and pleasant explanations of **some of the most popular and most powerful features of Guava**.

- Collections: Guava's extensions to the JDK collections ecosystem. These are some of the most mature and popular
- **Graphs**: a library for modeling **graph**-structured data, that is, entities and the relationships between them. Key features include:
  - **Strings**: A few extremely useful string utilities: splitting, joining, padding, and more.
  - **Primitives**: operations on primitive types, like `int` and `char`, not provided by the JDK, including unsigned variants for some types.
  - **Ranges**: Guava's powerful API for dealing with ranges on `Comparable` types, both continuous and discrete.
  - **I/O**: Simplified I/O operations, especially on whole I/O streams and files, for Java 5 and 6.
  - **Hashing**: Tools for more sophisticated hashes than what's provided by `Object.hashCode()`, including Bloom filters.
  - **EventBus**: Publish-subscribe-style communication between components without requiring the components to explicitly register with one another.
  - **Math**: Optimized, thoroughly tested math utilities not provided by the JDK.
  - **Reflection**: Guava utilities for Java's reflective capabilities.

# Introduction

## Stack Overflow Q&A

How do I convert a String to an int in Java? **Guava equivalent for IOUtils.toString(InputStream)**

How do I compare strings in Java?

Remove last character of a **StringBuilder**?

**Guava: Splitter and considering Escaping?**

**String utilities**

**initializing a Guava ImmutableMap**

**Flattening an Iterable<Iterable<T>> in Guava**

**How to directly initialize a HashMap (in a literal way)?**

**Google Guava isEmpty for collections**

**Google Guava “zip” two lists**

**Collection utilities**

**Merge ranges with guava**

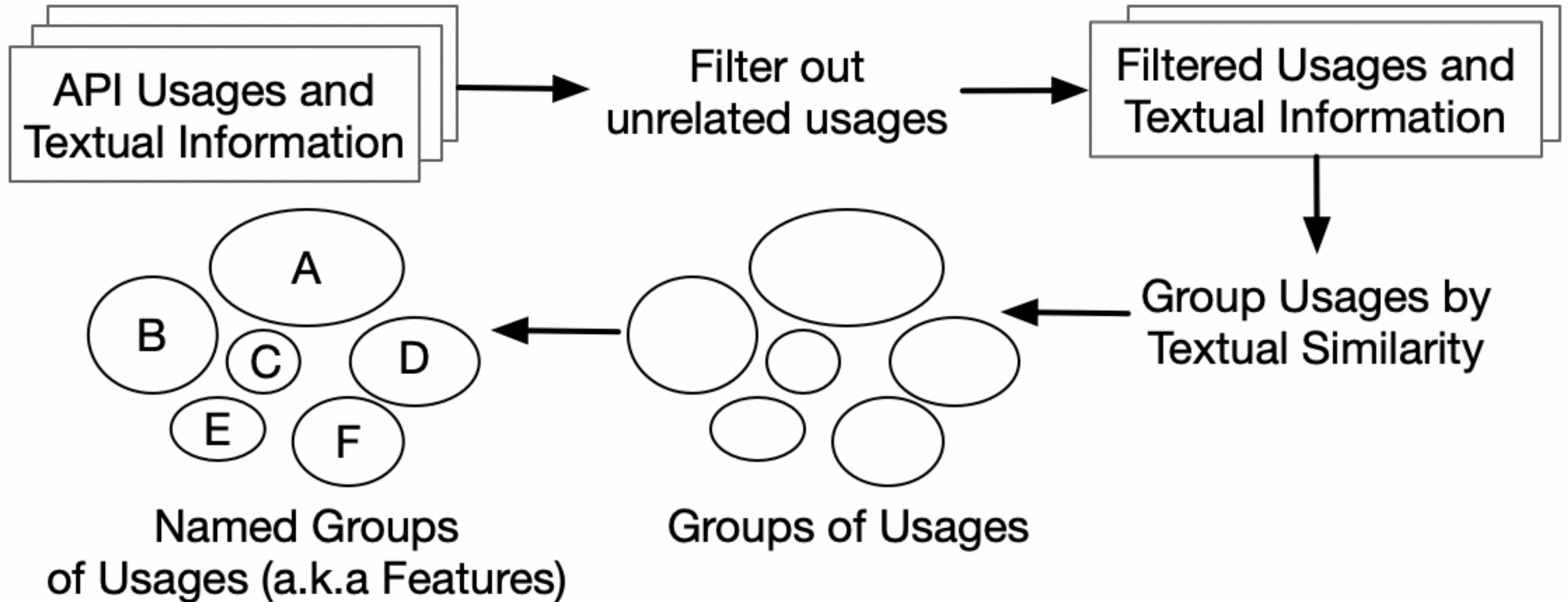
**Ranges**

**What is the simplest way to read a file into String?**

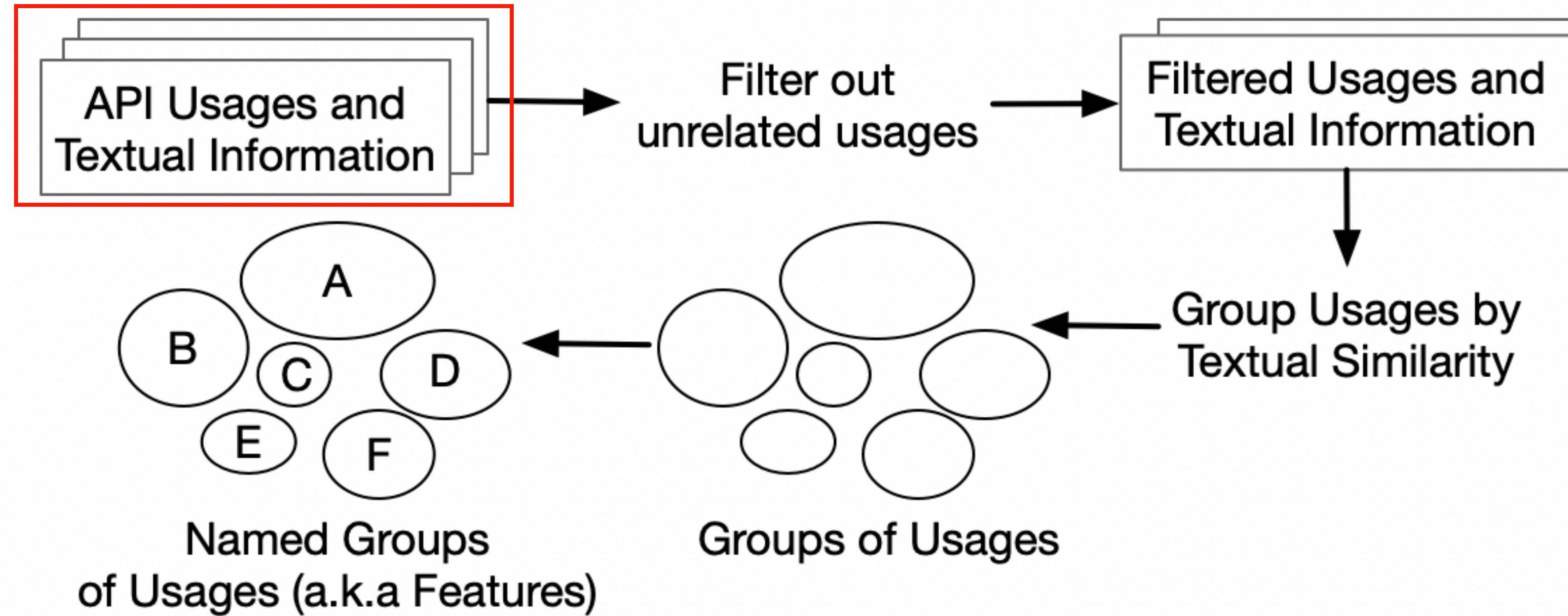
**Utils to read resource text file to String (Java)**

**I/O**

# Pipeline



# Extracting Usages

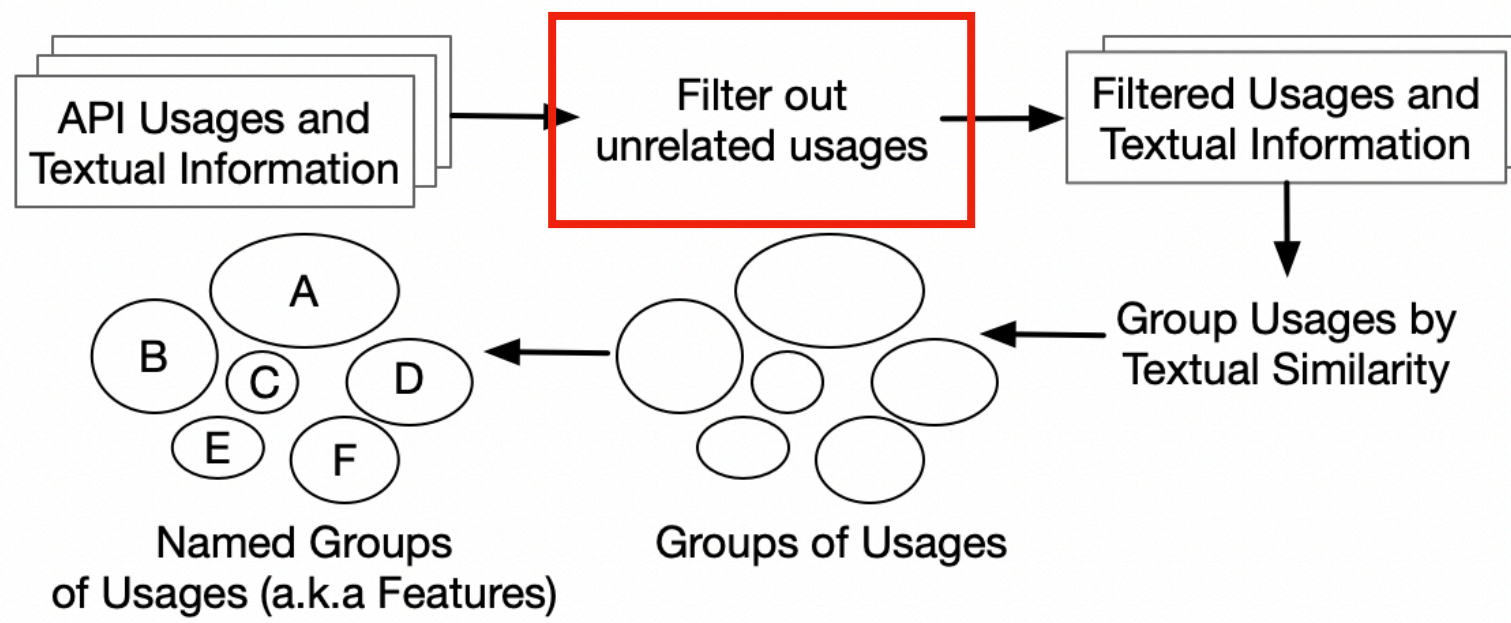


```
Gson gson = new Gson();  
PrintWriter out = new PrintWriter(  
String json = gson.toJson("{}")  
out.print(json)  
out.close();
```

```
[com.google.Gson: 0.91  
org.apache.http.NameValuePair: 0.04]  
Gson gson = new Gson();  
[java.io.PrintStream: 0.4  
java.io.BufferedReader: 0.4]  
PrintWriter out = new PrintWriter(  
[com.google.Gson: 0.98  
org.apache.http.NameValuePair: 0.01]  
String json = gson.toJson("{}")
```

```
[java.io.PrintStream: 0.95  
java.io.BufferedReader: 0.04]  
out.print(json)  
[java.io.PrintStream: 0.96  
java.io.BufferedReader: 0.03]  
out.close();
```

# Filtering



## Compute Covariance Matrix in Java

Here is a short example, how you can create it with Apache Commons Math (3.5):

10

```
RealMatrix mx = MatrixUtils.createRealMatrix(new double[][]{
    {1, 2, 3},
    {2, 4, 6}
});
RealMatrix cov = new Covariance(mx).getCovarianceMatrix();
```

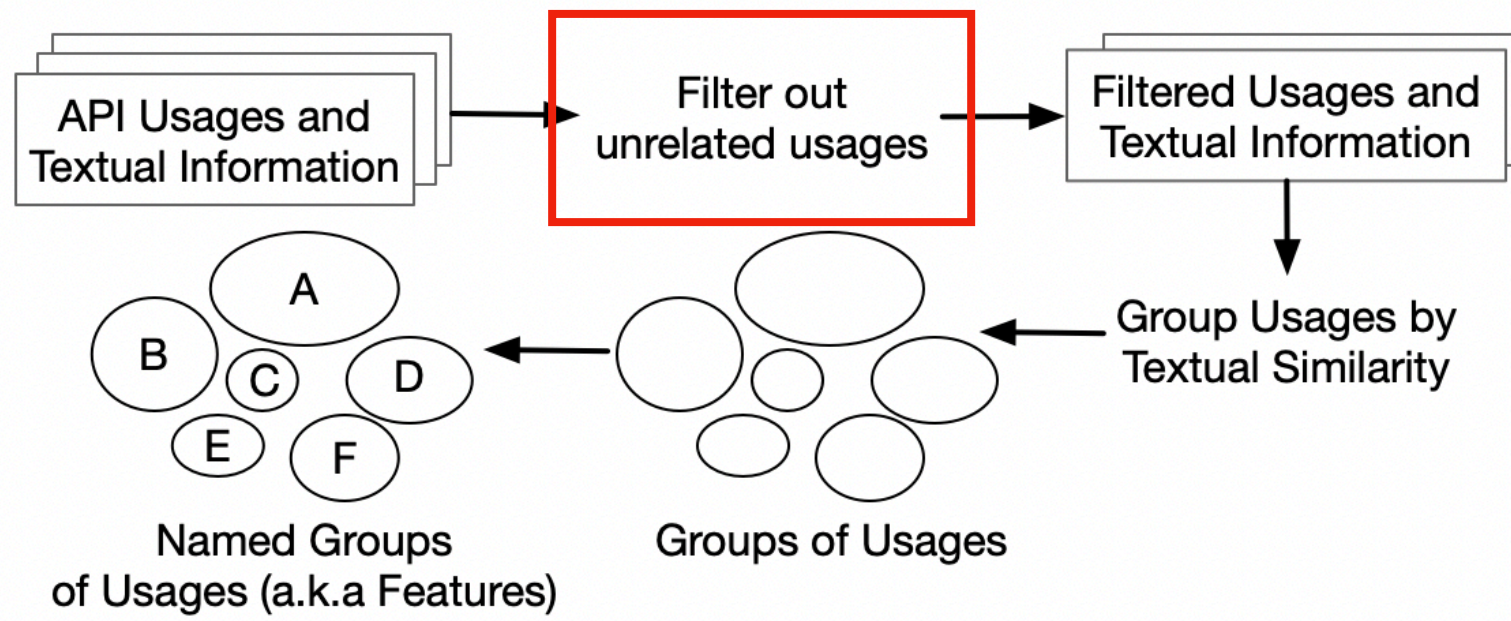
## How to average BigDecimals using Streams?

3

```
MutableDoubleList doubles = DoubleLists.mutable.with(1.0, 2.0, 3.0, 4.0);
List<BigDecimal> bigDecimals = doubles.collect(BigDecimal::new);
BigDecimal average =
    bigDecimals.stream()
        .collect(Collectors2.summarizingBigDecimal(e -> e))
        .getAverage(MathContext.DECIMAL32);

Assert.assertEquals(BigDecimal.valueOf(2.5), average);
```

# Filtering



## Compute **Covariance Matrix** in Java

Here is a short example, how you can create it with Apache Commons Math (3.5):

10

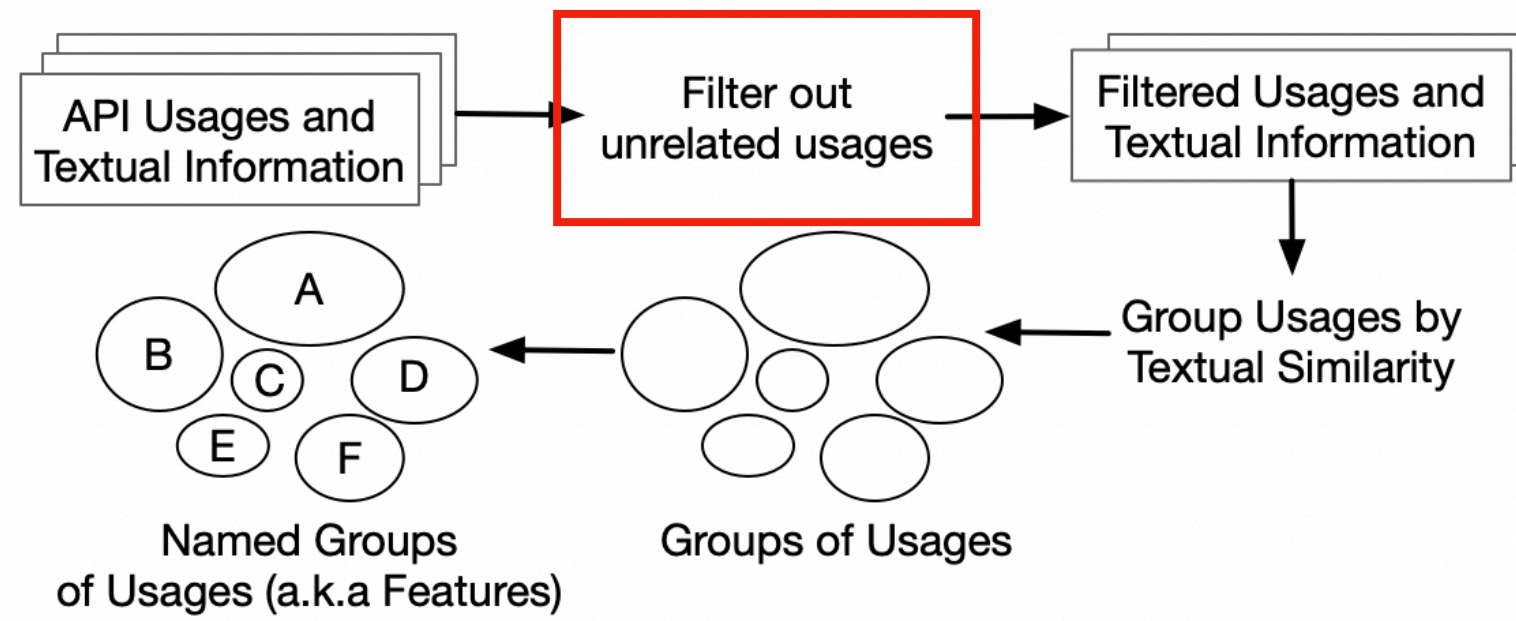
```
RealMatrix mx = MatrixUtils.createRealMatrix(new double[][]{  
    {1, 2, 3},  
    {2, 4, 6}  
});  
RealMatrix cov = new Covariance(mx).getCovarianceMatrix();
```

## How to **average BigDecimals** using **Streams**?

3

```
MutableDoubleList doubles = DoubleLists.mutable.with(1.0, 2.0, 3.0, 4.0);  
List<BigDecimal> bigDecimals = doubles.collect(BigDecimal::new);  
BigDecimal average =  
    bigDecimals.stream()  
        .collect(Collectors2.summarizingBigDecimal(e -> e))  
        .getAverage(MathContext.DECIMAL32);  
Assert.assertEquals(BigDecimal.valueOf(2.5), average);
```

# Filtering



1. Process post titles.

2. Process API usage names.

3. Load prebuilt embeddings model.

4. Compute semantic similarity.

## Compute Covariance Matrix in Java

Here is a short example, how you can create it with Apache Commons Math (3.5):

10

```
RealMatrix mx = MatrixUtils.createRealMatrix(new double[][]{  
    {1, 2, 3},  
    {2, 4, 6}  
});  
RealMatrix cov = new Covariance(mx).getCovarianceMatrix();
```

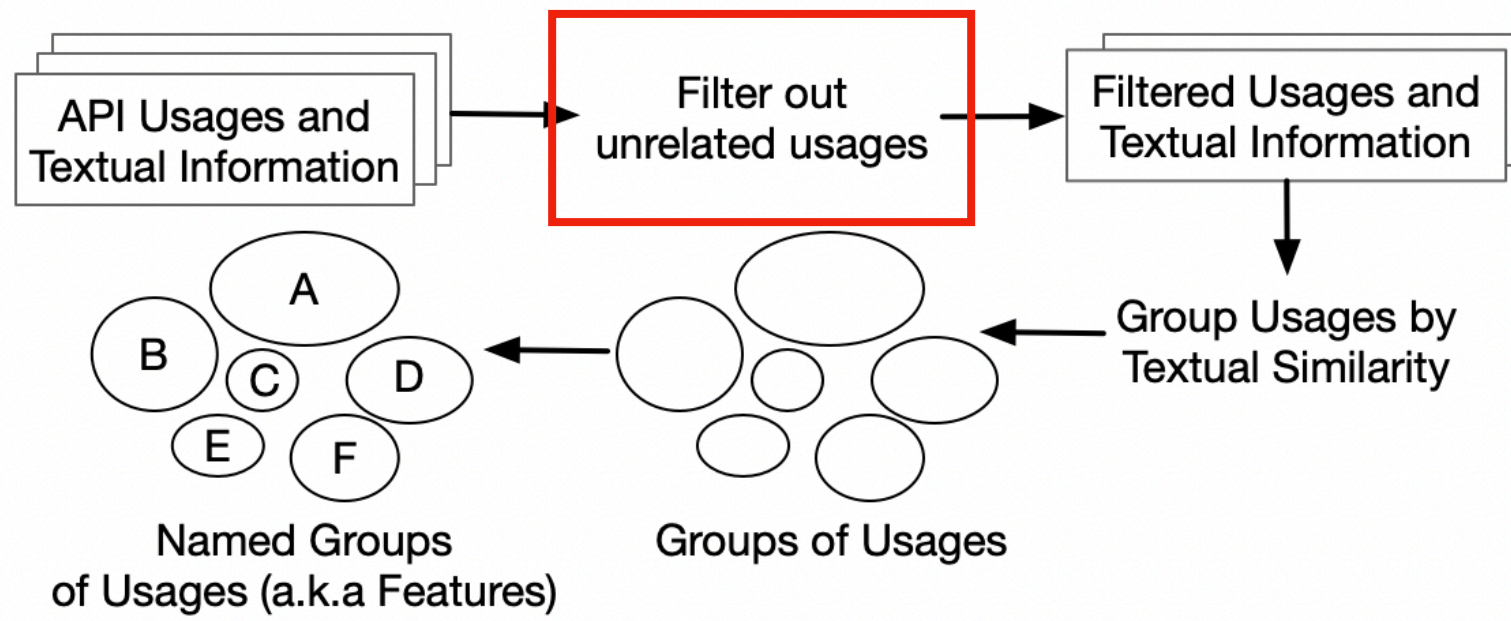
## How to average BigDecimals using Streams?

3

```
MutableDoubleList doubles = DoubleLists.mutable.with(1.0, 2.0, 3.0, 4.0);  
List<BigDecimal> bigDecimals = doubles.collect(BigDecimal::new);  
BigDecimal average =  
    bigDecimals.stream()  
        .collect(Collectors2.summarizingBigDecimal(e -> e))  
        .getAverage(MathContext.DECIMAL32);  
  
Assert.assertEquals(BigDecimal.valueOf(2.5), average);
```



# Filtering



1. Process post titles.
2. Process API usage names.
3. Load prebuilt embeddings model.
4. Compute semantic similarity.

0.86

## Compute Covariance Matrix in Java

Here is a short example, how you can create it with Apache Commons Math (3.5):

```
RealMatrix mx = MatrixUtils.createRealMatrix(new double[][]{{1, 2, 3}, {2, 4, 6}});  
RealMatrix cov = new Covariance(mx).getCovarianceMatrix();
```

10

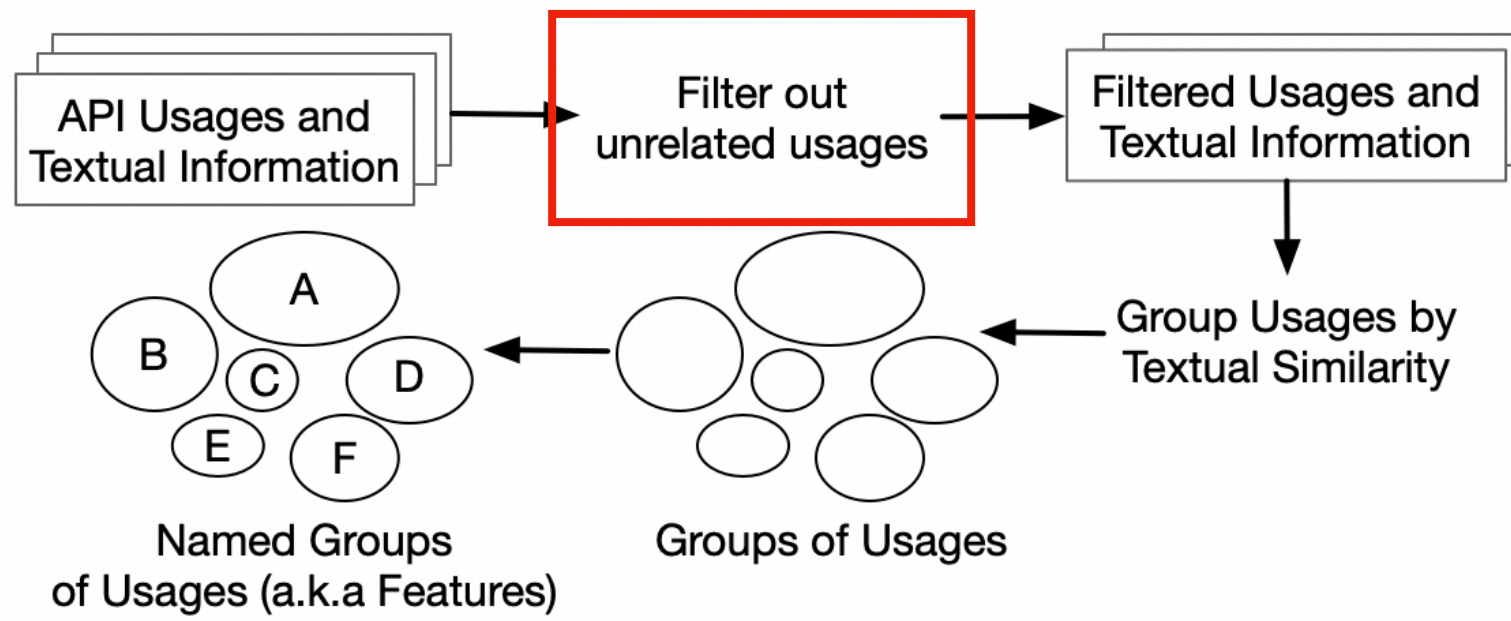
0.71

## How to average BigDecimals using Streams?

```
MutableDoubleList doubles = DoubleLists.mutable.with(1.0, 2.0, 3.0, 4.0);  
List<BigDecimal> bigDecimals = doubles.collect(BigDecimal::new);  
BigDecimal average =  
    bigDecimals.stream()  
        .collect(Collectors2.summarizingBigDecimal(e -> e))  
        .getAverage(MathContext.DECIMAL32);  
Assert.assertEquals(BigDecimal.valueOf(2.5), average);
```

3

# Filtering



1. Process post titles.
2. Process API usage names.
3. Load prebuilt embeddings model.
4. Compute semantic similarity.

0.86

## Compute Covariance Matrix in Java

Here is a short example, how you can create it with Apache Commons Math (3.5):

```
RealMatrix mx = MatrixUtils.createRealMatrix(new double[][]{{1, 2, 3}, {2, 4, 6}});  
RealMatrix cov = new Covariance(mx).getCovarianceMatrix();
```

10

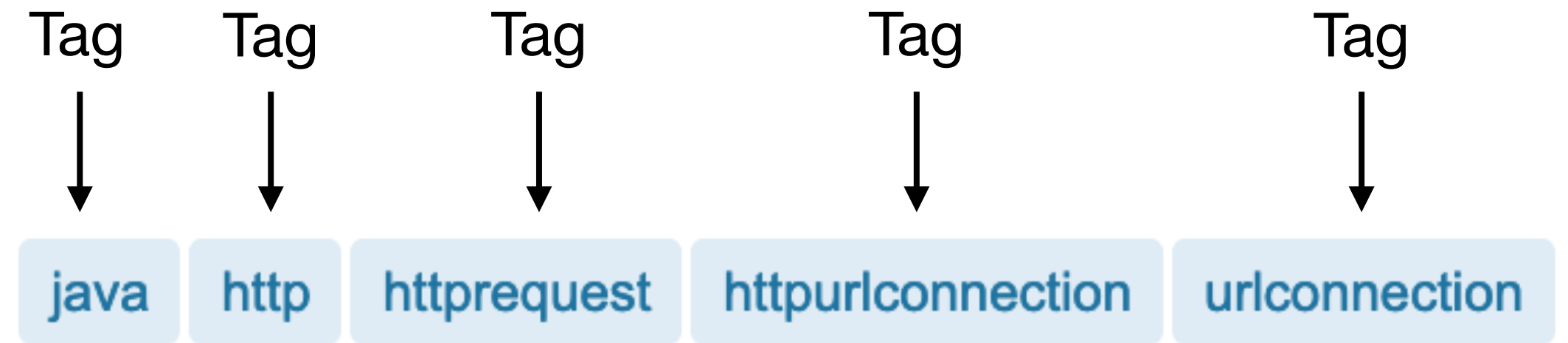
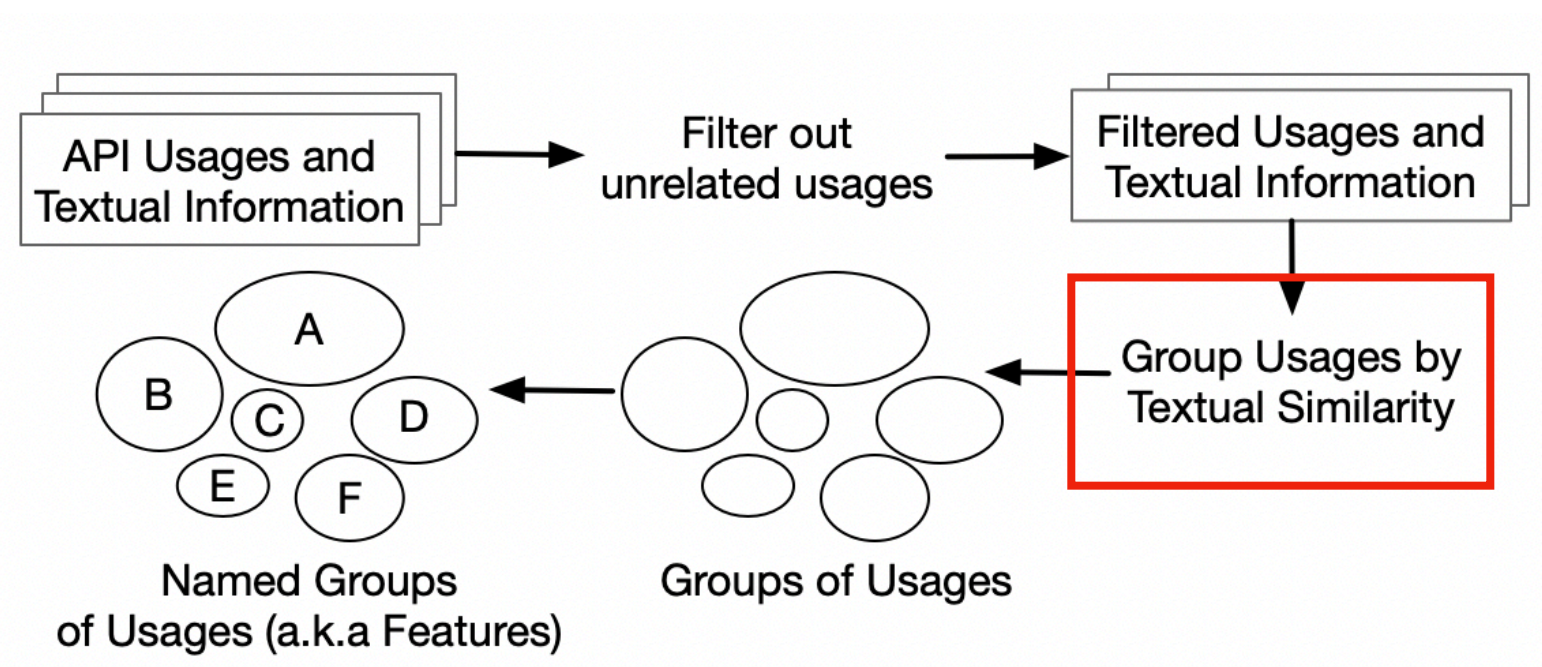
0.03

## How to average BigDecimals using Streams?

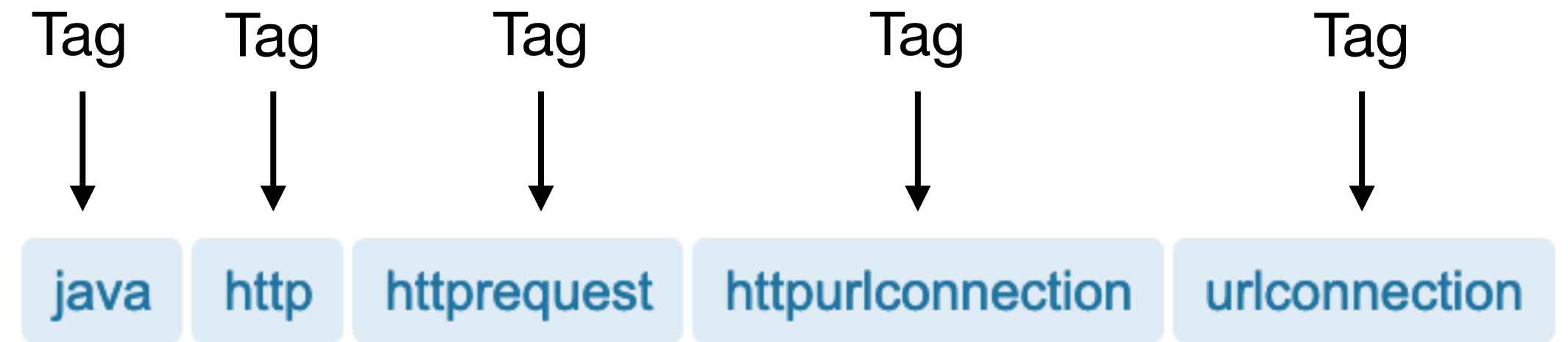
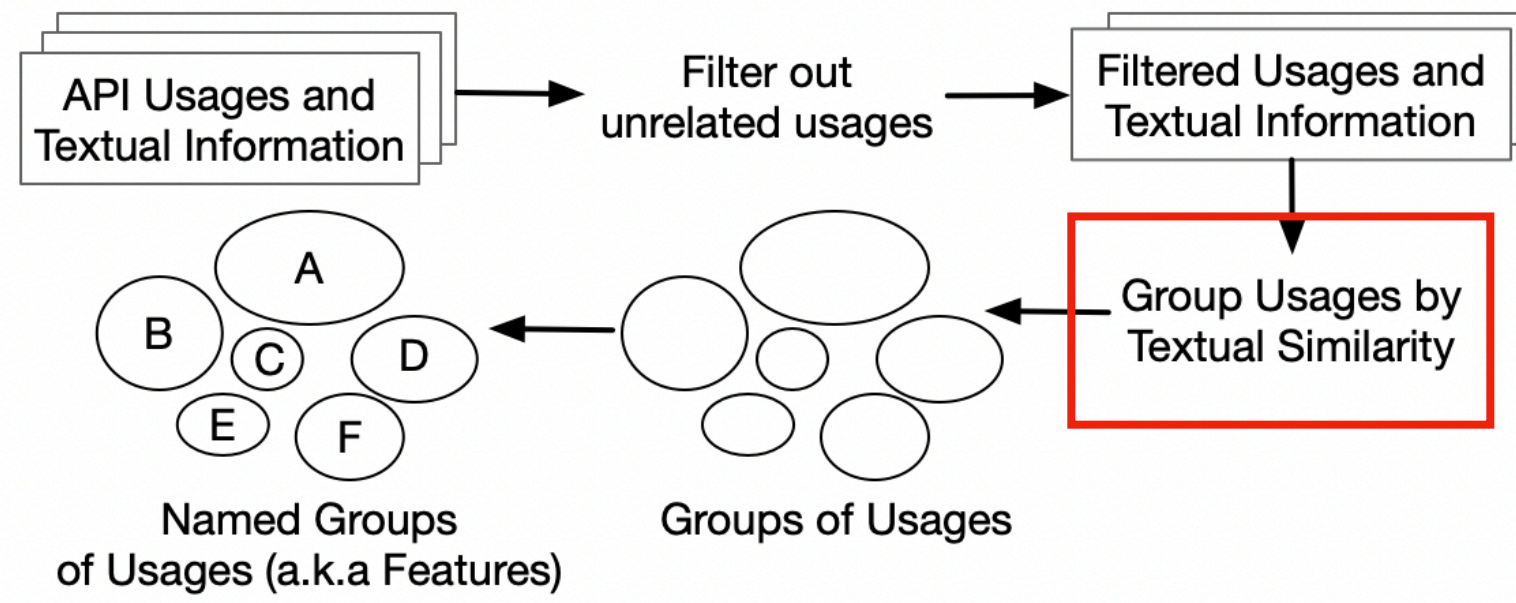
```
MutableDoubleList doubles = DoubleLists.mutable.with(1.0, 2.0, 3.0, 4.0);  
List<BigDecimal> bigDecimals = doubles.collect(BigDecimal::new);  
BigDecimal average =  
    bigDecimals.stream()  
        .collect(Collectors2.summarizingBigDecimal(e -> e))  
        .getAverage(MathContext.DECIMAL32);  
  
Assert.assertEquals(BigDecimal.valueOf(2.5), average);
```

3

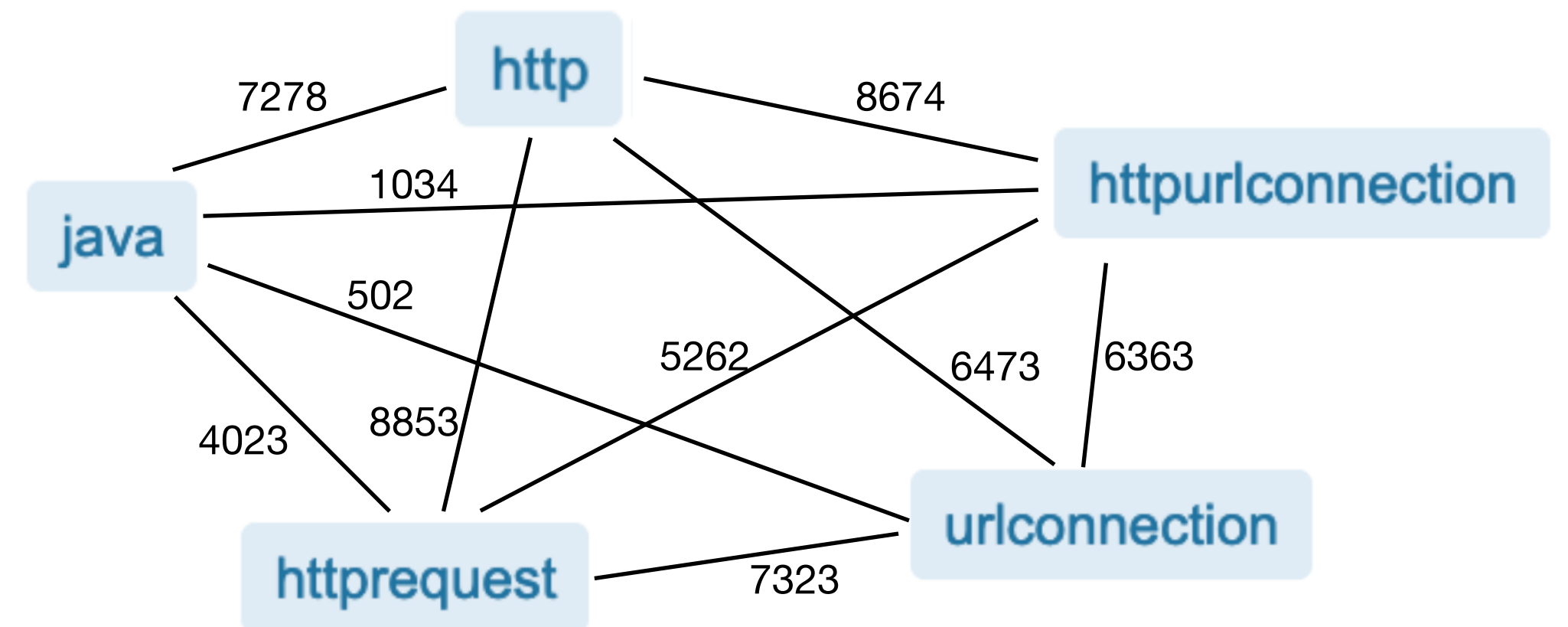
# Grouping



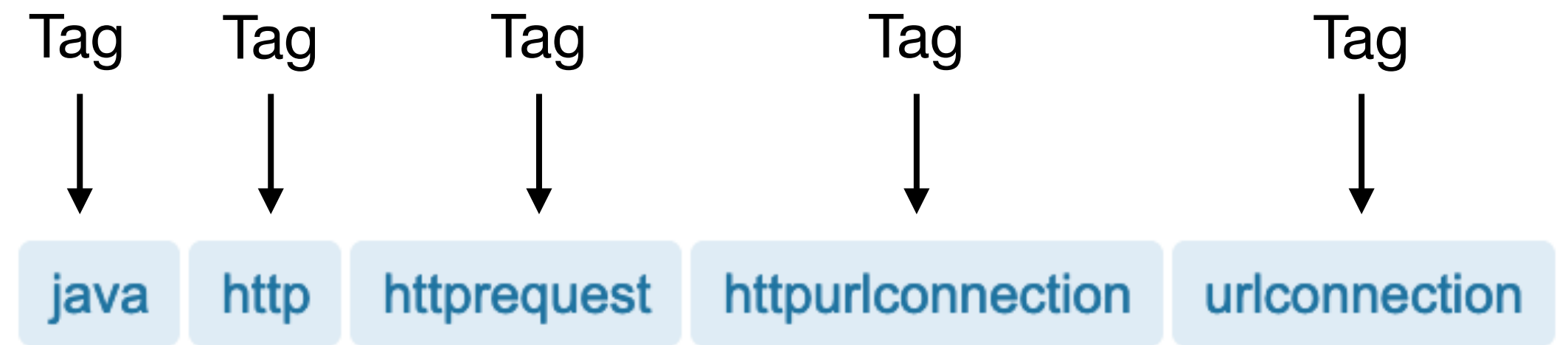
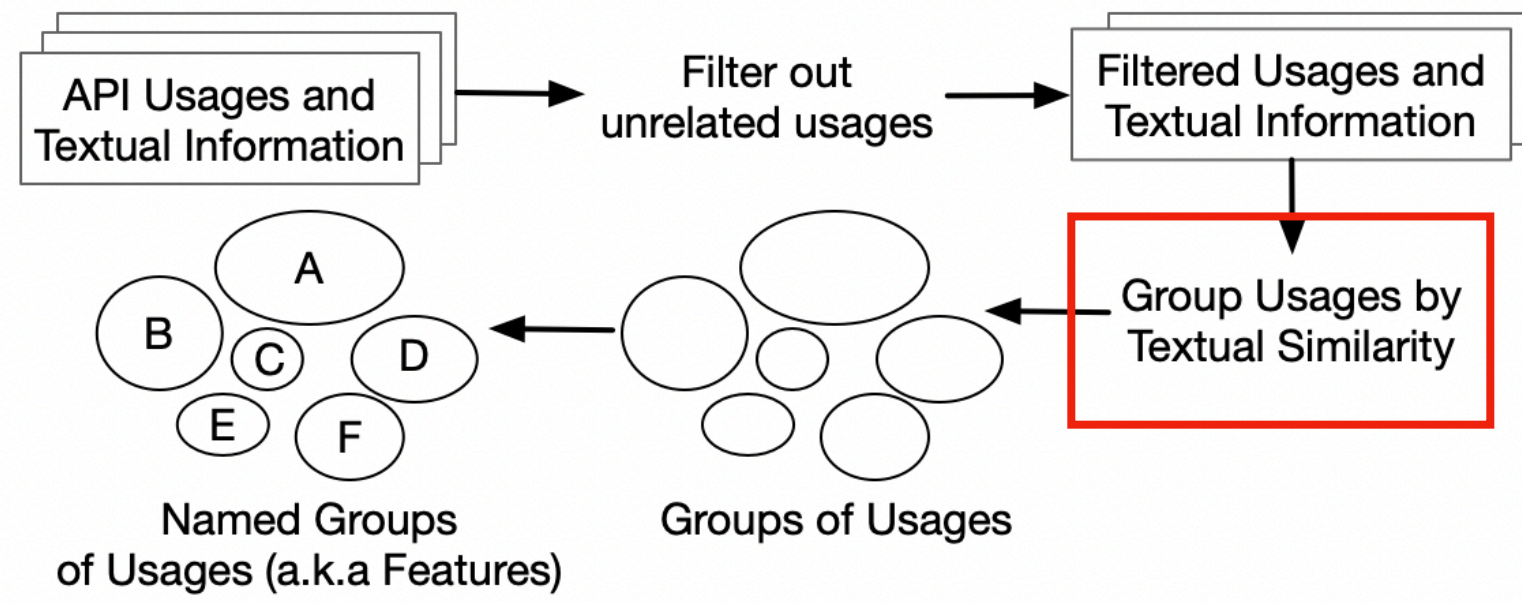
# Grouping



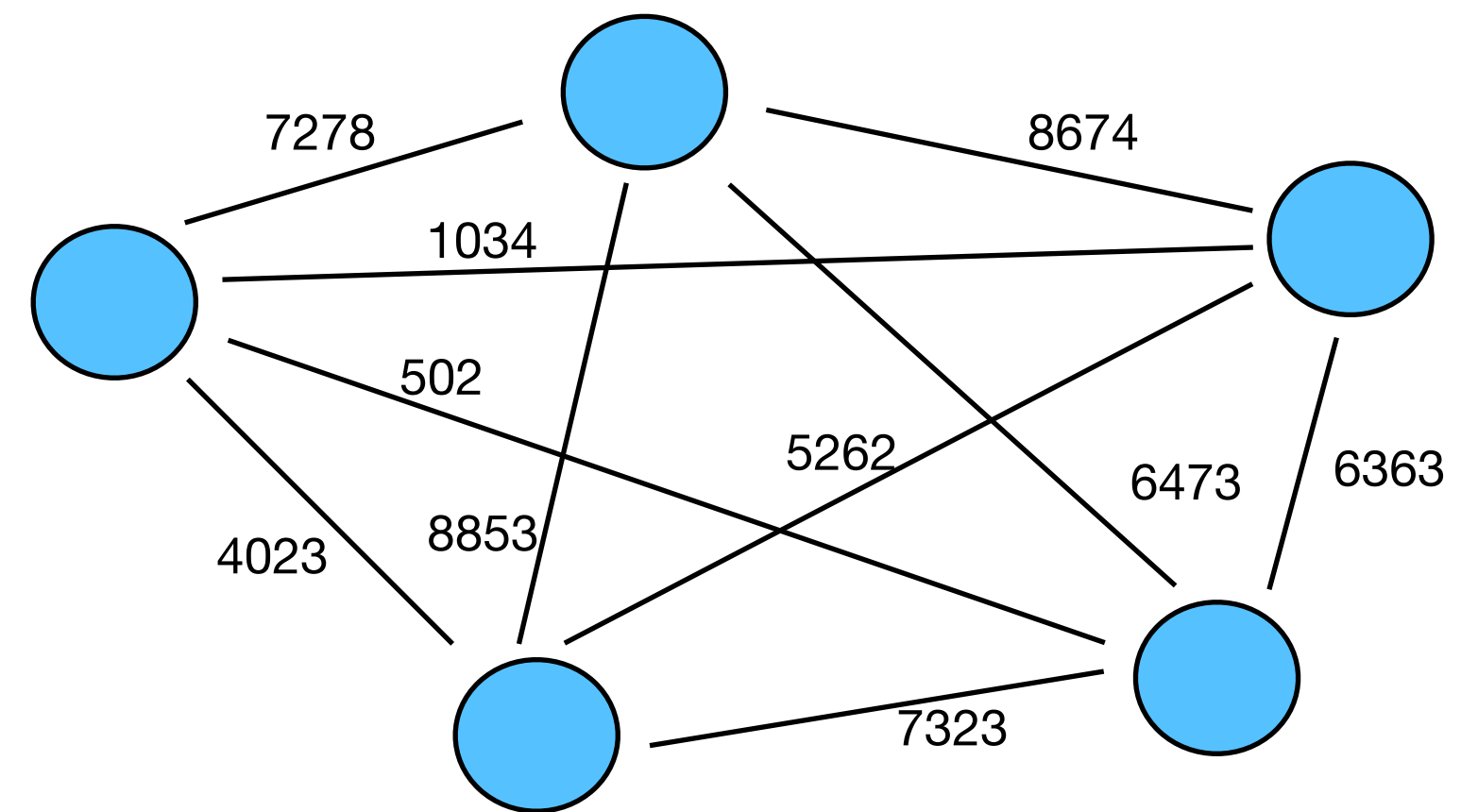
## Relation between Tags



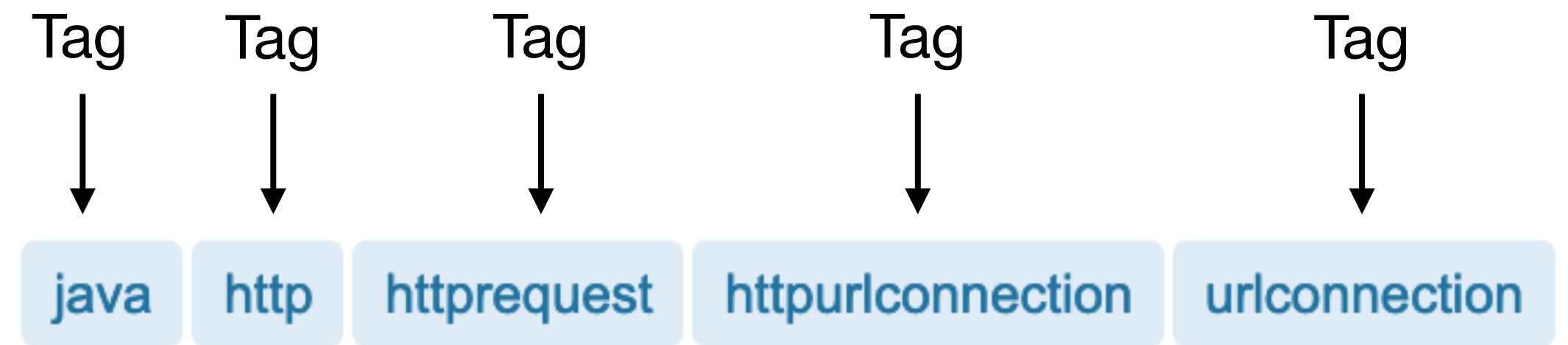
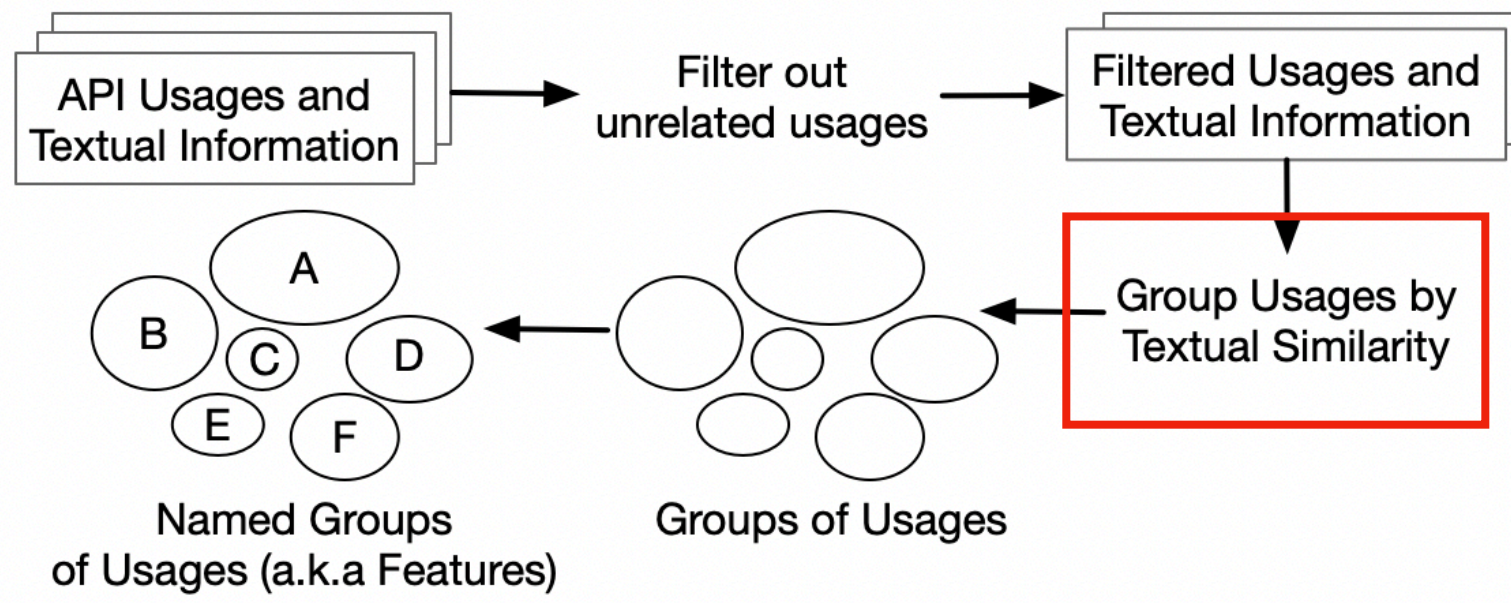
# Grouping



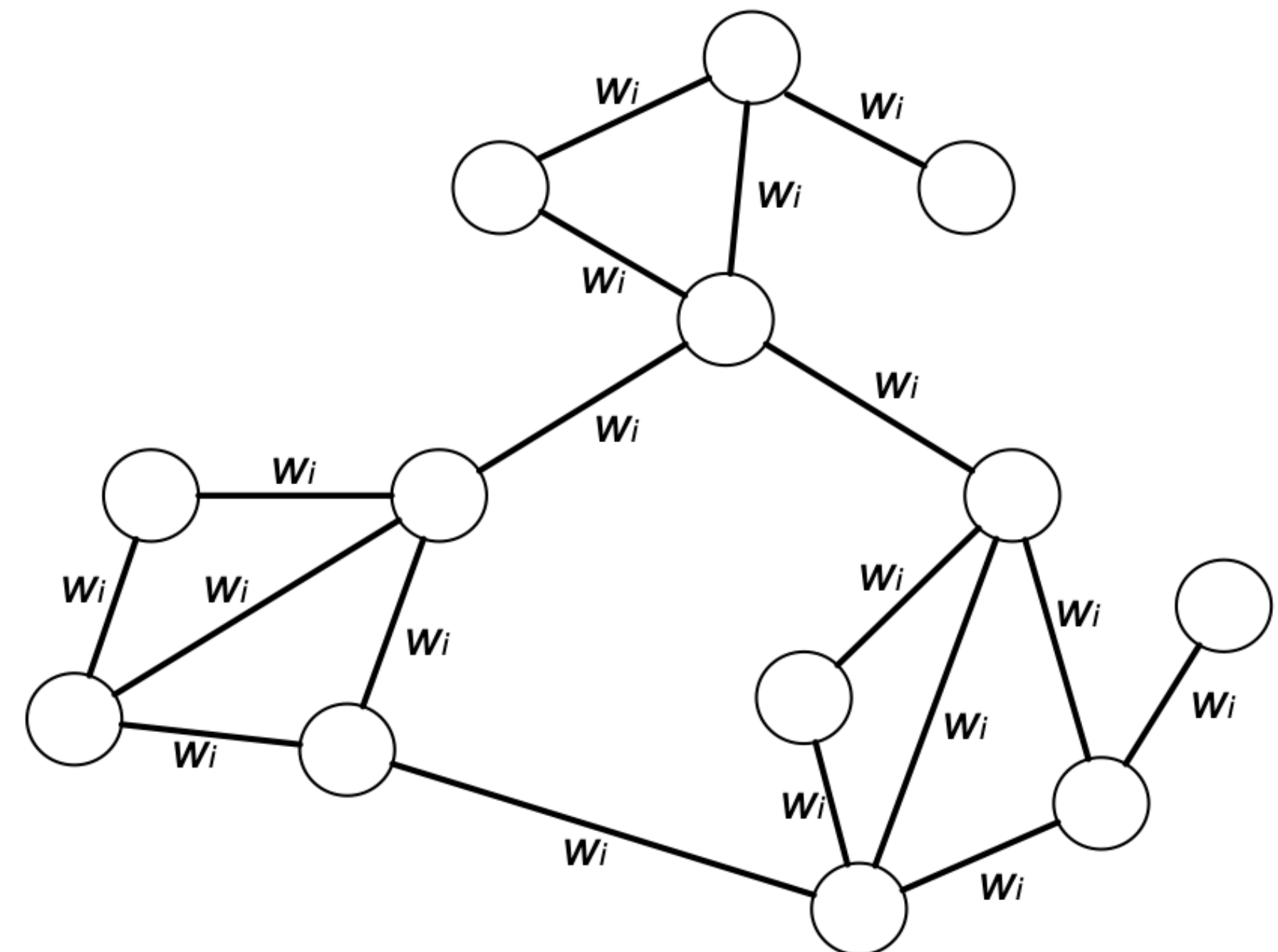
**Relation between Tags**



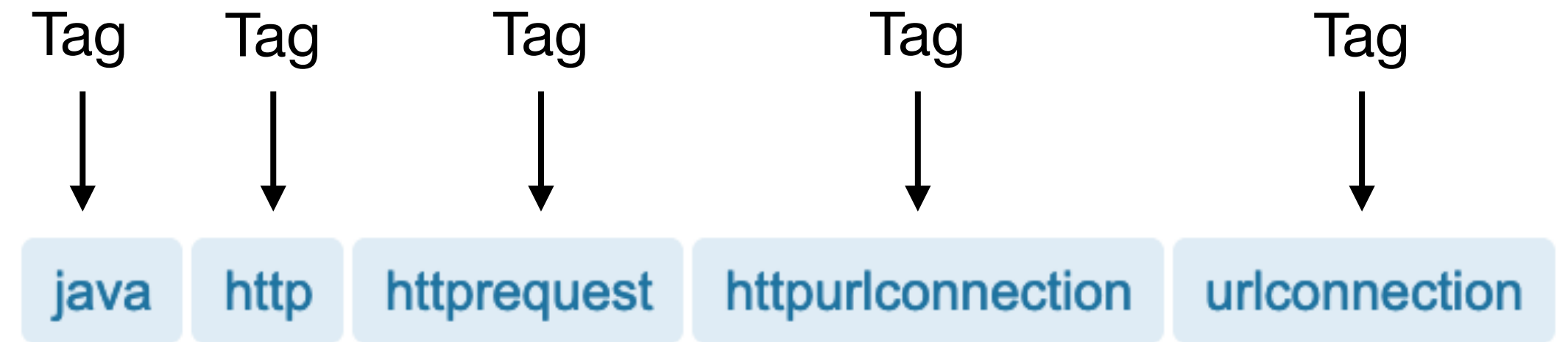
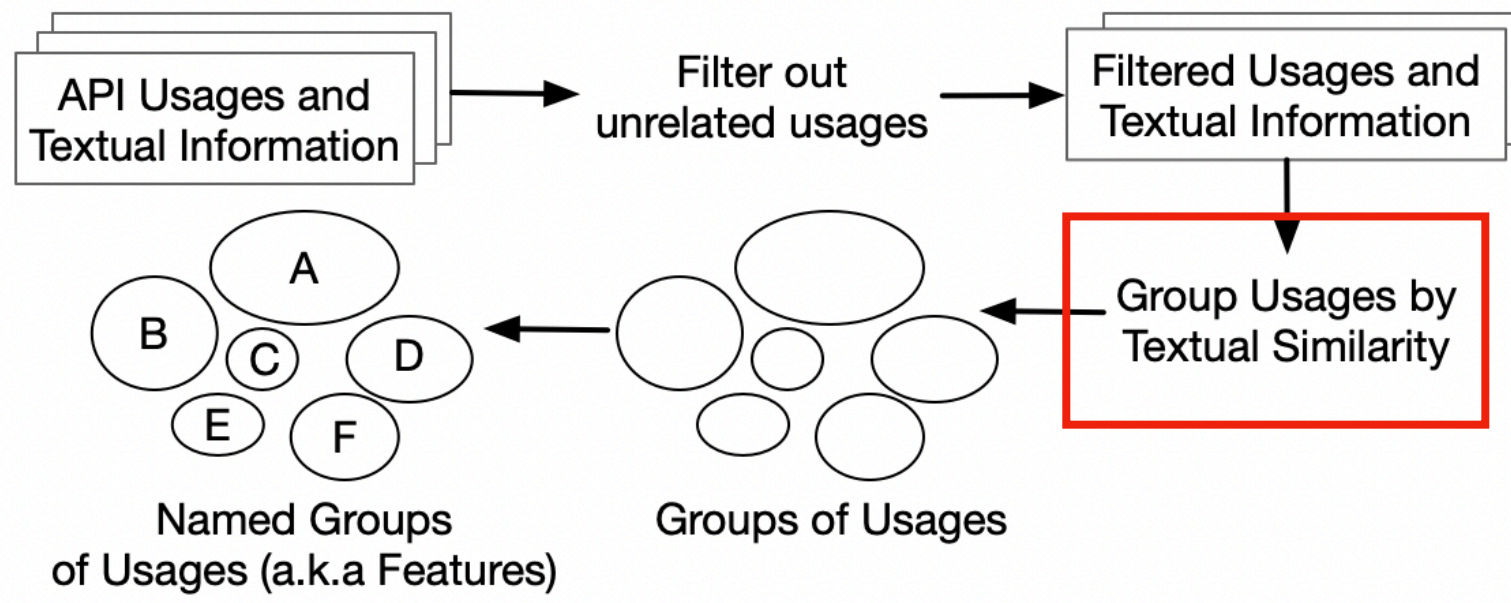
# Grouping



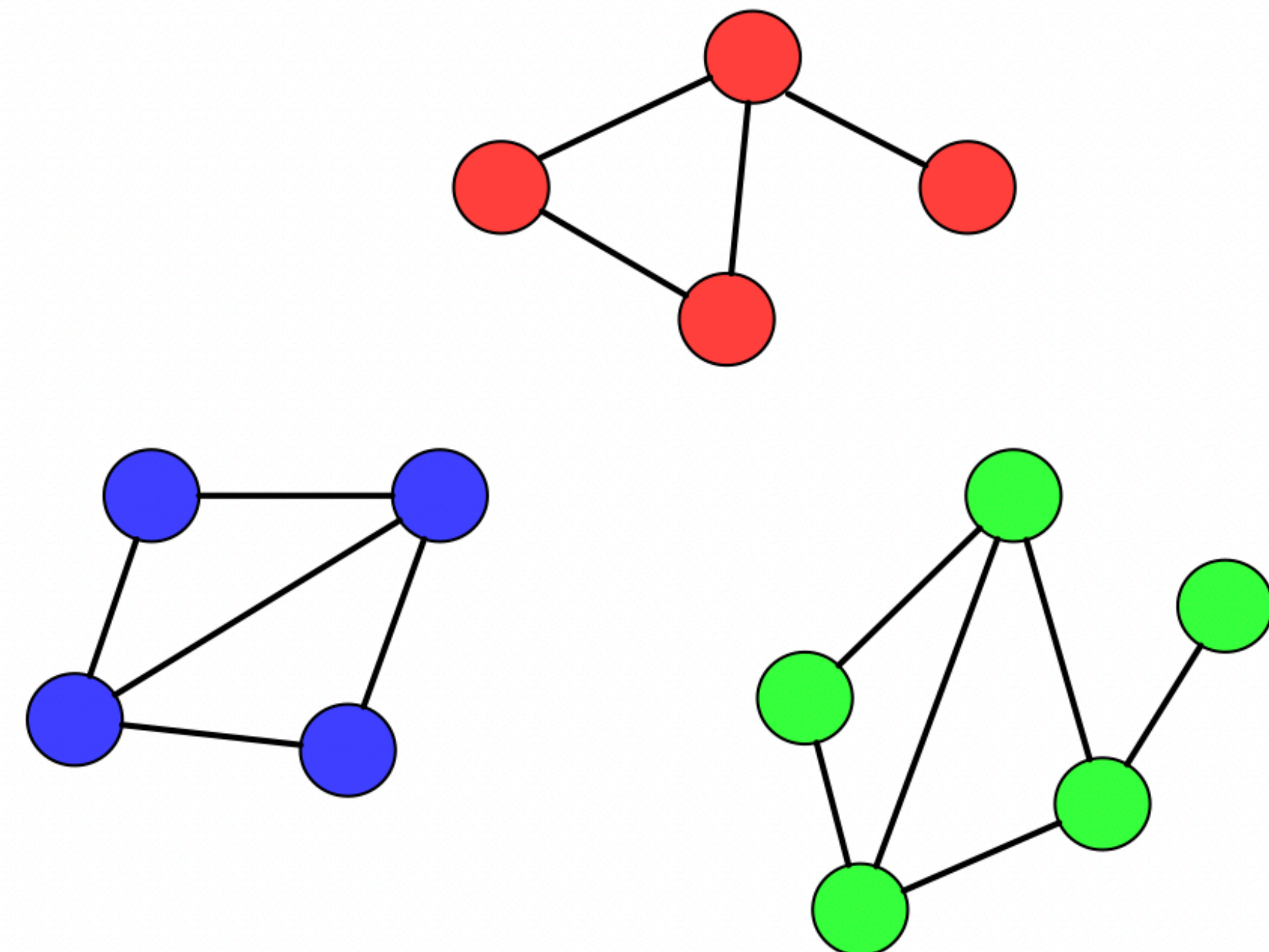
## Relation between Tags



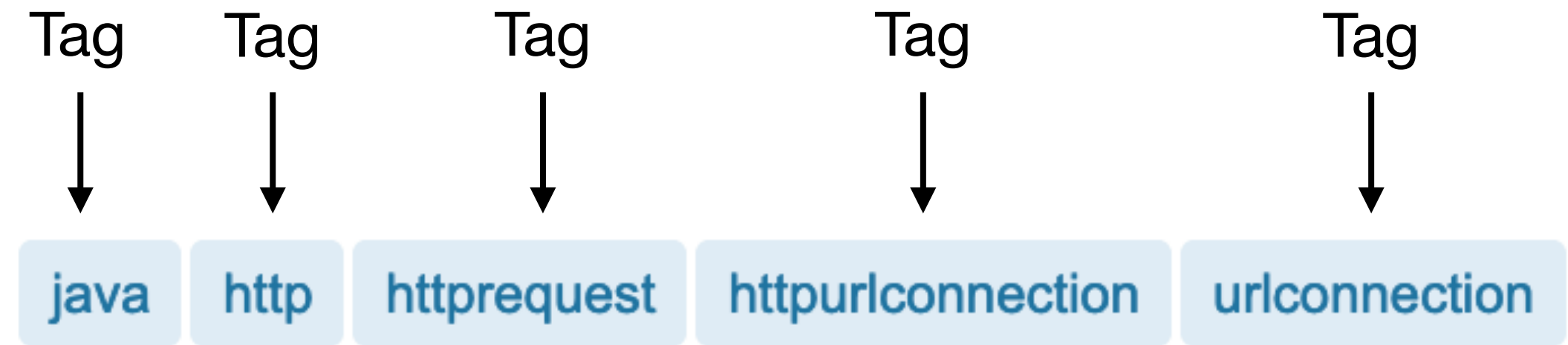
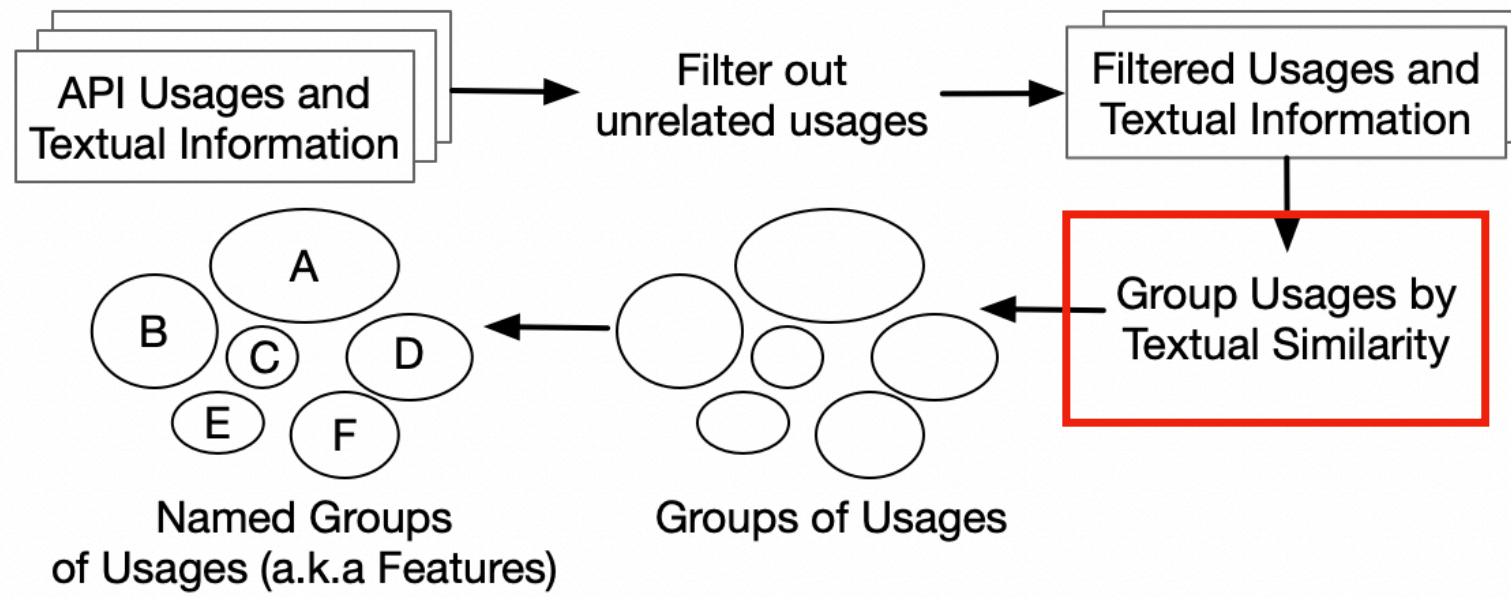
# Grouping



## Relation between Tags



# Grouping



## Identifying Features in Forks

Shurui Zhou  
Carnegie Mellon University

Ștefan Stănculescu  
IT University of Copenhagen

Olaf Leßenich  
University of Passau

Yingfei Xiong  
Peking University

Andrzej Wąsowski  
IT University of Copenhagen

Christian Kästner  
Carnegie Mellon University

### ABSTRACT

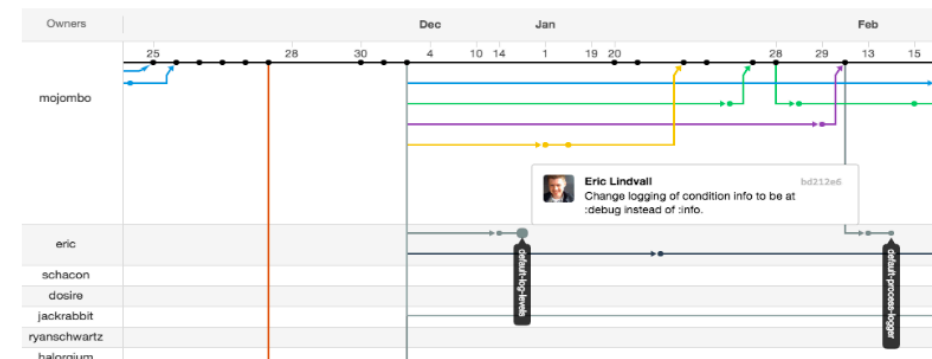
Fork-based development has been widely used both in open source communities and in industry, because it gives developers flexibility to modify their own fork without affecting others. Unfortunately, this mechanism has downsides: When the number of forks becomes large, it is difficult for developers to get or maintain an overview of activities in the forks. Current tools provide little help. We introduce INFOX, an approach to automatically identify non-merged features in forks and to generate an overview of active forks in a project. The approach clusters cohesive code fragments using code and network-analysis techniques and uses information-retrieval techniques to label clusters with keywords. The clustering is effective, with 90% accuracy on a set of known features. In addition, a human-subject evaluation shows that INFOX can provide actionable insight for developers of forks.

### ACM Reference Format:

Shurui Zhou, Ștefan Stănculescu, Olaf Leßenich, Yingfei Xiong, Andrzej Wąsowski, and Christian Kästner. 2018. Identifying Features in Forks. In *ICSE '18: 40th International Conference on Software Engineering*, May 27–June 3, 2018, Gothenburg, Sweden. ACM, New York, NY, USA, 12 pages. <https://doi.org/10.1145/3180155.3180205>

### 1 INTRODUCTION

Forking is a lightweight and easy mechanism that allows developers, both in open source and in industry, to start development from



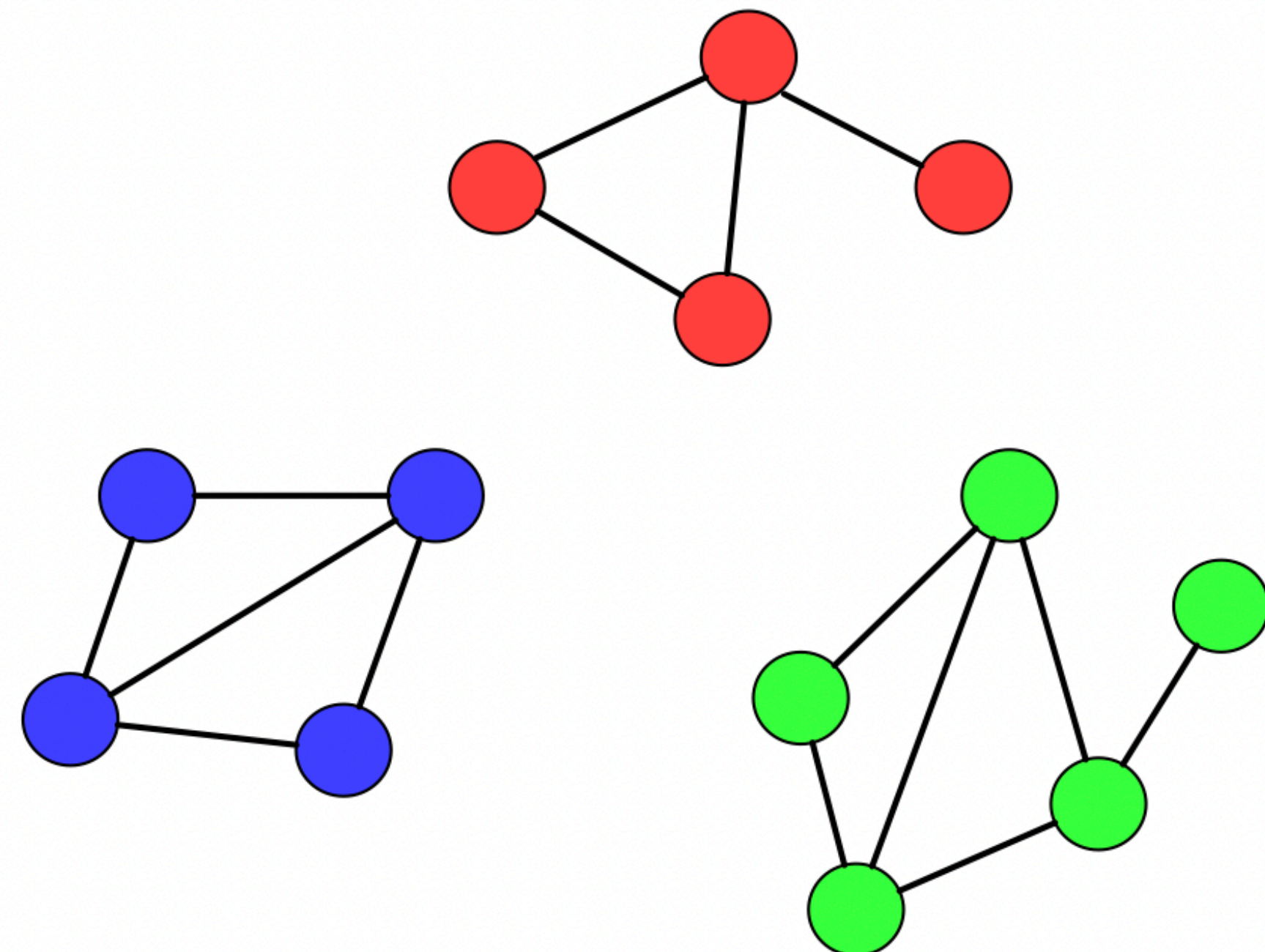
(a) GitHub network graph

Cluster	OpenBuilds/Smoothieware, last commit: 16 days ago	LOC
com.	[com, substr, smoothie, setcursor, build, lcd, c_str, openbuilds]	46
lcd.	[lcdwidth, shift, glcd_font, fb_size, framebuffer, memory, is_sh1106, size_sh1106, data, framebuffer]	10
con.	[contrast, apex, is_sh1106, support, contrast, reversed, oled, ol, sh1106]	5
ol.c.	[ol_checksum, new, variant, checksum, sh1106_ol, cksm, lcd_cksm, support, lcd, st7565]	2
Cluster	lciscon/Smoothieware, last commit: Jan 27	LOC
bed.	[bed, gcode, div, ha_letter, home_offset, gcode_receiv, nullstream, number, correct_checksum, st]	28
Cluster	arhi/Smoothieware, last commit: Jan 3	LOC
amp.	[ampmod1_pin, adc_valu, ampmod2_pin, thekernel, value, by_default, name_checksum, as_numbe]	92
add.	[added, example, config, function, pt100, class, streamoutput, snippet]	1

(b) INFOX overview

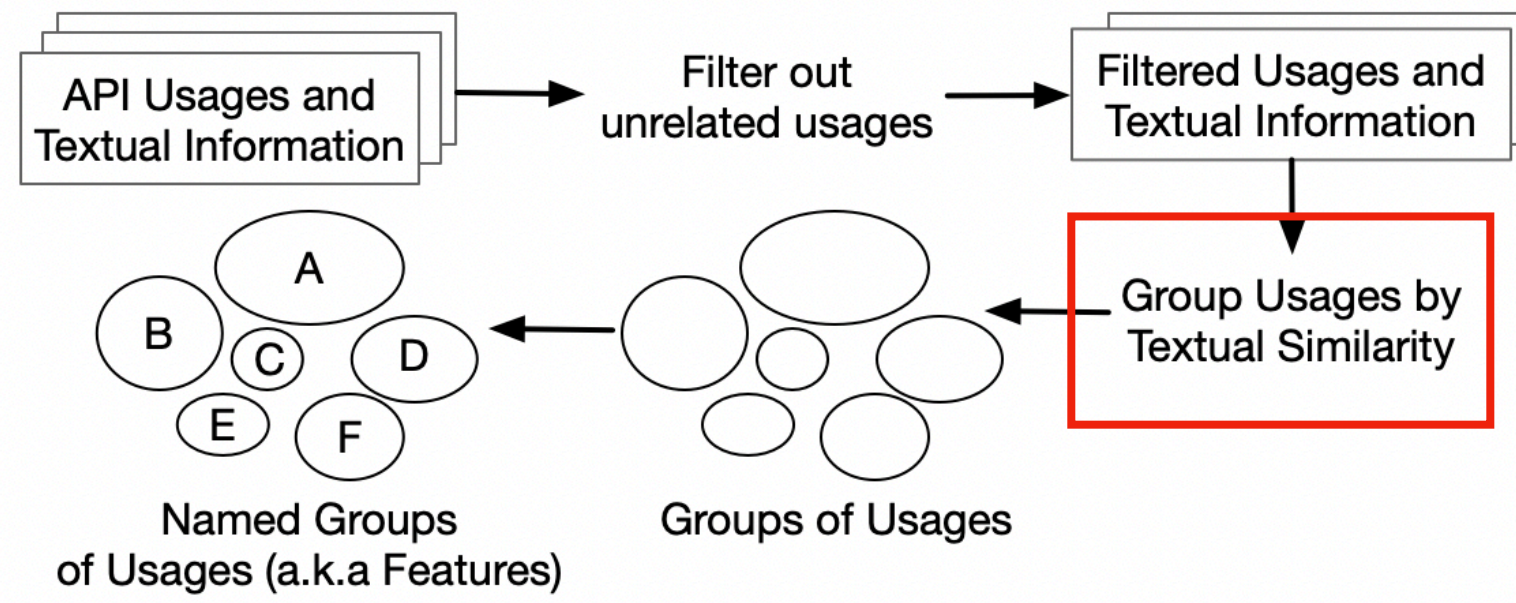
Figure 1: GitHub's network graph shows commits across known forks, but is difficult to use to gain an overview of activities in projects with many forks. INFOX's overview summarizes features in active forks.

## Relation between Tags





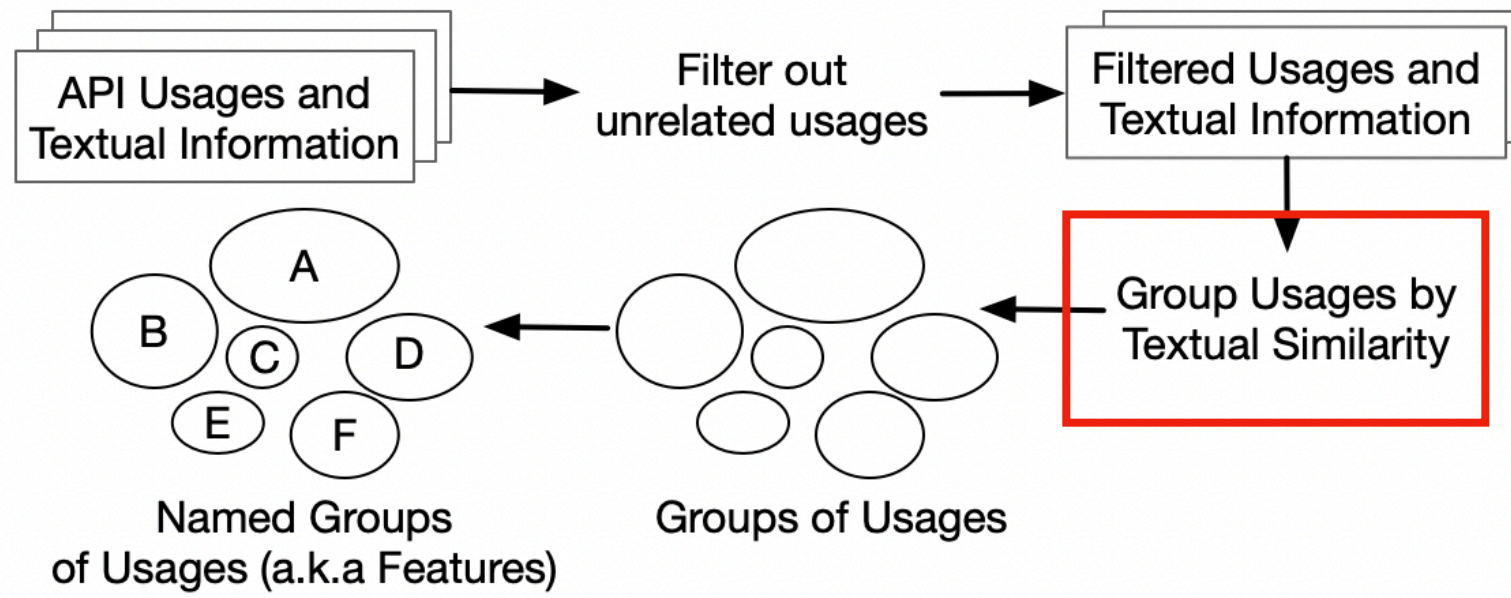
# Groups



Library: [org.apache.commons.math](https://org.apache.commons.math)

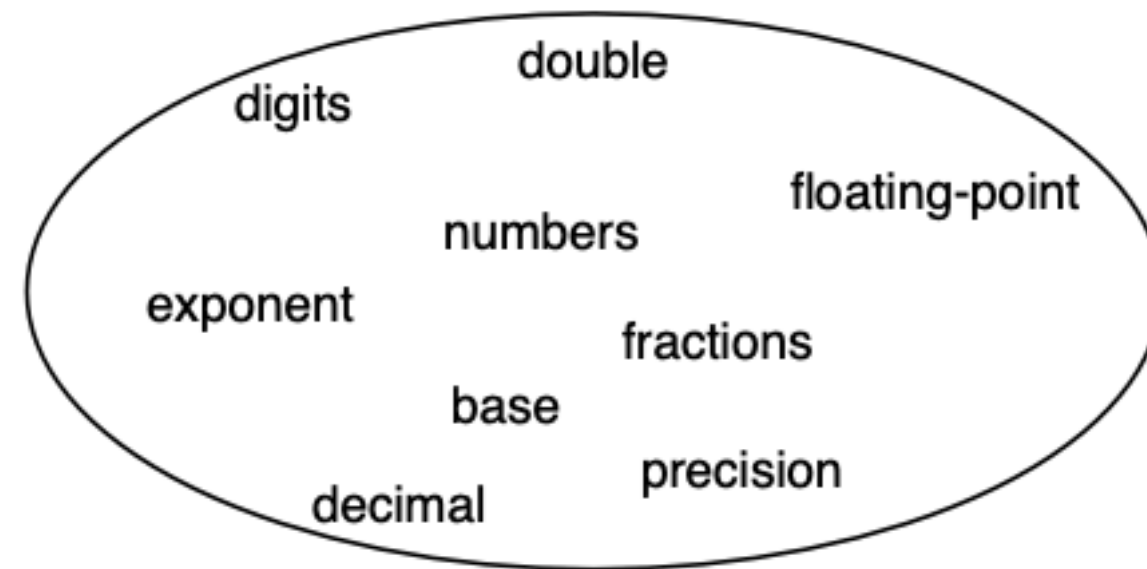
Mathematics Library from the Apache Foundation

# Groups

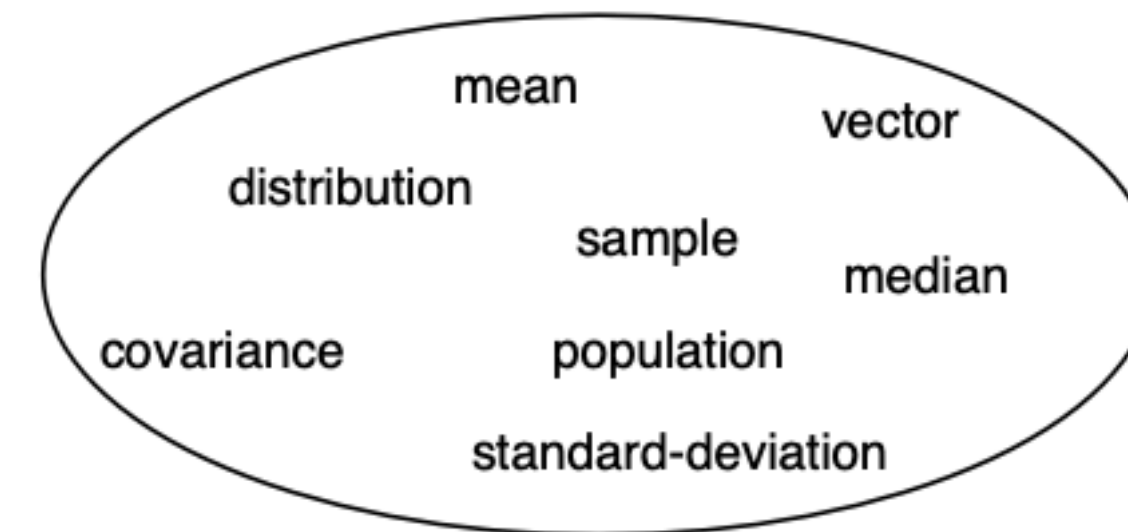


Library: `org.apache.commons.math`

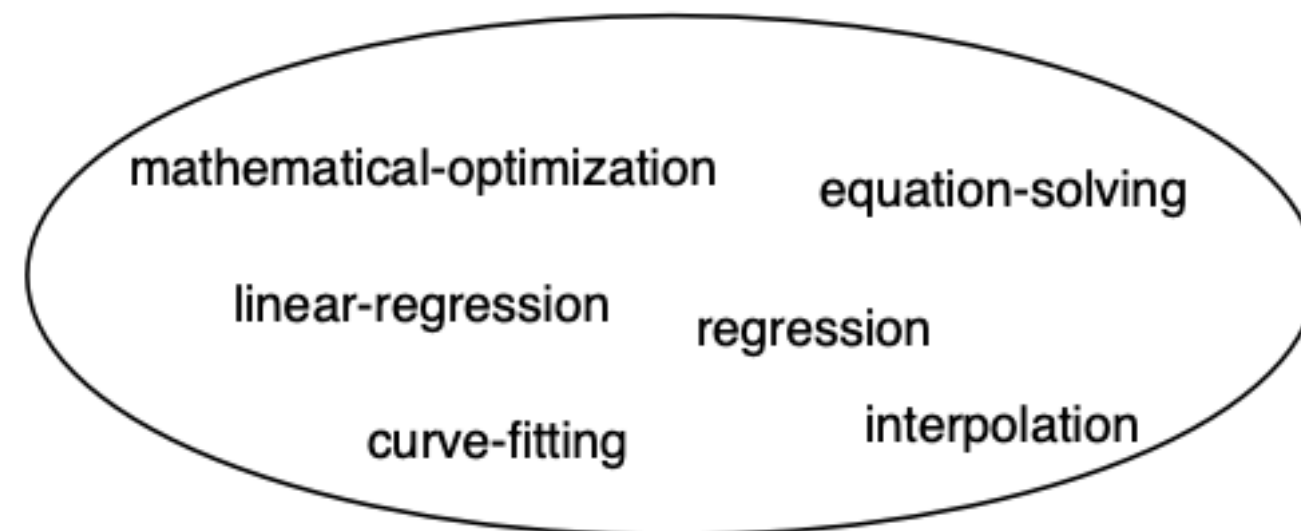
Mathematics Library from the Apache Foundation



`ComplexFormat.parse`  
`Complex.getReal`  
`Complex.getImaginary`  
`BigFraction.getNumerator`  
`BigFraction.getDenominator`  
`Precision.round`

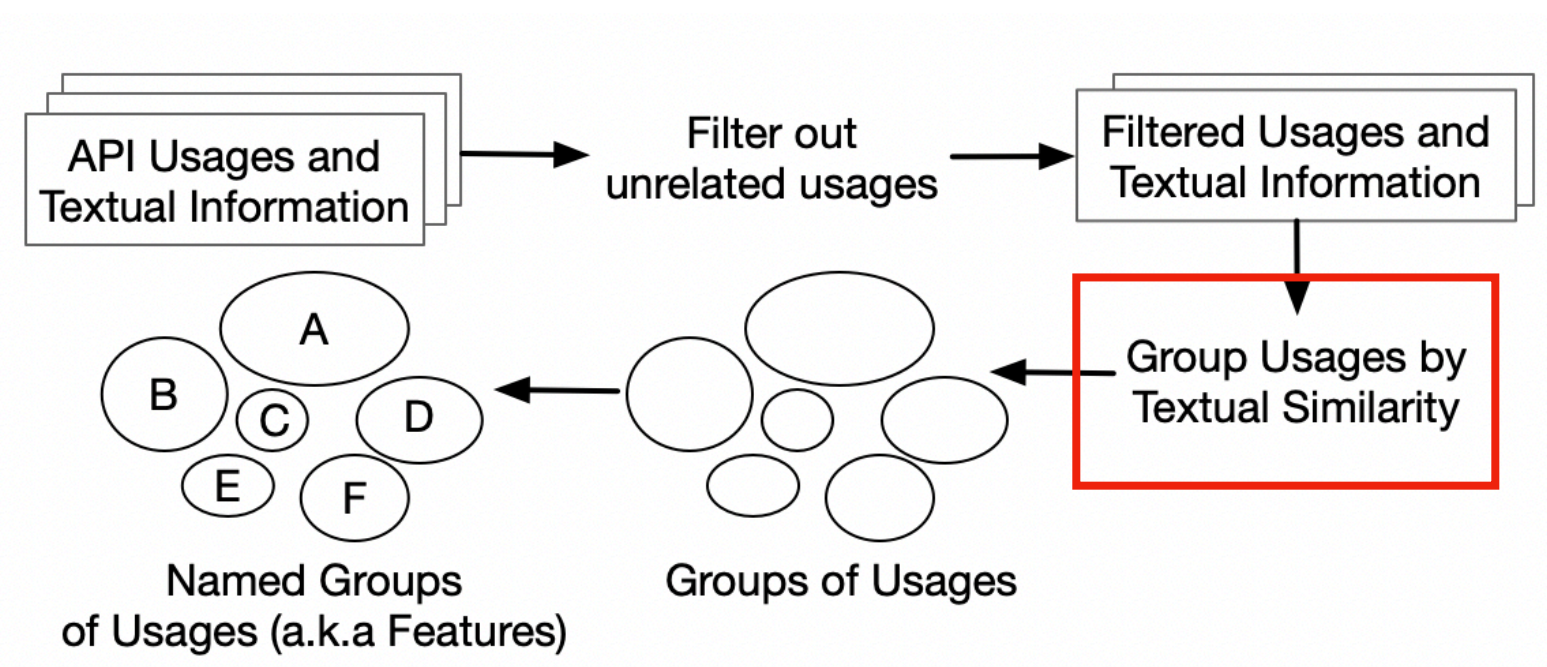


`SummaryStatistics.getMean`  
`SummaryStatistics.addValue`  
`SummaryStatistics.getStandardDeviation`  
`NormalDistribution.cumulativeProbability`  
`NormalDistribution.inverseCumulativeProbability`  
`LogNormalDistribution.sample`



`PolynomialSplineFunction.getPolynomials`  
`PolynomialSplineFunction.value`  
`LinearInterpolator.interpolate`  
`PointValuePair.getPoint`  
`SimplexSolver.optimize`  
`UnivariateOptimizer.optimize`  
`UnivariateOptimizer.getMaxEvaluations`  
`UnivariateOptimizer.getEvaluations`

# Groups



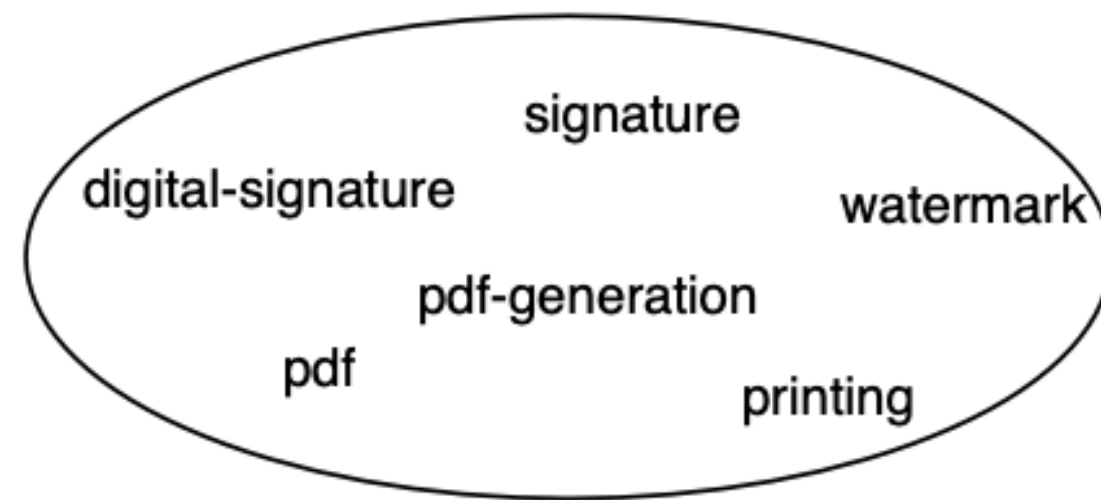
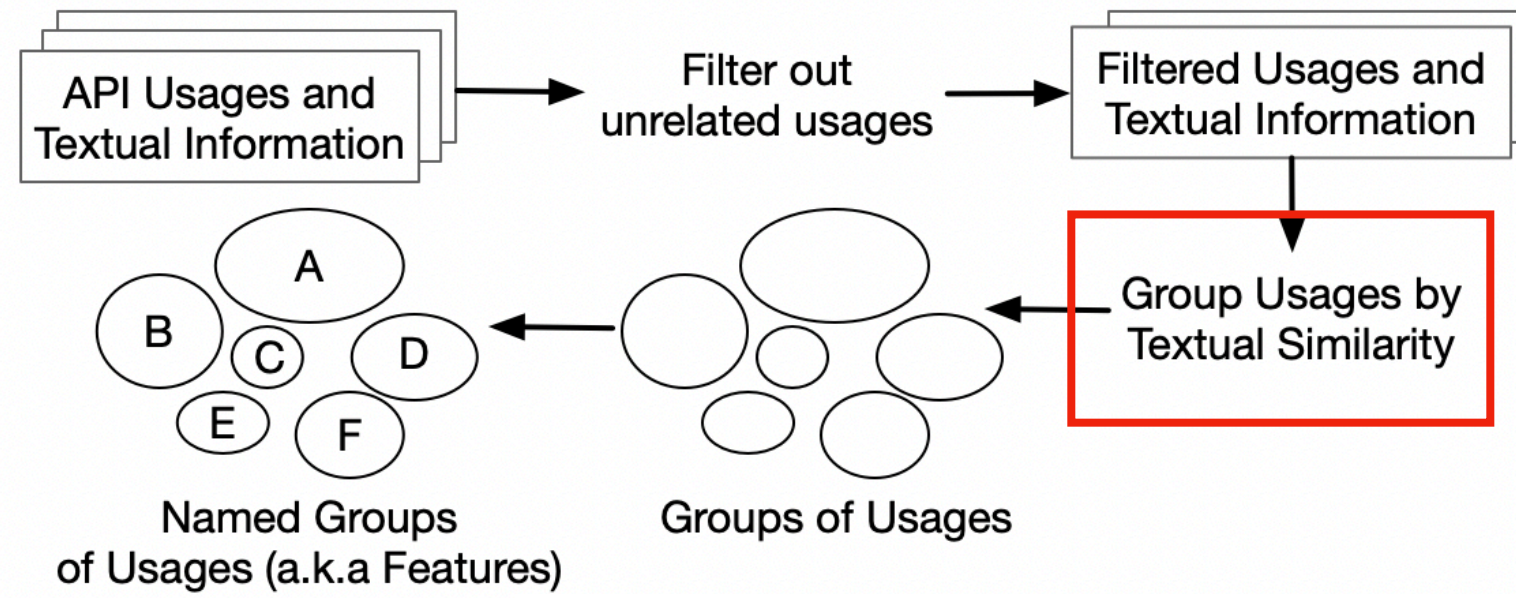
Library: com.itextpdf

PDF creation and transformation

# Groups

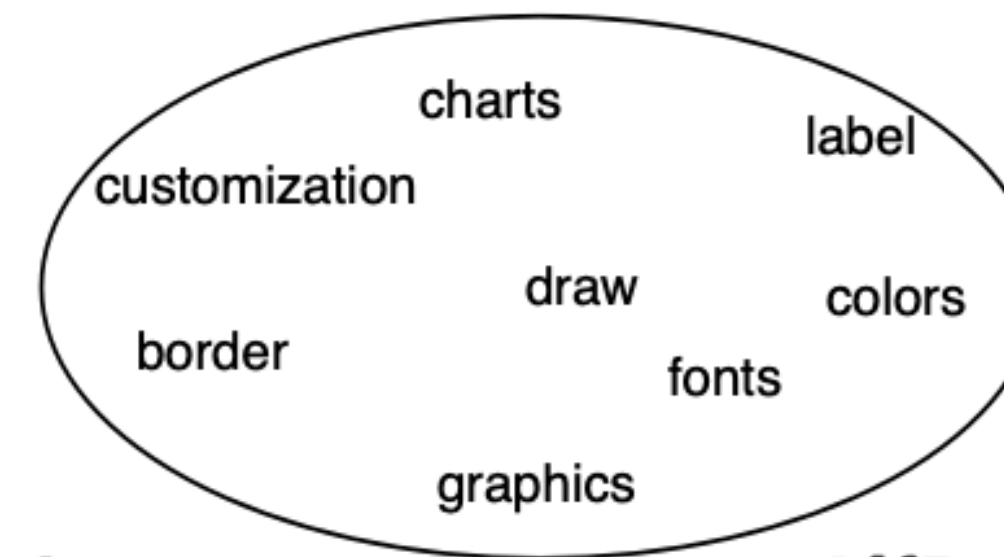
Library: com.itextpdf

PDF creation and transformation



**PdfStamper.setEncryption PdfStamper.close**  
**MakeSignature.signDeferred**

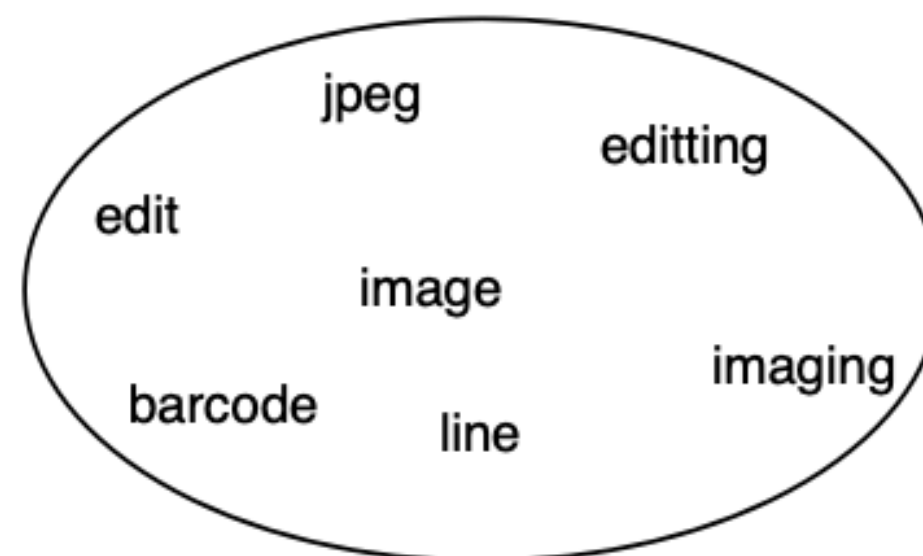
**PdfCopy.addDocument**  
**Document.open**  
**Document.isOpen**  
**Document.close**



**Document.close**  
**PdfContentByte.createGraphics**  
**Document.open**  
**DocumentException.getMessage**  
**Font.createFont**  
**PdfWriter.getInstance**

**Rectangle.setBackgroundColor**

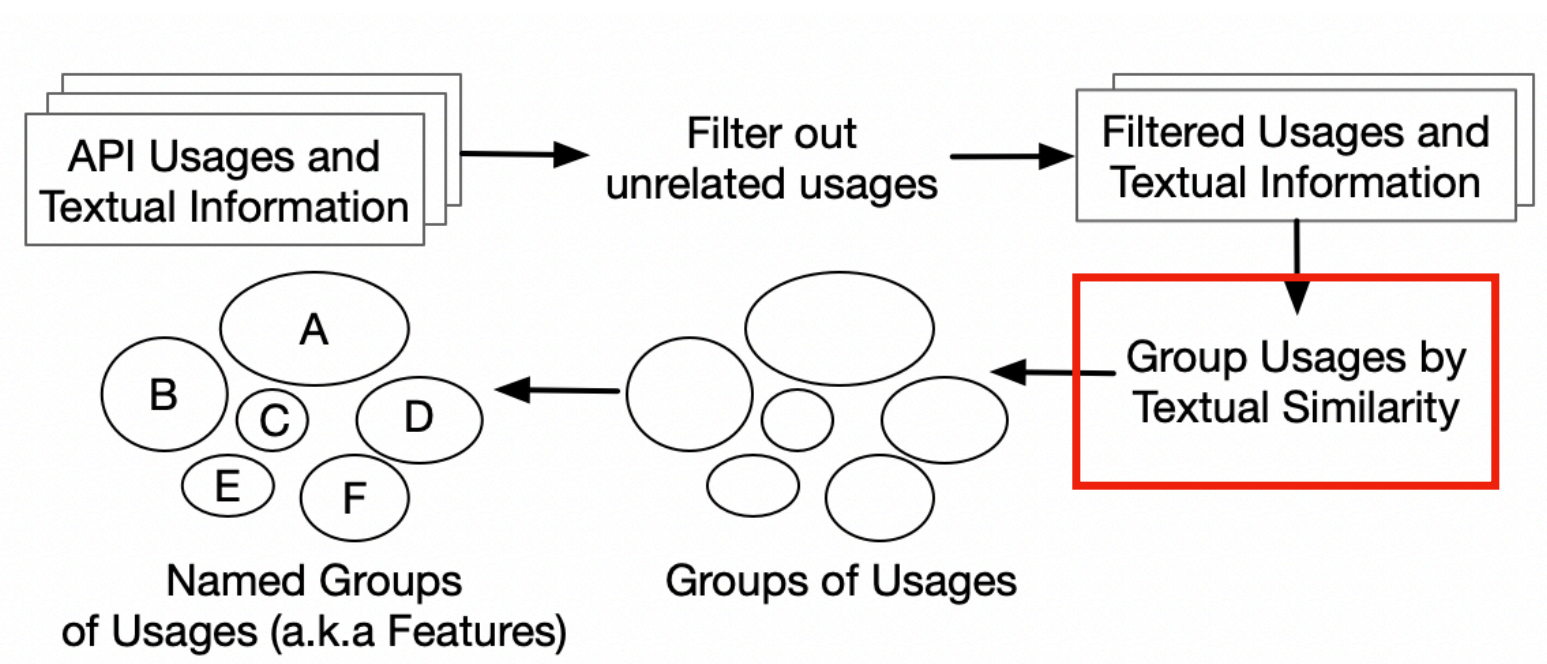
**PdfTemplate.fill**  
**PdfTemplate.setColorFill**  
**Image.getInstance**  
**PdfTemplate.rectangle**



**Image.getInstance**  
**PdfStamper.close**  
**PdfStamper.getOverContent**  
**ColumnText.showTextAligned**  
**PdfReader.close**

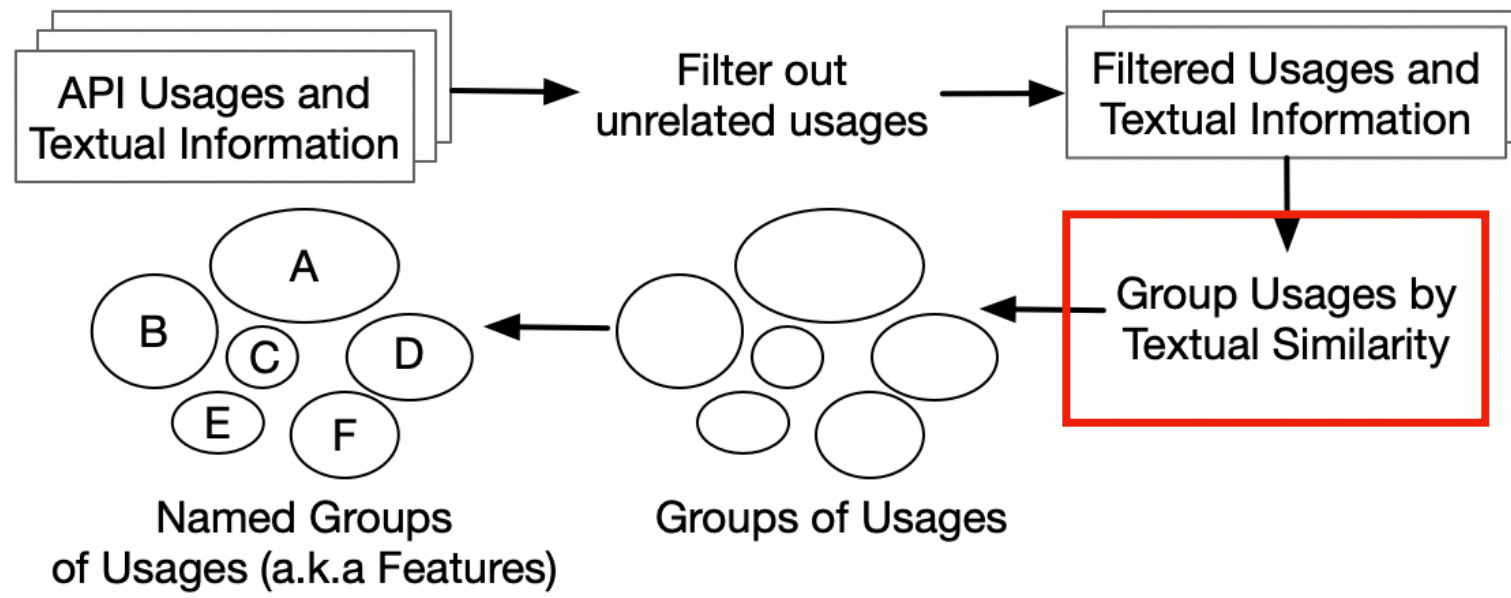
**Image.getWidth**  
**BarcodeQRCode.createAwtImage**

# Groups

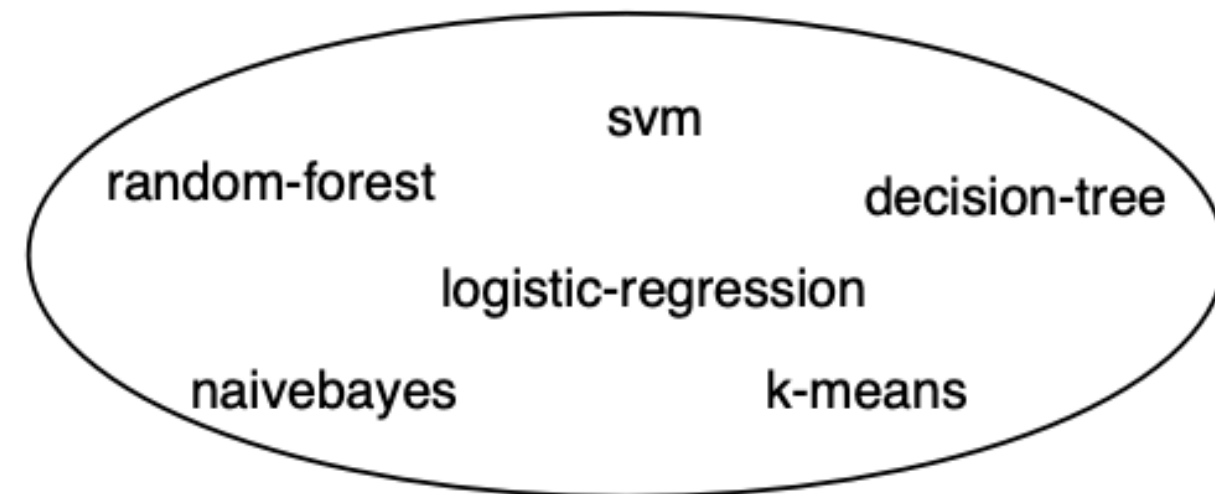


Library: [nz.ac.waikato.cms.weka](http://nz.ac.waikato.cms.weka)  
Data mining and Machine Learning Library

# Groups

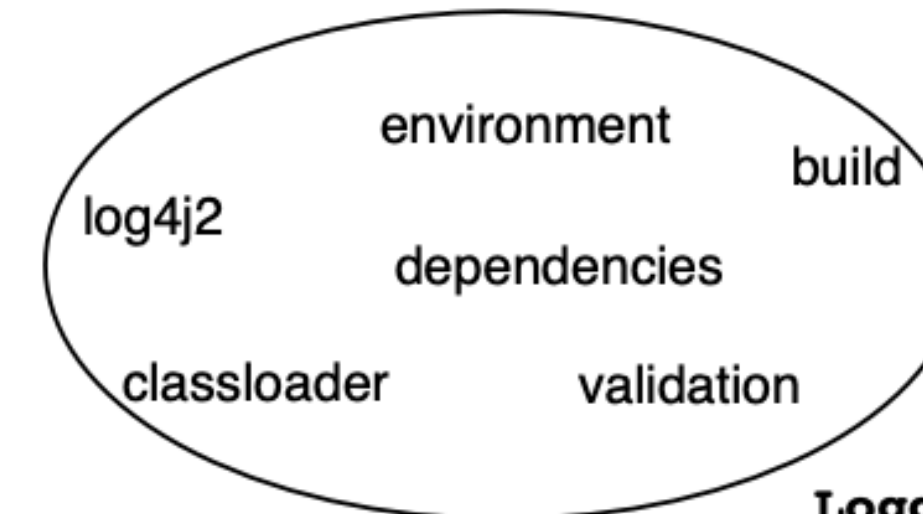


Library: `nz.ac.waikato.cms.weka`  
 Data mining and Machine Learning Library



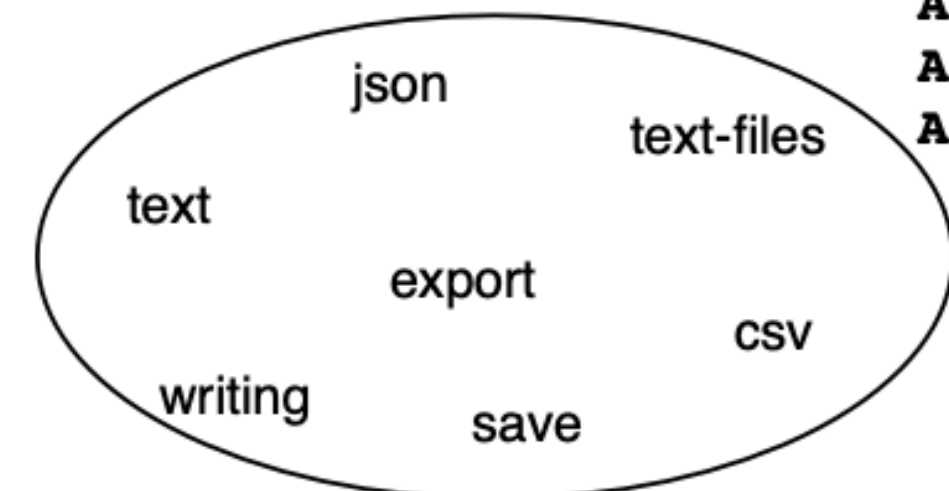
`J48.distributionForInstance`  
`DataSource.getDataSet`  
`J48.classifyInstance`

`SimpleKMeans.setNumClusters`  
`ClusterEvaluation.evaluateClusterer`  
`ClusterEvaluation.setClusterer`  
`SimpleKMeans.buildClusterer`  
`ClusterEvaluation.clusterResultsToString`



`Environment.getEnvironment`  
`Environment.getProperty`  
`Evaluation.evaluateModel`  
`Classifier.buildClassifier`

`Logger.getLogger`  
`Evaluation.evaluateModel`  
`Classifier.buildClassifier`

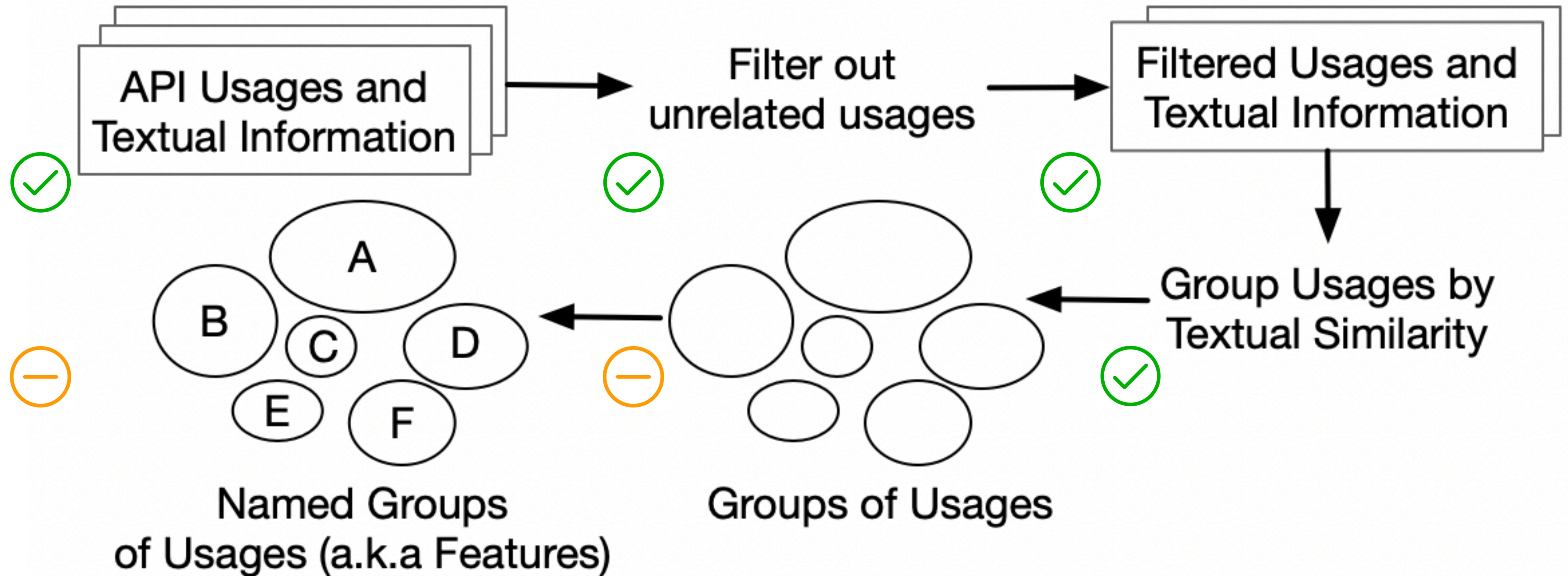


`ArffSaver.setDestination`  
`ArffSaver.writeBatch`  
`ArffSaver.setInstances`

`CSVLoader.getDataSet`  
`CSVLoader.setSource`  
`Instances.setClassIndex`  
`Instances.numAttributes`

`Evaluation.evaluateModel`  
`Evaluation.pctCorrect`  
`Evaluation.toSummaryString`

# Future Work



# Conclusion

## Introduction

Software products that evolve together in the same environment [1]



- I/O Utilities (17)
- HTTP Clients (81)
- Dependency Injection (51)
- Mocking (41)
- XML Processing (34)
- JSON Libraries (74)
- Collections (72)
- Reflection Libraries (47)
- Command Line Parsers (45)

**JCommander Library**

Version	Repository	Usages	Date
1.78	Central	20	Aug. 2019
1.77	Central	0	Aug. 2019
1.75	Central	3	Aug. 2019
1.72	Central	193	Jun. 2017
1.69	Central	18	Apr. 2017
1.64	Central	23	Mar. 2017
1.60	Central	30	Jan. 2017
1.58	Central	21	Sep. 2016
1.48	Central	209	Apr. 2015
1.47	Central	24	Dec. 2014
1.35	Central	123	Feb. 2014
1.32	Central	55	Sep. 2013
1.30	Central	54	Oct. 2012

[1] - M. Lungu, "Towards reverse engineering software ecosystems," in 2008 ICSM, IEEE, 2008, pp. 428-431.

## Introduction

### Selecting Third-Party Libraries: The Practitioners' Perspective

Enrique Larios Vargas  
Software Improvement Group  
The Netherlands  
elariosvargas@sig.eu

Mauricio Aniche  
Delft University of Technology  
The Netherlands  
M.FinavaroAniche@tudelft.nl

Christoph Treude  
University of Adelaide  
Australia  
christoph.treude@adelaide.edu.au

Magiel Bruntink  
Software Improvement Group  
The Netherlands  
m.bruntink@sig.eu

Georgios Gousios  
Delft University of Technology  
The Netherlands  
G.Gousios@tudelft.nl

**ABSTRACT**  
The selection of third-party libraries is an essential element of virtually any software development project. However, deciding which libraries to choose is a challenging practical problem. Selecting the wrong library can severely impact a software project in terms of cost, time, and development effort, with the severity of the impact depending on the role of the library in the software architecture, among others. Despite the importance of following a careful library selection process, in practice, the selection of third-party libraries is still conducted in an ad-hoc manner, where dozens of factors play an influential role in the decision.

In this paper, we study the factors that influence the selection process of libraries as perceived by industry developers. To that aim, we perform a cross-sectional interview study with 16 developers from 11 different businesses and survey 115 developers that are involved in the selection of libraries. We systematically devised a comprehensive set of 26 technical, human, and economic factors that developers take into consideration when selecting a software library. Eight of these factors are new to the literature. We explain each of these factors and how they play a role in the decision. Finally,

**1 INTRODUCTION**  
The use of third-party libraries<sup>1</sup> is a recurrent practice among practitioners to speed up the development of software systems and, as a consequence, to reduce its costs [10]. Software libraries provide developers with customized functionality which is useful while developing software [11]. Developers do not need to "reinvent the wheel", but rather seek libraries that suit their purpose [12]. The selection of the right library to adopt is a demanding task for developers. The number of factors that should be taken into account during the selection process is significant. For example, the adoption of a library that goes against the existing architectural decisions of the software system might increase its adoption time; a library that is easy to use but not performant enough might have an impact on the quality of the system; a library that has good documentation and active community, but with which the team has had a bad experience, might not be the best possible choice. On top

Survey of 16 participants from different companies

1. Technical Factors
    - Functionality
    - Quality
  2. Human Factors
    - Stakeholders
    - Organization
  3. Economic Factors
    - Total cost of ownership
    - Risk
- Type of Project
  - Release Process
  - Individual
  - Community

## Introduction

Stack Overflow Q&A

How do I convert a String to an int in Java? **Guava equivalent for IOUtils.toString(InputStream)**

How do I compare strings in Java? **Remove last character of a StringBuilder?**

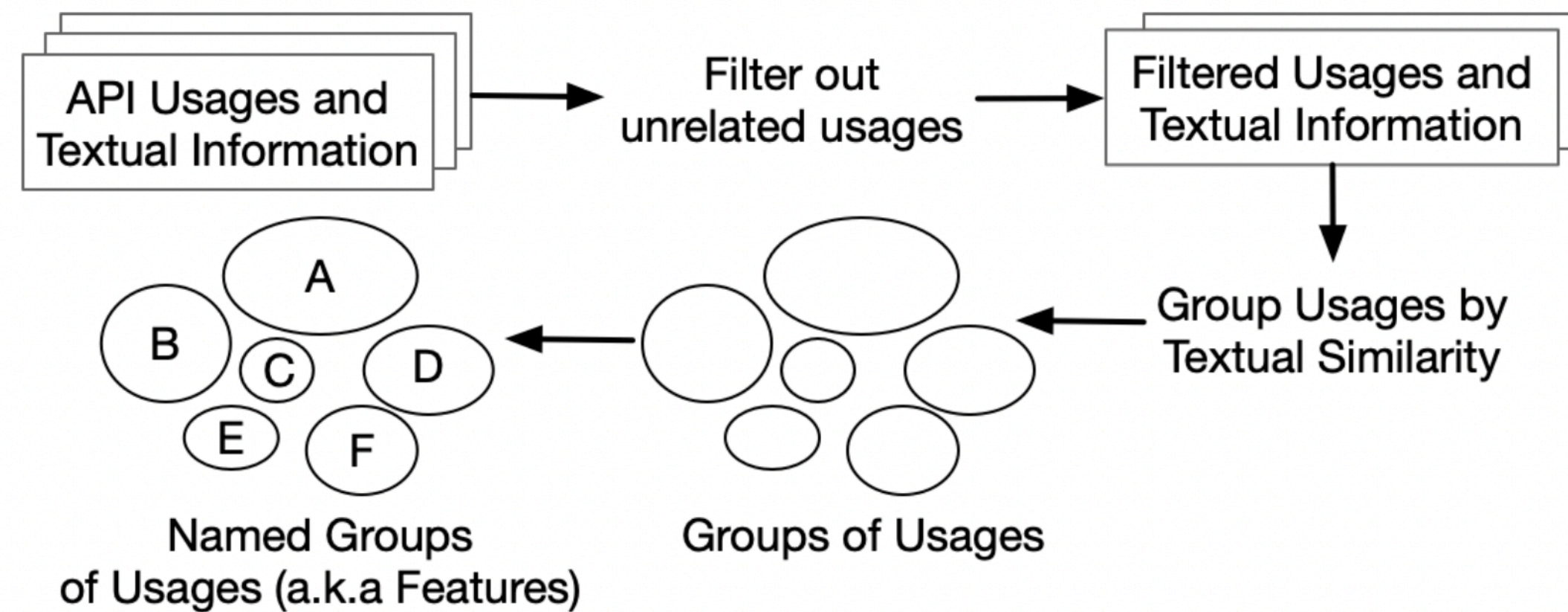
**Guava: Splitter and considering Escaping?** **String utilities**

**initializing a Guava ImmutableMap**  
**Flattening an Iterable<Iterable<T>> in Guava**  
**How to directly initialize a HashMap (in a literal way)?**  
**Google Guava isNullOrEmpty for collections**  
**Google Guava "zip" two lists** **Collection utilities**

**Merge ranges with guava** **Ranges**

**What is the simplest way to read a file into String?**  
**Utils to read resource text file to String (Java)** **I/O**

## Pipeline



## Groups

Library: org.apache.commons.math  
Mathematics Library from the Apache Foundation

**ComplexFormat.parse**  
**Complex.getReal**  
**Complex.getImaginary**

**BigFraction.getNumerator**  
**BigFraction.getDenominator**

**Precision.round**

**PolynomialSplineFunction.getPolynomials**  
**PolynomialSplineFunction.value**  
**LinearInterpolator.interpolate**

**PointValuePair.getPoint**  
**SimplexSolver.optimize**

**SummaryStatistics.getMean**  
**SummaryStatistics.addValue**  
**SummaryStatistics.getStandardDeviation**

**LogNormalDistribution.sample**

**NormalDistribution.cumulativeProbability**  
**NormalDistribution.inverseCumulativeProbability**

**UnivariateOptimizer.optimize**  
**UnivariateOptimizer.getMaxEvaluations**  
**UnivariateOptimizer.getEvaluations**