



# BIG DATA

## HÄUFIGKEITSBERECHNUNG

HTW Berlin | Programmierkonzepte und  
Algorithmen | 30 VI 2020 | Loparev & Nguyen



# INHALT

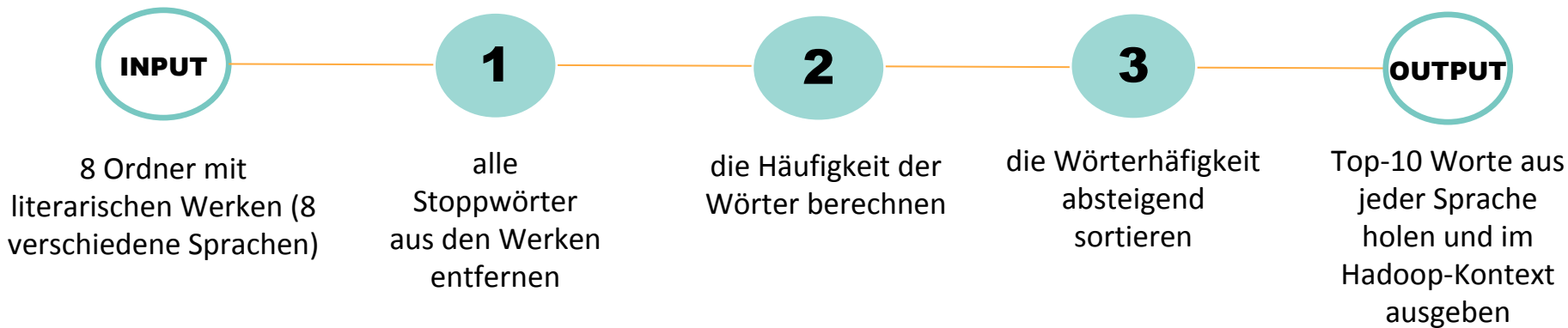
- 1 **Erklärung der Aufgabe**
- 2 **Planung und Lösungsbeschreibung**
- 3 **Ergebnis**
- 4 **Diagrammen für die Leistung und Laufzeit**
- 5 **Fazit**

# Aufgabe 09

## Häufigkeitsberechnung

- Berechnen Sie die Vorkommenshäufigkeit für alle Wörter in gegebenen Texten
- Stellen Sie die Ergebnisse als TOP-10 Liste getrennt für jede Sprache vor, mit Ausnahme der Stoppwörter (z.B. “a”, “an”, “the”, “of”, “und”, “mit”, “la”, “и”, “же” usw.)

# Planung



# Eingabe

# Lösungsbeschreibung

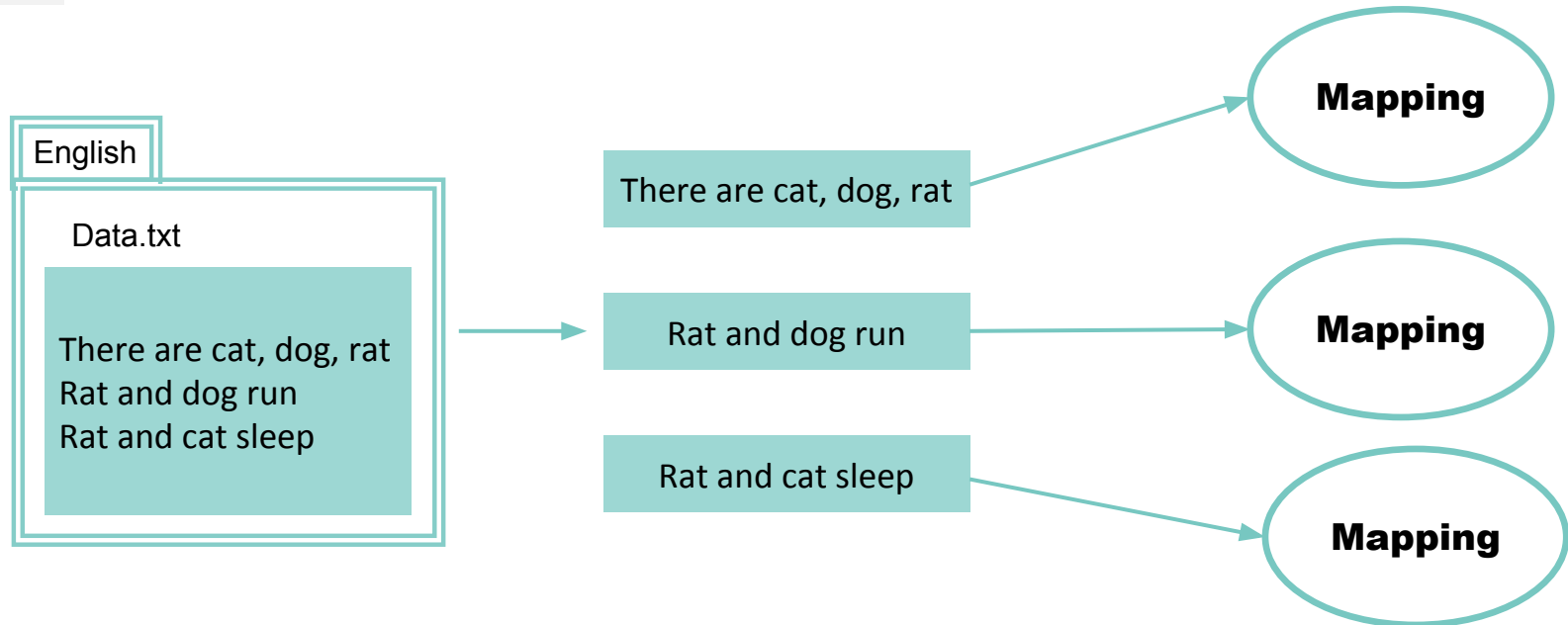
## Input

- English  
... Data.txt
- German
- Dutch
- Italian
- French
- Russian
- Ukrainian
- Spanish

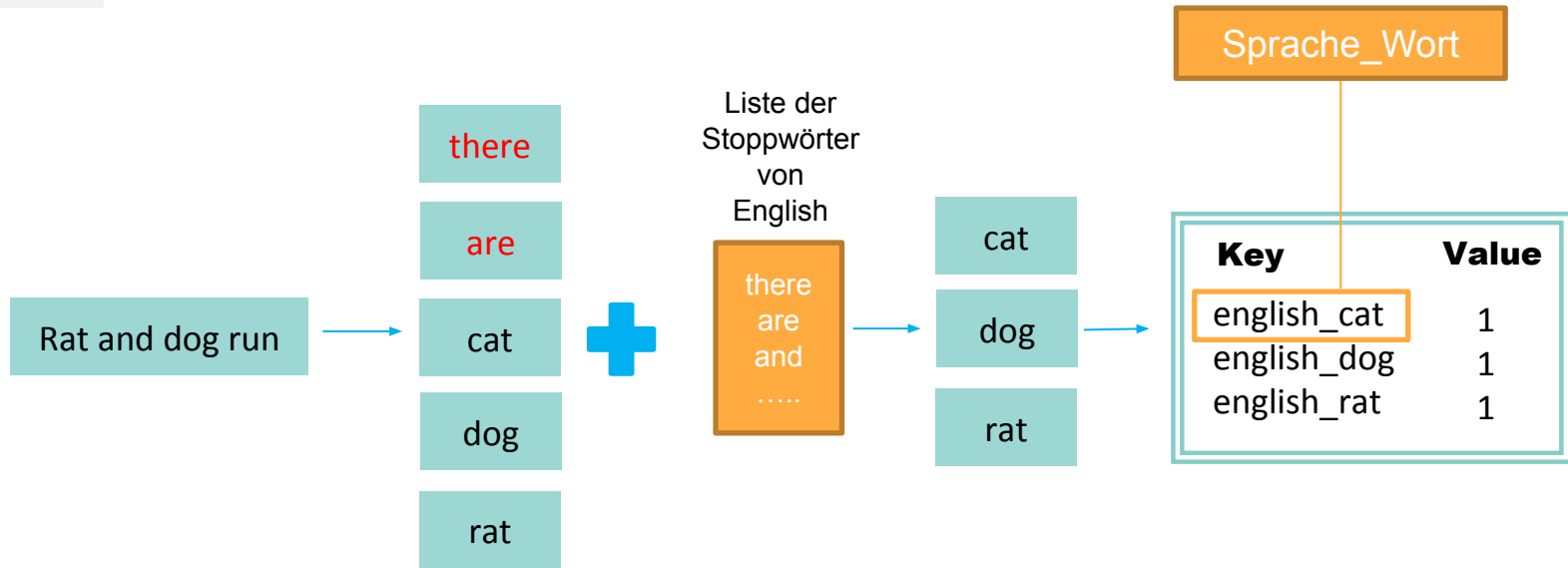
## Stopwords

English	German	Italian	Russian	Dutch	Ukrainian	French	Spanish
there are and ...	...	...	...	...	...	...	...

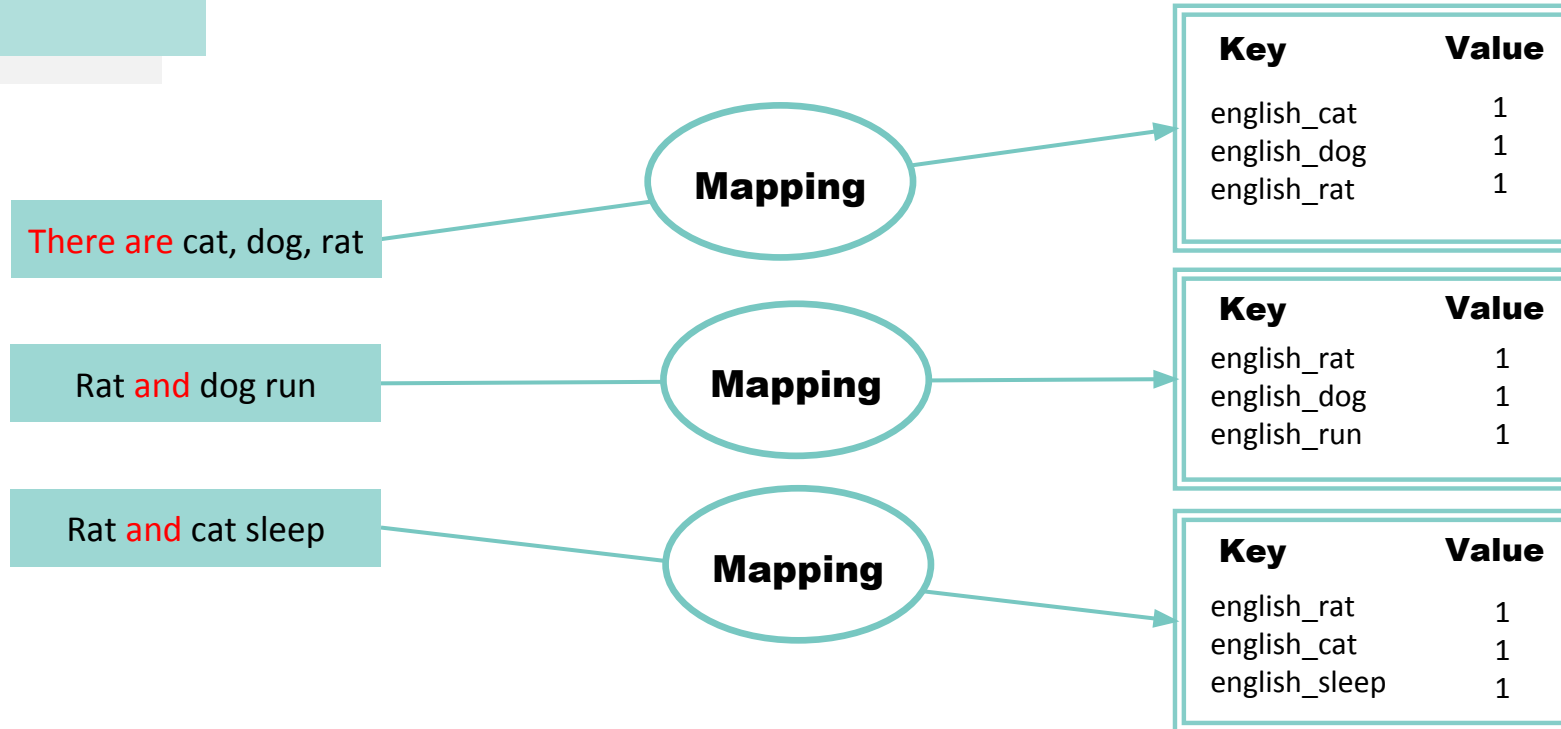
# Vorbereitung auf Eingabe für Mapping



# Mapping

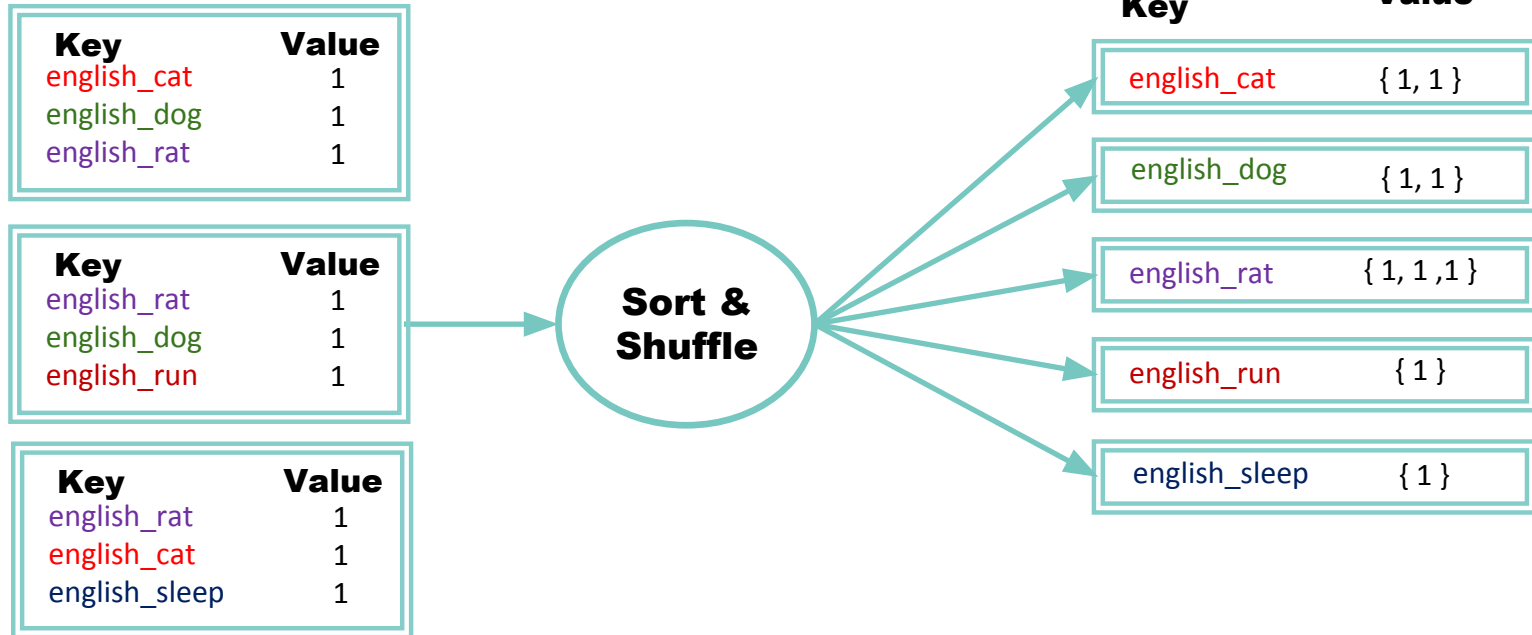


# Mapping

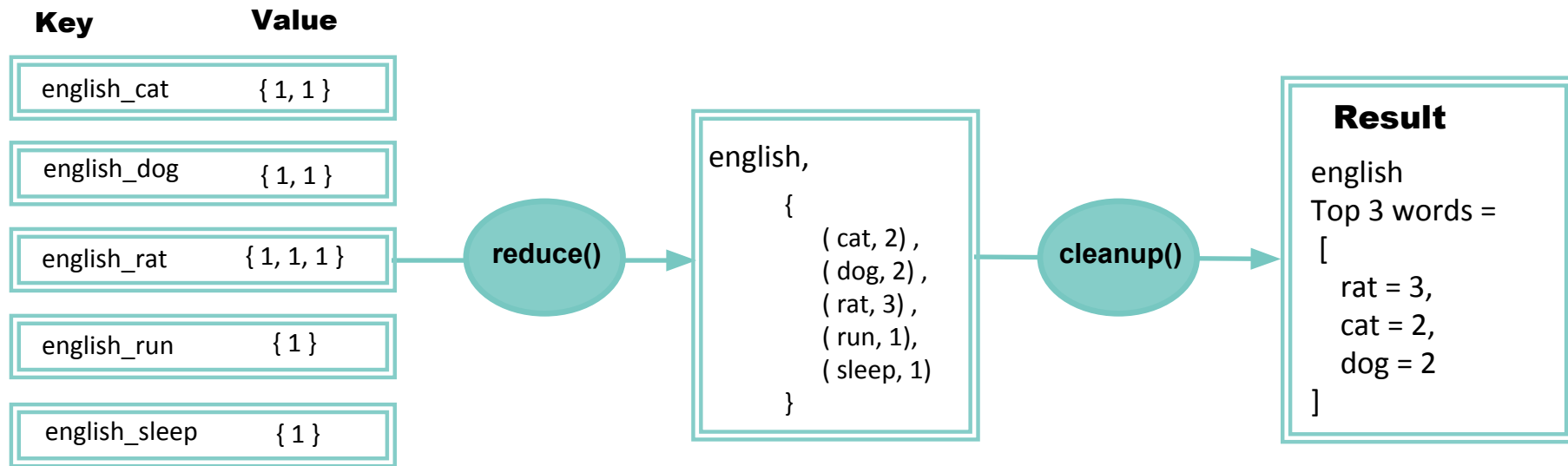




# Sort and Shuffle



# Reducer



# Mapping

```

01. public class MyMapper extends Mapper<LongWritable, Text, Text, Text> {
02.
03.     private final static Text one = new Text("1");
04.     private Text word = new Text();
05.     private HashMap<String, HashSet<String>> stopwordList = new HashMap<>();
06.
07.     public void map(LongWritable key, Text value, Context context) throws IOException, InterruptedException {
        String line = value.toString().toLowerCase();
        List<String> words = new ArrayList<>();
        Collections.addAll(words, line.split("\\p{L}+"));
        String filePathString = ((FileSplit) context.getInputSplit()).getPath().toString();
        String language = new Path(filePathString).getParent().getName().toLowerCase();
        if (!stopwordList.containsKey(language)) {
            HashSet<String> stopwords = Utils.readFileStopword(language + ".txt");
            if (stopwords == null) {
                return;
            }
            stopwordList.put(language, stopwords);
        }
        for (String w : words) {
            if (w.length() <= 2) {
                continue;
            }
            if (stopwordList.get(language).contains(w)) {
                continue;
            }
            word.set(language + "_" + w);
            context.write(word, one);
        }
    }
}

```

Eine Liste "stopwordList" als HashMap erstellen, in der die Stopwörter jeder Sprache gespeichert werden.

Die Kette in String und in Kleinbuchstaben konvertieren.

die Kette in Wörter teilen

Checken, ob die Sprache bereits in der Liste der Stopwörter ist oder nicht?

Den Name des Ordners holen, der die Textdatei der Kette enthält und auch der Name der Sprache ist

Laut der Sprache die Daten der Stopwörter in einer Liste laden

die Liste der Stopwörter in der gemeinsamen Liste der Stopwörter hinzufügen.

Wenn das Wort seine Länge <= 2 hat, wird es ignoriert

Die Daten in der Hadoop-Kontext als Paaren <"Sprache\_Wort", 1> schreiben

Checken, ob das Wort Stopwort ist oder nicht? Wenn Ja, wird es ignoriert

Eine Kette als Form "<Sprache>\_<Wort>" erstellen, um den Schlüssel zu erstellen

# Reducer

## reduce()

```

01. public class MyReducer extends Reducer<Text, Text, Text, Text> {
02.
03.     private HashMap<String, HashMap<String, Integer>> allData = new HashMap<>();
04.
05.     private final int MAX_TOP = 10;
06.
07.     @Override
08.     public void reduce(Text key, Iterable<Text> values, Context context) throws IOException, InterruptedException {
09.
10.         String[] language_key = key.toString().split("_");
11.         String language = language_key[0];
12.         String str = language_key[1];
13.
14.         int sum = 0;
15.
16.         for (Text val : values) {
17.             sum += Integer.parseInt(val.toString());
18.         }
19.
20.         if (allData.containsKey(language)) {
21.             allData.get(language).put(str, sum);
22.         } else {
23.             HashMap<String, Integer> hmap = new HashMap<>();
24.
25.             hmap.put(str, sum);
26.             allData.put(language, hmap);
27.         }
28.     }
29.
30. }

```

eine Liste "allData" als HashMap erstellen, in der die Wörter jeder Sprache gespeichert werden

Die Anzahl von TOP der Wörter, die ausgegeben werden

den "Key" in 2 Teilen schneiden. Der erste ist die Sprache und der zweite ist das Wort

die Häufigkeit des Wortes berechnen

Checken, ob die Sprache des Wortes bereits in der Liste "allData" ist oder nicht?

das Wort und seine Häufigkeit in der Liste hinzufügen.

## cleanup()

```

47.     @Override
48.     protected void cleanup(Context context) throws IOException, InterruptedException {
49.
50.         Comparator<Entry<String, Integer>> comparator = new Comparator<Entry<String, Integer>>() {
51.
52.             @Override
53.             public int compare(Entry<String, Integer> o1, Entry<String, Integer> o2) {
54.                 return o2.getValue().compareTo(o1.getValue());
55.             }
56.         };
57.
58.         for (Entry<String, HashMap<String, Integer>> entry : allData.entrySet()) {
59.
60.             List<Entry<String, Integer>> values = new ArrayList<>(entry.getValue().entrySet());
61.             Collections.sort(values, comparator);
62.
63.             StringBuilder sb = new StringBuilder("-----\n");
64.
65.             sb.append(entry.getKey());
66.             sb.append("\nTop " + MAX_TOP + " words=");
67.             sb.append(values.subList(0, values.size() > MAX_TOP ? MAX_TOP : values.size()));
68.
69.             context.write(new Text(sb.toString()), new Text("\n"));
70.         }
71.     }
72.
73. }

```

Erstellen einen Komparator, um die Häufigkeitszahl zwischen den Wörter zu vergleichen.

Sortieren die Wörter in absteigender Reihenfolge

Erstellen eine StringBuilder, um zu ausgehenden Informationen zu enthalten

Lesen jede Sprache in der Liste allData,

Fügen den Name der Sprache hinzu

Holen sich Top 10 Wörter

Schreiben die Informationen der Sprache in Hadoop-Kontext

# Ergebnis

-----  
dutch

Top 10 words=[den=1520, zoo=512, marten=330, baas=303, zien=230, jan=224, vrouw=201, oogen=194, riep=185, goed=178]

-----  
ukrainian

Top 10 words=[треба=1029, сказав=810, над=805, мати=783, перед=758, добре=749, серце=738, уже=738, чіпка=738, усе=731]

-----  
german

Top 10 words=[hand=5610, augen=5172, herr=4815, leben=3924, sprach=3701, stand=3630, vater=3553, ließ=3375, mußte=3243, alte=3102]

-----  
spanish

Top 10 words=[don=3073, casa=2451, ojos=1816, vida=1599, hombre=1565, dios=1508, señor=1497, padre=1395, mujer=1314, mano=1262]

-----  
russian

Top 10 words=[глаза=298, сердце=271, жизни=234, знаю=224, горский=224, руки=209, слова=200, вера=192, точно=191, руку=183]

-----  
english

Top 10 words=[man=11763, time=9585, good=6752, holmes=5899, long=5765, great=5534, day=5295, men=5018, thought=4828, night=4820]

-----  
italian

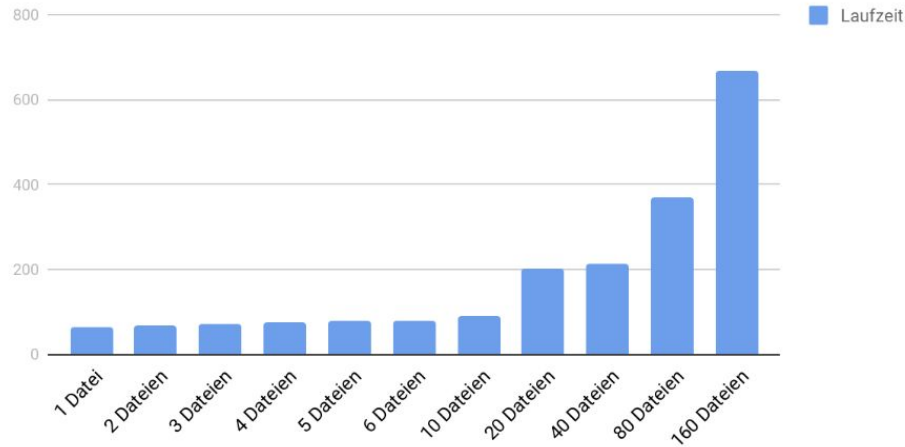
Top 10 words=[disse=8759, rispose=4023, occhi=3614, donna=3087, mano=3055, altra=2968, voce=2862, uomini=2851, signor=2710, più=2626]

-----  
french

Top 10 words=[homme=10266, faire=9207, monsieur=7084, point=6670, jeune=5891, femme=5866, fit=5731, yeux=5615, jamais=5572, oui=5526]

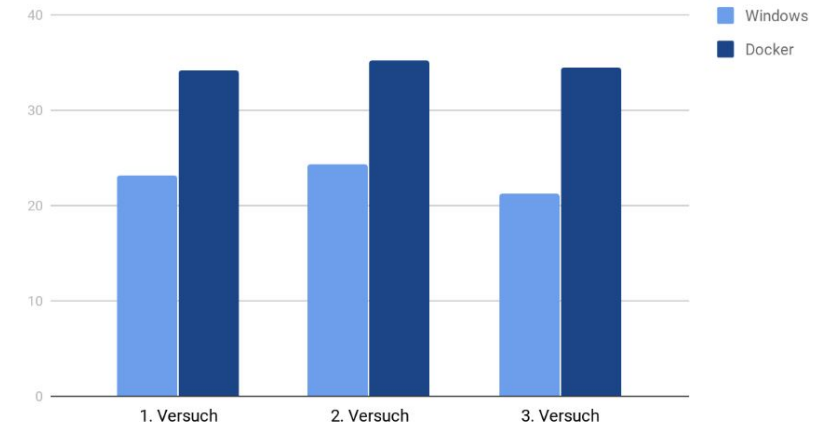
# Diagrammen

(Einheit : Sekunde)



4.1 Laufzeitvergleich in Abhängigkeit von der Inputgröße

(Einheit: Minute)



4.2 Laufzeitvergleich zwischen Hadoop und Docker

# Fazit

## Häufigkeitsberechnung

- Hadoop braucht relativ viel Zeit um einzelne Dateien zu öffnen
- die Verarbeitung an sich dauert verhältnismäßig schnell
- die Laufzeit kann durch Dateienvereinigung verbessert werden

# Vielen Dank!

für Ihre Merksamkeit