# Word Count и формальная модель MapReduce

**Драль Алексей**, study@bigdatateam.org
CEO at BigData Team, https://bigdatateam.org
https://www.facebook.com/bigdatateam

Apache Hadoop (/hə`du:p/) is an open-source software framework used for distributed storage and processing of dataset of big data using the MapReduce programming model. It consists of computer clusters built from commodity hardware.

All the modules in Hadoop are designed with a fundamental assumption that hardware failures are common occurrences and should be automatically handled by the framework...

Apache Hadoop (/həˈduːp/) is an open-source software framework used for distributed storage and processing of dataset of big data using the MapReduce programming model. It consists of computer clusters built from commodity hardware.

All the modules in Hadoop are designed with a fundamental assumption that hardware failures are common occurrences and should be automatically handled by the framework...

WIKIPEDIA
The Free Encyclopedia

```
'the': 3, 'of': 3, 'hadoop': 2, …
```

```
$ cat dataset.txt
```

```
    Apache Hadoop is a collection of open-source
software utilities that facilitates using a
network of many computers to solve problems
involving massive amounts of data and computation.
It provides a software framework for distributed
storage and processing of big data using the
MapReduce programming model...
```

```
$ cat dataset.txt | tr ' ' '\n'

    Apache
    Hadoop
    is
    a
    collection
    of
    ...
```

```
$ cat dataset.txt | tr ' ' '\n' | sort
```

```
    All
    Apache
    Hadoop
    Hadoop
    Hadoop
    It
    ...
```

```
$ cat dataset.txt | tr ' ' '\n' | sort | uniq -c
    1 All
    1 Apache
    3 Hadoop
    2 It
    1 MapReduce
    4 a
    ...
```

```
$ cat dataset.txt | tr ' ' '\n' | sort | uniq -c

    1 All
    1 Apache
    3 Hadoop
    2 It
    1 MapReduce
    4 a
    ...
```

```
$ cat dataset.txt | tr ' ' '\n' | sort | uniq -c
    1 All
    1 Apache
    3 Hadoop
    2 It
    1 MapReduce
    4 a
    ...
```

**BIGDATA TEAM**

```
$ cat dataset.txt | tr ' ' '\n' | sort | uniq -c
```

```
    1 All
    1 Apache
    3 Hadoop
    2 It
    1 MapReduce
    4 a
    ...
```

```
$ cat dataset.txt | tr ' ' '\n' | sort | uniq -c
```

❌ Фаза Map (нужна агрегация)

```
$ cat dataset.txt | tr ' ' '\n' | sort | uniq -c
```

❌ Фаза Map (нужна агрегация)
❌ Фаза Reduce (не хватит RAM / HDD)

# MapReduce =

# Map → Shuffle & Sort → Reduce

wikipedia.dump | tr ' ' ' \n' |     sort     |     uniq -c



**Wikipedia.dump**

**Block 1**
Apache Hadoop (/həˈduːp/) is an open-source software framework used for distributed storage

**Block 2**
and processing of dataset of big data using the MapReduce programming model. It consists of computer cluster built from

**Block M**
commodity hardware. All the modules in Hadoop are designed with a fundamental assumption that...

hash(word) % R

hash(word) % R

hash(word) % R

uniq -c

**a...**    **1**

uniq -c

**b...**    **2**

...

uniq -c

**z...**    **26**

$\Sigma$

wikipedia.dump -> map () -> word     shuffle & sort     reduce()

Фазы:

Фазы:

1. Map

Фазы:

1. Map
2. Shuffle & Sort

Фазы:

1. Map
2. Shuffle & Sort
3. Reduce

Фазы:

1. Map
2. Shuffle & Sort
3. Reduce

Worker'ы (контейнеры):

Фазы:

1. Map
2. Shuffle & Sort
3. Reduce

Worker'ы (контейнеры):

► Фаза Map → Mapper (использует функцию map)

Фазы:

1. Map
2. Shuffle & Sort
3. Reduce

Worker'ы (контейнеры):

► Фаза Map → Mapper (использует функцию map)
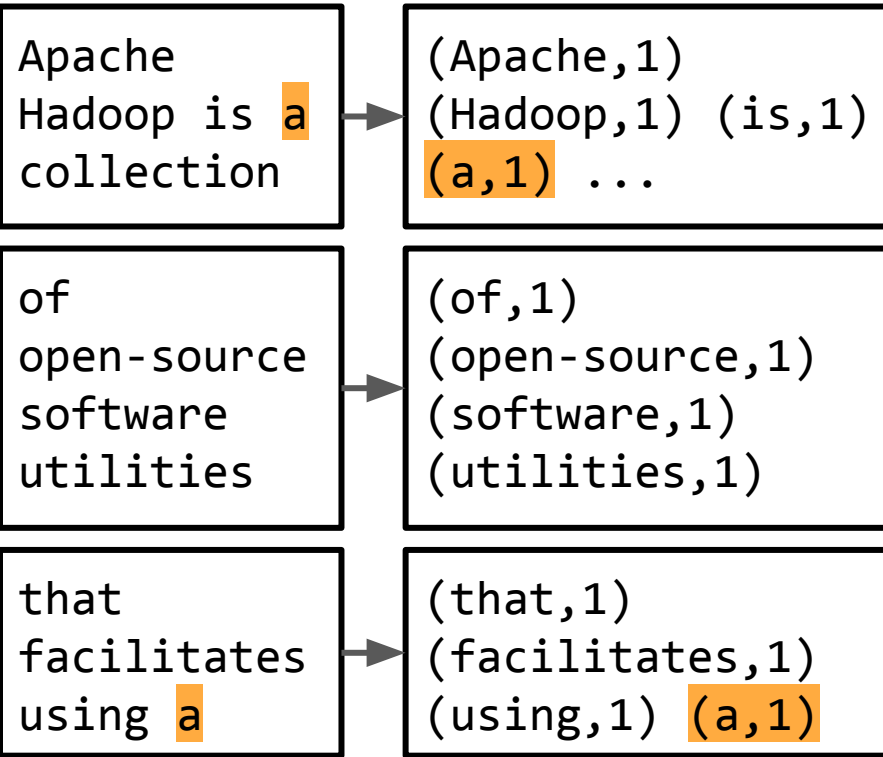► Фаза Reduce → Reducer (использует функцию reduce)

**BIGDATA TEAM**

Apache
Hadoop is a
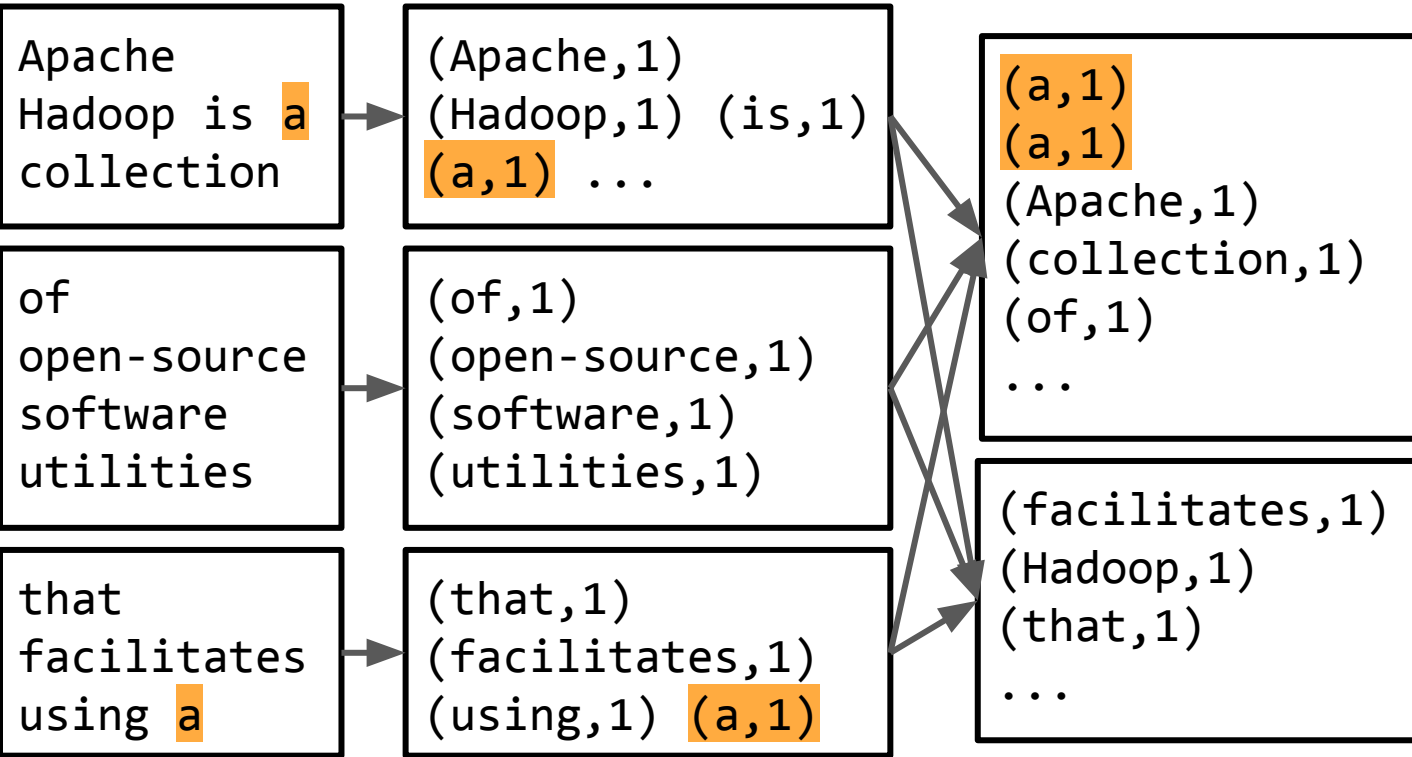collection

of
open-source
software
utilities

that
facilitates
using a

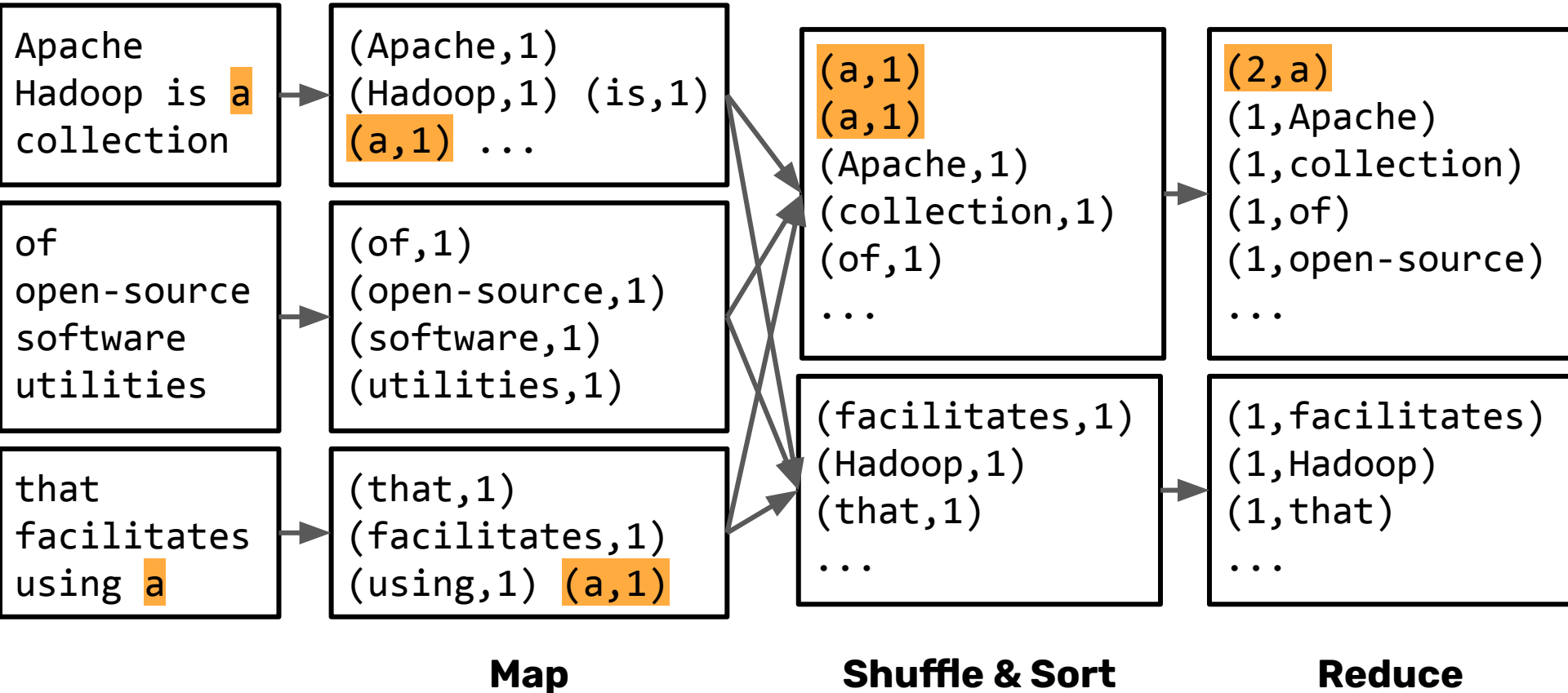Apache Hadoop is a collection

(Apache,1)
(Hadoop,1) (is,1)
(a,1) ...

of open-source software utilities

(of,1)
(open-source,1)
(software,1)
(utilities,1)

that facilitates using a

(that,1)
(facilitates,1)
(using,1) (a,1)

**Map**

**Map**

**Shuffle & Sort**

| Map | Shuffle & Sort | Reduce |

(k_in, v_in)

**map (функция)**

???

**Map (фаза)**

```
(k_in, v_in)
```

**map (функция)**

```
map = (tr ' ' '\n')
(-, line) → [(word, 1), ...]
```

```
[(k_interm, v_interm), ...]
```

**Map (фаза)**

**BIGDATA TEAM**

```
(k_in, v_in)
```

map (функция)

```
map = (tr ' ' '\n')
(-, line) → [(word, 1), ...]
```

```
[(k_interm, v_interm), ...]
```

**Map (фаза)**

sort and group by k_interm

**Shuffle & Sort**

```
(k_interm, [v_interm, ...])
```

reduce (функция)

```
???
```

**Reduce (фаза)**

# Формальная модель Word Count



(k_in, v_in)

**map (функция)**

```
map = (tr ' ' '\n')
(-, line) → [(word, 1), ...]
```

[(k_interm, v_interm), ...]

**Map (фаза)**

sort and group by k_interm

**Shuffle & Sort**

(k_interm, [v_interm, ...])

**reduce (функция)**

```
reduce = (uniq -c)
(word, [1,1,...]) → (7, word)
```

[(k_out, v_out), ...]

**Reduce (фаза)**

Теперь вы:

Теперь вы:

► Можете насчитать больше 2х фаз в MapReduce

Теперь вы:

► Можете насчитать больше 2х фаз в MapReduce

► Понимаете как решать Word Count при помощи MapReduce

Теперь вы:

► Можете насчитать больше 2х фаз в MapReduce

► Понимаете как решать Word Count при помощи MapReduce

► Знаете какие 3 типа пар (key, value) указываются при запуске (например) Java MapReduce приложения