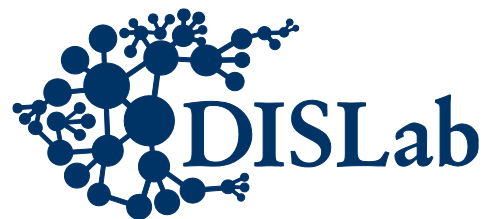


Параллельная обработка больших графов

Занятие 4

А.С. Семенов

dislab.org



Поиск всех кратчайших путей от заданной вершины в графе (Single Source Shortest Paths, SSSP)

$$G = (V, E, W), W : E \rightarrow R, \text{directed}$$

$$p(v_0, v_k) = \langle v_0, v_1, \dots, v_k \rangle$$

$$w(p) = \sum_{i=1}^k w(v_{i-1}, v_i)$$

$$\delta(u, v) = \begin{cases} \min\{w(p)\}, \exists p(u, v) \\ \infty, \text{otherwise} \end{cases}$$

Алгоритмы решения Single Source Shortest Paths

- Дейкстры, 1959
- Беллмана-Форда, 1969
- Дельта-степпинг, 2003

SSSP, Вспомогательные процедуры

InitializeSingleSource(G, s)

for all $v \in V$

$d[v] = \infty$

end for

$d[s] = 0$

Relax(u, v)

if $d[v] > d[u] + w(u, v)$

$d[v] = d[u] + w(u, v)$

$p[v] = u$

end if

SSSP, Алгоритм Дейкстры

Dijkstra (G, w, s)

InitializeSingleSource(G, s)

$S = \{ \}$

$Q = V$

while $Q \neq \{ \}$

$u = \text{GetMin}(Q)$

$S = S \cup \{u\}$

$Q = Q - \{u\}$

for all v **in** $\text{Adj}[u]$

 Relax (u, v)

end for

end while

- $G(V, E, W)$ – ориентированный, $w(e) \geq 0, w \in W$
- последовательный

Сложность $O(N^2 + M)$, $O(N \lg N + M)$

SSSP, Алгоритм Беллмана-Форда

```
BellmanFord (G, s)  
InitializeSingleSource(G, s)  
for i = 1 to |V| - 1  
    for all (u, v) in E  
        Relax(u, v)  
    end for  
end for
```

Сложность $O(N M)$

- $G(V, E, W)$ – ориентированный, $W: E \rightarrow R$
- может определять наличие циклов с отрицательным весом

SSSP, Алгоритм Беллмана-Форда

```
BellmanFord (G, s)
InitializeSingleSource(G, s)
for i = 1 to |V| - 1
    #pragma omp parallel for
    for all (u, v) in E
        Relax(u, v) // атомарно
    end for
end for
```

SSSP, Алгоритм delta-stepping (2003, U. Meyer, P. Sanders) (1/4)

- $G(V, E, W)$ – ориентированный, $w(e) \geq 0$, $w \in W$
- $\delta \in \mathbb{R}$, $\delta > 0$

DeltaStepping (G, s, δ)

InitializeDeltaStepping

$i = 0$

while true

 processBucket(i)

$i = \min(k > i : B_k \neq \{ \})$

if $i == \text{inf}$ **then break**

end while

SSSP, Алгоритм delta-stepping (2/4)

InitializeDeltaStepping (G, s)

for all u in V

 d(u) = inf

end for

d(s) = 0

B₀ = {s}

B_{inf} = V / {s}

B₁ = B₂ = ... = {}

$$B_i = \{v \in V : d[v] \in [i * \text{delta}; (i + 1) * \text{delta}]\}$$

SSSP, Алгоритм delta-stepping (3/4)

ProcessBucket (i)

$A = B_i$

while $A \neq \{\}$

for all u **in** A , **for all** $e = (u, v)$ **in** E

 Relax(u, v)

end for

$A' = \{x : d(x) \text{ изменен в предыдущем цикле}\}$

$A = B_i \cap A'$

end while

SSSP, Алгоритм delta-stepping (4/4)

Relax (u, v)

$i = \lfloor d(v) / \text{delta} \rfloor$ // старый bucket

$d(v) = \min(d(v), d(u) + w(u, v))$

$j = \lfloor d(v) / \text{delta} \rfloor$ // новый bucket

if $j < i$ **then** переместить v из B_i в B_j

SSSP, Алгоритм delta-stepping

DeltaStepping (G, s, delta)

InitializeDeltaStepping

$i = 0$

while true

 processBucket(i)

$i = \min(k > i : B_k \neq \{ \})$

if $i \neq \text{inf}$ **then break**

end while

- Для случайных графов «хороший» показатель $\text{delta} = 1/d$, где d – максимальная степень вершины, w равномерно распределены в $[0; 1]$
 - Сложность $O(M \log N)$
 - Количество итераций $O\left(\frac{d_c}{\text{delta}} * \frac{\log n}{\log \log n}\right)$
 - d_c - максимальная длина пути

SSSP, Алгоритм delta-stepping

ProcessBucket (i)

$A = B_i$

while $A \neq \{\}$

#pragma omp parallel for

for all u **in** A , **for all** $e = (u, v)$ **in** E

 Relax(u, v) // **атомарно**

end for

$A' = \{x : d(x) \text{ changed in previous for all}\}$

$A = B_i \cap A'$

end while

Варианты заданий

Реализовать алгоритм дельта-степпинг на распределенной памяти при помощи:

- Функции MPI-1
 - mpi4py
 - Односторонние коммуникации MPI-3
 - Библиотека OpenSHMEM
 - Charm++, charm4py
 - Chapel
 - Apache Spark, GraphX
 - Использование ускорителей
 - Combinatorial BLAS
 - Neo4j, ... ?
-
- Реализация должна быть масштабируемой **по памяти** и обладать **сильной масштабируемостью**
 - Для экзамена автоматом реализация должна обладать особенностью
-
- **Дедлайны:** 1.04 – корректность
 - 6.05 - масштабируемость

Big Data Landscape

Infrastructure

NoSQL / NewSQL Databases



Hadoop Related



MPP Databases



Crowdsourcing



Cluster Services



Management / Monitoring



Storage



Security



Monitoring



Analytics

Analytics Solutions



Data Visualization



Statistical Computing



Sentiment Analysis



Location / People / Events



Real-Time



Crowdsourced



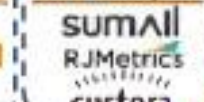
Social Media



IT Analytics



SMB Analytics



Applications

Ad Optimization



Publisher Tools



Marketing



Industry Applications



Data Sources

Data Marketplaces



Data Sources



Personal Data



Cross Infrastructure / Analytics



Open Source Projects

Framework



Programmability



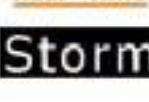
Data Access



Coordination / Workflow



Real - Time



Statistical Packages






Machine Learning

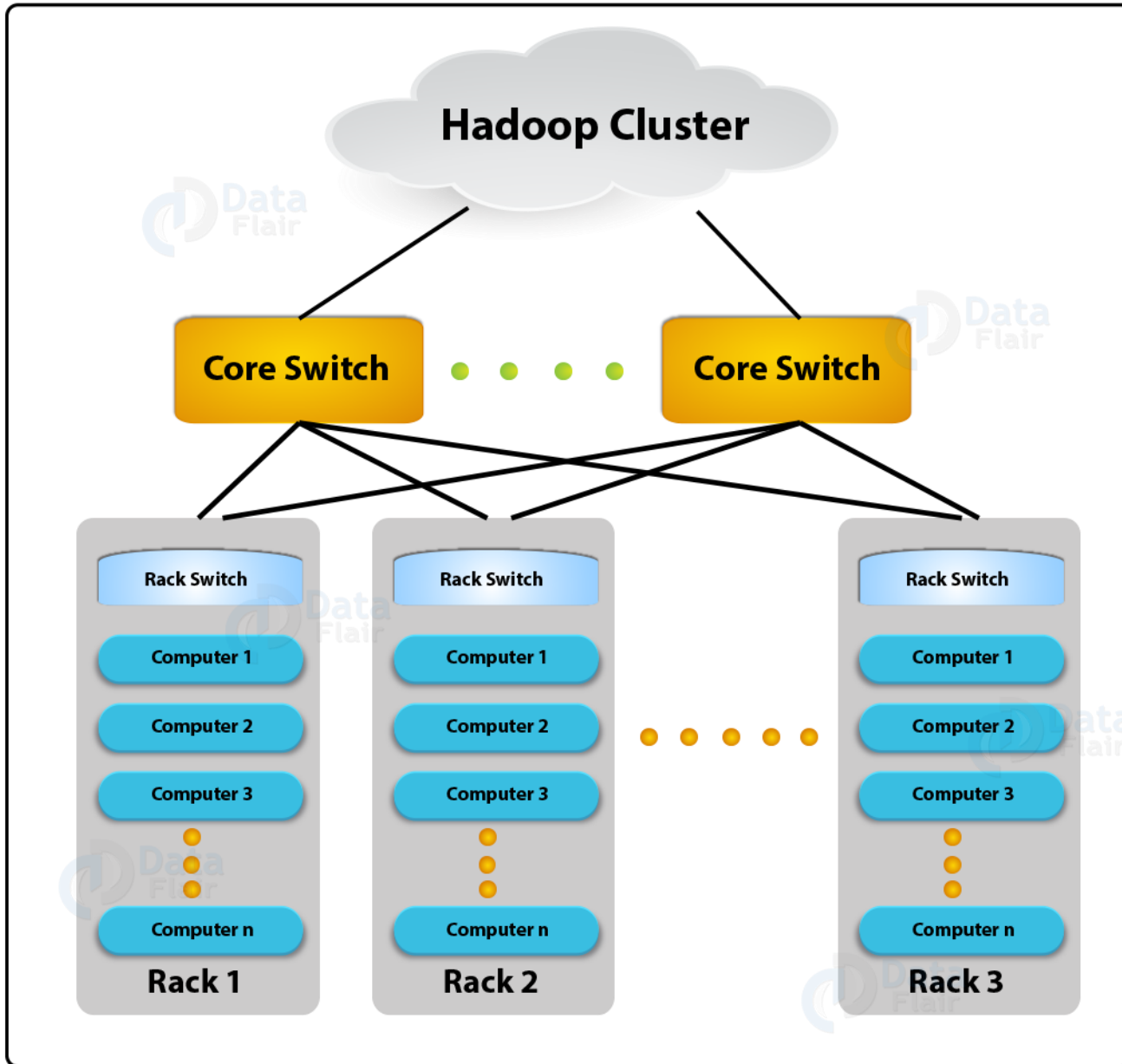


Apache Hadoop



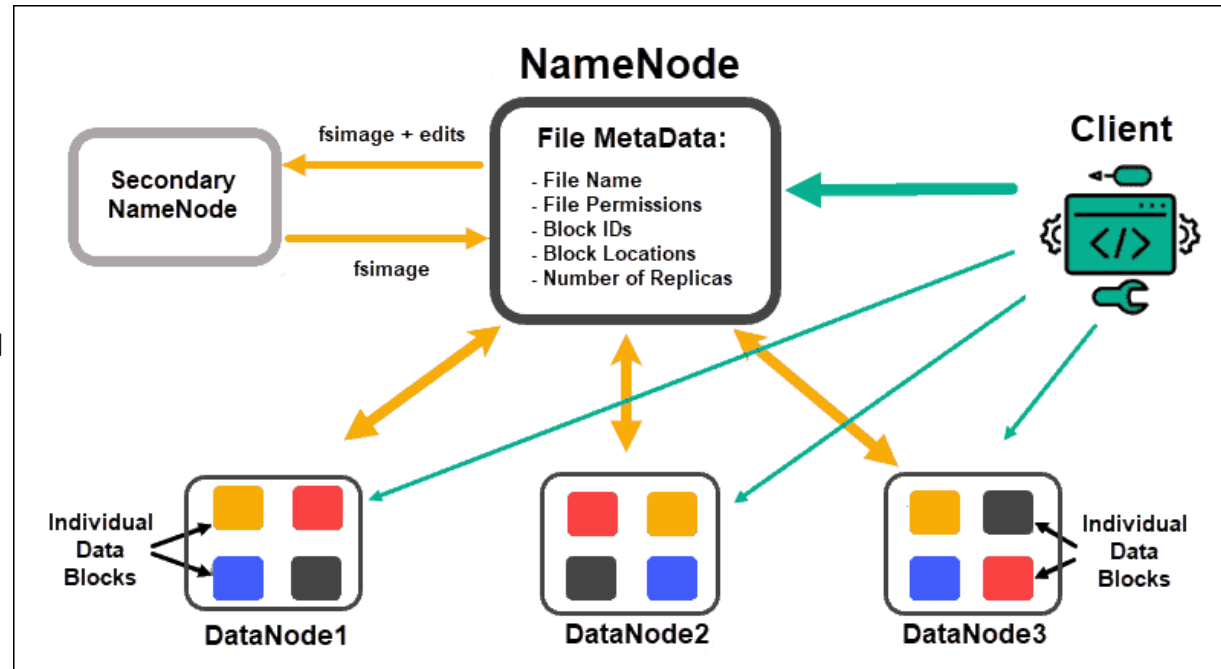
	1 PB	Hadoop/Spark Cluster
	1 TB	Hadoop/Spark Cluster
•  HDFS	100 GB	Postgres, Hadoop/Spark Cluster
•  Map Reduce	10 GB	pandas, Spark, Postgres
	GB	pandas, Spark, Postgres, CLI
•  APACHE Spark	MB	Excel, pandas, Postgres, CLI
	KB	Excel, CLI

Hadoop Cluster



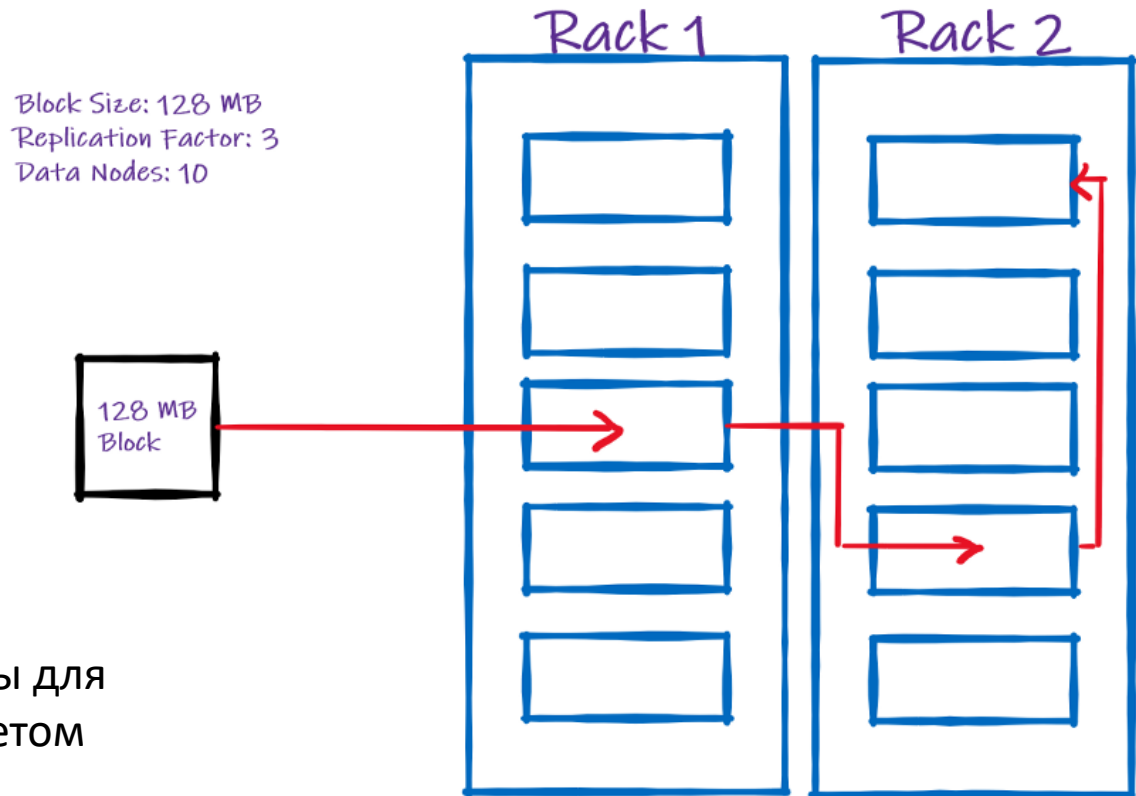
HDFS

- Работает полностью в user-space
- Отказоустойчивая
- Горизонтально высокомасштабируемая
- Разработана для обработки небольшого числа больших файлов
- Обеспечивает потоковую обработку: 1 запись, большое число чтений



- Разбиение на блоки по size (64 МБ), коэффициент репликации
- Узел имен: метаданные - информация о иерархии, о распределении
 - Номер стойки
- Standby NameNode (Hadoop 2.0)

HDFS – репликация



Алгоритм репликации:

- NameNode выбирает узлы для размещения реплик с учетом балансировки
- 1-я реплика размещается на первом узле из списка
- 2-я реплика копируется на другой узел
- И так далее

LISP

- 1958 John McCarty (MIT) Lisp (LISt Processor)
- Lisp – Функциональный язык программирования
- **(map f list)**
- **(reduce q list)**

- **(map square '(1 2 3 4))**
-> (1 4 9 16)
- **(reduce + '(1 4 9 16))**
-> 30
- **(reduce + (map square '(1 2 3 4)))**
-> 30

LISP

- 1958 John McCarty (MIT) Lisp (LISt Processor)
- Lisp – Функциональный язык программирования
- **(map f list)**
- **(reduce q list)**

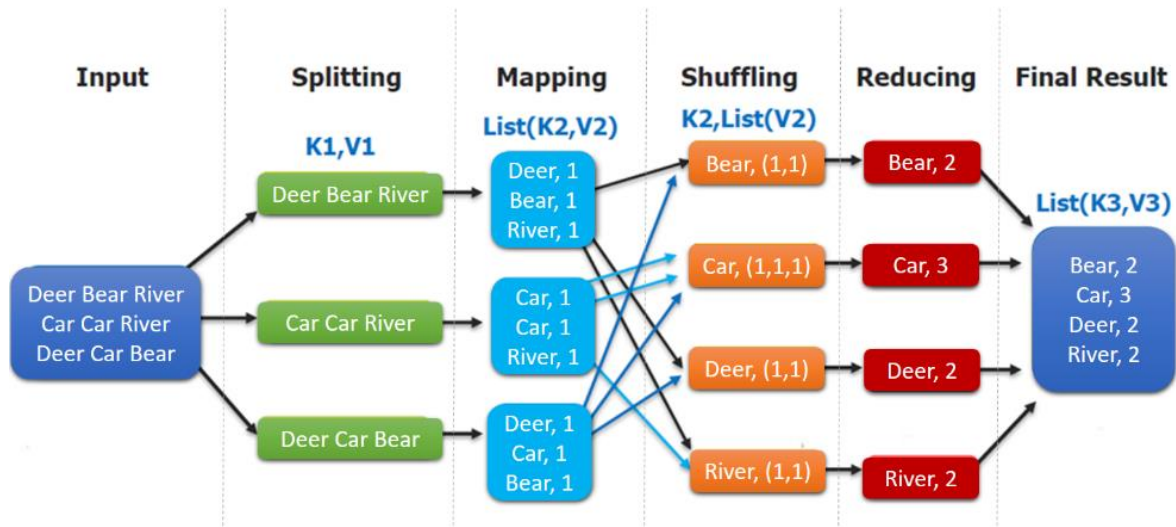
- **(map square '(1 2 3 4))**
-> (1 4 9 16)
- **(reduce + '(1 4 9 16))**
-> 30
- **(reduce + (map square '(1 2 3 4)))**
-> 30

Функциональное программирование:

- Функция – объект 1 класса (можно сохранять как значение, передавать как аргумент, возвращать значение)
- Отсутствие side-эффектов
- Неизменяемость значений переменных

MapReduce 2002 [Dean, 2004]

- 1. Input.** Получение входных данных на главном узле (master node)
- 2. Split.** Главный узел делит входные данные на части (K1) и передает рабочим узлам (worker node)
- 3. Map.** Каждый рабочий узел применяет функцию Map к каждому K1 (локальные данные) и записывает результат в формате «ключ-значение» (K2, V2) во временное хранилище
- 4. Shuffle (&Sort).** Рабочие узлы перераспределяют данные на основе ключей, чтобы все данные с одним ключом K2 лежали на одном рабочем узле.
- 5. Reduce.** Рабочие узлы применяют функцию Reduce к каждой группе результатов с одинаковым K2 и передает в главный узел

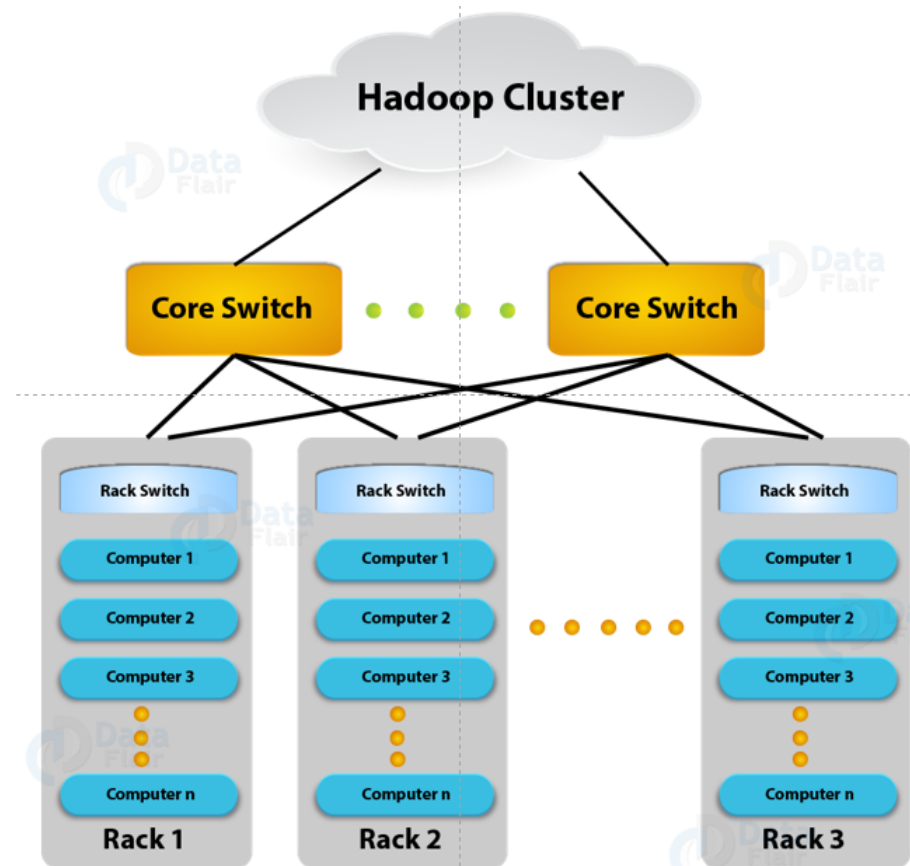


MapReduce

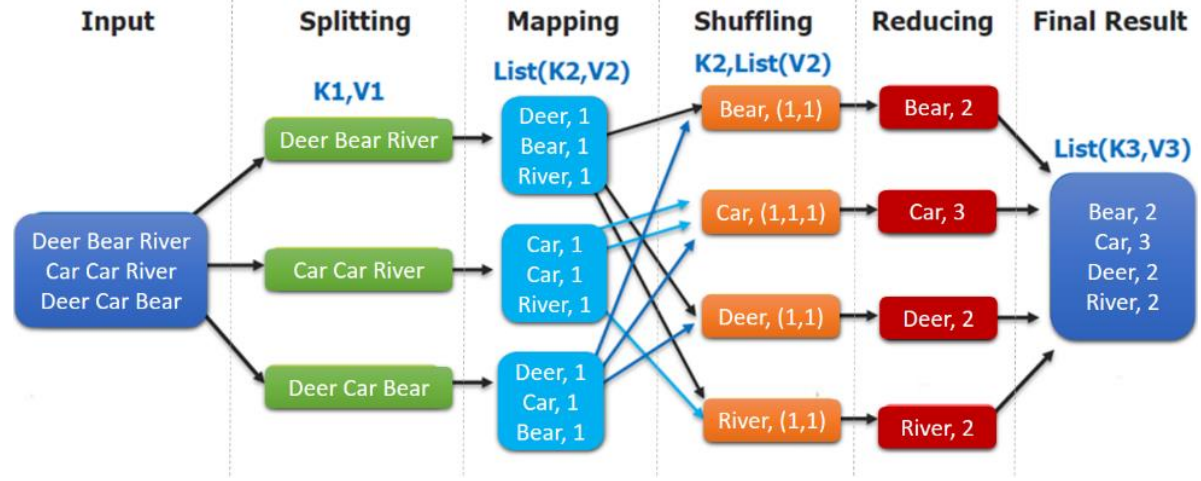
- **map (in_key, in_value)**
 - > (out_key, intermediate_value) list
- **reduce (out_key, intermediate_value list)**
 - > out_value list

Основные преимущества:

- Автоматическое распределение данных и распараллеливание вычислений
- Балансировка нагрузки
- Отказоустойчивость
 - Повтор выполнения операций в случае отказов



Пример

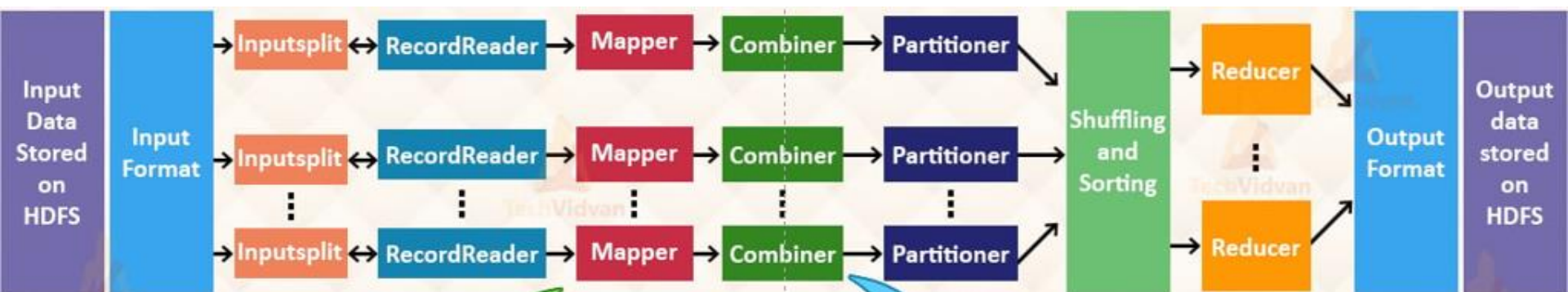


Java

```
public void map(LongWritable key, Text value, Context con) throws ... {  
    String line = value.toString();  
    String[] words=line.split(",");  
    for (String word: words ) {  
        Text outputKey = new Text(word.toUpperCase().trim());  
        IntWritable outputValue = new IntWritable(1);  
        con.write(outputKey, outputValue);  
    }  
}
```

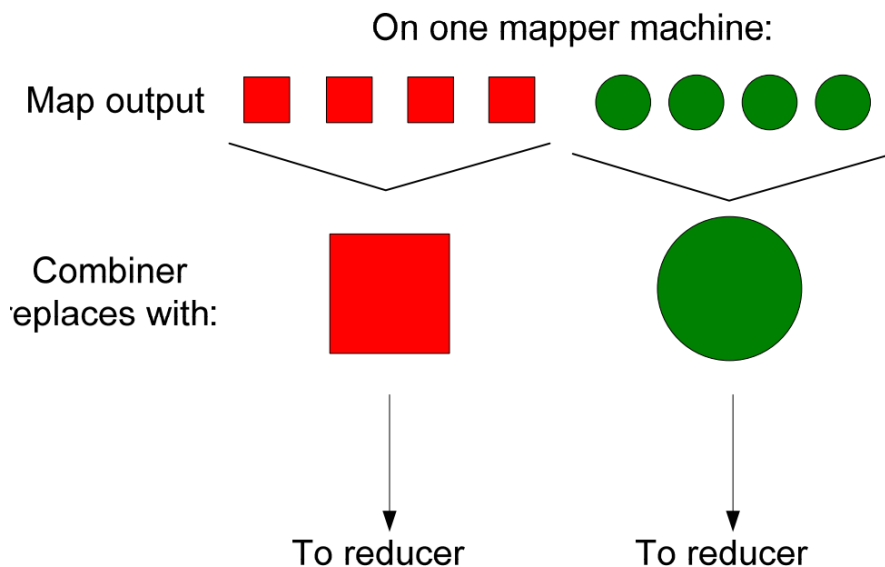
```
public void reduce(Text word, Iterable<IntWritable> values, Context con) throws ... {  
    int sum = 0;  
    for (IntWritable value : values) {  
        sum += value.get(); }  
    con.write(word, new IntWritable(sum));  
}
```


MapReduce



Оптимизации:

- Combiner – запуск на узлах mapper после фазы map (локальный mini reduce) для экономии пропускной способности



MapReduce

Реализации MapReduce:

- Google File System
- Apache Hadoop
- Yahoo
- Amazon AWS

Недостатки MapReduce:

- Трудоемкость программирования в чистой парадигме MapReduce
- Каждый раз результаты складываются на диск, производительность

MapReduce



Pregel

Giraph

Dremel

Drill

Tez

Impala

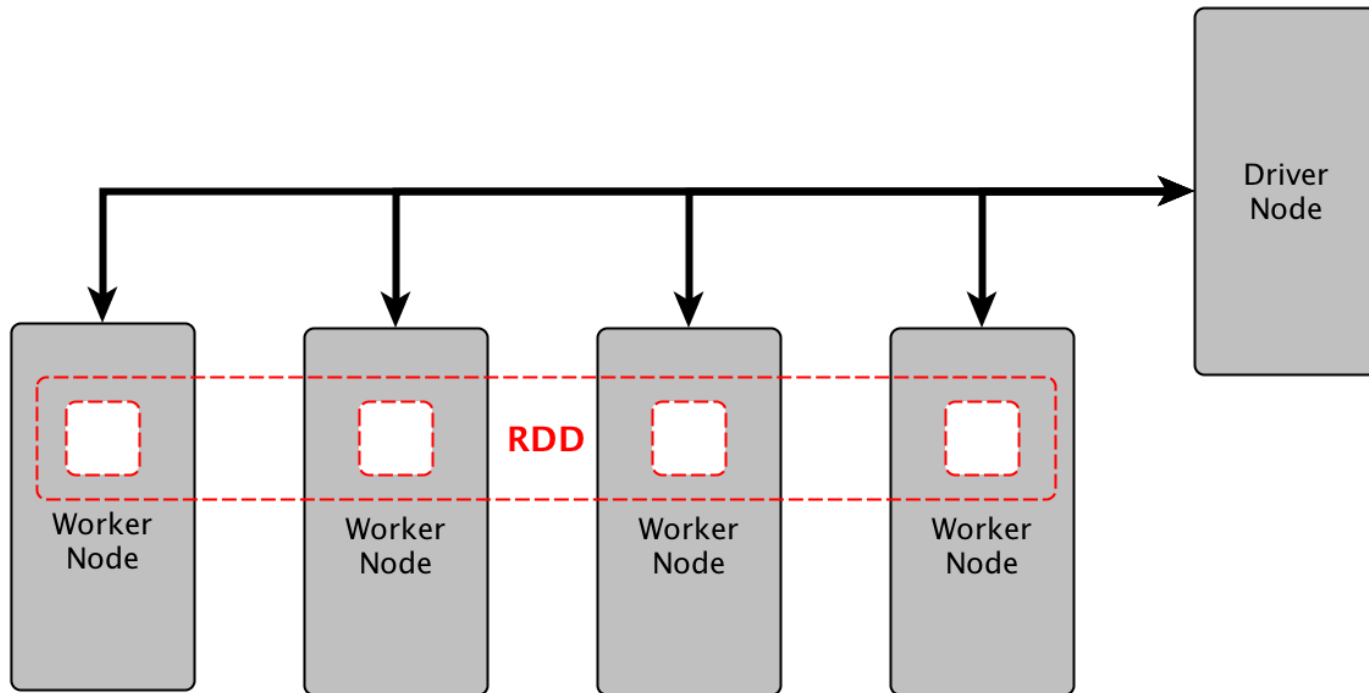
GraphLab

Storm

S4

Spark [Zaharia, 2010]

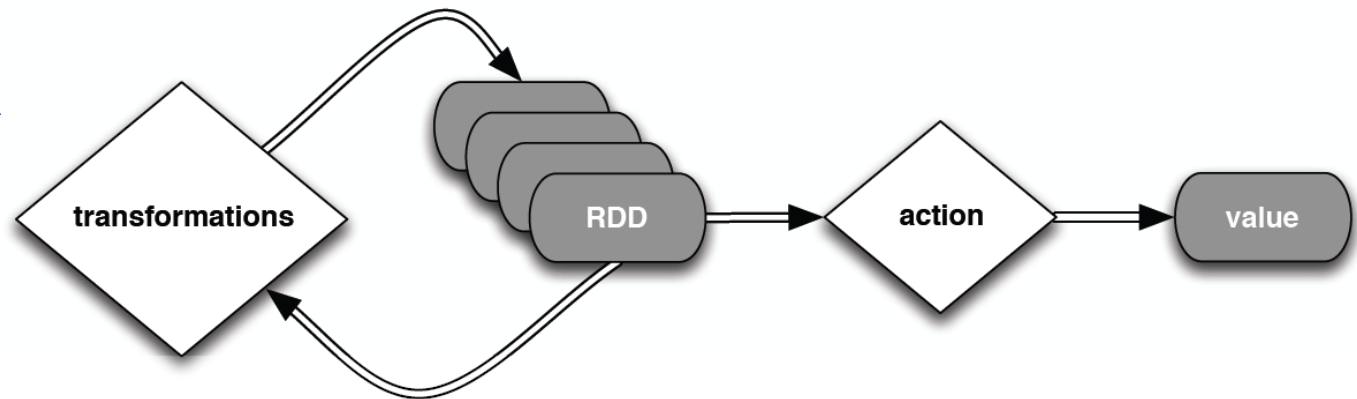
- RDD (Resilient Distributed Dataset) – главная абстракция в Spark:
 - коллекция элементов
 - отказоустойчивая
 - параллельная обработка
 - read-only



[Zaharia, 2010] Zaharia M. et al. Spark: Cluster computing with working sets //2nd USENIX Workshop on Hot Topics in Cloud Computing (HotCloud 10). – 2010.

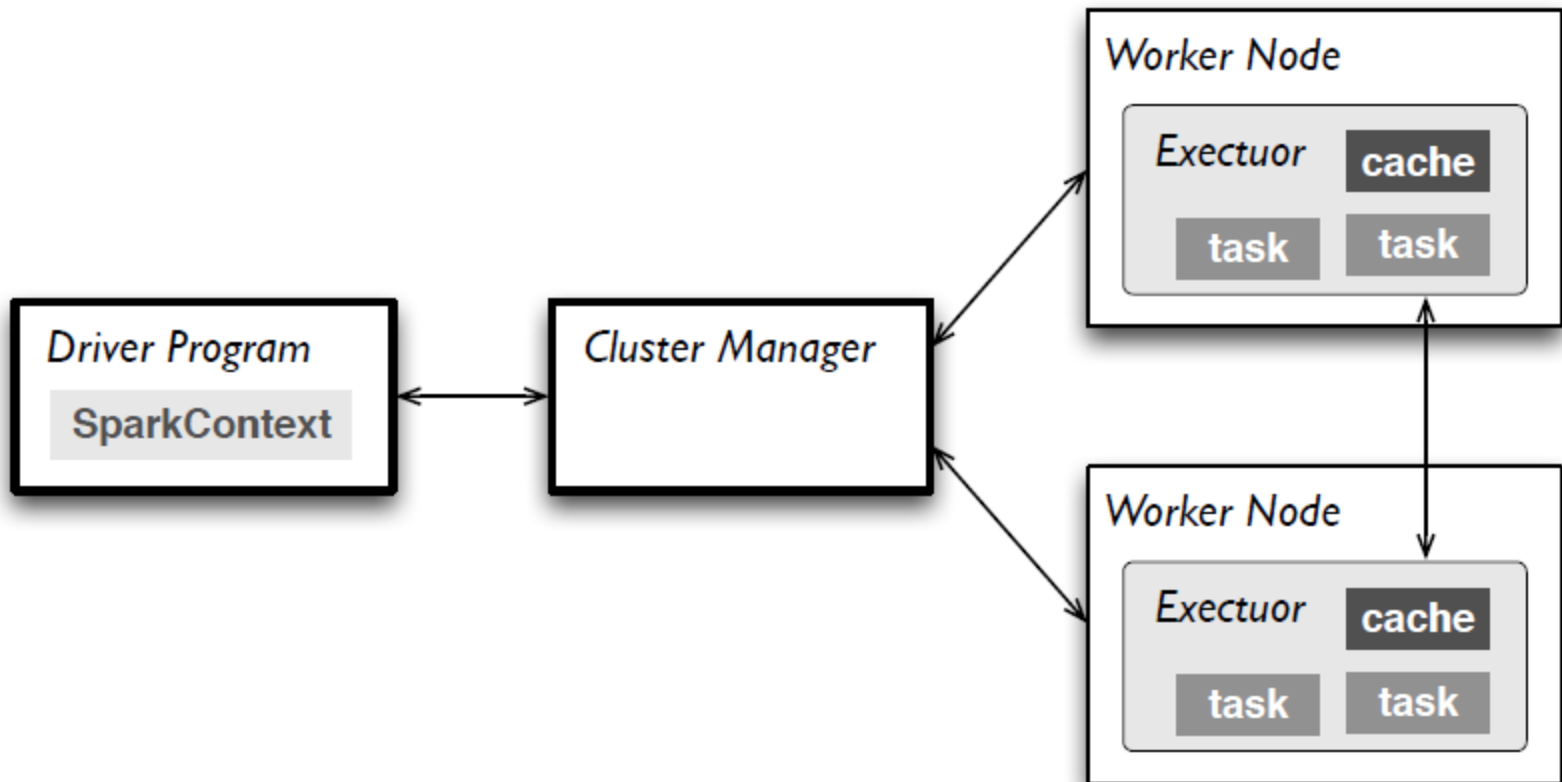
https://www.usenix.org/legacy/event/hotcloud10/tech/full_papers/Zaharia.pdf

Spark

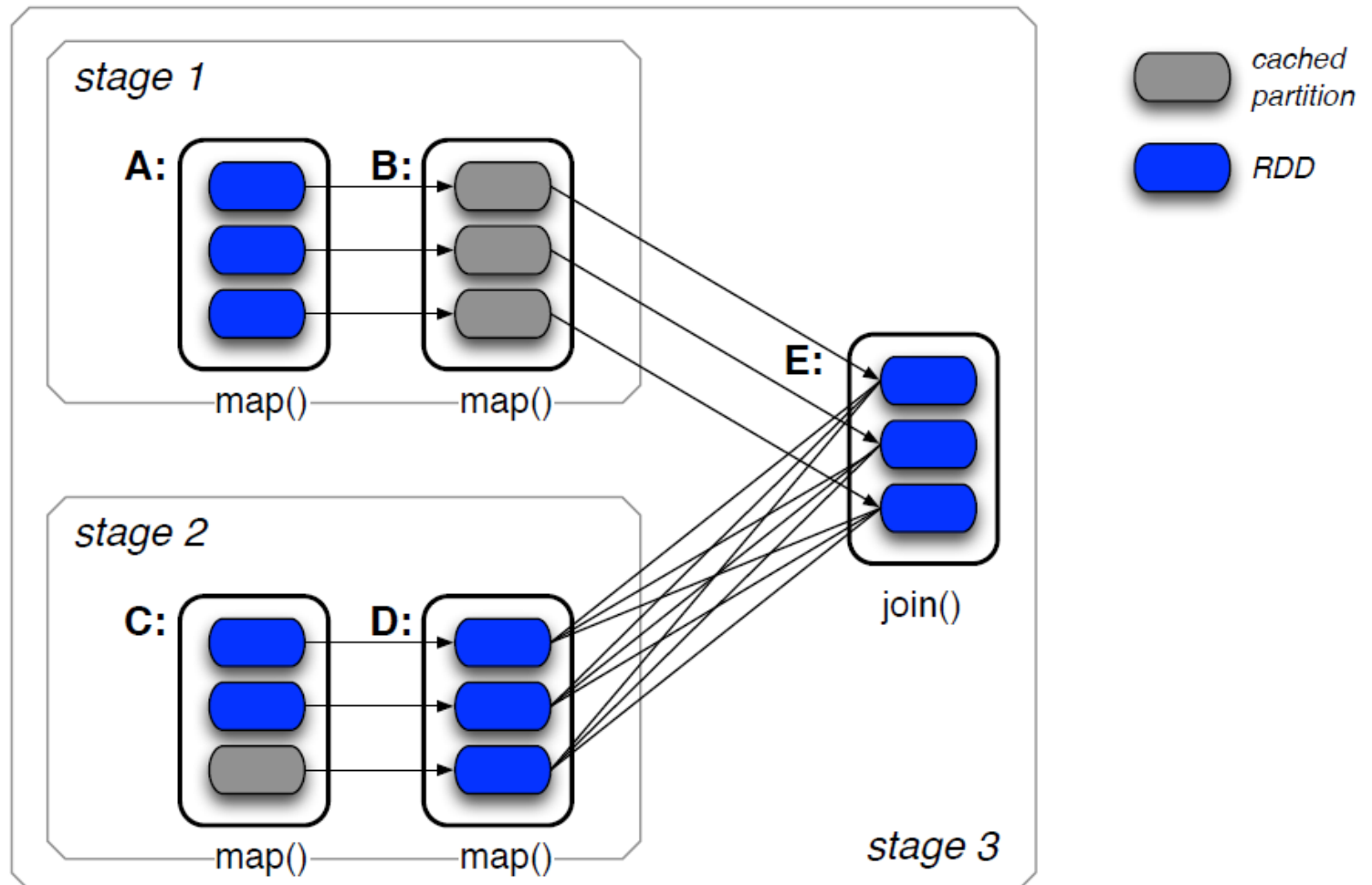


- Transformations
 - вычисления **ленивые** (оптимизация плана вычислений, восстановление)
 - map, filter, flatMap, union, distinct, groupByKey, sortByKey, reduceByKey, join
- Actions
 - reduce, count, collect, saveAs
- Persistence
 - **cache** – оставляет ленивым, но изменяет место хранения после выполнения вычислений: MEMORY_ONLY, MEMORY_AND_DISK, MEMORY_ONLY_SER, MEMORY_AND_DISK_SER, DISK_ONLY, MEMORY_ONLY_2, MEMORY_AND_DISK_2
- Разделяемые переменные
 - broadcast (read-only)
 - accumulators
(только модификация)

Spark: cluster



Spark: граф выполнения операций



Пример

Scala:

```
val f = sc.textFile("Input.file")
val wc = f.flatMap(l => l.split(" ")).map(word => (word, 1)).reduceByKey(_ + _)
wc.saveAsTextFile("wc_count.txt")
```


Spark

