

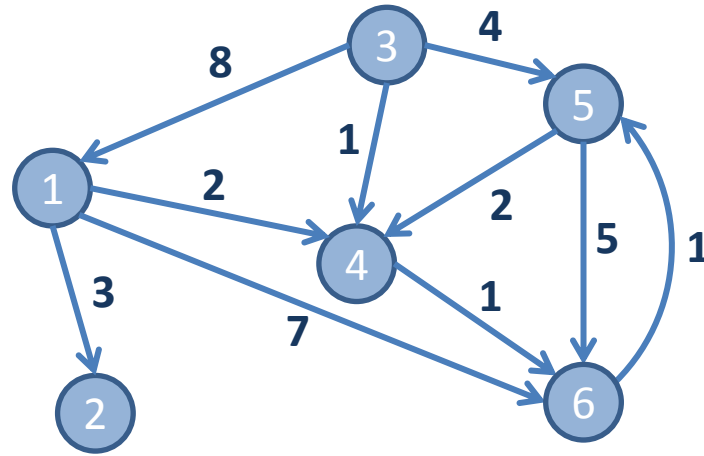
MapReduce в Hadoop

Графы

Граф как структура данных

$$G = (V, E)$$

- V представляет собой множество вершин (nodes)
- E представляет собой множество ребер (edges/links)
- Ребра и вершины могут содержать дополнительную информацию

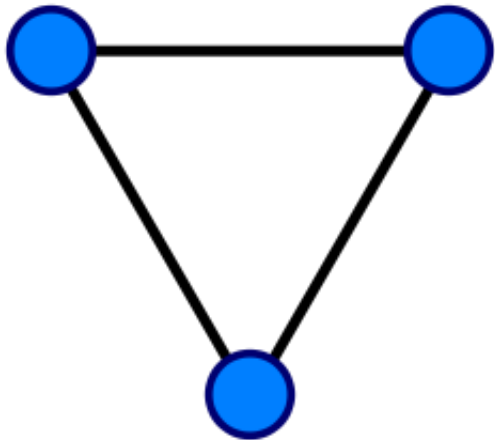


$$V = \{1, 2, 3, 4, 5, 6\}$$

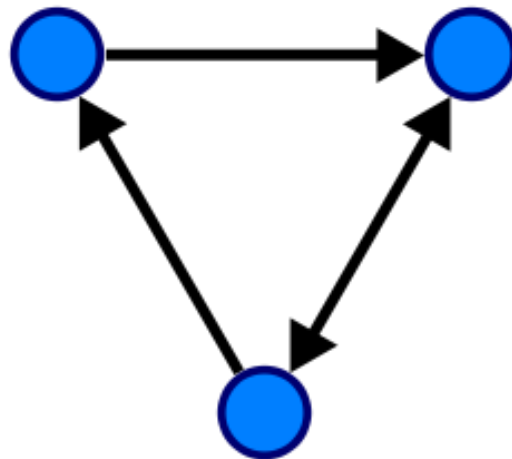
$$E = \{[1, 2], [1, 4], [1, 6], [3, 1], \dots\}$$

$$W_{1,2} = 3, W_{1,4} = 2, W_{1,6} = 7, W_{3,1} = 8, \dots$$

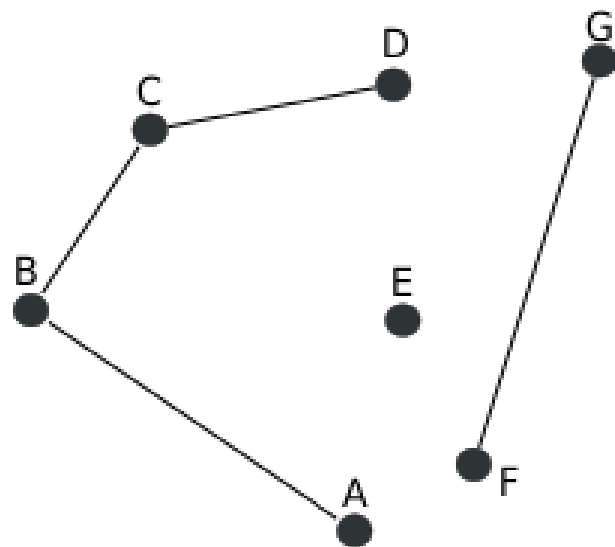
Виды графов



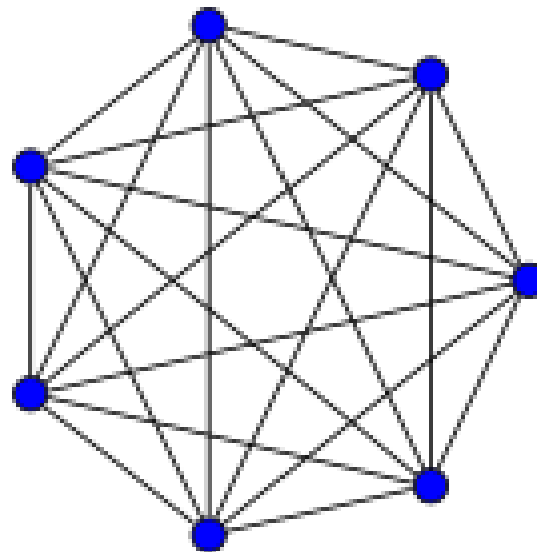
Неориентированный



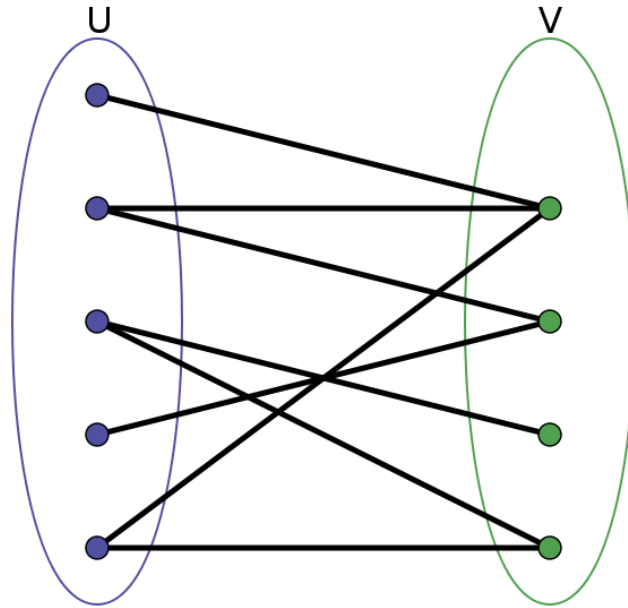
Ориентированный



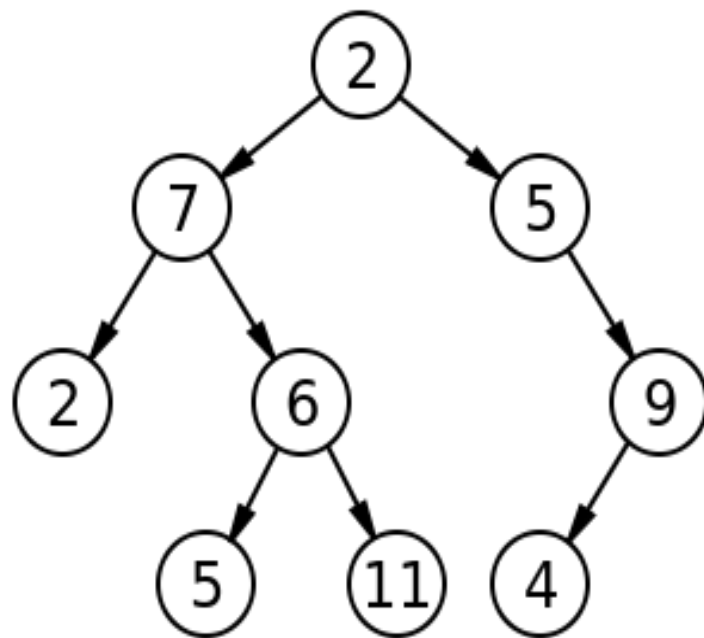
Несвязный



Полный

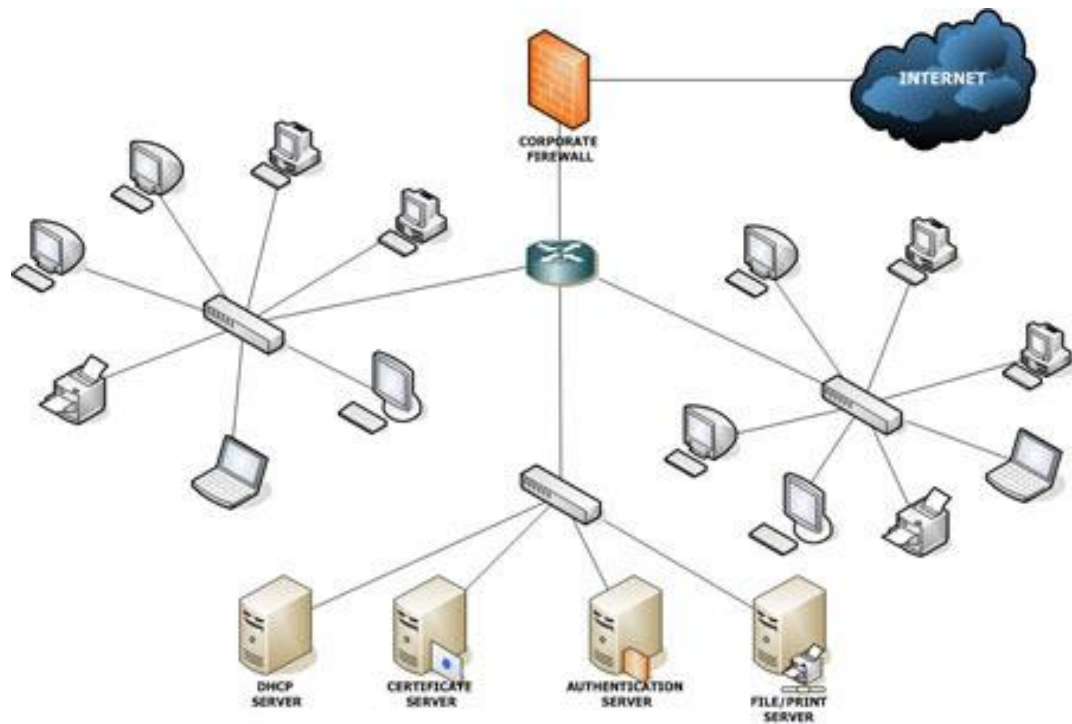


Двудольный

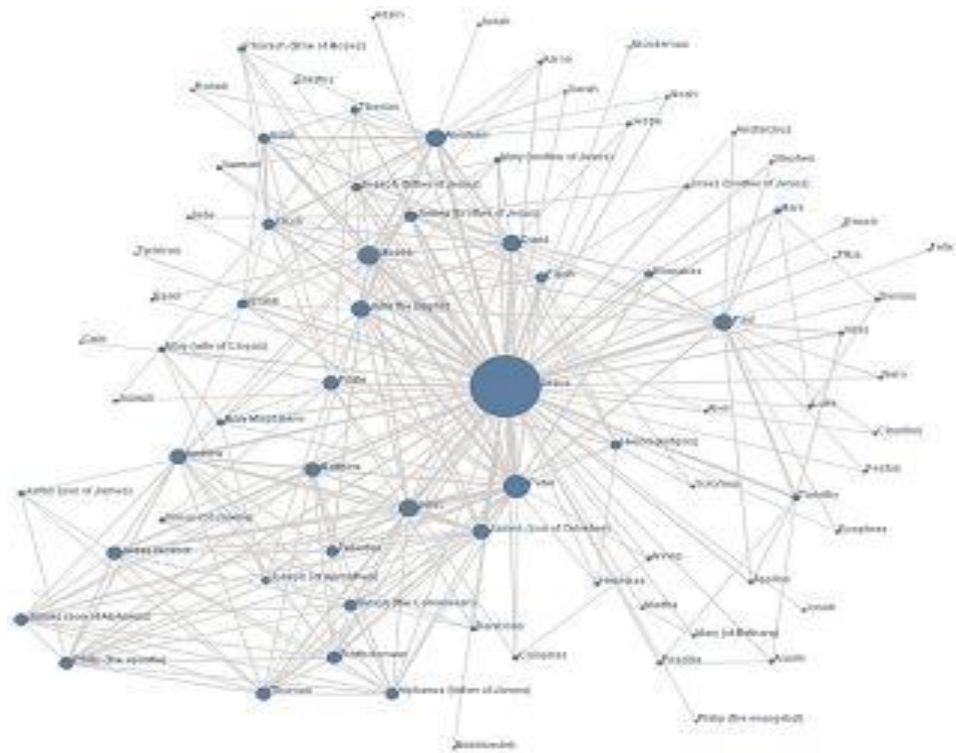


Дерево

Графы на практике



Структура компьютеров и серверов сети



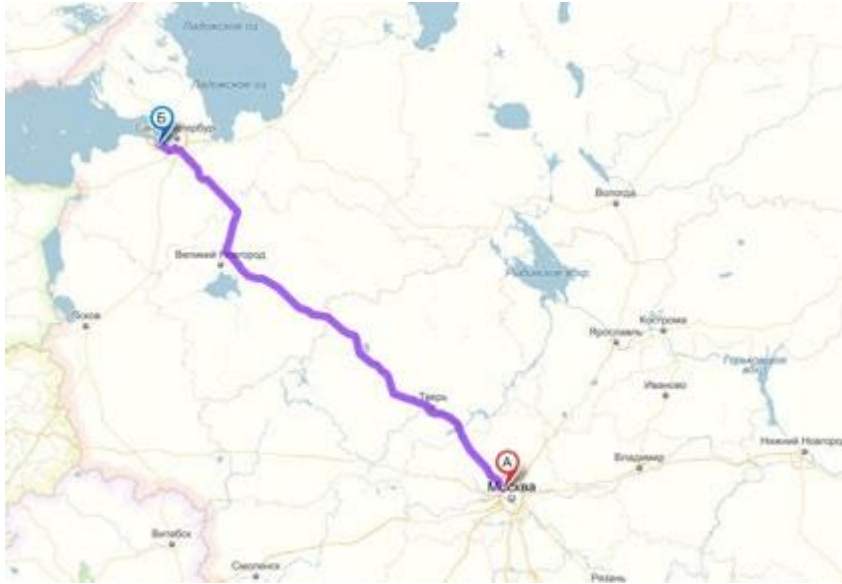
Социальные сети





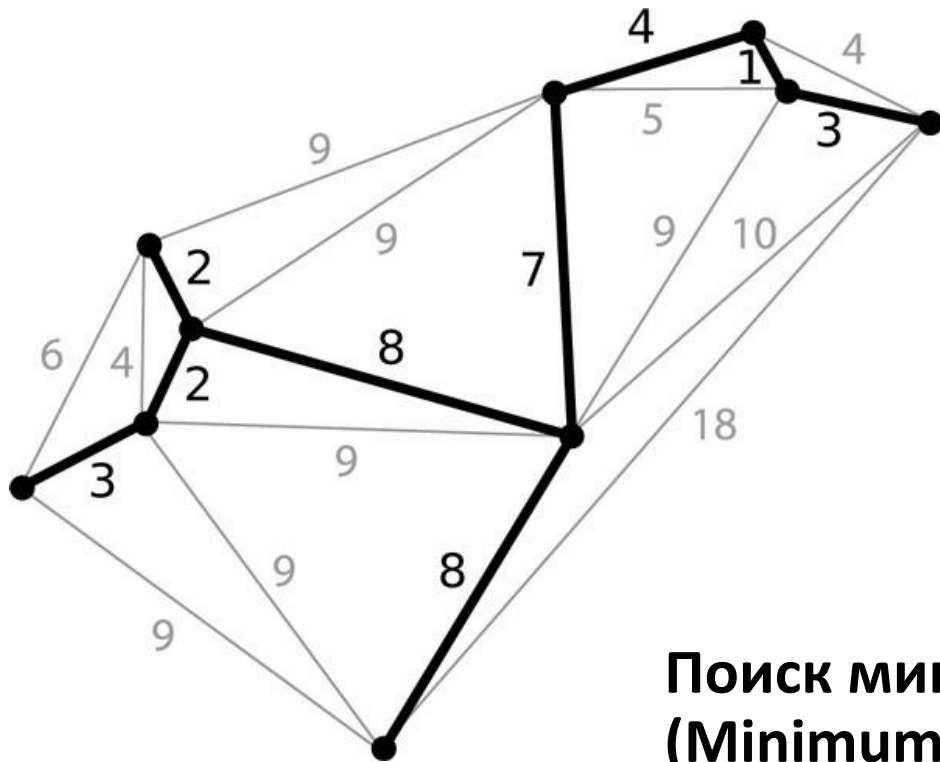
Структура дорог

Задачи и проблемы на графах



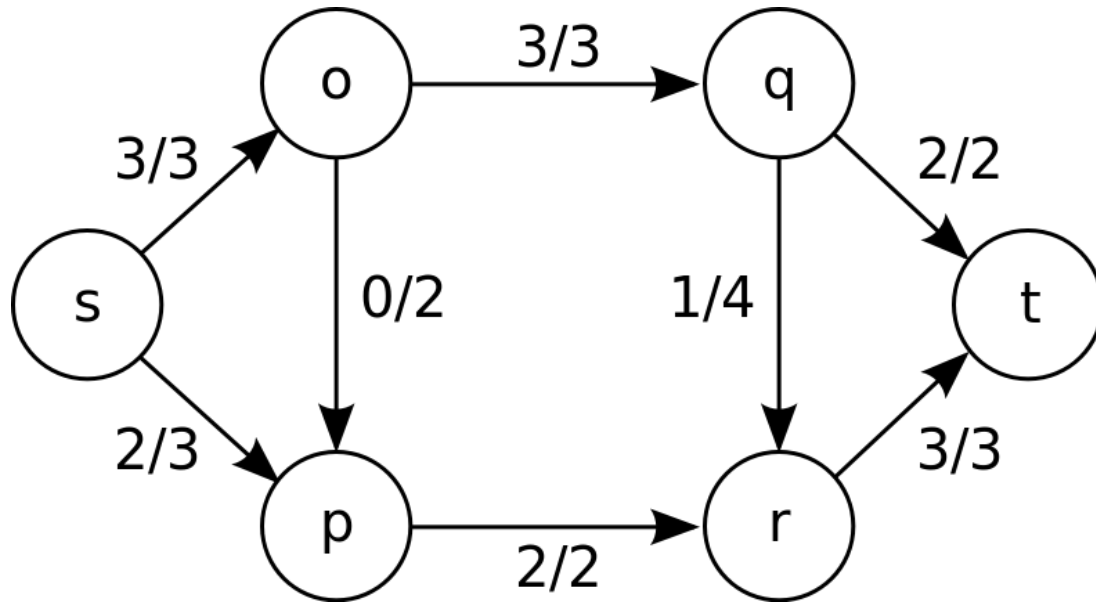
Поиск кратчайшего пути

- Роутинг траффика
- Навигация маршрута



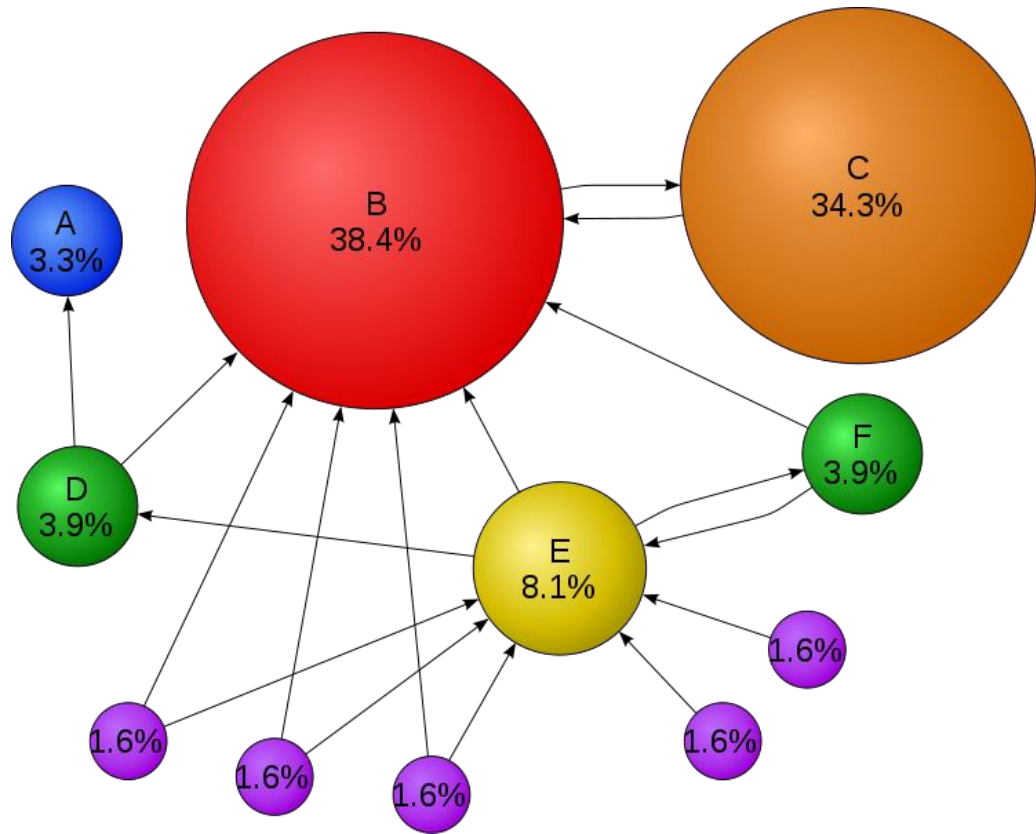
Поиск минимального остовного дерева (Minimum Spanning Tree)

- Телекоммуникационные компании



Поиск максимального потока (Max Flow)

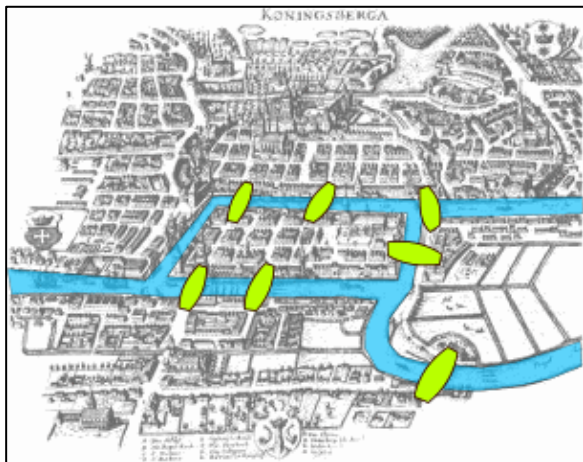
- Структура компьютеров и серверов Интернет



Алгоритмы ссылочного ранжирования

- PageRank
- HITS

Проблема семи мостов Кёнигсберга



Леонард Эйлер

Графы и MapReduce

- Вычисления на каждой вершине
- Обход графа

Ключевые вопросы:

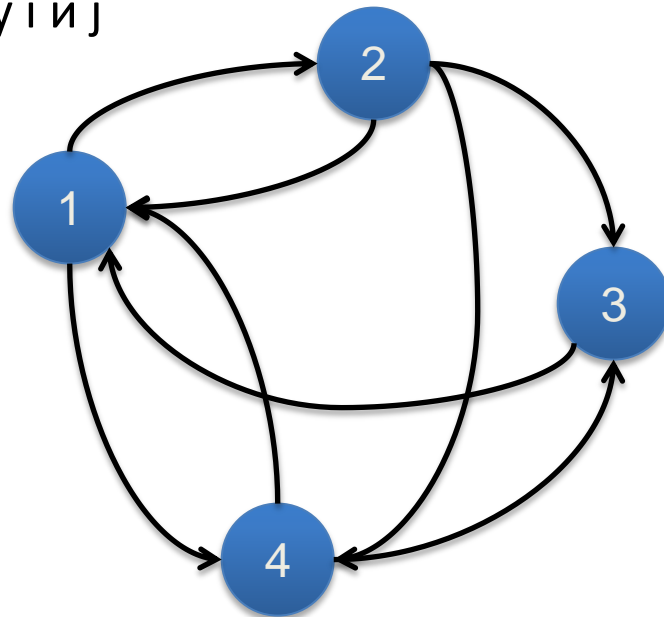
- Как представить граф в MapReduce?
- Как обходить граф в MapReduce?

Матрица смежности

Граф представляется как матрица M размером $n \times n$

- $n = |V|$
- $M_{ij} = 1$ означает наличие ребра между i и j

	1	2	3	4
1	0	1	0	1
2	1	0	1	1
3	1	0	0	0
4	1	0	1	0



Матрица смежности

Плюсы

- Удобство математических вычислений
- Перемещение по строкам и колонкам соответствует переходу по входящим и исходящим ссылкам

Минусы

- Матрица разреженная, множество лишних нулей
- Расходуется много лишнего места

	1	2	3	4
1	0	1	0	1
2	1	0	1	1
3	1	0	0	0
4	1	0	1	0

Списки смежности

Берем матрицу смежности и убираем все нули

	1	2	3	4
1	0	1	0	1
2	1	0	1	1
3	1	0	0	0
4	1	0	1	0



1: 2, 4

2: 1, 3, 4

3: 1

4: 1, 3

Списки смежности

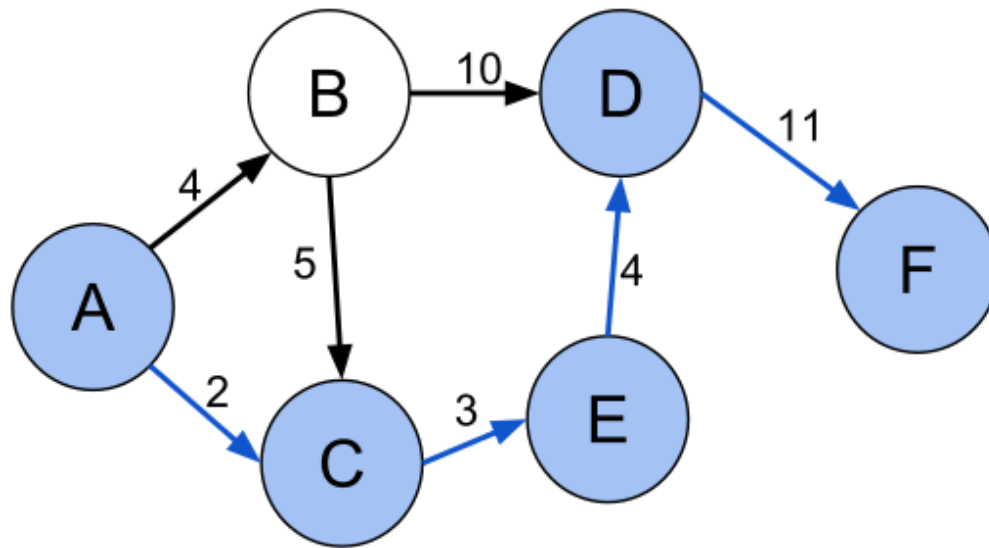
- **Плюсы**

- Намного более компактная реализация
- Легко найти все исходящие ссылки для вершины

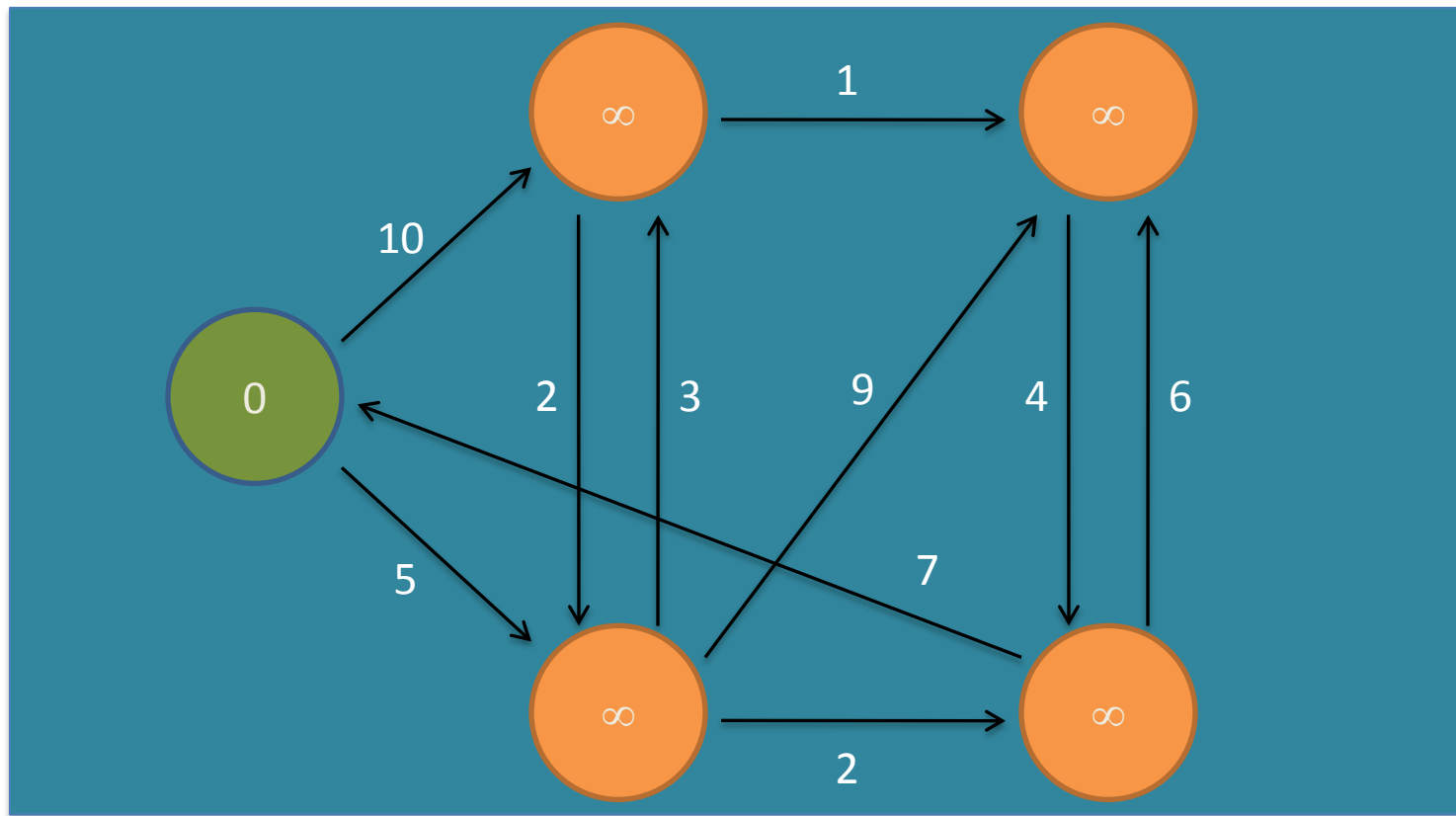
- **Минусы**

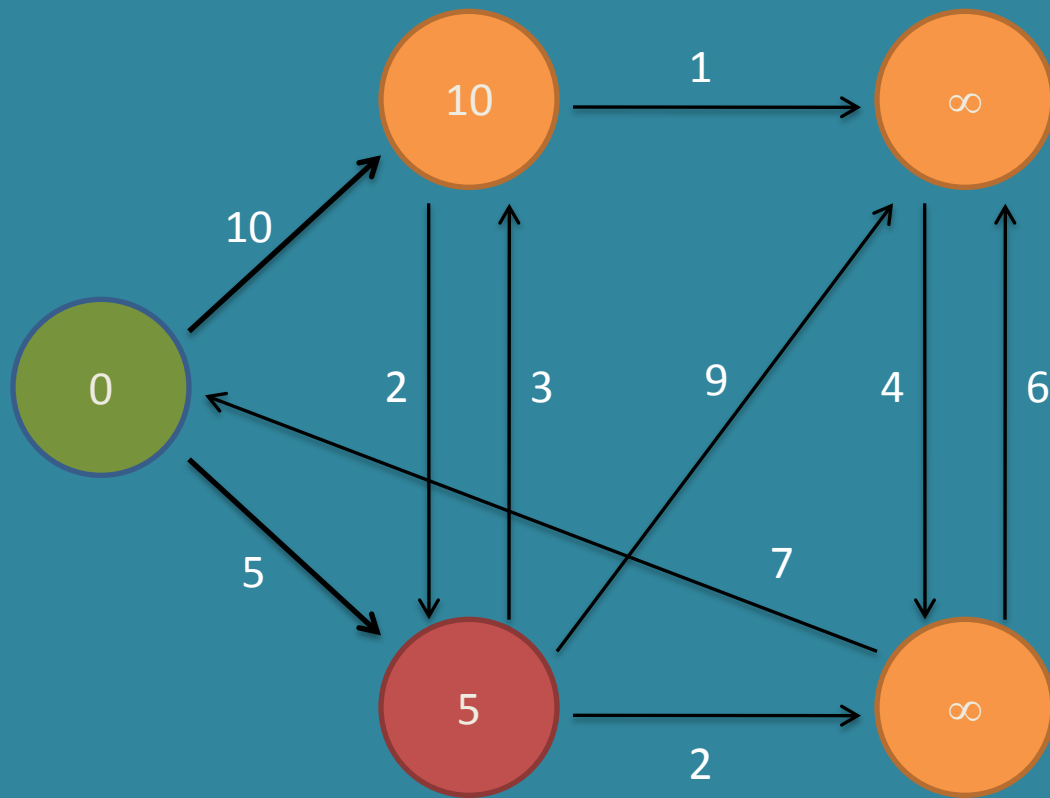
- Намного сложнее подсчитать входящие ссылки

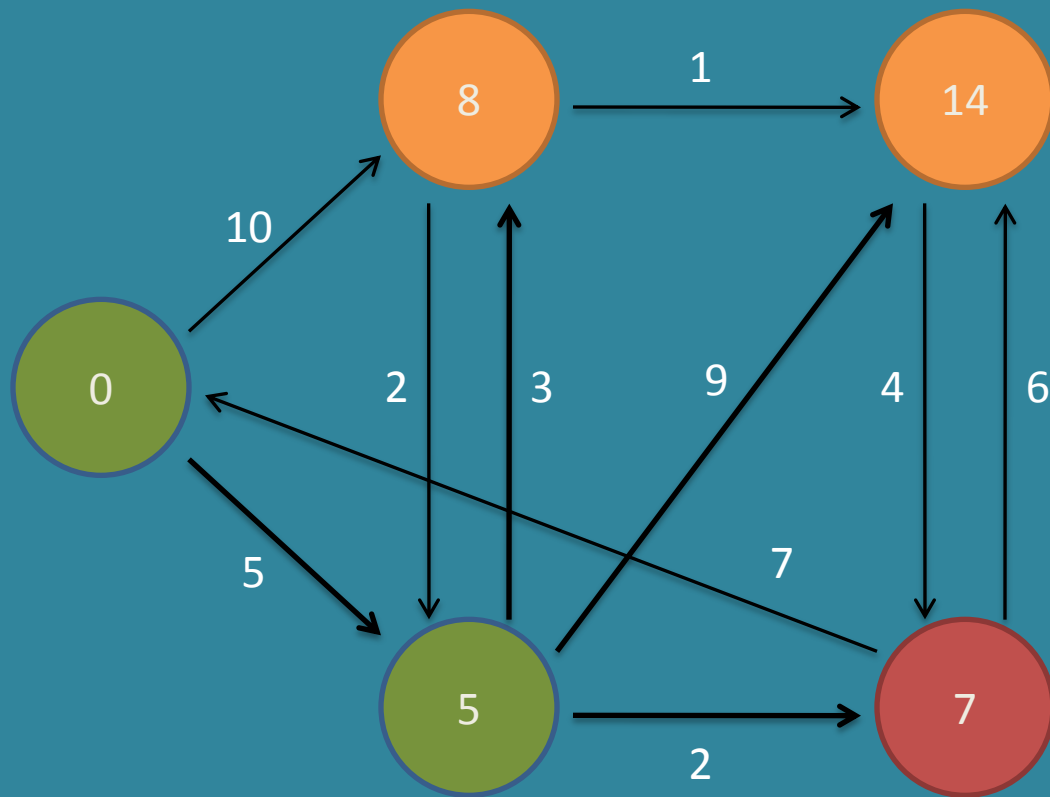
Поиск кратчайшего пути в графе

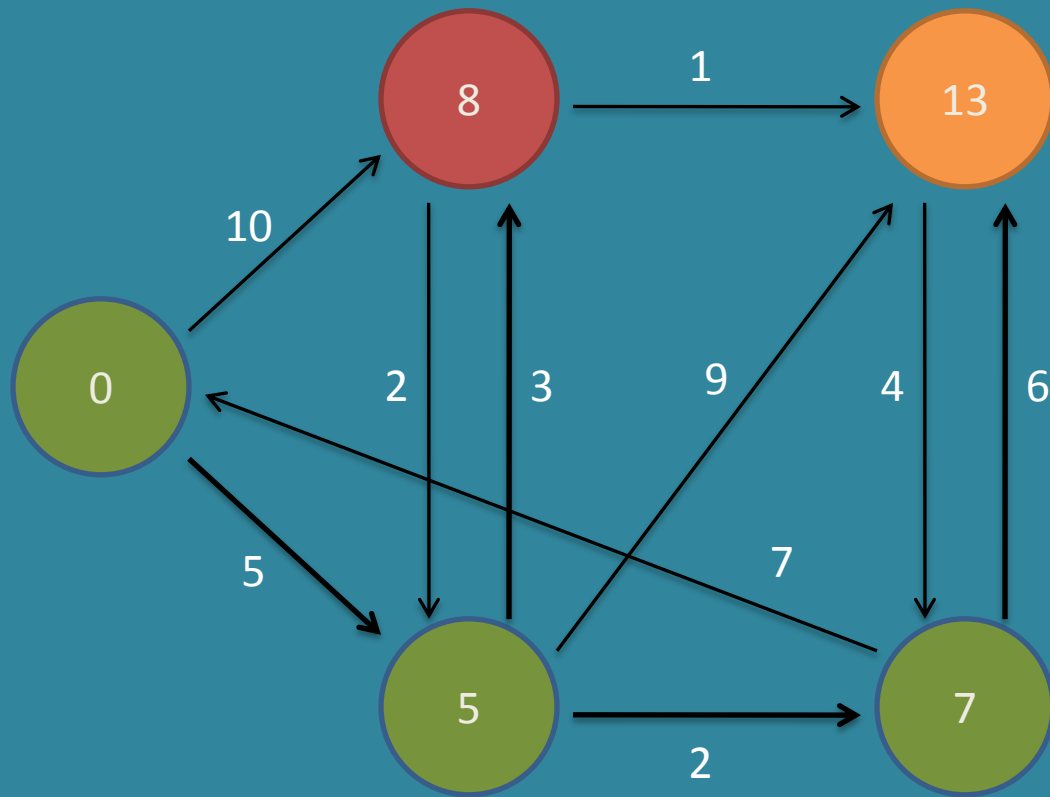


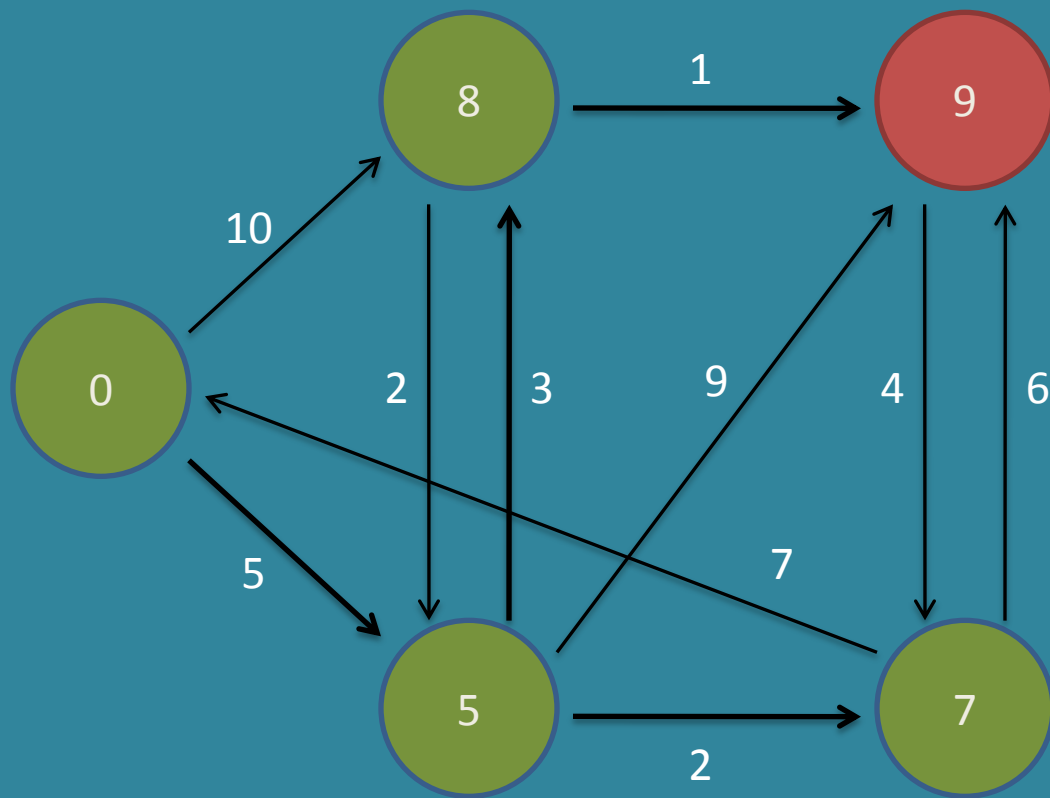
Алгоритм Дейкстры

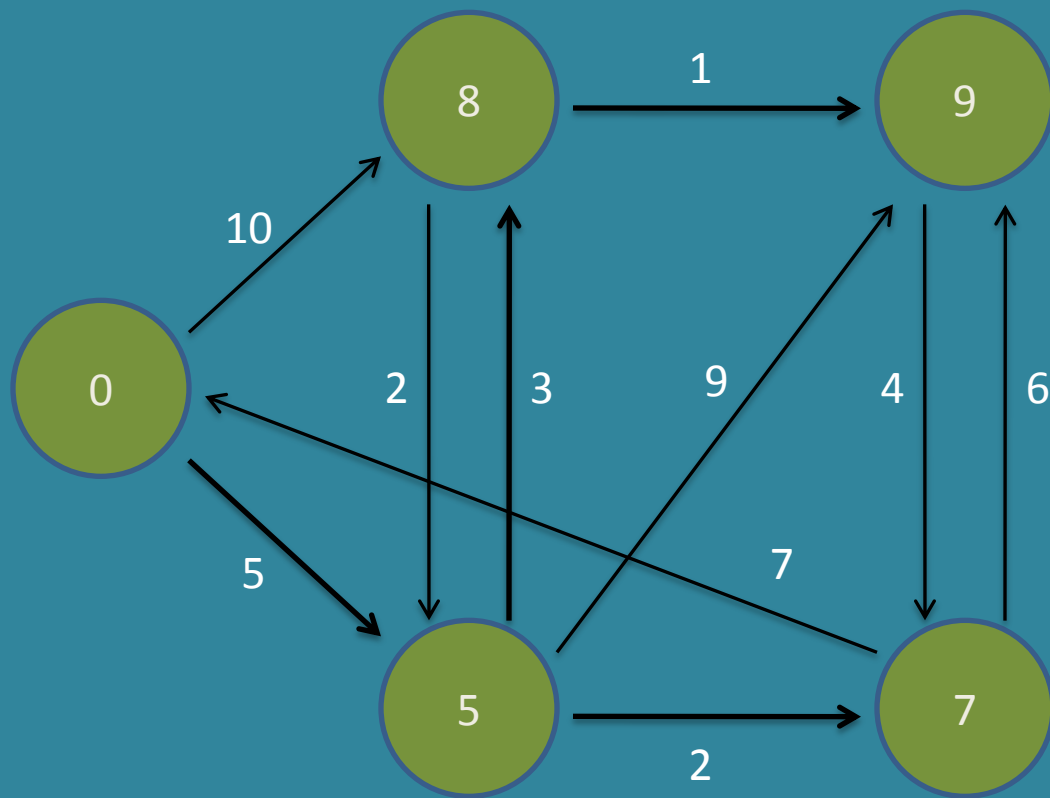










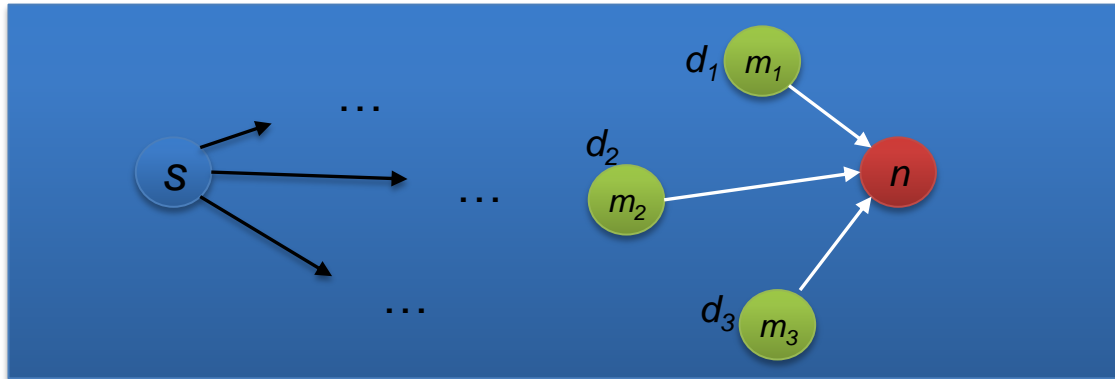


Алгоритм Дейкстры

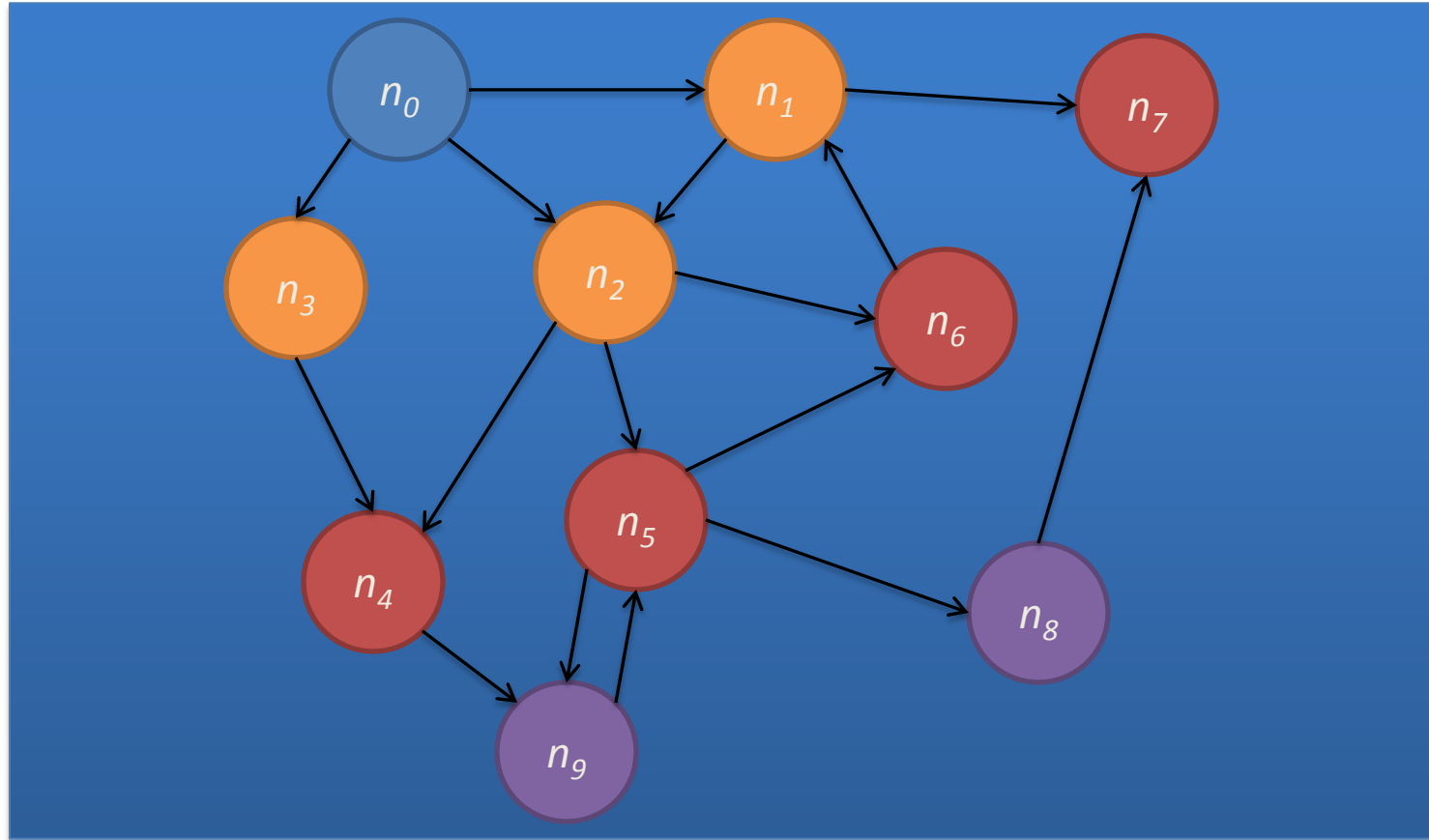
```
Dijkstra( $V, s, w$ )  
  for all vertex  $v \in V$  do  
     $d[v] \leftarrow \infty$   
   $d[s] \leftarrow 0$   
   $Q \leftarrow \{V\}$   
  while  $Q \neq \emptyset$  do  
     $u \leftarrow \text{ExtractMin}(Q)$   
    for all vertex  $v \in u.\text{AdjacencyList}$  do  
      if  $d[v] > d[u] + w(u, v)$  then  
         $d[v] \leftarrow d[u] + w(u, v)$ 
```

Поиск кратчайшего пути

- Пусть веса ребер равны 1
- Решение по индукции:
 - $DISTANCETo(s) = 0$
 - $DISTANCETo(s \rightarrow p) = 1$
 - $DISTANCETo(n) = 1 + \min(DISTANCETo(m), m \in M)$

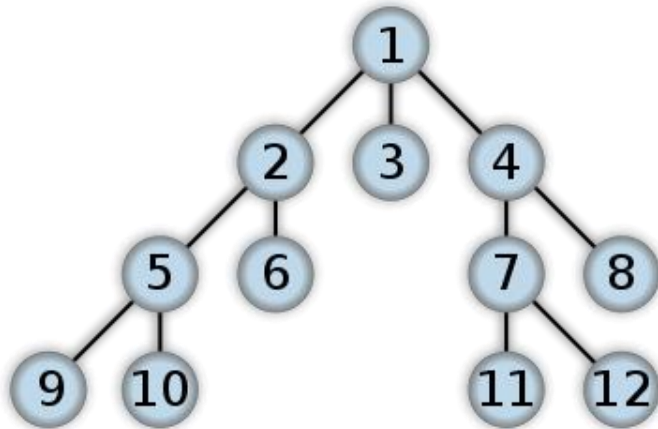


Параллельный поиск в ширину (BFS)



Breadth First Search: представление данных

- Key: вершина n
- Value: d (расстояние от начала), adjacency list (вершины, доступные из n)
- Инициализация: для всех вершин, кроме начальной, $d = \infty$



```
1 -> [0, {2, 3, 4}]
2 -> [∞, {5, 6} ]
3 -> [∞, {} ]
4 -> [∞, {7, 8} ]
5 -> [∞, {9, 10} ]
...
```

Breadth First Search: Mapper

```
mapper(key, value):
```

```
    emit(key, value)
```

```
     $\forall m \in \text{value.adjacency\_list}: \text{emit}(m, \text{value}.d + 1)$ 
```

Breadth First Search: Mapper

Mapper 1

1 \rightarrow [0, {2, 3, 4}]



1 \rightarrow [0, {2, 3, 4}]

2 \rightarrow [1, {}]

3 \rightarrow [1, {}]

4 \rightarrow [1, {}]

Mapper 2

2 \rightarrow [∞ , {5, 6}]



2 \rightarrow [∞ , {5, 6}]

5 \rightarrow [∞ , {}]

6 \rightarrow [∞ , {}]

Breadth First Search: Reducer

- **Sort/Shuffle**
 - Сгруппировать расстояния по достижимым вершинам
- **Reducer:**
 - Выбрать путь с минимальным расстоянием для каждой достижимой вершины
 - Сохранить структуру графа

Breadth First Search: Reducer

Reduce In:

$2 \rightarrow \{ [1, \{\}]], [\infty, \{5, 6\}]] \}$



Reduce Out:

$2 \rightarrow [1, \{5, 6\}]]$

BFS: псевдокод

```
class Mapper
  method Map(nid n, node N)
    d ← N.Distance
    Emit(nid n, N)      // Pass along graph structure
    for all nodeid m ∈ N.AdjacencyList do
      Emit(nid m, d + 1) // Emit distances to
reachable nodes
```

BFS: псевдокод

```
class Reducer
```

```
  method Reduce(nid m, [d1, d2, . . .])
```

```
    dmin  $\leftarrow \infty$ 
```

```
    M  $\leftarrow \emptyset$ 
```

```
    for all d  $\in$  counts [d1, d2, . . .] do
```

```
      if IsNode(d) then
```

```
        M  $\leftarrow$  d    // Recover graph structure
```

```
      else if d < dmin then
```

```
        dmin  $\leftarrow$  d
```

```
    M.Distance  $\leftarrow$  dmin    // Update shortest distance
```

```
    Emit(nid m, node M)
```

Input

```
1 -> [0, {2, 3, 4}]
2 -> [ $\infty$ , {5, 6} ]
3 -> [ $\infty$ , {} ]
4 -> [ $\infty$ , {7, 8} ]
5 -> [ $\infty$ , {9, 10} ]
...
```



Iteration 1

```
1 -> [0, {2, 3, 4}]
2 -> [1, {5, 6} ]
3 -> [1, {} ]
4 -> [1, {7, 8} ]
5 -> [ $\infty$ , {9, 10} ]
...
```



Iteration 2

```
1 -> [0, {2, 3, 4}]
2 -> [1, {5, 6} ]
3 -> [1, {} ]
4 -> [1, {7, 8} ]
5 -> [2, {9, 10} ]
...
```

...



Result

```
1 -> [0, {2, 3, 4}]
2 -> [1, {5, 6} ]
3 -> [1, {} ]
4 -> [1, {7, 8} ]
5 -> [2, {9, 10} ]
...
```

Breadth First Search: Итерации

- Каждая итерация задачи MapReduce смещает границу продвижения по графу (*frontier*) на один “hop”
 - Последующие операции включают все больше и больше посещенных вершин, т.к. граница (*frontier*) расширяется
 - Множество итераций требуется для обхода всего графа
- Сохранение структуры графа
 - Проблема: что делать со списком смежных вершин (*adjacency list*)?
 - Решение: Mapper также пишет (*n, adjacency list*)

BFS: критерий завершения

- Как много итераций нужно для завершения параллельного BFS?
- Когда первый раз посетили искомую вершину, значит найден самый короткий путь
- Равно диаметру графа
- Практическая реализация
 - Внешняя программа-драйвер для проверки оставшихся вершин с дистанцией ∞
 - Можно использовать счетчики из Hadoop MapReduce

BFS vs Дейкстра

- Алгоритм Дейкстры более эффективен
 - На каждом шаге используются вершины только из пути с минимальным весом
 - Нужна дополнительная структура данных (*priority queue*)
- MapReduce обходит все пути графа параллельно
 - Много лишней работы (brute-force подход)
 - Полезная часть выполняется только на текущей границе обхода

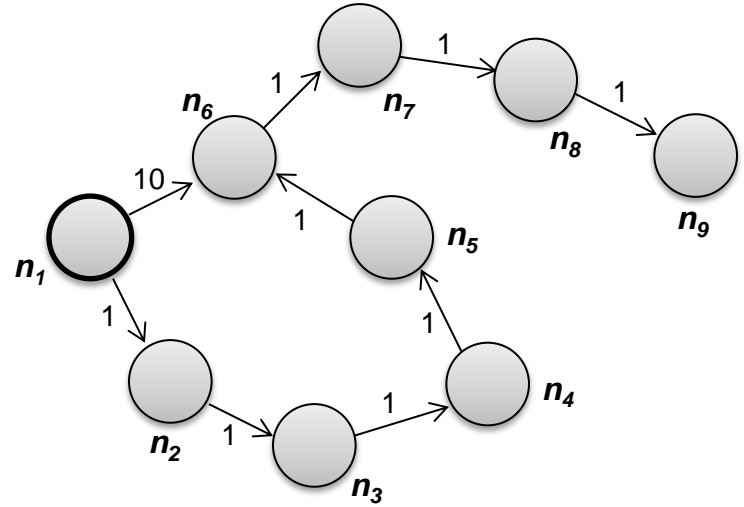
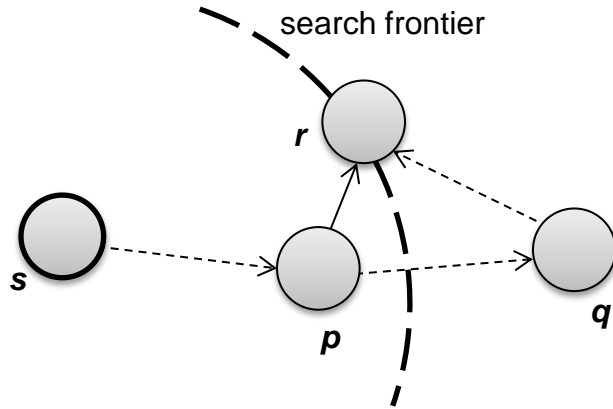
BFS: Weighted Edges

- Добавим положительный вес каждому ребру
- Простая доработка: добавим вес w для каждого ребра в список смежных вершин
 - В mapper, emit $(m, d + w_p)$ вместо $(m, d + 1)$ для каждой вершины m

BFS Weighted: критерий завершения

- Как много итераций нужно для завершения параллельного BFS (взвешенный граф)?
- Когда первый раз посетили искомую вершину, значит найден самый короткий путь
- **И это неверно!**

BFS Weighted: сложности



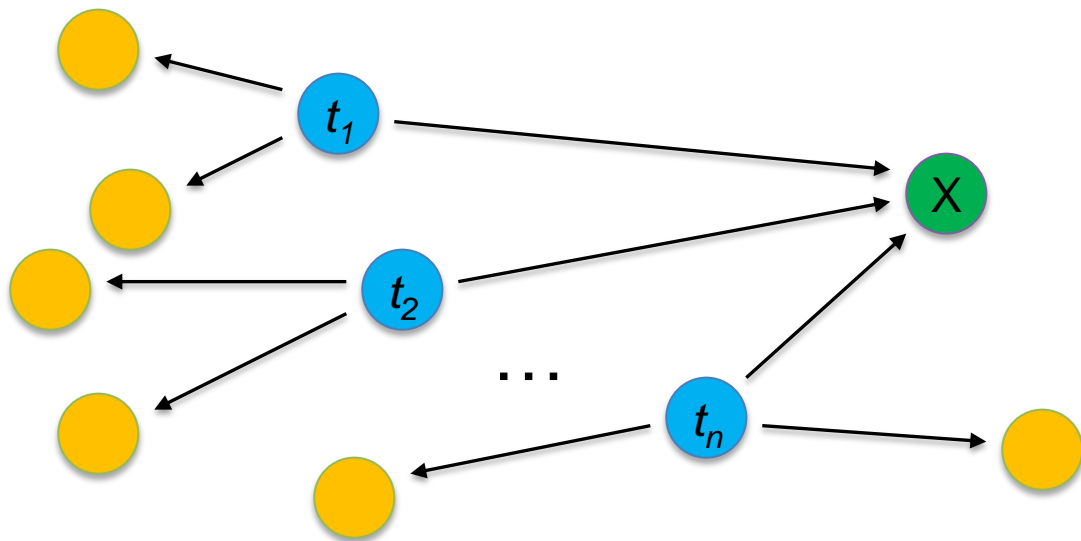
BFS Weighted: критерий завершения

- В худшем случае: $N - 1$
- В реальном мире \sim диаметру графа
- Практическая реализация
 - Итерации завершаются, когда минимальный путь у каждой вершины больше не меняется
 - Для этого можно также использовать счетчики в MapReduce

PageRank

PageRank

- Определяет важность страницы
- Характеризует кол-во времени, которое пользователь провел на данной странице
- Модель блуждающего веб-серфера
 - Пользователь начинает серфинг на случайной веб-странице
 - Пользователь произвольно кликает по ссылкам, тем самым перемещаясь от страницы к странице



$$PR(x) = \alpha \left(\frac{1}{N} \right) + (1 - \alpha) \sum_{i=1}^n \frac{PR(t_i)}{C(t_i)}$$

Вычисление PageRank

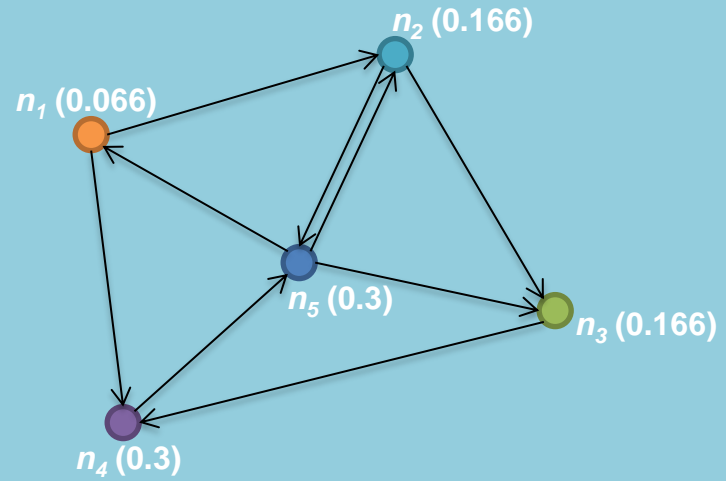
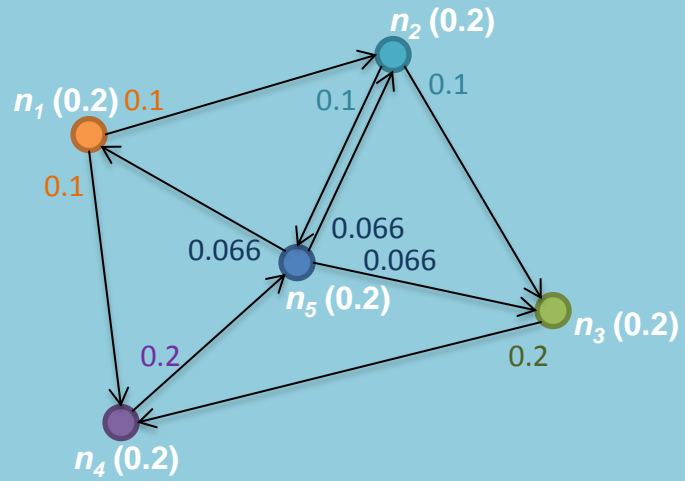
- Свойства PageRank'a
 - Может быть рассчитан итеративно
 - Локальный эффект на каждой итерации
- набросок алгоритма
 - Начать с некоторыми заданными значения PR_i
 - Каждая страница распределяет PR_i “кредит” всем страниц, на которые с нее есть ссылки
 - Каждая страница добавляет весь полученный “кредит” от страниц, которые на нее ссылаются, для подсчета PR_{i+1}
 - Продолжить итерации пока значения не сойдутся

Упрощения для PageRank

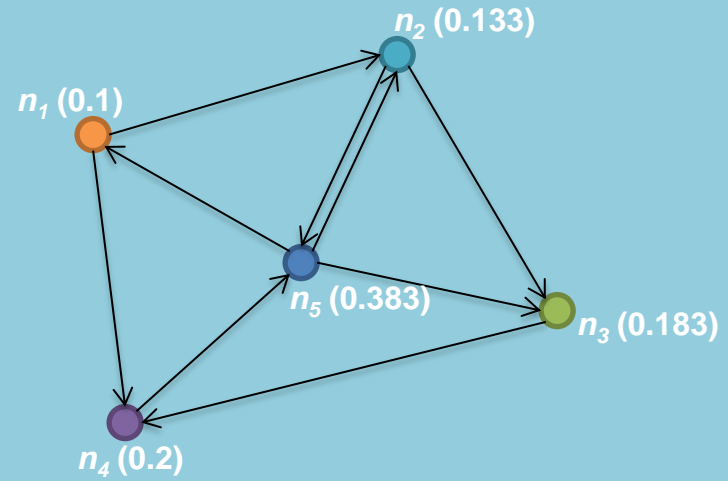
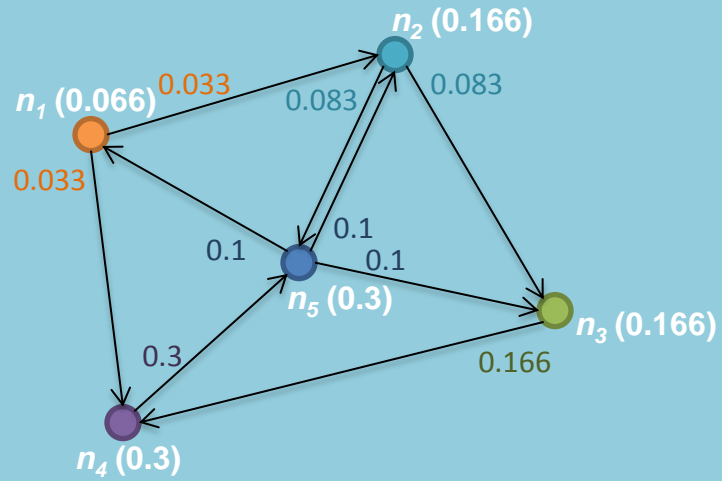
Рассмотрим простой случай

- Нет фактора случайного перехода (*random jump*)
- Нет “подвисших” вершин

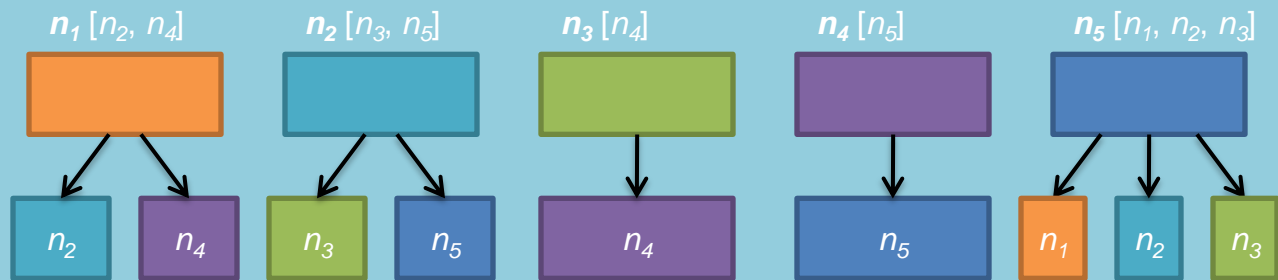
Iteration 1



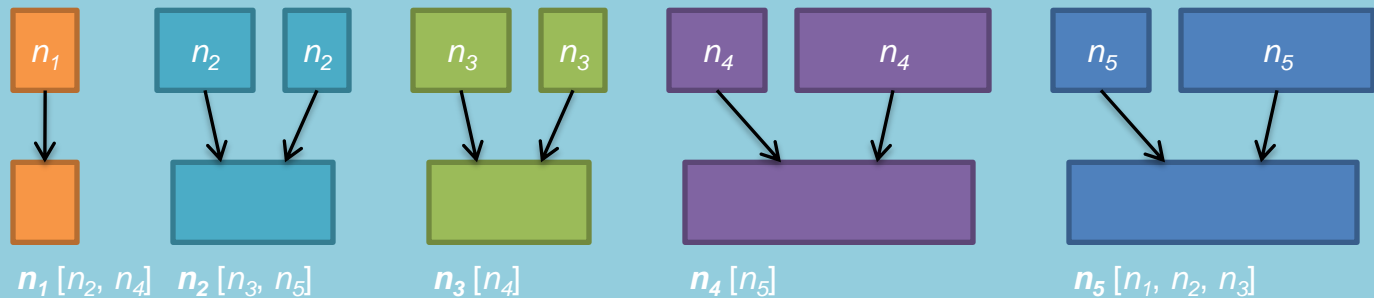
Iteration 2



Map



Reduce



PageRank: Mapper

```
class Mapper
  method Map(nid n, node N)
     $p \leftarrow N.\text{PageRank} / |N.\text{AdjacencyList}|$ 
    Emit(nid n, N)
    for all nodeid m  $\in N.\text{AdjacencyList}$  do
      Emit(nid m, p)
```

PageRank: Reducer

```
class Reducer
  method Reduce(nid m, [p1, p2, . . .])
     $M \leftarrow \emptyset$ 
    for all  $p \in \text{counts } [p1, p2, . . .]$  do
      if IsNode(p) then
         $M \leftarrow p$ 
      else
         $s \leftarrow s + p$ 
     $M.\text{PageRank} \leftarrow s$ 
    Emit(nid m, node M)
```

Полный PageRank

- Обработка “подвешенных” вершин
- Случайный переход (*random jump*)

$$p' = \alpha \left(\frac{1}{N} \right) + (1 - \alpha) \left(\frac{m}{N} + p \right)$$

Сходимость PageRank

- Продолжать итерации пока **значения** PageRank не перестанут изменяться
- Продолжать итерации пока **отношения** PageRank не перестанут изменяться
- Фиксированное число итераций

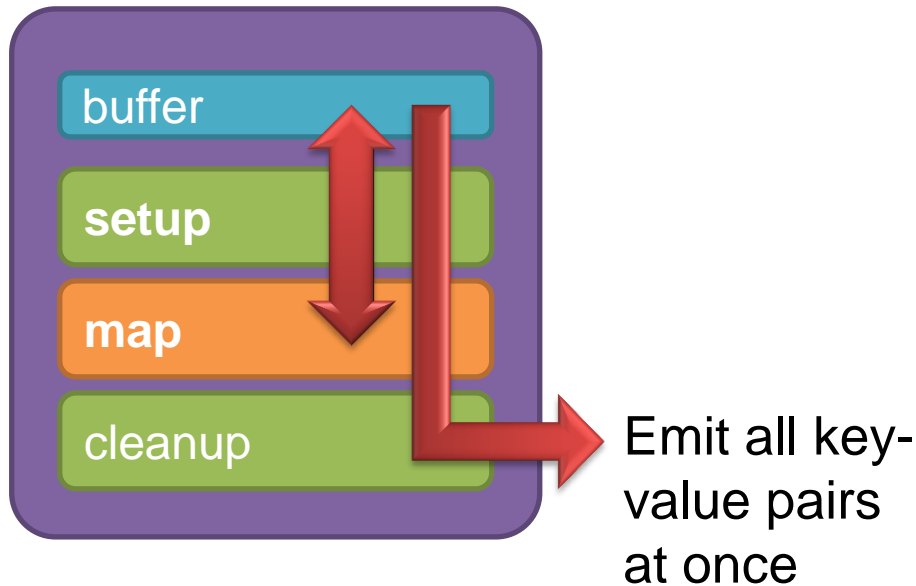
Проблемы MapReduce на графах

MapReduce на графах, проблемы

- Многословность Java
- Время запуска таска в Hadoop
- Медленные или зависшие таски
- Бесполезность фазы shuffle для графов
- Проверки на каждой итерации
- Итеративные алгоритмы на MapReduce неэффективны!

In-Mapper Combining

- Использование комбайнеров
 - Агрегирует данные на mapper
 - Но, промежуточные данные все равно обрабатываются
- In-mapper combining
 - Агрегируем сообщения в буфере
 - Но, требуется управление памятью

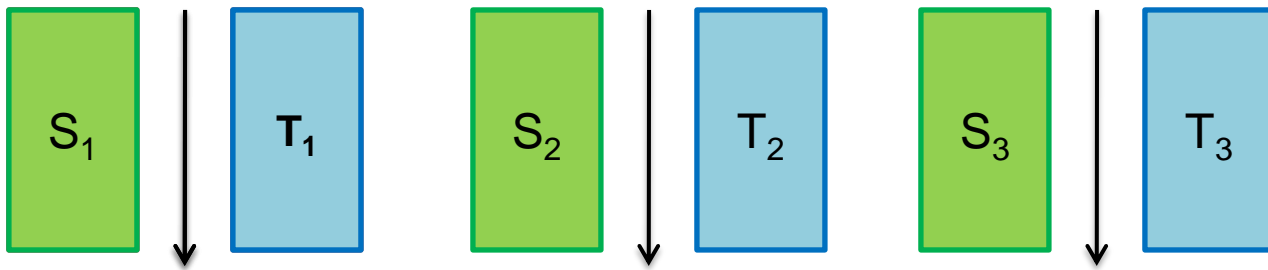


Улучшение партиционирования

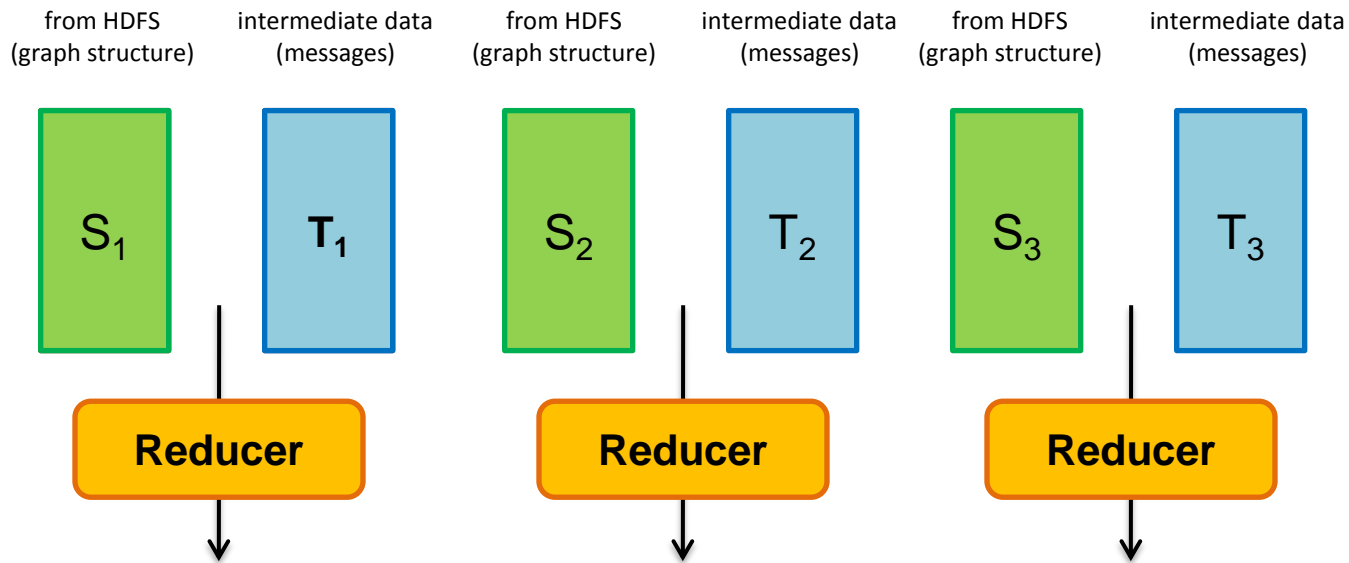
- По-умолчанию: hash partitioning
- Наблюдение: много графов имеют локальную структуру
 - Например, комьюнити в соц.сетях
 - Улучшение локальной агрегации
- Но, партиционирование довольно **сложно!**
 - Иногда простые эвристики помогают
 - Для веб-графа: использовать партиционирование на основе домена от URL

Schimmy Design Pattern

- Обычно два набора данных:
 - *Messages* (актуальные вычисления)
 - *Graph structure* (структура обрабатываемого графа)
- Schimmy: выполнять *shuffle* только для *messages*



Обе части (S и T) consistently partitioned and sorted by join key



Эксперимент

- Cluster setup:
 - 10 workers, each 2 cores (3.2 GHz Xeon), 4GB RAM, 367 GB disk
 - Hadoop 0.20.0 on RHEL 5.3
- Dataset:
 - Первый сегмент английского текста из коллекции ClueWeb09
 - 50.2m web pages (1.53 TB uncompressed, 247 GB compressed)
 - Extracted webgraph: 1.4 Млрд ссылок, 7.0 GB
 - Dataset сортирован в порядке краулинга
- Setup:
 - Измерялось время выполнения по каждой итерации (5 итераций)
 - 100 партиций

