

Масштабируемость в распределенных **in-memory** системах

Владимир Озеров

Яков Жданов

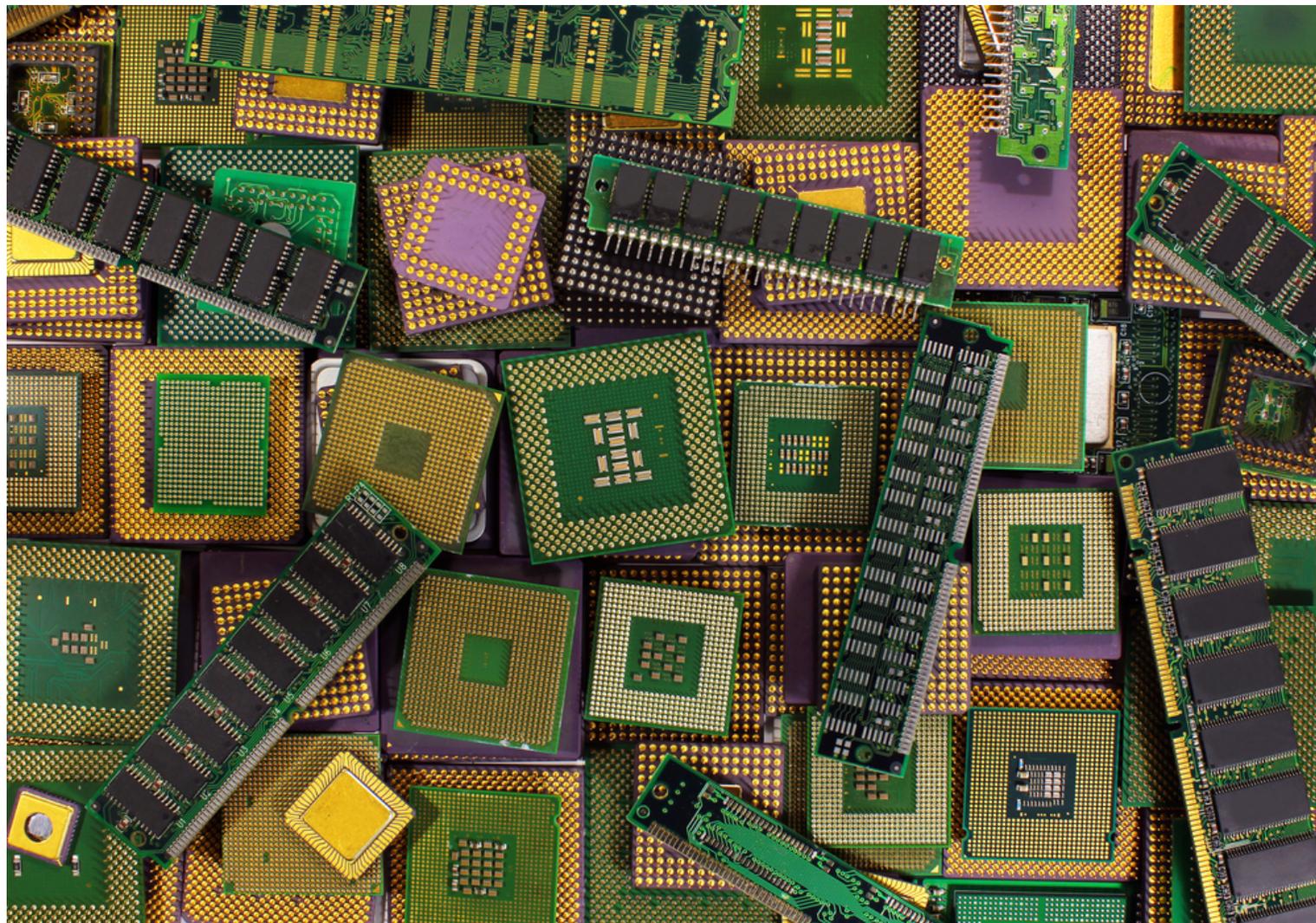
GridGain



КТО?

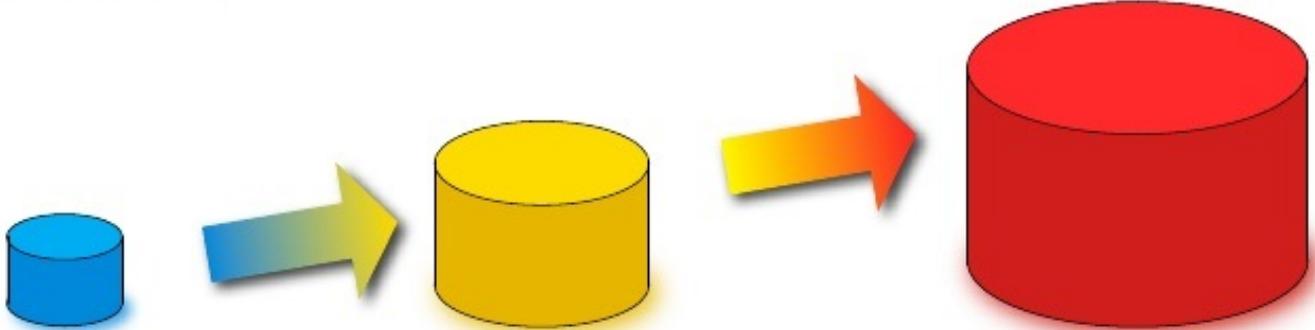


Почему in-memory?

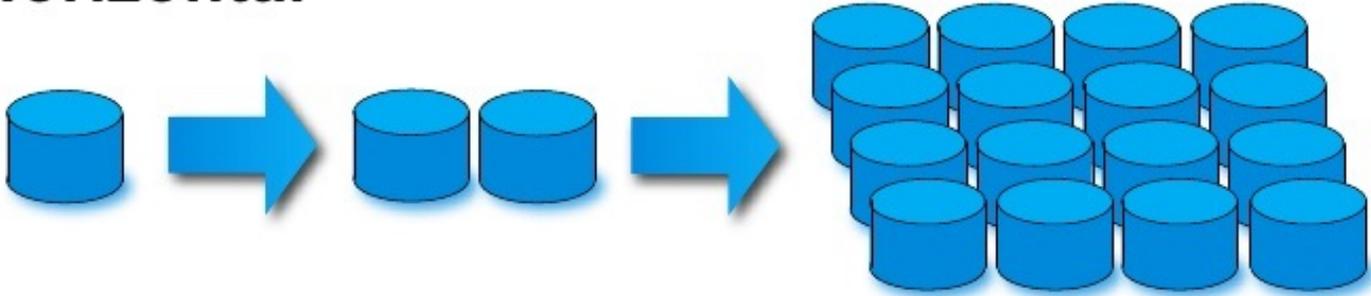


Масштабируем

Vertical



Horizontal



План

- Алгоритмы шардирования



План

- Алгоритмы шардирования
- Data co-location



План

- Алгоритмы шардирования
- Data co-location
- Синхронизация в кластере



План

- Алгоритмы шардирования
- Data co-location
- Синхронизация в кластере
- Локальная архитектура многопоточности



План

- Алгоритмы шардирования
- Data co-location
- Синхронизация в кластере
- Локальная архитектура многопоточности



Куда?

PUT(K, V)

?



Node 1



Node 2



Affinity

partition => node



Affinity

PUT(K, V)

?

0	2	4
6	8	10

Node 1

1	3	5
7	9	11

Node 2



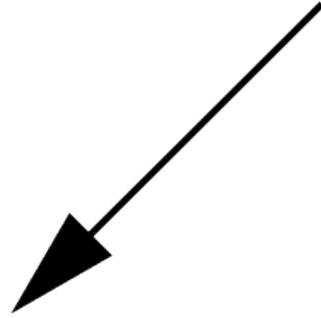
Affinity

key \Rightarrow partition \Rightarrow node



Affinity

PUT(K, V)



0	2	4
6	8	10

Node 1

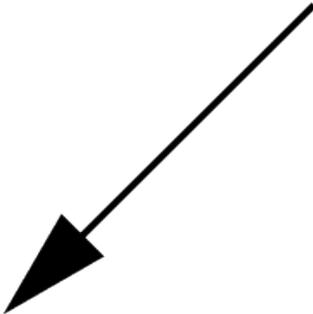
1	3	5
7	9	11

Node 2



Naïve affinity

PUT(K, V)



0	2	4
6	8	10

Node 1

50 TPS

1	3	5
7	9	11

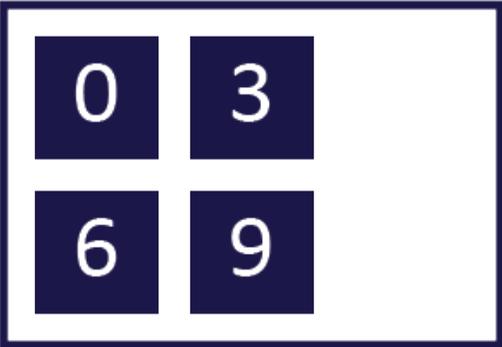
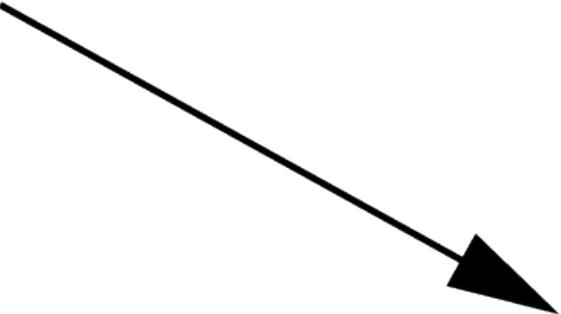
Node 2

50 TPS



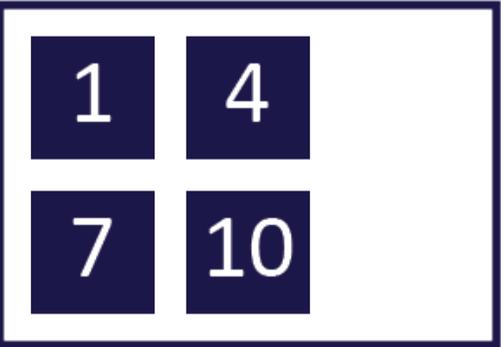
Naïve affinity

PUT(K, V)



Node 1

33 TPS



Node 2

33 TPS



Node 3

33 TPS



Naïve affinity

0	2	4
6	8	10

1	3	5
7	9	11

Node 1

Node 2



Naïve affinity

0	2	4
6	8	10

1	3	5
7	9	11

0	3
6	9

Node 1

1	4
7	10

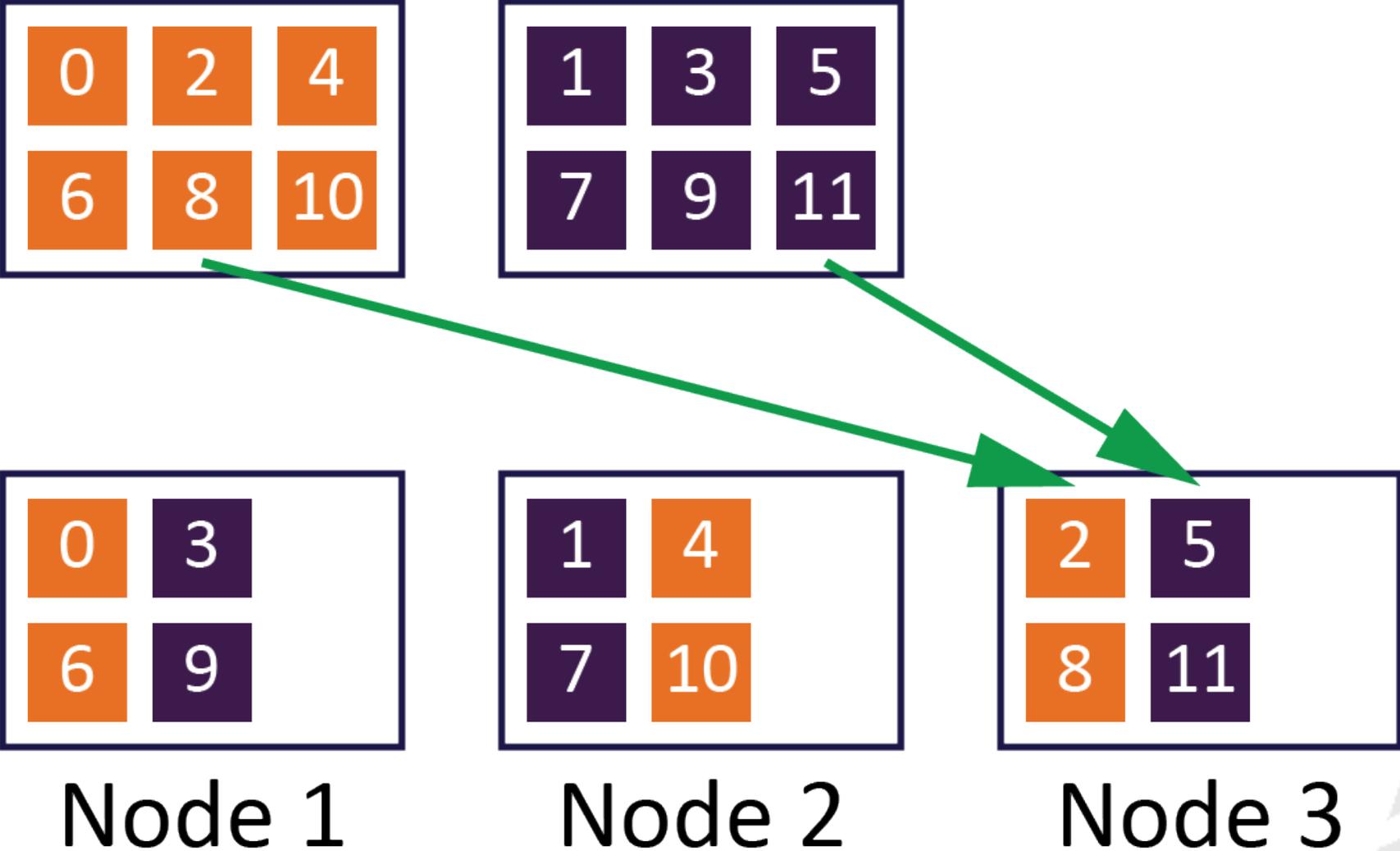
Node 2

2	5
8	11

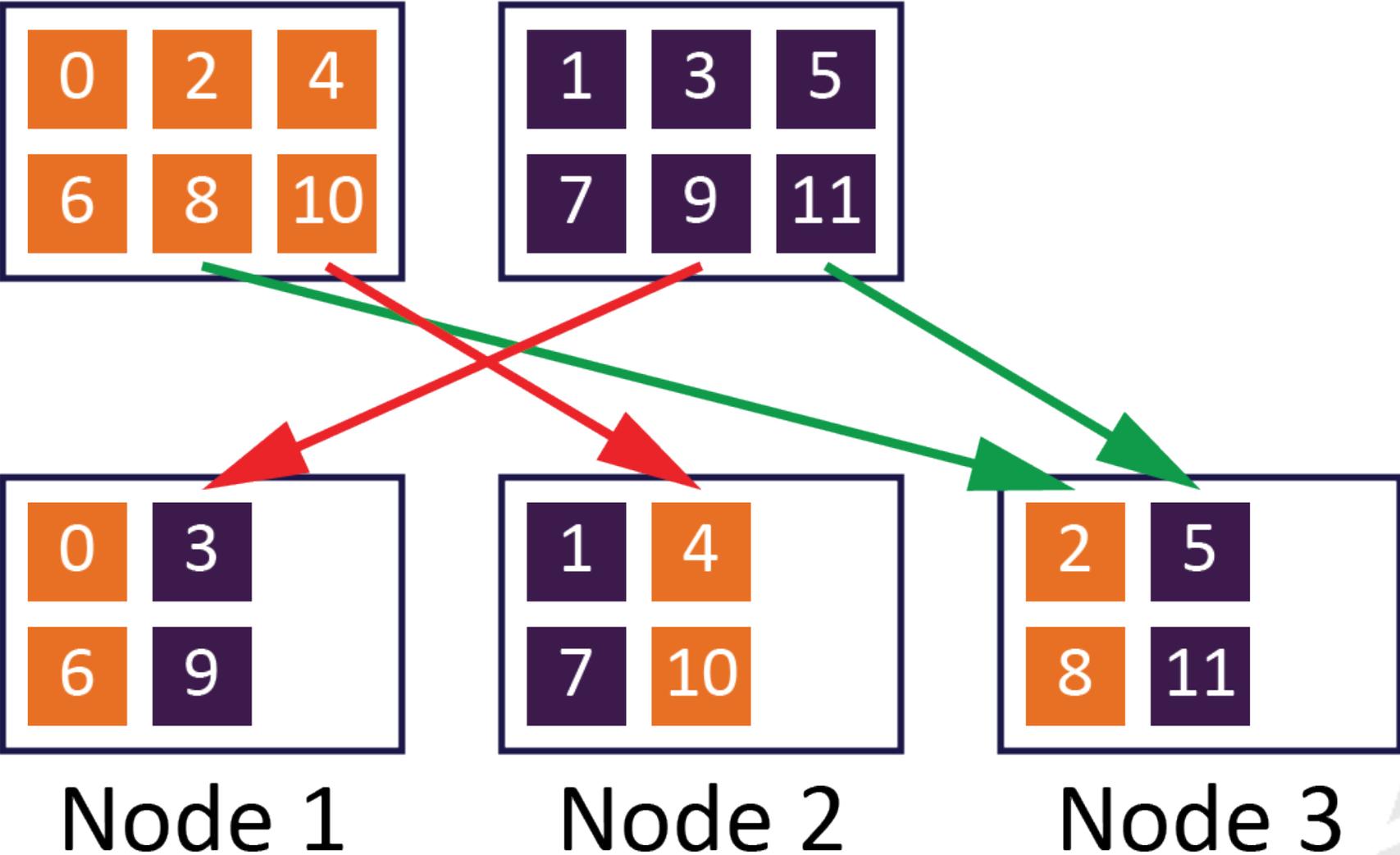
Node 3



Naïve affinity



Naïve affinity



Rendezvous affinity

- Consistent hashing [1]
- Rendezvous hashing (HRW) [2]

[1] https://en.wikipedia.org/wiki/Consistent_hashing

[2] https://en.wikipedia.org/wiki/Rendezvous_hashing

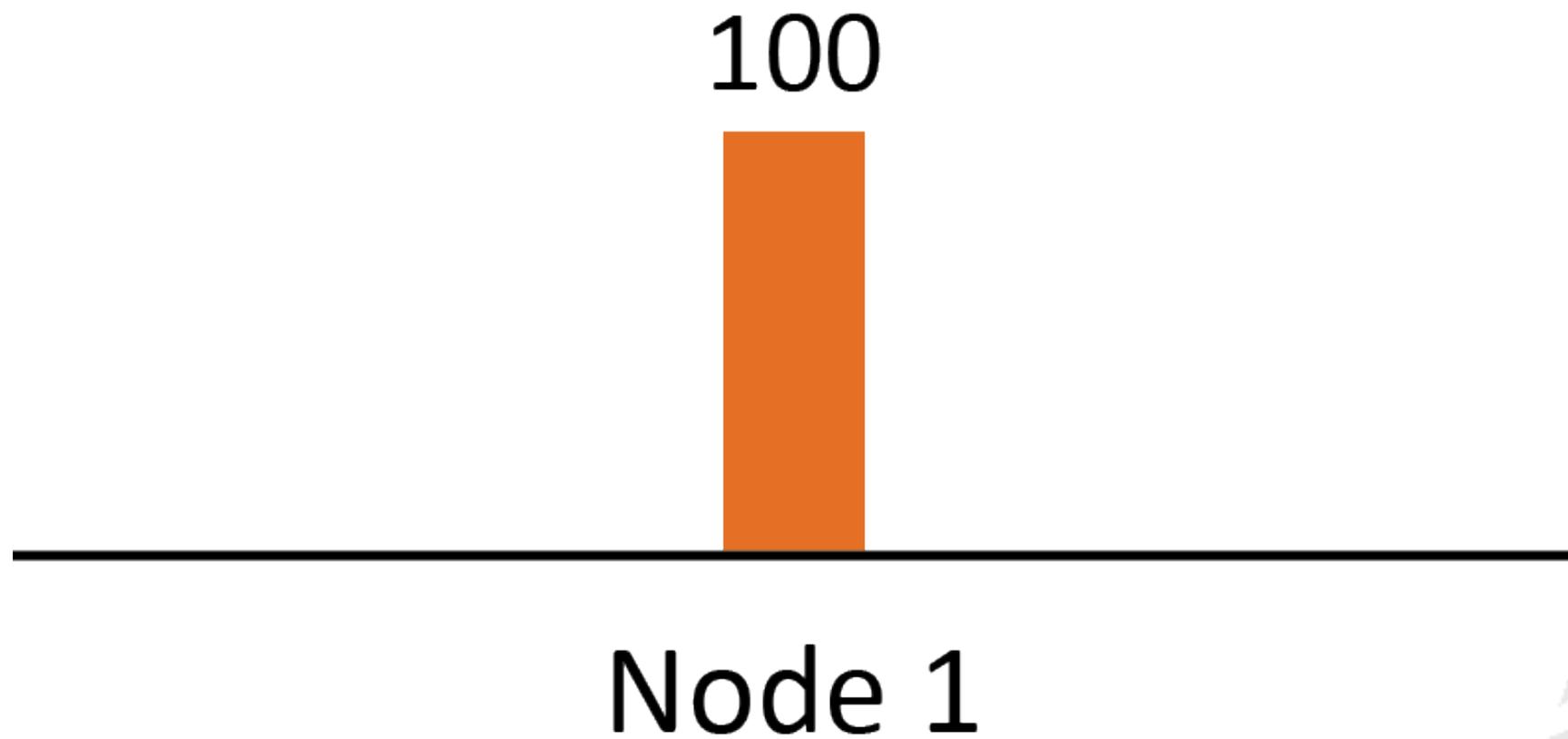


Rendezvous affinity

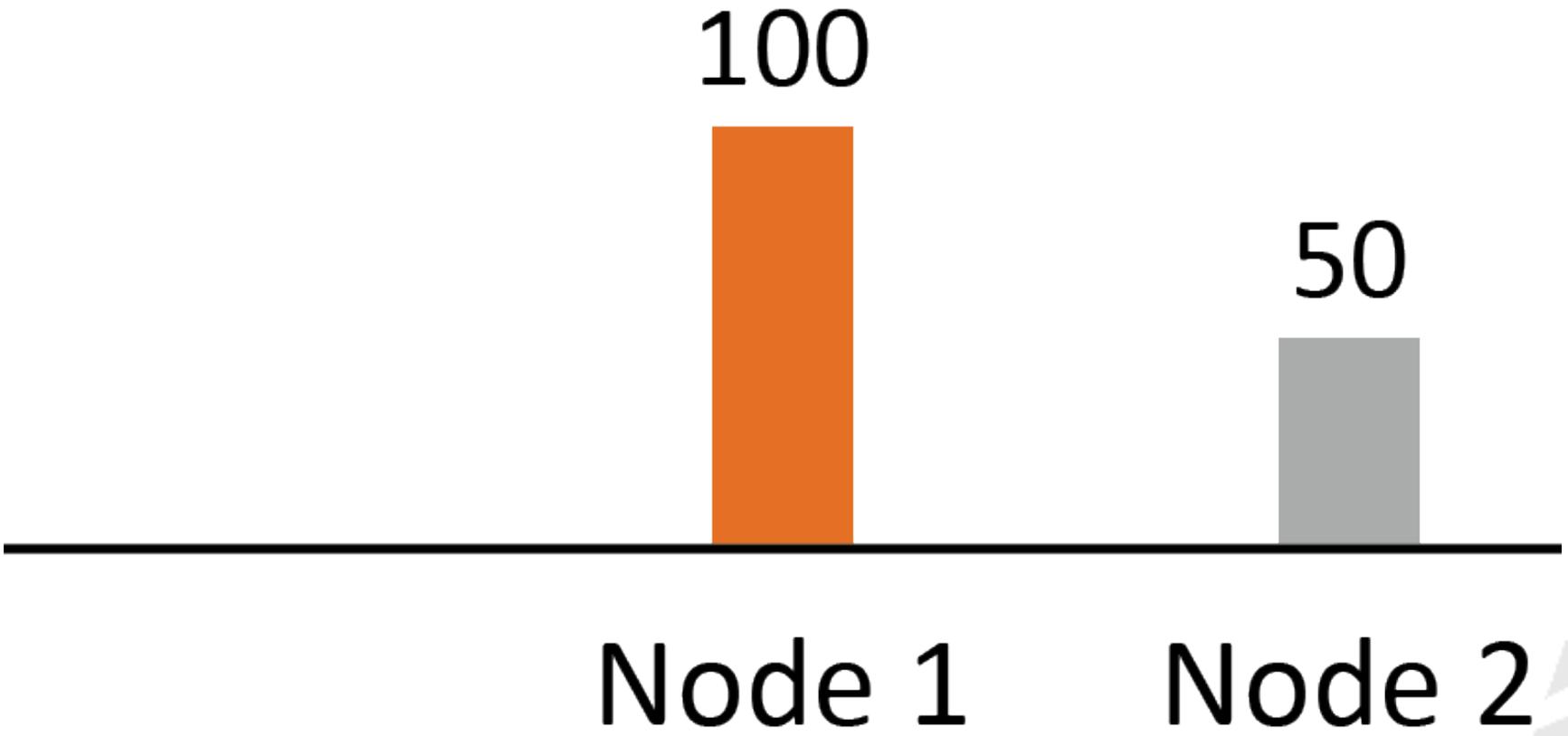
**WEIGHT (partition,
node)**



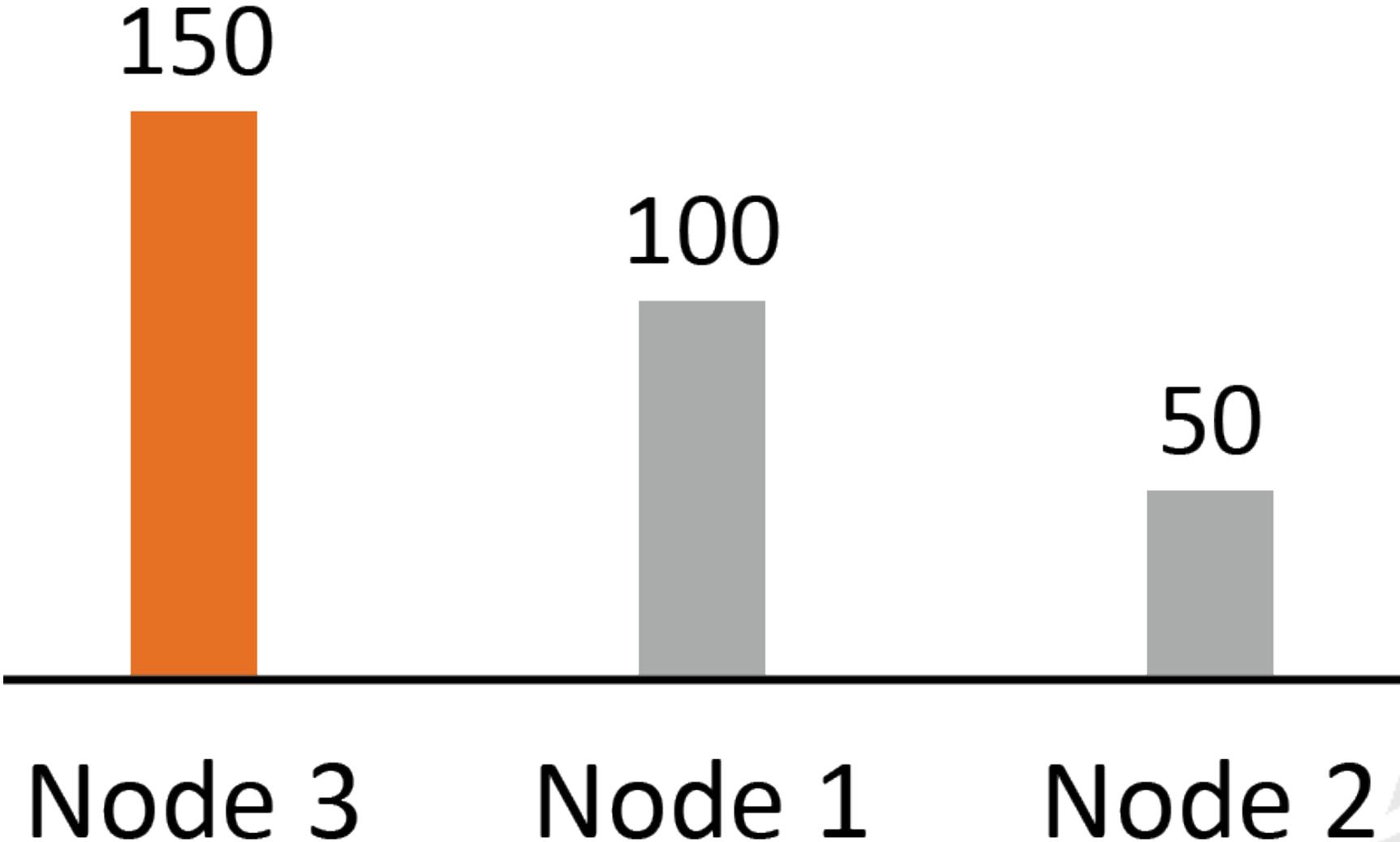
Rendezvous affinity



Rendezvous affinity



Rendezvous affinity



Rendezvous affinity

0	2	6
7	9	11

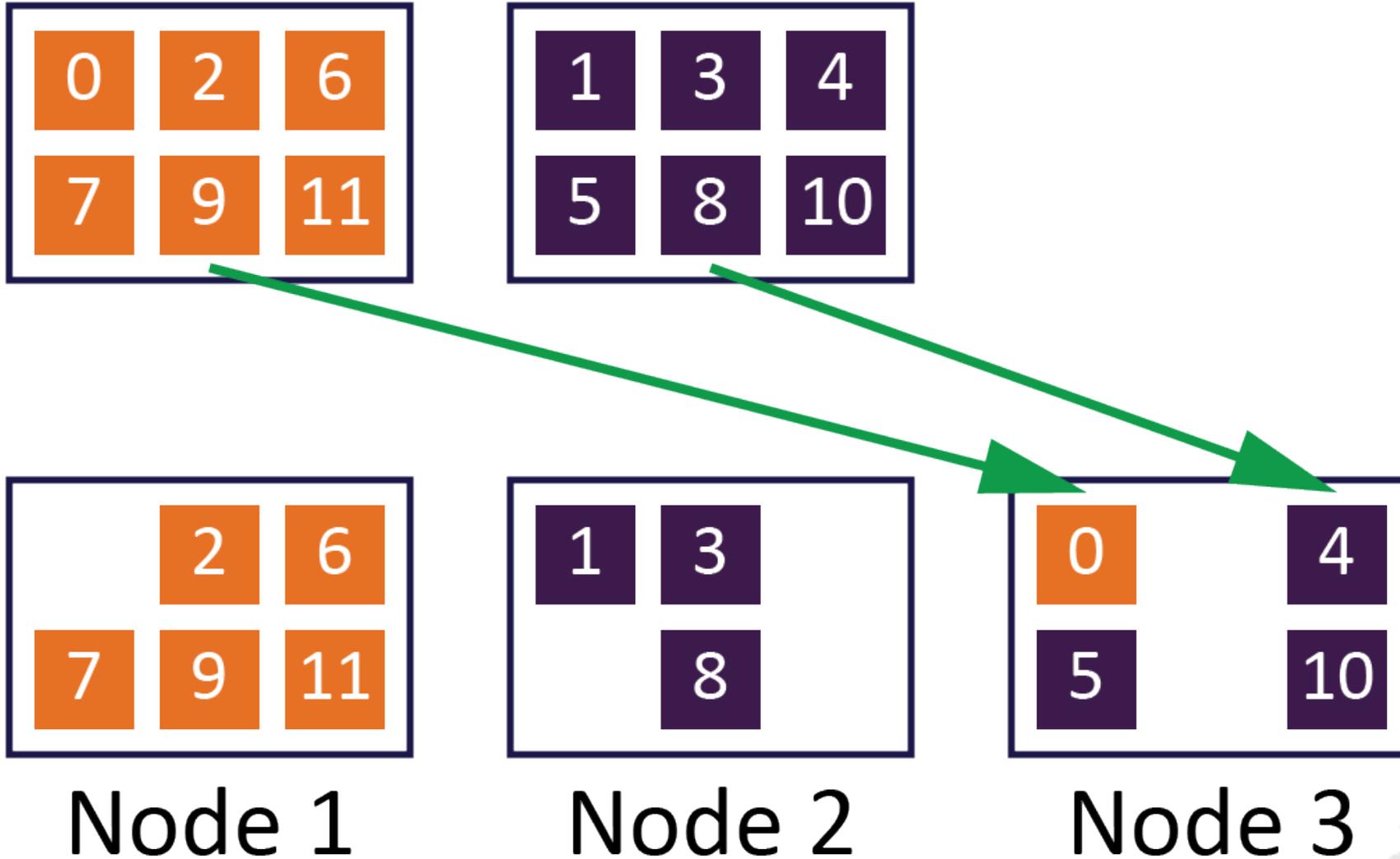
1	3	4
5	8	10

Node 1

Node 2

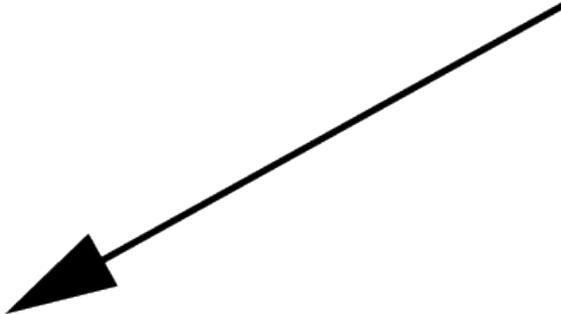


Rendezvous affinity



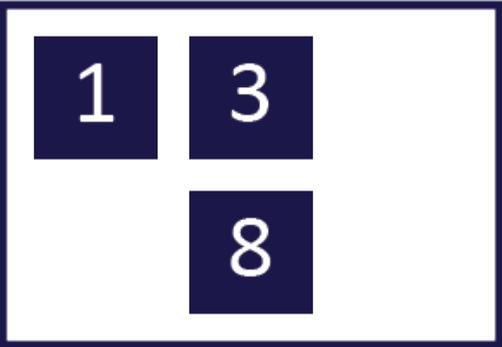
Rendezvous affinity

PUT(K, V)



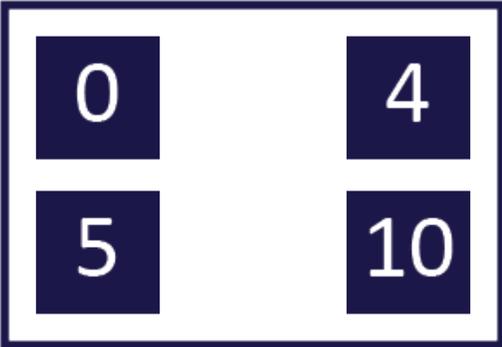
Node 1

42 TPS



Node 2

25 TPS



Node 3

33 TPS

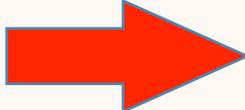
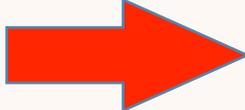


План

- Алгоритмы шардирования
- **Data co-location**
- Shared-nothing архитектура
- Синхронизация в кластере



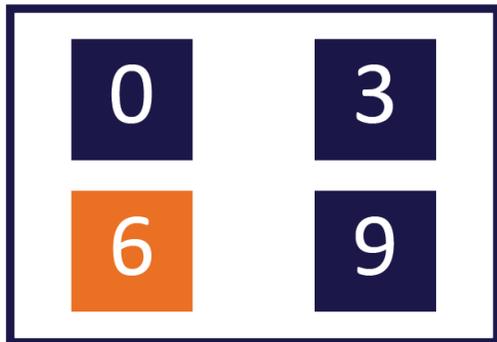
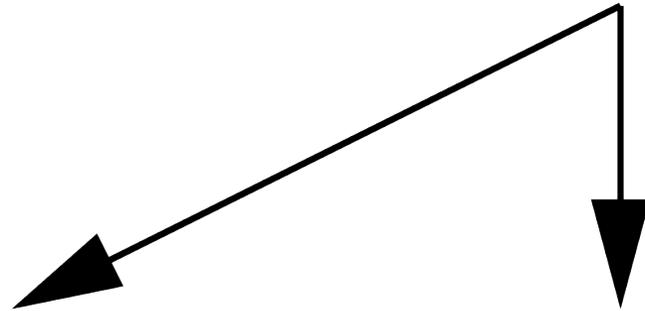
Без co-location: транзакции

```
1: class Customer {  
2:      long id;  
3:      City city;  
4: }
```

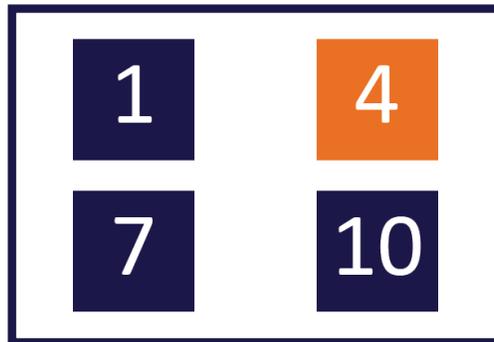
Без co-location: транзакции

UPDATE(id=1, city=MSK)

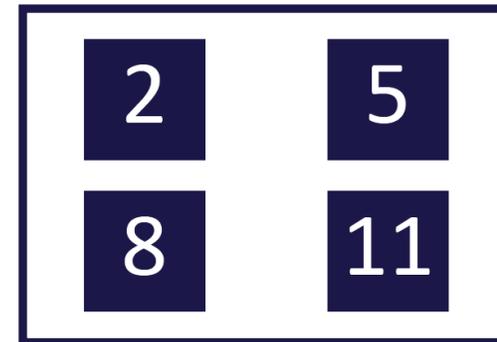
UPDATE(id=2, city=MSK)



Node 1

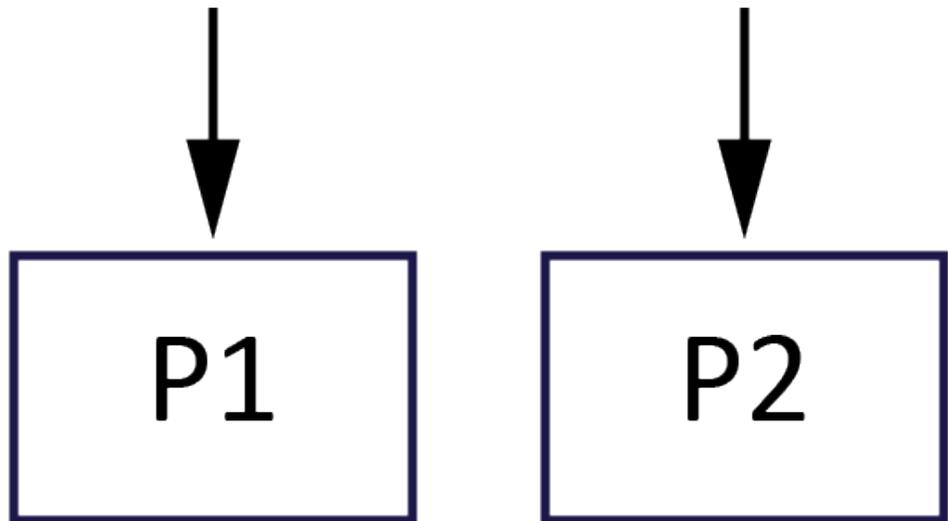


Node 2



Node 3

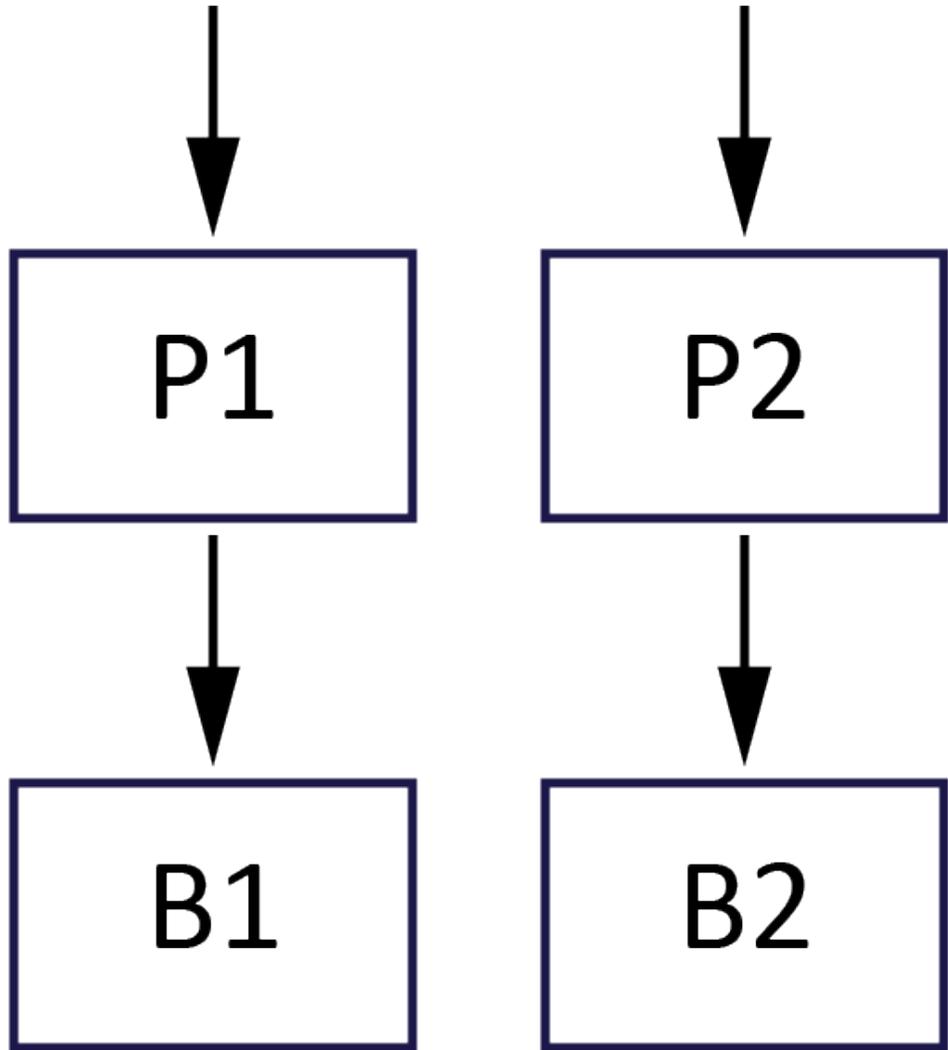
Без co-location: транзакции



2 (2 узла)



Без co-location: транзакции

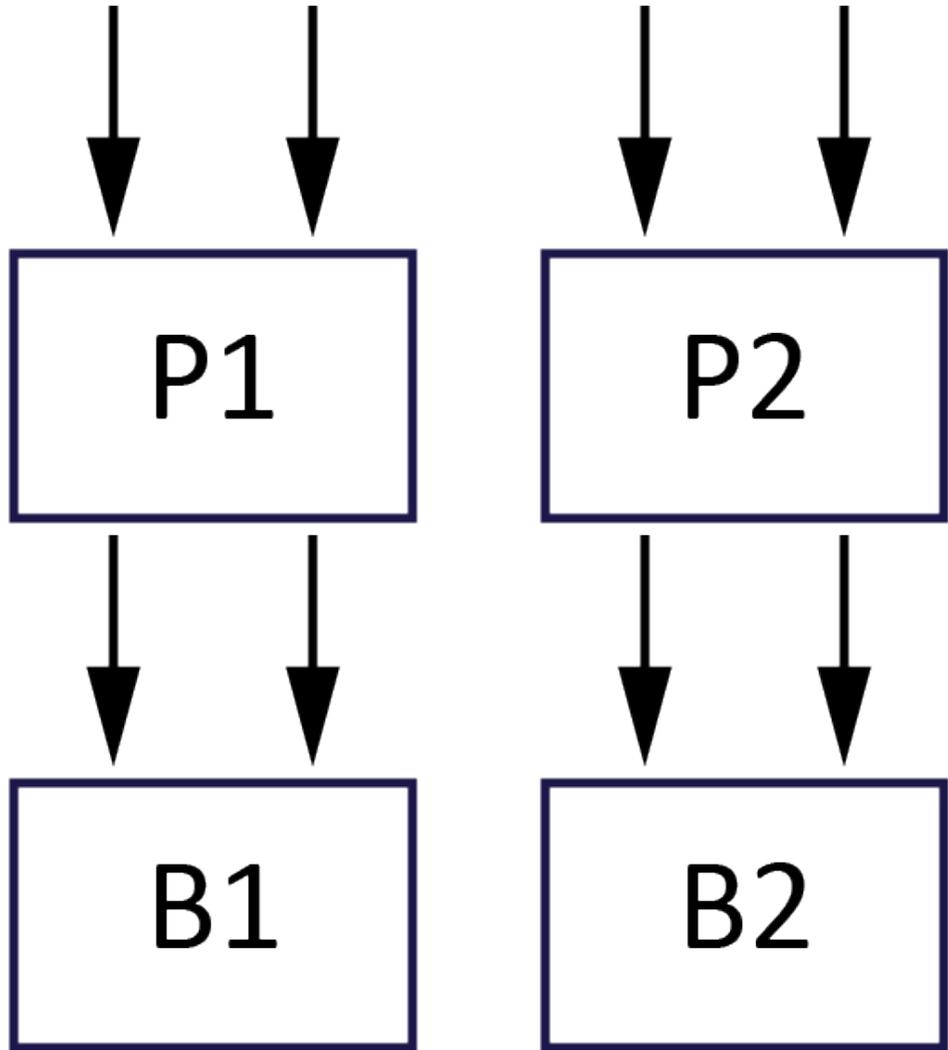


2 (2 узла)

2 (primary + backup)



Без co-location: транзакции



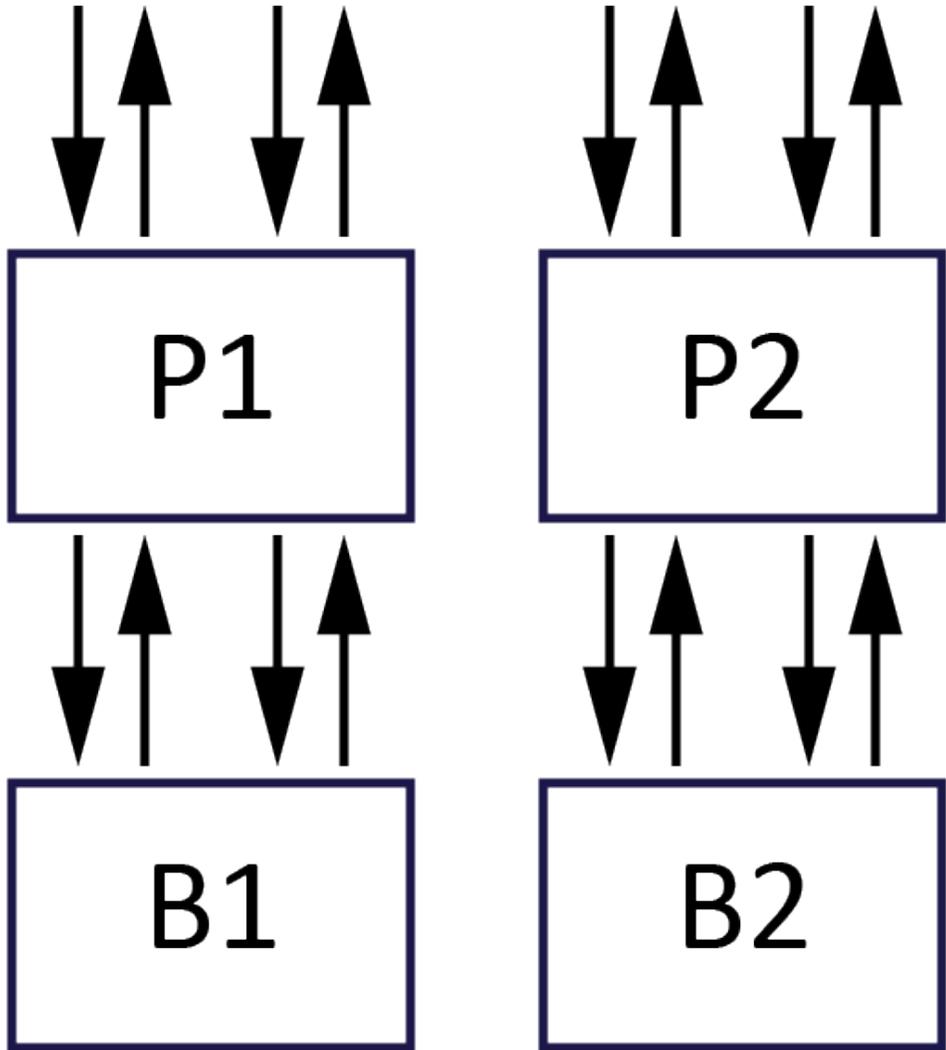
2 (2 узла)

2 (primary + backup)

2 (two-phase commit)



Без co-location: транзакции



2 (2 узла)

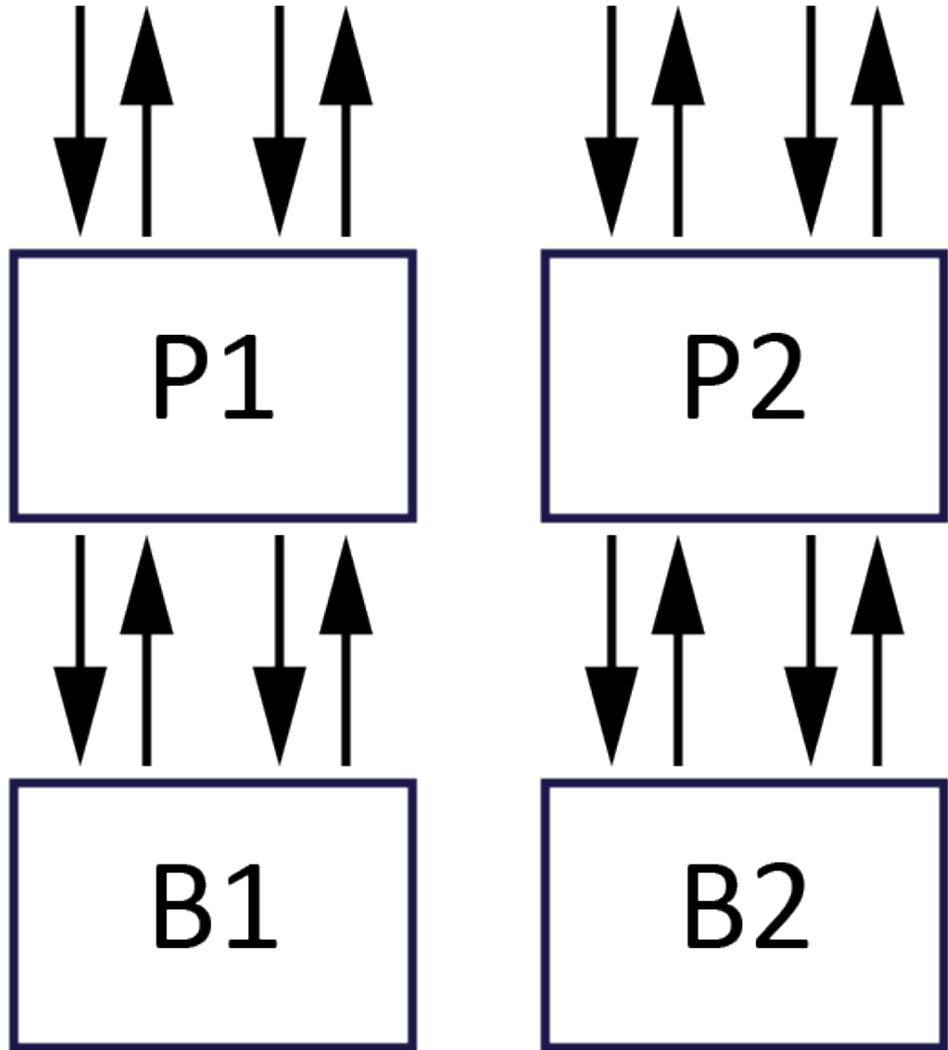
2 (primary + backup)

2 (two-phase commit)

2 (request-response)



Без co-location: транзакции



2 (2 узла)

2 (primary + backup)

2 (two-phase commit)

2 (request-response)

--

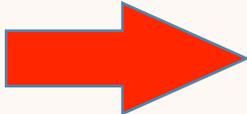
16 сообщений



Провод



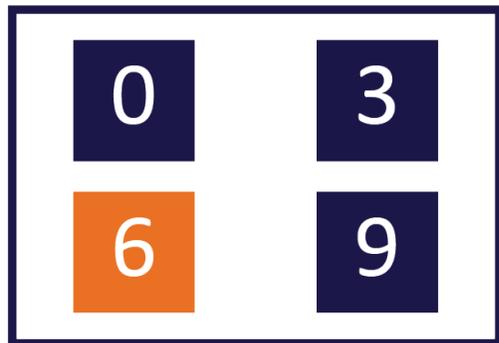
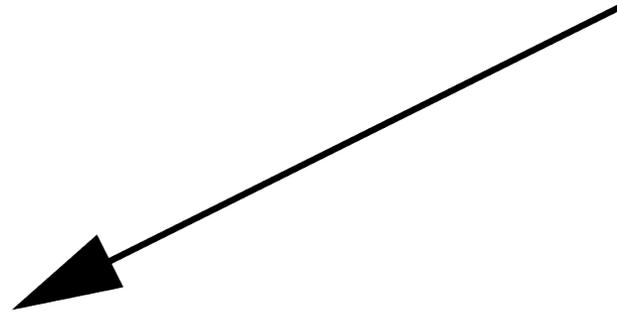
Co-location: транзакции

```
1: class Customer {  
2:     long id;  
3:  
4:     @AffinityKeyMapped  
5:      City city;  
6: }
```

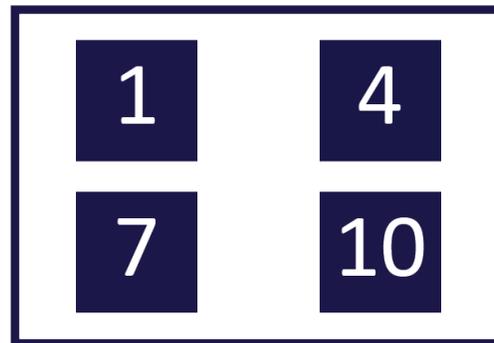
Co-location: транзакции

UPDATE(id=1, city=MSK)

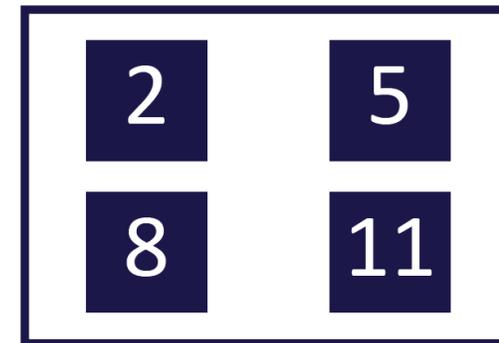
UPDATE(id=2, city=MSK)



Node 1

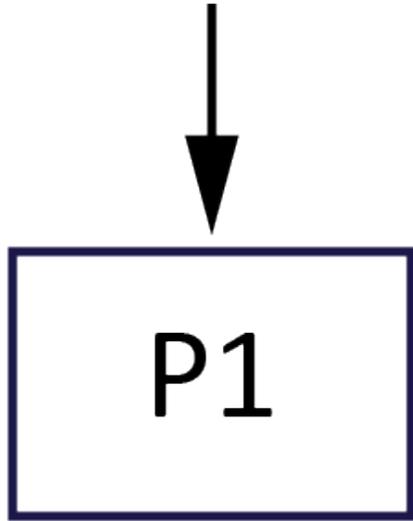


Node 2



Node 3

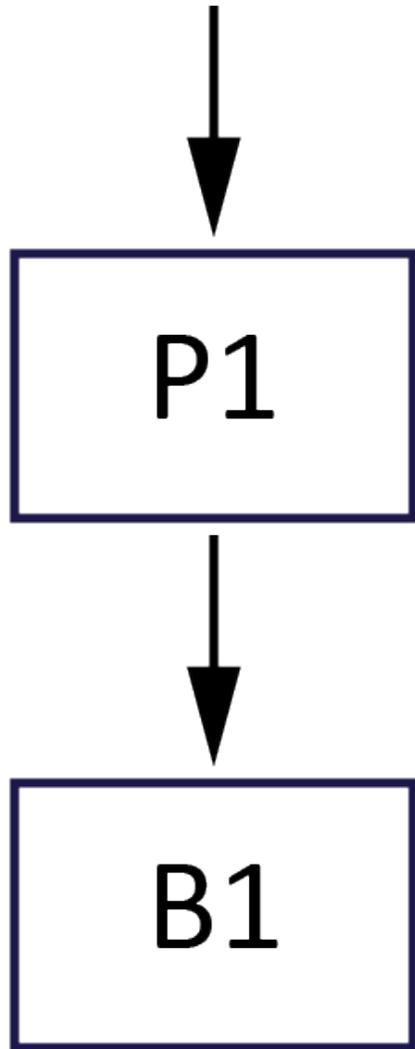
Co-location: транзакции



1 (1 узел)



Co-location: транзакции

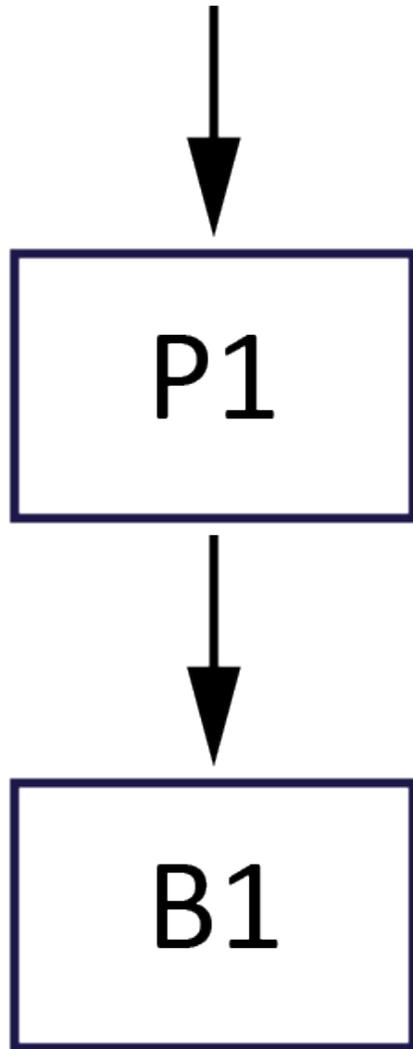


1 (1 узел)

2 (primary + backup)



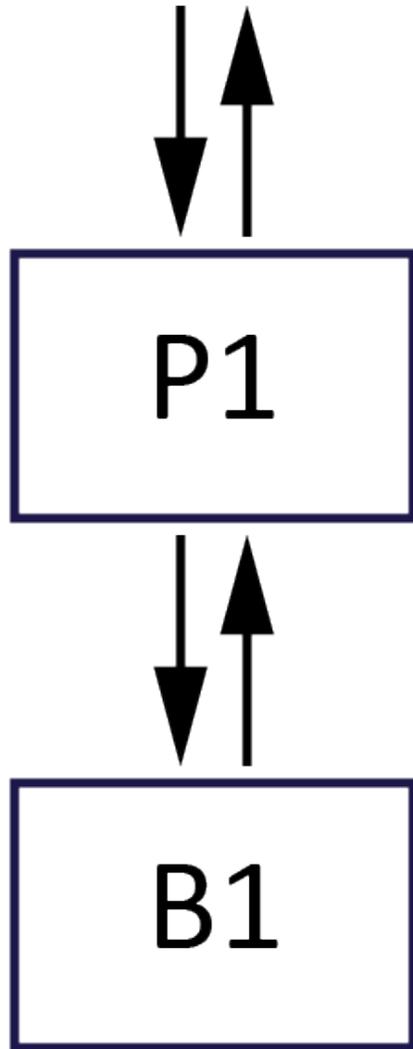
Co-location: транзакции



- 1** (1 узел)
- 2** (primary + backup)
- 1** (one-phase commit)



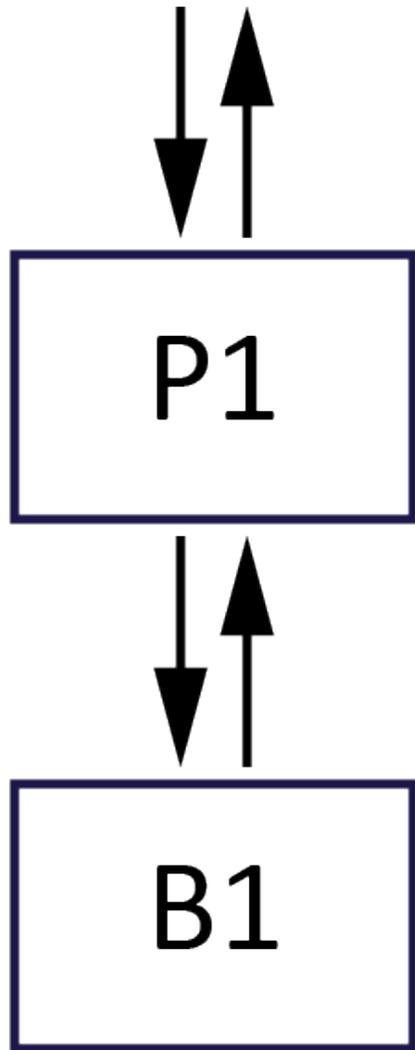
Co-location: транзакции



- 1** (1 узел)
- 2** (primary + backup)
- 1** (one-phase commit)
- 2** (request-response)



Co-location: транзакции



- 1** (1 узел)
- 2** (primary + backup)
- 1** (one-phase commit)
- 2** (request-response)

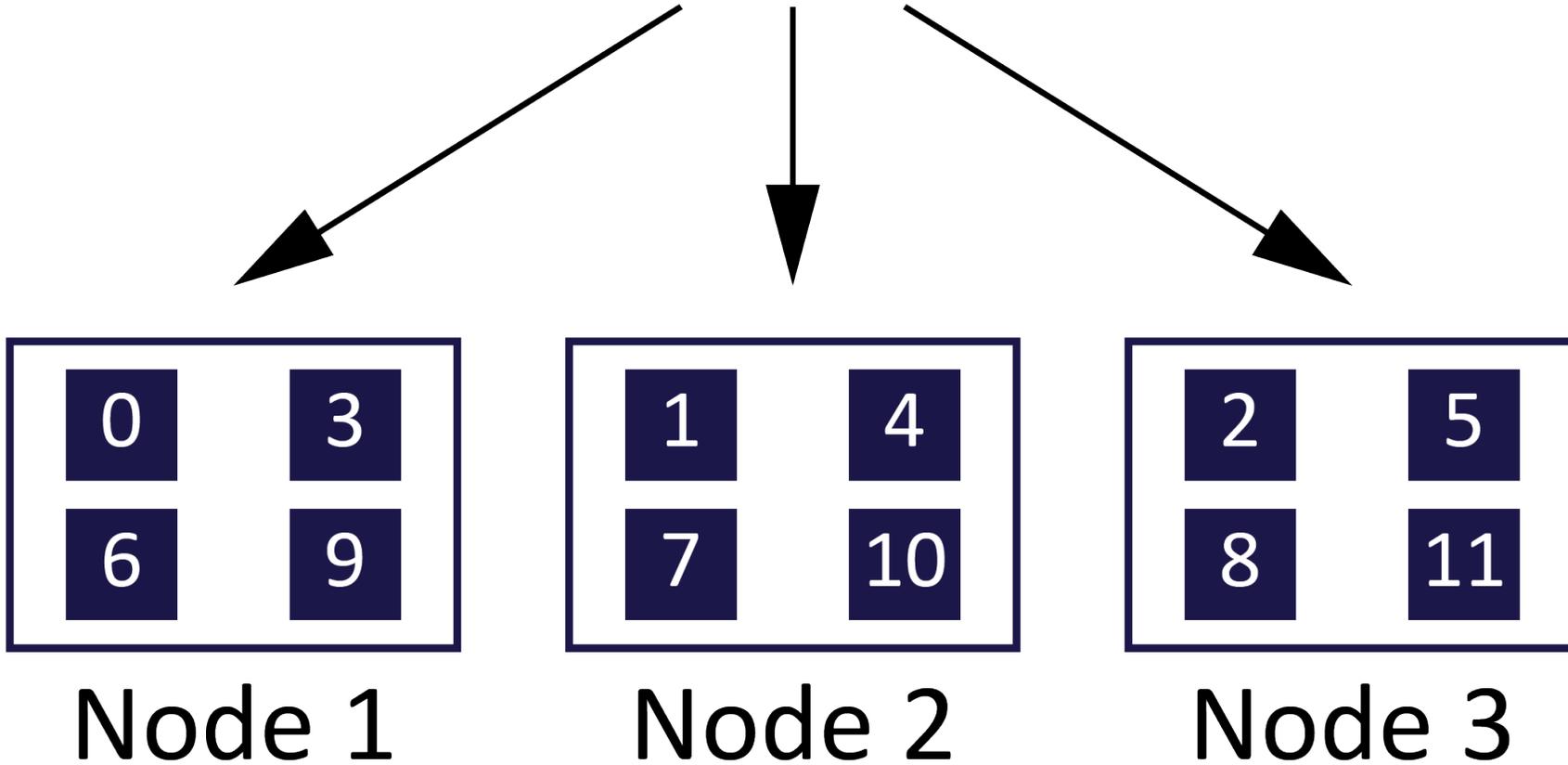
--

4 сообщения



Запустим SQL

SELECT ... WHERE city = MSK



Без co-location: full scan



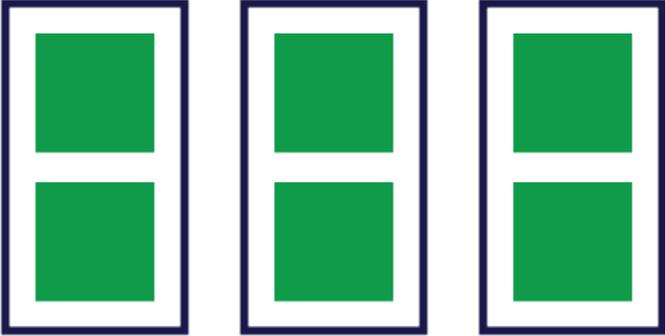
1 node



Без co-location: full scan



1 node



3 nodes



Без co-location: full scan

- **1/3x** latency!



Без co-location: full scan

- **1/3x** latency!
- **3x** capacity!



Индексы

А с индексами же еще круче
будет?



Индексы



1 node



Индексы



1 node



N nodes

Посчитаем сложность

$$\log 1_000_000 \approx 20$$



Посчитаем сложность

$$\log 1_000_000 \approx 20$$

vs

$$\log 333_333 \approx 18$$

$$\log 333_333 \approx 18$$

$$\log 333_333 \approx 18$$

Посчитаем сложность

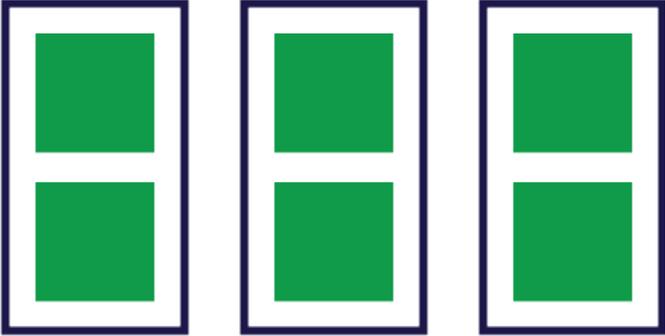
$$\log N < M * \log N / M$$



Без co-location: full scan



1 node



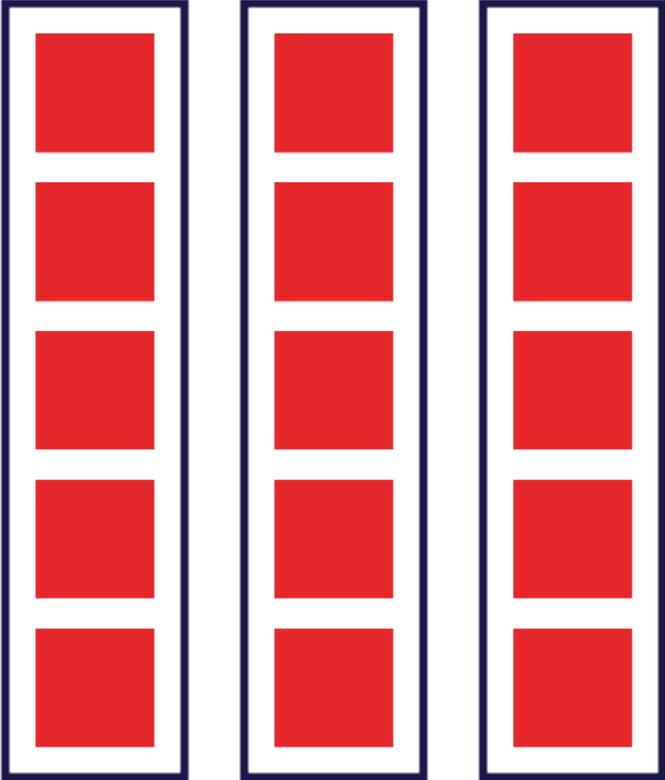
3 nodes



Без co-location: index seek



1 node



3 nodes



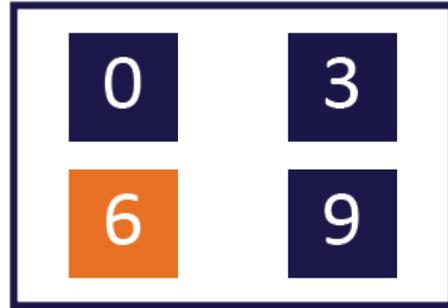
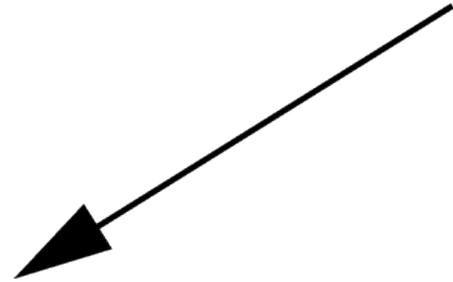
Без co-location: index seek

- **такой же** latency?
- **такой же** capacity?



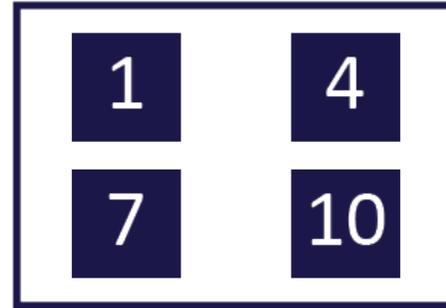
Co-location: index seek

SELECT ... WHERE city = **MSK**



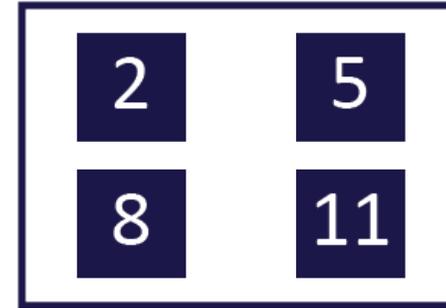
Node 1

100 TPS



Node 2

0 TPS



Node 3

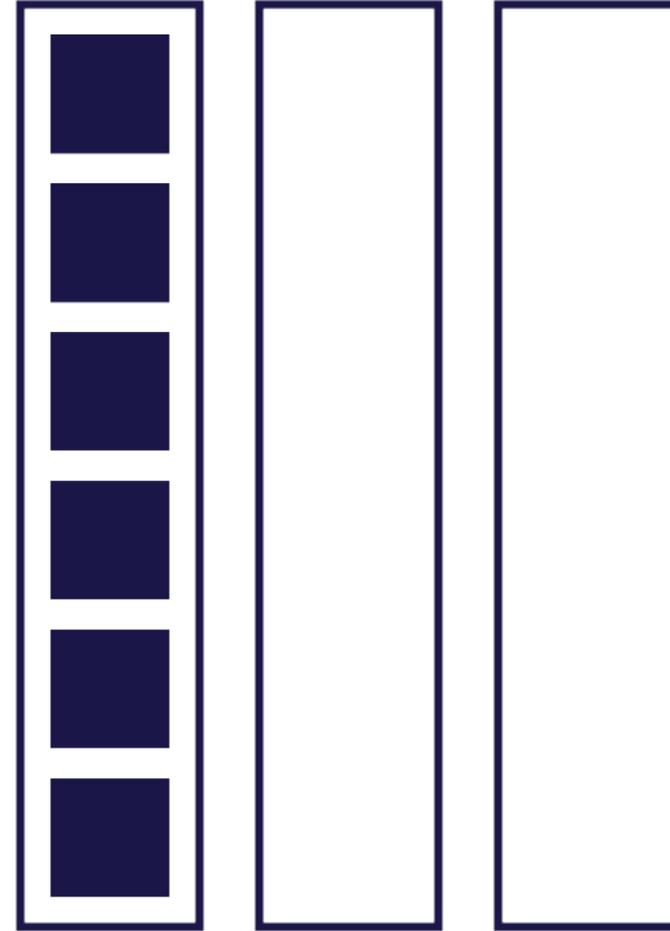
0 TPS



Co-location: index seek



1 node



3 nodes



Co-location: index seek

- **такой же** latency?
- **но 3x** capacity!



Дисбаланс



MSK

Node 1

100 TPS



SPB

Node 2

50 TPS



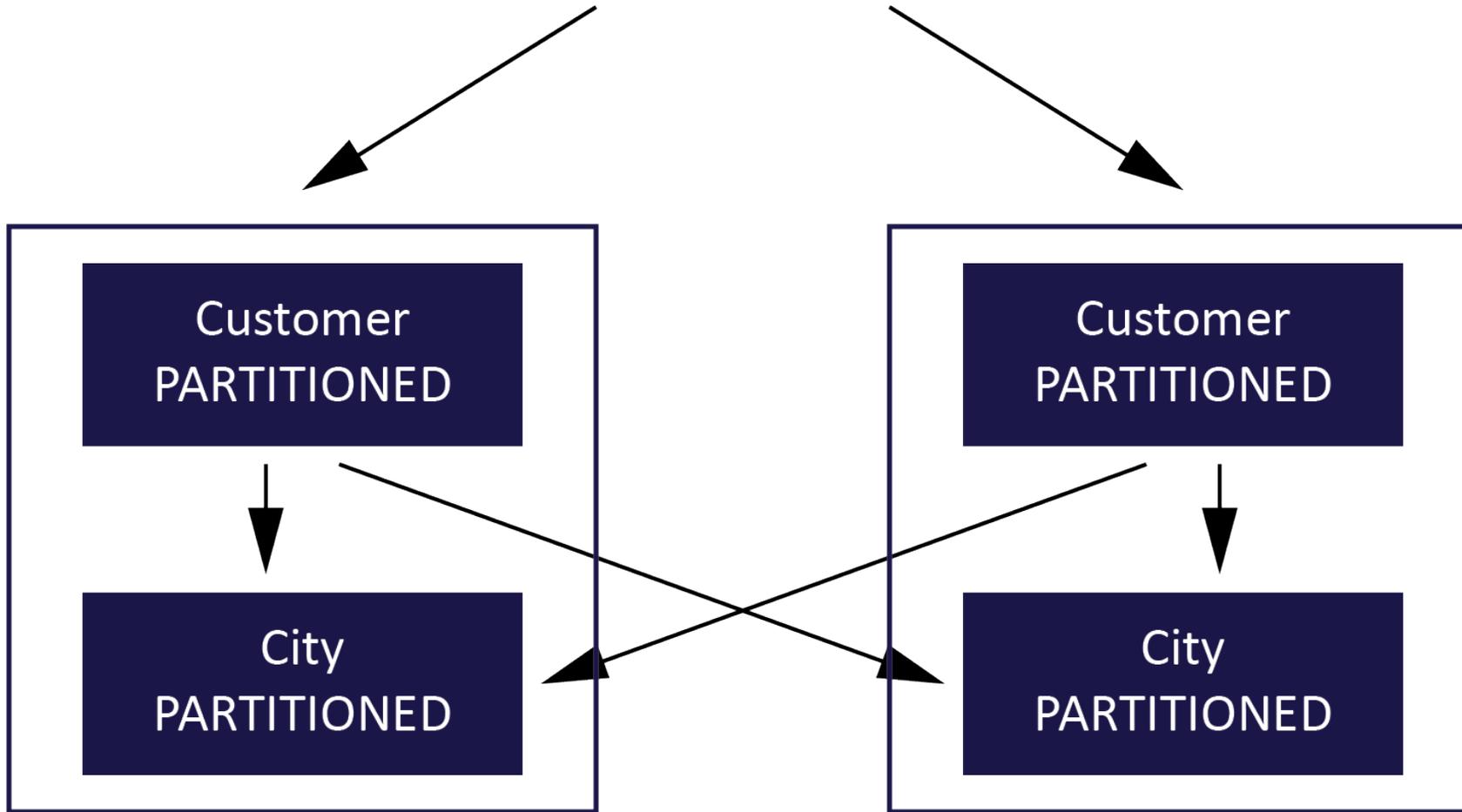
NSK

Node 3

10 TPS

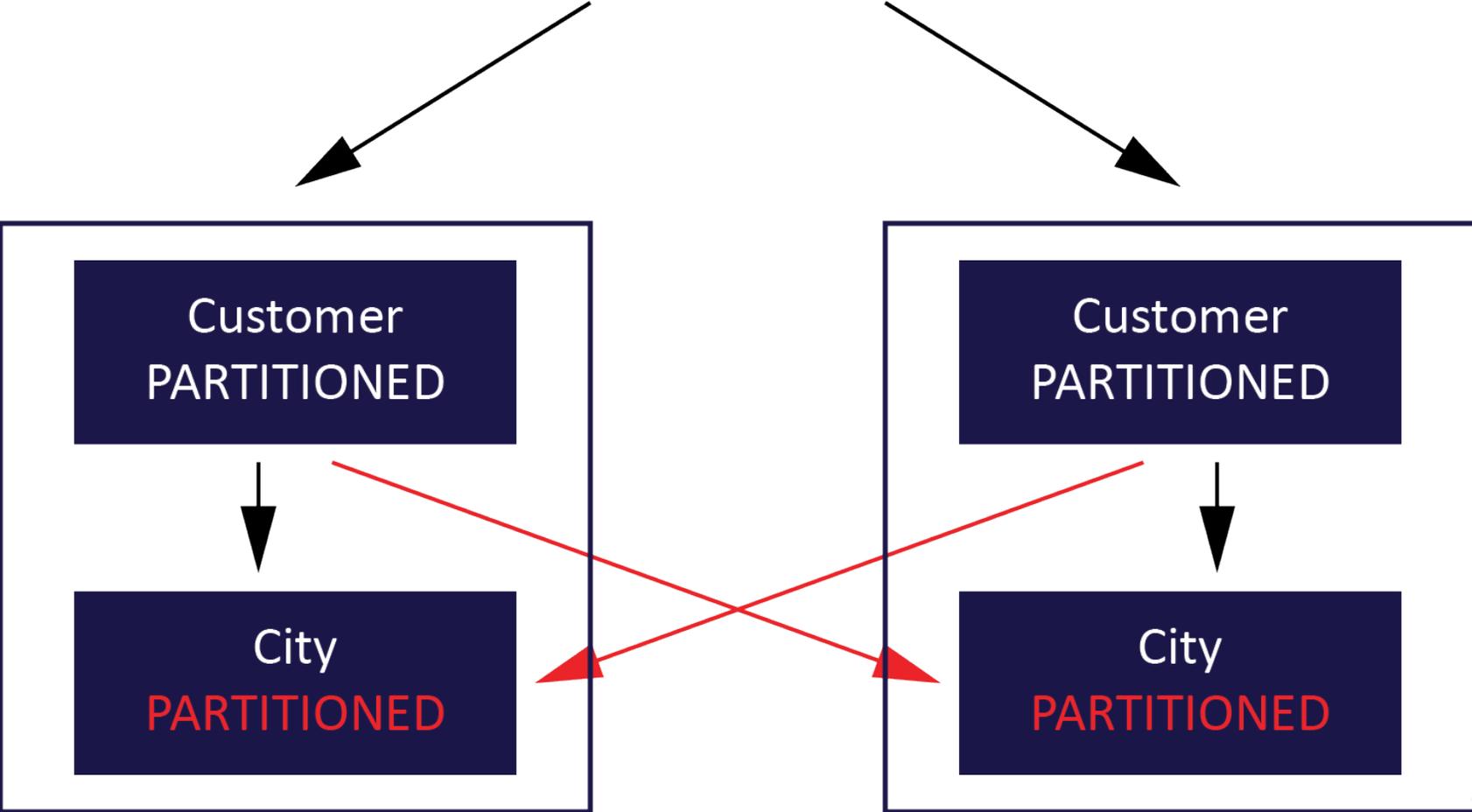
PARTITIONED cache

FROM Customer JOIN City



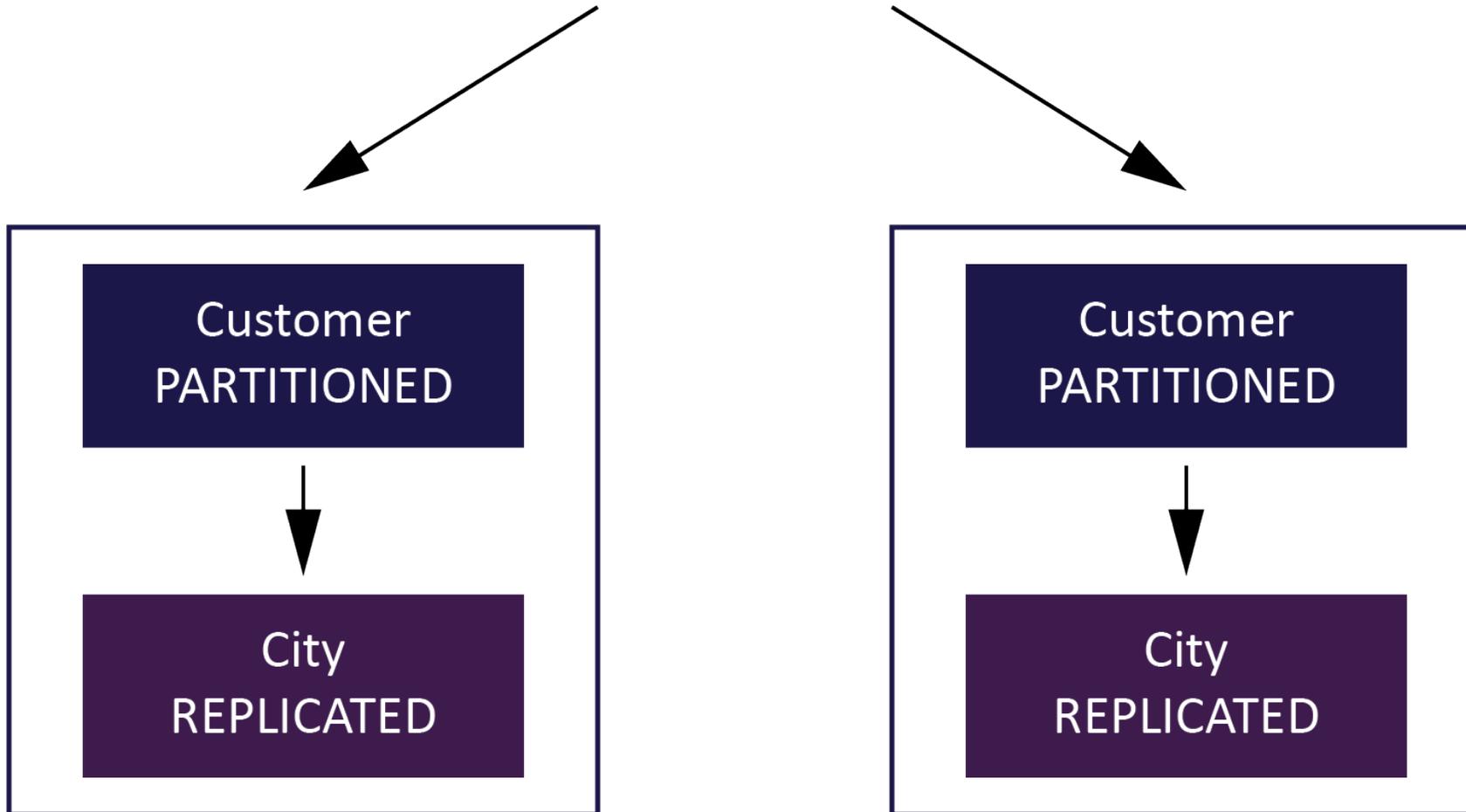
PARTITIONED cache

FROM Customer JOIN City



REPLICATED cache

FROM Customer JOIN City



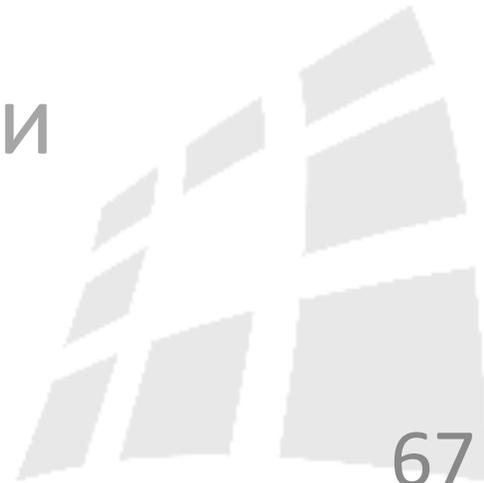
Адаптация модели

Local-distributed **impedance mismatch?**



План

- Алгоритмы шардирования
- Data co-location
- Синхронизация в кластере
- Локальная архитектура многопоточности

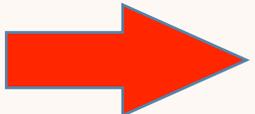


Счетчик: локально

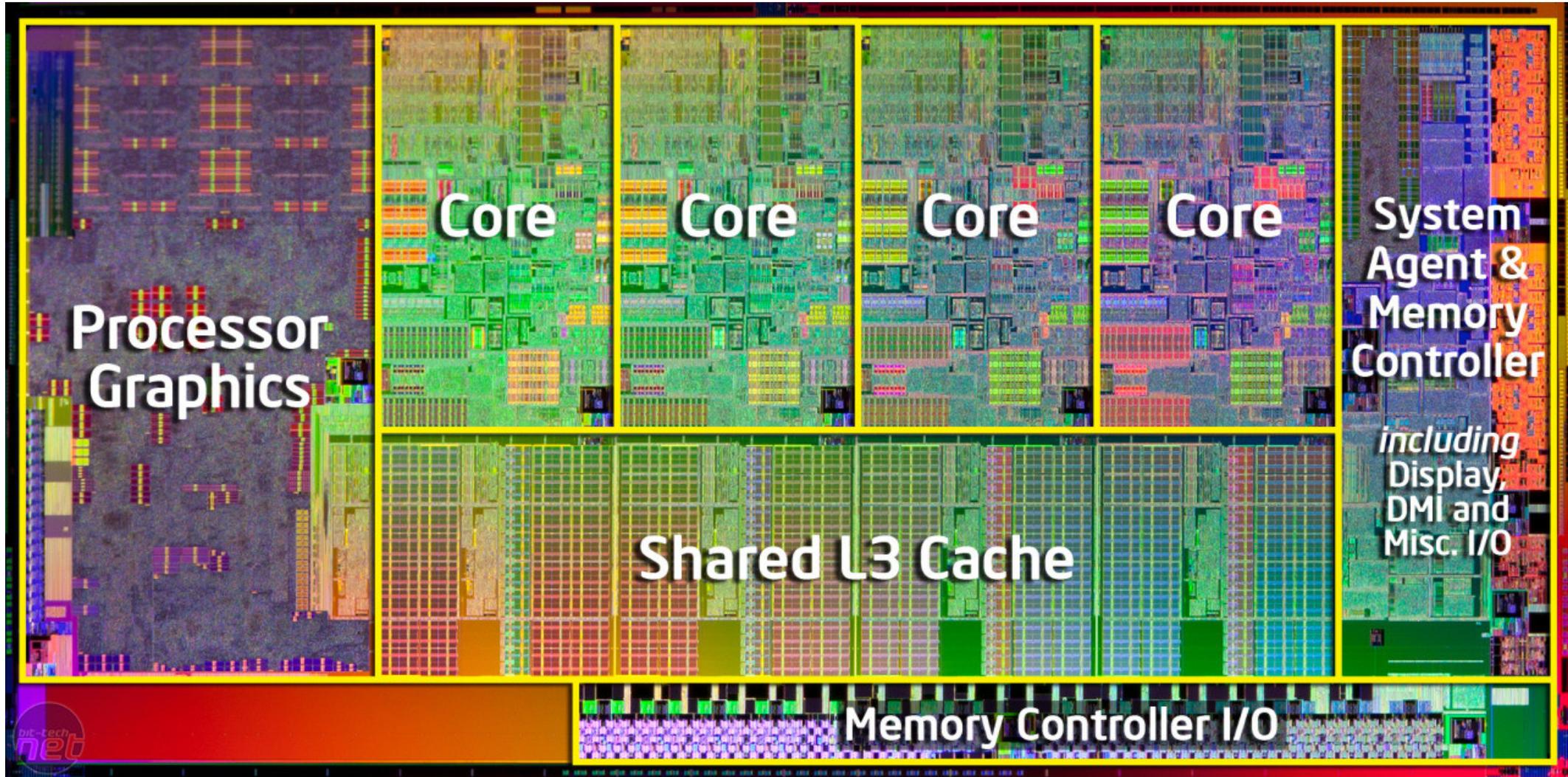
```
1: AtomicLong ctr;  
2:  
3: long getNext () {  
4:     return ctr.incrementAndGet ();  
5: }
```

Счетчик: локально, велосипедно

```
1: AtomicLong ctr;  
2: ThreadLocal<Long> localCtr;  
3:  
4: long getNext() {  
5:     long res = localCtr.get();  
6:  
7:     if (res % 1000 == 0)  
8:         res = ctr.getAndAdd(1000);  
9:  
10:    localCtr.set(++res);  
11:  
12:    return res;  
13: }
```



Процессор



Провод

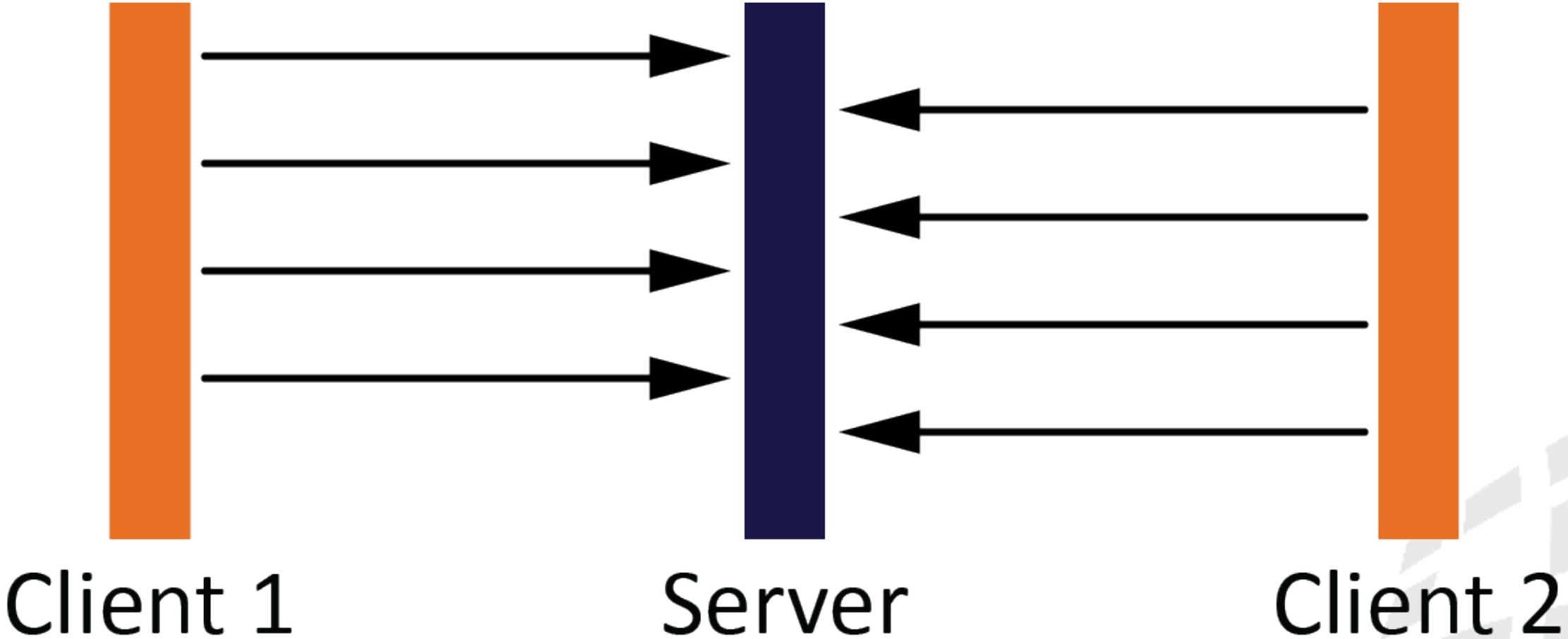


Счетчик: в кластере

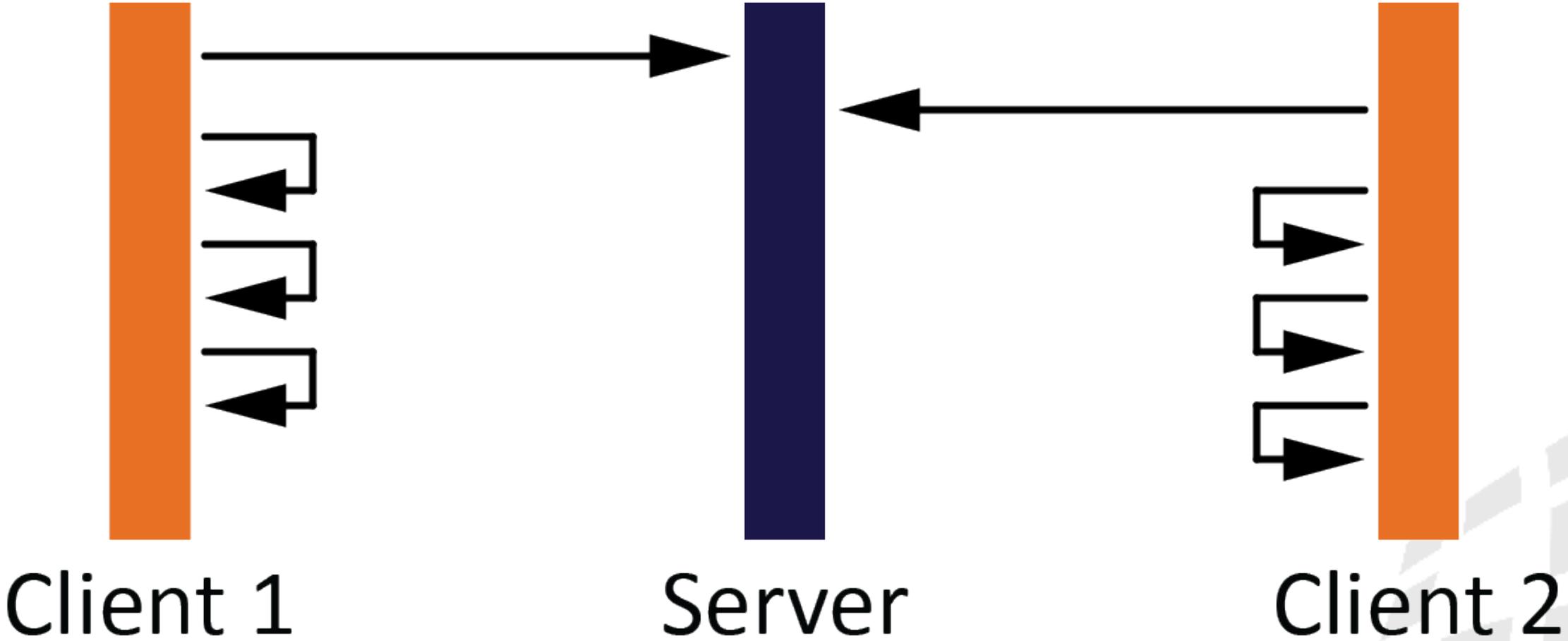
- Локально: **миллионы**
- Распределенно: **тысячи**



IgniteAtomicLong



IgniteAtomicSequence



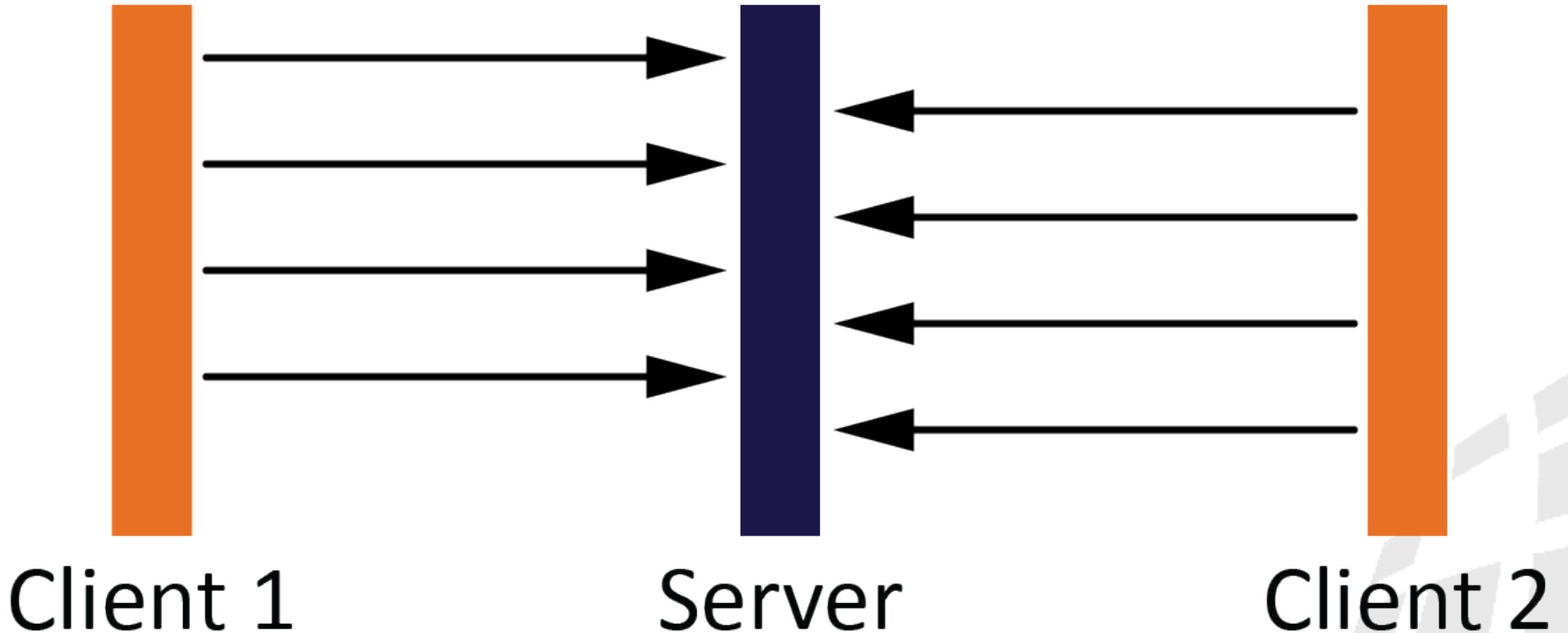
Разумные требования: счетчик

Уникальный, монотонно возрастающий, 8 байт!



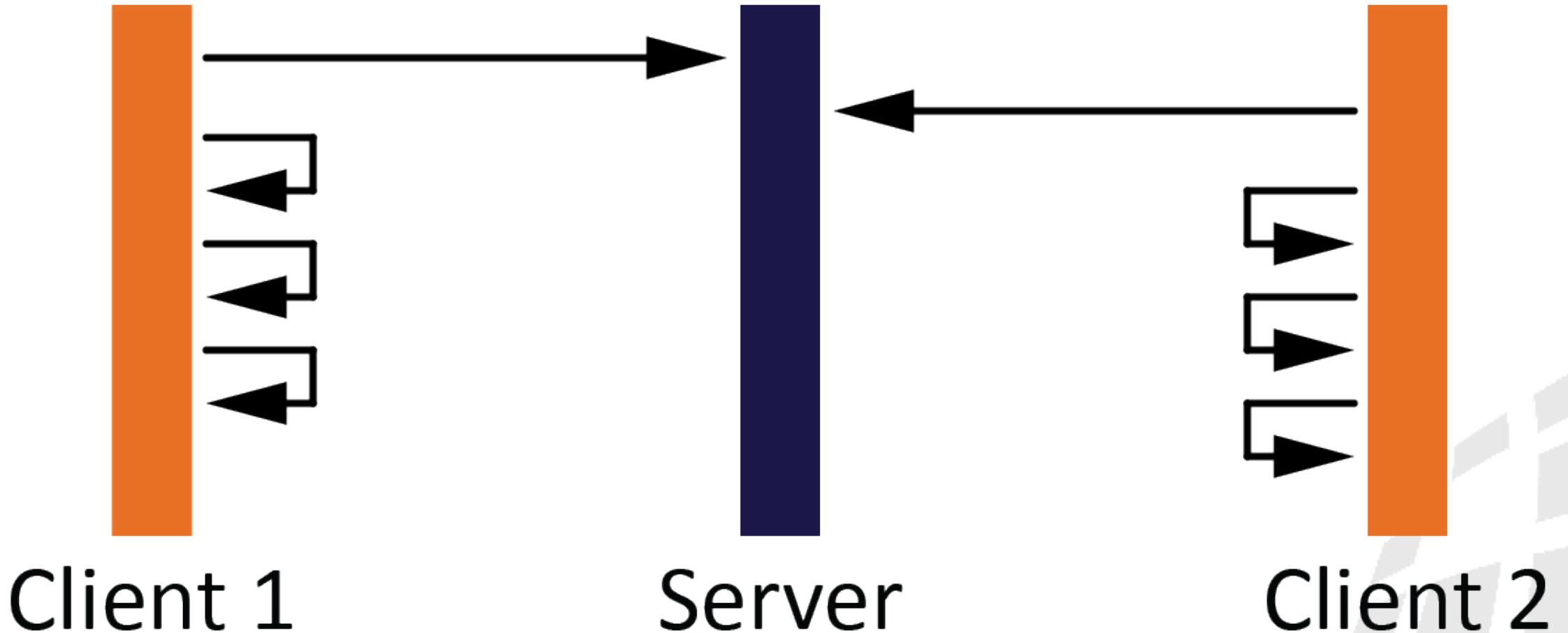
Разумные требования: счетчик

Уникальный, монотонно возрастающий, 8 байт!



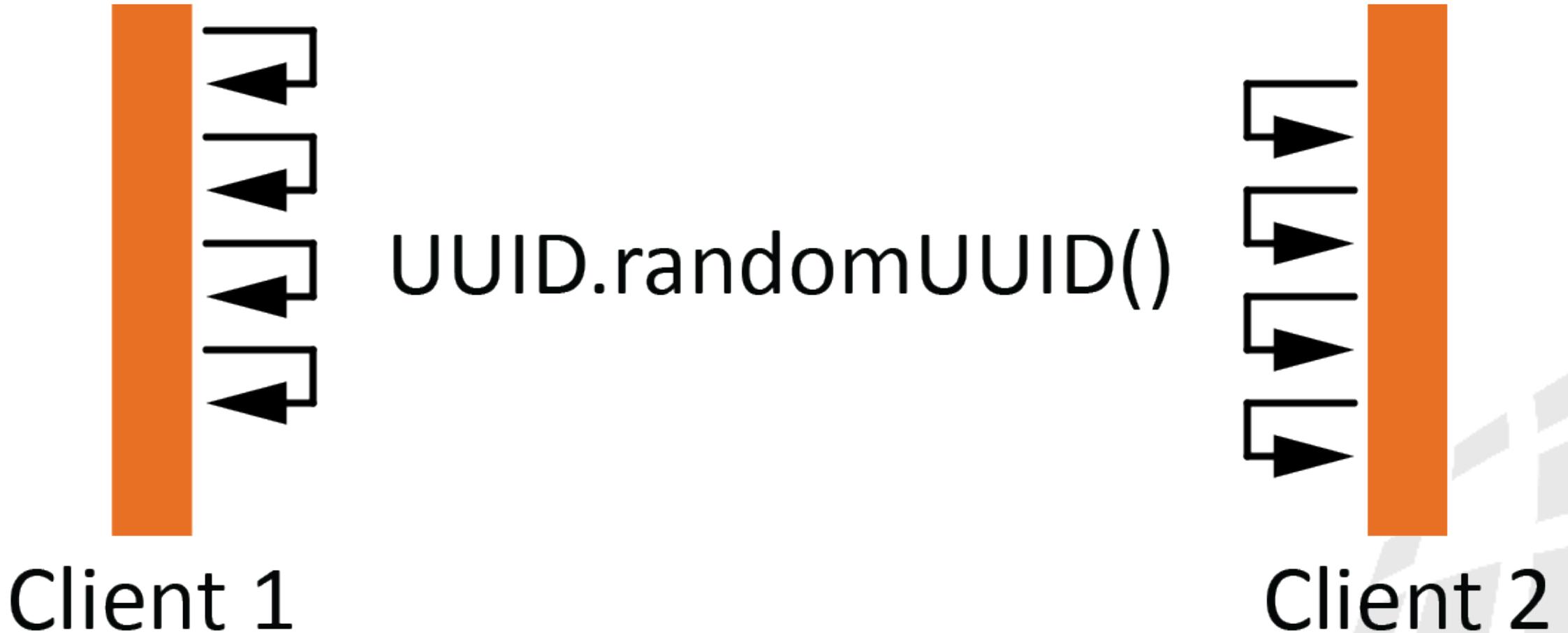
Разумные требования: счетчик

Уникальный, ~~монотонно-возрастающий~~, 8 байт!



Разумные требования: счетчик

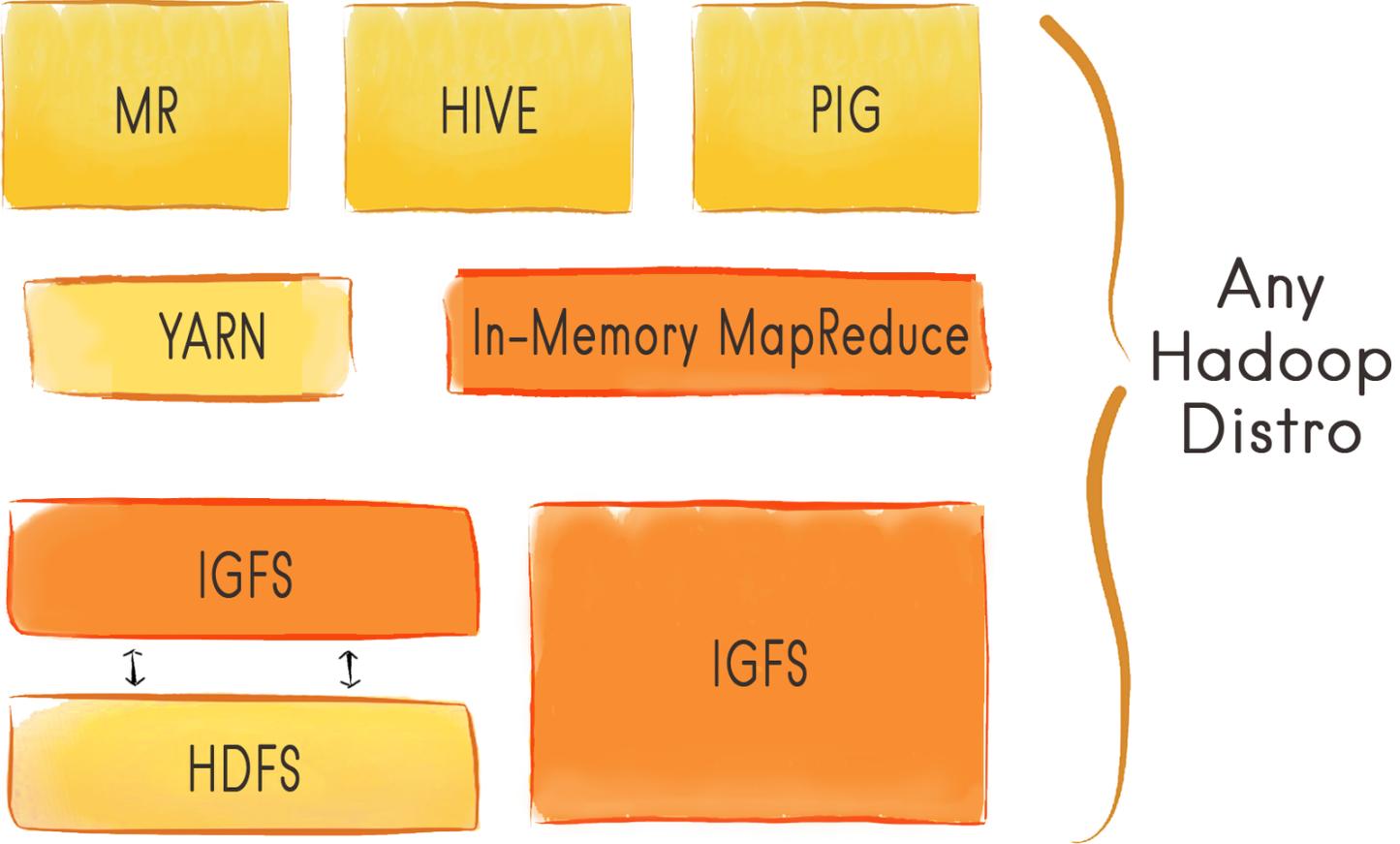
Уникальный, ~~монотонно-возрастающий~~, ~~8 байт!~~



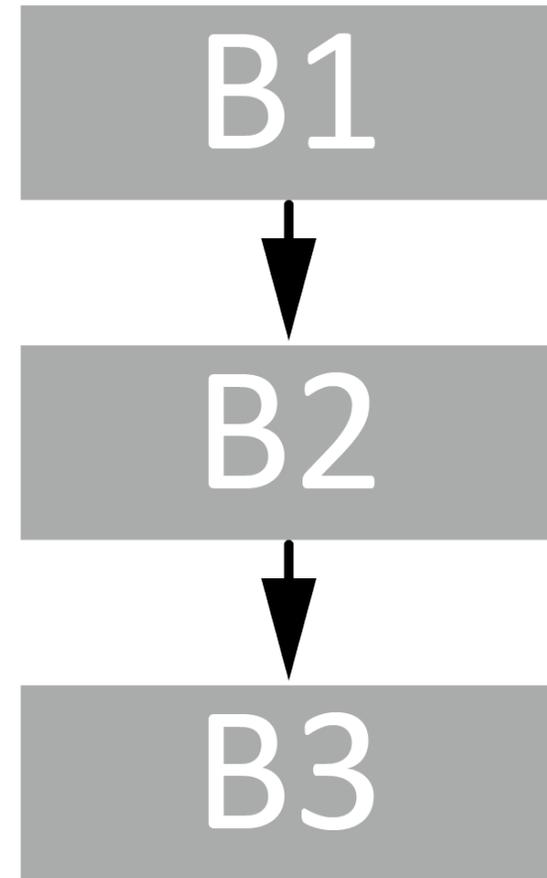
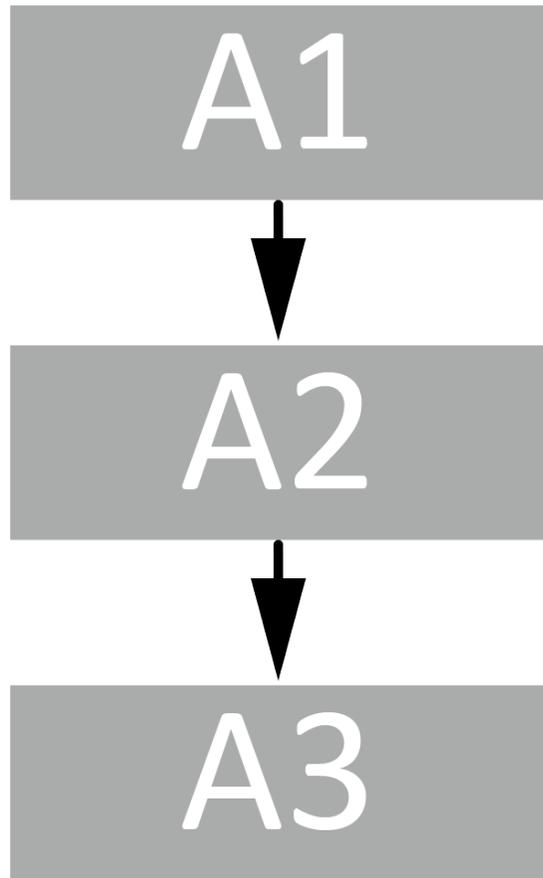
Координация



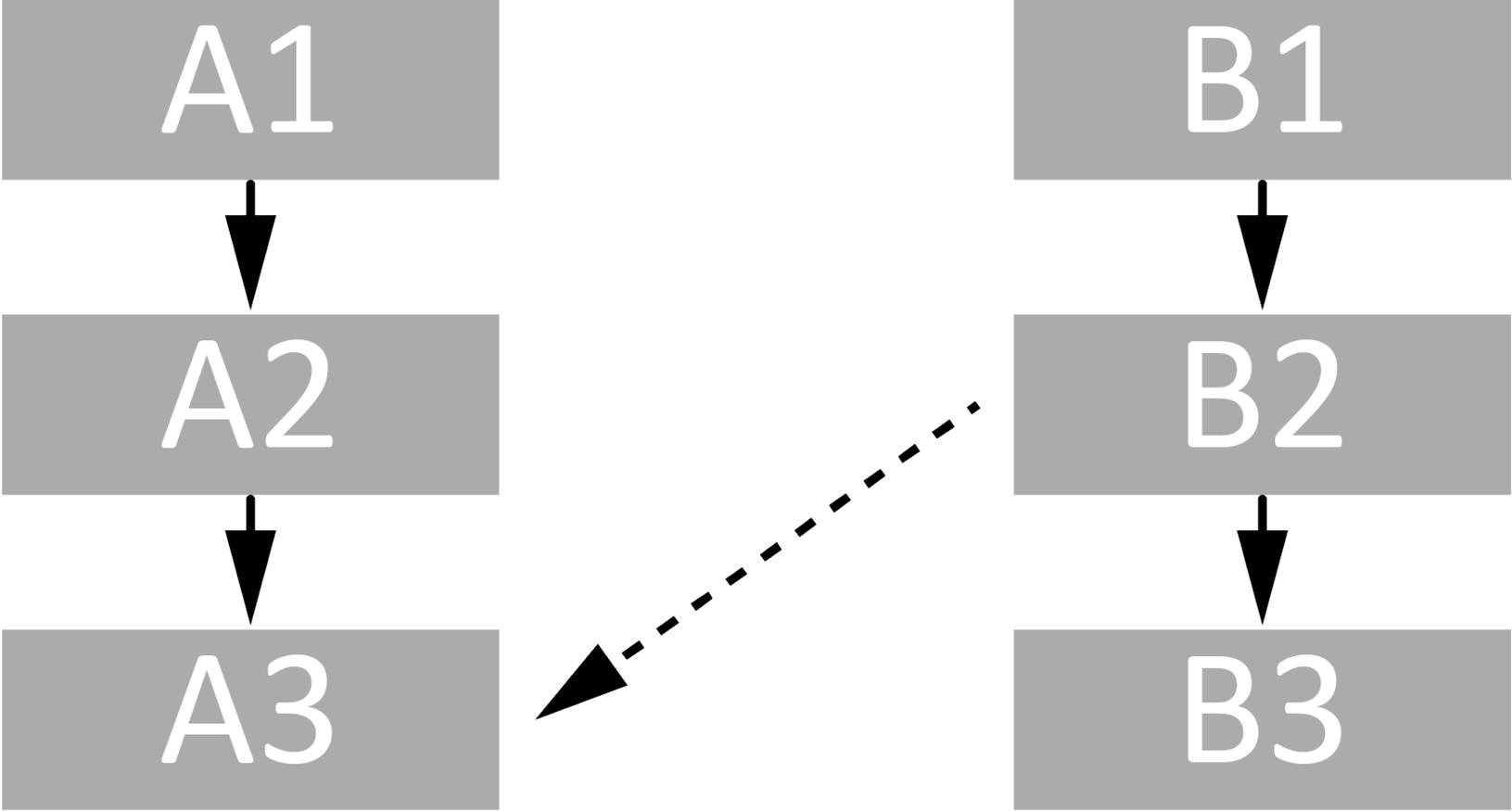
Разумные требования: IGFS



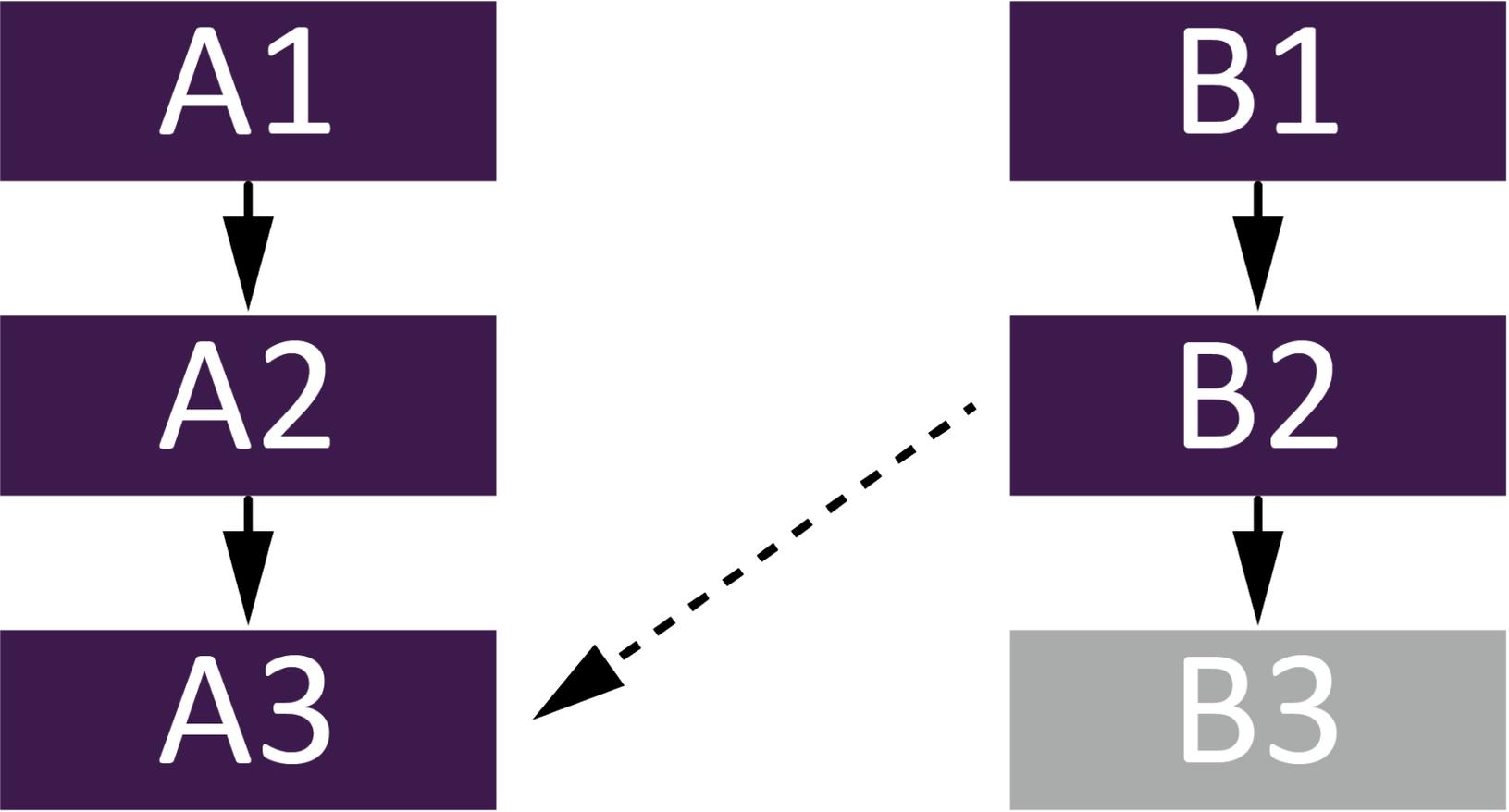
Разумные требования: IGFS



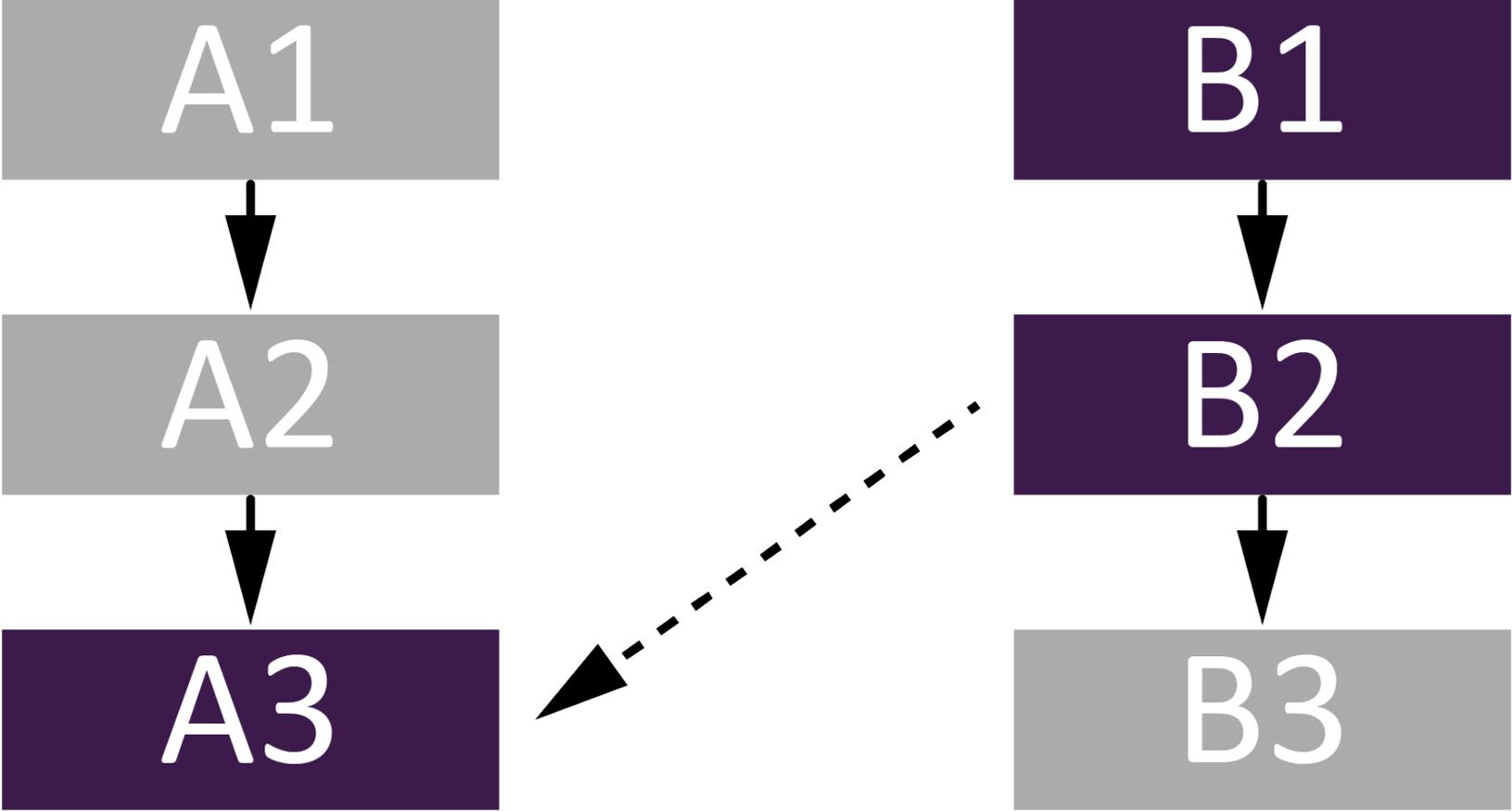
Разумные требования: IGFS



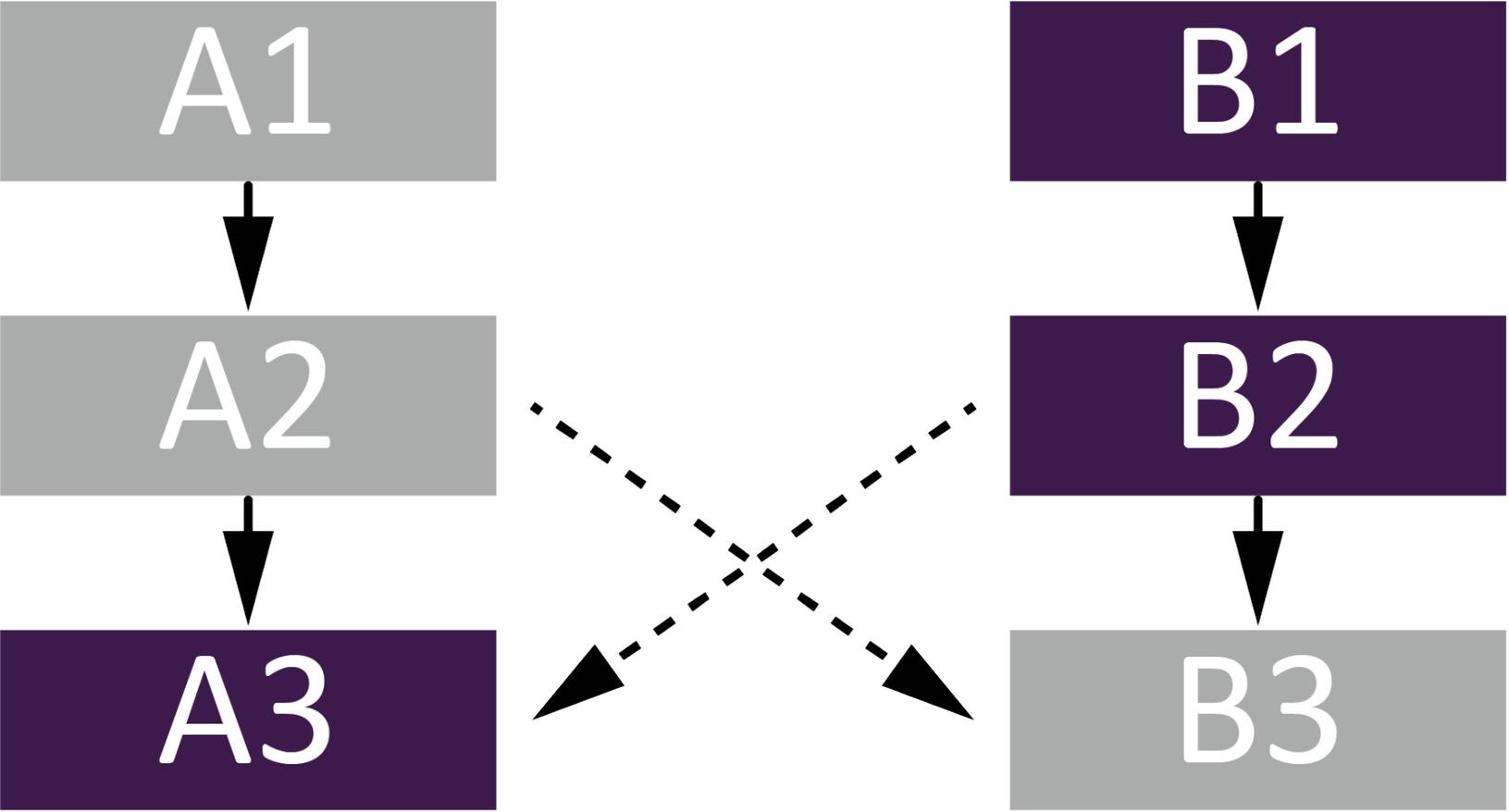
Разумные требования: IGFS



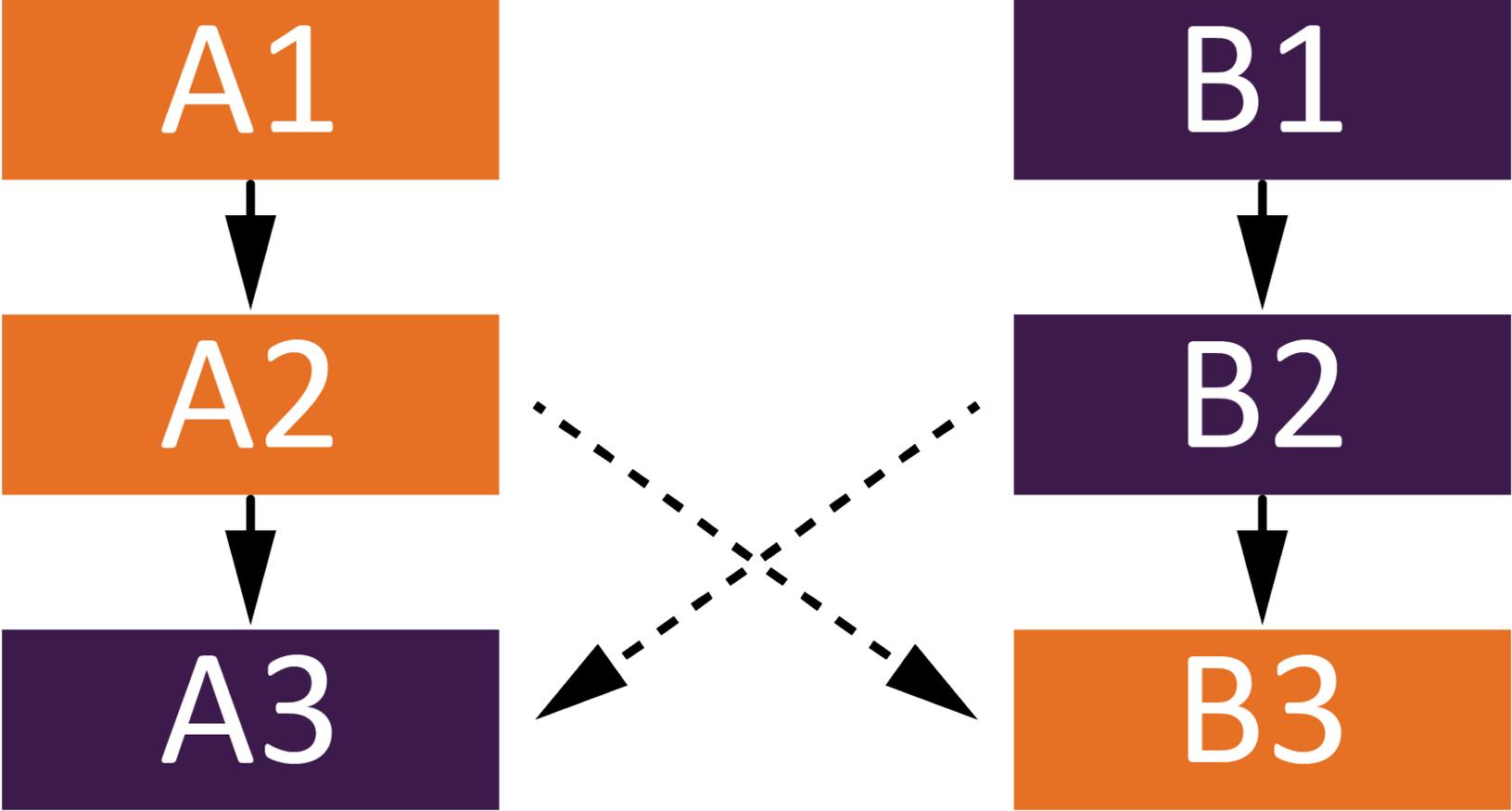
Разумные требования: IGFS



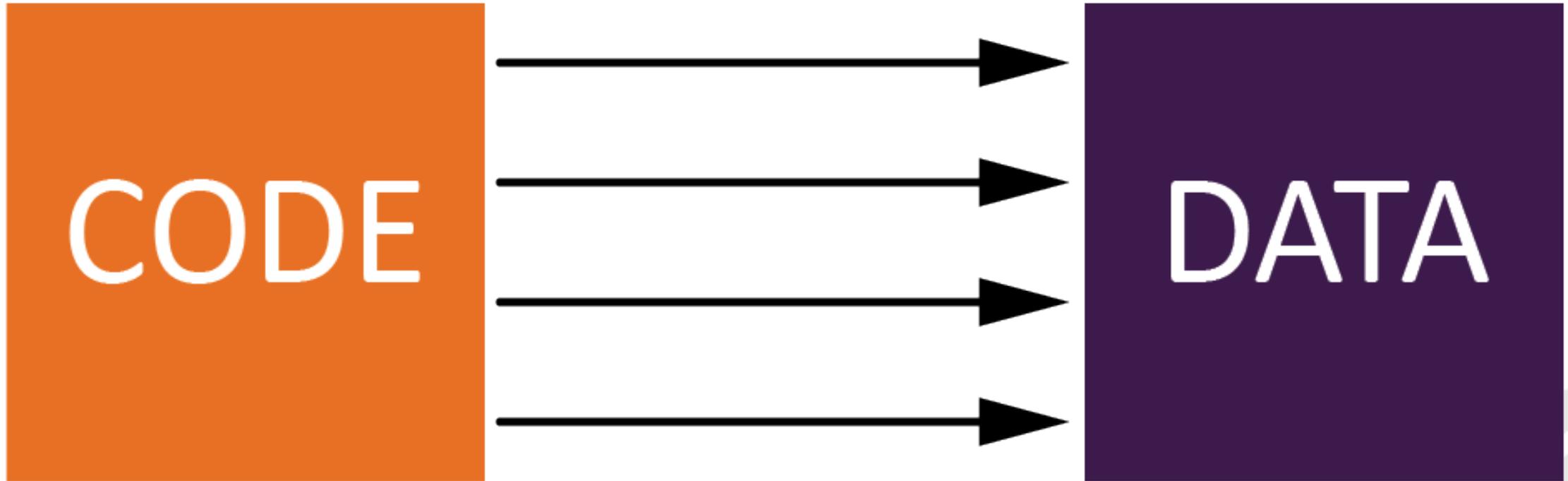
Разумные требования: IGFS



Разумные требования: IGFS



Данные к коду



Код к данным



Код к данным: а если баг?

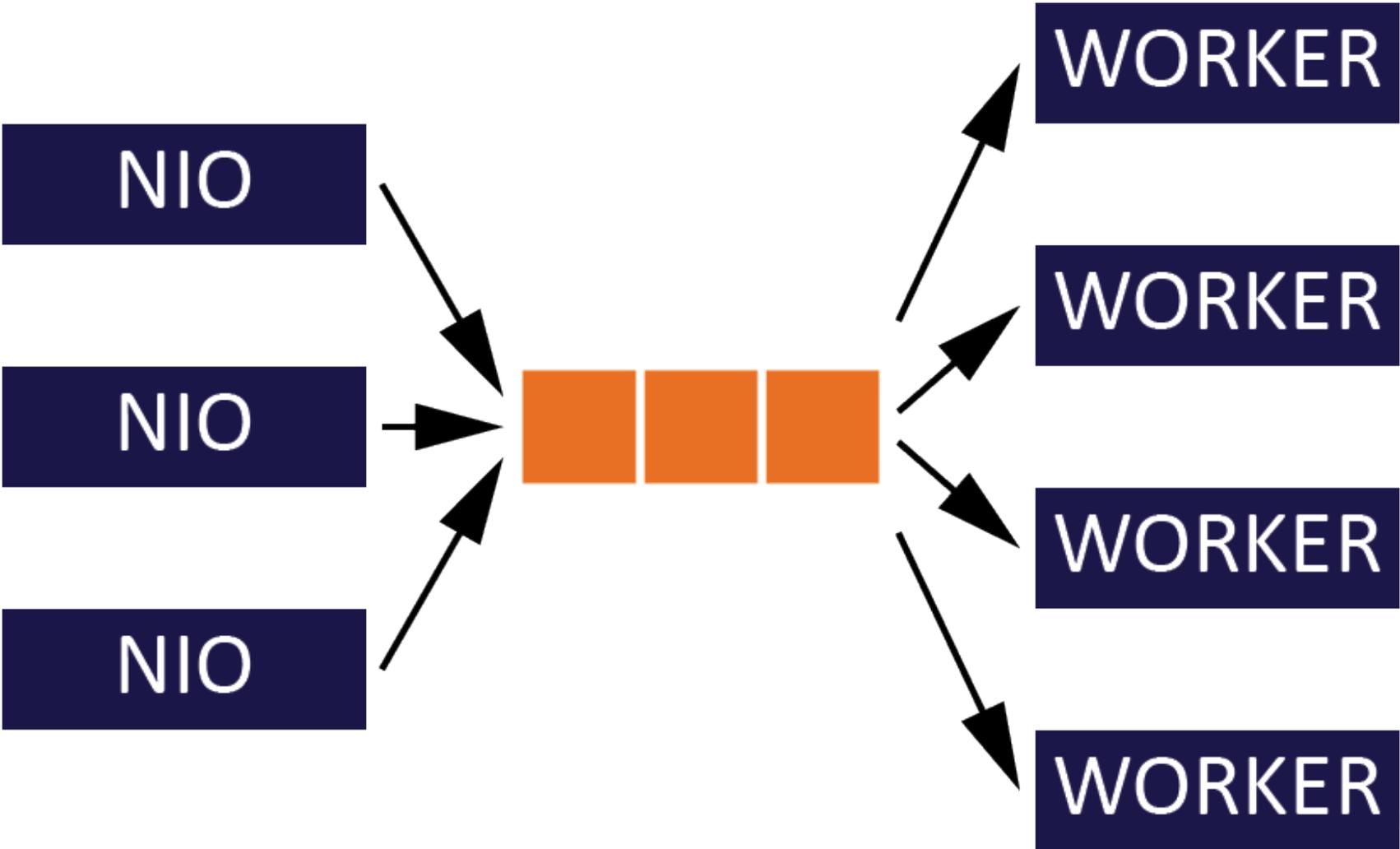


План

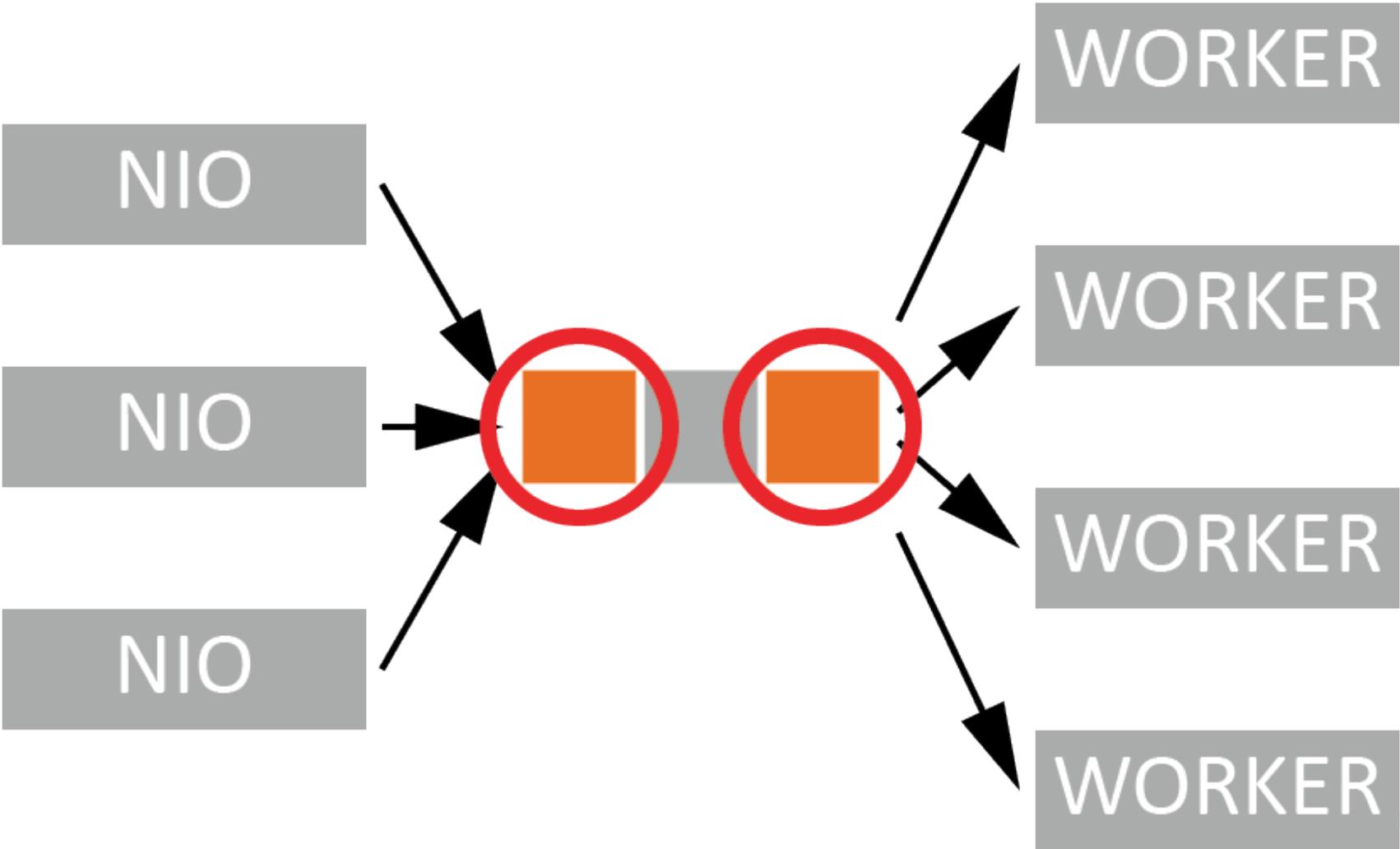
- Алгоритмы шардирования
- Data co-location
- Синхронизация в кластере
- **Локальная архитектура многопоточности**



Локальное распределение задач



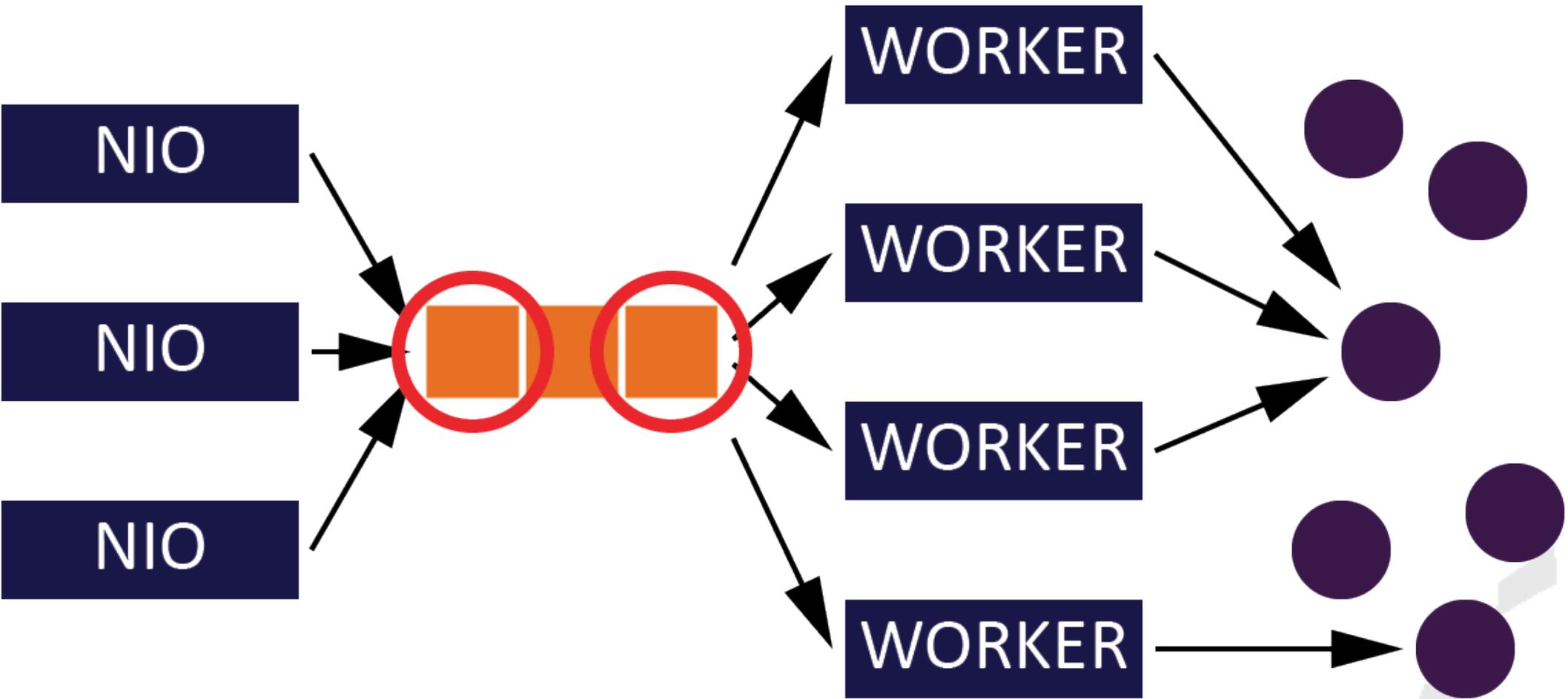
Локальное распределение задач



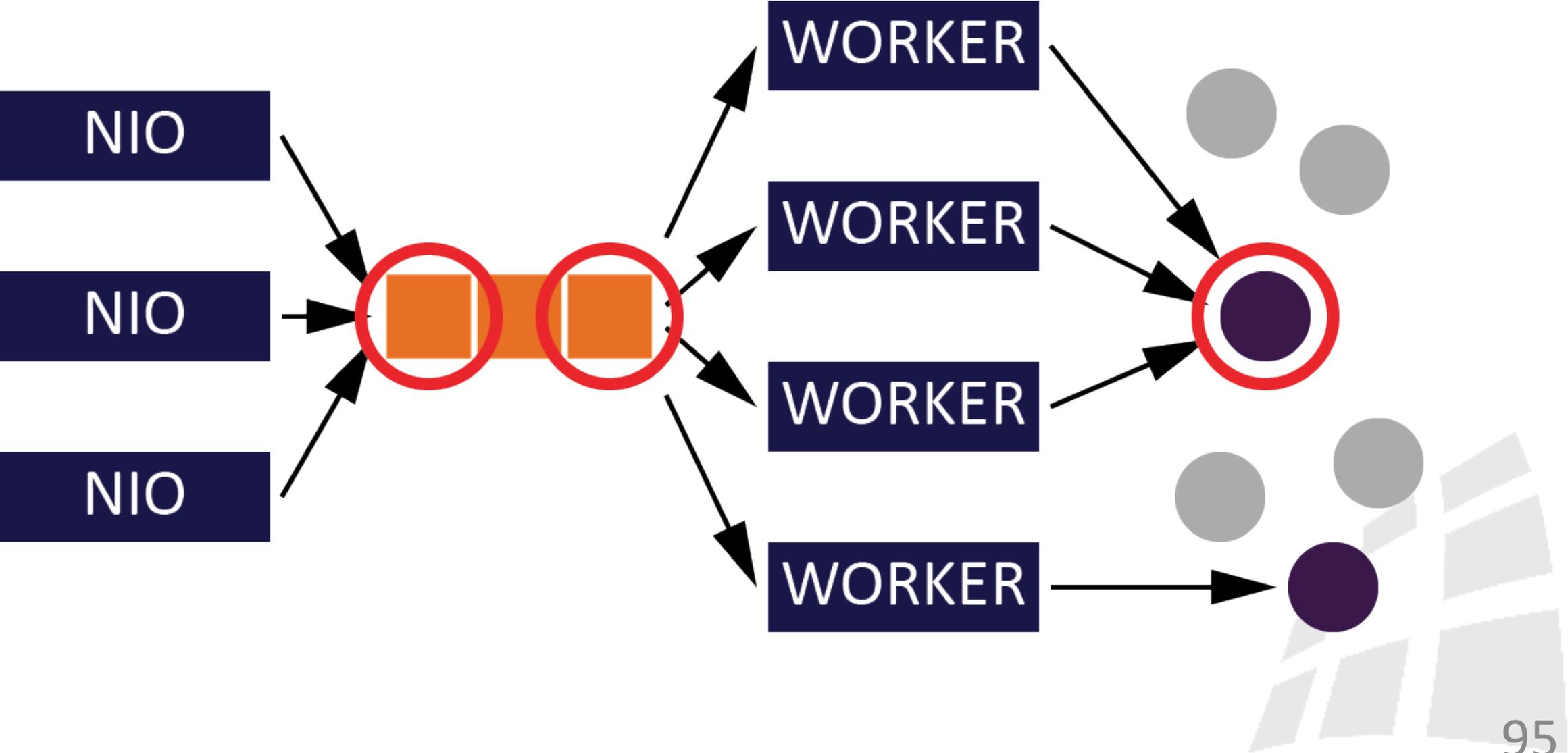
Локальное распределение задач



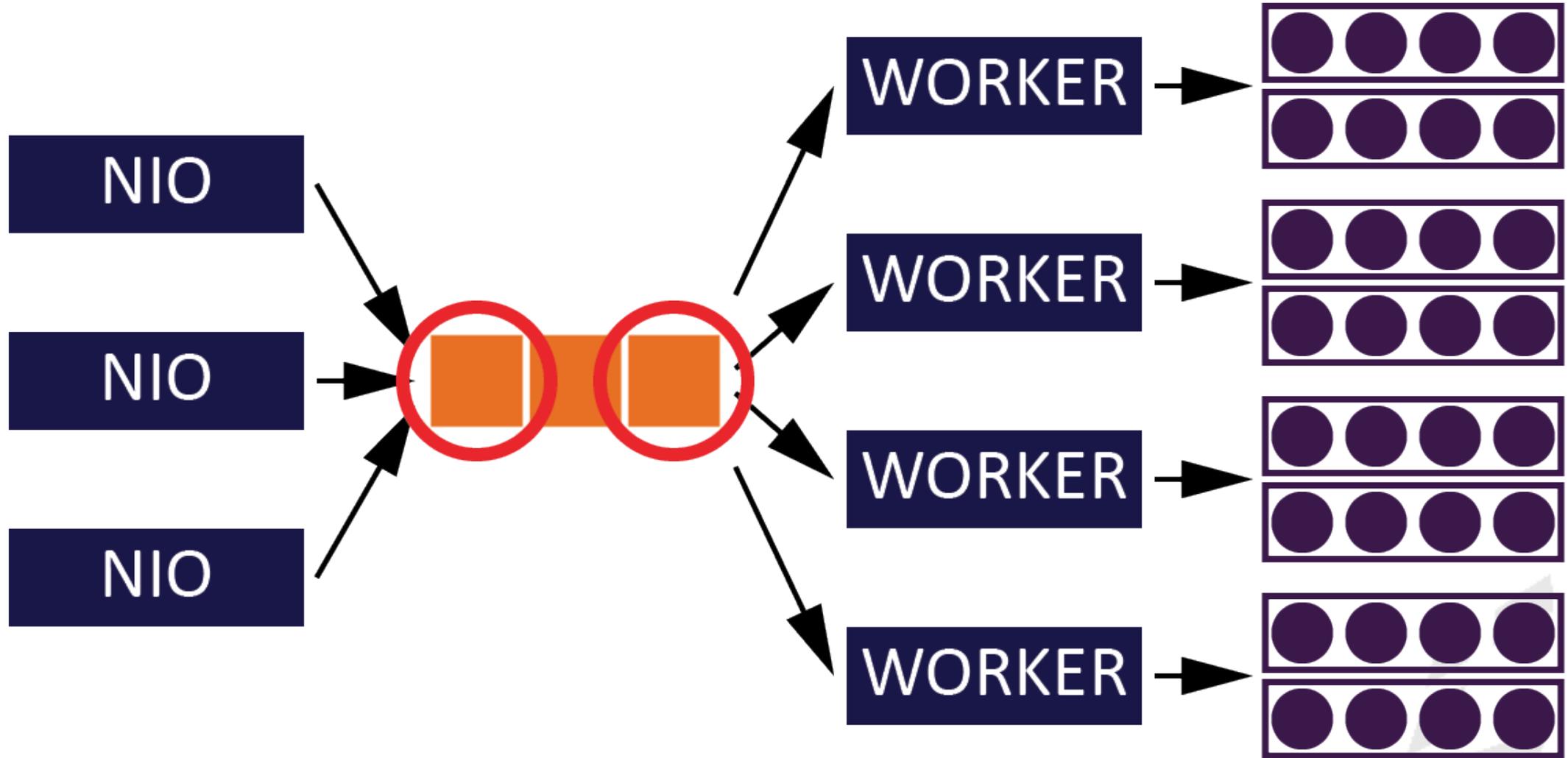
Локальное распределение задач



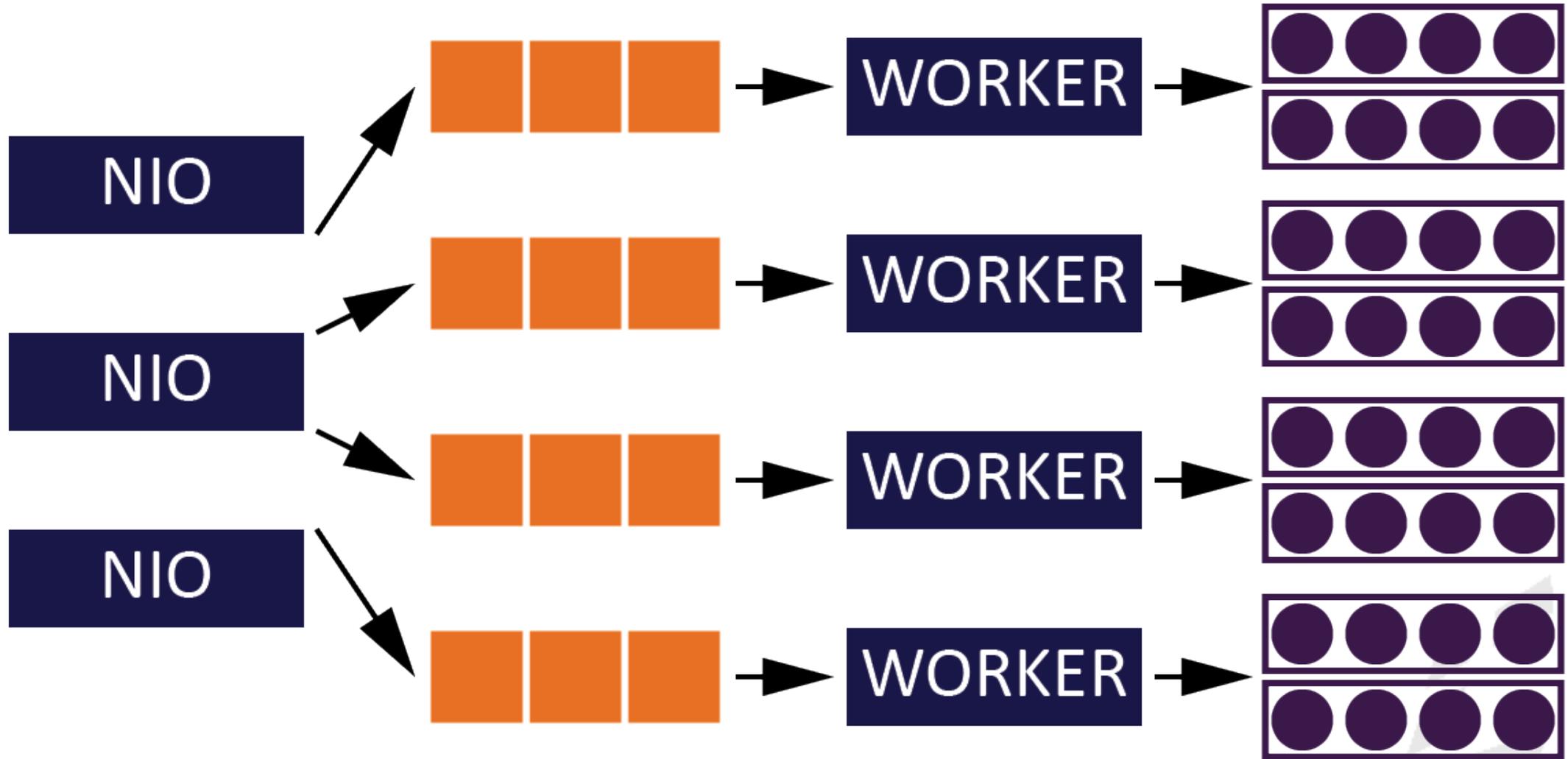
Локальное распределение задач



Thread-per-partition



Thread-per-partition



Thread-per-partition

A large red square containing the letters 'MSK' in white, bold, sans-serif font.

MSK

Node 1

100 TPS

A dark blue square containing the letters 'SPB' in white, bold, sans-serif font.

SPB

Node 2

50 TPS

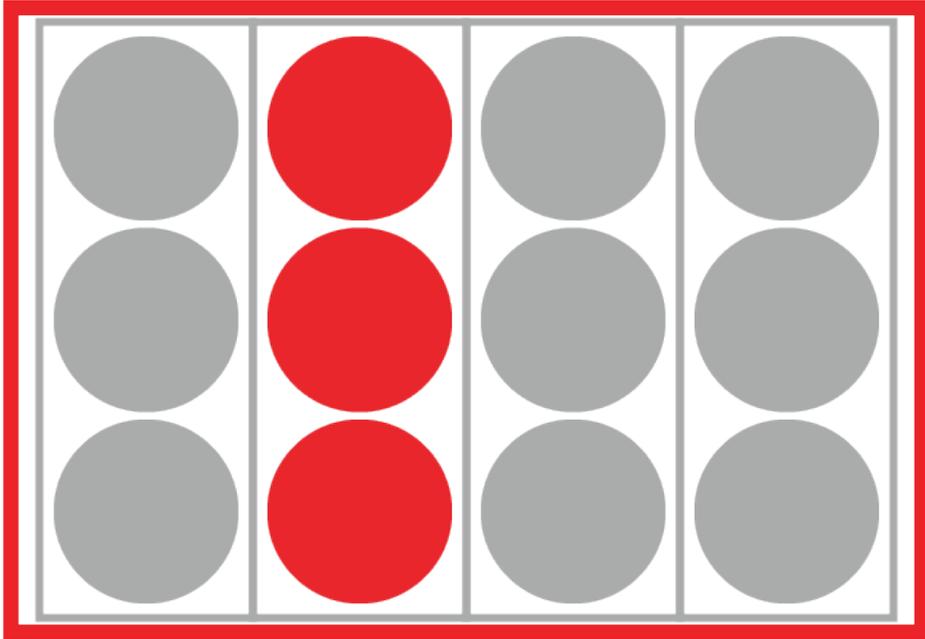
A dark blue square containing the letters 'NSK' in white, bold, sans-serif font.

NSK

Node 3

10 TPS

Thread-per-partition



Node 1

100 TPS



Node 2

50 TPS



Node 3

10 TPS

Итого

- Шардинг: важен не только баланс, но и стабильность



Итого

- Шардинг: важен не только баланс, но и стабильность
- Co-location: баланс или эффективность



Итого

- Шардинг: важен не только баланс, но и стабильность
- Co-location: баланс или эффективность
- Модель данных должна быть адаптирована



Итого

- Шардинг: важен не только баланс, но и стабильность
- Co-location: баланс или эффективность
- Модель данных должна быть адаптирована
- Координация – враг производительности



Итого

- Шардинг: важен не только баланс, но и стабильность
- Co-location: баланс или эффективность
- Модель данных должна быть адаптирована
- Координация – враг производительности
- Thread-per-partition – ускоряет простые сценарии, может замедлить сложные

Вопросы?

