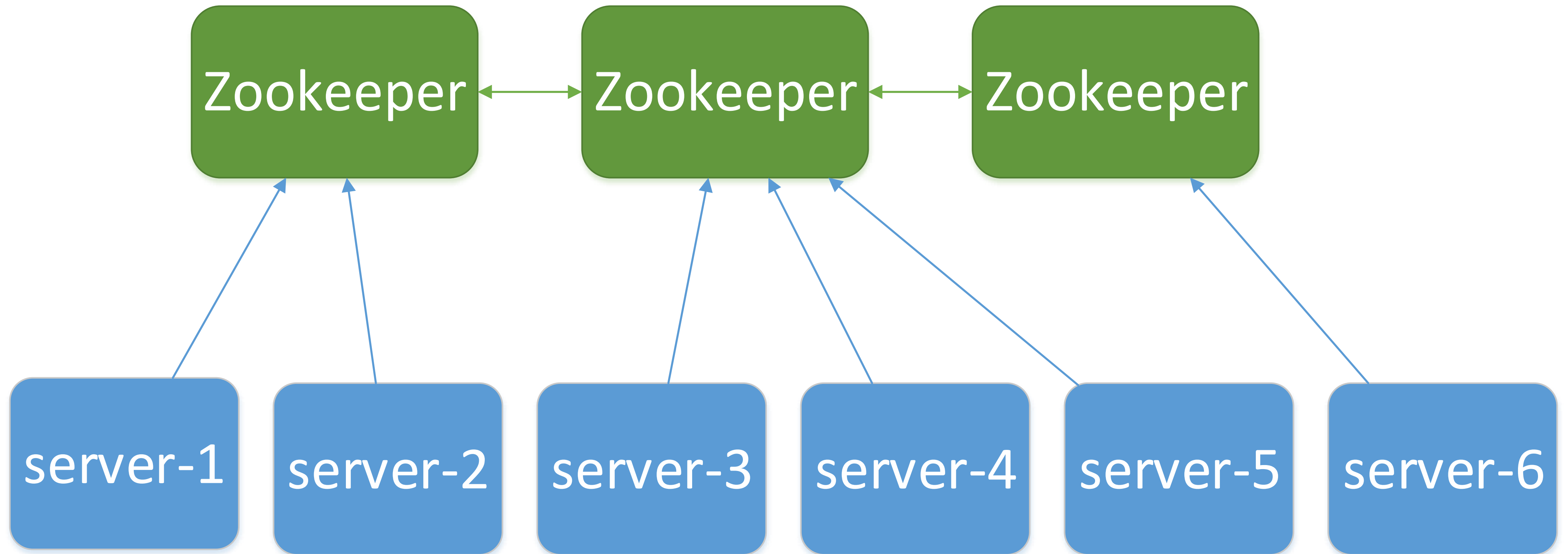


# Как не наступить на распределенные грабли или просто хорошие практики по работе с Kafka, Zookeeper и HBase

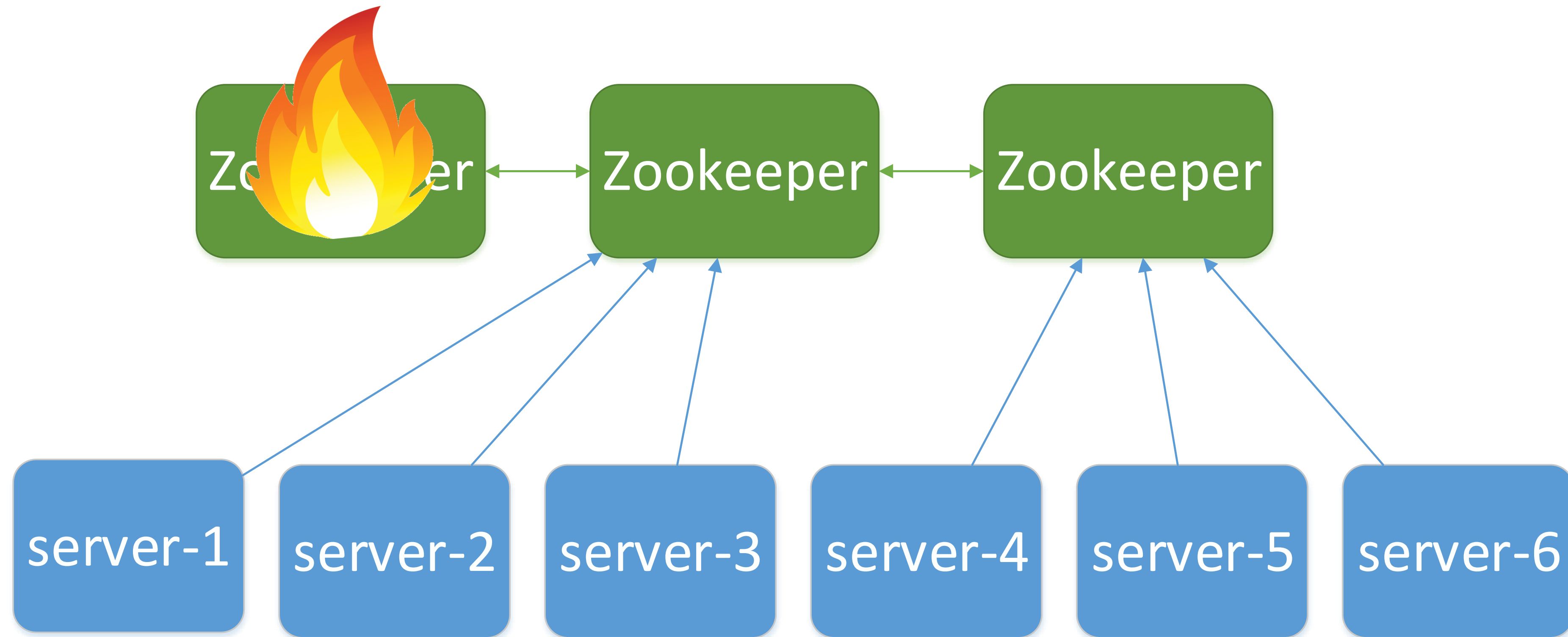
Работа с очередями или persistent хранилищами становится не такой очевидной, когда дело касается распределенных систем. Вопросы consistency и data availability требуют особого внимания, когда речь идет об обработке критичных данных. Мы расскажем о подводных камнях, которые нам попались, а также поделимся интересными практиками решения типовых задач.



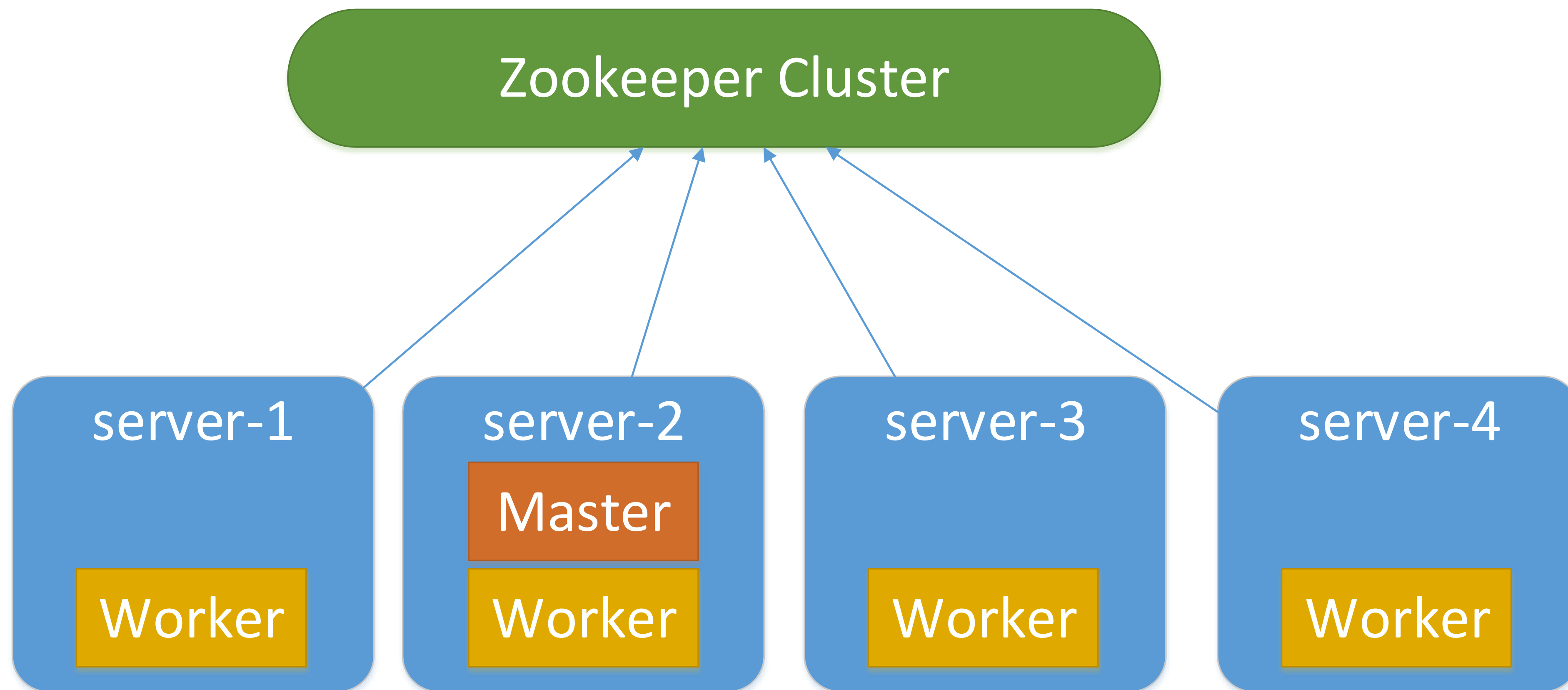




Consul



# distributed-job-manager





```
job-manager
├── locks
│   └── work-pooled
│       ├── async.report.building.job
│       ├── elasticsearch.upload.job
│       ├── fraud.user.event.export.job
│       ├── notification.processing.job
│       └── statistics.log.kafka.reader.job
│       ...
├── assignment-version
├── workers
│   ├── 20
│   ├── ...
│   └── 3
│       ├── available
│       │   └── work-pooled
│       │       ├── async.report.building.job
│       │       ├── cdr.upload.service
│       │       ├── deferred.sms.notification.job
│       │       ├── elasticsearch.upload.job
│       │       └── email.sender.job
│       │       ...
│       └── assigned
│           └── work-pooled
│               ├── async.report.building.job
│               └── elasticsearch.upload.job
```





## job-manager

```
└─ locks
   └─ work-pooled
      └─ async.report.building.job
      └─ elasticsearch.upload.job
      └─ fraud.user.event.export.job
      └─ notification.processing.job
      └─ statistics.log.kafka.reader.job
      ...
└─ assignment-version
└─ workers
   └─ 20
   └─ ...
   └─ 3
      └─ available
         └─ work-pooled
            └─ async.report.building.job
            └─ cdr.upload.service
            └─ deferred.sms.notification.job
            └─ elasticsearch.upload.job
            └─ email.sender.job
            ...
      └─ assigned
         └─ work-pooled
            └─ async.report.building.job
            └─ elasticsearch.upload.job
```



## job-manager

```
└─ locks
    └─ work-pooled
        └─ async.report.building.job
        └─ elasticsearch.upload.job
        └─ fraud.user.event.export.job
        └─ notification.processing.job
        └─ statistics.log.kafka.reader.job
        ...
└─ assignment-version
└─ workers
    └─ 20
    ...
    └─ 3
        └─ available
            └─ work-pooled
                └─ async.report.building.job
                └─ cdr.upload.service
                └─ deferred.sms.notification.job
                └─ elasticsearch.upload.job
                └─ email.sender.job
                ...
        └─ assigned
            └─ work-pooled
                └─ async.report.building.job
                └─ elasticsearch.upload.job
```

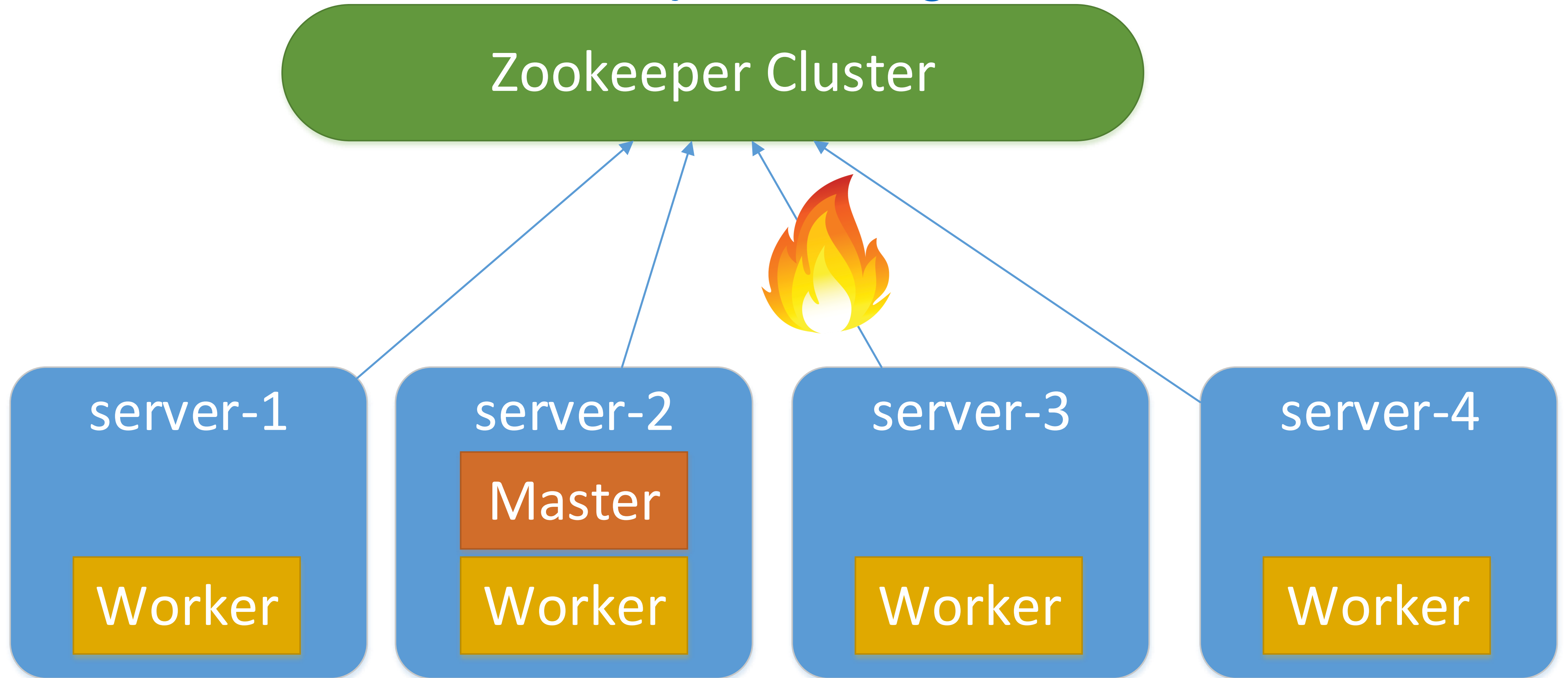




```
job-manager
├── locks
│   └── work-pooled
│       ├── async.report.building.job
│       ├── elasticsearch.upload.job
│       ├── fraud.user.event.export.job
│       ├── notification.processing.job
│       └── statistics.log.kafka.reader.job
│       ...
├── assignment-version
├── workers
│   ├── 20
│   ├── ...
│   └── 3
│       ├── available
│       │   └── work-pooled
│       │       ├── async.report.building.job
│       │       ├── cdr.upload.service
│       │       ├── deferred.sms.notification.job
│       │       ├── elasticsearch.upload.job
│       │       └── email.sender.job
│       │       ...
│       └── assigned
│           └── work-pooled
│               ├── async.report.building.job
│               └── elasticsearch.upload.job
```



## distributed-job-manager





```
class PersistentExpiringDistributedLock{
    public boolean expirableAcquire(long acquirePeriod, long timeout) throws Exception {...}
    public boolean checkAndProlong(long prolongationPeriod) throws Exception {...}
}
```

```
PersistentExpiringDistributedLock lock = getDistributedLock();
```

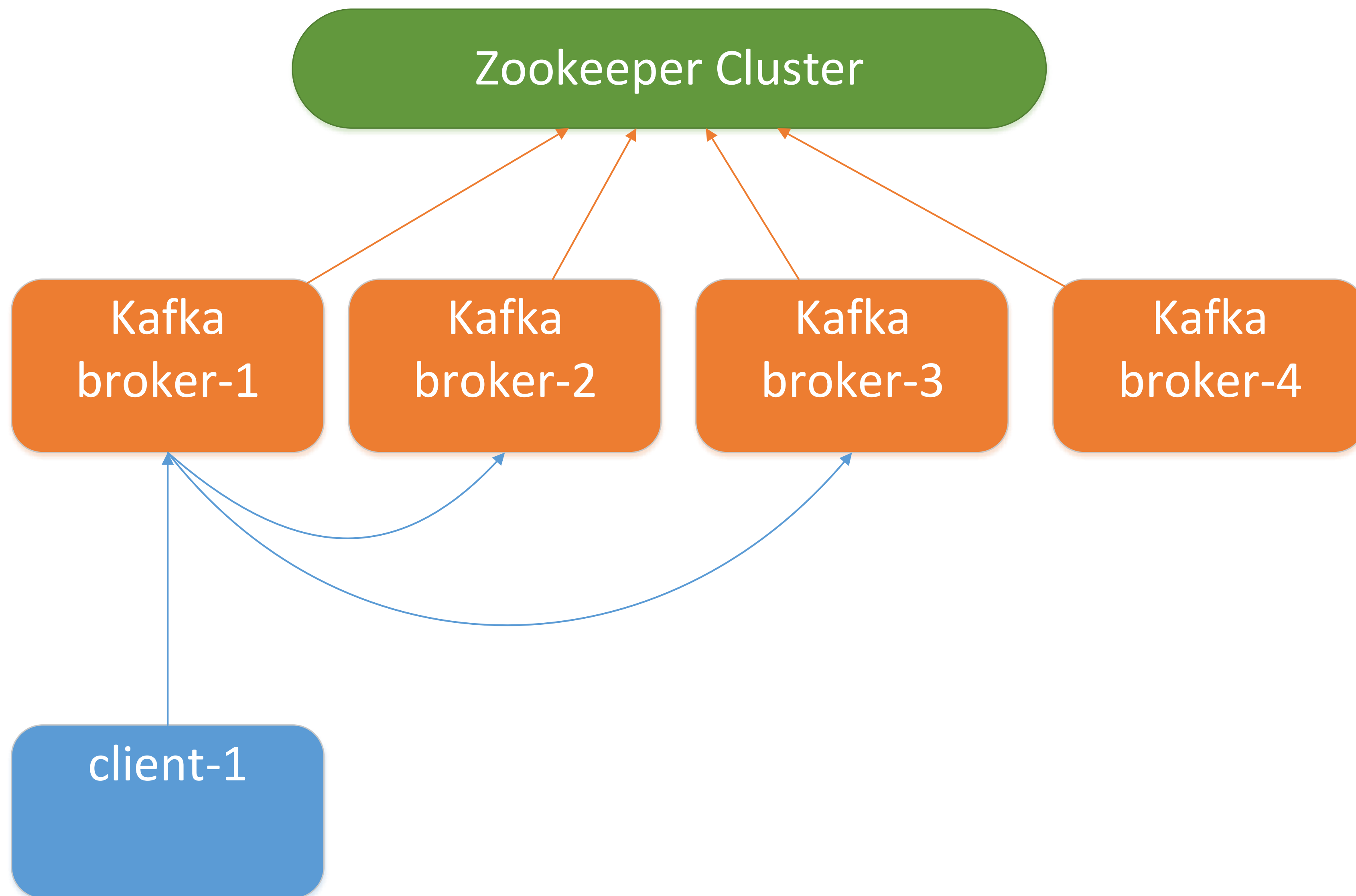
```
if (lock.expirableAcquire(MINUTES.toMillis(30), SECONDS.toMillis(5))) {
    while (haveWorkToDo()) {
        // do 10 minutes work
        doWork();
        // prolong persistent lock
        if (!lock.checkAndProlong(MINUTES.toMillis(30))) {
            //something gone wrong
            //gracefully stop task
            return;
        }
    }
}
```

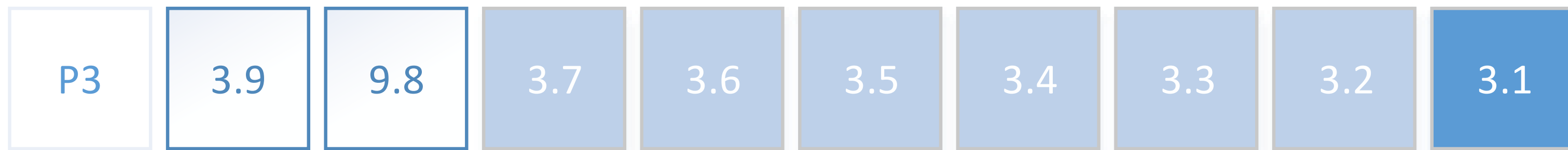
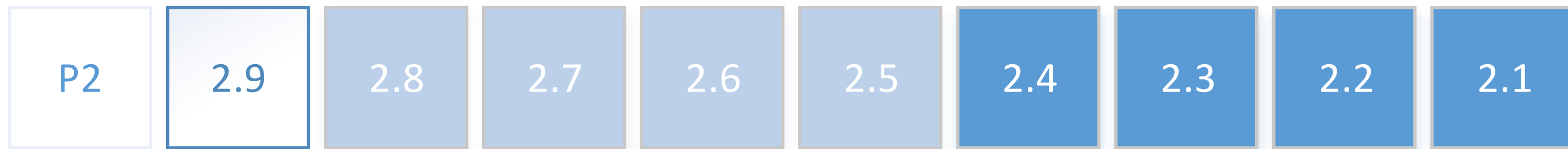
~~Ephemeral~~



```
job-manager
├── locks
│   └── ...
├── assignment-version
├── workers
│   ├── 20
│   ├── ...
│   └── 3
│       ├── available
│       │   └── work-pooled
│       │       ├── async.report.building.job
│       │       ├── cdr.upload.service
│       │       ├── deferred.sms.notification.job
│       │       ├── elasticsearch.upload.job
│       │       └── email.sender.job
│       │       ...
│       └── assigned
│           └── work-pooled
│               ├── async.report.building.job
│               └── elasticsearch.upload.job
```

```
public class ZooKeeper {
    ...
    public List<OpResult> multi(Iterable<Op> ops) throws Exception{...}
    ...
}
```





Last read offset

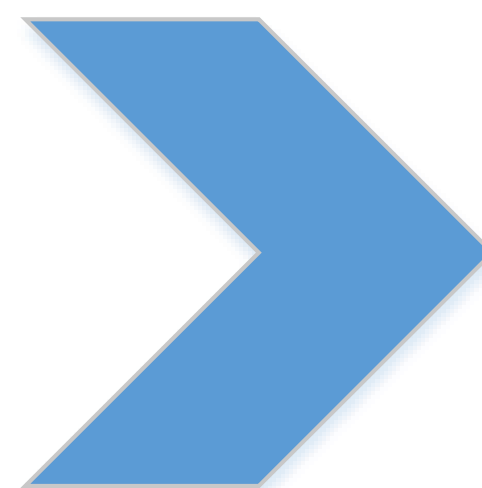


Last committed offset



2 – partition number  
8 – item number



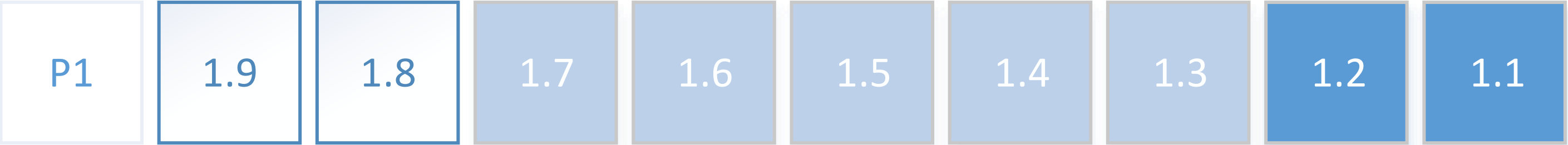


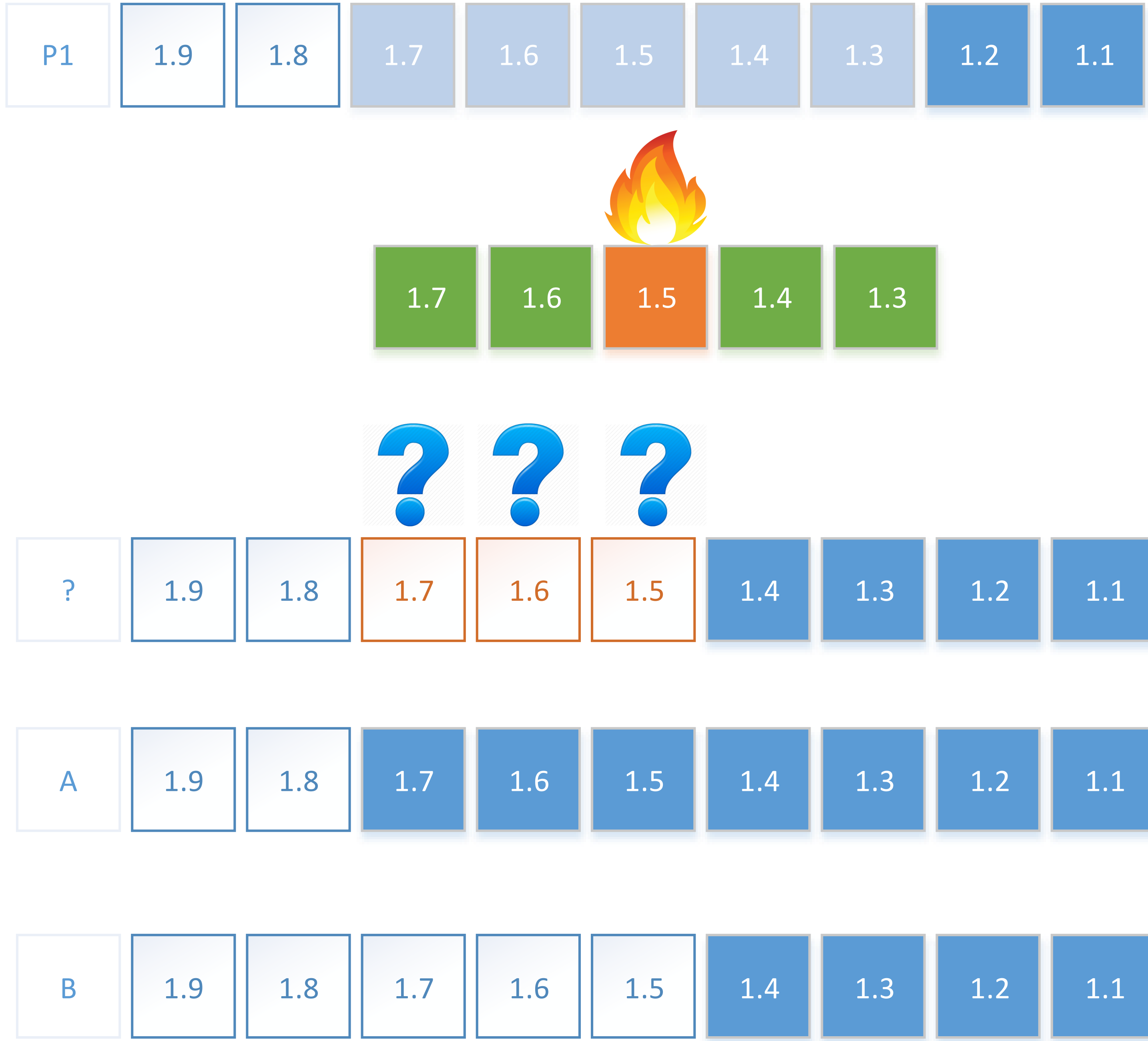
P1	1.9 25	1.8 22	1.7 19	1.6 16	1.5 13	1.4 10	1.3 7	1.2 4	1.1 1
----	-----------	-----------	-----------	-----------	-----------	-----------	----------	----------	----------

P2	2.9 26	2.8 23	2.7 20	2.6 17	2.5 14	2.4 11	2.3 8	2.2 5	2.1 2
----	-----------	-----------	-----------	-----------	-----------	-----------	----------	----------	----------

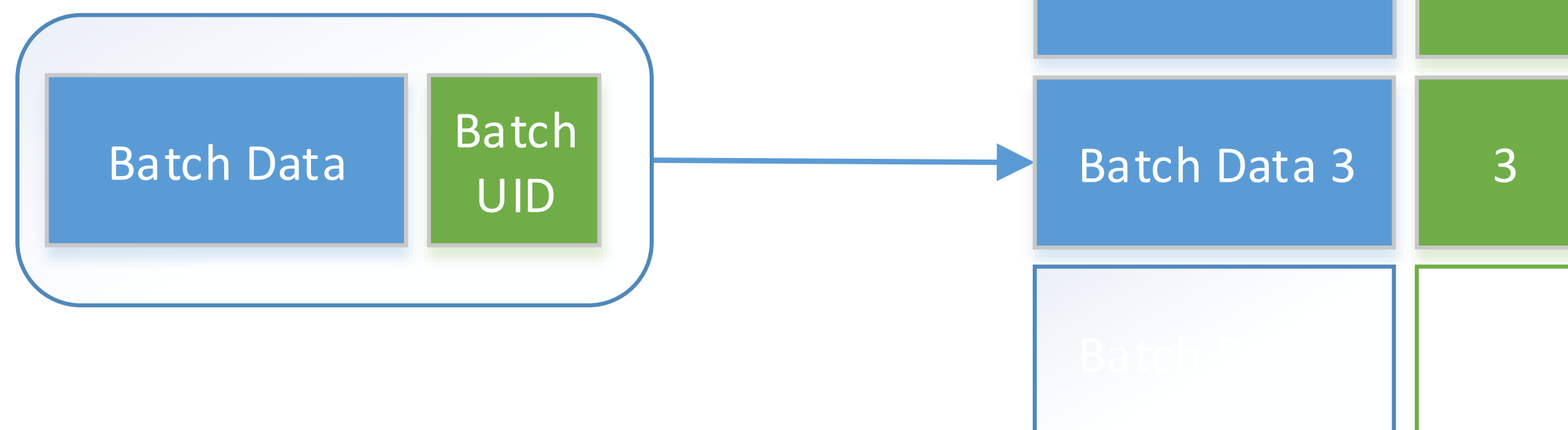
P3	3.9 27	9.8 24	3.7 21	3.6 18	3.5 15	3.4 12	3.3 9	3.2 6	3.1 3
----	-----------	-----------	-----------	-----------	-----------	-----------	----------	----------	----------

1.1 1	1.2 4	1.3 7	1.4 10	1.5 13	1.6 16	1.7 19	2.1 2	2.2 5	2.3 8
----------	----------	----------	-----------	-----------	-----------	-----------	----------	----------	----------





# Idempotency



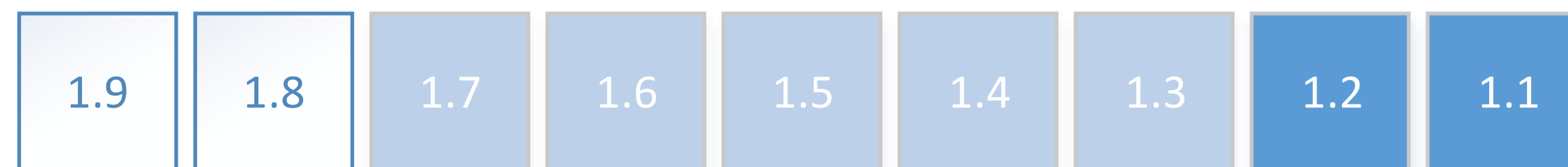


{userId:107, event: 'find-sword'}  
{userId:107, event: 'kill-monster-by-sword'}  
{userId:107, event: 'lose-sword'}



partition = hash(userId)





```
KafkaConsumer consumer = newConsumer();  
while (!isNeedToShutdown()) {  
    val records = consumer.poll(SECONDS.toMillis(10));  
    processRecords(records);  
    // if processRecords take too much time kafka session will die  
    consumer.commitSync();  
}
```



# KStreams?



KStms?

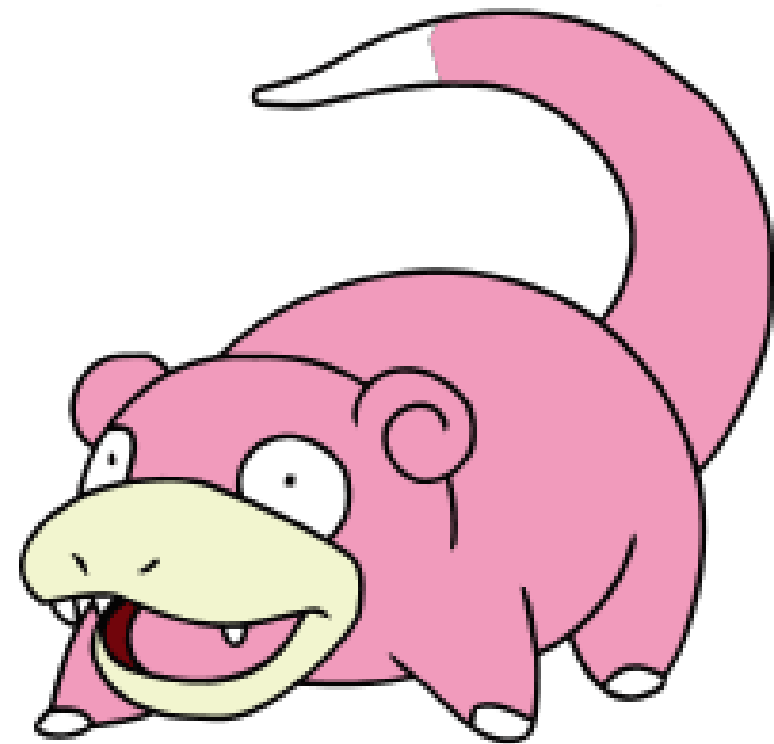




```
Map config = new HashMap();  
config.put(ConsumerConfig.GROUP_ID_CONFIG, group);  
config.put(ConsumerConfig.ENABLE_AUTO_COMMIT_CONFIG, "false");  
config.put(ConsumerConfig.AUTO_OFFSET_RESET_CONFIG, "earliest");  
...  
KafkaConsumer consumer = new KafkaConsumer(config);
```



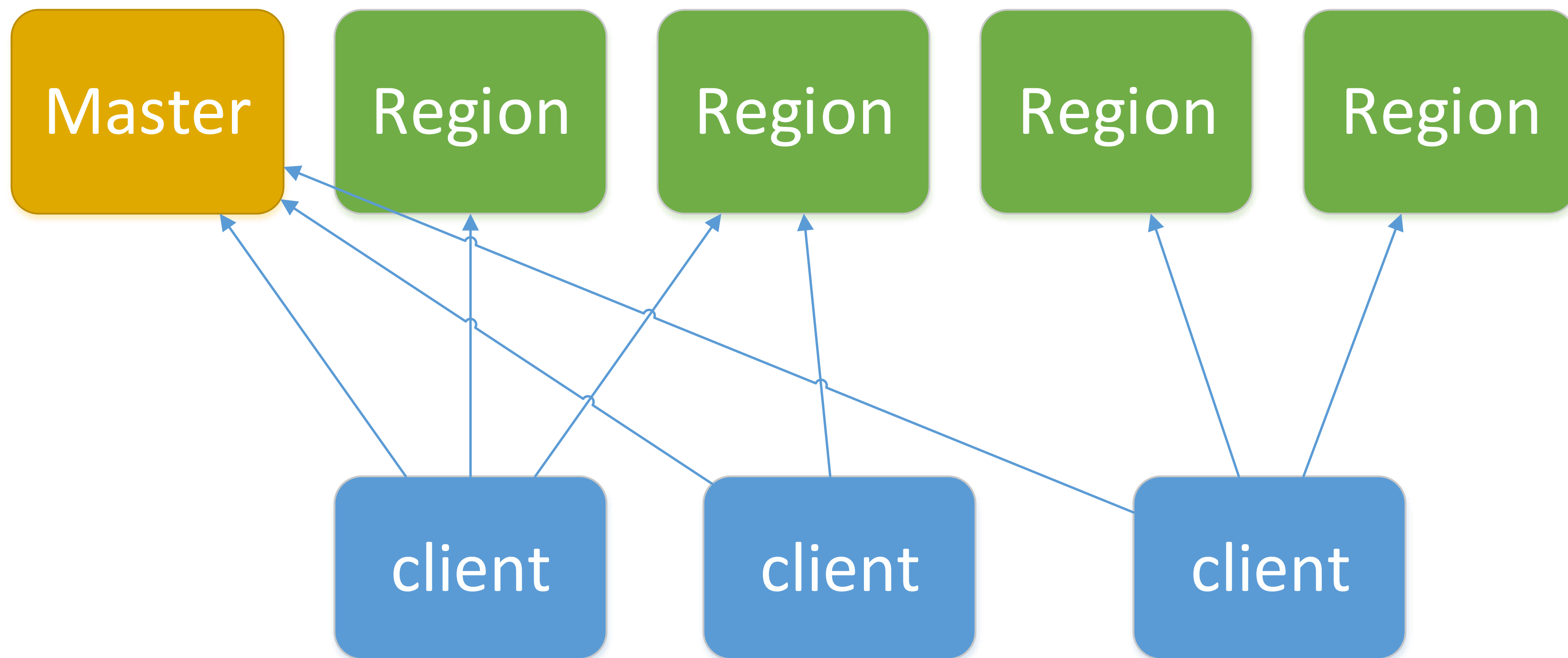
CAS!

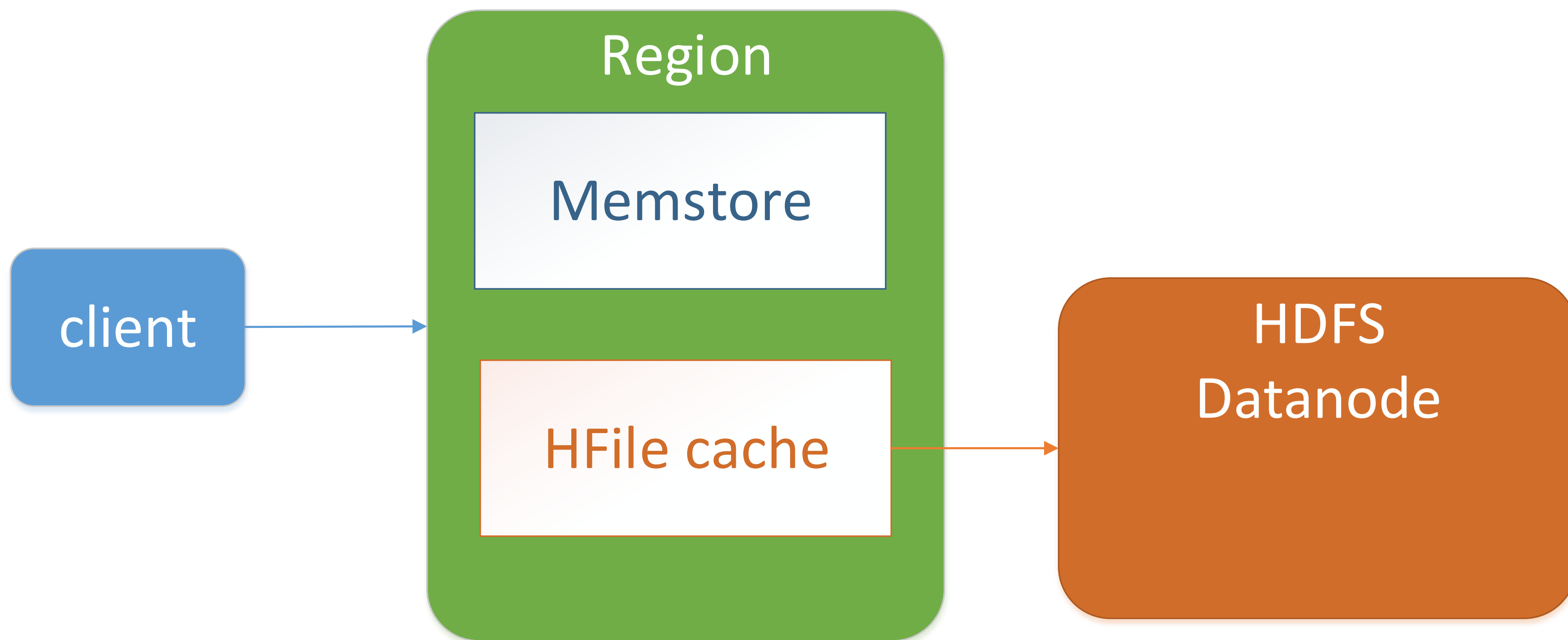


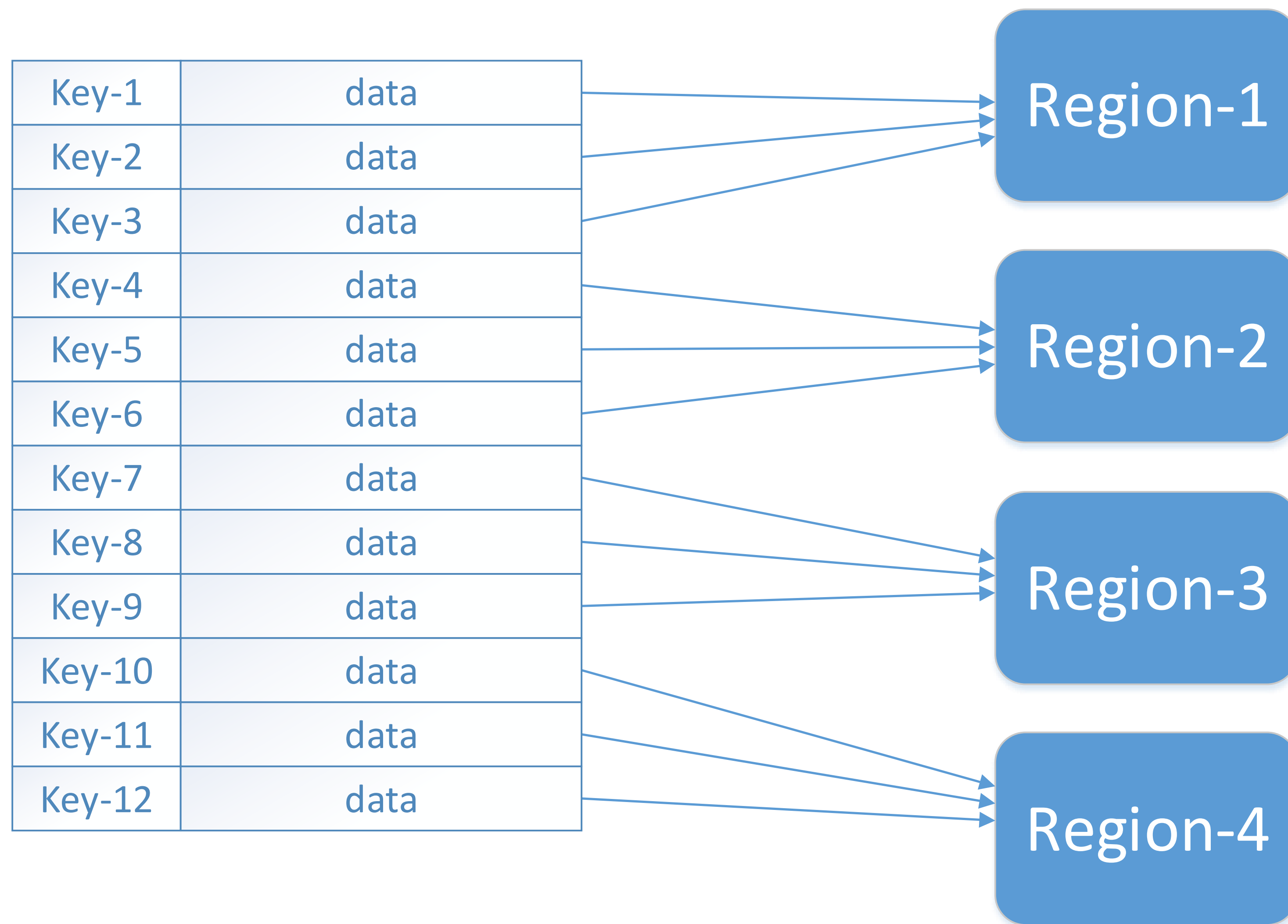
CAS!





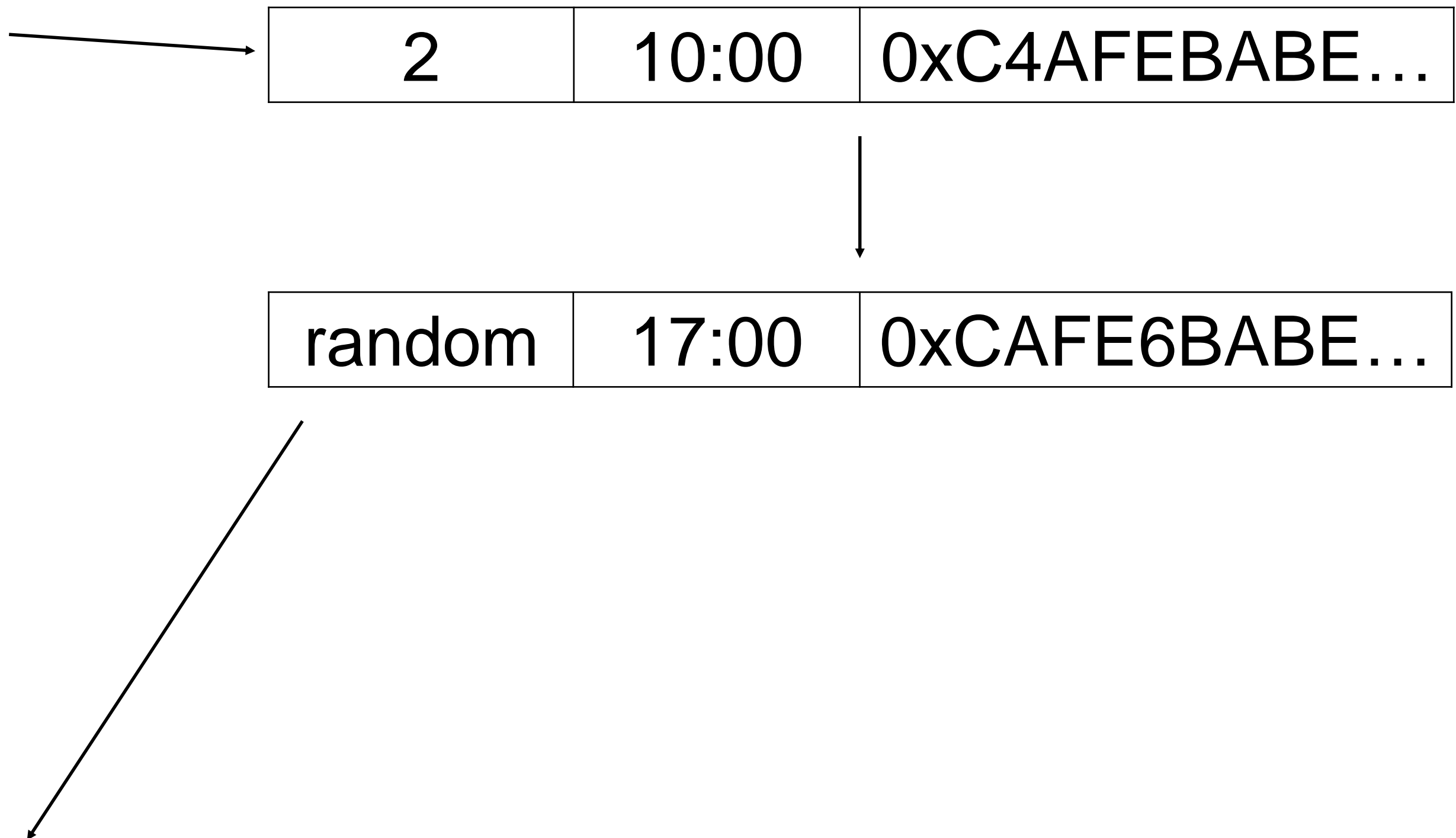








partition	timestamp	UID
5bit	43bit	64bit
1	10:00	0x3CAFEBABE...
1	11:17	0xC4AFEBABE...
1	23:05	0xCA5FEBABE...
2	10:00	0xCAF6EBABE...
2	10:01	0xCAFE3BABE...
2	11:45	0xCAFEB7ABE...
2	12:15	0xCAFEBAB9E...
2	13:10	0xCAFEBABE5...
3	10:00	0x8CAFEBABE...
3	10:03	0xC4AFEBABE...
3	12:17	0xCA6FEBABE...
3	15:22	0xCAF5EBABE...





Keep row key and column name as small as possible



event_type	timestamp	UID
11bit	43bit	64bit





9 21:10	8 20:19	7 20:18	6 20:12	5 19:53	4 12:45	3 11:17	2 10:05	1 10:00
------------	------------	------------	------------	------------	------------	------------	------------	------------



{ login: «Пицот метроф красной тряпки», name: «Иван», score: 45 },  
{ login: «Двести пийсят метроф красной тряпки», name: «Коля», score: 42 }

key = (crc32(login) << 32) | javaHash(login) // 64bit hash key  
Index = UUID

RowKey	Column		Value
Hash	type	index	
64bit	0x1	0x1CAFEBABE...	Пицот метроф красной тряпки
	0x1	0x2CAFEBABE...	Двести пийсят метроф красной тряпки
	0x2	0x1CAFEBABE...	{ name:«Иван», score: 45 }
	0x2	0x2CAFEBABE...	{ name:«Коля», score: 42 }



<http://github.com/ru-fix>

[kasfandiyarov@gmail.com](mailto:kasfandiyarov@gmail.com)



**СПАСИБО ЗА ВНИМАНИЕ**