



«Один раз в год сады цветут»

Разбор семантики exactly-once в Apache Kafka

Ты кто такой? Давай...

Solutions Architect
Developer Advocate

@gamussa in internetz

Ну погнущись, плав...



Тщеславие и Маркетинг ©



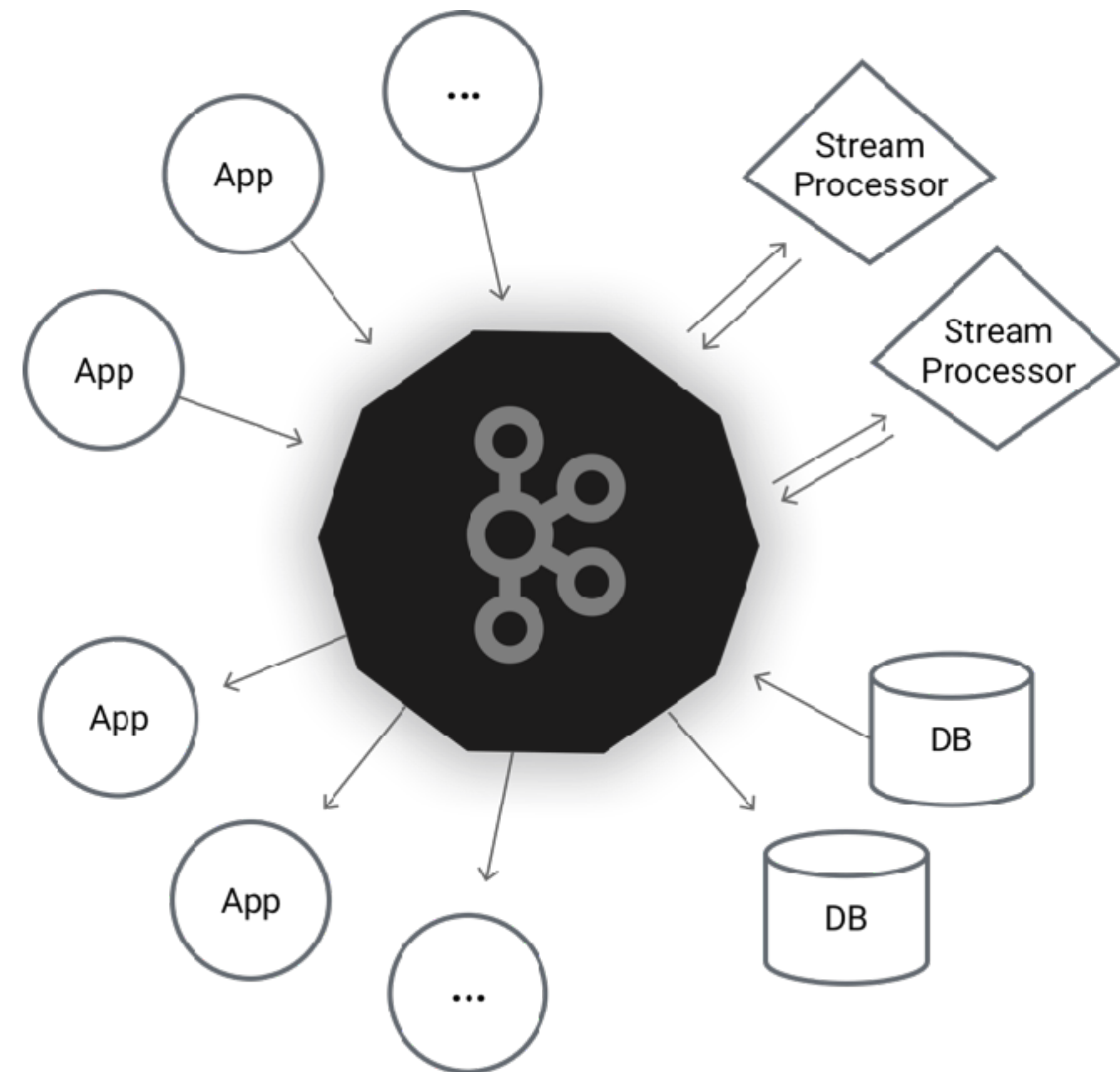


Разбор семантики Exactly-once в Apache Kafka

- Про Kafka по-быстрому
- Ограничения старой семантики
- Углубляемся в Exactly-once семантику
- Чему мы научились и будущие улучшения

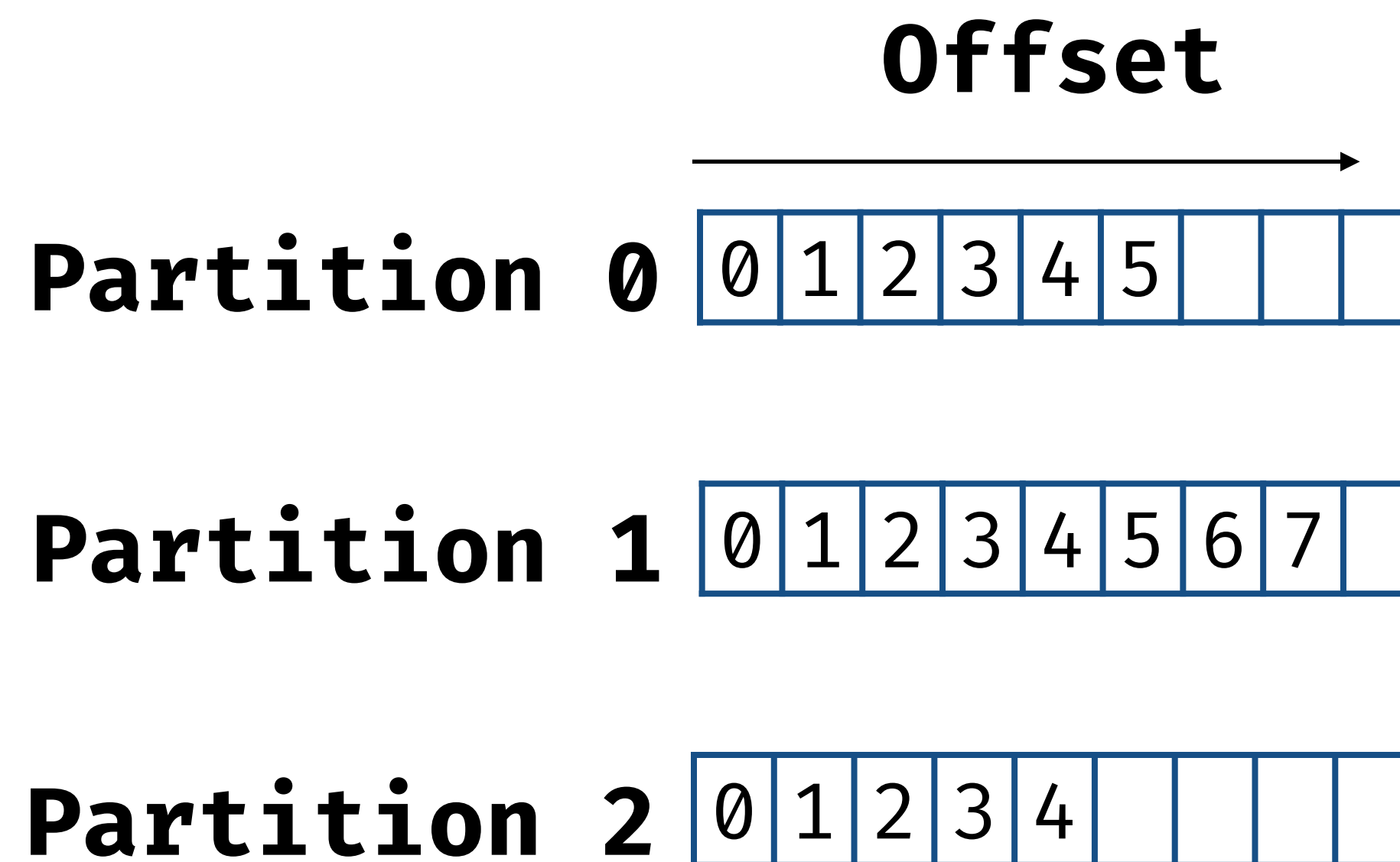
Про Kafka по-быстрому

- Kafka распределенная платформа обработки потоковой информации
- publish
- subscribe
- process
- store

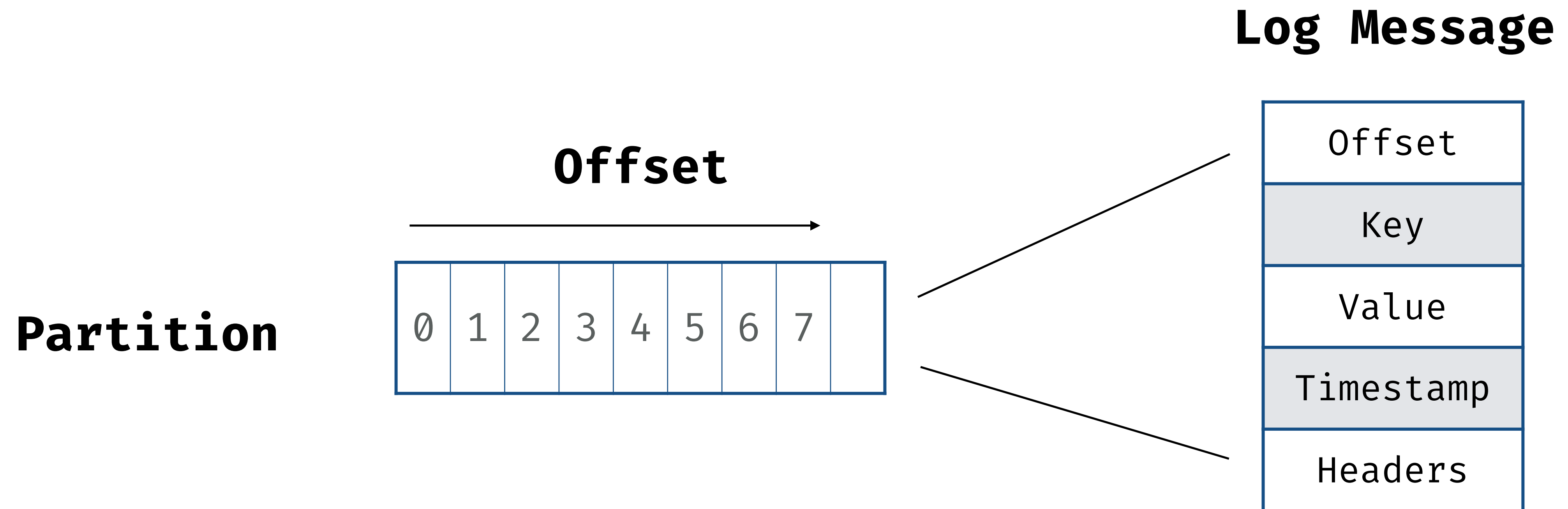


Про Kafka по-быстрому

- Kafka данные организованы в топиках
- Каждый топик состоит из партиций
- Партиция это реплицированный лог
- К сообщениям можно получать доступ по смещению (offset)



Про Kafka по-быстрому



Про Kafka по-быстрому

- Партици реплицируются
- Количество реплик настраивается
- После сбоя, реплики промотаются в лидеры

Broker 1

Partition 0



Partition 1



Broker 2

Partition 0

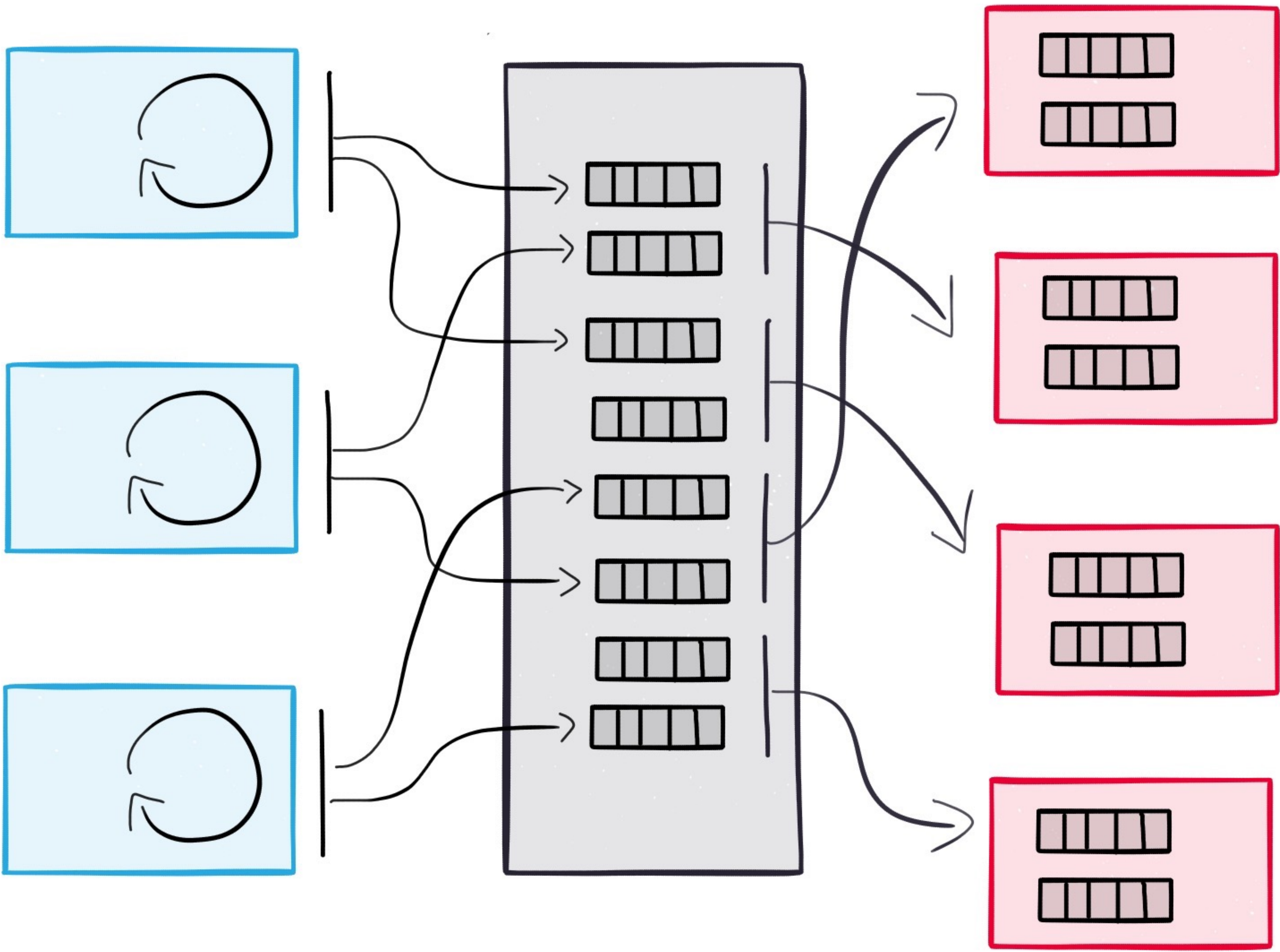


Partition 1



Про Kafka по-быстрому

- Клиенты Kafka
 - Producers записывают данные в партиции
 - Consumers читают из партиции



Kafka Producer

Producers

P1

P2

P3

Broker
1

Broker
2

Partition 0 (Leader)



Partition 1 (Follower)



Partition 0 (Follower)



Partition 1 (Leader)



Producers

P1

P2

P3



Broker
1

Broker
2

Partition 0 (Leader)



Partition 1 (Follower)



Partition 0 (Follower)



Partition 1 (Leader)



Producers

P1

P2

P3



Broker
1



Broker
2

Partition 0 (Leader)



Partition 1 (Follower)



Partition 0 (Follower)



Partition 1 (Leader)

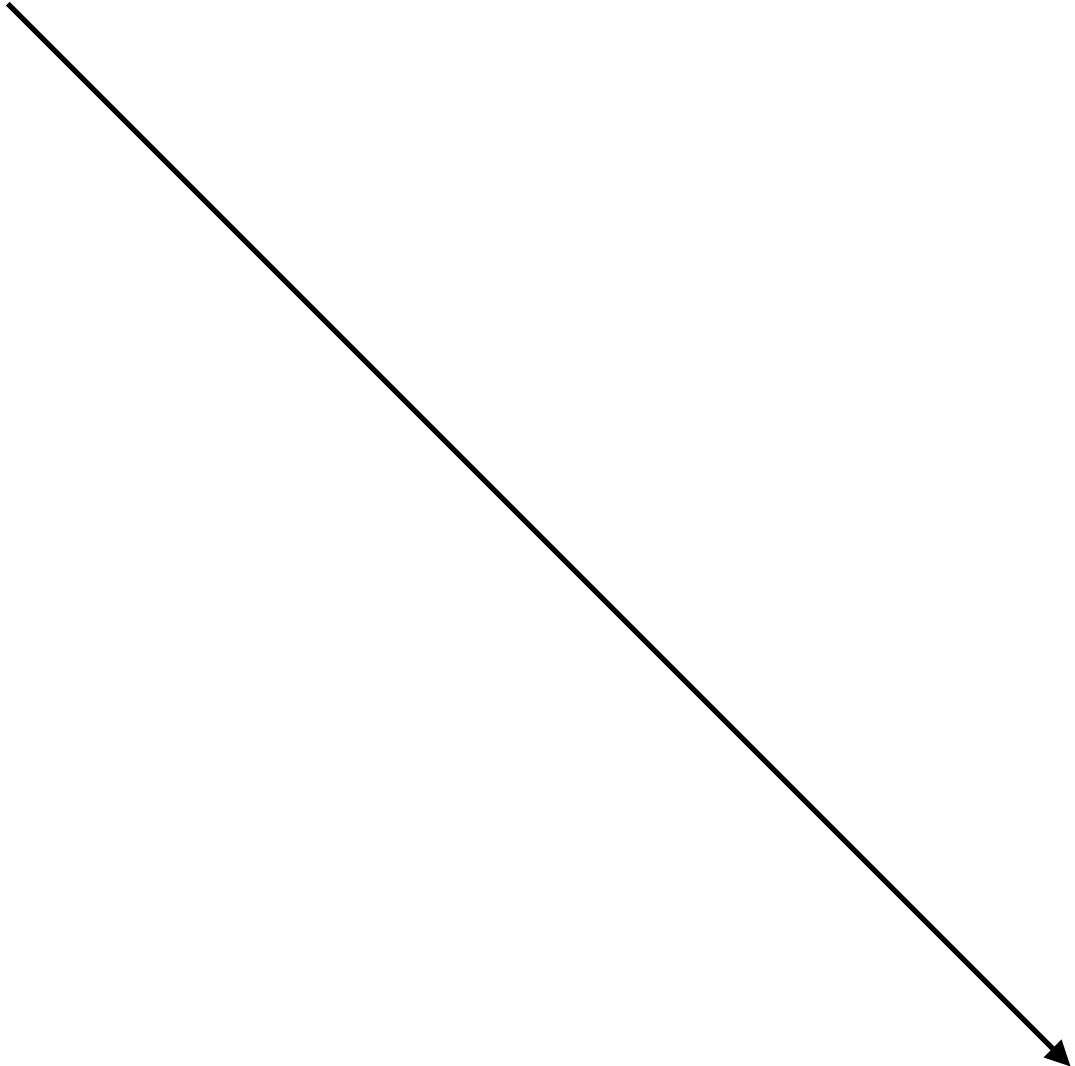


Producers

P1

P2

P3



Broker
1

Broker
2

Partition 0 (Leader)



Partition 1 (Follower)



Partition 0 (Follower)



Partition 1 (Leader)

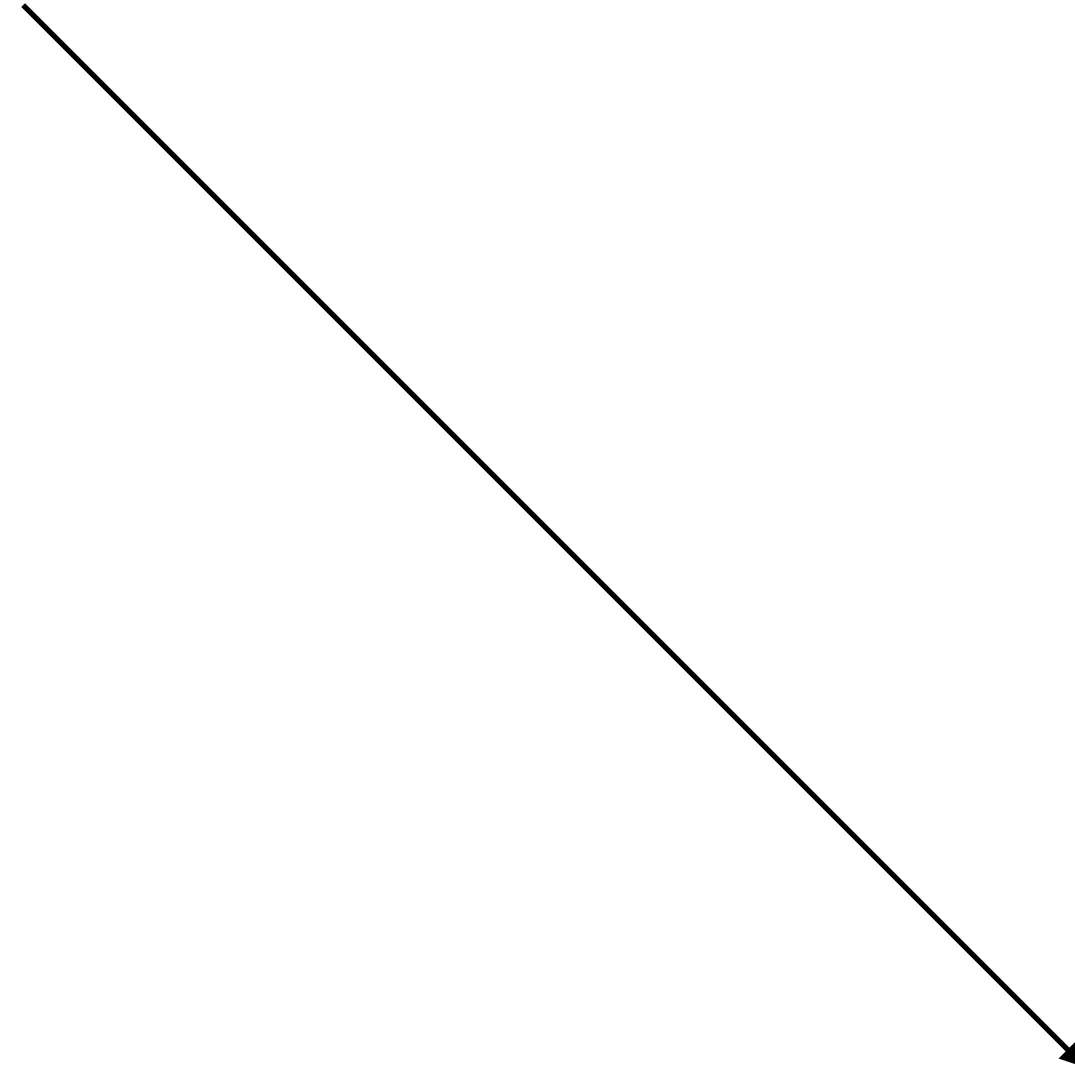


Producers

P1

P2

P3



Broker
1

Broker
2



Partition 0 (Leader)



Partition 1 (Follower)



Partition 0 (Follower)



Partition 1 (Leader)

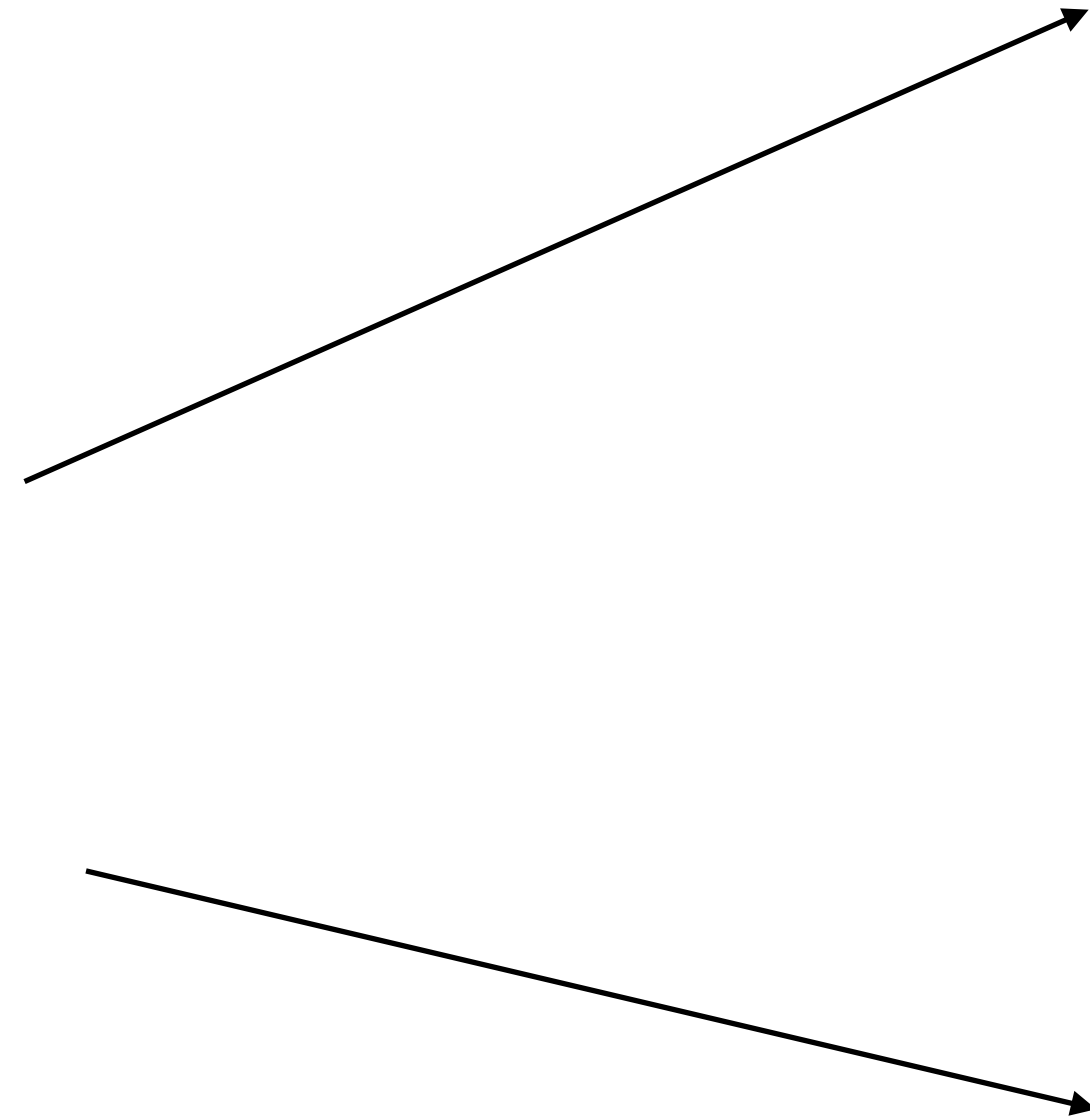


Producers

P1

P2

P3



Broker
1

Broker
2



Partition 0 (Leader)



Partition 1 (Follower)



Partition 0 (Follower)



Partition 1 (Leader)

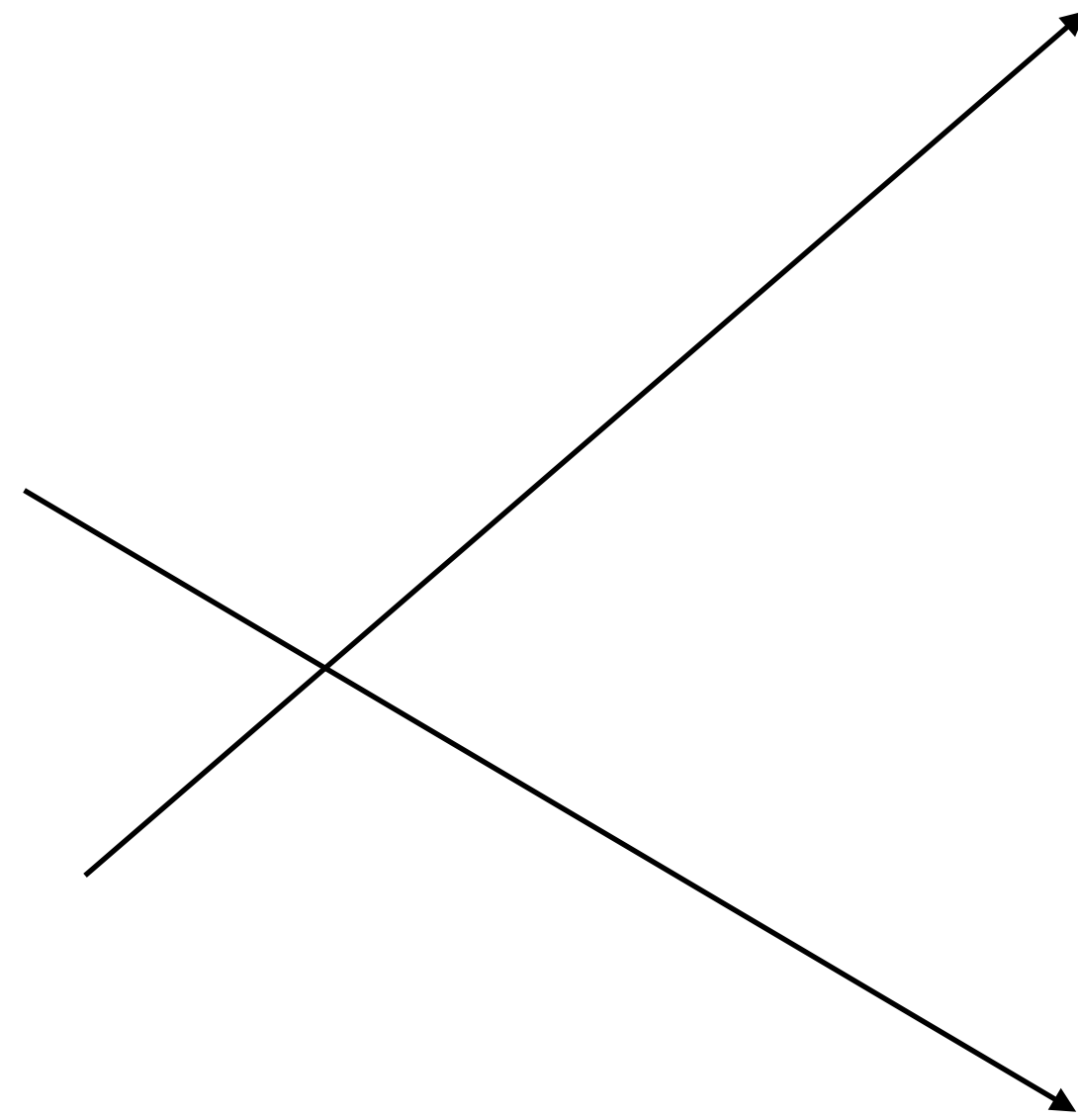


Producers

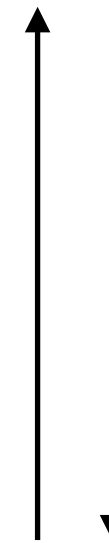
P1

P2

P3

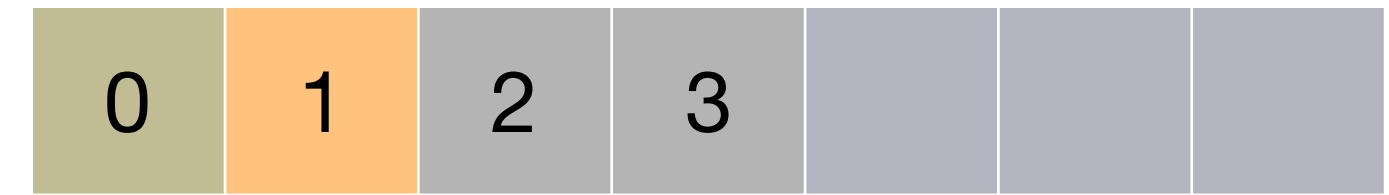


Broker
1



Broker
2

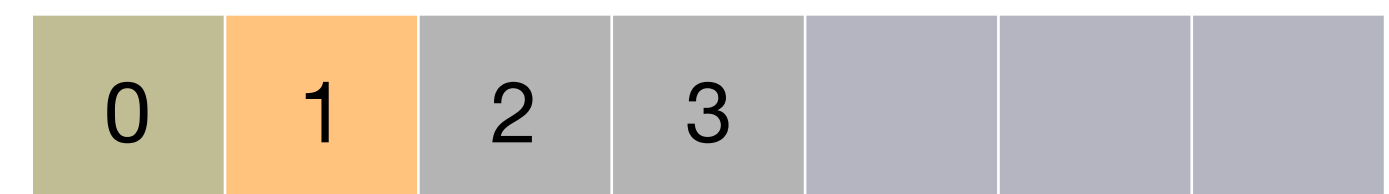
Partition 0 (Leader)



Partition 1 (Follower)



Partition 0 (Follower)



Partition 1 (Leader)



Producers

P1

P2

P3



Broker 2

Partition 0 (Leader)



Partition 1 (Follower)



Partition 0 (Follower)



Partition 1 (Leader)



Producers

P1

P2

P3



Partition 0 (Leader)



Partition 1 (Follower)



Partition 0 (Leader)



Partition 1 (Leader)

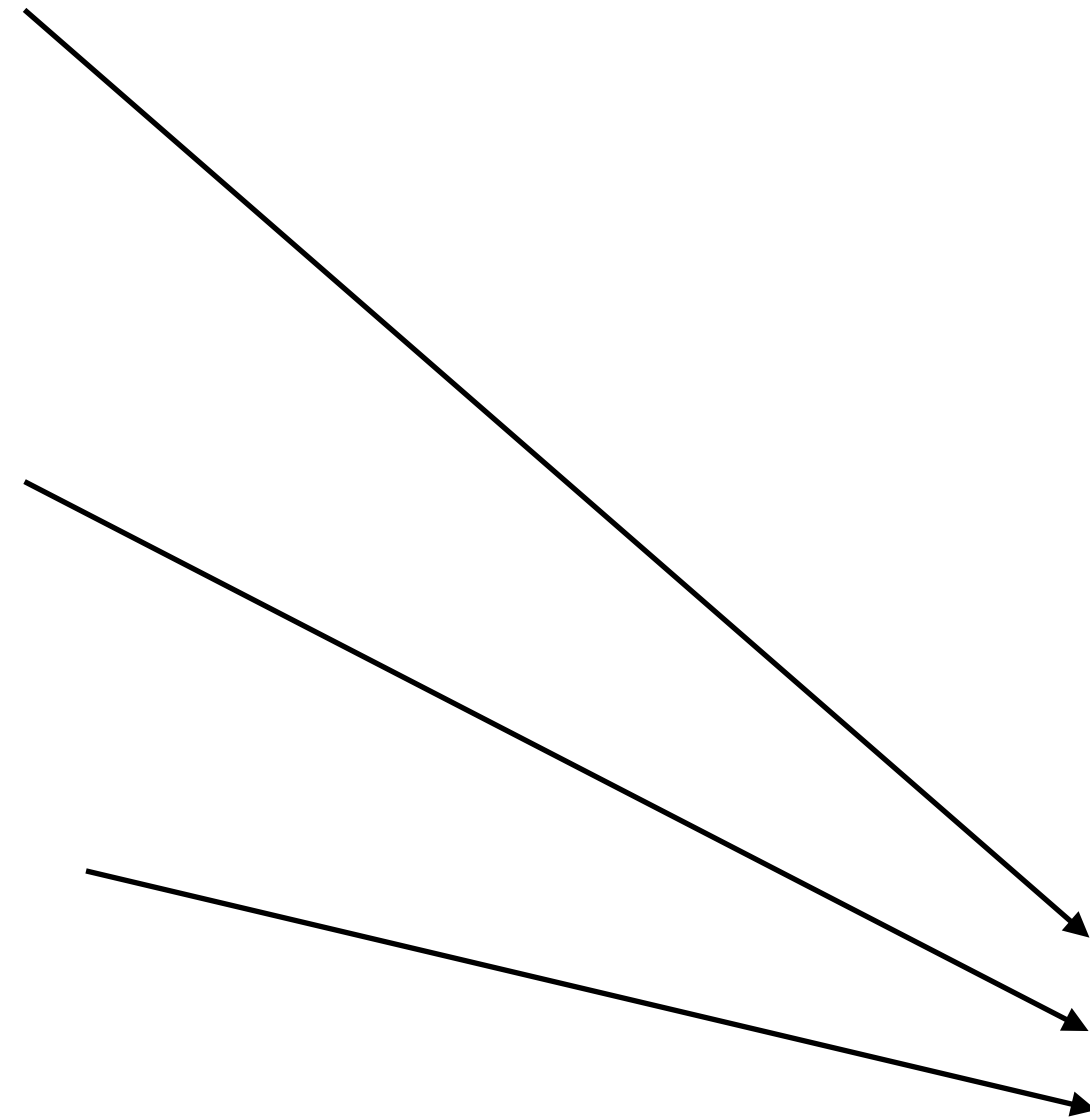


Producers

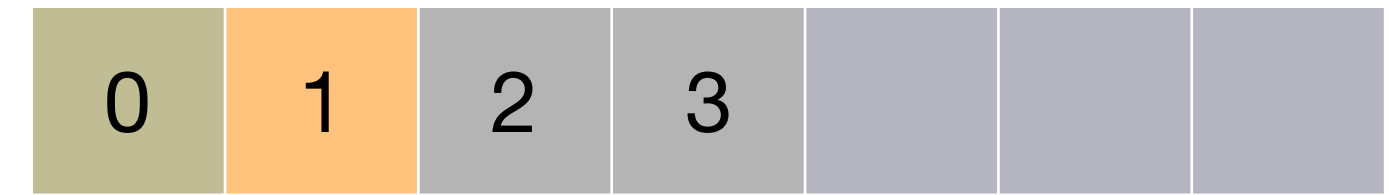
P1

P2

P3



Partition 0 (Leader)



Partition 1 (Follower)



Partition 0 (Leader)



Partition 1 (Leader)



Producers

P1

P2

P3

Broker
1

Broker
2

Partition 0 (Follower)



Partition 1 (Follower)



Partition 0 (Leader)



Partition 1 (Leader)



Producers

P1

P2

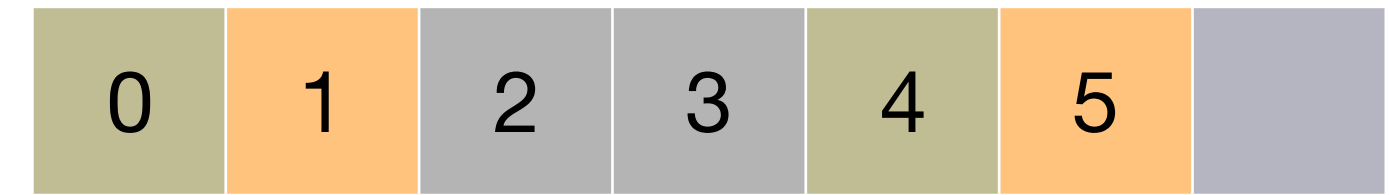
P3

Broker
1

Broker
2



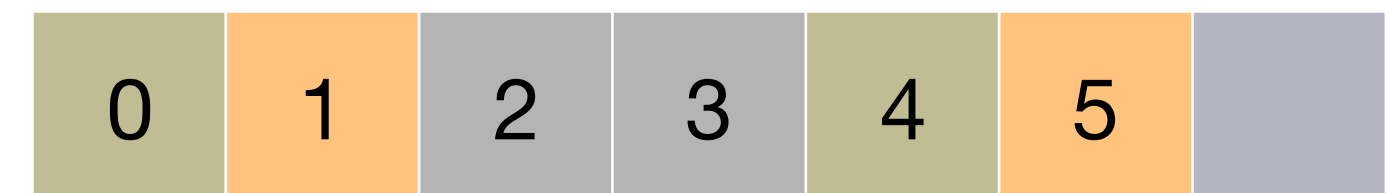
Partition 0 (Follower)



Partition 1 (Follower)



Partition 0 (Leader)



Partition 1 (Leader)



Producers

P1

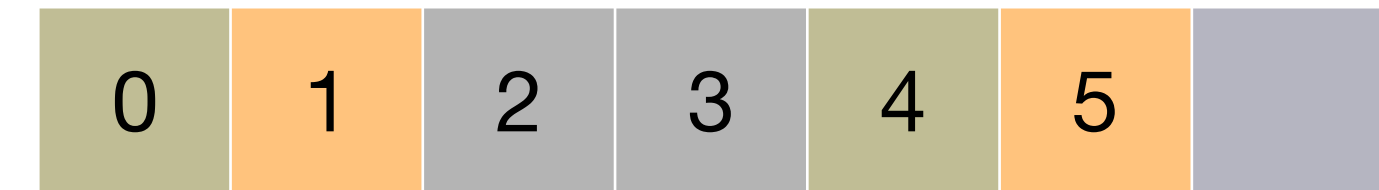
P2

P3

Broker
1

Broker
2

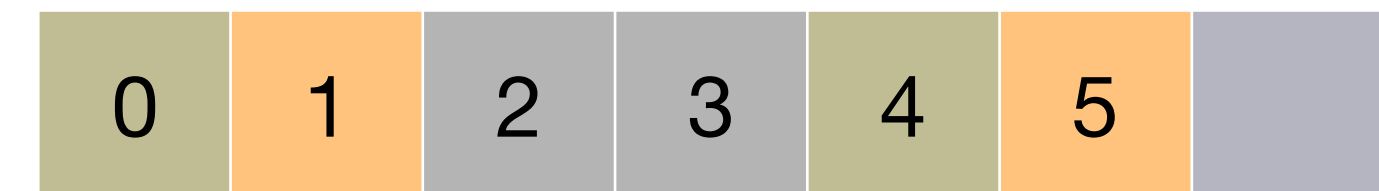
Partition 0 (Leader)



Partition 1 (Follower)



Partition 0 (Follower)



Partition 1 (Leader)

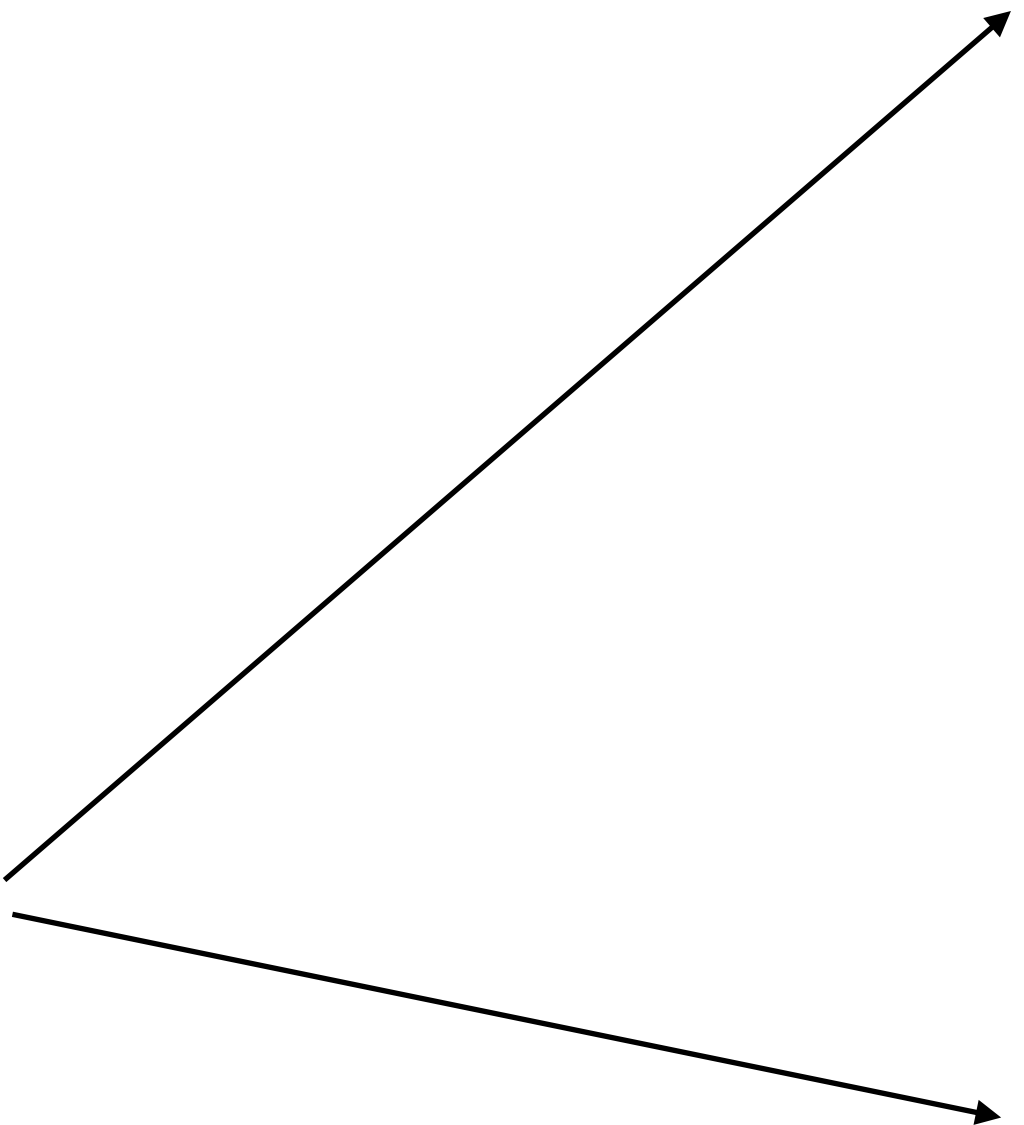


Producers

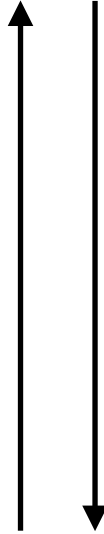
P1

P2

P3



Broker
1



Broker
2

Partition 0 (Leader)



Partition 1 (Follower)



Partition 0 (Follower)



Partition 1 (Leader)



Протокол прогюсера

- Для записи в партиции
- Ask означает успешный append

Client

Broker

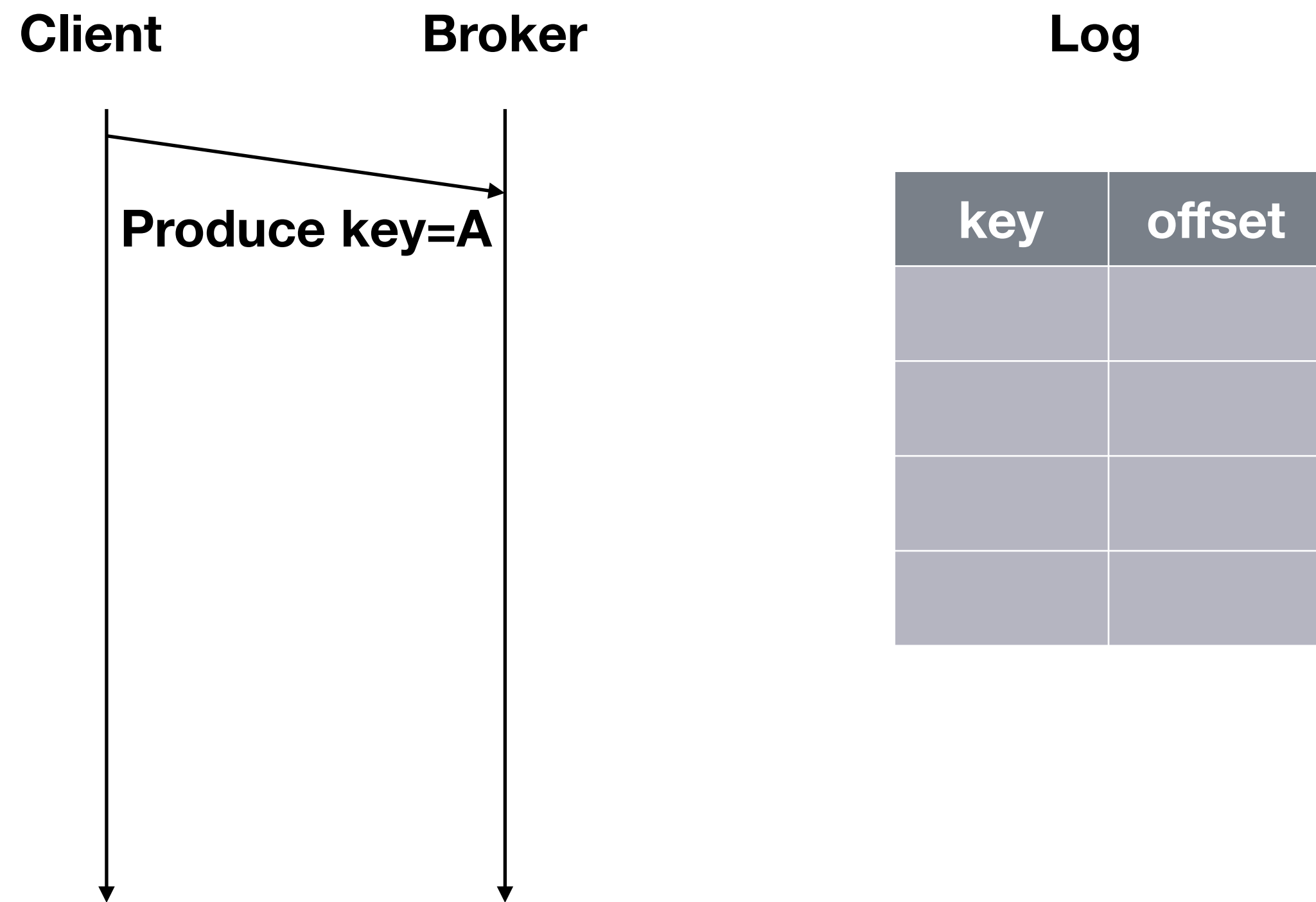
Log



key	offset

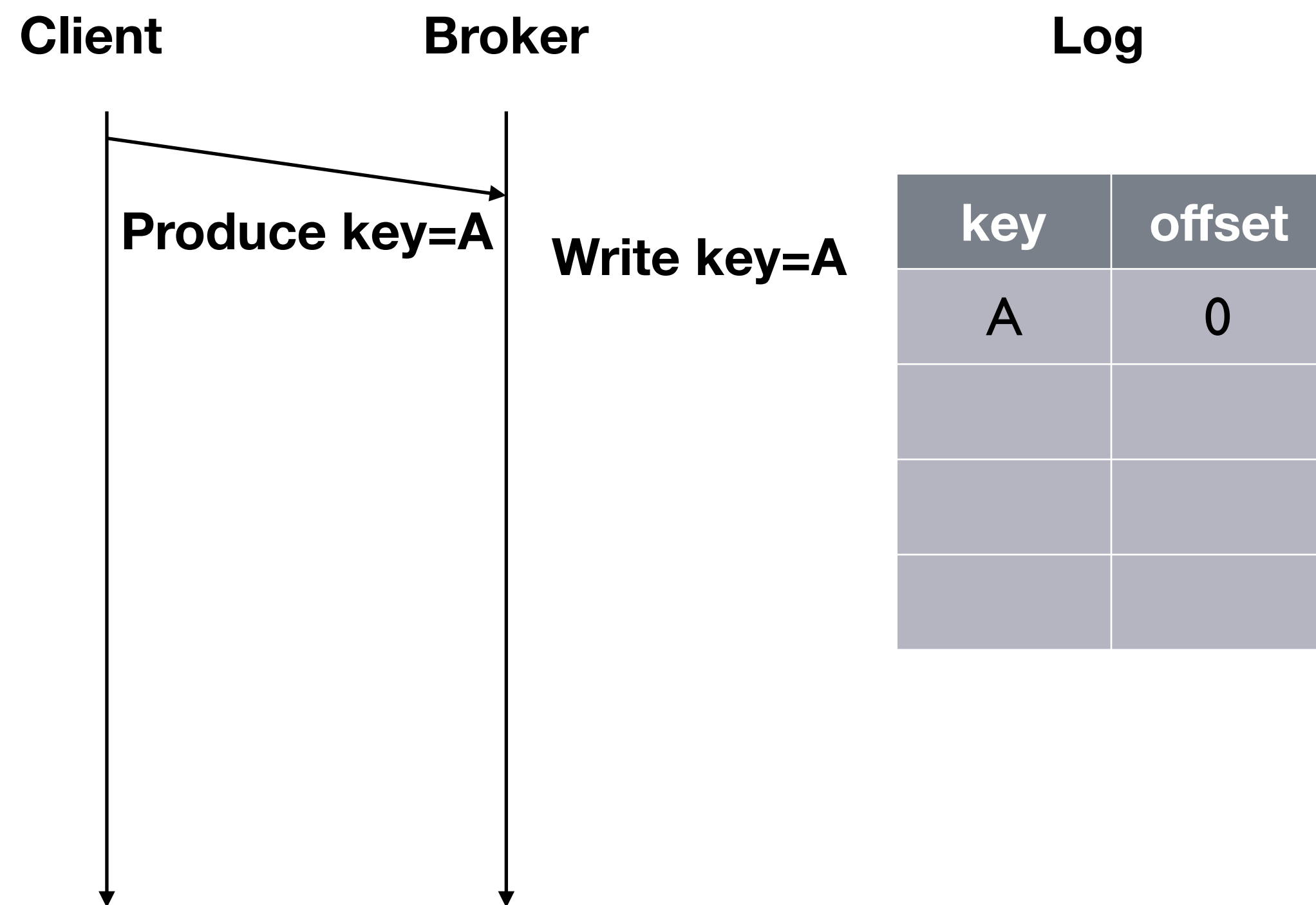
Протокол прогюсера

- Для записи в партиции
- Ask означает успешный append



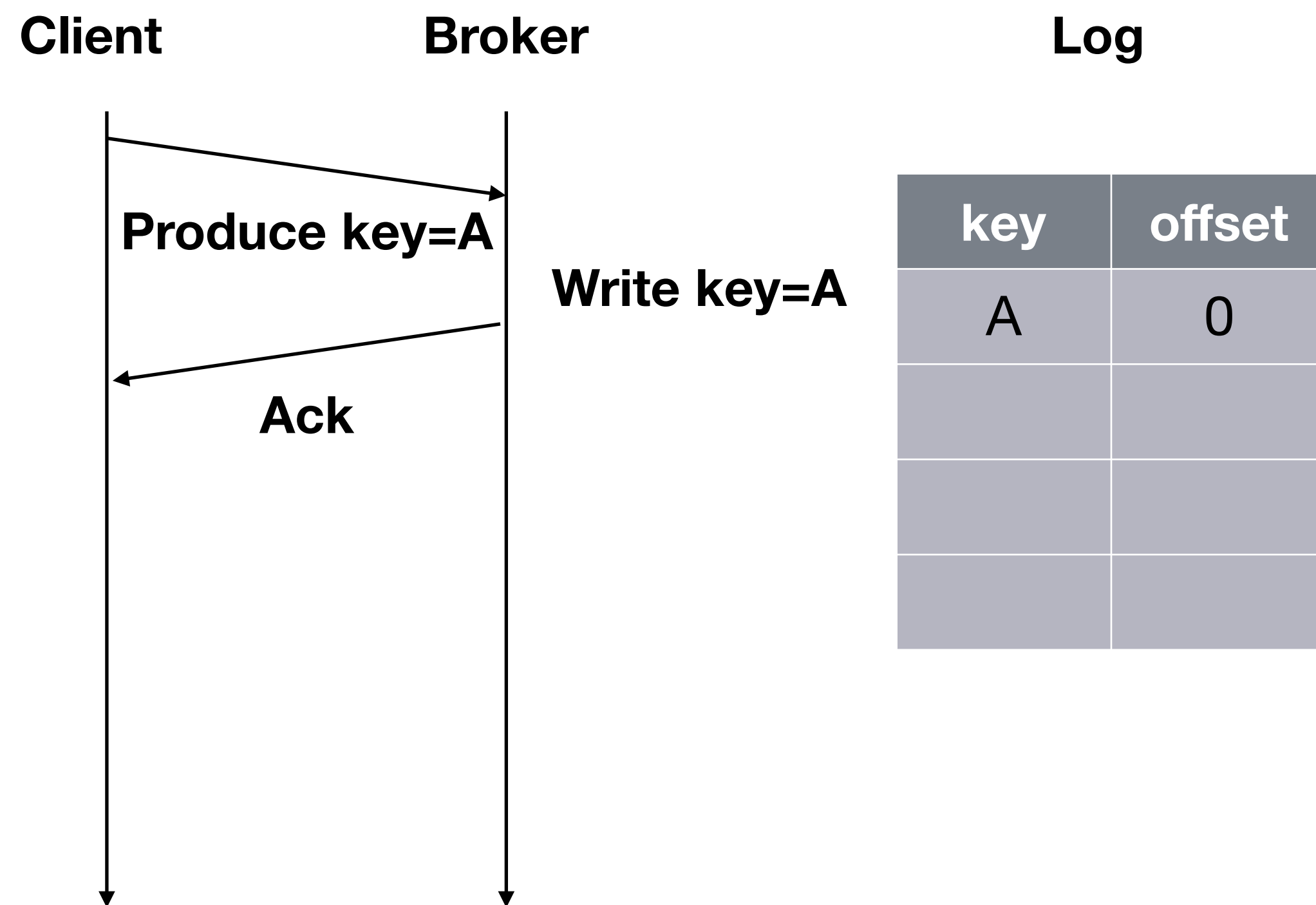
Протокол прогюсера

- Для записи в партиции
- Ask означает успешный append



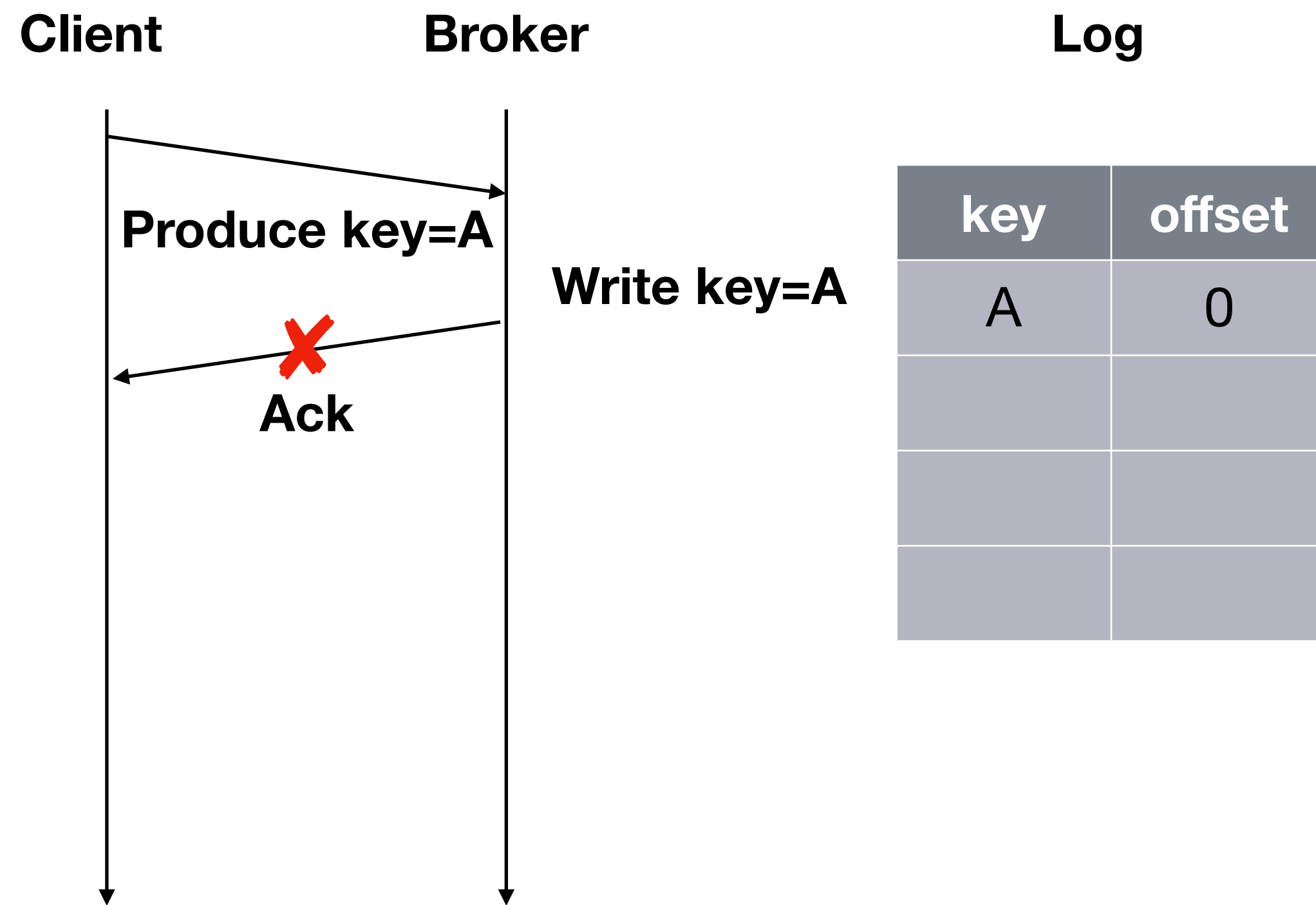
Протокол прогюсера

- Для записи в партиции
- Ack означает успешный append



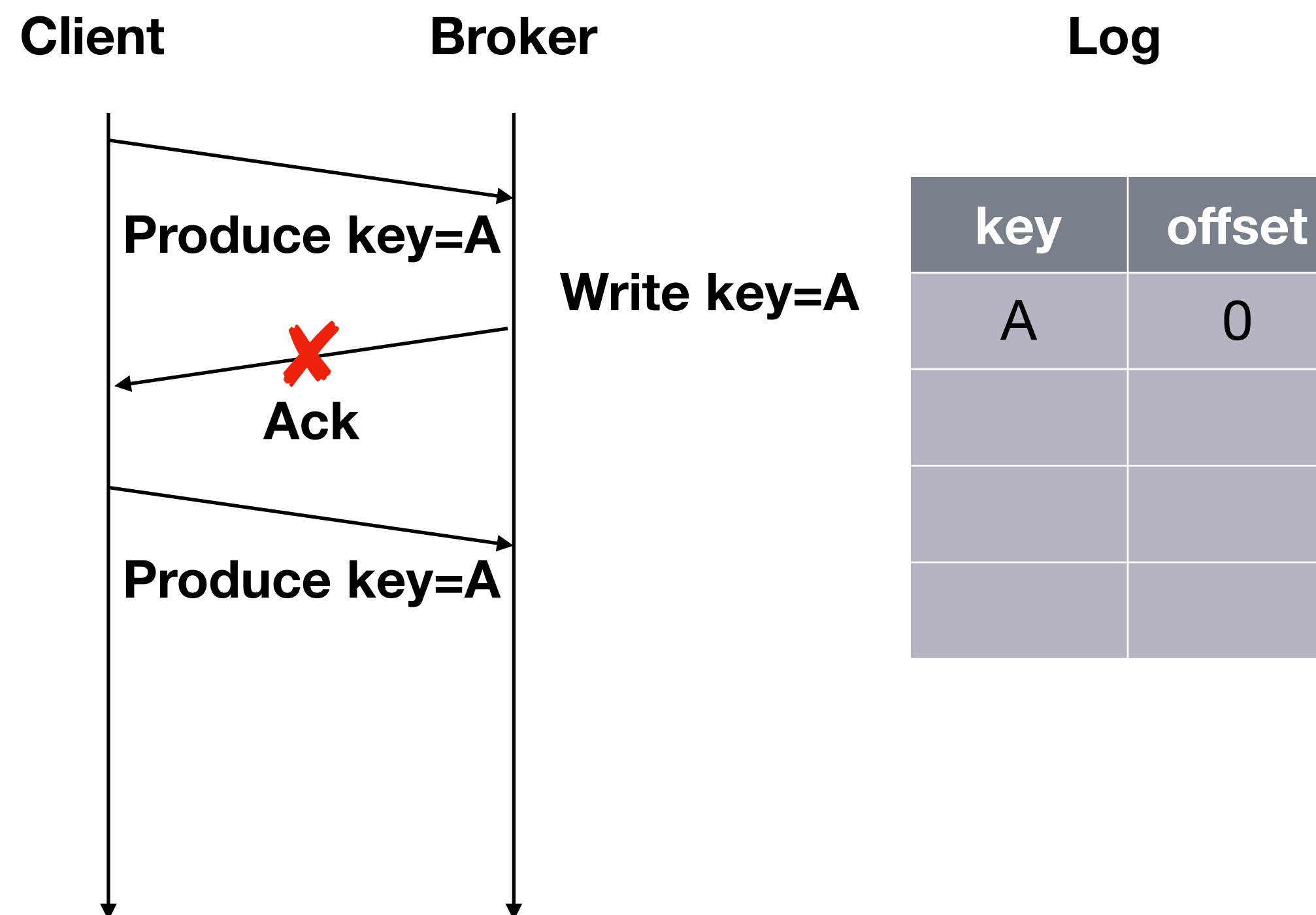
Протокол прогюсера

- Для записи в партиции
- Ack означает успешный append



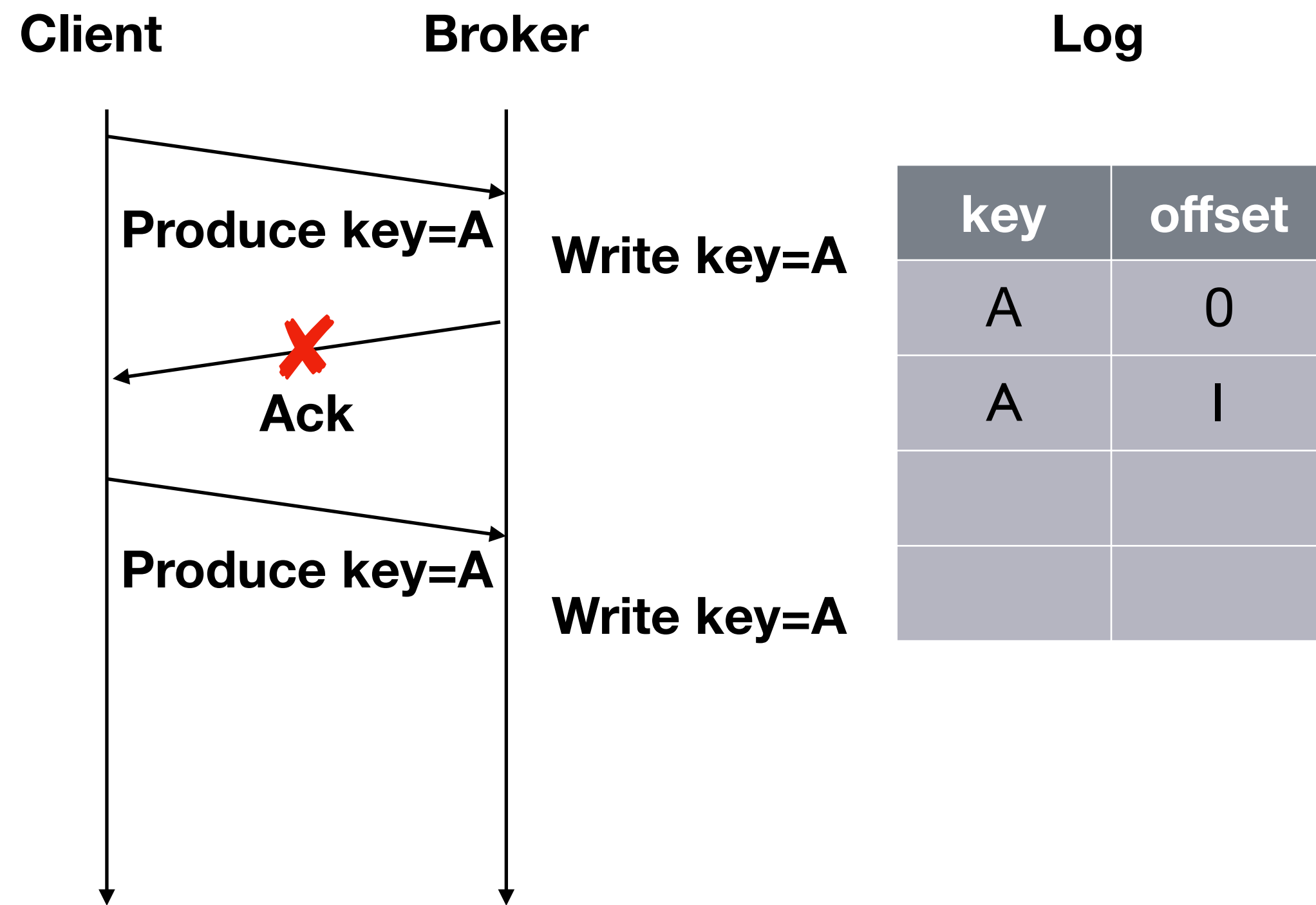
Протокол прогюсера

- Для записи в партиции
- Ack означает успешный append



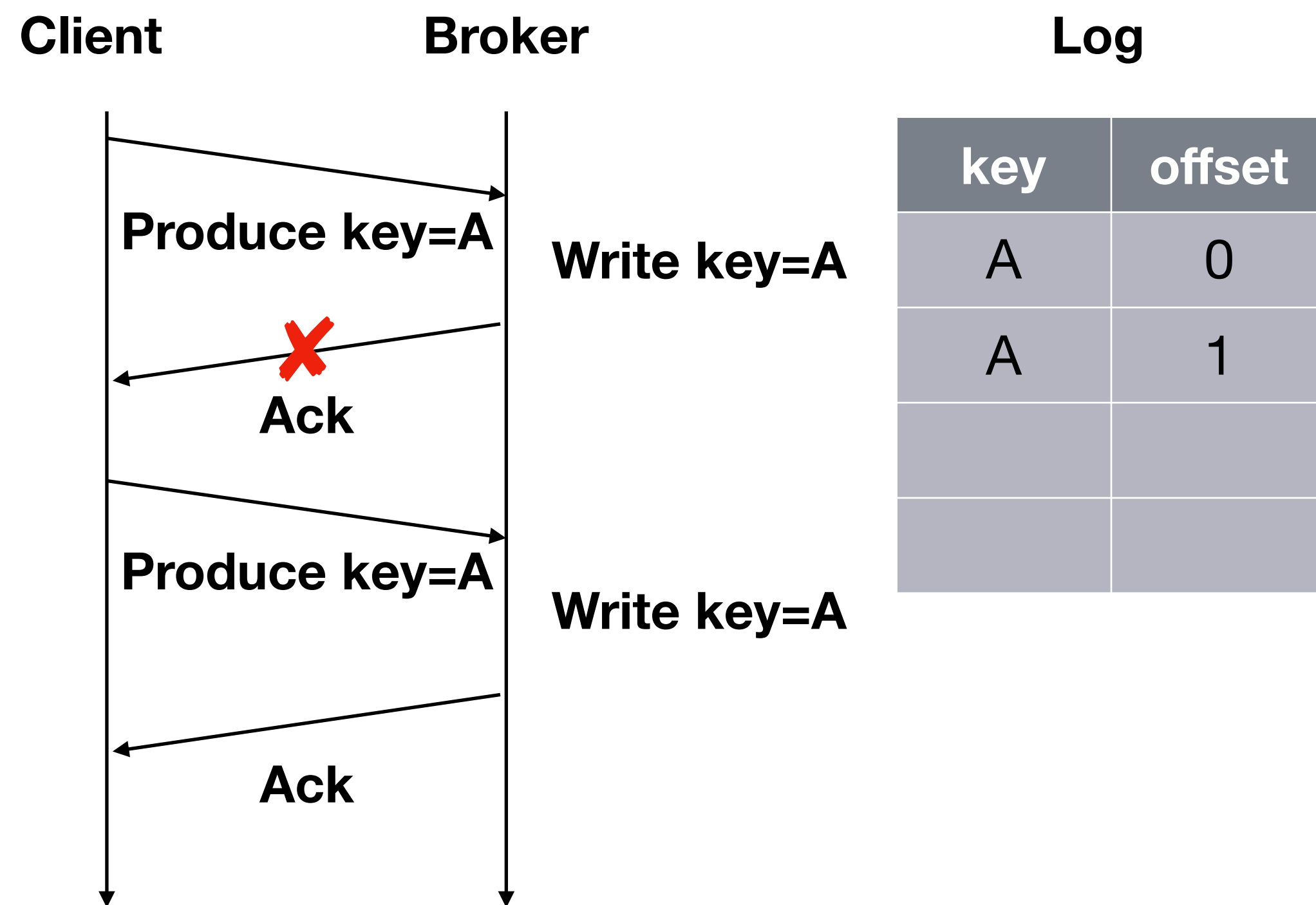
Протокол проглюсера

- Для записи в партиции
- Ack означает успешный append



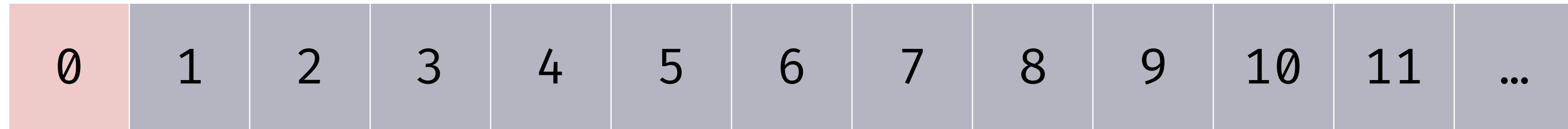
Протокол проглюсера

- Для записи в партиции
- Ack означает успешный append

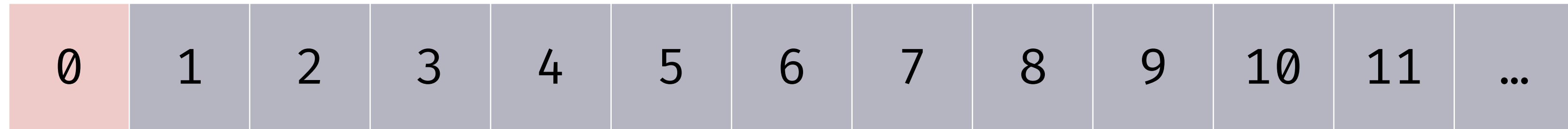


Kafka Consumer

Partition 0



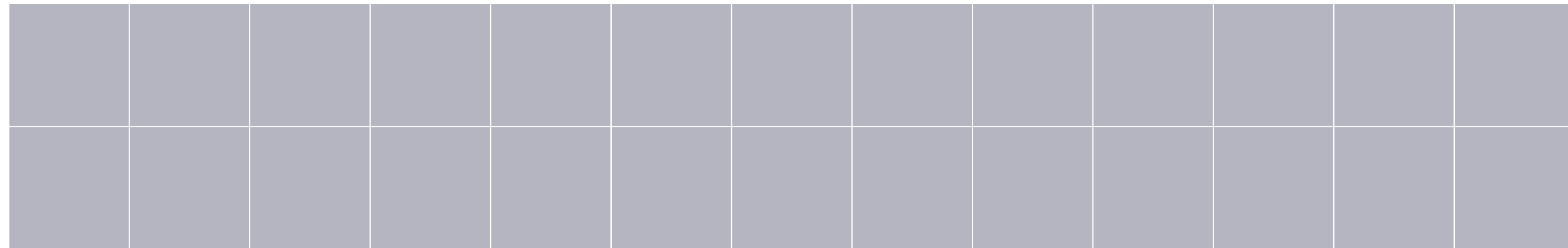
Partition 1



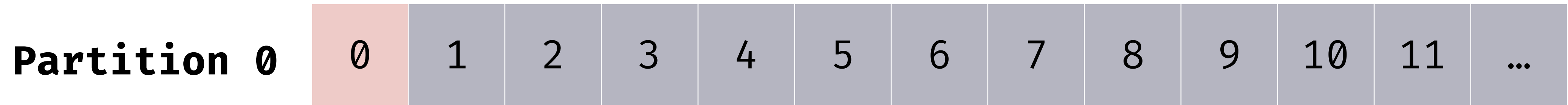
Consumed Partitions

Partition

Offset



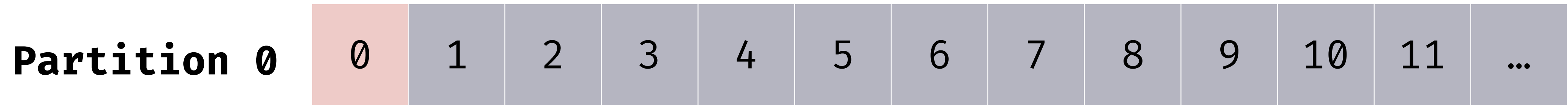
Committed Offsets



Consumed Partitions

Partition												
Offset												

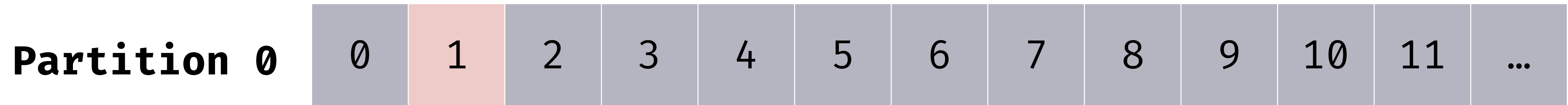
Committed Offsets



Consumed Partitions

Partition												
Offset												

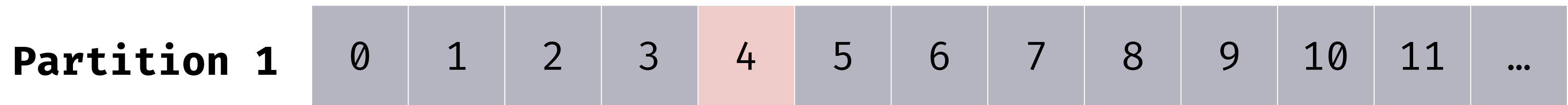
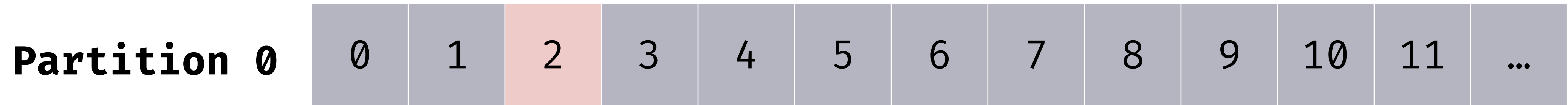
Committed Offsets



Consumed Partitions

Partition												
Offset												

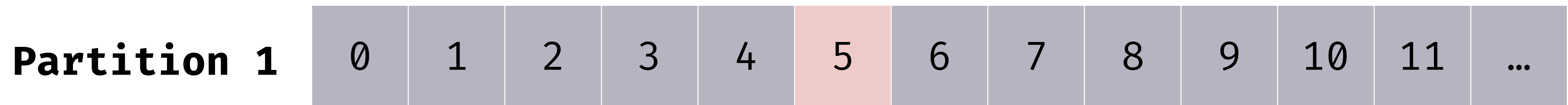
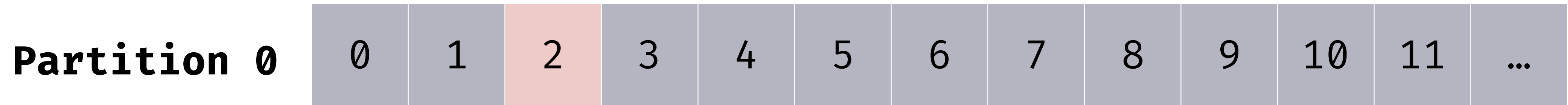
Committed Offsets



Consumed Partitions

Partition												
Offset												

Committed Offsets



Consumed Partitions

Partition												
Offset												

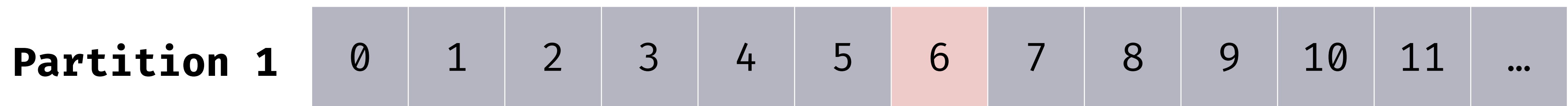
Committed Offsets



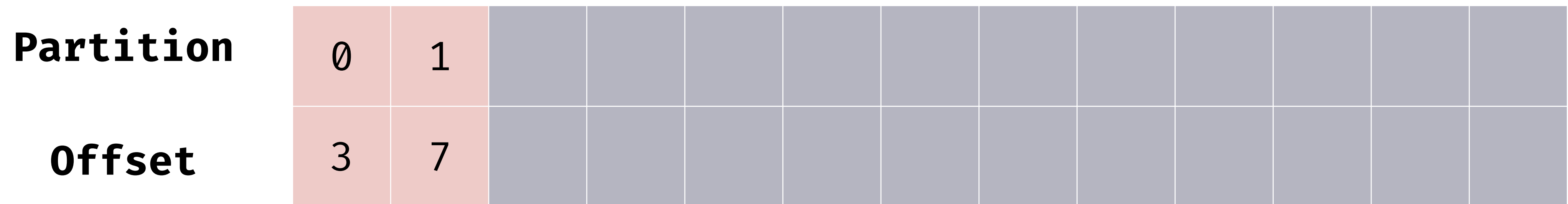
Consumed Partitions

Partition												
Offset												

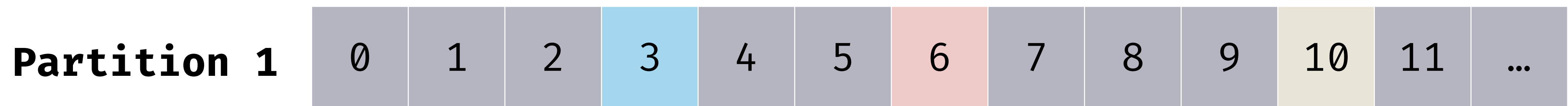
Committed Offsets



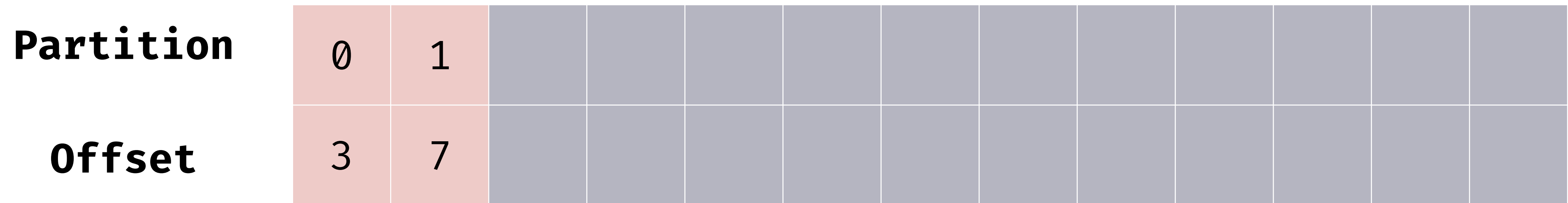
Consumed Partitions



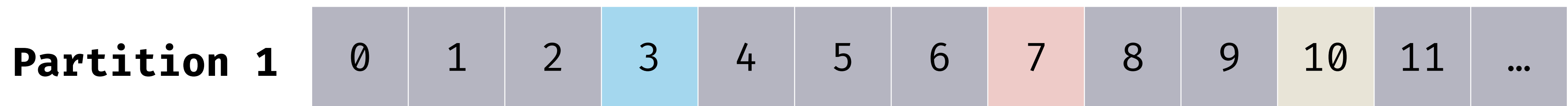
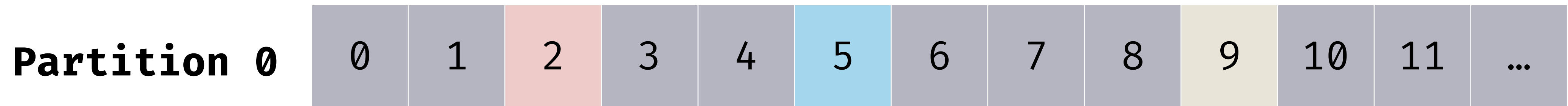
Committed Offsets



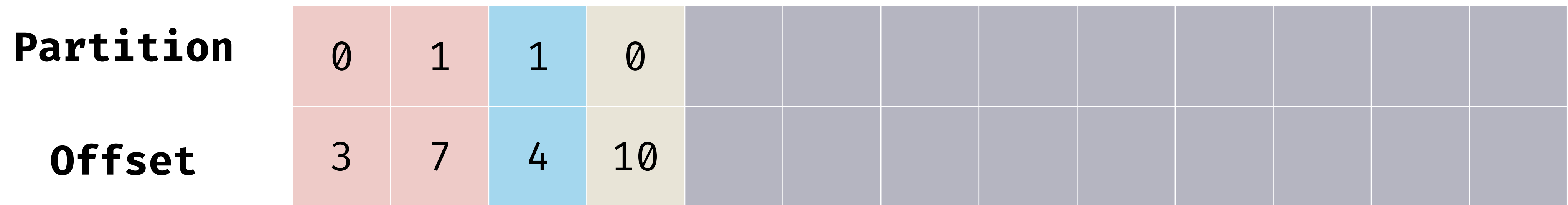
Consumed Partitions



Committed Offsets



Consumed Partitions



Committed Offsets

Partition 0

0	1	2	3	4	5	6	7	8	9	10	11	...
---	---	---	---	---	---	---	---	---	---	----	----	-----

Partition 1

0	1	2	3	4	5	6	7	8	9	10	11	...
---	---	---	---	---	---	---	---	---	---	----	----	-----

Consumed Partitions

Partition	0	1	1	0	1	1	0						
Offset	3	7	4	10	8	12	6						

Committed Offsets

Partition 0

0	1	2	3	4	5	6	7	8	9	10	11	...
---	---	---	---	---	---	---	---	---	---	----	----	-----

Partition 1

0	1	2	3	4	5	6	7	8	9	10	11	...
---	---	---	---	---	---	---	---	---	---	----	----	-----

Consumed Partitions

Partition	0	1	1	0	1	1	0	0	1	1			
Offset	3	7	4	10	8	12	6	3	9	6			

Committed Offsets

Протокол консьюмера

- Fetch API: Взять диапазон записей по смещению (offset)
- OffsetCommit API: Читать/записывать смещения (offset)

Client

Broker

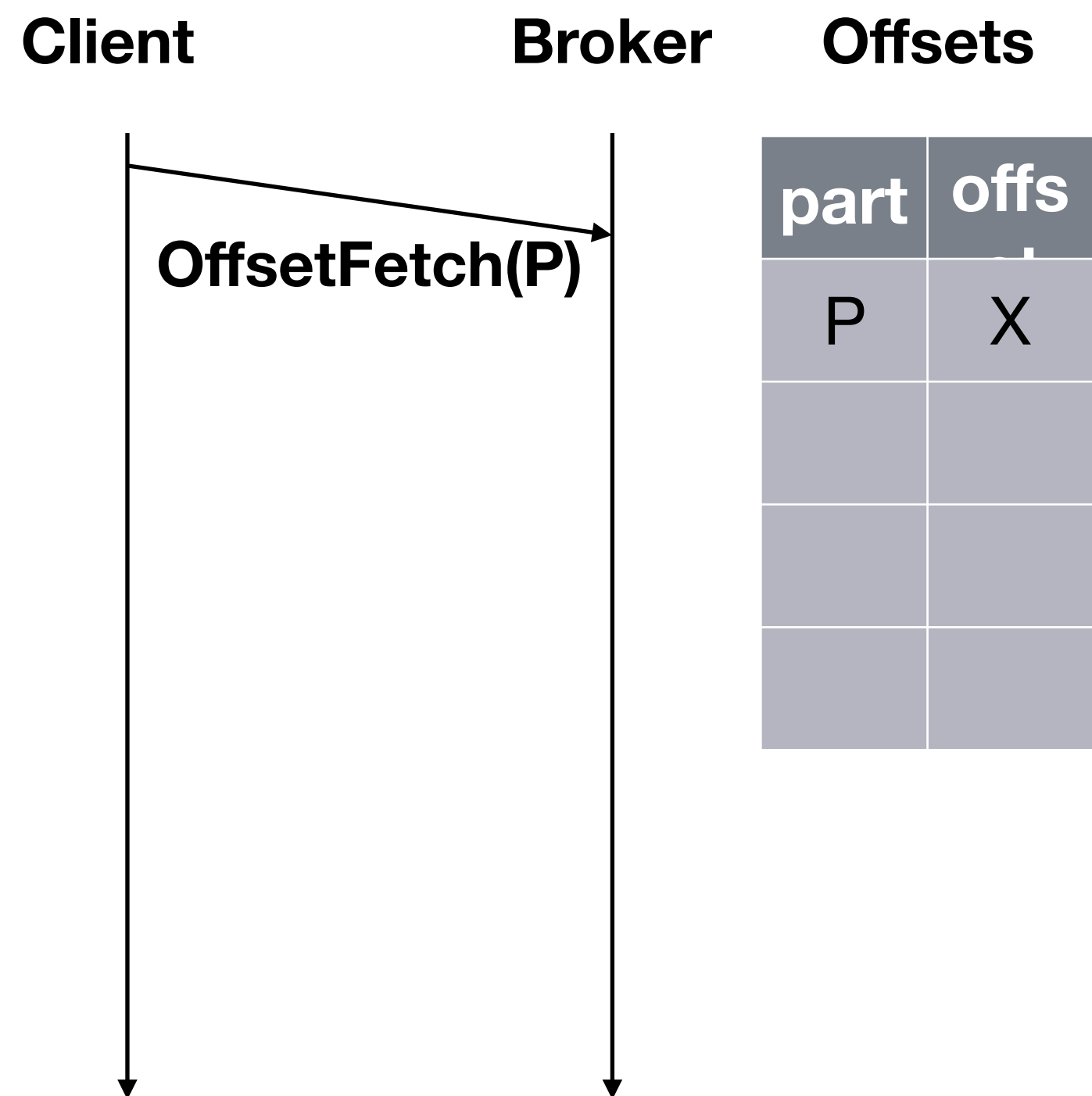
Offsets

part	offs
P	X



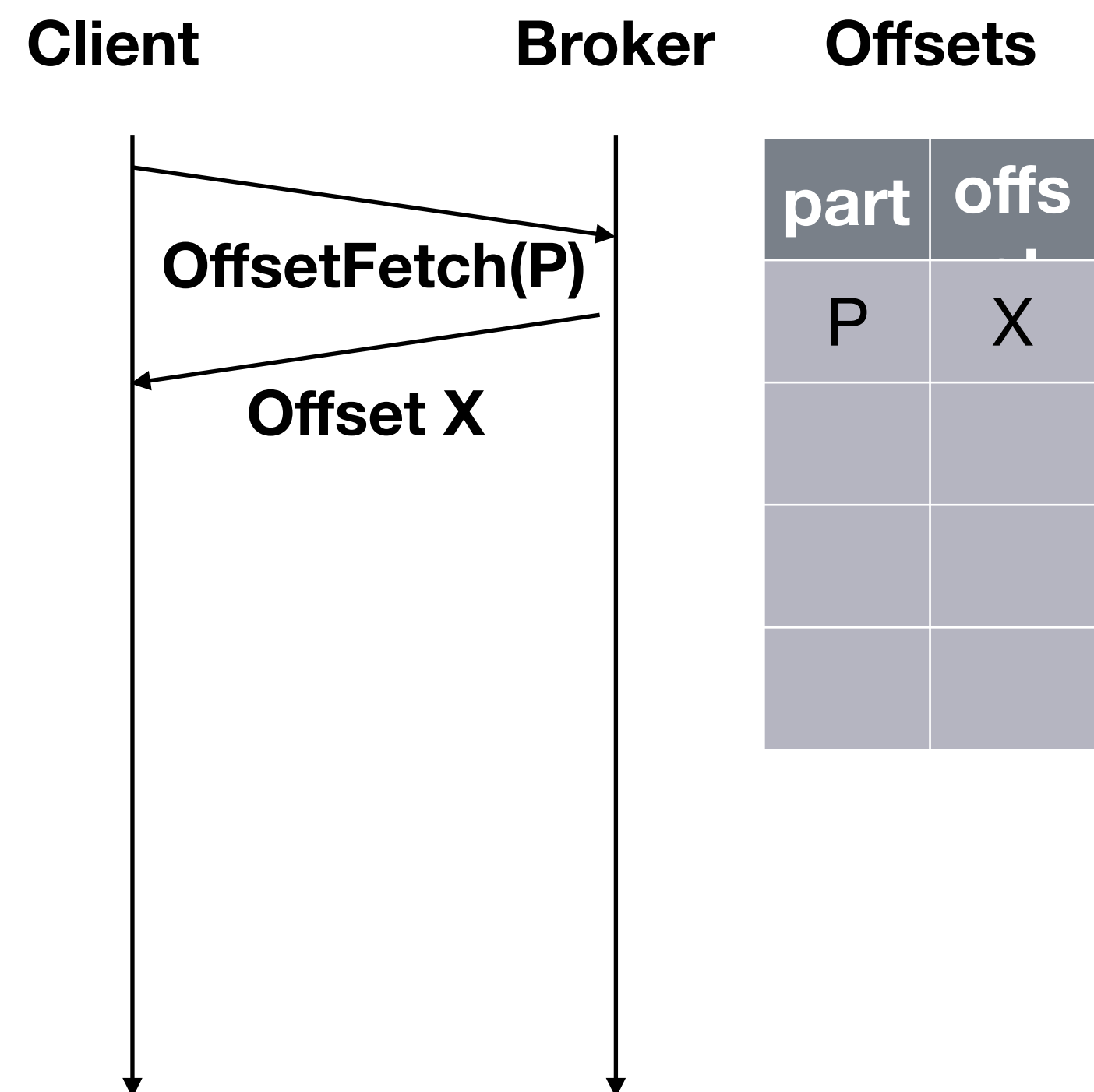
Протокол консьюмера

- Fetch API: Взять диапазон записей по смещению (offset)
- OffsetCommit API: Читать/записывать смещения (offset)



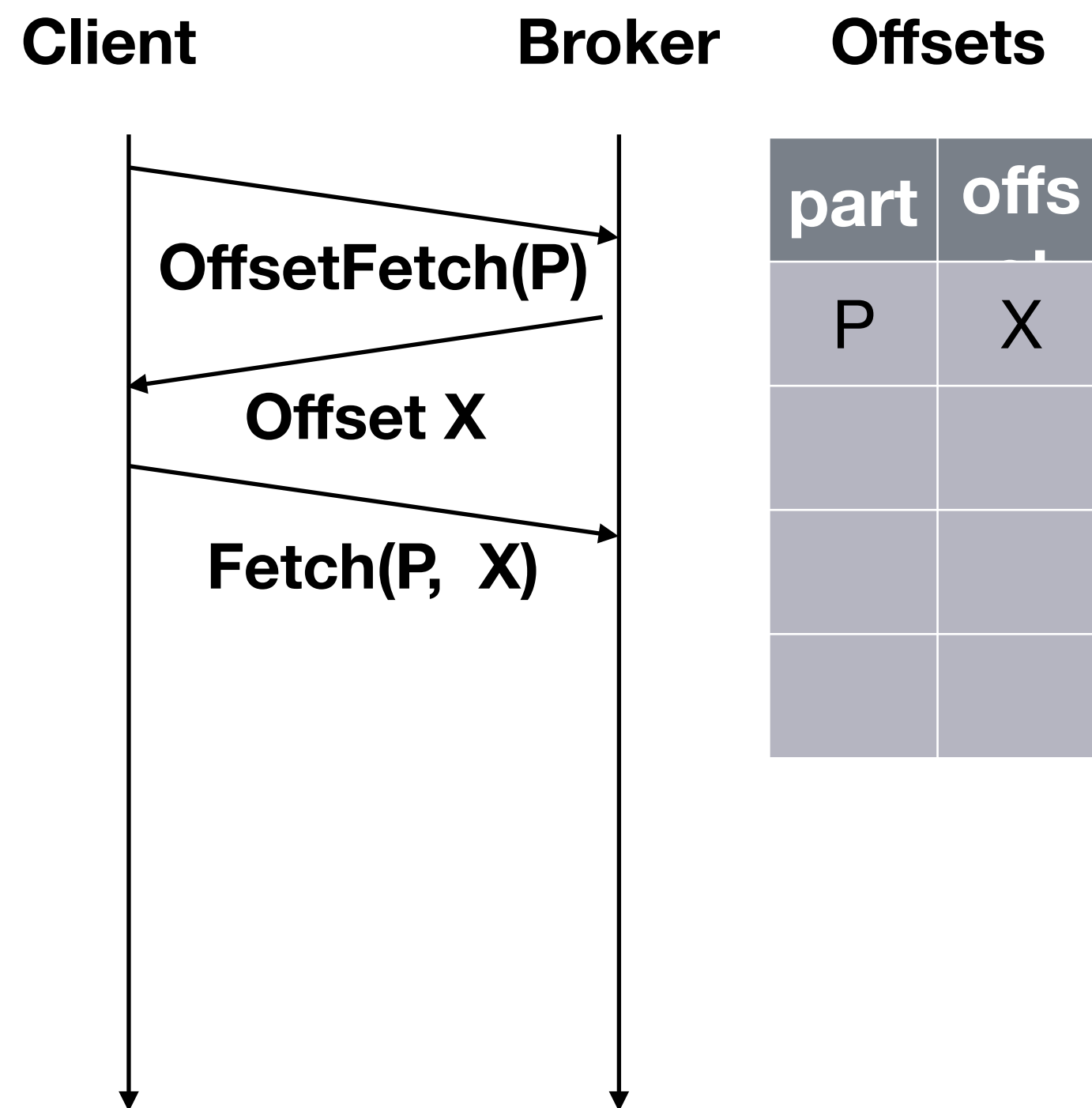
Протокол консьюмера

- Fetch API: Взять диапазон записей по смещению (offset)
- OffsetCommit API: Читать/записывать смещения (offset)



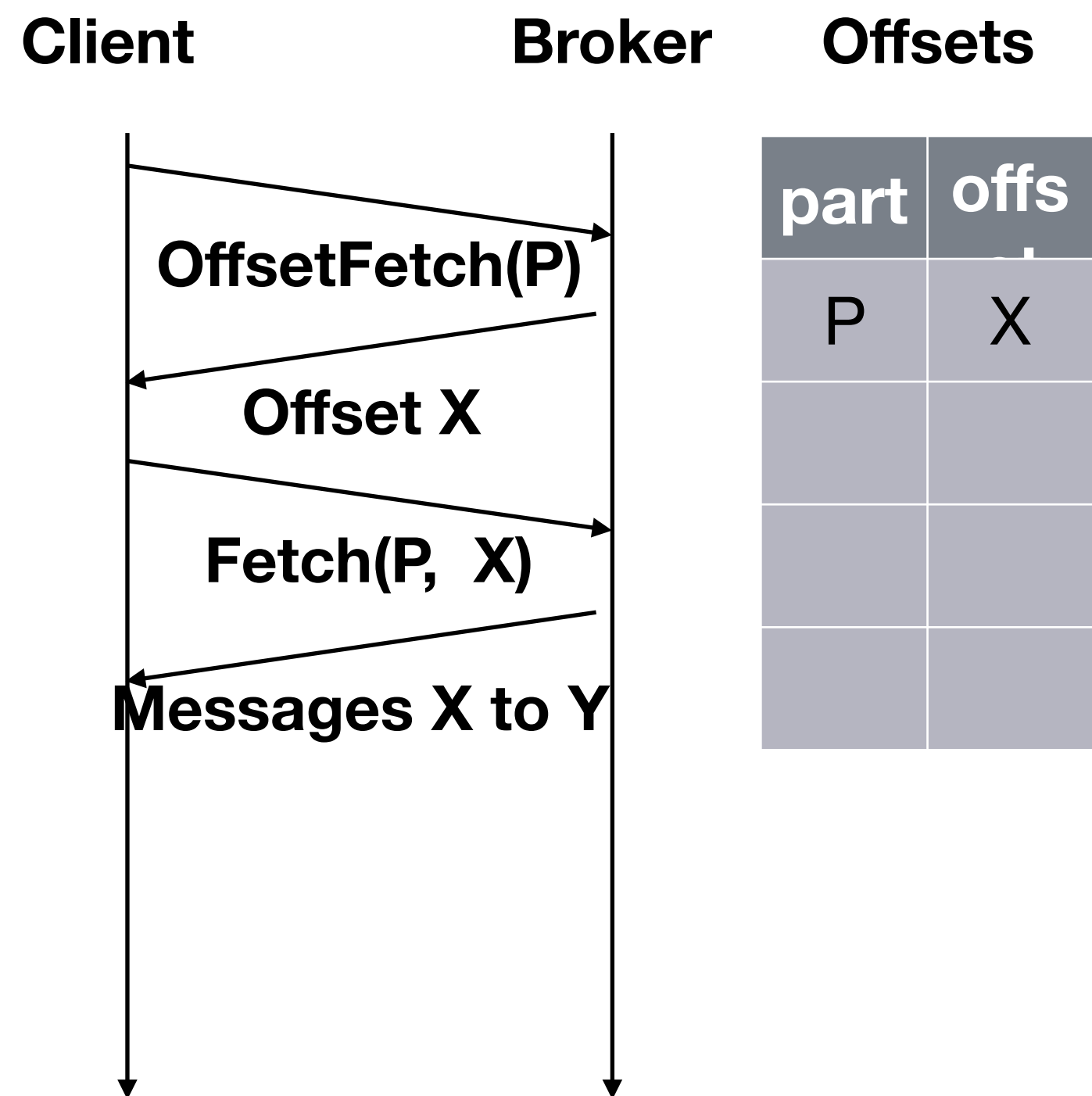
Протокол консьюмера

- Fetch API: Взять диапазон записей по смещению (offset)
- OffsetCommit API: Читать/записывать смещения (offset)



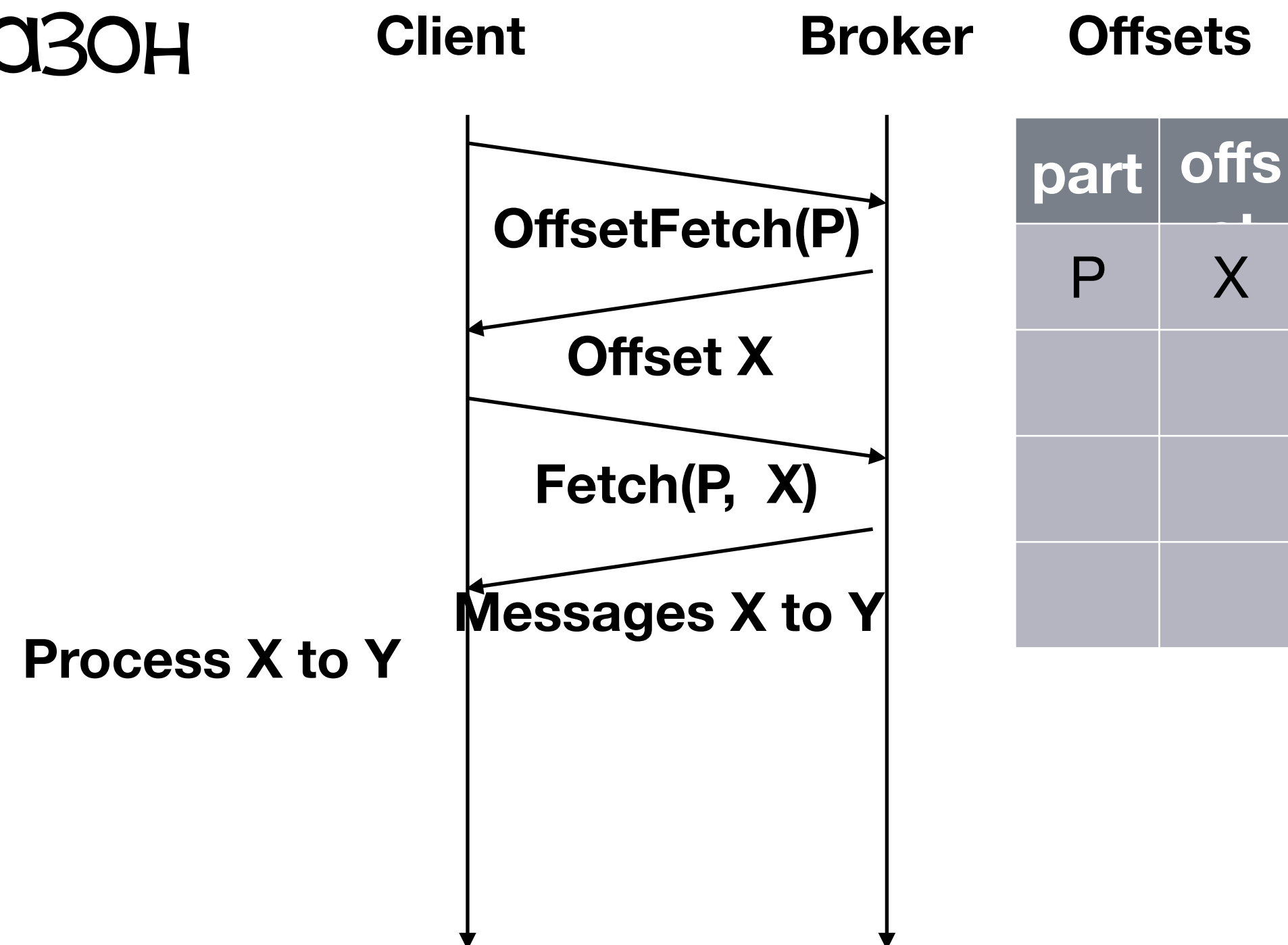
Протокол консьюмера

- Fetch API: Взять диапазон записей по смещению (offset)
- OffsetCommit API: Читать/записывать смещения (offset)



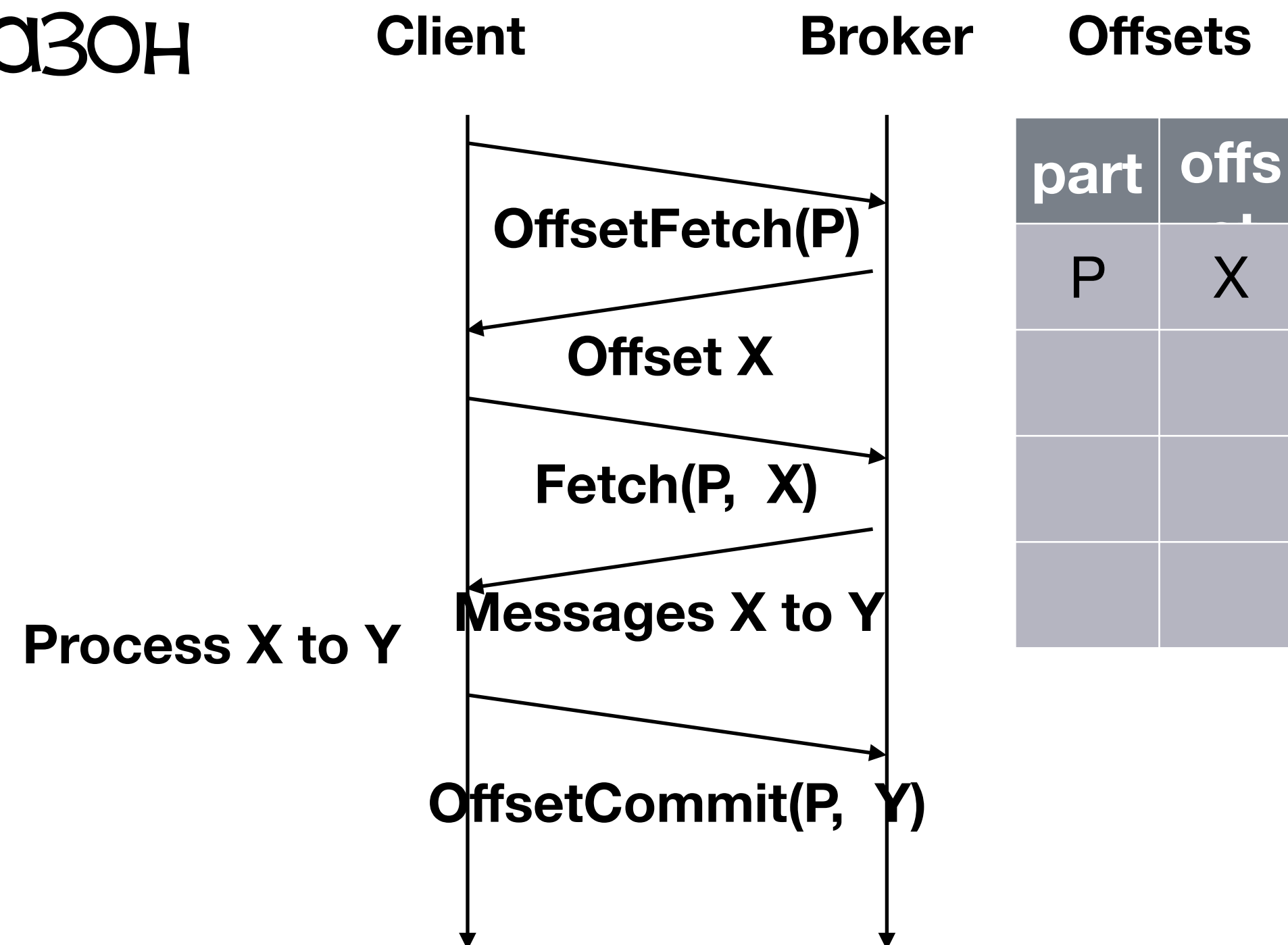
Протокол консьюмера

- Fetch API: Взять диапазон записей по смещению (offset)
- OffsetCommit API: Читать/записывать смещения (offset)



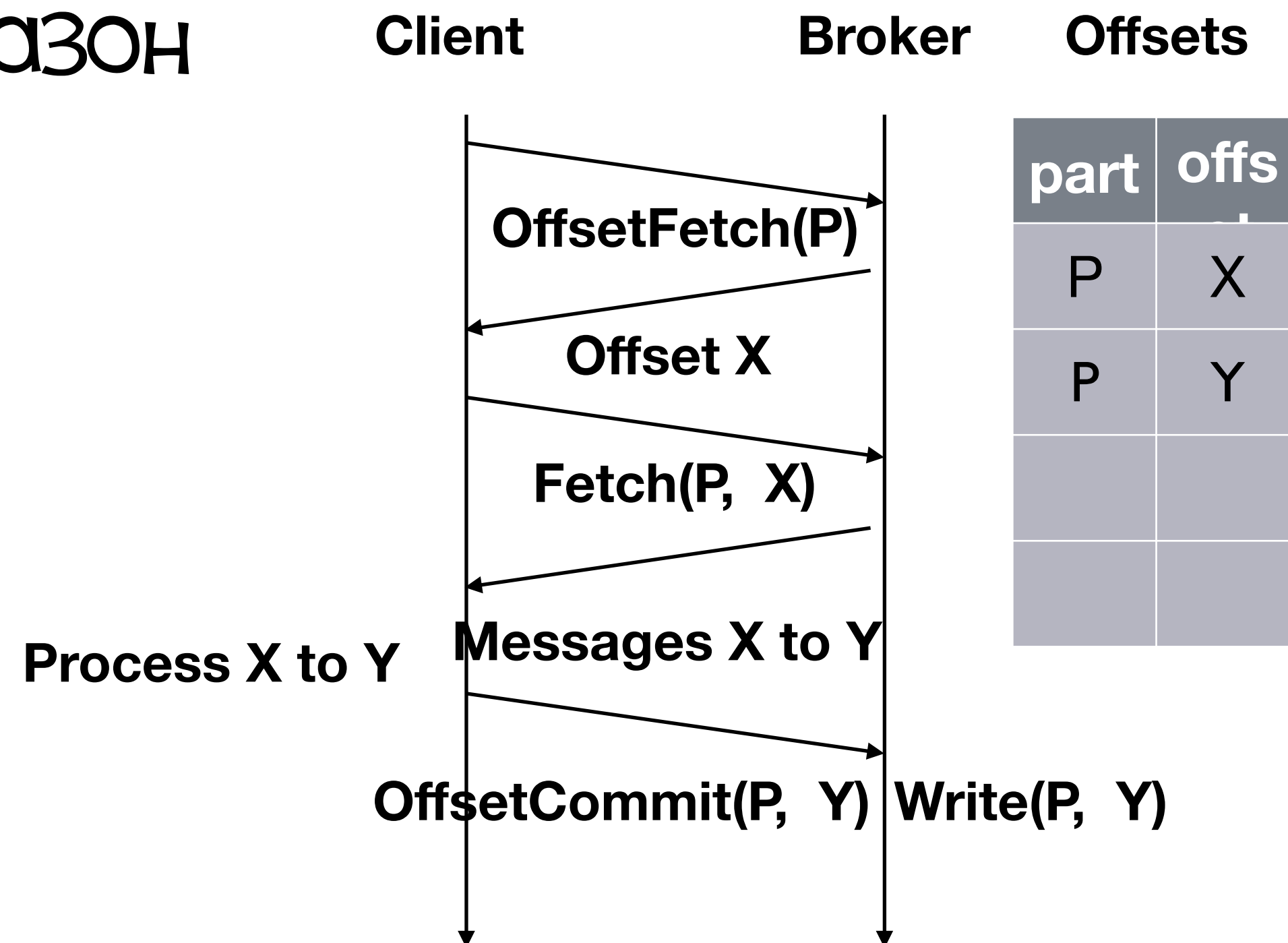
Протокол консьюмера

- Fetch API: Взять диапазон записей по смещению (offset)
- OffsetCommit API: Читать/записывать смещения (offset)



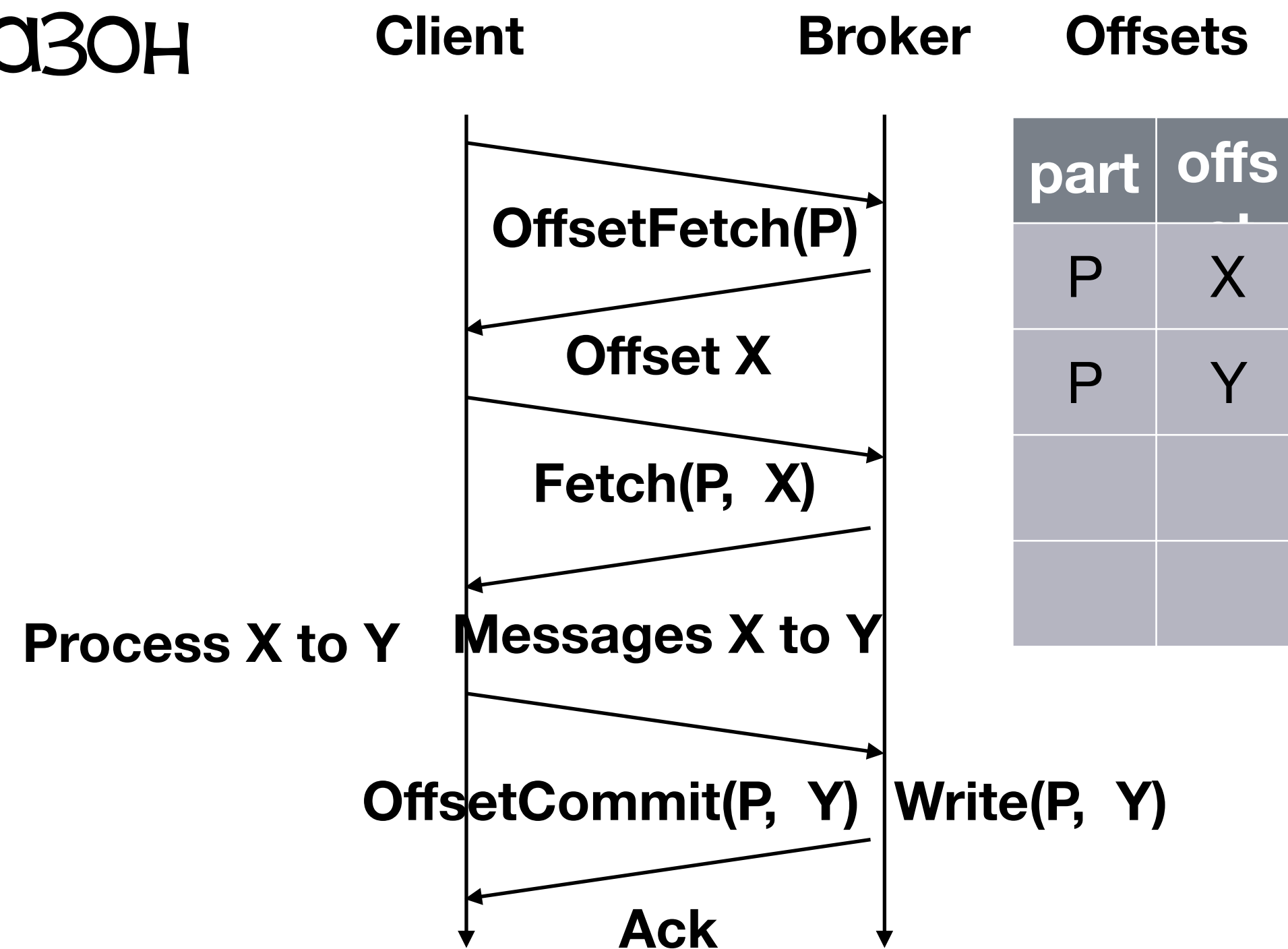
Протокол консьюмера

- Fetch API: Взять диапазон записей по смещению (offset)
- OffsetCommit API: Читать/записывать смещения (offset)



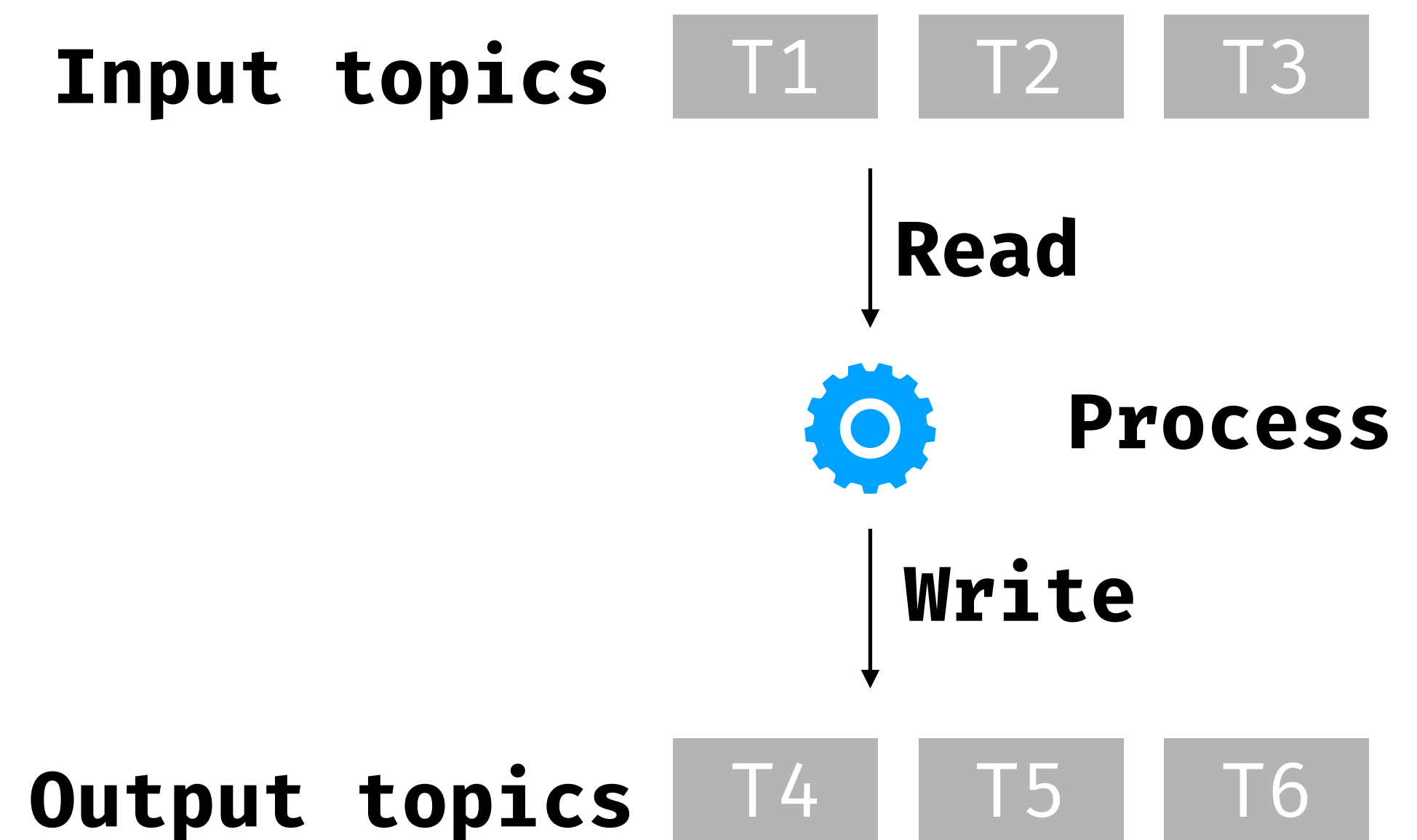
Протокол консьюмера

- Fetch API: Взять диапазон записей по смещению (offset)
- OffsetCommit API: Читать/записывать смещения (offset)



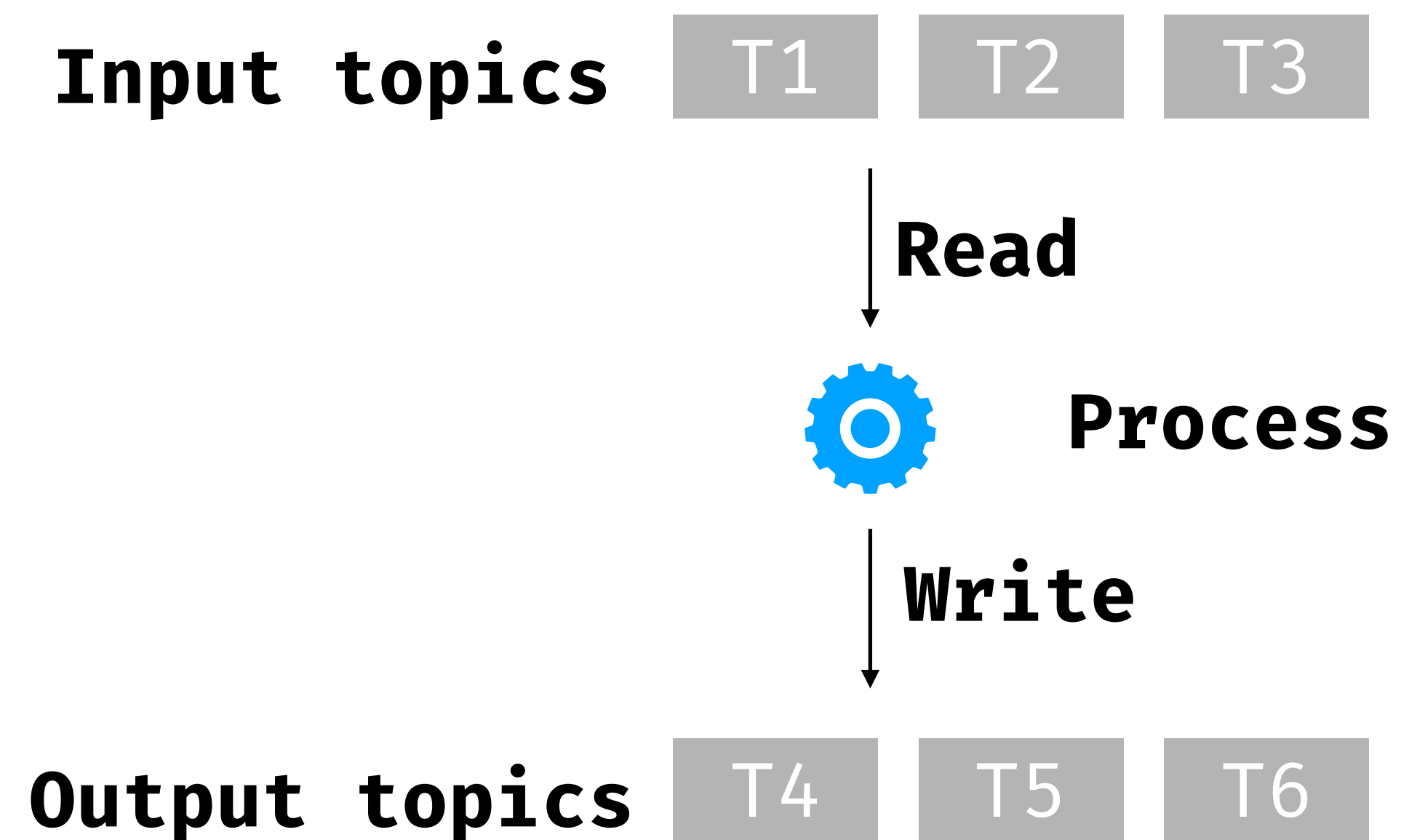
Модель обработки

- Прочитали из топика
- Что-то посчитали
- Записали результат в выходной топик



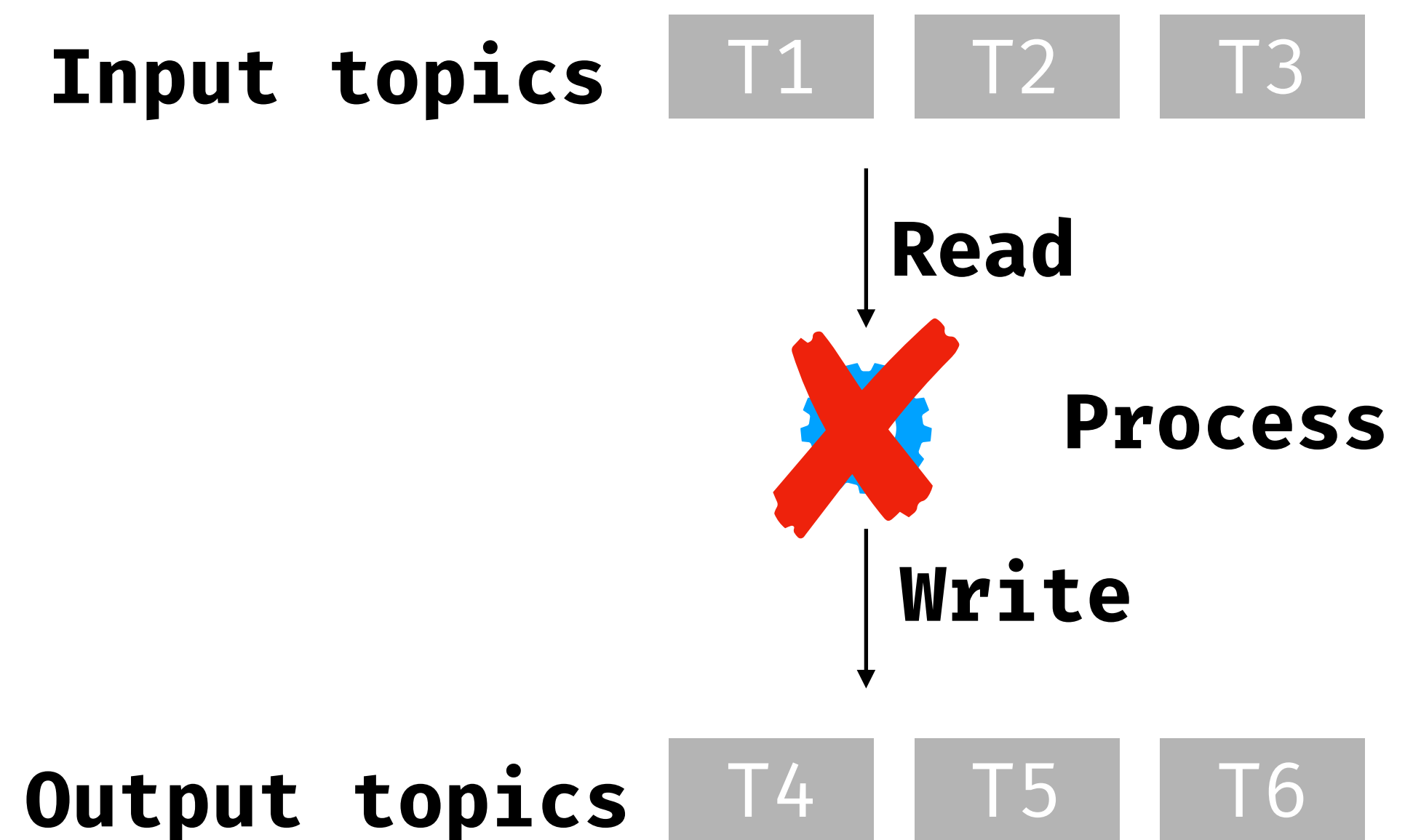
Модель обработки

- Вугы сбоев
- Все сломалось
- Зомби



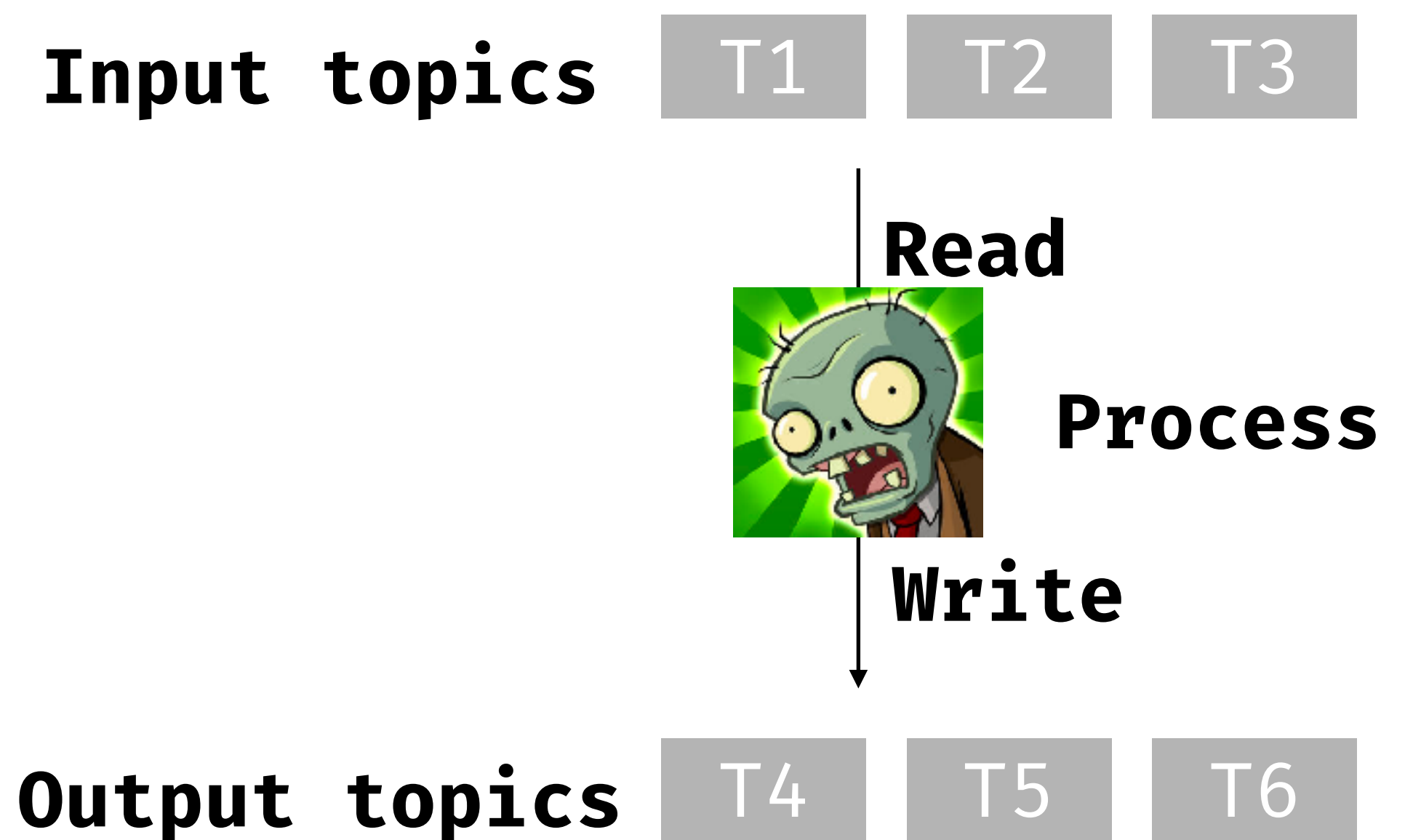
Модель обработки

- Вугы сбоев
- Все сломалось
- Зомби



Модель обработки

- Вугы сбоев
- Все сломалось
- Зомби



Семантики обработки сообщений

- Минимум один раз / At-least-once
- Максимум один раз / At-most-once: Вообще фигня какая-то!
- Точно один раз / Exactly-once: То что вы на самом деле хотите

Семантика «Как минимум один раз»

```
c.subscribe(inputTopics)
while (true) {
    input = c.poll()
    output = process(input)
    p.send(output)
    c.commitOffsets()
}
```

Семантика «Как минимум один раз»

```
c.subscribe(inputTopics)
while (true) {
    input = c.poll()
    output = process(input)
    p.send(output)
    c.commitOffsets()
}
```

```
void subscribe(Collection<String> topics);
```

Семантика «Как минимум один раз»

```
c.subscribe(inputTopics)
while (true) {
    input = c.poll()
    output = process(input)
    p.send(output)
    c.commitOffsets()
}
```

```
ConsumerRecords<K, V> poll(long timeout);
```


Семантика «Как минимум один раз»

```
c.subscribe(inputTopics)
while (true) {
    input = c.poll()
    output = process(input)
    p.send(output)
    c.commitOffsets()
}
```

(Kafka Streams)

Семантика «Как минимум один раз»

```
c.subscribe(inputTopics)
while (true) {
    input = c.poll()
    output = process(input)
    p.send(output)
    c.commitOffsets()
}
```

```
Future<RecordMetadata> send(ProducerRecord<K, V> r);
```

Семантика «Как минимум один раз»

```
c.subscribe(inputTopics)
while (true) {
    input = c.poll()
    output = process(input)
    p.send(output)
    c.commitOffsets()
}
```

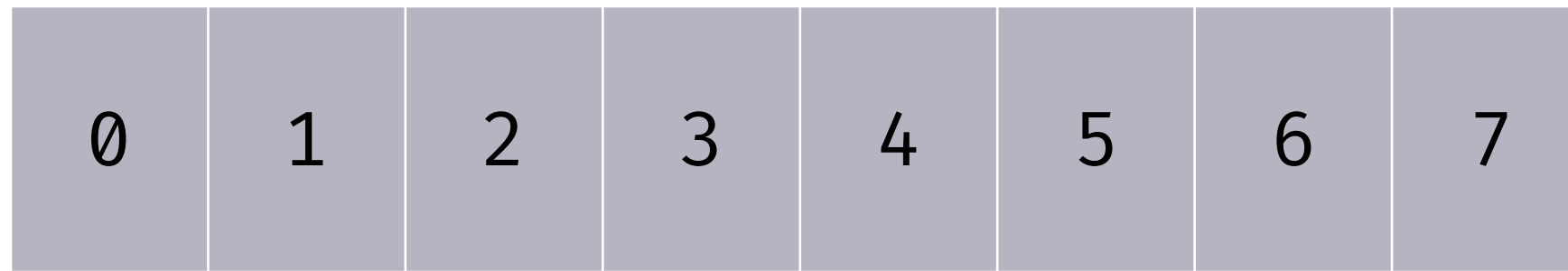
```
void commitSync();
```

Семантика «Как минимум один раз»

```
while (true) {  
    data = c.poll()  
    p.send(data)  
    c.commitOffsets()  
}
```

Например, когда просто надо скопировать ввод на вывод

Input



Output

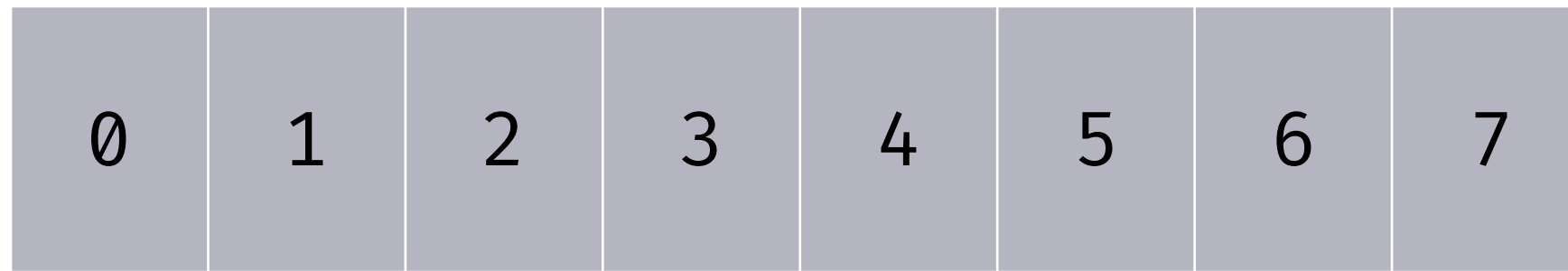


**Committed
Offsets**

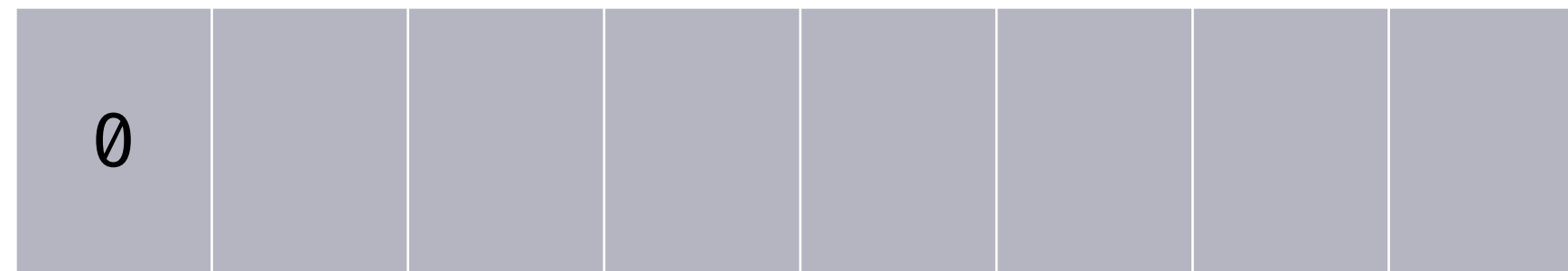


```
while (true) {  
    data = c.poll()  
    p.send(data)  
    c.commitOffsets()  
}
```

Input



Output

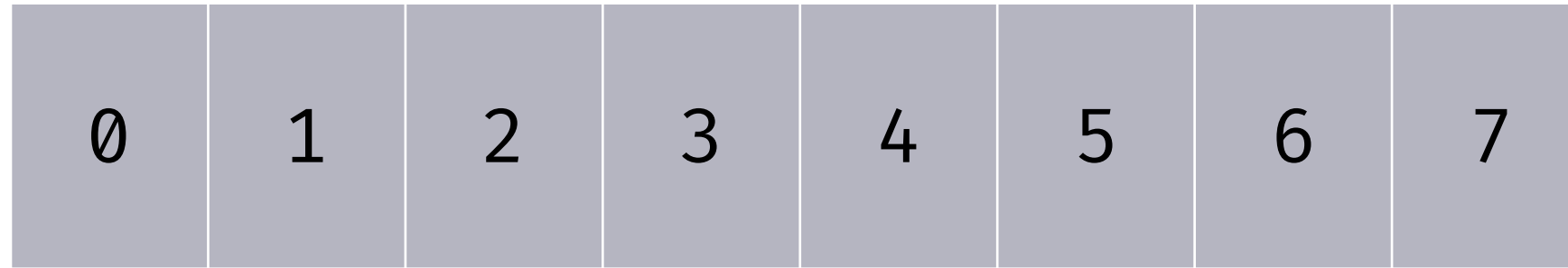


**Committed
Offsets**



```
while (true) {  
    data = c.poll()  
    p.send(data)  
    c.commitOffsets()  
}
```

Input



Output

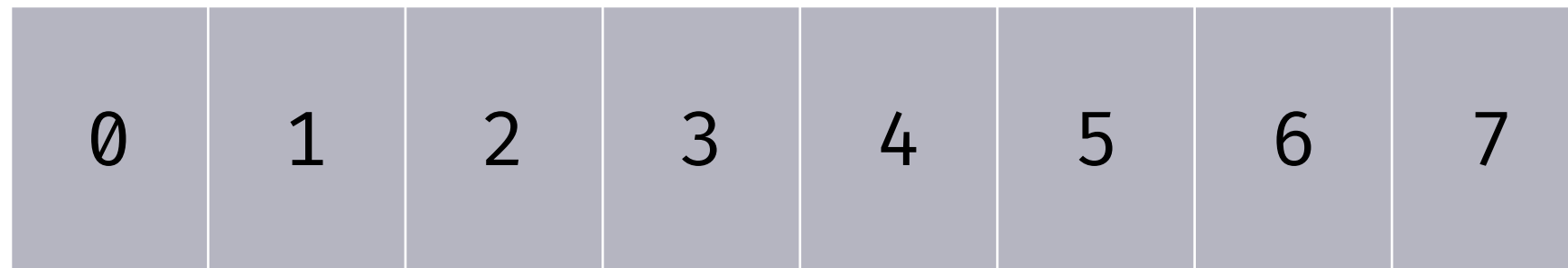


**Committed
Offsets**



```
while (true) {  
    data = c.poll()  
    p.send(data)  
    c.commitOffsets()  
}
```

Input



Output

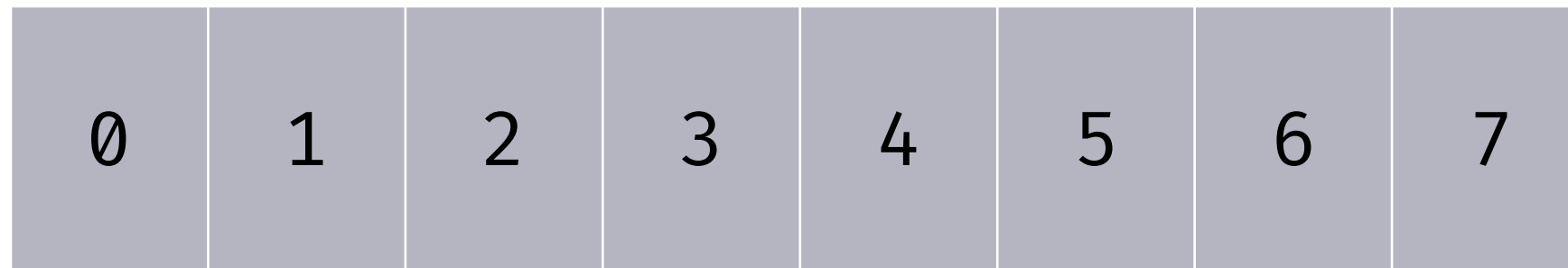


**Committed
Offsets**



```
while (true) {  
    data = c.poll()  
    p.send(data)  
    c.commitOffsets()  
}
```


Input



Output

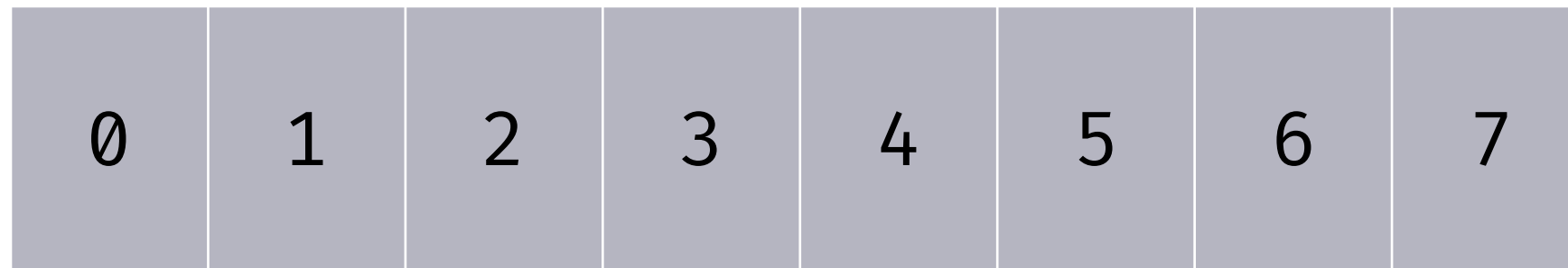


**Committed
Offsets**



```
while (true) {  
    data = c.poll()  
    p.send(data)  
    c.commitOffsets()  
}
```

Input



Output

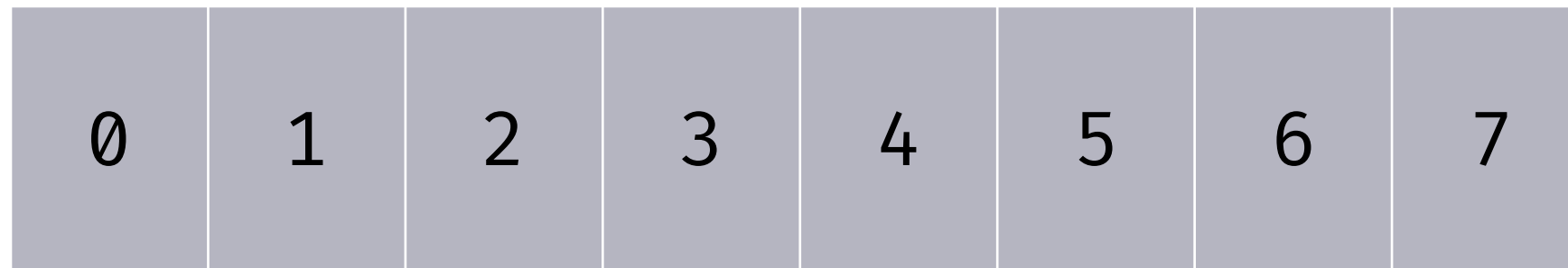


**Committed
Offsets**



```
while (true) {  
    data = c.poll()  
    p.send(data)  
    c.commitOffsets()  
} Ack lost, retry!
```

Input



Output

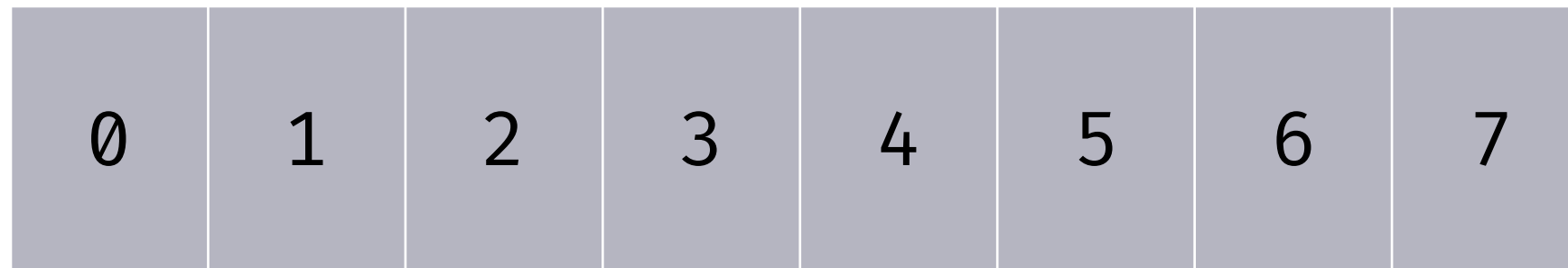


**Committed
Offsets**



```
while (true) {  
    data = c.poll()  
    p.send(data)  
    c.commitOffsets()  
}
```

Input



Output

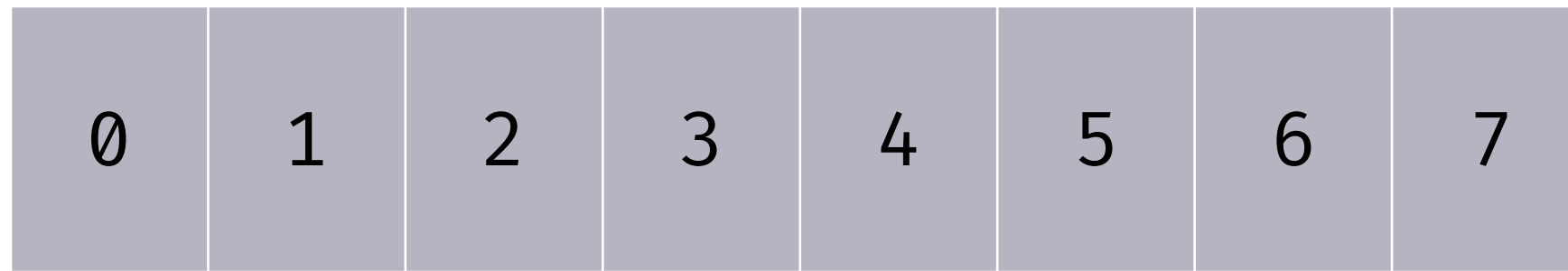


**Committed
Offsets**

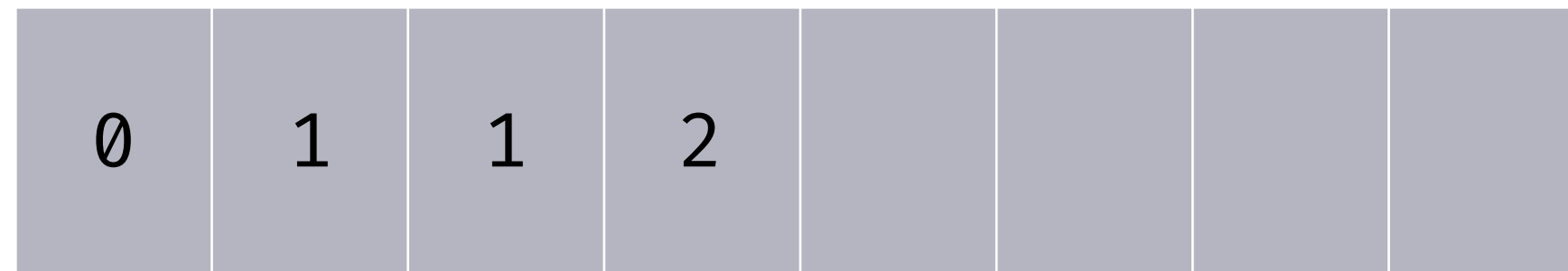


```
while (true) {  
    data = c.poll()  
    p.send(data)  
    c.commitOffsets()  
}
```

Input



Output

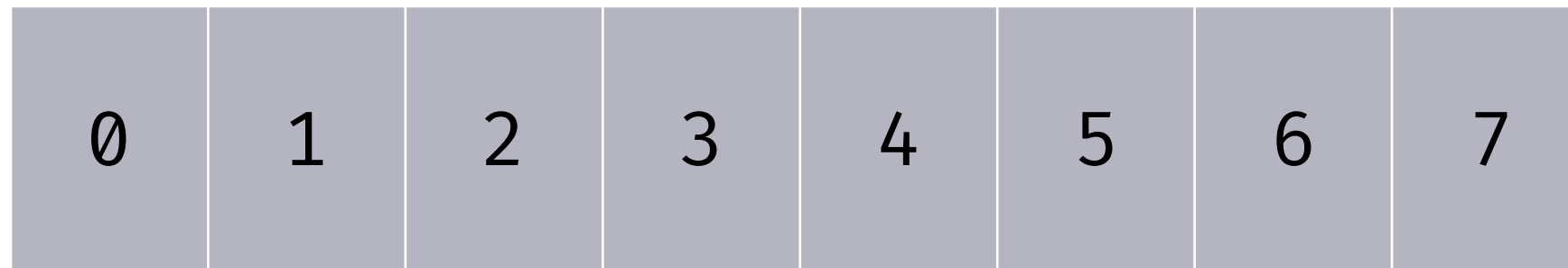


**Committed
Offsets**

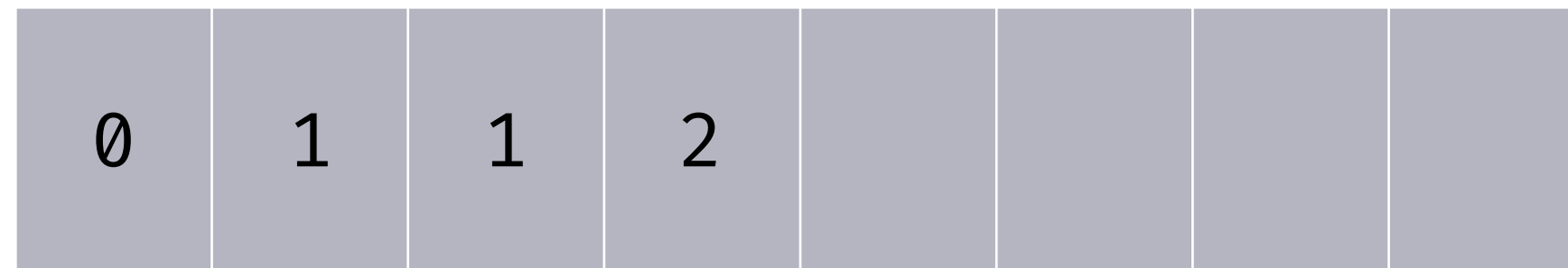


```
while (true) {  
    data = c.poll()  
    p.send(data)  
    c.commitOffsets()  
}
```

Input



Output



**Committed
Offsets**



```
while (true) {  
    data = c.poll()  
    p.send(data)  
    c.commitOffsets()  
}
```

Process crash!

Input

0	1	2	3	4	5	6	7
---	---	---	---	---	---	---	---

Output

0	1	1	2				
---	---	---	---	--	--	--	--

**Committed
Offsets**

1	2						
---	---	--	--	--	--	--	--

```
while (true) {  
    data = c.poll()  
    p.send(data)  
    c.commitOffsets()  
}
```

Restart!

Input

0	1	2	3	4	5	6	7
---	---	---	---	---	---	---	---

Output

0	1	1	2				
---	---	---	---	--	--	--	--

**Committed
Offsets**

1	2						
---	---	--	--	--	--	--	--

```
while (true) {  
    data = c.poll()  
    p.send(data)  
    c.commitOffsets()  
}
```


Input

0	1	2	3	4	5	6	7
---	---	---	---	---	---	---	---

Output

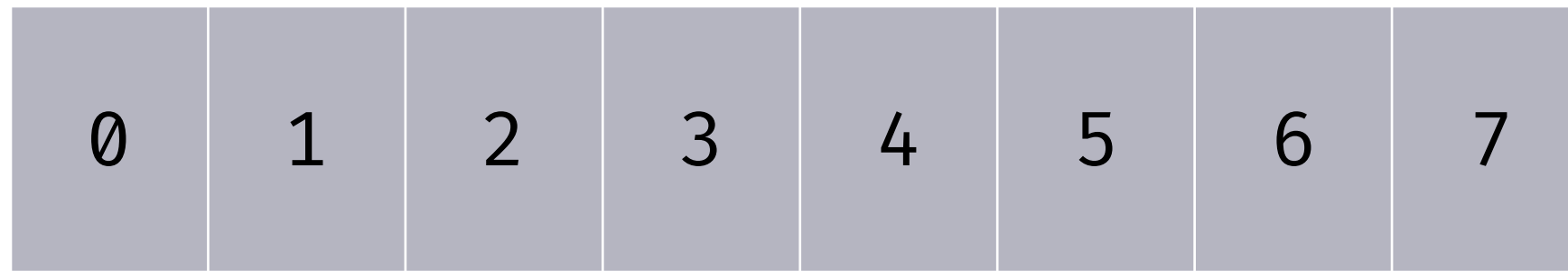
0	1	1	2	2			
---	---	---	---	---	--	--	--

**Committed
Offsets**

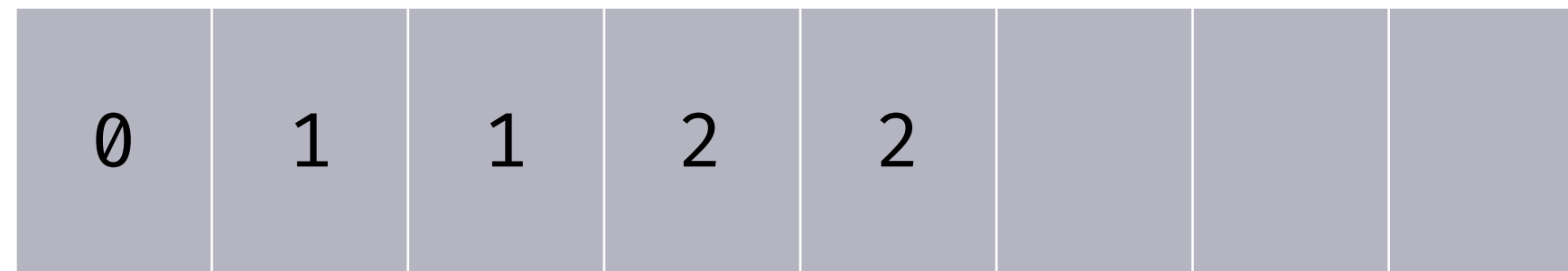
1	2						
---	---	--	--	--	--	--	--

```
while (true) {  
    data = c.poll()  
    p.send(data)  
    c.commitOffsets()  
}
```

Input



Output

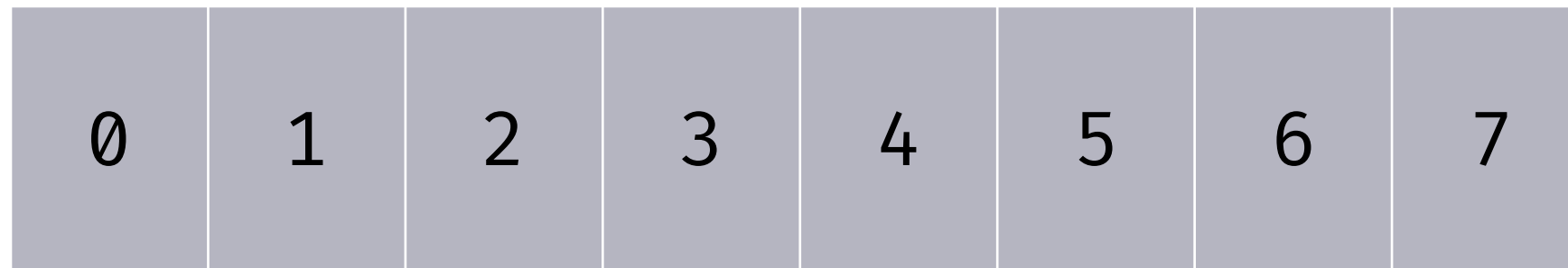


**Committed
Offsets**

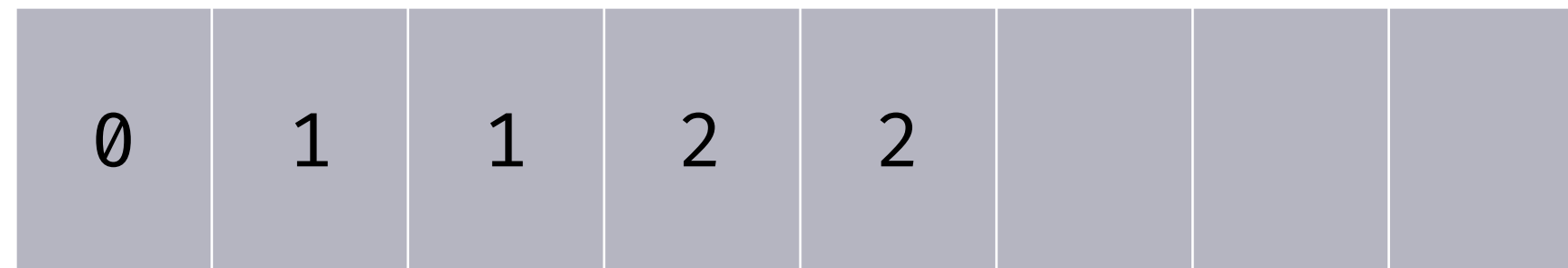


```
while (true) {  
    data = c.poll()  
    p.send(data)  
    c.commitOffsets()  
}
```

Input



Output



**Committed
Offsets**



```
while (true) {  
    data = c.poll()  
    p.send(data)  
    c.commitOffsets()  
}
```

Input

0	1	2	3	4	5	6	7
---	---	---	---	---	---	---	---


Output

0	1	1	2	2			
---	---	---	---	---	--	--	--

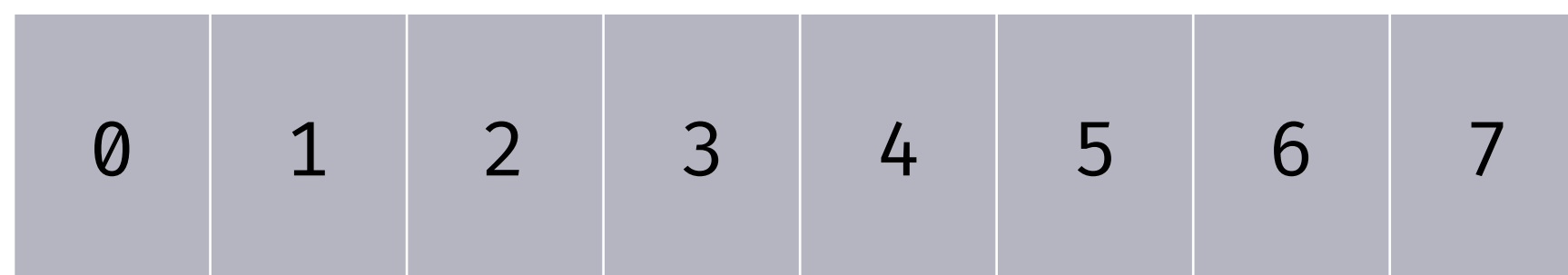
**Committed
Offsets**

1	2	3					
---	---	---	--	--	--	--	--

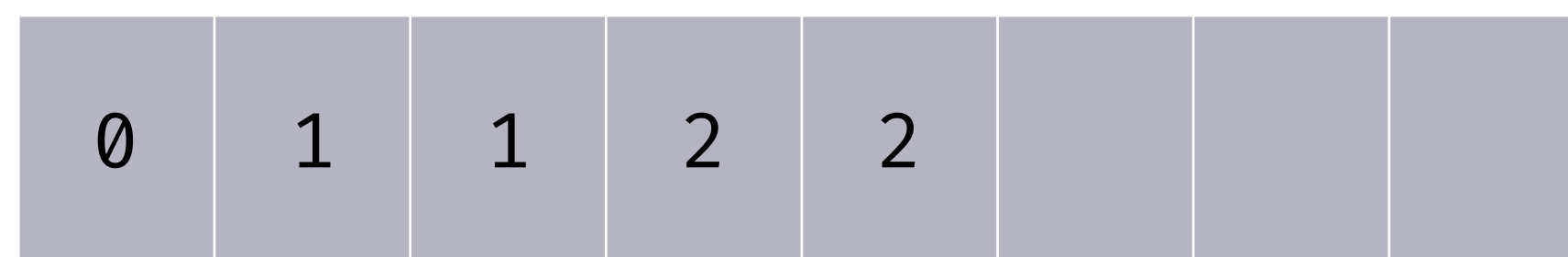
```
while {  
  dat ll()  
  p.s )  
  c.c sets()  
}
```



Input



Output



**Committed
Offsets**

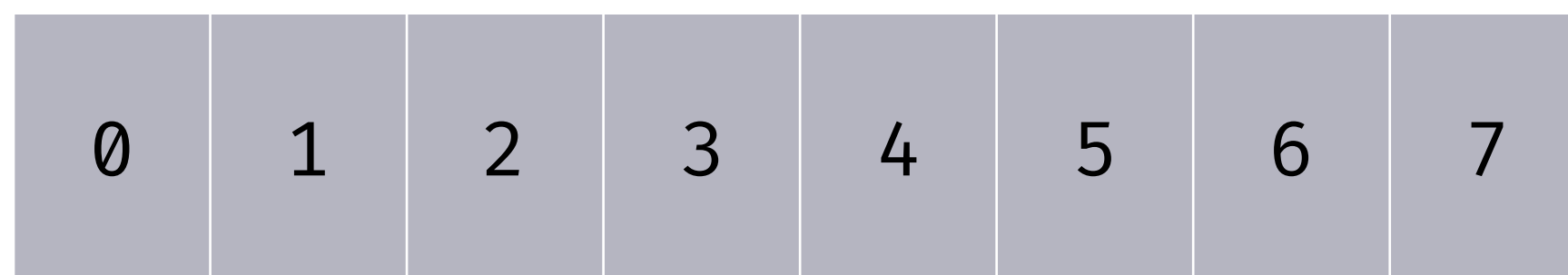


```
while (true) {  
  data = c.poll()  
  p.send(data)  
  c.commitOffsets()  
}
```

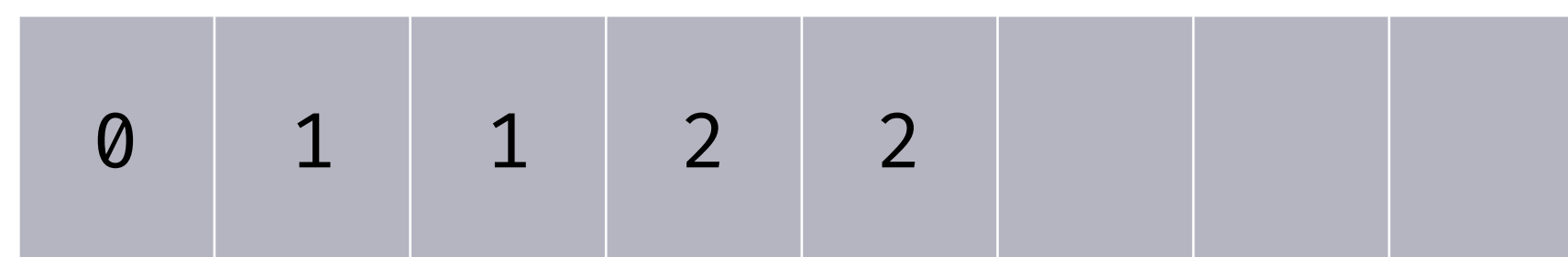


```
while (true) {  
  data = c.poll()  
  p.send(data)  
  c.commitOffsets()  
}
```

Input



Output



**Committed
Offsets**



```
while (true) {  
  data = c.poll()  
  p.send(data)  
  c.commitOffsets()  
}
```



```
while (true) {  
  data = c.poll()  
  p.send(data)  
  c.commitOffsets()  
}
```

Input

0	1	2	3	4	5	6	7
---	---	---	---	---	---	---	---

Output

0	1	1	2	2	3		
---	---	---	---	---	---	--	--

**Committed
Offsets**

1	2	3					
---	---	---	--	--	--	--	--

```
while (true) {  
  data = c.poll()  
  p.send(data)  
  c.commitOffsets()  
}
```



```
while (true) {  
  data = c.poll()  
  p.send(data)  
  c.commitOffsets()  
}
```

Input

0	1	2	3	4	5	6	7
---	---	---	---	---	---	---	---

Output

0	1	1	2	2	3		
---	---	---	---	---	---	--	--

**Committed
Offsets**

1	2	3	4				
---	---	---	---	--	--	--	--

```
while (true) {  
  data = c.poll()  
  p.send(data)  
  c.commitOffsets()  
}
```



```
while (true) {  
  data = c.poll()  
  p.send(data)  
  c.commitOffsets()  
}
```


Input

0	1	2	3	4	5	6	7
---	---	---	---	---	---	---	---


Output

0	1	1	2	2	3		
---	---	---	---	---	---	--	--

**Committed
Offsets**

1	2	3	4				
---	---	---	---	--	--	--	--

```
while (true) {  
  data = c.poll()  
  p.send(data)  
  c.commitOffsets()  
}
```



```
while (true) {  
  data = c.poll()  
  p.send(data)  
  c.commitOffsets()  
}
```

Input

0	1	2	3	4	5	6	7
---	---	---	---	---	---	---	---

Output

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

**Committed
Offsets**

1	2	3	4				
---	---	---	---	--	--	--	--

```
while (true) {  
  data = c.poll()  
  p.send(data)  
  c.commitOffsets()  
}
```



```
while (true) {  
  data = c.poll()  
  p.send(data)  
  c.commitOffsets()  
}
```

Input

0	1	2	3	4	5	6	7
---	---	---	---	---	---	---	---

Output

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

**Committed
Offsets**

1	2	3	4				
---	---	---	---	--	--	--	--

```
while (true) {  
    data = c.poll()  
    p.send(data)  
    c.commitOffsets()  
}
```

```
while (true) {  
    data = c.poll()  
    p.send(data)  
    c.commitOffsets()  
}
```

Input

0	1	2	3	4	5	6	7
---	---	---	---	---	---	---	---

Output

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

**Committed
Offsets**

1	2	3	4				
---	---	---	---	--	--	--	--

```
while (true) {  
    data = c.poll()  
    p.send(data)  
    c.commitOffsets()  
}  
  
while (true) {  
    data = c.poll()  
    p.send(data)  
    c.commitOffsets()  
}
```

Семантика «Как минимум один раз»

- Дубликаты

- Повторы прогьюсера

- Сбой перед записью смещений (offset)

- Зомби



Input

0	1	2	3	4	5	6	7
---	---	---	---	---	---	---	---

Output

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

**Committed
Offsets**

1	2	3	4				
---	---	---	---	--	--	--	--

Семантика «максимум один раз»

Вам это не надо

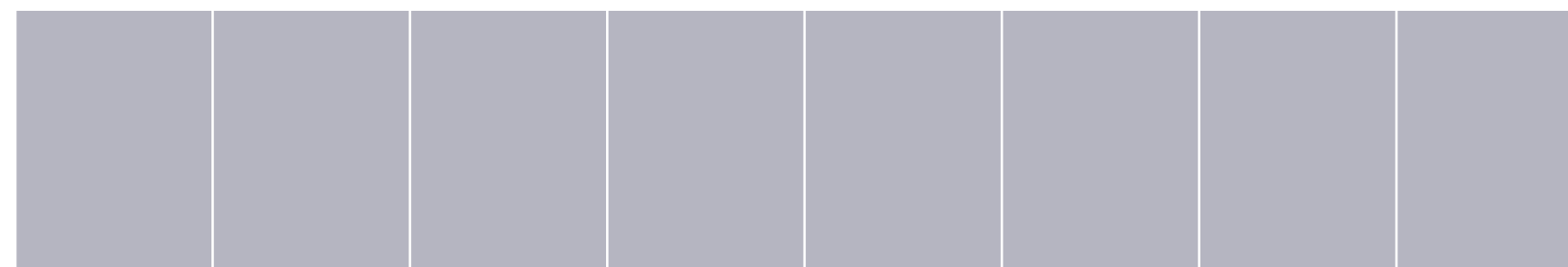
Input



Output



**Committed
Offsets**



```
while (true) {  
    data = c.poll()  
    c.commitOffsets()  
    p.send(data)  
}
```

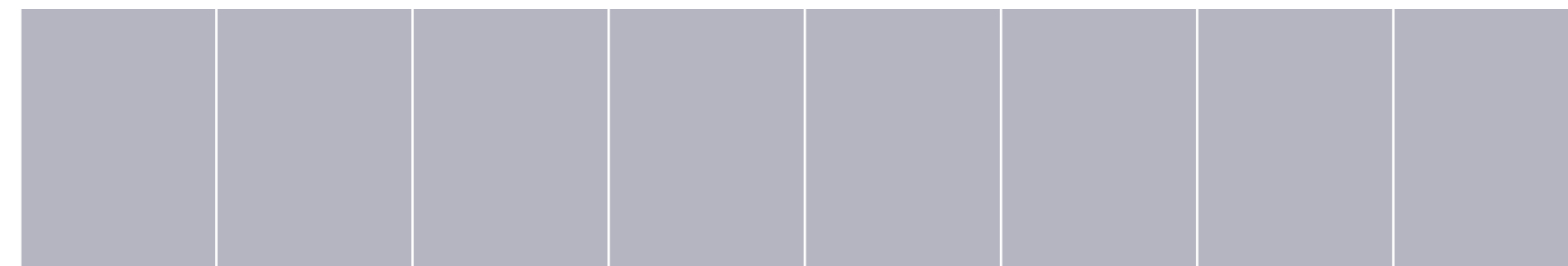
Семантика «максимум один раз»

Вам это не надо

Input



Output



**Committed
Offsets**

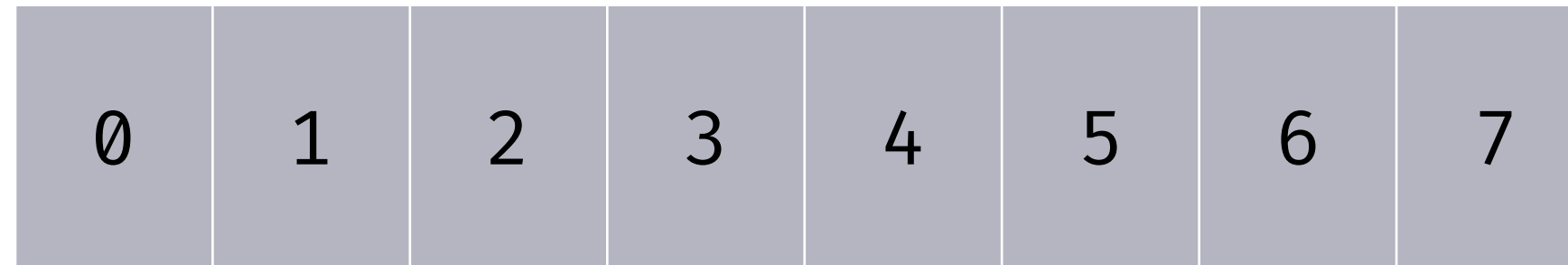


```
while (true) {  
    data = c.poll()  
    c.commitOffsets()  
    p.send(data)  
}
```

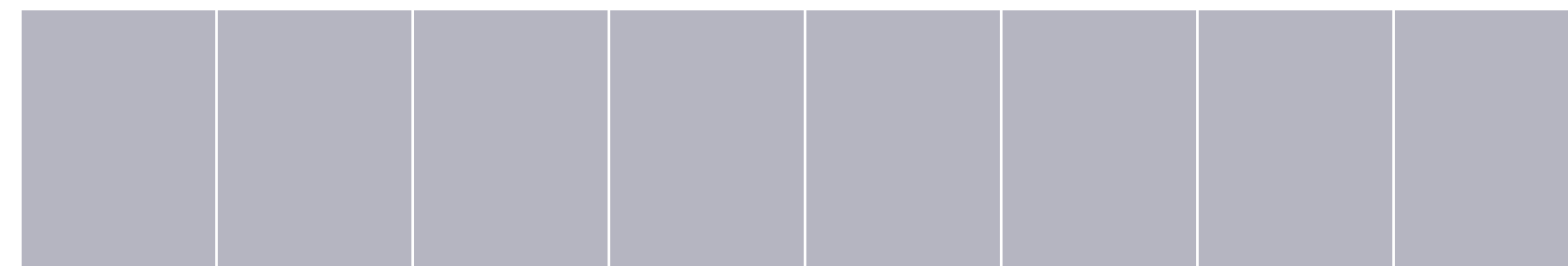
Семантика «максимум один раз»

Вам это не надо

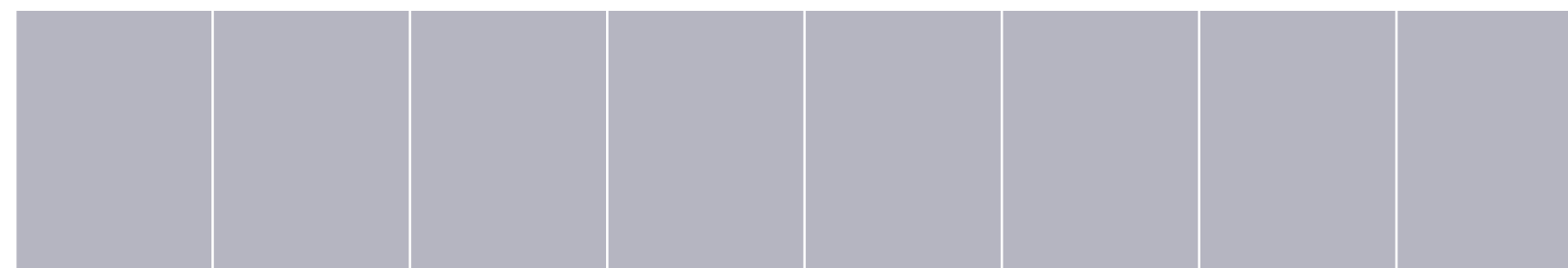
Input



Output



**Committed
Offsets**

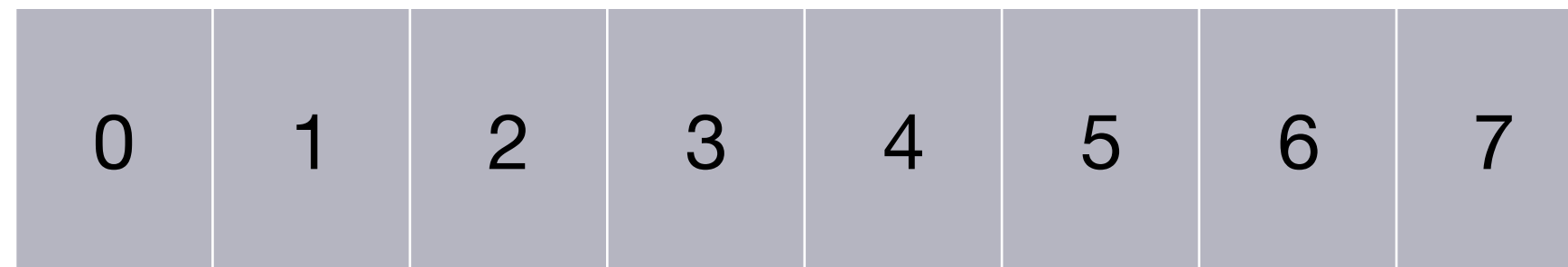


```
while (true) {  
    data = c.poll()  
    c.commitOffsets()  
    p.send(data)  
}
```


Семантика «максимум один раз»

Вам это не надо

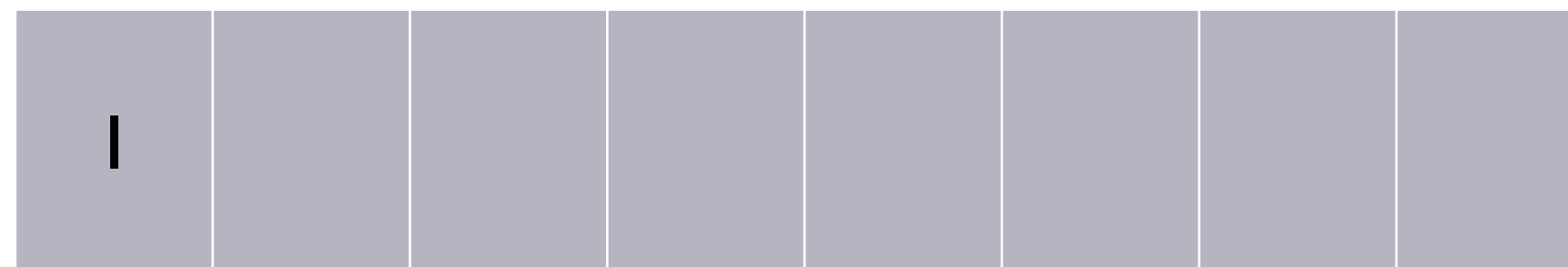
Input



Output



**Committed
Offsets**

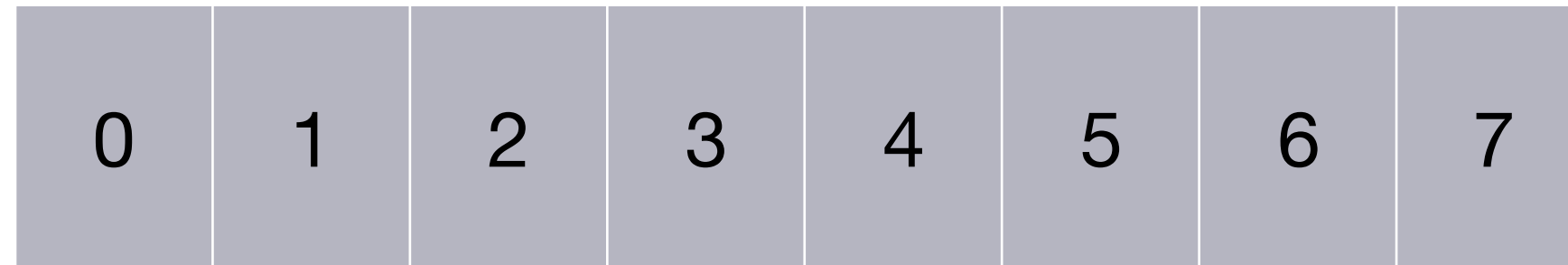


```
while (true) {  
    data = c.poll()  
    c.commitOffsets()  
    p.send(data)  
}
```

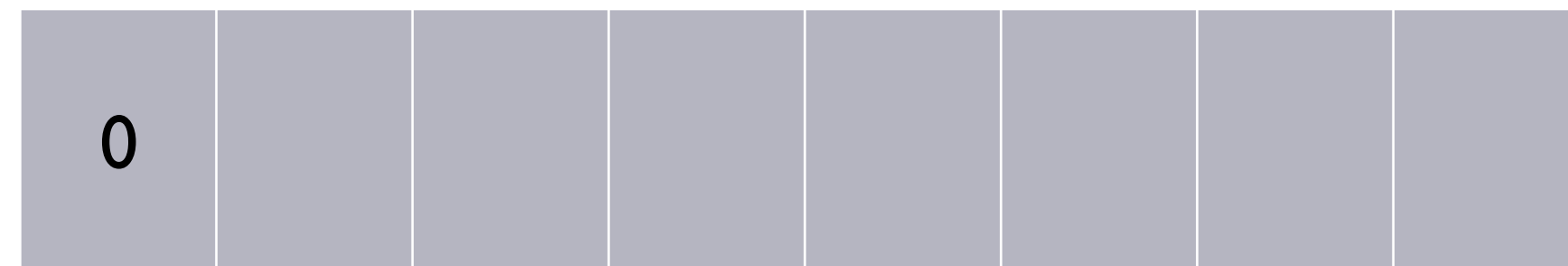
Семантика «максимум один раз»

Вам это не надо

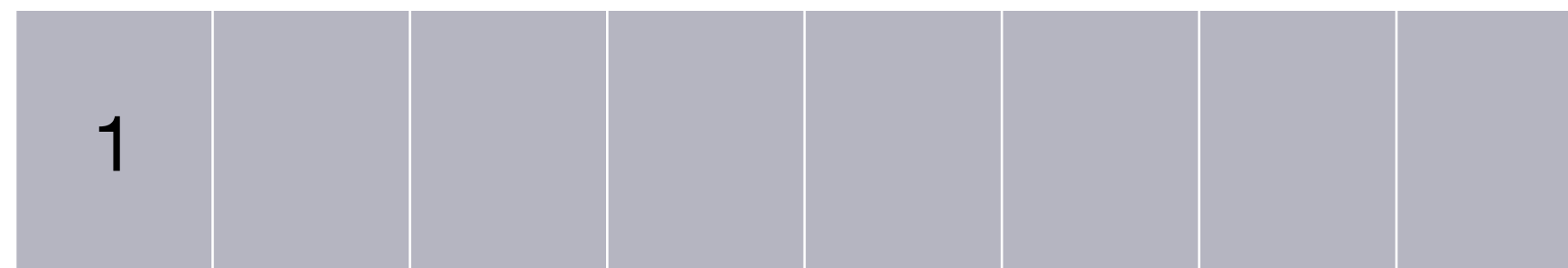
Input



Output



**Committed
Offsets**

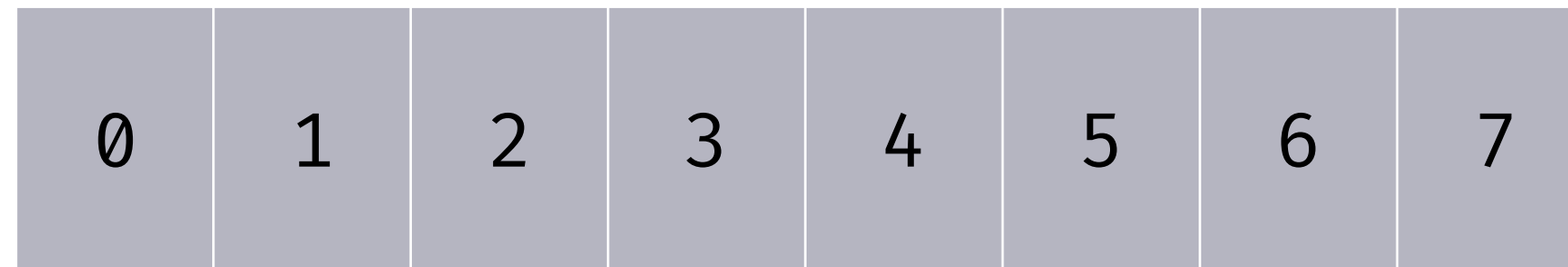


```
while (true) {  
    data = c.poll()  
    c.commitOffsets()  
    p.send(data)  
}
```

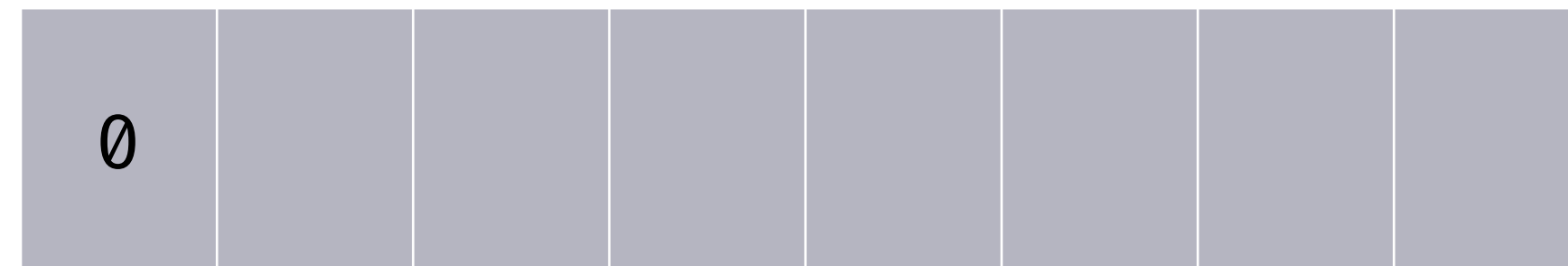
Семантика «максимум один раз»

Вам это не надо

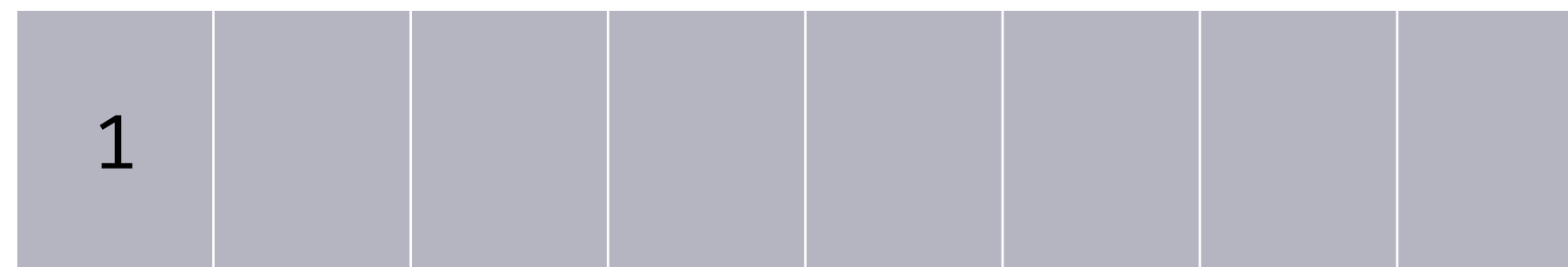
Input



Output



**Committed
Offsets**

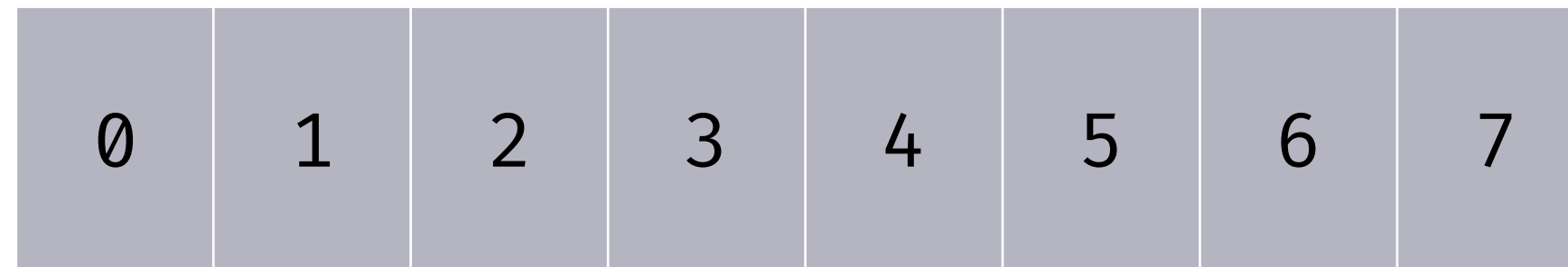


```
while (true) {  
    data = c.poll()  
    c.commitOffsets()  
    p.send(data)  
}
```

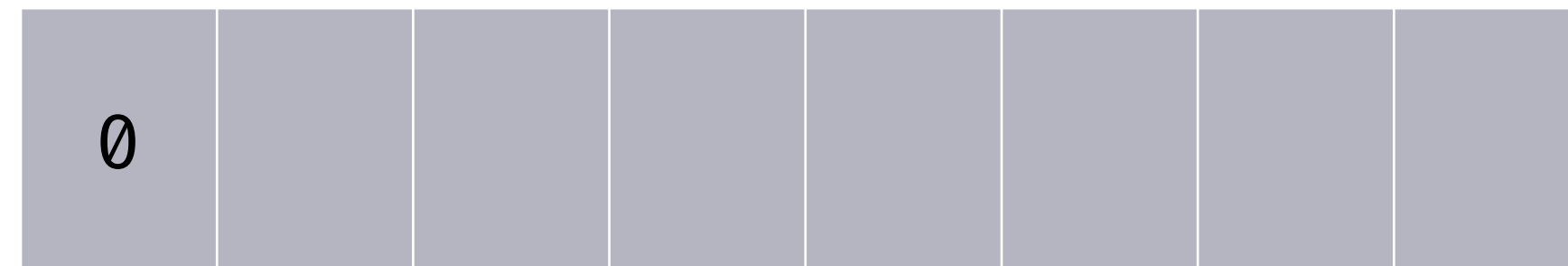
Семантика «максимум один раз»

Вам это не надо

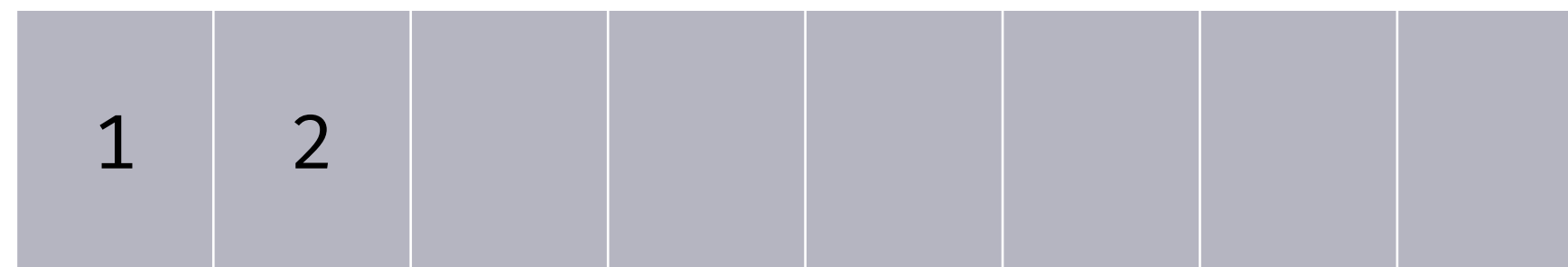
Input



Output



**Committed
Offsets**

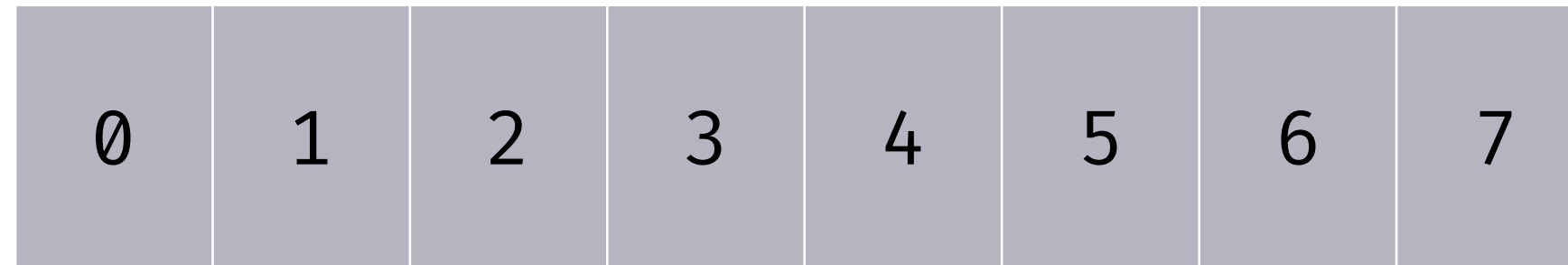


```
while (true) {  
    data = c.poll()  
    c.commitOffsets()  
    p.send(data)  
}
```

Семантика «максимум один раз»

Вам это не надо

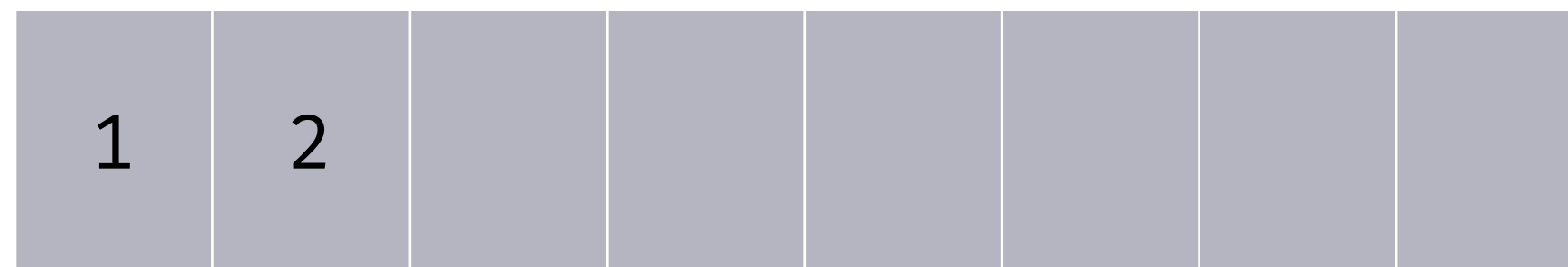
Input



Output



**Committed
Offsets**

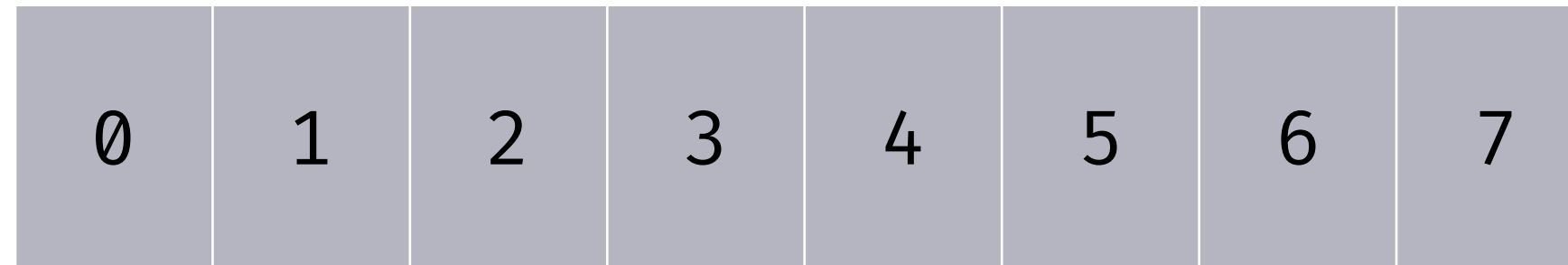


```
while (true) {  
    data = c.poll()  
    c.commitOffsets()  
    p.send(data)  
}
```

Семантика «максимум один раз»

Вам это не надо

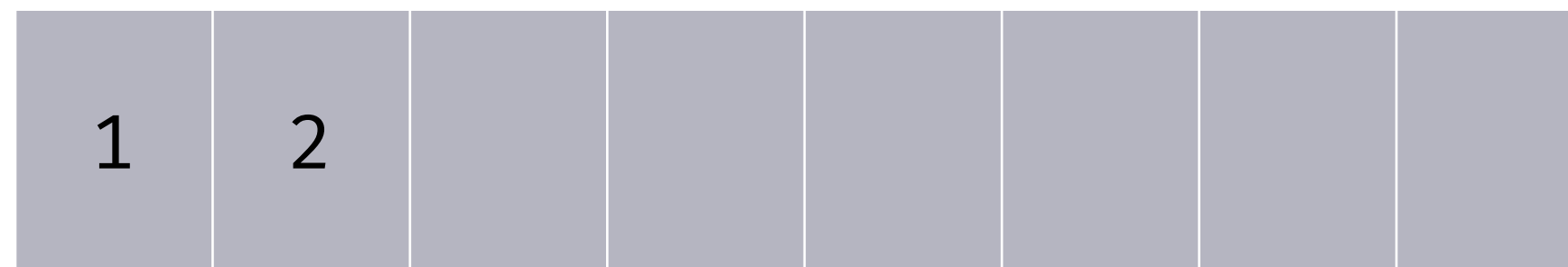
Input



Output



**Committed
Offsets**



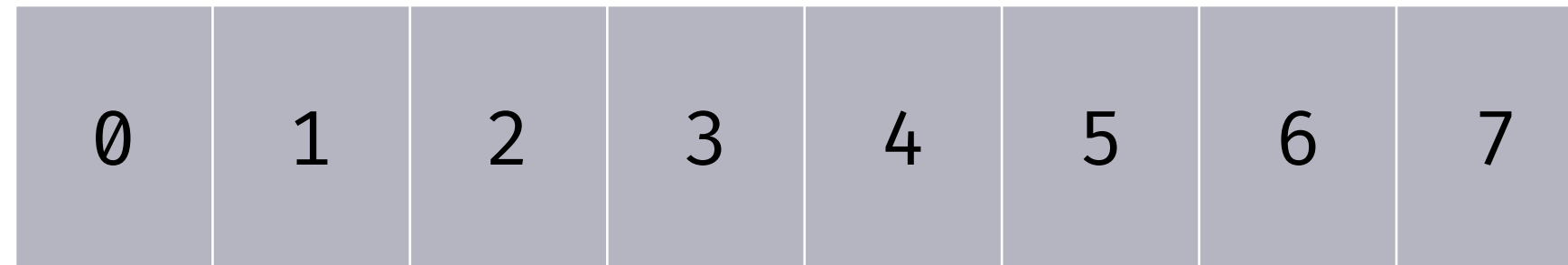
```
while (true) {  
    data = c.poll()  
    c.commitOffsets()  
    p.send(data)  
}
```

Ack lost, yolo!

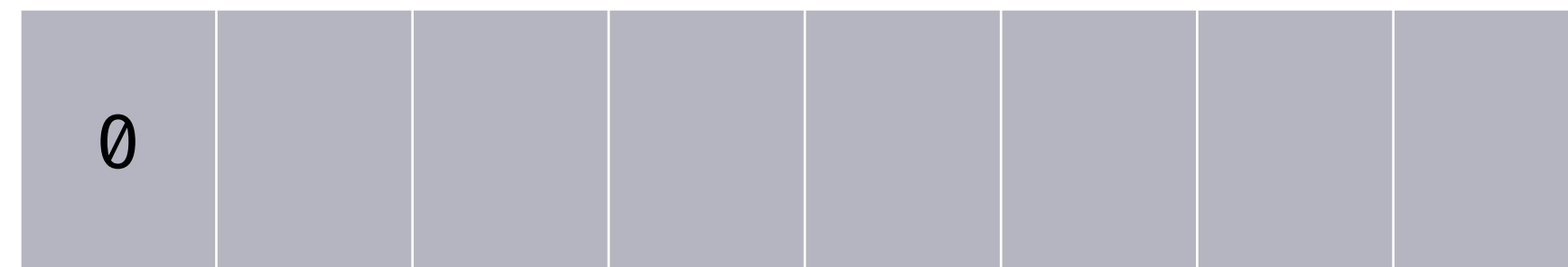
Семантика «максимум один раз»

Вам это не надо

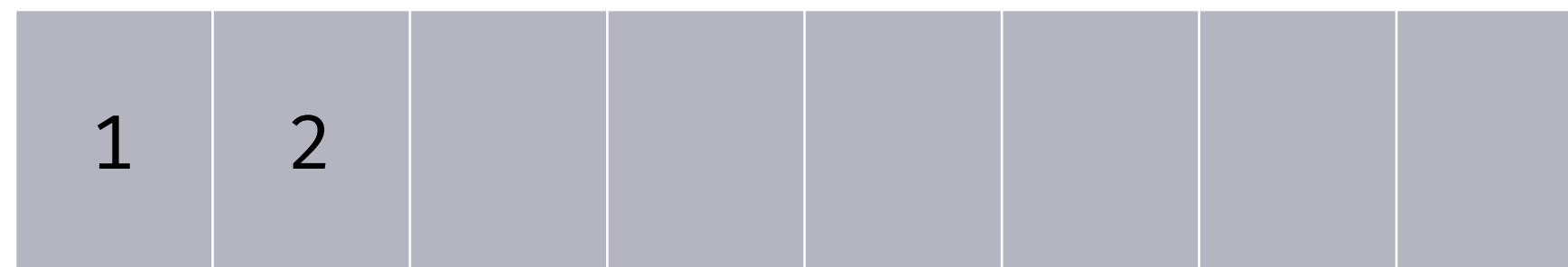
Input



Output



**Committed
Offsets**

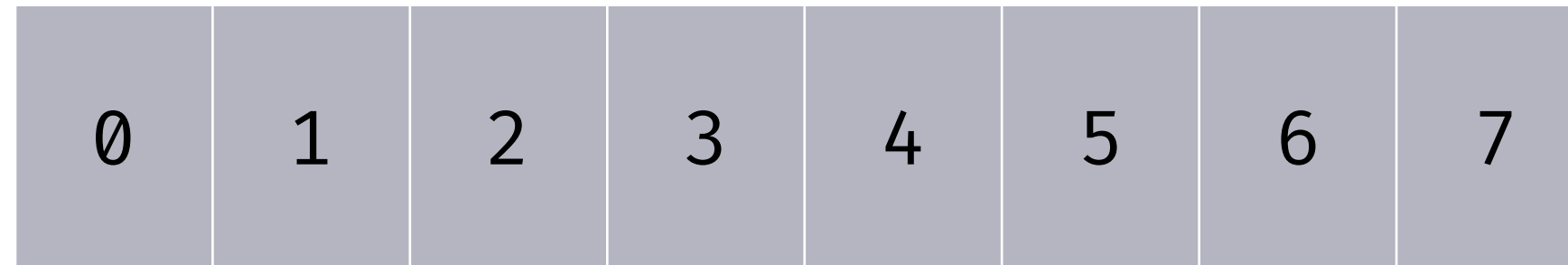


```
while (true) {  
    data = c.poll()  
    c.commitOffsets()  
    p.send(data)  
}
```

Семантика «максимум один раз»

Вам это не надо

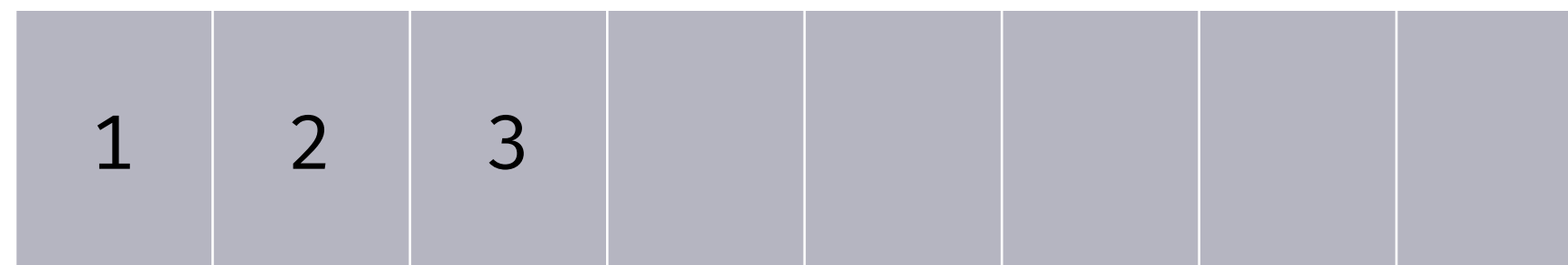
Input



Output



**Committed
Offsets**

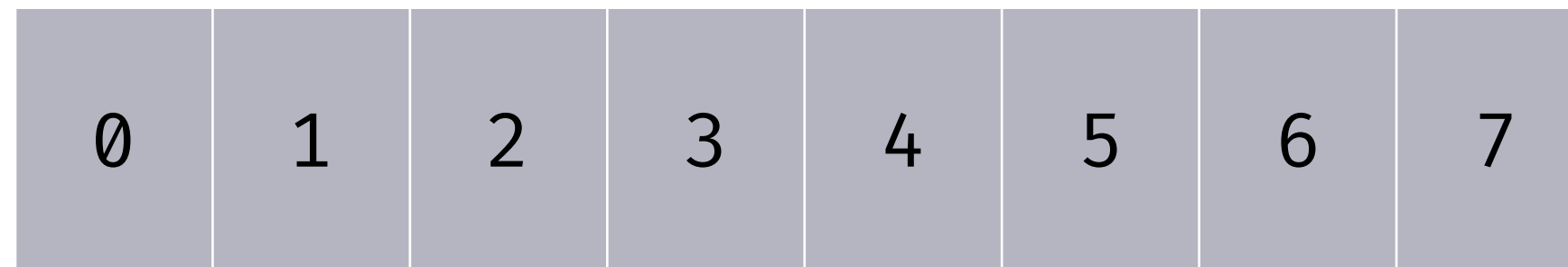


```
while (true) {  
    data = c.poll()  
    c.commitOffsets()  
    p.send(data)  
}
```


Семантика «максимум один раз»

Вам это не надо

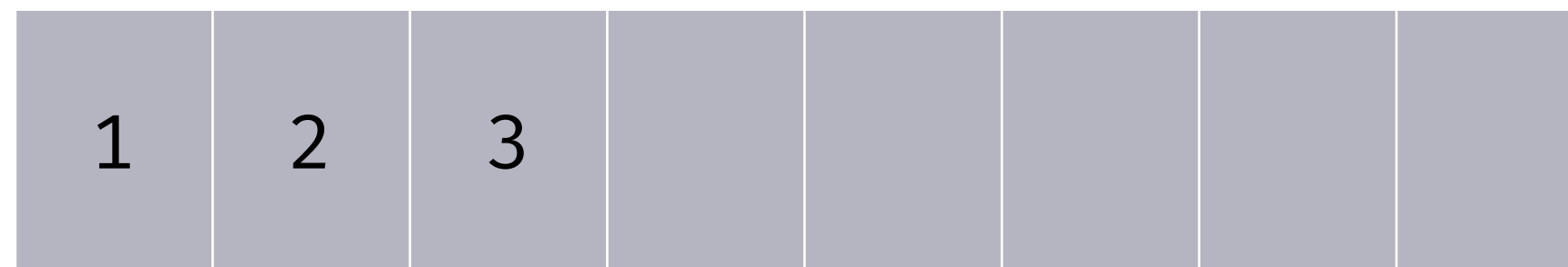
Input



Output



**Committed
Offsets**

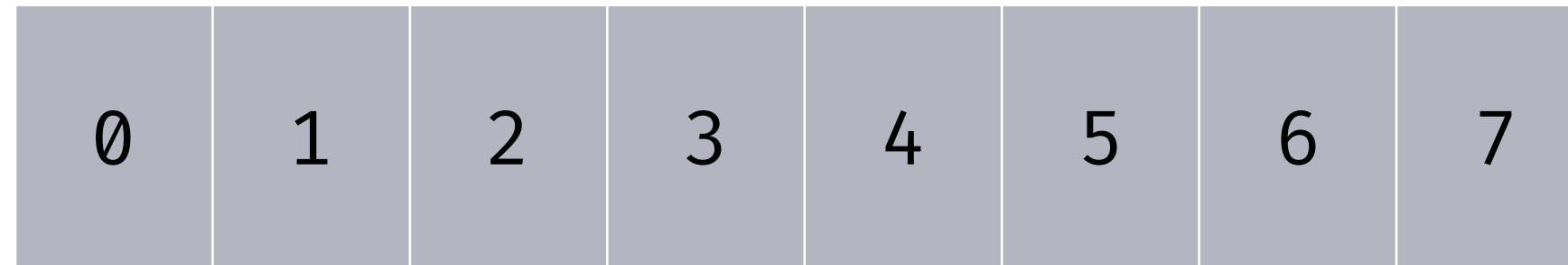


```
while (true) { Process crash!  
    data = c.poll()  
    c.commitOffsets()  
    p.send(data)  
}
```

Семантика «максимум один раз»

Вам это не надо

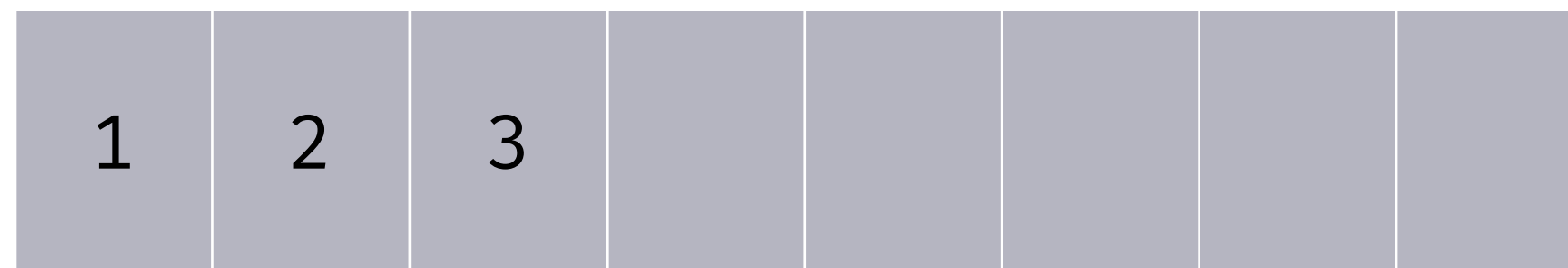
Input



Output



**Committed
Offsets**

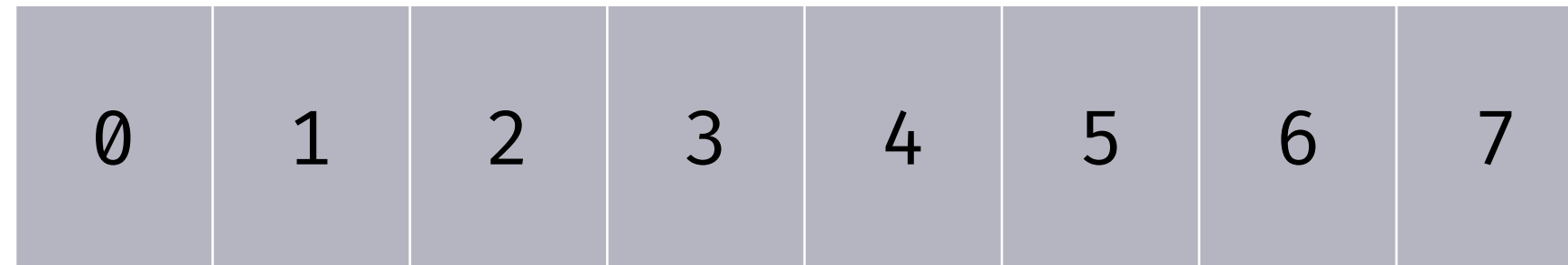


```
while (true) {  
    data = c.poll()  
    c.commitOffsets()  
    p.send(data)  
}  
Restart!
```

Семантика «максимум один раз»

Вам это не надо

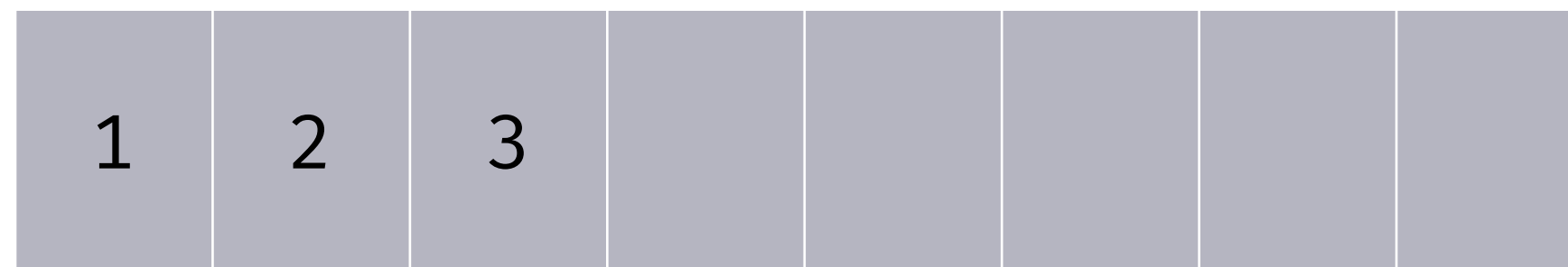
Input



Output



**Committed
Offsets**

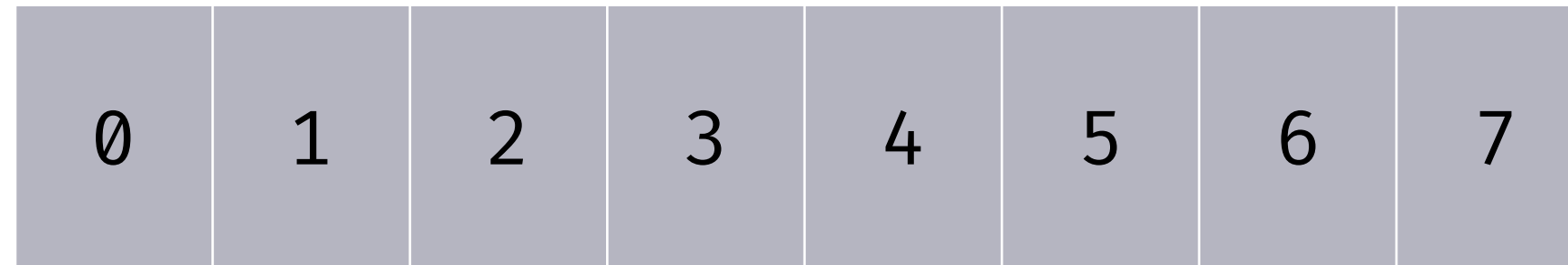


```
while (true) {  
    data = c.poll()  
    c.commitOffsets()  
    p.send(data)  
}
```

Семантика «максимум один раз»

Вам это не надо

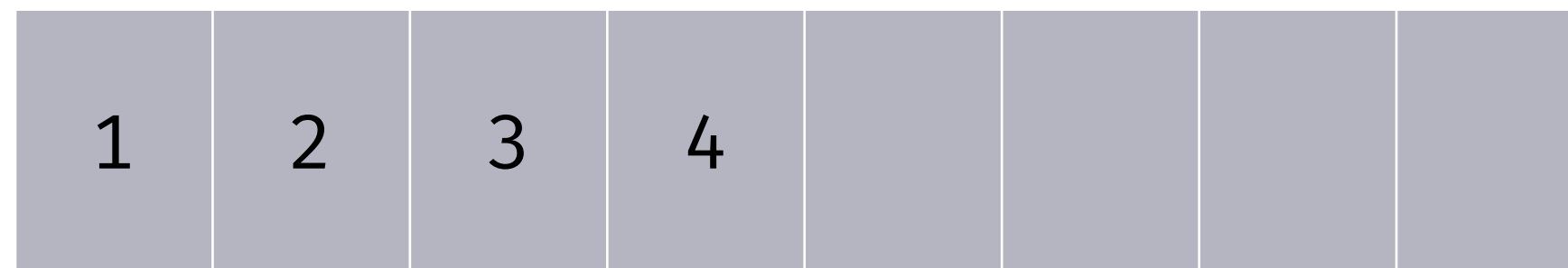
Input



Output



**Committed
Offsets**

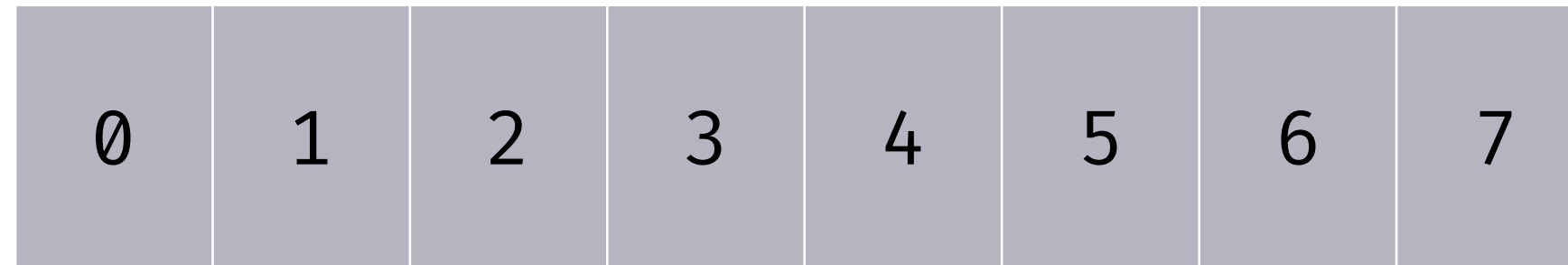


```
while (true) {  
    data = c.poll()  
    c.commitOffsets()  
    p.send(data)  
}
```

Семантика «максимум один раз»

Вам это не надо

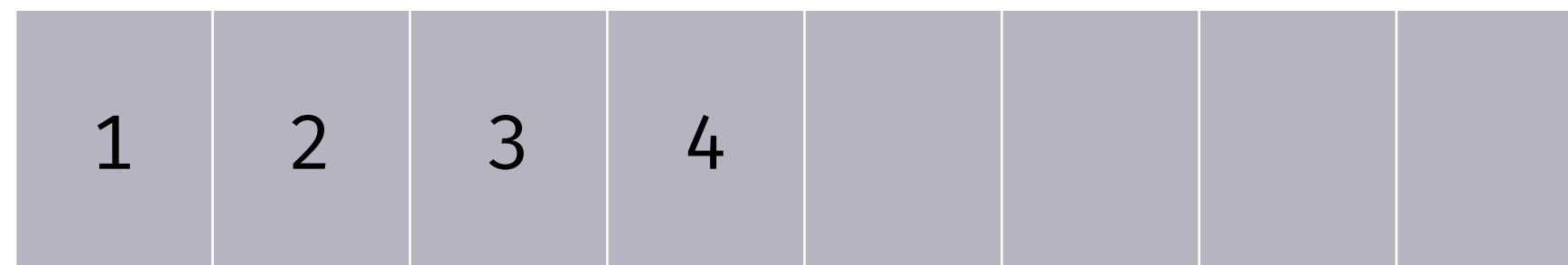
Input



Output



**Committed
Offsets**

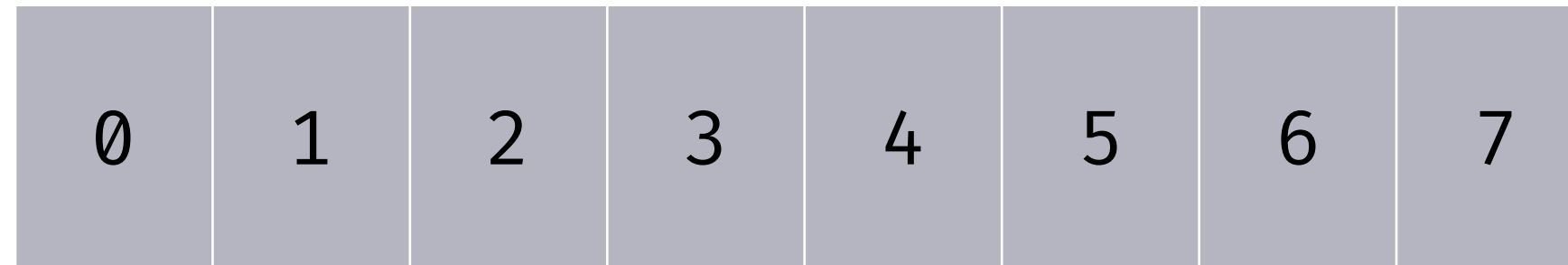


```
while (true) {  
    data = c.poll()  
    c.commitOffsets()  
    p.send(data)  
}
```

Семантика «максимум один раз»

Вам это не надо

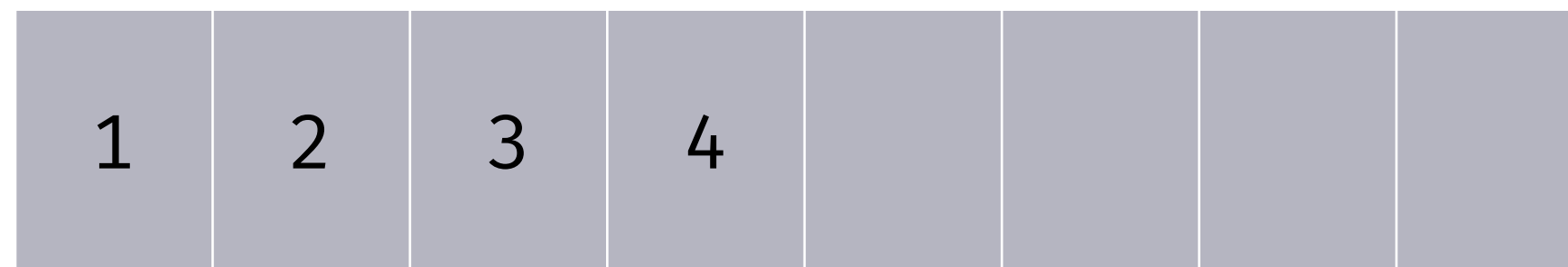
Input



Output



**Committed
Offsets**

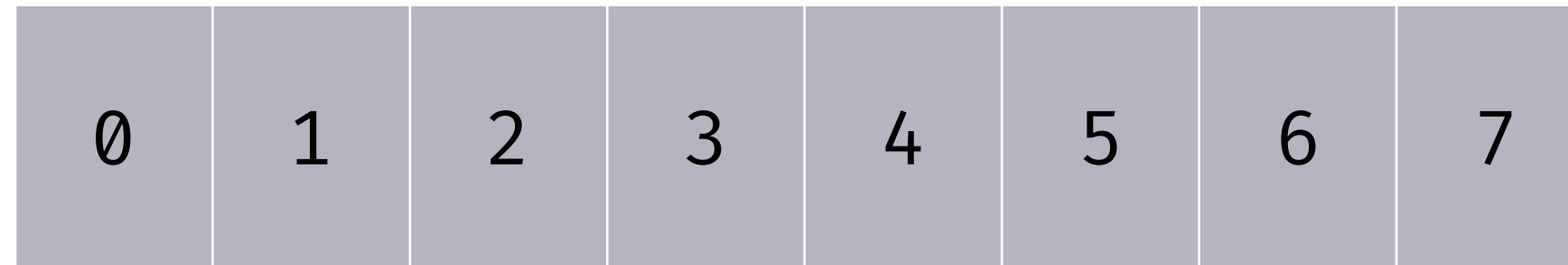


```
while (true) {  
    data = c.poll()  
    c.commitOffsets()  
    p.send(data)  
}
```

Семантика «максимум один раз»

Вам это не надо

Input



Output



**Committed
Offsets**

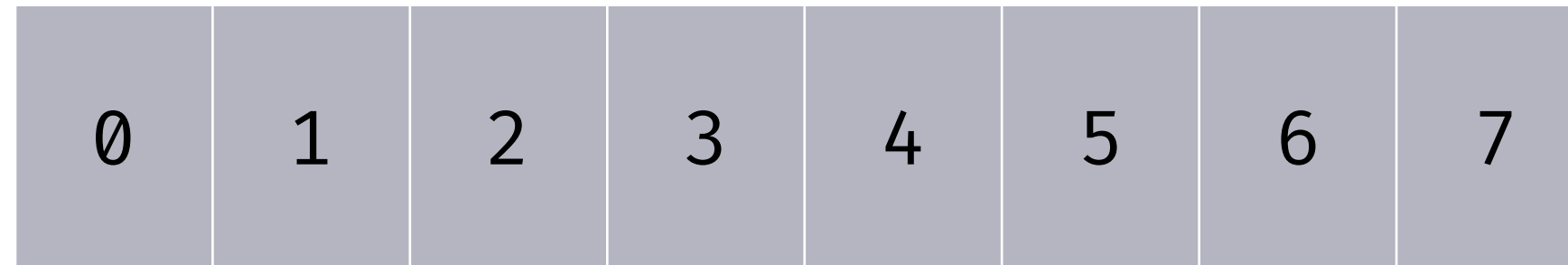


```
while (true) {  
    data = c.poll()  
    c.commitOffsets()  
    p.send(data)  
}
```

Семантика «максимум один раз»

Вам это не надо

Input



Output



**Committed
Offsets**



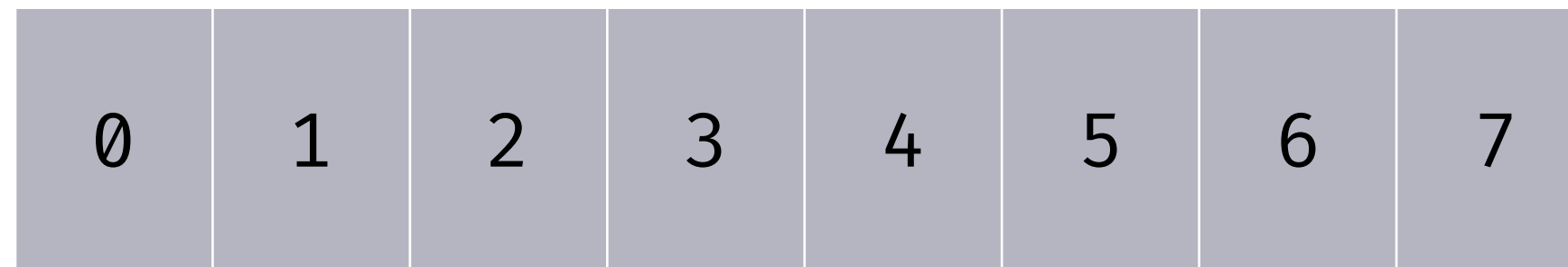
```
while (true) {  
    data = c.commits();  
    p.send();  
}
```



Семантика «максимум один раз»

Вам это не надо

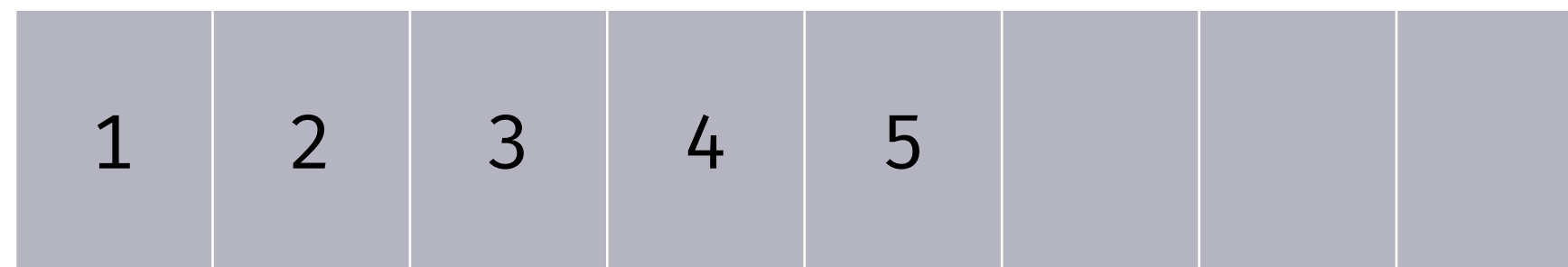
Input




Output



Committed Offsets



```
while (true) {  
  data = c.poll()  
  c.commitOffsets()  
  p.send(data)  
}  
while (true) {  
  data = c.poll()  
  c.commitOffsets()  
  p.send(data)  
}
```



Семантика «максимум один раз»

Вам это не надо

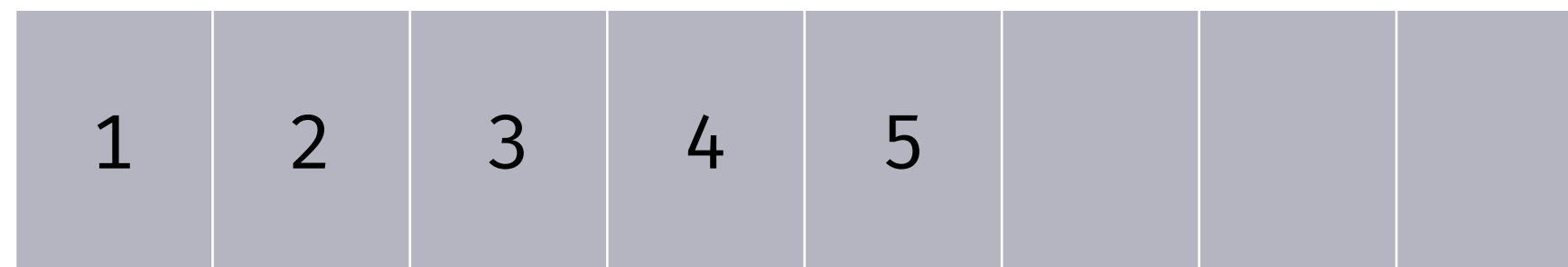
Input



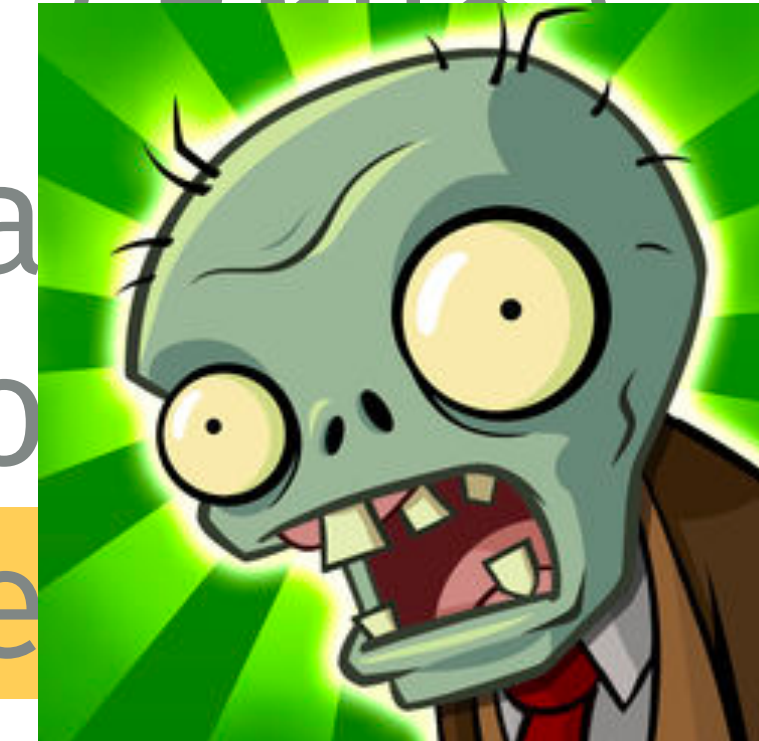
Output



**Committed
Offsets**



```
while (true) {  
  data = c.poll()  
  c.commitOffsets()  
  p.send(data)  
}
```

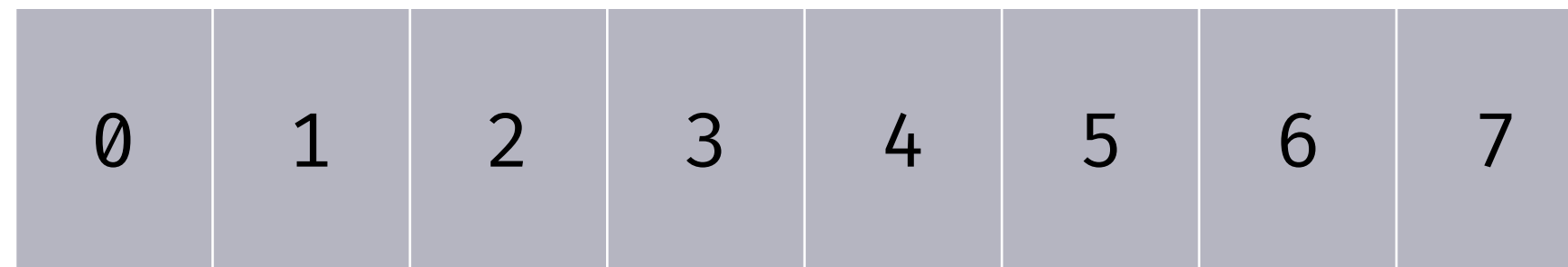


```
while (true) {  
  data = c.poll()  
  c.commitOffsets()  
  p.send(data)  
}
```

Семантика «максимум один раз»

Вам это не надо

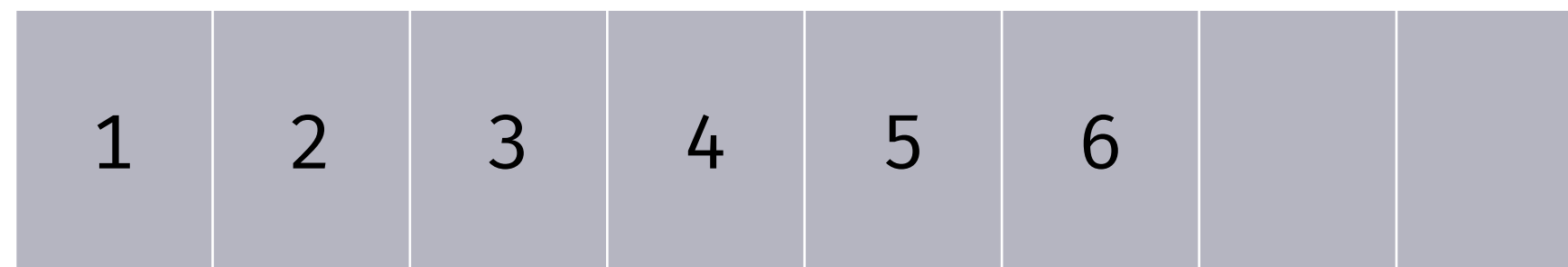
Input



Output



**Committed
Offsets**



```
while (true) {  
  data = c.poll()  
  c.commitOffsets()  
  p.send(data)  
}
```



```
while (true) {  
  data = c.poll()  
  c.commitOffsets()  
  p.send(data)  
}
```

Семантика «максимум один раз»

Вам это не надо

Input

0	1	2	3	4	5	6	7
---	---	---	---	---	---	---	---

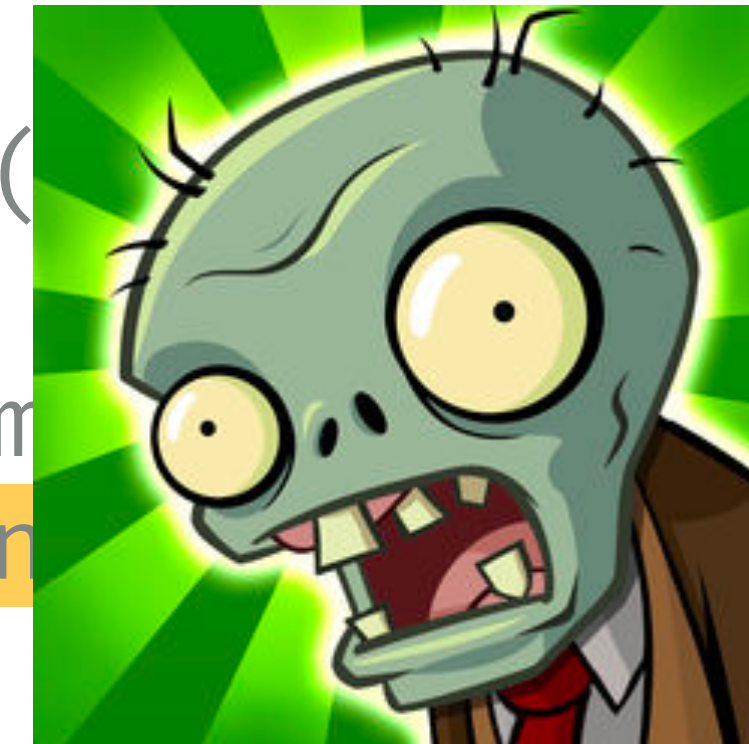
Output

0	3	5					
---	---	---	--	--	--	--	--

**Committed
Offsets**

1	2	3	4	5	6		
---	---	---	---	---	---	--	--

```
while (  
  data  
  c.com  
  p.send  
)  
{  
}
```

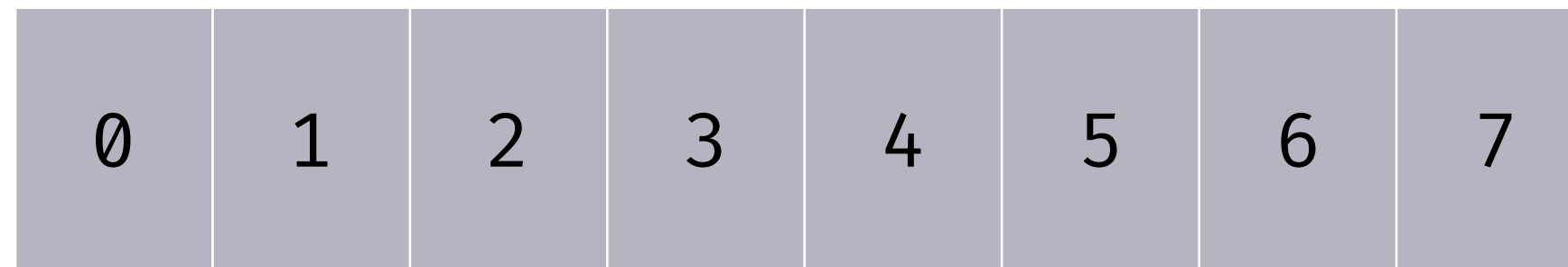


```
while (true) {  
  data = c.poll()  
  c.commitOffsets()  
  p.send(data)  
}
```

Семантика «максимум один раз»

Вам это не надо

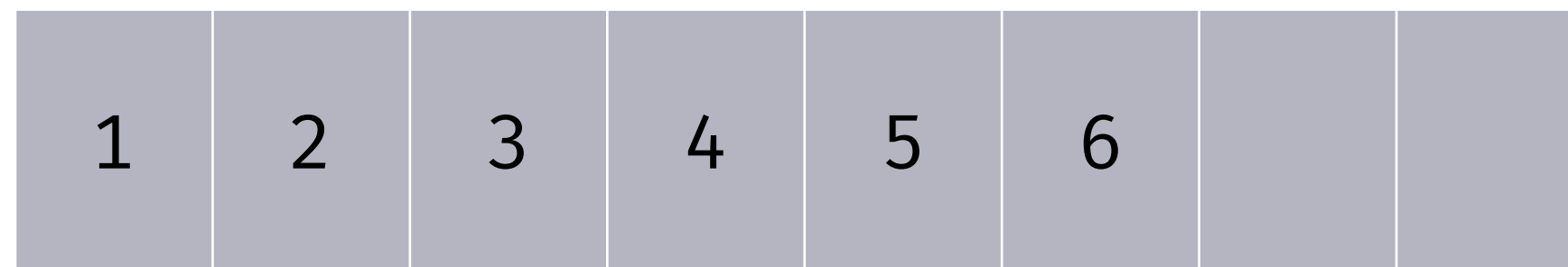
Input



Output



**Committed
Offsets**



```
while (true) {  
    data = c.poll()  
    c.commitOffsets()  
    p.send(data)  
}  
  
while (true) {  
    data = c.poll()  
    c.commitOffsets()  
    p.send(data)  
}
```

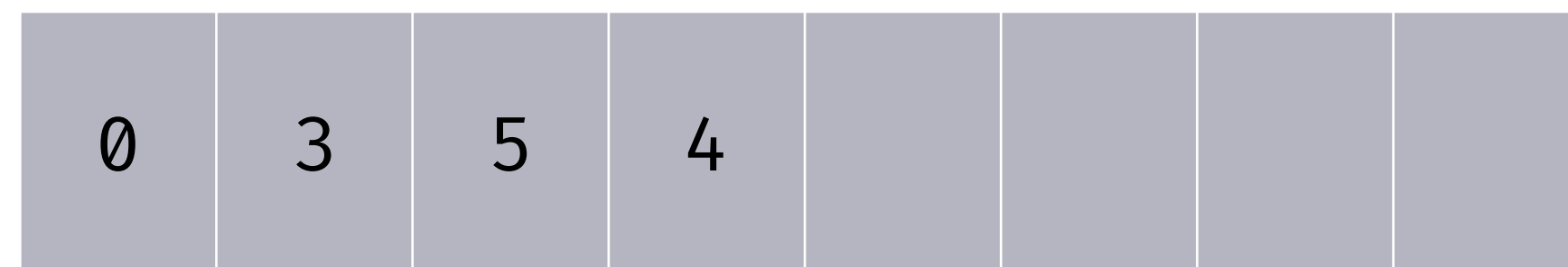
Семантика «максимум один раз»

Вам это не надо

Input



Output



**Committed
Offsets**

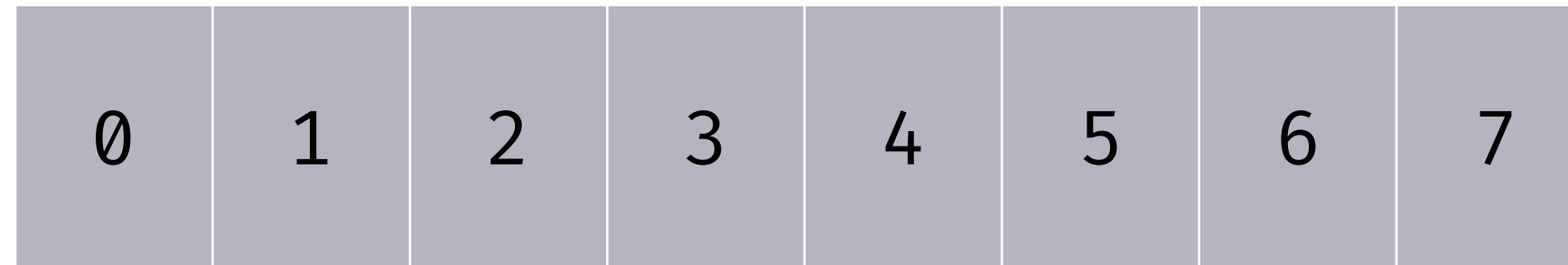


```
while (true) {  
    data = c.poll()  
    c.commitOffsets()  
    p.send(data)  
}  
  
while (true) {  
    data = c.poll()  
    c.commitOffsets()  
    p.send(data)  
}
```

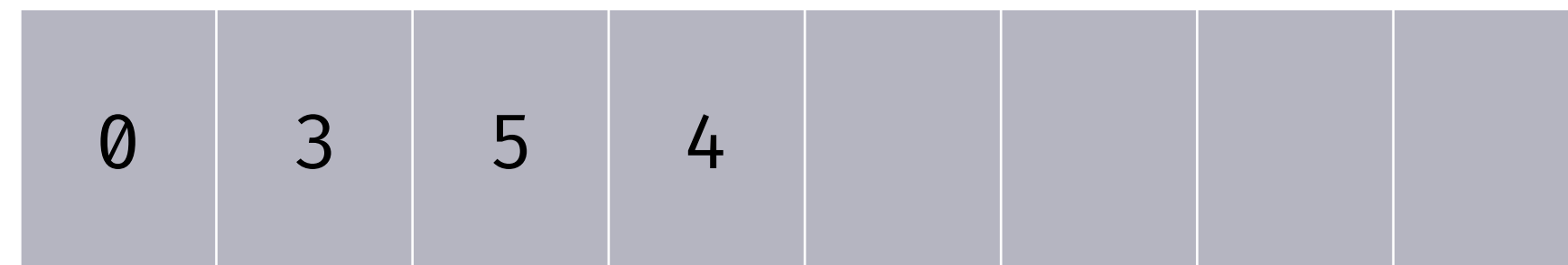
Семантика «максимум один раз»

Вам это не надо

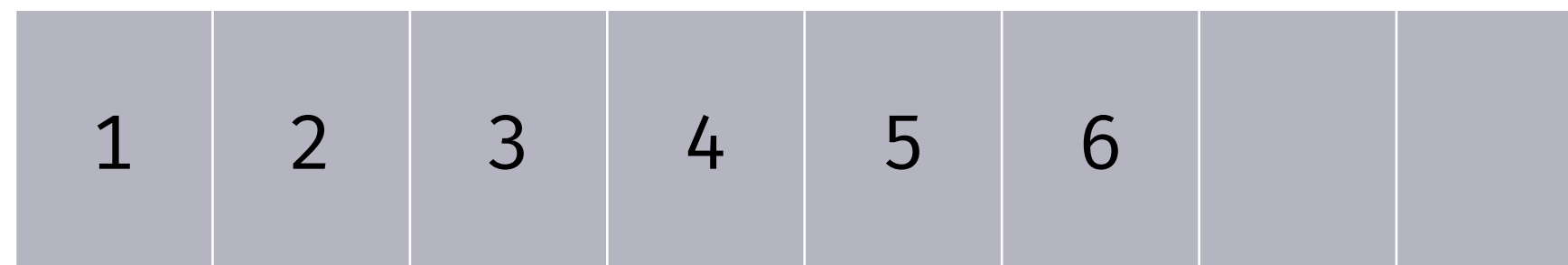
Input



Output



**Committed
Offsets**



- Потери сообщений
- Повторы продьюсера не разрешены
- Сбои перед тем как записываем данные

- Зомби могут
сообщений



порядок

Слабости семантики

- Повторные попытки продюсера небезопасны
- Данные не пишутся атомарно с их смещениями (offsets)
- Нет защиты от зомби

Слабости семантики

- Повторные попытки продюсера небезопасны
- Данные не пишутся атомарно с их смещениями (offsets)
- Нет защиты от зомби

Идемпотентный прогюсер

- `ProducerId` определяет экземпляр прогюсера
- `Sequence` служит для дедупликации
- `enable.idempotence=true`

Log Message

Offset

Key

Value

Timestamp

Headers

ProducerId

Sequence

Идемпотентный прогюсер

- `ProducerId` определяет экземпляр прогюсера
- `Sequence` служит для дедупликации

Client

Broker

Log

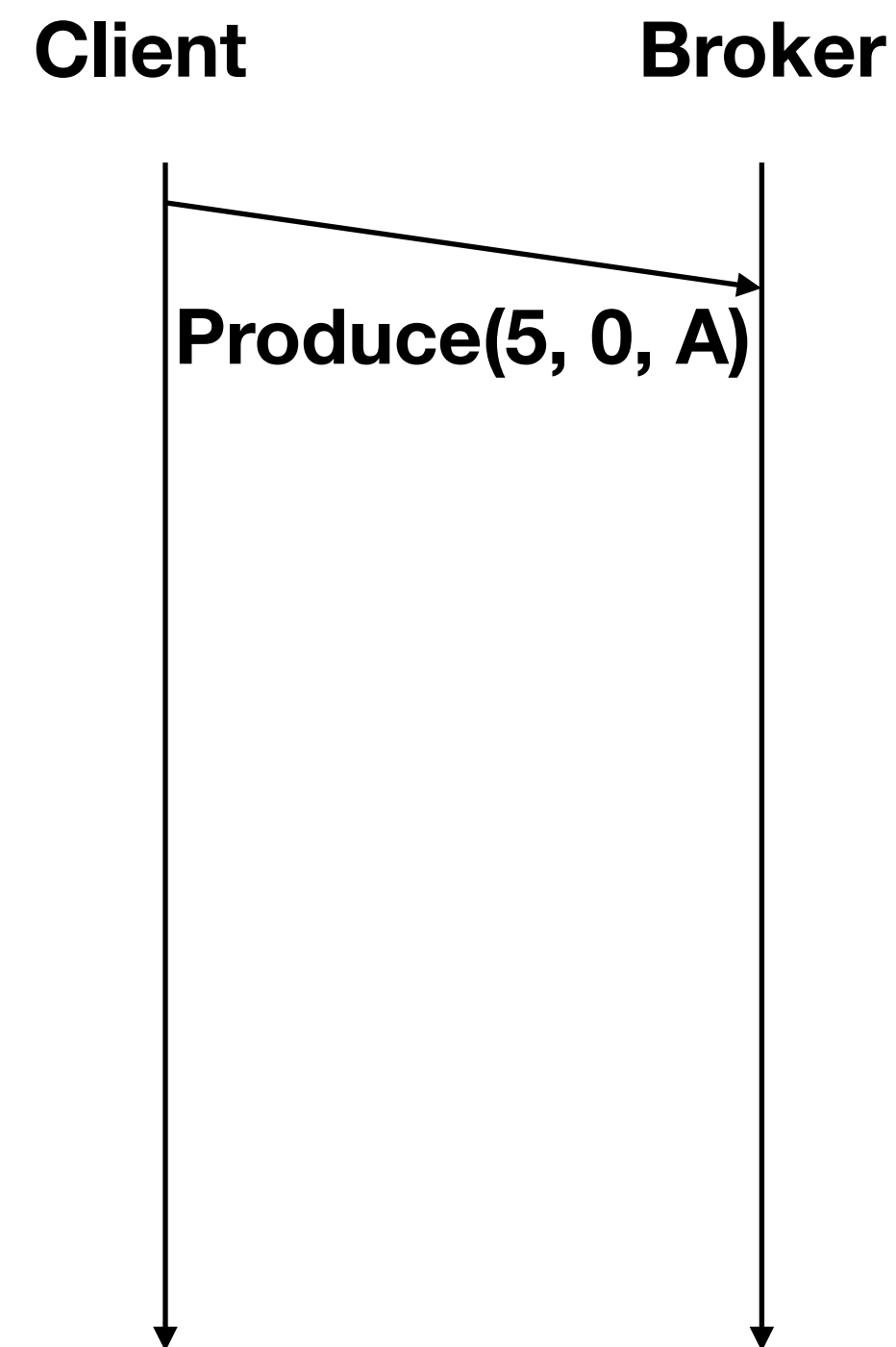
pid	seq	key

Produce(pid, seq, key)

Append(pid, seq, key)

Идемпотентный прогюсер

- `ProducerId` определяет экземпляр прогюсера
- `Sequence` служит для дедупликации



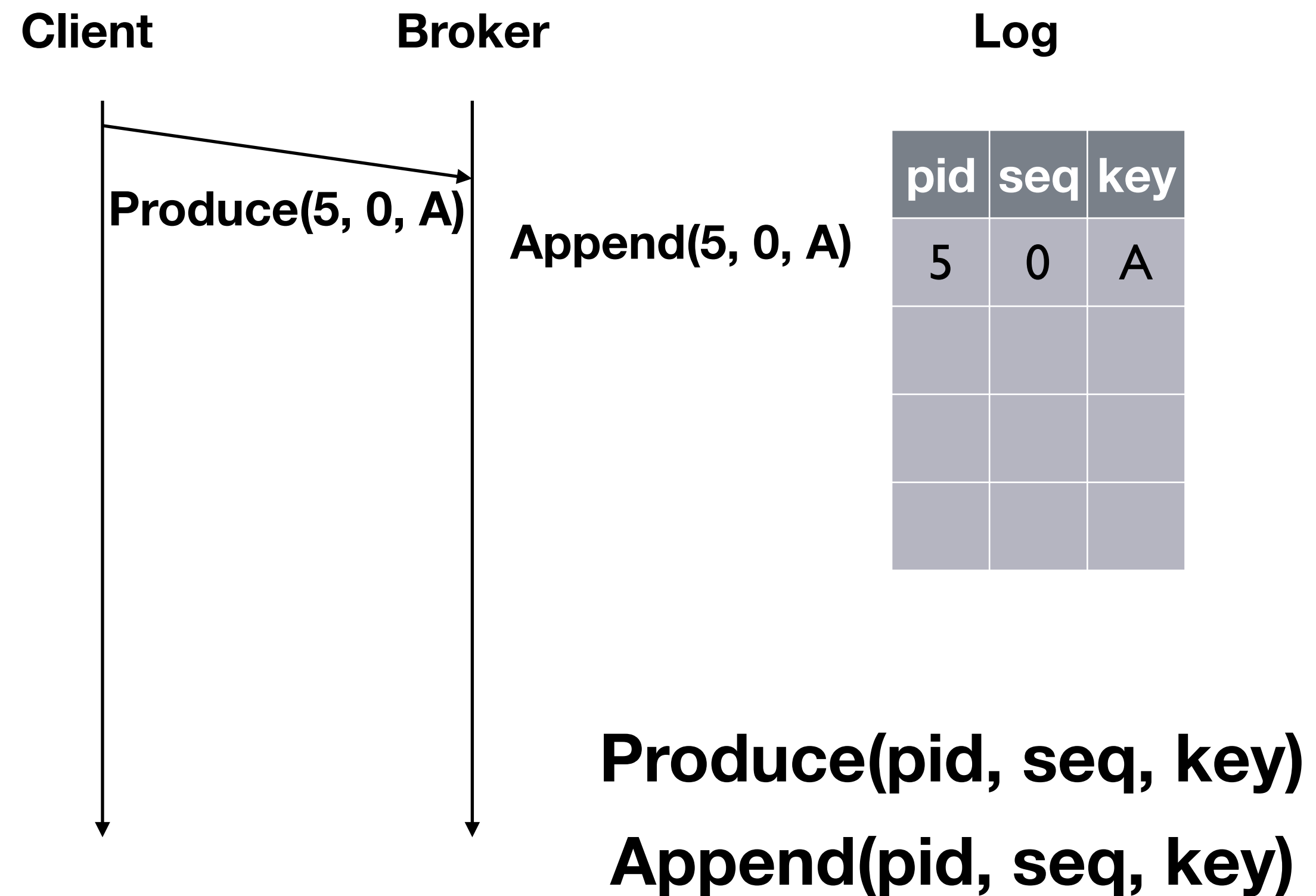
Log

pid	seq	key

`Produce(pid, seq, key)`
`Append(pid, seq, key)`

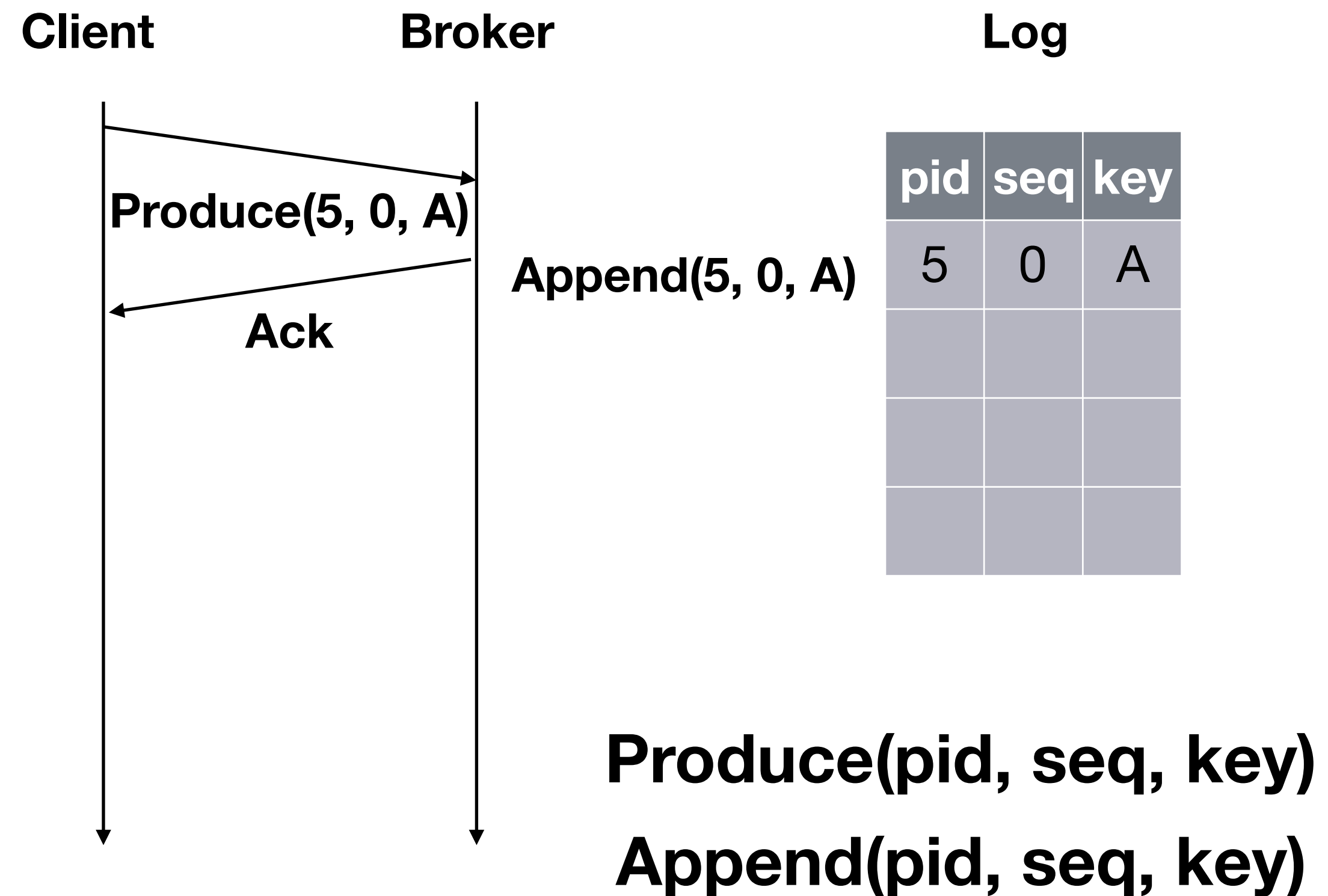
Идемпотентный прогюсер

- `ProducerId` определяет экземпляр прогюсера
- `Sequence` служит для дедупликации



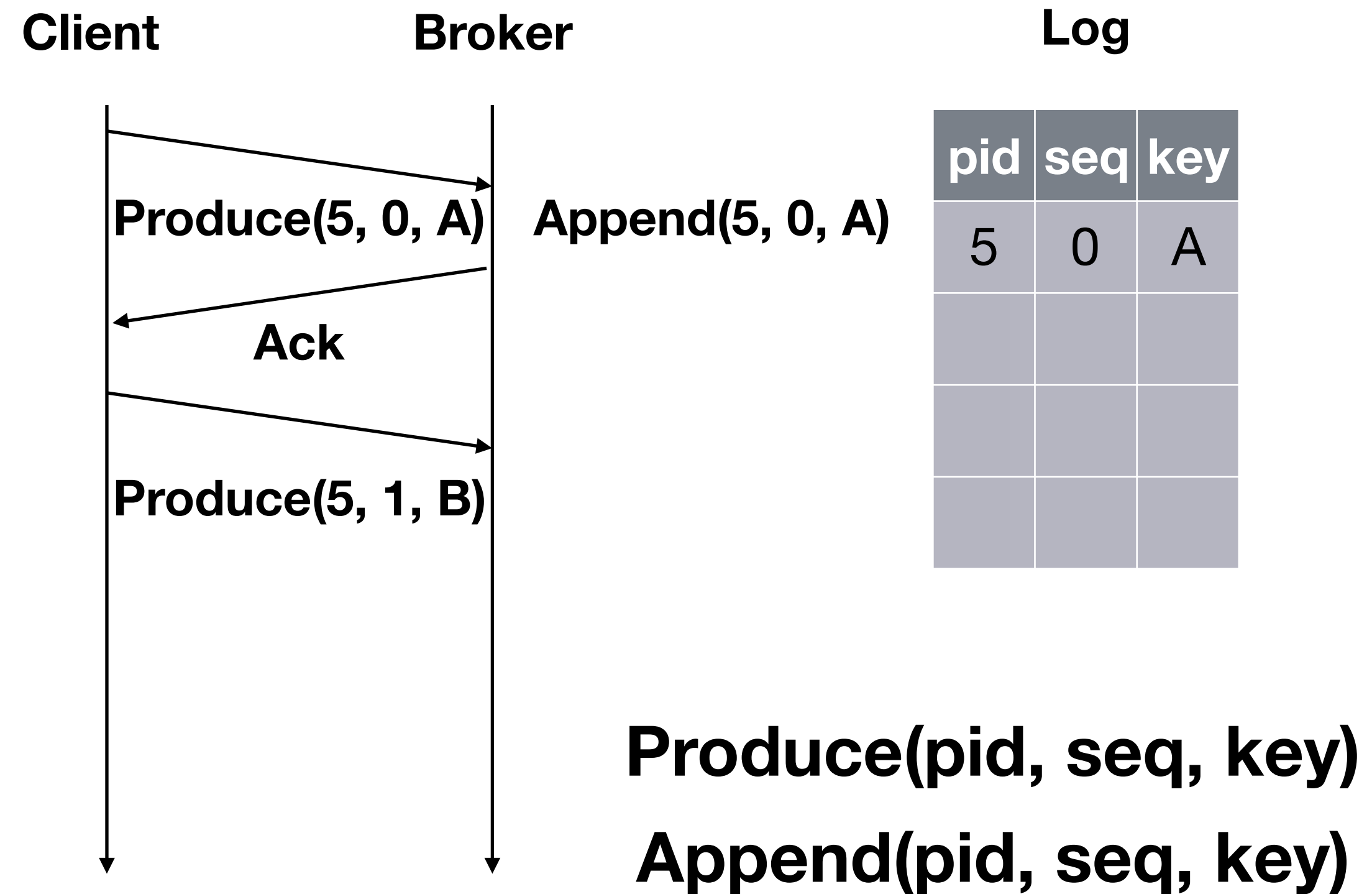
Идемпотентный прогюсер

- `ProducerId` определяет экземпляр прогюсера
- `Sequence` служит для дедупликации



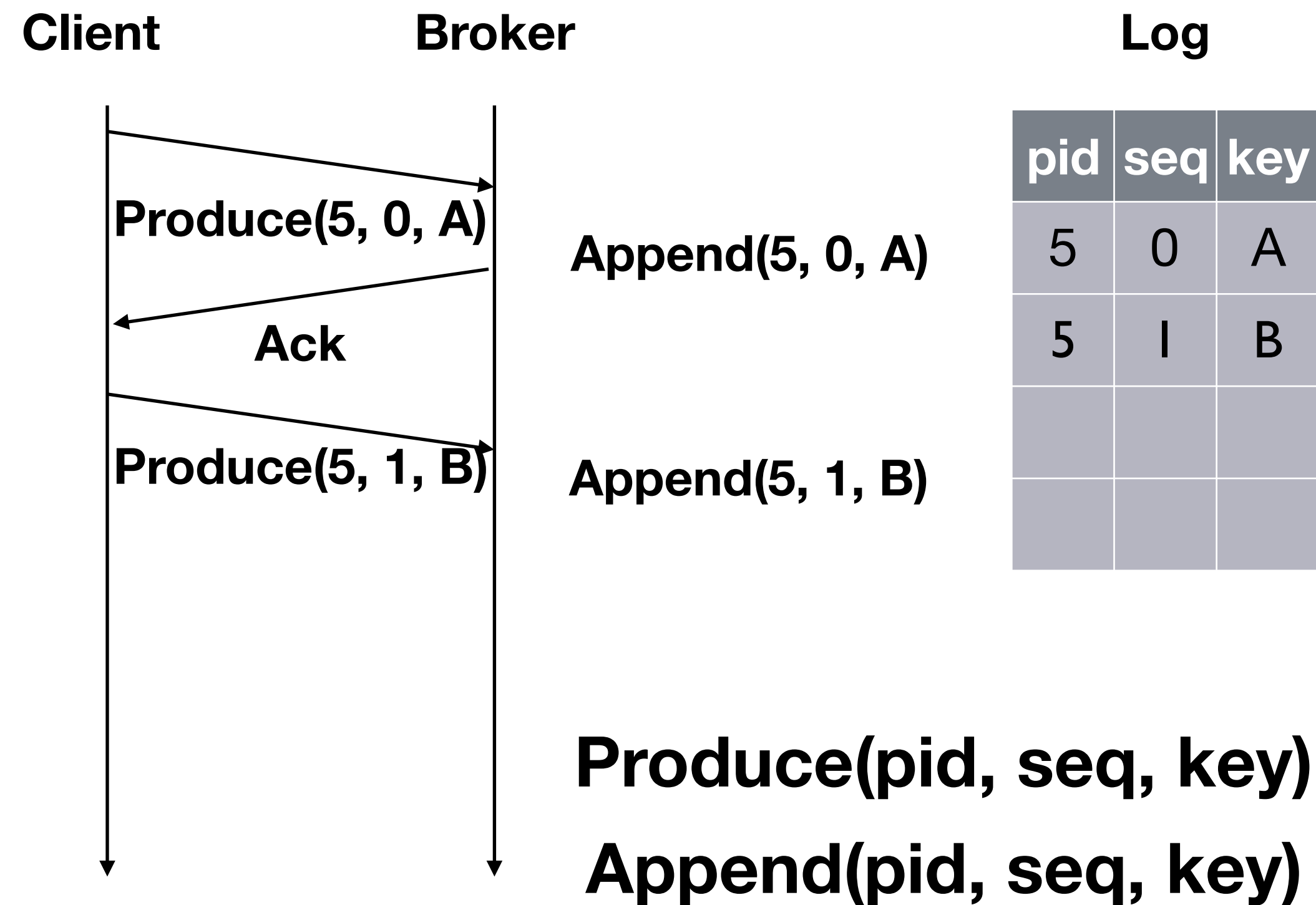
Идемпотентный прогюсер

- `ProducerId` определяет экземпляр прогюсера
- `Sequence` служит для дедупликации



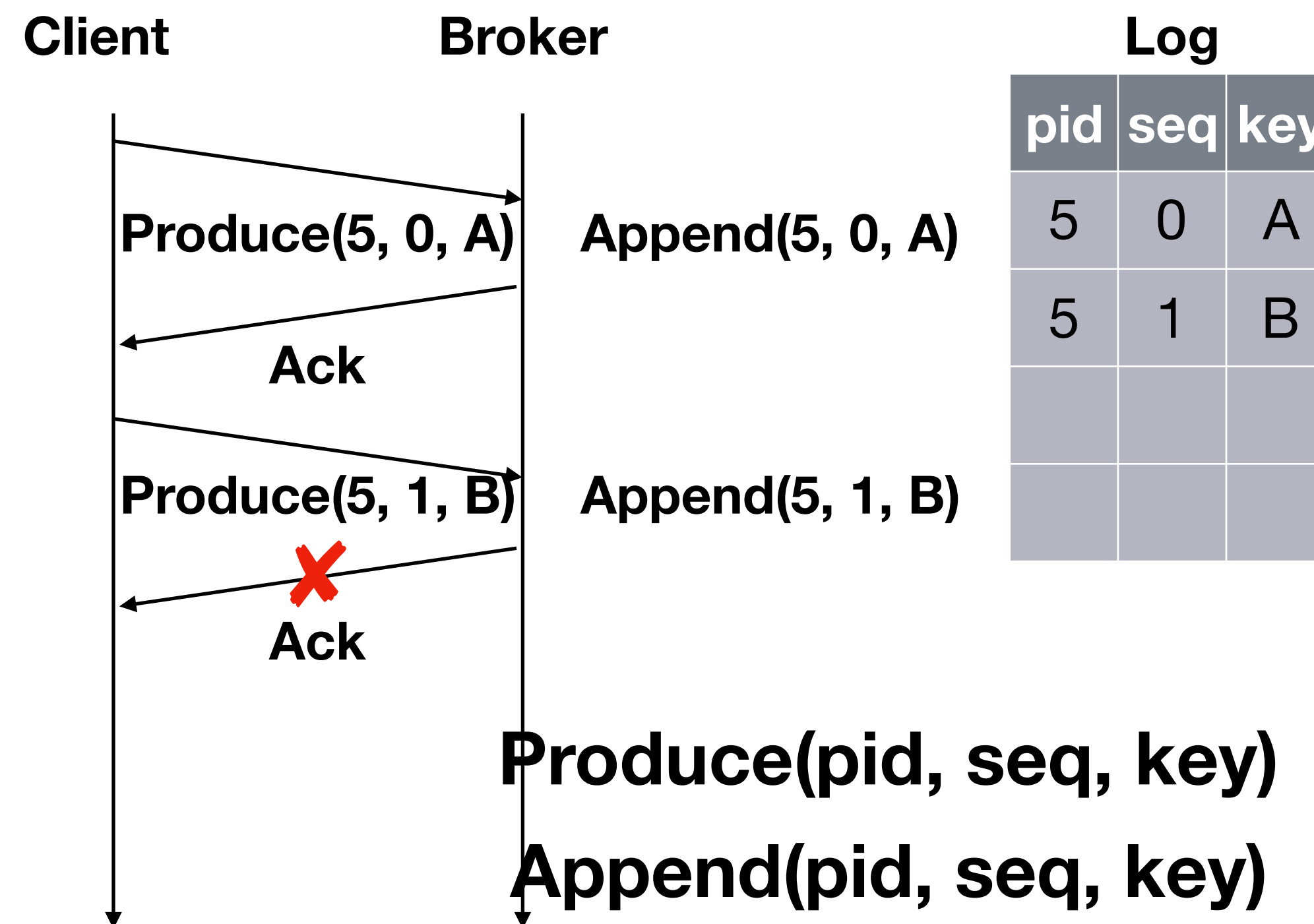
Идемпотентный прогюсер

- `ProducerId` определяет экземпляр прогюсера
- `Sequence` служит для дедупликации



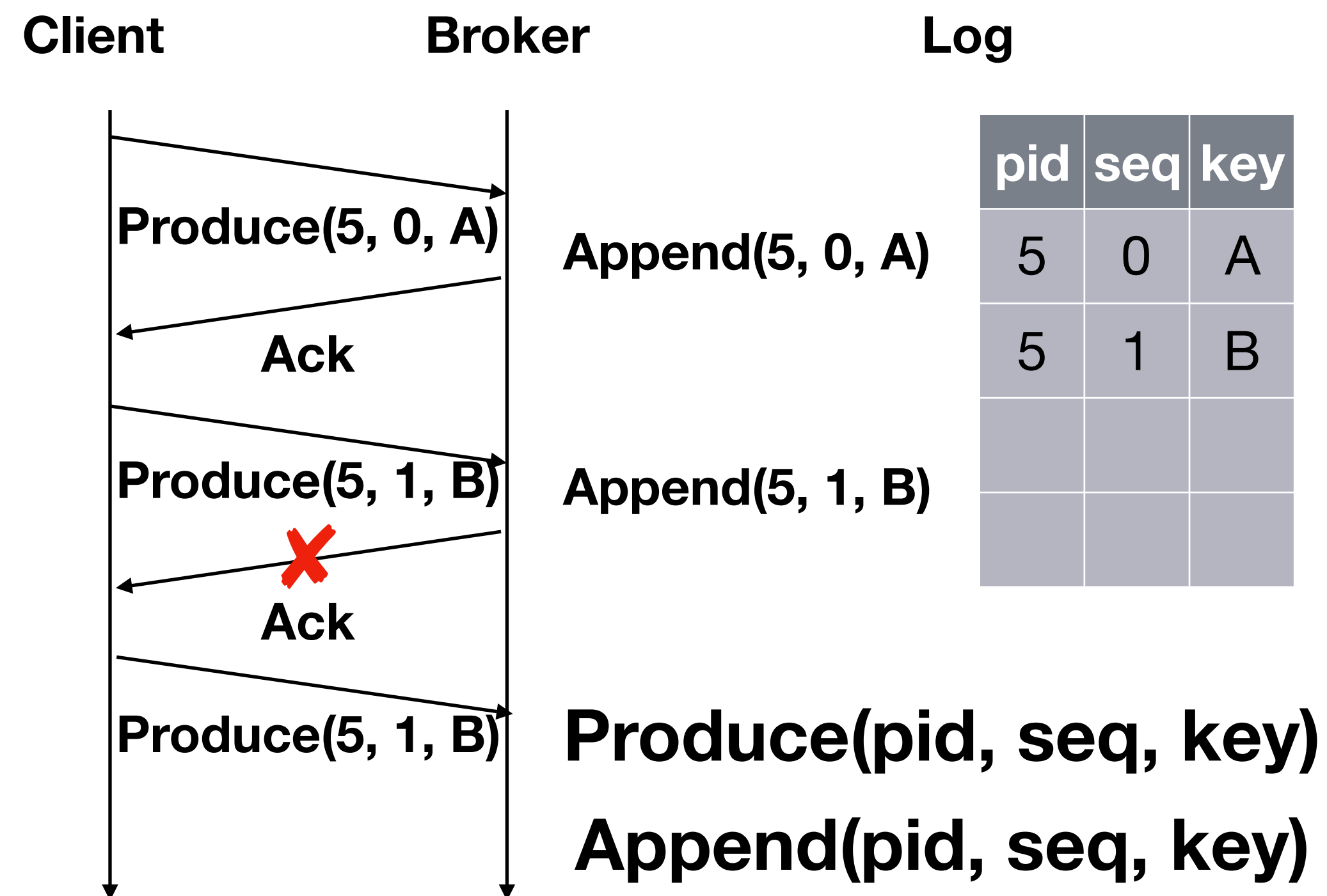
Идемпотентный прогюсер

- `ProducerId` определяет экземпляр прогюсера
- `Sequence` служит для дедупликации



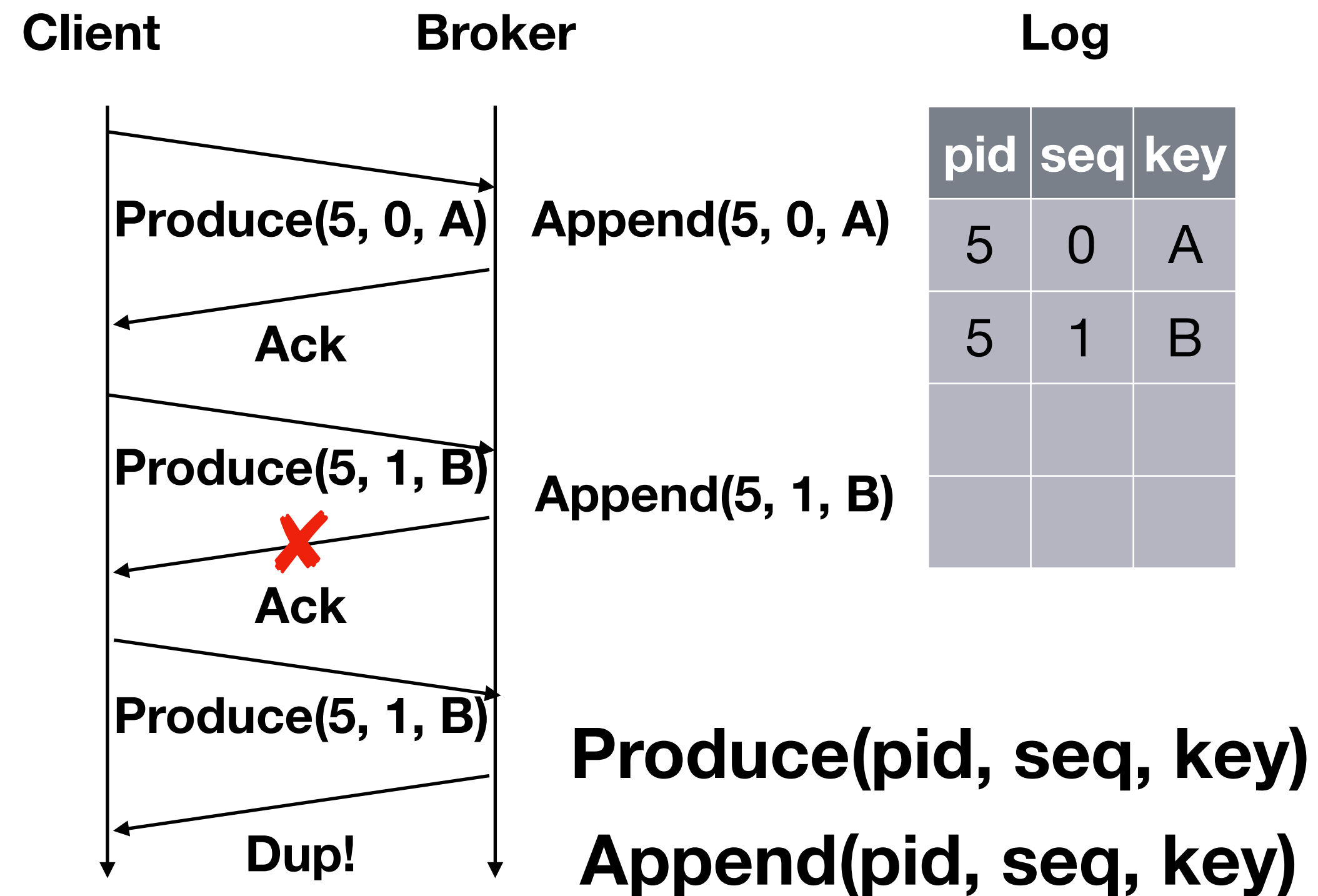
Идемпотентный прогюсер

- `ProducerId` определяет экземпляр прогюсера
- `Sequence` служит для дедупликации



Идемпотентный прогюсер

- `ProducerId` определяет экземпляр прогюсера
- `Sequence` служит для дедупликации



Слабости семантики

- Повторные попытки продюсера небезопасны
- Данные не пишутся атомарно с их смещениями (offsets)
- Нет защиты от зомби

Слабости семантики

- ~~Повторные попытки прогюсера небезопасны~~
- Данные не пишутся атомарно с их смещениями (offsets)
- Нет защиты от зомби

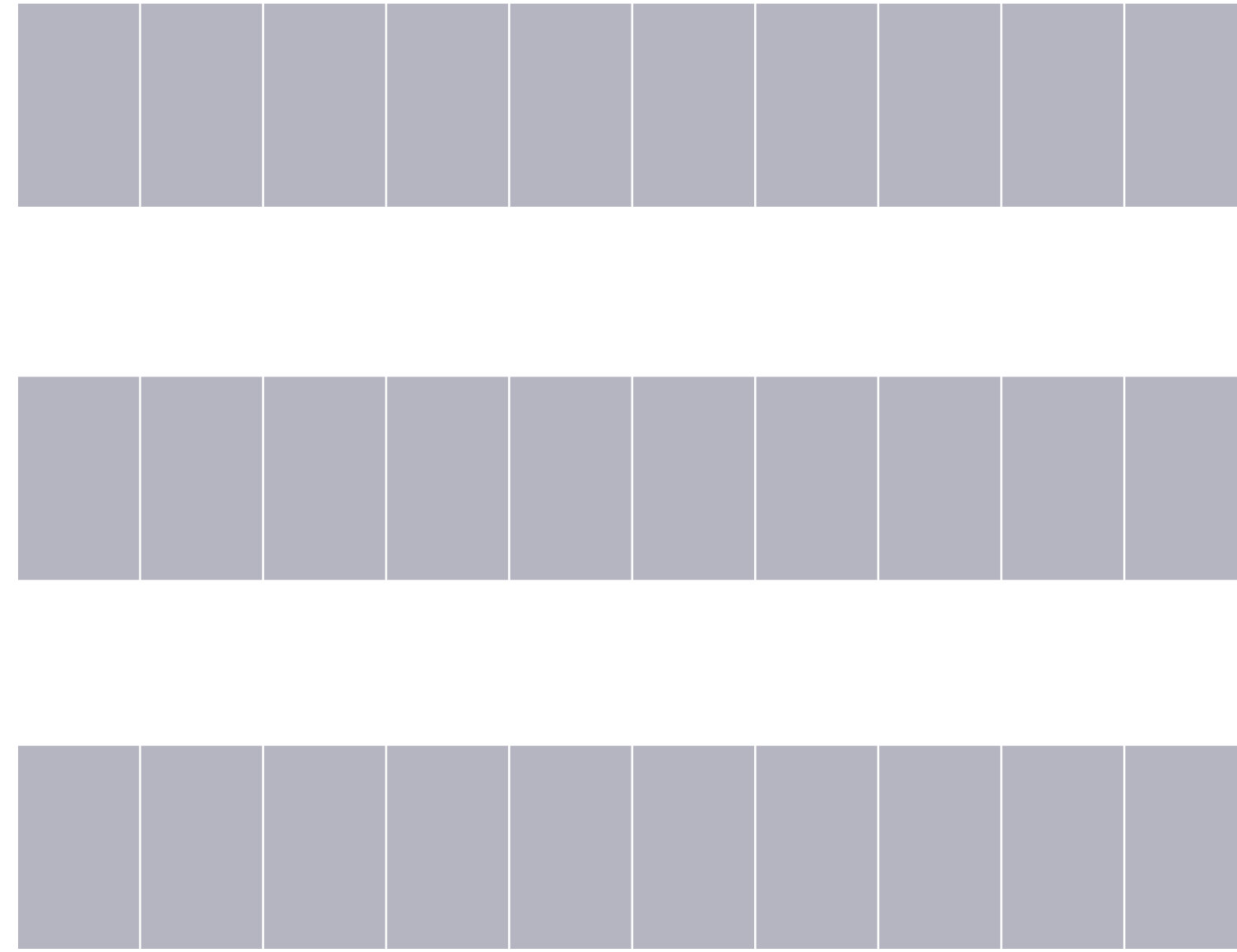
Слабости семантики

- ~~Повторные попытки прогюсера небезопасны~~
- Данные не пишутся атомарно с их смещениями (offsets)
- Нет защиты от зомби

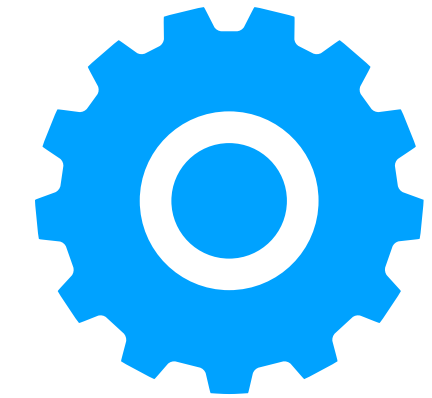
Снэгцшоты Чэгдү - Лэмпорт



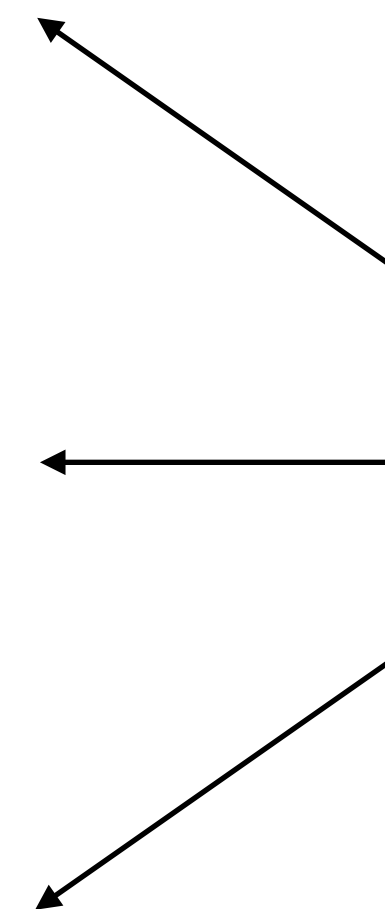
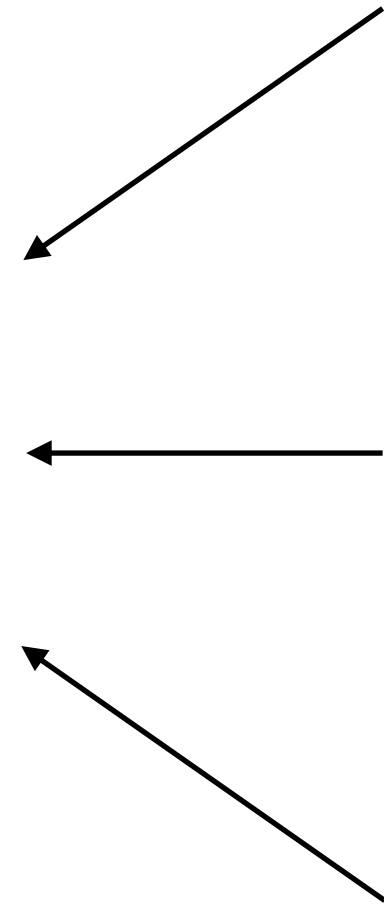
Observer



Output topics

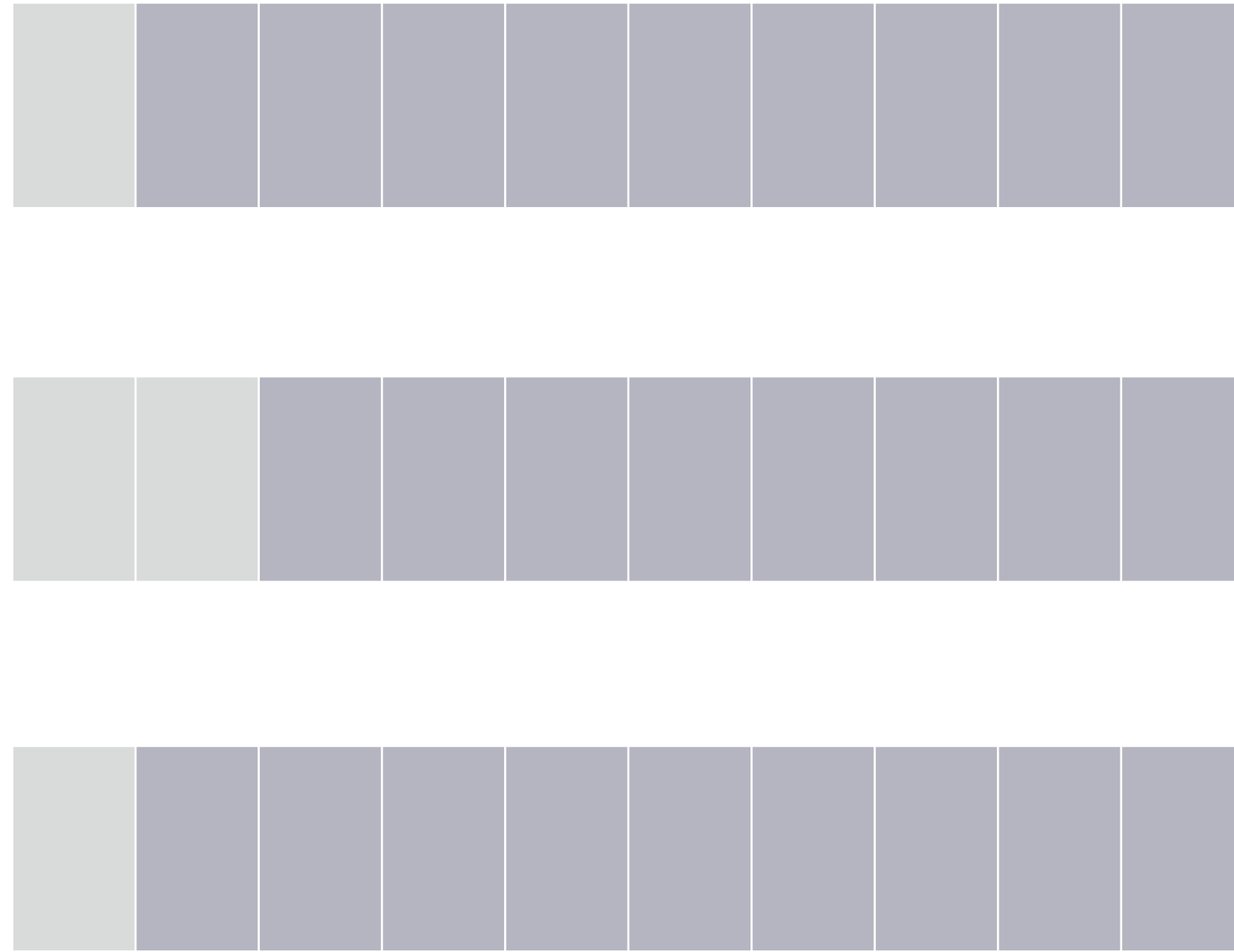
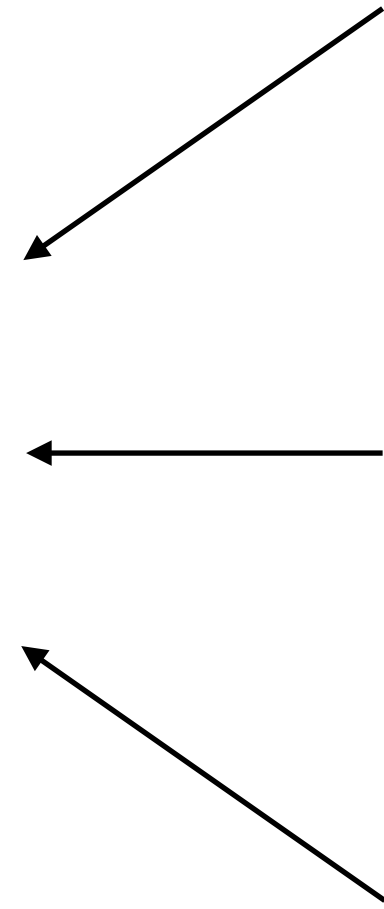


Processor

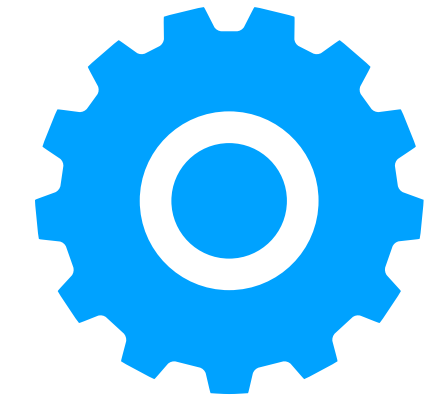
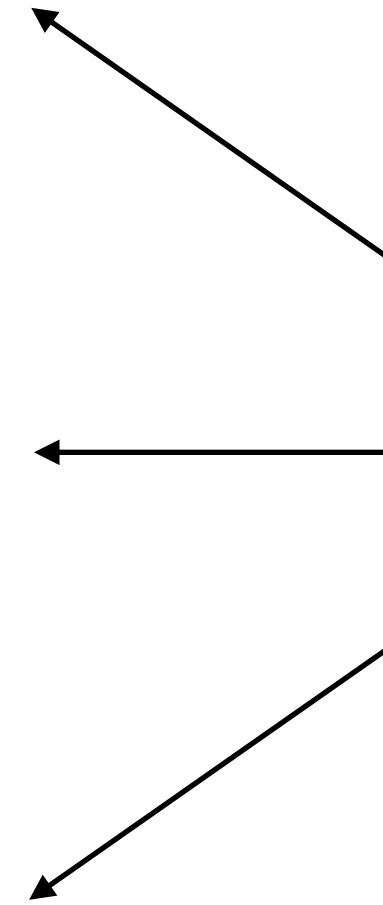




Observer



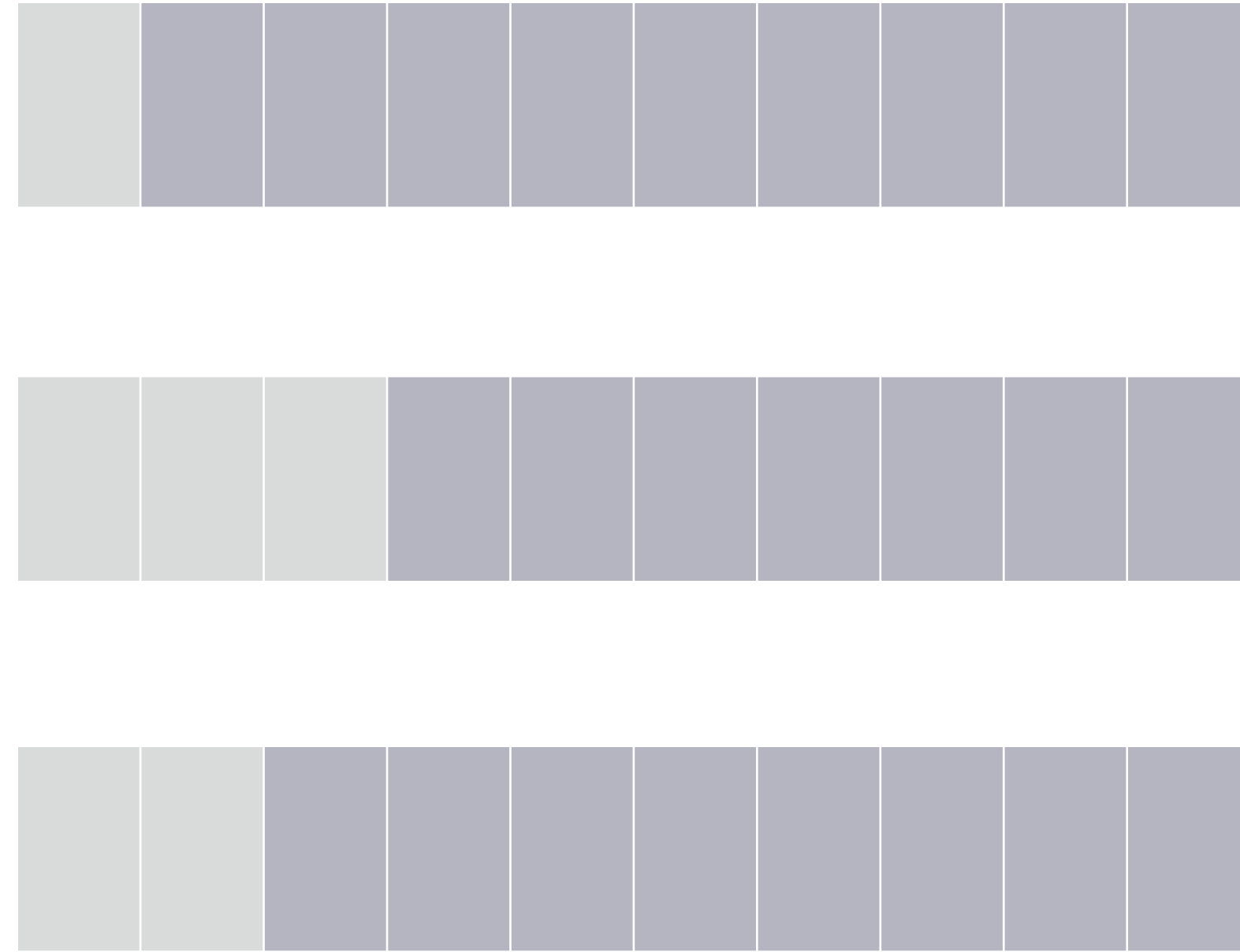
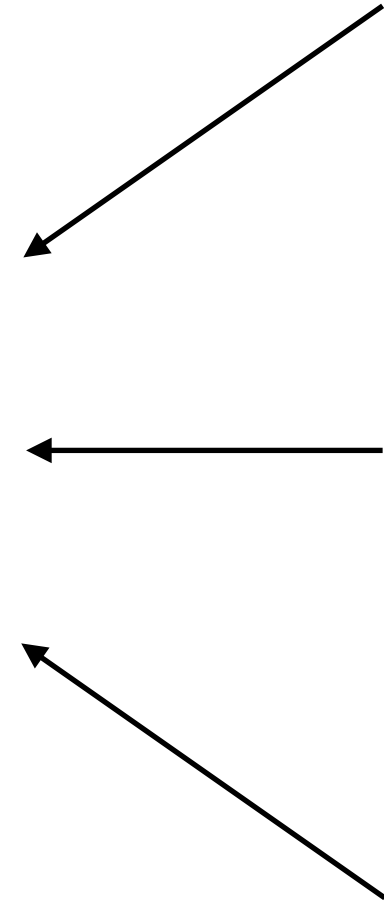
Output topics



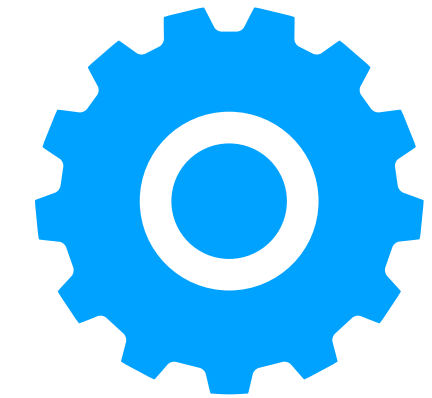
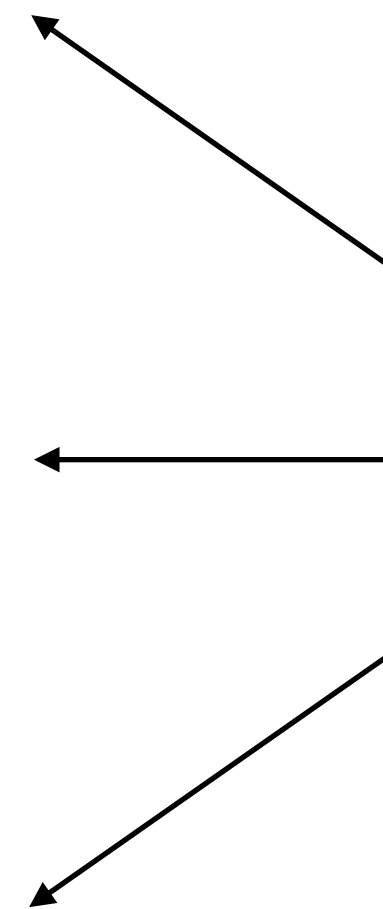
Processor



Observer



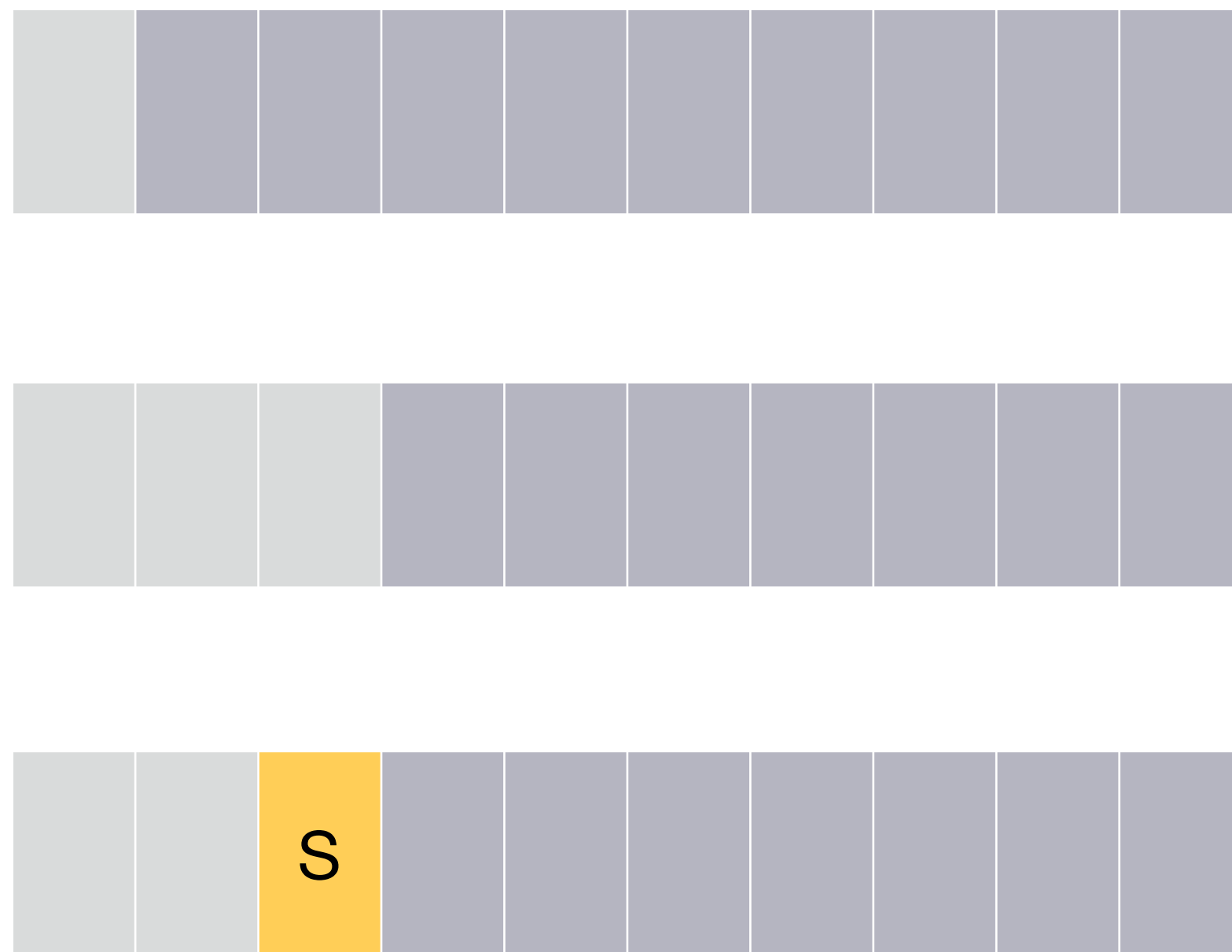
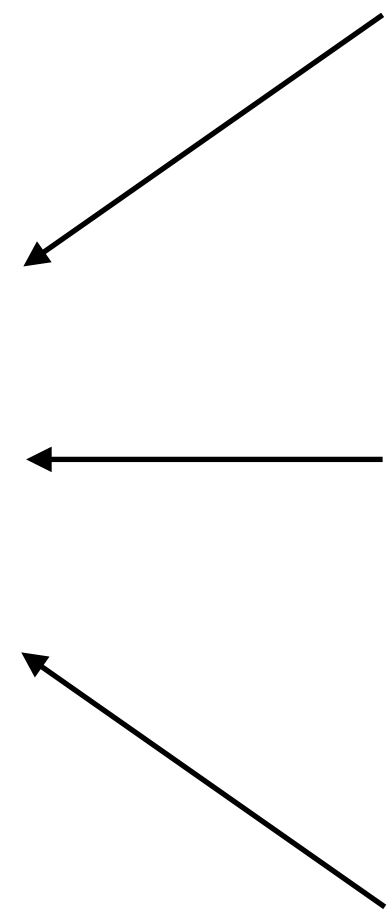
Output topics



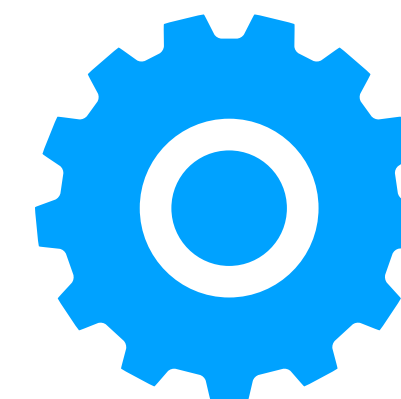
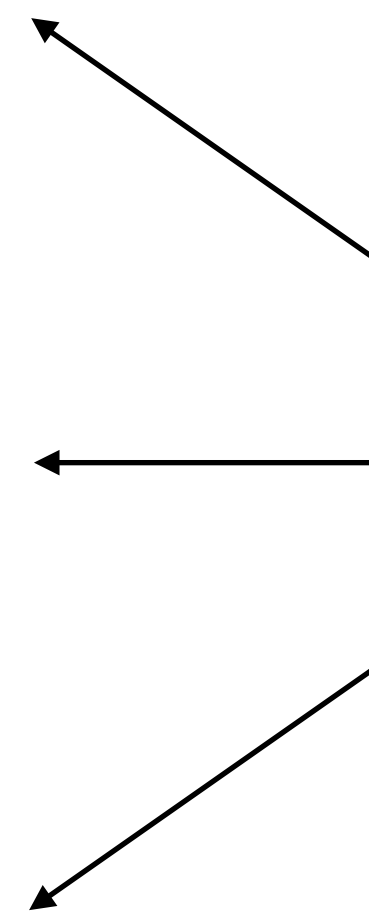
Processor



Observer



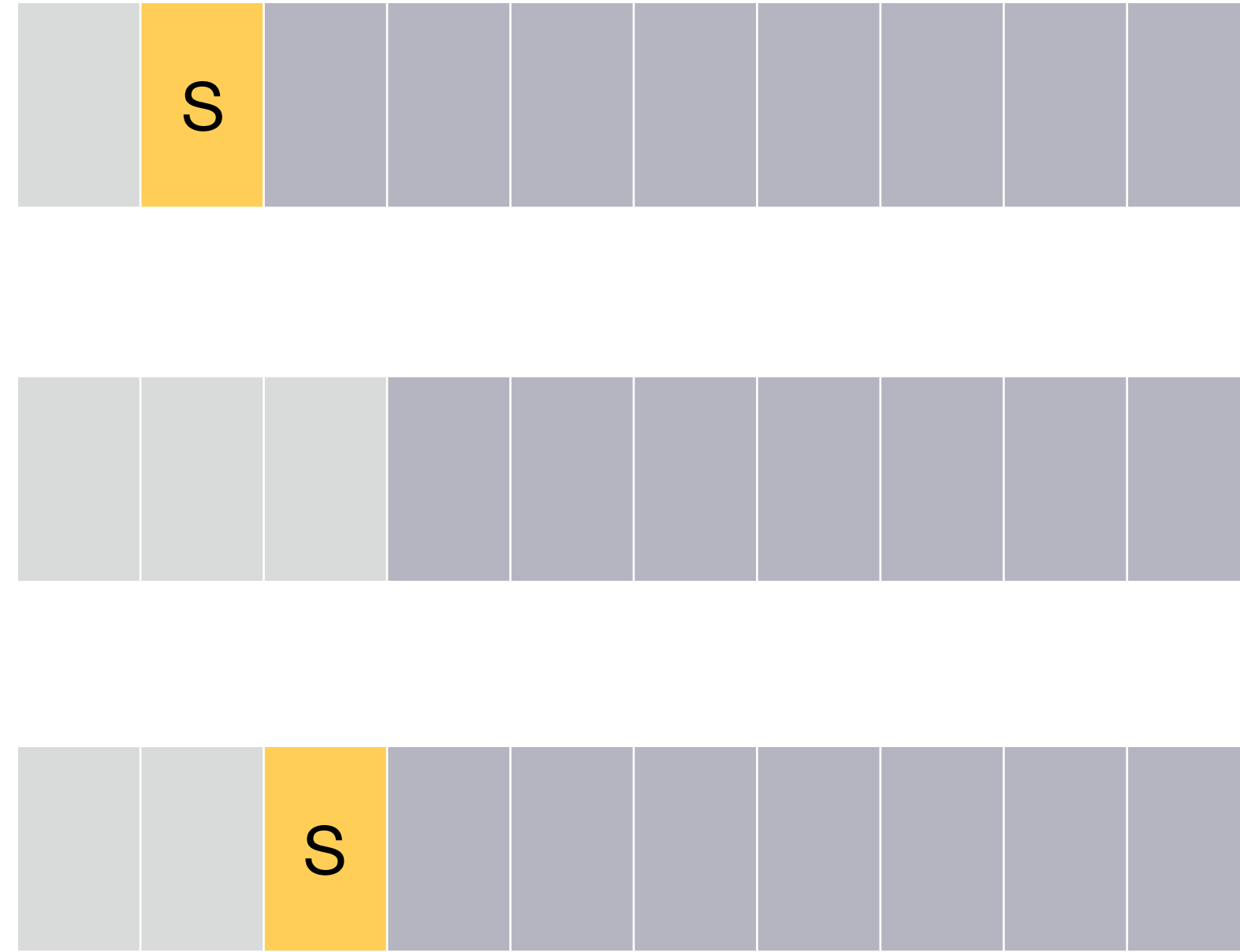
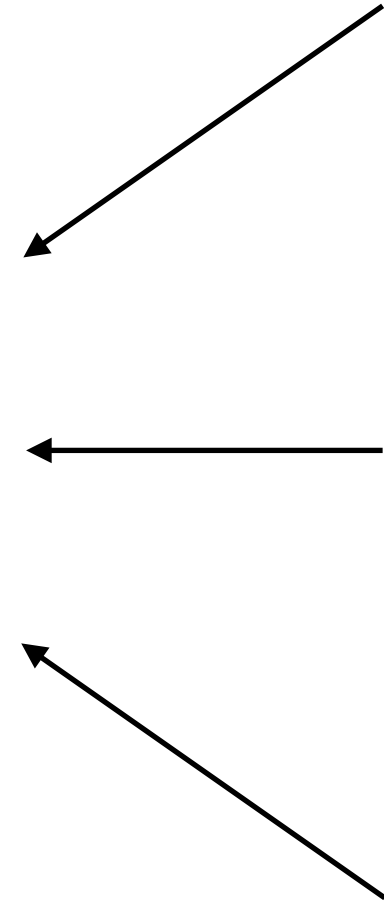
Output topics



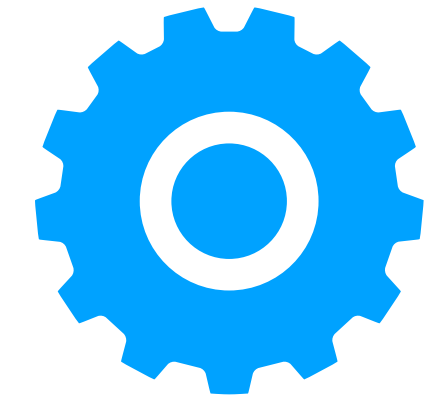
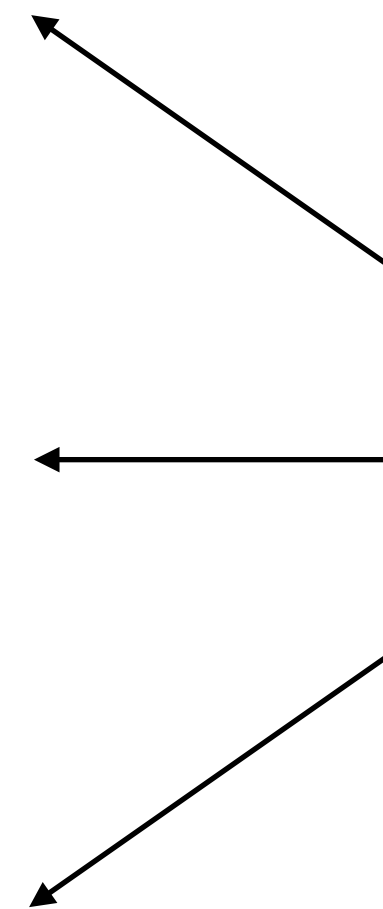
Processor



Observer



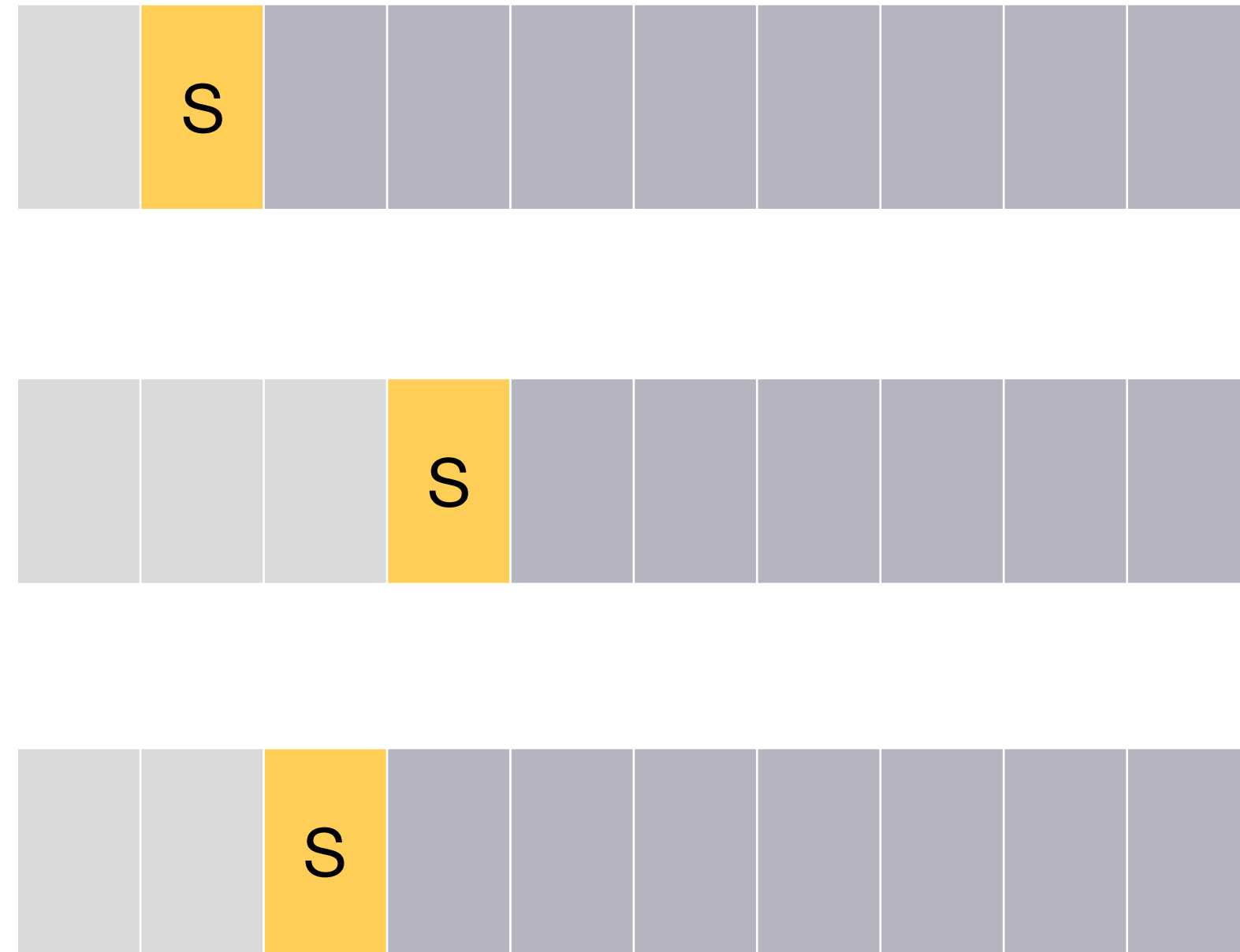
Output topics



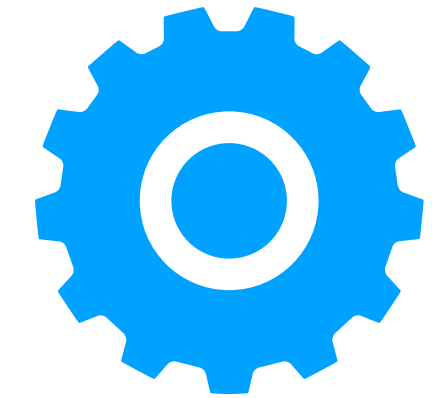
Processor



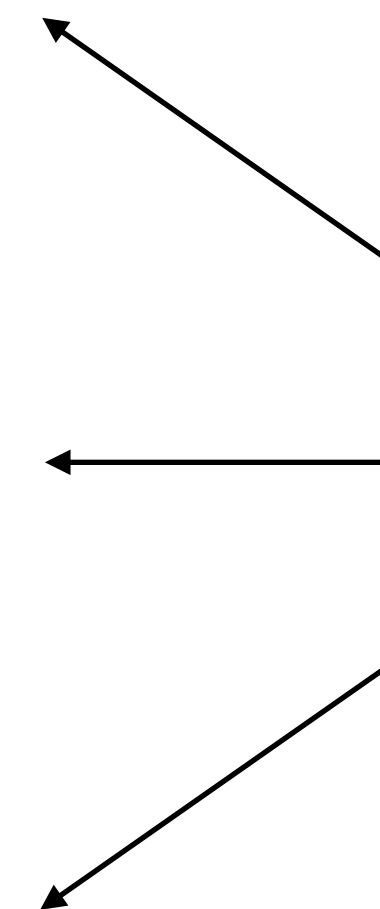
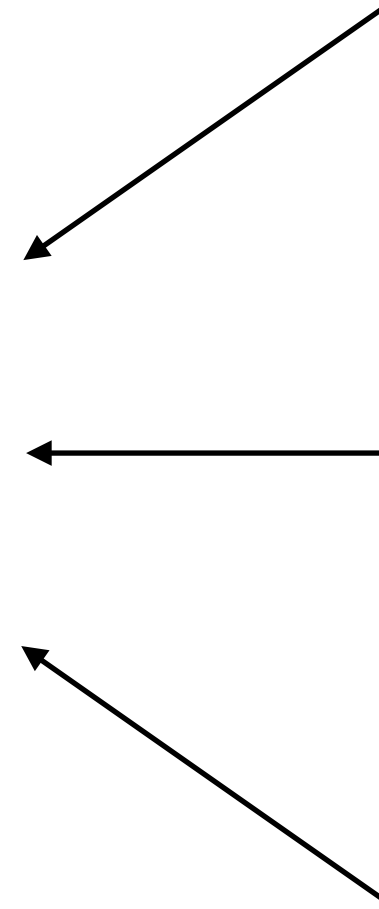
Observer



Output topics

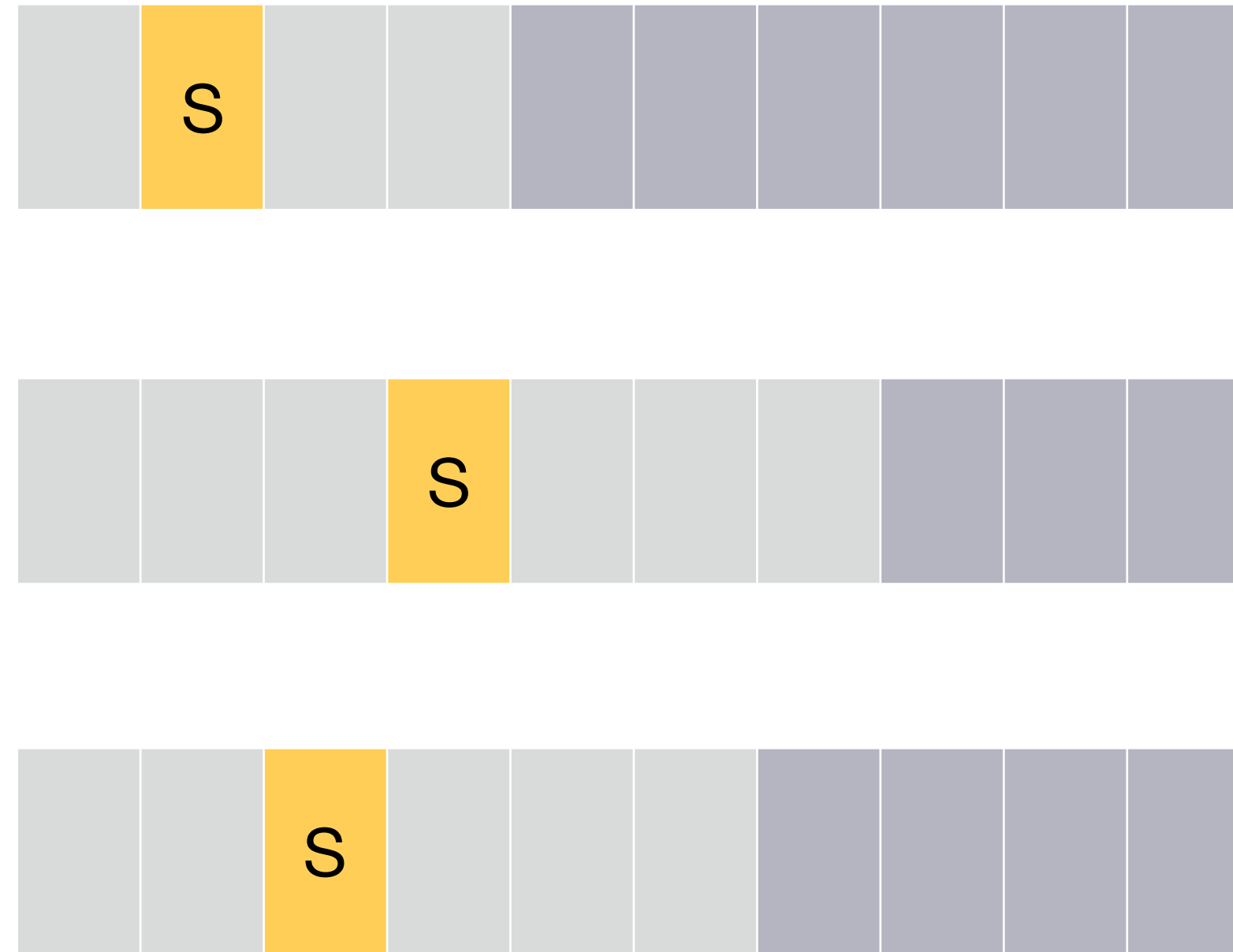


Processor

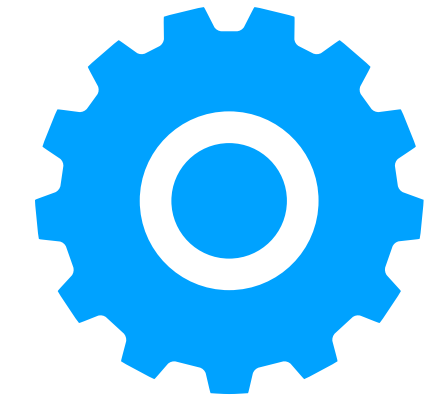




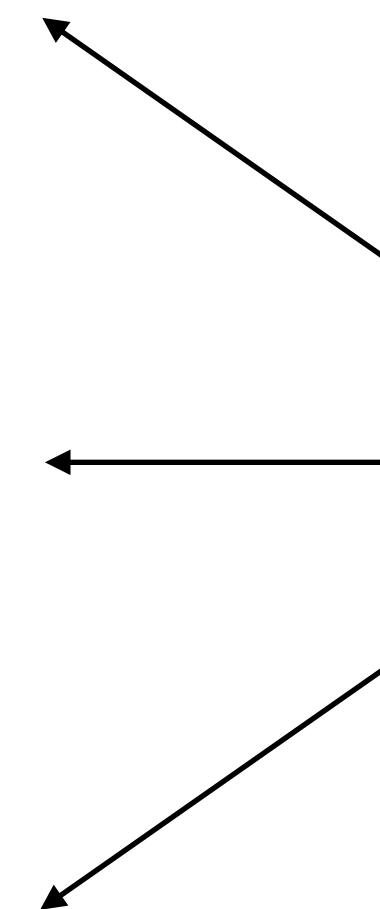
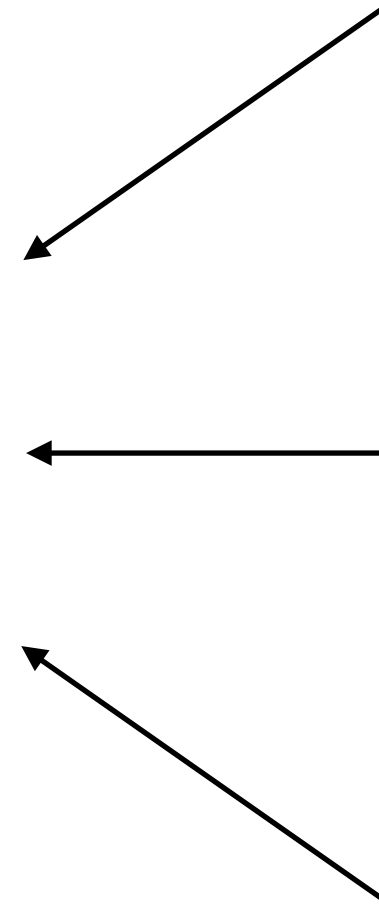
Observer



Output topics

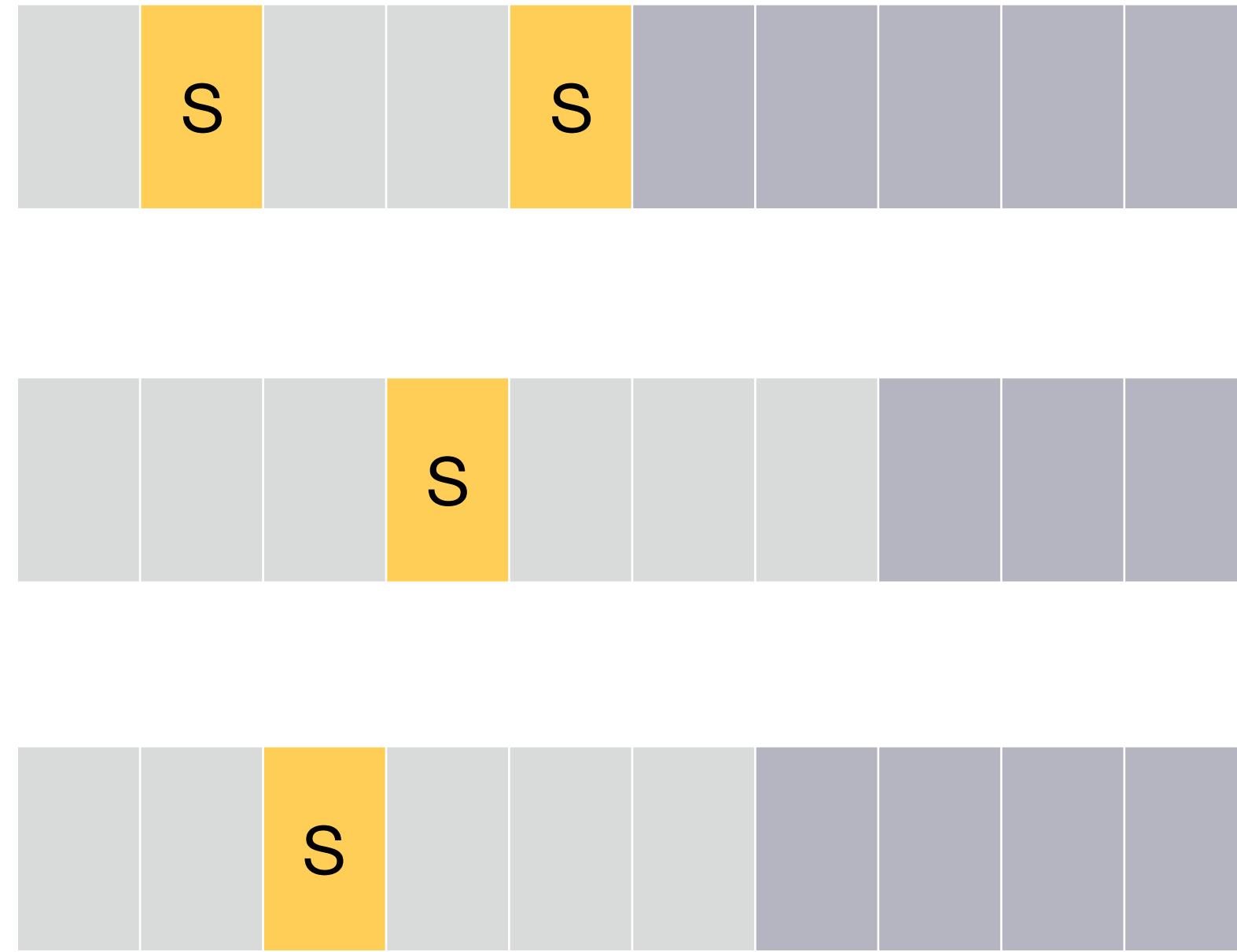


Processor

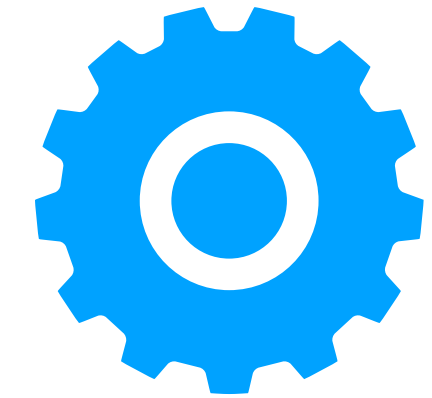




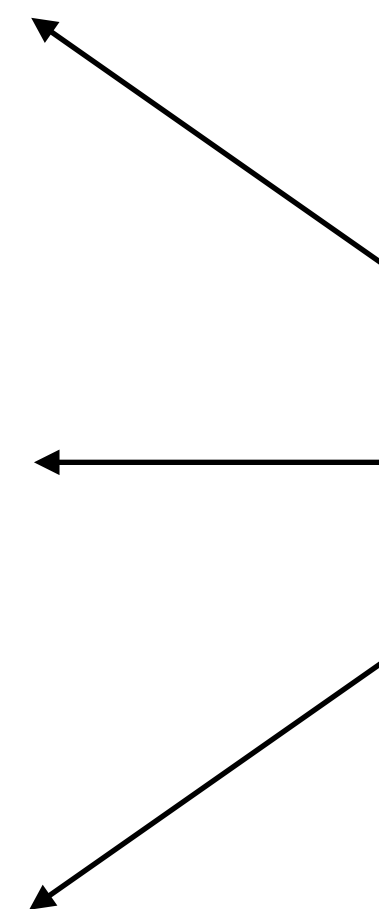
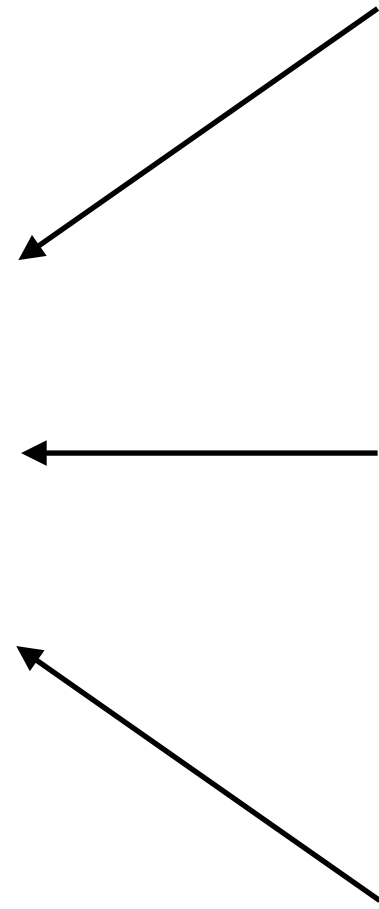
Observer



Output topics

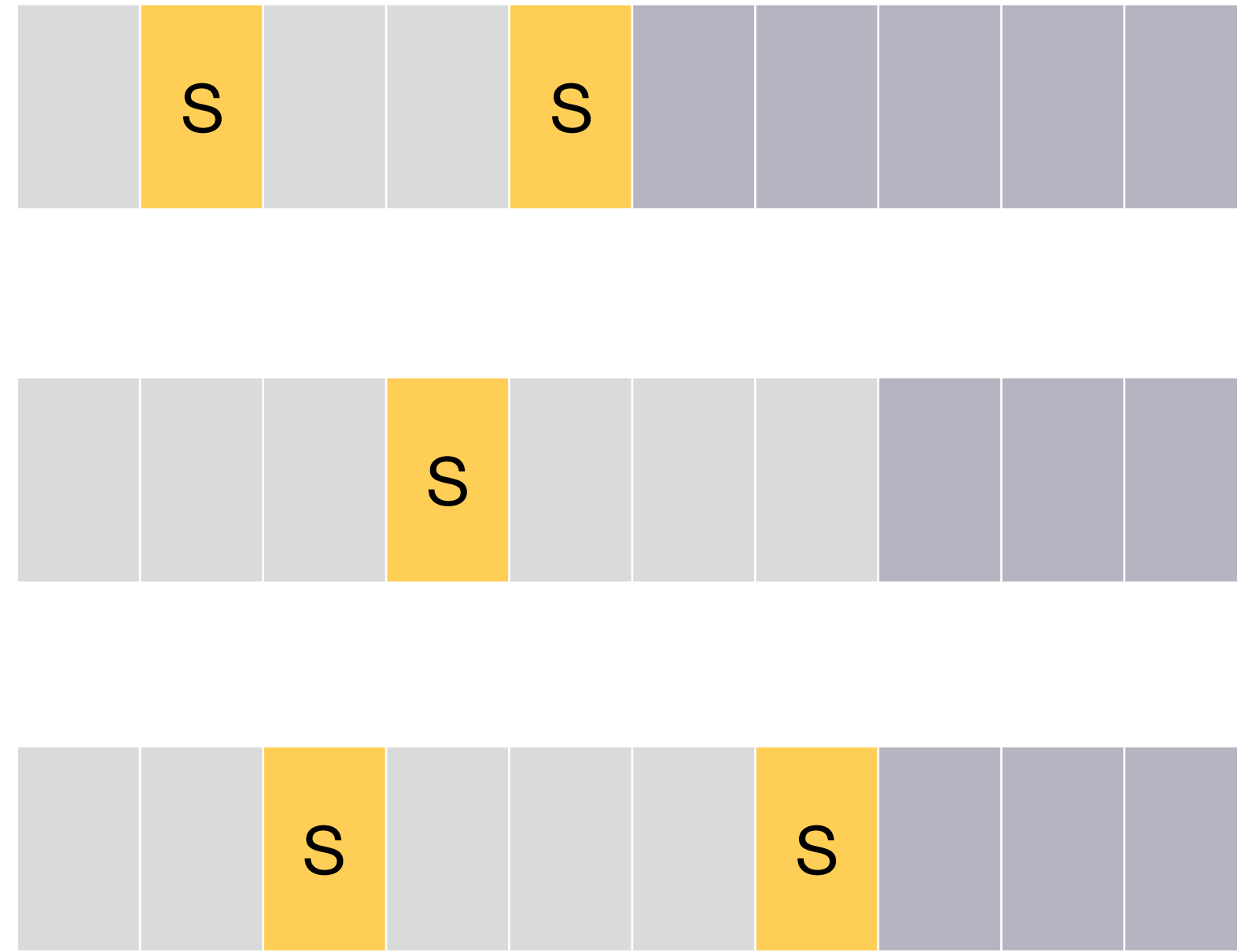


Processor

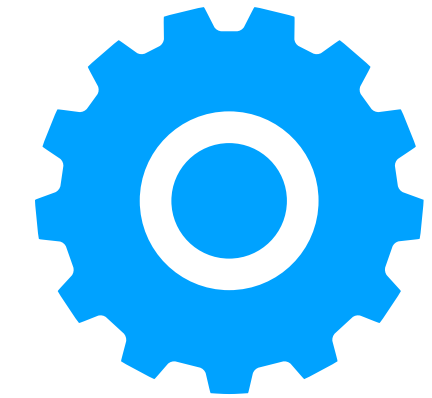




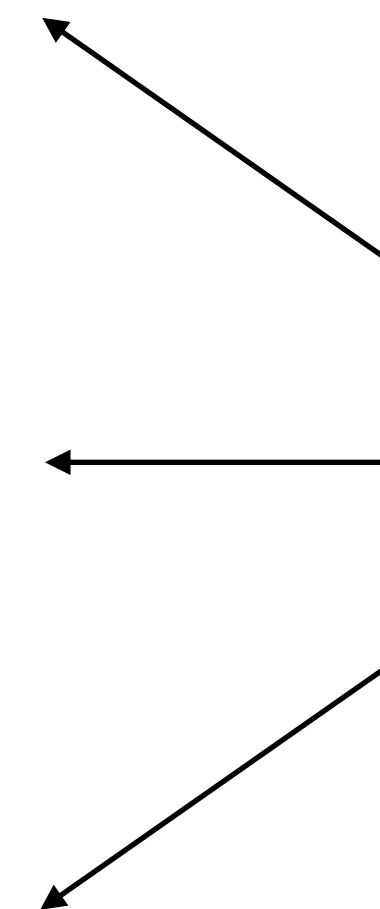
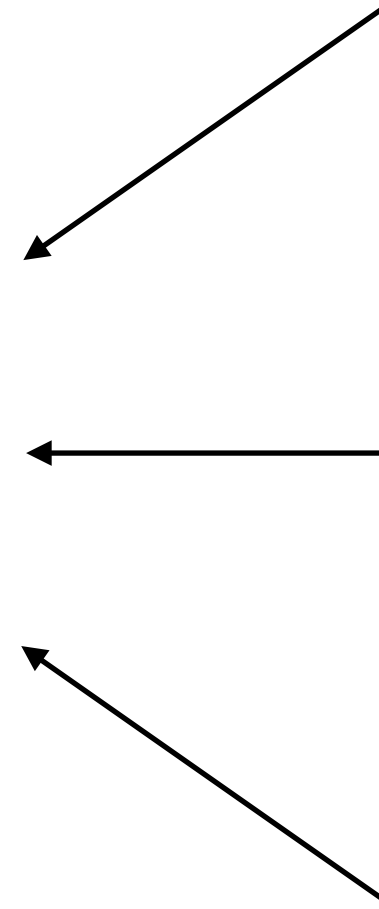
Observer



Output topics

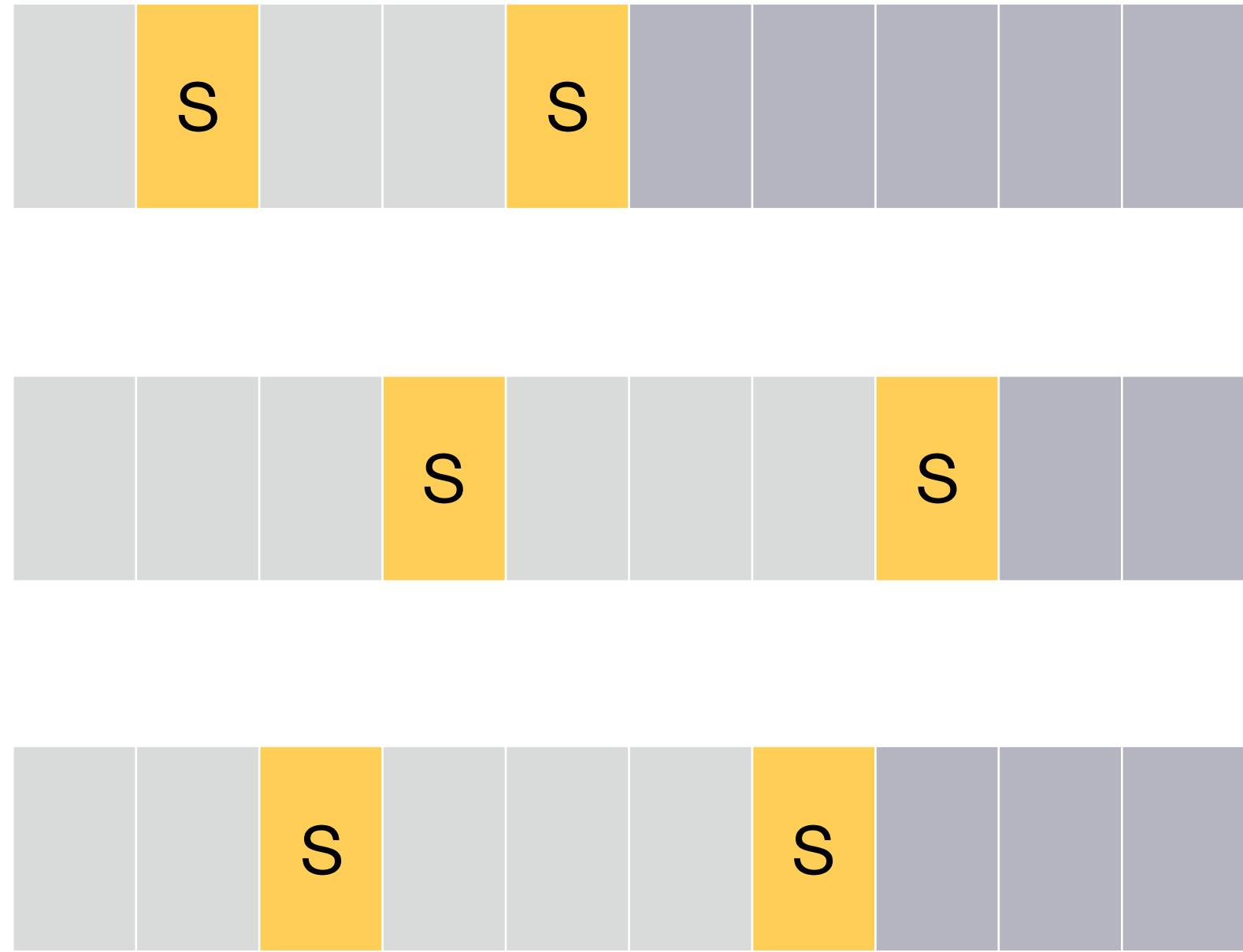
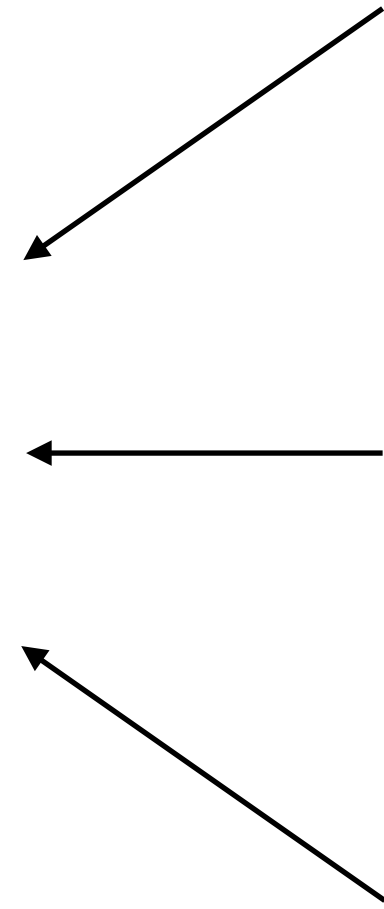


Processor

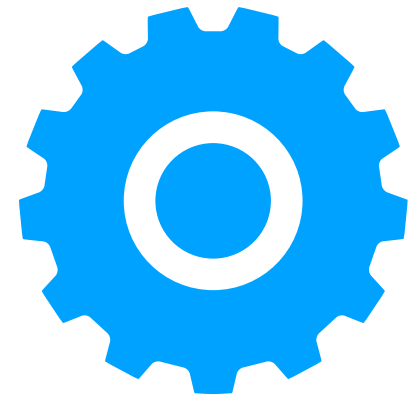
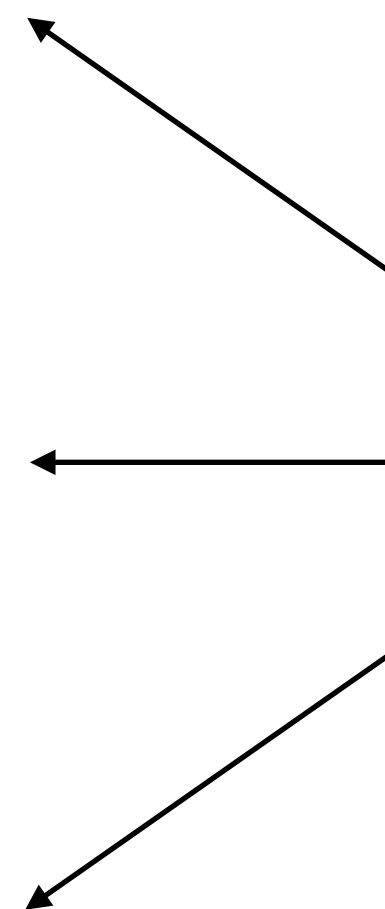




Observer



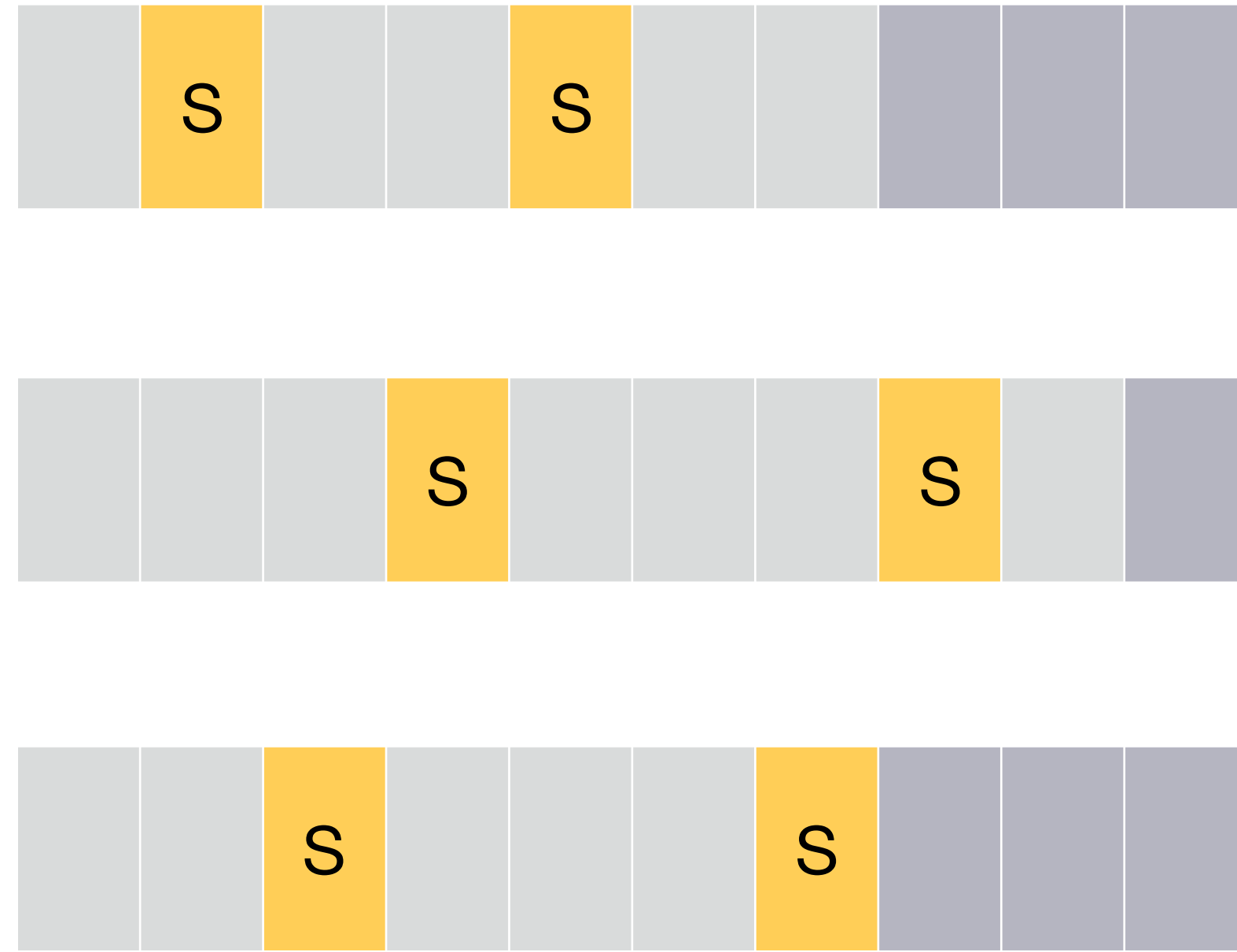
Output topics



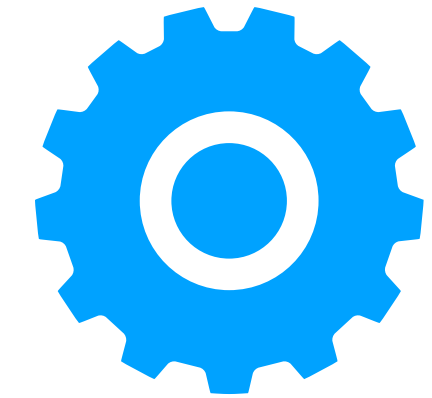
Processor



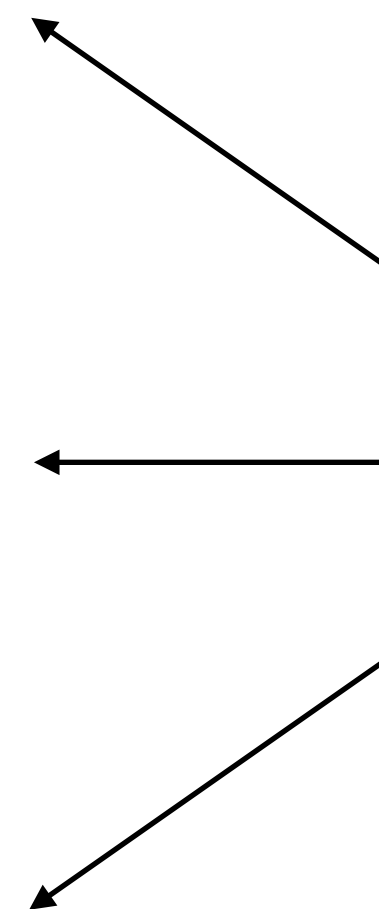
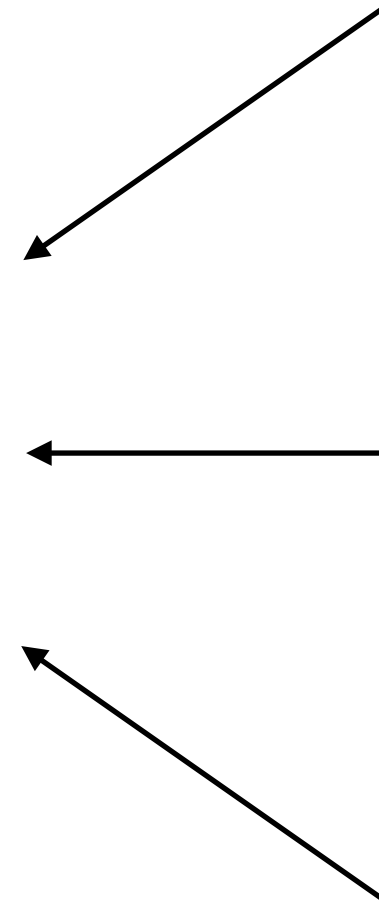
Observer



Output topics



Processor

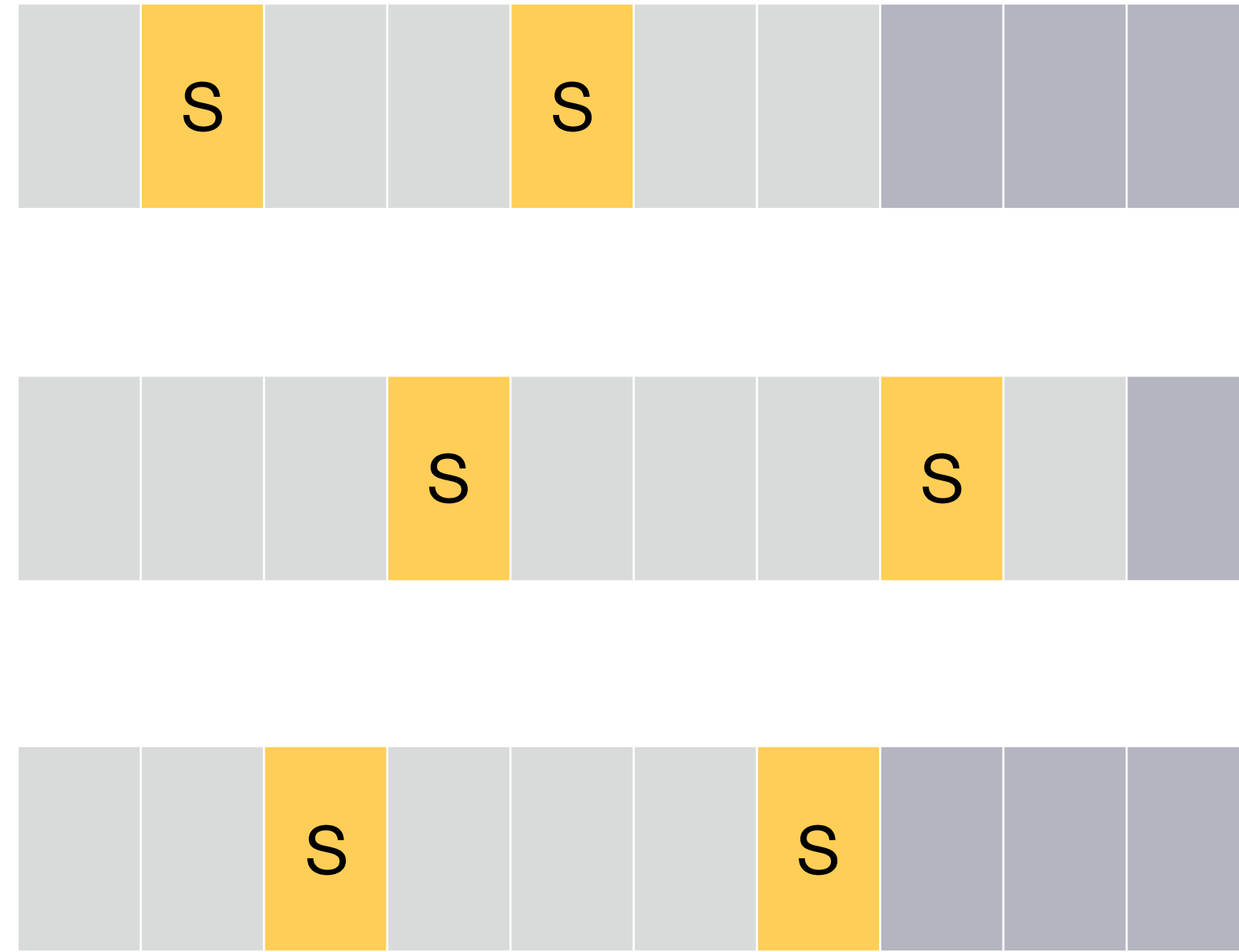
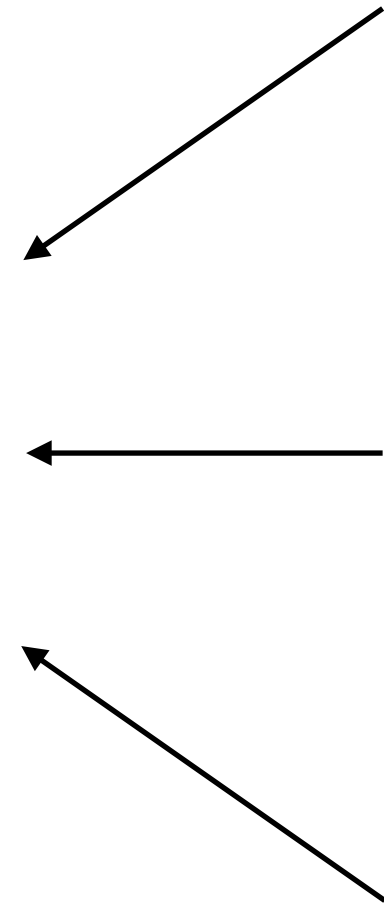


Снэпшоты Чэнду - Лэмпорт

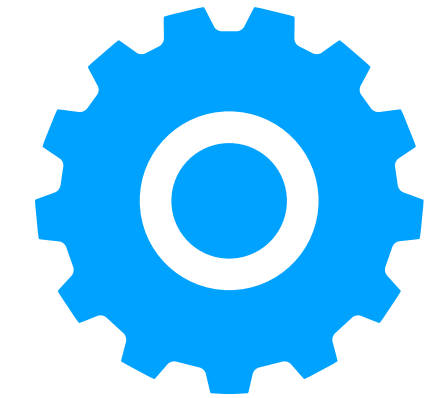
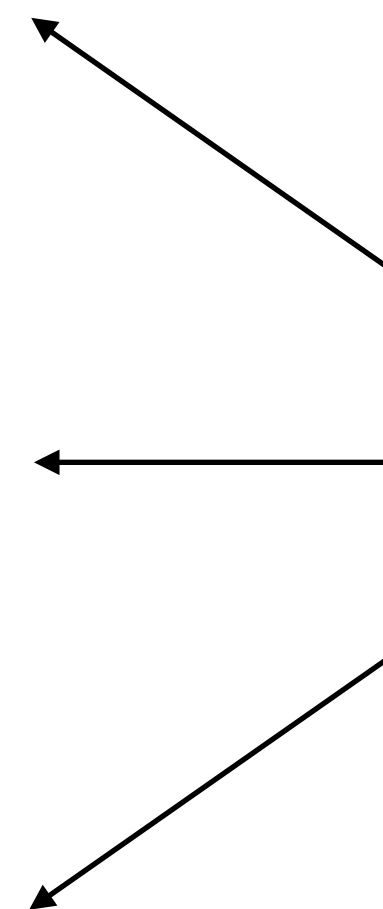
- А как же быть с со сбоями?
 - Soft: сбой во время обработки
 - Hard: сбой до или во время снятия снэпшота



Observer



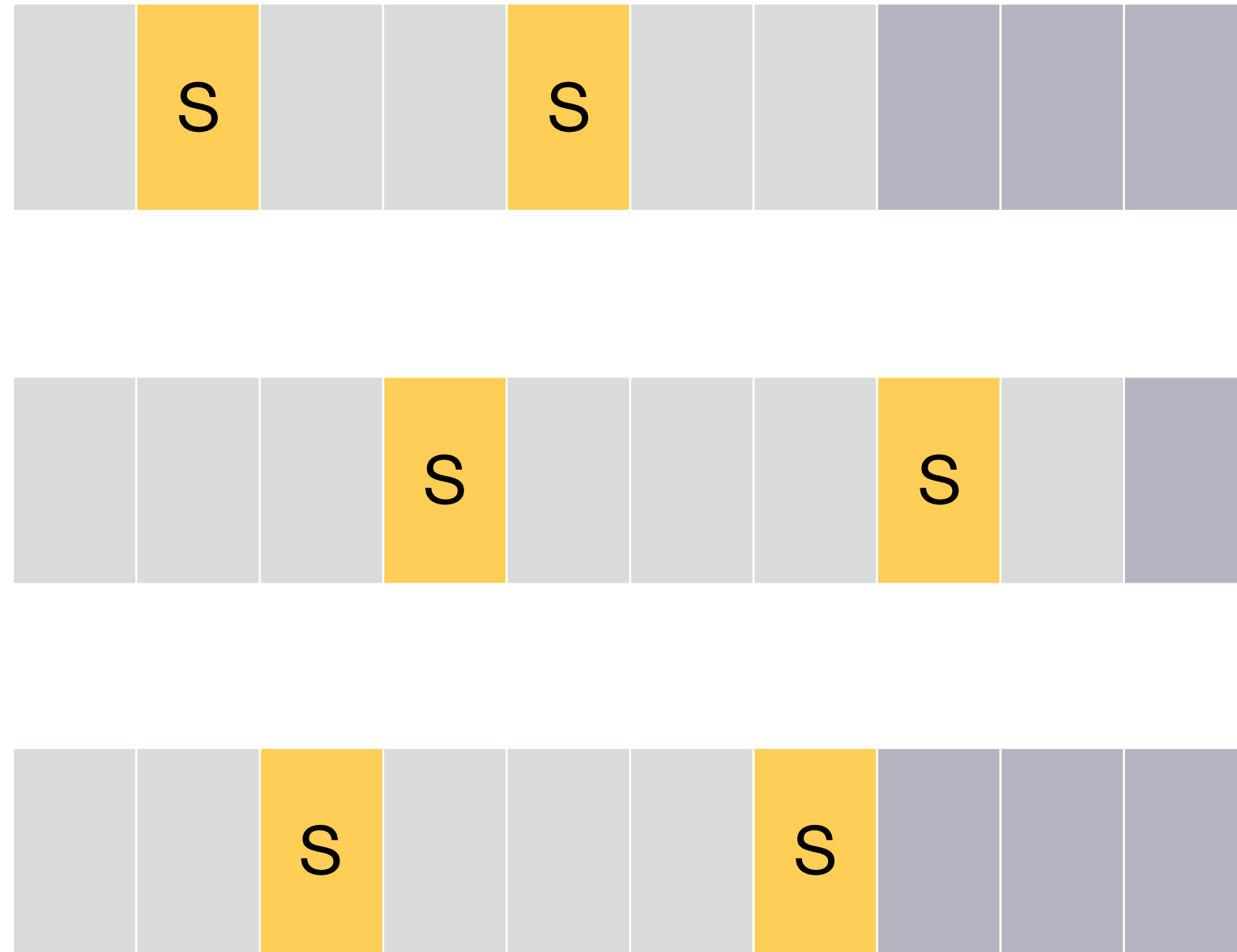
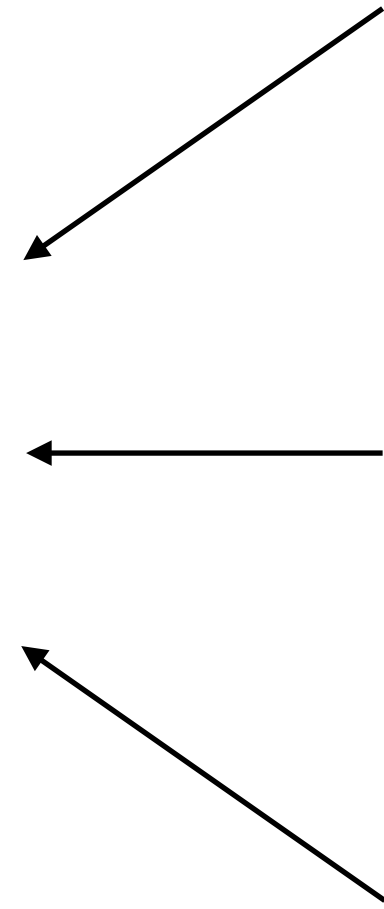
Output topics



Processor

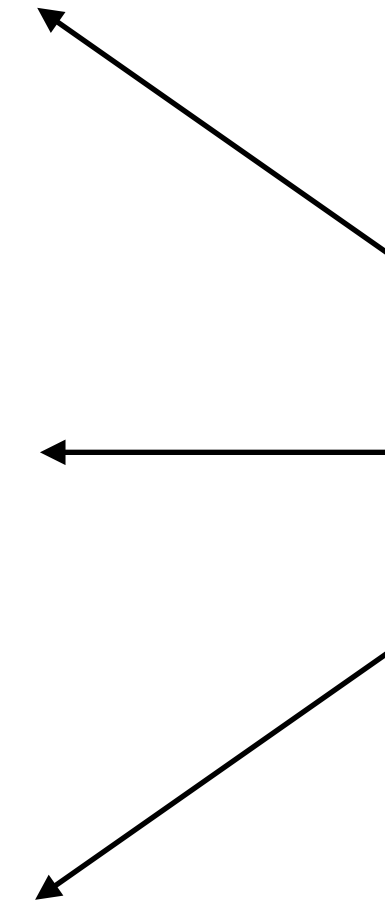


Observer



Output topics

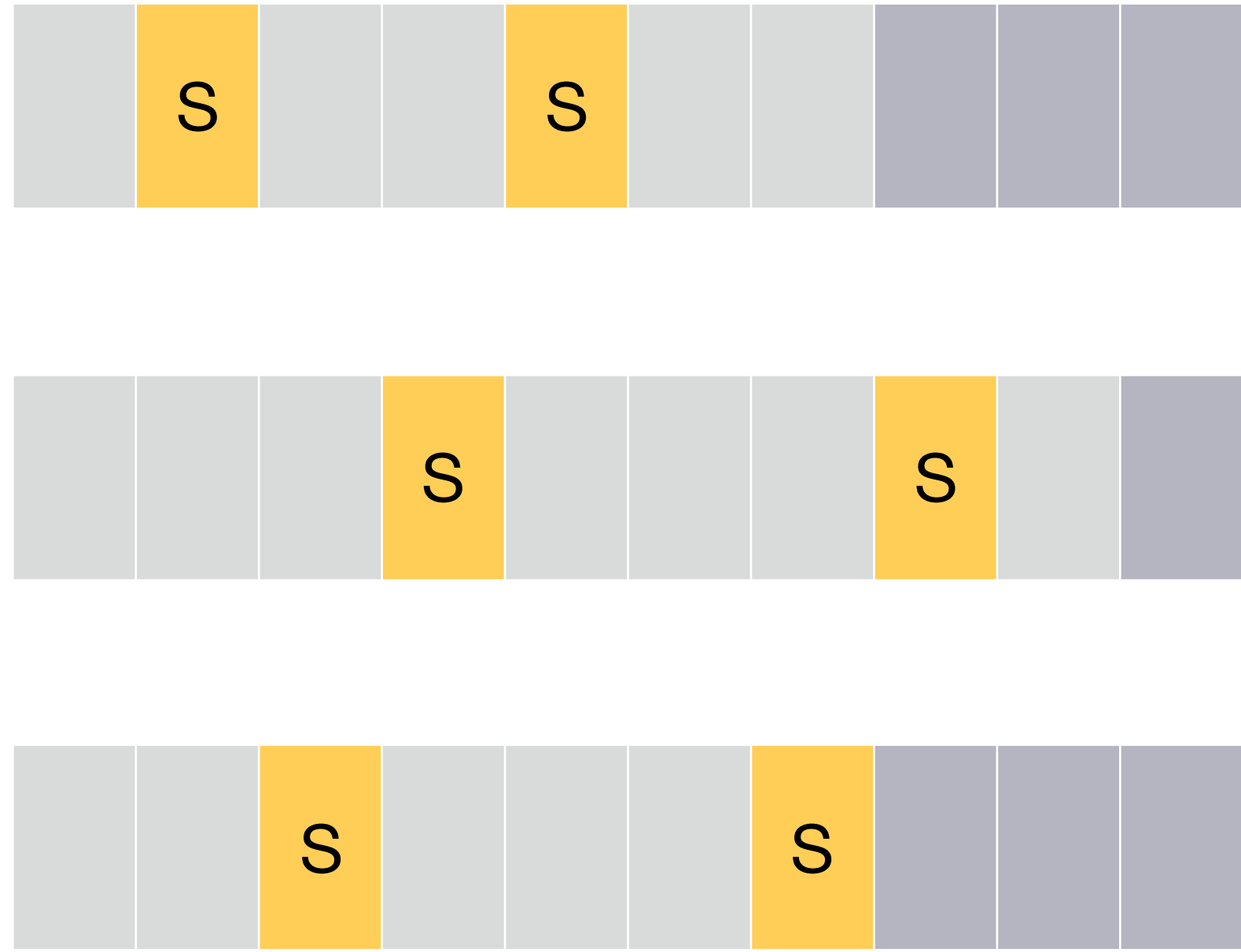
Processor crashed!



Processor

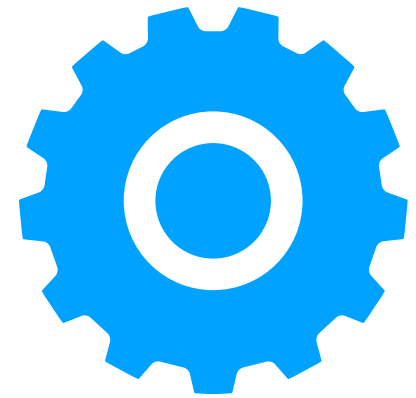


Observer

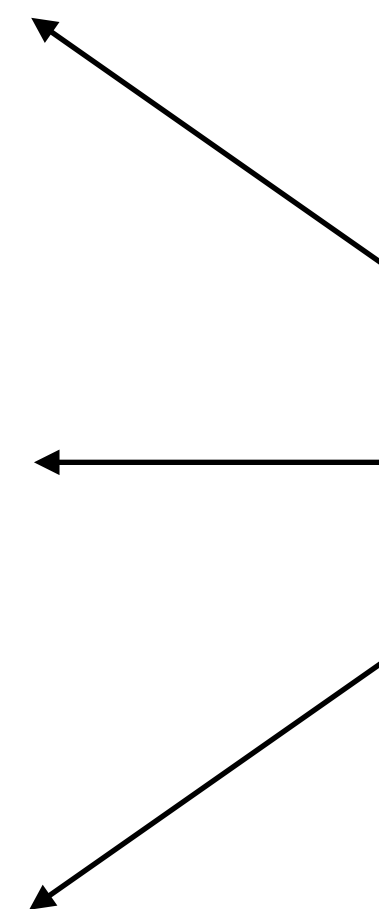
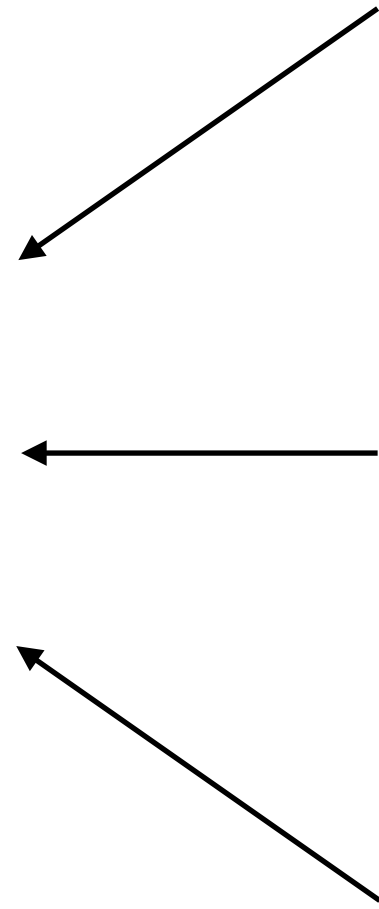


Output topics

Restart!

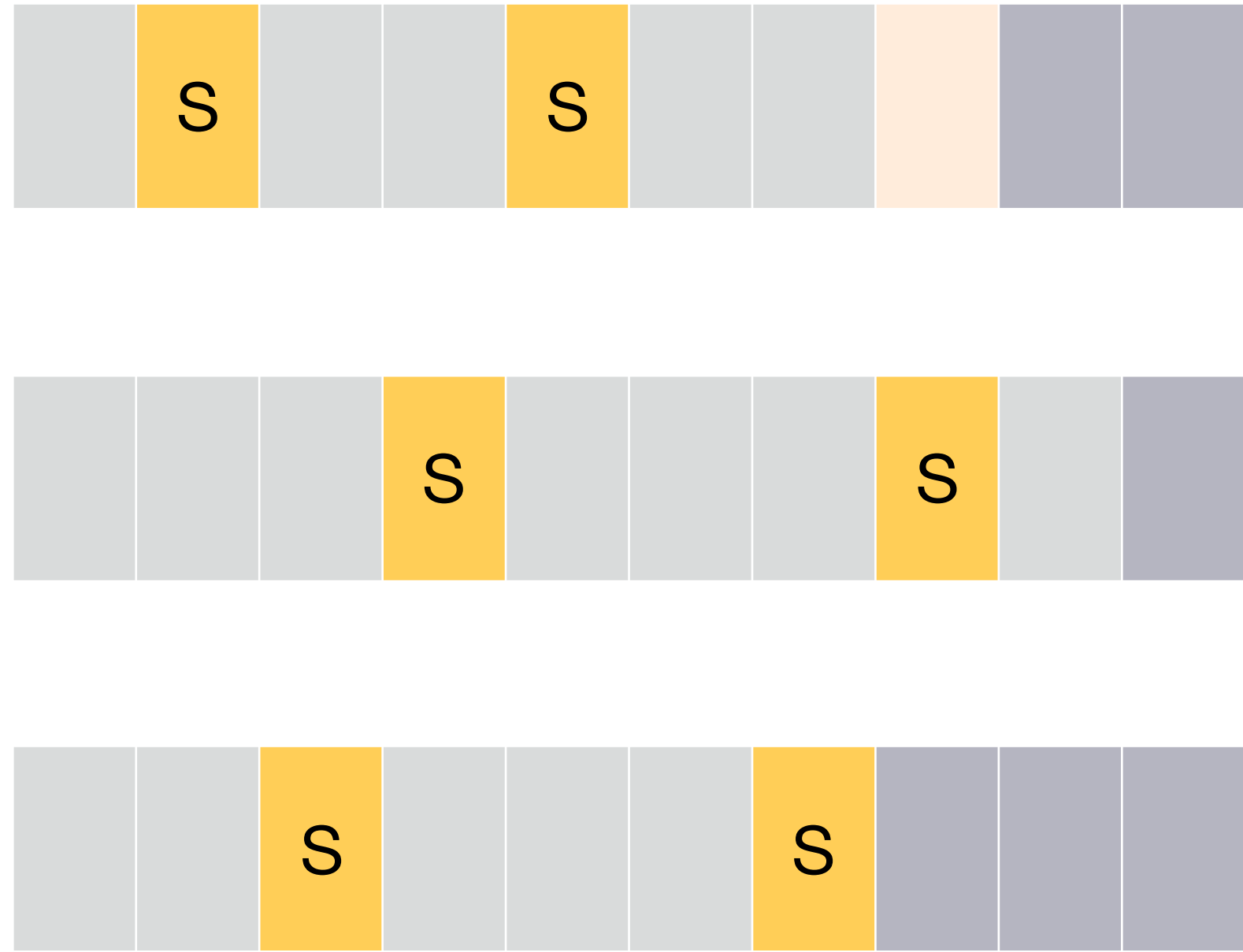
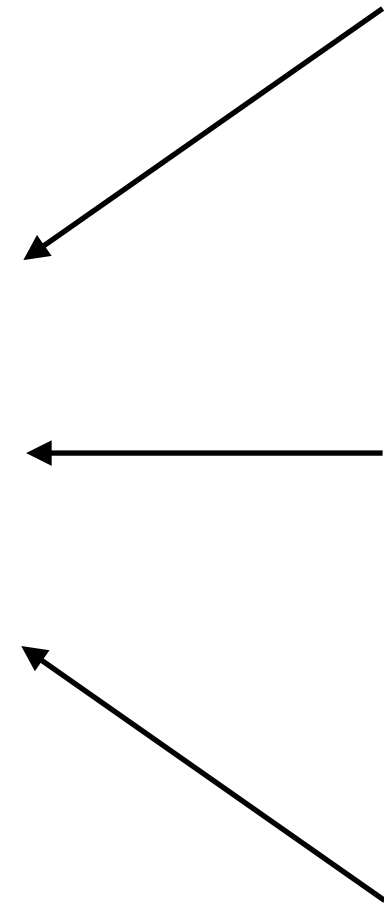


Processor

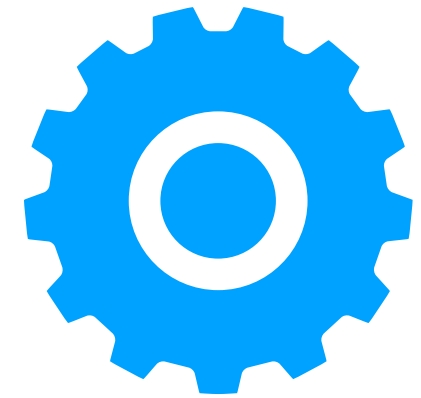
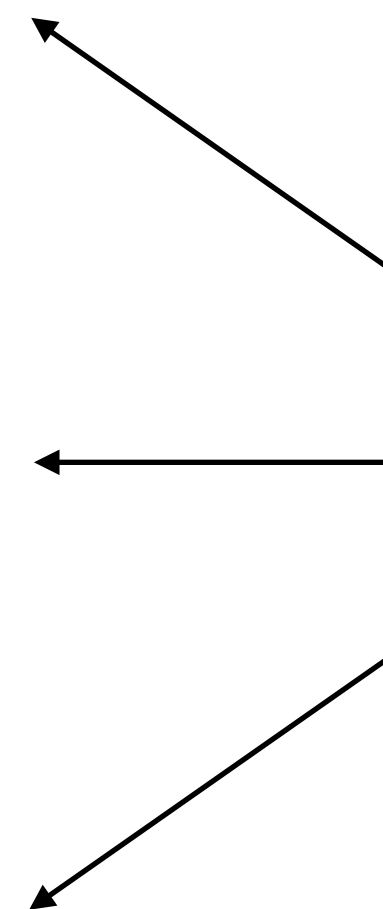




Observer



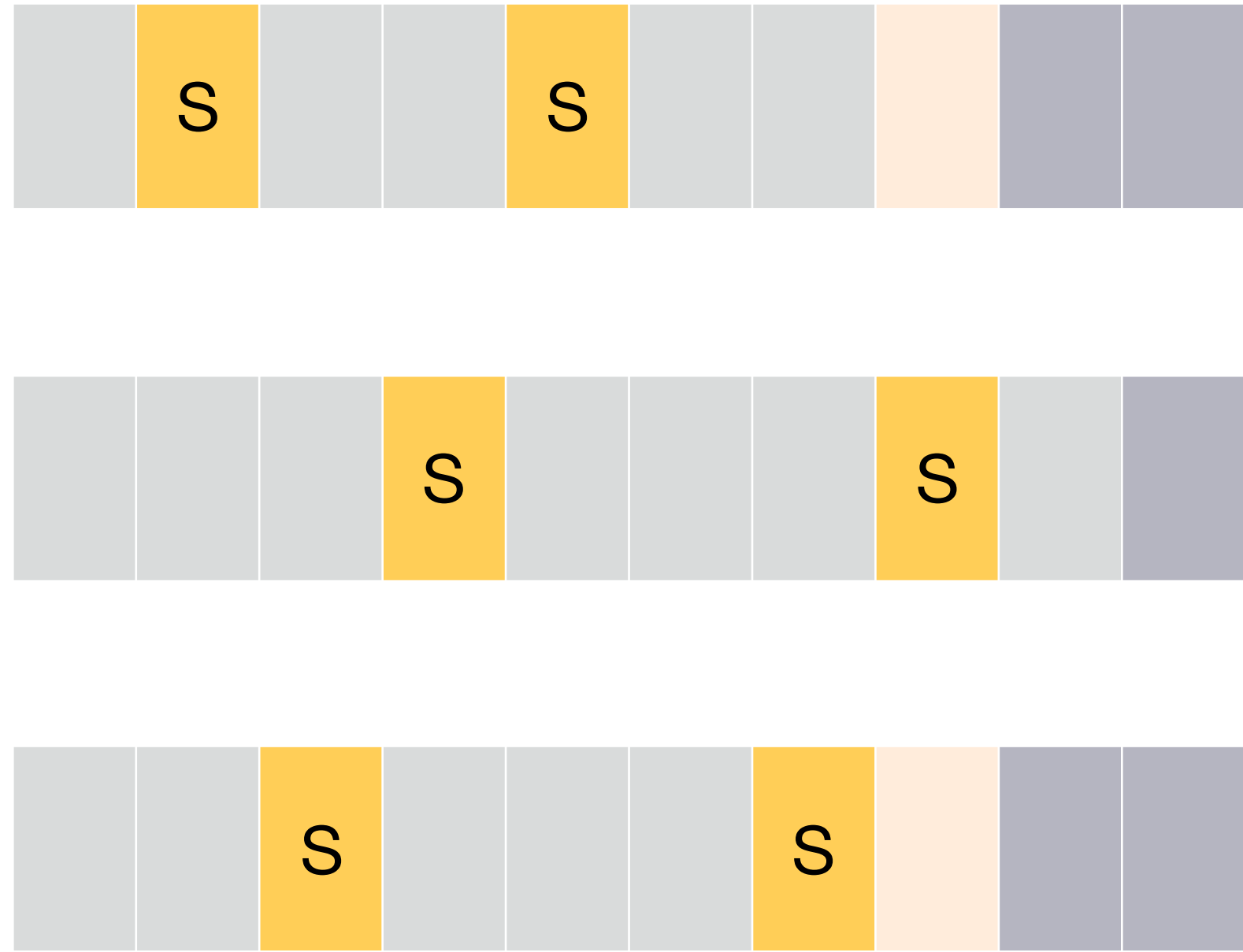
Output topics



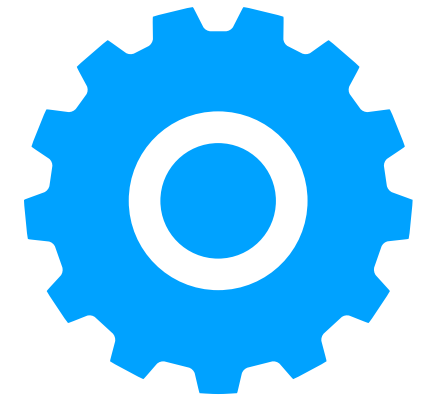
Processor



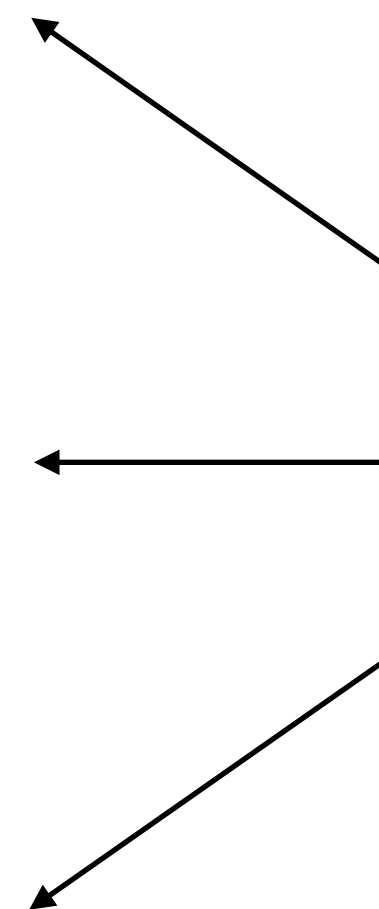
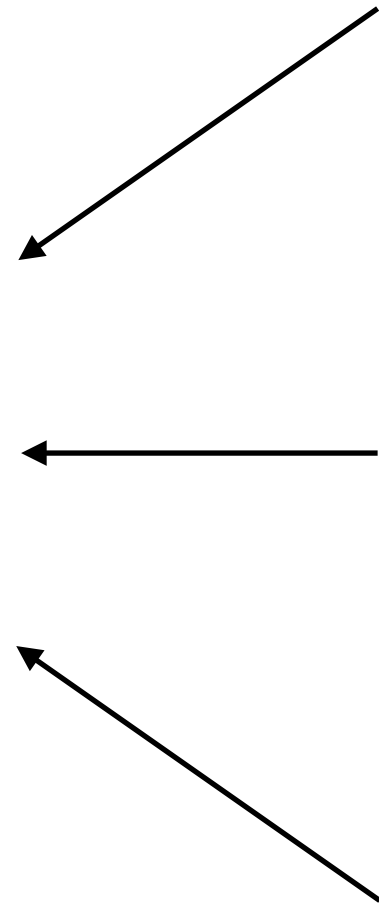
Observer



Output topics



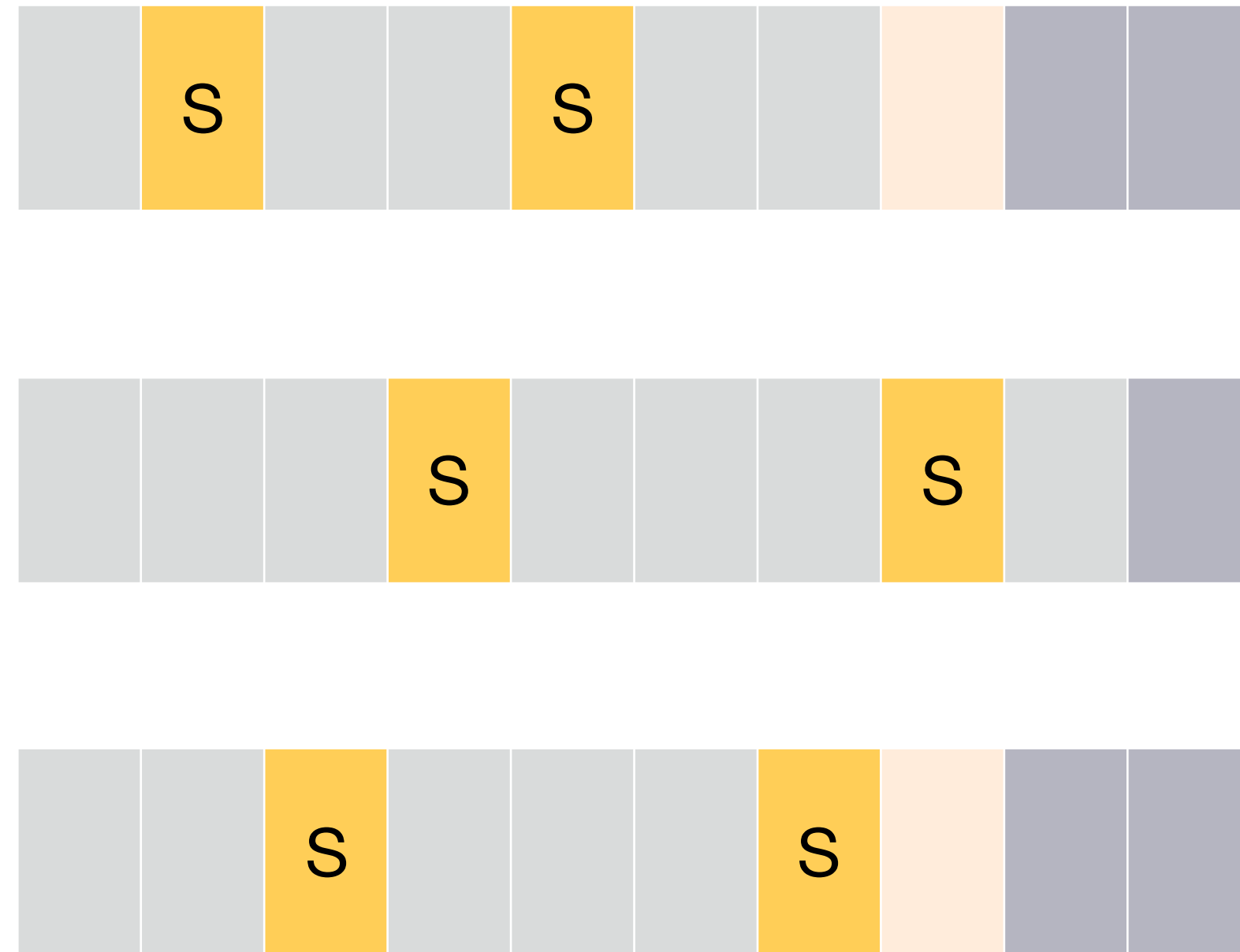
Processor



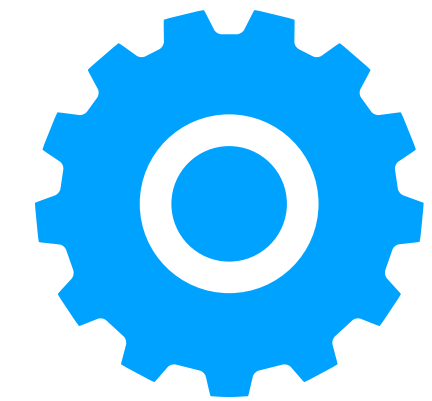
Hmm...



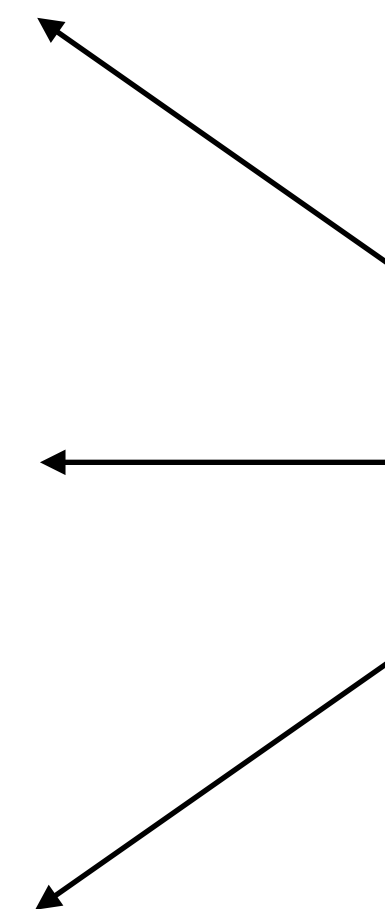
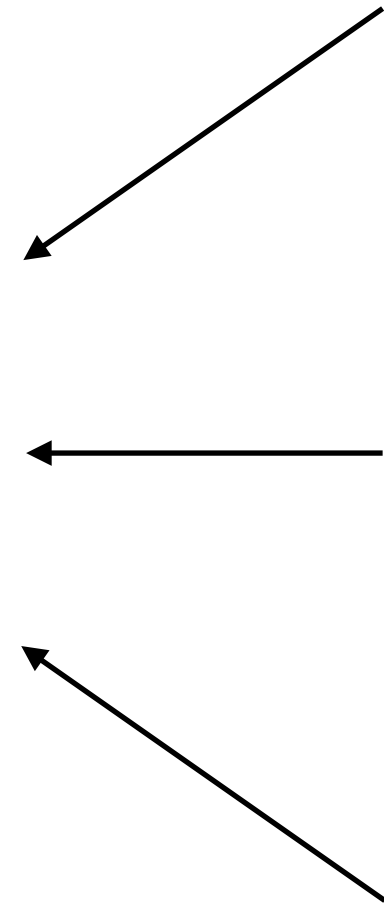
Observer



Output topics

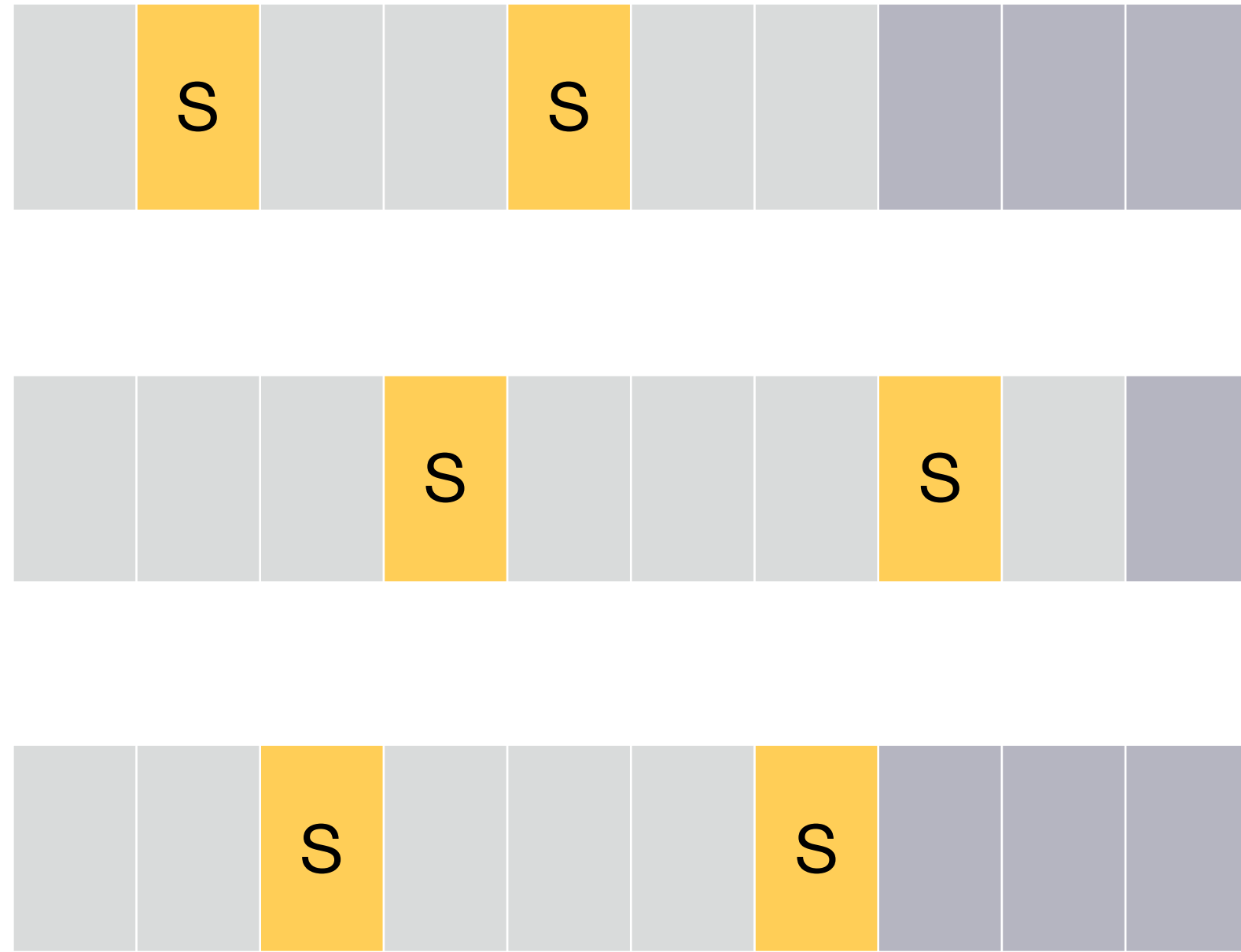
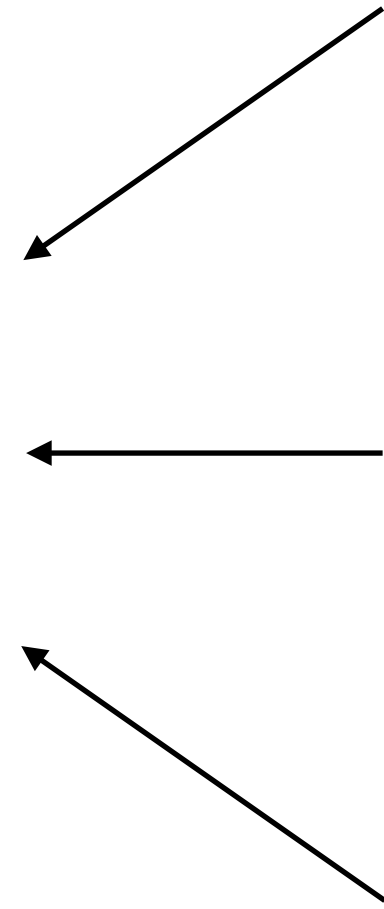


Processor

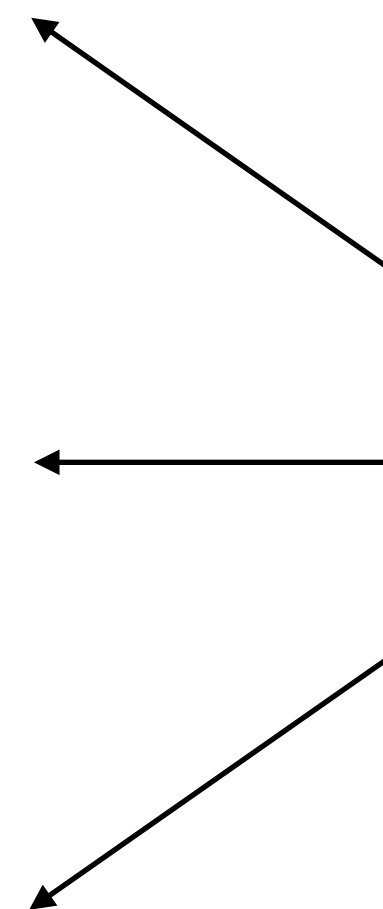




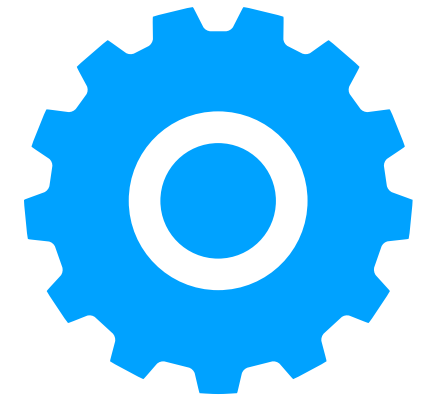
Observer



Output topics



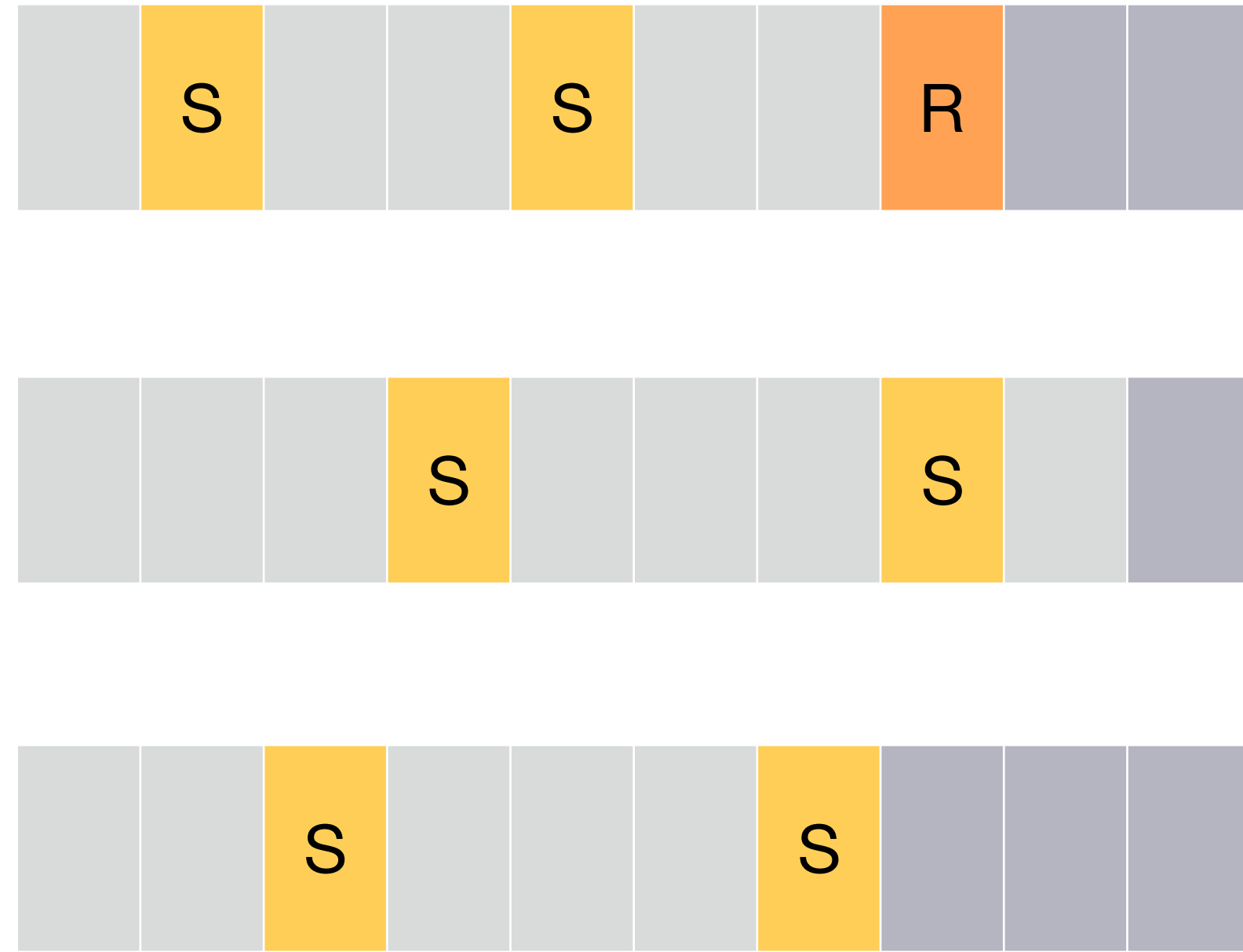
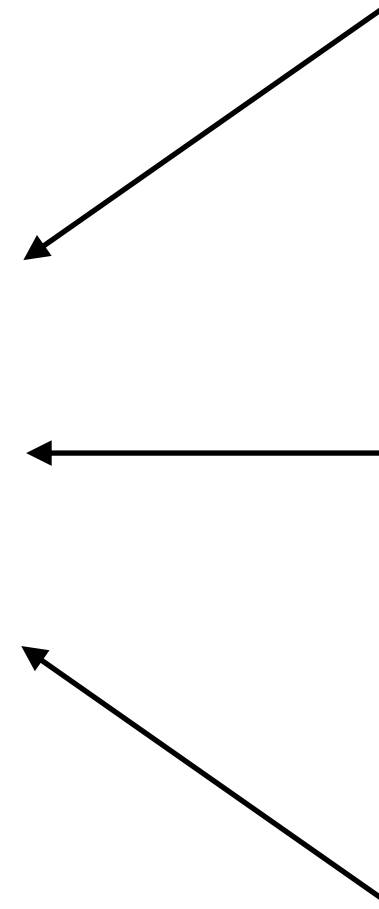
Restart!



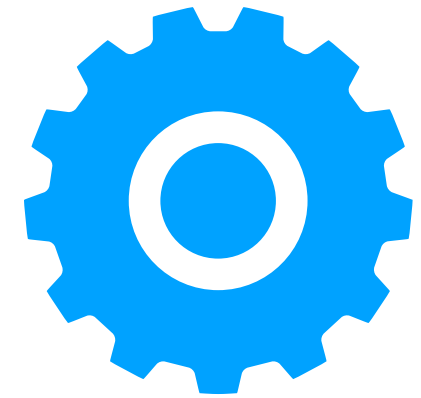
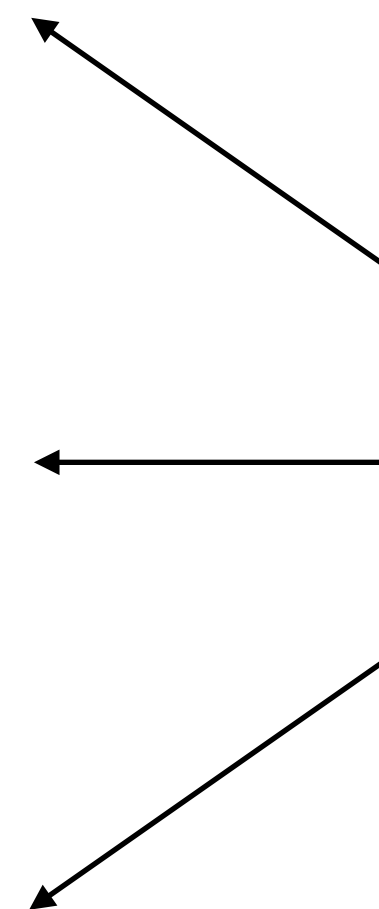
Processor



Observer



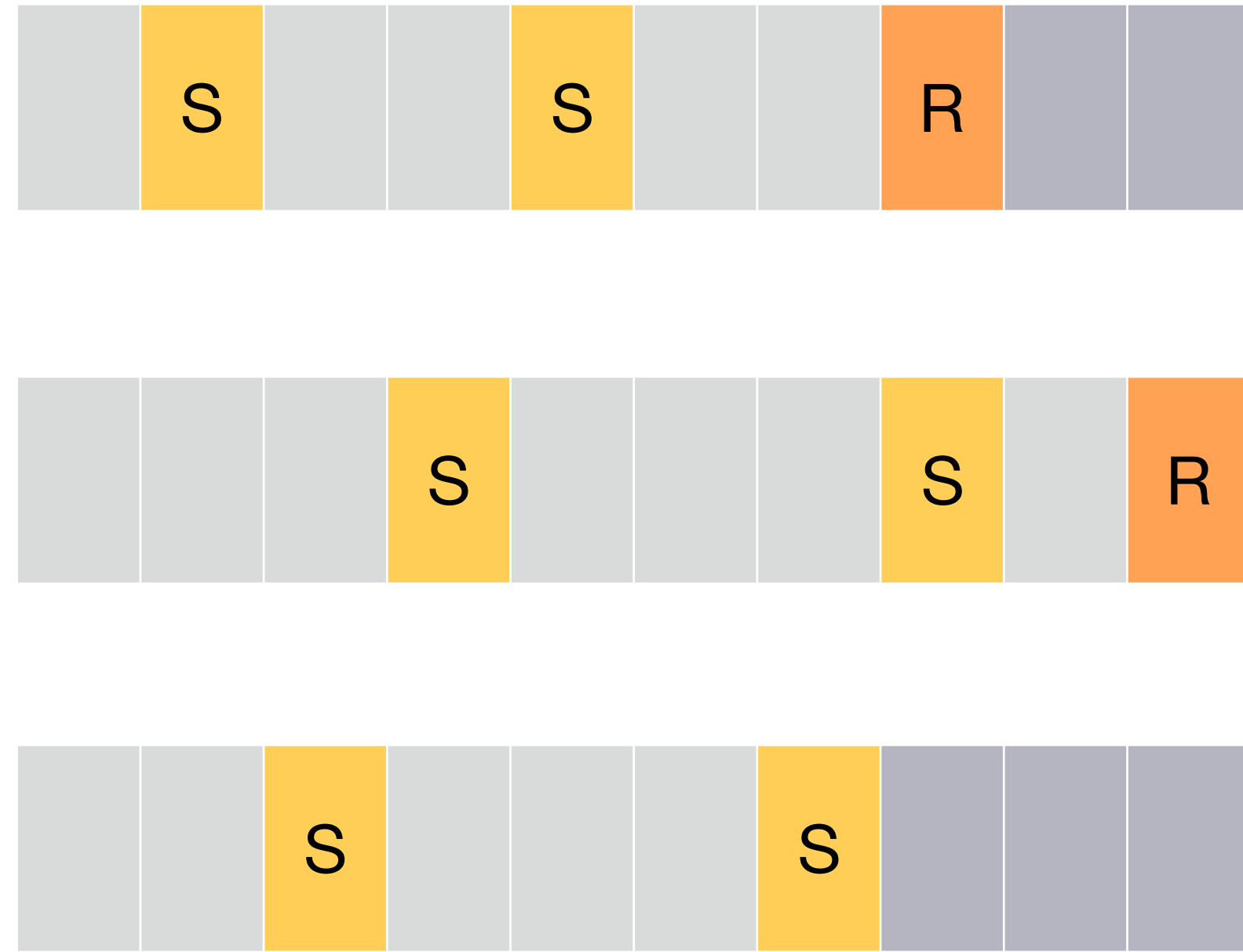
Output topics



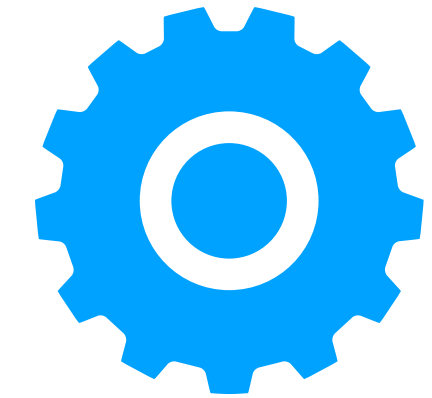
Processor



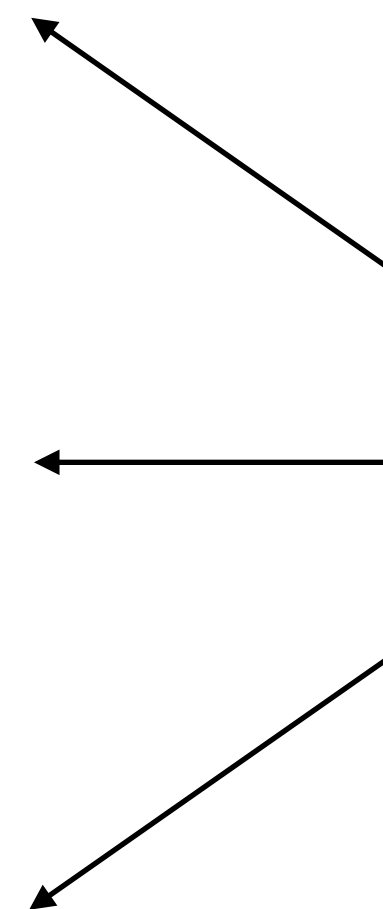
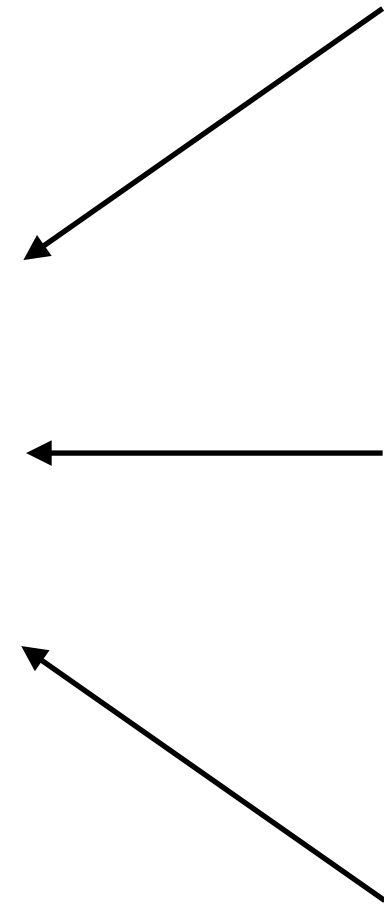
Observer



Output topics

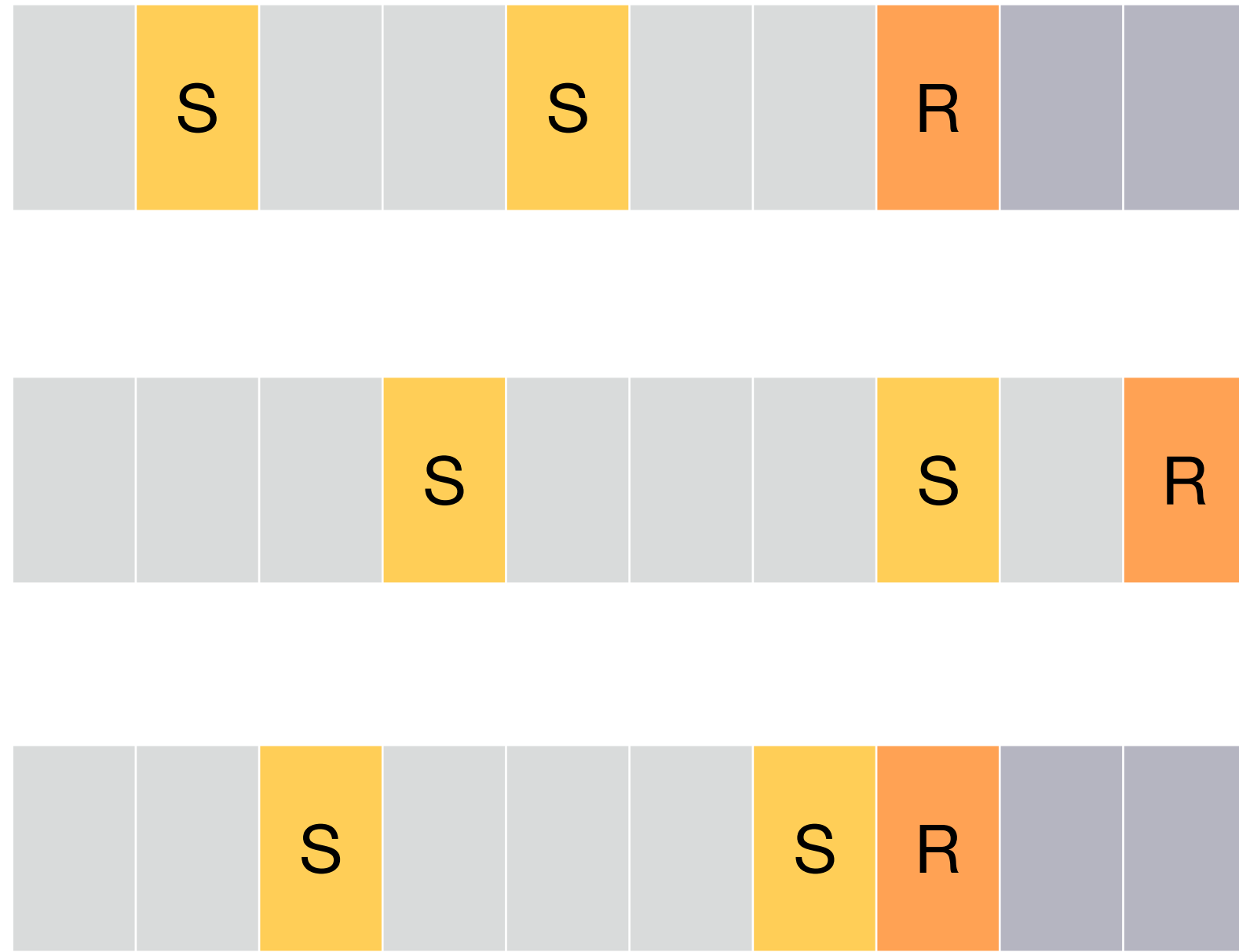


Processor

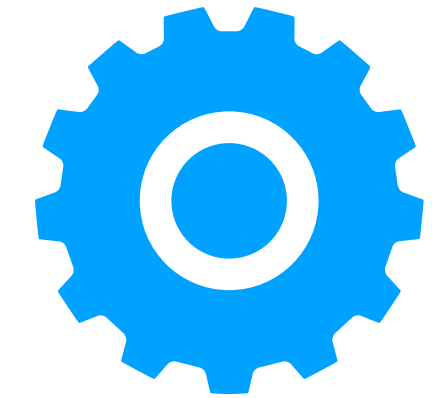




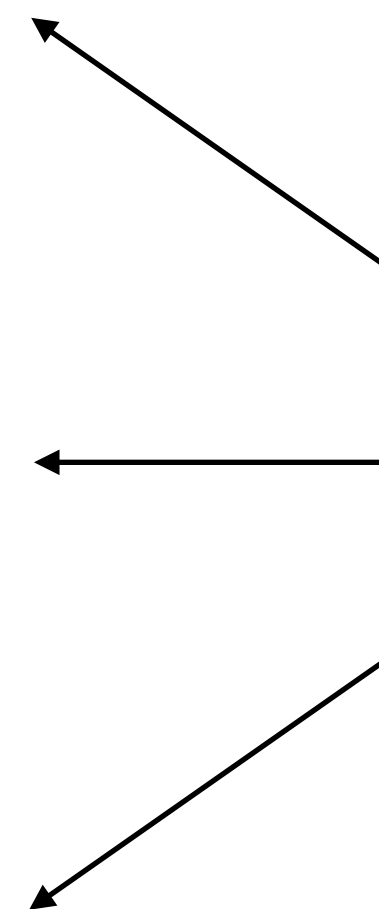
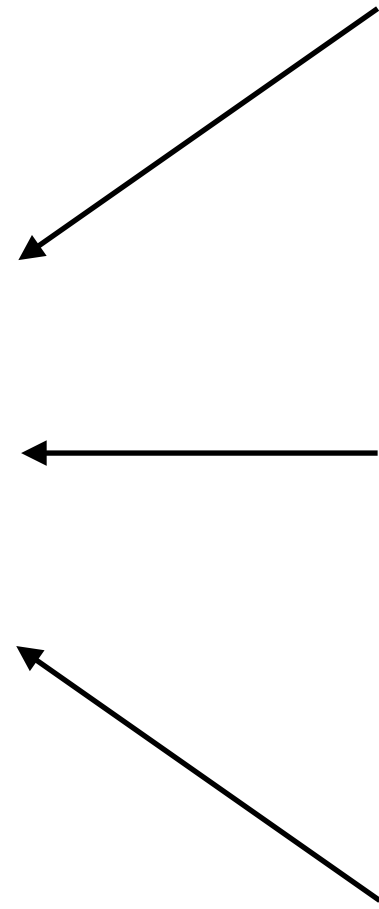
Observer



Output topics



Processor

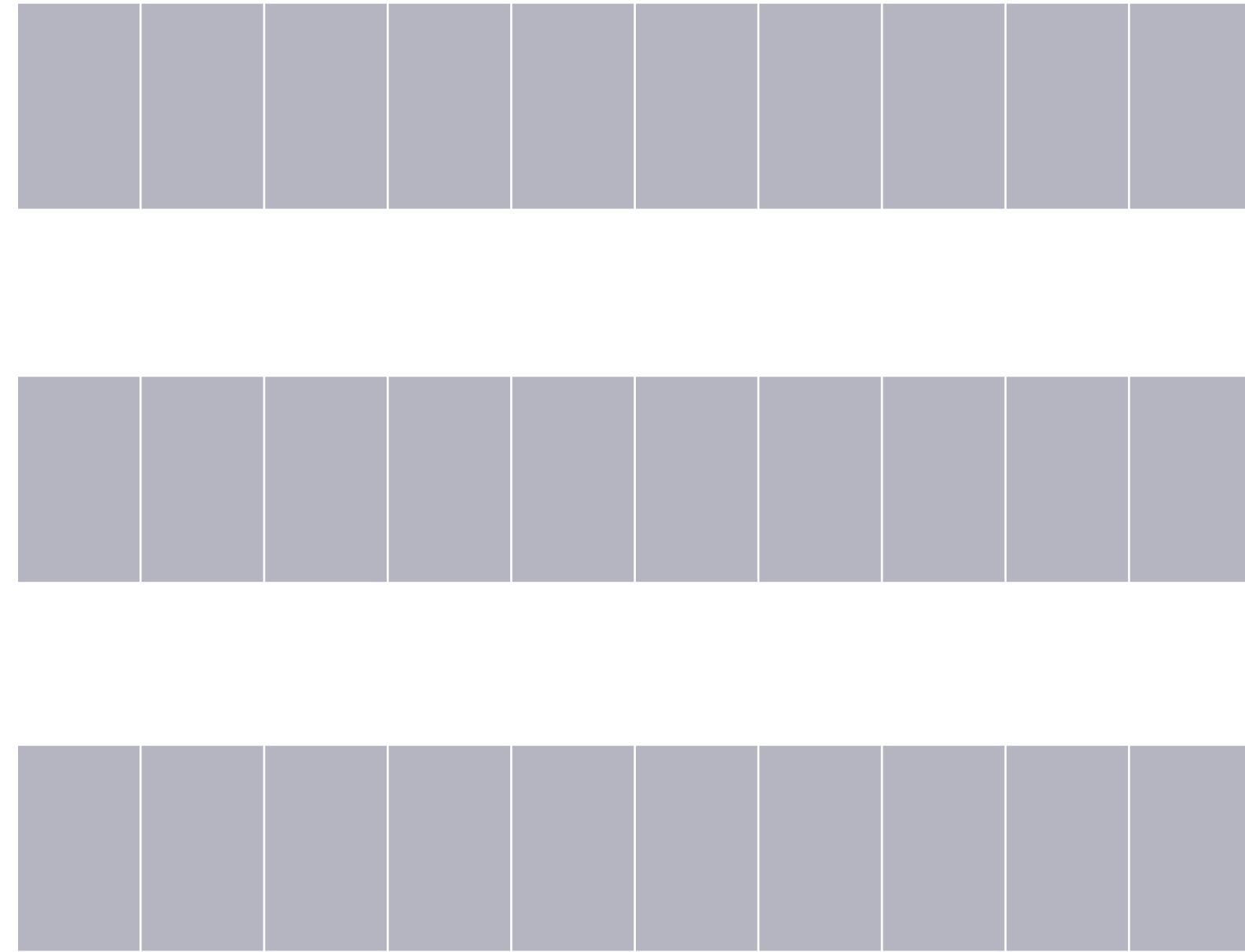
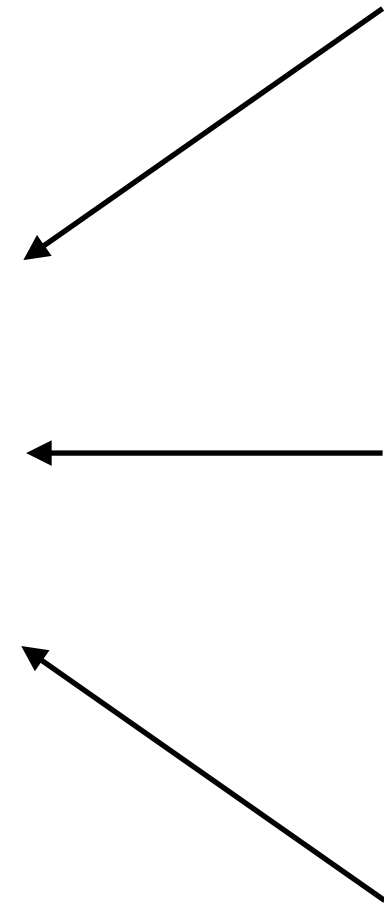


Транзакции в Кафке

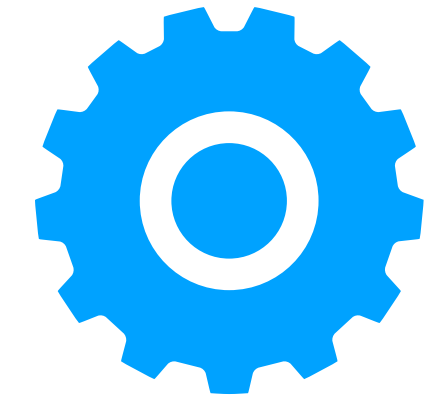
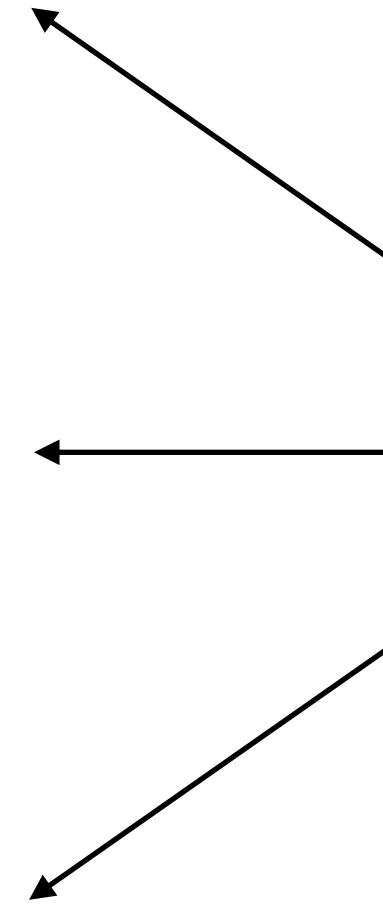
- Два типа маркеров COMMIT and ABORT
- Мягкие (Soft) типы сбоя обрабатываются записью ABORT маркера во все партиции



Observer



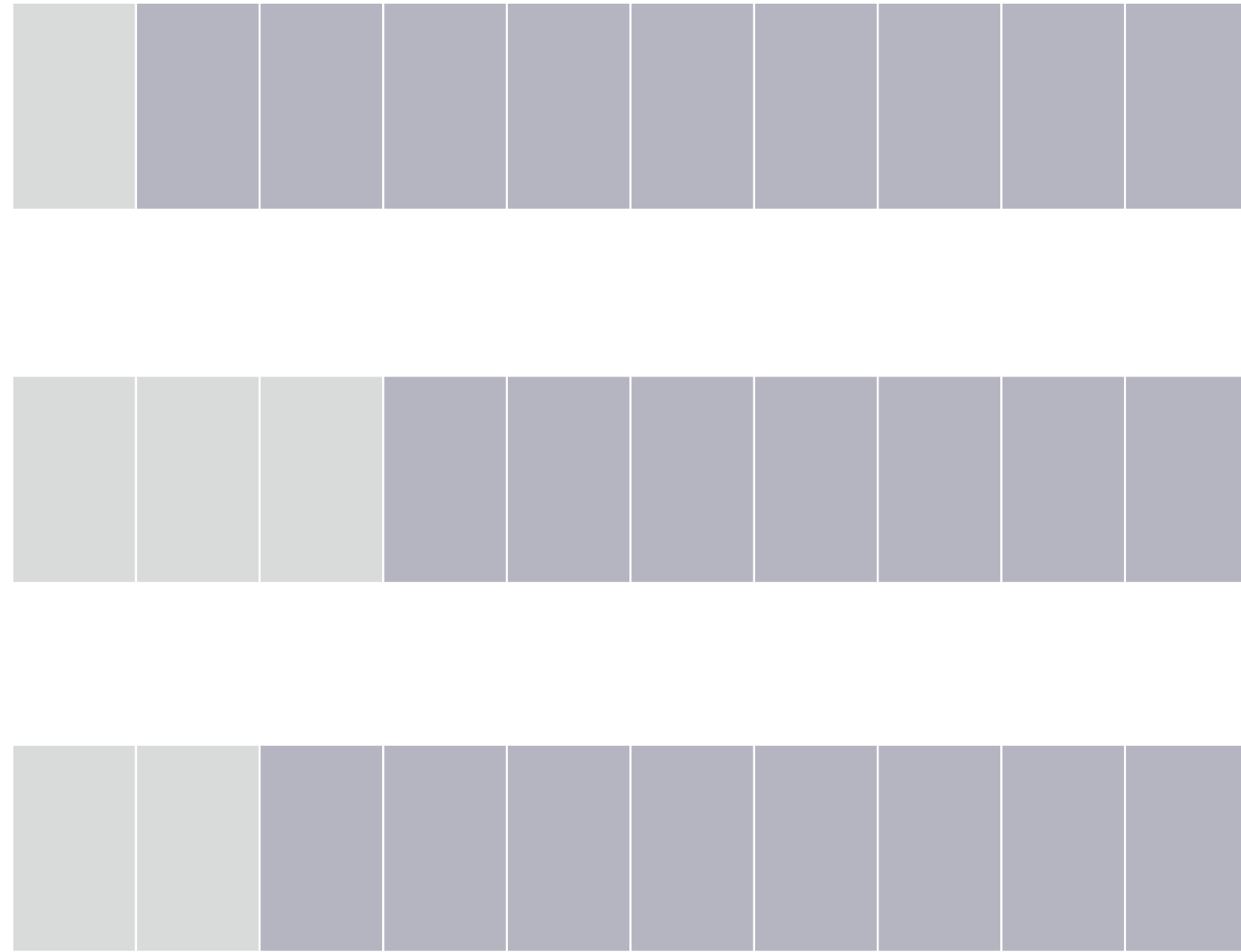
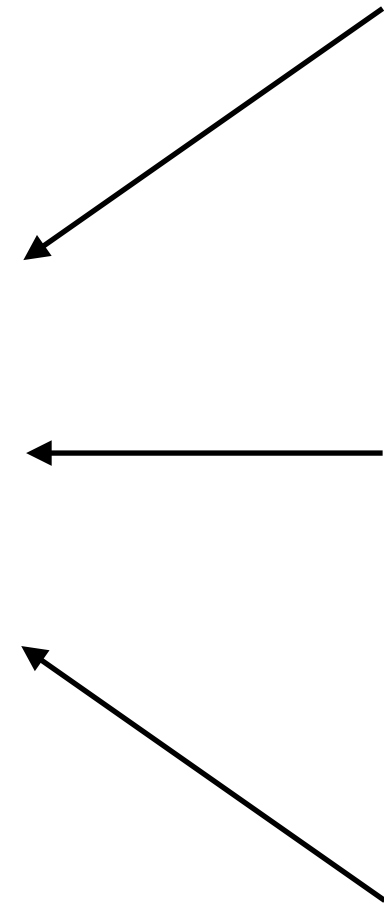
Output



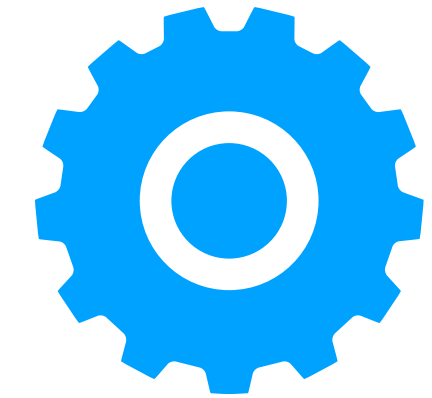
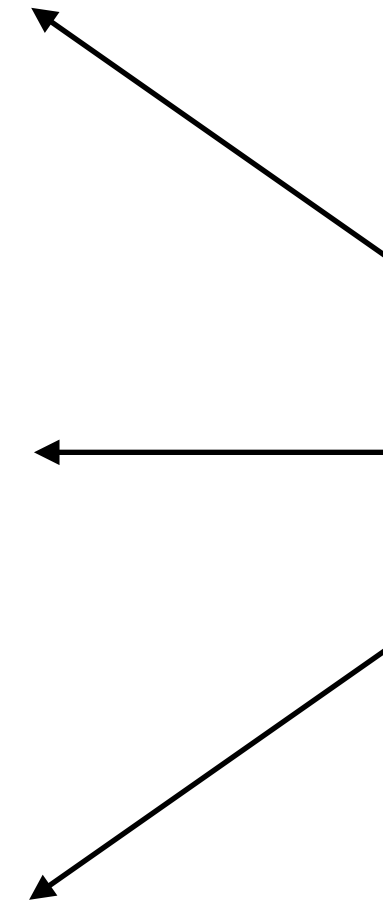
Processor



Observer



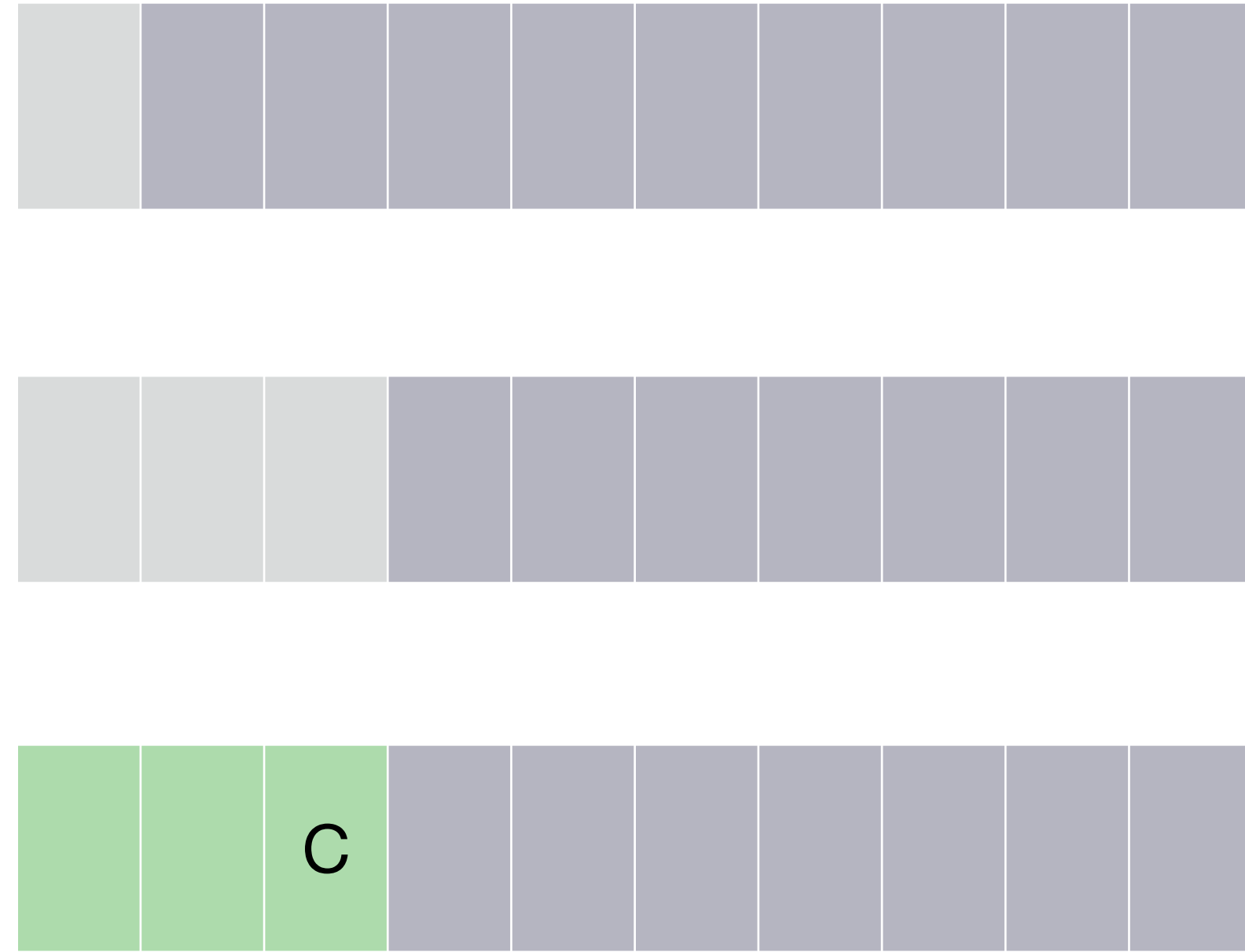
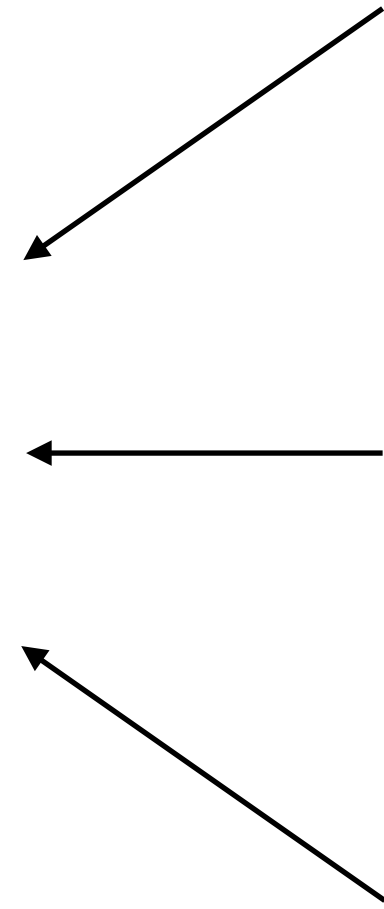
Output



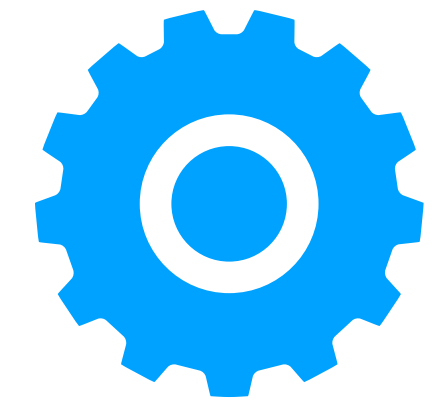
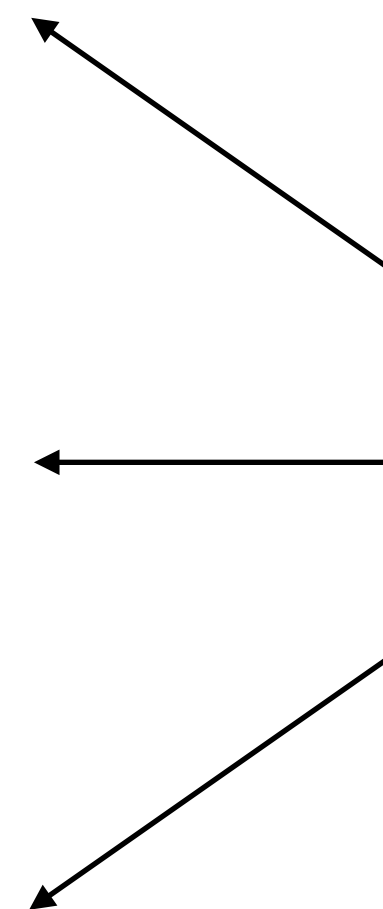
Processor



Observer



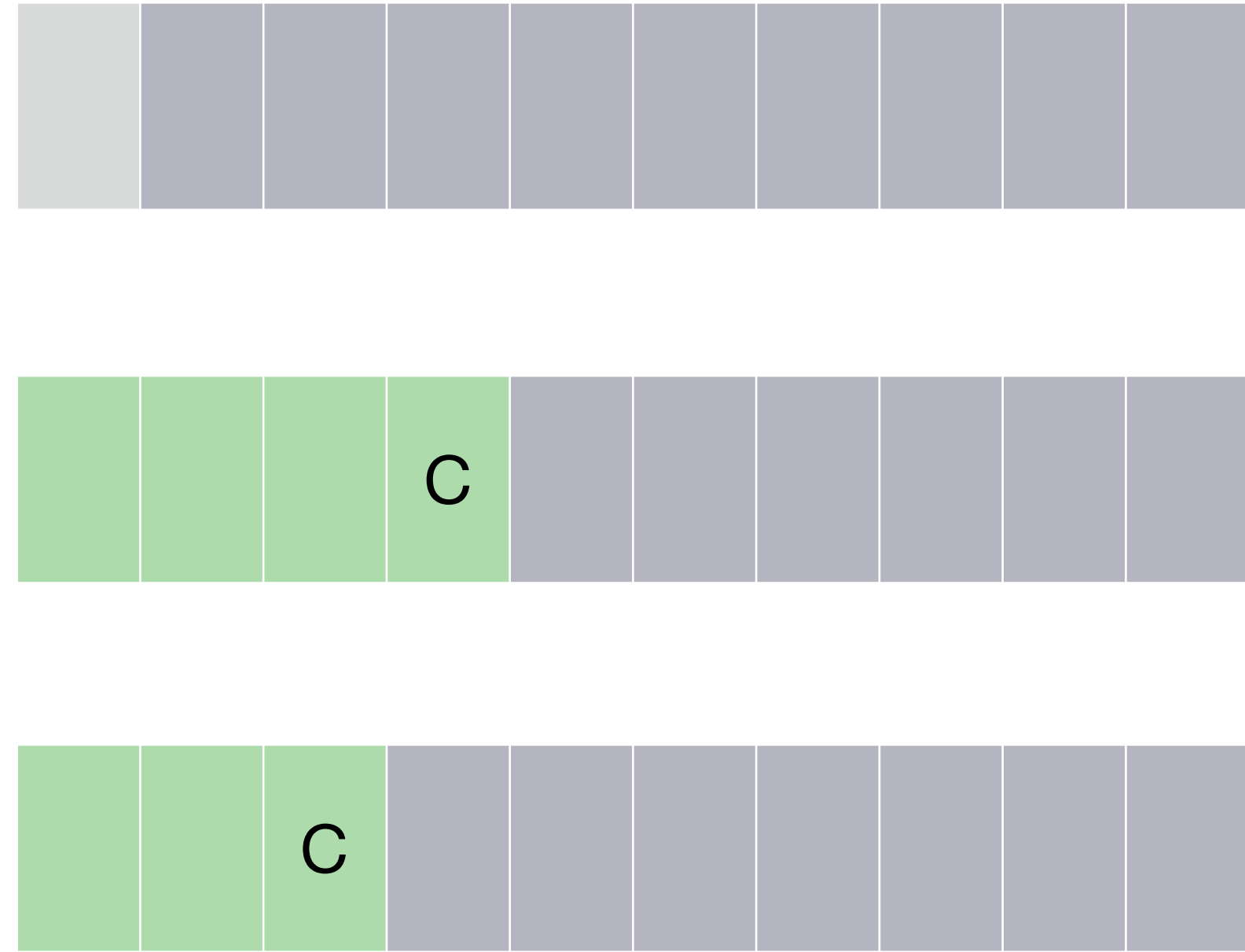
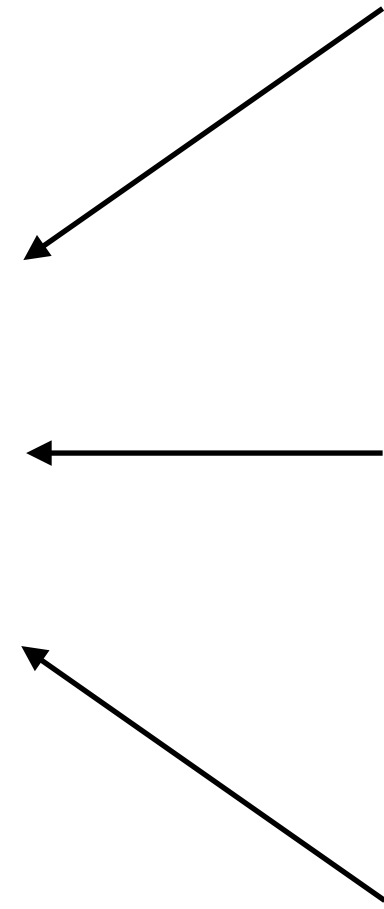
Output



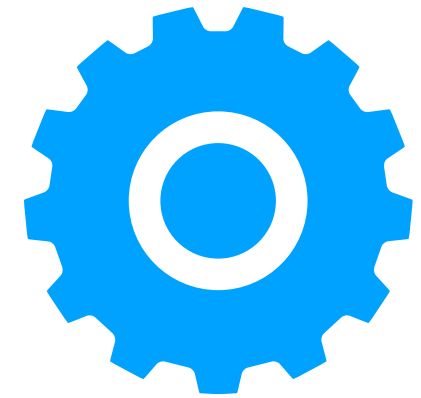
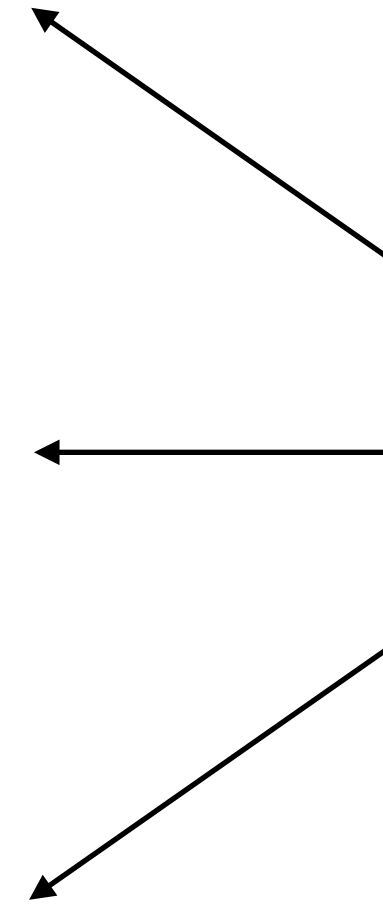
Processor



Observer



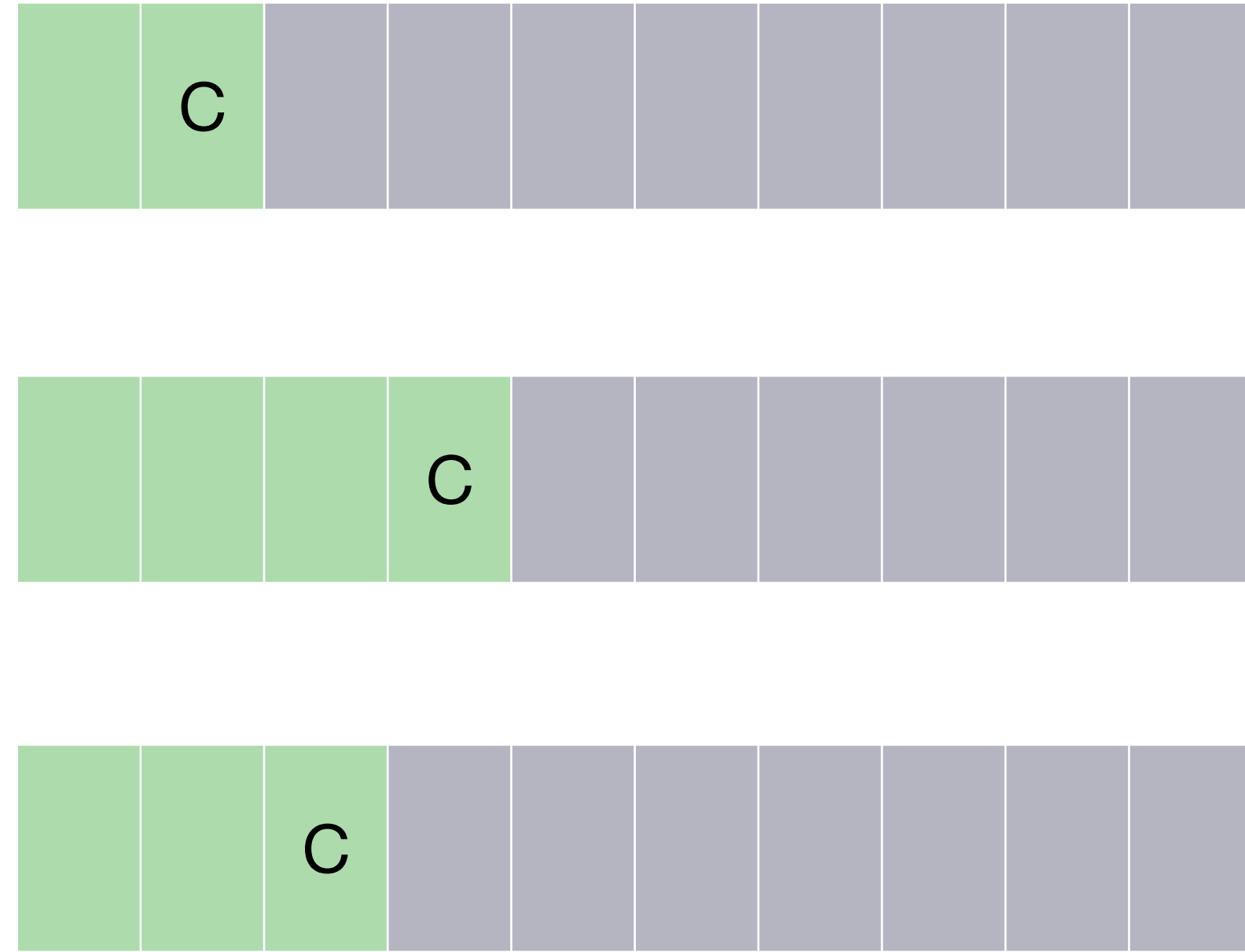
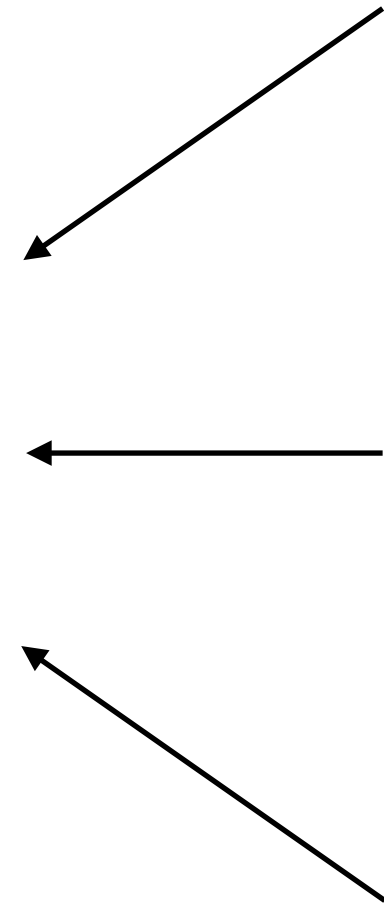
Output



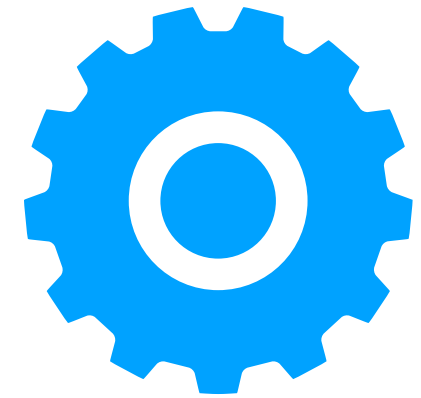
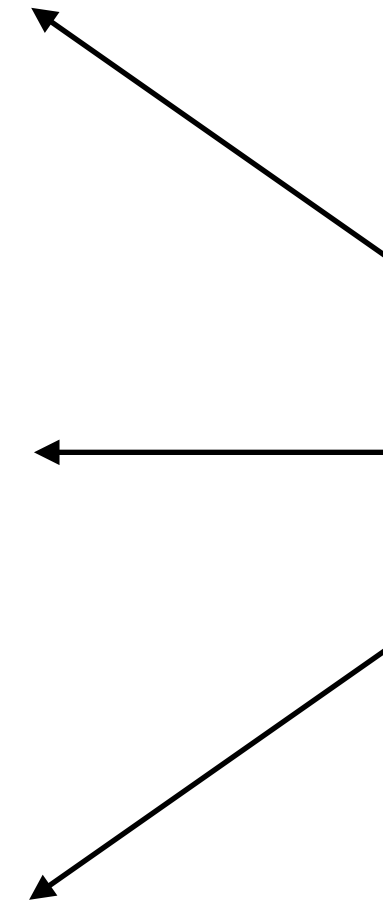
Processor



Observer



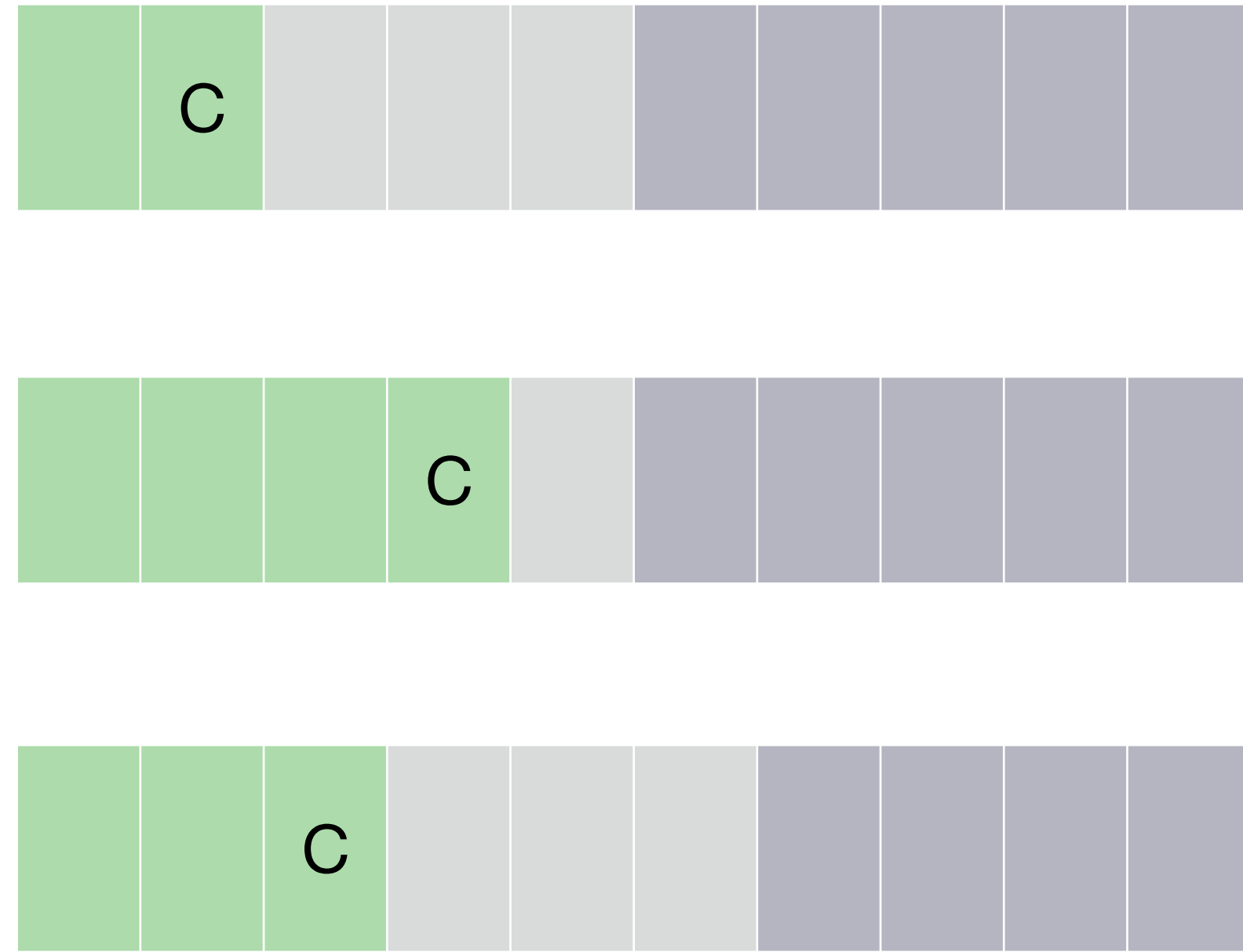
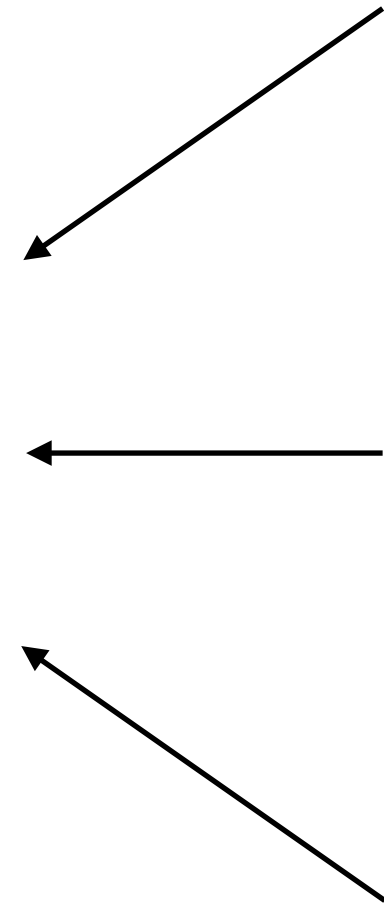
Output



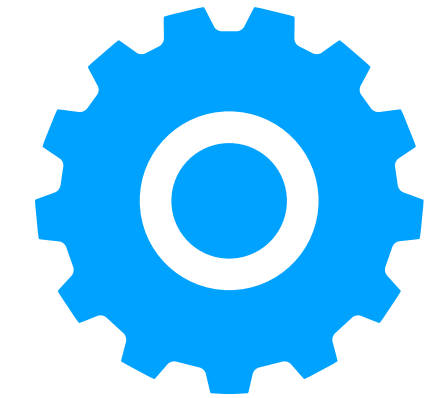
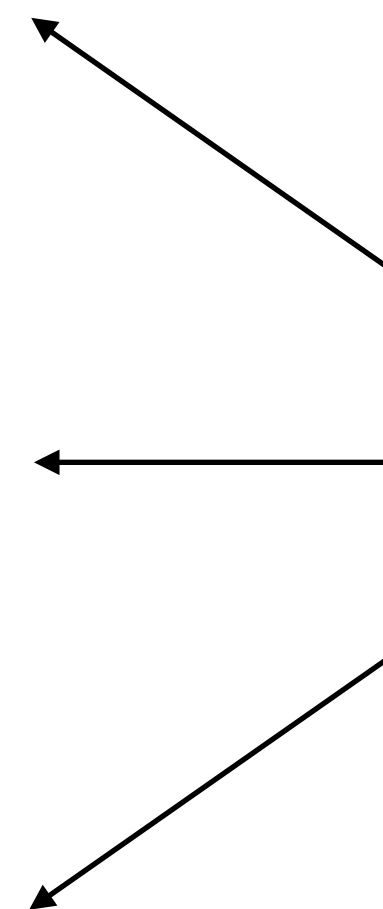
Processor



Observer



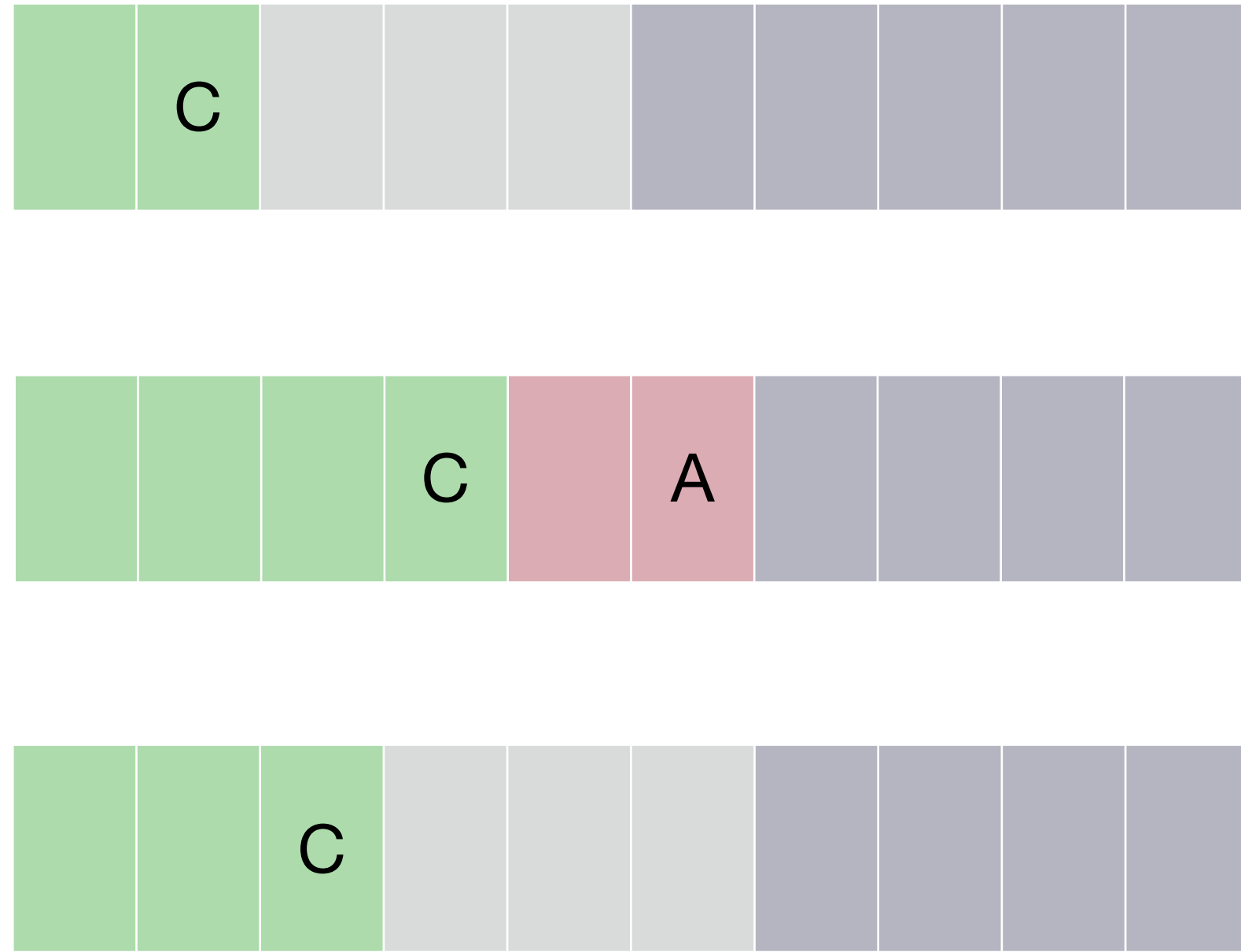
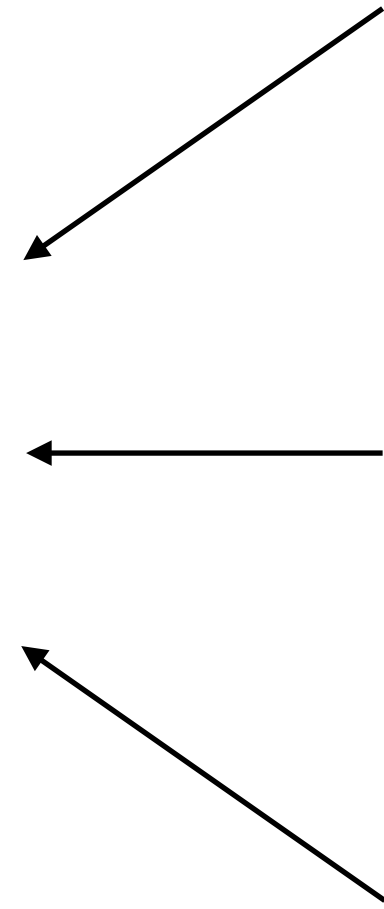
Output



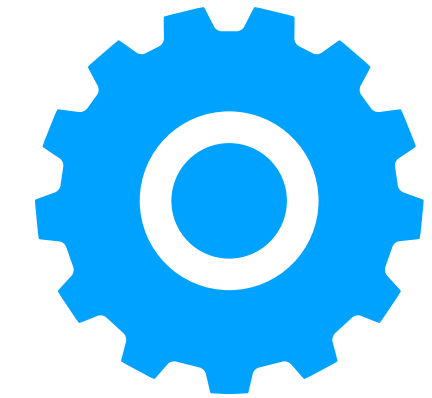
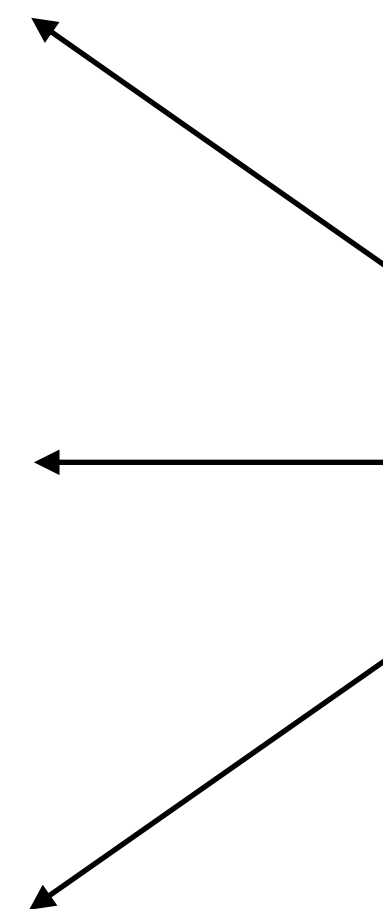
Processor



Observer



Output

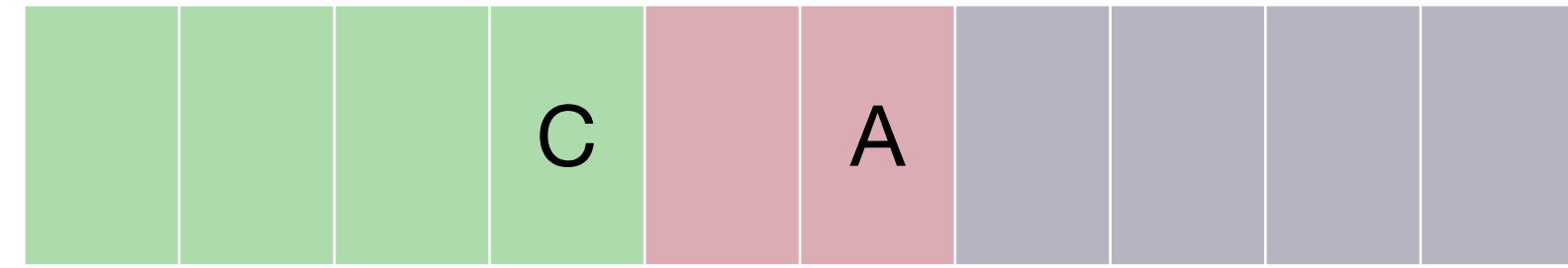


Processor

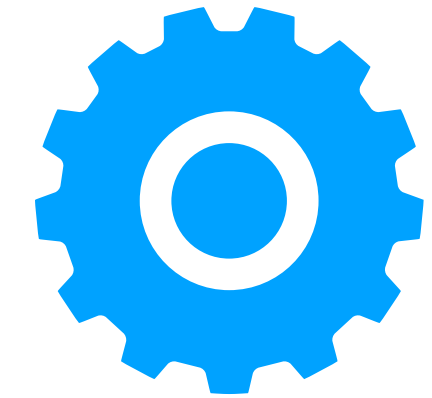
Error



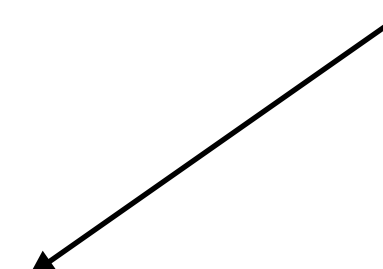
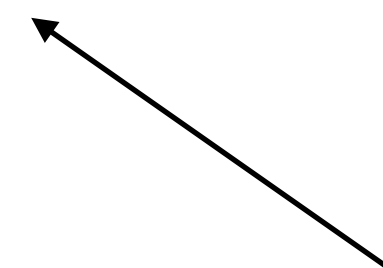
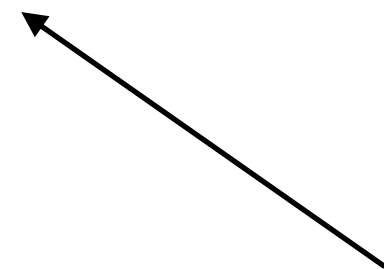
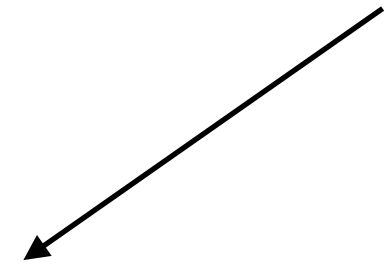
Observer



Output

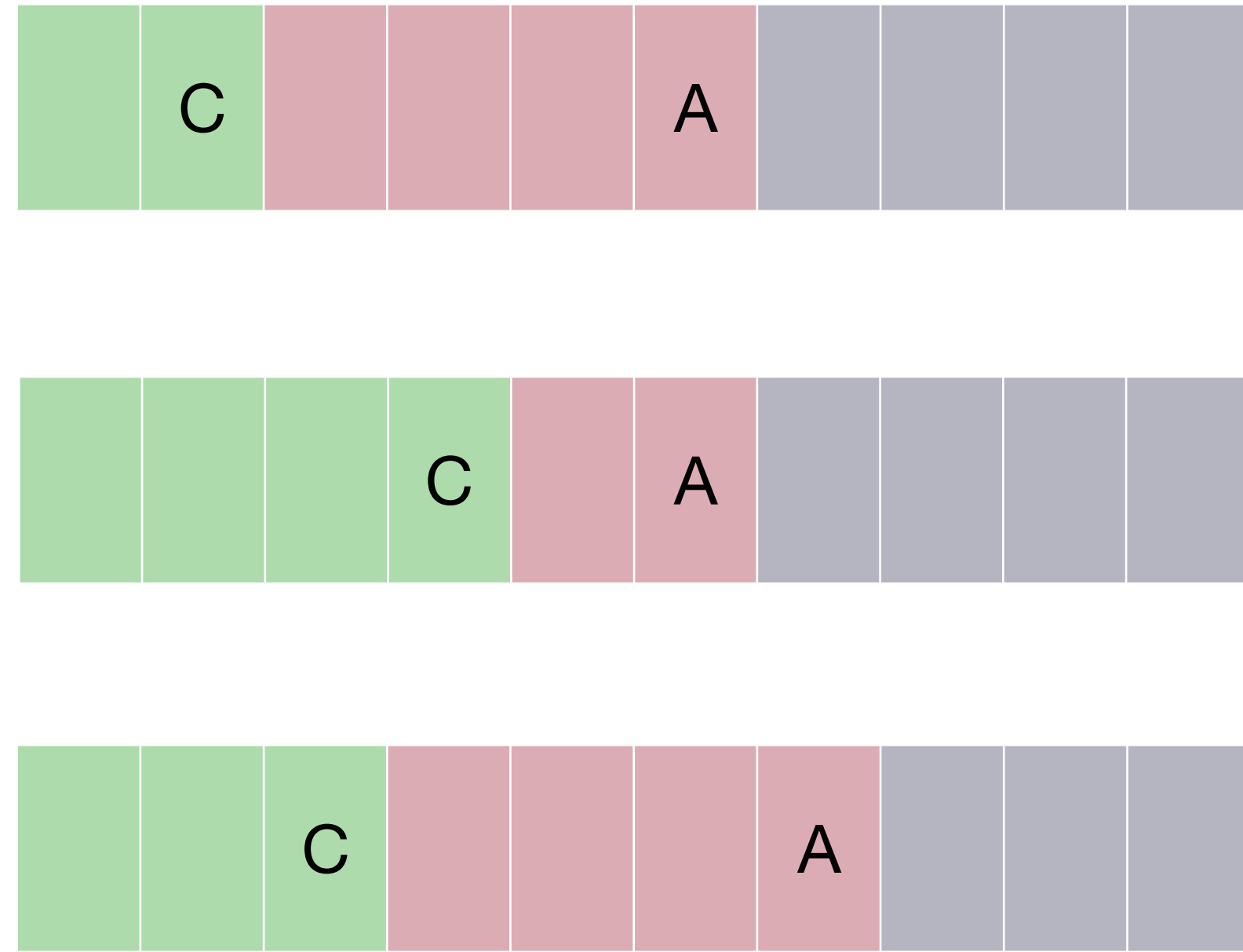
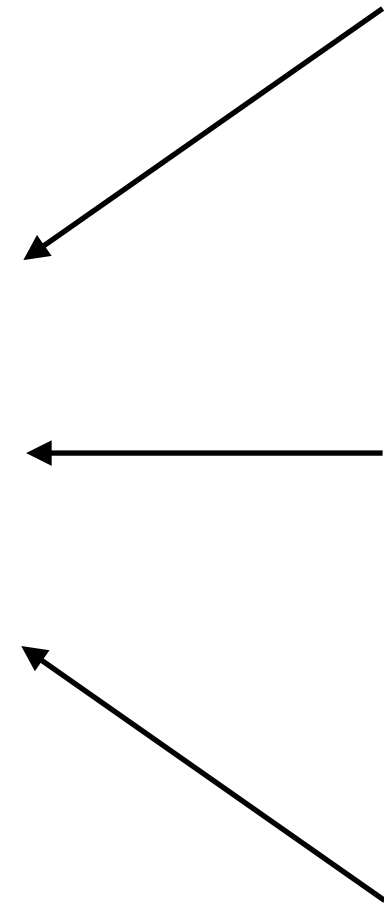


Processor

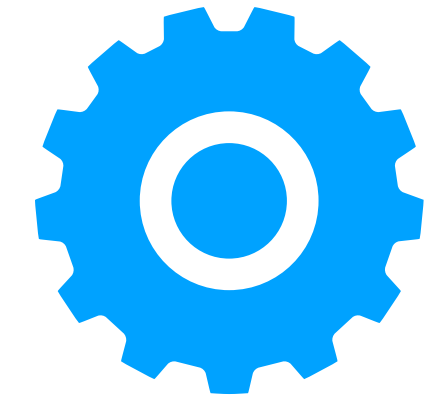
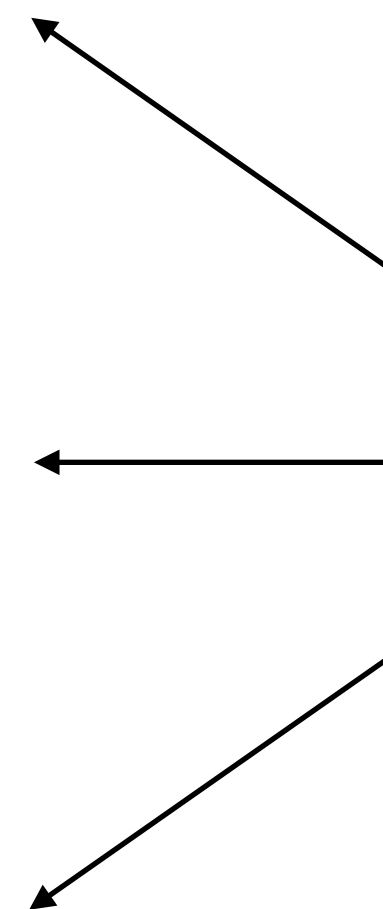




Observer



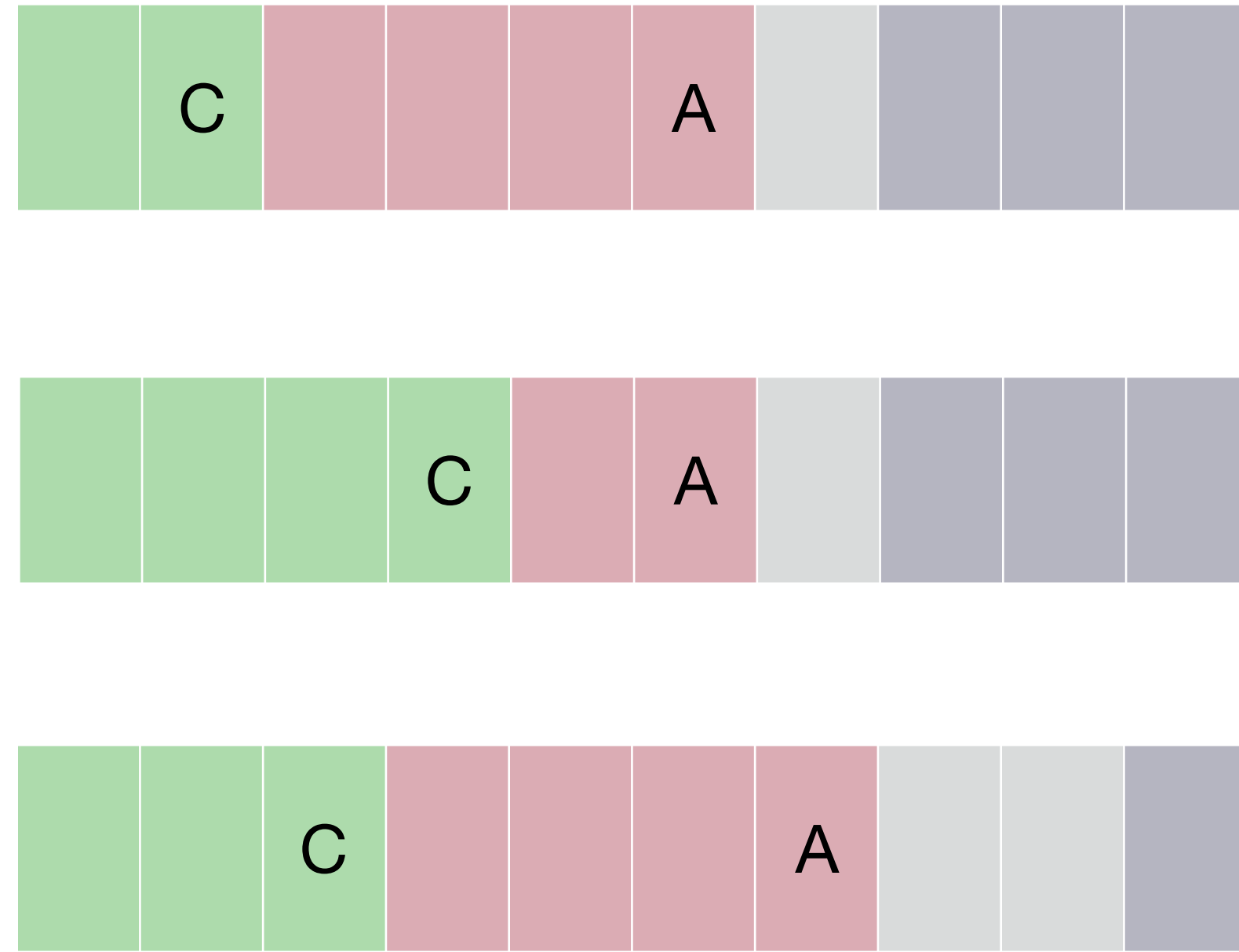
Output



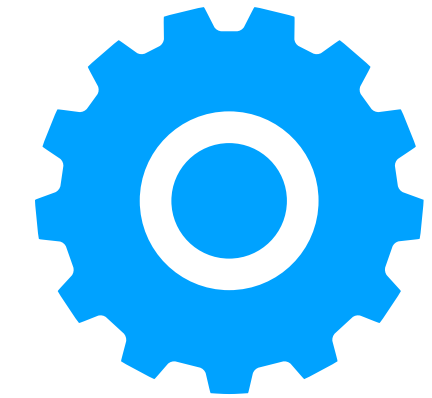
Processor



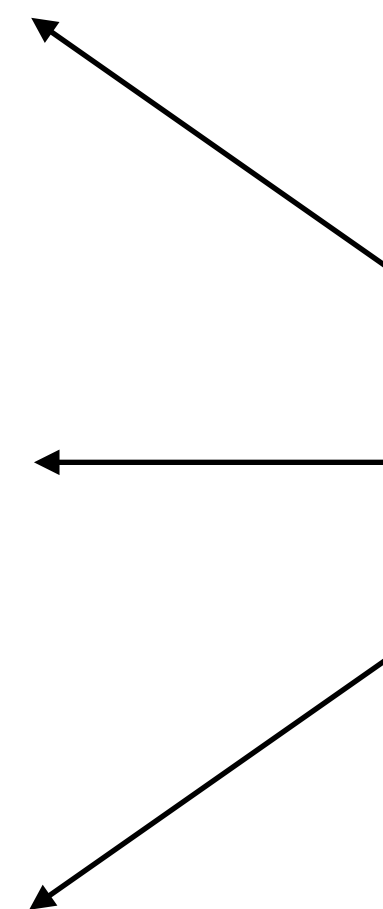
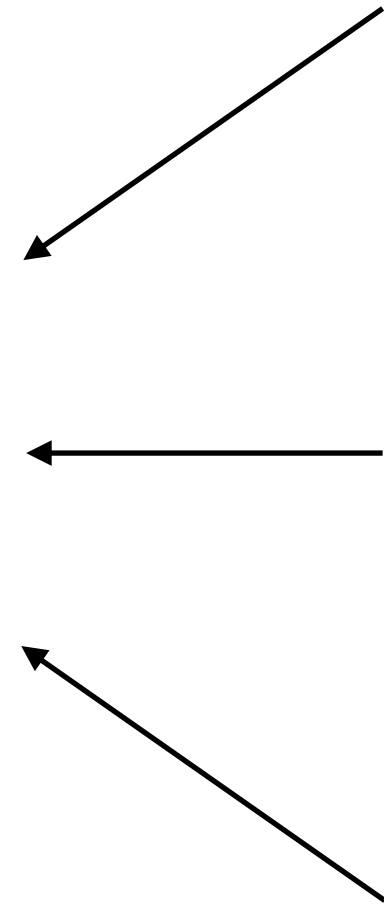
Observer



Output

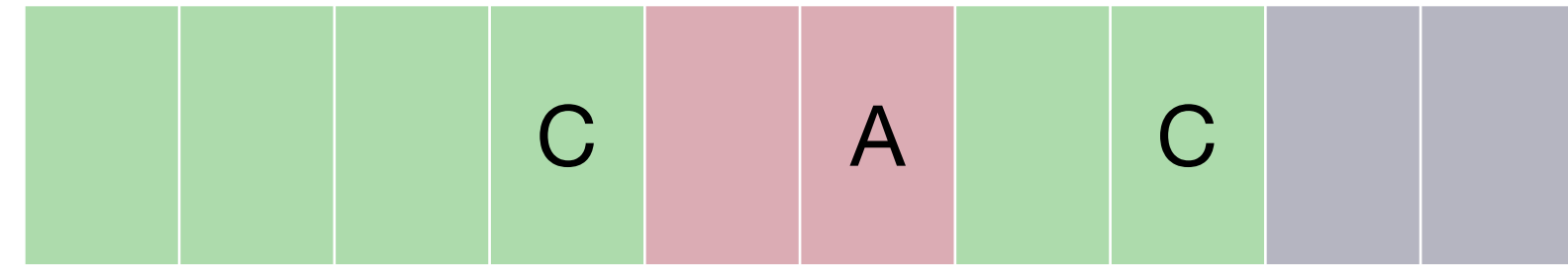
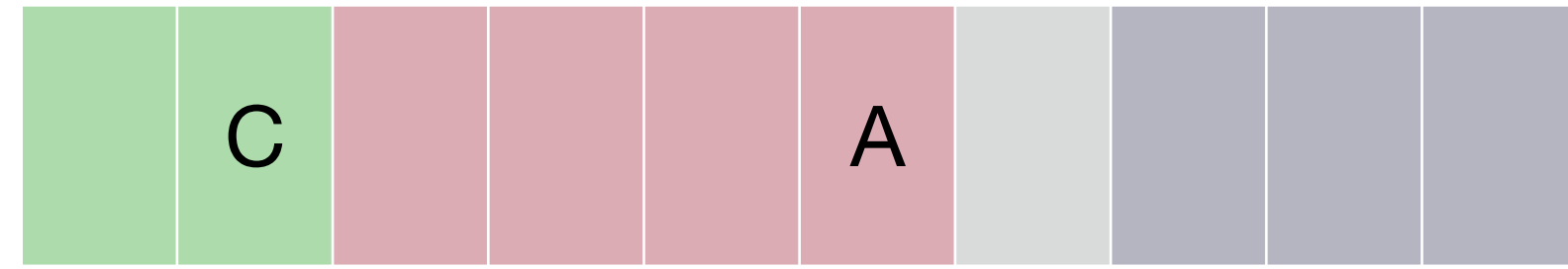


Processor

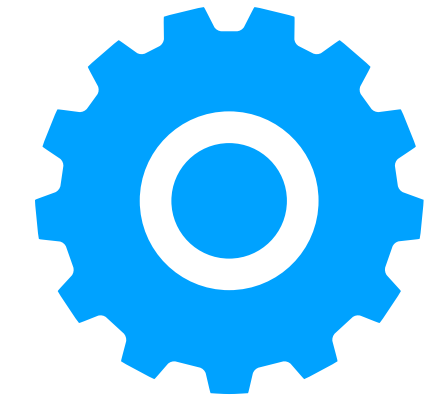




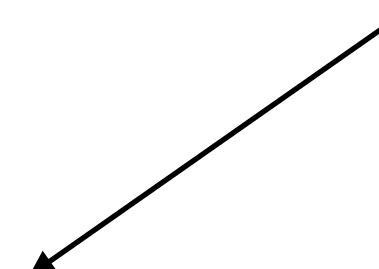
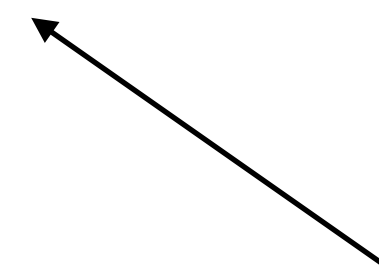
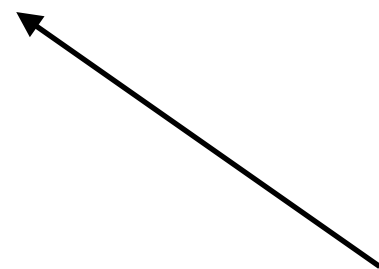
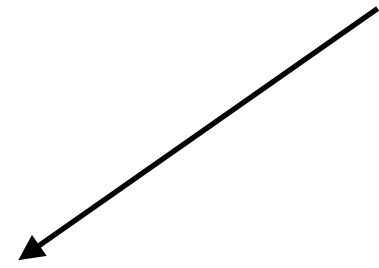
Observer



Output

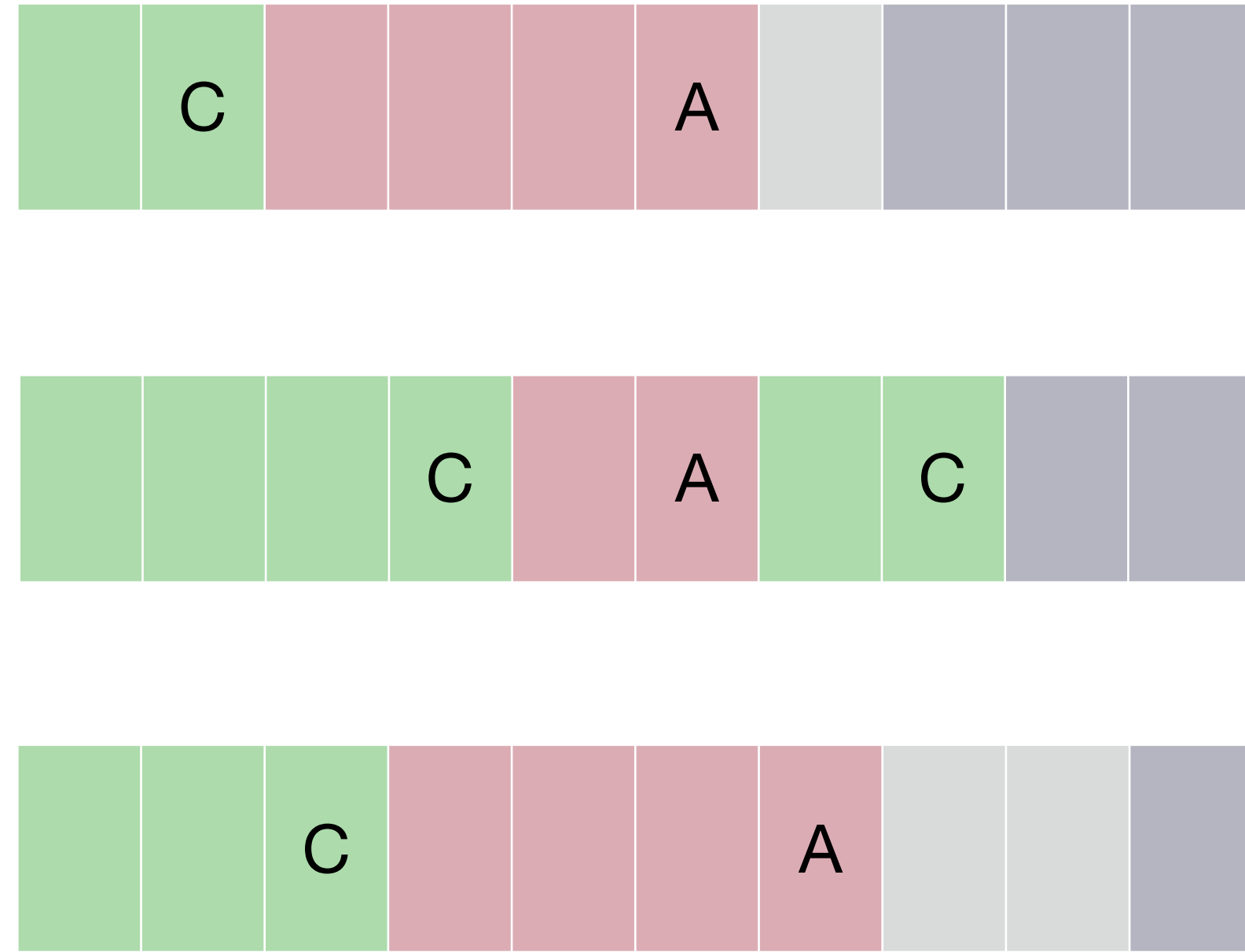
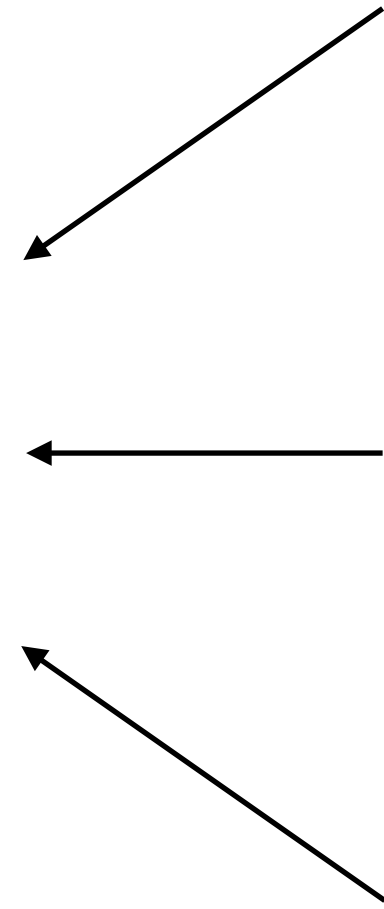


Processor

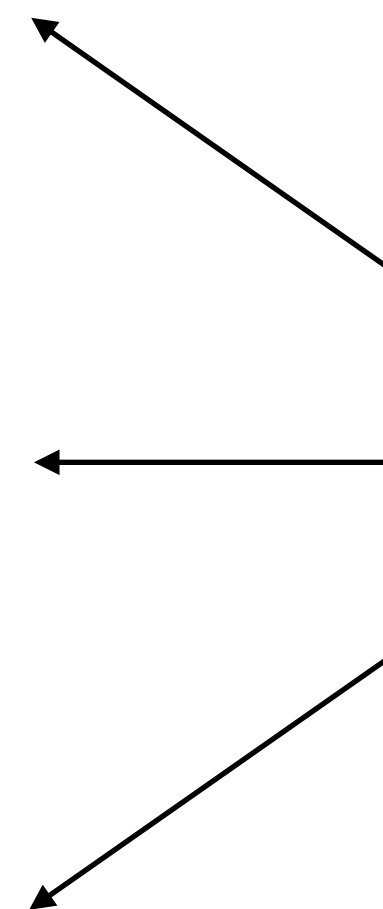




Observer



Output



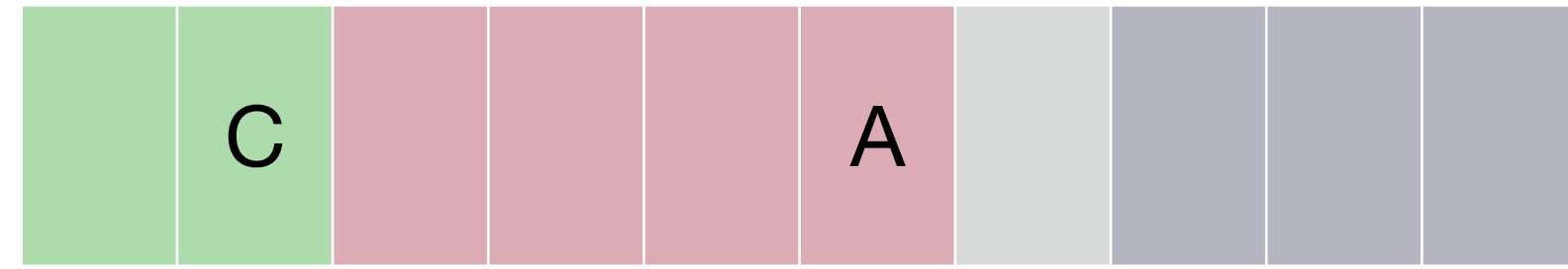
Error



Processor

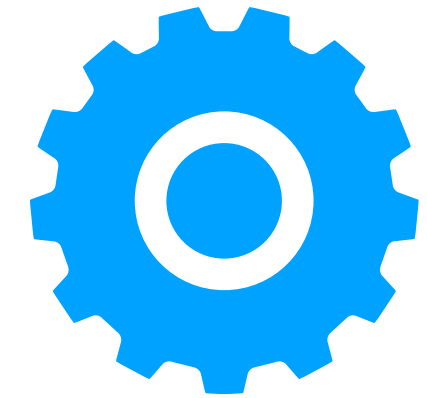


Observer

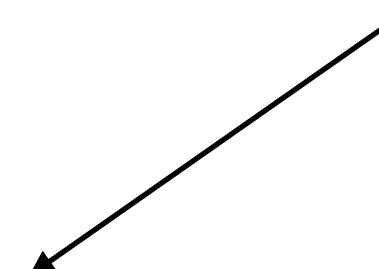
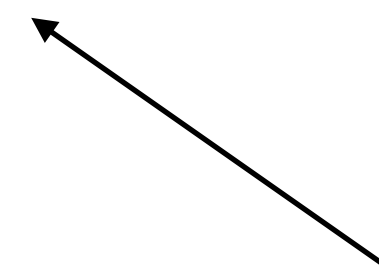
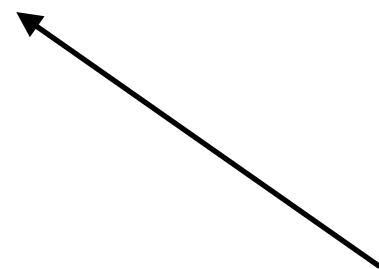
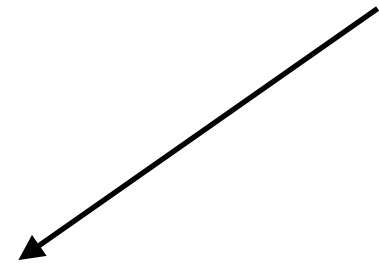


Output

Restart

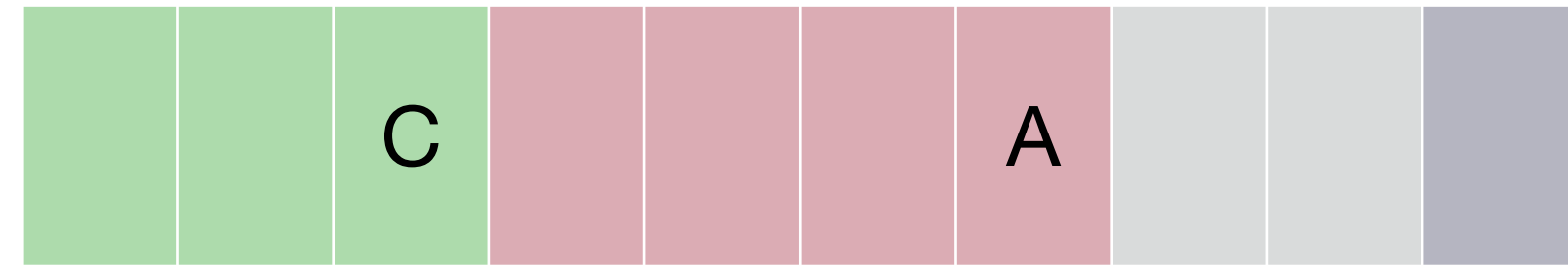
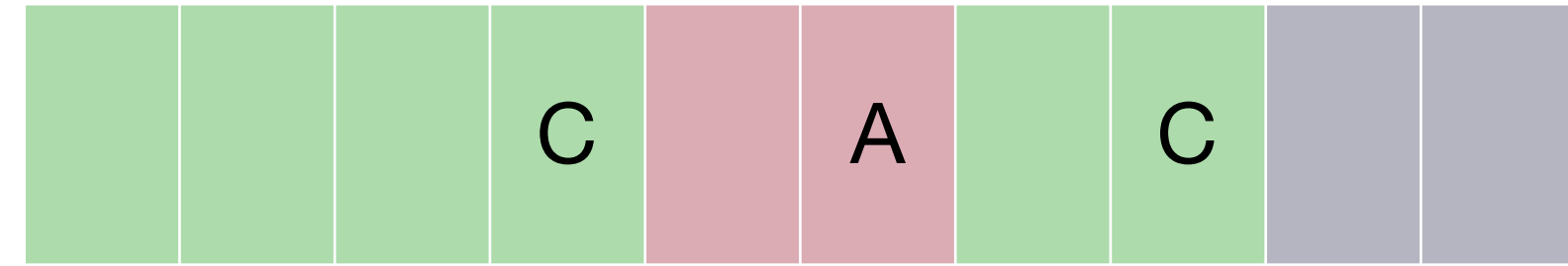
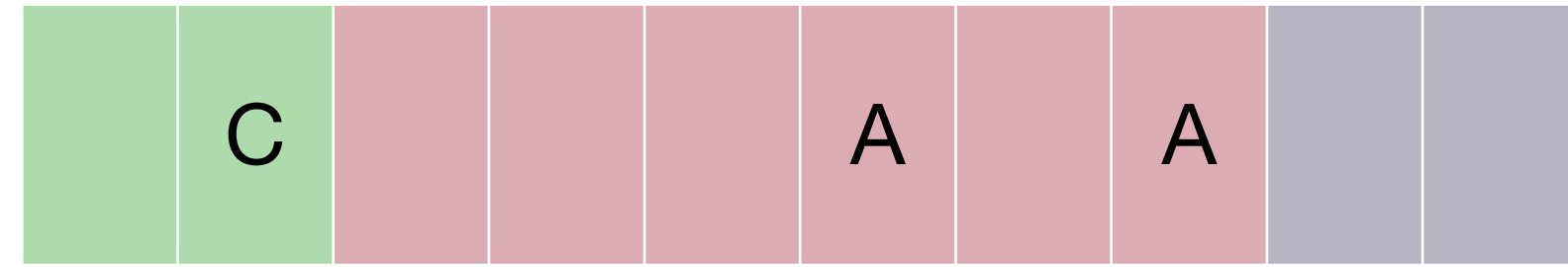


Processor



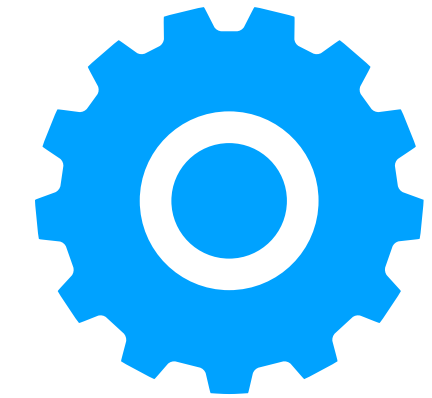


Observer

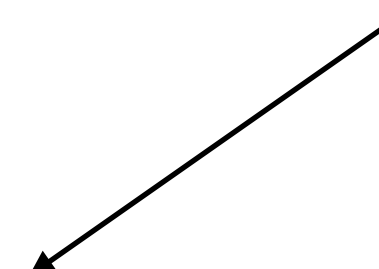
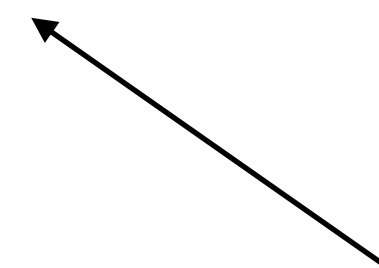
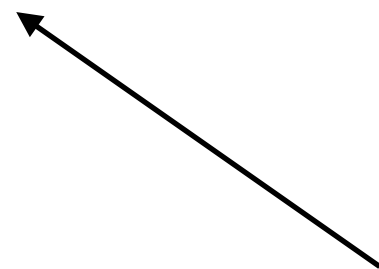
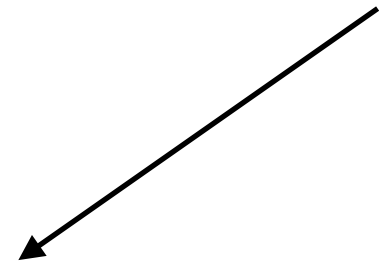


Output

Restart

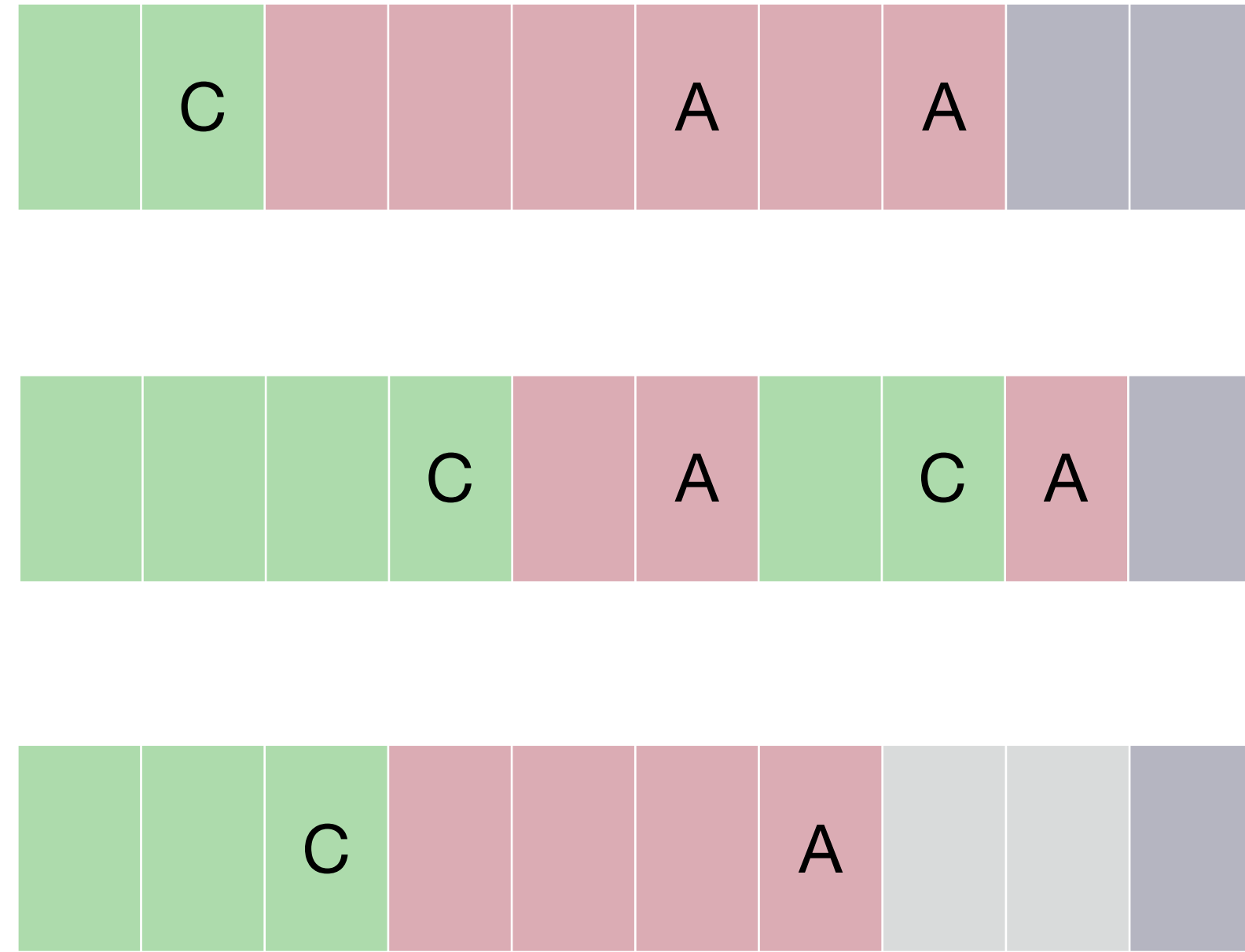


Processor



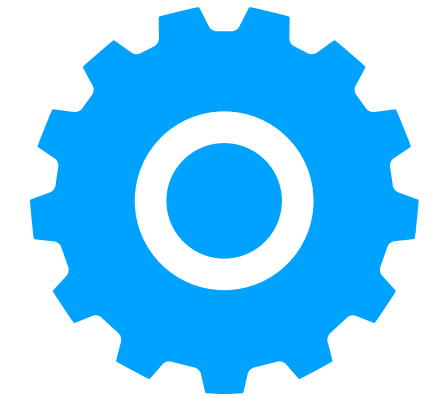


Observer

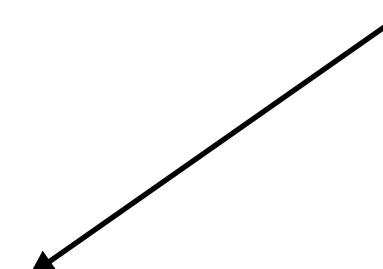
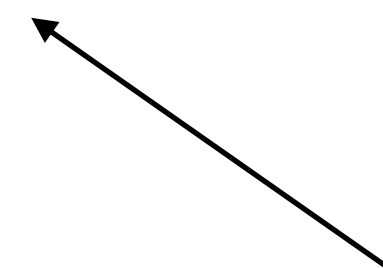
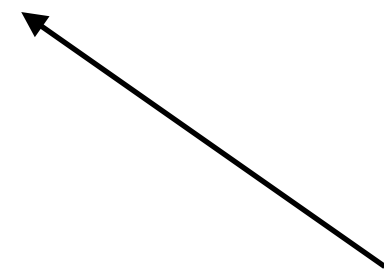
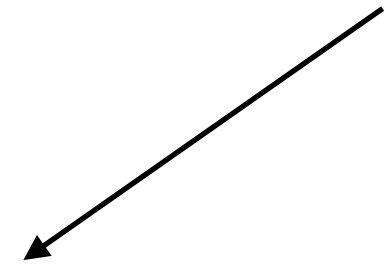


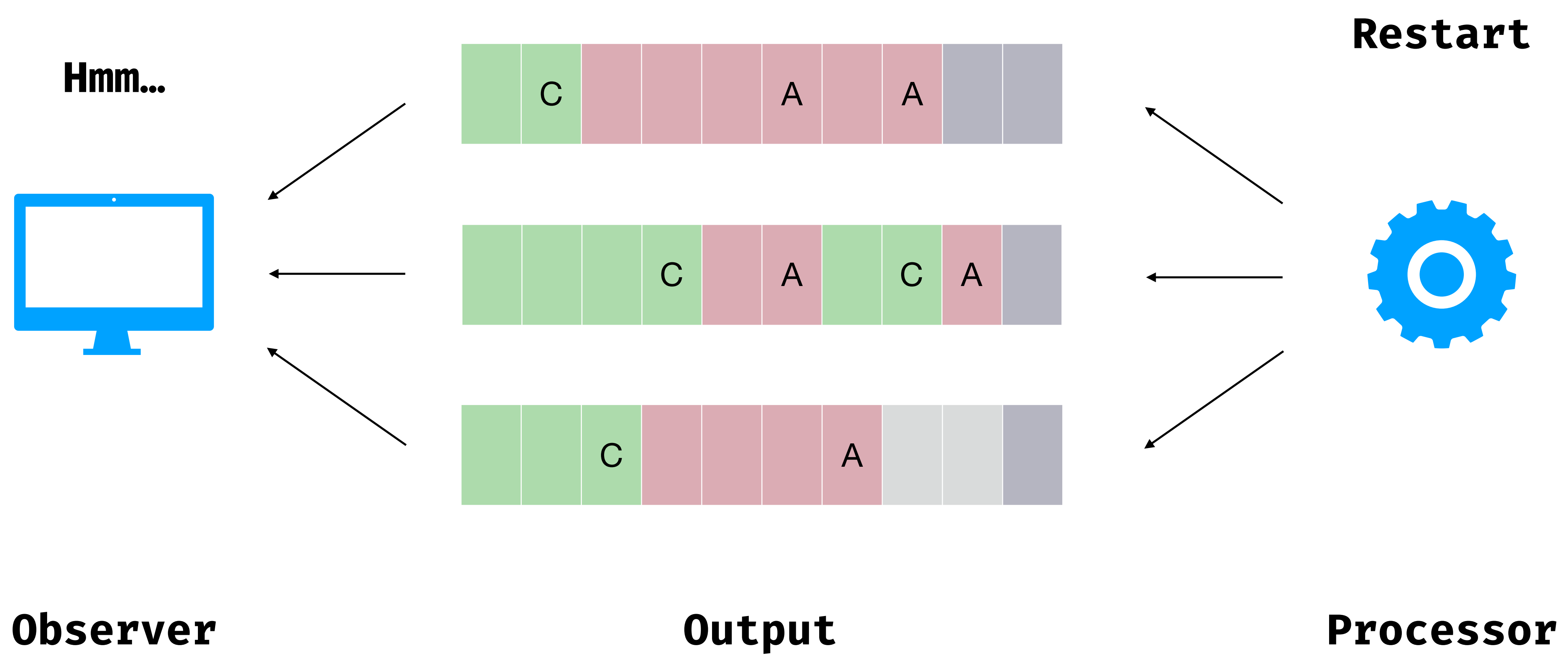
Output

Restart



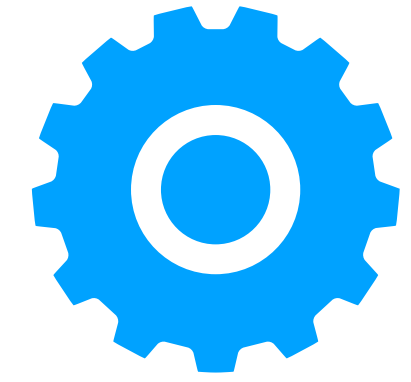
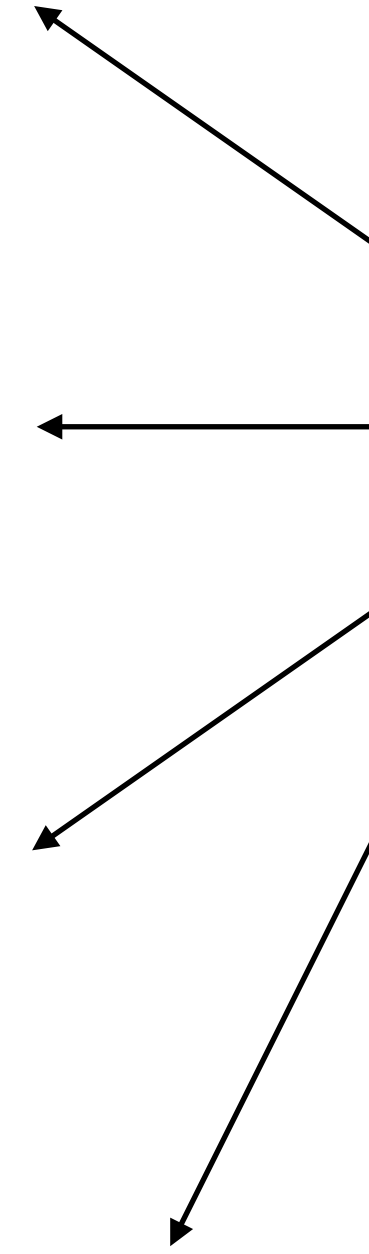
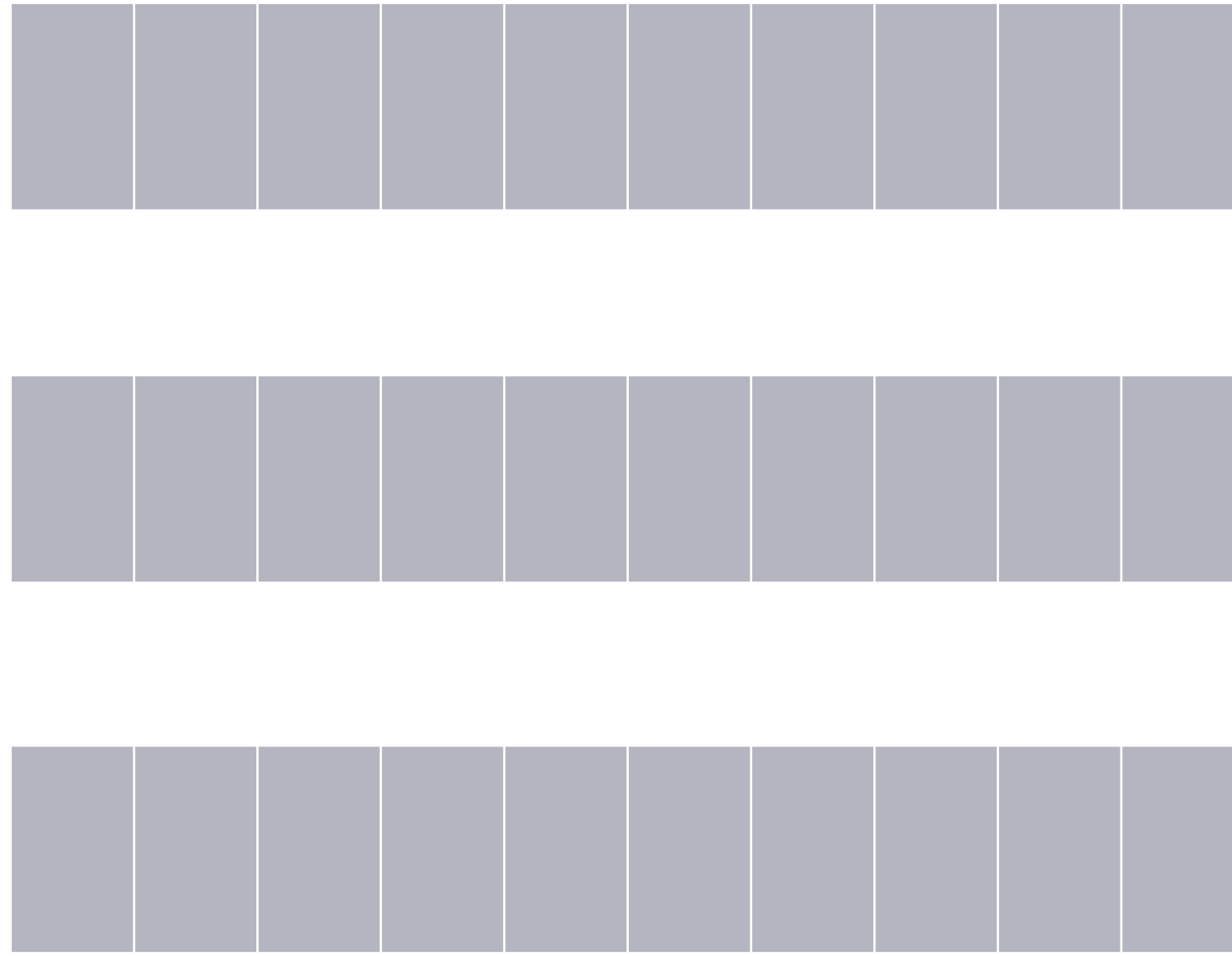
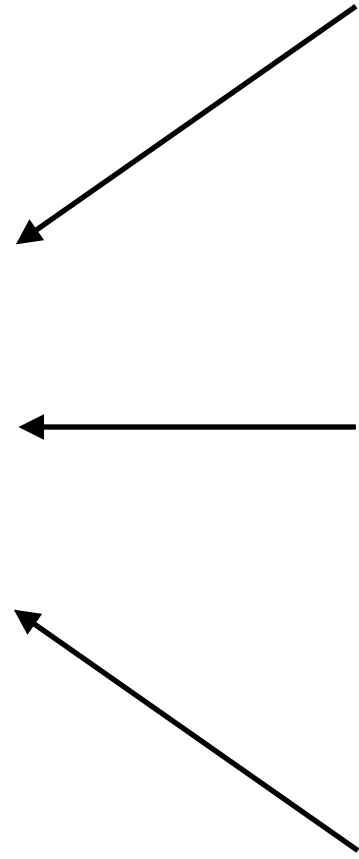
Processor





Транзакции

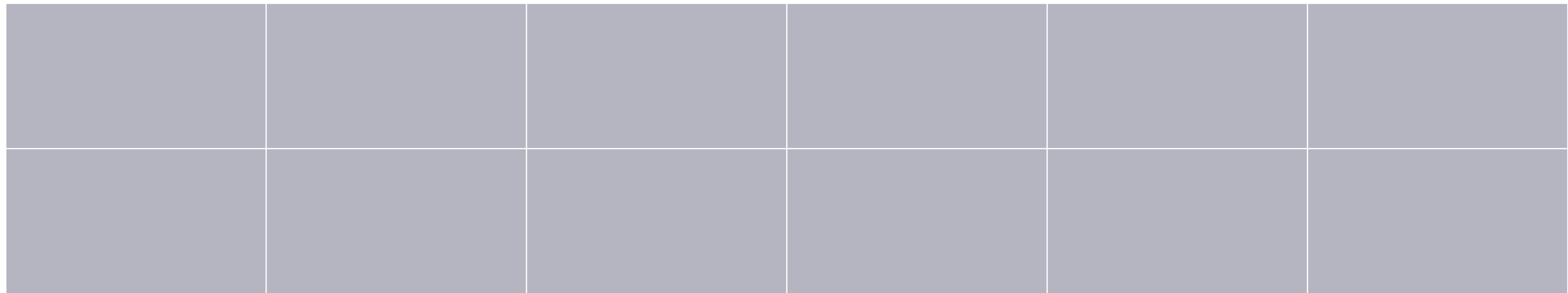
- 2-х фазный коммит и транзакционный лог
 - Запись «желания коммита» в лог
 - Запись маркеров в топичи
 - Запись маркера о завершении транзакции (commit / abort)

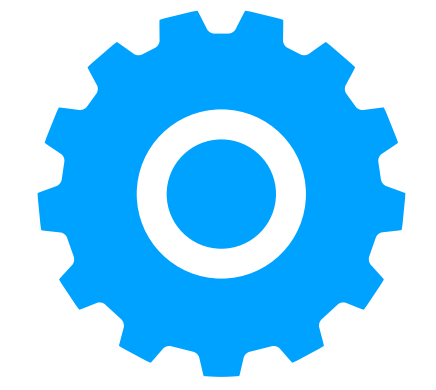
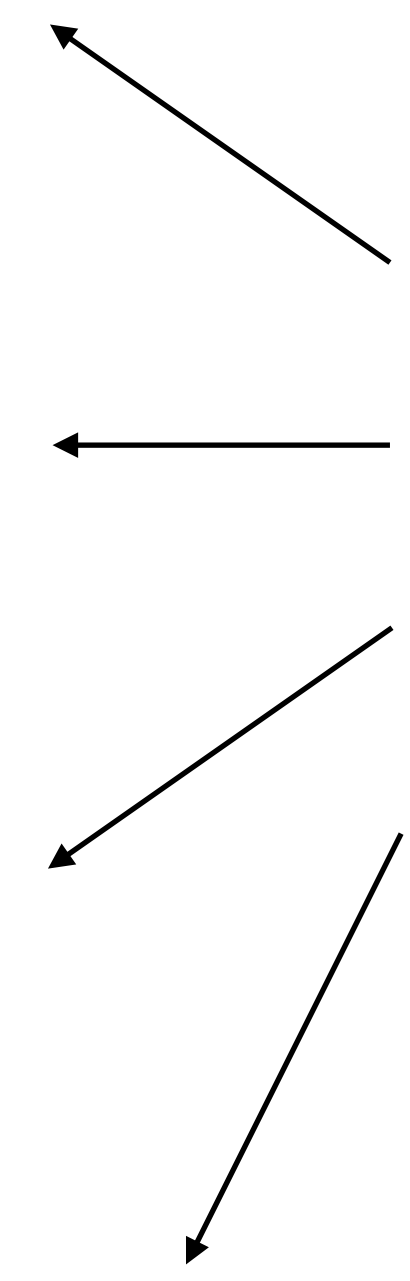
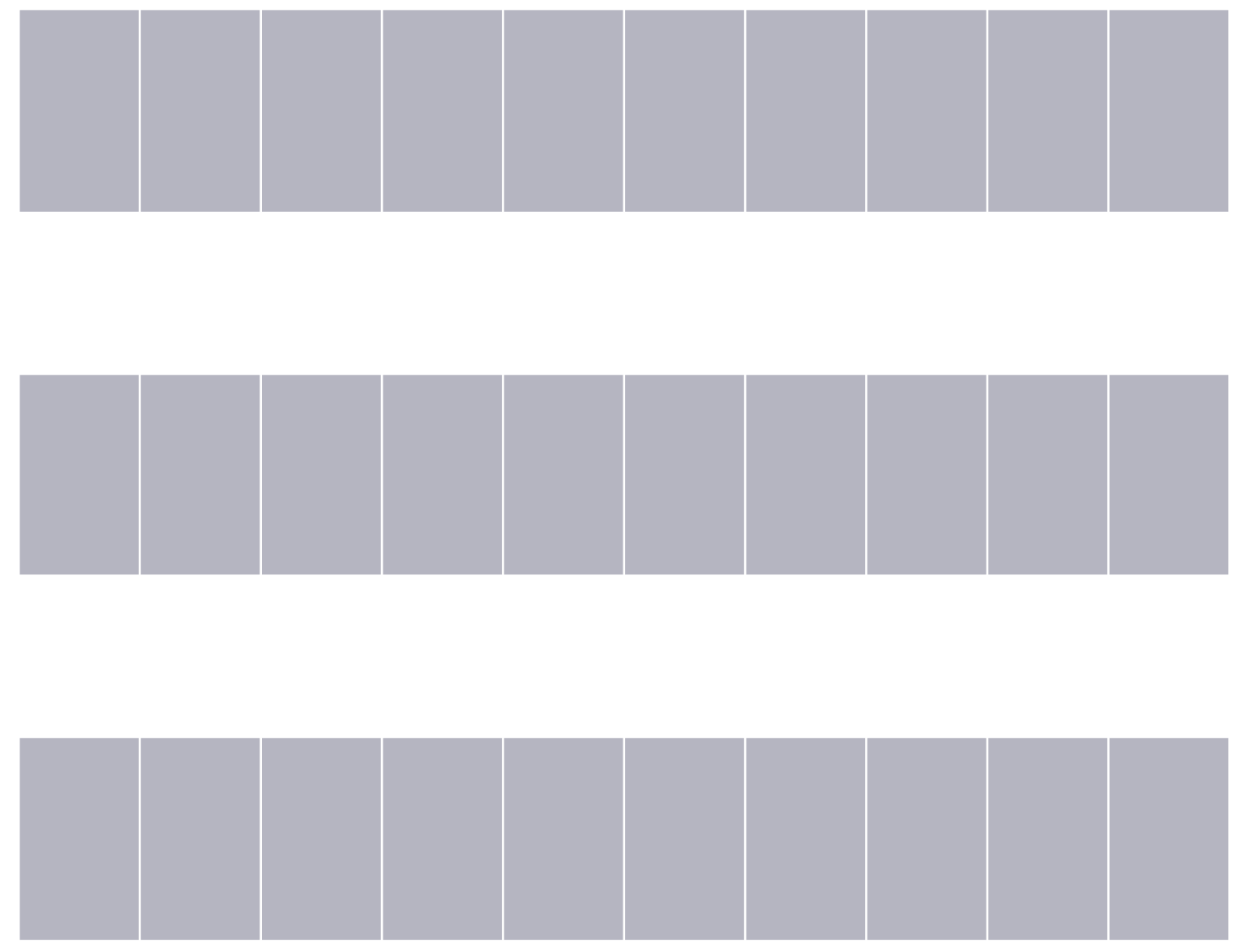
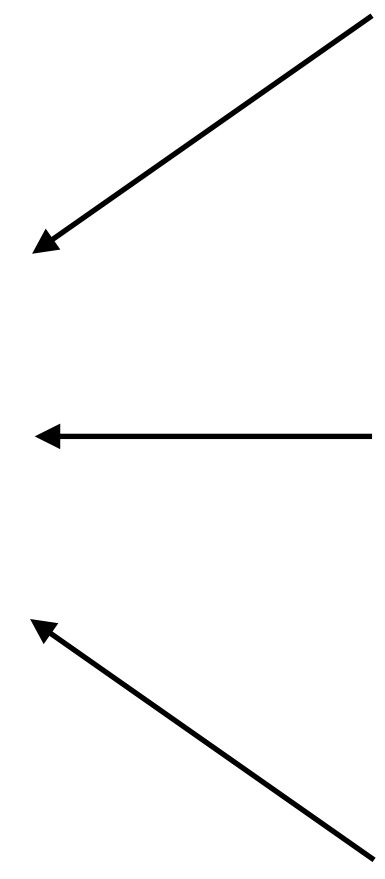


Transaction

Partitions

State





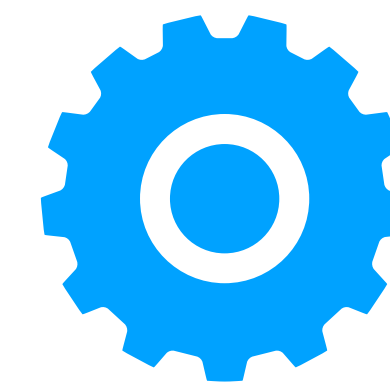
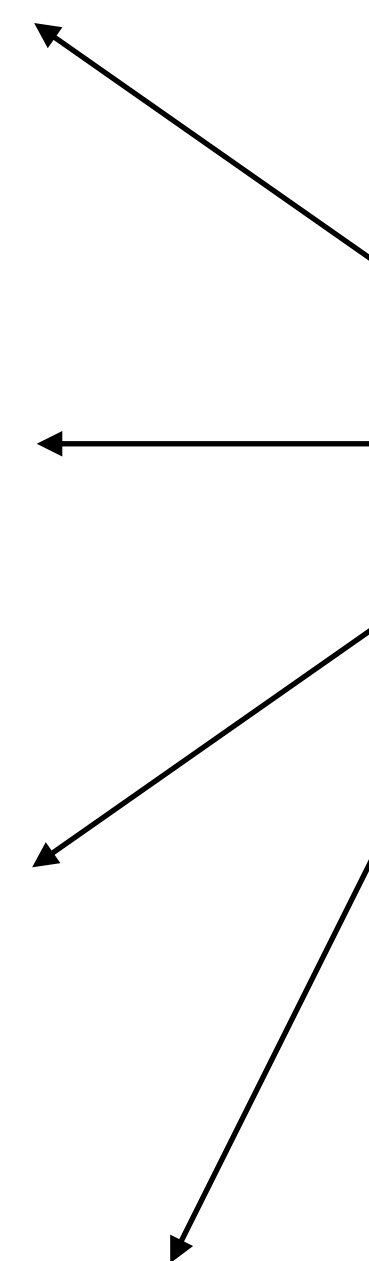
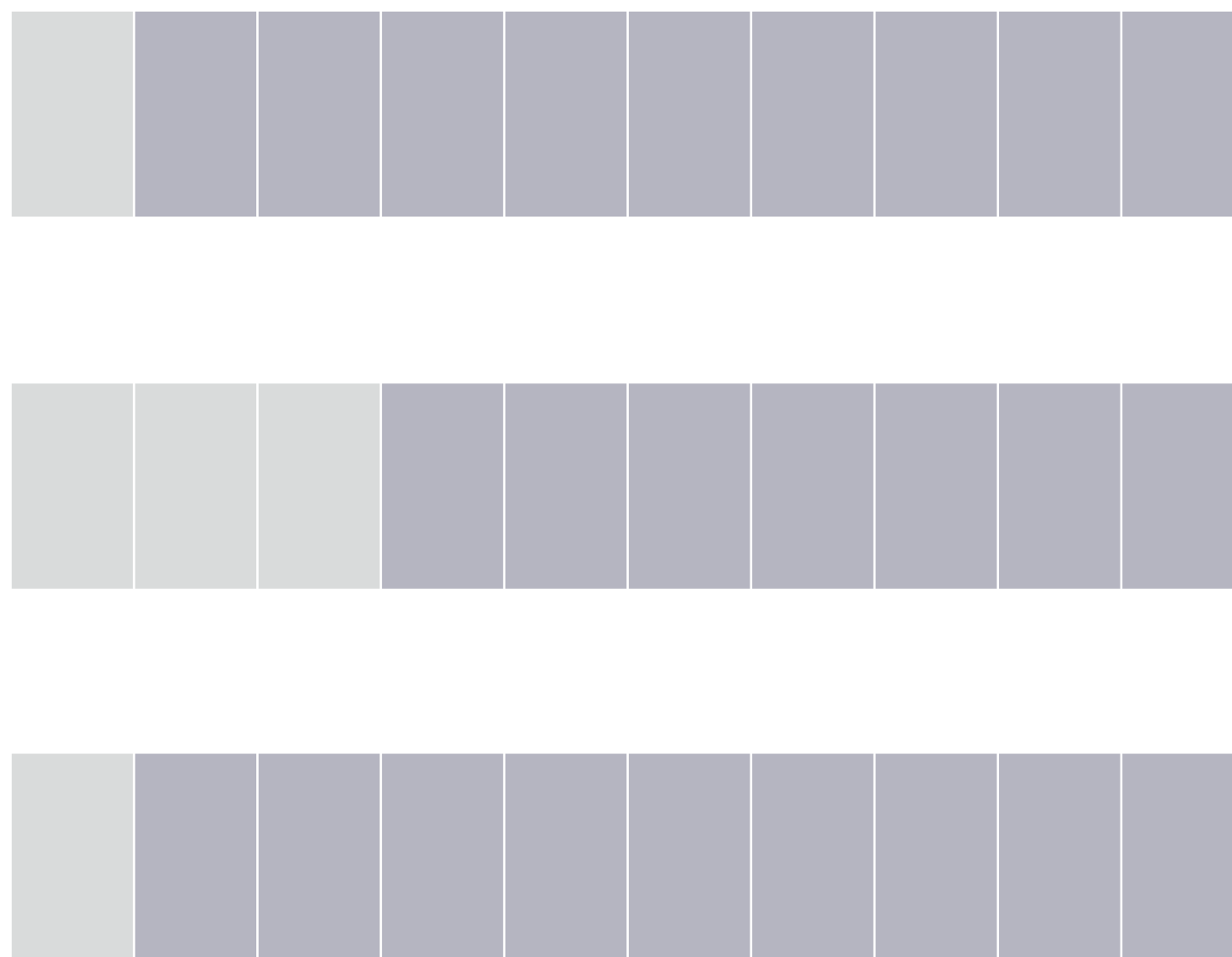
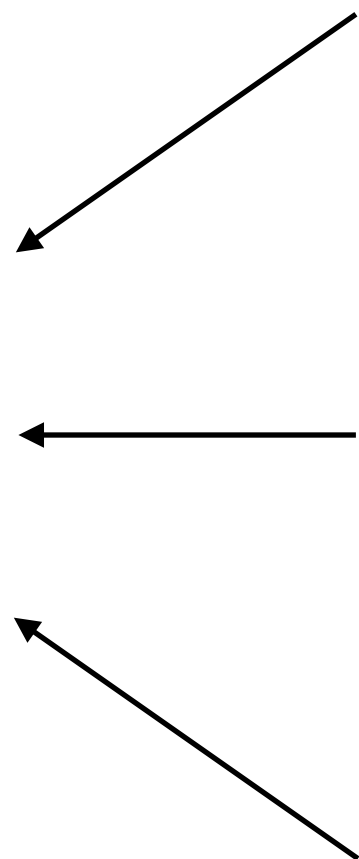
Transaction

Partitions

(0, 1, 2)					
-----------	--	--	--	--	--

State

ongoing					
---------	--	--	--	--	--

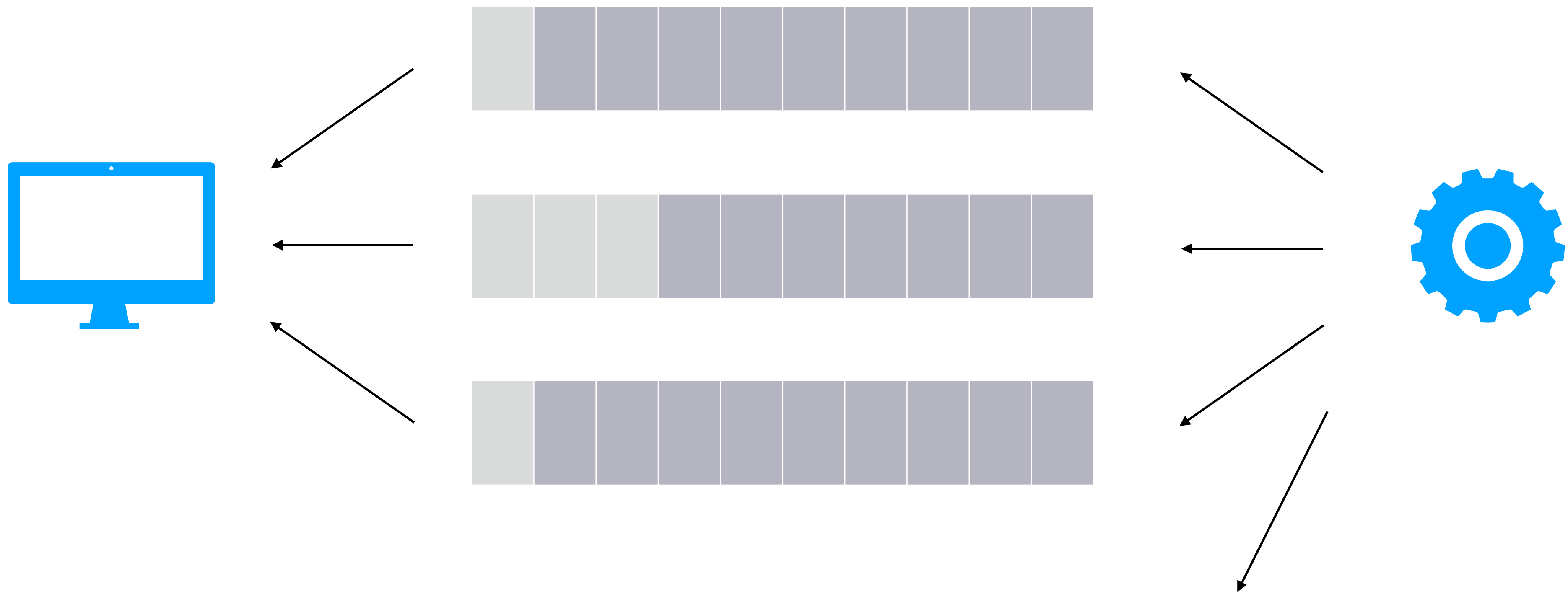


Transaction

Partitions

(0, 1, 2)					
ongoing					

State



Transaction

Partitions

(0, 1, 2)

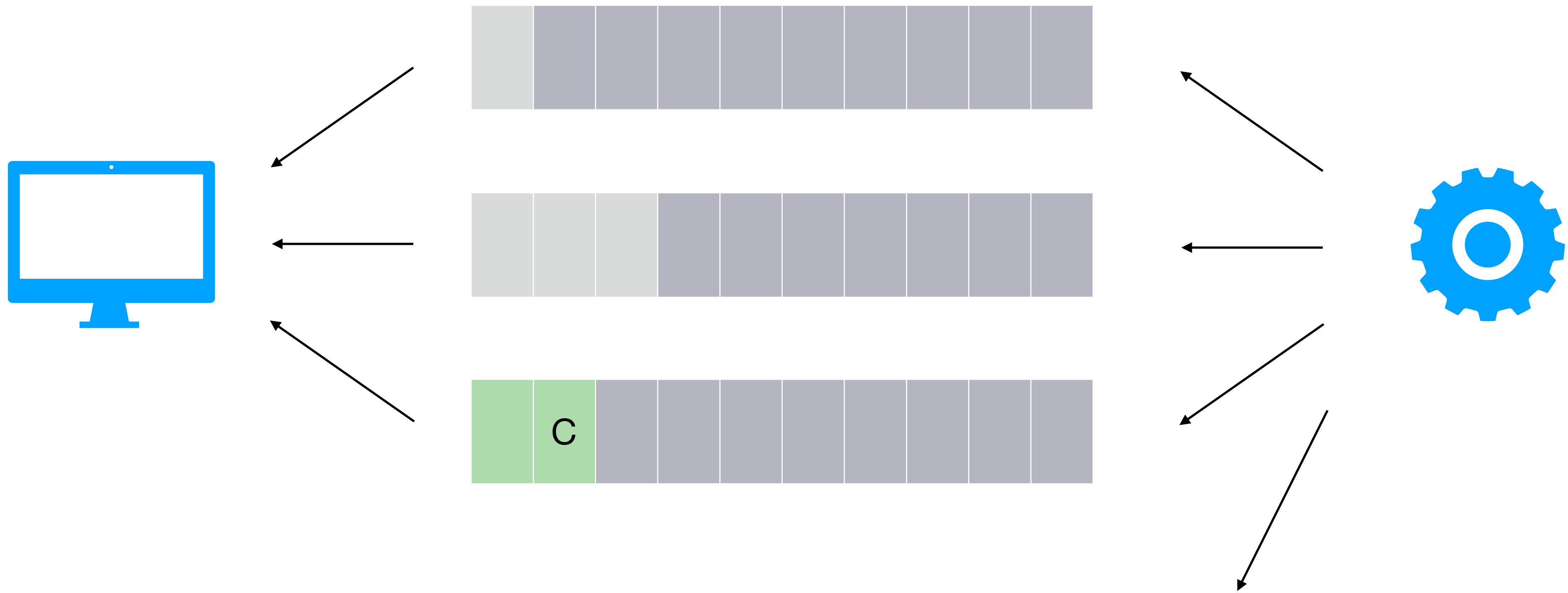
(0, 1, 2)

State

ongoing

prepare
commit

(0, 1, 2)	(0, 1, 2)				
ongoing	prepare commit				



Transaction

Partitions

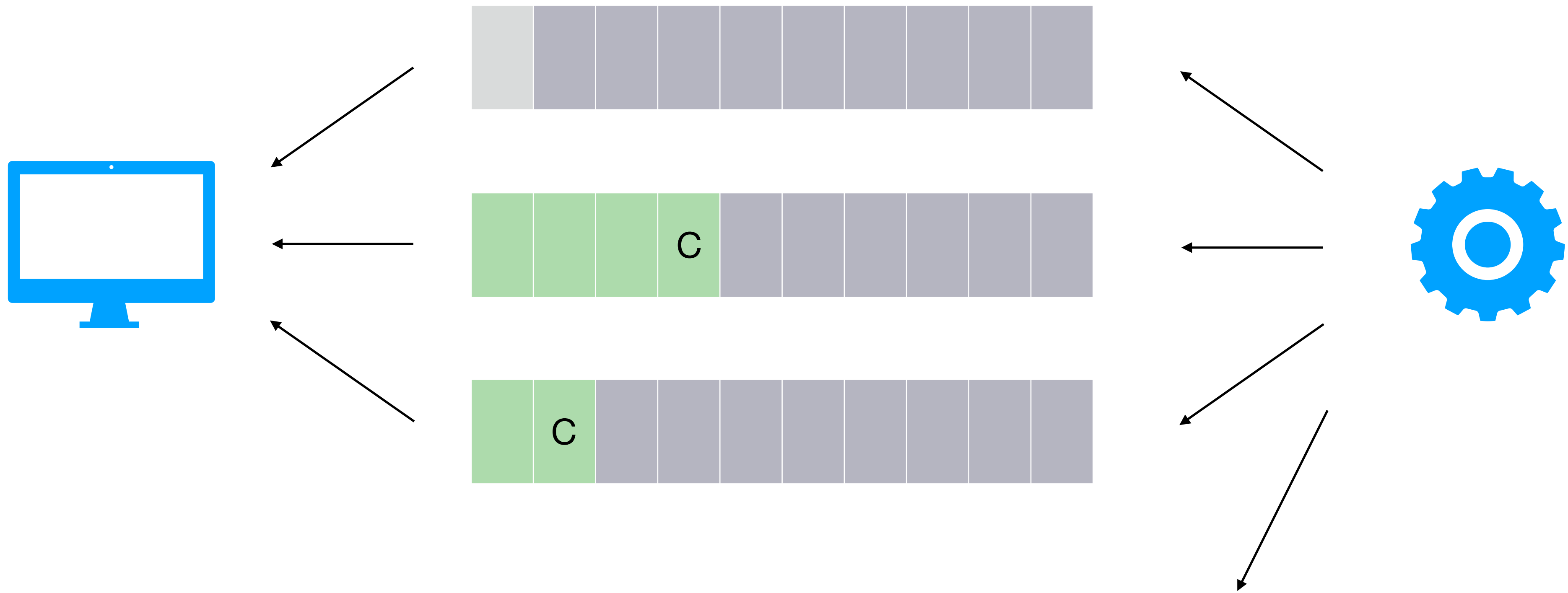
(0, 1, 2)

(0, 1, 2)

State

ongoing

prepare
commit



Transaction

Partitions

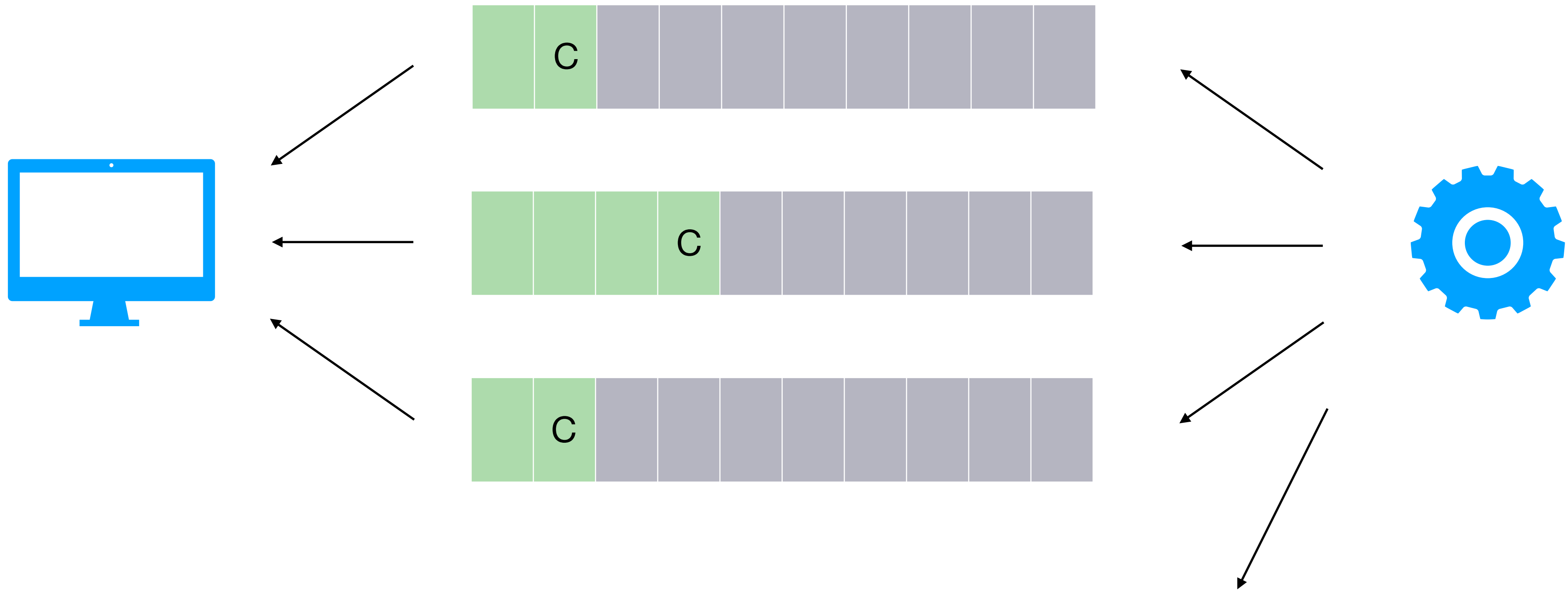
(0, 1, 2)

(0, 1, 2)

State

ongoing

prepare
commit



Transaction

Partitions

(0, 1, 2)

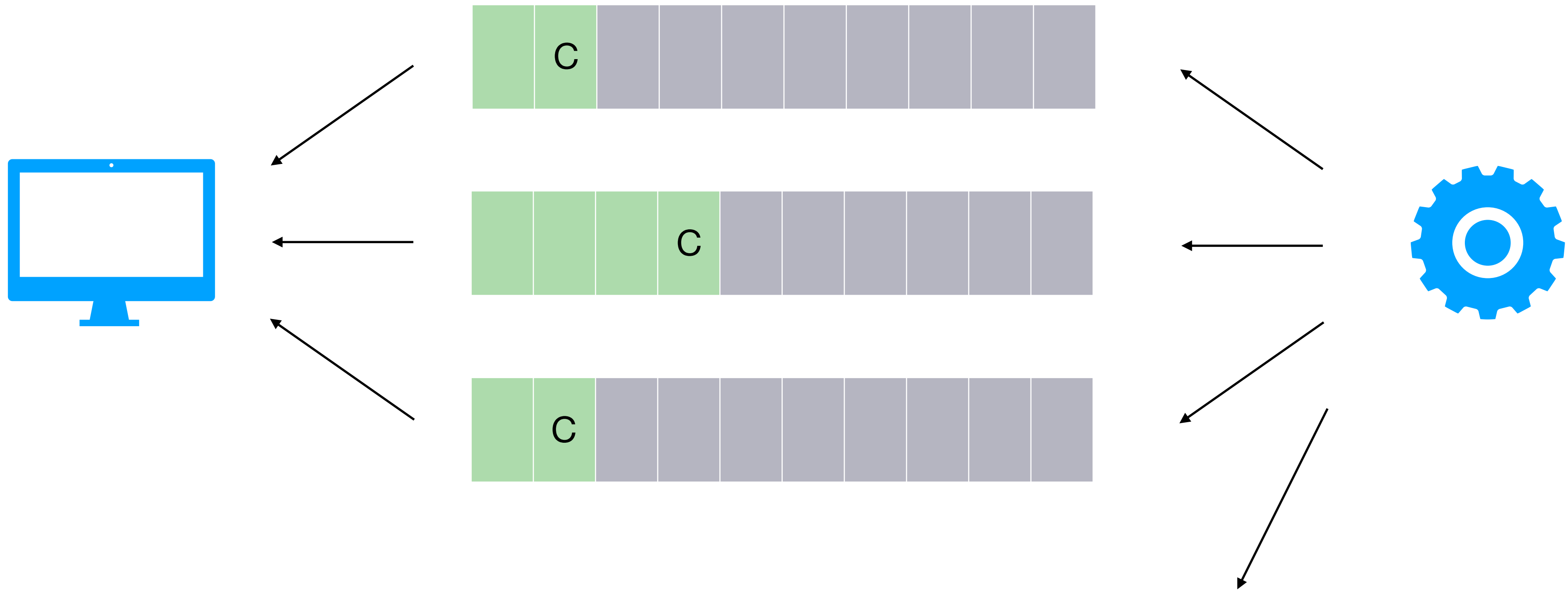
(0, 1, 2)

State

ongoing

prepare
commit

(0, 1, 2)	(0, 1, 2)				
ongoing	prepare commit				



Transaction

Partitions

(0, 1, 2)

(0, 1, 2)

(0, 1, 2)

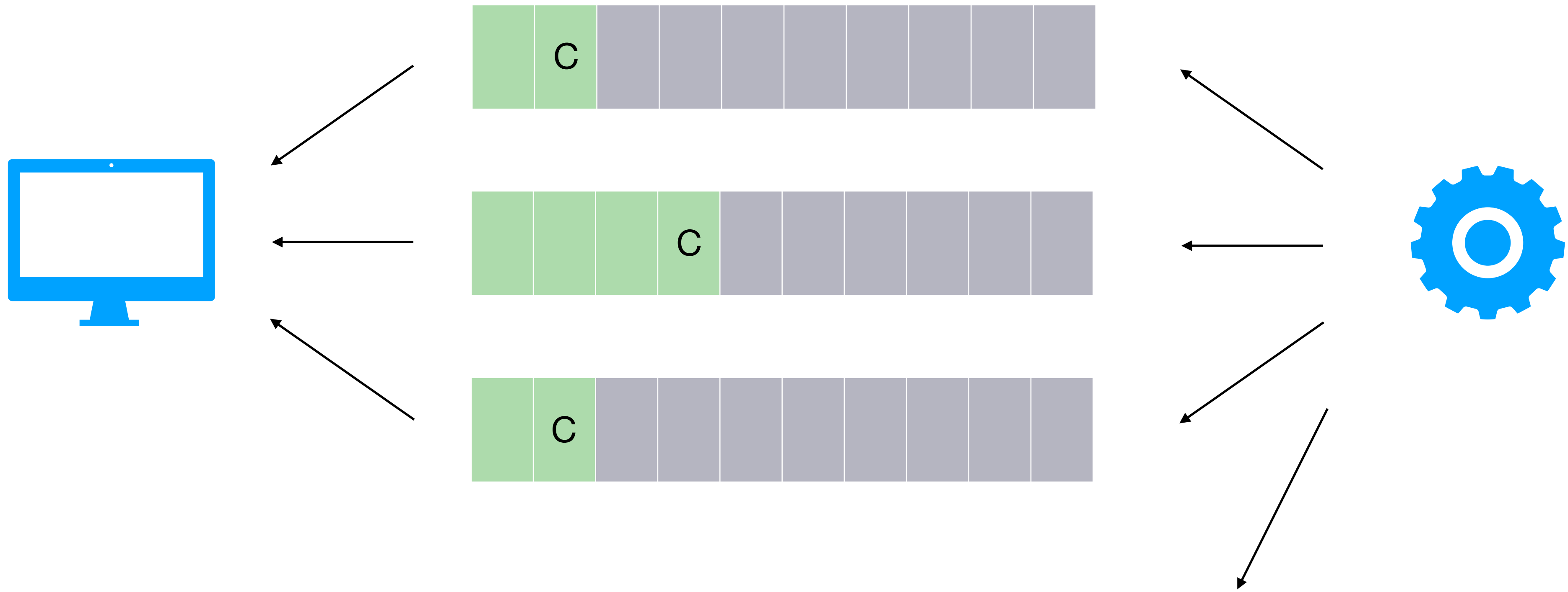
State

ongoing

prepare
commit

complete
commit

	(0, 1, 2)	(0, 1, 2)	(0, 1, 2)			
	ongoing	prepare commit	complete commit			



Transaction

Partitions

(0, 1, 2)

(0, 1, 2)

(0, 1, 2)

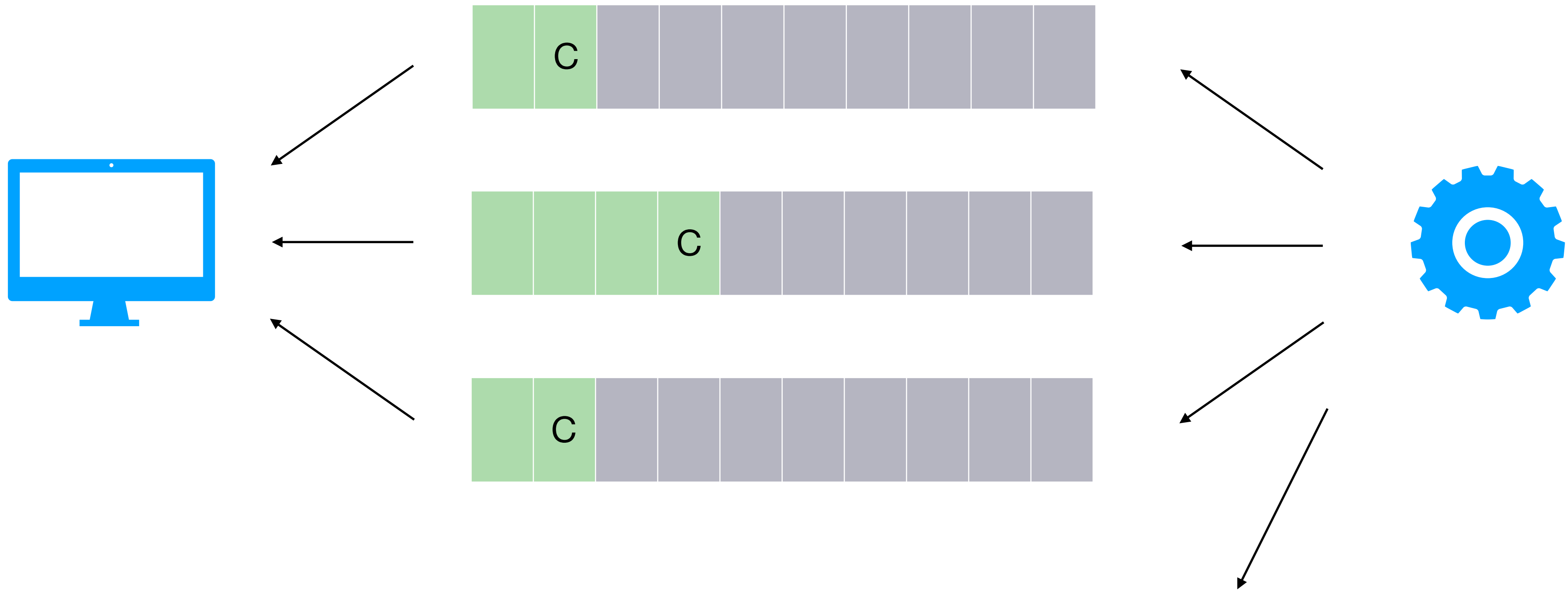
State

ongoing

prepare
commit

complete
commit

	(0, 1, 2)	(0, 1, 2)	(0, 1, 2)			
	ongoing	prepare commit	complete commit			



Transaction

Partitions

(0, 1, 2)

(0, 1, 2)

(0, 1, 2)

(0, 1)

State

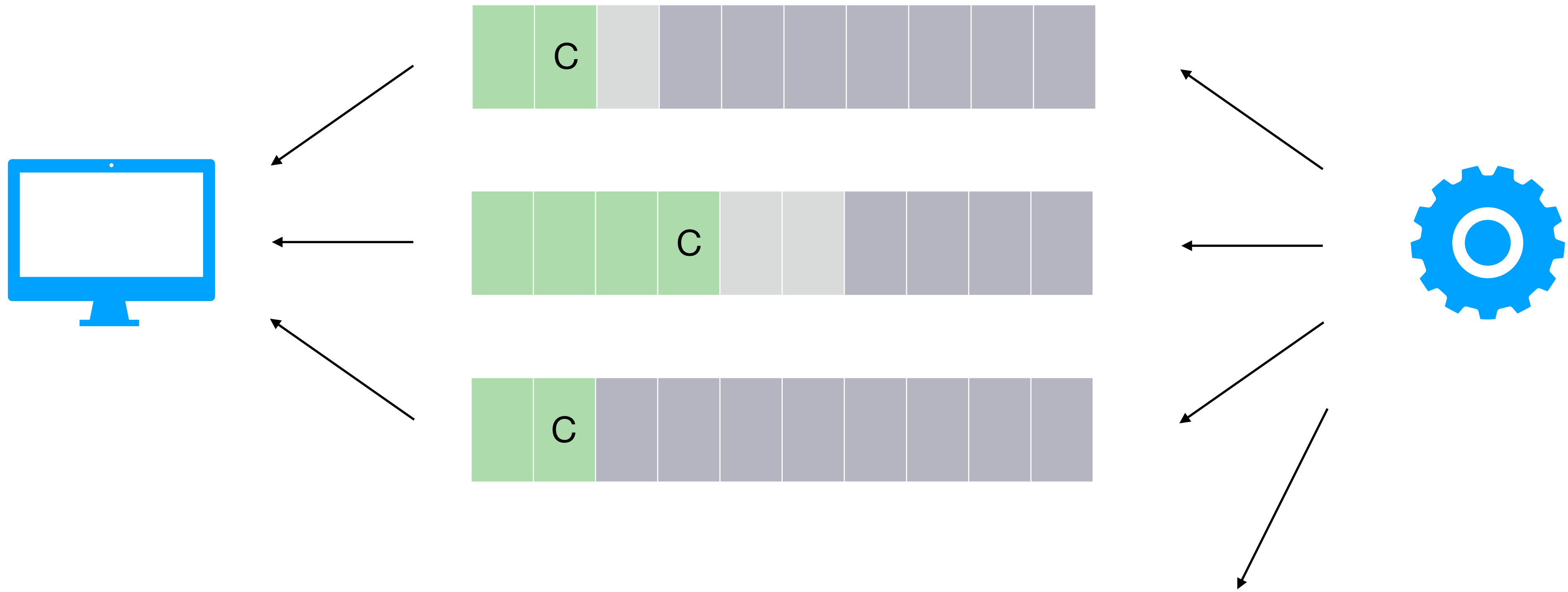
ongoing

prepare
commit

complete
commit

ongoing

	(0, 1, 2)	(0, 1, 2)	(0, 1, 2)	(0, 1)		
	ongoing	prepare commit	complete commit	ongoing		



Transaction

Partitions

(0, 1, 2)

(0, 1, 2)

(0, 1, 2)

(0, 1)

State

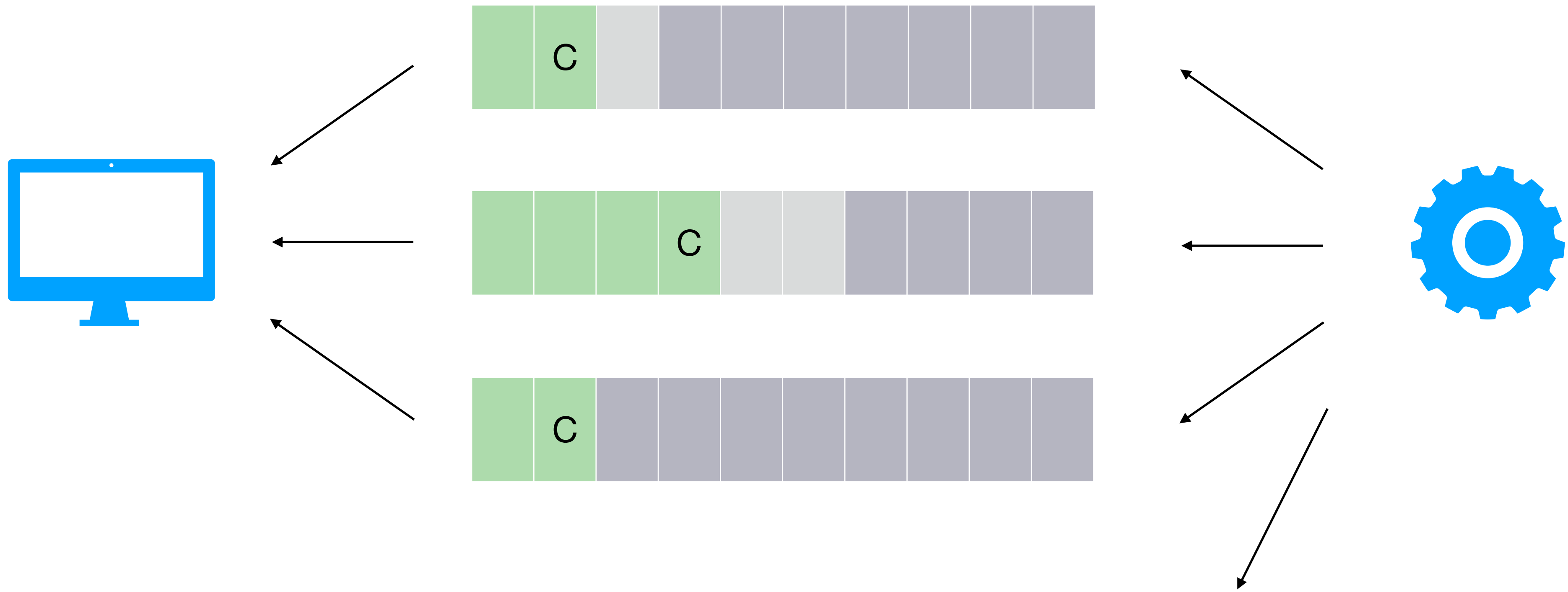
ongoing

prepare
commit

complete
commit

ongoing

	(0, 1, 2)	(0, 1, 2)	(0, 1, 2)	(0, 1)		
	ongoing	prepare commit	complete commit	ongoing		



Transaction

Partitions

(0, 1, 2)

(0, 1, 2)

(0, 1, 2)

(0, 1)

(0, 1)

State

ongoing

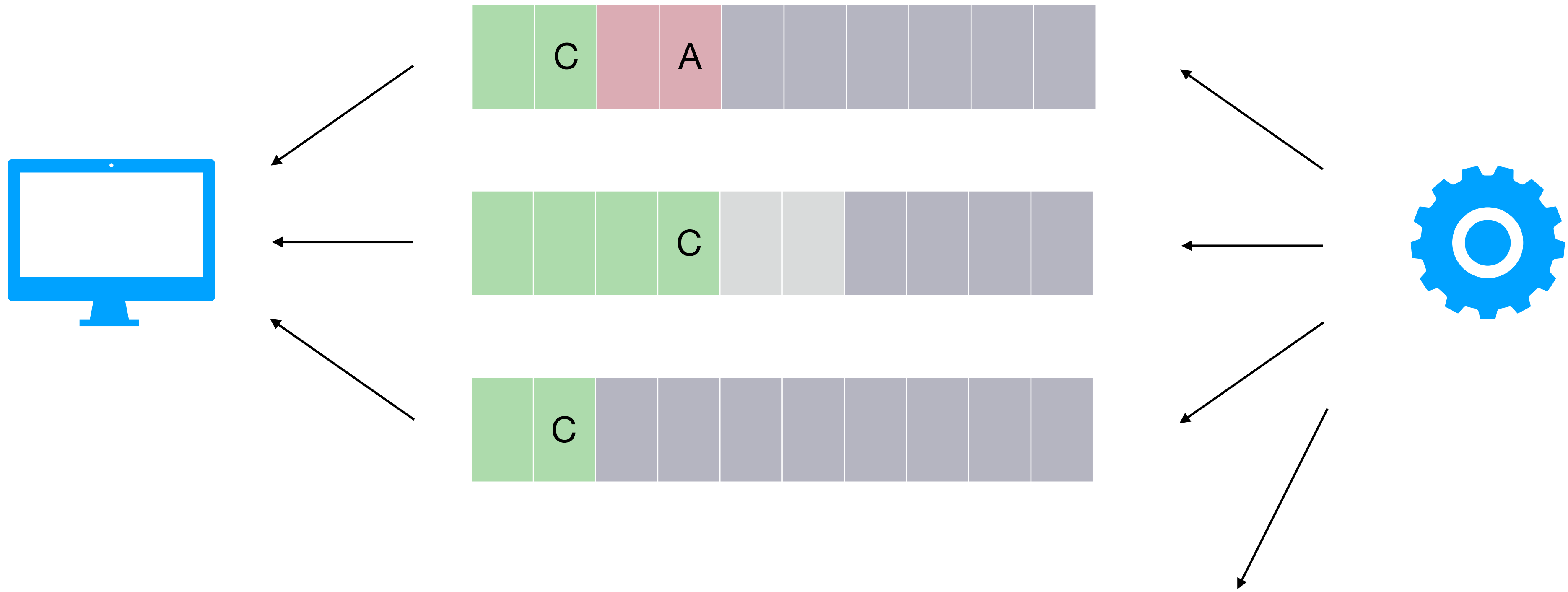
prepare
commit

complete
commit

ongoing

prepare
abort

	(0, 1, 2)	(0, 1, 2)	(0, 1, 2)	(0, 1)	(0, 1)	
	ongoing	prepare commit	complete commit	ongoing	prepare abort	



Transaction

Partitions

(0, 1, 2)

(0, 1, 2)

(0, 1, 2)

(0, 1)

(0, 1)

State

ongoing

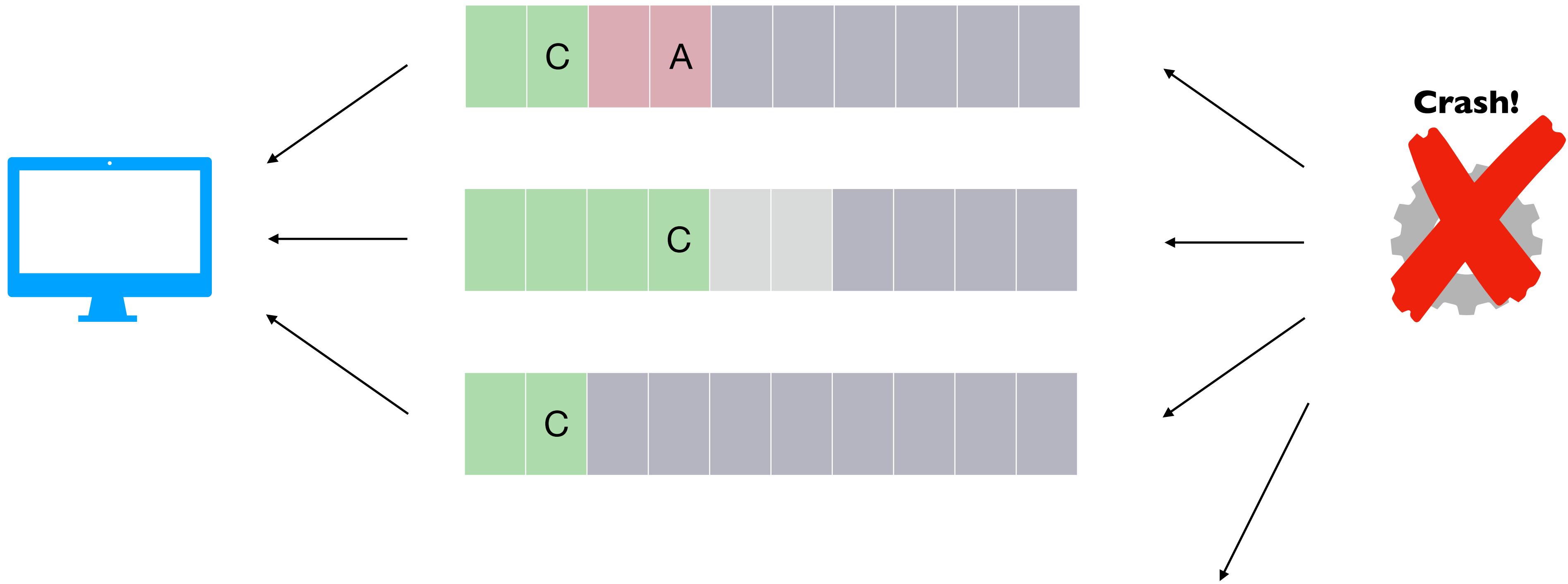
prepare
commit

complete
commit

ongoing

prepare
abort

	(0, 1, 2)	(0, 1, 2)	(0, 1, 2)	(0, 1)	(0, 1)	
	ongoing	prepare commit	complete commit	ongoing	prepare abort	

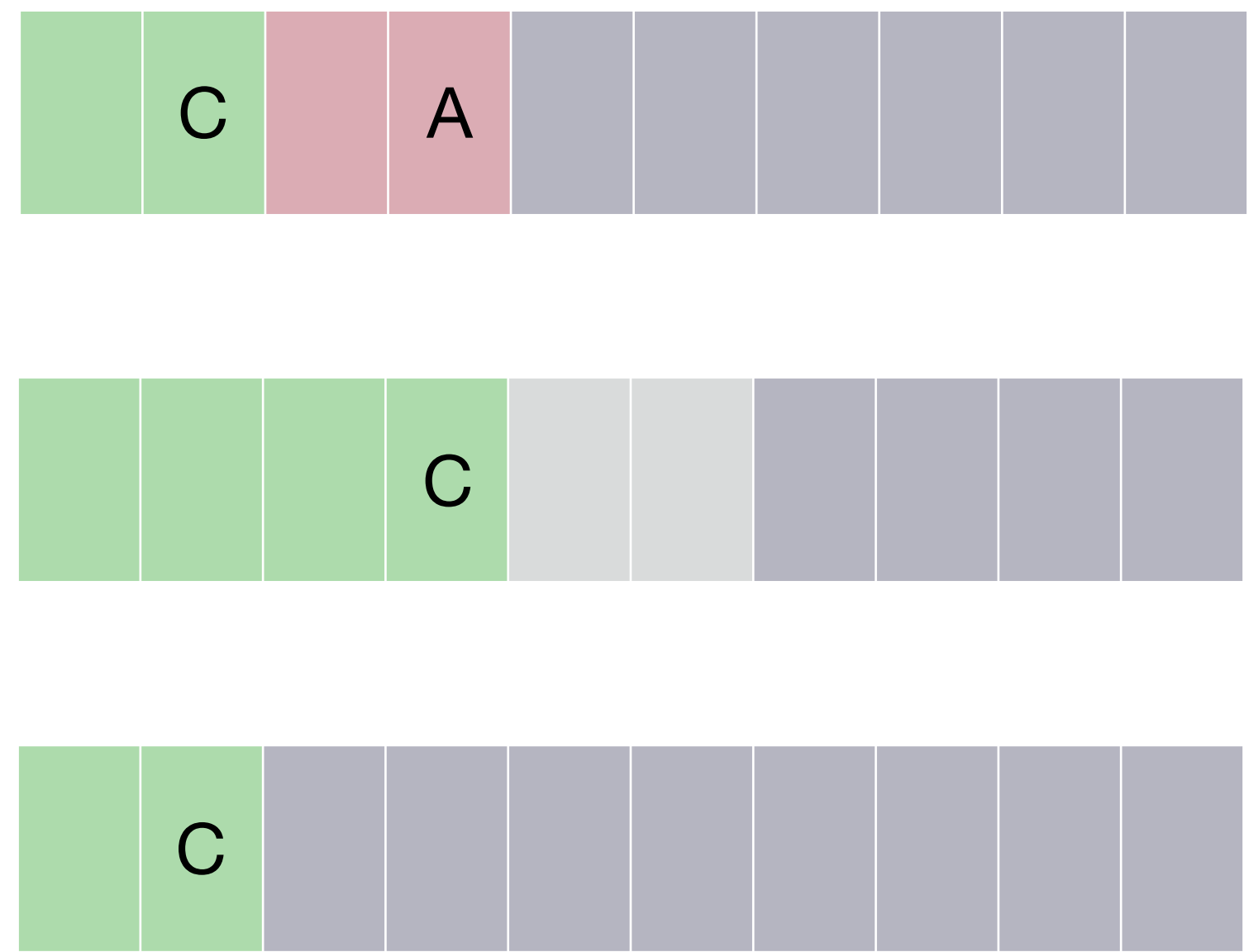
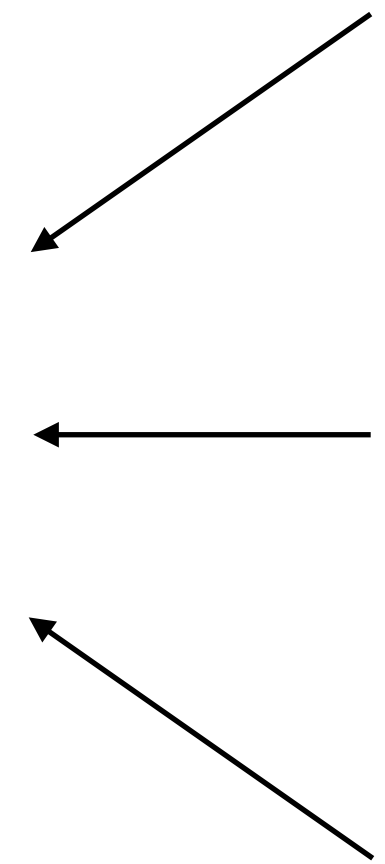


Transaction

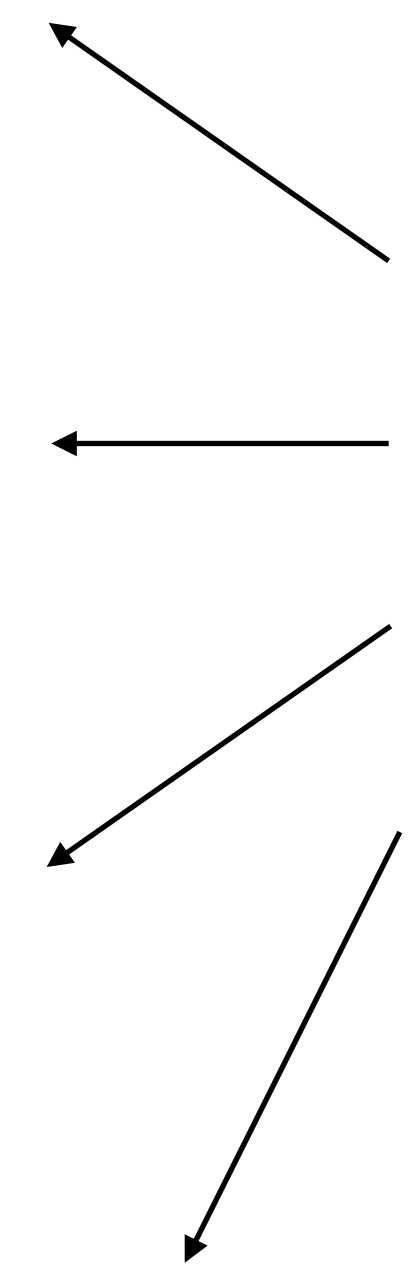
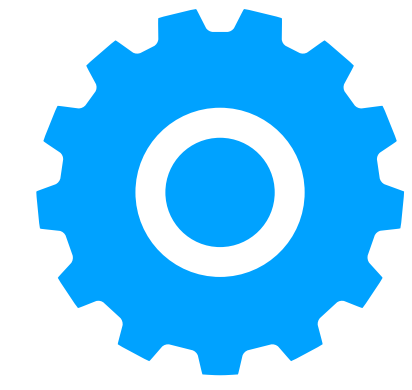
Partitions

State

$(0, 1, 2)$	$(0, 1, 2)$	$(0, 1, 2)$	$(0, 1)$	$(0, 1)$	
ongoing	prepare commit	complete commit	ongoing	prepare abort	



Restart!

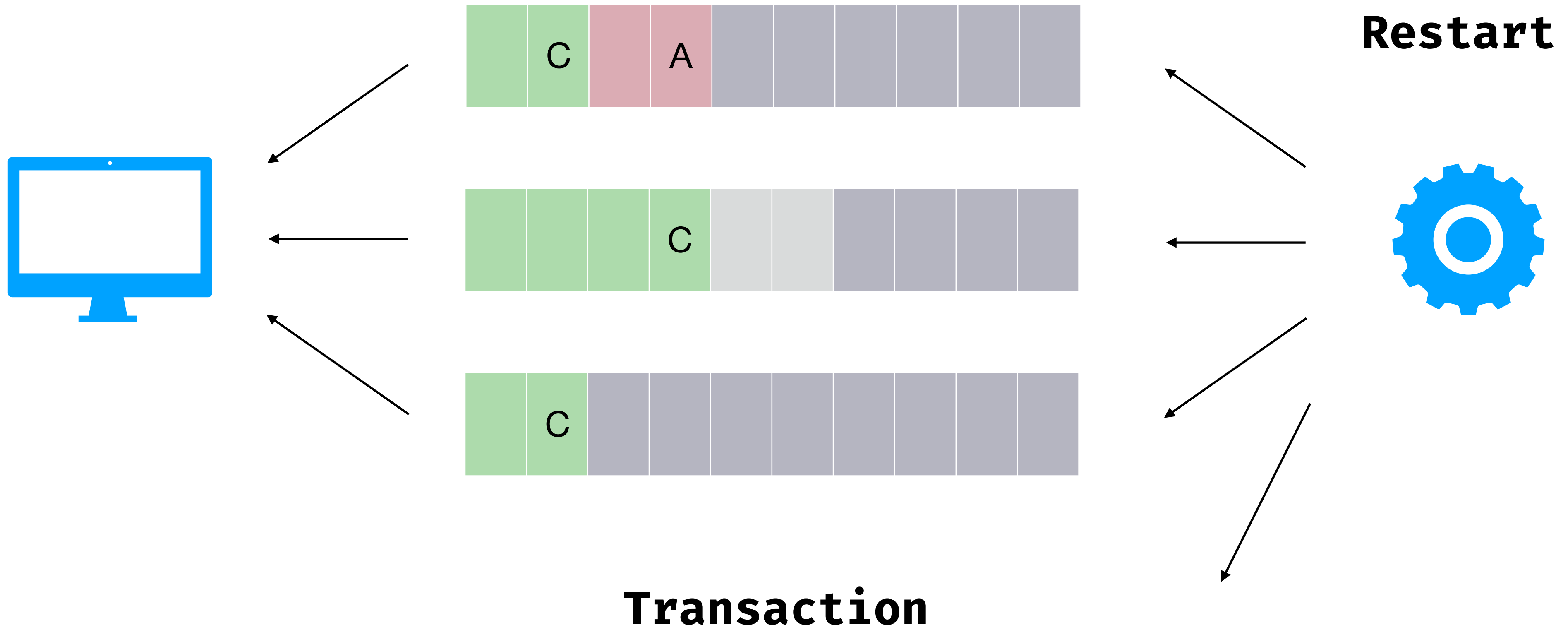


Transaction

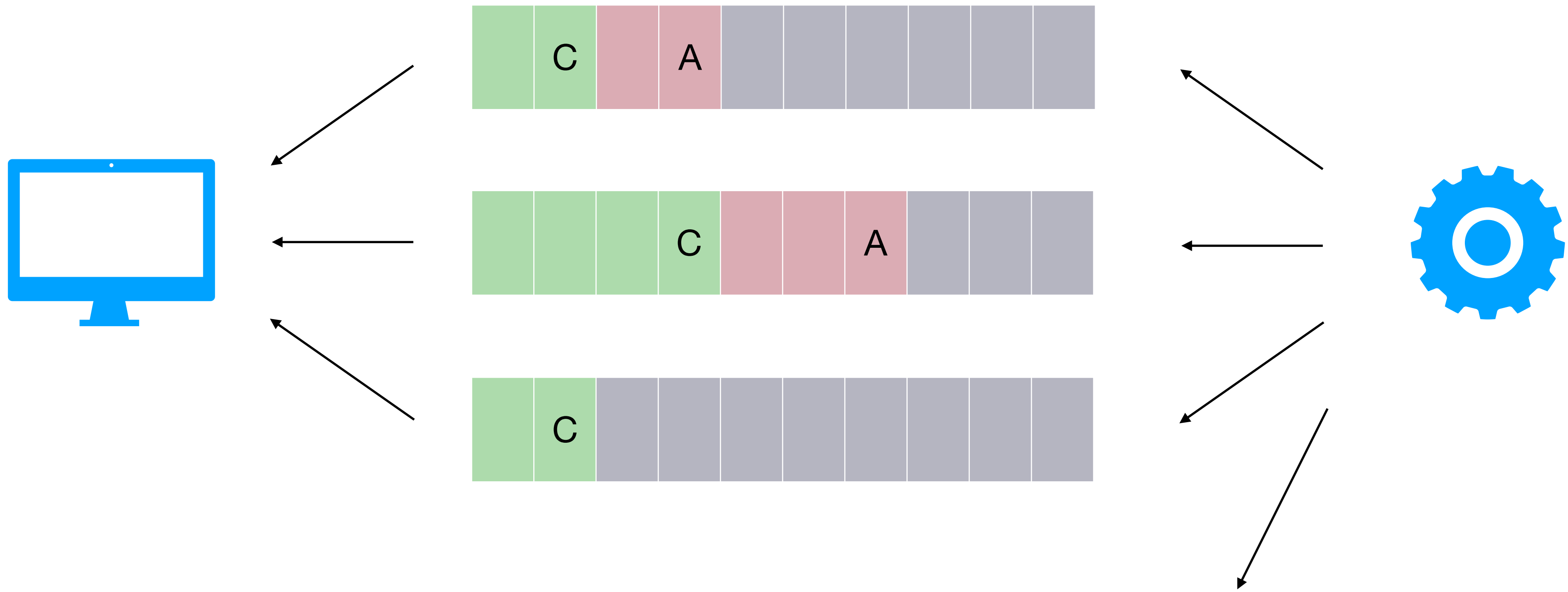
Partitions

State

$(0, 1, 2)$	$(0, 1, 2)$	$(0, 1, 2)$	$(0, 1)$	$(0, 1)$	
ongoing	prepare commit	complete commit	ongoing	prepare abort	



Partitions	(0, 1, 2)	(0, 1, 2)	(0, 1, 2)	(0, 1)	(0, 1)	
State	ongoing	prepare commit	complete commit	ongoing	prepare abort	



Transaction

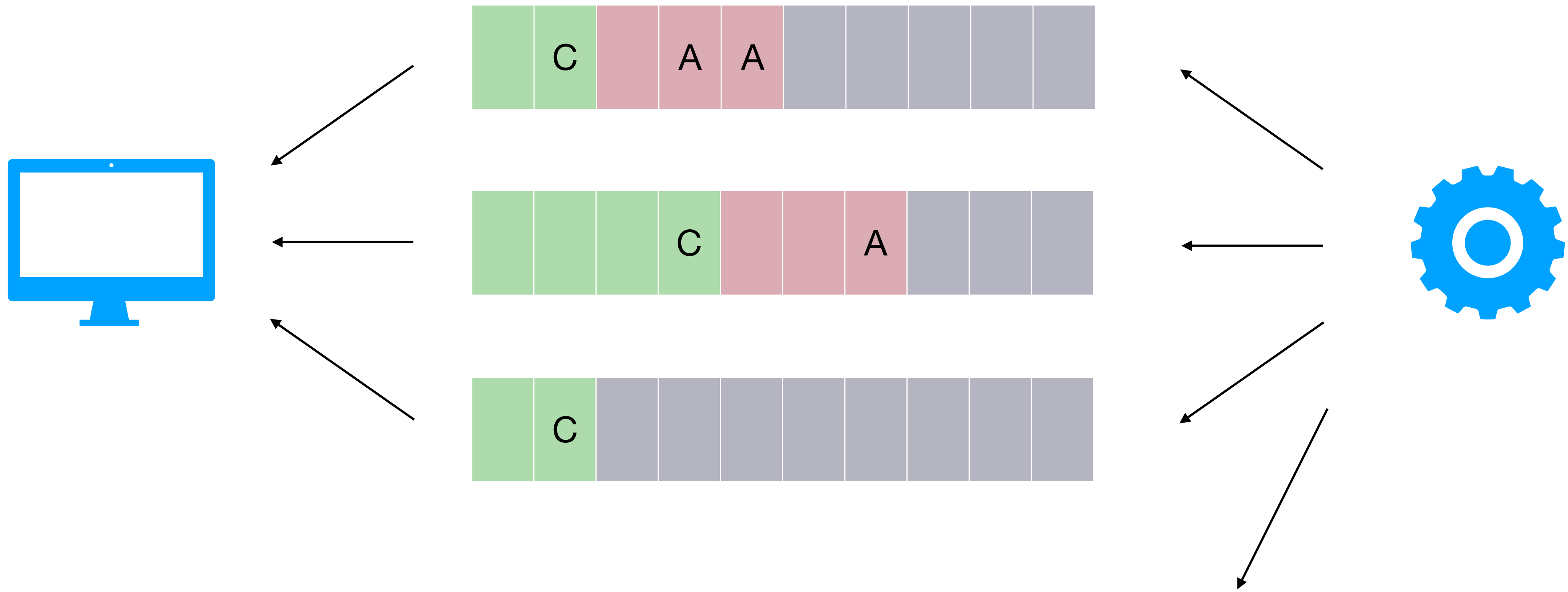
Partitions

(0, 1, 2) (0, 1, 2) (0, 1, 2) (0, 1) (0, 1)

State

ongoing prepare
commit complete
commit ongoing prepare
abort

(0, 1, 2)	(0, 1, 2)	(0, 1, 2)	(0, 1)	(0, 1)	
ongoing	prepare commit	complete commit	ongoing	prepare abort	



Transaction

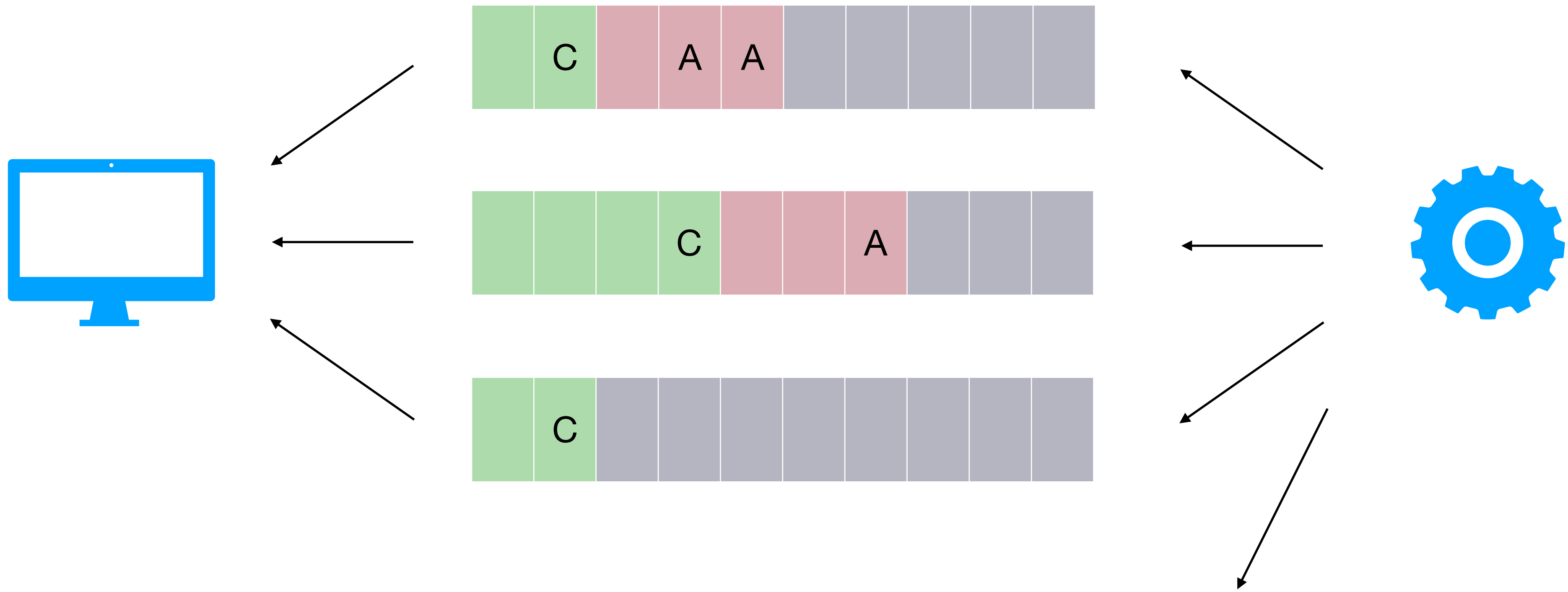
Partitions

(0, 1, 2) (0, 1, 2) (0, 1, 2) (0, 1) (0, 1)

State

ongoing prepare commit complete commit ongoing prepare abort

(0, 1, 2)	(0, 1, 2)	(0, 1, 2)	(0, 1)	(0, 1)	
ongoing	prepare commit	complete commit	ongoing	prepare abort	



Transaction

Partitions

(0, 1, 2) (0, 1, 2) (0, 1, 2) (0, 1) (0, 1) (0, 1)

State

	(0, 1, 2)	(0, 1, 2)	(0, 1, 2)	(0, 1)	(0, 1)	(0, 1)
	ongoing	prepare commit	complete commit	ongoing	prepare abort	complete abort

Транзакции

- Атомарная запись на множество партиций
- Включая `__consumer_offsets`

Слабости семантики

- ~~Повторные попытки прогюсера небезопасны~~
- Данные не пишутся атомарно с их смещениями (offsets)
- Нет защиты от зомби

Слабости семантики

- ~~Повторные попытки прогюсера небезопасны~~
- ~~Данные не пишутся атомарно с их смещениями (offsets)~~
- Нет защиты от зомби

Слабости семантики

- ~~Повторные попытки прогюсера небезопасны~~
- ~~Данные не пишутся атомарно с их смещениями (offsets)~~
- Нет защиты от зомби



Забор от Зомби



- Чтобы остановить зомби нам нужен «забор»



Забор от Зомби

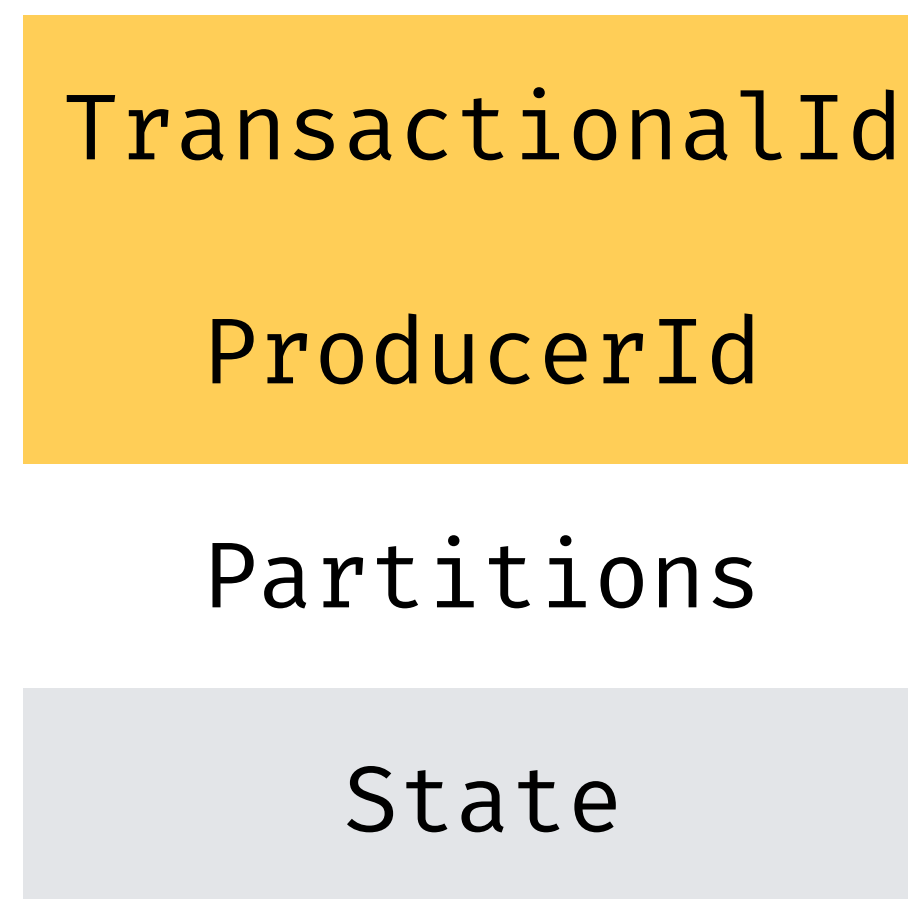


- Чтобы остановить зомби нам нужен «забор»
- Но для начала, надо каким-то образом идентифицировать экземпляр продьюсера между перезапусками

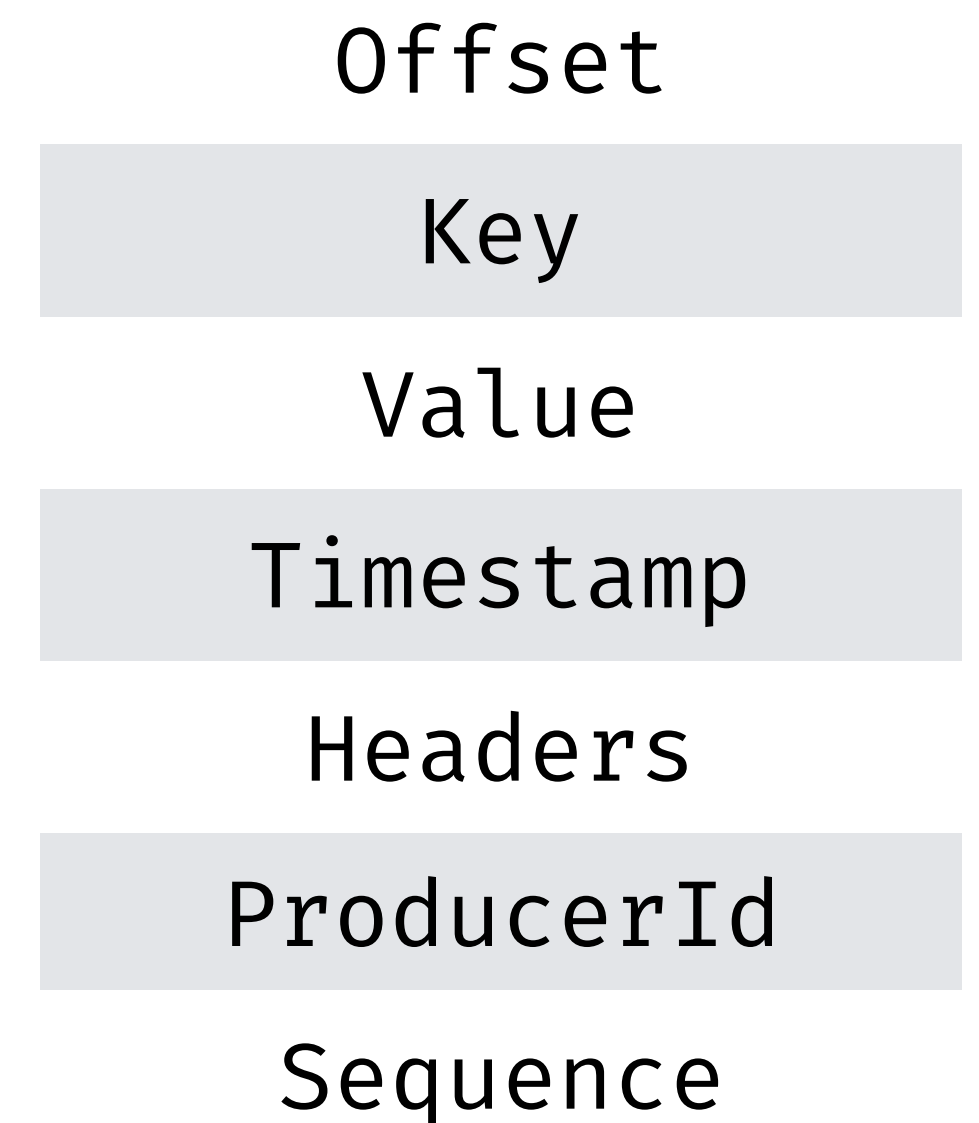
Zombie fencing

- “transactionalId” связан с producerId
- Не надо сохранять producerId
- `transactional.id={id}`

Transaction Log Message

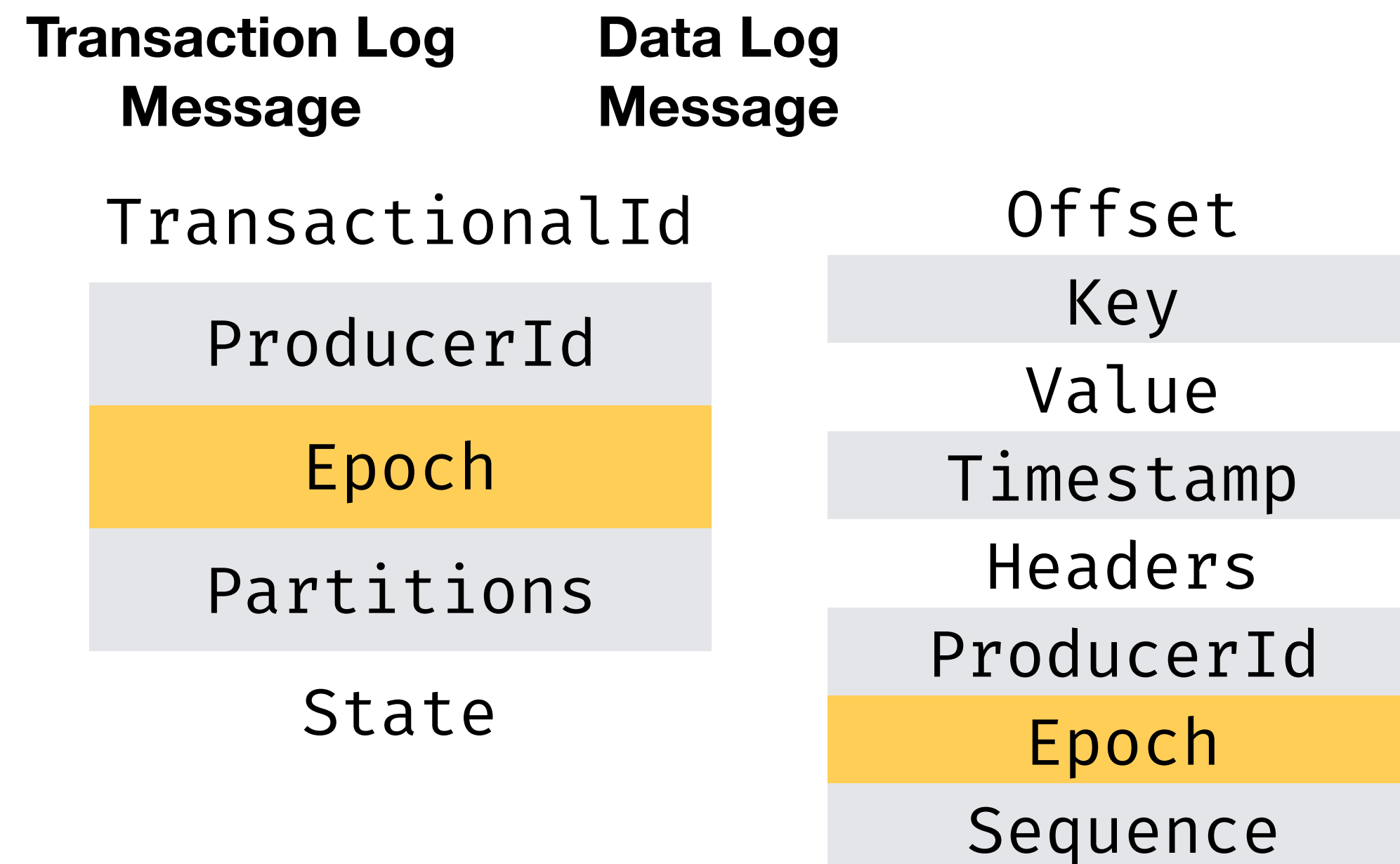


Data Log Message



Zombie fencing

- Мы добавили эпоху продюсера epoch
- Epoch инкрементируется при инициализации продюсера



Zombie fencing

Producer 1

txn.id=A

coordinator=1

Transaction Log 0 (Leader)

txn.id				
pid				
epoch				
partitions				
state				

Zombie fencing

Producer 1

txn1.id=A

coordinator=1

InitProducerId(txn1.id=A)



Transaction Log 0 (Leader)

txn1.id				
pid				
epoch				
partitions				
state				

Zombie fencing

Producer 1

txn.id=A

coordinator=1

InitProducerId(txn.id=A)



Transaction Log 0 (Leader)

txn.id	A			
pid	239			
epoch	0			
partitions	none			
state	init			

Zombie fencing

Producer 1

txnl.id=A

coordinator=1

pid=239, epoch=0



Transaction Log 0 (Leader)

txnl.id	A			
pid	239			
epoch	0			
partitions	none			
state	init			

Zombie fencing

Producer 1

txnl.id=A
coordinator=1
pid=239
epoch=0

Transaction Log 0 (Leader)

txnl.id	A			
pid	239			
epoch	0			
partitions	none			
state	init			

Zombie fencing

Producer 1

txnl.id=A
coordinator=1
pid=239
epoch=0

AddPartitions([0])



Transaction Log 0 (Leader)

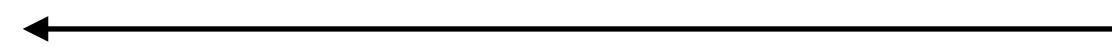
txnl.id	A	A		
pid	239	239		
epoch	0	0		
partitions	none	[0]		
state	init	ongoing		

Zombie fencing

Producer 1

txnl.id=A
coordinator=1
pid=239
epoch=0

OK



Transaction Log 0 (Leader)

txnl.id	A	A		
pid	239	239		
epoch	0	0		
partitions	none	[0]		
state	init	ongoing		

Zombie fencing

Producer 1

txnl.id=A
coordinator=1
pid=239
epoch=0

Produce(pid=239, epoch=0)



Transaction Log 0 (Leader)

txnl.id	A	A		
pid	239	239		
epoch	0	0		
partitions	none	[0]		
state	init	ongoing		

Data partition 0 (Leader)

offset	0							
pid	239							
epoch	0							

zombie fencing

Producer 1

txnl.id=A
coordinator=1
pid=239
epoch=0

Producer 2

txnl.id=A
coordinator=1

Transaction Log 0 (Leader)

txnl.id	A	A		
pid	239	239		
epoch	0	0		
partitions	none	[0]		
state	init	ongoing		

Data partition 0 (Leader)

offset	0							
pid	239							
epoch	0							

zombie fencing

Producer 1

txnl.id=A
coordinator=1
pid=239
epoch=0

Producer 2

txnl.id=A
coordinator=1

InitProducerId(txnl.id=A)

Transaction Log 0 (Leader)

txnl.id	A	A		
pid	239	239		
epoch	0	0		
partitions	none	[0]		
state	init	ongoing		

Data partition 0 (Leader)

offset	0							
pid	239							
epoch	0							

zombie fencing

Producer 1

txnl.id=A
coordinator=1
pid=239
epoch=0

Producer 2

txnl.id=A
coordinator=1

InitProducerId(txnl.id=A)

Transaction Log 0 (Leader)

txnl.id	A	A	A	
pid	239	239	239	
epoch	0	0	1	
partitions	none	[0]	[0]	
state	init	ongoing	pr_abort	

Data partition 0 (Leader)

offset	0							
pid	239							
epoch	0							

Zombie fencing

Producer 1

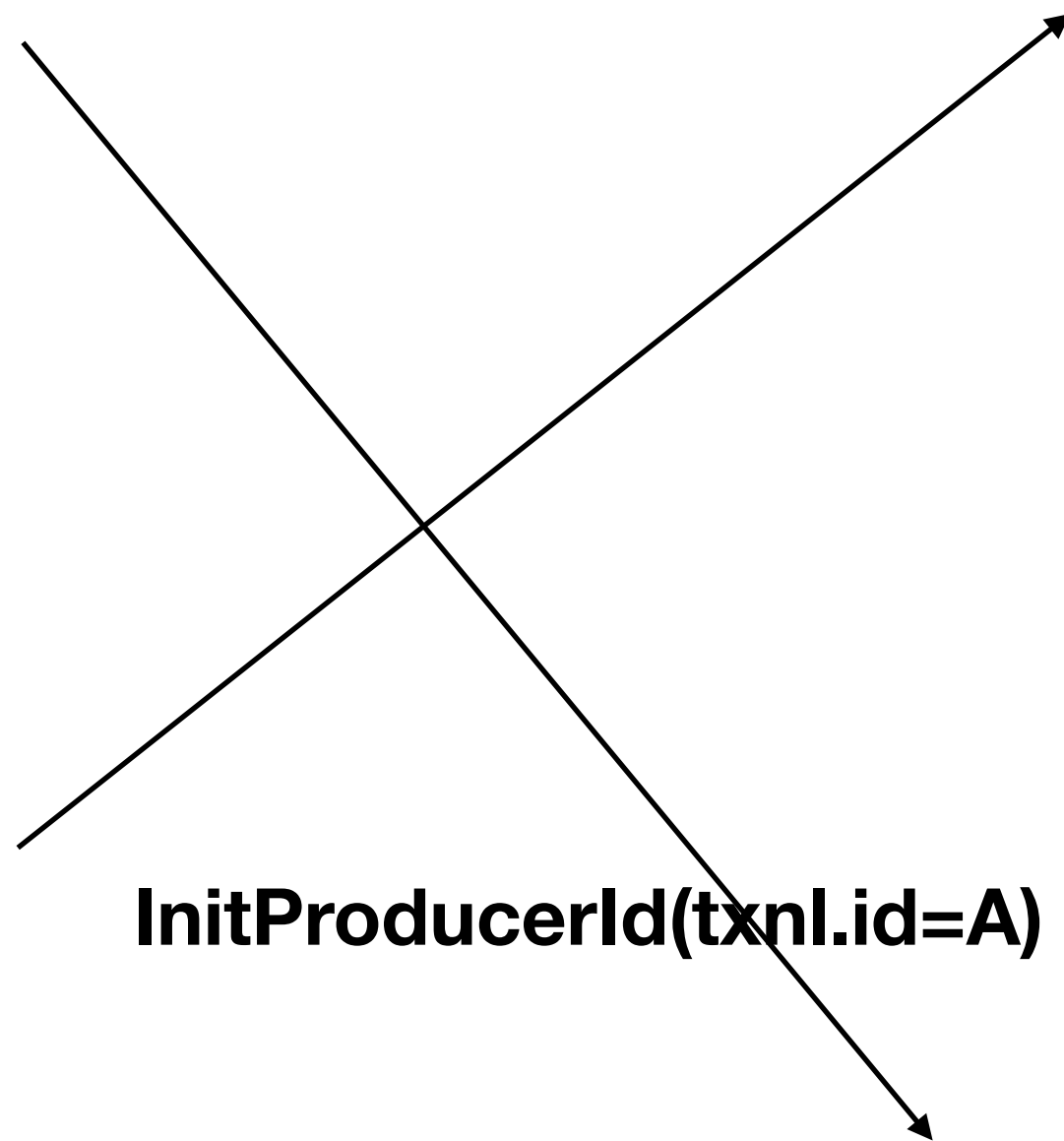
txnl.id=A
 coordinator=1
 pid=239
 epoch=0

Produce(pid=239, epoch=0)

Producer 2

txnl.id=A
 coordinator=1

InitProducerId(txnl.id=A)



Transaction Log 0 (Leader)

txnl.id	A	A	A	
pid	239	239	239	
epoch	0	0	1	
partitions	none	[0]	[0]	
state	init	ongoing	pr_abort	

Data partition 0 (Leader)

offset	0	1						
pid	239	239						
epoch	0	0						

zombie fencing

Producer 1

txnl.id=A
coordinator=1
pid=239
epoch=0

Producer 2

txnl.id=A
coordinator=1

InitProducerId(txnl.id=A)

Transaction Log 0 (Leader)

txnl.id	A	A	A	
pid	239	239	239	
epoch	0	0	1	
partitions	none	[0]	[0]	
state	init	ongoing	pr_abort	

Data partition 0 (Leader)

offset	0	1	2					
pid	239	239	239					
epoch	0	0	1					

Zombie fencing

Producer 1

txnl.id=A
coordinator=1
pid=239
epoch=0

Produce(pid=239, epoch=0)

Producer 2

txnl.id=A
coordinator=1

InitProducerId(txnl.id=A)

Transaction Log 0 (Leader)

txnl.id	A	A	A	
pid	239	239	239	
epoch	0	0	1	
partitions	none	[0]	[0]	
state	init	ongoing	pr_abort	

Data partition 0 (Leader)

offset	0	1	2					
pid	239	239	239					
epoch	0	0	1					

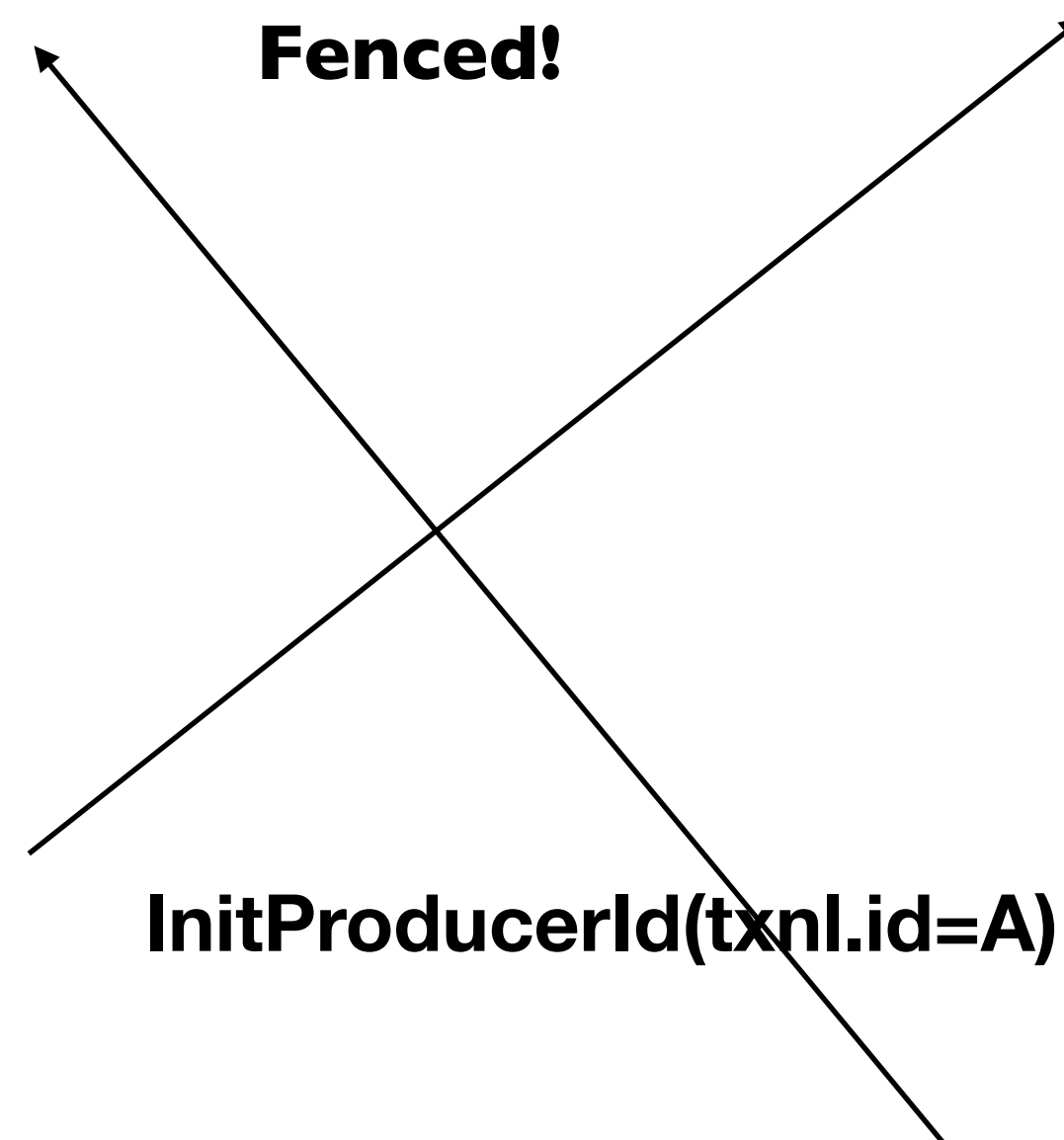
Zombie fencing

Producer 1

txnl.id=A
coordinator=1
pid=239
epoch=0

Producer 2

txnl.id=A
coordinator=1



Transaction Log 0 (Leader)

txnl.id	A	A	A	
pid	239	239	239	
epoch	0	0	1	
partitions	none	[0]	[0]	
state	init	ongoing	pr_abort	

Data partition 0 (Leader)

offset	0	1	2					
pid	239	239	239					
epoch	0	0	1					

Zombie fencing

Producer 1

txnl.id=A
coordinator=1
pid=239
epoch=0

Producer 2

txnl.id=A
coordinator=1

InitProducerId(txnl.id=A)

Transaction Log 0 (Leader)

txnl.id	A	A	A	A
pid	239	239	239	239
epoch	0	0	1	1
partitions	none	[0]	[0]	[0]
state	init	ongoing	p_abort	c_abort

Data partition 0 (Leader)

offset	0	1	2					
pid	239	239	239					
epoch	0	0	1					

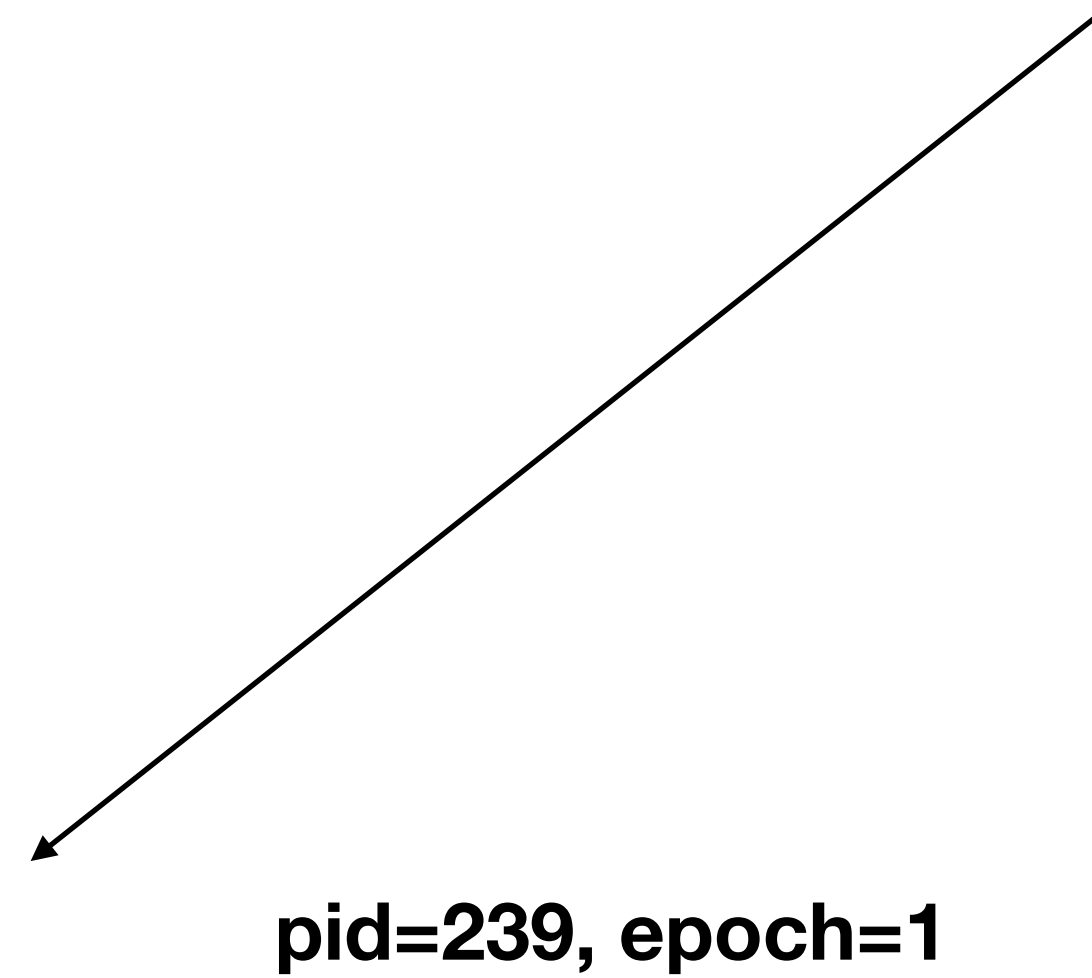
zombie fencing

Producer 1

txnl.id=A
coordinator=1
pid=239
epoch=0

Producer 2

txnl.id=A
coordinator=1



Transaction Log 0 (Leader)

txnl.id	A	A	A	A
pid	239	239	239	239
epoch	0	0	1	1
partitions	none	[0]	[0]	[0]
state	init	ongoing	p_abort	c_abort

Data partition 0 (Leader)

offset	0	1	2					
pid	239	239	239					
epoch	0	0	1					

Zombie fencing

Producer 1

txnl.id=A
coordinator=1
pid=239
epoch=0

Producer 2

txnl.id=A
coordinator=1
pid=239
epoch=1

Transaction Log 0 (Leader)

txnl.id	A	A	A	A
pid	239	239	239	239
epoch	0	0	1	1
partitions	none	[0]	[0]	[0]
state	init	ongoing	p_abort	c_abort

Data partition 0 (Leader)

offset	0	1	2					
pid	239	239	239					
epoch	0	0	1					

zombie fencing

Producer 1

txnl.id=A
coordinator=1
pid=239
epoch=0

CommitTxn(epoch=0)



Transaction Log 0 (Leader)

txnl.id	A	A	A	A
pid	239	239	239	239
epoch	0	0	1	1
partitions	none	[0]	[0]	[0]
state	init	ongoing	p_abort	c_abort

Producer 2

txnl.id=A
coordinator=1
pid=239
epoch=1

Data partition 0 (Leader)

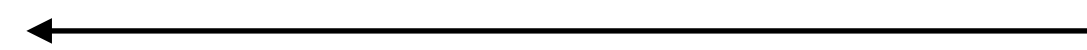
offset	0	1	2					
pid	239	239	239					
epoch	0	0	1					

Zombie fencing

Producer 1

txnl.id=A
coordinator=1
pid=239
epoch=0

Fenced!



Producer 2

txnl.id=A
coordinator=1
pid=239
epoch=1

Transaction Log 0 (Leader)

txnl.id	A	A	A	A
pid	239	239	239	239
epoch	0	0	1	1
partitions	none	[0]	[0]	[0]
state	init	ongoing	p_abort	c_abort

Data partition 0 (Leader)

offset	0	1	2					
pid	239	239	239					
epoch	0	0	1					

Слабости семантики

- ~~Повторные попытки прогюсера небезопасны~~
- ~~Данные не пишутся атомарно с их смещениями (offsets)~~
- Нет защиты от зомби

Слабости семантики

- ~~Повторные попытки прогюсера небезопасны~~
- ~~Данные не пишутся атомарно с их смещениями (offsets)~~
- ~~Нет защиты от зомби~~

Семантика «exactly-once»

```
p.initTxns()  
c.subscribe(inputTopics)  
while (true) {  
    offsets, input = c.poll()  
    output = process(input)  
    p.beginTxn()  
    p.send(output)  
    p.sendOffsets(offsets)  
    p.commitTxn()  
}
```

Семантика «exactly-once»

```
p.initTxns()  
c.subscribe(inputTopics)  
while (true) {  
    offsets, input = c.poll()  
    output = process(input)  
  
    p.beginTxn()                void beginTransaction();  
    p.send(output)  
    p.sendOffsets(offsets)  
    p.commitTxn()  
}
```

Семантика «exactly-once»

```
void sendOffsetsToTransaction(  
    Map<TopicPartition, OffsetAndMetadata> offsets, String consumerGroupId)  
throws ProducerFencedException;
```

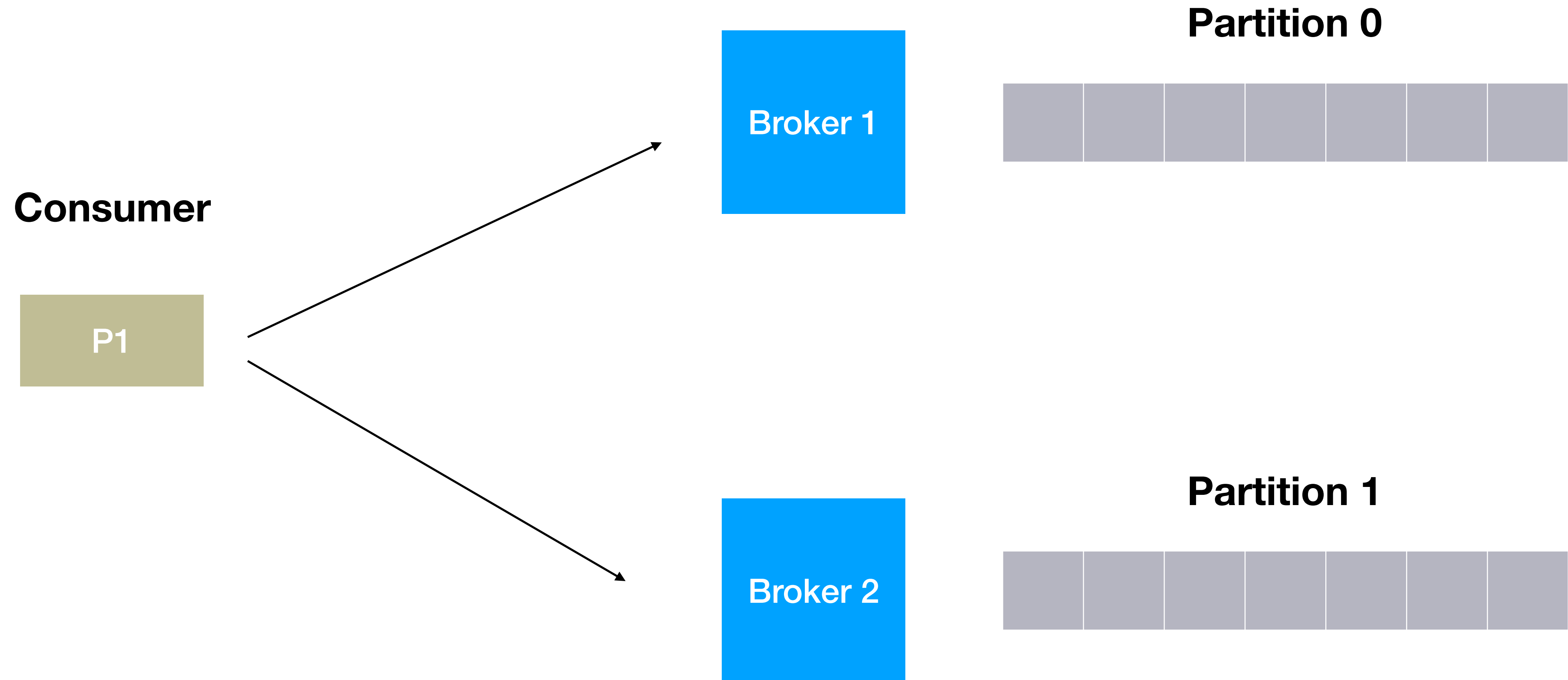
Семантика «exactly-once»

```
p.initTxns()  
c.subscribe(inputTopics)  
while (true) {  
    offsets, input = c.poll()  
    output = process(input)  
    p.beginTxn()                                void commitTransaction()  
    p.send(output)                               throws ProducerFencedException;  
    p.sendOffsets(offsets)  
    p.commitTxn()  
}
```

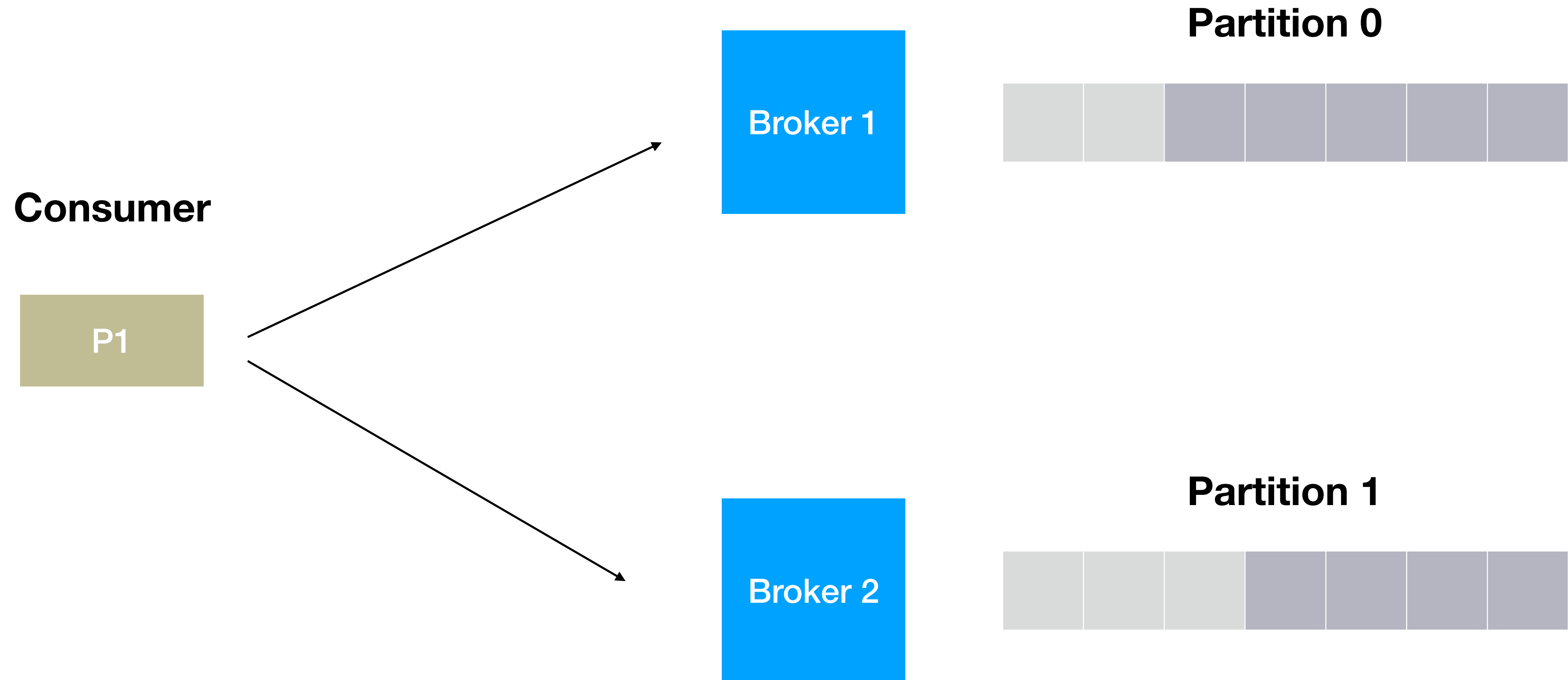
Изоляция Чтения у Consumer

- `Consumer isolation.level`:
 - `read_committed`: доступны только закоммиченные сообщения
 - `read_uncommitted`: Все сообщения доступны

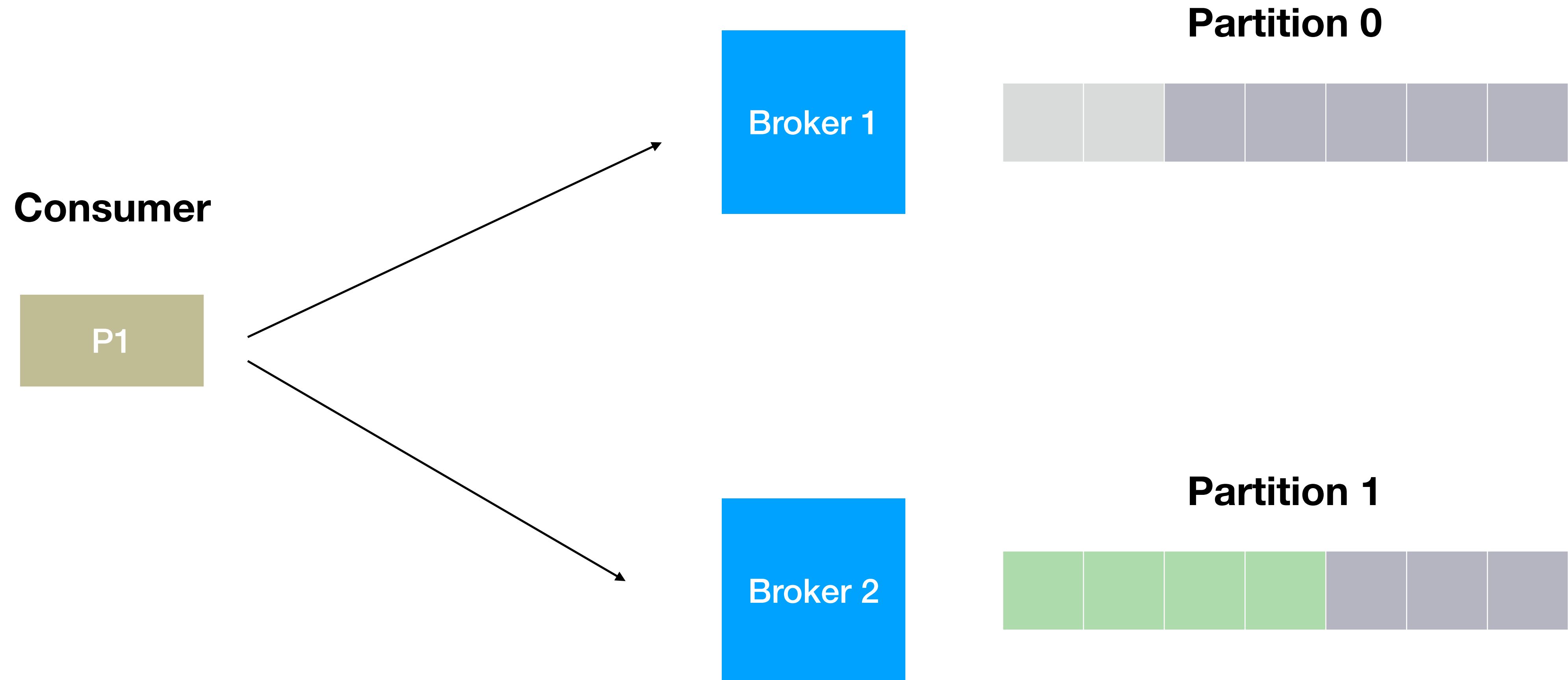
Изоляция чтения у Consumer



Изоляция чтения у Consumer

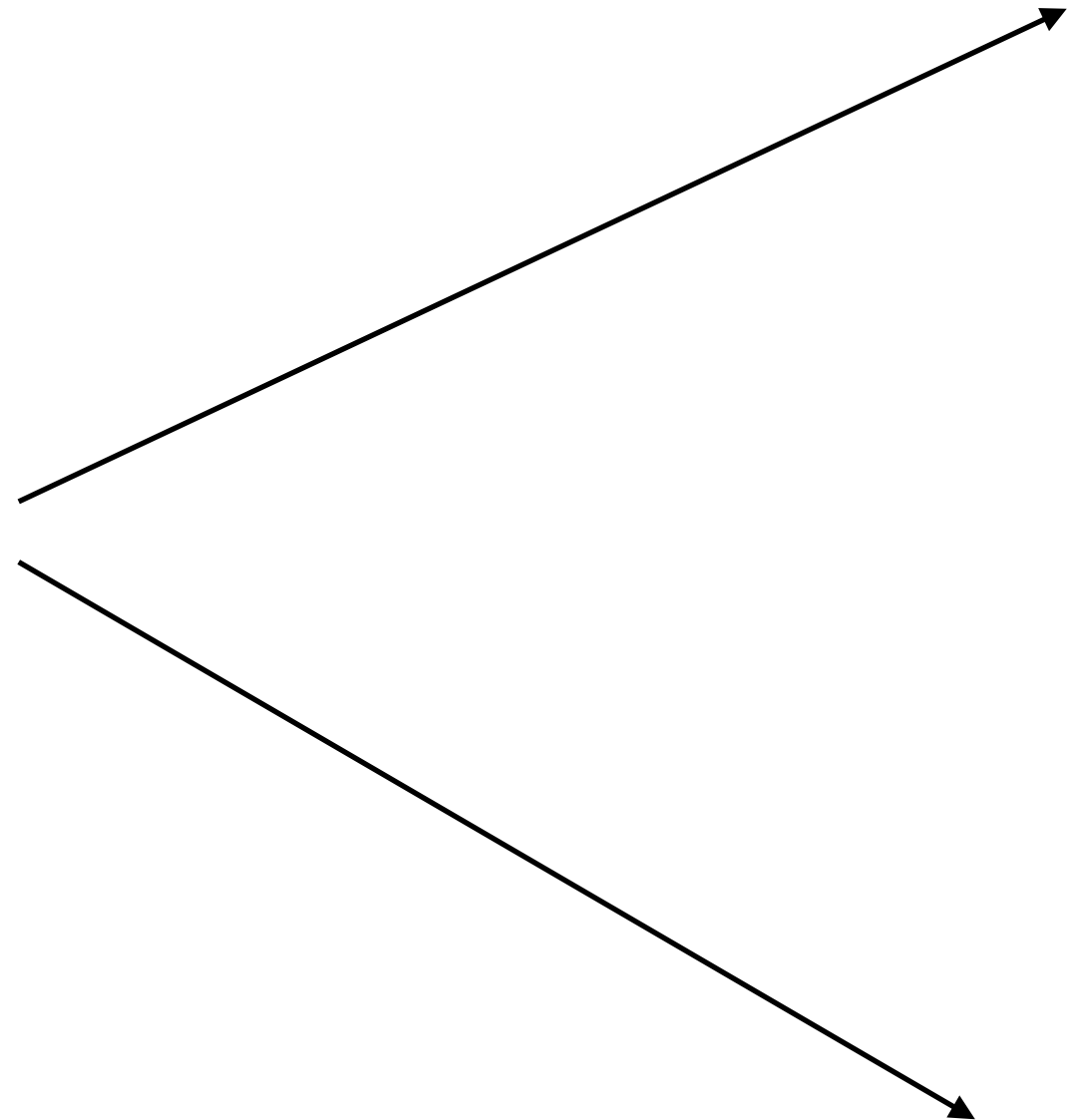


Изоляция чтения у Consumer



Изоляция чтения у Consumer

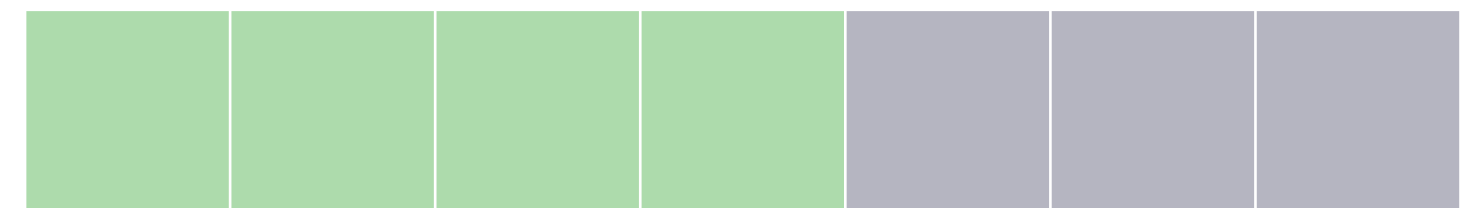
Consumer



Partition 0



Partition 1



Fetchable

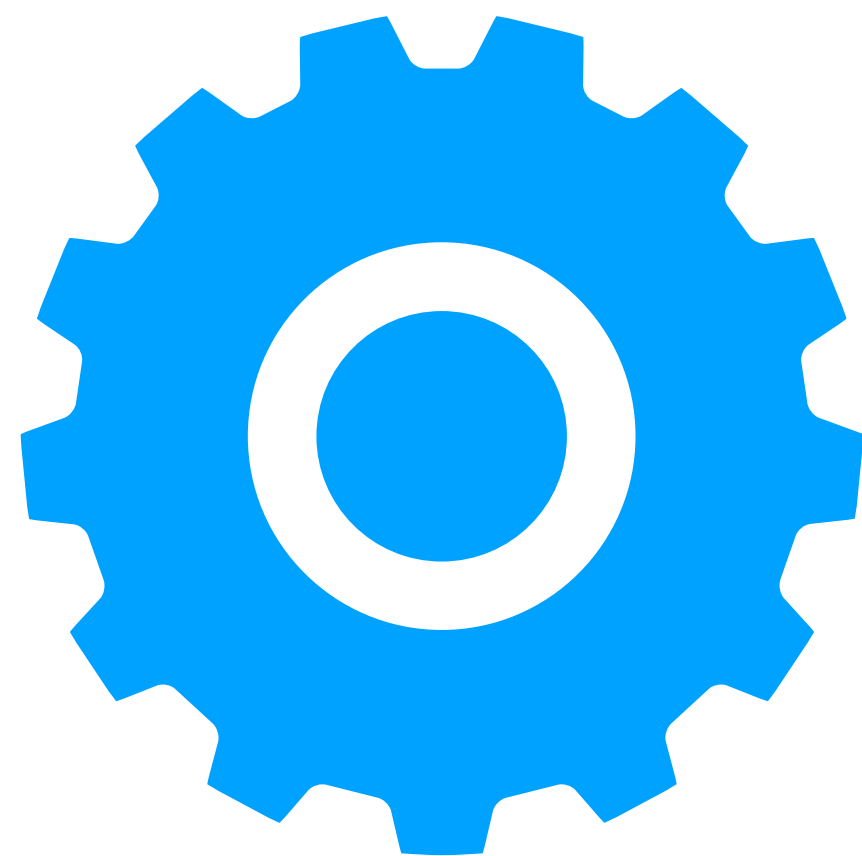
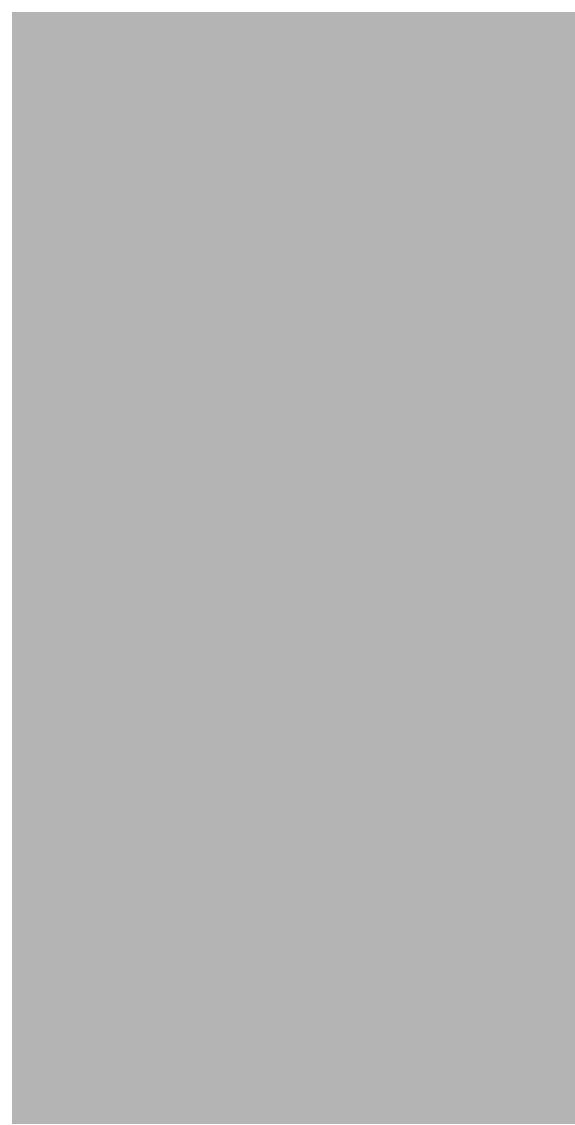
Семантика «Только один раз»

- Идемпотентный продюсер – доставка сообщений без дубликатов
- Транзакционный продюсер: возможность атомарной записи в несколько партиции (в том числе в хранилище offsets)
- Read-committed консьюмер: только законченные сообщения будут возвращены

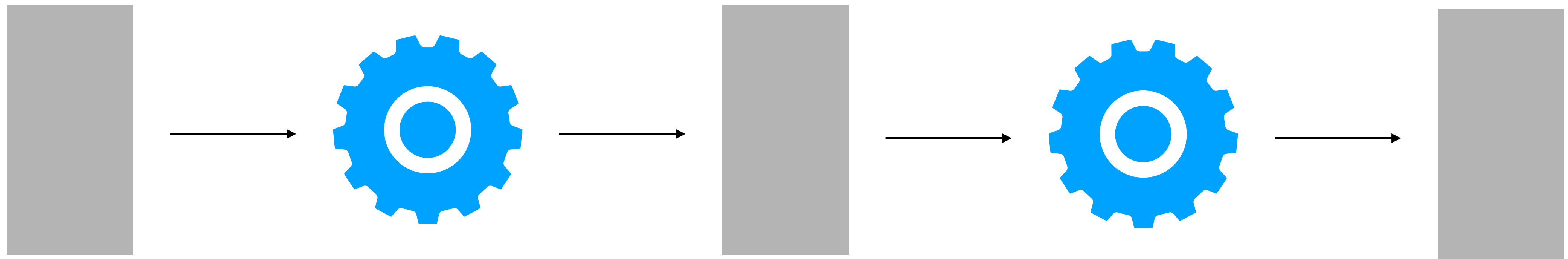
А тормозить не будет?

- Оверхед минимальный
- Чем продолжительней транзакция, тем выше end-to-end latency

end-to-end eos

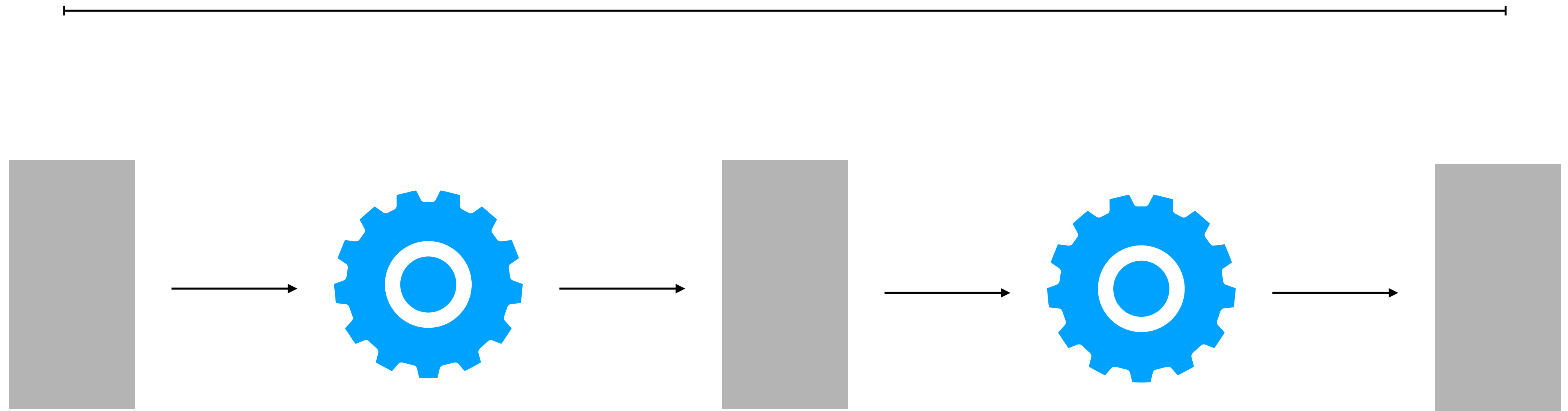


end-to-end eos



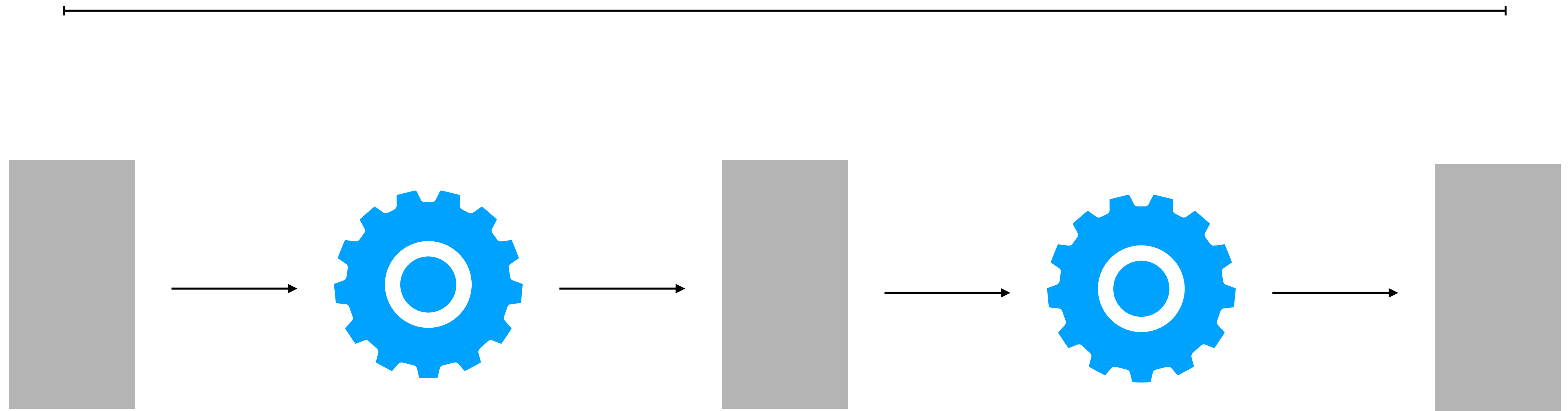
end-to-end eos

Kafka EoS



end-to-end eos

Kafka EoS



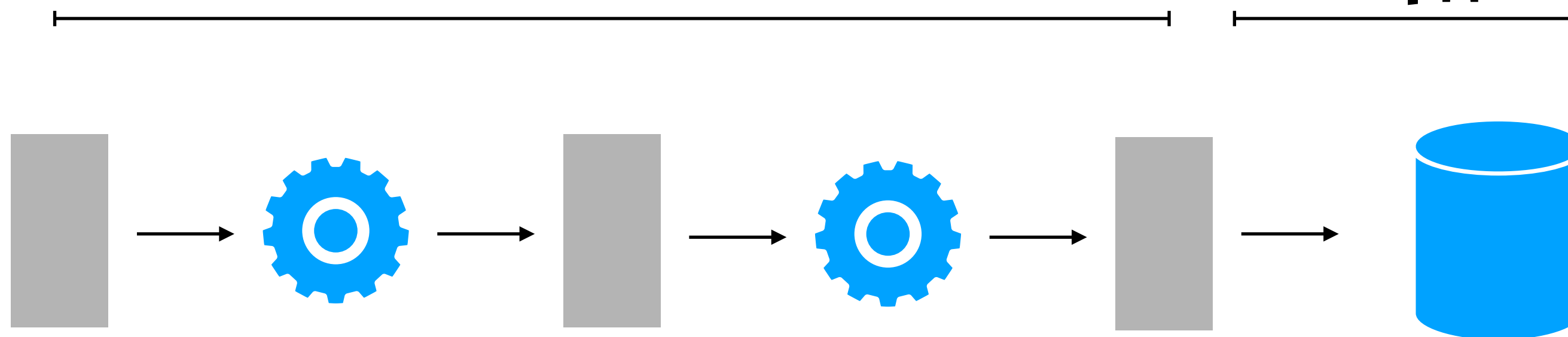
(Kafka Streams)

processing.guarantee=exactly_once

end-to-end eos

Kafka EoS

**Необходима
дополнительная
координация**



(Kafka Streams)

(Kafka Connect Sink)

`processing.guarantee=exactly_once`

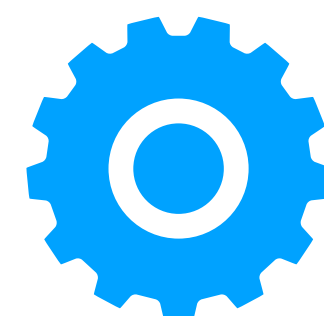
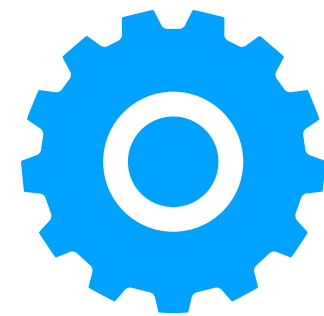
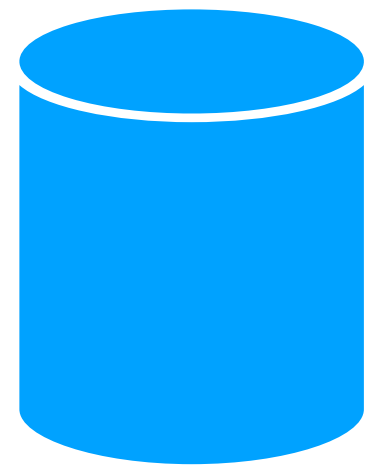
HDFS, Elasticsearch, S3,

end-to-end eos

Kafka EoS??

Kafka EoS

**Необходима
дополнительная
координация**



(Kafka Connect Source)

Work in progress

(Kafka Streams)

`processing.guarantee=exactly_once`

(Kafka Connect Sink)

**HDFS, Elasticsearch,
S3, etc.**

Команда EOS

пацаны ваще ребята...

- Apurva Mehta, Damian Guy, Flavio Junqueira, Guozhang Wang, Matthias J. Sax, Sriram Subramanian
- More info: [google KIP-98](#)

Приходите к
нам работать!



<https://www.confluent.io/careers/>

