

Высоконагруженные системы мониторинга

Владимир Гурьянов

Технический директор Deckhouse
Observability

- С **2011** занимаюсь эксплуатацией приложений
- С **2021** работаю в компании «Флант»
- С **2023** занимаюсь развитием направления Observability



Deckhouse Observability Platform

DOP

Единая платформа вместо зоопарка систем мониторинга



10+

тысяч серверов
на мониторинге

100+

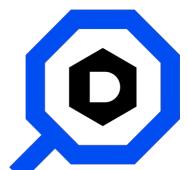
миллионов
метрик в минуту

22+

тысяч приложений
в production

600+

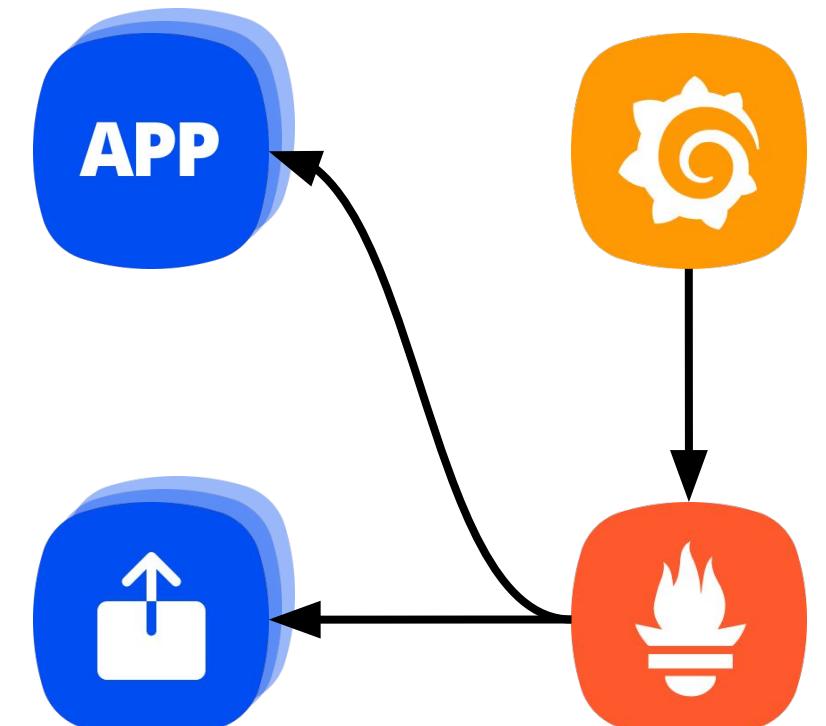
кластеров
на мониторинге



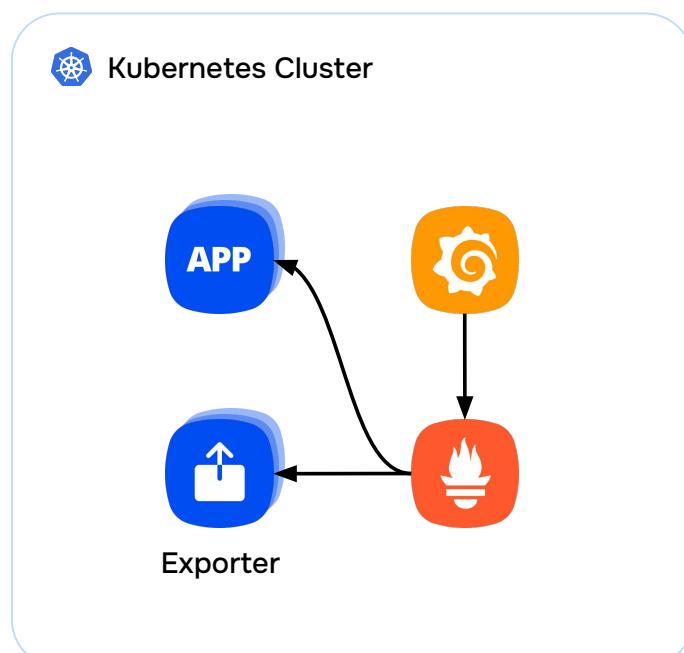
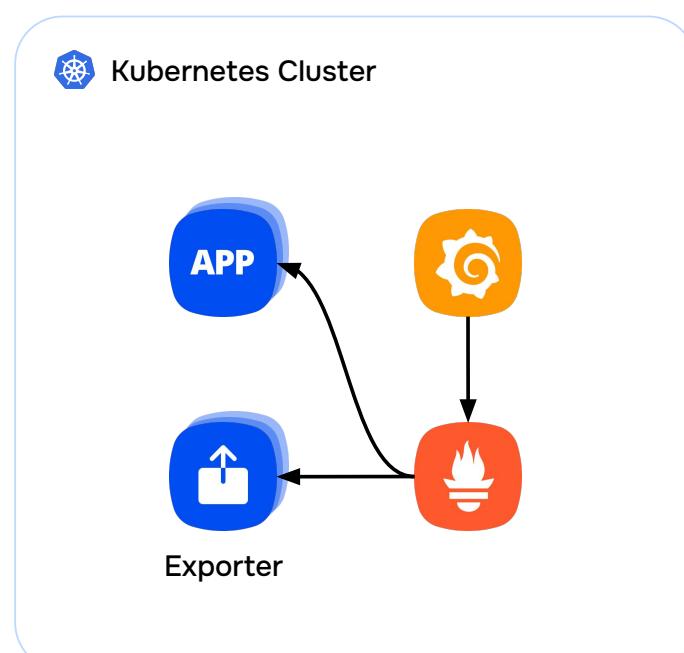
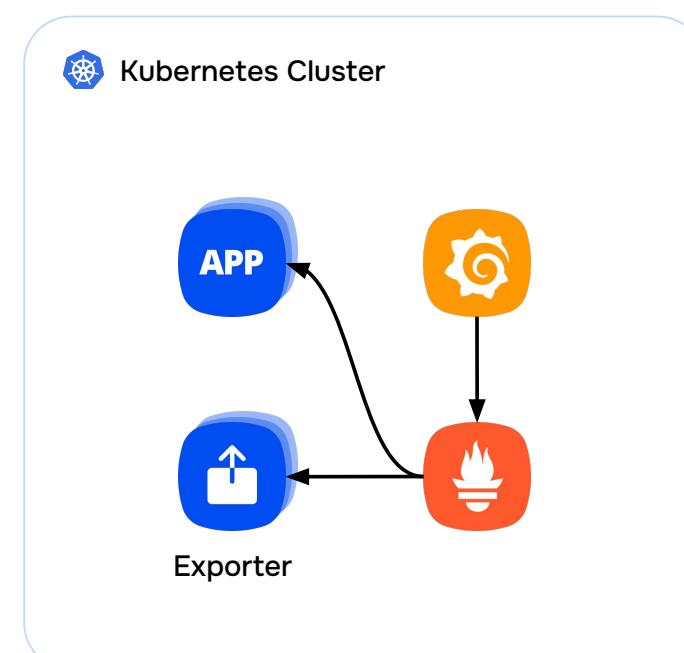
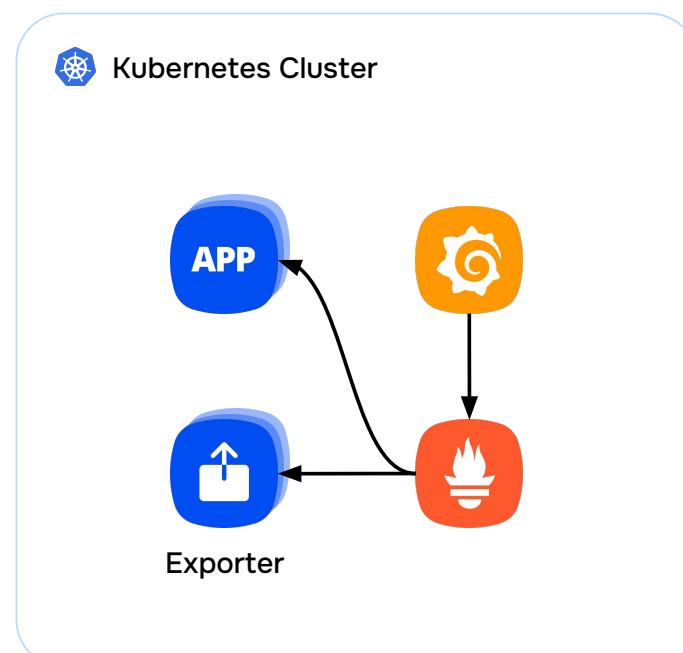
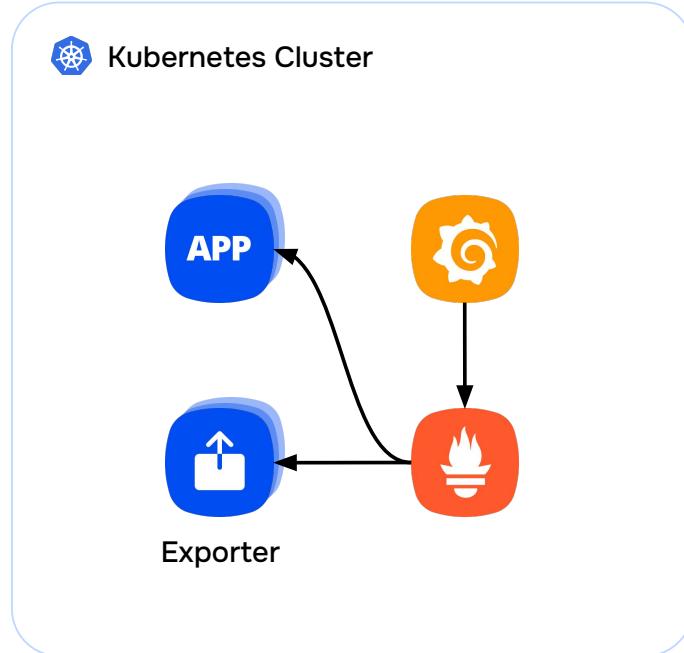
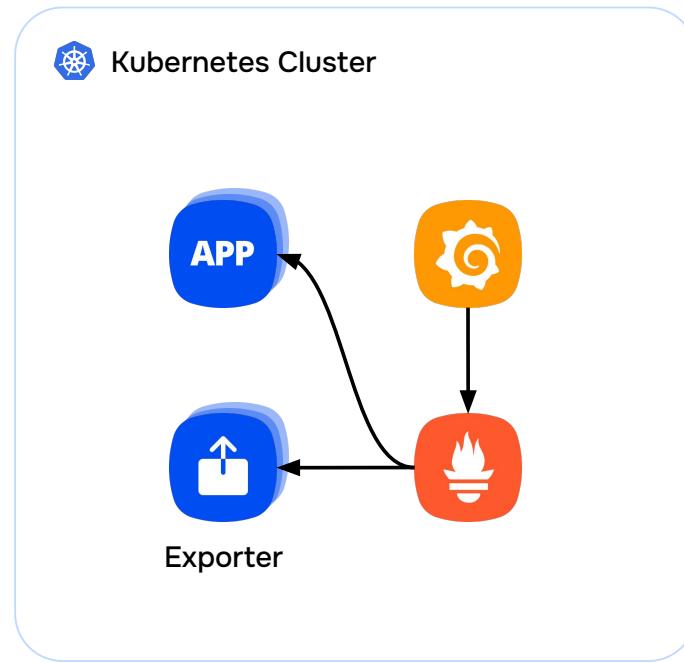
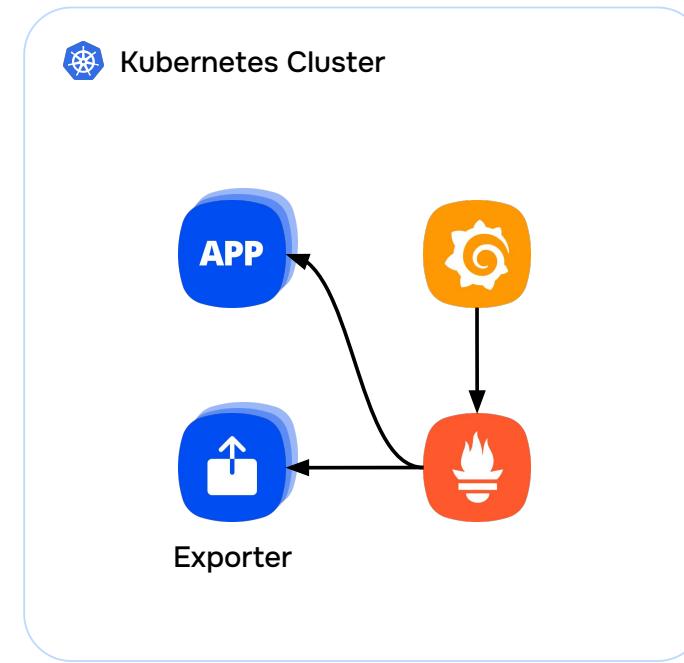
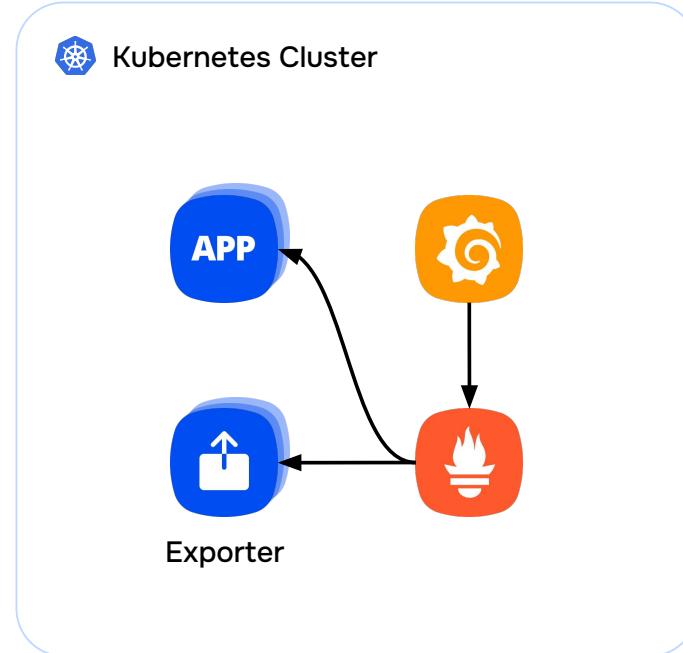
Deckhouse
Observability Platform

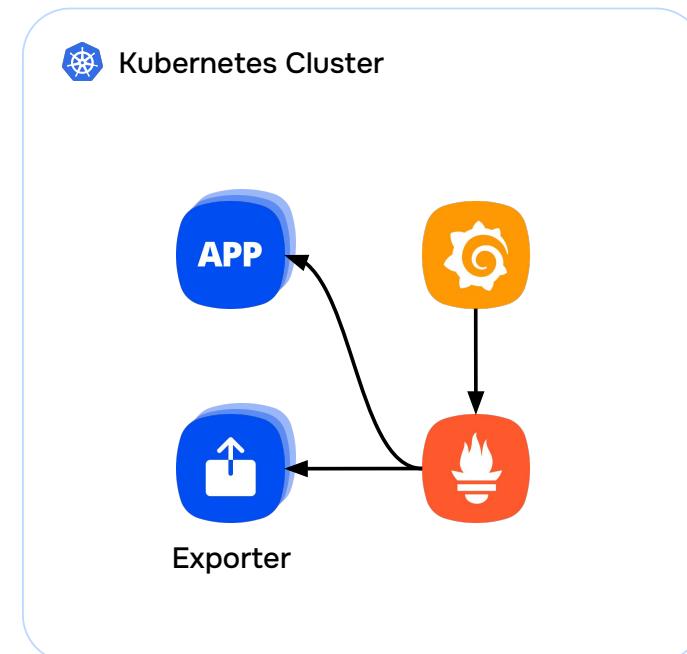
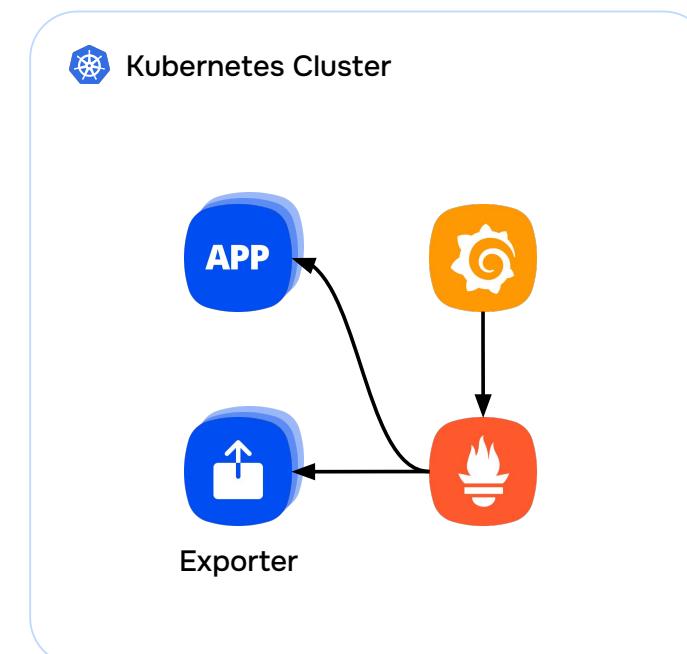
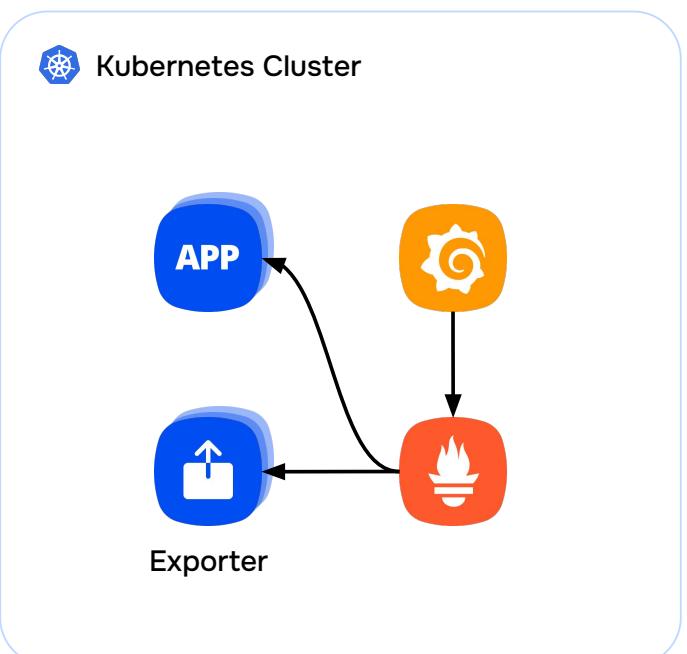
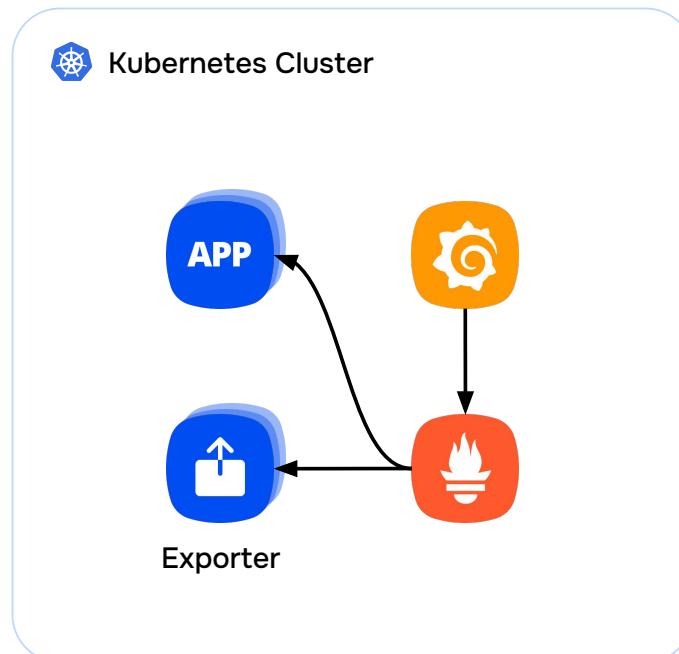
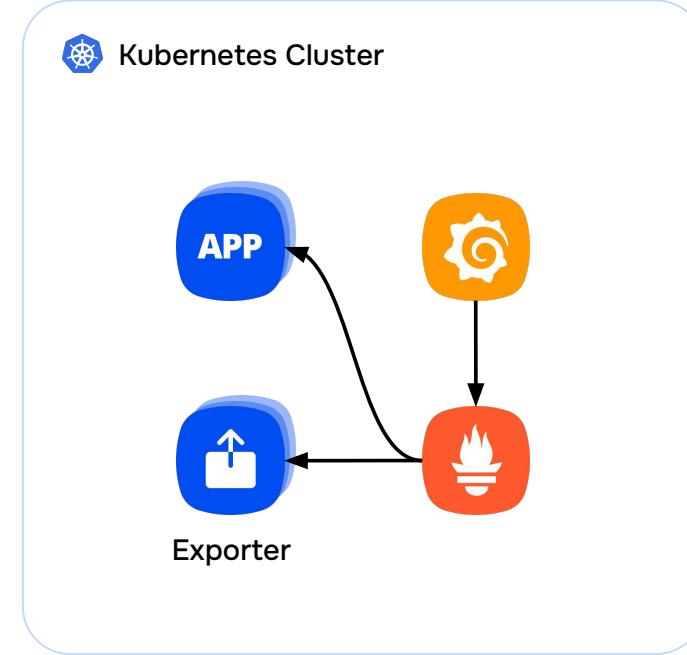
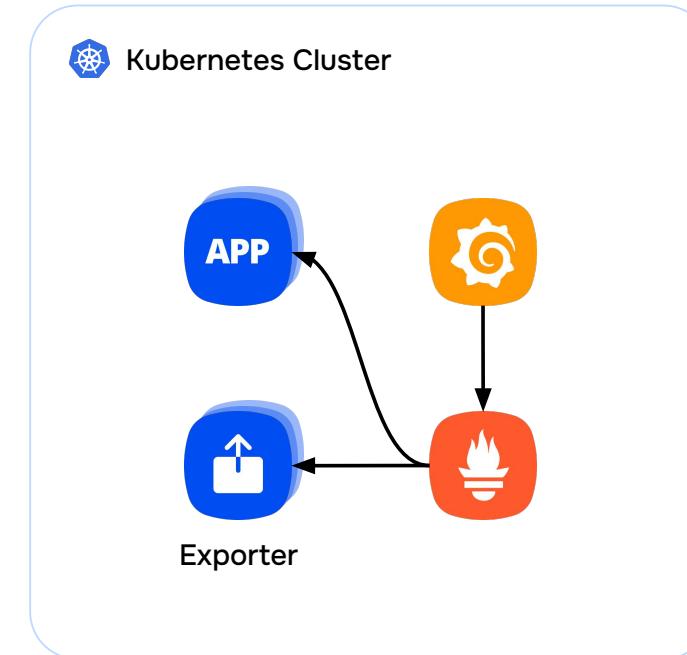
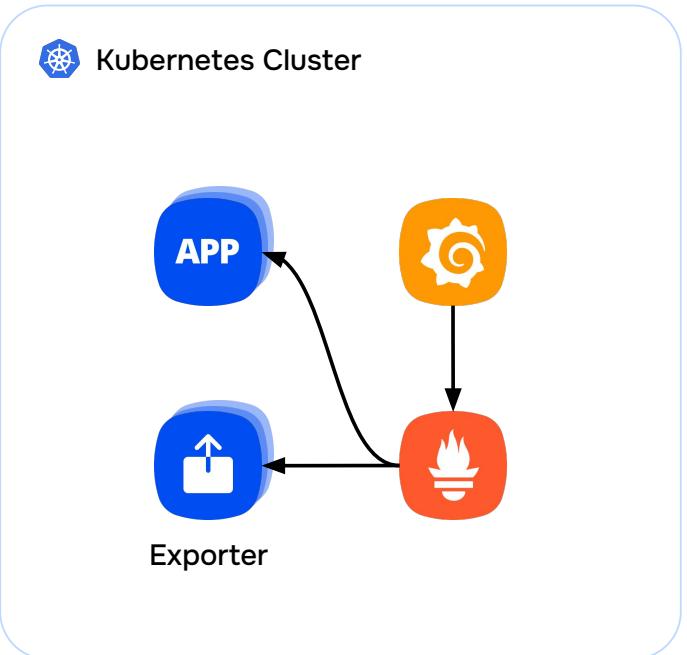
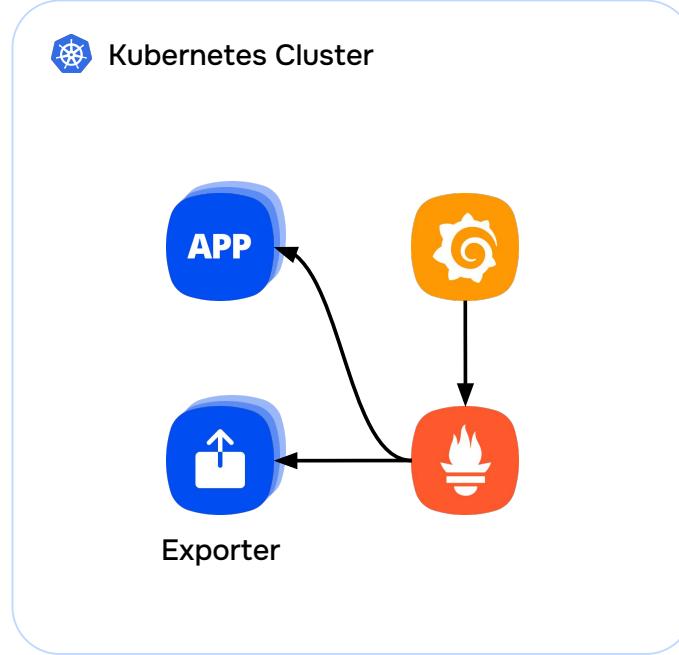


Kubernetes Cluster

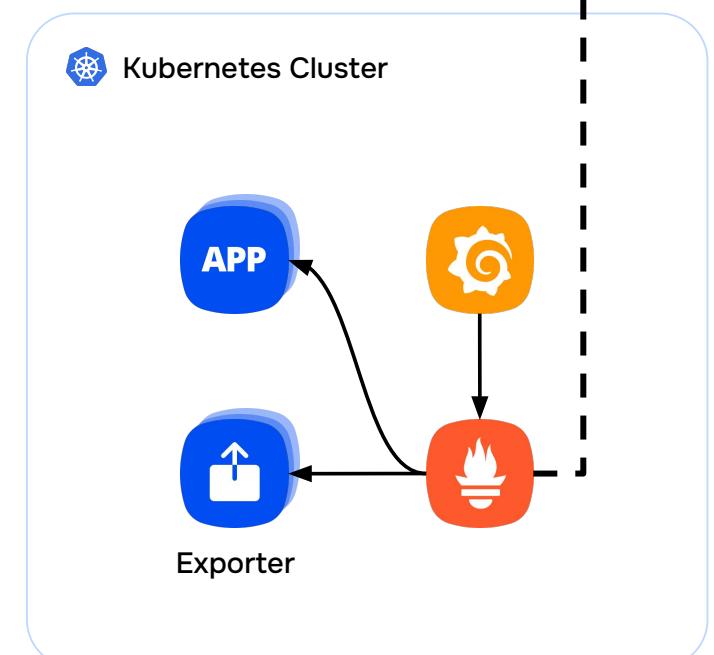
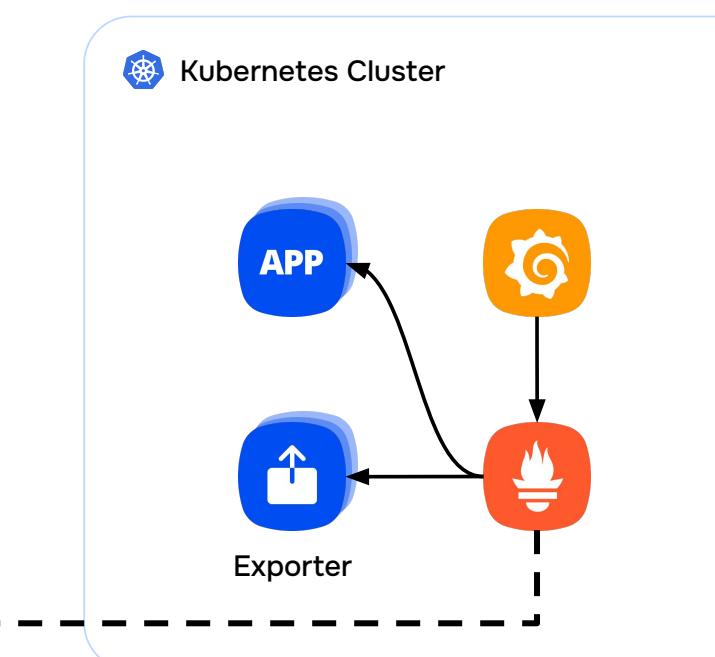
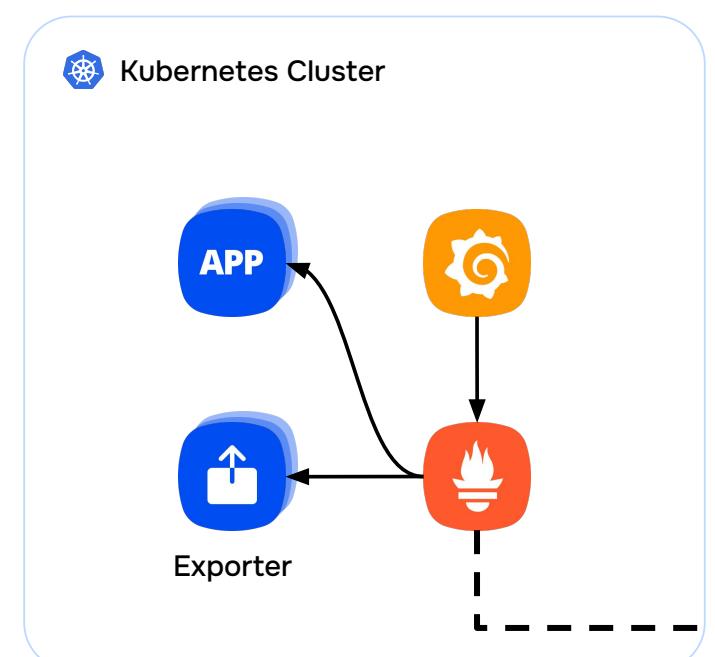
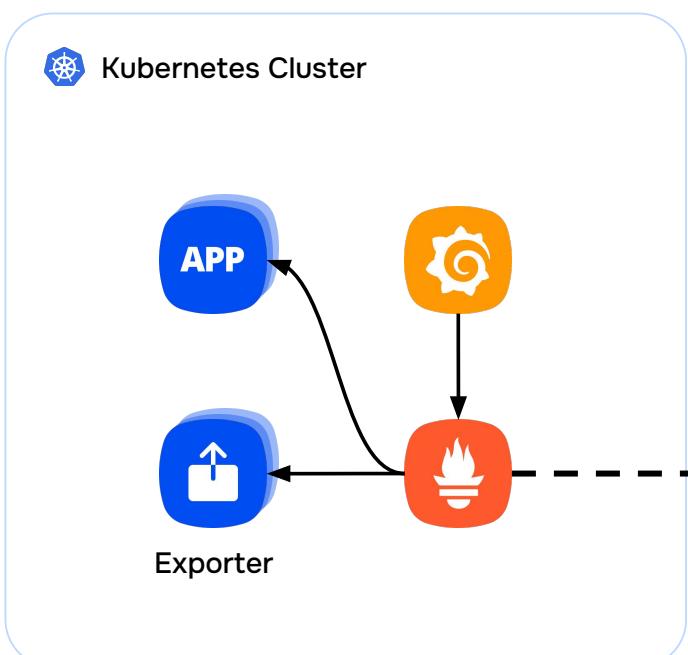
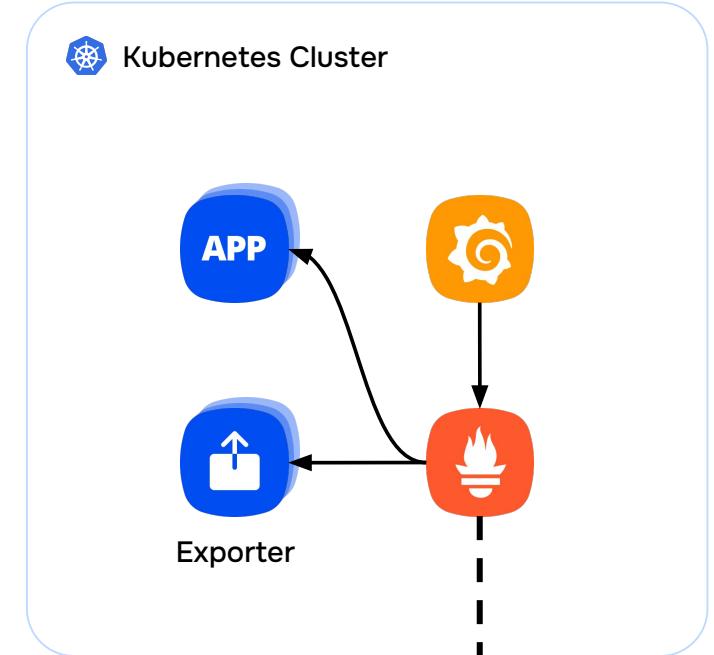
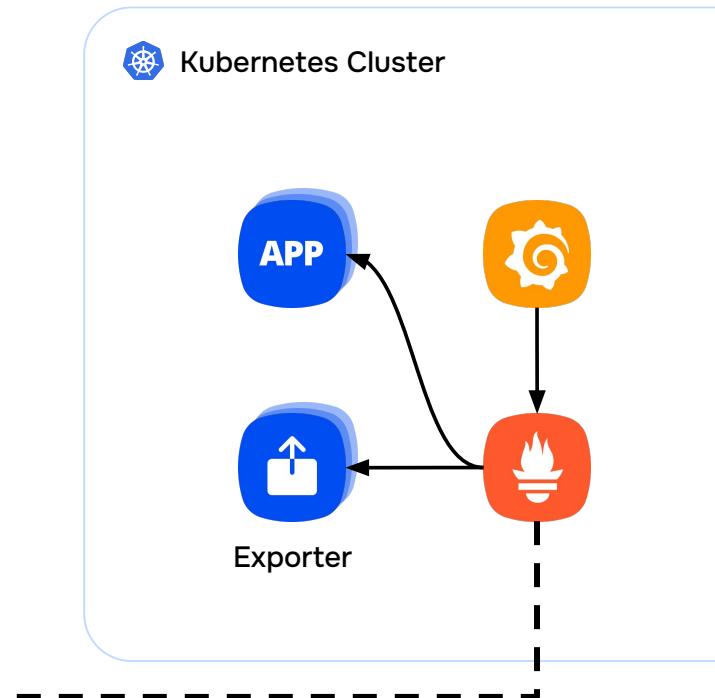
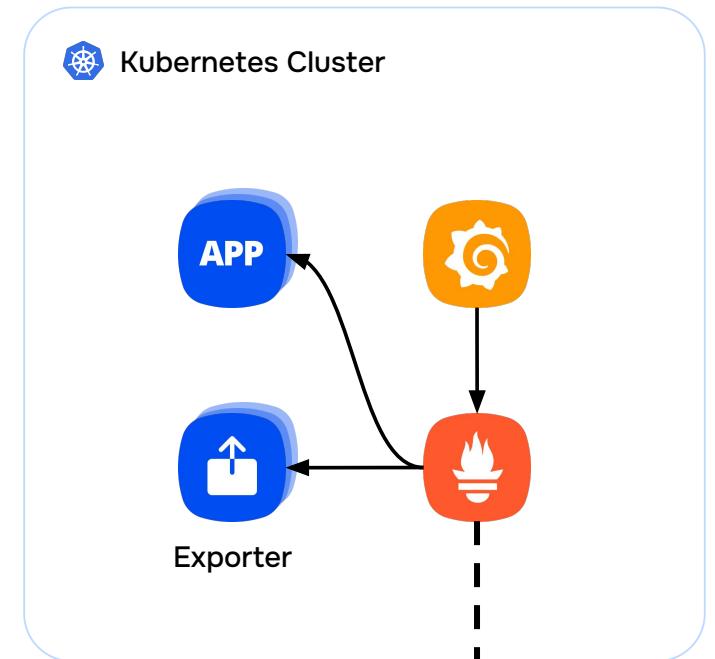
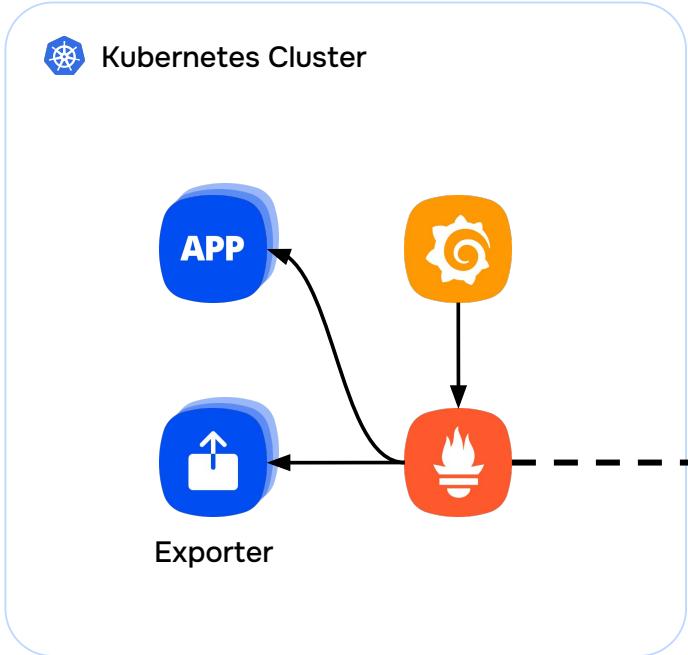


Exporter





Централизованное хранилище метрик



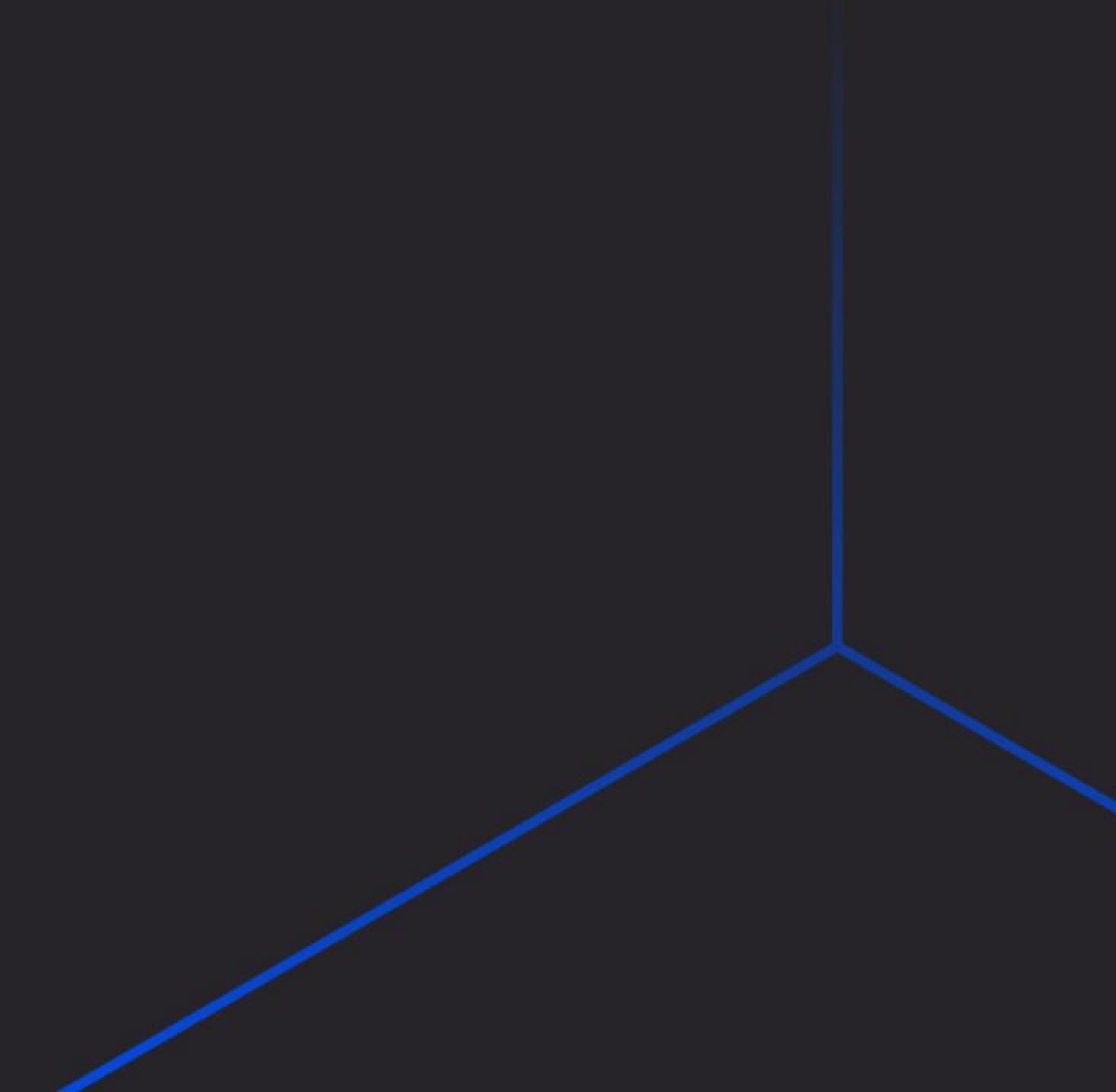
Централизованное
хранилище метрик

1 маленький кластер → 800 000+

1 маленький кластер → 800 000+

20 кластеров = $800\ 000 \times 20 = 16\ 000\ 000$

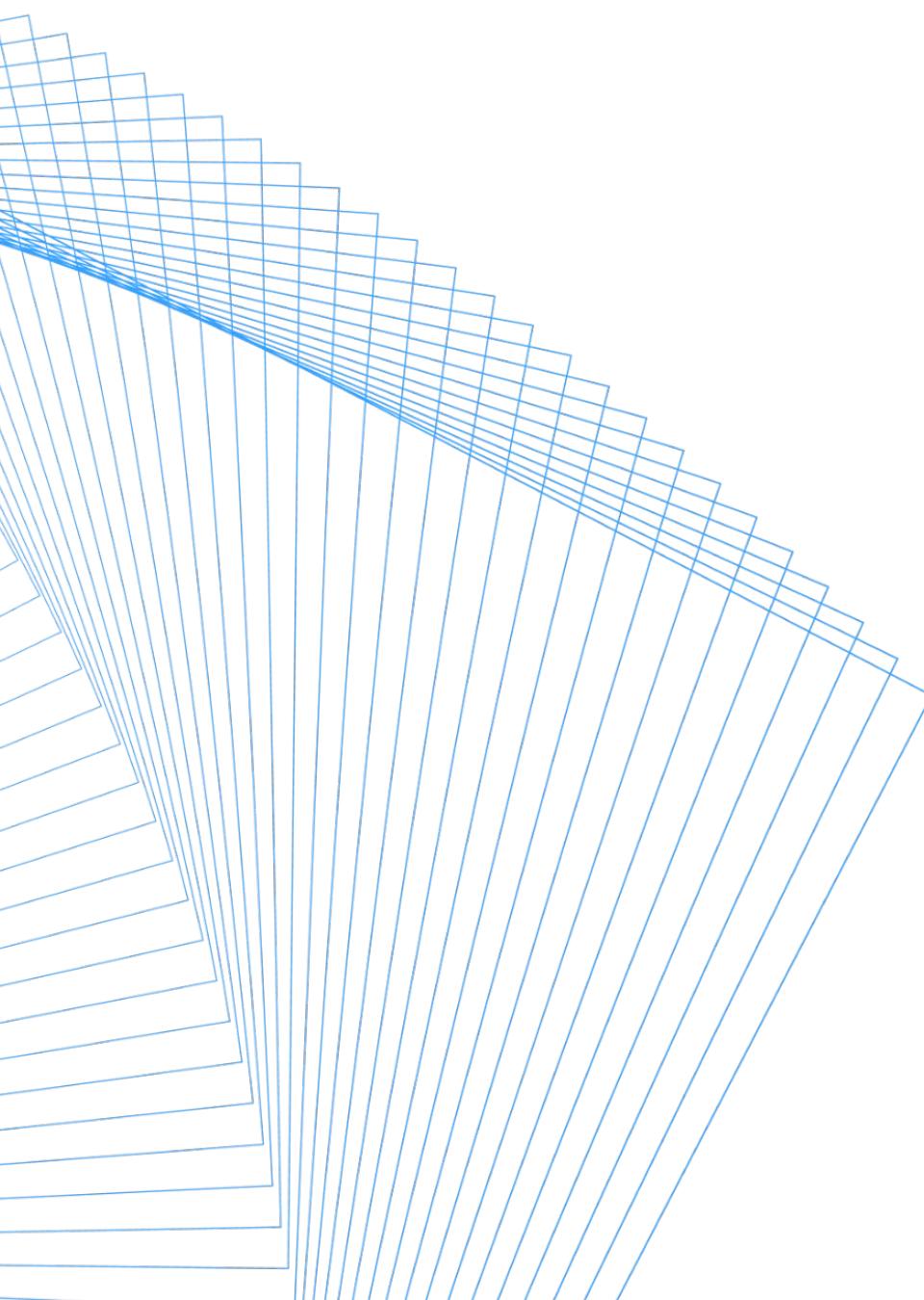
Требования к централизованной системе мониторинга



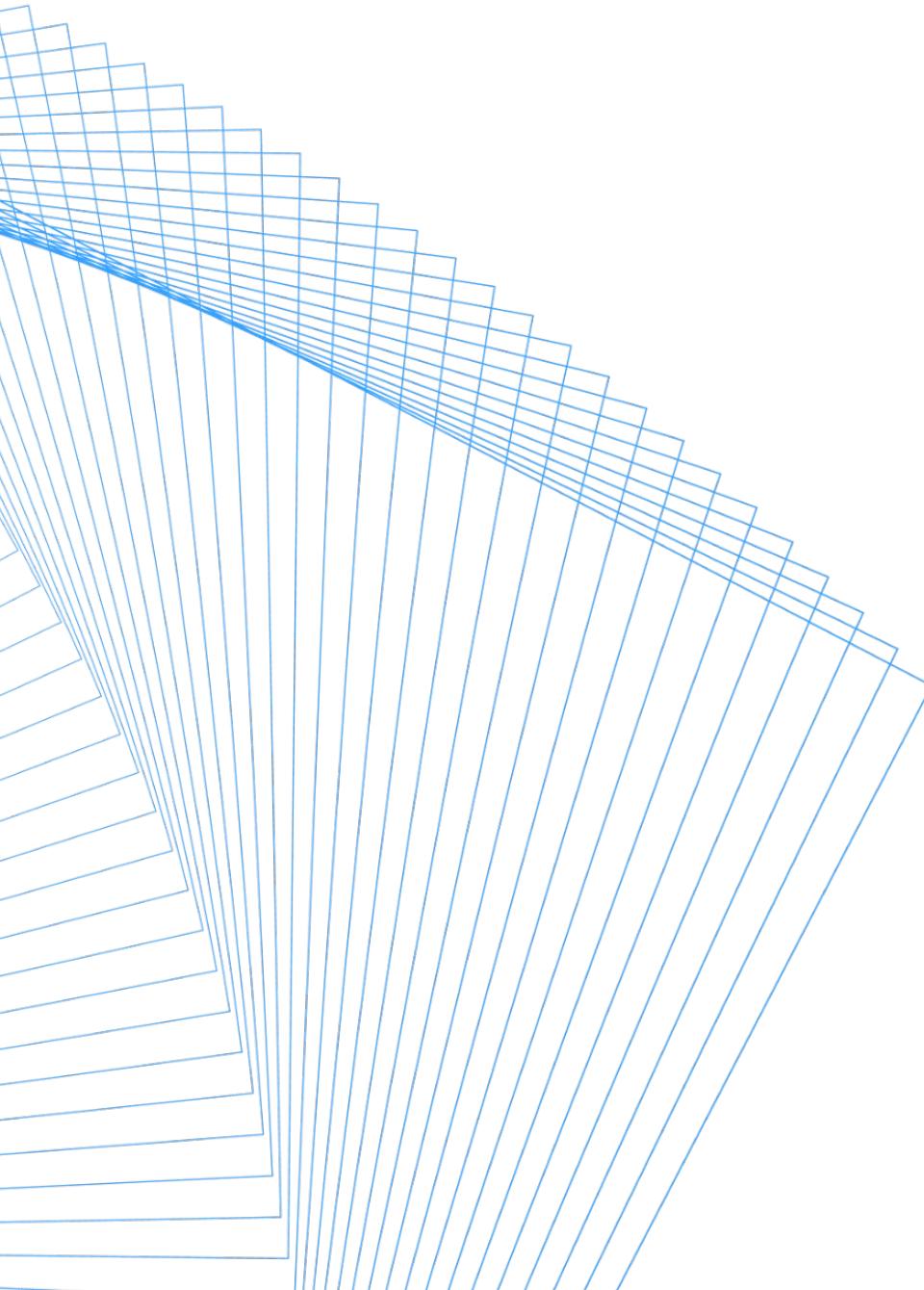
Требования

1

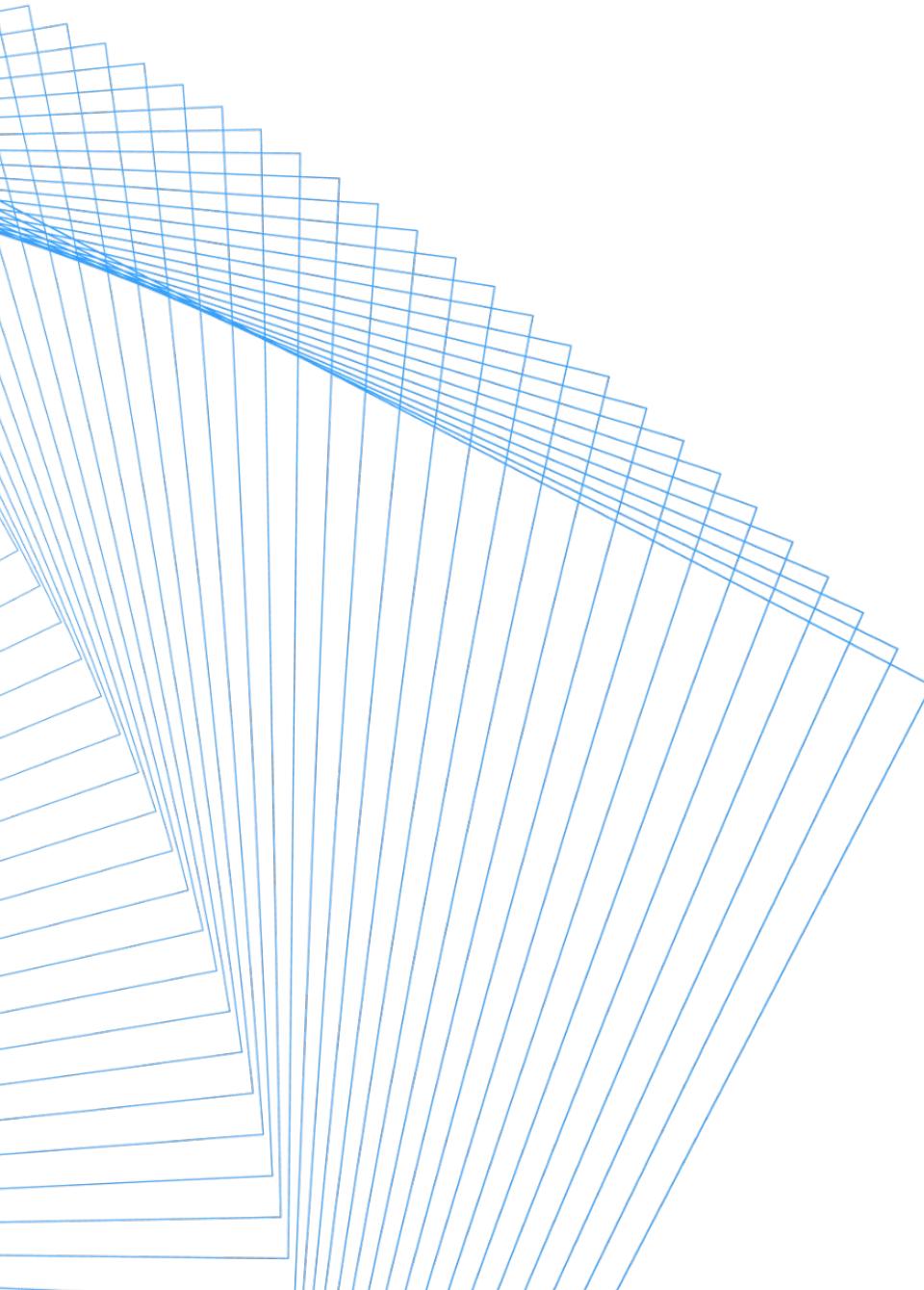
Отказоустойчивость



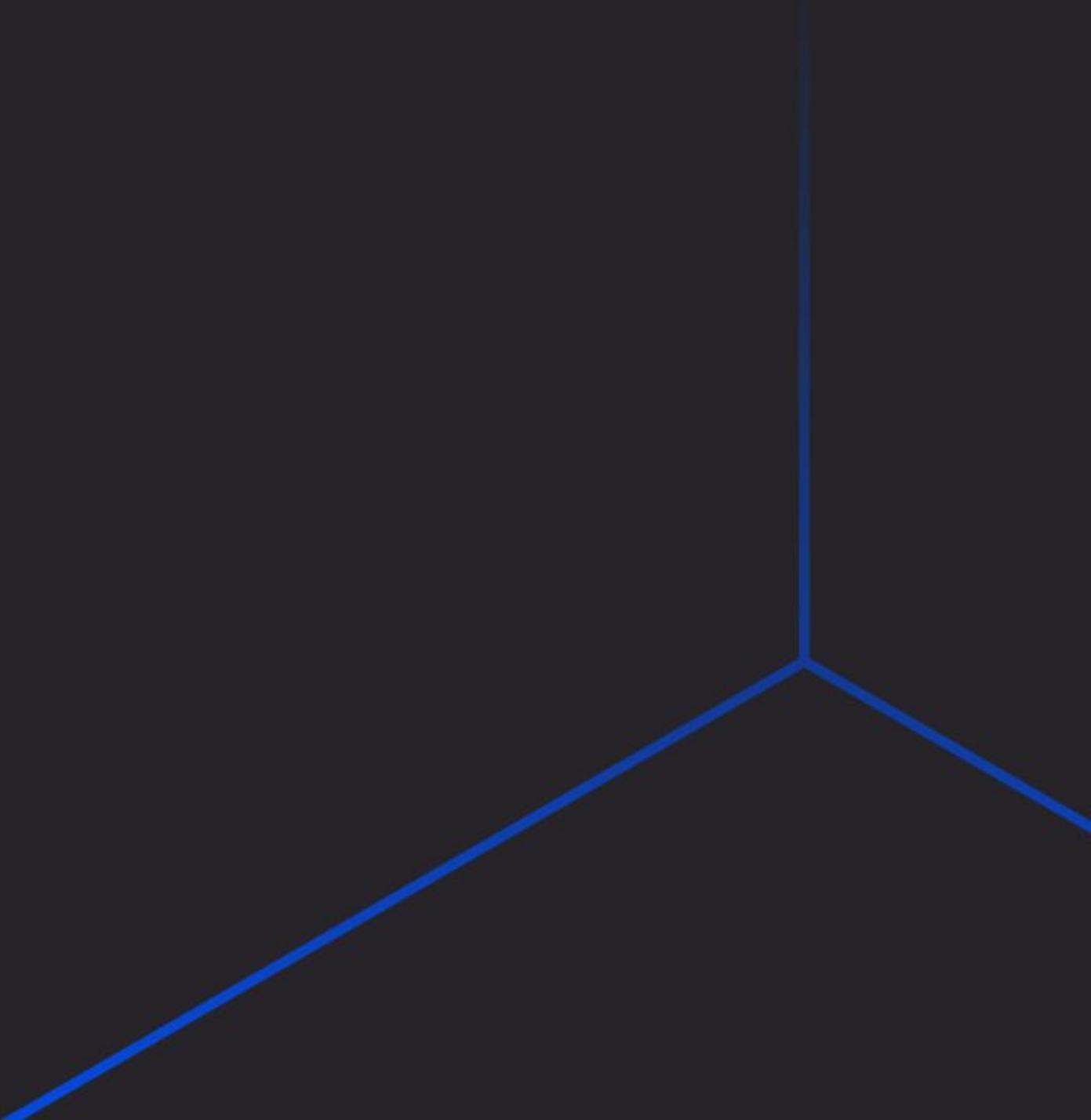
Требования

- 
- 1 Отказоустойчивость
 - 2 Шардирование данных

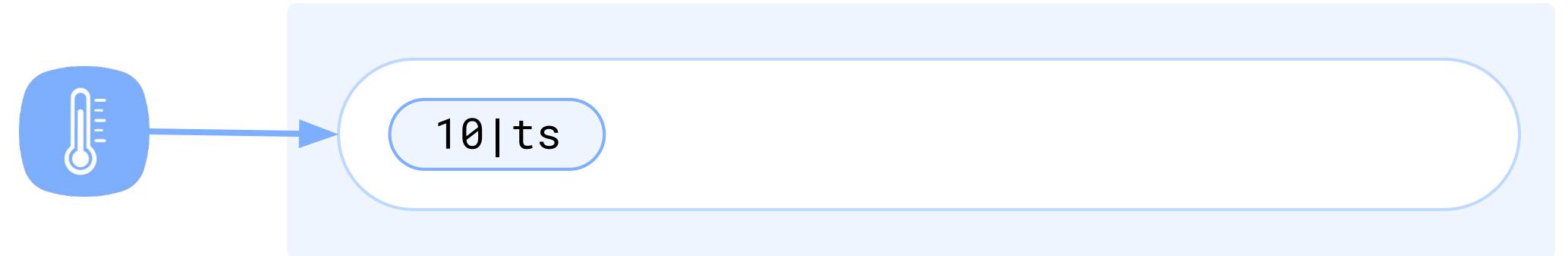
Требования

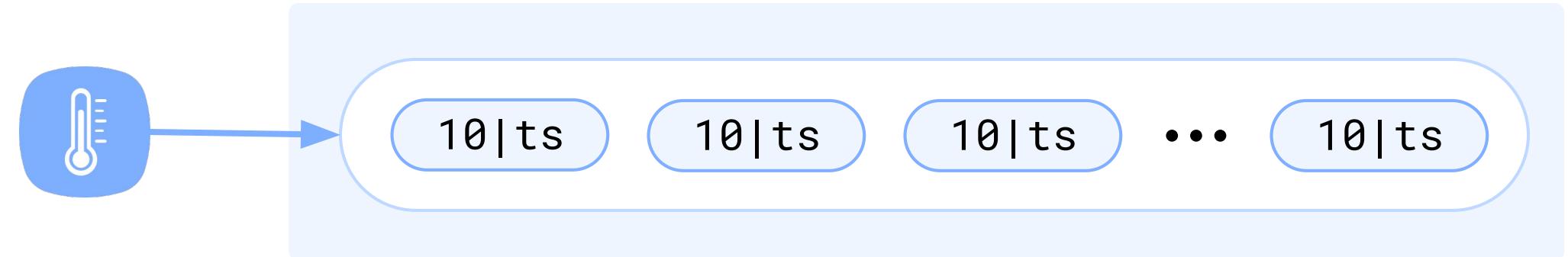
- 
- 1 Отказоустойчивость
 - 2 Шардирование данных
 - 3 Минимальная стоимость

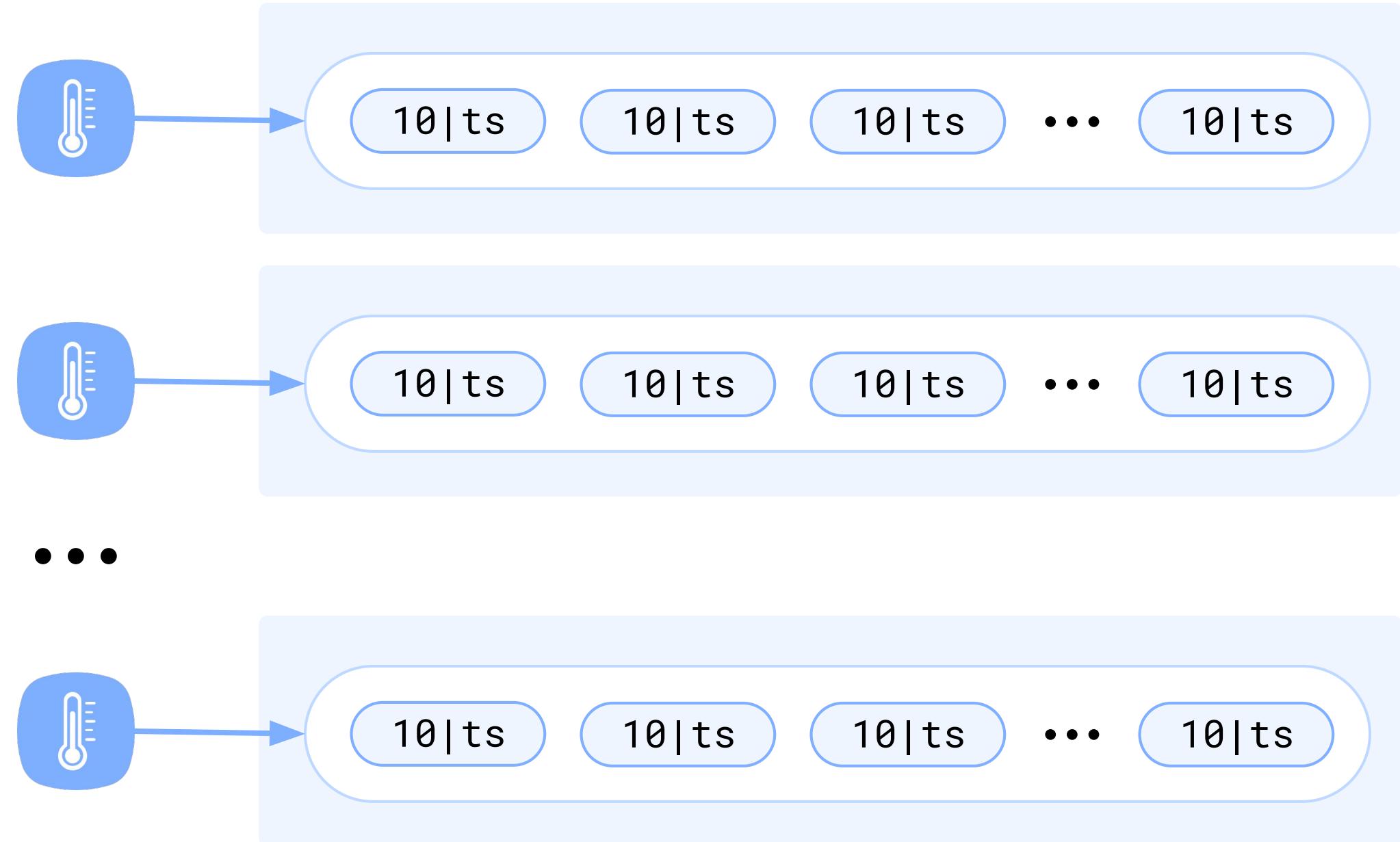
Что такое метрики и как они хранятся?

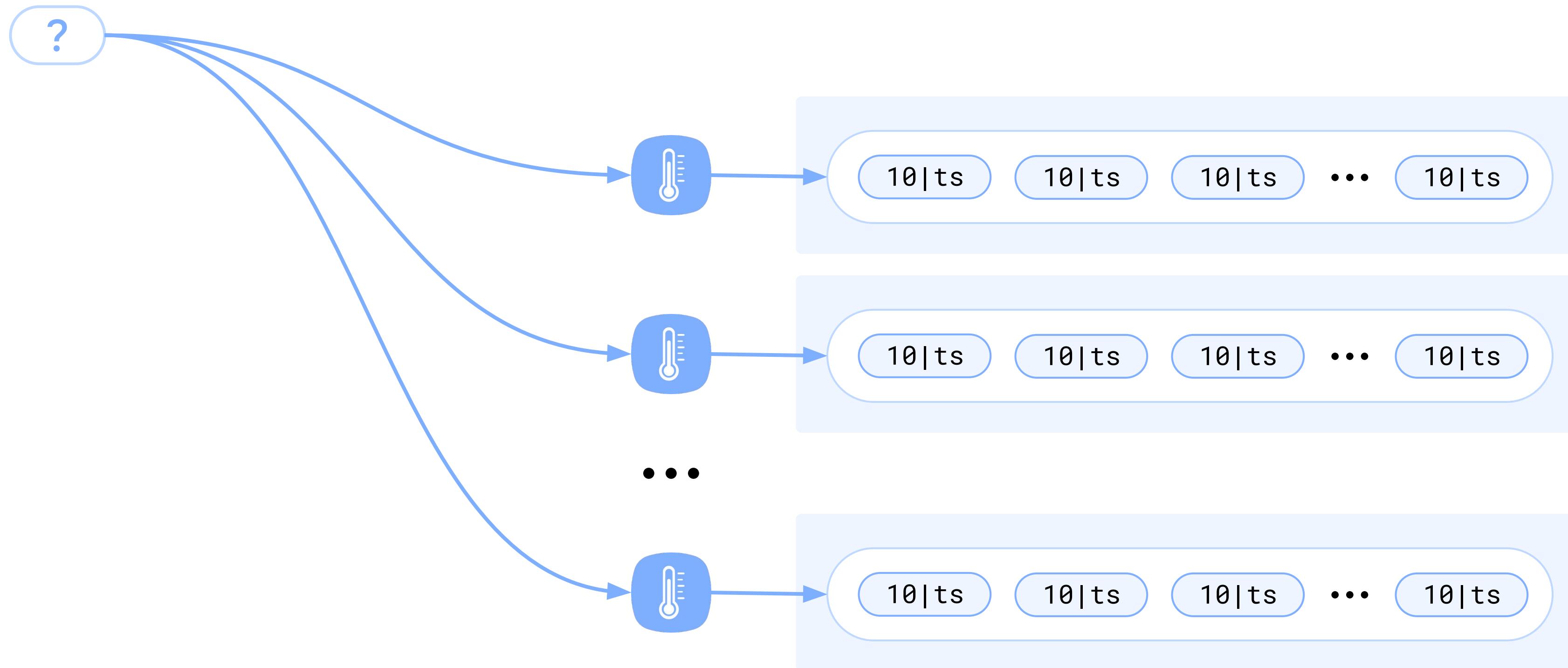












LabelSet



LabelSet

```
{  
    ключ: значение,  
    ключ: значение,  
    ...  
}
```

LabelSet

```
{  
    ключ: значение,  
    ключ: значение,  
    ...  
}
```

```
{  
    name: cpu_usage,  
    node: curiosity,  
    core: 4  
}
```

LabelSet

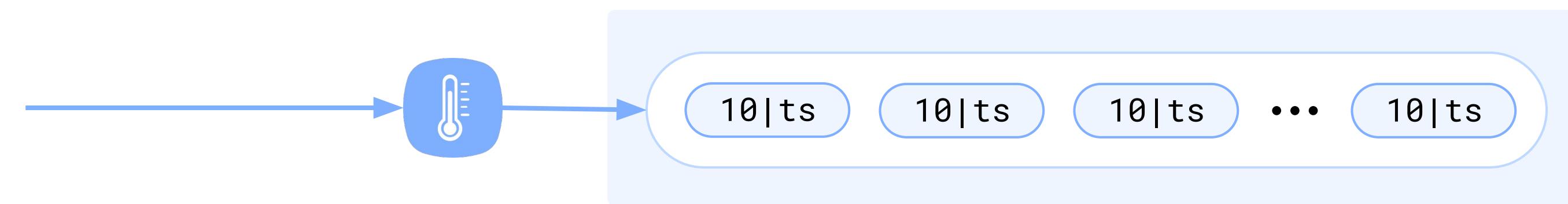


LabelSet

```
{  
  name:  cpu_usage,  
  node:  curiosity,  
  core:  0  
}
```

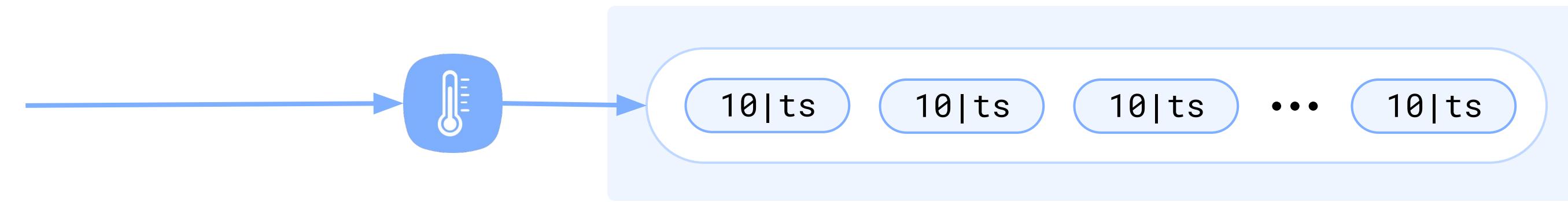


```
{  
  name:  cpu_usage,  
  node:  curiosity,  
  core:  1  
}
```



• • •

```
{  
  name:  cpu_usage,  
  node:  curiosity,  
  core:  999999  
}
```



{|s}

```
{  
  name:  cpu_usage,  
  node:  curiosity,  
  core:  0  
}
```



```
{  
  name:  cpu_usage,  
  node:  curiosity,  
  core:  1  
}
```



• • •

```
{  
  name:  cpu_usage,  
  node:  curiosity,  
  core:  999999  
}
```



```
{ls} → #t_id
```

```
{  
  name: cpu_usage,  
  node: curiosity,  
  core: 0  
}
```

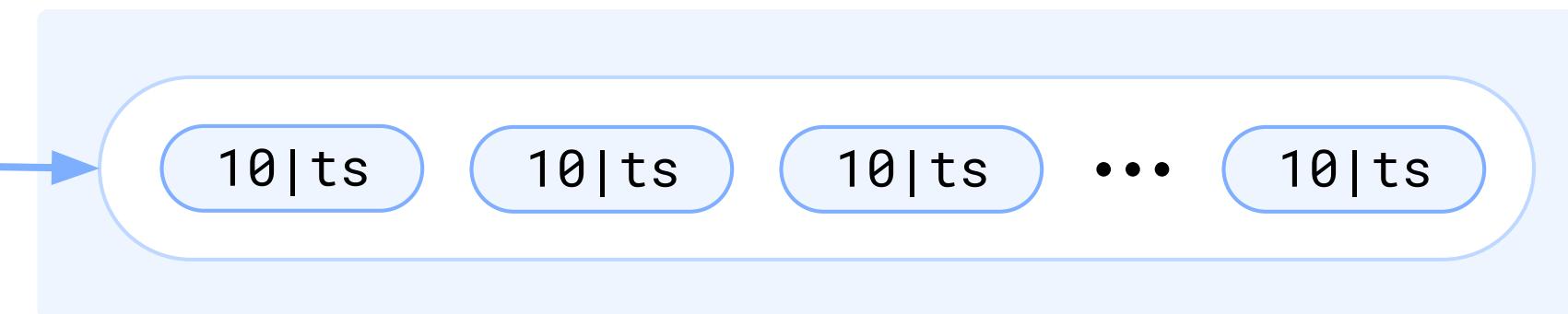


```
{  
  name: cpu_usage,  
  node: curiosity,  
  core: 1  
}
```



• • •

```
{  
  name: cpu_usage,  
  node: curiosity,  
  core: 999999  
}
```



{|s} → #t_id → Временной ряд

```
{  
  name: cpu_usage,  
  node: curiosity,  
  core: 0  
}
```



```
{  
  name: cpu_usage,  
  node: curiosity,  
  core: 1  
}
```

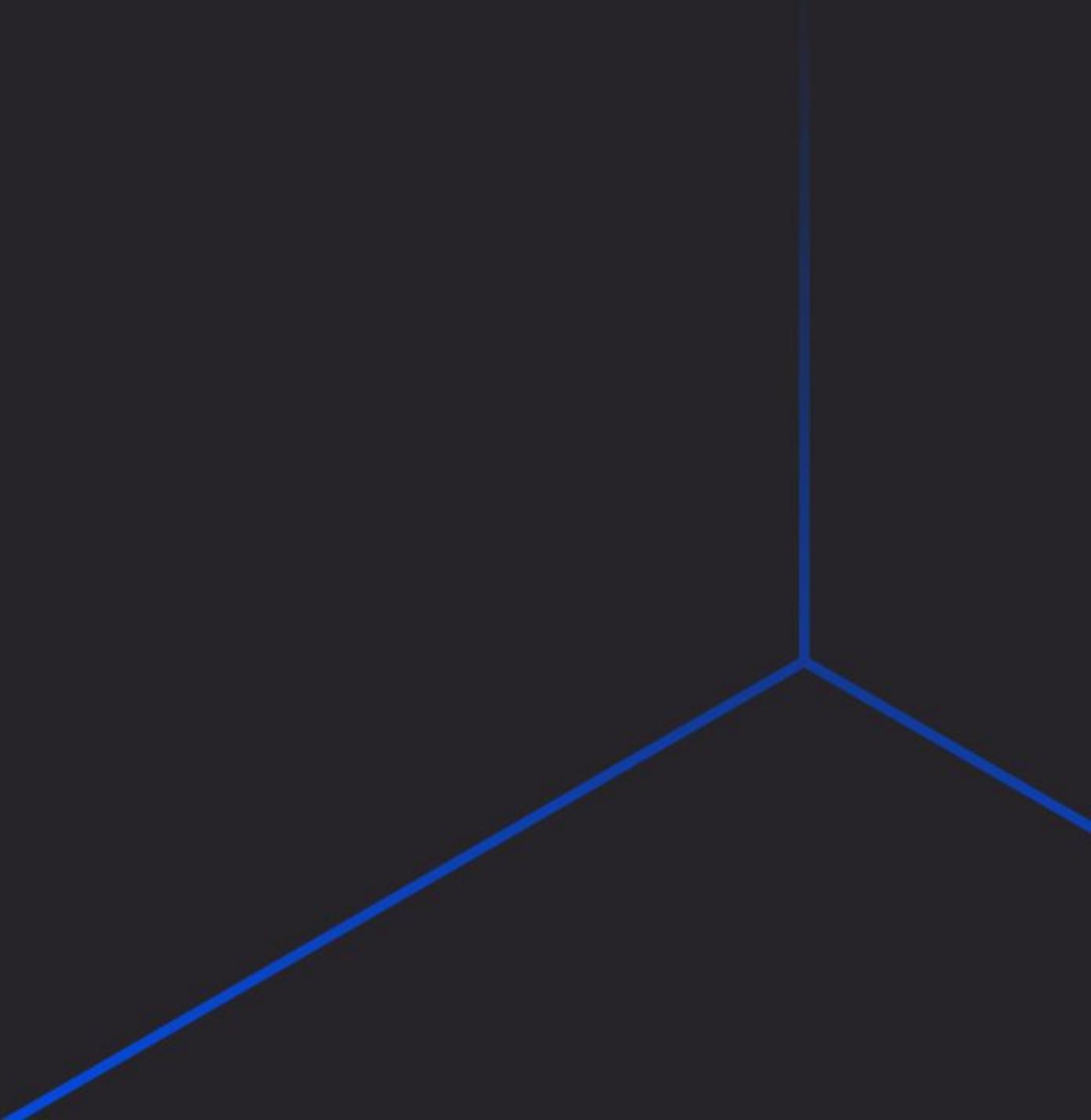


• • •

```
{  
  name: cpu_usage,  
  node: curiosity,  
  core: 999999  
}
```



Как хранятся метрики?





SSD

{ls}

{ls}

{ls}

1 000 000

...



RAM



SSD

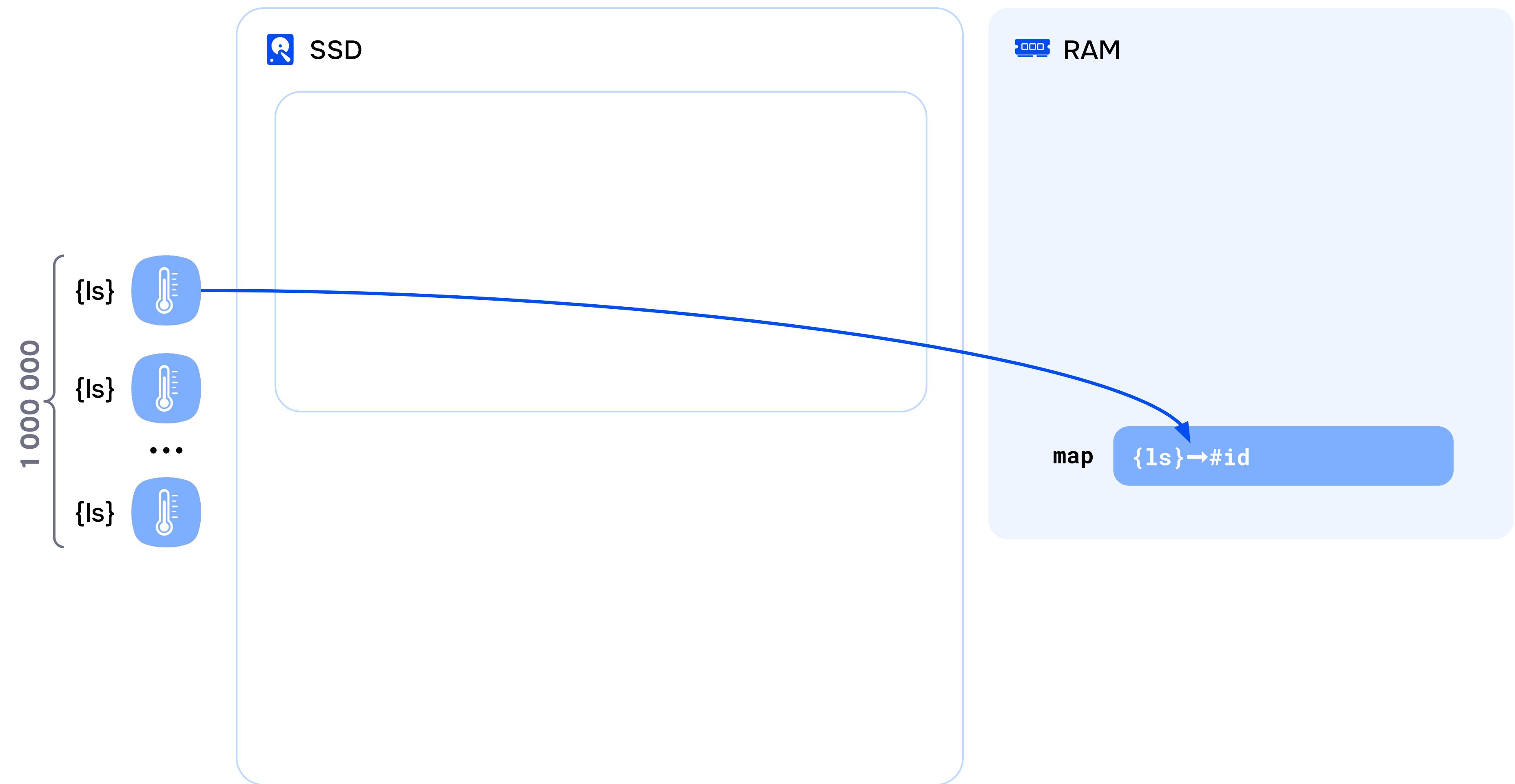
{ls}
{ls}
...
{ls}

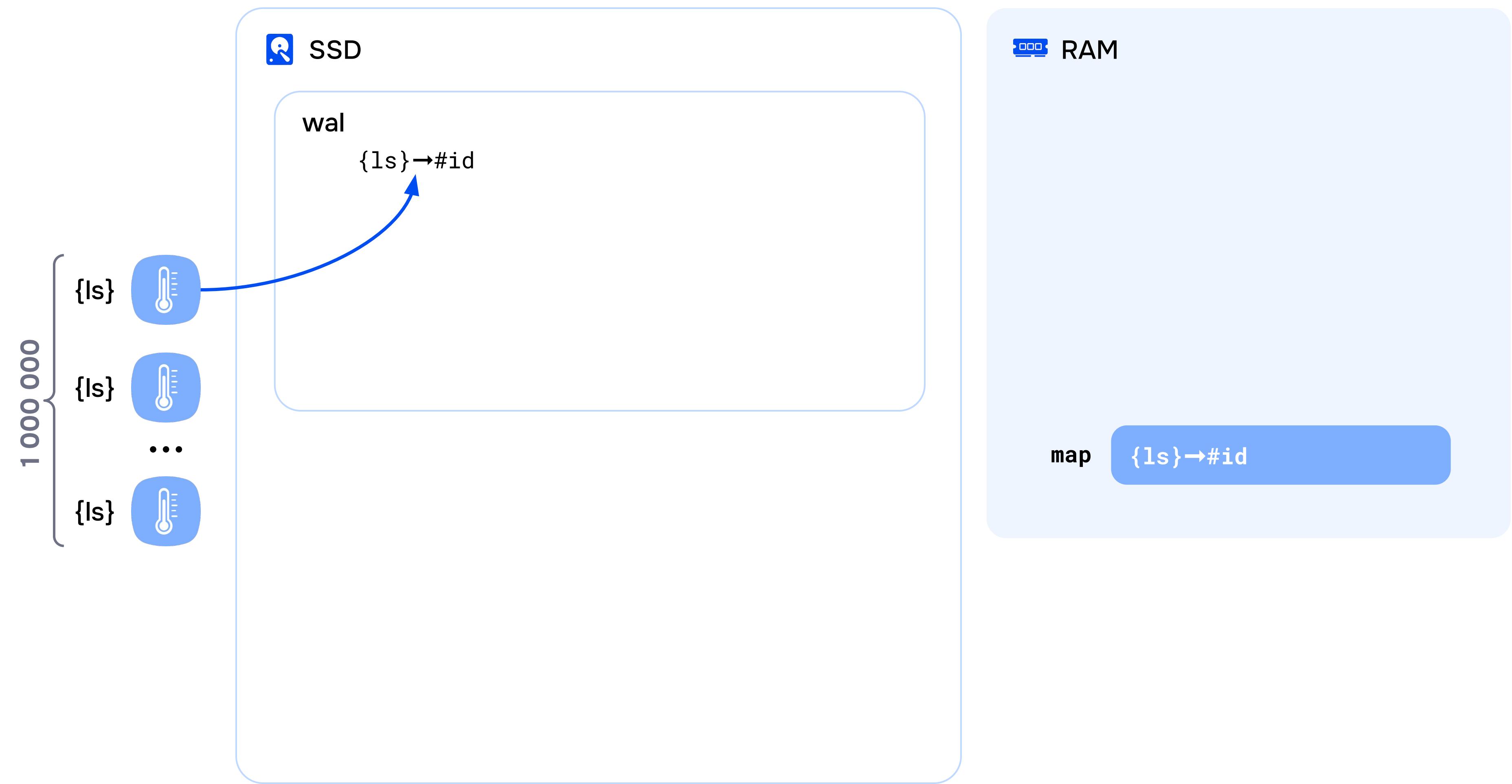
1 000 000



RAM

map





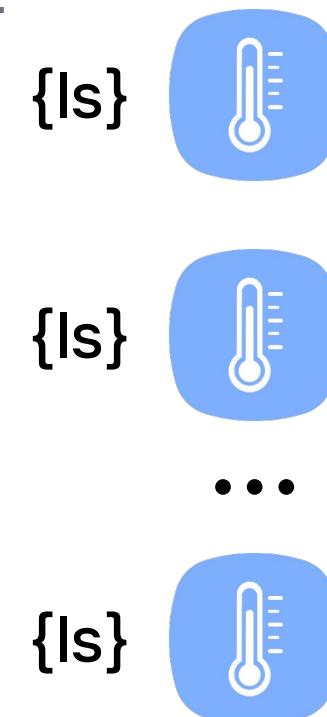


SSD

wal

{ls} → #id

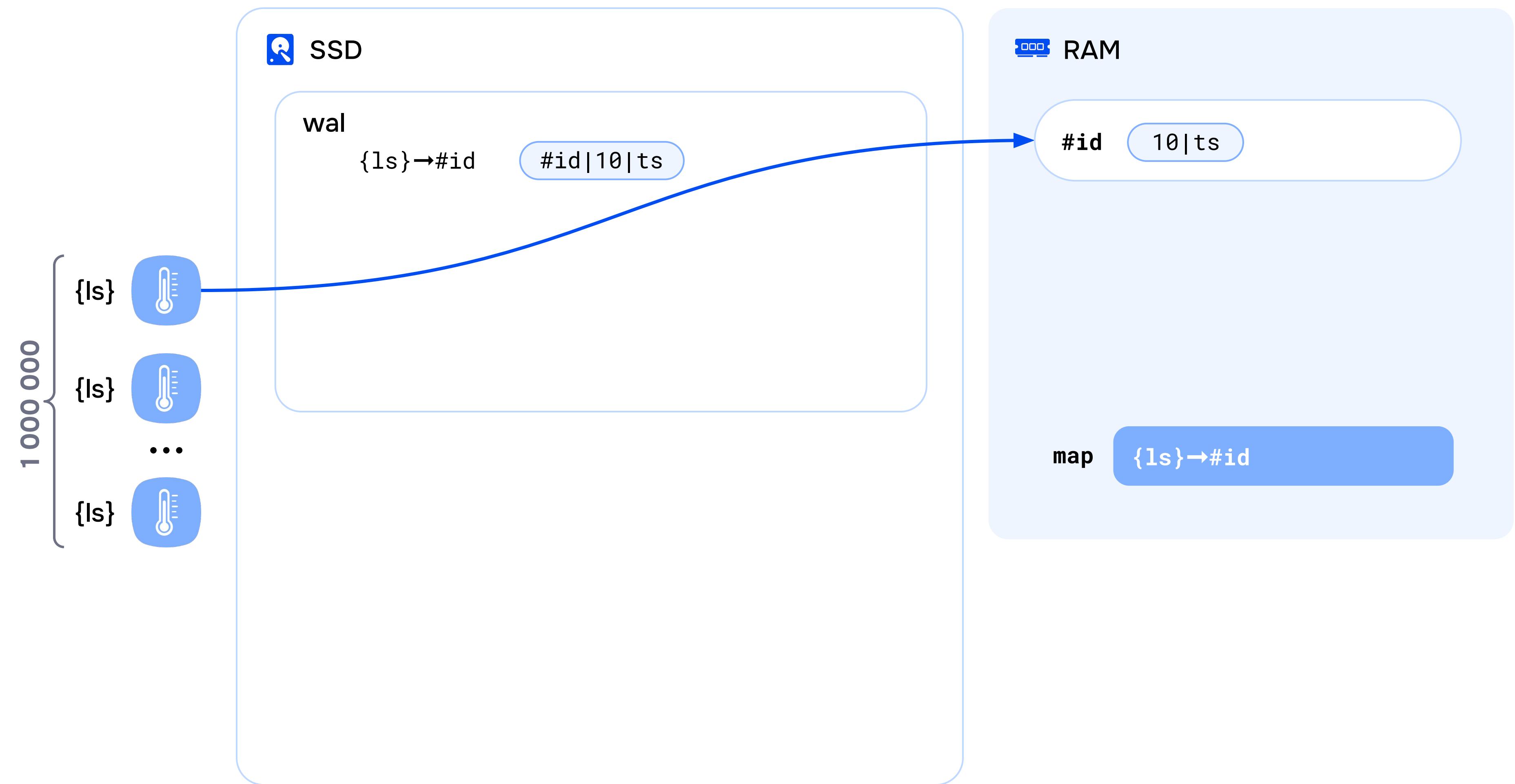
#id | 10 | ts



RAM

map

{ls} → #id



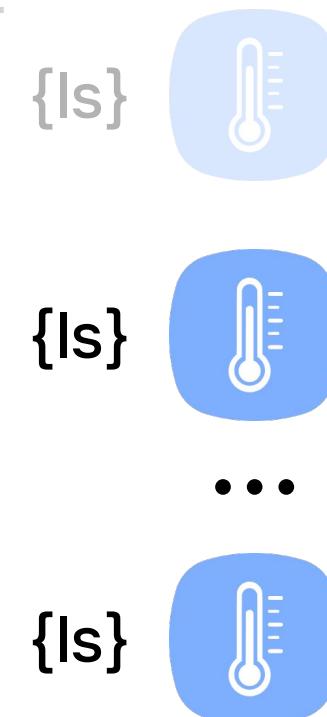


SSD

wal

{ls} → #id

#id | 10 | ts



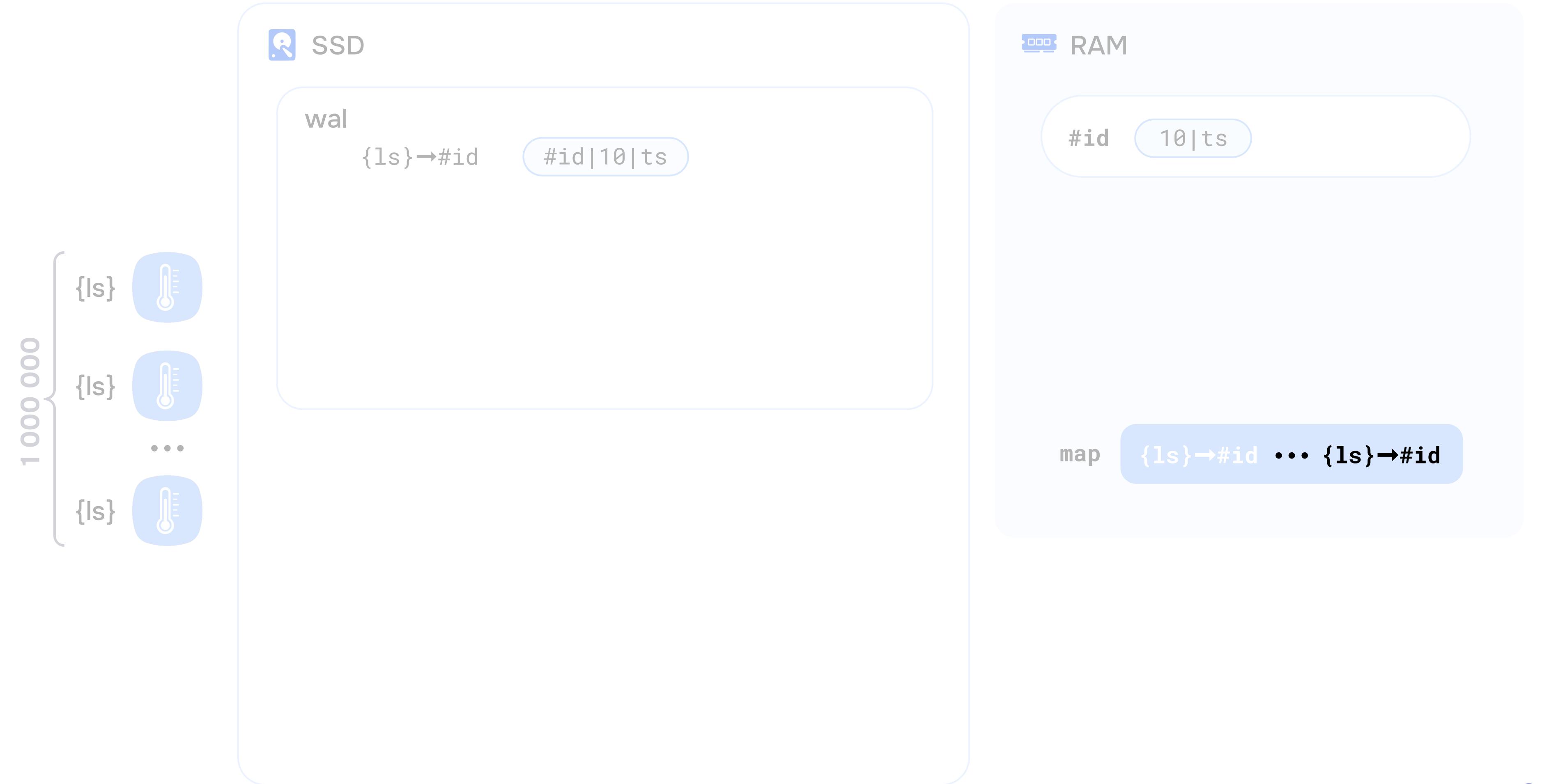
RAM

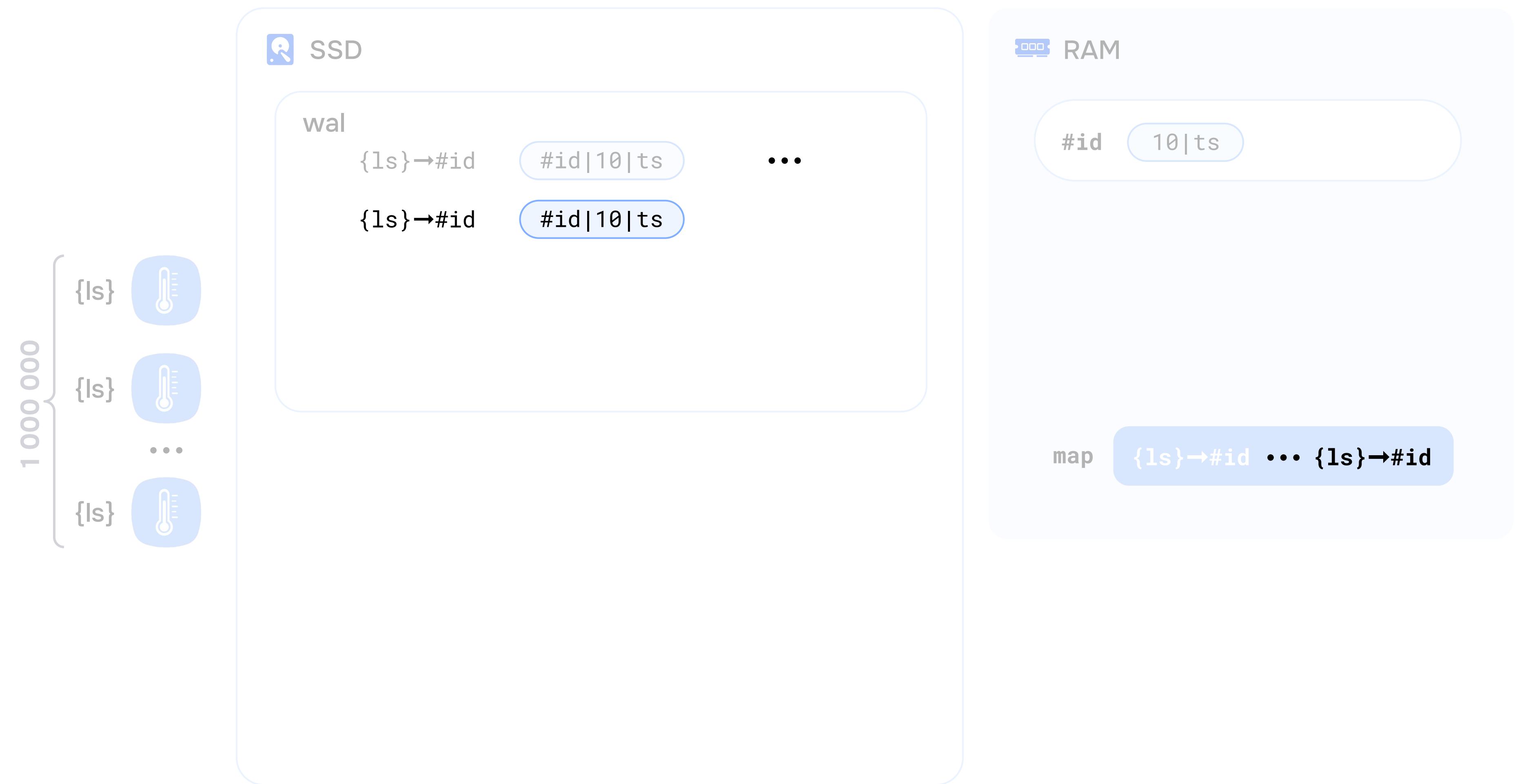
#id

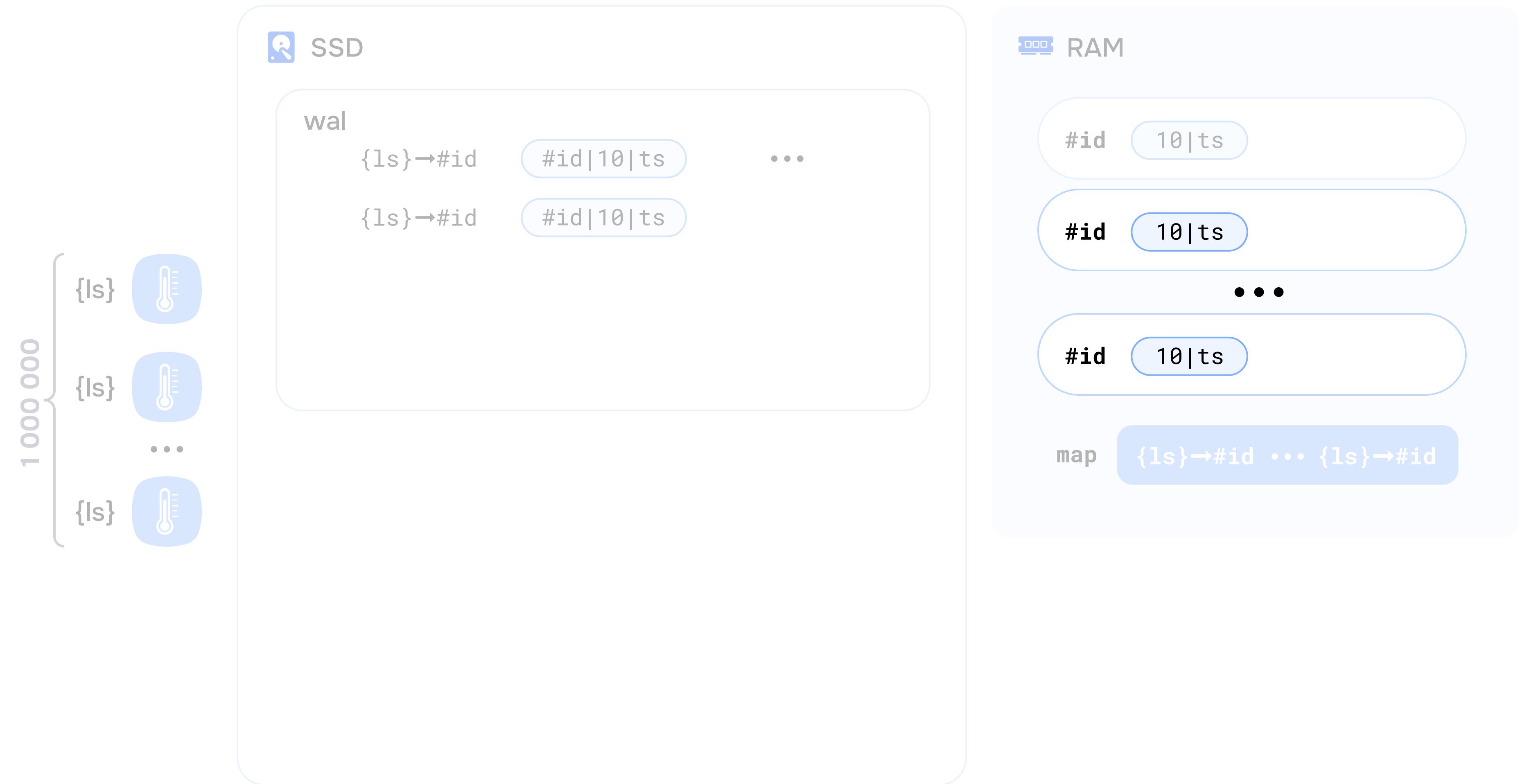
10 | ts

map

{ls} → #id









SSD

wal

{ls} → #id

#id|10|ts

...

{ls} → #id

#id|10|ts

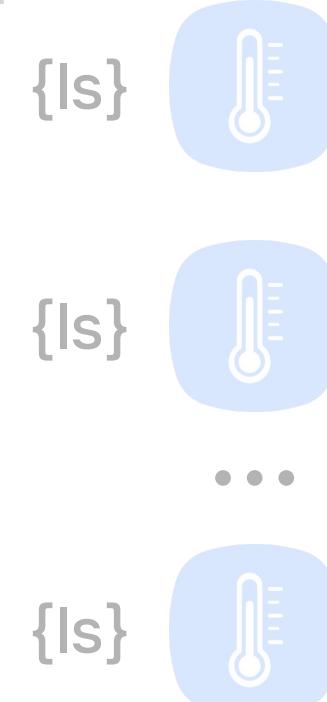
#id|10|ts

#id|10|ts

#id|10|ts

...

#id|10|ts



RAM

#id

10|ts

...

10|ts

#id

10|ts

...

10|ts

...

#id

10|ts

...

10|ts

map

{ls} → #id ... {ls} → #id



SSD

wal

{ls} → #id

#id|10|ts

...

{ls} → #id

#id|10|ts

#id|10|ts

#id|10|ts

#id|10|ts

...

#id|10|ts

1 000 000 000
{ls}
{ls}
...
{ls}



RAM

#id

10|ts

...

10|ts

#id

10|ts

...

10|ts

...

#id

10|ts

...

10|ts

map

{ls} → #id ... {ls} → #id



SSD

wal

{ls} → #id

#id|10|ts

...

{ls} → #id

#id|10|ts

#id|10|ts

...

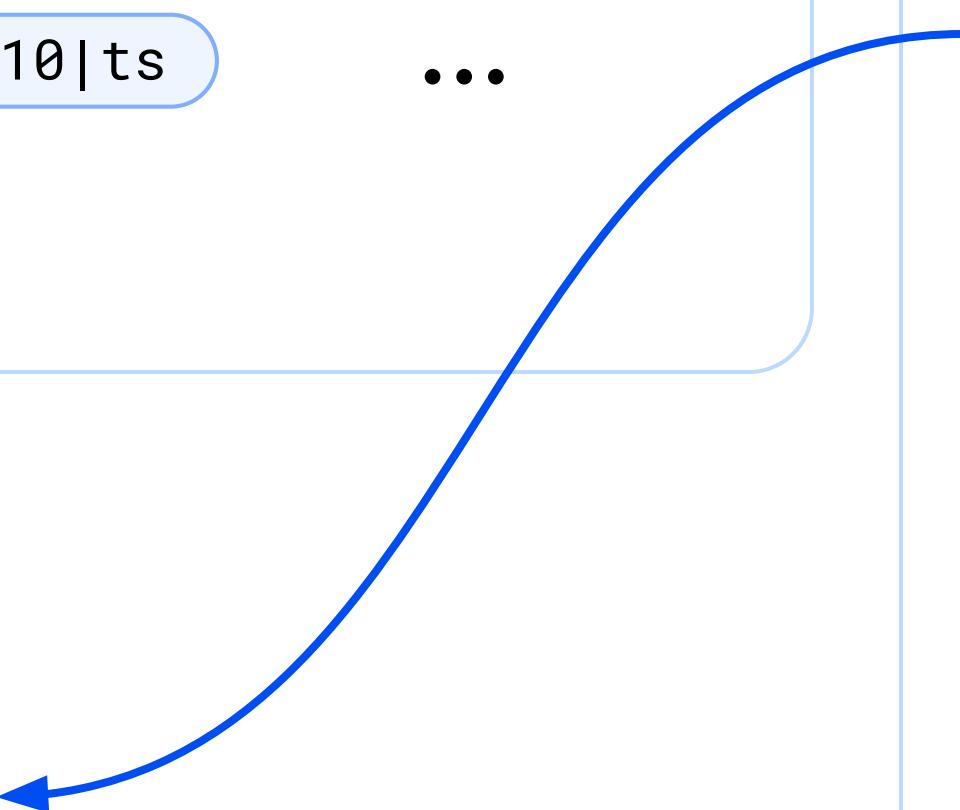
#id|10|ts

#id|10|ts

...

#id|10|ts

1 000 000
{ls} {ls} ... {ls}



RAM

#id

10|ts

...

10|ts

#id

10|ts

...

10|ts

...

#id

10|ts

...

10|ts

map

{ls} → #id ... {ls} → #id

#id 10|ts ... 10|ts

#id 10|ts ... 10|ts

...

#id 10|ts ... 10|ts

map {ls} → #id ... {ls} → #id



SSD

wal

{ls} → #id

#id|10|ts

...

{ls} → #id

#id|10|ts

#id|10|ts

#id|10|ts

#id|10|ts

...

#id|10|ts

{ls}



{ls}



...

{ls}



#id 10|ts ... 10|ts

#id 10|ts ... 10|ts

...

#id 10|ts ... 10|ts

map {ls} → #id ... {ls} → #id

#id 10|ts ... 10|ts

#id 10|ts ... 10|ts

...

#id 10|ts ... 10|ts

#id 10|ts ... 10|ts

#id 10|ts ... 10|ts

...

#id 10|ts ... 10|ts

map {ls} → #id ... {ls} → #id



RAM

#id

10|ts

...

10|ts

#id

10|ts

...

10|ts

...

#id

10|ts

...

10|ts

map

{ls} → #id ... {ls} → #id



SSD

wal

{ls} → #id

#id|10|ts

...

{ls} → #id

#id|10|ts

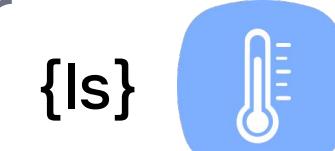
#id|10|ts

#id|10|ts

#id|10|ts

...

#id|10|ts



...



#id 10|ts ... 10|ts

#id 10|ts ... 10|ts

...

#id 10|ts ... 10|ts

map {ls} → #id ... {ls} → #id

#id 10|ts ... 10|ts

#id 10|ts ... 10|ts

...

#id 10|ts ... 10|ts

#id 10|ts ... 10|ts

#id 10|ts ... 10|ts

...

#id 10|ts ... 10|ts

map {ls} → #id ... {ls} → #id



RAM

#id

10|ts

...

10|ts

#id

10|ts

...

10|ts

...

#id

10|ts

...

10|ts

map

{ls} → #id ... {ls} → #id



SSD

wal

{ls} → #id

#id|10|ts

...

{ls} → #id

#id|10|ts

#id|10|ts

...

#id|10|ts

#id|10|ts

...

#id|10|ts

{ls} {ls} {ls} {ls}

#id|10|ts ... 10|ts

#id|10|ts ... 10|ts

...

#id|10|ts ... 10|ts

#id|10|ts ... 10|ts

...

#id|10|ts ... 10|ts

#id|10|ts ... 10|ts

...

map {ls} → #id ... {ls} → #id

map {ls} → #id ... {ls} → #id



RAM

#id

10|ts

...

10|ts

#id

10|ts

...

10|ts

...

#id

10|ts

...

10|ts

map {ls} → #id ... {ls} → #id

{ls} {

name: cpu_usage,
node: curiosity,
core: 4

}



SSD

wal

{ls} → #id

#id|10|ts

...

{ls} → #id

#id|10|ts

#id|10|ts

...

#id|10|ts

#id|10|ts

...

#id|10|ts

1 000 000 000

{ls}



{ls}



...

{ls}



#id|10|ts ... 10|ts

#id|10|ts ... 10|ts

...

#id|10|ts ... 10|ts

map {ls} → #id ... {ls} → #id



RAM

#id

10|ts

...

10|ts

#id

10|ts

...

10|ts

...

#id

10|ts

...

10|ts

map

{ls} → #id ... {ls} → #id

{ls} {

name: cpu_usage,
node: curiosity,
core: 4



SSD

wal

{ls} → #id

#id|10|ts

...

{ls} → #id

#id|10|ts

#id|10|ts

#id|10|ts

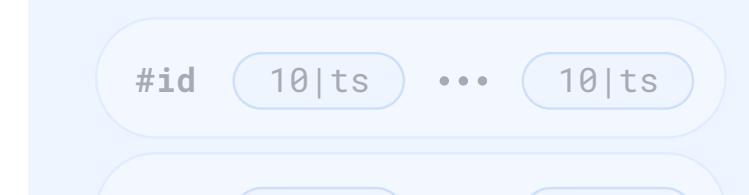
#id|10|ts

...

#id|10|ts



...



...



map {ls} → #id ... {ls} → #id



RAM

#id

10|ts

...

10|ts

#id

10|ts

...

10|ts

...

#id

10|ts

...

10|ts

map {ls} → #id ... {ls} → #id

{ls} {

name: cpu_usage,
node: curiosity,
core: 4

}



SSD

wal

{ls} → #id

#id|10|ts

...

{ls} → #id

#id|10|ts

#id|10|ts

...

#id|10|ts

#id|10|ts

...

{ls}



1 000 000 000

{ls}



...

{ls}



#id

10|ts

...

10|ts

#id

10|ts

...

10|ts

#id

10|ts

...

10|ts

map

{ls} → #id ... {ls} → #id

#id

10|ts

...

10|ts

#id

10|ts

...

10|ts

#id

10|ts

...

10|ts

map

{ls} → #id ... {ls} → #id



RAM

#id

10|ts

...

10|ts

#id

10|ts

...

10|ts

#id

10|ts

...

10|ts

map

{ls} → #id ... {ls} → #id

{ls} {

name: cpu_usage,
node: curiosity,
core: 4

}



SSD

wal

#id|10|ts

#id|10|ts

#id|10|ts

1 000 000
{ls}

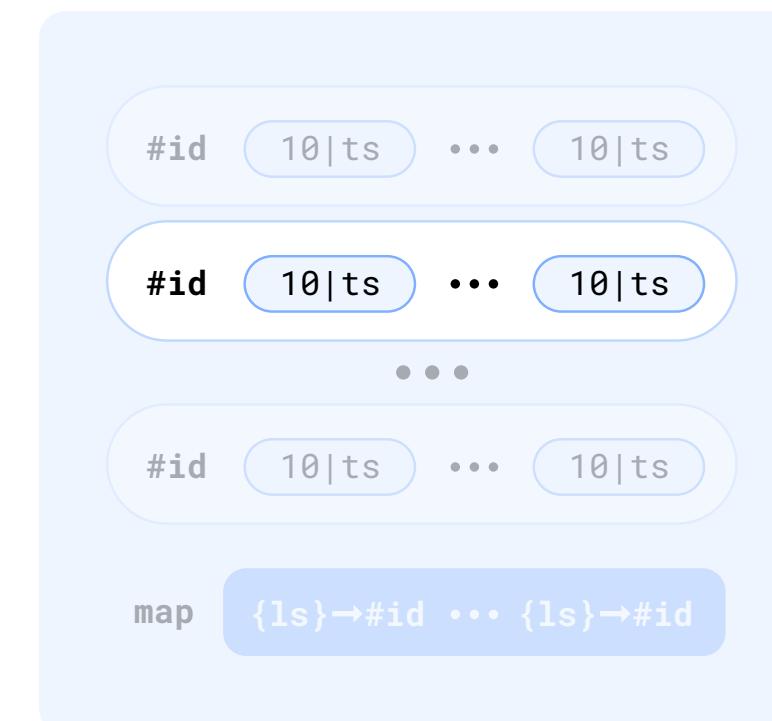


1 000 000
{ls}



...

1 000 000
{ls}



RAM

#id

10|ts

...

10|ts

#id

10|ts

...

10|ts

map

{ls} → #id ... {ls} → #id

{ls} {
name: cpu_usage,
node: curiosity,
core: 4
}



SSD

wal

#id|10|ts

#id|10|ts

#id|10|ts

{ls}



1 000 000

{ls}



...

{ls}



#id 10|ts ... 10|ts

10|ts

10|ts

10|ts

10|ts

10|ts

10|ts

10|ts

10|ts

{ls} {
name: cpu_usage,
node: curiosity,
core: 4
}



RAM

#id

10|ts

...

10|ts

#id

10|ts

...

10|ts

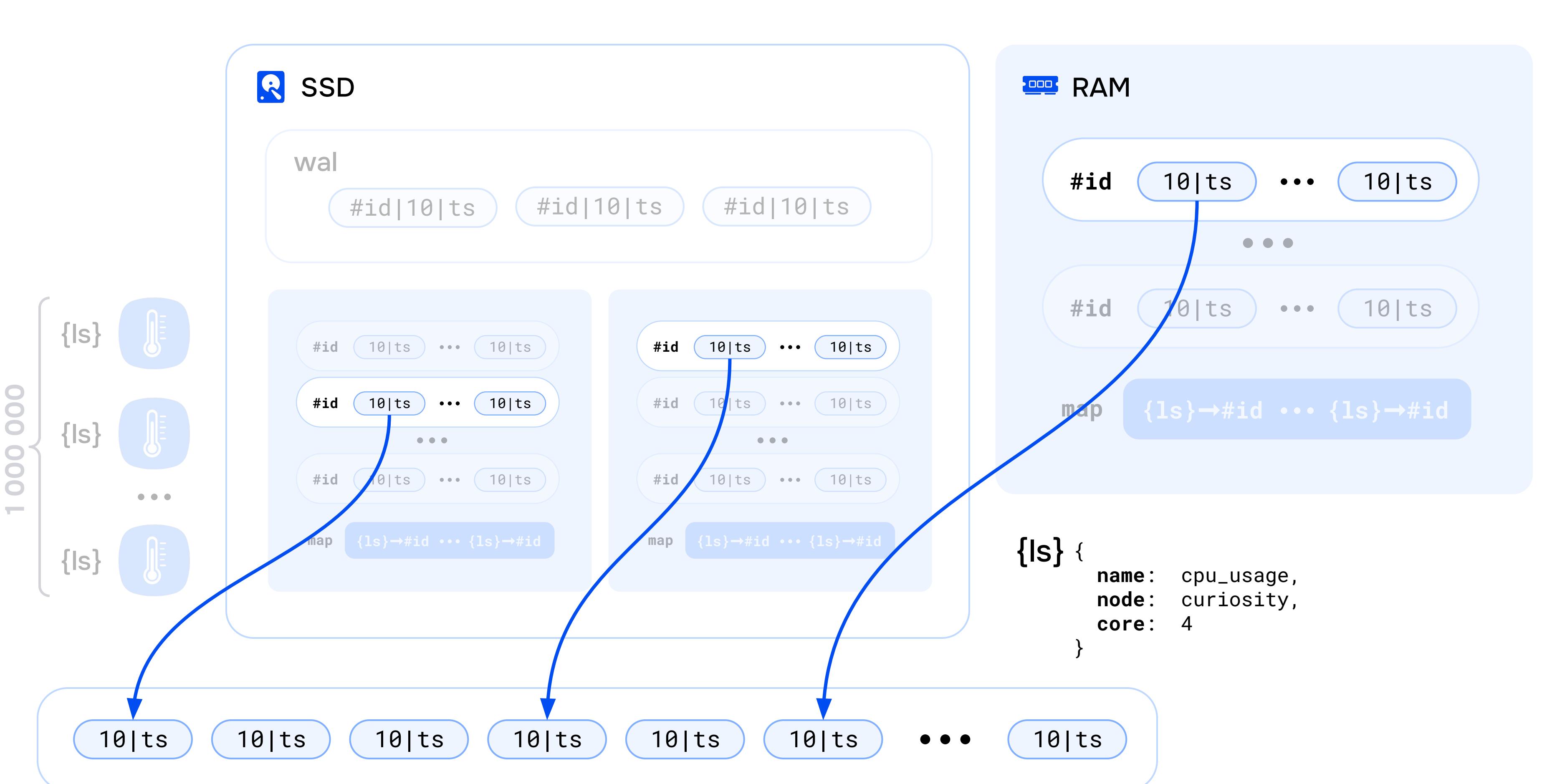
map

{ls}→#id ... {ls}→#id

map

{ls}→#id ... {ls}→#id

map {ls}→#id ... {ls}→#id





SSD

wal

#id|10|ts

#id|10|ts

#id|10|ts

{ls}



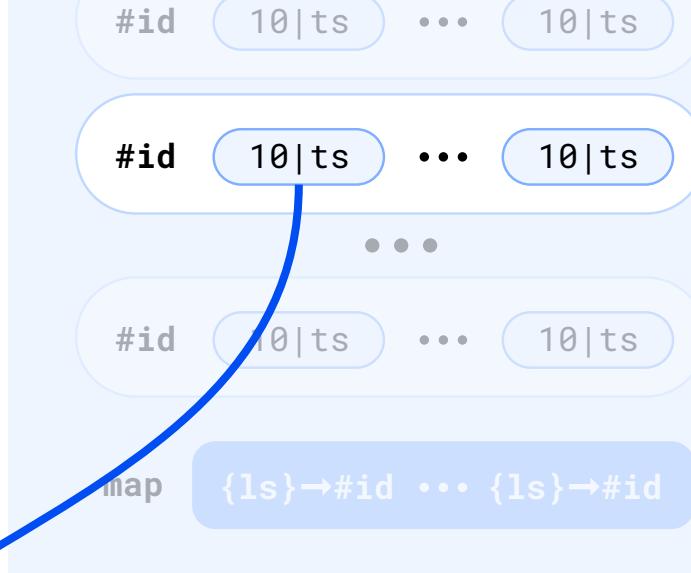
1 000 000

{ls}



...

{ls}



RAM

#id

10|ts

...

10|ts

#id

10|ts

...

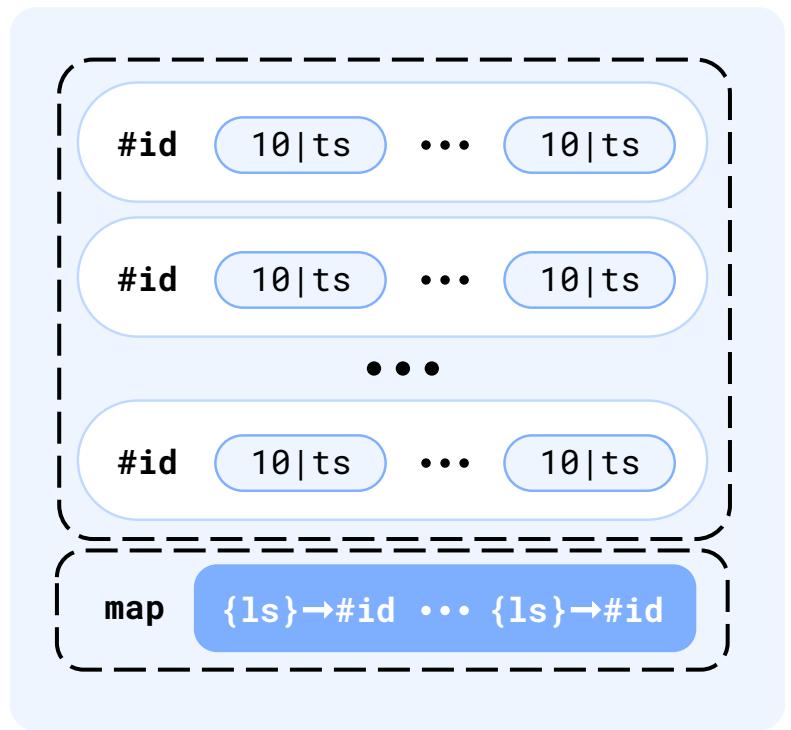
10|ts

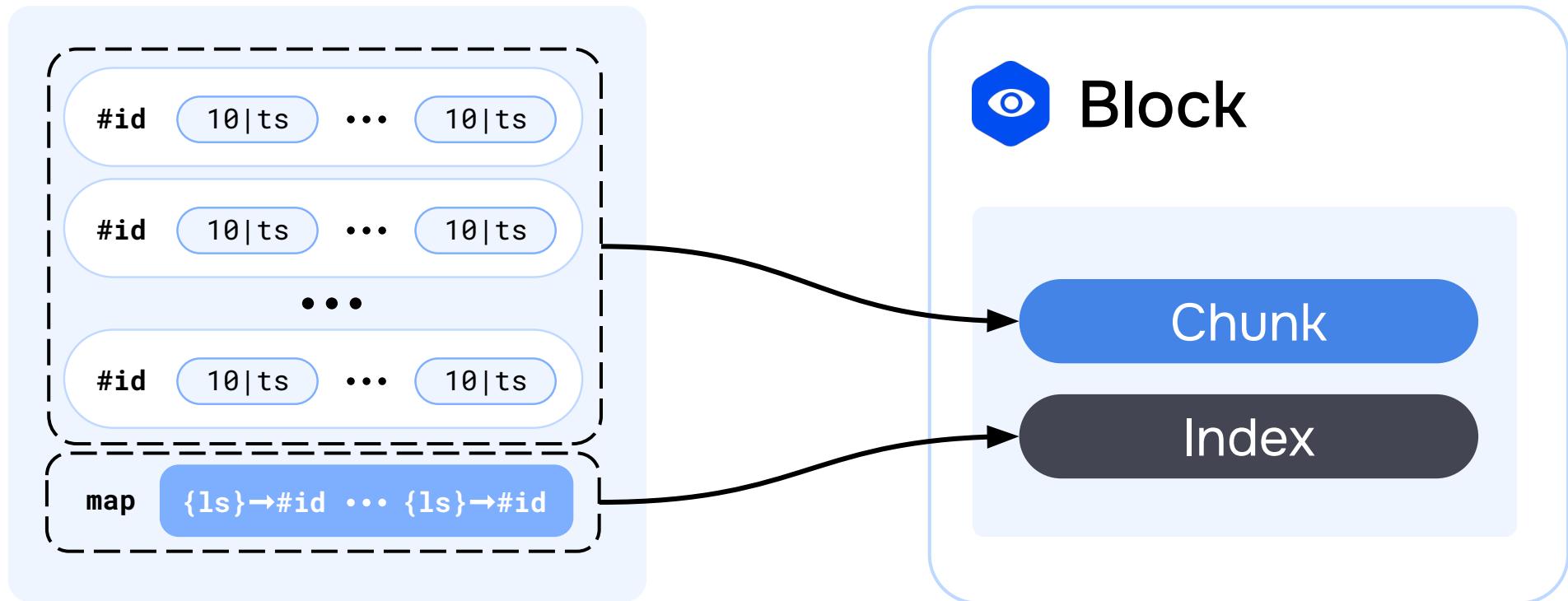
map

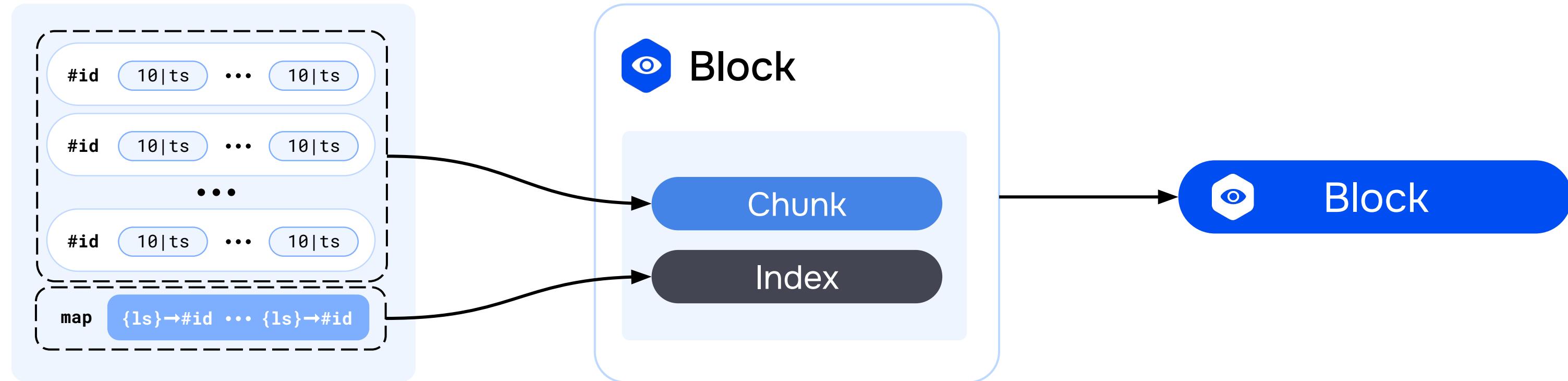
{ls} → #id ... {ls} → #id

{ls}

{
name: cpu_usage,
node: curiosity,
core: 4
}



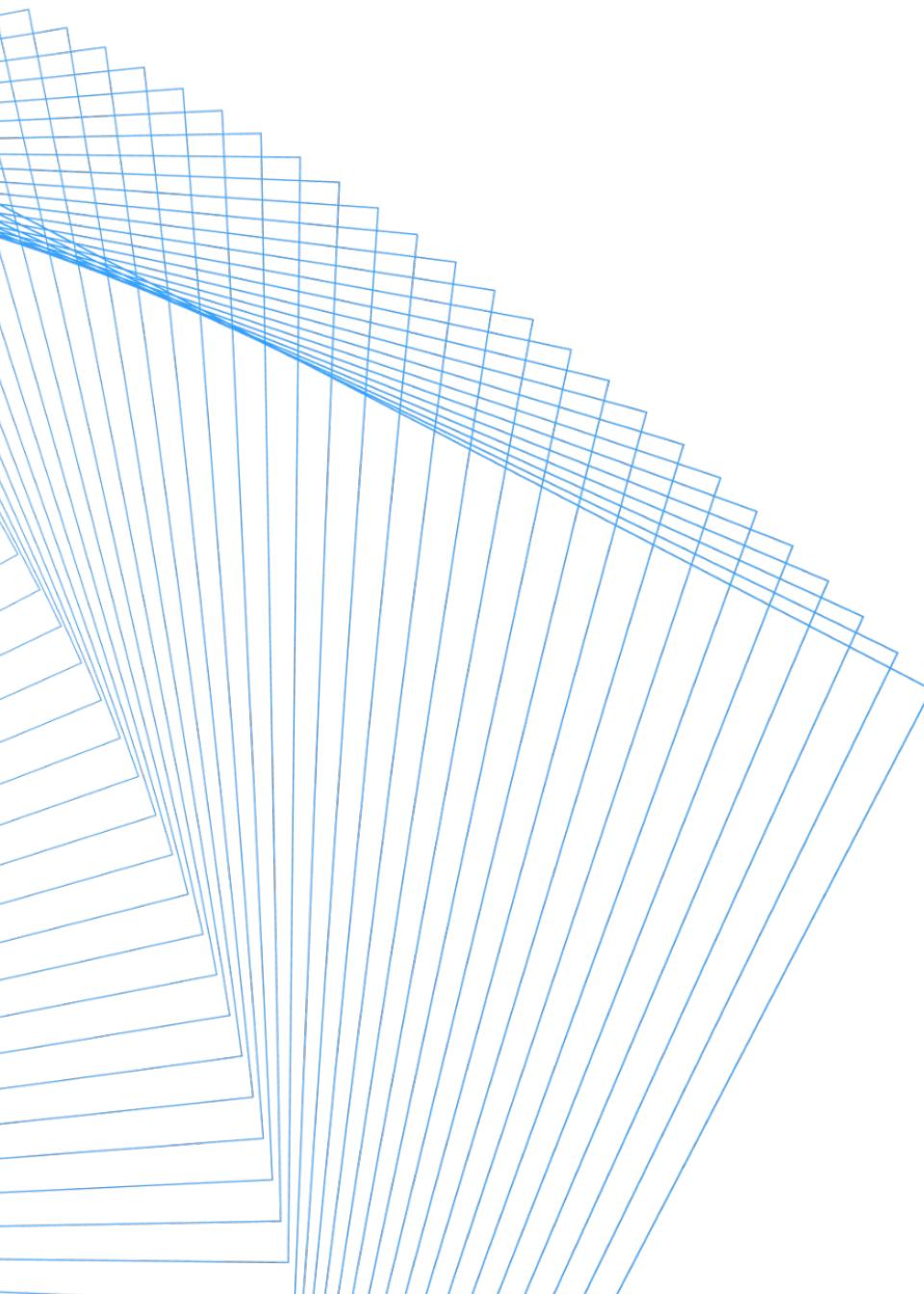




Важно

01

Блок не меняется после записи на диск



Важно

01

Блок не меняется после записи на диск

02

Блок состоит из двух частей:

- Index – содержит данные о LabelSet
- Chunk – содержит значения метрик

Важно

01

Блок не меняется после записи на диск

02

Блок состоит из двух частей:

- Index – содержит данные о LabelSet
- Chunk – содержит значения метрик

03

Index – примерно 10 % от размера
всего блока

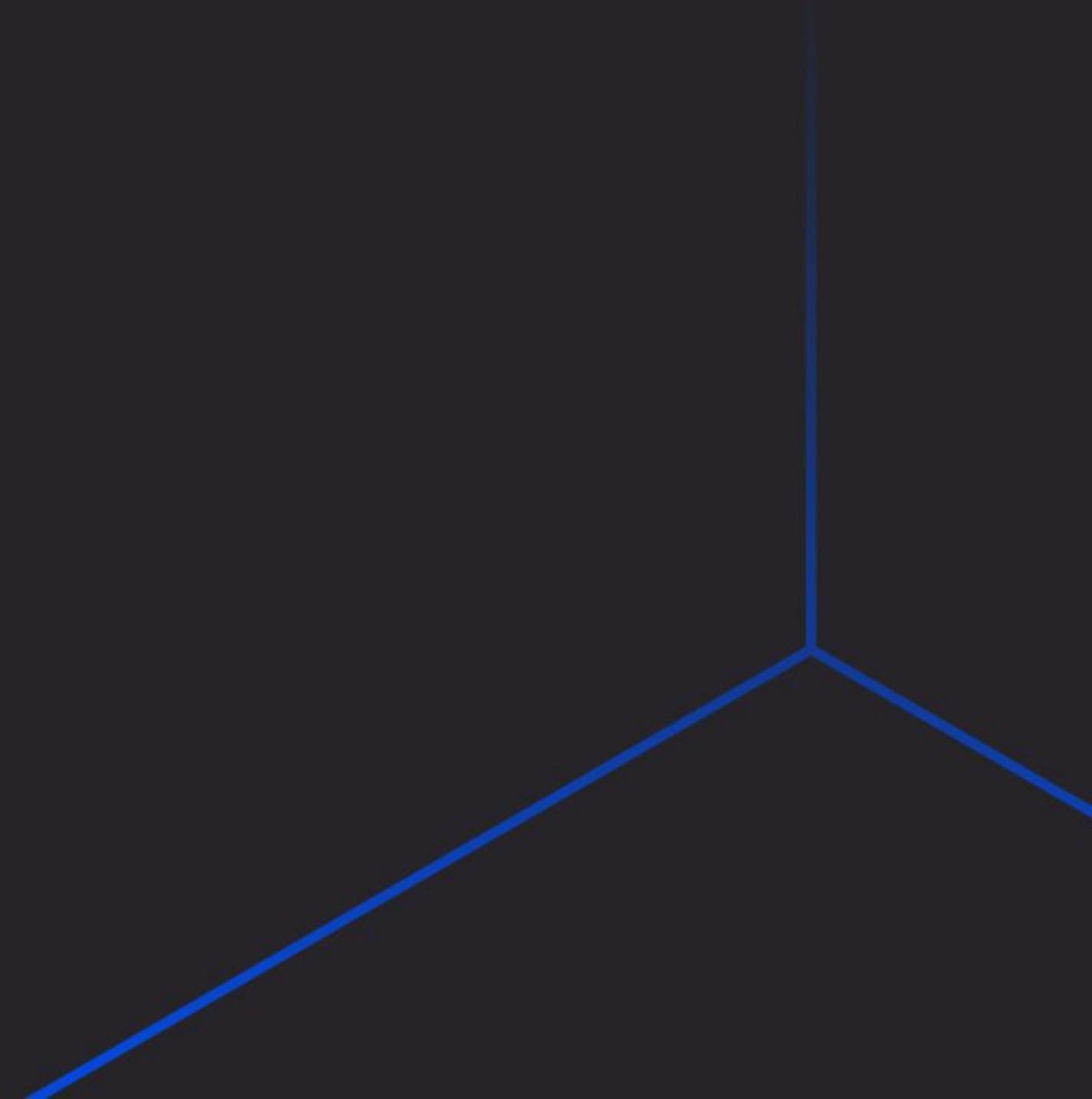
Мониторинг

Статья «Мимо тёщиного
дома я без метрик не хожу»

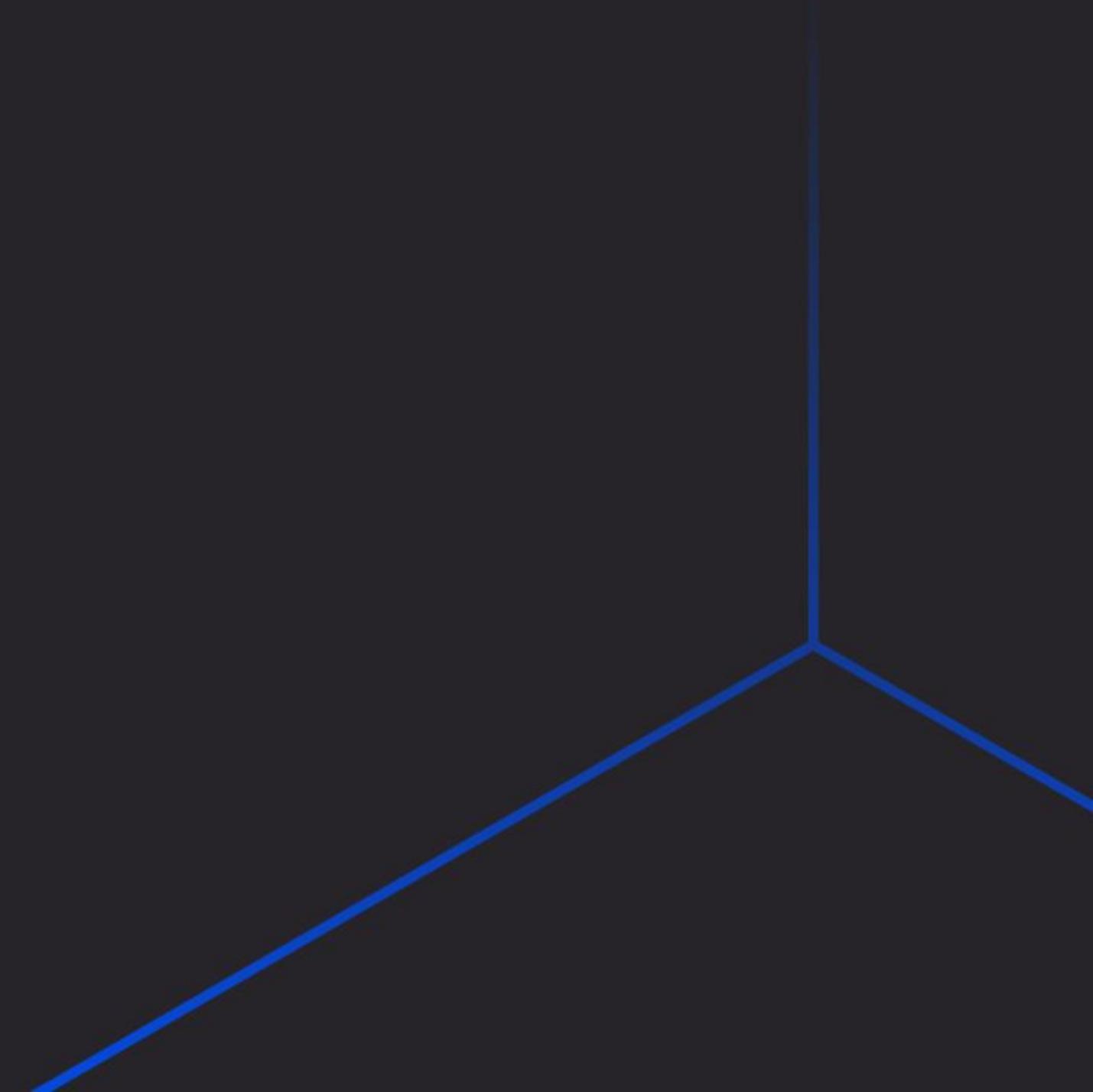


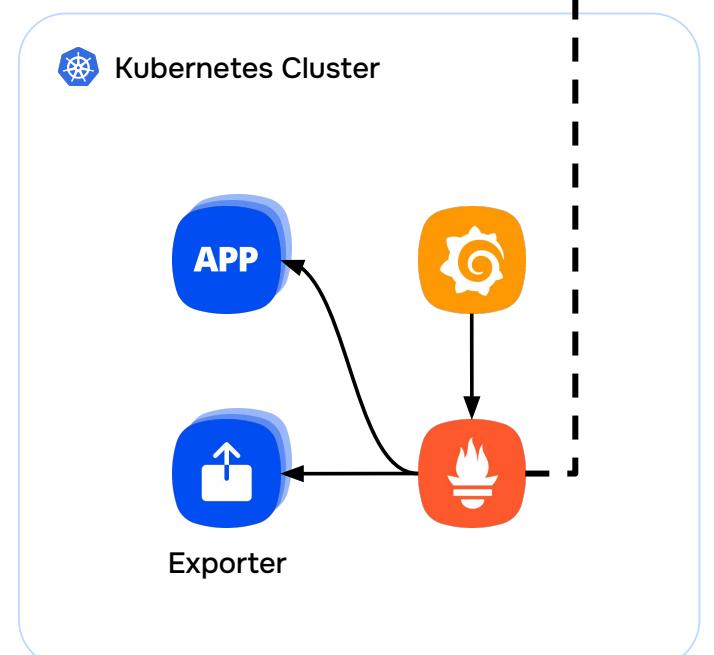
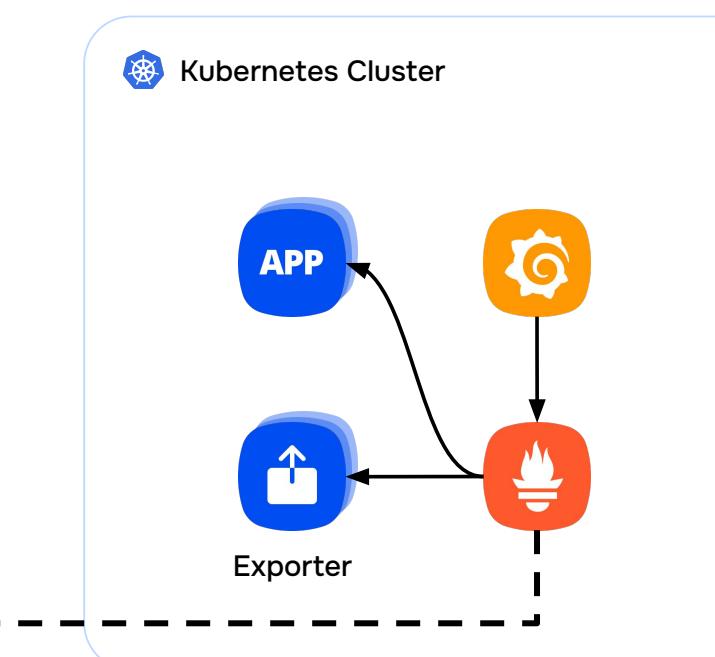
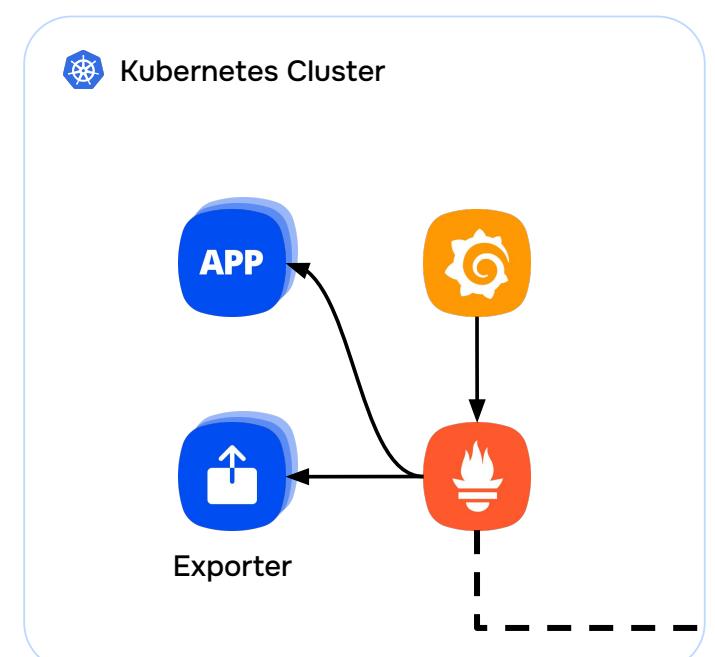
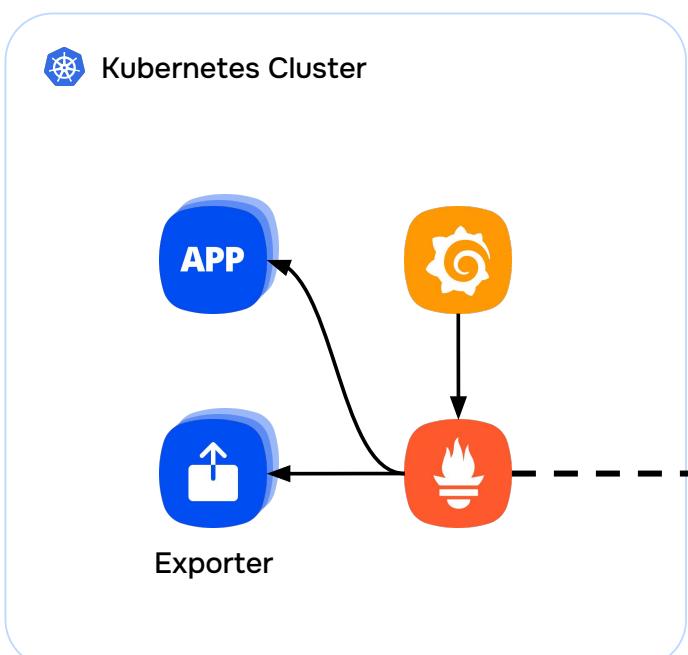
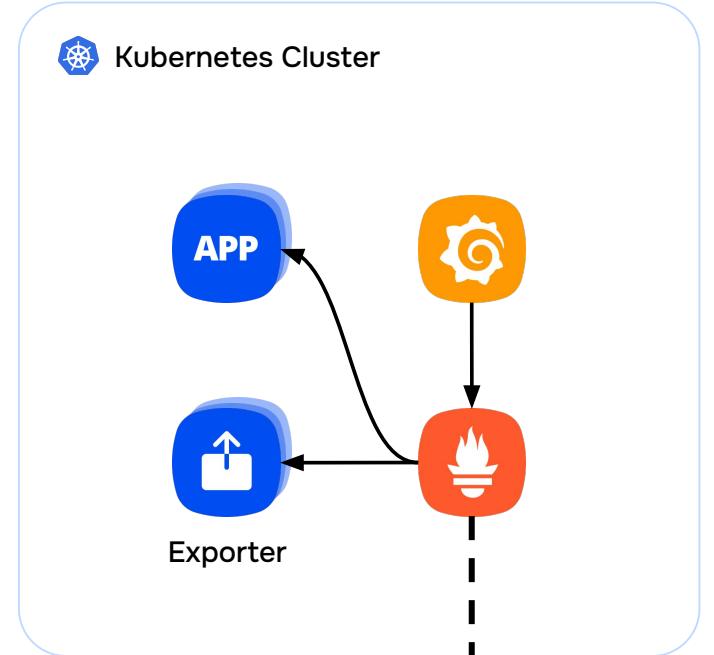
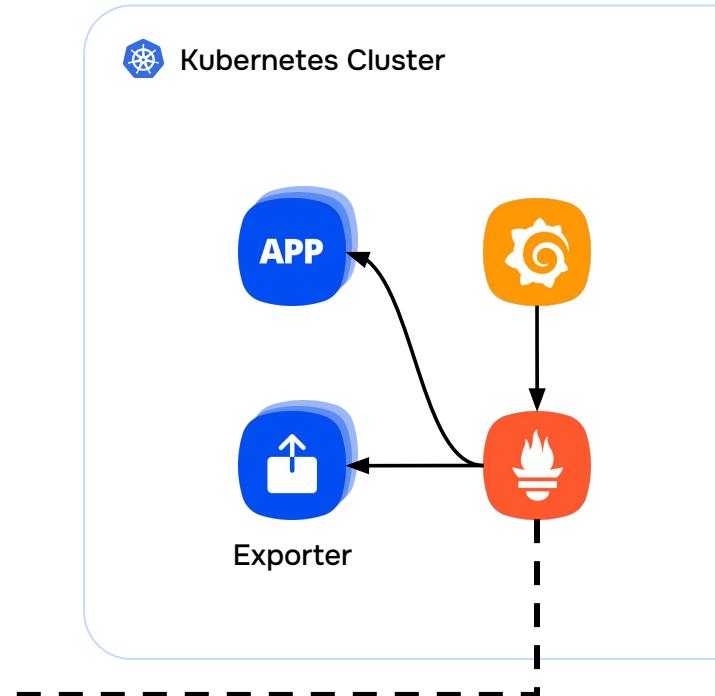
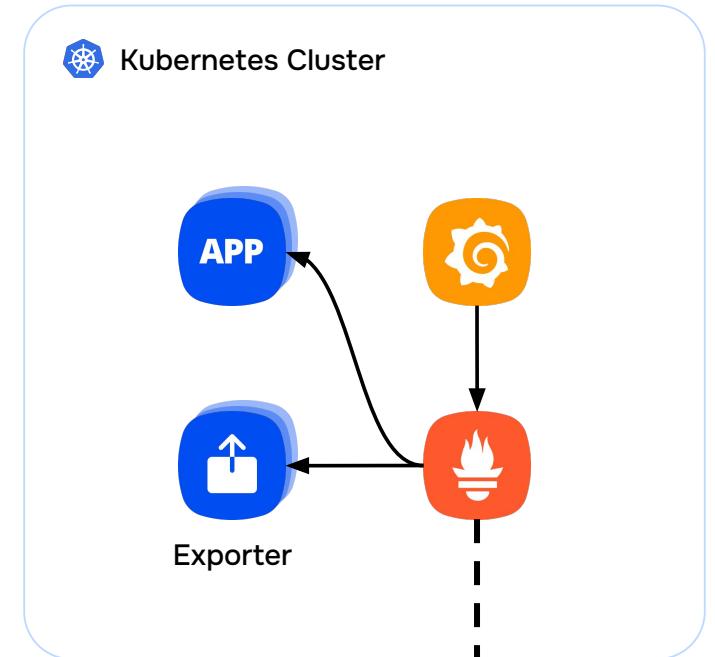
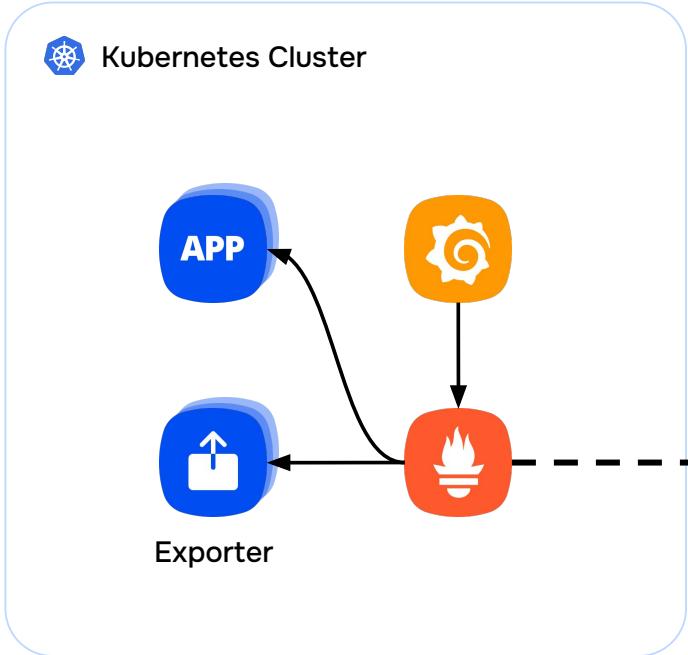
Подробнее

Запись метрик



СКОЛЬКО НУЖНО ДЛЯ ОТКАЗОУСТОЙЧИВОСТИ?





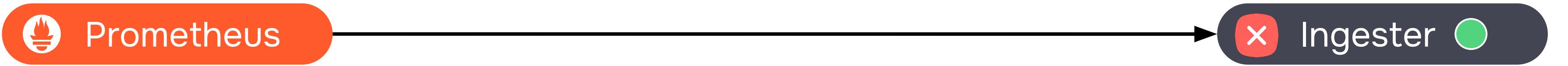
Централизованное
хранилище метрик



Prometheus









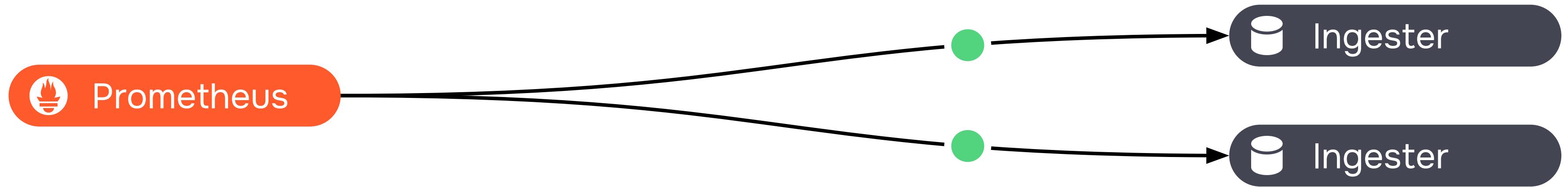
Prometheus

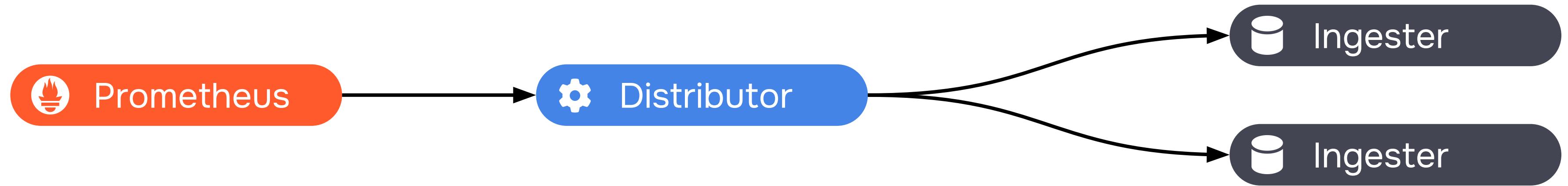


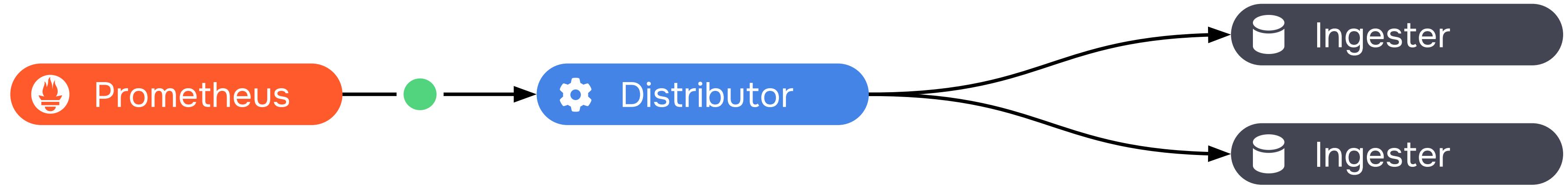
Ingestor

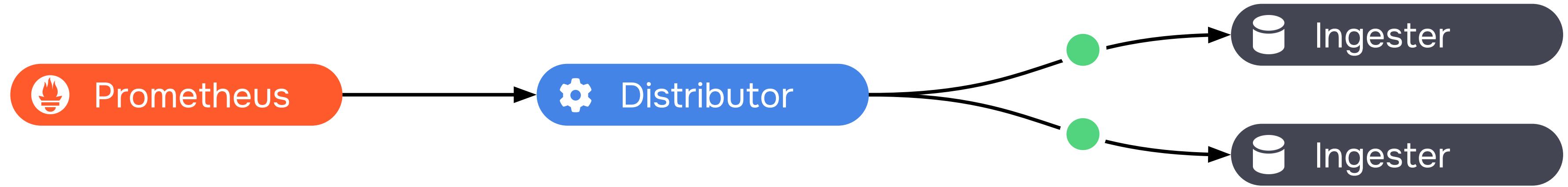


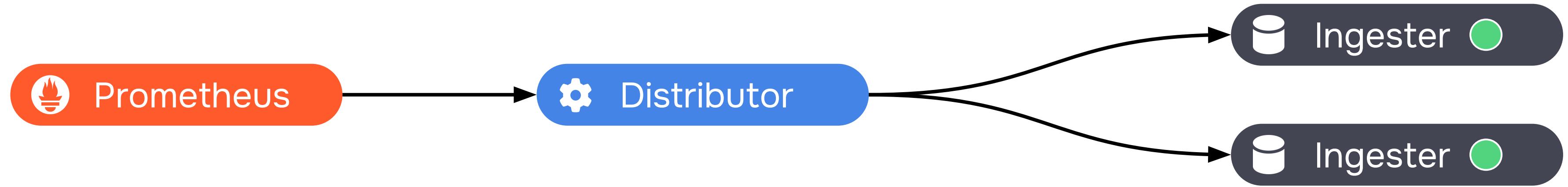
Ingestor

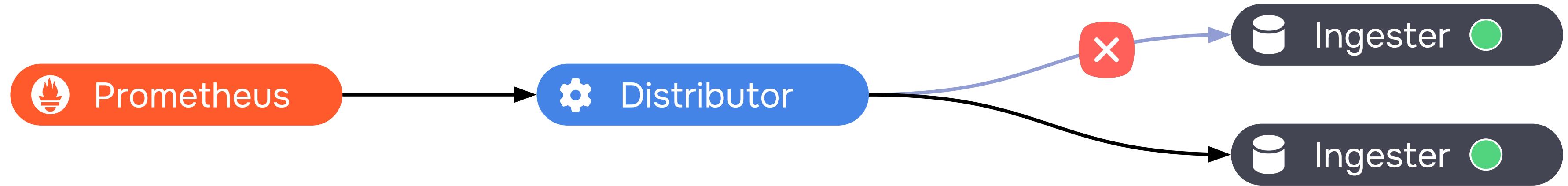


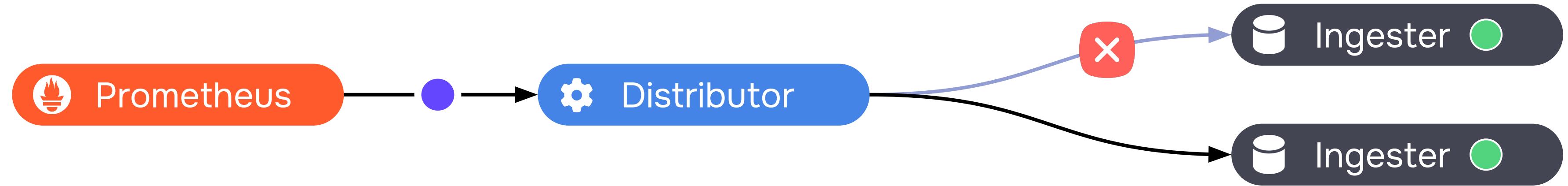


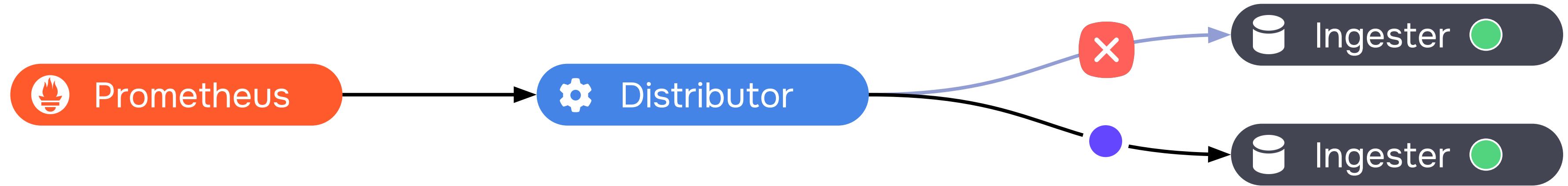


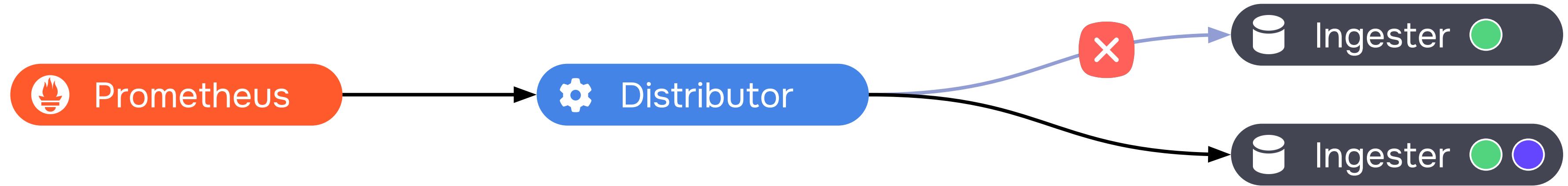


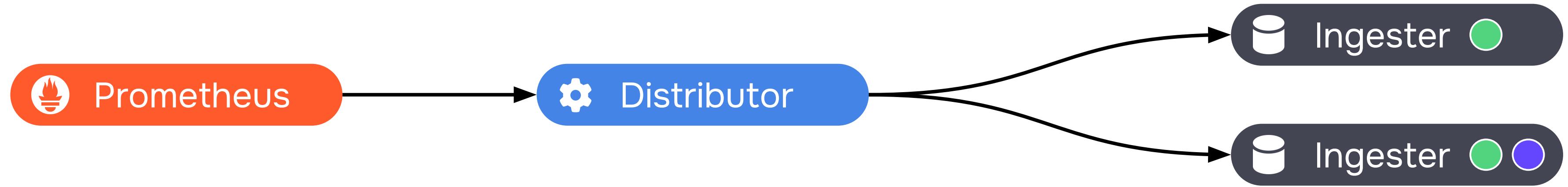


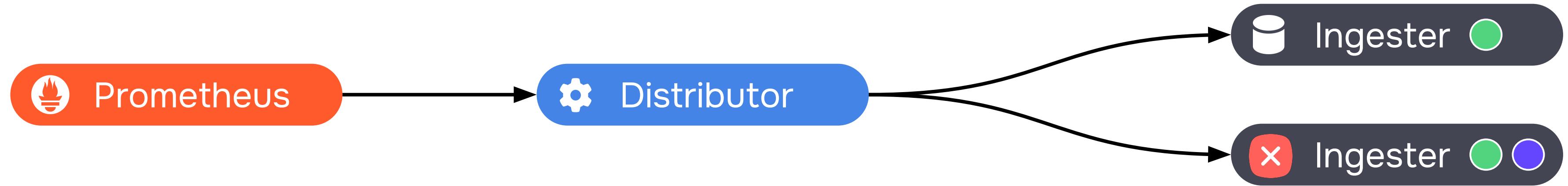


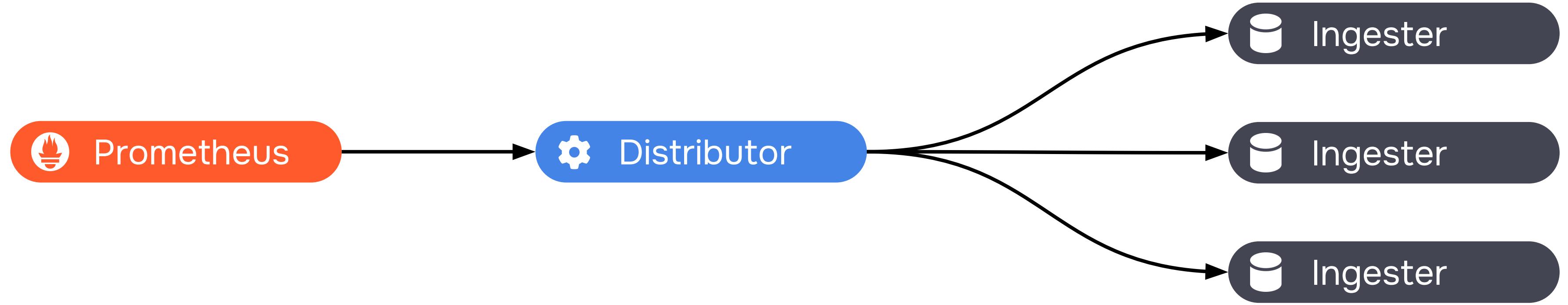


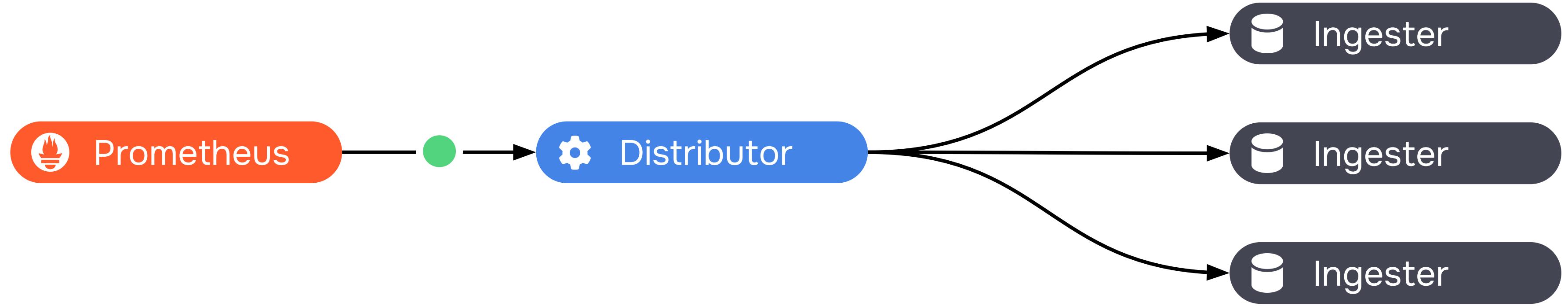


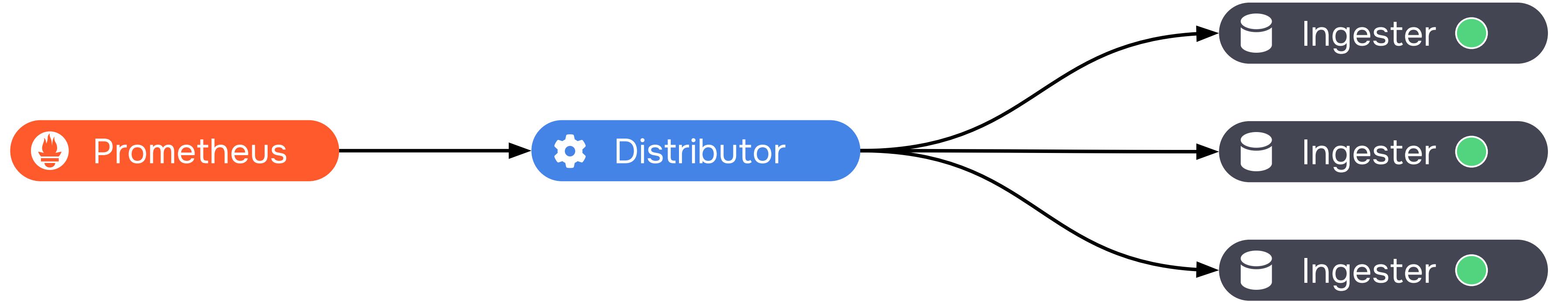


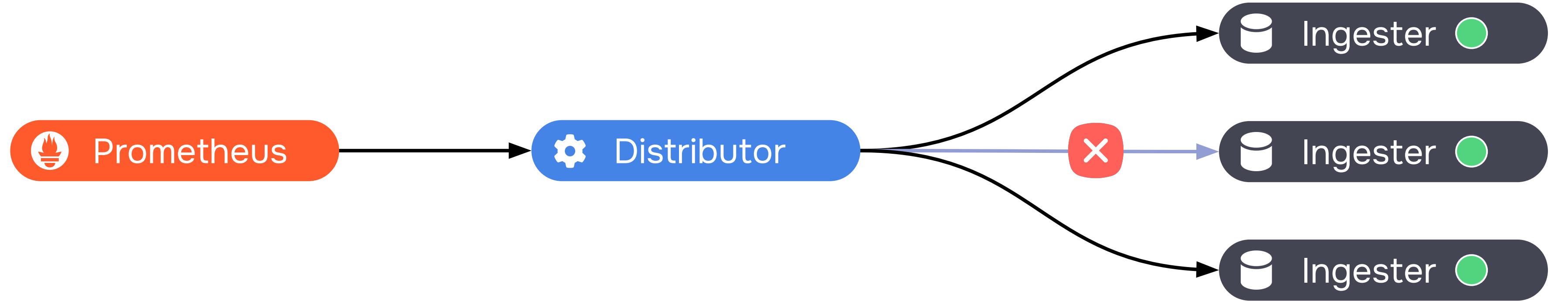


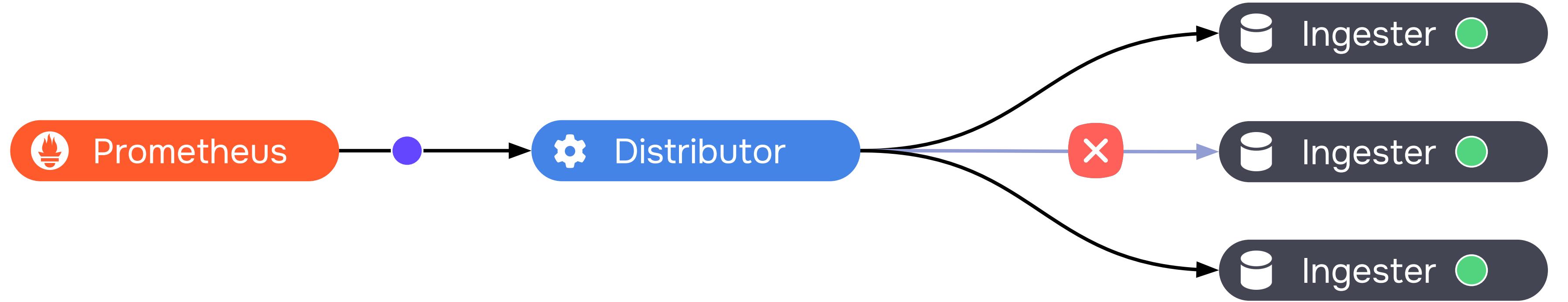


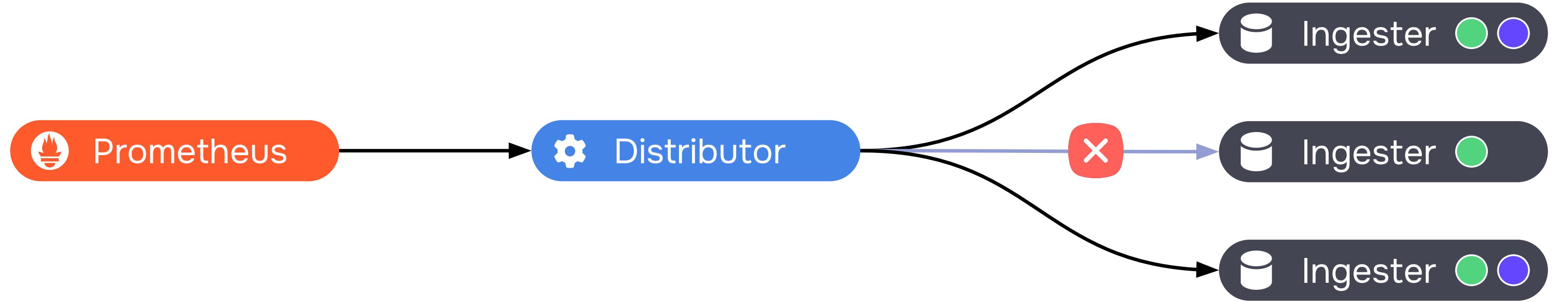


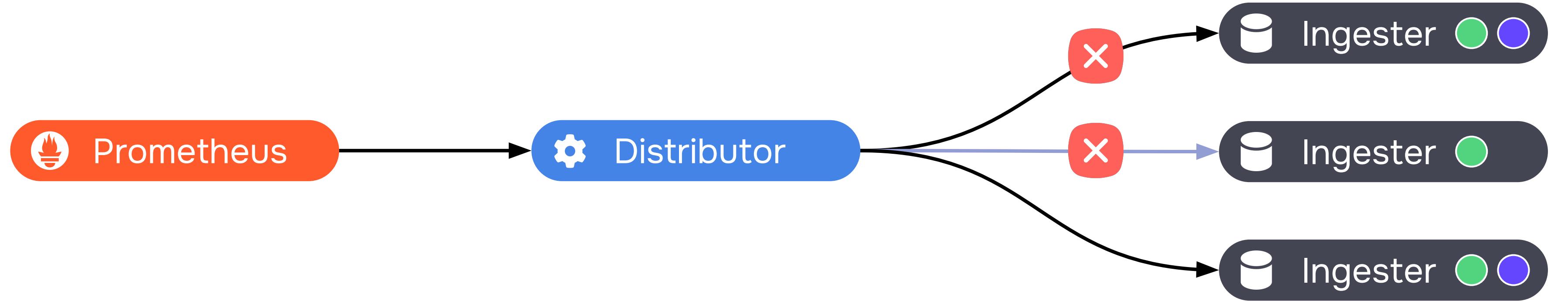


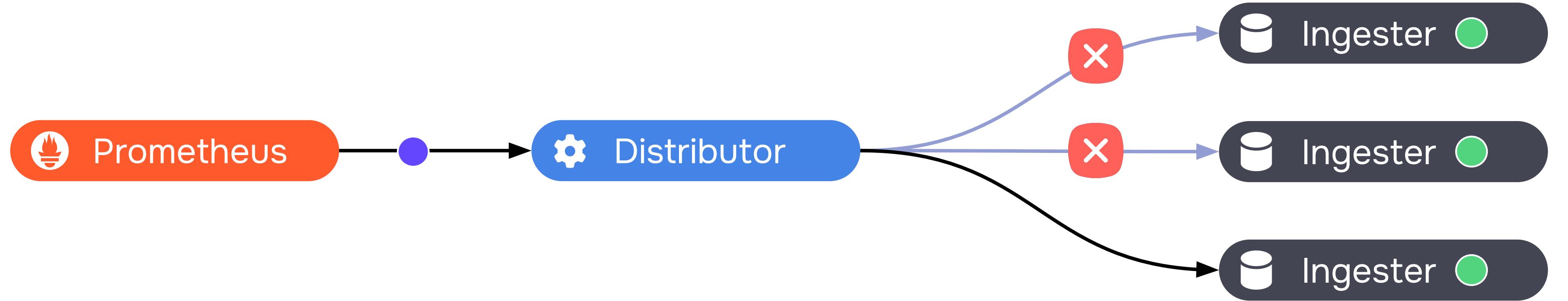


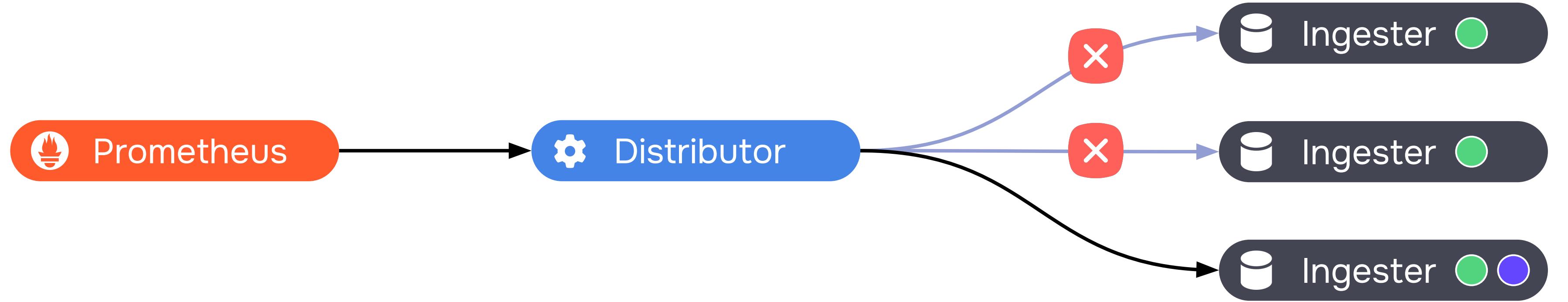


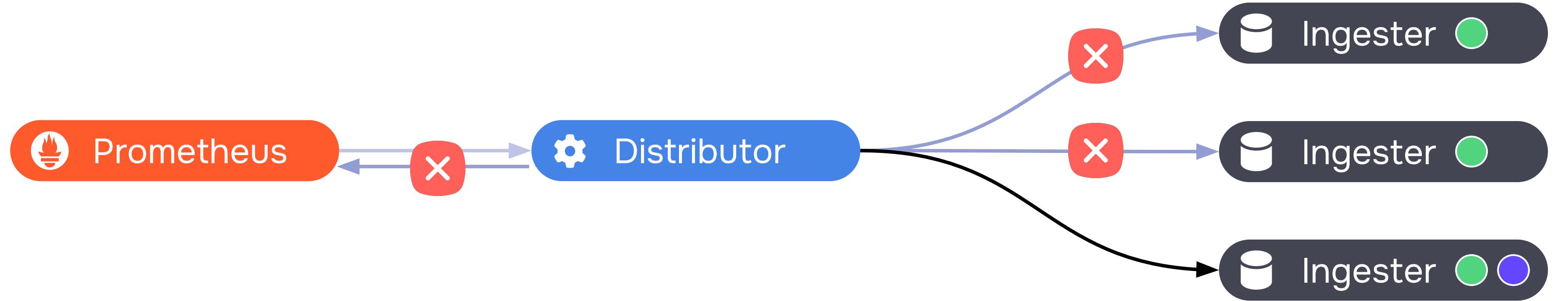


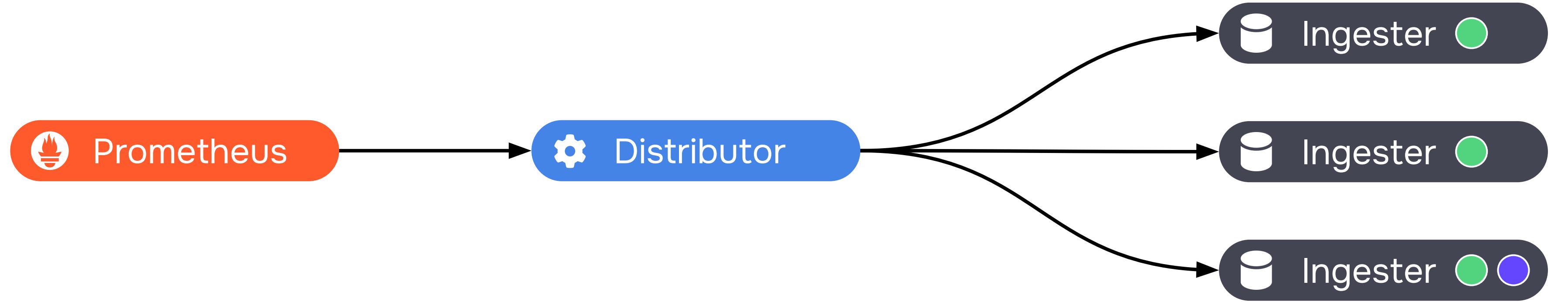


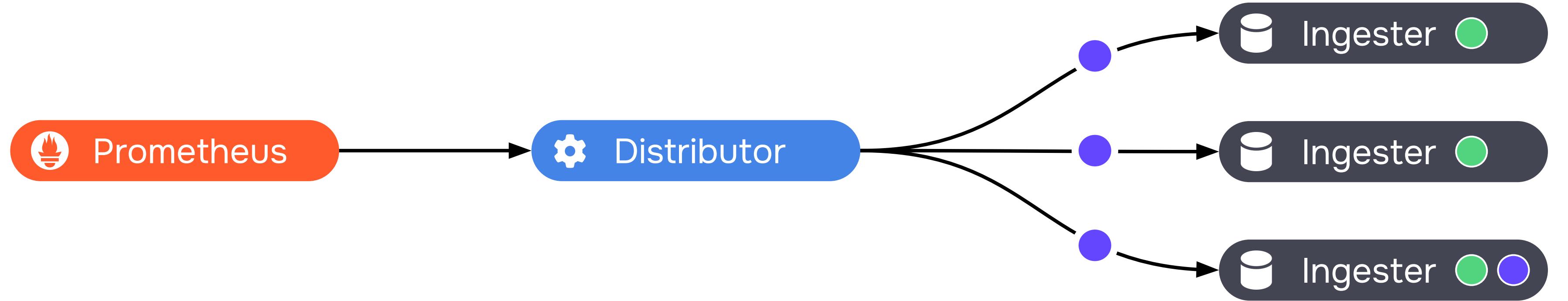


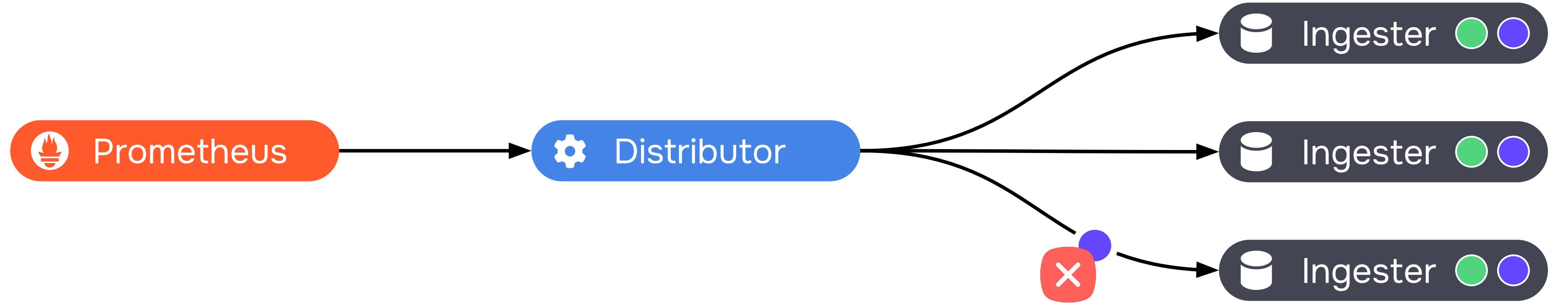


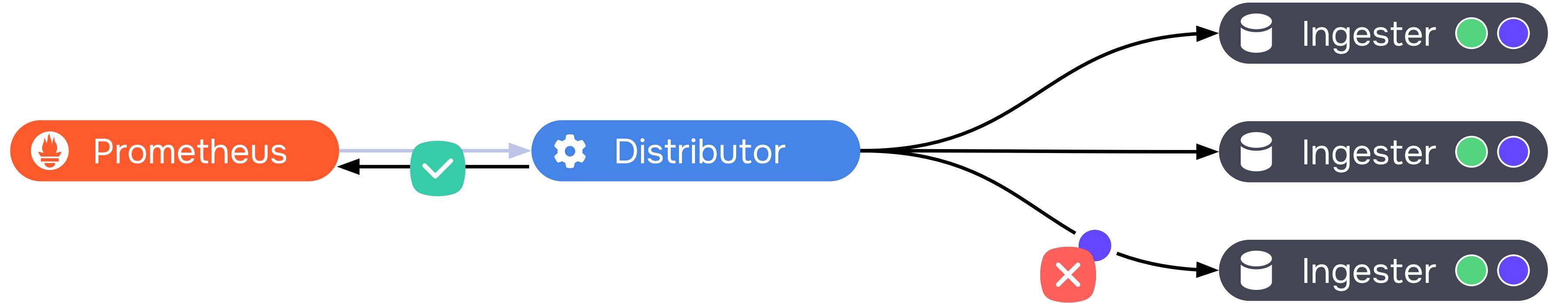






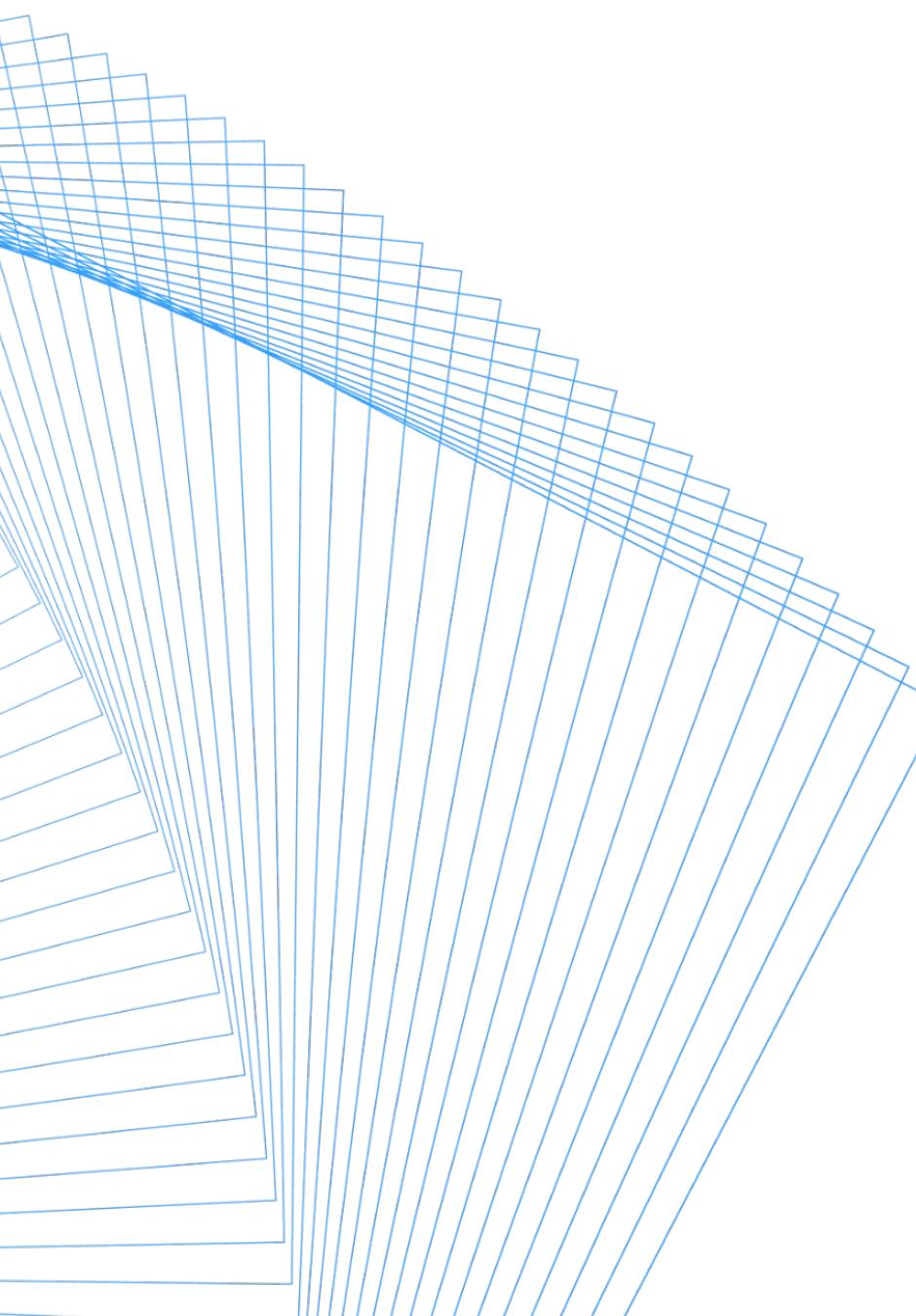






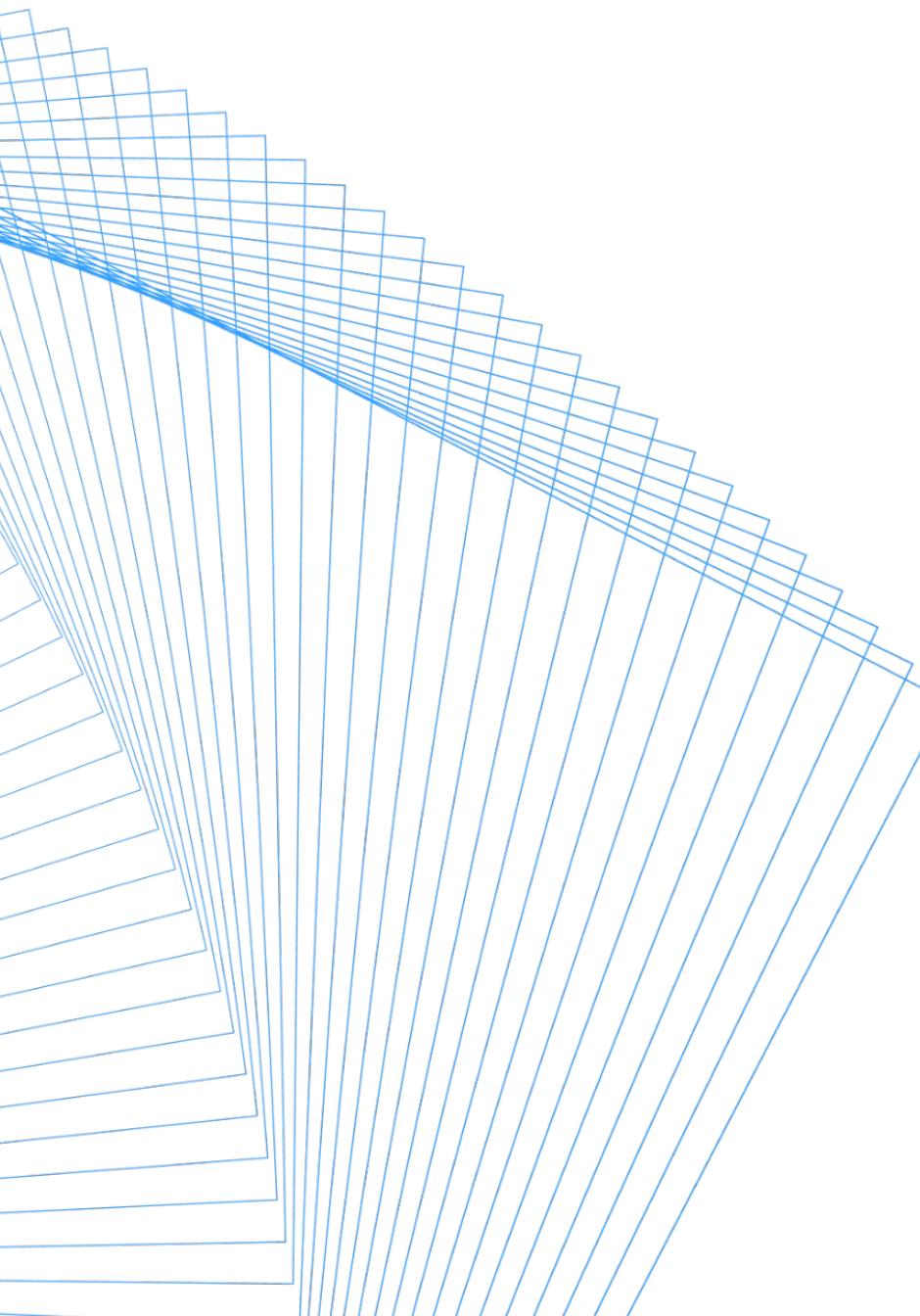
Централизованная система мониторинга

- Отказоустойчивость на запись: $WC = RF/2 + 1$

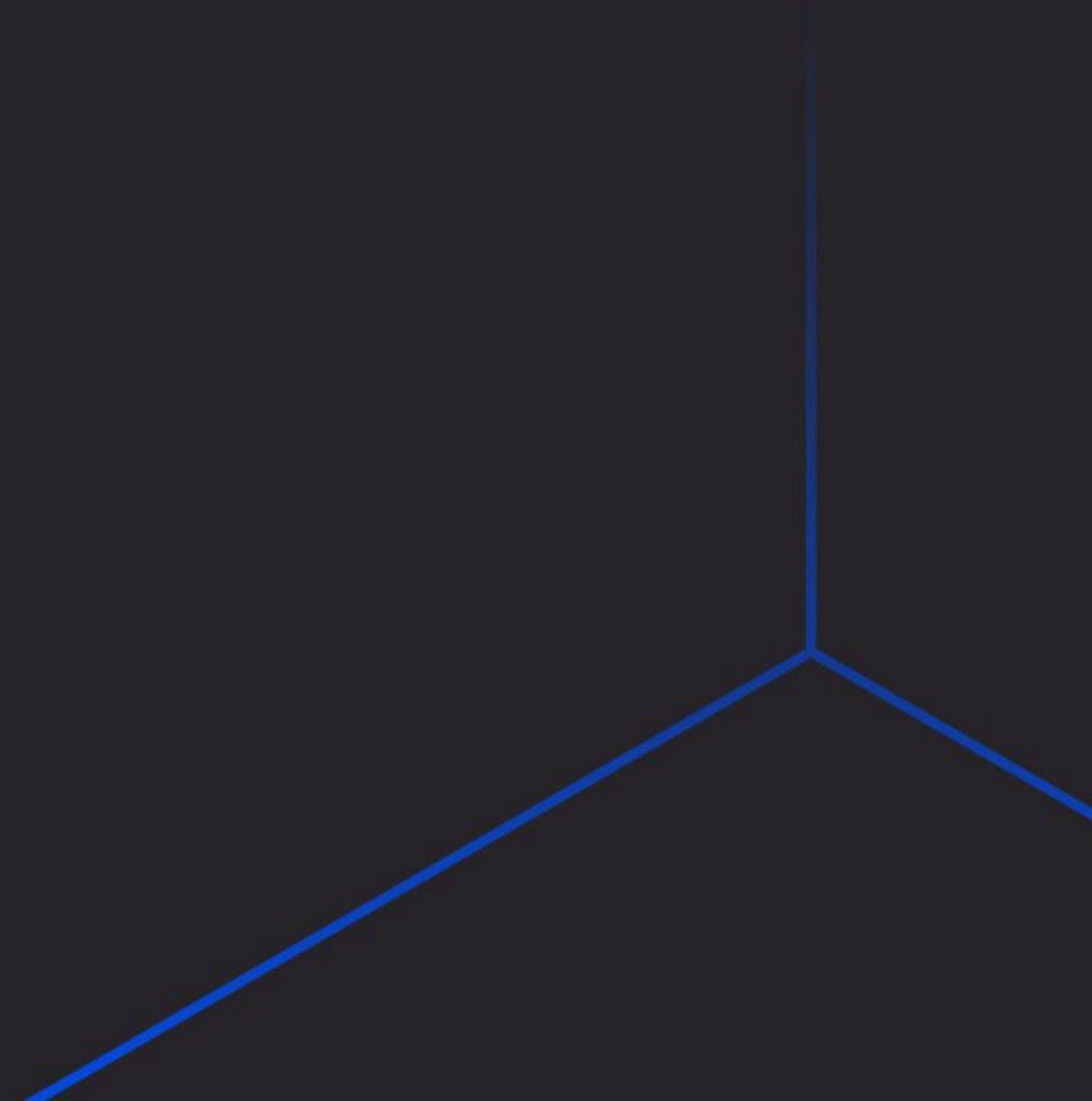


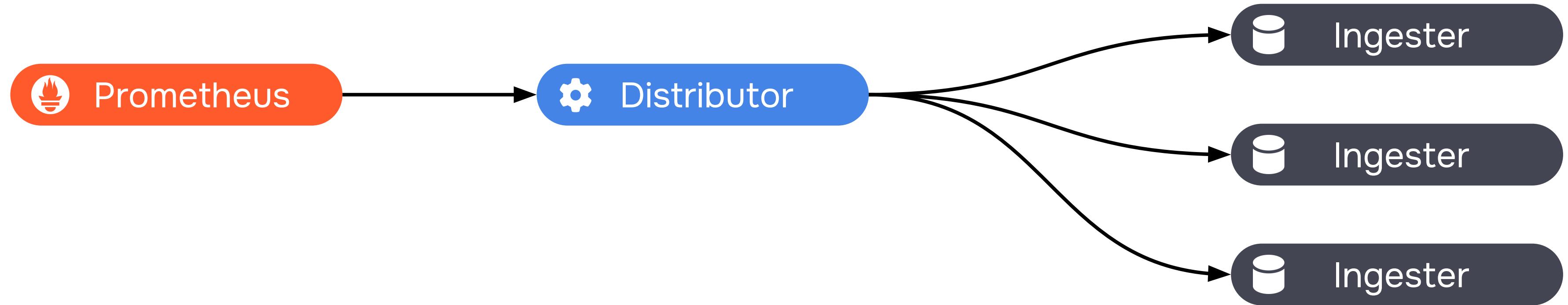
Централизованная система мониторинга

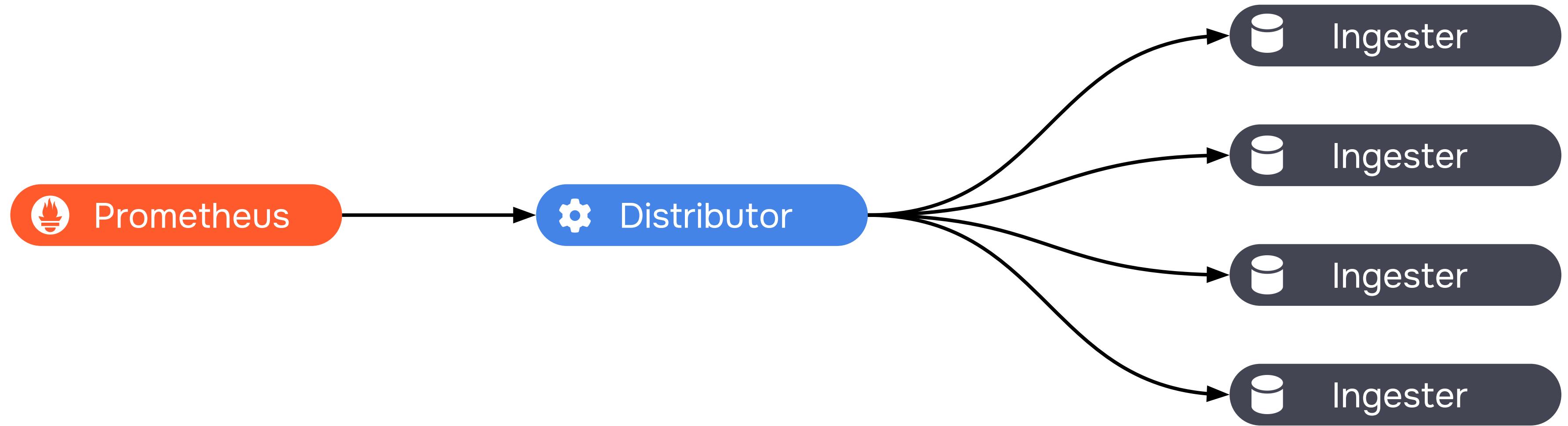
- Отказоустойчивость на запись: $WC = RF/2 + 1$
- Отказоустойчивость на чтение: $RC = RF - WC + 1$

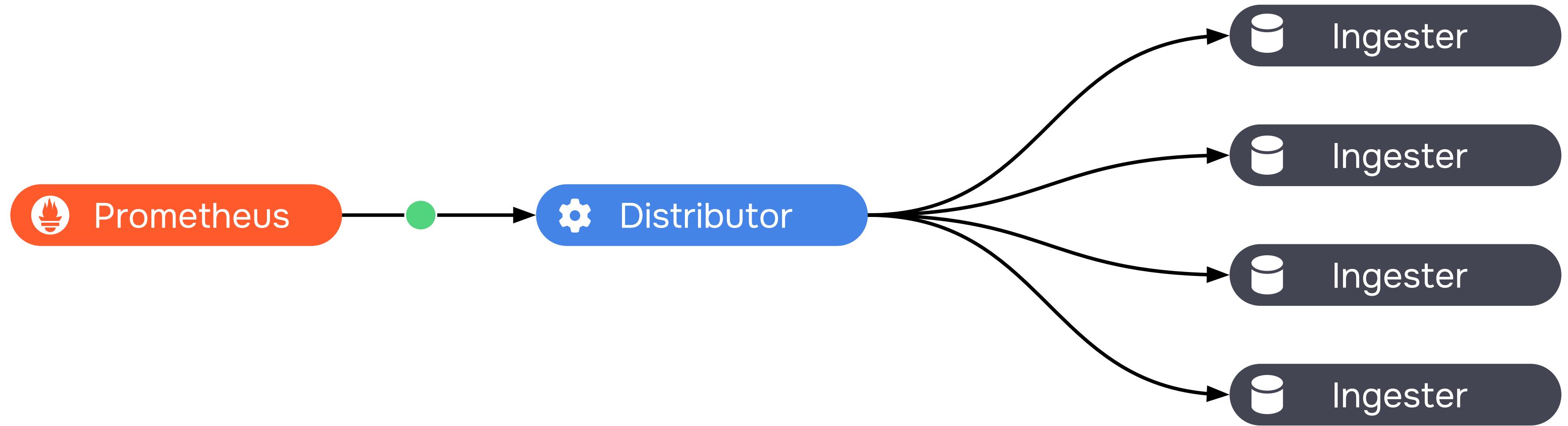


А как работает шардирование?



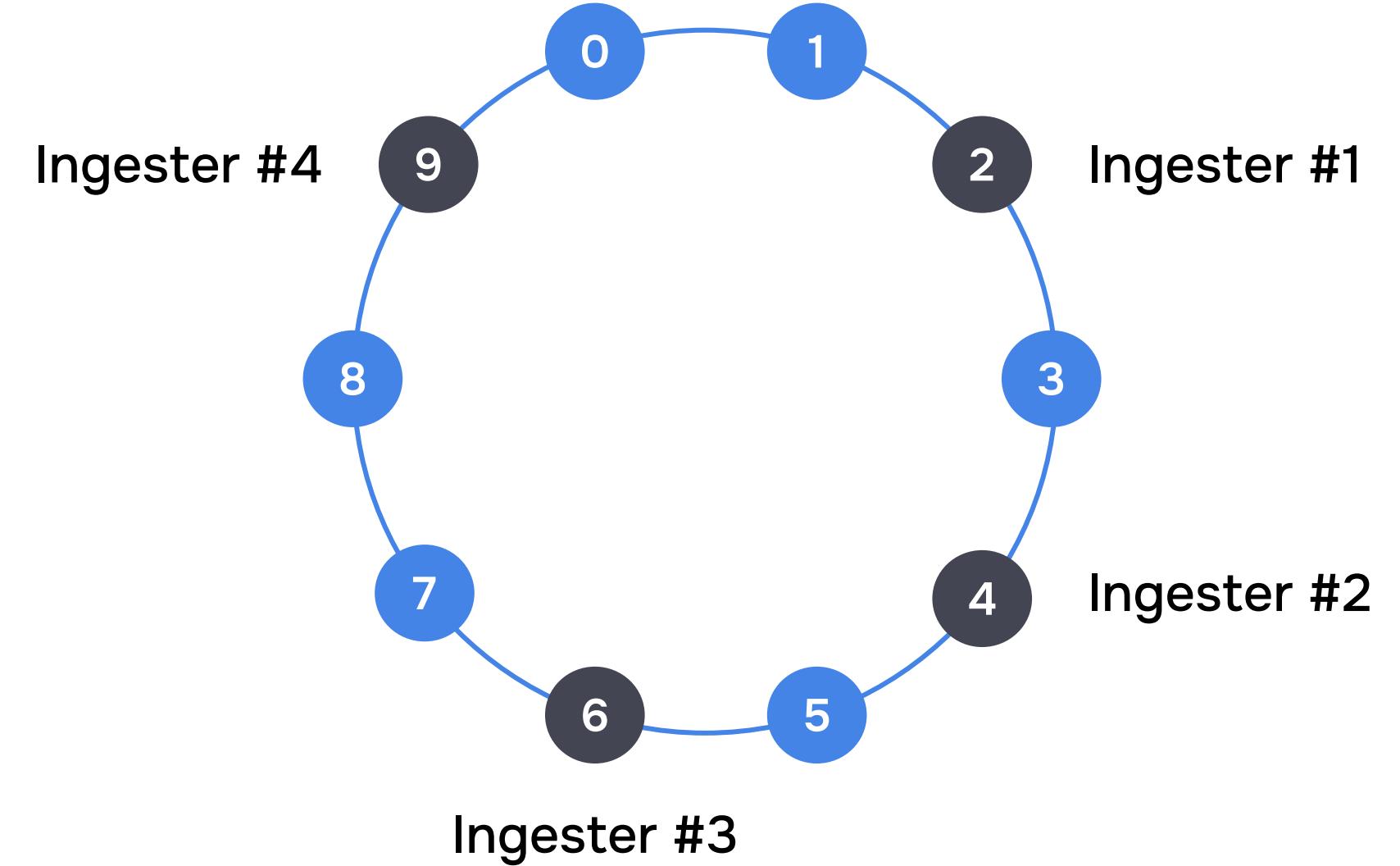






Хеш-ринг

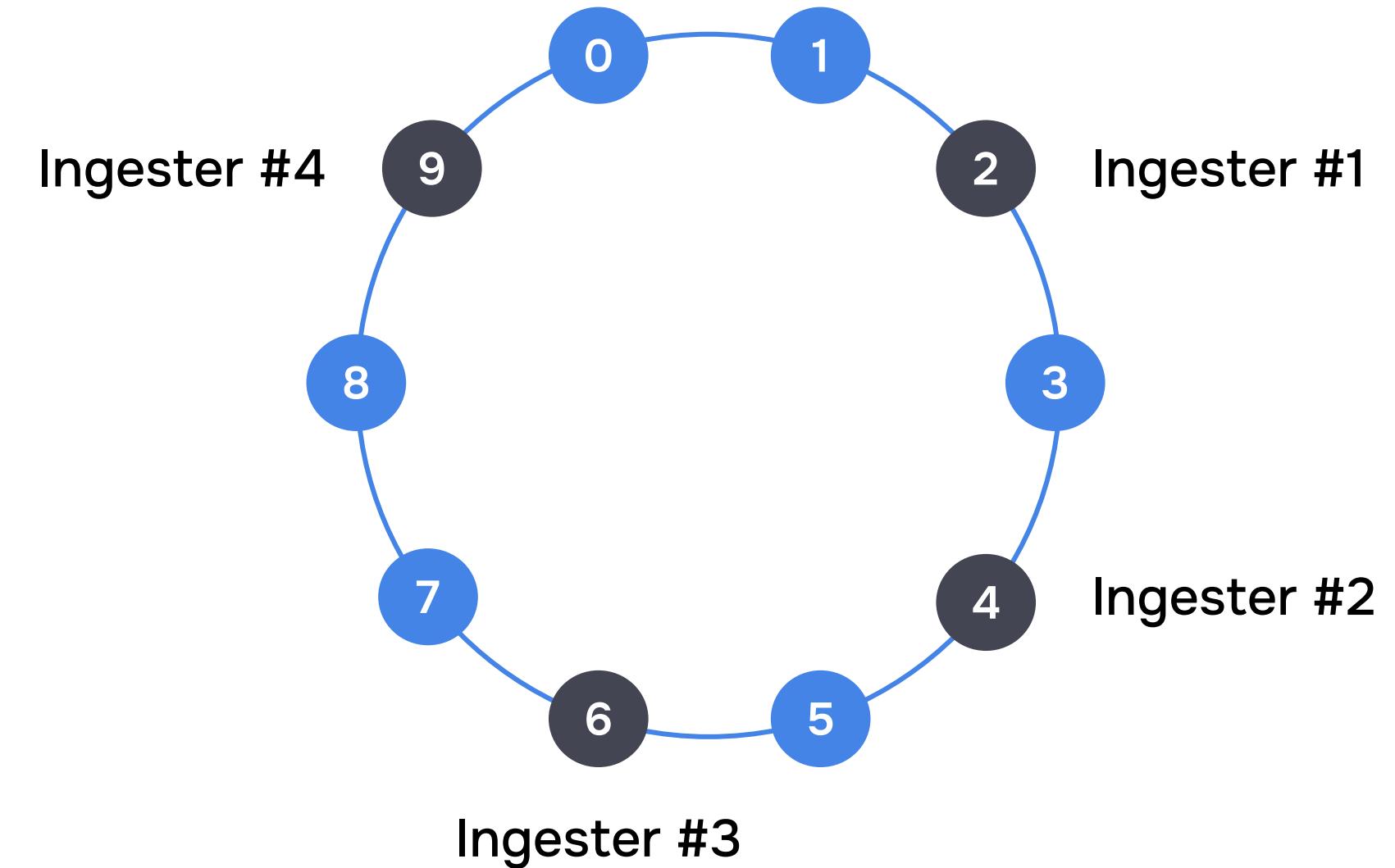
Хеш-ринг



Хеш-ринг

Series

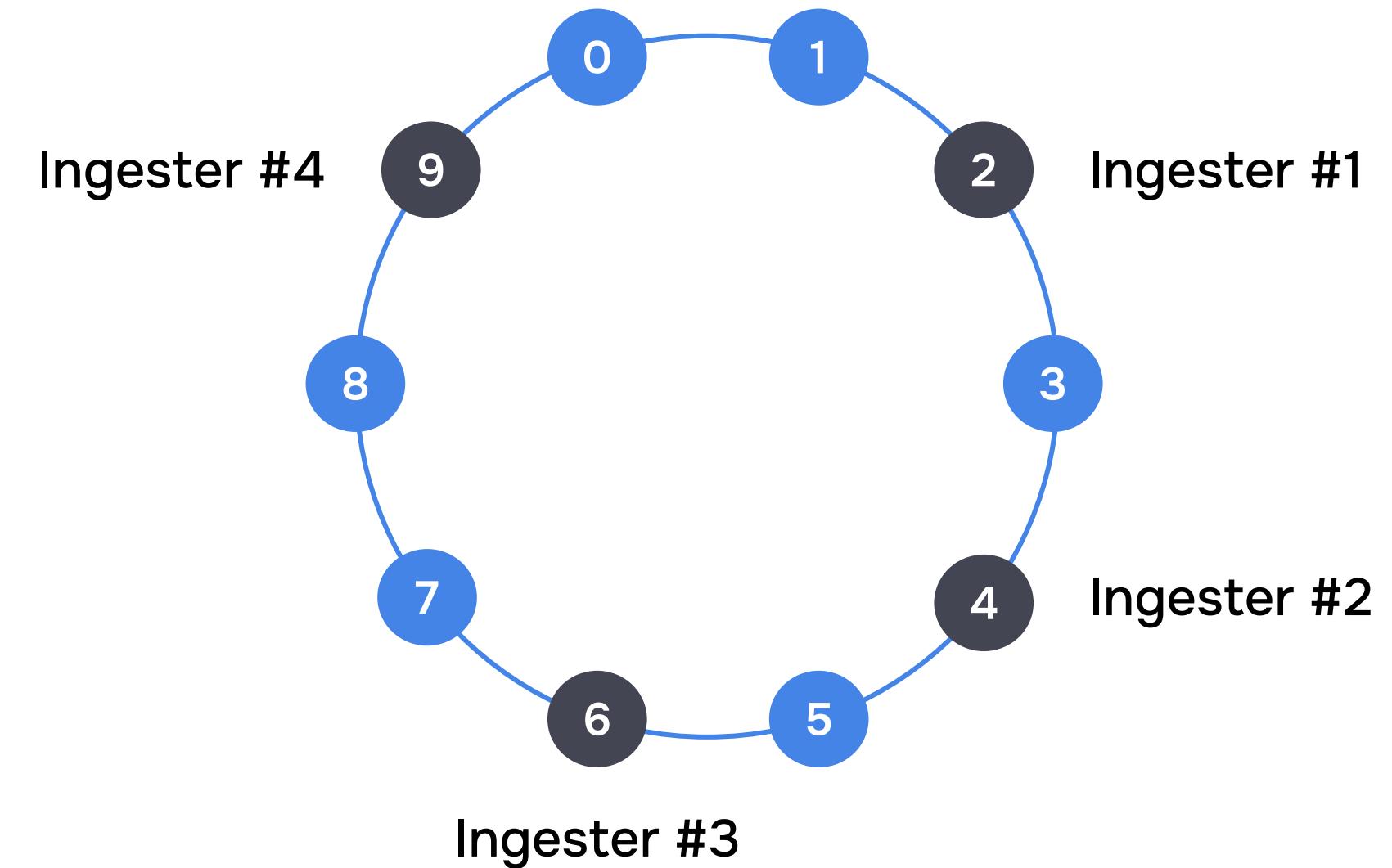
```
{  
    __name__="cpu_seconds_total,  
    instance: "1.1.1.1"  
}
```



Хеш-ринг

Series

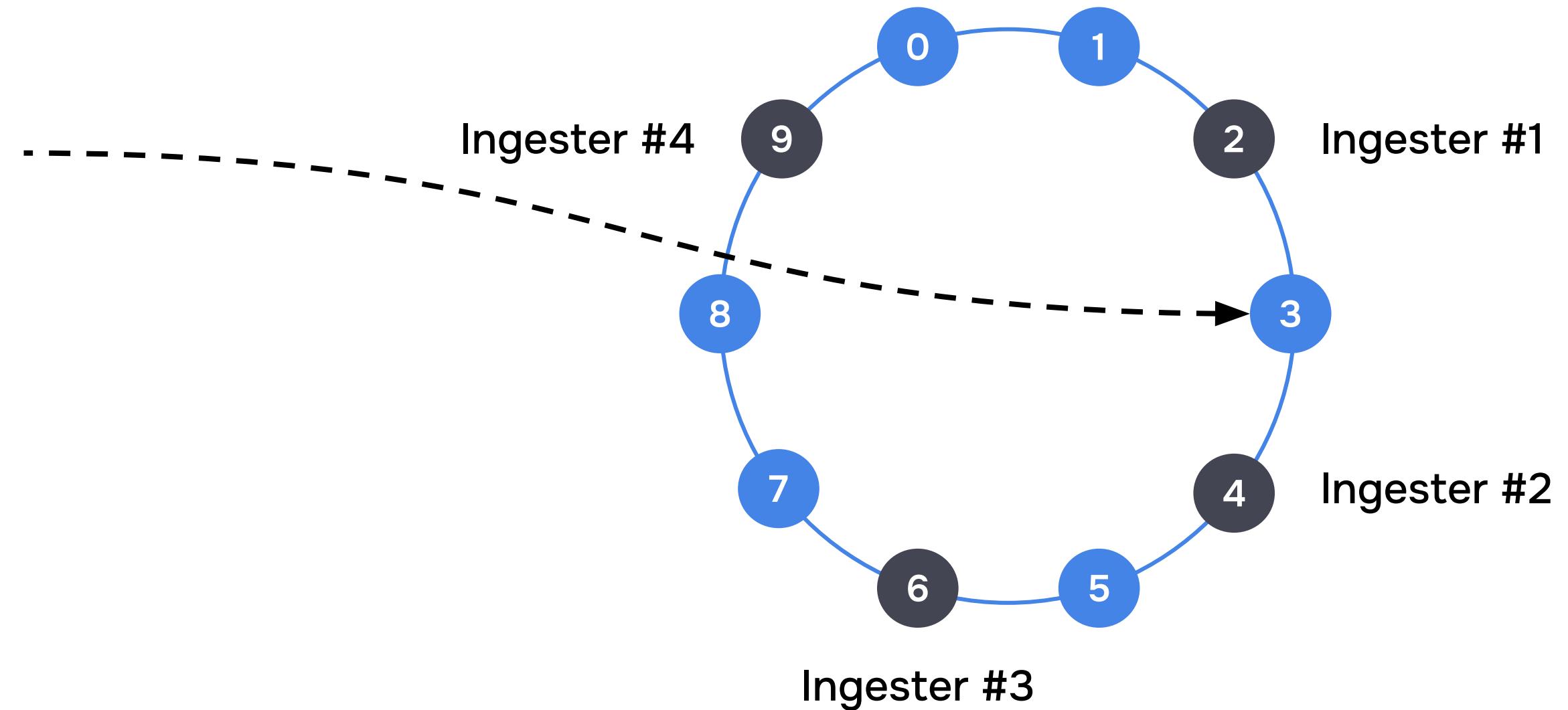
```
{  
    __name__="cpu_seconds_total,  
    instance: "1.1.1.1"  
}  
-> token = 3
```



Хеш-ринг

Series

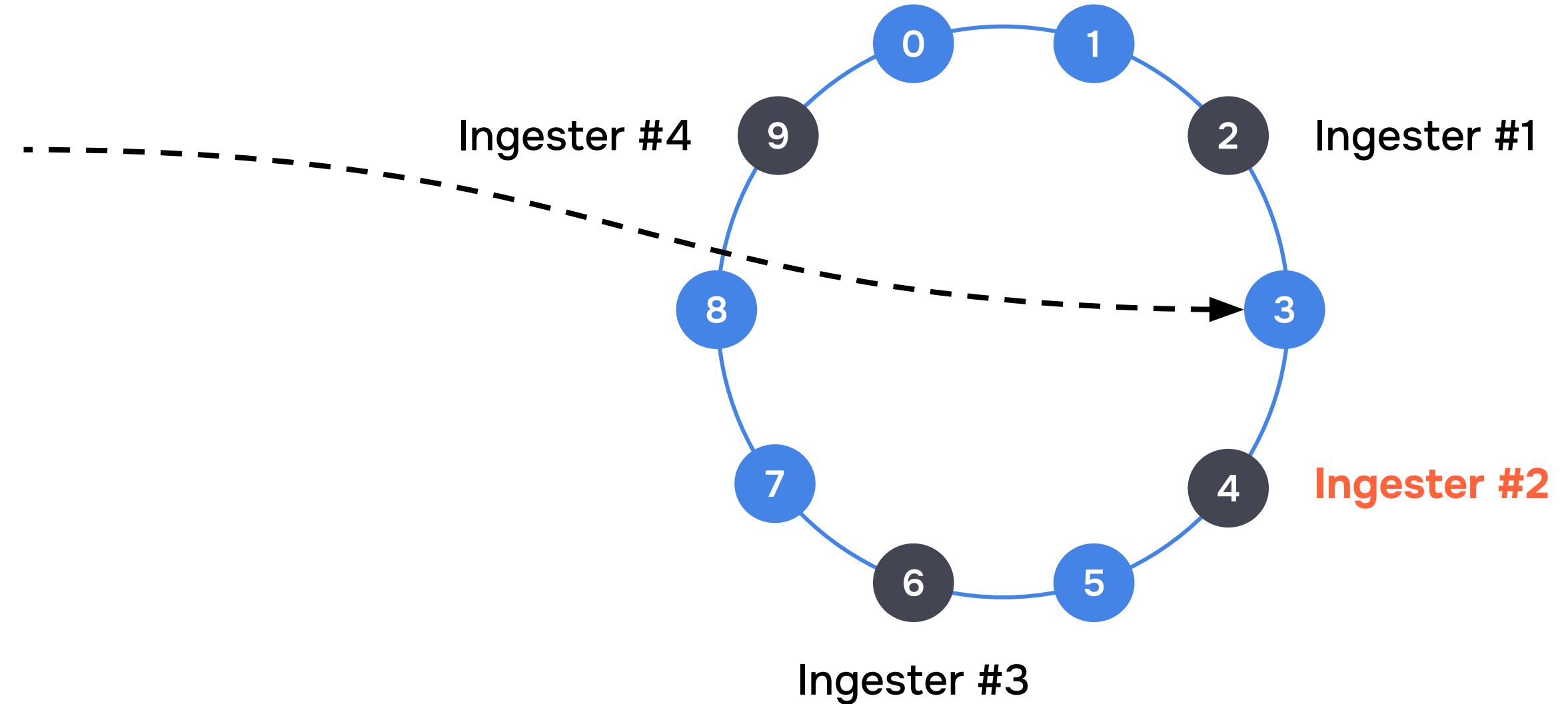
```
{  
    __name__="cpu_seconds_total,  
    instance: "1.1.1.1"  
}  
-> token = 3
```



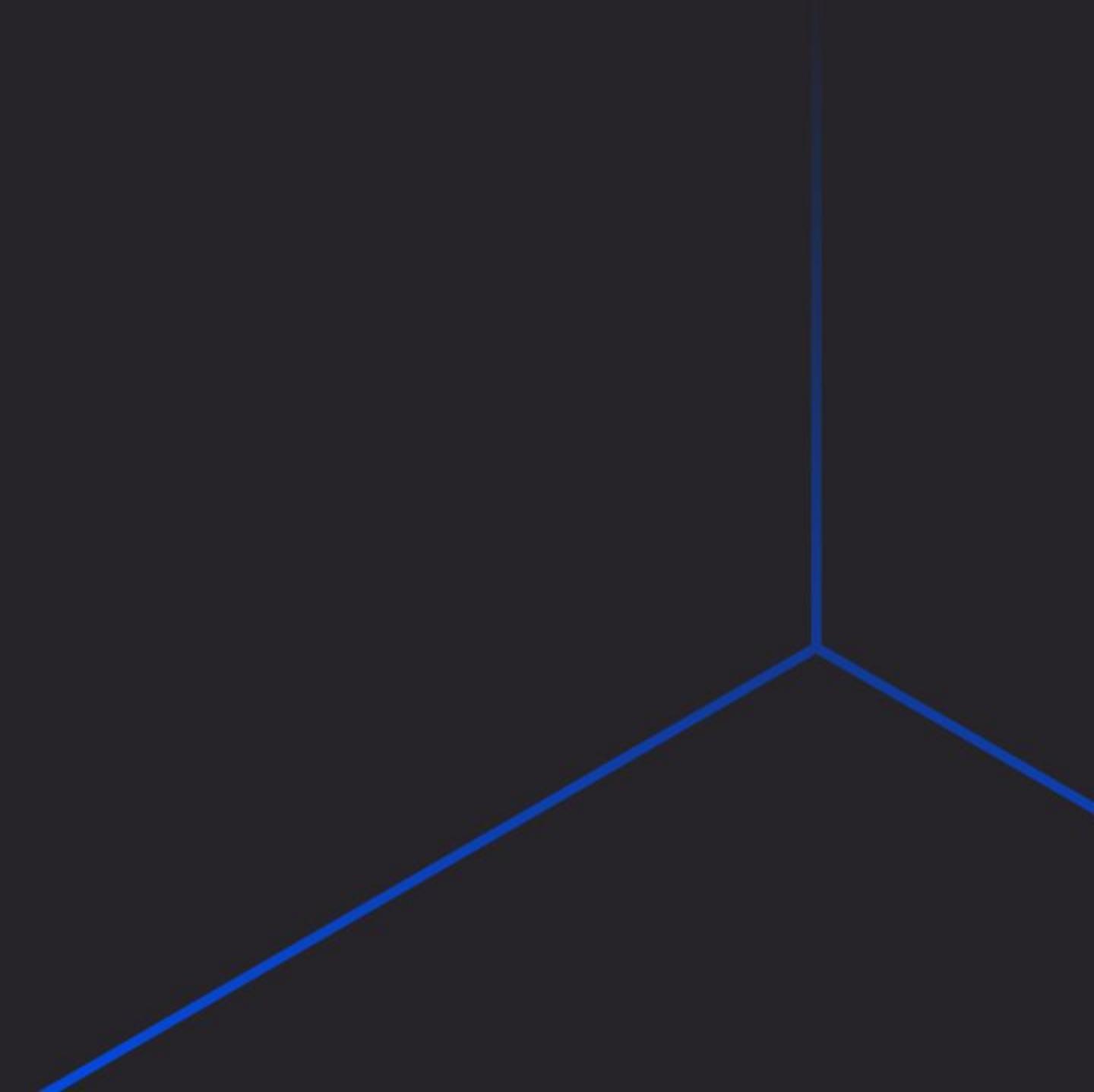
Хеш-ринг

Series

```
{  
    __name__="cpu_seconds_total,  
    instance: "1.1.1.1"  
}  
-> token = 3
```



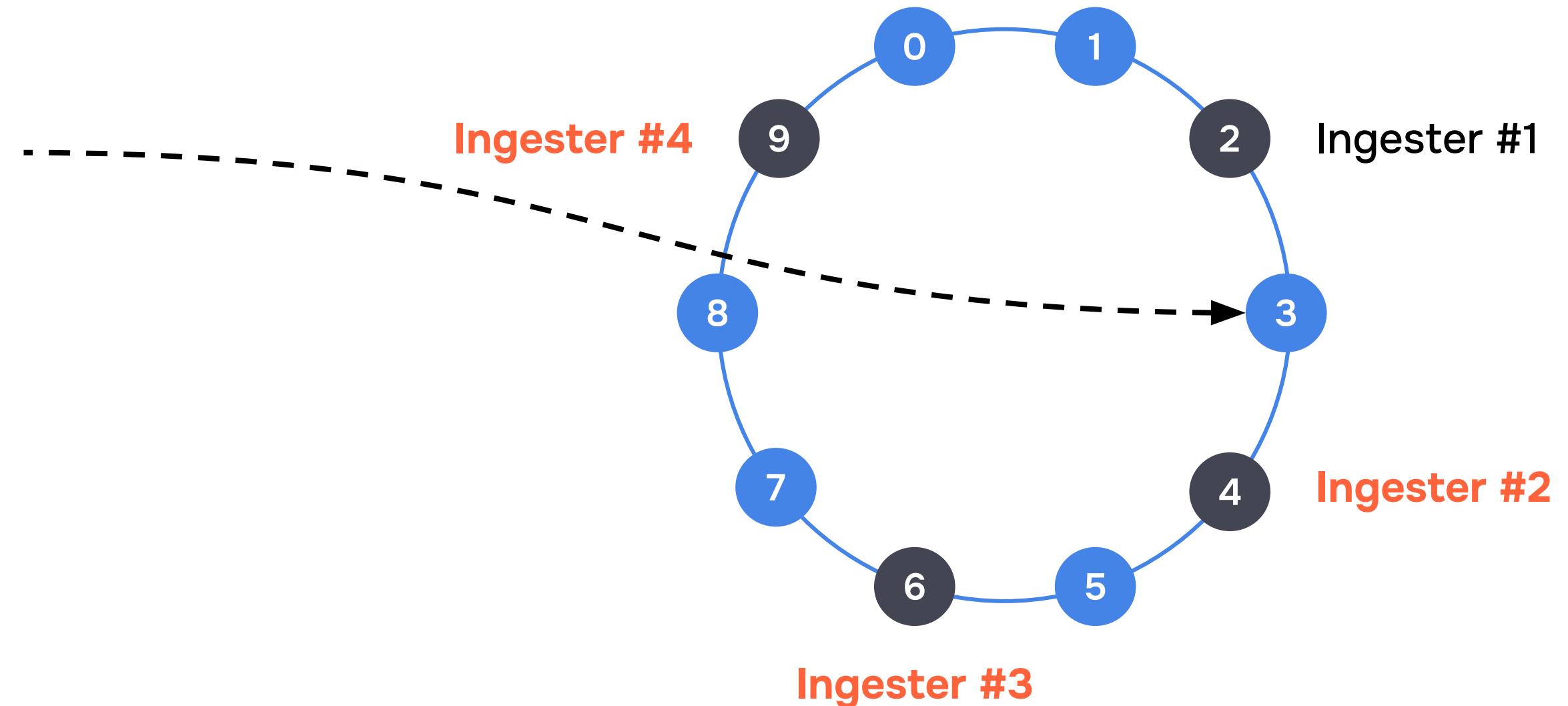
Репликация, а в какие ещё надо записать?



Хеш-ринг

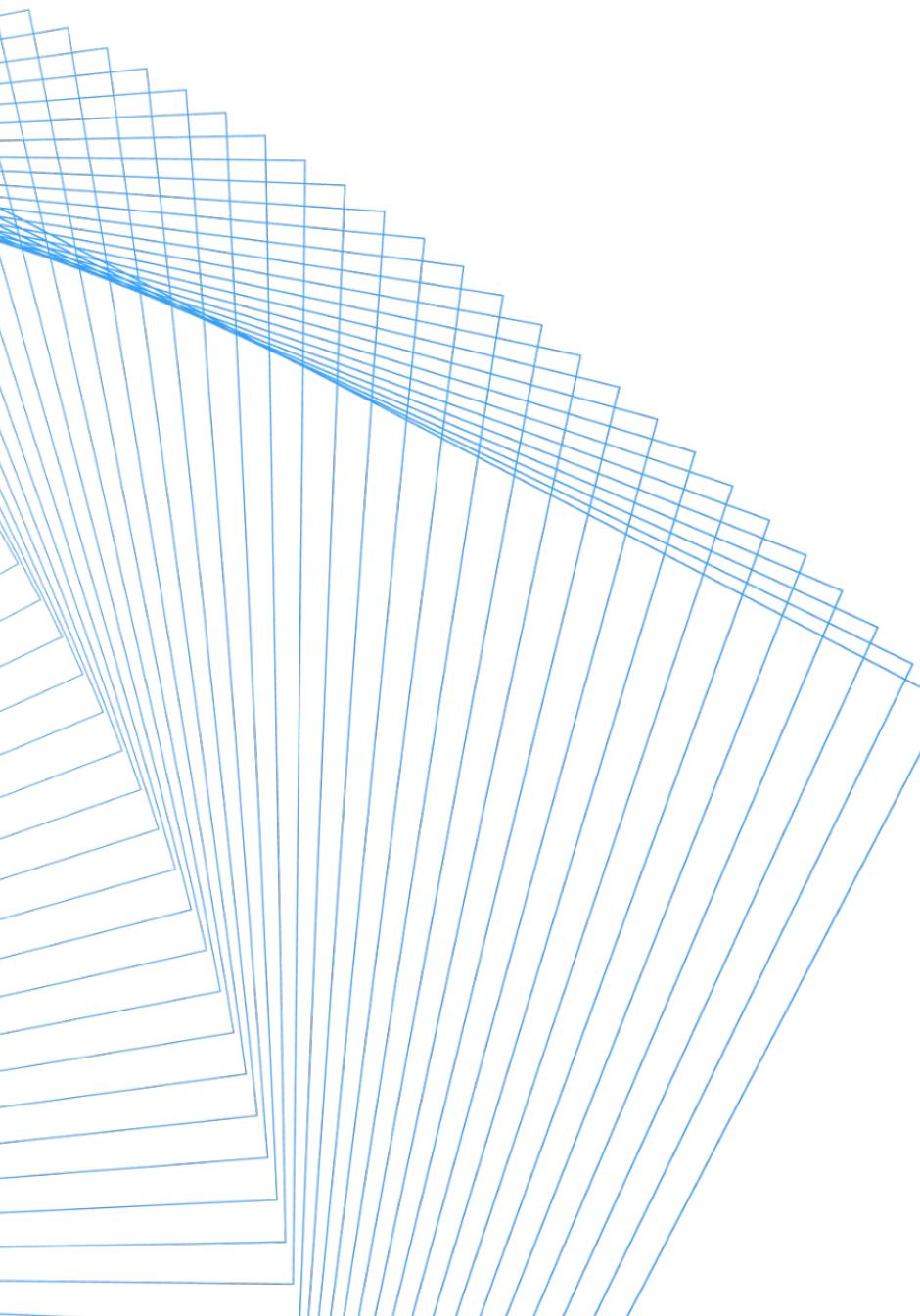
Series

```
{  
    __name__="cpu_seconds_total,  
    instance: "1.1.1.1"  
}  
-> token = 3
```

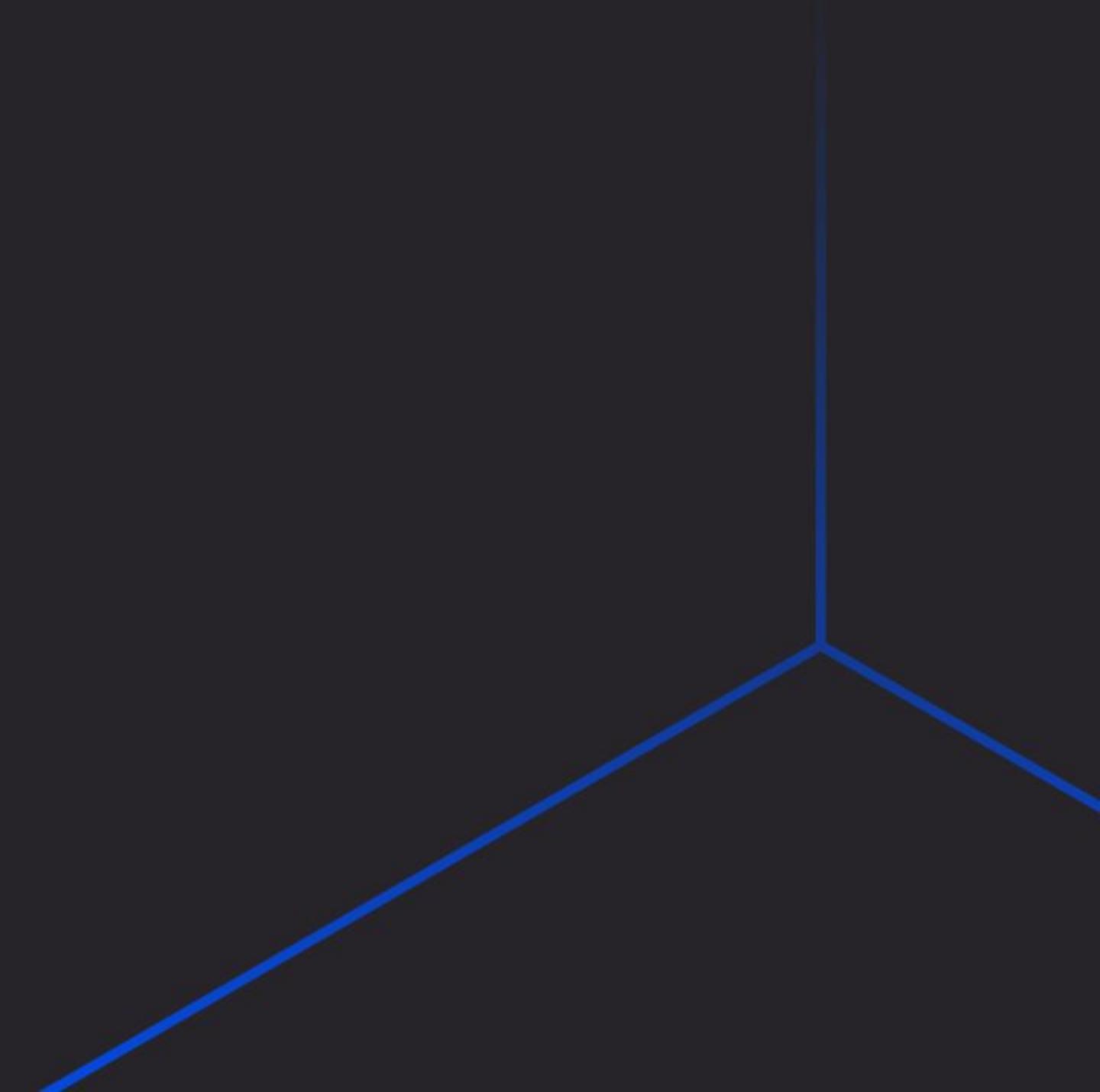


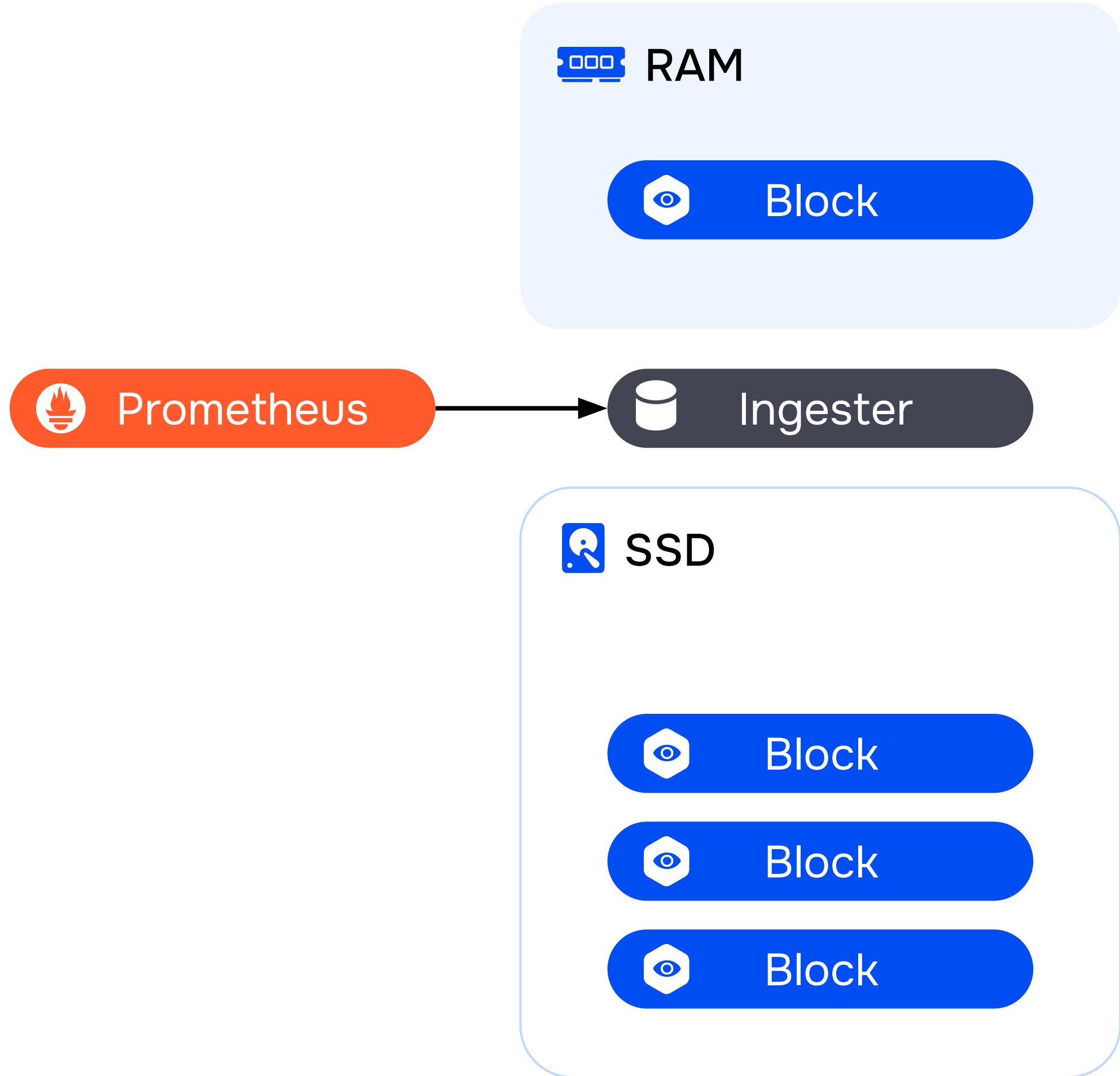
Централизованная система мониторинга

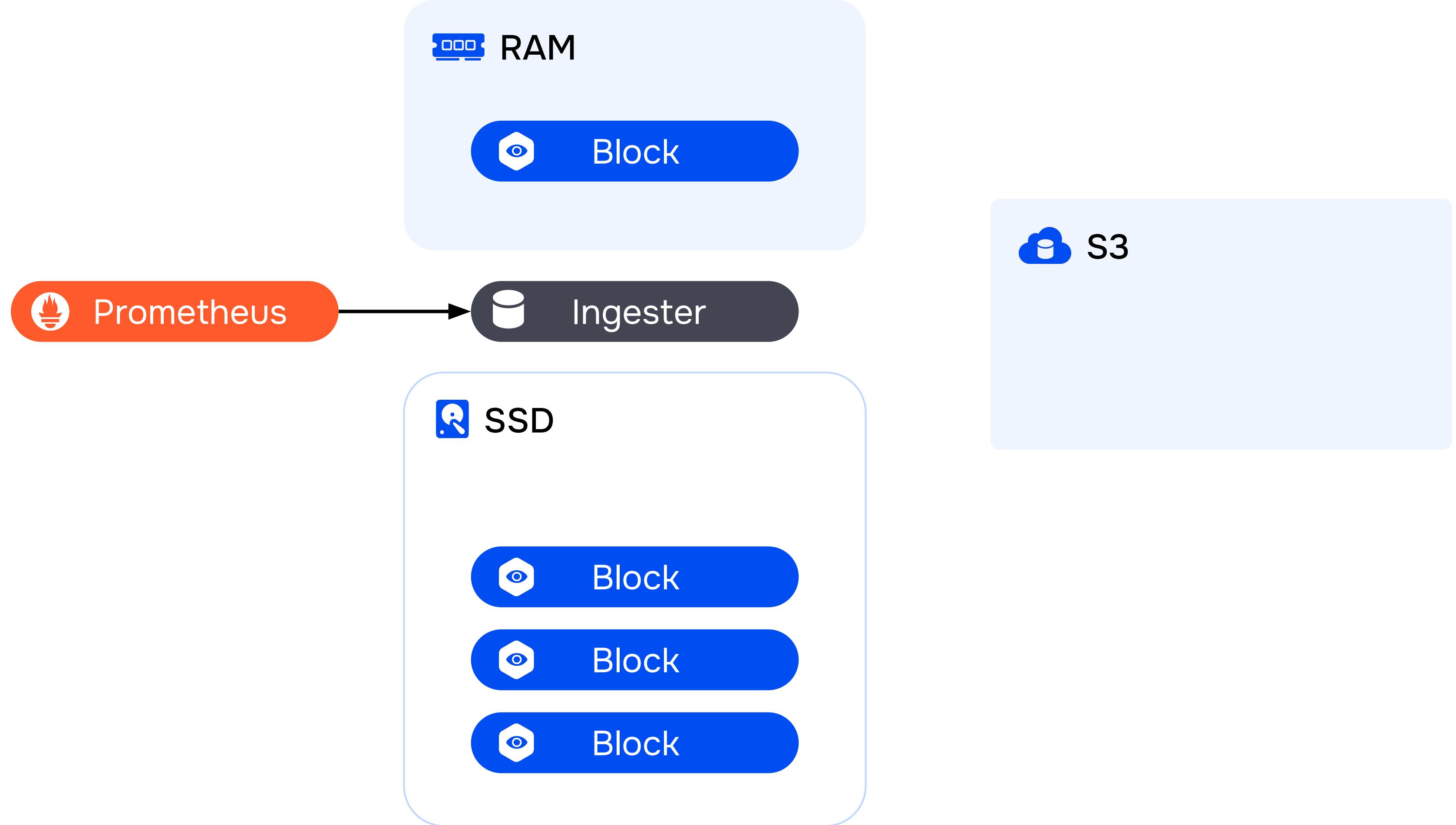
- Отказоустойчивость на запись: $WC = RF/2 + 1$
- Отказоустойчивость на чтение: $RC = RF - WC + 1$
- Шардирование на запись на основе хеш-ринга

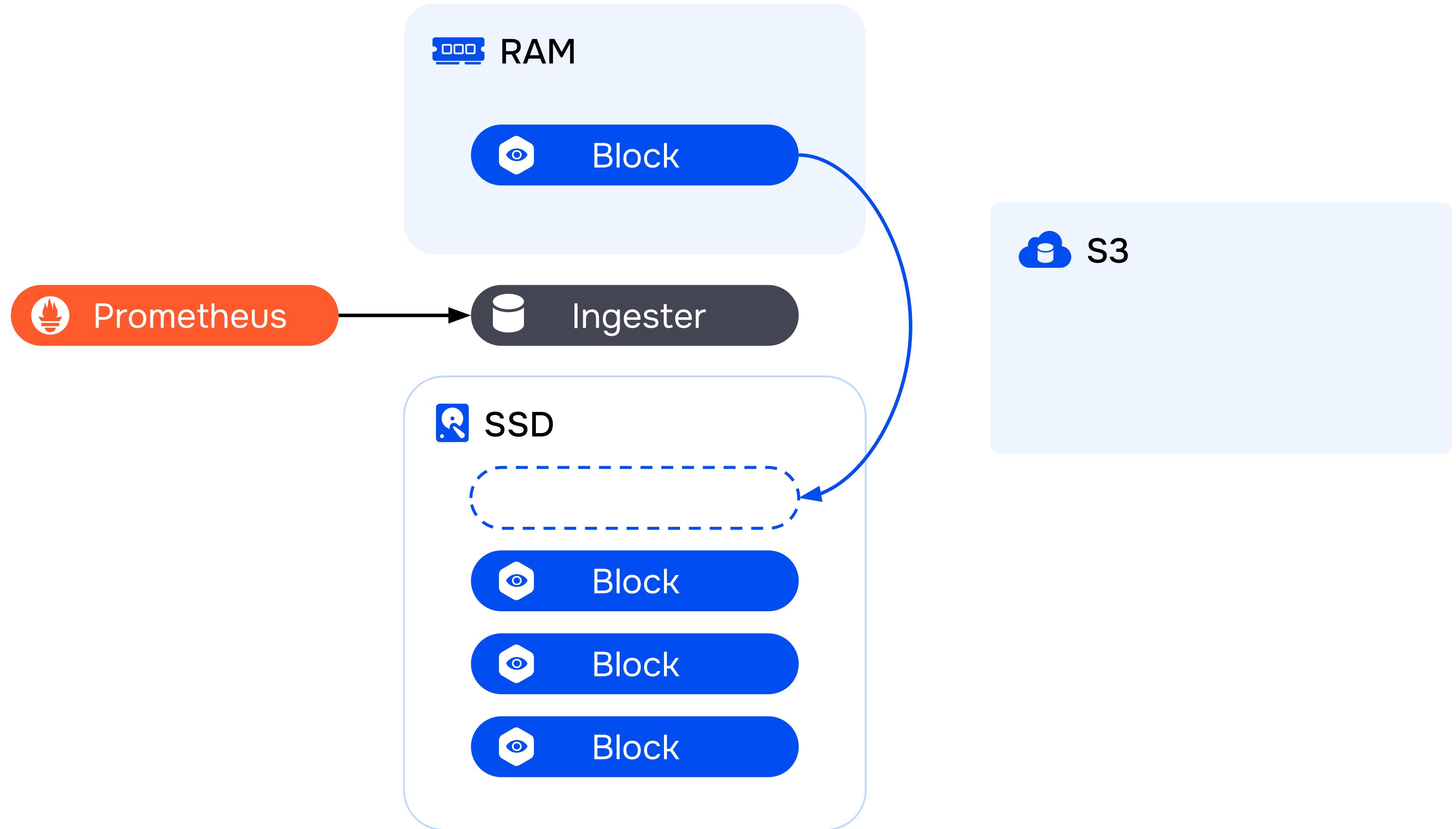


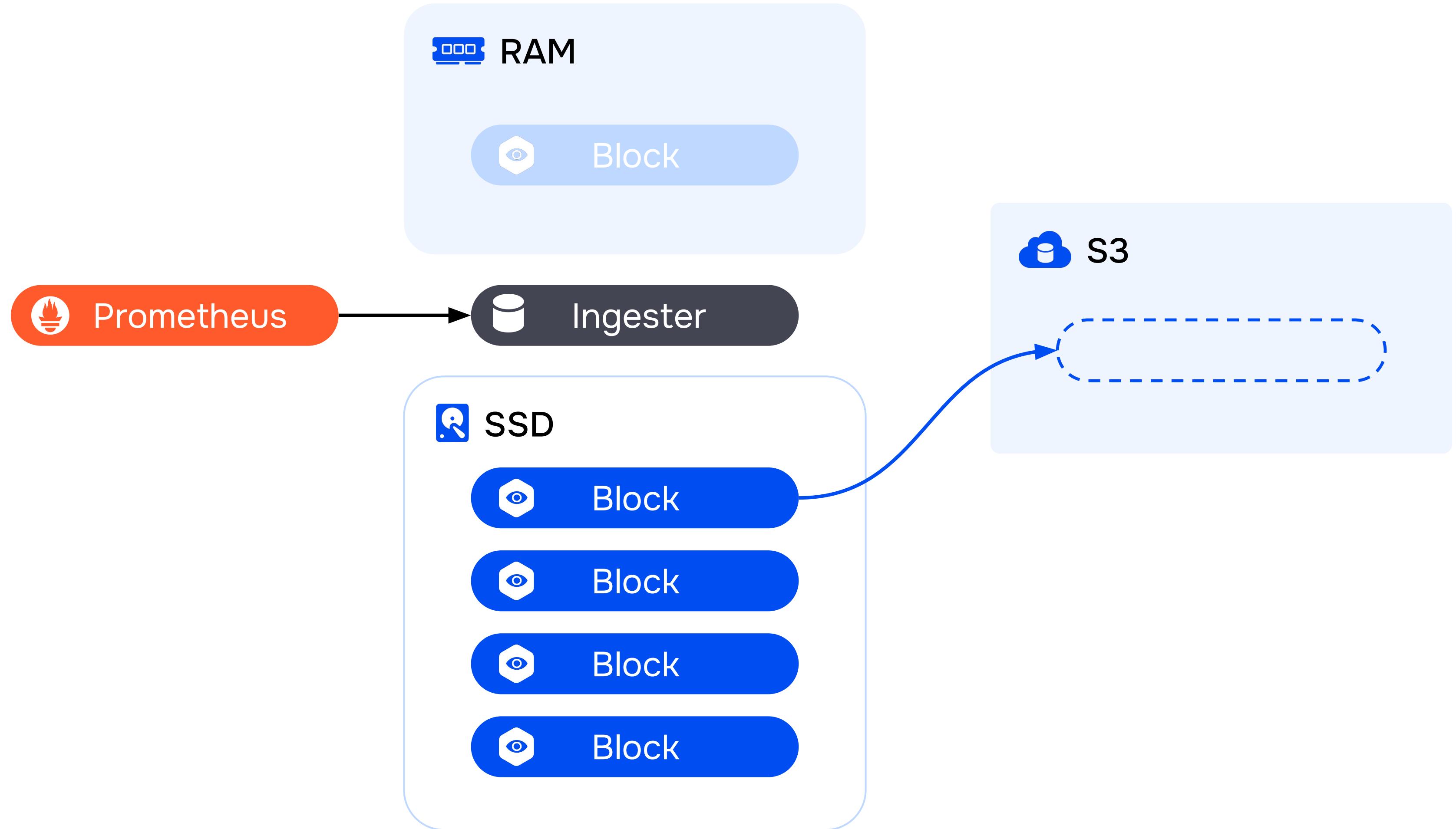
Но держать данные на Ingestor
неудобно!

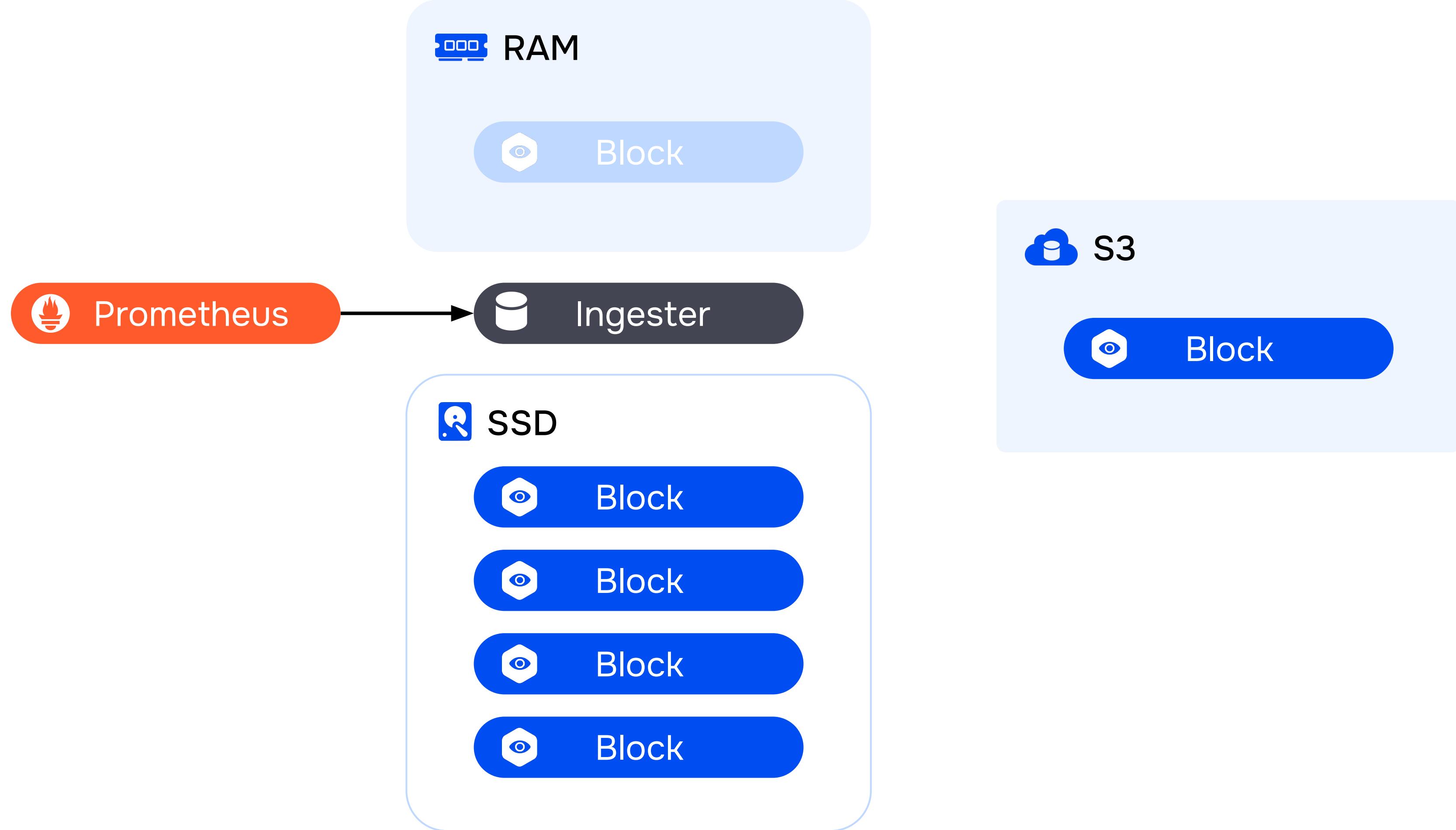


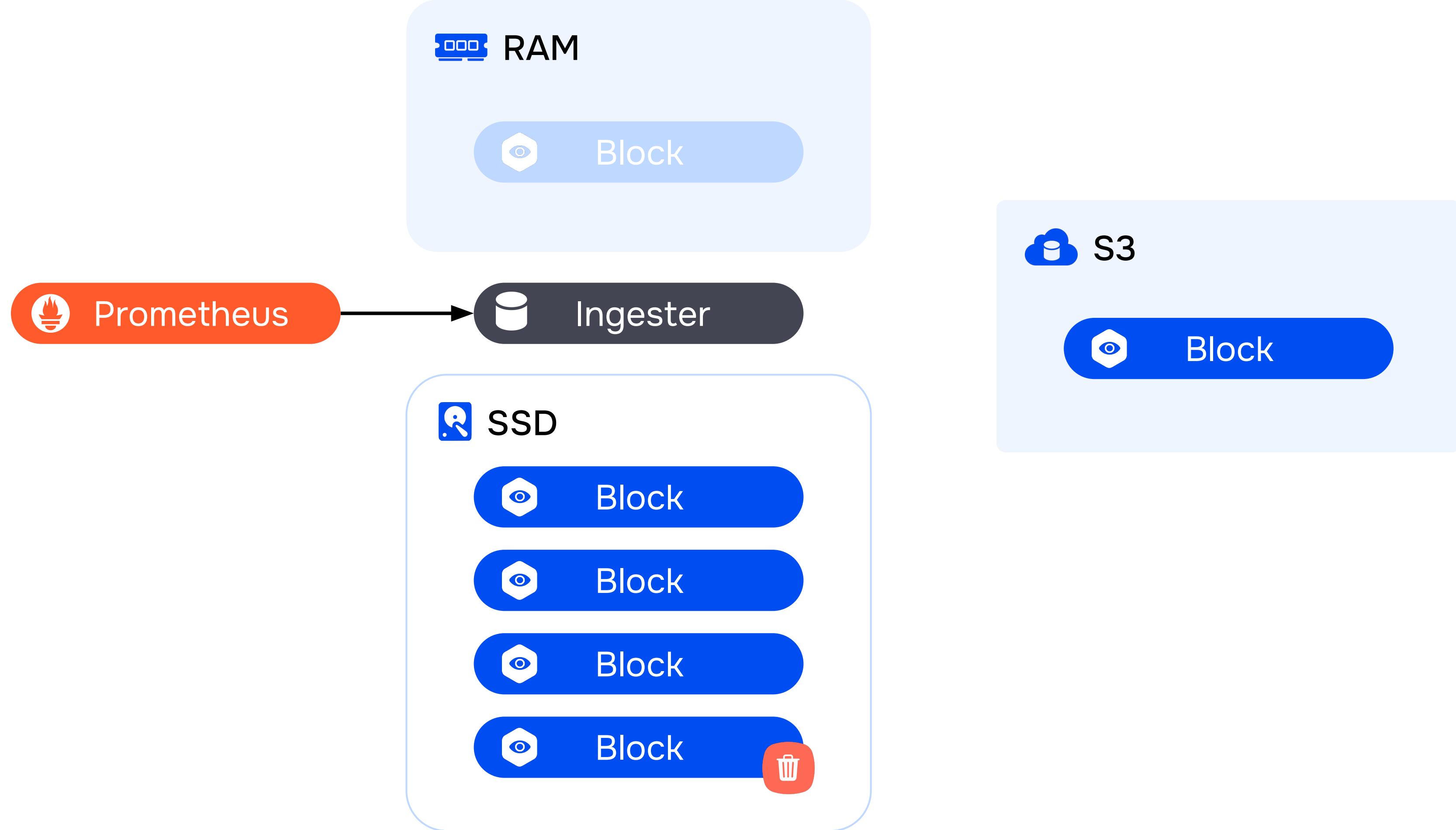


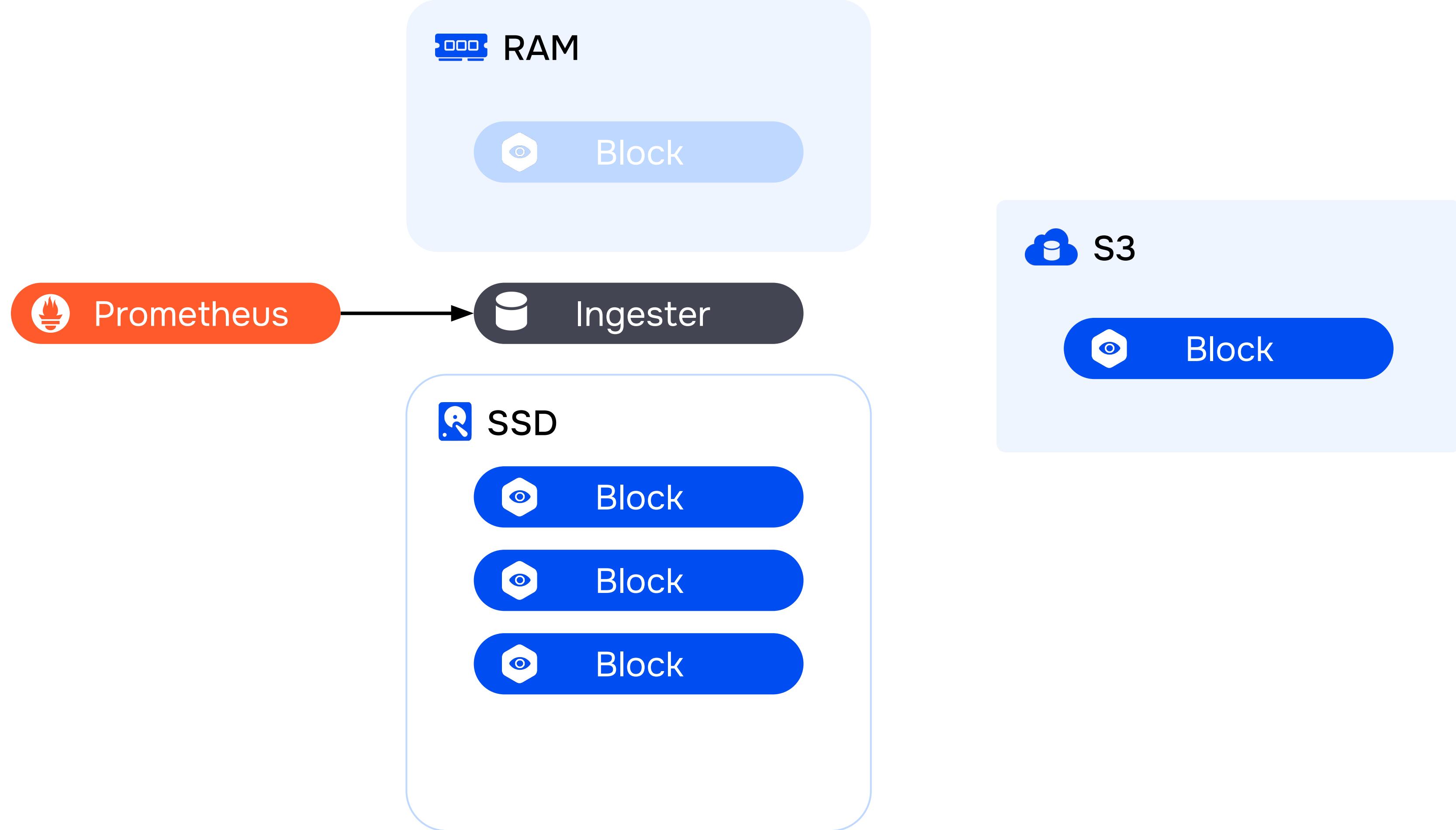


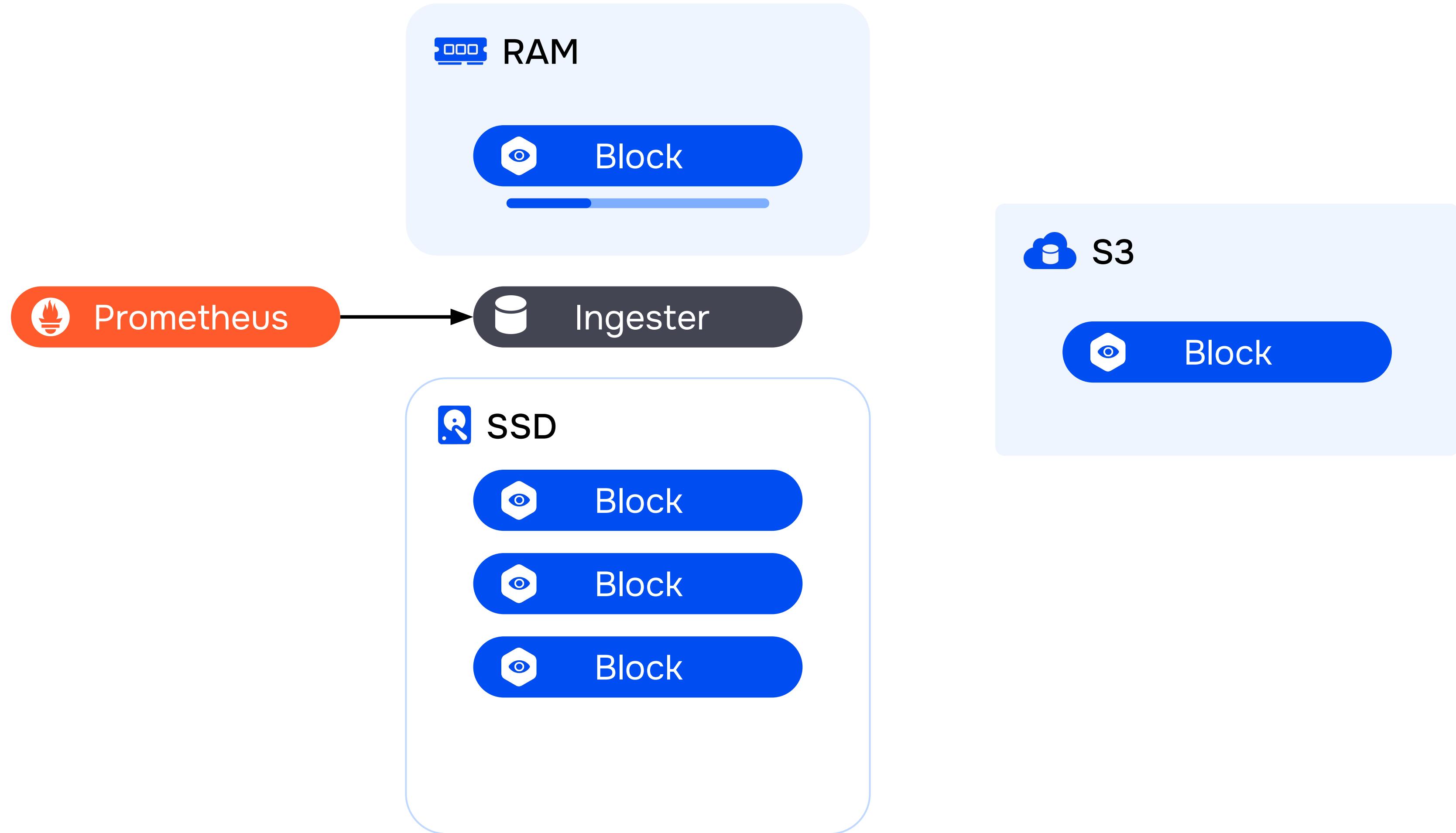


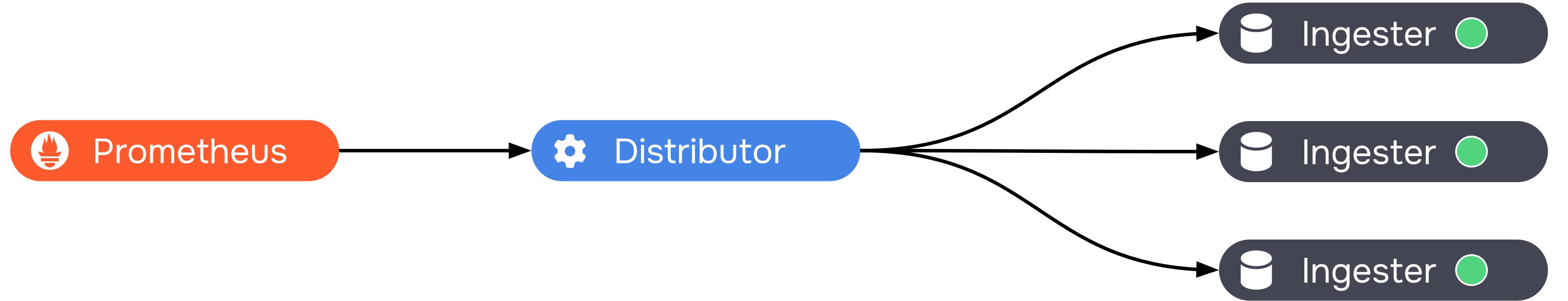














S3



Ingestor-1 Block-X

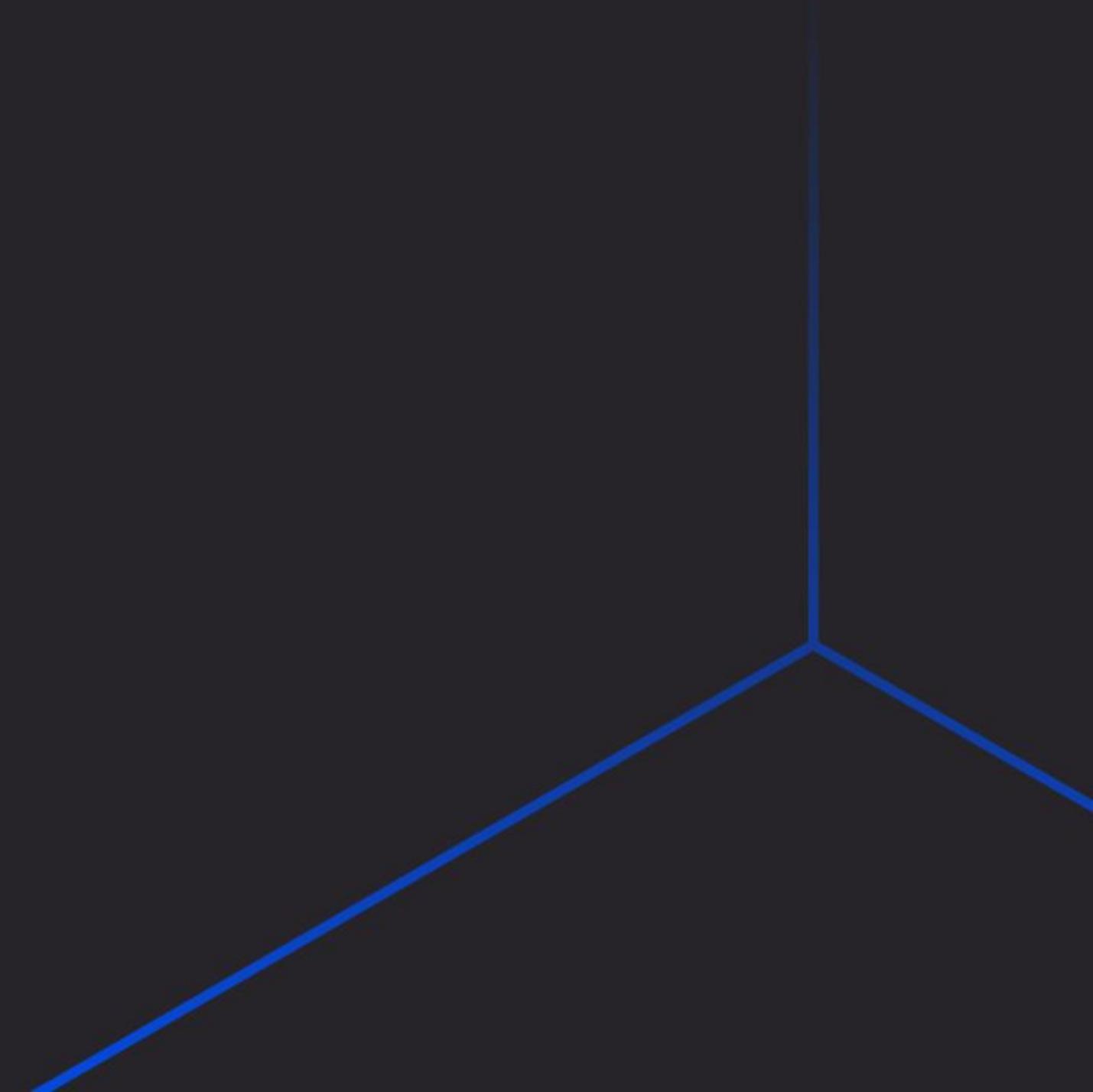


Ingestor-2 Block-X



Ingestor-3 Block-X

Устранием двойное резервирование





S3



Ingestor-1 Block-X



Ingestor-2 Block-X



Ingestor-3 Block-X

Compactor



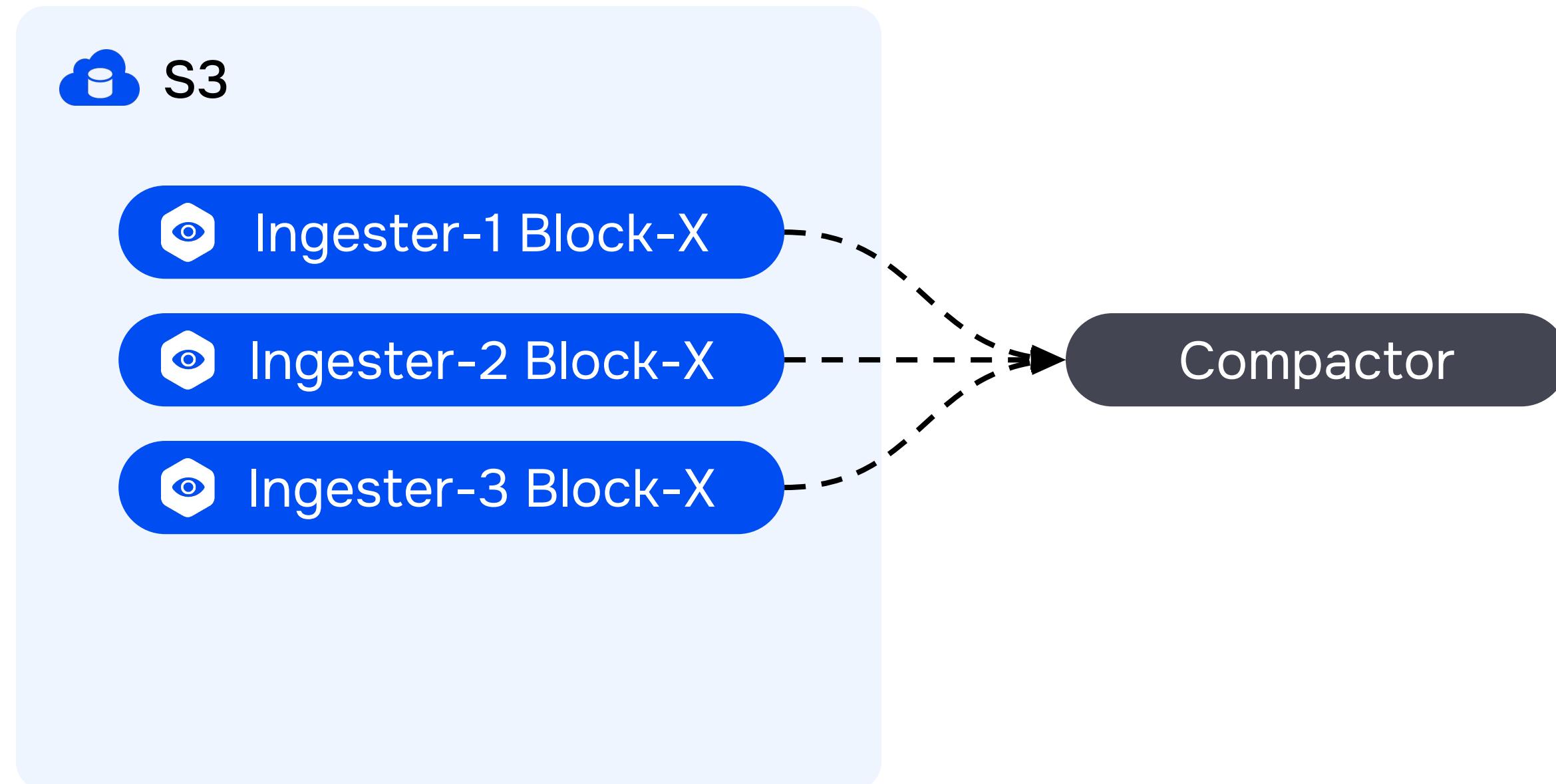
S3

Ingestor-1 Block-X

Ingestor-2 Block-X

Ingestor-3 Block-X

Compactor





S3

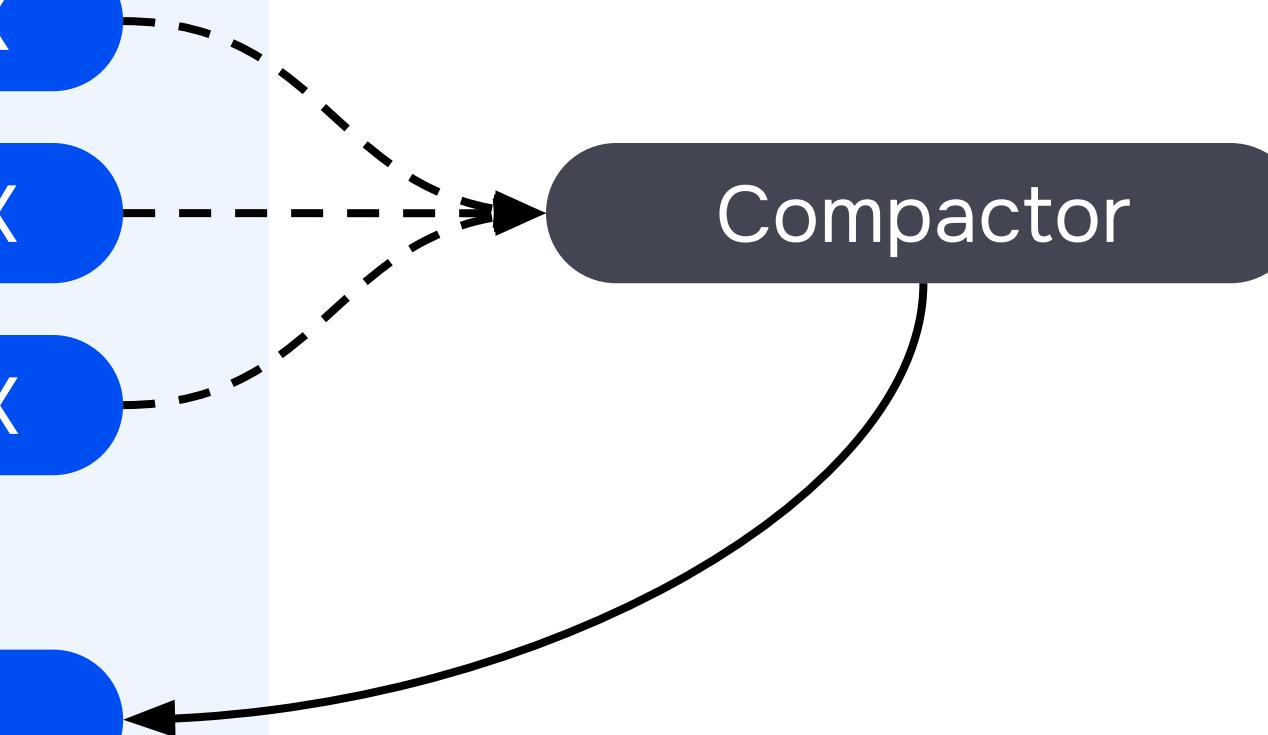
Ingestor-1 Block-X

Ingestor-2 Block-X

Ingestor-3 Block-X

Block-X

Compactor





S3



Ingestor-1 Block-X



Ingestor-2 Block-X



Ingestor-3 Block-X



Block-X

Compactor



S3

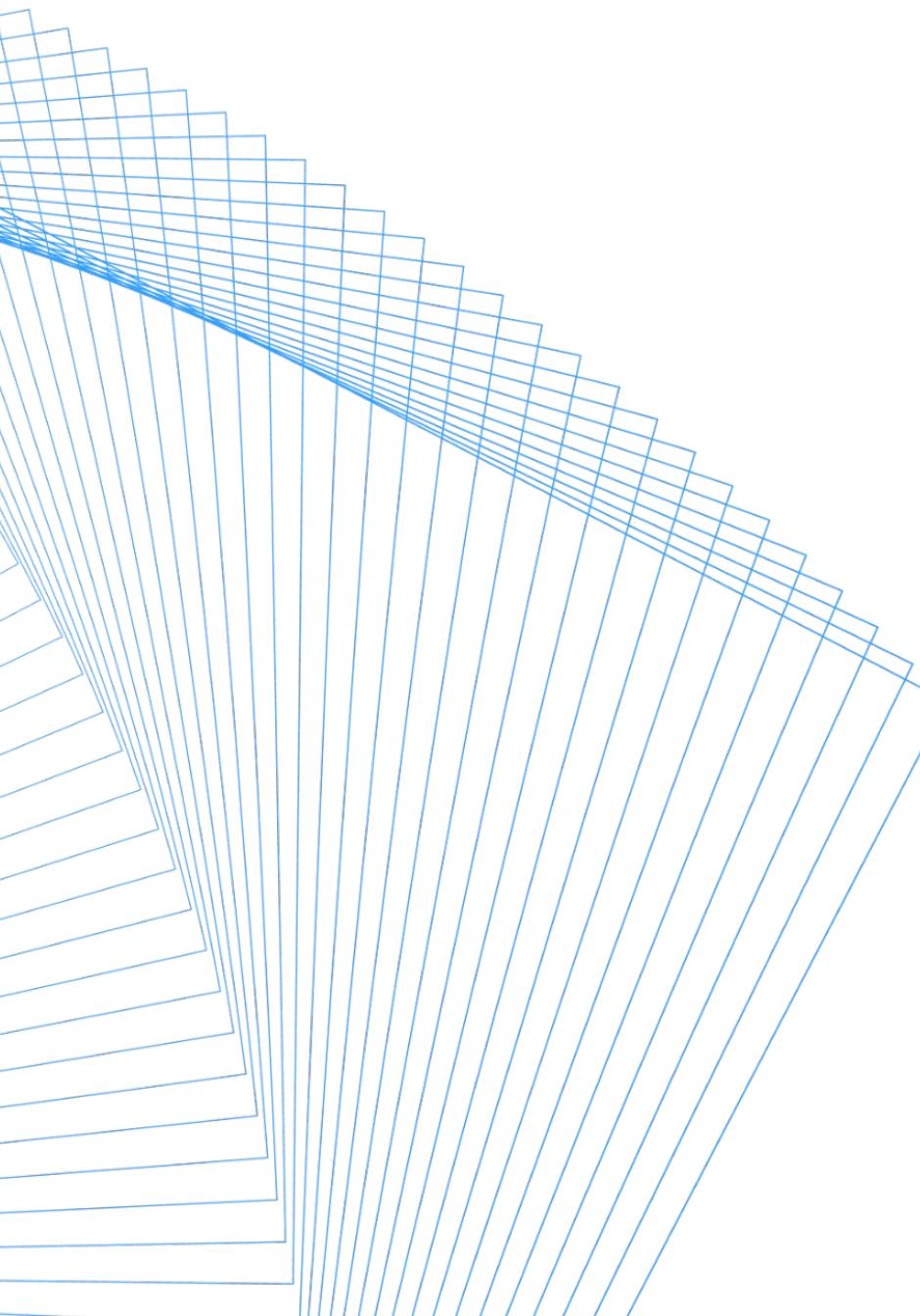


Block-X

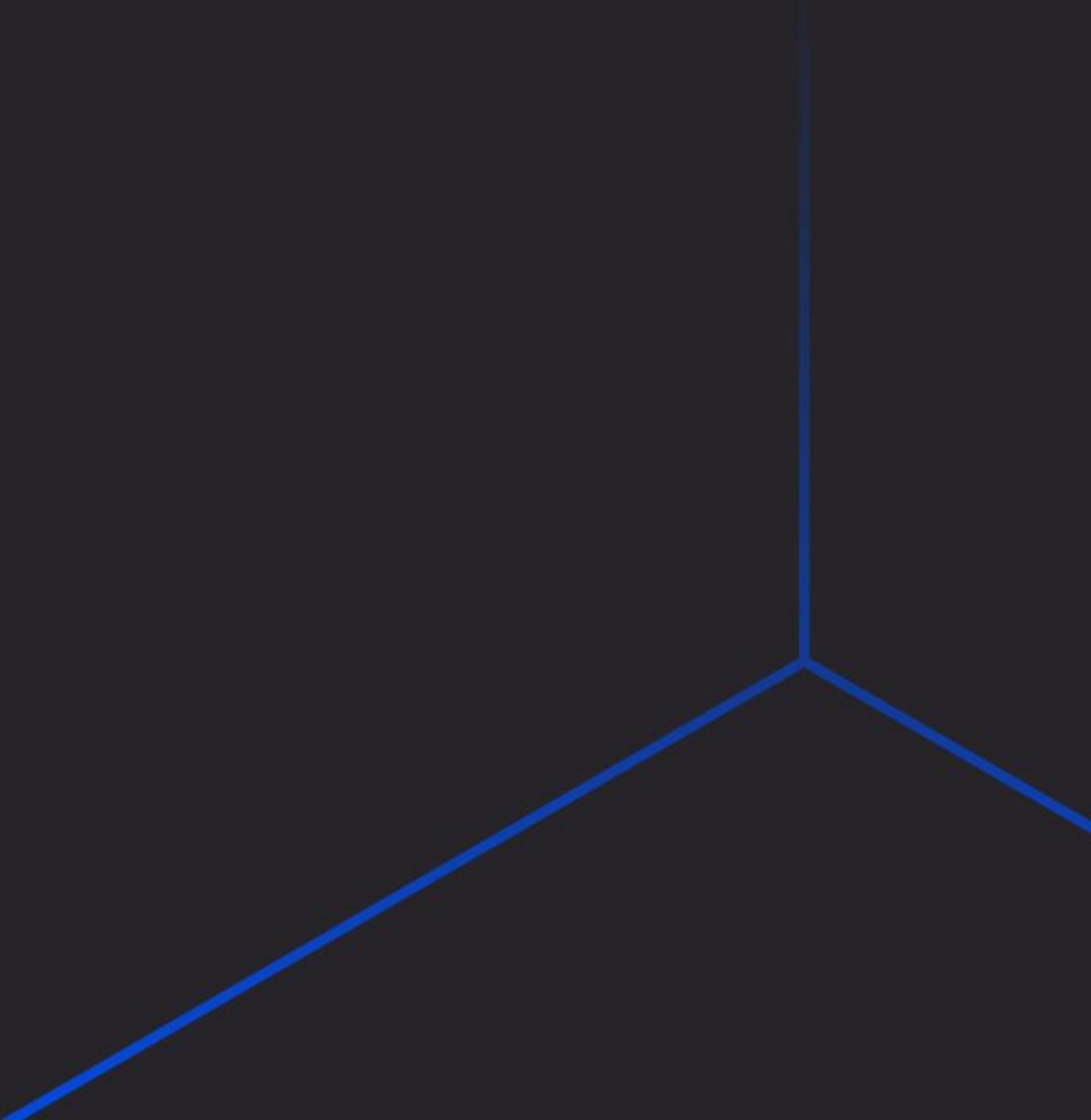
Compactor

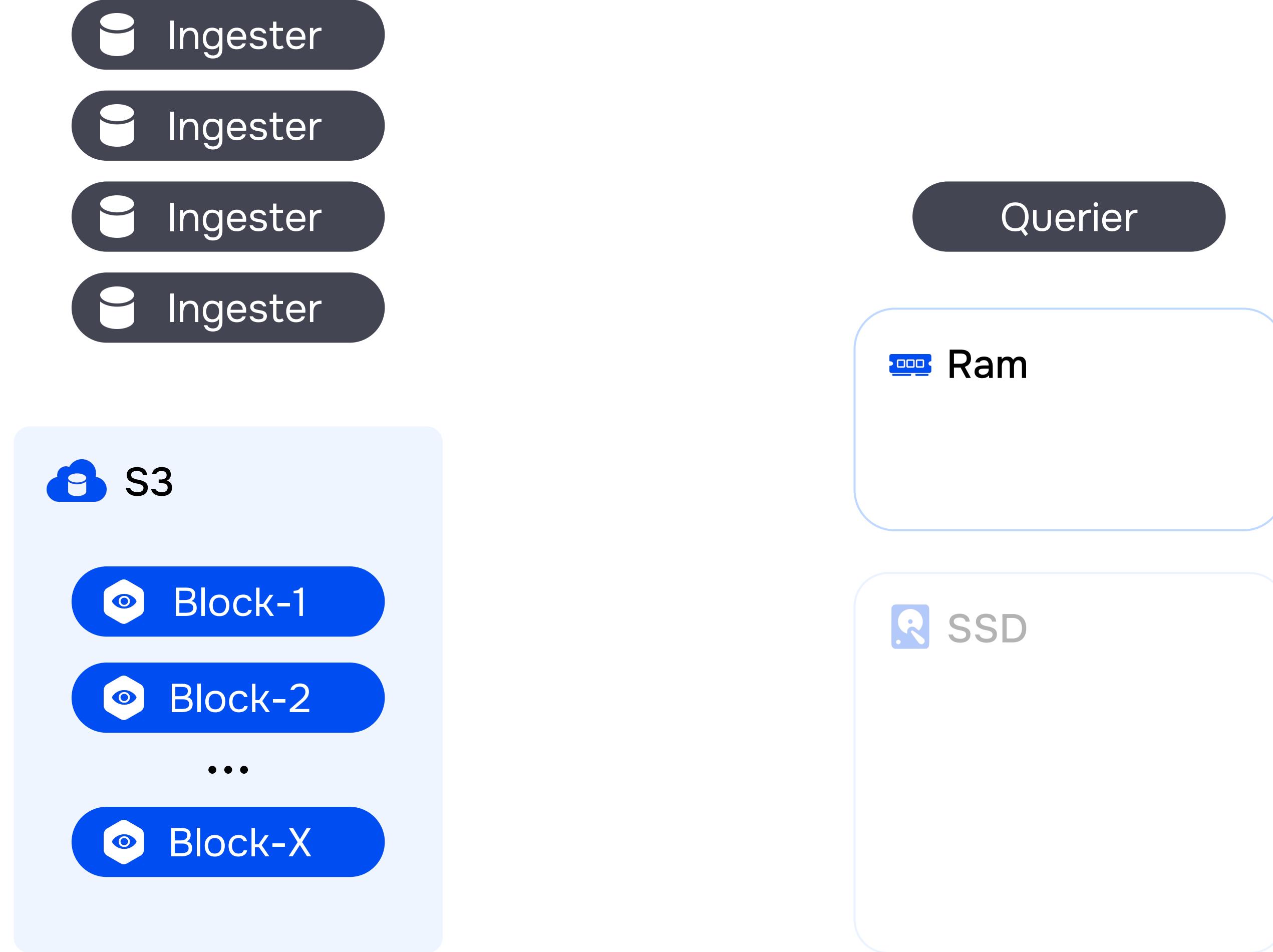
Централизованная система мониторинга

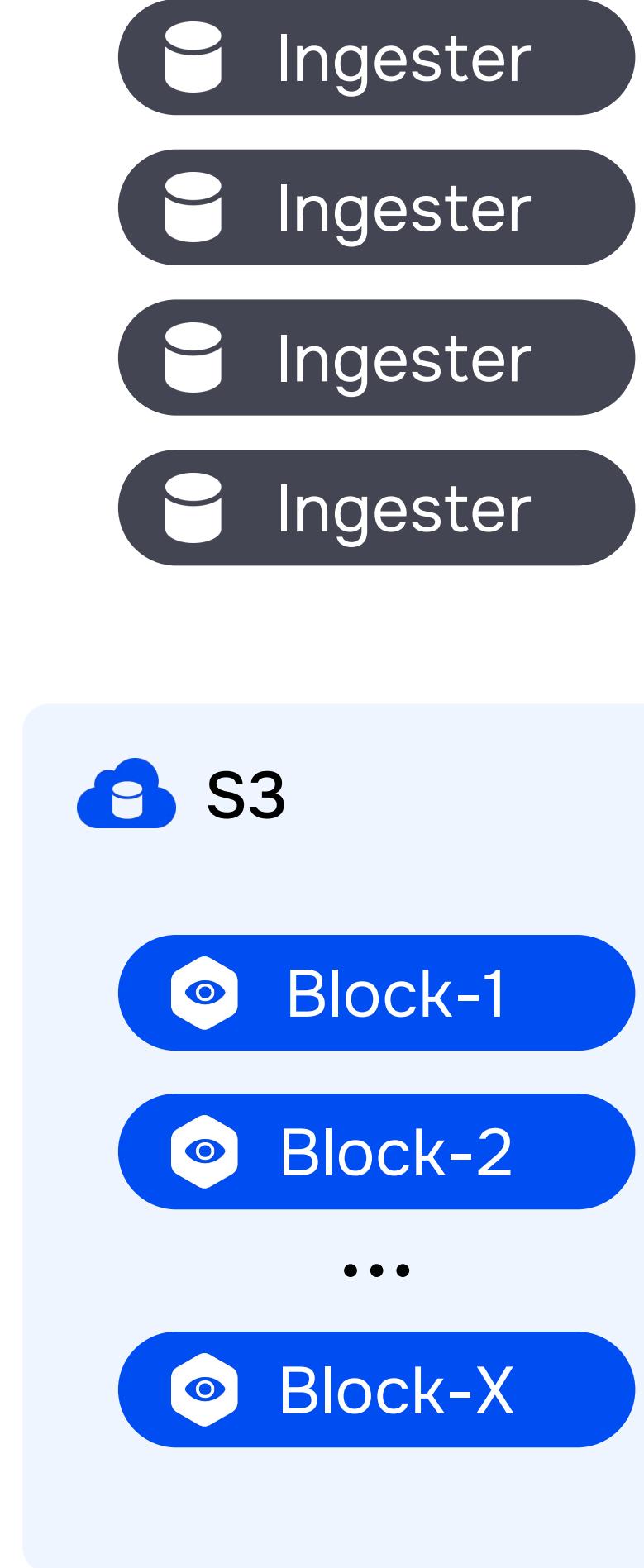
- Отказоустойчивость на запись: $WC = RF/2 + 1$
- Отказоустойчивость на чтение: $RC = RF - WC + 1$
- Шардирование на запись на основе хеш-ринга
- Исторические данные хранятся в S3



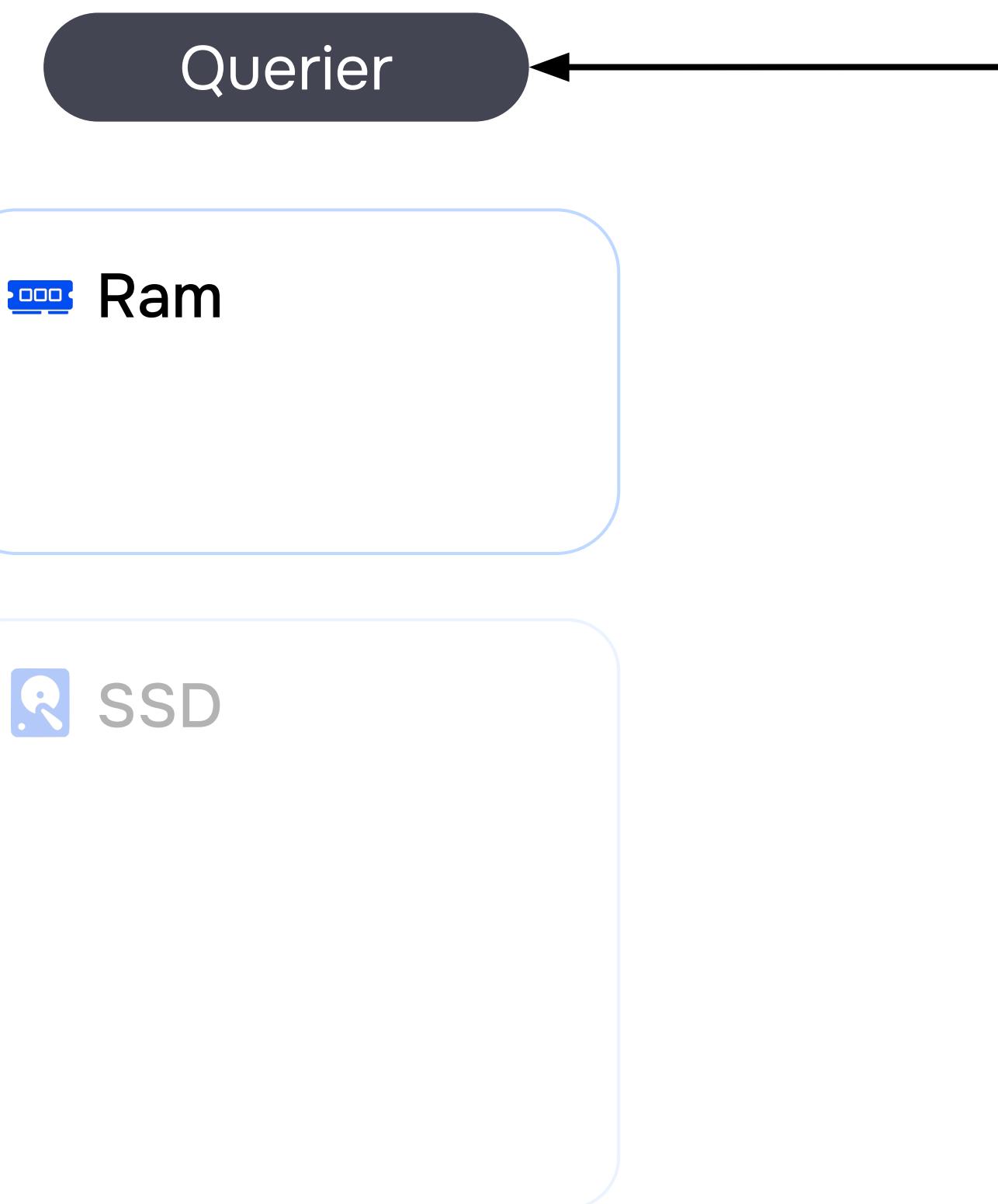
Как читать?

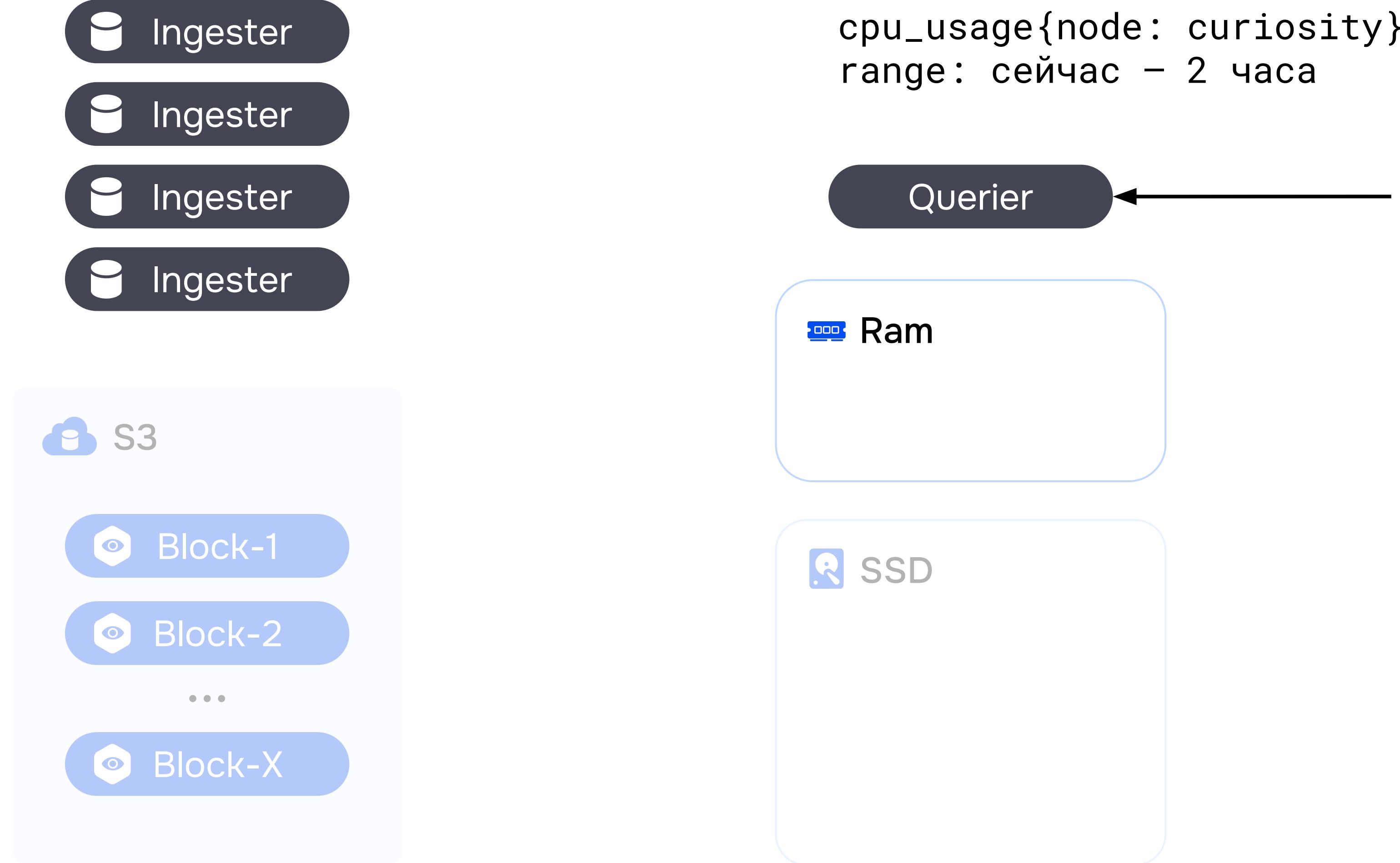


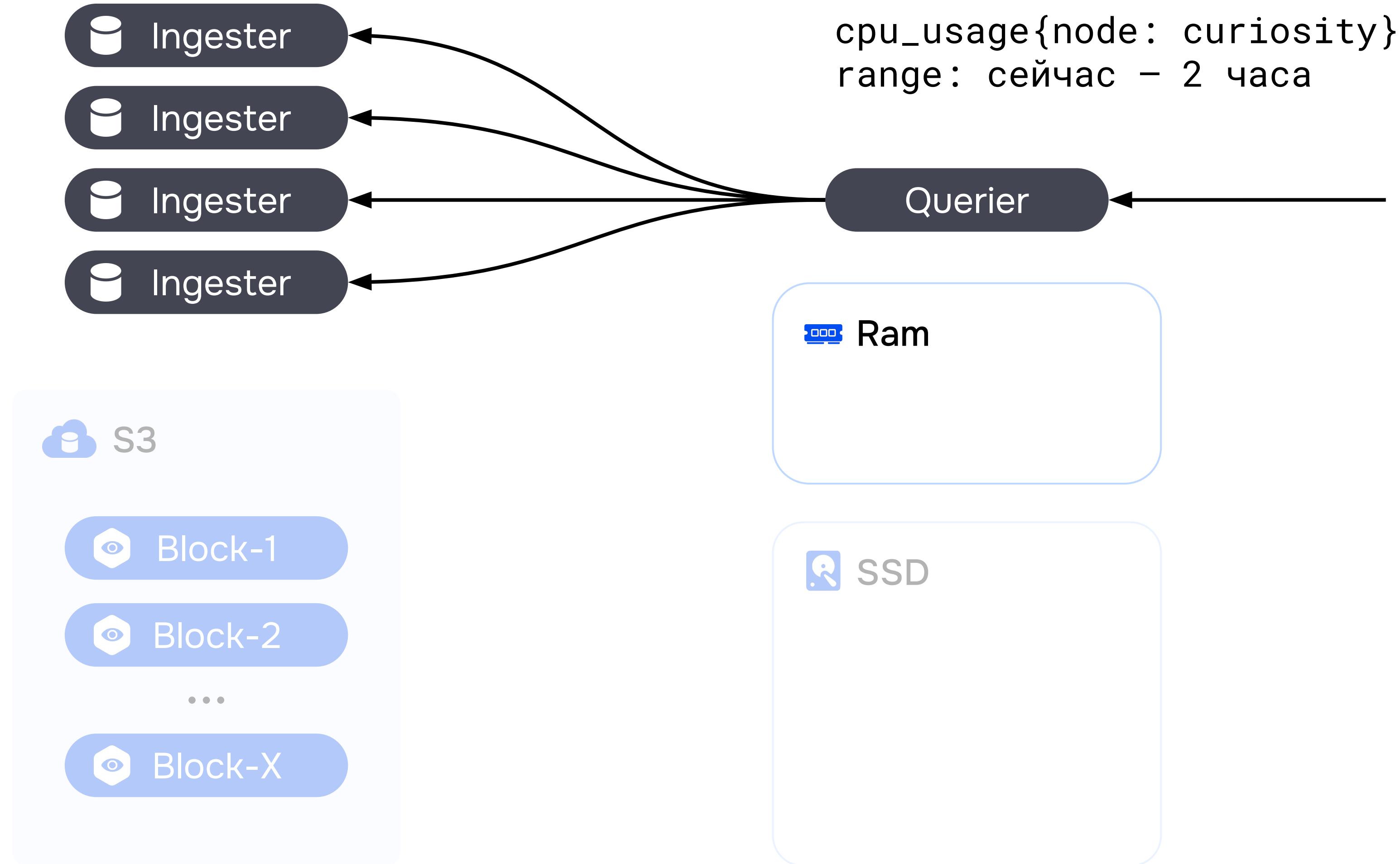


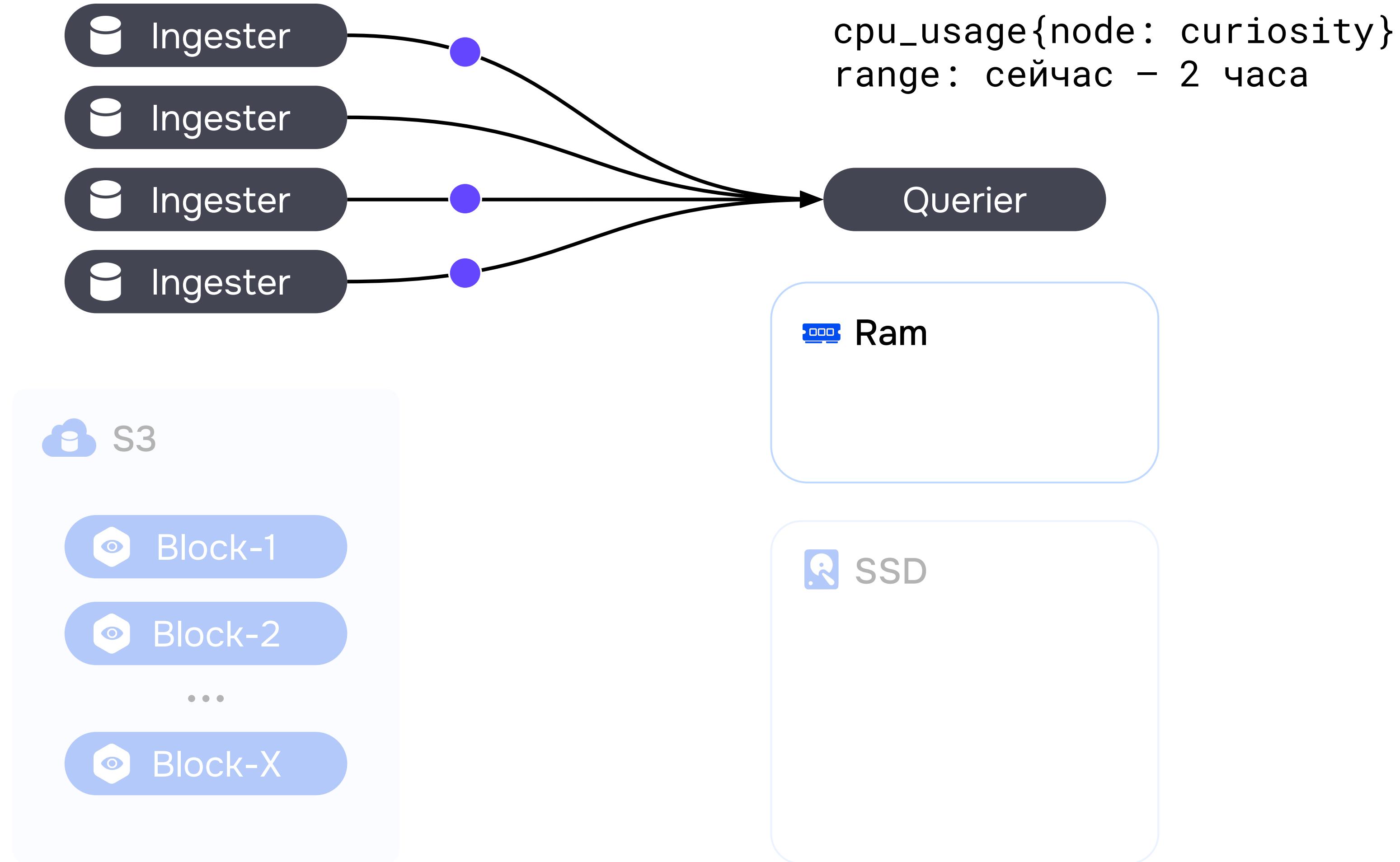


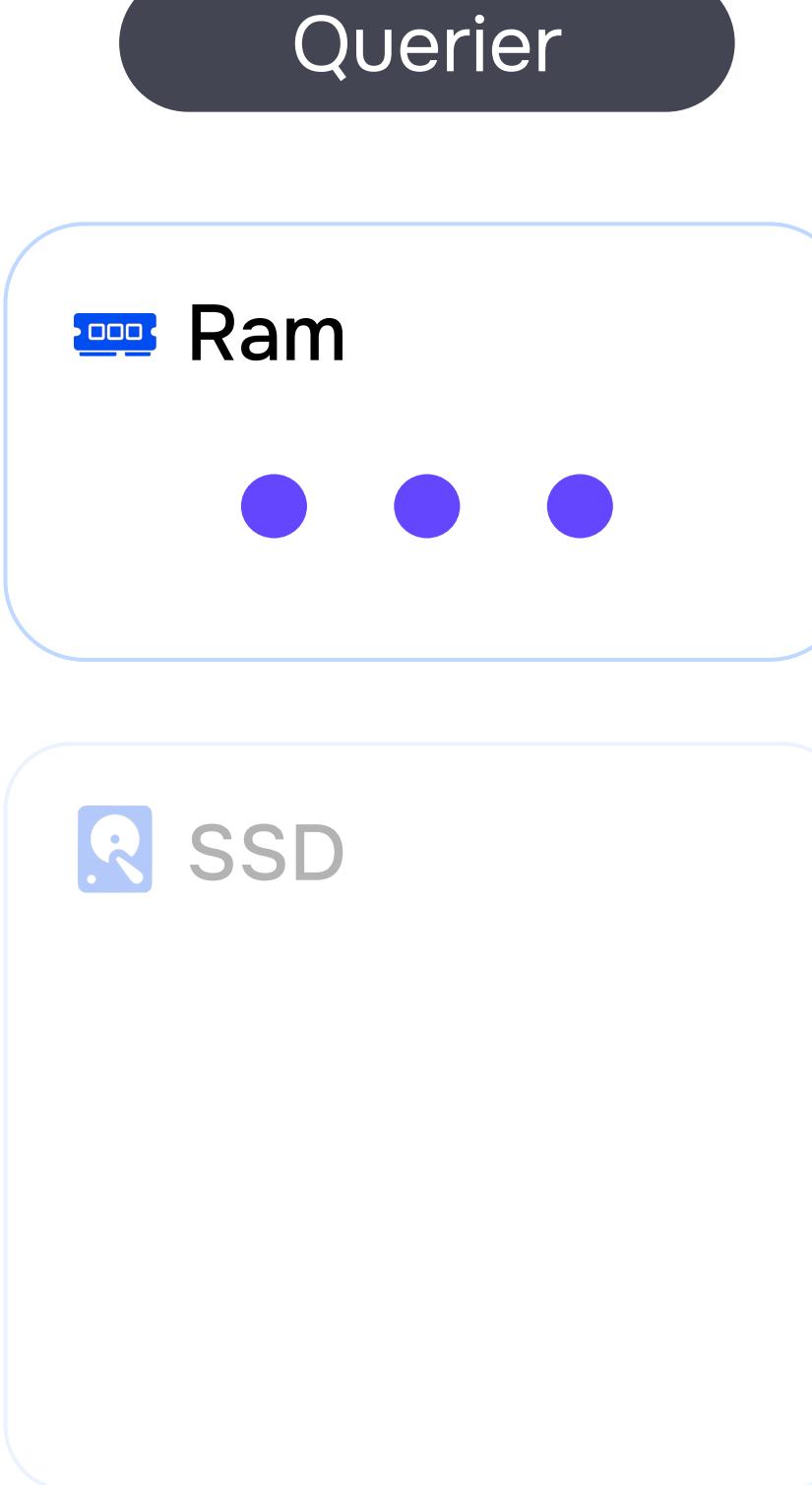
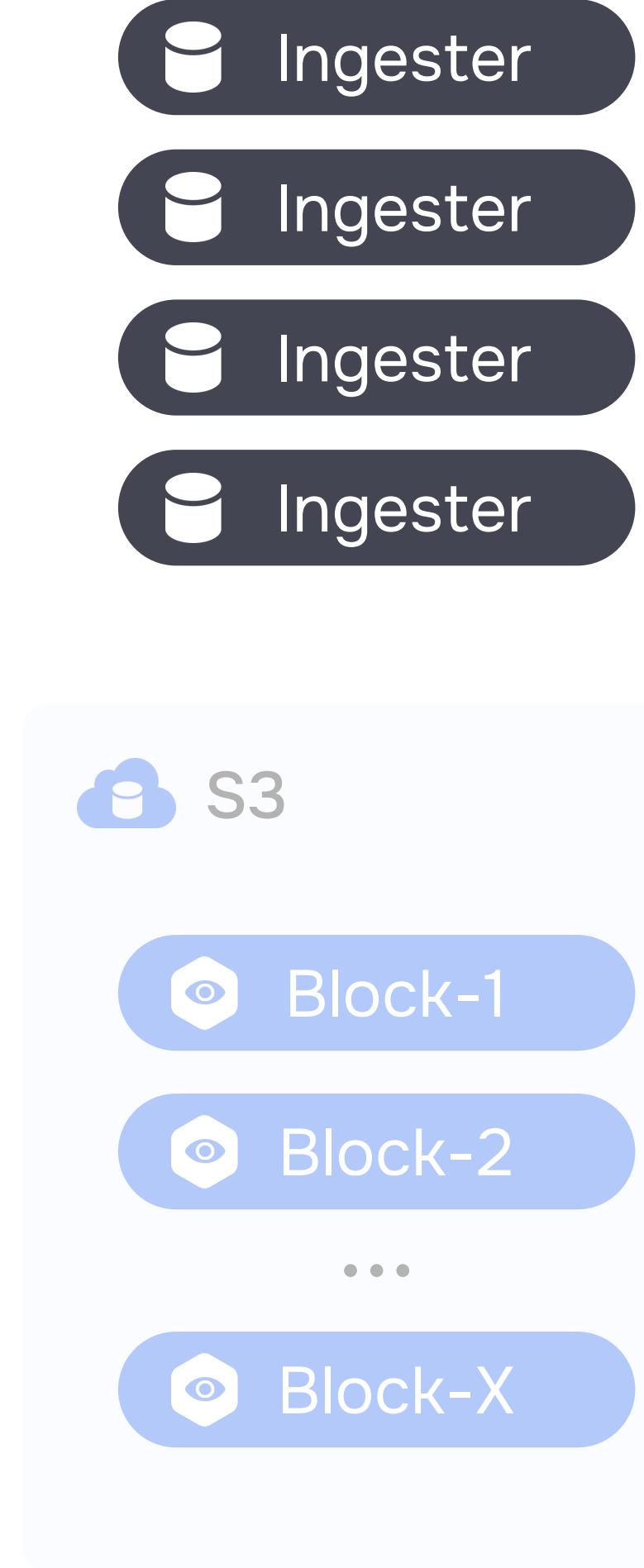
cpu_usage{node: curiosity}
range: сейчас – 2 часа

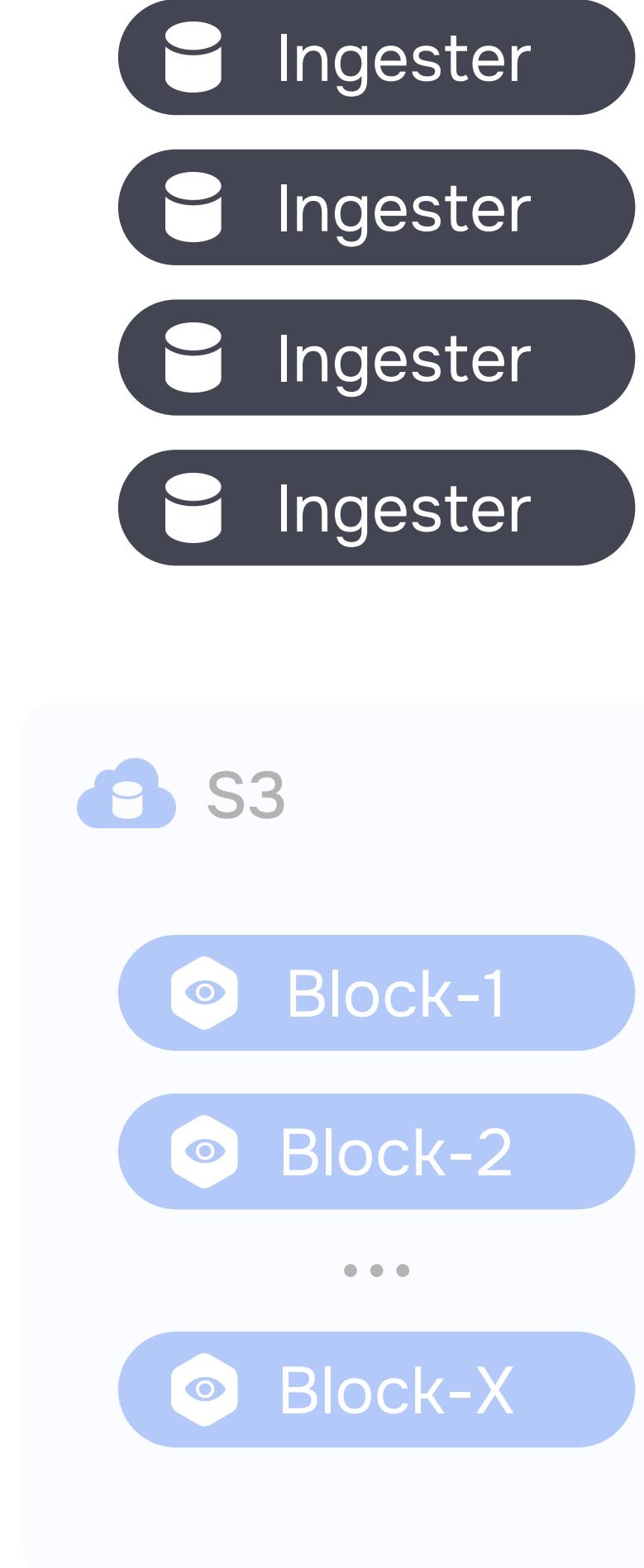




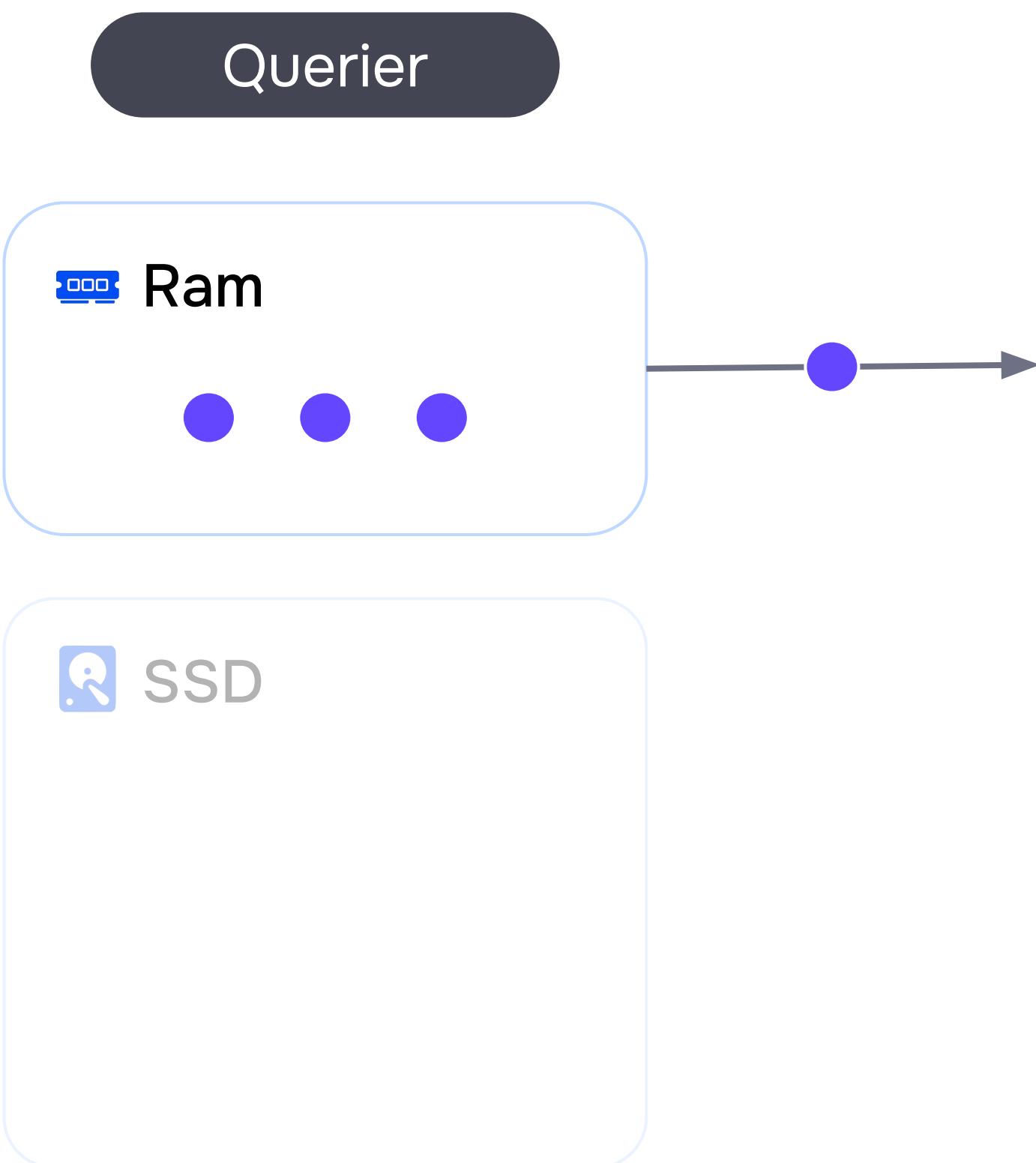


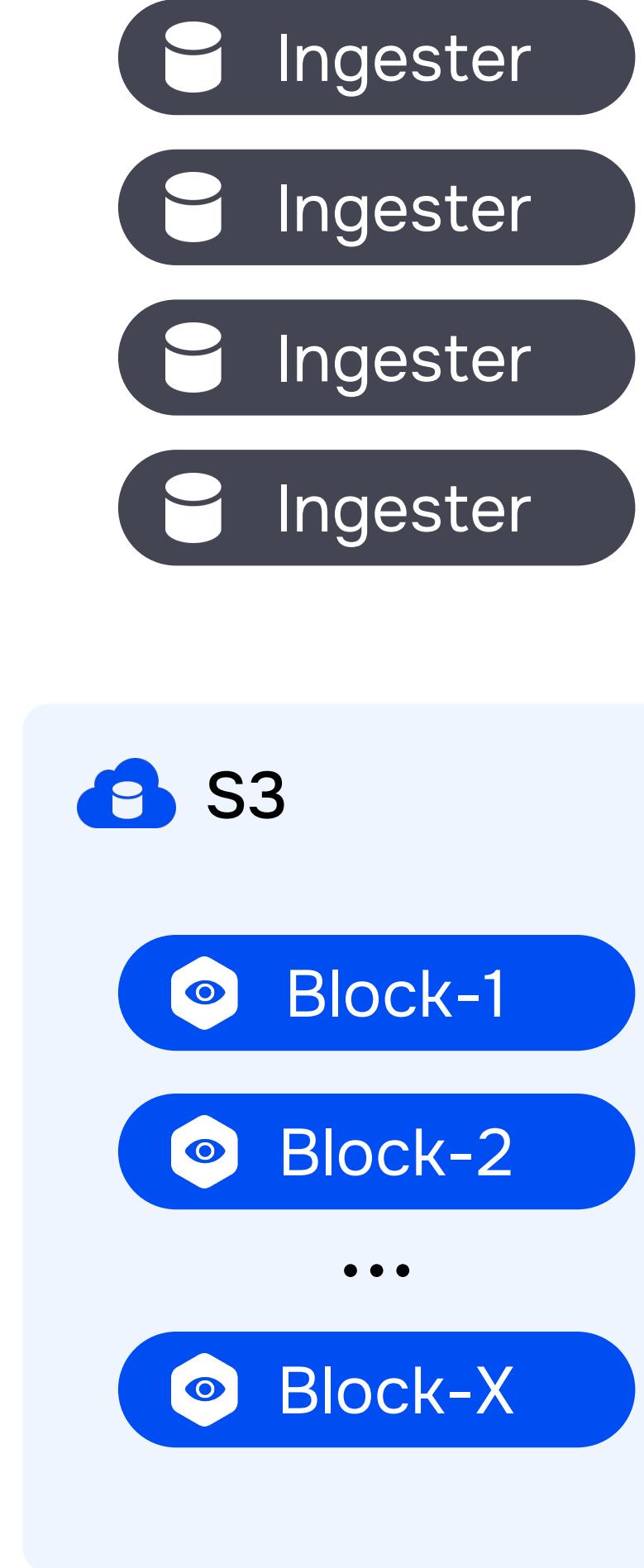




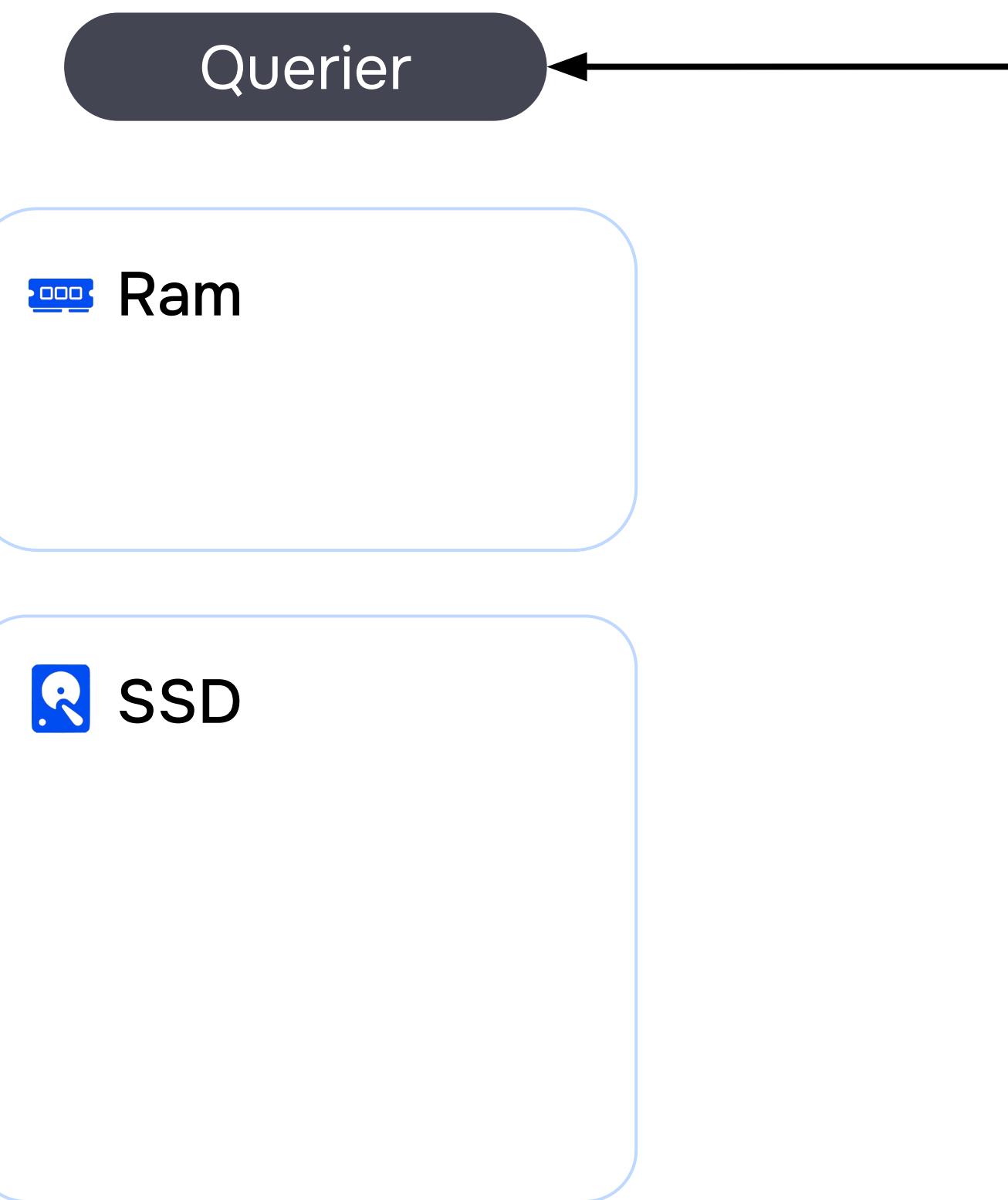


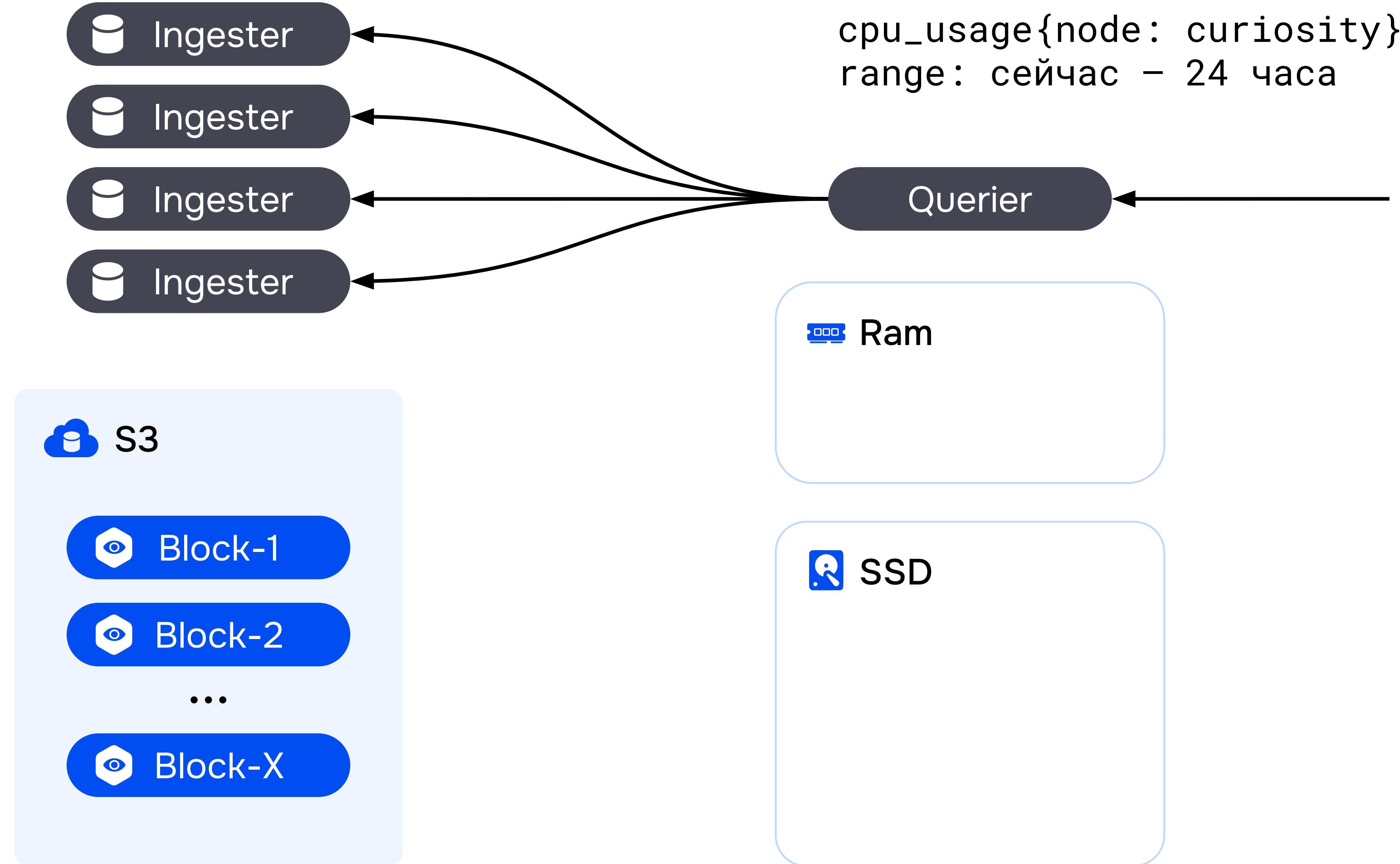
`cpu_usage{node: curiosity}`
range: сейчас – 2 часа

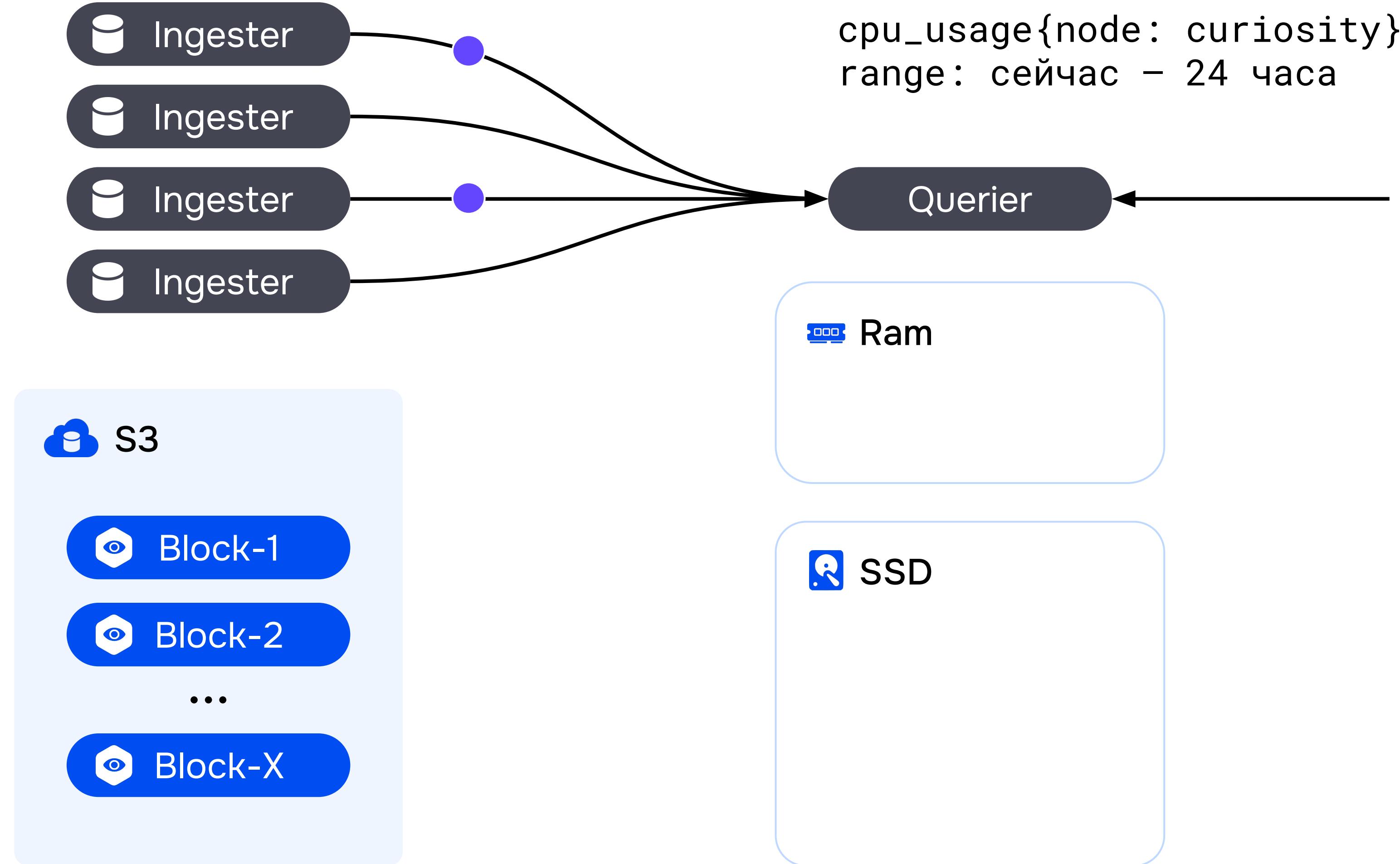


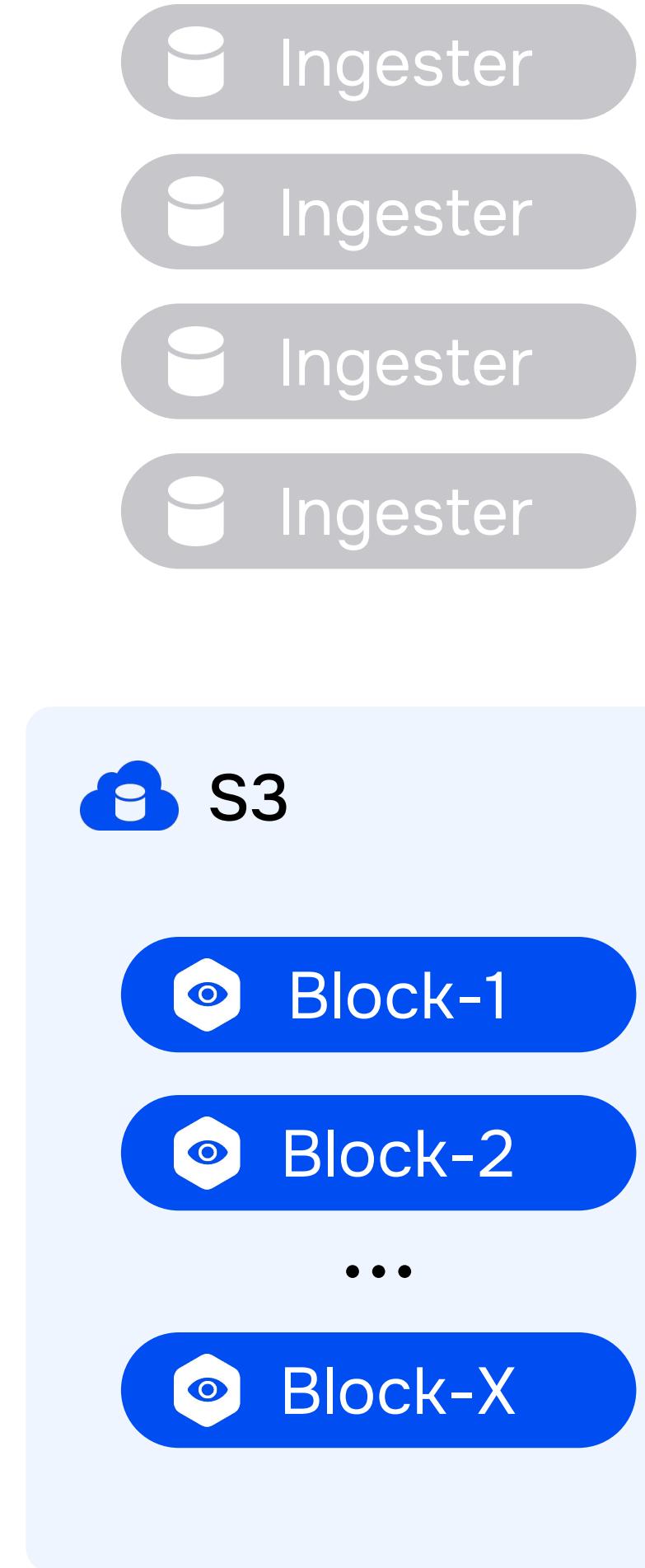


cpu_usage{node: curiosity}
range: сейчас – 24 часа

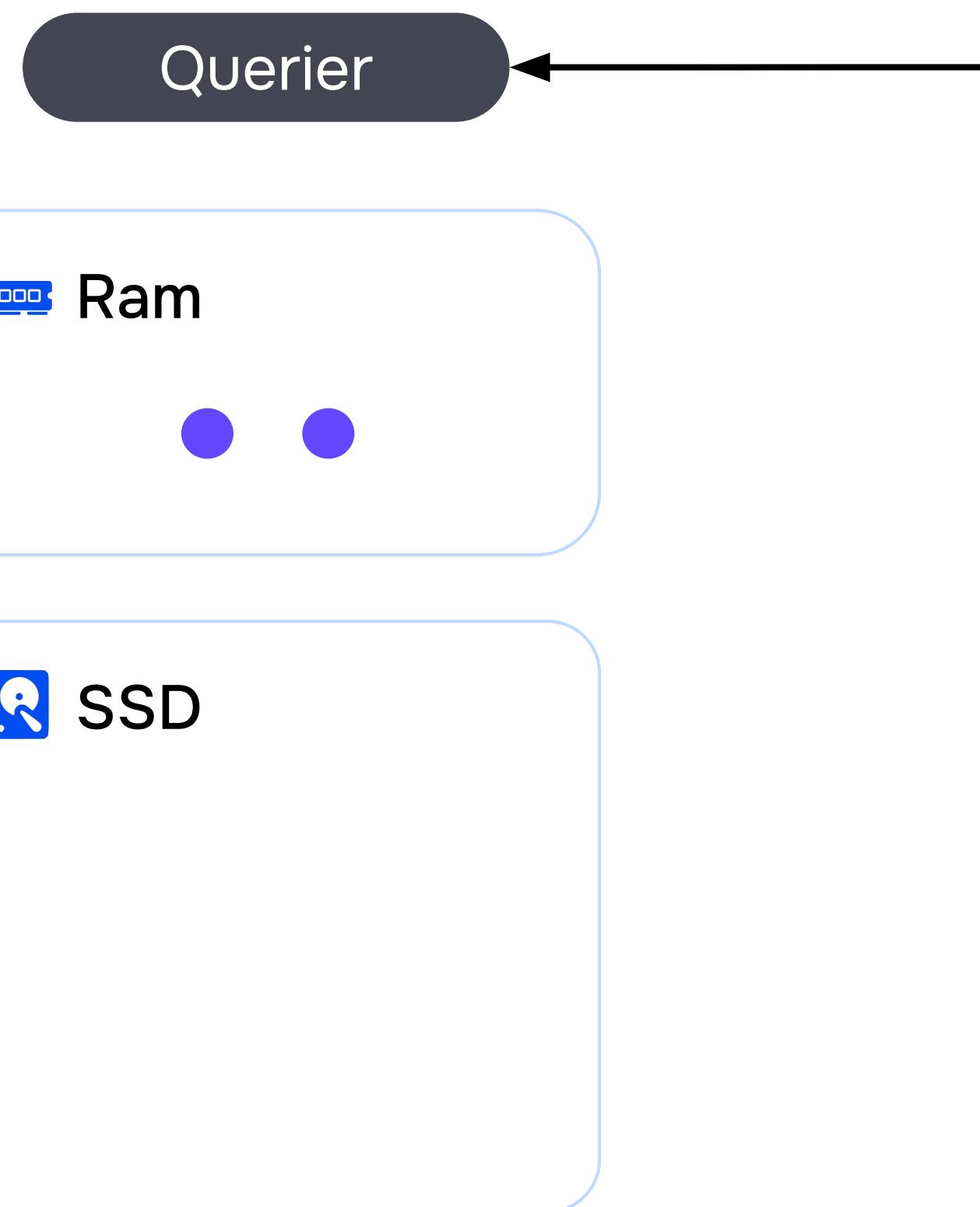


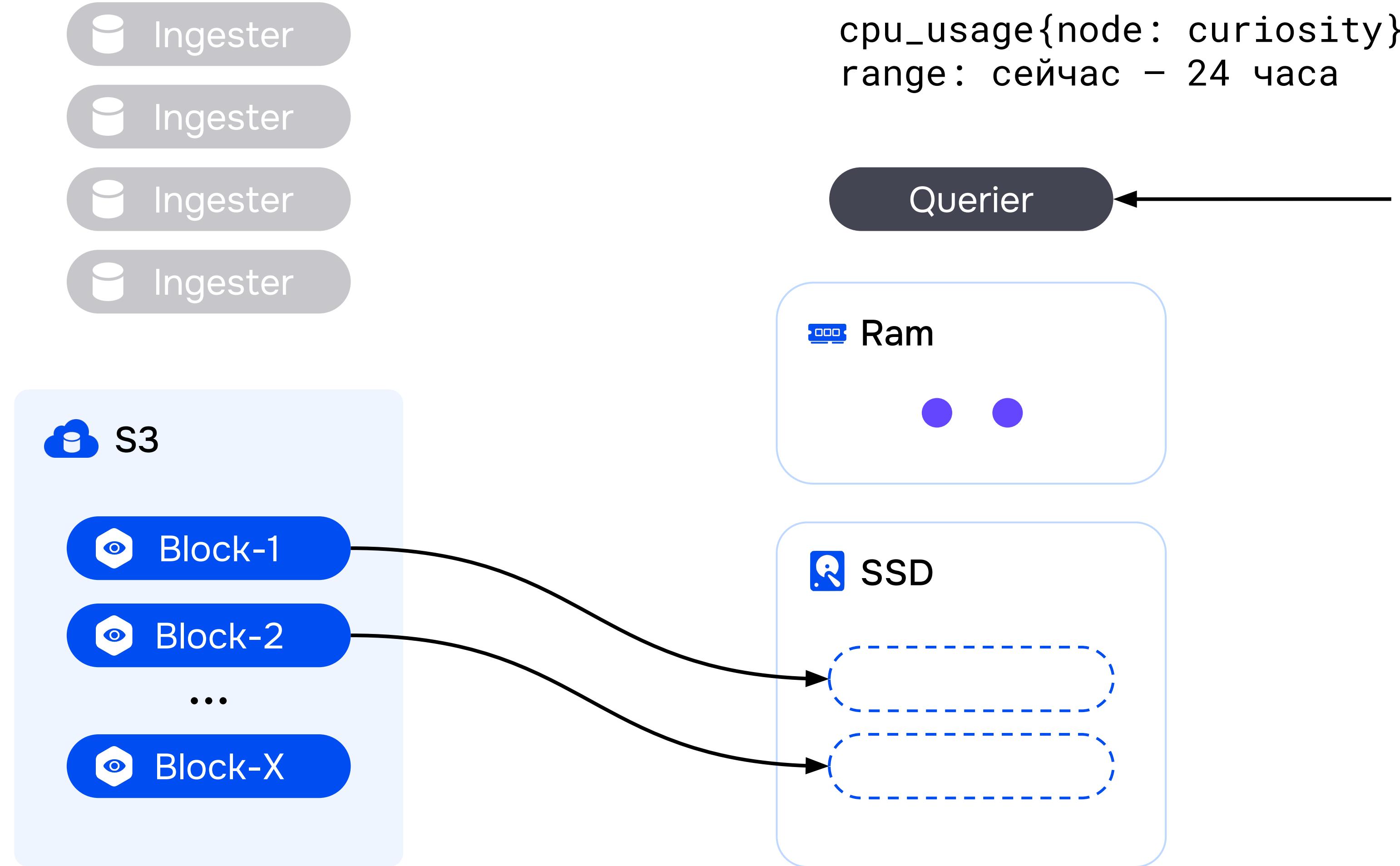


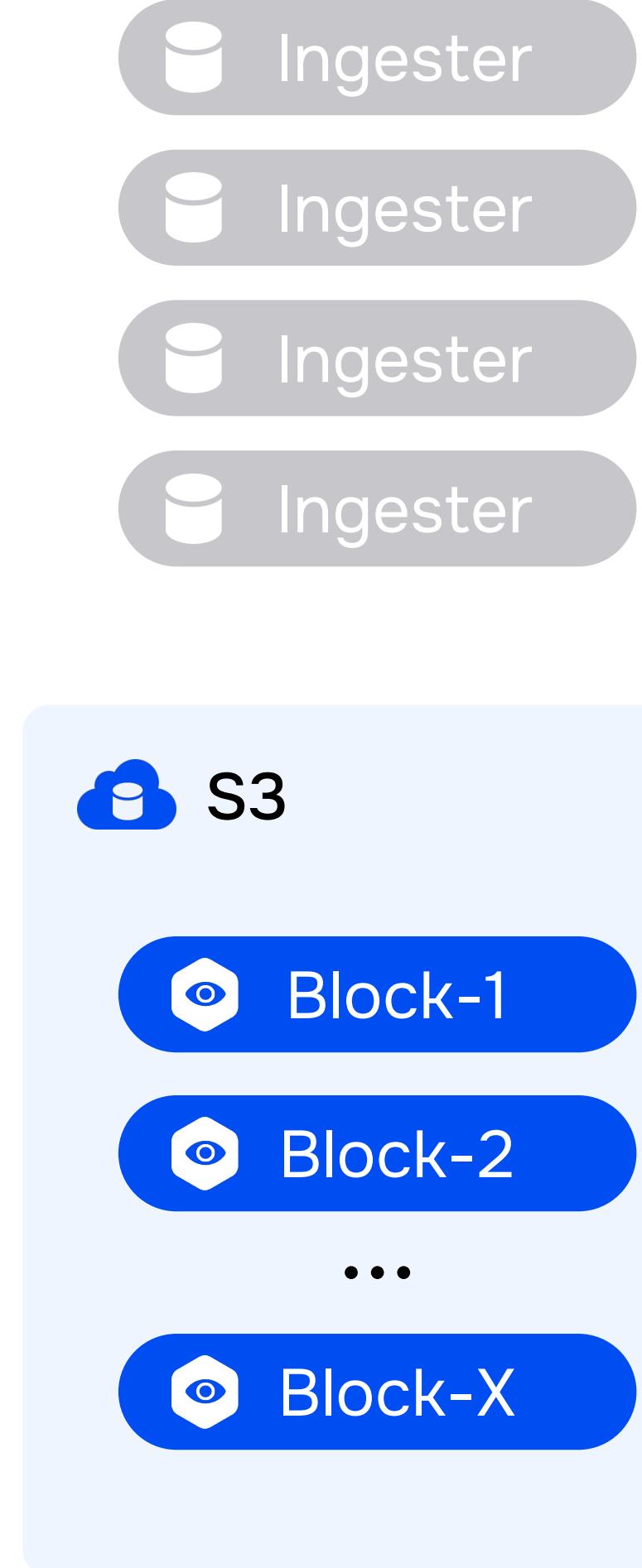




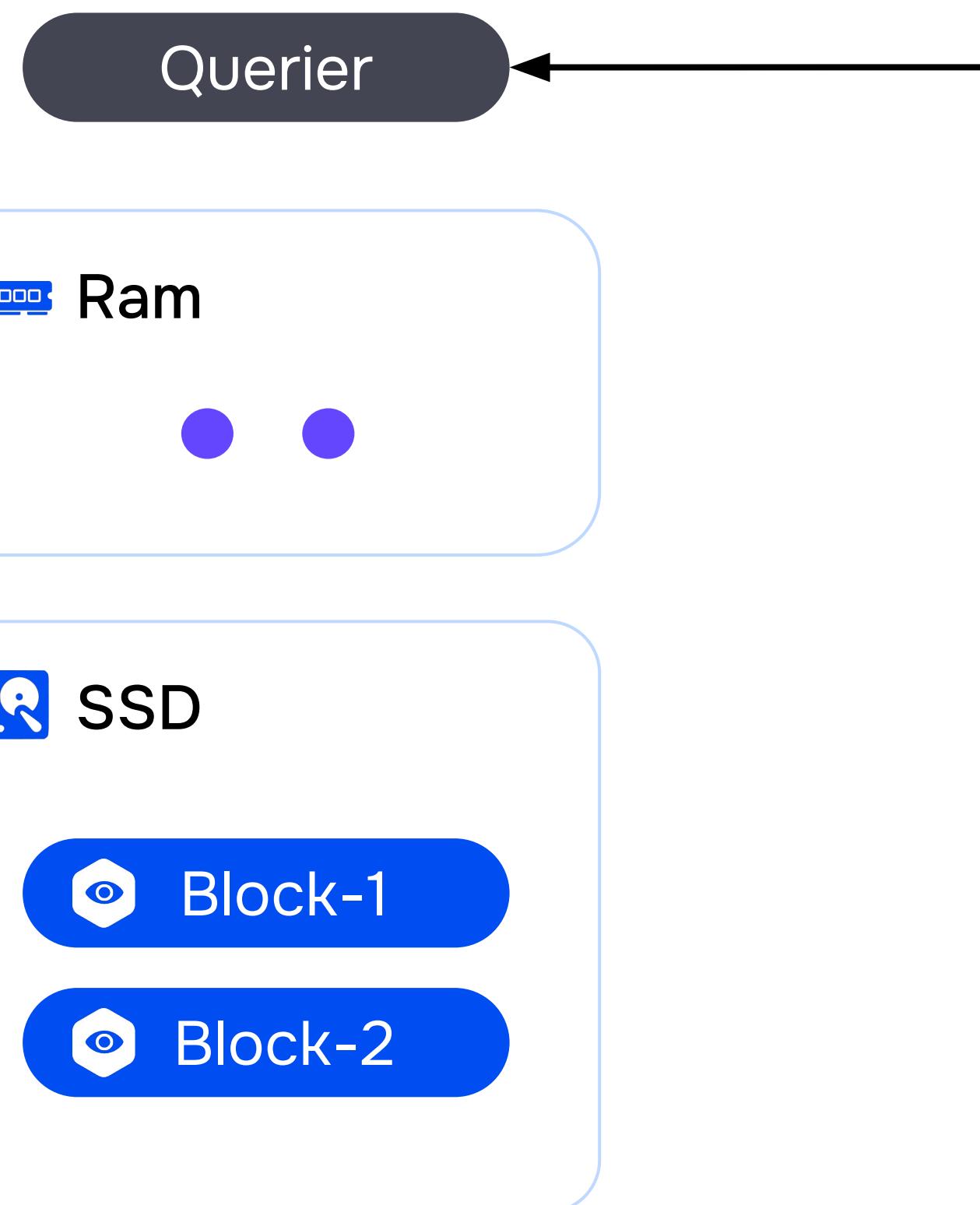
cpu_usage{node: curiosity}
range: сейчас – 24 часа

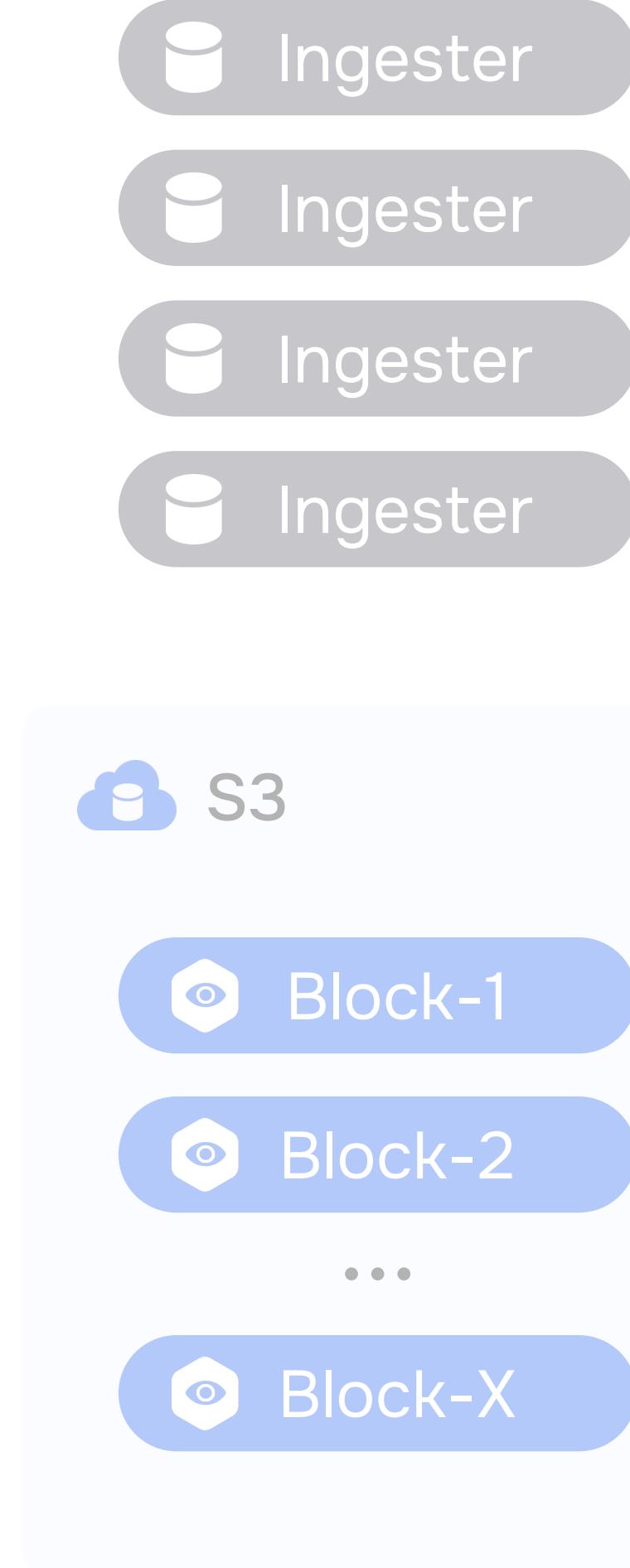




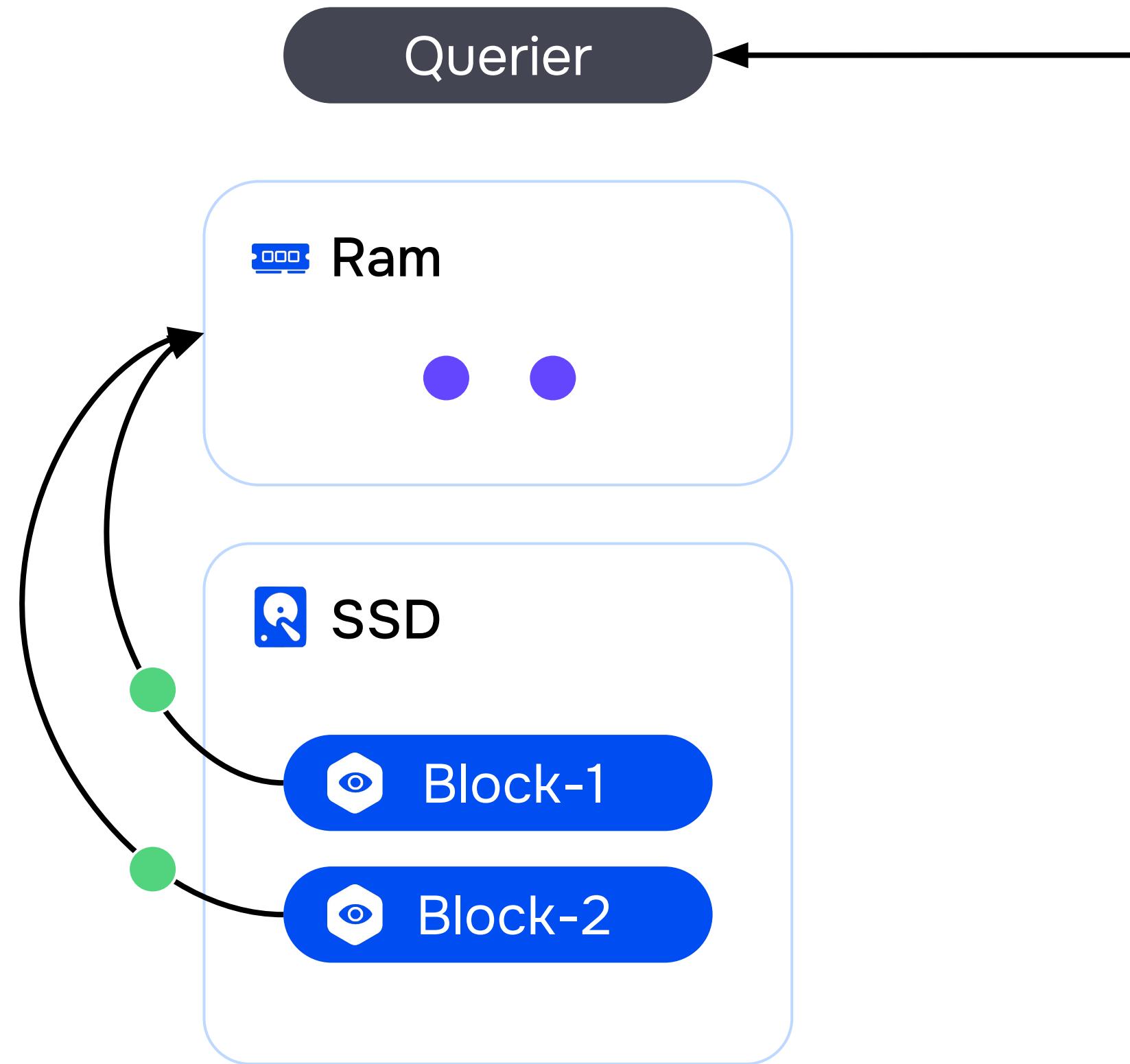


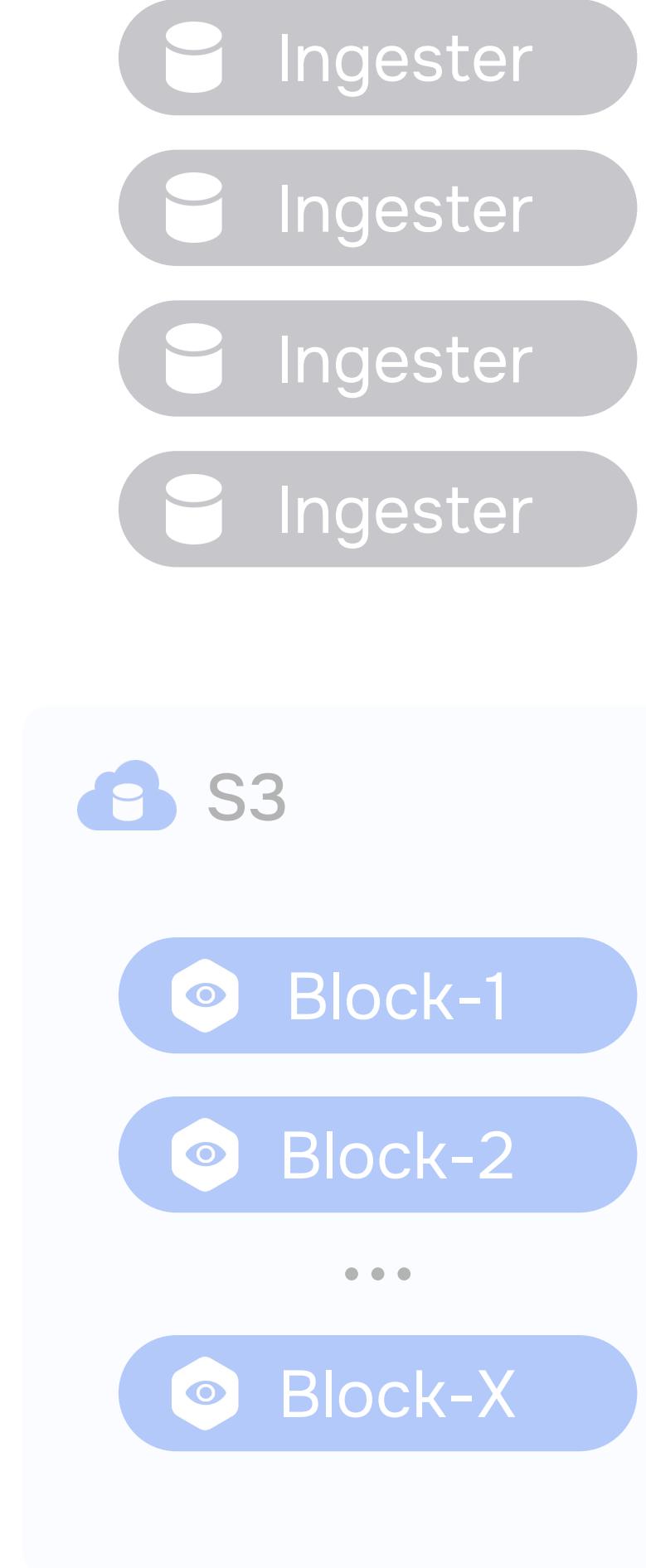
cpu_usage{node: curiosity}
range: сейчас – 24 часа



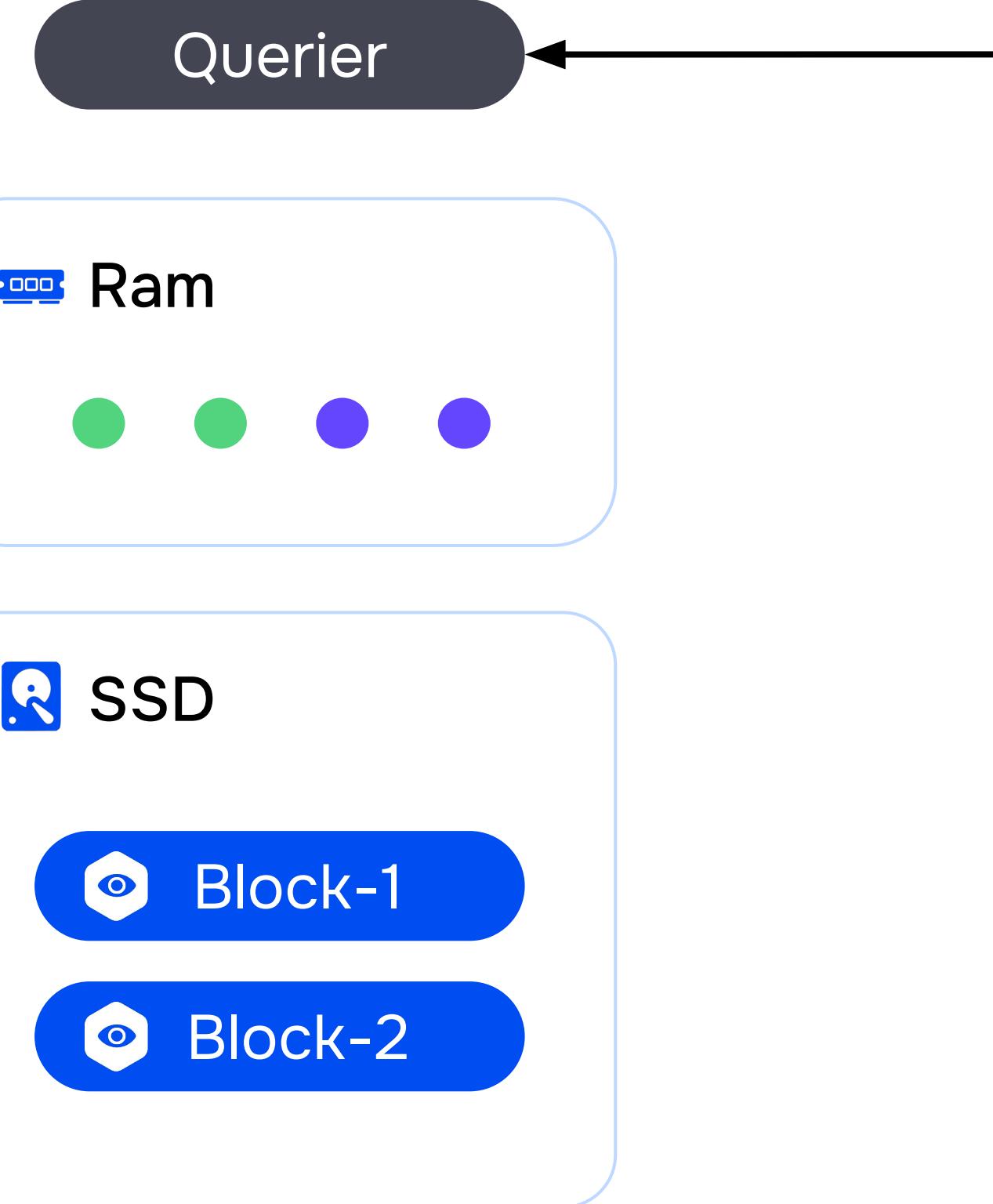


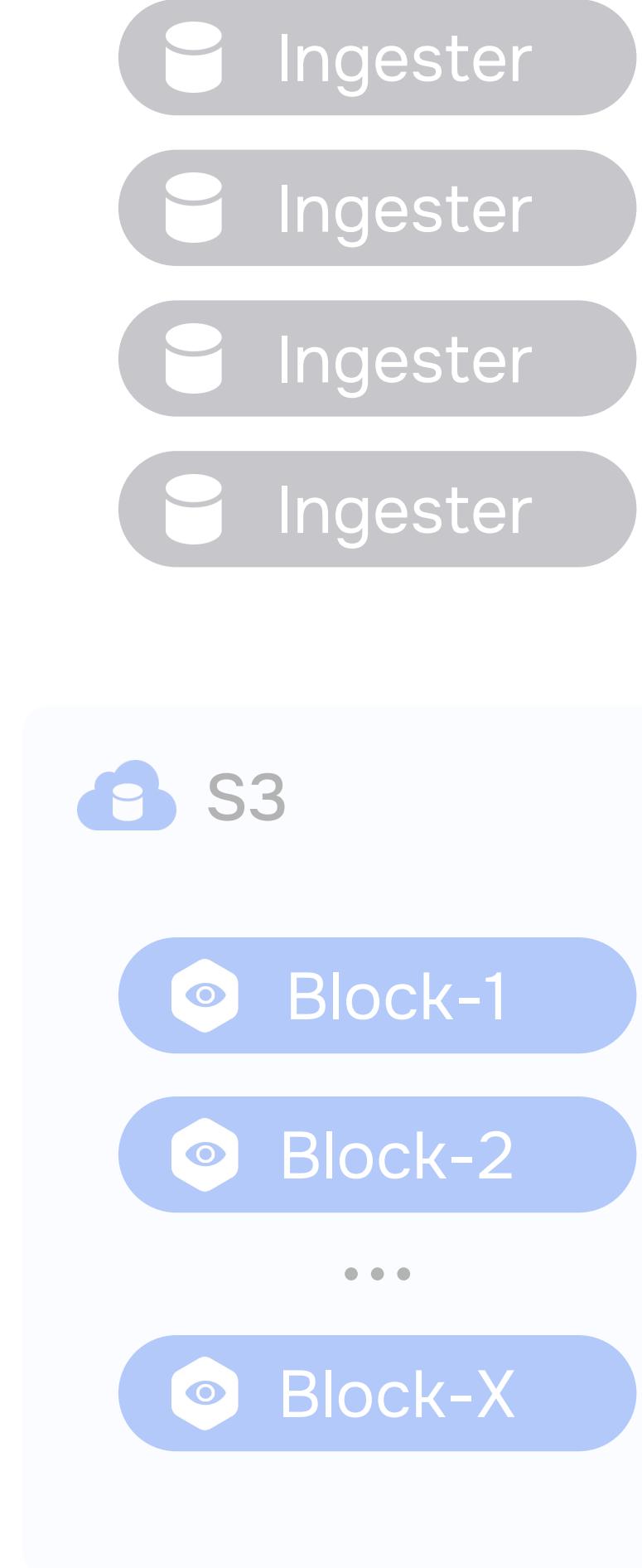
cpu_usage{node: curiosity}
range: сейчас – 24 часа



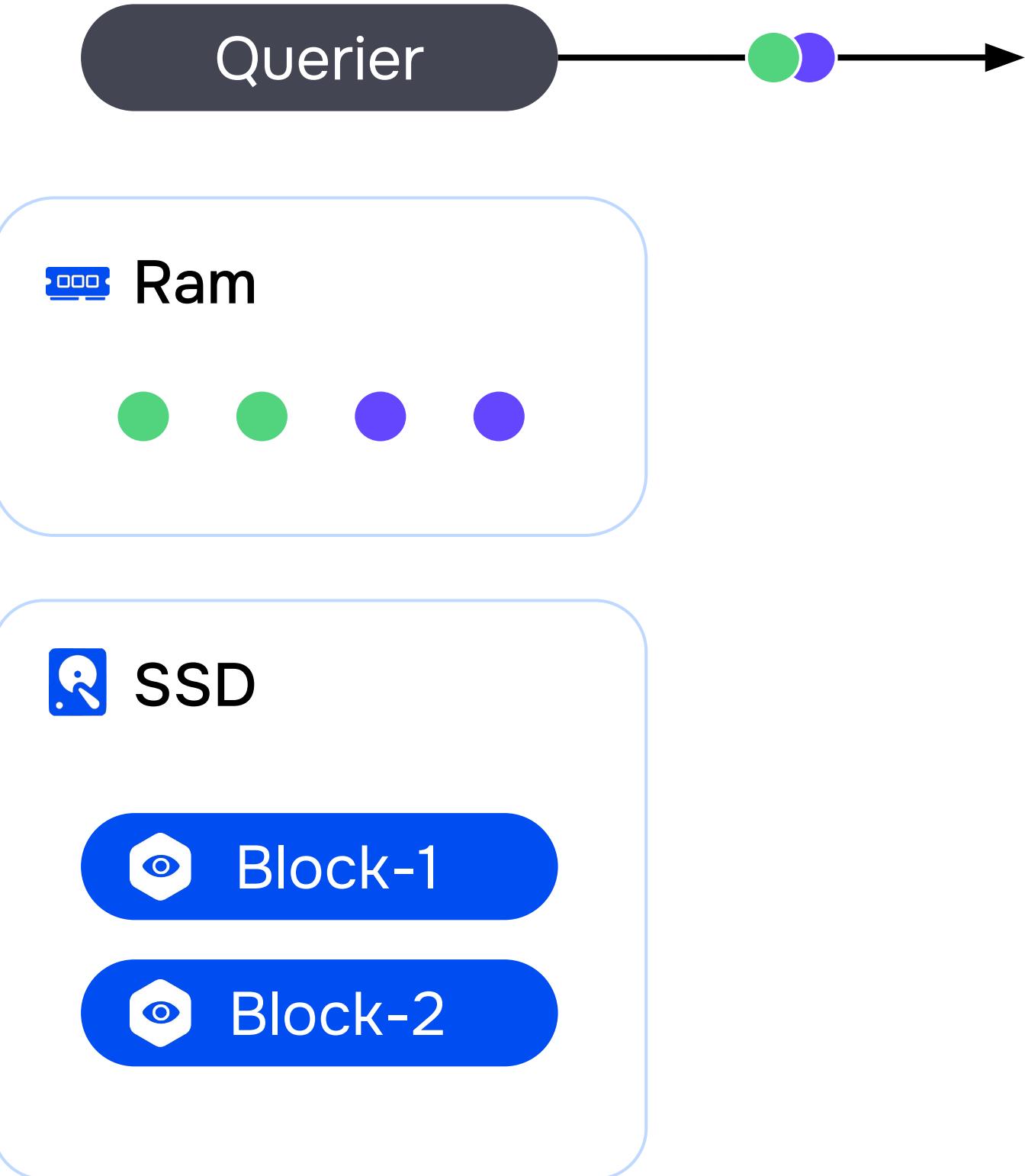


`cpu_usage{node: curiosity}`
range: сейчас – 24 часа

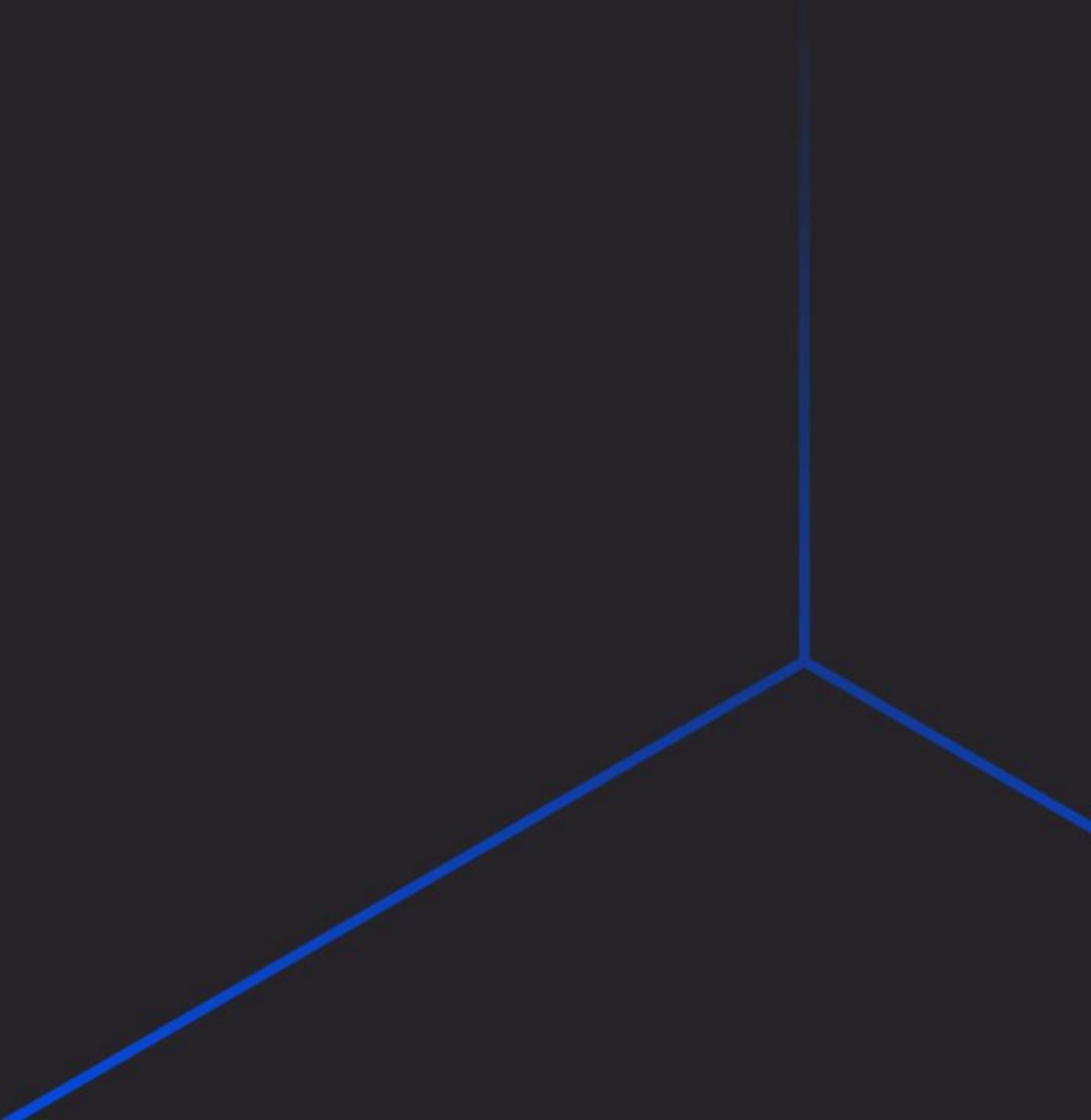


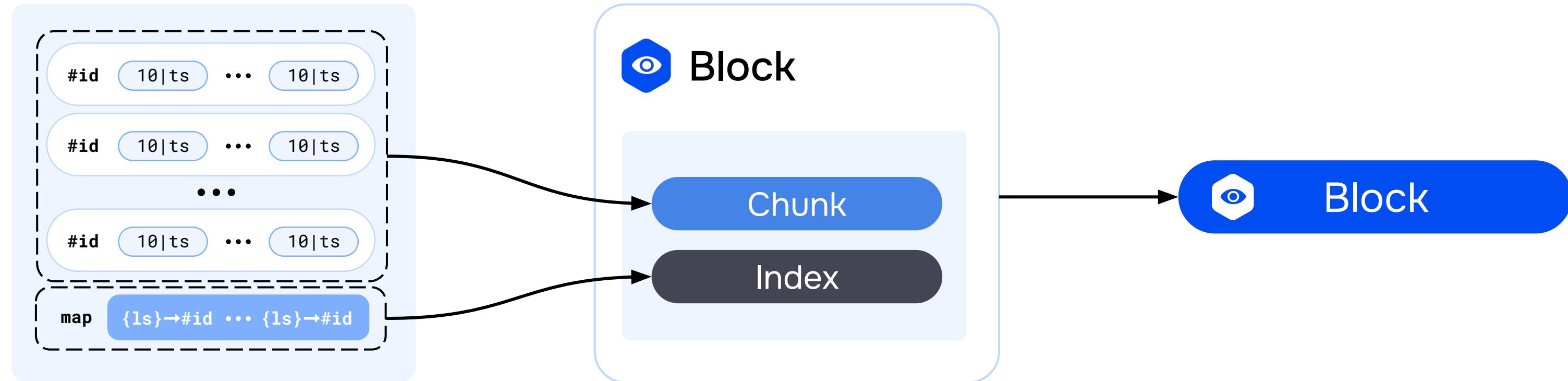


cpu_usage{node: curiosity}
range: сейчас – 24 часа



А дешевле можно?





Querier



Block -1

#id 10|ts ... 10|ts

#id 10|ts ... 10|ts

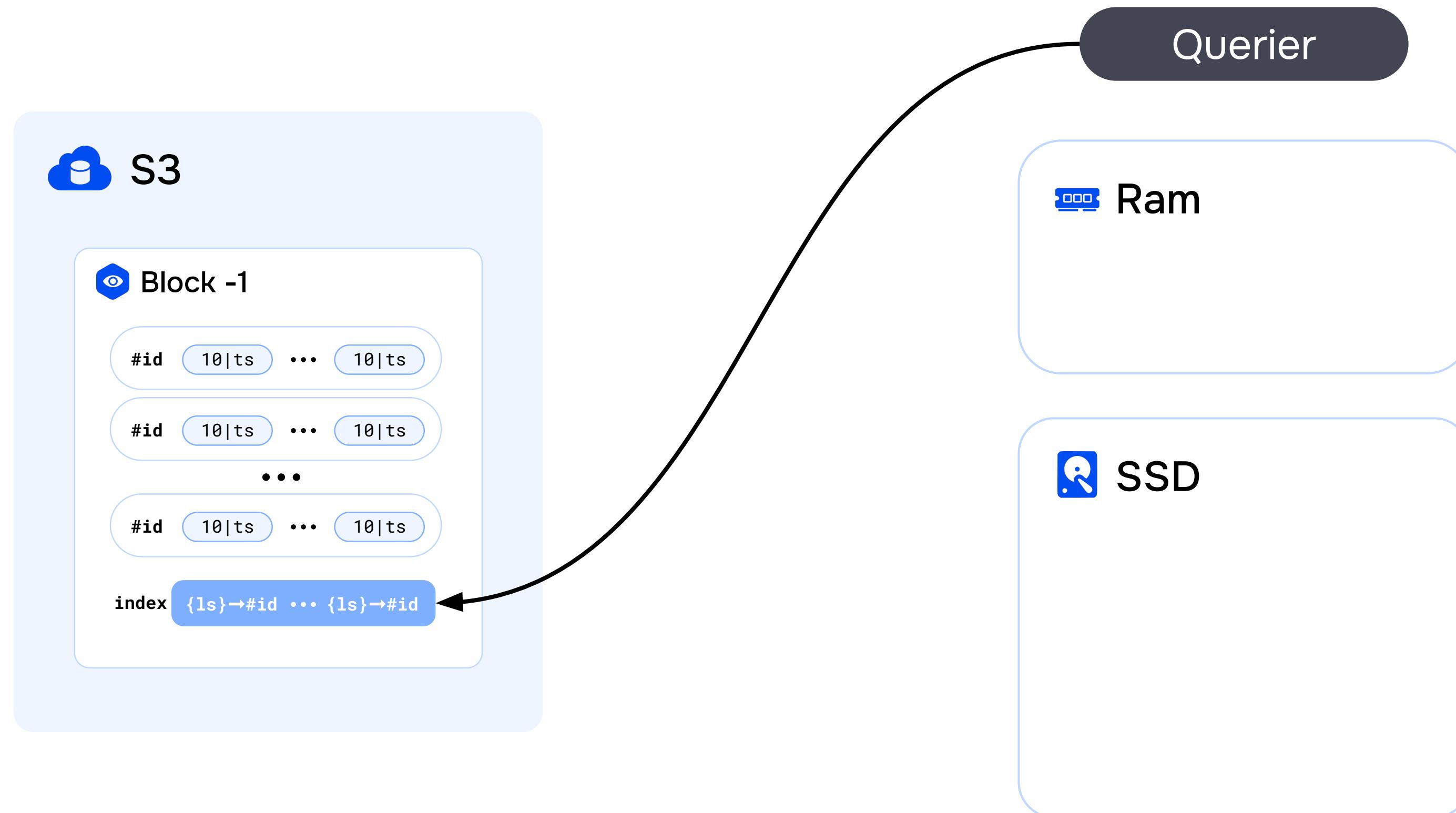
...

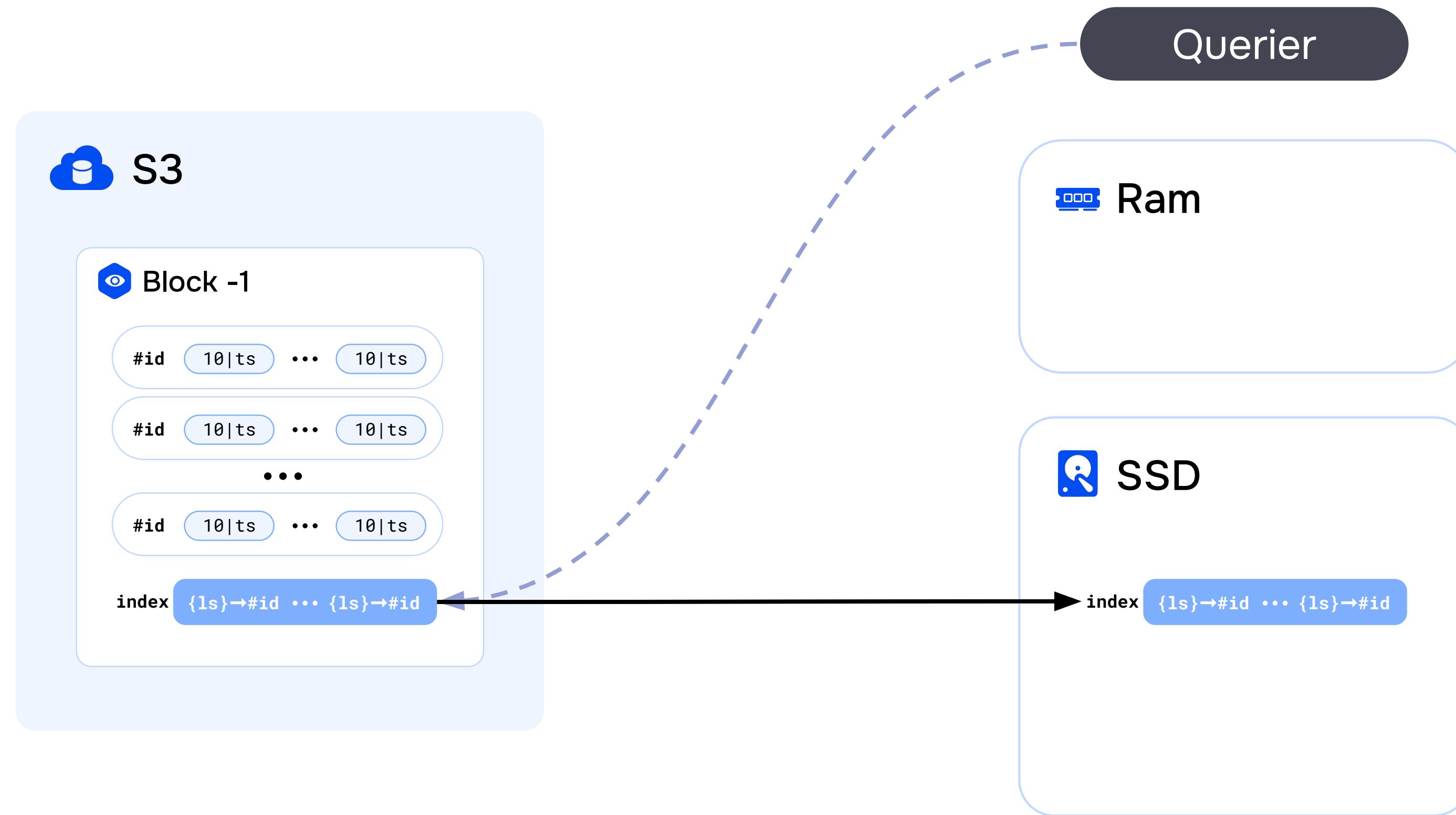
#id 10|ts ... 10|ts

index {ls}→#id ... {ls}→#id

Ram

SSD





Querier



Block -1

#id 10|ts ... 10|ts

#id 10|ts ... 10|ts

...

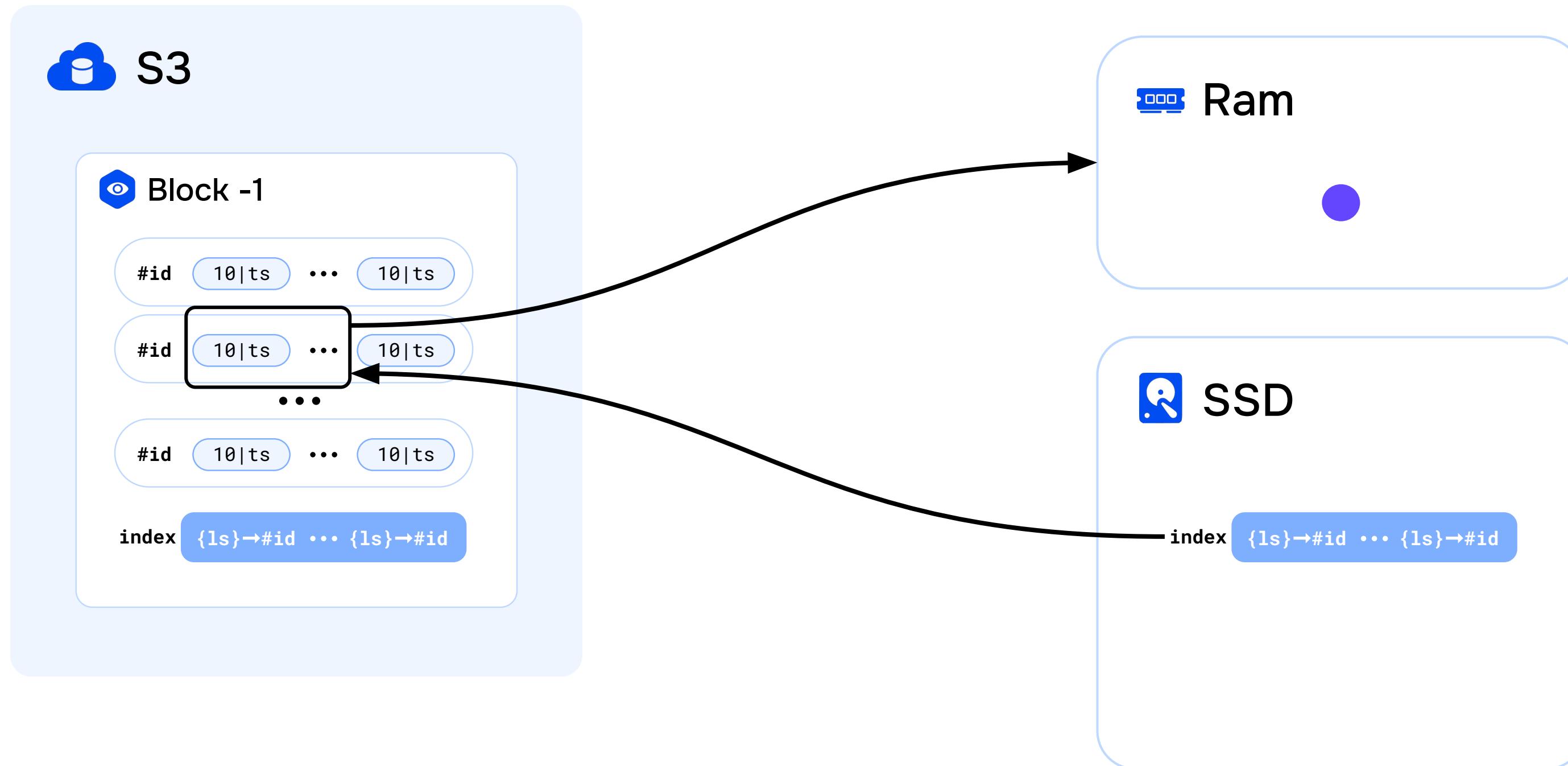
#id 10|ts ... 10|ts

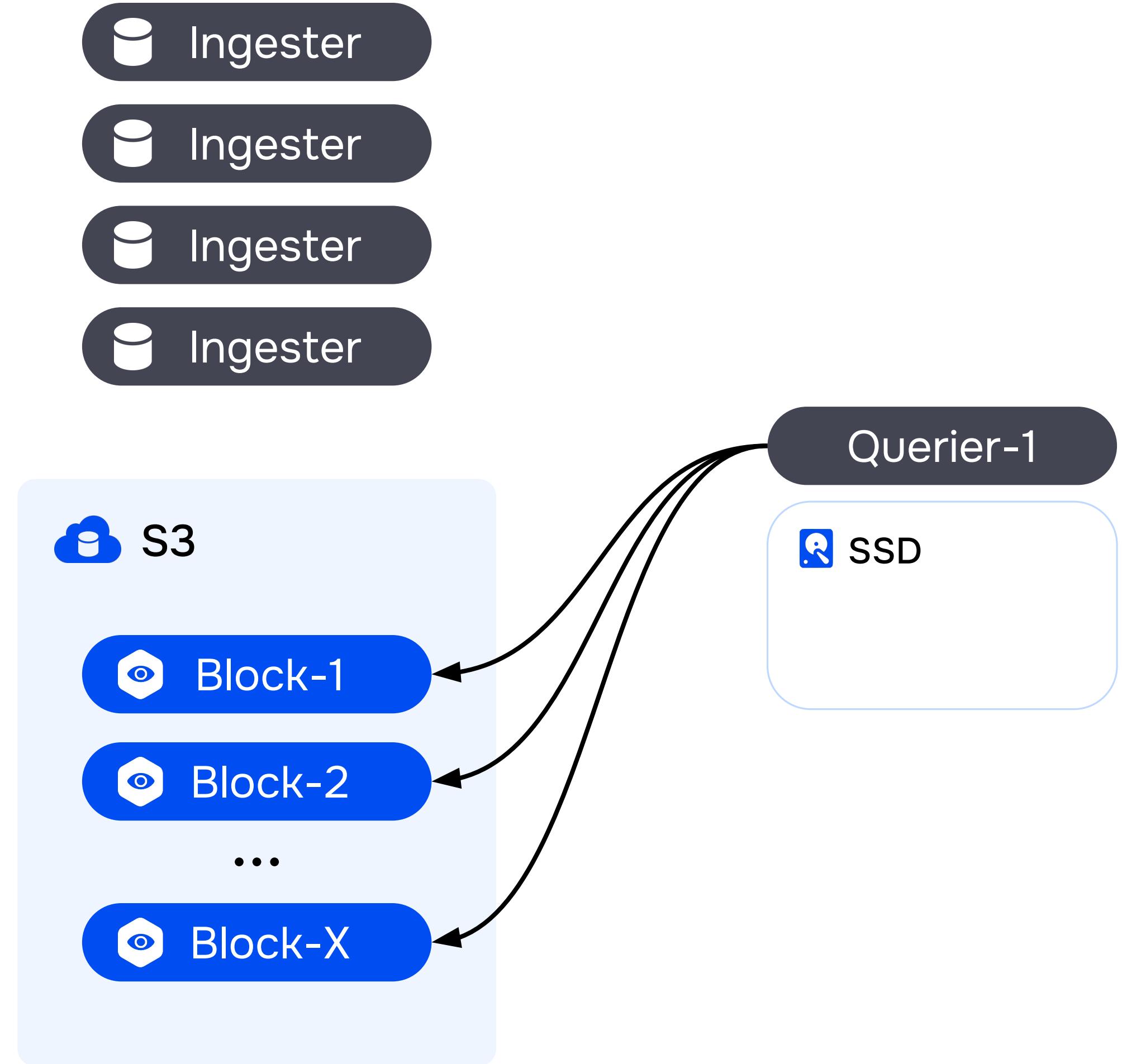
index {ls}→#id ... {ls}→#id



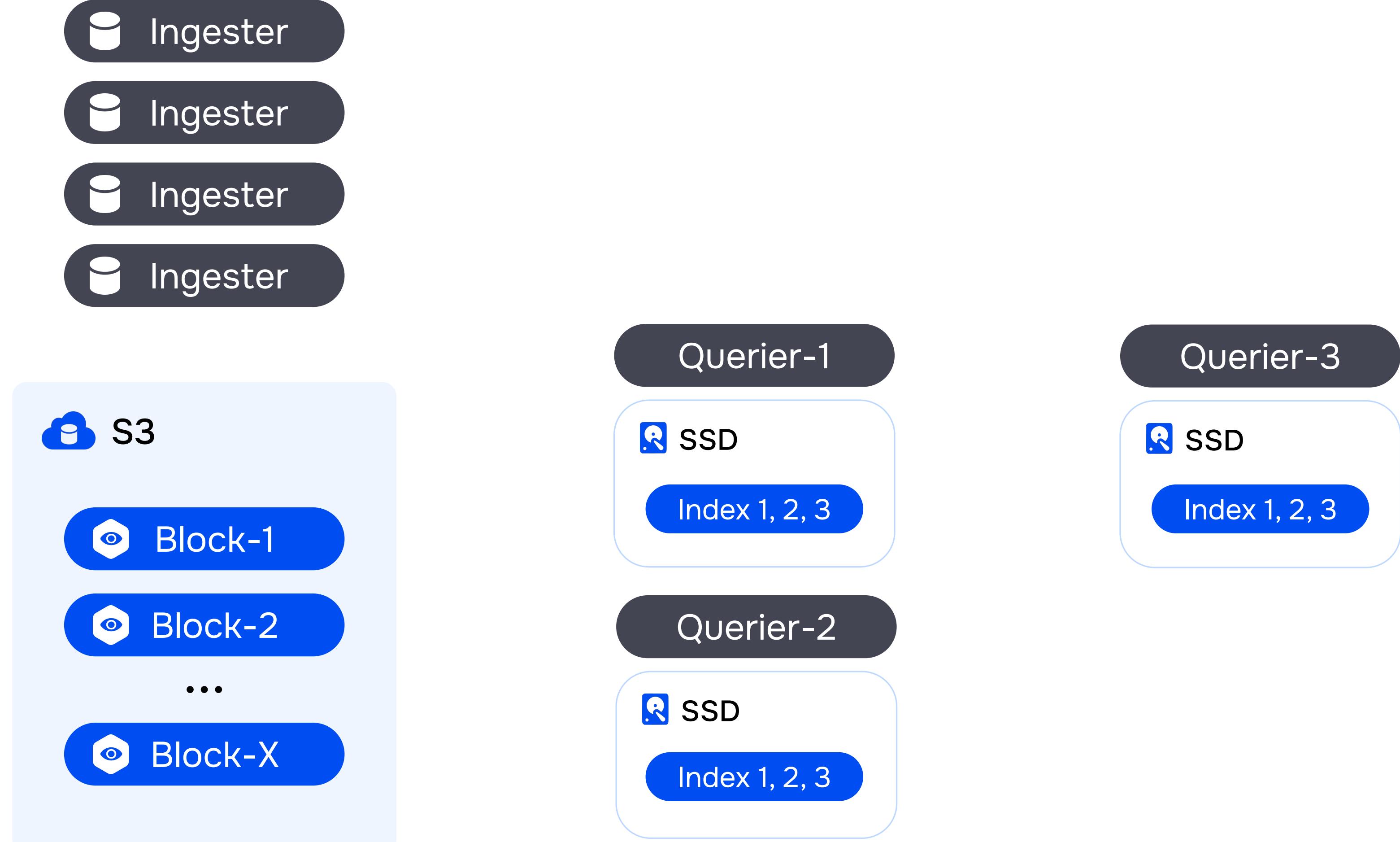
index {ls}→#id ... {ls}→#id

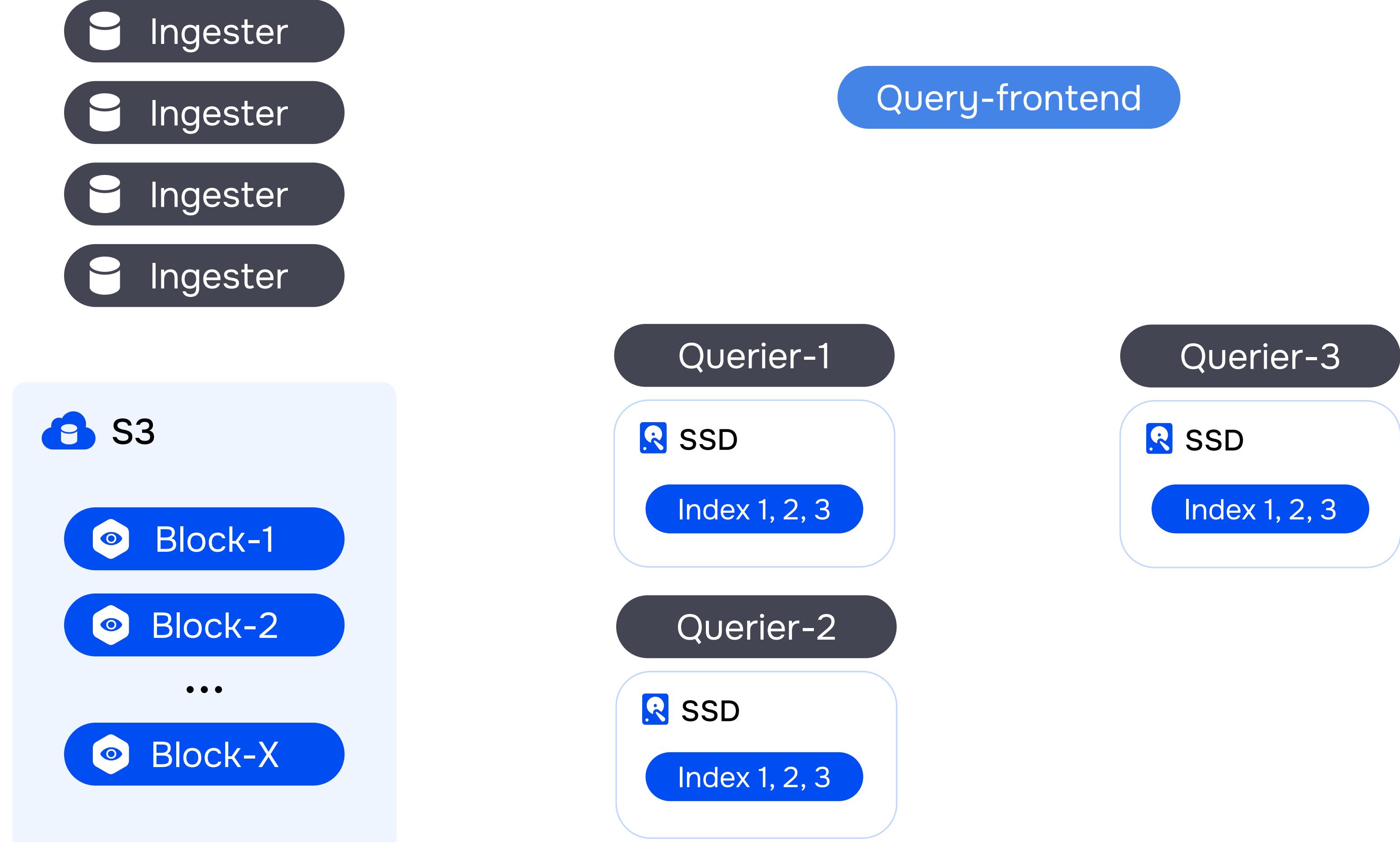
Querier

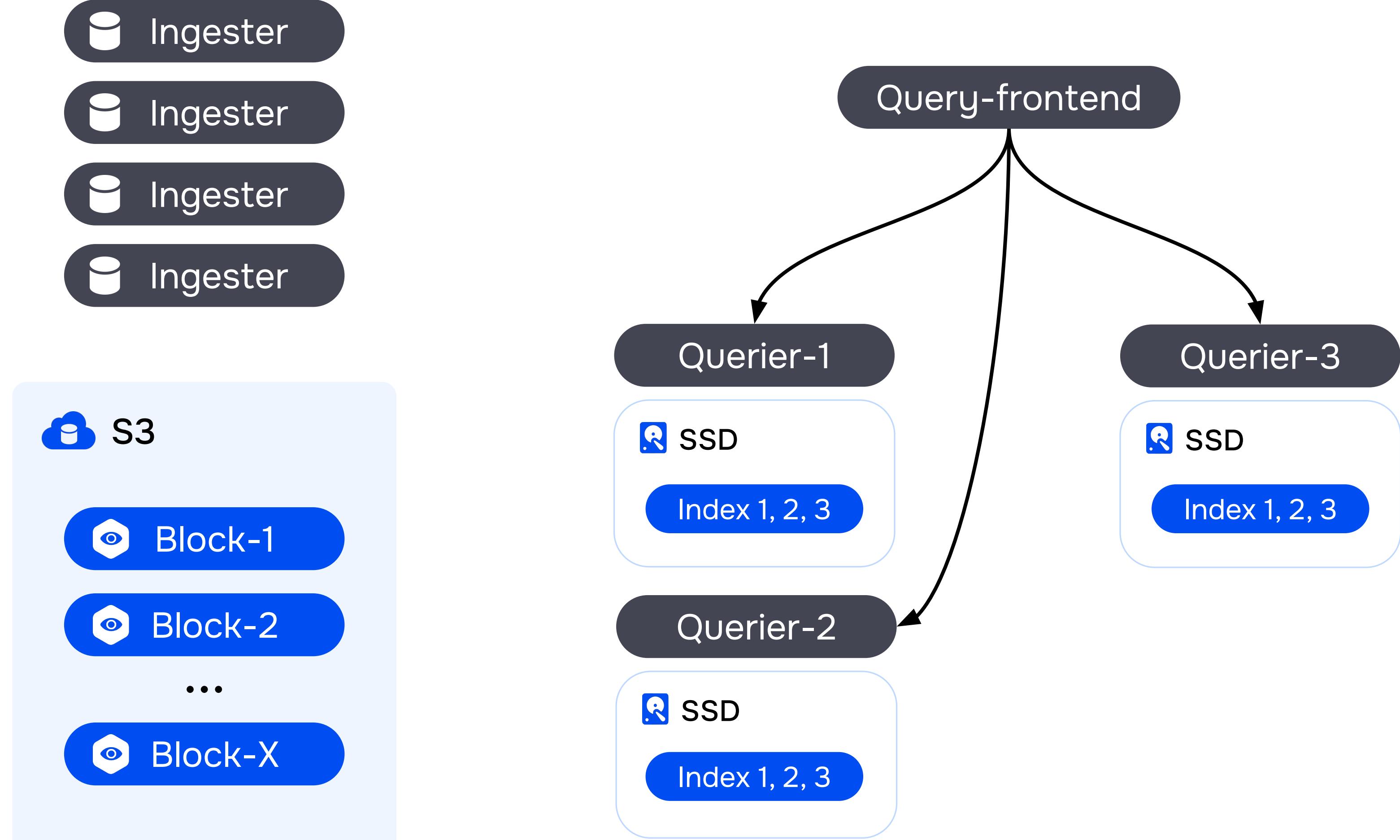






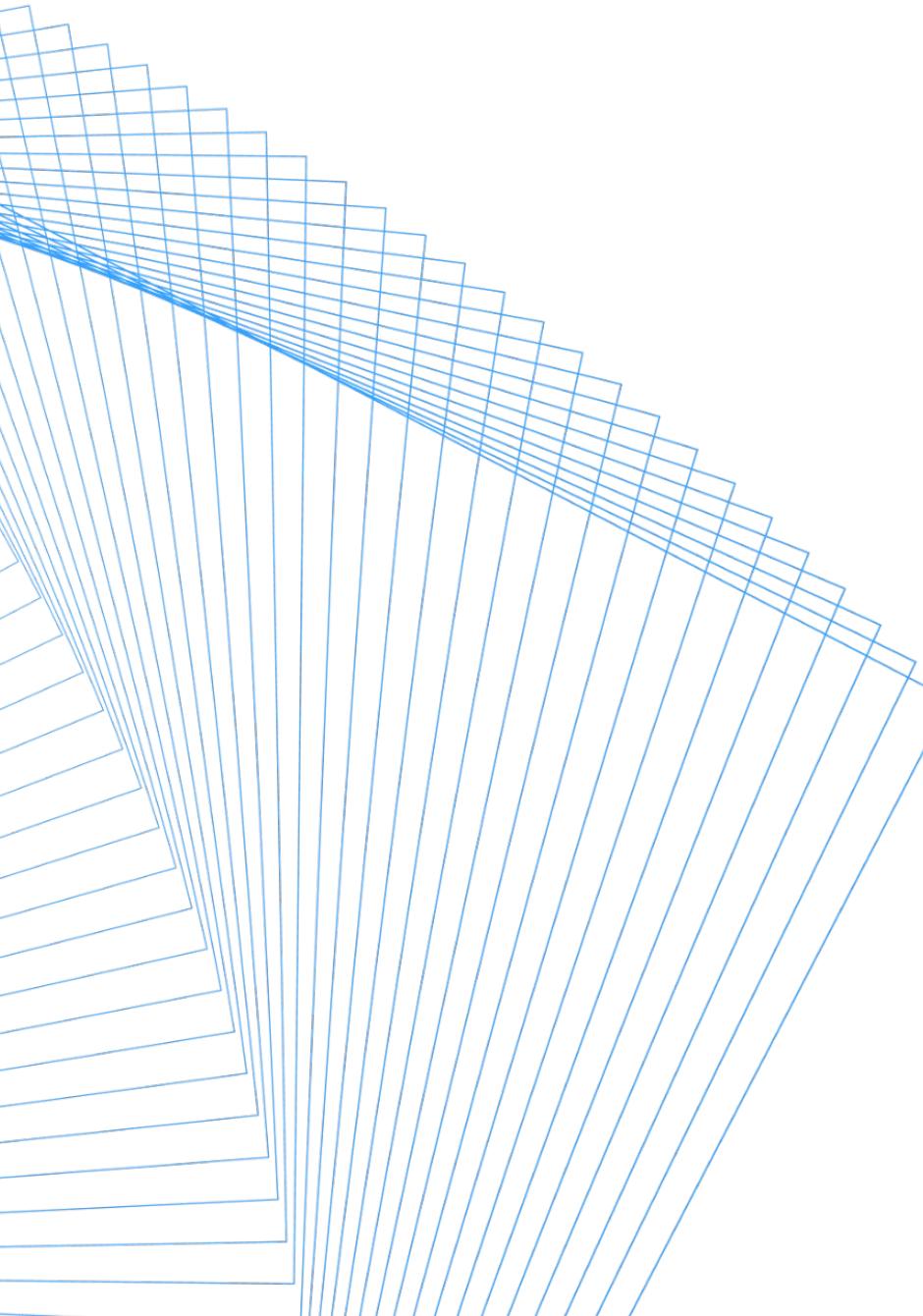






Централизованная система мониторинга

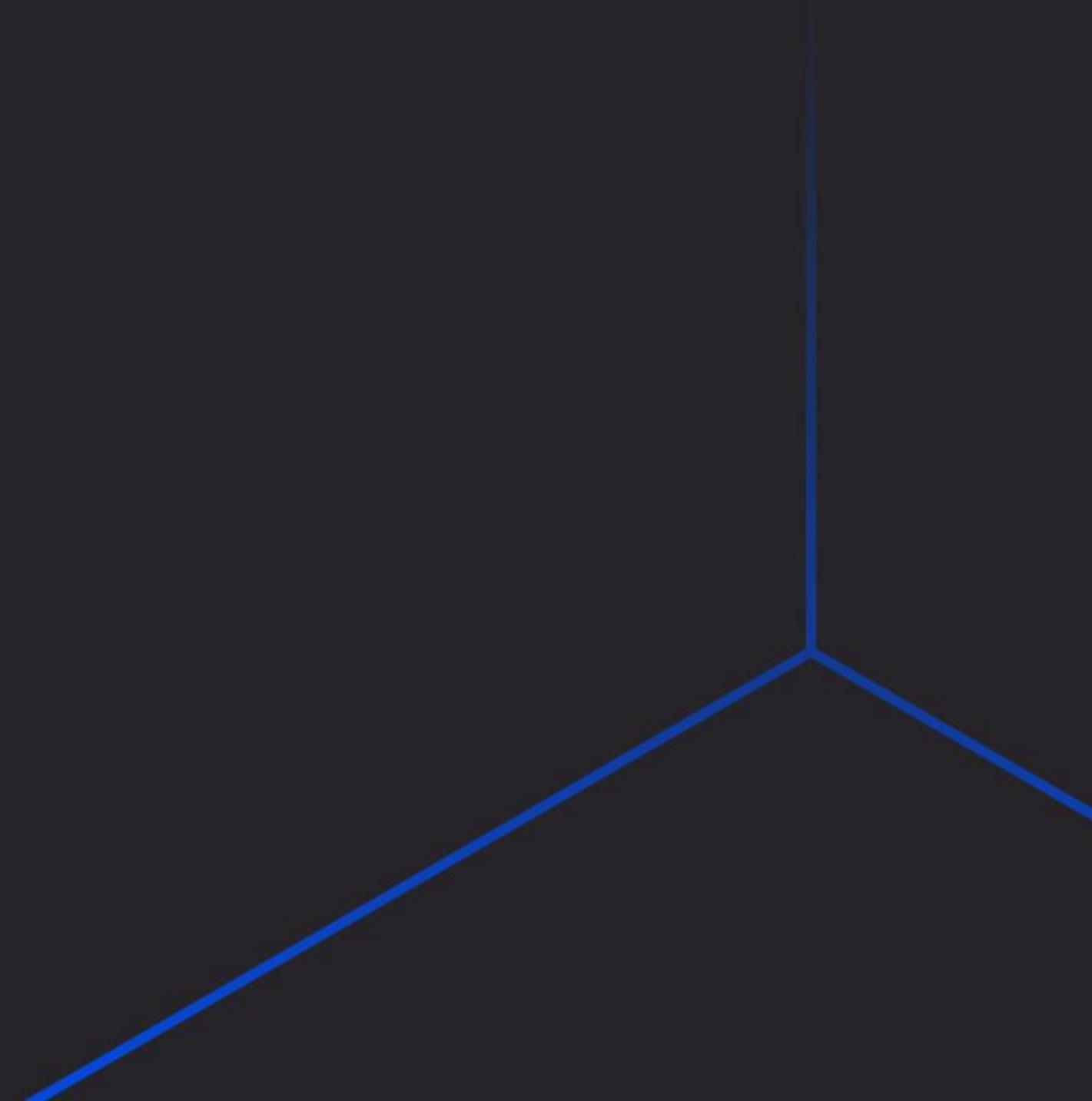
- Отказоустойчивость на запись: $WC = RF/2 + 1$
- Отказоустойчивость на чтение: $RC = RF - WC + 1$
- Шардирование на запись на основе хеш-ринга
- Исторические данные хранятся в S3
- При запуске читаем из S3 только индекс

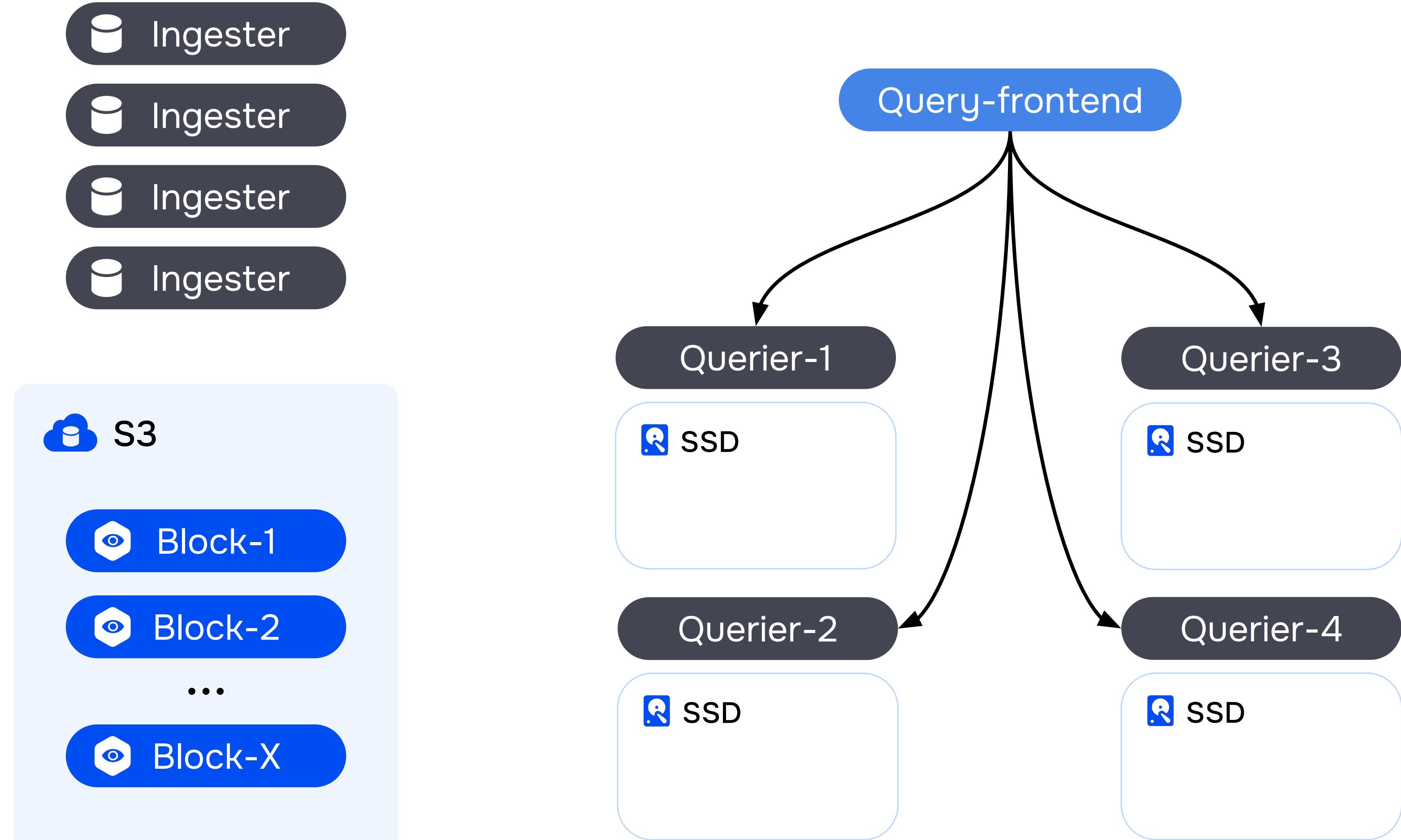


Централизованная система МОНИТОРИНГА

- Отказоустойчивость на запись: $WC = RF/2 + 1$
- Отказоустойчивость на чтение: $RC = RF - WC + 1$
- Шардирование на запись на основе хеш-ринга
- Исторические данные хранятся в S3
- При запуске читаем из S3 только индекс
- Читаем из чанков только нужные данные

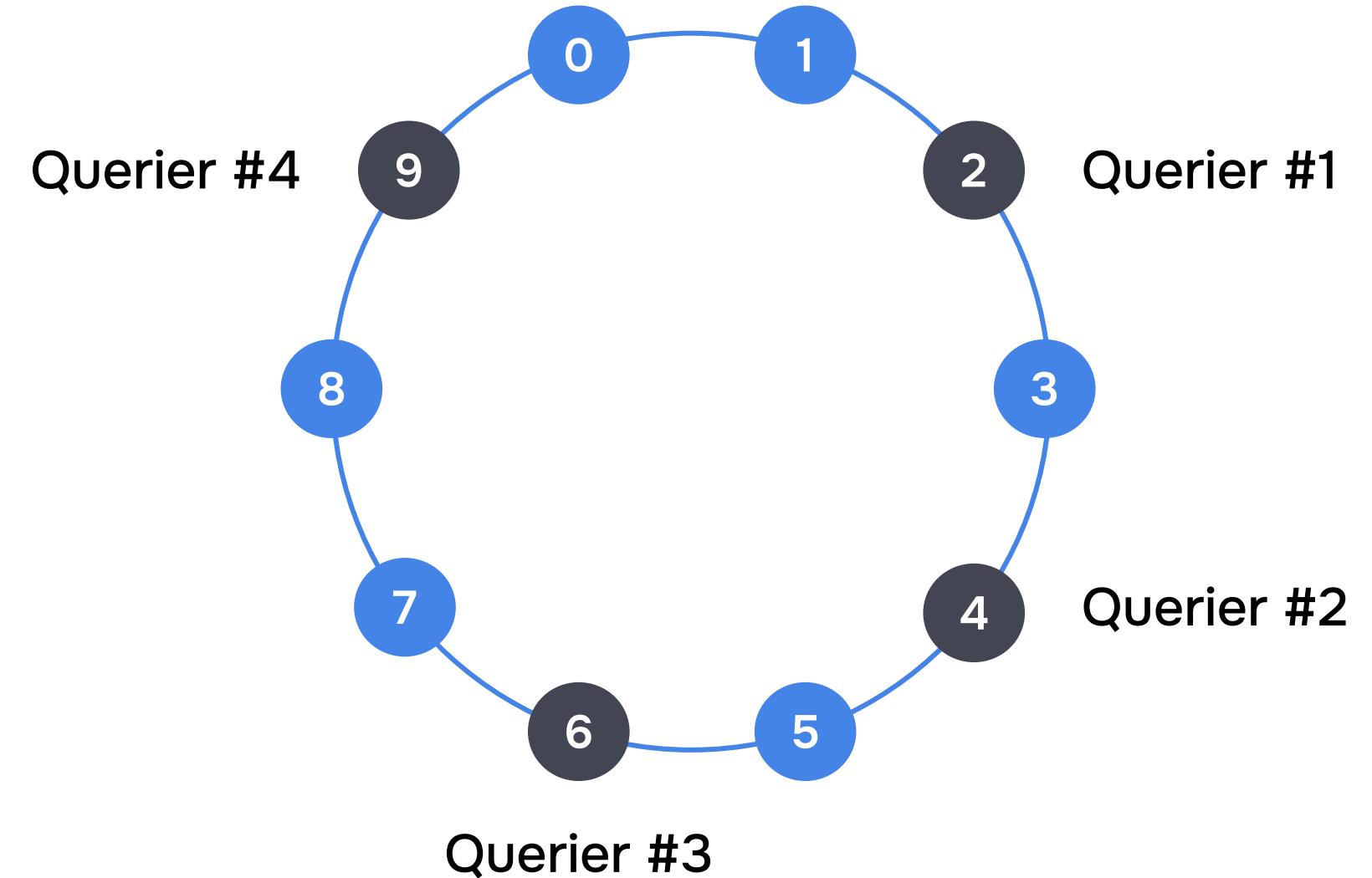
А как работает шардирование?





Хеш-ринг

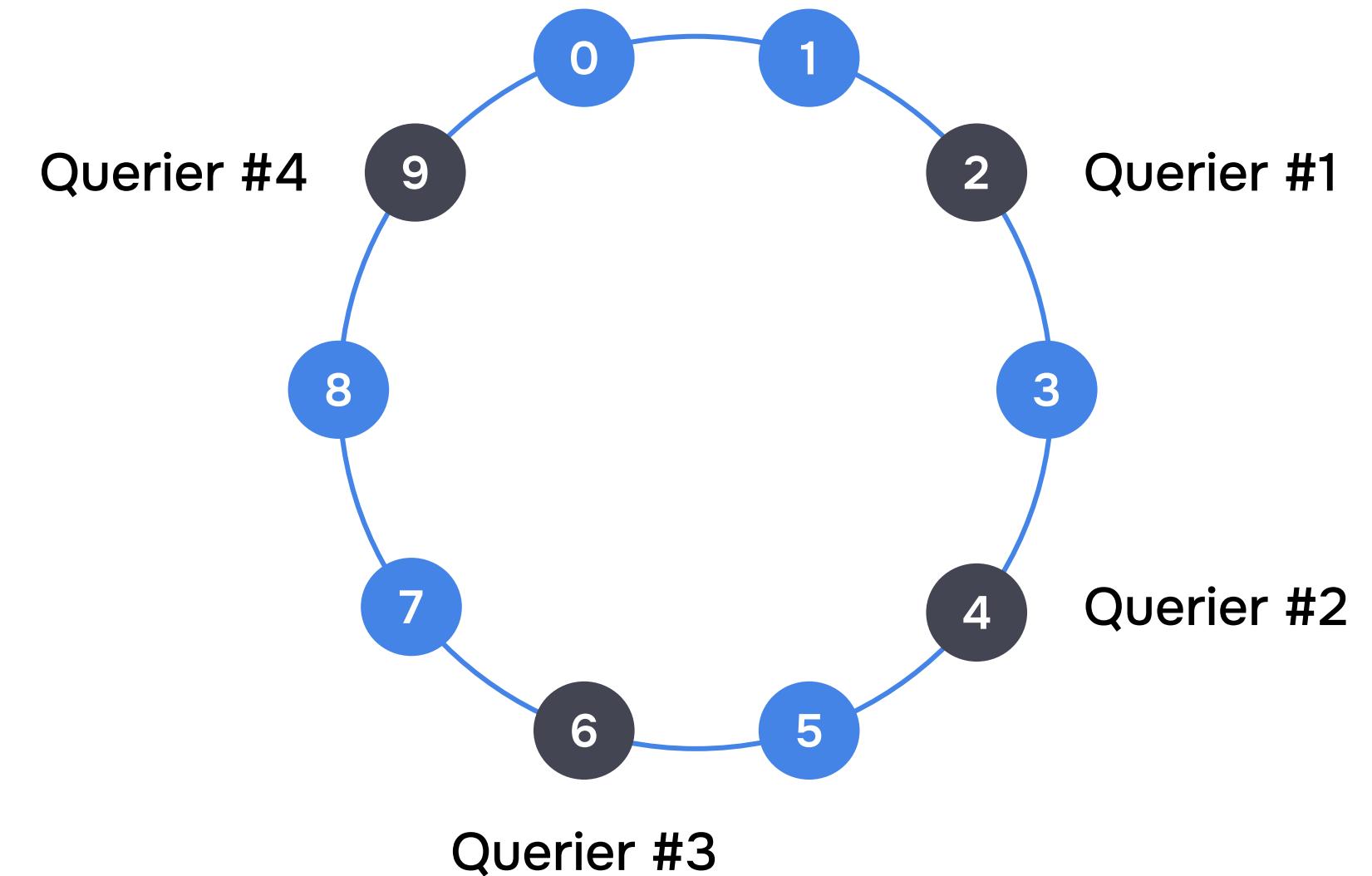
Хеш-ринг



Хеш-ринг

Block

UID: 01HZN9V7YZSC3F2Q0X8D0DYRZC

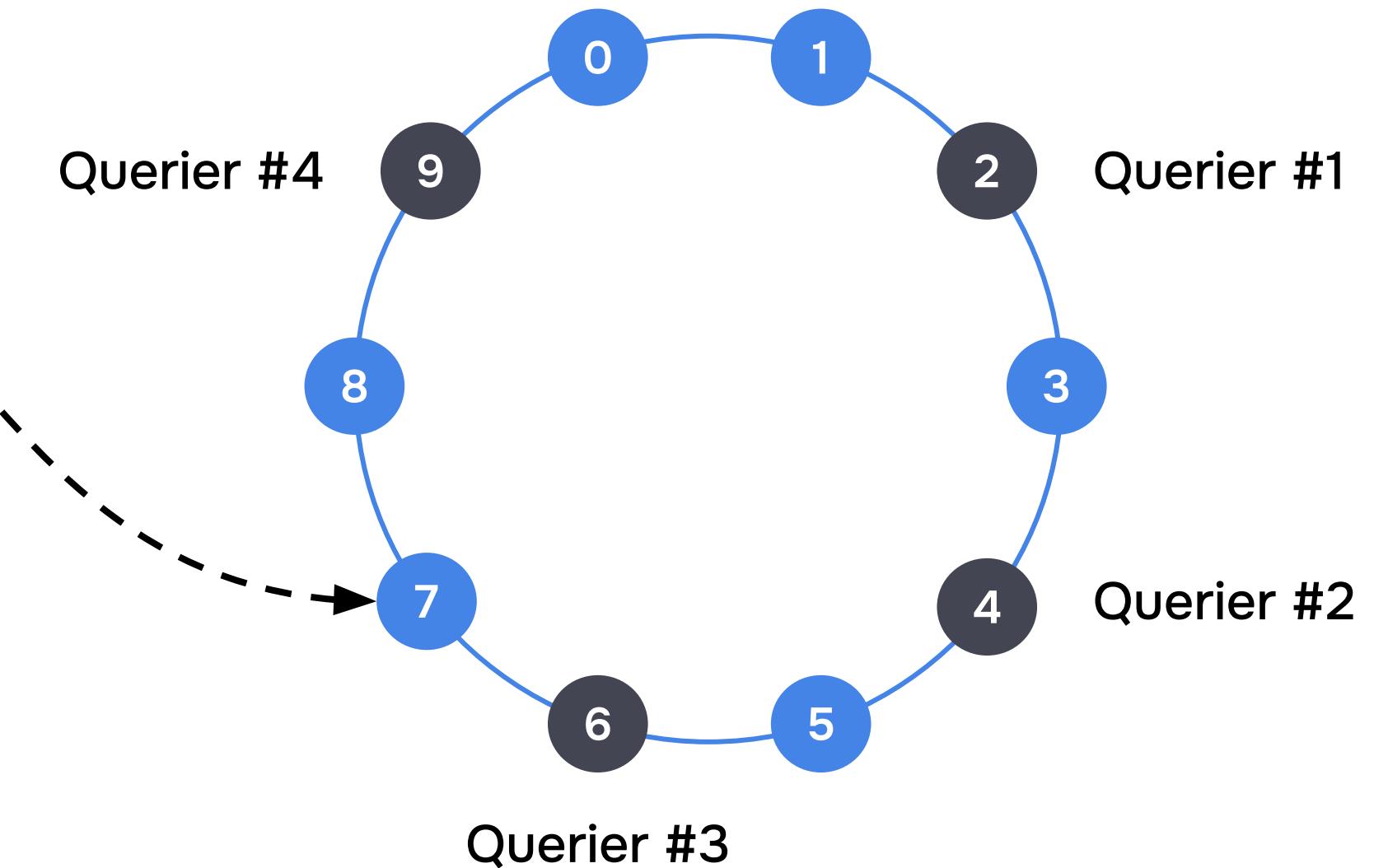


Хеш-риинг

Block

UID: 01HZN9V7YZSC3F2Q0X8D0DYRZC

-> **token = 7**

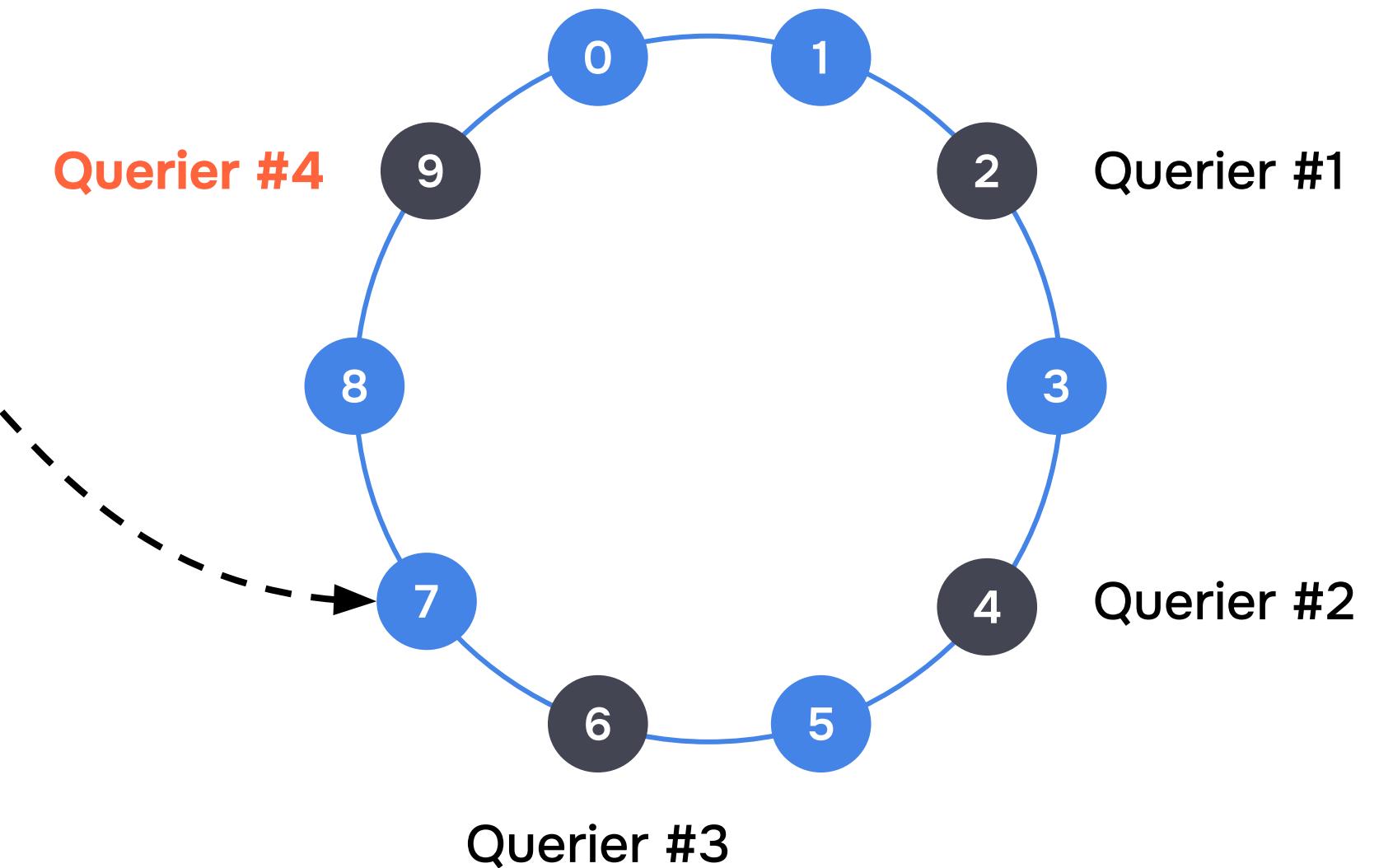


Хеш-ринг

Block

UID: 01HZN9V7YZSC3F2Q0X8D0DYRZC

-> **token = 7**

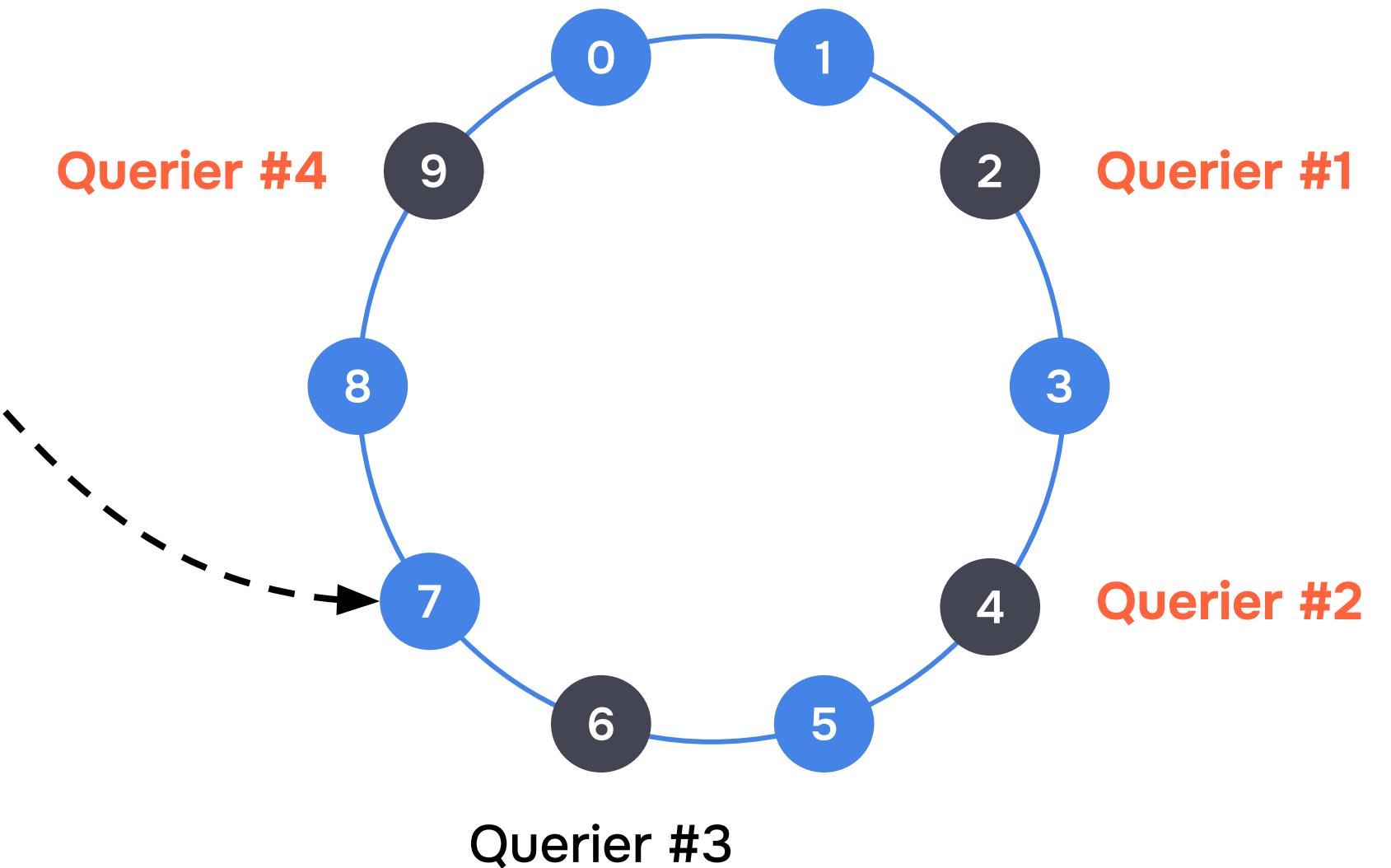


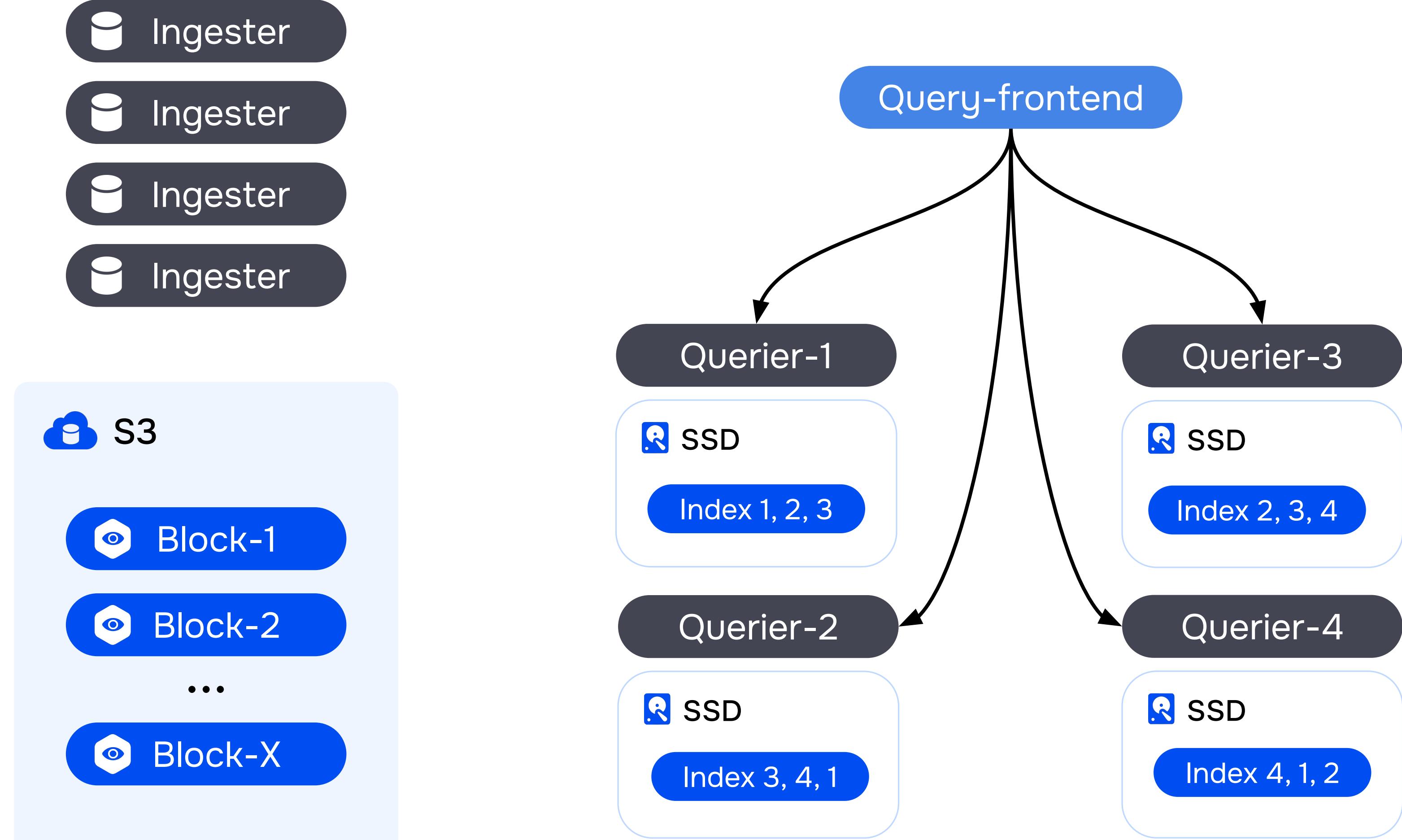
Хеш-ринг

Block

UID: 01HZN9V7YZSC3F2Q0X8D0DYZC

-> **token = 7**

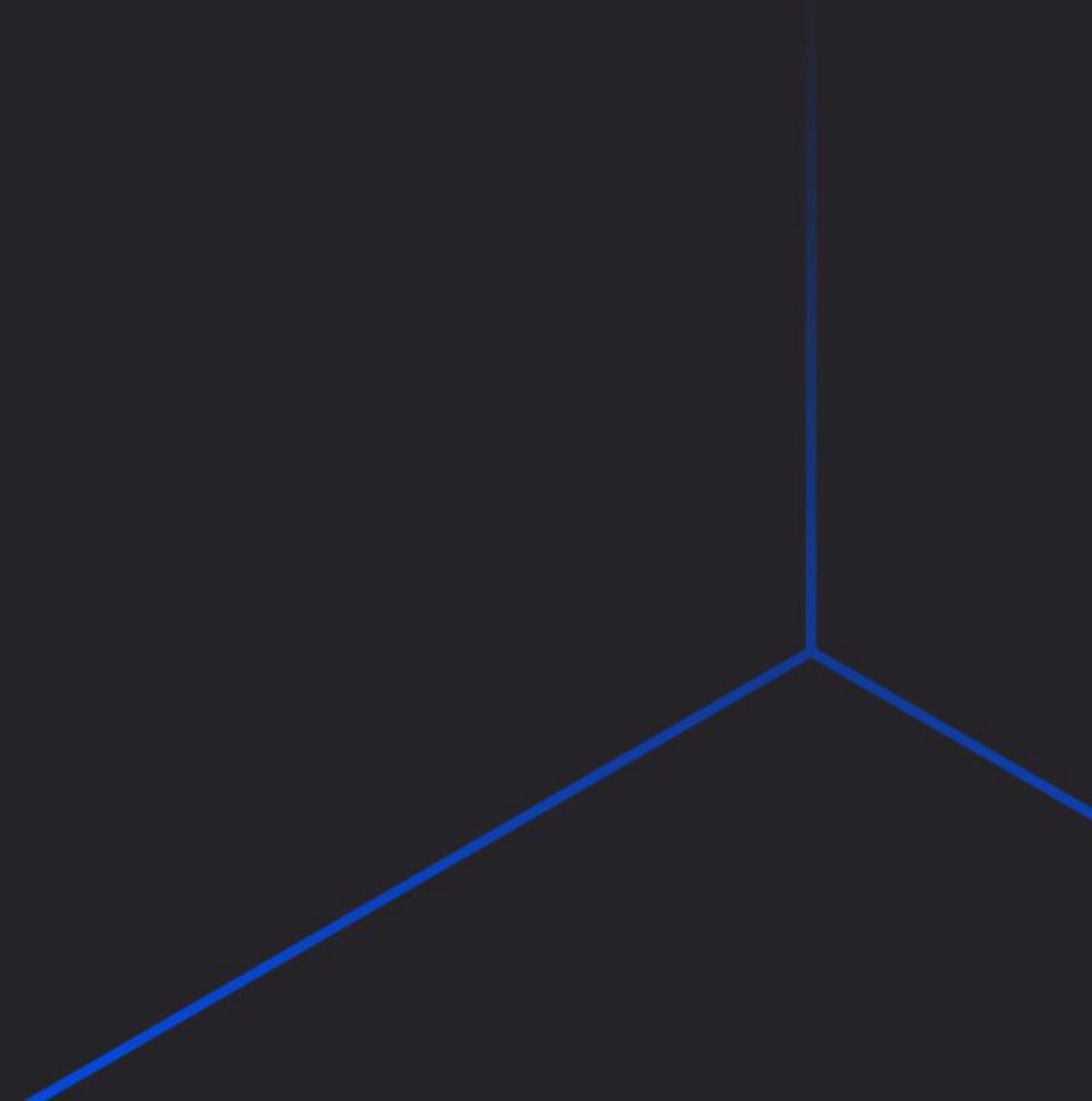


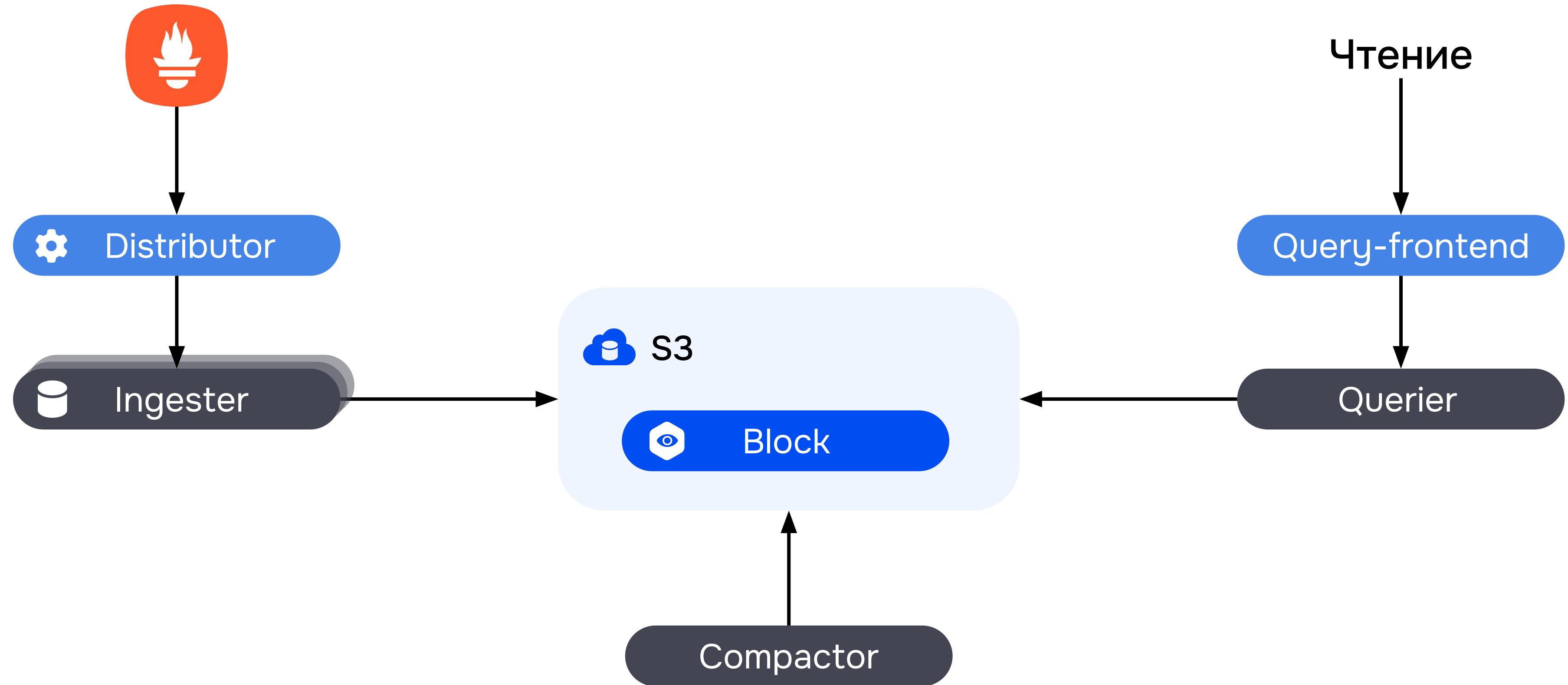


Централизованная система мониторинга

- Отказоустойчивость на запись: $WC = RF/2 + 1$
- Отказоустойчивость на чтение: $RC = RF - WC + 1$
- Шардирование на запись на основе хеш-ринга
- Исторические данные хранятся в S3
- При запуске читаем из S3 только индекс
- Читаем из чанков только нужные данные
- Шардирование на чтение на основе хеш-ринга

Итоги

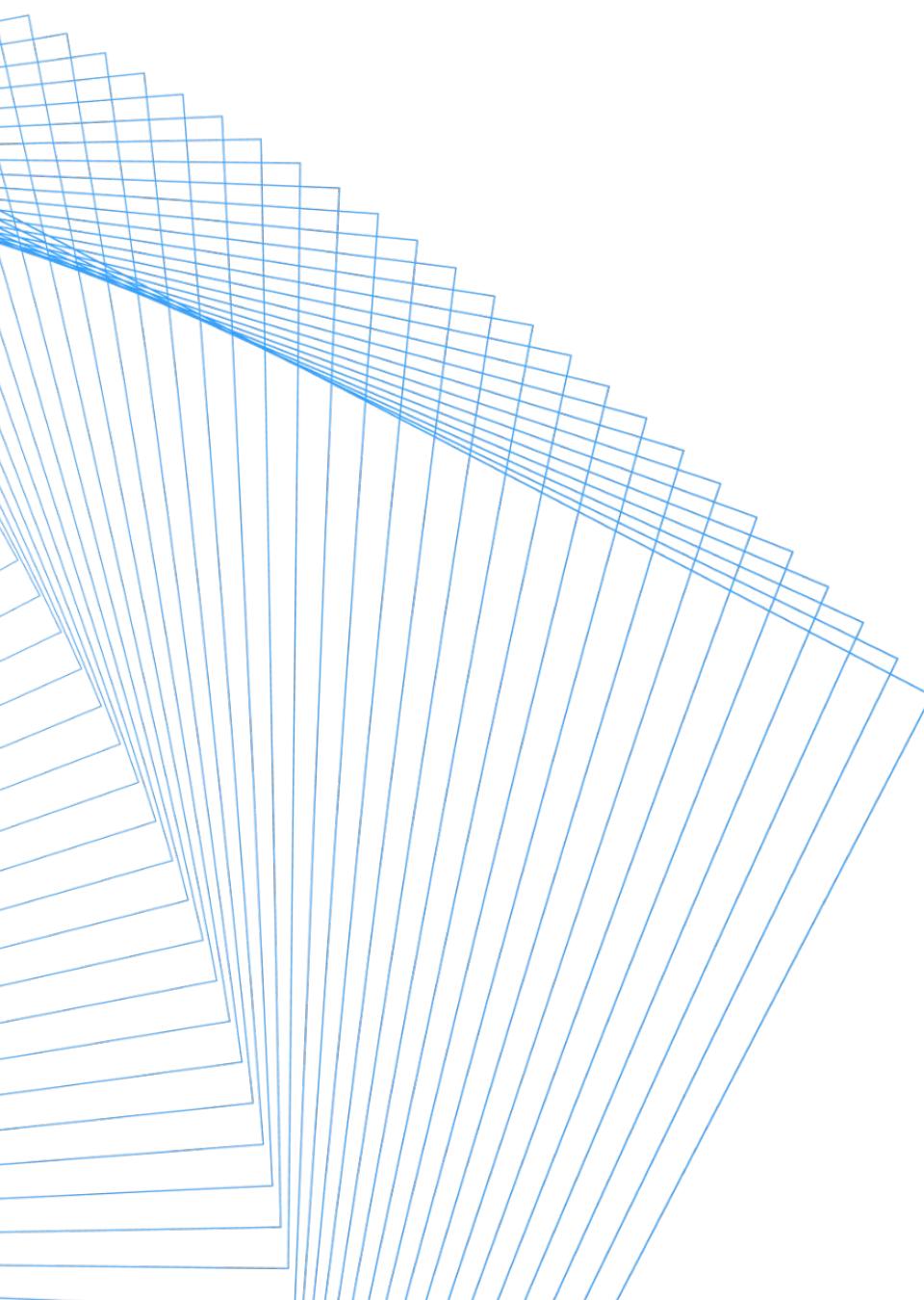




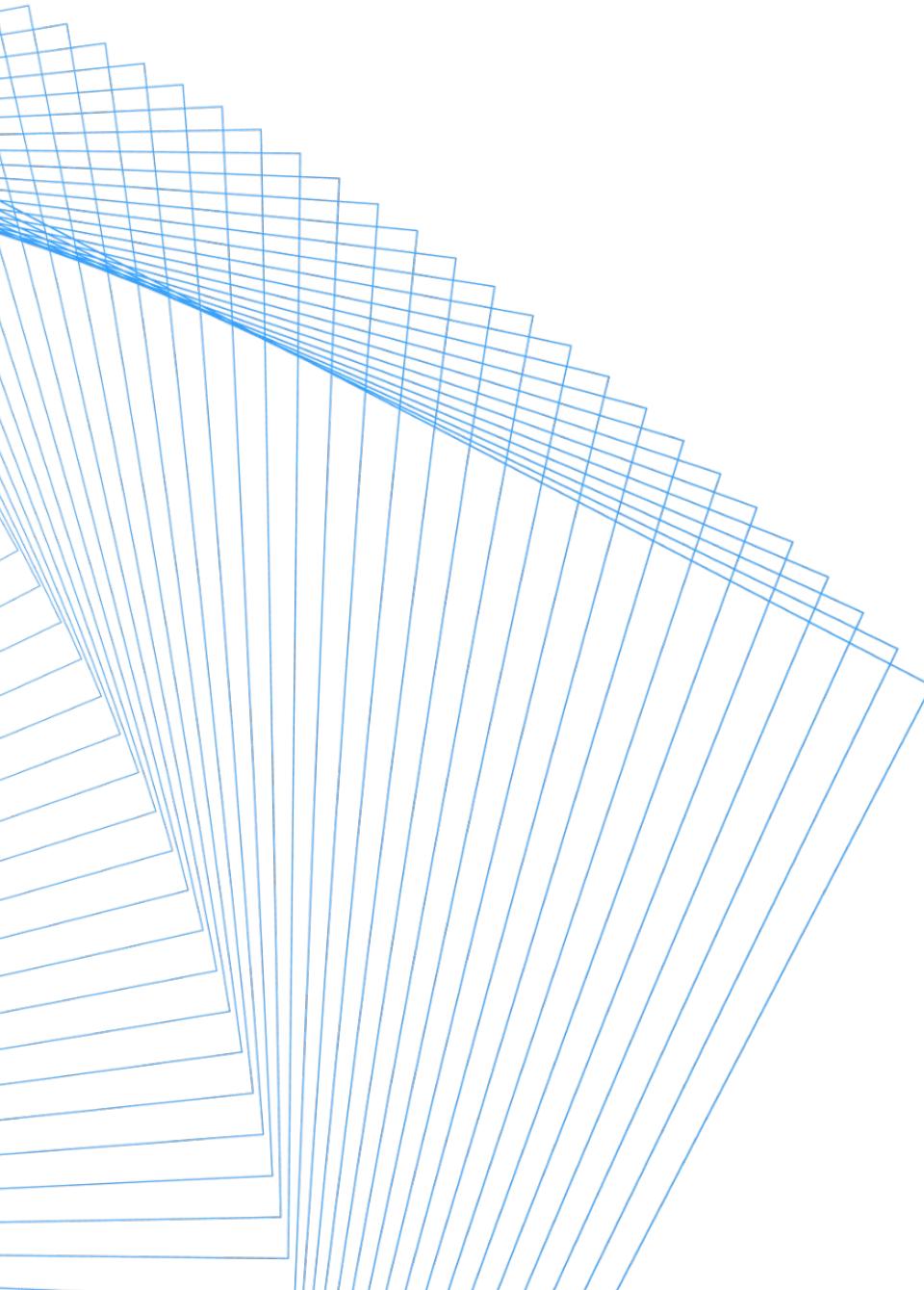
Итоги

1

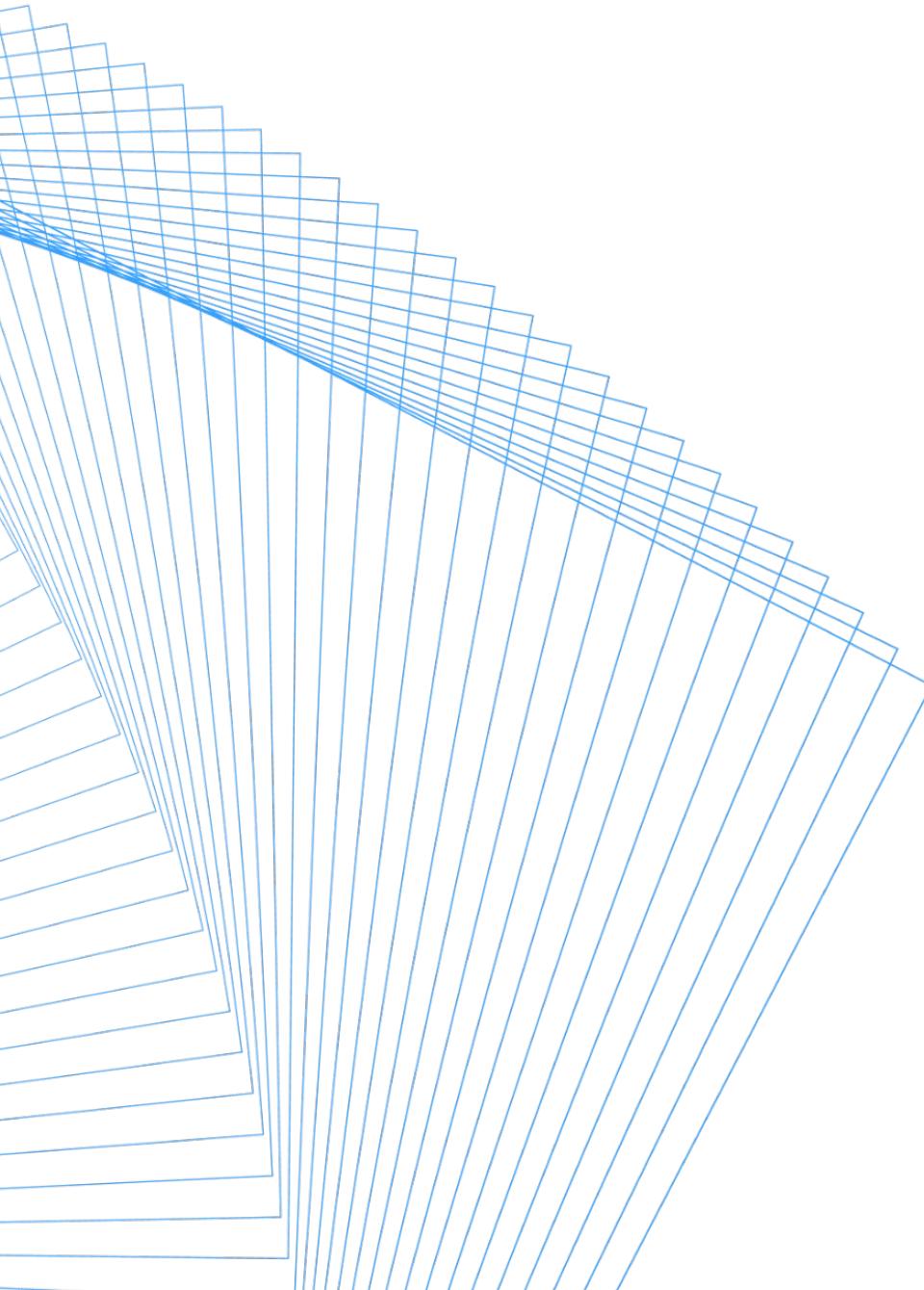
Отказоустойчивая система мониторинга



Итоги

- 
- 1 Отказоустойчивая система мониторинга
 - 2 Горизонтально масштабируемая система мониторинга

Итоги

- 
- 1 Отказоустойчивая система мониторинга
 - 2 Горизонтально масштабируемая система мониторинга
 - 3 С хранением данных в S3

Буду рад
ответить на ваши
вопросы!



[Telegram](#)



[YouTube](#)



[Habr](#)

 @Magvai69

 github.com/deckhouse/deckhouse

 github.com/deckhouse/prompp