



Off-heap вместо трех JEP'ов

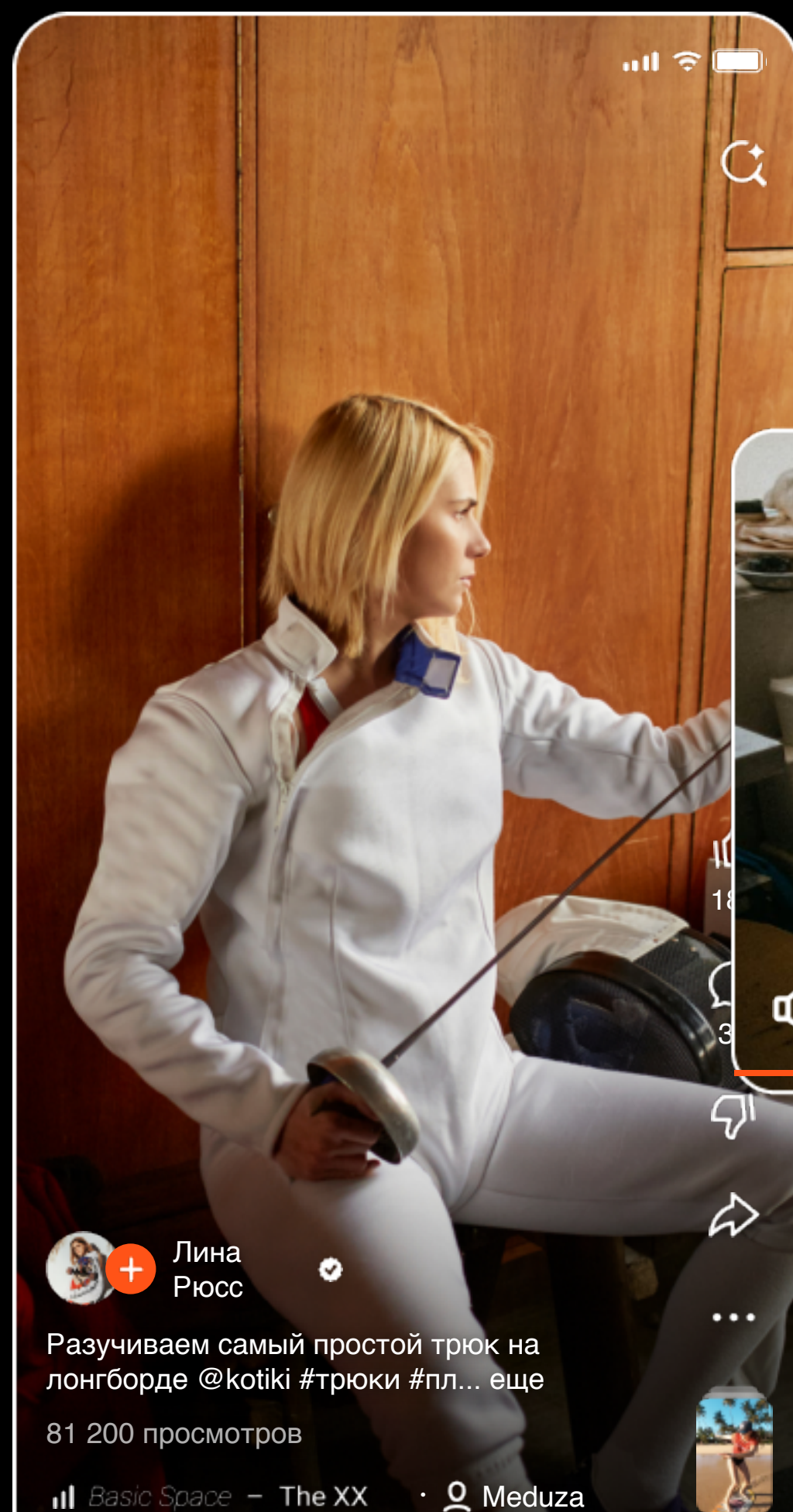
Дмитрий Погорелов, руководитель инфраструктуры рекомендательной системы Дзена
@rogorelych



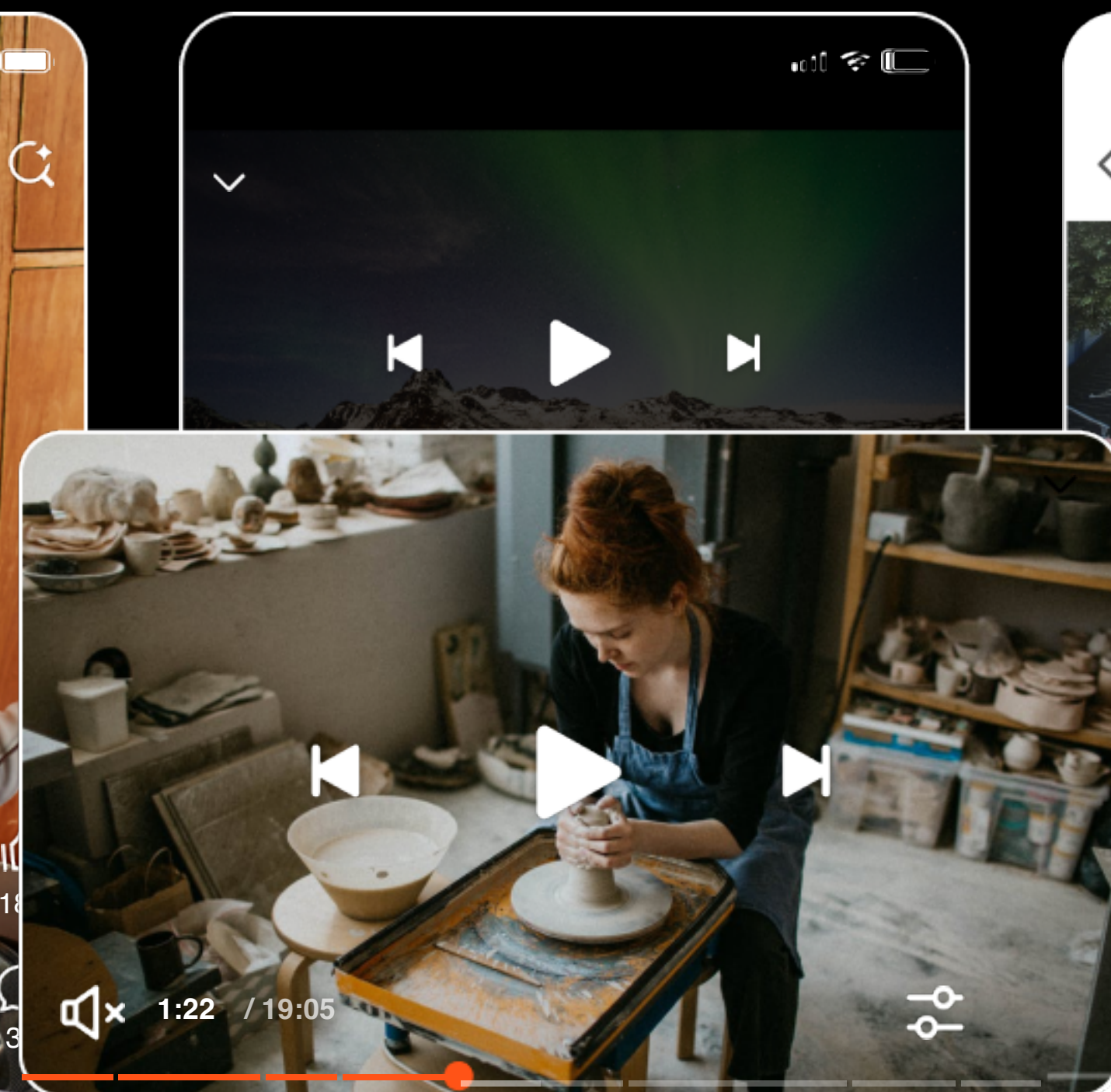
*Зачем C++ программист перешел
на другую сторону улицы?
Чтобы собрать мусор!*

Дзен — это большая контентная платформа

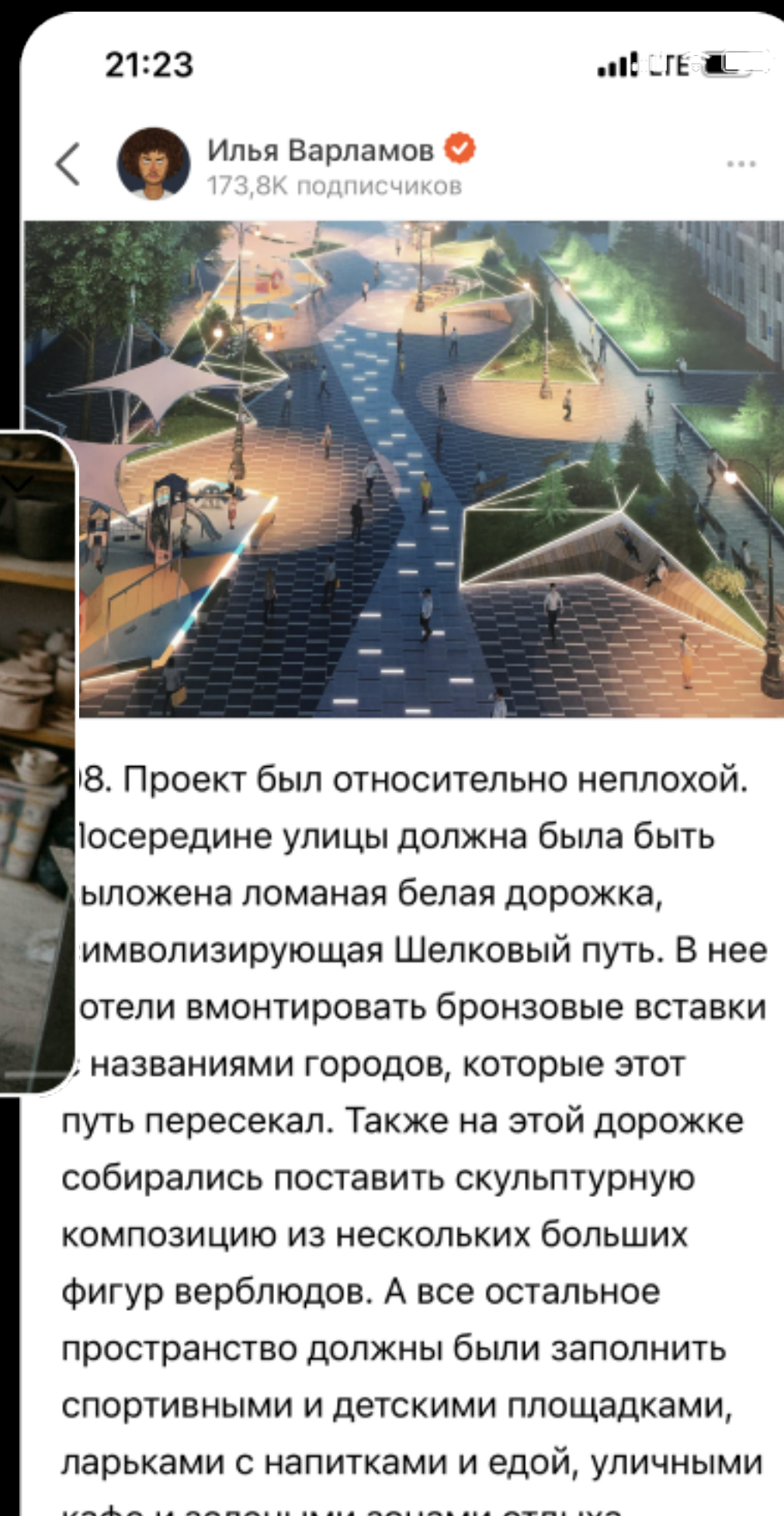
Вертикальные короткие ролики



Длинные видео



Тексты: лонгриды и посты



Дзен в цифрах

DAU 30M+
пользователей

MAU 70M+
пользователей

100K+
активных авторов

150K
RPS

100+
микросервисов

9999
отказоустойчивость

В сердце Дзена — рекомендательная система

5

Рекомендации работают в реальном времени

✦ 10 000 RPS
✦ 300 ms в 99_prc

10 млн документов

На запрос отскорим
50 тыс документов

... И МНОГО МАТЕМАТИКИ

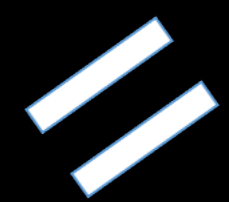
User \ Item	1	2	3	4	5	6	7	8	9	10
Петя	●		●	●	●			●		
Маша		●		●						
Вася	●		●		●		●			●
Катя	●			●	●			●		●



-0.82	-0.47	-1.05
-0.82	0.00	-0.00
-0.00	-1.28	0.77
0.82	-0.47	-1.05



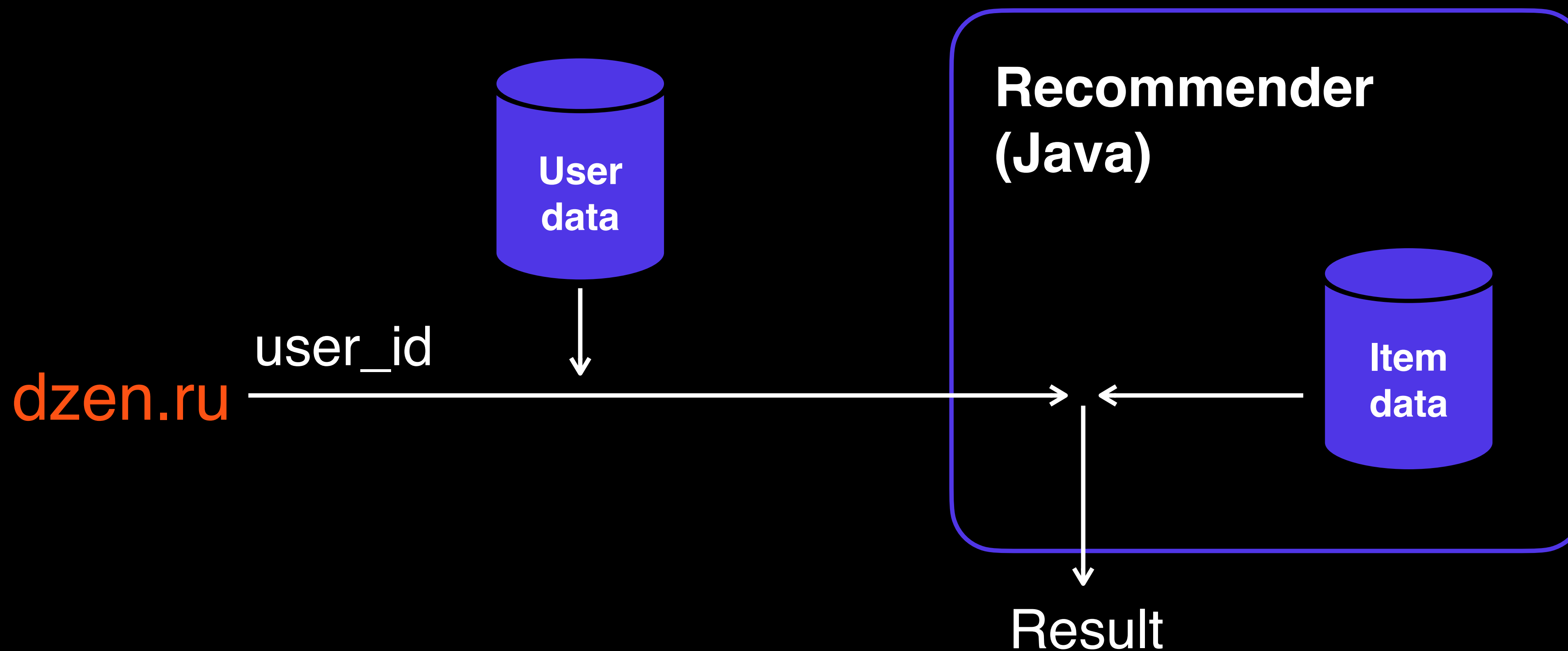
0.00	-0.41	-0.41	1.22	-0.00	0.00	0.00	0.00	0.00	0.00	-0.41
0.17	-0.00	0.39	0.00	1.07	0.00	-0.62	-0.45	0.00	0.00	-0.39
-1.03	0.00	-0.65	0.00	0.48	0.00	0.28	-0.75	0.00	0.00	0.65



User \ Item	1	2	3	4	5	6	7	8	9	10
Петя	1.00	0.33	0.83	-1.00	-1.00	0.00	0.00	1.00	0.00	-0.17
Маша	-0.00	0.33	0.33	-1.00	0.00	0.00	-0.00	-0.00	0.00	0.33
Вася	-1.00	-0.00	-1.00	0.00	-1.00	0.00	1.00	0.00	0.00	1.00
Катя	1.00	-0.33	0.17	1.00	-1.00	0.00	0.00	1.00	0.00	-0.83

Архитектура

Храним документы там же, где и ранжируем



Вызовы перед инфраструктурой

40+ гб

данных на один под

5 млрд

умножений векторов в секунду

Рецепт ранжирования

9

```
1 float[] user = {1f, 0f, 0.28f, .....}
2 float[] item = {0.3f, 1f, 0.121f, .....}
3
4 float feature_i = dotProduct(user, item)
5
6 float score(user,item) = matrixnet(float[] features)
7
8 float dotProduct(float[] a, float[] b) {
9     float sum = 0;
10    for (int i = 0; i < a.length; i++) {
11        sum += a[i] * b[i];
12    }
13    return sum;
14 }
```

Project Panama,
JEP-424: Foreign Function & Memory API

SIMD, AVX

JEP-448: Vector API

В off-heap проще

10

```
1 float* user = 0x434234
2 float* item = 0x434234
3 float feature = Eigen_JNI.dotProduct();
4 score = matrixnet(float*)
```

Много данных



Большой heap



Проблемы



- ✦ -XX:UseCompressedOops не работают
- ✦ Медленный GC
- ✦ Долгий старт



Off-heap коллекции

```
1 Map<Long, Ctrs> ctrs;  
2 record Ctrs (long clicks, long shows) {}
```

java.util.HashMap

14

```
1 Map<Long, Ctrs> ctrs = new HashMap ()
```

Node []

4 byte

92 byte

Из них 24 byte полезных

HashMap.Node

32 byte

- ✦ K key
- ✦ V value
- ✦ Node next
- ✦ int hash

Long
Ctrs

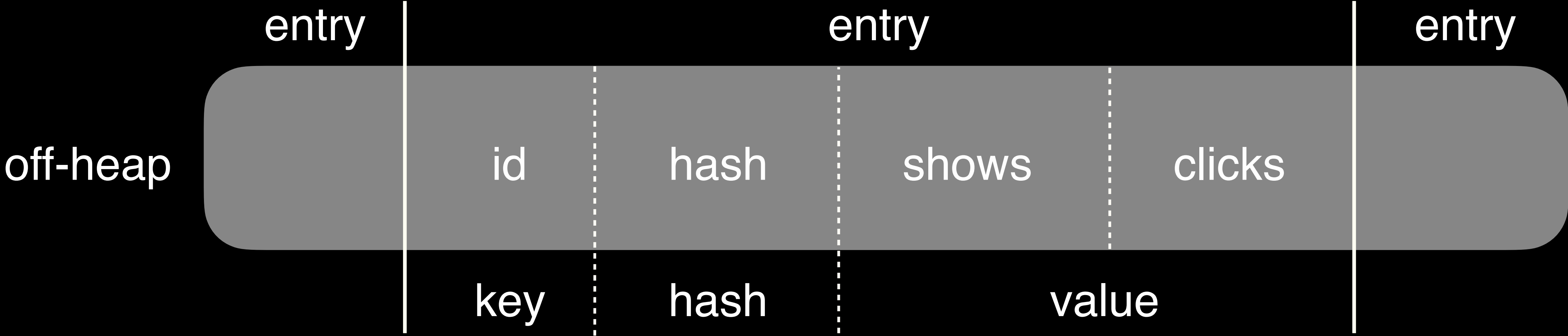
24 byte

32 byte

JEP 401: Value Classes and Objects

OffheapOpenHashMap

15



28 byte

Из них 24 byte полезных

```
1 V get (K key) {  
2     long addr = findPosition(key);  
3     V result = valueMemManager.get(addr);  
4 }
```

OffheapOpenHashMap.get(key) —> аллокация

OffheapOpenHashMap

17

- ✦ Маленький heap → проще GC
- ✦ Плотнo пакуем данные без overhead на заголовки
- ✦ Большие коллекции (больше 2 гб)
- ✦ Локальность данных



Мусорим при обращении

implementation `'ru.odnoklassniki:one-nio:1.7.1'`



github.com/odnoklassniki/one-nio

- ✦ `float[]` храним с выравниванием
- ✦ Ленивая десериализация `Map<Long, Map<Long, ...>>`
- ✦ Коллекции примитивов

```
1 LongLongHashMap  
2  
3 public long get(long key) {}  
4  
5 public void put(long key, long value) {}
```

Практическое применение

id документа



id автора

id документа



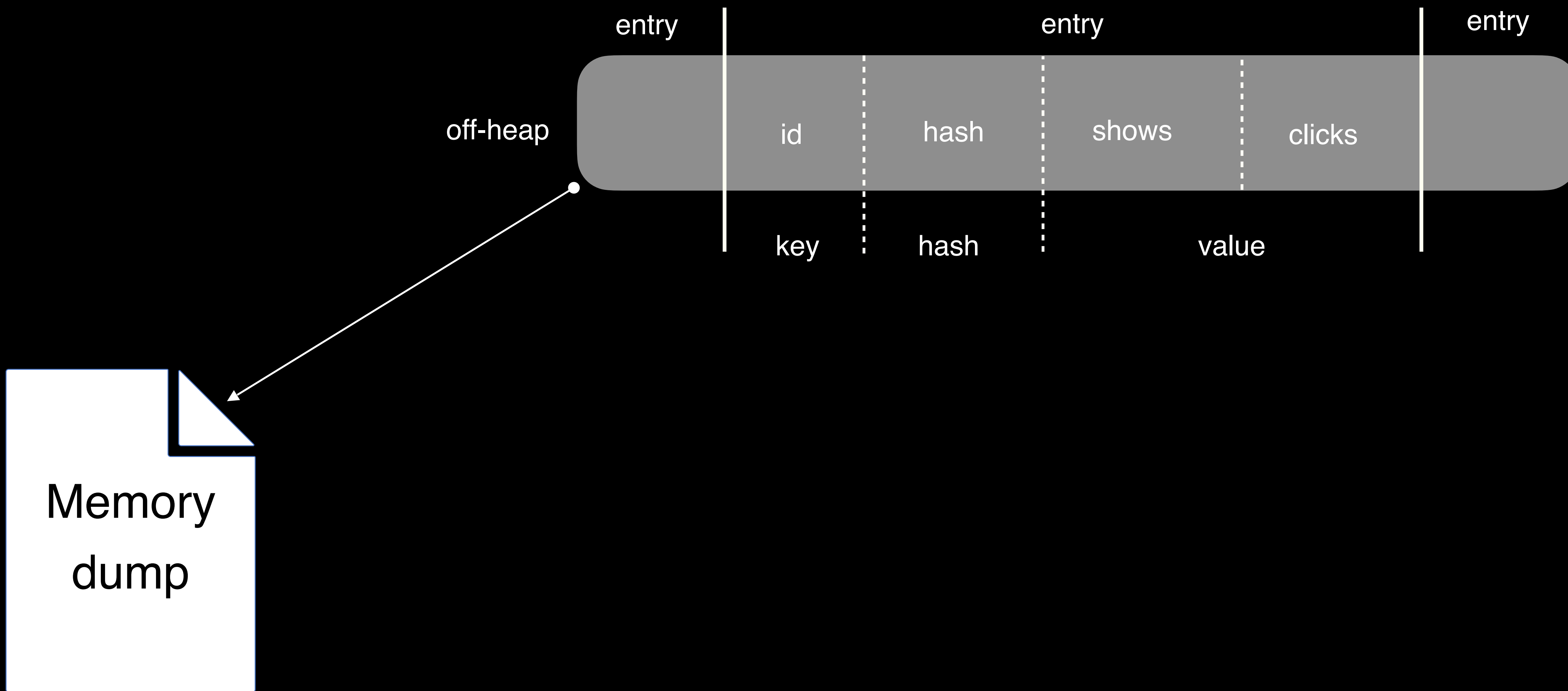
статус модерации

id документа

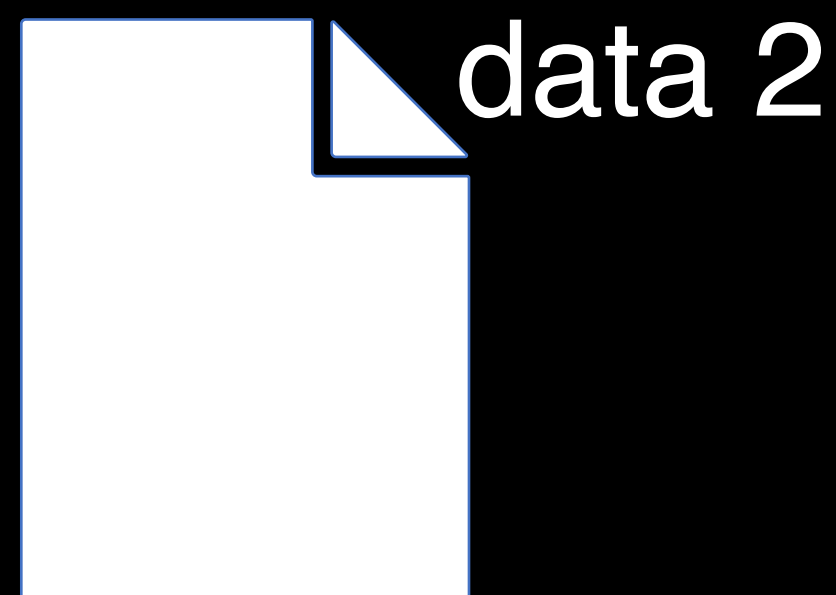
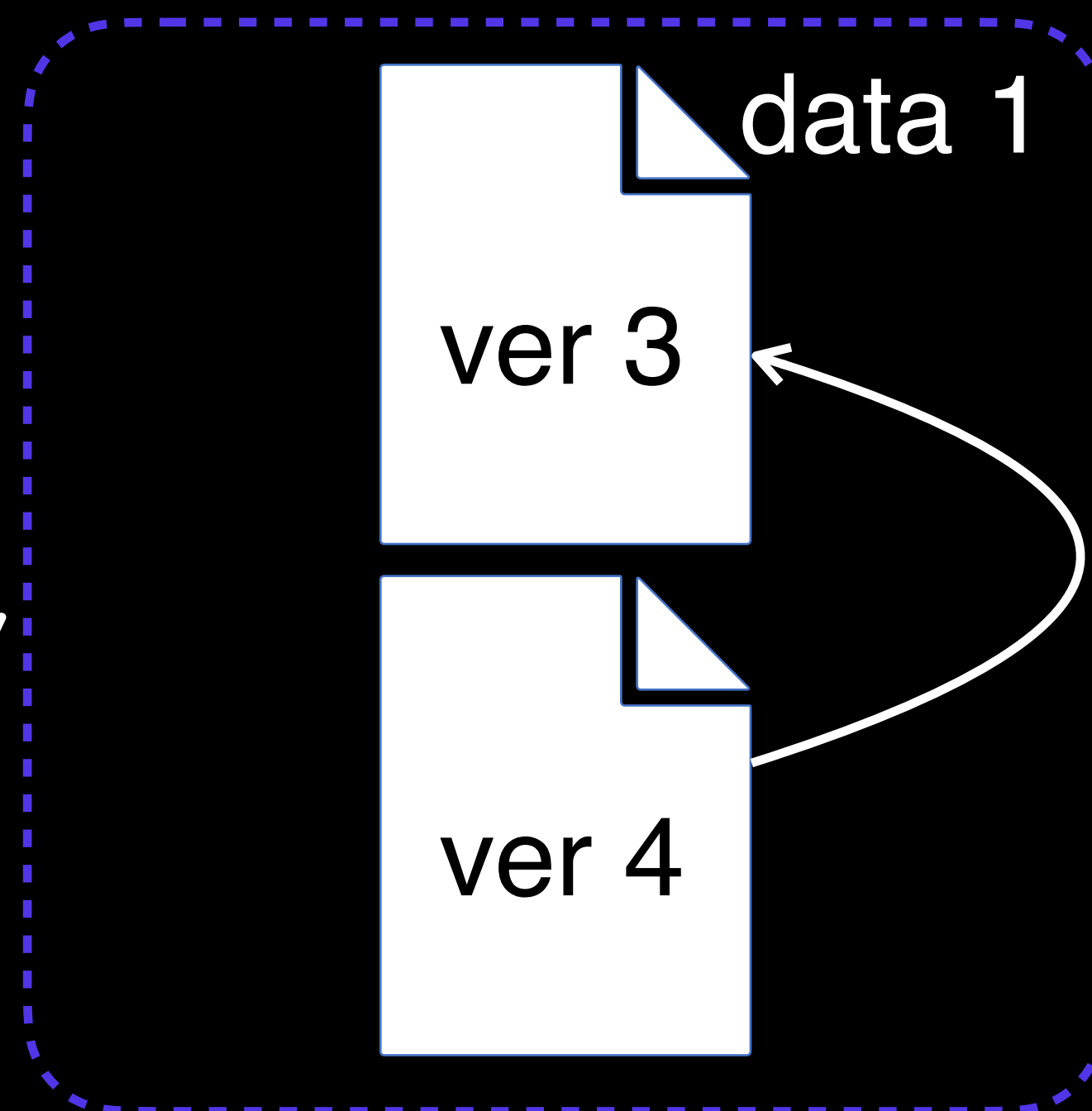


СПИСОК ПОХОЖИХ

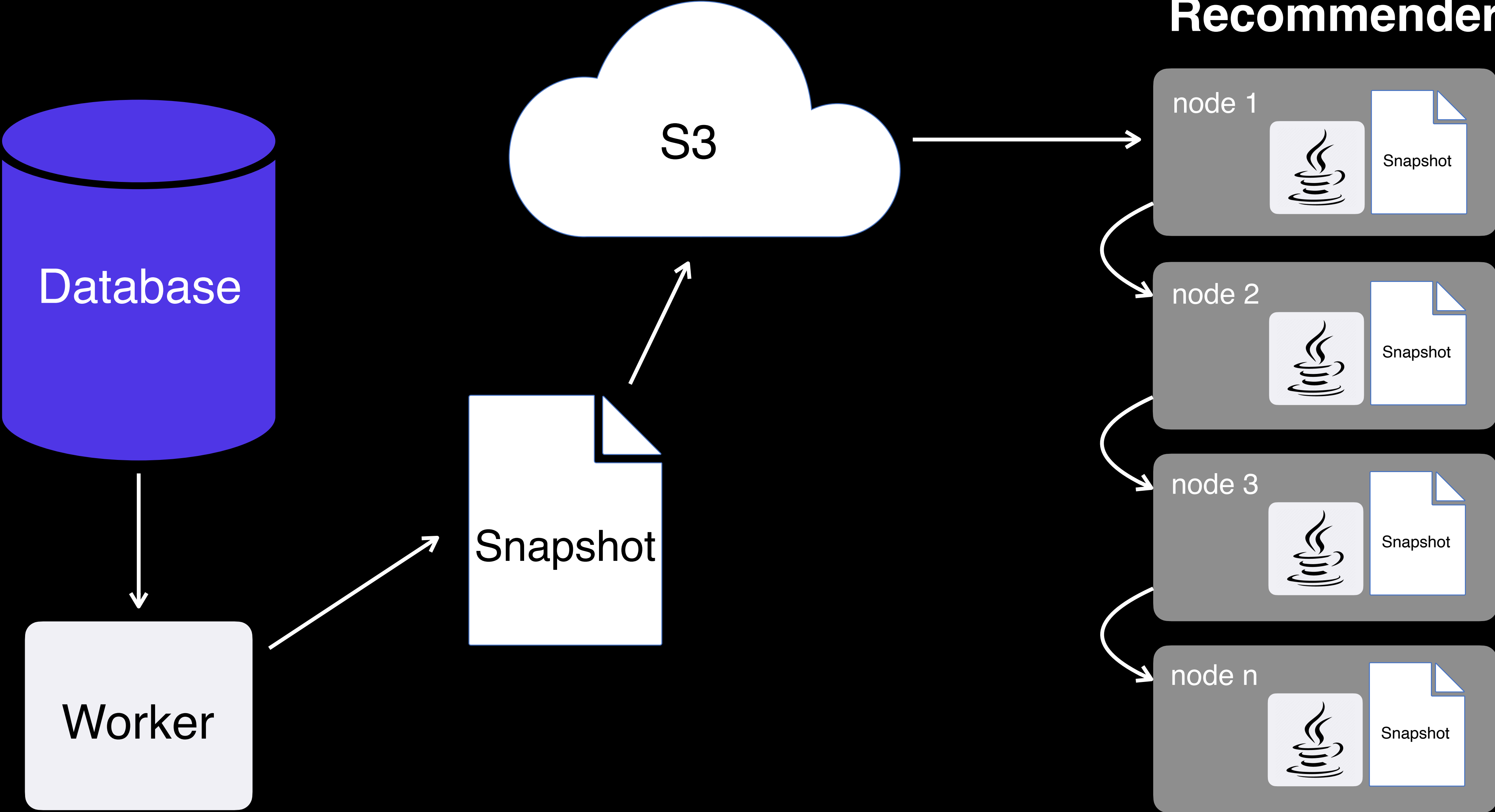
Dump памяти = бесплатная сериализация



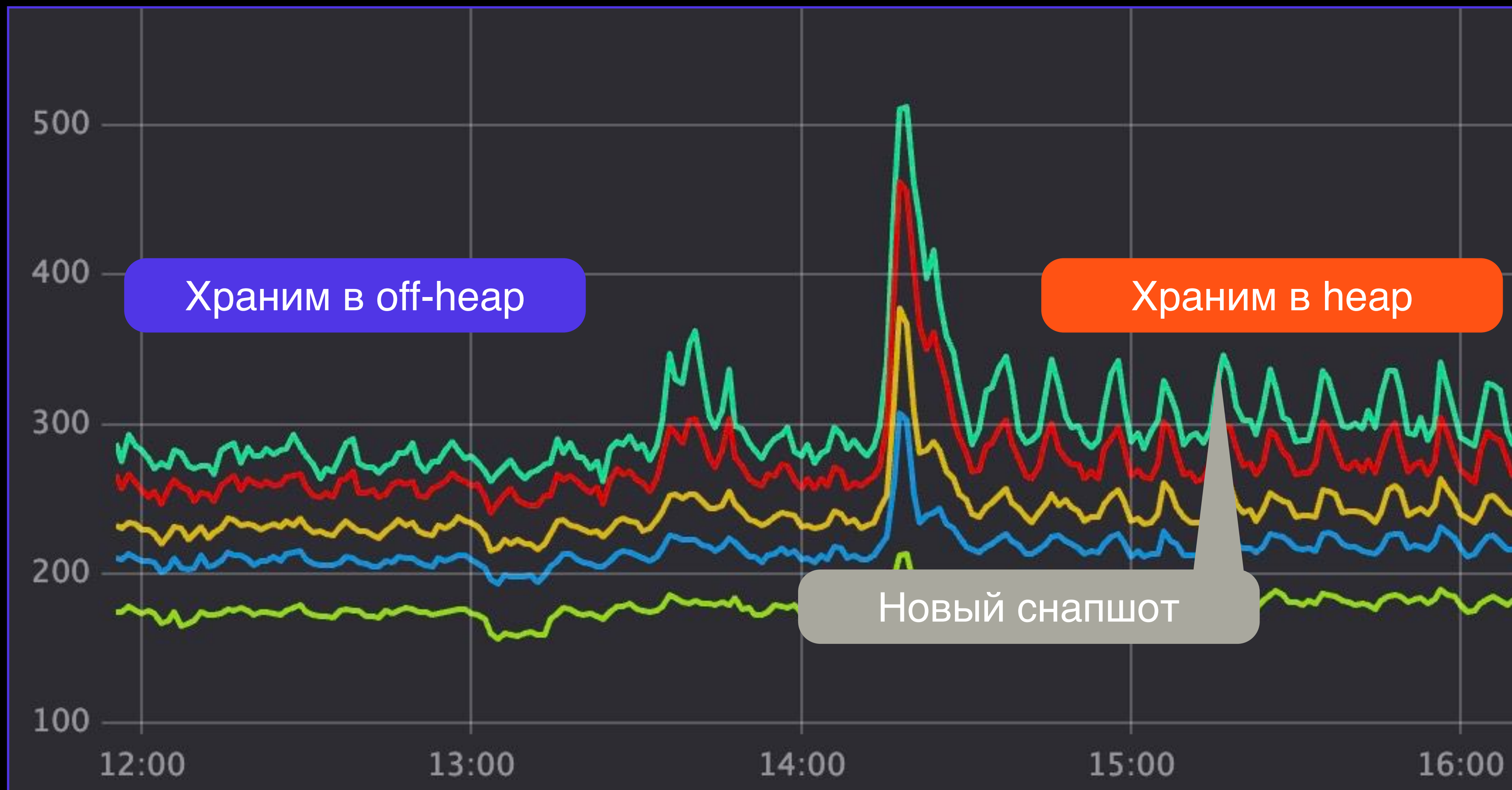
mmaped



Снапшоты на кластере



Переложили 3 гб в hear



- ✦ Бесплатно загружаем на нодах
- ✦ Моментальный рестарт с горячим кешом
- ✦ 40+ гб на 2 тыс нод за 15 минут...
или быстрее, упираемся только в сеть



Как работать с off-heap

Добываем off-heap

28

```
long addr = Unsafe.allocateMemory(size)
```

```
ByteBuffer buffer = ByteBuffer.allocateDirect(...)
```

```
MappedByteBuffer javamaped = file.getChannel().map(...)
```

```
MMapBuffer buf = new MMapBuffer(file, ...)
```

Добываем Unsafe

29

```
1 public static final sun.misc.Unsafe UNSAFE = initializeUnsafe();
2
3 private static sun.misc.Unsafe initializeUnsafe() {
4     try {
5         Field field = sun.misc.Unsafe.class.getDeclaredField("theUnsafe");
6         field.setAccessible(true);
7         return (sun.misc.Unsafe) field.get(null);
8     } catch (NoSuchFieldException | IllegalAccessException e) {
9         throw new Error(e);
10    }
11 }
```

```
1 long addr = UNSAFE.allocateMemory(size);  
2 UNSAFE.setMemory(addr, size, (byte) 0);  
3 UNSAFE.copyMemory(...);  
4 UNSAFE.getLong(addr);  
5 UNSAFE.putByte(addr);  
6 UNSAFE.freeMemory(addr);
```

DirectByteBuffer

31

```
1 int size = ...;
2 ByteBuffer buffer = ByteBuffer.allocateDirect(size);
3 byte a = buffer.get(index);
4 long b = buffer.getLong(index);
5 buffer.putLong(index, value);
```

MappedByteBuffer

32

```
1 RandomAccessFile file = new RandomAccessFile("file", "rw");
2
3 MappedByteBuffer javamaped =
4     file.getChannel().map(FileChannel.MapMode.READ_WRITE, 0, f.length());
5
6 javamaped.force();
```


DirectByteBuffer MappedByteBuffer

33



Подчищаются GC



Не больше 2 гб

XX:MaxDirectMemory

34

```
1 -XX:MaxDirectMemorySize
2 -XX:MaxRAM
3
4 Exception in thread "main" java.lang.OutOfMemoryError: Cannot reserve 200000000 bytes of direct buffer memory
5   at java.base/java.nio.Bits.reserveMemory(Bits.java:178)
6   at java.base/java.nio.DirectByteBuffer.<init>(DirectByteBuffer.java:121)
7   at java.base/java.nio.ByteBuffer.allocateDirect(ByteBuffer.java:332)
8   at Main.main(Main.java:42)
```

$$\text{MaxRam} = \text{Xmx} + \text{MaxMetaspaceSize} + \text{ReservedCodeCacheSize} + \text{MaxDirectMemorySize}$$

-XX:MaxRam -XX:MaxDirectMemory только для DirectByteBuffer

XX:MaxDirectMemory

36

-XX:MaxDirectMemorySize=1m
-XX:MaxRAM=100m

```
1 long size = 1000_000_000L;  
2 long addr = unsafe.allocateMemory(size);  
3 RandomAccessFile f = new RandomAccessFile("my_200mb_file", "rw");  
4 MappedByteBuffer javamaped =  
5     f.getChannel().map(FileChannel.MapMode.READ_WRITE, 0, f.length());
```

```
1 implementation 'com.indeed:util-mmap:1.0.52-3042601'  
2  
3 MMapBuffer buf = new MMapBuffer(file, FileChannel.MapMode.READ_ONLY,  
4 buf.mlock(0, bufferLength);  
5 long addr = buf.memory().getAddress()
```

com.indeed:util-mmap

38



Файлы любого размера



mlock

```
1 java.lang.OutOfMemoryError: null
2   at com.indeed.util.mmap.NativeMemoryUtils.mlock(NativeMemoryUtils.java:58)
3   at com.indeed.util.mmap.MMapBuffer.mlock(MMapBuffer.java:186) ~[util-mmap-1
```

Ставим капю на Java

```
setcap cap_ipc_lock=epi $JAVA_BINARY
```

Получаем ошибку

```
java: error while loading shared libraries: libjli.so: cannot  
open shared object file: No such file or directory
```

Решаем

```
echo "$JAVA_HOME/lib" > /etc/ld.so.conf.d/java.conf  
ldconfig
```



unix.stackexchange.com/a/616511

Что делать в k8s?

k8s использует rlimit с ноды, на которой запущен контейнер, и это не конфигурируется для контейнеров



unix.stackexchange.com/a/616511



Off-heap

GC

Он вам не Exception

43

```
#  
# A fatal error has been detected by the Java Runtime Environment:  
#  
# SIGSEGV (0xb) at pc=0x0000000102ac9260, pid=20968, tid=8707  
#  
# JRE version: OpenJDK Runtime Environment Corretto-17.0.6.10.1 (17.0.6+10) (build 17.0.6+10-LTS)  
# Java VM: OpenJDK 64-Bit Server VM Corretto-17.0.6.10.1 (17.0.6+10-LTS, mixed mode, sharing, tiered,  
# Problematic frame:  
# V [libjvm.dylib+0x9d5260] Unsafe_PutLong(JNIEnv_*, _jobject*, _jobject*, long, long)+0x118  
#  
# No core dump will be written. Core dumps have been disabled. To enable core dumping, try "ulimit -c  
#  
# If you would like to submit a bug report, please visit:  
# https://github.com/corretto/corretto-17/issues/
```

`-XX:ErrorFile=/var/log/my_app/java_error%p.log`

Try-with-resources

44

```
1 public class MyOffheapObject
2     implements Closeable {
3
4     private final long addr;
5
6     public MyOffheapObject() {
7         this.addr = UNSAFE.allocateMemory(size);
8     }
9
10    @Override
11    public void close() {
12        UnsafeUtils.UNSAFE.freeMemory(addr);
13    }
14 }
15 public static void main(String[] args) {
16     try (var obj = new MyOffheapObject()) {
17         //do some stuff
18     }
19 }
```

Try-with-resources — не лучший вариант, если...

45

- ✦ Объекты в глобальном контексте
- ✦ Хотим спрятать за интерфейс
- ✦ Лень про это помнить

~~Finalize~~ → ~~fanthom ref~~ → JDK9

```
1 public class MyOffheapObject implements AutoCloseable {
2     private final static Cleaner cleaner = Cleaner.create();
3
4     private final long addr;
5     private final Cleaner.Cleanable cleanable;
6
7     public MyOffheapObject(long size) {
8         this.addr = UNSAFE.allocateMemory(size);
9         this.cleanable = cleaner.register(
10             this, new Deallocator(addr));
11     }
12
13     @Override
14     public void close() {
15         cleanable.clean();
16     }
17 }
18 private static class Deallocator implements Runnable {
19     private final long addr;
20
21     private Deallocator(long addr) {
22         this.addr = addr;
23     }
24
25     @Override
26     public void run() {
27         UnsafeUtils.UNSAFE.freeMemory(addr);
28     }
29 }
```



... но GC ничего не знает про off-heap

Try-with-resources

49

```
1 try (resource1 = ... ; resource2 = ...) {  
2     ResourceContext context = new ResourceContext(resource1, resource2, ...)   
3     ...  
4  
5     doSomeStuff(context);  
6  
7     ...  
8     doAnotherStuff(context);  
9     ...  
10 }
```



“Отличный” способ получить связность в коде

Try-with-resources, но волшебный

50

```
1 try () {  
2     var resource1 = ...;  
3     ....  
4     doSomeStuff();  
5     ....  
6  
7 }  
8  
9 void doSomeStuff() {  
10  
11     var resource2 = ...;
```

Прикапываем созданное в `thread-local` или `spring request context`

Велосипед уместен, если он едет в нужную сторону :)

- ✦ Быстро работает математика
- ✦ Быстро доезжают данные
- ✦ Помещается большая база



Спасибо! Задавайте вопросы

Дмитрий Погорелов, руководитель инфраструктуры рекомендательной системы Дзена
@rogorelych