

Пишем эффективную многопоточную хеш-таблицу

Никита Коваль

ndkoval *at* ya *dot* ru
twitter.com/nkoval_

Осторожно! Этот доклад про многопоточность и алгоритмы!*

* В других залах вы можете послушать про связку Spring и Reactor (зал 2) или про эксплуатацию Hadoop (зал 4)

О себе

- Инженер-исследователь в лаборатории dxLab, Devexperts
- Преподаю курс по многопоточному программированию в ИТМО
- Заканчиваю магистратуру ИТМО



Зачем мне это?

- Кто использует HashMap?

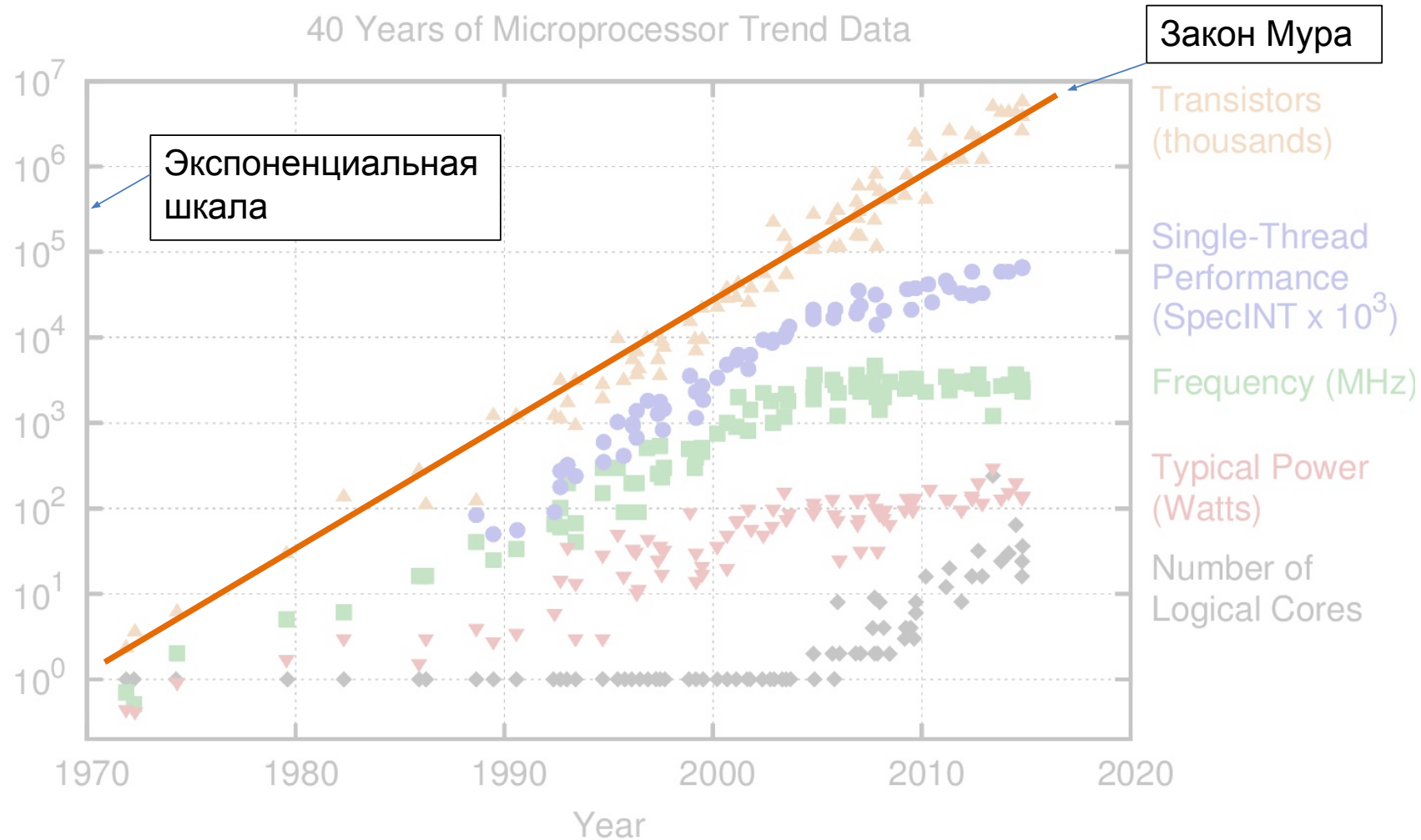
Зачем мне это?

- Кто использует HashMap?
 - Наверное, это самая популярная и полезная на практике структура данных

Зачем мне это?

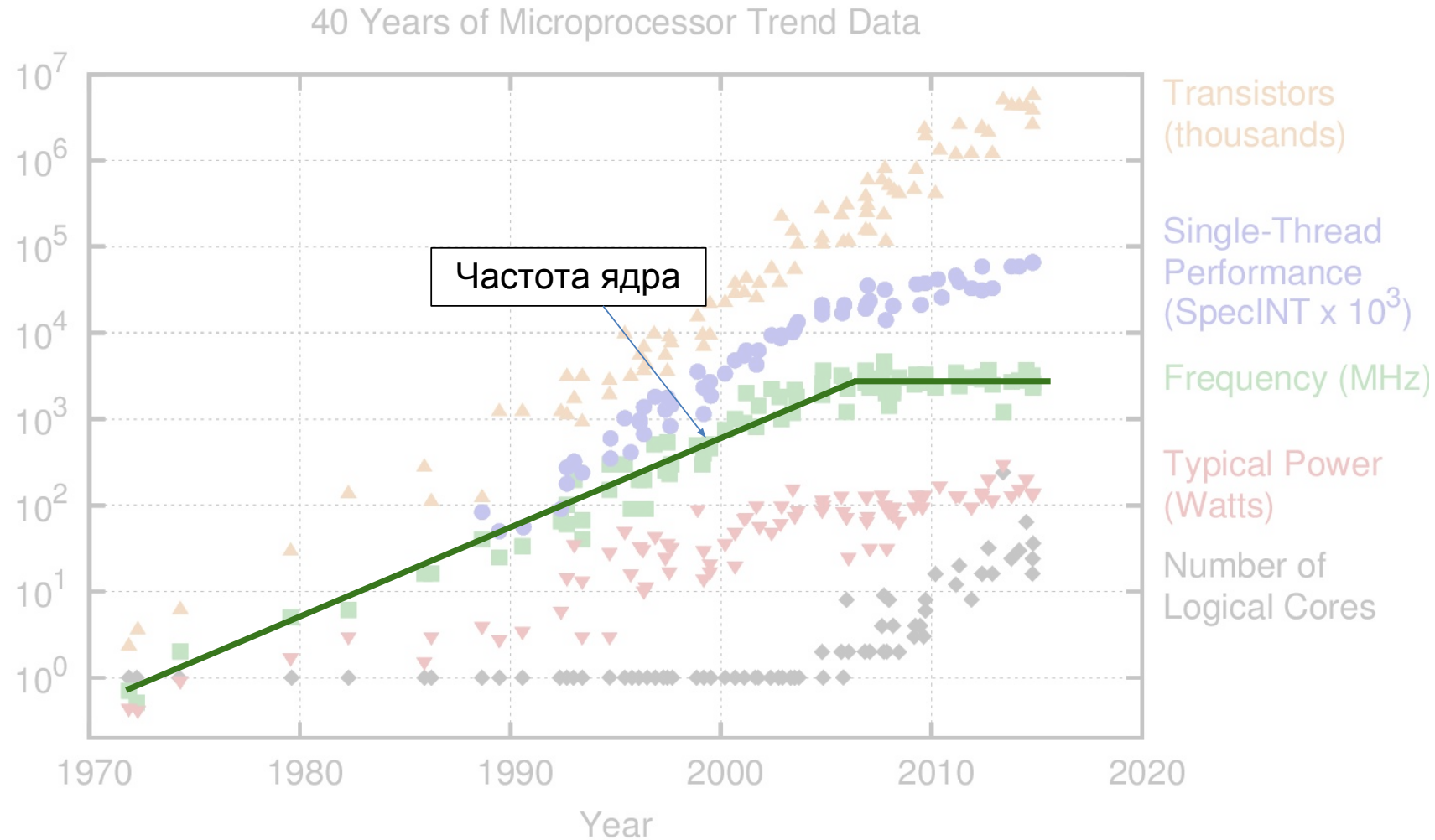
- Кто использует HashMap?
 - Наверное, это самая популярная и полезная на практике структура данных
- Но при чём тут многопоточность?

Многопоточность: а надо ли?



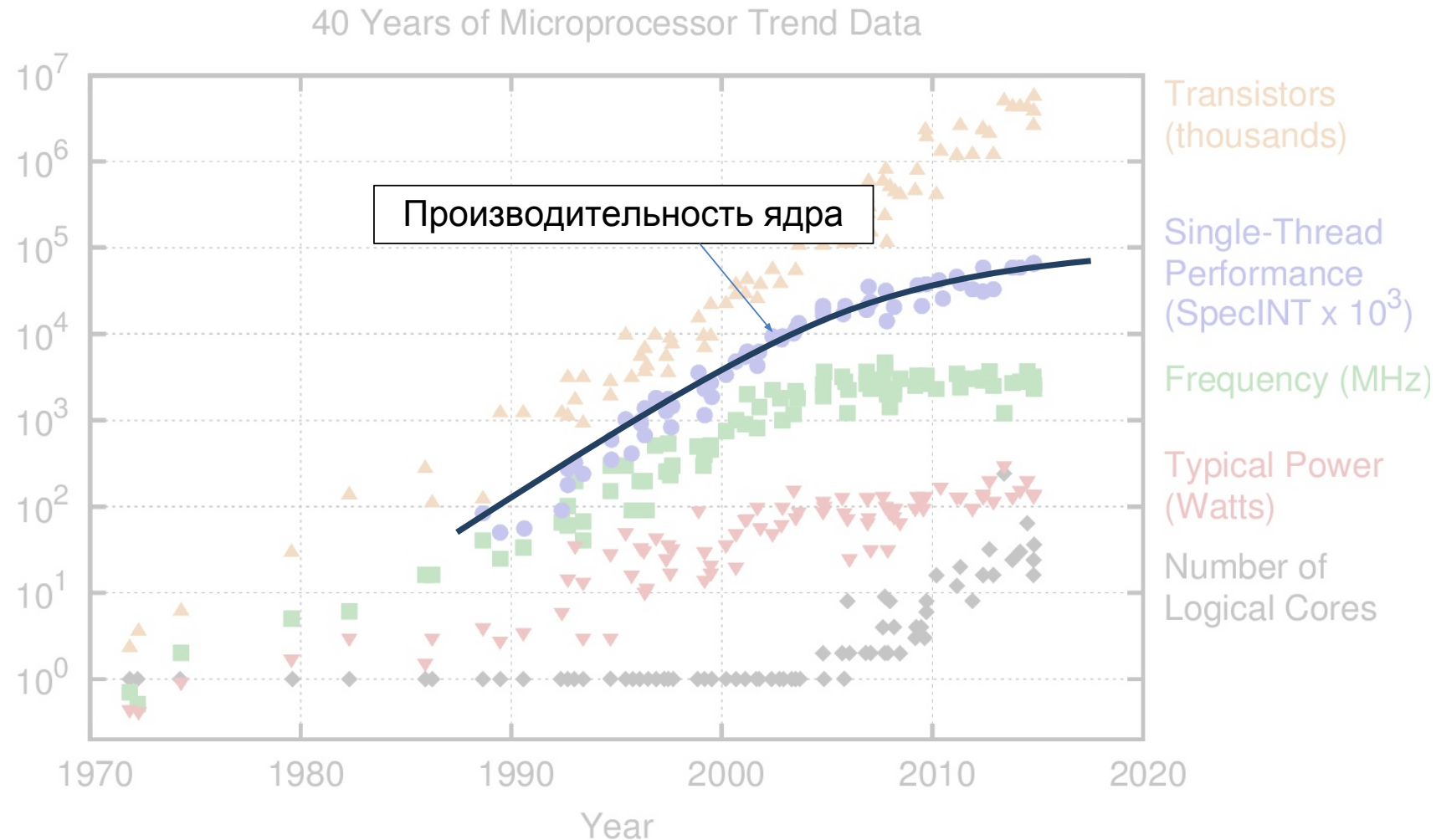
Original data up to the year 2010 collected and plotted by M. Horowitz, F. Labonte, O. Shacham, K. Olukotun, L. Hammond, and C. Batten
New plot and data collected for 2010-2015 by K. Rupp

Многопоточность: а надо ли?



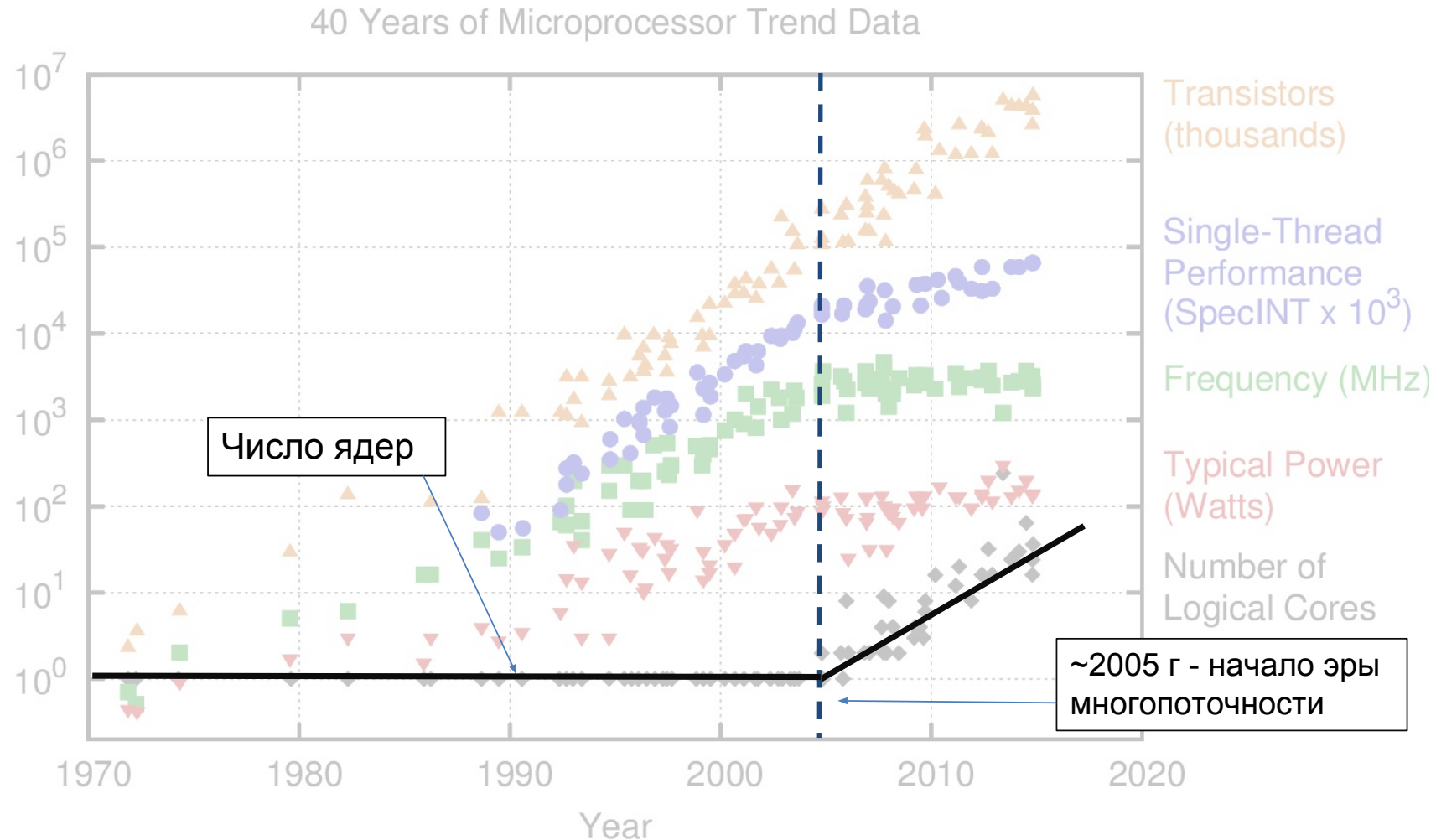
Original data up to the year 2010 collected and plotted by M. Horowitz, F. Labonte, O. Shacham, K. Olukotun, L. Hammond, and C. Batten
New plot and data collected for 2010-2015 by K. Rupp

Многопоточность: а надо ли?



Original data up to the year 2010 collected and plotted by M. Horowitz, F. Labonte, O. Shacham, K. Olukotun, L. Hammond, and C. Batten
New plot and data collected for 2010-2015 by K. Rupp

Многопоточность: а надо ли?



Original data up to the year 2010 collected and plotted by M. Horowitz, F. Labonte, O. Shacham, K. Olukotun, L. Hammond, and C. Batten
New plot and data collected for 2010-2015 by K. Rupp

И что?

- Хеш-таблица – одна из самых популярных структур данных в приложениях
- Для увеличения производительности нужно писать многопоточный код

И что?

- Хеш-таблица – одна из самых популярных структур данных в приложениях
- Для увеличения производительности нужно писать многопоточный код

Нужна быстрая многопоточная хеш-таблица!

Модель: кешируем Account-ы

- Часто необходимо сопоставлять id (*long*) с каким-то объектом (*Object*)
 - При этом большинство операций – чтение
 - Время от времени – добавление

Модель: кешируем Account-ы

```
class Account(val id: Int) {  
    val i1: Int = 0  
    val i2: Int = 0  
    val o1: Any? = null  
    val o2: Any? = null  
    val o3: Any? = null  
    val o4: Any? = null  
}
```

Модель: кешируем Account-ы

```
class Account(val id: Int) {  
    val i1: Int = 0  
    val i2: Int = 0  
    val o1: Any? = null  
    val o2: Any? = null  
    val o3: Any? = null  
    val o4: Any? = null  
}
```

```
interface AccountsCache {  
    fun addAccount(id: Int, account: Account)  
    fun getById(id: Int): Account?  
}
```

Как измерять?

Используем JMN:

- Стоимость отдельной операции ничего не значит, т.к. при перехешировании всё замедляется

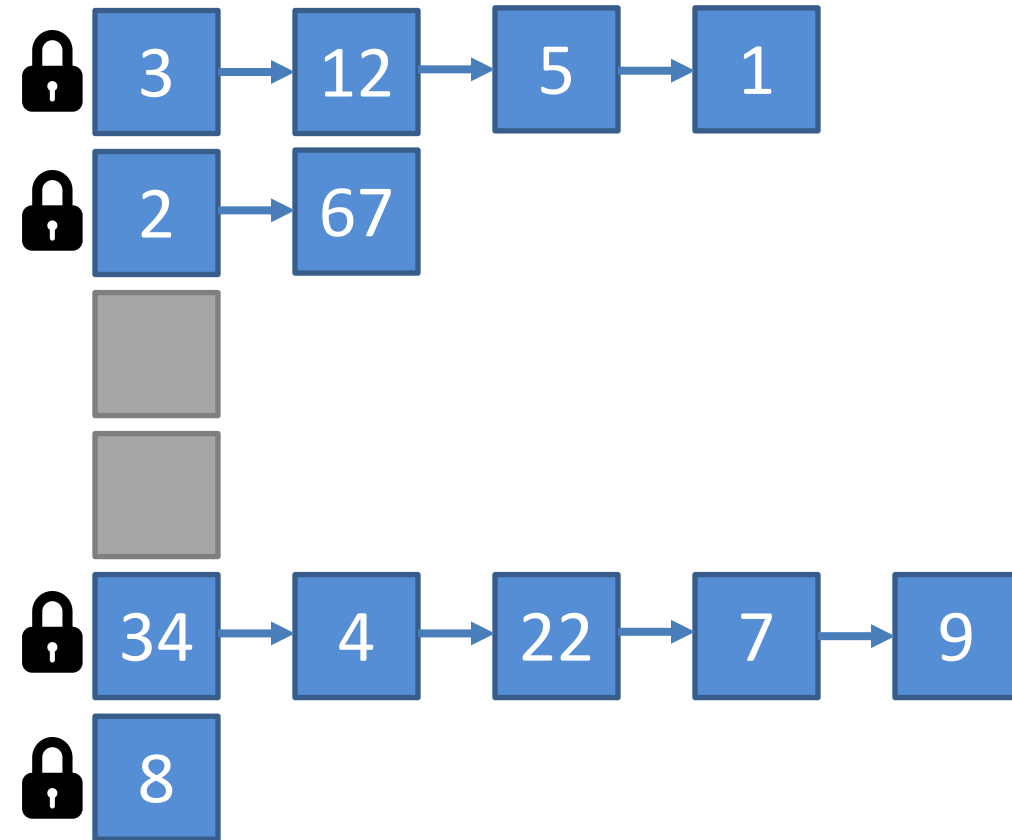
Как измерять?

Используем JMH:

- Стоимость отдельной операции ничего не значит, т.к. при перехешировании всё замедляется
- \Rightarrow имеет смысл считать стоимость пачки операций
`@Measurement(batchSize = 1000_000)`

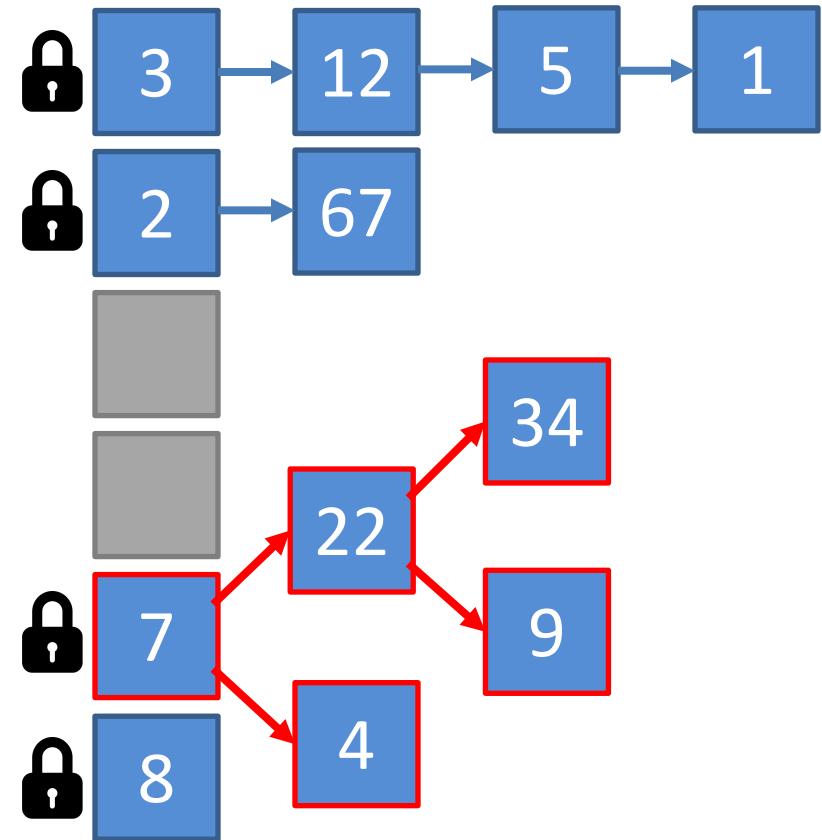
Что есть в Java: ConcurrentHashMap

- Работает за $O(1)$ в среднем
- Использует внутри блокировки
 - \Rightarrow не так хорошо масштабируется
- Хранит списки при коллизии
 - \Rightarrow лишние cache miss-ы



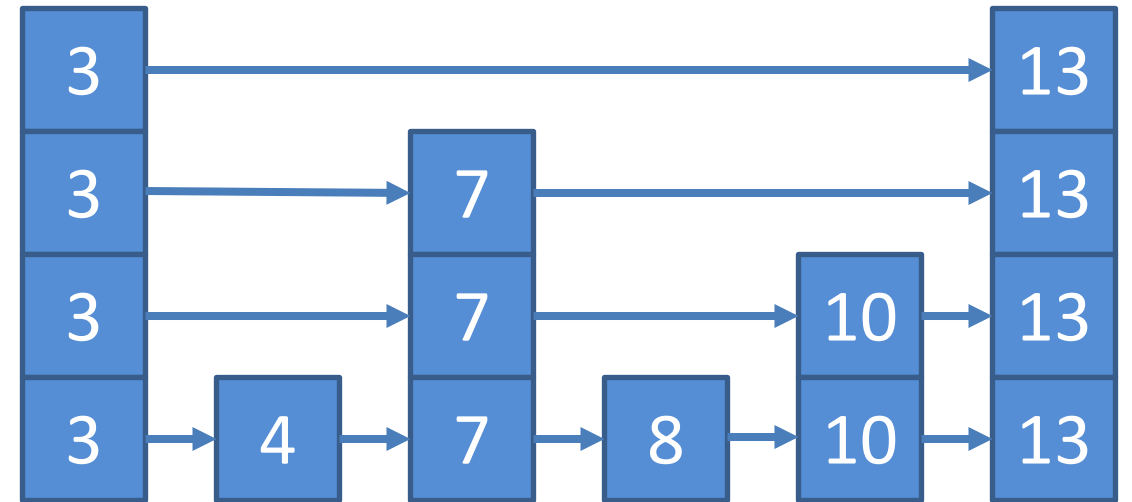
Что есть в Java: ConcurrentHashMap

- Работает за $O(1)$ в среднем
- Использует внутри блокировки
 - \Rightarrow не так хорошо масштабируется
- Хранит списки при коллизии
 - \Rightarrow лишние cache miss-ы
 - зато может хранить дерево при постоянных коллизиях*



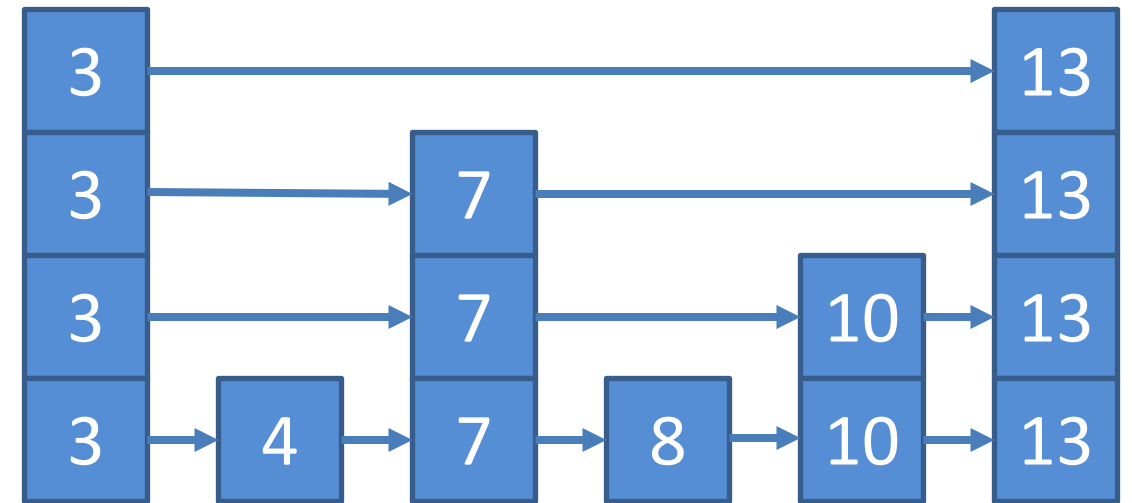
Что ещё есть в Java: ConcurrentSkipListMap

- Работает за $O(\log N)$ в среднем
- «Дерево» списков



Что ещё есть в Java: ConcurrentSkipListMap

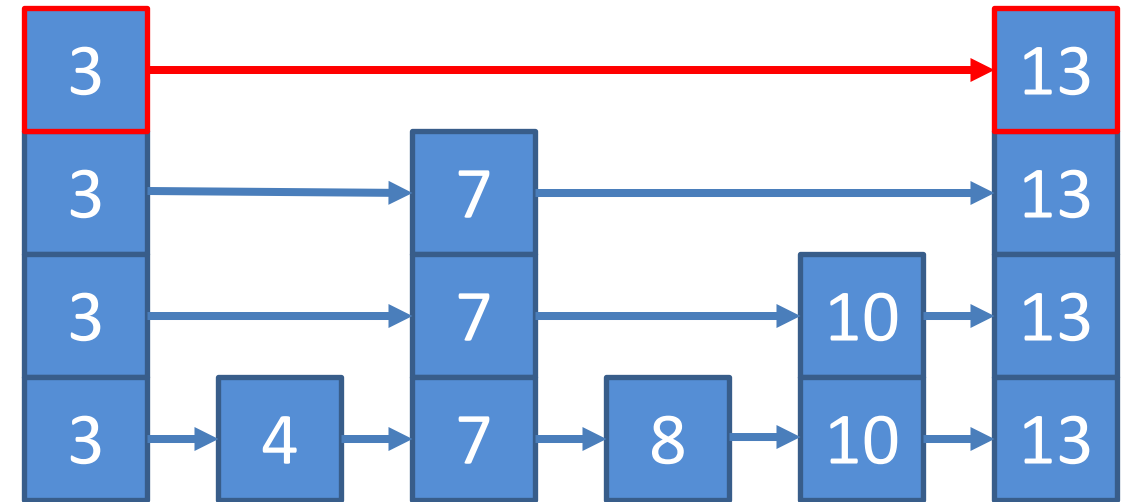
- Работает за $O(\log N)$ в среднем
- «Дерево» списков



Ищем ключ «8»

Что ещё есть в Java: ConcurrentSkipListMap

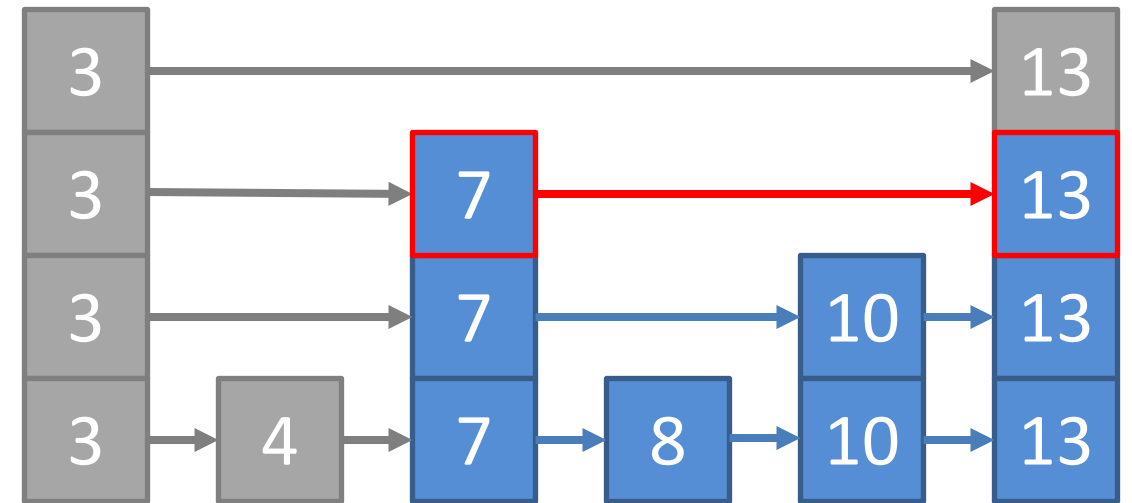
- Работает за $O(\log N)$ в среднем
- «Дерево» списков



Ищем ключ «8»

Что ещё есть в Java: ConcurrentSkipListMap

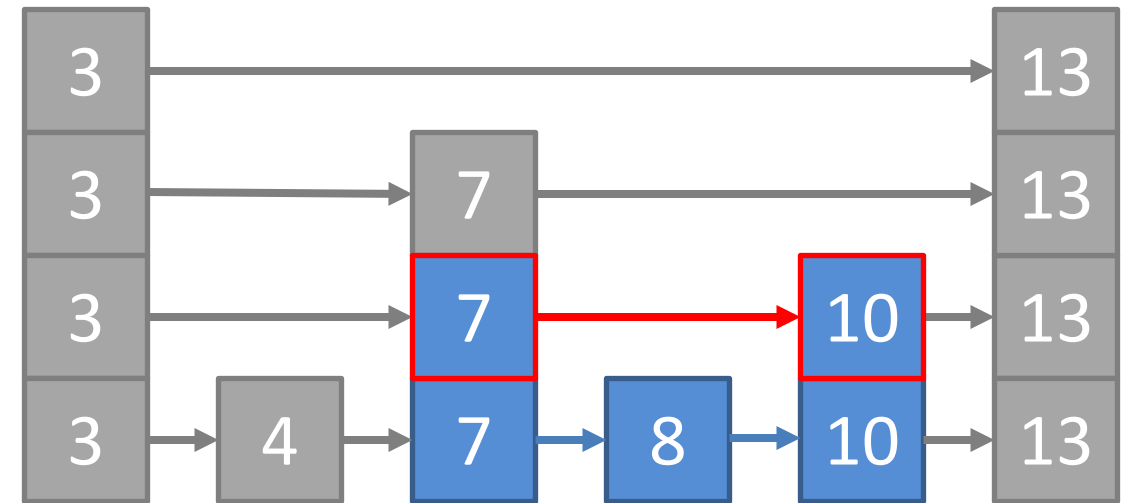
- Работает за $O(\log N)$ в среднем
- «Дерево» списков



Ищем ключ «8»

Что ещё есть в Java: ConcurrentSkipListMap

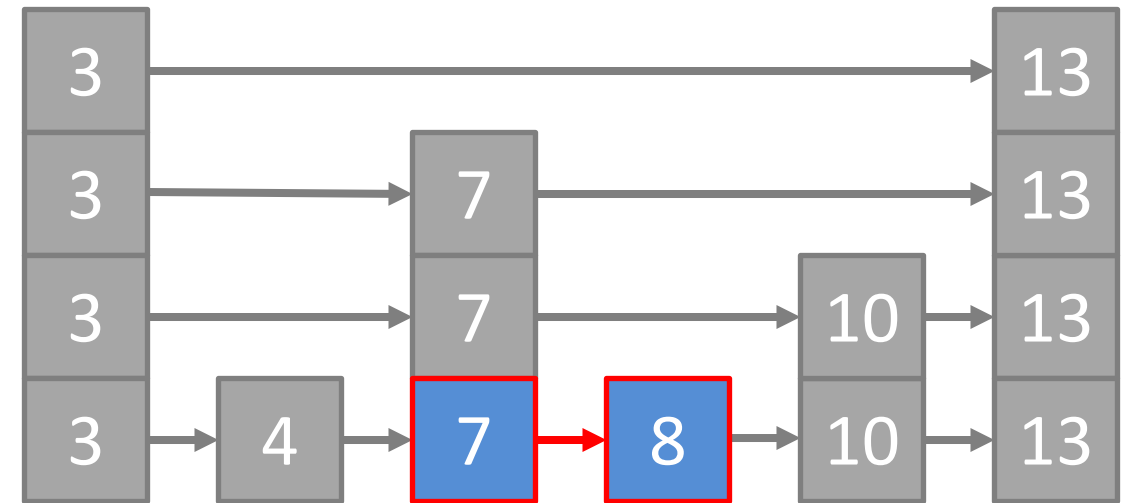
- Работает за $O(\log N)$ в среднем
- «Дерево» списков



Ищем ключ «8»

Что ещё есть в Java: ConcurrentSkipListMap

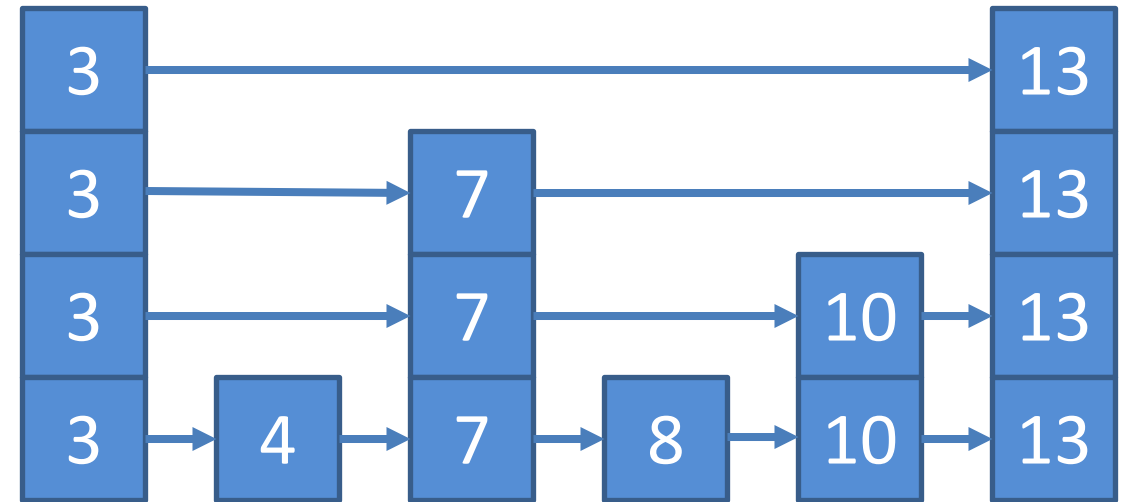
- Работает за $O(\log N)$ в среднем
- «Дерево» списков



Ищем ключ «8»

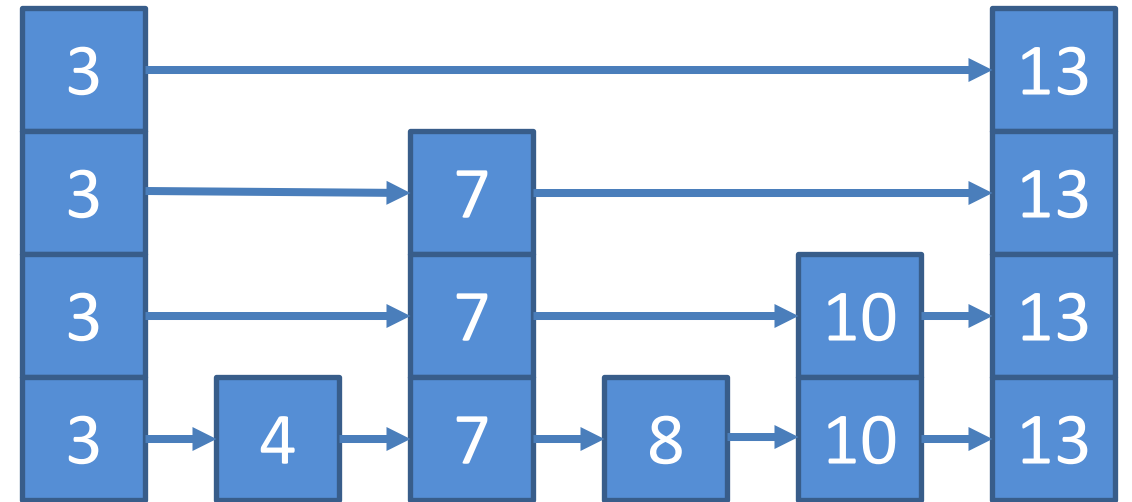
Что ещё есть в Java: ConcurrentSkipListMap

- Работает за $O(\log N)$ в среднем
- «Дерево» списков
 - Много лишних объектов
 - \Rightarrow постоянные cache miss-ы и нагружает GC



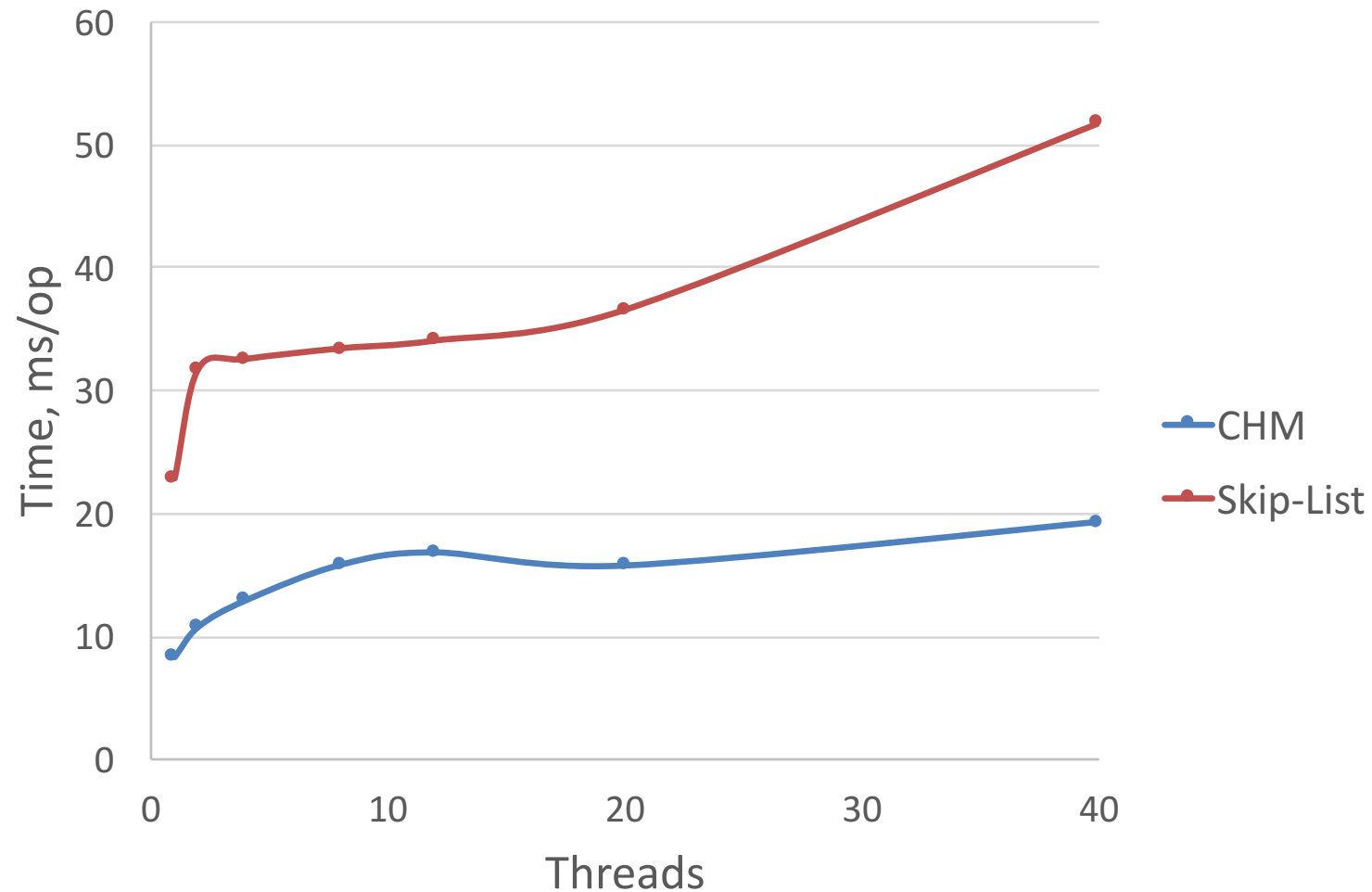
Что ещё есть в Java: ConcurrentSkipListMap

- Работает за $O(\log N)$ в среднем
- «Дерево» списков
 - Много лишних объектов
 - \Rightarrow постоянные cache miss-ы и нагружает GC

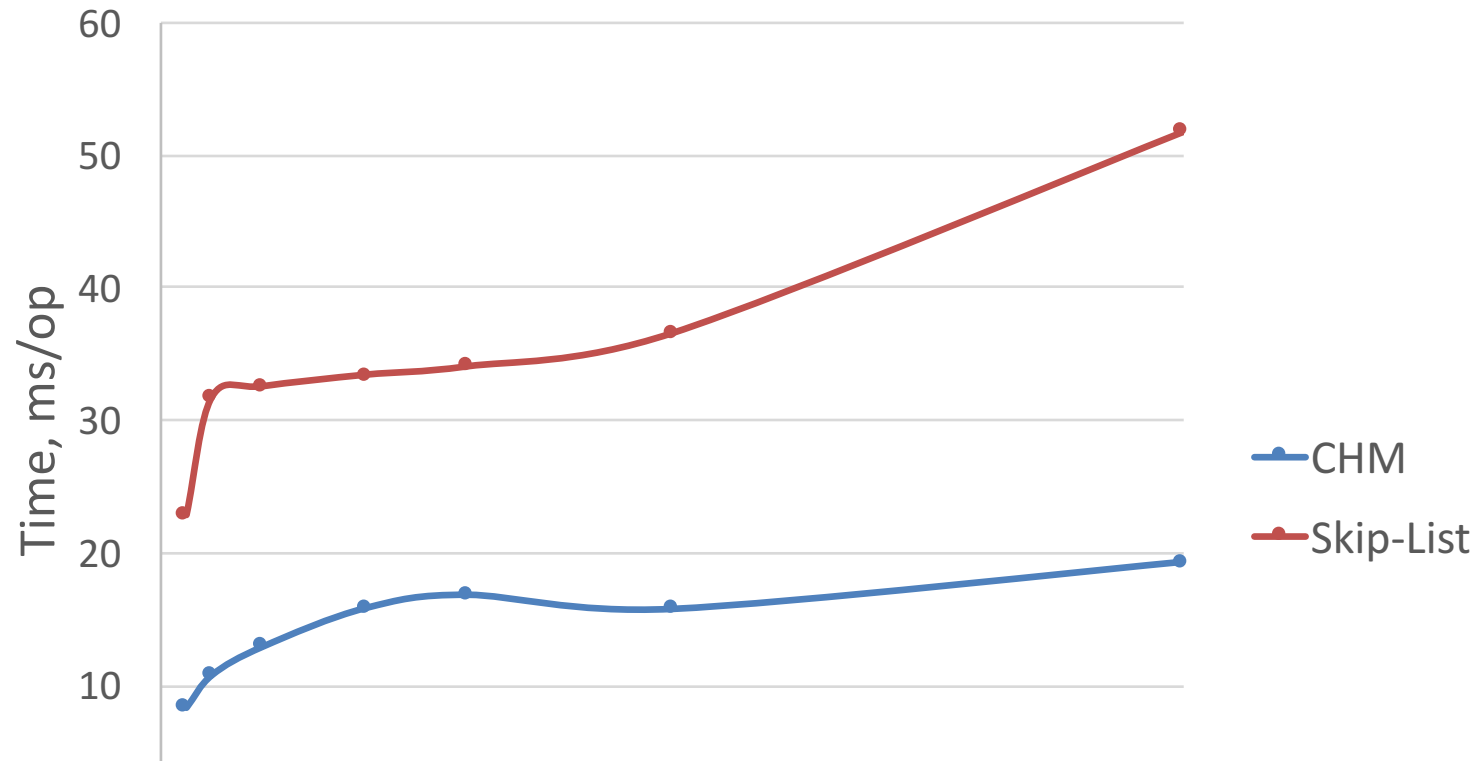


- И, вообще, это не хеш-таблица!

ConcurrentHashMap vs ConcurrentSkipListMap



ConcurrentHashMap vs ConcurrentSkipListMap



Даёшь больше производительности!

Потенциальные проблемы

- ConcurrentHashMap
 - Использует внутри блокировки \Rightarrow не lock-free

Потенциальные проблемы

- ConcurrentHashMap
 - Использует внутри блокировки \Rightarrow не lock-free



Важно для высоконагруженных систем!

Потенциальные проблемы

- ConcurrentHashMap
 - Использует внутри блокировки \Rightarrow не lock-free
 - Использует списки \Rightarrow cache miss-ы

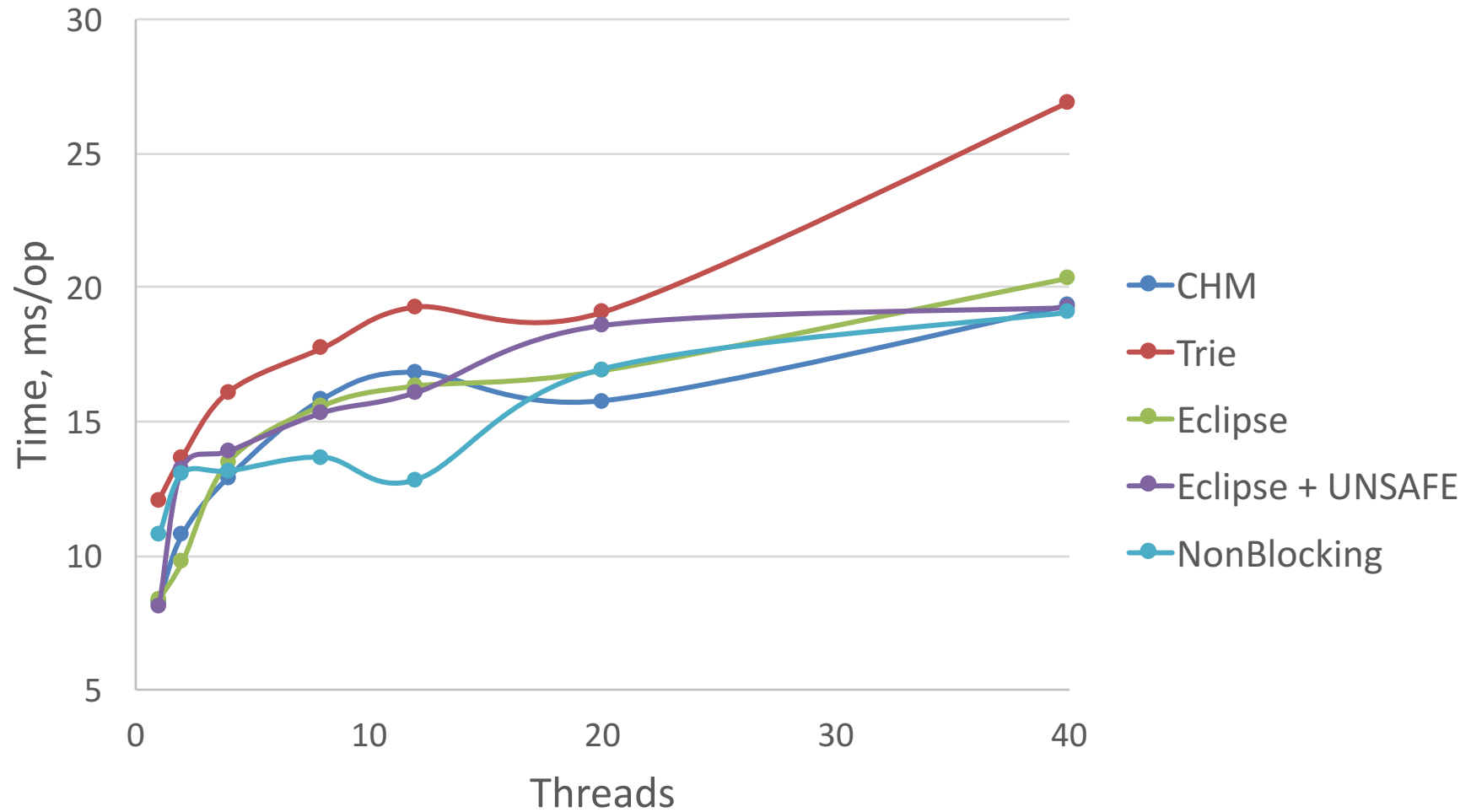
Потенциальные проблемы

- ConcurrentHashMap
 - Использует внутри блокировки \Rightarrow не lock-free
 - Использует списки \Rightarrow cache miss-ы
- ConcurrentSkipListMap
 - Много «лишних» объектов
 - Постоянные cache miss-ы

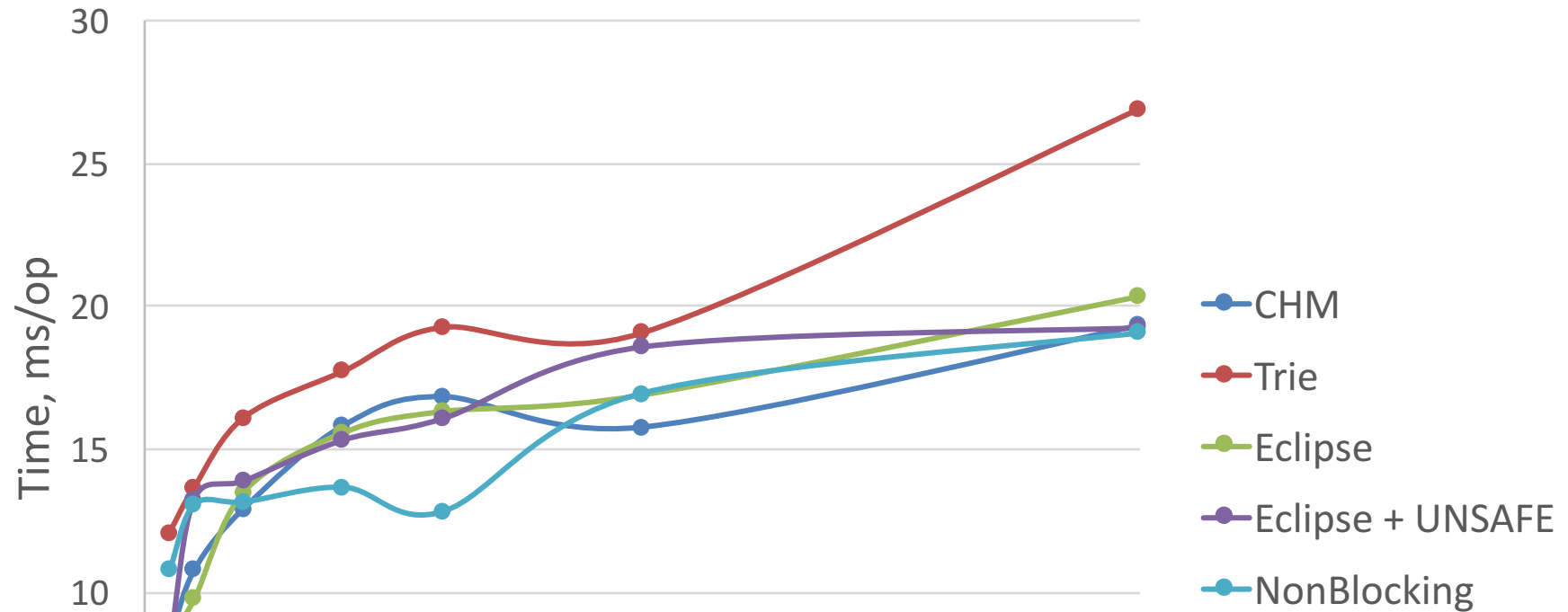
Альтернативы

- <https://github.com/edwardw/high-scale-java-lib>
 - NonBlockingHashMap
- <https://github.com/romix/java-concurrent-hash-trie-map>
- <https://github.com/eclipse/eclipse-collections>
 - ConcurrentHashMap и ConcurrentHashMapUnsafe
- <https://github.com/google/guava>
- <https://github.com/javolution/javolution>
- ...

Альтернативы: бенчмарки

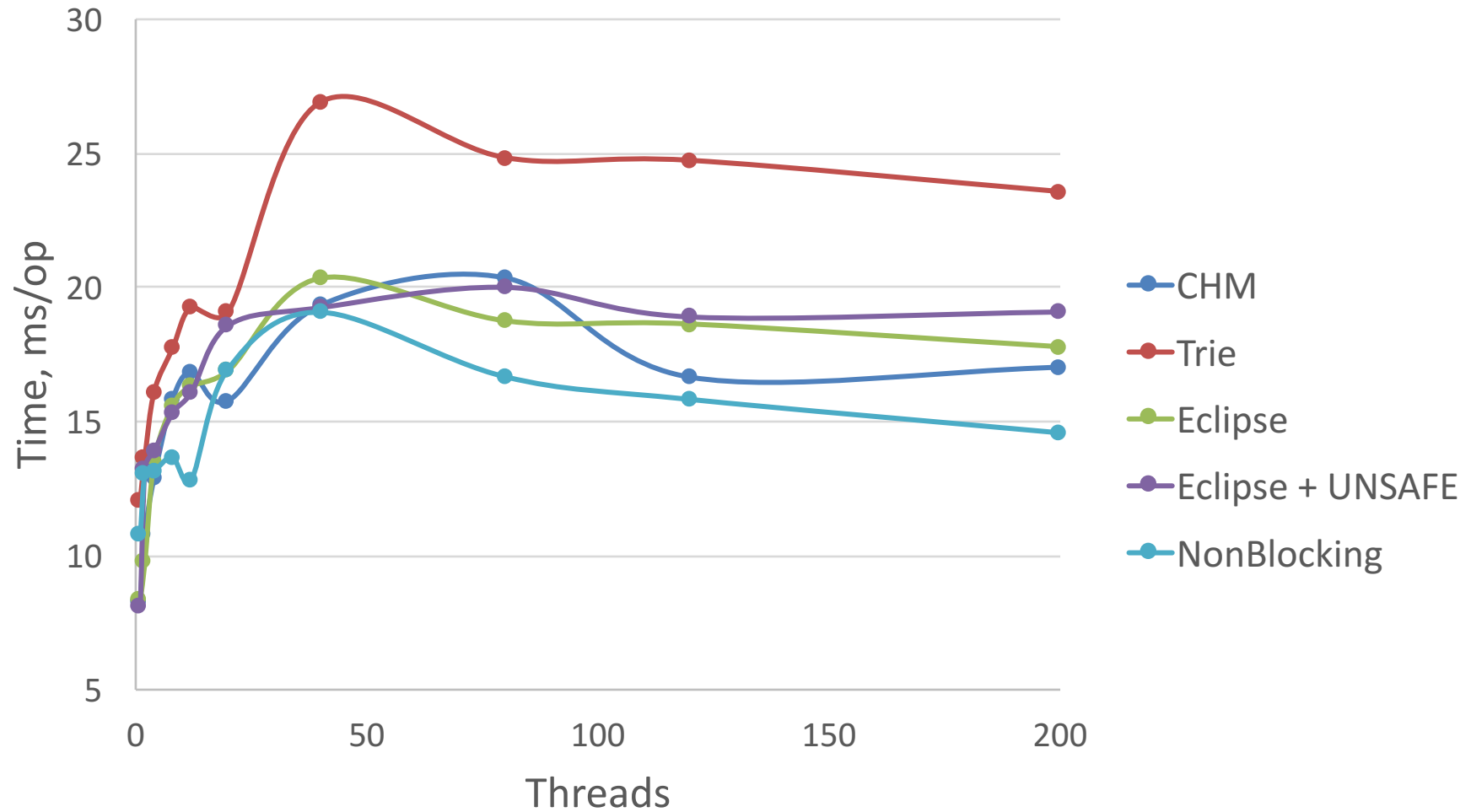


Альтернативы: бенчмарки

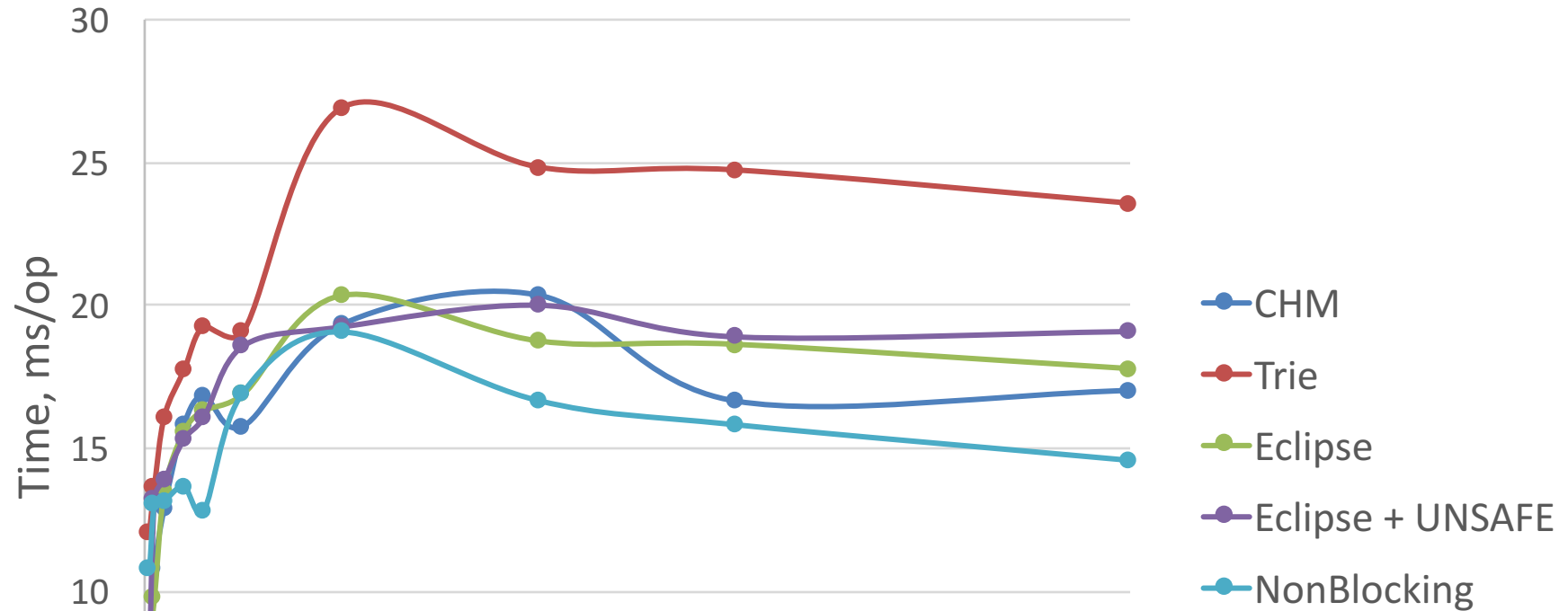


А если потоков больше, чем ядер?

Альтернативы: бенчмарки

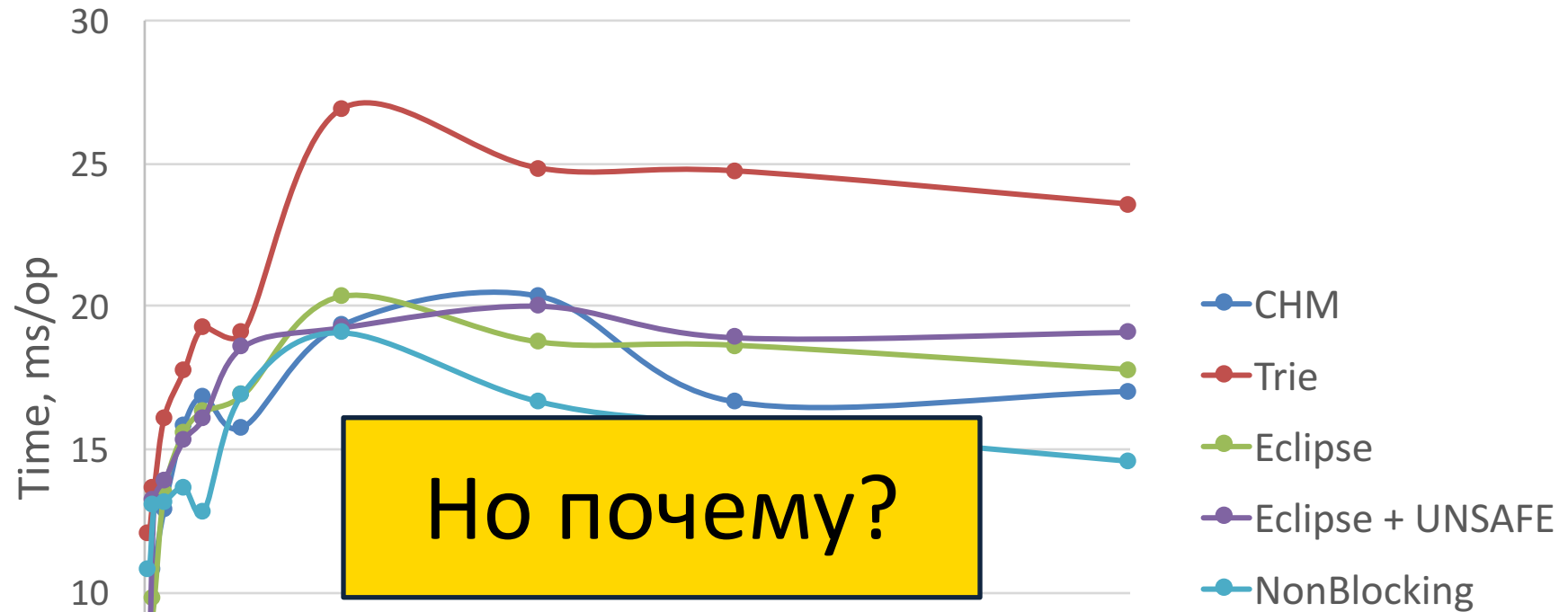


Альтернативы: бенчмарки



NonBlocking чуть-чуть всех рвёт

Альтернативы: бенчмарки



NonBlocking чуть-чуть всех рвёт

NonBlockingHashMap

Lock-free Dynamic Hash Tables with Open Addressing

Gao, H.¹, Groote, J.F.², Hesselink, W.H.¹

¹ Department of Mathematics and Computing Science, University of Groningen,
P.O. Box 800, 9700 AV Groningen, The Netherlands (Email: {hui,wim}@cs.rug.nl)

² Department of Mathematics and Computing Science, Eindhoven University of
Technology, P.O. Box 513, 5600 MB Eindhoven, The Netherlands and CWI,
P.O. Box 94079, 1090 GB Amsterdam, The Netherlands (Email: jfg@win.tue.nl)

Abstract

We present an efficient lock-free algorithm for parallel accessible hash tables with open addressing, which promises more robust performance and reliability than conventional lock-based implementations. “Lock-free” means that it is guaranteed that always at least one process completes its operation within a bounded number of steps. For a single processor architecture our solution is as efficient as sequential hash tables. On a multiprocessor architecture this is also the case when all processors have comparable speeds. The algorithm allows processors that have widely different speeds or come to a halt. It can easily be implemented using C-like languages and requires on average only constant time for insertion, deletion or accessing of elements. The algorithm allows the hash tables to grow and shrink when needed.

Lock-free algorithms are hard to design correctly, even when apparently straightforward. Ensuring the correctness of the design at the earliest possible stage is a major challenge in any responsible system development. In view of the complexity of the algorithm, we turned to the interactive theorem prover PVS for mechanical support. We employ standard deductive verification techniques to prove around 200 invariance properties of our algorithm, and describe how this is achieved with the theorem prover PVS.

CR Subject Classification (1991): D.1 Programming techniques

AMS Subject Classification (1991): 68Q22 Distributed algorithms, 68P20 Information storage and retrieval

Keywords & Phrases: Hash tables, Distributed algorithms, Lock-free, Wait-free



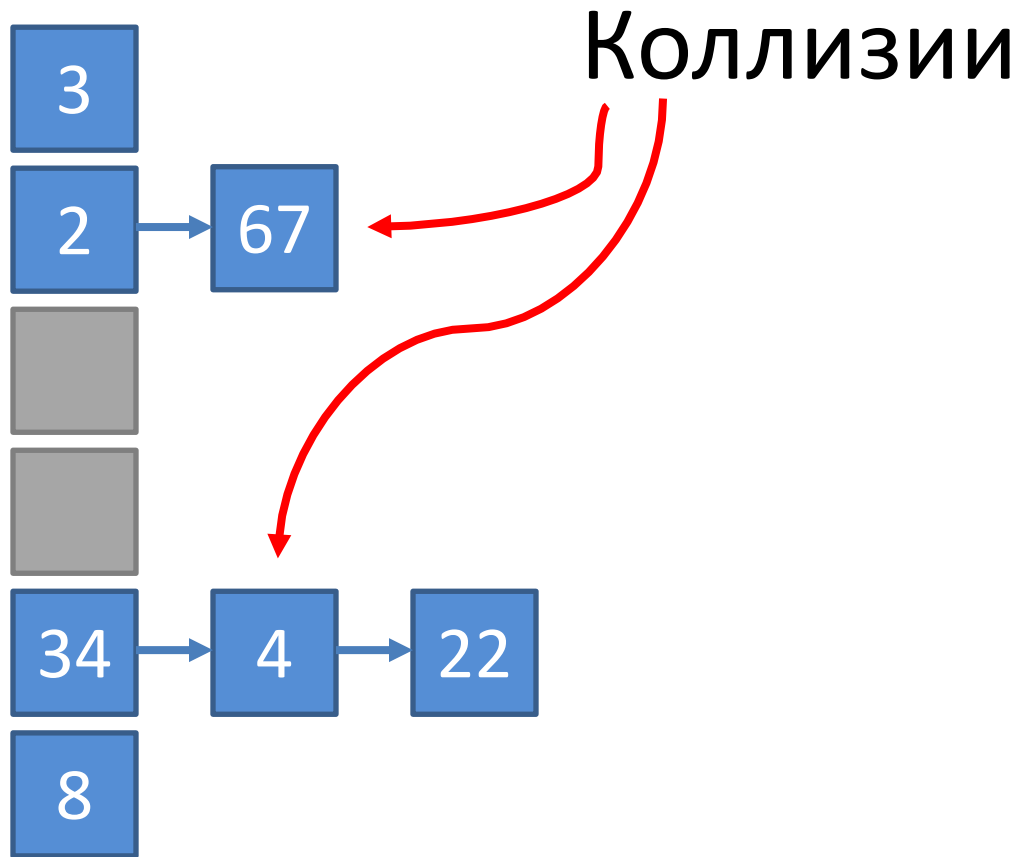
[Lock-free Dynamic Hash Tables
with Open Addressing](#) by Gao et al.

[A Lock-Free Wait-Free Hash Table](#) by C. Click

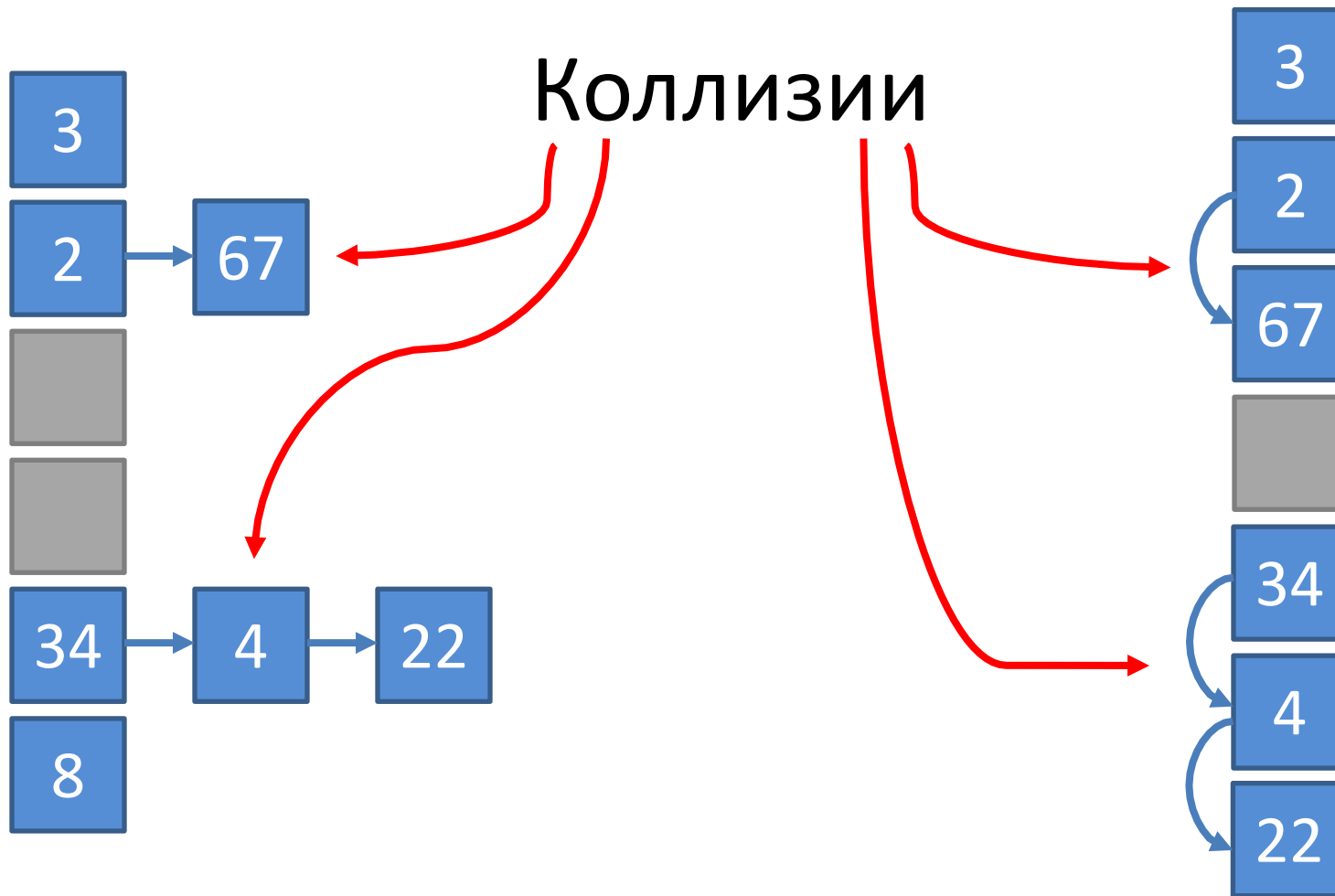
NonBlockingHashMap

- Использует открытую адресацию

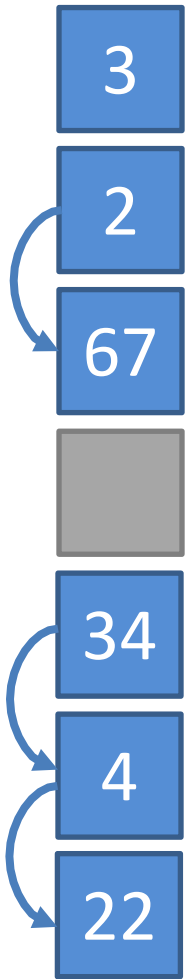
Открытая и закрытая адресации



Открытая и закрытая адресации

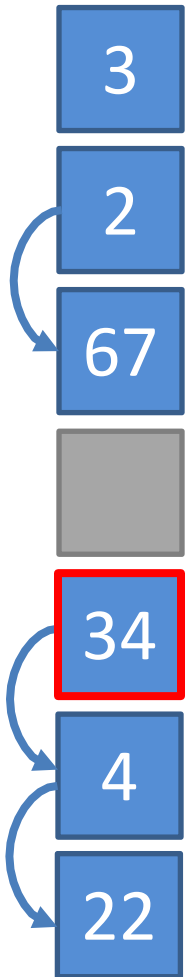


Открытая адресация: поиск



Ищем «22»

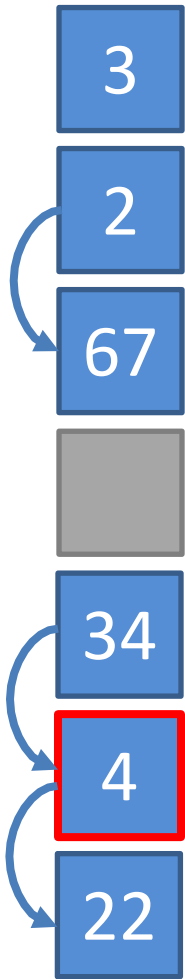
Открытая адресация: поиск



Ищем «22»:

1. Пришли в ячейку с «34»

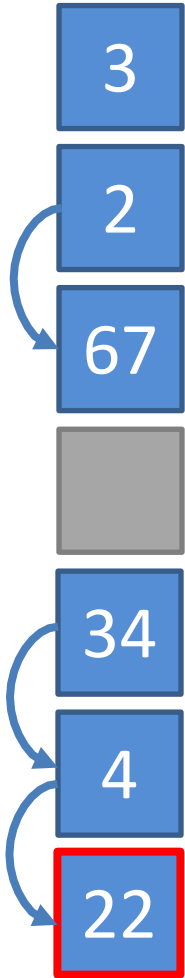
Открытая адресация: поиск



Ищем «22»:

1. Пришли в ячейку с «34»
2. Идём дальше пока не найдём «22»

Открытая адресация: поиск

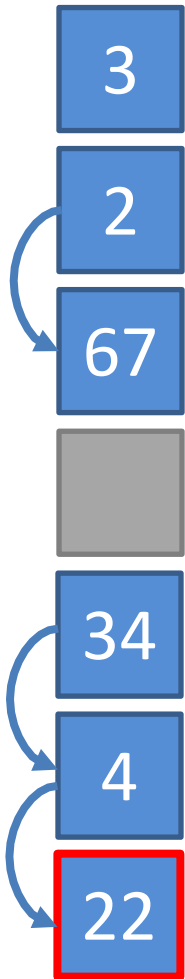


Ищем «22»:

1. Пришли в ячейку с «34»
2. Идём дальше пока не найдём «22»

Ура, нашли!

Открытая адресация: поиск



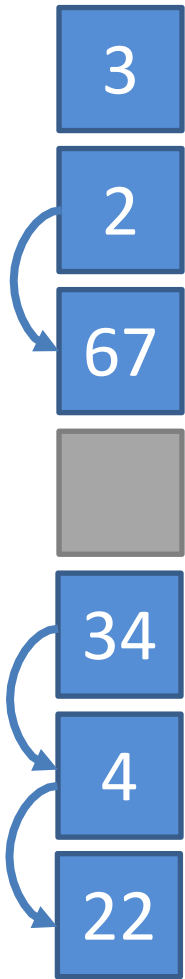
Ищем «22»:

1. Пришли в ячейку с «34»
2. Идём дальше пока не найдём «22»

Ура, нашли!

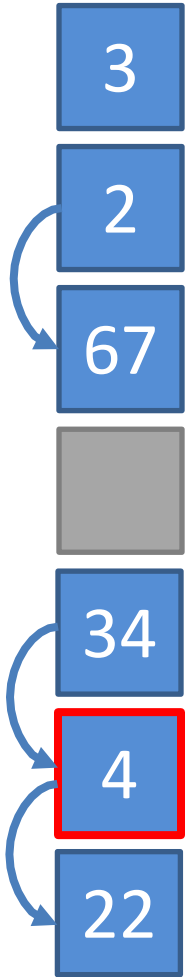
А если бы «4» была удалена?

Открытая адресация: удаление



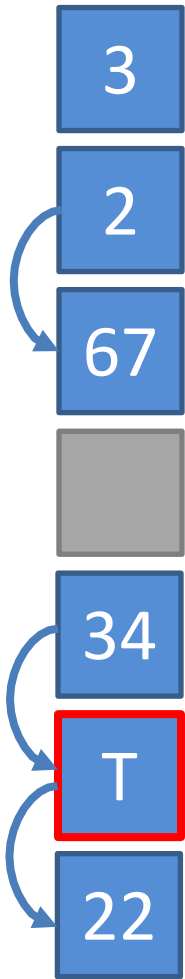
Удаляем «4»

Открытая адресация: удаление



Удаляем «4»:
1. Находим ячейку

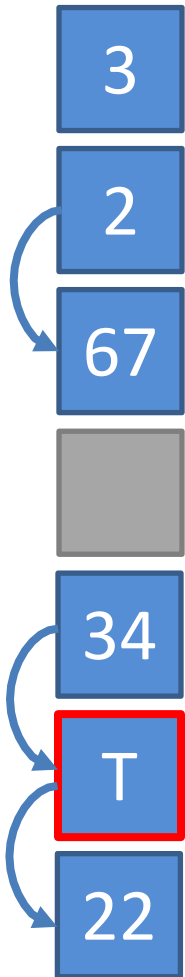
Открытая адресация: удаление



Удаляем «4»:

1. Находим ячейку
2. Меняем на специальное значение «раньше тут был элемент»

Открытая адресация: удаление



Удаляем «4»:

1. Находим ячейку
2. Меняем на специальное значение «раньше тут был элемент»

Теперь мы сможем
найти «22»!

Пара деталей: размер таблицы

- Как выбрать правильный размер таблицы?
- Теория говорит «давайте возьмём простое число»

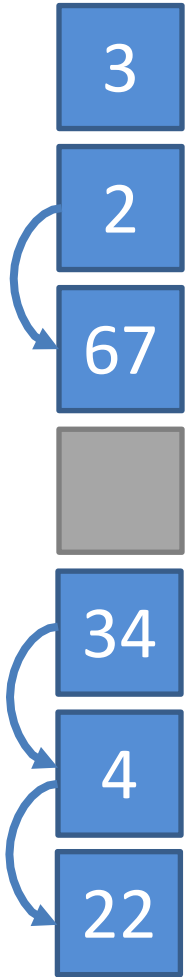
Пара деталей: размер таблицы

- Как выбрать правильный размер таблицы?
- Теория говорит «давайте возьмём простое число»
 - Тогда нужно всегда делать $\text{hash}(\text{key}) \% N$

Пара деталей: размер таблицы

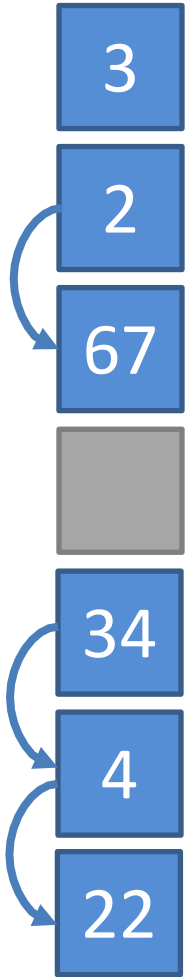
- Как выбрать правильный размер таблицы?
- Теория говорит «давайте возьмём простое число»
 - Тогда нужно всегда делать $\text{hash}(\text{key}) \% N$
 - Операция взятия по модулю слишком дорогая
 - На практике степень двойки лучше (это просто сдвиг)

Пара деталей: поиск элемента



- Можно смотреть не на следующий элемент
 - И в теории это лучше!

Пара деталей: поиск элемента



- Можно смотреть не на следующий элемент
 - И в теории это лучше!
- На практике посмотреть на следующий элемент дешевле, т.к. он скорее всего уже закеширован

Открытая адресация: get(key)

```
public T getInternal(long key) {  
    int i = index(key);  
    long k;  
    int probes = 0;  
    while ((k = keys.get(i)) != key) {  
        if (k == NULL_KEY)  
            return null;  
        if (++probes >= MAX_PROBES)  
            return null;  
        if (i == 0)  
            i = length;  
        i--;  
    }  
    return values.get(i);  
}
```

Открытая адресация: get(key)

```
public T getInternal(long key) {  
    int i = index(key);  
    long k;  
    int probes = 0;  
    while ((k = keys.get(i)) != key) {  
        if (k == NULL_KEY)  
            return null;  
        if (++probes >= MAX_PROBES)  
            return null;  
        if (i == 0)  
            i = length;  
        i--;  
    }  
    return values.get(i);  
}
```

Открытая адресация: get(key)

```
public T getInternal(long key) {  
    int i = index(key);  
    long k;  
    int probes = 0;  
    while ((k = keys.get(i)) != key) {  
        if (k == NULL_KEY)  
            return null;  
        if (++probes >= MAX_PROBES)  
            return null;  
        if (i == 0)  
            i = length;  
        i--;  
    }  
    return values.get(i);  
}
```

Открытая адресация: get(key)

```
public T getInternal(long key) {  
    int i = index(key);  
    long k;  
    int probes = 0;  
    while ((k = keys.get(i)) != key) {  
        if (k == NULL_KEY)  
            return null;  
        if (++probes >= MAX_PROBES)  
            return null;  
        if (i == 0)  
            i = length;  
        i--;  
    }  
    return values.get(i);  
}
```

NonBlockingHashMap

- Использует открытую адресацию
- Lock-free
 - Гарантирует, что система не стоит на месте даже при неудачном scheduling-e
 - Можно вставлять/читать во время перехеширования

NonBlockingHashMap

- Использует открытую адресацию
- Lock-free
 - Гарантирует, что система не стоит на месте даже при неудачном scheduling-e
 - Можно вставлять/читать во время перехеширования

Как этого добиться?

MRSW хеш-таблица

- Multiple Readers, Single Writer
 - Часто только один поток меняет данные!

MRSW хеш-таблица

- Как читать в процессе увеличения?

MRSW хеш-таблица

- Как читать в процессе увеличения?
- Увеличение в случае одного писателя тривиально:
 1. Создаём новую таблицу

MRSW хеш-таблица

- Как читать в процессе увеличения?
- Увеличение в случае одного писателя тривиально:
 1. Создаём новую таблицу
 2. Копируем элементы

MRSW хеш-таблица

- Как читать в процессе увеличения?
- Увеличение в случае одного писателя тривиально:
 1. Создаём новую таблицу
 2. Копируем элементы
 3. Меняем ссылку на таблицу

MRSW хеш-таблица

- Как читать в процессе увеличения?
- Увеличение в случае одного писателя тривиально:
 1. Создаём новую таблицу
 2. Копируем элементы
 3. Меняем ссылку на таблицу

MRSW хеш-таблица: get(key)

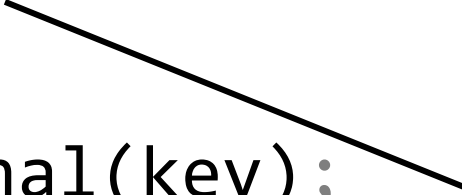
```
public class MRSWLongObjectHashMap<T> {  
    private volatile Core core = new Core(...);  
  
    public T get(long key) {  
        return core.getInternal(key);  
    }  
  
    private class Core {  
        final AtomicLongArray keys;  
        final AtomicReferenceArray<T> values;  
  
        public T getInternal(long key) { ... }  
    }  
}
```

MRSW хеш-таблица: get(key)

```
public class MRSWLongObjectHashMap<T> {  
    private volatile Core core = new Core(...);  
  
    public T get(long key) {  
        return core.getInternal(key);  
    }  
  
    private class Core {  
        final AtomicLongArray keys;  
        final AtomicReferenceArray<T> values;  
  
        public T getInternal(long key) { ... }  
    }  
}
```

MRSW хеш-таблица: get(key)

```
public class MRSWLongObjectHashMap<T> {  
    private volatile Core core = new Core(...);  
  
    public T get(long key) {  
        return core.getInternal(key);  
    }  
  
    private class Core {  
        final AtomicLongArray keys;  
        final AtomicReferenceArray<T> values;  
  
        public T getInternal(long key) { ... }  
    }  
}
```



Меняется после
перехеширования

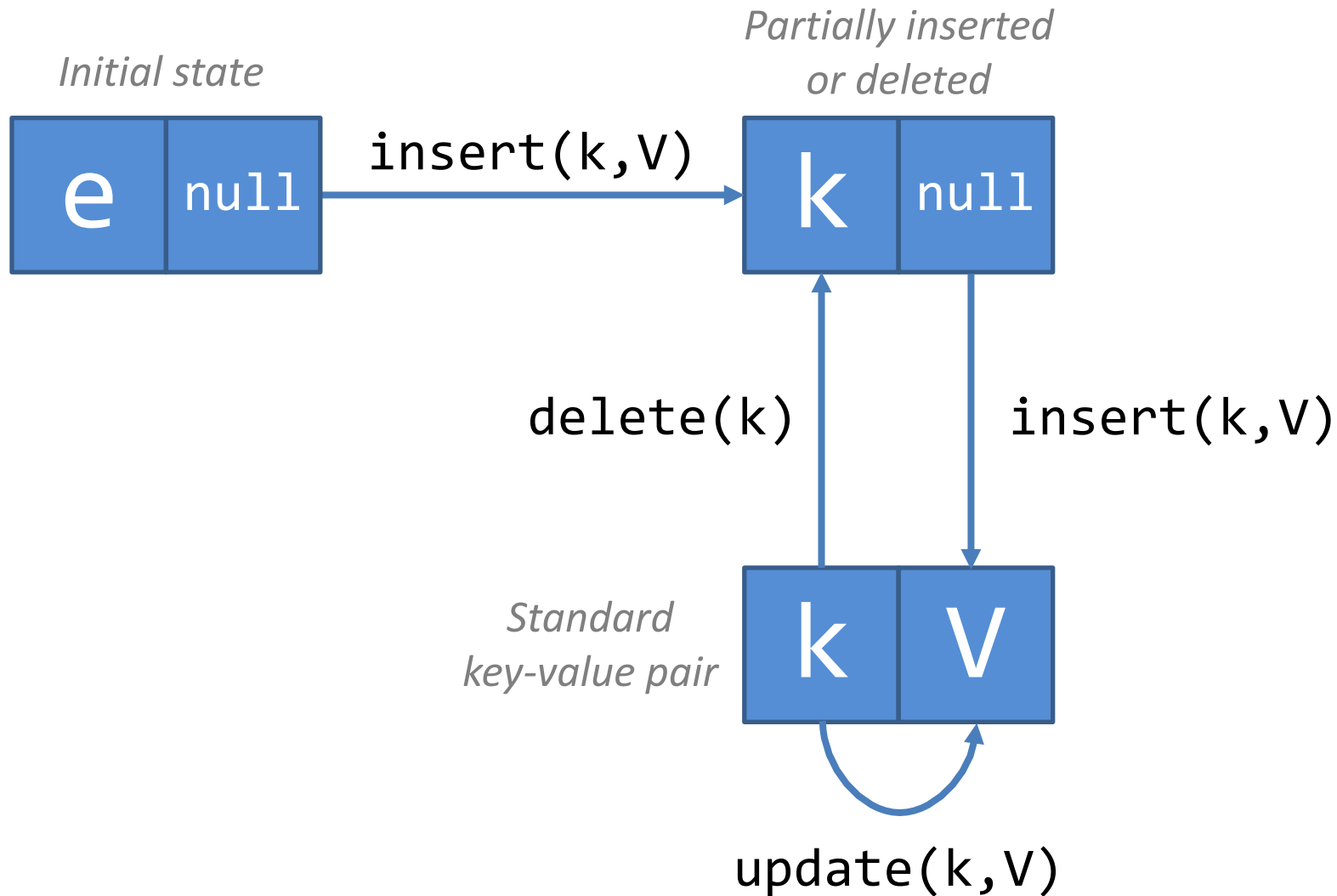
MRSW хеш-таблица: get(key)

```
public class MRSWLongObjectHashMap<T> {  
    private volatile Core core = new Core(...);  
  
    public T get(long key) {  
        return core.getInternal(key);  
    }  
  
    private class Core {  
        final AtomicLongArray keys;  
        final AtomicReferenceArray<T> values;  
  
        public T getInternal(long key) { ... }  
    }  
}
```

MRSW хеш-таблица: get(key)

```
public class MRSWLongObjectHashMap<T> {  
    private volatile Core core = new Core(...);  
  
    public T get(long key) {  
        return core.getInternal(key);  
    }  
  
    private class Core {  
        final AtomicLongArray keys;  
        final AtomicReferenceArray<T> values;  
  
        public T getInternal(long key) { ... }  
    }  
}
```

Жизнь ячейки



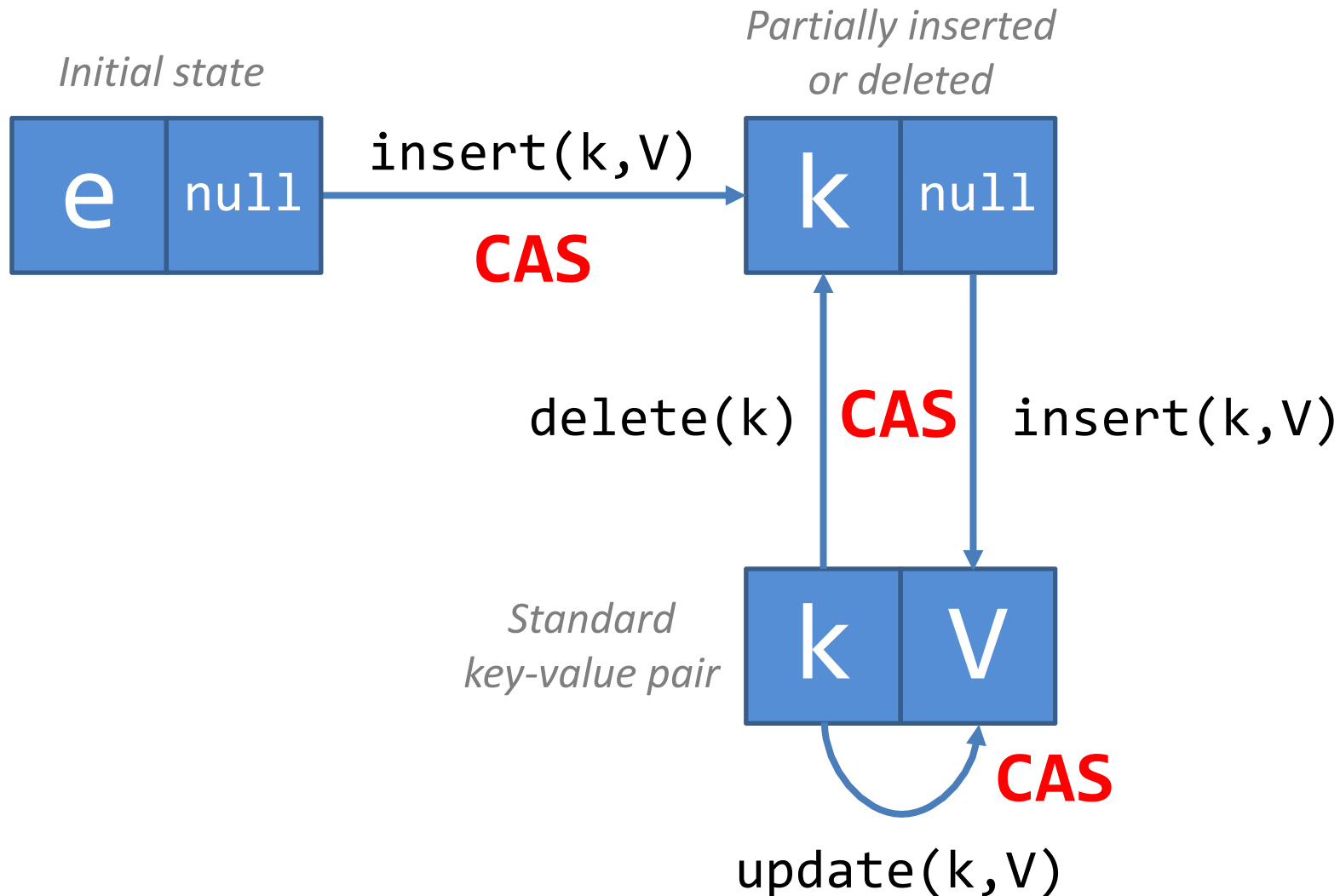
e = empty key
k = some key

V = some value
null = tombstone / empty

MRMW хеш-таблица

- Multiple Readers, Multiple Writers
 - То, чего мы хотим в общем случае
- Как организовать изменение из разных потоков?
- Как параллельно делать перенос элементов?

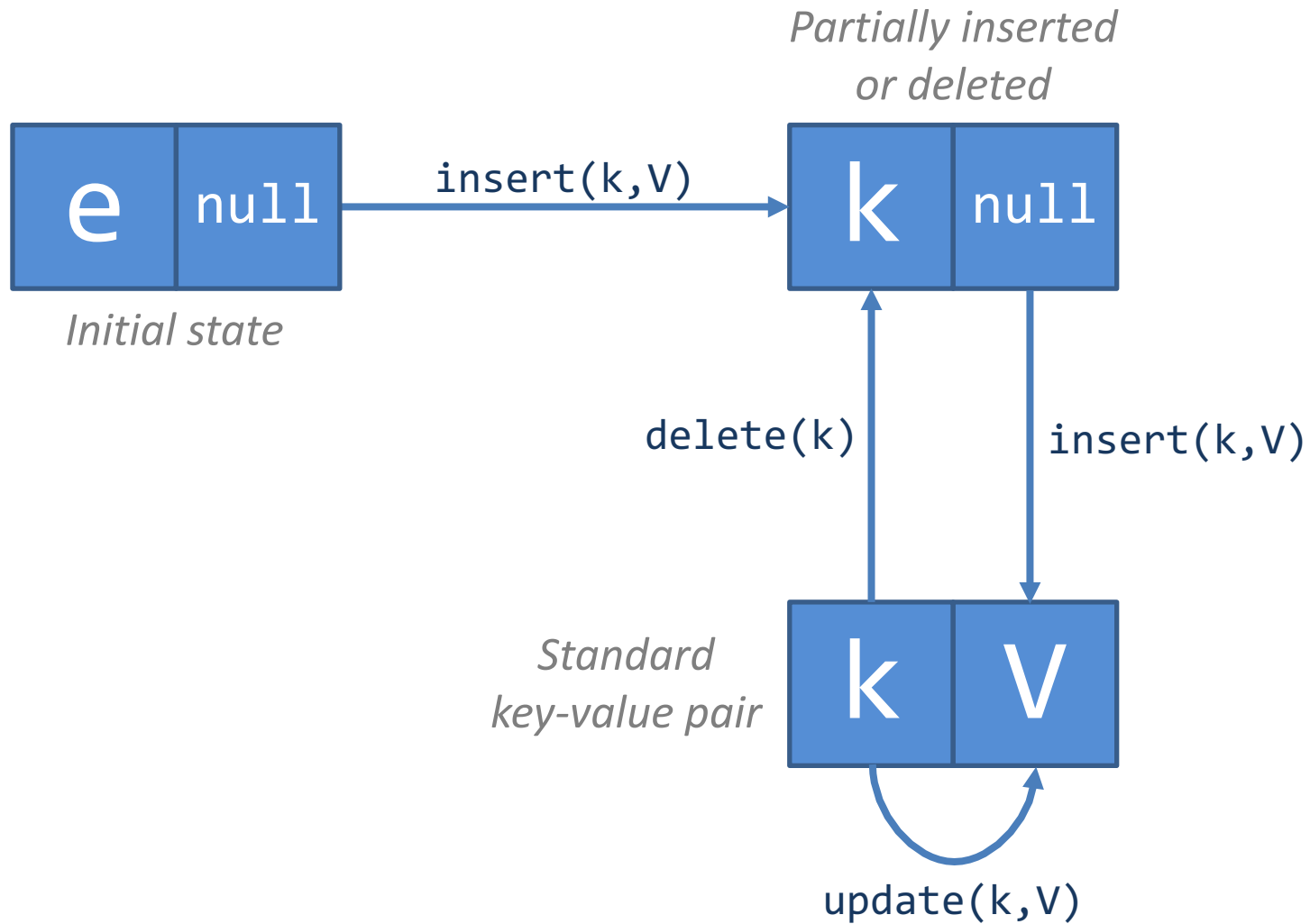
Жизнь ячейки: без переноса



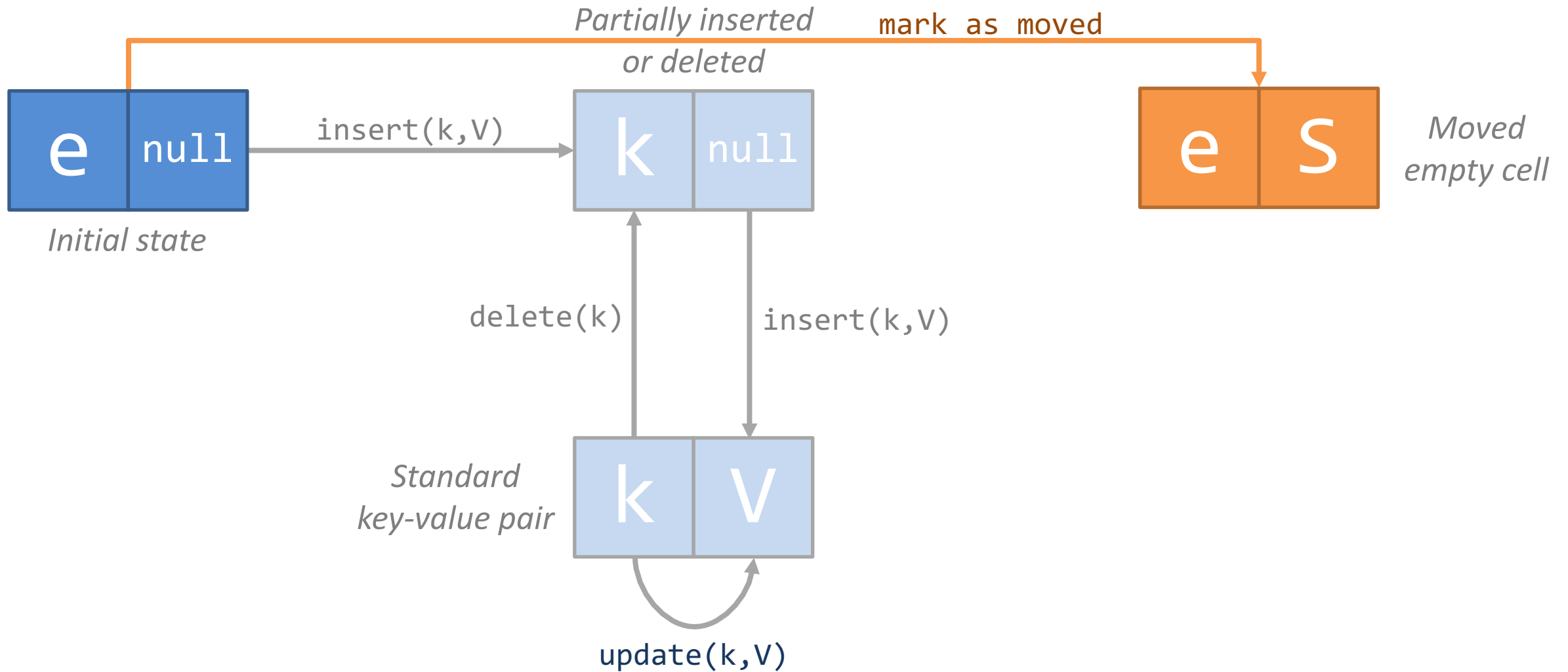
e = empty key
k = some key

V = some value
null = tombstone / empty

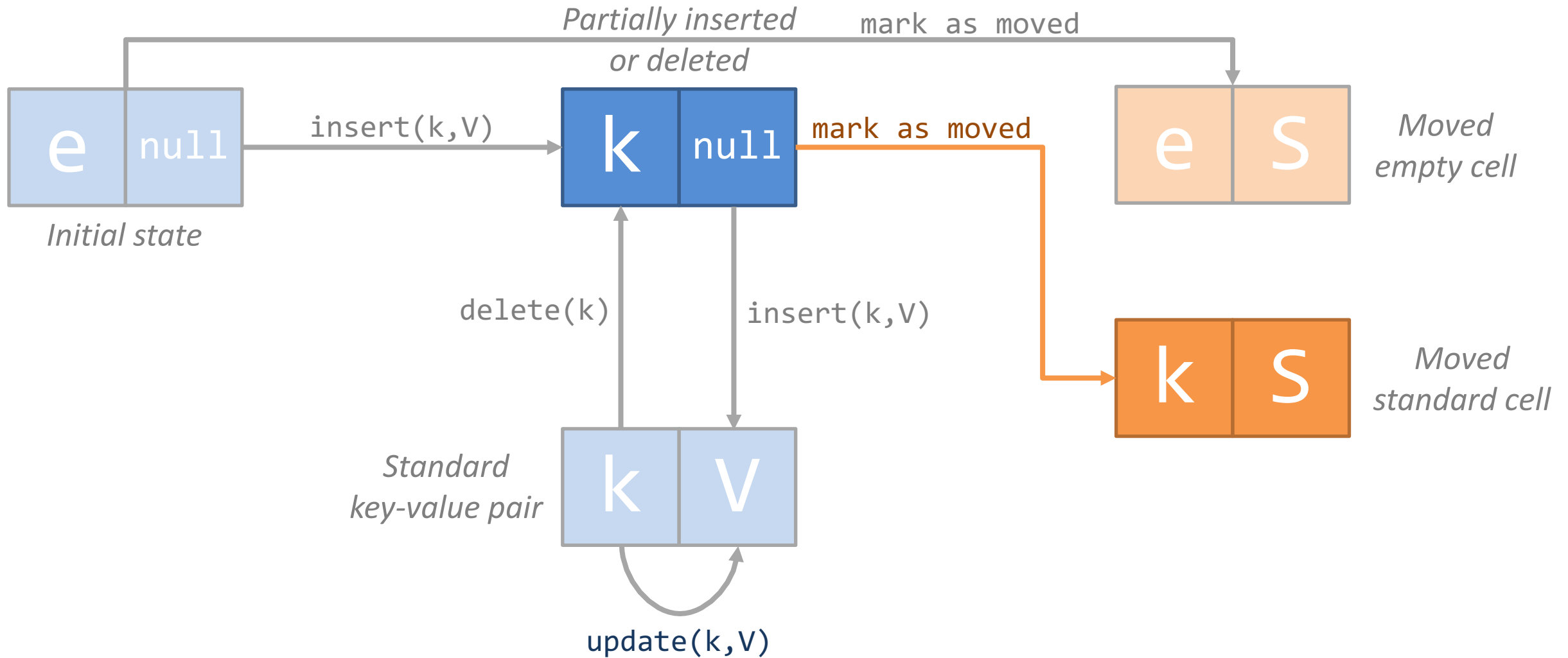
Жизнь ячейки: добавляем перенос



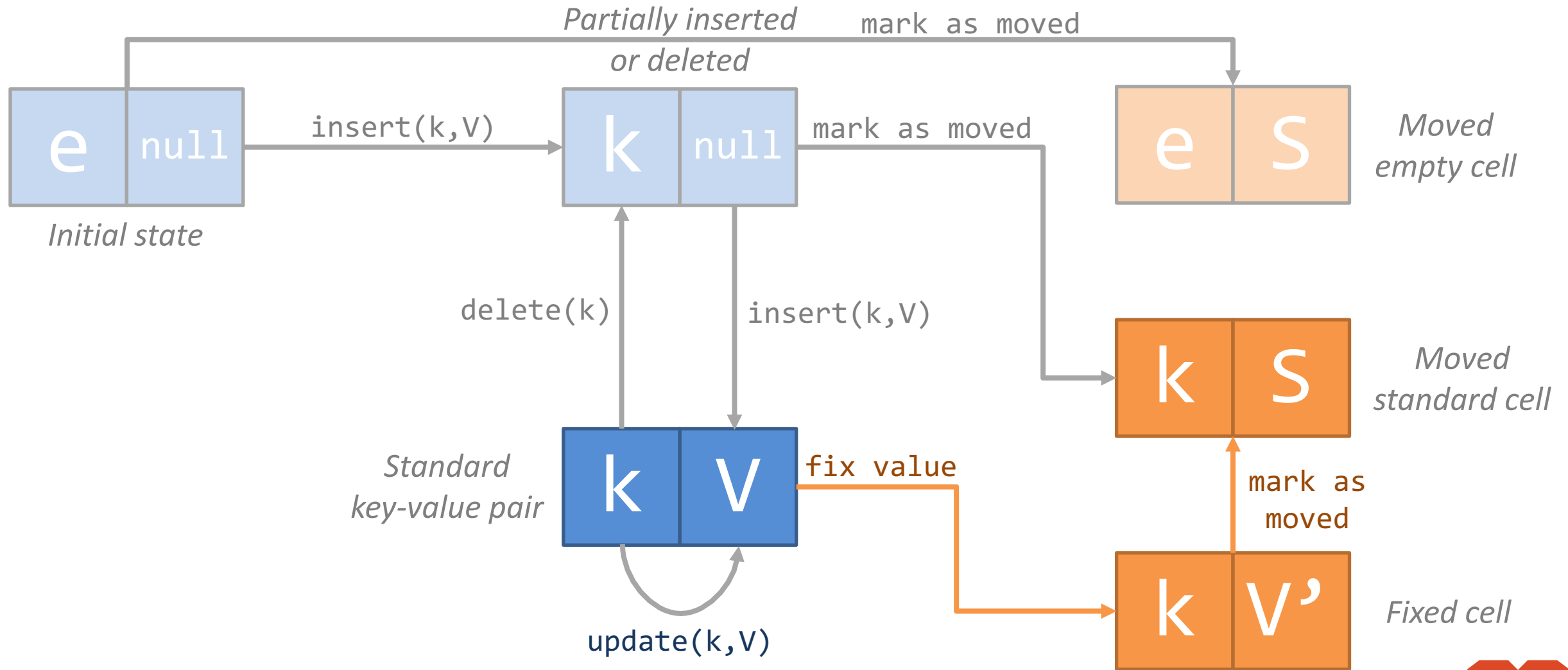
Жизнь ячейки: добавляем перенос



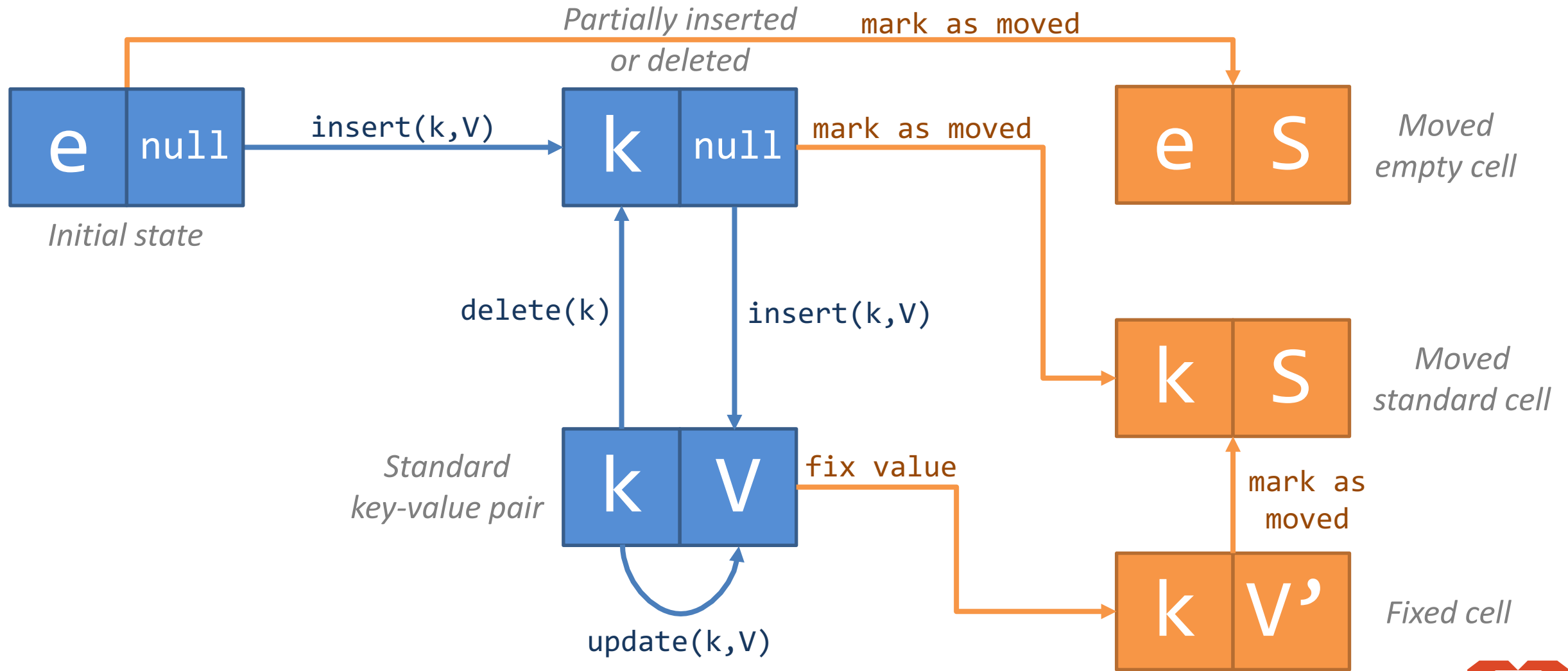
Жизнь ячейки: добавляем перенос



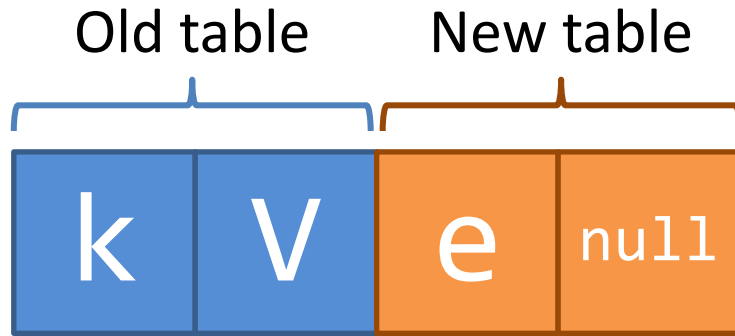
Жизнь ячейки: добавляем перенос



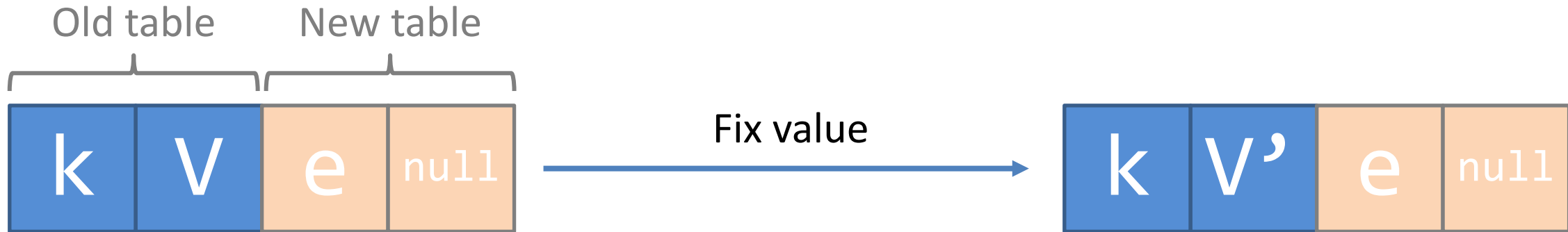
Жизнь ячейки: добавляем перенос



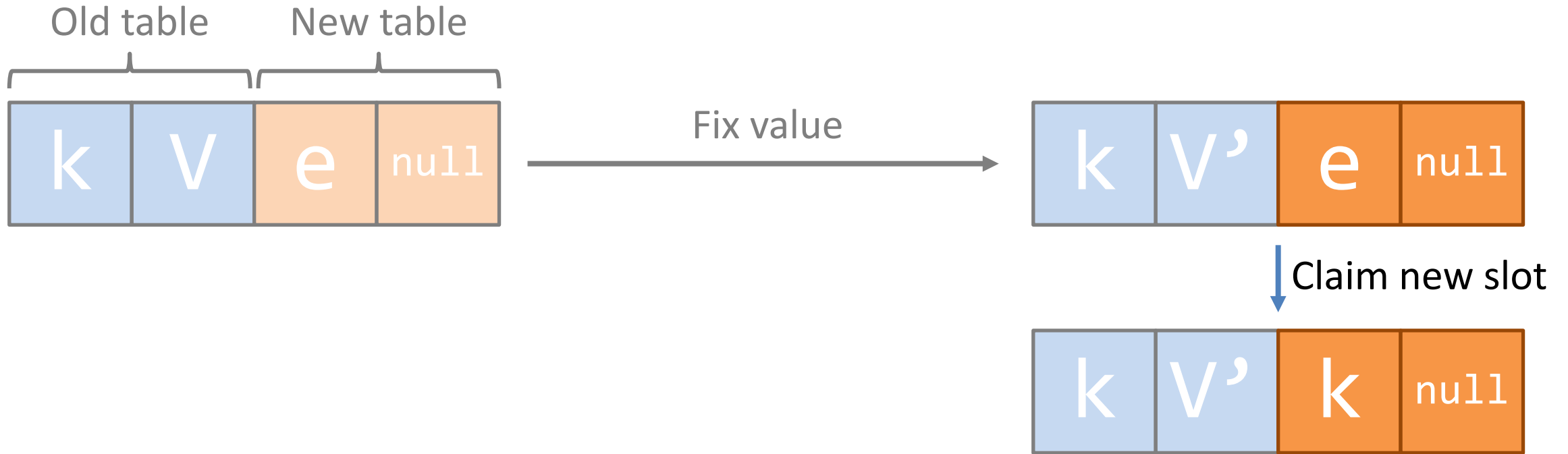
Перенос: с другого ракурса



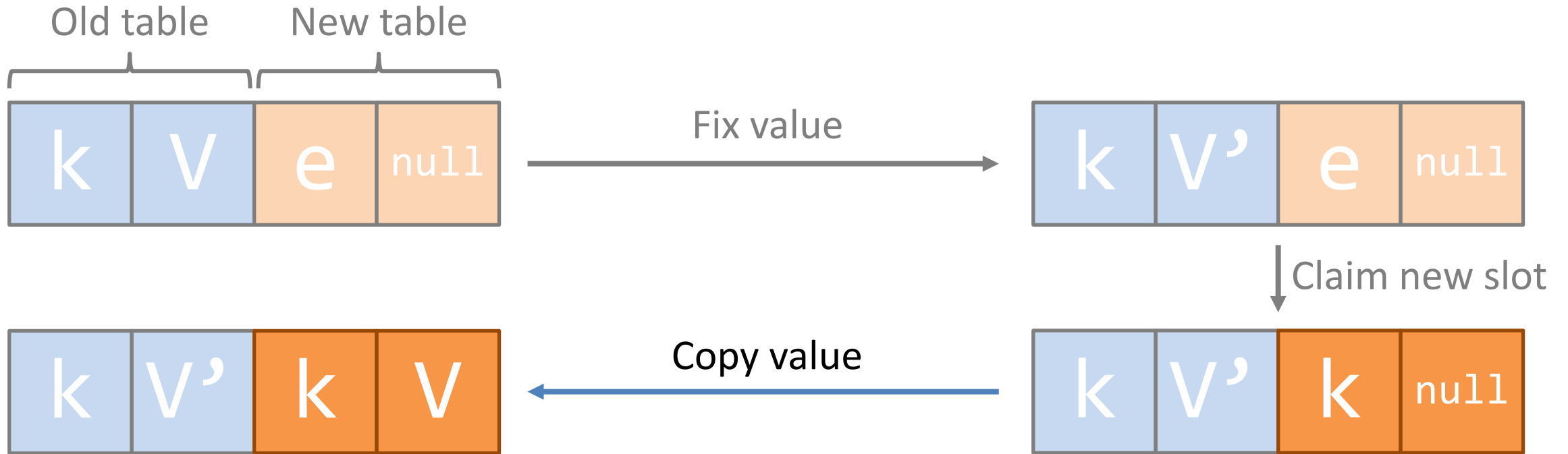
Перенос: с другого ракурса



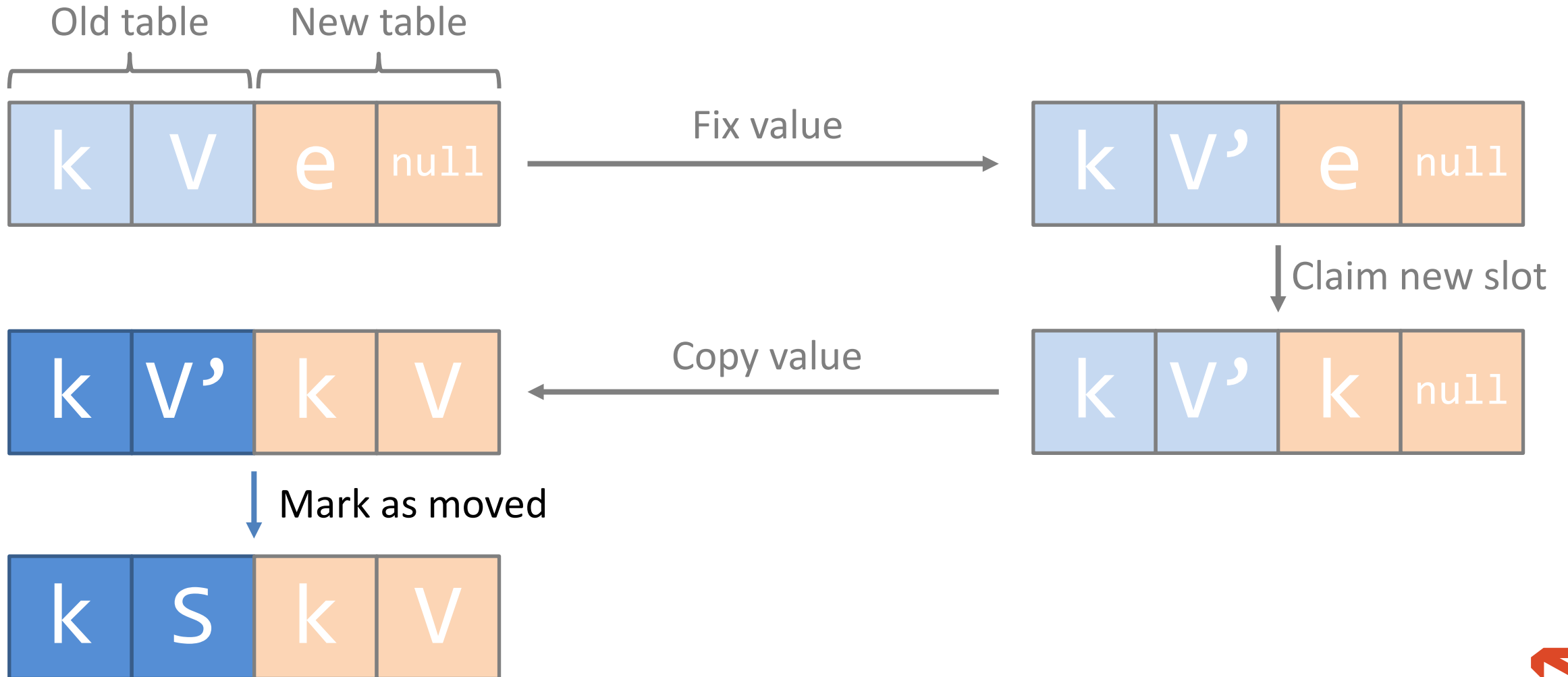
Перенос: с другого ракурса



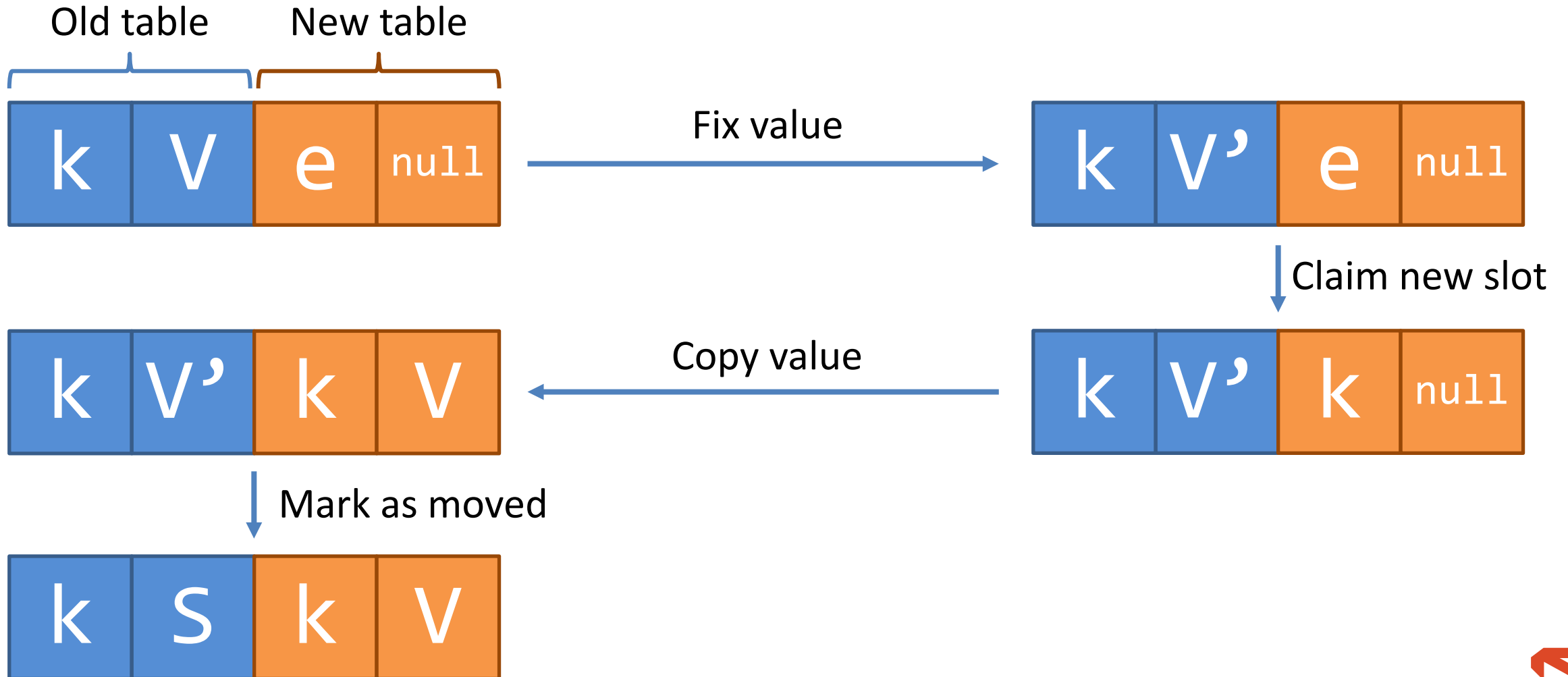
Перенос: с другого ракурса



Перенос: с другого ракурса



Перенос: с другого ракурса



Кооперация при переносе

- Как несколько потоков могут переносить элементы и не мешать друг другу?
 - Блоками по K элементов (так делает Cliff Click)
 - Переносом занимается выделенный поток

НУ И ЗАЧЕМ МНЕ ЭТО ВСЁ ЗНАТЬ???

Вспомним модель!

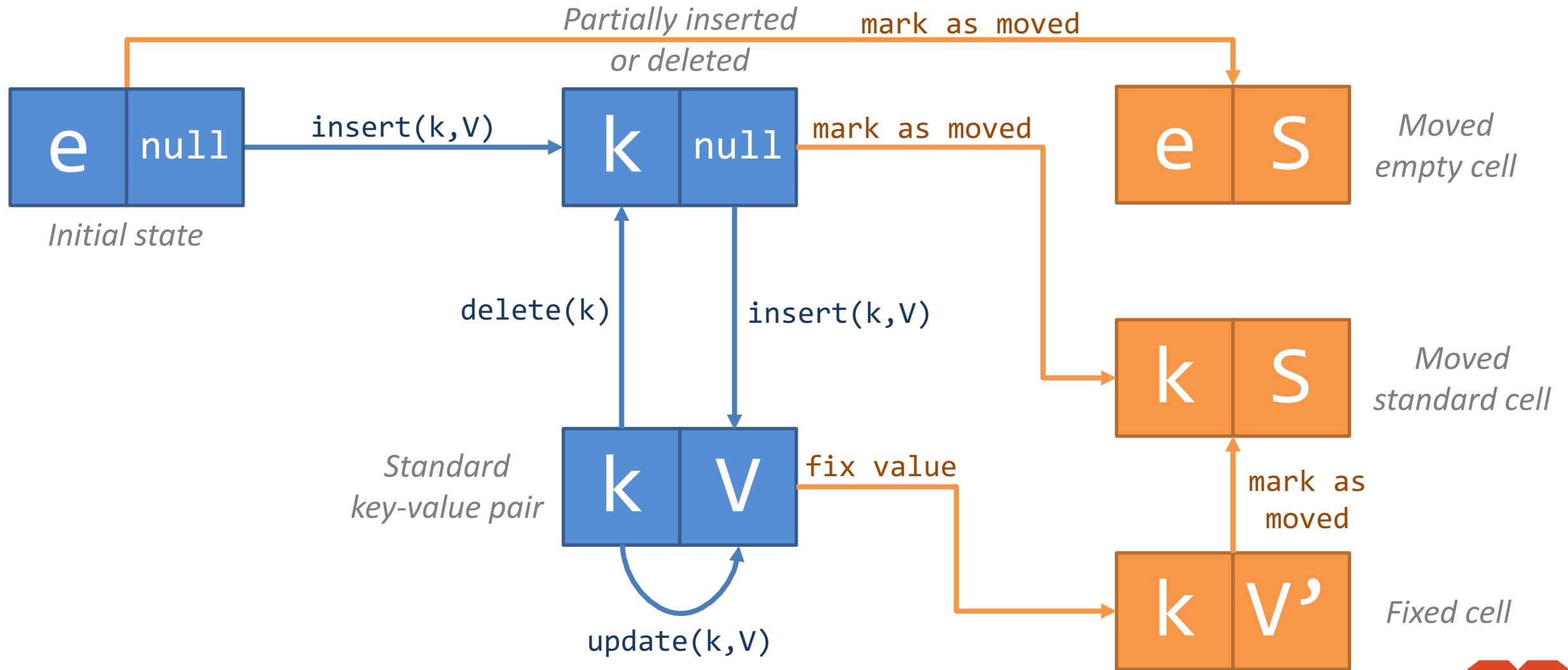
```
interface AccountsCache {  
    fun addAccount(id: Int, account: Account)  
    fun getById(id: Int): Account?  
}
```

Вспомним модель!

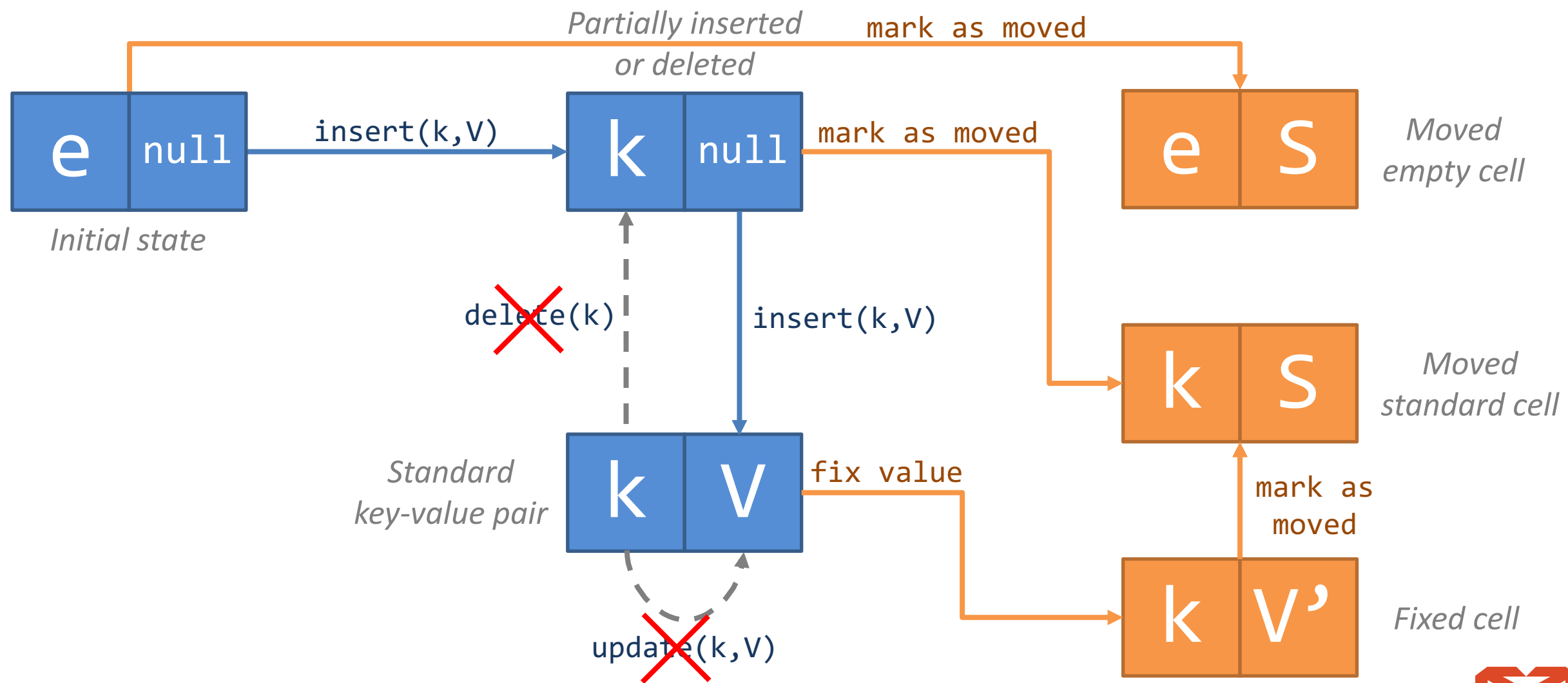
```
interface AccountsCache {  
    fun addAccount(id: Int, account: Account)  
    fun getById(id: Int): Account?  
}
```

- Account-ы никогда не меняются
- Удалений нет

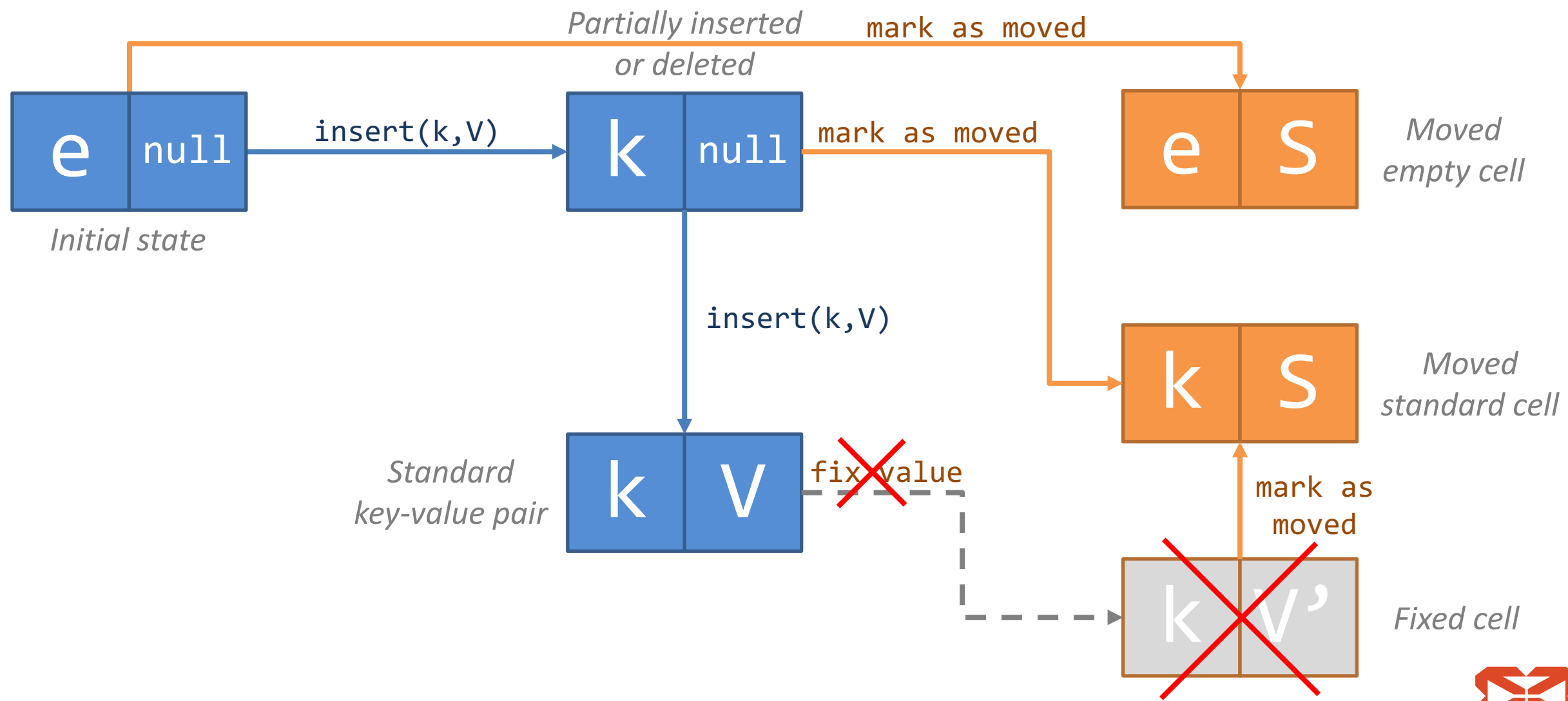
Упрощаем жизнь ячейки



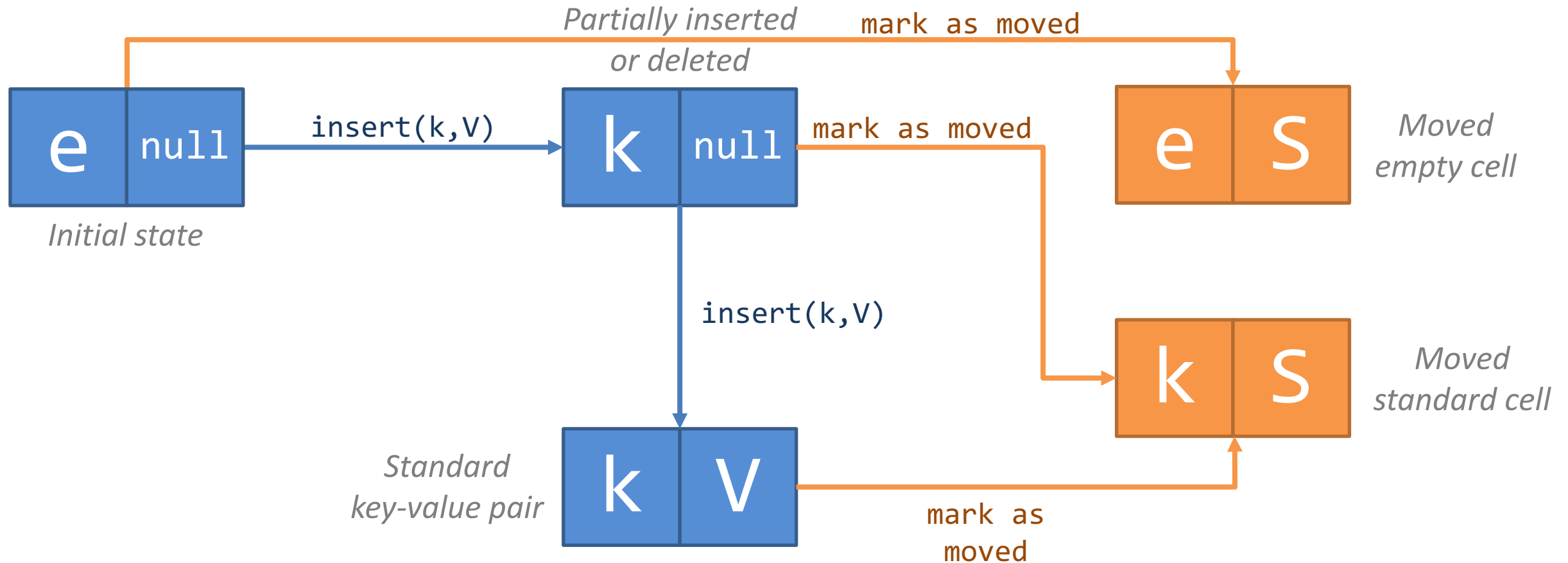
Упрощаем жизнь ячейки



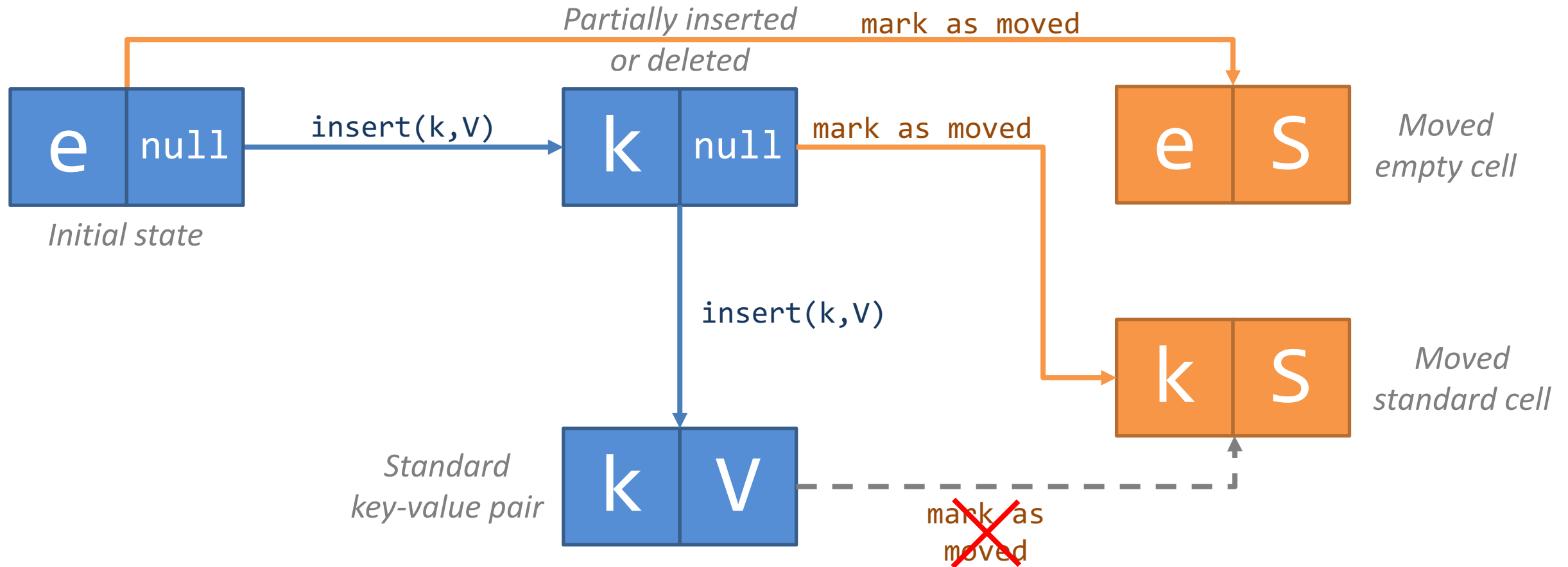
Упрощаем жизнь ячейки



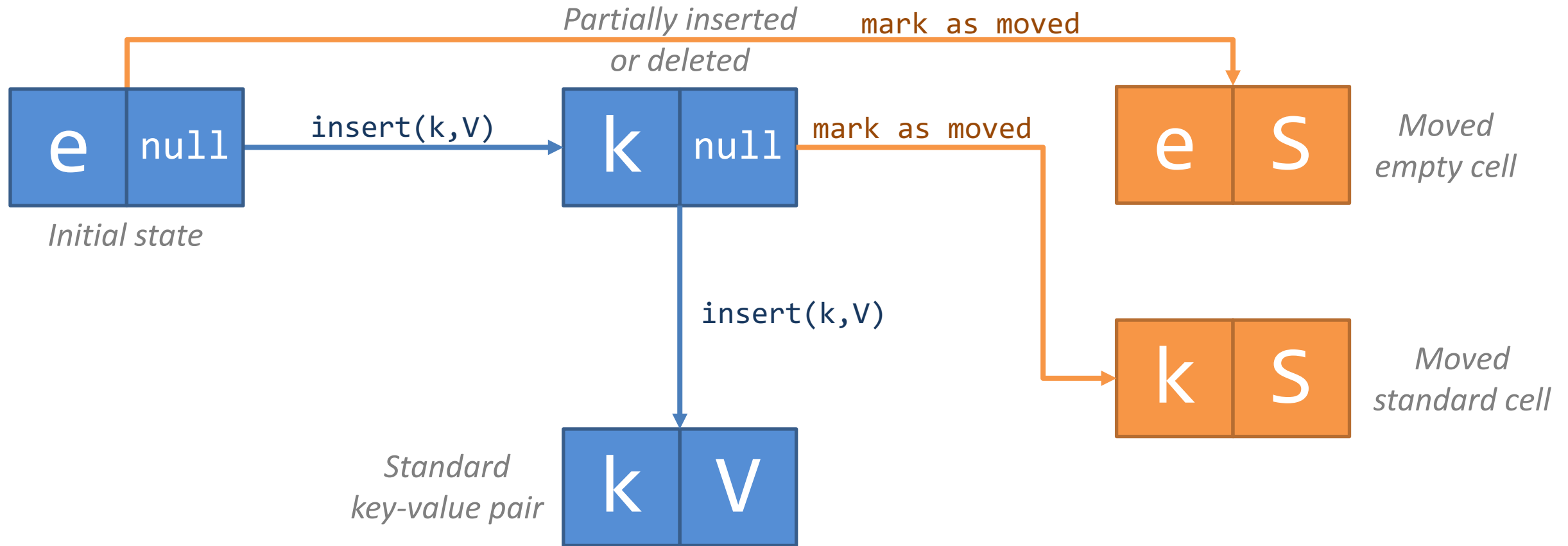
Упрощаем жизнь ячейки



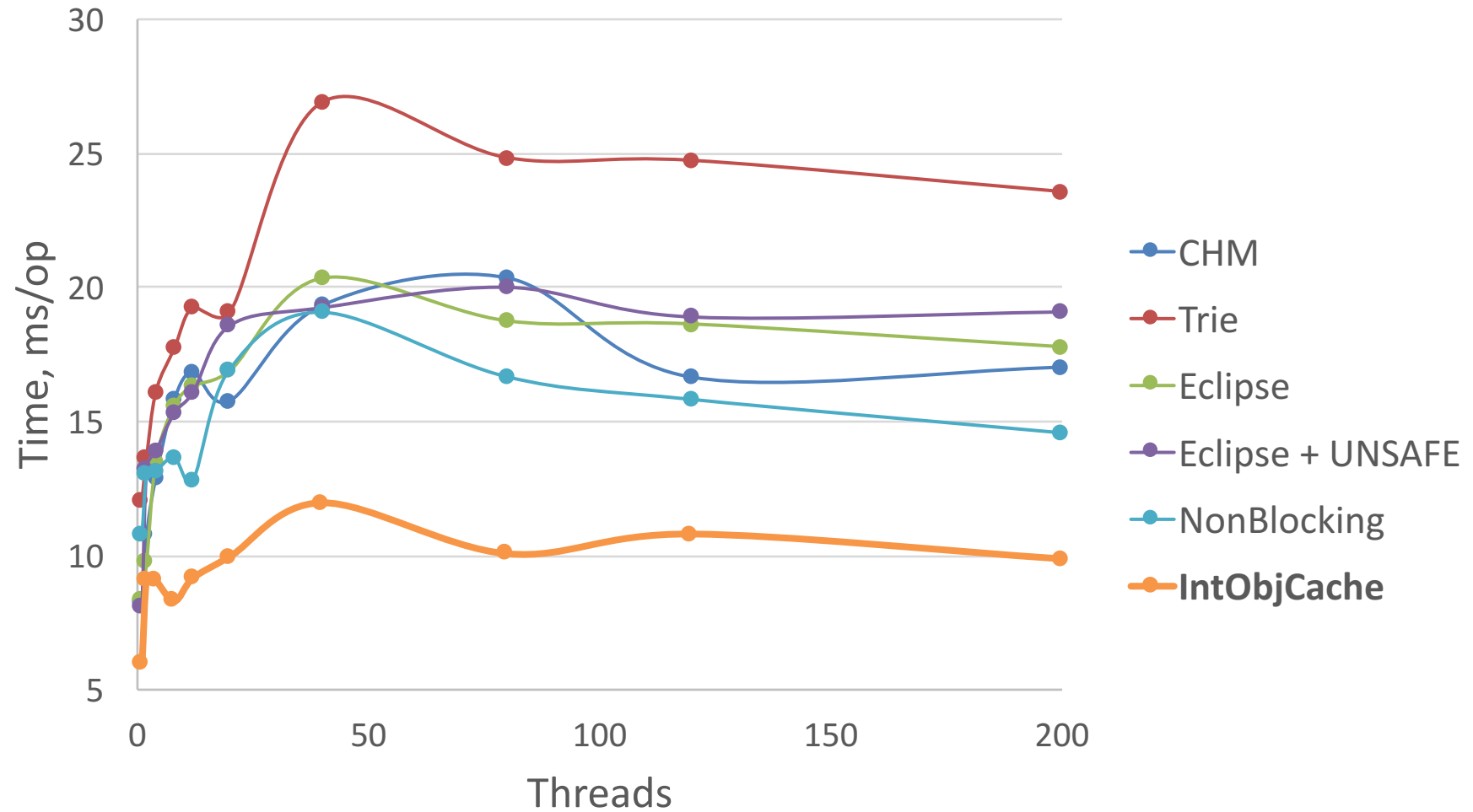
Упрощаем жизнь ячейки



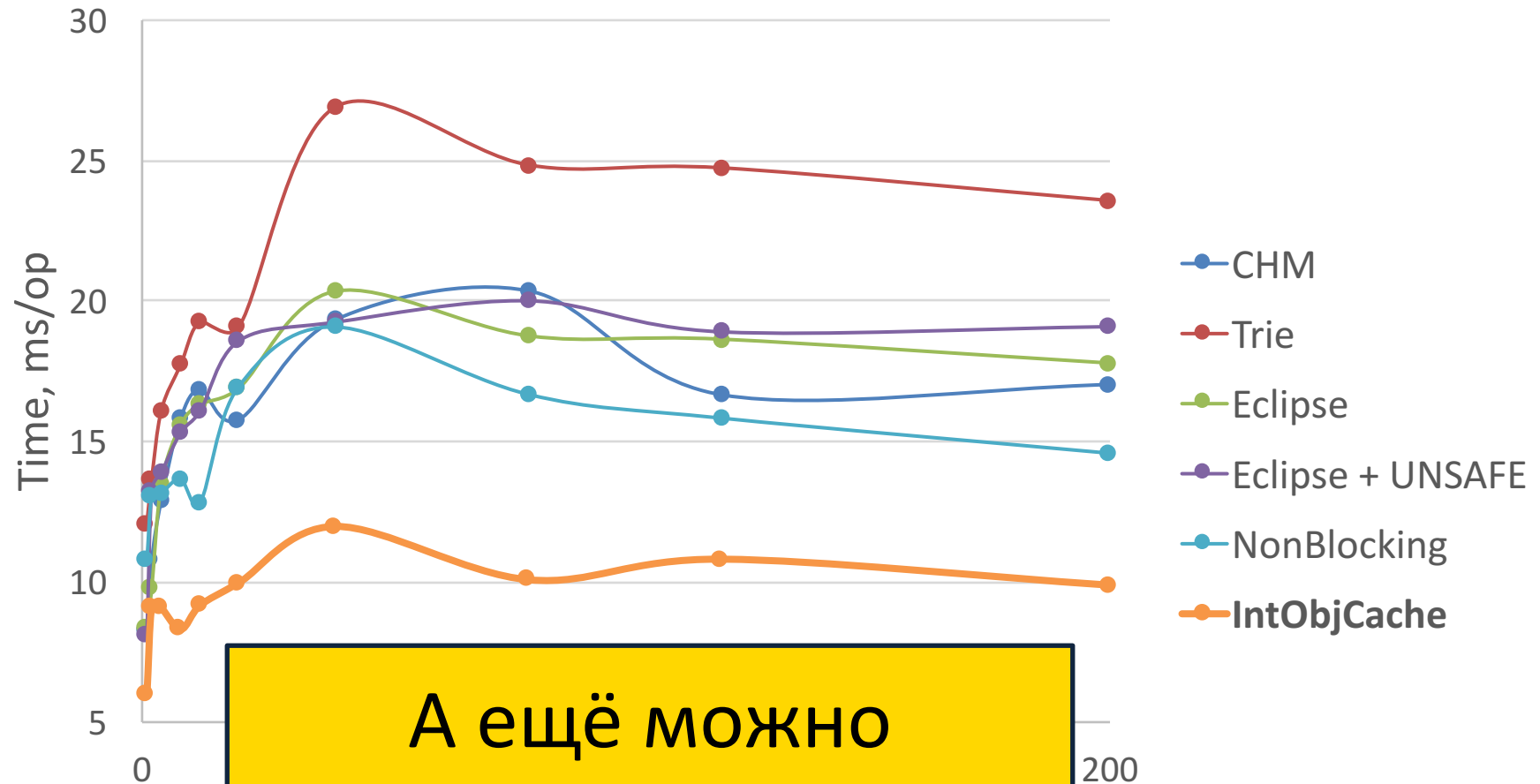
Упрощаем жизнь ячейки



Производительность



Производительность



А ещё можно
ускорить?

Посмотрим на модель ещё раз!

```
class Account(val id: Int) {  
    private val i1: Int = 0  
    private val i2: Int = 0  
    private val o1: Any? = null  
    private val o2: Any? = null  
    private val o3: Any? = null  
    private val o4: Any? = null  
}
```

Посмотрим на модель ещё раз!

```
class Account(val id: Int) {  
    private val i1: Int = 0  
    private val i2: Int = 0  
    private val o1: Any? = null  
    private val o2: Any? = null  
    private val o3: Any? = null  
    private val o4: Any? = null  
}
```

Уже храним id

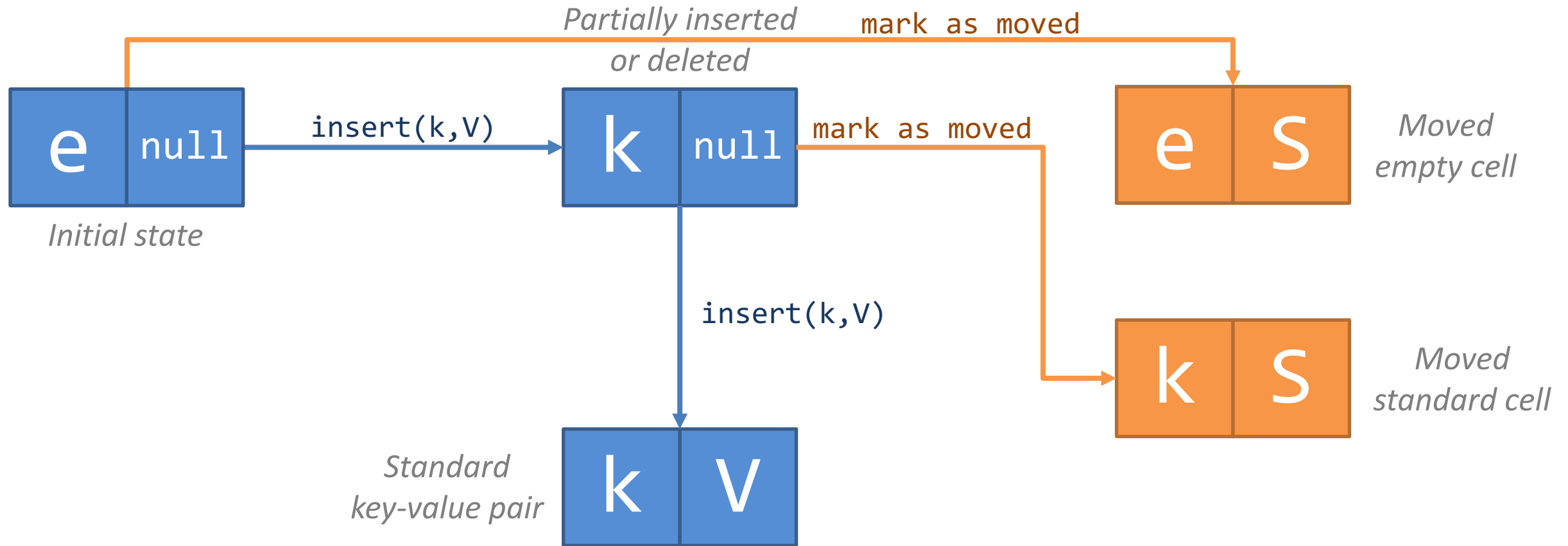
Посмотрим на модель ещё раз!

```
class Account(val id: Int) {  
    private val i1: Int = 0  
    private val i2: Int = 0  
    private val o1: Any? = null  
    private val o2: Any? = null  
    private val o3: Any? = null  
    private val o4: Any? = null  
}
```

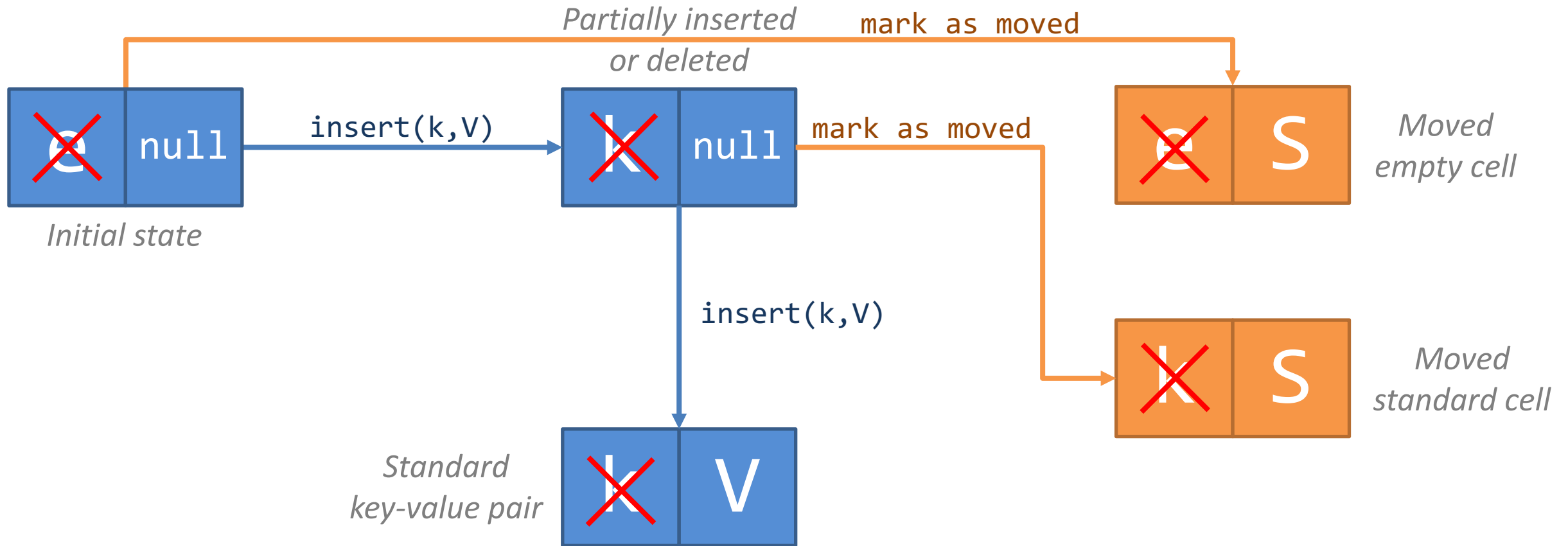
Уже храним id

Может не хранить
ключ отдельно?

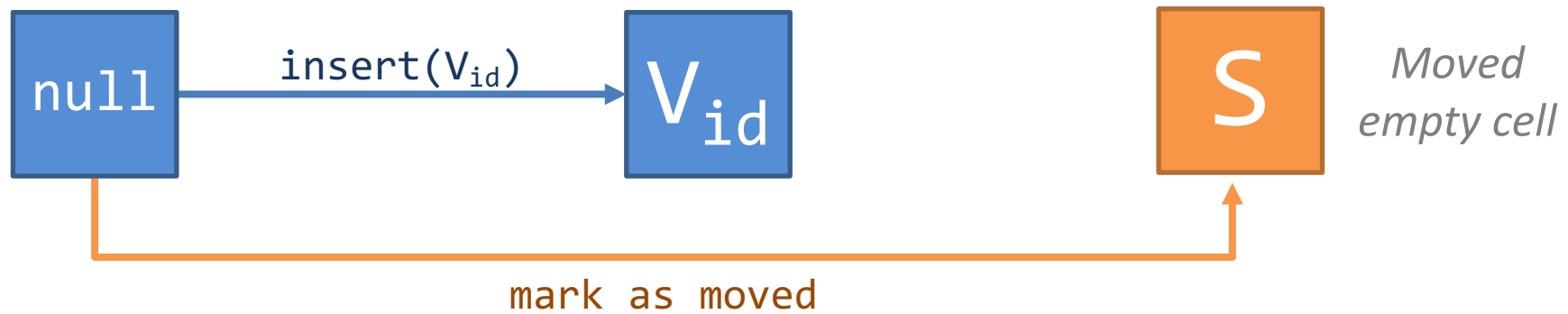
Упрощаем ячейку



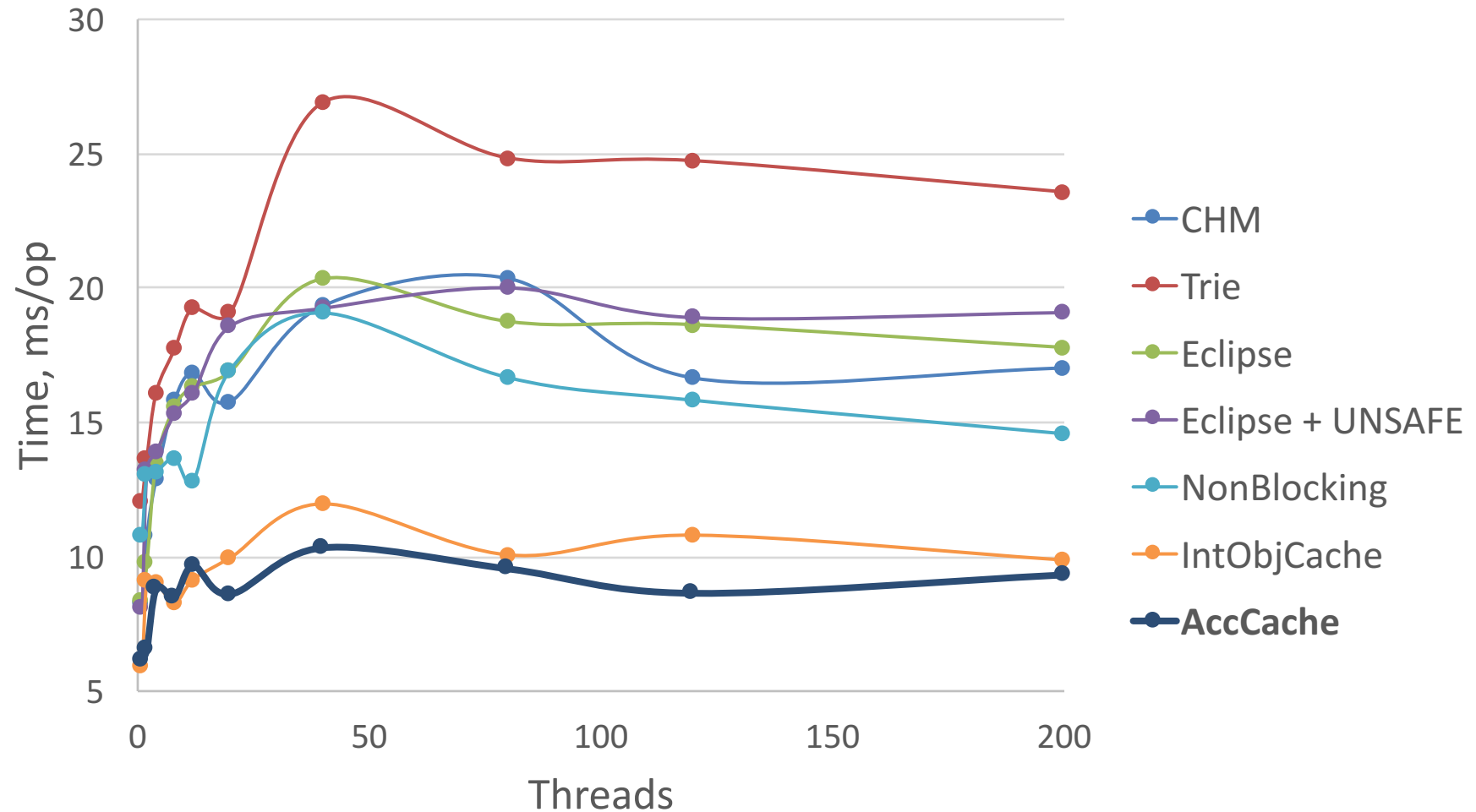
Упрощаем ячейку



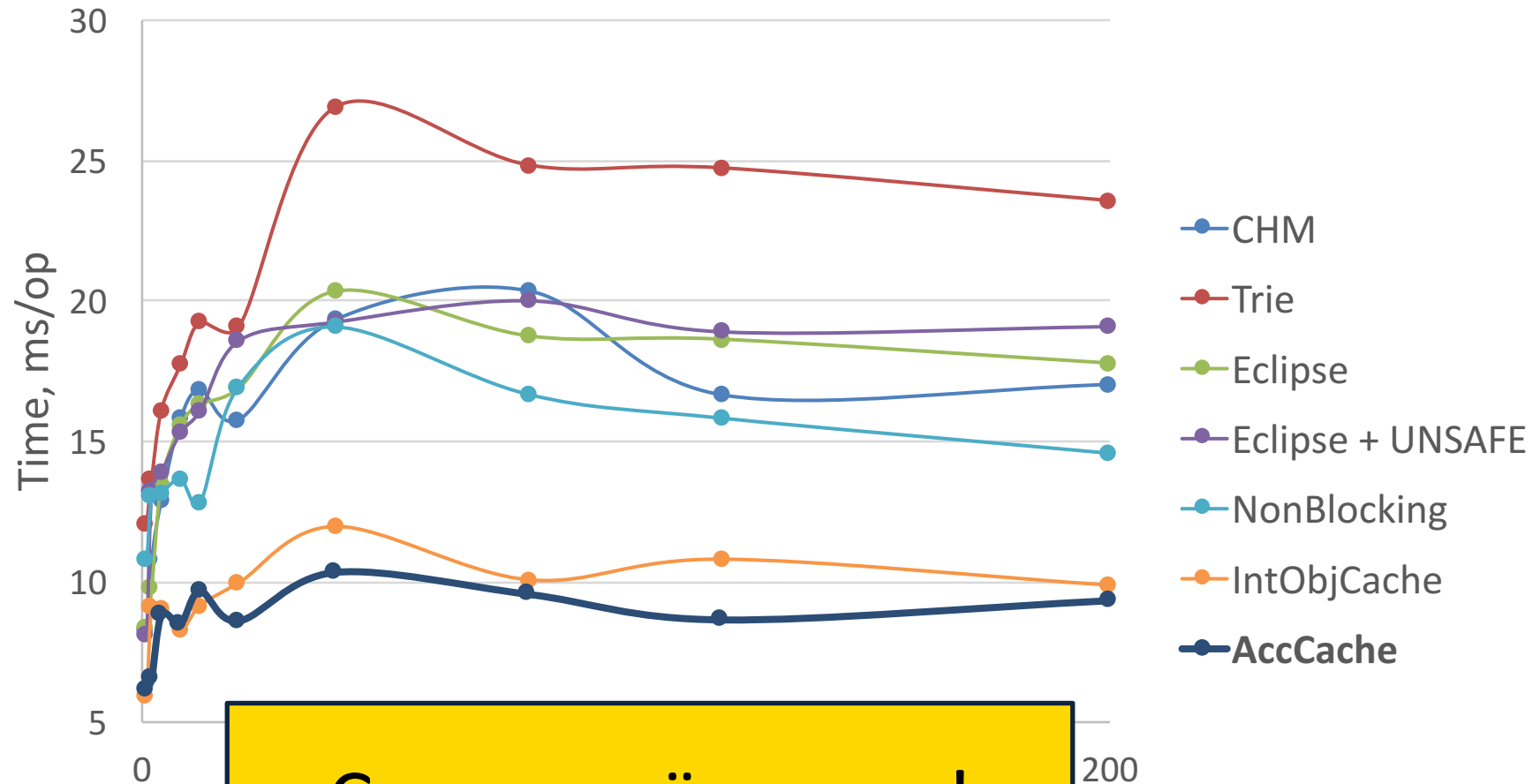
Упрощаем ячейку



Производительность



Производительность



Стало ещё лучше!

Можно ли ещё улучшить?

- IntObjCache и AccCache «плохо» делают rehash

Можно ли ещё улучшить?

- IntObjCache и AccCache «плохо» делают rehash
- VarHandle позволяет избежать лишних чтений и синхронизаций
 - compareAndExchange помогает не делать лишних чтений после неудачного CAS-а

Можно ли ещё улучшить?

- IntObjCache и AccCache «плохо» делают rehash
- VarHandle позволяет избежать лишних чтений и синхронизаций
 - compareAndExchange помогает не делать лишних чтений после неудачного CAS-а
 - Что-то вроде release-acquire semantics, но не специфицировано в JMM (есть «спека» от Doug Lea)
 - Высокий порог вхождения

Итоги

- Как правило достаточно стандартных реализаций СД
 - Если недостаточно стандартных – поищите альтернативы

Итоги

- Как правило достаточно стандартных реализаций СД
 - Если недостаточно стандартных – поищите альтернативы
- При разработке своих алгоритмов важно
 - Понимать существующие подходы
 - Иметь простую и понятную модель
 - Использовать вводимые моделью ограничения

Спасибо за внимание!

Никита Коваль

ndkoval *at* ya *dot* ru
twitter.com/nkoval_