

Методы параллельной обработки сверхбольших баз данных с использованием распределенных колоночных индексов

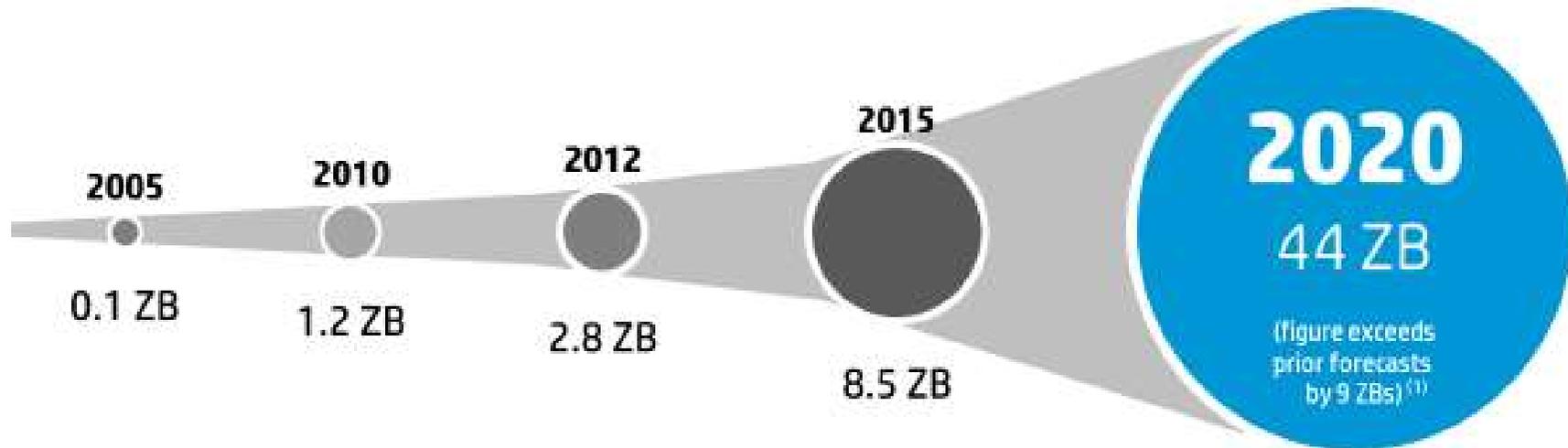
05.13.11 - математическое и программное обеспечение
вычислительных машин, комплексов и компьютерных сетей

Диссертация на соискание ученой степени кандидата физико-математических наук

Елена Владимировна Иванова

Научный руководитель:
СОКОЛИНСКИЙ Леонид Борисович,
доктор физ.-мат. наук, профессор

Проблема больших данных



- Объем хранимой информации удваивается каждые два года
- из всего объема существующих данных потенциально полезны 22%, из которых менее 5% были подвергнуты анализу

Тенденции в развитии технологий обработки больших данных

- Для обработки больших данных необходимо использовать СУБД (Майкл Стоунбрейкер)
- Кластерные вычислительные системы с большой суммарной оперативной памятью
- Базы данных в оперативной памяти
- Многоядерные ускорители
- Колоночное представление данных со сжатием

Цель диссертационной работы

Разработка и исследование эффективных методов параллельной обработки сверхбольших баз данных с использованием колоночного представления информации, ориентированных на кластерные вычислительные системы, оснащенные многоядерными ускорителями, и допускающих интеграцию с реляционными СУБД

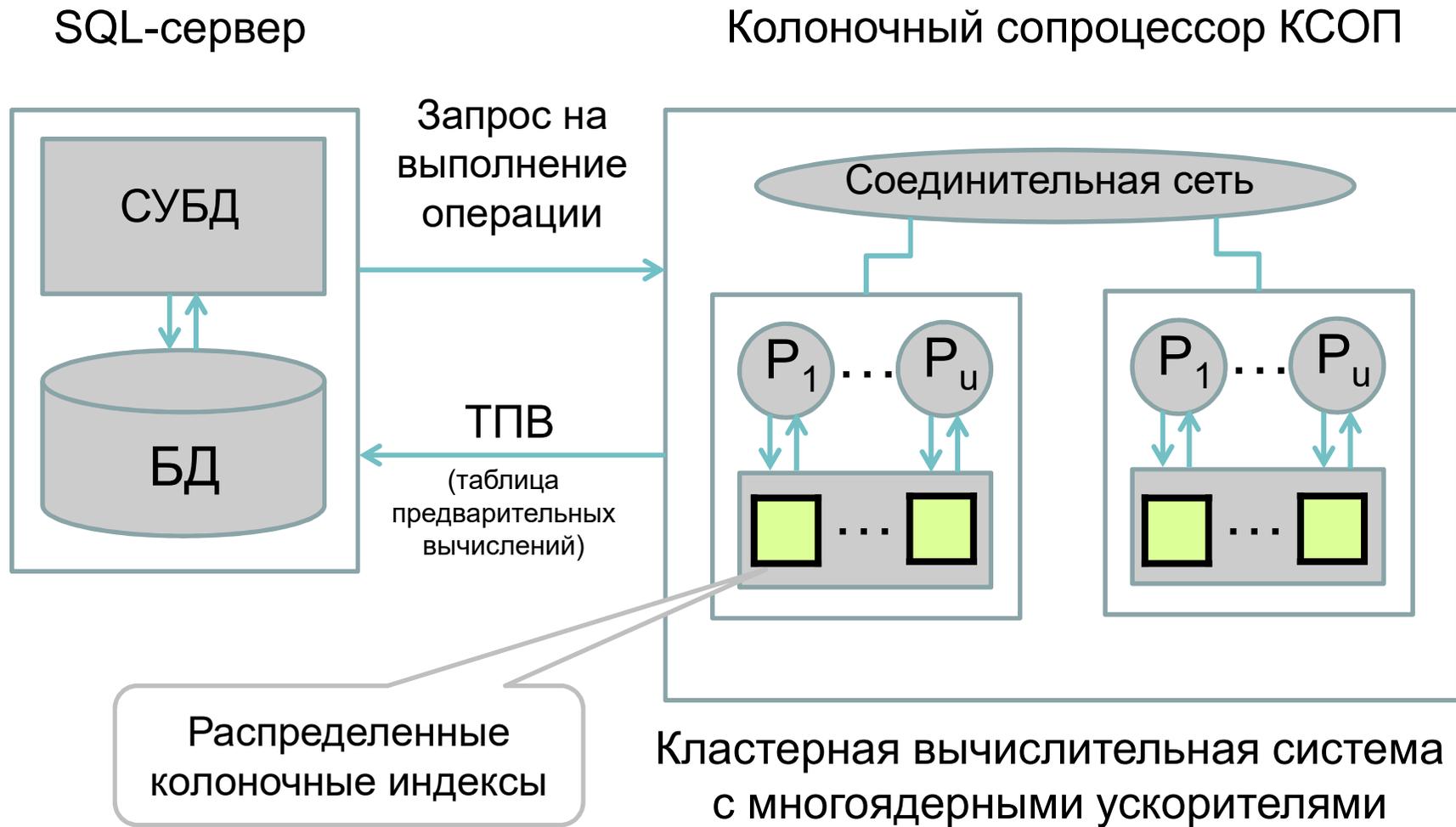
Основные задачи

1. Разработать колоночные индексы и методы их фрагментации
2. Разработать методы декомпозиции реляционных операций для распределенных колоночных индексов
3. Реализовать предложенные подходы и методы в виде колоночного сопроцессора КСОП для кластерных вычислительных систем
4. Провести вычислительные эксперименты с использованием КСОП

Работы по теме диссертации

| | | |
|---|---|--|
| 1 | Ramamurthy R., Dewitt D., Su Q. A case for fractured mirrors // Proceedings of the VLDB Endowment. 2002. Vol. 12, No. 2. P. 89-101. | «Разбитое зеркало» |
| 2 | Abadi D.J., Madden S.R., Hachem N. Column-Stores vs. Row-Stores: How Different Are They Really? // Proceedings of the 2008 ACM SIGMOD international conference on Management of data, June 9-12, 2008, Vancouver, BC, Canada. ACM, 2008. P. 967-980. | Эмуляция колоночного представления в строчной СУБД |
| 3 | Bruno N. Teaching an Old Elephant New Tricks // Online Proceedings of Fourth Biennial Conference on Innovative Data Systems Research (CIDR 2009), Asilomar, CA, USA, January 4-7, 2009. | C-таблицы |
| 4 | El-Helw A., Ross K.A., Bhattacharjee B., Lang C.A., Mihaila G.A. Column-oriented query processing for row stores // Proceedings of the ACM 14th international workshop on Data Warehousing and OLAP (DOLAP '11), October 28, 2011, Glasgow, United Kingdom. ACM, 2011. P. 67-74. | Только индексные планы |
| 5 | Larson P.-A., Clinciu C., Hanson E. N., Oks A., Price S. L., Rangarajan S., Surna A., Zhou Q. SQL server column store indexes // Proceedings of the 2011 ACM SIGMOD International Conference on Management of data (SIGMOD '11), June 12-16, 2011, Athens, Greece. ACM, 2011. P. 1177-1184. | Индексы колоночной памяти (column store indexes) |
| 6 | Jha S., He B., Lu M., Cheng X., Huynh H. P. Improving main memory hash joins on Intel Xeon Phi processors: an experimental approach // Proceedings of the VLDB Endowment. 2015. Vol. 8, No. 6. P. 642-653. | Использование Xeon Phi |

Предлагаемое решение



Колоночный индекс

$I_{R,B}$
(колоночный индекс
для B)

| A | B |
|-----|-----|
| 3 | 110 |
| 1 | 114 |
| 5 | 127 |
| 0 | 136 |
| 2 | 136 |
| 7 | 158 |
| 4 | 174 |
| 6 | 174 |
| 8 | 187 |

Таблица R

| A | B | C |
|-----|-----|-----|
| 0 | 136 | 17 |
| 1 | 114 | 10 |
| 2 | 136 | 25 |
| 3 | 110 | 10 |
| 4 | 174 | 15 |
| 5 | 127 | 99 |
| 6 | 174 | 97 |
| 7 | 158 | 63 |
| 8 | 187 | 55 |

$I_{R,C}$
(колоночный индекс
для C)

| A | C |
|-----|-----|
| 1 | 10 |
| 3 | 10 |
| 4 | 15 |
| 0 | 17 |
| 2 | 25 |
| 8 | 55 |
| 7 | 63 |
| 6 | 97 |
| 5 | 99 |

Формальное определение колоночного индекса

Пусть $R(A, B, \dots)$ – отношение R с суррогатным ключом A и атрибутом B . A состоит из целочисленных неотрицательных элементов.

\mathfrak{D}_B – домен атрибута B . На множестве \mathfrak{D}_B задано отношение линейного порядка. $T(R) = n$ – количество элементов в R .

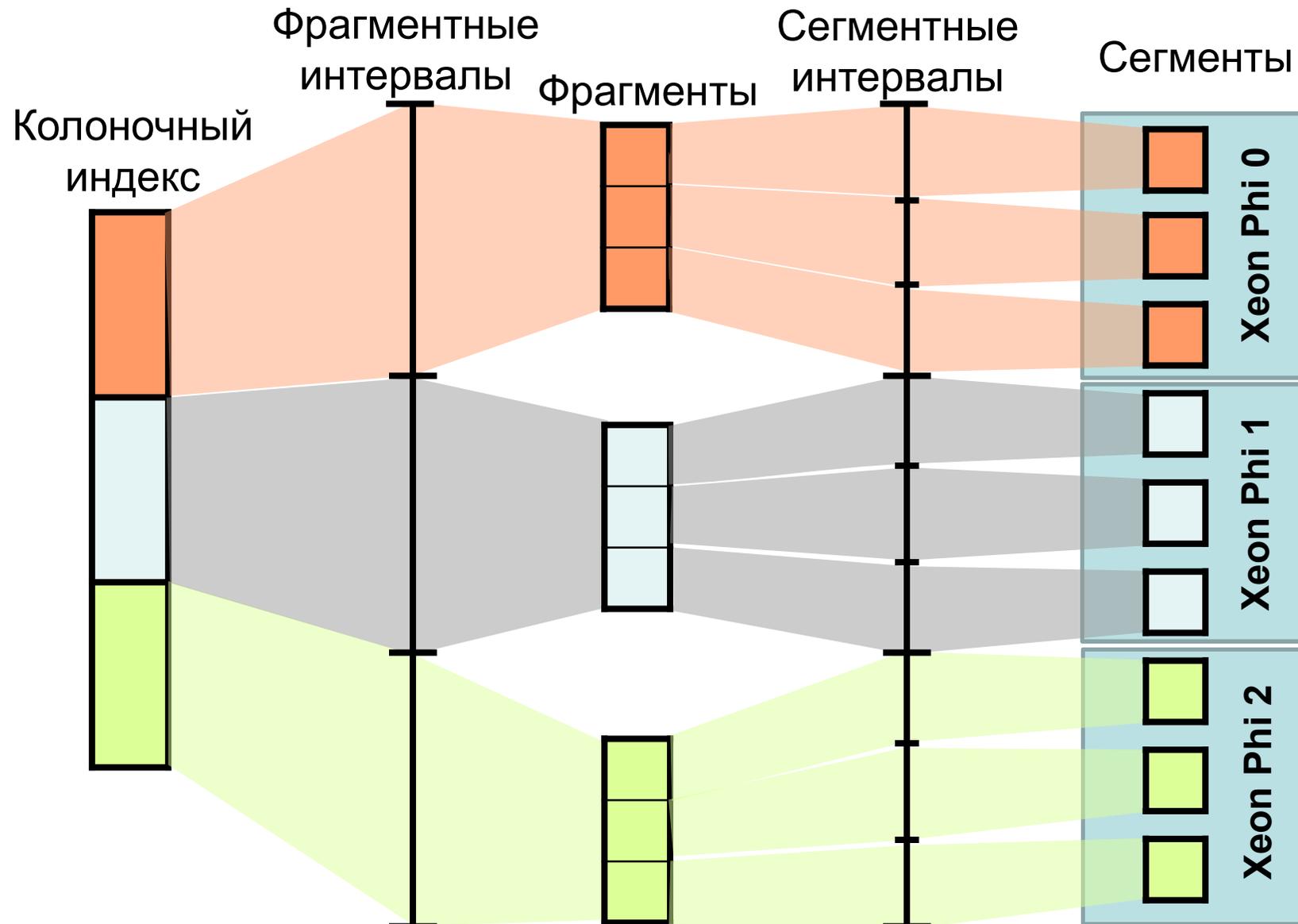
Колоночным индексом $I_{R.B}$ атрибута B отношения R называется упорядоченное отношение, удовлетворяющее следующим требованиям:

$$T(I_{R.B}) = n \quad \text{и} \quad \pi_A(I_{R.B}) = \pi_A(R); \quad (1)$$

$$\forall x_1, x_2 \in I_{R.B} (x_1 \leq x_2 \Leftrightarrow x_1.B \leq x_2.B); \quad (2)$$

$$\forall r \in R (\forall x \in I_{R.B} (r.A = x.A \Rightarrow r.B = x.B)). \quad (3)$$

Распределение данных



Формальное определение доменно-интервальной фрагментации

Разбиение домена \mathfrak{D}_B на k непересекающихся интервалов:

$$V_0 = [v_0; v_1); V_1 = [v_1; v_2); \dots; V_{k-1} = [v_{k-1}; v_k);$$
$$v_0 < v_1 < \dots < v_k;$$
$$\mathfrak{D}_B = \bigcup_{i=0}^{k-1} V_i$$

Доменная функция фрагментации

$$\varphi_{\mathfrak{D}_B}: \mathfrak{D}_B \rightarrow \{0, \dots, k-1\}$$

Доменно-интервальная функция фрагментации колоночного индекса: $\forall i \in \{0, \dots, k-1\} \left(\forall b \in \mathfrak{D}_B (\varphi_{\mathfrak{D}_B}(b) = i \Leftrightarrow b \in V_i) \right)$

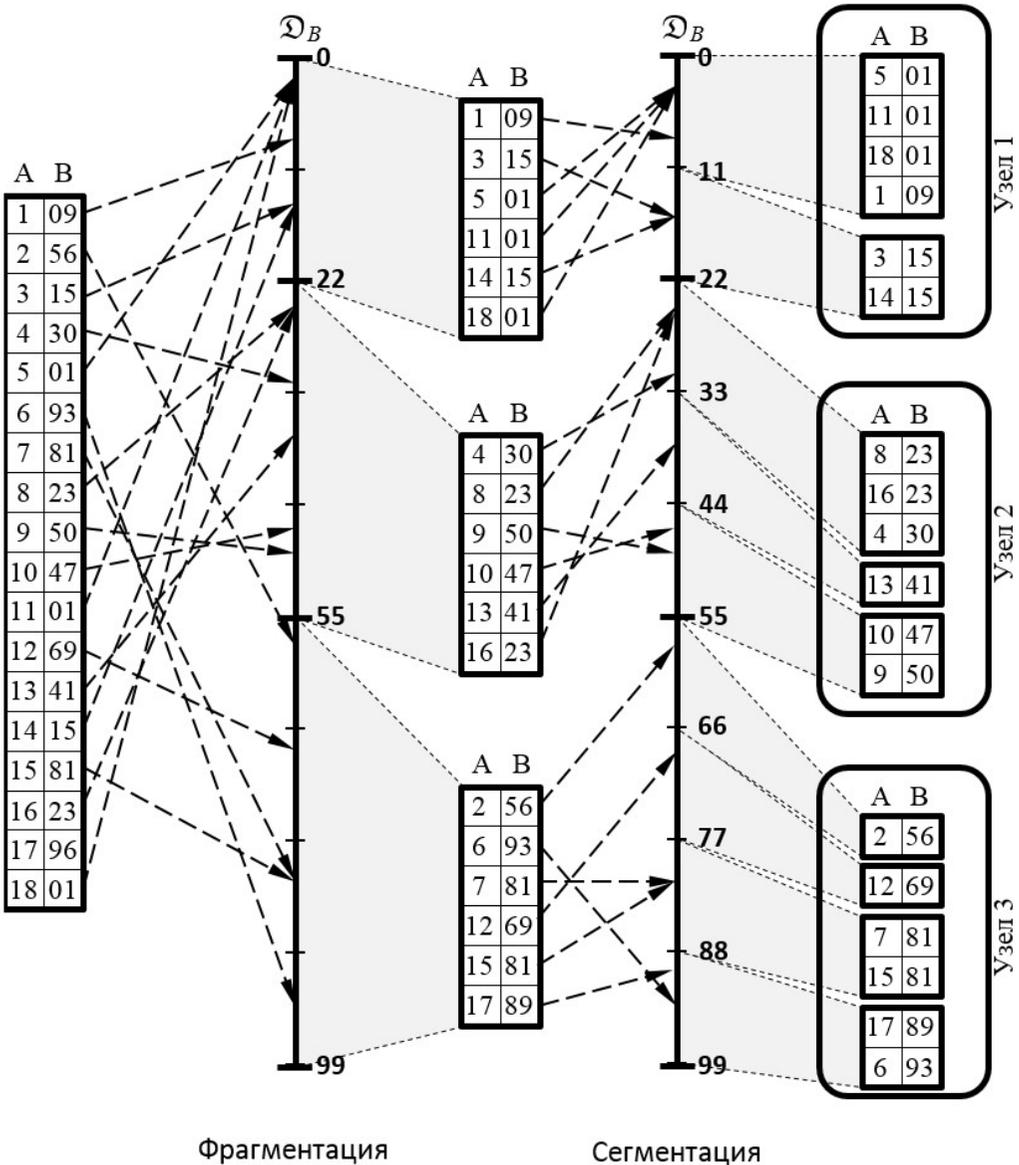
$$\varphi_{I_{R.B}}: I_{R.B} \rightarrow \{0, \dots, k-1\}$$

$$\forall x \in I_{R.B} \left(\varphi_{I_{R.B}}(x) = \varphi_{\mathfrak{D}_B}(x.B) \right)$$

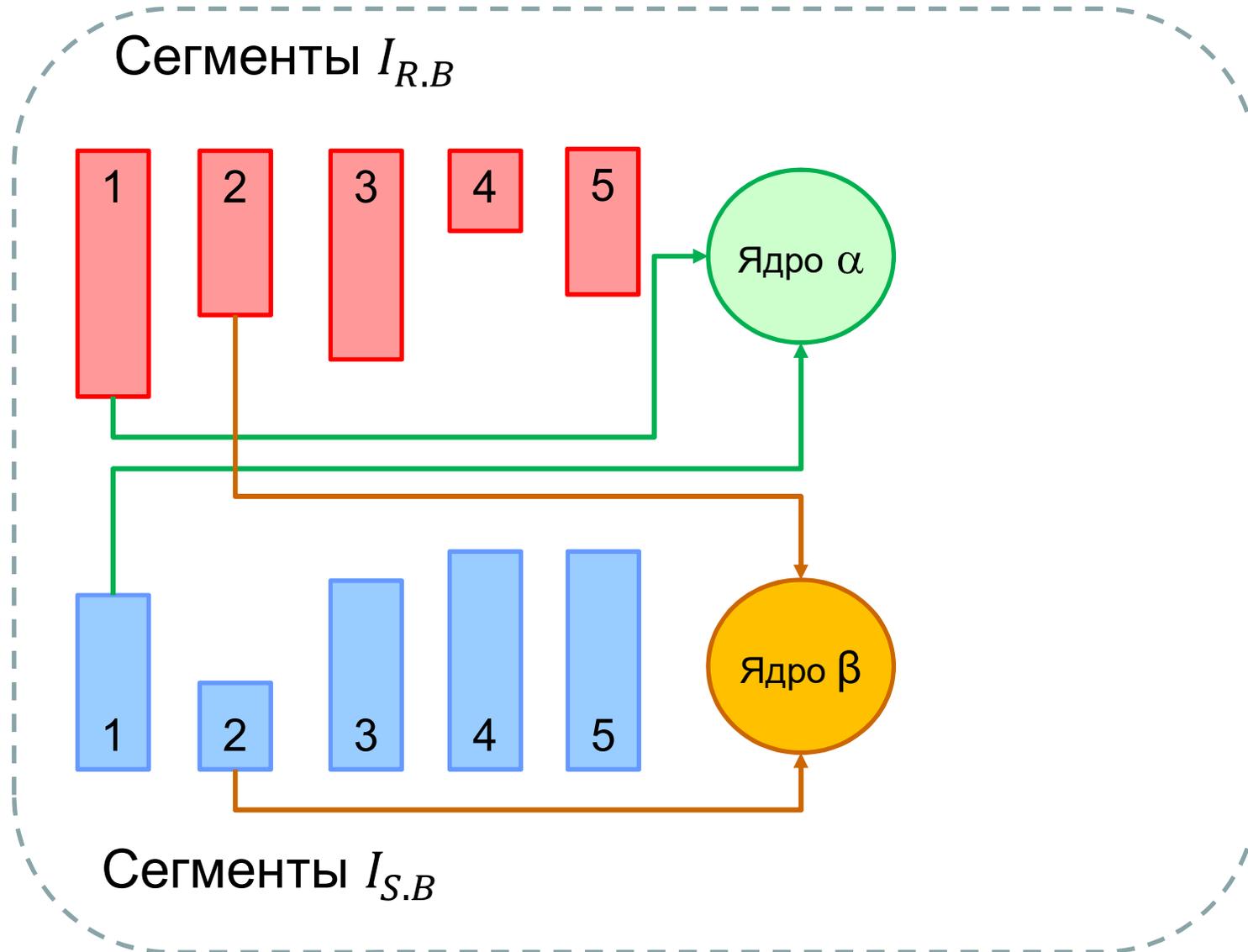
i -тый фрагмент колоночного индекса:

$$I_{R.B}^i = \left\{ x \mid x \in I_{R.B}; \varphi_{I_{R.B}}(x) = i \right\}$$

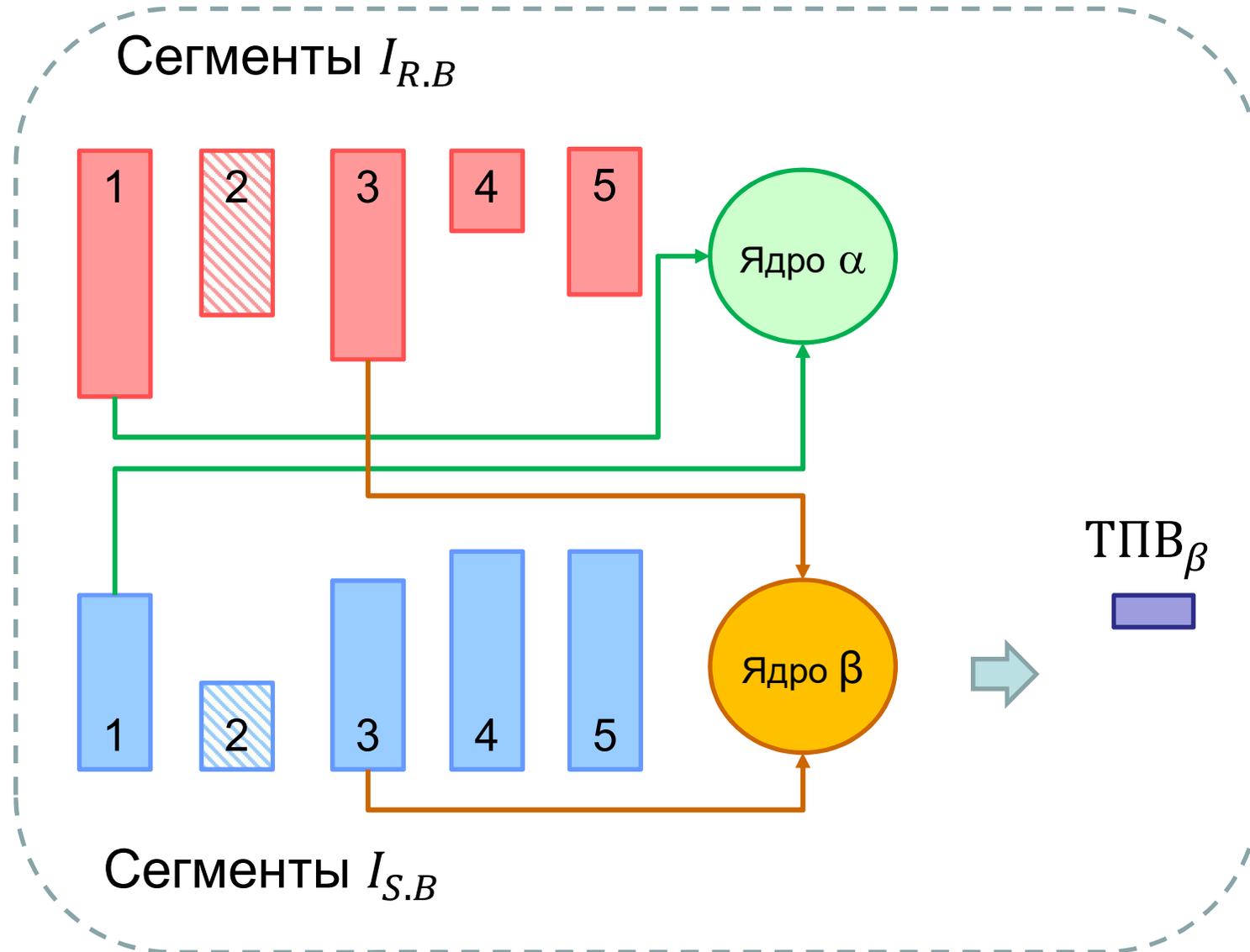
Пример двухуровневого разбиения колоночного индекса на фрагменты и сегменты



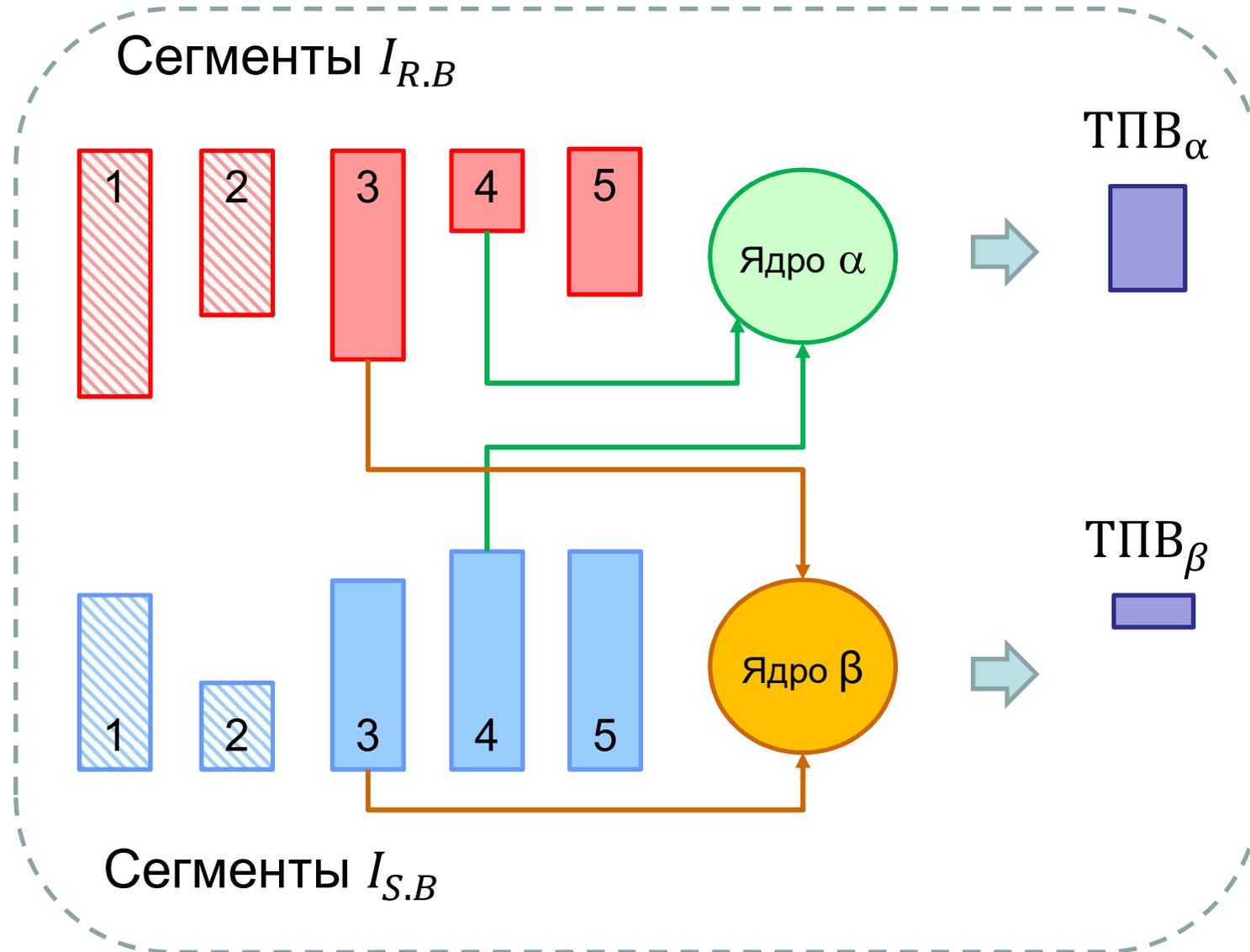
Балансировка загрузки ядер



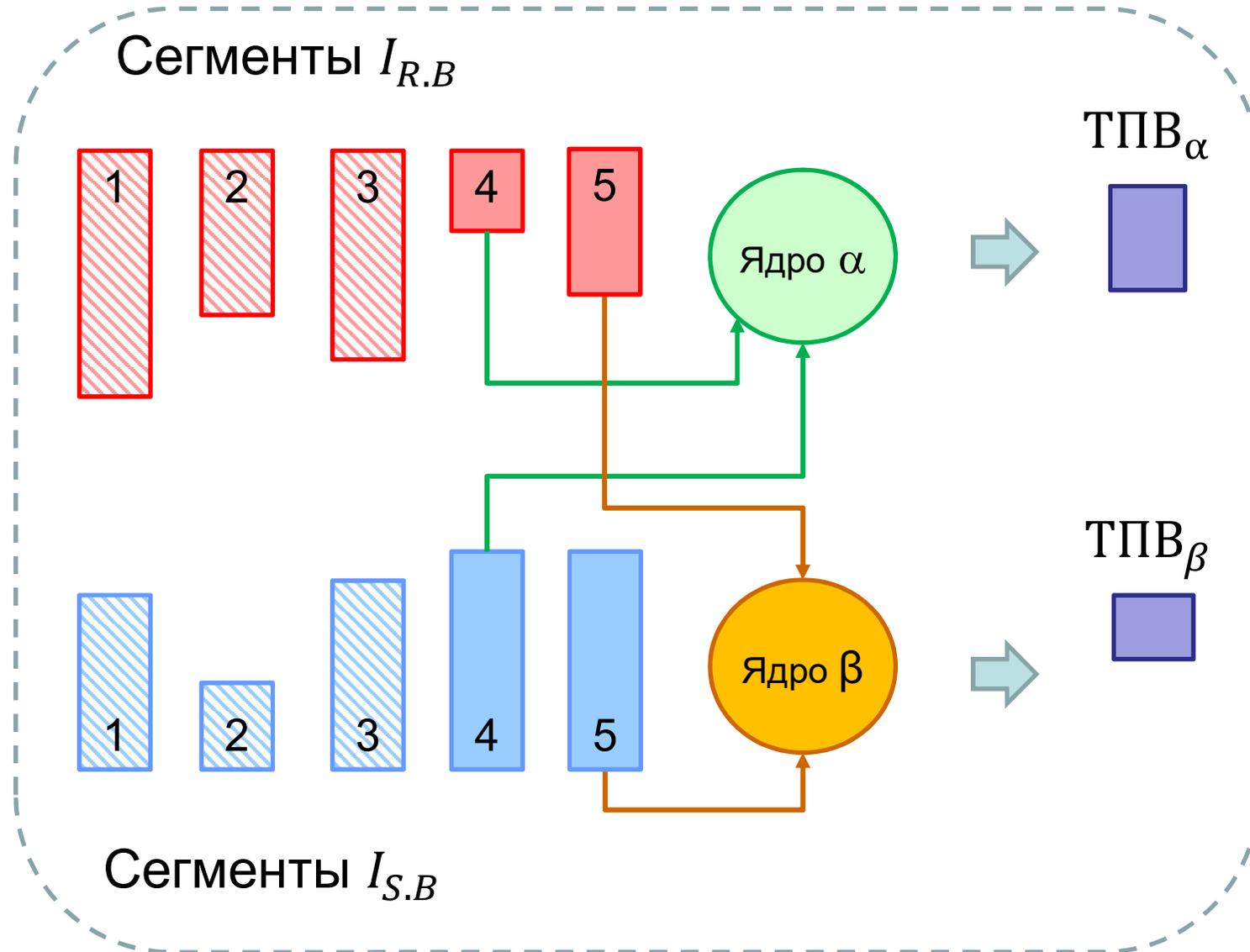
Балансировка загрузки ядер



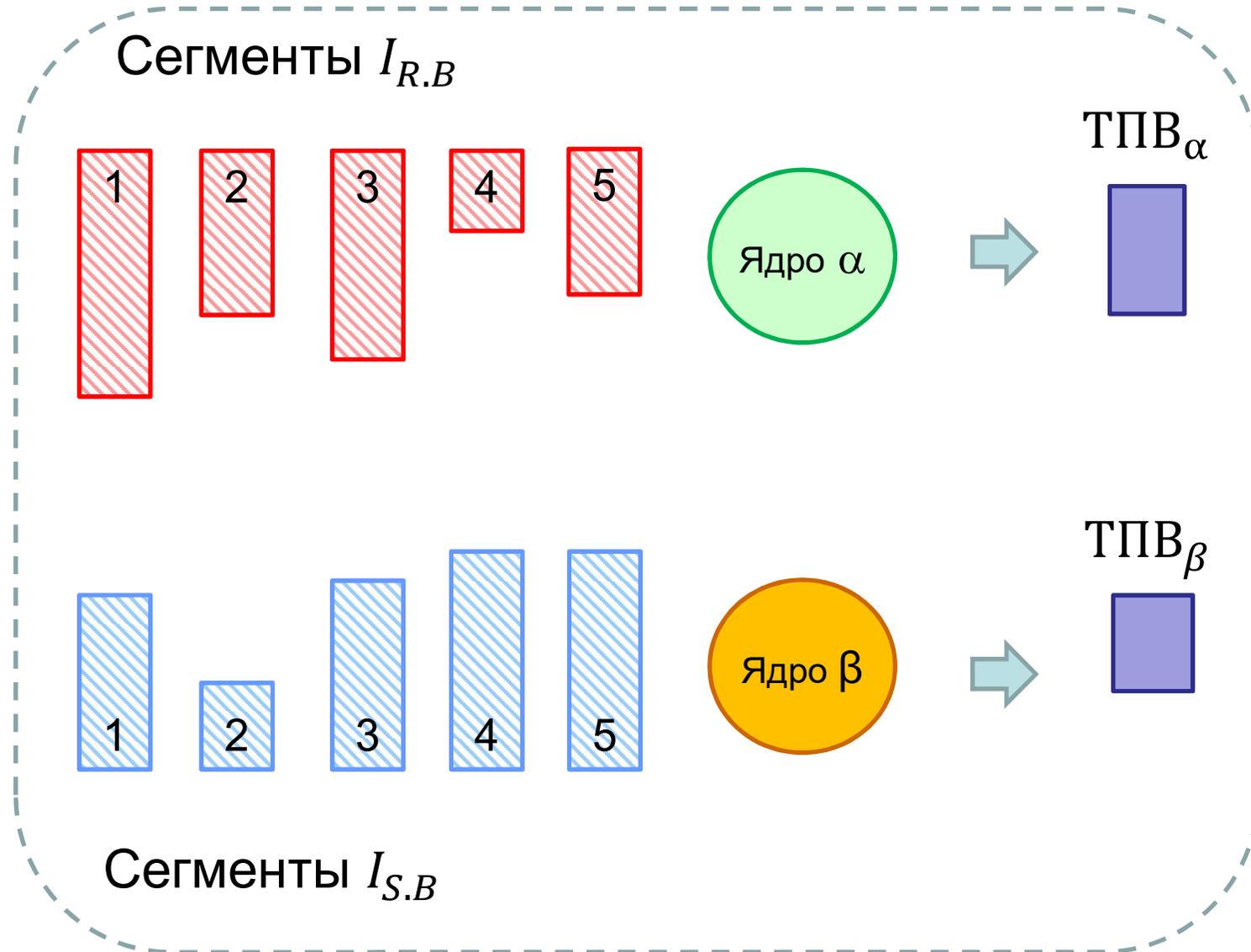
Балансировка загрузки ядер



Балансировка загрузки ядер



Балансировка загрузки ядер



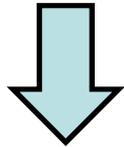
Параллельные алгоритмы выполнения реляционных операций

- Проекция
- Выбор
- Удаление дубликатов
- Пересечение
- Объединение
- Естественное соединение
- Группировка

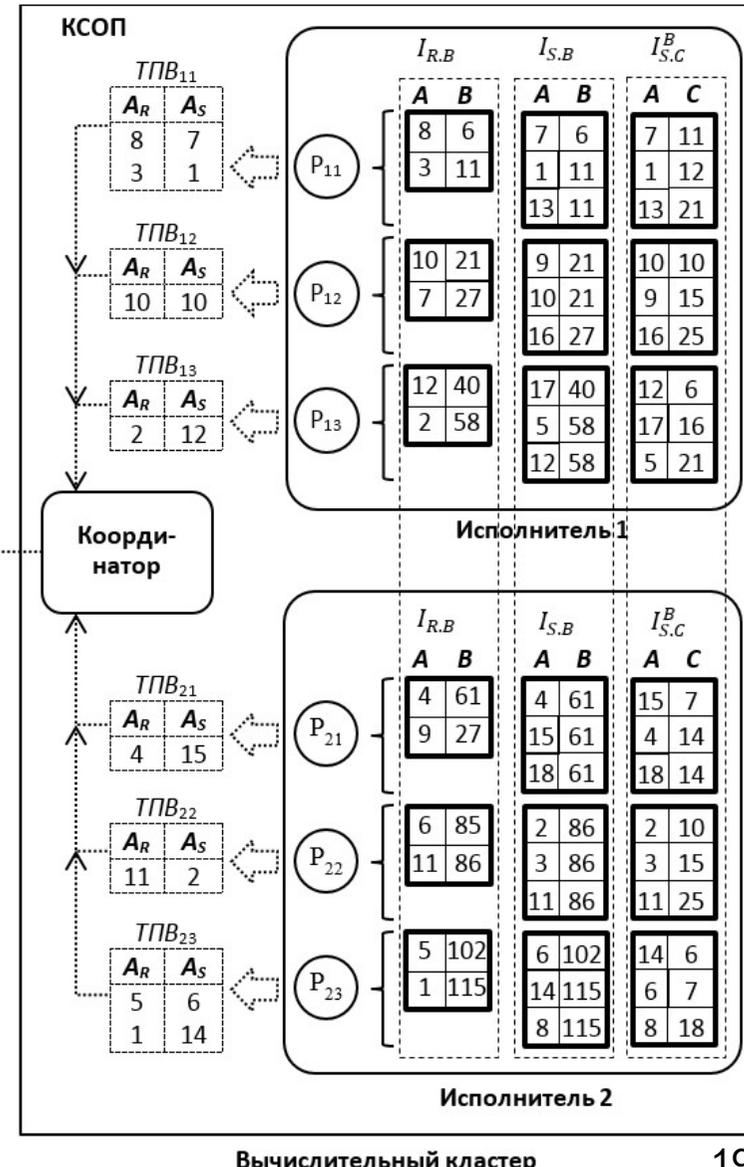
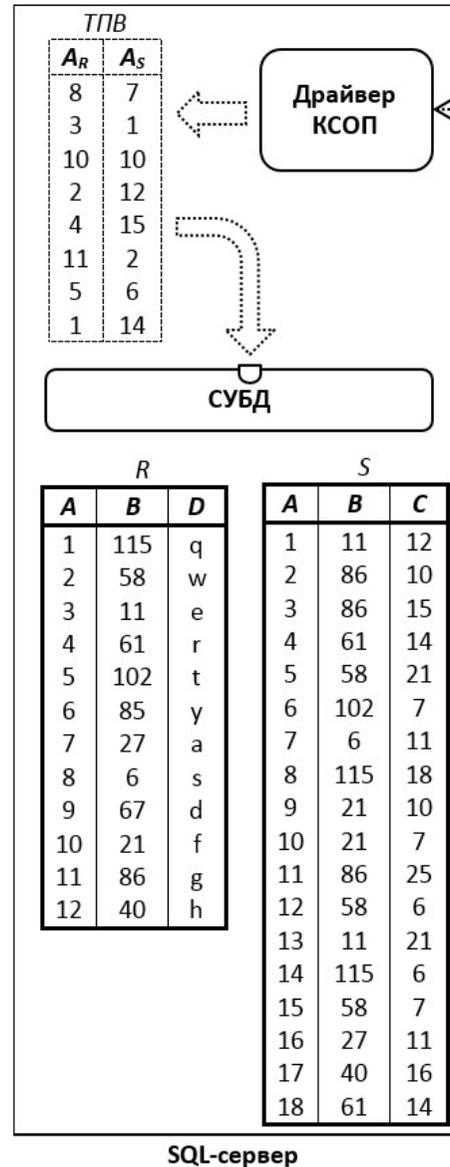
Взаимодействие SQL-сервера с колоночным сопроцессором КСОП

```

SELECT D, C
FROM R, S
WHERE
  R.B = S.B
  AND C < 13;
    
```



$$\pi_{I_{R,B} \cdot A \rightarrow A_R, I_{S,B} \cdot A \rightarrow A_S} \left(I_{R,B} \bowtie_{\begin{matrix} I_{R,B} \cdot B = \\ I_{S,B} \cdot B \end{matrix}} \left(I_{S,B} \bowtie \sigma_{C < 13} (I_{S,C}^B) \right) \right)$$



Реконструкция резултата на SQL-сервере

| <i>R</i> | | | <i>S</i> | | |
|----------|----------|----------|----------|----------|----------|
| <i>A</i> | <i>B</i> | <i>D</i> | <i>A</i> | <i>B</i> | <i>C</i> |
| 1 | 115 | q | 1 | 11 | 12 |
| 2 | 58 | w | 2 | 86 | 10 |
| 3 | 11 | e | 3 | 86 | 15 |
| 4 | 61 | r | 4 | 61 | 14 |
| 5 | 102 | t | 5 | 58 | 21 |
| 6 | 85 | y | 6 | 102 | 7 |
| 7 | 27 | a | 7 | 6 | 11 |
| 8 | 6 | s | 8 | 115 | 18 |
| 9 | 67 | d | 9 | 21 | 10 |
| 10 | 21 | f | 10 | 21 | 7 |
| 11 | 86 | g | 11 | 86 | 25 |
| 12 | 40 | h | 12 | 58 | 6 |
| | | | 13 | 11 | 21 |
| | | | 14 | 115 | 6 |
| | | | 15 | 58 | 7 |
| | | | 16 | 27 | 11 |
| | | | 17 | 40 | 16 |
| | | | 18 | 61 | 14 |

SELECT D, C

FROM

R INNER JOIN (

ТПВ INNER JOIN S ON (S.A = ТПВ.A_S)

) ON (R.A = ТПВ.A_R)

| <i>ТПВ</i> | | Резултат | |
|----------------------|----------------------|----------|----------|
| <i>A_R</i> | <i>A_S</i> | <i>D</i> | <i>C</i> |
| 8 | 7 | s | 11 |
| 3 | 1 | e | 12 |
| 10 | 10 | f | 7 |
| 2 | 12 | w | 6 |
| 4 | 15 | r | 7 |
| 11 | 2 | g | 10 |
| 5 | 6 | t | 7 |
| 1 | 14 | q | 6 |

Формальное определение декомпозиции естественного соединения $Q = \pi_{*\setminus A}(R) \bowtie \pi_{*\setminus A}(S)$

Схема БД: $R(A, B, D); S(A, B, C)$.

Колоночные индексы $I_{R.B}$ и $I_{S.B}$, для которых задана доменно-интервальная фрагментация степени k по атрибуту B :

$$I_{R.B} = \bigcup_{i=0}^{k-1} I_{R.B}^i \quad I_{S.B} = \bigcup_{i=0}^{k-1} I_{S.B}^i$$

$$P^i = \pi_{I_{R.B}^i \cdot A \rightarrow A_R, I_{S.B}^i \cdot A \rightarrow A_S} \left(I_{R.B}^i \bowtie_{I_{R.B}^i \cdot B = I_{S.B}^i \cdot B} I_{S.B}^i \right)$$

$$P = \bigcup_{i=0}^{k-1} P^i$$

$$Q = \left\{ \left(\&_R(p.A_R).B, \&_S(p.A_S).C, \&_R(p.A_R).D \right) \mid p \in P \right\}$$

Реализация КСОП

- Язык программирования: Си
- Технологии параллельного программирования: MPI и OpenMP
- Исходные тексты приложения свободно доступны в сети Интернет:
<https://github.com/elena-ivanova/columnindices/>

Характеристики вычислительных систем

| | «Торнадо ЮУрГУ» | «RSC PetaStream» |
|--------------------------|--|--|
| Количество узлов: | 384 | 64 |
| Тип процессоров: | 2 x Intel Xeon X5680 (12 ядер по 3.33 ГГц; 2 потока на ядро) | |
| Оперативная память узла: | 24 Гб | |
| Тип сопроцессора: | Intel Xeon Phi SE10X: (61 ядро по 1.1 ГГц; 4 потока на ядро) | Intel Xeon Phi 7120 (61 ядро по 1.24 ГГц) |
| Память сопроцессора: | 8 Гб | 16 Гб |
| Тип системной сети: | InfiniBand QDR | InfiniBand FDR |
| Тип управляющей сети: | Gigabit Ethernet | Gigabit Ethernet |
| Операционная система: | Linux CentOS 6.2 | Linux CentOS 7.0 |

Тестовая база данных

CUSTOMER
SF x 630 000

| |
|-------------|
| A |
| ID_CUSTOMER |
| NAME |
| ADDRESS |
| NATION |
| PHONE |
| ACCTBAL |
| MKTSEGMENT |
| COMMENT |

SF ∈ {1, 10}

ORDERS
SF x 63 000 000

| |
|-------------|
| A |
| ID_ORDER |
| ID_CUSTOMER |
| LINENUMBER |
| ORDERSTATUS |
| TOTALPRICE |
| ORDERDATE |
| PRIORITY |
| CLERK |
| SHIPRIORITY |
| ... |
| COMMENT |

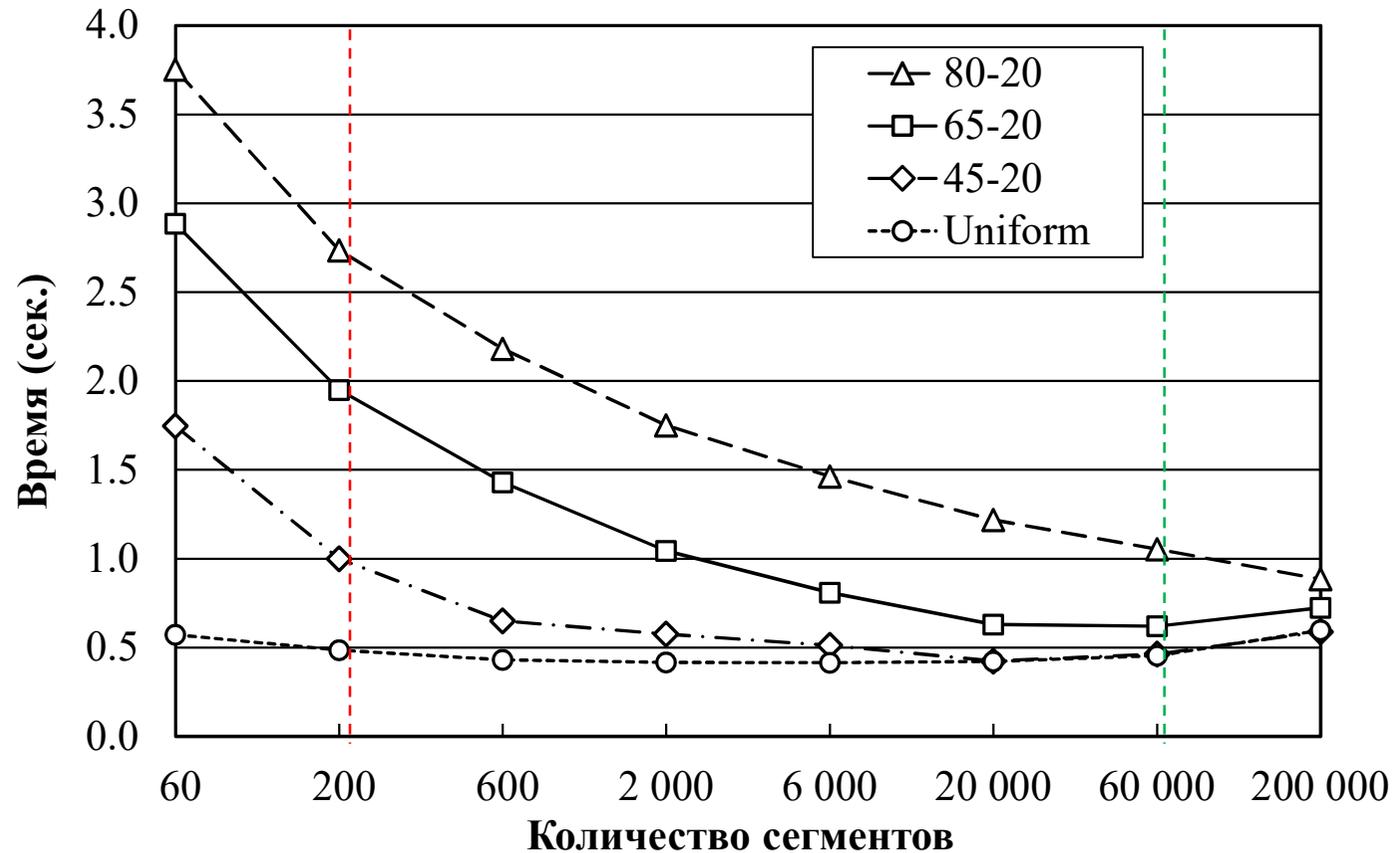
37 атрибутов

Тестовый запрос

Выдать заказы (с информацией о покупателе), общая стоимость которых не превышает величину $Se1 * 100\ 000$:

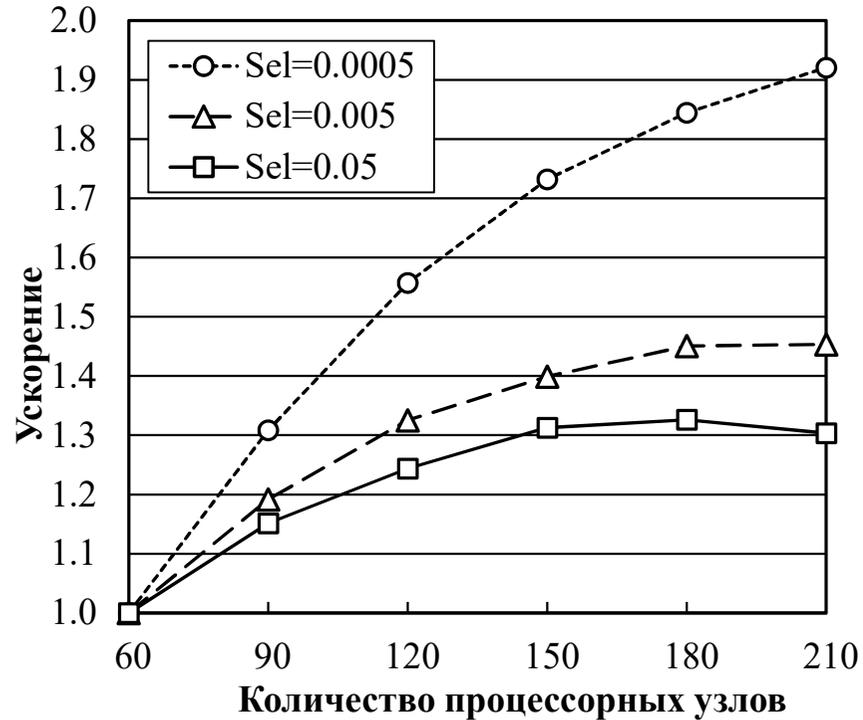
```
SELECT *  
FROM CUSTOMER, ORDERS  
WHERE (CUSTOMER.ID_CUSTOMER=ORDERS.ID_CUSTOMER)  
AND (ORDERS.TOTALPRICE <=  $Se1 * 100\ 000$ ).
```

Балансировка загрузки процессорных ядер Xeon Phi

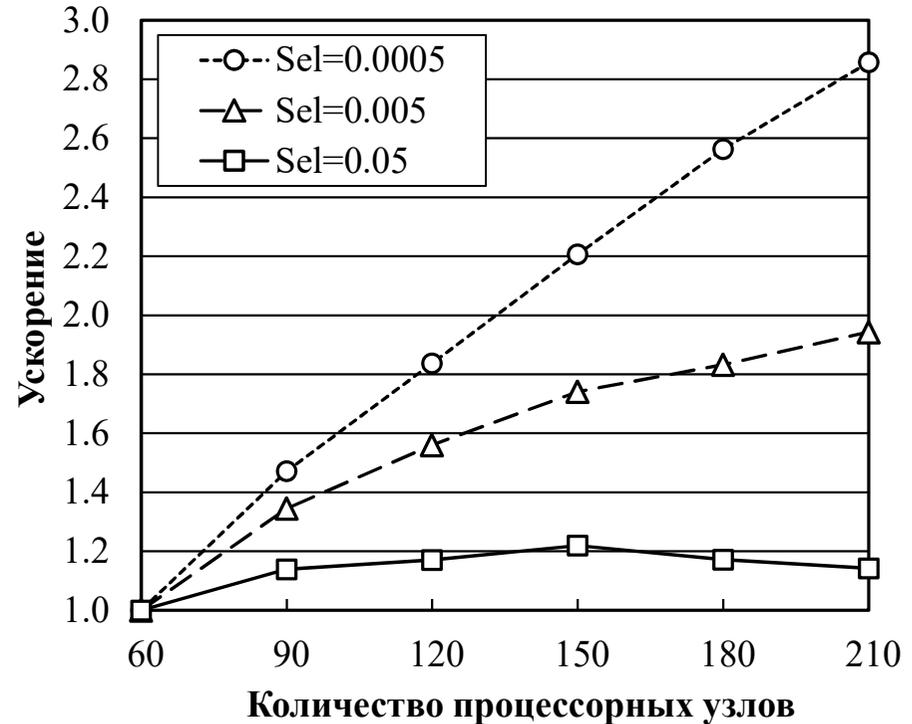


Размер хранилища: 63 млн. записей
Вычислительная система: «Торнадо ЮУрГУ»

Масштабируемость КСОП



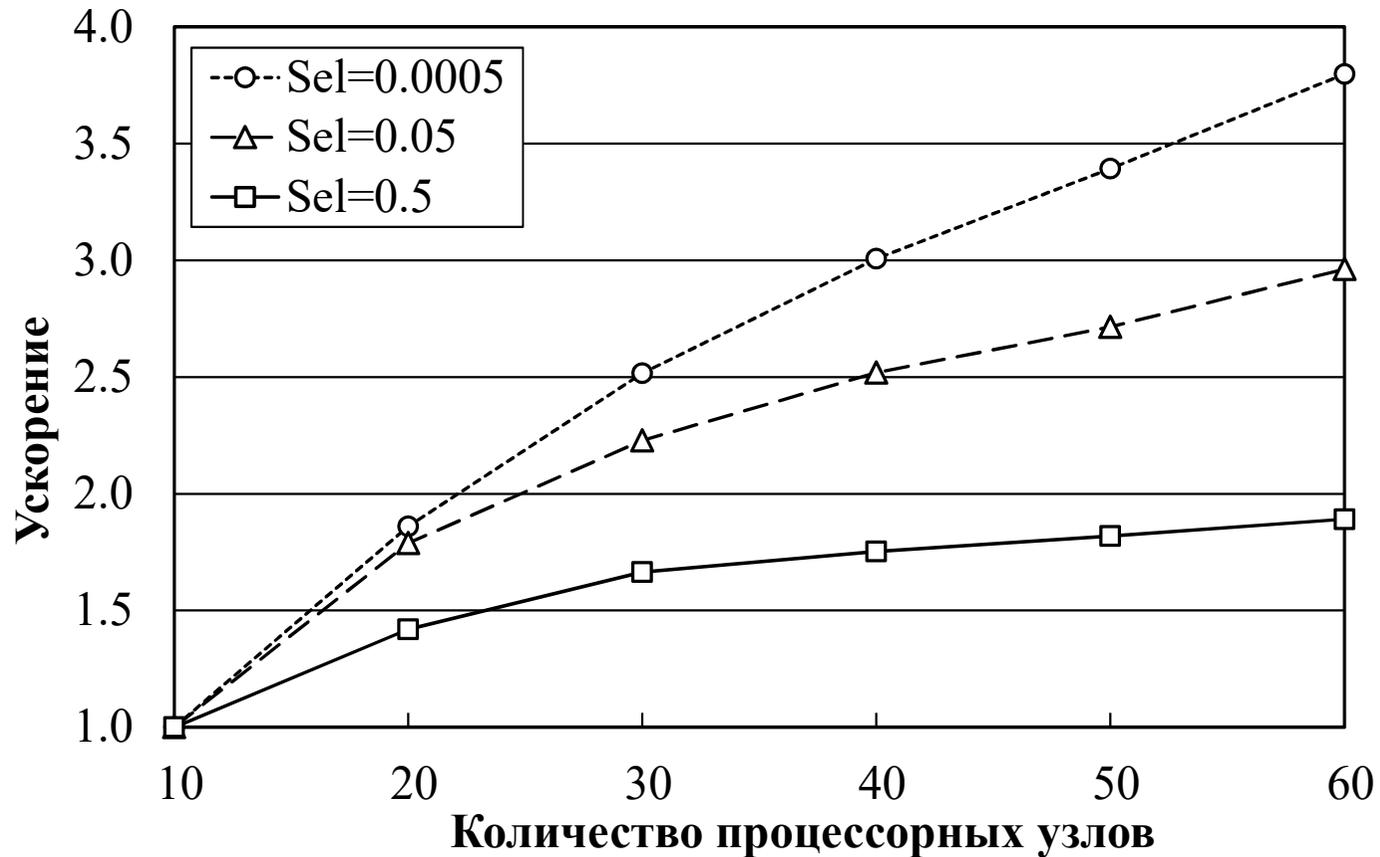
Размер хранилища:
63 млн. записей



Размер хранилища:
630 млн. записей

Вычислительная система: «Торнадо ЮУрГУ»

Масштабируемость КСОП



Размер хранилища: 63 млн. записей
Вычислительная система: «RSC PetaStream»

Использование КСОП при выполнении SQL-запросов

| Конфигурация | Время в минутах | | | | | |
|--|-----------------|------------|------------|------------|------------|------------|
| | Sel=0.0005 | | Sel=0.005 | | Sel=0.05 | |
| | 1-й запуск | 2-й запуск | 1-й запуск | 2-й запуск | 1-й запуск | 2-й запуск |
| PostgreSQL | 7.3 | 1.21 | 7.6 | 1.29 | 7.6 | 1.57 |
| PostgreSQL & B-Trees | 2.62 | 2.34 | 2.83 | 2.51 | 2.83 | 2.63 |
| PostgreSQL & CCOP | 0.073 | 0.008 | 0.65 | 0.05 | 2.03 | 1.72 |
| Ускорение | | | | | | |
| $\frac{t_{PostgreSQL}}{t_{PostgreSQL \& CCOP}}$ | 100 | 151 | 12 | 27 | 4 | 0.9 |
| $\frac{t_{PostgreSQL \& B-Trees}}{t_{PostgreSQL \& CCOP}}$ | 36 | 293 | 4 | 50 | 1.4 | 1.53 |

Основные результаты, выносимые на защиту

- 1) Разработана доменно-колоночная модель представления данных, на базе которой выполнена декомпозиция основных реляционных операций с помощью распределенных колоночных индексов.
- 2) Разработаны высокомасштабируемые параллельные алгоритмы выполнения основных реляционных операций, использующие распределенные колоночные индексы.
- 3) Выполнена реализация колоночного сопроцессора для кластерных вычислительных систем. Общий объем кода на языке Си составил около 2500 строк. Исходные тексты прототипа свободно доступны в Интернет по адресу: <https://github.com/elena-ivanova/colomnindices/>.
- 4) Проведены вычислительные эксперименты, подтверждающие эффективность предложенных подходов.