# Параллельный алгоритм поиска частых наборов для многоядерных процессорных систем

М.Л. Цымблер

#### Содержание

- о Задача поиска частых наборов
- о Алгоритм Apriori
- Алгоритм DIC (Dynamic Itemset Counting)
- Распараллеливание алгоритма DIC
- Вычислительные эксперименты

#### Примеры практических задач

Супермаркет: какие товары часто покупают совместно?



о Поиск побочных эффектов лекарств



#### Поиск частых наборов: определения

- ∘ Объект (*item*).
- ∘ Набор (*itemset*) конечное множество объектов.
- ∘ k-набор (k-itemset) множество из k≥1 объектов.
- $\circ$  *Транзакция* < *TID*, k-набор>.
- $\circ$  База данных транзакций D.
- Поддержка набора (itemset support count) количество транзакций, в которые входит набор.

#### D

<u> </u>	
TID	Набор
10	орешки, пиво, подгузники
20	кофе, пиво, подгузники
30	кофе, пиво, подгузники, хлеб
40	молоко, орешки, хлеб
50	кофе, молоко, орешки, подгузники, хлеб

#### о Пример:

- *supp*({пиво})=3
- *supp*({opeшки,пиво})=1
- *supp*({кофе,подгузники})=2
- *supp*({пиво,подгузники})=3
- *supp*({кофе,подгузники,хлеб})=2

#### Задача поиска частых наборов

- $\circ$  *Минимальная поддержка* minsup пороговое значение поддержки.
- $\circ$  Частый k-набор (k≥1) k-набор с поддержкой не ниже minsup (иначе набор  $pe \partial k u \ddot{u}$ ).
- $\circ L_k$  множество всех частых k-наборов.
- $\circ$  Задача: для заданных D и minsup найти все  $L_k$  (k≥1).

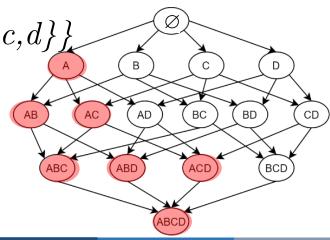
#### D

Набор	
орешки, пиво, подгузники	
кофе, пиво, подгузники	
кофе, пиво, подгузники, хлеб	
молоко, орешки, хлеб	
кофе, молоко, орешки, подгузники, хлеб	
	орешки, пиво, подгузники кофе, пиво, подгузники кофе, пиво, подгузники, хлеб молоко, орешки, хлеб

- При minsup=3
- $L_1$ ={{opeшки}: supp=3, {пиво}: supp=3, {подгузники}: supp=4, {хлеб}: supp=3, {кофе}: supp=3}
- *L*<sub>2</sub>={{пиво,подгузники}: supp=3}
- При *minsup*=2
  - L<sub>1</sub>={{opeшки}: supp=3, {пиво}: supp=3, {подгузники}: supp=4, {хлеб}: supp=3, {кофе}: supp=3, {молоко}: supp=2}
  - $L_2$ ={{кофе,пиво }: supp=2, {пиво,подгузники }: supp=3, ...}
  - L<sub>3</sub>=({кофе,подгузники,хлеб}: supp=2}

#### Алгоритм Apriori: идеи

- о Последовательность шагов
  - k:=1; Найти  $L_1$
  - Для всех *k*≥2
    - На основе  $L_{k-1}$  генерировать  $C_k$ , множество наборов-кандидатов в частые наборы
    - Отбросить заведомо редкие наборы из  $\,C_{\!k}\,$
    - Найти  $L_{k}$  сканируя D и вычисляя поддержку k-наборов из  $C_{k}$
    - Закончить, если  $L_k$ = $\varnothing$
- $\circ$  Генерация кандидатов: JOIN
  - $\bullet \ L_3 \!\!=\!\! \{\{a,b,c\}, \{a,b,d\}, \{a,c,d\}, \{a,c,e\}, \{b,c,d\}\}\}$
  - $C_4 = L_3 \ JOIN \ L_3 = \{\{a,b,c,d\},\{a,c,d,e\}\}$
- о Отбрасывание кандидатов: Apriori
  - Любое непустое подмножество частого набора должно быть частым набором.



02.05.2017

#### **D**, minsup=2

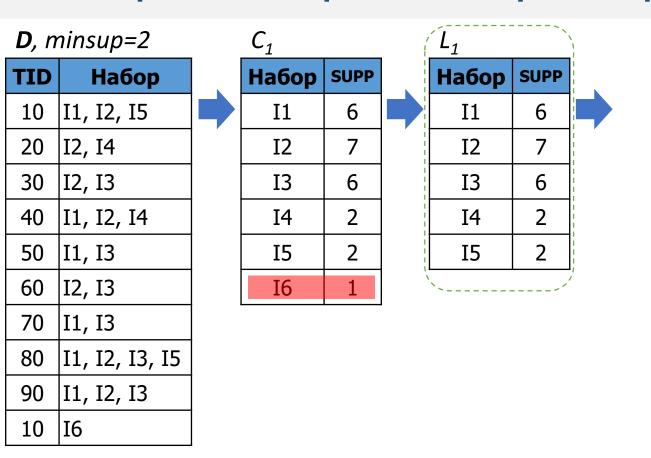
-,ep -		
TID	Набор	
10	I1, I2, I5	
20	I2, I4	
30	I2, I3	
40	I1, I2, I4	
50	I1, I3	
60	I2, I3	
70	I1, I3	
80	I1, I2, I3, I5	
90	I1, I2, I3	
100	I6	

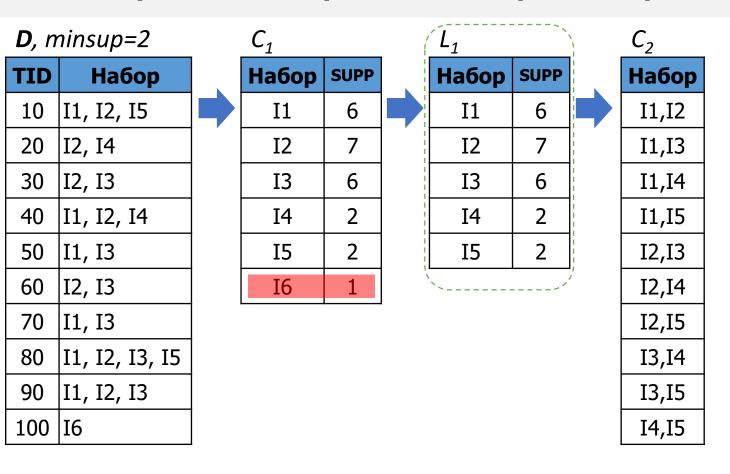
	•	$\sim$
,,	minciin.	_ ,
U.	minsup:	
_ ,		_

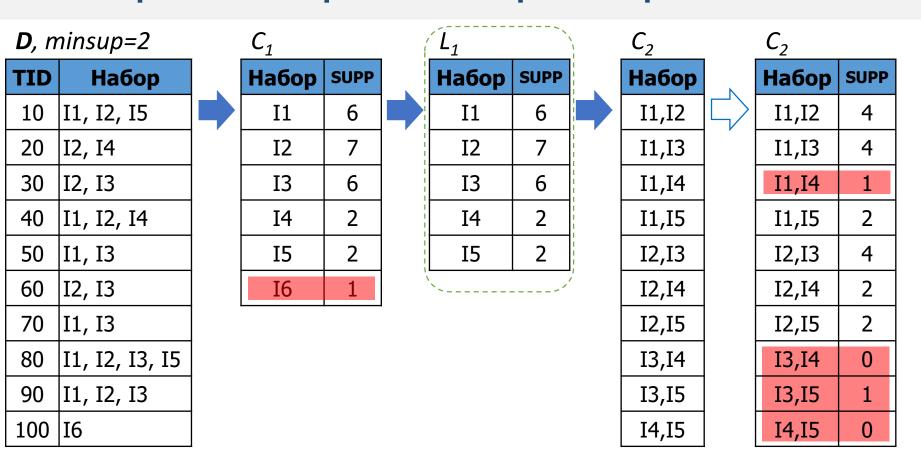
TID	Набор
10	I1, I2, I5
20	I2, I4
30	I2, I3
40	I1, I2, I4
50	I1, I3
60	I2, I3
70	I1, I3
80	I1, I2, I3, I5
90	I1, I2, I3
100	I6

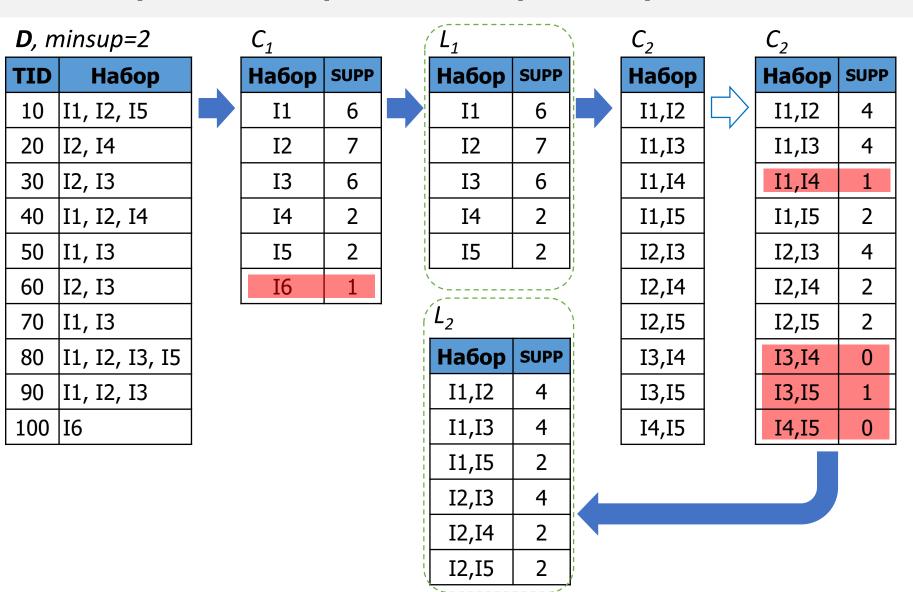
 $C_1$ 

Набор	SUPP
I1	6
I2	7
I3	6
I4	2
I5	2
I6	1









D.	minsup=	2
υ,	пппзир	_

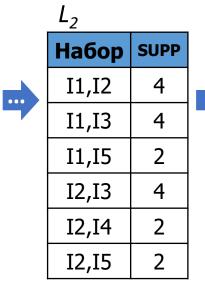
TID	Набор
10	I1, I2, I5
20	I2, I4
30	I2, I3
40	I1, I2, I4
50	I1, I3
60	I2, I3
70	I1, I3
80	I1, I2, I3, I5
90	I1, I2, I3
100	I6

 $L_2$ 

Набор	SUPP
I1,I2	4
I1,I3	4
I1,I5	2
I2,I3	4
I2,I4	2
I2,I5	2

	•	_
"	mincii	n-I
U.	minsu	レーと
_ ,		_

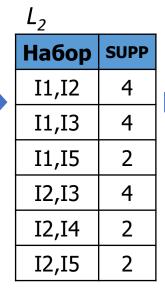
TID	Набор
10	I1, I2, I5
20	I2, I4
30	I2, I3
40	I1, I2, I4
50	I1, I3
60	I2, I3
70	I1, I3
80	I1, I2, I3, I5
90	I1, I2, I3
100	I6



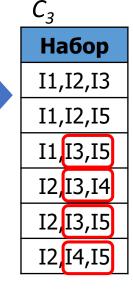
<i>C</i> <sub>3</sub>
Набор
I1,I2,I3
I1,I2,I5
I1, <mark>I3,I5</mark>
I2 <mark>,</mark> I3,I4
I2 <mark>,</mark> I3,I5
I2 <mark>,</mark> I4,I5

	•	$\sim$
"	minclin	<b>ー</b> ノ
U,	minsup	

TID	Набор
10	I1, I2, I5
20	I2, I4
30	I2, I3
40	I1, I2, I4
50	I1, I3
60	I2, I3
70	I1, I3
80	I1, I2, I3, I5
90	I1, I2, I3
100	I6



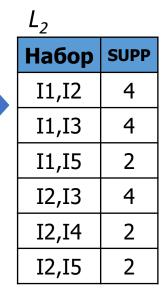
•••



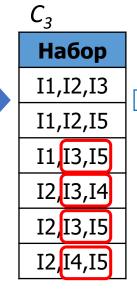
Набор	SUPP	
I1,I2,I3	3 2	
I1,I2,I5	5 2	

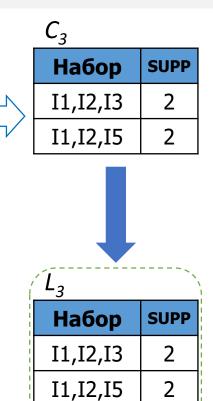


TID	Набор
10	I1, I2, I5
20	I2, I4
30	I2, I3
40	I1, I2, I4
50	I1, I3
60	I2, I3
70	I1, I3
80	I1, I2, I3, I5
90	I1, I2, I3
100	I6



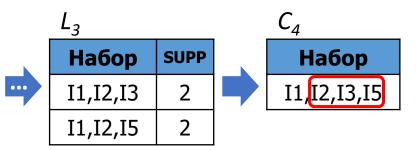
•••





	•	_
,,	minciir	ユーノ
v,	minsup	) — Z

TID	Набор
10	I1, I2, I5
20	I2, I4
30	I2, I3
40	I1, I2, I4
50	I1, I3
60	I2, I3
70	I1, I3
80	I1, I2, I3, I5
90	I1, I2, I3
100	I6



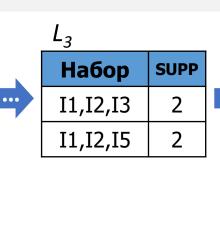


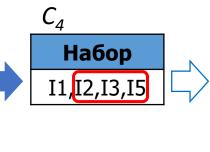
TID	Набор
10	I1, I2, I5
20	I2, I4
30	I2, I3
40	I1, I2, I4
50	I1, I3
60	I2, I3
70	I1, I3
80	I1, I2, I3, I5
90	I1, I2, I3
100	I6

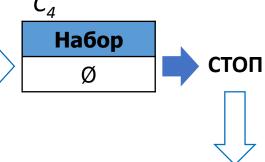




	<u> </u>
TID	Набор
10	I1, I2, I5
20	I2, I4
30	I2, I3
40	I1, I2, I4
50	I1, I3
60	I2, I3
70	I1, I3
80	I1, I2, I3, I5
90	I1, I2, I3
100	I6







 $L_3$ 

′	<u> - 1</u>		
	Набор	SUPP	
	I1	6	
	I2	7	
	I3	6	
	I4	2	
	I5	2	

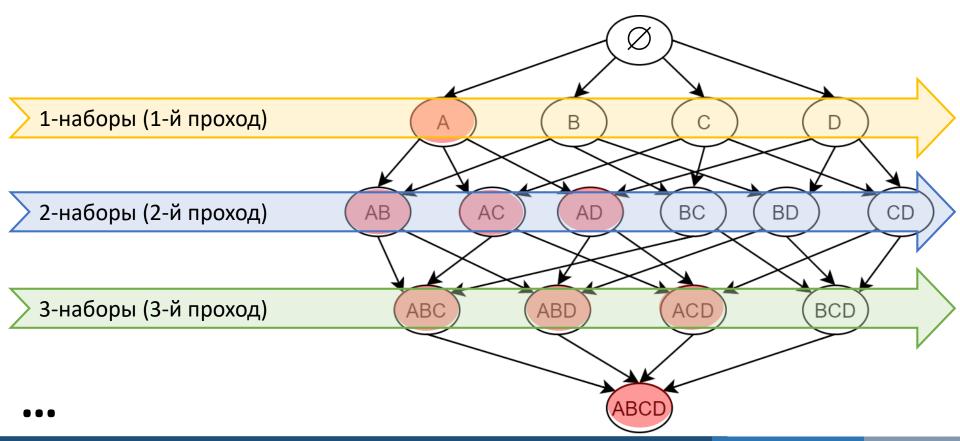
_		
Набор	SUPP	
I1,I2	4	
I1,I3	4	
I1,I5	2	
I2,I3	4	
I2,I4	2	
I2,I5	2	

 $L_2$ 

Набор	SUPP
I1,I2,I3	2
I1,I2,I5	2

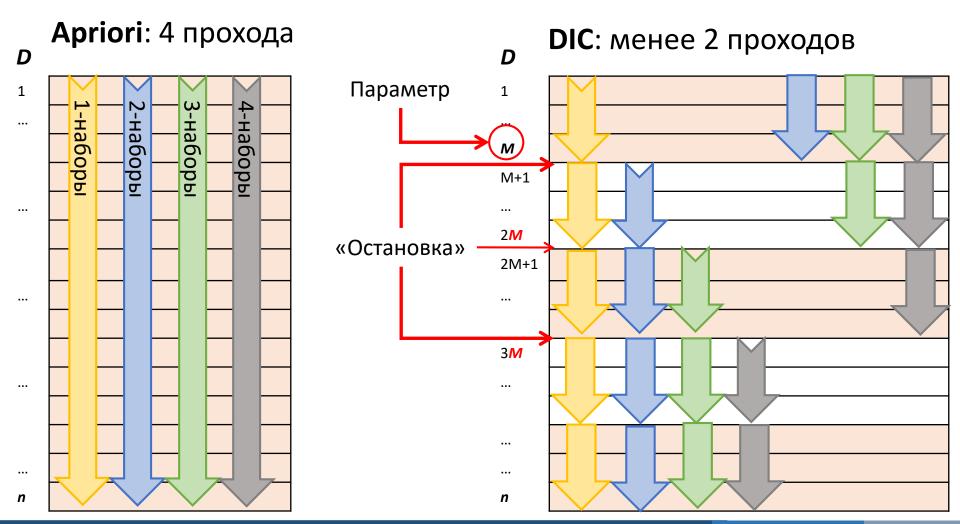
#### Алгоритм DIC (Dynamic Itemset Counting)

- Brin S., Motwani R., Ullman J.D., Tsur S. Dynamic Itemset Counting and Implication Rules for Market Basket Data. Proc. of the SIGMOD'1997, ACM Press, 1997, pp. 255–264.
- Основная идея улучшить Apriori путем сокращения количества проходов по базе данных транзакций



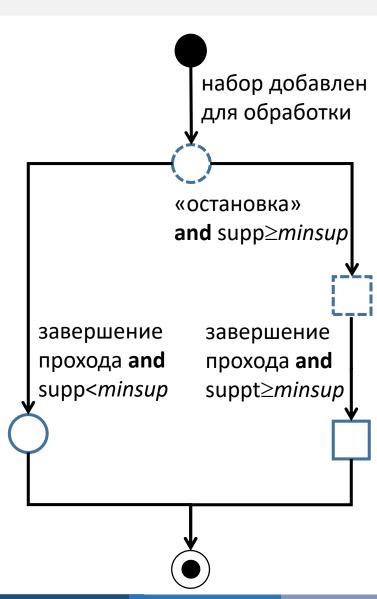
#### DIC: сокращение числа проходов

 $\circ$  Подсчет поддержки (k+1)-набора можно начать до окончания k прохода (если набор, возможно, частый)

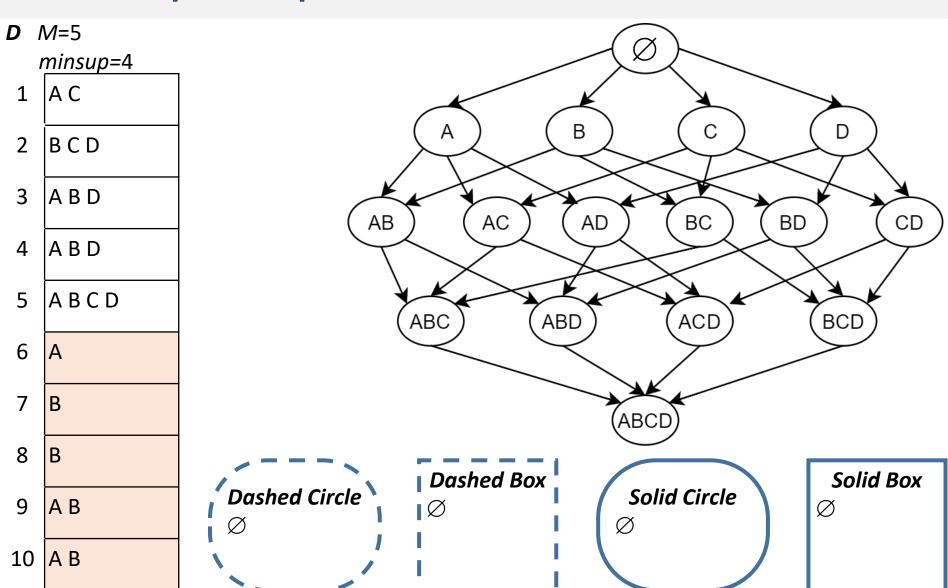


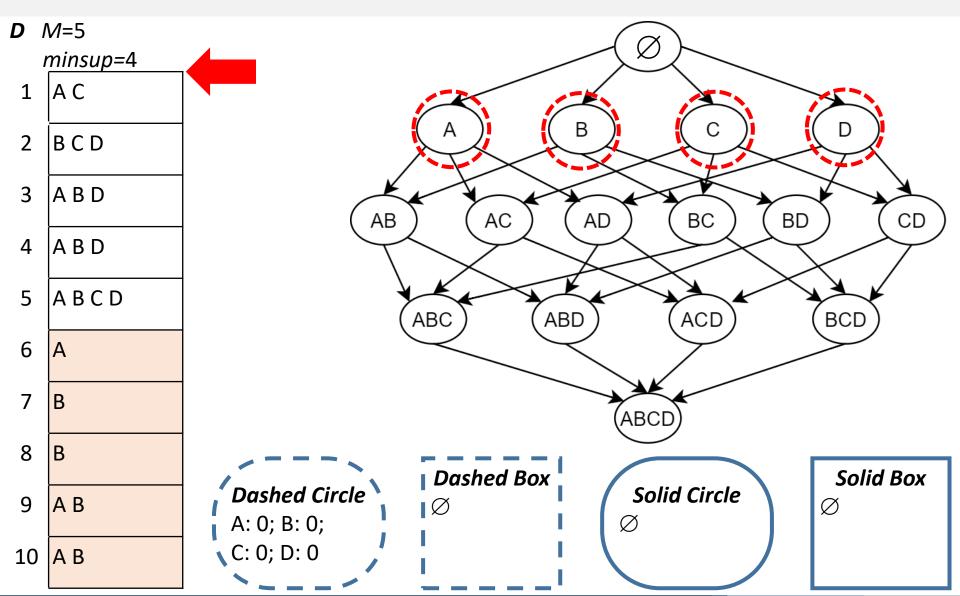
# DIC: разметка наборов

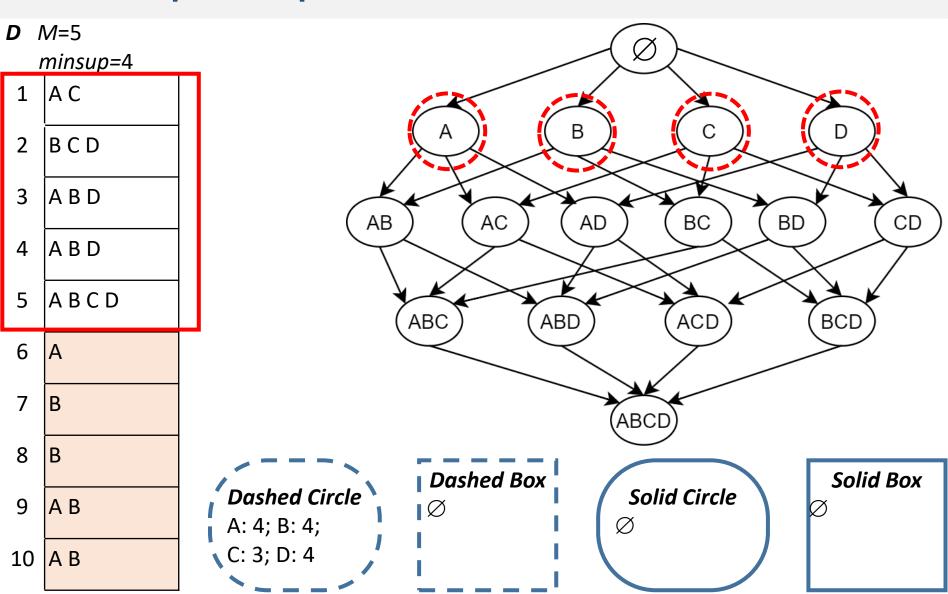
- о подтвержденные частые
  - конечный статус, результат
  - подсчет поддержки не нужен
- o подтвержденные редкие
  - конечный статус, не результат
  - подсчет поддержки не нужен
- неподтвержденные частые
  - промежуточный статус
  - нужен подсчет поддержки
  - 🔘 неподтвержденные редкие
    - начальный статус («кандидат»)
    - нужен подсчет поддержки

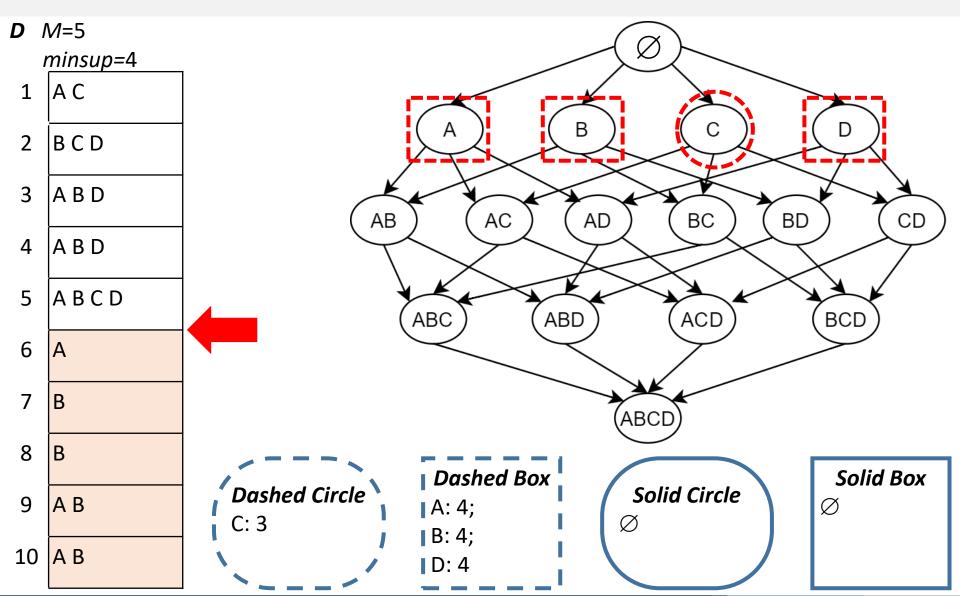


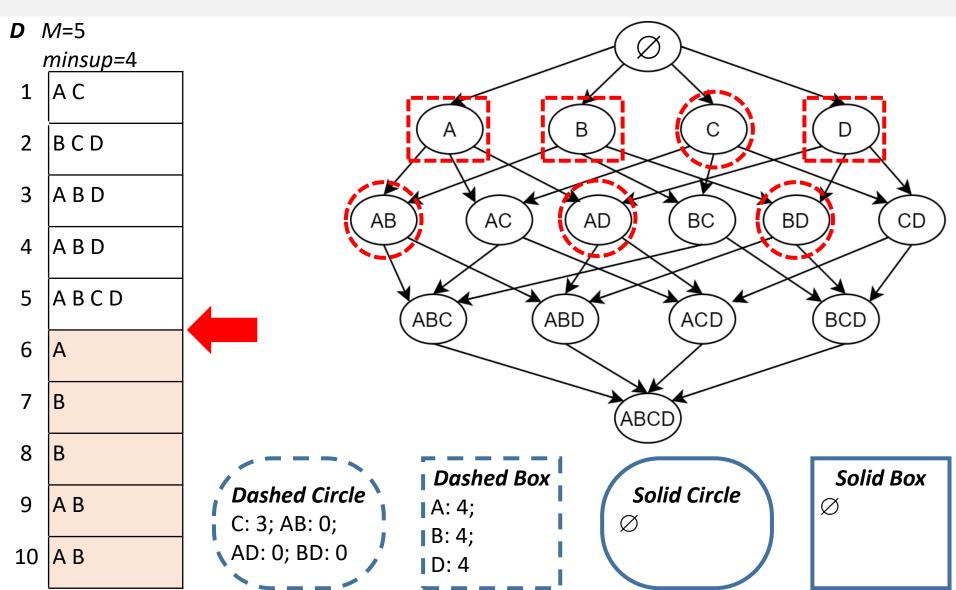
02.05.2017

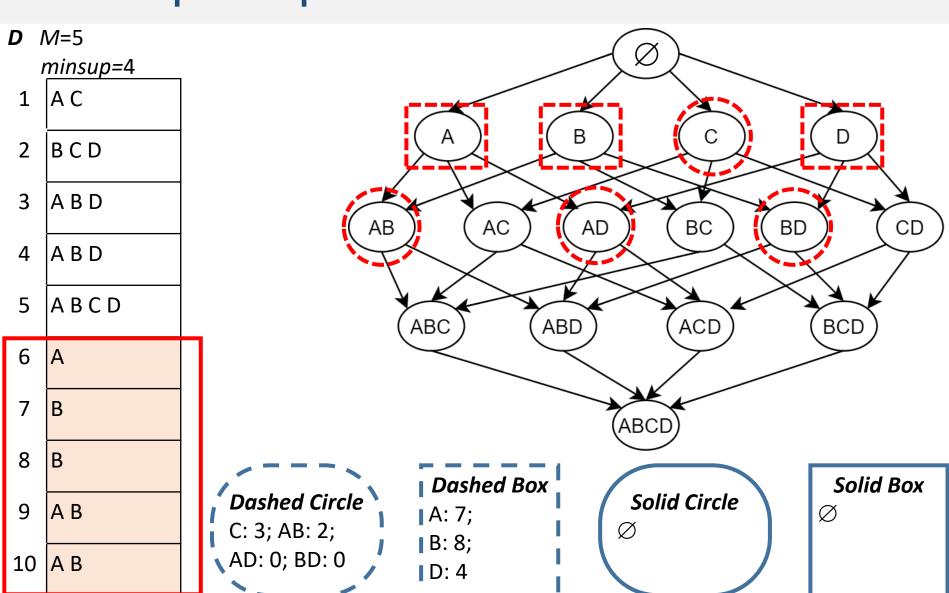


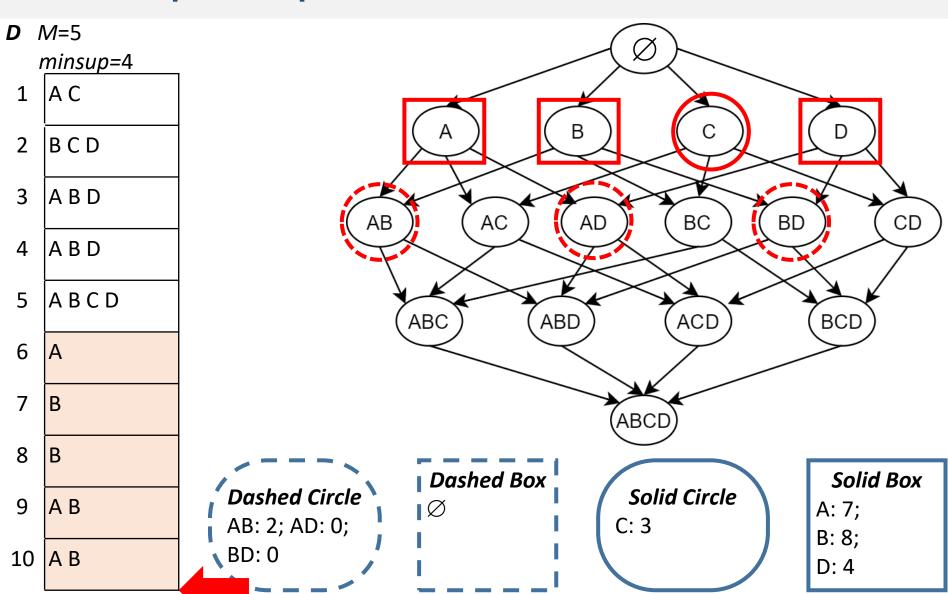


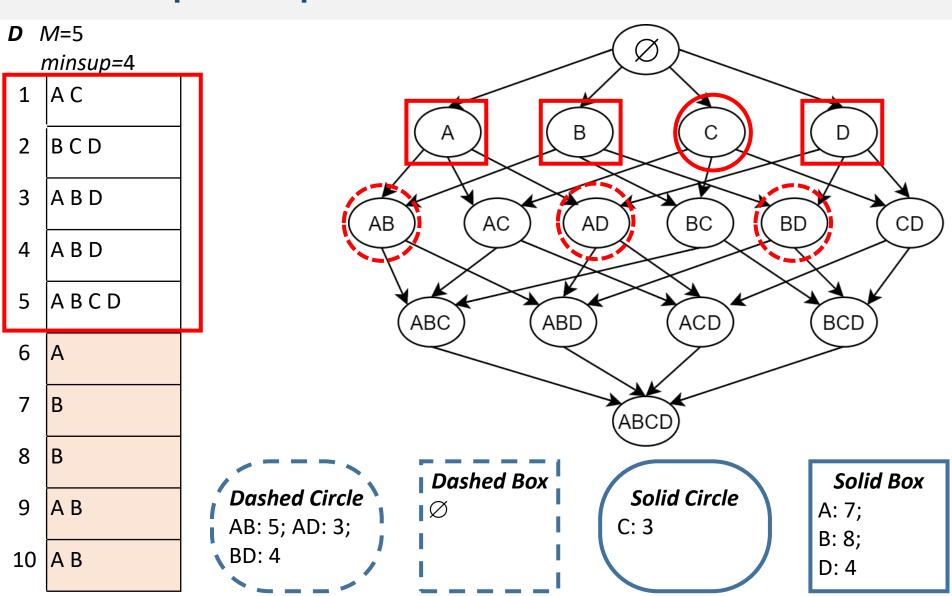


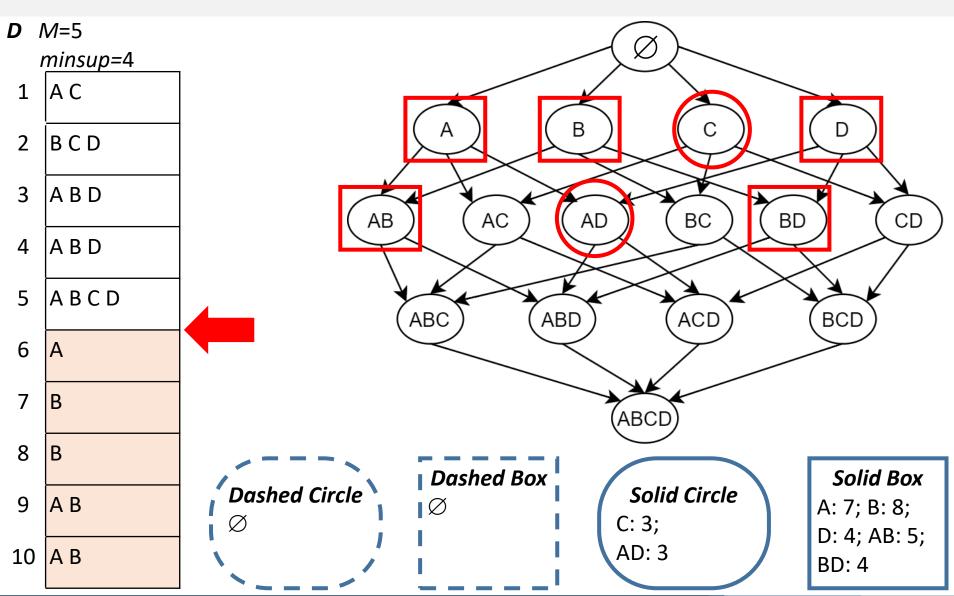












#### **Algorithm 1.** DIC(in $\mathcal{D}$ , in minsup, in M, out $\mathcal{L}$ )

- 2: SOLIDBOX  $\leftarrow \varnothing$ ; SOLIDCIRCLE  $\leftarrow \varnothing$ ; DASHEDBOX  $\leftarrow \varnothing$ DashedCircle  $\leftarrow \mathcal{I}$
- 4: while DashedCircle  $\cup$  DashedBox  $\neq \emptyset$  do

Scan database and rewind if necessary

 $Read(\mathcal{D}, M, Chunk)$ 6: if EOF(D) then

8:

10:

12:

14:

16:

18:

- Rewind( $\mathcal{D}$ ) for all  $T \in Chunk$  do
- - for all  $I \in \mathsf{DASHEDCIRCLE} \cup \mathsf{DASHEDBOX}$  do if  $I \subseteq T$  then
  - $support(I) \leftarrow support(I) + 1$

  - for all  $I \in DASHEDCIRCLE$  do
  - if support $(I) \geq minsup$  then
    - MoveItemset(I, DashedBox)
    - for all  $i \in \mathcal{I}$  do
      - $C \leftarrow I \cup i$
- if  $\forall s \subseteq C$   $s \in SOLIDBOX \cup DASHEDBOX$  then 20: MoveItemset(C, DashedCircle)
- 22: for all  $I \in \mathsf{DASHEDCIRCLE}$  do
- if IsPassCompleted(I) then 24:
  - MoveItemset(I, DASHEDBOX)
- for all  $I \in DASHEDBOX$  do 26: if IsPassCompleted(I) then
- MoveItemset(I, SoLidBox) 28:

 $\mathcal{L} \leftarrow SOLIDBOX$ 

#### ▶ Initialize sets of itemsets ANTOPHTM DIC

#### Работы по тематике DIC

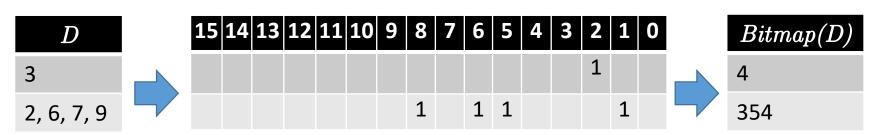
Год	Работа	Вклад
1997	Brin S., Motwani R., Ullman J.D., Tsur S. Dynamic Itemset Counting and Implication Rules for Market Basket Data. Proc. of the SIGMOD'1997, ACM Press, 1997, pp. 255–264.	Оригинальный последовательный алгоритм.
2001	Cheung D.W., Hu K., Xia S. An Adaptive Algorithm for Mining Association Rules on Shared-memory Parallel Machines. Distributed and Parallel Databases, vol. 9, no. 2, pp. 99–132, 2001.	<ul> <li>Параллельный алгоритм</li> <li>Для SMP-систем</li> <li>Техника полного первого прохода</li> <li>Техника отбрасывания кандидатов на основе виртуального разбиения</li> <li>Эксперименты на 12 CPU с БД до 10<sup>8</sup> транзакций; ускорение сублинейное (до 8 на 12 CPU) без учета I/O</li> </ul>
2011	Paranjape-Voditel P., Deshpande U., A DIC-based Distributed Algorithm for Frequent Itemset Generation. Journal of Software SW, vol. 6, no. 2, pp. 306–313, 2011.	<ul> <li>Распределенный алгоритм</li> <li>На базе MPI</li> <li>Обмены возможно частыми наборами между узлами на каждой «остановке»; старт подсчета поддержки наборов, возможно частых на других узлах, но не являющихся таковыми на данном узле</li> <li>Эксперименты на 6, 8, 12 узлах с БД до 10<sup>5</sup> транзакций; ускорение не измерялось; <i>ехр</i> рост количества обменов</li> </ul>
2015	Kumar P., Bhatt P., Choudhury R. Bitwise Dynamic Itemset Counting Algorithm. Proc. of the IEEE ICCIC'2015, IEEE, 2015, pp. 1–4.	Последовательный алгоритм • Вертикальное битовое представление (битовая карта транзакций, в которые входит набор) • Техника отбрасывания заведомо редких наборов • Слабые эксперименты

#### Параллельный DIC

- о Прямое битовое представление данных
- о Реализация структур данных
- о Оптимизация до распараллеливания
- о Распараллеливание
  - Порядок обработки данных
  - Балансировка загрузки

#### Битовое представление данных

- $\circ$  Объекты:  $\widetilde{I}$ =( $i_1$ ,  $i_2$ , ...,  $i_m$ )
- $\circ$  Транзакция (набор):  $T \subseteq \widetilde{I}$  ( $I \subseteq \widetilde{I}$ )
- $\circ$  База данных транзакций: D=(  $T_1$ ,  $T_2$ , ...,  $T_n$ ),  $T_i \subseteq \tilde{I}$
- $\circ$  Битовая маска транзакции: Bitmask(T)
  - $i_p$ -й бит установлен $\iff$   $i_p$   $\in$  T (1 $\le$  p  $\le$  m)
  - иначе  $i_p$ -й бит сброшен
- $\circ$  Битовая карта базы данных: Bitmap(D)
  - $Bitmap(D) = \{Bitmask(T_i): 1 \le j \le n\}$



#### Битовое представление данных

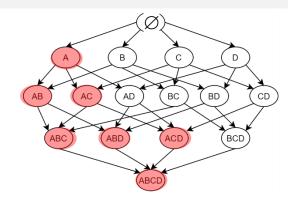
- о Сокращение размера базы данных
  - D можно целиком поместить в оперативную память
- о Ускорение вычисления поддержки
  - $I \subseteq T \iff Bitmask(I)$  **AND** Bitmask(T) = Bitmask(I)
  - 15 14 13 12 11 10 9 8 7 6 5 4 3 2 1 0
    I
    T
- о Усложнение генерации кандидата
  - $c \leftarrow Bitmask(a)$  **OR**  $Bitmask(b) \Leftrightarrow a$  и b отличны в одном бите
  - $(1,2,3,4) \leftarrow (1,2,3)$  *OR* (1,2,4), но (1,2,3) и  $(1,3,5) \rightarrow \emptyset$
- ∘ Случай *т*>64
  - объектов больше sizeof(unsigned long long int)
  - The GNU Multiple Precision Arithmetic Library, <a href="https://gmplib.org/">https://gmplib.org/</a>

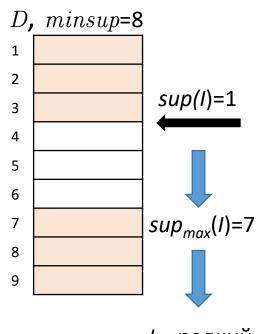
### Реализация структур данных

- База данных транзакций (битовая карта)
  - unsigned long long int \* BITMAP;
  - **BITMAP** = \_mm\_malloc(n\*sizeof(unsigned long long int)) // Выравнивание
- Hабор struct {
  - mask // Битовая маска
  - k // Количество объектов в наборе
  - stopNo // Номер «остановки»
  - supp // Поддержка набора
- о Множества наборов
  - C++ STL (Standard Template Library), класс вектор
  - vector<struct itemset> **DASHED** // Наборы
  - vector<struct itemset> **SOLID** // Наборы
  - vector<struct itemset> CAND // Наборы-кандидаты

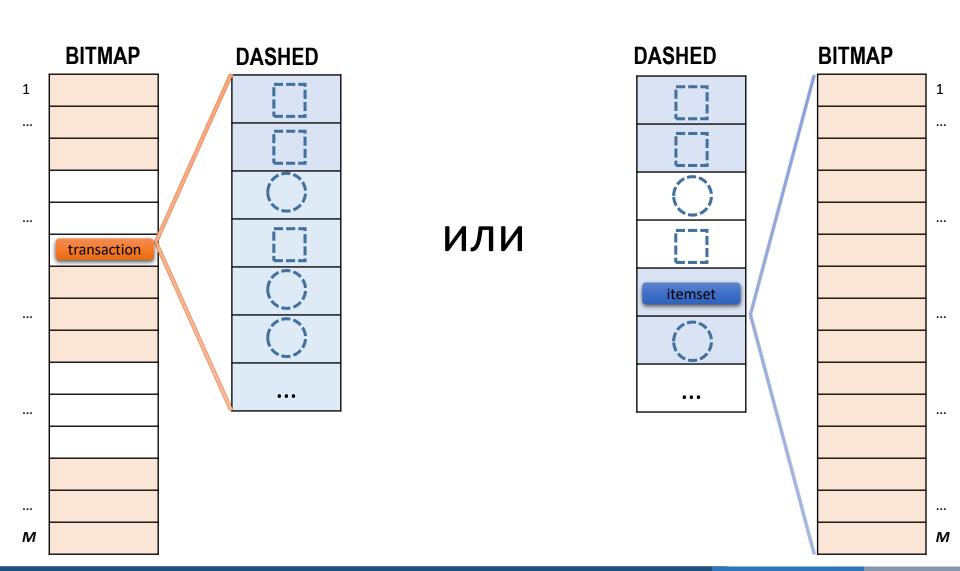
### Оптимизация до распараллеливания

- $\circ$  Первый проход по D полный
  - Найдем  $L_{\scriptscriptstyle 1}$
  - Отбросим редкие 1-наборы
  - Получим меньше 2-наборов кандидатов
- о Отсечение заведомо редких наборов
  - После обработки очередных M транзакций вычисляем максимально возможную поддержку набора  $sup_{max}$ 
    - предполагаем, что набор входит во все остальные транзакции
  - Если  $sup_{max}$  < minsup, то набор редкий
  - Отбрасываем также наборы, в которые входит данный редкий набор

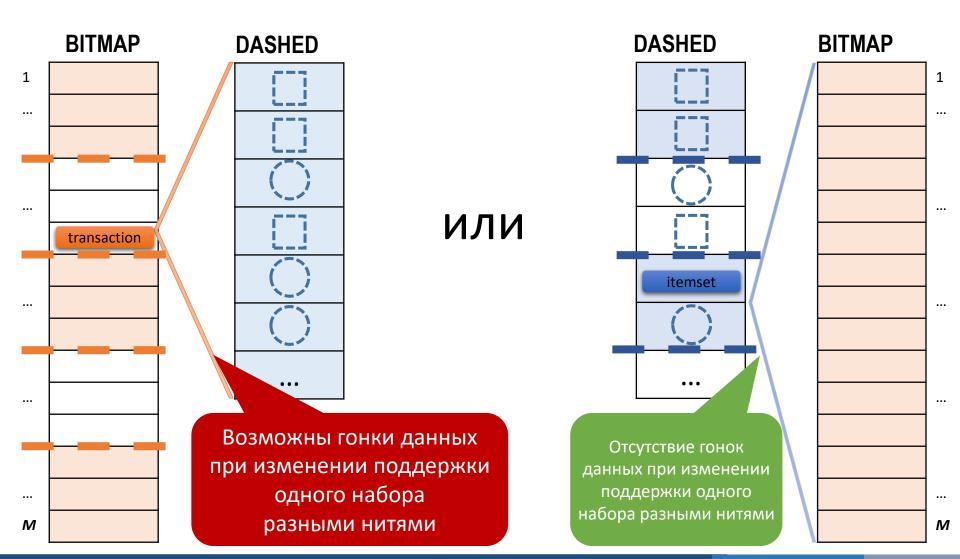




### Распараллеливание: порядок обработки

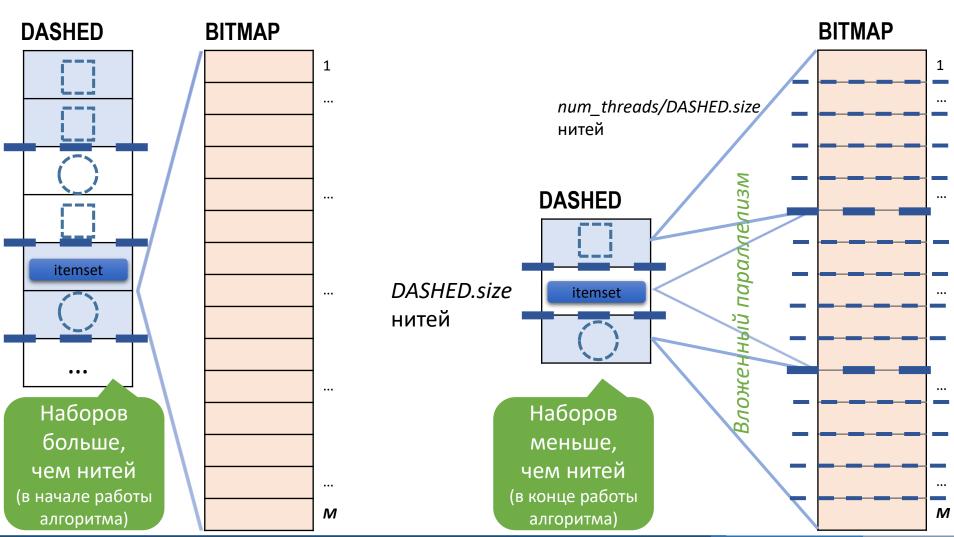


### Распараллеливание: порядок обработки



### Распараллеливание: балансировка

о Количество наборов в **DASHED** изменяется после очередной «остановки»



### Параллельный DIC: подсчет поддержки

```
Algorithm 3. CountSupport(in out DASHED)
    if DASHED.size() \ge num\_of\_threads then
         #pragma omp parallel for
         for all I \in DASHED do
             I.stop \leftarrow I.stop + 1
 4:
             for all T \in \mathcal{B}[first] .. \mathcal{B}[last] do
                 if I.mask AND T = I.mask then
 6:
                      I.supp \leftarrow I.supp + 1
 8: else
         omp_set_nested(true)
         #pragma omp parallel for
10:
         num_threads(DASHED.size())
         for all I \in DASHED do
12:
             I.stop \leftarrow I.stop + 1
             #pragma omp parallel for reduction(+:I.supp) num_threads(\lceil \frac{num\_of\_threads}{DASHED.size()} \rceil)
14:
             for all T \in \mathcal{B}[first] .. \mathcal{B}[last] do
16:
                 if I.mask AND T = I.mask then
                      I.supp \leftarrow I.supp + 1
18:
```

02.05.2017

Параллельный DIC

- 2: SOLID.init(); DASHED.init()  $k \leftarrow 1$
- 4: **for all**  $i \in 0..m 1$  **do**  $I.fig \leftarrow \texttt{NIL}; I.bitmask \leftarrow 0$
- 6:  $I.mask \leftarrow SetBit(I.mask, i)$  $I.stop \leftarrow 0; I.supp \leftarrow 0; I.k \leftarrow k$
- 8: SOLID. $push\_back(I)$  $stop_{max} \leftarrow \lceil \frac{n}{M} \rceil; stop \leftarrow 0$
- 10: FirstPass(SOLID, DASHED) while not DASHED.empty() do
- 12: ▷ Scan database and rewind if necessary
- 14:  $stop \leftarrow stop + 1$  $stop > stop_{max}$  then  $stop \leftarrow 1$
- 16:  $first \leftarrow (stop 1) \cdot M; \ last \leftarrow stop \cdot M 1 \\ k \leftarrow k + 1$
- 18: CountSupport(DASHED)
  CutDashedCircle(DASHED)
- 20: GenCandidates(DASHED)
  CheckFullPass(DASHED)
- 22:  $\mathcal{L} \leftarrow \{I \in \text{SOLID}, I.fig = \text{BOX}\}$

02.05.2017

### Параллельный DIC: отсечение кандидатов

```
Algorithm 4. CutDashedCircle(in out DASHED)
    #pragma omp parallel for
 2: for all I \in \mathsf{DASHED} and I.fig = \mathsf{CIRCLE} do
       if I.supp \ge minsup then
             4:
           I.fig \leftarrow BOX
       else
 6:
                           > Prune clearly infrequent itemset
           supp_{max} \leftarrow I.supp + M \cdot (stop_{max} - I.stop)
 8:
           if supp_{max} < minsup then
              I.fig \leftarrow \text{NIL}
10:
                      > Prune supersets of infrequent itemset
               for all J \in DASHED and J.fig = CIRCLE do
12:
                  if I.mask AND J.mask = I.mask then
                      J.fig \leftarrow \text{NIL}
14:
    DASHED.erase(\forall I, I.fig = NIL)
```

# Параллельный DIC:

#### проверка полного прохода для кандидатов

```
Algorithm 5. CheckFullPass(in out DASHED)
```

```
#pragma omp parallel for
2: for all I \in \text{DASHED} do
    if I.stop = stop_{max} then
4:    if I.supp \geq minsup then
    I.fig \leftarrow \text{BOX}
6:    SOLID.push\_back(I)
I.fig \leftarrow \text{NIL}
8: DASHED.erase(\forall I, I.fig = \text{NIL})
```

### Эксперименты

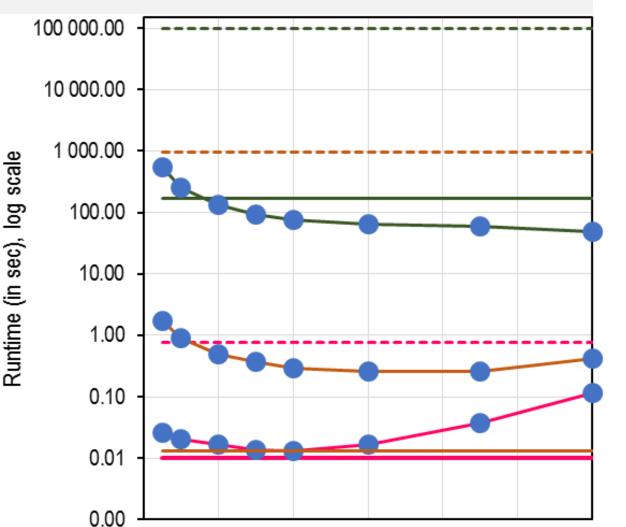
#### о Платформы

Спецификация	Процессор	Сопроцессор
Модель, Intel Xeon	X5680	Phi SE10X
Кол-во ядер	6	60
Тактовая частота, GGz	3.33	1.1
Гиперпоточность	2	4
Производительность, TFLOPS	0.371	1.076
Память, Gb	24	8

#### о Тестовые наборы

Набор	Кол-во транзакций	Кол-во объектов	Источник
SKIN	245 057	11	Dhall A., et. al. Adaptive Digital Makeup // ISVC 2009. LNCS, vol. 5876. Springer, 2009, pp. 728–736.
RECORDLINK	574 912	29	Sariyar M., et al. Controlling False Match Rates in Record Linkage Using Extreme Value Theory // Journal of Biomedical Informatics, vol. 44, no. 4, pp. 648–654, 2011.
20M	2*10 <sup>7</sup>	64	Синтез с помощью генератора IBM Quest (как у авторов DIC). Ср. длина транзакции 20, ср. длина частого набора 10.





- (1) By Cristian Borgelt
- (2) By Bart Goethals

Haбop *SKIN*—— Serial Apriori <sup>(1)</sup>

·---- Serial DIC (M=120 000)

—●— Parallel DIC (M=120 000)

Набор RECORDLINK

—— Serial Apriori<sup>(1)</sup>

---- Serial DIC (M=287 456)

—●— Parallel DIC (M=287 456)

Haбop 20M

——— Serial Apriori<sup>(1)</sup>

---- Serial DIC (M=10 000 000)

—●— Parallel DIC (M=10 000 000)

12

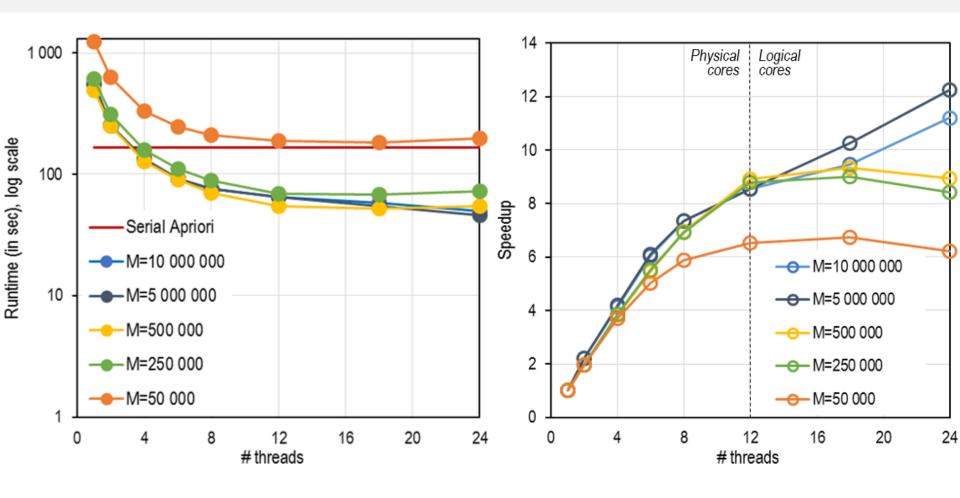
# threads

16

20

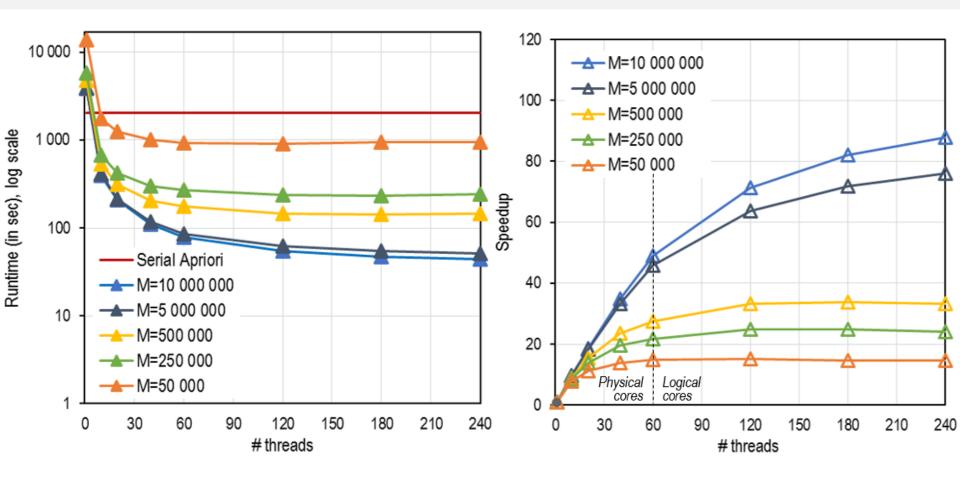
24

## Масштабируемость по M (CPU)



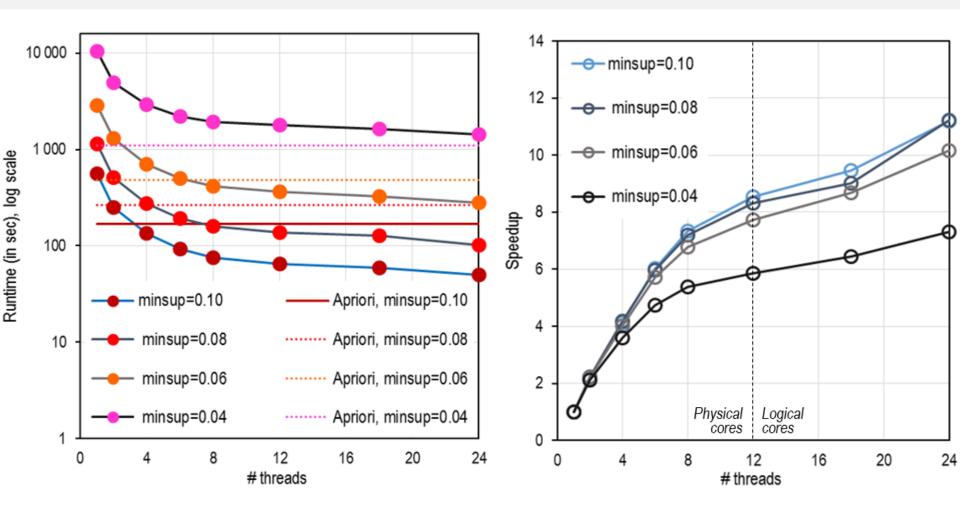
Hабор: 20M minsup=0.1

## Масштабируемость по M (Phi)



Набор: *20М minsup*=0.1

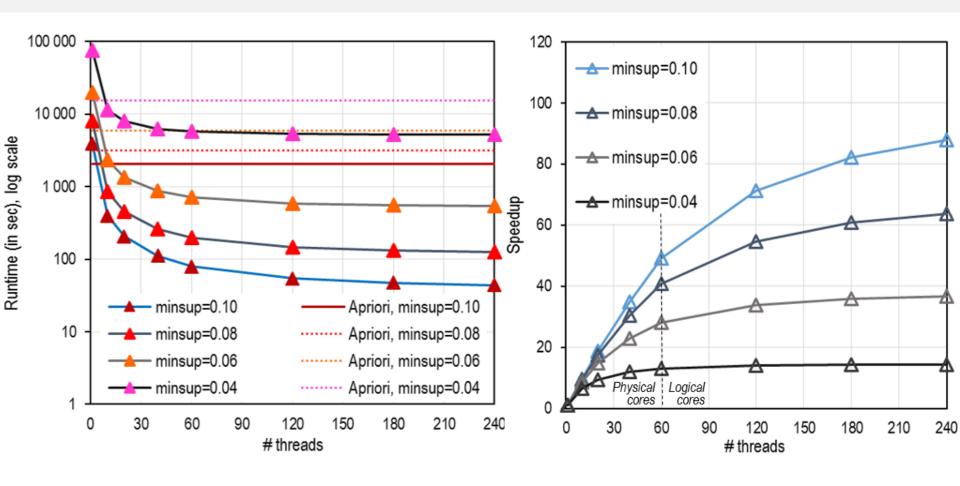
### Масштабируемость по minsup (CPU)



Набор: *20М* 

 $M=10^7$ 

## Масштабируемость по minsup (Phi)



Набор: *20М* 

 $M=10^7$