

Heterogeneous Parallel Programming

COMP4901D

Frequent Itemset Mining on the GPU

Overview

- Frequent item set mining
- GPU-based frequent item set mining

Frequent Itemset Mining (FIM)

Find groups of items, or *itemsets* where all items in a group frequently co-occur in a transaction.

Transaction ID	Items	<u>Minimum support: 2</u>
1	<i>A, B, C, D</i>	1-itemsets (frequent items): <i>A, B, C, D</i>
2	<i>A, B, D</i>	2-itemsets:
3	<i>A, C, D</i>	<i>AB, AC, AD, BD, CD</i>
4	<i>C, D</i>	3-itemsets: <i>ABD, ACD</i>

The Apriori Algorithm

Input: 1) Transaction Database 2) Minimum support

Output: All frequent itemsets

$L_1 = \{\text{All frequent 1-itemsets}\}$

$k = 2$

While ($L_{k-1} \neq \text{empty}$) {

//Generate candidate k-itemsets

Self-join //e.g. **ABC** JOIN **ABD** => **ABCD**

Subset test //e.g., Test all 3-subsets of ABCD

//Generate frequent k-itemsets

Support Counting

$k += 1$

}

Frequent 1-itemsets

Candidate 2-itemsets

Frequent 2-itemsets

Candidate 3-itemsets

Frequent 3-itemsets

...

Candidate (K-1)-itemsets

Frequent (K-1)-itemsets

Candidate K-itemsets

Frequent K-itemsets

GPU Acceleration for APriori

Input: 1) Transaction Database 2) Minimum support

Output: All frequent itemsets

$L_1 = \{\text{All frequent 1-itemsets}\}$

$k = 2$

While ($L_{k-1} \neq \text{empty}$) {

<u>Candidate Generation</u>	
<u>Pure Bitmap-based Impl. (PBI)</u> <u>on the GPU</u>	<u>Trie-based Impl. (TBI)</u> <u>on the CPU</u>
<u>Support Counting on the GPU</u>	

$k += 1$

}

Data Layout

Horizontal data layout

TID	Items
1	<i>A, B, C, D</i>
2	<i>A, B, D</i>
3	<i>A, C, D</i>
4	<i>C, D</i>

Vertical data layout

Item	Transactions
<i>A</i>	1,2,3
<i>B</i>	1,2
<i>C</i>	1,3,4
<i>D</i>	1,2,3,4

Horizontal layout requires scanning all transactions.

Pure Bitmap-based Implementation (PBI)

Bitmap representation for itemsets

	T1	T2	T3	T4
AB	1	1	0	0
AC	1	1	0	1
AD	1	0	0	1
BD	0	1	0	1
CD	0	0	1	1

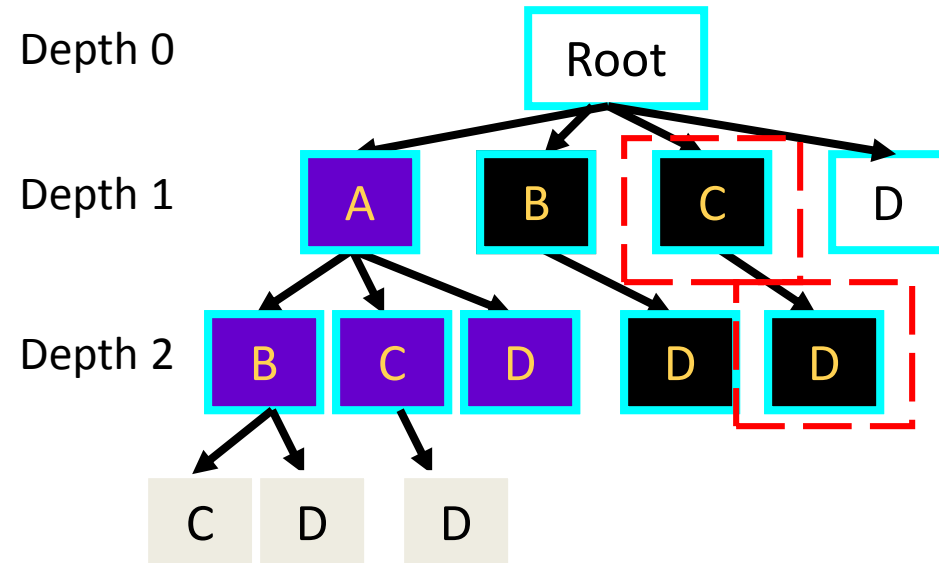
Bitwise OR In Join
(e.g., AB JOIN AD = ABD)

	T1	T2	T3	T4
ABD	1	1	0	1

Binary search
In Subset test
(e.g., 2-subsets {~~AB, AD, BD~~})

One GPU thread generates one candidate itemset.

Trie-based Impl. (TBI)



1-itemsets: {A, B, C, D}

2-itemsets: {AB, AC, AD, BD, **CD**}

AB JOIN AC = ABC

{AB, AC, ~~BC~~}

AB JOIN AD = ABD

{AB, AD, BD}

AC JOIN AD = ACD

{AC, AD, CD}

On CPU
Irregular memory access

Support Counting on the GPU

	T1	T2	T3	T4
AB	1	1	0	0
AC	1	0	1	0
AD	1	1	1	0
BD	1	1	1	0
CD	1	0	1	1

	T1	T2	T3	T4
ABD	1	1	0	0
ACD	1	0	1	0

Intersect = bitwise AND operation

1. Intersect
2. Count

Support Counting on the GPU (Cont.)

Bitmap representation for transactions

	T1	T2	T3	T4
ABD	1	1	0	0
ACD	1	0	1	0

Lookup table

Index	Count
0000	0
0001	1
...	
1000	1
1010	2
1011	3
1100	2

Constant memory

1. Cacheable
2. 64 KB
3. Shared by all GPU threads

Index

0000 0000 0000 0000 (0)
0000 0000 0000 0001 (1)
...
1111 1111 1111 1110 (65534)
1111 1111 1111 1111 (65535)

Count

0
1
...
15
16

1 byte

$$2^{16} = 65536$$

Support Counting on the GPU (Cont.)

	T1	T2	T3	T4
AB	1	1	0	0
AC	1	0	1	0
AD	1	1	1	0
BD	1	1	1	0
CD	1	0	1	1

	T1	T2	T3	T4
ABD	1	1	0	0
ACD	1	0	1	0

Thread block 1

Thread block 2

Support Counting on the GPU (Cont.)

Thread Block

Thread 1

Access vector type int4
In one instruction

Thread 2

Example:

int	int	int	int
-----	-----	-----	-----

AND AND AND AND

int	int	int	int
-----	-----	-----	-----

int	int	int	int
-----	-----	-----	-----

int	int	int	int
-----	-----	-----	-----

AND AND AND AND

int	int	int	int
-----	-----	-----	-----

int	int	int	int
-----	-----	-----	-----

1	1	0	0
---	---	---	---

1	0	1	0
---	---	---	---

1	0	0	0
---	---	---	---

LOOKUP TABLE

Counts of 1's for every 16-bit integer

Parallel Reduce

Support for this itemset

Counts: 1

Support:1

Experimental setup

Platform	Intel Core2 quad-core CPU	NV GTX 280 GPU
Processors	2.66 GHz * 4	1.35 GHz * 30 * 8
Memory Bandwidth (GB/sec)	10.4	141.7
Development Env.	Windows XP + Visual Studio 2005 + CUDA	

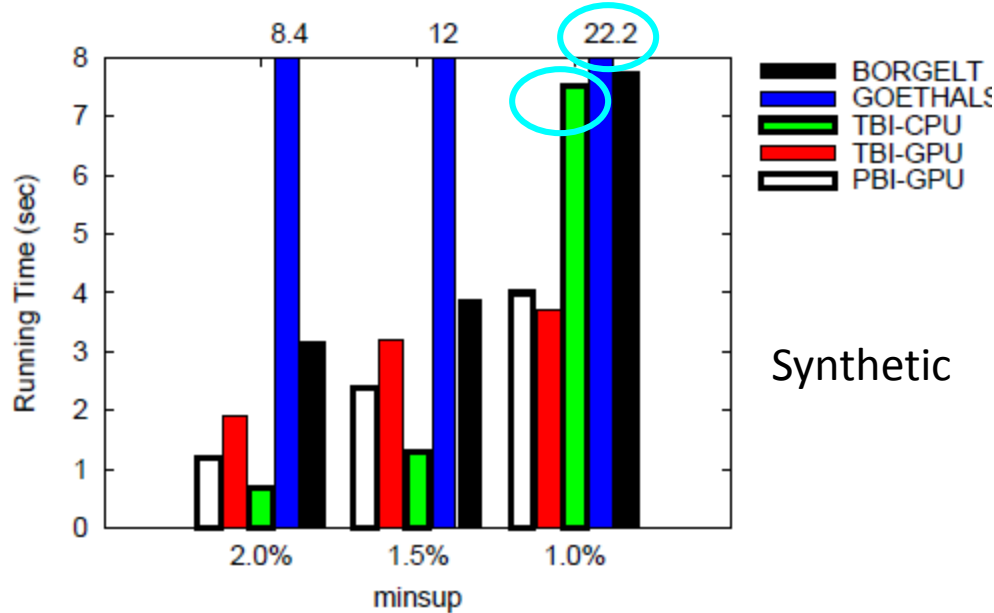
Dataset	#Item	Avg. Length	#Tran	Density	Original data size	Bitmap size
T40I10D100K (synthetic)	1,000	40	100,000	4%	~15MB	12MB
Retail	16,469	10.3	88,162	0.6%	~4MB	180MB
Chess	75	37	3196	49%	~335KB	30KB

Apriori Implementations

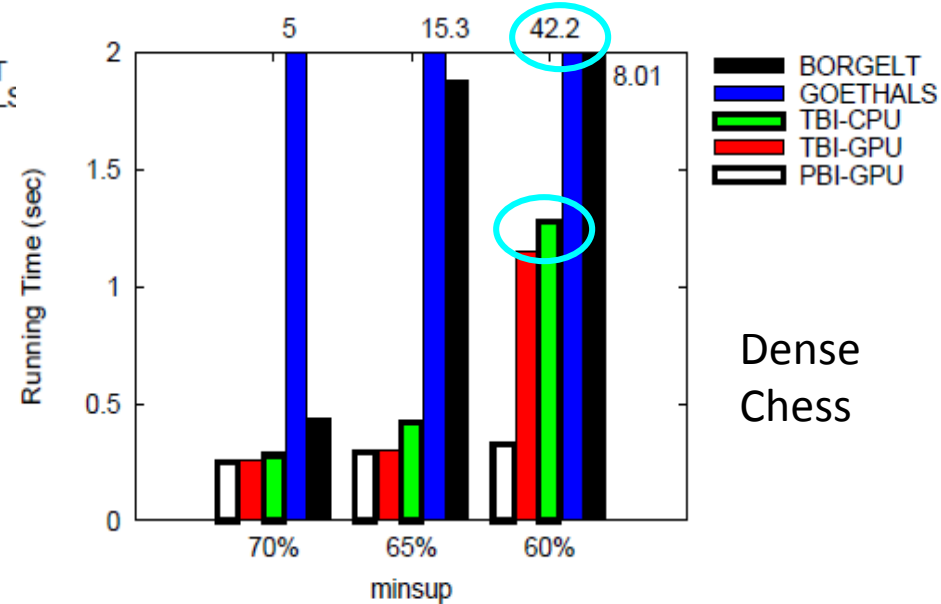
<u>Impl.</u>	<u>Candidate Generation</u>	<u>Support Counting</u>	<u>Itemset</u>	<u>Transactions</u>
PBI-GPU	Multi-threaded on the GPU	Multi-threaded on the GPU	Bitmap	Bitmap
TBI-GPU	Single-threaded on the CPU	Multi-threaded on the GPU	Trie	Bitmap
TBI-CPU	Single-threaded on the CPU	Multi-threaded on the CPU	Trie	Bitmap
GOETHALS	Single-threaded on the CPU	Multi-threaded on the CPU	Trie	Horizontal layout
BORGELT	Single-threaded on the CPU	Single-threaded on the CPU	Trie	Trie

Best Apriori implementation in FIMI repository.
(Frequent Itemset Mining Implementations Repository)

TBI-CPU vs GOETHALS

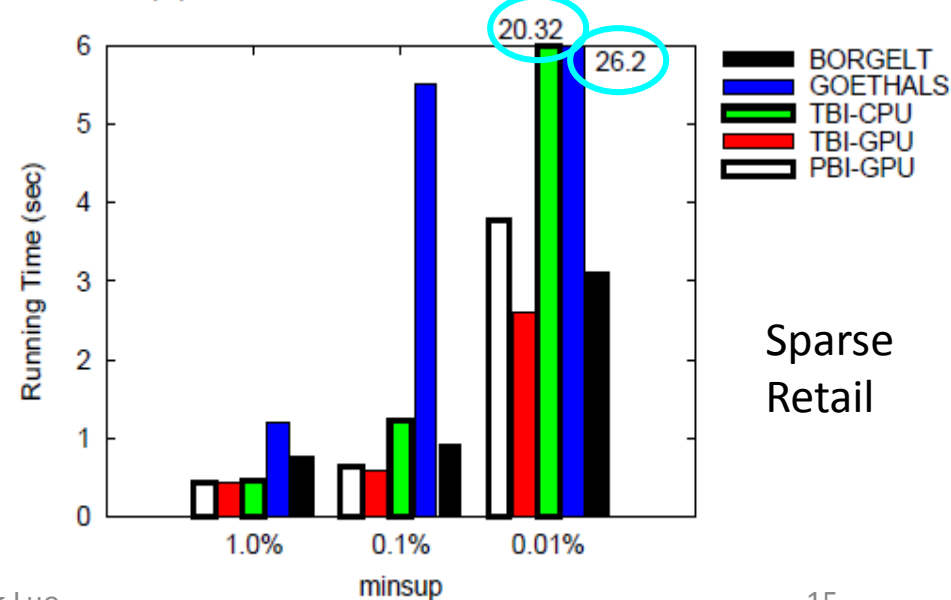


(a) Running time with various minsup



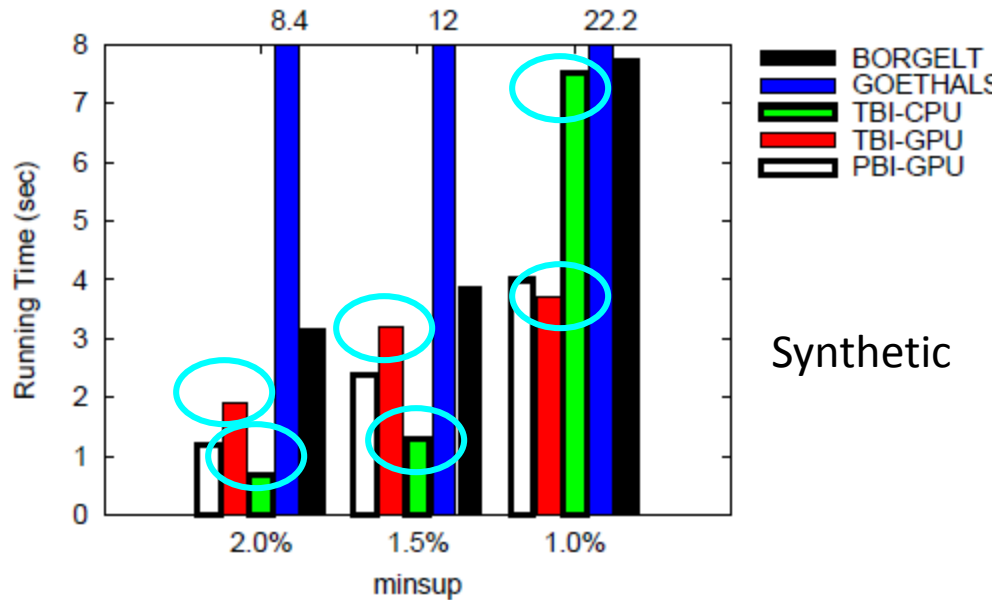
(a) Running time with various minsup

Impl.	Itemset/Candidate Generation	Transactions / Support Counting
TBI-CPU	Trie / CPU	Bitmap / CPU
GOETHALS	Trie / CPU	Horizontal layout / CPU

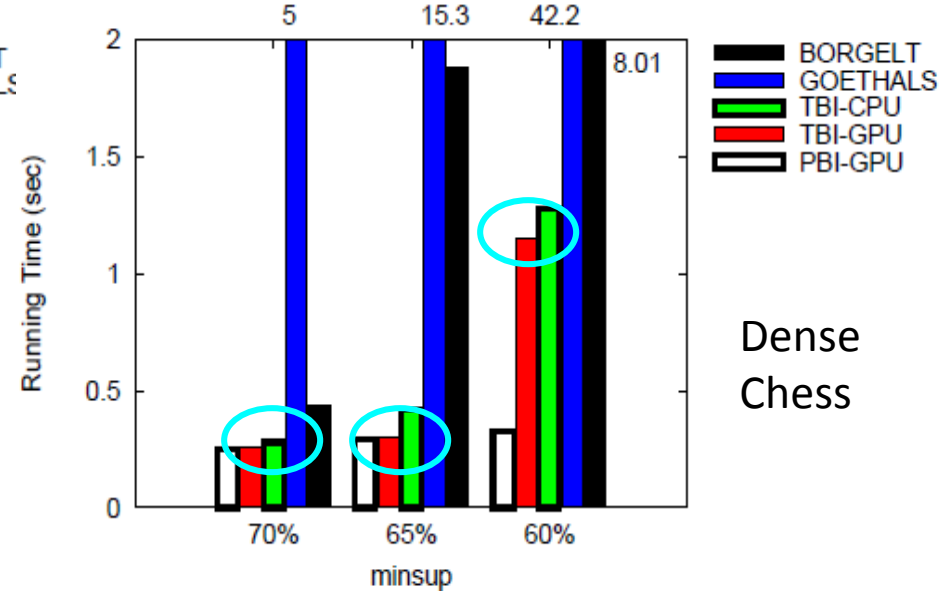


(a) Running time with various minsup

TBI-GPU vs TBI-CPU

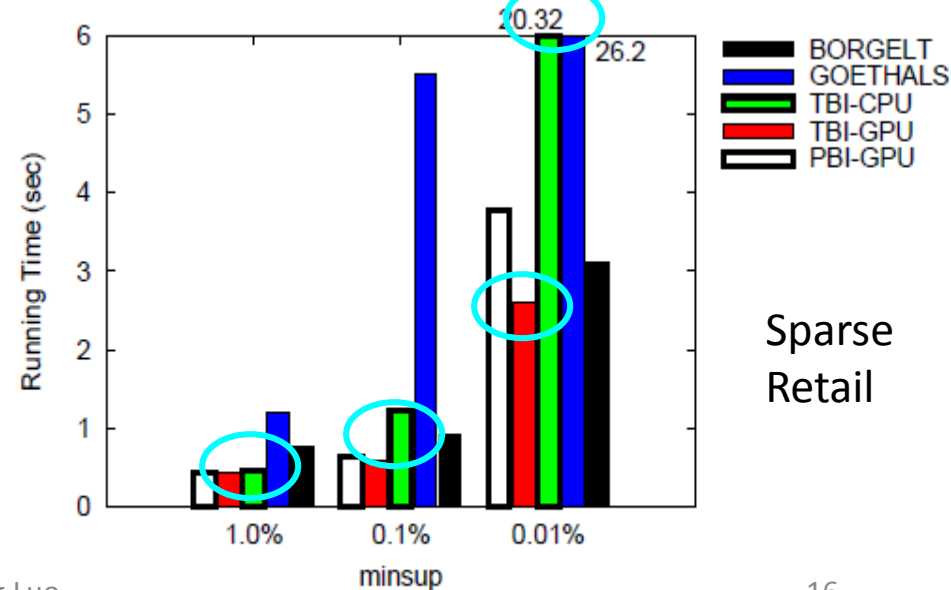


(a) Running time with various minsup



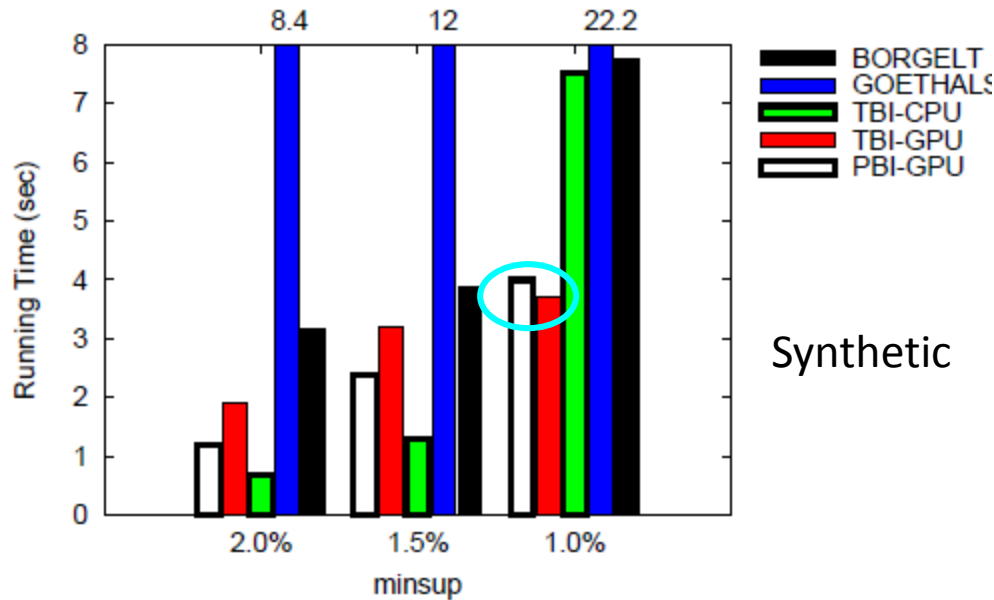
(a) Running time with various minsup

Impl.	Itemset/Candidate Generation	Transactions / Support Counting
TBI-GPU	Trie / CPU	Bitmap / GPU
TBI-CPU	Trie / CPU	Bitmap / CPU

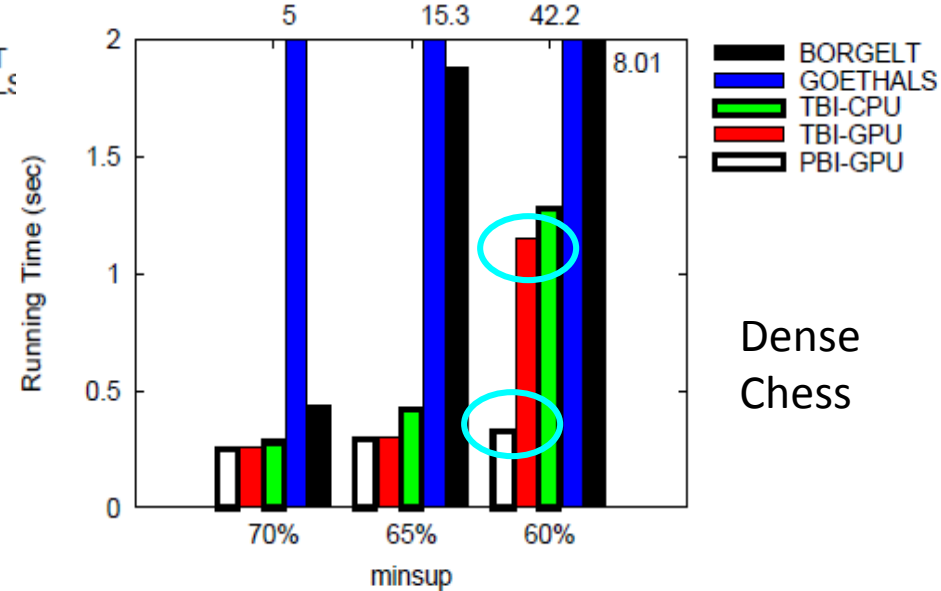


(a) Running time with various minsup

PBI-GPU vs TBI-GPU

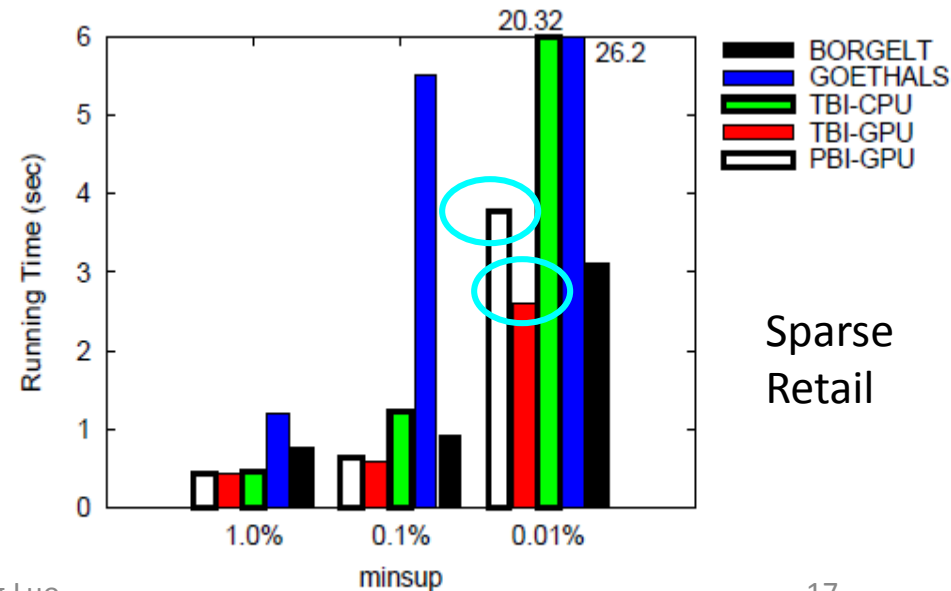


(a) Running time with various minsup



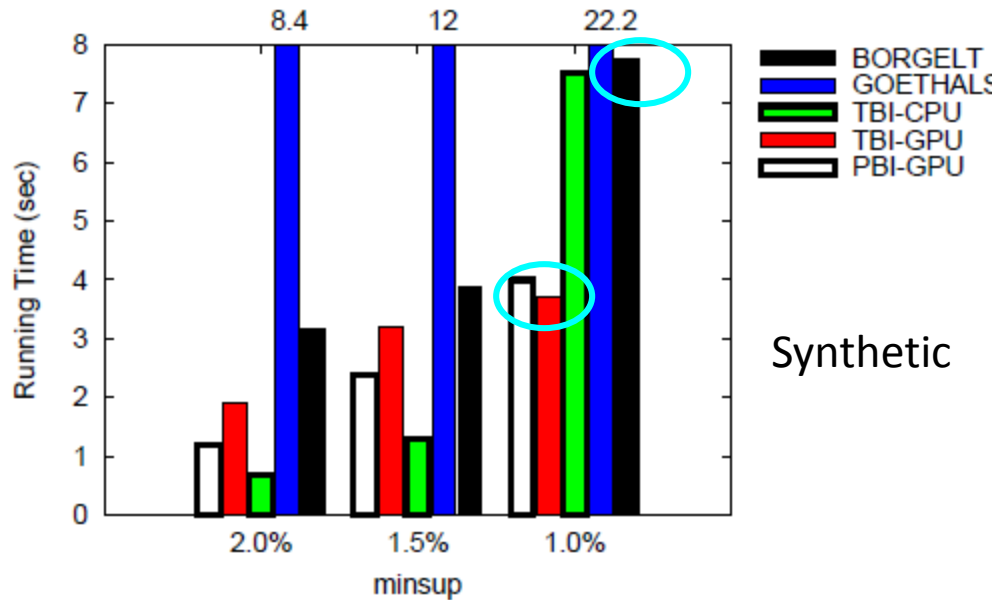
(a) Running time with various minsup

Impl.	Itemset/Candidate Generation	Transactions / Support Counting
PBI-GPU	Bitmap / GPU	Bitmap / GPU
TBI-GPU	Trie / CPU	Bitmap / GPU

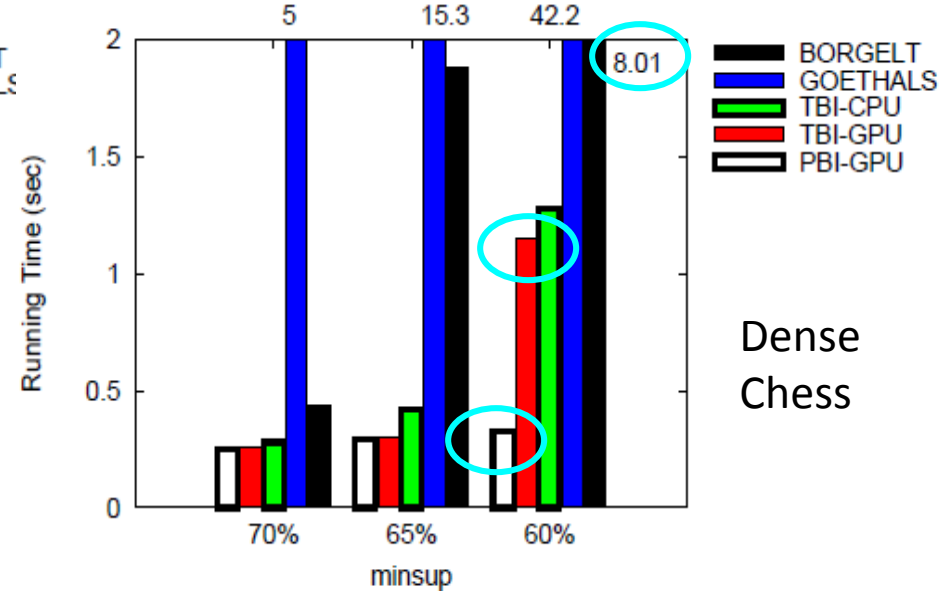


(a) Running time with various minsup

PBI-GPU/TBI-CPU vs BORGELT

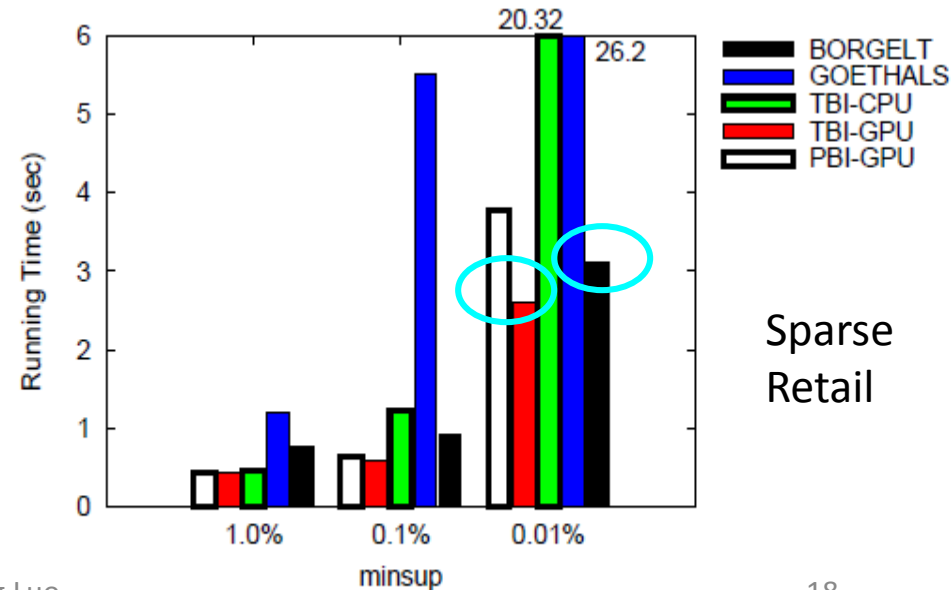


(a) Running time with various minsup



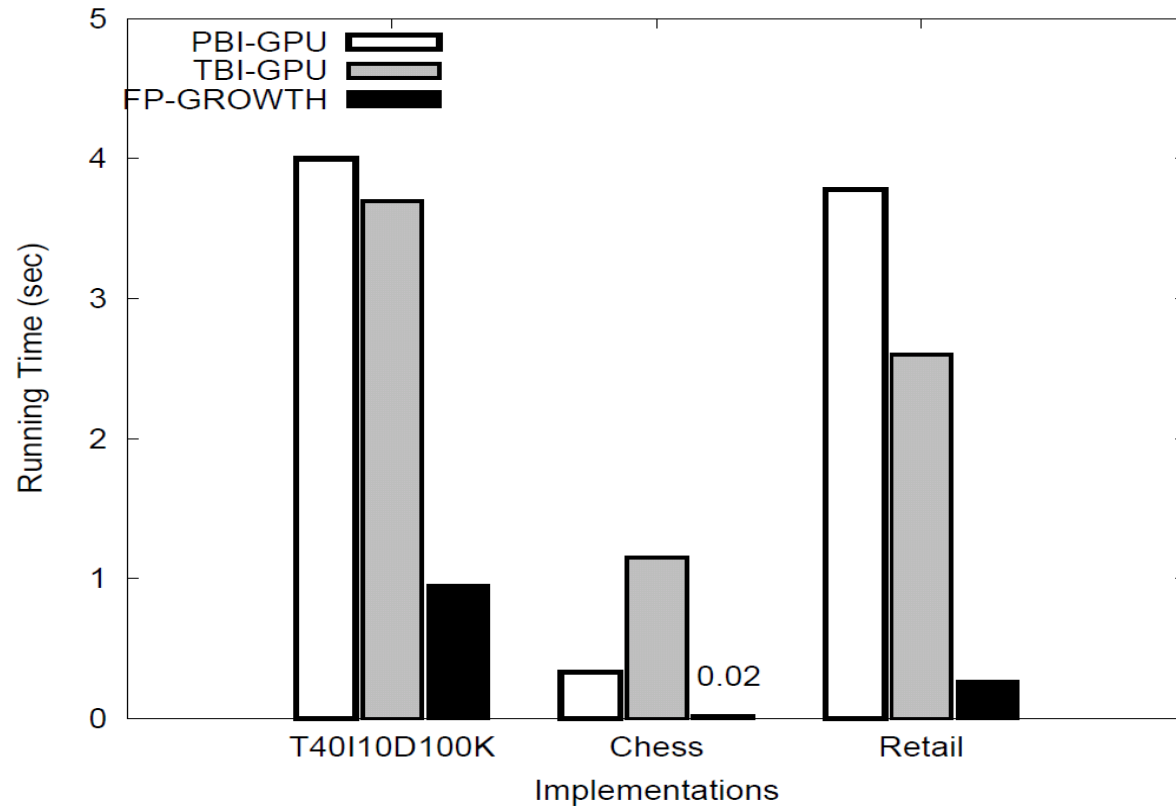
(a) Running time with various minsup

Impl.	Itemset/Candidate Generation	Transactions / Support Counting
PBI-GPU	Bitmap / GPU	Bitmap / GPU
TBI-GPU	Trie / CPU	Bitmap / GPU
BORGELT	Trie / CPU	Trie / CPU



(a) Running time with various minsup

Comparison with FP-growth



Summary

- GPU-based Apriori
 - Pure bitmap-based impl.
 - Bitmap representation for itemsets
 - Bitmap representation for transactions
 - GPU processing
 - Trie-based impl.
 - Trie Representation for itemsets
 - Bitmap Representation for transactions
 - GPU + CPU co-processing
- Better than CPU-based Apriori
- Still worse than CPU-based FP-growth