

A Parallel Pattern Mining Algorithm for Multi-Core Architectures

Benjamin Negrevergne

Jury:

- Marie-Christine Rousset
LIG/Grenoble University
- Alexandre Termier
LIG/Grenoble University
- Jean-François Méhaut
LIG/Grenoble University
- Hiroki Arimura
Graduate School of Science and
Technology/Hokkaido University
- Bruno Crémilleux
GREYC/University of Caen
- Anne Laurent
LIRMM/University of Montpellier

Pattern mining is a sub-topic of **Data Mining**
Knowledge discovery in databases [Fayyad, 96]

"Identifying valid, novel, potentially useful, and ultimately understandable patterns in data"

Pattern mining:

- Extract knowledge as patterns representing regularities (or irregularities) in data.

■ Mining frequent set of items purchased together

▷ _____

■ Mining frequent sub molecules in a molecular database

▷ _____

■ Mining correlations in sensors data

▷ _____

■ Mining frequent set of items purchased together



Possible frequent pattern: {*milk*, *cereals*}

- Dataset: set of receipts
- Patterns: set of items

■ Mining frequent sub molecules in a molecular database



■ Mining correlations in sensors data



■ Mining frequent set of items purchased together

▷ _____

■ Mining frequent sub molecules in a molecular database

▷ _____

■ Mining correlations in sensors data

▷ _____

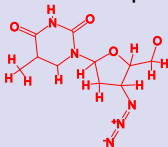
■ Mining frequent set of items purchased together



■ Mining frequent sub molecules in a molecular database



Krammer [2001] was able to identify the following molecular pattern in a set of anti-HIV molecules



- Dataset: set of molecules (represented as graphs)
- Patterns: graphs

■ Mining frequent set of items purchased together

▷ _____

■ Mining frequent sub molecules in a molecular database

▷ _____

■ Mining correlations in sensors data

▷ _____

■ Mining frequent set of items purchased together



■ Mining frequent sub molecules in a molecular database



■ Mining correlations in sensors data



Records from meteorological sensors

Frequent pattern:

Temperature ↓ Pressure ↓ Windspeed ↑

- Dataset: list of numerical sensor records
- Patterns: set of variations



■ Mining frequent set of items purchased together

▷ _____

■ Mining frequent sub molecules in a molecular database

▷ _____

■ Mining correlations in sensors data

▷ _____

The broad range of solutions

Nowadays: many different algorithms for each pattern mining problem

- **Frequent itemset mining**

Apriori	Eclat	FP-Growth	DCI-Closed	LCM
[Agrawal 94]	[Zaki 1997]	[Han 2004]	[Lucchese 2004]	[Uno 2004]

- **Graph Mining**

AGM	gSpan	FSG	CloseGraph	Gaston
[Inokuchi 2000]	[Yan 2002]	[Kuramochi 2001]	[Yan 2003]	[Nijssen 2004]

- **Tree Mining**

FreqT	TreeMiner	Dryade	CMTreeMiner
[Asai 2002]	[Zaki 2005]	[Termier 2004]	[Chi 2004]

- **Gradual itemset mining**

Grite	PGP-mc	PGLCM
[Di Jorio 2009]	[Laurent 2010]	[Do 2010]

- **Sequence**

Spade	PrefixSpan	CloSpan	LCM_seq
[Zaki 2000]	[Pei 2000]	[Yan 2003]	[Uno]

Challenge #1: Genericity

This algorithmic diversity

1. data owners refrain using pattern mining techniques
2. use inadequate algorithms

Challenge #1:

Use a **generic** approach to address different pattern mining problems with a unique algorithm.

Challenge #2: Build an efficient and generic pattern mining algorithm

Pattern mining is inherently **combinatorial**:

▶ very large number of possible patterns

Example

basket market analysis: 1000 items $\rightarrow 2^{1000}$ ($\sim 10^{300}$) possible patterns.

To be efficient algorithms exploit problem properties to reduce the search space ▶ efficient but non-generic

Challenge #2

Design an **efficient** generic and pattern mining algorithm

We proposed ParaMiner: both **generic** and **efficient** algorithm

How:

- Extend theoretical work on pattern enumeration [Boley 2007] and [Arimura 2009]
- Tackle large real world datasets through *dataset reduction*
- Exploit parallelism for multi-core architectures

We proposed ParaMiner: both **generic** and **efficient** algorithm

How:

- Extend theoretical work on pattern enumeration [Boley 2007] and [Arimura 2009]
- Tackle large real world datasets through *dataset reduction*
- Exploit parallelism for multi-core architectures

■ Results

ParaMiner:

- solves various pattern mining problems
- is time-efficient (compete with ad-hoc algorithms)

Generic framework and problem statement

ParaMiner

- Efficient exploration of the set of candidate patterns

- Speeding up candidate pattern testing

- Parallel execution of ParaMiner

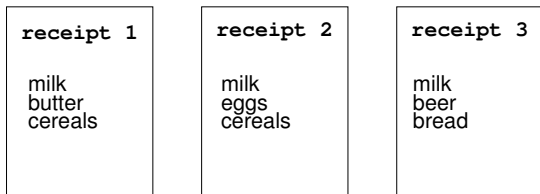
Experiments

- Parallel performance evaluation

- Comparative experiments

Conclusion and future work

Illustration: frequent itemset mining



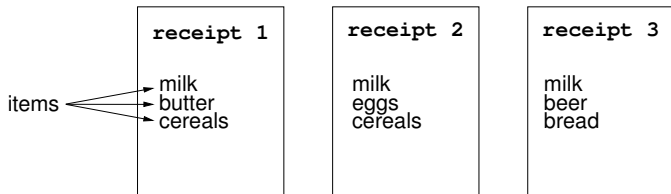
minimum frequency threshold = 50%

Illustration: frequent itemset mining

transactions: t_1

t_2

t_3



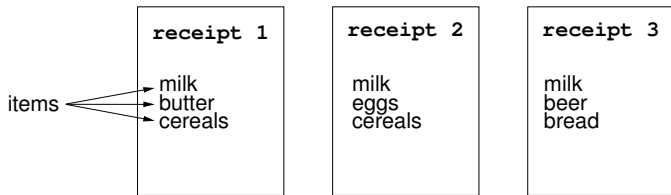
minimum frequency threshold = 50%

Illustration: frequent itemset mining

transactions: t_1

t_2

t_3



minimum frequency threshold = 50%

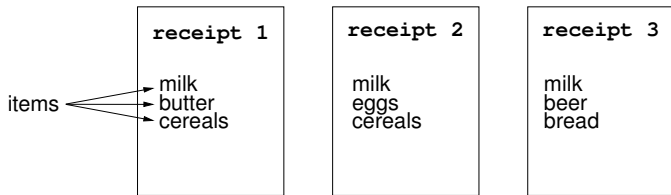
- $P_1 = \{\text{milk}\}$ is frequent (100%): t_1, t_2, t_3

Illustration: frequent itemset mining

transactions: t_1

t_2

t_3



minimum frequency threshold = 50%

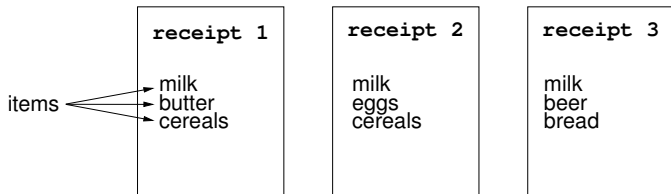
- $P_1 = \{\text{milk}\}$ is frequent (100%): t_1, t_2, t_3
- $P_2 = \{\text{cereals}\}$ is frequent (66%): t_1, t_2

Illustration: frequent itemset mining

transactions: t_1

t_2

t_3



minimum frequency threshold = 50%

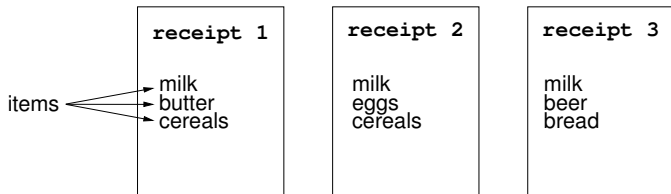
- $P_1 = \{\text{milk}\}$ is frequent (100%): t_1, t_2, t_3
- $P_2 = \{\text{cereals}\}$ is frequent (66%): t_1, t_2
- $P_3 = \{\text{milk, cereals}\}$ is frequent (66%): t_1, t_2

Illustration: frequent itemset mining

transactions: t_1

t_2

t_3



minimum frequency threshold = 50%

- $P_1 = \{\text{milk}\}$ is frequent (100%): $t_1, t_2, t_3 \rightarrow$ closed
- $P_2 = \{\text{cereals}\}$ is frequent (66%): t_1, t_2
- $P_3 = \{\text{milk}, \text{cereals}\}$ is frequent (66%): $t_1, t_2 \rightarrow$ closed

Closed pattern: maximal pattern in a set of transactions

- ▶ Every pattern represented as a set

- ▶ A pattern mining problem defined
 - A ground set E
 - A dataset \mathcal{D}_E
 - A selection criterion *Select*

■ Definition

- Set of all possible elements
- Every candidate pattern is a **subset** of the ground set

■ Definition

- Set of all possible elements
- Every candidate pattern is a **subset** of the ground set

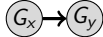
■ Examples

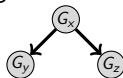
Frequent itemsets ► E : all possible items

$E = \{milk, cereals, eggs, \dots\}$ e.g. candidate: $\{milk, eggs\}$

Connected relational graphs ► E : all the possible arcs

$E = \{(G_1, G_2), (G_1, G_3), \dots, (G_5, G_4)\}$ e.g. candidate:

(G_x, G_y) represent the arc 



Gradual itemsets ► E : all the possible variations

$E = \{T \uparrow, T \downarrow, P \uparrow, \dots, W \downarrow\}$ e.g. candidate: $\{T \uparrow, W \uparrow\}$

■ Definition

- Sequence of *transactions* \mathcal{D}_E
- Each **transaction** is a subset of E

■ Definition

- Sequence of *transactions* \mathcal{D}_E
- Each **transaction** is a subset of E

■ Examples

Frequent itemsets ► in \mathcal{D}_E **each transaction = receipt**

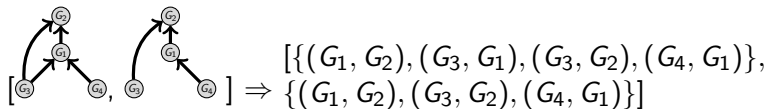
■ Definition

- Sequence of *transactions* \mathcal{D}_E
- Each **transaction** is a subset of E

■ Examples

Frequent itemsets ► in \mathcal{D}_E each transaction = receipt

Relational graphs ► in \mathcal{D}_E each transaction = input graph



Gradual itemsets ► in \mathcal{D}_E each transaction = variation of attributes between pairs of records

Date	T	P	W
Day 1	17.6	1021.20	57
Day 2	18.5	1021.30	56
Day 3	20.4	1018.20	51

⇒

(d_1, d_2)	$\{T \uparrow, P \uparrow, W \downarrow\}$
(d_1, d_3)	$\{T \uparrow, P \downarrow, W \downarrow\}$
(d_2, d_1)	$\{T \downarrow, P \downarrow, W \uparrow\}$
(d_2, d_3)	$\{T \uparrow, P \downarrow, W \downarrow\}$
(d_3, d_1)	$\{T \downarrow, P \uparrow, W \uparrow\}$
(d_3, d_2)	$\{T \downarrow, P \uparrow, W \uparrow\}$

■ Definition

- User-defined predicate $2^E \rightarrow \{true, false\}$
- $Select(P, \mathcal{D}_E) \equiv P$ is a **pattern** of interest in \mathcal{D}_E

■ Definition

- User-defined predicate $2^E \rightarrow \{true, false\}$
- $Select(P, \mathcal{D}_E) \equiv P$ is a **pattern** of interest in \mathcal{D}_E

■ Examples

Frequent itemsets ▶ $Select(P, \mathcal{D}_E) \equiv P$ is frequent in \mathcal{D}_E

Connected relational graphs ▶ $Select(P, \mathcal{D}_E) \equiv P$ is a connected graph **and** P is frequent in \mathcal{D}_E

Gradual itemsets ▶ $Select(P, \mathcal{D}_E) \equiv$ there exists a path $[(d_{x_1}, d_{x_2}), \dots, (d_{x_{n-1}}, d_{x_n})]$ with $n \geq minsup$

■ Closed patterns

Closed pattern: Any subset P of E such that

- P occurs in \mathcal{D}_E ($\mathcal{D}_E[P] \neq \emptyset$)
- P satisfies the selection criterion
- P is the maximal pattern in $\mathcal{D}_E[P]$

■ Problem statement

Given a **ground set** E , a **dataset** \mathcal{D}_E and a **selection criterion**
▶ extract all the **closed patterns**

Generic framework and problem statement

ParaMiner

- Efficient exploration of the set of candidate patterns

- Speeding up candidate pattern testing

- Parallel execution of ParaMiner

Experiments

- Parallel performance evaluation

- Comparative experiments

Conclusion and future work

■ The standard approach

generate and test.

1. **Generate** candidate patterns
2. **Test** if the candidate pattern is a pattern

■ The standard approach

generate and test.

1. **Generate** candidate patterns
2. **Test** if the candidate pattern is a pattern

■ ... is not naively applicable

- **Too many** candidate patterns
- **Each test** requires costly database accesses (e.g. to count frequency)

► Challenges

- Generic strategy to explore the set of candidate patterns
- Generic method to simplify the testing

Generic framework and problem statement

ParaMiner

Efficient exploration of the set of candidate patterns

Speeding up candidate pattern testing

Parallel execution of ParaMiner

Experiments

Parallel performance evaluation

Comparative experiments

Conclusion and future work

■ Structured search space

Most algorithms define a **pattern augmentation relation**

- ▶ search space is explored by repeatedly augmenting patterns

ParaMiner's augmentation relation

P and Q , two patterns:

- ▶ Q is the augmentation of $P \Leftrightarrow Q = P \cup \{e\}$

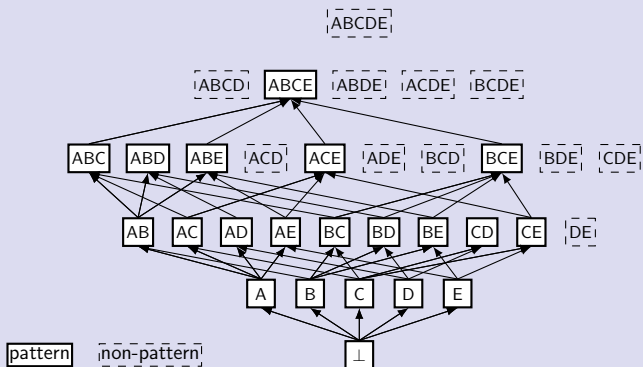
ACE \rightarrow ABCE

Set of patterns + **Augmentation relation** = **DAG¹ structure**

¹*directed acyclic graph*

Structured search space

DAG-structure of a set of patterns



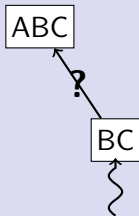
- several arcs lead to the same pattern

■ An enumeration strategy

Is required to:

- discover all the patterns
- avoid **duplicates generation**

Which augmentation path must be followed ?

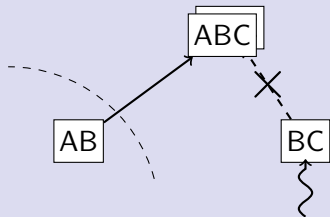


■ An enumeration strategy

Is required to:

- discover all the patterns
- avoid **duplicates generation**

Which augmentation path must be followed ?



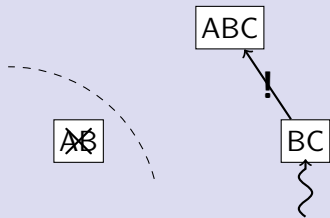
▶ **AB** is the **first_parent** of **ABC**

■ An enumeration strategy

Is required to:

- discover all the patterns
- avoid **duplicates generation**

Which augmentation path must be followed ?



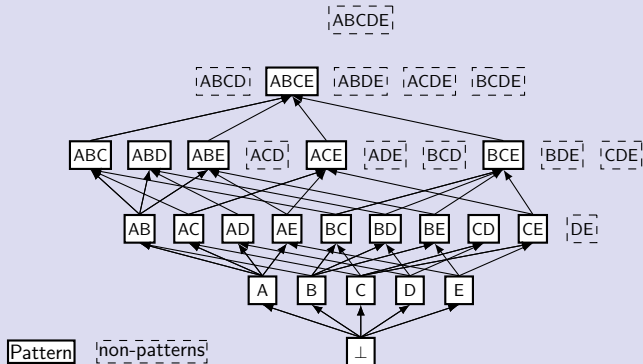
► **BC** is the **first_parent** of **ABC**

Enumeration strategy

Follow the augmentation $\boxed{P} \rightarrow \boxed{Q}$ if:

$\text{first_parent}(\boxed{Q}) = \boxed{P}$

DAG-structure following a tree search

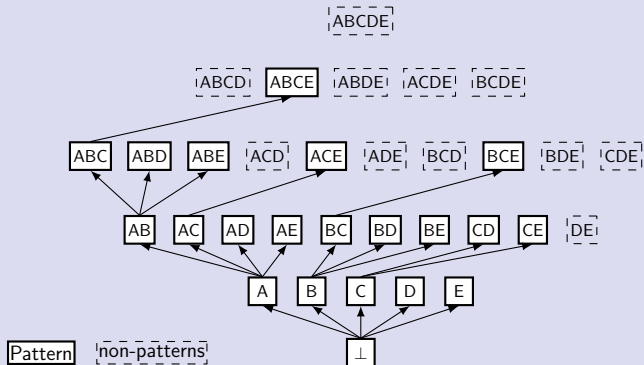


Enumeration strategy

Follow the augmentation $\boxed{P} \rightarrow \boxed{Q}$ if:

first_parent(\boxed{Q}) = \boxed{P}

DAG-structure following a tree search



■ Requirement

A method that is:

- **Generic**
 - ▶ Sound for all our pattern mining problems
- **Adequate to parallel exploration of the search space**
 - ▶ Computable without **global computations**

■ Requirement

A method that is:

- **Generic**
 - ▶ Sound for all our pattern mining problems
- **Adequate to parallel exploration of the search space**
 - ▶ Computable without **global computations**

■ State of the art:

Theoretical work closed pattern enumeration strategies

- [Arimura et al. 2009]
- [Boley et al. 2010]
- ▶ Poly-space enumeration strategies
 - ▶ `first_parent` do not rely on a global representation of the search space

Structural properties of the set of patterns

[Boley 2010] **Strong accessibility:**

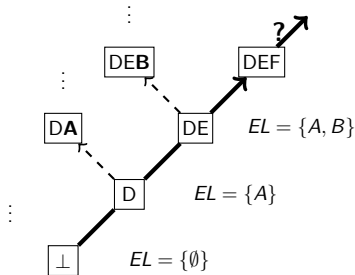
Idea: *We can reach every pattern by augmenting any subset that is a pattern*

Structural properties of the set of patterns

[Boley 2010] **Strong accessibility:**

Idea: *We can reach every pattern by augmenting any subset that is a pattern*

We can detect first parent by memorizing **only the root** of the forking branches:



- Branch can be explored independently → Parallel
- FIM, CRG and GRI satisfy this property (proved) → Generic

Conclusion on pattern enumeration

- Pattern augmentation
 - Structured search space
 - Enumeration strategy required
 - Problem of detecting the first parent
 - Strong accessibility property
 - Exclusion list
- ▶ generic and parallelizable enumeration strategy

Next: Identifying patterns

Generic framework and problem statement

ParaMiner

Efficient exploration of the set of candidate patterns

Speeding up candidate pattern testing

Parallel execution of ParaMiner

Experiments

Parallel performance evaluation

Comparative experiments

Conclusion and future work

■ Candidate pattern testing

- Check candidate pattern occurrence in the dataset
 - Computing the selection criterion
 - Computing pattern closure
- ▶ Intensive access to the dataset

Problem: How to reduce dataset accesses ?

- ▶ dataset reduction [Han 2000, Uno 2004]

Motivation: Every element in the dataset are not required to test each candidate pattern

Principle:

1. build a dataset **for each** new pattern P
2. **filter** unnecessary elements to P and its descendants
3. use it to test P and its descendants

In ParaMiner : EL-based filter ► generic

EL-based reduction (Proved)

Dataset:

\mathcal{D}_E

Pattern:

$P = \{D, E\}$

EL:

$EL = \{A, B, C\}$

► $\mathcal{D}_P^{\text{reduced}} ?$

\mathcal{D}_E				
t_1	A, B,	D, E,	F, G,	
t_2	A, D,	D, E,	F, G,	
t_3	A, B,	D, E,	F, H,	

EL-based reduction (Proved)

Dataset:

\mathcal{D}_E

Pattern:

$P = \{D, E\}$

EL:

$EL = \{A, B, C\}$

\mathcal{D}_E	$\in EL$	P	augm.	($\notin EL$)
t_1	A, B,	D, E,	F, G,	
t_2	A, D,	D, E,	F, G,	
t_3	A, B,	D, E,	F, H,	

► $\mathcal{D}_P^{reduced}$?

- Element required: Any e that can belongs to a closed pattern including P ► ensure that `first_parent` is sound

EL-based reduction (Proved)

Dataset:

\mathcal{D}_E

Pattern:

$P = \{D, E\}$

EL:

$EL = \{A, B, C\}$

\mathcal{D}_E	$\in EL$	P	augm.	($\notin EL$)
t_1	A, B,	D, E,	F, G,	
t_2	A, D,	D, E,	F, G,	
t_3	A, B,	D, E,	F, H,	

► $\mathcal{D}_P^{reduced}$?

- Element required: Any e that can belongs to a closed pattern including P ► ensure that `first_parent` is sound

EL-based reduction (Proved)

Dataset:

\mathcal{D}_E

Pattern:

$P = \{D, E\}$

EL:

$EL = \{A, B, C\}$

\mathcal{D}_E	$\in EL$	P	augm.	($\notin EL$)
t_1	A, B,	D, E,	F, G,	} \mathcal{P}_1
t_2	A, D,	D, E,	F, G,	
t_3	A, B,	D, E,	F, H,	} \mathcal{P}_2

► $\mathcal{D}_P^{reduced}$?

- Element required: Any e that can belongs to a closed pattern including P ► ensure that `first_parent` is sound

How to detect them:

- grouping transactions per augmentations supported
- Apply reduction

EL-based reduction (Proved)

Dataset:

\mathcal{D}_E

Pattern:

$P = \{D, E\}$

EL:

$EL = \{A, B, C\}$

\mathcal{D}_E	$\in EL$	P	augm.	($\notin EL$)
t_1	A, B,	D, E,	F, G,	} \mathcal{P}_1
t_2	A, D,	D, E,	F, G,	
t_3	A, B,	D, E,	F, H,	} \mathcal{P}_2
		↓		
$\mathcal{D}_P^{reduced}$	\cap	=	=	
$t'_1 \oplus t'_2$	A,	D, E,	F, G,	
t'_3	A, B,	D, E,	F, H,	

► $\mathcal{D}_P^{reduced}$?

- Element required: Any e that can belong to a closed pattern including P ► ensure that first_parent is sound

How to detect them:

- grouping transactions per augmentations supported
- Apply reduction

■ Metric

$$reduction_factor = \frac{|\mathcal{D}_E|}{|\mathcal{D}_P^{reduced}|}$$

$$average_reduction_factor = \frac{\sum_{P \in \mathcal{C}} |\mathcal{D}_E| / |\mathcal{D}_P^{reduced}|}{|\mathcal{C}|}$$

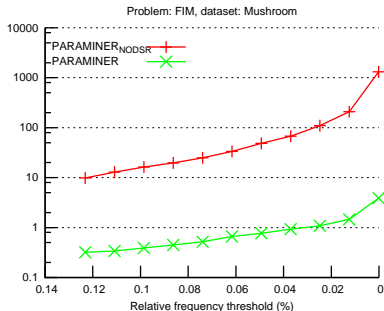
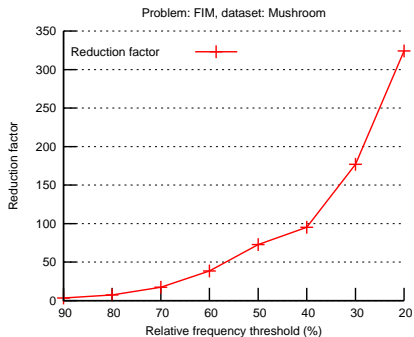
■ Frequent itemset mining, mushroom dataset

dataset name	ground set size	# transactions	dataset size
Mushroom	119	8,124	186,852

■ Metric

$$\text{reduction_factor} = \frac{|\mathcal{D}_E|}{|\mathcal{D}_P^{\text{reduced}}|}$$

$$\text{average_reduction_factor} = \frac{\sum_{P \in \mathcal{C}} |\mathcal{D}_E| / |\mathcal{D}_P^{\text{reduced}}|}{|\mathcal{C}|}$$



ParaMiner: algorithm

```
1: procedure expand( $P, \mathcal{D}_P^{\text{reduced}}, EL$ )
2:   for all  $e$  such that  $e$  occurs in  $\mathcal{D}_P^{\text{reduced}}$  do
3:     if  $\text{Select}(P \cup \{e\}, \mathcal{D}_P^{\text{reduced}})$  then
4:        $Q \leftarrow \text{Clo}(P \cup \{e\}, \mathcal{D}_P^{\text{reduced}})$ 
5:       if  $\text{is\_first\_parent}(P, EL, Q)$  then
6:         output  $Q$ 
7:          $\mathcal{D}_Q^{\text{reduced}} \leftarrow \text{reduce}(\mathcal{D}_P^{\text{reduced}}, e, EL)$ 
8:         spawn  $\text{expand}(Q, \mathcal{D}_Q^{\text{reduced}}, EL)$ 
9:          $EL \leftarrow EL \cup \{e\}$ 
10:      end if
11:    end if
12:  end for
13: end procedure
```

- **Proved sound and complete** under some constraints:
 - ▶ the set of patterns satisfying the **strong accessibility** property
 - ▶ the selection criterion is **decomposable**

Problems FIM, CRG and GRI satisfy these properties

- **Parallelized**: space exploration is divided into independent tasks
 - ▶ non-trivial for **closed** pattern mining

- **Proved sound and complete** under some constraints:
 - ▶ the set of patterns satisfying the **strong accessibility** property
 - ▶ the selection criterion is **decomposable**

Problems FIM, CRG and GRI satisfy these properties

- **Parallelized**: space exploration is divided into independent tasks
 - ▶ non-trivial for **closed** pattern mining
- Melinda a parallel execution engine for **data-driven** algorithms
 - ▶ can be extended with task scheduling strategies
 - ▶ execute tasks spawned by ParaMiner
 - ▶ also used in PLCM [Negrevergne 2010], PGLCM [Do 2010]

Generic framework and problem statement

ParaMiner

Efficient exploration of the set of candidate patterns

Speeding up candidate pattern testing

Parallel execution of ParaMiner

Experiments

Parallel performance evaluation

Comparative experiments

Conclusion and future work

■ Melinda: Main concepts

- **Tuples** (A, B, C) ▶ block of data with a fixed structure
- **Tuple Space** ▶ synchronized memory space
 - ▶ can contain tuples
 - ▶ supports **put()** and **get()**

■ Parallelizing ParaMiner with Melinda

ParaMiner parallelization scheme: A new tuple for each recursive call

Tuple structure: Parameters of recursive calls
(*Pattern*, *Reduce_dataset*, *Exclusion_list*)

■ Melinda: Main concepts

- **Tuples** (A, B, C) ▶ block of data with a fixed structure
- **Tuple Space** ▶ synchronized memory space
 - ▶ can contain tuples
 - ▶ supports **put()** and **get()**

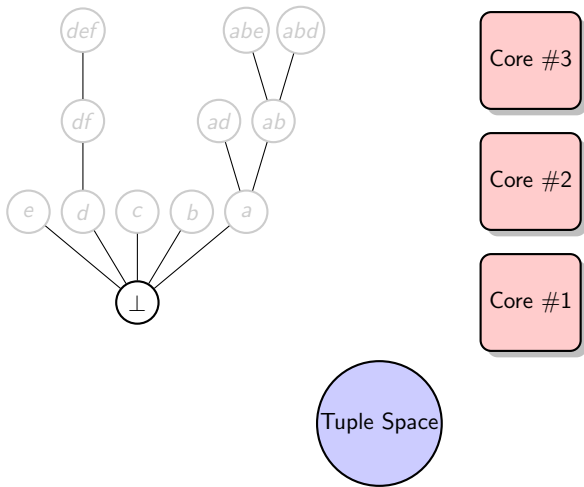
■ Parallelizing ParaMiner with Melinda

ParaMiner parallelization scheme: A new tuple for each recursive call

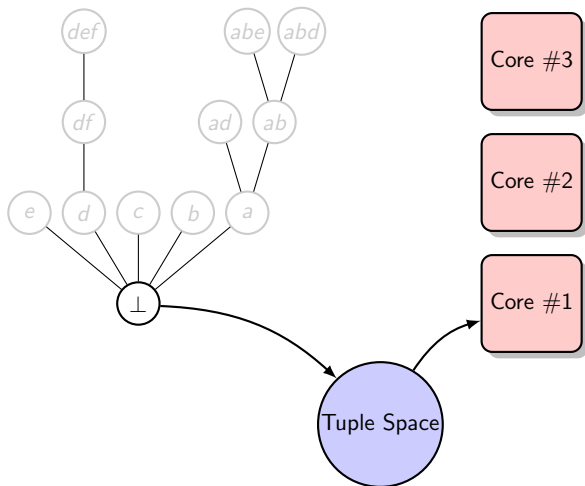
Tuple structure: Parameters of recursive calls
(*Pattern*, *Reduce_dataset*, *Exclusion_list*)

- Tuples are **produced** when a pattern is discovered
- Tuples are **consumed** when a core is idle

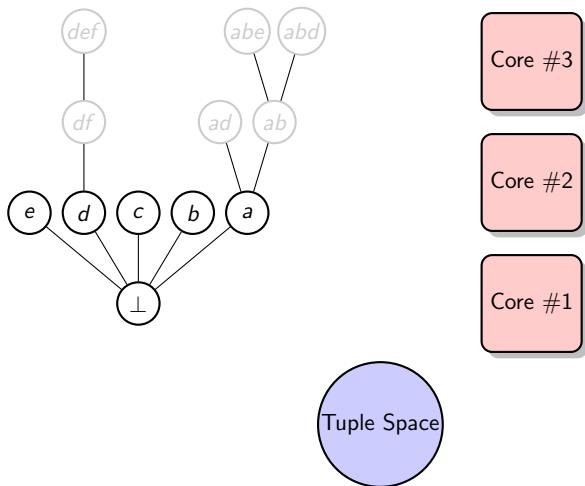
Parallel execution



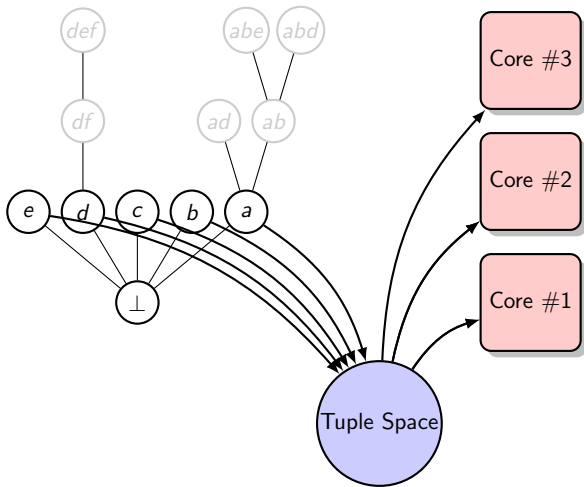
Parallel execution



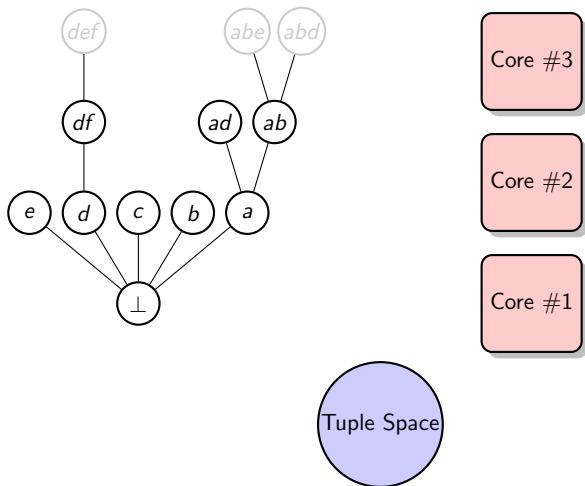
Parallel execution



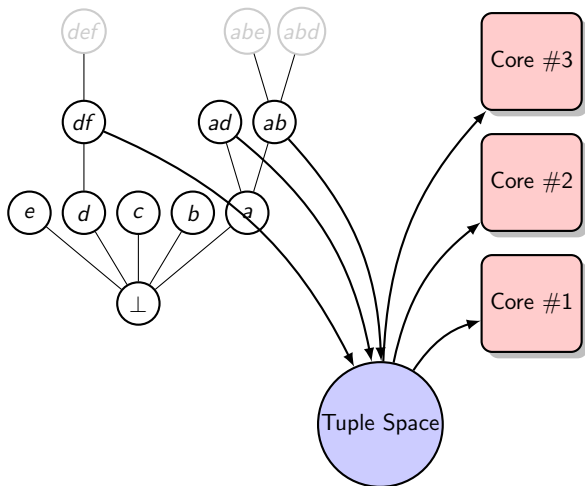
Parallel execution



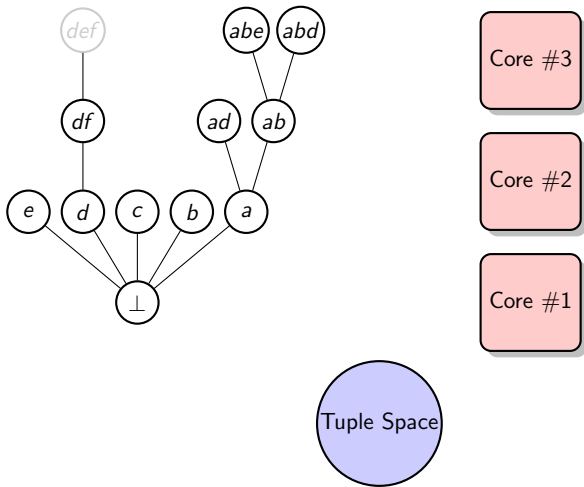
Parallel execution



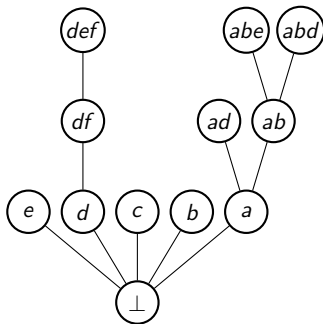
Parallel execution



Parallel execution



Parallel execution



Core #3

Core #2

Core #1

Tuple Space

Generic framework and problem statement

ParaMiner

- Efficient exploration of the set of candidate patterns

- Speeding up candidate pattern testing

- Parallel execution of ParaMiner

Experiments

- Parallel performance evaluation

- Comparative experiments

Conclusion and future work

Generic framework and problem statement

ParaMiner

- Efficient exploration of the set of candidate patterns

- Speeding up candidate pattern testing

- Parallel execution of ParaMiner

Experiments

- Parallel performance evaluation

- Comparative experiments

Conclusion and future work

■ ParaMiner instances

- frequent itemset mining
- frequent relational graph mining
- gradual itemset mining

■ Computing platforms

Laptop:

- 4-cores (4 core *i7*)
- 8 GB memory

Server:

- 32-cores (4 core *i7* with each core each)
- 64 GB memory

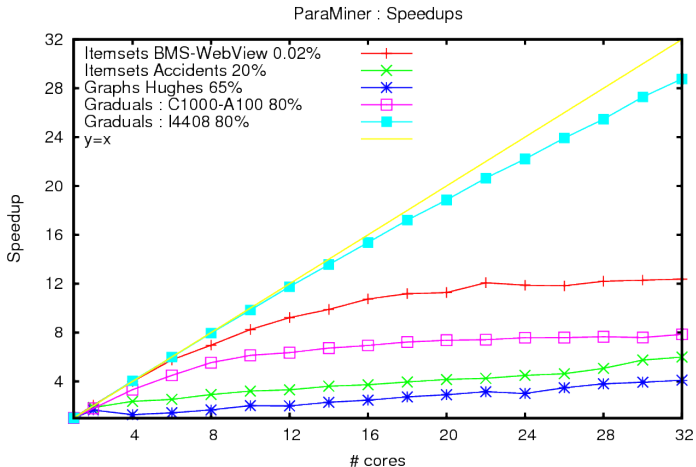
■ Metric

$$speedup = \frac{\text{time using one core}}{\text{time using n cores}}$$

■ 1. Laptop

- ▶ speedup on 4 cores: 3 to 4

2. Server



In [Buerhrer 2006, Tatikonda 2008], authors exhibit two important issues:

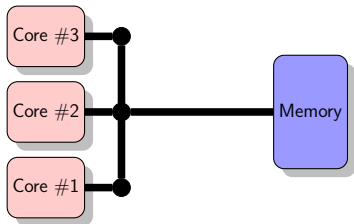
- load imbalance
- memory issues

In [Buerhrer 2006, Tatikonda 2008], authors exhibit two important issues:

- load imbalance
- memory issues

■ Bus contention

Memory bus: connect cores to memory



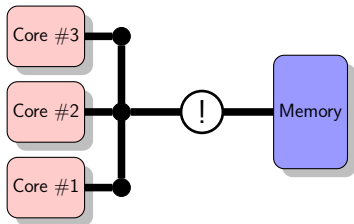
In [Buerhrer 2006, Tatikonda 2008], authors exhibit two important issues:

- load imbalance
- memory issues

■ Bus contention

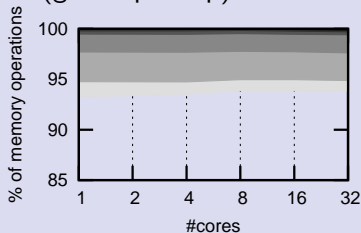
Memory bus: connect cores to memory

too many memory access ► subject to contention

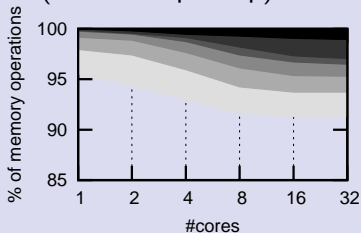


■ % hi-lat. memory operation VS # of cores

gradual itemset mining
(good speedup)



frequent itemset mining
(bounded speedup)



dark slices = high latencies memory operations

- **GRI**: more cores ► % hi-lat. memory operations remains constant
- **FIM**: more cores ► increasing % of hi-lat. memory operations

[Tatikonda 2008] proposed **architecture conscious** pattern mining algorithm

Requires deep modifications in the algorithm

Inadequate to ParaMiner:

▶ could penalize executions that perform well

■ **Melinda's** *tuple distribution strategies*

- distribute tuples to available cores
- user-defined
 - ▶ can exploit **algorithmic level** information (in tuples)
 - ▶ can exploit **architectural level** information (available in Melinda)

■ Structured tuplespace

Tuple space is divided into **internals**

▶ Internals can be used to classify tuples in the tuple space

■ Defining new strategies

- *distribute()* defines in which internal to put the tuple
- *retrieve()* from which internal the tuple has to be taken

▶ supports heavy usage

Architecture conscious tuple distribution strategy

1. **Distribute:** Tuples sharing datastructures go in the same internal
2. **Retrieve:** Cores of the same processor retrieve tuples from the same internal

Architecture conscious tuple distribution strategy

1. **Distribute**: Tuples sharing datastructures go in the same internal
2. **Retrieve**: Cores of the same processor retrieve tuples from the same internal

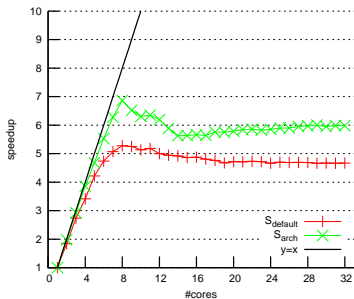
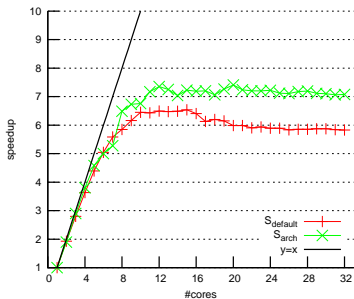
► **Improve cache usage** and **reduce bus contention**

■ **Frequent itemset mining, BMS-WebView/Accidents dataset**

dataset	# items	# transactions	dataset size	Density (%)
BMS-WebView-2	3,340	77,512	320,601	0.14
Accidents	468	340,184	11,500,870	7.22

Architecture conscious tuple distribution strategy

1. **Distribute:** Tuples sharing datastructures go in the same internal
2. **Retrieve:** Cores of the same processor retrieve tuples from the same internal



No source code modification ► 30% performance gain

Generic framework and problem statement

ParaMiner

Efficient exploration of the set of candidate patterns

Speeding up candidate pattern testing

Parallel execution of ParaMiner

Experiments

Parallel performance evaluation

Comparative experiments

Conclusion and future work

■ Algorithms

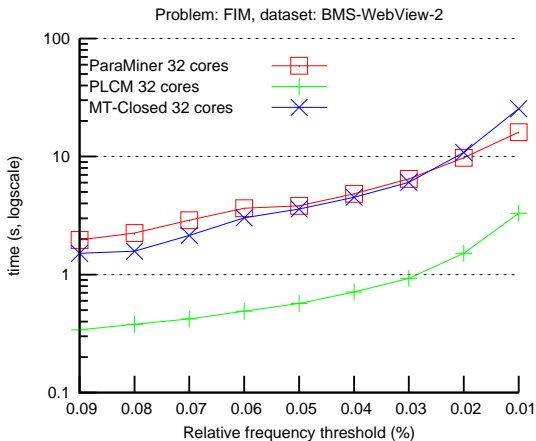
- **PLCM** [Negrevergne 2010] parallel implementation of LCM [Uno 2004] FIMI'04 Award
- **MT-Closed** [Lucchese 2007] parallel implementation of DCI-Closed [Lucchese 2004]

■ Datasets

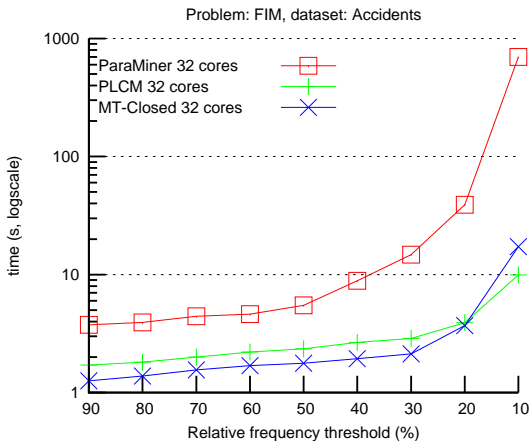
dataset	# items	# transactions	dataset size	Density (%)
BMS-WebView-2	3,340	77,512	320,601	0.14
Accidents	468	340,184	11,500,870	7.22

$$\text{density}(\mathcal{D}_E) = \frac{\#items \times \#transactions}{|\mathcal{D}_E|} (\times 100)$$

■ BMS-WebView-2 (sparse)



■ Accidents (dense)



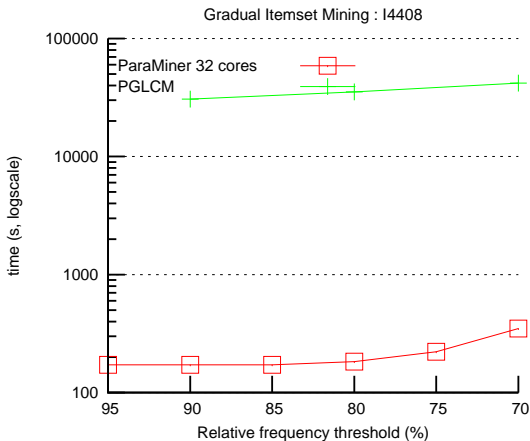
■ Algorithms

- **PGLCM** [Do 2010]
- **PGP-mc** (non-closed) [Laurent 2010]

■ Datasets

dataset name	ground set size	# transactions	dataset size	Density (%)
l4408	8824	11,556	50,985,072	50

■ I4408 (Real)



Generic framework and problem statement

ParaMiner

- Efficient exploration of the set of candidate patterns

- Speeding up candidate pattern testing

- Parallel execution of ParaMiner

Experiments

- Parallel performance evaluation

- Comparative experiments

Conclusion and future work

- **ParaMiner**
 - ▶ based on state of the art on pattern enumeration
 - ▶ useful to parallel pattern mining
 - ▶ designed an efficient generic dataset reduction
 - ▶ exploit multi-core architectures
- **Melinda**
 - ▶ parallel engine of ParaMiner (and other algorithms)
 - ▶ extensible through strategies

- **ParaMiner**

- ▶ based on state of the art on pattern enumeration
 - ▶ useful to parallel pattern mining
- ▶ designed an efficient generic dataset reduction
- ▶ exploit multi-core architectures

- **Melinda**

- ▶ parallel engine of ParaMiner (and other algorithms)
- ▶ extensible through strategies

■ Novelty

- ParaMiner: an efficient solution for new pattern mining problems
New pattern mining problem can benefit from 15 years of research
- ParaMiner/Melinda: tool to learn about parallel pattern mining

■ Extend ParaMiner's genericity

- Generalize strong accessibility properties to other structures
 - ▶ handle sequences, and general graphs
- Study partial strong accessibility of set of patterns
 - ▶ bridge with constraint pattern mining

■ Extend ParaMiner/Melinda's efficiency

- Improve Melinda's strategies
 - ▶ More flexible and expressive querying tuples
- Target larger computing platforms
 - ▶ cluster of computers

Thank you !