

# Mine 'Em All: A Note on (Complexity of) Mining All Graphs

Ondřej Kuželka<sup>1</sup> and Jan Ramon<sup>2</sup>

<sup>1</sup>Cardiff University, <sup>2</sup>KU Leuven

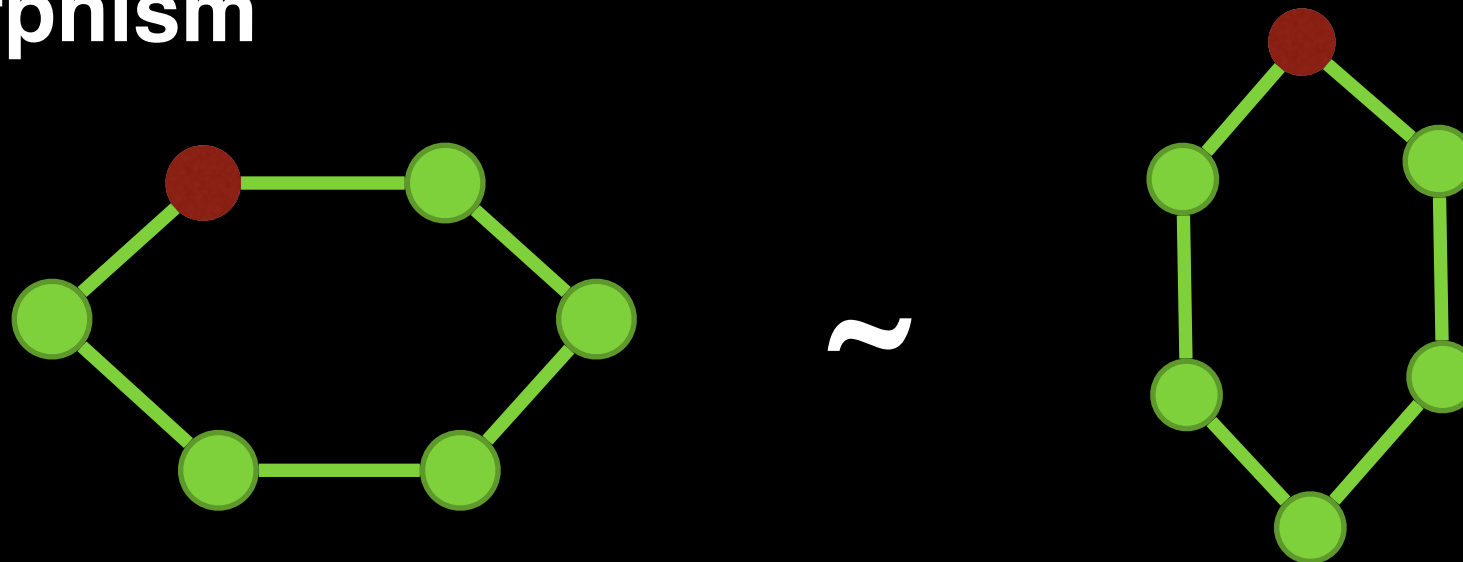


# Question

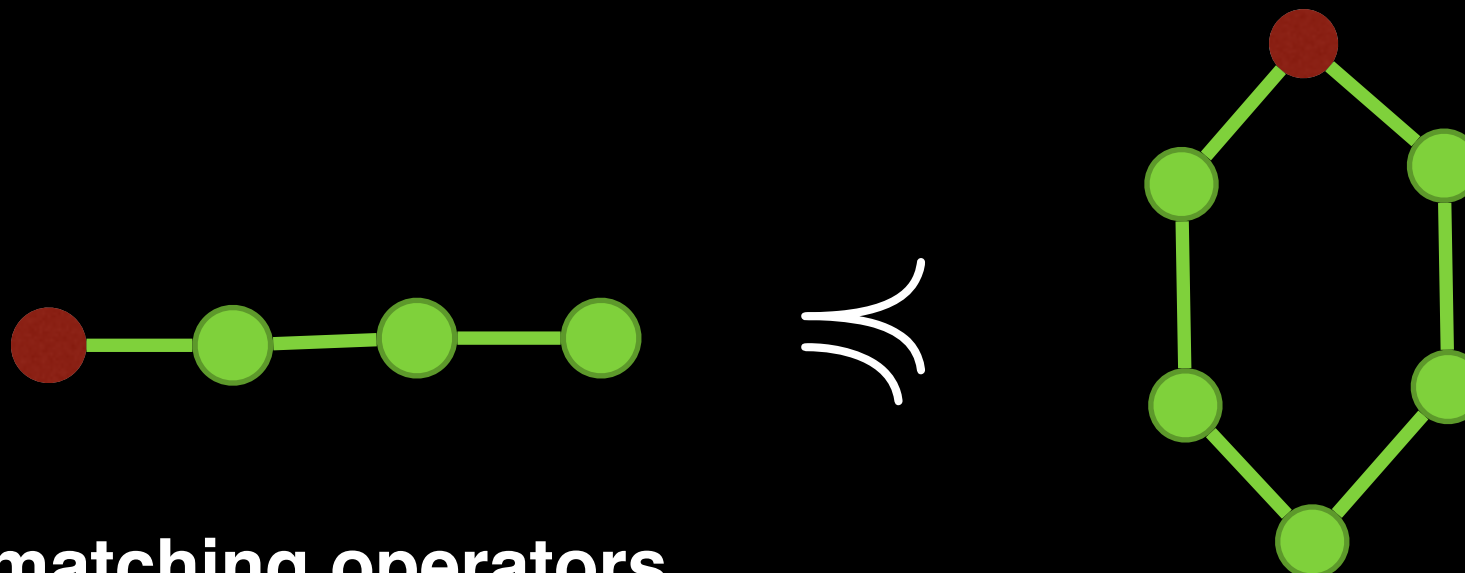
- **When can *graph mining* with an *intractable pattern matching operator* be fast?**
- Motivation: Horváth & Ramon have shown that frequent bounded-treewidth graphs can be mined in incremental-polynomial time even though subgraph isomorphism is NP-hard for them.

# Preliminaries

- **Isomorphism**



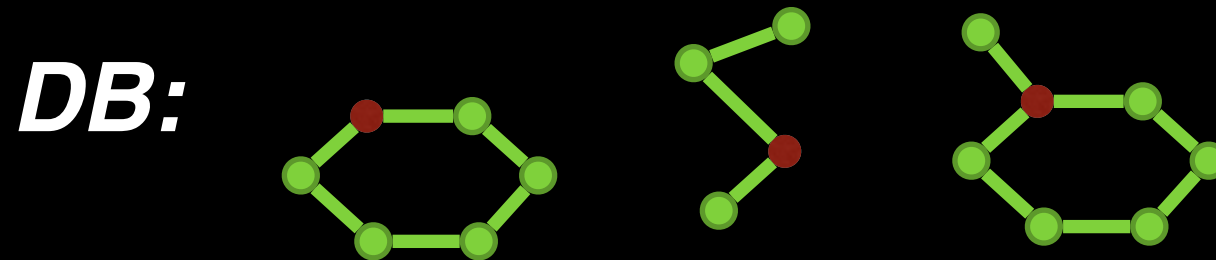
- **Subgraph isomorphism:**



- **+ other matching operators**  
(homeomorphism, minor embedding, induced operators...)

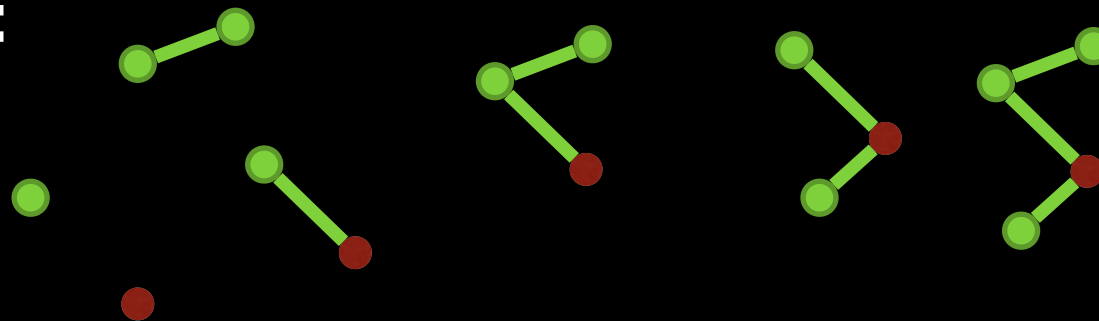
# Frequent Graph Mining

- **Given:** a database  $DB$  of graphs and a frequency threshold  $t$



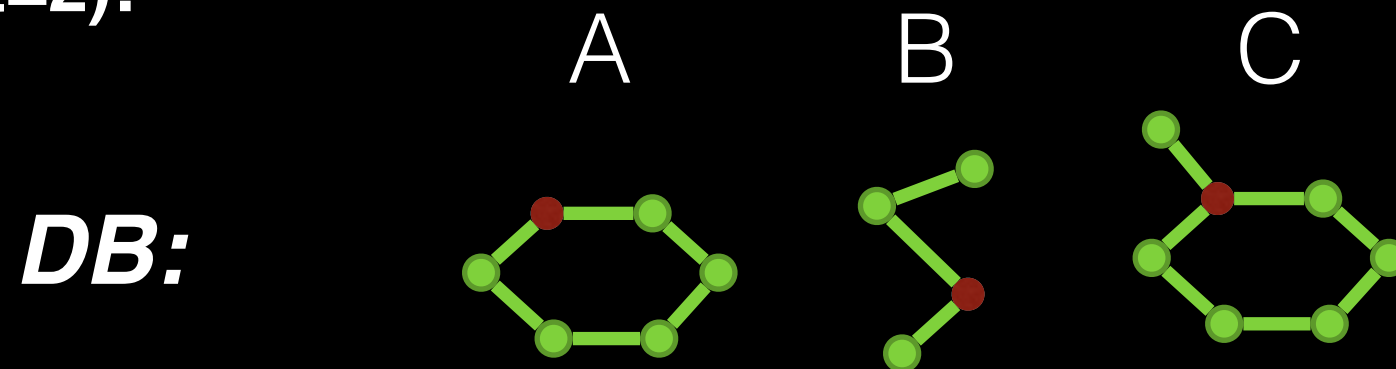
- **Task:** Output all nonisomorphic connected graphs subgraph isomorphic to at least  $t$  graphs from  $DB$ .

**Example ( $t=3$ ):**



# How Typical FGM Algos. Work

**Example (t=2):**



Candidates (1):

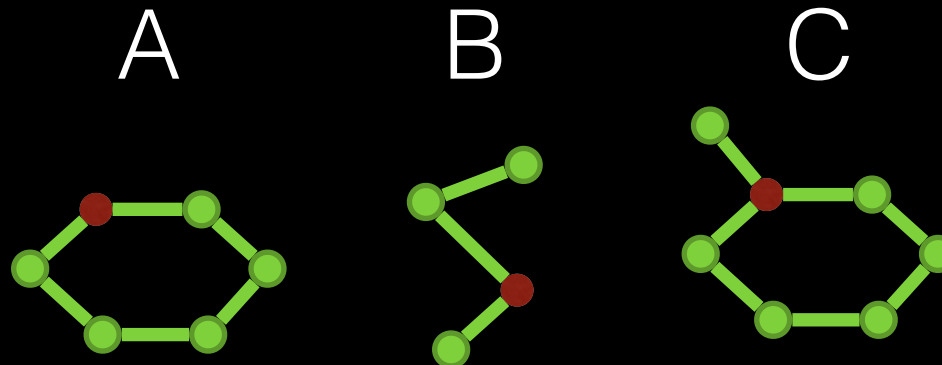


Occurrences: {A,B,C} and {A,B,C}

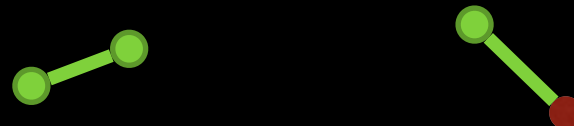
# How Typical FGM Algos. Work

**Example (t=2):**

***DB:***



Candidates (2):



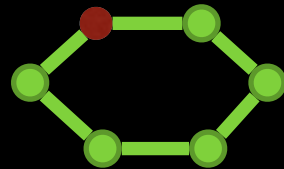
Occurrences: {A,B,C} and {A,B,C}

# How Typical FGM Algos. Work

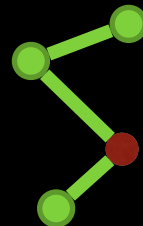
**Example (t=2):**

***DB:***

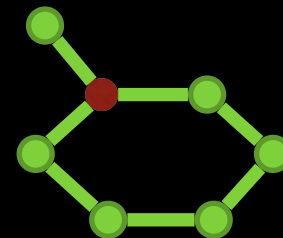
A



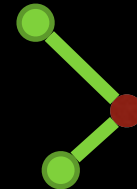
B



C



Candidates (3):



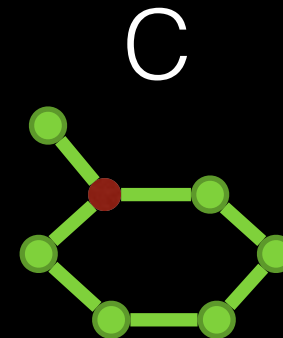
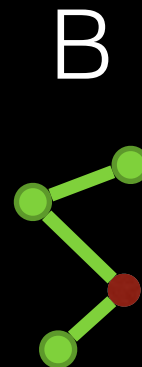
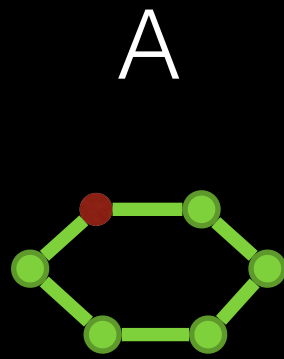
etc...

Occurrences: {A,B,C} and {A,B,C} and {A,C}

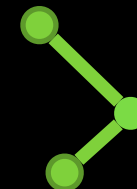
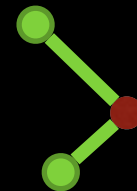
# How Typical FGM Algos. Work

**Example (t=2):**

***DB:***



Candidates (3):



etc...

Occurrences: {A,B,C} and {A,B,C} and {A,C}

Such an algorithm needs to be able to:

- remove isomorphic candidates (iso. not known to be in P)
- compute occurrences using subgraph isomorphism (NP-hard)



# Complexity of FGM

- **Complexity measures:**

- **Polynomial delay:** if the time between printing the next fr. graph (or terminating) is bounded by a polynomial of the size of input,
- **Incremental polynomial time:** if the time between printing next fr. graph (or terminating) is bounded by a polynomial of the size of input and of the size of output so far,
- **Output polynomial time:** if the algorithm finishes in time polynomial in the combined size of input and the entire output.

implies

implies

# Known Results

$\text{DB} \subseteq \text{All graphs}$

**NOT EVEN OUTPUT-POLY TIME POSSIBLE!**

Interesting cases

$\text{DB} \subseteq \text{Bounded-treewidth graphs}$   
**INCREMENTAL-POLY TIME!**  
[Horvath & Ramon, 2010]

Despite  
NP hard subgraph  
iso.

**Poly delay???**

**???**

**Open questions!**

$\text{DB} \subseteq \text{Hereditary graph classes with poly-time subgraph iso.}$

**POLY DELAY!**

# Change of Perspective

- **A more general problem (Ordered graph mining):**
  - Output all nonisomorphic connected graphs with frequency at least 1 and their occurrences in *DB* (*i.e. which DB graphs they match by subgraph iso.*):
    - **F  $\rightarrow$  I:** from frequent to infrequent (**generalizes FGM**)
    - **I  $\rightarrow$  F:** from infrequent to frequent (**generalizes IGM**)
    - **S  $\rightarrow$  L:** from smallest to largest
    - **L  $\rightarrow$  S:** from largest to smallest

*(If you cannot solve a problem, George Pólya in “How to Solve It” suggests studying a more general problem.)*

# Available Results

(From correspondence between FGM and  $F \rightarrow I$ )

	All Graphs	Planar Graphs	Bounded-Treewidth Graphs
$S \rightarrow L$	??	??	IncPoly [Horvath and Ramon, 2010]
$L \rightarrow S$	??	??	??
$F \rightarrow I$	Not IncPoly unless $P=NP$ [known]	??	IncPoly [Horvath and Ramon, 2010]
$I \rightarrow F$	??	??	??

# New Results and Corollaries

Corollaries of our theorems



	All Graphs	Planar Graphs	Bounded-Treewidth Graphs
S $\rightarrow$ L	Not IncPoly unless FPT = W[1]	??	IncPoly [Horvath and Ramon, 2010]
L $\rightarrow$ S	IncrPoly iff GI in P, Poly delay if CANON in P	Poly delay	Poly delay
F $\rightarrow$ I	Not IncPoly unless P=NP [known]	Not IncPoly unless P=NP	IncPoly [Horvath and Ramon, 2010]
I $\rightarrow$ F	Not IncPoly unless P=NP	Not IncPoly unless P=NP	Not IncPoly unless P=NP

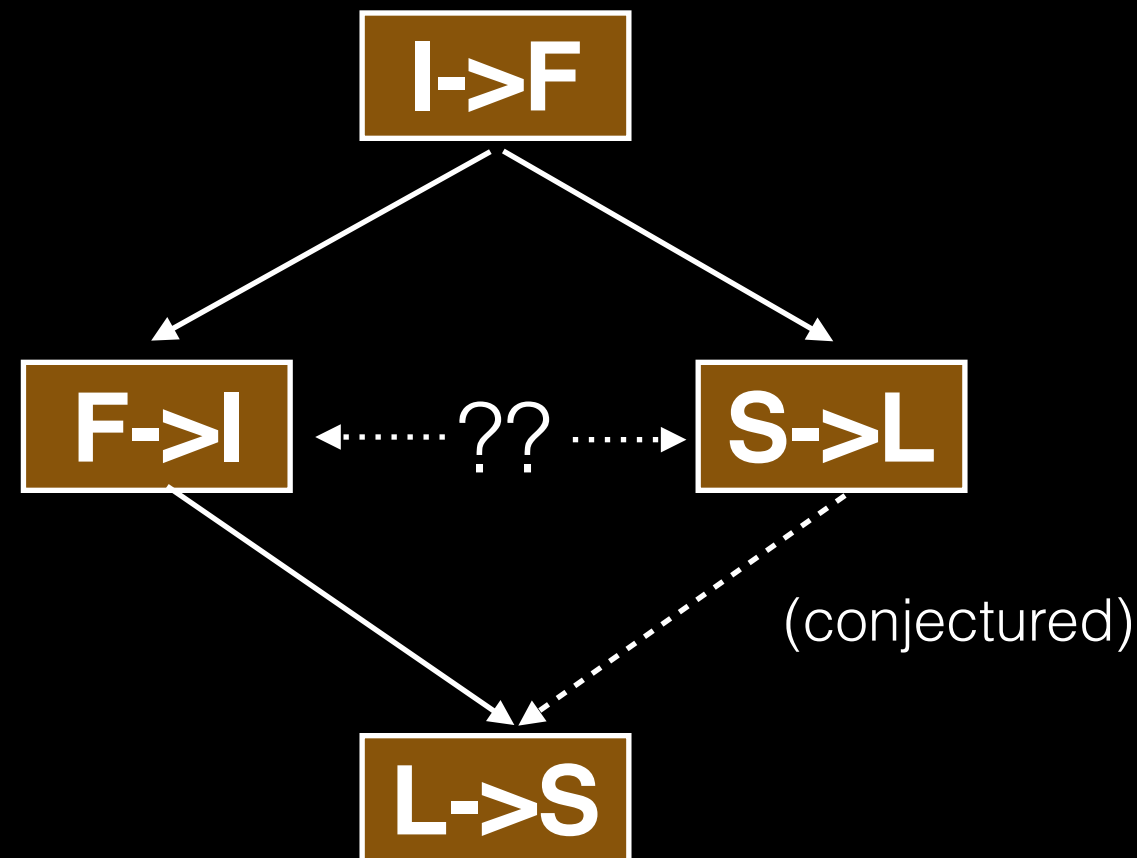
Positive

Negative

(More general results in the paper.)

# Relative Hardness

- Difficulty of the problems for the considered classes of graphs:



# Large to Small (Details)

**Require:** database  $DB$  of transaction graphs

**Ensure:** all connected (induced) subgraphs and their occurrences

```
1: let  $ALL$  be a data structure for storing graphs and their occurrences (as
   described in the main text).
2: for  $G \in DB$  do
3:   ADD( $G, \{ID(G)\}, ALL$ )
4: endfor
5: let  $m$  be the maximum order1 of a graph in  $DB$ .
6: for ( $l := m; l > 0; l := l - 1$ ) do
7:   for  $H \in KEYS(l, ALL)$  do
8:      $OCC \leftarrow GET(H, ALL)$ 
9:     PRINT( $H, OCC$ )
10:    for  $H' \in REFINE(H)$  do
11:      if  $H'$  is connected then
12:        ADD( $H', OCC, ALL$ )
13:      endif
14:    endfor
15:  endfor
16:  DELETE( $l, ALL$ )
17: endfor
```

- **Simple**, yet **poly-delay** algorithm for bounded TW graphs, planar graphs, ....
- It achieves poly-delay with NP-hard pattern matching operators and even if FGM cannot be solved in output-poly time (planar graphs).
- It may be combined with constraints such as maximum graph diameter which even leads to practical algorithms
- It can be generalised to **(induced) homeomorphism** and **(induced) minor emb.**

# Conclusions

- **Theory:**
  - **New results for complexity** of graph mining with NP-hard pattern matching operators (**some pretty surprising**).
  - We have proved analogical results for **induced subgraph isomorphism**, **(induced) homeomorphism** and **(induced) minor embedding**
- **Practice:**
  - Both the **positive** and **negative** results give guidelines e.g. for developing practical subgraph kernels.
  - **Larger-to-smaller** algorithm:
    - practically useful for mining subgraphs of bounded diameter
    - surprisingly also useful for mining **all** induced subgraphs of molecules of up to 25 non-hydrogen atoms (+ bigger molecules with additional hacks)



Thank you!