

Fast Graph Pattern Mining Made Easy

Xuhao Chen


Outline

- Graph Pattern Mining (GPM) Use Cases

- Pangolin Overview 

- Pangolin System Internals
 - How to Achieve Easy Programming?
 - How to Achieve High Performance?

- Evaluation

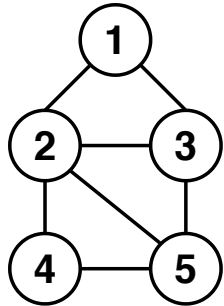
- Sandslash (optional) 

Graph Pattern Mining (GPM)

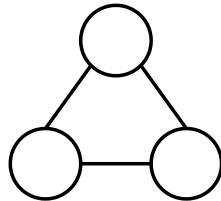
- **Inputs:** a graph G and a set of patterns $Q = \{P_1, P_2, \dots, P_n\}$
- **Goal:** count or list all subgraphs in G that matches $P_i \in Q$

We only consider
connected patterns

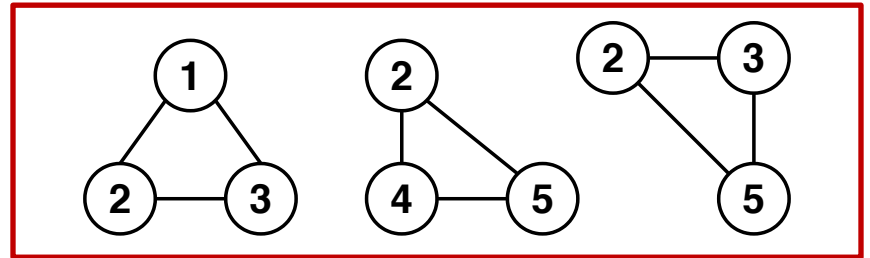
Input Graph G



Pattern P



Matched subgraphs

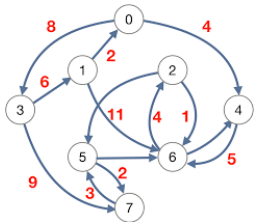


Graph Analytics vs. Graph Pattern Mining

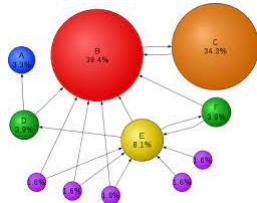
Graph Analytics

Updates vertex/edge labels of the input graph

Shortest Path



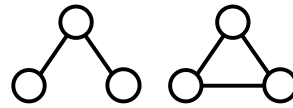
PageRank



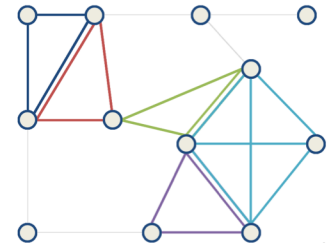
Graph Pattern Mining

Discovers structural patterns in the input graph

Motif Counting



Clique Listing



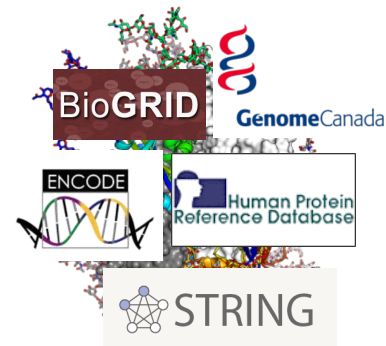
Why Graph Pattern Mining? (1/3)



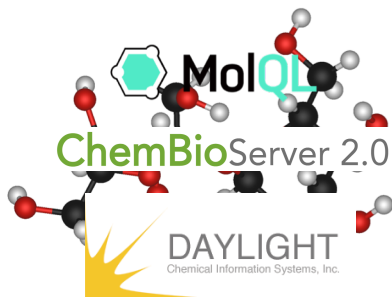
Social Networks



Recommender System



Bio-medicine



Chemical Engineering



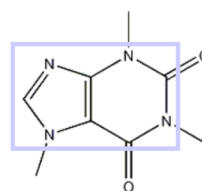
Graph Database

Why Graph Pattern Mining? (2/3)

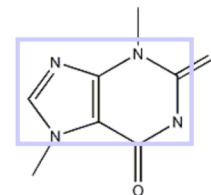
- Drug discovery



CHEMICAL COMPOUNDS

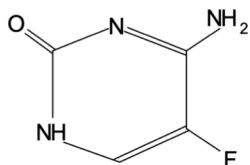
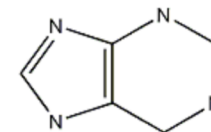


(a) caffeine

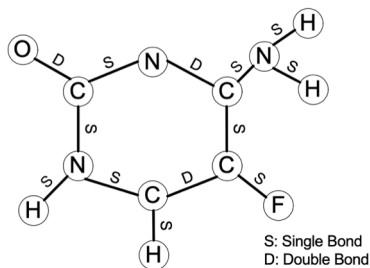


(b) diurobromine

FREQUENT SUBGRAPH



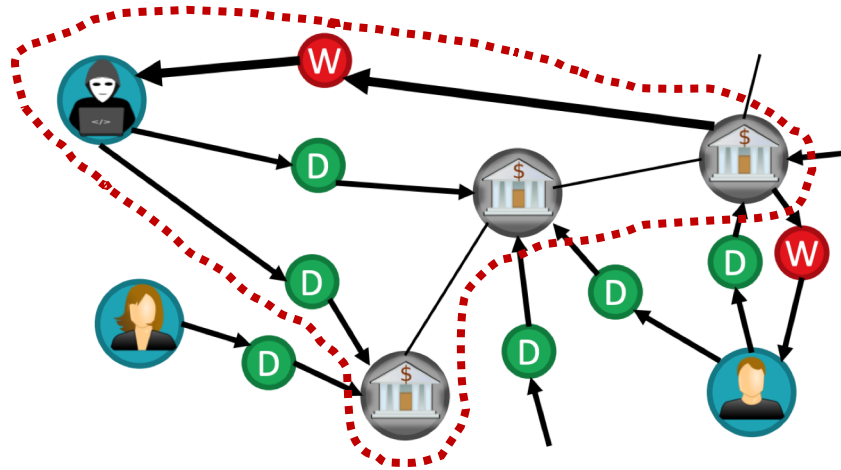
(a) NSC 103025 Flucytosine



(b) Graph Representation

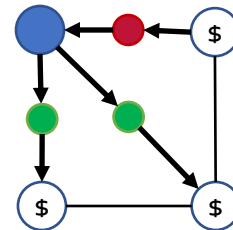
Figure modified from K. Borgwardt and X. Yan (KDD'08)

Why Graph Pattern Mining? (3/3)



Transactions and involved entities

Small deposits followed by large withdrawal



Outline

- Graph Pattern Mining (GPM) Use Cases
- **Pangolin Overview**
- Pangolin System Internals
 - How to Achieve Easy Programming?
 - How to Achieve High Performance?
- Evaluation



Pangolin: Efficient & Productive GPM Programming

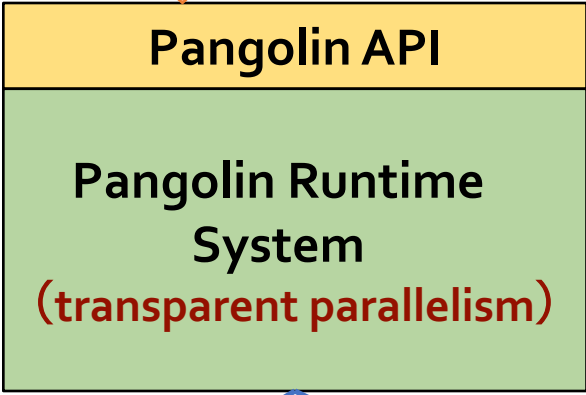
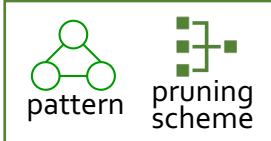
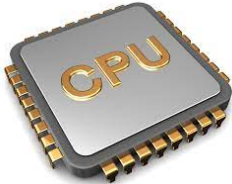


OpenMP



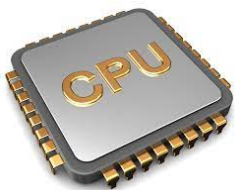
NVIDIA
CUDA

Too Difficult!



Sequential!
easy to code

run fast

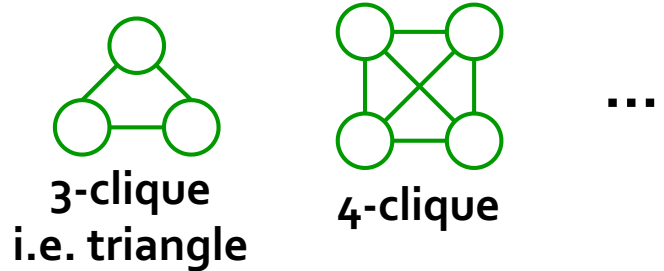


Easy Programming in Pangolin

◦ k -Clique Listing

```
1 bool toExtend(Subgraph sg, Vertex v) {  
2     return (sg.getLastVertex() == v);  
3 }  
  
4 bool toAdd(Subgraph sg, Vertex u) {  
5     for v ∈ sg.getVertices() except last:  
6         if (!isConnected(v,u)) return false;  
7     return true;  
8 }
```

Pangolin API
user defined functions

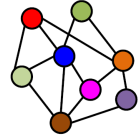


Hand Written	Pangolin
State-of-the-art: (KList: WWW'18) 395 Lines of Code !!!	8 Lines of code !!!

Outline

- Graph Pattern Mining (GPM) Use Cases
- Pangolin Overview
- **Pangolin System Internals**
 - How to Achieve Easy Programming?
 - How to Achieve High Performance?
- Evaluation

Graph Mining Challenges



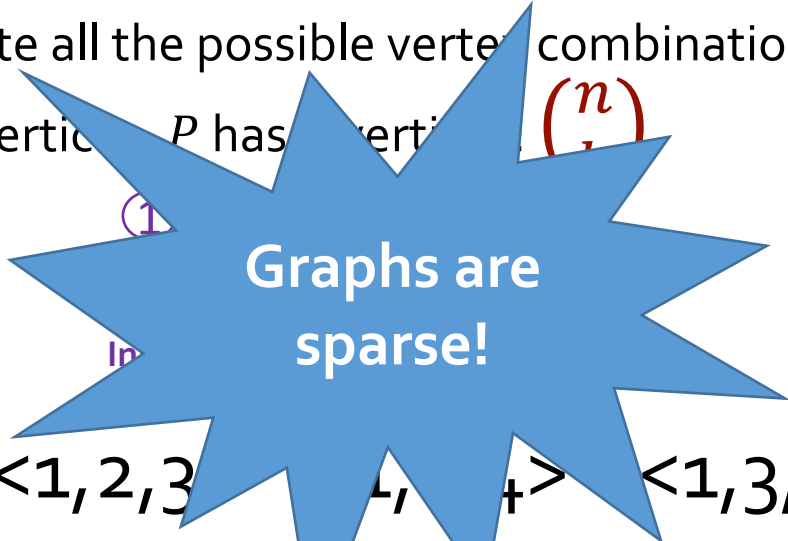
Easy Programming

High Performance



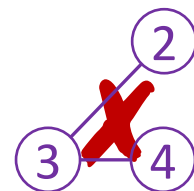
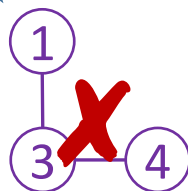
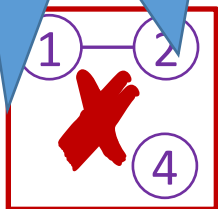
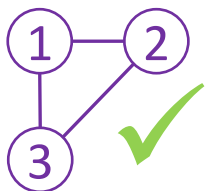
A "Naive" Abstraction

- Enumerate all the possible vertex combinations!!!
- G has n vertices, P has m vertices



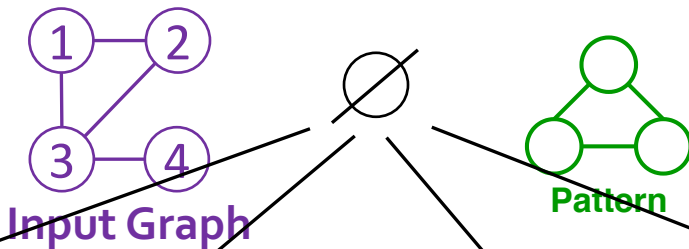
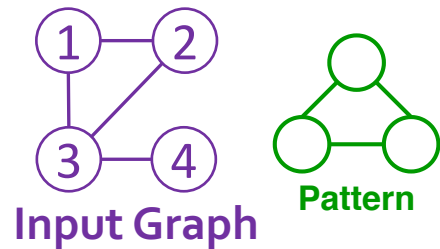
We only consider connected patterns

$$\binom{3}{4} = 4 : \langle 1, 2, 3 \rangle, \langle 1, 2, 4 \rangle, \langle 1, 3, 4 \rangle, \langle 2, 3, 4 \rangle$$



Isomorphism test

The "Subgraph Tree" Abstraction



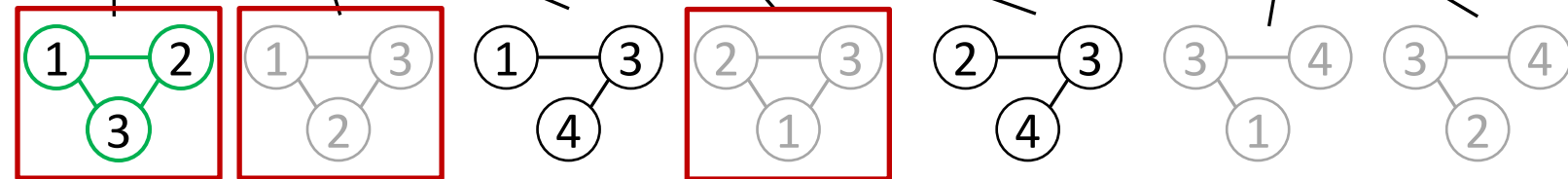
One abstraction for all the GPM problems

Level 1

Symmetry breaking



Level 2



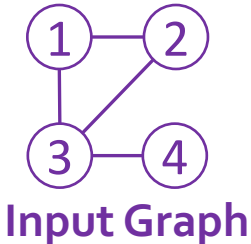
Level 3

Work Division between User & System

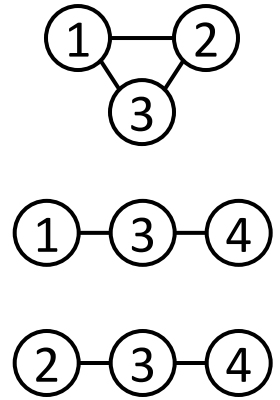
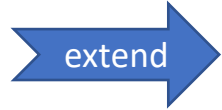
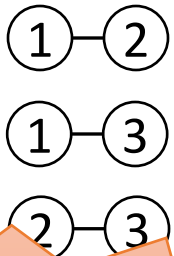
- Pangolin does a lot for you!

User responsibilities	Pangolin responsibilities	
<p>Pattern Specification</p>	<p>Subgraph Extension</p>	<p>Isomorphism Test</p>
<p>Search space Pruning</p>	<p>Symmetry Breaking</p>	<p>Parallelism Management</p>

Work Done By Pangolin



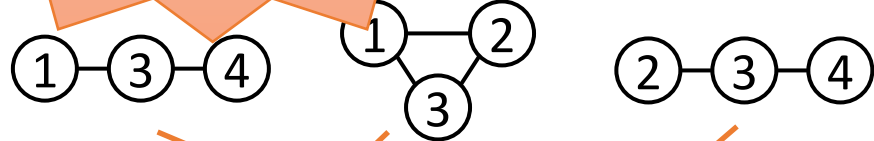
Extend Stage



Symmetry breaking

Transparently parallel!

Subgraphs



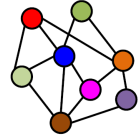
Reduce Stage

Isomorphism test

Pattern Map



Graph Mining Challenges



Easy Programming

High level abstraction
subgraph search tree

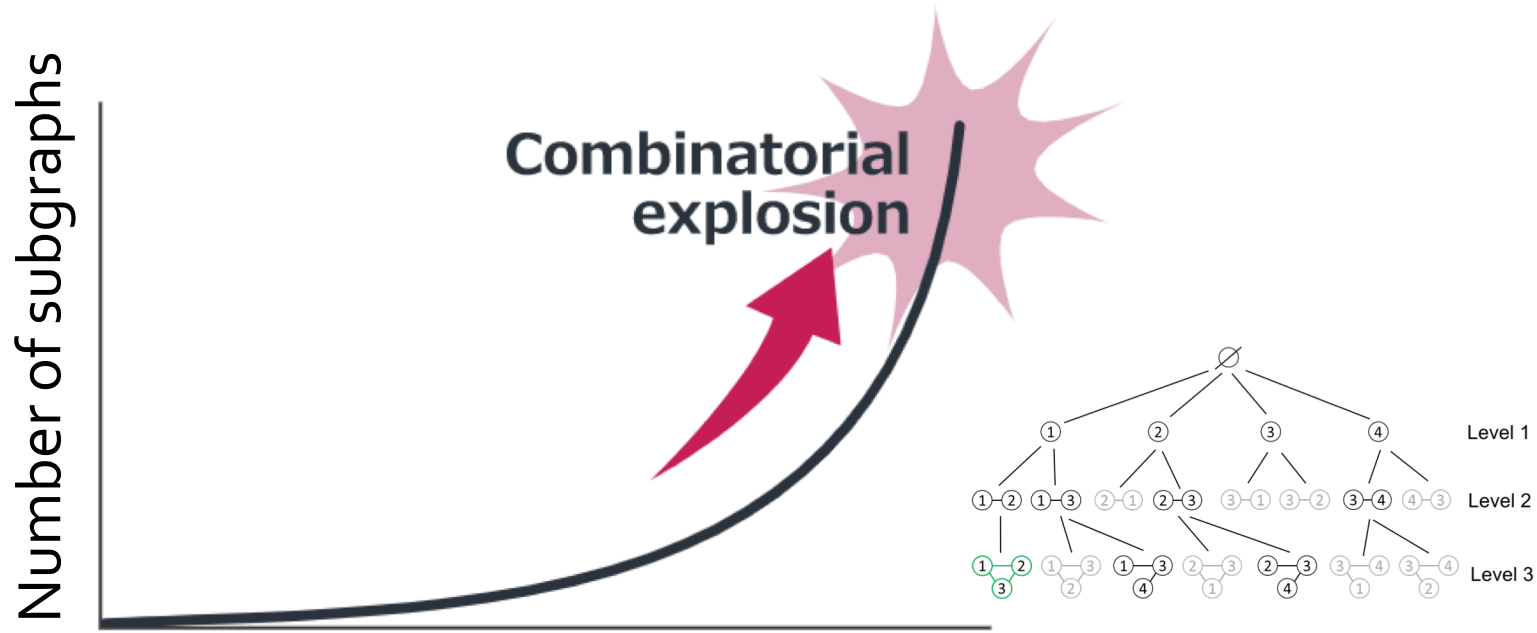
Work Division
Simple API

Transparent parallelism
code in sequential, not in parallel

High Performance



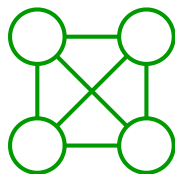
Search Space Pruning



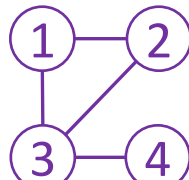
search space pruning is the key for GPM performance

Pattern-aware Search Space Pruning

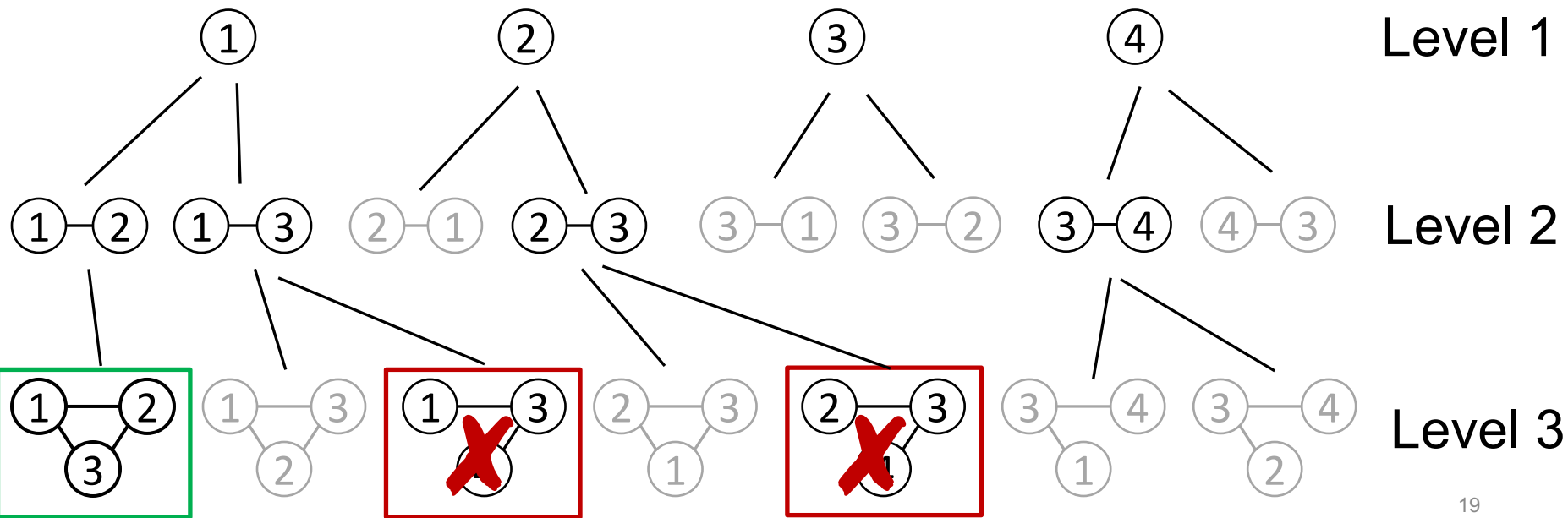
Problem: 4-Clique Listing



Pattern

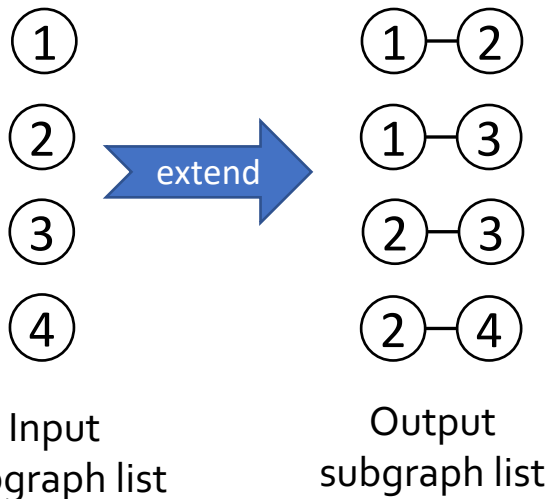


Input Graph



Exposing Interface for User Defined Pruning

Extend Stage



```
procedure Extend(SubgraphList in, SubgraphList out) {  
  for each subgraph  $sg \in in$   
    for each vertex  $v \in sg$   
      if (toExtend( $sg, v$ ) == true)  
        for each vertex  $u \in adj(v)$   
          if (toAdd( $sg, u$ ) == true)  
            out.insert( $sg \cup u$ )  
}
```

Two of the user defined functions*:

```
bool toExtend(Subgraph  $sg$ , Vertex  $v$ ) {  
  return true;  
}  
  
bool toAdd(Subgraph  $sg$ , Vertex  $u$ ) {  
  return true;  
}
```

* The full list of API functions can be found in the Pangolin paper: Chen et. al, VLDB 2020

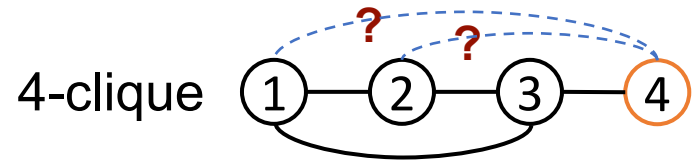
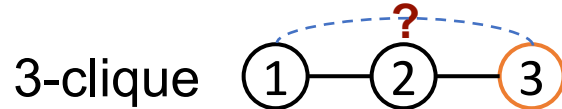
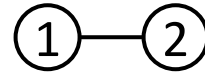
Easy Programming in Pangolin

◦ k -Clique Listing

- Extend only the last vertex v
- Check if new vertex u is connected to all previous vertices

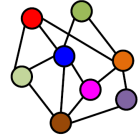
```
bool toExtend(Subgraph sg, Vertex v) {  
    return (sg.getLastVertex() == v);  
}
```

```
bool toAdd(Subgraph sg, Vertex u) {  
    for  $v \in$  sg.getVertices() except last:  
        if (!isConnected(v,u)) return false;  
    return true;  
}
```



Hand Written	Pangolin
State-of-the-art: (Klist: WWW'18) 395 Lines of Code !!!	8 Lines of code !!!

Graph Mining Challenges



Easy Programming

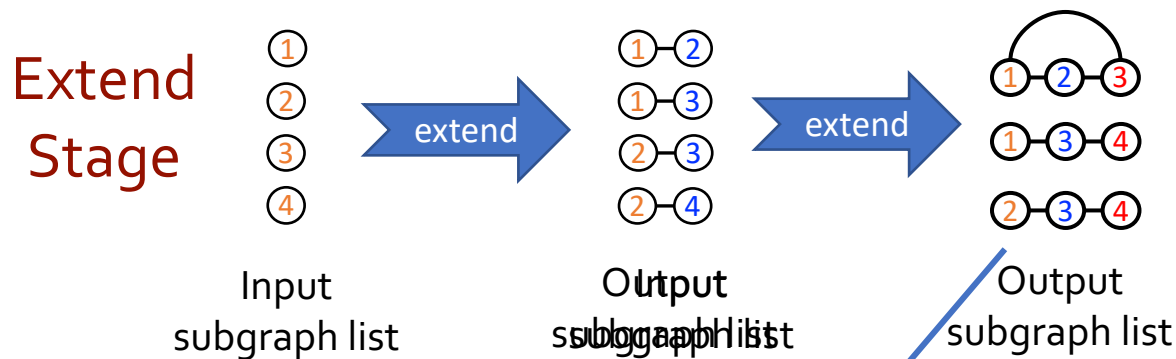
High Performance



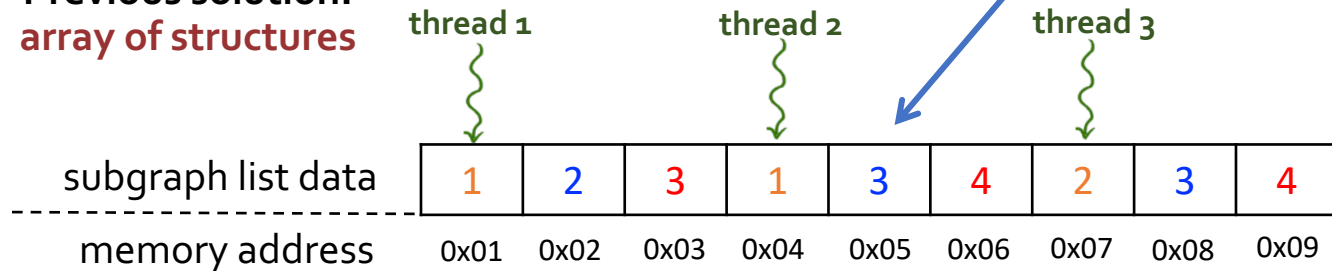
Search space
pruning

Architectural
Optimizations

Unlocking GPU Horsepower for GPM (1/3)



**Previous solution:
array of structures**



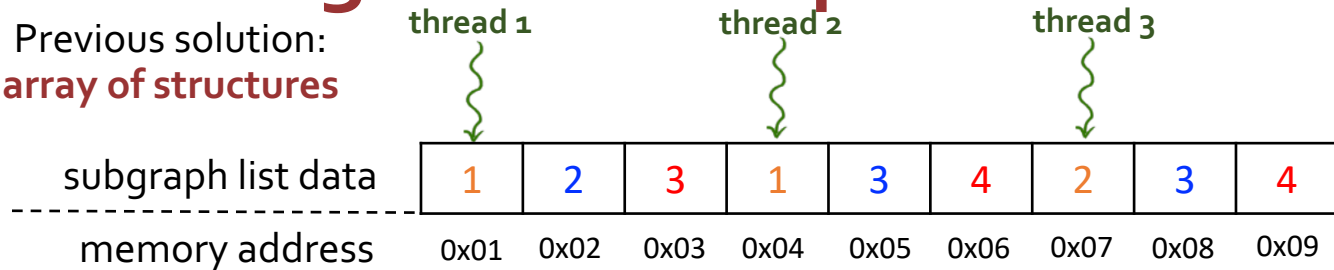
**Not GPU
friendly!**

Uncoalesced
memory accesses!!!



Unlocking GPU Horsepower for GPM (2/3)

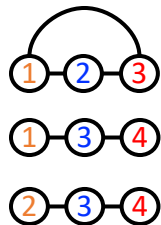
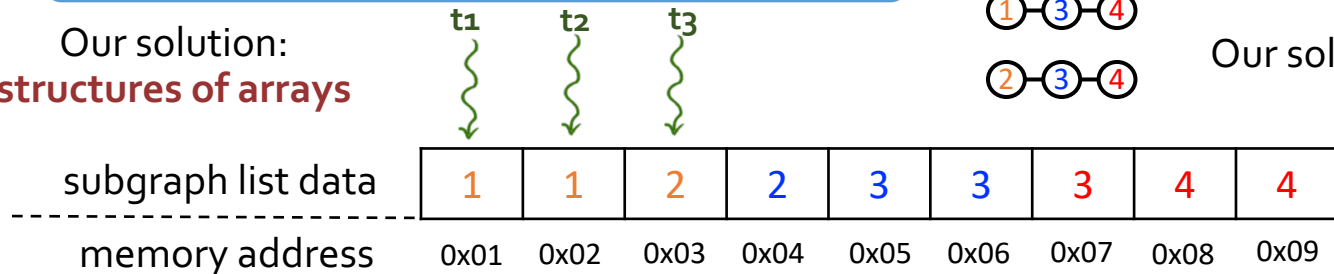
Previous solution:
array of structures



→ Not GPU friendly!

Key idea: Re-structure data layout

Our solution:
structures of arrays



Previous solution: one by one ❌

Our solution: level by level ✅

Coalesced memory accesses!!!

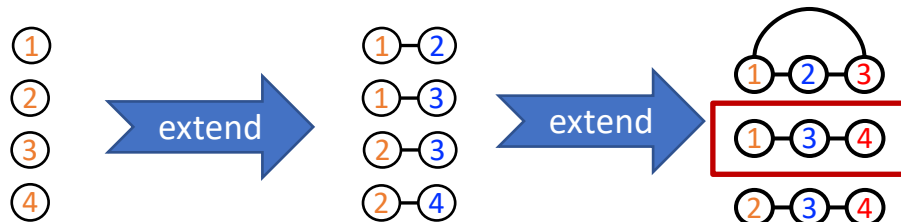


GPU friendly!

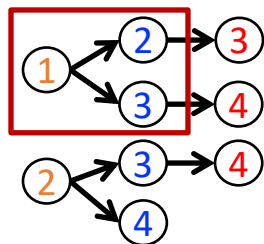


Unlocking GPU Horsepower for GPM (3/3)

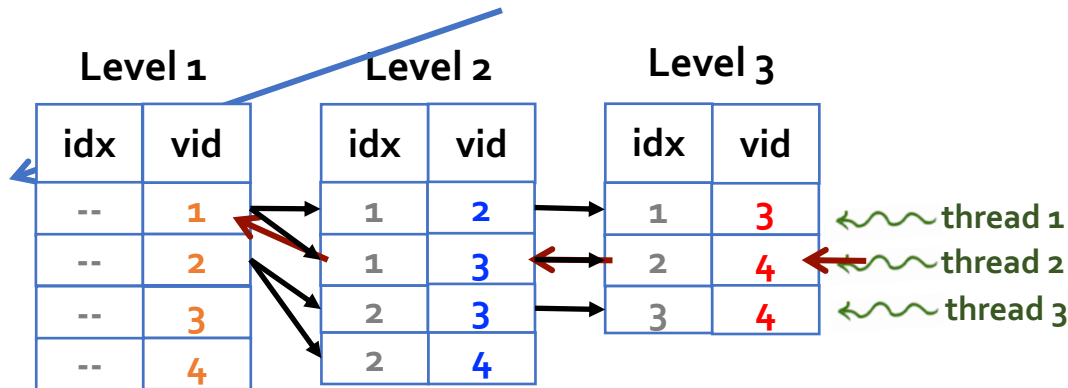
Extend Stage



Real implementation:
Build a **prefix tree**
level by level



Prefix tree

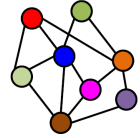


GPU friendly!

Coalesced memory accesses!!!



Graph Mining Challenges



Easy Programming

High Performance



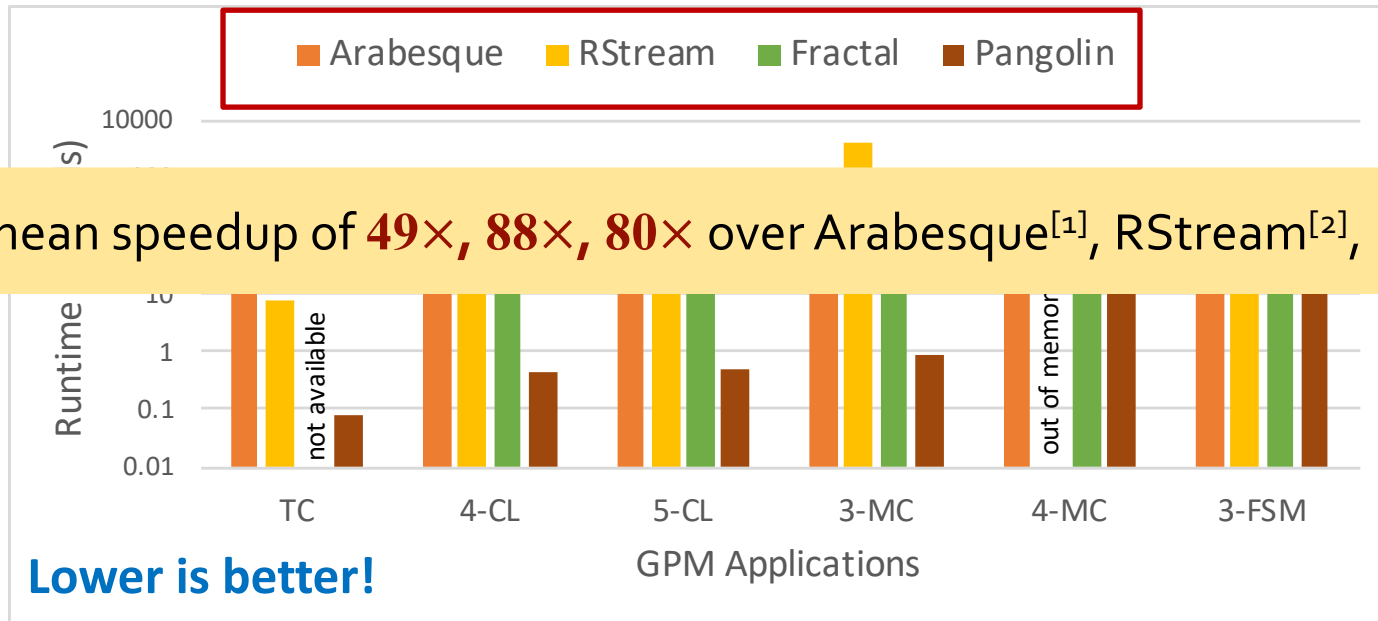
Search space
pruning

Architectural
Optimizations

Outline

- Graph Pattern Mining (GPM) Use Cases
- Pangolin Overview
- Pangolin System Internals
 - How to Achieve Easy Programming?
 - How to Achieve High Performance?
- **Evaluation**

Comparing with Existing Systems on CPU



TC: Triangle Counting

CL: Clique Listing

MC: Motif Counting

FSM: Frequent Subgraph Mining

Input graph: Patent-citations

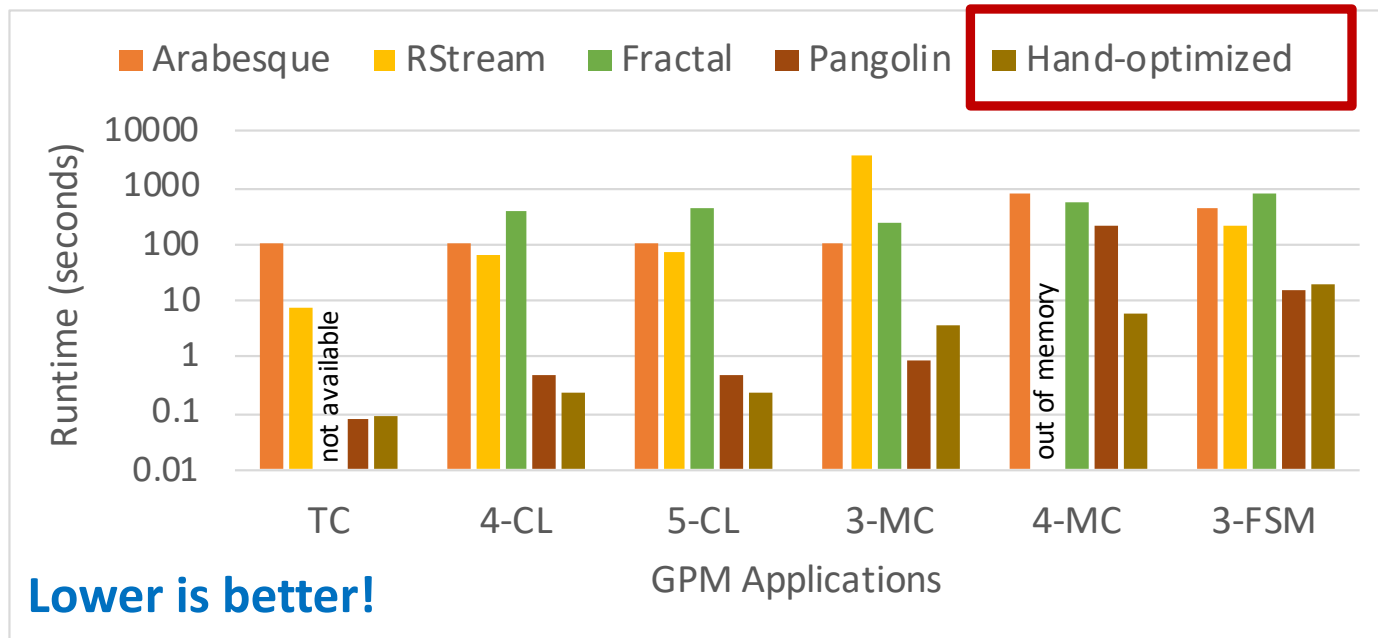
Intel Xeon Gold 5120 CPU (2.2GHz) 2 sockets
(14 cores each), 190GB memory, and 3TB SSD

[1] Arabesque SOSP'15

[2] RStream OSDI'18

[3] Fractal SIGMOD'19

Comparing with Hand-written Code on CPU



TC: Triangle Counting

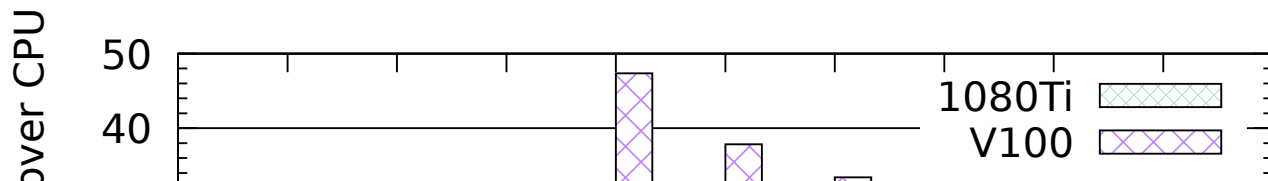
CL: Clique Listing

MC: Motif Counting

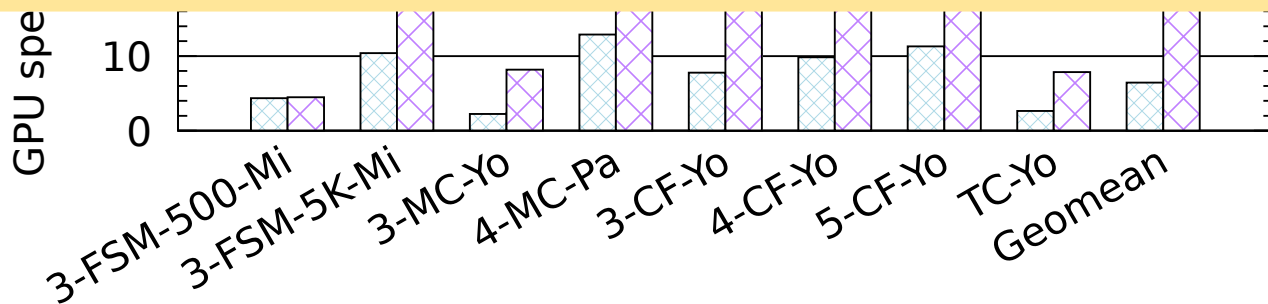
FSM: Frequent Subgraph Mining

Input graph: Patent-citations

Pangolin GPU Speedup over Pangolin CPU



Pangolin on V100 GPU gets **15×** speedup over Pangolin CPU



Speedup of Pangolin on GPU over Pangolin on 28-thread CPU

TC: Triangle Counting

CL: Clique Listing

MC: Motif Counting

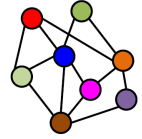
FSM: Frequent Subgraph Mining

Mi: Mico

Pa: Patent-citations

Yo: Youtube

Summary: GPM System Design



Easy Programming

Easy coding!

High level abstraction

Work Division

Transparent parallelism

**100x ~ 1000x
speedup!!!**

High Performance



Search space
pruning

Architectural
optimizations

Outline

- Graph Pattern Mining (GPM) Use Cases

- Pangolin Overview 

- Pangolin System Internals
 - How to Achieve Easy Programming?
 - How to Achieve High Performance?

- Evaluation

- **Sandlash (optional)** 

Sandslash*: A Two-level GPM System on CPU

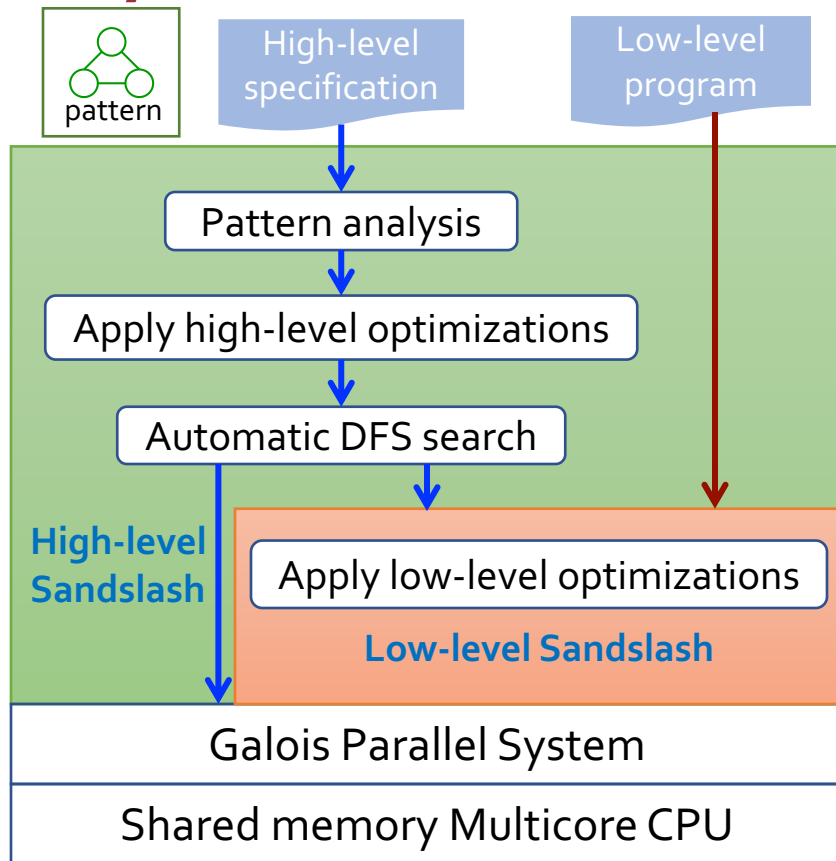


○ Easier programming

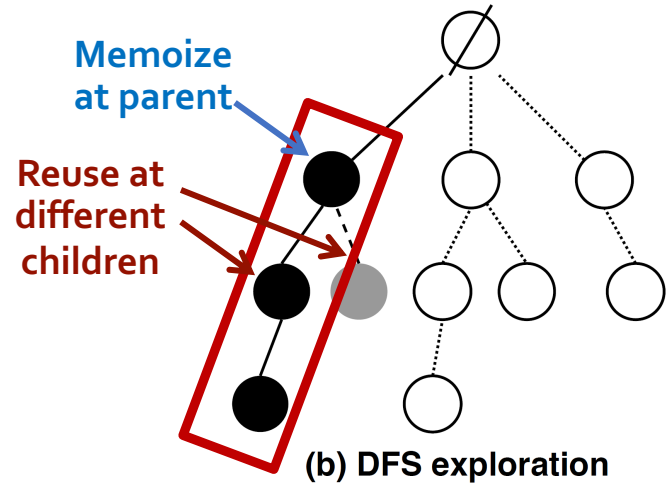
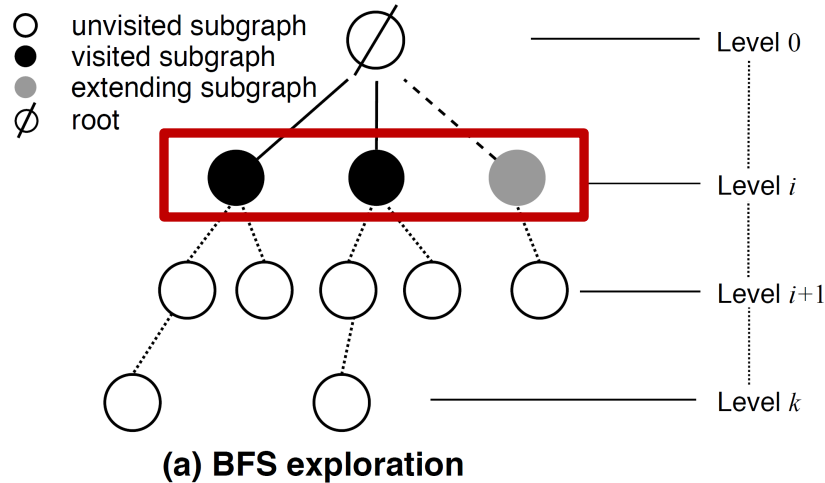
- High-level: pattern specification
- Low-level (**optional**): user defined pruning

○ Higher performance

- DFS exploration instead of BFS
- High-level optimizations: automated
- Low-level optimizations (**optional**)



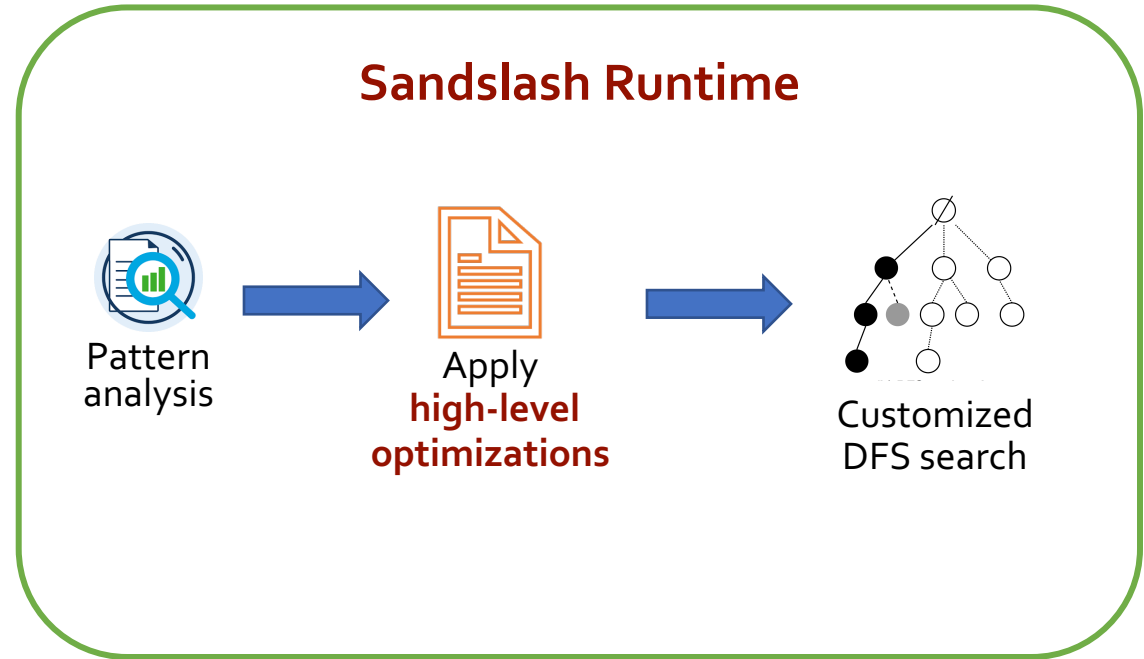
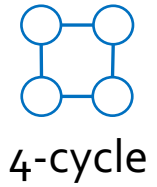
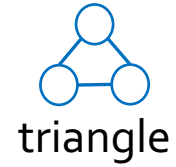
Search Tree Exploration: BFS vs. DFS



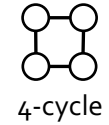
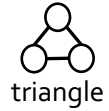
- Pangolin (both CPU and GPU) uses BFS
 - Good for massive parallelism on GPU
 - Large memory footprint for intermediate data
 - Lots of redundant computation

- Sandslash (CPU-only) uses DFS
 - Small and confined memory footprint
 - Avoid redundancy by **memoization**

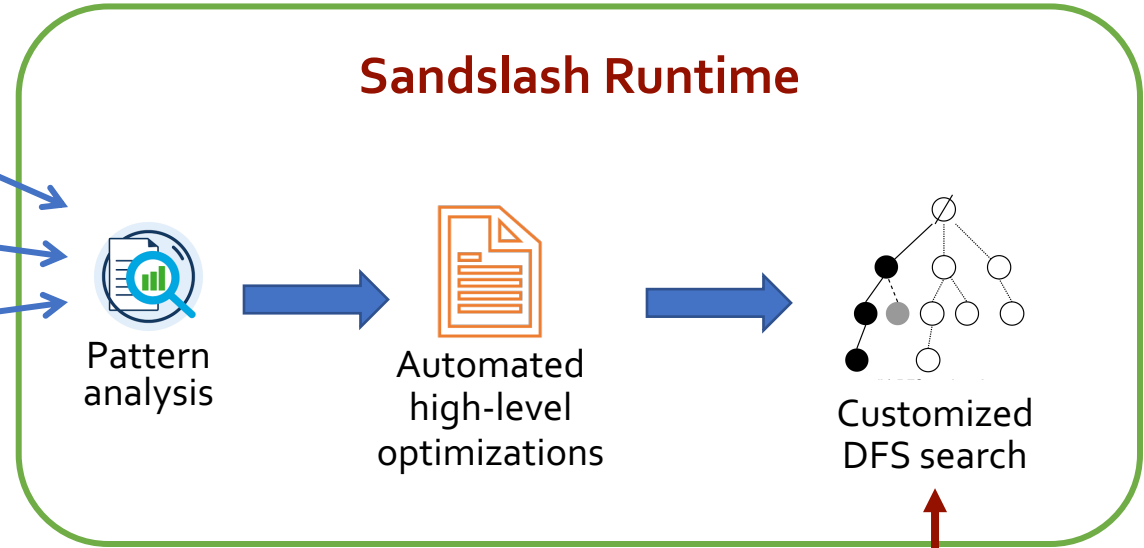
Automated High-level Optimizations



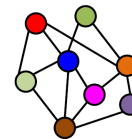
User Defined Low-level Optimizations




User specified
low-level
optimizations



Summary: GPM System on CPU



Easier Programming

Easier coding!

API separation

Automated high-level optimizations

On CPU **8×** and **14×**
speedup over Pangolin
and AutoMine*

Higher Performance



Avoid redundancy
with **memoization**

Low-level user
defined pruning

Thank You!
Q&A