# Graph Mining

# Graph Mining

A huge & complex graph dataset

# Graph Mining

A huge & complex graph dataset

Pattern counting (triangles, higher-order cliques, dense subgraphs, ...)
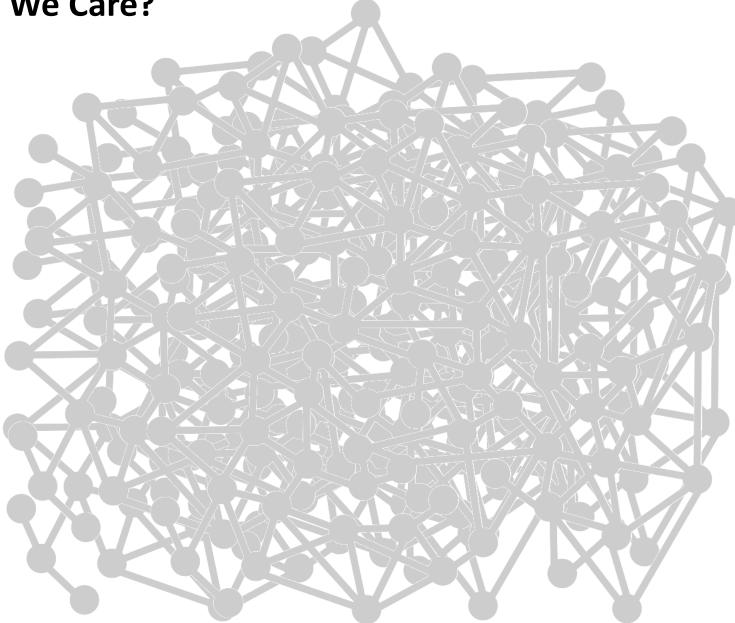
# Graph Mining

A huge & complex graph dataset

Pattern counting (triangles, higher-order cliques, dense subgraphs, ...)

Clustering, Link Prediction, Vertex Similarity, ...

# Graph Mining: Do We Care?

# Graph Mining: Do We Care?



Social sciences

# Graph Mining: Do We Care?

**Social sciences**

**Engineering**

# Graph Mining: Do We Care?



Social sciences

Biology

Chemistry

Engineering

Communication

Web graph analysis

Medicine

Cybersecurity

...even philosophy ☺

# Graph Mining: Do We Care?

Social sciences

Biology

Chemistry

Communication

Engineering

Web graph analysis

Me

...even philosop

## Challenges

Modeling a Philosophical Inquiry: from MySQL to a graph database

The short story of a long refactoring process

# Graph Mining & Graph Datasets: Challenges

# Graph Mining & Graph Datasets: Challenges

Huge



[1] Heng Lin et al.: ShenTu: Processing Multi-Trillion Edge Graphs on Millions of Cores in Seconds, SC18, Gordon Bell Finalist

4

# Graph Mining & Graph Datasets: Challenges



Huge

Irregular

[1] Heng Lin et al.: ShenTu: Processing Multi-Trillion Edge Graphs on Millions of Cores in Seconds, SC18, Gordon Bell Finalist

# Graph Mining & Graph Datasets: Challenges



Huge

Communication-heavy

Synchronization-heavy

Irregular

[1] Heng Lin et al.: ShenTu: Processing Multi-Trillion Edge Graphs on Millions of Cores in Seconds, SC18, Gordon Bell Finalist

# Graph Mining & Graph Datasets: Challenges

**Power-hungry**



**Huge**

**Communication-heavy**

**Irregular**

**Synchronization-heavy**

[1] Heng Lin et al.: ShenTu: Processing Multi-Trillion Edge Graphs on Millions of Cores in Seconds, SC18, Gordon Bell Finalist

# Graph Mining & Graph Datasets: Challenges

**Huge**

**Power-hungry**

**Communication-heavy**

**Synchronization-heavy**

**Irregular**

**Time complexities often O($n^k$) for k ≥ 2**

[1] Heng Lin et al.: ShenTu: Processing Multi-Trillion Edge Graphs on Millions of Cores in Seconds, SC18, Gordon Bell Finalist

# Goal: Making Graph Mining <u>Radically</u> Faster

**Goal: Making Graph Mining <u>Radically</u> Faster**



Do we need 100% accurate results in all cases?

# Goal: Making Graph Mining **Radically** Faster



Do we need 100% accurate results in all cases?

Let's say we can choose between…

Find all the patterns (e.g., cliques) in 1 day

Find ≥ 90% of all the patterns in 30 minutes

# Goal: Making Graph Mining <u>Radically</u> Faster

Do we need 100% accurate results in all cases?



Let's say we can choose between…

Find all the patterns (e.g., cliques) in 1 day
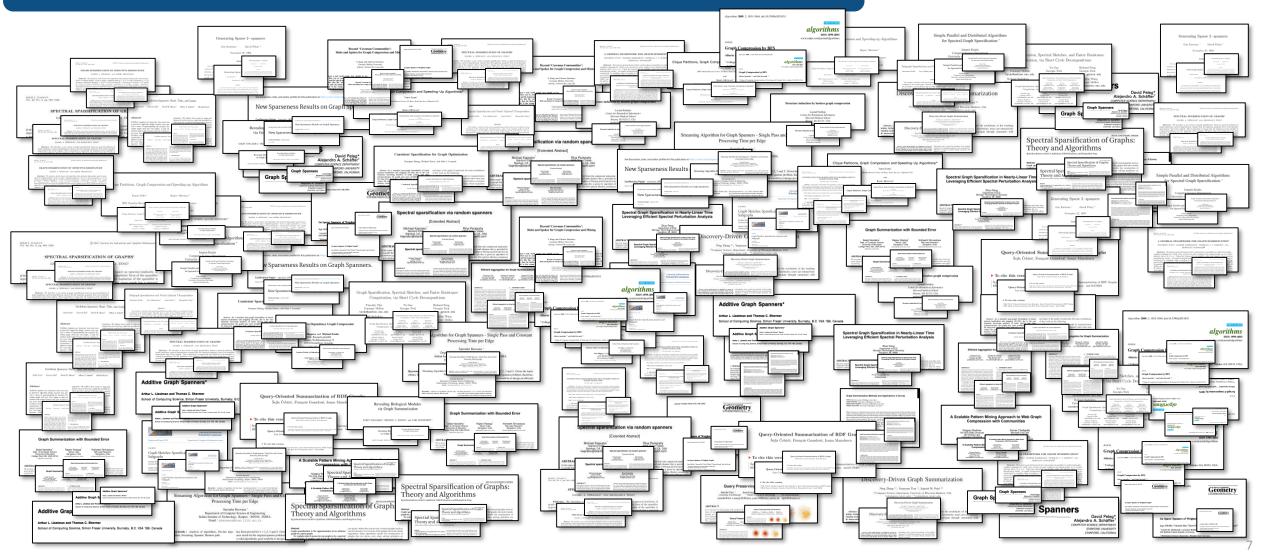
Find ≥ 90% of all the patterns in 30 minutes

# Approximate Graph Processing: State & Challenges

# Approximate Graph Processing: State & Challenges

We analyzed > 500 works and identified three classes of schemes…

# Approximate Graph Processing: State & Challenges

We analyzed > 500 works and identified three classes of schemes...
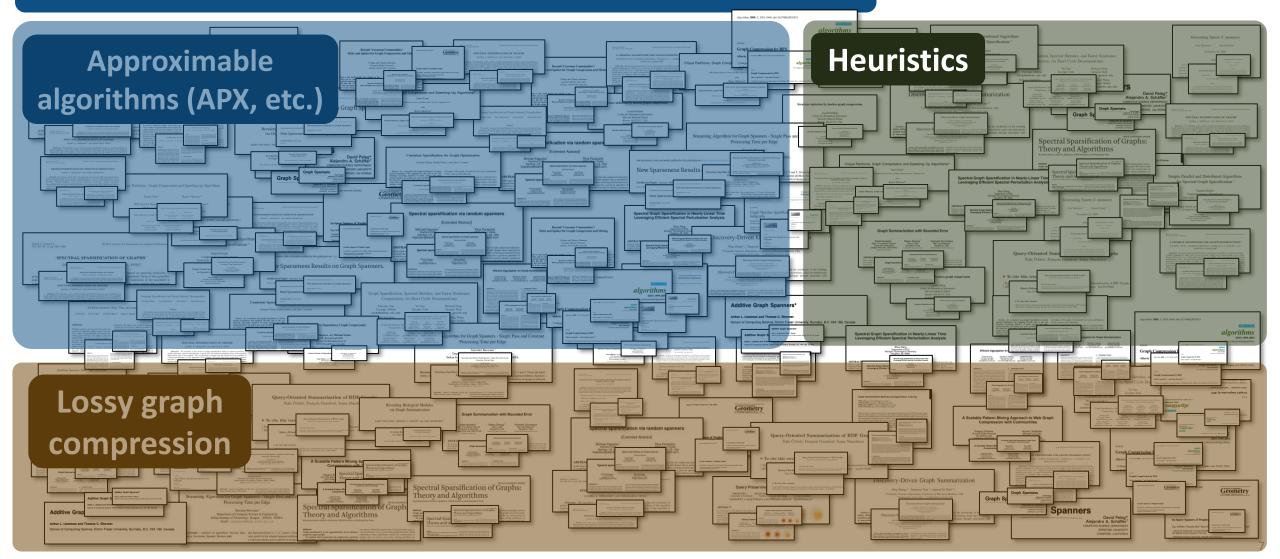
**Approximable algorithms (APX, etc.)**

**Heuristics**

**Lossy graph compression**

# Approximate Graph Processing: State & Challenges

**We analyzed > 500 works and identified three classes of schemes…**

**…they all have problems**

**Approximable algorithms (APX, etc.)**

**Heuristics**

**Lossy graph compression**

# Approximate Graph Processing: State & Challenges

We analyzed > 500 works and identified three classes of schemes…

…they all have problems

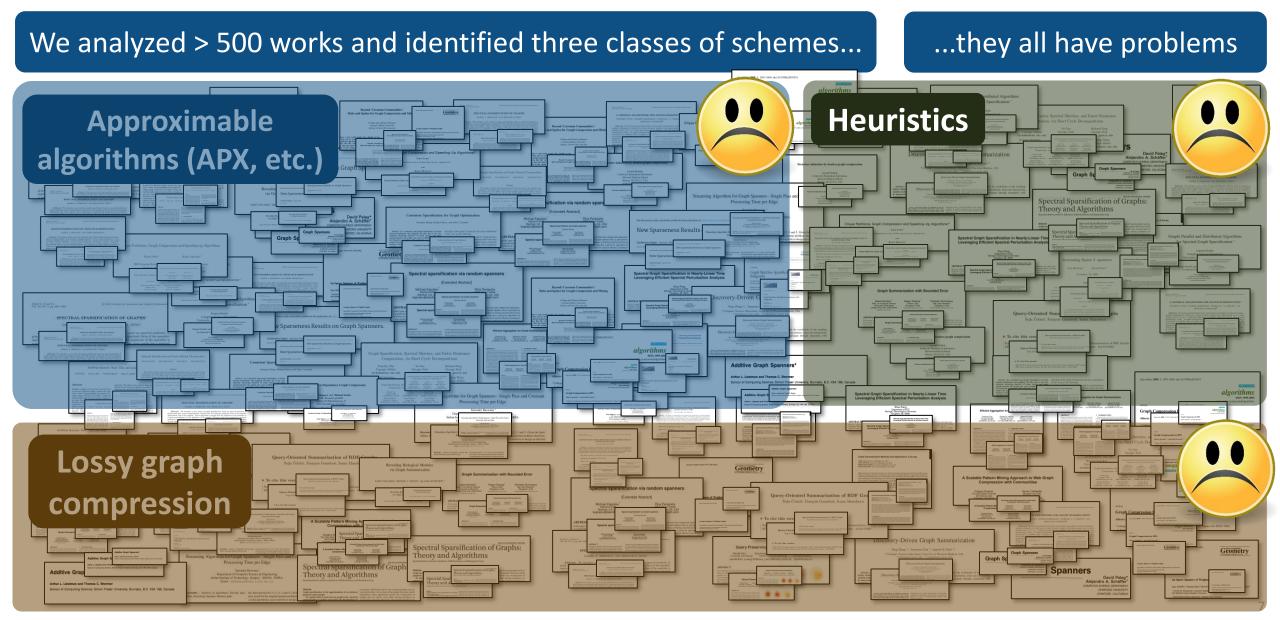Approximable algorithms (APX, etc.)

Heuristics

Little parallelism

Specific

Slow

Low accuracy

Lossy graph compression

# Approximate Graph Processing: State & Challenges

**We analyzed > 500 works and identified three classes of schemes...**

**...they all have problems**



**Approximable algorithms (APX, etc.)**

**Little parallelism**

**Specific**

**Slow**

**Low accuracy**

**Heuristics**

**Specific**

**No/loose accuracy guarantees**

**Lossy graph compression**

# Approximate Graph Processing: State & Challenges

We analyzed > 500 works and identified three classes of schemes...

...they all have problems

**Approximable algorithms (APX, etc.)**

Little parallelism

Specific

Slow

Low accuracy

**Heuristics**

Specific

No/loose accuracy guarantees

**Lossy graph compression**

Large memory overheads

No/loose accuracy guarantees

Slow

# Approximate Graph Processing: Current Issues & Our Objectives

# Approximate Graph Processing: Current Issues & Our Objectives



Little parallelism

Specific

No/loose accuracy guarantees

Slow

Low accuracy

Large memory overheads

# Approximate Graph Processing: Current Issues & Our Objectives

| | |
|---|---|
| Little parallelism | → Rich parallelism |
| Specific | |
| No/loose accuracy guarantees | |
| Slow | |
| Low accuracy | |
| Large memory overheads | |

# Approximate Graph Processing: Current Issues & Our Objectives

| | |
|---|---|
| Little parallelism | → Rich parallelism |
| Specific | → Wide applicability |
| No/loose accuracy guarantees | |
| Slow | |
| Low accuracy | |
| Large memory overheads | |

# Approximate Graph Processing: Current Issues & Our Objectives



| Little parallelism | → | Rich parallelism |
| Specific | → | Wide applicability |
| No/loose accuracy guarantees | → | Strong accuracy guarantees |
| Slow | | |
| Low accuracy | | |
| Large memory overheads | | |

# Approximate Graph Processing: Current Issues & Our Objectives



| Little parallelism | → | Rich parallelism |
| Specific | → | Wide applicability |
| No/loose accuracy guarantees | → | Strong accuracy guarantees |
| Slow | → | High performance |
| Low accuracy | | |
| Large memory overheads | | |

# Approximate Graph Processing: Current Issues & Our Objectives

| | | |
|---|---|---|
| Little parallelism | → | Rich parallelism |
| Specific | → | Wide applicability |
| No/loose accuracy guarantees | → | Strong accuracy guarantees |
| Slow | → | High performance |
| Low accuracy | → | High accuracy |
| Large memory overheads | | |

# Approximate Graph Processing: Current Issues & Our Objectives



| Little parallelism | → | Rich parallelism |
| Specific | → | Wide applicability |
| No/loose accuracy guarantees | → | Strong accuracy guarantees |
| Slow | → | High performance |
| Low accuracy | → | High accuracy |
| Large memory overheads | → | Low & controllable memory overheads |

How to achieve all these objectives in a single design?

No/loose accuracy guarantees → Strong accuracy guarantees

Slow → High performance

Low accuracy → High accuracy

Large memory overheads → Low & controllable memory overheads

**Approximate Graph Processing: Current Issues & Our Objectives**

How to achieve all these objectives in a single design?

No/loose accuracy guarantees

Strong accuracy guarantees

We develop **ProbGraph**: a graph representation that uses probabilistic set representations (aka sketches)

overheads

memory overheads

8

# High-Level Approach Taken in ProbGraph

# High-Level Approach Taken in ProbGraph

Keep the original graph

# High-Level Approach Taken in ProbGraph

Keep the original graph

Maintain a very small "sketch" of a graph



+

# High-Level Approach Taken in ProbGraph

Keep the original graph

Maintain a very small "sketch" of a graph



**+**

Use the sketch to answer performance critical queries

# High-Level Approach Taken in ProbGraph

**Keep the original graph**

**Maintain a very small "sketch" of a graph**

**+**

**What design to use for the sketch, to satisfy all the goals?**

**Use the sketch to answer performance critical queries**

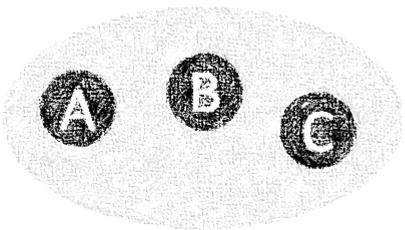# **ProbGraph key idea: Use probabilistic set representations (set sketches)**

# ProbGraph key idea: Use probabilistic set representations (set sketches)

A set = {A, B, C}

# **ProbGraph key idea**: Use probabilistic set representations (set sketches)

A set = {A, B, C}

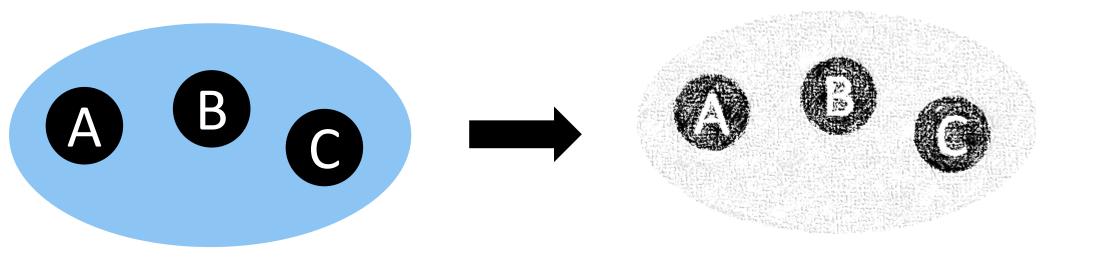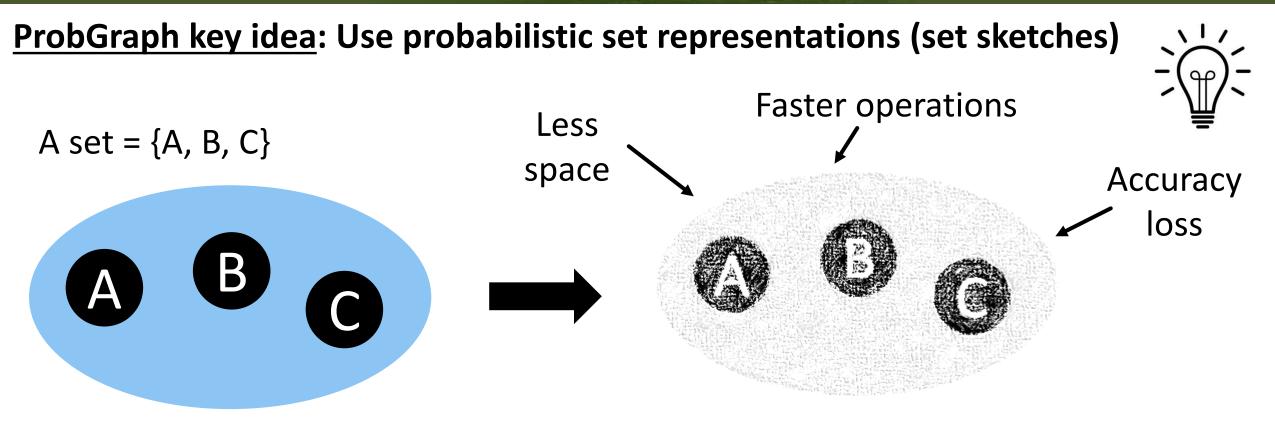# **ProbGraph key idea**: Use probabilistic set representations (set sketches)

A set = {A, B, C}

# ProbGraph key idea: Use probabilistic set representations (set sketches)

A set = {A, B, C}



**Bloom filters (BF) [1]**
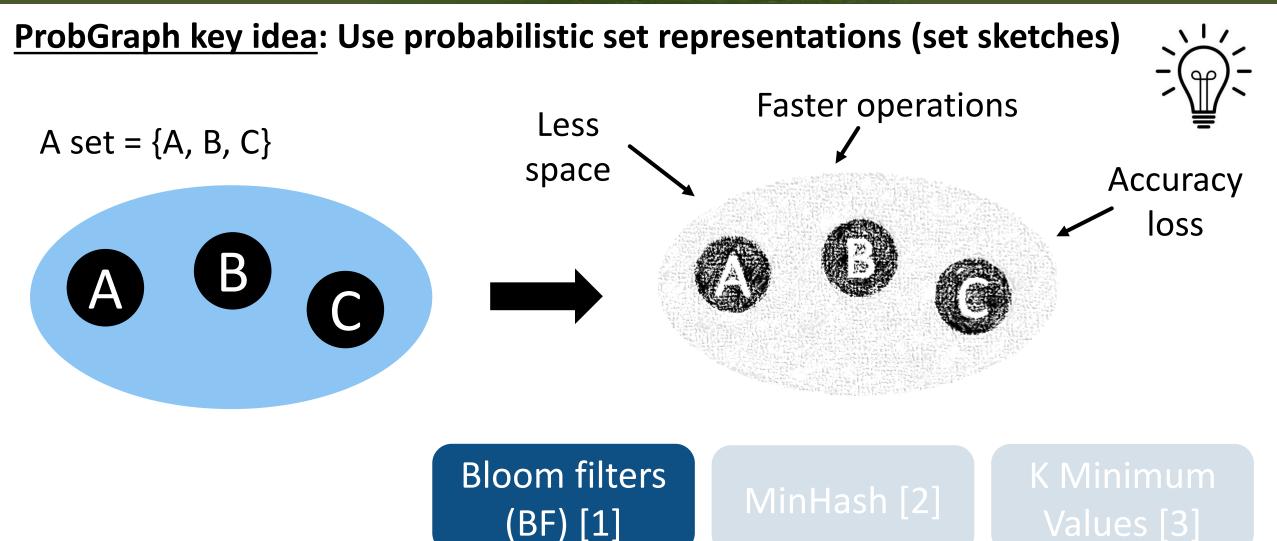
**MinHash [2]**

**K Minimum Values [3]**

[1] B. H. Bloom, "Space/time trade-offs in hash coding with allowable errors", CACM, 1970.
[2] A. Z. Broder, "On the resemblance and containment of documents", IEEE SEQUENCES, 1997.
[3] Z. Bar-Yossef et al., "Counting distinct elements in a data stream", in RANDOM, 2002.

# **ProbGraph key idea:** Use probabilistic set representations (set sketches)

A set = {A, B, C}

Less space

Faster operations

Accuracy loss



| Bloom filters (BF) [1] | MinHash [2] | K Minimum Values [3] |

[1] B. H. Bloom, "Space/time trade-offs in hash coding with allowable errors", CACM, 1970.
[2] A. Z. Broder, "On the resemblance and containment of documents", IEEE SEQUENCES, 1997.
[3] Z. Bar-Yossef et al., "Counting distinct elements in a data stream", in RANDOM, 2002.
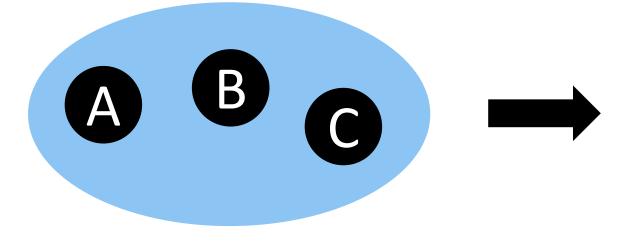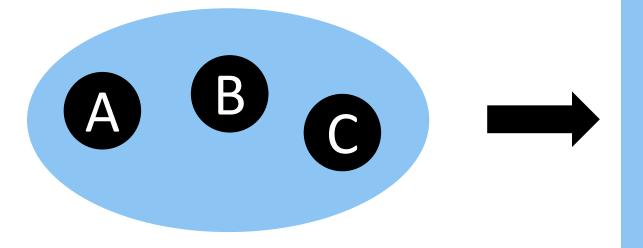
# **ProbGraph key idea: Use probabilistic set representations (set sketches)**

A set = {A, B, C}

Less space

Faster operations

Accuracy loss



| Bloom filters (BF) [1] | MinHash [2] | K Minimum Values [3] |

[1] B. H. Bloom, "Space/time trade-offs in hash coding with allowable errors", CACM, 1970.
[2] A. Z. Broder, "On the resemblance and containment of documents", IEEE SEQUENCES, 1997.
[3] Z. Bar-Yossef et al., "Counting distinct elements in a data stream", in RANDOM, 2002.

# Bloom Filters for Graph Mining

A set = {A, B, C}

# Bloom Filters for Graph Mining

A set = {A, B, C}



**Bloom filter $\mathcal{B}_X$ of X**

Bitvector of size $B_X$ [bits]

# Bloom Filters for Graph Mining

A set = {A, B, C}



**Bloom filter $\mathcal{B}_X$ of X**

Bitvector of size $B_X$ [bits]        $B_X = 12$



Hash functions $h_1, ..., h_b$
$h_i : X \rightarrow \{1, ..., B_X\}$

# Bloom Filters for Graph Mining

A set = {A, B, C}



**Bloom filter $\mathcal{B}_X$ of X**

Bitvector of size $B_X$ [bits]　　　$B_X = 12$

Hash functions $h_1, ..., h_b$
$h_i : X \rightarrow \{1, ..., B_X\}$

$b = 2$
$h_2, h_1 : X \rightarrow \{1, ..., 12\}$

# Bloom Filters for Graph Mining

A set = {A, B, C}



**Bloom filter $\mathcal{B}_X$ of X**

Bitvector of size $B_X$ [bits]          $B_X = 12$



Hash functions $h_1, ..., h_b$
$h_i : X \to \{1, ..., B_X\}$

$b = 2$

$h_2, h_1 : X \to \{1, ..., 12\}$

$h_1(A) = 3$

$h_2(A) = 5$

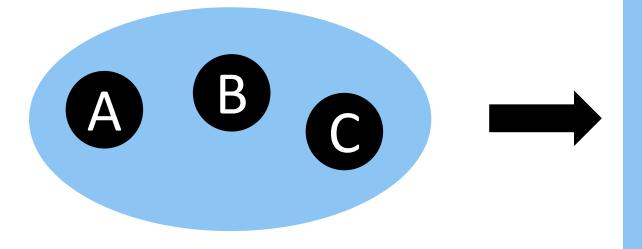# Bloom Filters for Graph Mining

A set = {A, B, C}



**Bloom filter $\mathcal{B}_X$ of X**

Bitvector of size $B_X$ [bits]          $B_X = 12$



Hash functions $h_1, ..., h_b$
$h_i : X \rightarrow \{1, ..., B_X\}$

$b = 2$

$h_2, h_1 : X \rightarrow \{1, ..., 12\}$

$h_1(\text{A}) = 3$

$h_2(\text{A}) = 5$
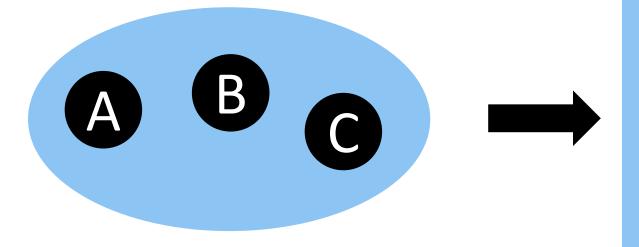
# Bloom Filters for Graph Mining

A set = {A, B, C}



## Bloom filter $\mathcal{B}_X$ of X

Bitvector of size $B_X$ [bits]

$B_X = 12$

Hash functions $h_1, ..., h_b$
$h_i : X \rightarrow \{1, ..., B_X\}$

$b = 2$

$h_2, h_1 : X \rightarrow \{1, ..., 12\}$

$h_1(A) = 3$    $h_1(B) = 1$

$h_2(A) = 5$    $h_2(B) = 8$

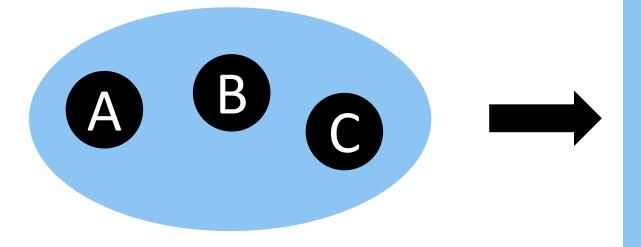# Bloom Filters for Graph Mining

A set = {A, B, C}



**Bloom filter $\mathcal{B}_X$ of X**

Bitvector of size $B_X$ [bits]        $B_X = 12$

Hash functions $h_1, ..., h_b$
$h_i : X \rightarrow \{1, ..., B_X\}$

$b = 2$

$h_2, h_1 : X \rightarrow \{1, ..., 12\}$

$h_1(A) = 3$      $h_1(B) = 1$

$h_2(A) = 5$      $h_2(B) = 8$

# Bloom Filters for Graph Mining
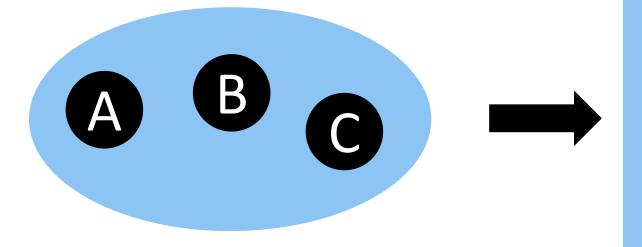
A set = {A, B, C}



## Bloom filter $\mathcal{B}_X$ of X

Bitvector of size $B_X$ [bits]          $B_X = 12$



Hash functions $h_1, ..., h_b$
$h_i : X \rightarrow \{1, ..., B_X\}$

$b = 2$

$h_2, h_1 : X \rightarrow \{1, ..., 12\}$

$h_1(A) = 3$          $h_1(B) = 1$          $h_1(C) = 4$

$h_2(A) = 5$          $h_2(B) = 8$          $h_2(C) = 11$

# Bloom Filters for Graph Mining

# Bloom Filters for Graph Mining

Each neighborhood $N_u$ is a set of vertices



$u$

$N_u$

# Bloom Filters for Graph Mining

Each neighborhood $N_u$ is a set of vertices

$u$

$N_u$

**Bloom Filters for Graph Mining**

Each neighborhood $N_u$ is a set of vertices ➡ „Sketch" each $N_u$ with a Bloom filter

u

$N_u$

# Bloom Filters for Graph Mining

Each neighborhood $N_u$ is a set of vertices

$\Rightarrow$

„Sketch" each $N_u$ with a Bloom filter

$u$

$N_u$

# Bloom Filters for Graph Mining

Each neighborhood $N_u$ is a set of vertices

➡️

„Sketch" each $N_u$ with a Bloom filter



u

$N_u$

11

# Bloom Filters for Graph Mining

Each neighborhood $N_u$ is a set of vertices → „Sketch" each $N_u$ with a Bloom filter

# ProbGraph: Summary of Design

# ProbGraph: Summary of Design

**Input graph G**

# ProbGraph: Summary of Design

**Input graph G**



**Standard graph representation (e.g., CSR)**

# ProbGraph: Summary of Design

**Input graph G**



**Standard graph representation (e.g., CSR)**

# ProbGraph: Summary of Design

**Input graph G**



**ProbGraph representation**

# ProbGraph: Summary of Design

**Input graph G**

**ProbGraph representation**

Bloom filters

# ProbGraph: Summary of Design

**Input graph G**



**ProbGraph representation**

Bloom filters

| 1 | → | 2 | 3 | | | |
| 2 | → | 1 | 3 | 4 | | |
| 3 | → | 1 | 2 | 4 | 5 | 6 |
| 4 | → | 2 | 3 | 5 | | |
| 5 | → | 3 | 4 | 6 | 7 | 8 |
| 6 | → | 3 | 5 | 7 | | |
| 7 | → | 5 | 6 | 8 | | |
| 8 | → | 5 | 7 | | | |

**+**

**Larger $B_X$** : more accuracy & more storage required. **Lower $B_X$** : vice versa.

$B_X$ [bits]

13

# ProbGraph: Summary of Design

## Input graph G



$B_X$ is often small → little storage

## ProbGraph representation

Bloom filters



$$+$$

**Larger $B_X$** : more accuracy & more storage required. **Lower $B_X$** : vice versa.

$B_X$ [bits]

13

# ProbGraph: Summary of Design

**Input graph G**

**ProbGraph representation**

Bloom filters

$B_X$ is often small → little storage

BFs have the same size → great load balancing

**Larger $B_X$** : more accuracy & more storage required. **Lower $B_X$** : vice versa.

$B_X$ [bits]

How does our idea compare to other Bloom filter use cases?

# How does our idea compare to other Bloom filter use cases?

Traditional BF use case: presence tracking

# How does our idea compare to other Bloom filter use cases?

## Traditional BF use case: presence tracking

**Data** stored somewhere

# How does our idea compare to other Bloom filter use cases?

## Traditional BF use case: presence tracking

**Data** stored somewhere

**A BF cache** tracking the presence of data

# How does our idea compare to other Bloom filter use cases?

## Traditional BF use case: presence tracking

**Data** stored somewhere

**A BF cache** tracking the presence of data

⬜⬜⬜⬜⬜⬜⬜⬜

# How does our idea compare to other Bloom filter use cases?

Traditional BF use case: presence tracking

**Data** stored somewhere

**A BF cache** tracking the presence of data

Insert an element

# How does our idea compare to other Bloom filter use cases?

## Traditional BF use case: presence tracking

**Data** stored somewhere

**A BF cache** tracking the presence of data

Set the appropriate BF bits

Insert an element

# How does our idea compare to other Bloom filter use cases?

Traditional BF use case: presence tracking

Is the data in question over there?

**A BF cache** tracking the presence of data

**Data** stored somewhere

Yes

# How does our idea compare to other Bloom filter use cases?

## Traditional BF use case: presence tracking

**Data** stored somewhere

A BF cache tracking
the presence of data

Is the data in
question over there?

Yes

Fetch the data

# How does our idea compare to other Bloom filter use cases?

Traditional BF use case: presence tracking



**Data** stored somewhere

A BF cache tracking the presence of data

This is a very fast operation

Is the data in question over there?

Yes

Fetch the data

This is usually a slow operation

Other Bloom filter use cases?

**The novelty of ProbGraph**

We use BFs as a sketch of the actual dataset

**Data** stored somewhere



Yes

Fetch the data

This is usually a slow operation

**The novelty of ProbGraph**

We use BFs as a sketch of the actual dataset

How do we exactly use these sketches to benefit graph mining?

**Data** stored somewhere

# Observation: Set Intersection Cardinality is Prevalent in Graph Mining

# Observation: Set Intersection Cardinality is Prevalent in Graph Mining

$$|X \cap Y|$$

# Observation: Set Intersection Cardinality is Prevalent in Graph Mining



$$|X \cap Y|$$

# Observation: Set Intersection Cardinality is Prevalent in Graph Mining

# Observation: Set Intersection Cardinality is Prevalent in Graph Mining



$$|X \cap Y|$$

We greatly accelerate $|X \cap Y|$ with BFs

## ProbGraph key idea, continued

## ProbGraph key idea, continued



u

v

# ProbGraph key idea, continued

# ProbGraph key idea, continued



$N_u$

u

v

$N_v$

# ProbGraph key idea, continued



$N_u$

u

v

$N_v$

# ProbGraph key idea, continued



$N_u$

u

v

$N_v$

# ProbGraph key idea, continued



$N_u$

u

v

$N_v$

**Bitwise**

**AND**

# ProbGraph key idea, continued



$N_u$

u

$N_v$

v

Bitwise

AND

$$|N_u \cap N_v|$$

# Approximate Graph Processing: Our Objectives

| Little parallelism | → | Rich parallelism |
|---|---|---|
| Specific | → | Wide applicability |
| No/loose accuracy guarantees | → | Strong accuracy guarantees |
| Slow | → | High performance |
| Low accuracy | → | High accuracy |
| Large memory overheads | → | Low & controllable memory overheads |

# Approximate Graph Processing: Our Objectives

| Little parallelism | → | Rich parallelism |

| Specific | → | Wide applicability |

| No/loose accuracy guarantees | → | Strong accuracy guarantees |

| Slow | → | High performance |

| Low accuracy | → | High accuracy |

| Large memory overheads | → | Low & controllable memory overheads |

# ProbGraph: Fast & Parallel Execution

# ProbGraph: Fast & Parallel Execution



(different colors indicate different workers)

# ProbGraph: Fast & Parallel Execution



Embarrassingly parallel, O(1) depth

$N_u$

u

v

$N_v$

Bitwise

AND

(different colors indicate different workers)

$$|\widehat{N_u \cap N_v}|$$

# Approximate Graph Processing: Our Objectives

| Little parallelism | → | Rich parallelism |
|---|---|---|
| Specific | → | Wide applicability |
| No/loose accuracy guarantees | → | Strong accuracy guarantees |
| Slow | → | High performance |
| Low accuracy | → | High accuracy |
| Large memory overheads | → | Low & controllable memory overheads |

# Approximate Graph Processing: Our Objectives

| Little parallelism → | Rich parallelism ✓ |
| --- | --- |
| Specific → | Wide applicability |
| No/loose accuracy guarantees → | Strong accuracy guarantees |
| Slow → | High performance |
| Low accuracy → | High accuracy |
| Large memory overheads → | Low & controllable memory overheads |

# Approximate Graph Processing: Our Objectives

| Little parallelism | → | Rich parallelism ✓ |
| Specific | → | Wide applicability |
| No/loose accuracy guarantees | → | Strong accuracy guarantees |
| Slow | → | High performance |
| Low accuracy | → | High accuracy |
| Large memory overheads | → | Low & controllable memory overheads |

# Approximate Graph Processing: Our Objectives

| Little parallelism | → | Rich parallelism ✓ |
| Specific | → | Wide applicability |
| No/loose accuracy guarantees | → | Strong accuracy guarantees |
| Slow | → | High performance |
| Low accuracy | → | H... |
| Large memory overheads | → | Low memory overheads |

Let's see 4 example use cases...

# Use Case 1: Link Prediction

Which links will appear?

Which links are missing?

# Use Case 1: Link Prediction

Which links will appear?

Which links are missing?

# Use Case 1: Link Prediction

Which links will appear?

Which links are missing?

# Use Case 1: Link Prediction

**Which links will appear?**

**Which links are missing?**

**Fixing missing data**

**Predict future data**

# Use Case 1: Link Prediction

Which links will appear?

Which links are missing?

Fixing missing data

Predict future data

# Use Case 2: Clique Counting

# Use Case 2: Clique Counting

# Use Case 2: Clique Counting

# Use Case 2: Clique Counting

Learning over higher-order networks

...

B-relaxation Distance
$B_0$ $B_1$ $B_2$ $B_3$ $B_4$

# Use Case 3: Clustering

**# Clusters?**

**Structure of clusters?**

# Use Case 3: Clustering

# Clusters?

Structure of clusters?

# Use Case 3: Clustering

# Clusters?

Structure of clusters?

Minibatch selection in Graph Neural Networks

# Use Case 4: Vertex Similarity

# Use Case 4: Vertex Similarity

# Use Case 4: Vertex Similarity

Enhancing graph embedding construction

# Use Case 4: Vertex Similarity

Enhancing graph embedding construction

# Use Case 4: Vertex Similarity

Enhancing graph embedding construction

# Approximate Graph Processing: Our Objectives

| Little parallelism | → | Rich parallelism ✓ |
| Specific | → | Wide applicability |
| No/loose accuracy guarantees | → | Strong accuracy guarantees |
| Slow | → | High performance |
| Low accuracy | → | High accuracy |
| Large memory overheads | → | Low & controllable memory overheads |

# Approximate Graph Processing: Our Objectives

| | | |
|---|---|---|
| Little parallelism | → | Rich parallelism ✅ |
| Specific | → | Wide applicability ✅ |
| No/loose accuracy guarantees | → | Strong accuracy guarantees |
| Slow | → | High performance |
| Low accuracy | → | High accuracy |
| Large memory overheads | → | Low & controllable memory overheads |

# Approximate Graph Processing: Our Objectives

| Little parallelism | → | Rich parallelism ✓ |
| --- | --- | --- |
| Specific | → | Wide applicability ✓ |
| No/loose accuracy guarantees | → | Strong accuracy guarantees |
| Slow | → | High performance |
| Low accuracy | → | High accuracy |
| Large memory overheads | → | Low & controllable memory overheads |

# ProbGraph: Summary of Theoretical Results

# ProbGraph: Summary of Theoretical Results

We want guarantees for
$$|ProbGraphEstimate - exactResult|$$

# ProbGraph: Summary of Theoretical Results

We want guarantees for
$$|ProbGraphEstimate - exactResult|$$

We incorporate statistical theory of estimators

# ProbGraph is <u>asymptotically unbiased</u>

# ProbGraph is <u>asymptotically unbiased</u>

# ProbGraph is <u>asymptotically unbiased</u>

# ProbGraph is <u>asymptotically unbiased</u>

# ProbGraph is __asymptotically unbiased__

# ProbGraph is **asymptotically unbiased**

Zero average error at some point...
but the variance can still go wild



Computation result

The difference goes to <u>zero</u>

$E[ProbGraphEstimate]$

$exactResult$

ProbGraph sketch size (storage needed)

# ProbGraph is consistent

# ProbGraph is consistent



Computation result

One can <u>always</u> find a ProbGraph sketch that delivers a <u>required</u> accuracy

*exactResult*

ProbGraph sketch size (storage needed)

# ProbGraph is consistent



Computation result

*ProbGraphEstimate*

One can always find a ProbGraph sketch that delivers a required accuracy

*exactResult*

ProbGraph sketch size (storage needed)

# ProbGraph is consistent

The variance also converges to zero with the increasing sketch size



One can <u>always</u> find a ProbGraph sketch that delivers a <u>required</u> accuracy

Computation result

*ProbGraphEstimate*

*exactResult*

ProbGraph sketch size (storage needed)

# ProbGraph is <u>asymptotically efficient</u>



Computation result

exactResult

ProbGraph sketch size (storage needed)

# ProbGraph is <u>asymptotically efficient</u>

[1] J. Tetek, "Approximate triangle counting via sampling and fast matrix ˇ multiplication", arXiv 2021.

[2] S. Assadi et al., "A simple sublinear-time algorithm for counting arbitrary subgraphs via edge sampling", arXiv 2018.

[3] T. Eden et al., "Approximately counting triangles in sublinear time", SIAM Journal on Computing, 2017.

[4] B. Bandyopadhyay et al., "Topological graph sketching for incremental and scalable analytics", CIKM, 2016.

[5] R. Pagh et al., "Colorful triangle counting and a MapReduce implementation", Information Processing Letters, 2012.

[6] O. Papapetrou et al., "Cardinality estimation and dynamic length adaptation for bloom filters", Distributed and Parallel Databases, 2010.

[7] C. E. Tsourakakis et al., "Doulion: counting triangles in massive graphs with a coin". ACM KDD. 2009.

[8] M. Besta et al., "Slim graph: Practical lossy graph compression for approximate graph processing, storage, and analytics". ACM/IEEE SC. 2019

Computation result

Other estimators [1-8]

*exactResult*

ProbGraph sketch size (storage needed)

# ProbGraph is **asymptotically efficient**



Computation result

*ProbGraphEstimate*

Other estimators
[1-8]

*exactResult*

ProbGraph sketch size (storage needed)

[1] J. Tetek, "Approximate triangle counting via sampling and fast matrix ˇ multiplication", arXiv 2021.

[2] S. Assadi et al., "A simple sublinear-time algorithm for counting arbitrary subgraphs via edge sampling", arXiv 2018.

[3] T. Eden et al., "Approximately counting triangles in sublinear time", SIAM Journal on Computing, 2017.

[4] B. Bandyopadhyay et al., "Topological graph sketching for incremental and scalable analytics", CIKM, 2016.

[5] R. Pagh et al., "Colorful triangle counting and a MapReduce implementation", Information Processing Letters, 2012.

[6] O. Papapetrou et al., "Cardinality estimation and dynamic length adaptation for bloom filters", Distributed and Parallel Databases, 2010.

[7] C. E. Tsourakakis et al., "Doulion: counting triangles in massive graphs with a coin". ACM KDD. 2009.

[8] M. Besta et al., "Slim graph: Practical lossy graph compression for approximate graph processing, storage, and analytics". ACM/IEEE SC. 2019

# ProbGraph is **asymptotically efficient**

Computation result

*ProbGraphEstimate*

Other estimators

[1-8]

*exactResult*

No other consistent estimator has lower MSE / variance

ProbGraph sketch size (storage needed)

[1] J. Tetek, "Approximate triangle counting via sampling and fast matrix ˇ multiplication", arXiv 2021.

[2] S. Assadi et al., "A simple sublinear-time algorithm for counting arbitrary subgraphs via edge sampling", arXiv 2018.

[3] T. Eden et al., "Approximately counting triangles in sublinear time", SIAM Journal on Computing, 2017.

[4] B. Bandyopadhyay et al., "Topological graph sketching for incremental and scalable analytics", CIKM, 2016.

[5] R. Pagh et al., "Colorful triangle counting and a MapReduce implementation", Information Processing Letters, 2012.

[6] O. Papapetrou et al., "Cardinality estimation and dynamic length adaptation for bloom filters", Distributed and Parallel Databases, 2010.

[7] C. E. Tsourakakis et al., "Doulion: counting triangles in massive graphs with a coin". ACM KDD. 2009.

[8] M. Besta et al., "Slim graph: Practical lossy graph compression for approximate graph processing, storage, and analytics". ACM/IEEE SC. 2019

# ProbGraph has strong concentration bounds

# ProbGraph has <u>strong concentration bounds</u>



Deviation $t$ from the real value

# ProbGraph has <u>strong concentration bounds</u>



$P(|ProbGraphEstimate$

Deviation $t$ from the real value

# ProbGraph has strong concentration bounds

$P(|ProbGraphEstimate$

This probability decreases __exponentially__ fast

Deviation $t$ from the real value

# ProbGraph has strong concentration bounds

ProbGraph is unlikely to deviate much from the true values



$P(|ProbGraphEstimate$

This probability decreases **exponentially** fast

Deviation $t$ from the real value

# Approximate Graph Processing: Our Objectives

| Little parallelism | → | Rich parallelism ✓ |
| Specific | → | Wide applicability ✓ |
| No/loose accuracy guarantees | → | Strong accuracy guarantees |
| Slow | → | High performance |
| Low accuracy | → | High accuracy |
| Large memory overheads | → | Low & controllable memory overheads |

# Approximate Graph Processing: Our Objectives



| Little parallelism | → | Rich parallelism ✓ |
| Specific | → | Wide applicability ✓ |
| No/loose accuracy guarantees | → | Strong accuracy guarantees ✓ |
| Slow | → | High performance |
| Low accuracy | → | High accuracy |
| Large memory overheads | → | Low & controllable memory overheads |

# Approximate Graph Processing: Our Objectives



| Little parallelism | → | Rich parallelism ✓ |
| Specific | → | Wide applicability ✓ |
| No/loose accuracy guarantees | → | Strong accuracy guarantees ✓ |
| Slow | → | High performance |
| Low accuracy | → | High accuracy |
| Large memory overheads | → | Low & controllable memory overheads |

**Evaluation: Used Machines & Objectives**

# Evaluation: Used Machines & Objectives



CSCS Cray Piz Daint,
64 GB per compute node

# Evaluation: Used Machines & Objectives

**Dell PowerEdge R910 server**

**CSCS Cray Piz Daint,**
**64 GB per compute node**

# Evaluation: Used Machines & Objectives

**Goal**: One design with…
<u>large</u> speedups +
<u>small & controlled</u> accuracy loss +
<u>small & controlled</u> memory requirements

**CSCS Cray Piz Daint,**
**64 GB per compute node**

**Dell PowerEdge R910 server**

29

# Considered Graph Datasets

# Considered Graph Datasets

67 graph datasets,
15 areas,
5 major graph
dataset repositories

# Considered Graph Datasets

67 graph datasets,
15 areas,
5 major graph
dataset repositories

**Real-world graphs**

Purchases

Gene functions



Social networks

Communication

Brain structure



Web graphs

Road nets

Citation graphs

Compute graphs

**Synthetic graphs**

Kronecker [1]

Erdös-Rényi [2]

Mathematics

Medicine

Chemistry

Economic nets

[1] J. Leskovec et al. Kronecker Graphs: An Approach to Modeling Networks. J. Mach. Learn. Research. 2010.
[2] P. Erdos and A. Renyi. On the evolution of random graphs. Pub. Math. Inst. Hun. A. Science. 1960.

## Considered Graph Datasets

67

5

dataset repositories

**Real-world graphs**

Purchases

Gene functions

Communication

Brain structure

Citation graphs

**Synthetic graphs**

Compute graphs

Kronecker [1]        Erdös-Rényi [2]

Mathematics        Medicine        Chemistry        Economic nets

**Highly irregular data**

**Lots of load imbalance**

[1] J. Leskovec et al. Kronecker Graphs: An Approach to Modeling Networks. J. Mach. Learn. Research. 2010.
[2] P. Erdos and A. Renyi. On the evolution of random graphs. Pub. Math. Inst. Hun. A. Science. 1960.

# Triangle Counting

**Cores/threads**: 32
**Max memory overhead**: 20%

# Triangle Counting

**Cores/threads**: 32
**Max memory**
**overhead**: 20%

# Triangle Counting

**Cores/threads**: 32
**Max memory overhead**: 20%



**Speedup** over the exact tuned baseline

# Triangle Counting

**Cores/threads**: 32
**Max memory overhead**: 20%



Accuracy: Relative count of a given pattern (y-axis, values 0, 1, 2)

**Speedup** over the exact tuned baseline (x-axis, values 0, 20)

# Triangle Counting

**Cores/threads**: 32
**Max memory overhead**: 20%

**Each data point:** the execution of a given scheme for a specific graph dataset

**Accuracy:** Relative count of a given pattern

2

1

0

0          20

**Speedup** over the exact tuned baseline

# Triangle Counting

**Cores/threads**: 32
**Max memory overhead**: 20%

**Each data point:** the execution of a given scheme for a specific graph dataset

★ Exact baseline [1]

[1] S. Beamer et al., „The GAP Benchmark Suite". 2015



**Accuracy**: Relative count of a given pattern (y-axis: 0, 1, 2)

**Speedup** over the exact tuned baseline (x-axis: 01, 20)

# Triangle Counting

**Each data point:** the execution of a given scheme for a specific graph dataset

**Cores/threads**: 32
**Max memory overhead**: 20%

ProbGraph

Exact baseline [1]

[1] S. Beamer et al., „The GAP Benchmark Suite". 2015



**Accuracy**: Relative count of a given pattern

**Speedup** over the exact tuned baseline

# Triangle Counting

**Each data point:** the execution of a given scheme for a specific graph dataset

**Cores/threads**: 32
**Max memory overhead**: 20%

ProbGraph

★ Exact baseline [1]

■ Heuristics, no formal guarantees [2]

[1] S. Beamer et al., „The GAP Benchmark Suite". 2015

[2] S. Singh et al., "Scalable and performant graph processing on GPUs using approximate computing". IEEE TMSCS. 2018

**Accuracy:** Relative count of a given pattern

**Speedup** over the exact tuned baseline

# Triangle Counting

**Each data point:** the execution of a given scheme for a specific graph dataset

**Cores/threads**: 32
**Max memory overhead**: 20%



ProbGraph

Exact baseline [1]

Heuristics, no formal guarantees [2]

Heuristics, formal guarantees [3-4]

[1] S. Beamer et al., „The GAP Benchmark Suite". 2015

[2] S. Singh et al., "Scalable and performant graph processing on GPUs using approximate computing". IEEE TMSCS. 2018

[3] R. Pagh et al., "Colorful triangle counting and a mapreduce implementation". Information Processing Letters. 2012

[4] Z. Shang et al., "Auto-approximation of graph computing". VLDB. 2014

**Accuracy:** Relative count of a given pattern

**Speedup** over the exact tuned baseline

# Triangle Counting

**Cores/threads**: 32
**Max memory overhead**: 20%

**Each data point:** the execution of a given scheme for a specific graph dataset



ProbGraph

★ Exact baseline [1]

■ Heuristics, no formal guarantees [2]

✖ Heuristics, formal guarantees [3-4]

▲ Lossy graph compression [5-6]

[1] S. Beamer et al., „The GAP Benchmark Suite". 2015

[2] S. Singh et al., "Scalable and performant graph processing on GPUs using approximate computing". IEEE TMSCS. 2018

[3] R. Pagh et al., "Colorful triangle counting and a mapreduce implementation". Information Processing Letters. 2012

[4] Z. Shang et al., "Auto-approximation of graph computing". VLDB. 2014

[5] C. E. Tsourakakis et al., "Doulion: counting triangles in massive graphs with a coin". ACM KDD. 2009.

[6] M. Besta et al., "Slim graph: Practical lossy graph compression for approximate graph processing, storage, and analytics". ACM/IEEE SC. 2019

**Accuracy:** Relative count of a given pattern

**Speedup** over the exact tuned baseline

# Triangle Counting

**Cores/threads**: 32

**Max memory overhead**: 20%

**Each data point:** the execution of a given scheme for a specific graph dataset

● ProbGraph

★ Exact baseline [1]

■ Heuristics, no formal guarantees [2]

✖ Heuristics, formal guarantees [3-4]

▲ Lossy graph compression [5-6]

[1] S. Beamer et al., „The GAP Benchmark Suite". 2015

[2] S. Singh et al., "Scalable and performant graph processing on GPUs using approximate computing". IEEE TMSCS. 2018

[3] R. Pagh et al., "Colorful triangle counting and a mapreduce implementation". Information Processing Letters. 2012

[4] Z. Shang et al., "Auto-approximation of graph computing". VLDB. 2014

[5] C. E. Tsourakakis et al., "Doulion: counting triangles in massive graphs with a coin". ACM KDD. 2009.

[6] M. Besta et al., "Slim graph: Practical lossy graph compression for approximate graph processing, storage, and analytics". ACM/IEEE SC. 2019



For most graphs, we have...

...very high speedups

**Accuracy:** Relative count of a given pattern

**Speedup** over the exact tuned baseline

31

# Triangle Counting

**Cores/threads**: 32
**Max memory overhead**: 20%

**Each data point:** the execution of a given scheme for a specific graph dataset

🔵 ProbGraph

⭐ **Exact baseline [1]**

🟩 **Heuristics, no formal guarantees [2]**

❌ **Heuristics, formal guarantees [3-4]**

🔺 **Lossy graph compression [5-6]**

[1] S. Beamer et al., „The GAP Benchmark Suite". 2015

[2] S. Singh et al., "Scalable and performant graph processing on GPUs using approximate computing". IEEE TMSCS. 2018

[3] R. Pagh et al., "Colorful triangle counting and a mapreduce implementation". Information Processing Letters. 2012

[4] Z. Shang et al., "Auto-approximation of graph computing". VLDB. 2014

[5] C. E. Tsourakakis et al., "Doulion: counting triangles in massive graphs with a coin". ACM KDD. 2009.

[6] M. Besta et al., "Slim graph: Practical lossy graph compression for approximate graph processing, storage, and analytics". ACM/IEEE SC. 2019

For most graphs, we have...

...very high speedups

...very good accuracy

80% accuracy

**Accuracy:** Relative count of a given pattern

**Speedup** over the exact tuned baseline

31

# Triangle Counting

**Each data point:** the execution of a given scheme for a specific graph dataset

**Cores/threads**: 32
**Max memory overhead**: 20%

● ProbGraph

★ Exact baseline [1]

■ Heuristics, no formal guarantees [2]

✖ Heuristics, formal guarantees [3-4]

▲ Lossy graph compression [5-6]

[1] S. Beamer et al., „The GAP Benchmark Suite". 2015

[2] S. Singh et al., "Scalable and performant graph processing on GPUs using approximate computing". IEEE TMSCS. 2018

[3] R. Pagh et al., "Colorful triangle counting and a mapreduce implementation". Information Processing Letters. 2012

[4] Z. Shang et al., "Auto-approximation of graph computing". VLDB. 2014

[5] C. E. Tsourakakis et al., "Doulion: counting triangles in massive graphs with a coin". ACM KDD. 2009.

[6] M. Besta et al., "Slim graph: Practical lossy graph compression for approximate graph processing, storage, and analytics". ACM/IEEE SC. 2019

For most graphs, we have...
...very high speedups
...very good accuracy
...mild memory requirements

**80% accuracy**

**Accuracy:** Relative count of a given pattern

**Speedup** over the exact tuned baseline

31

# 4-Clique Counting

**Cores/threads**: 32
**Max memory overhead**: 20%

**Each data point:** the execution of a given scheme for a specific graph dataset

★ Exact baseline [1]

● ProbGraph

For most graphs, we have...

...very high speedups

...very good accuracy

...mild memory requirements



[1] based on S. Beamer et al., „The GAP Benchmark Suite". 2015

# Clustering (Scaling)

**Max memory overhead**: 20%



Legend:
- Exact scheme [1]
- Heuristics [2]
- Compression [3-4]
- ProbGraph (A)
- ProbGraph (B)

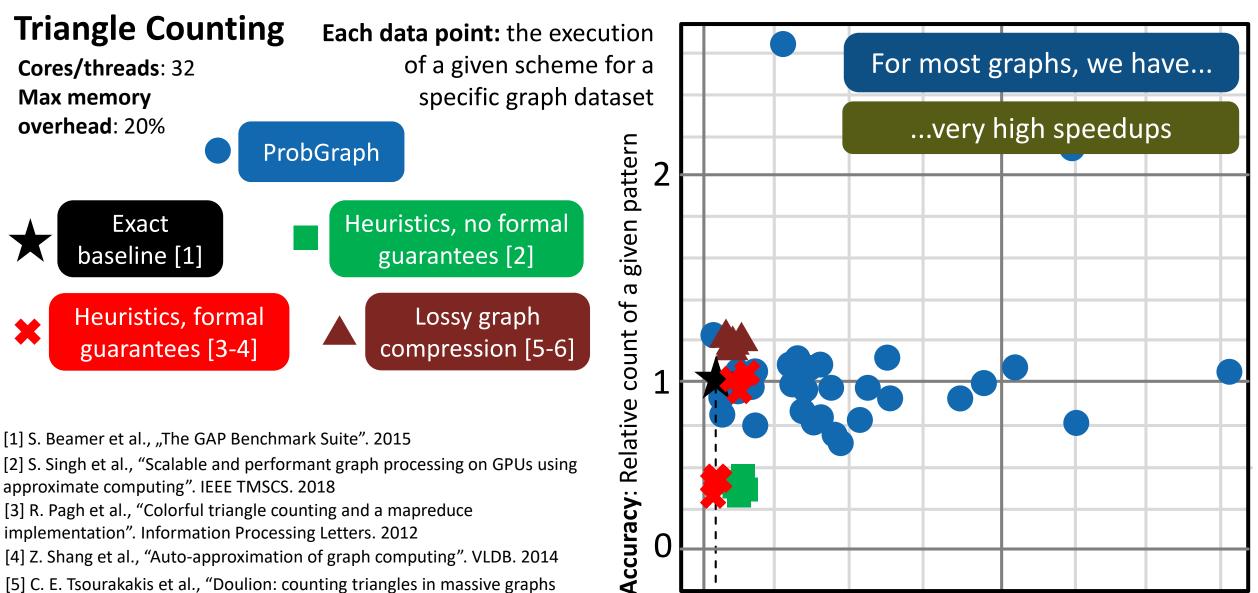[1] S. Beamer et al., „The GAP Benchmark Suite". 2015
[2] R. Pagh et al., "Colorful triangle counting and a mapreduce implementation". Information Processing Letters. 2012
[3] C. E. Tsourakakis et al., "Doulion: counting triangles in massive graphs with a coin". ACM KDD. 2009.
[4] M. Besta et al., "Slim graph: Practical lossy graph compression for approximate graph processing, storage, and analytics". ACM/IEEE SC. 2019

# Clustering (Scaling)

**Max memory overhead**: 20%



Legend:
- Exact scheme [1]
- Heuristics [2]
- Compression [3-4]
- ProbGraph (A)
- ProbGraph (B)

X-axis: # Workers ($2^0$, $2^1$, $2^2$, $2^3$, $2^4$, $2^5$)

Y-axis: Runtime [s] ($10^0$, $10^1$, $10^2$, $10^3$)

[1] S. Beamer et al., „The GAP Benchmark Suite". 2015
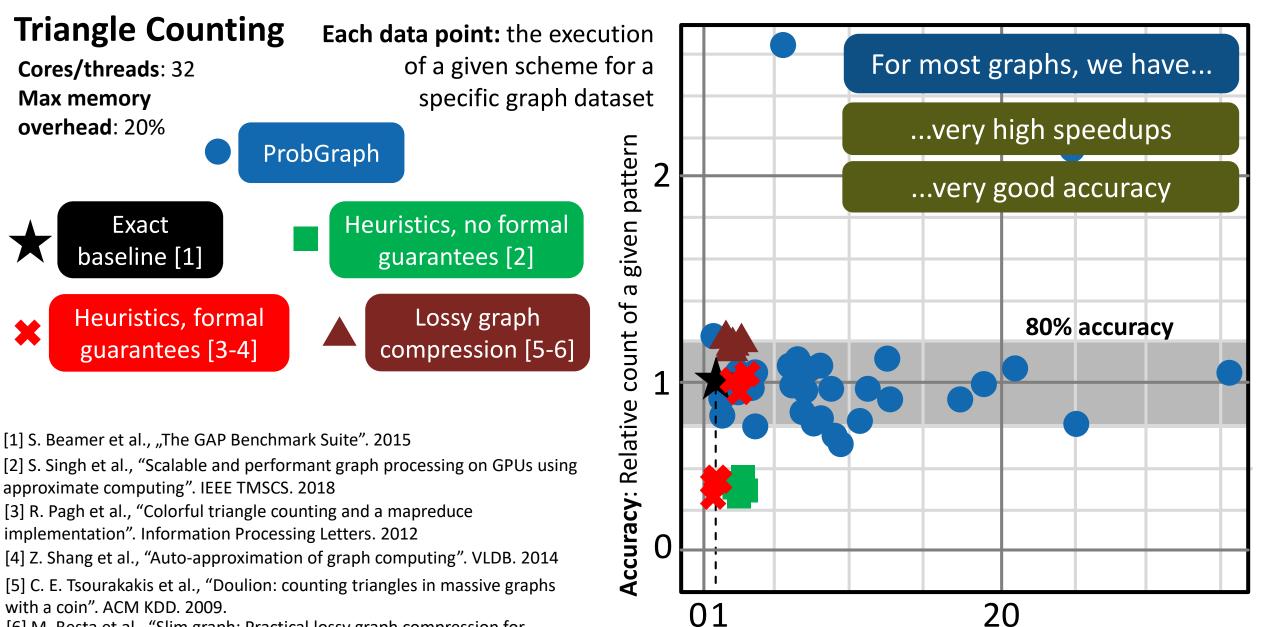[2] R. Pagh et al., "Colorful triangle counting and a mapreduce implementation". Information Processing Letters. 2012
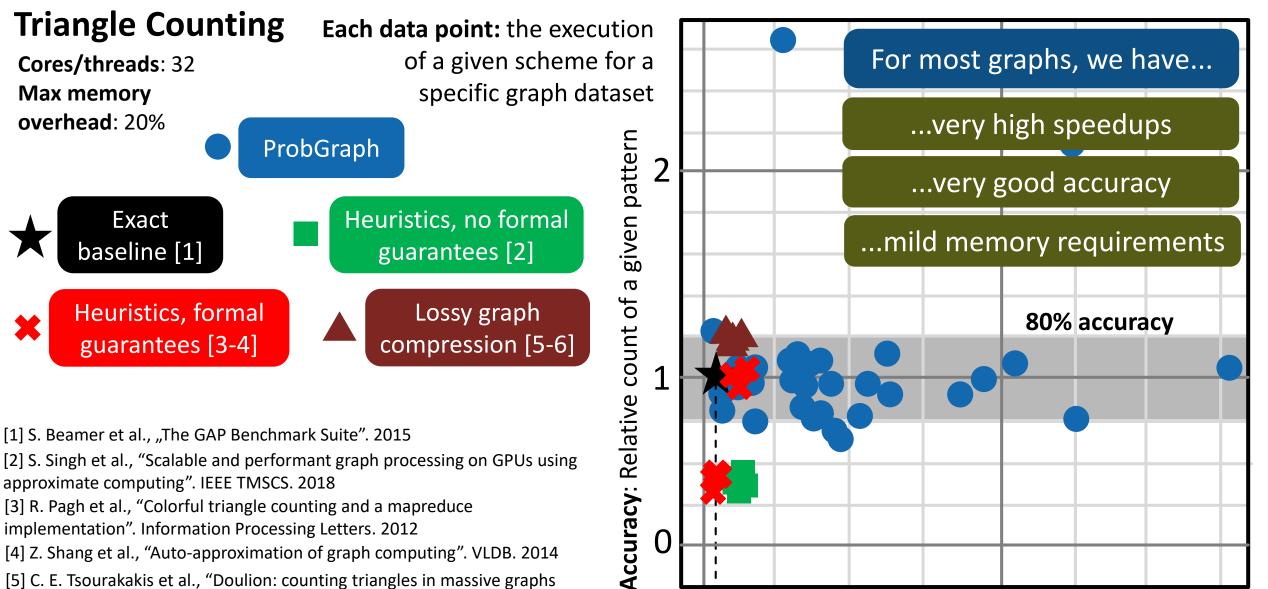[3] C. E. Tsourakakis et al., "Doulion: counting triangles in massive graphs with a coin". ACM KDD. 2009.
[4] M. Besta et al., "Slim graph: Practical lossy graph compression for approximate graph processing, storage, and analytics". ACM/IEEE SC. 2019

# Clustering (Scaling)

**Max memory overhead**: 20%

[1] S. Beamer et al., „The GAP Benchmark Suite". 2015
[2] R. Pagh et al., "Colorful triangle counting and a mapreduce implementation". Information Processing Letters. 2012
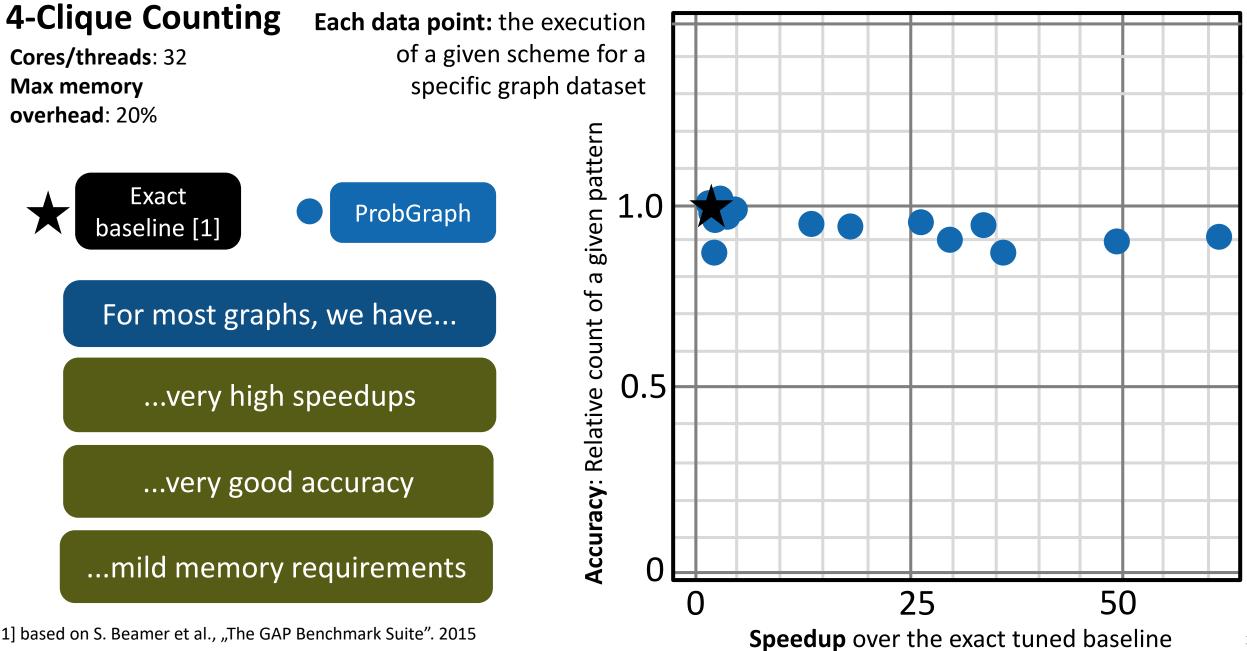[3] C. E. Tsourakakis et al., "Doulion: counting triangles in massive graphs with a coin". ACM KDD. 2009.
[4] M. Besta et al., "Slim graph: Practical lossy graph compression for approximate graph processing, storage, and analytics". ACM/IEEE SC. 2019

Legend:
- Exact scheme [1]
- Heuristics [2]
- Compression [3-4]
- ProbGraph (A)
- ProbGraph (B)

Runtime [s]

Graph densities

$2^2$  $2^4$  $2^6$  $2^8$  $2^{10}$  $2^{12}$

$2^0$  $2^1$  $2^2$  $2^3$  $2^4$  $2^5$

# Workers

# Clustering (Scaling)

**Max memory overhead**: 20%



**Graph densities**

$2^2$   $2^4$   $2^6$   $2^8$   $2^{10}$   $2^{12}$

$2^0$   $2^1$   $2^2$   $2^3$   $2^4$   $2^5$

# Workers

- Exact scheme [1]
- Heuristics [2]
- Compression [3-4]
- ProbGraph (A...
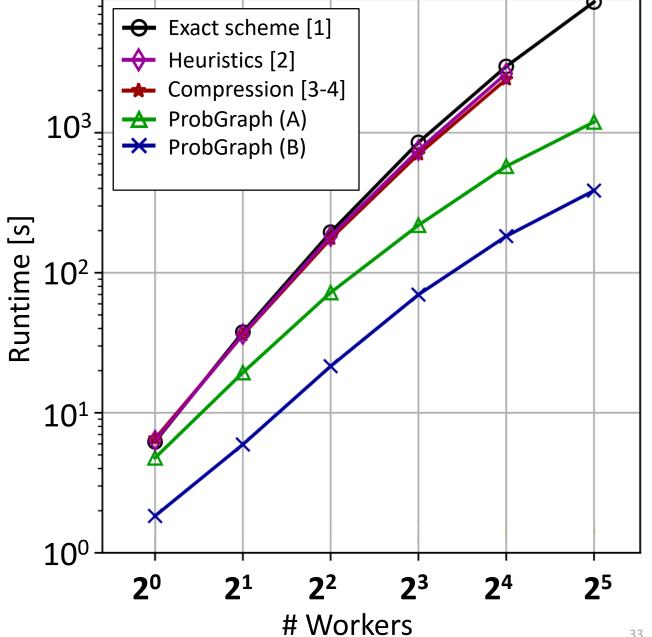- ProbGraph (...

$10^1$

$10^0$

[1] S. Beamer et al., „The GAP Benchmark Suite". 2015
[2] R. Pagh et al., "Colorful triangle counting and a mapreduce implementation". Information Processing Letters. 2012
[3] C. E. Tsourakakis et al., "Doulion: counting triangles in massive graphs with a coin". ACM KDD. 2009.
[4] M. Besta et al., "Slim graph: Practical lossy graph compression for approximate graph processing, storage, and analytics". ACM/IEEE SC. 2019

# Clustering (Scaling)

**Max memory overhead**: 20%



Exact scheme [1]
Heuristics [2]
Compression [3-4]
ProbGraph (A)
ProbGraph (B)

**Graph densities**

$2^2$    $2^4$    $2^6$    $2^8$    $2^{10}$    $2^{12}$

$10^1$

$10^0$

$2^0$    $2^1$    $2^2$    $2^3$    $2^4$    $2^5$

**# Workers**

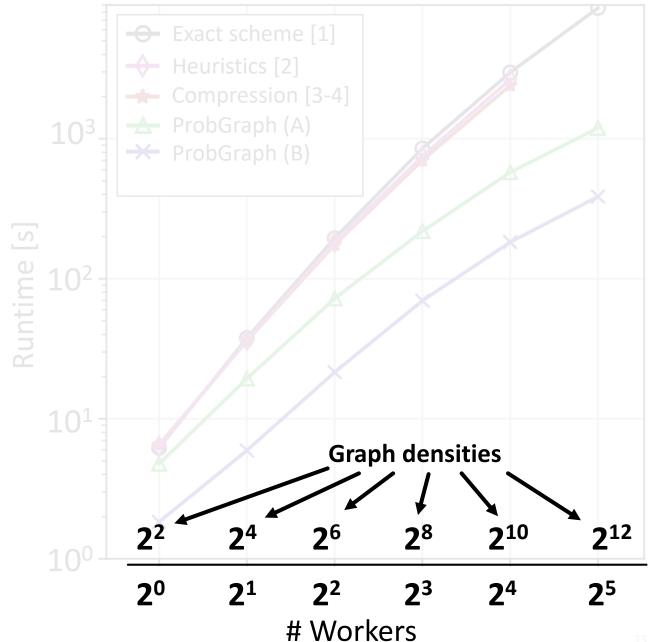[1] S. Beamer et al., „The GAP Benchmark Suite". 2015
[2] R. Pagh et al., "Colorful triangle counting and a mapreduce implementation". Information Processing Letters. 2012
[3] C. E. Tsourakakis et al., "Doulion: counting triangles in massive graphs with a coin". ACM KDD. 2009.
[4] M. Besta et al., "Slim graph: Practical lossy graph compression for approximate graph processing, storage, and analytics". ACM/IEEE SC. 2019
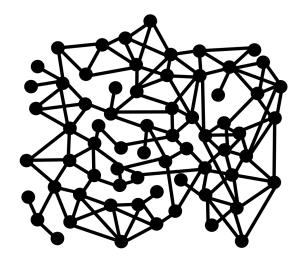
# Clustering (Scaling)

**Max memory overhead**: 20%

[1] S. Beamer et al., „The GAP Benchmark Suite". 2015

[2] R. Pagh et al., "Colorful triangle counting and a mapreduce implementation". Information Processing Letters. 2012

[3] C. E. Tsourakakis et al., "Doulion: counting triangles in massive graphs with a coin". ACM KDD. 2009.

[4] M. Besta et al., "Slim graph: Practical lossy graph compression for approximate graph processing, storage, and analytics". ACM/IEEE SC. 2019
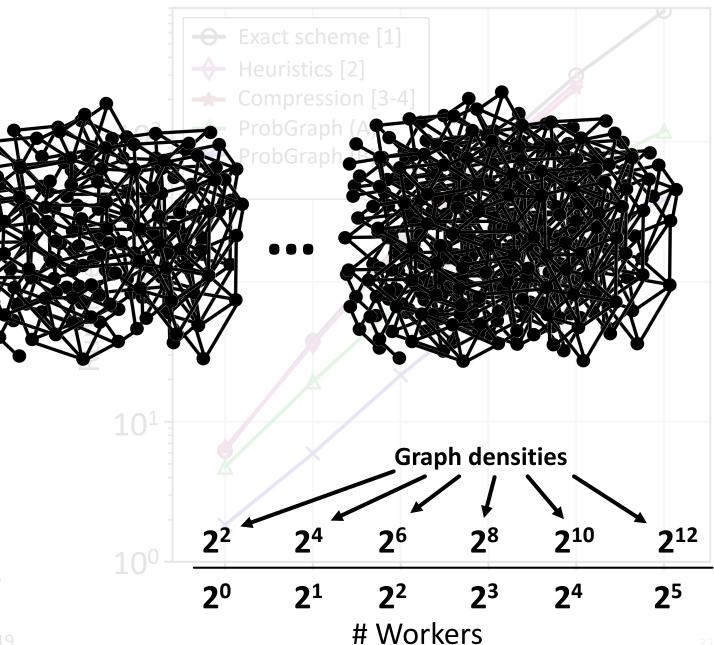
33

# Clustering (Scaling)

**Max memory overhead**: 20%

Why do we scale so well?



Graph densities

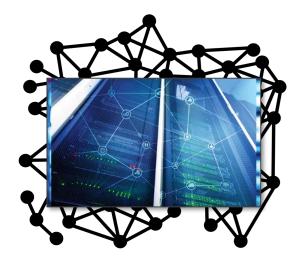[1] S. Beamer et al., „The GAP Benchmark Suite". 2015

[2] R. Pagh et al., "Colorful triangle counting and a mapreduce implementation". Information Processing Letters. 2012

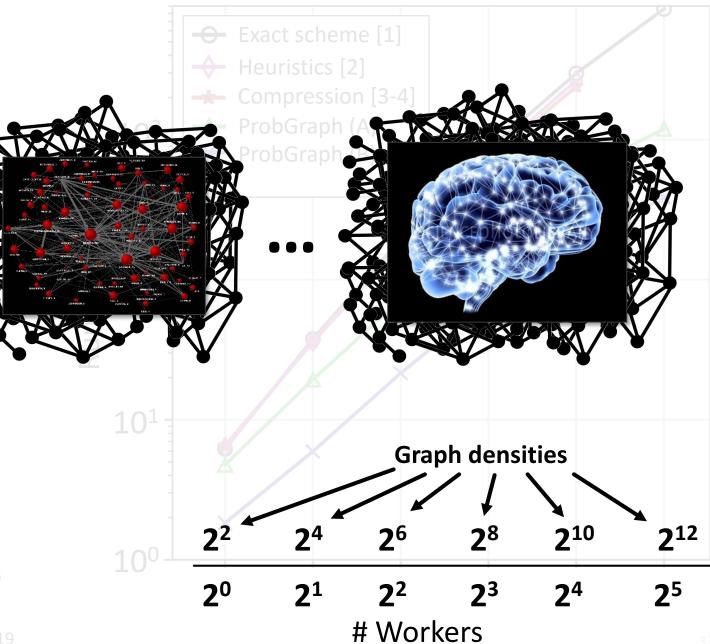[3] C. E. Tsourakakis et al., "Doulion: counting triangles in massive graphs with a coin". ACM KDD. 2009.

[4] M. Besta et al., "Slim graph: Practical lossy graph compression for approximate graph processing, storage, and analytics". ACM/IEEE SC. 2019

# Clustering (Scaling)

**Max memory overhead**: 20%

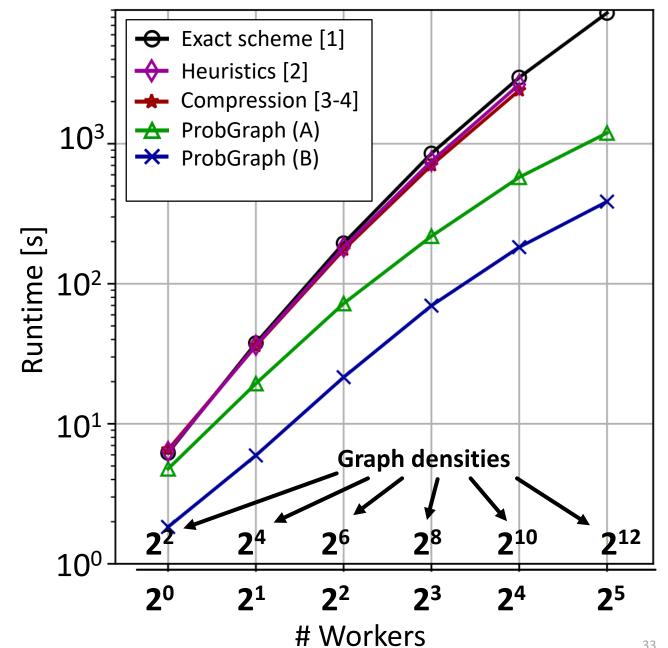Why do we scale so well?

Great load balancing properties

[1] S. Beamer et al., „The GAP Benchmark Suite". 2015
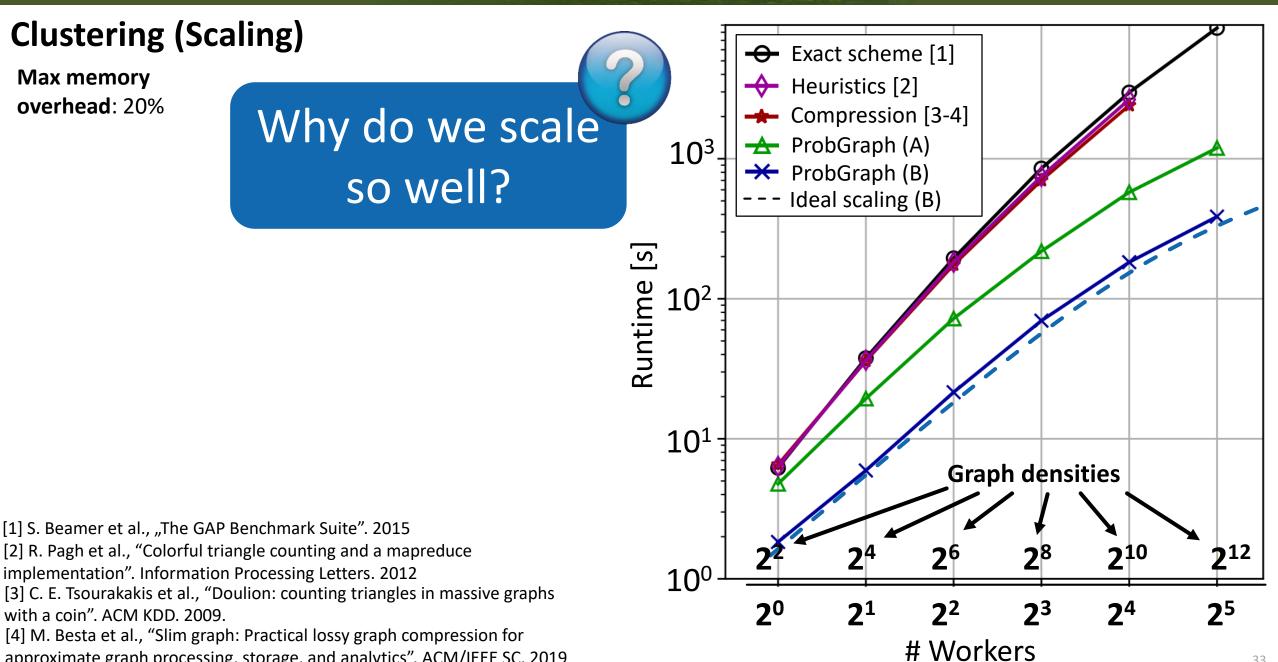[2] R. Pagh et al., "Colorful triangle counting and a mapreduce implementation". Information Processing Letters. 2012
[3] C. E. Tsourakakis et al., "Doulion: counting triangles in massive graphs with a coin". ACM KDD. 2009.
[4] M. Besta et al., "Slim graph: Practical lossy graph compression for approximate graph processing, storage, and analytics". ACM/IEEE SC. 2019
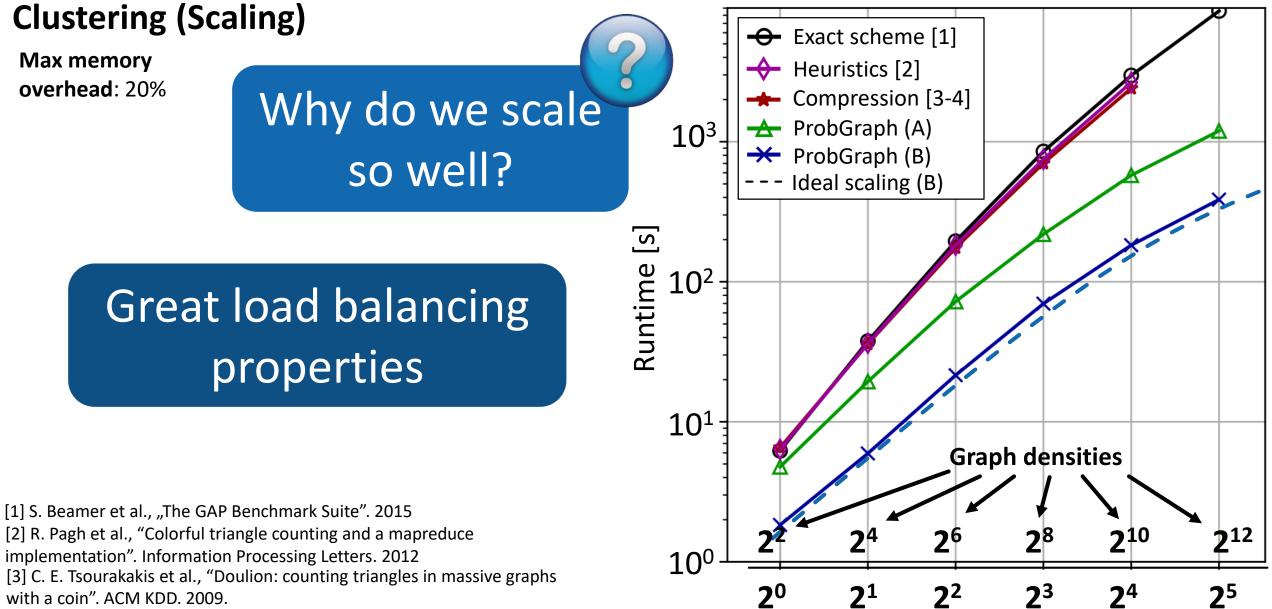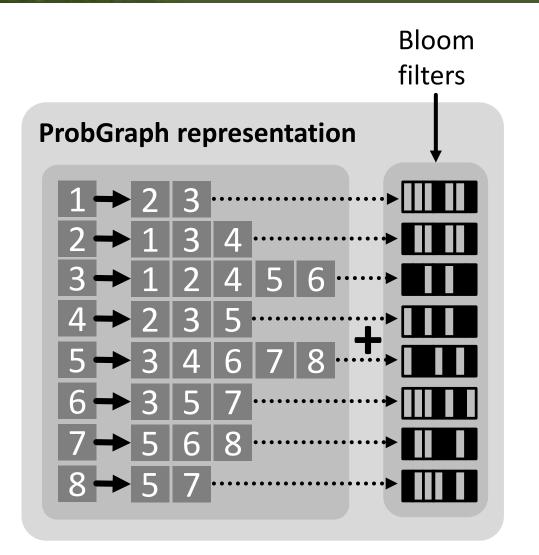
Bloom filters
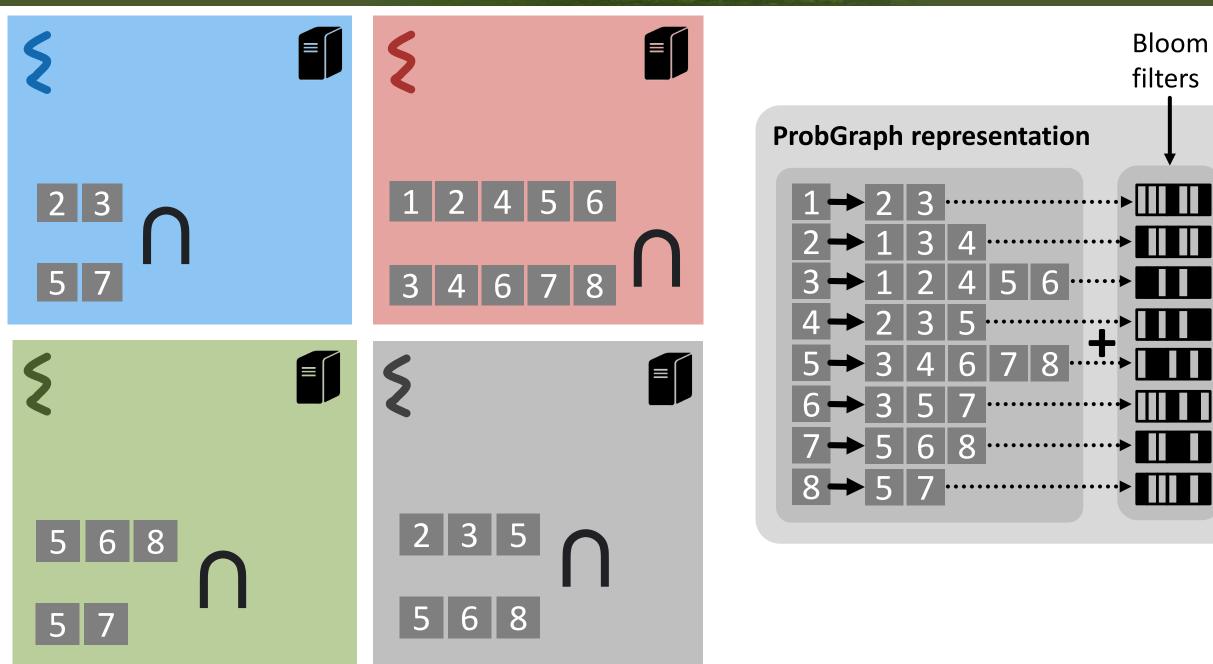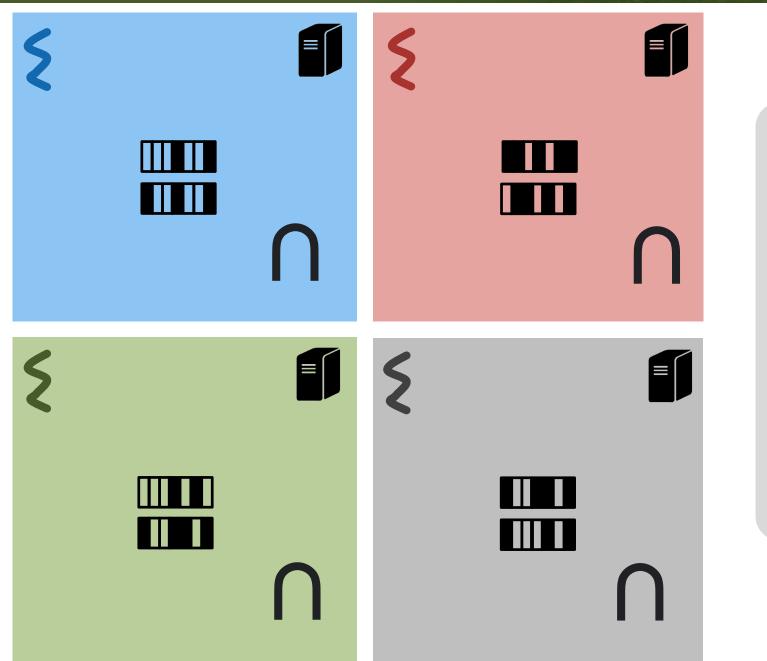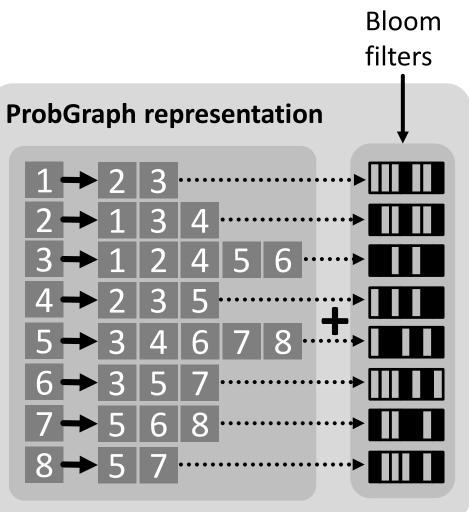
ProbGraph representation

ProbGraph representation

Bloom filters
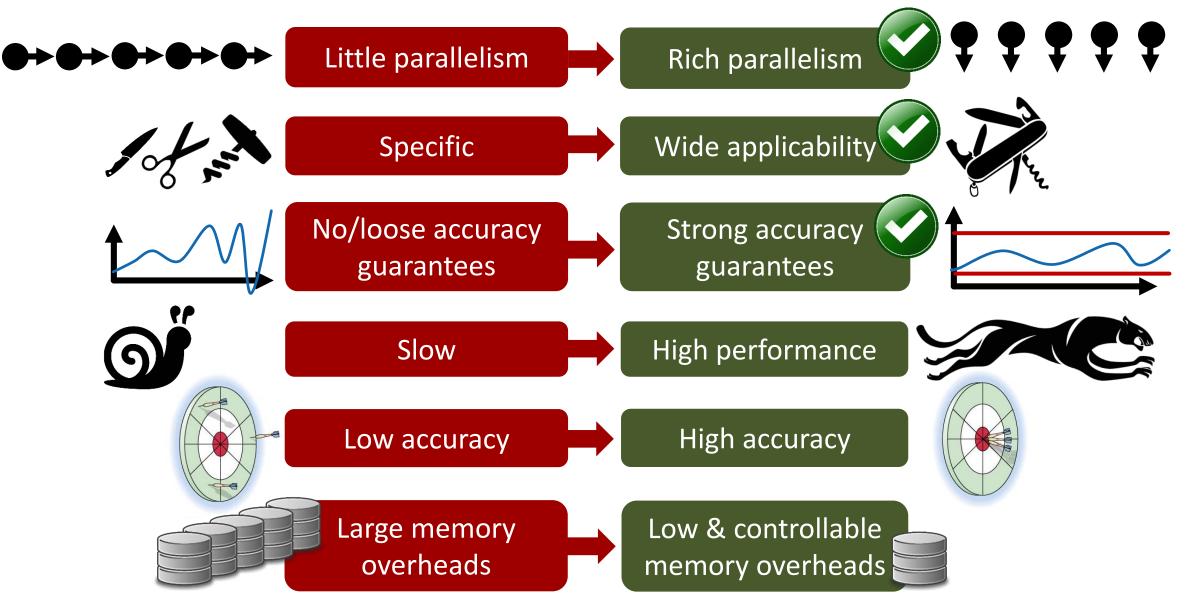
# ...Many more data & a lot of strong theory results!

$$P\left(\left|TC - \widehat{TC}_{1H}\right| \geq t\right) \leq 2\exp\left(-\frac{18\,k\,t^2}{\left(\sum_{v\in V} d(v)^2\right)^2}\right)$$

| Result | Where | Clas |
|---|---|---|
| $\widehat{|X|}_S$ | Eq. (1) | BF |
| $\widehat{|X\cap Y|}$ | | |
| $\widehat{|X\cap Y|}$ | | |
| $\widehat{|X\cap Y|}$ | | |

**Theorem A.6.** Let $Y_1 = $ and assume we partition $\mathcal{X}_1, \cdots, \mathcal{X}_\chi$ such that estim independent. Then for any $S = \sum_{i=1}^n C_i \frac{J}{1+J}$

| | | (work) | (depth) |
|---|---|---|---|
| | | $O(bd_v)$ | $O(\log(bd_v))$ |
| | | $O(kd_v)$ | $O(\log d_v)$ |
| | | **PG (BF)** | **PG (MH)** |
| | | $\frac{ndB_X}{W}$ | $O(ndk)$ |
| | | $\log\left(\frac{B_X}{W}\right)$ | $O(\log k)$ |
| | | $\frac{nd^2 B_X}{W}$ | $O(nd^2 k)$ |
| | | $\log d\log\left(\frac{B_X}{W}\right)$ | $O(\log^2 k)$ |

| Result | | |
|---|---|---|
| $\widehat{|X|}_S$ ★ | $P(|Y_1 - S| > t), P(|Y_k - S| >$ | |
| $|X\cap Y|_{AN}$ | | |
| $\widehat{|X\cap Y|}_L$ | | |
| $\widehat{|X\cap Y|}_{kH}$ | | |
| $\widehat{|X\cap Y|}_{1H}$ ★ | Eq. (7) | 1-Hash |

$$E[(\widehat{|X|} - |X|)^2] \tag{8}$$

$$= E[(\widehat{|X|} - |X|)^2 | \mathcal{E}]P(\mathcal{E}) + E[(\widehat{|X|} - |X|)^2 | \neg\mathcal{E}]P(\neg\mathcal{E}) \tag{9}$$

$$\leq (1+\varepsilon)E[(\widehat{|X|} - \kappa)^2 | \mathcal{E}] + \frac{1+\varepsilon}{\varepsilon}E[(\kappa - |X|)^2 | \mathcal{E}] + O(B_X^2 \log^2 B_X)\cdot\exp(-B_X^{\Omega(1)}) \tag{10}$$

$$\leq \frac{(1+\varepsilon)B_X^2}{b^2}E[(\log(B_{X,0}/B_X) - \log(1-1/B_X)^{b|X|})^2 | \mathcal{E}] + O((\kappa - |X|)^2) + \exp(-B_X^{\Omega(1)}) \tag{11}$$

$$\leq \frac{(1+\varepsilon)B_X^2}{b^2}E[(\log(B_{X,0}/B_X) - \log(1-1/B_X)^{b|X|})^2 | \mathcal{E}] + O(|X|/B_X) \tag{12}$$

$$\leq \frac{(1+\varepsilon)^2 B_X^2}{b^2}e^{2b|X|/B_X}E[(B_{X,0}/B_X - (1-1/B_X)^{b|X|})^2 | \mathcal{E}] + O(|X|/B_X) \tag{13}$$

$$\leq \frac{(1+\varepsilon)^2 B_X^2}{b^2}e^{2b|X|/(B_X-1)}\cdot E[(B_{X,0}/B_X - (1-1/B_X)^{b|X|})^2]/P[\mathcal{E}] + O(|X|/B_X) \tag{14}$$

$$= ((1+\varepsilon)^2 + o(1))\frac{B_X^2}{b^2}e^{2b|X|/(B_X-1)}\cdot E[(B_{X,0}/B_X - (1-1/B_X)^{b|X|})^2] + O(|X|/B_X) \tag{15}$$

$$= ((1+\varepsilon)^2 + o(1))\frac{e^{2b|X|/(B_X-1)}}{b^2}Var[B_{X,0}] + O(|X|/B_X) \tag{16}$$

$$\leq ((1+\varepsilon)^2 + o(1))e^{2b|X|/(B_X-1)}\cdot\left(e^{-\frac{b|X|}{B_X}}\frac{B_X}{b^2} - B_X/b^2 - |X|/b\right) + O(|X|/B_X) \tag{17}$$

$$\leq ((1+\varepsilon)^2 + o(1))\left(e^{|X|b/(B_X-1)}\frac{B_X}{b^2} - B_X/b^2 - |X|/b\right) + O(|X|/B_X) \tag{18}$$

$$\leq ((1+\varepsilon)^2 + o(1))\left(e^{|X|b/(B_X-1)}\frac{B_X}{b^2} - B_X/b^2 - |X|/b\right) \tag{19}$$

| | | | |
|---|---|---|---|
| | | | (30) |
| | | $|X||_i[E(\widehat{|X|}_j) - |X||_j]$ | (31) |
| | | $|X||_i\left||E(\widehat{|X|}_j) - |X||_j\right|$ | (32) |
| | | | (33) |
| | | $|_j) - |X||_j]^2$ | (35) |
| | | | (36) |
| | | $-\frac{2\Delta}{b}$ | (34) |
| | | | (37) |
| | | | (38) |

| Reference | Constr. time | Memory used |
|---|---|---|
| Doulion [46] | $O(m)$ | $O(pm)$ |
| Colorful [47] | $O(m)$ | $O(pm)$ |
| Sketching [48] | $O(km)$ | $O(kn)$ |
| ASAP [49] | n/a | $O(n+m)$ |
| GAP [50] | $O(m)$† | $O(m')$† |
| Slim Gr. [51] | $O(m)$ | $O(pm)$ |
| Eden et al. [52] | n/a | $O\left(\frac{n}{TC^{1/3}}\right)$ |
| Assadi et al. [53] | n/a | $O(1)$ |
| Tětek [54] | n/a | $\left(\frac{m^{1.41}}{TC^{0.82}}\right)$ |

| | | | |
|---|---|---|---|
| $\widehat{TC}_{AND}$ (BF) | $O(bm)$ | $O(n+m)$ |
| $\widehat{TC}_{kH}$ (MH) | $O(km)$ | $O(n+m)$ |
| $\widehat{TC}_{1H}$ (MH) | $O(km)$ | $O(n+m)$ |

**Relative difference:**

**Past results**

**PG**

**CSR (merge)**

**Work:** $O(d_u + d_v)$

**Depth:** $O(\log(d_u + d_v$

$$P\left(\left|TC - \widehat{T}\right.\right.$$

$$\frac{9\,t^2}{}$$

Number of Threads
$2^0$ $2^1$ $2^2$ $2^3$

$$P\left(\left|TC\right.\right.$$

$$\frac{B_X}{b^2} - \frac{2\Delta}{b}$$

$$\frac{}{4\,(\Delta+1)\sum_{v\in V}d(v)^3}$$

# Approximate Graph Processing: Our Objectives

| | |
|---|---|
| Little parallelism | Rich parallelism ✓ |
| Specific | Wide applicability ✓ |
| No/loose accuracy guarantees | Strong accuracy guarantees ✓ |
| Slow | High performance |
| Low accuracy | High accuracy |
| Large memory overheads | Low & controllable memory overheads |

# Approximate Graph Processing: Our Objectives



| Little parallelism | → | Rich parallelism ✓ |
| Specific | → | Wide applicability ✓ |
| No/loose accuracy guarantees | → | Strong accuracy guarantees ✓ |
| Slow | → | High performance ✓ |
| Low accuracy | → | High accuracy ✓ |
| Large memory overheads | → | Low & controllable memory overheads ✓ |

# Conclusion: ProbGraph Enables Approximate Graph Mining with...

Rich parallelism

Wide applicability

Strong accuracy guarantees

High performance

High accuracy

Low & controllable memory overheads

# Conclusion: ProbGraph Enables Approximate Graph Mining with...

- Rich parallelism ✓
- Wide applicability ✓
- Strong accuracy guarantees ✓
- High performance ✓
- High accuracy ✓
- Low & controllable memory overheads ✓

# Conclusion: ProbGraph Enables Approximate Graph Mining with...

Rich parallelism

Wide applicability

Strong accuracy guarantees

High performance

High accuracy

Low & controllable memory overheads

# Thank you

**Want to know more?**

youtube.com/@spcl

twitter.com/spcl_eth

spcl.inf.ethz.ch

github.com/spcl