## PhD Dissertation Defense: Systems Abstractions for Big Data Processing on a Single Machine

Steffen Maass

Advisor: Taesoo Kim

April 3, 2019

## Large-scale data processing is ubiquitous

#### One Trillion Edges: Graph Processing at Facebook-Scale

Avery Ching Facebook 1 Hacker Lane Menlo Park, California aching@fb.com Sergey Edunov Facebook 1 Hacker Lane Menlo Park, California edunov@fb.com Maja Kabiljo Facebook 1 Hacker Lane Menlo Park, California majakabiljo@fb.com

Dionysios Logothetis Facebook 1 Hacker Lane Menlo Park, California dionysios@fb.com Sambavi Muthukrishnan Facebook 1 Hacker Lane Menlo Park, California sambavim@fb.com

#### ABSTRACT

Analyzing large graphs provides valuable insights for oxial networking and web comparison is content ranking and recommendations. While numerous graph processing systems graphs of up to 608 edges, they often face significant difficulties in scaling to much large graphs. Industry graphs can be two order of magnitude larger - hundreds of bhlizons or up to one trillion edges. In addition to scalability combine graph processing workflows than previously evalu-

#### Social networks

a project to run Facebook-scale graph applications in the summer of 2012 and is still the case today.

Table 1: Popular benchmark graphs.

Graph	Vertices	Edges
LiveJournal [9]	4.8M	69M
Twitter 2010 [31]	42M	1.5B
UK web graph 2007 [10]	109M	3.7B
Yahoo web [8]	1.4B	6.6B

## Large-scale data processing is ubiquitous

#### One Trillion Edges:



Dionysios Fac 1 Hac Menlo Pa dionysic

#### ABSTRACT

Analyzing large graphs provides valuable i networking and web companies in content ommendations. While numerous graph pr have been developed and evaluated on avagraphs of up to 6.6B edges, they often fac ficulties in scaling to much larger graphs. can be two orders of magnitude larger lions or up to one trillion edges. In addit challenges, real world applications often re complex graph processing workflows than

#### Social networks



Home Current issue News & comment Research Archive V Authors & referees V About the journal

home > archive > issue > computational biology > primer > full text

NATURE BIOTECHNOLOGY | COMPUTATIONAL BIOLOGY | PRIMER

< 🔒

#### How to apply de Bruijn graphs to genome assembly

Phillip E C Compeau, Pavel A Peyzner & Glenn Tesler

Affiliations | Corresponding author

Nature Biotechnology 29, 987-991 (2011) | doi:10.1038/nbt.2023 Published online 08 November 2011



A mathematical concept known as a de Bruijn graph turns the formidable challenge of assembling a contiguous genome from billions of short seguencing reads into a tractable computational problem.

#### Genome analysis

## Large-scale data processing is ubiquitous

#### **One Trillion Edges:**

Avery Ching Facebook 1 Hacker Lane Menlo Park, California aching@fb.com

Dionysios Fac 1 Hac Menlo Pa dionysic

#### ABSTRACT

Analyzing large graphs provides valuable i networking and web companies in content ormnerolations. While numerous graph pr have been developed and evaluated on ava graphs of up to 6.6B edges, they often far ficulties in scaling to much larger graphs. can be two orders of magnitude larger -lions or up to one trillion edges. In addit challenges, real world applications often re complex ranho processing worldworks than

#### Social networks



#### Home | Current issue | News & cor

home > archive > issue > computa Graph-powered Machine Learning at Google

NATURE BIOTECHNOLOGY | CO

# How to apply de E assembly

#### Phillip E C Compeau, Pavel A P

#### Affiliations | Corresponding aut

Nature Biotechnology 29, 987–99 Published online 08 November 20



A mathematical concept known assembling a contiguous genor computational problem.

### Genome analysi



image recognition system in Google Photos.

Thursday, October 06, 2016

Posted by Suiith Ravi, Staff Research Scientist, Google Research



Recently, there have been significant advances in Machine Learning that enable computer systems

to solve complex real-world problems. One of those advances is Google's large scale, graph-based

machine learning platform, built by the Expander team in Google Research. A technology that is

behind many of the Google products and features you may use everyday, graph-based machine

learning is a powerful tool that can be used to power useful features such as reminders in Inbox

and smart messaging in Allo, or used in conjunction with deep neural networks to power the latest



#### Learning with Minimal Supervision

### Graphs enable Machine Learning

Steffen Maass

Abstractions for Big Data Processing

April 3, 2019 2 / 46



#### More sockets. More memory. More SAP HANA.

by Cori Pasinetti on July 29, 2015



SGI UV 300H 20-Socket Appliance Certified by SAP to Run SAP HANA® Under Controlled Availability Announcing the first 20-socket SAP HANA-certified in-memory server!

SGI announced today that the SGIB UV<sup>™</sup> 300H is now SAPB-entified to run the SAP HANAB platform in controlled availability at 20-acckets-delivering up to 15 terabytes (TB) of in-memory computing capacity in a single node. Assetting the value of key enhancements in support package stack. 10 (SPSID) for SAP HANA and SAP's doce collaboration with system providers, SGI UV 300H delivers outstanding single-node performance and simplicity for entryrises moving to SAP HANA to gain business breakthroughs.

SGI UV 3000 is a specialized offering in the SGI8 UV<sup>®</sup> event line for in-merry comparing the tables in the server line for in-merry comparing the tables matching boot involution, and lower IT costs with S40 advantages for business running S40% advantages for business advantages for business running S40% advantess Sufe advantages for business running S40% advantess advantages for business running S40% advantess advantages for business running S40% advantess advantages for business running server history advantages advantages for advantages advantag

Integrated with the recently announced SAP HANA SP510, SGI UV 300H capitalizes on deep collaboration between SAP, Intel and SGI to optimize SAP HANA-based workloads on multicore NUMA (non-uniform memory More sockets, More memory, More SAP HANA Now SQI offers industry-leading controlled availability – up to 20-sockets!

access) systems. This enables enterprises to leverage single-node systems with very large memory capacity and

### Gigabytes of RAM

Learning



The Trusted Leader in High Performance Computing

Intel Unveils Plans for Knights Mill, a Xeon Phi for Deep

#### More sockets, More

by Corl Pasinetti on July 29, 2015 Michael Feldman | August 18, 2016 01:33 CEST

Tweet 0 🖬 Share 0

SGI UV 300H 20-Socket Applian Announcing the first 20-socket

At the Intel Developer Forum (IDF) this week in San Francisco, Intel revealed it is working on a new Xeon Phi processor aimed at deep learning applications, Diane Bryant, executive VP and GM of Intel's Data Center Group, unveiled the new chip, known SGI announced today that the SGI as Knights Mill, during her IDF keynote address on Wednesday.

controlled availability at 20-sockets single node. Asserting the value of SAP's close collaboration with syste simplicity for enterprises moving to

SGI UV 300H is a specialized offering server line for in-memory computin enterprises to further unlock value real-time, boost innovation, and low HANA. Featuring a highly differentia architecture, the system delivers sig advantages for businesses running 4 SAP HANA (SAP S/4HANA) and cc extreme scale. The single-node sim enterprises eliminate overhead asso environments, streamline high avai seamlessly as data volumes grow w performance.

Integrated with the recently annour SPS10, SGI UV 300H capitalizes on between SAP, Intel and SGI to optin workloads on multicore NUMA (non access) systems. This enables enter

### Gigabytes of



#### Powerful many-core coprocessors



Fast, large-capacity non-volatile Memory



 $\Rightarrow$  Enable tera-scale data processing on a single machine However: design and systems-level mechanisms required!

workloads on multicore NUMA (non access) systems. This enables enter

Gigabytes of



Intel is introducing a new family of enterprise PCIe SSDs with the aim of outperforming their existing DC P3600 series and even beating the DC P3700 series in many metrics. To do this, they've essentially put two P3600 SSDs on to one expansion card and widened the interface to 8 lanes of PCIe 3.0. While this does Powerful ma come across as a bit of a quick and dirty solution, it is a very straightforward way for Intel to deliver higher performance, albeit at the cost of sharply increased power consumption.

Fast, large-capacity non-volatile Memory

Abstractions for Big Data Processing

Large-scale big data analytics is possible on a single machine using systems and design-level abstractions for commodity and heterogeneous single machines. Large-scale big data analytics is possible on a single machine using systems and design-level abstractions for commodity and heterogeneous single machines.

Approaches:

- **Systems-level** analysis and mechanism for improving virtual memory performance with LATR
- **Design-level** abstractions for trillion-scale graph analytics with MOSAIC
- **Design-level** abstractions and data structures for billion-scale evolving graphs with CYTOM

## High-level structure



## High-level structure



 $\Rightarrow$  Focus of this talk: Processing evolving graphs with  $\rm Cytom$ 

## High-level structure



- $\bullet~\mathrm{Cytom}$  builds on lessons learned with  $\mathrm{MOSAIC}$
- However: Different challenges for setting of evolving graphs

### Thesis Statement and Approaches

### **2** CYTOM: Processing Billion-Scale Evolving Graphs

- Applications & Challenges
- CYTOM's Cell-based Graph Representation
- Handling Deletions
- Programming Interface

### 3 Conclusion

# New Tweets per second record, and how!

By ©raffi Friday, 16 August 2013 ₩ f in Ø

Recently, something remarkable happened on Twitter: On Saturday, August 3 in Japan, people watched an airing of Castle in the Sky, and at one moment they took to Twitter so much that we hit a one-second peak of 143,199 Tweets per second. (August 2 at 7:21:50 PUT), August 3 at 11:21:50 JST)

To give you some context of how that compares to typical numbers, we normally take in more than 500 million Tweets a day which means about 5,700 Tweets a second, on average. This particular spike was around 25 times greater than our steady state.



#### Social networks

# New Tweets per second record, and how!



Live traffic & route calculation

# New Tweets per second record, and how!



Analysis of next-gen Cellular Data, Connected Cars, ...

# New Tweets per second record, and how!



- $\Rightarrow$  Applications enable by processing evolving graphs
  - Online route calculation
  - Online Anomaly detection
  - Online Ranking
  - Online Community detection

Live traffic & route



Analysis of next-gen Cellular Data, Connected Cars, ...

- Constant stream of updates
- Optimizing graph storage format not possible on every update
  - Graph data structures supporting updates
- Co-running graph updates with algorithm
- Special challenges for single machines (compared to distributed systems):
  - Limited memory & compute
  - Fault tolerance

## Background: Graph data structures for evolving graphs

- Goal: Efficient update support and computational efficiency
- Four options:
  - Adjacency matrix
  - Adjacency lists
  - Compressed sparse rows (CSR)
  - Edge lists

## Background: Graph data structures for evolving graphs

- Goal: Efficient update support and computational efficiency
- Four options:
  - Adjacency matrix
  - Adjacency lists
  - Compressed sparse rows (CSR)
  - Edge lists

Sample graph:



## Background: Adjacency matrix

- Minimal storage per edge (1 bit)
- Great for dense graphs
- Limitation: Most real-world graphs are sparse

Example:



- For best case, storage cost per edge: 8 bytes
- Relatively simple construction
- Limitation: Overhead on traversal and allocation of new edges
- Used by Stinger



- For best case, storage cost per edge: 8 bytes
- Efficient access and lower storage overhead
- Limitation: Changes require rebuilding or complicated overflow areas
- Used by Llama



- For best case, storage cost per edge: 8 bytes
- Efficient access and lower storage overhead
- Limitation: Changes require rebuilding or complicated overflow areas
- Used by Llama

Example: Add edge 4



- For best case, storage cost per edge: 8 bytes
- Efficient access and lower storage overhead
- Limitation: Changes require rebuilding or complicated overflow areas
- Used by Llama

Problem: Array has to be shifted



- Simple construction
- Adding new edges possible
- Limitation: Overhead for storage cost, 16 bytes per edge
- Used by GraphIn, GraphOne



### • Comparison of presented data structures

Representation	Sparse graphs	Updates	Traversals	Storage
Adj. matrix	-	$\checkmark$	$\checkmark$	1 bit
Adj. lists	$\checkmark$	$\checkmark$	-	8 bytes
CSR	$\checkmark$	-	$\checkmark$	8 bytes
Edge list	$\checkmark$	$\checkmark$	$\checkmark$	16 bytes

### • Comparison of presented data structures

Representation	Sparse graphs	Updates	Traversals	Storage
Adj. matrix	-	$\checkmark$	$\checkmark$	1 bit
Adj. lists	$\checkmark$	$\checkmark$	-	8 bytes
CSR	$\checkmark$	-	$\checkmark$	8 bytes
Edge list	$\checkmark$	$\checkmark$	$\checkmark$	16 bytes
$Edge\ list_{\mathrm{CYTOM}}$	$\checkmark$	$\checkmark$	$\checkmark$	4 bytes

- Our choice for evolving graphs: *edge lists*
- However: Mitigation for storage overhead needed

From raw input graph to algorithmic output:



## High-level architecture

### Step 1 : Update to input graph



### Step **2** : Update CYTOM's internal structure



## High-level architecture

### Step 3 : Distribute update



### Step **4** : Algorithmic processing of update



## High-level architecture: Cells



 $\Rightarrow$  First focus: CYTOM's cells

• Lightweight partitioning

Step 1: Input graph:


- Lightweight partitioning
- Step 2: Adjacency Matrix



- Lightweight partitioning
- Step 3: Construction of cell  $C_{11}$



- Lightweight partitioning
- Step 4: Construction of all other cells



- Lightweight partitioning
- Step 5: Overlay current cell size



# Construction of $\operatorname{CYTOM}\nolimits$ 's cells: Adding new edges

- Adding new edges lightweight
- Division to obtain partition



# Construction of CYTOM's cells: Adding new edges

- Adding new edges lightweight
- Division to obtain partition



- Subgraph-centric
- Difference to *tiles* (MOSAIC) or *partitions* (GridGraph)?
  - Both designed for static execution
  - Immutable structures
  - Changes only possible coupled with preprocessing step
  - $\bullet~\mathrm{CYTOM}$  's cells can grow (and shrink) as required without preprocessing
- Enables short identifiers (16 bits)
  - $2 \times$  to  $4 \times$  less memory (32 or 64 bit identifiers)
  - $\bullet$  Compared to  $\operatorname{MOSAIC:}$  No metadata needed

# Benefit of Cytom's cells

- Enables short identifiers (16 bits)
  - 32 bits (4 bytes) per edge
  - $2 \times$  to  $4 \times$  less memory (32 or 64 bit identifiers)
  - No metadata needed
- Usage of edge lists
  - Insert: Append at end of vector

Base case with global edge lists



# Benefit of Cytom's cells

- Enables short identifiers (16 bits)
  - 32 bits (4 bytes) per edge
  - $2 \times$  to  $4 \times$  less memory (32 or 64 bit identifiers)
  - No metadata needed
- Usage of edge lists
  - Insert: Append at end of vector

Idea: Translate global IDs into local ones, store in CYTOM's edge list



# Benefit of Cytom's cells

- Enables short identifiers (16 bits)
  - 32 bits (4 bytes) per edge
  - $2 \times$  to  $4 \times$  less memory (32 or 64 bit identifiers)
  - No metadata needed
- Usage of edge lists
  - Insert: Append at end of vector

Idea: Translate global IDs into local ones, store in CYTOM's edge list



## Impact of short identifiers - lower memory footprint



 $\Rightarrow$  Up to 24% higher overall performance, 47% fewer cache references

## High-level architecture: Cell distribution

How to operate on active cells?



- Problem: Many more cells than threads
  - Millions of cells, dozens of threads
- Skewed cell size distribution

## How to load balance cells?



- Problem: Many more cells than threads
  - Millions of cells, dozens of threads
- Skewed cell size distribution
- Four choices implemented:
  - Static partitioning



static

- Problem: Many more cells than threads
  - Millions of cells, dozens of threads
- Skewed cell size distribution
- Four choices implemented:
  - Static partitioning
  - Work pool



- Problem: Many more cells than threads
  - Millions of cells, dozens of threads
- Skewed cell size distribution
- Four choices implemented:
  - Static partitioning
  - Work pool
  - Cell distributor



- Problem: Many more cells than threads
  - Millions of cells, dozens of threads
- Skewed cell size distribution
- Four choices implemented:
  - Static partitioning
  - Work pool
  - Cell distributor
  - Hierarchical cell distributor



- Problem: Many more cells than threads
  - Millions of cells. dozens of threads
- Skewed cell size distribution
- Eour choices implemented:
  - Static partitioning
  - Work pool
  - Cell distributor
  - Hierarchical coll
  - $\Rightarrow$  Hierarchical cell distributor improves throughput by 2.1 $\times$ Further optimizations enabled by 2D structure?



Abstractions for Big Data Processing

# 2D-partitioning: Opportunity for different traversals

 $\bullet\,$  Multiple traversal strategies are possible with  $\rm CYTOM\,$ 

- Evaluated options:
  - Column first
    - Read optimized



# 2D-partitioning: Opportunity for different traversals

 $\bullet\,$  Multiple traversal strategies are possible with  $\rm CYTOM\,$ 

- Evaluated options:
  - Column first
    - Read optimized
  - Row first
    - Write optimized



# 2D-partitioning: Opportunity for different traversals

 $\bullet\,$  Multiple traversal strategies are possible with  $\rm CYTOM\,$ 

- Evaluated options:
  - Column first
    - Read optimized
  - Row first
    - Write optimized
  - Hilbert order
    - Mix between read and write optimization



- With MOSAIC: Hilbert order best for locality and performance
  - Optimize for reads and writes
- MOSAIC's scenario:
  - Batch-processing of large amounts of edges
  - Entire graph being processed at a time
- However: CYTOM's scenario is different:
  - Small sets of changes, e.g., one million edges (maximum: 12 MB of vertex and edge data)
  - Only parts of the graph being processed
- For CYTOM: Optimizing writes beneficial!

#### Traversal strategies with $\operatorname{Cytom}$



 $\Rightarrow$  Row-first, write-optimized strategy performs up to 80% better

# Optimization: Selective scheduling

- Important in evolving graphs:
  - Not all vertices active all the time
- Effect: Many inactive edges
- Common in many algorithms
- Idea: Only operate on cells with active vertices

# Optimization: Selective scheduling

- Important in evolving graphs:
  - Not all vertices active all the time
- Effect: Many inactive edges
- Common in many algorithms
- Idea: Only operate on cells with active vertices

Target vertex



# Optimization: Selective scheduling

- Important in evolving graphs:
  - Not all vertices active all the time
- Effect: Many inactive edges
- Common in many algorithms
- Idea: Only operate on cells with active vertices

Target vertex



## Impact of selective scheduling



 $\Rightarrow$  Up to 5.6× improvement due to reduced number cells

## High-level architecture: Algorithmic interface

#### What is CYTOM's algorithmic interface?



### Example connected components

- Iterative algorithm
- Minimize connected component over incoming edges

Step 1: Initialize with local vertex ID



- Iterative algorithm
- Minimize connected component over incoming edges

Step 2: Minimize over all incoming edges with Pull



- Iterative algorithm
- Minimize connected component over incoming edges

Step 3: Update local connected component value



- Iterative algorithm
- Minimize connected component over incoming edges

Step 4: Run apply and mark active vertices



- Iterative algorithm
- Minimize connected component over incoming edges

Step 5: New iteration, minimize with Pull



- Iterative algorithm
- Minimize connected component over incoming edges

Step 6: Run Apply on all active vertices



- Iterative algorithm
- Minimize connected component over incoming edges

Step 7: Convergence reached after one more iteration


- Add new vertex and edge
- Algorithm has to update connected component

- Add new vertex and edge
- Algorithm has to update connected component

Step 1: Add edge ④ and vertex ④, init connected component



- Add new vertex and edge
- Algorithm has to update connected component

Step 2: Pull incoming edge for ④



- Add new vertex and edge
- Algorithm has to update connected component

Step 3: Convergence after next iteration



- Potentially devastating impact on algorithmic result
- Example with connected components:

- Potentially devastating impact on algorithmic result
- Example with connected components:



- Potentially devastating impact on algorithmic result
- Example with connected components:



- Potentially devastating impact on algorithmic result
- Example with connected components:



- Potentially devastating impact on algorithmic result
- Example with connected components:



- "Safe" solution: Re-execute algorithm
- CYTOM's approach: Mark critical edges carrying results
- If critical edge deleted or updated: Re-execute



Delete non-critical edge **2** :



### Delete *critical* edge ① :



Re-execute algorithm for correct result:



- Scatter-Gather-Apply (GAS) variant
- Algorithm developer not exposed to subgraph-centric aspects
- Special addition for evolving graphs: edgeChanged callback
  - React to edge changes
  - Trigger algorithm re-execution if necessary (e.g., deletions)

## Example: Connected components

#### Edge processing callbacks:

#### **Edge-centric** operation

IK	1	// Inside a Cell (local)	
Local graph processing on Ce	2	// Edge e = (Vertex src, Vertex tgt)	
	3	def Pull(Vertex src, Vertex tgt):	
	4	if src.value < tgt.value:	
	5	markCritial(src, tgt)	
	6	return src.value	
	7	// Collecting Cell Results (local & global)	
	8	def Reduce(Vertex $v_1$ , Vertex $v_2$ ):	aph ing
	9	<b>return</b> $min(v_1.value, v_2.value)$	
	10	// On Global Vertices	l gi essi
	11	def Apply(Vertex v):	oba
	12	<pre>if current(v.value) != next(v.value):</pre>	<sup>D</sup> G
	13	activate(v) <i>Vertex-centric op</i>	eration

### Example: Connected components

Special for evolving graphs: edgeChanged callback

- Filter updates without implied changes
- Handle deletions



#### **Edge-centric operation**

## Impact of edgeChanged callback for insertions

• Benefit: Filter out uninteresting updates



 $\Rightarrow$  Up to 2.8× improvement in overall throughput due to reduced work for algorithm

Steffen Maass

## Impact of edgeChanged callback for deletions



Re-execution vs. critical edges approach with CC:

- How does CYTOM perform compared to other systems?
- Metrics:
  - Throughput of graph updates
  - Overall throughput with updates and algorithm processing

#### • Features offered by each system

Characteristics	LLAMA	GraphIn	GraphOne	Stinger	Cytom
Incremental proc.	-	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$
Synchronous proc.	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$
Asynchronous proc.	-	-	-	-	$\checkmark$
Approximation	-	-	-	-	$\checkmark$
Disk persistence	$\checkmark$	-	$\checkmark$	-	$\checkmark$
Deletions	$\checkmark$	$\checkmark$	$\checkmark$	-	$\checkmark$
Snapshot view	$\checkmark$	-	-	-	$\checkmark$

- $\bullet$  Compared to GraphOne,  $\mathrm{Cytom}$  offers 2.8× the throughput
  - $\bullet\,$  Uses CSR for static snapshots, edge lists for new data
  - Log-structured approach
- $\bullet$  Compared to Differential Dataflow,  $\mathrm{CYTOM}$  offers 4.8 $\times$  the throughput
  - Generalized data flow system
  - Graph support with edge lists

## Comparison to GraphIn

- GraphIn uses edge lists
- Only scales to pprox 100M edges



 $\Rightarrow$  Up to 5.5× improvement in throughput

## Comparison to STINGER

- STINGER's focus: Use atomics to insert into adjacency lists
- Earliest systems of all evaluated ones



 $\Rightarrow$  At least 60× speedup due to more efficient graph updates and APIs

Steffen Maass

 $\bullet~\mathrm{Cytom},$  a billion-scale engine for evolving graphs

- Cell-based design enables high update rate
- Selective scheduling enabled by cell-based design
- edgeChanged callback to quickly react to graph updates

 $\bullet~\mathrm{CYTOM},$  a billion-scale engine for evolving graphs

- Cell-based design enables high update rate
- Selective scheduling enabled by cell-based design
- edgeChanged callback to quickly react to graph updates
- Additional components of dissertation:
  - $\bullet~\mathrm{MOSAIC},$  a trillion-scale engine for static graphs
  - $\bullet~{\rm LATR},$  an OS-level approach to reduce overheads of synchronous TLB shootdowns

• Outcome: We can do large-scale data analytics on a single machine

 $\bullet~\mathrm{CYTOM},$  a billion-scale engine for evolving graphs

- Cell-based design enables high update rate
- Selective scheduling enabled by cell-based design
- edgeChanged callback to quickly react to graph updates
- Additional components of dissertation:
  - $\bullet~\mathrm{MOSAIC},$  a trillion-scale engine for static graphs
  - $\bullet~{\rm LATR},$  an OS-level approach to reduce overheads of synchronous TLB shootdowns

• Outcome: We can do large-scale data analytics on a single machine

### Thanks!