# $1000^5$ Petascale graphs: The end-to-end challenge



**Full Internet Map
[Lumeta]**



**Social Graph
[Facebook]**

- GraphLab is indeed promising
- But we struggled with feeding it and other practicalities
- Set out to study potential approaches...

Cluster Computing Architecture

# But first...
# a bit on how we got here.



Haijie (Jay) Gu
PhD Candidate
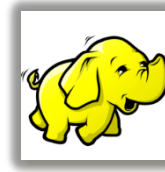CMU
and Summer 2012
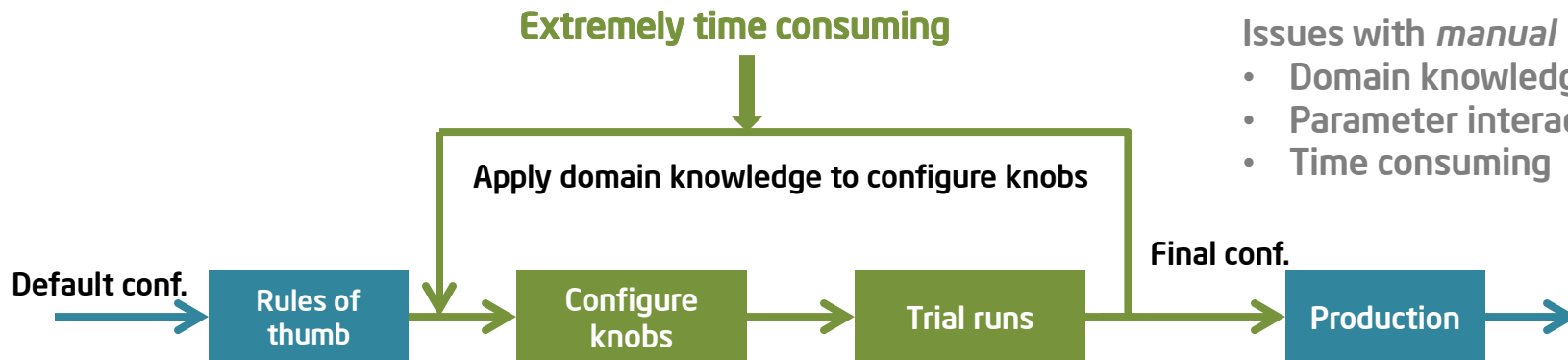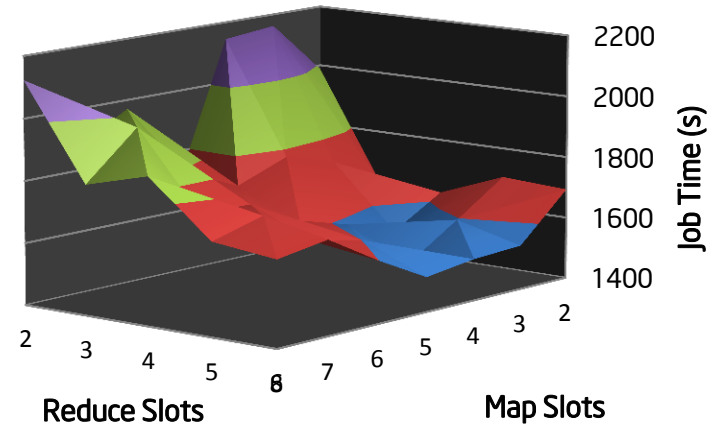Intern

Danny Bickson
Postdoc CMU

Joseph Gonzalez
Postdoc AMPLab

Yucheng Low
PhD Candidate CMU

# Hadoop Research

- Evaluating new cluster technologies requires solid baselines
- But an exhaustive search for the best Hadoop configuration would take **7,257,600,000 lengthy trials!**
- Solve the challenge and accelerate our other work at the same time?

Job Time (s): 2200, 2000, 1800, 1600, 1400

Reduce Slots: 2, 3, 4, 5, 6, 7, 8

Map Slots: 2, 3, 4, 5, 6, 7

**Extremely time consuming**

**Issues with _manual tuning_:**
- Domain knowledge
- Parameter interaction
- Time consuming

Apply domain knowledge to configure knobs

Default conf. → **Rules of thumb** → **Configure knobs** → **Trial runs** → Final conf. → **Production**

**We spend weeks tuning clusters for a few days of experiments.**

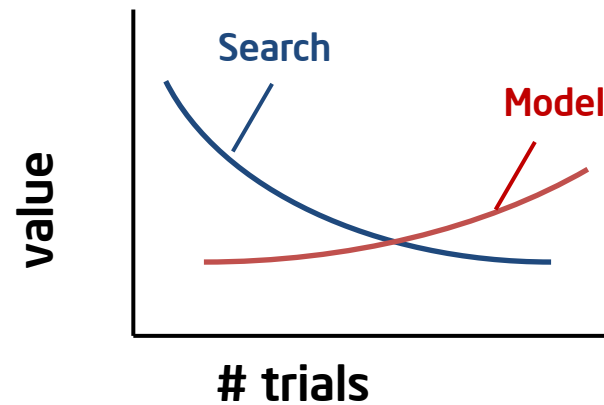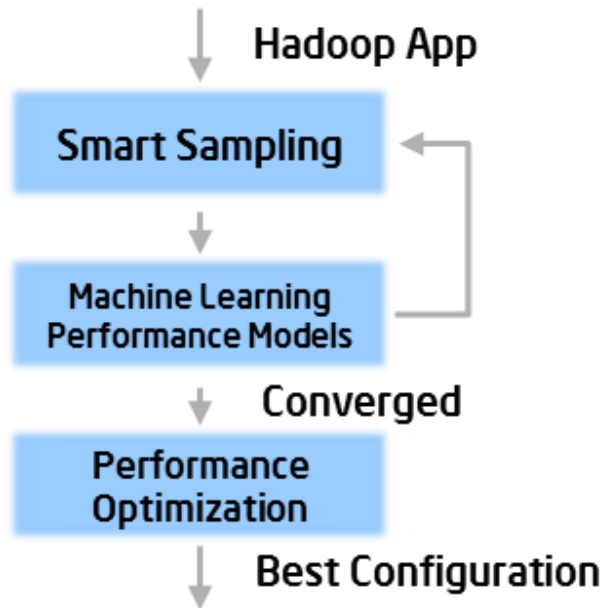Cluster Computing Architecture

4

# Our Approach

1. Focus on the most important parameters for each circumstance*

2. Apply generalized search algorithms to efficiently explore the parameter space

3. Model the system to reason about unexplored space

* Future work

Cluster Computing Architecture

# Gunther: The Elephant Trainer

## An Auto-tuner for Hadoop MapReduce

Hadoop App

Smart Sampling

Machine Learning
Performance Models

Converged

Performance
Optimization
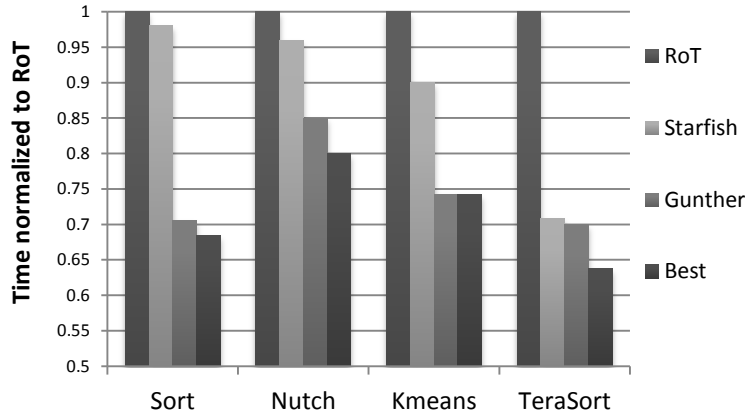
Best Configuration

Search
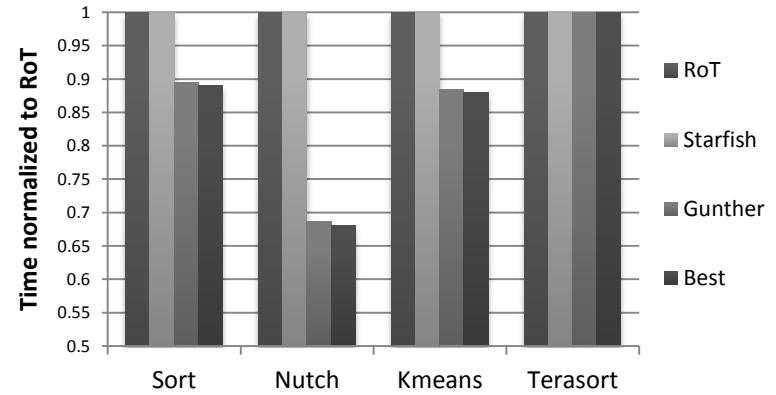
Model

value

# trials

*Key benefits:*
- Little domain knowledge required
- Easily adapts to new datasets, workloads, frameworks, & clusters
- More effective and faster than manual approaches

Cluster Computing Architecture

(intel)

# Search + Model

## Genetic search algorithm is ~95% *effective* in <30 trials



**Storage bottlenecked cluster**



**Network bottlenecked cluster**

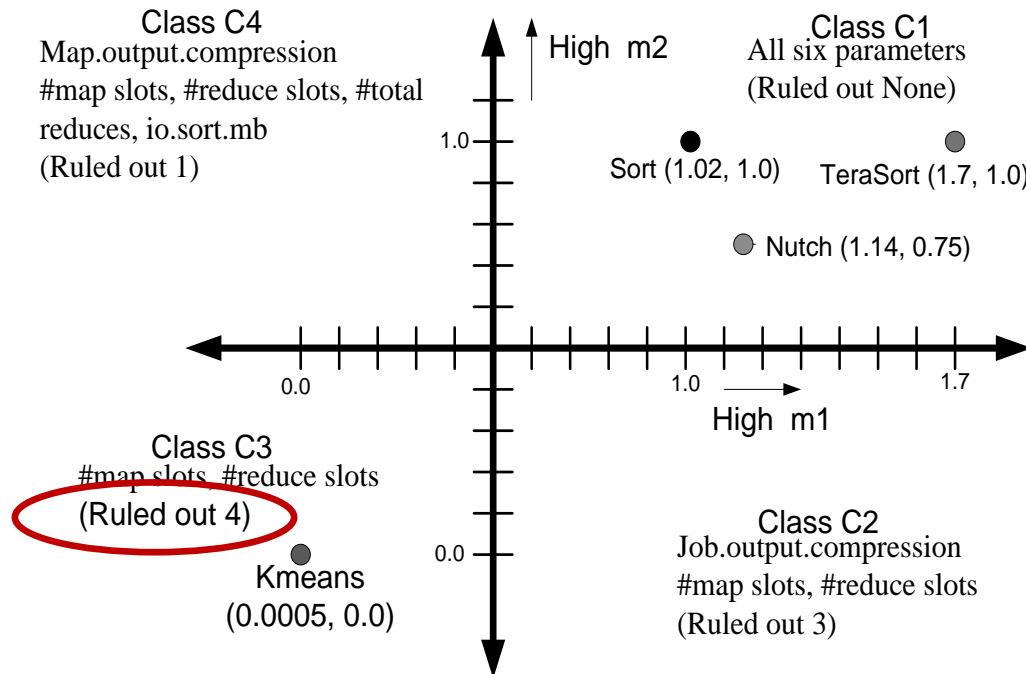## SVR model is very *accurate* but requires hundreds of trials (320 in this case)

| Modeling Approach | SNB Cluster | | | | | | | ZT Cluster | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | min | Q1 | median | mean | Q3 | max | IQR | min | Q1 | median | mean | Q3 | max | IQR |
| MLR | 0 | 7 | 15 | 17 | 23 | 68 | 16 | 0 | 8 | 18 | 23 | 33 | 117 | 25 |
| MLR-I | 0 | 6 | 14 | 16 | 22 | 69 | 16 | 0 | 8 | 17 | 22 | 31 | 129 | 23 |
| MLR-Q | 0 | 6 | 11 | 15 | 22 | 66 | 16 | 0 | 7 | 16 | 19 | 27 | 92 | 20 |
| MLR-IQ | 0 | 6 | 11 | 14 | 20 | 67 | 14 | 0 | 7 | 15 | 18 | 25 | 87 | 18 |
| ANN | 0 | 4 | 10 | 12 | 17 | 61 | 13 | 0 | 4 | 10 | 12 | 17 | 61 | 13 |
| M5Tree | 0 | 5 | 10 | 12 | 17 | 65 | 12 | 0 | 5 | 10 | 14 | 19 | 71 | 14 |
| SVR | 0 | 2 | 4 | 8 | 10 | 73 | 8 | 0 | 3 | 6 | 10 | 13 | 64 | 10 |

## Apply model to predict perf and inform future searches.

Cluster Computing Architecture

# Dimensionality Reduction

Rule out parameters *up front* that primarily affect resources that aren't likely to bottleneck



Class C4
Map.output.compression
#map slots, #reduce slots, #total
reduces, io.sort.mb
(Ruled out 1)

High m2

Class C1
All six parameters
(Ruled out None)

Sort (1.02, 1.0)    TeraSort (1.7, 1.0)

Nutch (1.14, 0.75)

1.0

0.0    1.0    1.7
High m1

Class C3
#map slots, #reduce slots
(Ruled out 4)

Kmeans
(0.0005, 0.0)

0.0

Class C2
Job.output.compression
#map slots, #reduce slots
(Ruled out 3)

$$m1= \frac{spilled}{map+reduce\ inputs}$$

$$m2= \frac{HDFS\ bytes\ written}{HDFS\ bytes\ read}$$

**Direction:**
1. Incorporate node- & cluster-level utilization observations (*m*) into model
2. Apply EV-based MV analysis offline to determine *what* params matter *when*
3. Use 1st run to collect *m* and apply to search

# But tuned MapReduce is still MapReduce.

# "Chance favors the *connected* mind."
## --Steven Johnson



**A new country**

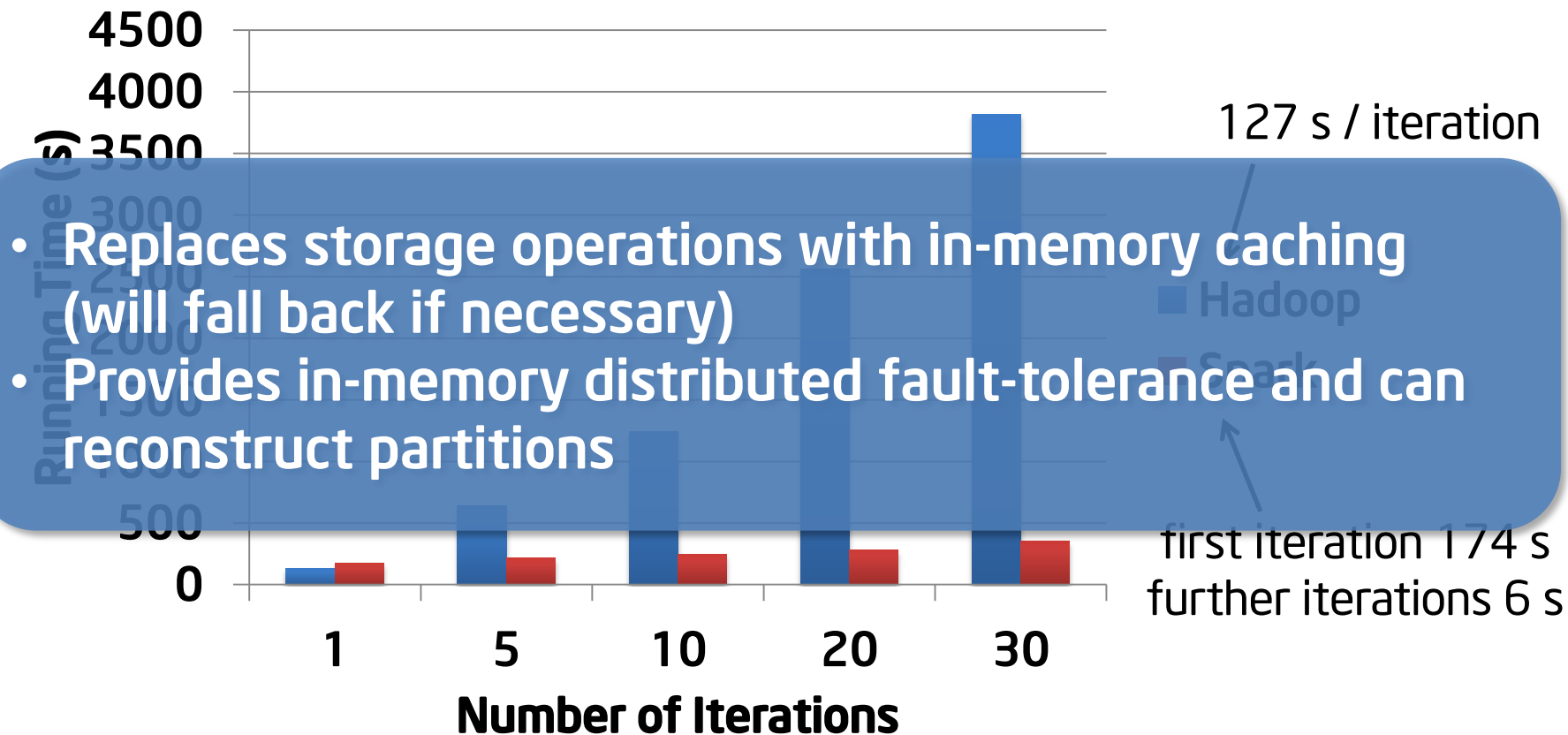**A new planet**

# Garth made the connection at the December 2011 ISTC retreat!

# Spark*

## Fast, Interactive, Language-Integrated Cluster Computing



- **Replaces storage operations with in-memory caching (will fall back if necessary)**
- **Provides in-memory distributed fault-tolerance and can reconstruct partitions**

127 s / iteration

Hadoop

Spark

first iteration 174 s
further iterations 6 s

**Running Time (s)**

4500
4000
3500
3000
2500
2000
1500
1000
500
0

1    5    10    20    30

**Number of Iterations**

*Zaharia et al.  UC Berkeley.  Retrieved from www.spark-project.org.*

Cluster Computing Architecture

(intel)

11

# PageRank on 64 DP Xeons

**40M Webpages,  1.4 Billion Links**



Hadoop — 5.5 hrs

Twister — 1 hr

GraphLab — **8 min**

**GraphLab is 41X faster. What's going on here?**

Hadoop results from [Kang et al. '11]
Twister (in-memory MapReduce) [Ekanayake et al. '10]

# Big graphs are a big deal!

| | | |
|---|---|---|
| 100B Neuron 100T Relationships | 1B Users 140B Friendships | 1 Trillion Pages 100s T Links |
| Human Brain | Social Network | Internet |
| Millions of Products & Users | 27M Users 70K Movies | Large Biological Cell Networks |
| e-commerce | Online Services | Science |

**Many problems involve irregular data structures most naturally expressed as graphs, trees, and arbitrary sets**

**(paraphrasing Keshav Pingali)**

(intel)

# Collaborative Filtering: Mining Relationships

## Customers Who Bought This Item Also Bought What?

Dog Food

**|P|**

**|C|**

Dog Food
#1

Milo's Meatball Treats
#2

Greenies for Teeth
#3

# Time to find product similarities is

| Days | worst case |
| Minutes | if algorithm **exploits data dependency structure** |
| Seconds | if **ideally partitioned** across M machines |

intel

# Graph processing:
# An extremely short history

**Data-parallel**
▼
Data-parallel + Iterative
▼
Graph-parallel + Iterative
▼
Asynchronous graph-parallel

Ship the **entire** graph structure...
... over and over ...
or, better yet, pass the results ...
... whenever you want.

(intel)

# So we're done, right?

intel

"I spend more than half of my time integrating, cleansing and transforming data without doing any actual analysis. Most of the time I'm lucky if I get to do any *analysis* at all."

Anonymous Data Scientist
from Jeff Heer's (Stanford) interview study, 2012

(intel)

Cluster Computing Architecture

# Taking a Broader Perspective



* Adaptation of GraphLab team material

How do we construct the graph?
How do we store it? Query it?
Analyze it? _____ it?

Many of these challenges are solved for small problems... but what about Internet scale?

# Challenges for Emerging Area

1. Few people skilled in the apps and algos
2. App frameworks emerging and evolving rapidly
3. Lack of tools to deploy systems and analyze system behavior

**Lightly charted territory offers big opportunities for Intel and other companies.**

Cluster Computing Architecture

(intel)

# Parallel Machine Learning (ML):
## Joint work with the ISTC for CC (UW/CMU)

**intel**

# Machine Learning Pipeline



**Data** → **Extract Features** → **Graph Formation** → **Structured Machine Learning Algorithm** → **Value from Data**

**Graph Ingress**
mostly data-parallel

**Graph-Structured Computation**
graph-parallel

# Hadoop for Graph Construction

- **Intuitive Map and Reduce programming model (in Java)**
- **Framework takes care of resource provisioning**
- **Provides redundant storage and fault recovery**

**Map**

SHUFFLE

**Reduce**

Cluster Computing Architecture

23

**People**     **Interests**

Kushal

Diana

Nilesh

Danny

Ted

Frank

Ivy

Jay

Cluster Computing Architecture

24

# Building Graphs for Practical Apps

| | Raw Data | Pre-processing | Graph Formation | Add Network Information |
|---|---|---|---|---|
| What **words** are most associated with what **(hidden) topics**? | XML Docs | Extract Doc Names and Words | Bipartite (Docs, Words) | Count Word Frequency |
| What does context tell me about the **type** (person, place, thing) of this noun? | News Feeds | Extract Noun Phrases and Contexts | Bipartite (NP, Context) | Count NP Frequency & Initialize **type** Distribution |
| What are the highest **ranked pages**? | Web Pages | Extract Page URLs and Links | Directed Graph | N/A |

(intel)

# And, in practice and at scale we must:

| Raw Data | Pre-processing | Graph Formation | Add Network Information | Finalize for Parallel Computation |

- Minimize the use of system resources, like memory, storage, etc.

- Ensure GL's computational effort is load balanced for *power-law graphs*

- Do our best to ensure the graph we generated is the one we intended to

## ... but the application programmer shouldn't be responsible for this domain expertise!

(intel)

Large-Scale Graph Construction using Apache™ Hadoop™

## GraphBuilder makes it easy.

- Fills a hole in the ecosystem
- Written in Java for easy use in Hadoop MR and apps
- Offloads domain expertise

Cluster Computing Architecture

# *GraphBuilder* Data flow

## **E**xtract
Graph formation from data source(s)

HDFS

DB

XML Docs

Feature Extraction and Tabulation

## **T**ransform
Apply cleaning and transformation

Graph Checks and Transformation

## **L**oad
Prepare for graph analytics

Graph Compression, Partitioning, and Serialization

**App-Specific Code**

***GraphBuilder* Library**

Cluster Computing Architecture

# Extract - Graph Formation

## Extract features from data to construct relationships

**Read** → **Tokenize** → *Optional* Reduce (●—●) → f(x)

```
conf.set(XMLInputFormat.START_TAG_KEY, START_TAG);
conf.set(XMLInputFormat.END_TAG_KEY, END_TAG);
new XMLRecordReader((FileSplit) split, conf);
```
**Read Records**

```
Document doc = builder.parse(new
InputSource(new StringReader(s)));
title = xpath.evaluate("//page/title/text()",
doc);
```
**Extract Features**

```
title = title.replaceAll("\\s", "_");
id = xpath.evaluate("//page/id/text()", doc);
String text =
xpath.evaluate("//page/revision/text/text()",
doc);
parseLinks(text);
```
**Parse Elements**

- Write simple data-specific functions.
- Program sequential, not parallel!

# Extract - Tabulation

Built-in tabulation functions for TF, TFIDF, WC, ADD, MUL, DIV. Interface for custom tabulation on source and/or target vertex

Example: Term Frequency



*User Defined:*

▶ Reduce ( ⬤—⬤ ) → f(x)

▶ Apply (f(x)) → ⬤—⬤

$$tf(t,d) = \frac{f(t,d)}{\max\{f(w,d): w \in d\}}$$

# Transform – Graph Transforms & Checks

- Would like the ability to:
  - Optionally filter duplicate, dangling and/or self edges
  - Transform a directed graph into an undirected graph
  - Calculate graph statistics, compute sub-graphs, etc.
- The library provides:
  - Functions to perform self-, dangling- and duplicate-edge removal
  - Directionality transformation
- Solutions are based on a distributed hashing algorithm



Steering function

# Load - Graph Compression

- We can save memory if we compress/normalize graph
- But, seems to call for global lookups in a framework that prefers *independent subproblems*
- A simple, scalable solution is to "shard" ordered lists:

Dictionary

(Aaron,0)
(AMD,4)
(Brad,1)
(CMU,2)
(Dan,5)
(Dave,3)
(IBM,6)
(Intel,7)

**Machine 1**

Dictionary Shard 1

(Aaron,0)
(AMD,4)
(Brad,1)
(CMU,2)

Unconverted Edge List

(Aaron,IBM)
(Brad,Intel)

(AMD,5)
(CMU,3)

Converted Edge List

(5,4)
(3,2)
(0,6)
(1,7)

Dictionary Shard 2

(Dan,5)
(Dave,3)
(IBM,6)
(Intel,7)

(Dan,AMD)
(Dave,CMU)

(IBM,0)
(Intel,1)

(Source Sorted)

(Dest Sorted)

**Machine 2**

# Load - Graph Partitioning

"Cut quality varies inversely with cut balance." [Kevin Lang, '04]

- Minimize communications by minimizing the number of machines *v* spans

- Place about the same number of *edges* on each machine

# Load - Graph Partitioning

"Cut quality varies inversely with cut balance." [Kevin Lang, '04]

- Minimize communications by minimizing the number of machines *v* spans

- Place about the same number of *edges* on each machine

# Heuristic-Based Partitioning Strategies

- **Random edge placement:** Edges are placed randomly by each system

- **Greedy edge placement:** Global coordination for edge placement to minimizes the vertex spanned

- **Oblivious greedy placement:** Implements a local version of the Greedy without global coordination

# Oblivious Algorithm

# Machine 1's Shard



## Partition 1

CASE 1:
Both end points have never been seen before

→ Randomly assign

## Partition 2

# Machine 1's Shard



**Partition 1**

**Partition 2**

**CASE 2:**
Both end points have been seen before on the same partition

→ Assign to a partition which contains both endpoints

# Machine 1's Shard



## Partition 1



**CASE 3:**
Both end points have been seen before but on different partitions

→ Assign to any partition that contains an endpoint

## Partition 2

# Machine 1's Shard



**Partition 1**



## CASE 3:
Both end points have been seen before but on different partitions

→ Assign to any partition that contains an endpoint

**Partition 2**

(intel)

# Machine 1's Shard



## Partition1

CASE 4:
Only one end point
has been seen
before

→ Assign to a
partition that
contains the
endpoint

## Partition 2

Cluster Computing Architecture

# Partitioning Quality

## Twitter Graph: 41M vertices, 1.4B edges



Greedy yields a quality cut and the best performance....

# Performance for Partitioning



Performance is inversely proportional to replication.

*Gonzalez et al., "PowerGraph: Distributed Graph-Parallel Computation on Natural Graphs," [OSDI'12]

# Load - Graph Serialization

**Partitioning**

↓

**JSON Encoding**

↓                    ↓

**Edge Lists**      **Vertex Lists**

```
{
   "src_id": 34,
   "dest_id": 45
   "e-data": 30
}
```

```
{
   "ver_id": 34,
   "v-data": 56,
   "mirror": [1,2,3],
   "owner": 1
}
```

- Self-describing data format
  - JSON +/- compression
- Extensible
  - Easy to connect with Graph Databases
  - Plug-in Graph Visualizers

Cluster Computing Architecture

# GraphBuilder Stack



**Built-in Parser/Tabulator**

**Custom Parser/Tabulator**

Extract

Transform

Load

Hadoop MapReduce

Distributed Graph

Hadoop/HDFS

**GraphBuilder app**

GraphLab app

**GraphBuilder**

MapReduce

Distributed Graph Computation (GLv2)

HDFS

(intel)

# GraphBuilder Demo

**WIKIPEDIA**

Partitioned
Bipartite graph

Latent Dirichlet Allocation
(LDA) Algorithm

*GraphBuilder*

GraphLab

*Knowledge Extraction*

## WordCloud Visualizer

| NTopics | NWords | NDocs | NTokens | Alpha | Beta |
|---|---|---|---|---|---|
| 50 | 296248 | 4003417 | 824408165 | 0.5 | 0.1 |

van africa
dutch african
united south swedish
germany references
european sweden
netherlands list
air aircraft airport
force flight squadron wing base
flying international united space fighter
aviation group pilot training mission
operations airlines states raf service world
control missile ii military unit ground

norwegian norway danish
europe republic france kingdom
international denmark states
external usa belgium finland see finnish
countries

war army
military general battle
forces division regiment
commander corps command
infantry force st officer service troops
battalion british cross men world nd major
soldiers brigade attack german chief
lieutenant

## Topic Modeling

age
population
people living average
income years census
median city family town
households families total size
square area white county township
household united density females mile
references made present american

india indian
pakistan sri temple khan ali tamil
singh islamic muslim iran
references state known muhammad
delhi hindu ibn arabic shah lanka name islam
see iranian people arab persian malaysia

music musical theatre
opera piano orchestra jazz
dance performed composer
concert symphony performance major songs
folk works classical violin composed played
string play solo performances musicians
sound york recorded song

(intel)

# Our Wikipedia Graphs

**LDA**

Top 1% of vertices are adjacent to 49% of the edges!

*(plot: Number of Vertices vs Degree, LDA)*

**PageRank**

Power Law

*(plot: Number of Vertices vs Degree, PageRank)*

| Graph | |V| | |E| | α |
|---|---|---|---|
| LDA | 4.9M | 478M | 2.23 |
| PageRank | 9.7M | 107M | 2.41 |

Cluster Computing Architecture

(intel)

# Prototype Overview

- **Hardware: 8 node cluster**
  - 1U Dual CPU (Intel SNB) Amazon build ZT systems
  - 64 GB Memory, Four SATA Hard Drives
  - Intel 10G Adapter and Switch
- **Software:**
  - Apache Hadoop 1.0.1
  - GraphLab v2.1
  - *GraphBuilder* beta

# Preliminary Results

| Graph | Custom plug-in code | Graph Compression | Partitioning Improvement (vs. Random) |
|---|---|---|---|
| PageRank | 100 lines | 60% | 17% |
| LDA | 130 lines | 5% | 32% |

PageRank Graph

**45 min**
$|V| = 54M, |E| = 1.4B$

LDA Graph

**13 min**
$|V| = 20M, |E| = 128M$

■ **Extract**  ■ **Transform**  ■ **Load**

Cluster Computing Architecture

# Scaling Experiment



Stacked bar chart. X-axis: Dataset Size (4x, 2x, 1x). Y-axis: Overall Time (Seconds), 0 to 1000. Legend: Finalization (green), Partitioning (red), Pre-processing/Graph Formation (blue).

# Collaboration ahead!



Sam Madden    Carlos Guestrin    Ted Willke

# All Together Now

**Parallel Machine Learning Pipeline**

| ML and Analytics Toolkits |
| --- |

| Parallel ML Cluster API |
| --- |

GraphBuilder

Hadoop MapReduce

Distributed GraphLab

Data Parallel

Graph Ingress

Graph Parallel

| Distributed FS and/or Graph DB | Local Store |
| --- | --- |

## Future areas for ISTC collaboration:

1. Improve usability and data wrangling
2. Research GL fault tolerance and local storage support
3. Advance GB + GL for streaming and time-evolving apps

Cluster Computing Architecture

(intel)

# Launches today!



Intel open source portal at
http://www.01.org

GraphLab2 at
http://graphlab.org

Both under Apache 2.0 licensing.

*Looking for research partners and committers.*

Contacts:
nilesh.jain@intel.com
theodore.l.willke@intel.com

# How many people are pointing to you and what's their relative importance?



Loops in graph - Must iterate!

# Properties of Graph-Structured Computation

Dependency
Graph

Local
Update

Iterative
Computations



## Similar properties for many other problems!

# Data Dependencies

- **MapReduce does not efficiently express data dependencies**
  - **User must code substantial data transformations**
  - **Costly data replication**



- **MapReduce does not efficiently express iterative algorithms**

Cluster Computing Architecture

(intel)

# Approaches to Graph-Structured Computation

- **Bulk Synchronous Processing (BSP)**
  - Giraph on Hadoop (Inspired by Google Pregel)
  - Dryad (Microsoft Research)
  - Apache Hama on Hadoop (Twitter)

- **Asynchronous Graph-Parallel**
  - Galois (UT Austin) → Edge partitioning
  - GraphLab (CMU) → Vertex partitioning

## GraphLab has an edge.



Runtime in Seconds vs Number of CPUs — BSP, Async



Program Like This → Run Like This
Machine 1  Machine 2
Master  Slave
Split **High-Degree** Vertices

(intel)

# GraphLab Goals

- ## Designed specifically for ML
  - Graph dependencies
  - Iterative
  - Asynchronous
  - Dynamic

- ## Simplifies design of parallel programs
  - Abstracts away hardware issues
  - Automatic data synchronization
  - Addresses multiple hardware architectures

# The GraphLab Framework

Graph Based
*Data Representation*

Update Functions
*User Computation*



Scheduler

Consistency Model



(intel)

# GraphBuilder makes it easy.

- Fills a hole in the ecosystem
- Written in Java for easy use in Hadoop MapReduce and applications
- Offloads domain expertise

Raw Data

Parsing | Tokenization or Feature Extraction | Edge List Generation | E, V Data Tabulation | Graph Normalization | Graph Checks & Transforms | Graph Partitioning & Serialization

To GraphLab

**App-Specific Code**

*GraphBuilder* **Library**

Cluster Computing Architecture

(intel)

# Prototype Overview

- **Hardware: 8 node cluster**
  - 1U Dual CPU (Intel SNB) Amazon build ZT systems
  - 64 GB  Memory, Four SATA Hard Drives
  - Intel 10G Adapter and Switch
- **Software:**
  - Apache Hadoop 1.0.1
  - GraphLab v2.1
  - *GraphBuilder* beta



**Cluster Computing Architecture**

# Example Topics Discovered from Wikipedia