Declarative Abstractions and Scalable Platforms for Big Data Analytics

Amol Deshpande

Department of Computer Science and UMIACS
University of Maryland at College Park

Big Data

- Explosion of data, in pretty much every domain
 - Sensing devices and sensor networks that can monitor everything from temperature to pollution to vital signs 24/7
 - Increasingly sophisticated smart phones
 - Internet, social networks making it very easy to publish data
 - Scientific experiments and simulations
 - Internet of Things
 - Many aspects of life being turned into data ("dataification")
- "Big Data" (= extracting knowledge and insights from data) becoming fundamental
 - Science, business, politics -- largely driven by data and analytics
 - Many others (Education, Social Good) are slowly being

Four V's of Big Data

- Big data not just about "Volume"
 - Large scale of data certainly poses many problems
 - But most datasets are pretty small...
- Variety and heterogeneity in both data and applications
 - Text, networks, time series, nested/hierarchical, multimedia, ...
 - Increasingly complex and specialized analysis tasks

Velocity

 Data generated at very high rates and often needs to be processed in real time

Veracity

- What/who to trust? How to reason about data quality issues?
- Easy to draw wrong statistical conclusions from large datasets
- Issues becoming more important with increasing automation...

My Research

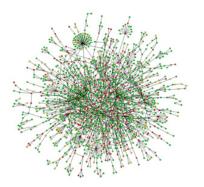
- Building data management systems to address challenges in managing and analyzing big data by..
 - Designing intuitive, formal, and <u>declarative abstractions</u> to empower users, and
 - Developing <u>scalable platforms</u> and algorithms to support those abstractions over large volumes of data
- Major research thrusts over the last 10 years
 - Uncertain and probabilistic data management
 - Graph data management
 - Data management in the cloud
 - Collaborative data analytics
 - Query processing and optimization

Outline

- Graph Data Management
 - Declarative Graph Cleaning
 - A Framework for Distributed Graph Analytics
- DataHub: A platform for collaborative data science
 - Recreation/Storage Tradeoff in Version Management (or, why git/svn are not good at managing datasets)

Graph Data

 Increasing interest in querying and reasoning about the underlying graph (network) structure in a variety of disciplines



A protein-protein interaction network

Citation networks

Communication networks

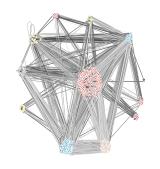
Disease transmission networks



Social networks



Knowledge Graph



Financial transaction networks

World Wide Web



Stock Trading Networks

Different types of "queries"

Subgraph pattern matching: Given a "query" graph, find where it occurs in a given "data" graph

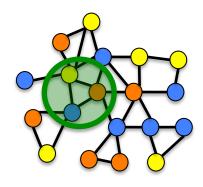
Reachability; Shortest path; Keyword search; ...

networks

Historical or Temporal queries: "Find most important nodes in a communication network in 2002?"



Query Graph



Data Grapi

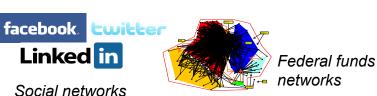


A protein-protein interaction network

Citation networks

Communication networks

Disease transmission networks



Knowledge Graph World Wide Web



Stock Trading Networks

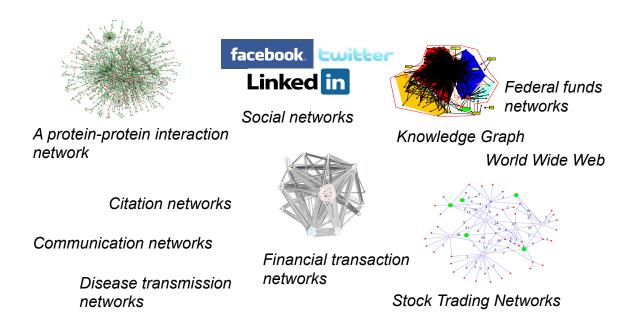
Different types of "queries"

Subgraph pattern matching; Reachability; Shortest path; Keyword search; Historical or Temporal queries...

Continuous "queries" and Real-time analytics

Online prediction in response to new data Monitoring: "Tell me when a topic is suddenly trending in my friend circle"

Anomaly/Event detection: "Alert me if the communication activity around a node changes drastically"



Different types of "queries"

Subgraph pattern matching; Reachability; Shortest path; Keyword search; Historical or Temporal queries...

Continuous "queries" and Realtime analytics

Online prediction; Monitoring; Anomaly/Event detection

Batch analysis tasks/Network science

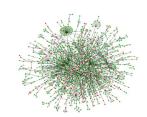
Centrality analysis: Find the most central nodes in a network

Community detection: Partition vertices into groups with dense interactions

Network evolution: Build models for network formation and evolution

Network measurements: Measure statistical properties

Graph cleaning/inference: Remove noise in the observed network data





Social networks

A protein-protein interaction network

Citation networks

Communication networks

Disease transmission networks



Financial transaction networks

Stock

Know

Different types of "queries"

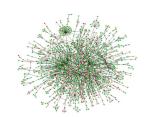
Subgraph pattern matching; Reachability; Shortest path; Keyword search; Historical or Temporal queries...

Continuous "queries" and Realtime analytics

Online prediction; Monitoring; Anomaly/Event detection

Batch analysis tasks

Centrality analysis; Community detection; Network evolution; Network measurements; Graph cleaning/inference



facebook. twitter
Linked in
Social networks

Federal funds networks

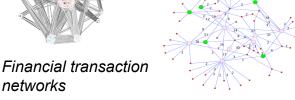
A protein-protein interaction network

Citation networks

Communication networks

Disease transmission networks





Stock Trading Networks

Machine learning tasks

Many algorithms can be seen as message passing in specially constructed graphs

Graph Data Management: State of the Art

- Much prior and ongoing work most of it outside, or on top of, general-purpose data management systems
 - Specialized indexes or algorithms for specific types of queries
 - Stand-alone prototypes for specific analysis tasks
- Emergence of specialized graph databases in recent years
 - Neo4j, Titan, OrientDB, DEX, AllegroGraph, ...
 - Rudimentary declarative interfaces/query languages
- Several "vertex-centric" frameworks in recent years
 - Pregel, Giraph, GraphLab, GRACE, GraphX, ...
 - Only work well for a very limited set of tasks
- Little work on continuous/real-time query processing, or on supporting evolutionary or temporal analytics

What we are doing

- Goal: A graph data management system with unified declarative abstractions for graph queries and analytics
- Work so far
 - Declarative graph cleaning [GDM'11, SIGMOD Demo'13]
 - NScale: a distributed programming framework [VLDB Demo'14]
 - Real-time continuous queries [SIGMOD'12, ESNAM'14, SIGMOD'14]
 - Techniques for continuous query processing over large dynamic graphs
 - Expressive query language for specifying anomaly detection queries
 - Historical graph data management [ICDE'13, SIGMOD Demo'13]
 - New indexing structure for retrieving historical snapshots
 - A high-level temporal/evolutionary analytics framework
 - Subgraph pattern matching and counting [ICDE'12, ICDE'14]
 - GraphGen: graph analytics over relational data [VLDB Demo'15]

Outline

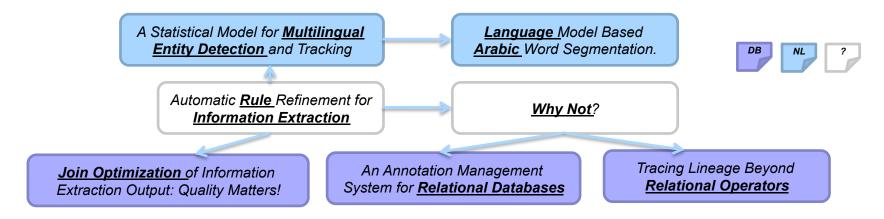
- Graph Data Management
 - Declarative Graph Cleaning
 - A Framework for Distributed Graph Analytics
- DataHub: A platform for collaborative data science
 - Recreation/Storage Tradeoff in Version Management

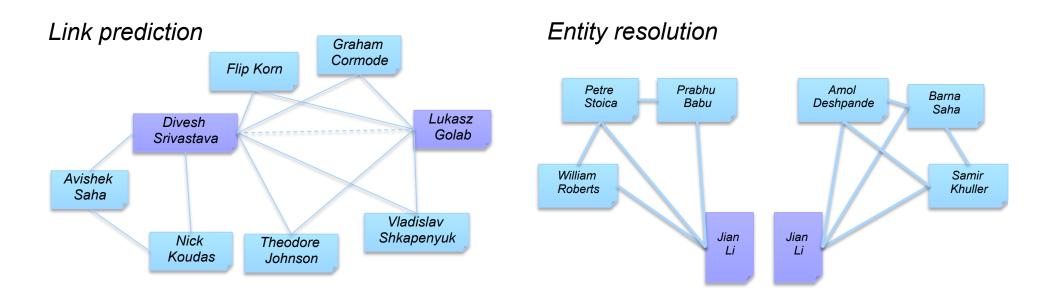
Motivation

- The observed, automatically-extracted information networks are often noisy and incomplete
- Need to extract the underlying true information network through:
 - Attribute Prediction: to predict values of missing attributes
 - Link Prediction: to infer missing links
 - Entity Resolution: to decide if two references refer to the same entity

Collective (relational) Inference

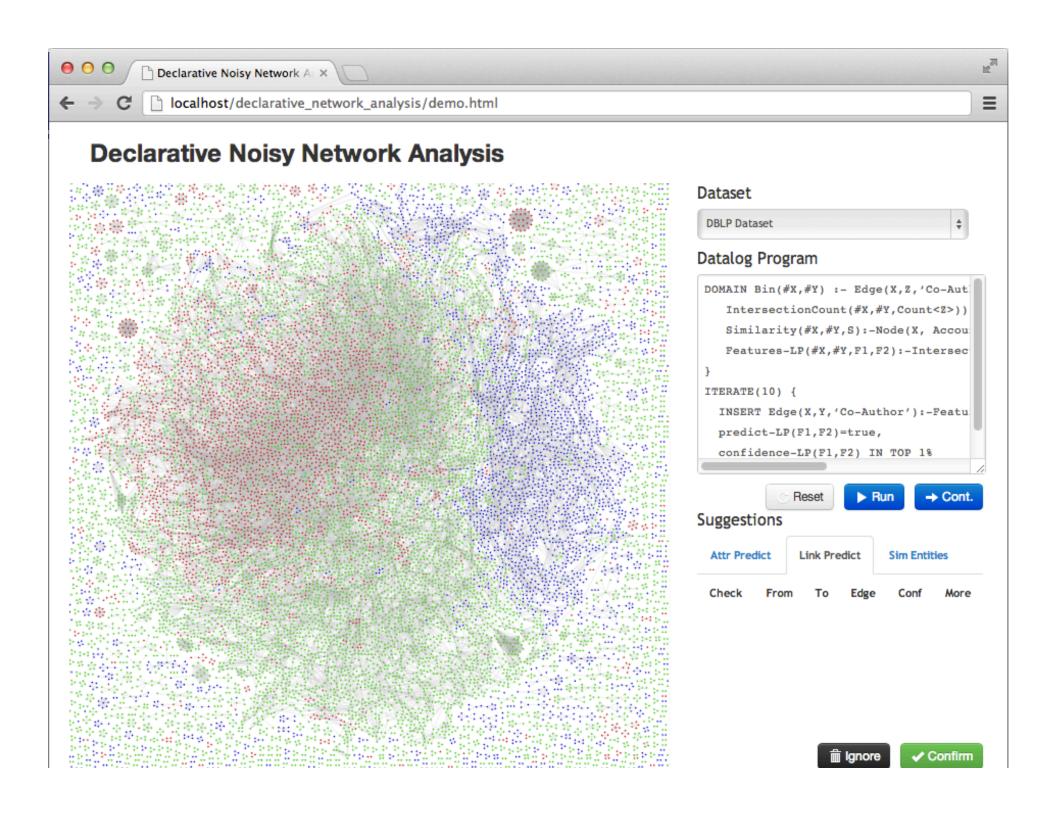
Attribute prediction: Predict topic of the paper

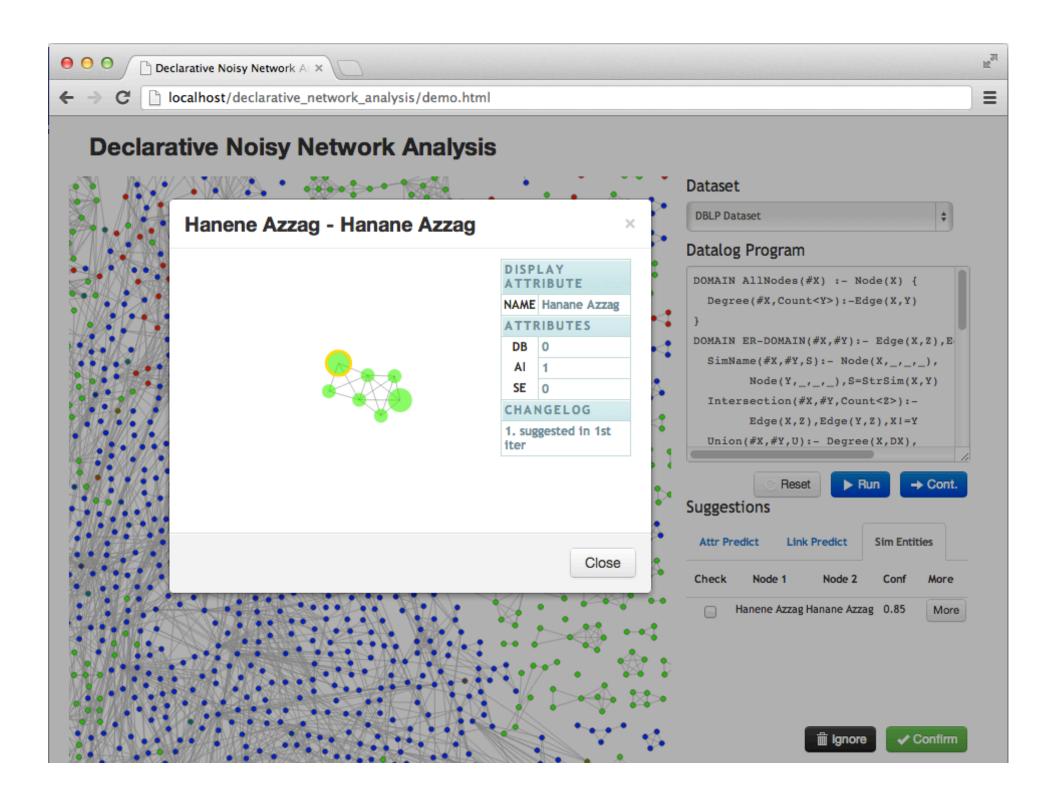




Motivation

- The observed, automatically-extracted information networks are often noisy and incomplete
- Need to extract the underlying true information network through:
 - Attribute Prediction: to predict values of missing attributes
 - Link Prediction: to infer missing links
 - Entity Resolution: to decide if two references refer to the same entity
- Typically iterative and interleaved application of the techniques
 - Use results of one to improve the accuracy of other operations
 - Significant benefits to using graph structure ("collective" inference)
- Numerous techniques developed for the tasks in isolation
 - No support from data management systems
 - Hard to evaluate new techniques, especially for joint inference

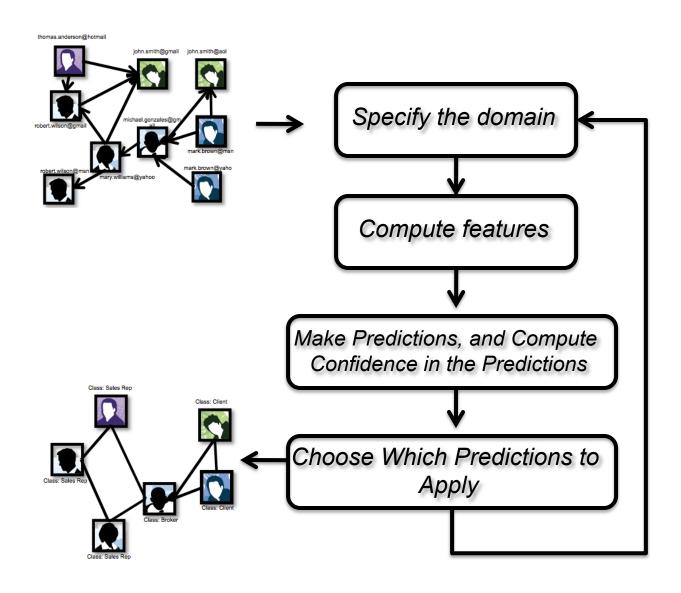




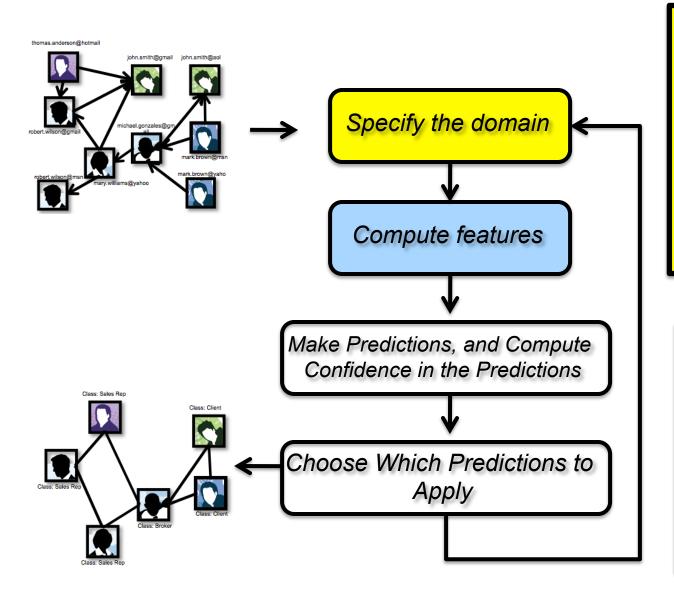
Overview of the Approach

- Declarative specification of the cleaning task
 - Datalog-based language for specifying --
 - Prediction features (including local and relational features)
 - The details of how to accomplish the cleaning task
 - Arbitrary interleaving or pipelining of different tasks

Task Specification Framework



Task Specification Framework



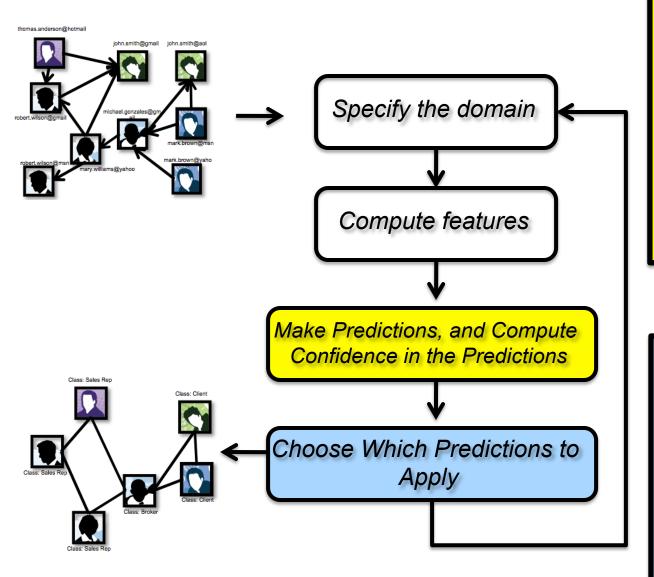
For attribute prediction, the domain is a subset of the graph nodes.

For link prediction and entity resolution, the domain is a subset of pairs of nodes.

Local: word frequency, income, etc.

Relational: degree, clustering coeff., no. of neighbors with each attribute value, common neighbors between pairs of nodes, etc.

Task Specification Framework



Attribute prediction: the missing attribute

Link prediction: add link or not?

Entity resolution: merge two nodes or not?

After predictions are made,

the graph changes:

Attribute prediction changes local attributes.

Link prediction changes the graph links.

Entity resolution changes both local attributes and graph links.

Overview of the Approach

- Declarative specification of the cleaning task
 - Datalog-based language for specifying --
 - Prediction features (including local and relational features)
 - The details of how to accomplish the cleaning task
 - Arbitrary interleaving or pipelining of different tasks
- A mix of declarative constructs and user-defined functions to specify complex prediction functions
- Prototype implementation using Java BerkeleyDB
 - Datalog rules converted into SQL
 - Optimized the execution through caching, incremental evaluation, and pre-computed data structures

Example

- Real-world PubMed graph
 - Set of publications from the medical domain, their abstracts, and citations
- 50,634 publications, 115,323 citation edges
- Task: Attribute prediction
 - Predict if the paper is categorized as Cognition, Learning, Perception or Thinking
- Choose top 10% predictions after each iteration, for 10 iterations

Outline

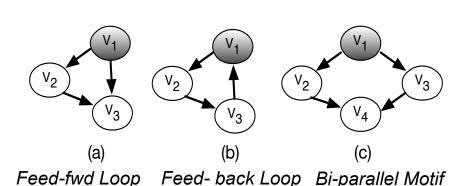
- Graph Data Management
 - Declarative Graph Cleaning
 - A Framework for Distributed Graph Analytics
- DataHub: A platform for collaborative data science
 - Recreation/Storage Tradeoff in Version Management

Scaling Graph Analysis Tasks

- Graph analytics/network science tasks too varied
 - Centrality analysis; evolution models; community detection
 - Link prediction; belief propagation; recommendations
 - Motif counting; frequent subgraph mining; influence analysis
 - Outlier detection; graph algorithms like matching, max-flow

Friends in

An active area of research in itself...



CS dept

Colleagues

Colleagues

Colleagues

Colleagues

Friends in friends

database lab in CS dept

Friends members

Office

Work place friends

High school

Counting network motifs

Identify Social circles in a user's ego network

Scaling Graph Analysis Tasks

- Graph analytics/network science tasks too varied
 - Centrality analysis; evolution models; community detection
 - Link prediction; belief propagation; recommendations
 - Motif counting; frequent subgraph mining; influence analysis
 - Outlier detection; graph algorithms like matching, max-flow
 - An active area of research in itself...
- Hard to build general platforms like MapReduce
 - What is a good programming abstraction to provide?
 - Needs to cover a large fraction of use cases, and be easy to use
 - MR not a good fit for graph analytics
 - No clear winner yet, so little progress on systems
 - Especially on distributed or parallel systems
 - Application developers largely doing their own thing

"Vertex-centric" Frameworks

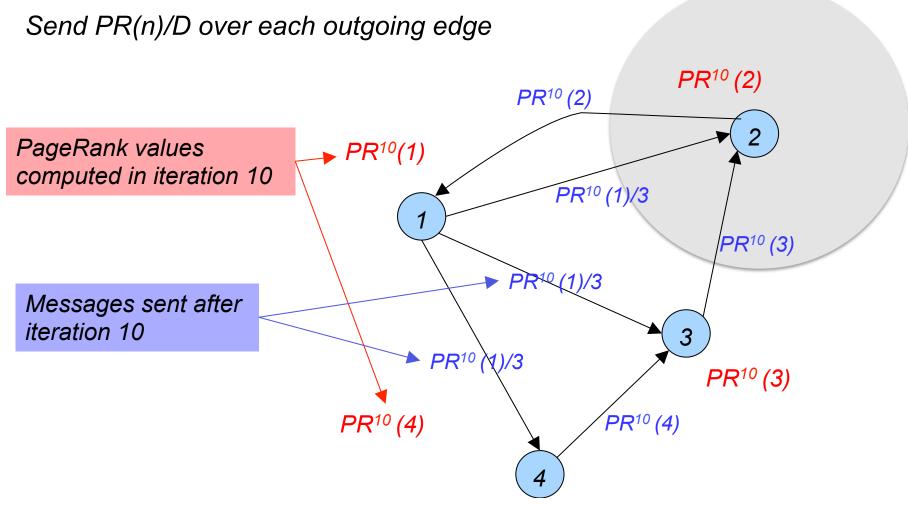
- Introduced by Google in a system called "Pregel"
 - Inspired by BSP (Bulk Synchronous Protocol)
- Adopted by many other systems
 - GraphLab, Apache Giraph, GraphX, Xstream, ...
 - Most of the research, especially in databases, focuses on it
- "Think like a vertex" paradigm
 - User provides a single compute() function that operates on a vertex
 - Executed in parallel on all vertices in an iterative fashion
 - Exchange information at the end of each iteration through message passing

Example: PageRank

Compute() at Node n:

PR(n) = sum up all the incoming weights

Let the outDegree be D



Programming Frameworks

- Vertex-centric framework
 - Works well for some applications
 - Pagerank, Connected Components, ...
 - Some machine learning algorithms can be mapped to it
 - However, the framework is very restrictive
 - Most analysis tasks or algorithms cannot be written easily
 - Simple tasks like counting neighborhood properties infeasible
 - Fundamentally: Not easy to decompose analysis tasks into vertex-level, independent local computations
- Alternatives?
 - Galois, Ligra, GreenMarl: Not sufficiently high-level
 - Some others (e.g., Socialite) restrictive for different reasons

Example: Local Clustering Coefficient

LCC(n) = neighborhood density around n Compute() at Node n:

Need to count the no. of edges between neighbors

But does not have access to that information

Option 1: Each node transmits its list of

neighbors to its neighbors

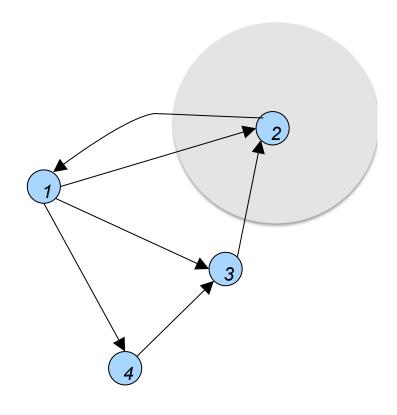
Huge memory consumption

Option 2: Allow access to neighbors' state

Neighbors may not be local

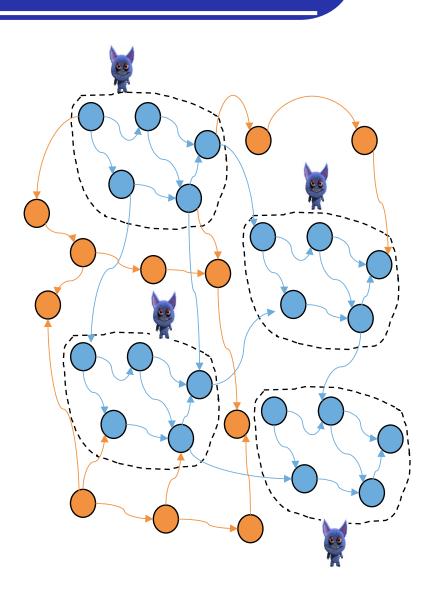
What about computations that require 2-

hop information?



NScale Programming Framework

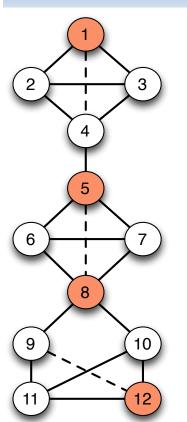
- An end-to-end distributed graph programming framework
- Users/application programs specify:
 - Neighborhoods or subgraphs of interest
 - A kernel computation to operate upon those subgraphs
- Framework:
 - Extracts the relevant subgraphs from underlying data and loads in memory
 - Execution engine: Executes user computation on materialized subgraphs
 - Communication: Shared state/ message passing



NScale programming model

Underlying graph data on HDFS

Subgraph extraction query:



```
Compute (LCC) on 
Extract ({Node.color=orange} {k=1} 
{Node.color=white} 
{Edge.type=solid} 
)
```

Query-vertex predicate

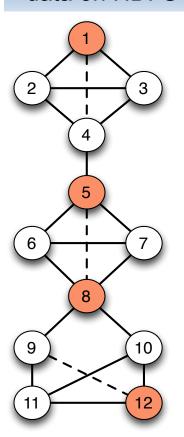
Neighborhood Size

Neighborhood vertex predicate

Neighborhood edge predicate

NScale programming model

Underlying graph data on HDFS

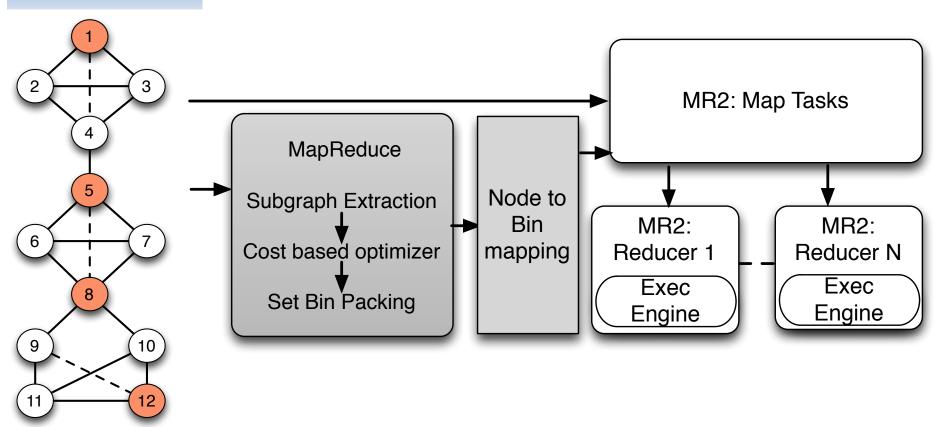


Specifying Computation: BluePrints API

Program cannot be executed as is in vertex-centric programming frameworks.

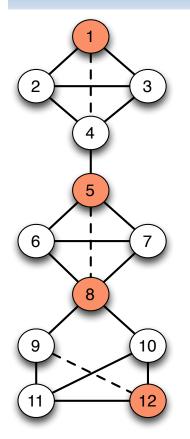
GEP: Graph extraction and packing

Underlying graph data on HDFS

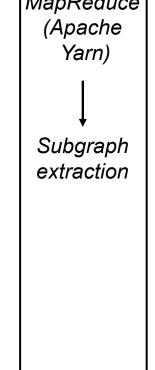


GEP: Graph extraction and packing

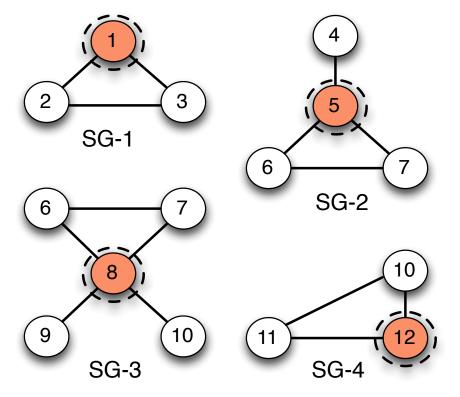
Underlying graph data on HDFS



Graph Extraction and Loading MapReduce



Extracted Subgraphs



GEP: Graph extraction and packing

Goal:

- Group graphs with high similarity
- Minimizes memory consumption

Techniques explored

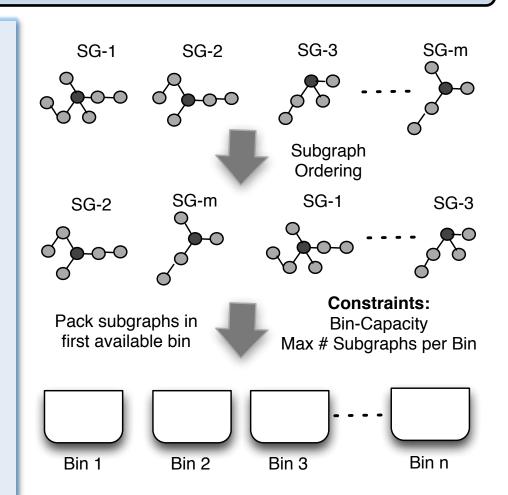
 Set bin packing, graph partitioning, clustering

Shingle based set bin packing

- Min-hash signatures based sorting
- Grouping based on Jaccard similarity

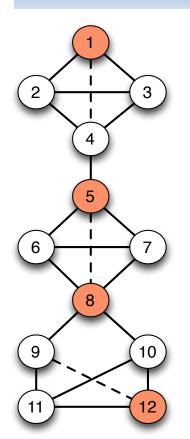
Bin Packing

- Set union operation
- Bin Capacity: Elastic resource allocation
- Max # Subgraphs: Handles Skew

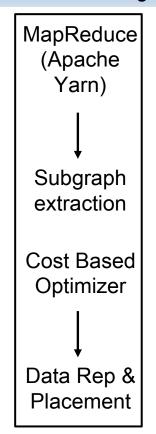


GEP: Graph extraction and packing

Underlying graph data on HDFS

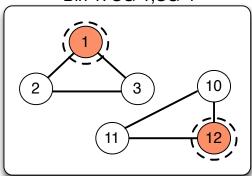


Graph Extraction and Loading

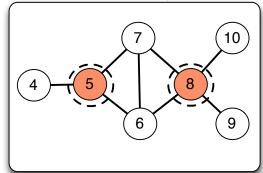


Sample bin packing using Shingles

Bin 1: SG-1,SG-4

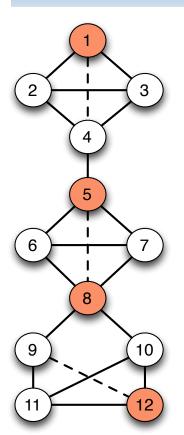


Bin 2: SG-2, SG-3



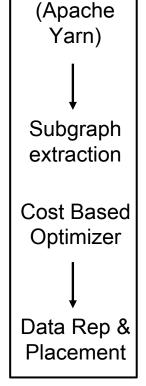
GEP: Graph extraction and packing

Underlying graph data on HDFS

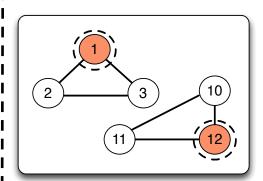


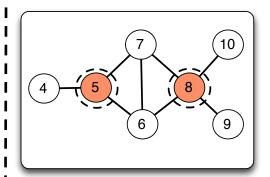
Graph Extraction I and Loading

MapReduce



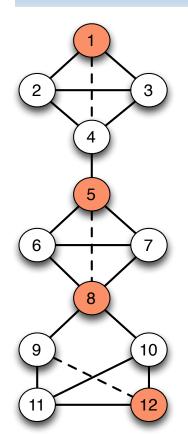
Subgraphs in Distributed Memory





Distributed execution of user computation

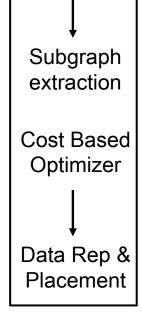
Underlying graph data on HDFS



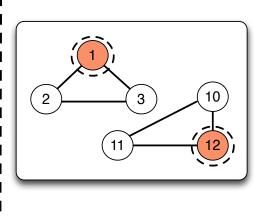
Graph Extraction and Loading

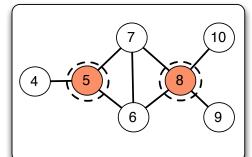
MapReduce (Apache

Yarn)

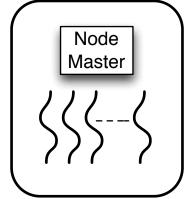


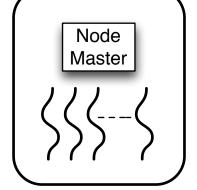
Subgraphs in Distributed Memory





Distributed Execution Engine





NScale: Summary

- Users write programs at the abstraction of a graph
 - More intuitive for graph analytics
 - Captures mechanics of common graph analysis/cleaning tasks
- Generalization: Flexibility in subgraph definition
 - Subgraph = vertex and associated edges: vertex-centric programs
 - Subgraph = an entire graph: global programs
- Scalability
 - Only relevant portions of the graph data loaded into memory
 - User can specify subgraphs of interest, and select nodes or edges based on properties
 - Carefully partition (pack) nodes across machines so that:
 - Every subgraph is entirely in memory on a machine, while using very few machines

Experimental Evaluation

Datasets

- Web graphs
- Communication/interaction graphs
- Social networks

Graph applications

- Local Clustering Coefficient
- Motif counting
- Identifying weak ties
- Triangle Counting
- Personalized Page Rank

Baselines

- Apache Giraph
- GraphLab
- GraphX

Evaluation Metrics

- Computational Effort
- Execution Time
- Cluster Memory

Cluster Setup

- 16 Node Cluster
- Apache YARN (MRv2)
- Each Node:
 - 2 x 4-core Intel Xeon
 - 24GB RAM, 3 x 2 TB disks

Experimental Evaluation

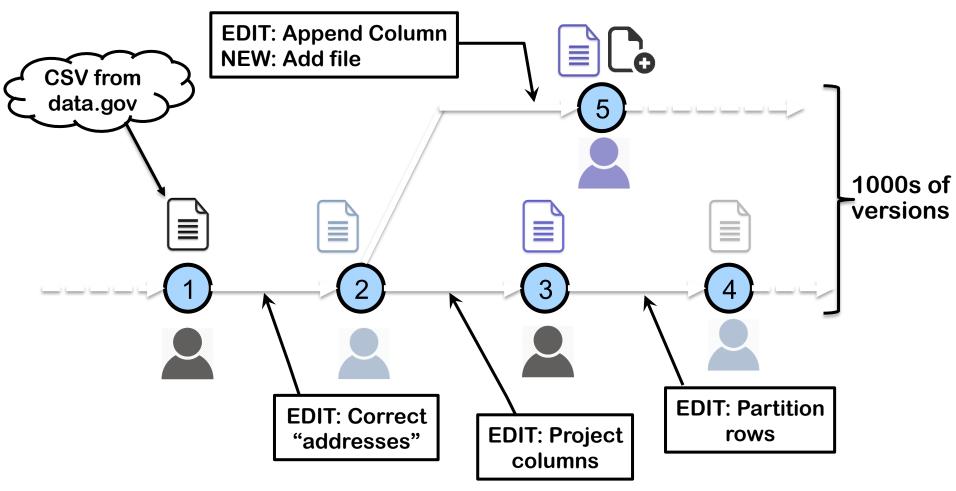
		Local Clustering Coefficient													
Dataset	1	NScale			Giraph				GraphLab				GraphX		
	CE (Node Secs)	- Cluster Mem (G	- , -				er CE (GB) Secs		,		Cluster (Mem (GB)		de-	Cluster Mem (GB)	
EU Email	377	9.00	-	1150		26.17		365		20.10		225		4.95	
NotreDame	620	19.07	1	1564		30.14		550		21.40		340		9.75	
Google Web	658	25.82	2	2024		35.35		600		33.50		1485		21.92	
WikiTalk	726	24.16	[DNC		ООМ		1125		37.	37.22			32.00	
LiveJournal	1800	50.00	[DNC		ООМ		550	5500		3.62	4515		84.00	
Orkut	2000	62.00	DNC		ООМ			DNC		00	OOM 2			125.00	
		Personalized Page Rank on 2-Hop Neighborhood													
Dataset		NSo	NScale		Gira		aph		G	GraphLab			GraphX		
	#Source Vertices	CE (Node- Secs)	Cluster Mem (GB)		CE (Node- Secs)		Cluster Mem (GB)		CE (Node- Secs)		Cluster Mem (GB)	•	Node-)	Cluster Mem (GB)	
EU Email	3200	52	3.35		782		17.10		710		28.87	9975	5	85.50	
NotreDame	3500	119	9.56		1058		31.76		870		70.54	5059	95	95.00	
Google Web	4150	464	21.52		10482		64.16		1080		108.28	DNC		-	
WikiTalk	12000	3343	79.43		DNC		ООМ		DNC		ООМ	DNC		-	
LiveJournal	20000	4286	84.94		DNC		ООМ		DNC		ООМ	DNC		-	
Orkut	20000	4691	93.07		DNC		OOM		DNC		ООМ	DNC		-	

Outline

- Graph Data Management
 - Declarative Graph Cleaning
 - A Framework for Distributed Graph Analytics
- DataHub: A platform for collaborative data science
 - Recreation/Storage Tradeoff in Version Management

Collaborative Data Science

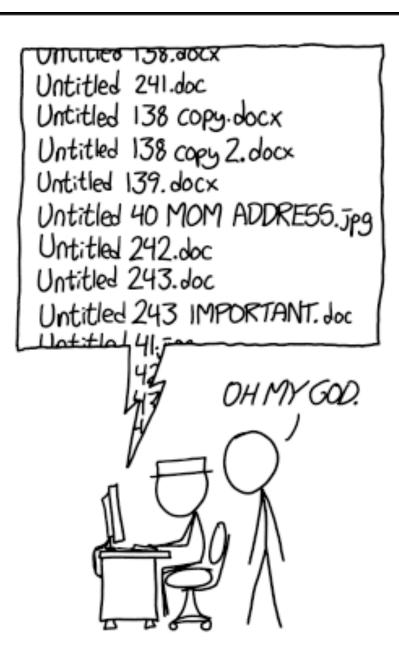
Widespread use of "data science" in many many domains



A typical data analysis workflow

Collaborat

- Widespread
- Increasingly especially d
 - Many prival
 - No easy wa
 - Manual int
 - No efficient
 - No way to
- Ad hoc data
 - Much of th
 - The proces
 - Scientists/



PROTIP: NEVER LOOK IN SOMEONE. ELSE'S DOCUMENTS FOLDER. hy many domains g the *process*,

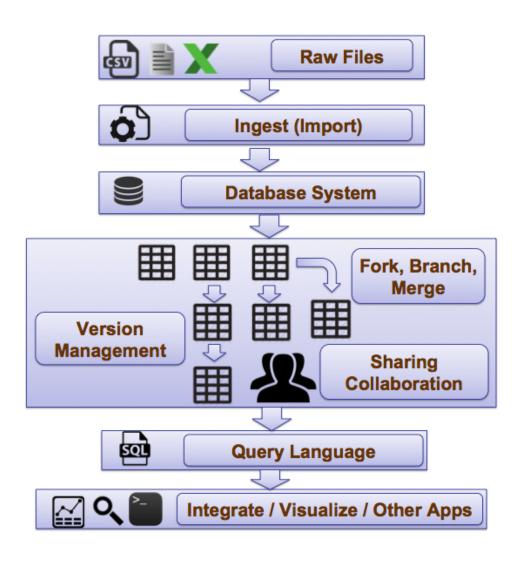
sive redundancy between datasets pnflicts datasets of a dataset

Dropbox) used y can't use DBs hoc and exploratory uch on their own

DataHub: A Collaborative Data Science Platform

The one-stop solution for collaborative data science and dataset version management

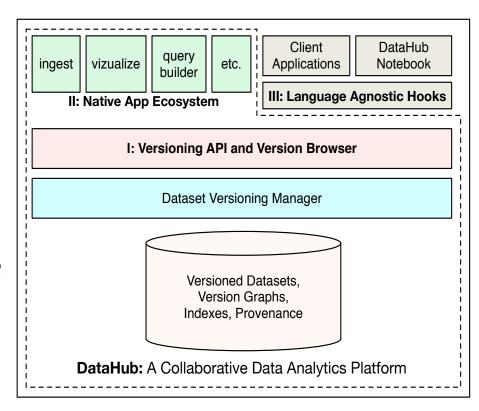
http://data-hub.org



DataHub: A Collaborative Data Science Platform

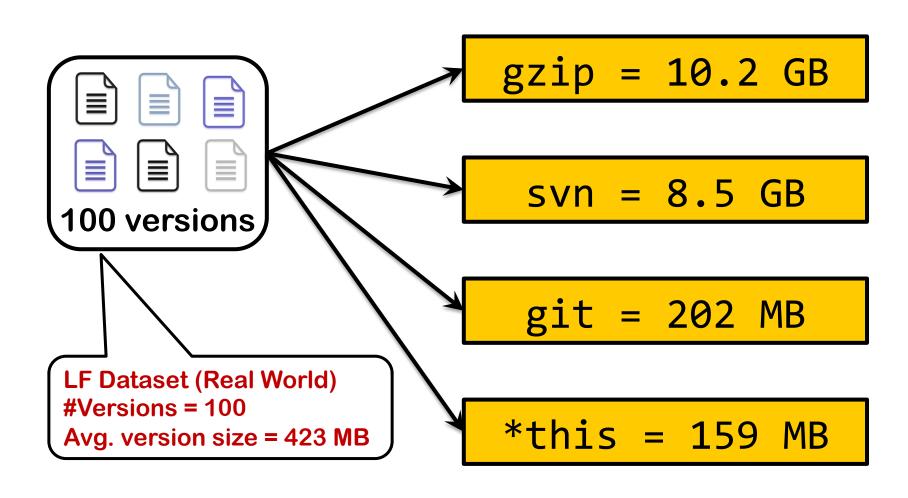
The one-stop solution for collaborative data science and dataset version management

- a dataset management system import, search, query, analyze a large number of (public) datasets
- a dataset version control system –
 branch, update, merge, transform large structured or unstructured datasets
- an app ecosystem and hooks for external applications (Matlab, R, iPython Notebook, etc)

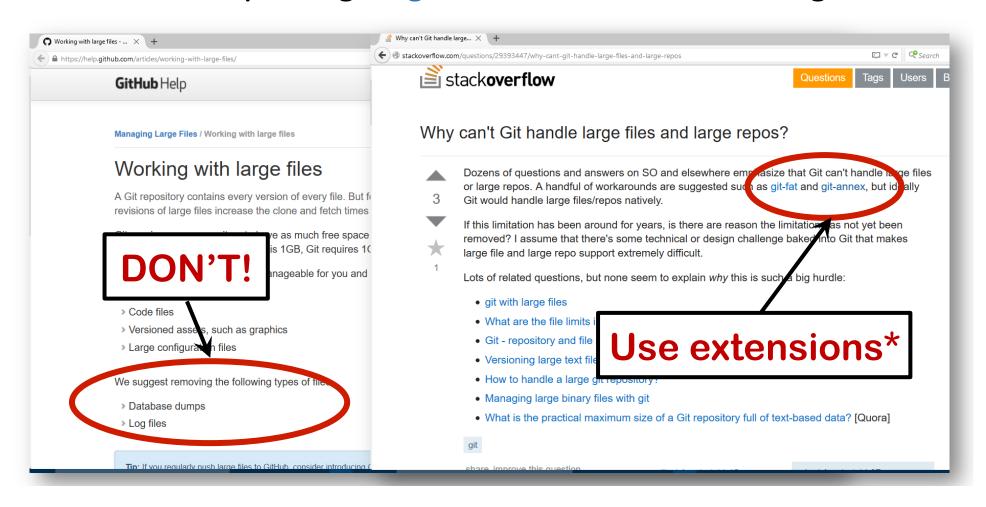


DataHub Architecture

No, because they typically use fairly simple algorithms and are optimized to work for code-like data



- No, because they typically use fairly simple algorithms and are optimized to work for code-like data
- Git ends up using large amounts of RAM for large files



- No, because they typically use fairly simple algorithms and are optimized to work for code-like data
- ★ Git ends up using large amounts of RAM for large files
- Querying and retrieval functionalities are primitive, and revolve around single version and metadata retrieval
- No way to specify queries like:
 - identify all datasets derived of dataset A that satisfy property P
 - identify all predecessor versions of version A that differ from it by a large number of records
 - rank a set of versions according to a scoring function
 - find the version where the result of an aggregate query is above a threshold
 - find parent records of all records in version A that satisfy certain property

- No, because they typically use fairly simple algorithms and are optimized to work for code-like data
- **✗** Git ends up using large amounts of RAM for large files
 - VQuel: A Unified Query Language for querying versioning and derivation information [USENIX TAPP'15]
 - Example: What commits did Alice make after January 01, 2015?

- Illia the version where the result of an aggregate query is above a threshold
- find parent records of all records in version A that satisfy certain property

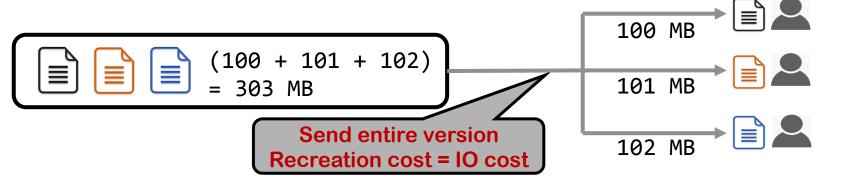
Outline

- Graph Data Management
 - Declarative Graph Cleaning
 - A Framework for Distributed Graph Analytics
- DataHub: A platform for collaborative data science
 - Recreation/Storage Tradeoff in Version Management

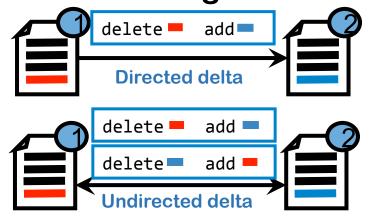
Storage cost is the space required to store a set of versions



Recreation cost is the time* required to access a version



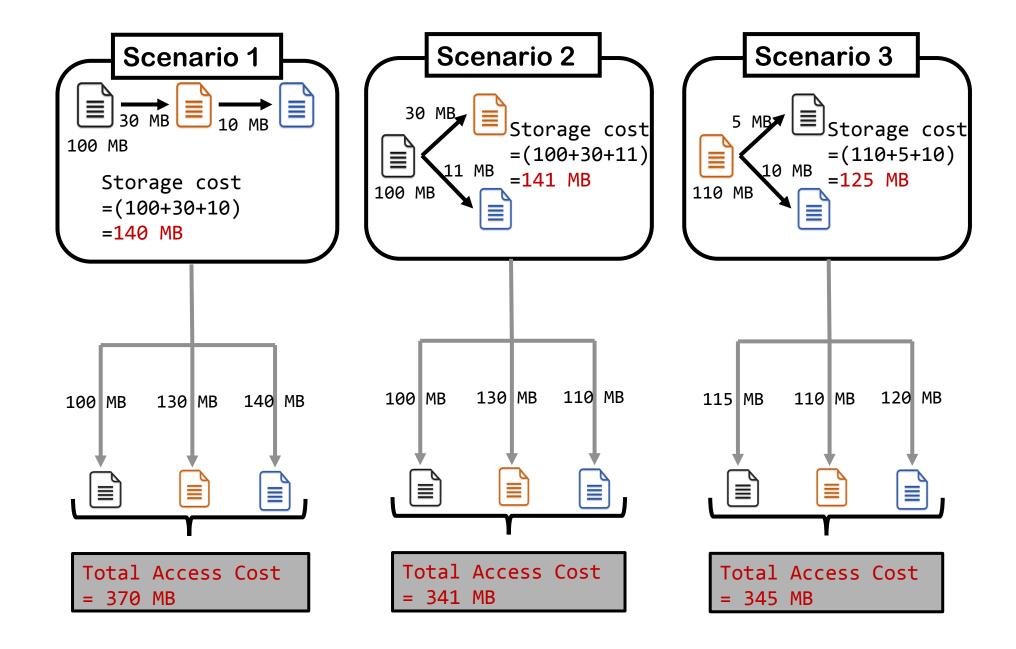
A delta between versions is a file which allows constructing one version given the other



Example: Unix diff, xdelta, XOR, etc.

A delta has its own storage cost and recreation cost, which, in general, are independent of each other

Storage-Recreation Tradeoff



Storage-Recreation Tradeoff

Given

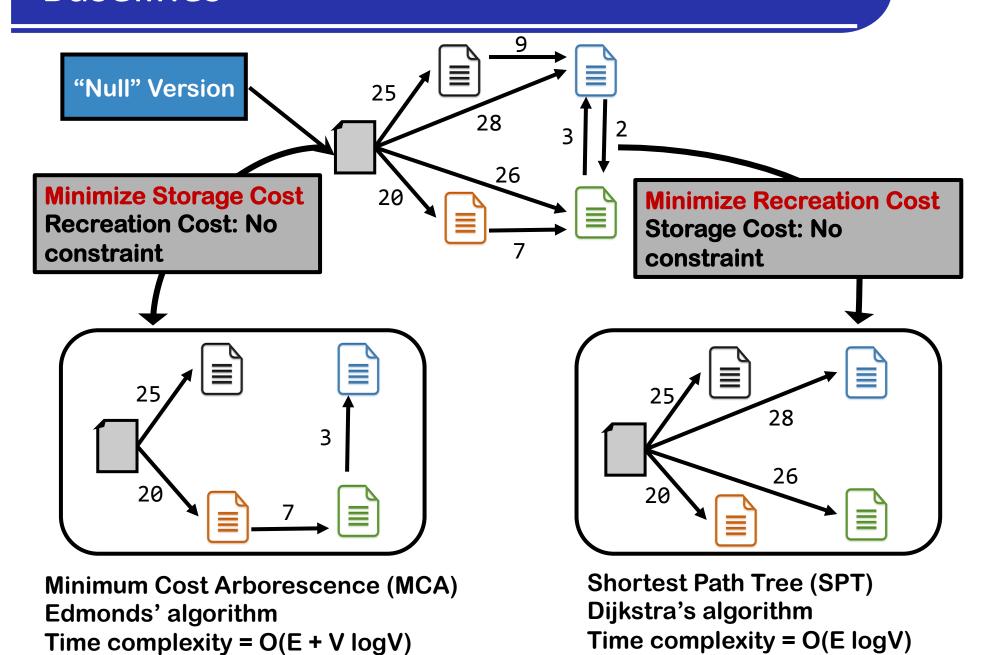
- a set of versions
- 2) partial information about deltas between versions

Find a Storage Solution that:

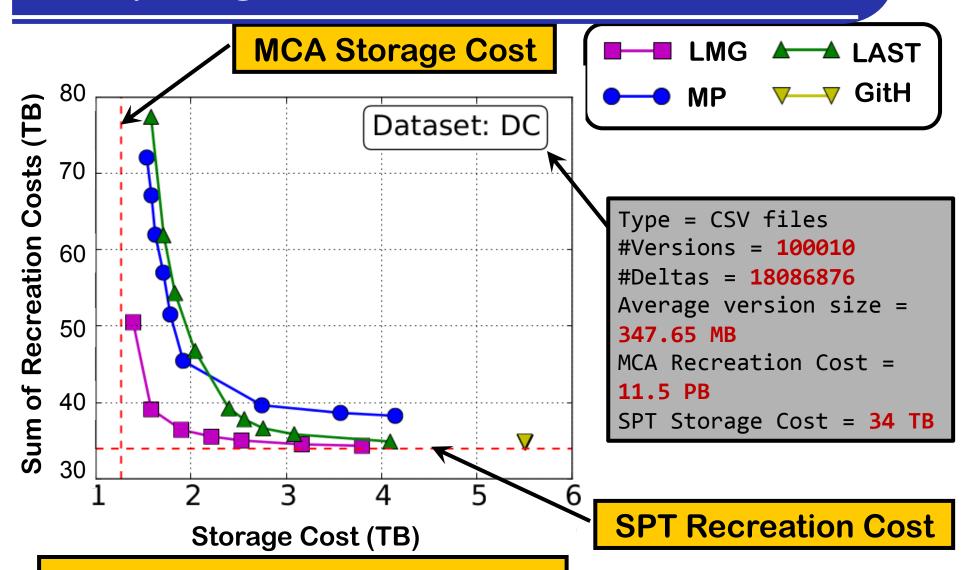
- minimizes total recreation cost given a storage budget, or
- minimizes max recreation cost given a storage budget

	Storage Cost	Recreation Cost	Undirected Case, $\Delta = \Phi$	Directed Case, $\Delta = \Phi$	Directed Case, Δ≠Φ			
P1	min C	R _i < ∞, ∀ i	PTime, Minimum Cost Arborescence (MCA)					
P2	C < ∞	$\min \left\{ \max \left\{ R_i \mid 1 \le i \le n \right\} \right\}$	PTime, Shortest Path Tree (SPT)					
Р3	C≤β	$\min\left\{\sum_{i=1}^{n}R_{i}\right\}$	NP-hard,	NP-hard, LMG Algorithm				
P4	C≤β	$\min \left\{ \max \left\{ R_i \mid 1 \le i \le n \right\} \right\}$	LAST* Alg	NP-hard, MP Algorithm				
P5	min C	$\sum_{i=1}^{n} R_i \leq \theta$	NP-hard,	NP-hard, LMG Algorithm				
P6	min C	$\max \{R_i \mid 1 \le i \le n\} \le \theta$	LAST* Alg	NP-hard, MP Algorithm				

Baselines



Comparing Different Solutions



Storage budget of 1.1X the MCA reduces total recreation cost by 1000X