# Part 1: Systems Infrastructure for Big-Data Analytics

#### Ananth Grama

Center for Science of Information Dept. of Computer Science, Purdue University

With help from Giorgos Kollias, Naresh Rapolu, Karthik Kambatla, and Adnan Hassan

Thanks to NSF, DoE, the Center for Science of Information, Microsoft, and Intel.

- Motivating applications PageRank, Graph Alignment
- Case study: single *mat-vec* in MapReduce
- Asynchrony and speculation: Optimizations across iterations
  - Asynchronous algorithms through Relaxed Synchronization
  - Speculative parallelism through TransMR (transactional *MapReduce*)
  - Locking techniques for efficient distribution transactions
- Online Analytics
  - Motivating applications
  - Streams (Storm) and Machine Learning (Vowpal Wabbit)
  - Runtime and optimizations

# Motivating Example: Functional PageRank (PR)

#### Computing PageRank (PR)

- PageRank as a random surfer process: Start surfing from a random node and keep following links with probability μ restarting with probability 1 μ; the node for restarting will be selected based on a personalization vector v. The ranking value x<sub>i</sub> of a node i is the probability of visiting this node during surfing.
- PR can also be cast in power series representation as  $x = (1 \mu) \sum_{j=0}^{k} \mu^{j} S^{j} v$ ; S encodes column-stochastic adjacencies.

#### Functional rankings

- A general method to assign ranking values to graph nodes as  $x = \sum_{j=0}^{k} \zeta_j S^j v$ . PR is a functional ranking,  $\zeta_j = (1 \mu)\mu^j$ .
- Terms attenuated by outdegrees in S and damping coefficients  $\zeta_j$ .

ヘロン 人間 とくほと 人 ほとう

# Motivating Example: Graph Matching



- Node similarity: Two nodes are similar if they are linked by other similar node pairs. By pairing similar nodes, the two graphs become *aligned*.
- Let à and B̃ be the normalized adjacency matrices of the graphs (normalized by columns), H<sub>ij</sub> be the independently known similarity scores (preferences matrix) of nodes i ∈ V<sub>B</sub> and j ∈ V<sub>A</sub>, and μ be the fractional contribution of topological similarity.
- To compute X, IsoRank iterates:

$$X \leftarrow \mu \tilde{B} X \tilde{A}^T + (1 - \mu) H$$

## MapReduce: Basics

Execute in parallel user-defined functions on individual data-items distributed across machines.

- Simple programming model map and reduce functions
- Scalable, distributed execution of these functions on massive amounts of data on commodity hardware



э

#### MapReduce: DataFlow

input



6 / 53

Ξ

# Example: Shotest Path (*mat-vec*) in MapReduce

```
map(node, adjList) {
  for each arc in adjList {
    output(arc.dst, node.distance + arc.weight)
  }
}
```

```
reduce(node, newDistList) {
   node.distance = min(newDistList)
}
```

while(not converged) {
 runMapReduceJob(map, reduce, Ax)
}

 We call this the Naive mat-vec — map takes an adjacency list as input

• Pegasus (CMU) implementation reads one edge per map

#### Naive *mat-vec*: Resource Utilizaton



Ananth Grama (Purdue University)

3 / 53

## Optimizing the *mat-vec* further



- Stages overlap; hence, sum of percentages > 100
- I/O time > Computation time
- For performance:
  - Batch read data
  - Each *map* processes more data (as much as can fit in memory)

#### Partitioned mat-vec

- Each *map* operates on a *graph partition* a set of adjacency lists
- Partition size is constrained by the heap size
- On our setup, maximum partition size was 2<sup>19</sup> nodes



#### Performance: Pegasus vs Naive vs Partitioned



Algorithmic optimizations:

- Asynchronous algorithms: Algorithms that allow asynchrony

   ordering of updates doesn't affect the correctness of the
   algorithm
  - e.g., PageRank, Alignments
- **Speculative Parallelism**: Algorithms where concurrent computations can have potential conflicts, the conflicts are rare and can only be detected at runtime e.g., Boruvka's MST, Single Source Shortest Path

- Improve performance in parallel environments.
  - Infrequent synchronization reduces communication
  - Examples
    - Graph algorithms, Numerical methods, ML kernels, etc.
- More pronounced gains in distributed environments
  - Higher communication and data-movement costs
  - $\bullet\,$  Read once, process multiple times allows more computation for the same I/O cost

- Synchronize once every few iterations
- Approach: Two levels of MapReduce
  - Global MapReduce: The regular MapReduce
    - Requires global synchronization
  - Local MapReduce: MapReduce within a global map
    - Each global map runs a few iterations of local MapReduce
    - Partial synchronization of data of a single global map task
  - Input data partitioning for fewer dependencies across partitions

## PageRanks Using Traditional MapReduce



## PageRank with Relaxed Synchronization



æ

Image: Image:

- ∢ ∃ ▶

-∢∃>

#### Realizing Relaxed Synchronization Semantics

#### • Code gmap, greduce, Imap, Ireduce

- Imap, Ireduce use EmitLocalIntermediate() and Emit Local()
- Synchronized hashtables for local storage

```
gmap(xs : X list) {
  while(no-local-convergence-intimated) {
    for each element x in xs {
        lmap(x); // emits lkey, lval
    }
    lreduce(); // operates on the output of lmap functions
  }
  for each value in lreduce-output{
    EmitIntermediate(key, value);
  }
}
```

## Evaluation — Relaxed Synchronization

- Sample applications
  - Single Source Shortest Path (MST, transitive closure, etc.)
  - PageRank (mat-vec: eigen value and linear system solvers)
- Experimental Testbed
  - 8 Amazon EC2 Large Instances
    - 64-bit compute units with 15 GB RAM, 4x 420 GB storage
    - Hadoop 0.20.1; 4 GB heap space per slave

#### Single Source Shortest Path



э

- Input: Partitioned using METIS (less than 10 seconds)
- Damping factor = 0.85

|       | GraphA    | GraphB    |
|-------|-----------|-----------|
| Nodes | 280,000   | 100,000   |
| Edges | 3 million | 3 million |

æ

イロト イポト イヨト イヨト

#### PageRank Performance: GraphA



э

#### PageRank Performance: GraphB



э

#### Speculative parallelism

- Most of the data can be operated on in parallel.
- Some executions conflict. These can only be detected at runtime. [Pingali et.al., PLDI'11]
- Online algorithms/ Pipelined workflows
  - MapReduce Online [Condie'10] is an approach needing havy checkpointing.
- Software Transactional Memory (STM)

- Goal: Exploit speculative data-parallelism
- Support data-sharing across concurrent computations to detect and resolve conflicts at runtime
- Solution:
  - Use distributed key-value stores as shared address space across computations
  - Address inconsistencies arising due to the disparate fault-tolerance mechanisms
  - Transactional execution of map and reduce functions

#### TransMR: System Architecture

• Distributed key-value store provides a shared-memory abstraction to the distributed execution-layer.



- Data-centric function scope Map/Reduce/Merge etc, termed as a Computation Unit (CU), is executed as a transaction.
- Optimistic reads and write-buffering. Local Store (LS) forms the write-buffer of a CU.
  - Put (K, V): Write to LS, which is later atomically committed to GS.
  - Get (K, V): Return from LS, if already present; otherwise, fetch from GS and store in LS.
  - Other Op: Any thread local operation.
- The output of a CU is always committed to the GS before being visible to other CUs of the same or different type.
  - Eliminates the costly shuffle phase of MapReduce.

# **Design Principles**

- Optimistic concurrency control over pessimistic locking
  - Locks are acquired at the end of the transaction. Write-buffer and read-set is validated against those of concurrent Trx assuring serializability.
  - Client is potentially executing on the slowest node in the system; in this case, pessimistic locking hinders parallel transaction execution.
- Consistency (C) and Tolerance to Network Partitions (P) over Availability (A) in CAP Theorem for Distributed transactions.
  - Application correctness mandates strict consistency of execution. Relaxed consistency models are application-specific optimizations.
  - Intermittent non-availability is not too costly for batch-processing applications, where client is fault-prone in itself.

イロト 不得下 イヨト イヨト

- We show performance gains on two applications, which are hitherto implemented sequentially without transactional support; both exhibit Optimistic data-parallelism.
- Boruka's MST
  - Each iteration is coded as a Map function with input as a node. Reduce is an identity function. Conflicting maps are serialized while others are executed in parallel.
  - After n iterations of coalescing, we get the MST of an n node graph.
  - A graph of 100 thousand nodes, with average degree of 50, generated based on the forest-fire model.

## Boruvka's MST

• Speedup of 3.73 on 16 nodes, with less than 0.5 % re-executions due to aborts.



# Maximum Flow Using Push-Relabel Algorithm

- Each Map function executes a Push or a Relabel operation on the input node, depending on the constraints on its neighbors.
- Push operation increases the flow to a neighboring node and changes their "Excess".
- Relabel operation increases the height of the input node if it is the lowest among its neighbors.
- Conflicting Maps operating on neighboring nodes get serialized due to their transactional nature.
- Only sequential implementation possible without support for runtime conflict detection.

#### Maximum flow using Push-Relabel algorithm

• Speedup of 4.5 is observed on 16 nodes with 4% re-executions on a windown of 40 iterations.



#### TransMR: Intermediate Lessons

- TransMR programming model enables data-sharing in data-centric programming models for enhanced applicability.
- Similar to other data-centric programming models, the programmer only specifies operation on the individual data-element without concerning about its interaction with other operations.
- Prototype implementation shows that many important applications can be expressed in this model while extracting significant performance gains through increased parallelism.

#### BUT: What about the locking operations!

- Transactions are costly in a large scale distributed settings
  - two-phase locking (concurrency control)
  - two-phase commit (atomicity)
- Careful examination of the protocols and optimizations crucial to performance of TransMR-like systems
- These optimizations also useful for general purpose transactions on databases using key-value store as the underlying storage

## Lock Management in Distributed Transactions

- Lock management the major bottleneck affecting the latency of distributed transactions.
- Consider Strong Strict two phase locking (SS2PL) waiting case: The lock-acquiring stage is the only sequential stage. The other stages can be parallelized to finish in a single round-trip.
- Holds true even in optimistic-concurrency techniques where the locks are acquired at the end.



### Workload Aware Lock Management

- **Contention based Lock-ordering**: Order the locks so as to decrease the total amount of waiting time.
- For the waiting case, the lock with the least contention should be acquired first. This increases pipelining while decreasing lock-holding times.
- Contention order is a runtime characteristic, and is updated consistently. All clients should adhere to the same order to avoid deadlocks.



#### Constrained k-way Graph Partitioning

- Graph partitioning algorithm to split the locking into k non-overlapping partitions, minimizing the sum of weights on cut-edges, while approximately balancing the total weight (sum of node-weights) of individual partitions.
- The result of the partitioning algorithm is the load-balanced-partitioning of locks among k storage nodes.



- A cluster of 20 machines was used for all evaluations. Each machine had a Quad-core xeon processor with 8 GB of RAM. HBase is the underlying key-value store.
- The YCSB benchmark was extended with the atomic multi-put operation. A client transaction involves an atomic Read-Modify-Write operation on a set of keys.
- The keys for the atomic operation are generated using a Zipfian generator with variable zipfian parameter. Each transaction updates 15 keys out of 50K keys.

# Two Phase Locking: Waiting Version



Figure : Performance of Lock Partitioning

- Ordering of keys in increasing-order of their contention significantly better than the decreasing order.
- The increasing-order reduced lock holding time for highest contended locks reducing waiting time for other transactions.

# Two Phase Locking: Waiting Version



Figure : Performance of Lock Partitioning

• Partitioning is done using Metis and partitions are placed at separate nodes. Lock-partitioning improves the throughput by reducing the number of network-roundtrips needed for sequential locking by the client.

# Optimistic Concurrency Control – No-waiting Version



Figure : Performance of Lock Ordering

- Smaller improvement for OCC mainly due to the shorter duration of locking.

# Optimistic Concurrency Control – No-waiting Version



Figure : Lock wastage due to restarts

 Optimistic techniques not suitable at high contention levels as the time spent in reading and local updating gets wasted due to conflicts during commit.

Ananth Grama (Purdue University)

## Lock Optimization: Conclusions

- The waiting version of SS2PL with increasing-contention-order and partitioning outperforms the other protocols significantly.
- Restarts due to conflicts constitute a major overhead in distributed transactions. Reducing restarts by busy-waiting for locks is an important step towards increasing performance.
- Understanding the workload even simple statistics on contention is enough to achieve significant gains (up to 200%).
- Lock-partitioning through graph clustering and partitioning techniques can be performed dynamically to achieve performance gains.

# Towards Online Learning

Traditional off-line learning models:



э

(日) (同) (三) (三)

# Motivating Applications of Online Learning

- Ad-Servers (learning from click-streams, spatial parameters)
- Real-time sentiment classification (learning from twitter feeds, blogs, etc.)
- Read-time content recommendation (correlating tweets, hashtags, etc.)

```
Randomly shuffle training samples
         repeat {
              for i := 1, 2, ..., m {
                  T_{-j} := T_{-j} - a(ht(x(i) - y(i)) x(j))
                        (for every j = 0, 1, ..., n)
Can be used in applications like Supervised Semantic Indexing (e.g.,
Minimize loss function F = 1 - q^* W^+ d^+ + q^* W^+ d^-, where q is the
query document, W is a weight matrix, d^+ corresponds to a positive
document and d^{-} corresponds to a negative document.)
```

#### Stream Processing Solutions

Storm is a stream processing envine built at Twitter. It executes a DAG of operators consisting of Spouts and Bolts.



Image: Image:

# Building Online Applications: Storm and Vowpal Wabbit



э

(日) (同) (三) (三)

## Limitations of Simple Integration

Long reduction times, network bandwidth bottlenecks, static topologies and dynamic system state.



# Limitations of Simple Integration

Impedance mismatch between synchronous reductions and asynchronous execution engine (LMAX Disruptor).



# A dynamically orchestrated online learning framework

Controller dynamically schedules reductions and forces rate control and routing on input data streams



< A</li>

# Controller Dexign: Asynchronous Reductions

Schedules asynchronous reductions (staggered butterfly) while forcing rate control (blacklisting reducers)



< A</li>

# Controller Design: Dynamic Mapping of Virtual Overlay to Physical Operators

Controller uses operator metrics to decide on dynamic mapping and schedule reductions. Metrics: Model sparsity triggers new reduction; Latency and throughput of individual operators (scale-in/scale-out)



TransMR and Concurrency Control Modules available on request. Online analytics framework currently being validated. Available in limited release.

B ▶ < B ▶