# Web-Scale Data Management Systems

Prof. Oliver Kennedy

http://www.cse.buffalo.edu/~okennedy/courses/cse704fa2012.html

### <u>okennedy@buffalo.edu</u>

# About Me

- New to UB (super excited to be here)
- 4+ years working in databases
  - And 3+ years in distributed systems before that
- Worked with analytics teams at Microsoft Azure

# Office Hours

- Davis Hall 338H
  - Monday: 2:00-4:00 PM
  - Thursday: 2:00-4:00 PM
  - or by appointment

# Course Goals

- Familiarize you with systems for storing and analyzing really really big data
  - Real world systems (if you want to go into industry)
  - Interesting problems (if you want to go into academia)
- Past, present, and future systems
- Give (some of) you experience working with these systems.

# Course Requirements (Everyone)

- Submit paper overviews **Each Week** 
  - I-2 paragraphs per paper
  - What's the message of the paper?
  - What are the pros and cons of the paper's approach?
- Everyone gets 2 weeks of missed overviews

# Course Requirements (2+ Credits)

- Present 2 papers
  - 30-40 minute presentation
  - 20-30 minutes of discussion
  - Arrange to go over your presentation with me by the prior Thursday.

# Course Requirements (3+ Credits)

- A simple implementation project
  - Hadoop-based Join Algorithms
  - A Ring DHT
  - Standalone Distributed Join Algorithms
- Or develop your own project!

# Course Requirements (3+ Credits; Continued)

- Available Resources for Projects
  - A 12 core development testbed
  - \$100 of Amazon Cloud Credit
    - Be careful with this credit!
    - Avoid running out early!

# Course Requirements (3+ Credits; Continued)

- Project Requirements
  - A I-page milestone report on Nov I
  - A short (~4 page) final report
  - A short (10-15 minute) presentation
- Confirm your project with me by Oct. I

**Data-Flow Computation** Dryad Map Reduce & HDFS Map/Reduce Hive & HadoopDB SQL on M/R Pig & Dremel Other Languages on M/R MonetDB & DataCyclotron **Column Stores** Cassandra & BigTable/HBase Semistructured Databases Zookeeper & Percolator **Distributed Consistency** Chord & Dynamo **Distributed Hash Tables** PIQL **Enforced Scalability** Lipstick **Workflow Provenance Borealis & DBToaster Stream Processors** 

**Data-Flow Computation** Dryad Map Reduce & HDFS Map/Reduce Hive & HadoopDB SQL on M/R Pig & Dremel Other Languages on M/R MonetDB & DataCyclotron **Column Stores** Cassandra & BigTable/HBase Semistructured Databases Zookeeper & Percolator **Distributed Consistency Distributed Hash Tables** Chord & Dynamo **Enforced Scalability** PIQL Very recent (presented last month) Lipstick Workflow Provenance **Borealis & DBToaster Stream Processors** 

# Questions about the Course?

# Dryad

### Michael Isard, Mihai Budiu, Yuan Yu, Andrew Birrell, Dennis Fetterly Microsoft Research: Silicon Valley

(presented by Oliver Kennedy)





### Hard Problem!



### Hard Problem!

# Deadlock! IPC!

Error Handling!







### The programmer knows what's data-parallel.

The **computer** understands the **infrastructure**.

cat users.dat

cat users.dat

| sed 's/\([^ ]\*\).\*/\1/'
| grep 'buffalo'
| wc -l

# Unix pipes allow programmers to compose 'nuggets' of computation

Use the same metaphor!

Join nuggets of code into a graph

How about an Example?

# SELECT distinct p.objID FROM photoObjAll p JOIN neighbors n ON p.objID = n.objID AND n.objID < n.neighborObjID JOIN photoObjAll l ON l.objID = n.neighborObjID AND SimilarColor(l.rgb,n.rgb)</pre>















### Read pMn Emit n.neighborID : pMn



Monday, September 10, 12











# Read n.neighborID : pMn Emit sorted pMn by n.neighborId






# Graph Programming



# Graph Programming



#### • Graph Programming

- Execution Model
- Evaluation

#### < V, E, I, O >

#### < V, E, I, O >





#### < {A}, {}, {A}, {A} >

A^k



#### < {A}, {}, {A}, {A} >

A^k



< {A1,A2,...}, {}, {A1,A2,...}, {A1,A2,...} >





< {A1,A2,...}, {}, {A1,A2,...}, {A1,A2,...} >







< VA $\oplus$ VB, EAUEBUENEW, IA, OB >









< VA $\oplus$ VB, EAUEBUENEW, IA, OB >

#### $A^k >> B^k$



< VA $\oplus$ VB, EAUEBUENEW, IA, OB >



 $< V_{A} \oplus V_{B}$ ,  $E_{A} \cup E_{B}$ ,  $I_{A} \cup I_{B}$ ,  $O_{A} \cup O_{B} >$ 



#### $< V_{A} \oplus V_{B}$ , $E_{A} \cup E_{B}$ , $I_{A} \cup I_{B}$ , $O_{A} \cup O_{B} >$



Monday, September 10, 12

## Summary

- A<sup>k</sup> Parallelize A, with k replicas
- A >= B Connect A's outputs to B's inputs (one-to-one)
- A >> B Connect A's outputs to B's inputs (all-to-any)

A | | B (deduplicating nodes in both A and B)

• Graph Programming

#### • Execution Model

• Evaluation



Job Manager



#### The Job Manager coordinates graph execution

















The programmer doesn't need to worry about the channels







For In-Memory FIFOs, downstream nodes are spawned immediately.
















































What about aggregation?



Aggregation is Expensive!



Aggregation is Expensive!



We can use runtime information to optimize!



Once the Job Manager assigns nodes to machines...



... we can apply a **user-provided** function to pre-aggregate.



... we can apply a **user-provided** function to pre-aggregate.

- Graph Programming
- Execution Model
- Evaluation

SELECT distinct p.objID
FROM photoObjAll p
JOIN neighbors n
ON p.objID = n.objID
AND n.objID < n.neighborObjID
JOIN photoObjAll l
ON l.objID = n.neighborObjID
AND SimilarColor(l.rgb,n.rgb)</pre>



- P: Read/Parse Input
- **D:** Distribute Copies
- S: In-Memory Sort
- MS: Merge Sort
- C: Count



Time: ?? (really bad)



Time: II min 30 sec



# Conclusions

- The hard part of distributed programming is infrastructure.
- Programmers like thinking in code nuggets.
- Graphs are a natural representation of distributed computations.
- Dryad isn't used much, but production dataflow systems use virtually identical techniques.
  - E.g., Storm (Twitter's analytics infrastructure)

00		nathanmarz/storm - GitHub			
	GitHub, Inc. 🔒 github.com/n	athanmarz/storm			C Reader 🔘
6-∂ ∰ Bills ▼ School ▼ Lang	guages T Comics T Info P	rojects * Gaming *	Personal * TV * Tools * Macros *		_
			nathanmarz/storm - GitHub		+
	github		Signup and Pricing	Explore GitHub Features Blog Sign in	
	-				
PUBLIC	nathanmarz / storr	n		★ Star < 3,980	
	Code	Network	Pull Requests 10 Issues 78	Wiki Graphs	
	Distributed and fault-tolerant realtime computation: stream processing, continuous computation, distributed RPC, and more — Read more http://storm-project.net				
	Clone in Mac 🗘 Z	IP HTTP Git Re	ad-Only https://github.com/nathanmarz/storm.gi	Read-Only access	
	🎾 branch: master 🝷	iles Commits	Branches 33	Tags Downloads 17	
	C Latest commit to the master brand	h			
	update changelog				
	nathanmarz authored 3 days	ago		Commit 24ae774832	
	storm /				
	name	age	message	history	
	ill bin	8 days ago	Fix directory bug in install_zmq.sh [minghan]		
	iii conf	4 days ago	Implemented pluggable spout wait strategies [nathann	narz]	
	iim log4j	7 months ago	add storm dir for logs folder to use absolute path for lo	ogs folder [haitaoyao]	
	in src	4 days ago	added ISchemableSpout interface [nathanmarz]		
	in test	5 days ago	formatting [nathanmarz]		
	.gitignore	a month ago	Adding IntelliJ files to ignore [sjoerdmulder]		
	CHANGELOG.md	3 days ago	update changelog [nathanmarz]		

release commit [nathanmarz]

a year ago

LICENSE.html

# Conclusions

- The hard part of distributed programming is infrastructure.
- Programmers like thinking in code nuggets.
- Graphs are a natural representation of distributed computations.
- Dryad isn't used much, but production dataflow systems use virtually identical techniques.
  - E.g., Storm (Twitter's analytics infrastructure)

# Conclusions

- The hard part of distributed programming is infrastructure.
- Programmers like thinking in code nuggets.
- Graphs are a natural representation of distributed computations.
- Dryad isn't used much, but production dataflow systems use virtually identical techniques.
  - E.g., Storm (Twitter's analytics infrastructure)

#### Questions?