

CS535 BIG DATA

PART 2. SCALABLE FRAMEWORKS FOR REAL-TIME BIG DATA ANALYTICS
2. SERVING LAYER: CASSANDRA

Sangmi Lee Pallickara
Computer Science, Colorado State University
<http://www.cs.colostate.edu/~cs535>

11/3/2016 CS535 Big Data, Fall 2016 W11.B.1

FAQs

- Assignment 2 has been posted
 - URL for zookeeper has been updated
 - Use the same machines and ports assignment
 - Use only English message

11/3/2016 CS535 Big Data, Fall 2016 W11.B.2

Today's topics

- Document-oriented storage system - continued
 - Apache Cassandra

11/3/2016 CS535 Big Data, Fall 2016 W11.B.3

ByteOrderPartitioner

- This partitioner orders rows lexically by key bytes
- The ordered partitioner allows ordered scans by primary key
 - If your application has user names as the partition key, you can scan rows for users whose names fall between Jake and Joe
- Disadvantage of this partitioner
 - Difficult load balancing
 - Sequential writes can cause hot spots
 - Uneven load balancing for multiple tables

11/3/2016 CS535 Big Data, Fall 2016 W11.B.4

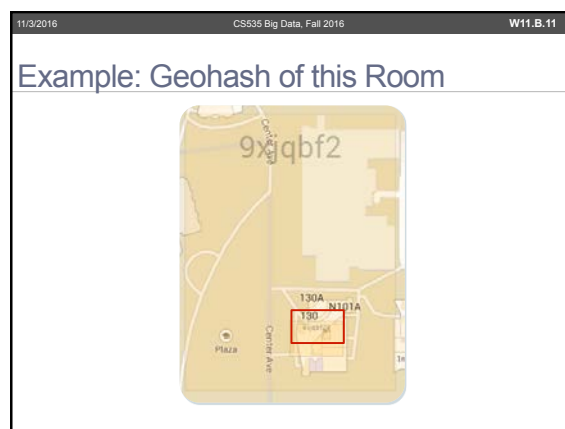
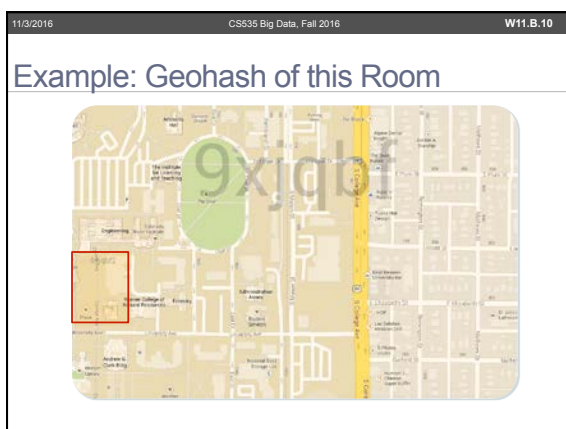
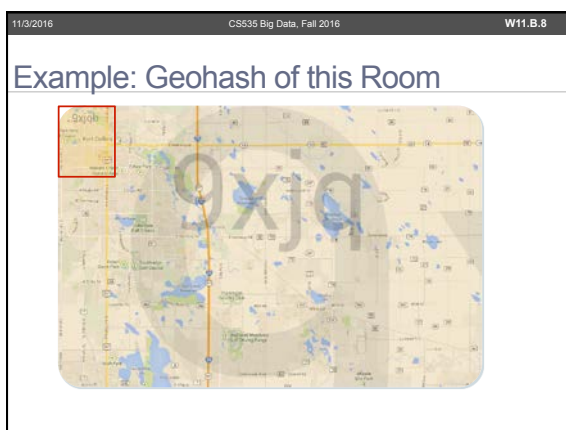
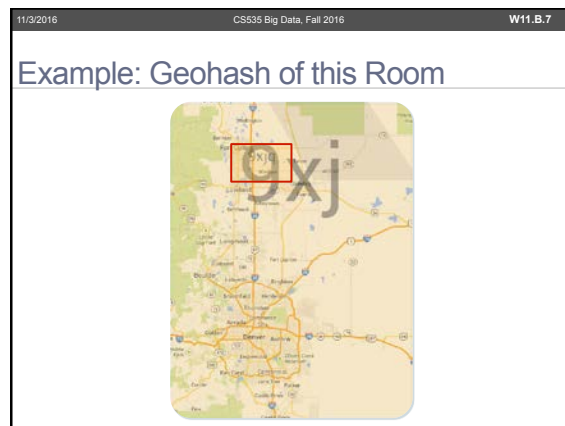
Geohashes
(2-dimensional geospatial data to DHT)

- Used in Galileo, MongoDB
 - Proximity search
- Subdivides the globe into a **hierarchy** represented by strings
- (40.573879, -105.084282) → **9XJQBDJK4XUT**
- Longer strings represent more precise coordinates
- Strings with similar prefixes are **geographically close**

11/3/2016 CS535 Big Data, Fall 2016 W11.B.5

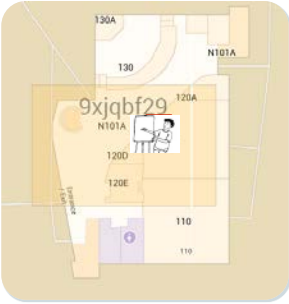
Example: Geohash of this Room





11/3/2016 CS535 Big Data, Fall 2016 W11.B.12

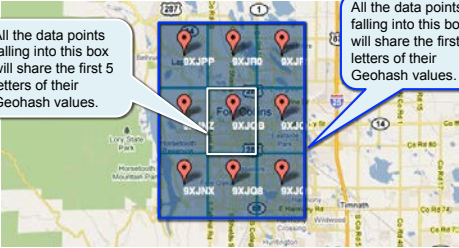
Example: Geohash of this Room



11/3/2016 CS535 Big Data, Fall 2016 W11.B.13

Adjacency with the geohash algorithm

- Geohash for Colorado State University:
(40.573879, -105.084282) → 9XJQB DJK4XUT

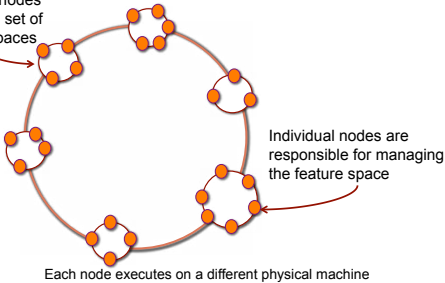


All the data points falling into this box will share the first 5 letters of their Geohash values.

All the data points falling into this box will share the first 3 letters of their Geohash values.

11/3/2016 CS535 Big Data, Fall 2016 W11.B.14

Nodes are organized in a ring of rings



A group of nodes manages a set of geohash spaces

Individual nodes are responsible for managing the feature space

Each node executes on a different physical machine

11/3/2016 CS535 Big Data, Fall 2016 W11.B.15

Apache Cassandra Replication

11/3/2016 CS535 Big Data, Fall 2016 W11.B.16

Replication

- Provides high availability and durability
- For a replication factor (replication degree) of N
 - The coordinator replicates these keys at N-1 nodes
 - Client can specify the replication scheme
 - Rack-aware/Rack-unaware/"Datacenter-aware"
- There is no master or primary replica
- Two replication strategies are available
 - SimpleStrategy
 - Use for a single data center only
 - NetworkTopologyStrategy
 - Multi-data center setup

11/3/2016 CS535 Big Data, Fall 2016 W11.B.17

SimpleStrategy

- Used only for a single data center
- Places the first replica on a node determined by the partitioner
- Places additional replicas on the next nodes clockwise in the ring without considering topology
 - Does not consider rack or data center location

11/3/2016 CS535 Big Data, Fall 2016 W11.B.18

NetworkTopologyStrategy (1/3)

- For the data cluster deployed across multiple data centers
 - This strategy specifies how many replicas you want in each data center
- Places replicas in the same data center by walking the ring clockwise until it reaches the first node in another rack
 - Attempts to place replicas on distinct racks
 - Nodes in the same rack (or similar physical grouping) often fail at the same time due to power, cooling, or network issues.

11/3/2016 CS535 Big Data, Fall 2016 W11.B.19

NetworkTopologyStrategy (2/3)

- When deciding how many replicas to configure in each data center, you should consider:
 - being able to satisfy reads locally, without incurring cross data-center latency
 - failure scenario
- The two most common ways to configure multiple data center clusters
 - Two replicas in each data center
 - This configuration tolerates the failure of a single node per replication group and still allows local reads at a consistency level of ONE.
 - Three replicas in each data center
 - This configuration tolerates either the failure of one node per replication group at a strong consistency level of LOCAL_QUORUM or multiple node failures per data center using consistency level ONE.

11/3/2016 CS535 Big Data, Fall 2016 W11.B.20

NetworkTopologyStrategy (3/3)

- Asymmetrical replication groupings
 - For example, you can maintain 4 replicas
 - Three replicas in one data center to serve real-time application requests
 - A single replica elsewhere for running analytics.

11/3/2016 CS535 Big Data, Fall 2016 W11.B.21

Apache Cassandra Virtual Node

11/3/2016 CS535 Big Data, Fall 2016 W11.B.22

What are Vnodes?

- With consistent hashing, a node owns exactly one contiguous range in the ring-space
- Vnodes change from one token or range per node, to many per node
 - Within a cluster these can be randomly selected and be non-contiguous, giving us many smaller ranges that belong to each node

11/3/2016 CS535 Big Data, Fall 2016 W11.B.23

The diagram illustrates the concept of Virtual Nodes (Vnodes) in Apache Cassandra. It shows two scenarios for a ring with 6 nodes (Node 1 to Node 6):

- Ring without VNodes:** Each node owns a single contiguous range of tokens. For example, Node 1 owns tokens A, B, C, D, E, F.
- Ring with VNodes:** Each node owns multiple non-contiguous ranges of tokens. For example, Node 1 owns tokens A, B, C, D, E, F, G, H, I, J, K, L, M, N, O, P, Q, R, S, T, U, V, W, X, Y, Z.

The diagram shows that with Vnodes, the number of tokens per node increases significantly, allowing for more granular distribution of data across the cluster.

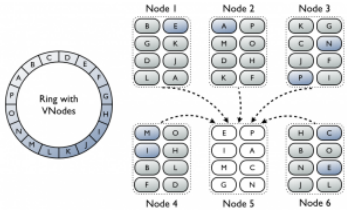
<http://www.datastax.com/dev/blog/virtual-nodes-in-cassandra-1-2>

11/3/2016 CS535 Big Data, Fall 2016 W11.B.24

Advantages of Vnodes

- **Example**
 - 30 nodes and replication factor of 3
 - A node dies completely, and we need to bring up a replacement
 - A replica for 3 different ranges to reconstitute
 - 1 set of the first natural replica
 - 2 sets of replica for replication factor of 3
 - Since our RF is 3 and we lost a node, we logically only have 2 replicas left, which for 3 ranges means there are up to 6 nodes we can stream from
 - With the setup of RF3, data will be streamed from 3 other nodes total

11/3/2016 CS535 Big Data, Fall 2016 W11.B.25



- If vnodes are spread throughout the entire cluster
 - Data transfers will be distributed on more machines

11/3/2016 CS535 Big Data, Fall 2016 W11.B.26

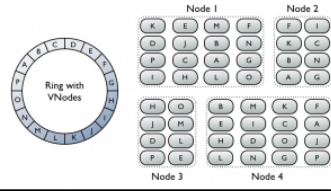
Restoring a new disk with vnodes

- **Process of restoring a disk**
 - Validating all the data and generating a Merkle tree
 - This might take an hour
 - Streaming when the actual data that is needed is sent
 - This phase takes a few minutes
- **Advantage of using Vnodes**
 - Since the ranges are smaller, data will be sent to the damaged node in a more incremental fashion
 - Instead of waiting until the end of a large validation phase
 - The validation phase will be parallelized across more machines, causing it to complete faster

11/3/2016 CS535 Big Data, Fall 2016 W11.B.27

The use of heterogeneous machines with vnodes

- **Newer nodes might be able to bear more load immediately**
 - You just assign a proportional number of vnodes to the machines with more capacity
 - e.g. If you started your older machines with 64 vnodes per node and the new machines are twice as powerful, give them 128 vnodes each and the cluster remains balanced even during transition



11/3/2016 CS535 Big Data, Fall 2016 W11.B.28

Apache Cassandra

Gossip (Internode communications)

11/3/2016 CS535 Big Data, Fall 2016 W11.B.29

Use of Gossip in Cassandra

- **Peer-to-peer communication protocol**
 - Periodically exchange state information about nodes themselves and about other nodes they know about
- Every node talks to up to three other nodes in the cluster
- A gossip message has a version associated with it
 - During a gossip exchange, older information is overwritten with the most current state for a particular node

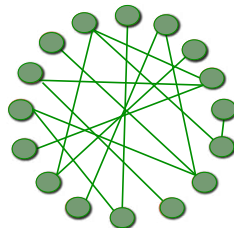
11/3/2016 CS535 Big Data, Fall 2016 W11.B.30

What is gossip?

- Broadcast protocol for disseminating data
- Decentralized, peer-to-peer networks
- 'epidemic'
- Fault tolerant
- Epidemic broadcast protocol provides a resilient and efficient mechanism for data dissemination
- Cassandra uses gossip for peer discovery and metadata propagation

11/3/2016 CS535 Big Data, Fall 2016 W11.B.31

Broadcast protocol for disseminating data



11/3/2016 CS535 Big Data, Fall 2016 W11.B.32

Why gossip for Cassandra?

- Reliably disseminate node metadata to peers
 - Cluster membership
 - Heartbeat
 - Node status
- Each node maintains a view of all peers

11/3/2016 CS535 Big Data, Fall 2016 W11.B.33

What gossip is not for in Cassandra?

- Streaming
- Repair
- Reads/write
- Compaction
- Hint
- CQL query parsing/execution

11/3/2016 CS535 Big Data, Fall 2016 W11.B.34

Data structure

- HeartBeatState
- ApplicationState
- EndpointState
 - Wrapper of a heartbeat state and a set of application state

11/3/2016 CS535 Big Data, Fall 2016 W11.B.35

HeartBeatState

- Generation
- Heartbeat
 - Periodically update monotonically increasing value

11/3/2016 CS535 Big Data, Fall 2016 W11.B.36

Application state

- {enum_name, value, version}
- Contained as a map in `EndpointState` per peer

11/3/2016 CS535 Big Data, Fall 2016 W11.B.37

ApplicationState enum

- DC/RACK
 - Where you are
- SCHEMA
- LOAD
 - Updated every 60 seconds
- SEVERITY
 - I/O load
- STATUS

11/3/2016 CS535 Big Data, Fall 2016 W11.B.38

Status (AppState)

- Bootstrap
 - For new nodes
- Hibernate
- Normal
- Leaving/Left
- Removing/Removed

11/3/2016 CS535 Big Data, Fall 2016 W11.B.39

Gossip messaging

- Every second, each node starts a new round
- Peer selection (1-3 peers)
 - Live peer
 - Seed (maybe)
 - Unreachable peer(maybe)

11/3/2016 CS535 Big Data, Fall 2016 W11.B.40

Gossip Exchange

- SYN/ACK/ACK2
 - Similar to TCP 3-way handshake
 - Add anti-entropy to gossiping

11/3/2016 CS535 Big Data, Fall 2016 W11.B.41

SYN: GossipDigestSynMessage

- Initiator sends a digest of all the nodes it knows about to a peer
- {ipAddr, generation, heartbeat}

11/3/2016 CS535 Big Data, Fall 2016 W11.B.42

GossipDigestSynMessage

```
# Suppose we are in 10.0.0.1
EndPointState 10.0.0.1
HeartBeatState: generation 1259909635, version 325
ApplicationState "load-information": 5.2, generation 1259909635, version 45
ApplicationState "bootstrapping": bXlpasf3XD8KyKs, generation 1259909635, version 56
ApplicationState "normal": bXlpasf3XD8KyKs, generation 1259909635, version 87
EndPointState 10.0.0.2
HeartBeatState: generation 1259911052, version 61
ApplicationState "load-information": 2.7, generation 1259911052, version 2
ApplicationState "bootstrapping": AuJDMftpyUvebnn, generation 1259911052, version 31
EndPointState 10.0.0.3
HeartBeatState: generation 1259912238, version 5
ApplicationState "load-information": 12.0, generation 1259912238, version 3
EndPointState 10.0.0.4
HeartBeatState: generation 1259912942, version 18
ApplicationState "load-information": 6.7, generation 1259912942, version 3
ApplicationState "normal": bJ05IV0lVwKw2xH, generation 1259912942, version 7
```

Max version number for this endpoint: 325, 61, 5, and 18

GossipDigestSynMessage

```
10.0.0.1:1259909635:325
10.0.0.2:1259911052:61
10.0.0.3:1259912238:5
10.0.0.4:1259912942:18
```

Source: <https://wiki.apache.org/cassandra/ArchitectureGossip>

11/3/2016 CS535 Big Data, Fall 2016 W11.B.43

ACK: GossipDigestActMessage

- Peer receives GossipDigestSynMessage
- Sort gossip digest list according to the difference in max version number between sender's digest and own information in descending order
 - Handle those digests first that differ mostly in version number
- Produces a diff and sends back an ACK
- Diff contains
 - Map of APPStates (for any node) that the peer has which the initiator does not
 - Digest of nodes (and their corresponding metadata) which a peer needs from an initiator

11/3/2016 CS535 Big Data, Fall 2016 W11.B.44

GossipDigestActMessage

```
# Suppose we are in 10.0.0.2
# The EndPointState in 10.0.0.2 will be
EndPointState 10.0.0.1
HeartBeatState: generation 1259909635, version 324
ApplicationState "load-information": 5.2, generation 1259909635, version 45
ApplicationState "bootstrapping": bXlpasf3XD8KyKs, generation 1259909635, version 56
ApplicationState "normal": bXlpasf3XD8KyKs, generation 1259909635, version 87
EndPointState 10.0.0.2
HeartBeatState: generation 1259911052, version 63
ApplicationState "load-information": 2.7, generation 1259911052, version 2
ApplicationState "bootstrapping": AuJDMftpyUvebnn, generation 1259911052, version 31
ApplicationState "normal": AuJDMftpyUvebnn, generation 1259911052, version 62
EndPointState 10.0.0.3
HeartBeatState: generation 1259812143, version 2142
ApplicationState "load-information": 16.0, generation 1259812143, version 1803
ApplicationState "normal": W20IXVUC3hggpC7, generation 1259812143, version 6
```

```
# GossipDigestSynMessage
From 10.0.0.1
10.0.0.1:1259909635:325
10.0.0.2:1259911052:61
10.0.0.3:1259912238:5
10.0.0.4:1259912942:18

# The GossipDigestActMessage from 10.0.0.2 is
10.0.0.1:1259909635:324
10.0.0.3:1259912238:0
10.0.0.4:1259912942:0
10.0.0.2:ApplicationState "normal": AuJDMftpyUvebnn, generation 1259911052, version 62, [HeartBeatState, generation 1259911052, version 63]
```

11/3/2016 CS535 Big Data, Fall 2016 W11.B.45

ACK2: GossipDigestAct2Message

- Initiator receives ACK
- Applies any AppState and sends back an ACK2
- ACK2 has a map of AppStates which the peer does not have

```
# The GossipDigestAct2Message from 10.0.0.1 is
10.0.0.1:[ApplicationState "load-information": 5.2, generation 1259909635, version 45],
[ApplicationState "bootstrapping": bXlpasf3XD8KyKs, generation 1259909635, version 56],
[ApplicationState "normal": bXlpasf3XD8KyKs, generation 1259909635, version 87],
[HeartBeatState, generation 1259909635, version 325]
10.0.0.3:[ApplicationState "load-information": 12.0, generation 1259912238, version 3],
[HeartBeatState, generation 1259912238, version 3]
10.0.0.4:
[ApplicationState "load-information": 6.7, generation 1259912942, version 3],
[ApplicationState "normal": bJ05IV0lVwKw2xH, generation 1259912942, version 7],
[HeartBeatState: generation 1259912942, version 18]
```

11/3/2016 CS535 Big Data, Fall 2016 W11.B.46

AppState Reconciliation

- Generation
- Heartbeat
- AppState based on comparing version

11/3/2016 CS535 Big Data, Fall 2016 W11.B.47

Reconciliation example

	A	
A	gen:1234 Hb: 994 Status: normal {4}	gen:1234 Hb: 990 Status: normal {4}
B	Gen:2345 Hb: 10 Status: bootstrap {1}	Gen:2345 Hb:17 Status: normal {2}
C	Gen:5555 Hb: 1111 Status: normal {5}	
D	Gen:2222 Hb: 4444 status: normal {3}	Gen:3333 Hb: 11 Status: normal {3}

11/3/2016 CS535 Big Data, Fall 2016 W11.B.48

Messaging summary

- Each node starts a gossip round every second
- 1-3 peers per round
- 3 messages passed
- Constant amount of network traffic

11/3/2016 CS535 Big Data, Fall 2016 W11.B.49

Practical implications

- Who is in the cluster?
- How are peers judged UP or DOWN?
- When does a node stop sending a peer traffic?
- When is one peer preferred over another?
- When does a node leave the cluster?

11/3/2016 CS535 Big Data, Fall 2016 W11.B.50

Cluster membership

- Gossip with a seed upon startup
- Learn about all peers
- Gossip
- Lather, rinse, repeat

11/3/2016 CS535 Big Data, Fall 2016 W11.B.51

UP/DOWN?

- Local to each node
 - Not shared via gossip
- Determined via heartbeat

11/3/2016 CS535 Big Data, Fall 2016 W11.B.52

Failure Detection

- Glorified heartbeat listener
- Records timestamp when heartbeat update is received for each peer
- Keeps backlog of timestamp intervals between updates
- Periodically checks all peers to make sure that we've heard from them recently

11/3/2016 CS535 Big Data, Fall 2016 W11.B.53

UP/DOWN affects

- Stop sending writes (hints)
- Sending reads
- Gossip
 - It is down
 - This node is treated as an unavailable node
- Repair/stream sessions are terminated

11/3/2016 CS535 Big Data, Fall 2016 W11.B.54

What if a peer is really slow?

- Peer is NOT marked down
 - We will try to avoid it

11/3/2016 CS535 Big Data, Fall 2016 W11.B.55

Dynamic “Snitch”

- Determine when to avoid a slow node
- Scoring peers based on response times
 - Scores recalculated every 100ms (default)
 - Scores reset every 10m (default)

11/3/2016 CS535 Big Data, Fall 2016 W11.B.56

How do nodes leave?

- STATUS = LEAVING
- Stream data
- Stream hints
- STATUS = LEFT, expiryTime

11/3/2016 CS535 Big Data, Fall 2016 W11.B.57

Decomission

- STATUS = LEAVING
- Stream data
- Stream hints
- STATUS = LEFT, expiryTime

11/3/2016 CS535 Big Data, Fall 2016 W11.B.58

Remove node

- STATUS = REMOVING
- Rebalance cluster
 - Notify coordinator
- Delete hint
- STATUS = REMOVED, expiryTime

11/3/2016 CS535 Big Data, Fall 2016 W11.B.59

Replace node

- Cassandra.replace_address
- “shadow gossip”
- Take tokens/hostID(hints)
 - Check that previous owner hasn't gossiped
- Stream data

11/3/2016 CS535 Big Data, Fall 2016 W11.B.60

"Assassinate!"

- Managing hanging non-functional nodes
- `unsafeAssassinateEndpoint(ipAddr)`
 - Use with caution
- Forces change to peer

11/3/2016 CS535 Big Data, Fall 2016 W11.B.61

Failure detection: Φ Accrual Failure Detector

11/3/2016 CS535 Big Data, Fall 2016 W11.B.62

Failure detection (1/3)

- Φ Accrual Failure Detector
- Accrual Failure Detection does not emit a Boolean value stating a node is up or down
 - Emits a value which represents a suspicion level for each of the monitored nodes
 - This value is defined as Φ
- Dynamically adjusts to reflect network and load conditions at the monitored nodes

11/3/2016 CS535 Big Data, Fall 2016 W11.B.63

Failure detection (2/3)

- Given some threshold Φ , and assuming that we decide to suspect a node A
 - e.g. when $\Phi = 1$, then the likelihood that we will make a mistake is about 10%.
 - The likelihood is about 1% with $\Phi = 2$
 - The likelihood is about 0.1% with $\Phi = 3$

11/3/2016 CS535 Big Data, Fall 2016 W11.B.64

Failure detection (3/3)

- Every node maintains a sliding window of inter-arrival times of gossip messages from other nodes in the cluster
- Exponential Distribution
 - The nature of gossip channel and its impact on latency

11/3/2016 CS535 Big Data, Fall 2016 W11.B.65

Bootstrapping and persistence

11/3/2016 CS535 Big Data, Fall 2016 W11.B.66

Bootstrapping

- When a node joins the ID ring, the mapping is persisted to the disk locally and in Zookeeper
 - Then the token information is gossiped around the cluster
- With bootstrapping, a node joins with a configuration file that contains a list of a few contact points
 - Seeds of the cluster
- Seeds can be provided by a configuration service (e.g. Zookeeper)

11/3/2016 CS535 Big Data, Fall 2016 W11.B.67

Local persistence: Write Operation

- Write into a commit log
 - Durability and recoverability
 - Dedicated disk for each node
- Write into an in-memory data structure
 - When in-memory data structure crosses a certain threshold, it dumps itself to disk
- Write into disk
 - Generates an index for efficient lookup based on row key
- Similar to Bigtable (compaction)

11/3/2016 CS535 Big Data, Fall 2016 W11.B.68

Local persistence: Read Operation

- First queries the in-memory data structure
- Disk lookup
 - Look-up a key
- To narrow down the lookup process
 - a bloom filter is stored in each data file and memory

11/3/2016 CS535 Big Data, Fall 2016 W11.B.69

Apache Cassandra Replication

11/3/2016 CS535 Big Data, Fall 2016 W11.B.70

Replication

- Provides high availability and durability
- For a replication factor (replication degree) of N
 - The coordinator replicates these keys at N-1 nodes
 - Client can specify the replication scheme
 - Rack-aware/Rack-unaware/"Datacenter-aware"
- There is no master or primary replica
- Two replication strategies are available
 - SimpleStrategy
 - Use for a single data center only
 - NetworkTopologyStrategy
 - Multi-data center setup

11/3/2016 CS535 Big Data, Fall 2016 W11.B.71

SimpleStrategy

- Used only for a single data center
- Places the first replica on a node determined by the **partitioner**
- Places additional replicas on the next nodes clockwise in the ring without considering topology
 - Does not consider rack or data center location

11/3/2016 CS535 Big Data, Fall 2016 W11.B.72

NetworkTopologyStrategy (1/3)

- For the data cluster deployed across multiple data centers
 - This strategy specifies how many replicas you want in each data center
- Places replicas in the same data center by walking the ring clockwise until it reaches the first node in another rack
 - Attempts to place replicas on distinct racks
 - Nodes in the same rack (or similar physical grouping) often fail at the same time due to power, cooling, or network issues.

11/3/2016 CS535 Big Data, Fall 2016 W11.B.73

NetworkTopologyStrategy (2/3)

- When deciding how many replicas to configure in each data center, you should consider:
 - being able to satisfy reads locally, without incurring cross data-center latency
 - failure scenario
- The two most common ways to configure multiple data center clusters
 - Two replicas in each data center
 - This configuration tolerates the failure of a single node per replication group and still allows local reads at a consistency level of ONE.
 - Three replicas in each data center
 - This configuration tolerates either the failure of one node per replication group at a strong consistency level of LOCAL_QUORUM or multiple node failures per data center using consistency level ONE.

11/3/2016 CS535 Big Data, Fall 2016 W11.B.74

NetworkTopologyStrategy (3/3)

- Asymmetrical replication groupings
 - For example, you can maintain 4 replicas
 - Three replicas in one data center to serve real-time application requests
 - A single replica elsewhere for running analytics.

11/3/2016 CS535 Big Data, Fall 2016 W11.B.75

Apache Cassandra Virtual Node

11/3/2016 CS535 Big Data, Fall 2016 W11.B.76

What are Vnodes?

- With consistent hashing, a node owns exactly one contiguous range in the ring-space
- Vnodes change from one token or range per node, to many per node
 - Within a cluster these can be randomly selected and be non-contiguous, giving us many smaller ranges that belong to each node

11/3/2016 CS535 Big Data, Fall 2016 W11.B.77

The diagram illustrates the transition from a ring without Vnodes to a ring with Vnodes. The top part shows a ring with 6 nodes (Node 1 to Node 6) and their respective token ranges. The bottom part shows the same ring with Vnodes, where each node has multiple non-contiguous token ranges.

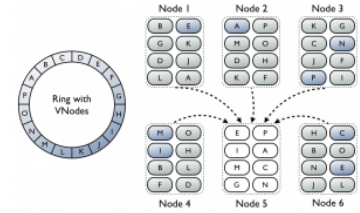
<http://www.datastax.com/dev/blog/virtual-nodes-in-cassandra-1-2>

11/3/2016 CS535 Big Data, Fall 2016 W11.B.76

Advantages of Vnodes

- Example
 - 30 nodes and replication factor of 3
 - A node dies completely, and we need to bring up a replacement
 - A replica for 3 different ranges to reconstitute
 - 1 set of the first natural replica
 - 2 sets of replica for replication factor of 3
 - Since our RF is 3 and we lost a node, we logically only have 2 replicas left, which for 3 ranges means there are up to 6 nodes we can stream from
 - With the setup of RF3, data will be streamed from 3 other nodes total

11/3/2016 CS535 Big Data, Fall 2016 W11.B.79



- If vnodes are spread throughout the entire cluster
 - Data transfers will be distributed on more machines

11/3/2016 CS535 Big Data, Fall 2016 W11.B.80

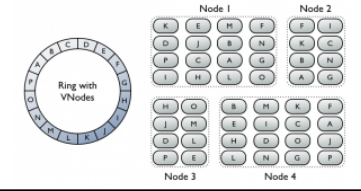
Restoring a new disk with vnodes

- Process of restoring a disk
 - Validating all the data and generating a Merkle tree
 - This might take an hour
 - Streaming when the actual data that is needed is sent
 - This phase takes a few minutes
- Advantage of using Vnodes
 - Since the ranges are smaller, data will be sent to the damaged node in a more incremental fashion
 - Instead of waiting until the end of a large validation phase
 - The validation phase will be parallelized across more machines, causing it to complete faster

11/3/2016 CS535 Big Data, Fall 2016 W11.B.81

The use of heterogeneous machines with vnodes

- Newer nodes might be able to bear more load immediately
 - You just assign a proportional number of vnodes to the machines with more capacity
 - e.g. If you started your older machines with 64 vnodes per node and the new machines are twice as powerful, give them 128 vnodes each and the cluster remains balanced even during transition



11/3/2016 CS535 Big Data, Fall 2016 W11.B.82

Apache Cassandra

Gossip (Internode communications)

11/3/2016 CS535 Big Data, Fall 2016 W11.B.83

Use of Gossip in Cassandra

- Peer-to-peer communication protocol
 - Periodically exchange state information about nodes themselves and about other nodes they know about
- Every node talks to up to three other nodes in the cluster
- A gossip message has a version associated with it
 - During a gossip exchange, older information is overwritten with the most current state for a particular node

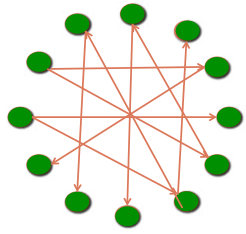
11/3/2016 CS535 Big Data, Fall 2016 W11.B.84

What is gossip?

- Broadcast protocol for disseminating data
- Decentralized, peer-to-peer networks
- 'epidemic'
- Fault tolerant
- Epidemic broadcast protocol provides a resilient and efficient mechanism for data dissemination
- Cassandra uses gossip for peer discovery and metadata propagation

11/3/2016 CS535 Big Data, Fall 2016 W11.B.85

Broadcast protocol for disseminating data



11/3/2016 CS535 Big Data, Fall 2016 W11.B.86

Why gossip for Cassandra?

- Reliably disseminate node metadata to peers
 - Cluster membership
 - Heatbeat
 - Node status
 - **Each node maintains a view of ALL peers**

11/3/2016 CS535 Big Data, Fall 2016 W11.B.87

What gossip is not for in Cassandra?

- Streaming
- Repair
- Reads/write
- Compaction
- Hint
- CQL query parsing/execution

11/3/2016 CS535 Big Data, Fall 2016 W11.B.88

Data structure

- HeartBeatState
- ApplicationState
- EndpointState
 - Wrapper of a heartbeat state and a set of application state

11/3/2016 CS535 Big Data, Fall 2016 W11.B.89

HeartBeatState

- Generation
- Heartbeat
 - Periodically update monotonically increasing value

11/3/2016 CS535 Big Data, Fall 2016 W11.B.90

Application state

- {enum_name, value, version}
- Contained as a map in `EndpointState` per peer

11/3/2016 CS535 Big Data, Fall 2016 W11.B.91

ApplicationState enum

- DC/RACK
 - Where you are
- SCHEMA
- LOAD
 - Updated every 60 seconds
- SEVERITY
 - I/O load
- STATUS

11/3/2016 CS535 Big Data, Fall 2016 W11.B.92

Status (AppState)

- Bootstrap
 - For new nodes
- Hibernate
- Normal
- Leaving/Left
- Removing/Removed

11/3/2016 CS535 Big Data, Fall 2016 W11.B.93

Gossip messaging

- Every second, each node starts a new round
- Peer selection (1-3 peers)
 - Live peer
 - Seed (maybe)
 - Unreachable peer(maybe)

11/3/2016 CS535 Big Data, Fall 2016 W11.B.94

Gossip Exchange

- SYN/ACK/ACK2
 - Similar to TCP 3-way handshake
 - Add anti-entropy to gossiping

11/3/2016 CS535 Big Data, Fall 2016 W11.B.95

SYN: GossipDigestSynMessage

- Initiator sends a **digest** of all the nodes it knows about to a peer
- {ipAddr, generation, heartbeat}

11/3/2016 CS535 Big Data, Fall 2016 W11.B.96

GossipDigestSynMessage

Suppose we are in 10.0.0.1

```
EndPointState 10.0.0.1
HeartBeatState: generation 1259909635, version 325
ApplicationState "load-information": 5.2, generation 1259909635, version 45
ApplicationState "bootstrapping": bXlpasF3XD8KyKs, generation 1259909635, version 56
ApplicationState "normal": bXlpasF3XD8KyKs, generation 1259909635, version 87
EndPointState 10.0.0.2
HeartBeatState: generation 1259911052, version 61
ApplicationState "load-information": 2.7, generation 1259911052, version 2
ApplicationState "bootstrapping": AuJDMtppUvEbnn, generation 1259911052, version 31
EndPointState 10.0.0.3
HeartBeatState: generation 1259912238, version 5
ApplicationState "load-information": 12.0, generation 1259912238, version 3
EndPointState 10.0.0.4
HeartBeatState: generation 1259912942, version 18
ApplicationState "load-information": 6.7, generation 1259912942, version 3
ApplicationState "normal": bJ05IVc0lvRWx2Xh, generation 1259912942, version 7
```

Max version number for this endpoint: 325, 61, 5, and 18

GossipDigestSynMessage

```
10.0.0.1:1259909635:325
10.0.0.2:1259911052:61
10.0.0.3:1259912238:5
10.0.0.4:1259912942:18
```

Source: <https://wiki.apache.org/cassandra/ArchitectureGossip>

11/3/2016 CS535 Big Data, Fall 2016 W11.B.97

ACK: GossipDigestActMessage

- Peer receives GossipDigestSynMessage
- Sort gossip digest list according to the difference in max version number between sender's digest and own information in descending order
 - Handle those digests first that differ mostly in version number
- Produces a diff and sends back an ACK
- Diff contains
 - Map of APPStates (for any node) that the peer has which the initiator does not
 - Digest of nodes (and their corresponding metadata) which a peer needs from an initiator

11/3/2016 CS535 Big Data, Fall 2016 W11.B.98

GossipDigestActMessage

Suppose we are in 10.0.0.2

The EndPointState in 10.0.0.2 will be

```
EndPointState 10.0.0.1
HeartBeatState: generation 1259909635, version 324
ApplicationState "load-information": 5.2, generation 1259909635, version 45
ApplicationState "bootstrapping": bXlpasF3XD8KyKs, generation 1259909635, version 56
ApplicationState "normal": bXlpasF3XD8KyKs, generation 1259909635, version 87
EndPointState 10.0.0.2
HeartBeatState: generation 1259911052, version 63
ApplicationState "load-information": 2.7, generation 1259911052, version 2
ApplicationState "bootstrapping": AuJDMtppUvEbnn, generation 1259911052, version 31
ApplicationState "normal": AuJDMtppUvEbnn, generation 1259911052, version 62
EndPointState 10.0.0.3
HeartBeatState: generation 1259812143, version 2142
ApplicationState "load-information": 16.0, generation 1259812143, version 1803
ApplicationState "normal": W20IXVCUwMgpc77, generation 1259812143, version 6
```

Compares

Generates ACT

```
# GossipDigestSynMessage From 10.0.0.1
10.0.0.1:1259909635:325
10.0.0.2:1259911052:61
10.0.0.3:1259912238:5
10.0.0.4:1259912942:18

# The GossipDigestActMessage from 10.0.0.2 is
10.0.0.1:1259909635:324
10.0.0.3:1259912238:0
10.0.0.4:1259912942:0
10.0.0.2:ApplicationState "normal": AuJDMtppUvEbnn, generation 1259911052, version 62, [HeartBeatState, generation 1259911052, version 63]
```

11/3/2016 CS535 Big Data, Fall 2016 W11.B.99

ACK2: GossipDigestAct2Message

- Initiator receives ACK
- Applies any AppState and sends back an ACK2
- ACK2 has a map of AppStates which the peer does not have

```
# The GossipDigestAct2Message from 10.0.0.1 is
10.0.0.1:[ApplicationState "load-information": 5.2, generation 1259909635, version 45],
[ApplicationState "bootstrapping": bXlpasF3XD8KyKs, generation 1259909635, version 56],
[ApplicationState "normal": bXlpasF3XD8KyKs, generation 1259909635, version 87],
[HeartBeatState, generation 1259909635, version 325]
10.0.0.3:[ApplicationState "load-information": 12.0, generation 1259912238, version 3],
[HeartBeatState, generation 1259912238, version 3]
10.0.0.4:
[ApplicationState "load-information": 6.7, generation 1259912942, version 3],
[ApplicationState "normal": bJ05IVc0lvRWx2Xh, generation 1259912942, version 7],
[HeartBeatState: generation 1259912942, version 18]
```

11/3/2016 CS535 Big Data, Fall 2016 W11.B.100

AppState Reconciliation

- Generation
- Heartbeat
- AppState based on comparing version

11/3/2016 CS535 Big Data, Fall 2016 W11.B.101

Reconciliation example

	A	B	After reconciliation
A	gen:1234 Hb: 994 Status: normal {4}	gen:1234 Hb: 990 Status: normal {4}	Heartbeats are different. B will update its Hb to A's status
B	Gen:2345 Hb: 10 Status: bootstrap {1}	Gen:2345 Hb:17 Status: normal {2}	The version numbers are different A will update the state to B's status
C	Gen:5555 Hb: 1111 Status: normal {5}		B will know that C exists in this cluster
D	Gen:2222 Hb: 4444 status: normal {3}	Gen:3333 Hb: 11 Status: normal {3}	Node D has been re-started A will take all of the metadata in B

4 nodes cluster A, B, C and D

11/3/2016 CS535 Big Data, Fall 2016 W11.B.102

Messaging summary

- Each node starts a gossip round every second
- 1-3 peers per round
- 3 messages passed
- Constant amount of network traffic
 - Broadcast same number of gossip messages

CS535 BIG DATA

Sangmi Lee Pallickara
Computer Science, Colorado State University
<http://www.cs.colostate.edu/~cs535>

11/3/2016 CS535 Big Data, Fall 2016 W11.B.104

FAQs

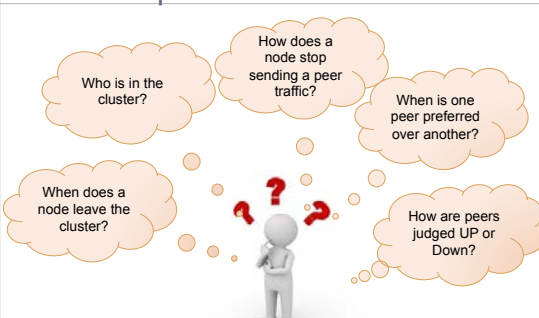
11/3/2016 CS535 Big Data, Fall 2016 W11.B.105

Today's topics

- Document-oriented storage system - continued
 - Apache Cassandra
- Graph storage and processing models

11/3/2016 CS535 Big Data, Fall 2016 W11.B.106

Practical implications



- Who is in the cluster?
- When does a node leave the cluster?
- How does a node stop sending a peer traffic?
- When is one peer preferred over another?
- How are peers judged UP or Down?

11/3/2016 CS535 Big Data, Fall 2016 W11.B.107

Cluster membership

- Gossip with a seed upon startup
- Learn about all peers
- Gossip
- Lather, rinse, repeat

11/3/2016 CS535 Big Data, Fall 2016 W11.B.108

UP/DOWN?

- Local to each node
 - Not shared via gossip
- Determined via heartbeat

11/3/2016 CS535 Big Data, Fall 2016 W11.B.109

Failure Detector

- Glorified heartbeat listener
- Records timestamp when heartbeat update is received for each peer
- Keeps backlog of timestamp intervals between updates
- Periodically checks all peers to make sure that we've heard from them recently

11/3/2016 CS535 Big Data, Fall 2016 W11.B.110

UP/DOWN affects

- Stop sending writes (hints)
- Sending reads
- Gossip
 - It is down
 - This node is treated as an unavailable node
- Repair/stream sessions are terminated
 - Terminate sockets

11/3/2016 CS535 Big Data, Fall 2016 W11.B.111

What if a peer is really slow?

- Peer is NOT marked down
 - We will try to avoid it

11/3/2016 CS535 Big Data, Fall 2016 W11.B.112

Dynamic "Snitch"

- Determine when to avoid a slow node
- Scoring peers based on response times
 - Scores recalculated every 100ms (default)
 - The updates are capped at a maximum of 10,000 per scoring interval
 - Scores reset every 10 minutes (default)
 - Uses statistically significant random sampling

<http://www.datastax.com/dev/blog/dynamic-snitching-in-cassandra-past-present-and-future>

11/3/2016 CS535 Big Data, Fall 2016 W11.B.113

How do nodes leave?

11/3/2016 CS535 Big Data, Fall 2016 W11.B.114

Decommission

- STATUS = LEAVING
- Stream data
- Stream hints
- STATUS = LEFT, expiryTime (hard coded as 3 days)

11/3/2016 CS535 Big Data, Fall 2016 W11.B.115

Remove node

- STATUS = REMOVING
- Rebalance cluster
 - Notify coordinator
- Delete hint
- STATUS = REMOVED, expiryTime (hard coded as 3 days)

11/3/2016 CS535 Big Data, Fall 2016 W11.B.116

Replace node

- Cassandra.replace_address
- "shadow gossip"
- Take tokens/hostID(hints)
 - Check that previous owner hasn't gossiped
- Stream data

11/3/2016 CS535 Big Data, Fall 2016 W11.B.117

"Assassinate!"

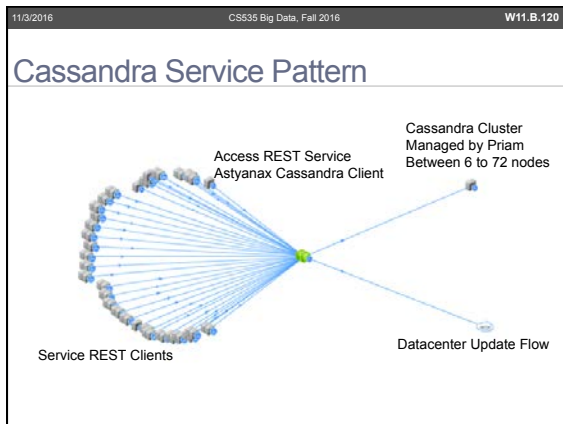
- Managing hanging non-functional nodes
- unsafeAssassinateEndpoint(ipAddr)
 - Use with caution
- Forces change to peer

11/3/2016 CS535 Big Data, Fall 2016 W11.B.118

Running Netflix on Cassandra in the Cloud

11/3/2016 CS535 Big Data, Fall 2016 W11.B.119

2009	2009	2010	2010	2010	2011
Content	Logs	Play	WWW	API	Customer Service
Content management	S3 terabytes	DRM	Sign-up	Metadata	International CS lookup
EC2 Encoding	EMR	CDN routing	Search	Device config	Diagnostics and Actions
S3 Petabytes	Hive and Pig	Bookmarks	Movie Choosing	TV Movie Choosing	Customer Call log
Business intelligence	Logging	Ratings	Social Facebook	CS analytics	
CDN's ISPs Terabits Customers	Hadoop	Cassandra	MySQL	Cassandra	Cassandra



11/3/2016 CS535 Big Data, Fall 2016 W11.B.121

Production deployment

- Over 50 Cassandra Clusters
- Over 500 nodes
- Over 30 TB of daily backups
- Biggest cluster 72 nodes
- 1 cluster over 250 K Writes/s

11/3/2016 CS535 Big Data, Fall 2016 W11.B.122

High Availability

- Cassandra stores 3 local copies, 1 per zone
 - Synchronous access, durable, highly available
 - Read/Write One fastest, use for fire and forget
 - Read/Write Quorum 2 of 3, user for read-after-write
- AWS Availability Zones
 - Separate buildings
 - Separate power etc.
 - Fairly close together

11/3/2016 CS535 Big Data, Fall 2016 W11.B.123

Graph Storage systems

11/3/2016 CS535 Big Data, Fall 2016 W11.B.124

Graphs around us?

- Social graph
 - People you may know
- Ratings graph
 - Products you might like
- Social+ratings graph
 - Movies you should watch and the friends you should watch them with

11/3/2016 CS535 Big Data, Fall 2016 W11.B.125

Graph-like structure in RDBMS

- It can be stored and also handle queries
 - A single type of relationship
 - "who is my manager"
- Adding another relationship to the mix?
 - Requires schema changes and data movement
 - e.g. add preferences to lunch menu for each of the employees who have lunch at cafeteria
 - Add new table or column for new menu

EmployeeID	Name	ManagerID	EmployeeID	Likes_caf_menu_A
0001	ABC	5545	0001	1
0002	BCD	5567	0002	0
0003	CDE	6677	0003	1

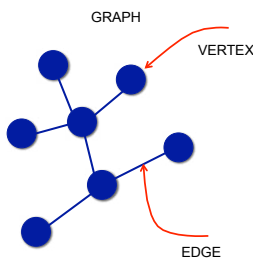
11/3/2016 CS535 Big Data, Fall 2016 W11.B.126

Traversal in Graph data

- A query on the graph is known as traversing the graph
- Traversing can be changed without having to change the nodes or edges
- Data can be traversed as it needs to be
- Traversing the joins or relationship is very fast
 - Relationship between nodes is not calculated at query time
 - It is actually persisted as a relationship

11/3/2016 CS535 Big Data, Fall 2016 W11.B.127

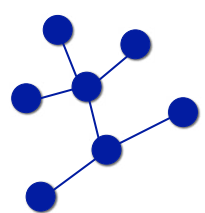
Graphs



$G=(V,E)$ where,
 V is a set of vertices
 E is a set of edges

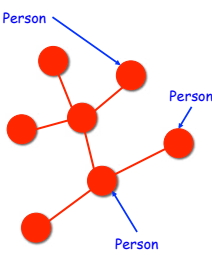
11/3/2016 CS535 Big Data, Fall 2016 W11.B.128

Homogeneous set of vertices connected by a homogeneous set of edges



11/3/2016 CS535 Big Data, Fall 2016 W11.B.129

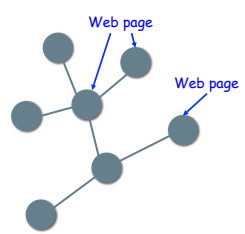
Homogeneous set of vertices connected by a homogeneous set of edges



People follows relationships

11/3/2016 CS535 Big Data, Fall 2016 W11.B.130

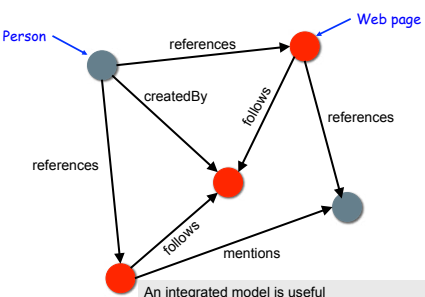
Homogeneous set of vertices connected by a homogeneous set of edges



Web pages contain links

11/3/2016 CS535 Big Data, Fall 2016 W11.B.131

Multi-relational property graph



An integrated model is useful
 Allows for more interesting/novel algorithms
 Allows for a universal model of things and their relationships

11/3/2016 CS535 Big Data, Fall 2016 W11.B.132

The property graph

- $G=(V,E,\lambda)$
- Directed, attributed, edge-labeled graph
- Multi-relational graph with key/value pairs on the elements

11/3/2016 CS535 Big Data, Fall 2016 W11.B.133

Graph query language

- Gremlin
 - A graph traversal language
 - Expresses complex graph traversals
 - Supports traversal patterns
 - Backtrack pattern, Path pattern, Loop pattern, Split/Merge Pattern, etc.

```
gremlin>
g.V('name','hercules').out('father').out('father').name
→ cronos
```

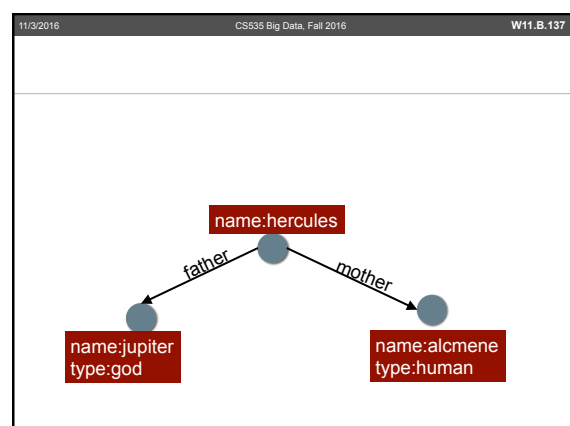
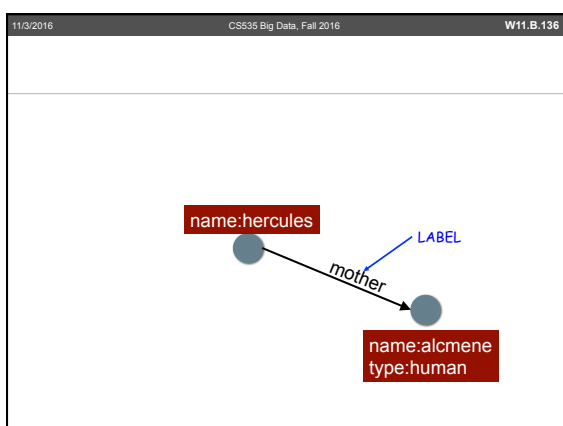
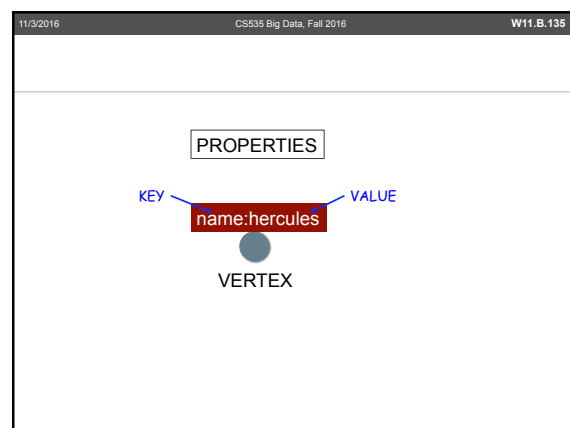
- g: for the current graph
- v('name','hercules'): get all vertices with name property "hercules"
- out('father'): traverse outgoing father edge's from Hercules
- out('father'): traverse outgoing father edge's from Hercules' father's vertex
- name: get the name property of the "hercules" vertex's grandfather

11/3/2016 CS535 Big Data, Fall 2016 W11.B.134

Traversing with Functions

- Adding vertex and edge

```
gremlin> theseus = g.addVertex([name:'theseus',type:'human'])
==>v[302]
gremlin> cerberus = g.V('name','cerberus').next()
==>v[48]
gremlin> g.addEdge(theseus,cerberus,'battled')
==>e[151200009:302:36028797018964038][302-battled->48]
```



11/3/2016 CS535 Big Data, Fall 2016 W11.B.138

What is Hercules' type?

```

graph TD
    H((name:hercules)) -- father --> J((name:jupiter type:god))
    H -- mother --> A((name:alcmene type:human))
  
```

* This example is from <http://thinkarelius.github.io/titan/>

11/3/2016 CS535 Big Data, Fall 2016 W11.B.139

```

graph TD
    H((name:hercules)) -- father --> J((name:jupiter type:god))
    H -- mother --> A((name:alcmene type:human))
  
```

```

gremlin> hercules.out('mother','father')
==>v[1]
==>v[2]
  
```

11/3/2016 CS535 Big Data, Fall 2016 W11.B.140

```

graph TD
    H((name:hercules)) -- father --> J((name:jupiter type:god))
    H -- mother --> A((name:alcmene type:human))
  
```

```

gremlin> hercules.out('mother','father').type
==>human
==>god
  
```

11/3/2016 CS535 Big Data, Fall 2016 W11.B.141

Hercules is 'demigod'

```

graph TD
    H((name:hercules type:demigod)) -- father --> J((name:jupiter type:god))
    H -- mother --> A((name:alcmene type:human))
  
```

```

gremlin> hercules.type='demigod'
==>demigod
  
```

11/3/2016 CS535 Big Data, Fall 2016 W11.B.142

Who else might Hercules know?

```

graph TD
    H((0)) -- knows --> C((1))
    H -- knows --> N((2))
    H -- knows --> Hy((3))
    C -- knows --> P((4))
    N -- knows --> Ne((5))
    Hy -- knows --> J((6))
  
```

```

gremlin> hercules
==>v[0]
  
```

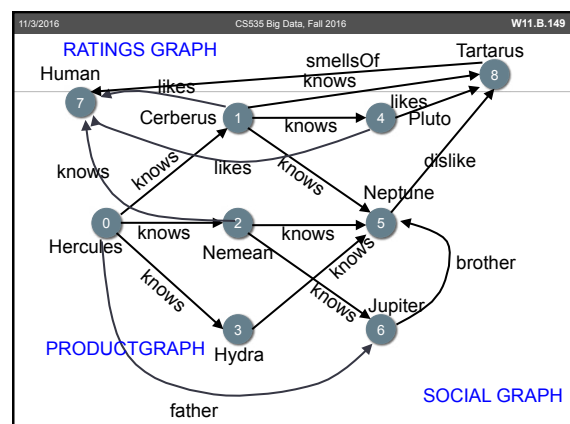
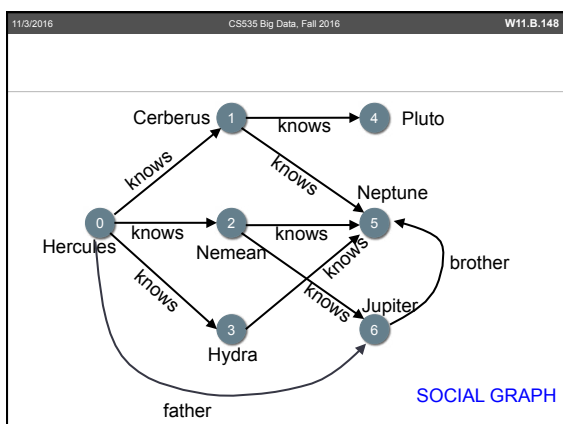
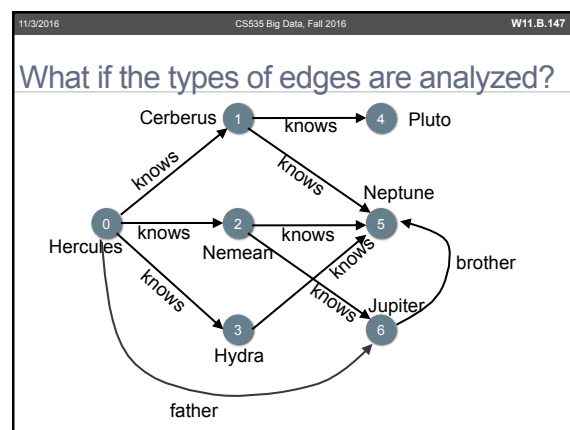
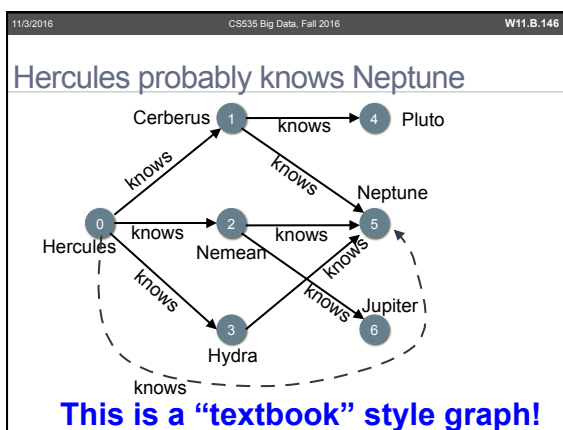
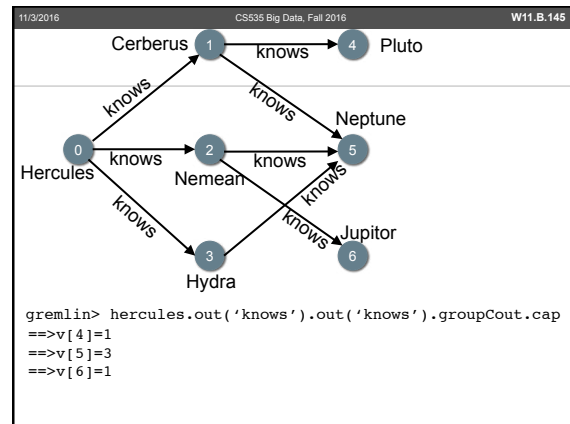
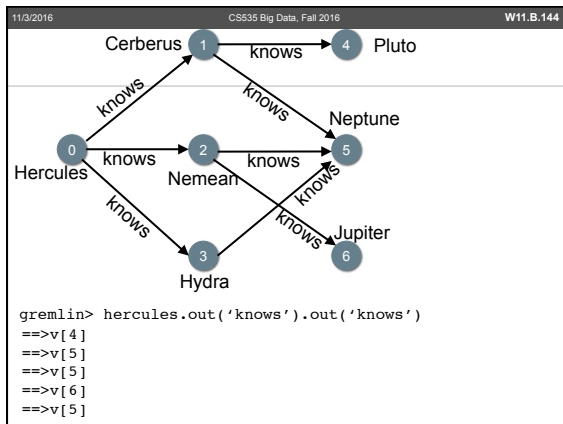
11/3/2016 CS535 Big Data, Fall 2016 W11.B.143

```

graph TD
    H((0)) -- knows --> C((1))
    H -- knows --> N((2))
    H -- knows --> Hy((3))
    C -- knows --> P((4))
    N -- knows --> Ne((5))
    Hy -- knows --> J((6))
  
```

```

gremlin> hercules.out('knows')
==>v[1]
==>v[2]
==>v[3]
  
```

11/3/2016 CS535 Big Data, Fall 2016 W11.B.150

Finding paths

- Movie graph
 - How is this person related to this film?
- Book graph
 - Which authors of this book also wrote a New York Times Bestseller?
- Movie+Book graph
 - Which movies are based on a book by a New York Times Bestseller?

11/3/2016 CS535 Big Data, Fall 2016 W11.B.151

Who played Hercules in What movie?

```
gremlin> hercules.out('depictedIn').as('movie').
  out('hasActor').out('role').retain(hercules).nack(2).
  Out('actor').as('star').select{it.name}
  ==>[movie:hercules in new york, star"Arnold sch
  Warzenegger]
```

Some part of the graph is removed for the space.

11/3/2016 CS535 Big Data, Fall 2016 W11.B.152

Queries generates social influence

- Who are the most influential people in
 - Java, mathematics, art, surreal art, politics?
- Which region of the social graph will propagate this advertisement the furthest?
- Which 3 experts should review this submitted article?
- Which people should I talk to at the upcoming conference and what topics should I talk to them about?

11/3/2016 CS535 Big Data, Fall 2016 W11.B.153

Pattern identification

- This connectivity pattern is a sign of a financial fraud.
 - Transaction graph
- Healthy discourse is typified by a discussion board with a branch factor in this range and a concept clique score in this range
 - Discussion graph

11/3/2016 CS535 Big Data, Fall 2016 W11.B.154

Graph computing engines: Memory-based graph framework

- Graph size is constrained by local machine's RAM
- Rich graph algorithm and visualization packages
- Applications
 - iGraph
 - <http://igraph.sourceforge.net>
 - NetworkX
 - <http://networkx.lanl.gov>
 - JUNG
 - <http://jung.sourceforge.net>

11/3/2016 CS535 Big Data, Fall 2016 W11.B.155

Graph computing engines: Disk-based graphs

- Graph size is constrained by the local disk
- Optimized for local graph algorithms
- Oriented towards property graphs
- Graph database
 - Neo4J
 - <http://neo4j.org>
 - InfiniteGraph
 - <http://objectivity.com>
 - OrientDB
 - <http://orientdb.org>
 - DEX
 - <http://www.sparsity-technologies.com/dex>

11/3/2016 CS535 Big Data, Fall 2016 W11.B.156

Graph computing engines: Cluster-based graphs

- Graph size is constrained by cluster's total RAM
- Optimized for global graph algorithm
- Bulk synchronous parallel processing
 - Pregel/Hama
 - <http://incubator.apache.org/hama>
 - Giraph
 - <http://incubator.apache.org/giraph>
 - GoldenOrb
 - <http://goldenorbos.org>

11/3/2016 CS535 Big Data, Fall 2016 W11.B.157

Titan: a Distributed Graph storage

11/3/2016 CS535 Big Data, Fall 2016 W11.B.158

Demands on Titan

- Processes graphs that have a 100+ billion edge scale with thousands of concurrent transactions
- Local graph traversals and batch graph processing
- Extremely scalable in the presence of:
 - Graph volume
 - Number of concurrent access to the graph

11/3/2016 CS535 Big Data, Fall 2016 W11.B.159

Key features

- Scalable graph
- Scalable access
- Real-time local traversals
- Global batch processing via Hadoop
- Open-source with Apache2 license

11/3/2016 CS535 Big Data, Fall 2016 W11.B.160

Options for the backend storage system

- Hbase
- Cassandra
- Oracle's BerkeleyDB

11/3/2016 CS535 Big Data, Fall 2016 W11.B.161

Inherited features from Cassandra

- High availability
- No single point of failure
- Replication scheme
- Caching
- Elastic scalability

11/3/2016 CS535 Big Data, Fall 2016 W11.B.162

Distinctive challenges

- Storage model
- Graph Partitioning
- Edge compression
- Vertex-centric indices

11/3/2016 CS535 Big Data, Fall 2016 W11.B.163

Titan Storage Model

11/3/2016 CS535 Big Data, Fall 2016 W11.B.164

Titan Storage Model

- Adjacency list in one column family
- Row key= vertex id
- Each property and edge in one column
 - Denormalized (i.e. stored twice)
- Direction and label/key as column prefix
 - Use slice predicate for quick retrieval

11/3/2016 CS535 Big Data, Fall 2016 W11.B.165

vertex	edge	edge	edge	edge
ID/Props	ID/Props/Label/ID	ID/Label/ID	ID/Props/Label/ID	ID/Label/ID
1 e:f	4 c:d A 2	5 B 0	6 g:h A 3	7 C 3

11/3/2016 CS535 Big Data, Fall 2016 W11.B.166

Type check (1/2)

- Type check for each of the data types
 - `time:"twelve"` is not acceptable
- Edge label signatures
 - Create association using signature of particular key
- TitanLabel


```
battled=g.makeType().name("battled").signature(timeKey)
```

11/3/2016 CS535 Big Data, Fall 2016 W11.B.167

Type check (2/2)

- Functional Declarations


```
TitanLabel father = g.makeType().name("father").functional()
```

11/3/2016 CS535 Big Data, Fall 2016 W11.B.168

Locking system

- Ensures consistency over non-consistent storage backends

1. Acquire lock at the end of the transaction
 - Locking mechanism depends on storage layer consistency guarantees
2. Verify original read
3. Fail transaction if any precondition is violated

GRAPH PARTITIONING

11/3/2016 CS535 Big Data, Fall 2016 W11.B.170

Graph Partitioning (1/2)

- Maximize locality of vertices
 - Co-locate vertices
- Titan maintains multiple ID rings
- OrderedPartitioner (from Cassandra) in storage backend
- Dynamically determines good partition and allocates corresponding IDs
- Graph Partitioning is NP complete problem
- Ongoing effort

11/3/2016 CS535 Big Data, Fall 2016 W11.B.171

Graph Partitioning (2/2)

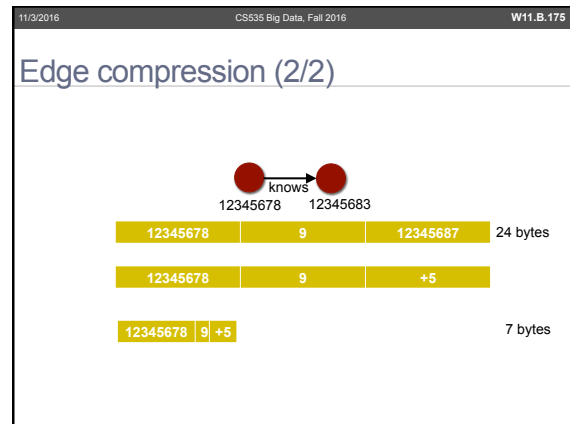
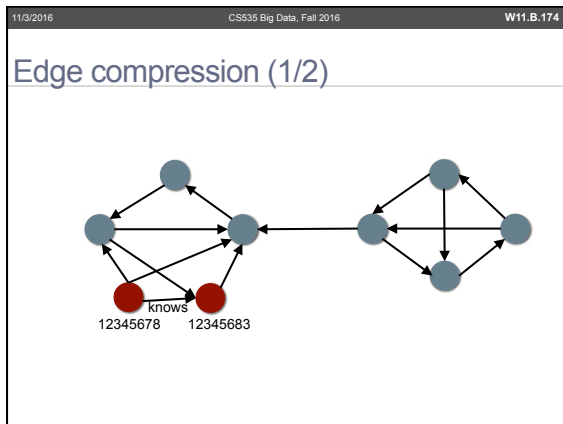
EDGE COMPRESSION

11/3/2016 CS535 Big Data, Fall 2016 W11.B.173

Small-world networks

- Natural graphs have a small world, community/cluster property

- High intra-connectivity within a community and low inter-connectivity between communities
- Watts, D.J., Strogatz, S. H., "Collective Dynamics of 'Small-World' Networks," Nature 393(6684), pp. 440-442, 1998



11/3/2016 CS535 Big Data, Fall 2016 W11.B.176

VERTEX-CENTRIC INDICES

11/3/2016 CS535 Big Data, Fall 2016 W11.B.177

Top 5 Twitters based on Followers

- 5) 34,952,307
 - Taylor Swift
- 4) 37,849,134
 - Barack Obama
- 3) 40,242,656
 - Lady Gaga
- 2) 44,951,137
 - Katy Perry
- 1) 45,728,072
 - Justin Bieber

<http://twitaholic.com/top100/followers>

11/3/2016 CS535 Big Data, Fall 2016 W11.B.178

The super node problem (1/2)

- The super node problem
 - Natural, real-world graphs contain vertices of high degree
 - Even if rare, their degree ensures that they exist on many paths
 - Traversing a high degree vertex means touching numerous incident edges and potentially touching most of the graph in only a few steps

11/3/2016 CS535 Big Data, Fall 2016 W11.B.179

The super node problem (2/2)

- A "super node" only exists from the vantage point of classic "textbook style" graphs
- In the world of property graphs, intelligent disk-based filtering can interpret a "super node" as a more manageable low-degree vertex
- Vertex-centric querying utilizes B-Trees and sort orders for speedy lookup of incident edges with particular qualities

11/3/2016 CS535 Big Data, Fall 2016 W11.B.180

Vertex-centric Indices

- Sort and index the incident edges (adjacent vertices)
 - Per vertex by primary key
 - Based on the edge's labels and properties
- Enables efficient focused traversals
 - Only retrieve edges that matter
- Uses push down predicates for quick, index-driven retrieval

<https://github.com/thinkaurelius/titan/wiki/Vertex-Centric-Indices>

11/3/2016 CS535 Big Data, Fall 2016 W11.B.181

Vertex-centric indices

```
TitanLabel battled =  
g.makeType().name("battled")  
.primaryKey(time)
```

11/3/2016 CS535 Big Data, Fall 2016 W11.B.182

Vertex-centric indices

```
Vertex.query().group("family")..
```

11/3/2016 CS535 Big Data, Fall 2016 W11.B.183

Vertex-centric Indices

Incident edge size	No vertex-centric indices	Vertex-centric indices
1000	0.82 ms	0.50 ms
10000	7.06 ms	0.52 ms
100000	70.90 ms	0.46 ms
1000000	778.94 ms	0.46 ms

11/3/2016 CS535 Big Data, Fall 2016 W11.B.184

Pushdown predicates

- Vertex-centric indices enables pushdown predicates efficiently
- A vertex has following information about the incident edges
 - Labels
 - Properties
 - Directions
- Vertex-centric query provides query from the perspective of a vertex.
- Edges that meet the query criteria are selected from the underlying graph representation

11/3/2016 CS535 Big Data, Fall 2016 W11.B.185

Pushdown predicates

`vertex.query()`

8 edges

11/3/2016 CS535 Big Data, Fall 2016 W11.B.186

Pushdown predicates

```
vertex.query().direction(out)
```

7 edges

11/3/2016 CS535 Big Data, Fall 2016 W11.B.187

Pushdown predicates

```
vertex.query().direction(OUT).labels("likes")
```

5 edges

11/3/2016 CS535 Big Data, Fall 2016 W11.B.188

Pushdown predicates

```
vertex.query().direction(OUT).labels("likes").has("stars",5)
```

1 edges

11/3/2016 CS535 Big Data, Fall 2016 W11.B.189

Pushdown predicates

- Query Query.direction(Direction)
- Query Query.labels(String... labels)
- Query Query.has(String, Object, Compare)
- Query Query.has(String, Object)
- Query Query.range(String, Object, Object)
- Iterable<Vertex> Query.vertices()
- Iterable<Edge> Query.edges()

INTERACTING WITH GRAPH ANALYTICS

11/3/2016 CS535 Big Data, Fall 2016 W11.B.191

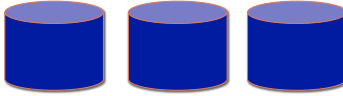
Faunus

- Hadoop-based Graph computing framework
- Graph analytics
- Breath-first traversals
- Global graph computations
- Batch big graph data

11/3/2016 CS535 Big Data, Fall 2016 W11.B.192

An Adjacency list + Cluster


0
1
2
3
4
5
6
7
8
9
10
11



11/3/2016 CS535 Big Data, Fall 2016 W11.B.193


A Distributed adjacency list

0	4	8
1	5	9
2	6	10
3	7	11



11/3/2016 CS535 Big Data, Fall 2016 W11.B.194

Fanus workflow



Compressed HDFS Graphs

- Stored in sequence files
- Variable length encoding
- Prefix compression

11/3/2016 CS535 Big Data, Fall 2016 W11.B.195

Graph analytics with Titan

- Titan
 - Stores a massive-scale property graph, allowing real time traversals and updates
- Analysis results
 - Batch processing of large graphs with Hadoop
- Fulgora
 - Runs global graph algorithms on large, compressed in-memory graphs

Handwritten notes: "Bulk load" and "MapReduce" with arrows pointing to the Titan and Fulgora sections respectively.

PREGEL:
A SYSTEM FOR LARGE-SCALE
GRAPH PROCESSING

11/3/2016 CS535 Big Data, Fall 2016 W11.B.197

This material is built based on,

- Grzegorz Malewicz, Matthew H. Austern, Aart J.C. Bik, Names C. Dehnert, Ilan Horn, Naty Leiser, Grzegorz Czajkowski, "Pregel: a system for large-scale graph processing", Proceedings of the 2010 ACM SIGMOD International Conference on Management of Data, pp. 135-146
- Apache's Hama
 - Open source project inspired by Pregel
 - <http://hama.apache.org>

11/3/2016 CS535 Big Data, Fall 2016 W11.B.198

Graph analysis at Google?

- MapReduce tasks
 - Google's 80% of data analysis
 - Large-scale web search indexing
 - Clustering problems for Google News
 - Produce reports for popular queries (e.g. Google Trend)
 - Processing of satellite imagery data
 - Language model processing for statistical machine translations
 - Large-scale machine learning problems
 - Back-up/restore
- The other 20%?

11/3/2016 CS535 Big Data, Fall 2016 W11.B.199

Graph analysis at Google?

- Large graph analysis
 - Graph algorithms
 - PageRank
 - Shortest path
 - Connected components
 - Clustering techniques
- Graph data
 - Web graph
 - Transportation routes
 - Citation relationships
 - Social networks

11/3/2016 CS535 Big Data, Fall 2016 W11.B.200

MapReduce is NOT great for graph processing

- Many iterations are needed for parallel graph processing
- Materializations of intermediate results at every MapReduce iteration causes performance bottleneck

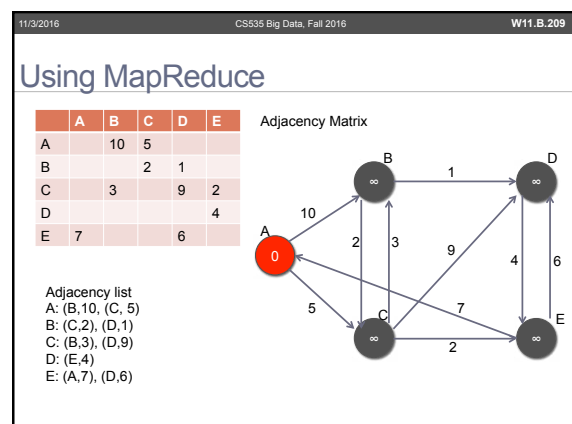
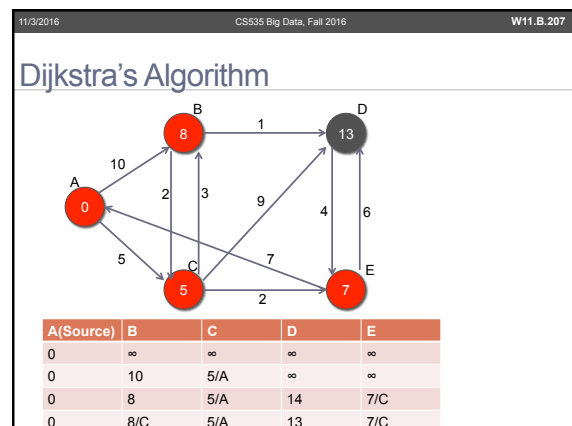
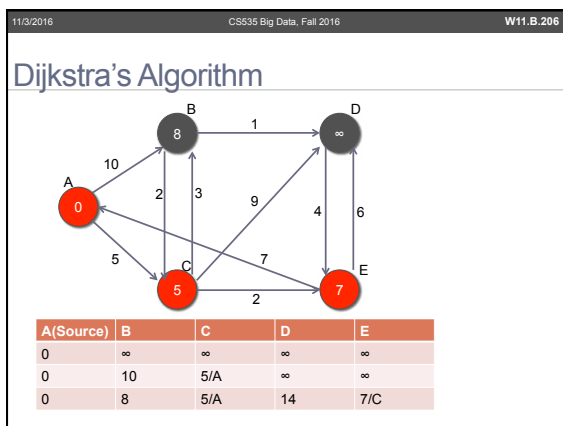
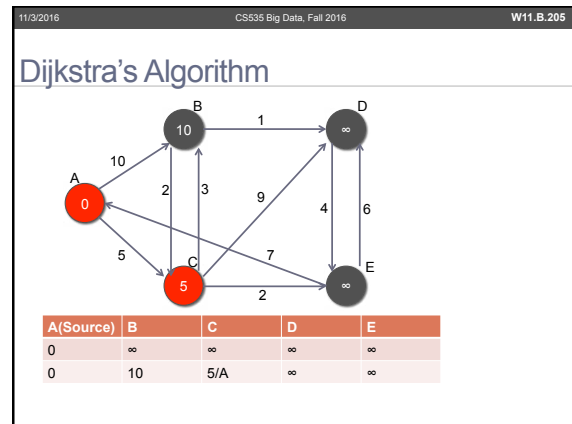
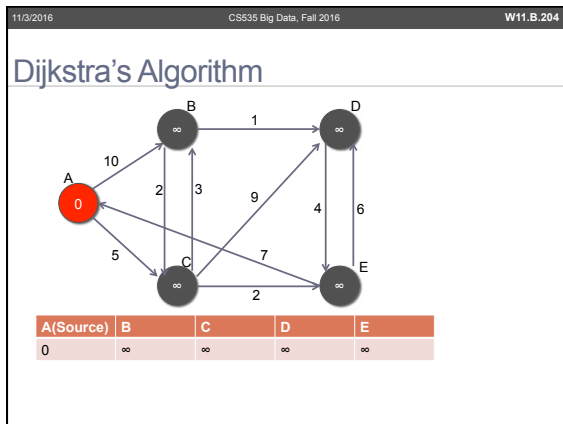
11/3/2016 CS535 Big Data, Fall 2016 W11.B.201

Single Source Shortest Path (SSSP)

- Find shortest path from a source node to all target nodes
- If you have a single processor machine?
 - Dijkstra's algorithm

Finding SSSP using Dijkstra's Algorithm

11/3/2016 CS535 Big Data, Fall 2016 W11.B.203



11/3/2016 CS535 Big Data, Fall 2016 W11.B.210

Using MapReduce

Map input: <nodeID, :dist, adj list>>

```
<A, <0, <(B,10), (C, 5)>>>
<B, <inf, <(C,2), (D,1)>>>
<C, <inf, <(B,3), (D,9)>>>
<D, <inf, <(E,4)>>>
<E, <inf, <(A,7), (D,6)>>>
```

Map output: <dest node ID, dist>

```
<B,10><C,5>
<C, inf><D, inf>
<B, inf><D, inf> <E, inf>
<E, inf>
<A, inf><D, inf>
```

Flushed to local DFS

11/3/2016 CS535 Big Data, Fall 2016 W11.B.211

Using MapReduce

Reduce input: <nodeID, dist>

```
<A, <0, <(B,10), (C, 5)>>>
<A, inf>
```

```
<B, <inf, <(C,2), (D,1)>>>
<B,10><B, inf>
```

```
<C, <inf, <(B,3), (D,9)>>>
<C,5><C, inf>
```

```
<D, <inf, <(E,4)>>>
<D, inf><D, inf><D, inf>
```

```
<E, <inf, <(A,7), (D,6)>>>
<E, inf><E, inf>
```

11/3/2016 CS535 Big Data, Fall 2016 W11.B.212

Using MapReduce

Reduce input: <nodeID, dist>

```
<A, <0, <(B,10), (C, 5)>>>
<A, inf>
```

```
<B, <inf, <(C,2), (D,1)>>>
<B,10><B, inf>
```

```
<C, <inf, <(B,3), (D,9)>>>
<C,5><C, inf>
```

```
<D, <inf, <(E,4)>>>
<D, inf><D, inf><D, inf>
```

```
<E, <inf, <(A,7), (D,6)>>>
<E, inf><E, inf>
```

11/3/2016 CS535 Big Data, Fall 2016 W11.B.213

Using MapReduce

Reduce input: <nodeID, <dist,adj list>>

= Map input for next iteration

```
<A, <0, <(B,10), (C, 5)>>>
<B, <10, <(C,2), (D,1)>>>
<C, <5, <(B,3), (D,9)>>>
<D, <inf, <(E,4)>>>
<E, <inf, <(A,7), (D,6)>>>
```

Map output: <dest node ID, dist>

```
<B,10><C,5> <A, <0, <(B,10), (C, 5)>>>
<C,12><D,11> <B, <10, <(C,2), (D,1)>>>
<B,8><D,14><C, <5, <(B,3), (D,9)>>>
<E, inf> <D, <inf, <(E,4)>>>
<A, inf><C,inf> <E, <inf, <(A,7), (D,6)>>>
```

Flushed to local DFS

11/3/2016 CS535 Big Data, Fall 2016 W11.B.214

Using MapReduce

Reduce input: <nodeID, <dist,adj list>>

= Map input for next iteration

```
<A, <0, <(B,10), (C, 5)>>>
<A, inf>
```

```
<B, <10, <(C,2), (D,1)>>>
<B,10><B, 8>
```

```
<C, <5, <(B,3), (D,9)>>>
<C,5><C,12>
```

```
<D, <inf, <(E,4)>>>
<D,11><D,14><D,inf>
```

```
<E, <inf, <(A,7), (D,6)>>>
<E, inf><E,7>
```

11/3/2016 CS535 Big Data, Fall 2016 W11.B.215

Using MapReduce

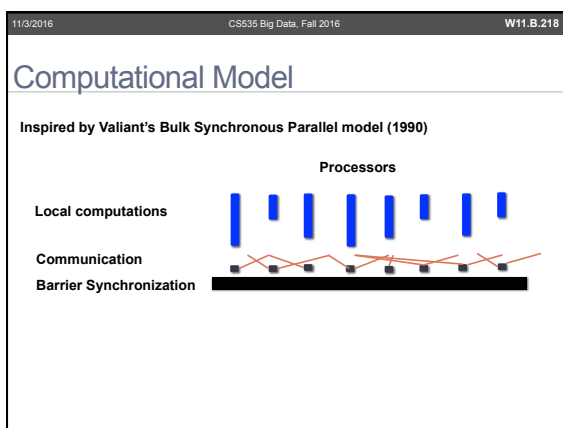
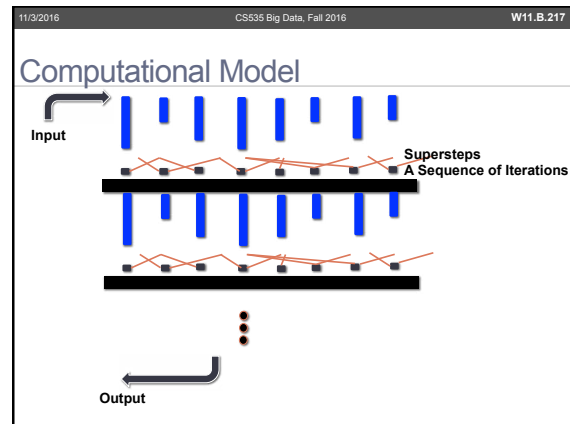
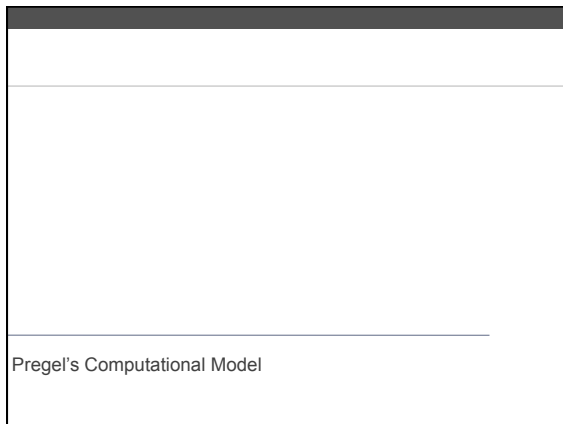
Reduce input: <nodeID, <dist,adj list>>

= Map input for next iteration

```
<A, <0, <(B,10), (C, 5)>>>
<B, <8, <(C,2), (D,1)>>>
<C, <5, <(B,3), (D,9)>>>
<D, <11, <(E,4)>>>
<E, <7, <(A,7), (D,6)>>>
```

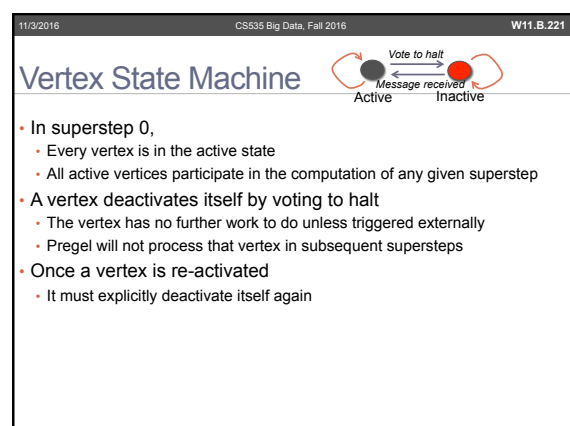
Flushed to local DFS

Keep going..



- Computational Model (1/2)
- Superstep: the vertices compute in parallel
 - Each vertex
 - Receives messages from the previous superstep
 - Executes the same user-defined function
 - Sends messages to other vertices
 - Mutates the topology of the graph if need be
 - Votes to halt if it has no further work to do
 - When to terminate?
 - All vertices are simultaneously inactive
 - There are no messages in transit

- Computation Model (2/2)
- Input to the Pregel computation
 - A directed graph
 - Vertex
 - String vertex ID
 - Associated user defined value
 - Edge
 - Associated with their source vertices
 - User defined value and a target vertex ID
 - Computation in the vertex
 - Executes the same user-defined function
 - Modifies the state
 - Sometimes changes the outgoing edges
 - Receive/send message
 - Mutate topology
 - There is no computation associated with the edges



11/3/2016 CS535 Big Data, Fall 2016 W11.B.222

Output of a Pregel program

- Set of values explicitly output by the vertices
 - Often a directed graph isomorphic to the input
 - E.g. clustering algorithm
 - Creates small set of disconnected vertices selected from a large graph
 - E.g. graph mining algorithm
 - Generates aggregated statistics mined from the graph

11/3/2016 CS535 Big Data, Fall 2016 W11.B.223

Message Passing (1/2)

- Vertices communicate directly with one another by sending messages
 - Message value
 - Name of the destination vertex
- A vertex can send any number of messages in a superstep
- There is no guaranteed order of messages in the iterator. However,
 - Message is delivered reliably
 - There will be no duplicate

11/3/2016 CS535 Big Data, Fall 2016 W11.B.224

Message Passing (2/2)

- Common usage pattern
 - A vertex V to iterate over its outgoing edges and sending a message to the destination vertex of each edge
- Destination vertex need not be a neighbor of V
 - E.g.: A vertex can learn the identifier of a non-neighbor from a message received earlier
 - E.g.: implicitly vertex info is distributed
- If destination does not exist, user-defined handler will be executed.
 - Create the missing vertex or remove the dangling edge

11/3/2016 CS535 Big Data, Fall 2016 W11.B.225

SSSP using parallel BFS in Pregel

Legend: Active vertex sends message to inactive vertex. Inactive vertex becomes active upon receiving a message.

11/3/2016 CS535 Big Data, Fall 2016 W11.B.226

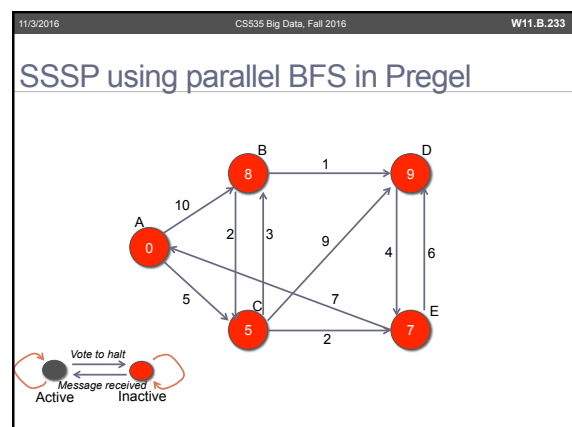
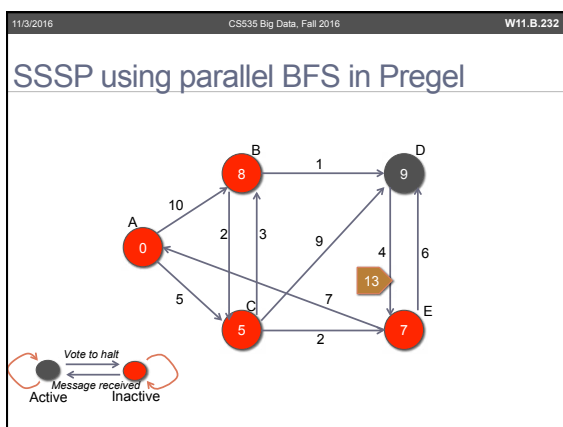
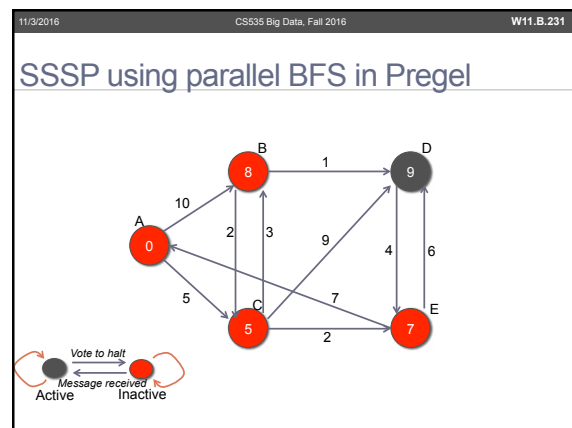
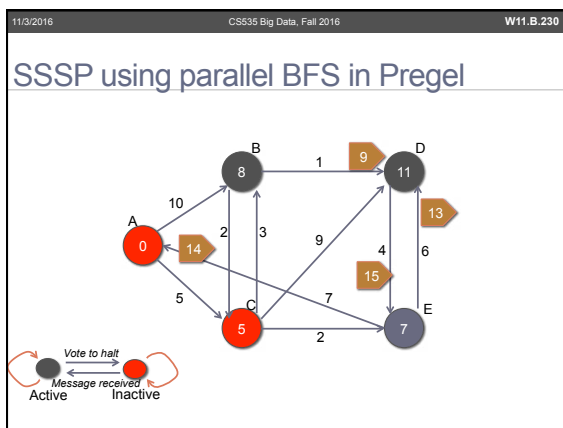
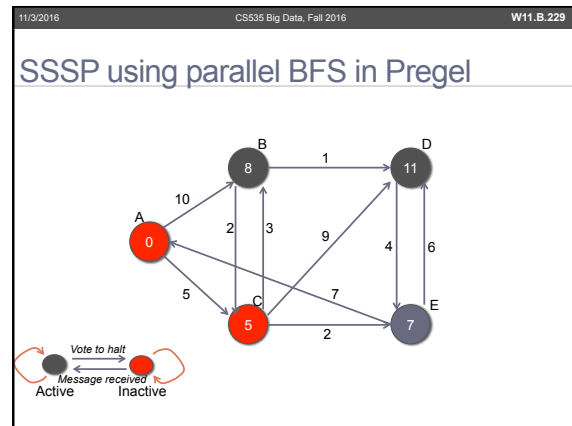
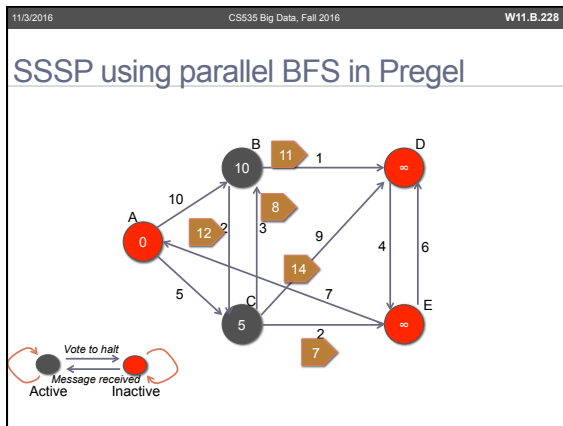
SSSP using parallel BFS in Pregel

Legend: Active vertex sends message to inactive vertex. Inactive vertex becomes active upon receiving a message.

11/3/2016 CS535 Big Data, Fall 2016 W11.B.227

SSSP using parallel BFS in Pregel

Legend: Active vertex sends message to inactive vertex. Inactive vertex becomes active upon receiving a message.



11/3/2016 CS535 Big Data, Fall 2016 W11.B.235

Combiners

- Sending message incurs overhead
- System can combine several messages intended for a vertex V into a single message
 - E.g. calculating sum: values intended for V will be combined as the local sum and then passed to the vertex V
- Combiner class overrides a virtual `Combine()` method
- For SSSP, more than 4x reduction in message traffic

Advanced Message Passing Features

11/3/2016 CS535 Big Data, Fall 2016 W11.B.236

Aggregators

- A mechanism for global communication
- Each vertex can provide a value to an aggregator in superstep S
 - System combines these values using a reduction operator
- Resulting value is available to all vertices in the superstep $S+1$

11/3/2016 CS535 Big Data, Fall 2016 W11.B.237

Aggregators

- Pre-defined aggregator
 - min, max or sum operations
- User-defined aggregator
 - Specify the process in the Aggregator class
- Sticky aggregator
 - Uses input values from all supersteps
 - E.g. maintaining global edge count

11/3/2016 CS535 Big Data, Fall 2016 W11.B.238

Topology Mutations

- Some graph algorithms need to change the graph's topology
 - Clustering algorithm
 - Might replace each cluster with a single vertex
 - Minimum spanning tree
 - Might remove all but the tree edges
- Pregel allows the edges to be removed/added

11/3/2016 CS535 Big Data, Fall 2016 W11.B.239

Resolving conflicting requests (1/2)

- Multiple vertices may issue conflicting requests in the same superstep
 - E.g. two requests to add a vertex V , with different initial values
- Partial ordering
 - Mutations become effective in the superstep after the requests were issued
 - Within that superstep, removals are performed first
 - Edge removal before the vertex removal
 - After removal, addition will be performed
 - Vertex addition before the edge addition

11/3/2016 CS535 Big Data, Fall 2016 W11.B.240

Resolving conflicting requests (2/2)

- If there are multiple requests to create the same vertex in the same superstep
 - By default the system just picks one arbitrarily
 - Users can specify their own conflict resolution policy
 - By defining a handler method in their `Vertex` subclass

11/3/2016 CS535 Big Data, Fall 2016 W11.B.241

Global mutation

- Lazy coordination mechanism
 - Global mutation does not require coordination until the point when they are applied
 - A vertex `V` will handle its own modifications

Contrasting Pregel & MapReduce

11/3/2016 CS535 Big Data, Fall 2016 W11.B.243

Pregel vs. MapReduce

- Many of graph algorithms can be written as a series of chained MapReduce invocations
- Pregel
 - Once the vertices and edges are loaded into computing nodes, they will stay on that machine
 - Only messages will be transferred through the network
- MapReduce
 - Passes the entire state of graph for every iteration
 - External coordinator is required to create a "chain" of MapReduce jobs

System Architecture

11/3/2016 CS535 Big Data, Fall 2016 W11.B.245

System Architecture

- Master/worker model
 - Worker
 - Processes user-defined tasks
 - Communicates with other workers (messaging)
 - Master
 - Maintains information about workers
 - No portion of graph assigned
 - Recovers from faults
 - Uses monitoring tools
- Underlying persistent data storage: GFS or BigTable
- Temporary data is stored on local disk

11/3/2016 CS535 Big Data, Fall 2016 W11.B.246

Step-by-step execution (1/4)

1. A client launches a Pregel job
 - Many copies of the user program begin executing on a cluster of machines
 - One of these copies acts as the master
 - Workers use the cluster management system's name service to discover the master's location
 - Send registration messages to the master
2. The master assigns a partition of the input to each worker
 - Each worker loads the vertices and marks them as active

11/3/2016 CS535 Big Data, Fall 2016 W11.B.247

Step-by-step execution (2/4)

2. The master assigns a partition of the input to each worker
 - Worker:
 - Loads the vertices and marks them as active
 - Maintains the state of its section of the graph
 - Executes user's `compute()` method on its vertices
 - Manages messages to and from other workers

11/3/2016 CS535 Big Data, Fall 2016 W11.B.248

Step-by-step execution (3/4)

3. The master instructs each worker to perform a superstep
 - Performs user-defined function on the active vertices
 - Messages are sent *asynchronously*
 - Before the end of the superstep
 - This step is repeated until: *(a) all of the vertices are inactive simultaneously && (b) no messages are transferred*

11/3/2016 CS535 Big Data, Fall 2016 W11.B.249

Step-by-step execution (4/4)

- 4 After the computation halts, the master may instruct each worker to save its portion of the graph

11/3/2016 CS535 Big Data, Fall 2016 W11.B.250

Fault tolerance (1/2)

- System maintains checkpoints
- The master periodically requests the workers to save the state of their partitions to persistent storage
 - State is saved as checkpoints, and includes ...
 - Vertex values, edge values, incoming messages

11/3/2016 CS535 Big Data, Fall 2016 W11.B.251

Fault tolerance (2/2)

- Failure detection
 - Regular "ping" message
- Recovery
 - The master reassigns graph partitions to the current available workers
 - The workers all reload their partition state from most recent available checkpoint

PageRank Algorithm in Pregel

PageRank Algorithm

- Link analysis algorithm
- Probability distribution
 - Represents the likelihood that a person randomly clicking on links will arrive at any particular page
- Probability
 - Between 0 and 1
 - PageRank of 0.5
 - There is a 50% chance that a person clicking on a random link will be directed to the document with the 0.5 PageRank

Iterative approach

- A link to a page counts as a vote of support
- At $t=0$,
- $PR(p_i;0)=1/N$
- At each time step, the computation yields,


$$PR(p_j;t+1) = \frac{1-d}{N} + d \sum \frac{PF(p_j;t)}{L(p_j)}$$

Iterative approach


- Damping factor, d
 - An imaginary surfer who is randomly clicking on links and he/she will eventually stop
 - The probability that the imaginary person will continue in that step
 - Generally assumed as around 0.85

$$PR(p_i;t+1) = \frac{1-d}{N} + d \sum \frac{PF(p_j;t)}{L(p_j)}$$

- Where, $PR(p,t)$: PageRank for the page p at time step t
- $L(p)$: number of links from page p




- $PR(A)=1, PR(B)=1, d=0.85$
- $PR(A) = (1-d) + (0.85) \times 1 = 1$
- $PR(B) = (1-d) + (0.85) \times 1 = 1$
- $PR(A)=0, PR(B)=0, d=0.85$
- $PR(A) = (1-d) + (0.85) \times 0 = 0.15$
- $PR(B) = (1-d) + (0.85) \times 0.15 = 0.2775$
- $PR(A)=0.15, PR(B)=0.2775, d=0.85$
- $PR(A) = (1-d) + (0.85) \times 0.2775 = 0.385875$
- $PR(B) = (1-d) + (0.85) \times 0.385875 = 0.47799375$



- $PR(A)=0.385875, PR(B)=0.47799375, d=0.85$
- $PR(A) = (1-d) + (0.85) \times 0.47799375 = 0.5562946875$
- $PR(B) = (1-d) + (0.85) \times 0.5562946875 = 0.622850484375$
- The numbers just keep going up.
- But will the numbers stop increasing when they get to 1.0?
- What if a calculation over-shoots and goes above 1.0?

11/3/2016 CS535 Big Data, Fall 2016 W11.B.258



• $PR(A)=40, PR(B)=40, d=0.85$
• $PR(A) = (1-d) + (0.85) \times 40 = 34.25$
• $PR(B) = (1-d) + (0.85) \times 34.25 = 29.1775$

• $PR(A)=34.25, PR(B)=29.1775, d=0.85$
• $PR(A) = (1-d) + (0.85) \times 29.1775 = 24.950875$
• $PR(B) = (1-d) + (0.85) \times 34.25 = 21.35824375$

• The numbers are heading down.
• The numbers will get to 1.0 and stop

11/3/2016 CS535 Big Data, Fall 2016 W11.B.259

In Pregel

```
class PageRankVertex:public Vertex<double, void, double>{
public:
    virtual void Compute(MessageIterator* msgs){
        if (superstep()>=1){
            double sum = 0;
            for(; !msgs->Done(); msgs->Next())
                sum +=msgs->Value();
            *MutableValue()=0.15/NumVertices(+0.85*sum;
        }
        if (superstep() < 30) {
            const int64 n = GetOutEdgeIterator().size();
            SendMessageToAllNeighbors(GetValue()/n);
        }
        else{
            VoteToHalt();
        }
    }
};
```

11/3/2016 CS535 Big Data, Fall 2016 W11.B.260