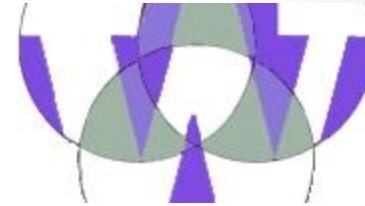


# Introduction to Data Science

Lecture 9; December 7<sup>th</sup>, 2016

Ernst Henle  
[ErnstHe@UW.edu](mailto:ErnstHe@UW.edu)  
Skype: ernst-henle

# Agenda



- Announcements
  - Social Interactions on LinkedIn: Continue to discuss statistics, job prospects, etc.
- HIVE Lab with Hue, Impala, and File Browser
- MapReduce
- Break
- MapReduce Lab
- Break
- Graph Data Intro
- Quiz Graph Data (Measures)
- Graph Data Popularity
- Quiz Graph Data (Popularity)
- Google Matrix
- Assignment. See assignment slides at the end of the deck. Complete all assignments items from all assignment slides. Submit by this Saturday at 11:57 PM.

# Hadoop-2 Labs

( 3 )

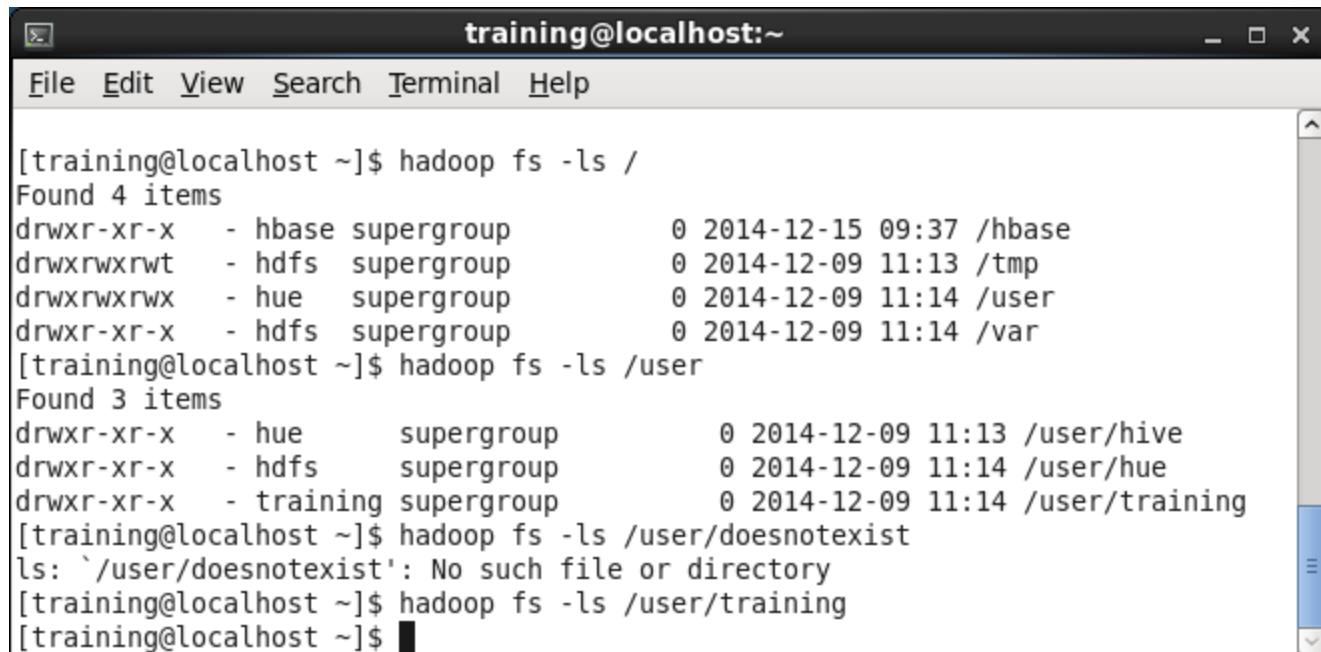
# HDFS Lab (0)



"So you want to hire me as a Data Scientist  
for Intelligent Virtualized Deep Machine  
Learning Real-time Big Data in the Cloud  
for Social Networks? Ok, but if you also  
want Hadoop, increase my salary by 50%."

# HDFS Lab (1)

- Enter these commands into the console to list some folders and files inside HDFS:
  - \$ hadoop fs -ls /
  - \$ hadoop fs -ls /user
  - \$ hadoop fs -ls /user/doesnotexist
  - \$ hadoop fs -ls /user/training
- “fs” stands for file system (HDFS). “ls” is an argument to list HDFS files.

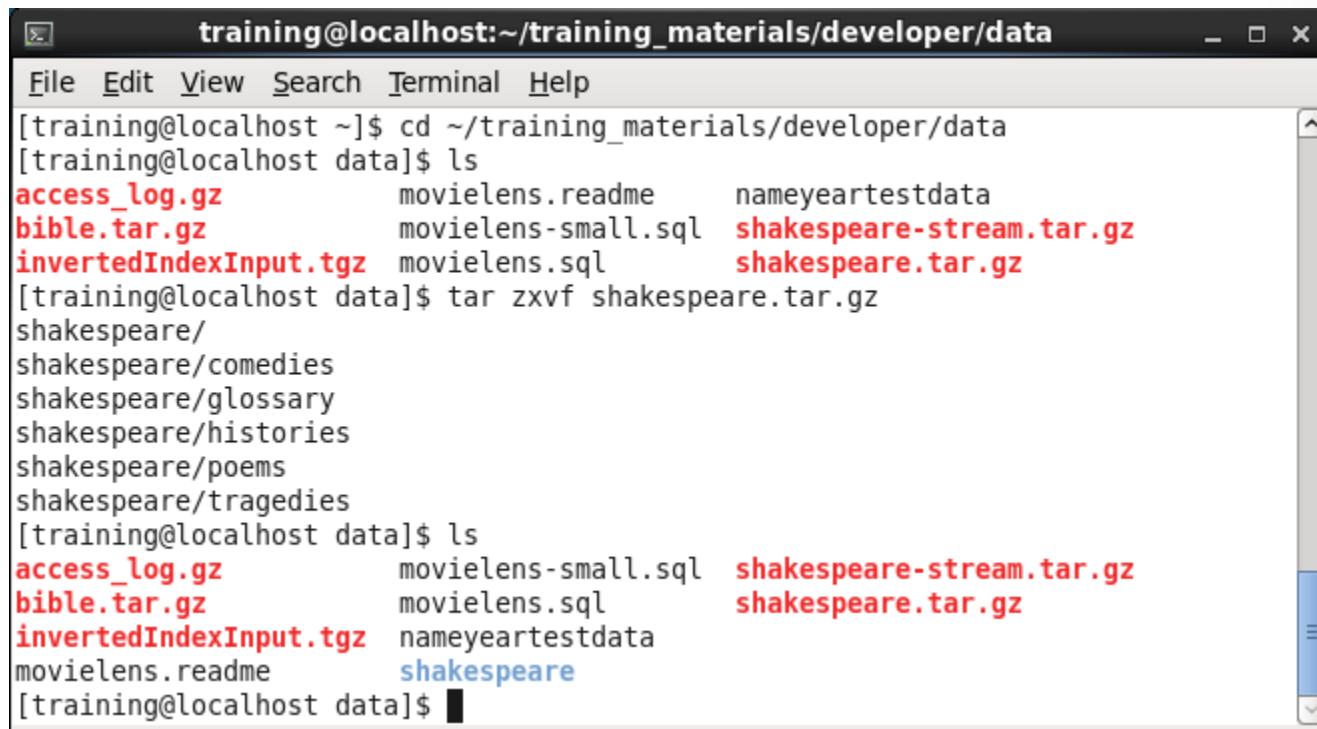


A screenshot of a terminal window titled "training@localhost:~". The window shows the following command-line session:

```
[training@localhost ~]$ hadoop fs -ls /
Found 4 items
drwxr-xr-x  - hbase supergroup          0 2014-12-15 09:37 /hbase
drwxrwxrwt  - hdfs supergroup          0 2014-12-09 11:13 /tmp
drwxrwxrwx  - hue  supergroup          0 2014-12-09 11:14 /user
drwxr-xr-x  - hdfs supergroup          0 2014-12-09 11:14 /var
[training@localhost ~]$ hadoop fs -ls /user
Found 3 items
drwxr-xr-x  - hue      supergroup      0 2014-12-09 11:13 /user/hive
drwxr-xr-x  - hdfs     supergroup      0 2014-12-09 11:14 /user/hue
drwxr-xr-x  - training  supergroup     0 2014-12-09 11:14 /user/training
[training@localhost ~]$ hadoop fs -ls /user/doesnotexist
ls: `/user/doesnotexist': No such file or directory
[training@localhost ~]$ hadoop fs -ls /user/training
[training@localhost ~]$ █
```

# HDFS Lab (2)

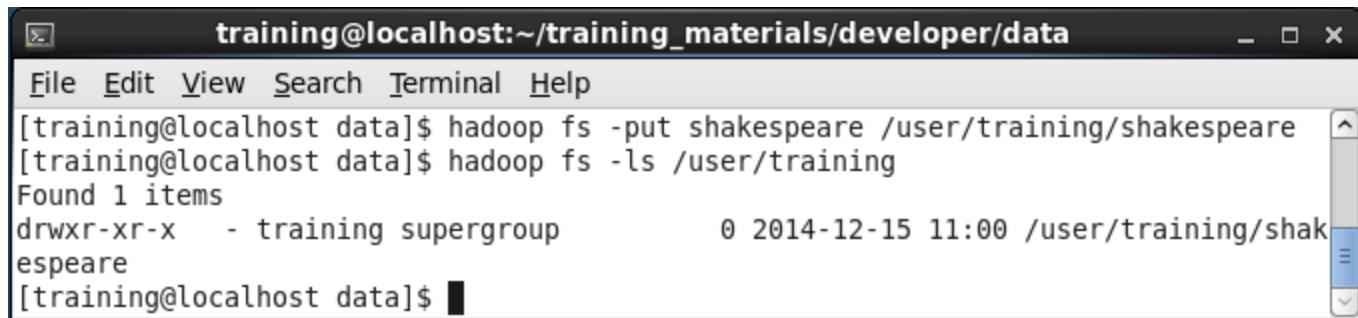
- Enter these commands into the console to find and expand the shakespeare tar on the Linux OS :
  - \$ cd ~/training\_materials/developer/data
  - \$ ls
  - \$ tar zxvf shakespeare.tar.gz
  - \$ ls



```
training@localhost:~/training_materials/developer/data
File Edit View Search Terminal Help
[training@localhost ~]$ cd ~/training_materials/developer/data
[training@localhost data]$ ls
access_log.gz      movielens.readme    nameyeartestdata
bible.tar.gz       movielens-small.sql   Shakespeare-stream.tar.gz
invertedIndexInput.tgz movielens.sql     Shakespeare.tar.gz
[training@localhost data]$ tar zxvf shakespeare.tar.gz
shakespeare/
shakespeare/comedies
shakespeare/glossary
shakespeare/histories
shakespeare/poems
shakespeare/tragedies
[training@localhost data]$ ls
access_log.gz      movielens-small.sql   Shakespeare-stream.tar.gz
bible.tar.gz       movielens.sql        Shakespeare.tar.gz
invertedIndexInput.tgz nameyeartestdata
movielens.readme    Shakespeare
[training@localhost data]$
```

# HDFS Lab (3)

- Enter these commands into the console to place the shakespeare directory into HDFS under training and then verify that the directory is in HDFS:
  - \$ hadoop fs -put shakespeare /user/training/shakespeare
  - \$ hadoop fs -ls /user/training

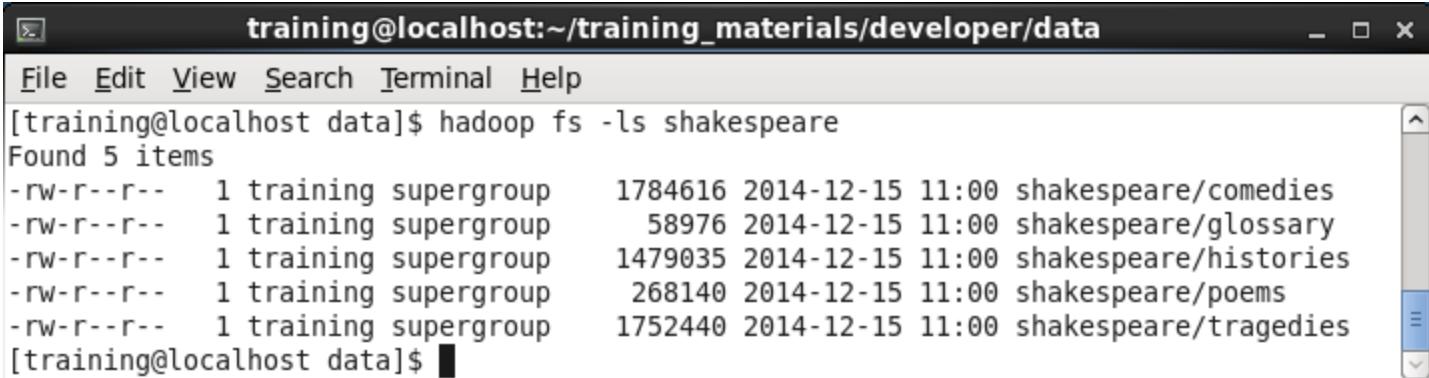


A screenshot of a terminal window titled "training@localhost:~/training\_materials/developer/data". The window contains the following text:

```
[training@localhost data]$ hadoop fs -put shakespeare /user/training/shakespeare
[training@localhost data]$ hadoop fs -ls /user/training
Found 1 items
drwxr-xr-x  - training supergroup          0 2014-12-15 11:00 /user/training/shak
espeare
[training@localhost data]$
```

# HDFS Lab (4)

- Enter this command into the console to list the 5 files in the `shakespeare` directory of HDFS:
  - `$ hadoop fs -ls /user/training/shakespeare`

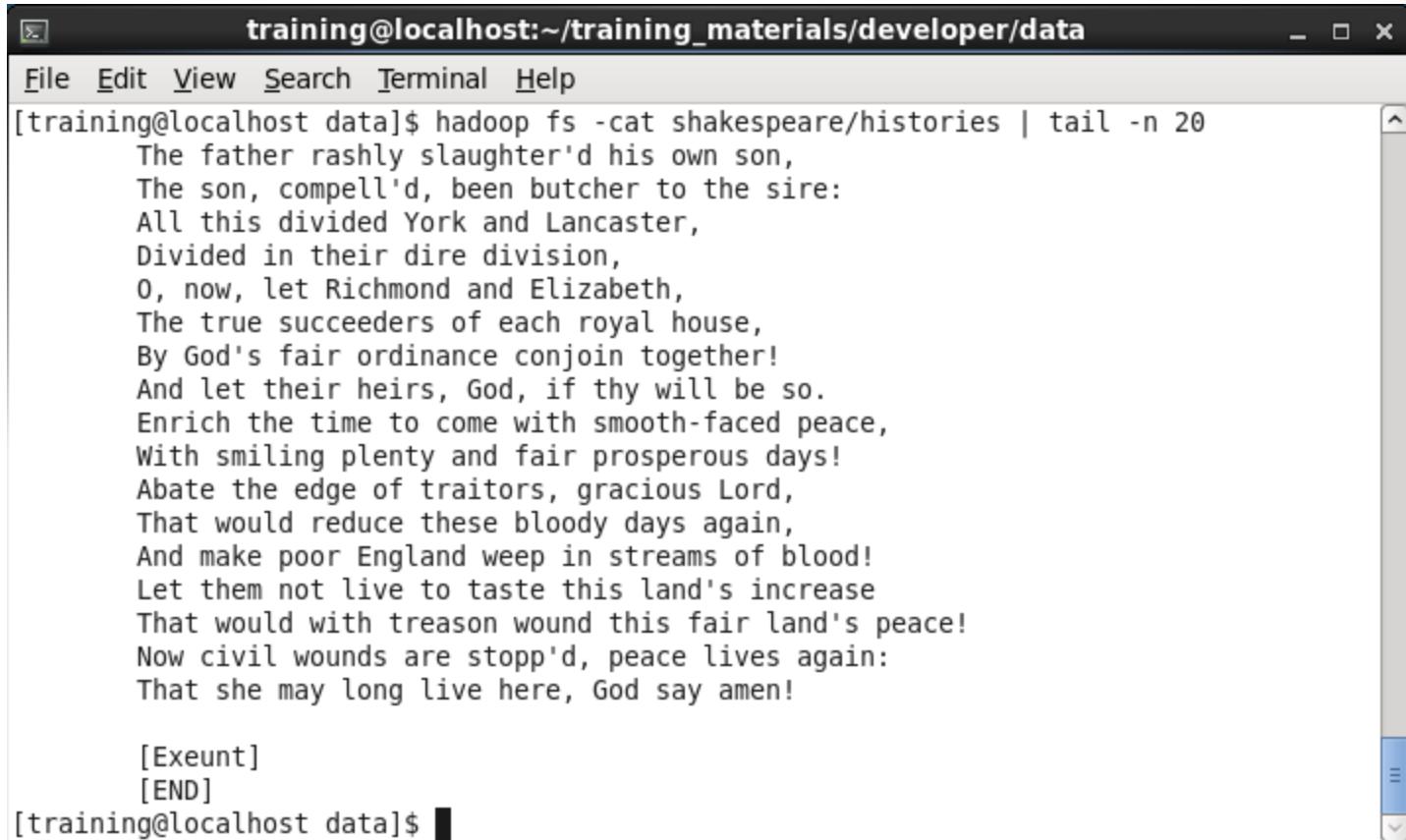


A screenshot of a terminal window titled "training@localhost:~/training\_materials/developer/data". The window contains the following text:

```
[training@localhost data]$ hadoop fs -ls shakespeare
Found 5 items
-rw-r--r-- 1 training supergroup 1784616 2014-12-15 11:00 shakespeare/comedies
-rw-r--r-- 1 training supergroup 58976 2014-12-15 11:00 shakespeare/glossary
-rw-r--r-- 1 training supergroup 1479035 2014-12-15 11:00 shakespeare/histories
-rw-r--r-- 1 training supergroup 268140 2014-12-15 11:00 shakespeare/poems
-rw-r--r-- 1 training supergroup 1752440 2014-12-15 11:00 shakespeare/tragedies
[training@localhost data]$
```

# HDFS Lab (5)

- Enter this command into the console to read the last 20 lines of the histories file in HDFS:
  - \$ hadoop fs -cat **shakespeare/histories** | tail -n 20



The screenshot shows a terminal window titled "training@localhost:~/training\_materials/developer/data". The window contains the following text:

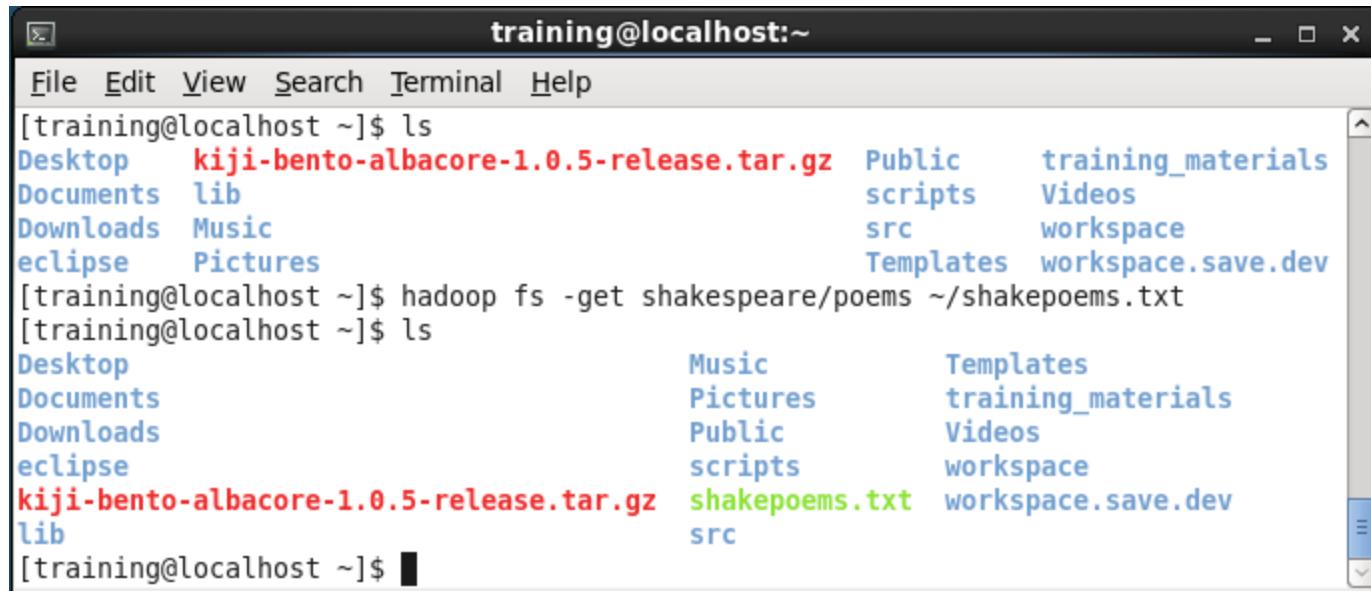
```
File Edit View Search Terminal Help
[training@localhost data]$ hadoop fs -cat shakespeare/histories | tail -n 20
The father rashly slaughter'd his own son,
The son, compell'd, been butcher to the sire:
All this divided York and Lancaster,
Divided in their dire division,
O, now, let Richmond and Elizabeth,
The true succeeders of each royal house,
By God's fair ordinance conjoin together!
And let their heirs, God, if thy will be so.
Enrich the time to come with smooth-faced peace,
With smiling plenty and fair prosperous days!
Abate the edge of traitors, gracious Lord,
That would reduce these bloody days again,
And make poor England weep in streams of blood!
Let them not live to taste this land's increase
That would with treason wound this fair land's peace!
Now civil wounds are stopp'd, peace lives again:
That she may long live here, God say amen!

[Exeunt]
[END]
[training@localhost data]$ █
```

A green box containing the number "9" is located in the bottom right corner of the slide.

# HDFS Lab (6)

- Enter the following commands to retrieve files from HDFS and list them in the Linux OS:
  - \$ cd ~
  - \$ ls
  - \$ hadoop fs -get **shakespeare/poems** ~/shakepoems.txt
  - \$ ls



The screenshot shows a terminal window titled "training@localhost:~". The window contains the following command-line session:

```
[training@localhost ~]$ ls
Desktop      kiji-bento-albacore-1.0.5-release.tar.gz  Public      training_materials
Documents    lib                                scripts    Videos
Downloads   Music                               src        workspace
eclipse      Pictures                           Templates  workspace.save.dev
[training@localhost ~]$ hadoop fs -get shakespeare/poems ~/shakepoems.txt
[training@localhost ~]$ ls
Desktop          Music          Templates
Documents        Pictures        training_materials
Downloads        Public         Videos
eclipse          scripts        workspace
kiji-bento-albacore-1.0.5-release.tar.gz shakepoems.txt workspace.save.dev
lib              src
[training@localhost ~]$
```

The terminal window has a standard Linux-style interface with a title bar, menu bar, and scroll bars.

# Sqoop Lab (0)

- SQL + Hadoop → Sqoop
- Use Sqoop to share between a RDBMS and Hadoop



*"Finance here - we're not sure  
about this Hadoop thing...  
Could you just dump it all  
into Excel for us?"*

# Sqoop Lab (1)

- Use this Sqoop command to familiarize yourself with Sqoop:
  - `$ sqoop help`
- Use these sqoop commands to list mysql databases on localhost and then tables in one of those databases:
  - `$ sqoop list-databases --connect jdbc:mysql://localhost --username training --password training`
  - `$ sqoop list-tables --connect jdbc:mysql://localhost/movielens --username training --password training`

# Sqoop Lab (2)

- Use this Sqoop command to import the movie table from mysql to an HDFS file. The fields of a table record (row) will be separated by tabs in the corresponding record of the HDFS file:
  - `$ sqoop import --connect jdbc:mysql://localhost/movielens --username training --password training --fields-terminated-by '\t' --table movie`
- Note that Sqoop constructs some SQL statements for MySQL and that it converts the SQL statement into a MapReduce job that only has a Map component:
  - Test SQL: `SELECT t.* FROM `movie` AS t LIMIT 1`
  - `SELECT MIN(`id`), MAX(`id`) FROM `movie``
- Use this HDFS command to list the newly imported file(s) and then view the last part of one of the files:
  - `$ hadoop fs -ls movie`
  - `$ hadoop fs -tail movie/part-m-00000`

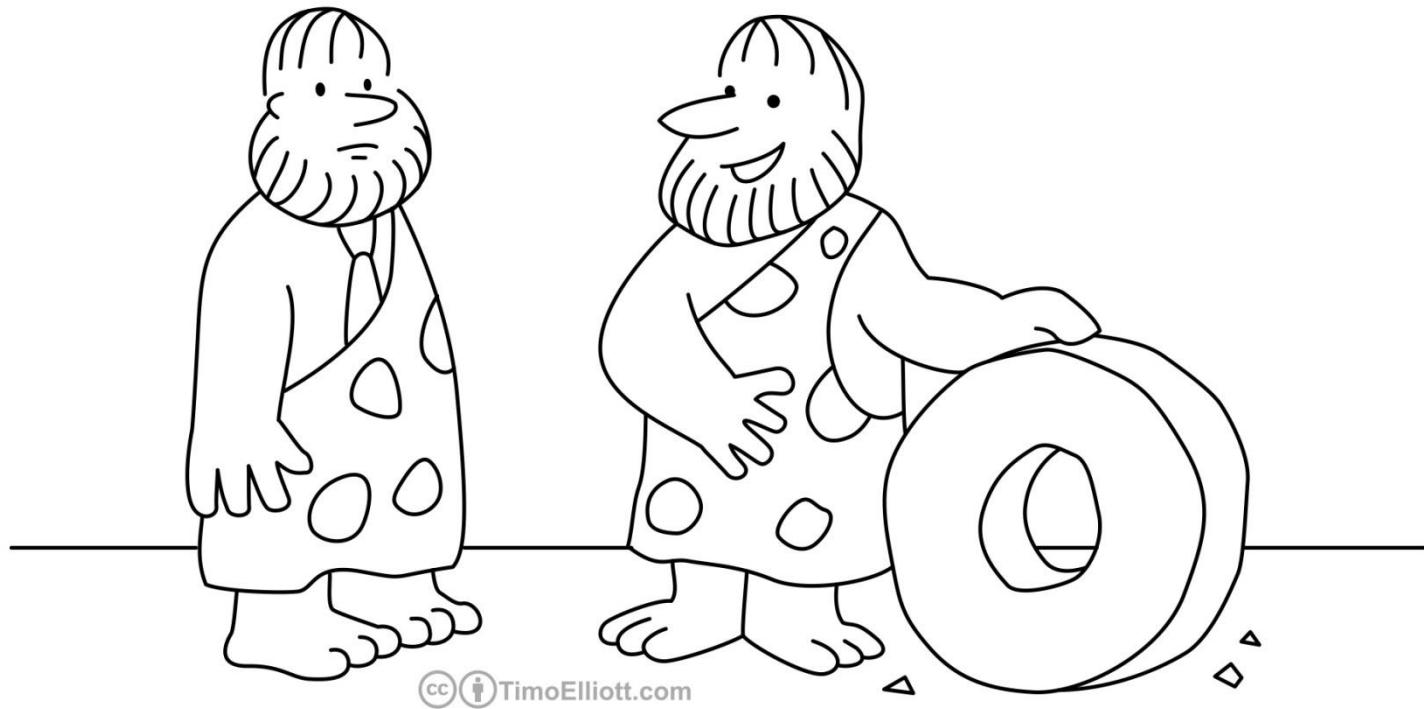
# Sqoop Lab (3)

- Use this Sqoop command to import the movierating table from mysql to an HDFS file. The fields of a table record (row) will be separated by tabs in the corresponding record of the HDFS file:
  - `$ sqoop import --connect jdbc:mysql://localhost/movielens --username training --password training --fields-terminated-by '\t' --table movierating`
- Use these HDFS commands to list the newly imported file(s) and then view the last part of one of the files:
  - `$ hadoop fs -ls movierating`
  - `$ hadoop fs -tail movierating/part-m-00000`

# Break



# Hive Lab (0)



*“It does look similar—but this one  
is powered by Hadoop”*

# Hive Lab (1)

- Check that we have movie and movierating in HDFS
  - `$ hadoop fs -cat /movie/part-m-00000 | head`
  - `$ hadoop fs -cat /movierating/part-m-00000 | head`



- Start Hive
  - `$ hive`

- Tell Hive that the HDFS files, movie and movierating, are tables:
  - `hive> CREATE EXTERNAL TABLE movie (id INT, name STRING, year INT) ROW FORMAT DELIMITED FIELDS TERMINATED BY '\t' LOCATION '/user/training/movie';`
  - `hive> CREATE EXTERNAL TABLE movierating (userid INT, movieid INT, rating INT) ROW FORMAT DELIMITED FIELDS TERMINATED BY '\t' LOCATION '/user/training/movierating';`

# Hive Lab (2)

- Ask Hive to show tables and provide metadata on these tables
  - `hive> SHOW TABLES;`
  - `hive> DESCRIBE movie;`
  - `hive> DESCRIBE movierating;`
- Use Hive to find information from your tables
  - `hive> SELECT * FROM movie WHERE year < 1925;`
  - `hive> SELECT * FROM movie WHERE year < 1925 AND year != 0 ORDER BY name;`
  - `hive> SELECT * FROM movierating WHERE userid=149;`
  - `hive> SELECT AVG(rating) FROM movierating WHERE userid=149;`
  - `hive> SELECT userid, COUNT(userid), AVG(rating) FROM movierating WHERE userid < 10 GROUP BY userid;`

# Hive Lab (3)

- Use Hive to create a new table
  - `hive> CREATE TABLE MovieRated (MovieID INT, NumRatings INT, AvgRating FLOAT);`
  - `hive> insert overwrite table MovieRated SELECT movieid,COUNT(movieid),AVG(rating) FROM movierating GROUP BY movieid;`
- Query the new table
  - `hive> SELECT COUNT(*) FROM MovieRated;`
  - `hive> SELECT * FROM MovieRated WHERE AvgRating > 4.4 and NumRatings > 1000;`

# Hive Lab (4)

- Query with Join
  - `hive> select name, NumRatings , AvgRating from MovieRated join movie on MovieRated.movieid=movie.id where AvgRating > 4.4 and NumRatings > 500 ORDER BY AvgRating;`
- Quit Hive
  - `hive> QUIT;`
- Take-home Exercise
  - What is the name of the oldest movie in the database that has a top rating?

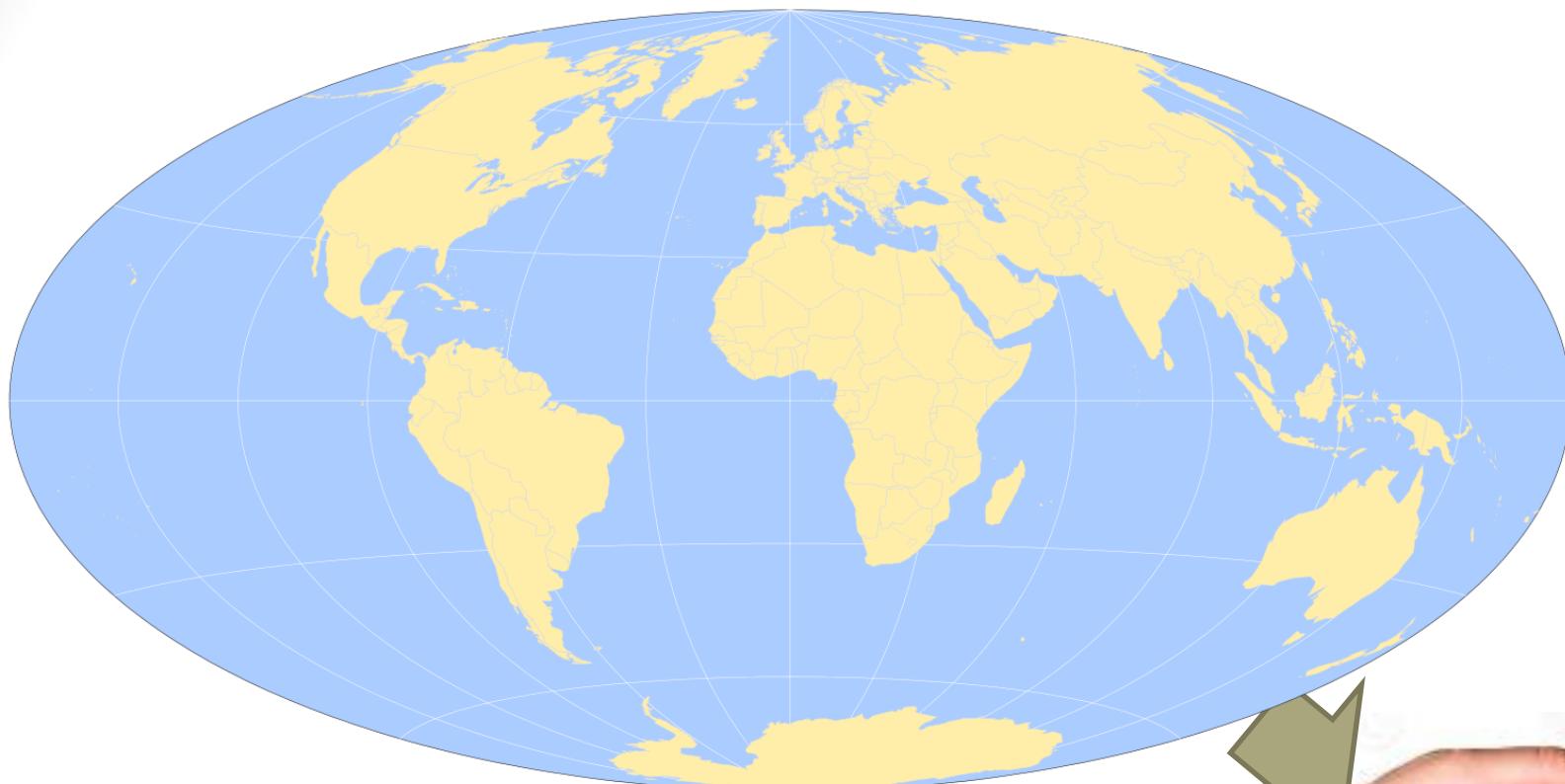
# Hive Lab (5) Impala

- Use Impala to execute “HiveQL”

1. Open Firefox
2. Start Hue by clicking on the Hue link in the Bookmarks toolbar (Username: training, Password: training)
3. Select Impala and Query Editor Tab
4. Enter Query into Query text box for Impala:
  - `select name, NumRatings , AvgRating from MovieRated join movie on MovieRated.movieid=movie.id where NumRatings > 500 ORDER BY AvgRating DESC limit 20;`
5. Click on Execute



# MapReduce (0)



# MapReduce (1)

- The purpose of Hadoop is to facilitate scale-out computing.
- MapReduce is Hadoop's answer to scale-out processing (Like HDFS is Hadoop's answer to scale-out persistence).
- The ideas behind distributed processing:
  - Send the program to the data as opposed to the data to the program
  - Make use of distributed data where each chunk of data already has its own CPUs
  - Allow independent, asynchronous processing of data
- MapReduce is a pattern (programming model) designed to make coding of distributed processes easy.

# MapReduce (2)

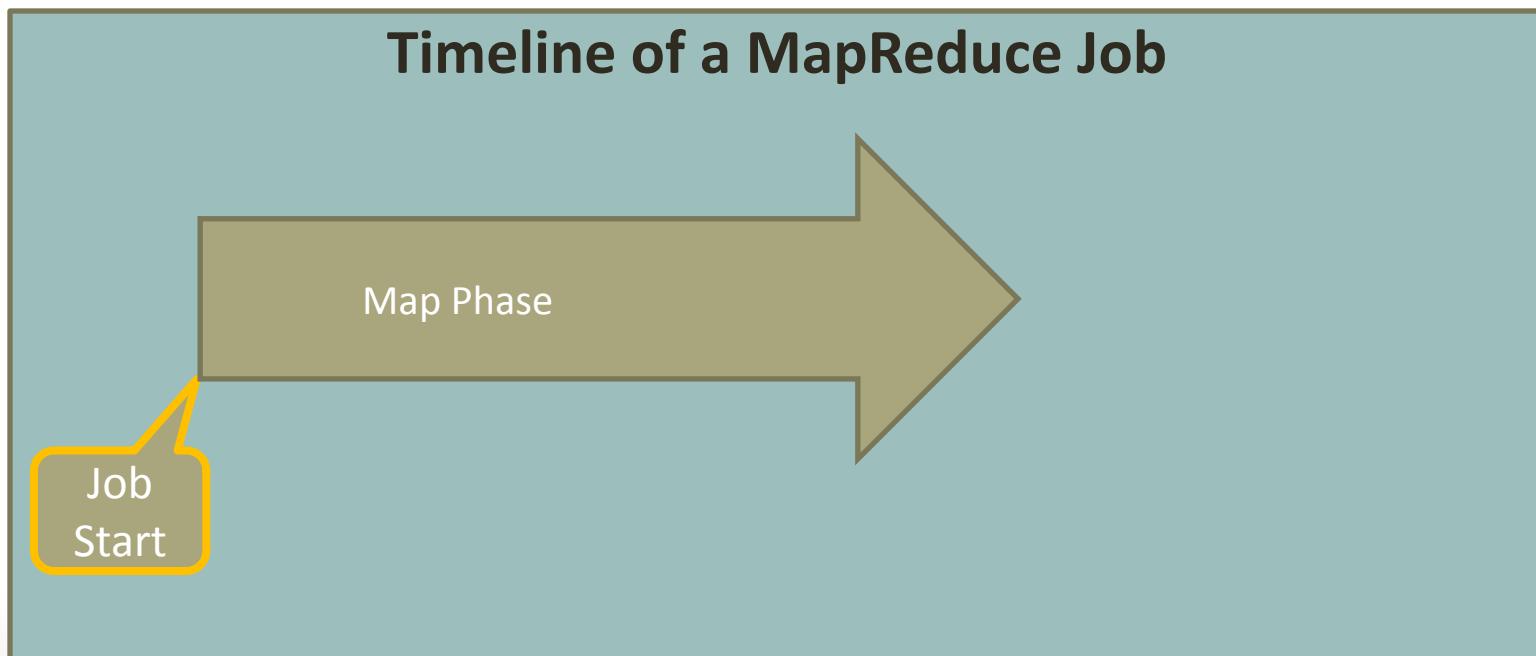
- The word MapReduce is made from **Map** and **Reduce**
- When a MapReduce job is started, then Hadoop sends **Map** and **Reduce** jobs to the data nodes. Hadoop manages all the details of data passing among data nodes.
  - The **Map** portion of the computation occurs on the individual data nodes where the data already reside. There is no need for the data to travel.
  - The **Reduce** portion of the computation occurs on some of the data nodes (not necessarily where the data reside). In some cases part of the **Reduce** operation can occur on the data node where the data already reside.

# MapReduce (3)

- MapReduce uses three stages:
  1. Map
  2. Shuffle
  3. Reduce
- Then, why don't we call MapReduce: MapShuffleReduce?
  - Because: In most cases, the programmer need only be concerned with Map and Reduce. The Shuffle and Sorting is taken care of by Hadoop.
  - To achieve scalable programs, Hadoop takes care off:
    - Creating tasks on the various data nodes
    - Tracking Tasks
    - Moving data among data nodes
    - File I/O, networking, process synchronization, recovery from failures, re-running jobs, splitting input, Moving Key-Value-Pairs from Map to Reduce
    - Outputs results

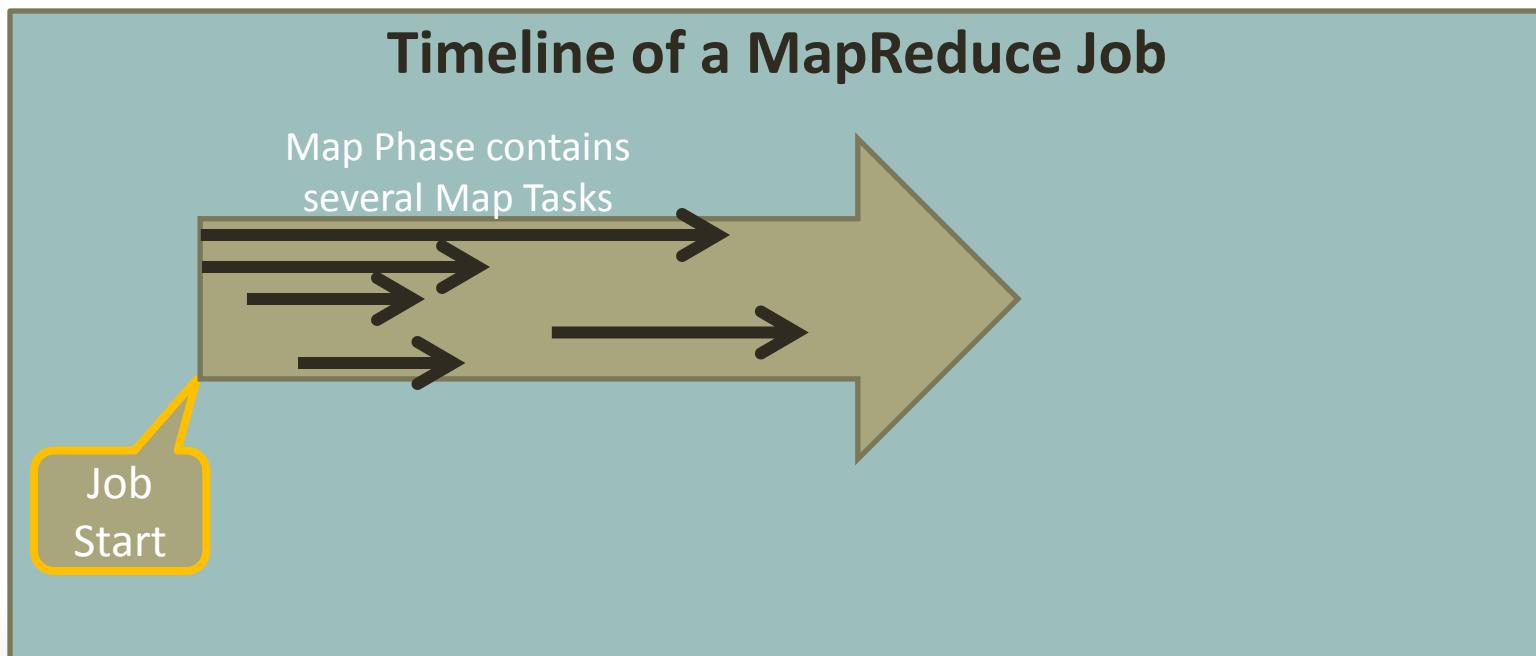
# MapReduce (4)

- Parallelization Rules in MapReduce



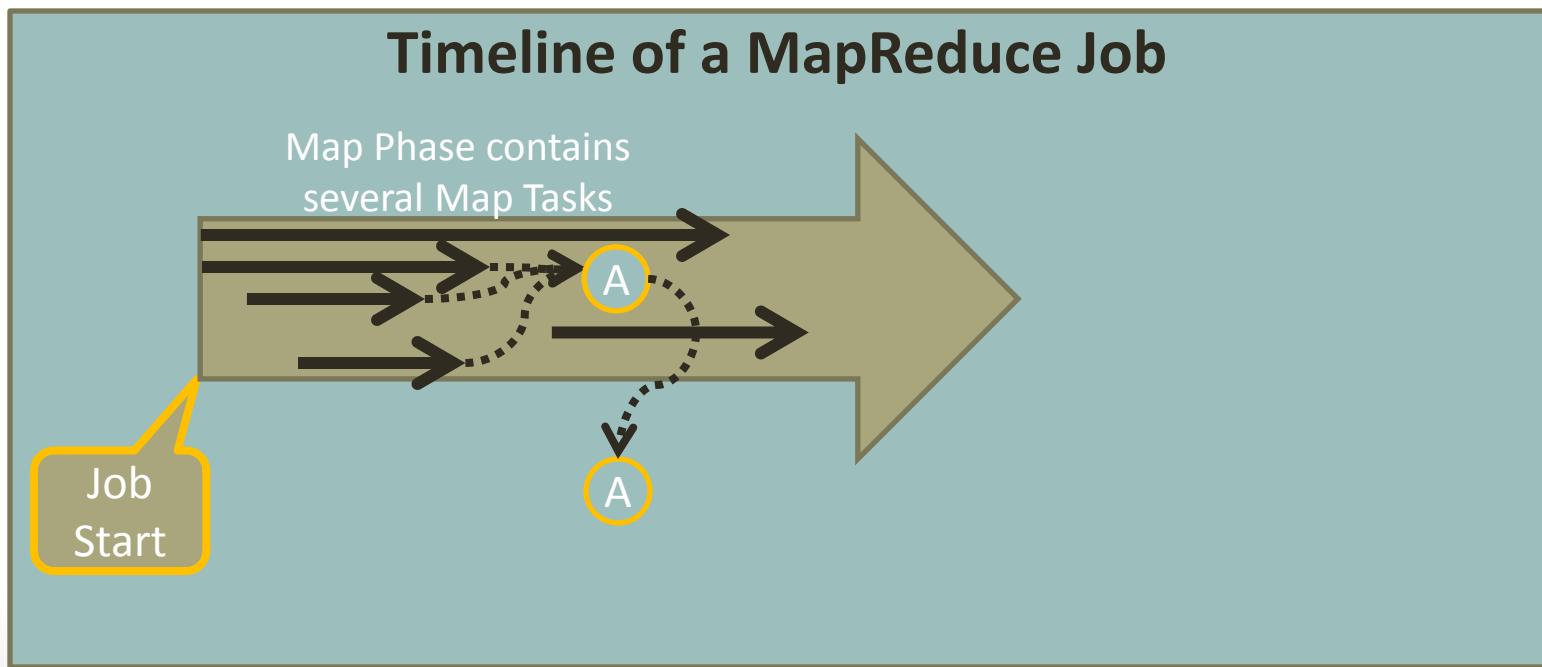
# MapReduce (5)

- Parallelization Rules in MapReduce
  - Map tasks occur in parallel independent of each other



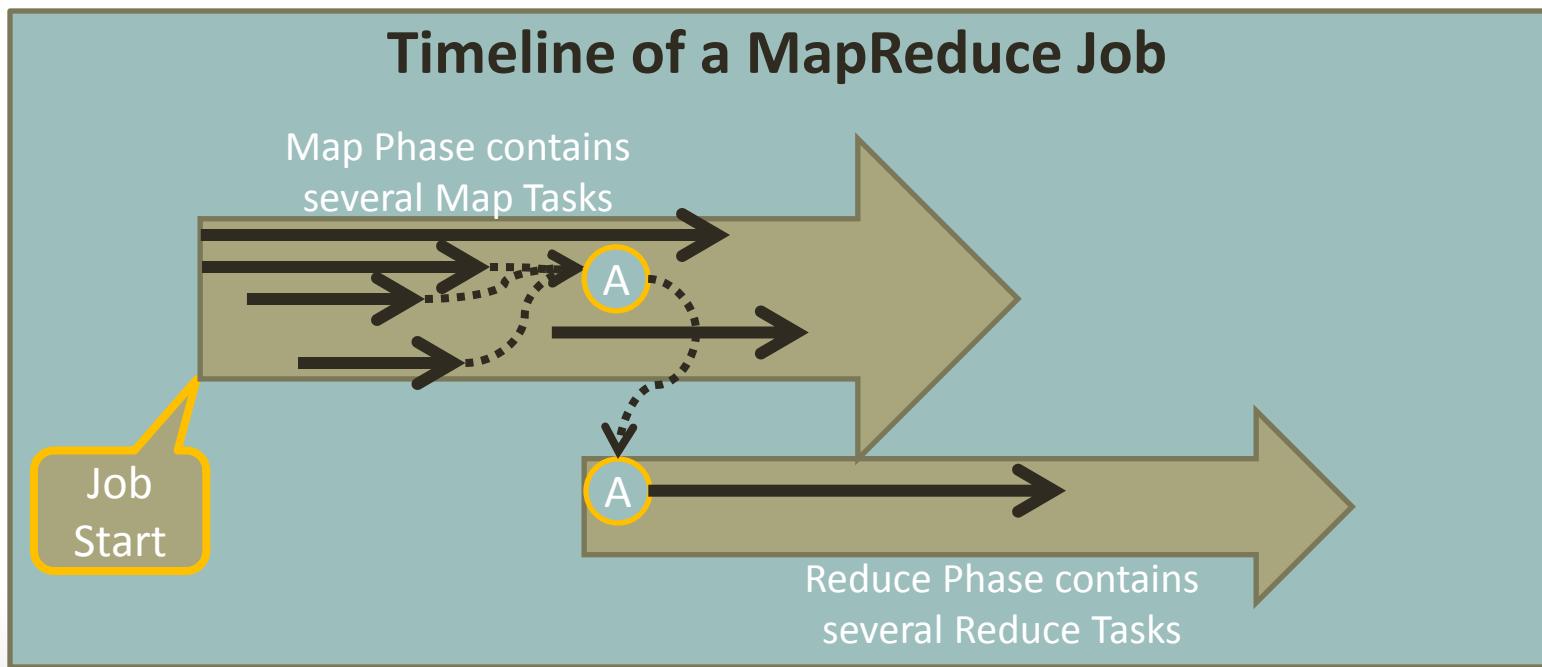
# MapReduce (6)

- Parallelization Rules in MapReduce
  - Map tasks occur in parallel independent of each other
  - Map tasks terminate in “Shuffle and Sort”



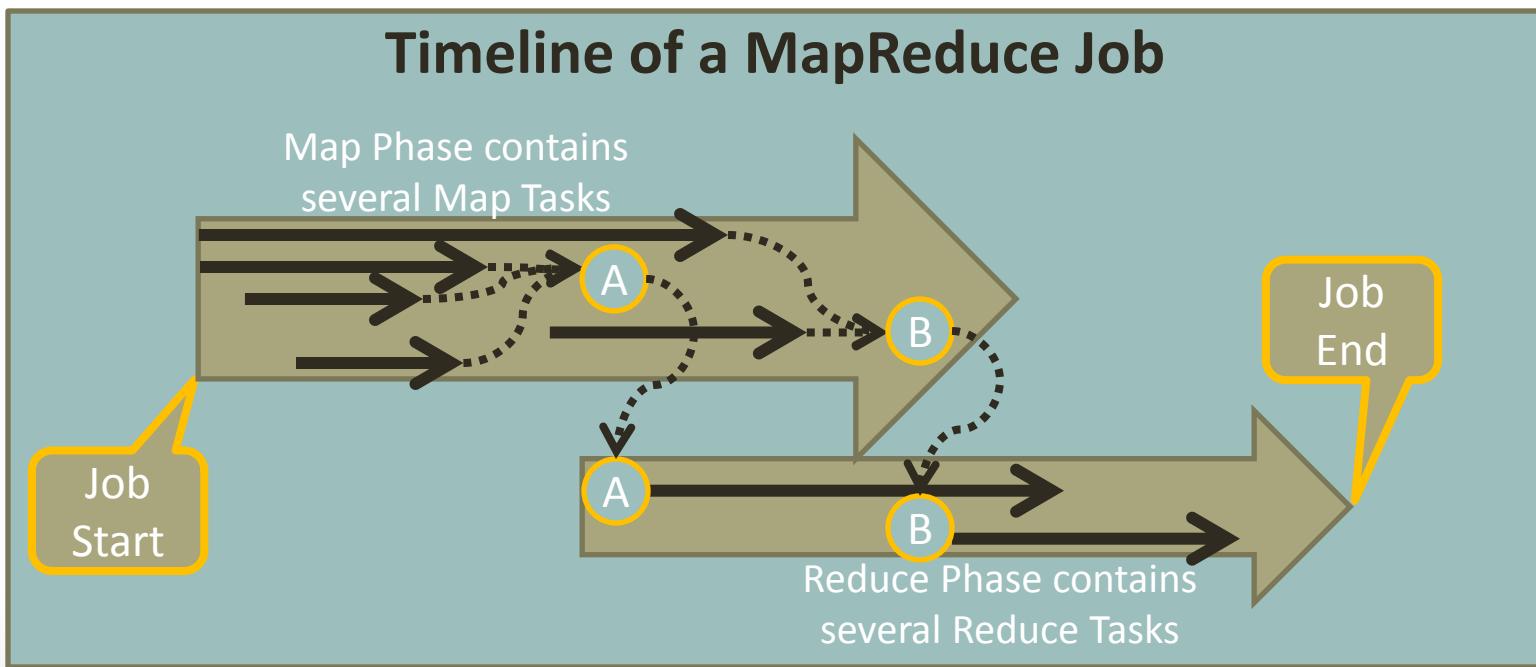
# MapReduce (7)

- Parallelization Rules in MapReduce
  - Map tasks occur in parallel independent of each other
  - Map tasks terminate in “Shuffle and Sort”
  - Some Reduce tasks can start before all Mappers are done.



# MapReduce (8)

- Parallelization Rules in MapReduce
  - Map tasks occur in parallel independent of each other
  - Map tasks terminate in “Shuffle and Sort”
  - Some Reduce tasks can start before all Mappers are done.
  - Reduce tasks occur in parallel independent of each other
  - Reducers cannot complete before all Mapping and Shuffle & Sort have completed.



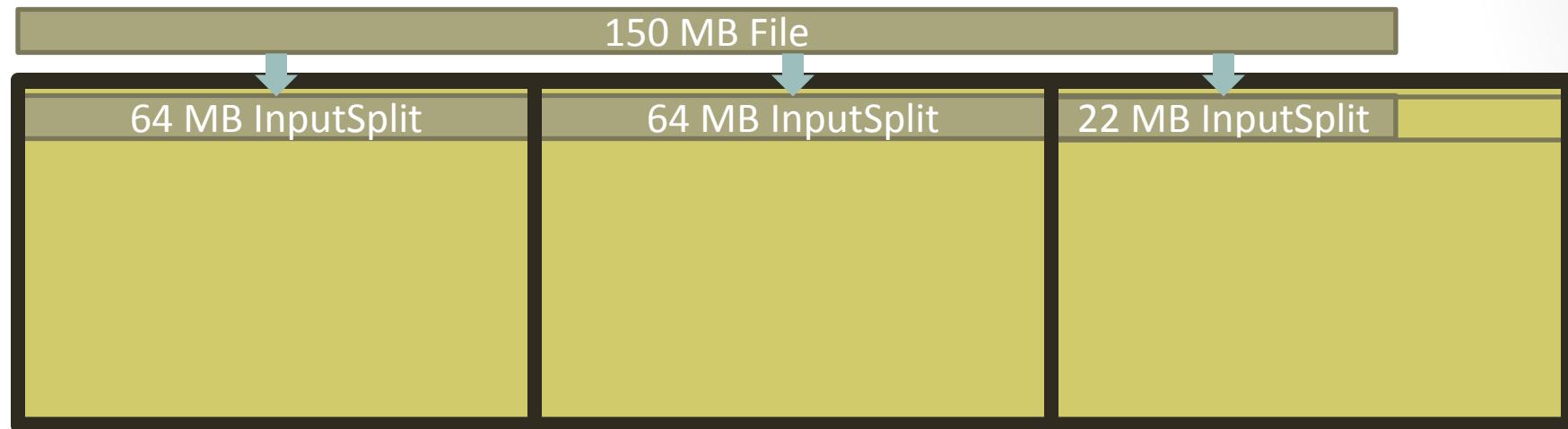
# MapReduce (9)

150 MB File

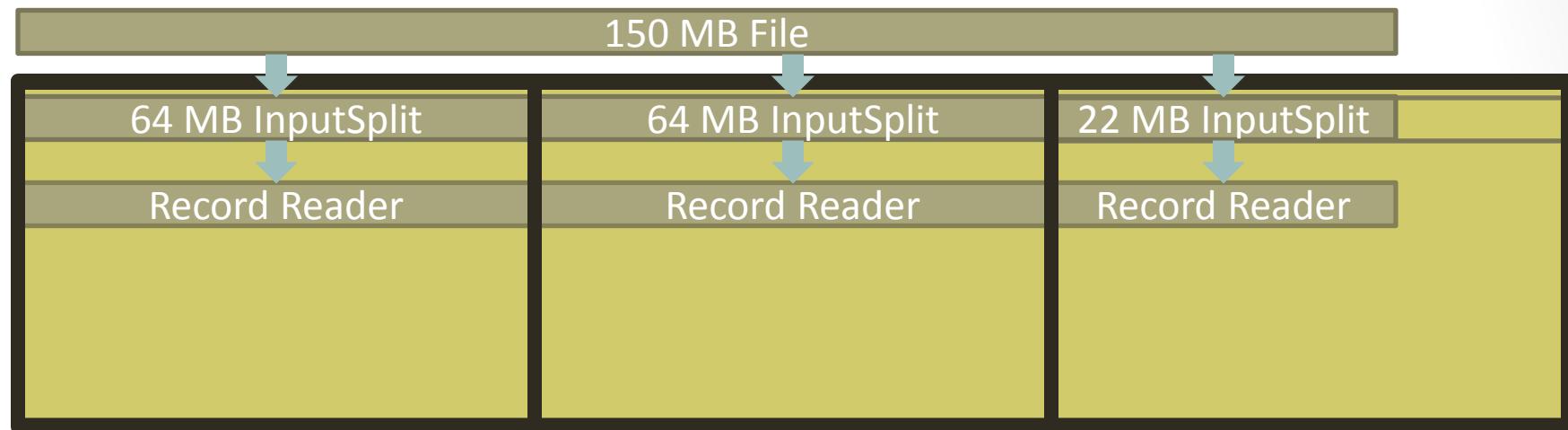
# MapReduce (10)



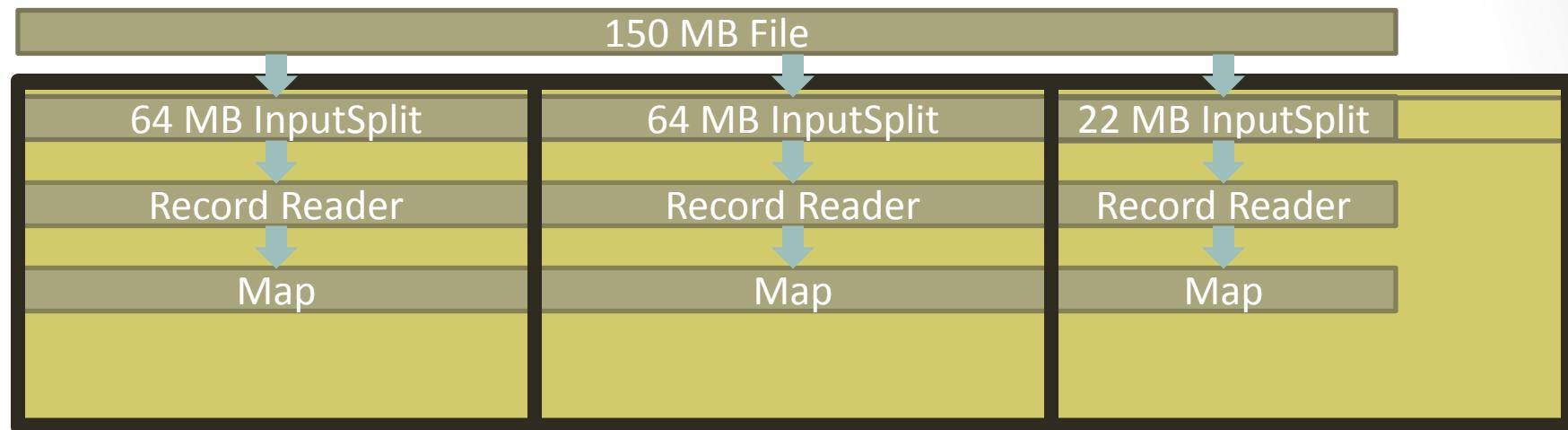
# MapReduce (11)



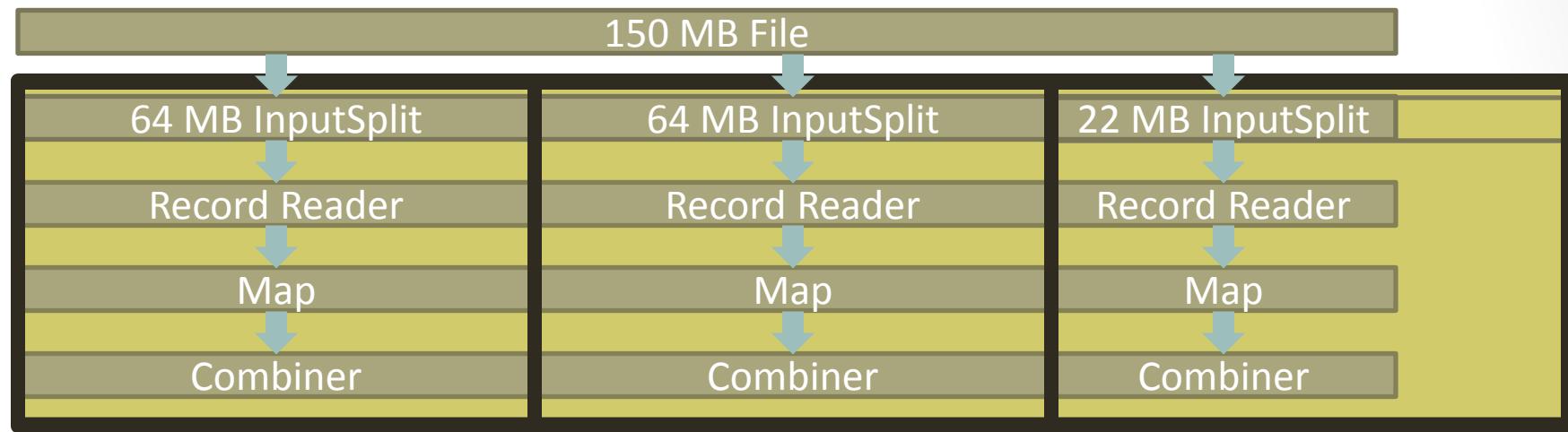
# MapReduce (12)



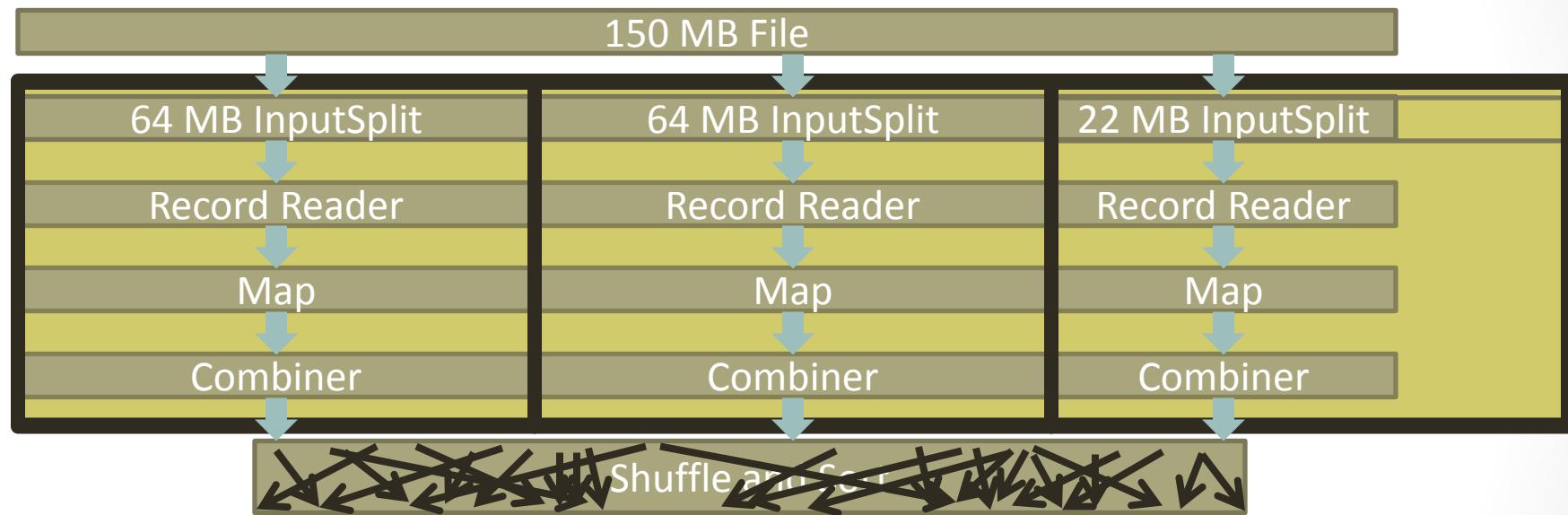
# MapReduce (13)



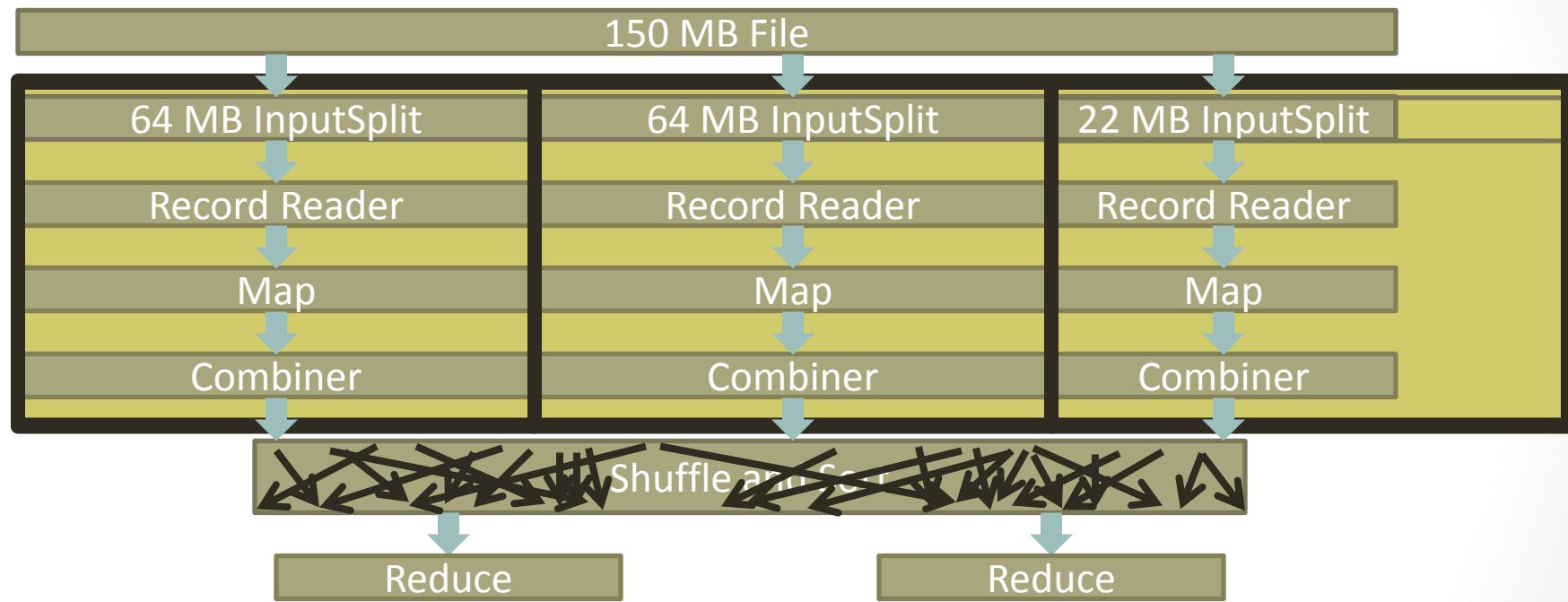
# MapReduce (14)



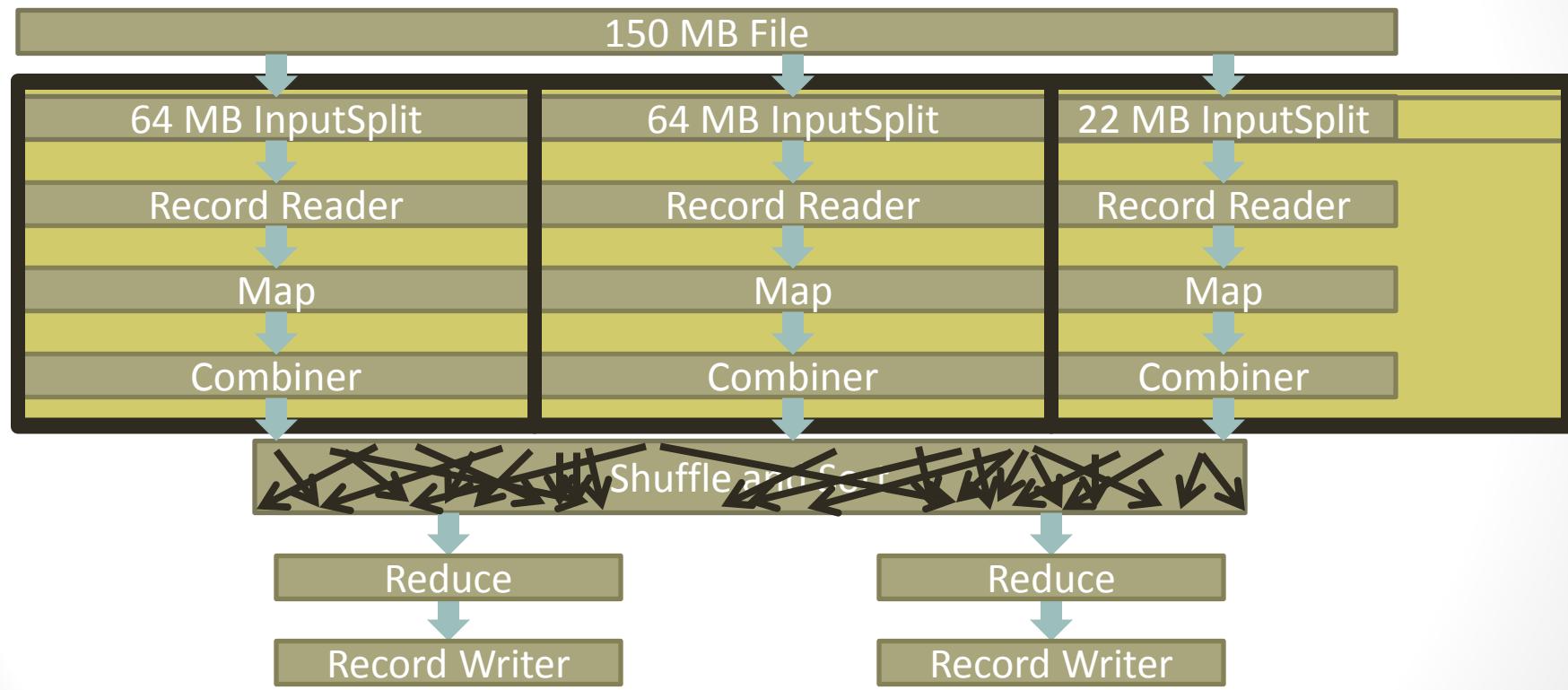
# MapReduce (15)



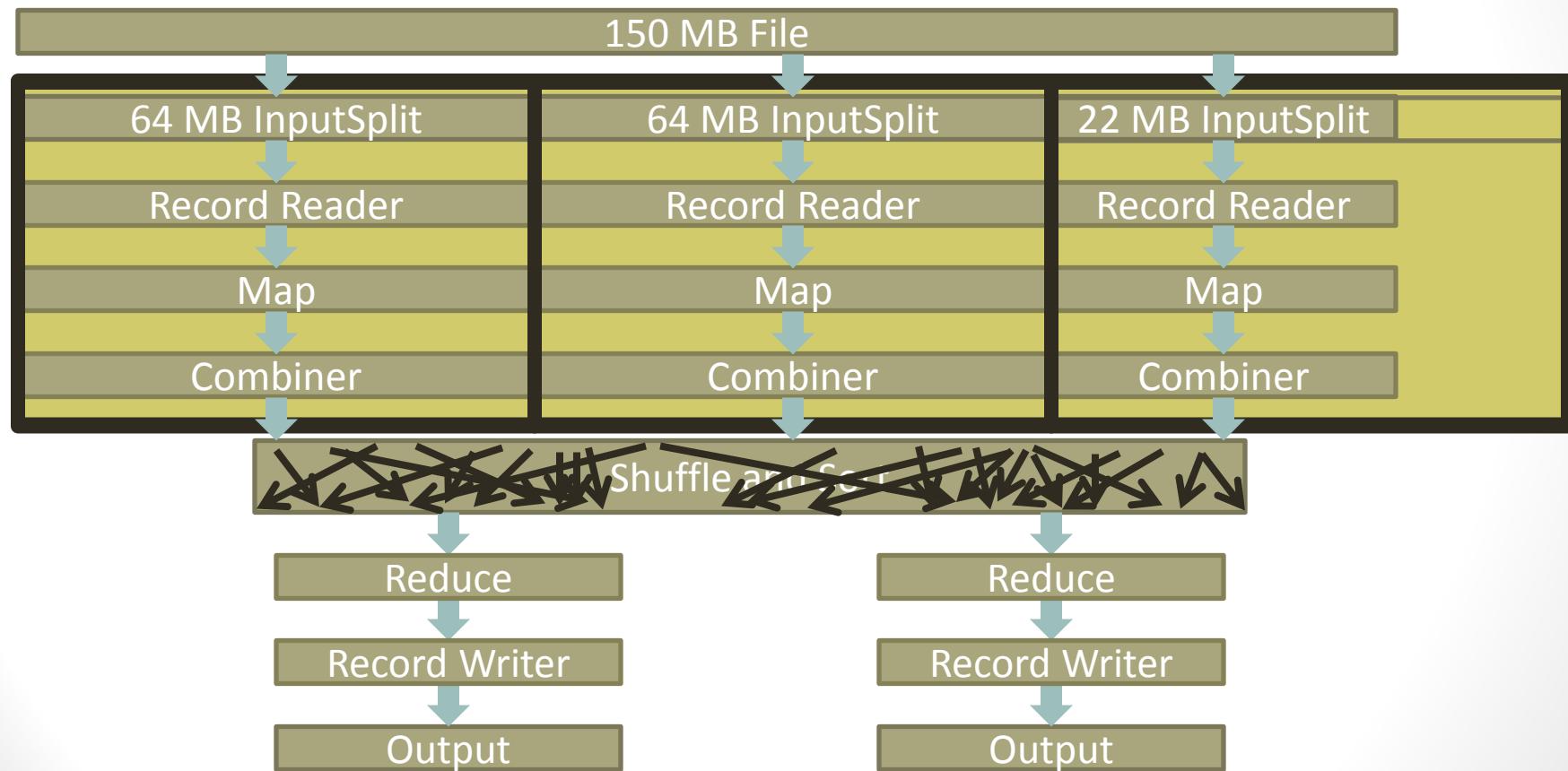
# MapReduce (16)



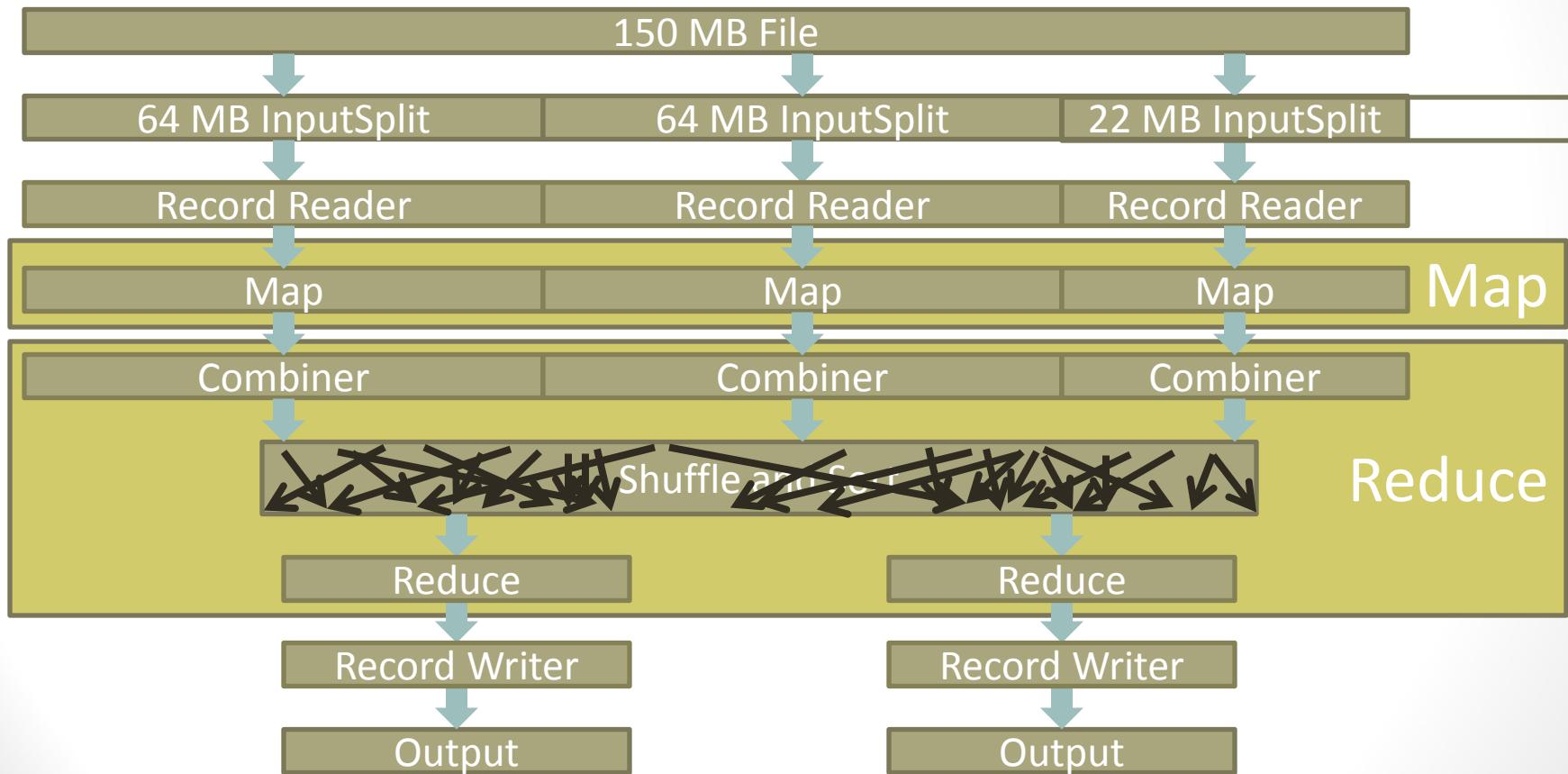
# MapReduce (17)



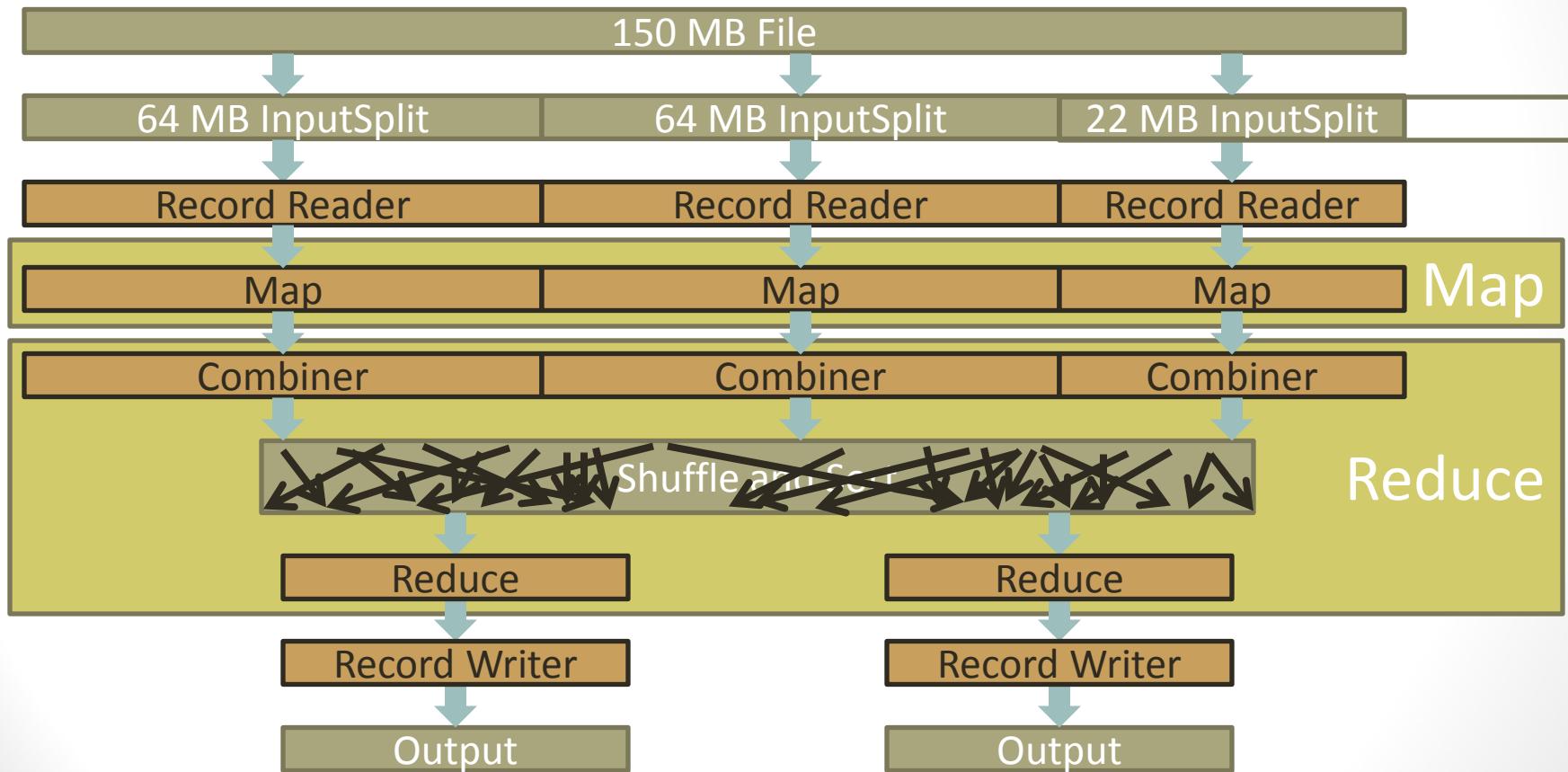
# MapReduce (18)



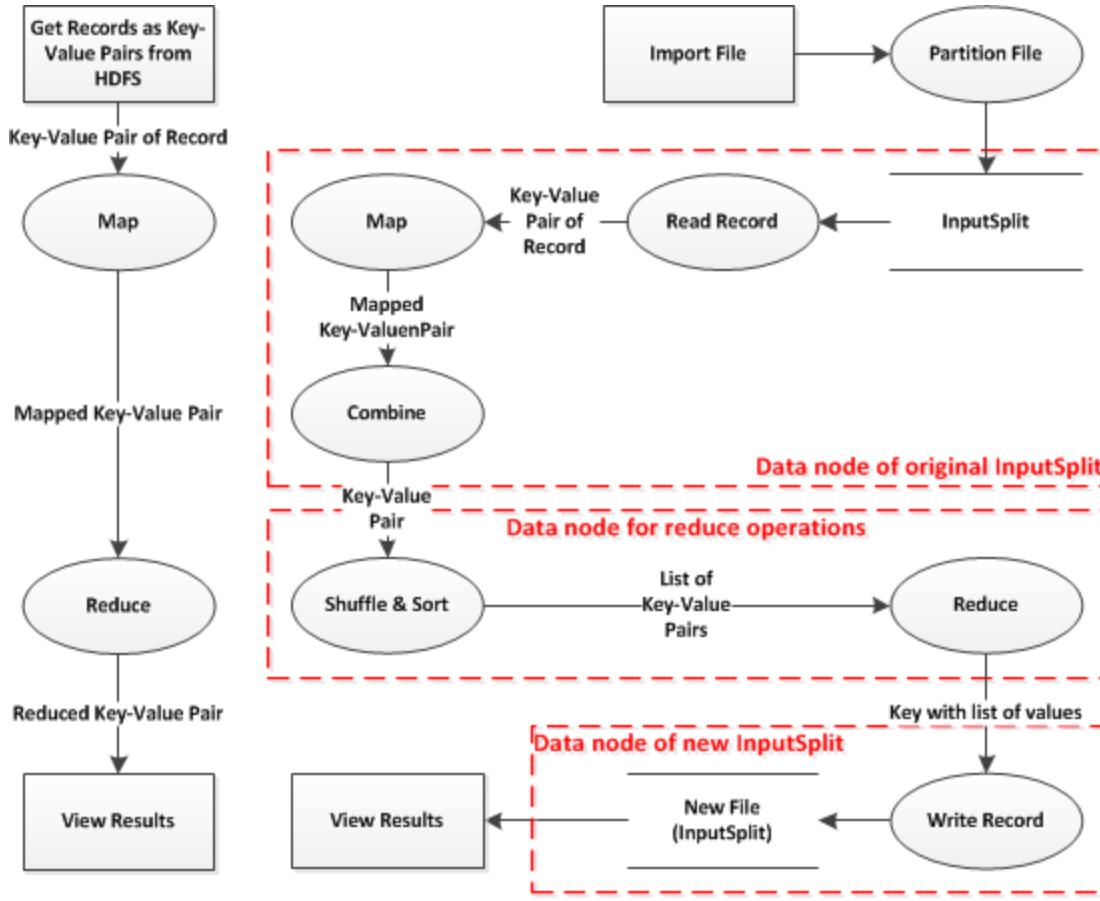
# MapReduce(19)



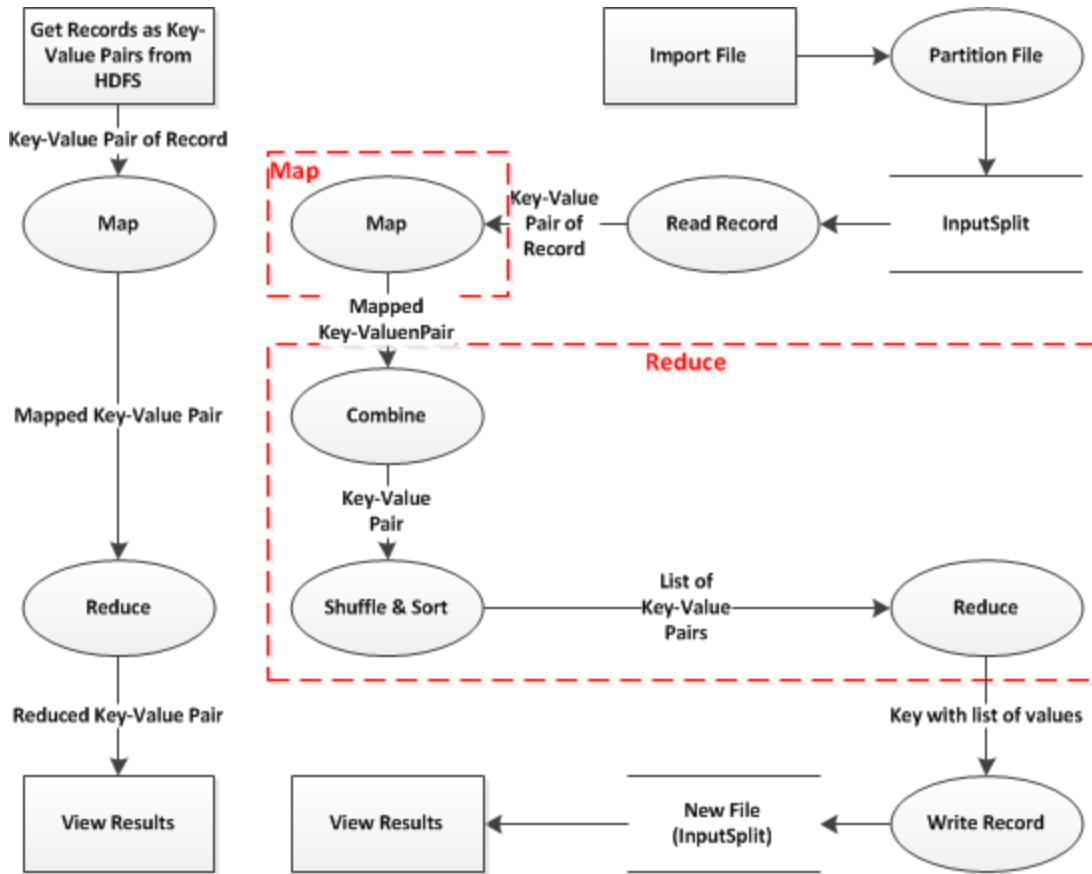
# MapReduce (20)



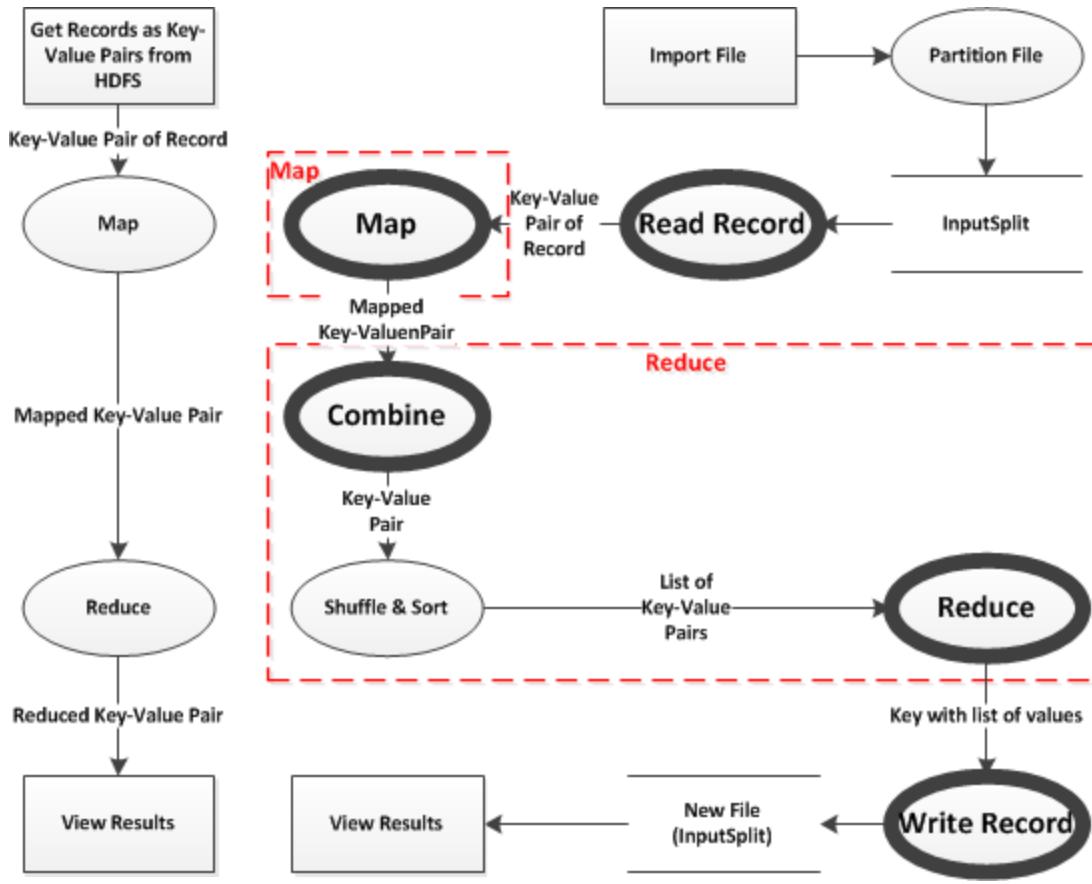
# MapReduce (21)



# MapReduce (22)



# MapReduce (23)

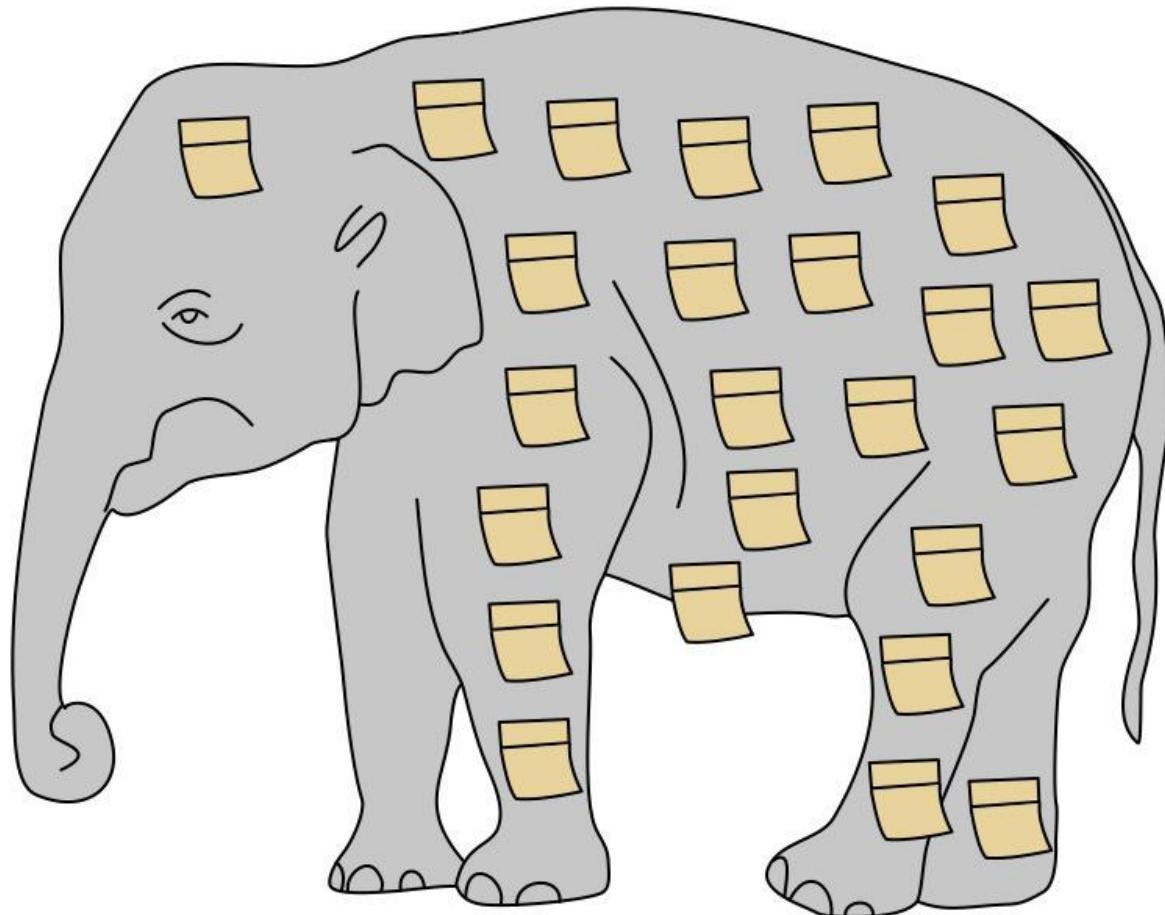


# MapReduce (24)

- Record Read
  - $\text{Record} \rightarrow (\text{K1}, \text{V1})$
  - $\text{to be or not to be} \rightarrow (0, \text{"to be or not to be"})$
- Map
  - $(\text{K1}, \text{V1}) \rightarrow \text{list}(\text{K2}, \text{V2})$
  - $(0, \text{"to be or not to be"}) \rightarrow [(\text{to}, 1), (\text{be}, 1), (\text{or}, 1), (\text{not}, 1), (\text{to}, 1), (\text{be}, 1)]$
- Shuffle and Sort
  - $\text{list}(\text{K2}, \text{V2}) \rightarrow (\text{K2}, \text{list}(\text{V2}))$
  - $[(\text{to}, 1), (\text{be}, 1), (\text{or}, 1), (\text{not}, 1), (\text{to}, 1), (\text{be}, 1)] \rightarrow (\text{to}, [1, 1]), (\text{be}, [1, 1]), (\text{or}, [1]), (\text{not}, [1])$
- Reduce
  - $(\text{K2}, \text{list}(\text{V2})) \rightarrow \text{list}(\text{K3}, \text{V3})$
  - $(\text{to}, [1, 1]), (\text{be}, [1, 1]), (\text{or}, [1]), (\text{not}, [1]) \rightarrow [(\text{to}, 2), (\text{be}, 2), (\text{or}, 1), (\text{not}, 1)]$
- Record Write
  - $(\text{K3}, \text{V3}) \rightarrow \text{Records}$
  - $[(\text{to}, 2), (\text{be}, 2), (\text{or}, 1), (\text{not}, 1)] \rightarrow$ 

to,2
be,2
or,1
not,1

# Break



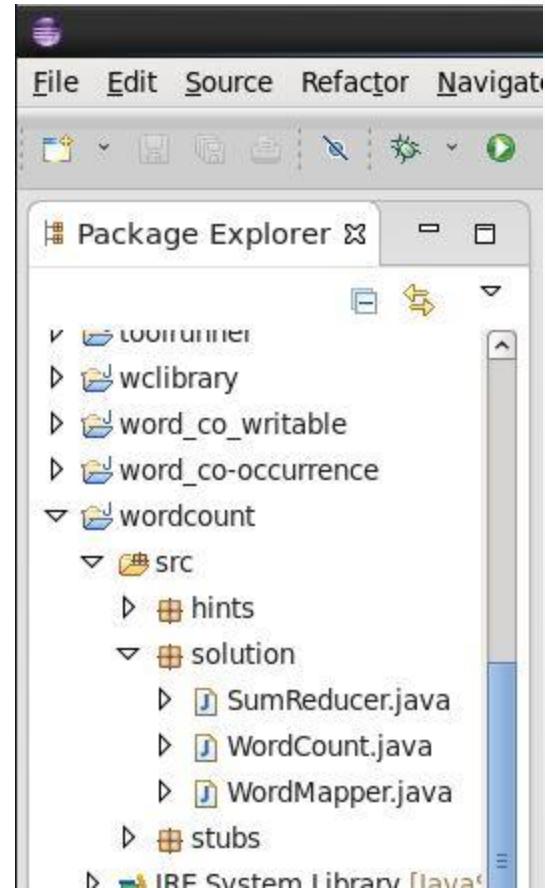
TimoElliott.com

*An early Hadoop prototype...*

# MapReduce Lab (1)



- Open Eclipse by clicking on the Eclipse icon. Then, use the Package Explorer to navigate to wordcount.
- See the java files in wordcount\src\solution\
  - WordCount.java
  - WordMapper.java
  - SumReducer.java
- Study these java files



# MapReduce Lab (2)

**Input File**

ByteOffset

Record

**Mapper Class**  
extends Mapper<  
MapperInputKeyType,  
MapperInputValueType,  
MapperEmitKeyType,  
MapperEmitValueType>

map(  
MapperInputKeyType  
MapperInputKey,  
  
MapperInputValueType  
MapperInputValue,  
Context context)

    context.write(  
        MapperEmitKey,  
        MapperEmitValue);

**Reducer Class**  
extends Reducer<  
MapperEmitKeyType,  
MapperEmitValueType,  
OutputKeyType,  
OutputValueType>

reduce(  
MapperEmitKeyType  
MapperEmitKey,  
Iterable<  
MapperEmitValueType>  
MapperEmitValue,  
Context context)

    context.write(  
        OutputKey,  
        OutputValue);

**main**  
job.setOutputKeyClass(  
    OutputKeyType.class);  
job.setOutputValueClass(  
    OutputValueType.class);

**Output File**  
OutputKey  
OutputValue

# MapReduce Lab (3)

**Mapper input key is typically just the file byte offset. Ignore it.**

**Mapper Class**  
extends Mapper<  
**MapperInputKeyType**,  
MapperInputValueType,  
MapperEmitKeyType,  
MapperEmitValueType>

```
map(  
  MapperInputKeyType  
  MapperInputKey,  
  
  MapperInputValueType  
  MapperInputValue,  
  Context context)  
  
  context.write(  
    MapperEmitKey,  
    MapperEmitValue);
```

**Reducer Class**  
extends Reducer<  
MapperEmitKeyType,  
MapperEmitValueType,  
OutputKeyType,  
OutputValueType>

```
reduce(  
  MapperEmitKeyType  
  MapperEmitKey,  
  Iterable<  
    MapperEmitValueType>  
  MapperEmitValue,  
  Context context)
```

**main**  
job.setOutputKeyClass(  
 OutputKeyType.class);  
job.setOutputValueClass(  
 OutputValueType.class);

**Input File  
ByteOffset**

Record

```
map(  
  MapperInputKeyType  
  MapperInputKey,  
  
  MapperInputValueType  
  MapperInputValue,  
  Context context)
```

```
  context.write(  
    MapperEmitKey,  
    MapperEmitValue);
```

**Output File**  
OutputKey  
OutputValue

# MapReduce Lab (4)

Mapper input value is typically a file record. It is of type Text.

**Mapper Class**  
extends Mapper<  
MapperInputKeyType,  
**MapperInputValueType**,  
MapperEmitKeyType,  
MapperEmitValueType>

map(  
MapperInputKeyType  
MapperInputKey,  
**MapperInputValueType**  
**MapperInputValue**,  
Context context)

context.write(  
MapperEmitKey,  
MapperEmitValue);

**Reducer Class**  
extends Reducer<  
MapperEmitKeyType,  
MapperEmitValueType,  
OutputKeyType,  
OutputValueType>

reduce(  
MapperEmitKeyType  
MapperEmitKey,  
Iterable<  
MapperEmitValueType>  
MapperEmitValue,  
Context context)

context.write(  
OutputKey,  
OutputValue);

**main**  
job.setOutputKeyClass(  
OutputKeyType.class);  
job.setOutputValueClass(  
OutputValueType.class);

**Input File**  
ByteOffset

Record

[ 52 ]

# MapReduce Lab (5)

Mapper output or emit key is the same the reducer input key.

**Mapper Class**  
extends Mapper<  
MapperInputKeyType,  
MapperInputValueType,  
**MapperEmitKeyType**,  
MapperEmitValueType>

map(  
MapperInputKeyType  
MapperInputKey,  
  
MapperInputValueType  
MapperInputValue,  
Context context)  
  
context.write(  
**MapperEmitKey**,  
MapperEmitValue);

**Reducer Class**  
extends Reducer<  
**MapperEmitKeyType**,  
MapperEmitValueType,  
OutputKeyType,  
OutputValueType>

reduce(  
**MapperEmitKeyType**  
**MapperEmitKey**,  
Iterable<  
MapperEmitValueType>  
MapperEmitValue,  
Context context)  
  
context.write(  
OutputKey,  
OutputValue);

**main**  
job.setOutputKeyClass(  
OutputKeyType.class);  
job.setOutputValueClass(  
OutputValueType.class);

**Input File**  
ByteOffset

Record

**Output File**  
OutputKey  
OutputValue

# MapReduce Lab (6)

**Mapper output or emit value is the reducer input value.**

**Mapper Class**  
extends Mapper<  
MapperInputKeyType,  
MapperInputValueType,  
MapperEmitKeyType,  
**MapperEmitValueType**>

map(  
MapperInputKeyType  
MapperInputKey,

MapperInputValueType  
MapperInputValue,  
Context context)

```
context.write(  
MapperEmitKey,  
MapperEmitValue)  
;
```

**Reducer Class**  
extends Reducer<  
MapperEmitKeyType,  
**MapperEmitValueType**,  
OutputKeyType,  
OutputValueType>

reduce(  
MapperEmitKeyType  
MapperEmitKey,  
Iterable<  
**MapperEmitValueType**>  
**MapperEmitValue**,  
Context context)

```
context.write(  
OutputKey,  
OutputValue);
```

**main**  
job.setOutputKeyClass(  
OutputKeyType.class);  
job.setOutputValueClass(  
OutputValueType.class);

**Input File**  
ByteOffset

Record

**Output File**  
OutputKey  
OutputValue

# MapReduce Lab (7)

**Reducer output or emit key is the output key from main and the key of the Output File**

**Mapper Class**  
extends Mapper<  
MapperInputKeyType,  
MapperInputValueType,  
MapperEmitKeyType,  
MapperEmitValueType>

map(  
MapperInputKeyType  
MapperInputKey,  
  
MapperInputValueType  
MapperInputValue,  
Context context)

context.write(  
MapperEmitKey,  
MapperEmitValue);

**Reducer Class**  
extends Reducer<  
MapperEmitKeyType,  
MapperEmitValueType,  
**OutputKeyType**,  
OutputValueType>

reduce(  
MapperEmitKeyType  
MapperEmitKey,  
Iterable<  
MapperEmitValueType>  
MapperEmitValue,  
Context context)

context.write(  
**OutputKey**,  
OutputValue);

**main**  
job.setOutputKeyClass(  
**OutputKeyType**.class);  
job.setOutputValueClass(  
OutputValueType.class);

**Input File**  
ByteOffset

Record

**Output File**  
**OutputKey**  
OutputValue

# MapReduce Lab (8)

**Reducer output or emit value is the output value from main and the value of the Output File**

**Mapper Class**  
extends Mapper<  
MapperInputKeyType,  
MapperInputValueType,  
MapperEmitKeyType,  
MapperEmitValueType>

map(  
MapperInputKeyType  
MapperInputKey,  
  
MapperInputValueType  
MapperInputValue,  
Context context)

context.write(  
MapperEmitKey,  
MapperEmitValue);

**Reducer Class**  
extends Reducer<  
MapperEmitKeyType,  
MapperEmitValueType,  
OutputKeyType,  
**OutputValueType**>

reduce(  
MapperEmitKeyType  
MapperEmitKey,  
Iterable<  
MapperEmitValueType>  
MapperEmitValue,  
Context context)

context.write(  
OutputKey,  
**OutputValue**);

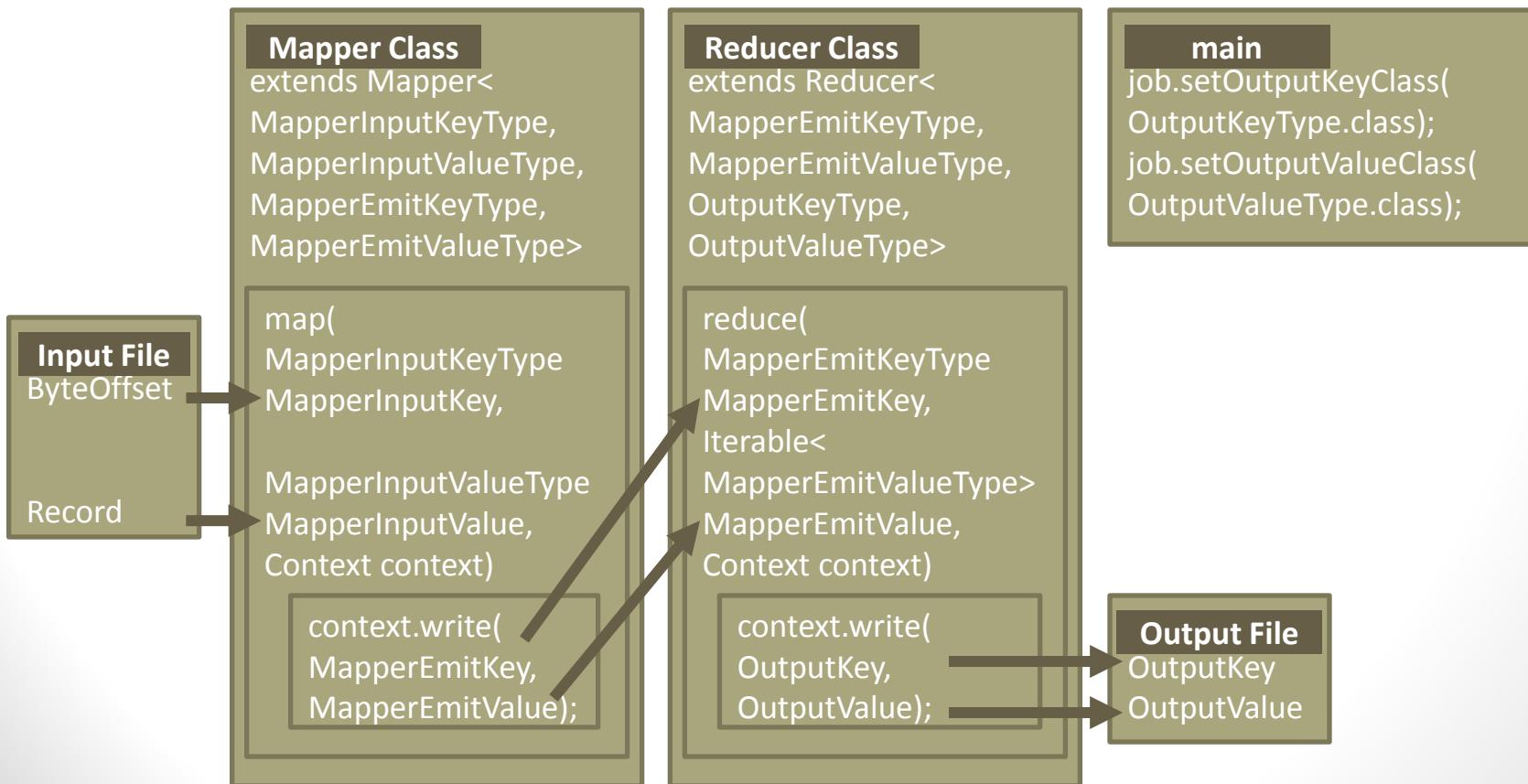
**main**  
job.setOutputKeyClass(  
OutputKeyType.class);  
job.setOutputValueClass(  
**OutputValueType**.class);

**Input File**  
ByteOffset

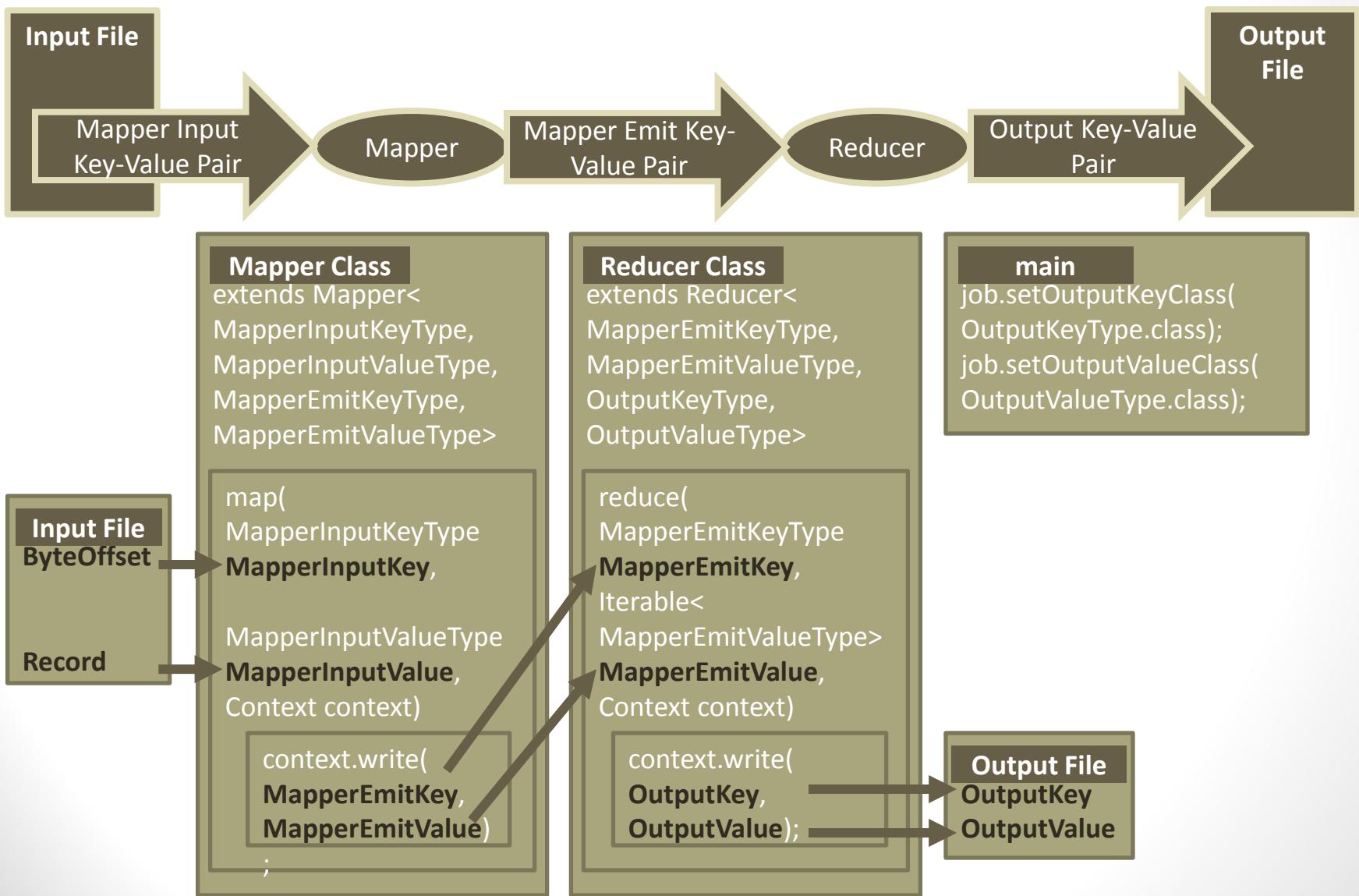
Record

**Output File**  
OutputKey  
**OutputValue**

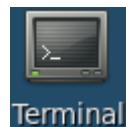
# MapReduce Lab (9)



# MapReduce Lab (10)



# MapReduce Lab (11)



Open a terminal

In the terminal change the current directory:

- `$ cd ~/workspace/wordcount/src`

List the directories and files in the current directory:

- `$ ls`

List the files in the solution directory:

- `$ ls solution`

Show the classpath to get the Hadoop Jar and other required libraries:

- `$ hadoop classpath`

Compile the java files to class files:

- `$ javac -classpath `hadoop classpath` solution/*.java`

List the files in the solution directory (note the new class files):

- `$ ls solution`

# MapReduce Lab (12)



In the terminal:

Create a JAR file containing the class files:

- `$ jar cvf wordcount.jar solution/*.class`

List the files in the current directory (note the new JAR file):

- `$ ls`

In Hadoop, run the wordcount program on all of Shakespeare's works:

- `$ hadoop jar wordcount.jar solution.WordCount shakespeare wordcounts`

In HDFS list the files that contain the results:

- `$ hadoop fs -ls wordcounts`

In HDFS list the results:

- `$ hadoop fs -cat wordcounts/part-r-00000 | less`

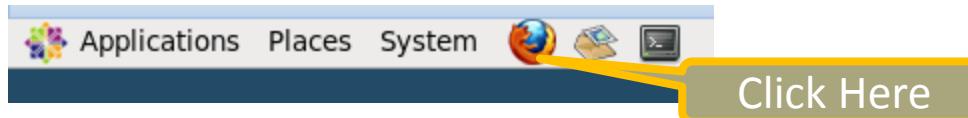
In Hadoop, run the wordcount program on Shakespeare's poems:

- `$ hadoop jar wc.jar solution.WordCount shakespeare/poems pwords`

# MapReduce Lab (13)

- Use Hue's HDFS file browser to browse and read the results

1. Open Firefox



2. Start Hue by clicking on the Hue link in the Bookmarks toolbar (Username: training, Password: training)



3. Start File Browser (Username: training, Password: training)



4. Find wordcounts/part-r-00000



5. View Contents of part-r-00000



# MapReduce Lab (14)

- The next two slides describe the assignment for this week.

# MapReduce Lab (15)

- Lab: Word Length Count 1
- This is part of today's homework assignment.
- Create a MapReduce program that lists word lengths and how often words of that length occur. Modify package wordcount.solution:
  - Calculate word length (`word.length()`) and cast to a `Text`.
  - Change map method to use the `Text` word length as the key.
  - Compile and run the program on shakespeare's works:
    - `hadoop jar wordcount.jar solution.WordCount shakespeare lengthcounts1`
  - Copy the modified java files (E.g. `MyFile.java` -> `MyFile_1.java`)
  - View the results using the file browser
    - `hadoop fs -cat lengthcounts1/part-r-00000 | less`
  - Take a screen shot of the results (approx. 18 rows). Similar to this:

1	58839
10	10342
11	3830
12	1366
13	475
14	261

# MapReduce Lab (16)

- Lab: Word Length Count 2
- This is part of today's homework assignment.
- Create a MapReduce program that lists word lengths and how often words of that length occur. Modify package wordcount.solution:
  - Calculate word length (`word.length()`) and cast to IntWritable
  - Change map method to use the word length as the key
  - Compile and run the program on Shakespeare's works:
    - `hadoop jar wordcount.jar solution.WordCount shakespeare lengthcounts2`
  - Copy the modified java files (E.g. `MyFile.java` -> `MyFile_2.java`)
  - View the results using the file browser
    - `hadoop fs -cat lengthcounts2/part-r-00000 | less`
  - Take a screen shot of the results (approx. 18 rows). Similar to this:

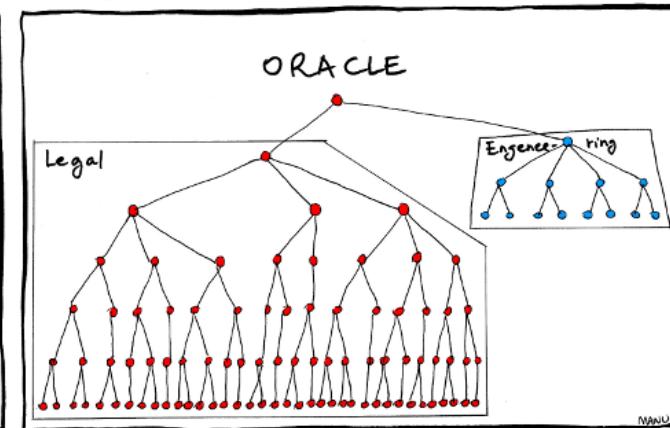
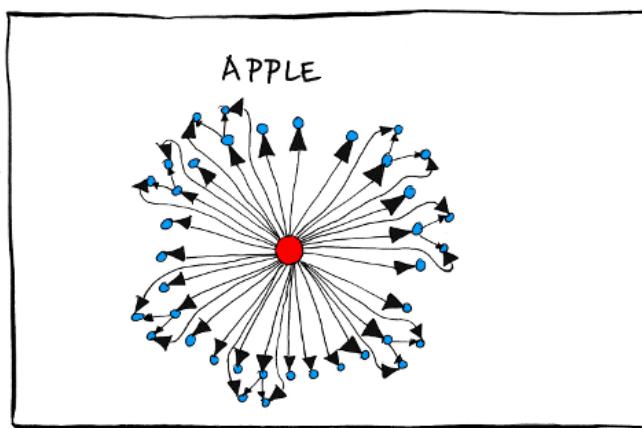
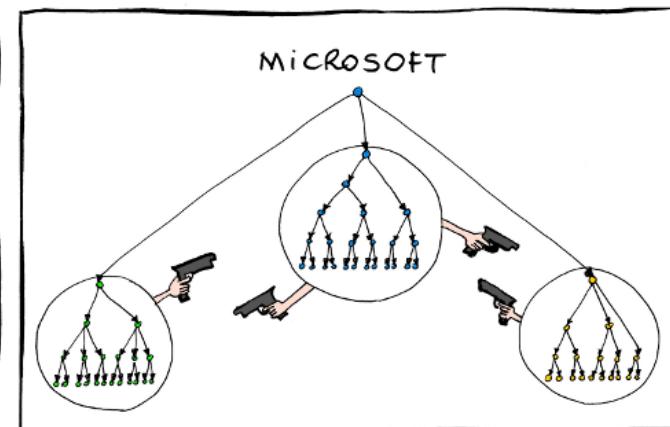
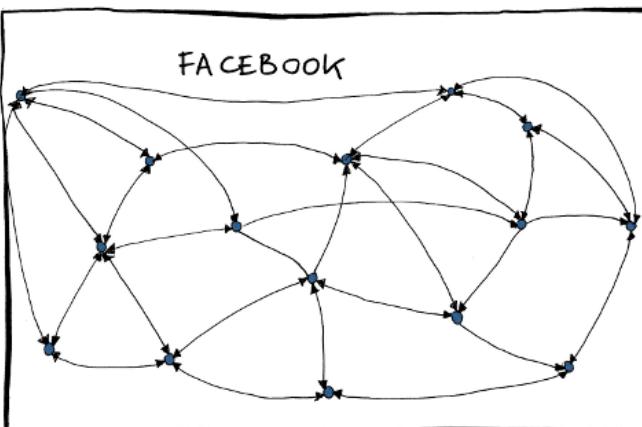
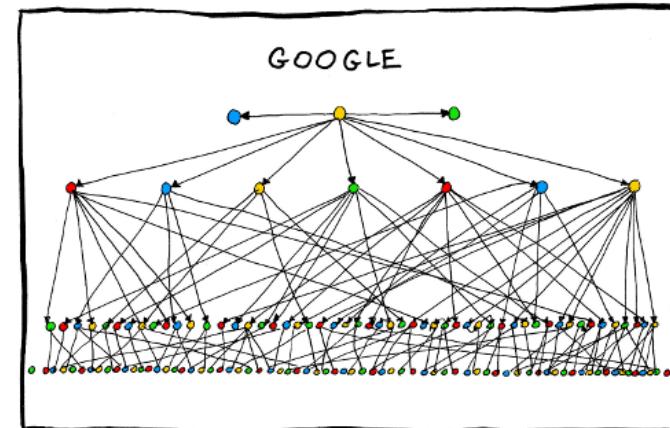
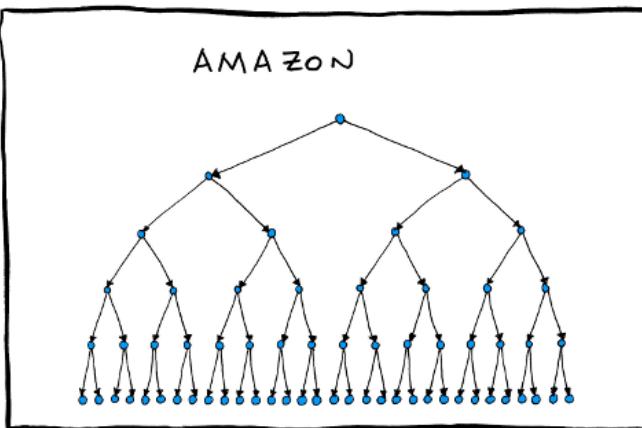
1	58839
2	163405
3	199082
4	219685
5	120151
6	79584



# Hadoop-2 Labs

[ 66 ]

# Break



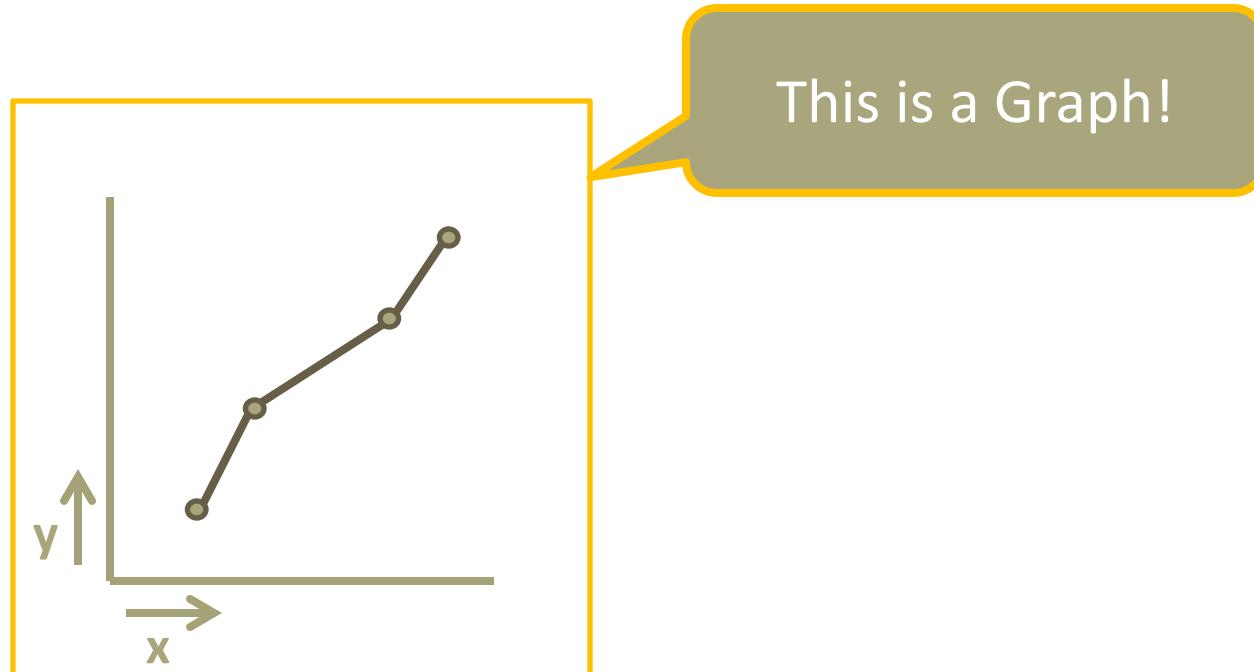
# Introduction to Graph Data

( 68 )

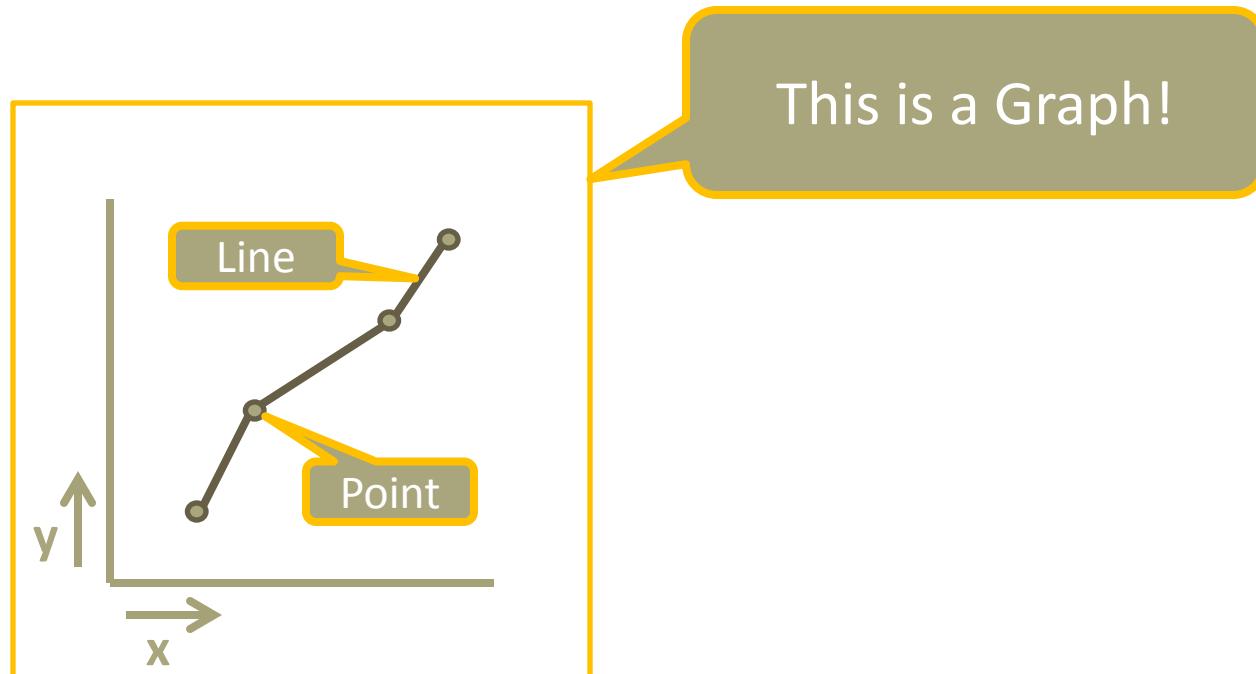
# Graph Data: What is it? (0)

What is a Graph?

# Graph Data: What is it? (1)

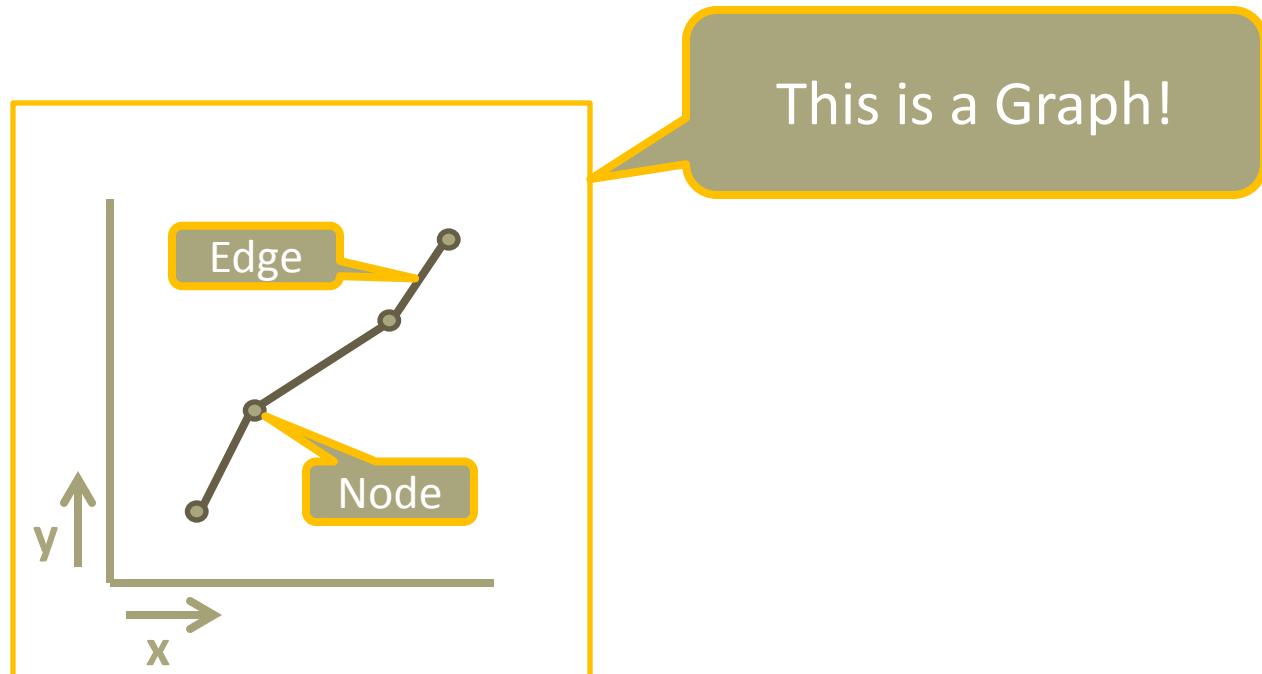


# Graph Data: What is it? (2)

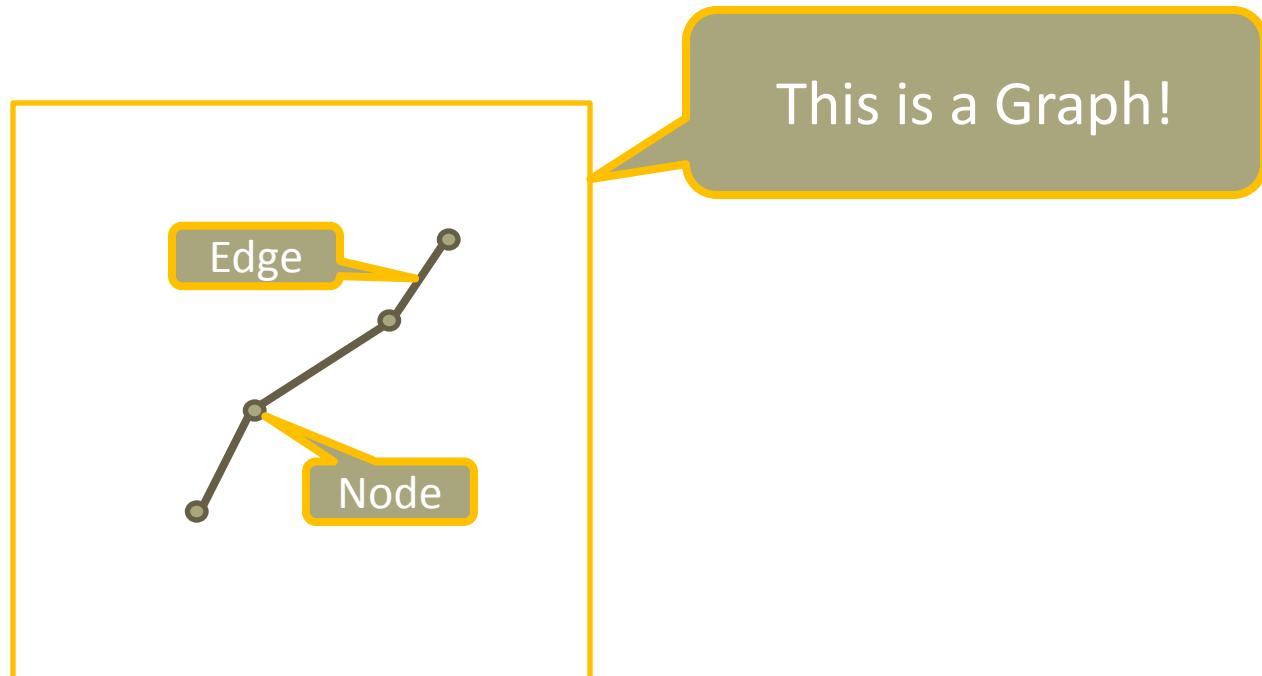


This is a Graph!

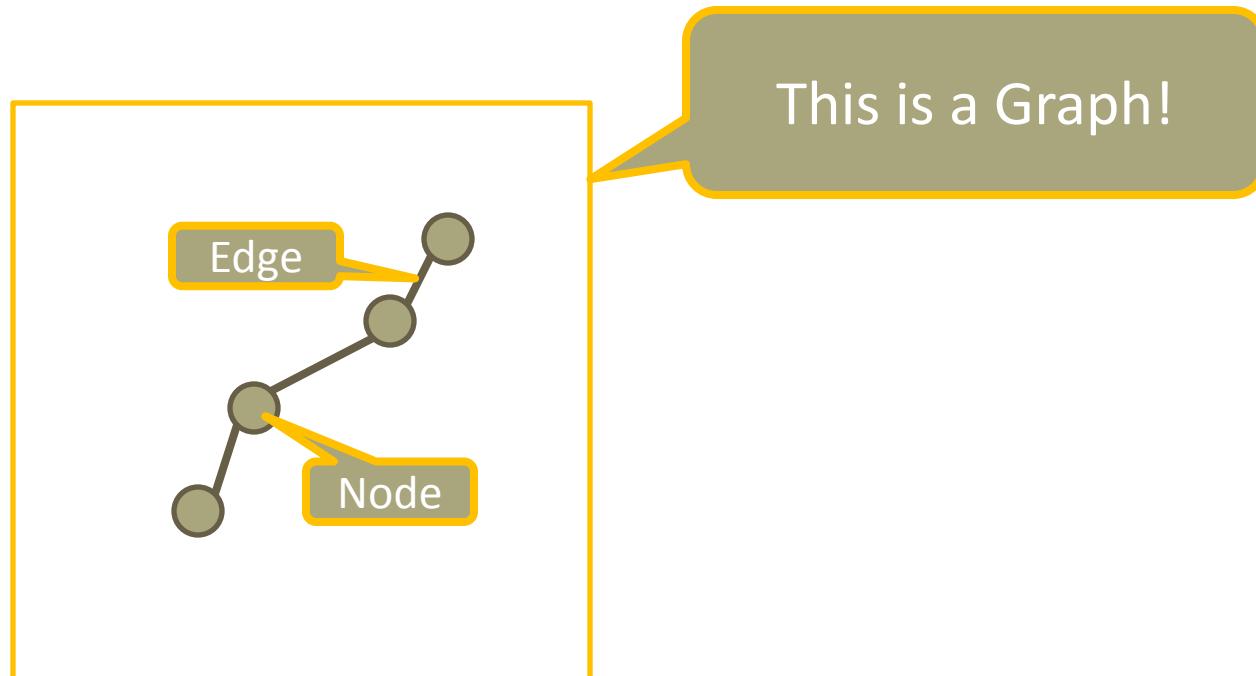
# Graph Data: What is it? (3)



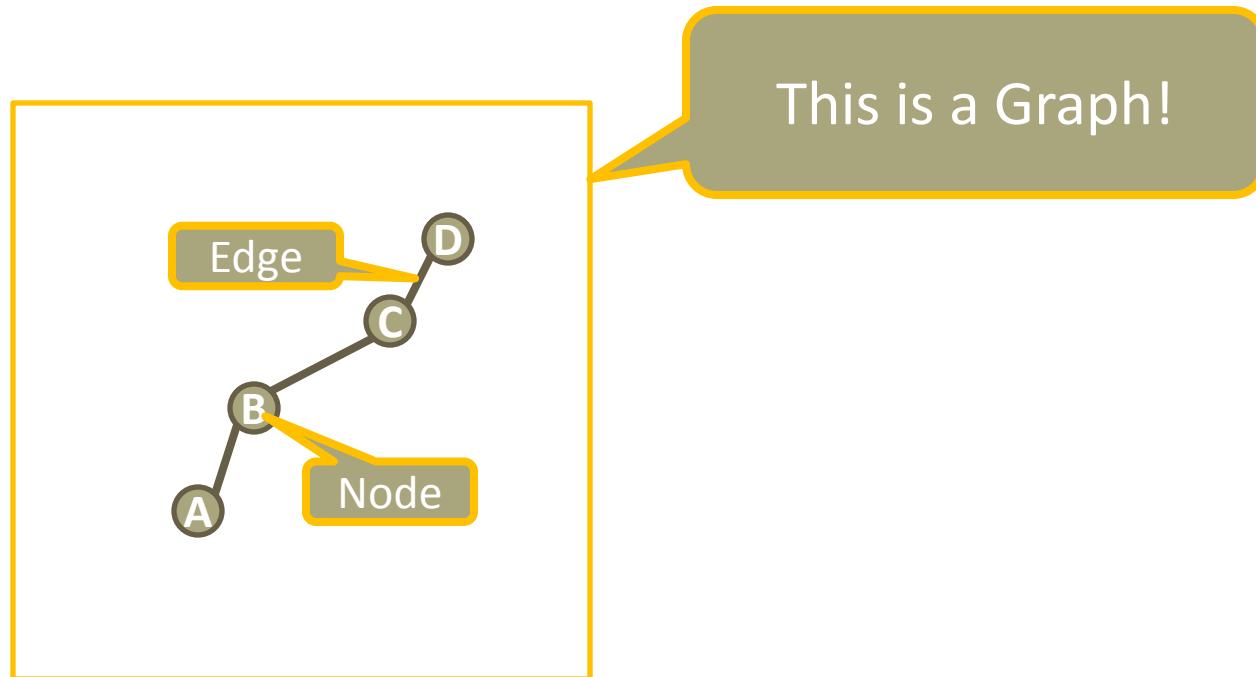
# Graph Data: What is it? (4)



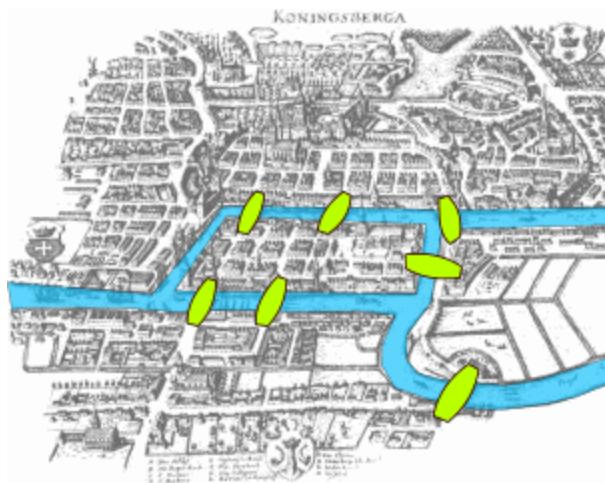
# Graph Data: What is it? (5)



# Graph Data: What is it? (6)



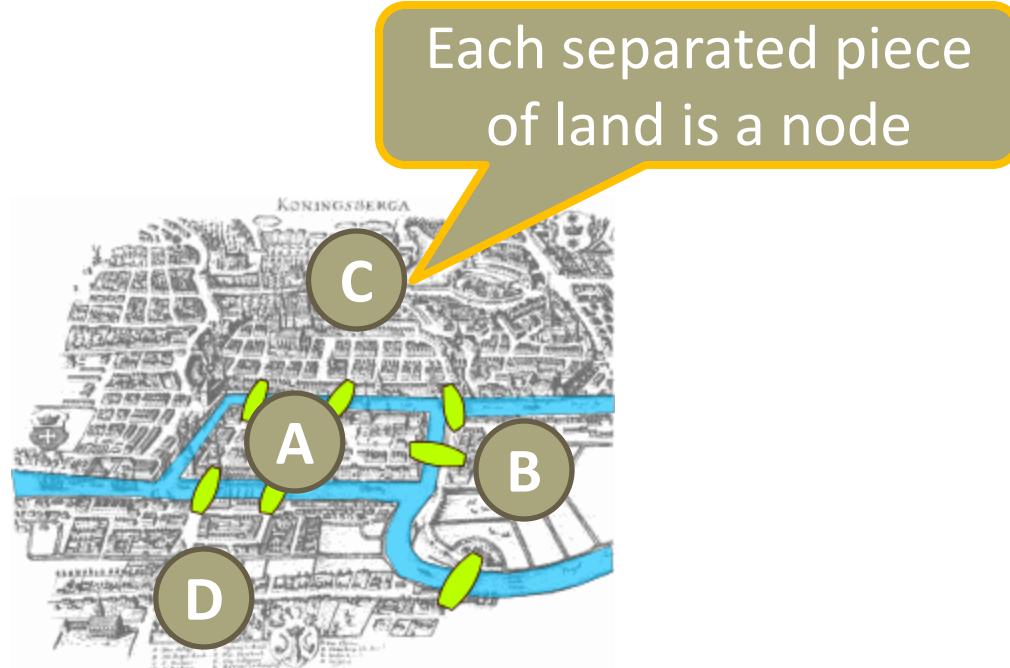
# Graph Data: Euler's Seven Bridges of Königsberg (0)



“find a walk through the city that would cross each bridge once and only once”

[http://en.wikipedia.org/wiki/Seven\\_Bridges\\_of\\_K%C3%B6nigsberg](http://en.wikipedia.org/wiki/Seven_Bridges_of_K%C3%B6nigsberg)

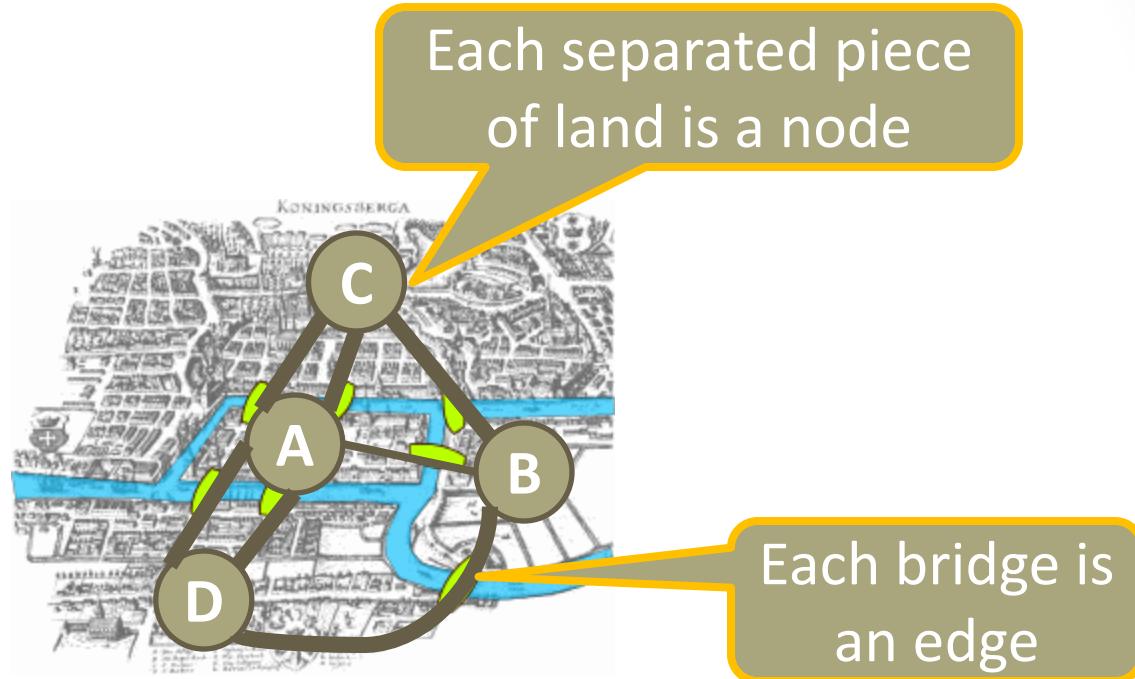
# Graph Data: Euler's Seven Bridges of Königsberg (1)



“find a walk through the city that would cross each bridge once and only once”

[http://en.wikipedia.org/wiki/Seven\\_Bridges\\_of\\_K%C3%B6nigsberg](http://en.wikipedia.org/wiki/Seven_Bridges_of_K%C3%B6nigsberg)

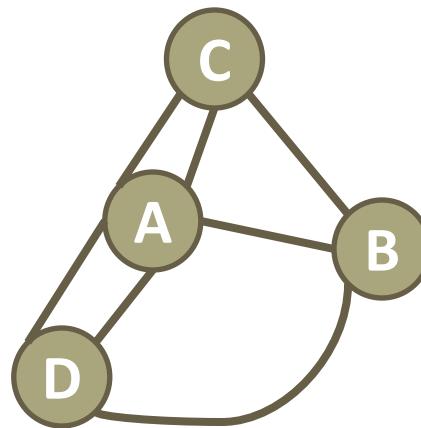
# Graph Data: Euler's Seven Bridges of Königsberg (2)



“find a walk through the city that would cross each bridge once and only once”

[http://en.wikipedia.org/wiki/Seven\\_Bridges\\_of\\_K%C3%B6nigsberg](http://en.wikipedia.org/wiki/Seven_Bridges_of_K%C3%B6nigsberg)

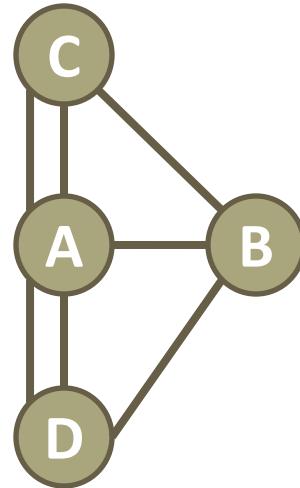
# Graph Data: Euler's Seven Bridges of Königsberg (3)



“find a walk through the city that would cross each bridge once and only once”

[http://en.wikipedia.org/wiki/Seven\\_Bridges\\_of\\_K%C3%B6nigsberg](http://en.wikipedia.org/wiki/Seven_Bridges_of_K%C3%B6nigsberg)

# Graph Data: Euler's Seven Bridges of Königsberg (4)

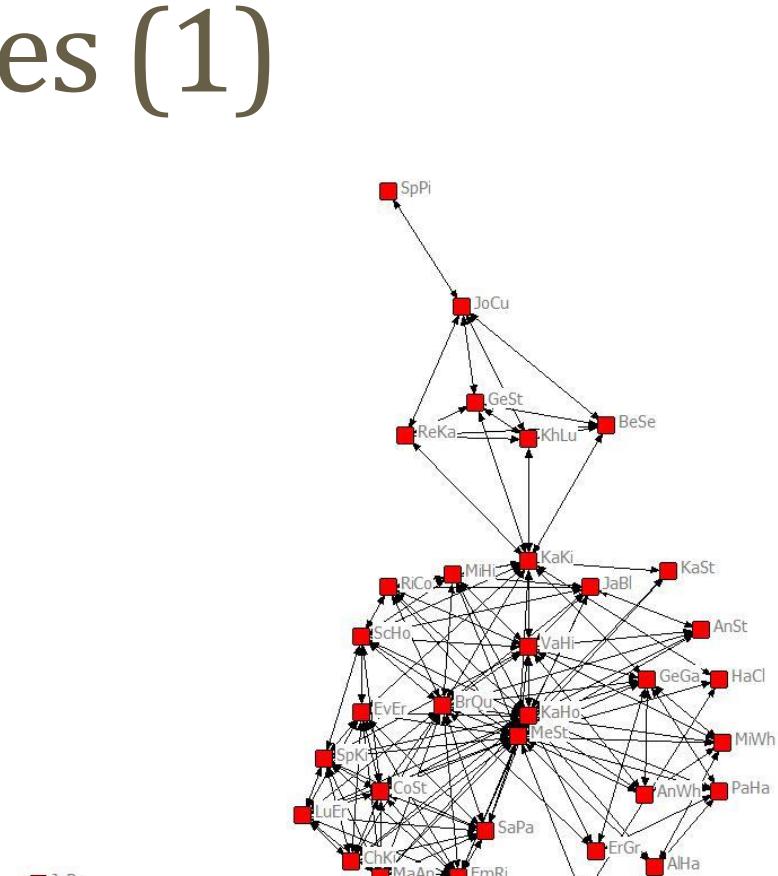
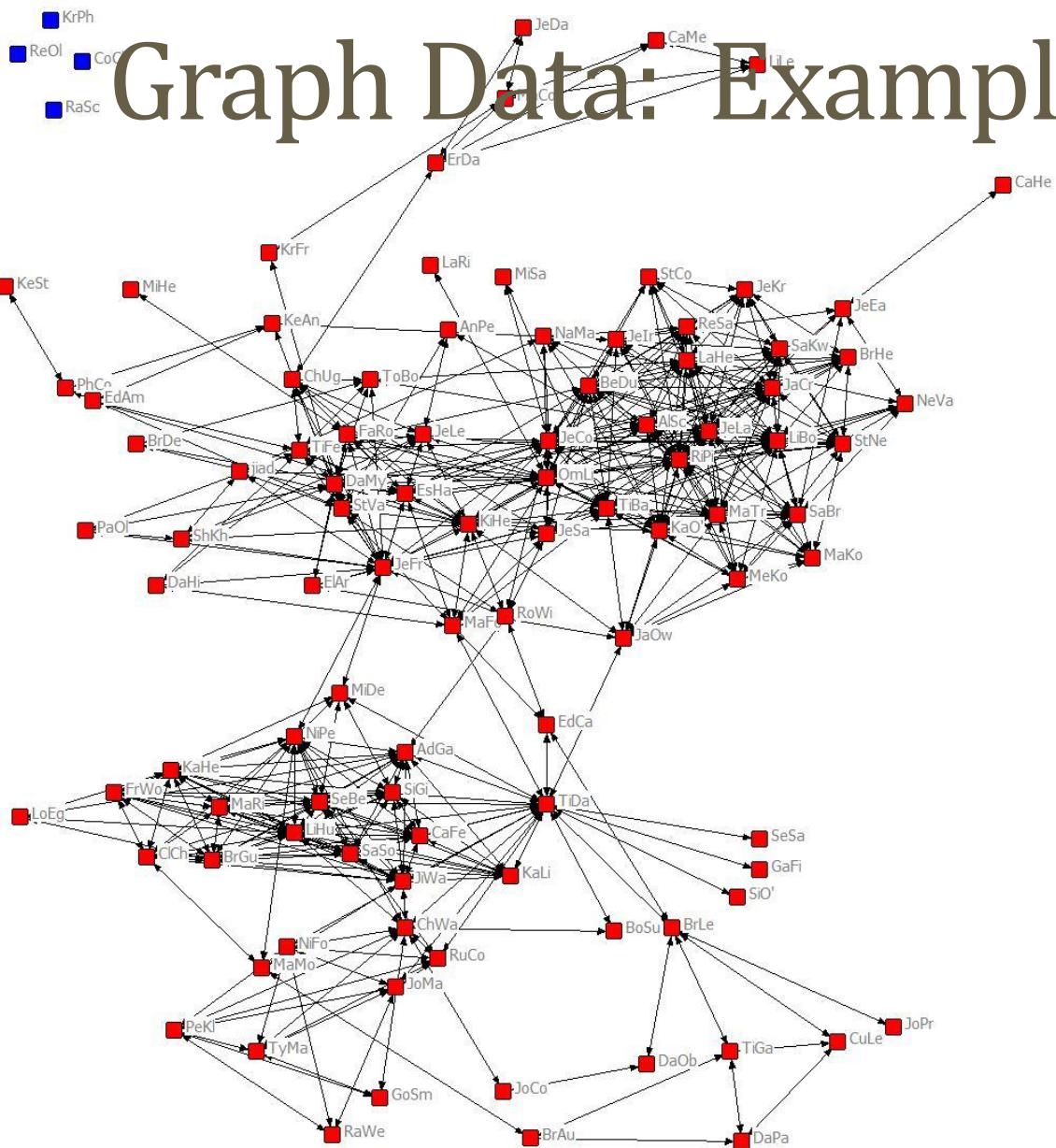


“find a walk through the city that would cross each bridge once and only once”

[http://en.wikipedia.org/wiki/Seven\\_Bridges\\_of\\_K%C3%B6nigsberg](http://en.wikipedia.org/wiki/Seven_Bridges_of_K%C3%B6nigsberg)

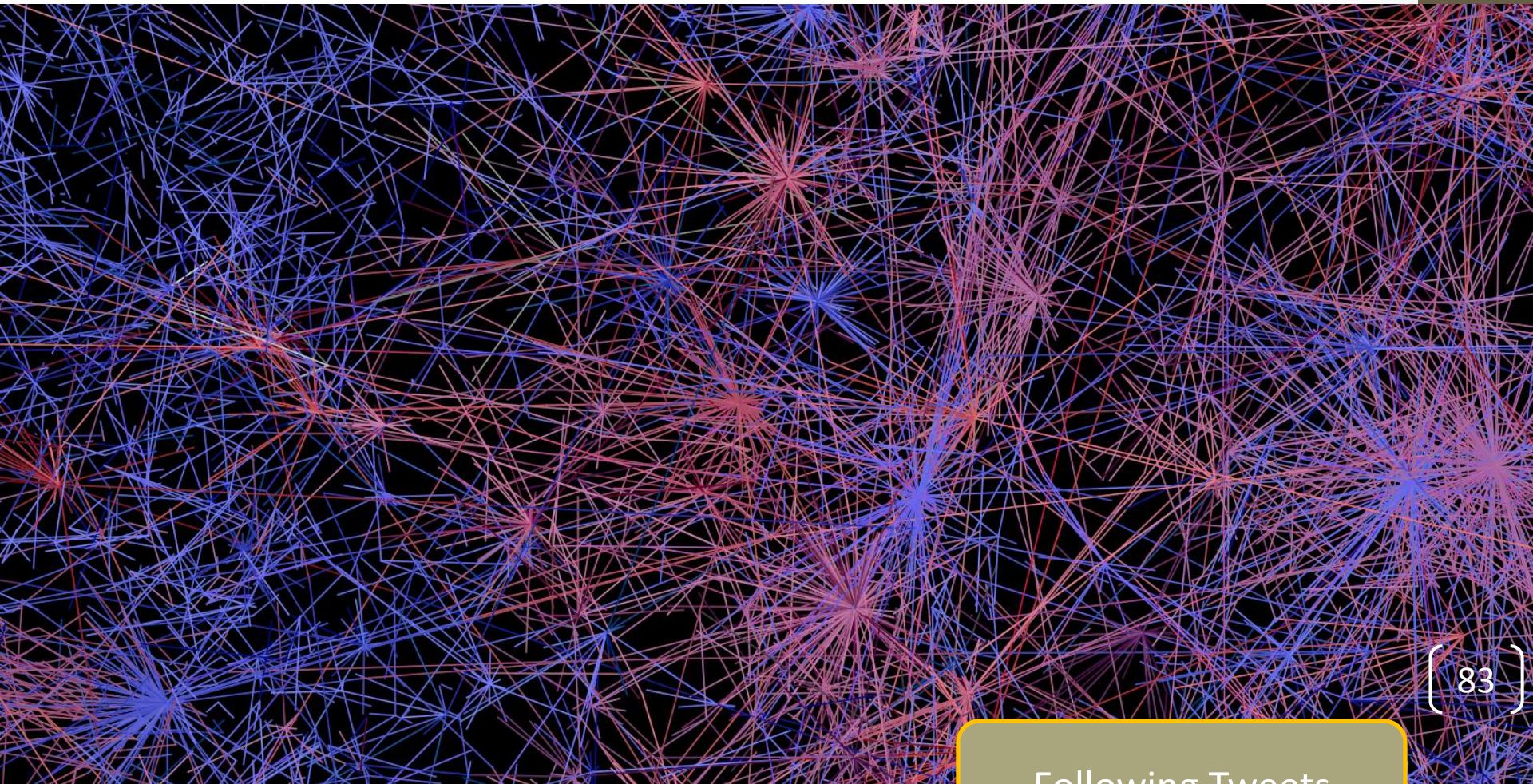
# Graph Data: Examples (0)

- Graphs represent many things like Associations and Networks.
- Graphs represent the world wide web
- Graphs represent Facebook networks
- Graphs represent metabolic pathways



Facebook Connections

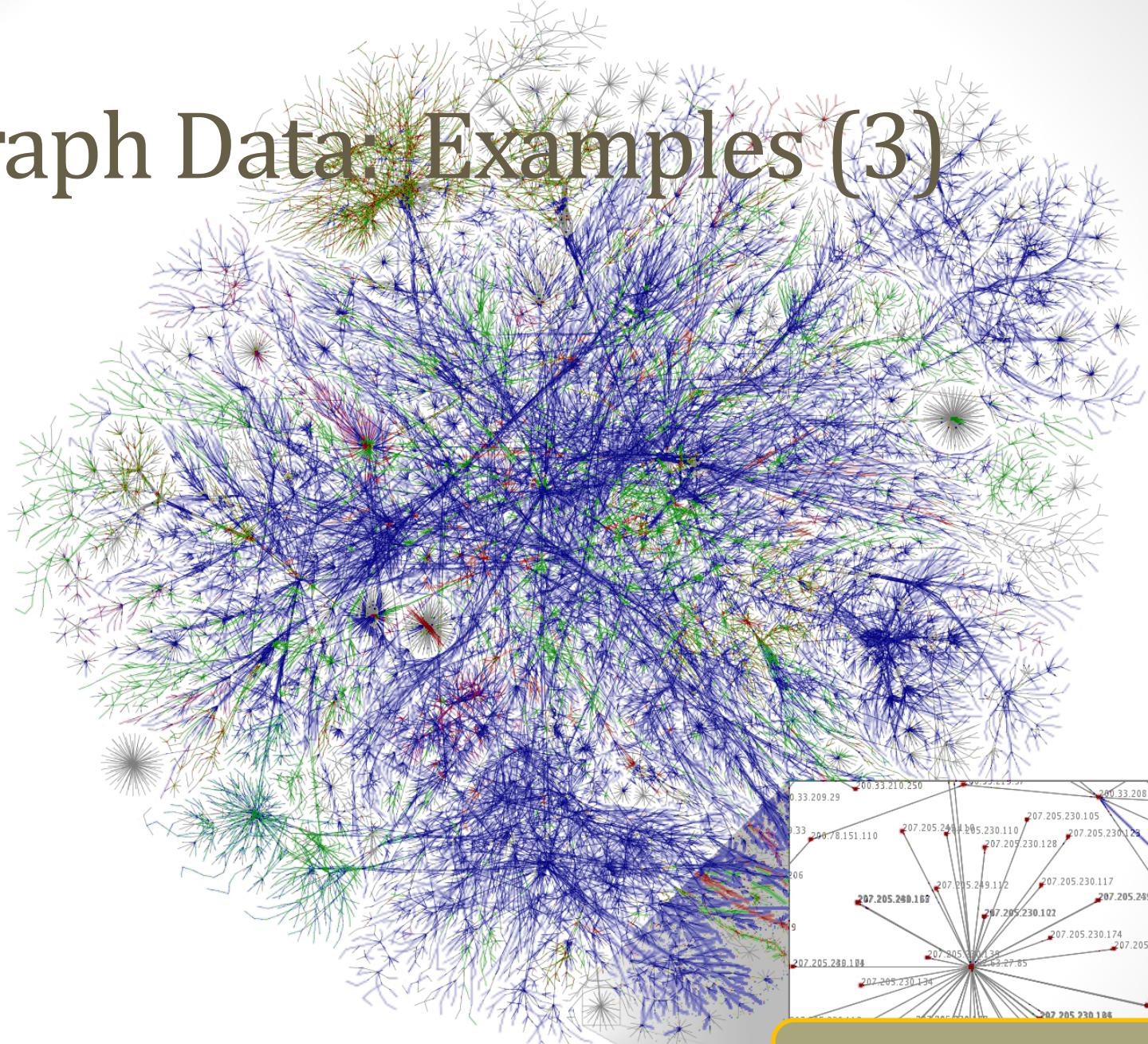
# Graph Data: Examples (2)



[ 83 ]

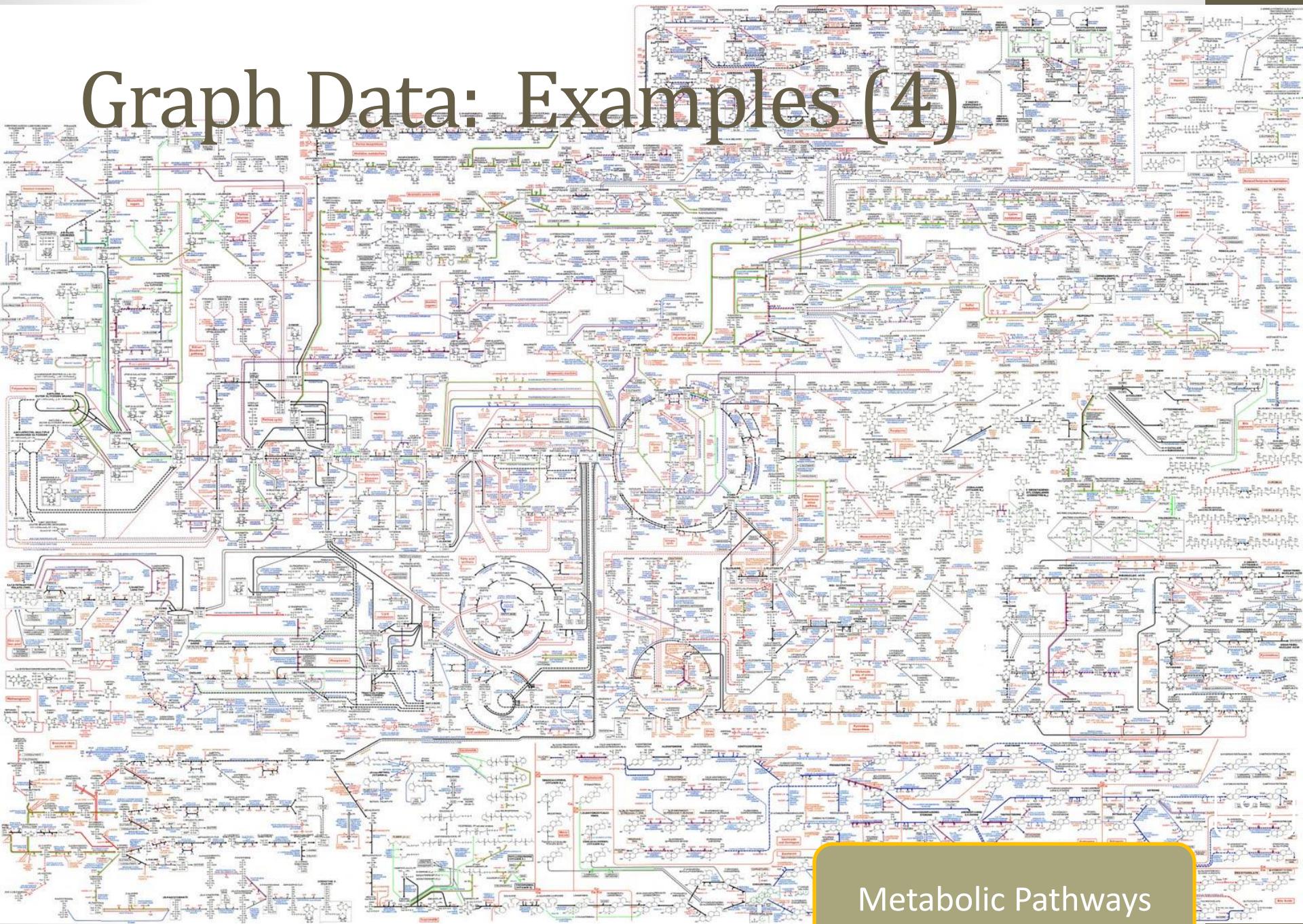
Following Tweets

# Graph Data: Examples (3)

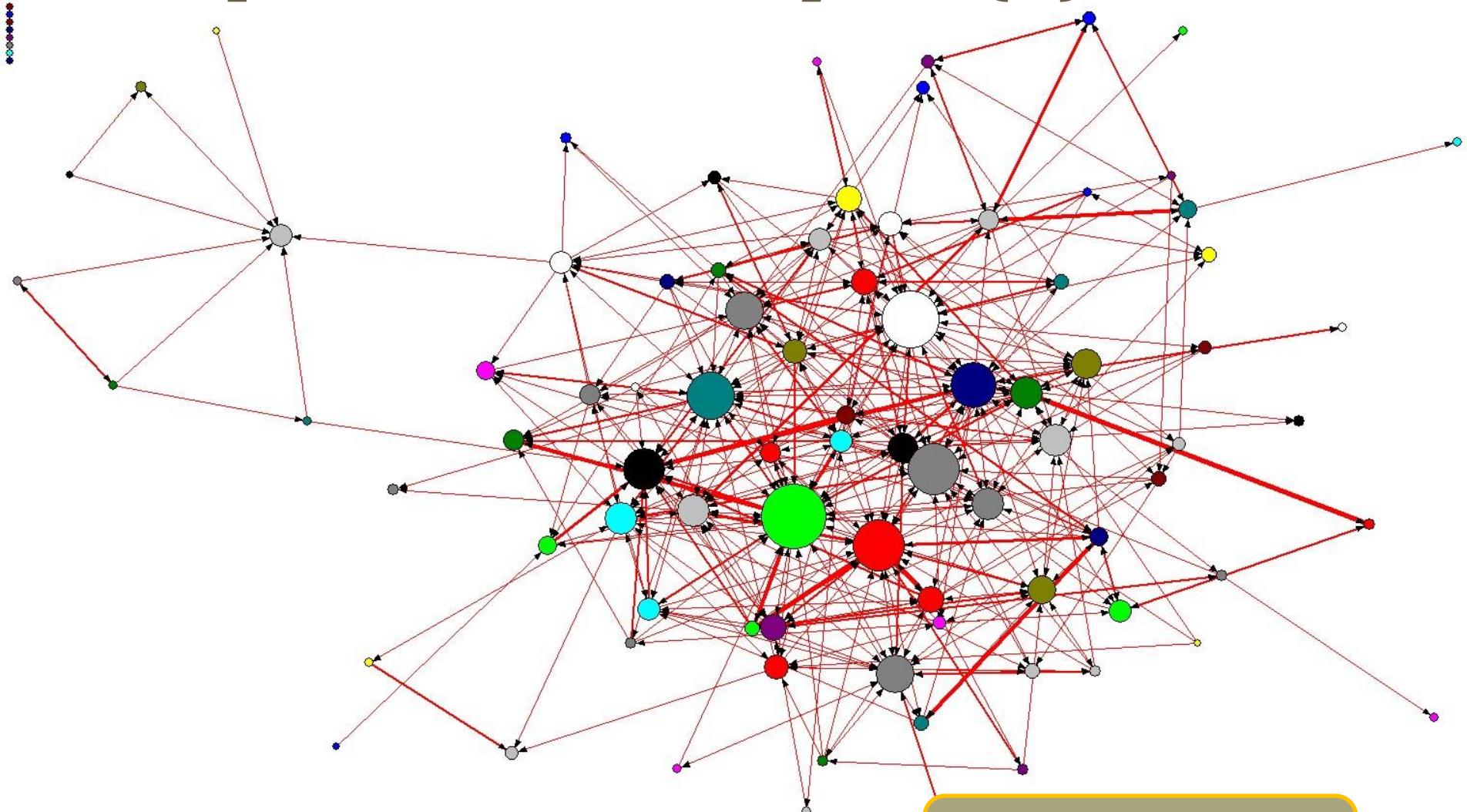


# Section of the Internet

# Graph Data: Examples (4)

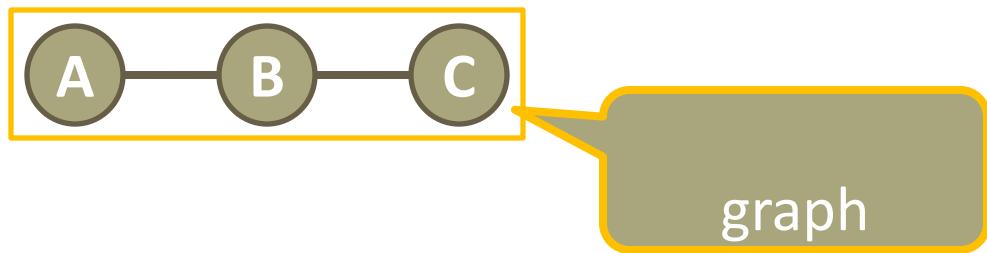


# Graph Data: Examples (5)

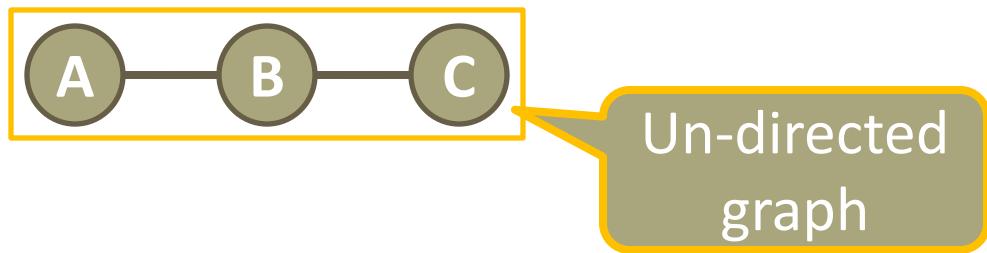


# Ranked Web Pages

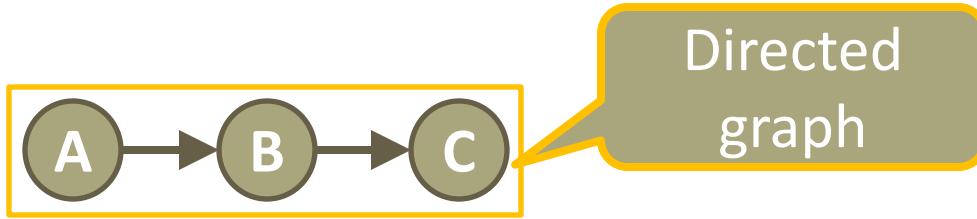
# Graph Data: Directed and Un-directed Graphs (0)



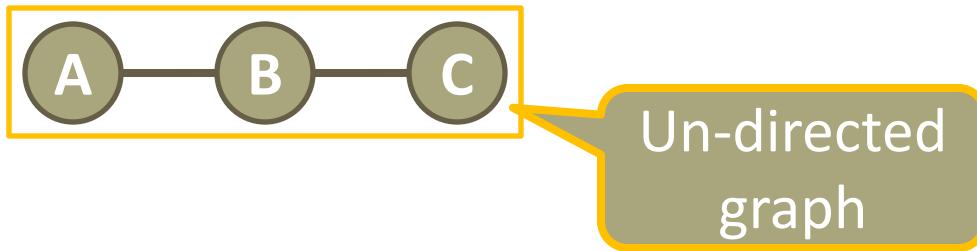
# Graph Data: Directed and Un-directed Graphs (1)



# Graph Data: Directed and Un-directed Graphs (2)

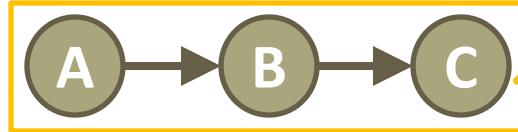


Directed  
graph



Un-directed  
graph

# Graph Data: Directed and Un-directed Graphs (3)



Directed  
graph

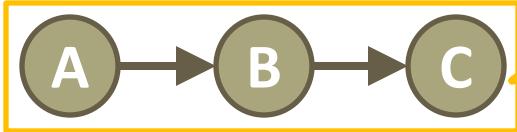
“The link goes from  
A to B”



Un-directed  
graph

“The link is  
between A and B”

# Graph Data: Directed and Un-directed Graphs (4)



Directed graph

"The link goes from A to B"

Example of directed graphs:  
World Wide Web

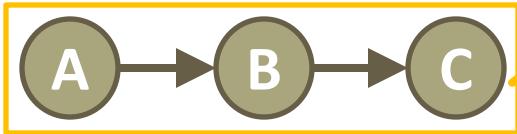


Un-directed graph

"The link is between A and B"

Examples of un-directed graphs:  
Facebook, LinkedIn

# Graph Data: Directed and Undirected Graphs (5)



Directed graph

"The link goes from A to B"

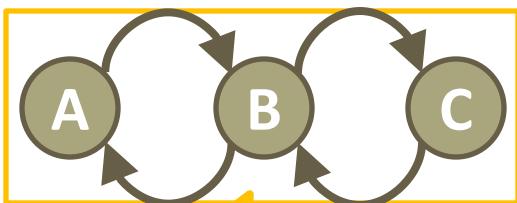
Example of directed graphs:  
World Wide Web



Un-directed graph

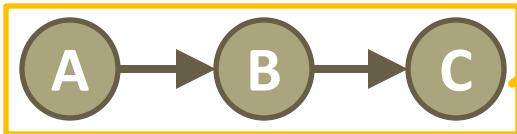
"The link is between A and B"

Examples of un-directed graphs:  
Facebook, LinkedIn



graph with bi-directional links

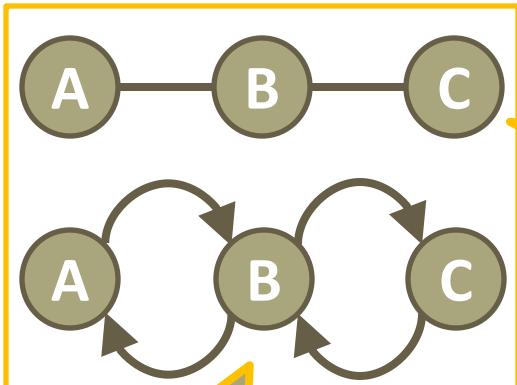
# Graph Data: Directed and Undirected Graphs (6)



Directed graph

"The link goes from A to B"

Example of directed graphs:  
World Wide Web



Un-directed graph

"The link is between A and B"

Examples of un-directed graphs:  
Facebook, LinkedIn

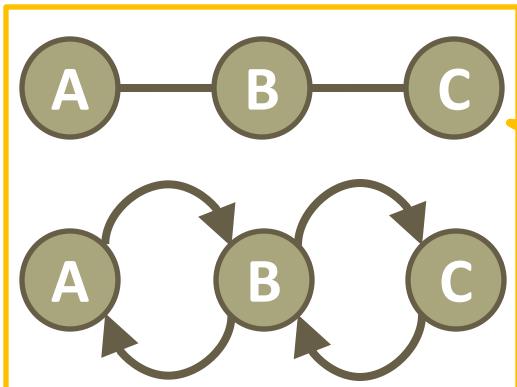
graph with bi-directional links

A graph with bi-directional links is like an undirected graph

# Graph Data: Formalize as Rectangular Data (0)



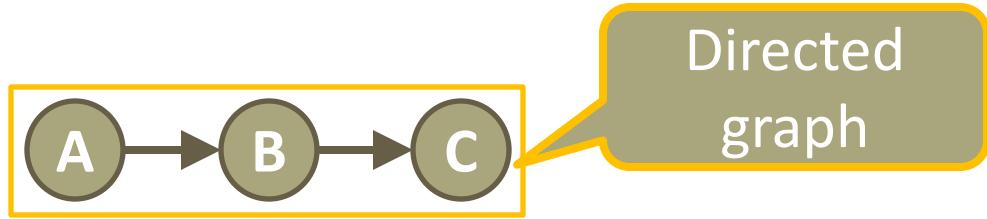
Directed  
graph



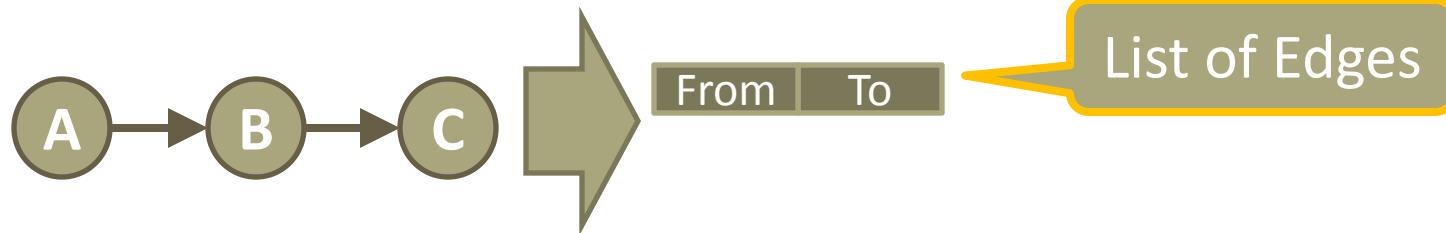
Un-directed  
graph

- How can we perform operations on graphs?
- How can we formalize graphs as rectangular data?

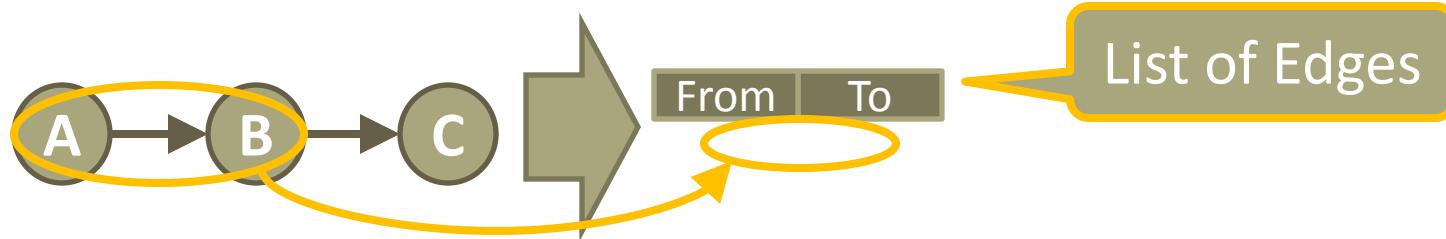
# Graph Data: Formalize as Rectangular Data (1)



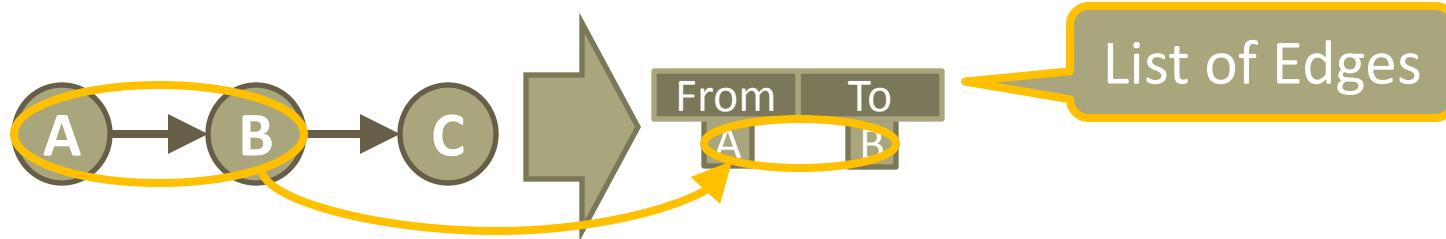
# Graph Data: Formalize as Rectangular Data (2)



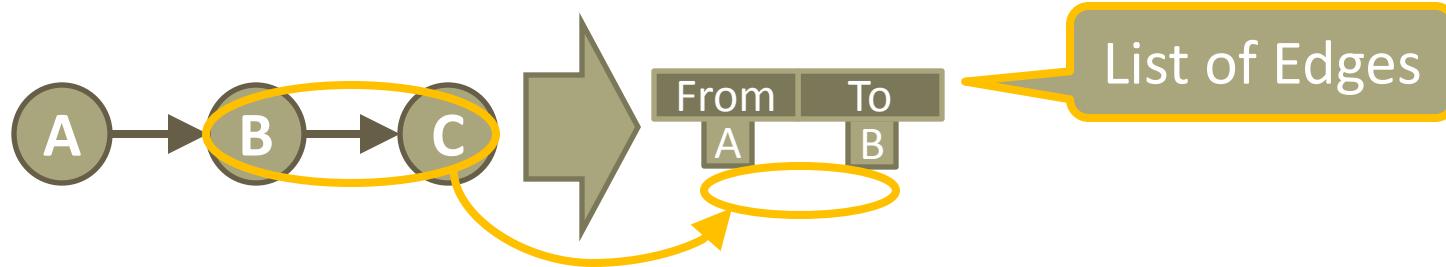
# Graph Data: Formalize as Rectangular Data (3)



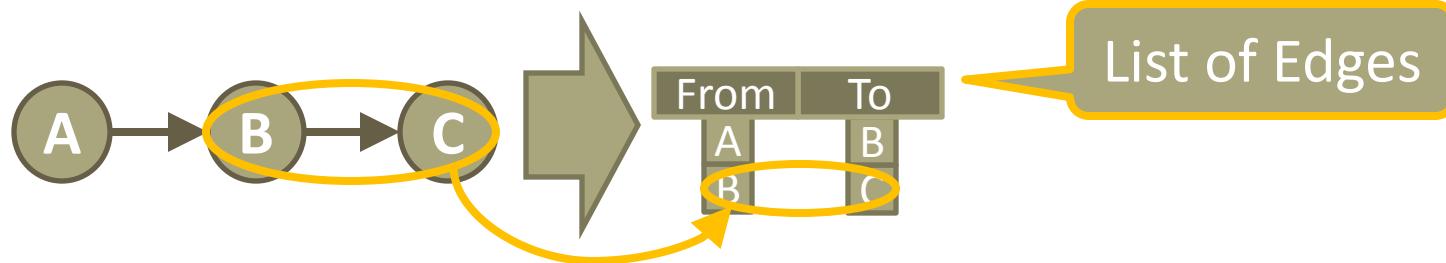
# Graph Data: Formalize as Rectangular Data (4)



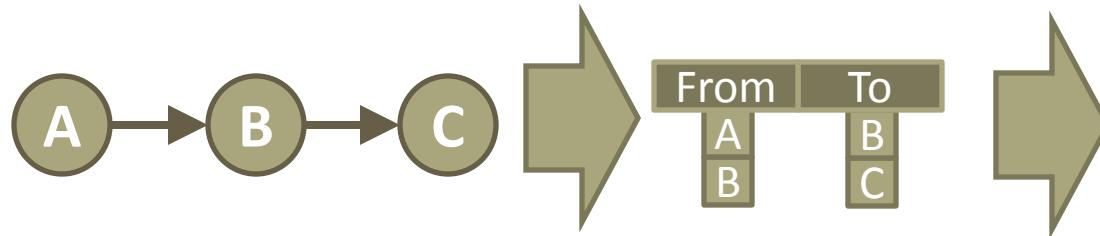
# Graph Data: Formalize as Rectangular Data (5)



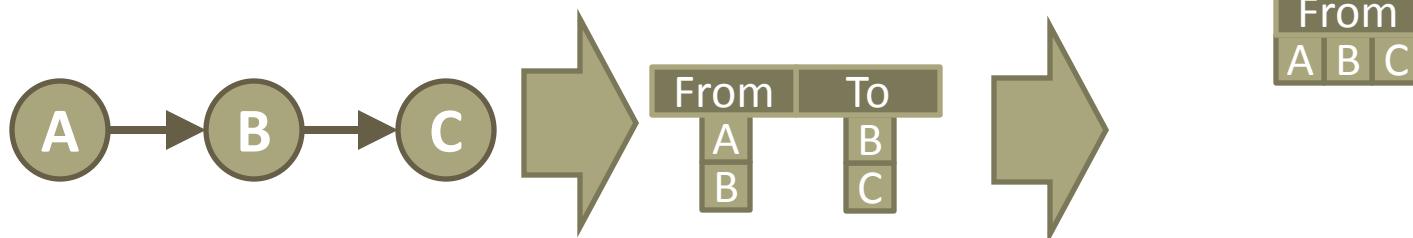
# Graph Data: Formalize as Rectangular Data (6)



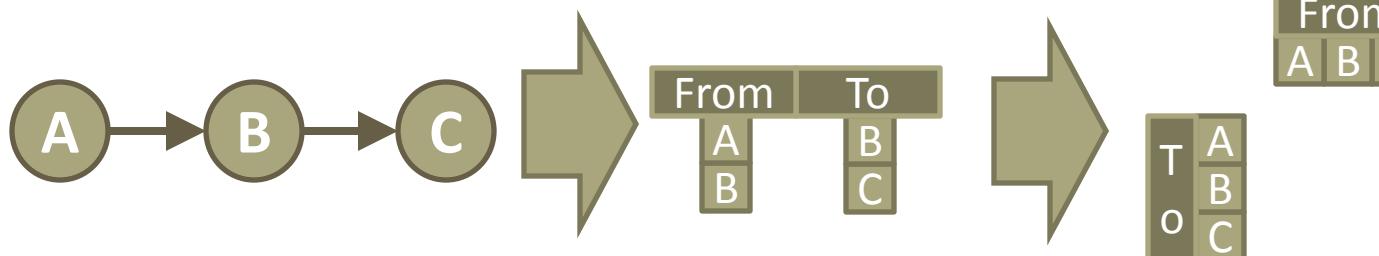
# Graph Data: Formalize as Rectangular Data (7)



# Graph Data: Formalize as Rectangular Data (8)



# Graph Data: Formalize as Rectangular Data (9)



Edges start from  
these nodes

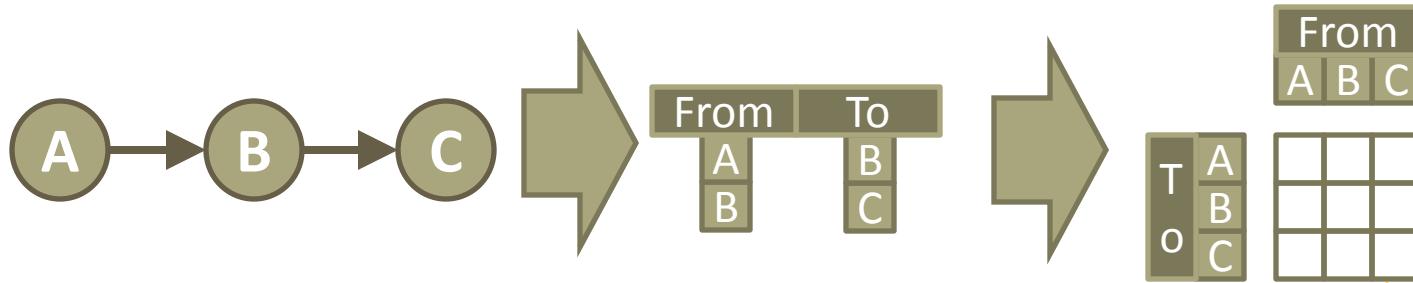
From  
A B C

Edges go to these  
nodes

To  
B C  
A

Every node could emanate and receive  
links.

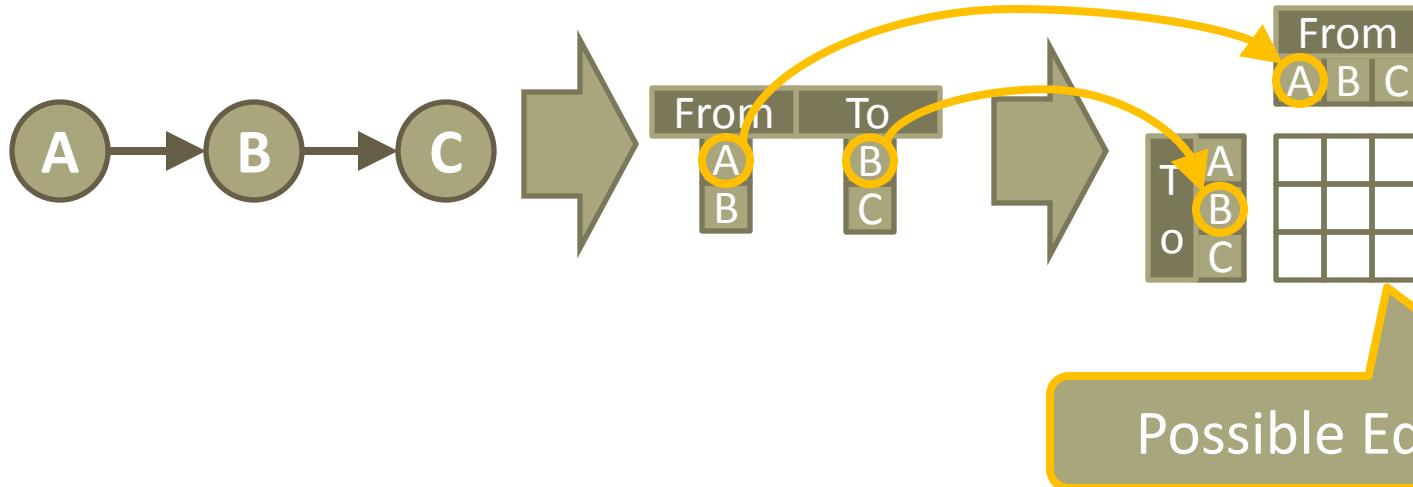
# Graph Data: Formalize as Rectangular Data (10)



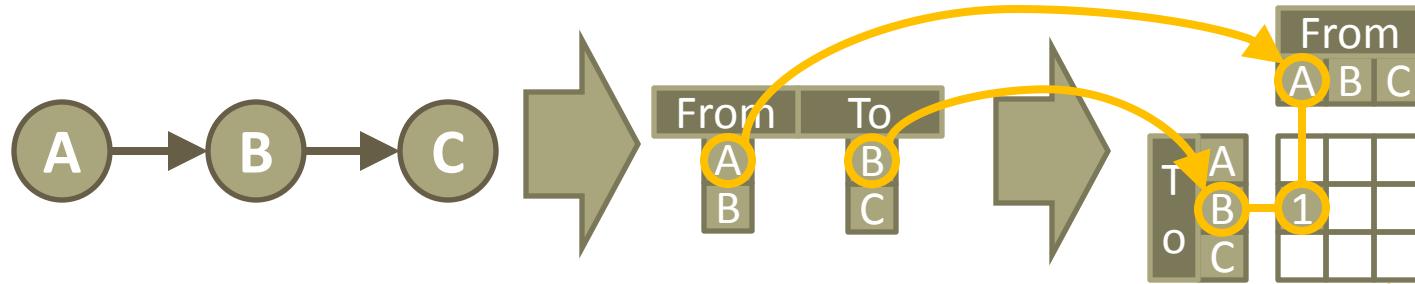
Possible Edges

Every node could emanate and receive links. Therefore the matrix is square.

# Graph Data: Formalize as Rectangular Data (11)

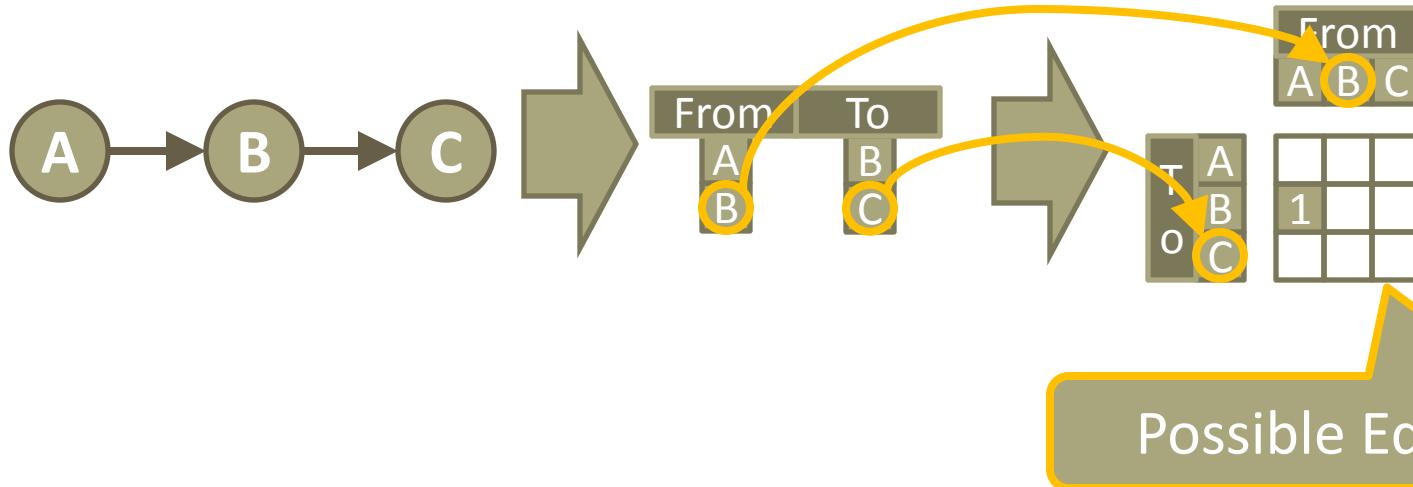


# Graph Data: Formalize as Rectangular Data (12)

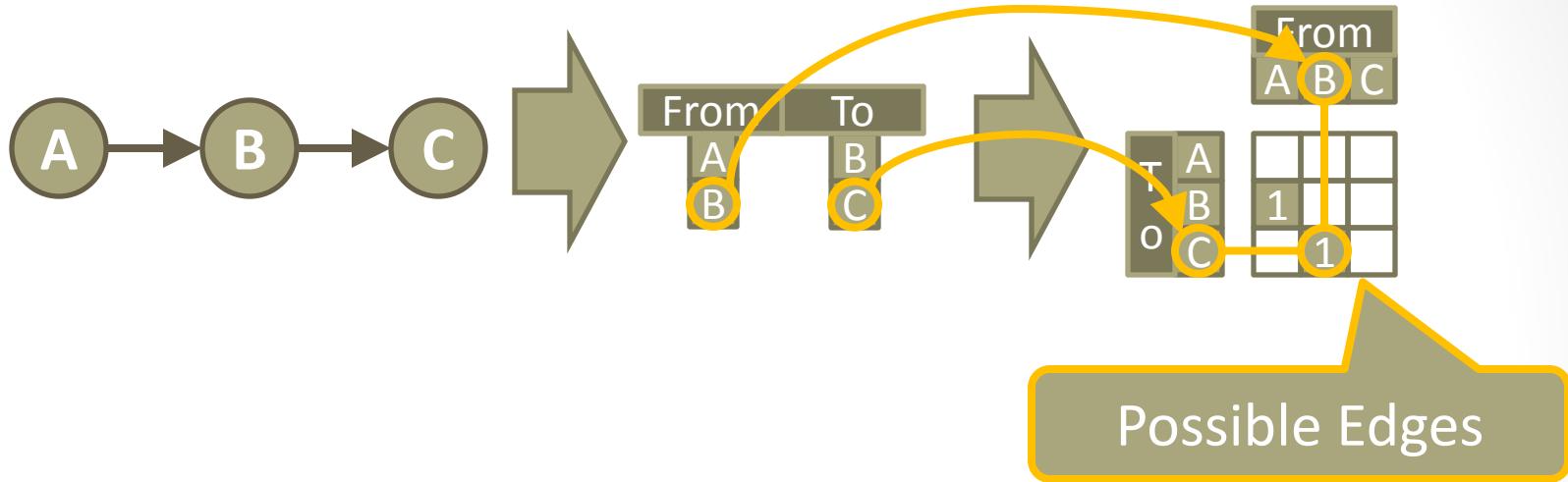


Possible Edges

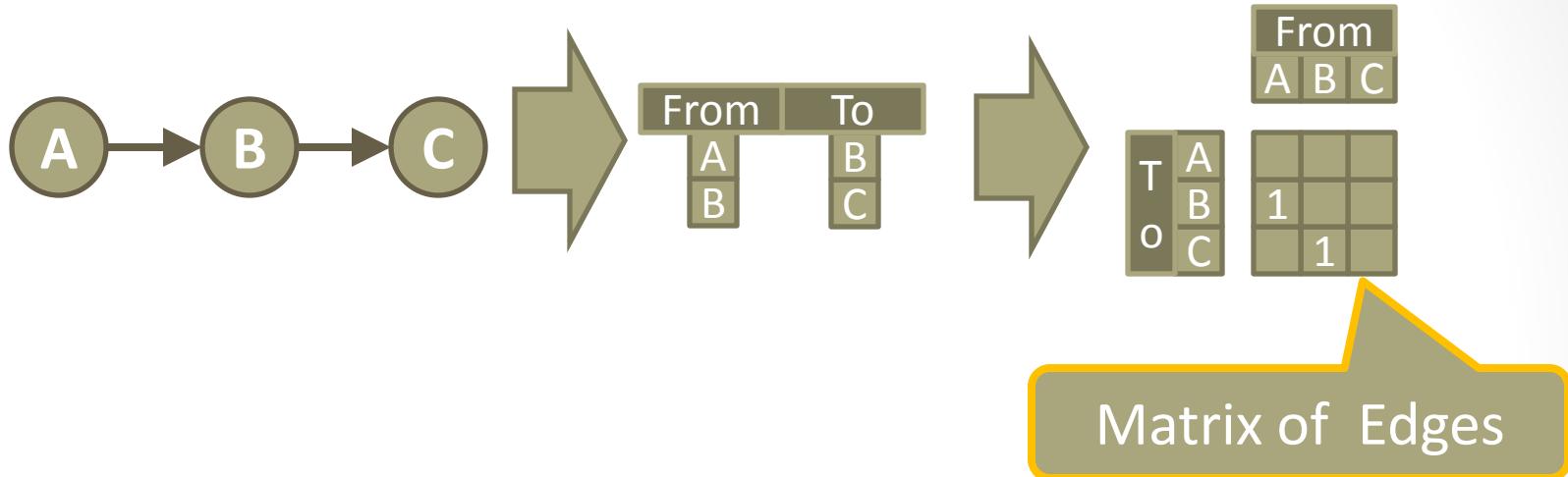
# Graph Data: Formalize as Rectangular Data (13)



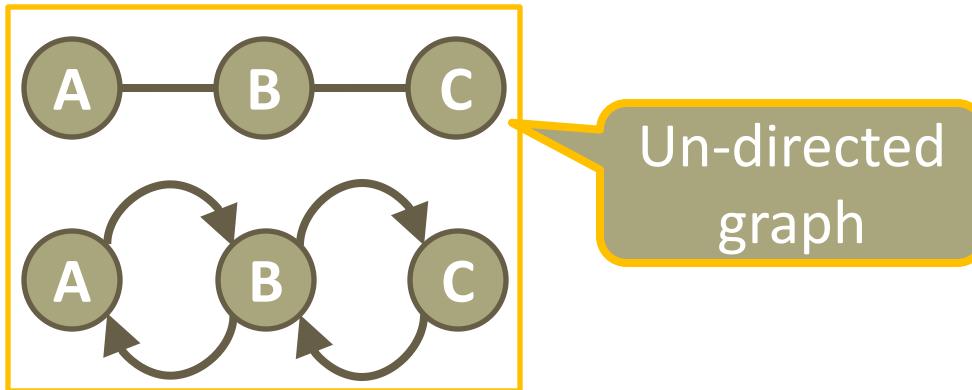
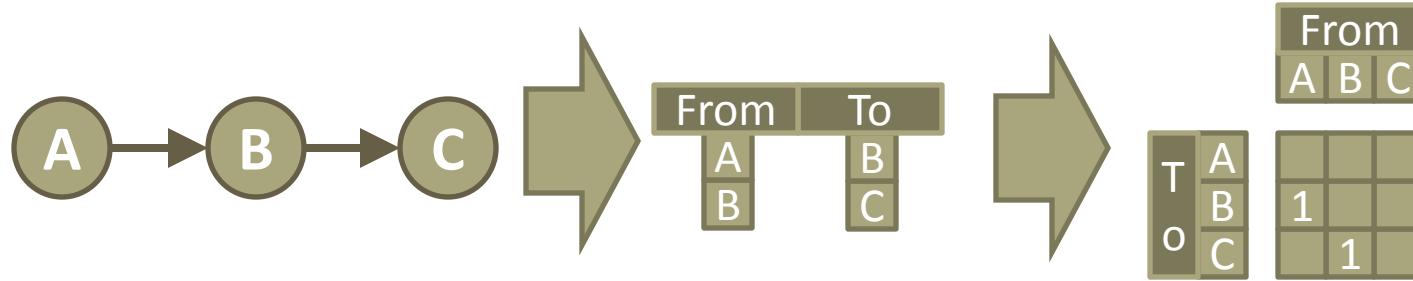
# Graph Data: Formalize as Rectangular Data (14)



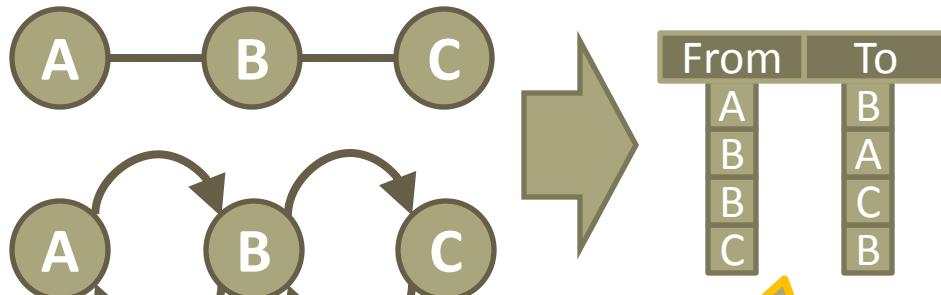
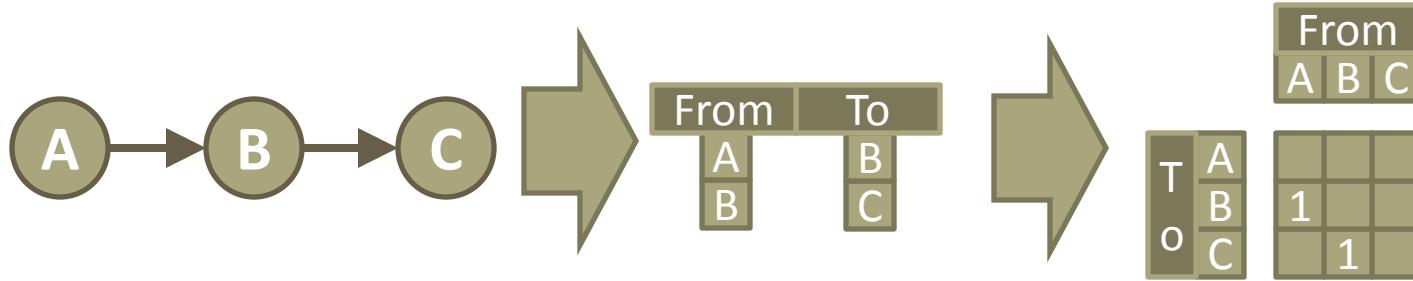
# Graph Data: Formalize as Rectangular Data (15)



# Graph Data: Formalize as Rectangular Data (16)



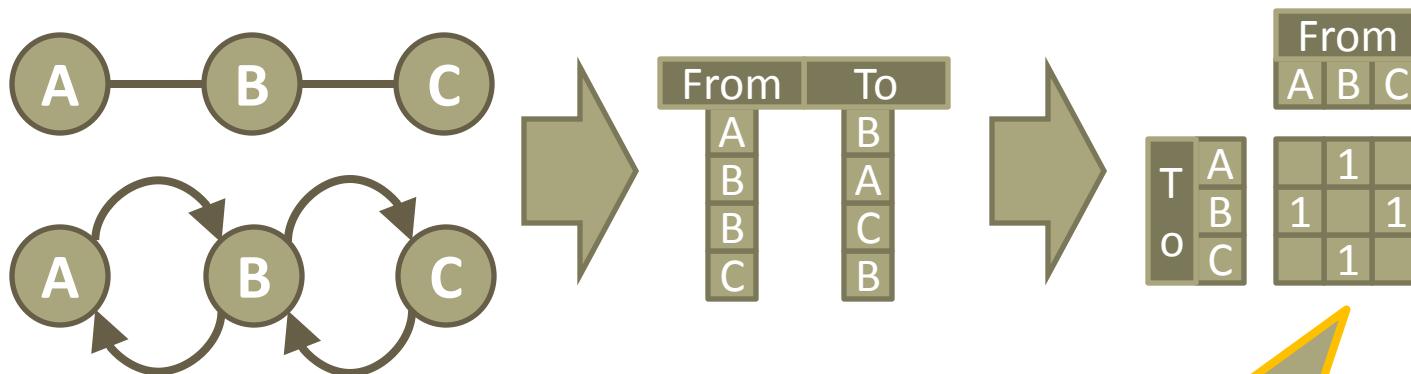
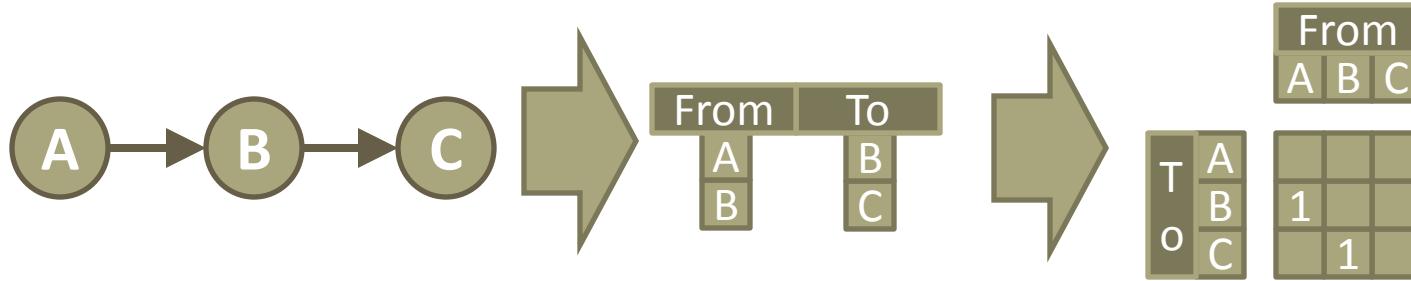
# Graph Data: Formalize as Rectangular Data (17)



List of Edges

[ 111 ]

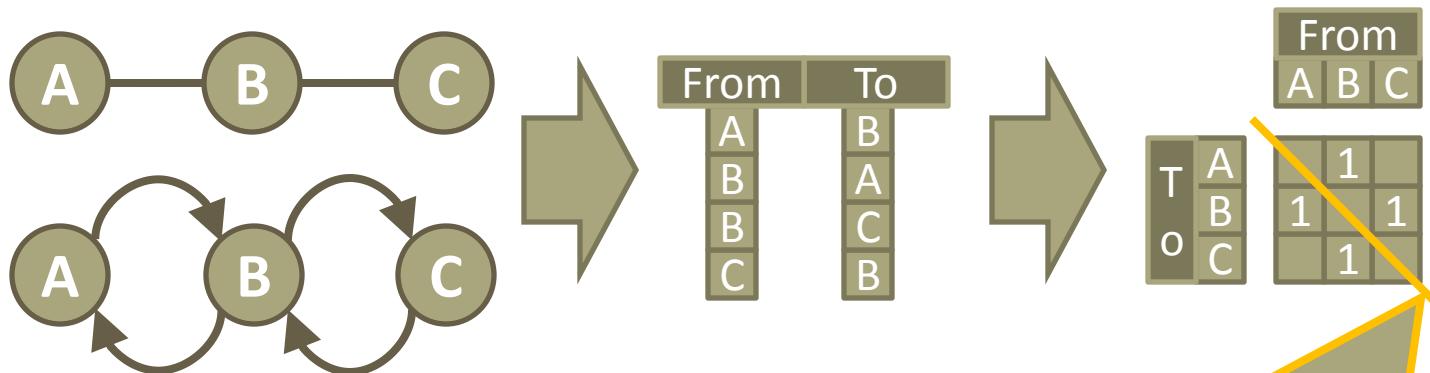
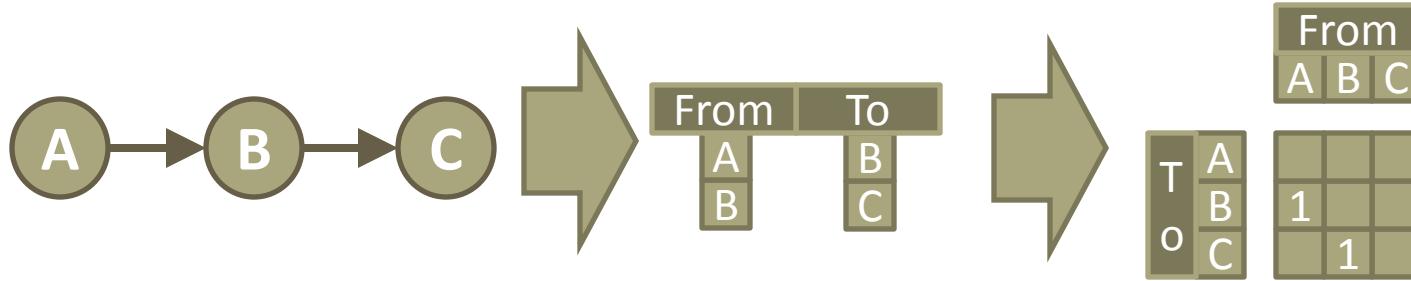
# Graph Data: Formalize as Rectangular Data (18)



Matrix of Edges

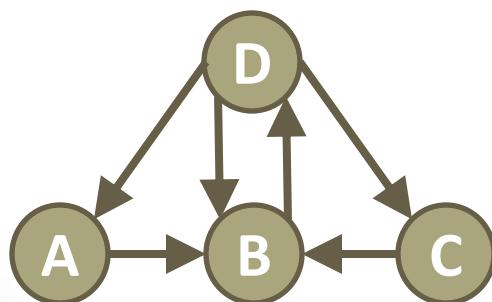
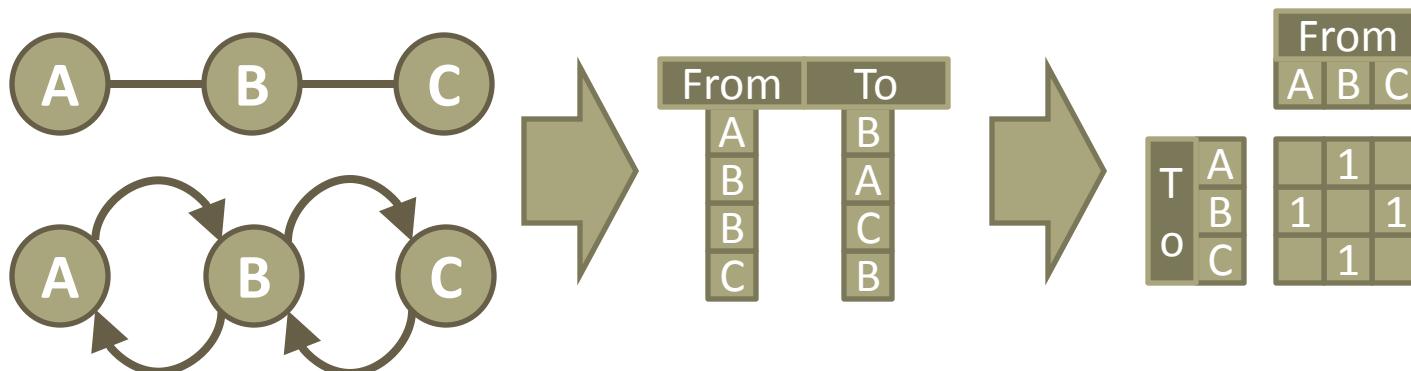
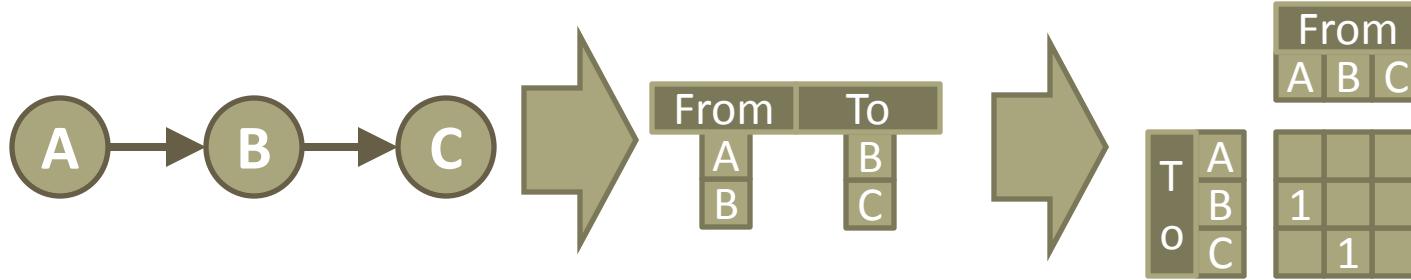
(112)

# Graph Data: Formalize as Rectangular Data (19)

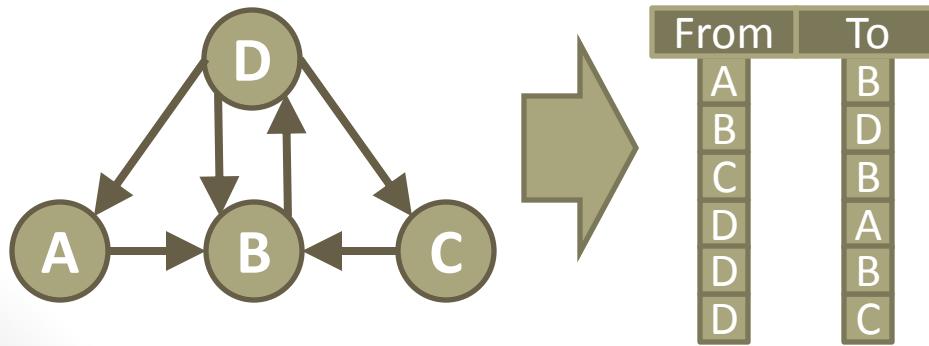
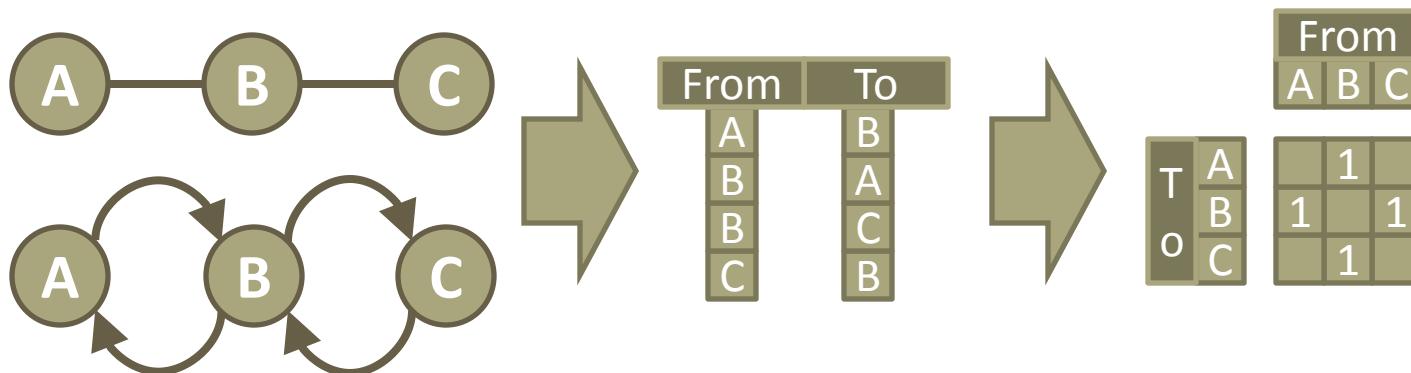
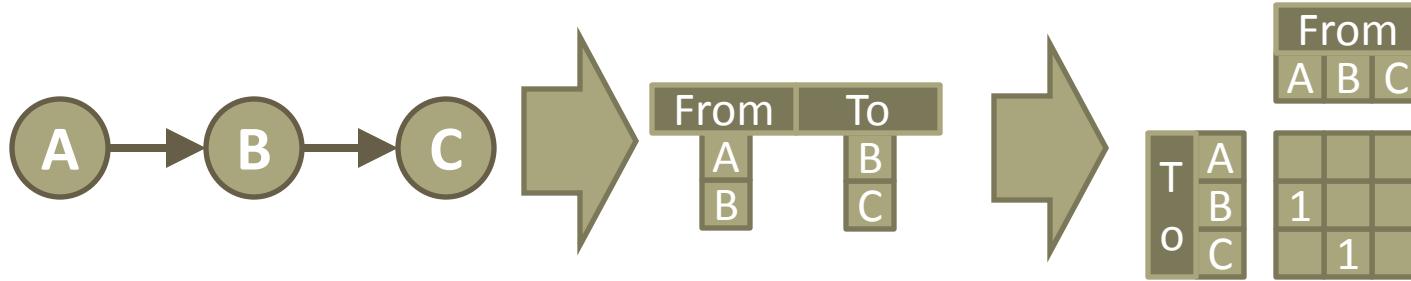


Note: undirected or bi-directional graphs Have symmetric matrices

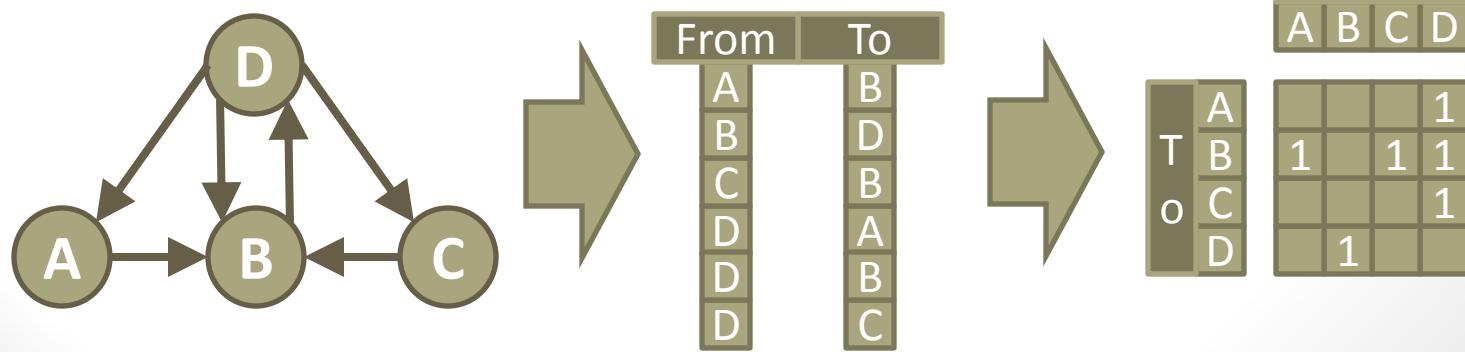
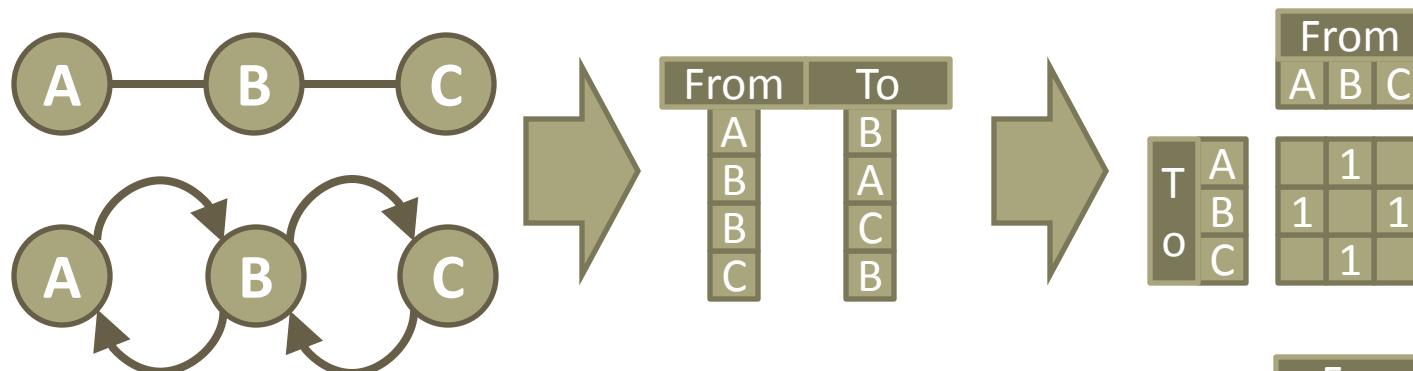
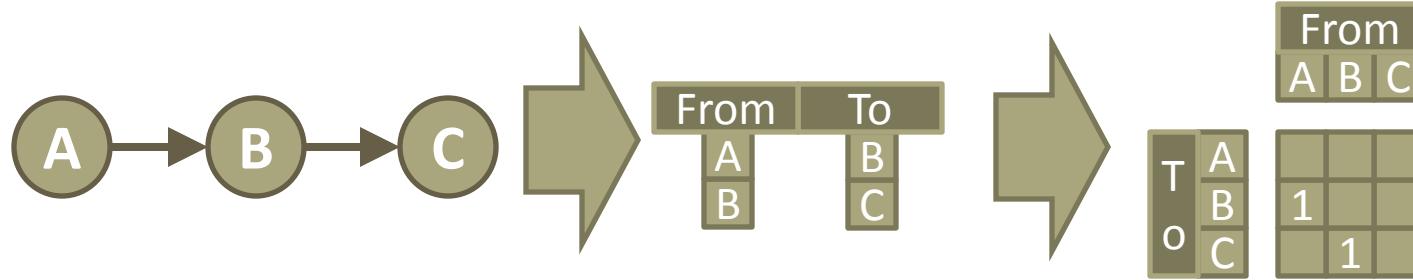
# Graph Data: Formalize as Rectangular Data (20)



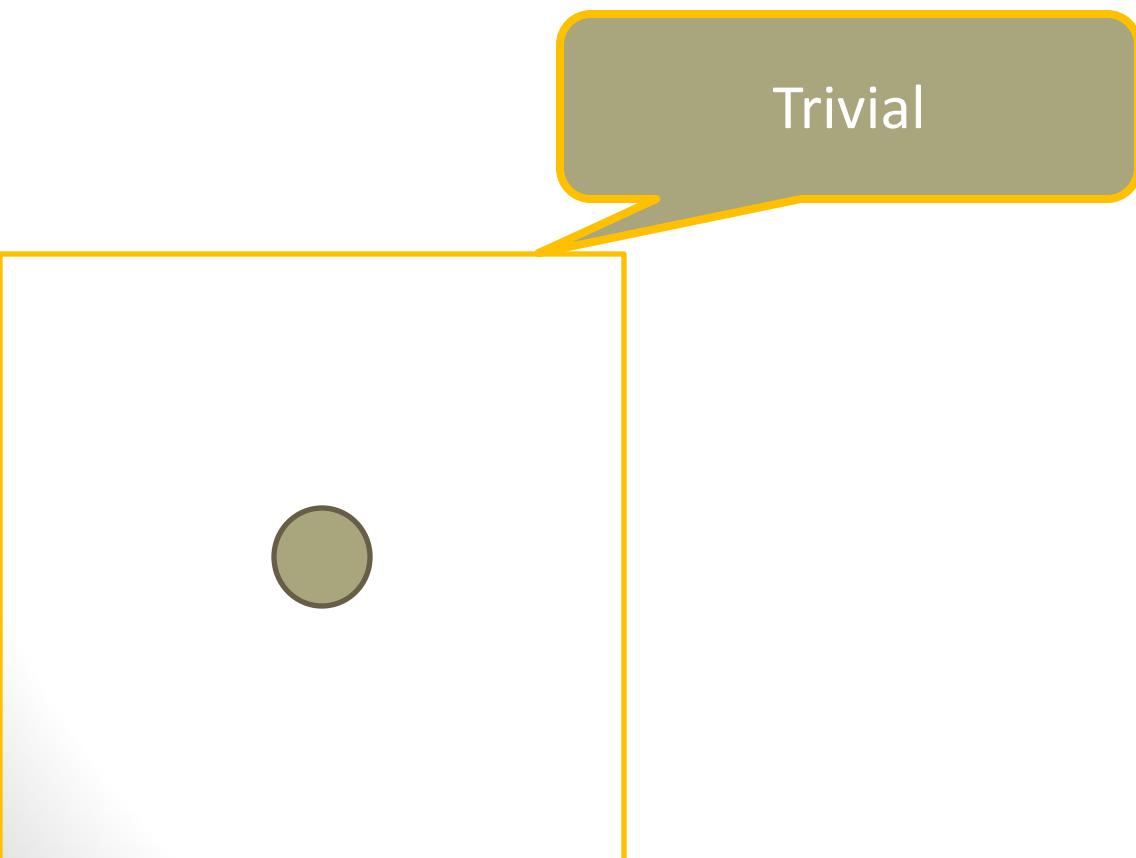
# Graph Data: Formalize as Rectangular Data (21)



# Graph Data: Formalize as Rectangular Data (22)

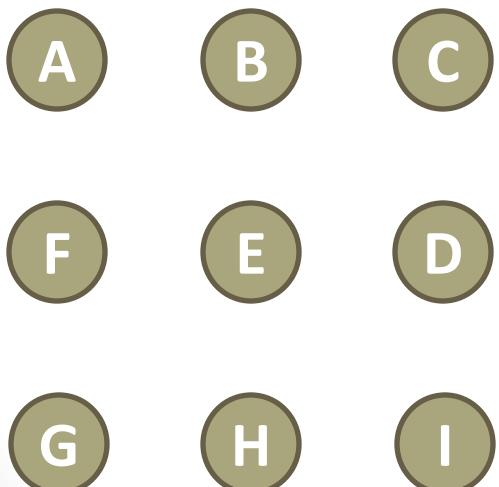


# Graph Data: Connectedness and Density (0)

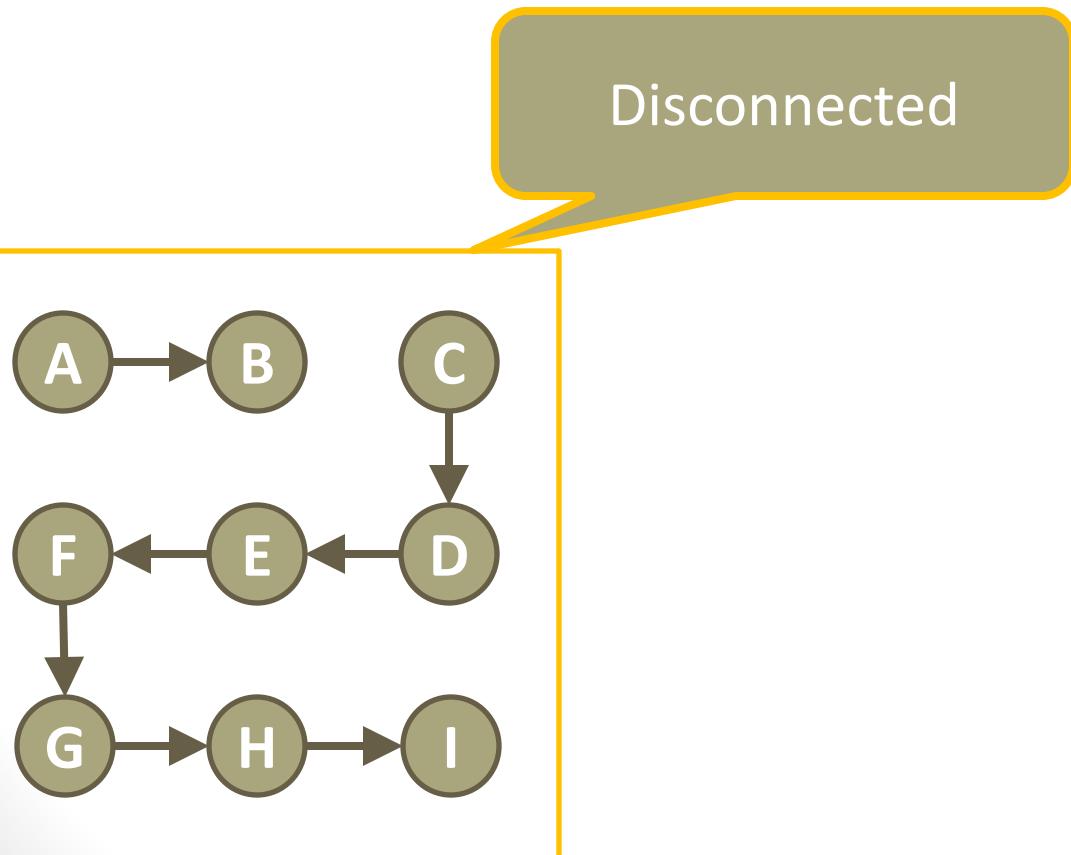


# Graph Data: Connectedness and Density (1)

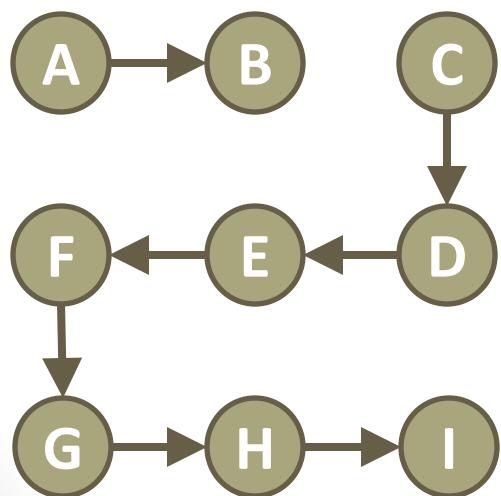
Edgeless



# Graph Data: Connectedness and Density (2)

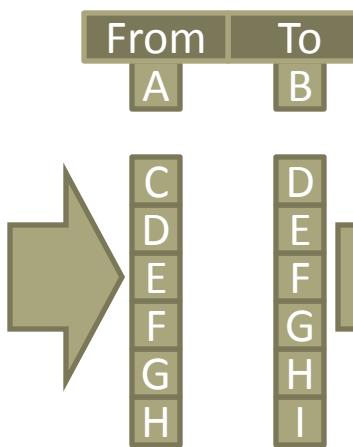
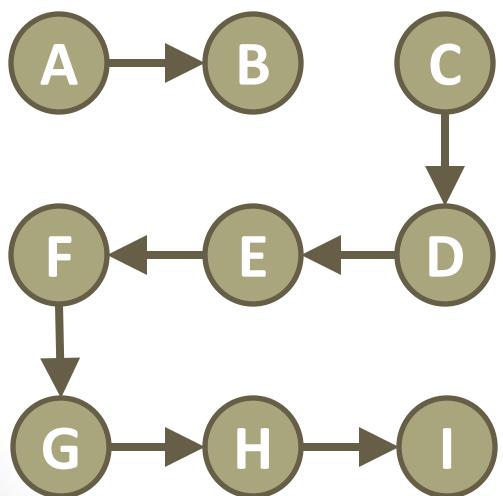


# Graph Data: Connectedness and Density (3)



From	To
A	B
C	D
D	E
D	F
E	G
F	G
G	H
H	I

# Graph Data: Connectedness and Density (4)

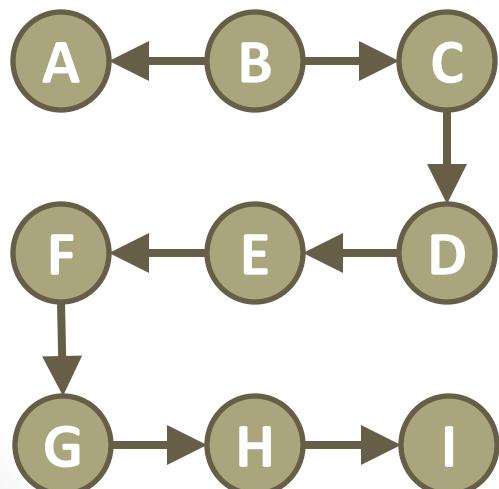


From	A	B	C	D	E	F	G	H	I
To	A	B	C	D	E	F	G	H	I
1									
	1								
		1							
			1						
				1					
					1				
						1			
							1		
								1	
									1

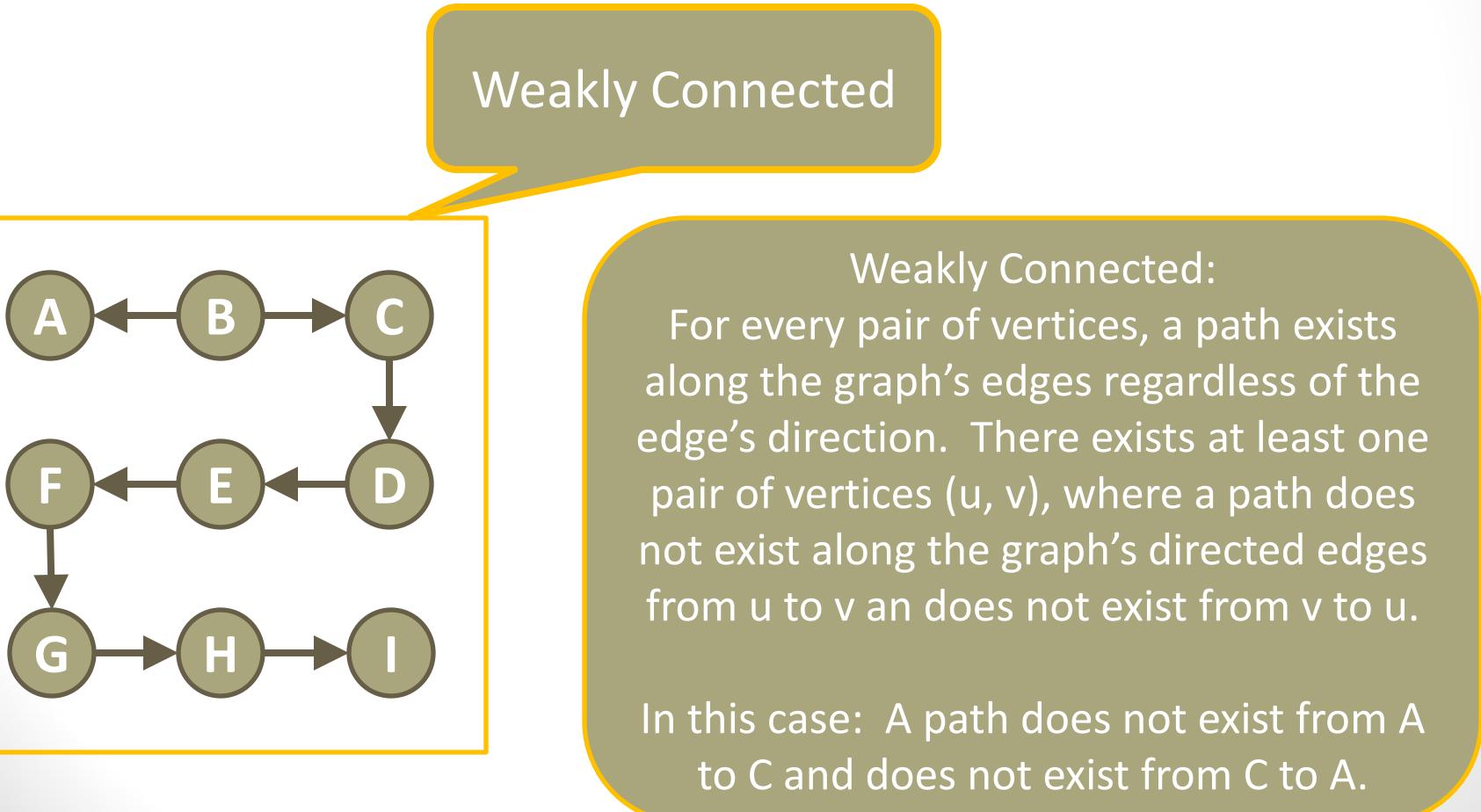
(121)

# Graph Data: Connectedness and Density (5)

Weakly Connected

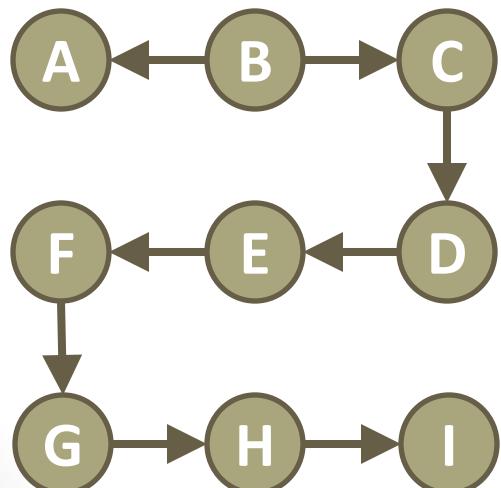


# Graph Data: Connectedness and Density (6)



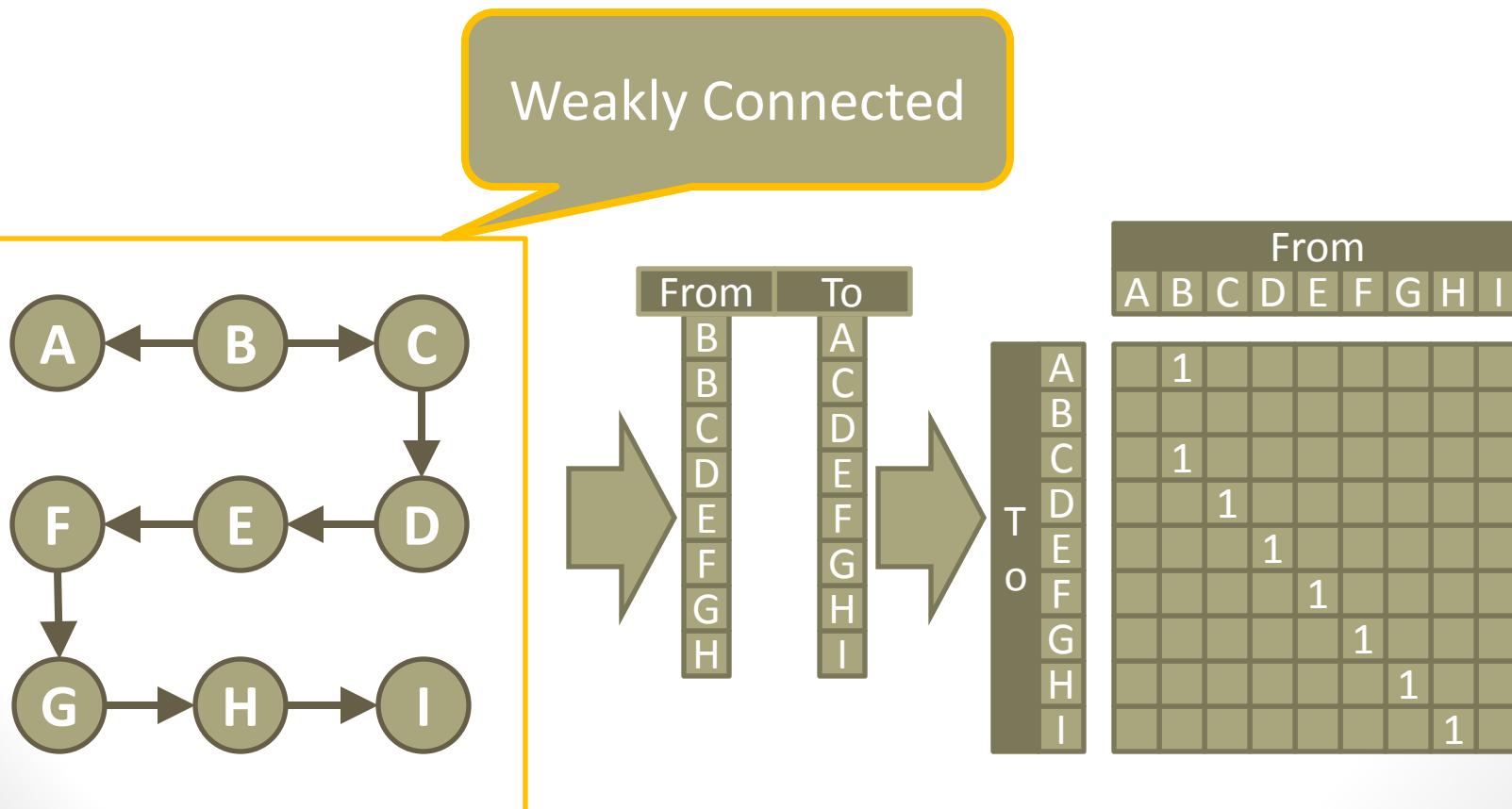
# Graph Data: Connectedness and Density (7)

Weakly Connected

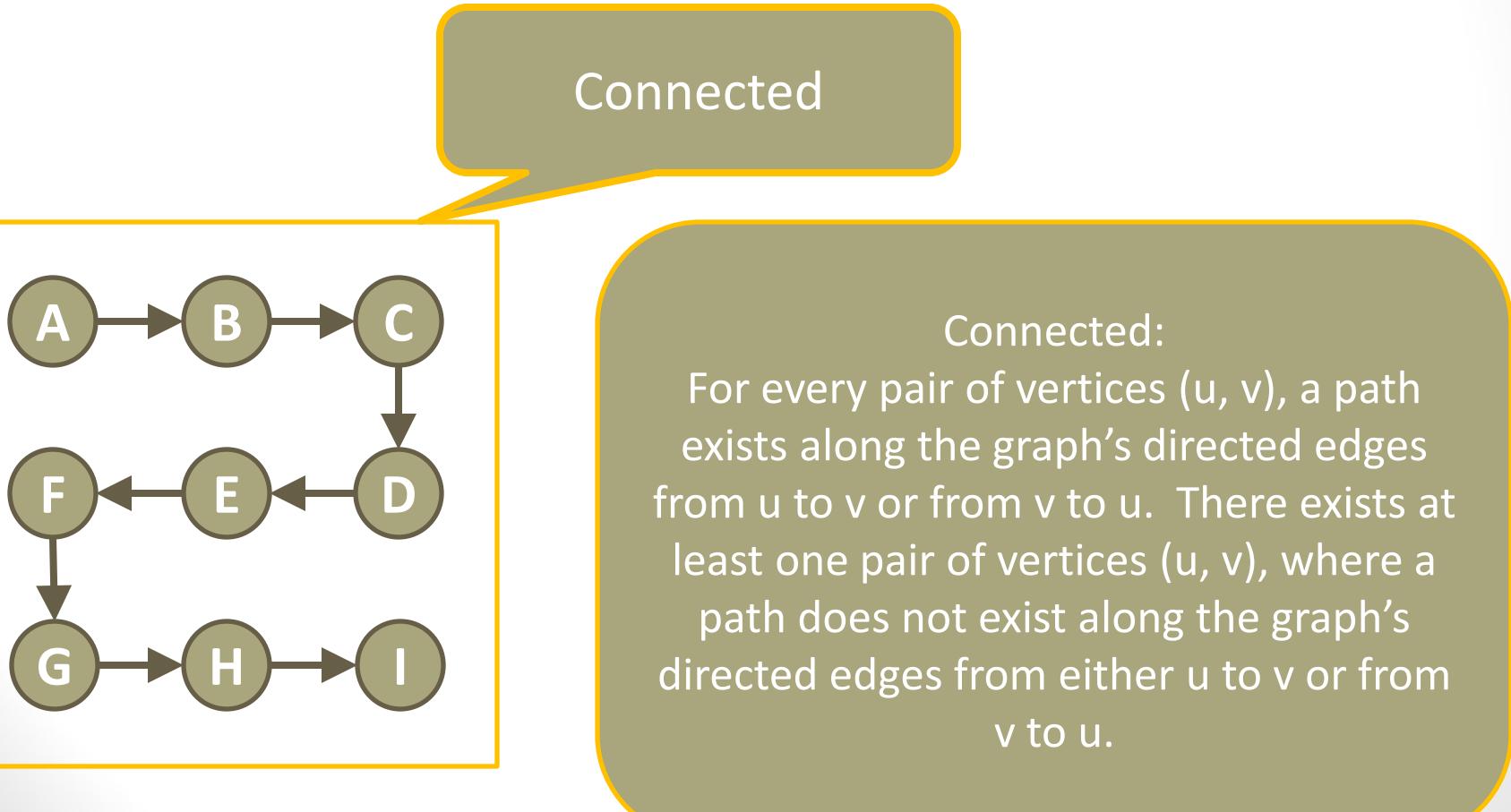


From	To
B	A
B	C
C	D
D	E
E	F
F	G
G	H
H	I

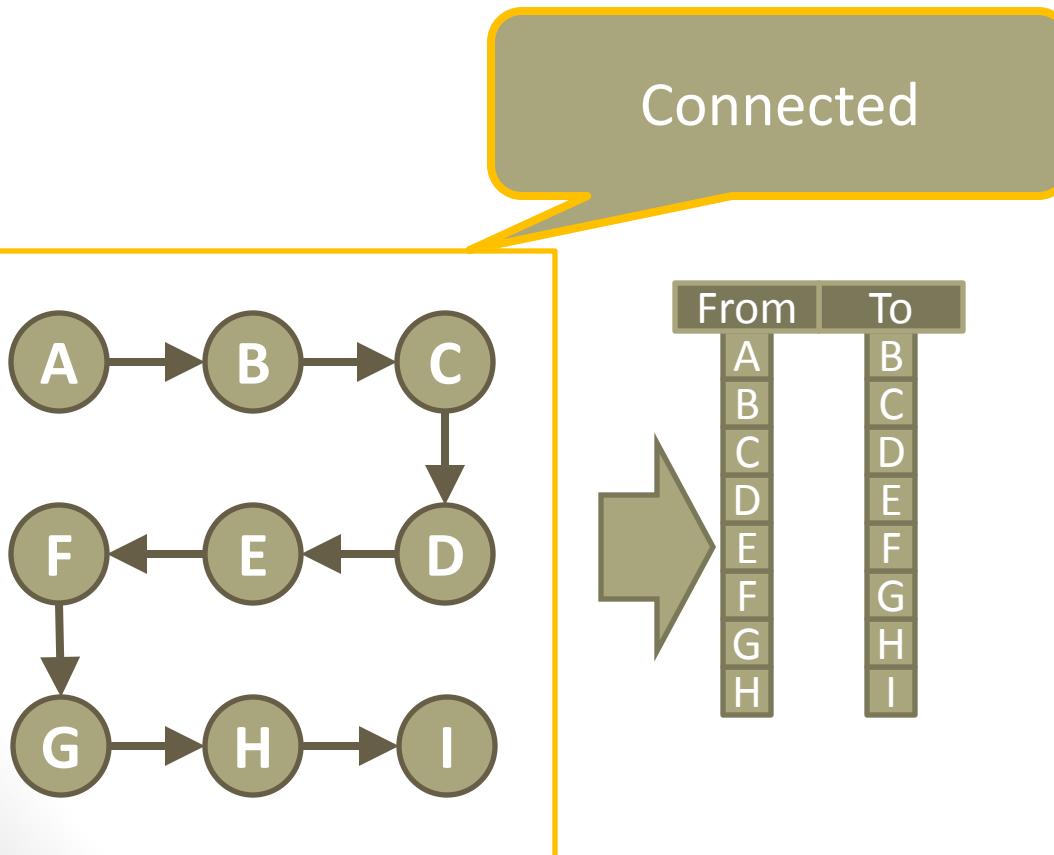
# Graph Data: Connectedness and Density (8)



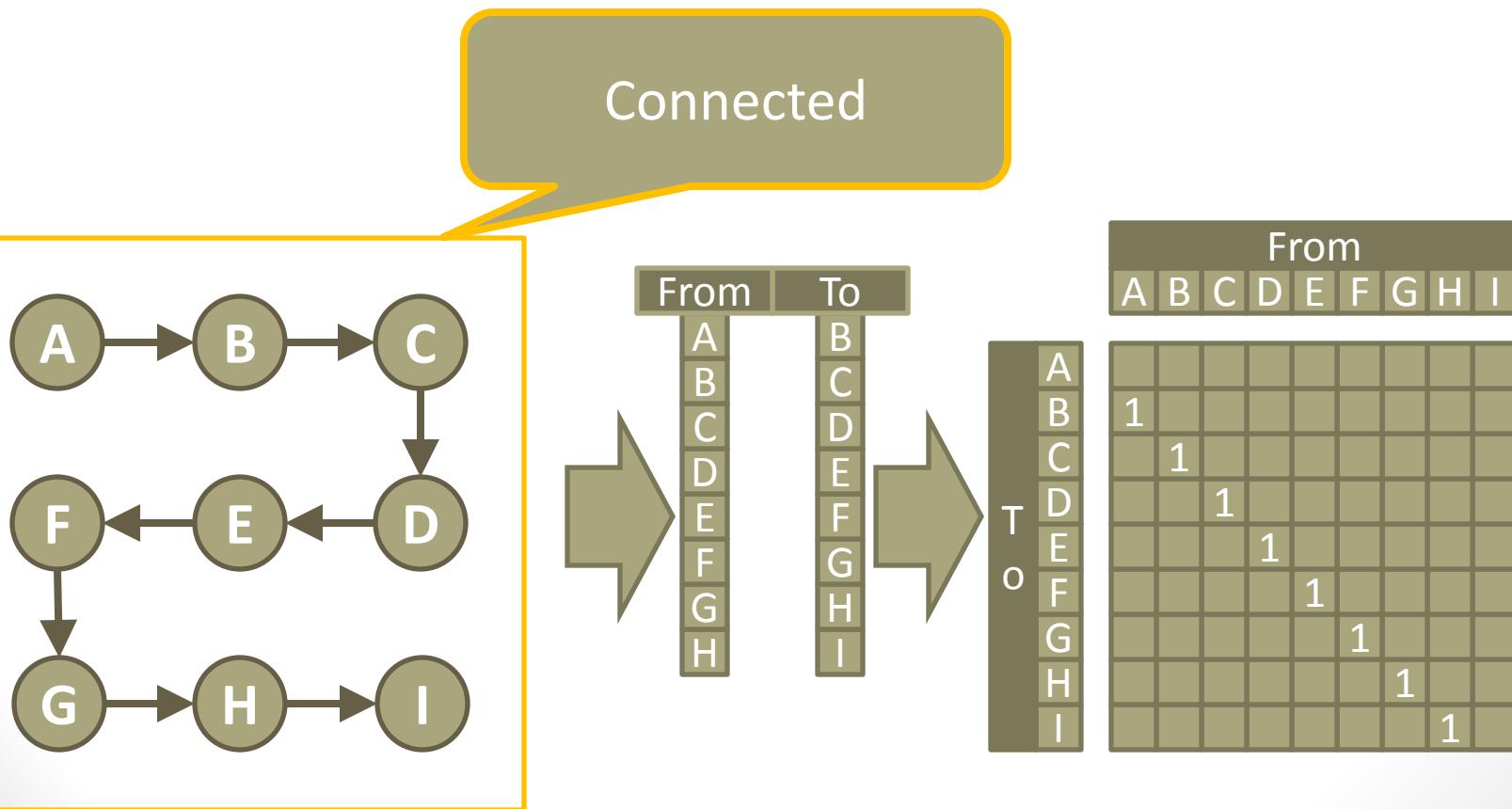
# Graph Data: Connectedness and Density (9)



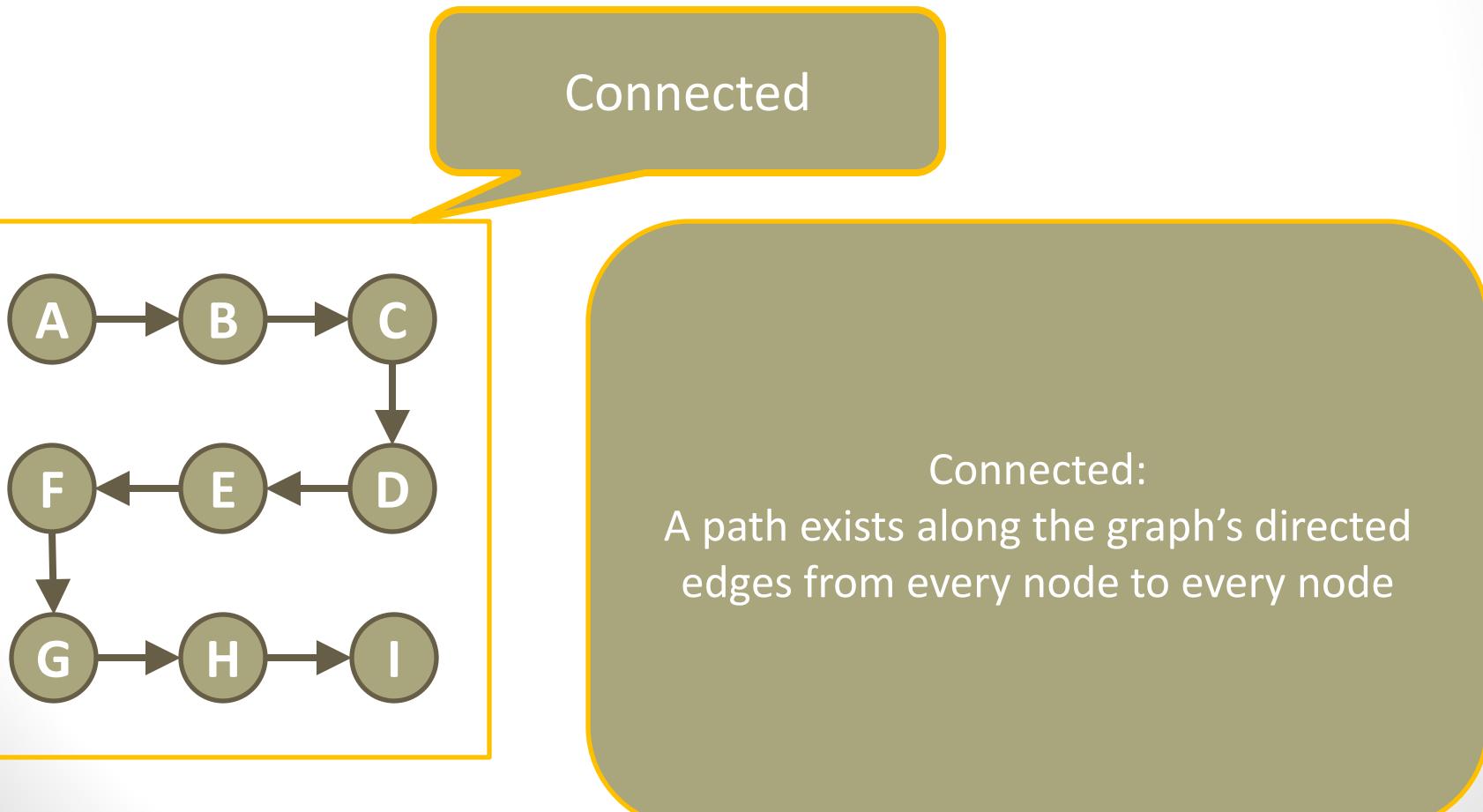
# Graph Data: Connectedness and Density (10)



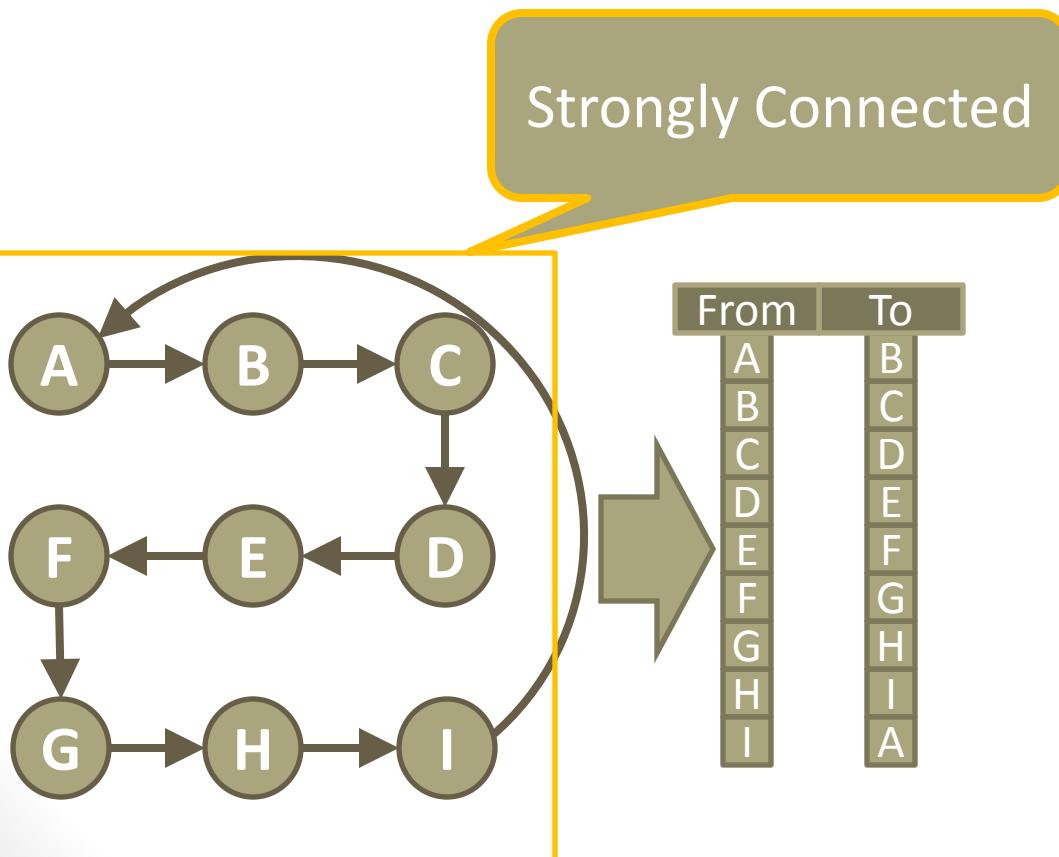
# Graph Data: Connectedness and Density (11)



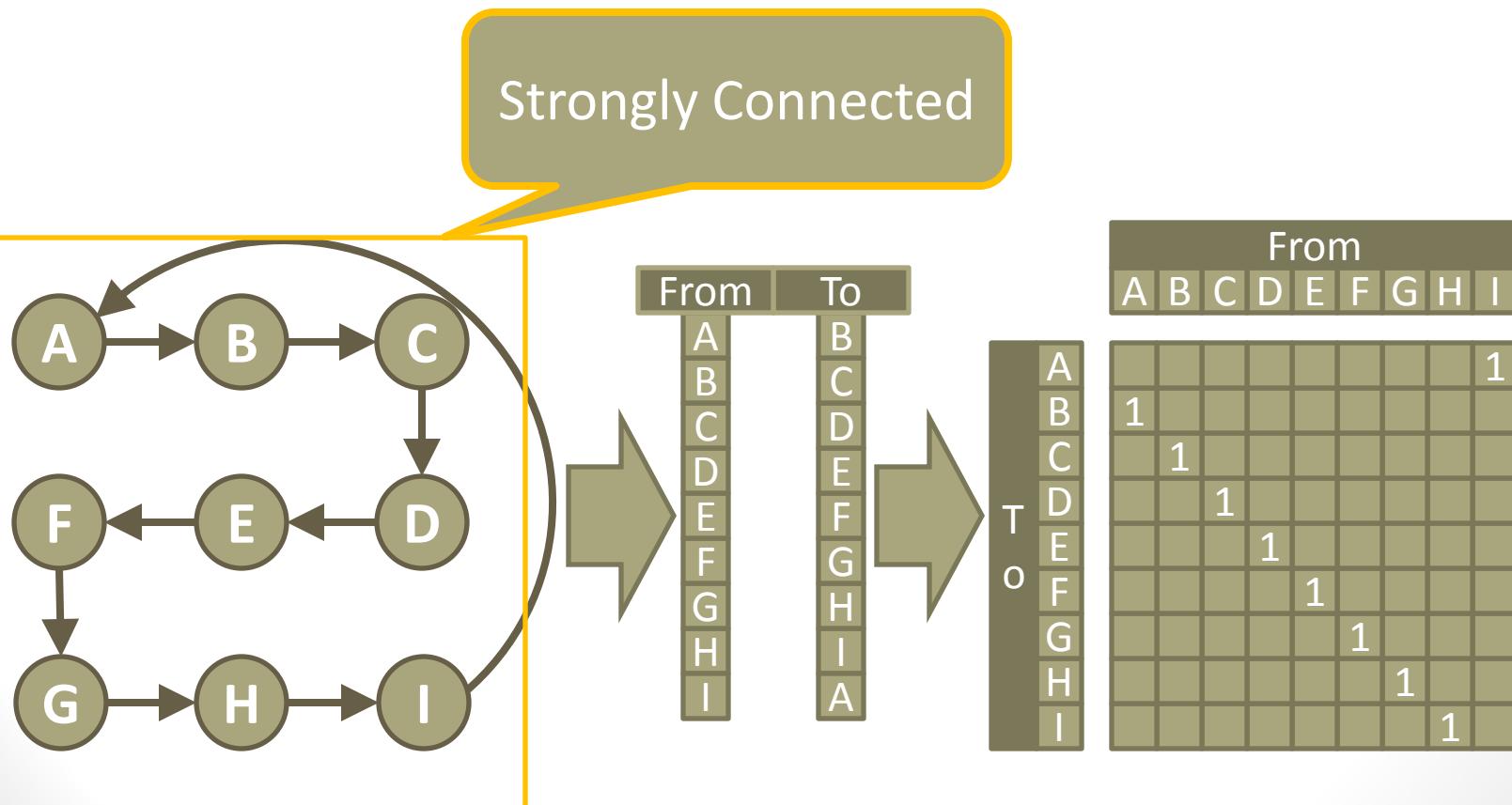
# Graph Data: Connectedness and Density (12)



# Graph Data: Connectedness and Density (13)

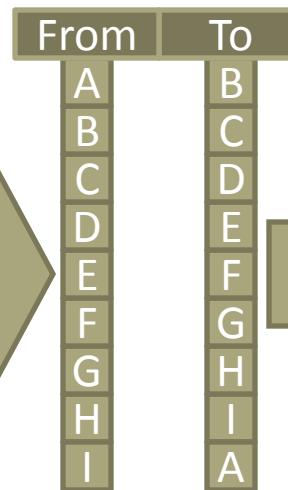
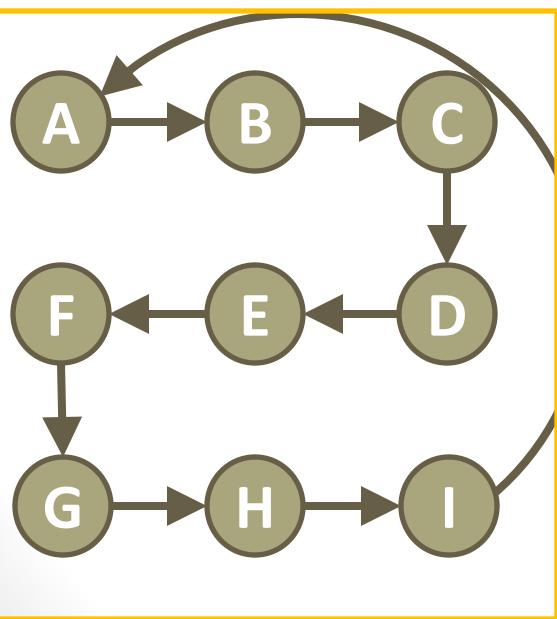


# Graph Data: Connectedness and Density (14)



# Graph Data: Connectedness and Density (15)

Sparse: Mean row or column sum is close to 1

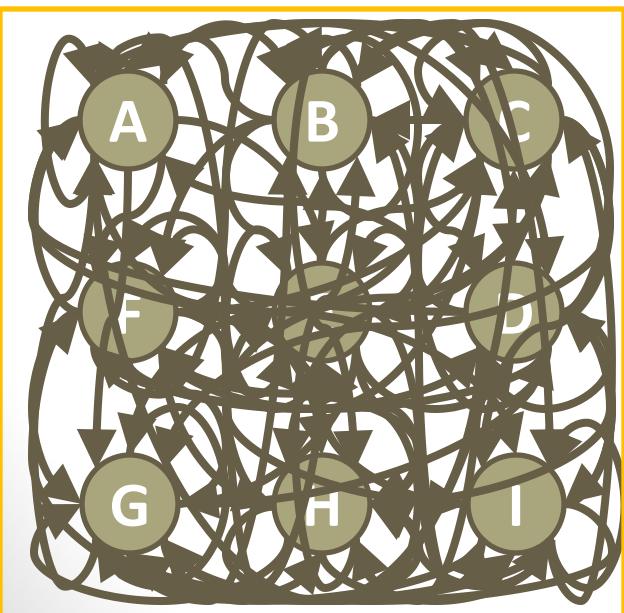


From	A	B	C	D	E	F	G	H	I	Sum:
A	1									1
B		1								1
C			1							1
D				1						1
E					1					1
F						1				1
G							1			1
H								1		1
I									1	1

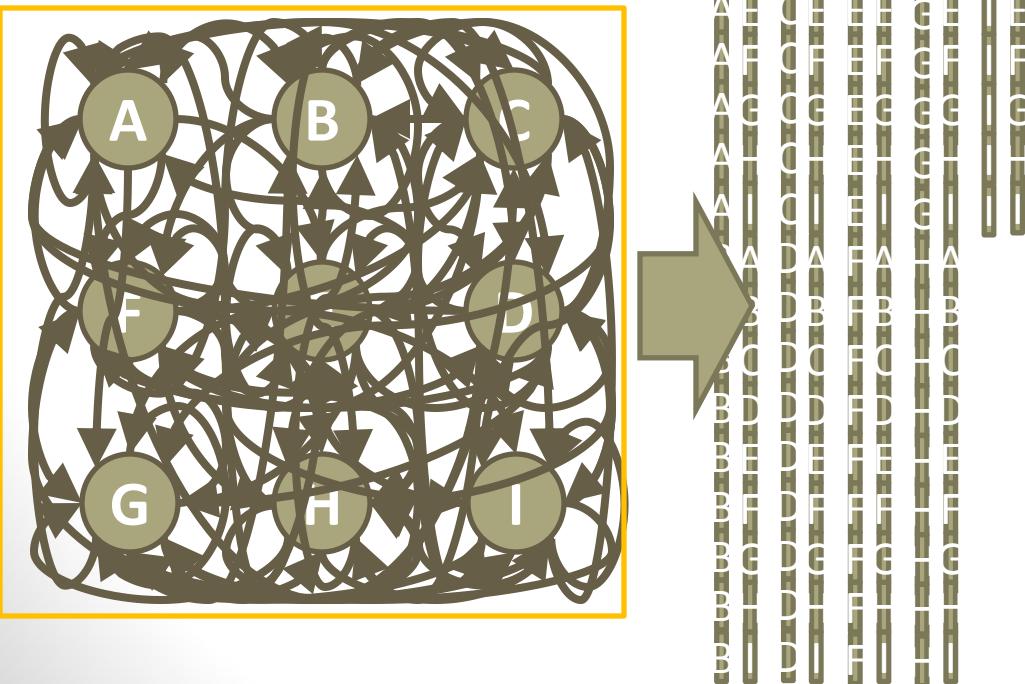
Sum: [1 1 1 1 1 1 1 1 1]

[132]

# Graph Data: Connectedness and Density (16)



# Graph Data: Connectedness and Density (17)



# Graph Data: Connectedness and Density (18)

From	To
A	A
A	B
A	C
A	D
A	E
A	F
A	G
A	H
A	I
B	A
B	B
B	C
B	D
B	E
B	F
B	G
B	H
B	I

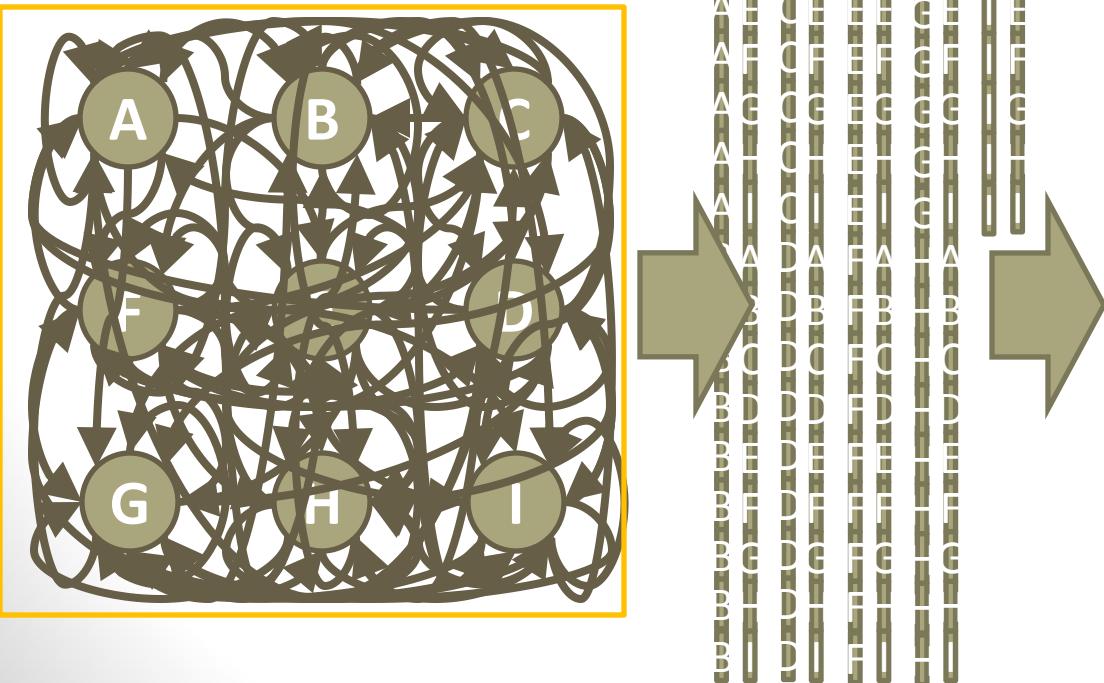
From	To
C	A
C	B
C	C
C	D
C	E
C	F
C	G
C	H
C	I
D	A
D	B
D	C
D	D
D	E
D	F
D	G
D	H
D	I

From	To
E	A
E	B
E	C
E	D
E	E
E	F
E	G
E	H
F	I
F	A
F	B
F	C
F	D
F	E
F	F
F	G
F	H
F	I

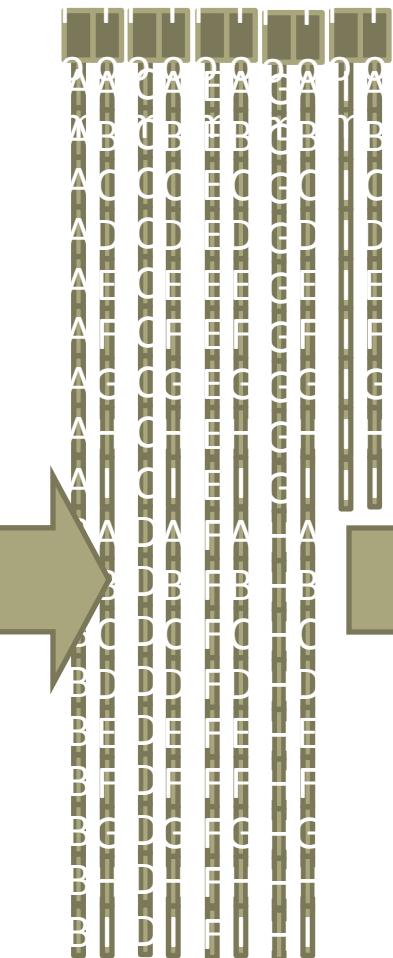
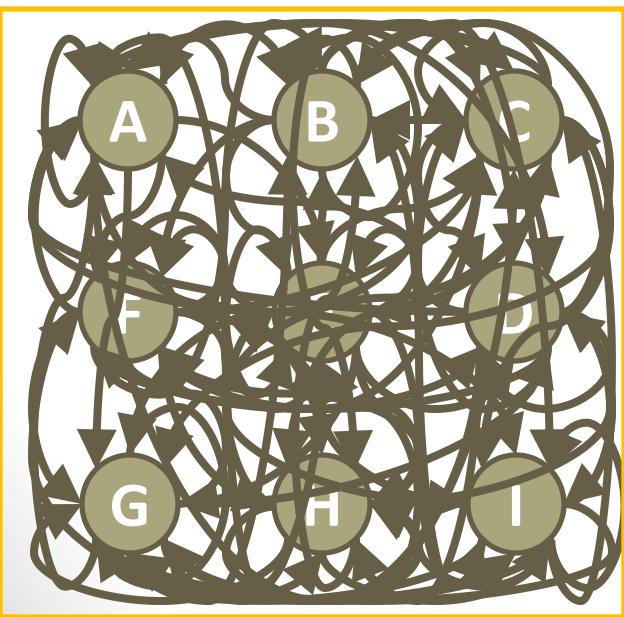
From	To
G	A
G	B
G	C
G	D
G	E
G	F
G	G
G	H
H	I
H	A
H	B
H	C
H	D
H	E
H	F
H	G
H	H
H	I

From	To
I	A
I	B
I	C
I	D
I	E
I	F
I	G
I	H
I	I

# Graph Data: Connectedness and Density (19)

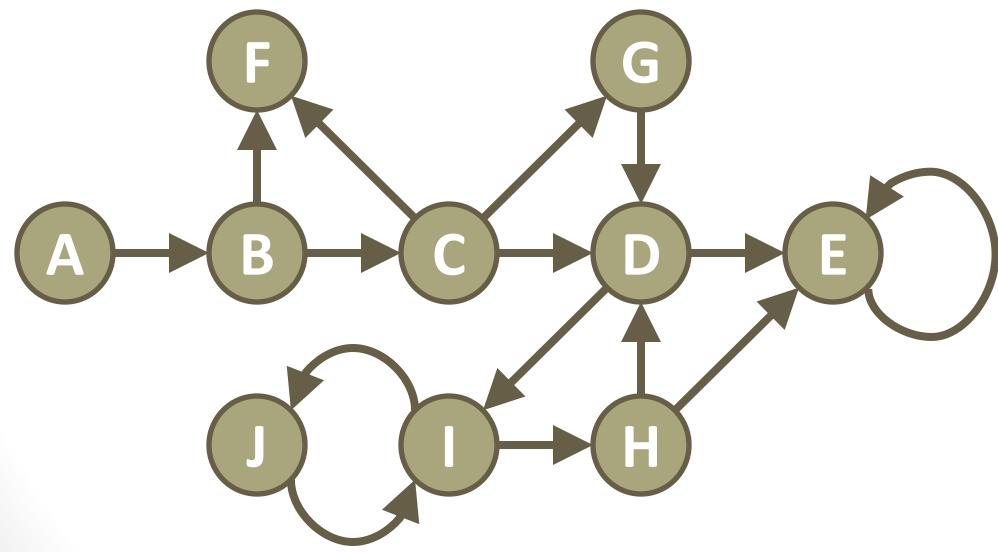


# Graph Data: Connectedness and Density (20)



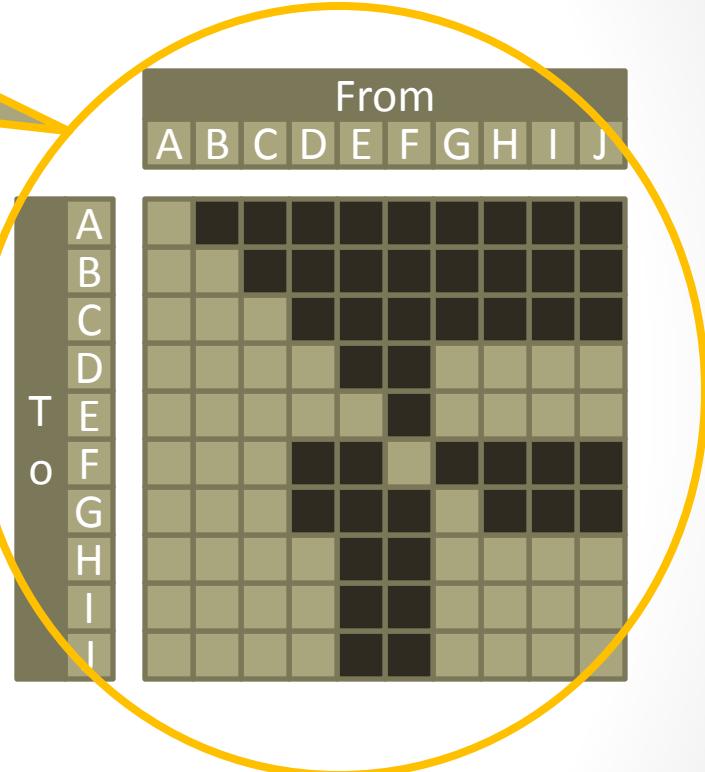
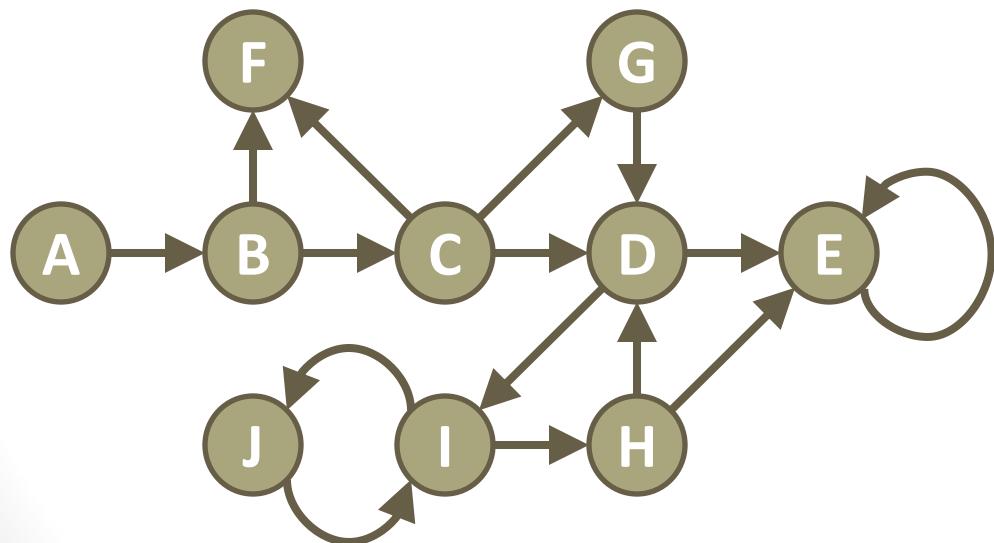
Dense: Mean row or column sum is close to number of nodes

# Graph Data: Allowed Paths, Distance, and Diameter (0)

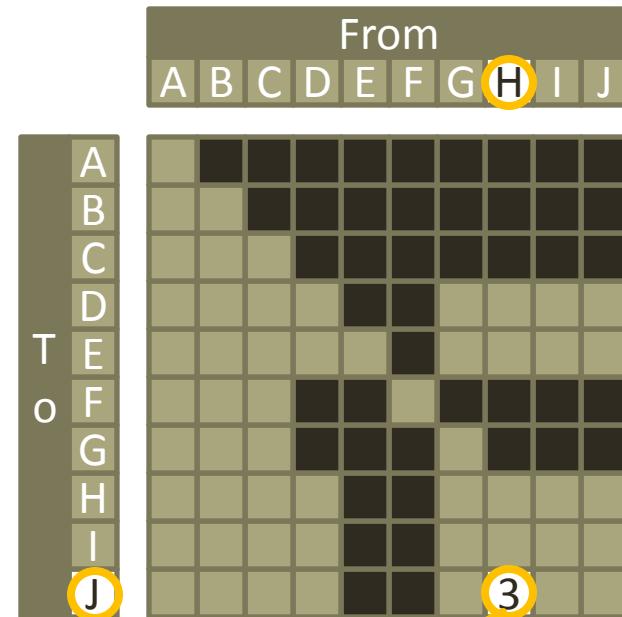
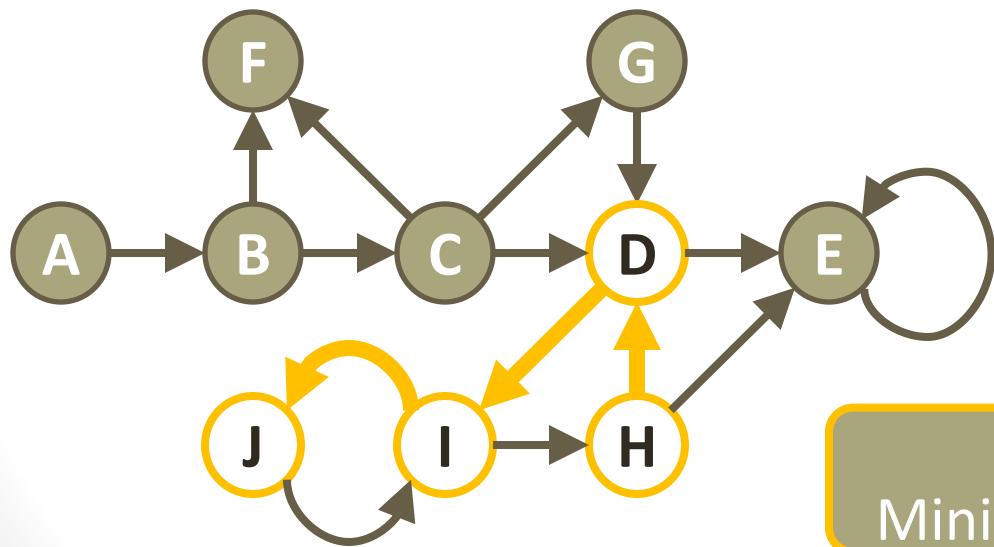


# Graph Data: Allowed Paths, Distance, and Diameter (1)

Allowed Paths



# Graph Data: Allowed Paths, Distance, and Diameter (2)

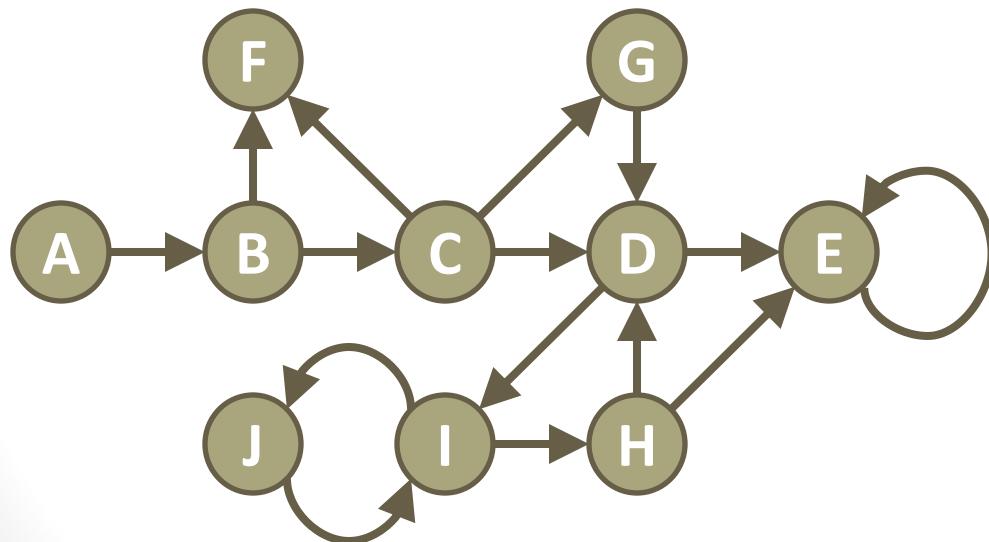


Path Distance:  
Minimum Distance from H to J

# Graph Data: Allowed Paths, Distance, and Diameter (3)

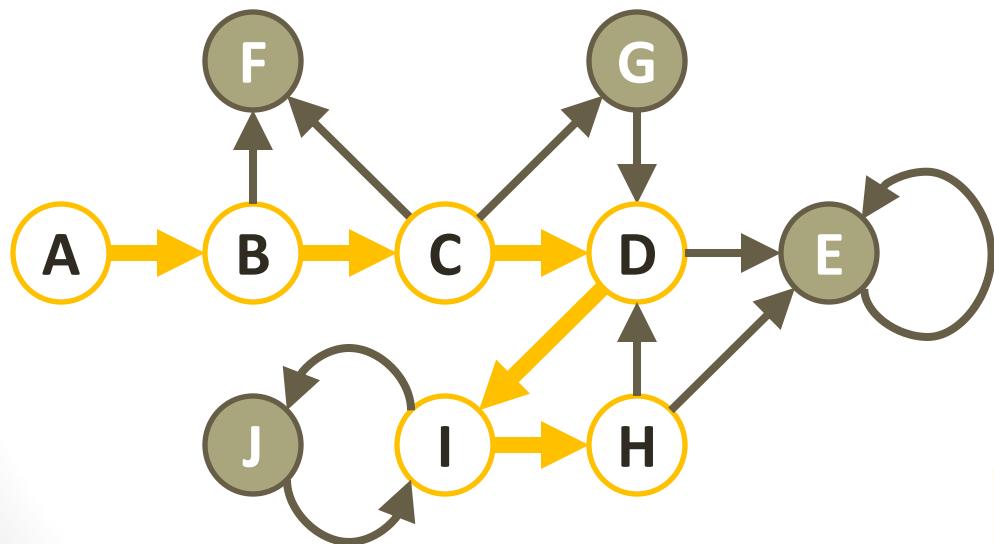
Distance Matrix:

Minimum distance between any two vertices along a directed edge



		From									
		A	B	C	D	E	F	G	H	I	J
To	A	0									
	B	1	0								
	C	2	1	0							
	D	3	2	1	0				1	1	2
	E	4	3	2	1	0			2	1	2
	F	2	1	1			0				
	G	4	2	1				0			
	H	5	4	3	2				3	0	1
	I	4	3	2	1				2	2	0
	J	5	4	3	2				3	3	1

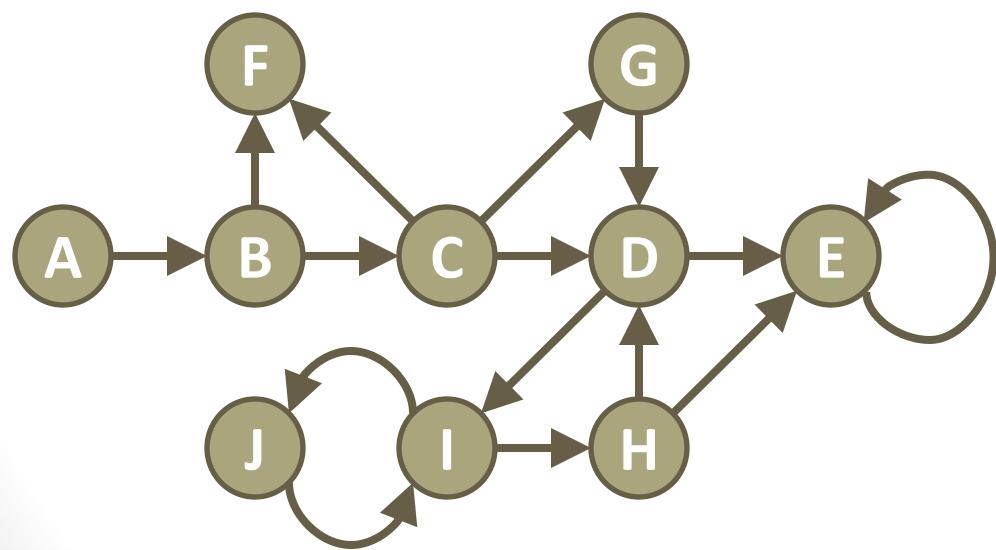
# Graph Data: Allowed Paths, Distance, and Diameter (4)



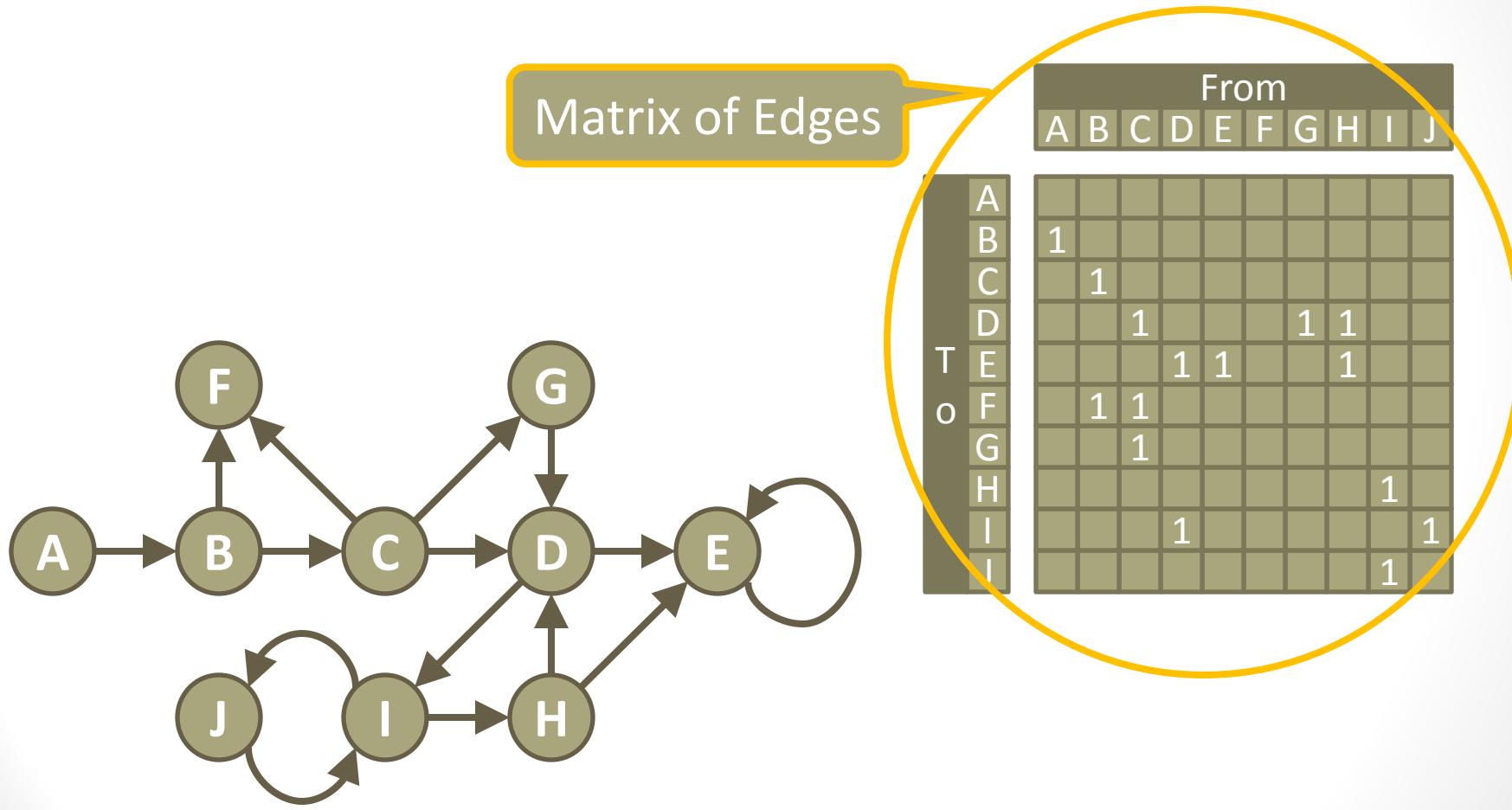
From									
A	B	C	D	E	F	G	H	I	J
A	0								
B	1	0							
C	2	1	0						
D	3	2	1	0				1	1
E	4	3	2	1	0			2	1
F	2	1	1		0			2	3
G	4	2	1			0			
H	4	3	2	2		3	0	1	2
I	4	3	2	1		2	2	0	1
J	4	3	2	2		3	3	1	0

Diameter:  
Max Path Distance

# Graph Data: Order and Size (0)

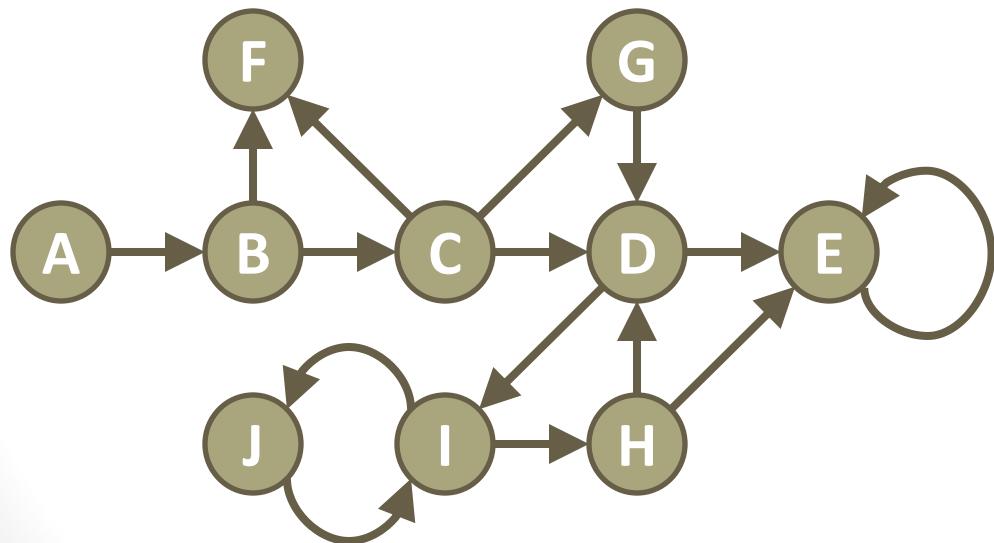


# Graph Data: Order and Size (1)



# Graph Data: Order and Size (2)

Order of Graph: 10

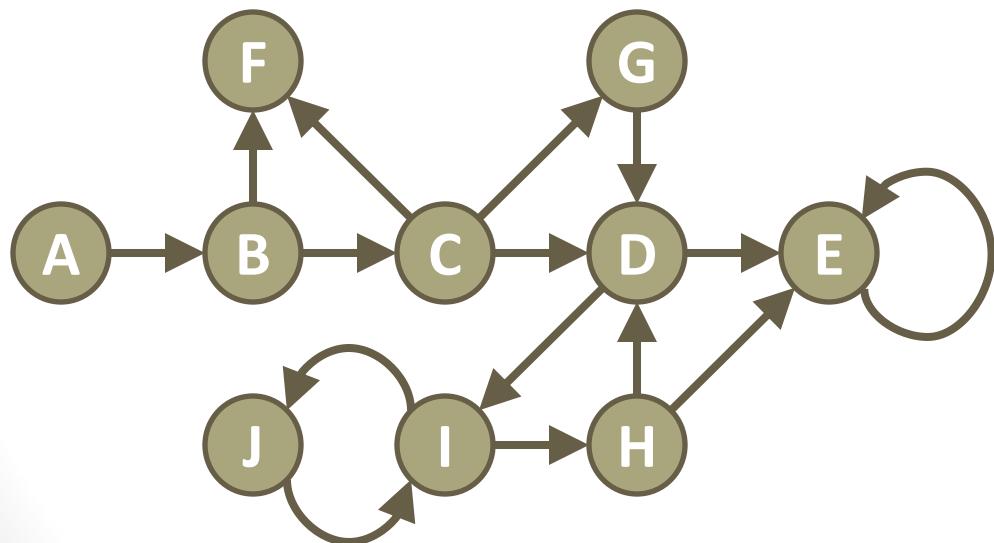


From									
A	B	C	D	E	F	G	H	I	J
T	1	1	1	1	1	1	1	1	1
C		1	1	1	1	1	1	1	1
O			1	1	1	1	1	1	1
R				1	1	1	1	1	1
E					1	1	1	1	1
N						1	1	1	1
S							1	1	1
P								1	1
L									1

Order of Graph:  
Number of Nodes

# Graph Data: Order and Size (3)

Order of Graph: 10  
Size of Graph: 15



From									
A	B	C	D	E	F	G	H	I	J
T	1	1	1	1	1	1	1	1	1
O	1	1	1	1	1	1	1	1	1
R	1	1	1	1	1	1	1	1	1
E	1	1	1	1	1	1	1	1	1
N	1	1	1	1	1	1	1	1	1
S	1	1	1	1	1	1	1	1	1
P	1	1	1	1	1	1	1	1	1
U	1	1	1	1	1	1	1	1	1

Size of Graph:  
Number of Edges

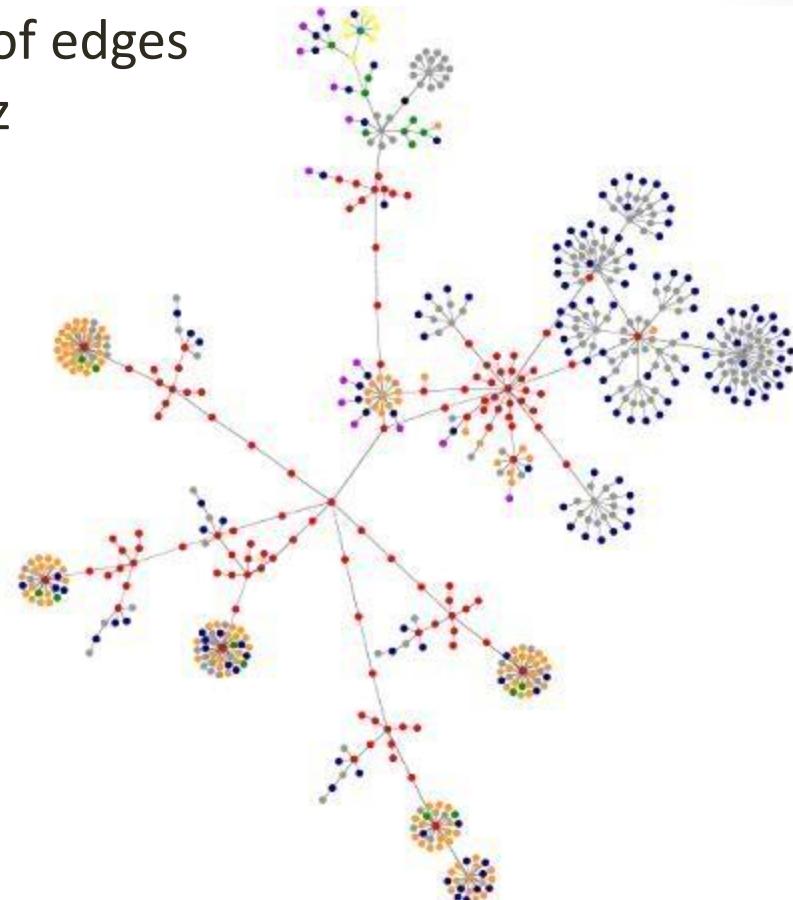
(146)

# Graph Data: Links (0)

- Graphs In General:
  - [http://en.wikipedia.org/wiki/Graph\\_theory](http://en.wikipedia.org/wiki/Graph_theory)
  - [https://en.wikipedia.org/wiki/Graph\\_\(discrete\\_mathematics\)](https://en.wikipedia.org/wiki/Graph_(discrete_mathematics))
  - [https://en.wikipedia.org/wiki/Graph\\_\(abstract\\_data\\_type\)](https://en.wikipedia.org/wiki/Graph_(abstract_data_type))
  - [https://en.wikipedia.org/wiki/List\\_of\\_graph\\_theory\\_topics](https://en.wikipedia.org/wiki/List_of_graph_theory_topics)
- Graph Diameter
  - <http://people.hofstra.edu/geotrans/eng/methods/diameter1.html>
- Distance
  - [http://en.wikipedia.org/wiki/Distance\\_\(graph\\_theory\)](http://en.wikipedia.org/wiki/Distance_(graph_theory))
- Resistance
  - [http://en.wikipedia.org/wiki/Resistance\\_distance](http://en.wikipedia.org/wiki/Resistance_distance)
- Centrality
  - [http://en.wikipedia.org/wiki/Betweenness#Betweenness\\_centrality](http://en.wikipedia.org/wiki/Betweenness#Betweenness_centrality)
- Connectivity
  - [http://en.wikipedia.org/wiki/Connectivity\\_\(graph\\_theory\)](http://en.wikipedia.org/wiki/Connectivity_(graph_theory))
- Hypergraph: One edge can connect many nodes
  - <http://en.wikipedia.org/wiki/Hypergraph>
- Degree Distribution: Histogram of in-degrees, out-degrees, in-degrees normalized by out-degrees
- Popularity (Google Matrix)
- Density: #edges/(#nodes<sup>2</sup>-#nodes)
- Clustering Coefficient
  - [http://en.wikipedia.org/wiki/Clustering\\_coefficient](http://en.wikipedia.org/wiki/Clustering_coefficient)

# Quiz Graph Data (Measures)

- The presented slide contains a list of edges that describes the graph in the quiz



# Introduction to Graph Data

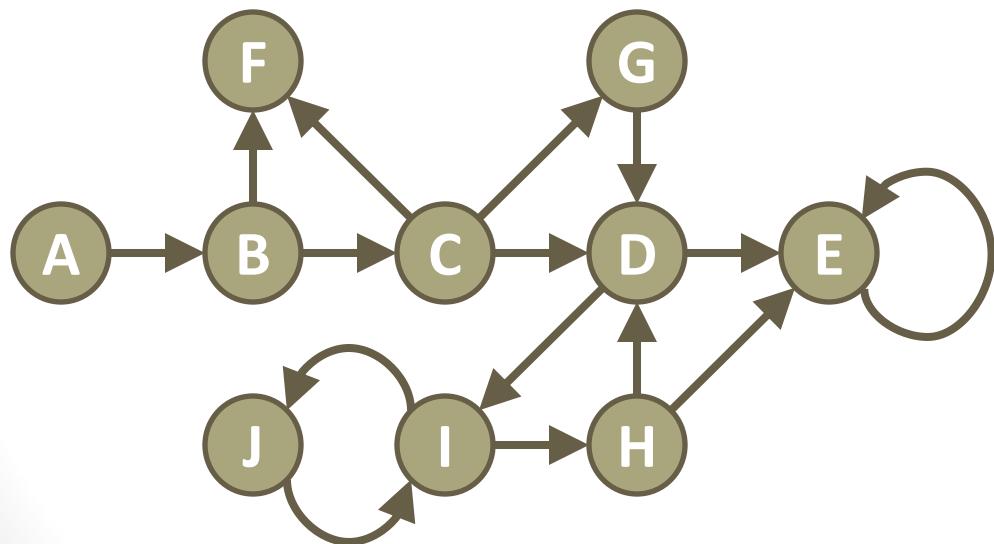
(150)

# Graph Data: Google Matrix

# Google Matrix

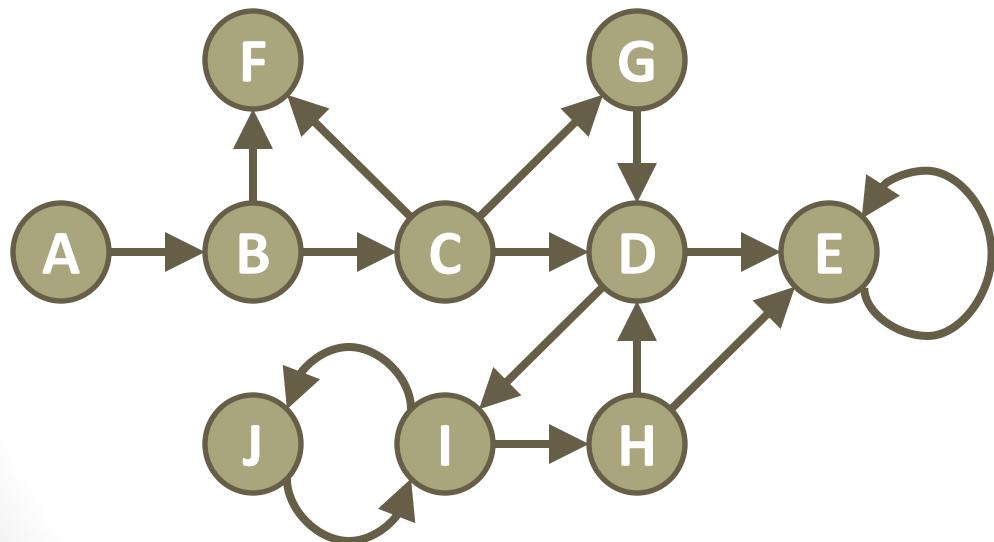
- [http://en.wikipedia.org/wiki/Google\\_matrix](http://en.wikipedia.org/wiki/Google_matrix)

# Graph Data: Popularity (0)



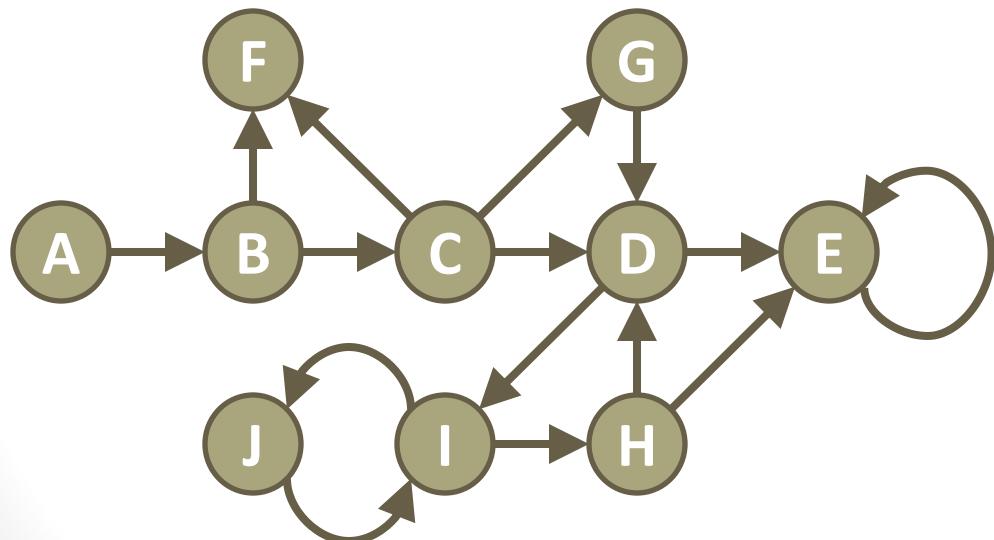
# Graph Data: Popularity (1)

# In-degree of a node: Popularity



# Graph Data: Popularity (2)

In-degree of a node: Popularity



From										Sum:	
A	B	C	D	E	F	G	H	I	J		
To	1									0	
		1								1	
			1							1	
				1						3	
					1					3	
						1				2	
							1			1	
								1		1	
									1	2	
										1	
Sum:										Sum:	
1 2 3 2 1 0 1 2 2 1										Sum:	

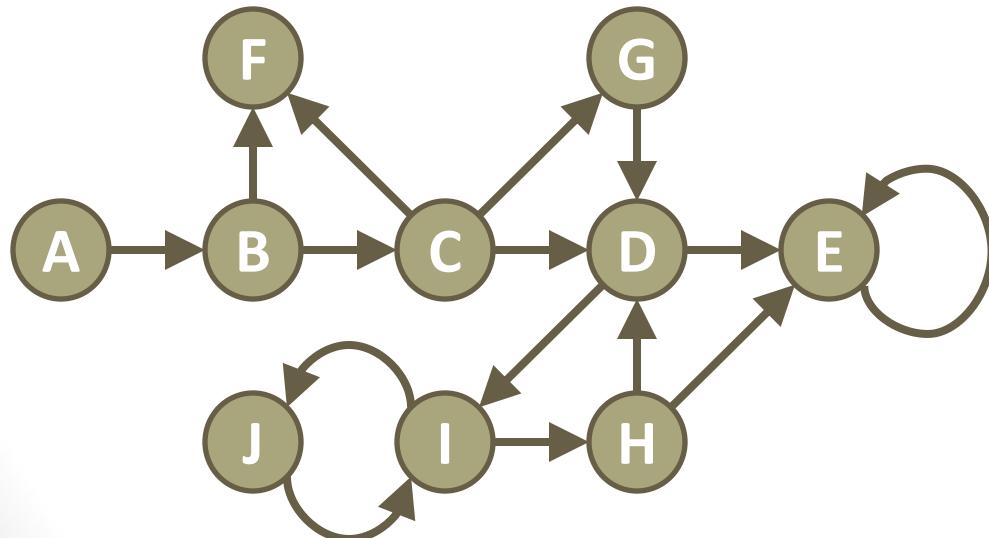
Out-degree of a node:  
Gregariousness

# Quiz Graph Data (Popularity)

- Graph Data (Popularity)
- You need to view the projected slide to answer the questions

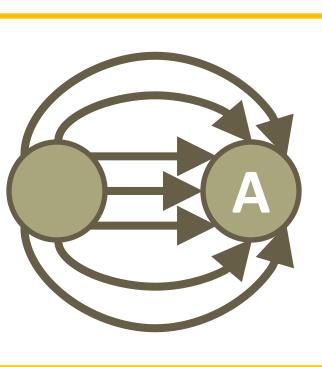
# Graph Data: Popularity (3)

But, not every link is equally important.  
Otherwise, you could create a website with  
a thousand links that all point to your site.

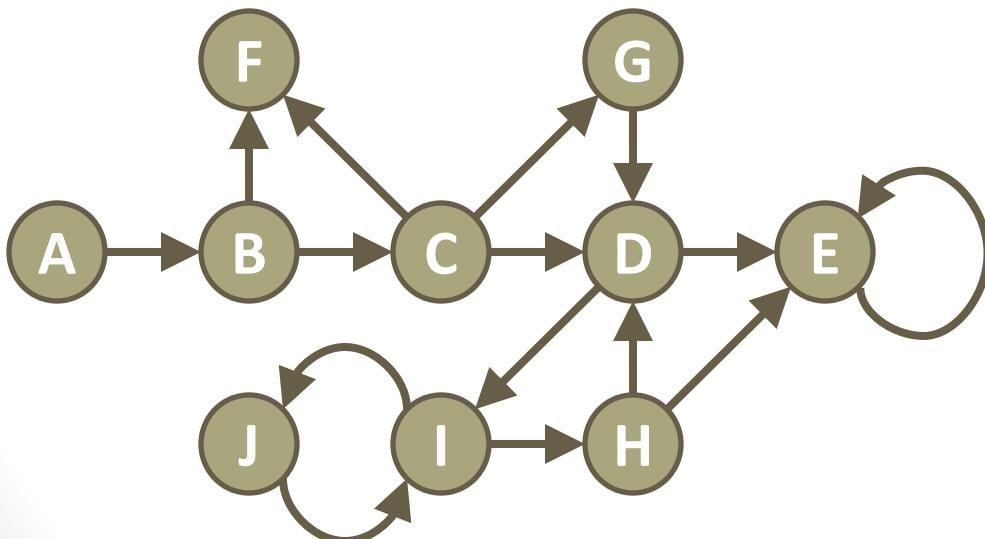


From										Sum:
A	B	C	D	E	F	G	H	I	J	
To	1									0
		1								1
			1							1
				1						3
					1					3
						1				2
							1			1
								1		1
									1	2
										1
Sum:										1 2 3 2 1 0 1 2 2 1

# Graph Data: Popularity (4)

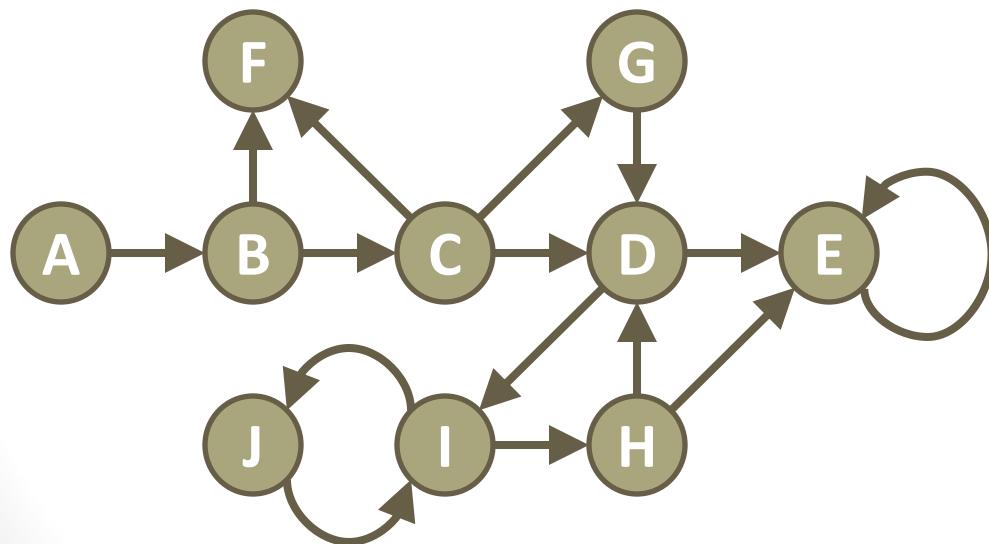


One very gregarious node  
could make another node  
very popular



From										Sum:	
	A	B	C	D	E	F	G	H	I	J	
A											0
B	1										1
C		1									1
D			1					1	1		3
T				1	1					1	3
E					1	1				1	3
O	F		1	1							2
G				1							1
H									1		1
I					1					1	2
J									1		1
Sum:											
1	2	3	2	1	0	1	2	2	1		

# Graph Data: Popularity (5)



From										Sum:
A	B	C	D	E	F	G	H	I	J	
To	1									0
		1								1
			1							1
				1						3
					1					3
						1				2
							1			1
								1		1
									1	2
										1

Sum: 1 2 3 2 1 0 1 2 2 1

The amount of popularity that a single node can send out needs to be tempered by Gregariousness

# Graph Data: Popularity (6)

This matrix needs some  
adjustments

From										Sum:
A	B	C	D	E	F	G	H	I	J	
To	1									0
		1								1
			1							1
				1						3
					1	1				3
						1	1			2
							1			1
								1		2
									1	1
										1

Sum: 1 2 3 2 1 0 1 2 2 1

# Graph Data: Popularity (7)

	From										
	A	B	C	D	E	F	G	H	I	J	
A											0
B	1										1
C		1									1
D			1					1	1		3
E				1	1				1		3
F		1	1								2
G			1								1
H				1					1		1
I					1					1	2
J								1			1
sum:	1	2	3	2	1	0	1	2	2	1	

# Graph Data: Popularity (8)

	From										
	A	B	C	D	E	F	G	H	I	J	
A											sum: 0
B	1										1.00
C		1/2									0.50
D			1/3					1	1/2		1.83
E				1/2	1				1/2		2.00
F		1/2	1/3								0.83
G			1/3								0.33
H				1/2						1/2	0.50
I					1/2						1.50
J									1/2		0.50
sum:	1	1	1	1	1	0	1	1	1	1	

# Graph Data: Popularity (9)

	From											
	A	B	C	D	E	F	G	H	I	J	sum:	Popularity:
A											0	0
B	1										1.00	0.111
C		1/2									0.50	0.056
D			1/3				1	1/2			1.83	0.204
E				1/2	1			1/2			2.00	0.222
F		1/2	1/3								0.83	0.092
G			1/3								0.33	0.037
H									1/2		0.50	0.056
I				1/2						1	1.50	0.167
J									1/2		0.50	0.056
sum:												
	1	1	1	1	1	0	1	1	1	1		

# Graph Data: Popularity (10)

- Rules of a Popularity Network
  - Rule 0: Popularity transmits from one node to a downstream connected node.
  - Rule 1: Popularity transmissions from multiple nodes to a single node are additive.
  - Rule 2: A more popular node transmits more popularity.
  - Rule 3: A node splits up its popularity among its recipients.
- Continue with **EigenVectorOfGraphMatrix.R**
- Solve (or Iterate until convergence)
  - Modified Popularity is the Popularity divided by Gregariousness (Out-Degree)
  - Popularity of a website is the sum of the Modified Popularities that point to the website.

# Graph Data: Google Matrix

# Assignment

(167)

# Assignment

1. Comment in our LinkedIn Group on a discussion started by yourself or a fellow student about statistics. Copy the post to a text file called post.txt. This assignment item is similar to last week's assignment.
2. Create a MapReduce program that lists word lengths next to the number of times a word of that length occurs. The word lengths are sorted as texts, even though they are numbers. Specifically, modify package wordcount.solution as described in the MapReduce Lab for "Word Length Count 1". In this lab, the output keys are of type text. Save all code of the modified java file(s) and a screen shot of all the results (approx. 18 rows).
3. Create a MapReduce program that lists word lengths next to the number of times a word of that length occurs. The word lengths are sorted numerically. Specifically, modify package wordcount.solution as described in the MapReduce Lab for "Word Length Count 2". In this lab, the output keys are of type IntWritable (numeric). Save all code of the modified java file(s) and a screen shot of all the results (approx. 18 rows).
4. Submit to Canvas by Saturday 11:57 PM the following: post.txt. The code as \*.java files from "Word Length Count 1". The screenshot of the results from "Word Length Count 1". The code as \*.java files from "Word Length Count 2". The screenshot of the results from "Word Length Count 2".
5. Look through the Preview section, especially the SPARQL Exercises, prior to next week's lecture. (There will be a quiz).

# Assignment

(169)

# Introduction to Data Science

(170)