

Welcome (Again!) to MATH 4100/COMP 5360– Introduction to Data Science

Network Analysis

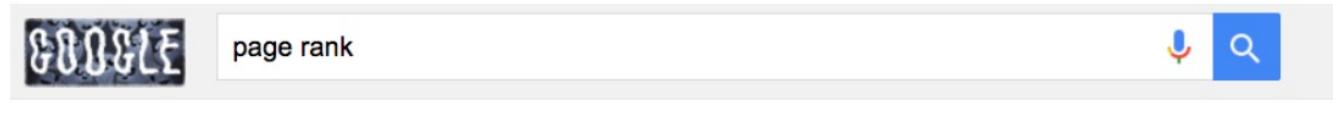
April 18, 2024

*Based in part on prior lectures from
Alex Lex and Bei Wang Phillips
+ others where noted*



Applications of Graphs

Without graphs, there would be none of these:



About 431,000,000 results (0.86 seconds)

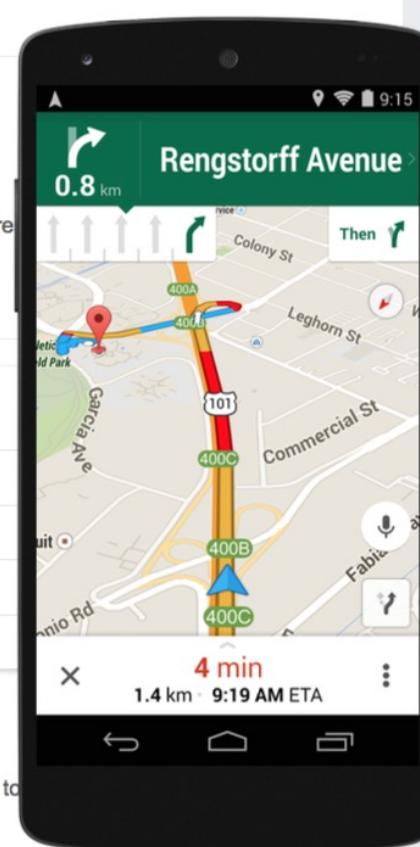
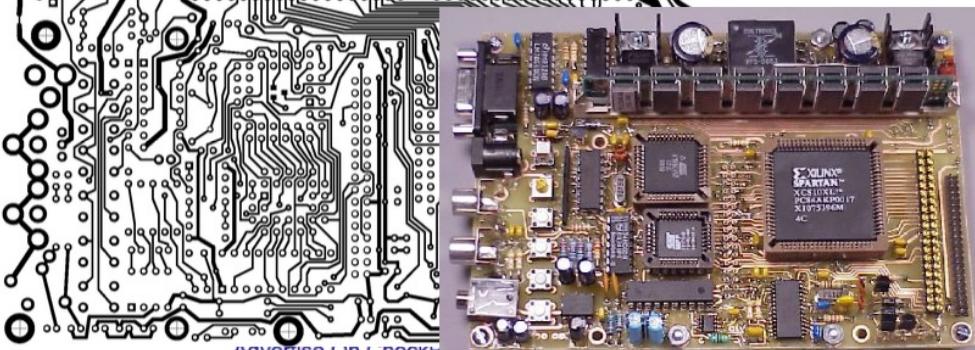
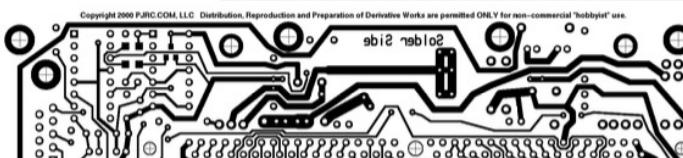
[PageRank - Wikipedia](#)

<https://en.wikipedia.org/wiki/PageRank> ▾ Wikipedia ▾

PageRank is an algorithm used by Google Search to rank websites in their search engine re-

PageRank was named after Larry Page, one of the founders of ...

[Description](#) · [History](#) · [Algorithm](#) · [Variations](#)



Follow Mark to get his public posts in your News Feed.



83,820,410 Followers



Mark Zuckerberg

October 28 at 1:26pm ·

Happy birthday, Bill! Thanks for your friendship and here of changing the world.



facebook

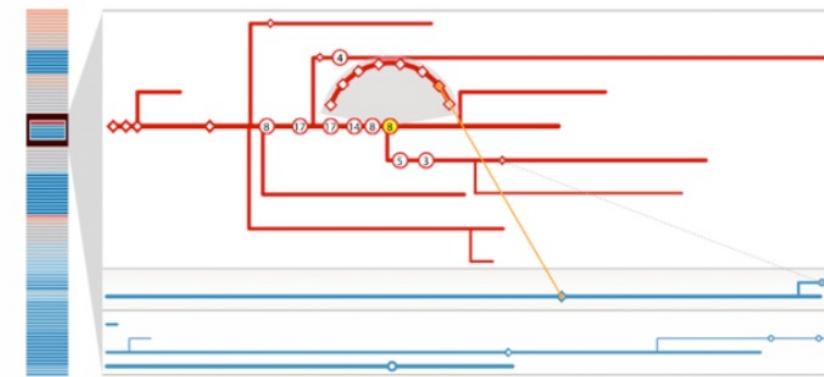
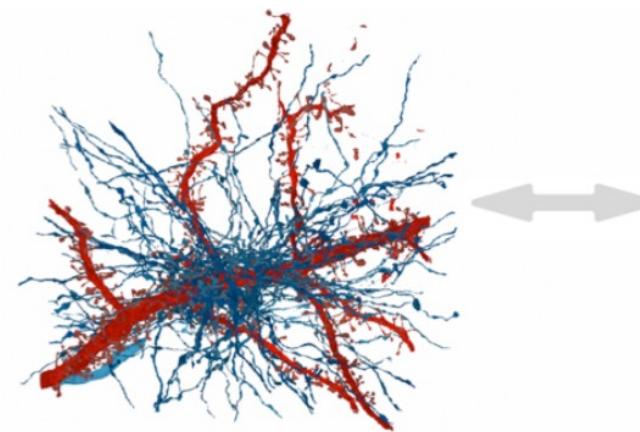
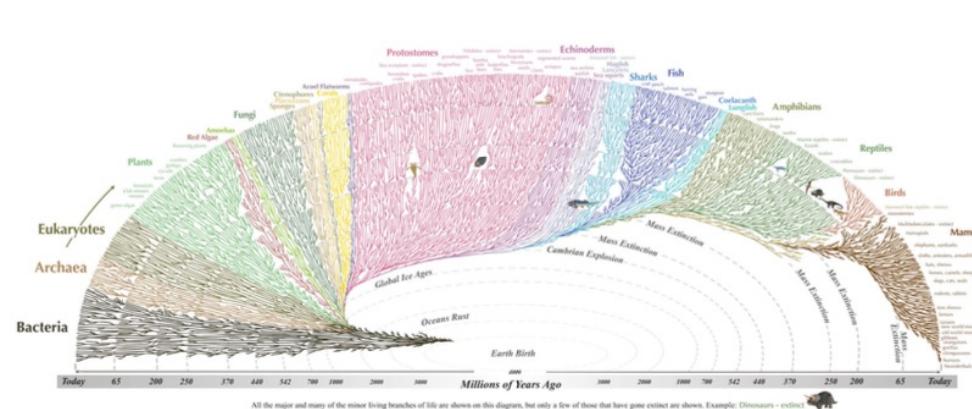
Biological Networks

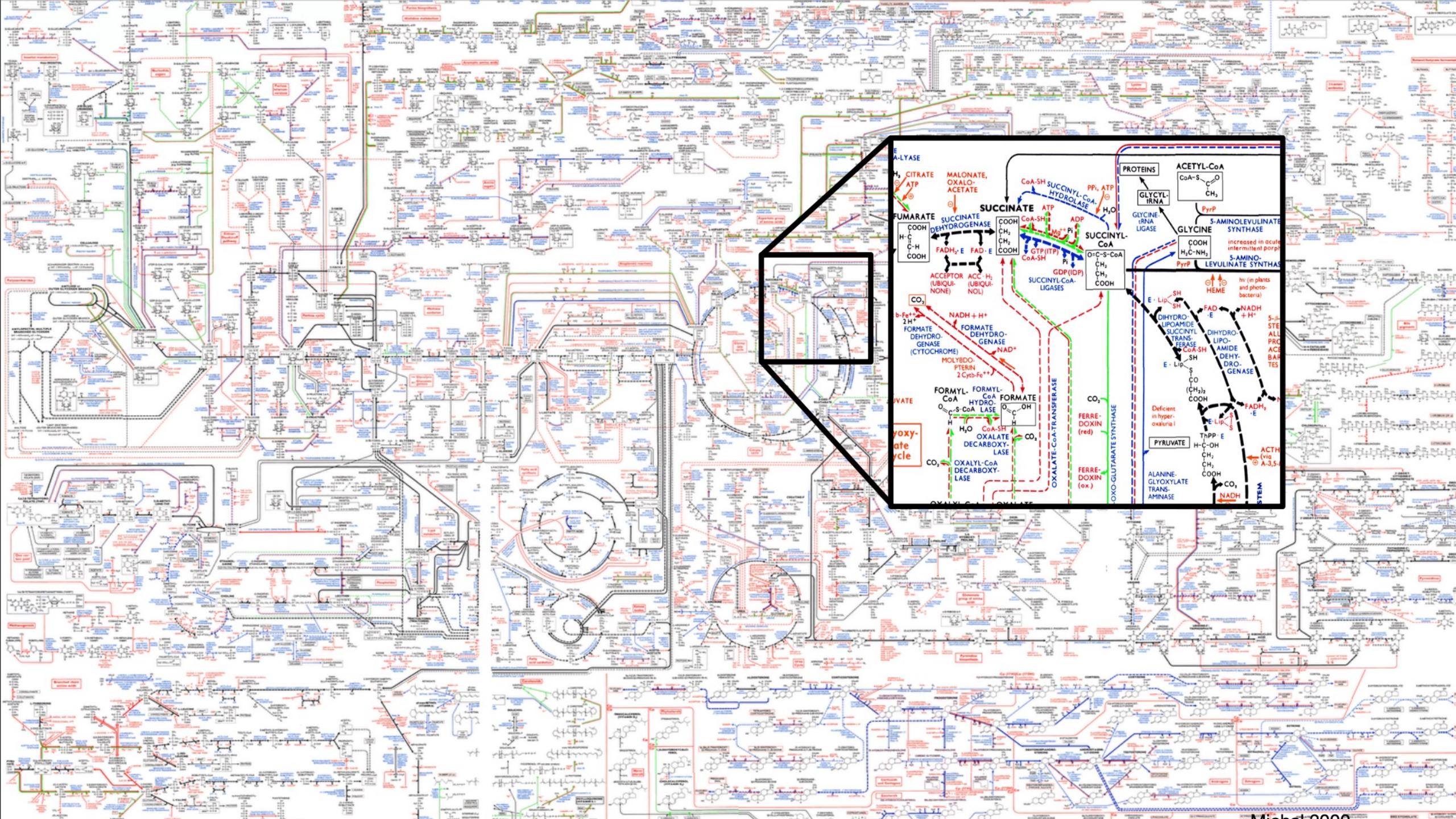
Interaction between genes, proteins and chemical products

The brain: connections between neurons

Your ancestry: the relations between you and your family

Phylogeny: the evolutionary relationships of life



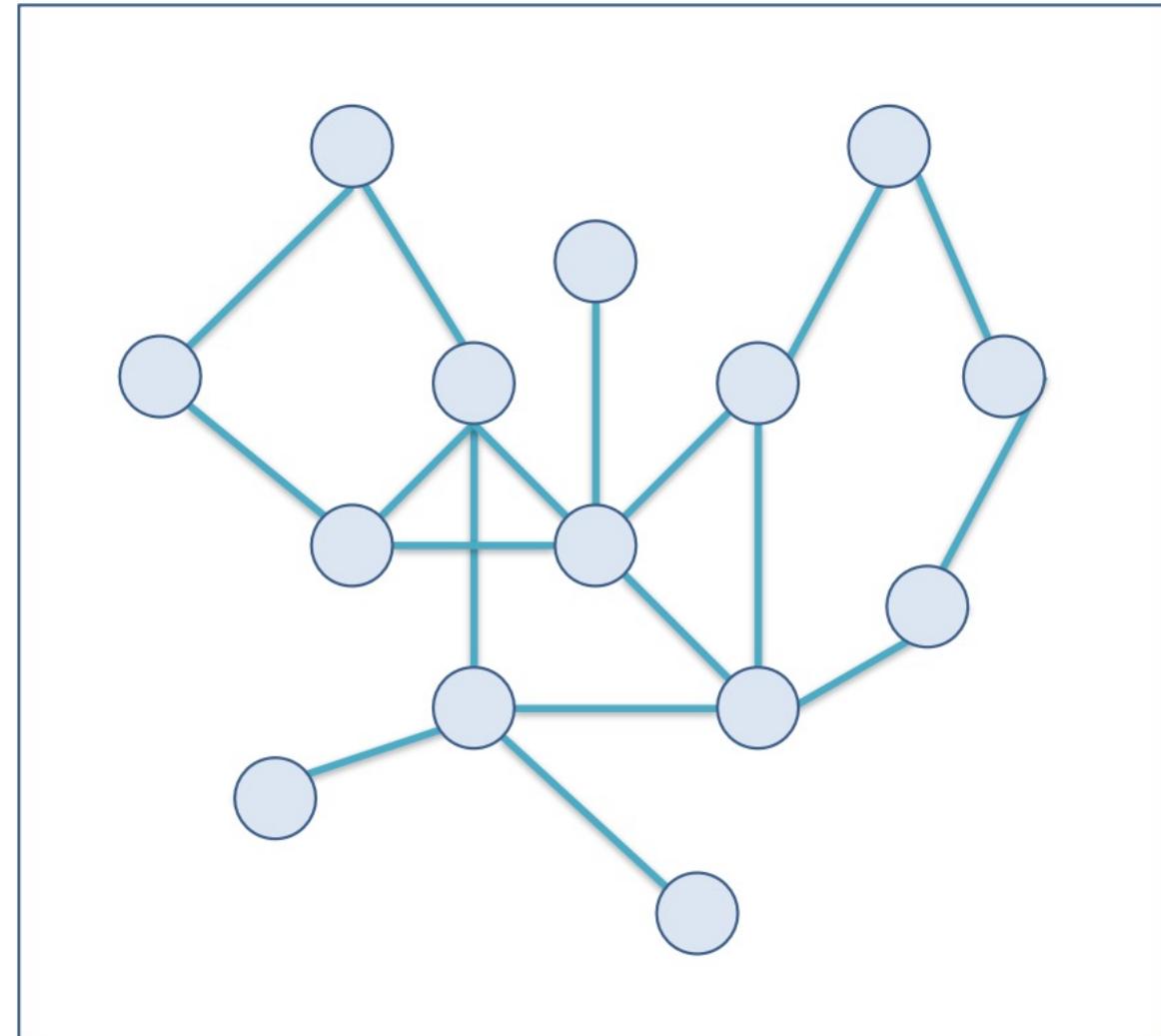


Network (Graph) Terminology

Mathematics Terminology

A graph $\mathbf{G(V,E)}$ denotes a graph with vertices **V** and edges **E**.

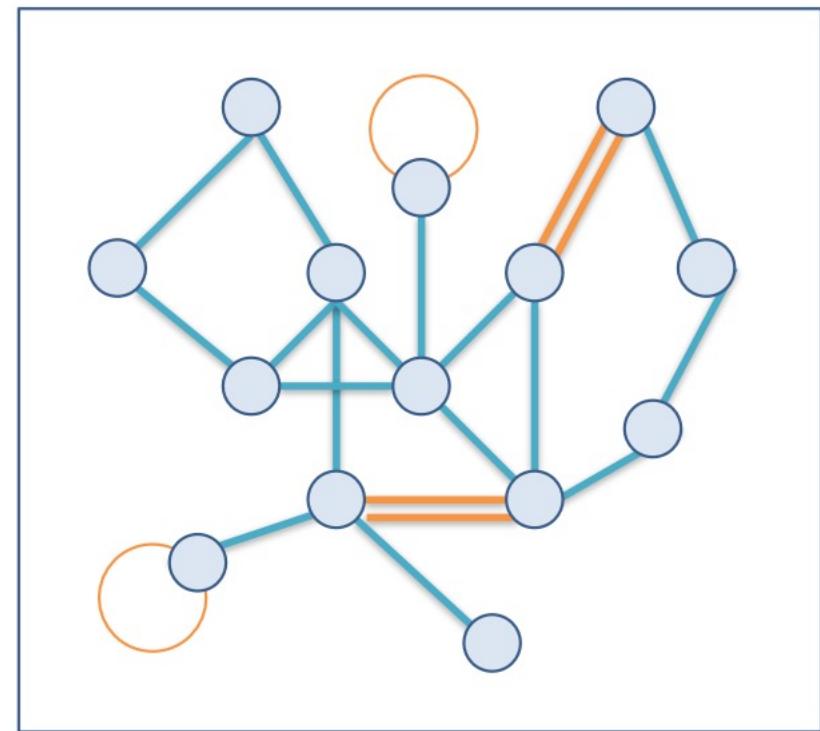
- Vertices and edges are sometimes referred to as **nodes** and **links**.
- Graphs are sometimes referred to as **networks**.



Mathematics Terminology

Properties of a *simple* graph:

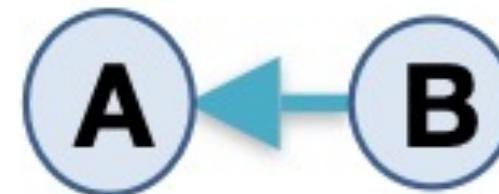
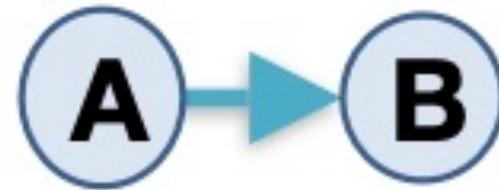
- No self-loops
 - No multi-edges



Not a simple graph!
→ A *general graph*

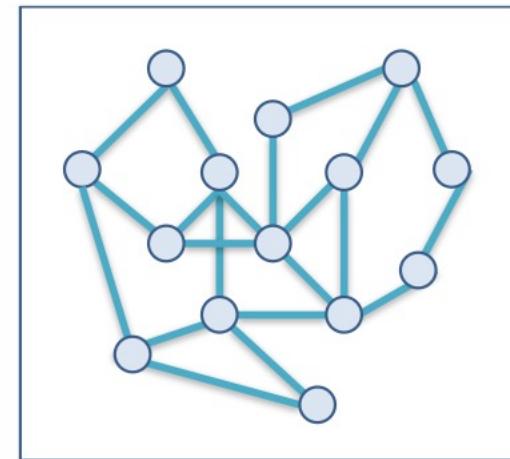
Mathematics Terminology

A *directed* graph specifies the order of the source and target of each edge.



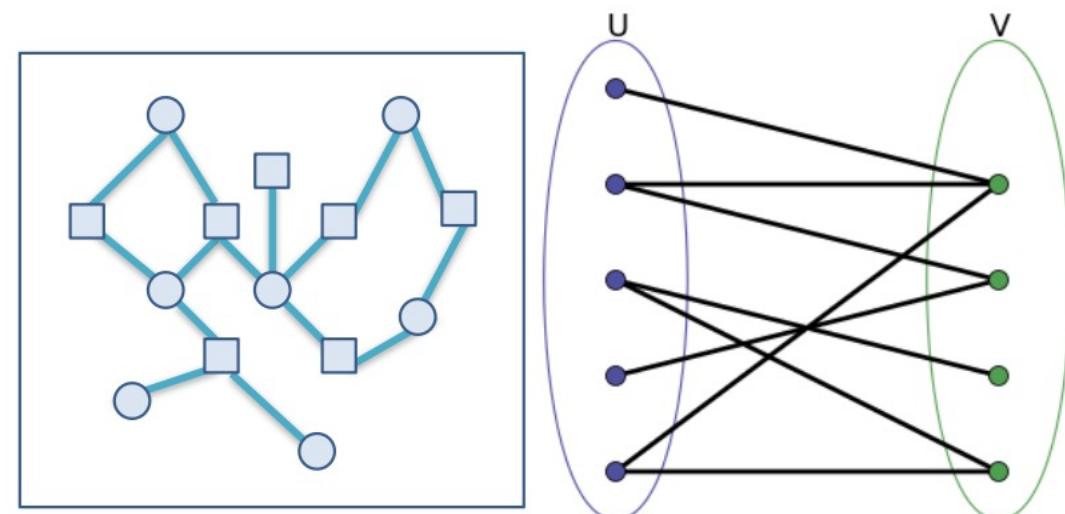
Mathematics Terminology

A ***biconnected*** graph is one that is still connected when you remove any one node.



Biconnected Graph

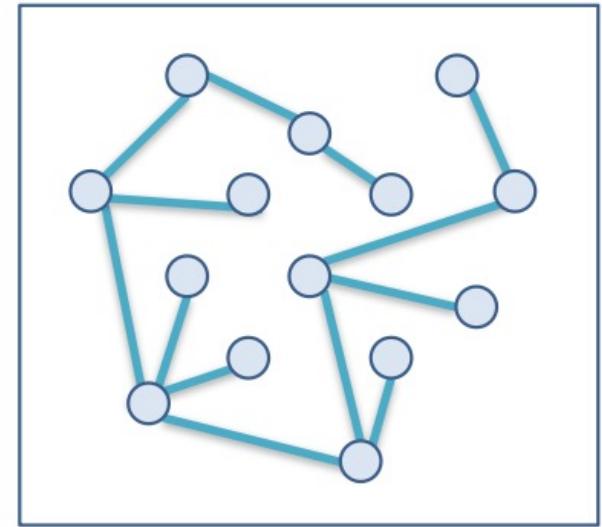
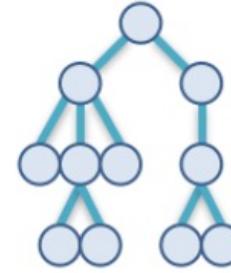
A ***bipartite*** graph can be partitioned into two sets of unconnected vertices.



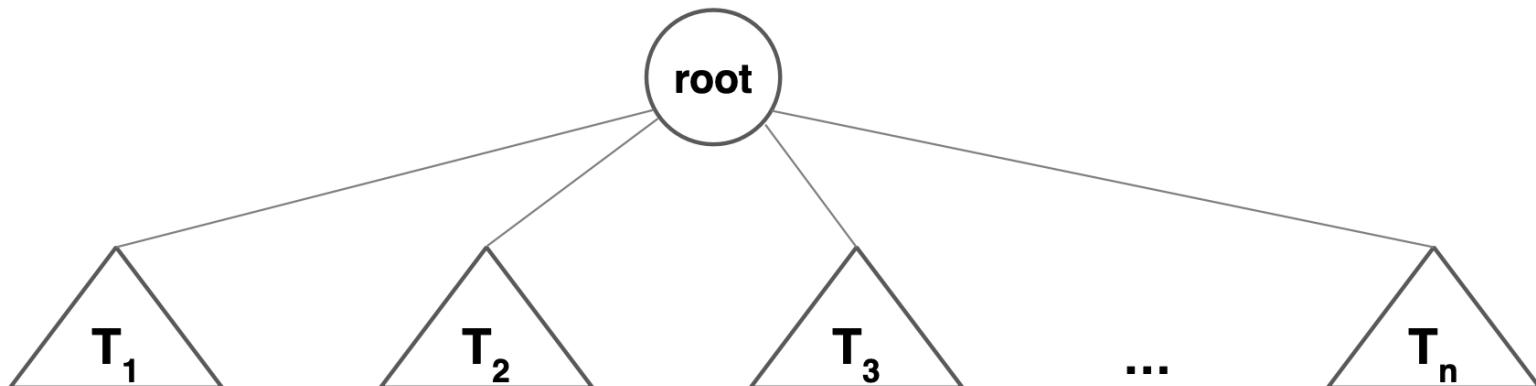
Bipartite Graph

Mathematics Terminology

A ***tree*** is a connected graph without any cycles.

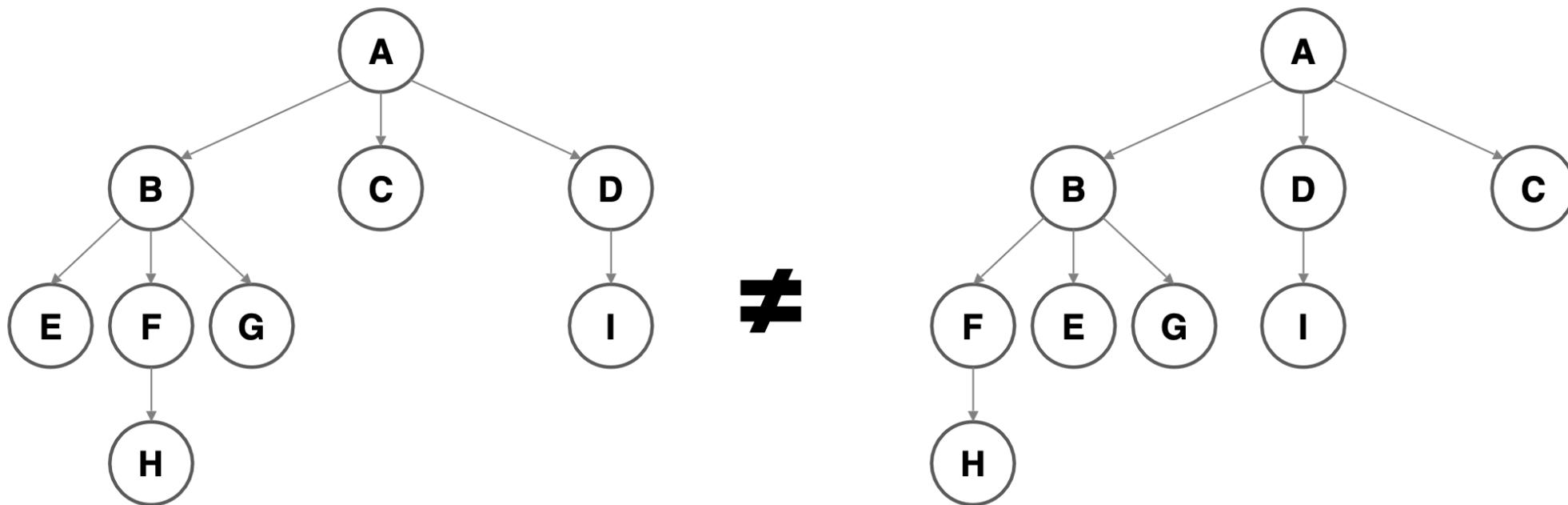


You may choose a single vertex as the ***root*** of a tree. Often it is helpful to consider sub-trees of the rooted tree.



Mathematics Terminology

An *ordered tree* defines an order to nodes coming from the same parent.

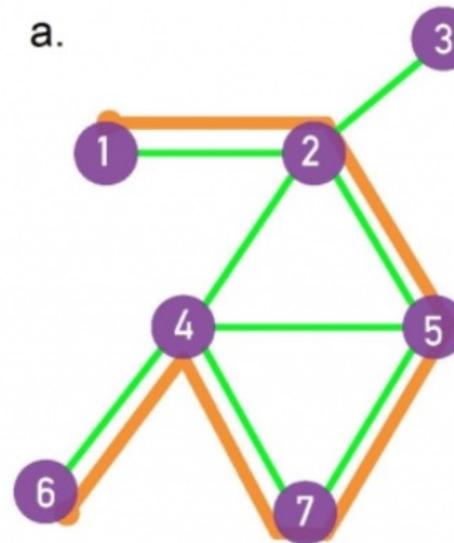


Mathematics Terminology

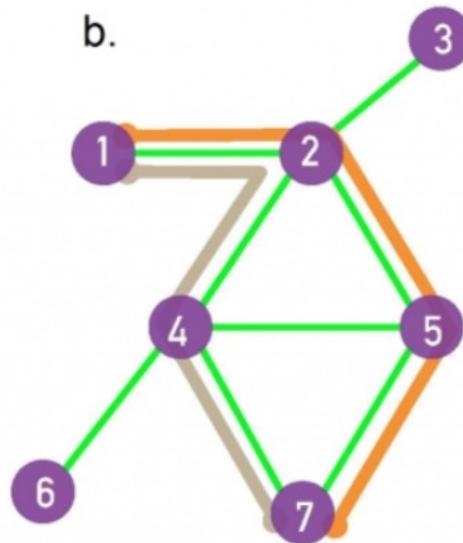
An ***path*** is an ordered list of unrepeated, connected vertices from one vertex to another.

The ***length*** of the path is its number edges.

This is often used in routing.



A path from 1 to 6



Shortest paths (two) from 1 to 7

Mathematics Terminology

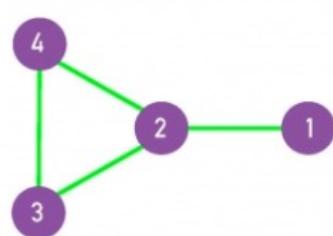
The *degree* of a node is the number of edges connected to it.

Some analyses involve average degree:

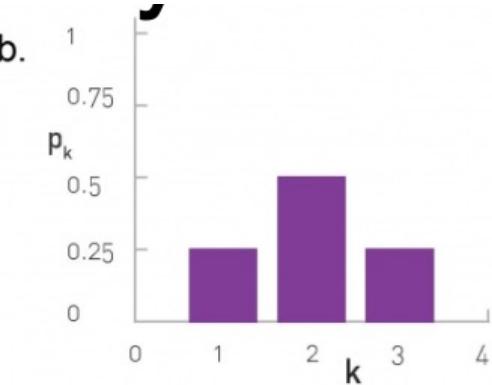
$$\langle k \rangle = \frac{1}{N} \sum_{i=1}^N k_i = \frac{2L}{N}$$

Others examine the degree distribution.

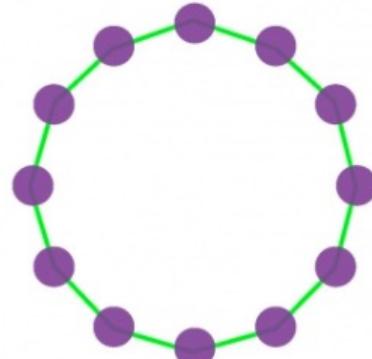
a.



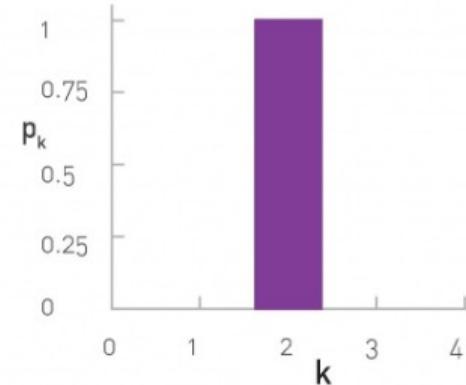
b.



c.

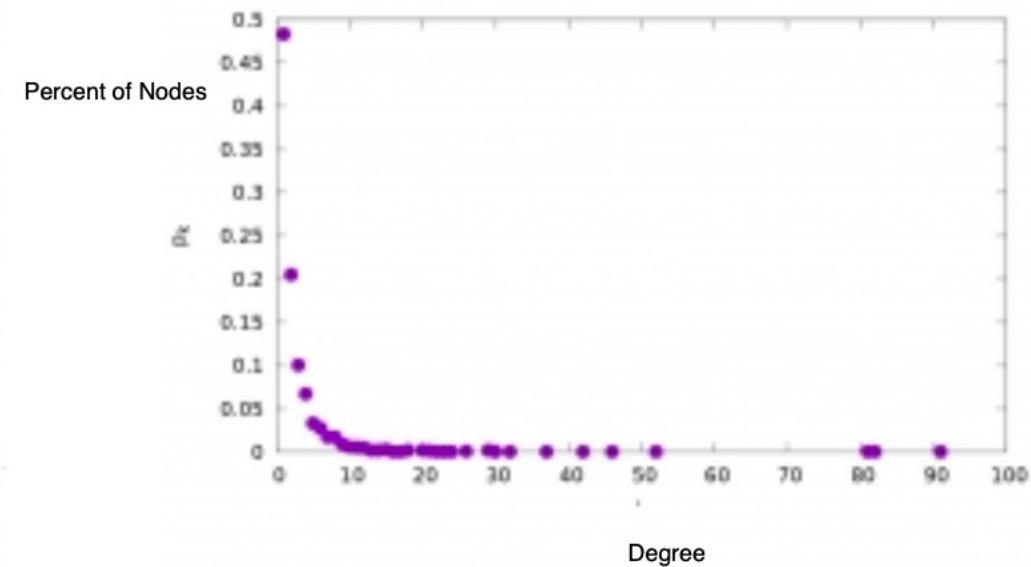


d.



Degree distribution example

a.

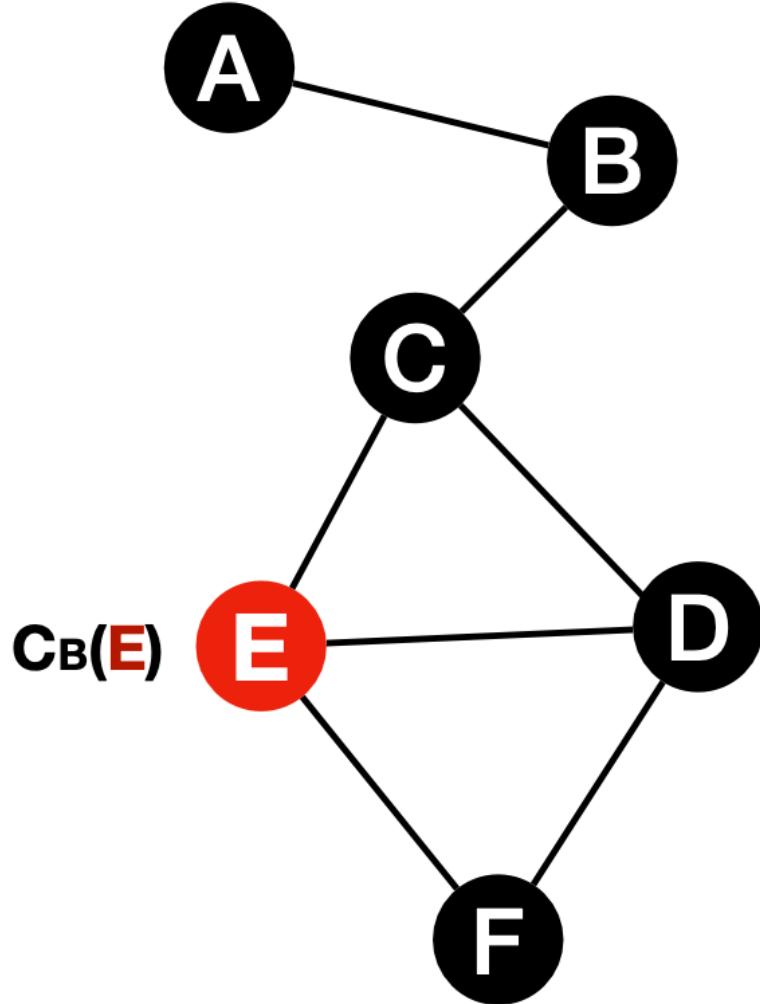


Protein Interaction Network

Betweenness Centrality

$$c_B(v) = \sum_{s \neq v \neq t \in V} \frac{\sigma(s,t|v)}{\sigma(s,t)}$$

Number of shortest paths
Total number of paths



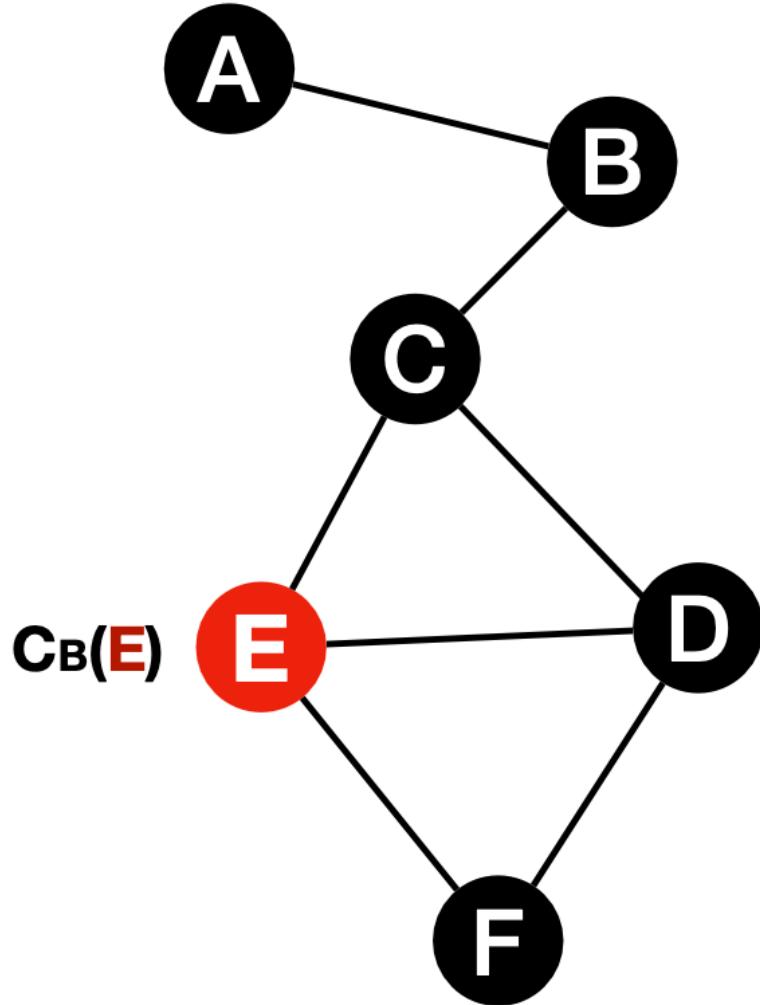
Betweenness centrality measures how many shortest paths pass through a given node.

It is often analyzed to determine a node's importance.

Betweenness Centrality

$$c_B(v) = \sum_{s \neq v \neq t \in V} \frac{\sigma(s,t|v)}{\sigma(s,t)}$$

Number of shortest paths
Total number of paths

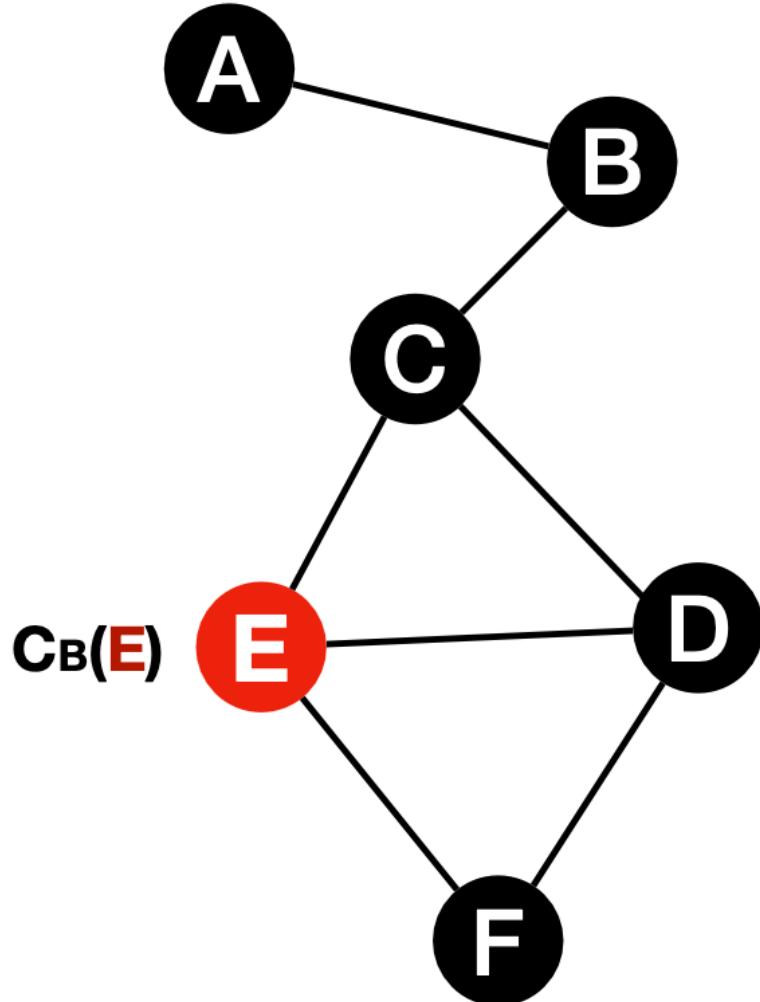


Shortest Path through E?	Number of shortest paths
A, B	0
A, C	0
A, D	0
A, F	1
B, C	0
B, D	0
B, F	1
C, D	0
C, F	1
D, F	0

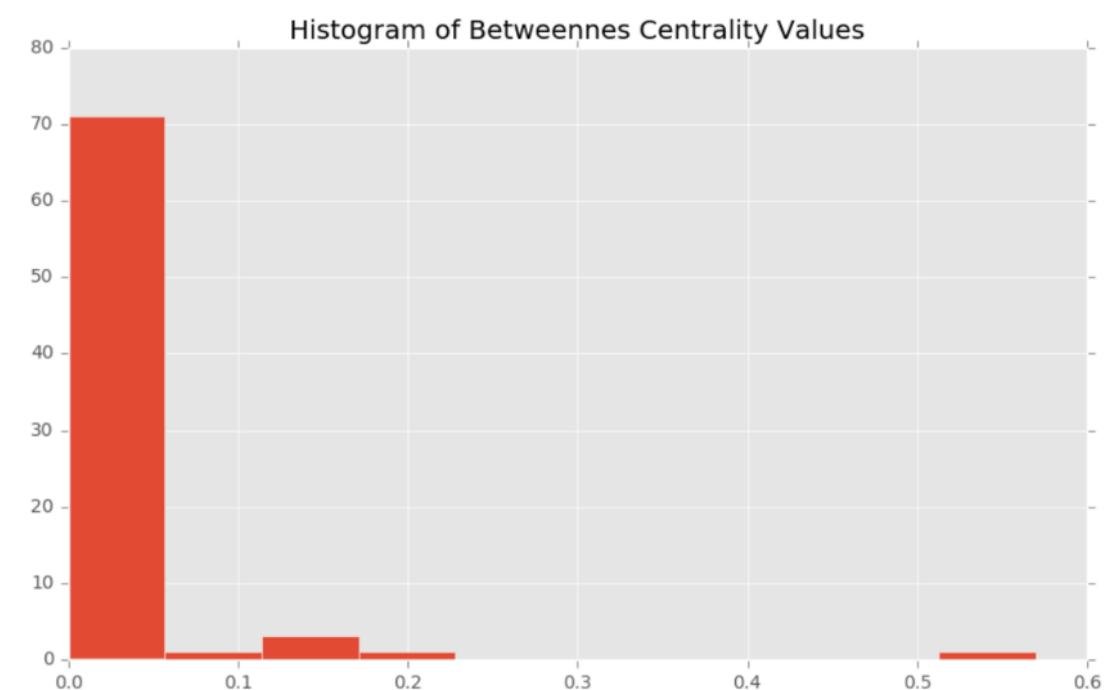
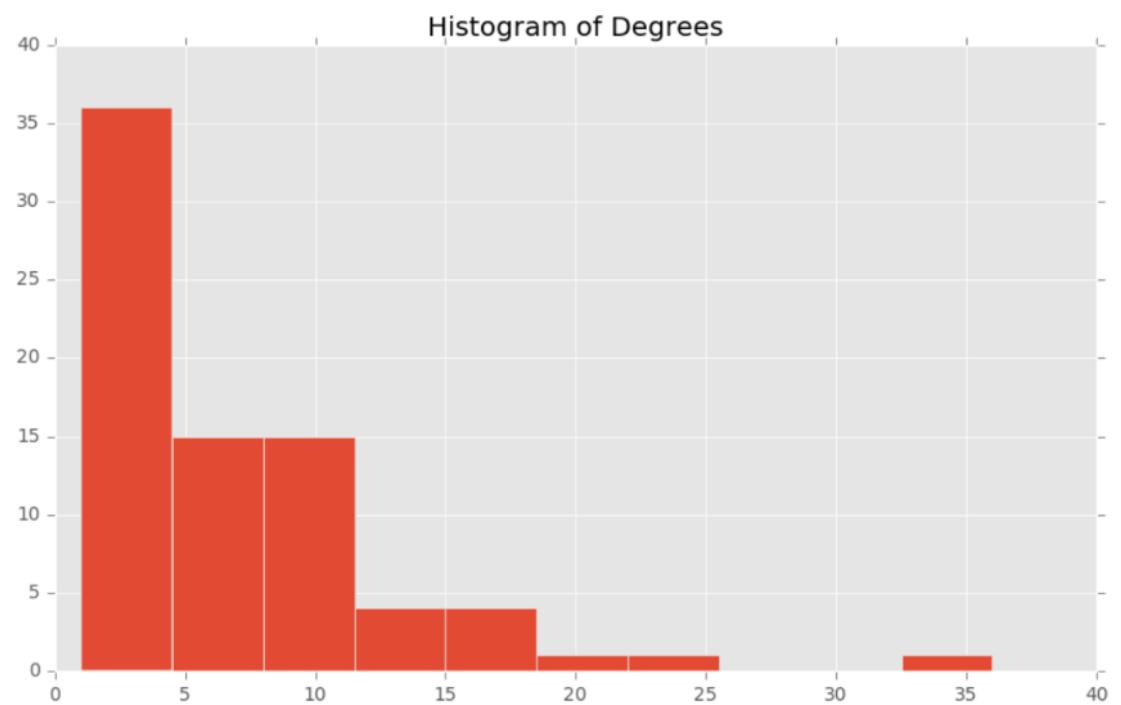
Betweenness Centrality

$$c_B(v) = \sum_{s \neq v \neq t \in V} \frac{\sigma(s,t|v)}{\sigma(s,t)}$$

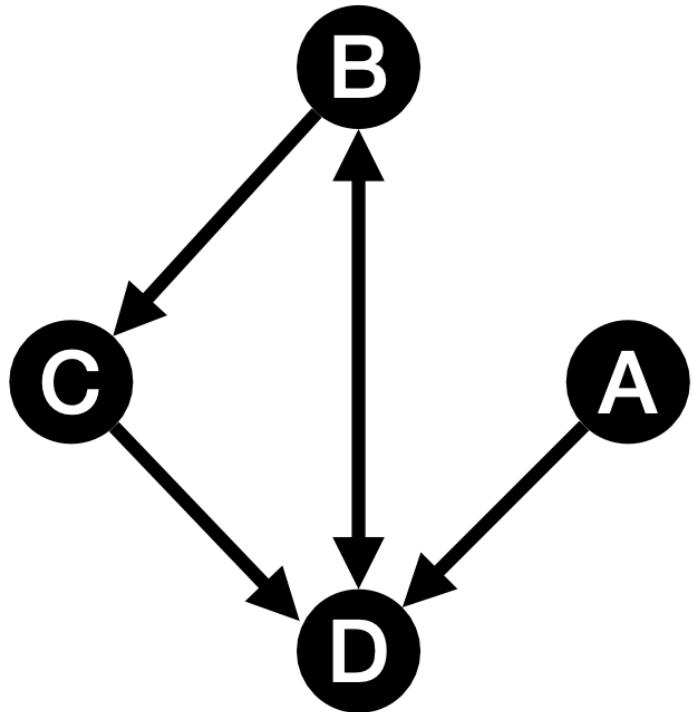
Number of shortest paths
Total number of paths



	Shortest Path through E?	Number of shortest paths
A, B	0	1
A, C	0	1
A, D	0	1
A, F	1	2
B, C	0	1
B, D	0	1
B, F	1	2
C, D	0	1
C, F	1	2
D, F	0	1
		= 1.5



PageRank

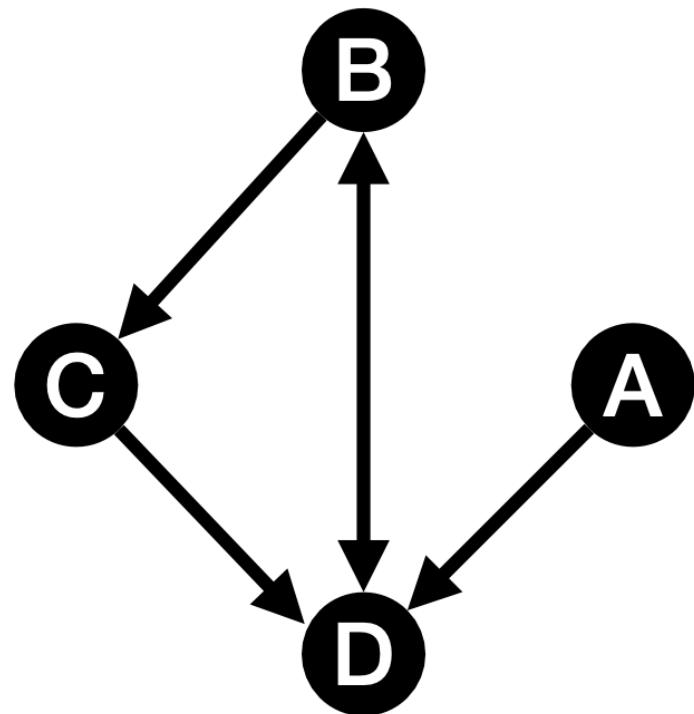


PageRank is the algorithm that gave Google early dominance in search engines in the late 1990s.

Like Betweenness Centrality, it calculates the importance of nodes based on their links.

Page Rank

...and continues until settled...



A	0.25	0	0	0
B	0.25	$0 + 0.25$	$0 + 0.625$	$0 + 0.1875$
C	0.25	$0 + 0.125$	$0 + 0.125$	$0 + 0.3125$
D	$0 + 0.25$ $+ 0.125$ $+ 0.25$	$0 + +$ $0.125 +$ 0.0625	$0 + +$ $0.3125 +$ 0.0625	

Exploratory (Visual) Analysis of Networks

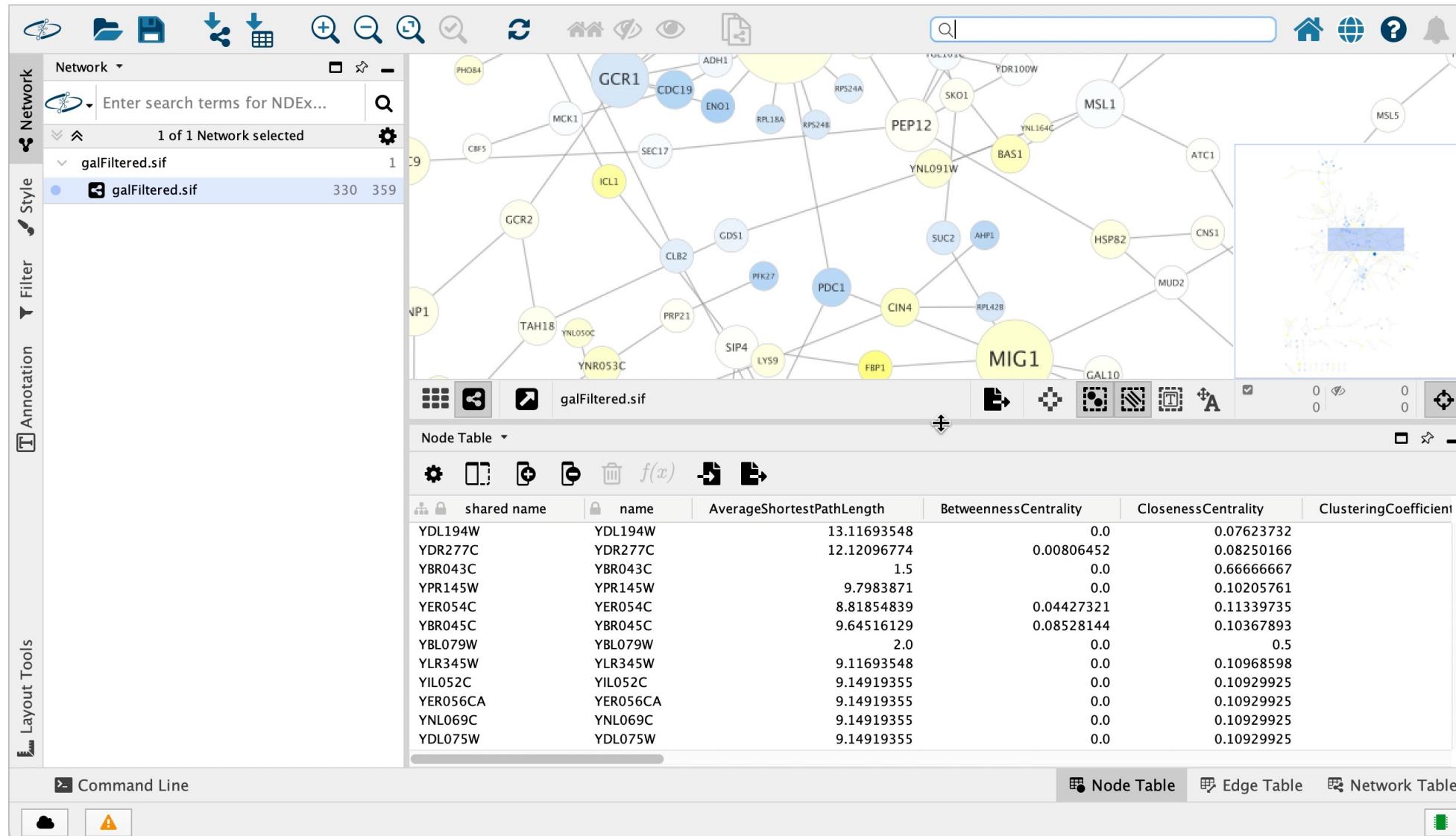
A good place to start for network-centric analysis and visualization is a graph drawing framework.

This slide:

- Cytoscape

Next slide:

- Gephi



 Overview Data Laboratory PreviewWorkspace 1 Appearance Nodes Edges    

Unique Partition Ranking

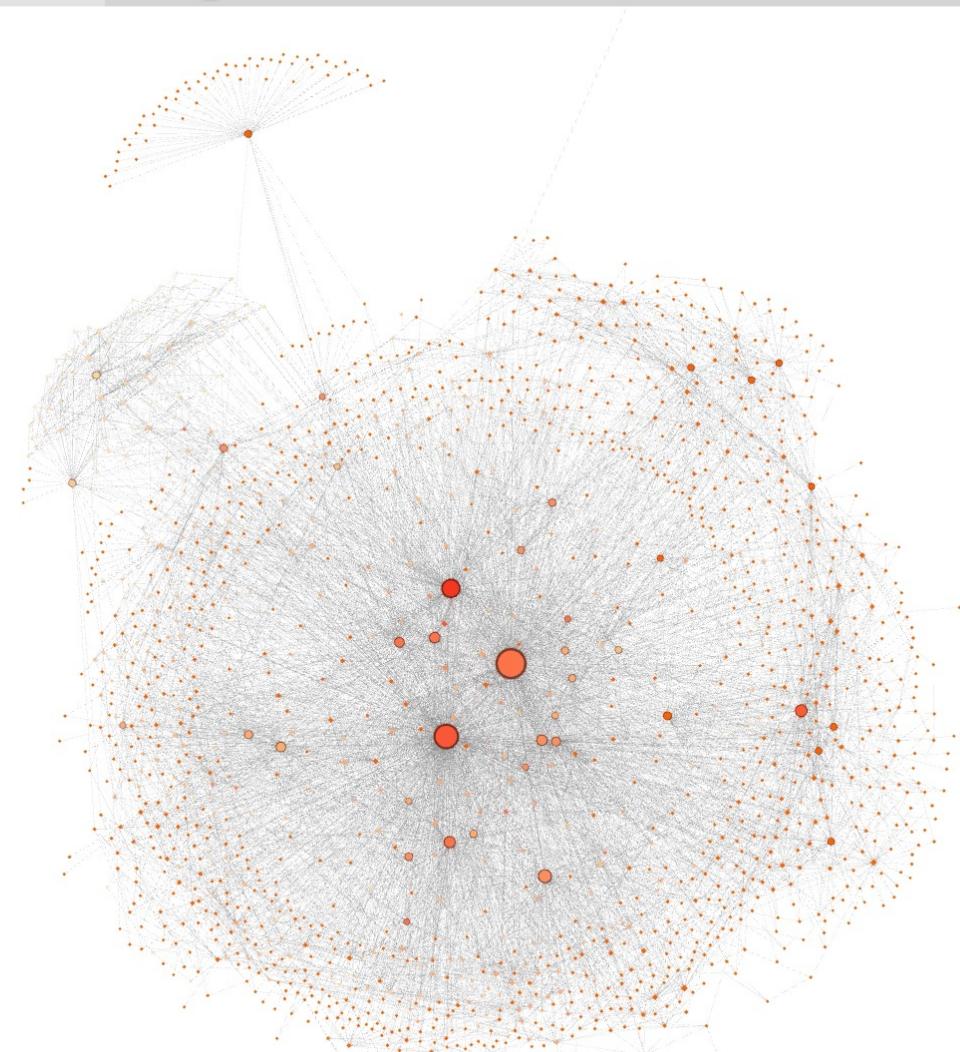
#c0c0c0

Layout Fruchterman Reingold  Apply

Fruchterman Reingold

Area	10000.0
Gravity	10.0
Speed	1.0

Fruchterman Reingold

Graph Mouse selection Color: Context 

Nodes: 1538

Edges: 8032

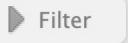
Directed Graph

Filters Statistics Reset Library 

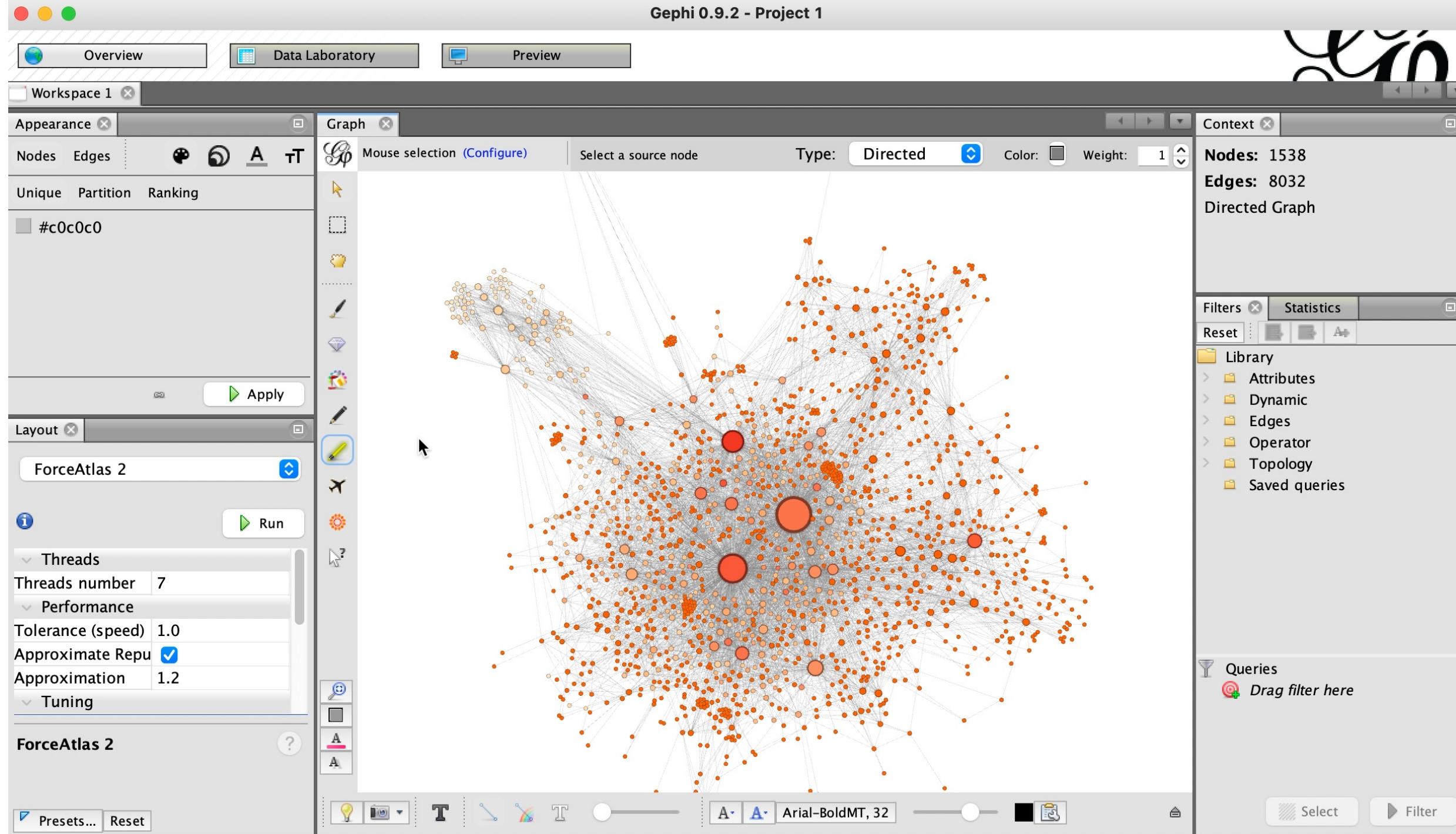
- > Attributes
- > Dynamic
- > Edges
- > Operator
- > Topology
- > Saved queries

Queries 

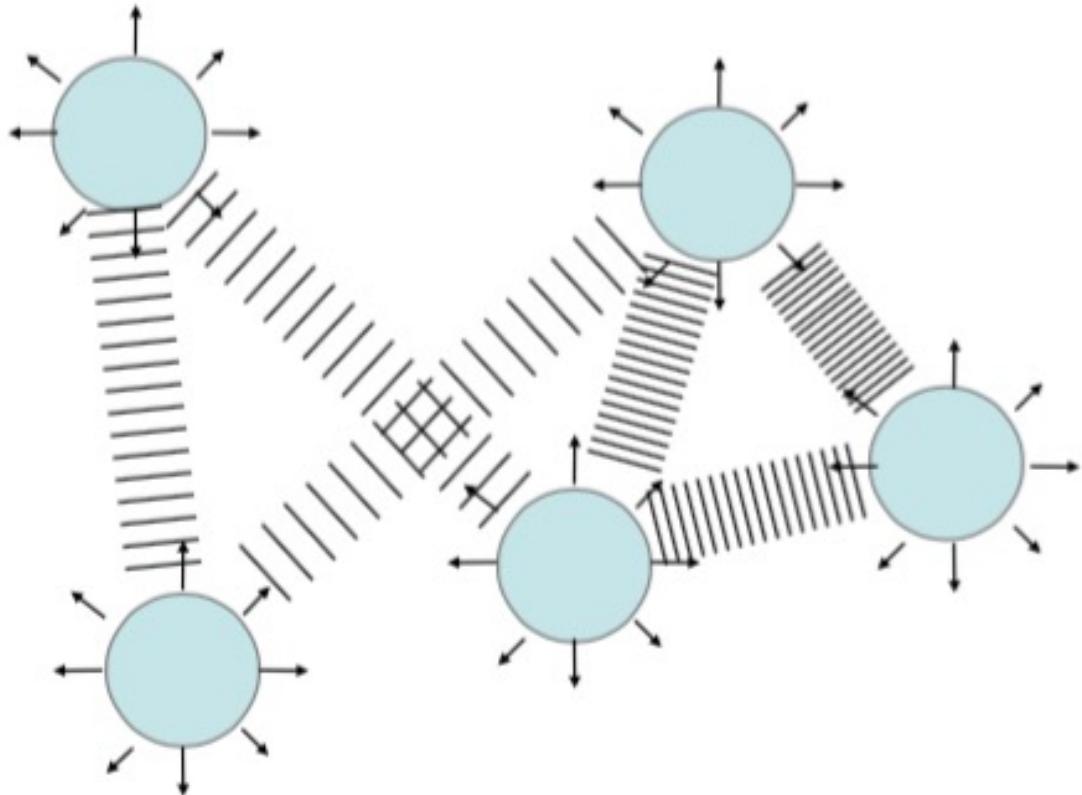
Drag filter here

 Presets... 

A- A+ Arial-BoldMT, 32



Force-directed layouts model networks as a physical system where nodes repel but edges attract



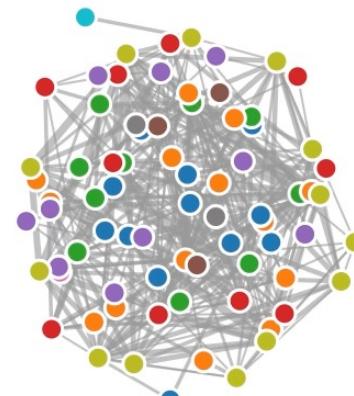
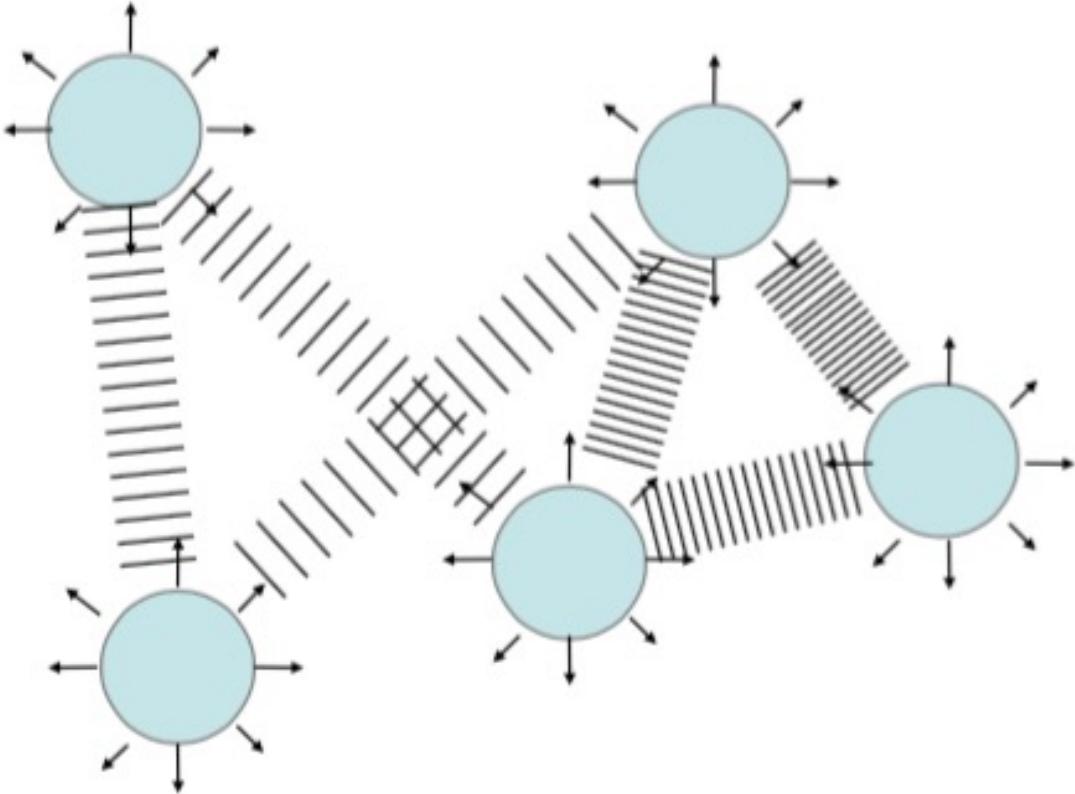
These layouts “work” on a wide variety of graphs.

They also scale better than other approaches.

Issues:

- Lots of tweakable settings
- Are unstable - may not be repeatable

Force-directed layout animation for Les Miserable graph



**If you know more about your network,
another layout may perform “better.”**

Hierarchical networks assume data has a layered, ordered, or otherwise hierarchical structure.

Examples:

- Workflows
- Computing
- Genealogy

Well-known hierarchical layout:

- dot in GraphViz
(pygraphviz)

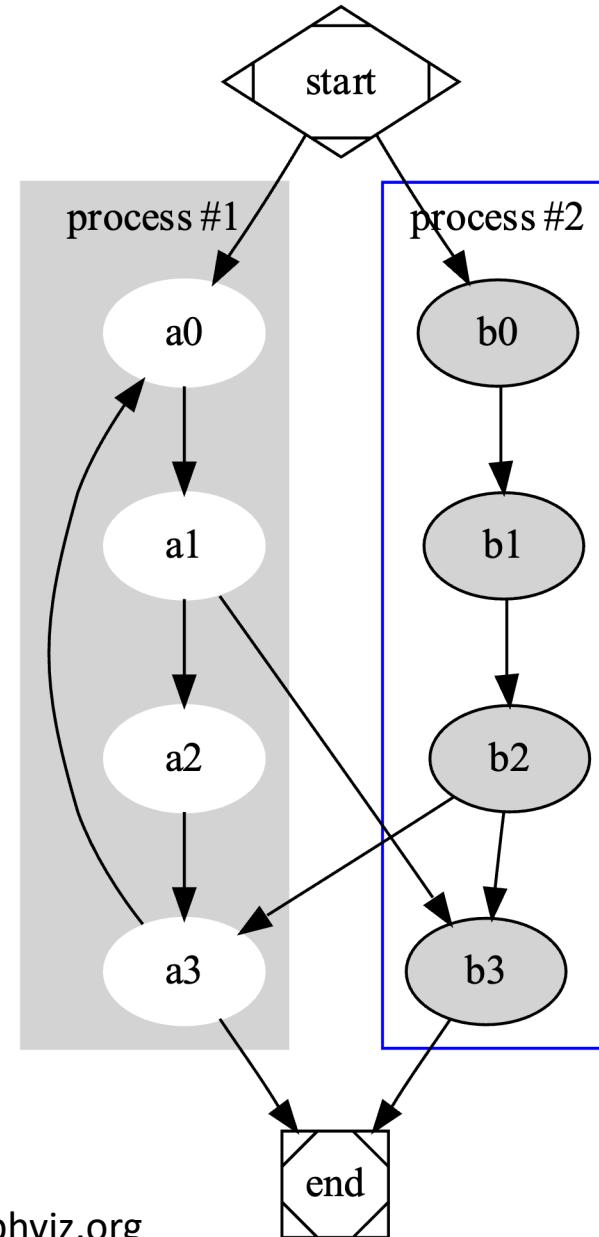


Image from graphviz.org

Dimensionality

All

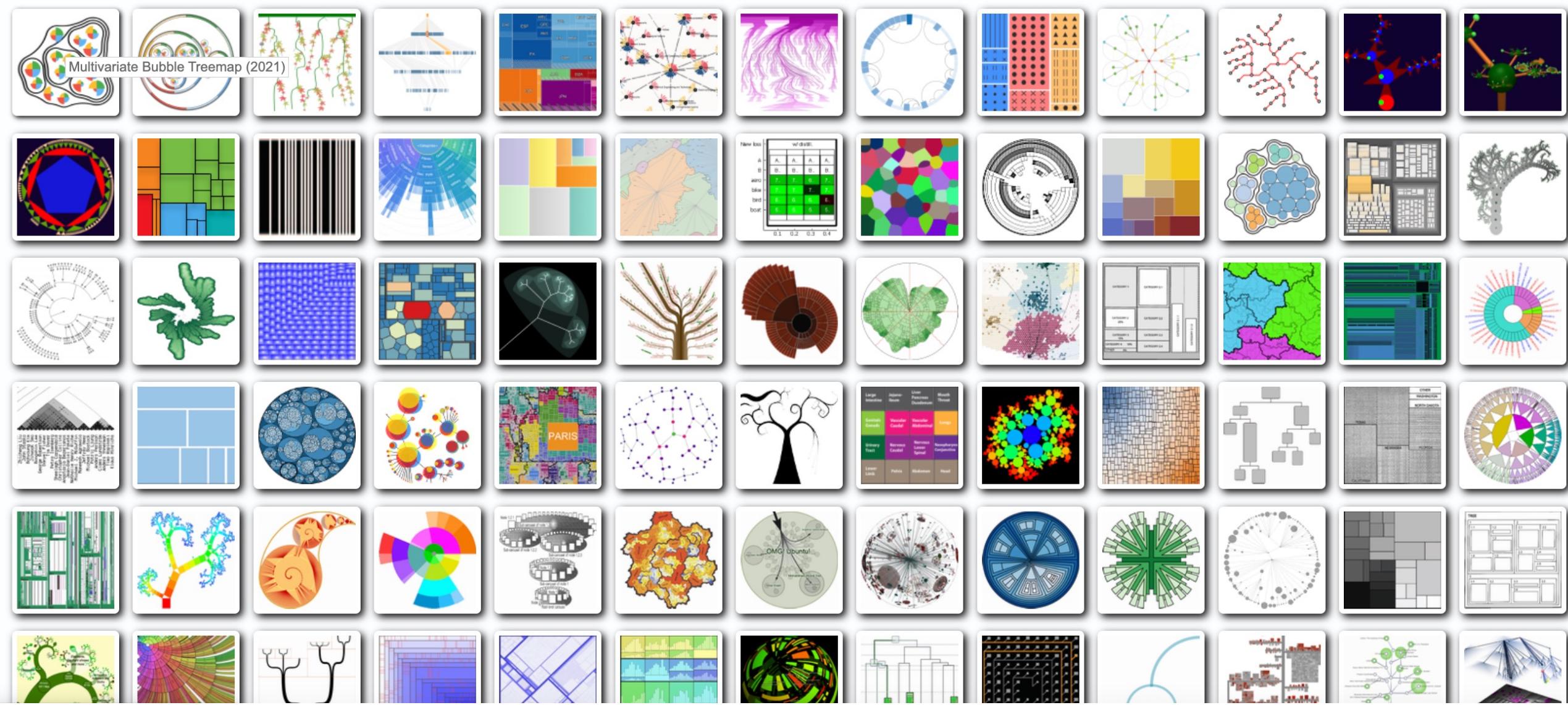
Representation

Alignment

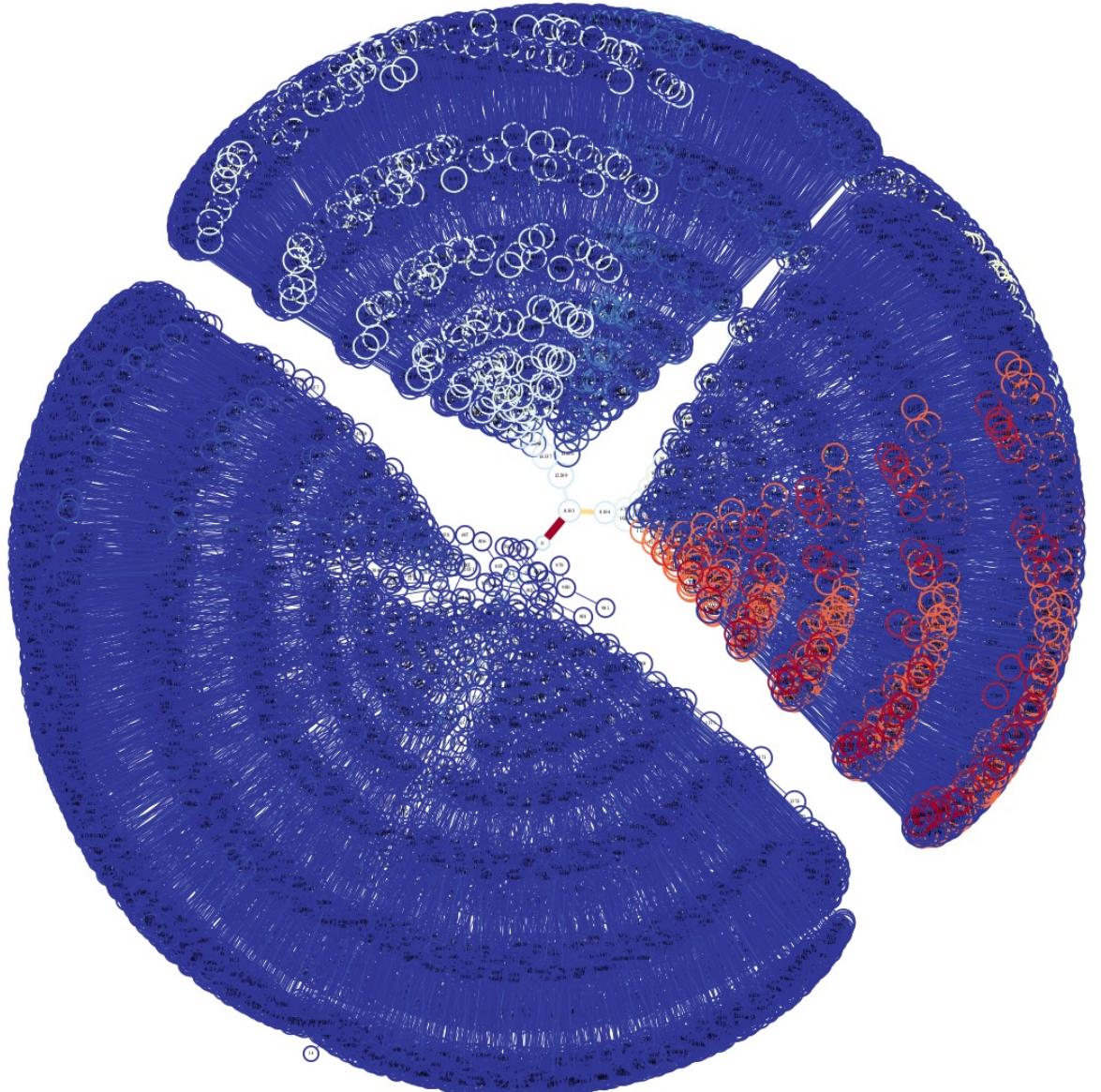
Fulltext Search

Techniques Shown

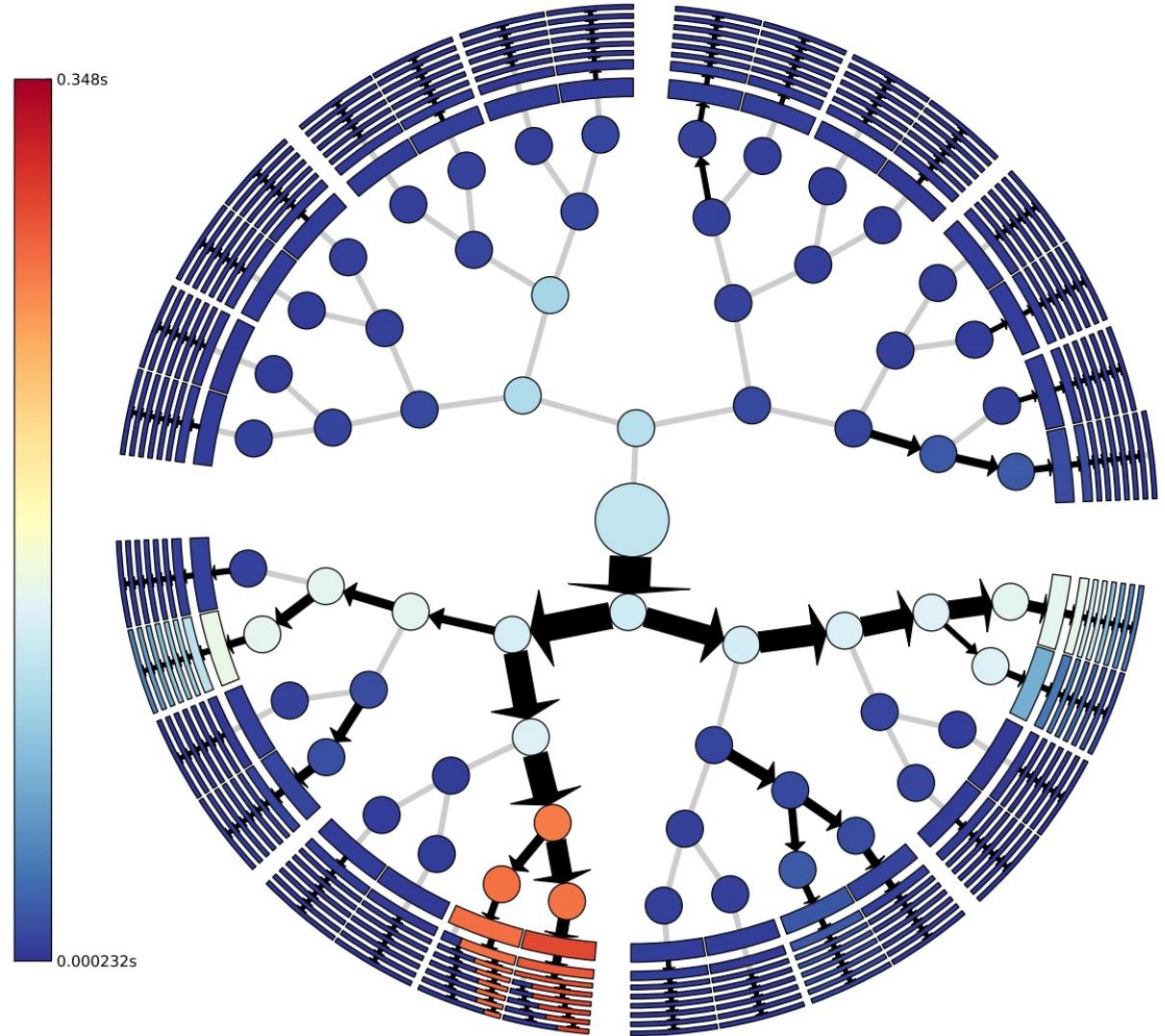
333



When networks get large, aggregation can help.

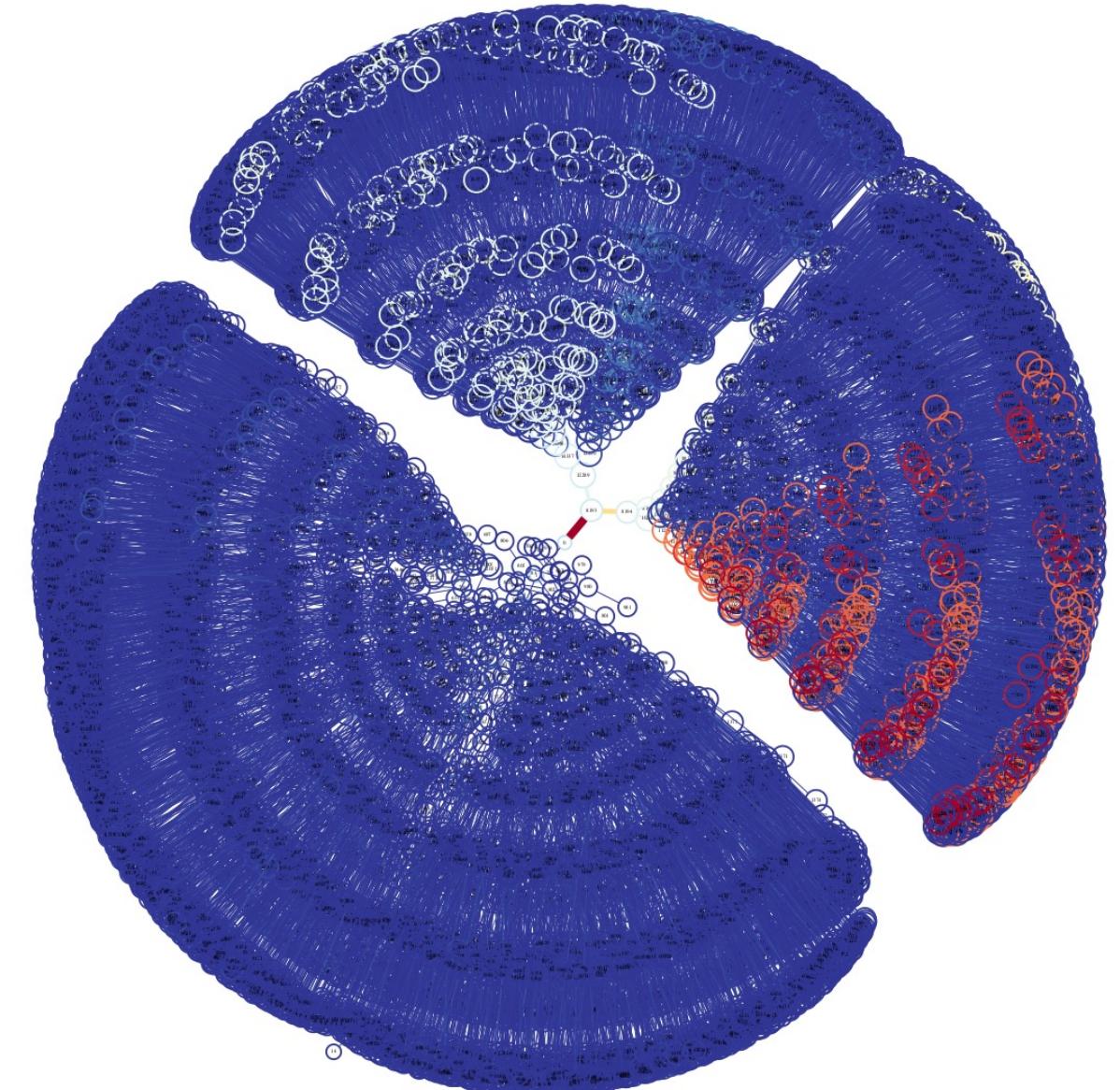


16,384 nodes, force-directed, GraphViz

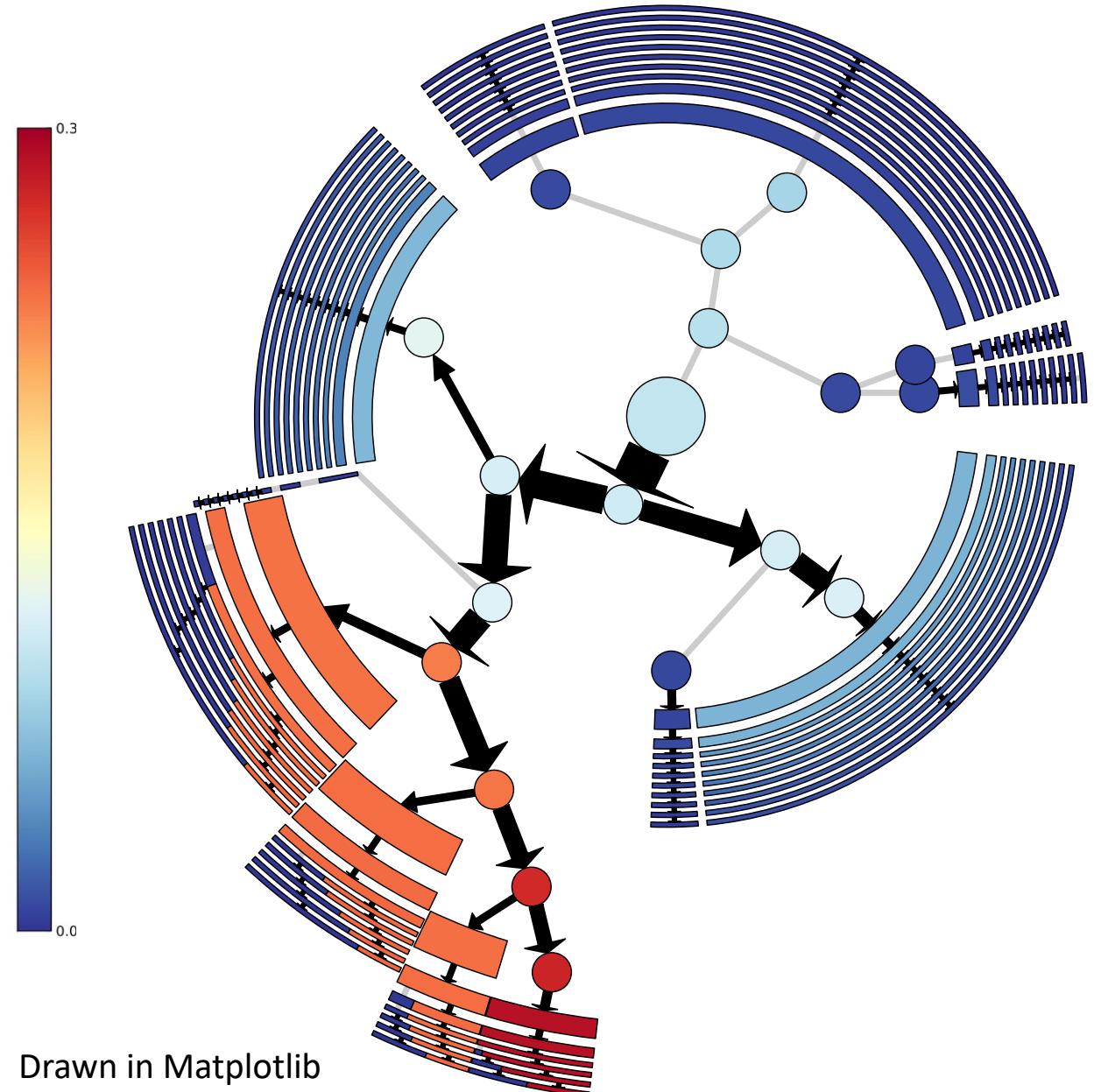


Drawn in Matplotlib

When networks get large, aggregation can help.



16,384 nodes, force-directed, GraphViz



Drawn in Matplotlib

jupyter Workflow Example Last Checkpoint: 16 hours ago (unsaved changes)

File Edit View Insert Cell Kernel Widgets Help

Logout Trusted Python 3

File + % Run Code

```
In [1]: from IPython.display import HTML, display  
import hatchet as ht  
  
display(HTML("<style>.container { width:60% !important; }</style>"))  
%load_ext hatchet.vis.loader
```

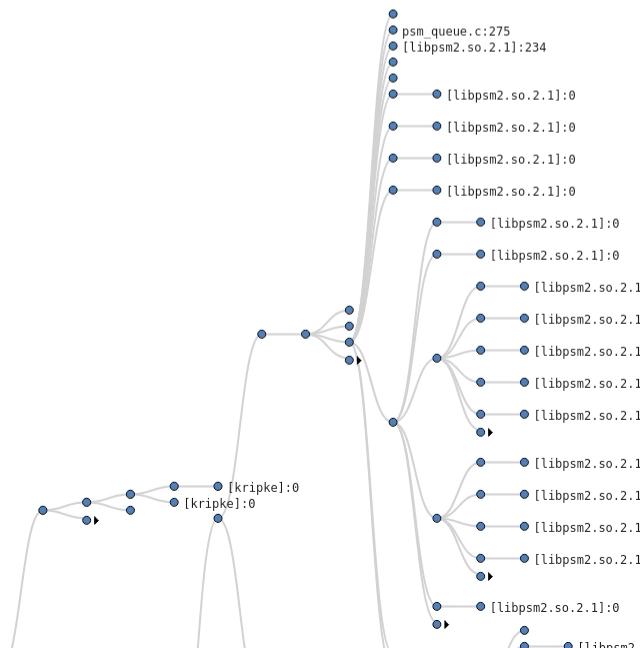
Identifying Optimization Targets

```
In [2]: gf_64 = ht.GraphFrame.from_hpctoolkit('cct-vis-eval/datasets/kripke-scaling/hpctoolkit-kripke-64-cores/')  
gf_128 = ht.GraphFrame.from_hpctoolkit('cct-vis-eval/datasets/kripke-scaling/hpctoolkit-kripke-128-cores/')
```

```
In [3]: %cct gf_64
```

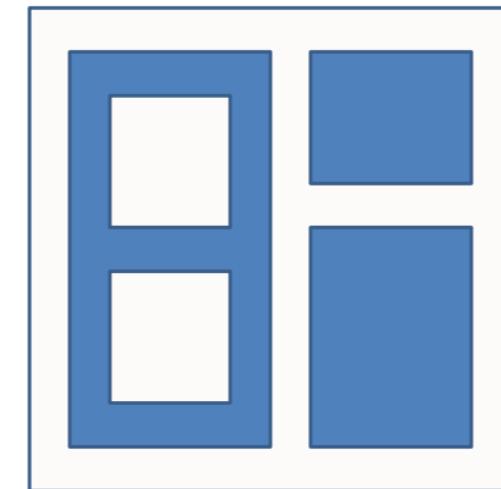
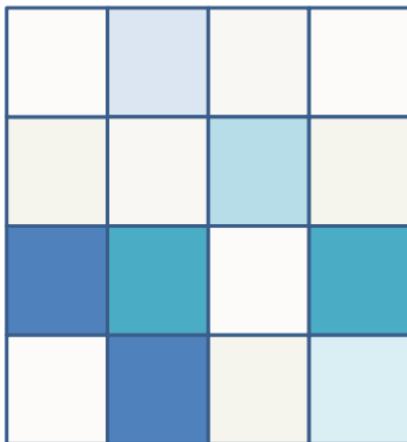
Metrics Display Query Interactive Calling Context Tree

Legend for metric: time inc
64.32M - 77.19M
51.46M - 64.32M
38.59M - 51.46M
25.73M - 38.59M
12.86M - 25.73M
0.00 - 12.86M
0.00 - 31.98M
159.89M - 191.87M
127.91M - 159.89M
95.94M - 127.91M
63.96M - 95.94M
31.98M - 63.96M
0.00 - 31.98M

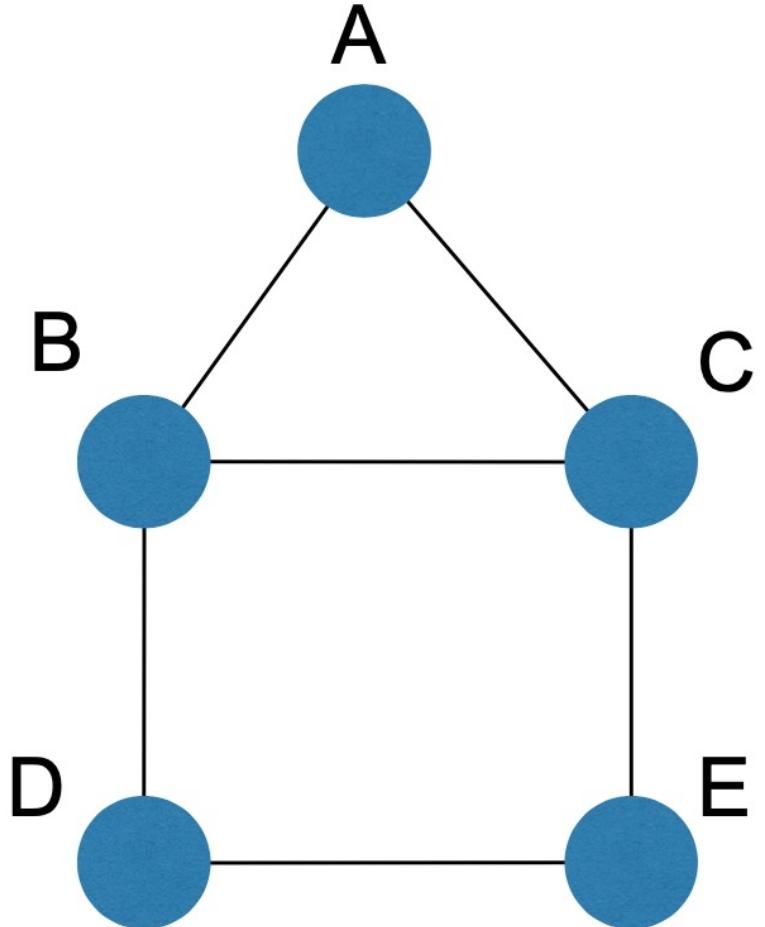


Eliding or filtering can help as well.

Networks don't have to be drawn as node-link diagrams.

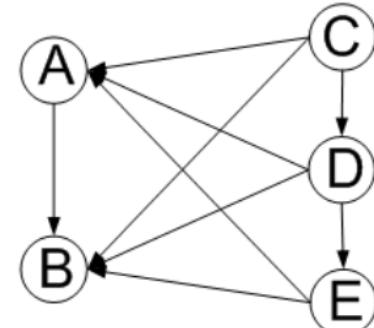
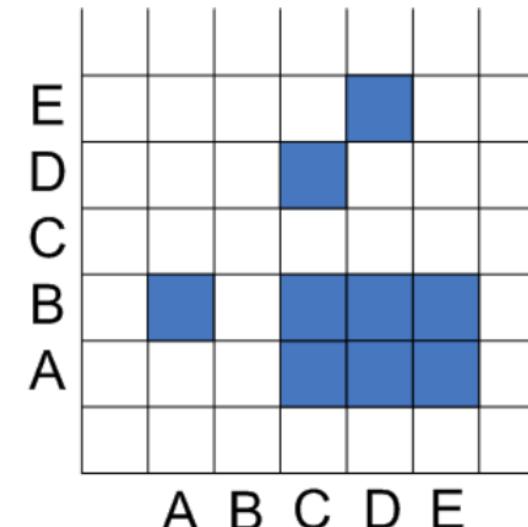
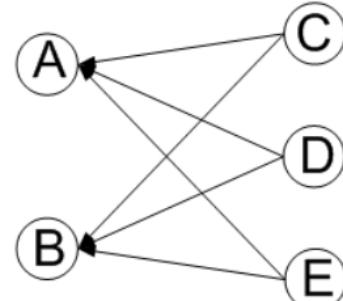
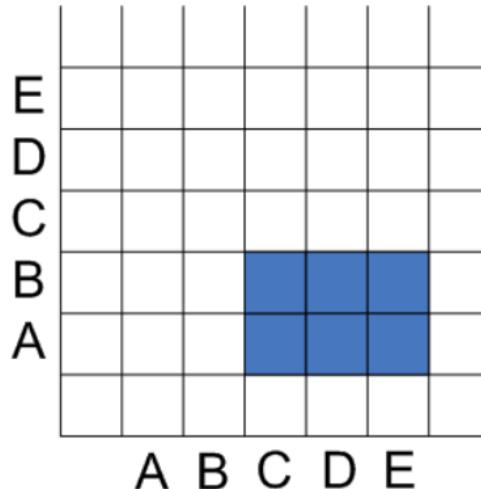
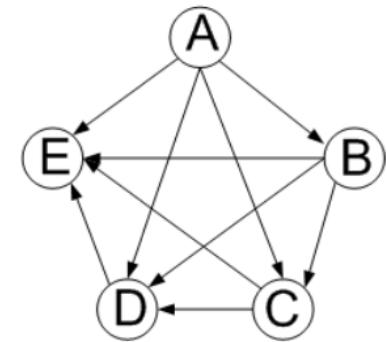
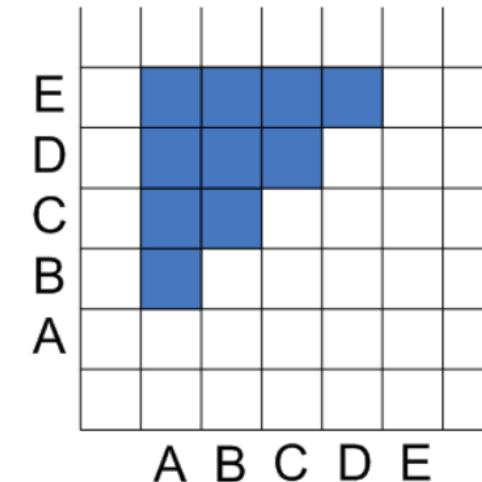
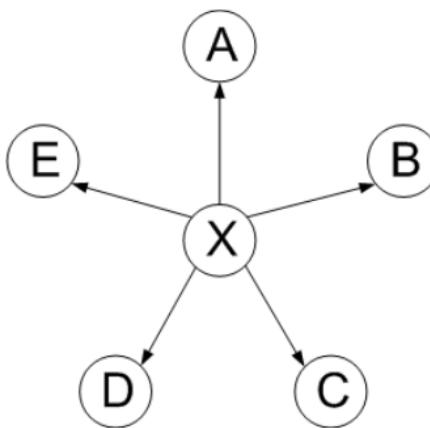
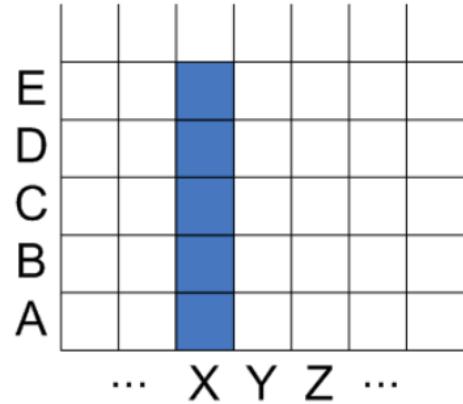


Graphs can be represented as adjacency matrices.

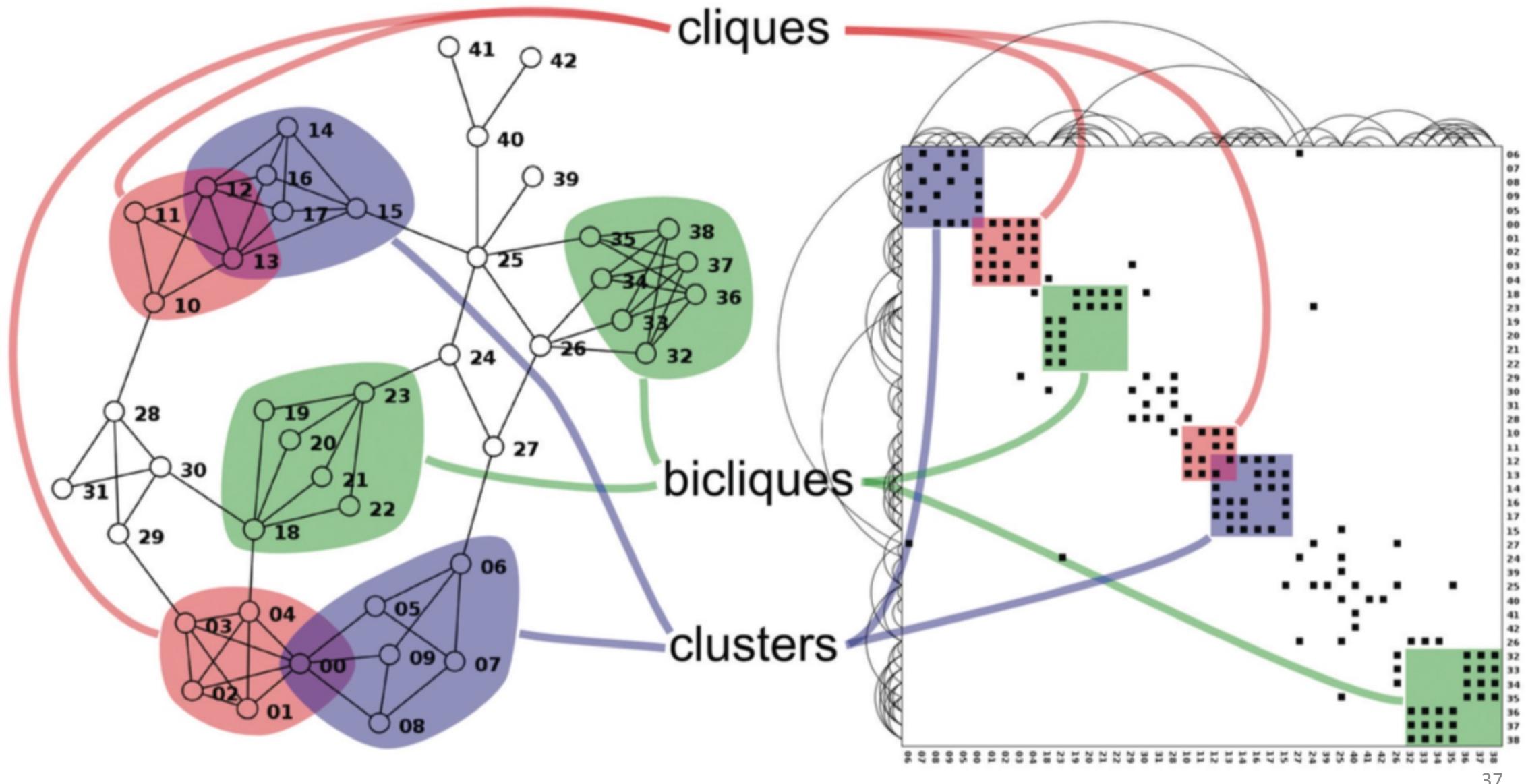


	A	B	C	D	E
A					
B					
C					
D					
E					

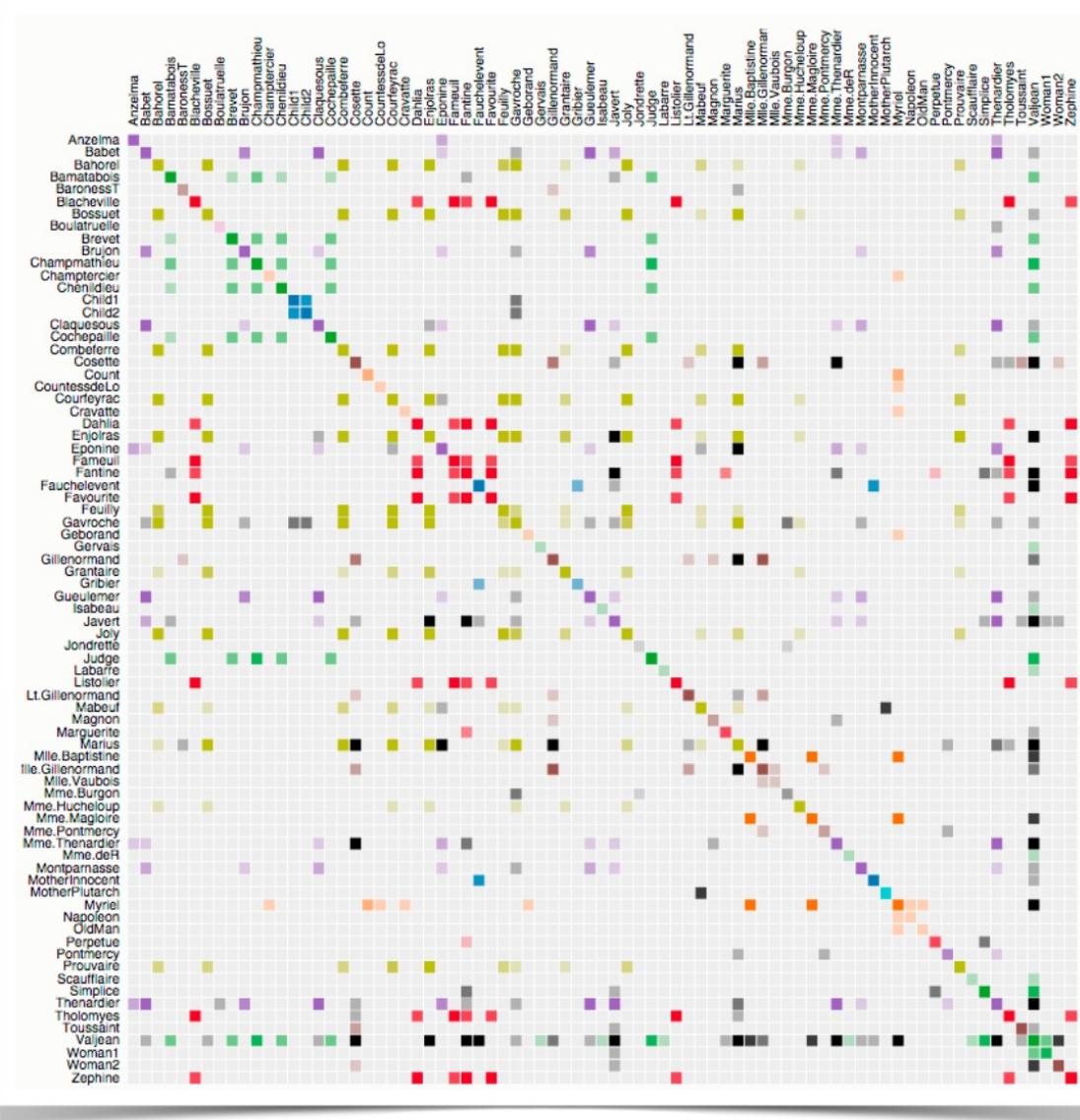
Patterns in directed adjacency matrices



Patterns in directed adjacency matrices



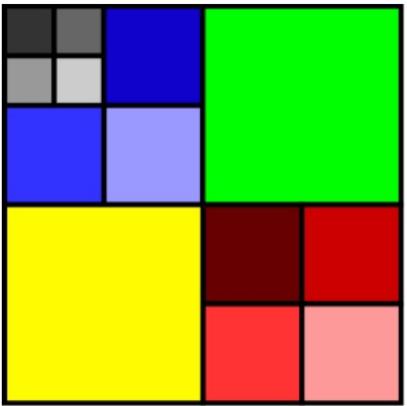
But patterns are only revealed with order...



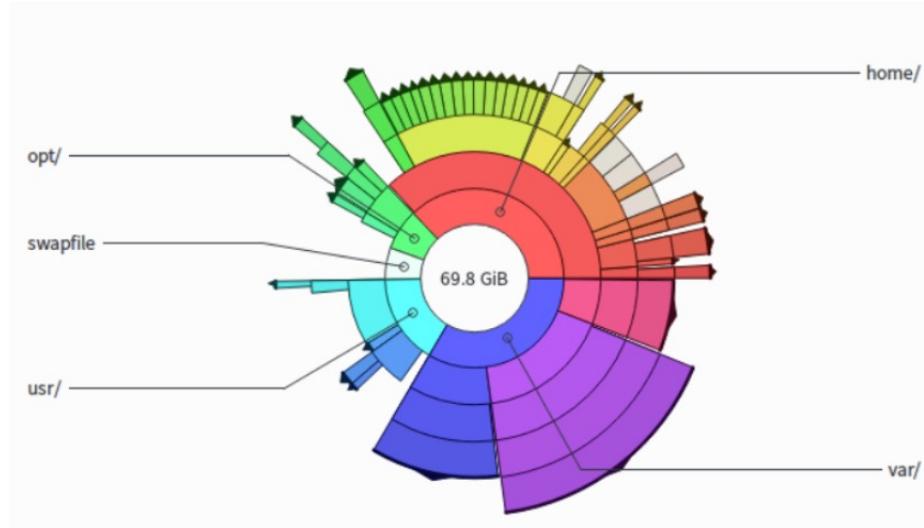
Les Miserable graph

Trees have additional options

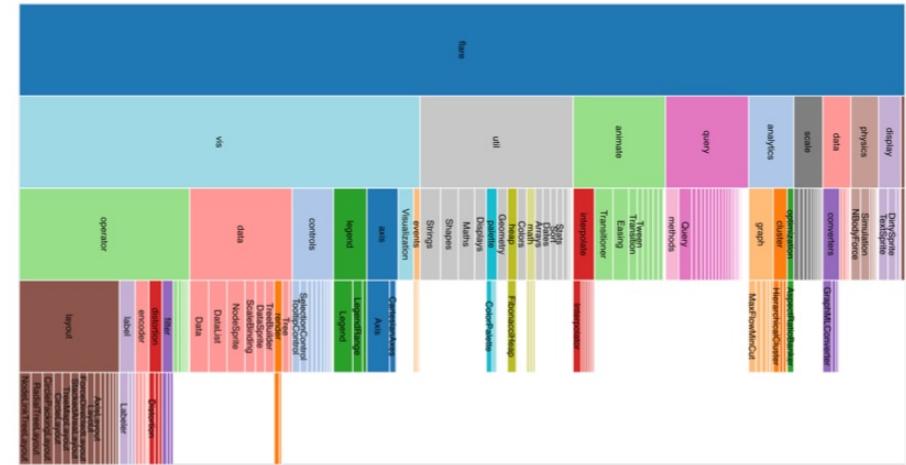
Treemap



Sunburst



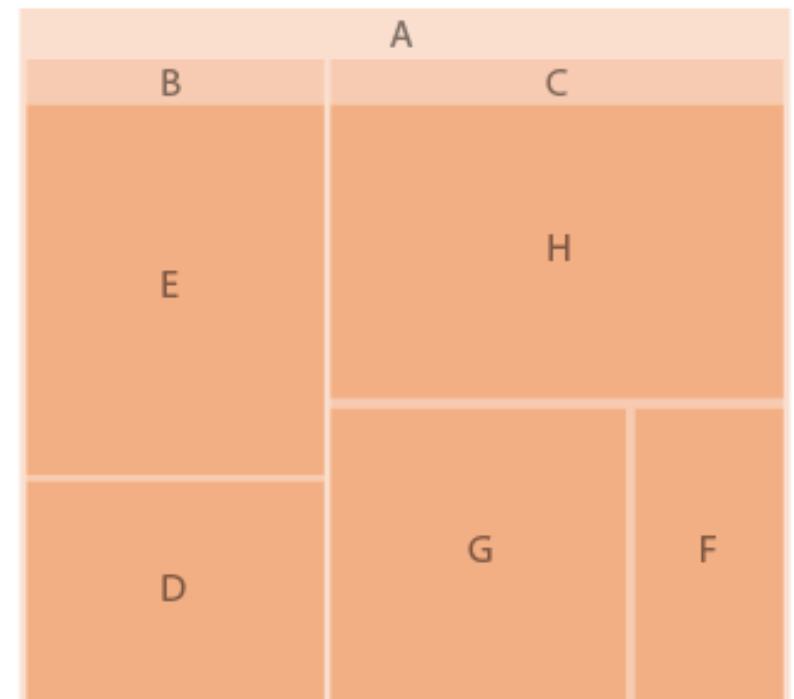
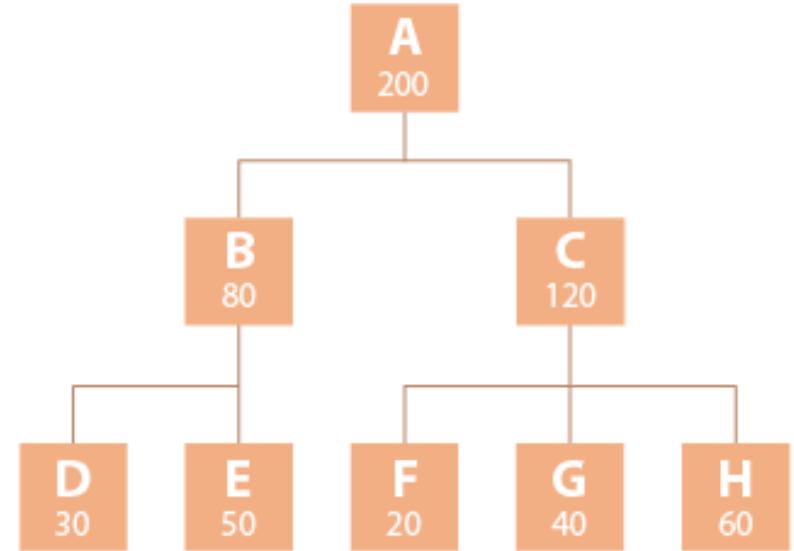
Icicle Plot



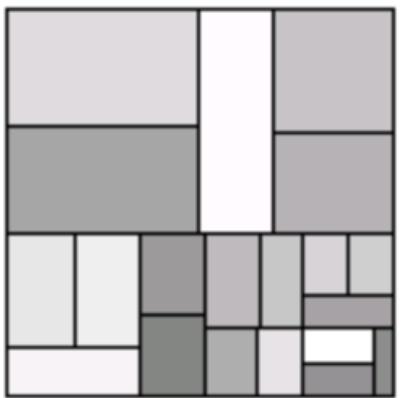
Treemaps

Recursively divide space based on some size metric of the nodes.

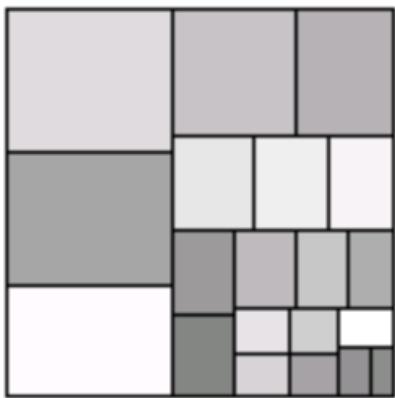
- What metrics?
 - Size of subtree
 - Attribute of node
- How to divide space?
 - Cluster
 - Pivot
 - Slice & Dice
 - Squarified
 - Strip
 - More...



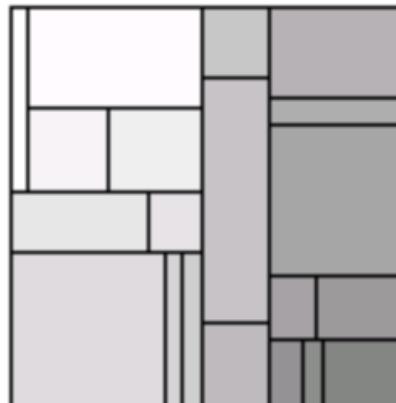
Treemap Layout Algorithms



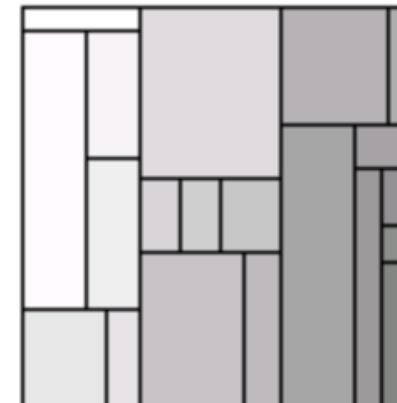
Cluster



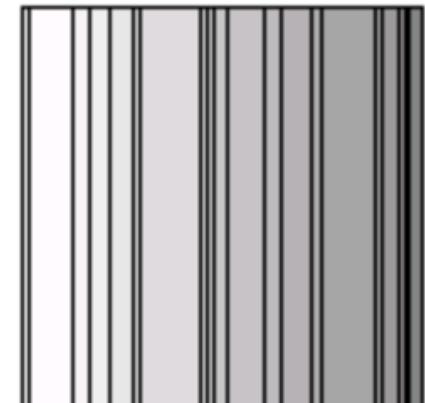
Squareified



Pivot-by-Middle

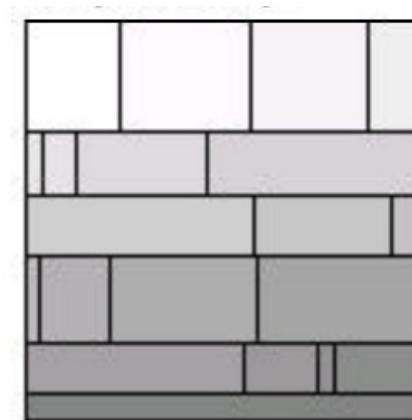


Pivot-by-Size



Slice-and-Dice

Grayscale indicates index order of nodes.



Strip

TECHNOLOGY

INTERNET INFORMATION



APPLICATIONS



SEMICONDUCT



BUSINESS SOFTWARE & S

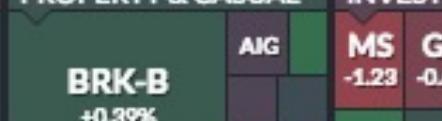


FINANCIAL

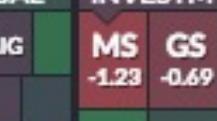
MONEY CENTER BANKS



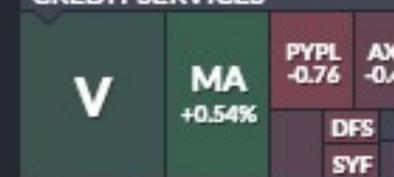
PROPERTY & CASUAL



INVESTM



CREDIT SERVICES



SERVICES

CATALOG & MAIL O



DISCOUNT



ENTERTAI



CONSUMER GOODS



BASIC MATERIALS

MAJOR INT



INDEPE



SPECIAL



CVX

HEALTHCARE

DRUG MANUFACTURE



HEALTH CARE PLANS



INDUSTRIAL GOODS



Use mouse wheel to zoom in and out. Drag zoomed map to pan it.

Double-click a ticker to display detailed information in a new window.

Hover mouse cursor over a ticker to see its main competitors in a stacked view with a 3-month history graph.

-3%

-2%

-1%

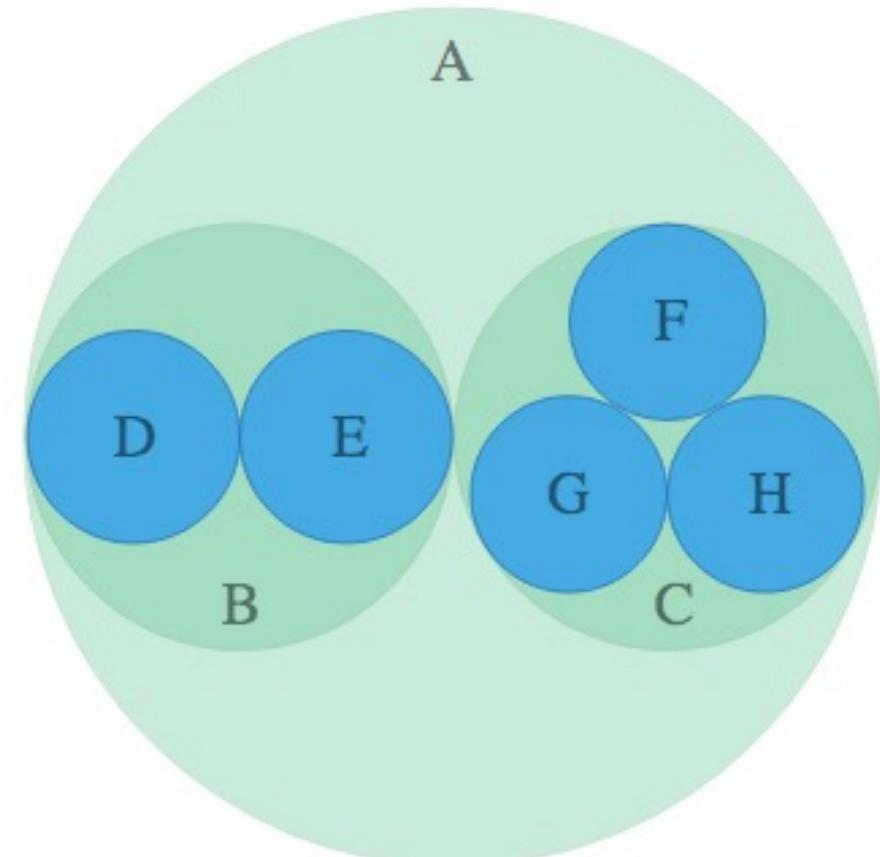
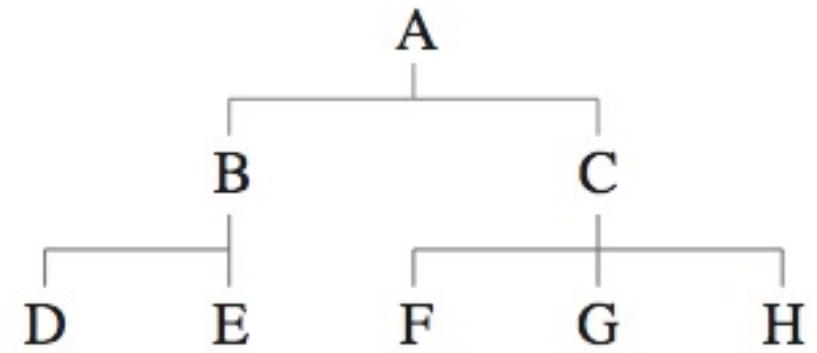
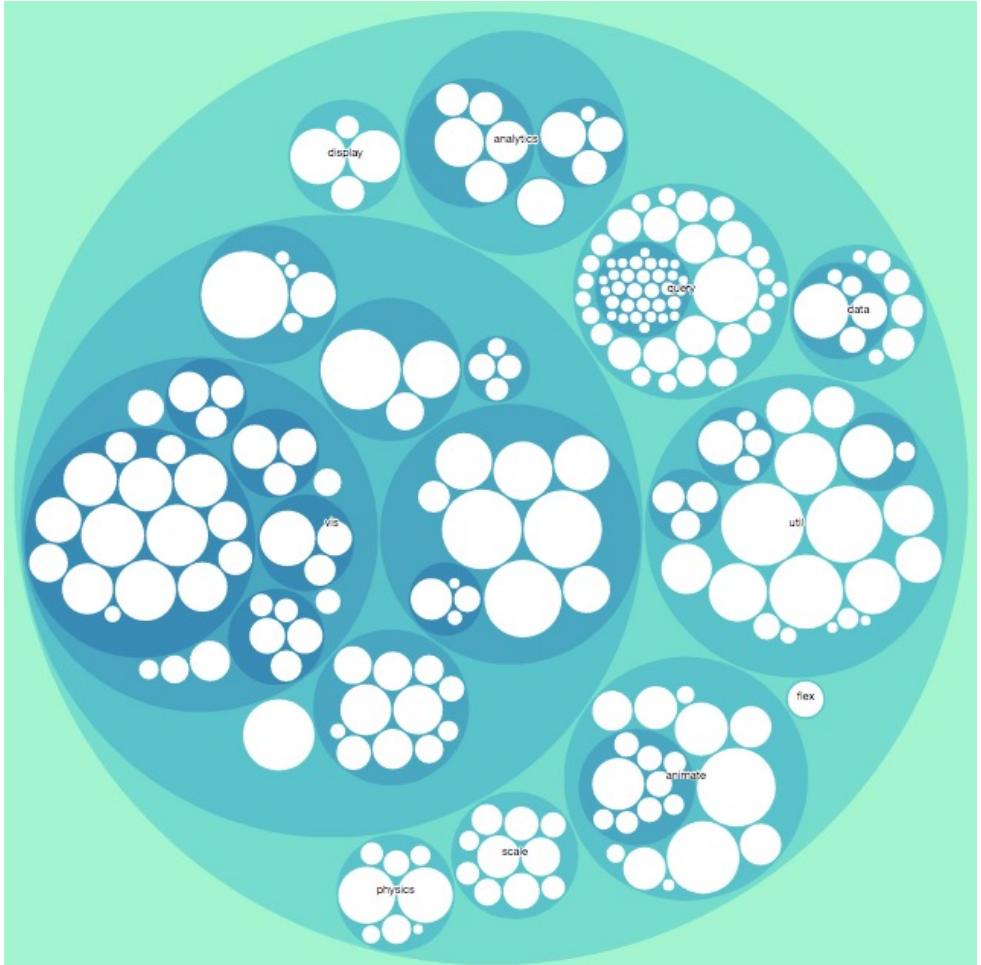
0%

+1%

+2%

+3%

Treemaps need not be rectangular

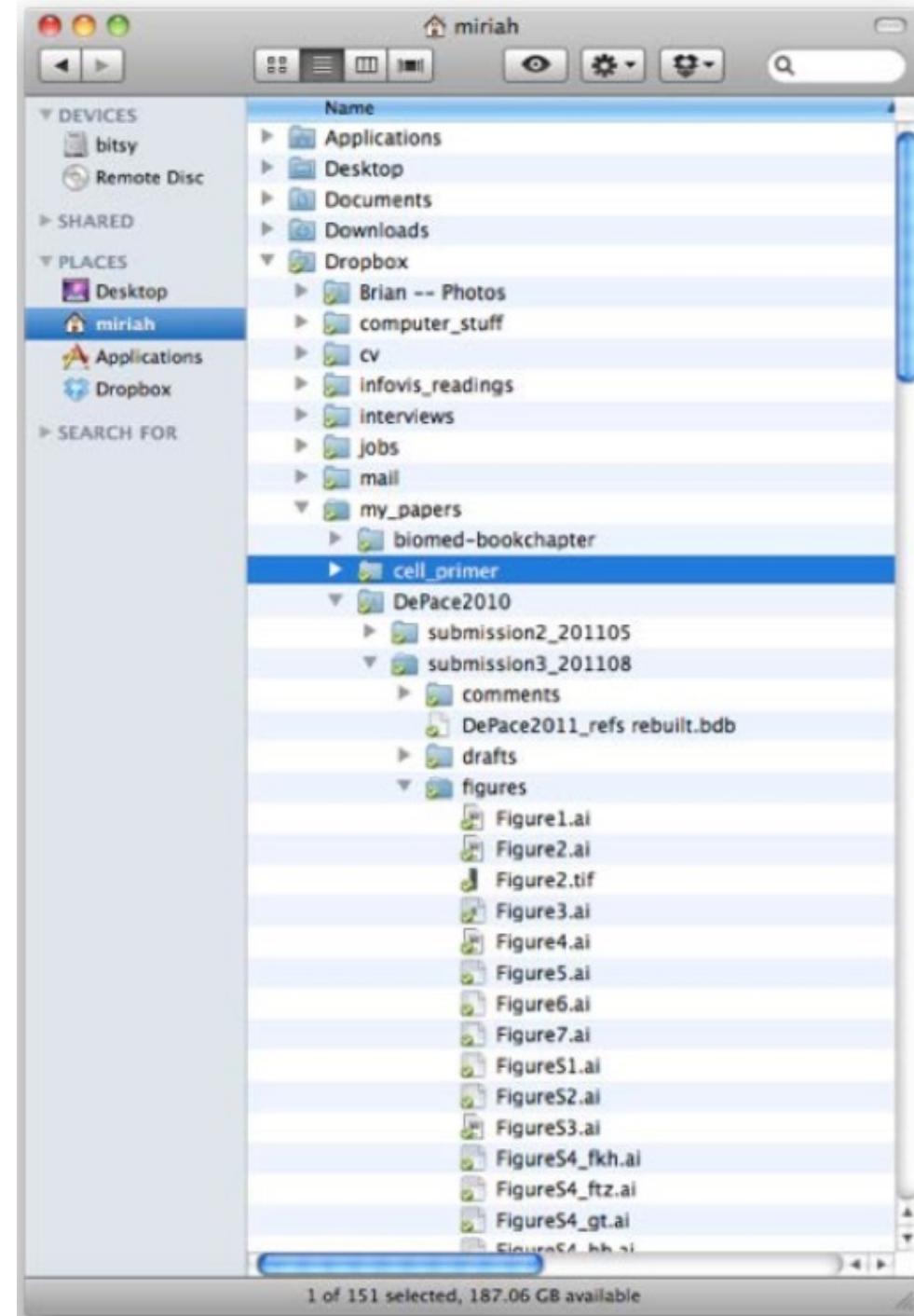


Also known as circle packing, Images from DataVizCatalogue.com and
<https://bl.ocks.org/mbostock/7607535>

Indented Trees

- Parent-child relationships shown via indentation
- Trade-off between breadth (one level) and depth (to the leaves)

```
root
├── dir1
├── dir2
│   └── file1
└── dir3
    ├── file2
    ├── file3
    └── dir4
        └── file4
└── file5
```

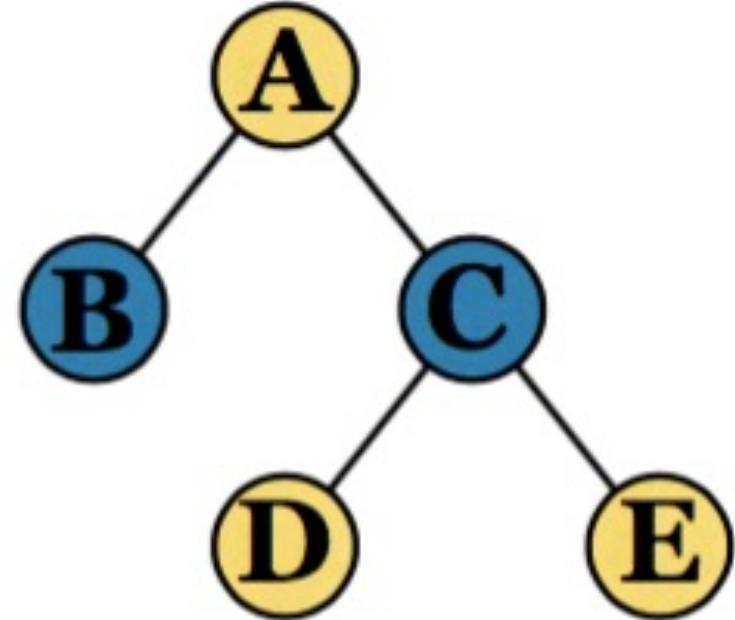


Layered Trees

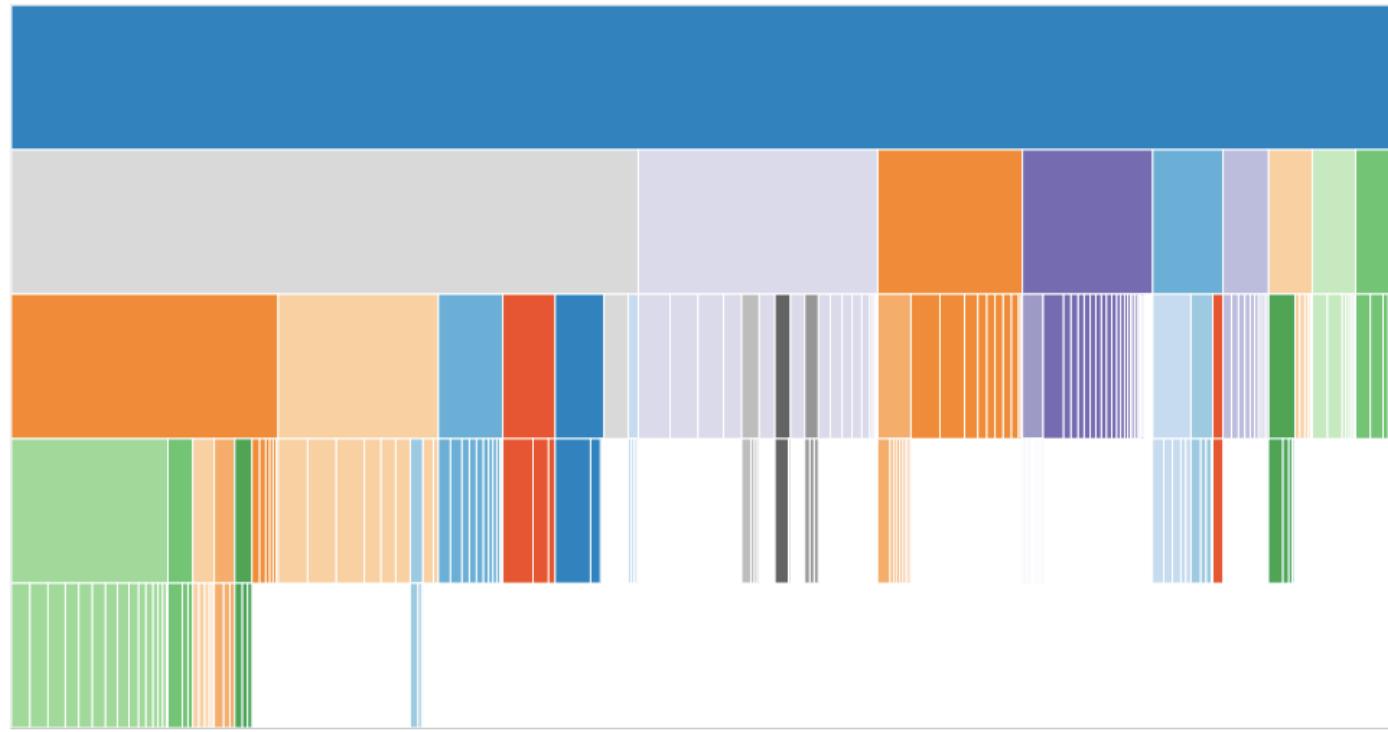
Parent-child relationships shown via layering, adjacency, & alignment

Layout similar to node-link without edges

Extent of parent (e.g., length, angle) constrains extent of children, similar to enclosure



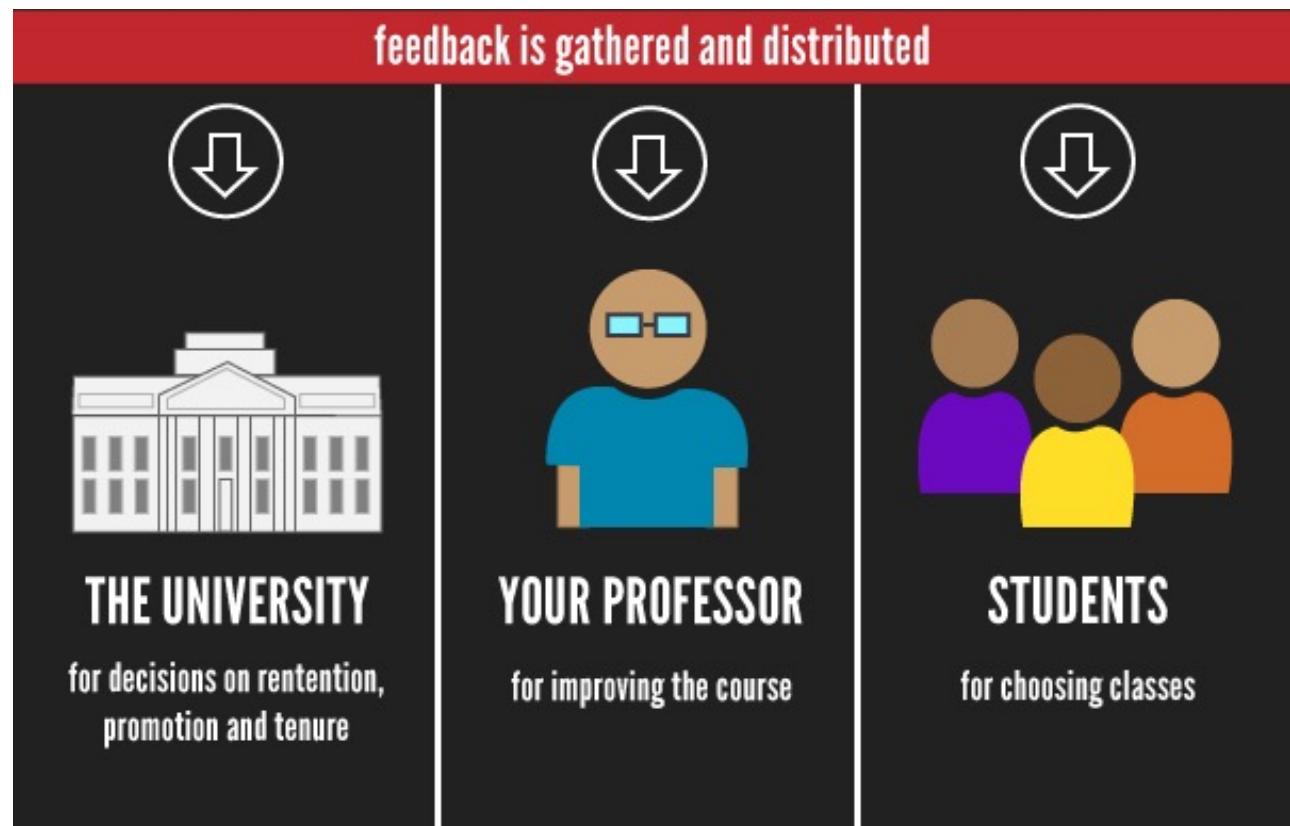
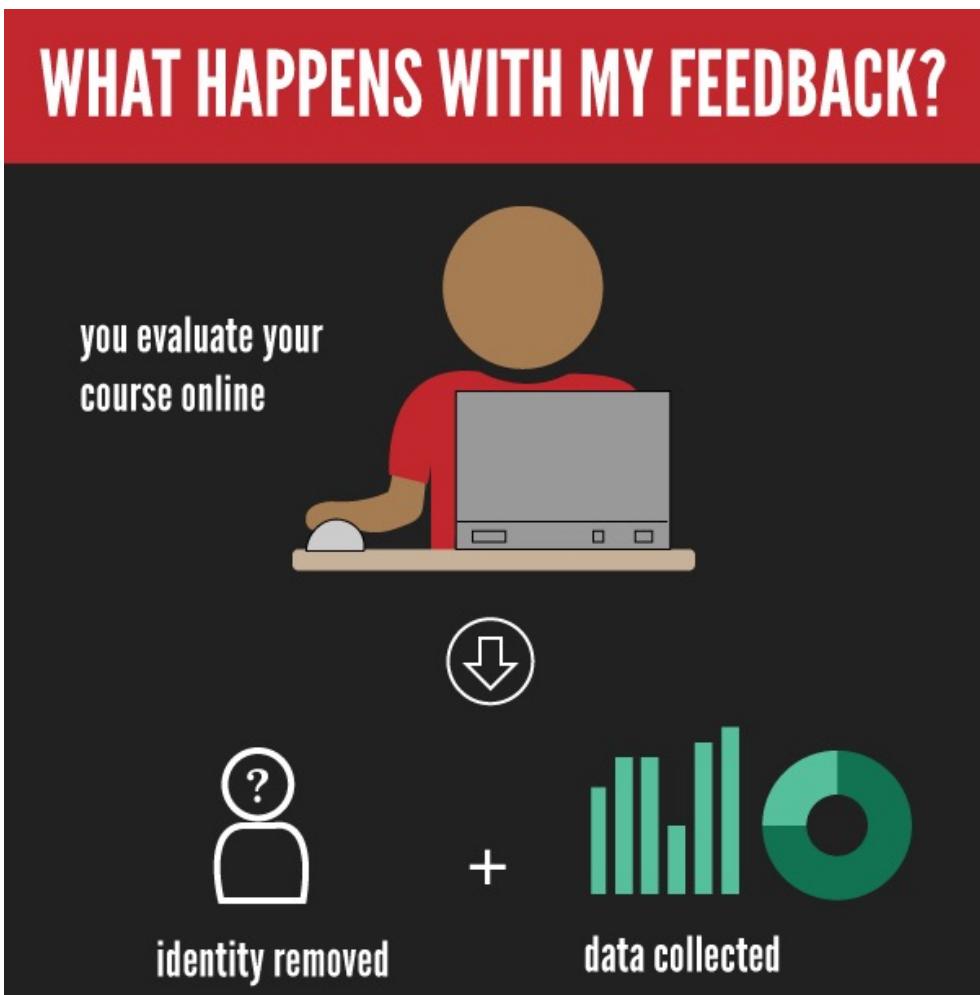
Icicle charts and sunbursts use layering



Images from <https://bl.ocks.org/mbostock/1005873>,
<https://bl.ocks.org/mbostock/4348373>

Student Course Feedback

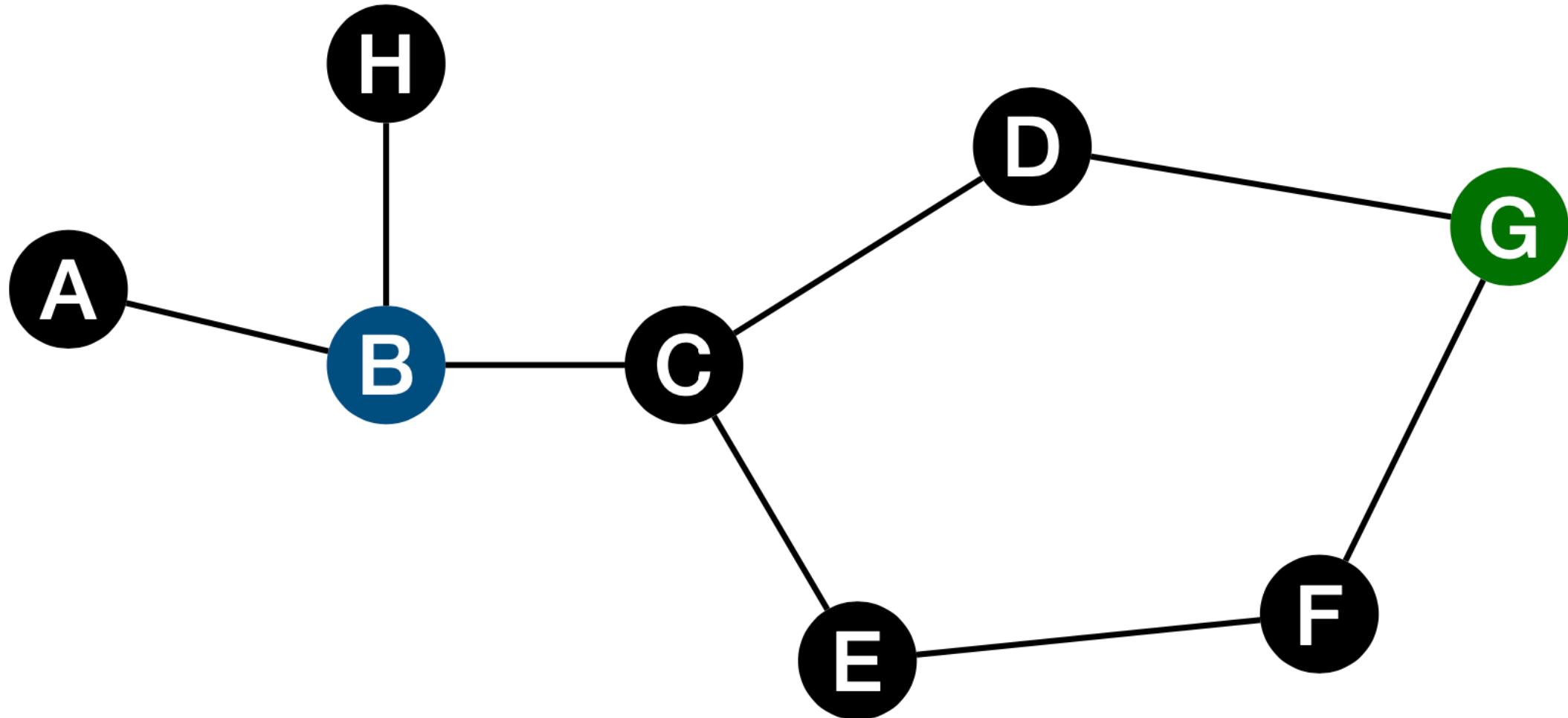
<https://scf.utah.edu>



Shortest Path Algorithms

Breadth-First Search Example

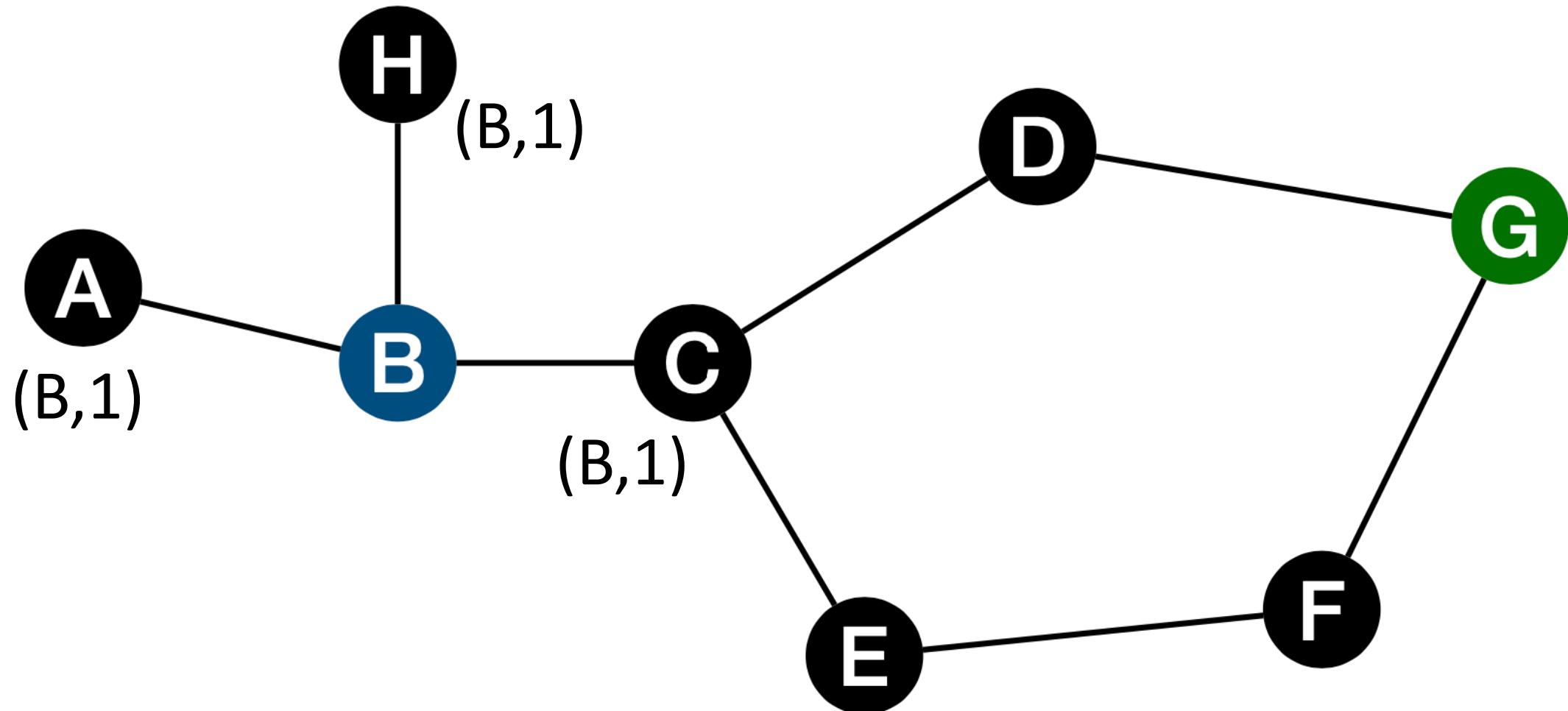
Shortest path from B to G



[B]

Breadth-First Search Example

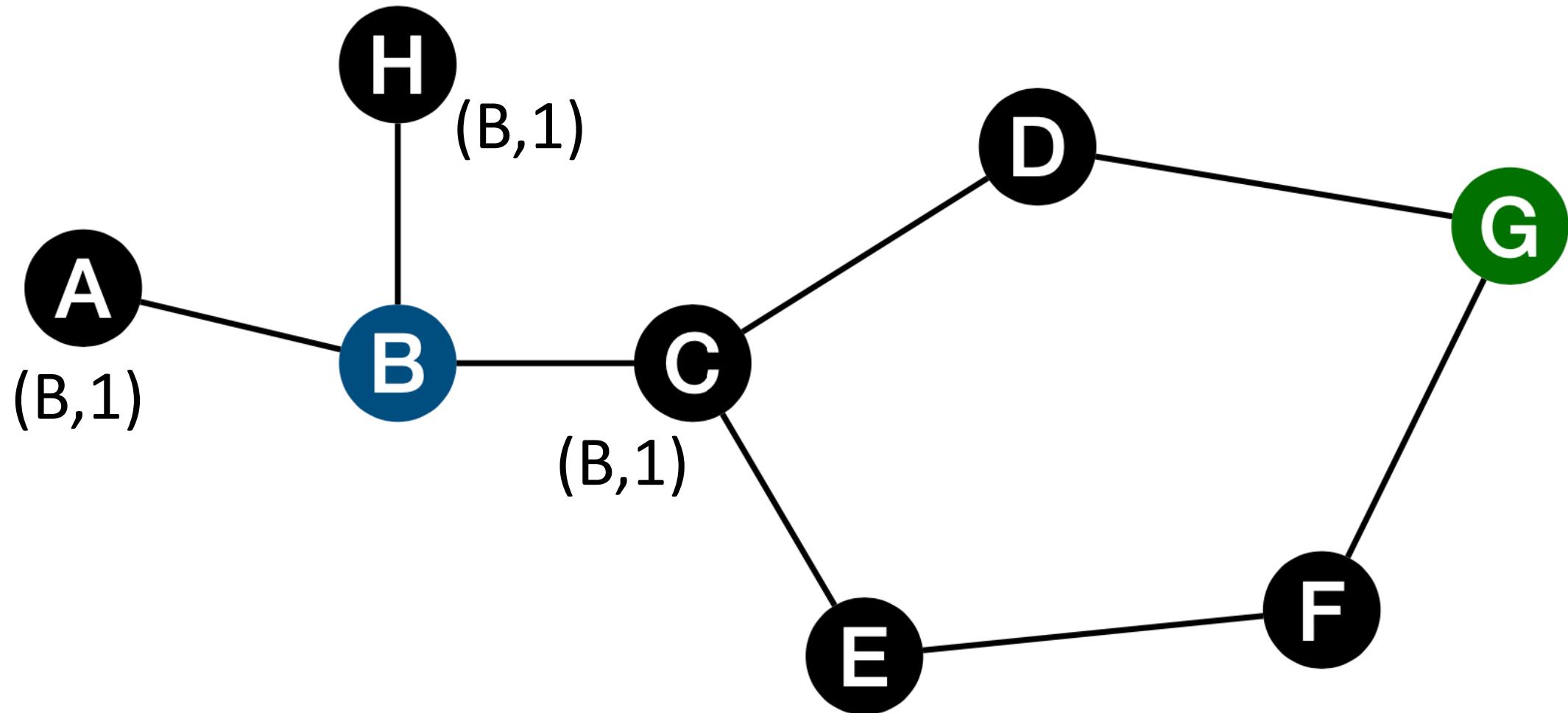
Shortest path from B to G



[B, A, H, C]

Breadth-First Search Example

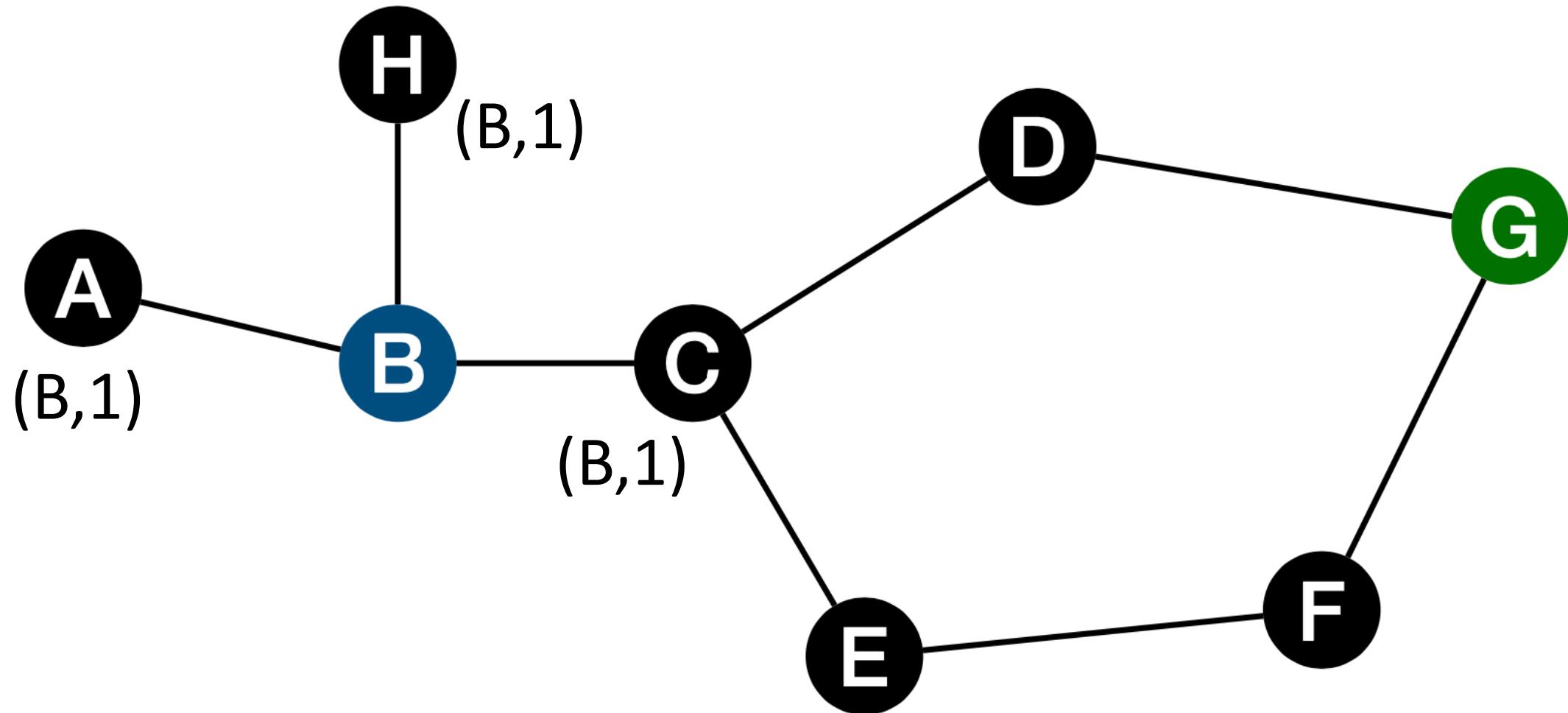
Shortest path from B to G



[B, A, H, C]

Breadth-First Search Example

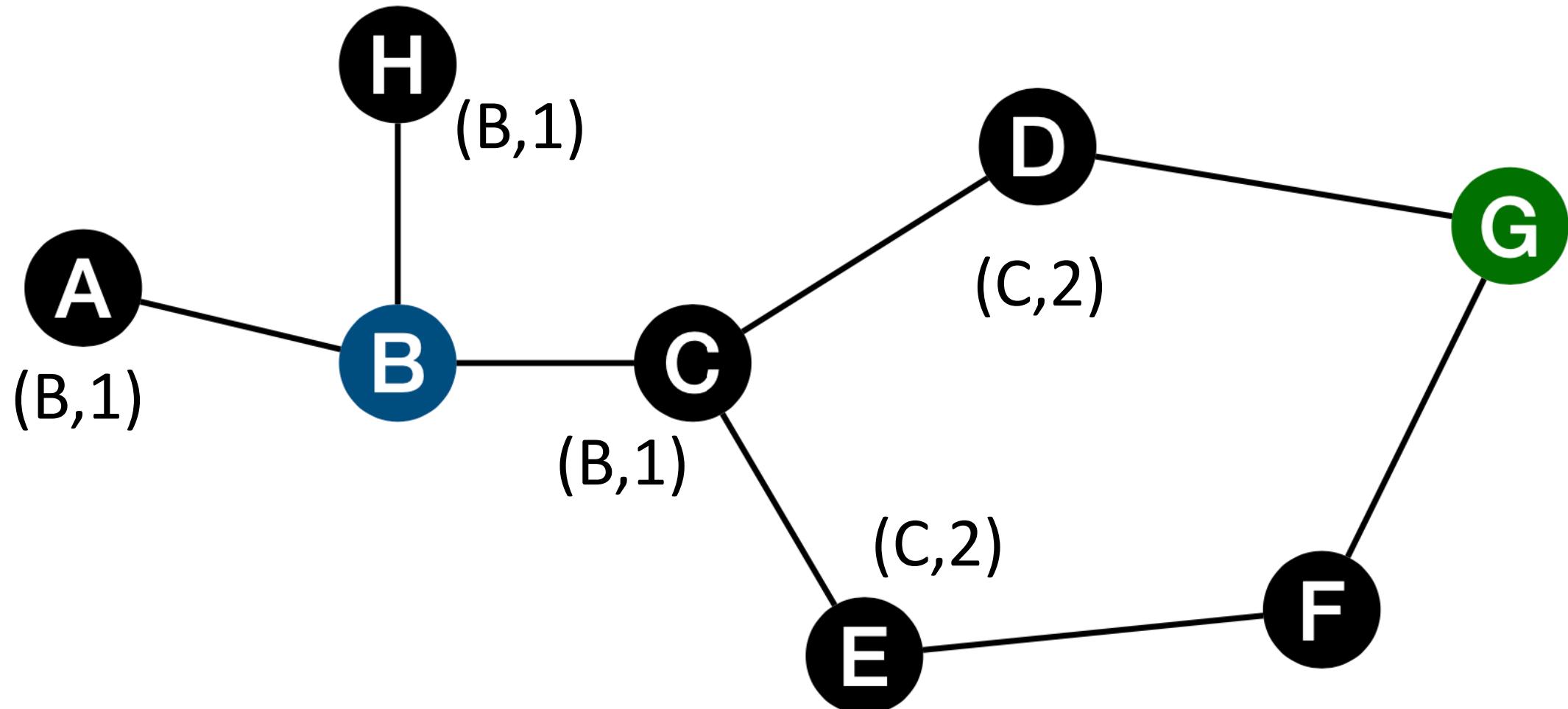
Shortest path from B to G



[B, A, H, C]

Breadth-First Search Example

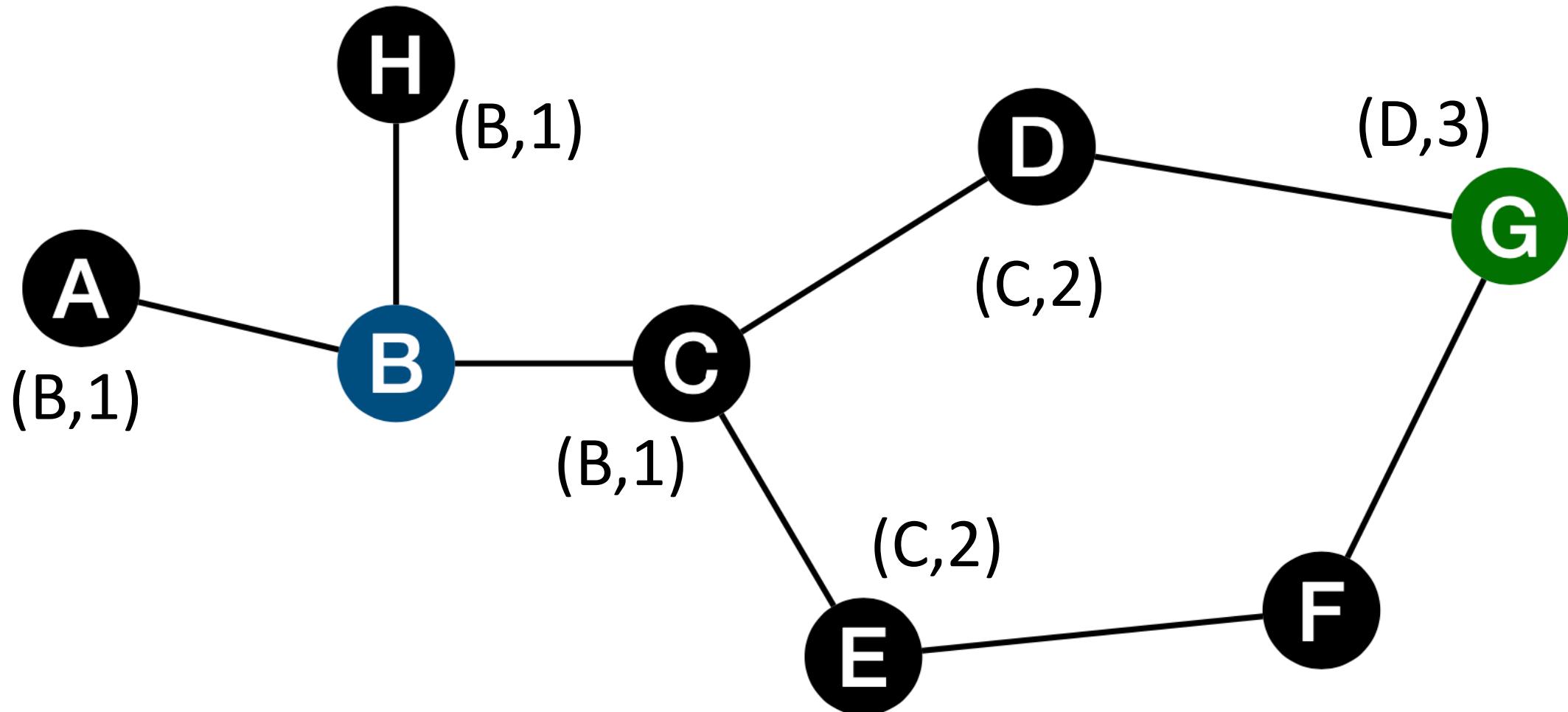
Shortest path from B to G



[B, A, H, C, D, E]

Breadth-First Search Example

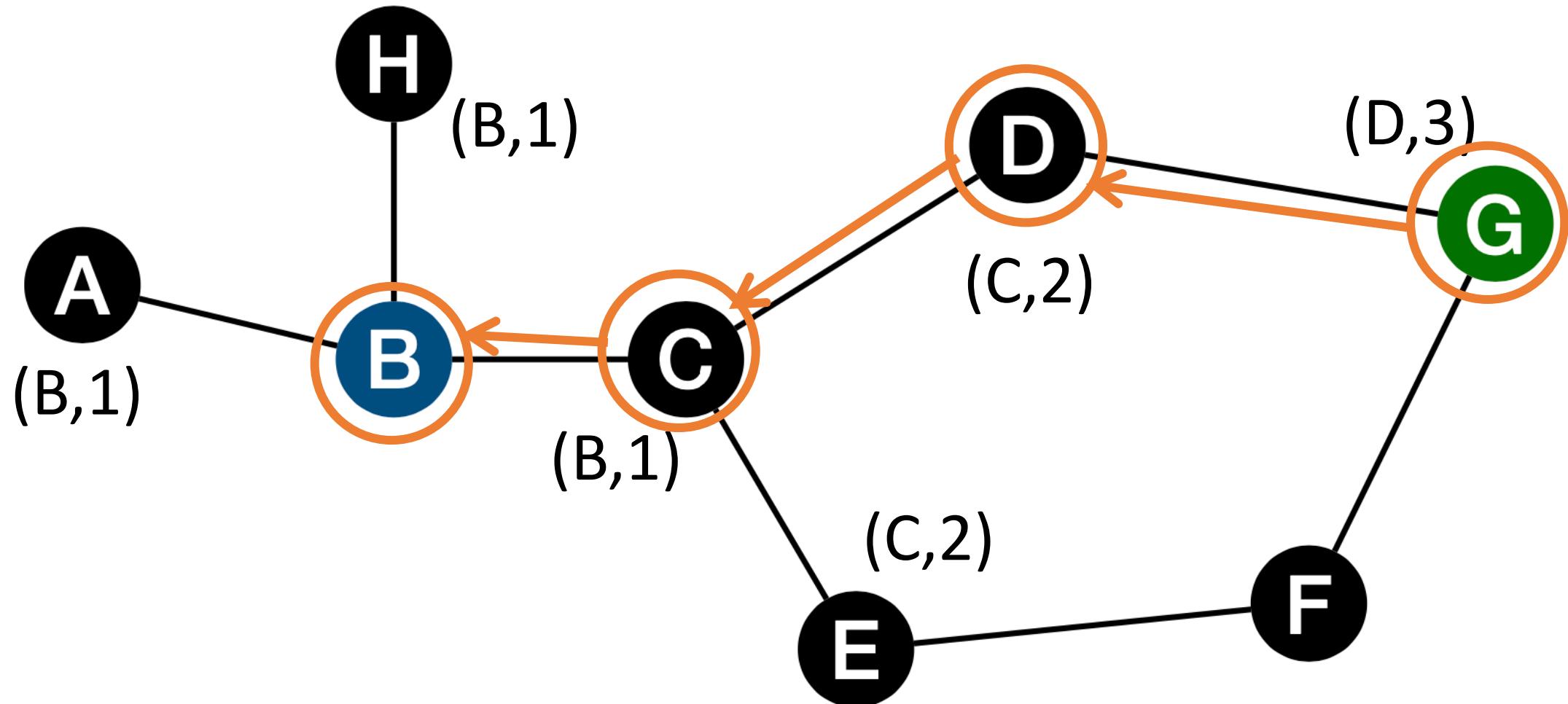
Shortest path from B to G



[B, A, H, C, D, E, G]

Breadth-First Search Example

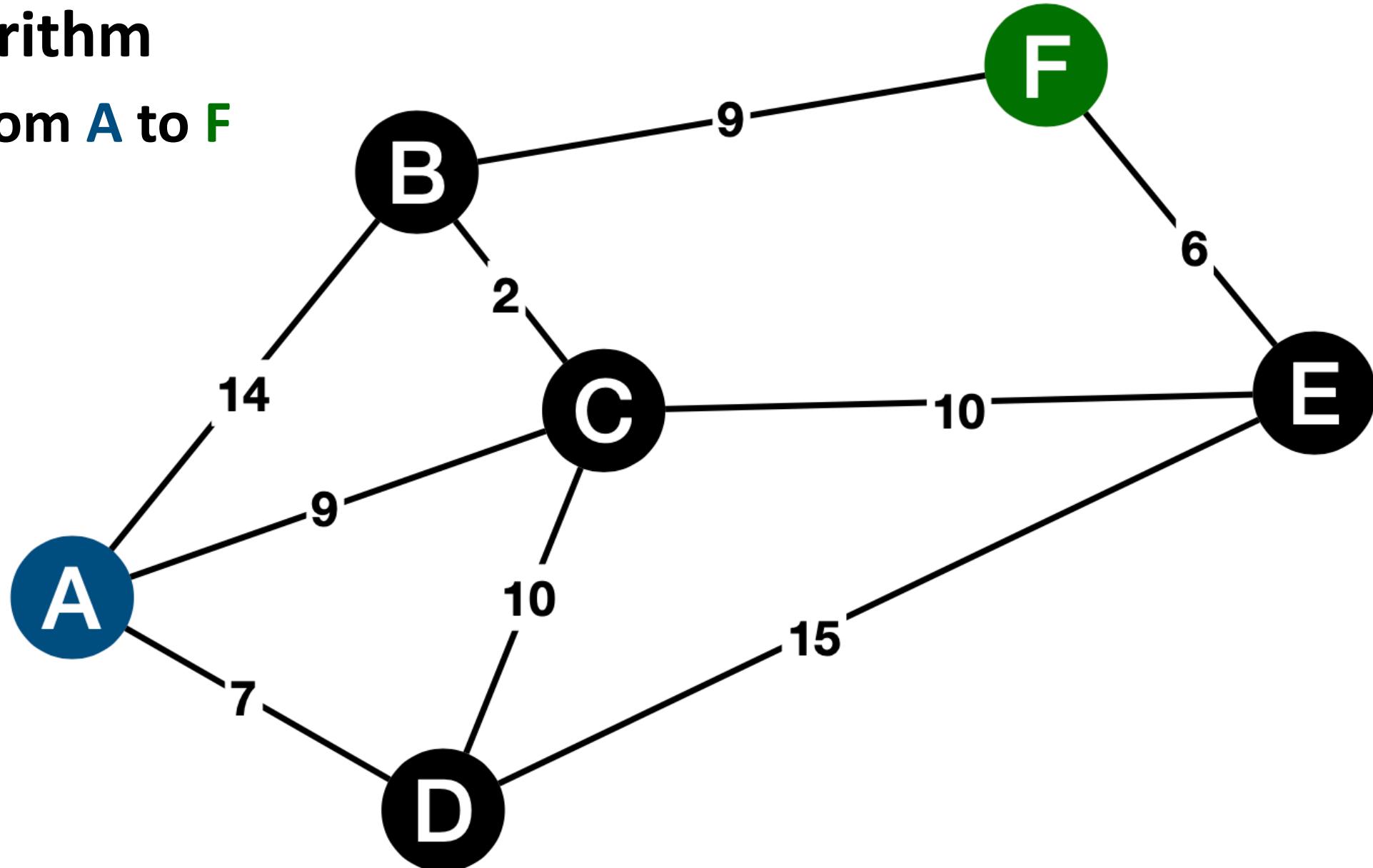
Shortest path from B to G



[B, A, H, C, D, E, G]

Dijkstra's Algorithm

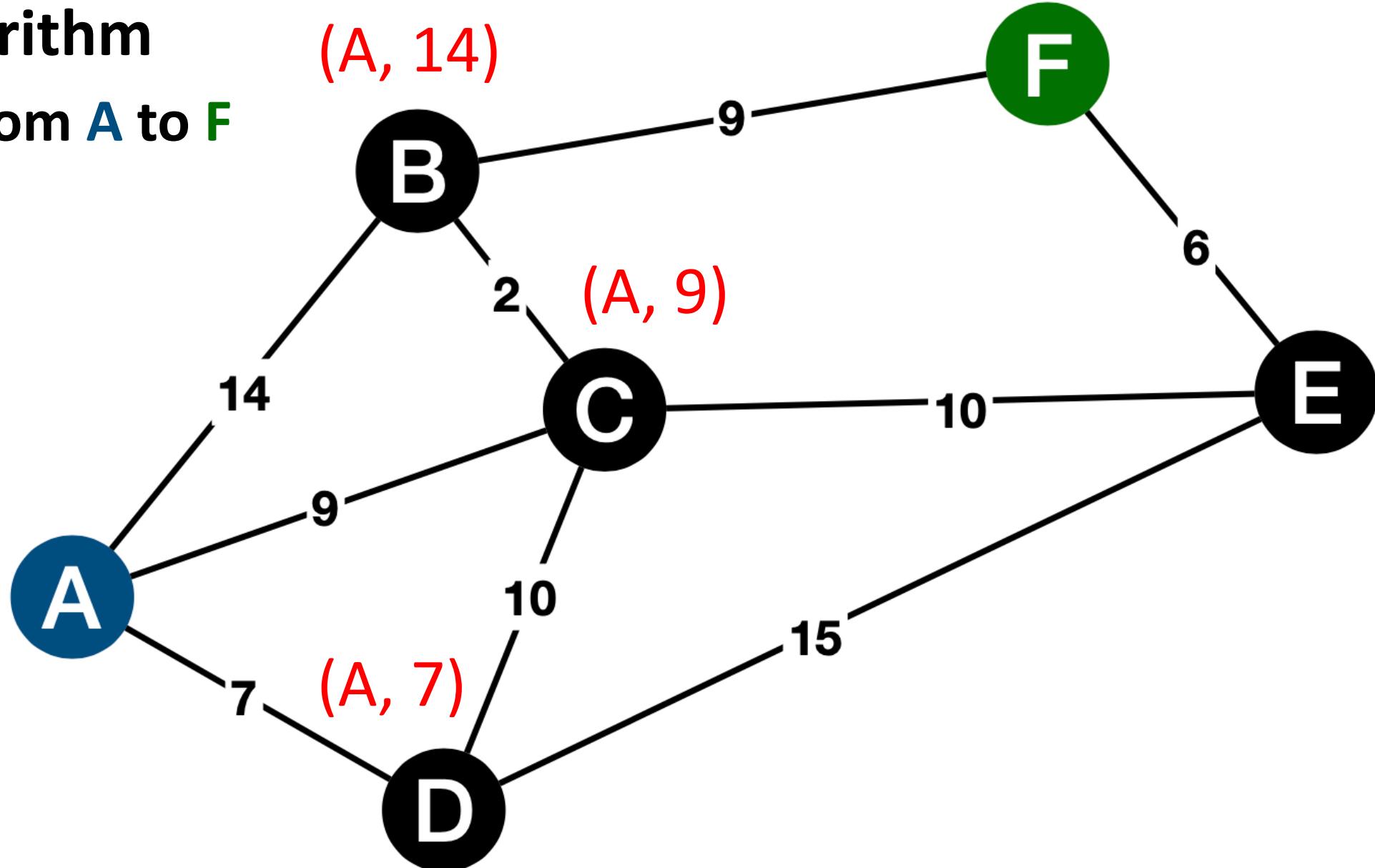
Shortest path from A to F



[A]

Dijkstra's Algorithm

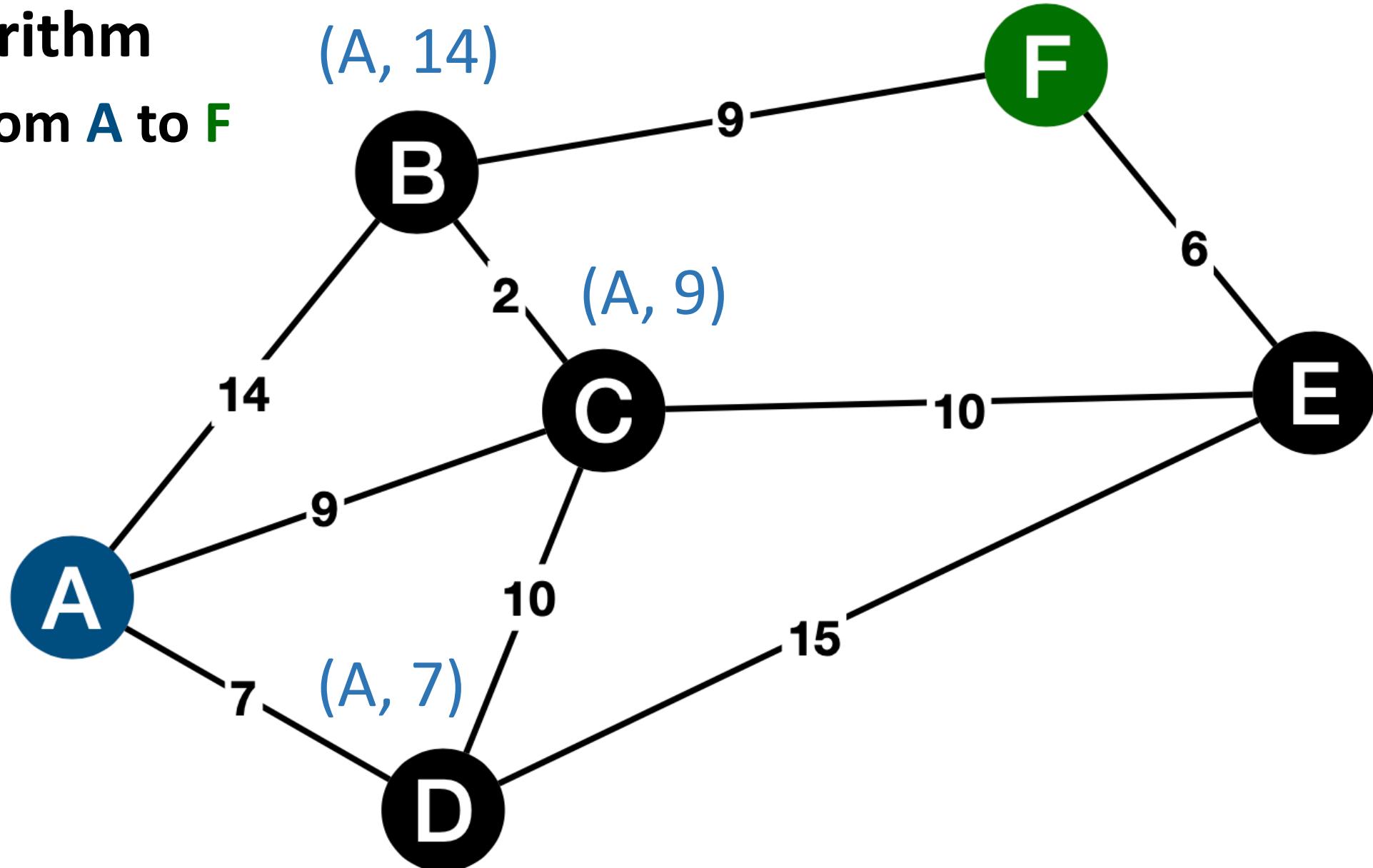
Shortest path from A to F



[A, D(7), C(9), B(14)]

Dijkstra's Algorithm

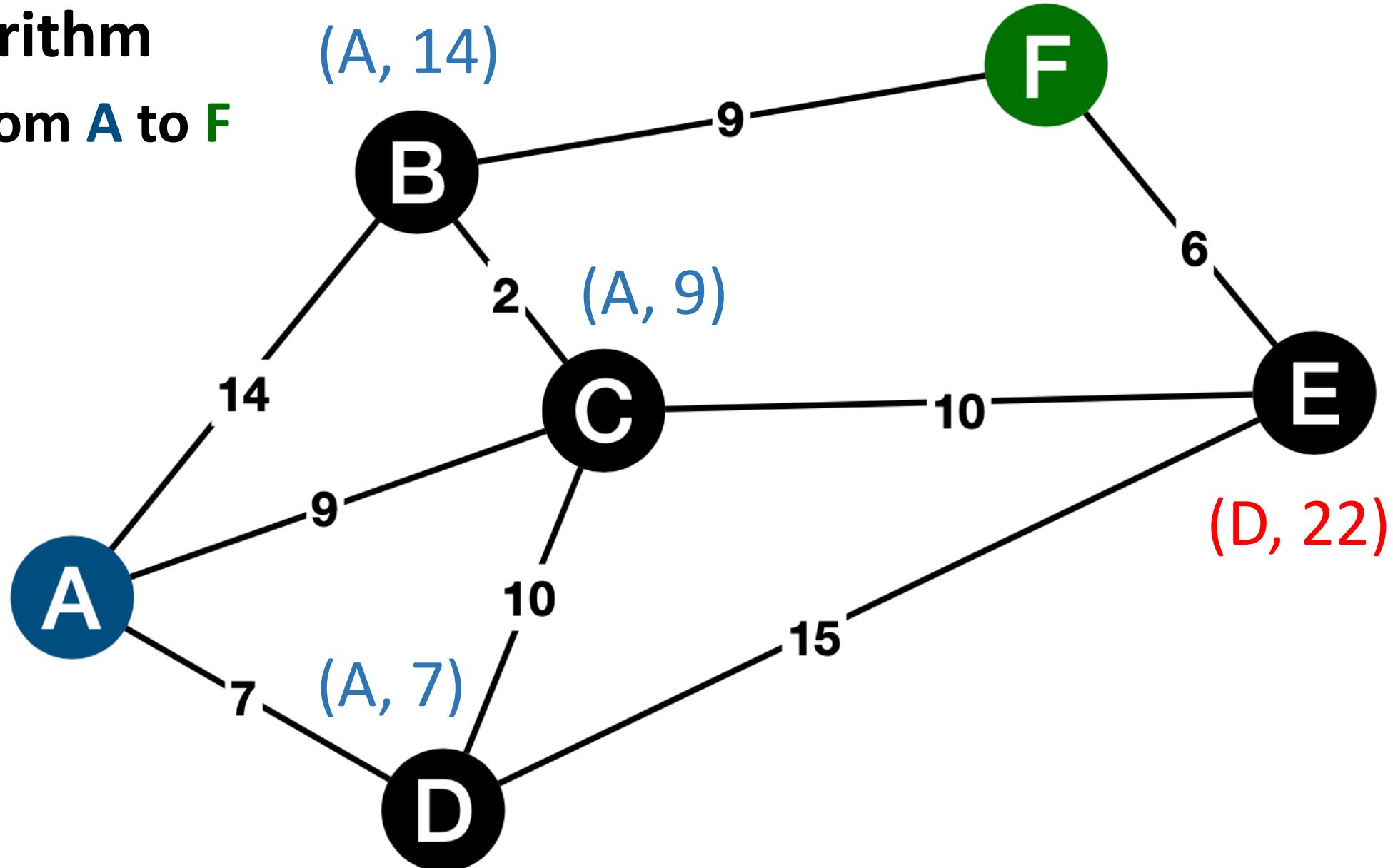
Shortest path from A to F



[A, D(7), C(9), B(14)]

Dijkstra's Algorithm

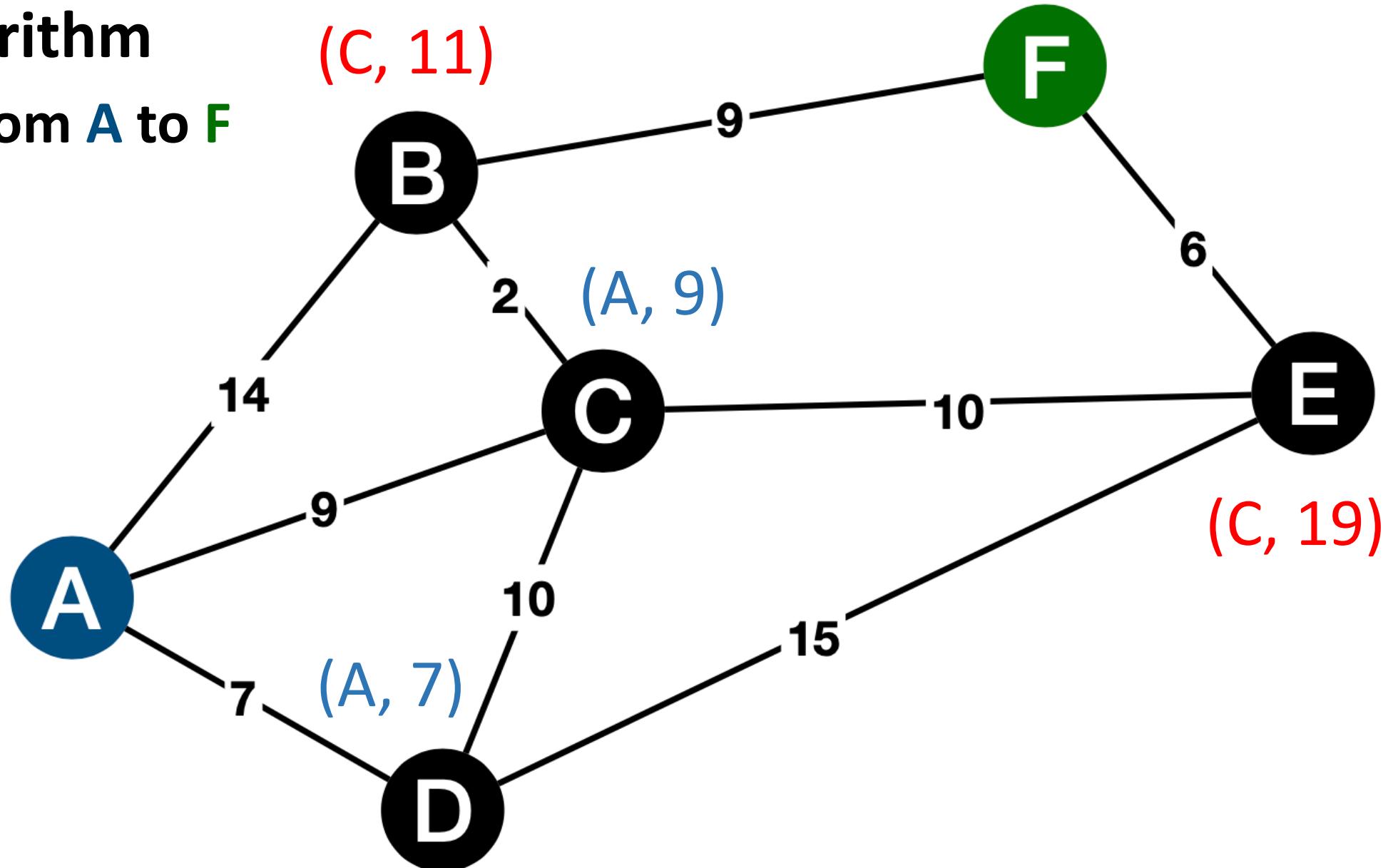
Shortest path from A to F



[A, D(7), C(9), B(14), E(22)]

Dijkstra's Algorithm

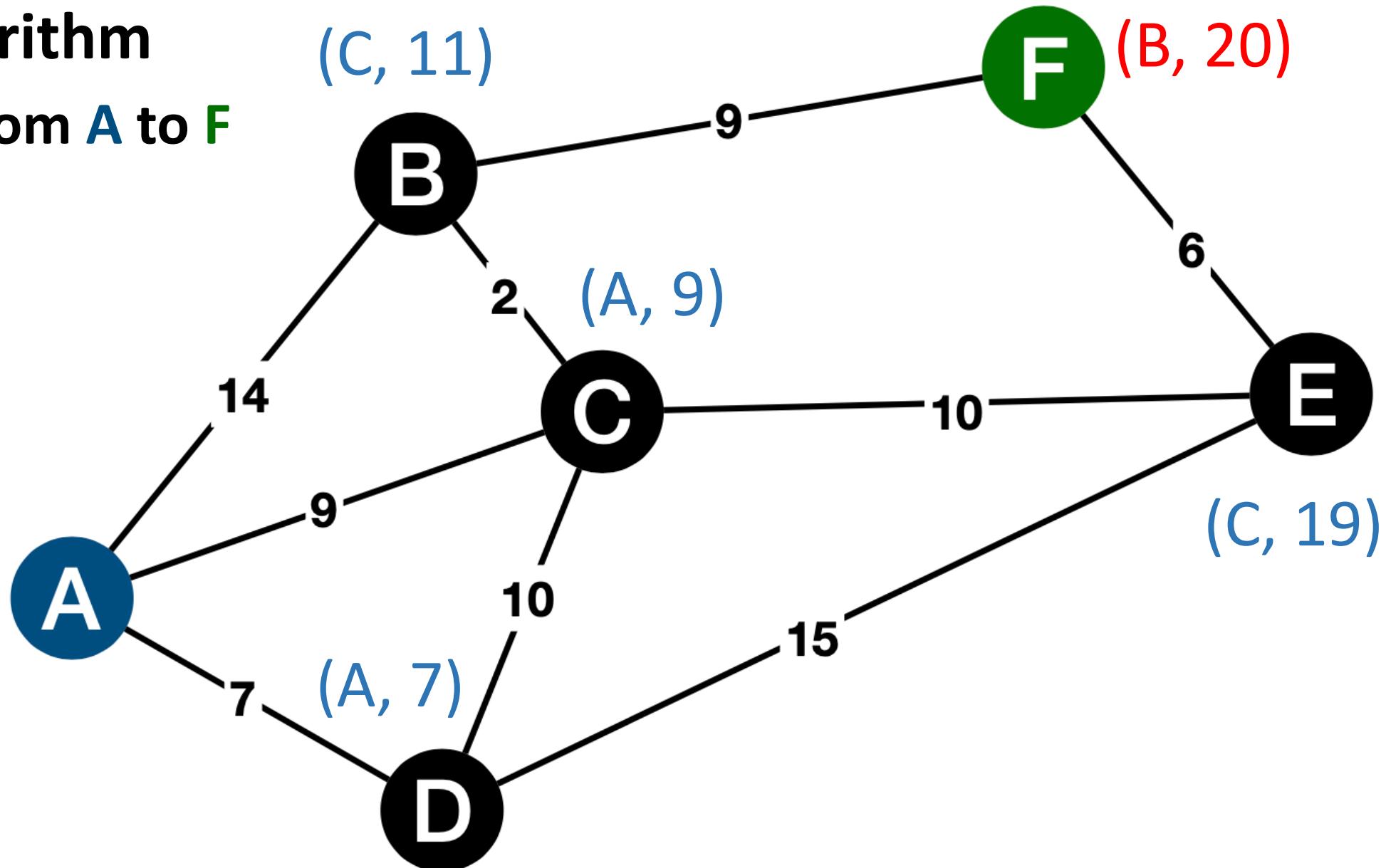
Shortest path from A to F



[A, D(7), C(9), B(11), E(19)]

Dijkstra's Algorithm

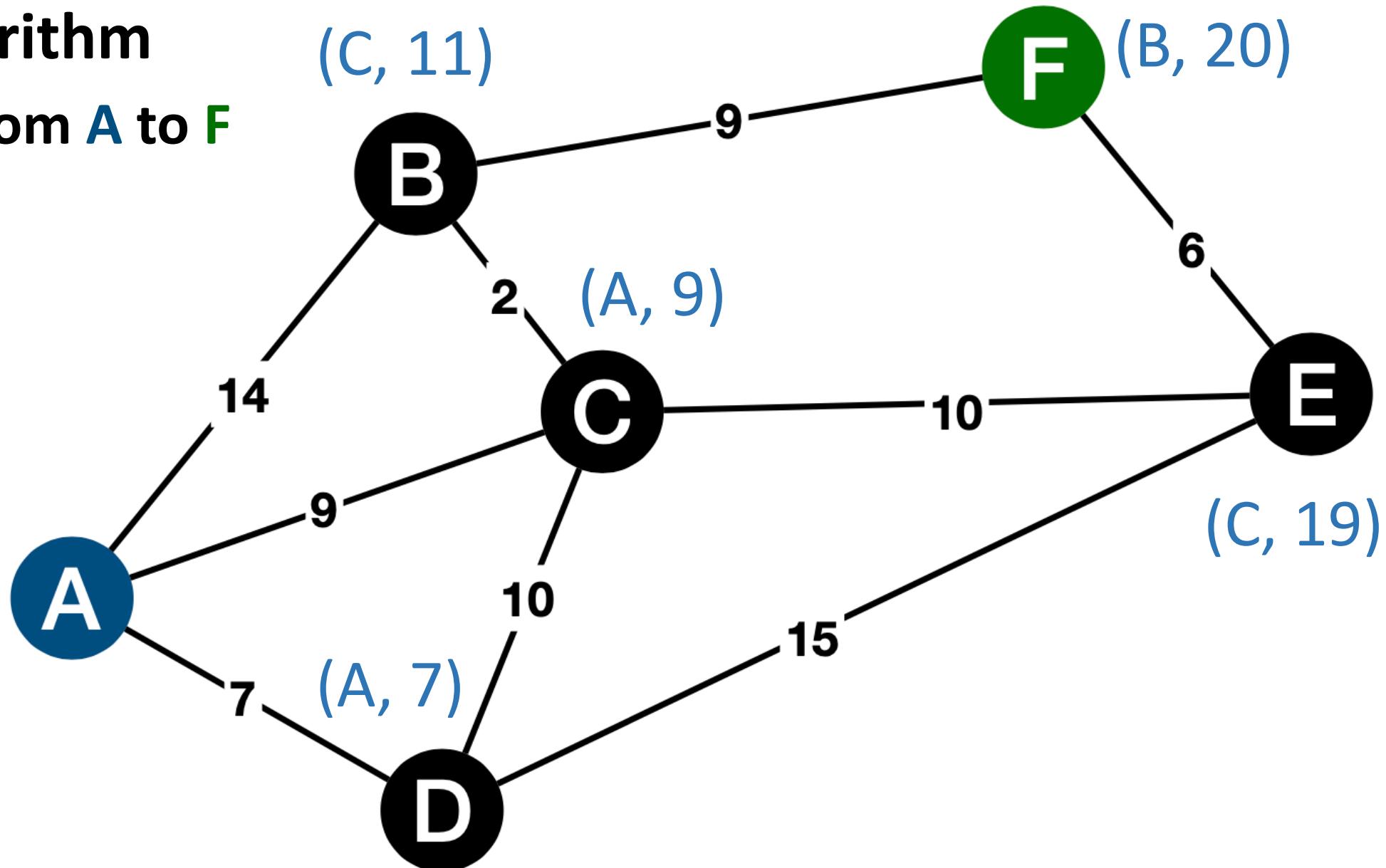
Shortest path from A to F



[A, D(7), C(9), B(11), E(19), F(20)]

Dijkstra's Algorithm

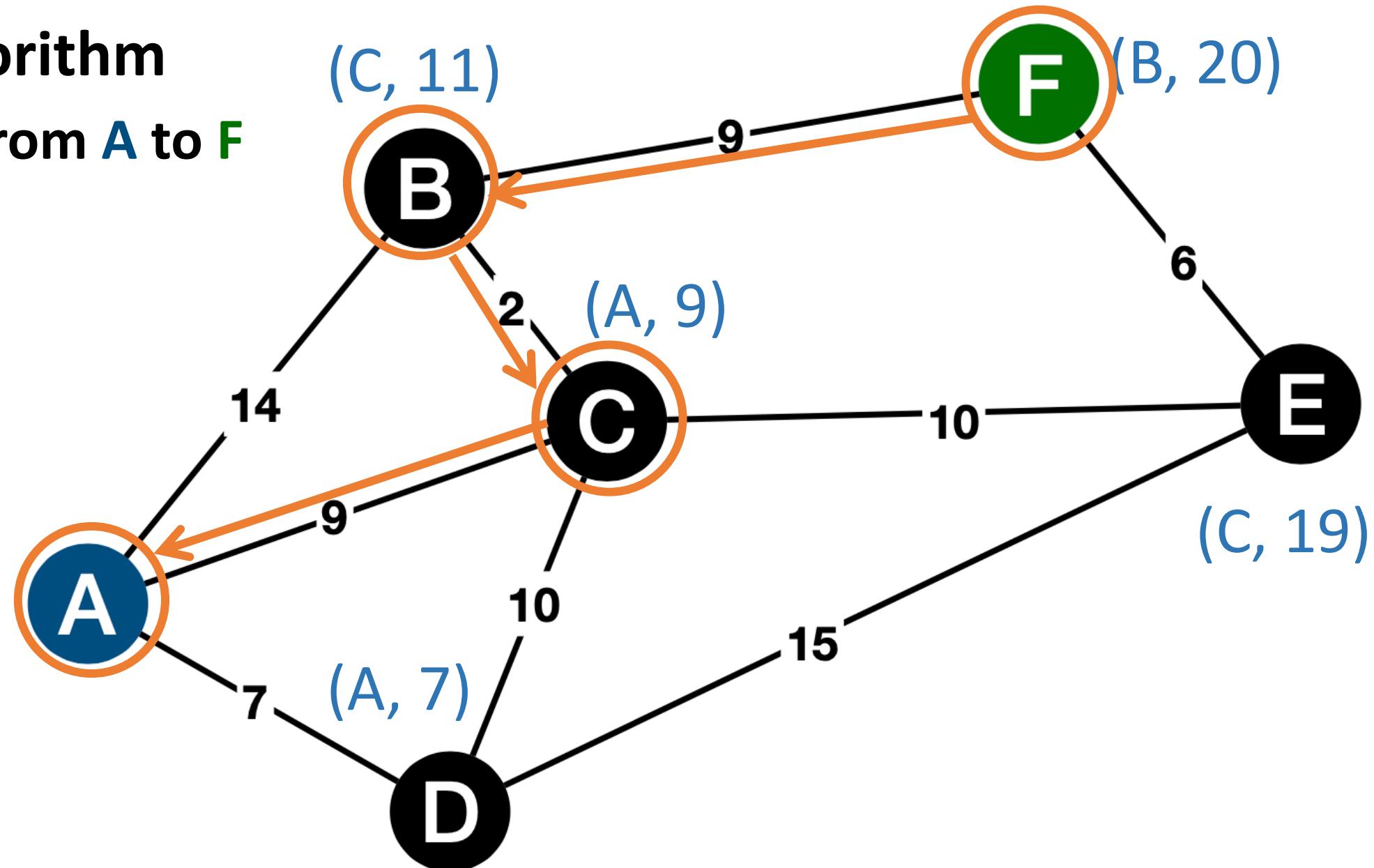
Shortest path from A to F



[A, D(7), C(9), B(11), E(19), F(20)]

Dijkstra's Algorithm

Shortest path from A to F



[A, ~~D(7)~~, ~~C(9)~~, ~~B(11)~~, ~~E(19)~~, ~~F(20)~~]