

# COMP9313: Big Data Management

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro

**Lecturer: Xin Cao**

Course web site: <http://www.cse.unsw.edu.au/~cs9313/>

# Chapter **Assignment Project Exam Help** **rocessing**

<https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro

# What's a Graph?

- $G = (V, E)$ , where
  - $V$  represents the set of vertices (nodes)
  - $E$  represents the set of edges (links)
  - Both vertices and edges may contain additional information
- Different types of graphs
  - Directed vs.
  - Presence or
- Graphs are everywhere:
  - Hyperlink structure of the Web
  - Physical structure of computers on the Internet
  - Interstate highway system
  - Social networks

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro

# Graph Data: Social Networks

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro

**Facebook social graph**

4-degrees of separation [Backstrom-Boldi-Rosa-Ugander-Vigna, 2011]

# Graph Data: Technological Networks

Assignment Project Exam Help

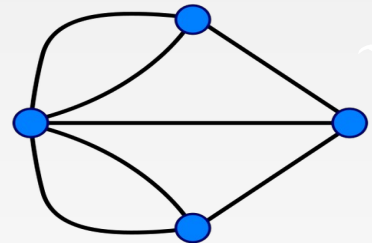
<https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro

## Seven Bridges of Königsberg

[Euler, 1735]

Return to the starting point by traveling each link of the graph once and only once.



# Some Graph Problems

- Finding shortest paths
    - Routing Internet traffic and UPS trucks
  - Finding minimum spanning trees
    - Telco laying down fiber
  - Finding Max Flow
  - Identify “special”
  - Bipartite matching
  - And of course... PageRank
- Assignment Project Exam Help  
<https://eduassistpro.github.io/>  
Add WeChat edu\_assist\_pro

# Graph Analytics

- General Graph
  - Count the number of nodes whose degree is equal to 5
  - Find the diameter of the graphs
- Web Graph
  - Rank each web page in the web graph or each user in the twitter graph using **Assignment Project Exam Help**
- Transportation **<https://eduassistpro.github.io/>**
  - Return the shortest or cheapest **Add WeChat edu\_assist\_pro** one city to another
- Social Network
  - Detect a group of users who have similar interests
- Financial Network
  - Find the path connecting two suspicious transactions;
- ... ..

# Graphs and MapReduce

- Graph algorithms typically involve:
  - Performing computations at each node: based on node features, edge features, and local link structure
  - Propagating computations: “traversing” the graph
- Key questions:
  - How do you reduce?
  - How do you <https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro



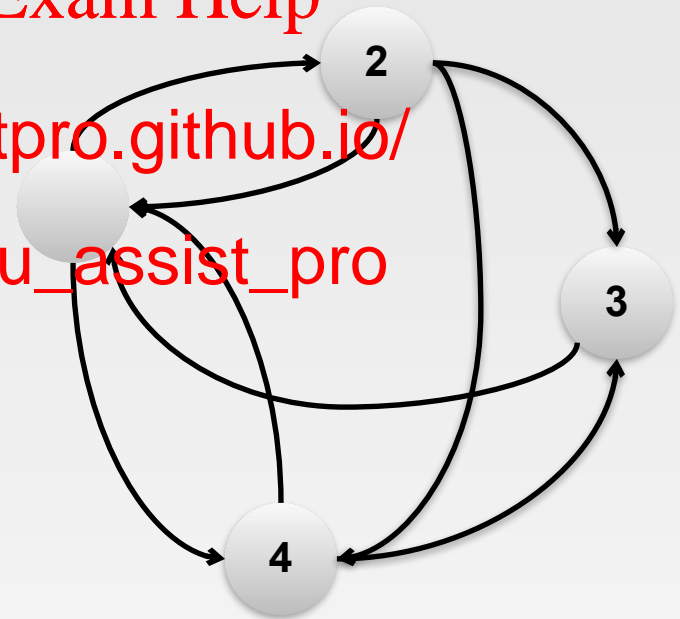
# Representing Graphs

- Adjacency Matrices: Represent a graph as an  $n \times n$  square matrix  $M$ 
  - $n = |V|$
  - $M_{ij} = 1$  means a link from node  $i$  to  $j$

## Assignment Project Exam Help

	1	2		
1	0	1	0	1
2	1	0	1	1
3	1	0	0	0
4	1	0	1	0

<https://eduassistpro.github.io/>  
Add WeChat edu\_assist\_pro



# Adjacency Matrices: Critique

- Advantages:
  - Amenable to mathematical manipulation
  - Iteration over rows and columns corresponds to computations on outlinks and inlinks
- Disadvantages:
  - Lots of zero
  - Lots of wast

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro

# Representing Graphs

- Adjacency Lists: Take adjacency matrices... and throw away all the zeros

## Assignment Project Exam Help

	1	2	3	4
1	0	1	0	1
2	1	0	1	1
3	1	0	0	0
4	1	0	1	0

<https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro



4  
, 3, 4  
3: 1  
4: 1, 3

# Adjacency Lists: Critique

- Advantages:
  - Much more compact representation
  - Easy to compute over outlinks
- Disadvantages:
  - Much more difficult to compute over inlinks

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro

Assignment Project Exam Help

**Singl** <https://eduassistpro.github.io/> **Path**

Add WeChat edu\_assist\_pro

# Single-Source Shortest Path (SSSP)

- **Problem:** find shortest path from a source node to one or more target nodes
  - Shortest might also mean lowest weight or cost
- Dijkstra's Algorithm:
  - For a given source node in the graph, the algorithm finds the shortest path between that node and every other

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro

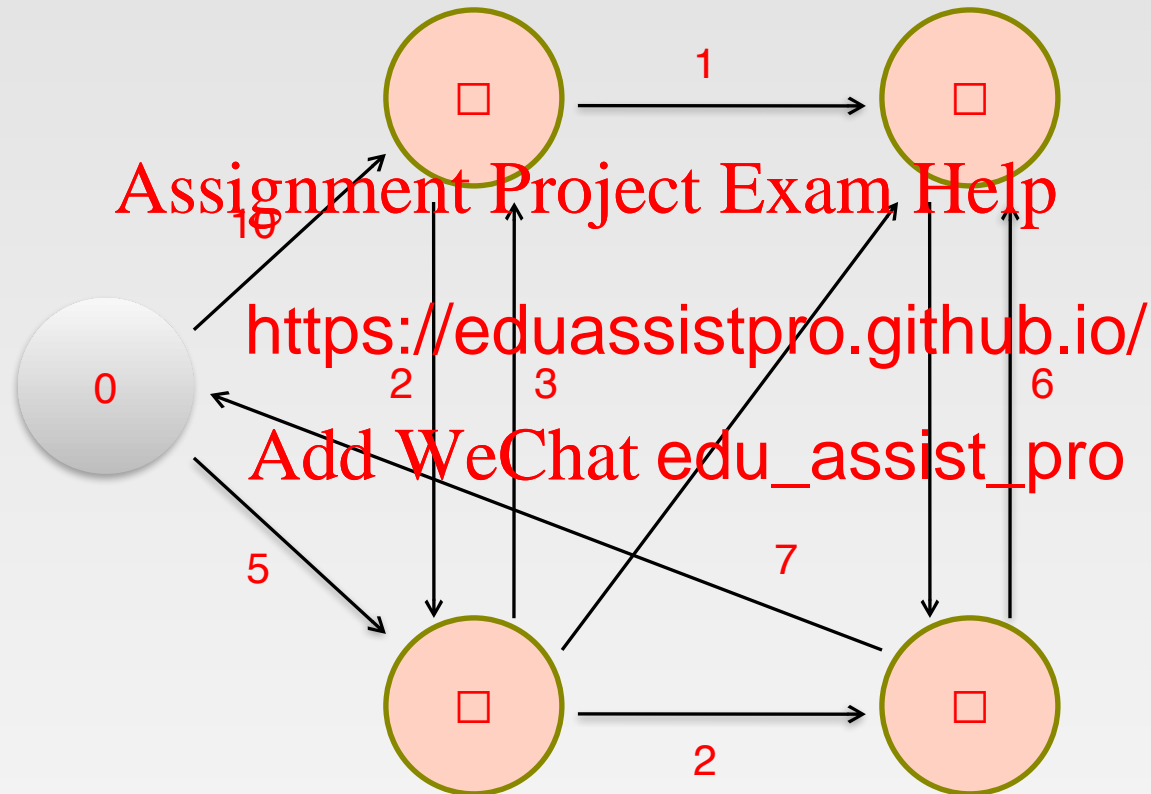
# Dijkstra's Algorithm

Assignment Project Exam Help

<https://eduassistpro.github.io/>

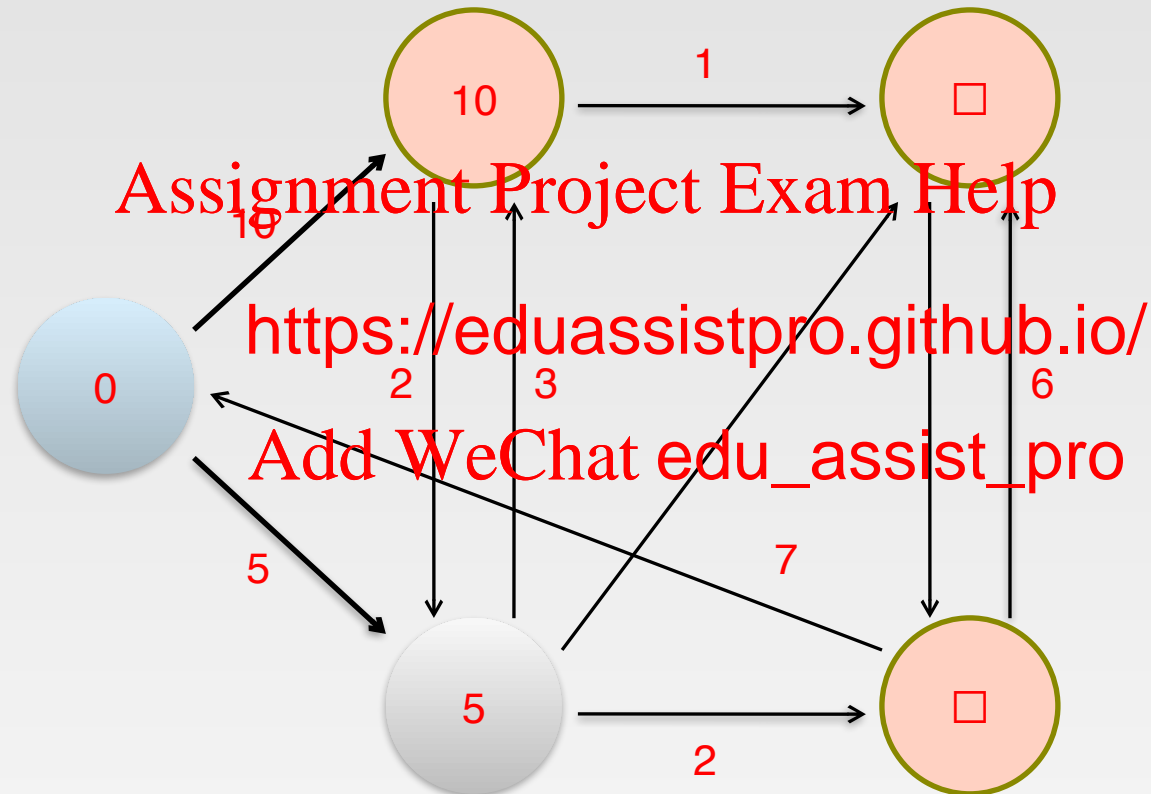
Add WeChat edu\_assist\_pro

# Dijkstra's Algorithm Example

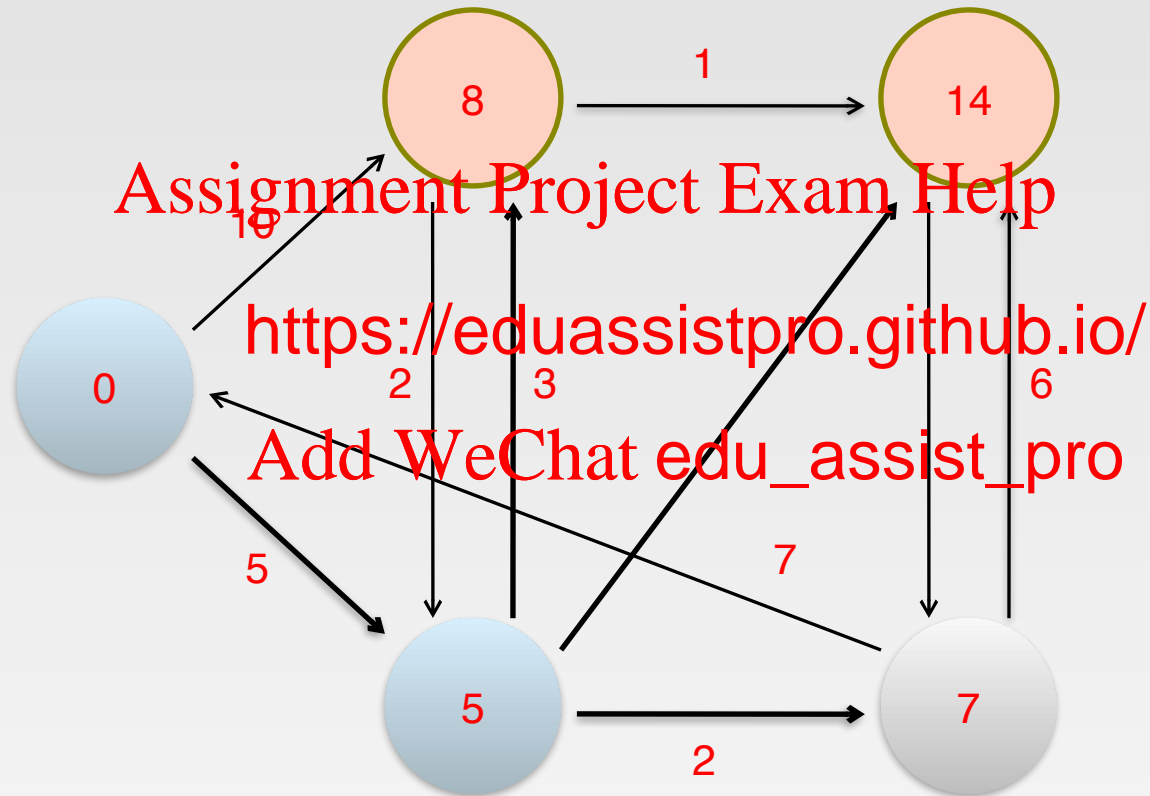




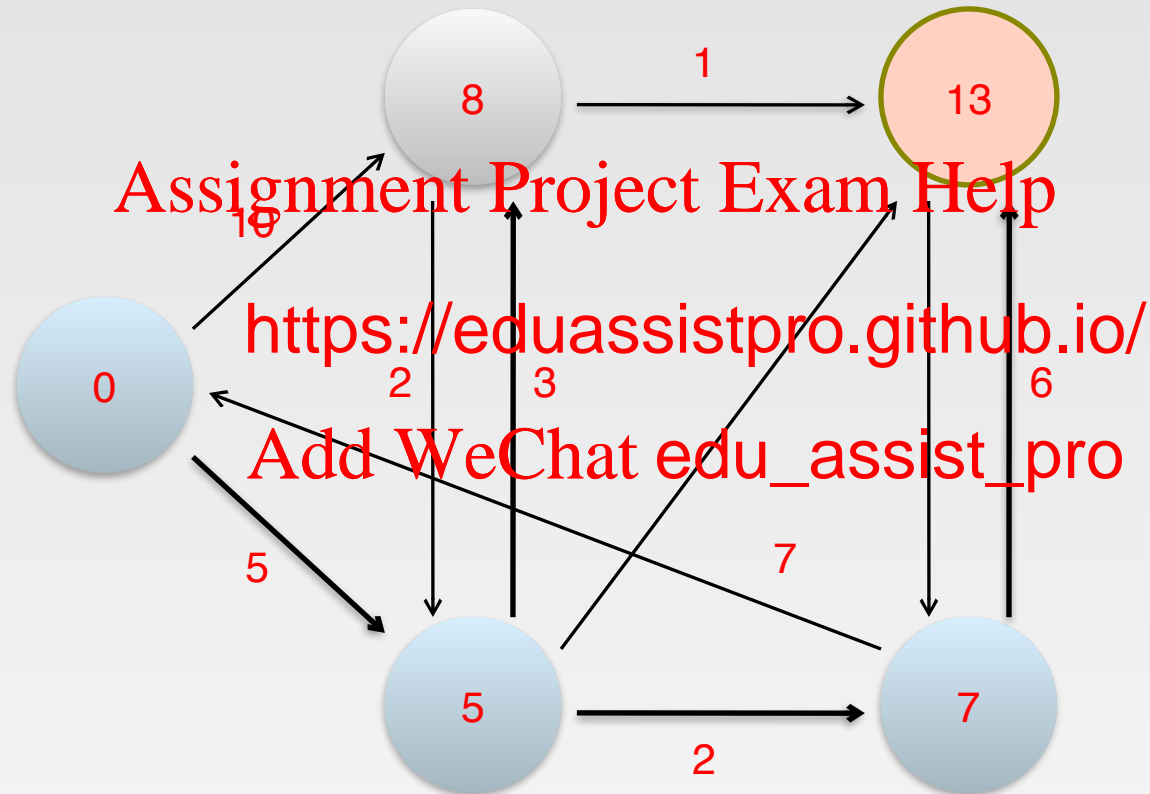
# Dijkstra's Algorithm Example



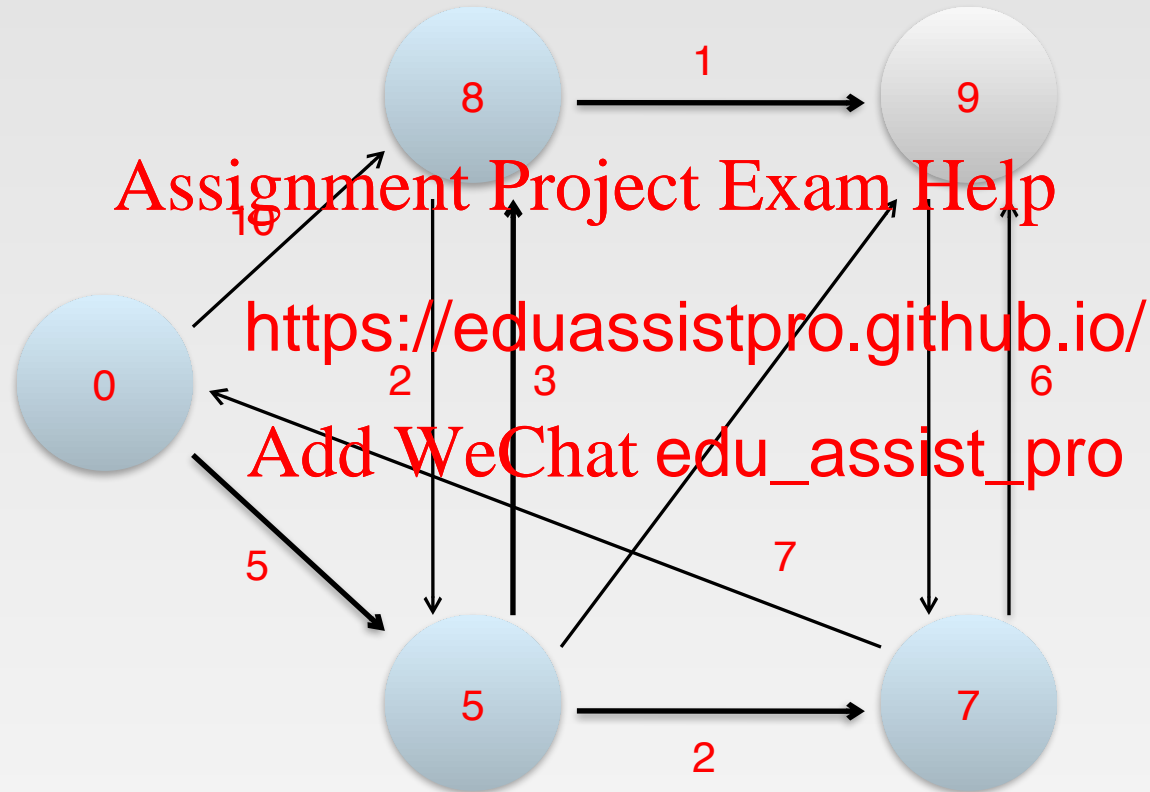
# Dijkstra's Algorithm Example



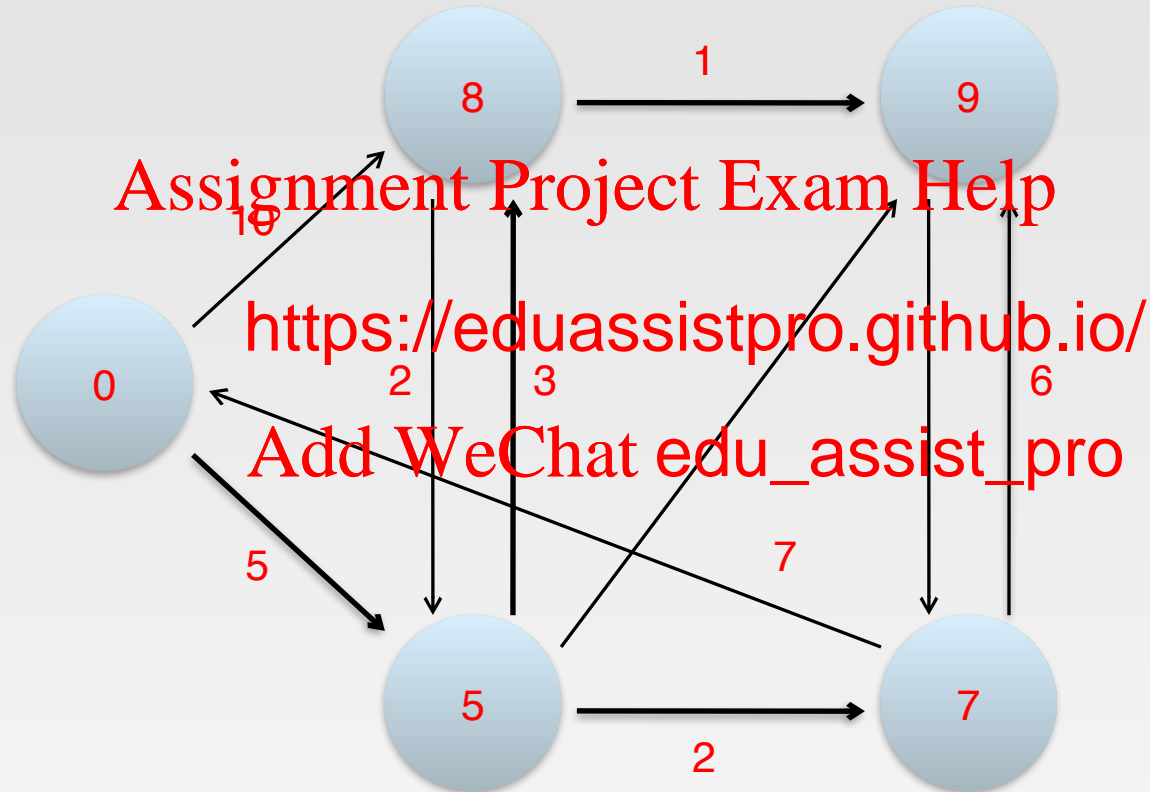
# Dijkstra's Algorithm Example



# Dijkstra's Algorithm Example



# Dijkstra's Algorithm Example



**Finish!**

# Single Source Shortest Path

- **Problem:** find shortest path from a source node to one or more target nodes
  - Shortest might also mean lowest weight or cost
- Single processor machine: Dijkstra's Algorithm
- MapReduce: parallel Breadth First Search (BFS)

Assignment Project Exam Help

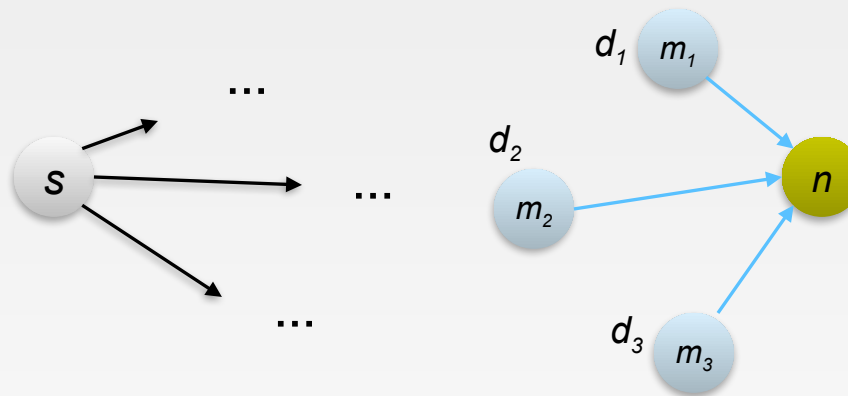
<https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro

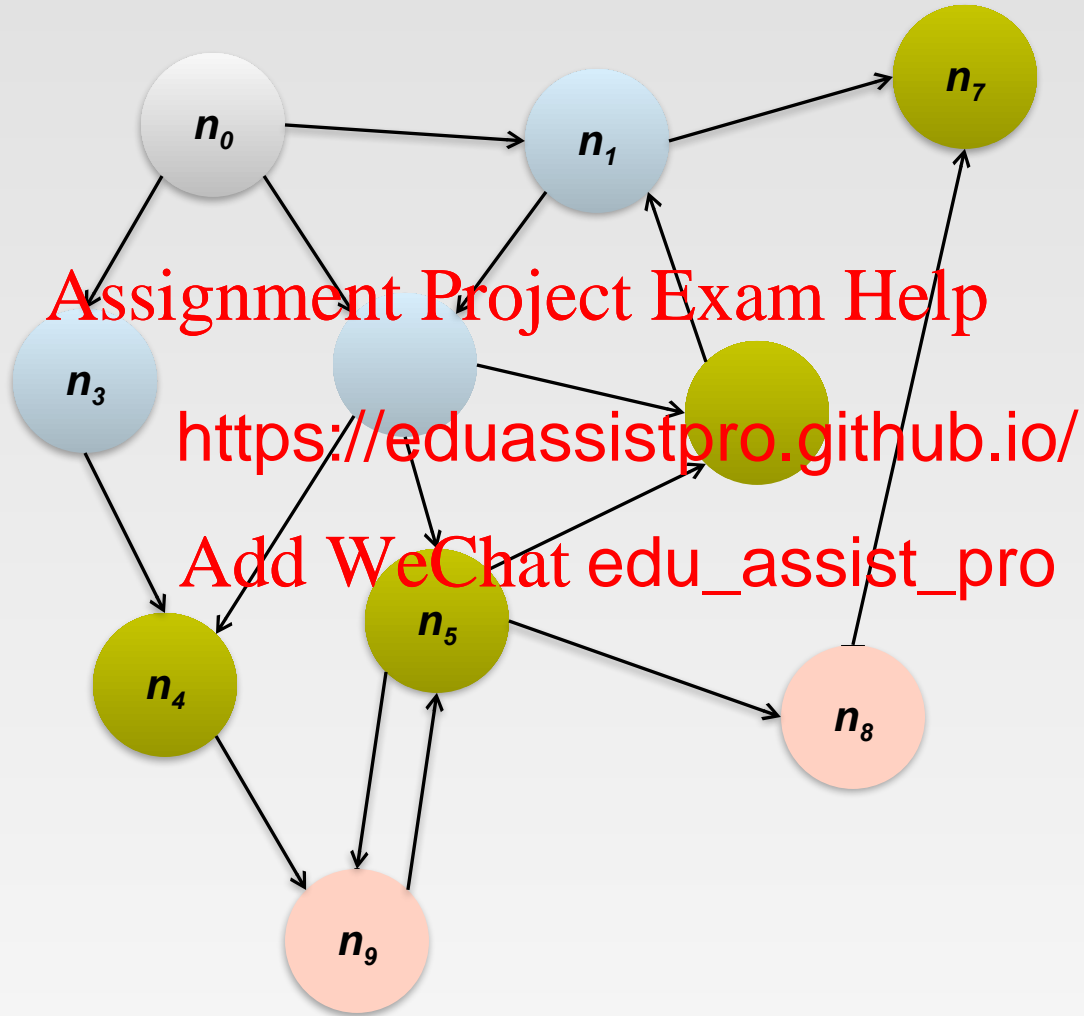
# Finding the Shortest Path

- Consider simple case of equal edge weights
- Solution to the problem can be defined inductively
- Here's the intuition:
  - Define:  $b$  is reachable from  $a$  if  $b$  is on adjacency list of  $a$
  - $\text{DISTANCE}(s) = 0$
  - For all node  $n$ 

$$\text{DISTANCE}(n) = 1 + \min_{m \in M} (\text{DISTANCE}(m))$$



# Visualizing Parallel BFS





# From Intuition to Algorithm

- Data representation:
  - Key: node  $n$
  - Value:  $d$  (distance from start), adjacency list (list of nodes reachable from  $n$ )
  - Initialization: for all nodes except for start node,  $d = \infty$
- Mapper:
  - $\forall m \in \text{adjac}$
- Sort/Shuffle
  - Groups distances by reachable
- Reducer:
  - Selects minimum distance path for each reachable node
  - Additional bookkeeping needed to keep track of actual path

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro

# Multiple Iterations Needed

- Each MapReduce iteration advances the “known frontier” by one hop
  - Subsequent iterations include more and more reachable nodes as frontier expands
  - The input of Mapper is the output of Reducer in the previous iteration
  - Multiple iterations are needed to explore entire graph

Assignment Project Exam Help

<https://eduassistpro.github.io/>

- Preserving graph structure:
  - Problem: Where did the adjacency
  - Solution: mapper emits  $(n, \text{adjacency list})$  as well

Add WeChat edu\_assist\_pro

# BFS Pseudo-Code

- Equal Edge Weights (how to deal with weighted edges?)
- Only distances, no paths stored (how to obtain paths?)

```
class Mapper
  method Map(nid n, node N)
    d ← N.Distance
    Emit(nid n, N) // Pass along graph structure
    for all nodeid m ∈ N.
      Emit(nid m, d+reachable nodes)
```

```
class Reducer
  method Reduce(nid m, [d1, d2, ...])
    dmin ← ∞
    M ← ∅
    for all d ∈ counts [d1, d2, ...] do
      if IsNode(d) then
        M ← d //Recover graph structure
      else if d < dmin then //Look for shorter distance
        dmin ← d
    M.Distance ← dmin //Update shortest distance
    Emit(nid m, node M)
```

# Stopping Criterion

- How many iterations are needed in parallel BFS (equal edge weight case)?
- Convince yourself: when a node is first “discovered”, we’ve found the shortest path
- Now answer the question.
  - The diameter of the graph, or the greatest distance between any pair of node
  - Six degrees
    - ▶ If this is indeed true, then a first search on the global social network would x MapReduce iterations.

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro

# Implementation in MapReduce

- The actual checking of the termination condition must occur outside of MapReduce.
- The driver (main) checks to see if a termination condition has been met, and if not, repeats.
- Hadoop provides a lightweight API called “counters”.
  - It can be used for counting events that occur during execution, e.g., number of times a certain condition is met.
  - Counters can be designed to count the number of nodes that have distances of  $\infty$  at the end of the program. The driver can access the final counter value and check to see if another iteration is necessary.

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro

# How to Find the Shortest Path?

- The parallel breadth-first search algorithm only finds the shortest distances.
- Store “back-pointers” at each node, as with Dijkstra's algorithm
  - Not efficient to recover the path from the back pointers
- A simpler approach is to store the shortest path distances in the graph at all times, so that each node will have an easily accessible path to the source.
  - The additional space requirement is  $O(V^2)$ .

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro

# Mapper

```
public static class TheMapper extends Mapper<LongWritable, Text, LongWritable, Text> {
    public void map(LongWritable key, Text value, Context context) throws IOException, InterruptedException {
        Text word = new Text();
        String line = value.toString();//looks like 1 0 2:3:
        String[] sp = line.split(" ");//splits on space
        int distanceadd = Integer.parseInt(sp[1]) + 1;
        String[] PointsTo = sp[2].split(":");
        for(int i=0; i<PointsTo.length; i++){
            word.set("VALUE "+distanceadd+PointsTo[i]);
            context.write(new LongWritable(1), word);
            word.clear(); }
        //pass in current node's distance (if it is the lowest distance)
        word.set("VALUE "+sp[1]);
        context.write( new LongWritable( Integer.parseInt( sp[1] ) ), word );
        word.clear();
        word.set("NODES "+sp[2]);//tells me to append on the final tally
        context.write( new LongWritable( Integer.parseInt( sp[0] ) ), word );
        word.clear();
    }
}
```

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro

# Reducer

```
public static class TheReducer extends Reducer<LongWritable, Text, LongWritable, Text> {
    public void reduce(LongWritable key, Iterable<Text> values, Context context) throws
IOException, InterruptedException {
    String nodes = "UNMODED";
    Text word = new Text();
    int lowest = INFINITY; //start at infinity

    for (Text val : values)
        //looks like NODES/VALUES 1 0 2:3:, we need to use the first as a key
        String[] sp = val.toString().split(" ");
        //look at first value
        if(sp[0].equals("NODES"))
            nodes = null;
            nodes = sp[1];
        }else if(sp[0].equals("VALUES")){
            int distance = Integer.parseInt(sp[2]);
            lowest = Math.min(distance, lowest);
        }
    }
    word.set(lowest+" "+nodes);
    context.write(key, word);
    word.clear();
}
```

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro

<https://github.com/himank/Graph-Algorithm-Map-Reduce/blob/master/src/DijkstraAlgo.java>



# BFS Pseudo-Code (Weighted Edges)

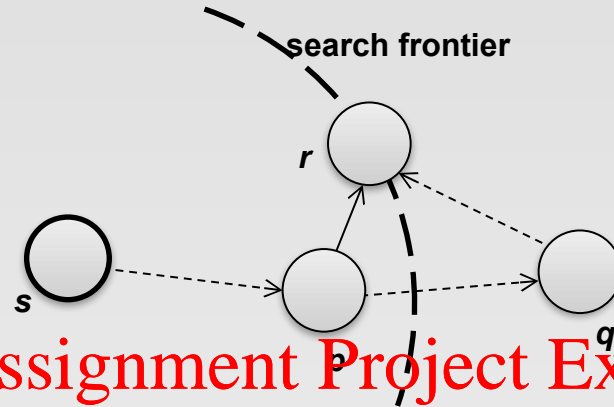
- The adjacency lists, which were previously lists of node ids, must now encode the edge distances as well
  - Positive weights!
- In line 6 of the mapper code, instead of emitting  $d+1$  as the value, we must now emit  $d + w$ , where  $w$  is the edge distance
- **The termination behaviour is very different**
  - How many iterations are needed? ( $\infty$  for positive edge weight case)
  - Convince yourself: when a node is first “discovered”, we’ve found the shortest path

<https://eduassistpro.github.io/>

Add WeChat: edu\_assist\_pro

Not true!

# Additional Complexities

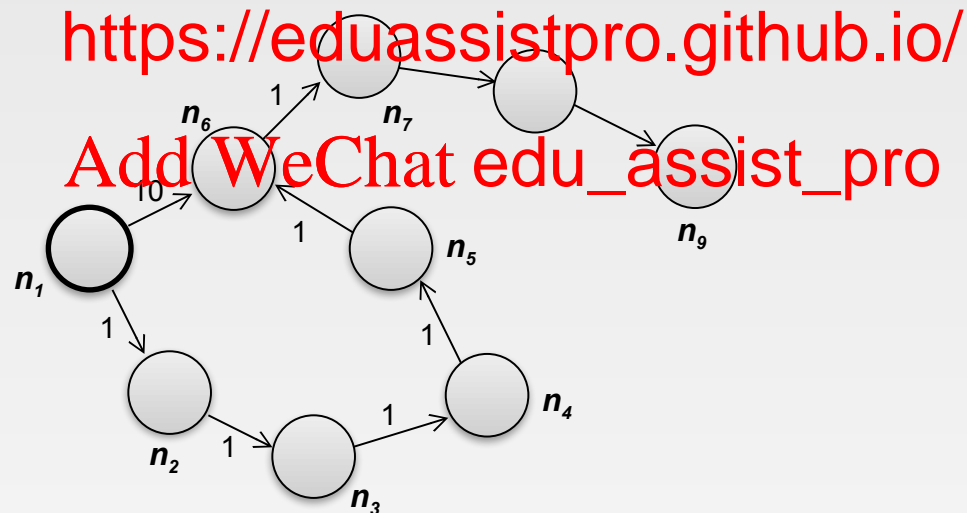


## Assignment Project Exam Help

- Assume that  $p$  is the very first node for the very first time.
  - In the current time.
  - We've already discovered the shortest distance to  $r$  through  $p$ .
  - Is  $s \rightarrow p \rightarrow r$  the shortest path from  $s$  to  $r$ ?
- The shortest path from source  $s$  to node  $r$  may go outside the current search frontier
  - It is possible that  $p \rightarrow q \rightarrow r$  is shorter than  $p \rightarrow r$ !
  - We will not find the shortest distance to  $r$  until the search frontier expands to cover  $q$ .

# How Many Iterations Are Needed?

- In the worst case, we might need as many iterations as there are nodes in the graph minus one
  - A sample graph that elicits worst-case behaviour for parallel breadth-first search.
  - Eight iterations are required to discover shortest distances to all nodes from  $n_1$ .



# Example (only distances)

□ Input file:

s --> 0 | n1: 10, n2: 5

n1 -->  $\infty$  | n2: 2, n3:1

n2 -->  $\infty$  | n1: 3, n3:9 , n4:2

n3 -->  $\infty$  | n4:4

n4 -->  $\infty$  | s:7, n3:6

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro

# Iteration 1

□ Map:

Read s --> 0 | n1: 10, n2: 5

Emit: (n1, 10), (n2, 5), and the adjacency list (s, n1: 10, n2: 5)

*The other lists will also be read and emit, but they do not contribute, and thus ignored* Assignment Project Exam Help

□ Reduce:

Receives: (n1, 10), <https://eduassistpro.github.io/>

*The adjacency list of each node will also be ignored in example*

Emit: Add WeChat edu\_assist\_pro

s --> 0 | n1: 10, n2: 5

n1 --> 10 | n2: 2, n3:1

n2 --> 5 | n1: 3, n3:9 , n4:2

# Iteration 2

□ Map:

Read: n1 --> 10 | n2: 2, n3:1

Emit: (n2, 12), (n3, 11), (n1, <10, (n2: 2, n3:1)>)

Read: n2 --> 5 | n1: 3, n3:9 , n4:2

Emit: (n1, 8), (n3, 14), (n4, 7), (n2, <5, (n1:3, n3:9, n4:2)>)

*Ignore the processi*

□ Reduce:

Receives: (n1, (8, <10, (n2: 2, n3:1)>)), (n3, (11, 14)), (n4, 7), (n2, <5, (n1: 3, n3:9, n4:2)>))

Emit:

n1 --> 8 | n2: 2, n3:1

n2 --> 5 | n1: 3, n3:9 , n4:2

n3 --> 11 | n4:4

n4 --> 7 | s:7, n3:6

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro

# Iteration 3

□ Map:

Read: n1 --> 8 | n2: 2, n3:1

Emit: (n2, 10), (n3, 9), (n1, <8, (n2: 2, n3:1)>)

Read: n2 --> 5 | n1: 3, n3:9 , n4:2 (**Again!**)

Emit: (n1, 8), (n3, 14), (n4, 7), (n2, <5, (n1:3, n3:9, n4:2)>)

Read: n3 --> 11 | n

Emit: (n4, 15) , (n <https://eduassistpro.github.io/>

Read: n4 --> 7 | s:7, n3:6

Emit: (s, 14), (n3, 13), (n4, <7, (s:7, n3: [Add WeChat edu\\_assist\\_pro](#)

□ Reduce:

Emit:

n1 --> 8 | n2: 2, n3:1

n2 --> 5 | n1: 3, n3:9 , n4:2

n3 --> 9 | n4:4

n4 --> 7 | s:7, n3:6

# Iteration 4

□ Map:

Read: n1 --> 8 | n2: 2, n3:1 (**Again!**)

Emit: (n2, 10), (n3, 9), (n1, <8, (n2: 2, n3:1)>)

Read: n2 --> 5 | n1: 3, n3:9 , n4:2 (**Again!**)

Emit: (n1, 8), (n3, 14), (n4, 7), (n2, <5, (n1:3, n3:9, n4:2)>)

Read: n3 --> 9 | n4:

Emit: (n4, 13) , (n <https://eduassistpro.github.io/>

Read: n4 --> 7 | s:7, n3:6 (**Again!**)

Emit: (s, 14), (n3, 13), (n4, <7, (s:7, n3:

□ Reduce:

Emit:

n1 --> 8 | n2: 2, n3:1

n2 --> 5 | n1: 3, n3:9 , n4:2

n3 --> 9 | n4:4

n4 --> 7 | s:7, n3:6

**In order to avoid duplicated computations, you can use a status value to indicate whether the distance of the node has been modified in the previous iteration.**

**No updates. Terminate.**



# Comparison to Dijkstra

- Dijkstra's algorithm is more efficient
  - At any step it only pursues edges from the minimum-cost path inside the frontier
- MapReduce explores all paths in parallel
  - Lots of "waste"
  - Useful work
- Why can't we do <https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro

# Graphs and MapReduce

- Graph algorithms typically involve:
  - Performing computations at each node: based on node features, edge features, and local link structure
  - Propagating computations: “traversing” the graph
- Generic recipe:
  - Represent  $g$
  - Perform local computation at each node
  - Pass along partial results via outgoing edges to destination node
  - Perform aggregation in reducer node
  - Iterate until convergence: controlled by external “driver”
  - Don’t forget to pass the graph structure between iterations

Assignment Project Exam Help

<https://eduassistpro.github.io/>

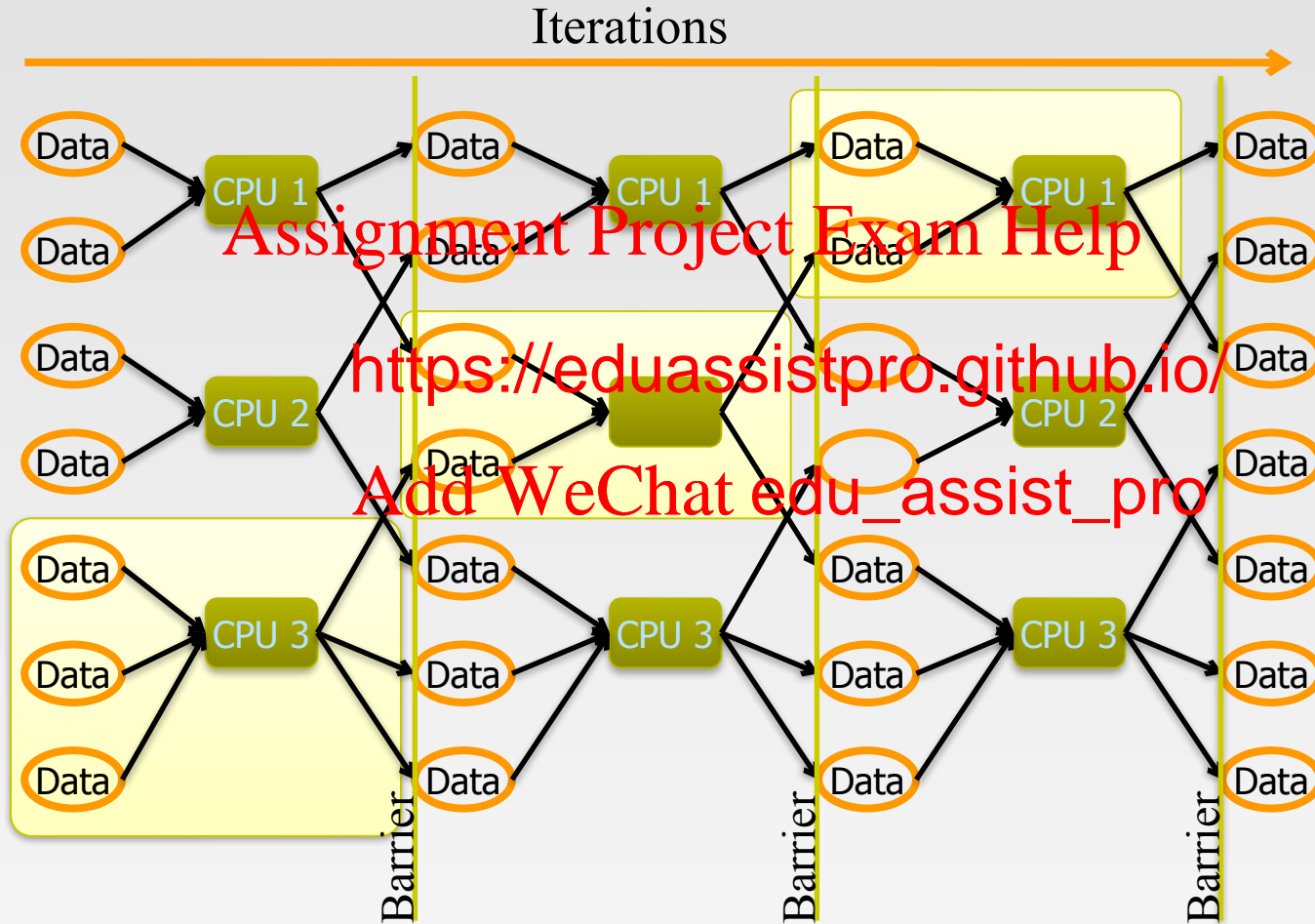
Add WeChat edu\_assist\_pro

# Issues with MapReduce on Graph Processing

- MapReduce Does not support iterative graph computations:
  - External driver. Huge I/O incurs
  - No mechanism to support global data structures that can be accessed and updated by all mappers and reducers
    - ▶ Passing information is only possible within the local graph structure – through adjacency list
    - ▶ Dijkstra's algorithm in Hadoop : a global priority queue that guid
    - ▶ Dijkstra's algorithm in Hadoop : a global priority queue available. Do some “wasted” computation
- MapReduce algorithms are often impractical on large, dense graphs.
  - The amount of intermediate data generated is on the order of the number of edges.
  - For dense graphs, MapReduce running time would be dominated by copying intermediate data across the network.

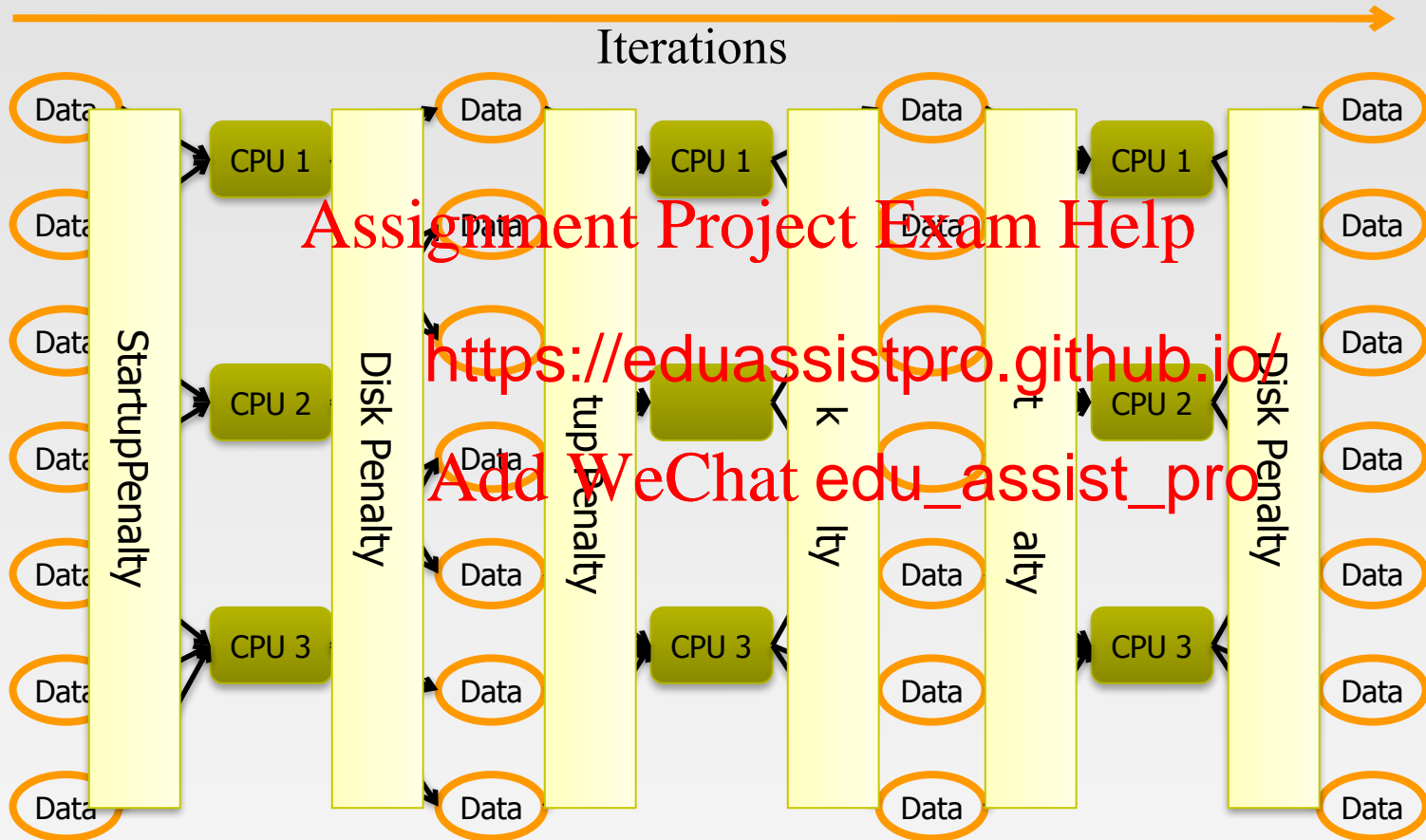
# Iterative MapReduce

- Only a subset of data needs computation:



# Iterative MapReduce

- System is not optimized for iteration:



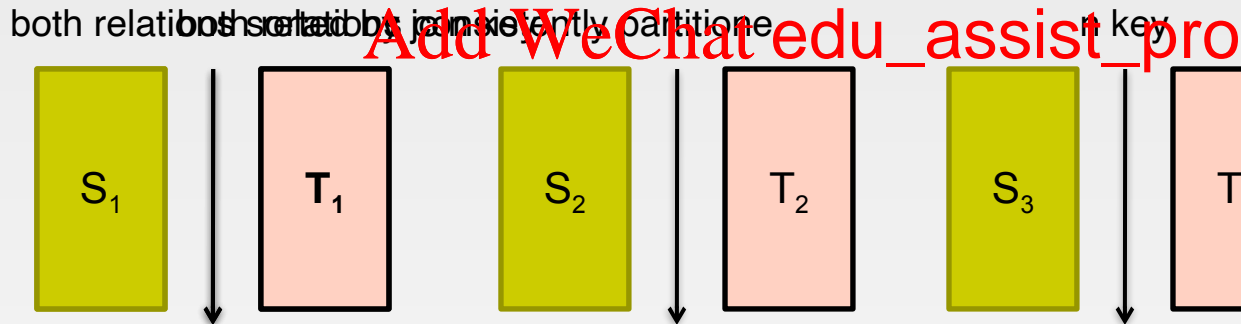
# Better Partitioning

- Default: hash partitioning
  - Randomly assign nodes to partitions
- Observation: many graphs exhibit local structure
  - E.g., communities in social networks
  - Better partitioning creates more opportunities for local aggregation
- Unfortunately, p
  - Sometimes, <https://eduassistpro.github.io/>
  - But cheap heuristics sometimes
  - For webgraphs: range partition ted URLs

# Schimmy Design Pattern

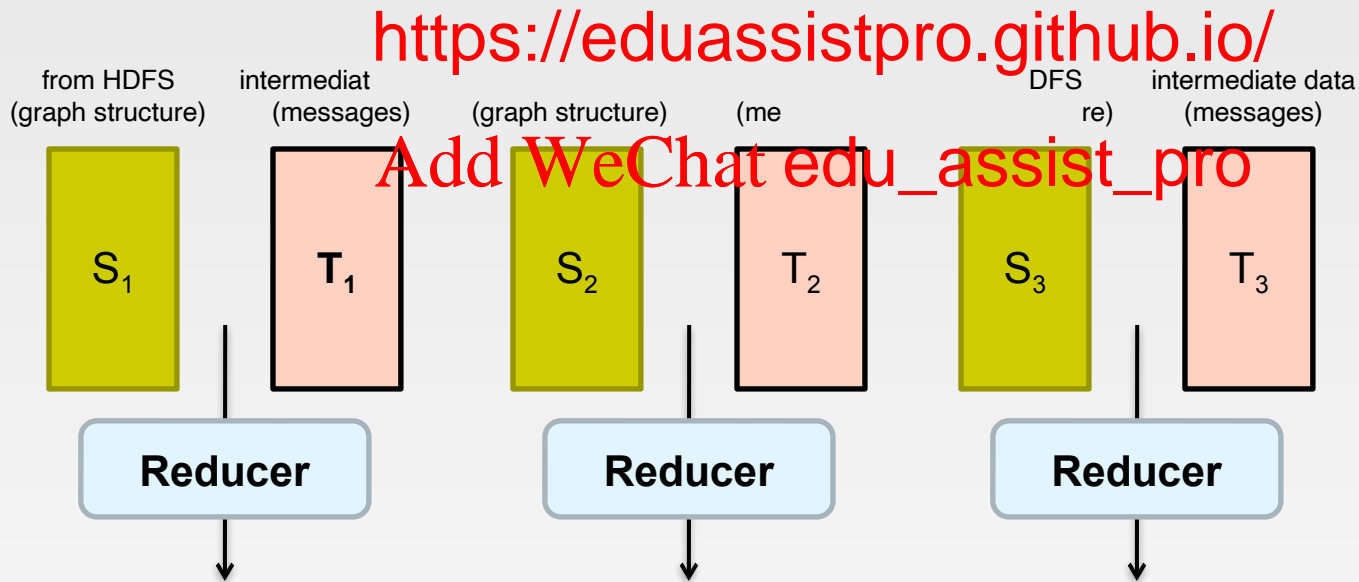
- Basic implementation contains two dataflows:
  - Messages (actual computations)
  - Graph structure (“bookkeeping”)
- Schimmy: separate the two dataflows, shuffle only the messages
  - Basic idea: merge join between graph structure and messages

<https://eduassistpro.github.io/>



# Do the Schimmy!

- Schimmy = reduce side parallel merge join between graph structure and messages
  - Consistent partitioning between input and intermediate data
  - Mappers emit only messages (actual computation)
  - Reducers read graph structure directly from HDFS





Assignment Project Exam Help

**In** <https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro

# Motivation of Pregel

- Many practical computing problems concern large graphs

## Large graph data

Web graph  
Transportation routes  
Citation relationships  
Social networks

## Graph algorithms

PageRank  
Shortest path  
Connected components  
Mining techniques



Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro

- Single computer graph library does
- MapReduce is ill-suited for graph processing
  - Many iterations are needed for parallel graph processing
  - Materializations of intermediate results at every MapReduce iteration harm performance

# Pregel

- **Pregel**: A System for Large-Scale **Graph** Processing (Google) - Malewicz et al. SIGMOD 2010.

- Scalable and Fault-tolerant platform

## Assignment Project Exam Help

- API with flexibility

<https://eduassistpro.github.io/>

- Inspired by Valiant's Bulk Synchronous Parallel Model
  - Leslie G. Valiant: A Bridging Model for Parallel Computation. Commun. ACM 33 (8): 103-111 (1990)

- Vertex centric computation (Think like a vertex)

# Bulk Synchronous Parallel Model (BSP)

analogous to MapReduce rounds

- Processing: a series of **supersteps**
- Vertex**: computation is defined to run on each vertex
- Superstep S**: *all vertices compute in parallel; each vertex v may*
  - receive **messages** sent to v from superstep S – 1;
  - perform some computation: modify its states and the states of its outgoing ed
  - Send **messages** (to be received in the next superstep)

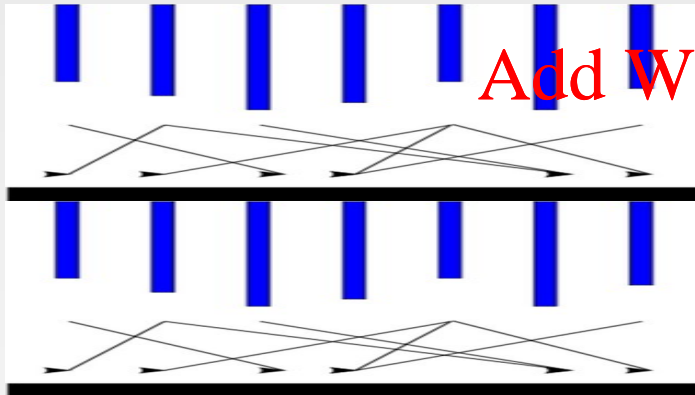
<https://eduassistpro.github.io/>  
Add WeChat: edu\_assist\_pro

*Vertex-centric, message passing*

# Pregel Computation Model

- Based on Bulk Synchronous Parallel (BSP)
  - Computational units encoded in a directed graph
  - Computation proceeds in a series of supersteps
  - Message passing architecture

Input



Output

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro

Supersteps  
(a sequence of iterations)

# Pregel Computation Model (Cont')

- Concurrent computation and Communication need not be ordered in time
- Communication through message passing

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro

# Pregel Computation Model (Cont')

□ Superstep: the vertices compute in parallel

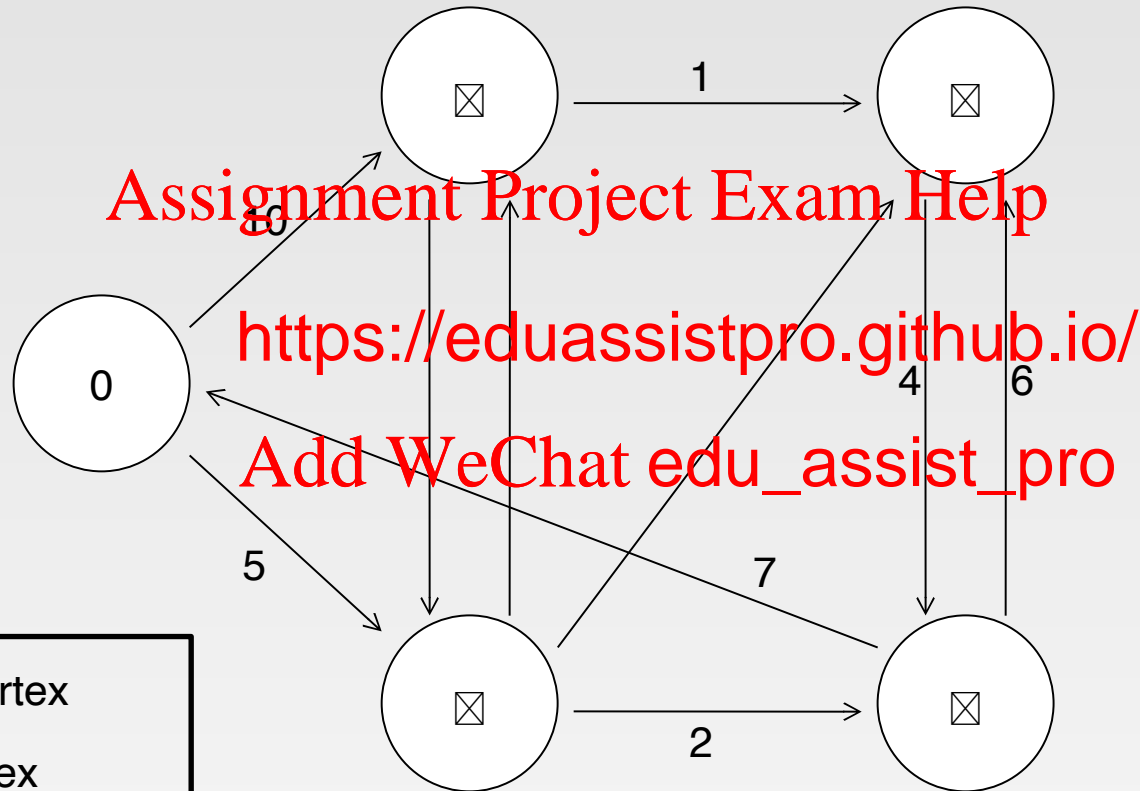
□ Each vertex

- ▶ Receives messages sent in the previous superstep
- ▶ Executes the same user-defined function
- ▶ Modifies edges
- ▶ Sends messages to other vertices (received in the next superstep)
- ▶ Votes to halt if it has no further work to do

□ Termination condition

- ▶ All vertices are simultaneously inactive
- ▶ A vertex can choose to deactivate itself
- ▶ Is “woken up” if new messages received

# Example: SSSP – Parallel BFS in Pregel



● Inactive Vertex

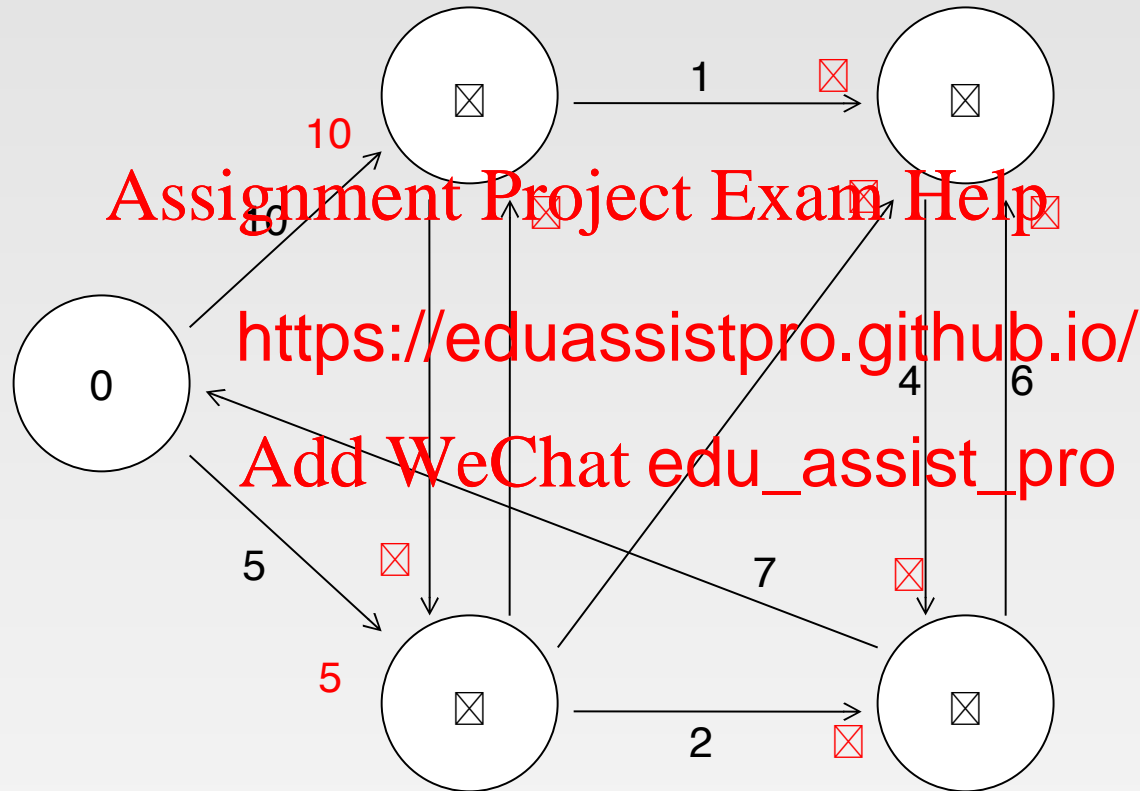
○ Active Vertex

$\xrightarrow{x}$  Edge weight

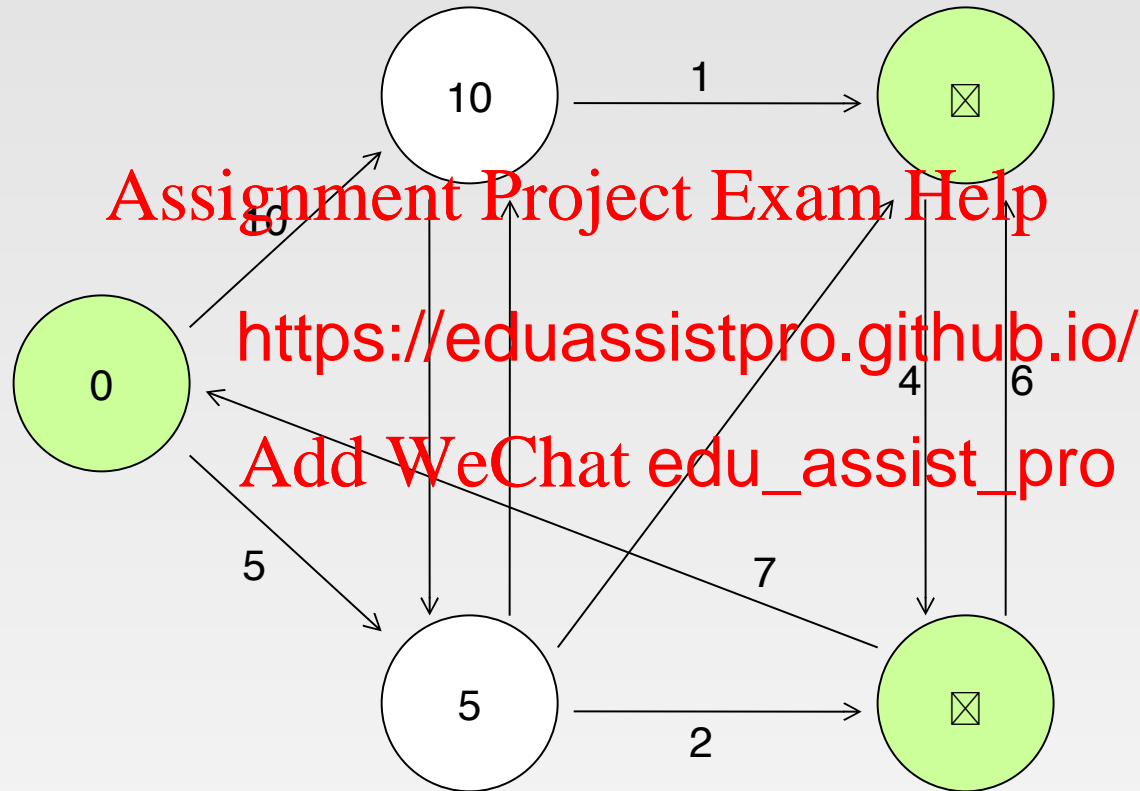
$\times$  Message



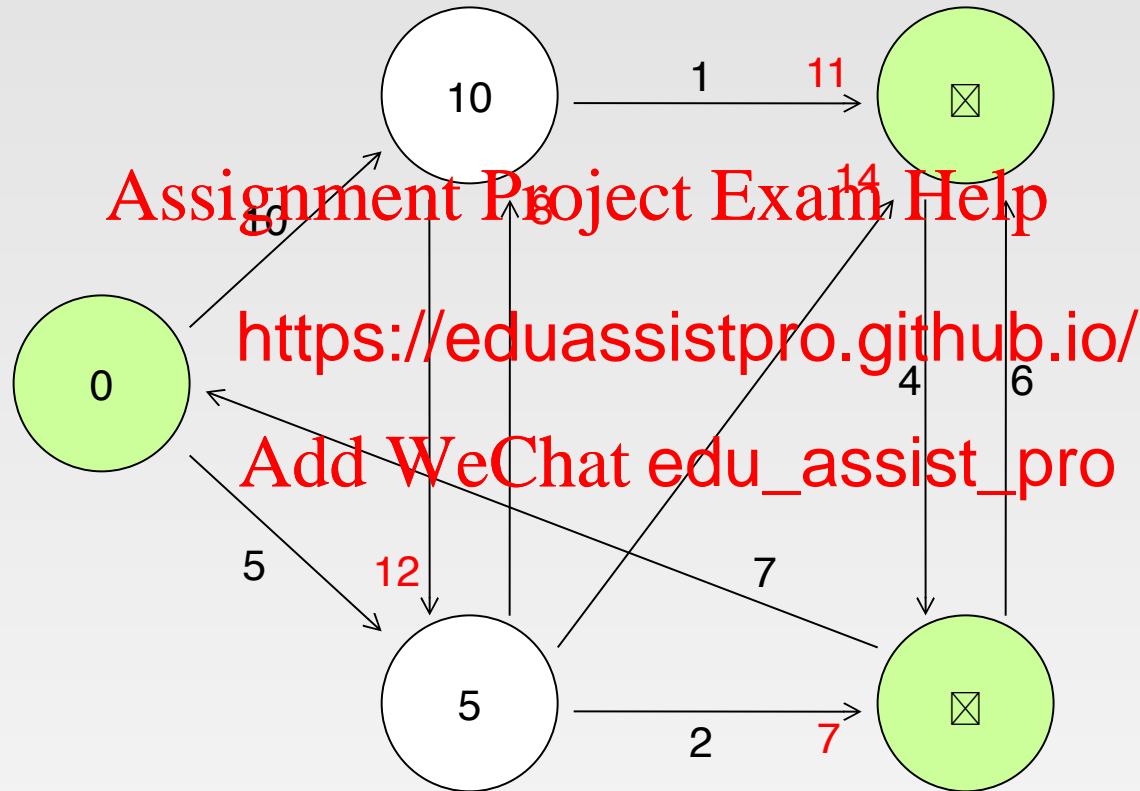
# Example: SSSP – Parallel BFS in Pregel



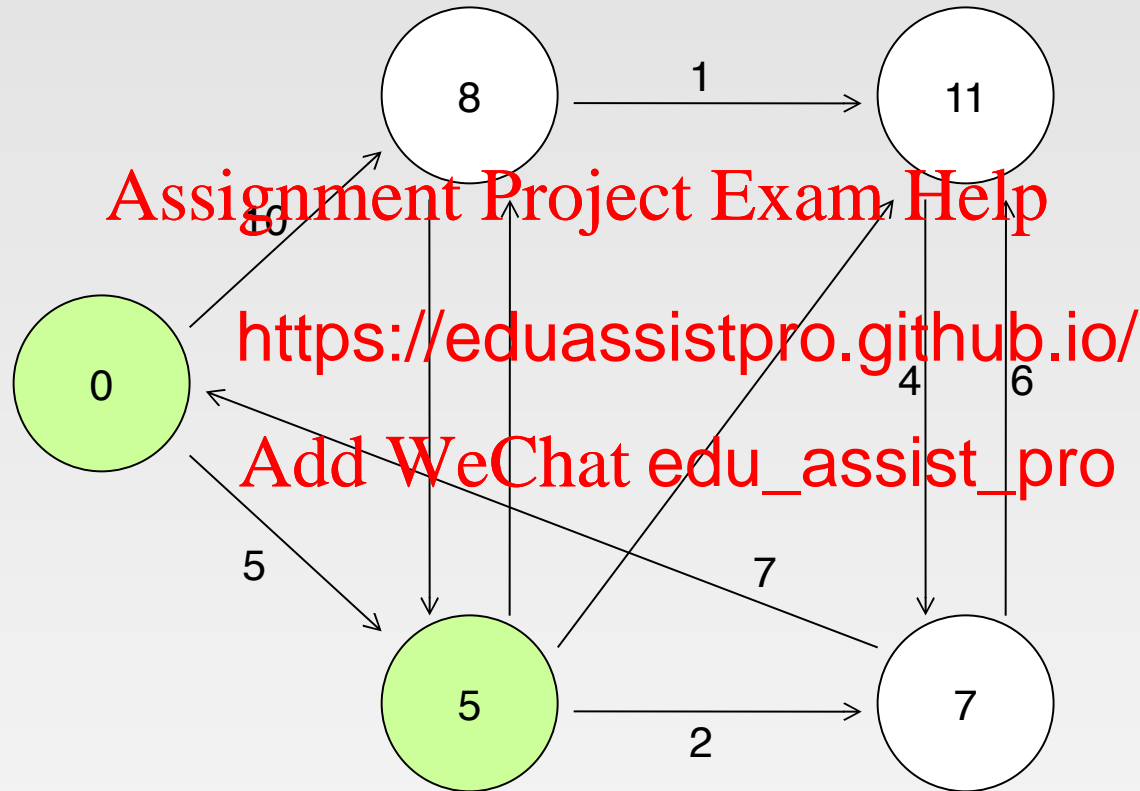
# Example: SSSP – Parallel BFS in Pregel



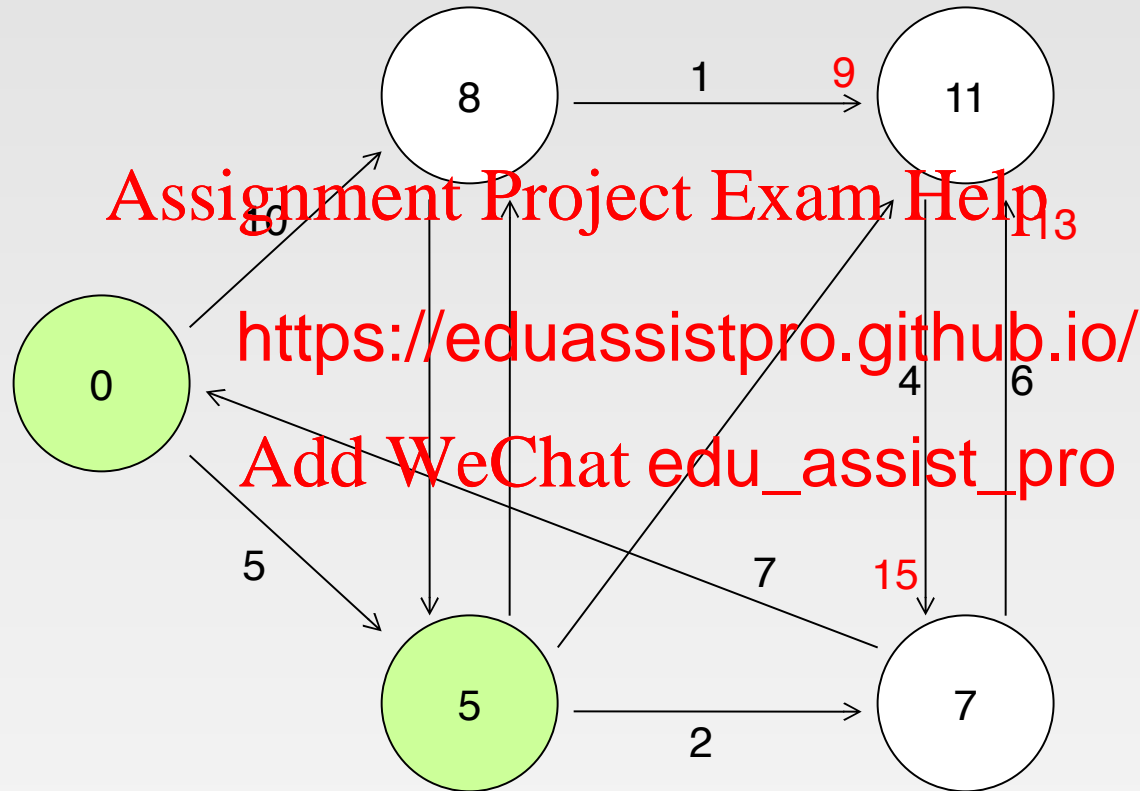
# Example: SSSP – Parallel BFS in Pregel



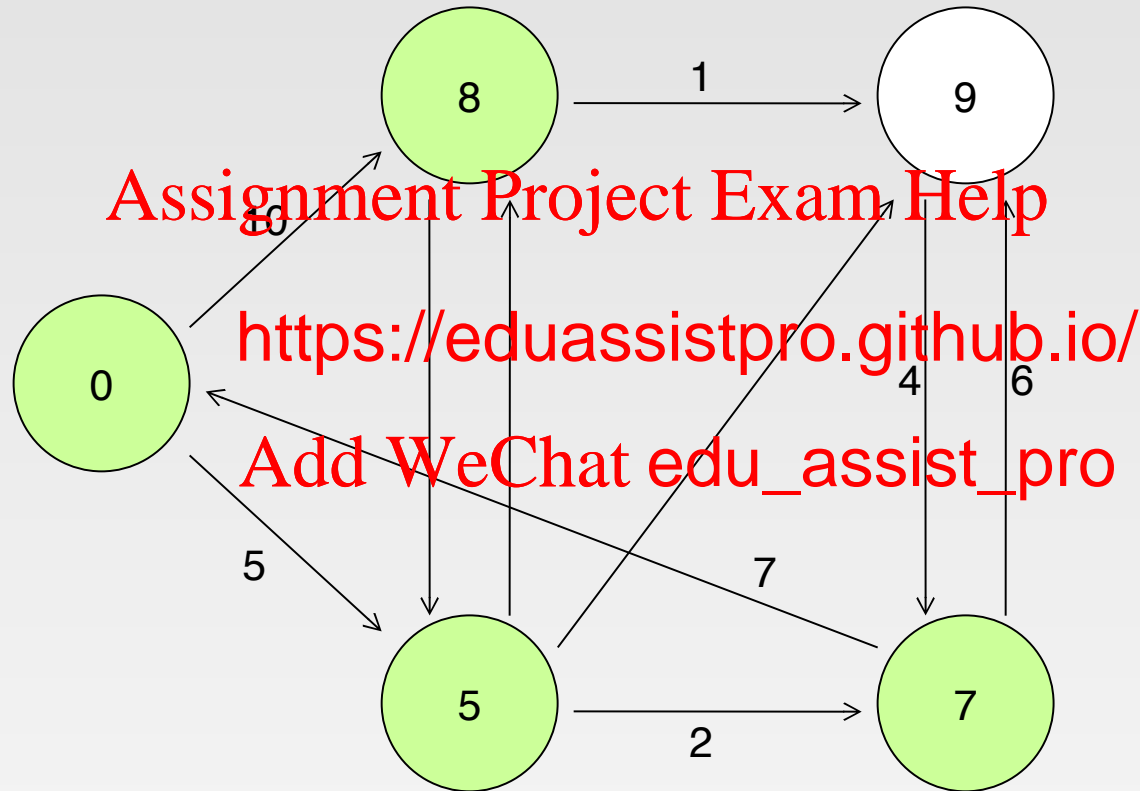
# Example: SSSP – Parallel BFS in Pregel



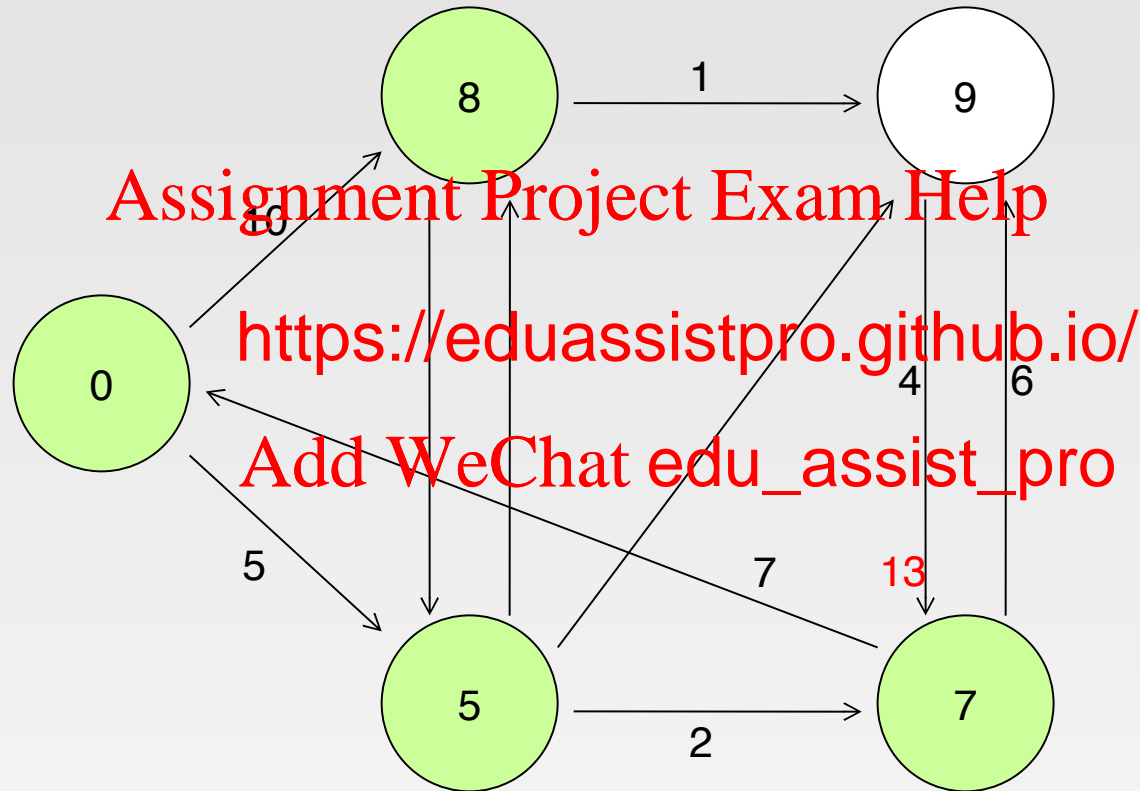
# Example: SSSP – Parallel BFS in Pregel



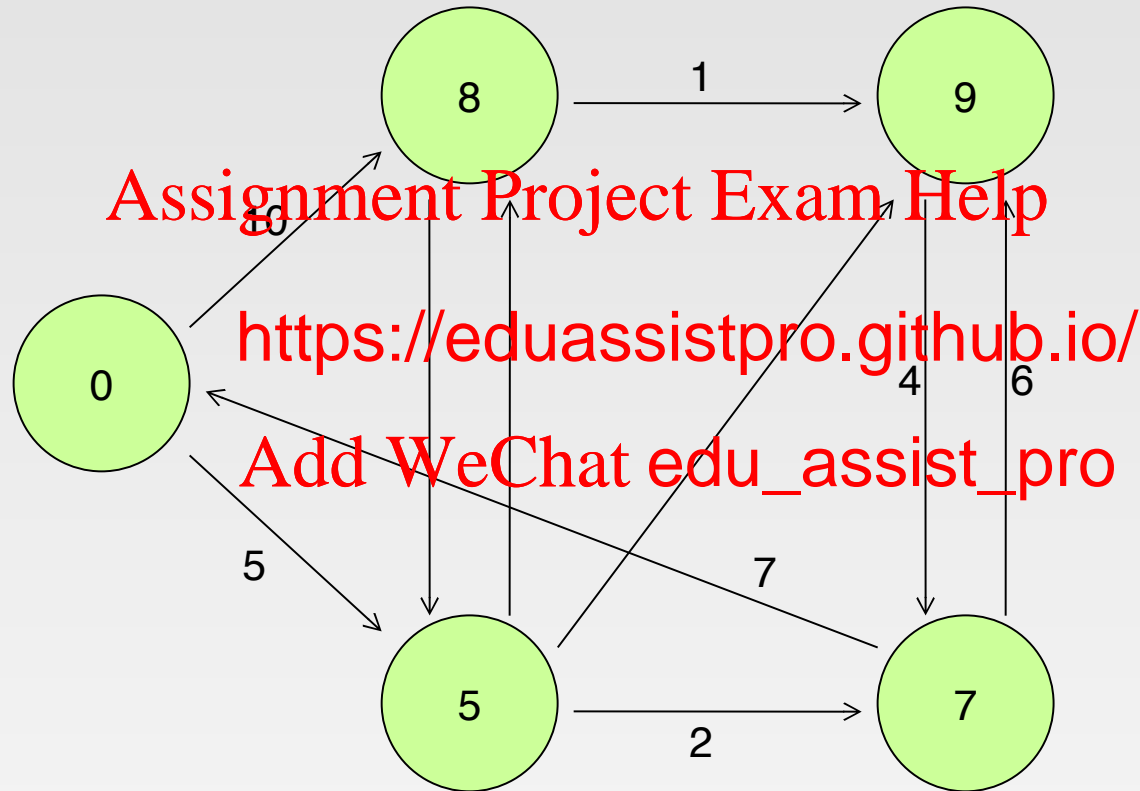
# Example: SSSP – Parallel BFS in Pregel



# Example: SSSP – Parallel BFS in Pregel



# Example: SSSP – Parallel BFS in Pregel





# Differences from MapReduce

- Graph algorithms can be written as a series of chained MapReduce jobs
- Pregel
  - Keeps vertex computation local, while MapReduce performs computation on the entire graph
  - Uses network transfers only for vertex-to-vertex communication
- MapReduce
  - Passes the entire state of the graph from one stage to the next
  - Needs to coordinate the steps of a chained MapReduce

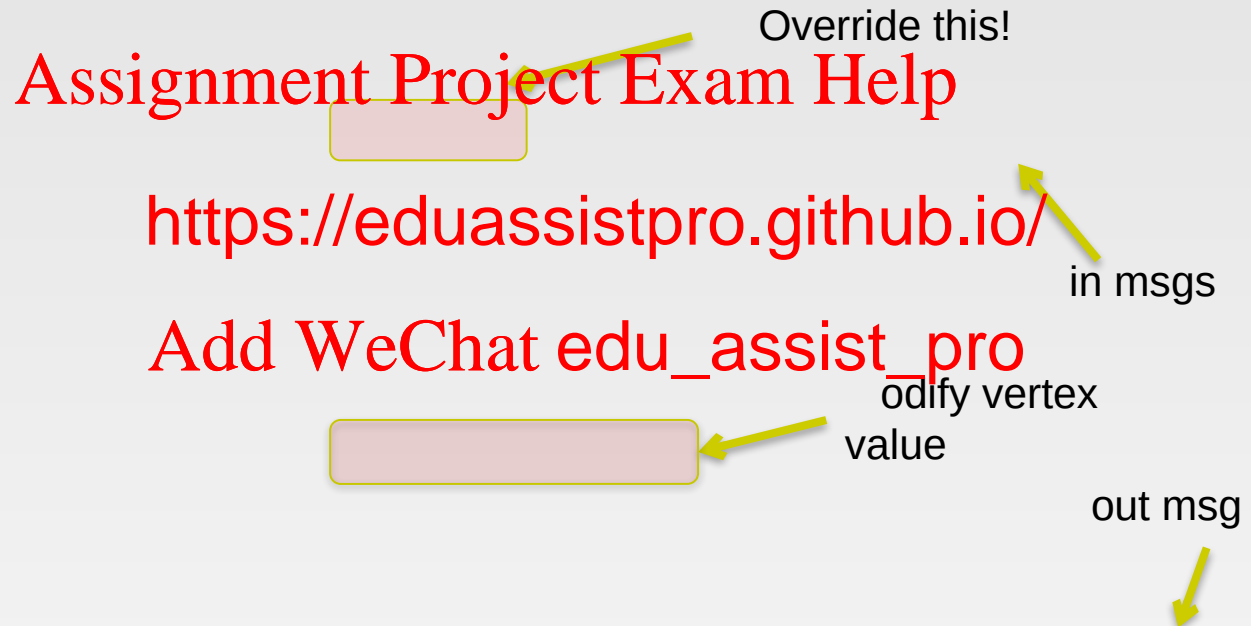
Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat: edu\_assist\_pro

# Writing a Pregel Program (C++)

- Subclassing the predefined **Vertex** class



# Pregel: SSSP (C++)

Refer to the current node as  $u$

```
class ShortestPathVertex : public Vertex<int, int, int>
{
    void Compute(MessageIterator* msgs) {
        int mindist = IsSource(vertex_id()) ? 0 : INF;
        for (; !msgs.Done(); msgs.Next()) {
            int e = msgs.GetValue();
            if (mindist > e) {
                *MutableValue() = e;
                OutEdgeIterator iter = GetOutEdgeIterator();
                for (; !iter.Done(); iter.Next()) {
                    SendMsgTo(iter.Target(), mindist + iter.GetWeight());
                }
            }
        }
        VoteToHalt();
    }
};
```

## aggregation

Messages: distances to u

Value of the current  
e

# Add WeChat edu\_assist\_pro

Pass revised distance to its neighbors

# More Tools on Big Graph Processing

## □ Graph databases: Storage and Basic Operators

- [http://en.wikipedia.org/wiki/Graph\\_database](http://en.wikipedia.org/wiki/Graph_database)
- Neo4j (an open source graph database)
- InfiniteGraph
- VertexDB
- ... ..

Assignment Project Exam Help

<https://eduassistpro.github.io/>

## □ Distributed Graph Processing (most nly)

- Google's Pregel (vertex centere)
- Giraph (Apache)
- GraphX (Spark)
- GraphLab
- ... ..

Add WeChat edu\_assist\_pro

# References

- Chapter 5. Data-Intensive Text Processing with MapReduce
- Chapter 5. Mining of Massive Datasets.

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro

Assignment Project Exam Help

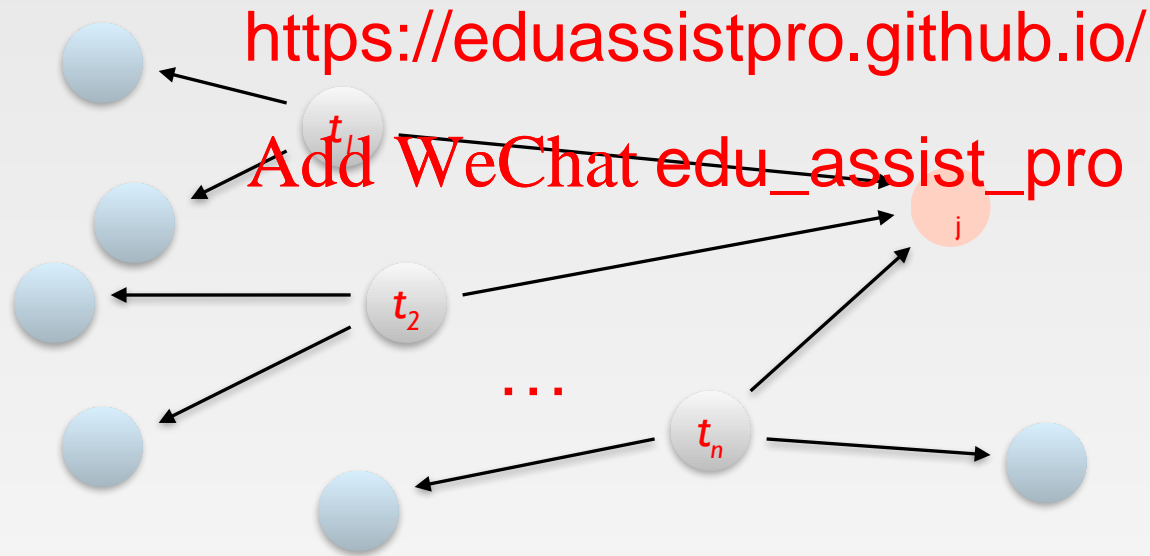
<https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro

# PageRank Review

- Given page  $t_j$  with in-coming neighbors  $t_1 \dots t_n$ , where
  - $d_i$  is the out-degree of  $t_i$
  - $\beta$  is the teleport probability
  - $N$  is the total number of nodes in the graph

Assignment Project Exam Help



# Computing PageRank

- Properties of PageRank

- Can be computed iteratively
- Effects at each iteration are local

- Sketch of algorithm:

- Start with seed  $r$  values
- Each page  $i$  it links to
- Each target  $j$   $\frac{r_i}{\text{in-degree}(j)}$  multiple in-bound links to
- compute  $r_j$
- Iterate until values converge

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro



# Simplified PageRank

- First, tackle the simple case:
  - No teleport
  - No dangling nodes (dead ends)
- Then, factor in these complexities...
  - How to deal with the teleport probability?
  - How to deal

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro

# Sample PageRank Iteration (1)

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro

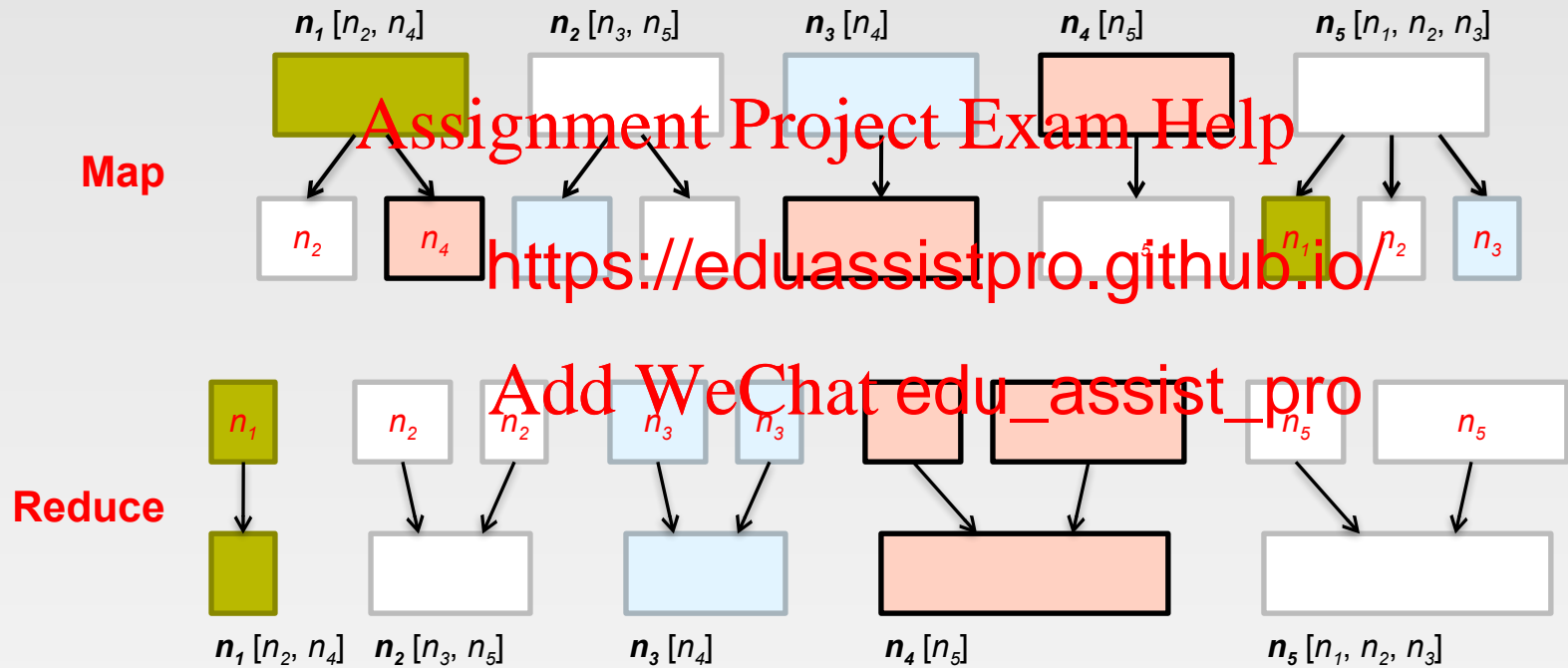
# Sample PageRank Iteration (2)

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro

# PageRank in MapReduce (One Iteration)



# PageRank Pseudo-Code

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro

# PageRank vs. BFS

	PageRank	BFS
Map	$\sum_i \frac{w_{ji}}{d_j}$	$d + w$
Reduce	sum	min

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro

# PageRank in Pregel

□ **Superstep 0:** Value of each vertex is `1/NumVertices()`

```
virtual void Compute(MessageIterator* msgs) {  
    if (superstep() >= 1) {  
        double sum = 0;  
        for (; !msgs->done(); msgs->Next())  
            sum += msgs->Value();  
        Value() = 1.0 / n * sum;  
    }  
    if (supersteps() < 30) {  
        const int64 n = G.rator().size();  
        SendMessageToAllNeighbors(GetValue() / n);  
    } else {  
        VoteToHalt();  
    }  
}
```

Assignment Project Exam Help  
<https://eduassistpro.github.io/>  
Add WeChat edu\_assist\_pro