# Application 1: Analysis of Citation Graphs
# Algorithmic Thinking (Part1)

Shamsuddin Rehmani

July 17, 2019

## Application 1 Description

In the Module 1 Application, we will combine the mathematical analysis that we began in the Homework with the code that you have written in the Project to analyze a real-world problem: How do scientific papers get cited? This part of the module will probably be much more unstructured than you are accustomed to in an on-line class. Our goal is to provide a more realistic simulation of how the concepts that you are learning are actually used in practice. Your key task in this part of the module is to think about the problem at hand as you answer each question.

## Citation graphs

Our task for this application is to analyze the structure of graphs generated by citation patterns from scientific papers. Each scientific paper cites many other papers, say 20-40, and sometimes (e.g., review papers) hundreds of other papers. But, let's face it: It is often the case that the authors of a paper are superficially familiar with some (many?) of the papers they cite. So, the question is: Are the cited papers chosen randomly (from within the domain of the paper) or is there some "hidden pattern"?

Given that we will be looking at "paper $i$ cites paper $j$" relationships, it makes sense to represent the citation data as a directed graph (a citation graph) in which the nodes correspond to papers, and there is an edge from node $i$ to node $j$ if the paper corresponding to node $i$ cites the paper corresponding to node $j$. Since we're interested in understanding how papers get cited, we will analyze the in-degree distribution of a specific graph, and contrast it to those of graphs generated by two different random processes
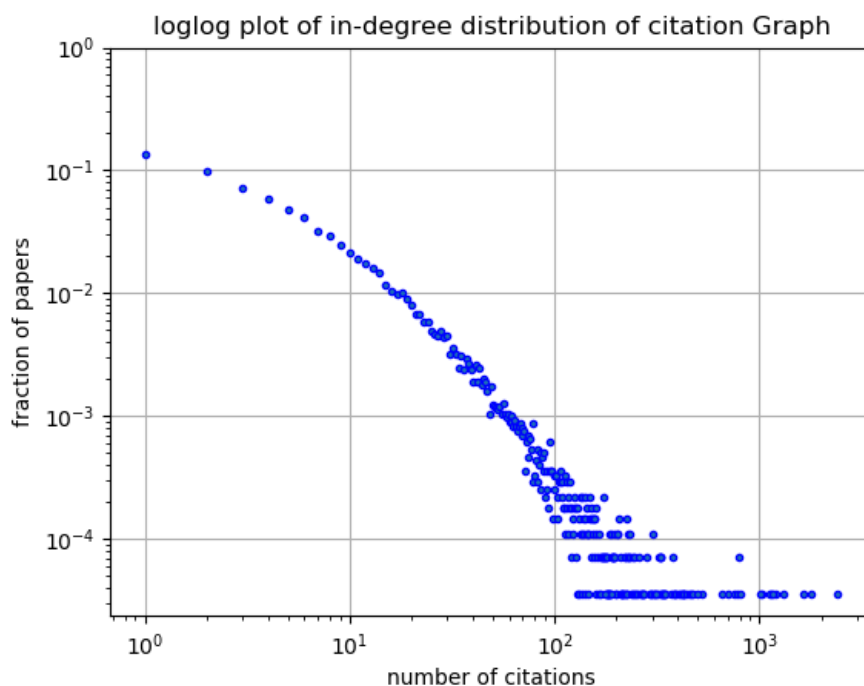
## Answer to Question 1



Figure 1: in-dgree distribution of citation graph

# Answer to Question 2

**Q2.1: Is the expected value of the in-degree the same for every node in an ER graph? Please answer yes or no and include a short explanation for your answer.**

Ans: yes it is the same for all nodes since the presence of an edge is independent of all other edges i.e it is independent of the current structure of the graph The expected value of in-degree is given by $p(n-1)$

**Q2.2: What does the in-degree distribution for an ER graph look like? Provide a short written description of the shape of the distribution.**

Ans: we know that the probability that a given node has degree $k$ is given by a binomial distribution as seen in the homework. Thus as $p \to 0$ (probability $p$ becomes smaller), we see more nodes with smaller in-degree and thus the in-degree distribution shape looks like a bell curve skewed towards the left i.e near in-degree 0. As $p \to 1$, we get more nodes with higher in-degree and the shape is increasing curve with most points near the higher in-degree region. For large number of nodes, and small p, this becomes a symmetric bell shaped curve and approaches a normal distribution

**Q2.3: Does the shape of the in-degree distribution plot for ER look similar to the shape of the in-degree distribution for the citation graph? Provide a short explanation of the similarities or differences. Focus on comparing the shape of the two plots as discussed in the class page on "Creating, formatting, and comparing plots".**

Ans: As mentioned in answer of Q2.2, the shape for the ER in-degree approaches a bell-shaped curve for large $N$ values and small $p$. However, for the citation graph it is a decreasing curve with majority of point located near in-degree of zero.

# Answer to Question 3

Value of $n$ and $m$ that yield a DPA graph whose number of nodes and edges is roughly the same to those of the citation graph are as follows:

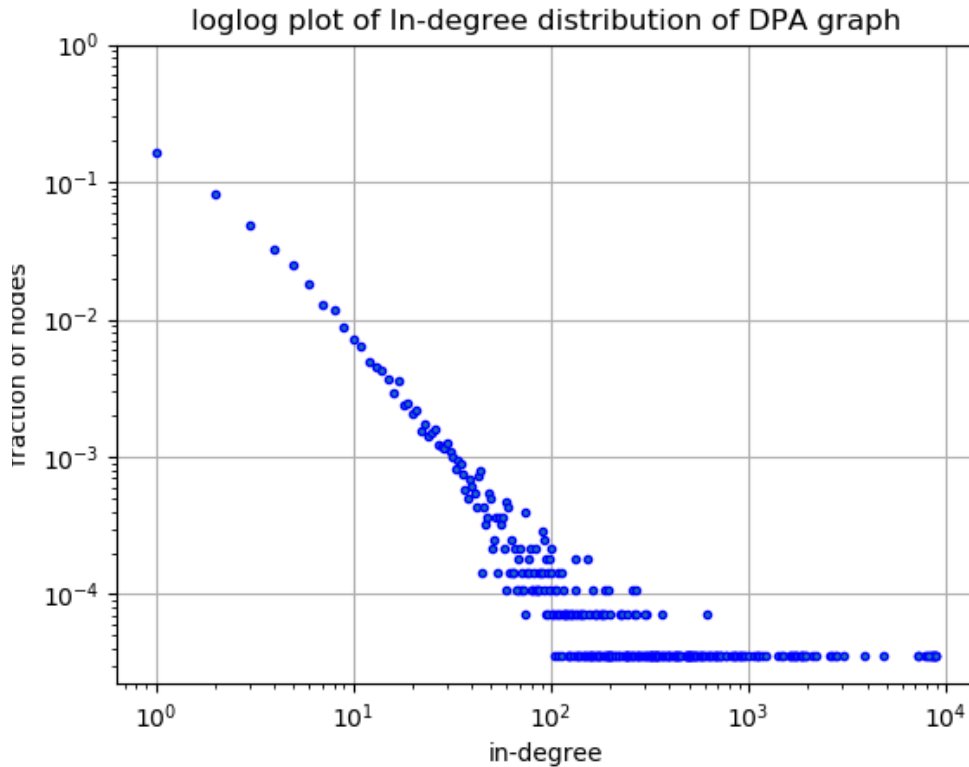- $m = 27770$
- $m = 13$

# Answer to Question 4



Figure 2: in-degree distribution for DPA graphs

# Answer to Question 5

**Q5.1: Is the plot of the in-degree distribution for the DPA graph similar to that of the citation graph? Provide a short explanation of the similarities or differences. Focus on the various properties of the two plots as discussed in the class page on "Creating, formatting, and comparing plots".**

Ans: Yes they are similar since both follow a linear log-log decreasing trend i.e they both follow the power law distribution and the point are spread out more as the in-degree increases

**Q5.2: Which one of the three social phenomena listed above mimics the behavior of the DPA process? Provide a short explanation for your answer.**

Ans: DPA process mimics the "rich get richer" or the "preferential attachment" phenomena since every new node that is added to the graph is most likely to be connected to the neighbor with highest in-degree.

**Q5.3: Could one of these phenomena explain the structure of the physics citation graph? Provide a short explanation for your answer.**

Ans: The citation graph also mimics the "rich get richer" phenomena as the paper with higher citations i.e higher degree tend to be more likely used in other papers as well due to being more visible

# A    Python code used to answer the Application Questions

```python
 1  """
 2  Analyze the structure of graphs generated by citation patterns from scientific papers
 3  """
 4
 5  import matplotlib.pyplot as plt
 6  import numpy as np
 7  import parse_graph
 8  import alg_project1_solution as alg_proj1_sol
 9  import alg_dpa_trial as dpa
10
11  ##### Q1 Solution #####
12  # For this question, your task is to load a provided citation graph for 27,770
13  # high energy physics theory papers. This graph has 352,768 edges. You should
14  # use the following code to load the citation graph as a dictionary. In
15  # CodeSkulptor, loading the graph should take 5-10 seconds. (For an extra
16  # challenge, you are welcome to write your own function to create the citation
17  # graph by parsing this text representation of the citation graph.)
18  #
19  # Your task for this question is to compute the in-degree distribution for this
20  # citation graph. Once you have computed this distribution, you should normalize
21  # the distribution (make the values in the dictionary sum to one) and then
22  # compute a log/log plot of the points in this normalized distribution. How you
23  # create this point plot is up to you. You are welcome to use a package such as
24  # matplotlib for desktop Python, use the simpleplot module in CodeSkulptor, or
25  # use any other method that you wish
26
27  # load the graph from the text file
28  cit_graph = parse_graph.load_graph("citation_graph.txt")
29
30  # get the unnormalized in degree distribution
31  in_deg_dist = alg_proj1_sol.in_degree_distribution(cit_graph)
32
33  # normalize the in degree distribution
34  sum_val = sum(in_deg_dist.values())
35  in_deg_dist.update((degree, freq / float(sum_val))  for degree, freq in in_deg_dist.items())
36
37  # draw the loglog plot of the normalized in degree distribution of the citation graphe
38  plt.figure(0)
39  plt.loglog(in_deg_dist.keys(), in_deg_dist.values(), basex=10, basey=10, linestyle='None',
40          marker='.', markeredgecolor='blue')
41  plt.title('loglog plot of in-degree distribution of citation Graph')
42  plt.xlabel('number of citations')
43  plt.ylabel('fraction of papers')
44  plt.grid()
45  plt.ylim(None, 1)
46  plt.show()
47  #plt.savefig("Q1_loglog_degree_dist_citgraph.png")
48
49  ##### Q2 Solution #####
50  # In Homework 1, you saw Algorithm ER for generating random graphs and reasoned
51  # analytically about the properties of the ER graphs it generates. Consider the
52  # simple modification of the algorithm to generate random directed graphs: For
53  # every ordered pair of distinct nodes (i, j), the modified algorithm adds the
54  # directed edge from i to j with probability p.
55  #
56  # For this question, your task is to consider the shape of the in-degree
57  # distribution for an ER graph and compare its shape to that of the physics
58  # citation graph. In the homework, we considered the probability of a specific
59  # in-degree, k, for a single node. Now, we are interested in the in-degree
60  # distribution for the entire ER graph. To determine the shape of this
61  # distribution, you are welcome to compute several examples of in-degree
62  # distributions or determine the shape mathematically.
63  #
64  # Once you have determined the shape of the in-degree distributions for ER graphs,
65  # compare the shape of this distribution to the shape of the in-degree distribution
66  # for the citation graph. When answering this question, make sure to address the
67  # following points:
68  #
69  # Q2.1: Is the expected value of the in-degree the same for every node in an ER graph?
70  # Please answer yes or no and include a short explanation for your answer.
71
72  # Ans: yes it same for all nodes since the presence of an edge is independent of
73  # all other edges i.e it is independent of the current structure of the graph
74  # The expected value of in-degree is given by p * (n-1)
75
76  # Q2.2: What does the in-degree distribution for an ER graph look like?
77  # Provide a short written description of the shape of the distribution.
78
79  # Ans: we know that the probability that a given node has degree k is given by
80  # a binomial distribution as seen in the homework. Thus as p -> 0 (probability
81  # p becomes smaller), we see more nodes with smaller in-degree and thus
82  # the in-degree distribution shape looks like a bell curve skewed towards the
83  # left i.e near in-degree 0. As p -> 1, we get more nodes with higher in-degree
84  # and the shape is increasing curve with most points near the higher
85  # in-degree region. For large number of nodes, and small p, this becomes
86  # a symmetric bell shaped curve and approaches a normal distribution
87
88  # Q2.3: Does the shape of the in-degree distribution plot for ER look similar
89  # to the shape of the in-degree distribution for the citation graph?
90  # Provide a short explanation of the similarities or differences.
91  # Focus on comparing the shape of the two plots as discussed in the class page on
92  # "Creating, formatting, and comparing plots".
93
94  # Ans: As mentioned in answer of Q2.2, the shape for the ER in-degree approaches
```

```python
 95    # a bell−shaped curve for large N values and small p. However, for the citation
 96    # graph it is a decreasing curve with majority of point located near in−degree
 97    # of zero.
 98
 99    ##### Q3 Solution #####
100    # We next consider a different process for generating synthetic directed graphs.
101    # In this process, a random directed graph is generated iteratively, where in
102    # each iteration a new node is created, added to the graph, and connected to a
103    # subset of the existing nodes. This subset is chosen based on the in−degrees
104    # of the existing nodes. More formally, to generate a random directed graph in
105    # this process, the user must specify two parameters: nn, which is the final
106    # number of nodes, and m (where m <= n), which is the number of existing
107    # nodes to which a new node is connected during each iteration. Notice that m
108    # is fixed throughout the procedure.
109    #
110    # The algorithm starts by creating a complete directed graph on mm nodes.
111    # (Note, you've already written the code for this part in the Project.) Then,
112    # the algorithm grows the graph by adding n−m nodes, where each new node is
113    # connected to m nodes randomly chosen from the set of existing nodes. As an
114    # existing node may be chosen more than once in an iteration, we eliminate
115    # duplicates (to avoid parallel edges); hence, the new node may be connected
116    # to fewer than m existing nodes upon its addition.
117
118    # The algorithm is called Algorithm DPA (note that the m in the input is a
119    # parameter that is specified to this algorithm, and it does not denote the
120    # total number of edges in the resulting graph).
121
122    # For this question, we will choose values for n and m that yield a DPA
123    # graph whose number of nodes and edges is roughly the same to those of the
124    # citation graph. For the nodes, choosing n to be the number of nodes as
125    # the citation graph is easy. Since each step in the DPA algorithm adds m
126    # edges to the graph, a good choice for m is an integer that is close to
127    # the average out−degree of the physics citation graph.
128
129    # For this question, provide numerical values for n and m that you will
130    # use in your construction of the DPA graph.
131
132    # calculate n, i.e the number of nodes in citation graph for DPA algorithm
133    n_nodes = len(cit_graph.keys())
134    # calculate m, i.e the average out−degree in citation graph for DPA algorithm
135    m_nodes = int(round(np.mean([len(neighbors) for  neighbors in cit_graph.values()])))
136
137    print "Q3 Solution:"
138    print "n = ", n_nodes
139    print "m = ", m_nodes
140
141    ##### Q4 Solution #####
142    # Your task for this question is to implement the DPA algorithm, compute a DPA
143    # graph using the values from Question 3, and then plot the in−degree distribution
144    # for this DPA graph. Creating an efficient implementation of the DPA algorithm
145    # from scratch is surprisingly tricky. The key issue in implementing the algorithm
146    # is to avoid iterating through every node in the graph when executing Line 6
147    # of the provided pseudocode. Using a loop to implement Line 6 leads to implementations
148    # that require on the order of 30 minutes in desktop Python to create a DPA graph with
149    # 28000 nodes.
150    #
151    # To avoid this bottleneck, you are welcome to use this provided code that implements
152    # a DPATrial class.
153    #
154    # Once you have created a DPA graph of the appropriate size, compute a (normalized)
155    # log/log plot of the points in the graph's in−degree distribution
156
157    # write the function for generating DPA graphs
158    def alg_dpa(n_num_nodes, m_num_nodes):
159        """
160        Uses the DPA algorithm provided in Q3 of the Application
161        to generates a random directed graph iteratively, where
162        each iteration a new node is created, added to the graph,
163        and connected to the subset of the existing node
164
165        Arguments:
166            n_nodes {integer} −− final number of nodes in the generated graph
167            m_nodes {integer} −− number of existing nodes to which a new node is connected
168                                 during each iteration
169
170        Returns:
171            dictionary −− the generated graph based on DPA algorithm
172        """
173
174        # create a complete graph of m_nodes noes
175        graph = alg_proj1_sol.make_complete_graph(m_num_nodes)
176
177        # create the DPA trial object corresponding to complete graph
178        dpa_trial = dpa.DPATrial(m_num_nodes)
179
180        # add each new ode to m_nodes from the existing graph randomly
181        # chosen with probability:
182        # (in−degree of new_node + 1) / (in−degree of all nodes +
183        # total number of existing nodes)
184        # simulated by the run_trial of the DPATrial class
185        for new_node in range(m_num_nodes, n_num_nodes):
186            # randomly select m_nodes from the existing graph that
187            # the new_node will be connected to. Remove if any
188            # duplicate nodes in the m_nodes selected
189            new_node_neighbors = dpa_trial.run_trial(m_num_nodes)
190
```

```python
191            # update the existing graph to add this new node and its
192            # neighbors
193            graph[new_node] = new_node_neighbors
194
195
196        return graph
197
198 # create the graph using the DPA algorithm
199 dpa_graph = alg_dpa(n_nodes, m_nodes)
200
201 # get the in-degree distribution for the DPA graph
202 in_deg_dist_dpa = alg_proj1_sol.in_degree_distribution(dpa_graph)
203
204 # normalize the in degree distribution for the DPA graph
205 sum_val = sum(in_deg_dist_dpa.values())
206 in_deg_dist_dpa.update((degree, freq / float(sum_val))  for degree, freq in in_deg_dist_dpa.items())
207
208 # draw the loglog plot of the normalized in-degree distribution of the DPA graph
209 plt.figure(1)
210 plt.loglog(in_deg_dist_dpa.keys(), in_deg_dist_dpa.values(), basex=10, basey=10,
211            linestyle='None', marker='.', markeredgecolor='blue')
212 plt.title('loglog plot of In-degree distribution of DPA graph')
213 plt.xlabel('in-degree')
214 plt.ylabel('fraction of nodes')
215 plt.ylim(None, 1)
216 plt.grid()
217 plt.show()
218 #plt.savefig("Q4_loglog_indegree_dist_dpa.png")
219
220 ##### Q5 Solution #####
221 # In this last problem, we will compare the in-degree distribution for the citation graph
222 # to the in-degree distribution for the DPA graph as constructed in Question 4. In
223 # particular, we will consider whether the shape of these two distributions are similar
224 # and, if they are similar, what might be the cause of the similarity.
225 #
226 # To help you in your analysis, you should consider the following three phenomena:
227 # - The "six degrees of separation" phenomenon,
228 # - The "rich gets richer" phenomenon, and
229 # - The "Hierarchical structure of networks" phenomenon.
230 #
231 # Your task for this problem is to consider how one of these phenomena might explain
232 # the structure of the citation graph or, alternatively, how the citations patterns
233 # follow one of these phenomena.
234 #
235 # When answering this question, please include answers to the following:
236 #
237 # Q5.1: Is the plot of the in-degree distribution for the DPA graph similar to that of the
238 # citation graph? Provide a short explanation of the similarities or differences.
239 # Focus on the various properties of the two plots as discussed in the class page on
240 # "Creating, formatting, and comparing plots".
241
242 # Ans: Yes they are similar since both follow a linear log-log decreasing trend i.e
243 # they both follow the power law distribution and the point are spread out more
244 # as the in-degree increases
245
246 # Q5.2: Which one of the three social phenomena listed above mimics the behavior of the DPA
247 # process? Provide a short explanation for your answer.
248
249 # Ans: DPA process mimics the "rich get richer" or the "preferential attachment" phenomena
250 # since every new node that is added to the graph is most likely to be connected to
251 # the neighbor with highest in-degree.
252
253 # Q5.3: Could one of these phenomena explain the structure of the physics citation graph?
254 # Provide a short explanation for your answer.
255
256 # Ans: The citation graph also mimics the "rich get richer" phenomena as the paper with
257 # higher citations i.e higher degree tend to be more likely used in other papers
258 # as well due to being more visible
```

# B   All functions for project 4 used in the application

```python
1  """
2  Functions for Prject #1: "Degree Distribution for Graphs". These functions will be
3  used in the Application #1: "Analysis of Citation Graphs"
4  """
5
6  # define directed graph constants for testing
7  EX_GRAPH0 = {
8      0 : set([1, 2]),
9      1 : set(),
10     2 : set(),
11     }
12
13 EX_GRAPH1 = {
14     0 : set([1, 4, 5]),
15     1 : set([2, 6]),
16     2 : set([3]),
17     3 : set([0]),
18     4 : set([1]),
19     5 : set([2]),
20     6 : set([]),
21     }
22
23 EX_GRAPH2 = {
```

```python
24        0 : set([1, 4, 5]),
25        1 : set([2, 6]),
26        2 : set([3, 7]),
27        3 : set([7]),
28        4 : set([1]),
29        5 : set([2]),
30        6 : set(),
31        7 : set([3]),
32        8 : set([1, 2]),
33        9 : set([0, 3, 4, 5, 6, 7]),
34        }

36    def make_complete_graph(num_nodes):
37        """
38        create and return a complete graph with nodes from
39        0 to num_nodes - 1 for num_nodes > 0. Otherwise
40        the function returns a dictionary corresponding to
41        the empty graph
42
43        Arguments:
44            num_nodes {integer} -- number of nodes for the graph
45
46        Returns:
47            dictionary -- returns a dictionary corresponding to a complete directed
48            graph with the specified number of nodes.
49        """
50        # local variable for the complete graph
51        graph = {}
52
53        # return an empty graph if num_nodes is not positive
54        if num_nodes <= 0:
55            return graph
56
57        for node in range(num_nodes):
58            # create an adjacency list for a directed complete graph with no
59            # self loops or parallel edges
60            graph[node] = set([val for val in range(num_nodes) if val != node])
61
62        return graph

64    def compute_in_degrees(digraph):
65        """
66        computes the in-degree of the nodes in a graph
67
68        Arguments:
69            digraph {dictionary} -- a directed graph with no self loop or parallel edges
70
71        Returns:
72            dictionary -- returns a dictionary with same set of keys(nodes) as digraph
73                          whose corresponding values are the number of edges whose
74                          head matches a particular node
75
76        """
77
78        # initialize the in degree for the nodes of digraph to 0
79        in_degree = dict(zip(digraph.keys(), len(digraph) * [0]))
80
81        for tail_node in digraph:
82            for head_node in digraph[tail_node]:
83                in_degree[head_node] += 1
84
85        return in_degree

87    def in_degree_distribution(digraph):
88        """
89        computes the unnormalized distribution of the in-degrees of the graph
90
91        Arguments:
92            digraph {dictionary} -- a directed graph with no self loops or parallel
93                                    edges
94
95        Returns:
96            dictionary -- unnormalized distribution of the in-degrees of the graph
97                          with key being the in-degree of nodes in the graph and
98                          the value associated with each particular in-degree is
99                          the number of nodes with that in-degree. In-degrees with
100                         no corresponding nodes in the graph are not included in
101                         the dictionary.
102        """
103        # initialize the dictionary to store
104        in_degree_dist = {}
105        # get the in-degree for each node
106        in_degrees = compute_in_degrees(digraph)
107
108        for degree_vals in in_degrees.values():
109            if in_degree_dist.has_key(degree_vals):
110                in_degree_dist[degree_vals] += 1
111            else:
112                in_degree_dist[degree_vals] = 1
113
114        return in_degree_dist
```

# C   Code to load the graphs from text file

```
1   """
2   Common functions used for the both Application #1 and #2
3   """
4
5   def load_graph(graph_file):
6       """
7       Helper function to solve Q1 Application #1: Analysis of
8       Citation Graphs
9       converts the text representation of a graph from
10      a text file to dictionary representation.
11
12      Arguments:
13          graph_file {string} — a file name of the file with text representation of a graph
14
15      Returns:
16          dictionary — returns a dictionary representation of the graph
17      """
18      # will store the dictionary representation of the graph
19      graph = {}
20
21      with open(graph_file) as grh_file:
22          for line in grh_file:
23              # get the tail node and corresponding head nodes in the current
24              # line of the adjacency list text representation of graph
25              nodes = line.split()
26              # convert the head node string to integer
27              head_nodes = map(int, nodes[1:])
28              # add the key, the tail node, and value, the head nodes to the
29              # dictionary representation of the graph
30              graph[int(nodes[0])] = set(head_nodes)
31
32      return graph
```