

# Edge Replacement Grammars: A Formal Language Approach for Generating Graphs

Revanth Reddy<sup>1\*</sup>, Sarath Chandar<sup>2\*</sup>, Balaraman Ravindran<sup>1,3</sup>

<sup>1</sup>Department of Computer Science and Engineering, Indian Institute of Technology Madras

<sup>2</sup>Mila, Université de Montréal

<sup>3</sup>Robert Bosch Centre for Data Science and AI, Indian Institute of Technology Madras

## Motivation

- ▶ A generative model that closely mimics the structural properties of a given set of graphs has utility in a variety of domains.
- ▶ The model can be used to anonymize graph data, by generating graphs similar to the original graphs.
- ▶ If we are able to fit the model more accurately, we can just save the model instead of the entire graph data.
- ▶ We can do graph classification by determining the notion of likelihood of the test graph as per the given model.
- ▶ Our graph model is based on a variant of Probabilistic Edge Replacement Grammar (PERG) called Restricted PERG (RPERG).

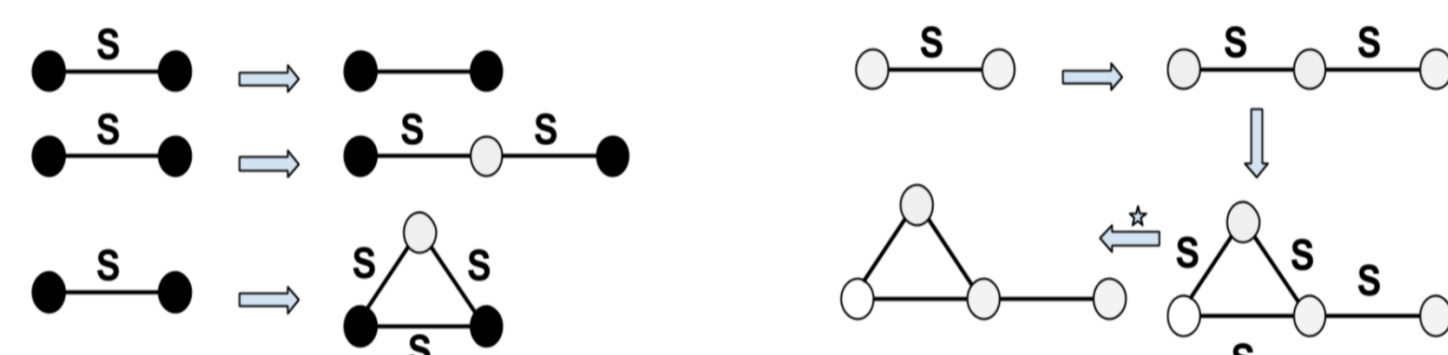
## Definitions

**Definition 1:** An Edge Replacement Grammar (ERG) is a tuple  $G = \langle N, T, P, S \rangle$  where

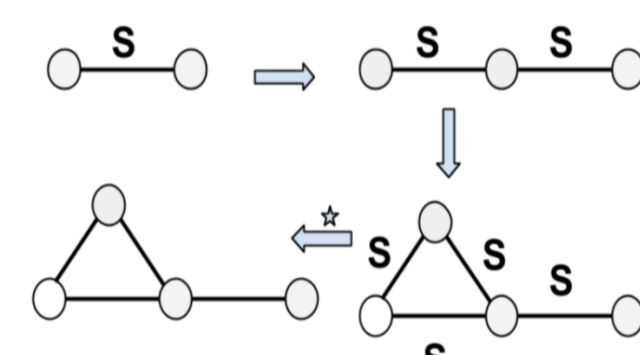
- ▶  $N$  and  $T$  are finite disjoint set of non-terminal and terminal edge labels.
- ▶  $S \in N$  is the start edge label.
- ▶  $P$  is a finite set of productions of the form  $A \rightarrow R$ , where  $A \in N$  and  $R$  is a graph fragment with edge labels drawn from  $N \cup T$ .

**Definition 2:** A Probabilistic Edge Replacement Grammar (PERG) consists of

- ▶ An edge replacement grammar  $G = \langle N, T, P, S \rangle$
- ▶ A parameter  $p(A \rightarrow R)$  for each rule  $A \rightarrow R \in P$ . This parameter is the conditional probability of choosing this rule given that the non-terminal being expanded is  $A$  with the following constraint, for any  $X \in N$ ,  $\sum_{A \rightarrow R: A=X} p(A \rightarrow R) = 1$



(a) Sample ERG



(b) Sample derivation

**Definition 3:** Let  $u, v$  be a pair of vertices in the graph  $G$ . Let  $g_1, g_2, \dots, g_t$  be the connected components obtained by removing  $u, v$  from  $G$ . A squeezing operation with respect to  $u, v$  is an operation where one of the components  $g_i$  is replaced by an edge between  $u, v$ .

When  $t = 1$ , the entire graph is squeezed into a single edge. If  $t \geq 3$  and  $g_1, \dots, g_t$  are isolated vertices, then squeeze operation replaces entire graph with the edge  $u, v$ . If a graph can be squeezed into a single edge, it is a **trivial squeeze**.

**Definition 4:** A non-squeezable graph is a graph in which the only squeeze operation that is possible is the trivial squeeze. Star graphs and triangle are considered as the degenerate cases.

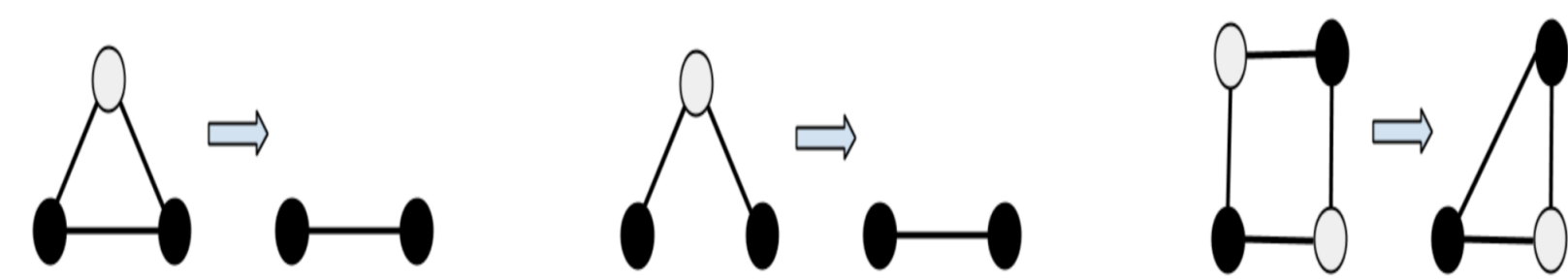


Figure: Examples for Squeezing

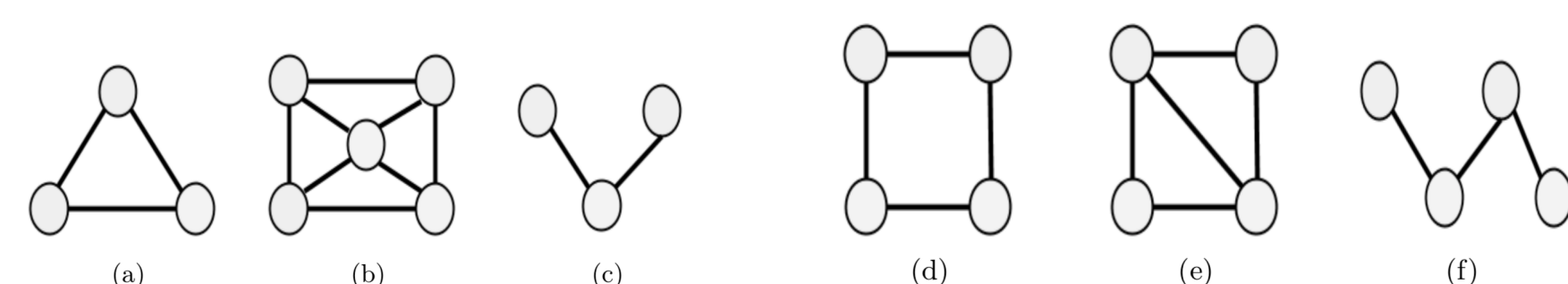


Figure: a,b,c are non-squeezable while d,e,f are squeezable

**Definition 5:** A Restricted Probabilistic Edge Replacement Grammar (RPERG) is a PERG such that for every rule  $A \rightarrow R \in \text{RPERG}$ ,  $R$  is a non-squeezable graph fragment.

## Learning the Grammar

If  $c_D(A \rightarrow R)$  is the count of the occurrences of the sub-graph  $R$  in the data  $D$ , then the maximum likelihood estimation of the parameters of the model is given by

$$p_{ML}^{A \rightarrow R} = \frac{c_D(A \rightarrow R)}{\sum_{R': A \rightarrow R'} c_D(A \rightarrow R')}$$

The learning problem now has been reduced to getting the counts of non-squeezable components in a graph.

### Algorithm 1 Learn RPERG

**Input:** Set of Graphs  $D = \{g_1, g_2, \dots, g_n\}$ .

**Output:** RPERG

```

1: function MAIN(Set of Graphs D)
2:   Stack ← empty stack
3:   for each graph  $g_i$  do
4:     GET_COMPONENTS( $g_i$ )
5:     while Stack is not empty do
6:        $g \leftarrow$  Stack.pop()
7:       find a split pair (a,b) in  $g$ 
8:       if  $\exists$  no split pair then
9:          $C(A \rightarrow g) + 1$ 
10:      else
11:         $g_1, g_2 \leftarrow$  Obtained by splitting  $g$  at (a,b)
12:        if edge(a,b)  $\notin$   $g$  then
13:          Add edge(a,b) to  $g_2$ 
14:        end if
15:        for  $g'$  in  $g_1, g_2$  do
16:          GET_COMPONENTS( $g'$ )
17:        end for
18:      end if
19:    end while
20:  end for
21: end function

```

```

1: function GET_COMPONENTS(Graph g)
2:   CV ← cut vertices in g
3:   for each  $v_k$  in CV do
4:      $n \leftarrow$  no. of bi-connected components connected by  $v_k$ 
5:      $C(A \rightarrow star(n)) + 1$ 
6:   end for
7:   S ← set of all bi-connected components in g
8:   for each  $s_i$  in S do
9:     Stack.push( $s_i$ )
10:  end for
11: end function

```

## Generative Model

We assume that the network is homogeneous and the links are un-weighted. Since we have only one type of link, number of nonterminal labels is one. The learning algorithm will consider all edges in the given graph to be non-terminal edges.

### Algorithm 2 Generative Model

**Input:** RPERG

**Output:** A Graph

```

1: Graph G = NULL
2: Add a non-terminal edge to G
3: while desired graph size is not reached do
4:   Randomly pick a non-terminal edge A in G.
5:   Sample a rule  $A \rightarrow R$  from RPERG and replace A with R in G.
6: end while
7: Convert all non-terminal edges in G to terminal edges.

```

## Results

We tested the proposed model by fitting it onto several real life graph datasets. The networks vary not only in the number of vertices and edges, but also in the clustering coefficient, diameter, degree distribution and many other graph properties. We compare the properties of the approximate graphs generated from RPERG, HRG, Chung-Lu and Kronecker graph models.

Dataset	Nodes	Edges	Diameter	Clust. Coeff.
Arxiv	5242	14496	17	0.529
Routers	6474	13895	9	0.252
Enron	36692	183831	11	0.497
DBLP	317080	1049866	21	0.632

(a) Dataset Statistics

Model	Color
Original Graph	Black
RPERG	Red
HRG	Blue
Chung-Lu	Green
Kronecker	Brown

(b) Color Coding

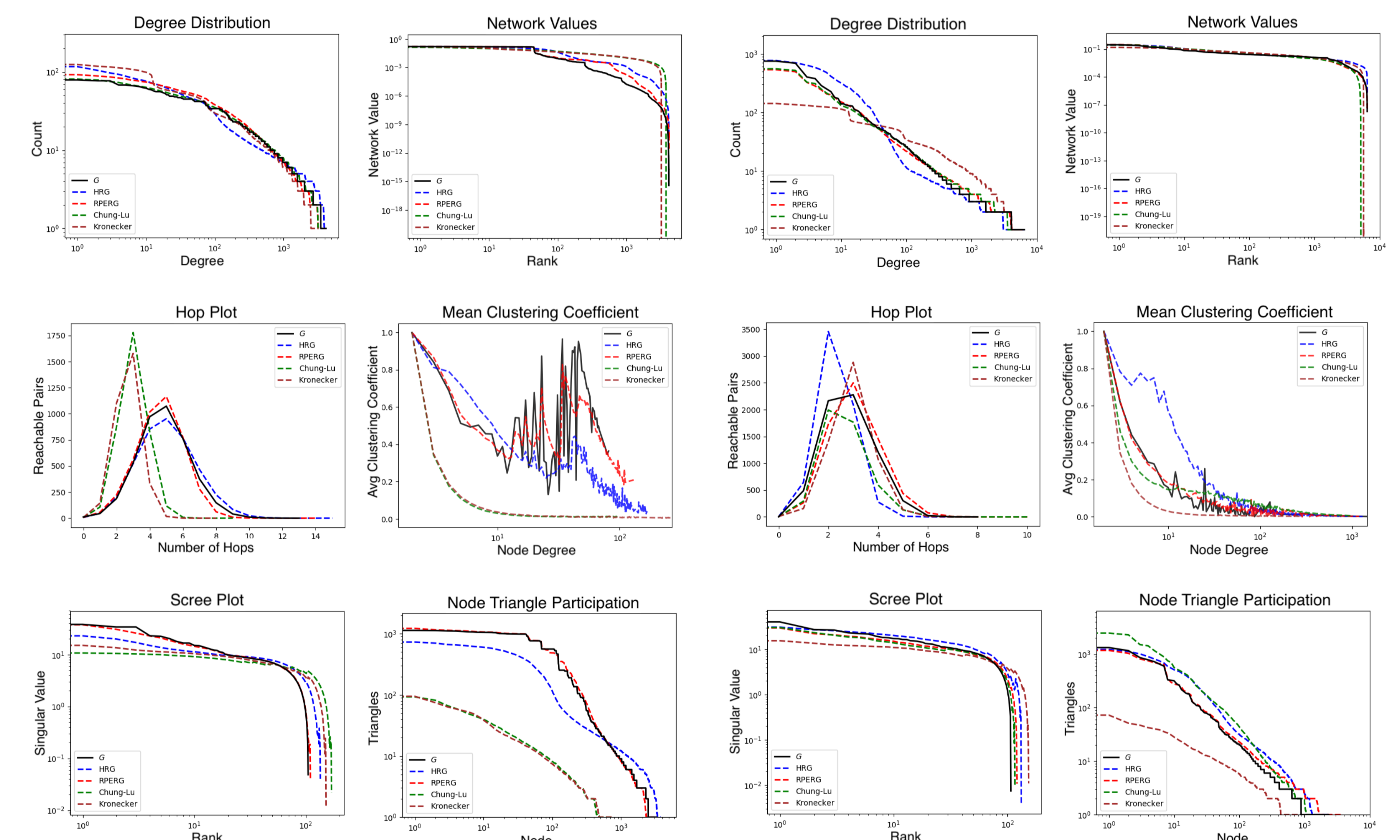


Figure: Arxiv GR-QC

Figure: Internet Routers

To more concretely compare the eigenvectors, the cosine distance between the eigenvector centrality of the original graph and the models' generated graphs is calculated.

Graphlet Correlation Distance (GCD) computes the distance between two graphlet correlation matrices by measuring the frequency of the various graphlets present in each graph, i.e the number of edges, wedges, triangles, squares, 4-cliques, etc., and comparing the graphlet frequencies between two graphs.

Dataset	RPERG	HRG	Chung-Lu	Kronecker
Arxiv	<b>0.0025</b>	0.0161	0.3496	0.3406
Routers	<b>0.0247</b>	0.0411	0.0379	0.0614
Enron	<b>0.00007</b>	0.0002	0.0052	0.0676
DBLP	<b>0.0079</b>	0.0649	0.5854	0.4997

(a) Cosine Distance

Dataset	RPERG	HRG	Chung-Lu	Kronecker
Arxiv	<b>1.086</b>	1.094	1.792	2.071
Routers	<b>1.293</b>	1.404	1.975	2.776
Enron	<b>0.487</b>	0.525	1.319	2.83
DBLP	<b>0.409</b>	1.602	1.738	2.821

(b) GCD Values

## Conclusion

- ▶ Even though the grammar is context free, it is able to capture most of the statistical properties of the graph.
- ▶ We observe that our algorithm is easily parallelizable as we can run the algorithm simultaneously on multiple graphs.
- ▶ In future work, we can try to model preferential attachment by converting the grammar into a context sensitive grammar.
- ▶ Tackling graphs with multiple types of links (heterogeneous links) is also a challenging problem.