

Growing Attributed Networks through Local Processes

Harshay Shah, Suhansanu Kumar, Hari Sundaram
University of Illinois at Urbana-Champaign
{hrshah4,skumar56,hs1}@illinois.edu

ABSTRACT

This paper proposes an attributed network growth model. Despite the knowledge that individuals use limited resources to form connections to similar others, we lack an understanding of how local and resource-constrained mechanisms explain the emergence of rich structural properties found in real-world networks. We make three contributions. First, we propose an interpretable and accurate model of attributed network growth that jointly explains the emergence of in-degree distribution, local clustering, clustering-degree relationship and attribute mixing patterns. Second, we make use of biased random walks to develop a model that forms edges locally, without recourse to global information. Third, we account for multiple sociological phenomena: bounded rationality; structural constraints; triadic closure; attribute homophily; preferential attachment. Our experiments show that the proposed Attributed Network Growth (ARW) model accurately preserves network structure and attribute mixing patterns of six real-world networks; it improves upon the performance of eight well-known models by a significant margin of 2.5–10 \times .

CCS CONCEPTS

• **Information systems** \rightarrow **Web applications**; *Data mining*; *Web mining*; • **Applied computing** \rightarrow *Sociology*.

KEYWORDS

Network growth; Network Structure; Attributed networks

ACM Reference Format:

Harshay Shah, Suhansanu Kumar, Hari Sundaram. 2019. Growing Attributed Networks through Local Processes. In *Proceedings of the 2019 World Wide Web Conference (WWW'19)*, May 13–17, 2019, San Francisco, CA, USA. ACM, New York, NY, USA, 8 pages. <https://doi.org/10.1145/3308558.3313640>

1 INTRODUCTION

We present a network growth model that explains how distinct structural properties of attributed networks can emerge from a local edge formation process. In real-world networks, individuals form edges with limited information and partial network access. Moreover, phenomena such as triadic closure and homophily {simultaneously influence individuals' decisions to form connections. Over time, these decisions cumulatively shape real-world networks to exhibit rich structural properties: heavy-tailed in-degree distribution, skewed local clustering and diverse attribute mixing patterns. However, we lack an understanding of local, resource-constrained

mechanisms that incorporate sociological factors to jointly explain the emergence of multiple structural properties. Additionally, accurate network growth models are useful for synthesizing networks and extrapolating existing real-world networks.

Well-known models of network growth tend to make unrealistic assumptions about how individuals form edges. Consider a simple stylized example: the process of finding a set of papers to cite when writing an article. In preferential attachment [3] or fitness [5, 10, 44] based models, a node making m citations would pick papers from the *entire* network in proportion to their in-degree or fitness respectively. This process assumes that individuals possess complete knowledge of in-degree or fitness of every node in the network. An equivalent formulation—vertex copying [25]—induces preferential attachment: for every citation, a node would pick a paper uniformly at random from *all* papers, and either cite it or copy its citations. Notice that vertex copying assumes individuals have complete access to the network and forms each edge independently. Although these models explain the emergence of power law degree distributions, they are unrealistic: preferential attachment and vertex copying require global node-level knowledge or complete network access respectively. Additionally, they do not account for the role of assortative mixing [35] via nodal attributes (e.g., venue of paper, political interests of Facebook users) in network formation.

Recent papers tackle resource constraints [32, 45, 47] as well as nodal attributes [12, 17]. However, the former disregard attributes and the latter do not provide a realistic representation of edge formation under resource constraints. Furthermore, both sets of models do not jointly preserve multiple structural properties. Developing an interpretable and accurate model of attributed network growth that accounts for resource constraints and observed sociological phenomena is non-trivial.

We make three key contributions. First, we propose a simple and accurate model of attributed network growth. Second, our model is based on local processes to form edges, without recourse to global network information. Third, our model unifies multiple sociological phenomena—bounded rationality; structural constraints; triadic closure; attribute homophily; preferential attachment—to jointly model global network structure and attribute mixing patterns.

The proposed model—Attributed Random Walk (ARW)—jointly explains the emergence of in-degree distribution, local clustering, clustering-degree relationship and attribute mixing patterns through a resource constrained mechanism based on random walks (see Figure 1). In particular, the model relies entirely on local information to grow the network, without access to information of all nodes. In ARW, incoming nodes select a seed node based on attribute similarity and initiate a biased random walk: at each step of the walk, the incoming node either jumps back to its seed or chooses an outgoing link or incoming link to visit another node; it links to each visited node with some probability and halts after it has exhausted its budget to form connections. Our experiments on six

This paper is published under the Creative Commons Attribution 4.0 International (CC-BY 4.0) license. Authors reserve their rights to disseminate the work on their personal and corporate Web sites with the appropriate attribution.

WWW '19, May 13–17, 2019, San Francisco, CA, USA

© 2019 IW3C2 (International World Wide Web Conference Committee), published under Creative Commons CC-BY 4.0 License.

ACM ISBN 978-1-4503-6674-8/19/05.

<https://doi.org/10.1145/3308558.3313640>

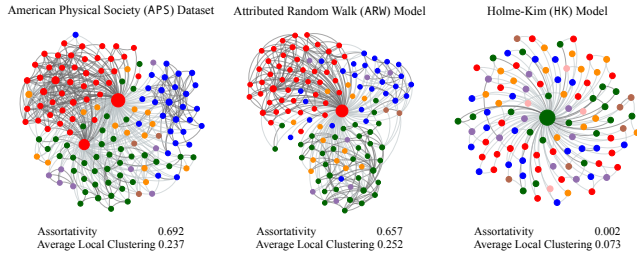


Figure 1: We contrast our proposed model, Attributed Random Walk (ARW), with a non-attributed growth model [20] to underscore the importance of using attributes for network growth.

large-scale network datasets indicate that the proposed growth model outperforms eight state-of-the-art network growth models, including attributed growth models, by a statistically significant margin of 2.5–10 \times .

The rest of the paper is organized as follows. We begin by defining the problem statement in Section 2. In Section 3, we outline six network datasets, describe key structural properties of real-world networks and discuss insights from sociological studies. Then, in Section 4, we describe the network growth model. Then, We present experiments in Section 5, discuss related work in Section 6 and conclude in Section 7.

2 PROBLEM STATEMENT

Consider an attributed directed network $G = (V, E, B)$, where V & E are sets of nodes & edges and each node has an attribute value $b \in B$. The goal is to develop a directed network growth model that preserves structural and attribute based properties observed in G . The growth model should be normative, accurate and parsimonious:

- (1) **Normative:** The model should account for normative behavior. In real-world networks, multiple sociological phenomena influence how individuals form edges under constraints of limited global information and partial network access.
- (2) **Accurate:** The model should preserve key structural and attribute based properties such as heavy tailed degree distribution, skewed local clustering, negatively correlated degree-clustering relationship and attribute mixing patterns.
- (3) **Interpretability:** The model should be expressive enough to generate networks with varying structural properties, while having as few parameters as possible.

Next, we present empirical analysis on real-world datasets to motivate our attributed random walk model.

3 EMPIRICAL ANALYSIS

We begin by describing six large-scale network datasets that we use in our analysis and experiments. Then, we describe global network properties, insights from empirical studies in the Social Science and common assumptions in network modeling. Finally, we discuss the role of structural proximity in edge formation.

3.1 Datasets

We consider six citation networks of different scales (size, time) from diverse sources: research articles, utility patents and judicial cases. We list the summary statistics and global network properties of these datasets in Table 1. Three of the six datasets are attributed networks; that is, each node has a categorical attribute value.

We focus on citation networks for two reasons. First, since nodes in citation networks form all outgoing edges to existing nodes at the time of joining the network, citation networks provide a clean basis to study edge formation mechanisms in attributed social networks. Second, citation network span long periods of time (e.g., the USSC judicial citation network span several hundred years). As a result, identifying local edge formation processes that accurately model growth for this duration is non-trivial. Next, we study the structural and content properties of these networks.

3.2 Global Network Properties

Compact statistical descriptors of global network properties [33] such as degree distribution, local clustering, and attribute assortativity quantify the extent to which local edge formation phenomena shape global network structure.

Heavy tailed degree distribution: Real-world networks tend to exhibit heavy tailed degree distributions. These distributions can emerge from the well-known preferential attachment process [3, 40], where incoming nodes connect with nodes in proportion to their degree. Log-normal fits, with parameters listed in Table 1, well describe the in-degree distribution of all network datasets, consistent Broido and Clauset’s [9] observation that scale-free, real-world networks are rare

High Local Clustering: Real-world networks exhibit high local clustering (LCC), as shown in Table 1. Local clustering can arise from triadic closure [34, 39], where nodes with common neighbor(s)

Network	Description	$ V $	$ E $	T	$A, A $	$\text{LN}(\mu, \sigma)$	DPL α	Avg. LCC	AA r
USSC [14]	U.S. Supreme Court cases	30,288	216,738	1754-2002	-	(1.19, 1.18)	2.32	0.12	-
HEP-PH [15]	ArXiv Physics manuscripts	34,546	421,533	1992-2002	-	(1.32, 1.41)	1.67	0.12	-
Semantic [2]	Academic Search Engine	7,706,506	59,079,055	1991-2016	-	(1.78, 0.96)	1.58	0.06	-
ACL [37]	NLP papers	18,665	115,311	1965-2016	VENUE, 50	(1.93, 1.38)	1.43	0.07	0.07
APS [1]	Physics journals	577,046	6,967,873	1893-2015	JOURNAL, 13	(1.62, 1.20)	1.26	0.11	0.44
Patents [27]	U.S. NBER patents	3,923,922	16,522,438	1975-1999	CATEGORY, 6	(1.10, 1.01)	1.94	0.04	0.72

Table 1: Summary statistics & global properties of six network datasets: $|V|$ nodes join the networks and form edges $|E|$ over time period T . In attributed networks, each node has a categorical attribute value that belongs to set A of size $|A|$. The networks exhibit lognormal (LN) in-degree distribution with mean μ and standard deviation σ , high average local clustering (LCC) & attribute assortativity (AA) coefficient and densify over time with power law (DPL) exponent α .

have an increased likelihood of forming a connection. The coefficient of node i equals the probability with which two randomly chosen neighbors of the node i are connected. In directed networks, the neighborhood of a node i can refer to the nodes that link to i , nodes that i links to or both. We define the neighborhood to be the set of all nodes that link to node i . In Figure 2, we show that (a) average local clustering is not a representative statistic of the skewed local clustering distributions and (b) real-world networks exhibit a negative correlation between in-degree and clustering. That is, low in-degree nodes have small, tightly knit neighborhoods and high in-degree nodes tend to have large, star-shaped neighborhoods.

Homophily: Attributed networks tend to exhibit homophily [29], the phenomenon where similar nodes are more likely to be connected than dissimilar nodes. The assortativity coefficient [35] $r \in [-1, 1]$, quantifies the level of homophily in an attributed network. Intuitively, assortativity compares the observed fraction of edges between nodes with the same attribute value to the expected fraction of edges between nodes with same attribute value if the edges were rewired randomly. In Figure 3, we show that attributed networks ACL, APS and Patents exhibit varying level of homophily with assortativity coefficient ranging from 0.07 to 0.72.

Increasing Out-degree over Time: The out-degree of nodes that join real-world networks tends to increase as functions of network size and time. This phenomenon densifies networks and can shrink effective diameter over time. Densification tends to exhibit a power law relationship [27] between the number of edges $e(t)$ and nodes $n(t)$ at time t : $e(t) \propto n(t)^\alpha$. Table 1 lists the densification power law (DPL) exponent α of the network datasets.

To summarize, citation networks tend to be homophilic networks that undergo accelerated network growth and exhibit regularities in structural properties: heavy tailed in-degree distribution, skewed local clustering distribution, negatively correlated degree-clustering relationship, and varying attribute mixing patterns.

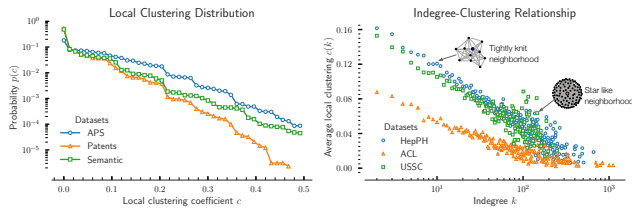


Figure 2: Local clustering in real-world networks have common characteristics: skewed local clustering distribution (left subplot) and a negatively correlated relationship between in-degree and average local clustering (right subplot).



Figure 3: Attributed networks exhibit varying levels of homophily. The subplots illustrate the mixing patterns in ACL, APS and Patents w.r.t. attributes Venue ($r = 0.07$), Journal ($r = 0.44$) and Category ($r = 0.72$) respectively.

3.3 Insights from Sociological Studies

Sociological studies on network formation seek to explain how individuals form edges in real-world networks.

Interplay of Triadic Closure and Homophily: Empirical studies [6, 24] that analyze the interplay between triadic closure and homophily indicate that *both* structural proximity and homophily are statistically significant factors that simultaneously influence edge formation. Homophilic preferences [29] induce edges between similar nodes, whereas structural factors such as network distance limit edge formation to proximate nodes (e.g. friend of a friend).

Bounded Rationality: Extensive work [16, 28, 41] on decision making shows that individuals are boundedly rational actors; constraints such as limited information, cognitive capacity and time impact decision making. This suggests that resource-constrained individuals that join networks are likely to employ simple rules to form edges using limited information and partial network access.

Current preferential attachment and fitness-based models [3, 13, 42] make two assumptions that are at variance with these findings. First, by assuming that successive edge formations are independent, these models disregard the effect of triadic closure and structural proximity. Second, these models implicitly require incoming nodes to have complete network access (e.g., be able to connect to any node) or explicit knowledge of one or more properties (e.g., fitness, degree) of every node. For example, a preferential attachment model, by making connections in proportion to degree, requires non-local information: the degree distribution of the entire network.

To summarize, insights from sociological studies indicate that edge formation in real-world networks comprises biases towards nodes that are similar, well-connected or structurally proximate. Next, we analyze the role of structural proximity in edge formation.

To summarize, empirical analyses and insights from the Social Sciences motivate the need to model how resource-constrained edge formation processes collectively shape well-defined global network properties of large-scale networks over time.

4 ATTRIBUTED RANDOM WALK MODEL

We propose an Attributed Random Walk (ARW) model to explain the emergence of key structural properties of real-world networks through entirely local edge formation mechanisms.

Consider a stylized example of how a researcher might go about finding relevant papers to cite. First, the researcher broadly identifies one or more relevant papers, possibly with the help of external information (e.g. Google Scholar). These initial set of papers act as seed nodes. Then, acting under time and information constraints, she will examine papers cited by the seed and papers that cite the seed. Thus, she navigates a chain of backward and forward references to identify similar, relevant papers. Next, through careful analysis, she will cite a subset of these papers. Similarly, users in online social networks might form new friendships by navigating their social circle (e.g., friends of friends) to find similar others.

ARW grows a directed network as new nodes join the network. The mechanism is motivated by the stylized example: an incoming node selects a seed node and initiates a random walk to explore the network by navigating through neighborhoods of existing nodes. It halts the random walk after connecting to a few visited nodes.

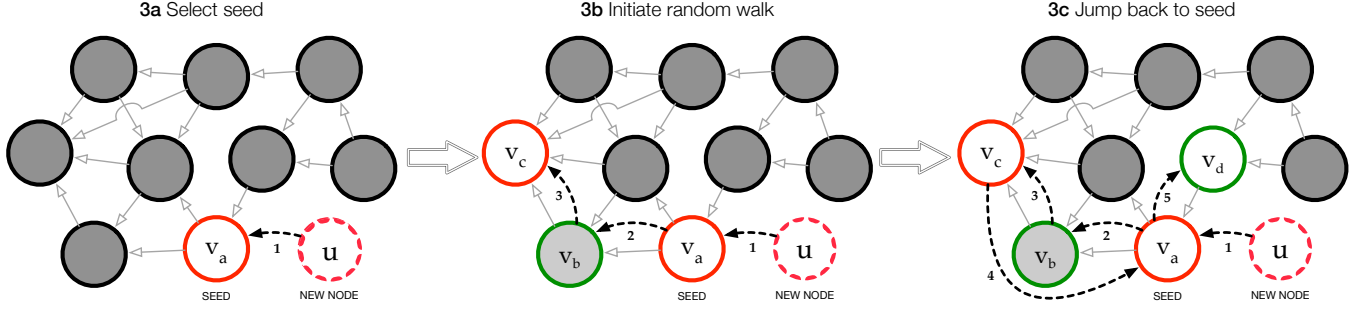


Figure 4: Edge formation in ARW: consider an incoming node u with outdegree $m = 3$ and attribute value $B(u) = \text{RED} \in \{\text{RED}, \text{GREEN}\}$. In fig. 3a, u joins the network and selects seed v_a via SELECT-SEED. Then, in fig. 3b, u initiates a RANDOM-WALK and traverses from v_a to v_b to v_c . Finally, u jumps back to its seed v_a and restarts the walk, as shown in fig. 3c. Node u halts the random walk after linking to v_a , v_c & v_d .

In this section, we describe the edge formation mechanisms underlying ARW, explain how ARW unifies multiple sociological phenomena, discuss model interpretability and summarize methods required to fit ARW to network data.

4.1 Model Description

The Attributed Random Walk (ARW) model grows a directed network $\{\hat{G}_t\}_{t=1}^T$ in T time steps. More formally, at every discrete time step t , a new node u , with attribute value $B(u)$, joins the network \hat{G}_t . After joining the network, node u forms $m(t)$ edges to existing nodes.

The edge formation mechanism consists of two components: SELECT-SEED and RANDOM-WALK. As shown in Figure 4, an incoming node u with attribute value $B(u)$ that joins the network at time t first selects a seed node using SELECT-SEED:

SELECT-SEED

- (1) With probability $p_{\text{same}}/p_{\text{same}}+p_{\text{diff}}$, randomly select a seed node from existing nodes that have the same attribute value, $B(u)$.
- (2) Otherwise, with probability $p_{\text{diff}}/p_{\text{same}}+p_{\text{diff}}$, randomly select a seed node from existing nodes that do *not* have the same attribute value, $B(u)$.

SELECT-SEED accounts for homophilic preferences of incoming nodes using parameters p_{same} and p_{diff} , which incorporate attribute preferences of incoming nodes. As shown in Figure 4, after selecting the seed node, u initiates a random walk using RANDOM-WALK to form $m(t)$ links. The RANDOM-WALK mechanism consists of four parameters: attribute-based parameters p_{same} & p_{diff} model edge formation decisions and the jump parameter p_{jump} & out-link parameter p_{out} characterize random walk traversals:

RANDOM-WALK

- (1) At each step of the walk, new node u visits node v_i .
 - If $B(u) = B(v_i)$, u links to v_i with probability p_{same}
 - Otherwise, u links to v_i with probability p_{diff}
- (2) Then, with probability p_{jump} , u jumps back to seed s_u .
- (3) Otherwise, with probability $1 - p_{\text{jump}}$, u continues to walk. It picks an outgoing edge with prob. p_{out} or an incoming edge with prob. $1 - p_{\text{out}}$ to visit a neighbor of v_i .
- (4) Steps 1-3 are repeated until u links to $m(t)$ nodes.

When attribute data is absent, ARW simplifies further. A single link parameter p_{link} replaces both attribute parameters p_{same} & p_{diff} . SELECT-SEED reduces to uniform seed selection and in RANDOM-WALK, the probability of linking to visited nodes equals p_{link} .

Note that ARW has two exogenous parameters: the out-degree $m(t)$ and attribute $B(u)$ of incoming nodes. The attribute distribution varies with time as new attribute values (e.g., journals) crop up, necessitating an exogenous parameter. The parameter $m(t)$ is the mean-field value of out-degree m at time t in the observed network. While it is straightforward to model $m(t)$ endogenously by incorporating a densification power-law DPL exponent, exogenous factors (e.g., venue, topic) may influence node out-degree.

Next, we explain how each parameter is necessary to conform to normative behavior of individuals in evolving networks.

4.2 ARW and Normative Behavior

The Attributed Random Walk model unifies multiple sociological phenomena into its edge formation mechanisms.

Phenomenon 1. (Limited Resources) *Individuals are boundedly rational [16, 28, 41] actors that form edges under constraints of limited information, partial network access and finite cognitive capacity.*

ARW uses random walk traversals to incorporate constraints of limited information and partial network access. A new node u selects a seed node from which it initiates a biased random walk. Then, u uses simple rules to connect to each visited nodes probabilistically and halts the walk after forming $m(t)$ edges, as shown in Figure 4. Random walks require information only about the 1-hop neighborhood of a few visited nodes, thereby accounting for the constraints of limited information and partial network access.

Phenomenon 2. (Structural Constraints) *Structural factors such as network distance act as constraints that limit edge formation to proximate nodes. [24]*

We incorporate structural constraints into ARW using p_{jump} , the probability with which a new node jumps back to its seed node after each step of the random walk. This implies that the probability with which the new node is at most k steps from its seed node is $(1 - p_{\text{jump}})^k$; as a result, p_{jump} controls the extent to which nodes' random walks explore the network to form edges.

Phenomenon 3. (Triadic Closure) *Nodes with common neighbors have an increased likelihood of forming a connection. [39]*

When attribute data is absent, ARW controls the effect of triadic closure on link formation using p_{link} . This is because a new node u closes a triad through its random walk by linking to both, a visited node and its neighbor, with probability proportional to p_{link}^2 . Similarly, in attributed networks, the probability of triad completion equals pq , where p and q can equal p_{same} or p_{diff} , depending on the attribute values of u and the visited nodes.

Phenomenon 4. (*Attribute Homophily*) *Nodes that have similar attributes are more likely to form a connection.* [29]

The attribute parameters p_{same} and p_{diff} modulate attribute assortativity. When $p_{\text{same}} > p_{\text{diff}}$, nodes are more likely to connect if they share the same attribute value, thereby resulting in a homophilic network over time. Similarly, $p_{\text{same}} < p_{\text{diff}}$ and $p_{\text{same}} = p_{\text{diff}}$ make edge formation heterophilic and attribute agnostic respectively.

Phenomenon 5. (*Preferential Attachment*) *Nodes tend to link to high degree nodes that have more visibility.* [3]

ARW controls preferential attachment by adding structural bias to the random walk traversal using outlink parameter p_{out} , instead of relying on the global degree distribution. Random walks that traverse outgoing edges only (i.e., $p_{\text{out}} = 1$) eventually visit old nodes that tend to have high in-degree. Similarly, random walks that traverse incoming edges only (i.e., $p_{\text{out}} = 0$) visit recently joined nodes that tend to have low indegree. As a result, we use p_{out} to adjust the effect of preferential attachment on edge formation.

To summarize: ARW incorporates five well-known sociological phenomena— bounded rationality; structural constraints; triadic closure; attribute homophily; preferential attachment—into a single edge formation mechanism based on random walks.

4.3 Model Fitting

We now briefly describe methods to estimate model parameters, initialize \hat{G} , densify \hat{G} over time and sample nodes' attribute values.

Parameter Estimation. The parameter estimation task consists of finding the set of parameters values for $(p_{\text{same}}, p_{\text{diff}}, p_{\text{jump}}, p_{\text{out}})$ that best preserve the structural properties of an observed network G . We use a straightforward grid search method to estimate the four parameters using evaluation metrics and selection criterion described in Subsection 5.1.

Initialization. ARW is sensitive to a large number of weakly connected components (WCCs) in initial network \hat{G}_0 because incoming nodes only form edges to nodes in the same WCC. To ensure that \hat{G}_0 is weakly connected, we perform an undirected breadth-first search on the observed, to-be-fitted network G that starts from the oldest node and halts after visiting 0.1% of the nodes. The initial network \hat{G}_0 is the small WCC induced from the set of visited nodes.

Node Out-degree. Node out-degree increases non-linearly over time in real-world networks. We coarsely mirror the growth rate of observed network G as follows. Each incoming node u that joins \hat{G} at time t corresponds to some node that joins the observed network G in year $y(t)$; the number of edges $m(t)$ that u forms is equal to the average out-degree of nodes that join G in year $y(t)$.

Sampling Attribute Values. The distribution over nodal attribute values $P_G(B)$ tends to change over time. The change in the attribute distribution over time is an exogenous factor and varies for every network. Therefore, we sample the attribute value $B(u)$ of node

Model	Abbreviation	Type	Attributed
Dorogovtsev et al. [13]	DMS	PA	No
Relay Linking [42]	RL	PA	No
Kim-Altman [22]	KA	PA	Yes
Social Attribute Network [17]	SAN	PA+TC	Yes
Holme-Kim [20]	HK	PA+TC	No
Herera-Zufria [19]	HZ	RW	No
Saramaki-Kaski [38]	SK	RW	No
Forest Fire [27]	FF	RW	No

Table 2: We evaluate the performance of our model ARW relative to 3 preferential attachment (PA) models, 2 pref. attachment & triangle closing (PA+TC) models and 3 random walk (RW) models.

u , that joins \hat{G} at time t , from $P_G(B \mid \text{year} = y(t))$, the observed attribute distribution conditioned on the year of arrival of node u .

To summarize, ARW intuitively describes how individuals form edges under resource constraints. ARW uses four parameters— p_{same} , p_{diff} , p_{jump} , p_{out} —to incorporate individuals' biases towards similar, proximate and high degree nodes. Next, we discuss our experiments on the performance of ARW in accurately preserving multiple structural and attribute properties of real networks.

5 MODELING NETWORK STRUCTURE

In this section, we evaluate ARW's performance in preserving real-world network structure relative to well-known growth models.

5.1 Setup

In this subsection, we introduce eight representative growth models and describe evaluation metrics used to fit models to the datasets.

State-of-the-art Growth Models. We compare ARW to eight state-of-the-art growth models representative of the key edge formation mechanisms: preferential attachment, fitness, triangle closing and random walks. Two of the eight models account for attribute homophily and preserve attribute mixing patterns, as listed in Table 2.

Ensuring Fair Comparison. To ensure fair comparison, we modify existing models in three ways. First, for DMS, SAN, KA do not have an explicitly defined initial graph, so we use initialization method used for ARW, described in subsection 4.3. Second, we extend models that use constant node outdegree m by increasing outdegree over time $m(t)$ using the method described in subsection 4.3. In the absence of model-specific parameter estimation methods, we use grid search to estimate the parameters of every network model, including ARW, using evaluation metrics and selection criterion described below.

Evaluation Metrics. We evaluate the network model fit by comparing four structural properties of G & \hat{G} : degree distribution, local clustering distribution, degree-clustering relationship and attribute assortativity. We use Kolmogorov-Smirnov (KS) statistic to compare in-degree & local clustering distributions. We compare the degree-clustering relationship in G and \hat{G} using Weighted Relative Error (WRE), which aggregates the relative error between the average local clustering $c(k)$ and $\hat{c}(k)$ of nodes with in-degree k in G and \hat{G} respectively; The relative error between $c(k)$ and $\hat{c}(k)$ is weighted in proportion to the number of nodes with in-degree k in G .

		Significance level $\alpha < 0.001$ $\alpha < 0.01$																			
		A: INDEGREE DISTRIBUTION (KS STAT)						B: LOCAL CLUSTERING DISTRIBUTION (KS STAT)						C: INDEGREE & CLUSTERING RELATIONSHIP (WRE)							
PREFERENTIAL ATTACHMENT	■	0.03	0.03	0.05	0.09	0.04	0.02	0.80	0.82	0.56	0.63	0.83	0.50	1.00	1.00	1.00	1.00	1.00	1.00	DMS	✗
	■	0.11	0.19	0.22	0.26	0.13	0.06	0.80	0.82	0.56	0.63	0.82	0.50	1.00	1.00	1.00	1.00	1.00	1.00	KA	✓
	■	0.12	0.12	0.17	0.15	0.07	0.15	0.79	0.82	0.56	0.62	0.83	0.50	0.99	1.00	1.00	0.99	1.00	1.00	RL	✗
TRIANGLE CLOSING	■	0.11	0.19	0.22	0.26	0.13	0.05	0.39	0.55	0.15	0.08	0.52	0.05	0.59	0.74	0.08	0.25	0.73	0.17	HK	✗
	■	0.12	0.18	0.19	0.24	0.11	0.05	0.12	0.05	0.12	0.16	0.05	0.19	0.13	0.14	0.34	0.31	0.15	1.28	SAN	✓
RANDOM WALK	■	0.16	0.17	0.14	0.12	0.46	0.32	0.53	0.54	0.33	0.69	0.19	0.40	1.64	1.74	0.54	4.11	0.15	0.73	FF	✗
	■	0.19	0.22	0.25	0.27	0.13	0.13	0.15	0.29	0.26	0.34	0.34	0.11	0.14	0.46	0.74	0.41	0.51	0.38	SK	✗
	■	0.18	0.22	0.23	0.26	0.13	0.13	0.08	0.29	0.10	0.07	0.34	0.03	0.18	0.45	0.21	0.22	0.51	0.04	HZ	✗
	■	0.07	0.06	0.07	0.09	0.07	0.08	0.08	0.04	0.05	0.05	0.05	0.09	0.14	0.10	0.05	0.13	0.08	0.08	ARW	✓
		USSC	HepPH	Semantic	ACL	APS	Patents	USSC	HepPH	Semantic	ACL	APS	Patents	USSC	HepPH	Semantic	ACL	APS	Patents	Assortativity $ r - \bar{r} < \epsilon$	

Figure 5: Modeling network structure. We assess the extent to which network models fit key structural properties of six real-world networks. Tables 5A, 5B and 5C measure the accuracy of eight models in fitting the in-degree distribution, local clustering distribution, in-degree & clustering relationship respectively and global attribute assortativity. Existing models tend to underperform because they either disregard the effect of factors such as triadic closure and/or homophily or are unable to generate networks with varying structural properties. Our model, ARW, jointly preserves all three properties accurately and often performs considerably better than existing models: the cells are shaded gray or dark gray if the proposed model ARW performs better at significance level $\alpha = 0.01$ (■) or $\alpha = 0.001$ (■) respectively.

Jointly preserving multiple structural properties is a multi-objective optimization problem; model parameters that accurately preserve the degree distribution (i.e. low KS statistic) may not preserve the clustering distribution. Therefore, for each model, the selection criterion for the grid search parameter estimation method chooses the model parameters that minimizes the ℓ^2 -norm of the aforementioned evaluation metrics. Since the metrics have different scales, we normalize the metrics before computing the ℓ^2 -norm to prevent unwanted bias towards any particular metric. We note that the parameter sensitivity of the Forest Fire (FF) model necessitates a manually guided grid search method.

5.2 Results

Now, we evaluate the performance of ARW relative to eight well-known existing models on the datasets introduced in Subsection 3.1. Figure 5 tabulates the evaluation metrics for every pair of model and dataset. These metrics measure the accuracy with which the fitted models preserve key global network properties: degree distribution, local clustering distribution, and in-degree-clustering relationship.

To evaluate the performance of these models, we first fit each model to all network datasets G in Subsection 3.1. Thereafter, we compare the structural properties of network dataset G and network \hat{G} generated by the fitted model using evaluation metrics in Subsection 5.1. We average out fluctuations in \hat{G} over 100 runs.

We use one-sided permutation tests [18] to evaluate the relative performance of ARW. If ARW performs better than a model on a dataset with significance level $\alpha = 0.01$ or $\alpha = 0.001$, the corresponding cells in Figure 5 are shaded gray (■) or dark gray (■) respectively. We also group models that have similar edge formation mechanisms by color-coding the corresponding rows in Figure 5. We use green ticks in Figure 5 to annotate models that preserve assortativity up to two decimal places.

Figure 5 shows that existing models fail to jointly preserve multiple structural properties in an accurate manner. This is because existing models either disregard important mechanisms such as

triadic closure and homophily or are not flexible enough to generate networks with varying structural properties.

Preferential attachment models: DMS, RL and KA preserve in-degree distributions but disregard clustering. DMS outperforms other models in accurately modeling degree distribution (Figure 5A) because its “initial attractiveness” parameter can be tuned to adjust preference towards low degree nodes. Unlike KA, however, DMS cannot preserve global assortativity. However, by assuming that successive edge formations are independent, both models disregard triadic closure and local clustering. (Figure 5B & Figure 5C).

Triangle Closing Models: HK and SAN are preferential attachment models that use triangle closing mechanisms to generate scale-free networks with high average local clustering. While triangle closing leads to considerable improvement over DMS and KA in modeling local clustering, HK and SAN are not flexible enough to preserve local clustering in all datasets (see Figure 5B & Figure 5C).

Existing random walk models: FF, SK, and HZ cannot accurately preserve structural properties of real-world network datasets. The recursive approach in FF considerably overestimates local clustering, because nodes perform a probabilistic breadth-first search and link to *all* visited/burned nodes. SK and HZ can control local clustering to some extent, as nodes perform a single random walk and link to each visited node with tunable probability μ . However, both models lack control over the in-degree distribution. Furthermore, existing random walk models disregard attribute homophily and do not account for attribute mixing patterns.

Attributed Random Walk model: Figure 5 clearly indicates the effectiveness of ARW in jointly preserving multiple global network properties. ARW can generate networks with tunable in-degree distribution by adjusting nodes’ bias towards high degree nodes using p_{out} . As a result, ARW accurately preserves in-degree distributions (Figure 5A), often significantly better than all models except DMS. Similarly, ARW matches the local clustering distribution (Figure 5B) and in-degree & clustering relationship (Figure 5C) with high accuracy using p_{jump} and p_{link} . Similarly, ARW preserves attribute

assortativity using the attribute parameters p_{same} and p_{diff} . Barring one to two datasets, ARW preserves all three properties significantly better ($\alpha < 0.001$) than existing random walk models.

To summarize, ARW unifies five sociological phenomena into a single mechanism to jointly preserve real-world network structure.

6 RELATED WORK

Network growth models seek to explain a subset of structural properties observed in real networks. Below, we discuss relevant work on modeling network growth.

Preferential attachment and fitness-based models [4, 5, 10, 30] can preserve heavy-tailed degree distribution, small diameter [8] and temporal dynamics [44] of real-world networks. Furthermore, while extensions of preferential attachment [32, 45, 47], they disregard network properties such as clustering and mixing patterns.

A set of models [20, 23, 26] incorporate triadic closure using triangle closing mechanisms. While this increases average local clustering by forming edges between nodes with one or more common neighbors, as shown in Section 5, it cannot accurately preserve distributional properties of local clustering.

Models [12, 17, 21, 48] that account for attribute mixing patterns can be broadly categorized as (a) fitness-based model that define fitness as a function of attribute similarity and (b) microscopic models of network evolution that require complete temporal information about edge arrivals & deletion. Our experiment results in Subsection 5.2 show that well-known attributed network models SAN and KA cannot jointly preserve mixing patterns and multiple structural properties of real-world networks.

First introduced by Vazquez [43], random walk models are inherently local. Models [7] in which new nodes only link to terminal nodes of short random walks generate networks with power law degree distributions [11] and small diameter [31] but do not preserve clustering. Models such as SK [38] and HZ [19], in which new nodes probabilistically link to each visited nodes incorporate triadic closure but are not flexible enough to preserve skewed local clustering of real-world networks, as shown in Subsection 5.2. Recursive random walk models such as FF [27] preserve temporal properties such as shrinking diameter but considerably overestimate local clustering. Furthermore, existing random walk models do not account for nodal attributes.

To summarize, existing models do not accurately explain how resource constrained and local processes *jointly* preserve multiple global network properties of attributed networks.

7 CONCLUSION

In this paper, we proposed a simple, interpretable model of attributed network growth. ARW grows a directed network in the following manner: an incoming node selects a seed node based on attribute similarity, initiates a biased random walk to explore the network by navigating through neighborhoods of existing nodes, and halts the random walk after connecting to a few visited nodes. To the best of our knowledge, ARW is the first model that unifies multiple sociological phenomena—bounded rationality; structural constraints; triadic closure; attribute homophily; preferential attachment—into a single local process to model global network structure *and* attribute mixing patterns. We explored the parameter space of the model to show how each parameter intuitively controls one or more key structural properties. Our experiments on six large-scale citation networks showed that ARW outperforms relevant and recent existing models by a statistically significant factor of 2.5–10 \times .

We plan to extend the ARW model in three ways: modeling undirected, social networks, understanding the emergence of higher-order clustering [46] and modeling the effect of homophily on the formation of temporal motifs [36]

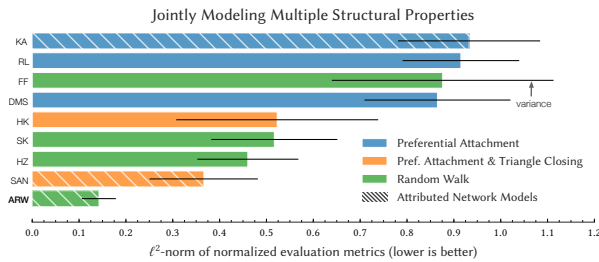


Figure 6: ARW outperforms existing network models in jointly preserving key structural properties—in-degree distribution, local clustering distribution and degree-clustering relationship— by a significant margin of 2.5x-10x.

REFERENCES

- [1] [n. d.]. APS Datasets for Research. ([n. d.]). <https://journals.aps.org/datasets>
- [2] Waleed Ammar, Dirk Groeneveld, Chandra Bhagavatula, Iz Beltagy, Miles Crawford, Doug Downey, Jason Dunkelberger, Ahmed Elgohary, Sergey Feldman, Vu Ha, Rodney Kinney, Sebastian Kohlmeier, Kyle Lo, Tyler Murray, Hsu-Han Ooi, Matthew Peters, Joanna Power, Sam Skjonsberg, Lucy Lu Wang, Chris Wilhelm, Zheng Yuan, Madeleine van Zuylen, and Oren Etzioni. 2018. Construction of the Literature Graph in Semantic Scholar. In *NAACL*.
- [3] Albert-László Barabási and Réka Albert. 1999. Emergence of scaling in random networks. *science* 286, 5439 (1999), 509–512.
- [4] Michael Bell, Supun Perera, Mahendrarajah Piraveenan, Michiel Bliemer, Tanya Latty, and Chris Reid. 2017. Network growth models: A behavioural basis for attachment proportional to fitness. *Scientific Reports* 7 (2017), 42431.
- [5] Ginestra Bianconi and Albert-László Barabási. 2001. Bose-Einstein condensation in complex networks. *Physical review letters* 86, 24 (2001), 5632.
- [6] Per Block and Thomas Grund. 2014. Multidimensional homophily in friendship networks. *Network Science* 2, 2 (2014), 189–212.
- [7] Avrim Blum, TH Hubert Chan, and Mugizi Robert Rwebangira. 2006. A random-surfer web-graph model. In *2006 Proceedings of the Third Workshop on Analytic Algorithmics and Combinatorics (ANALCO)*. SIAM, 238–246.
- [8] Béla Bollobás and Oliver Riordan. 2004. The diameter of a scale-free random graph. *Combinatorica* 24, 1 (2004), 5–34.
- [9] Anna D Broido and Aaron Clauset. 2018. Scale-free networks are rare. *arXiv preprint arXiv:1801.03400* (2018).
- [10] Guido Caldarelli, Andrea Capocci, Paolo De Los Rios, and Miguel A Munoz. 2002. Scale-free networks from varying vertex intrinsic fitness. *Physical review letters* 89, 25 (2002), 258702.
- [11] Prasad Chebolu and Páll Melsted. 2008. PageRank and the random surfer model. In *Proceedings of the nineteenth annual ACM-SIAM symposium on Discrete algorithms*. Society for Industrial and Applied Mathematics, 1010–1018.
- [12] Maurício L de Almeida, Gabriel A Mendes, G Madras Viswanathan, and Luciano R da Silva. 2013. Scale-free homophilic network. *The European Physical Journal B* 86, 2 (2013), 38.
- [13] Sergey N Dorogovtsev, Jose Ferreira F Mendes, and Alexander N Samukhin. 2000. Structure of Growing Networks: Exact Solution of the Barabási–Albert’s Model. *arXiv preprint cond-mat/0004434* (2000).
- [14] James H Fowler and Sangick Jeon. 2008. The authority of Supreme Court precedent. *Social networks* 30, 1 (2008), 16–30.
- [15] Johannes Gehrke, Paul Ginsparg, and Jon Kleinberg. 2003. Overview of the 2003 KDD Cup. *ACM SIGKDD Explorations Newsletter* 5, 2 (2003), 149–151.
- [16] Gerd Gigerenzer and Daniel G Goldstein. 1996. Reasoning the fast and frugal way: models of bounded rationality. *Psychological review* 103, 4 (1996), 650.
- [17] Neil Zhenqiang Gong, Wenchang Xu, Ling Huang, Prateek Mittal, Emil Stefanov, Vyas Sekar, and Dawn Song. 2012. Evolution of social-attribute networks: measurements, modeling, and implications using google+. In *Proceedings of the 2012 Internet Measurement Conference*. ACM, 131–144.
- [18] Phillip Good. 2013. *Permutation tests: a practical guide to resampling methods for testing hypotheses*. Springer Science & Business Media.
- [19] Carlos Herrera and Pedro J Zufiria. 2011. Generating scale-free networks with adjustable clustering coefficient via random walks. In *Network Science Workshop (NSW), 2011 IEEE*. IEEE, 167–172.
- [20] Petter Holme and Beom Jun Kim. 2002. Growing scale-free networks with tunable clustering. *Physical review E* 65, 2 (2002), 026107.
- [21] Fariba Karimi, Mathieu Géniois, Claudia Wagner, Philipp Singer, and Markus Strohmaier. 2017. Visibility of minorities in social networks. *arXiv preprint arXiv:1702.00150* (2017).
- [22] Kibae Kim and Jörn Altmann. 2017. Effect of homophily on network formation. *Communications in Nonlinear Science and Numerical Simulation* 44 (2017), 482–494.
- [23] Konstantin Klemm and Victor M Eguiluz. 2002. Highly clustered scale-free networks. *Physical Review E* 65, 3 (2002), 036123.
- [24] Gueorgi Kossinets and Duncan J. Watts. 2009. Origins of Homophily in an Evolving Social Network. *Amer. J. Sociology* 115 (2009), 405–450. <http://www.journals.uchicago.edu/doi/abs/10.1086/599247>
- [25] Ravi Kumar, Prabhakar Raghavan, Sridhar Rajagopalan, D Sivakumar, Andrew Tomkins, and Eli Upfal. 2000. Stochastic models for the web graph. In *Foundations of Computer Science, 2000. Proceedings. 41st Annual Symposium on*. IEEE, 57–65.
- [26] Jure Leskovec, Lars Backstrom, Ravi Kumar, and Andrew Tomkins. 2008. Microscopic evolution of social networks. In *Proceedings of the 14th ACM SIGKDD international conference on Knowledge discovery and data mining*. ACM, 462–470.
- [27] Jure Leskovec, Jon Kleinberg, and Christos Faloutsos. 2005. Graphs over time: densification laws, shrinking diameters and possible explanations. In *Proceedings of the eleventh ACM SIGKDD international conference on Knowledge discovery in data mining*. ACM, 177–187.
- [28] Barton L Lipman. 1995. Information processing and bounded rationality: a survey. *Canadian Journal of Economics* (1995), 42–67.
- [29] Miller McPherson, Lynn Smith-Lovin, and James M Cook. 2001. Birds of a feather: Homophily in social networks. *Annual review of sociology* 27, 1 (2001), 415–444.
- [30] Matúš Medo, Giulio Cimini, and Stanislao Gualdi. 2011. Temporal effects in the growth of networks. *Physical review letters* 107, 23 (2011), 238701.
- [31] Abbas Mehrabian and Nick Wormald. 2016. It’s a small world for random surfers. *Algorithmica* 76, 2 (2016), 344–380.
- [32] Stefano Mossa, Marc Barthélemy, H Eugene Stanley, and Luis A Nunes Amaral. 2002. Truncation of power law behavior in scale-free network models due to information filtering. *Physical Review Letters* 88, 13 (2002), 138701.
- [33] Mark Newman. 2010. *Networks: an introduction*. Oxford university press.
- [34] Mark EJ Newman. 2001. Clustering and preferential attachment in growing networks. *Physical review E* 64, 2 (2001), 025102.
- [35] Mark EJ Newman. 2002. Assortative mixing in networks. *Physical review letters* 89, 20 (2002), 208701.
- [36] Ashwin Paranjape, Austin R Benson, and Jure Leskovec. 2017. Motifs in temporal networks. In *Proceedings of the Tenth ACM International Conference on Web Search and Data Mining*. ACM, 601–610.
- [37] Dragomir R. Radev, Pradeep Muthukrishnan, Vahed Qazvinian, and Amjad Abu-Jbara. 2013. The ACL anthology network corpus. *Language Resources and Evaluation* (2013), 1–26. <https://doi.org/10.1007/s10579-012-9211-2>
- [38] Jari Saramäki and Kimmo Kaski. 2004. Scale-free networks generated by random walkers. *Physica A: Statistical Mechanics and its Applications* 341 (2004), 80–86.
- [39] Georg Simmel. 1950. *The sociology of georg simmel*. Vol. 92892. Simon and Schuster.
- [40] Herbert A Simon. 1955. On a class of skew distribution functions. *Biometrika* 42, 3/4 (1955), 425–440.
- [41] Herbert A Simon. 1972. Theories of bounded rationality. *Decision and organization* 1, 1 (1972), 161–176.
- [42] Mayank Singh, Rajdeep Sarkar, Pawan Goyal, Animesh Mukherjee, and Soumen Chakrabarti. 2017. Relay-linking models for prominence and obsolescence in evolving networks. In *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. ACM, 1077–1086.
- [43] Alexei Vazquez. 2000. Knowing a network by walking on it: emergence of scaling. *arXiv preprint cond-mat/0006132* (2000).
- [44] Dashun Wang, Chaoming Song, and Albert-László Barabási. 2013. Quantifying long-term scientific impact. *Science* 342, 6154 (2013), 127–132.
- [45] Li-Na Wang, Jin-Li Guo, Han-Xin Yang, and Tao Zhou. 2009. Local preferential attachment model for hierarchical networks. *Physica A: Statistical Mechanics and its Applications* 388, 8 (2009), 1713–1720.
- [46] Hao Yin, Austin R Benson, and Jure Leskovec. 2018. Higher-order clustering in networks. *Physical Review E* 97, 5 (2018), 052306.
- [47] Jianyang Zeng, Wen-Jing Hsu, and Suiping Zhou. 2005. Construction of scale-free networks with partial information. *Lecture notes in computer science* 3595 (2005), 146.
- [48] Elena Zheleva, Hossam Sharara, and Lise Getoor. 2009. Co-evolution of social and affiliation networks. In *Proceedings of the 15th ACM SIGKDD international conference on Knowledge discovery and data mining*. ACM, 1007–1016.