
DAMNETS: A Deep Autoregressive Model for Generating Markovian Network Time Series

Anonymous Author(s)

Anonymous Affiliation

Anonymous Email

Abstract

Generative models for network time series (also known as dynamic graphs) have tremendous potential in fields such as epidemiology, biology and economics, where complex graph-based dynamics are core objects of study. Designing flexible and scalable generative models is a very challenging task due to the high dimensionality of the data, as well as the need to represent temporal dependencies *and* marginal network structure. Here we introduce DAMNETS, a scalable deep generative model for network time series. DAMNETS outperforms competing methods on all of our measures of sample quality, over real and synthetic data sets.

1 Introduction

Temporal networks (also known as dynamic graphs) arise naturally in many fields of study such as the spread of disease [1], molecular interaction networks [2], interbank liability networks [3] and online social [4] and citation networks [5]. Accurate data-driven generating modelling of these processes could have a profound wide-reaching impact, for example in simulating the trajectories of future pandemics or financial contagion risk in economic crash scenarios.

In contrast to generating static networks (i.e., networks that do not evolve over time), generating *time series of networks* has received relatively little attention in the literature. While static networks usually include complex dependencies, network time series contain complex dependencies also across time. As an example, in a time series of social contact networks, the interest may lie in replicating not only the degree distribution but also the clustering behaviour, to capture the interplay between these summary statistics over different times of the day. This complexity is further exacerbated due to the high dimensional nature of network time series; a dataset with N network time series on n nodes each, and of length T each, has size $N \times T \times n^2$. Building a generative model that faithfully replicates both network topology and dependence between graph snapshots is an extremely challenging task.

Data-driven generative models of other types of sequential data, such as natural language, commonly follow an *encoder-decoder* structure, e.g. Sequence2Sequence [6] and Transformer [7] models. We combine ideas from the static network generation and sequence modelling literatures in DAMNETS, an efficient and high quality generator for Markovian network time series. We leverage the insight that the *delta matrix*, that is the difference between subsequent adjacency matrices, is very sparse for most networks of interest. **The key novelty of DAMNETS is that it uses** a GNN to encode the current state of the network, and an efficient sparse matrix sampler to generate delta matrices conditioned on the node embeddings computed by the GNN.

In this paper, we restrict our attention to time series G_0, G_1, \dots, G_T of simple, undirected, labelled graphs on a fixed node set $V = \{1, \dots, n\}$ with edge set $E_t \subseteq \{(i, j) : i, j \in V\}$. An element of the sequence $G_t = (V, E_t)$ has a random edge set E_t drawn from a time-dependent probability distribution $p_t(V \times V)$ over the set of node pairs on V , and emits adjacency matrix $A^{(t)}$.

The remainder of this paper is structured as follows. Section 2 is a review of related work. Section 3 introduces the DAMNETS algorithmic pipeline. Section 4 details the outputs of numerical experiments for representative generative models from the network literature as well as real world networks. Section 5 summarises our main findings and proposes future avenues of investigation. The DAMNETS code is available at [this link](#).

42 **2 Related Work**

43 **2.1 Static Network Generation**

44 Static graph generation involves learning a probability distribution $p(G)$ over an observed set of net-
 45 works. Recently, several machine learning approaches have shown good performance on generating
 46 arbitrary sets of networks, including DeepGMG [8], GraphRNN [9], GRAN [10] and BiGG [11].
 47 Our paper continues this progression to the network time series setting.

48 **BiGG.** BiGG is a scaleable model for generating static networks that we will introduce briefly here,
 49 as our approach shares some similarities. Popular frameworks such as GraphRNN, GRAN and BiGG
 50 all employ the following high-level pattern for sampling the adjacency matrix; they sample each
 51 row of the adjacency matrix one at a time, using a row-wise auto-regressive model to capture the
 52 topological structure of the sampled graph and a second auto-regressive model to capture within-row
 53 edge-level correlations. GraphRNN uses a hierarchical RNN structure, GRAN uses a graph neural
 54 network with a conditional mixture of Bernoulli likelihood and BiGG uses a binary tree type structure,
 55 which is particularly suited to sparse graphs.

56 The major innovation introduced in BiGG is an improvement upon the naive $O(n)$ time complexity
 57 for sampling a row of the adjacency matrix. Instead of sampling each of the n entries using a
 58 linear-time autoregressive model (such as a RNN), the authors propose to sample each row using a
 59 binary tree. Each node u is associated with a random binary tree \mathcal{T}_u which is constructed as follows.
 60 Each tree node k corresponds to an interval of graph nodes $[v_l, v_r]$. The process starts from the root
 61 $[1, n]$ and terminates at leaf nodes $[v, v]$. At each decision step the model decides whether the tree
 62 has a left child (lch), with probability $p(\text{lch}(k))$, and right child (rch), with probability $p(\text{rch}(k))$,
 63 and if so descends further down the tree until it reaches a leaf node. The probability of this tree being
 64 a particular realisation $\mathcal{T}_u = \tau_u$ is thus

$$p(\tau_u) = \prod_{k \in \tau_u} p(\text{lch}(k))p(\text{rch}(k)). \tag{1}$$

65 The tree τ_u is then represented as a row vector of length n of an adjacency matrix, with position
 66 v having entry 1 if τ_u contains the leaf $[v, v]$, and 0 otherwise. The algorithmic advantage stems
 67 from setting all entries $[v_l, \frac{v_l+v_r}{2}]$ to 0 in row u as soon as at tree node $k = [v_l, v_r]$ the left child is not
 68 generated (and similarly if a right child is not generated). Thus for a node u , the corresponding row
 69 of the adjacency matrix can be sampled in $O(|\mathcal{T}_u|)$ decision steps. Since $|\mathcal{N}_u|$, the size of the graph
 70 neighbourhood of u , equals the number of leaf nodes and $\log n$ is the maximum depth of the binary
 71 tree, the upper bound $|\mathcal{T}_u| \leq |\mathcal{N}_u| \log n$ follows. Moreover, significantly larger time savings can be
 72 made in practice if the model decides to not descend further into the tree in the upper levels.

73 To include dependence between entries within the row of the adjacency matrix, BiGG augments the
 74 process to produce state variables that track the decisions made, both above and below in the tree. At
 75 each tree node k , one always decides first whether to generate the left child conditionally on the state
 76 of the tree above, which is denoted $h_u^{\text{top}}(k)$, with the decision sampled from $p(\text{lch}(k)|h_u^{\text{top}}(k))$. If the
 77 model decides to descend into the left child, the entire left subtree is generated before returning to t
 78 and making a decision about whether to generate the right child. The left subtree that was generated
 79 is summarised by a *bottom-up* state variable, denoted $h_u^{\text{bot}}(k)$, and this is used to decide whether to
 80 sample a right child (rch) for the subtree. The model for \mathcal{T}_u therefore becomes

$$p(\mathcal{T}_u) = \prod_{k \in \mathcal{T}_u} p(\text{lch}(k)|h_u^{\text{top}}(k)) p(\text{rch}(k)|h_u^{\text{top}}(k), h_u^{\text{bot}}(\text{lch}(k))), \tag{2}$$

81 where the exact equations for h_u^{top} and h_u^{bot} are given in Algorithm 2. The child probabilities are
 82 finally created via two MLPs, denoted $\text{MLP}_x : \mathbb{R}^F \rightarrow \mathbb{R}$ for $x = L, R$, via

$$p(\text{lch}(k) | h_u^{\text{top}}(k)) = \text{Bernoulli}(\text{MLP}_L(h_u^{\text{top}}(k))), \tag{3}$$

$$p(\text{rch}(k) | h_u^{\text{top}}(k), h_u^{\text{bot}}(\text{lch}(k))) = \text{Bernoulli}(\text{MLP}_R(h_u^{\text{top}}(k), h_u^{\text{bot}}(\text{lch}(k)))). \tag{4}$$

83 **2.2 Network Time Series (NTS) Generation**

84 There are classical models for generating time series of networks designed to capture a specific
 85 set of NTS characteristics, such as the *forest fire* process [5], which can produce power-law degree
 86 distributions and shrinking effective diameter (i.e., the largest shortest path length in the graph). These
 87 classical models, while very effective at re-creating certain types of behaviour, are not data-driven and

88 require the network to obey a pre-defined set of characteristics to be effective. Approaches attempting
 89 to generate arbitrary network time series have appeared in the machine learning literature, such as the
 90 TagGen model [12], which uses a self-attention mechanism to learn from temporal random walks on
 91 a NTS, from which new NTSs are subsequently generated. Another recent algorithm is DYMOND
 92 [13], which is a simpler approach that models the arrival times of 3-node motifs, then samples these
 93 subgraphs to generate the NTS. It is important to note that both DYMOND and TagGen attempt to
 94 solve a slightly different problem to DAMNETS; they take as input a single time series G_0, \dots, G_T and
 95 pre-defined network statistics, and aim to generate an entire time series with these network statistics
 96 similar to this single realisation. Instead of specifying the network statistics of interest, DAMNETS
 97 aims to learn a probability distribution $p(G_t|G_{t-1})$ such that given an arbitrary graph G_{t-1} (not in
 98 the training set), one can draw many samples for G_t and reason about the future trajectory of the
 99 network. This requires a different set of evaluation metrics and datasets, see Section 4 for discussion.

100 **AGE.** The approach most similar to our own is the Attention-Based Graph Evolution (AGE) model
 101 [14]. AGE uses a model very similar to a Transformer [7] (only omitting the positional encoding
 102 step), where a self-attention mechanism is applied to the rows of $A^{(t-1)}$ to learn node embeddings,
 103 and a source target attention module is sequentially applied to generate the rows of $A^{(t)}$. AGE has
 104 two clear shortcomings; the first one is that it does not explicitly account for graph connectivity,
 105 which is left to the attention mechanism to deduce. The second is that it does not capture edge-level
 106 correlations on the sampled rows. To give a simple example of why this is important, suppose we
 107 were considering a NTS where in every graph snapshot, each node has exactly two neighbours; the
 108 model should have some mechanism to condition on the edges it has sampled for a node so that it
 109 can stop once it has generated two edges. Furthermore AGE operates directly between two adjacency
 110 matrices rather than generating only differences, which does not allow it to utilise sparsity, limiting
 111 the scalability of the method. In contrast, DAMNETS explicitly utilises graph connectivity in the
 112 model pipeline and has the capacity to model edge correlations within rows of the adjacency matrix.

113 3 DAMNETS Architecture

114 Our goal is to learn a generative model $p(\cdot|G_{t-1})$ for the next network in a NTS, given a set of
 115 training network time series $\{\{G_t^1\}_{t=0}^{T_1}, \dots, \{G_t^N\}_{t=0}^{T_N}\}$. Our model has a Markovian structure and
 116 hence for generating G_t all relevant information about the past is assumed to be contained in G_{t-1} .

117 For a description of our model we first introduce the *delta matrix* $\Delta^{(t)} \in \{-1, 0, 1\}^{n \times n}$ defined as

$$\Delta_{ij}^{(t)} = A^{(t)} - A^{(t-1)} = \begin{cases} 1 & \implies \text{add edge } (i, j) \\ 0 & \implies \text{no change in } (i, j) \\ -1 & \implies \text{remove edge } (i, j). \end{cases}$$

118 When conditioned on $A^{(t-1)}$, each entry $\Delta_{ij}^{(t)}$ can only take *two values*, namely $\Delta_{ij}^{(t)}$ can only be
 119 0 or 1 if $A_{ij}^{(t-1)} = 0$, and $\Delta_{ij}^{(t)}$ can only be -1 or 0 if $A_{ij}^{(t-1)} = 1$. Learning a generative model
 120 $p(\Delta^{(t)}|G_{t-1})$ is equivalent to learning $p(G_t|G_{t-1})$. Thus, this model only has to learn to produce
 121 the temporal update, rather than to reproduce the current graph *and* apply the temporal update.

122 As we consider only undirected graphs, we only model the lower triangular part of the delta matrix.
 123 As our approach is an encoder-decoder framework, we first summarise the previous network G_{t-1}
 124 by computing node embeddings using a GNN as an encoder, then combine these with a modified
 125 version of the very efficient sparse graph sampler BiGG [11] to act as a decoder for the delta matrix.

126 3.1 The Encoder

127 The first step is to compute node embeddings for G_{t-1} , using a GNN. We employ a Graph Attention
 128 Network (GAT) [15], although any GNN layer is applicable. We use $GAT(X, A)$ to represent
 129 the application of a GAT network to a graph with node feature matrix X and adjacency matrix A .
 130 and in the absence of other node features we use the identity matrix as node features (which here
 131 corresponds to a one-hot encoding of the nodes). Node or edge-level features, whenever available,
 132 can be incorporated into the pipeline. The embedding of G_{t-1} is given by

$$H^{(t-1)} = GAT(X, A^{(t-1)}), \quad (5)$$

133 where $X \in \mathbb{R}^{n \times p}$ is the node feature matrix, and $H^{(t-1)} \in \mathbb{R}^{n \times q}$ is the node embedding matrix.

134 3.2 The Decoder

Starting with the first node according to the given node ordering, conditioning gives

$$p(\Delta) = p(\{\Delta_u\}_{u \in V}) = \prod_{u \in V} p(\Delta_u \mid \{\Delta_w : w < u\}).$$

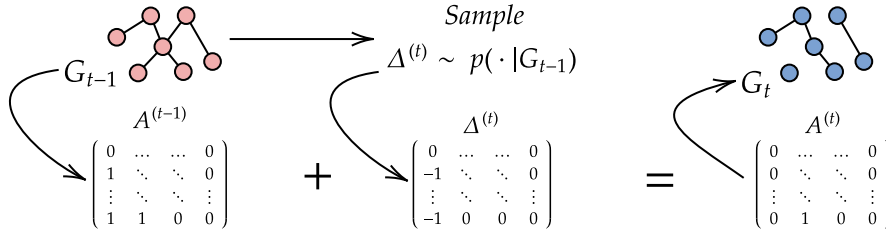
135 We sample each row of Δ using Algorithm 2, a modified version of the BiGG row sampling algorithm.
 136 We enhance the procedure, allowing it to distinguish between a tree leaf which would be an edge
 137 addition and a tree leaf which would be an edge deletion. If the left (resp. right) child at level k is a
 138 leaf node corresponding to entry $\Delta_{ij}^{(t)}$, instead of (3) we sample the leaf node using

$$p(\text{lch}(k) \mid h) = \begin{cases} \text{Bernoulli}(\text{MLP}_+(h)) & \text{if } A_{ij}^{(t)} = 0, \\ \text{Bernoulli}(\text{MLP}_-(h)) & \text{if } A_{ij}^{(t)} = 1, \end{cases} \quad (6)$$

139 where $h \in \mathbb{R}^q$ is the corresponding state variable. Each application of Algorithm 2 returns an
 140 embedding, namely $g_u = h_u^{\text{bot}}(\text{root})$ which depends on every entry in the row. As is done in the
 141 static setting we apply an auto-regressive model across these row embeddings to capture dependencies
 142 between rows. The bottom-up embeddings of each tree have no other computational dependencies, so
 143 can be efficiently pre-computed during training. We chose to use a standard Transformer self-attention
 144 layer [7] (which we call TFEncoder) with sinusoidal positional embedding for this auto-regressive
 145 component; this was chosen to provide similar representation power to the baseline model AGE. Self
 146 attention does not scale to very long sequences however, so for very large graphs with many nodes,
 147 this could be replaced by either an LSTM or the *Fenwick Tree* structure proposed in [16].

148 3.3 The DAMNETS model architecture

Figure 1: An overview of our approach to generating Markovian transitions in a network time series.
 We learn a generative model of the lower triangular part of the *delta matrix* given the previous graph
 G_{t-1} . We then draw a sample $\Delta^{(t)}$ and add this to $A^{(t-1)}$ to produce a sample G_t .



149 With the two key components of our model defined, we now explain how these models are combined
 150 to generate delta matrices given an input graph. As stated in Equation (5), we first compute node
 151 embeddings $H^{(t-1)} \in \mathbb{R}^{n \times F}$, with $H_i^{(t-1)} \in \mathbb{R}^F$ representing the node embedding computed for
 152 node i in G_{t-1} . When generating the row tree \mathcal{T}_u for node u , (which corresponds to generating the
 153 row of the delta matrix for node u), we combine the node embedding from the previous network with
 154 the row-wise auto-regressive term h_{u-1}^{row} computed by TFEncoder via an MLP

$$h_u^{\text{top}}(\text{root}) = \text{MLP}_{\text{cat}}(h_{u-1}^{\text{row}}, H_u^{(t-1)}). \quad (7)$$

155 where $\text{MLP}_{\text{cat}} : \mathbb{R}^{2F} \rightarrow \mathbb{R}^F$. The full procedure is described in Algorithm 1, with a detailed version
 156 Algorithm 2 in the SI, and is visualised in Figures 1 and 2. The model is trained via maximum
 157 likelihood over the entries of the delta matrix using gradient descent. The advantage of this framework
 158 is twofold; firstly the delta matrix is usually much sparser than the full adjacency matrix, allowing us
 159 to well utilise sparse sampling methods. This is a very natural assumption: one does not expect *most*
 160 of the network to change at each timestep, but rather just a small subset of the edges. The second
 161 is that differencing a time series makes learning easier. It is very common in traditional time series
 162 analysis to perform differencing transformations on data, as differencing may alleviate trends in the
 163 time series.

164 4 Experiments

165 Evaluating a generative model usually follows the following recipe: fit the generative model on
 166 the training data, draw samples from the model and then compare the distribution of these samples

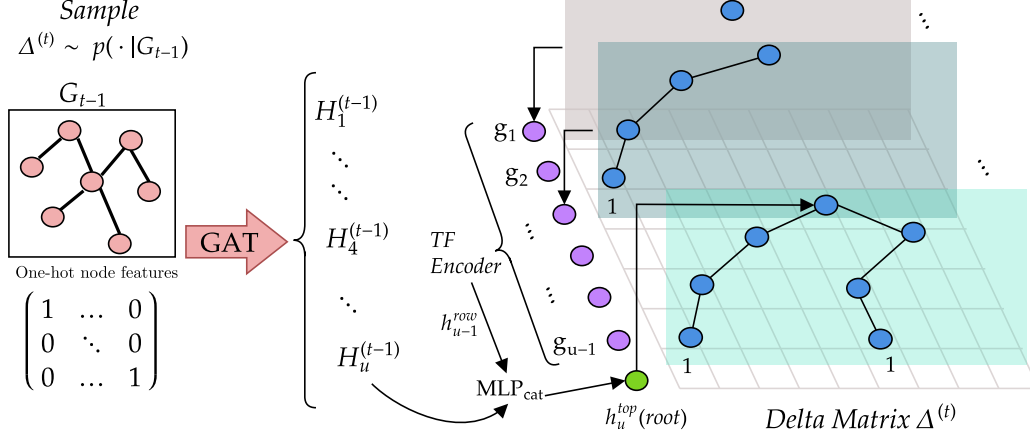


Figure 2: A visualisation of the generation of the u -th row of the delta matrix $\Delta^{(t)}$ using the DAMNETS model architecture. The nodes shown in red indicate the graph G_{t-1} . We use a GAT to compute node embeddings $H^{(t)}$ for each node in G_t . Nodes shown in blue belong to the binary tree generated for each row; each tree is generated by combining the node embedding in the previous graph with an auto-regressive term computed using a Transformer (TF) Encoder across the rows of the delta matrix to produce h_{u-1}^{row} , which is used in Equation (7) to initialise the top-down descent of each tree.

Algorithm 1: Algorithm for generating the the delta matrix $\Delta^{(t)}$ using DAMNETS

Input: Input graph $G_{t-1} = (V, E_{t-1})$, node features X

$H^{(t-1)} \leftarrow GAT(X, A^{(t-1)})$

$h_0^{row} \leftarrow \emptyset$

for $u \leftarrow 1$ **to** n **do**

 Let $k = \{1, \dots, u-1\}$ be the root of tree \mathcal{T}_u .

$h_u^{top}(k) = MLP_{cat}(h_{u-1}^{row}, H_u^{(t-1)})$.

$g_u, \mathcal{N}_u \leftarrow \text{Recursive}(u, k, h_u^{top}(k))$

 /* Algorithm 2 */

 /* Only non-zero indices are returned in \mathcal{N}_u */

$\Delta_u \leftarrow \text{Determine sign of entries using } A^{(t-1)} \text{ and transform into a vector.}$

$h_u^{row} \leftarrow \text{TFEncoder}(g_u; g_{1:u-1})$

end

Return $\Delta^{(t)}$ with rows $\Delta_u, u = 1, \dots, n$.

167 to some held out test data using some kind of statistical test or metric on the space of probability
 168 distributions. For static graphs, there exist a number of *graph kernels* [17] from which a Maximum
 169 Mean Discrepancy (MMD) [18] type metric can be derived. However these are very computational
 170 costly (some scaling as $O(n^4)$ for a graph with n nodes). It is therefore common to define a set of
 171 *summary statistics* over the graphs, such as the degree distribution or clustering coefficient distribution,
 172 and compare the distributions of these summary statistics computed over the sampled and test graphs.

173 We adopt a similar approach applied to the marginal distributions of the network time series. We
 174 choose to compare six different network statistics, three local and three global (see [19] for a
 175 background on network statistics). Our three local properties are the degree distribution, clustering
 176 coefficient distribution and the eigenvalue distribution of the graph Laplacian as introduced in [10].
 177 For each graph, we compute a histogram of these properties over the nodes in the graph, and use a
 178 Gaussian kernel with the total-variation metric to compute the MMD. Our three global measures are
 179 transitivity, assortativity and closeness centrality. Each of these metrics produces one scalar value per
 180 graph, and we again use a Gaussian kernel with the ℓ^2 metric to compute the MMD.

181 For each time point t and statistic $S(\cdot)$, we compute $MMD_t(S(G_t^{test}), S(G_t^{sampled}))$, and use as
 182 final metric the sum $\overline{MMD}(S) = \sum_t MMD_t(S(G_t^{test}), S(G_t^{sampled}))$. If the marginal distributions
 183 match exactly, $\overline{MMD}(S)$ will equal 0, and smaller values indicate better agreement between the
 184 distributions. We display all \overline{MMD} scores to three significant figures. Comparing the marginal

185 distributions alone does not suffice as a comparison metric, so we also provide summary plots of
 186 these network statistics through time to verify that the evolution of these statistics match. In addition,
 187 we have designed several synthetic-data experiments to verify specific time-series properties observed
 188 in real-world networks which we would like to capture.

189 A difficulty for graph generative model evaluation is that proper comparison of a network time
 190 series generator requires many realisations of this time series drawn from the same distribution to
 191 facilitate learning and subsequent comparison. Papers such as TagGen [12] and DYMOND [13]
 192 utilise datasets that comprise of one realisation of a real world temporal network, and aim to simply
 193 produce “*surrogate*” networks that closely resemble that single realisation. We aim to assess whether
 194 our model is able to generalise to new examples, in the sense that given a new graph G_{t-1} drawn from
 195 the same distribution as the training distribution, we can draw samples from $G_t \sim p(\cdot|G_{t-1})$. We are
 196 therefore unable to use the same data sets as these papers, and instead design a new experimental
 197 setup in line with our objective.

198 Our general experimental framework is as follows: we are given a set of realisations
 199 $\{\{G_t^1\}_{t=0}^{T_1}, \dots, \{G_t^N\}_{t=0}^{T_N}\}$. For DAMNETS and AGE, we split this up into a set of training time
 200 series and test time series, and fit each model on the training set, then evaluate the performance
 201 on the test set. As DYMOND and TagGen can only learn from one time series at a time and produce
 202 realisations from that specific time series, we instead train an instance of these models separately
 203 on each time series in the test set and sample one time series from each trained model. This might
 204 seem like a large advantage for these models, as they have direct access to the test set. However
 205 our experimental results show that the aggregated behaviour of these samples does not match the
 206 underlying distribution well, suggesting these methods are not suitable for learning the true underlying
 207 process that a given sample was drawn from. Due to the fact that DYMOND and TagGen have to be
 208 re-trained on every single time series, we provide two sets of results for some datasets, with a smaller
 209 dataset chosen such that DYMOND and TagGen converge within 24 hours.

210 4.1 The Barabási–Albert Model

211 The family of Barabási–Albert (B-A) models [20] was designed to capture the so-called *scale-*
 212 *free* property observed in many real world networks through a preferential attachment mechanism.
 213 Formally a scale-free network is one whose degree distribution follows a power-law; if $\text{deg}(i)$
 214 represents the degree of node i in a random network model, then the network is scale free if
 215 $\mathbb{P}(\text{deg}(i) = d) \propto \frac{1}{d^\gamma}$, for some constant $\gamma \in \mathbb{R}$. Degree distributions with a power-law tail have been
 216 observed in many real networks of interest, such as hyperlinks on the World-Wide Web or metabolic
 217 networks, although the ubiquity of power law degree distributions has been disputed [21].

218 The B-A model has two integer parameters, the number of nodes n and the number of edges m to be
 219 added at each iteration. The network is initialised with m initial connected nodes. At each iteration t ,
 220 a new node is added and is connected to m existing nodes, with probability proportional to the current
 221 degree $p_u = \frac{\text{deg}(u)}{\sum_{v \in V} \text{deg}(v)}$. Here, the standard NetworkX [22] implementation is used. Constructing
 222 a B-A network in this way yields a network time series of length $T = n - m$, where each graph G_t is
 223 the graph after node $m + t$ has the first edges attached to it. Nodes with a many existing connections
 224 (known as *hubs*) will likely accumulate more links; this is the *preferential attachment* property which,
 225 in the B-A model, leads to a power-law degree distribution with scale parameter $\gamma = 3$.

226 For the B-A experiments, we take $N = 200$ time series with parameters $n = 100$ and $m = 4$,
 227 yielding time series of length $T = 96$. The results are displayed in Table 1 and Figure 3. We see
 228 DAMNETS produces samples with orders of magnitude lower $\overline{\text{MMD}}$ than the baseline methods, and is
 229 the only model to correctly replicate the power law degree distribution.

Table 1: The $\overline{\text{MMD}}$ on the B-A dataset for each network statistic. Lower is better.

Model	Degree	Clustering	Spectral	Transitivity	Assortativity	Closeness
DYMOND	14.01	61.20	8.78	7.28	4.76	3.19
TagGen	16.33	16.55	2.29	2.06	23.95	0.10
AGE	15.08	25.15	9.45	3.42	6.37	2.36
DAMNETS	$8e^{-3}$	0.78	0.14	0.01	0.01	$5e^{-6}$

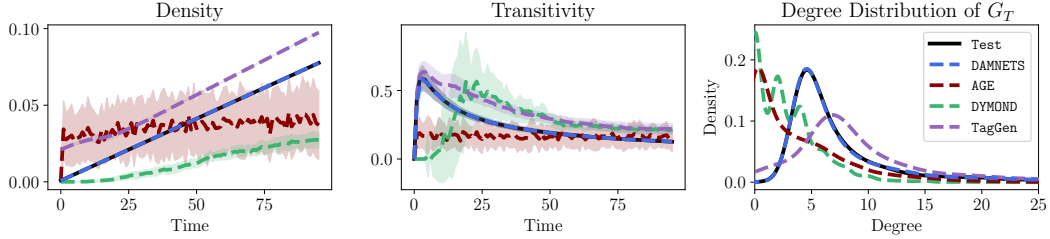
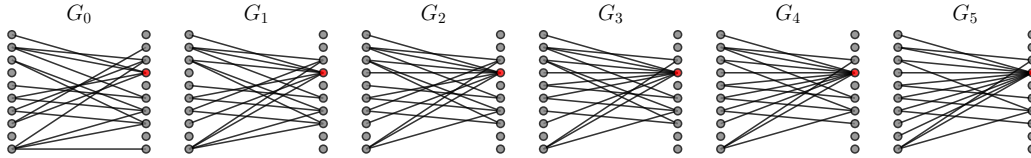


Figure 3: Plots for the B-A model. Left: density against time; middle: transitivity against time; right: the average degree distribution of the final network G_T produced by the models. Only DAMNETS correctly replicates the power law degree distribution.

230 4.2 Bipartite Concentration

Figure 4: A sample from the bipartite concentration model with 10 nodes in each partition, with an initial connection probability of $p = 0.2$ and a concentration proportion $p^{con} = 0.3$. The highest degree node is shown in red; links concentrate on this node over time.



231 This dataset is designed to simulate behaviour in rating systems where objects with many links tend
 232 to accumulate more recommendations [23]. For example in a data set consisting of users and movies,
 233 movies with many existing recommendations are likely to accumulate more over time. The graph
 234 G_0 is initialised as a random bipartite graph with connection probability p . At each timestep, we
 235 select the node in the right-hand partition with the most links (ties broken at random) and re-wire a
 236 proportion p^{con} of non-adjacent edges to that node.

237 For the experiments we set $p = 0.5$ and $p^{con} = 0.1$. For the smaller data set (S), we place 30 nodes
 238 in each partition (so $n = 60$) and iterate for $T = 10$ timesteps. For the larger dataset (L) we place
 239 250 nodes in each partition ($n = 500$) and iterate for $T = 15$ timesteps. To measure the extent to
 240 which the different generators replicate this bipartite structure, in addition to our standard summaries
 241 we also compute the mean Spectral Bipartivity (SB) [24] through time, which takes values in $[0, 1]$,
 242 with 0 indicating the network is not bipartite and 1 indicating the network is fully bipartite. The
 243 results are displayed in Table 2 and Figure 9. DAMNETS consistently outperforms all the baseline
 models across all summary statistics.

Table 2: The $\overline{\text{MMD}}$ for each network statistic (lower is better) and Spectral Bipartivity (closer to 1 is better) across the small (S) and large (L) bipartite contraction test datasets.

Model	Deg.		Clust.		Spec.		Trans.		Assort.		Closeness		SB	
	(S)	(L)	(S)	(L)	(S)	(L)	(S)	(L)	(S)	(L)	(S)	(L)	(S)	(L)
DYMOND	1.06	—	9.55	—	0.12	—	1.67	—	$9e^{-4}$	—	0.14	—	0.50	—
TagGen	0.81	—	1.73	—	0.29	—	$5e^{-4}$	—	0.07	—	$2e^{-4}$	—	0.56	—
AGE	0.92	2.75	9.46	15.3	0.13	0.25	1.48	3.71	0.72	4.81	0.16	0.36	0.55	0.52
DAMNETS	0.01	$4e^{-3}$	0.11	$3e^{-3}$	0.03	$5e^{-4}$	$7e^{-6}$	$8e^{-8}$	$1e^{-4}$	$7e^{-6}$	$4e^{-7}$	$1e^{-7}$	0.99	0.99

244

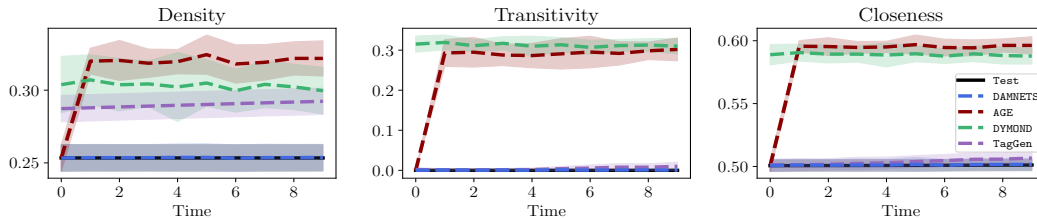
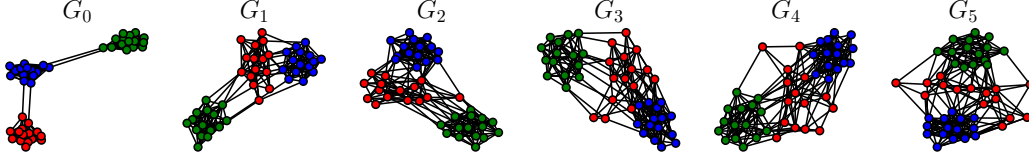


Figure 5: Plots for the bipartite contraction model. Left: density against time; middle: transitivity against time; right: closeness against time. Only DAMNETS shows good performance in all statistics.

245 4.3 Community Evolution and Decay

Figure 6: A sample from the community decay model of length $T = 5$ on $V = \{1, \dots, 45\}$, with 15 nodes in each of the $Q = 3$ communities, connection probabilities $p_{int} = 0.7$, $p_{ext} = 0.005$, decay community $D = 3$ (coloured red) and decay proportion $p_{dec} = 0.2$.



246 Our next network time series benchmark considers a dynamic community structure model. We
 247 initialise a three-community stochastic block model on n nodes. At each time step, we re-wire a fixed
 248 proportion f_{dec} of the third community (which we call the decay community), replacing them with a
 249 random outgoing edge to a node in one of the other communities. A sample from the model is shown
 250 in Figure 6, and a full description of the model is given in Appendix A.2.

251 For the experiments we use inter-community connection probability $p_{int} = 0.9$, intra-community
 252 $p_{ext} = 0.01$, decay fraction $f_{dec} = 0.2$ and iterate for $T = 20$ timesteps. For the small (S) dataset
 253 we place 20 nodes in each community (for a total of $n = 60$ nodes) and for the large (L) dataset we
 254 place 400 nodes in each community ($n = 1200$ in total). The non-decay communities should have
 255 constant density, and the decay community should have density decaying exponentially at rate f_{dec} .
 256 The results are displayed in Table 3 and Figure 10. DAMNETS is the best performing model overall,
 257 although AGE also shows strong performance on this dataset.

Table 3: The $\overline{\text{MMD}}$ for each network statistic across the small (S) and large (L) community decay test datasets, with a (–) when the model did not converge within 24 hours. A lower $\overline{\text{MMD}}$ is better.

Model	Deg.		Clust.		Spec.		Trans.		Assort.		Closeness	
	(S)	(L)	(S)	(L)	(S)	(L)	(S)	(L)	(S)	(L)	(S)	(L)
DYMOND	1.95	–	3.20	–	0.66	–	0.88	–	1.02	–	0.33	–
TagGen	10.99	–	2.91	–	2.18	–	0.26	–	2.37	–	1.04	–
AGE	0.15	0.17	2.00	2.06	0.43	0.42	0.02	0.03	0.07	0.06	0.01	0.03
DAMNETS	0.19	0.21	1.90	1.91	0.39	0.40	0.01	0.01	0.03	0.04	0.01	0.02

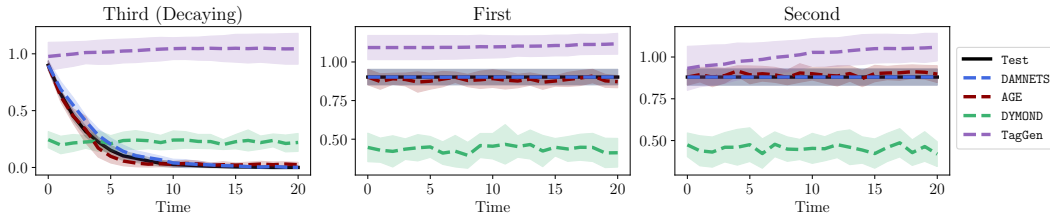


Figure 7: The density of each community through time in the 3-community dataset.

258 4.4 Correlation Networks

259 This data set consists of financial correlation networks built from time series of asset prices from the
 260 Wharton CRSP database [25]. We consider a set of 49 liquid stocks from the US equity market, for
 261 which we have available minutely prices data. We construct a graph by assigning each stock to a
 262 node. We then estimate the correlation matrix of their 5-minute returns each day, and threshold these
 263 correlations at 1 standard deviation in order to construct the edges (so stocks are connected by an
 264 edge if they are strongly correlated). The data set spans $N = 97$ weeks, with each week giving a
 265 time series of length $T = 5$.

266 One issue with this dataset is that correlations between financial instruments are known to be unstable
 267 over time (hence different realisations may not drawn from the same distribution). To mitigate this
 268 we did not split the data chronologically, but have rather drawn the training and test splits randomly
 269 (which correspond to selecting random weekly time series from the dataset). We repeat this procedure
 270 over 5 seeds and compute the average $\overline{\text{MMD}}$. The results are displayed in Table 4 and Figure 11.
 271 DAMNETS is the only model to show good performance across all statistics.

Table 4: The $\overline{\text{MMD}}$ for each network statistic across the correlation test dataset. Lower is better.

Model	Degree	Clustering	Spectral	Transitivity	Assortativity	Closeness
DYMOND	0.16	0.58	0.27	0.17	0.04	0.06
TagGen	0.95	0.56	0.85	$4e^{-3}$	0.08	0.48
AGE	0.14	1.07	0.31	0.26	0.08	0.10
DAMNETS	0.13	0.21	0.25	0.04	0.02	0.01

272 4.5 The MIT Reality Mining Dataset

273 This is a contact network between students and faculty at the MIT media lab recorded between August
 274 2004 to May 2005 [26]. Each contact between two different people (recorded via a bluetooth device
 275 on the subjects’ mobile phones) forms a timestamped edge. We aggregated all daily contacts into
 276 networks, and evaluate our procedure on generating weekly time series of these contact networks.
 277 We dropped weeks without observations for all the days, giving a set of 32 weekly time series. We
 278 used 16 weeks for training, 6 for validation and 10 for testing. As with the correlation dataset,
 279 we randomly sample weeks to form the train and test set, and repeat the experiment across three
 280 seeds. The results in Table 5 show that DAMNETS performs best on all statistics except closeness,
 281 even compared to DYMOND and TagGen which have access to the test data at training. The strong
 282 performance of DAMNETS is particularly evident across the local summary statistics, which suggests
 283 DAMNETS is particularly well suited to represent fine-grained local structure.

Table 5: The $\overline{\text{MMD}}$ for each network statistic across the MIT test dataset. Lower is better.

Model	Degree	Clustering	Spectral	Transitivity	Assortativity	Closeness
DYMOND	1.24	2.21	1.28	0.39	0.18	0.04
TagGen	2.57	2.99	2.42	0.54	0.32	0.48
AGE	2.02	2.75	2.17	0.37	0.38	0.73
DAMNETS	0.41	1.42	0.46	0.34	0.10	0.09

284 4.6 Ablation Study

285 We see that DAMNETS outperforms all the baseline models on each dataset under consideration, in
 286 particular the AGE model, which is the most similar in that it also follows a Sequence2Sequence
 287 framework. DAMNETS differs from AGE in two major ways, namely the formulation in terms of the
 288 delta matrix and the model architecture adapted for sampling this sparse matrix. We provide an
 289 ablation study in Appendix B where we modify AGE to generate delta matrices, and also a version
 290 where we add positional encodings. We find that the delta matrix formulation significantly improves
 291 the performance of AGE, while positional encodings do not change the performance much, with
 292 neither variant of AGE able to match the performance of DAMNETS. This suggests it is the combination
 293 of our re-formulation of the problem combined with a model architecture suited to sample sparse delta
 294 matrices that provides such strong performance. We also provide separate experiments to examine
 295 the influence of the GNN layer type, GNN depth and type of recurrent module in the decoder.

296 5 Discussion and Conclusion

297 DAMNETS provides a novel approach to generating network time series, with the ability to have
 298 fine-grained edge-level conditioning while maintaining scalability by generating delta matrices
 299 rather than entire graphs and efficiently utilising the sparsity of these matrices. We have shown
 300 through extensive experiments that DAMNETS is able to learn a variety of important network models
 301 that existing methods simply cannot. DAMNETS can learn to generate long time series, re-produce
 302 power-law degree distributions, bipartite structure and maintains very strong performance on larger
 303 networks, while none of the baseline models are able to capture all of these properties.

304 In future work, the Markovian assumption underlying DAMNETS could be relaxed to incorporate time
 305 series with long range dependencies, using techniques such as *node memory* introduced in the TGNN
 306 model [27]. The model could also be extended to handle graphs of varying size: node deletion could
 307 be performed by adding a step before the sampling of each row-tree wherein the model makes a
 308 decision about whether the node should persist to the current timestep. Node additions could be
 309 handled by allowing optional rows to be appended to the end of the delta matrix (and only sampling
 310 ones for these rows, as a new node could not have any edge deletions). It would also be interesting
 311 and fairly straightforward to extend DAMNETS to generate node attributes, along the lines of [28].

References

- 312 [1] Naoki Masuda and Petter Holme. *Temporal Network Epidemiology*. 01 2017.
- 313 [2] Teresa Przytycka and Yoo-Ah Kim. Network integration meets network dynamics. *BMC*
314 *biology*, 8:48, 04 2010.
- 315 [3] John Leventides, Kalliopi Loukaki, and Vassilios G. Papavassiliou. Simulating financial conta-
316 gion dynamics in random interbank networks. *Journal of Economic Behavior & Organization*,
317 158:500–525, 2019.
- 318 [4] Jure Leskovec, Lars Backstrom, Ravi Kumar, and Andrew Tomkins. Microscopic evolution
319 of social networks. In *Proceedings of the 14th ACM SIGKDD International Conference on*
320 *Knowledge Discovery and Data Mining*, KDD '08, page 462–470, New York, NY, USA, 2008.
321 Association for Computing Machinery.
- 322 [5] Jure Leskovec, Jon Kleinberg, and Christos Faloutsos. Graphs over time: Densification laws,
323 shrinking diameters and possible explanations. In *Proceedings of the Eleventh ACM SIGKDD*
324 *International Conference on Knowledge Discovery in Data Mining*, KDD '05, page 177–187,
325 New York, NY, USA, 2005. Association for Computing Machinery.
- 326 [6] Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. Neural machine translation by jointly
327 learning to align and translate. *CoRR*, abs/1409.0473, 2015.
- 328 [7] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,
329 Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In I. Guyon, U. V. Luxburg,
330 S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural*
331 *Information Processing Systems*, volume 30. Curran Associates, Inc., 2017.
- 332 [8] Yujia Li, Oriol Vinyals, Chris Dyer, Razvan Pascanu, and Peter Battaglia. Learning deep
333 generative models of graphs. *arXiv preprint arXiv:1803.03324*, 2018.
- 334 [9] Jiaxuan You, Rex Ying, Xiang Ren, William L. Hamilton, and Jure Leskovec. GraphRNN:
335 Generating Realistic Graphs with Deep Auto-regressive Models. *35th International Conference*
336 *on Machine Learning, ICML 2018*, 13:9072–9081, 2 2018.
- 337 [10] Renjie Liao, Yujia Li, Yang Song, Shenlong Wang, Charlie Nash, William L Hamilton, David
338 Duvenaud, Raquel Urtasun, and Richard Zemel. Efficient Graph Generation with Graph
339 Recurrent Attention Networks. In *NeurIPS*, 2019.
- 340 [11] Hanjun Dai, Azade Nazi, Yujia Li, Bo Dai, and Dale Schuurmans. Scalable deep generative
341 modeling for sparse graphs. In *Proceedings of the 37th International Conference on Machine*
342 *Learning*, ICML'20. JMLR.org, 2020.
- 343 [12] Dawei Zhou, Lecheng Zheng, Jiawei Han, and Jingrui He. A data-driven graph generative model
344 for temporal interaction networks. In *Proceedings of the 26th ACM SIGKDD International*
345 *Conference on Knowledge Discovery and Data Mining*, pages 401–411, 2020.
- 346 [13] Giselle Zeno, Timothy La Fond, and Jennifer Neville. DYMOND: DYnamic MOTif-NoDes
347 Network Generative Model. In *Proceedings of the Web Conference 2021 (WWW '21)*, page 12,
348 Ljubljana, Slovenia, 2021. ACM, New York, NY, USA.
- 349 [14] Shuangfei Fan and Bert Huang. Attention-Based Graph Evolution. *Advances in Knowledge*
350 *Discovery and Data Mining*, 12084:436, 2020.
- 351 [15] Petar Veličković, Arantxa Casanova, Pietro Liò, Guillem Cucurull, Adriana Romero, and Yoshua
352 Bengio. Graph attention networks. *6th International Conference on Learning Representations*,
353 *ICLR 2018 - Conference Track Proceedings*, pages 1–12, 2018.
- 354 [16] Hanjun Dai, Azade Nazi, Yujia Li, Bo Dai, and Dale Schuurmans. Scalable Deep Generative
355 Modeling for Sparse Graphs. 2020.
- 356 [17] Giannis Nikolentzos, Giannis Siglidis, and Michalis Vazirgiannis. Graph kernels: A survey. *J.*
357 *Artif. Int. Res.*, 72:943–1027, jan 2022.
- 358 [18] Arthur Gretton, Karsten M. Borgwardt, Malte J. Rasch, Bernhard Schölkopf, and Alexander
359 Smola. A kernel two-sample test. *Journal of Machine Learning Research*, 13(25):723–773,
360 2012.
- 361

- 362 [19] Mark Newman. *Networks: an Introduction*. Oxford University Press, second edition, 2018.
- 363 [20] Reka Albert and Albert-Laszlo Barabasi. Statistical mechanics of complex networks. *Reviews*
364 *of Modern Physics*, 74(1):47–97, 6 2001.
- 365 [21] Aaron Clauset, Cosma Rohilla Shalizi, and M E J Newman. Power-Law Distributions in
366 Empirical Data. *SIAM Review*, 51(4):661–703, 7 2009.
- 367 [22] Aric A. Hagberg, Daniel A. Schult, and Pieter J. Swart. Exploring network structure, dynamics,
368 and function using networkx. In Gaël Varoquaux, Travis Vaught, and Jarrod Millman, editors,
369 *Proceedings of the 7th Python in Science Conference*, pages 11 – 15, Pasadena, CA USA, 2008.
- 370 [23] Massimiliano Zanin, Pedro Cano, Javier M. Buldú, and Oscar Celma. Complex networks in
371 recommendation systems. In *Proceedings of the 2nd WSEAS International Conference on*
372 *Computer Engineering and Applications*, CEA’08, page 120–124, Stevens Point, Wisconsin,
373 USA, 2008. World Scientific and Engineering Academy and Society (WSEAS).
- 374 [24] Ernesto Estrada and Juan A. Rodríguez-Velázquez. Spectral measures of bipartivity in complex
375 networks. *Phys. Rev. E*, 72:046105, Oct 2005.
- 376 [25] The University of Chicago Booth School of Business. Center for Research in Security Prices
377 (CRSP), 2022. Data retrieved from <https://wrds-www.wharton.upenn.edu>.
- 378 [26] Nathan Eagle and Alex (Sandy) Pentland. Reality Mining: Sensing complex social systems.
379 *Personal Ubiquitous Comput.*, 10(4):255–268, 2006.
- 380 [27] Emanuele Rossi, Ben Chamberlain, Fabrizio Frasca, Davide Eynard, Federico Monti, and
381 Michael Bronstein. Temporal graph networks for deep learning on dynamic graphs. In *ICML*
382 *2020 Workshop on Graph Representation Learning*, 2020.
- 383 [28] Ehsan Hajiramezani, Arman Hasanzadeh, Nick G. Duffield, Krishna R. Narayanan, Mingyuan
384 Zhou, and Xiaoning Qian. Variational graph recurrent neural networks. In *NeurIPS*, 2019.
- 385 [29] Kai Sheng Tai, Richard Socher, and Christopher D. Manning. Improved semantic representations
386 from tree-structured long short-term memory networks. In *Proceedings of the 53rd Annual*
387 *Meeting of the Association for Computational Linguistics and the 7th International Joint*
388 *Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 1556–1566,
389 Beijing, China, July 2015. Association for Computational Linguistics.
- 390 [30] Sepp Hochreiter and Jürgen Schmidhuber. Long Short-Term Memory. *Neural Computation*,
391 9(8):1735–1780, 11 1997.
- 392 [31] Marina Knight, Kathryn Leeming, Guy Nason, and Matthew Nunes. Generalized network
393 autoregressive processes and the GNAR package. *Journal of Statistical Software*, 96(5):1–36,
394 2020.
- 395 [32] Daniele Zambon, Daniele Grattarola, Lorenzo Livi, and Cesare Alippi. Autoregressive models
396 for sequences of graphs. pages 1–8, 07 2019.
- 397 [33] Benjamin Paaßen, Daniele Grattarola, Daniele Zambon, Cesare Alippi, and Barbara Hammer.
398 Graph edit networks. In Shakir Mohamed, Katja Hofmann, Alice Oh, Naila Murray, and Ivan
399 Titov, editors, *Proceedings of the Ninth International Conference on Learning Representations*
400 *(ICLR 2021)*, 2021.
- 401 [34] Thomas N. Kipf and Max Welling. Semi-supervised classification with graph convolutional
402 networks. In *International Conference on Learning Representations (ICLR)*, 2017.
- 403 [35] William L. Hamilton, Rex Ying, and Jure Leskovec. Inductive representation learning on large
404 graphs. *Advances in Neural Information Processing Systems*, 2017-Decem:1025–1035, 6 2017.
- 405 [36] Kenta Oono and Taiji Suzuki. Graph neural networks exponentially lose expressive power for
406 node classification. In *International Conference on Learning Representations*, 2020.
- 407 [37] Diederik P. Kingma and Jimmy Lei Ba. Adam: A method for stochastic optimization. In
408 *3rd International Conference on Learning Representations, ICLR 2015 - Conference Track*
409 *Proceedings*, 2015.

- 410 [38] Matthias Fey and Jan E. Lenssen. Fast graph representation learning with PyTorch Geometric.
411 In *ICLR Workshop on Representation Learning on Graphs and Manifolds*, 2019.

- 412 [39] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan,
413 Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas
414 Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy,
415 Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-
416 performance deep learning library. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-
417 Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems 32*,
418 pages 8024–8035. Curran Associates, Inc., 2019.

419 **A Supplementary Information for DAMNETS: A Deep Autoregressive Model**
 420 **for Generating Markovian Network Time Series**

421 **A.1 The DAMNETS Row Generation Algorithm**

Algorithm 2: Algorithm for generating the the u^{th} row of the delta matrix

```

Function Sample_Leaf( $u, k, h$ ):
   $e \leftarrow (u, k)$ 
  if  $e$  is Edge Addition then
    |  $has\_leaf \sim \text{Bernoulli}(\text{MLP}_+(h))$ 
  else
    |  $has\_leaf \sim \text{Bernoulli}(\text{MLP}_-(h))$            /* Edge deletion */
  end
  if  $has\_leaf$  then
    | return  $\vec{1}, e$ 
  else
    | return  $\vec{0}, \emptyset$ 
  end
End Function

Function Recursive( $u, k, h_u^{top}(k)$ ):
  if is_leaf( $lch_u(k)$ ) then
    |  $h_u^{bot}(lch_u(k)), \mathcal{N}_u^{k,left} \leftarrow \text{Sample\_Leaf}(u, lch_u(k), h_u^{top}(k))$ 
  else
    |  $has\_left \sim \text{Bernoulli}(\text{MLP}_L(h_u^{top}(k)))$ 
    if  $has\_left$  then
      |  $h_u^{top}(lch_u(k)) \leftarrow \text{LSTMCell}(h_u^{top}(k), \text{embed}(\text{left}))$ 
      |  $h_u^{bot}(lch_u(k)), \mathcal{N}_u^{k,left} \leftarrow \text{Recursive}(u, lch_u(k), h_u^{top}(lch_u(k)))$ 
    else
      |  $h_u^{bot}(lch_u(k)), \mathcal{N}_u^{k,left} \leftarrow \vec{0}, \emptyset$ 
    end
  end
   $\hat{h}_u^{top}(rch_u(k)) \leftarrow \text{TreeCell}^{top}(h_u^{bot}(lch_u(k)), h_u^{top}(lch_u(k)))$ 
  if is_leaf( $rch_u(k)$ ) then
    |  $h_u^{bot}(rch_u(k)), \mathcal{N}_u^{k,right} \leftarrow \text{Sample\_Leaf}(u, rch_u(k), \hat{h}_u^{top}(rch_u(k)))$ 
  else
    |  $has\_right \sim \text{Bernoulli}(\text{MLP}_L(\hat{h}_u^{top}(rch_u(k))))$ 
    if  $has\_right$  then
      |  $h_u^{top}(rch_u(k)) \leftarrow \text{LSTMCell}(\hat{h}_u^{top}(rch_u(k)), \text{embed}(\text{right}))$ 
      |  $h_u^{bot}(rch_u(k)), \mathcal{N}_u^{k,right} \leftarrow \text{Recursive}(u, rch_u(k), h_u^{top}(rch_u(k)))$ 
    else
      |  $h_u^{bot}(rch_u(k)), \mathcal{N}_u^{k,right} \leftarrow \vec{0}, \emptyset$ 
    end
  end
   $h_u^{bot}(k) \leftarrow \text{TreeCell}^{bot}(h_u^{bot}(lch_u(k)), h_u^{bot}(rch_u(k)))$ 
   $\mathcal{N}_u^k \leftarrow \mathcal{N}_u^{k,left} \cup \mathcal{N}_u^{k,right}$ 
  return  $h_u^{bot}(k), \mathcal{N}_u^k$ 
End Function

```

423 First we provide details for the DAMNETS row generation algorithm given in Algorithm 2. Here,
 424 TreeCell^{bot} and TreeCell^{top} are two TreeLSTM [29] cells, $\text{embed}(\text{left})$ and $\text{embed}(\text{right})$ are learned
 425 embeddings for the binary values "left" and "right", and LSTMCell is a standard LSTM [30]. The top
 426 down cell summarises decisions made above t in the tree, and the bottom up cell summarises lower
 427 levels of the tree (if they exist), where $h_u^{bot}(\emptyset) = 0$. Notice that that h_u^{bot} is computed independently
 428 of h_u^{top} .

429 A.2 The Community Decay Model

430 The three community decay model is formally defined as follows. The initial network $G_0 = (V, E_0)$
 431 with node set $V = \{1, \dots, n\}$ is equipped with a surjective community membership function
 432 $C : \{1, \dots, n\} \rightarrow \{1, \dots, Q\}$ that encodes which of the Q communities a given node i belongs to (a
 433 node can only belong to one community). Here we assume that the community memberships are
 434 known. The initial graph G_0 is then fully described by the *interior* (within community) and *exterior*
 435 (across communities) edge probabilities $p_{ij} := \mathbb{P}((i, j) \in E_0)$, given by

$$p_{ij} = \begin{cases} p_{int} & \text{if } C(i) = C(j) \\ p_{ext} & \text{if } C(i) \neq C(j). \end{cases} \quad (8)$$

436 A network time series G_1, \dots, G_T is then constructed as follows; we fix a community $D \in$
 437 $\{1, \dots, Q\}$ as the *decay* community. We define the set of *internal* edges for community D as

$$D_t^{int} := \{(i, j) \in E_t \mid C(i) = C(j) = D\}. \quad (9)$$

439 At each iteration t , (i.e time step), a fixed proportion f_{dec} of the *internal* edges D_t^{int} are replaced
 440 with *external* edges. This is achieved by selecting a random internal edge (i, j) and removing it from
 441 the edge set E_t , then selecting a node u uniformly from $\{i, j\}$. We then select a random endpoint
 442 k uniformly from $\{v \in V \mid C(v) \neq D, (u, v) \notin E_t\}$, the set of nodes not in community D and not
 443 connected to u , and finally add the edge (u, k) to the edge set E_t . We repeat this procedure T times
 444 to generate our network time series.

445 The model can be interpreted as starting with a network with Q densely connected communities,
 446 decaying in time to have only $Q - 1$ clear communities; the decay community D will appear as
 447 noise around those left unperturbed. A sample from the model can be seen in Figure 6; for ease of
 448 visualisation, each initial community has only 15 nodes.

449 A.3 Graph Attention Networks

450 For the encoder step in DAMNETS we compute node embeddings for G_{t-1} , using a GNN. We employ
 451 a Graph Attention Network (GAT) [15], although any GNN layer is applicable. Given node features
 452 $X_1, \dots, X_n, X_i \in \mathbb{R}^F$, a GAT layer produces a new set of node features $h_i \in \mathbb{R}^{F'}$ according to

$$h_i = \sigma \left(\sum_{j \in \mathcal{N}_i} \alpha_{ij} W X_j \right), \quad (10)$$

453 where $W \in \mathbb{R}^{F' \times F}$ is a learnable weight matrix, $\sigma(\cdot)$ is a non-linear function applied element-wise,
 454 and $\alpha_{ij} \in \mathbb{R}$ are normalised attention coefficients computed as

$$e_{ij} = a(W X_i \parallel W X_j), \quad (11)$$

$$\alpha_{ij} = \frac{\exp(e_{ij})}{\sum_{k \in \mathcal{N}_i} \exp(e_{ik})}, \quad (12)$$

455 where \parallel represents the concatenation operation, and $a(\cdot)$ is a single layer MLP with the LeakyReLU
 456 activation function. These layers are stacked to produce a GAT network. GAT layers can also employ
 457 *multi-head* attention [7]. We write $GAT(X, A)$ to represent the application of a GAT network to a
 458 graph with node feature matrix X and adjacency matrix A .

459 A.4 Further Related Work

460 **Network Time Series Forecasting.** A distinct but related area of study is network time series
 461 forecasting, where the goal is to predict node attributes or links in a graph at a future time point.
 462 Classical approaches include the GNAR model [31] which assumes a simple linear model based on
 463 lagged network attributes. Many approaches have appeared in the deep learning literature, such as the
 464 Graph AR model [32] and Variational Graph Recurrent Neural networks [28]. Markov models have
 465 been considered in this literature before, in particular the Graph Edit Network model [33], which
 466 uses a GNN encoder to predict a list of insertions and deletions for both nodes and edges at the next
 467 timestep, which can be viewed as an estimate of the delta matrix at the next time step. It remains to
 468 note that forecasting focuses on the next time (points) and does not produce a whole time series of
 469 networks which resembles the observed time series.

B Ablation Studies

B.1 The Delta Parameterisation

In many of our experiments, DAMNETS outperforms the AGE model [14]. One may conjecture that it is the use of the delta matrix which is the main driver of this difference in performance. To assess this hypothesis we perform the following ablation study: we re-formulate the AGE model to generate delta matrices instead of entire adjacency matrices. Recall that AGE is a transformer model as described in [7], but without the positional encodings. The transformer is trained via maximum likelihood to generate rows of the adjacency matrix $A^{(t+1)}$ from the rows of the previous adjacency $A^{(t)}$ using the standard Sequence2Sequence framework.

We instead re-formulate AGE to generate delta matrices, which we call AGE-D. To ensure valid delta matrices, we propose to train a transformer to generate $|\Delta^{(t)}|$ row-wise from the rows of $A^{(t)}$, then construct $A^{(t+1)}$ as

$$A^{(t+1)} = \left\{ A^{(t)} + \left| \Delta^{(t)} \right| \right\} \pmod{2}.$$

Note that this always produces a valid adjacency matrix. We also include a variant with positional encodings on both the input and output rows, which we title AGE-DPE. We compare the performance of AGE and the two proposed variants on the BA and the Bipartite Contraction (L) datasets, as there was a particularly large gap in performance between DAMNETS and AGE on these datasets. We also include the $\overline{\text{MMD}}$ for DAMNETS again for ease of reference. The results are displayed in Tables 6 and 7.

Table 6: The $\overline{\text{MMD}}$ for each network statistic on the BA dataset. Lower is better.

Model	Degree	Clustering	Spectral	Transitivity	Assortativity	Closeness
AGE	15.08	25.15	9.45	3.42	6.37	2.36
AGE-D	0.76	2.45	0.69	0.51	4.52	$2e^{-3}$
AGE-DPE	0.76	2.37	0.71	0.49	4.31	$2e^{-3}$
DAMNETS	$8e^{-3}$	0.78	0.14	0.01	0.01	$5e^{-6}$

Table 7: The first block shows the $\overline{\text{MMD}}$ for each network statistic on the Bipartite Contraction (L) dataset, for which lower is better. The last column shows the spectral bipartivity, for which a value closer to 1 is better.

Model	Degree	Clustering	Spectral	Transitivity	Assortativity	Closeness	SB
AGE	2.75	15.3	0.25	3.71	4.81	0.36	0.52
AGE-D	0.15	6.14	0.15	0.04	0.02	$1e^{-2}$	0.85
AGE-DPE	0.13	6.07	0.17	0.03	0.02	$1e^{-2}$	0.87
DAMNETS	$4e^{-3}$	$3e^{-3}$	$5e^{-4}$	$8e^{-8}$	$7e^{-6}$	$1e^{-7}$	0.99

We see that re-formulating AGE to generate delta matrices significantly improves the performance of the method on these datasets, whilst adding the positional encodings provided little to no gain in performance. We also note that while AGE-D performs better than AGE, it still does not match the performance of DAMNETS on these datasets, indicating that it is the combination of our formulation of the problem *and* specific choice of architecture that leads to such strong performance.

B.2 GNN Layer Type

Here we study the impact on the performance of DAMNETS when using different types of GNN layers in the encoder. The three layers we consider are GAT [15] used in the main text, GCN [34] and GraphSAGE [35] with mean aggregation. We repeat the BA experiment from the main text using a single GNN layer of each type, and report the $\overline{\text{MMD}}$ statistic for each variation in Table 8. We see that DAMNETS is not particularly sensitive to the choice of GNN layer.

B.3 GNN Depth

Table 8: The $\overline{\text{MMD}}$ for each network statistic on the BA dataset for samples generated by DAMNETS using different GNN encoder layers. In each case we use a single layer GNN. Lower is better. We see that DAMNETS is not particularly sensitive to the choice of GNN.

Model	Degree	Clustering	Spectral	Transitivity	Assortativity	Closeness
GAT	$8e^{-3}$	0.78	0.14	0.01	0.01	$5e^{-6}$
GCN	$8e^{-3}$	0.81	0.17	0.02	0.01	$7e^{-6}$
GraphSAGE	$7e^{-3}$	0.80	0.16	0.01	0.02	$7e^{-6}$

501 We repeat the BA experiment from the main text varying the depth of the GNN. We use GAT [15]
 502 layers for the encoder with a depth of 1, 2 and 4 layers respectively. The results are displayed in
 503 Table 9. We see that performance slightly degrades with deeper networks, although the difference is
 504 not substantial.

Table 9: The $\overline{\text{MMD}}$ for each network statistic on the BA dataset for samples generated by DAMNETS using different different numbers of GAT layers. Lower is better. We see that the performance is similar across GNN depths, with shallower networks performing slightly better.

Number of Layers	Degree	Clustering	Spectral	Transitivity	Assortativity	Closeness
1	$8e^{-3}$	0.78	0.14	0.01	0.01	$5e^{-6}$
2	$9e^{-3}$	0.77	0.15	0.02	0.01	$7e^{-6}$
4	0.04	0.91	0.22	0.01	0.06	$4e^{-3}$

505 It is important to keep in mind that the GNN embeddings are just one component of the model
 506 pipeline, with the downstream task being generating samples at the next timestep. Therefore intuition
 507 that applies in other applications of GNNs such as deeper networks being better (up to the point of
 508 oversmoothing) may not apply here. In particular DAMNETS has no information other than the present
 509 state of the graph G_t . We hypothesise that if the GNN embeddings computed at time G_t are similar
 510 to those computed at some other time G_{t+k} (which may occur in a deeper network that oversmooths
 511 the node features), then the model may become "confused" and sample the wrong type of transition,
 512 hence shallower networks perform better.

513 B.4 Choice of Encoder

514 In Section 3.2 we mentioned that the self-attention layer in the decoder could be replaced with other
 515 recurrent modules. Here we study the performance of DAMNETS when using an LSTM instead of
 516 self-attention (which as before we call TFEncoder. Again we repeat the BA experiment, using a
 3-layer LSTM or TFEncoder, the results of which are displayed in Table 10.

Table 10: The $\overline{\text{MMD}}$ for each network statistic on the BA dataset for samples generated by DAMNETS using different decoder layers. Lower is better. We used three layers for both the TFEncoder and LSTM. We see that the performance only degrades mildly when moving from the self attention layer through to the LSTM.

Layer Type	Degree	Clustering	Spectral	Transitivity	Assortativity	Closeness
TFEncoder	$8e^{-3}$	0.78	0.14	0.01	0.01	$5e^{-6}$
LSTM	$9e^{-3}$	0.91	0.20	0.02	0.03	$4e^{-5}$

517

518 We see a slight degradation in performance when moving from self-attention to LSTM. This is to
 519 be expected - the self attention layer can directly attend over all the tokens (here node-embeddings)
 520 in it's receptive field, whereas the LSTM only models dependencies via the hidden state. LSTM
 521 layers use significantly less memory however, so for large graphs with many nodes this may be
 522 advantageous. Another alternative is the Fenwick Tree structure introduced in [11].

523 C Experimental Details

524 C.1 Model Specification and Training Details

525 For the experiments in this paper we used a hidden size of $F = 256$ for all experiments, as this was
526 the default hidden size used in BiGG [16]. We used a single layer GAT [15] for all experiments.
527 The rationale for this is as follows: GNNs are known to suffer from an *oversmoothing* problem
528 [36], whereby node embeddings all become similar when using many stacked GNN layers. This
529 would be particularly problematic in our case, as the model would not be able to distinguish between
530 different states in the Markov chain and would likely perform very poorly. We therefore chose to use
531 a very simple model with one GNN layer. It is possible this could be improved upon. All the LSTM
532 networks in the BiGG decoder used 2-layers.

533 We used the Adam [37] optimiser for all experiments, with learning rate 0.001 and weight decay
534 parameter 0.0005. We have not made an effort to optimise these parameters. We used early stopping
535 based on the log-likelihood of the validation set, which was comprised of 30% of the training data,
536 chosen randomly. We used a batch size of 32 graphs (using gradient accumulation for the larger
537 graphs to keep this consistent) and clipped gradients at a norm of 5. We found that training to 0
538 training loss was very harmful for out of sample performance, and that early stopping is necessary for
539 good performance. All numerical results are averaged over five seeds.

540 We implemented the GAT using Torch Geometric [38], and used PyTorch [39] for the other deep
541 learning functionality. We used Networkx [22] for processing the network data. We modified the
542 original BiGG implementation to combine this with the encoder, which can be found at [this link](#).

543 C.2 Baseline Model Information

544 We used the publically released versions of DYMOND and TagGen. Both of these had fatal errors in
545 their implementation, which we have fixed and released as a part of our source code. There is no
546 available code for AGE, so we implemented this using standard PyTorch Transformer modules. We
547 used all the default hyperparameters given in the respective papers. For training AGE we also used
548 early stopping with the same validation log-likelihood criterion, batch size and optimiser settings. As
549 AGE is a Transformer model, we experimented with many "tricks" that are commonly used to train
550 Transformers, such as warmup learning rates as described in [7], but found they did not improve the
551 performance of the model.

552 C.3 Hardware and Running Time

553 All the experiments in this paper were carried out on a single Nvidia GeForce RTX 3090 GPU with
554 an Intel(R) Xeon(R) Silver 4210R CPU @ 2.40GHz. For the smaller experiments we capped the
555 run time of DAMNETS and AGE at one hour, and 24 hours for the larger datasets (although in all cases
556 both these models early stopped before the cap). DYMOND is also fast to run despite needing to be
557 re-trained on each NTS due to its simplicity. TagGen required 24 hours to complete its experimental
558 run on all datasets.

559 D Further Experimental Plots

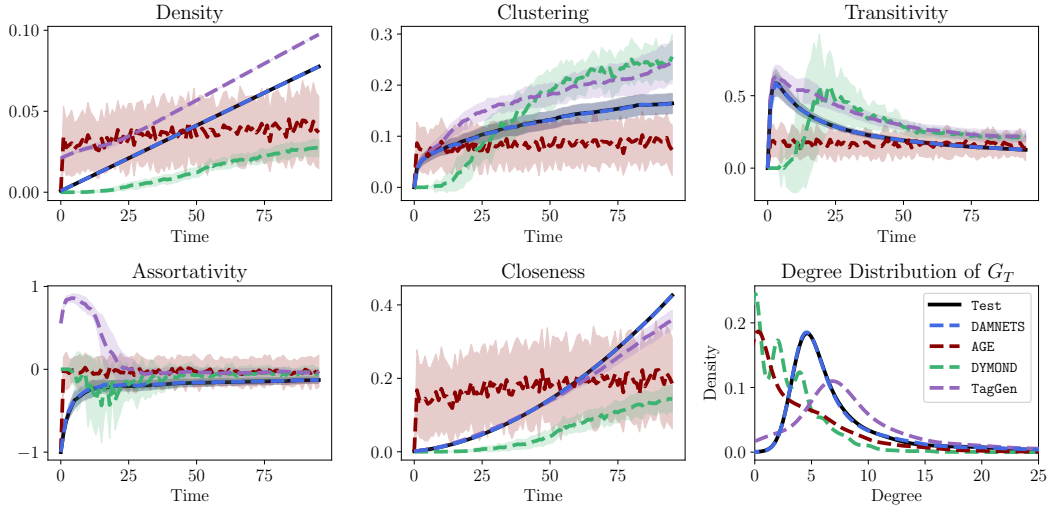


Figure 8: The first five plots show the mean and standard deviation of the network statistics computed through time for the B-A dataset. We see that DAMNETS produces samples that are very similar to the test set across all metrics, whereas the baseline methods fail to do so. The final plot shows the average degree distribution of the final network G_T produced by the models. Only DAMNETS correctly replicates the power law degree distribution.

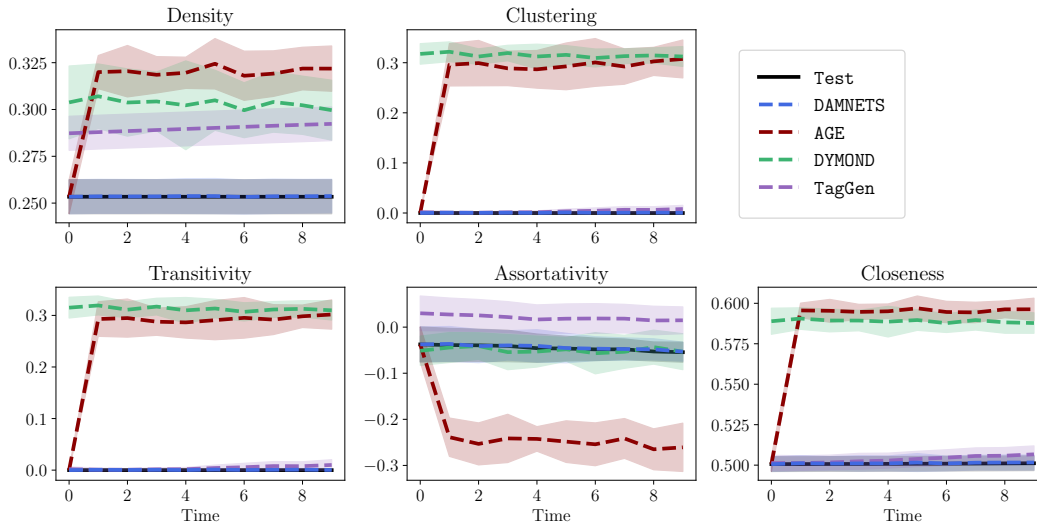


Figure 9: The network statistics computed through time for the bipartite contraction model. We see that DAMNETS shows excellent performance on all statistics, whereas the other models are not able to learn the dynamics.

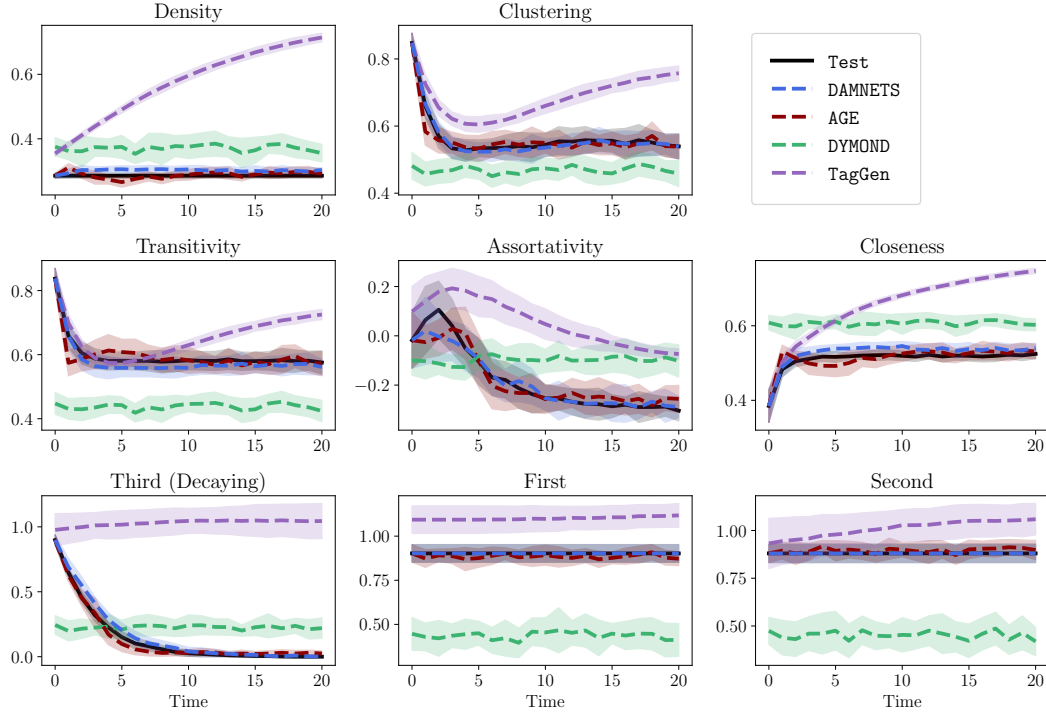


Figure 10: Statistics computed through time on the test set for the three-community decay model. First two rows: The average networks statistics computed across time. Final row: the density of each community through time. We see that both AGE and DAMNETS both show strong performance on this model, while DYMOND performs poorly.

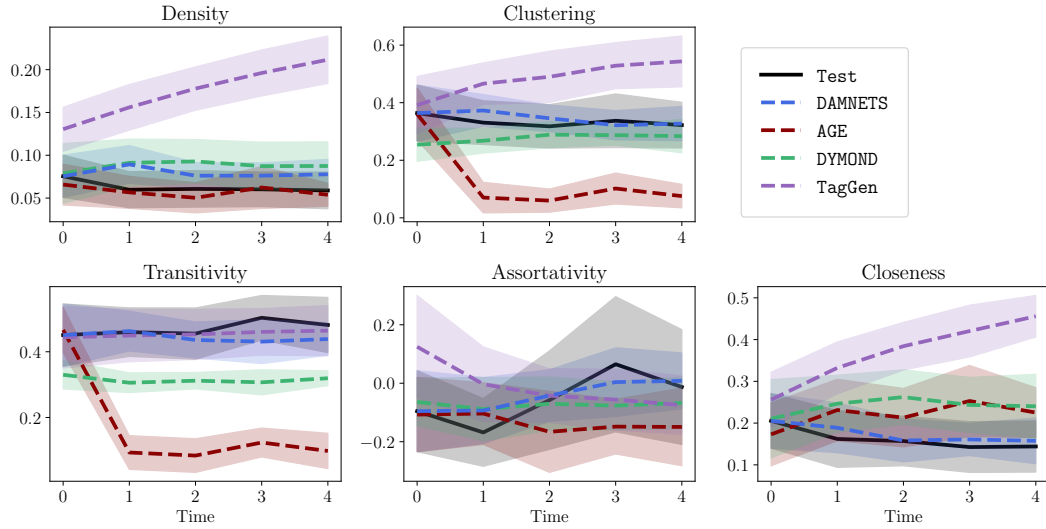


Figure 11: The average network statistics computed through time for the test correlation networks. We see DAMNETS closely tracks the test distribution on all statistics other than density.