

## 712 **A Key Information about GraphGT**

### 713 **A.1 Dataset Documentation**

714 We provide detailed documentation of dataset collection, processing, task for each dataset both in  
715 section C and in our website. We provide statistics, taxonomy, detailed description, and task for each  
716 dataset and can be tracked in our website <https://graphgt.github.io/>.

### 717 **A.2 Intended Use**

718 GraphGT is intended for the deep graph learning as well as specific domain (e.g. physics, biology,  
719 chemistry, etc.) community to use and develop machine learning algorithms to advance applications  
720 in various domains.

### 721 **A.3 URLs**

722 Official website (<https://graphgt.github.io/>) contains all references of GraphGT, including  
723 dataset taxonomy, task, evaluation, visualization, tutorials, papers, GitHub, and other useful resources.  
724 GitHub repository (<https://github.com/yuanqidu/GraphGT>) hosts all source codes, installation  
725 instructions, and tutorials of GraphGT.

### 726 **A.4 Hosting and Maintenance Plan**

727 Our GraphGT Python library is regularly maintained and version-tracked via GitHub. All datasets are  
728 currently hosted on Dropbox and will be transferred to Emory University server soon. Our dataset is  
729 both directly downloadable with a Dropbox link or from our Python APIs. Our core team commit  
730 to maintain this initiative for at least five years. In the meantime, we will expand the community in  
731 multiple dimensions and attract external contributors from the whole community. We will regularly  
732 update new dataset, task, evaluation and visualization methods to GraphGT.

### 733 **A.5 Limitations**

734 Graph generation and transformation is a fast-growing, vast, and promising field and their applications  
735 cover a wide range of applications. We start this initiative to build the infrastructure for the community  
736 which includes most of the mainstream datasets in the graph generation and transformation field and  
737 many more new datasets. However, it is an ongoing effort and we strive to continuously include more  
738 datasets, evaluation and visualization methods to advance the field.

### 739 **A.6 Potential Negative Societal Impacts**

740 Graph generation and transformation are motivated by generating novel graph-structured data and  
741 understanding the graph-structured data; thus, they have vast applications, such as drug discovery,  
742 protein design, mobility synthesis, etc., which could potentially lead to better designed drug, traffic  
743 network, etc., and save lives, time, etc. We envision that GraphGT can facilitate algorithmic  
744 and scientific advances in various domains across subjects and accelerate machine learning model  
745 development and application for real-world use. GraphGT neither involves human subject research  
746 nor contains personally identifiable information.

## 747 **B Dataset Format**

748 We store each of the dataset in a Numpy<sup>3</sup> array format. For different datasets with different information  
749 available as shown in Table 1. For all the datasets, each has at most five types of features available  
750 including adjacency matrices, node features, edge features, spatial features, and labels. Among  
751 all the features, *adjacency matrices* denote the edge connections between pairs of nodes, *node*  
752 *features* denote features attaching to each node, *edge features* denote features attaching to each edge  
753 connection, *spatial features* denote the spatial geometry of a graph (in most of the cases, they are  
754 coordinates attaching to each node), *labels* denote either node-level or graph-level labels of a graph.  
755 For temporal graphs, we store two versions of the graphs, which one flattens and shuffles all the  
756 snapshots of the temporal graphs, and the other one keeps the temporal dimension and order. For  
757 graph transformation datasets, we store both the source and the target graph and available features  
758 separately.

## 759 **C Dataset Details**

760 We list detailed information for each of the datasets in GraphGT.

---

<sup>3</sup><https://numpy.org/doc/>

## 761 C.1 Molecules

762 We have 6 molecule datasets, in which 4 (QM9 [44], ZINC250K [45], MOSES [46], ChEMBL [48])  
763 for graph generation and 2 (MolOpt [47], ChemReact [31]) for graph transformation. For all of the  
764 molecule datasets, we store adjacency matrix, node feature (i.e. atoms), edge feature (i.e. bonds),  
765 spatial feature (i.e. geometry), and smiles (i.e. string representation). There are in total 4 types of  
766 atoms in QM9, 0 = H, 1 = C, 2 = N, 3 = O, 4 = F. There are in total 14 types of atoms in ZINC250K  
767 dataset, MOSES, and ChEMBL dataset, 0 = Br, 1 = C, 2 = Cl, 3 = F, 4 = H, 5 = I, 6 = N, 7 = N, 8 =  
768 N, 9 = O, 10 = O, 11 = S, 12 = S, 13 = S. There are in total 4 types of bonds in all the datasets, and  
769 we represent them as follows: 0 = Single, 1 = Double, 2 = Triple, 3 = Aromatic.

770 **QM9** [44] dataset is an enumeration of around 134,000 stable organic molecules with up to 9 heavy  
771 atoms (carbon, oxygen, nitrogen and fluorine). As no filtering is applied, the molecules in this dataset  
772 only reflect basic structural constraints. In QM9 dataset, each graph contains approximately 9 nodes  
773 and 19 edges. A node in QM9 represents an atom with atom type as the node feature. An edge  
774 in QM9 dataset represents a bond in the molecule with bond type as the edge feature. Moreover,  
775 QM9 dataset contains the 3D spatial feature for each graph. In GraphGT, the QM9 dataset has been  
776 reformatted as adj.npy, edge\_feat.npy, label.npy, node\_feat.npy and spatial.npy that contain molecular  
777 structure information, node features, edge features and spatial features.

778 The information of QM9 is initially stored in .xyz files separately for each molecule. We use Python  
779 to process the SMILE of each molecule and convert the molecule graph to Numpy formats.

780 **ZINC250K** [45] dataset is a curated set of 250k commercially available drug-like chemical com-  
781 pounds. On average, these molecules are bigger (about 23 heavy atoms) and structurally more  
782 complex than the molecules in QM9 dataset. Each graph in ZINC250K dataset contains approxi-  
783 mated 23 nodes and 50 edges. In ZINC250K dataset, each node represents an atom, with atom type  
784 as the node feature. An edge in 250K dataset represents a bond in the molecule with bond type as  
785 the edge feature. 250K dataset also contains 3D spatial feature for each graph. In GraphGT, the  
786 ZINC250K dataset has been reformatted as adj.npy, edge\_feat.npy, label.npy, node\_feat.npy and  
787 spatial.npy that contain molecular structure information, node features, edge features and spatial  
788 features.

789 ZINC dataset is stored in one .csv file including 249,455 molecules. After reading the data by Python,  
790 we process the SMILE of each molecule to convert the data to a graph. And all the graphs are saved  
791 in .npy format.

792 **Molecular Sets (MOSES)** [46] is a benchmark platform for distribution learning based molecule  
793 generation. Within this benchmark, MOSES provides a cleaned dataset of molecules that are ideal of  
794 optimization. It is processed from the ZINC Clean Leads dataset, and contains 193,696 molecules in  
795 total. Each graph in the dataset contains around 22 nodes and 47 edges. In MOSES dataset, each  
796 node represents an atom, with atom type as the node feature. An edge in MOSES dataset represents  
797 the bond in the molecule with bond type as the edge feature. MOSES datasets also contains 3D  
798 spatial features.

799 The data is originally stored in a .txt file. We first read the data and then process the SMILE of the  
800 molecule based on the Python rdkit library. The final data format is saved as .npy files.

801 **ChEMBL** [48] dataset is a manually curated database of bioactive molecules with drug-like properties.  
802 It brings together chemical, bioactivity and genomic data to aid the translation of genomic information  
803 into effective new drugs. ChEMBL contains 1,799,433 graphs in total. Each graph in the dataset  
804 contains around 27 nodes and 58 edges. In ChEMBL dataset, each node represents an atom, with  
805 atom type as the node feature. An edge in ChEMBL dataset represents the bond in the molecule with  
806 bond type as the edge feature. This datasets also contains 3D spatial features.

807 ChEMBL is originally stored in a .txt file containing all the molecules. We first read the data and  
808 then process the SMILE of the molecule based on the Python rdkit library. The final data format is  
809 saved as .npy files.

810 **MolOpt** [47] dataset extracts translation pairs from the ZINC database in terms of three molecular  
811 properties, Penalized logP, Drug-likeness, and Dopamine Receptor. MolOpt contains 229,473 pairs  
812 of graphs in total. Each graph in the dataset contains around 24 nodes and 53 edges. In MolOpt  
813 dataset, each node represents an atom, with atom type as the node feature. An edge in ChEMBL  
814 dataset represents the bond in the molecule with bond type as the edge feature. This datasets also  
815 contains 3D spatial features.

816 This dataset is originally stored in several .csv files and the format of the dataset has been preprocessed.  
817 We read the .csv files and convert the SMILE molecules to graphs and then save them as .npy files.

818 **ChemReact** [31] dataset has totally 7180 pairs of reactant and product molecule graph in the dataset  
819 derived from USPTO dataset. Each graph in the dataset contains around 20 nodes and 16 edges. In  
820 ChemReact dataset, each node represents an atom, with atom type as the node feature. An edge in  
821 ChemReact dataset represents the bond in the molecule with bond type as the edge feature. This  
822 datasets also contains 3D spatial features. [95].

823 Chemical Reaction dataset is originally stored in several .txt files. The first step for processing the  
824 data is to aggregate data from different sources. Then we convert the SMILE of molecules to graph  
825 formats, and then save them in .numpy files.

### 826 C.1.1 License

827 **QM9**: CC BY-NC-SA 4.0.

828 **ZINC250K**: Free to use for everyone.

829 **MOSES**: The dataset is generated by [46], which is under MIT License. The license of the dataset is  
830 not specified.

831 **ChEMBL**: CC BY-NC-SA 3.0.

832 **MolOpt**: Extracted from ZINC Database.

833 **ChemReact**: Not specified.

## 834 C.2 Proteins

835 We have three protein datasets available in GraphGT, which includes protein structures, Enzyme and  
836 dynamic protein folding process.

837 **Protein** [30] dataset contains 918 protein graphs. Each protein is represented by a graph in Protein  
838 dataset, where nodes are amino acids and two nodes are connected if they are less than 6 Angstroms  
839 apart. Proteins dataset contains 1,113 graphs in total. Each graph in the dataset contains around 39  
840 nodes and 73 edges. Node feature is contained in the dataset representing the type of amino acids.  
841 Protein dataset can be used for attributed graph generation.

842 Protein dataset is originally stored in several .txt files with the unit of node. We read all .txt files to  
843 generate graphs, convert them to Numpy arrays and save them in .numpy format.

844 **Enzyme** [28] dataset contains protein tertiary structures representing 600 Enzyme. Nodes in a graph  
845 (protein) represent secondary structure elements, and two nodes are connected if the corresponding  
846 elements are interacting. The node labels indicate the type of secondary structure, which is either  
847 helices, turns, or sheets. Each graph in the dataset contains around 33 nodes and 62 edges. The node  
848 features in the graph represent type of amino acids. This dataset can be employed for attributed graph  
849 generation.

850 Enzyme dataset is originally stored in several .txt files with the unit of node. We read all .txt files to  
851 generate graphs, convert them to Numpy arrays and save them in .numpy format.

852 **ProFold** [29] dataset contains dynamic folding processes of a protein peptide with sequence  
853 AGAAAAGA in 38 steps. ProFold contains 76,000 graphs in total. Each graph has 8 nodes  
854 and around 40 edges. The node represents amino acid of the protein, and the edge represent the bond  
855 between amino acids. The node feature of each protein is the sequence (AGAAAAGA) along with  
856 the spatial locations of each amino acid, and the edge feature of each protein is an adjacency matrix  
857 constructed by connecting all pairs of nodes with distance  $< 8 \text{ \AA}$ . This dataset can be used for either  
858 attributed graph generation or temporal graph generation.

### 859 C.2.1 License

860 **Enzyme**: CC-BY-4.0.

861 **ProFold**: The dataset is collected by [29]. The license is not specified.

862 **Protein**: CC-BY-4.0.

## 863 C.3 Brain Networks

864 The Brain dataset comes from the human connectome project (HCP) [31] and has a few branches:  
865 restingstate, emotion, gambling, language, motor, relational, social and wm according to different  
866 tasks. In this dataset, the source graphs reflect the structural connectivity (SC), and the target graphs  
867 represent the functional connectivity [31]. Specifically, both types of connectivities are processed  
868 from the magnetic resonance imaging (MRI) data from HCP. SC is obtained by applying probabilistic  
869 tracking on the diffusion MRI data by Protrackx tool from the FMRIB Software Library [96] with  
870 68 regions of interest (ROI). The edge attributes of FC are defined as Pearson's correlation between  
871 two ROIs blood oxygen level-dependent time obtained from the resting-state functional MRI data.

872 Node attributes is a one-hot vector representing index of each node. In total, 823 pairs of SC and  
873 FC samples are enrolled in the dataset. The dataset has been splitted into 8 categories for 8 specific  
874 domains, including Brain-restingstate, Brain-emotion, Brain-gambling, Brain-language, Brain-motor,  
875 Brain-relational, Brain-social and Brain-wm. All of these datasets can be employed for eight weighted  
876 graph transformation or signed graph transformation tasks.

877 Originally, data is a group of .npz files, containing the structural connectivities for each subject,  
878 functional connectivities for each subject, and list of subject IDs for each task using different  
879 correlations. Unfortunately, the subjects used are not universal for all tasks, and so we eliminate  
880 all but those that appeared in every single task. From there, we simply concatenate all of the  
881 functional connectivities from all of the various tasks using FC correlation, and concatenated all  
882 of the structural connectivities from all of the various tasks using FC correlation, thus creating  
883 FC\_concatenated\_edge\_feat and SC\_concatenated\_edge\_feat. For the adjacency matrix containing  
884 .npy arrays, we encounter a small issue; the adjacency matrix is required to be formatted with a  
885 specific shape, but that shape is not compatible with the edge feature shape, and so we make the  
886 adjacency matrix a placeholder basically. For details please refer to readme.txt.

### 887 C.3.1 License

888 **Brain:** This dataset comes from the human connectome project. Data collection and sharing for this  
889 project was provided by the MGH-USC Human Connectome Project (HCP; Principal Investigators:  
890 Bruce Rosen, M.D., Ph.D., Arthur W. Toga, Ph.D., Van J. Weeden, MD). HCP funding was provided  
891 by the National Institute of Dental and Craniofacial Research (NIDCR), the National Institute of  
892 Mental Health (NIMH), and the National Institute of Neurological Disorders and Stroke (NINDS).  
893 HCP data are disseminated by the Laboratory of Neuro Imaging at the University of Southern  
894 California.

### 895 C.4 Physical Simulations

896 **N-body-charged** [49] dataset simulates a system containing 5 particles with positive or negative  
897 charges. Particles are located in 2D coordinates without any external forces except attracting force  
898 and repelling force. The quantity of electrical charges is sampled from uniform probability. Each  
899 particle interacts via Coulomb forces. Every two particles interact, either attract or repel each other.  
900 The temporal length of each sequence is 49, which obtains from sub-sampling every 100 steps in  
901 a trajectory. N-body-charged dataset contains 3,430,000 graphs in total, each of which contains 25  
902 nodes with around 3 edges. Each node represents a particle and each edge represents interaction  
903 between nodes. Node attribute represents node input. 2d spatial features and temporal are included  
904 in the dataset. N-body-charged can be used for either attributed graph transformation, spatial graph  
905 transformation or temporal graph transformation.

906 Originally, for the charged dataset, there are separate numpy files for the velocities, edges, and  
907 locations of each particle for train, validation, and testing. Then, all velocity arrays(train, valid, test)  
908 for the charged dataset were merged into a single one, and the same was done for all of the location  
909 arrays, and all of the edge arrays. To convert the charged edge features into adjacency matrices, all  
910 nonzero values were turned to ones, and since all particles had some form of connection, that meant  
911 all adjacency matrices ended up being all ones for the charged dataset. Then, for each new temporal  
912 array we had here, we created a new version: a non-temporal one, where we concatenated the first  
913 two dimensions of the array, as the second dimension represented the different temporal instances.  
914 For details information, please refer to readme.txt.

915 **N-body-spring** [49] dataset simulates a system containing 5 particles connected by springs. Particles  
916 are located in 2D coordinates without any external forces except elastic collisions. Particles are  
917 connected via springs with probability of 0.5, and interactions between springs follow Hooke's law.  
918 The initial location of each particle is sampled from a Gaussian distribution and the initial velocity of  
919 each particle is a random vector of norm 0.5. The trajectories of all springs are calculated by solving  
920 Newton's equations of motion PDE. The temporal length of each sequence is 49, which obtains from  
921 sub-sampling every 100 steps in a trajectory. N-body-spring dataset contains 3,430,000 graphs in  
922 total, each of which contains 5 nodes with around 10 edges. Each node represents a particle and each  
923 edge represents interaction between nodes. Node attribute represents node input. 2D spatial features  
924 and temporal features are included in the dataset. N-body-spring can be used for either attributed  
925 graph transformation, spatial graph transformation or temporal graph transformation.

926 Originally, for the spring dataset, there were separate numpy files for the velocities, edges, and  
927 locations of each particle for train, validation, and testing. There are 5 particles, 5 springs in each  
928 graph. Then, all velocity arrays(train, valid, test) for the spring dataset were merged into a single one,

929 and the same was done for all of the location arrays, and all of the edge arrays. For the springs dataset,  
930 we had only ones and zeroes in the edges: connection or no connection, and so we simply took this  
931 as our adjacency matrix as well for each matrix in the springs dataset. Then, for each new temporal  
932 array we had here, we created a new version: a non-temporal one, where we concatenated the first  
933 two dimensions of the array, as the second dimension represented the different temporal instances.

#### 934 C.4.1 License

935 **N-body-charged:** The dataset is simulated by [49], which is under MIT License. The license of the  
936 dataset is not specified.

937 **N-body-spring:** The dataset is simulated by [49], which is under MIT License. The license of the  
938 dataset is not specified.

### 939 C.5 Collaboration Networks

940 **CollabNet** [55] dataset is collected from DBLP-Citation-network V12, which contains around 4.9  
941 million papers and 45 million citation relationships. We construct graphs by selecting authors as  
942 nodes and co-authorships as edges during the time period from 1990 to 2019. To cut the graphs into  
943 pieces, we generate sub-graphs based on the Fields of Study attribute from papers. For each field, we  
944 generate one spatio-temporal graph. We generate 2361 spatio-tempora graphs with a total of 303,308  
945 nodes and a total of 207,632 of edges. This dataset contains temporal and GCS spatial features, so  
946 that the dataset can be used for spatial graph generation and temporal graph generation.

#### 947 C.5.1 License

948 **CollabNet:** The dataset is collected from DBLP-Citation-network V12. The license is not specified.

### 949 C.6 Traffic Networks

950 **METR-LA** [53] dataset is collected by Los Angeles Metropolitan Transportation Authority (LA-  
951 Metro), and processed by University of Southern California’s Integrated Media Systems Center. This  
952 dataset contains traffic information collected from 207 loop detectors in the highway of Los Angeles  
953 County for 4 months (from Mar 1st 2012 to Jun 30th 2012). Each sensor records traffic speed value  
954 per 5 minutes. The dataset contains 34,272 graphs, each of which has 325 nodes and 2,369 edges. In  
955 METR-LA, each node represent a speed sensor and each edge represents a road. The node features  
956 of the dataset represent the traffic speed captured by the sensor. The dataset contains GCS spatial  
957 features and temporal features. METR-LA can be used for spatial graph generation, temporal graph  
958 generation, attributed graph generation and weighted graph generation.

959 The information of the METR-LA dataset is stored in three files with different formats. We borrow  
960 Python to read these data, and convert them to Numpy formats. We then save the data in .npy format.

961 **PeMS-BAY** [54] dataset is collected by California Transportation Agencies (CalTrans) Performance  
962 Measurement System (PeMS). PeMS-BAY dataset collects traffic information in the Bay Area. The  
963 dataset contains traffic information of 325 sensors within 5 months (From Jan 1st 2017 to May 31st  
964 2017). Each sensor records traffic speed value per 5 minutes. The dataset contains 50,221 graphs,  
965 each of which has 207 nodes and 1,515 edges. In PeMS-BAY, each node represent a speed sensor and  
966 each edge represents a road. The node features of the dataset represent the traffic speed captured by  
967 the sensor. The dataset contains GCS spatial features and temporal features. PeMS-BAY can be used  
968 for spatial graph generation, temporal graph generation, attributed graph generation and weighted  
969 graph generation.

970 The information of the PeMS-BAY dataset is stored in three files with different formats. We borrow  
971 Python to read these data, and convert them to Numpy formats. We then save the data in .npy format.

#### 972 C.6.1 License

973 **METR-LA:** The dataset is collected by Los Angeles Metropolitan Transportation Authority (LA-  
974 Metro), and processed by University of Southern California’s Integrated Media Systems Center. The  
975 license is not specified.

976 **PeMS-BAY:** The dataset is collected by California Transportation Agencies (CalTrans) Performance  
977 Measurement System (PeMS). The license is not specified.

### 978 C.7 Authentication Networks

979 **AuthNet** dataset includes the authentication activities of users on their computers and servers in their  
980 enterprise computer network and is published by Los Alamos National Laboratory (LANL). [97, 41].  
981 There are two subsets of different sizes of graphs (e.g., 50 and 300) in AuthNet dataset with 114 and  
982 412 graphs, respectively. For each subset, we train and test folder separately. Train set contains the



983 graph pairs (one-to-one) which are just used for training. Test set contains data for each user. For  
984 each user, there are several input graphs (e.g., regular user authentication activity graph) and several  
985 target graphs (e.g., malware user authentication activity graph). Input and target graphs in test set  
986 are not one-to-one, which can be tested by indirect evaluation. There are no node attributes for this  
987 dataset, and only edge attribute is considered. For each graph, the value of the  $i$ -th row and the  
988  $j$ -th column refers to the edge attribute of node  $i$  and  $j$  (0 refers to no links). This dataset can be  
989 employed for weighted graph generation.

### 990 C.7.1 License

991 **AuthNet**: The dataset is publically released by LANL [97]. To the extent possible under law,  
992 LANL has waived all copyright and related or neighboring rights to User-Computer Authentication  
993 Associations in Time. This work is published from: United States.

994 We collect this dataset from DBLP-Citation-network V12. We chose authors with affiliations, papers  
995 with more than one authors, and the time period from 1990 to 2019. To cut the graphs into pieces,  
996 we generate sub-graphs based on the fields of study of papers. For each field, we generate one  
997 spatio-temporal graph. Then we concatenate and pad all graphs, and save them into Numpy arrays.  
998 We save the graphs in .npy format.

## 999 C.8 IoT Networks

1000 **IoTNet** is the malware dataset collected for malware confinement prediction [31]. There are three  
1001 sets of IoT nodes at different amounts (20, 40 and 60) encompassing temperature sensors connected  
1002 with Intel ATLASEDGE Board and Beagle Boards (BeagleBone Blue), communicating via Bluetooth  
1003 protocol. Benign and malware activities are executed on these devices to generate the initial attacked  
1004 networks (i.e., the Internet of Things) as input graphs. Benign activities include MiBench [98] and  
1005 SPEC2006 [99], Linux system programs, and word processors. The nodes represent devices and node  
1006 attribute is a binary value referring to whether the device is compromised or not. Edge represents the  
1007 connection of two devices and the edge attribute is a continuous value reflecting the distance of two  
1008 devices. The real target graphs are generated by the classical malware confinement method: stochastic  
1009 controlling with malware detection [100, 101, 102]. We collect 334 pairs of input and target graphs  
1010 with different contextual parameters (infection rate, recovery rate and decay rate) for each of the three  
1011 datasets. In this dataset, there are both nodes attributes and edge attributes considered. IoTNet can be  
1012 used for attributed graph generation and weighted graph generation.

1013 The original format of IoTNet contains 1,029 .csv files, we convert them to .npy files, input\_adj.npy,  
1014 input\_edge.npy, input\_node.npy, target\_adj.npy, target\_edge.npy, target\_node.npy, Iot\_20\_labels.npy,  
1015 Iot\_40\_labels.npy and Iot\_60\_labels.npy, to contain structure, node features, edge features and  
1016 labels and to be easily read by Python. The detailed information of the data can be found in the  
1017 corresponding readme.txt file. To reformat the data, we use glob to read in all .csv files from the  
1018 directory, and separate the original .csv files into input data and target data; For both input and target  
1019 data, we get edge feature from the original .csv files, get node feature(0 or 1 for IoTNet) from the  
1020 diagonals of each file, and get adjacent matrix from the edge feature while setting the diagonals to  
1021 be 0. For IoTNet, we also split the name and get labels from the name of each .csv file. We then  
1022 reshaped all arrays into the required dimensions and converted them to NumPy files.

### 1023 C.8.1 License

1024 **IoTNet**: The dataset is generated by [31]. The license is not specified.

## 1025 C.9 Skeleton Graphs

1026 **Kinetics** [51] dataset is a large-scale human action dataset with 300000 videos clips in 400 classes.  
1027 Those video clips are from YouTube with a great variety. The raw Kinetics dataset doesn't contain  
1028 skeleton data, and [51] uses OpenPose toolbox to generate skeleton with 18 joints on every frame.  
1029 Kinetics-Skeleton contains 240000 clips of training data and 20000 clips of test data. This dataset  
1030 does not contain node or edge attributes, but contain temporal and 2D spatial features to be used in  
1031 spatial graph generation and temporal graph generation tasks.

1032 The raw Kinetics dataset is stored in a few .json files, and each json file contains information of a  
1033 single video clip. We traverse all .json files, and concatenate their contents into several Numpy arrays  
1034 with paddings for short video clips. We then remove extra skeletons, and leave each video clip only  
1035 one skeleton. Finally, we save the data in .npy array.

1036 **NTU-RGB+D** [52] dataset is a large and widely used action recognition dataset with 56000 action  
1037 clips in 60 classes. These clips are performed by 40 volunteers captured in a constrained lab  
1038 environment, with three camera views recorded simultaneously. The dataset provides 3D joint

1039 locations of each frame and 25 joints for each subject. NTU-RGB+D does not contain node or edge  
1040 attributes, but contain temporal and 3D spatial features to be used in 3D spatial graph generation and  
1041 temporal graph generation tasks.

1042 We process this dataset by the code from github. The dataset is originally stored in a few files, and  
1043 each contains information of one single video clip. After the same processing process as we do for  
1044 Kinetics dataset, we save the data in .numpy format.

#### 1045 **C.9.1 License**

1046 **Skeleton (Kinectics):** CC BY 4.0.

1047 **Skeleton (NTU-RGB+D):** Not specified.

#### 1048 **C.10 Social Networks**

1049 **Ego:** Ego dataset contains 757 3-hop ego networks extracted from the Citeseer [103]. The number of  
1050 nodes of the graph in Ego dataset ranges from 50 to 399, and 145 in average. Each graph in Edo has  
1051 around 335 edges. Nodes represent documents and edges represent citation relationships [34]. Ego  
1052 does not contain node or edge attributes, and can be used for graph generation tasks.

1053 **TwitterNet:** The dataset is processed by [56] and obtained from 5 different countries in Latin  
1054 America, namely Brazil, Colombia, Mexico, Paraguay, and Venezuela. Data sources from Twitter are  
1055 adopted as the model inputs. In each case the data for the period from July 1, 2013 to February 9,  
1056 2014 is used for training and validation, where the validation set consists of a randomly chosen 30%  
1057 of the data, and the rest is used for training; the data from February 10, 2014 to December 31, 2014 is  
1058 used for the performance evaluation. TwitterNet contains 2,580 graphs in total, each of which has  
1059 300 nodes and 0.5 edges in average. This dataset can be employed in graph transformation tasks.

#### 1060 **C.10.1 License**

1061 **Ego:** This dataset is extracted from Citeseer [103]. Citeseer is under CC BY-NC-SA 3.0.

1062 **TwitterNet:** The dataset is obtained from [104]. The license is not specified.

#### 1063 **C.11 Scene Graphs**

1064 **CLEVR** [50] dataset provides a dataset for visual question answer, which can be formalized as a  
1065 spatial-graph dataset. CLEVR dataset contains 85,000 graphs in total. There are 10 objects in the  
1066 image with different 3D locations. Each object is identified by its shape, such as sphere, cylinder, and  
1067 cube. The relationship between two objects can be categorized into four types: right, behind, front,  
1068 left, with directions. Thus, each image can be formalized as a labeled directed graph with different  
1069 edge types and node types. Thus, the spatial information of each nodes is closely correlated with the  
1070 edge types between each pair of nodes. As a result, CLEVR dataset can be employed for attributed  
1071 graph generation, weighted graph generation and spatial graph generation.

#### 1072 **C.11.1 License**

1073 **CLEVR:** CC BY 4.0.

#### 1074 **C.12 Synthetic Graphs**

1075 **Barab’asi-Albert Graphs:** This dataset is generated by the Barab’asi-Albert model [31]. It fits the  
1076 "one-to-one" mapping problem of graph translation. It contains pairs of input and target graphs. The  
1077 target graph topology is the 2-hop connection of the input graph, where each edge in the target graph  
1078 refers to the 3-hop reachability in the input graph (e.g., if node  $i$  is 3-hop reachable to node  $j$  in the  
1079 input graph, then they are connected in the target graph). There are edge and node attributes for graphs  
1080 in this dataset: the edge attribute  $E_{(i,j)}$  denotes the existence of the edge, and the node attributes  
1081 are continuous values computed following the polynomial function:  $f(x) : y = ax^2 + bx + c$   
1082 ( $a = 0; b = 1; c = 5$ ), where  $x$  is the node degree and  $f(x)$  is the node attribute. Here we provide the  
1083 datasets with three different node sizes. Barab’asi-Albert Graphs dataset can be used for attributed  
1084 graph transformation.

1085 The original Barab’asi-Albert Graphs dataset contains 3,000 .csv files. We reformat them into .numpy  
1086 files, including input\_adj.npy, input\_edge.npy, input\_node.npy, target\_adj.npy and target\_edge.npy,  
1087 target\_node.npy for the community to use. To reformat the data, we use glob to read in all .csv files  
1088 from the directory, and separate the original .csv files into input data and target data; For both input  
1089 and target data, we get edge feature from the original .csv files, get node feature from the diagonals  
1090 of each file, and get adjacent matrix from the edge feature while setting the diagonals to be 0. We  
1091 then reshaped all arrays into the required dimensions and converted them to NumPy files.

1092 **Community:** This dataset is generate by [34] and contains 3,000 two-community graphs, each of  
1093 which has 64 nodes and around 340 edges. Each community is generated by the Erdos-Renyi model  
1094 (E-R) [105] with  $\frac{|V|}{2}$  nodes and the edge probability of 0.3. Then add  $0.05|V|$  inter-community edges  
1095 are added with uniform probability. This dataset does not have node or edge attributes. Community  
1096 can be used for graph generation tasks.

1097 **Erdos-Renyi Graphs:** This dataset is generated by the Erdos-Renyi model with the edge probability  
1098 of 0.2 [31]. It fits the "one-to-one" mapping problem of graph translation. It contains pairs of (input,  
1099 target) graphs. The target graph topology is the 2-hop connection of the input graph, where each  
1100 edge in the target graph refers to the 2-hop reachability in the input graph (e.g., if node  $i$  is 2-hop  
1101 reachable to node  $j$  in the input graph, then they are connected in the target graph). There are  
1102 edge and node attributes for graphs in this dataset: the edge attribute  $E_{(i,j)}$  denotes the existence of  
1103 the edge, and node attributes are continuous values computed following the polynomial function:  
1104  $f(x) : y = ax^2 + bx + c$  ( $a = 0; b = 1; c = 5$ ), where  $x$  is the node degree and  $f(x)$  is the node  
1105 attribute. This dataset contains 1,000 graphs in total, and can be used for attributed graph generation.

1106 The original Erdos-Renyi Graphs dataset contains 3,000 .csv files. We reformat them into .numpy  
1107 files, including input\_adj.npy, input\_edge.npy, input\_node.npy, target\_adj.npy and target\_edge.npy,  
1108 target\_node.npy for the community to use. Detailed information can be found in ER\_Readme.rtf. To  
1109 reformat the data, we use glob to read in all .csv files from the directory, and separate the original  
1110 .csv files into input data and target data; For both input and target data, we get edge feature from  
1111 the original .csv files, get node feature from the diagonals of each file, and get adjacent matrix from  
1112 the edge feature while setting the diagonals to be 0. We then reshaped all arrays into the required  
1113 dimensions and converted them to NumPy files.

1114 **Scale-free:** This dataset is generated as a directed scale-free network [41], which is a network  
1115 whose degree distribution follows power-law property [83]. It fits the "one-to-many" mapping graph  
1116 translation problem. There are no node features in this dataset, and the goal is to learn the mapping  
1117 from the input graph's topology to the target graph's topology. To generate a target graph, a node  
1118 will be selected as target node with probability proportional to its in-degree, which will be linked to  
1119 a new source node with probability of 0.41. Similarly, a node will be selected as the source node  
1120 with the probability proportional to its out-degree, which will be linked to a new target node with  
1121 the probability of 0.54. Then, a corresponding target graph is generated by adding  $m$  (number of  
1122 nodes of the input graph) edges between two nodes. Thus, both input and target graphs are directed  
1123 scale-free graphs. This dataset contains 10,000 graphs in total, and can be splitted into subsets that  
1124 contains 10, 20, 50, 100, 150 nodes along with 20, 40, 100, 200 and 320 edges, respectively.

1125 The original Scale-free dataset contains 10,000 .csv files and we convert it to .numpy files for peo-  
1126 ple to read in Python. The detailed information of the data can be found in the corresponding  
1127 scale\_free\_Readme.rtf. To reformat the data, we use glob to read in all .csv files from the directory,  
1128 and separate the original .csv files into input data and target data; For both input and target data, we  
1129 get edge feature from the original .csv files, get node feature from the diagonals of each file, and get  
1130 adjacent matrix from the edge feature while setting the diagonals to be 0. Due to the massive .csv files  
1131 in the Scale-free Graphs, we optimize to reduce the time complexity in order to process the dataset  
1132 faster. We then reshaped all arrays into the required dimensions and converted them to NumPy files.

1133 **Waxman Graphs:** This dataset contains graphs generated by the Waxman random graph model that  
1134 places  $n$  nodes uniformly at random in a rectangular domain [106, 29]. There are three types of  
1135 factors that are related to the generation of Waxman graphs: the independent graph factor  $b$  that  
1136 controls node attributes, the independent spatial factor  $p$  that controls the overall node positions,  
1137 and the graph-spatial correlated factor  $s$  that controls both graph and spatial density [29]. There are  
1138 80,000 samples for training and 80,000 for testing. Each graph in the dataset contains 25 nodes and  
1139 around 250 edges. Waxman Graphs dataset can be used for a few tasks, including attributed graph  
1140 generation, spatial graph generation and temporal graph generation.

1141 The original Waxman Graphs dataset contains 96,000 graph files saved in Numpy array.  
1142 We reformat them into .numpy files, including adj.npy, edge\_feat.npy, label.npy, node\_feat.npy,  
1143 spatial.npy, temporal\_adj.npy, temporal\_edge.npy, temporal\_label.npy, temporal\_node.npy and  
1144 temporal\_spatial.npy. The detailed information can be found in waxman\_Readme.rtf. To reformat  
1145 these files, we load the testing and training dataset and converted the sparse matrices to dense matrices.  
1146 we concatenate the testing and training datasets and reshape them into the required dimensions. To  
1147 get the version of datasets with temporal dimension, we flattened the NumPy arrays. All datasets  
1148 were saved as NumPy files eventually.



1149 **Random Geometric Graphs:** This dataset contains graphs generated by the random geometric graph  
1150 model that places  $n$  nodes uniformly at random in a rectangular domain [29]. Two nodes are joined  
1151 by an edge if their distance is larger than a threshold  $\beta = 12$ . The node attributes among a graph  
1152 are generated in the same rule as that for generating Waxman graphs. There are 8,000 samples for  
1153 training and 1,600 for testing in this dataset. Each graph in the dataset contains 25 nodes and around  
1154 350 edges. Random Geometric Graphs dataset can be used for a few tasks, including attributed graph  
1155 generation, spatial graph generation and temporal graph generation.

1156 The original Random Geometric Graphs dataset contains 96,000 graph files saved in Numpy ar-  
1157 ray. We reformat them into .numpy files, including adj.npy, edge\_feat.npy, label.npy, node\_feat.npy,  
1158 spatial.npy, temporal\_adj.npy, temporal\_edge.npy, temporal\_label.npy, temporal\_node.npy and  
1159 temporal\_spatial.npy. The detailed information can be found in random\_geo\_Readme.rtf. To  
1160 reformat these files, we load the testing and training dataset and converted the sparse matrices to  
1161 dense matrices. we concatenate the testing and training datasets and reshape them into the required  
1162 dimensions. To get the version of datasets with temporal dimension, we flattened the NumPy arrays.  
1163 All datasets were saved as NumPy files eventually.

#### 1164 C.12.1 License

1165 **Barab’asi-Albert Graphs:** The dataset is generated by [31]. The license is not specified.

1166 **Community:** The dataset is generated by [34], which is under MIT License. The license of the  
1167 dataset is not specified.

1168 **Erdos-Renyi graphs:** The dataset is generated by [31]. The license is not specified.

1169 **Scale-free:** The dataset is generated by [41]. The license is not specified.

1170 **Waxman graphs:** The dataset is generated by [29]. The license is not specified.

1171 **Random geometric:** The dataset is generated by [29]. The license is not specified.

## 1172 D Benchmark Results

1173 We benchmark all the datasets with graph generation and transformation models. For graph generation  
1174 task, we conduct experiments on three models, GraphRNN [34], GraphVAE [18], GraphGMG [8].  
1175 For graph transformation task, we conduct experiments on two models, Interaction Networks [38]  
1176 and NEC-DGT [31].

### 1177 D.1 Molecule Generation Results.

1178 As mentioned above, graph generation task could be very domain-specific, meaning that each domain  
1179 has specific expectations over the generative tasks. Our first benchmark focuses on one of the most  
1180 developed areas, molecular graph generation, which is motivated by drug and material discovery. For  
1181 molecule generation task, we utilize the above mentioned self-quality based evaluation, where the  
1182 validity, uniqueness and novelty are measured. We survey a list of state-of-the-art deep generative  
1183 models on molecules and report the performance regarding validity, novelty, and uniqueness on  
1184 two popular benchmark datasets (QM9 [44] and ZINC250K [45]) in the original paper as shown  
1185 in Table 4. In Table 4, it is clearly to observe that the state-of-the-art models, such as MoFlow,  
1186 GraphEBM, GraphDF, almost perform perfectly on the two common benchmarked datasets. As  
1187 described in the following section, one key point to generate good molecular graphs is to handle the  
1188 valency constraints. Some models utilize sequential generation, some utilize valency check, some  
1189 design regularization, but overall, the best-performing models handle the valency constraint properly.  
1190 However, it is not the end of the area. The molecule space being searched currently is small with  
1191 very limited set of atoms and bonds and small size of molecules. Thus, benchmark datasets with  
1192 larger molecules and molecules with more diverse atom and bond types are urgent to advance the  
1193 field. From another perspective, it is important to generate molecules with desired properties which  
1194 more domain-specific analyses and explorations could be done.

### 1195 D.2 Baseline Models

1196 **GraphRNN [34].** GraphRNN represents graph generation as an auto-regressive process and builds  
1197 an generative RNN model to generate nodes and edges sequentially.

1198 **GraphVAE [18].** GraphVAE represents each graph by its adjacency matrix and feature vectors  
1199 and utilizes graph neural network to encode the graphs into a vector space. Then, the model learns  
1200 the distribution of the graphs via a VAE setting which minimizes the distance between the latent  
1201 distribution and Gaussian distribution. Finally, the model decodes the latent vectors to reconstruct  
1202 graphs.

**Table 4:** Quantitative evaluation and comparison on molecular graph generation tasks by different deep generative models on graphs (“Valid.” is short for validity. “Novel.” is short for novelty. “Unique.” is short for uniqueness.).

Method→ Dataset↓	QM9			ZINC250K		
	Valid.	Novel.	Unique.	Valid.	Novel.	Unique.
GrammarVAE [107]	31.00%	<b>100.00%</b>	10.76%	30.00%	95.44%	9.30%
GraphVAE [18]	14.00%	<b>100.00%</b>	31.60%	61.00%	85.00%	40.90%
CGVAE [108]	<b>100.00%</b>	<b>100.00%</b>	99.82%	<b>100.00%</b>	94.35%	98.54%
GraphNVP[86]	74.30%	<b>100.00%</b>	94.80%	90.10%	54.00%	97.30%
GRF [109]	73.40%	<b>100.00%</b>	53.70%	84.50%	58.60%	66.00%
GraphAF[35]	<b>100.00%</b>	<b>100.00%</b>	99.10%	<b>100.00%</b>	88.83%	94.51%
CGSVAE [110]	34.90%	<b>100.00%</b>	-	96.60%	97.50%	-
JT-VAE [20]	<b>100.00%</b>	<b>100.00%</b>	99.80%	-	-	-
GCPN [93]	<b>100.00%</b>	<b>100.00%</b>	<b>99.97%</b>	-	-	-
MolecularRNN [36]	<b>100.00%</b>	<b>100.00%</b>	99.89%	-	-	-
MolGAN [19]	-	-	-	98.10%	94.10%	10.40%
MPGVAE [111]	-	-	-	91.00%	54.00%	68.00%
SCAT [112]	-	-	-	47.40%	92.00%	98.30%
MoFlow [32]	<b>100.00%</b>	98.03%	99.20%	<b>100.00%</b>	<b>100.00%</b>	<b>99.99%</b>
GraphEBM [33]	<b>100.00%</b>	97.01%	97.90%	99.96%	<b>100.00%</b>	98.79%
GraphDF [37]	<b>100.00%</b>	98.10%	97.62%	<b>100.00%</b>	<b>100.00%</b>	99.55%

1203 **GraphGMG [8].** GraphGMG first learns a node-level embedding of a given graph, then learns  
 1204 a probability distribution over possible outcomes for each generation step. During the generation  
 1205 process, the model sequentially connects nodes and edges to a new graph.

1206 **GrammarVAE [107].** GrammarVAE is one of the first deep generative models that learn to generative  
 1207 novel molecules with a string representation.

1208 **GraphVAE [18].** GraphVAE is a VAE-based graph generative models that generates graphs in an  
 1209 one-shot fashion.

1210 **CGVAE [108].** CGVAE is a VAE-based graph generative model that formulates the generation  
 1211 process as an iterative process.

1212 **GraphNVP [86].** GraphNVP first introduces the idea of invertible normalizing flow-based methods  
 1213 to molecular graph generation in an one-shot generation way.

1214 **GRF [109].** GRF introduces residual flows for molecular graph generation which circumvents the  
 1215 requirement of partitioning of the latent vector in GraphNVP.

1216 **GraphAF [35].** GraphAF takes one step further than GraphNVP to formulate the problem as a  
 1217 sequential generation problem.

1218 **CGSVAE [110].** CGSVAE is a VAE-based graph generative models that proposes a regularization  
 1219 method that encourages the model to generate valid molecules.

1220 **JT-VAE [20].** JT-VAE is motivated to explicitly model substructures in the generative models that  
 1221 introduces an extra junction tree encoder-decoder part which each node denotes a substructure rather  
 1222 than an atom in a molecule.

1223 **GCPN [93].** GCPN formulates molecular graph generation as a reinforcement learning problem  
 1224 where each state is a generation step, every step, it takes the action to connect two atoms and labels  
 1225 the edges by bond types. It stops when no atoms are connected.

1226 **MolecularRNN [36].** MolecularRNN follows the idea of GraphRNN and adopts it for the molecular  
 1227 graph generation task.

1228 **MolGAN [19].** MolGAN is a GAN-based molecular graph generation method that implements a  
 1229 GAN model to generate molecular graphs in an one-shot fashion.

1230 **MPGVAE [111].** MPGVAE designs a VAE-based which follows Graphite [113] and generate  
 1231 molecular graphs in an one-shot way.

1232 **SCAT [112].** SCAT takes a scattering transform and gaussianization as an encoder and utilizes a  
 1233 MLP as a decoder to generate novel molecular graphs in an one-shot way.

1234 **MoFlow [32].** MoFlow improves over GraphNVP by introducing a valency correction mechanism in  
 1235 the framework.

**Table 5:** Hyper-parameters for graph generation benchmark.

Method	Learning Rate	Epoch	Batch Size	# graphs
GraphRNN	$3 \times 10^{-3}$	1,000	32	1,000
GraphVAE	$1 \times 10^{-3}$	10	1	1,000
GraphGMG	$1 \times 10^{-3}$	10	1	1,000

**Table 6:** Hyper-parameters for graph transformation benchmark. Due the capacity of our memory, for the graph transformation task, we sampled a subset from a few datasets for evaluation. The size of the subset depends on the graph size and total number of graphs contained in the dataset.

Dataset	Interaction Networks			NEC-DGT		
	Learning Rate	Epoch	#graphs	Learning Rate	Epoch	#graphs
AuthNet	$1 \times 10^{-2}$	100	412	$1 \times 10^{-4}$	500	412
Barab’asi-Albert Graphs	$1 \times 10^{-2}$	100	1,000	$1 \times 10^{-4}$	500	1,000
Brain-restingstate	$1 \times 10^{-2}$	100	823	$1 \times 10^{-4}$	500	823
Brain-emotion	$1 \times 10^{-2}$	100	811	$1 \times 10^{-4}$	500	811
Brain-grambling	$1 \times 10^{-2}$	100	818	$1 \times 10^{-4}$	500	818
Brain-language	$1 \times 10^{-2}$	100	816	$1 \times 10^{-4}$	500	816
Brain-motor	$1 \times 10^{-2}$	100	816	$1 \times 10^{-4}$	500	816
Brain-relational	$1 \times 10^{-2}$	100	808	$1 \times 10^{-4}$	500	808
Brain-social	$1 \times 10^{-2}$	100	816	$1 \times 10^{-4}$	500	816
Brain-wm	$1 \times 10^{-2}$	100	812	$1 \times 10^{-4}$	500	812
Scale-free	$1 \times 10^{-2}$	100	250	$1 \times 10^{-4}$	500	250
TwitterNet	$1 \times 10^{-2}$	100	250	$1 \times 10^{-4}$	500	250
N-body-charged	$1 \times 10^{-2}$	100	150	$1 \times 10^{-4}$	500	150
N-body-spring	$1 \times 10^{-2}$	100	150	$1 \times 10^{-4}$	500	150
ChemReact	$1 \times 10^{-2}$	100	1,000	$1 \times 10^{-4}$	500	1,000
IoTNet	$1 \times 10^{-2}$	100	343	$1 \times 10^{-4}$	500	343
MolOpt	$1 \times 10^{-2}$	100	500	$1 \times 10^{-4}$	500	500

1236 **GraphEBM [33].** GraphEBM is an energy-based generative model that utilizes Langevin Dynamics  
 1237 to sample novel molecules.

1238 **GraphDF [37].** GraphDF improves over GraphAF by learning discrete latent variables rather than  
 1239 continuous latent variables as in most of the Flow and VAE-based methods.

1240 **Interaction Network [38].** Physical domain is the target for Interaction Networks, the input of which  
 1241 is a graph that represents a system of objects and relations. Interaction Networks instantiates the  
 1242 pairwise interaction and compute its effects via a relational model. The effects are then aggregated  
 1243 and combined with the objects and external effects to generate the input for an object model, which  
 1244 predicts how the interactions and dynamics influence the objects.

1245 **NEC-DGT [31].** In NEC-DGT, the node and edge attributes of input graphs are inputted to the  
 1246 model. The model outputs node attributes and edges attributes of the generated target graphs via  
 1247 several blocks, which have edge and node translation paths co-evolved and combined by a graph  
 1248 regularization during training process.

### 1249 D.3 Hyper-parameters

1250 All experiments are conducted on a 64-bit machine with a 6 core Intel CPU i9-9820X, 32GB RAM,  
 1251 and an NVIDIA GPU (GeForce RTX 2080ti, 1545MHz, 11GB GDDR6). The detailed hyper-  
 1252 parameters can be found in Table 5 and Table 6. For the molecular graph generation benchmark, we  
 1253 take experiment results from the original reports.

## 1254 E Tutorials

1255 We provide data processors, evaluators, as well as visualizers which simplify the pipeline for graph  
 1256 generation and transformation, as shown in Fig. 4, 5 and 6, respectively.

```

import graphgt

qm9_data_loader = graphgt.DataLoader(name='qm9', save_path='./', format='numpy')

Downloading node feature...
100% ██████████ | 31.1M/31.1M [00:03<00:00, 9
.13MiB/s]
Done!
Downloading edge feature...
100% ██████████ | 49.5M/49.5M [00:04<00:00, 1
0.1MiB/s]
Done!
Downloading spatial feature...
100% ██████████ | 35.8M/35.8M [00:03<00:00, 9
.19MiB/s]
Done!
Downloading adjacency matrix...
100% ██████████ | 49.5M/49.5M [00:04<00:00, 1
0.2MiB/s]
Done!
Downloading smiles string...
100% ██████████ | 15.0M/15.0M [00:02<00:00, 6
.97MiB/s]
Done!

adj, node_feat, edge_feat, spatial, smile = qm9_data_loader.get_data()

```

Figure 4: Loading generation dataset.

```

import graphgt

ER20_data_loader = graphgt.DataLoader(name='ER_20', save_path='./', format='numpy')

Downloading input node feature...
100% ██████████ | 80.1k/80.1k [00:00<00:00,
230kiB/s]
Done!
Downloading input edge feature...
100% ██████████ | 1.60M/1.60M [00:01<00:00, 1
.51MiB/s]
Done!
Downloading input adjacency matrix...
100% ██████████ | 800k/800k [00:00<00:00,
902kiB/s]
Done!
Downloading target node feature...
100% ██████████ | 80.1k/80.1k [00:00<00:00, 1
.20MiB/s]
Done!
Downloading target edge feature...
100% ██████████ | 1.60M/1.60M [00:01<00:00, 1
.54MiB/s]
Done!
Downloading adjacency matrix...
100% ██████████ | 800k/800k [00:00<00:00,
883kiB/s]
Done!

input_adj, input_node_feat, input_edge_feat, input_spatial, target_adj, target_node_feat, target_edge_feat, target_s

```

Figure 5: Loading transformation dataset.

```

import graphgt
import numpy as np

batch = 1000
x = np.random.rand(batch,1)
y_baseline = np.random.rand(batch,1)
y_pred = np.zeros((batch,1))

print('MMD baseline', graphgt.compute_mmd(x,y_baseline))
print('MMD prediction', graphgt.compute_mmd(x,y_pred))
print('KLD', graphgt.compute_kld(x,y_baseline))
print('EMB', graphgt.compute_emd(x,y_baseline))

MMD baseline 9.684740112247958e-05
MMD prediction 0.3751574658037742
KLD [0.51577211]
EMB 0.01009273634128826

```

Figure 6: Evaluation APIs.