756 FURTHER DETAILS ON MODEL А

758 We learn node embeddings using a graph neural network (GNN) (Kipf & Welling, 2017; Veličković 759 et al., 2018), which is a deep learning model that leverages graph-structured data by iteratively 760 aggregating and transforming feature information from neighboring nodes. The GNN takes in as 761 inputs the graph G and a set of features for each node *i*. We use one-hot node features which is 762 common in settings like ours without natural node features (Cui et al., 2022). We learn the type embeddings using a linear layer with one-hot type feature inputs. 763

764 The details of our model architecture are as follows: We use a 2 layer GNN where each layer consists 765 of a graph convolution, leaky ReLU activation, and batch normalization. We use an intermediate 766 dimension equal to the number of nodes n = 2292 and an embedding dimension of $E_n = E_{\tau} = 50$.

767 We batch our data such that data points with observed ratings and unobserved ratings are batched 768 separately. During training, we freeze the reporting model for batches for which there are no observed 769 ratings (i.e., we learn the reporting coefficients only from types for which ratings are observed). 770

We conduct a hyperparameter search over the loss weights $\gamma_1, \gamma_2, \gamma_3$, embedding dimension sizes, 771 number of layers, batch size, and learning rate using Weights and Biases on a validation set. We select 772 the set of hyperparameters that maximize the correlation of predicted reports and ratings. Based on 773 the hyperparameter search, we run experiments with a learning rate of 0.01 and a batch size of 16000. 774 Our full model uses weights $\gamma_1 = 20, \gamma_2 = 1, \gamma_3 = 10^{-6}$. All experiments are conducted on a cluster 775 with access to NVIDIA A100 and A6000 GPUs. Our model can comfortably train on one GPU. 776

In our experiments, we wish to assess the effect of using reports and ratings. Thus, we compare 777 inferences from models with (i) both reports and ratings (full model), (ii) only reports (reports-only 778 *model*), and (iii) only ratings (*ratings-only model*). The full model uses both reporting and rating 779 data and all demographic coefficients. Its hyperparameters are set to the specifications listed above. The reports-only and ratings-only models are identical to the full model, except for their loss. The 781 reports-only model sets a weight of 0 on the loss terms that evaluate against ground truth reports 782 $\mathcal{L}_{report unobserved}, \mathcal{L}_{report observed}$. The ratings-only model sets a weight of 0 on the loss term that evaluates 783 against ground truth ratings \mathcal{L}_{rating} . 784

785 В FURTHER DETAILS ON SEMI-SYNTHETIC EXPERIMENTS 786

787

B.1 SEMI-SYNTHETIC DATA

789 We generate synthetic inspection ratings r_{ikt} using equation 11. We separately generate ratings 790 for the train and test split. For example, for the train split, $\mathbb{E}_t(T_{ikt})$ is defined as the empirical 791 frequency of T_{ikt} over all weeks in the train time period. We draw α_k and θ_k from a Gaussian. 792 The mean of the Gaussian is calculated as follows: We take our real rating data, and separately for 793 each type fit a logisitic regression predicting reports from demographics and the ground truth rating $(T_{ikt} \sim X_i, r_{ikt})$. We set the mean α_k and θ_k to be the mean coefficients predicted across these 794 type-specific logisitc regressions. We set the intercept such that the ratings are zero mean. Thus, our generated and real inspection ratings take on both negative and nonnegative values. 796

798

801 802

804 805

797

788

B.2 Semi-synthetic Results

1	9	9	
8	0	0	

	Full model	Reports-only model	Ratings-only model
RMSE on predicted reports	0.08	0.06	-
RMSE on predicted ratings	1.01	_	1.01

Table 4: Semisynthetic data RMSE results. We compare our full model to a reports-only and 807 a ratings-only model. Compared to both baselines, our full model can estimate ratings without compromising accuracy in predicting reports. We calculate the RMSE between our predicted 808 probabilities of reports and the true probabilities for all node/type pairs. We calculate the RMSE between our predicted ratings and the true ratings for all nodes and for all types with observed ratings. We report the median correlation across 5 synthetic datasets.

We evaluate our predicted reports and ratings
using both *correlation* and *root mean squared error (RMSE)*.

813 814

815 **Correlation results:** We evaluate reports by 816 calculating the correlation between our models 817 predicted probability of a report and the average 818 true report across each node/type pair. In other 819 words, we calculate $\operatorname{corr}(\hat{P}(T_{ikt}), \mathbf{E}_t[T_{ikt}])$. 820 We evaluate ratings by calculating the correla-821 tion between our models predicted rating and the 822 average true rating across each node/type pair. In other words, we calculate corr(\hat{r}_{ikt} , $\mathbf{E}_t[r_{ikt}]$). 823

824 In Table 1, we calculate the average corre-825 lation on reports across all node/type pairs 826 and the average correlation on ratings across 827 node/type pairs with observed ratings. The 828 ratings-only model only predicts ratings, so we cannot evaluate its performance on predicting 829 reports. Similarly, the reports-only model only 830 predicts probabilities of reports. Thus in or-831 der to estimate the reports-only model's correla-832 tion on ratings we use the predicted probability 833 of a report as a proxy for rating and evaluate 834 $\operatorname{corr}(P(T_{ikt}), \mathbf{E}_t[r_{ikt}]).$ 835



Figure 5: We evaluate our model's performance in predicting ratings across type frequencies. We measure type frequency as $\mathbb{E}_{it}[T_{ikt}]$. Particularly for rare types, compared to the reports-only model, our full model, which uses both reporting and rating data, predicts more correlated ratings. We show results for all types that the model *does not* observe ratings for. We plot the median across 5 synthetic datasets.

In Figure 2a, we calculate the correlation on reports for each type with observed ratings separately.
 In Figure 5, we calculate the correlation on reports for each type with unobserved ratings. In both cases, compared to the reports-only model, we find that our full model predicts ratings that are more correlated with the ground truth.

840 841

858 859

RMSE results: We evaluate reports by calculating the RMSE between our models predicted probability of a report and the average true report across each node/type pair. In other words, we calculate RMSE($\hat{P}(T_{ikt})$, $\mathbf{E}_t[T_{ikt}]$). We evaluate ratings by calculating the RMSE between our models predicted rating and the average true rating across each node/type pair. In other words, we calculate RMSE(\hat{r}_{ikt} , $\mathbf{E}_t[T_{ikt}]$).

In Table 4, we calculate the average RMSE on reports across all node/type pairs and the average RMSE on ratings across node/type pairs with observed ratings. We report the median RMSE across 5 synthetic datasets. We note that the ratings-only model only predicts ratings. Therefore, we cannot evaluate its performance on predicting reports. Similarly, the reports-only model only predicts probabilities of reports. Therefore, we cannot evaluate its performance on predicting ratings. Note that unlike for correlation, when calculating RMSE, we *cannot* use a proxy for rating (e.g., $\hat{P}(T_{ikt})$).

	Full model	Reports-only model	Ratings-only model
RMSE on predicted reports	0.11	0.06	_
RMSE on predicted ratings	0.83	-	0.84

Table 5: Real data RMSE results. We compare our full model to a reports-only and a ratings-only model. Compared to both baselines, our full model can estimate ratings without compromising accuracy in predicting reports. We calculate the RMSE between our predicted probabilities of reports and the true probabilities for all node/type pairs. We calculate the RMSE between our predicted ratings and the true ratings for all nodes and for all types with observed ratings.

864 С FURTHER DETAILS ON REAL DATA 865

866 **Details on processing real reporting data:** We use reports T_{ikt} from New York City 311 data (NYC Open Data, 2024a). We analyze all Census tracts with valid demographic information (n = 2292) 868 nodes), complaint types with a reporting frequency greater than 0.1% ($\tau = 141$ types), and all weeks in the two years from 2022 - 2023. $T_{ikt} \in \{0,1\}$ denotes whether at least one report of type k was 870 made in node i during week t. In total we analyze more than 55 million reports.

871

867

872 Feature processing: We include demographic features collected for each Census tract. The full list of features that we include is: log population density, percentage of population with a bachelors 873 degree, percentage of households occupied by a renter, log median income, percentage of population 874 that is white, and median age. We normalize all features to have mean 0 and standard deviation 1. 875

876

Details on processing real rating data: 877 We collect ratings from government in-878 spection data for five complaint types: (i)

879 street conditions (NYC Open Data, 2023), 880 (ii) park maintenance or facility conditions (NYC Open Data, 2024c), (iii) rodents (NYC Open Data, 2024e), (iv) food es-883 tablishment/mobile food vendor/food poi-884 soning (NYC Open Data, 2024d), and (v) 885 DCWP consumer complaints (NYC Open Data, 2024b). Each rating is for a fine-886 grained unit within a Census tract. Street 887 ratings are for street segments; park ratings are for parks; rodent ratings are averaged 889 over each Borough-Block-Lot (BBL); food 890 ratings are averaged over each BBL; and 891 DCWP ratings are averaged over each Cen-892 sus block. We match each fine-grained rat-

Covariate	Mean coefficient
Bachelors degree population	0.28
Households occupied by renter	0.24
log(Population density)	0.20
Median age	0.13
White population	-0.08
log(Median income)	-0.10
True inspection rating	-0.20

Table 6: Multivariate reporting coefficients. We report the average predicted multivariate demographic coefficients across types with observed ratings. The estimated coefficients capture known demographic factors: tracts that are more dense, more educated, or are older are more likely to report incidents. We also report the coefficient on the true inspection rating. Tracts that have lower ratings are more likely to be reported.

893 ing to its corresponding fine-grained report (i.e., reports in that same street segment). For rodents, 894 food, and DCWP the matching is done directly (i.e., we match the aggregated rodent rating for a 895 BBL to the aggregated report for the same BBL). For streets and parks, we run a distance heuristic to complete the matching. We match each rating with its nearest report. If the nearest report is above a 896 certain distance threshold, we filter out the rating (consider it unobserved). Within the same tract, 897 all fine-grained ratings and reports are provided to the model and are mapped to the same node's 898 embedding, as well as the corresponding type's embedding. 899

900 We also process the inspection data to remove any inspections triggered by 311 reports. The rodent 901 inspection data dictionary states that DOHMH inspectors run both random inspections and inspections 902 triggered by 311 reports NYC Open Data (2024e). It is also stated that the random inspections occur block by block. The inspection data is not labeled as random versus 311 initiated, thus we run a 903 heuristic to identify inspections triggered by 311 reports. We calculate the number of inspections that 904 occur each week in each Census tract. We filter out all inspections that fall in tracts under the 50th 905 percentile. Inspection data for the other types are described to be purely random. 906

907 908

909 910

911

912

916

D FURTHER DETAILS ON THE REAL-WORLD CASE STUDY

Real data results We report our model's correlation on predicted reports and ratings in Table 2. In Table 5, we report our model's RMSE on predicted reports and ratings.

Predicted demographic coefficients: In Table 3 we report the demographic coefficients predicted 913 by univariate models. In Table 6 we report the demographic coefficients predicted by a multivariate 914 model. 915

Clustered nodes are demographically distinct. For each node *i*, we create a vector $\mathbf{r}_i = \{r_{ikt}\}_{i=1}^{T}$ 917 of ratings over all types k. We use each node's \mathbf{r}_i vector to cluster the nodes into 4 groups. We

918	Cluster	0	1	2	3
919	Race:Non-Hispanic White	55%	29%	34%	12%
920	Race:Asian	16%	18%	18%	5%
921	Race:African-American	8%	19%	22%	35%
922	Households occupied by renter	72%	60%	46%	87%
923	Bachelors degree	71%	33%	36%	26%
924	Population	5,500	3,900	2,600	5,200
025	Median income	120,000	71,000	73,000	47,000
925	Median age	37	38	40	35

Table 7: Clustering ratings for each node. We find that the clustering correlates with differences in demographics. All differences between clusters are statistically significant (p < 0.001, ANOVA test). The largest value in each row is shown in bold.

 find that the predicted clusters are spatially correlated and demographically distinct. We report the statistically significant differences in demographics for each cluster in Table 7.

Clustering ratings for each type: For each type k, we create a vector $\mathbf{r}_k = \{r_{ikt}\}_{i=1}^n$ of ratings over all nodes i to cluster the types into 8 groups. We find that each group contains a coherent cluster of types, and in Table 8 we describe and list the types captured by each cluster. Additionally, Figure 6 shows that the dimension of highest variability (i.e., first PCA dimension) of the \mathbf{r}_k vectors captures type frequency (i.e., $\mathbb{E}_{it}[T_{ikt}]$).



Figure 6: Each type's learned ratings over nodes capture type frequency information. In particular, the dimension of highest variance of our type ratings (PC1) has a high correlation with type reporting frequency $\mathbb{E}_k[T_{ikt}]$.

972 Table 8: Ratings capture correlations between 311 complaint types: Using each type's vector of 973 learned ratings over nodes, we cluster types into 8 groups using a k means clustering algorithm. We 974 manually assign a succinct cluster description to each group. We find that the clusters group similar 975 types together.

978	Cluster Des
979	
980	
981	Natao and D
982	Noise and P
983	ASSIStance
984	
985	
986	
987	
988	
989	Residential
990	Parking Vic
991	
992	
993	
994	
995	
996	
997	Housing Mai
998	
999	
1000	
1001	
1002	
1003	
1004	Stroot and
1005	Conditions
1006	condicions
1007	
1008	
1009	
1010	

976 977

Cluster Description	Complaint Types		
	Consumer Complaint (DCA)		
	Noise (DEP)		
	Homeless Person Assistance (DHS)		
Noise and Public	Traffic Signal Condition (DOT)		
Assistance Issues	Encampment (NYPD)		
	Noise - Commercial (NYPD)		
	Noise - Vehicle (NYPD)		
	For Hire Vehicle Complaint (TLC)		
	Dirty Condition (DSNY)		
	Missed Collection (DSNY)		
	Heat/Hot Water (HDD)		
Residential and	Unsanitary Condition (HPD)		
Parking Violations	Placked Drivoway (NYPD)		
Farking violations	Illogal Darking (NYDD)		
	IIIEgal Parking (NIPD)		
	Noise - Residential (NYPD)		
	Noise - Street/Sidewalk (NYPD)		
	Appliance (HPD)		
	Door/Window (HPD)		
	Electric (HPD)		
Housing Maintenanc	Flooring/Stairs (HPD)		
inclusing marmeenane	General (HPD)		
	Paint/Plaster (HPD)		
	Plumbing (HPD)		
	Water Leak (HPD)		
	Sewer (DEP)		
	Water System (DEP)		
	General Construction/Plumbing (DOB)		
	Rodent (DOHMH)		
Street and Vehicle	Sidewalk Condition (DOT)		
Conditions	Street Light Condition (DOT)		
	Damaged Tree (DPR)		
	Derelict Vehicles (DSNY)		
	Illegal Dumping (DSNY)		
	Abandoned Vehicle (NYPD)		
	Air Quality (DEP)		
	Lead (DEP)		
	Building/Use (DOB)		
	Elevator (DOB)		
	Curb Condition (DOT)		
	Street Sign - Damaged (DOT)		
	Dead/Duing Tree (DPR)		
Environmental and	New Tree Request (DPR)		
Building Operation	New The Neguest (DIN) Overgrown Tree/Branches (DDD)		
Concerns	Doot /Sever/Sidewalk Condition (DDD)		
CONCETIIS	Root/Sewer/Staewark Collattion (DFR)		
	Deau Allillai (DSNI) Electropico Nocto Appointment (DCNV)		
	Electronics waste Appointment (DSNI)		
	ODSUPUCTION (DENY)		
	Residential Disposal Complaint (DSNY)		
	Street Sweeping Complaint (DSNY)		
	Salety (HPD)		
	Non-Emergency Police Matter (NYPD)		

1026		Asbestos (DEP)
1027		AHV Inspection Unit (DOB)
1028		BFST/Site Safety (DOB)
1029		Scaffold Safety (DOB)
1030		Achastas (DOUMU)
1021		Roach/Rool/Sauna Complaint (DOUMU)
1031		Construction Load Dust (DOUMH)
1032		Construction Lead Dust (DOHMH)
1033		Drinking Water (DOHMH)
1034		Illegal Animal Kept as Pet (DOHMH)
1035		Indoor Sewage (DOHMH)
1036		Mold (DOHMH)
1027		Mosquitoes (DOHMH)
1007		Pet Shop (DOHMH)
1038		Poison Ivy (DOHMH)
1039	Uselth and Cafeta	Tattooing (DOHMH)
1040	Health and Salety	Bike Rack Condition (DOT)
1041		Bus Stop Shelter Placement (DOT)
1042		DEP Street Condition (DOT)
10/13		E-Scooter (DOT)
10//		Uprooted Stump (DPR)
1044		Wood Pile Remaining (DPP)
1045		Nood IIIE Nemaining (DER) Ndort-N-Rackat (DSNV)
1046		AUUPL-A-DASKEL (DSNI)
1047		Seasonal Collection (DSNY)
1048		Outside Building (HPD)
1049		Sewer (NYC311-PRD)
1050		Water System (NYC311-PRD)
1050		Disorderly Youth (NYPD)
1051		For Hire Vehicle Report (TLC)
1052		Green Taxi Complaint (TLC)
1053		Taxi Report (TLC)
1054		Consumer Complaint (DCWP)
1055		Encampment (DHS)
1056		Boilers (DOB)
1050		Emergency Response Team (ERT) (DOB)
1057		Real Time Enforcement (DOB)
1058		Special Projects Inspection Team (SPIT) (DOB)
1059		Indoor Air Quality (DOHMH)
1060		Smoking (DOHMH)
1061		Broken Barking Meter (DOT)
1062		Outdeen Dining (DOT)
1062		Ctreat Gian Dangling (DOT)
1005		Street Sign - Dangling (DOI)
1004		Street Sign - Missing (DUI)
1065		Animai in a Park (DPR)
1066	Public Space and	IIIegal Tree Damage (DPR)
1067	Community	Maintenance or Facility (DPR)
1068	Violations	Violation of Park Rules (DPR)
1069	* 1010010115	Commercial Disposal Complaint (DSNY)
1070		Graffiti (DSNY)
1070		Litter Basket Request (DSNY)
1071		Noise - Helicopter (EDC)
1072		Animal-Abuse (NYPD)
1073		Bike/Roller/Skate Chronic (NYPD)
1074		Drug Activity (NYPD)
1075		Graffiti (NYPD)
1076		Illegal Fireworks (NYPD)
1077		Noise - Park (NYPD)
1077		Dephandling (NYDD)
1078		raimanutting (NIFD)
1079		II dIIIC (NIPD)
		LOST Property (TLC)
		Taxi Complaint (TLC)

1080		
1081		Hazardous Materials (DEP)
1082		Industrial Waste (DEP)
1083		Water Conservation (DEP)
1084		Water Quality (DEP)
1085		Electrical (DOB)
1086		Investigations and Discipline (IAD) (DOB)
1000		Plumbing (DOB)
1007		School Maintenance (DOE)
1088		Day Care (DOHMH)
1089		Face Covering Violation (DOHMH)
1090		Homboring Deeg (Means (DOHMH)
1091		Non-Posidontial Hoat (DOHMH)
1092		Standing Water (DOHMH)
1093	Sanitation and	Unleashed Dog (DOHMH)
1094	Water	Unsanitary Animal Put Property (DOHMH)
1095	Water	Unsanitary Pigeon Condition (DOHMH)
1096		Bus Stop Shelter Complaint (DOT)
1097		Street Condition (DOT)
1098		Abandoned Bike (DSNY)
1099		Dumpster Complaint (DSNY)
1100		Illegal Posting (DSNY)
1101		Litter Basket Complaint (DSNY)
1102		Lot Condition (DSNY)
1103		Sanitation Worker or Vehicle Complaint (DSNY)
1104		Snow or Ice (DSNY)
1105		Elevator (HPD)
1106		Drinking (NYPD)
1107		Noise - House of Worship (NYPD)
1108		Urinating in Public (NYPD)
1100		
1110		
1110		
1111		
1112		
1113		
1114		
61115		
1116		
1117		
1118		
1119		
1120		
1121		
1122		
1123		
1124		
1125		
1126		
1127		
1128		
1129		
1130		
1131		
1132		
1133		