

Zero-shot drug repurposing with geometric deep learning and clinician centered design

Kexin Huang^{1,†,*}, Payal Chandak^{2,*}, Qianwen Wang¹, Shreyas Havaladar³, Akhil Vaid^{3,4}, Jure Leskovec⁵, Girish Nadkarni⁴, Benjamin S. Glicksberg^{3,4}, Nils Gehlenborg¹, and Marinka Zitnik^{1,6,7,8,‡}

¹Department of Biomedical Informatics, Harvard Medical School, Boston, MA 02115

²Harvard-MIT Program in Health Sciences and Technology, Cambridge, MA 02139

³Hasso Plattner Institute for Digital Health, Icahn School of Medicine at Mount Sinai, NY 10029

⁴Charles Bronfman Institute for Personalized Medicine, Icahn School of Medicine at Mount Sinai, NY 10029

⁵Department of Computer Science, Stanford University, Stanford, CA 94305

⁶Broad Institute of MIT and Harvard, Cambridge, MA 02142

⁷Harvard Data Science Initiative, Cambridge, MA 02138

⁸Kempner Institute for the Study of Natural and Artificial Intelligence, Harvard University, MA 02134

† Present address: Department of Computer Science, Stanford University

* Equal contribution

‡ Corresponding author: marinka@hms.harvard.edu

1 **Historically, drug repurposing – identifying new therapeutic uses for approved drugs – has**
2 **been attributed to serendipity. While recent advances have leveraged knowledge graphs and**
3 **deep learning to identify potential therapeutic candidates, their clinical utility remains lim-**
4 **ited because they focus on diseases with available existing treatments and rich molecular**
5 **knowledge. Here, we introduce TxGNN, a geometric deep learning approach designed for**
6 **“zero-shot” drug repurposing, identifying therapeutic candidates even for diseases with no**
7 **existing medicines. Trained on a medical knowledge graph, TxGNN utilizes a graph neural**
8 **network and metric-learning module to rank therapeutic candidates as potential indications**
9 **and contraindications across 17,080 diseases. When benchmarked against eight methods,**
10 **TxGNN significantly improves prediction accuracy for indications by 49.2% and contraindi-**
11 **cations by 35.1% under stringent zero-shot evaluation. To facilitate interpretation and anal-**
12 **ysis of the model’s predictions, TxGNN’s Explainer module offers transparent insights into**
13 **the multi-hop paths that form TxGNN’s predictive rationale. Our pilot human evaluation**
14 **of TxGNN’s Explainer showed that TxGNN’s novel predictions and explanations perform**
15 **encouragingly on multiple axes of model performance beyond accuracy. Many of TxGNN’s**
16 **novel predictions are aligned with off-label prescriptions made by clinicians within a large**
17 **healthcare system, affirming their potential clinical utility. TxGNN provides drug repurpos-**
18 **ing predictions that are more accurate than existing methods, are consistent with off-label**
19 **prescription decisions made by clinicians, and can be investigated by human experts through**
20 **multi-hop interpretable explanations.**

21 Introduction

22 There is a pressing need to develop therapies for many diseases that currently lack treatments^{1,2}.
23 Of over 7,000 rare diseases worldwide, only 5-7% of rare diseases have FDA-approved drugs³.
24 Leveraging existing therapies and expanding their use by identifying new therapeutic indications
25 via drug repurposing can alleviate the global disease burden. By using safety and efficacy data for
26 existing drugs, drug repurposing can expedite translation to the clinic and lower development costs
27 than designing drugs from scratch⁴ (Figure 1a). The fundamental premise behind repurposing is
28 that drugs can have pleiotropic effects beyond the mechanism of action of their direct targets⁵.
29 Approximately 30% of FDA-approved drugs are issued at least one post-approval new indication,
30 and many drugs have accrued over ten indications over the years⁶. However, most repurposed
31 drugs are the result of serendipity^{7,8} – either observed through off-label prescriptions written by
32 clinicians, as with gabapentin and bupropion⁸ or discovered through patient experience, as with
33 sildenafil⁶. The relationships between drug candidates and their potential new applications have
34 not been studied systematically because the underlying mechanism ‘connecting’ them is often
35 intricate and dispersed through the biomedical literature⁷.

36 Owing to technological advances, the effects of drugs can now be prospectively matched to
37 new indications by systematically analyzing medical knowledge graphs^{5,9}. The new strategies rely
38 on identifying therapeutic candidates based on their effects on cell signalling, gene expression, and
39 disease phenotypes^{5,10-12}. Machine learning has been used to analyze high-throughput molecular
40 interactomes to unravel genetic architecture perturbed in disease^{12,13} and help design therapies to
41 target them¹⁴. To provide therapeutic predictions, geometric deep learning models optimized on
42 large medical knowledge graphs¹⁵ can match disease signatures to therapeutic candidates based on
43 networks perturbed in disease¹⁵⁻¹⁹.

44 Although computational approaches have identified promising repurposing candidates for
45 complex diseases^{16,20,21}, there remain two key challenges that could significantly enhance the
46 clinical relevance of repurposing predictions made by machine learning models. (1) First, ex-
47 isting methods assume that diseases for which we would like to make therapeutic predictions are
48 well-understood and likely to have existing therapies. While this is the case for more widespread
49 diseases⁹, a long tail of diseases does not satisfy this assumption – 92% of 17,080 diseases ex-
50 amined in our study have no indications. Moreover, around 95% of rare diseases have no FDA-
51 approved drugs, and up to 85% of rare diseases do not have even one drug developed that would
52 show promise in rare disease treatment, diagnosis, or prevention²². This long tail of diseases with

53 few or no therapies and limited molecular understanding presents a clinically fruitful challenge
54 for drug repurposing models to prioritize. (2) Second, a repurposed indication for a therapeutic
55 candidate can be unrelated to the indication for which the drug was initially studied. Originally
56 proposed to help with morning sickness during pregnancy, Thalidomide was repurposed in 1964
57 for an autoimmune complication of leprosy and again in 2006 for multiple myeloma⁸. Collec-
58 tively, we refer to these challenges as the zero-shot drug repurposing problem (Figure 1b). To be
59 clinically useful, machine learning models must make “zero-shot” predictions; that is, they need
60 to extend therapeutic predictions to diseases whose understanding is incomplete and, further, to
61 diseases with no approved drugs. Unfortunately, the ability of existing machine learning models to
62 identify therapeutic candidates for diseases with incomplete, sparse data and zero known therapies
63 drops drastically^{16,23} (as we demonstrate across eight benchmarks in Figures 2c and 2d).

64 Here, we introduce TXGNN, a geometric deep learning approach for zero-shot drug repur-
65 posing that can prioritize therapeutic candidates for diseases with no therapies (Figure 1c). Foun-
66 dation models like TxGNN are transforming deep learning: instead of training disease-specific
67 models for every disease, TXGNN is a single pretrained model that adapts across many diseases.
68 TXGNN is trained on a medical knowledge graph that collates decades of biological research
69 across 17,080 diseases (Figure 1d). TXGNN uses a graph neural network model to embed ther-
70 apeutic candidates and diseases into a latent representation space and is optimized to reflect the
71 geometry of TXGNN’s medical knowledge graph. To make therapeutic predictions under zero-
72 shot settings, TXGNN implements a metric learning module to learn similarities between diseases
73 with indications and diseases without indications to transfer knowledge between these diseases
74 and make zero-shot predictions. Once trained, TXGNN performs zero-shot inference on new
75 diseases without additional parameters or fine-tuning. To facilitate interpretation and analysis of
76 the therapeutic candidates that TXGNN ranks highly, we develop a TXGNN Explainer module
77 that offers transparent insights into the multi-hop pathways that form TXGNN’s predictive ratio-
78 nale. TXGNN’s predictions and explanations are available at <http://txgnn.org>. Our pilot human
79 evaluation of TXGNN’s Explainer showed that TXGNN’s explanations perform encouragingly on
80 multiple axes of model performance such as accuracy, trust, usefulness, and time efficiency (Figure
81 4). Moreover, many of TXGNN’s novel predictions have shown alignment with off-label prescrip-
82 tions made by clinicians within a large healthcare system and TXGNN’s explanatory rationales
83 have demonstrated consistency with medical reasoning in selected case studies, encouraging the
84 potential real-world clinical utility of TXGNN.

85 Results

86 **Overview of TxGNN zero-shot drug repurposing model.** A problem not previously considered
87 in biomedical deep learning research, zero-shot drug repurposing involves predicting therapeutic
88 candidates for diseases that do not have any existing indications (Figure 1b). Mathematically, the
89 model takes a query drug-disease pair as input and provides the likelihood of the drug acting on
90 the disease as output. The gold standard labels for evaluating such a model come from our previ-
91 ously curated and validated a large-scale medical knowledge graph⁹ (Figure 1d, Tables S2 and S3)
92 that consists of 9,388 indications and 30,675 contraindications²⁴. The knowledge graph covers a
93 vast range of 17,080 diseases where 92% have no FDA-approved drugs, including rare diseases and
94 less-understood complex diseases. The knowledge graph also comprises 7,957 potential candidates
95 for drug repurposing, ranging from FDA-approved drugs to experimental drugs investigated in on-
96 going clinical trials. Our model for zero-shot drug repurposing, TxGNN operates on the principle
97 that effective drugs can target disease-perturbed and disease-associated networks of biomolecules,
98 and it has two modules: (1) the TxGNN *Predictor* module enables the accurate prediction of
99 indications and contraindications in the zero-shot setting and (2) the TxGNN *Explainer* module
100 provides interpretable multi-hop pathways that connect the drug to the disease (Figure 1c).

101 **TxGNN Predictor** The Predictor module consists of a graph neural network (GNN) optimized on
102 the relationships within the biomedical knowledge graph (Methods 2.2). Through large-scale self-
103 supervised pre-training, the GNN produces biologically meaningful representations for any entity
104 in this knowledge graph. Then, this GNN is finetuned to predict relationships between therapeutic
105 candidates and diseases. TxGNN leverages metric learning for zero-shot prediction. TxGNN
106 capitalizes on the insight that diseases are intrinsically related^{10,14} by leveraging molecular mecha-
107 nisms of well-annotated diseases to enhance predictions on diseases with limited annotations (Fig-
108 ure 2a, Figure S1). This is achieved by creating a disease signature vector for each disease based
109 on its neighbors in the knowledge graph. The similarity between a pair of diseases is measured
110 by the normalized dot product of their signature vectors. Since most disease pairs do not share
111 underlying pathologies, they have low similarity scores. In contrast, a relatively high similarity
112 score (>0.2) between diseases suggests similar mechanisms (Figure 2b). A detailed description of
113 the model and its architecture can be found in Methods 2 and Figure S2.

114 When querying a specific disease, TxGNN retrieves similar diseases, generates embeddings
115 for them, and then adaptively aggregates them based on their similarity to the queried disease. The
116 aggregated output embedding summarizes knowledge borrowed from similar diseases fused with

117 the query disease embedding. This step can also be interpreted as a graph rewiring technique in
118 the geometric machine learning literature (Figure S3). TxGNN processes different downstream
119 therapeutic tasks, such as indication and contraindication prediction, in a unified manner using
120 shared drug and disease embeddings (Methods 2.3). Given a query disease, TxGNN ranks drugs
121 based on their predicted likelihood scores, offering a prioritized list of therapeutic candidates with
122 potential for repurposing.

123 **TxGNN Explainer** While TxGNN Predictor provides likelihood scores for therapeutic candi-
124 dates, these scores alone are insufficient for trustworthy model deployment. Clinicians and scien-
125 tists seek to understand the reasoning behind these predictions to validate the model’s hypotheses
126 and better understand disease pathology. To this end, TxGNN Explainer delves into the knowl-
127 edge graph to pinpoint and succinctly present relevant biological pathways for the drug-disease
128 pair of interest (Figure 4a). This conceptual subgraph mirrors the analytical process clinical re-
129 searchers use to examine relationships between therapeutic candidates and disease and how the
130 drug perturbs local biological networks to produce a therapeutic effect on disease.

131 TxGNN uses a self-explaining approach called GraphMask²⁵ (Methods 2.6). For a particu-
132 lar therapeutic use prediction, GraphMask generates a sparse yet sufficient subgraph of biological
133 entities considered critical by TxGNN for making the prediction. Particularly, it yields an im-
134 portance score between 0 and 1 for every edge in the subgraph between the drug and disease,
135 with 1 indicating the edge is vital for prediction and 0 suggesting it is irrelevant. TxGNN Ex-
136 plainer combines the drug-disease subgraph and edge importance scores to produce multi-hop
137 explanations connecting the disease to the predicted therapeutic candidate. Unlike widely recog-
138 nized explainability techniques such as SHAP²⁶ that generate feature attribution maps, TxGNN
139 Explainer offers granular and straightforward explanations that are, as we show in a pilot human
140 study, aligned with clinician/scientist’s intuition.

141 We developed a human-centered graphical user interface that presents these subgraph ex-
142 planations proposed by TxGNN Explainer (Figure 4b). Amongst a range of designs, as shown
143 in Figures S4 and S5, we focused on visual path-based reasoning because our pilot human study
144 demonstrated that this design choice enhanced clinician comprehension and satisfaction²⁷. This
145 interface with TxGNN’s predictions and explanations is openly accessible at <http://txgnn.org>.

146 **Comparative assessment of TxGNN against existing methods.** We evaluated model perfor-
147 mance in drug repurposing across various hold-out datasets. We generated a hold-out dataset by
148 sampling diseases from the knowledge graph. These diseases were deliberately omitted during the

149 training phase and later served as test cases to gauge the model’s ability to generalize its insights
150 to previously unseen diseases. These held-out diseases were either chosen randomly, following a
151 standard evaluation strategy, or specifically selected to evaluate zero-shot prediction. In our study,
152 we used both types of hold-out datasets to thoroughly evaluate methods. We compared TxGNN to
153 eight established methods in predicting therapeutic use. They included network medicine statistical
154 techniques, including KL and JS divergence¹⁶, graph-theoretic network proximity approach²⁰, and
155 diffusion state distance (DSD)²⁸, state-of-the-art graph neural network methods, including rela-
156 tional graph convolutional networks (RGCN)^{19,29}, heterogeneous graph transformer (HGT)³⁰, and
157 heterogeneous attention networks (HAN)³¹, and a natural language processing model, BioBERT³²
158 (Supplementary Note S4).

159 Initially, we followed the standard evaluation strategy where drug-disease pairs were ran-
160 domly shuffled, and a subset of these pairs was set aside as a hold-out set (testing set; Figure 2c).
161 Under this strategy, the diseases being evaluated as hold-outs may already have had indications
162 and contraindication relationships with drugs in the training set. Therefore, the learning objective
163 was to identify additional therapeutic candidates for well-studied diseases. This evaluation method
164 aligns with the approach predominantly used in literature¹⁹. We use the area under the precision-
165 recall curve (AUPRC) as the evaluation metric as it measures the recall and precision tradeoff of a
166 model at different thresholds. Our experimental results in this setting concur, with 3 of 8 existing
167 methods achieving AUPRC greater than 0.80, and HAN as the best at 0.873 AUPRC. TxGNN
168 also had a comparable performance as established methods. In predicting indications, TxGNN
169 achieved a 4.3% increase in AUPRC (0.913) over the strongest baseline, HAN.

170 As shown by the above experiments, machine learning methods can help identify repurpos-
171 ing opportunities for diseases that already have some FDA-approved drugs^{12–16,20,21}. However,
172 Duran et al.³³ reason that many methods simply retrieve additional therapeutic candidates that are
173 similar to existing ones across biological levels. This suggests the standard evaluation strategy is
174 unsuitable for evaluating diseases that have no FDA-approved drugs (Figure 1b). Given this limi-
175 tation, we evaluate models under zero-shot drug repurposing. We began by holding out a random
176 set of diseases and then moved all their associated drugs to the hold-out set (Figure 2d). From a
177 biological standpoint, the model was required to predict therapeutic candidates for diseases that
178 lacked treatments, meaning it had to operate without any available data on drug similarities. In this
179 scenario, TxGNN outperformed all existing methods by a large margin. TxGNN significantly
180 improves over the next best baseline in predicting both indications (19.0% AUPRC gain) and con-

181 traindications (23.9% AUPRC gain). While established methods achieved satisfactory results in
182 conventional drug repurposing evaluations, they often fell short on more challenging zero-shot
183 drug repurposing scenarios. TXGNN was the only method that achieved consistent performance
184 in both settings.

185 **TXGNN’s zero-shot drug repurposing performance across disease areas.** Diseases with bio-
186 logical similarities often share therapeutic candidates¹⁰. For instance, beta-blockers are effective
187 in treating a multitude of cardiovascular issues, including heart failure, cardiac arrest, and hyper-
188 tension. Likewise, selective serotonin reuptake inhibitors (SSRIs) can address various psychiatric
189 conditions such as major depressive disorder, anxiety disorder, and obsessive-compulsive disorder.
190 If, during training, a model learns that an SSRI is indicated for major depressive disorder, it does
191 not take a large leap to suggest that the same SSRI could be effective for obsessive-compulsive dis-
192 order during testing²³. This phenomenon is known as shortcut learning^{34,35} and underlies many of
193 deep learning’s failures^{36,37}. Shortcut decision rules tend to perform well on standard benchmarks
194 but typically fail to transfer to challenging testing conditions³⁸, such as the real-world scenario of
195 predicting therapeutic candidates for rare or neglected diseases.

196 To evaluate drug repurposing models for these challenging diseases, we curated a stringent
197 hold-out dataset that contained a group of biologically related diseases that we refer to as a dis-
198 ease area. Given the diseases in a specific disease area, all their indications and contraindications
199 were removed from the training dataset. Further, a fraction of the connections from medical en-
200 tities to these diseases were excluded from the training dataset. For diseases in the chosen area,
201 these conditions simulated limited molecular characterization and lack of existing treatments (Fig-
202 ure 3a). Under this setup, we observe that diseases in the hold-out evaluation set have a signif-
203 icantly smaller number of neighbors compared to the training set (Figure S6). In this study, we
204 considered nine disease area hold-out datasets characterized in Table 1 and listed here in order of
205 increasing disease area size: (1) diabetes-related diseases such as Gestational diabetes and Lipoat-
206 rophic diabetes; (2) ‘adrenal gland’ diseases like Addison and ectopic crushing syndrome; (3) ‘au-
207 toimmune’ diseases like Celiac disease and Graves disease; (4) ‘anemia’ with conditions such as
208 thalassemia and hemoglobin C disease; (5) ‘neurodegenerative’ diseases include pick disease and
209 Neuroferritinopathy; (6) ‘mental health’ disorders like anorexia nervosa and depressive disorder;
210 (7) ‘metabolic disorder’ such as Macroglobulinemia and Gilbert syndrome; (8) ‘cardiovascular’
211 diseases, including long QT syndrome and mitral valve stenosis; (9) ‘cancerous’ diseases such as
212 neurofibroma and Leydig cell tumors. These cover a wide range of diverse disease areas.

213 We benchmarked the performance of TxGNN and all methods above on these rigorous hold-
214 out datasets in Figure 3b-f and S7. We found that TxGNN consistently improved predictive per-
215 formance over existing methods. For indications, TxGNN had 26.1%, 59.3%, 32.2%, 42.3%,
216 13.6%, 36.2%, 11.1%, 10.2%, 0.5% relative gain in AUPRC over the next best baseline across dia-
217 betes, adrenal glands, autoimmune, anemia, neurodegenerative, mental health, metabolic disorder,
218 cancer, and cardiovascular disease hold-outs respectively. For contraindications, TxGNN robustly
219 improved over the next best baseline, with relative gains ranging from 11.8% to 35.6%. For in-
220 dication prediction, the natural language processing method, BioBERT, had the best performance
221 (in 7/9 disease area hold-outs) amongst the group of established methods. For contraindication
222 prediction, the graph-based method, RGCN, was the best baseline across 8 of 9 hold-out datasets,
223 and BioBERT's performance gain observed for indication prediction disappeared. TxGNN was
224 consistently the best-performing method across all nine disease area hold-outs for both indication
225 and contraindication prediction tasks. These rigorous benchmarks demonstrate that TxGNN was
226 broadly generalizable and produced accurate predictions in zero-shot drug repurposing settings.

227 TxGNN demonstrated higher performance in eight of nine disease area hold-outs; how-
228 ever, its performance was equivalent to existing methods in the cardiovascular hold-out. This
229 equivalence may be due to an absence of related disease knowledge in the training dataset when
230 entire disease areas are excluded. Visualization of the latent representations of TxGNN Predictor
231 revealed that it supports knowledge transfer from unrelated diseases to those with limited infor-
232 mation (Figure S8). Additional evaluation metrics, including AUROC and recall, are detailed in
233 Figures S9, S10, and S11. Ablation analyses confirmed that each component of TxGNN Predictor
234 is critical for the model's predictive performance (Figure S12). Additional data splits were con-
235 ducted to stress test the model, including evaluations on diseases with minimal connections to the
236 knowledge graph (Figure S13), evaluations with certain percentages of disease local neighborhood
237 masked (Figure S14), and evaluations on various knowledge graph configurations (Figure S15).
238 These evaluations showed that TxGNN maintains robust and strong predictive performance.

239 **TxGNN's multi-hop explanations reflect model's predictive rationale.** TxGNN's Explainer
240 extracts multi-hop explanations as sequences of associations between predicted drugs and dis-
241 eases in the knowledge graph to substantiate TxGNN's predictions. This tool identifies maxi-
242 mally predictive subgraphs within the knowledge graph, connecting the query drug to the query
243 disease through multiple hops, following relationships in the graph. The performance of these
244 subgraphs is nearly equivalent to that of the entire knowledge graph. To assess the quality of ex-

245 planations, we first compared the AUPRC of TxGNN’s predictions using the entire knowledge
246 graph against the AUPRC derived from only the predictive subgraphs. A strong correlation in-
247 dicates that TxGNN’s Explainer effectively identifies key associations³⁹ and that explanations
248 accurately reflect TxGNN’s internal reasoning⁴⁰. Focusing on the most predictive relationships
249 (i.e., edges with importance scores above 0.5, representing an average of 14.9% of edges from the
250 knowledge graph), the model’s performance showed a slight reduction from AUPRC=0.890 (STD:
251 0.006) to AUPRC=0.886 (STD: 0.005). Conversely, when excluding edges deemed predictive by
252 TxGNN and considering the remaining irrelevant relationships (i.e., edges with importance scores
253 below 0.5, accounting for an average of 85.1% of edges), the predictive performance significantly
254 dropped from AUPRC=0.890 (STD: 0.006) to AUPRC=0.628 (STD: 0.026).

255 To assess the quality of TxGNN’s explanations, we employed three established metrics:
256 (1) insertion, which measures predictive performance using only the top K% of edges ranked
257 highest by explanation weight; (2) deletion, which assesses performance after removing the top
258 K% of edges considered most explainable; (3) stability, which evaluates the consistency of ex-
259 planation weights through Pearson’s correlation before and after introducing random perturba-
260 tions to the knowledge graph. We included experiments with three graph explainability methods:
261 GNNExplainer⁴¹, Integrated Gradients⁴², and Information Bottleneck⁴³. As shown in Figure S16,
262 the top-ranked explainable edges are crucial, significantly impacting performance when either re-
263 moved from or inserted into a graph. The performance remained consistent across all insertion
264 and deletion percentages. Additionally, TxGNN Explainer demonstrated the most stable expla-
265 nation weights under various levels of knowledge graph perturbation. These analyses confirm
266 that TxGNN’s multi-hop explanations capture elements of the knowledge graph most critical for
267 making accurate predictions.

268 **TxGNN Explainer supports the human-centric evaluation of therapeutic candidates.** To ex-
269 amine the utility of TxGNN’s multi-hop interpretable explanations for human expert evaluations,
270 we conducted a pilot human study with clinicians and scientists (see Figure S17 for the study inter-
271 face). The study participants included five clinicians, five clinical researchers, and two pharmacists
272 (7 males, 5 females, mean age=34.3, Figure 4c). The user study took around 65 minutes in average,
273 including the assessment of drug-disease indication predictions from TxGNN, a usability ques-
274 tionnaire, and a semi-structured interview. For assessing drug-disease indication predictions, these
275 participants were asked to assess 16 predictions from TxGNN, 12 of which were accurate. For
276 each prediction, we recorded participants’ assessment accuracy, exploration time, and confidence

277 scores, totaling 192 trials (16 predictions \times 12 participants).

278 In evaluating the drug repurposing candidates, participants reported a significant improve-
279 ment in both accuracy (+46%, $p = 0.0443 < 0.05$) and confidence (+49%, $p = 0.0041 < 0.05$)
280 when provided with explanations. Participants took more time to think ($p = 0.0014$) to contextu-
281 alize TxGNN's explanations with their domain expertise, which led to more confident decisions
282 (confidence +49%, $p = 0.0041 < 0.05$). When using TxGNN Explainer, participants are more
283 accurate in evaluating the correctness of drug repurposing predictions than using TxGNN predic-
284 tions alone (accuracy +46%, $p = 0.0443 < 0.05$; Tables S6 and S7).

285 In the post-task questionnaires and interviews, participants reported greater satisfaction when
286 using TxGNN Explainer compared to the baseline (Figure 4e), with 11/12 (91.6%) agreeing
287 or strongly agreeing that the predictions and explanations provided by TxGNN were valuable.
288 In contrast, without explanations, 8/12 (75.0%) disagreed or strongly disagreed with relying on
289 TxGNN's predictions. Participants expressed significantly more confidence in correct predic-
290 tions made by TxGNN when the TxGNN Explainer was included ($t(11) = 3.64, p < 0.01$,
291 using a two-sided Tukey's honestly significant difference test⁴⁴). Some participants indicated that
292 multi-hop interpretable explanations were helpful when examining molecular target interactions
293 identified by TxGNN Explainer and guiding evaluations of potential adverse drug events.

294 **Alignment between TxGNN's drug repurposing predictions and medical evidence.** For three
295 rare diseases, we investigated whether predicted drugs and their multi-hop explanations align with
296 medical reasoning. The evaluation protocol was structured into three stages (Figure 5a). Initially,
297 a human expert queried TxGNN Predictor to identify drugs potentially repurposable for a spe-
298 cific disease. The TxGNN Predictor provided a candidate drug, specifying the confidence in the
299 prediction and its comparative ranking against other candidates. Subsequently, the TxGNN Ex-
300 plainer was queried to elucidate why the selected drug was considered for repurposing. This model
301 revealed its rationale through multi-hop interpretable paths linking the disease to the drug via in-
302 termediate biological interactions. In the final stage, independent medical evidence was collected
303 and analyzed to verify the model's predictions and explanations.

304 First, we examined TxGNN's predictions for Kleefstra syndrome, a disease with a preva-
305 lence of less than one in a million. The condition is attributed to mutations in the EHMT1 gene,
306 leading to pronounced speech development delays, autism spectrum disorder, and childhood hy-
307 potonia. Kleefstra syndrome often features underdeveloped brains with many dormant neuronal
308 pathways. On querying TxGNN Predictor, it recommended Zolpidem as the number one drug

309 repurposing candidate (Figure 5b). At first, this seemed like it would worsen the underdeveloped
310 brains since Zolpidem is commonly used as a sedative and has an inhibitory effect on GABA-A
311 receptors (gene GABRG2) in the brain. TXGNN Explainer's pathways proposed that Zolpidem's
312 action on GABRG2 could reduce autism susceptibility and enhance prefrontal cortex function-
313 ing. Surprisingly, we found that Zolpidem has also demonstrated unexpected stimulative effects
314 in various neurological conditions. For various neurodevelopmental disorders, Zolpidem has been
315 observed to temporarily awaken underactive neurons, offering a potential therapeutic avenue⁴⁵.
316 This paradoxical improvement in neuronal activity can lead to enhancements in speech, motor
317 skills, and alertness in individuals with severe brain injuries or neurodevelopmental disorders, as
318 supported by anecdotal evidence and a handful of clinical studies^{46,47}. TXGNN's prediction and
319 explanatory rationale are both aligned with medical evidence about the paradoxical mechanism of
320 action for Zolpidem, despite none of these clinical cases being directly encountered by the model
321 during training.

322 Next, we explored TXGNN's prediction of Tretinoin for Ehlers-Danlos syndrome, a rare
323 connective tissue disorder that affects 1-9 individuals per 100,000. This disorder arises from muta-
324 tions in collagen-coding genes (such as COL1A1 and COL1A2) and is marked by impaired wound
325 healing and the development of atypical scars. TXGNN Predictor ranks Tretinoin as the number
326 one drug repurposing candidate for Ehlers-Danlos syndrome. Tretinoin, a vitamin A derivative
327 commonly used for acne treatment, is transported by albumin (ALB) and targets ALDH1A2 to
328 mitigate collagen loss and inflammation. Both of these members of Tretinoin's mechanism of ac-
329 tion occur in TXGNN's predictive rationale for this prediction (seen in Figure 5c), indicating that
330 TXGNN's predictive rationale is aligned with medical reasoning. Tretinoin may help in Ehlers-
331 Danlos syndrome by potentially enhancing wound healing and improving the appearance of scars
332 due to its ability to stimulate collagen production in the skin. Further, some subtypes of Ehlers-
333 Danlos syndrome have been associated with a pathogenic mutation in the ALB gene in Landrum
334 et al.⁴⁸ and weakly linked to ALDH1A1 in Javed et al.⁴⁹. In this case, TXGNN Explainer's rea-
335 soning about the pathways that connect Tretinoin to Ehlers-Danlos syndrome was congruent with
336 contemporary clinical evidence.

337 In the final example, we looked at a rare condition, nephrogenic syndrome of inappropriate
338 antidiuresis (NSIAD). This disease is characterized by water and sodium imbalance caused by a
339 mutation in the AVPR2 gene. Patients with congestive heart failure face similar fluid retention
340 challenges, and congestive heart failure has been strongly associated with both AVPR2 and NPR1

341 genes⁵⁰⁻⁵². TxGNN Predictor identified Amyl Nitrite among the top 5 therapeutic candidates
342 (Figure 5d). TxGNN Explainer proposed that the relationship between NSIAD and Amyl Nitrite
343 passes through AVPR2, congestive heart failure, and NPR1. As per medical literature, the AVPR2
344 and NPR1 genes play pivotal roles in regulating fluid and electrolyte balance via complementary
345 but distinct pathways. AVPR2 contributes to water retention and urine concentration, whereas
346 NPR1 facilitates vasodilation, lowers blood pressure, and enhances water excretion⁵³. Enhancing
347 NPR1 activity could counteract the excessive water reabsorption caused by the malfunctioning
348 AVPR2 receptors in NSIAD patients. Amyl Nitrite, which targets the NPR1 gene, emerges as a
349 potential therapeutic option for NSIAD, confirming consistency of TxGNN's explanations with
350 medical evidence. We share TxGNN drug repurposing predictions and explanations for 17,080
351 diseases at <http://txgnn.org>.

352 **Evaluation of TxGNN's predictions using medical records from a large healthcare system.**

353 TxGNN's remarkable performance in previous evaluations suggests that its novel predictions—*i.e.*,
354 therapies not yet FDA-approved for a disease but ranked highly by TxGNN—may hold significant
355 clinical value. As these therapies have not yet been approved for treatment, there is no established
356 gold standard against which to validate them. Recognizing the longstanding clinical practice of
357 off-label drug prescription, we used the enrichment of disease-drug pair co-occurrence in a health
358 system's electronic health records as a proxy measure of being a potential indication. From the
359 Mount Sinai Health System medical records, we curated a cohort of 1,272,085 adults with at least
360 one drug prescription and one diagnosis each (Figure 6a). This cohort was 40.1 percent male, and
361 the average age was 48.6 years (STD: 18.6 years). The demographic breakdown is in Figure 6b-c.
362 Diseases were included if at least one patient was diagnosed with it, and drugs were included if
363 prescribed to a minimum of ten patients (Table 2 and Methods 4), resulting in a broad spectrum of
364 480 diseases and 1,290 drugs as illustrated in Figure 6d.

365 Across these medical records, we measured disease-drug co-occurrence enrichment as the
366 ratio of the odds of using a specific drug for a disease to the odds of using it for other diseases.
367 We derived 619,200 log-odds ratios (log-ORs) for each drug-disease pair. We found that FDA-
368 approved drug-disease pairs exhibited significantly higher log-ORs than other pairs (Figure 6e).
369 Contraindications represented a potential confounding factor in this analysis because adverse drug
370 events could increase the co-occurrence between drug-disease pairs. However, in our study of
371 contraindications, we found no significant enrichment in the co-occurrence of drug-disease pairs,
372 which suggested that adverse drug effects were not a major confounding factor.

373 For each disease in the electronic health records, TXGNN produced a ranked list of potential
374 therapeutic candidates. We omitted drugs already linked to the disease, categorized the remain-
375 ing novel candidates into top-1, top-5, top-5%, and bottom-50%, and calculated their respective
376 mean log-ORs (Figure 6f). We found that the top-1 novel TXGNN prediction had, on average, a
377 107% higher log-OR than the mean log-OR of the bottom-50% predictions. This suggested that
378 TXGNN's top candidate had much higher enrichment in the medical records and, thereby, had
379 a greater likelihood of being an appropriate indication. In addition, the log-OR increased as we
380 broadened the fraction of retrieved candidates, suggesting that TXGNN's prediction scores were
381 meaningful in capturing the likelihood of indication. Although the average log-OR stands at 1.09,
382 the top-1 therapeutic candidate predicted by TXGNN had a log-OR of 2.26, approaching the av-
383 erage log-OR of 2.92 for FDA-approved indications, indicating the enrichment of off-label drug
384 prescriptions among TXGNN's top-ranked predictions.

385 Examining TXGNN's predicted drugs for Wilson's disease, a rare disease causing excessive
386 copper accumulation that frequently instigates liver cirrhosis in children (Figure 3g), we observed
387 that TXGNN predicts likelihoods close to zero for most drugs, with only a select few drugs highly
388 likely to be indications. TXGNN ranked Deferasirox as the most promising candidate for Wilson's
389 disease. Wilson's disease and Deferasirox had a log-OR of 5.26 in the medical records, and litera-
390 ture indicates that Deferasirox may effectively eliminate hepatic iron⁵⁴. In a separate analysis, we
391 evaluated TXGNN on ten recent FDA approvals introduced after the knowledge cutoff date (Table
392 S1). TXGNN consistently ranked newly introduced drugs favorably and, in two instances, placed
393 the newly approved drugs within the top 5% of predicted drugs.

394 Discussion

395 Drug repurposing has been embraced as a drug discovery approach to address the major produc-
396 tivity issues of cost, time to market, and the inherent risks of developing entirely new drugs. While
397 the conventional 'one disease—one model' approach has been utilized in drug repurposing efforts
398 to enhance success rates, the majority of successful drug repurposing cases have resulted from
399 unexpected findings in clinical and preclinical in vivo settings. We propose that a comprehensive
400 way to reposition drugs is to find new indications through multi-disease predictive models. Yet,
401 existing predictive models are based on the assumption that, for a disease, some drugs already
402 exist for it or that drugs already exist for closely related diseases. This overlooks the vast array
403 of diseases—92% of the 17,080 diseases we analyzed—lacking such pre-existing indications and

404 known molecular target interactions. Addressing the needs of these diseases, many of which are
405 complex, neglected, or rare, is a top clinical priority⁵⁵⁻⁵⁷. We define this challenge as zero-shot
406 drug repurposing.

407 We introduce TxGNN, a geometric deep learning model that addresses this problem head-
408 on, specifically targeting diseases with limited molecular understanding and no treatment av-
409 enues. TxGNN achieves state-of-the-art performance in drug repurposing by leveraging a network
410 medicine principle that focuses on disease-treatment mechanisms¹⁵. When asked to suggest ther-
411 apeutic candidates for a disease, TxGNN identifies diseases with shared pathways, phenotypes,
412 and pathologies, extracts relevant knowledge, and fuses it back into the disease of interest. By
413 effectively capturing these latent relationships between diseases, TxGNN can generalize to dis-
414 eases with few treatment options and perform zero-shot inference for unseen diseases. The design
415 behind TxGNN that enables effective zero-shot drug repurposing can be adapted to a wide range
416 of problems, such as disease-target identification and phenotype modeling.

417 TxGNN Predictor is a unified model for indication and contraindication prediction across
418 17,080 diseases. It satisfies an early drug repurposing approach as a high-capacity model that
419 is not limited to a single therapeutic area. Our findings suggest that evaluating a large number
420 of approved or development-stage drugs through multi-disease predictive models should yield a
421 larger number of repositioned drug candidates than approaches limited to a single therapeutic area
422 that can produce infrequent hits. It was found that predicted drug candidates are consistent with
423 off-label prescription rates in a large healthcare system. In the limited evaluation using clinical
424 prescription data and human expert assessment, it was found that predicted drugs were aligned with
425 scientific and clinical consensus. While these estimates suggest beneficial therapeutic potential for
426 existing drugs, predicted drugs would need to undergo extensive screening to establish safety and
427 efficacy as well as determine other drug parameters, such as drug dosage and the sequence and
428 timing of treatments.

429 TxGNN Explainer generates multi-hop interpretable explanations, offering rationales for
430 predicted drugs. These rationales can be analyzed to assess if predicted drugs might elicit ad-
431 ditional biological responses, considering the original indication or molecular target interactions
432 identified by TxGNN Explainer. A pilot human evaluation showed that experts could examine pre-
433 dicted drugs and identify failure points more effectively with multi-hop explanations compared to
434 alternative explanation visualizations. These findings confirm the importance of considering clini-
435 cal needs and explainability when integrating machine learning models into discovery workflows⁵⁸.

436 While TxGNN demonstrates promising performance for zero-shot drug repurposing, its ca-
437 pabilities depend on the quality of medical knowledge graphs. These graphs may lack comprehen-
438 sive data on host-pathogen interactions, essential for predicting drug repurposing in infectious dis-
439 eases (Table S1), and information on the pathogenicity of genetic variants, crucial for identifying
440 repurposing opportunities for genetic diseases⁵⁹. Additionally, challenges such as data biases and
441 the potential for outdated information within the knowledge graph must be addressed. Strategies
442 for overcoming these issues include using techniques for continual learning and model editing⁶⁰,
443 and utilizing easily updatable knowledge graphs, as the one used in this study⁹. Another fruitful
444 future direction is using uncertainty quantification techniques to evaluate the reliability of model
445 predictions⁶¹. We also envision integrating patient information with medical knowledge graphs to
446 provide personalized drug repurposing predictions. Our pilot human evaluation engaged a small
447 sample size (N=12) of clinicians and scientists, prioritizing an in-depth analysis with a smaller,
448 more qualified group over a broader study with a larger, potentially less specialized participant
449 pool. While the results were statistically significant and this participant number is considered
450 a common practice for evaluating highly specialized tools^{62,63}, a larger study could incorporate a
451 greater diversity of user expertise. Despite the promising performance of TxGNN's predictions on
452 tests using medical records, confounders might have biased the enrichment scores measured. We
453 conducted a comprehensive evaluation across multiple axes of model performance beyond accu-
454 racy, including evaluation across diverse hold-out datasets, a pilot evaluation with human experts,
455 and a large-scale enrichment analysis using medical records.

456 TxGNN zero-shot drug repurposing model predicts drugs for diseases without FDA-approved
457 treatments and with minimal available knowledge. TxGNN's Explainer enhances the transparency
458 of TxGNN's predictions, fostering trust and aiding human expert evaluations. TxGNN stream-
459 lines drug repurposing prediction, especially when the limited availability of disease-specific datasets
460 hinders drug development. In the quest for cost-effective therapeutic innovations, models like
461 TxGNN highlight the computational potential for novel therapeutic avenues.

462 **Data availability.** TxGNN's website is at <https://zitniklab.hms.harvard.edu/projects/TxGNN>.
463 The knowledge graph dataset is available at [Harvard Dataverse](https://dataverse.harvard.edu/dataset.xhtml?persistentId=doi:10.7910/DVN/IXA7BM) under a persistent identifier <https://doi.org/10.7910/DVN/IXA7BM>. All clinical and electronic medical record data were deidenti-
464 fied, and the Institutional Review Board at Mount Sinai, New York City, U.S., approved the study.
465

466 **Code availability.** Python implementation of the methodology developed and used in the study
467 is available via the project website at <https://zitniklab.hms.harvard.edu/projects/TxGNN>. The code
468 to reproduce results, documentation, and usage examples are at [https://github.com/mims-harvard/](https://github.com/mims-harvard/TxGNN)
469 [TxGNN](https://github.com/mims-harvard/TxGNN). To facilitate the usage of the algorithm, we developed a TxGNN Explainer, a web-based
470 app available at <http://txgnn.org> to access TxGNN's predictions.

471 **Acknowledgements.** K.H., P.C., and M.Z. gratefully acknowledge the support of NIH R01-
472 HD108794, NSF CAREER 2339524, US DoD FA8702-15-D-0001, awards from Harvard Data
473 Science Initiative, Amazon Faculty Research, Google Research Scholar Program, AstraZeneca
474 Research, Roche Alliance with Distinguished Scientists, Sanofi iDEA-iTECH Award, Pfizer Re-
475 search, Chan Zuckerberg Initiative, John and Virginia Kaneb Fellowship award at Harvard Medical
476 School, Aligning Science Across Parkinson's (ASAP) Initiative, Biswas Computational Biology
477 Initiative in partnership with the Milken Institute, and Kempner Institute for the Study of Natural
478 and Artificial Intelligence at Harvard University. P.C. was supported, in part, by the Harvard Sum-
479 mer Institute in Biomedical Informatics. Any opinions, findings, conclusions or recommendations
480 expressed in this material are those of the authors and do not necessarily reflect the views of the
481 funders.

482 **Authors contribution.** P.C. retrieved, processed, and analyzed the knowledge graph. K.H. and
483 P.C. developed and implemented new machine learning methods, benchmarked machine learning
484 models, and analyzed model behavior, all together with M.Z. Q.W. and N.G. implemented the
485 clinician-centered visual explorer of model predictions and performed a user study to evaluate its
486 usability. S.H., A.V., G.N. and B.S.G. performed a validation study examining new predictions
487 of therapeutic use through the electronic health record system. K.H., P.C., Q.W., S.H., A.V., J.L.,
488 G.N., B.S.G., N.G., and M.Z. contributed new analytic tools and wrote the manuscript. All authors
489 discussed the results and contributed to the final manuscript. M.Z. designed the study.

490 **Competing interests.** The authors declare no competing interests.

491 **Inclusion and ethics statement in global research.** We have complied with all relevant ethical
492 regulations. Our research team represents a diverse group of collaborators. Roles and responsibili-
493 ties were clearly defined and agreed upon among collaborators before the start of the research. All
494 researchers were included in the study design, study implementation, data ownership, intellectual
495 property, and authorship of publications. Our research did not face severe restrictions or prohibi-
496 tions in the setting of the local researchers, and no specific exceptions were granted for this study
497 in agreement with local stakeholders. Animal welfare regulations, environmental protection and
498 risk-related regulations, transfer of biological materials, cultural artifacts, or associated traditional
499 knowledge out of the country do not apply to our research. Our research does not result in stigmati-
500 zation, incrimination, discrimination, or personal risk to participants. Appropriate provisions were
501 taken to ensure the safety and well-being of all participants involved. Our team was committed to
502 promoting equitable access to resources and benefits resulting from the research.

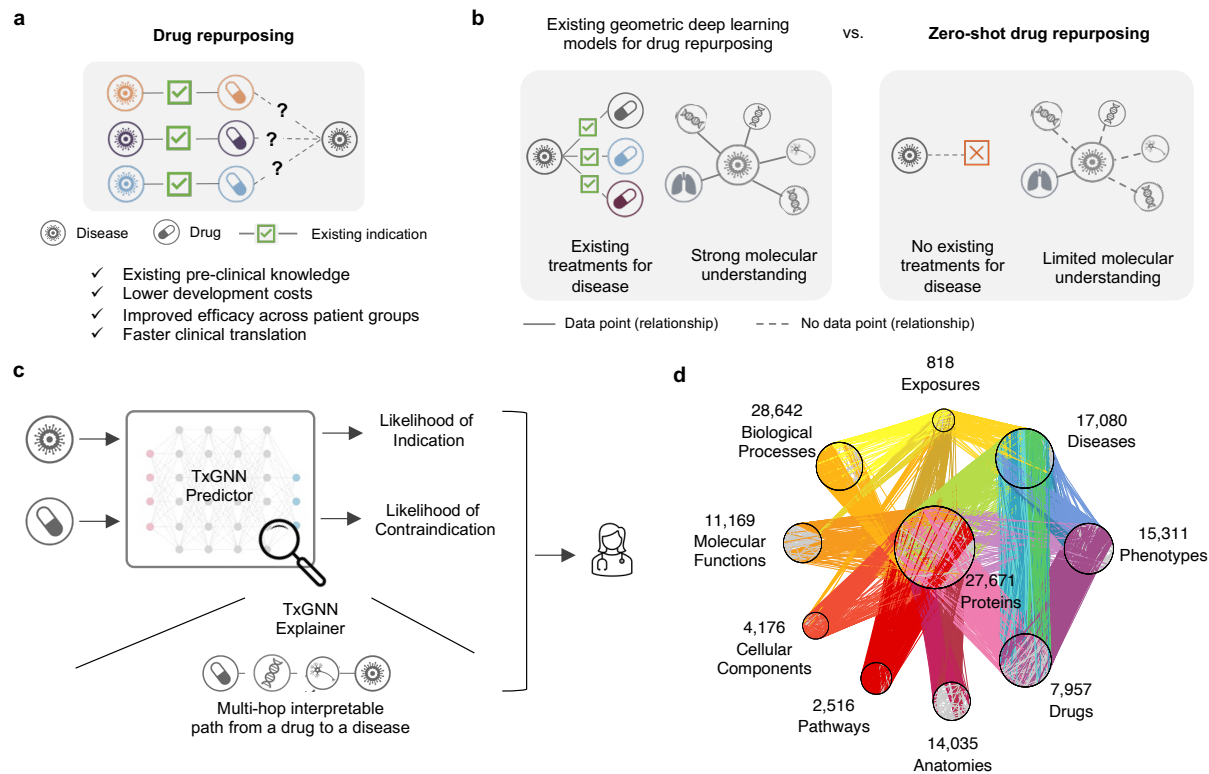


Figure 1: TxGNN is a geometric deep learning approach for drug repurposing across challenging diseases with no known treatments and limited molecular understanding. **a.** Drug repurposing involves exploring new therapeutic applications for existing drugs to treat different diseases. By capitalizing on abundant pre-existing safety and efficacy data, it can dramatically cut down the cost and time to deliver life-saving therapeutics. **b.** Although AI-based drug repurposing has shown promise, its success has been primarily evaluated on diseases with approved treatments and well-understood molecular mechanisms. However, many diseases of critical pharmaceutical interest lack any available treatments (i.e., zero-shot) and exhibit unclear disease mechanisms. These inherent constraints pose challenges to existing AI methods. In this work, we tackle this problem head-on by formulating it as a zero-shot drug repurposing challenge. **c.** TxGNN presents a novel AI framework that generates actionable predictions for zero-shot drug repurposing. TxGNN geometric deep learning model incorporates a vast and comprehensive biological knowledge graph to accurately predict the likelihood of indication or contraindication for any given disease-drug pair. Additionally, TxGNN generates explainable multi-hop paths, facilitating a scientist-friendly understanding of how the prediction is grounded in biological mechanisms in the KG. The combined power of rich predictions and path-based explanations empowers practitioners to prioritize the most promising drug repurposing candidates. **d.** To support our drug repurposing efforts, we develop a large-scale therapeutics-driven knowledge graph that integrates 17 primary data sources. This knowledge graph paints a comprehensive landscape of biological mechanisms across 17,080 diseases and 7,957 repurposable drugs, compiling scientific knowledge for zero-shot drug repurposing endeavors.

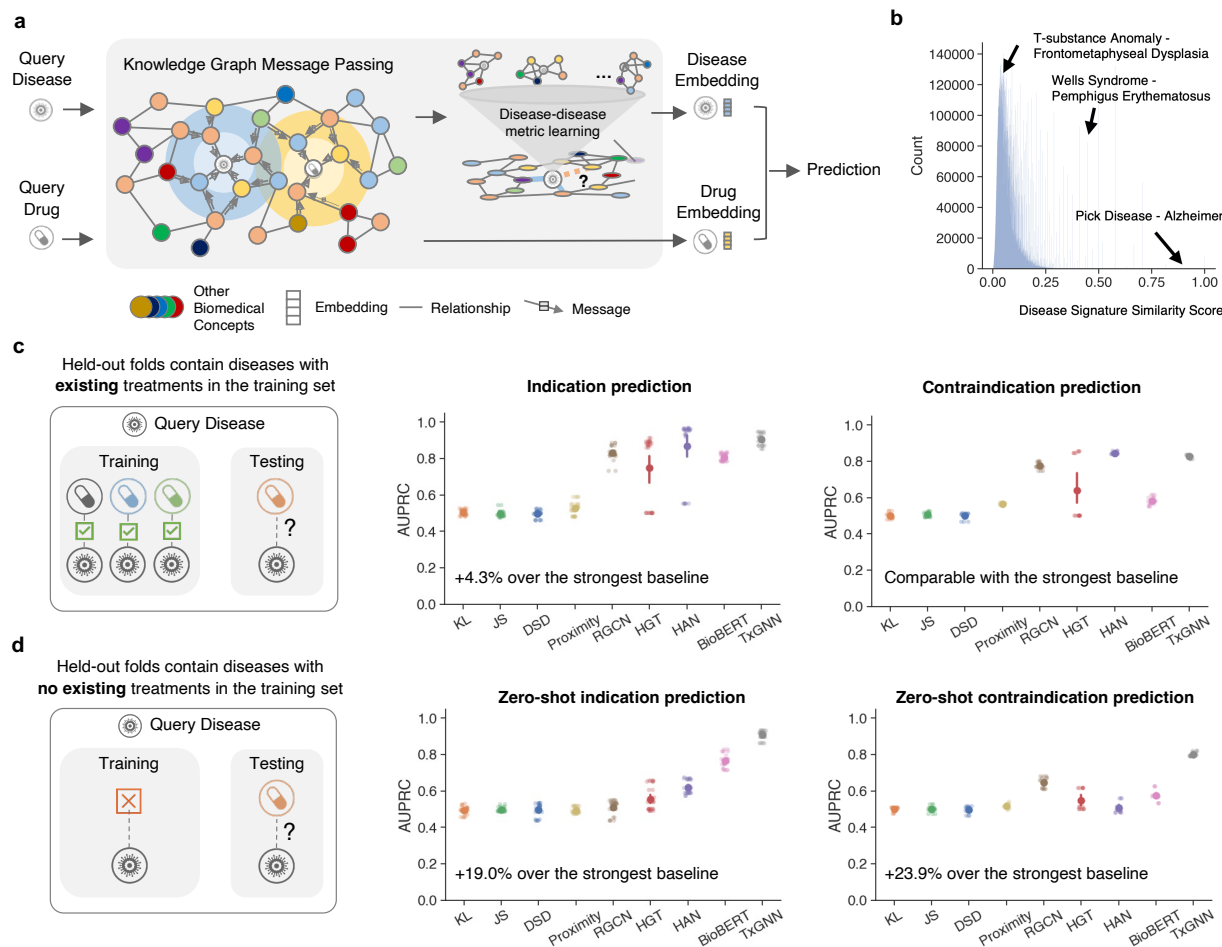


Figure 2: TXGNN predicts indications and contraindications for diseases of no known treatments with high precision. a.

TXGNN is a deep learning model that learns to reason over large-scale knowledge graph on predicting the relationship between drug and disease. In zero-shot repurposing, there is limited indication and mechanism information available for the query disease. Our key insight revolves around the interconnectedness of biological systems. We recognize that diseases, despite their distinctiveness, can exhibit partial similarities and share multiple underlying mechanisms. Based on this motivation, we have developed a specialized module known as disease pooling, which harnesses the power of network medicine principles. This module identifies mechanistically similar diseases and employs them to enhance the information available for the query disease. The disease pooling module has demonstrated significant improvements in the prioritization of repurposing candidates within zero-shot settings. **b.** The TXGNN disease similarity score provides a nuanced and meaningful measure of the relationship between diseases. For instance, disease pairs with low similarity scores, such as T-substance anomaly and frontometaphyseal dysplasia (score: 0.084), indicate a lack of shared mechanisms. Conversely, significant similarity is observed when two diseases receive relatively high scores (>0.2). For instance, Wells syndrome and pemphigus erythematosus exhibit a similarity score of 0.433. Both diseases are skin disorders caused by autoimmune dysregulation, although they differ in phenotypic manifestations, with Wells syndrome characterized by redness and swelling and pemphigus erythematosus characterized by blisters. Moreover, certain disease pairs display exceptionally high similarity scores, such as Pick's disease and Alzheimer's (similarity: 0.909), due to their shared neurological causes. This metric empowers TXGNN to discover similar diseases that can inform and enrich the understanding of query diseases lacking treatment and mechanistic information. **c.** The conventional AI-based repurposing evaluates indication predictions on diseases where the model may have seen other approved drugs during training. In this scenario, we show that TXGNN achieves good performance along with existing methods. **d.** To provide a more realistic evaluation, we introduce a novel setup for assessing zero-shot repurposing, where the model is evaluated on diseases that have no approved drugs available during training. In this challenging setting, we observe a significant degradation in performance for baseline methods. In contrast, TXGNN consistently exhibits robust performance, surpassing the best baseline by up to 19% for indications and 23.9% for contraindications. These results highlight the advanced reasoning capabilities of TXGNN when confronted with query diseases lacking treatment options. The evaluation utilizes the area under the precision-recall curve (AUPRC) and is conducted with five random data splits. The mean performance is highlighted, while the 95% confidence intervals are represented by error bars.

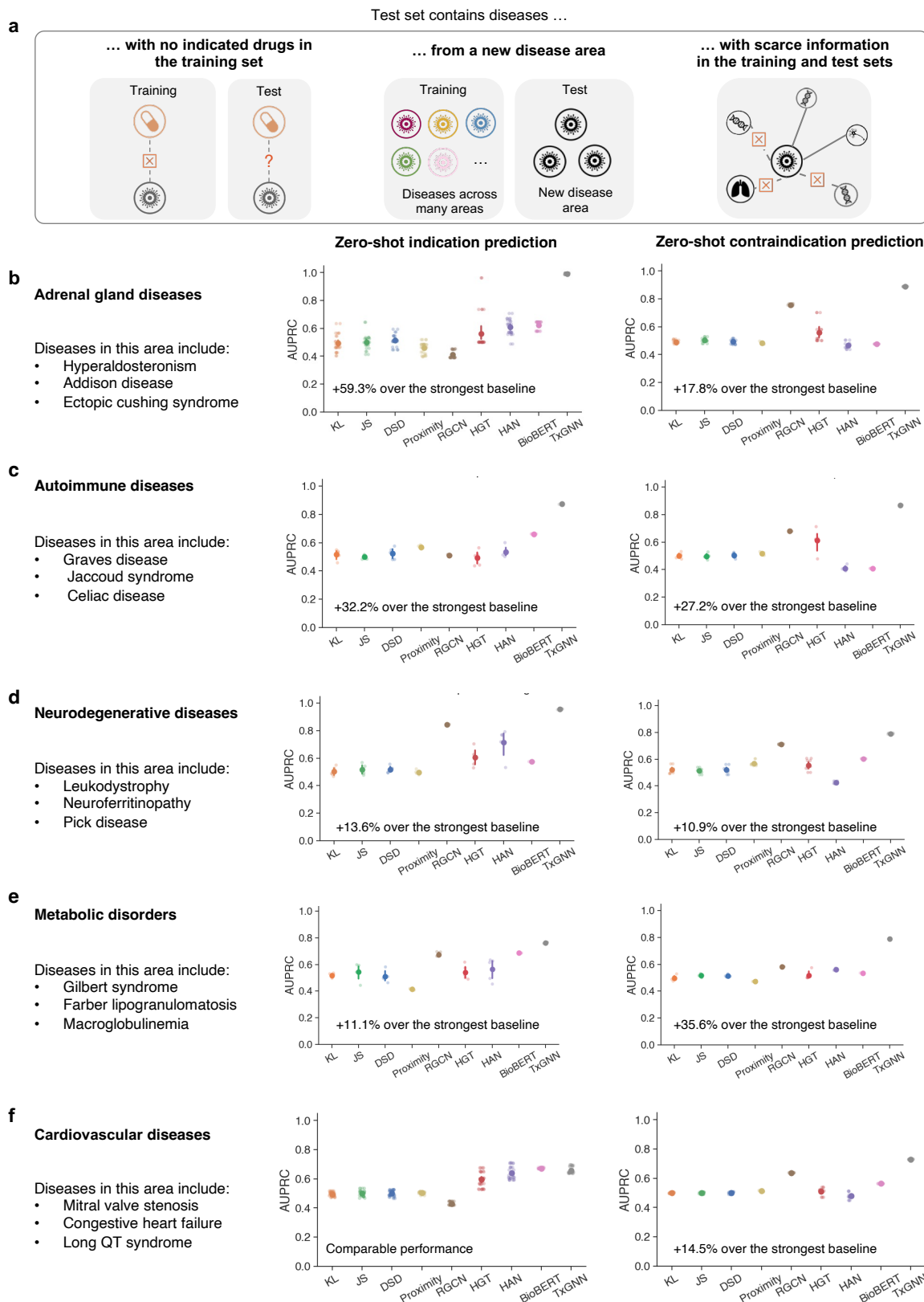


Figure 3: TXGNN accurately predicts therapeutics indications and contraindications across challenging disease areas with limited mechanism understanding. **a.** Zero-shot drug repurposing addresses diseases without any existing treatments and with a dearth of prior biomedical knowledge. We construct a set of ‘disease area’ splits to simulate these conditions. The diseases in the holdout set have (1) no approved drugs in training, (2) limited overlap with the training disease set because we exclude similar diseases, and (3) lack molecular data because we deliberately remove their biological neighbors from the training set. These data splits constitute challenging but realistic evaluation scenarios that mimic zero-shot drug repurposing settings. **b-f.** Holdout folds evaluate diseases related to adrenal glands, autoimmune diseases, neurodegenerative diseases, metabolic disorders, and cardiovascular diseases. Additional four disease areas in anemia, diabetes, cancer, and mental health are provided in Figure S7. Raw scores are provided in Tables S4 and S5. TXGNN shows up to 59.3% improvement over the next best baseline in ranking therapeutic candidates, measured by area under the precision-recall curve.

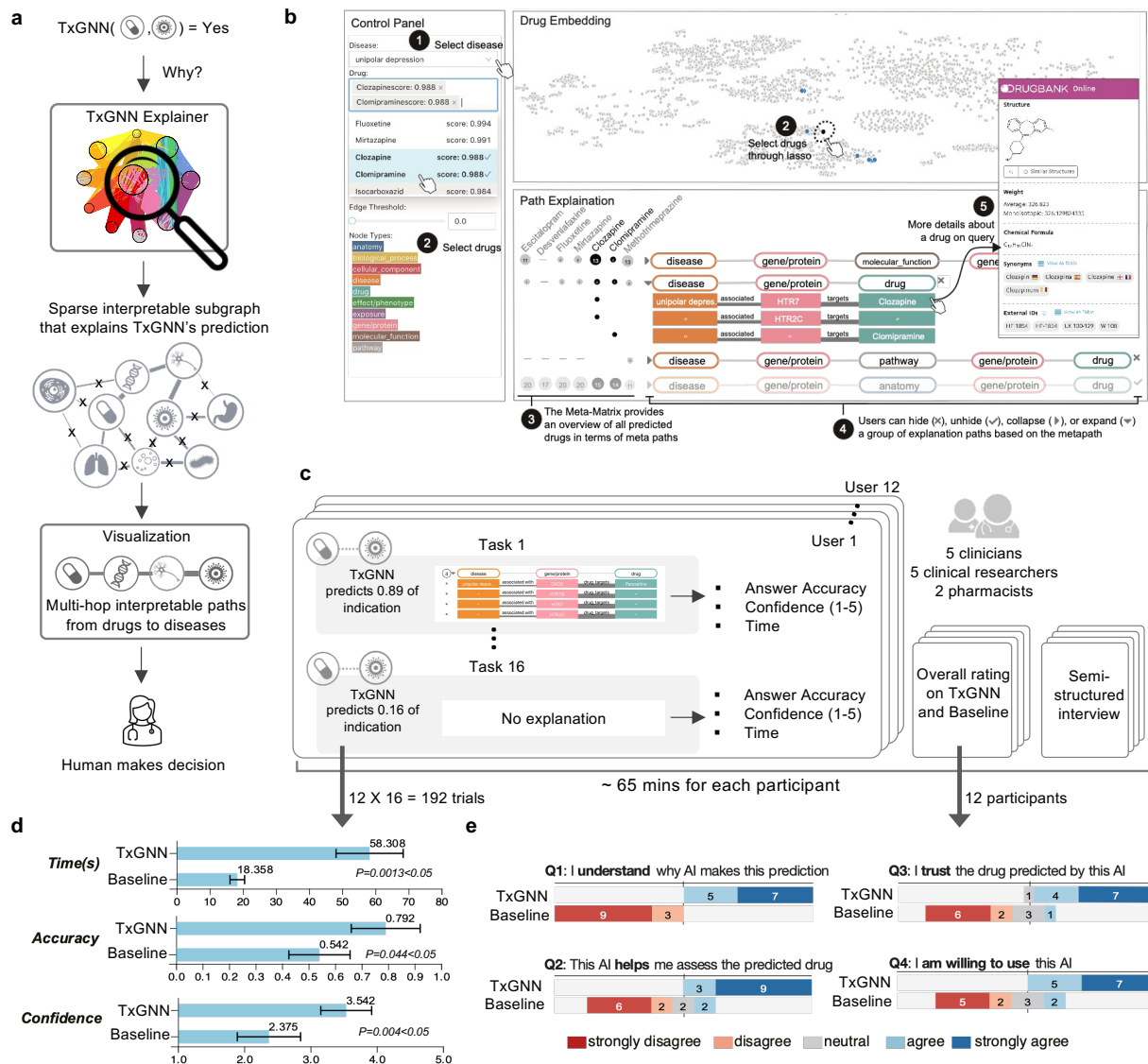


Figure 4: Development, visualization, and evaluation of explanations provided by TxGNN. **a.** Since prediction scores alone are often insufficient for trustworthy deployment of machine learning models, we develop TxGNN Explainer to facilitate adoption by clinicians and scientists. TxGNN Explainer uses state-of-the-art graph explainability techniques to identify a sparse interpretable subgraph that underlies the model's predictions. For each therapeutic candidate, TxGNN Explainer generates a multi-hop pathway composed of various biomedical entities that connect the query disease to the proposed therapeutic candidate. We develop a visualization module that transforms the identified subgraph into these multi-hop paths in a manner that aligns with the cognitive processes of clinicians and scientists. **b.** We design a web-based graphical user interface to support clinicians and scientists in exploring and analyzing the predictions and explanations generated by TxGNN. The 'Control Panel' allows users to select the disease of interest and view the top-ranked TxGNN predictions for the query disease. The 'edge threshold' module enables users to modify the sparsity of the explanation and thereby control the density of the multi-hop paths displayed. The 'Drug Embedding' panel allows users to compare the position of a selected drug relative to the entire repurposing candidate library. The 'Path Explanation' panel displays the biological relations that have been identified as crucial for TxGNN's predictions regarding therapeutic use. **c.** To evaluate the usefulness of TxGNN explanations, we conducted a user study involving 5 clinicians, 5 clinical researchers, and 2 pharmacists. These participants were shown 16 drug-disease combinations with TxGNN's predictions, where 12 predictions were accurate. For each pairing, participants indicated whether they agreed or disagreed with TxGNN's predictions using the explanations provided. **d.** We compared the performance of TxGNN Explainer with a no-explanation baseline in terms of user answer accuracy, task completion time, and user confidence. The results are aggregated on 192 trials (12 participants \times 16 tasks) and reveal a significant improvement in accuracy (+46%) and confidence (+49%) when explanations were provided. Error bars represent 95% confidence intervals. **e.** At the conclusion of the user study, participants were asked qualitative usability questions. Clinicians and scientists agreed that the explanations provided by TxGNN were helpful in assessing the predicted drug-disease relationships and instilled greater trust in the TxGNN's predictions.

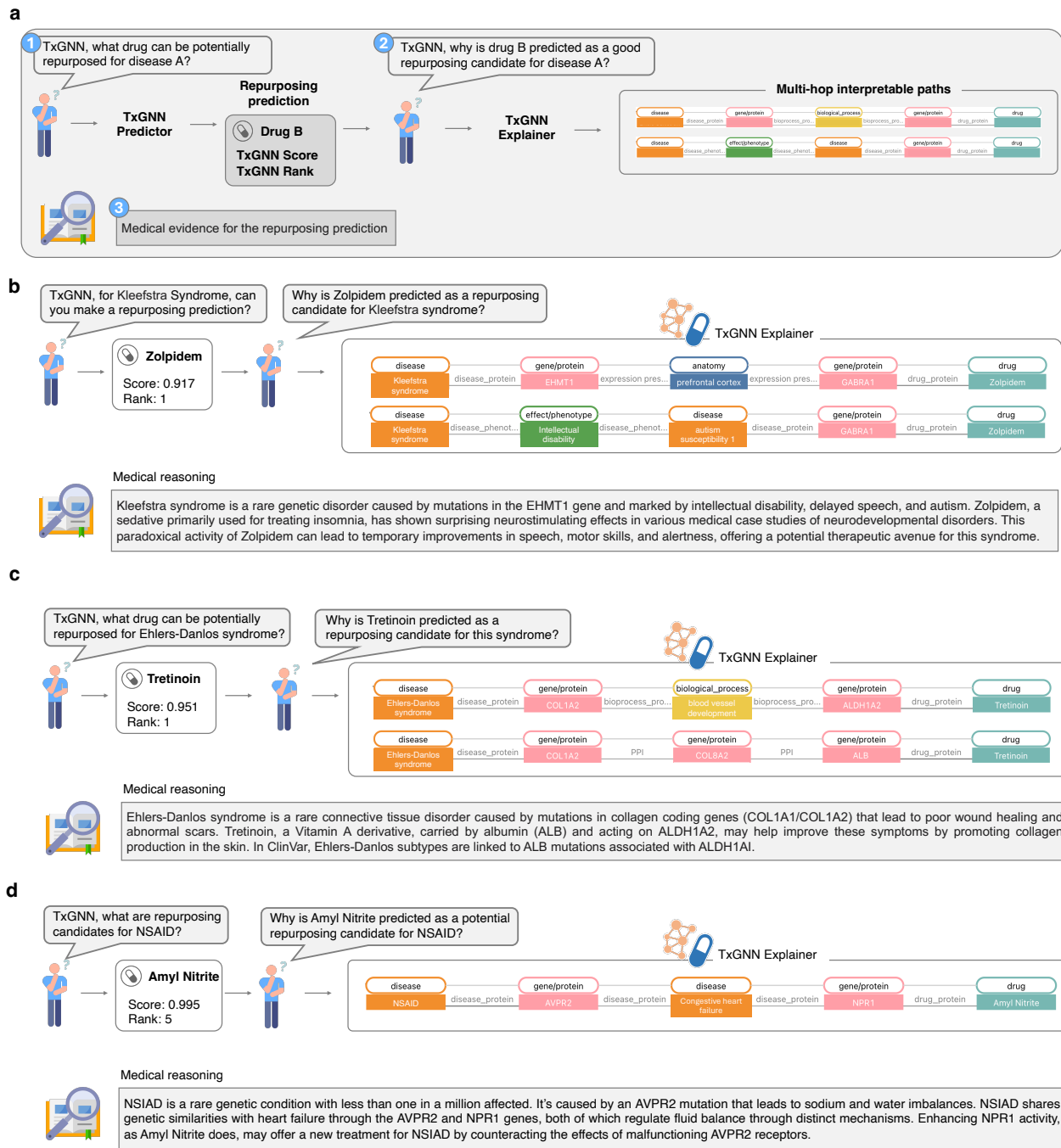


Figure 5: Highlighted cases where interpretable paths produced by TxGNN Explainer align with clinical evidence a. We assess the alignment of drug repurposing candidates identified by TxGNN with established medical reasoning across three rare diseases. The process begins with the TxGNN Predictor, which selects potential drugs for repurposing based on a disease query, and continues with the TxGNN Explorer, which provides interpretable paths explaining the selection. Our case studies conclude with an independent verification of the TxGNN's predictions against clinical knowledge, showcasing the congruence between the TxGNN's recommendations and medical insights. **b.** TxGNN predicts Zolpidem, typically used as a sedative, as a repurposing candidate for Kleefstra syndrome, characterized by developmental delays and neurological symptoms. Despite Zolpidem's conventional inhibitory effects on the brain, TxGNN Explainer suggests its potential to enhance prefrontal cortex activity and improve cognitive functions in those with Kleefstra syndrome. TxGNN's counterintuitive recommendation aligns with emerging clinical evidence of Zolpidem's ability to "awaken" dormant neurons, thereby potentially aiding in speech, motor skills, and alertness in individuals with neurodevelopmental disorders. **c.** TxGNN identifies Tretinoin as the top candidate for treating Ehlers-Danlos syndrome. TxGNN's predictive rationale is rooted in the drug's interactions with albumin (ALB) and ALDH1A2, which aligns with medical insights about Ehlers-Danlos syndrome regarding collagen loss and inflammation mitigation. **d.** TxGNN identifies Amyl Nitrite as a therapeutic option for nephrogenic syndrome of inappropriate antidiuresis (NSIAD). In NSIAD, an AVPR2 mutation leads to water and sodium imbalances. TxGNN Explorer points out the connection between NSIAD and Amyl Nitrite through congestive heart failure, a condition with similar fluid retention issues, by exploring gene interactions (AVPR2 and NPR1) that regulate electrolyte balance.

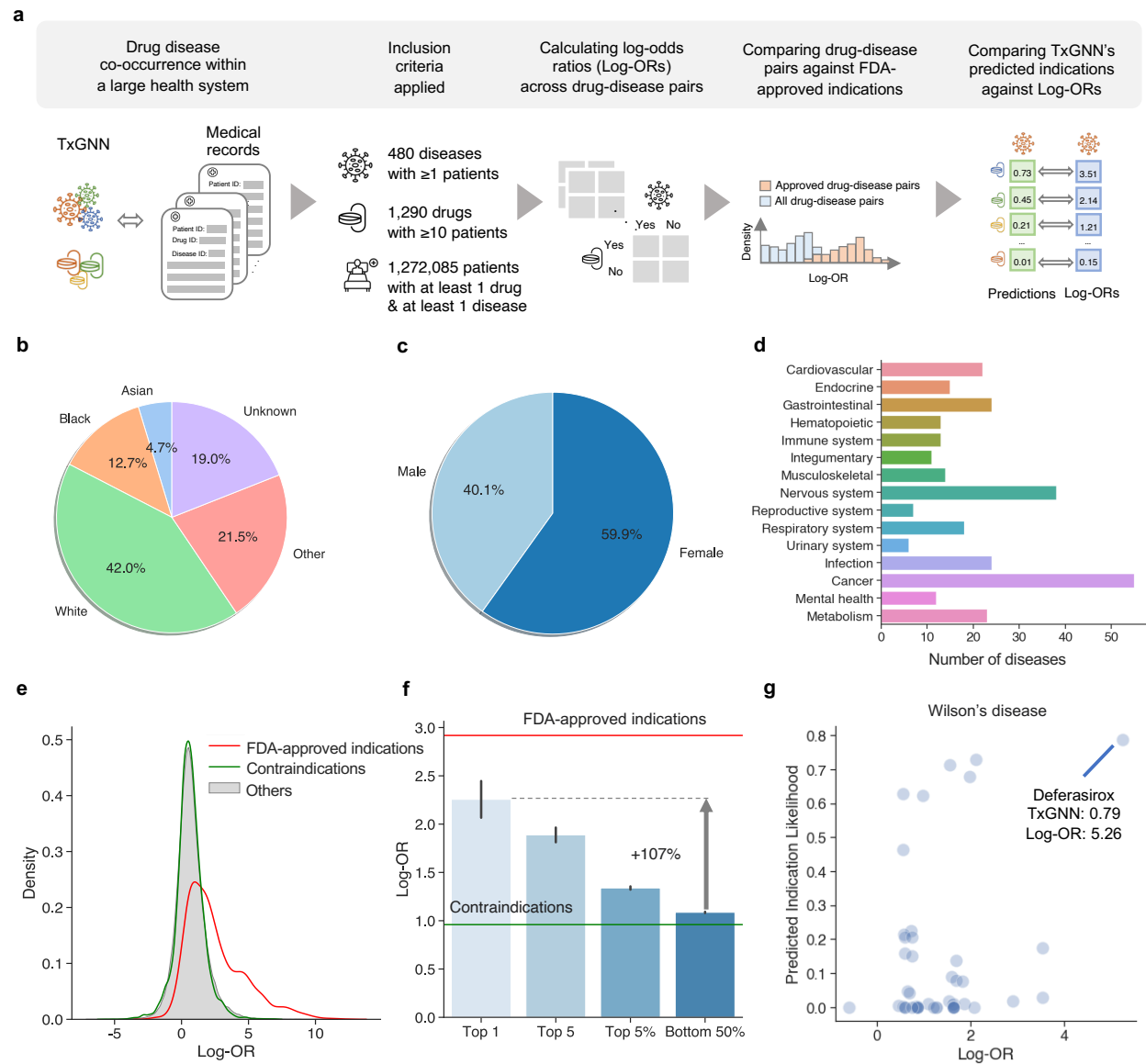


Figure 6: Evaluating TxGNN's novel predictions in a large healthcare system. **a.** We illustrate the steps taken to evaluate TxGNN's novel indications predictions in Mount Sinai's electronic health record (EHR) system. First, we matched the drugs and diseases in the TxGNN knowledge graph to the EHR database, resulting in a curated cohort of 1.27 million patients spanning 480 diseases and 1,290 drugs. Next, we calculated the log-odds ratio (log-OR) for each drug-disease pair, which served as an indicator of the usage of a particular drug for a specific disease. We then validated the log-OR metric as a proxy for clinical usage by comparing drug-disease pairs against FDA-approved indications. Finally, we evaluated TxGNN's novel predictions to determine if their Log-ORs exhibited enrichment within the medical records. **b.** The racial diversity within the patient cohort. **c.** The sex distribution of the patient cohort. **d.** The medical records encompassed a diverse range of diseases spanning major disease areas, ensuring comprehensive coverage and representation. **e.** In validating log-ORs as a proxy metric for clinical prescription, we observed that while the majority of drug-disease pairs exhibited low log-OR values, there was a significant enrichment of log-OR values for FDA-approved indications. Additionally, we noted that contraindications displayed similar log-OR values to the general non-indicated drug-disease pairs, minimizing potential confounders such as adverse drug effects. **f.** We evaluated Log-ORs for the novel indications proposed by TxGNN. The y-axis represents the Log-OR of the disease-drug pairs, serving as a proxy for clinical usage. For each disease, we ranked TxGNN's predictions and extracted the average Log-OR values for the top 1, top 5, top 5%, and bottom 50% of novel drug candidates. The red horizontal line represents the average Log-OR for FDA-approved indications, while the green horizontal line represents the average Log-OR for contraindications. We observed a remarkable enrichment in the clinical usage of TxGNN's novel predictions. The error bar is 95% confidence interval. **g.** We provide a case study of TxGNN's predicted scores plotted against the Log-OR for Wilson's disease. Each point on the plot represents a therapeutic candidate. The top 1 most probable candidate suggested by TxGNN is highlighted, indicating its associated TxGNN score and Log-OR.

503 Online Methods

504 The Methods are structured as follows: 1) curation of knowledge graph dataset (Section 1), 2)
505 description of machine learning approach (Section 2), 3) pilot human evaluation and usability
506 study (Section 3), and 4) evaluation of novel predictions against medical records within a large
507 healthcare system (Section 4).

508 1 Training dataset

509 The knowledge graph is heterogeneous, with 10 types of nodes and 29 types of undirected edges.
510 It contains 123,527 nodes and 8,063,026 edges. Tables S2 and S3 show a breakdown of nodes by
511 node type and edges by edge type, respectively. The knowledge graph and all auxiliary data files are
512 available via Harvard Dataverse at <https://doi.org/10.7910/DVN/IXA7BM>. Supplementary Note
513 S1 provides detailed information about datasets and curation of the knowledge graph.

514 2 Geometric deep learning approach

515 **Notation.** We are given a heterogeneous knowledge graph (KG) $G = (\mathcal{V}, \mathcal{E}, \mathcal{T}_R)$ with nodes in the
516 vertex set $v_i \in \mathcal{V}$, edges $e_{i,j} = (v_i, r, v_j)$ in the edge set \mathcal{E} , where $r \in \mathcal{T}_R$ indicates the relation
517 type, v_i is called the head/source node and v_j is called the tail/target node. Each node also belongs
518 to a node type set \mathcal{T}_V . Each node also has an initial embedding, which we denote as $\mathbf{h}_i^{(0)}$.

519 **Problem definition.** Given a disease i and drug j , we want to predict the likelihood of the drug
520 being (1) indicated for the disease or (2) contraindicated for the disease. Our approach is to induce
521 inductive priors in the model by incorporating factual knowledge from the KG into the model.
522 This process enhances the model's reasoning capabilities to form hypotheses and make predictions
523 about disease treatments.

524 **Experimental setup.** We describe detailed experimental protocols, including data split curation,
525 negative sampling scheme, hyperparameter tuning, and implementation details in Supplementary
526 Note S4.

527 2.1 Overview of TxGNN approach

528 TxGNN is a deep learning approach for mechanistic predictions in drug discovery based on molec-
529 ular networks perturbed in disease and targeted by therapeutics. TxGNN is composed of four
530 modules: (1) a heterogeneous graph neural network-based encoder to obtain biologically mean-
531 ingsful network representation for each biomedical entity; (2) a disease similarity-based metric

532 learning decoder to leverage auxiliary information to enrich the representation of diseases that lack
533 molecular characterization; (3) an all-relation stochastic pre-training followed by a drug-disease
534 centric full-graph fine-tuning strategy; (4) a graph explainability module to retain a sparse set of
535 edges that are crucial for prediction as a post-training step. Next, we expand each module in detail.

536 **2.2 Heterogeneous graph neural network encoder**

537 Our objective is to learn a general encoder of a biomedical knowledge graph by learning a numer-
538 ical vector (embedding) for each node, encapsulating the biomedical knowledge contained within
539 its neighboring relational structures. This involves transforming initial node embeddings using
540 a sequence of local graph-based non-linear function transformations to refine embeddings^{29,64}.
541 These transformations are subject to iterative optimization, guided by a loss function aimed at
542 minimizing the error in therapeutic use predictions. Through this process, the system converges to
543 an optimized set of node embeddings.

544 **Step 1: Initializing latent representations.** We denote the input node embedding \mathbf{X}_i for each
545 node i , which is initialized using Xavier uniform initialization⁶⁵. For every layer l of message-
546 passing, there are the following three stages:

547 **Step 2: Propagating relation-specific neural messages.** For every relation type, first calculates
548 a transformation of node embedding from the previous layer $\mathbf{h}^{(l-1)}$, where the first layer $\mathbf{h}^{(0)} =$
549 \mathbf{X} . This is achieved via applying a relation-specific weight matrix $\mathbf{W}_{r,M}^{(l)}$ on the previous layer
550 embedding:

$$\mathbf{m}_{r,i}^{(l)} = \mathbf{W}_{r,M}^{(l)} \mathbf{h}_i^{(l-1)}. \quad (1)$$

551 **Step 3: Aggregating local network neighborhoods.** For each node v_i , we aggregate on the
552 incoming messages $\{\mathbf{m}_{r,j}^{(l)} | j \in \mathcal{N}_r(i)\}$ from neighboring nodes of each relation r denoted as $\mathcal{N}_r(i)$
553 by taking the average of these messages:

$$\widetilde{\mathbf{m}}_{r,i}^{(l)} = \frac{1}{|\mathcal{N}_r(i)|} \sum_{j \in \mathcal{N}_r(i)} \mathbf{m}_{r,j}^{(l)}. \quad (2)$$

554 **Step 4: Updating latent representations.** We then combine the node embedding from the last

555 layer and the aggregated messages from all relations to obtain the new node embedding:

$$\mathbf{h}_i^{(l)} = \mathbf{h}_i^{(l-1)} + \sum_{r \in \mathcal{T}_R} \widetilde{\mathbf{m}}_{r,i}^{(l)}. \quad (3)$$

556 After L layers of propagation, we arrive at our encoded node embeddings \mathbf{h}_i for each node i .

557 **2.3 Predicting drug-disease relationships**

558 TXGNN employs disease and drug embeddings to predict indications, contra-indications, and off-
559 label use for each disease-drug pair. Considering the three relation types that need prediction,
560 a trainable weight vector \mathbf{w}_r is assigned to each type. The interaction likelihood for a specific
561 relation is then determined using the DistMult approach⁶⁶. Formally, for a disease i , drug j , and
562 relation r , the predicted likelihood p is calculated as follows:

$$p_{i,j,r} = \frac{1}{1 + \exp(-\text{sum}(\mathbf{h}_i \cdot \mathbf{w}_r \cdot \mathbf{h}_j))}. \quad (4)$$

563 **2.4 Embedding-based disease similarity search**

564 Research on diseases varies widely based on factors such as their prevalence and complexity. For
565 instance, the molecular basis of many rare diseases remains poorly understood^{67,68}. Despite this,
566 rare diseases often offer significant opportunities for therapeutic advancements⁶⁹. The limited
567 knowledge surrounding these diseases has heightened the importance of machine learning predic-
568 tions. This shortage of research is evident in the biological knowledge graph, where rare diseases
569 are characterized by a lack of relevant nodes and edges, leading to lower-quality graph embed-
570 dings. For example, diseases without any connections in the knowledge graph are assigned a
571 random initialization for their embedding. Empirical evidence indicates that GNN models exhibit
572 substantially reduced predictive performance on disease-centric splits designed to reflect the sparse
573 nature of knowledge on these diseases, as opposed to random splits (Figure 1g).

574 We posit that the network embeddings generated for these diseases lack significance due to
575 the sparse prior information in the KG. Consequently, there is a necessity for a model to enhance
576 and supplement the network embeddings for these diseases. The underlying principle is that human
577 physiology represents an interconnected system wherein diseases exhibit similarities across various
578 dimensions—e.g., lung cancer and brain cancer are analogous within the cancer disease dimension,
579 while lung cancer and asthma are comparable within the lung disease dimension. Leveraging this

580 concept by utilizing a model to extract relevant information from a group of similar but better-
 581 characterized diseases in the KG, it is possible to enrich the embedding of a target disease, thereby
 582 improving its predictive accuracy.

583 To achieve this, TxGNN employs a three-step procedure: (1) it constructs a disease sig-
 584 nature vector to capture the complex similarities among diseases; (2) it utilizes an aggregation
 585 mechanism to combine the embeddings of similar diseases into a comprehensive auxiliary embed-
 586 ding, which supplements the original disease embedding; (3) it introduces a gating mechanism to
 587 modulate the influence between the original disease embedding and the auxiliary disease embed-
 588 ding, acknowledging that many well-characterized diseases possess adequate embeddings and do
 589 not require supplementation. Each of these steps is elaborated upon in the sections that follow.

590 **Disease signature vectors.** The primary objective of this module is to derive a signature vector
 591 \mathbf{p}_i for each disease i . Given the insufficiency of disease representations produced solely by graph
 592 neural networks in fully capturing the nuances of diseases, these representations are not ideal for
 593 direct similarity computations. Instead, we employ graph theoretical methods¹⁴ to calculate disease
 594 similarities. Additionally, variations of signature vectors are detailed in Supplementary Note S2.
 595 Specifically, we generate a vector that encapsulates the local neighborhoods surrounding a disease.
 596 For disease i , the signature vector is formally defined as follows:

$$\mathbf{p}_i^{\text{AT}} = [p_1 \cdots p_{|\mathcal{V}_P|} \text{ep}_1 \cdots \text{ep}_{|\mathcal{V}_{\text{EP}}|} \text{ex}_1 \cdots \text{ex}_{|\mathcal{V}_{\text{EX}}|} \text{ep}_1 \cdots \text{ep}_{|\mathcal{V}_{\text{EP}}|} d_1 \cdots d_{|\mathcal{V}_D|}], \quad (5)$$

597 where

$$p_j = \begin{cases} 1 & \text{if } j \in \mathcal{N}_i^{\text{P}} \\ 0 & \text{otherwise} \end{cases}, \text{ep}_j = \begin{cases} 1 & \text{if } j \in \mathcal{N}_i^{\text{EP}} \\ 0 & \text{otherwise} \end{cases}, \text{ex}_j = \begin{cases} 1 & \text{if } j \in \mathcal{N}_i^{\text{EX}} \\ 0 & \text{otherwise} \end{cases}, d_j = \begin{cases} 1 & \text{if } j \in \mathcal{N}_i^{\text{D}} \\ 0 & \text{otherwise} \end{cases} \quad (6)$$

598 and $\mathcal{N}_i^{\text{P}}, \mathcal{N}_i^{\text{EP}}, \mathcal{N}_i^{\text{EX}}, \mathcal{N}_i^{\text{D}}$ is the set of gene/protein, effect/phenotype, exposure, diseases nodes lie
 599 in the 1-hop neighborhood of disease i . We also adopt the dot product as the similarity measure,
 600 which means the similarity is the sum of all shared nodes across the four node types:

$$\text{sim}^{\text{AT}}(i, j) = \mathbf{p}_i^{\text{AT}} \cdot \mathbf{p}_j^{\text{AT}} = |\mathcal{N}_i^{\text{P}} \cap \mathcal{N}_j^{\text{P}}| + |\mathcal{N}_i^{\text{EP}} \cap \mathcal{N}_j^{\text{EP}}| + |\mathcal{N}_i^{\text{EX}} \cap \mathcal{N}_j^{\text{EX}}| + |\mathcal{N}_i^{\text{D}} \cap \mathcal{N}_j^{\text{D}}|. \quad (7)$$

601 Given the selected signature for diseases and calculated similarities among the diseases, for

602 a query disease, we can then obtain k most similar diseases for a query disease i :

$$\mathcal{D}_{\text{sim},i} = \operatorname{argmax}_{j \in \mathcal{V}_D}^k \operatorname{sim}(i, j). \quad (8)$$

603 **Disease metric learning.** Given a set of similar diseases, TXGNN generates disease embeddings
604 that integrate various measures of disease similarity into a unified embedding, capable of augment-
605 ing the representation of a query disease that may be sparsely annotated. To achieve this, we adopt
606 a weighted scheme, wherein each disease is weighted according to its similarity score, as follows:

$$\mathbf{h}_i^{\text{sim}} = \sum_{j \in \mathcal{D}_{\text{sim}}} \frac{\operatorname{sim}(i, j)}{\sum_{k \in \mathcal{D}_{\text{sim}}} \operatorname{sim}(i, k)} * \mathbf{h}_j. \quad (9)$$

607 **Gating disease embeddings.** The final stage involves updating the original disease embedding \mathbf{h}_i
608 with the disease-disease metric learning embedding $\mathbf{h}_i^{\text{sim}}$ via a gating mechanism. This mechanism
609 employs a scalar $c \in [0, 1]$ to modulate the influence between these two embeddings. Special con-
610 sideration is needed here because, for diseases that are well-documented in the knowledge graph,
611 the disease-disease metric learning embedding might not be necessary and could potentially skew
612 the final embedding. Conversely, for diseases lacking characterization, the disease-disease met-
613 ric learning embedding is invaluable due to the original embedding’s inadequacy in representing
614 molecular mechanisms. The use of a learnable attention mechanism for deciding whether to pri-
615 oritize the original or augmented embedding is not effective, as it tends to overvalue the original
616 embeddings for well-characterized diseases, thereby neglecting the supplementary embedding. Al-
617 ternatively, we introduce a heuristic algorithm that determines weighting based on the degree of
618 node connectivity $|\mathcal{N}_i^r|$ within the drug-disease relationship being analyzed. A higher degree in-
619 dicates a well-characterized disease, suggesting a reduced reliance on the disease-disease metric
620 learning embedding and vice versa. The scalar’s value is designed to be significantly high for min-
621 imal node degrees (0 or 1) and to decrease rapidly with increasing node degrees. To achieve this
622 gradient, we use an inflated exponential distribution density function with $\lambda = 0.7$:

$$c_i = 0.7 * \exp(-0.7 * |\mathcal{N}_i^r|) + 0.2. \quad (10)$$

623 We observe the result is not sensitive to λ (Figure S6). Finally, we use parameter search and find

624 optimal $\lambda = 0.7$. Then, we can finally obtain an augmented disease embedding:

$$\widehat{\mathbf{h}}_i = c_i * \mathbf{h}_i^{sim} + (1 - c_i) * \mathbf{h}_i. \quad (11)$$

625 Finally, TXGNN uses augmented disease embeddings as input to the latent decoder described in
626 Section 2.3 to produce drug repurposing predictions.

627 2.5 Training TXGNN deep graph models

628 **Objective function.** The objective of the training process is to predict the presence of a relation
629 between two entities within a knowledge graph, which can be viewed as a binary classification
630 task for each relation type. The dataset for positive samples, denoted as \mathcal{D}_+ , comprises all pairs
631 (i, j) across various relation types $r \in \mathcal{T}_R$, with the label $y_{i,r,j} = 1$ indicating the presence of a
632 relation. To generate the dataset for negative samples, \mathcal{D}_- , we use a sampling technique detailed
633 in Supplementary Note S4.3, creating counterparts for each positive pair. For a given pair i, j and
634 relation type r , the model estimates the probability $p_{i,r,j}$ of a relation's existence. The training loss
635 is then calculated using the binary cross-entropy loss formula:

$$\mathcal{L} = \sum_{(i,r,j) \in \mathcal{D}_+ \cup \mathcal{D}_-} y_{i,r,j} * \log(p_{i,r,j}) + (1 - y_{i,r,j}) * \log(1 - p_{i,r,j}). \quad (12)$$

636 Previous research has emphasized knowledge graph completion, optimizing models across the
637 entire spectrum of relations within a knowledge graph⁷⁰. This approach, however, may dilute
638 the model's capacity to capture specific knowledge, particularly when the interest lies solely in
639 drug-disease relations. Given that drug-disease interactions are governed by complex biological
640 mechanisms, the extensive range of biomedical relations in a knowledge graph can offer a com-
641 prehensive view of biological systems. The primary challenge lies in optimizing performance on
642 a select group of relations while beneficially leveraging the broader set of relations for knowledge
643 transfer, avoiding catastrophic forgetting of general knowledge.

644 To address this challenge, TXGNN adopts a pre-training strategy. Initially, during pre-
645 training, TXGNN learns to predict relations across the entire KG using stochastic mini-batching,
646 which helps to encapsulate biomedical knowledge within enriched node embeddings. Subse-
647 quently, in the fine-tuning phase, TXGNN focuses specifically on drug-disease relations. This tar-
648 getted training sharpens the model's ability to generate drug-disease-specific embeddings, thereby

649 optimizing the quality of drug repurposing predictions.

650 **Pre-training.** TxGNN initially undergoes pre-training on millions of biomedical entity pairs
651 spanning the entire array of relations. Given the extensive number of edges, training on the full
652 graph is not computationally viable. Therefore, stochastic mini-batching is employed, allowing
653 for the training on a subset of pairs at each step. This process ensures that each epoch covers
654 all data pairs within the training knowledge graph. During this phase, degree-adjusted disease
655 augmentation is deactivated and all relation types are treated equally. The weights from the pre-
656 trained encoder are subsequently utilized to initialize the encoder model for the fine-tuning phase.
657 It is important to note that the weights in the decoder, specifically for DistMult w_r , are reinitialized
658 prior to fine-tuning to mitigate the risk of negative knowledge transfer.

659 **Fine-tuning.** After the pre-training phase, the model initialization encapsulates a broad spec-
660 trum of biological knowledge. The subsequent phase concentrates on refining the prediction of
661 drug-disease relations. This refinement is achieved by exclusively considering samples of drug-
662 disease pairs (i, j) , where the relation types r fall within the set {indication, contraindication,
663 rev_indication, rev_contraindication}. Other relation types, while not directly included in the train-
664 ing objective, remain part of the knowledge graph to facilitate information flow between drug and
665 disease nodes. During the fine-tuning phase, the model activates the degree-adjusted inter-disease
666 embedding feature. The TxGNN model undergoes both pre-training and fine-tuning in an end-to-
667 end process. The variant that exhibits the highest performance on the validation set is selected for
668 evaluation on the test set and is used for downstream analyses.

669 2.6 Generating multi-hop interpretable explanations

670 In a trained drug repurposing prediction model, consider a target node j and a neighboring source
671 node i connected by an edge $e_{i,j}$ at layer l . For each relation r , intermediate messages $\mathbf{m}_{r,i}^{(l)}$ and
672 $\mathbf{m}_{r,j}^{(l)}$ are computed. These embeddings are concatenated and input into a relation-specific, single-
673 layer neural network parameterized by $\mathbf{W}_{g,r}^{(l)}$. This network predicts the probability of masking the
674 message from source node i during the computation of the target node j 's embedding. The output
675 is processed through a gate, which includes a sigmoid layer to constrain the probability to the range
676 $[0, 1]$, followed by an indicator function that determines whether the edge should be dropped:

$$z_{i,j,r}^{(l)} = \mathbf{1}_{[\mathbb{R}>0.5]} \left(\text{sigmoid} \left(\mathbf{W}_{g,r}^{(l)} \left(\mathbf{m}_{r,i}^{(l)} \parallel \mathbf{m}_{r,j}^{(l)} \right) \right) \right), \quad (13)$$

677 such that $z_{i,j,r}^{(l)} \in [0, 1]$. In practice, a location bias of 3 is added to the sigmoid function during
 678 initialization to ensure that its outputs are initially close to 1. This means that at the start, the gates
 679 remain open, allowing the model to adaptively close the gates and mask edges within the subgraph
 680 as needed. This approach is essential because starting with random initialization, which drops
 681 edges randomly, creates a significant discrepancy between the original and updated predictions.
 682 Consequently, the model’s primary focus shifts towards minimizing this discrepancy rather than
 683 balancing the two objectives. To refine this mechanism, when a gate outputs 0, the corresponding
 684 message is not simply removed. Instead, it is substituted with a learnable baseline vector $\mathbf{b}_r^{(l)}$ for
 685 each relation r and layer l . Therefore, the revised message from source node i to target node j is
 686 represented as follows:

$$\hat{\mathbf{m}}_{i,r}^{(l)} = z_{i,j,r}^{(l)} \cdot \mathbf{m}_{i,r}^{(l)} + (1 - z_{i,j,r}^{(l)}) \cdot \mathbf{b}_r^{(l)}. \quad (14)$$

687 Following the modification of messages with the learnable baseline vector, the process con-
 688 tinues with the standard steps of message aggregation and node embedding updates as described
 689 in Section 2.2. This updated node embedding is then utilized in inter-disease augmentation (Sec-
 690 tion 2.4) and to generate the updated predictions \hat{p} for the interaction between a drug and a disease
 691 (Section 2.3). The optimization of the GraphMask gate weights is guided by two objectives. The
 692 first, faithfulness, aims to ensure that the updated predictions, after applying the mask, align closely
 693 with the initial prediction outcomes. The second objective encourages the model to apply as exten-
 694 sive a masking as feasible. These objectives inherently entail a trade-off: increasing the extent of
 695 masking tends to enlarge the discrepancy between the updated and original predictions. This sce-
 696 nario is addressed through constrained optimization, employing Lagrange relaxation to balance the
 697 objectives. Specifically, the optimization seeks to maximize the Lagrange multiplier λ to enforce
 698 the constraint while simultaneously minimizing the primary objective. The loss function employed
 699 for this purpose is formulated as follows:

$$\max_{\lambda} \min_{\mathbf{W}_g} \sum_{k=1}^L \sum_{(i,j,r) \in \mathcal{D}_+ \cup \mathcal{D}_-} \mathbf{1}_{[\mathbb{R} \neq 0]} z_{i,j,r}^{(k)} + \lambda (\|\hat{p}_{i,j,r} - p_{i,j,r}\|_2^2 - \beta), \quad (15)$$

700 where β is the margin between the updated and original prediction. After the training process is
 701 complete, edges (i, j, r) for which $z_{i,j,r}^{(k)} = 0$ can be eliminated. The remaining edges serve as ex-
 702 planations for the model’s predictions. Additionally, the value computed prior to the application of
 703 the indicator function can be employed to quantify each edge’s contribution to the prediction. This

704 facilitates the adjustment of granular differences in the contributions. More detailed adaptations of
705 the GraphMask approach are discussed in Supplementary Note S3.

706 **3 Pilot usability evaluation of TXGNN with medical experts**

707 The TXGNN Explorer was developed following a user-centric design study process, as outlined in
708 a prior study²⁷. This process involved comparing three visual presentations of GNN explanations
709 from the user's perspective. The findings from this comparison motivated the adoption of path-
710 based explanations, which were preferred based on user feedback. The usability of the TXGNN
711 Explorer was assessed through a comparison with a baseline that only displayed drug predictions
712 and their associated confidence scores.

713 For this usability study, twelve medical experts (7 males and 5 females, average age 34.25,
714 referred to as P1-12) were recruited through personal contacts, Slack channels, and email lists from
715 collaborating institutions, with all participants providing informed consent. The group comprised
716 five clinical researchers (P1-3, P11-12) and five practicing physicians (P4, P7-10), all holding
717 M.D. degrees, and two medical school students with prior experience as pharmacists (P5, P6).
718 Each participant had at least five years of experience in various medical specialties.

719 The study was conducted remotely via Zoom in compliance with COVID-19-related restric-
720 tions. Participants accessed the study system (as shown in Figure S5) using their own computers
721 and shared their screens with the interviewer. The sequence in which predictions were presented,
722 along with the conditions (TXGNN Explorer or the baseline approach), was randomized and coun-
723 terbalanced across participants and tasks.

724 In the drug assessment tasks, participants' accuracy, confidence levels, and task completion
725 times were evaluated across 192 trials (16 tasks × 12 participants). Specifically, participants were
726 tasked with 1) determining the correctness of a drug prediction (i.e., if the drug could potentially
727 be used to treat the disease) and 2) rating their confidence in their decision on a 5-point Likert scale
728 (1=not confident at all, 5=completely confident). The system automatically logged the time taken
729 to evaluate each prediction.

730 Upon completing all predictions, participants provided subjective ratings for both tasks re-
731 garding *Trust*, *Helpfulness*, *Understandability*, and *Willingness to Use*, using a 5-point Likert scale
732 (1=strongly disagree, 5=strongly agree). Subsequent semi-structured interviews yielded insights
733 and feedback on the tool's predictions, explanations, and overall user experience. Each session of
734 the user study lasted approximately 65 minutes.

735 **4 Analysis of medical records from a large healthcare system**

736 Patient data from the Mount Sinai Health System’s electronic health records (EHR) in New York
737 City, U.S., were utilized to examine patterns from predictions in clinical practice. The Mount Sinai
738 Institutional Review Board approved the study, ensuring all clinical data were de-identified. The
739 initial cohort included over 10 million patients, refined to those over 18 years of age with at least
740 one drug and one diagnosis on record, resulting in 1,272,085 patients. This refined cohort com-
741 prised 40.1% males, with an average age of 48.6 years (SD: 18.6 years). The racial composition of
742 the dataset is detailed in Table 2.

743 Disease and medication data were structured according to the Observational Medical Out-
744 comes Partnership (OMOP) standard data model^{71,72}. Predictions were generated for 1,363 dis-
745 eases, identified by training a knowledge graph on 5% of randomly selected drug-disease pairs,
746 serving as a validation set for early stopping. This methodology does not extend to zero-shot per-
747 formance evaluation across all 17,080 diseases, focusing instead on conditions with established
748 indications. Disease names in the prediction dataset were aligned with SNOMED or ICD-10
749 codes and then mapped to OMOP concepts within the Mount Sinai data system. The analysis
750 was restricted to diseases diagnosed in at least one patient, narrowing the focus to 480 conditions.
751 Similarly, medication names were matched to DrugBank IDs, then to RxNorm IDs and OMOP
752 concepts, limiting the scope to medications prescribed to at least one patient, resulting in 1,290
753 medications. Drug-disease pairs were further refined to those with at least one recorded instance
754 of a patient being prescribed the drug for the disease, leading to a final count of 1,236 drugs and 470
755 diseases. Contingency tables were created for each drug-disease pair, and the Fisher exact func-
756 tion from the SciPy library⁷³ was employed to calculate 2-sided odds ratios and p-values for each
757 pair. A two-sided Bonferroni correction was applied to the p-values using the statsmodels Python
758 library’s multi-test function⁷⁴, identifying statistically significant drug-disease pairs as those with
759 $p < 0.005$.

Disease area	Number of diseases	Number of indications	Number of Contraindications
Diseases of cell proliferation	183	854	1007
Mental health diseases	56	213	1038
Cardiovascular diseases	104	300	3131
Diseases of anemia	15	55	545
Adrenal gland diseases	6	33	303
Autoimmune diseases	18	75	319
Metabolic disorders	54	68	523
Diabetes	3	102	364
Neurodegenerative diseases	16	123	134

Table 1: Statistics on disease-area-based dataset splits used to evaluate the zero-shot prediction of therapeutic use. Given all diseases in a given disease area, all indications and contraindications were removed from the dataset used to train machine learning models. Additionally, a fraction (5%) of the connections between biomedical entities to these diseases were removed from the therapeutics-centered knowledge graph. Disease-area splits were curated to evaluate model performance on diseases with limited molecular understanding and no existing treatments.

Racial group	Number of patients	Percent (%)
Asian	60,041	4.7
Black	162,102	12.7
White	534,305	42.0
Unknown	241,998	19.0
Other	273,639	21.5
Total number of patients	1,272,085	100.0

Table 2: Demographics of the electronic health record dataset at Mount Sinai Health System in New York City used to validate TxGNN’s hypotheses on therapeutic use prediction.

References

- 760
761
762 1. Feigin, V. L. *et al.* Burden of neurological disorders across the us from 1990-2017: a global
763 burden of disease study. *JAMA neurology* **78**, 165–176 (2021).
- 764 2. Vetter, N. Editor's choice. *British Medical Bulletin* **93**, 1–5 (2010).
- 765 3. Food, U. & Administration, D. Rare Disease Day 2021. [https://www.fda.gov/news-](https://www.fda.gov/news-events/fda-voices/rare-disease-day-2021-fda-shows-sustained-support-rare-disease-product-development-during-public)
766 [events/fda-voices/rare-disease-day-2021-fda-shows-sustained-support-rare-disease-product-](https://www.fda.gov/news-events/fda-voices/rare-disease-day-2021-fda-shows-sustained-support-rare-disease-product-development-during-public)
767 [development-during-public](https://www.fda.gov/news-events/fda-voices/rare-disease-day-2021-fda-shows-sustained-support-rare-disease-product-development-during-public) (2023). [Online; accessed 19-September-2023].
- 768 4. Pushpakom, S. *et al.* Drug repurposing: progress, challenges and recommendations. *Nature*
769 *Reviews Drug Discovery* **18**, 41–58 (2019).
- 770 5. Abdelsayed, M., Kort, E. J., Jovinge, S. & Mercola, M. Repurposing drugs to treat cardio-
771 vascular disease in the era of precision medicine. *Nature Reviews Cardiology* **19**, 751–764
772 (2022).
- 773 6. Sahragardjoonegani, B., Beall, R. F., Kesselheim, A. S. & Hollis, A. Repurposing existing
774 drugs for new uses: a cohort study of the frequency of FDA-granted new indication exclusivi-
775 ties since 1997. *Journal of Pharmaceutical Policy and Practice* **14** (2021).
- 776 7. Sardana, D. *et al.* Drug repositioning for orphan diseases. *Briefings in Bioinformatics* **12**,
777 346–356 (2011).
- 778 8. Jourdan, J.-P., Bureau, R., Rochais, C. & Dallemagne, P. Drug repositioning: a brief overview.
779 *Journal of Pharmacy and Pharmacology* **72**, 1145–1151 (2020).
- 780 9. Chandak, P., Huang, K. & Zitnik, M. Building a knowledge graph to enable precision
781 medicine. *Scientific Data* **10**, 67 (2023).
- 782 10. Menche, J. *et al.* Uncovering disease-disease relationships through the incomplete interac-
783 tome. *Science* **347** (2015).
- 784 11. Zitnik, M., Feldman, M. W., Leskovec, J. *et al.* Evolution of resilience in protein interactomes
785 across the tree of life. *Proceedings of the National Academy of Sciences* **116**, 4426–4433
786 (2019).
- 787 12. Ruiz, C., Zitnik, M. & Leskovec, J. Identification of disease treatment mechanisms through
788 the multiscale interactome. *Nature Communications* **12**, 1–15 (2021).
- 789 13. Goh, K.-I. *et al.* The human disease network. *Proceedings of the National Academy of Sciences*
790 **104**, 8685–8690 (2007).
- 791 14. Barabási, A.-L., Gulbahce, N. & Loscalzo, J. Network medicine: a network-based approach
792 to human disease. *Nature Reviews Genetics* **12**, 56–68 (2011).
- 793 15. Li, M. M., Huang, K. & Zitnik, M. Graph representation learning in biomedicine and health-
794 care. *Nature Biomedical Engineering* 1–17 (2022).
- 795 16. Gysi, D. M. *et al.* Network medicine framework for identifying drug-repurposing opportuni-
796 ties for covid-19. *Proceedings of the National Academy of Sciences* **118** (2021).
- 797 17. Cao, M. *et al.* Going the distance for protein function prediction: A new distance metric for
798 protein interaction networks. *PLoS ONE* **8**, e76339 (2013).

- 799 18. Cheng, F. *et al.* Network-based approach to prediction and population-based validation of in
800 silico drug repurposing. *Nature Communications* **9** (2018).
- 801 19. Zitnik, M., Agrawal, M. & Leskovec, J. Modeling polypharmacy side effects with graph
802 convolutional networks. *Bioinformatics* **34**, i457–i466 (2018).
- 803 20. Guney, E., Menche, J., Vidal, M. & Barábasi, A.-L. Network-based in silico drug efficacy
804 screening. *Nature Communications* **7**, 1–13 (2016).
- 805 21. Cheng, F., Kovács, I. A. & Barabási, A.-L. Network-based prediction of drug combinations.
806 *Nature Communications* **10**, 1–11 (2019).
- 807 22. Fermaglich, L. J. & Miller, K. L. A comprehensive study of the rare diseases and conditions
808 targeted by orphan drug designations and approvals over the forty years of the orphan drug
809 act. *Orphanet Journal of Rare Diseases* **18**, 1–8 (2023).
- 810 23. Guney, E. Reproducible drug repurposing: When similarity does not suffice. In *Pacific Sym-*
811 *posium on Biocomputing 2017*, 132–143 (World Scientific, 2017).
- 812 24. Avram, S. *et al.* DrugCentral 2021 supports drug discovery and repositioning. *Nucleic Acids*
813 *Research* **49**, D1160–D1169 (2021).
- 814 25. Schlichtkrull, M. S., De Cao, N. & Titov, I. Interpreting graph neural networks for NLP with
815 differentiable edge masking. *International Conference on Learning Representations* (2021).
- 816 26. Lundberg, S. M. & Lee, S.-I. A unified approach to interpreting model predictions. *NeurIPS*
817 **30** (2017).
- 818 27. Wang, Q., Huang, K., Chandak, P., Zitnik, M. & Gehlenborg, N. Extending the nested model
819 for user-centric xai: A design study on gnn-based drug repurposing. *IEEE Transactions on*
820 *Visualization and Computer Graphics* **29**, 1266–1276 (2023).
- 821 28. Cao, M. *et al.* Going the distance for protein function prediction: a new distance metric for
822 protein interaction networks. *PloS one* **8**, e76339 (2013).
- 823 29. Schlichtkrull, M. *et al.* Modeling relational data with graph convolutional networks. In *ESWC*,
824 593–607 (Springer, 2018).
- 825 30. Hu, Z., Dong, Y., Wang, K. & Sun, Y. Heterogeneous graph transformer (2020).
- 826 31. Wang, X. *et al.* Heterogeneous graph attention network (2019).
- 827 32. Lee, J. *et al.* BioBERT: a pre-trained biomedical language representation model for biomedical
828 text mining. *Bioinformatics* btz682 (2019).
- 829 33. Duran-Frigola, M. *et al.* Extending the small-molecule similarity principle to all levels of
830 biology with the chemical checker. *Nature Biotechnology* **38**, 1087–1096 (2020).
- 831 34. Bickel, S., Brückner, M. & Scheffer, T. Discriminative learning under covariate shift. *Journal*
832 *of Machine Learning Research* **10** (2009).
- 833 35. Schölkopf, B. *et al.* On causal and anticausal learning. *ICML* 1255–1262 (2012).
- 834 36. Niven, T. & Kao, H.-Y. Probing neural network comprehension of natural language argu-
835 ments. *Proc. 57th Annual Meeting of the Association of Computational Linguistics* 4658–4664
836 (2019).

- 837 37. Zech, J. R. *et al.* Variable generalization performance of a deep learning model to detect
838 pneumonia in chest radiographs: a cross-sectional study. *PLoS medicine* **15**, e1002683 (2018).
- 839 38. Geirhos, R. *et al.* Shortcut learning in deep neural networks. *Nature Machine Intelligence* **2**,
840 665–673 (2020).
- 841 39. Agarwal, C., Queen, O., Lakkaraju, H. & Zitnik, M. Evaluating explainability for graph neural
842 networks. *Scientific Data* **10** (2023).
- 843 40. Agarwal, C., Zitnik, M. & Lakkaraju, H. Probing GNN explainers: A rigorous theoretical
844 and empirical analysis of gnn explanation methods. In *International Conference on Artificial
845 Intelligence and Statistics*, 8969–8996 (2022).
- 846 41. Ying, Z., Bourgeois, D., You, J., Zitnik, M. & Leskovec, J. Gnnexplainer: Generating expla-
847 nations for graph neural networks. *NeurIPS* **32** (2019).
- 848 42. Sundararajan, M., Taly, A. & Yan, Q. Axiomatic attribution for deep networks. In *ICML*,
849 3319–3328 (PMLR, 2017).
- 850 43. Wang, J. *et al.* Empower post-hoc graph explanations with information bottleneck: A pre-
851 training and fine-tuning perspective. In *KDD*, 2349–2360 (2023).
- 852 44. Tukey, J. W. Comparing individual means in the analysis of variance. *Biometrics* 99–114
853 (1949).
- 854 45. Bomalaski, M. N., Claffin, E. S., Townsend, W. & Peterson, M. D. Zolpidem for the treatment
855 of neurologic disorders: a systematic review. *Jama Neurology* **74**, 1130–1139 (2017).
- 856 46. Boisgontier, J. *et al.* Case report: Zolpidem’s paradoxical restorative action: A case report of
857 functional brain imaging. *Frontiers in Neuroscience* **17**, 1127542 (2023).
- 858 47. Sripad, P. *et al.* Effect of zolpidem in the aftermath of traumatic brain injury: an meg study.
859 *Case reports in neurological medicine* **2020** (2020).
- 860 48. Landrum, M. J. *et al.* Clinvar: improvements to accessing data. *Nucleic acids research* **48**,
861 D835–D844 (2020).
- 862 49. Javed, S. *et al.* Aldh1 & cd133 in invasive cervical carcinoma & their association with the
863 outcome of chemoradiation therapy. *The Indian journal of medical research* **154**, 367 (2021).
- 864 50. Ghoussaini, M. *et al.* Open targets genetics: systematic identification of trait-associated genes
865 using large-scale genetics and functional genomics. *Nucleic acids research* **49**, D1311–D1320
866 (2021).
- 867 51. Goltsman, I. *et al.* Rosiglitazone treatment restores renal responsiveness to atrial natriuretic
868 peptide in rats with congestive heart failure. *Journal of Cellular and Molecular Medicine* **23**,
869 4779–4794 (2019).
- 870 52. Bryan, P. M., Xu, X., Dickey, D. M., Chen, Y. & Potter, L. R. Renal hyporesponsiveness to
871 atrial natriuretic peptide in congestive heart failure results from reduced atrial natriuretic pep-
872 tide receptor concentrations. *American Journal of Physiology-Renal Physiology* **292**, F1636–
873 F1644 (2007).
- 874 53. Wishart, D. S. *et al.* DrugBank 5.0: a major update to the DrugBank database for 2018.
875 *Nucleic Acids Research* **46**, D1074–D1082 (2018).

- 876 54. Seetharaman, J. & Sarma, M. S. Chelation therapy in liver diseases of childhood: Current
877 status and response. *World Journal of Hepatology* **13**, 1552 (2021).
- 878 55. Alsentzer, E. *et al.* Deep learning for diagnosing patients with rare genetic diseases. *medRxiv*
879 2022–12 (2022).
- 880 56. O’Connell, D. Neglected diseases. *Nature* **449**, 157–157 (2007).
- 881 57. Tambuyzer, E. *et al.* Therapies for rare diseases: therapeutic modalities, progress and chal-
882 lenges ahead. *Nature Reviews Drug Discovery* **19**, 93–111 (2020).
- 883 58. Zhang, A., Xing, L., Zou, J. & Wu, J. C. Shifting machine learning for healthcare from
884 development to deployment and from models to data. *Nature Biomedical Engineering* 1–16
885 (2022).
- 886 59. Duffy, Á. *et al.* Development of a human genetics-guided priority score for 19,365 genes and
887 399 drug indications. *Nature Genetics* 1–9 (2024).
- 888 60. Cheng, J., Dasoulas, G., He, H., Agarwal, C. & Zitnik, M. GNNDelete: a general strategy for
889 unlearning in graph neural networks. *International Conference on Learning Representations*
890 (2023).
- 891 61. Huang, K., Jin, Y., Candes, E. & Leskovec, J. Uncertainty quantification over graph with
892 conformalized graph neural networks. *Advances in Neural Information Processing Systems*
893 **36** (2024).
- 894 62. Cai, C. J. *et al.* Human-centered tools for coping with imperfect algorithms during medical
895 decision-making. In *Proceedings of the 2019 chi conference on human factors in computing*
896 *systems*, 1–14 (2019).
- 897 63. Macefield, R. How to specify the participant group size for usability studies: a practitioner’s
898 guide. *Journal of usability studies* **5**, 34–45 (2009).
- 899 64. Gilmer, J., Schoenholz, S. S., Riley, P. F., Vinyals, O. & Dahl, G. E. Neural message passing
900 for quantum chemistry. In *ICML*, 1263–1272 (PMLR, 2017).
- 901 65. Glorot, X. & Bengio, Y. Understanding the difficulty of training deep feedforward neural
902 networks. In *AISTATS*, 249–256 (2010).
- 903 66. Yang, B., Yih, W.-t., He, X., Gao, J. & Deng, L. Embedding entities and relations for learning
904 and inference in knowledge bases. *ICLR* (2015).
- 905 67. Griggs, R. C. *et al.* Clinical research for rare disease: opportunities, challenges, and solutions.
906 *Molecular Genetics and Metabolism* **96**, 20–26 (2009).
- 907 68. Boycott, K. M., Vanstone, M. R., Bulman, D. E. & MacKenzie, A. E. Rare-disease genetics in
908 the era of next-generation sequencing: discovery to translation. *Nature Reviews Genetics* **14**,
909 681–691 (2013).
- 910 69. Thomas, S. & Caplan, A. The orphan drug act revisited. *Jama* **321**, 833–834 (2019).
- 911 70. Lin, Y., Liu, Z., Sun, M., Liu, Y. & Zhu, X. Learning entity and relation embeddings for
912 knowledge graph completion. In *AAAI* (2015).

- 913 71. Stang, P. E. *et al.* Advancing the science for active surveillance: rationale and design for
914 the Observational Medical Outcomes Partnership. *Annals of Internal Medicine* **153**, 600–606
915 (2010).
- 916 72. Klann, J. G., Joss, M. A., Embree, K. & Murphy, S. N. Data model harmonization for the All
917 Of Us Research Program: Transforming i2b2 data into the OMOP common data model. *PloS*
918 *ONE* **14**, e0212463 (2019).
- 919 73. Virtanen, P. *et al.* SciPy 1.0: fundamental algorithms for scientific computing in Python.
920 *Nature Methods* **17**, 261–272 (2020).
- 921 74. Seabold, S. & Perktold, J. Statsmodels: Econometric and statistical modeling with python. In
922 *Proceedings of the 9th Python in Science Conference*, vol. 57 (2010).