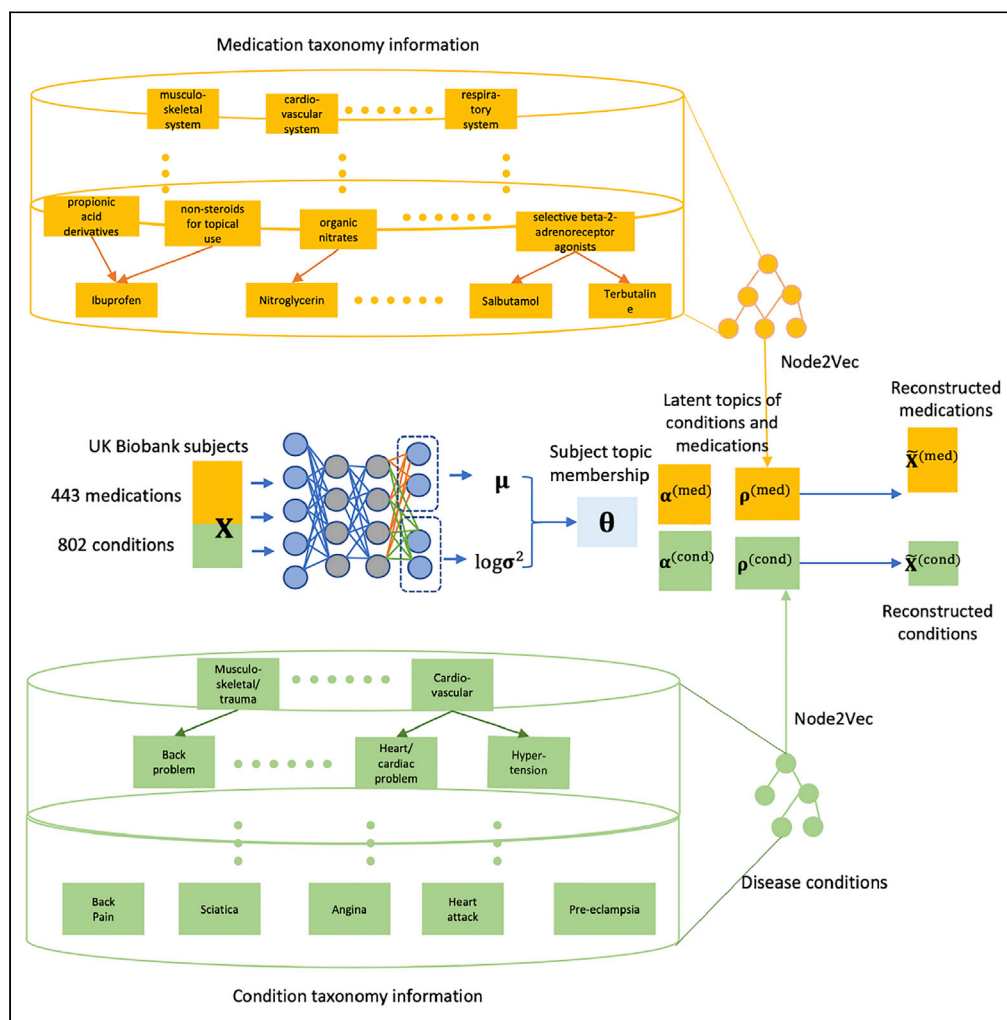


Article

A graph-embedded topic model enables characterization of diverse pain phenotypes among UK biobank individuals



Yuening Wang,
Rodrigo Benavides, Luda Diatchenko, Audrey V. Grant, Yue Li

audrey.grant@mcgill.ca (A.V.G.)
yueli@cs.mcgill.ca (Y.L.)

Highlights

Interpretable deep learning to integrate knowledge graphs and patient data

Modeling phenotypes from self-reports of 457,461 individuals from the UK Biobank

Predicting and characterizing chronic pain phenotypes using latent phenotypes

Potential link between cardiovascular conditions or medications and chronic pain

Wang et al., iScience 25, 104390
June 17, 2022 © 2022 The Author(s).
<https://doi.org/10.1016/j.isci.2022.104390>



Article

A graph-embedded topic model enables characterization of diverse pain phenotypes among UK biobank individuals

Yuening Wang,¹ Rodrigo Benavides,² Luda Diatchenko,^{3,4,5} Audrey V. Grant,^{3,4,5,*} and Yue Li^{1,6,*}

SUMMARY

Large biobank repositories of clinical conditions and medications data open opportunities to investigate the phenotypic disease network. We present a graph embedded topic model (GETM). We integrate existing biomedical knowledge graph information in the form of pre-trained graph embedding into the embedded topic model. Via a variational autoencoder framework, we infer patient phenotypic mixture by modeling multi-modal discrete patient medical records. We applied GETM to UK Biobank (UKB) self-reported clinical phenotype data, which contains 443 self-reported medical conditions and 802 medications for 457,461 individuals. Compared to existing methods, GETM demonstrates good imputation performance. With a more focused application on characterizing pain phenotypes, we observe that GETM-inferred phenotypes not only accurately predict the status of chronic musculoskeletal (CMK) pain but also reveal known pain-related topics. Intriguingly, medications and conditions in the cardiovascular category are enriched among the most predictive topics of chronic pain.

INTRODUCTION

The advent of electronic health records (EHR) has started to transform the way healthcare data are recorded and used in clinical practice and in research settings. Besides free-form clinical text, most modern healthcare centers now routinely collect structured EHR data describing comprehensive aspects of care, including diagnosis, medications, treatments, laboratory test results, and other measures. Deriving coherent phenotypes from these EHR data is crucial in downstream phenome-wide association analyses (PheWAS) and may greatly improve the power of detecting genetic associations using the genome-wide association (GWA) approach (McCoy et al., 2017). Besides better characterizing known phenotypes (e.g., through comorbidities or the demographics of age-of-onset), mining phenomic data also has the potential to reveal novel combinations of diseases and other variables of potential etiological interest. This will help identify specific strata of study subjects most at risk for disease or targeted for specific drug recommendations. Despite these promises, clinical phenotype data sources remain underused (Jensen et al., 2012). As genotype and deep phenotype data become increasingly available through consortia or large government-funded cohorts such as the UK Biobank (UKB) (Bycroft et al., 2018) and the 100,000 Genomes Project (Turnbull et al., 2018), there is an urgent need for an automatic and accurate phenotyping tool to accelerate novel disease comorbidity discoveries and improve the yield of GWA studies for complex phenotypes and diseases in humans.

Among many machine learning approaches, topic models (Blei et al., 2003a) stand out as a particularly well-adapted framework for automatic phenotyping. They are extremely efficient at modeling sparse and discrete data such as text documents. Topic models were originally developed to discover patterns of word usages from corpuses of text documents by accomplishing two related tasks: (1) inferring a set of latent categorical distributions over the vocabulary (i.e., topics); (2) using these latent topic distributions to infer topic mixture memberships of each document, thereby connecting them under similar topical themes. In our context, we consider EHR as our documents and the diagnostic and medication codes as our vocabulary.

Several topic methods were developed recently for effectively mining EHR data (Li et al., 2020; Song et al., 2021). However, most existing topic models are unable to incorporate existing biomedical knowledge

¹School of Computer Science, McGill University, Canada

²Department of Anesthesiology, Centro Nacional de Rehabilitación, San Jose, Costa Rica

³Department of Anesthesia, McGill University, Canada

⁴Faculty of Dentistry, McGill University, Canada

⁵Alan Edwards Centre for Research on Pain, McGill University, Canada

⁶Lead contact

*Correspondence: audrey.grant@mcgill.ca (A.V.G.), yueli@cs.mcgill.ca (Y.L.)
<https://doi.org/10.1016/j.isci.2022.104390>



graphs, which manifest in several forms such as disease taxonomy and drug classification systems. Knowledge Graph Embedding LDA (KGE-LDA) (Yao et al., 2017) models the distribution of the word embedding learned from TransE (Bordes et al., 2013) on a words-by-words relational graph. Latent-feature LDA (Nguyen et al., 2015) and Embedded Topic Model (ETM) (Dieng et al., 2019) use the word embedding to compose the topic distribution. These methods were applied to standard benchmark corpus data and only works with one data modality as opposed to the multimodal patient electronic medical record data (e.g., disease conditions and medications). In addition, except for ETM, the aforementioned topic models use traditional inference algorithms (e.g., Gibbs sampling or mean-field variational inference) to infer topic distributions. Therefore, these models have limited flexibility to capture the non-linear connections between observed EHR codes and the underlying patient phenotypes.

In this paper, we present a graph-embedded topic model (GETM) for learning phenotypes from heterogeneous EHR data by leveraging biomedical graph information. As the main method contribution, GETM seamlessly integrates two existing models: ETM (Dieng et al., 2019) and node2vec (Grover and Leskovec, 2016). Briefly, we first use node2vec to learn the condition and medication embeddings based on their taxonomic information; then, we incorporate these embeddings into the ETM, which tri-factorizes the individuals-by-conditions/medications matrix into individuals-by-topics, topics-by-embedding, and embedding-by-conditions/medications matrices. The distribution of the individuals-by-topics is approximated by the output of a feedforward neural network using the amortized variational inference technique while fixing the embedding-by-conditions/medications matrices to the node2vec-learned node embedding from conditions/medications taxonomical graphs, respectively.

As a proof-of-concept, we applied GETM to UKB phenotype data, where 457,461 individuals of European descent from across the United Kingdom were deeply phenotyped through extensive self-report based questionnaires for about 443 well-defined phenotype conditions and 802 active ingredients of medications (Bycroft et al., 2018). We then turned to a more focused analysis on predicting and characterizing different pain phenotypes. Chronic pain is the result of dysfunction of the nociceptive circuitry leading to sustained perception of pain. Chronic pain is highly prevalent in aging populations affecting up to 50 % of older adults (> 65 years old) (Fayaz et al., 2016a) and decreases the overall mental and emotional health of affected individuals (Dueñas et al., 2016). Using GETM, we provide a refreshing view of pain phenotypes, considered as CMK pain, chronic pain by body site, non-specific acute and chronic pain, and the transition from acute to chronic pain by making use of phenotype data from subsequent visits on a subset of the UKB study population. By correlating the inferred GETM-topics and pain phenotypes across the UKB subjects, we discover not only the known pain-related conditions and medications among the highly predictive topics but also novel combinations after removing labels of known pain-related conditions and analgesics.

RESULTS

Graph embedded topic model overview

GETM models the distribution of medications and conditions as discrete clinical features for each individual. For a given study subject, the expected rate of each feature is determined by both the logistic-normal latent subject's topic mixture and the point-estimate latent topic distributions over the features. The goal of GETM is to approximately infer the distributions of these latent variables. To this end, we carry out an amortized variational inference (Kingma and Welling, 2014) in two steps (Figure 1). In the first step, to infer the topic mixture of a given patient, we provide to a feedforward neural network (i.e., the encoder) the binary vector of the individual's observed discrete features. We then sample the topic mixture from a variational Gaussian distribution with the mean and SD computed by the encoder. In the second step, we decode the sampled topic mixture back to the original conditions and medications. We used two linear decoders each with separate topic/feature embeddings for medications and conditions, respectively. Specifically, each decoder tri-factorizes the individuals-by-features matrix respectively into individuals-by-topics (θ), topics-by-embeddings ($\alpha^{(t)}$), and embeddings-by-features (medication or condition) ($\rho^{(t)}$) matrices ($t \in \{med, cond\}$). Notably, the two linear decoders share the patient-level topic mixture θ , whereas $\alpha^{(t)}$ and $\rho^{(t)}$ are learned separately. Importantly, the medication embedding $\rho^{(med)}$ and condition embedding $\rho^{(cond)}$ are pretrained by node2vec (Grover and Leskovec, 2016) from the taxonomic tree-structured graphs of conditions and medications while the two sets of topic embedding counterparts ($\alpha^{(med)}$, $\alpha^{(cond)}$) are directly learned from the UKB participant data by GETM. This tri-factorization design allows for exploring topics, study subjects and relationships among conditions and medications in a highly interpretable way.

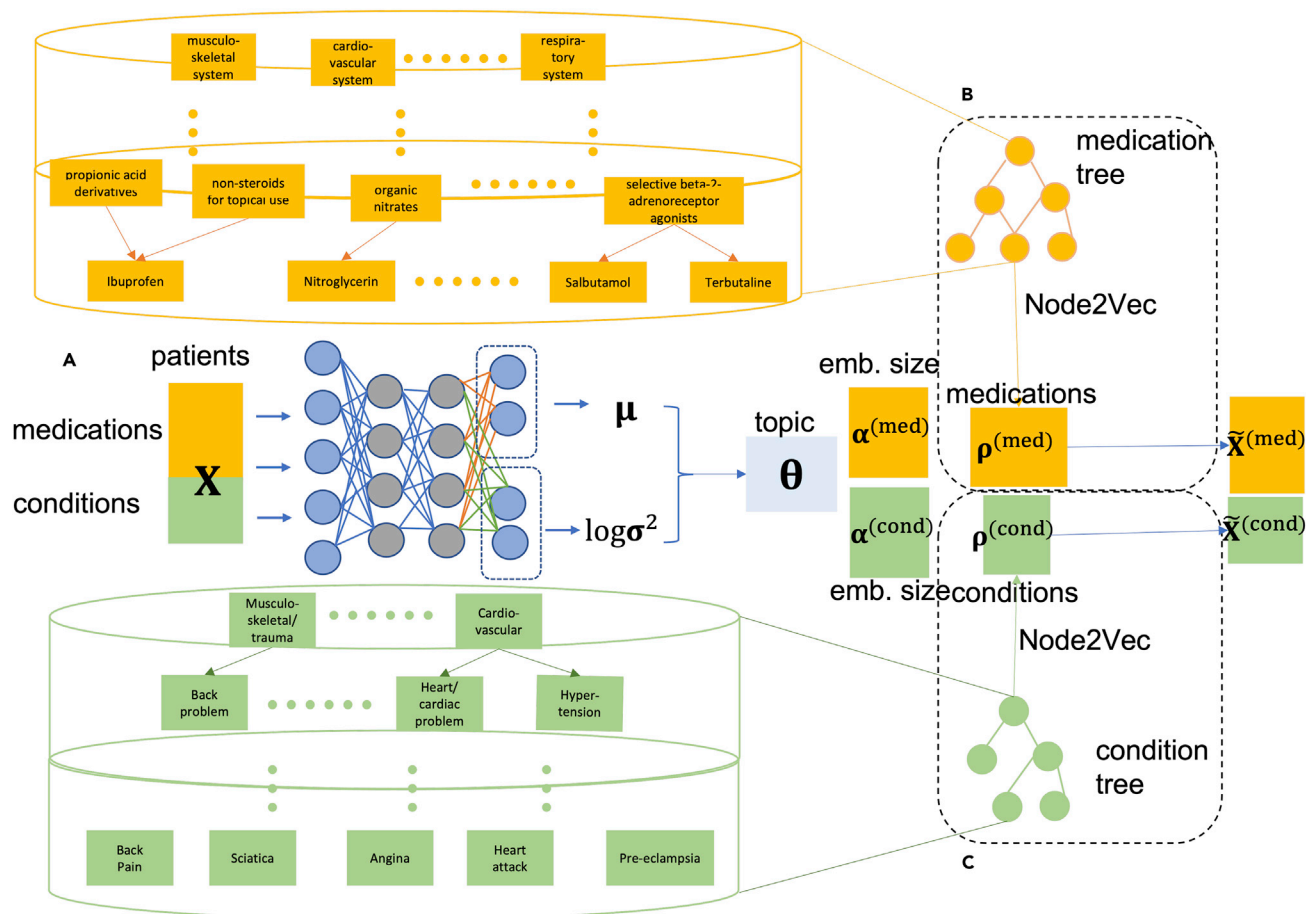


Figure 1. Overview of Graph-embedded topic model and its application on UKB phenotype data

The UKB data consists of 443 conditions and 802 medications for 457,461 individuals. We developed GETM to model these data although our GETM can be applied to other datasets as well.

(A) GETM training. GETM is a variational autoencoder (VAE) model. The neural network encoder takes individuals' condition and medication information as input and produces the variational mean μ and variance σ^2 for the patient topic mixtures θ . The decoder is linear and consists of two tri-factorizations. One learns medication-defined topic embedding $\alpha^{(med)}$ and medication embedding $\rho^{(med)}$. The other learns condition-specific topic embedding $\alpha^{(cond)}$ and the condition embedding $\rho^{(cond)}$. We separately pre-train (B) the embedding of medications $\rho^{(med)}$ and (C) the embedding of conditions $\rho^{(cond)}$ using node2vec (Grover and Leskovec, 2016) based on their structural meta-information. This is done learning the node embedding that maximizes the likelihood of the tree-structured relational graphs of conditions and medications.

Topic quality evaluation

Using data from the UKB on 457,461 individuals of European descent from the baseline visit, we trained GETM to obtain the topic embedding and conditions/medications embedding (i.e., $\alpha^{(t)}$, $\rho^{(t)}$, respectively). As a qualitative exploratory analysis, we visualized these embeddings using UMAP (Figure 2). For illustration purposes, we picked five topics representing diverse conditions and medications and observed that the top features under these topics belong to coherent categories of conditions or medications. For example, top medications atenolol, bisoprolol, metoprolol, nebivolol, and carvedilol in topic 32 all belong to cardiovascular-system medications (Figure 2D), while topic 32 is assigned to the cluster mainly composed of medications from the same category (Figure 2B). Moreover, the top five conditions from topic eight are all from the musculoskeletal/trauma category, whereas the top five medications from topic eight are all from the dermatological category. In contrast, ETM without using the Knowledge-Graph (KG)-informed embedding produced less interpretable topics, each covering heterogeneous categories (Figure S1). Thus, GETM allows for the identification of related features from sparse, heterogeneous data in a data-driven manner.

Given that we are leveraging embeddings of conditions and medications learned from their taxonomic graphs, we expected a higher quality of topics to be inferred by GETM compared to baseline models

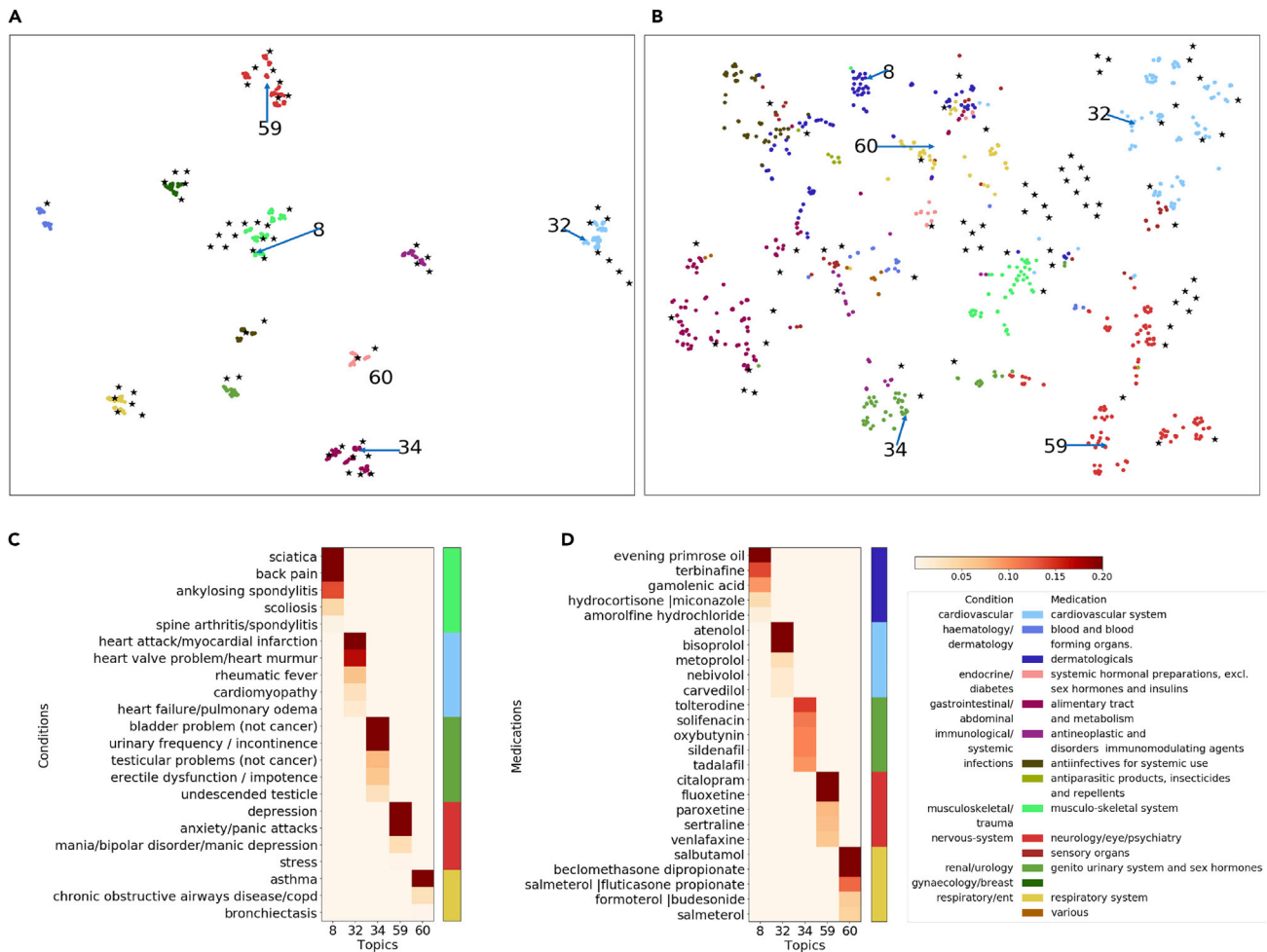


Figure 2. Visualizing the topic embedding and feature embedding learned by GETM on the UK Biobank data

The analysis was based on results from the GETM model with both condition and medication embedding that are KG-informed using 75 topics (Table S1). (A) Visualizing the embedding of topics and conditions. Because of the shared embedding space, we applied a single UMAP to project and visualize the condition embedding $\mathbf{p}^{(cond)}$ and topic embedding $\boldsymbol{\alpha}^{(cond)}$.

(B) Visualizing the embedding of topics and medications. Similarly, we applied a single UMAP to visualize medication embedding $\mathbf{p}^{(med)}$ and topic embedding $\boldsymbol{\alpha}^{(med)}$. The solid dots on the UMAP plot are the features, and the asterisks are the topics. The points are colored by corresponding category of condition and medication or annotated by its topic number. Visualizing the UMAP in this way allows us to examine (1) the similarity among topics, (2) the similarity among features, and (3) the similarity between topics and features. As indicated by the arrows, we identified five topics on each UMAP and displayed their top features in panel c and d.

(C and D) Heatmap visualization of select topics. We generated heatmaps for the top five conditions and the top five medications with the highest probability $\beta_k^{(cond)}$ and $\beta_k^{(med)}$ under each of the five topics. The color intensity displays proportionally the probabilities. The color bars on the right indicate the categories of the conditions and medications. All four panels share the same color legend of the categories.

that do not use or use only partial graph embeddings. Topic quality was quantified as the product of topic coherence and topic diversity (Dieng et al., 2019) (STAR Methods). For ease of reference, we listed the nine models that we compared in Table S1. We first evaluated the quality of the medication-defined topic (Table S4). To evaluate the medication topic coherence unbiasedly, we used an external set of 59 expert-curated medication categories (STAR Methods) that were not part of the medication taxonomy we used in training the graph embedding for our GETM model.

We repeated our experiments 5 times each time with a different random initialization for each model and evaluated the resulting topic coherence, topic diversity, and topic quality. Overall, we observed the highest average topic quality (Figure 3 and Table S4) for the 50-topic GETM that used the graph embeddings for both the conditions and medications (i.e., GETM in Table S1). Statistically, GETM is significantly better than

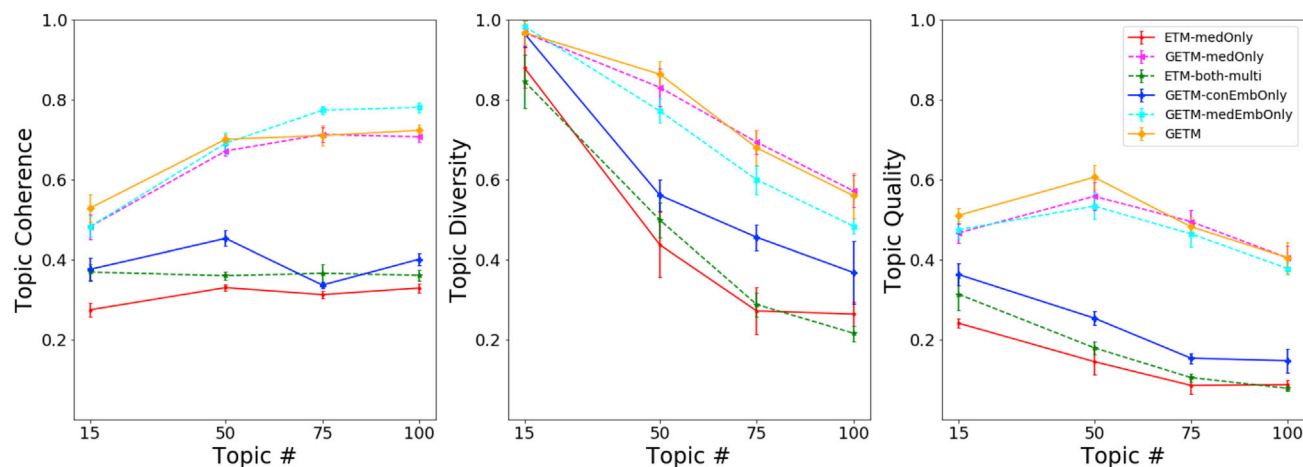


Figure 3. Medication-specific topic quality evaluation

We experimented with GETM as well as five baseline models (Table S1) with four predefined number of topics. To compute statistical significance between GETM and the baseline methods, we ran each model 5 times on the full UK Biobank data each with a different random initialization. The line plot displays the topic coherence, topic diversity, and topic quality, where the error bar indicates the standard deviations over the five experiments.

all of the baseline methods at p value < 0.01 (one-sided t-test) except for GETM-medOnly, compared to which GETM is better at p value < 0.05 (one-sided t-test). Therefore, the results suggest the benefits of using KG-informed medication embeddings in improving topic quality.

GETM-medEmbOnly confers slightly higher topic coherence but much poorer topic diversity for medications compared to GETM. This is possibly because the condition embedding directly learned from the UKB data provides complementary supports to inferring more coherent medication topics compared to the node2vec-pre-trained condition embedding from the taxonomy. Note that topic coherence measures each topic separately, whereas topic diversity measures the differences among topics. Both metrics are important. For instance, a model that generates a set of very similar topic distributions each with high topic coherent score is much less useful than a model that generates a set of diverse topic distributions each with slightly lower topic coherence score.

We evaluated condition-derived topic quality (Figure S2). Owing to the lack of external condition categories, we calculated topic coherence based on whether the top conditions from each topic were present in the same patients. One caveat of this approach is that a patient rarely exhibits multiple conditions from the same category and that many conditions can be mutually exclusive within the same patient. For example, asthma and COPD as shown in our GETM topic 60 in Figure 2B are rarely observed in the same patient despite the common physiological connection between them. This is reflected in the low topic coherence among the methods. Nonetheless, we found that the GETM with pre-trained embeddings of both medications and conditions (Table S1) dominate other models among three out of four model settings in terms of the overall topic quality scores (Table S5). GETM-medEmbOnly confers much poorer topic coherence for conditions (Table S5). This suggests an overall benefit of transferring the pre-trained embedding for both conditions and medications to the GETM when modeling the UKB data.

GETM reveals known or potentially novel condition-medication relations from UKB data

We compared the total number of unique known pairs between medications and conditions that were identified by the five models (Tables 1 and S1). The full GETM extracted most pairs of correlated conditions and medications illustrated by topic number 50 (161 pairs), 75 (175 pairs), and 100 (203 pairs). We examined some of the topics with high condition-medication associations based on pharmacological knowledge (Figures 2B and 2D). For instance, in topic 32, the top medication bisoprolol is known to be used to treat the top condition heart failure under that topic. Under topic 60, the top medication salmeterol is often prescribed to treat asthma and chronic obstructive airways (COPD) <https://go.drugbank.com/>, which are the top conditions under that topic. Interestingly, although the medication solifenacin in topic 59 is not known to be indicative of depression according to DrugBank, recent research shows potential for solifenacin,

Table 1. Number of known pairs between conditions and medications

Algorithm	Topic #			
	15	50	75	100
	Number of matched pairs			
ETM-both-flat	53	84	119	132
ETM-both-multi	63	86	105	126
GETM-condEmbOnly	62	118	159	162
GETM-medEmbOnly	65	135	171	178
GETM	61	161	175	203

We mapped our medications and conditions identifiers to CTD and DrugBank databases to obtain known connections between them. For each topic inferred, we generated nine condition-medication pairs from the top three conditions and the top three medications. Among these pairs, we calculated the number of known pairs (Table S1). Higher numbers imply more clinically meaningful topics inferred by the algorithm.

along with other muscarinic antagonists, in treating patients with depression through drug repurposing (Wróbel et al., 2020; Oliveira et al., 2020).

GETM accurately imputes UKB conditions and medications

To further ascertain the overall utility of GETM, we used GETM to reconstruct randomly masked 50% of the medications or conditions of each test individual. As a baseline, we used the observed condition or medication prevalence from the training data to impute the masked feature in the test individuals. We have also experimented with the ablated models listed in Table S1. Among all models, the full GETM conferred the lowest reconstruction errors for both conditions (Table 2) and medications (Table S6) using 100 topics with ~ 5% improvement over the second best method. The improvements are largely attributed to the use of KG-informed embeddings compared to the ablated ETM model that learns these embeddings from the data.

Moreover, we also evaluated GETM in its ability to recapitulate the medication data using only the condition information by masking all medication information of the test individuals. The full GETM outperformed all other models and gave the lowest reconstruction error (14.6125), the highest precision at top five predicted hits (precision@5) (0.2612), and the highest recall@5 (0.5787) (Tables 3 and 4). The upper bounds, which were obtained from the GETM trained on unmasked test data, achieved reconstruction error of 11.4936, precision@5 of 0.4226, and recall@5 of 0.8027. Interestingly, we found out that on average around 44 % of the unobserved medications from the top 10 predicted medications by GETM in fact have a treatment effect based on condition-medication association information extracted from Comparative Toxicogenomics Database (CTD) <http://ctdbase.org/> and DrugBank <https://go.drugbank.com/>. In particular, we took a closer look at the three best predicted patients and three worse predicted patients by GETM (Figure 4). As expected, among the best predicted patients, the top 10 predicted medications by GETM contain mostly observed medications, which are known to match the subset of the conditions of these three patients. Although most top predicted medications were not observed for the worst three predicted patients, they were known to treat the observed conditions for these patients.

CMK pain prediction

We sought to investigate the predictive ability of reported past conditions and current medications on pain experience in the UKB using GETM. We first focused on CMK pain because it has the highest number of positive cases. Using logistic regression (LR), we evaluated the predictive accuracy of the inferred topic mixture on CMK pain in terms of area under the receiver operating characteristic (AUROC) curve and area under the precision-recall curve (AUPRC). We used GETM with 128 topics for this experiment for the reason of its high performance on the validation dataset (Figure S3). As a baseline (i.e., the raw model), we trained another LR model that directly used all of the 443 conditions and 802 medications. We also sought to investigate the relative improvements of GETM over standard topic modeling after removing obvious conditions or medications for CMK pain. We experimented with six different ways of filtering out features based on ORs or expert knowledge (STAR Methods).

Table 2. Reconstructing 50% masked conditions

Algorithm	Topic #			
	15	50	75	100
	Reconstruction Error			
ETM-condOnly	5.89	5.56	5.39	5.80
GETM-condOnly	5.76	5.33	5.05	4.93
ETM-both_multi	5.08	5.09	4.80	4.96
GETM-condEmbOnly	5.12	4.84	4.55	4.40
GETM-medEmbOnly	5.06	4.86	4.75	4.58
GETM	4.69	4.84	4.45	4.34
Proportion	8.69			

We split the UKB data into 80% training and 20% test. For the test data, we randomly masked 50% of the values such that each test patient would have 50% of their conditions and medications observed. We then reconstructed the matrix with learned θ , α and ρ . The reconstruction error (i.e., negative log-likelihood) was calculated for the held-out data. Same condition data was used for all algorithms. For ETM-both-multi, GETM-condEmbOnly, GETM-medEmbOnly and GETM, the same models and data were used as for the medication reconstruction Table S6. Description of the algorithm names are in Table S1. As another baseline method, we evaluated the performance of filling in masked conditions based on their overall proportion over all of the UKB population.

Predictive performance using the original features (i.e., unfiltered) and the filtered features is summarized in Figure 5. From these results, we observed that (1) The GETM topic mixture (GETM, Table S1) achieved larger AUROC and larger AUPRC across all feature filtering regimes (i.e., no filter plus the six filtering rules; Table S2); (2) As we removed more pain-related signature conditions and medications, the performance of using raw features dropped more drastically than that using the patient topic mixtures. In other words, the relative improvement of using patient topic mixtures over using raw data increased as we removed more indicative conditions and medications. In particular, GETM conferred a larger than 40% improvement over the baseline when using the fewest conditions and medications (i.e., m579c322: filtered feature set 7) in contrast to less than 5% improvement over the baseline when using all features (m802c443: filtered feature set 1; Figure 5 and Table S2). These results echo the superior imputation performance of GETM we observed previously. This is because GETM uses the pre-trained graph embedding and inferred topic mixture to compensate for the information loss from the feature removals. Nonetheless, the absolute values of AUROC and AUPRC are not high. We discussed these results as limitations of the study in the Discussion section.

Table 3. Medication imputation

Algorithm	Topic #			
	15	50	75	100
	Medication Recovery Error			
Upper bound	12.25	11.24	11.50	11.50
ETM-both-flat	18.51	19.85	19.81	20.06
GETM-both-multi	14.99	14.88	14.94	15.05
GETM-condEmbOnly	15.24	14.95	14.95	14.96
GETM-medEmbOnly	15.18	14.88	14.90	14.88
GETM	14.82	14.65	14.61	14.65

The medication data was masked for test individuals. In contrast to the medication reconstruction experiments (Table S6), we masked the entire medications for the test patients. We imputed their medication data using the inferred θ from their condition data only and the learned embedding $\alpha^{(med)}$ and $\rho^{(med)}$. The reconstruction error (i.e., negative log-likelihood) was calculated. The upper bounds were obtained using reconstruction errors calculated from unmasked test data using GETM (GETM, Table S1). Description of the algorithm names are in Table S1.

Table 4. Medication imputation accuracy

Algorithm	Topic #							
	15	50	75	100	15	50	75	100
	recall@5				precision@5			
Upper bound	0.6672	0.7943	0.8027	0.7862	0.3342	0.4084	0.4226	0.4049
ETM-both-flat	0.4397	0.4486	0.4667	0.4614	0.1901	0.2109	0.2111	0.2162
GETM	0.5543	0.5639	0.5568	0.5388	0.2479	0.2516	0.2504	0.2417
GETM-condEmbOnly	0.5479	0.5664	0.5670	0.5647	0.2373	0.2524	0.2493	0.2486
GETM-medEmbOnly	0.5519	0.5722	0.5668	0.5716	0.2440	0.2533	0.2521	0.2532
GETM	0.5692	0.5787	0.5732	0.5753	0.2504	0.2612	0.2606	0.2578

The imputation procedures were the same as in Table 3. We calculated the precision and recall at the top five imputed medications for each individual. The upper bounds were calculated from unmasked test data using GETM (GETM, Table S1).

CMK pain-related conditions and medications

We then investigated the most pain-related conditions and medications based on LR coefficients and calculated overlapping proportions with physician-curated lists (Figure 6). In comparison with the overlapping proportions from ETM and odds ratio calculation, GETM identified a much greater proportion of known pain-related conditions among the top conditions and medications under these topics (36.7% from top 10 conditions, 33.3% from top 30 conditions and 30.0% from top 50 conditions) and medications (60.0% from top 10 medications, 33.3% from top 30 medications and 32.0% from top 50 medications) in the provided lists. This suggests that GETM improves the ability to extract otherwise hidden associations, which could identify pain-related comorbidities. Therefore, GETM does not only confer superior performance but higher model explainability from its inferred topics.

We also had a close look at the three most positively associated topics and three most negatively associated topics to CMK pain based on learned ω (Figure 7A). This analysis showed that topics 56, 34, and 51 are strongly positively associated with CMK pain and topics 73, 68, 89 are strongly negatively associated with CMK pain, respectively. For each topic, we examined their semantic meaning according to the top conditions and top medications (Figures 7B and 7C). Particularly, topics 56 and 34 contained musculoskeletal system conditions and medications, which are clinically meaningful because they are highly related to CMK pain. In particular, prolapsed disc or slipped disc as a condition is painful, and ibuprofen is an analgesic in the NSAID (non-steroidal anti-inflammatory drug) class.

Topic 51 is in the cardiovascular category, and the top medication acetylsalicylic acid (aspirin), also an NSAID, is prescribed more frequently and typically at lower doses for its cardioprotective properties in prevention of stroke and heart attacks; it acts as a “blood thinner”. Dipyridamole inhibits blood clot formation and therefore prevents potential consequences of blood clotting. The top condition under topic 51 is high cholesterol, which is a known and common risk factor for atherosclerosis. Atherosclerosis is a process of deposition of fatty material in the walls of arteries, and this thickening leads to an increase in stroke and heart attack risk. Thus, although the two medications are not directly used as a cure for the condition of high cholesterol, by way of atherosclerosis, high cholesterol leads to higher risk for other cardiovascular outcomes and the medications are used to prevent those outcomes (Al-Ghamdi et al., 2021).

Contrary to the risk topics, the protective topics identified combinations of medications and conditions that did not yield as straightforward pairings or links. Topic 73 is one of the top three negatively predictive topics of CMK pain. Its top medications identified were ramipril, lisinopril and enalapril, which are angiotensin-converting enzyme (ACE) inhibitors, used to treat high blood pressure and may be used in response to heart failure or heart attack. Incidentally, all conditions under topic 73 may be categorized as allergic or atopic with hayfever contributing the most.

This finding suggests that a particular subset of individuals suffering from allergic conditions, whose immunity is therefore skewed, and who are also undergoing cardiovascular treatment, are at lower risk for CMK pain. The implication of immune mechanisms in chronic pain manifestation is well known, where

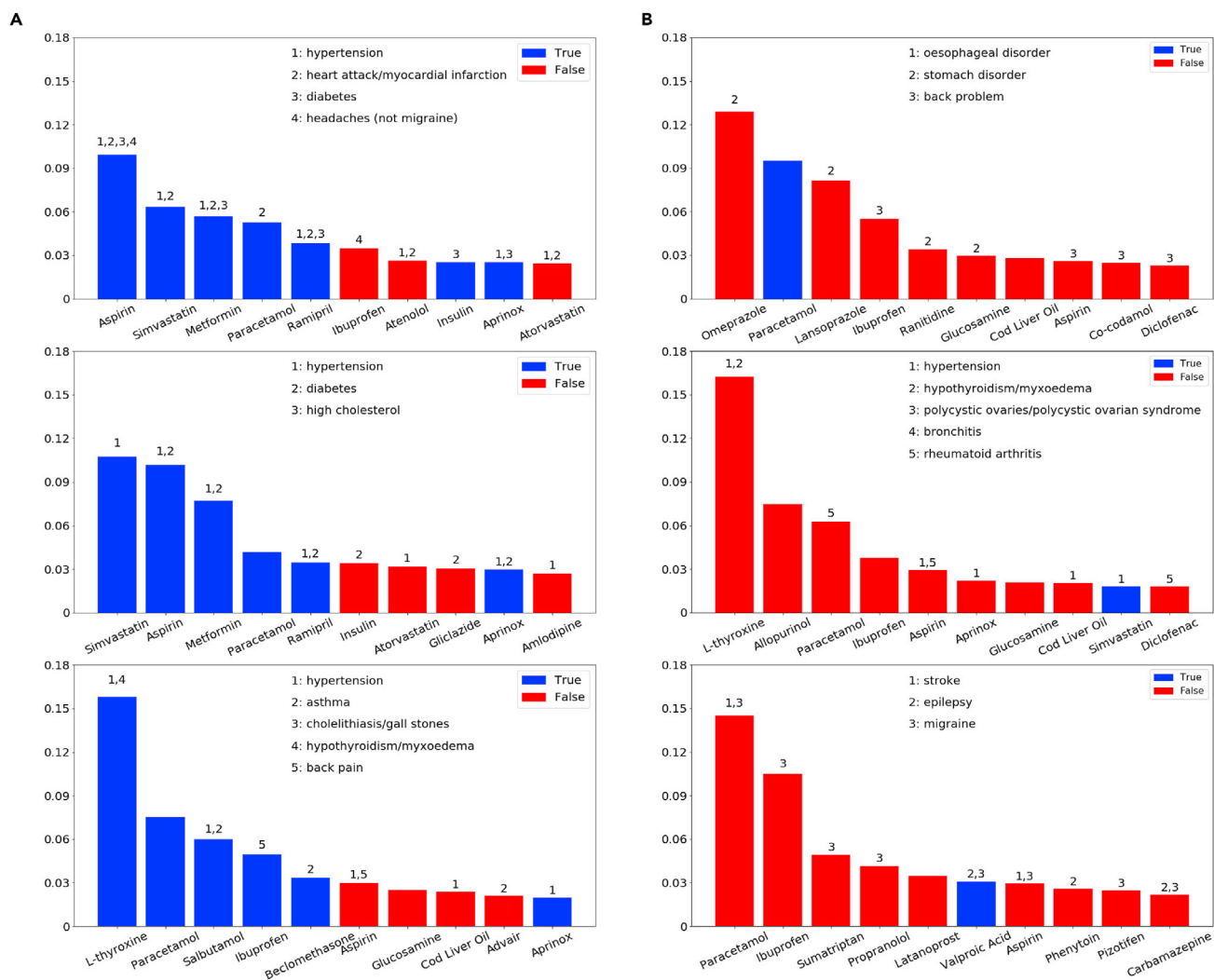


Figure 4. Illustration of GETM for medication imputation through examples of study subjects

(A) The three best-matched study subjects displayed as three barplot panels on the left. We chose three best-matched individuals for whom our imputed medications matched well with the observed medications.

(B) The three worst-matched individuals displayed as three barplot panels on the right. The y axis represents the predicted probability of the medications on the x axis. The numbers on the top of each bar represents the patient's observed conditions that are known to be treated by the corresponding medication under the bar. Inset in each panel lists the conditions names corresponding to the numbers on top of each bar. Blue bars indicate observed medications in the patient, and red bars indicate unobserved medications.

proinflammatory states are associated with chronic pain (Baral et al., 2019). In addition, high blood pressure is the most common indication for ACE inhibitors (Heran et al., 2008). It is known to have analgesic effects that persist with treatment (Ghione, 1996; Makovac et al., 2020). Thus, hypertension related analgesia would help explain this negative predictive topic.

Another negatively predictive topic of CMK pain, topic 89, is limited to women given the top medication, conjugated estrogens. The top condition, hypertension, as seen previously, is known to have analgesic effects. Sex hormones including estrogens modulate pain, and estrogens have been found to have neuroprotective, antinociceptive properties. These could explain the protective association observed here on CMK pain. Indeed, estrogens are known to influence chronic pain conditions (Chen et al., 2021). These combinations of medications and conditions within topics could lead to new CMK pain etiology hypotheses for further exploration.

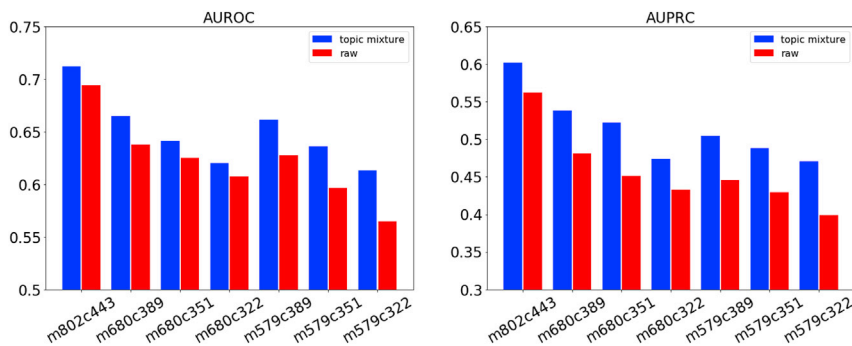


Figure 5. Performance of logistic regression (LR) for CMK pain

We trained LR models using θ obtained from GETM (Table S1) with 128 topics as input to predict CMK pain. The baseline LR model directly used the raw condition and medication data as input. We iterated through seven feature filtering schemes with different filtered condition sets and medication sets as indicated by x axis (details in Table S2 and described in STAR Methods Section 8). Barplot displays the AUROC and AUPRC across these experiments.

Characterizing diverse pain types by topic analysis

We extended our analysis from CMK to predicting several other definitions of pain phenotypes, including acute pain, the acute to chronic pain transition for a subset of UKB individuals, chronic pain at specific body sites (neck/shoulder, hip, back, stomachache/abdominal, knee, headache, face) and chronic pain all over body. We used three different feature filtering regimes to progressively remove obvious pain-type predictors (STAR Methods). Logistic regression using the GETM's topic mixture conferred larger AUROC and larger AUPRC values across all three feature filtering regimes (Figure 8 and Table S1). We then examined the relative contributions of the medication/condition categories for each pain type (Figure S5; STAR Methods). The results we have presented in the following sections came from at least one of the three feature filtering regimes as highlighted in Figure S5. In addition to the defined categories, we used the Anatomical Therapeutic Chemical (ATC) classification system to consider the contribution of the medication subclass analgesics, which has a direct relationship with pain. By considering the implication of condition and medication categories in predicting different pain phenotypes, we were able to identify expected trends and to consider and interpret unexpected findings.

Overall, the acute and acute-to-chronic transition displayed similar trends. Here the acute-to-chronic transition is defined such that non-site specific MSK acute pain observed at the first visit turned into chronic pain as diagnosed at the following visit of the UKB individuals. As an example, low positive prediction from analgesics and high negative contribution were seen for acute pain (Figure 9).

As expected, the category from the literature that displays the highest levels of chronic pain comorbidities, neurology/eye/psychiatry (conditions) and nervous system (medications) contributed the most to prediction of chronic pain, particularly headache and face pain (Figures S5A-I and S5B-I). Endocrine/diabetes (conditions) exhibited a highly prominent contribution for the acute phenotype classification (Figure S5A-II). Immunological/systemic disorders showed a strong contribution for protection from acute pain (Figure S5C-I). Gastrointestinal/abdominal (conditions) and alimentary tract and metabolism (medications) showed an expected higher positive proportion of contribution toward stomach/abdominal pain prediction compared to other chronic pain types (Figure S5A-III).

Knee pain compared to face pain and headache exhibited a sharp contrast in predictive compositions (Figures S6A and S6B). Interestingly, obesity-related medication categories such as systemic hormonal preparations excluding sex hormones and insulins exhibited the highest protective contribution for chronic knee pain (Figures S7D–S7I). Chronic knee pain also had the highest positive contribution toward the prediction of cardiovascular condition and no negative contribution (Figures S6A-I and S6C-I). A similar but less striking pattern was observed for hip pain (Figures S6A-II, S6C-II). For headache and face pain, this trend was reversed with the highest protective contribution and lowest positive contribution by cardiovascular conditions (Figures S6A-III and S6C-III).

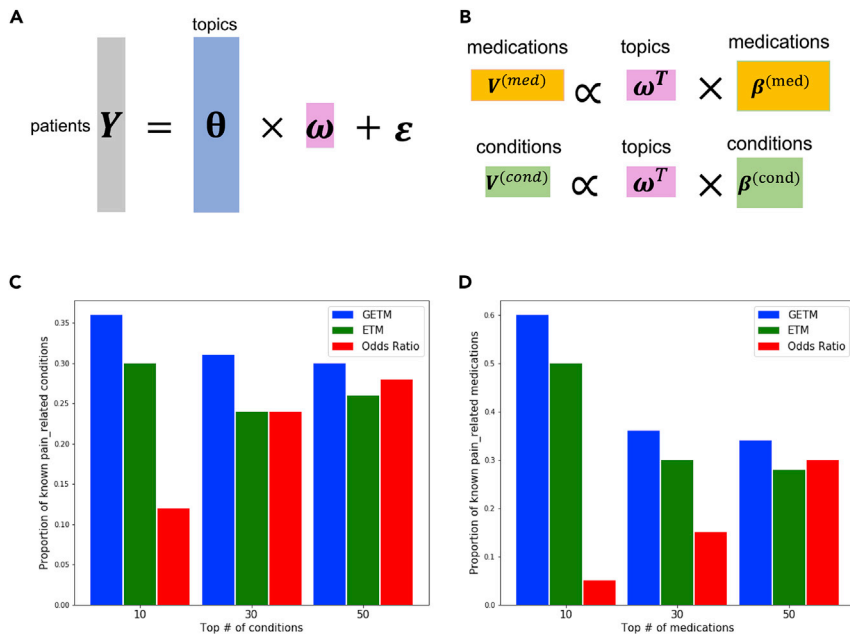


Figure 6. Analysis of CMK-pain-related conditions and medications

(A) Logistic regression using patient topic mixture $\theta \in \mathbb{R}^{D \times K}$ of D individuals and K topics to predict CMK pain as a binary outcome.

(B) Importance score computation for medications and conditions. Taking the inner product of the regression coefficients $\omega^T \in \mathbb{R}^{1 \times K}$ and $\beta^{(med)} \in \mathbb{R}^{K \times M}$ from GETM (GETM, Table S2), we obtained the importance scores of predicting CMK pain as $v^{(med)} \in \mathbb{R}^{1 \times M}$ and $v^{(cond)} \in \mathbb{R}^{1 \times C}$ for medications and conditions, respectively.

(C) Proportions of known CMK-related conditions based on a physician-compiled list.

(D) Proportions of known CMK-related medications based on a physician-compiled list of analgesic medications.

Comparisons relate to two baselines: (1) Using ETM which treated conditions and medications as the same type of feature (i.e., ETM-both-flat, Table S1). For this baseline, we selected the top medications and conditions from the resulting $v \in \mathbb{R}^{1 \times (M+C)}$. (2) Odds Ratio (STAR Methods Section 8).

Intriguingly, the cardiovascular medications category exhibited a strong negative prediction of 0.7 across all chronic pain phenotypes but not acute pain (Figure 9). This suggests that topics containing cardiovascular medications, other than analgesics, had a protective effect against chronic pain. Headache, and more specifically migraine, has a known vascular etiological component. Nonetheless, the patterns seen across all chronic pain phenotypes are similar, indicating a more general chronic pain protective effect for cardiovascular medications.

DISCUSSION

Large biobanks with clinical data such as the UKB are valuable resources enabling greater understanding of factors impacting the manifestation and treatment regimens of complex diseases such as chronic pain. However, the sparsity, heterogeneity and sheer size of these data pose challenges when restricted to a classical statistical toolbox. The typical biomedical researcher is hindered from taking full advantage of these invaluable data resources. As a result, most investigations focus on one or a small subset of related diseases (Groen et al., 2020; Song et al., 2020; Singh et al., 2019). In the present study, we developed GETM to model all self-reported conditions and medications among 457,461 UKB subjects. Our main focus is to understand the comorbidity among the UKB subjects in relation to pain-related phenotypes. By introducing the knowledge graph and simultaneously training different types of features, GETM was able to infer more coherent topics compared to the ablated models without using the KG embedding (Tables 1 and S4). In contrast to the topic modeling without the graph embedding, this allows for better interpretation of the disease topics and any findings related to certain topics with clearer clinical grounds.

As demonstrated using UKB data, GETM achieved superior performance in imputing 50% of the observed conditions (Table 2) and 50% of the observed medications over all test individuals (Table 3) as well as in

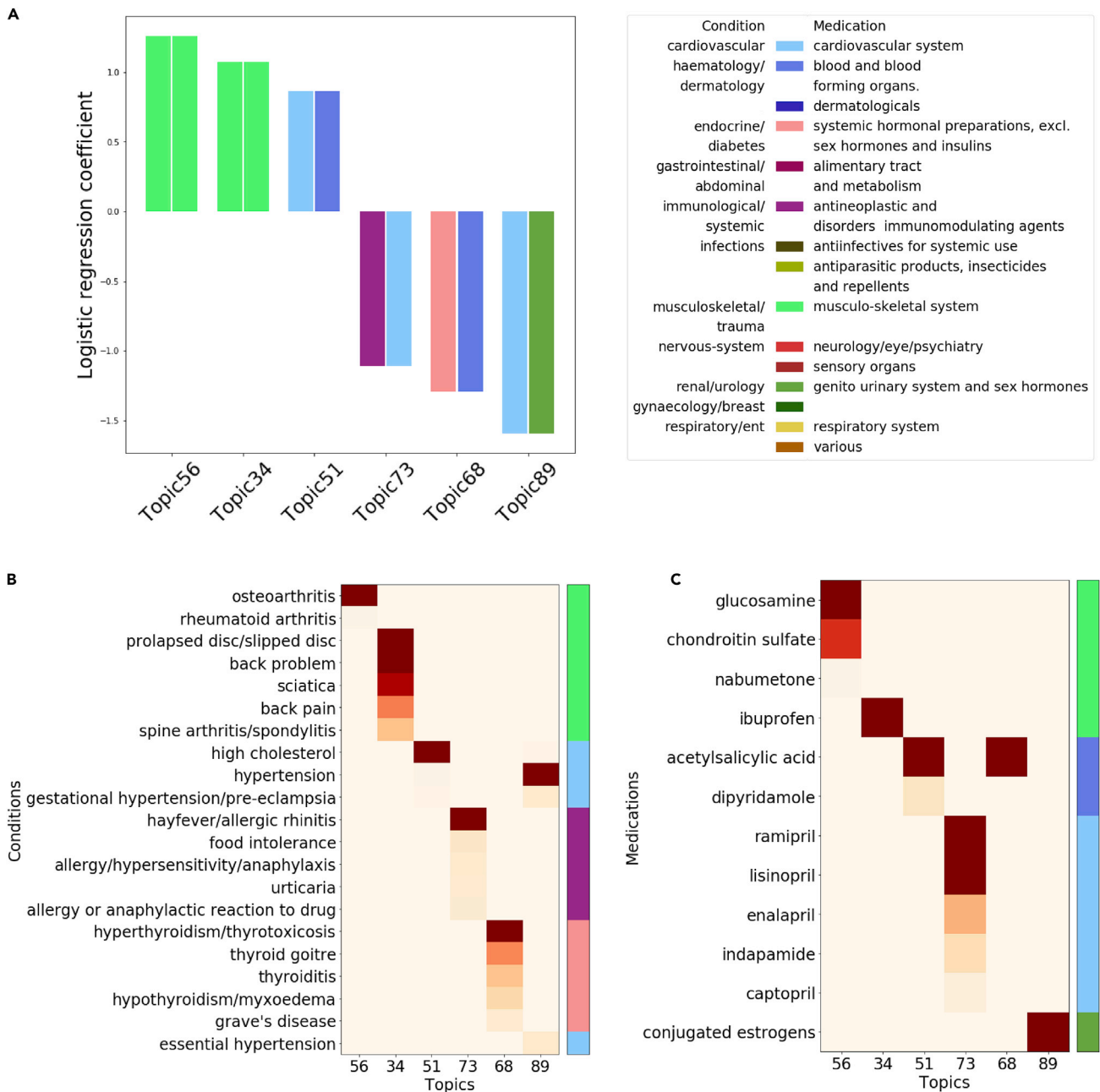


Figure 7. Topic analysis for CMK pain

(A) The most predictive CMK pain topics. Based on the logistic regression coefficients of predicting CMK pain ω (Figure 6A), we chose three topics with the highest coefficients and three topics with the most negative coefficients. Each bar is composed of condition and medication category. (B and C) The top conditions and medications under the six most predictive topics. Same as in Figure 2, the color bars on the right indicate the categories of the conditions and medications.

imputing the entire medication records based only on the conditions records of test individuals (Table 4). Many top unobserved medications predicted by GETM also exhibit known connections with the observed conditions of the subject (Figure 4). Importantly, GETM offers excellent model interpretability and can be used to discover meaningful disease comorbidities and disease-medication combinations via topic analysis (Figure 2).

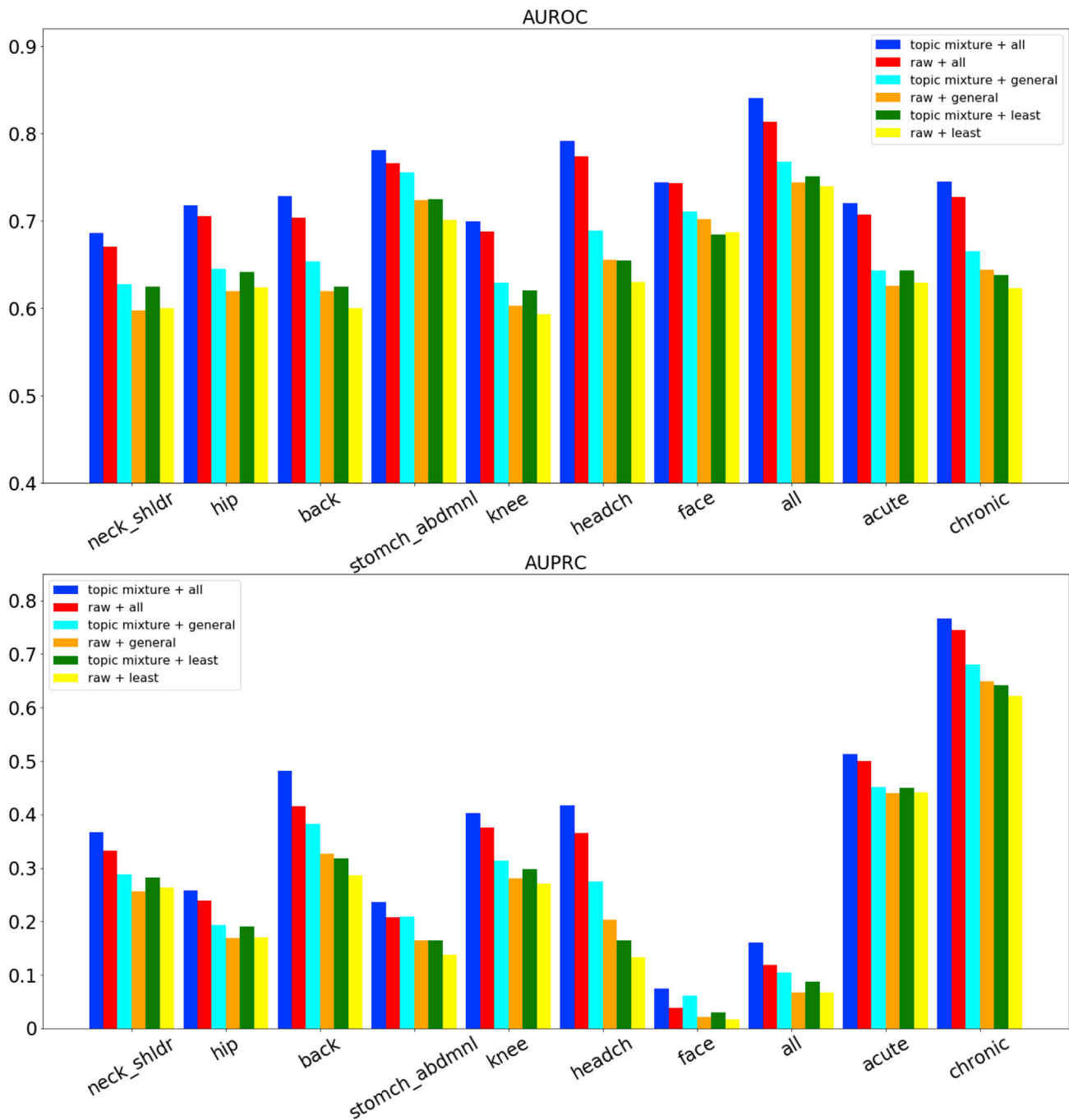


Figure 8. Prediction performance of 10 pain types

Logistic regression was trained to predict 10 pain types as indicated in the x axis. We experimented with three different filtered feature sets of conditions and medications: (1) all (m802c443): non-filtered conditions and medications, (2) general (m680c351): physician-curated general pain-related conditions and medications were filtered, and (3) least: physician-curated and conditions and medications based on with odds ratios were filtered. Details of the data filtering were described in [Table S2](#).

In a focused analysis to predict chronic musculoskeletal (CMK) pain, the same logistic regression (LR) classifier using GETM-inferred topic mixtures conferred robustly and consistently higher AUROC and AUPRC compared to using the raw features ([Figure 5](#)). The top conditions and medications under the most predictive topics are enriched for elements from the physician-curated condition list and medication list,

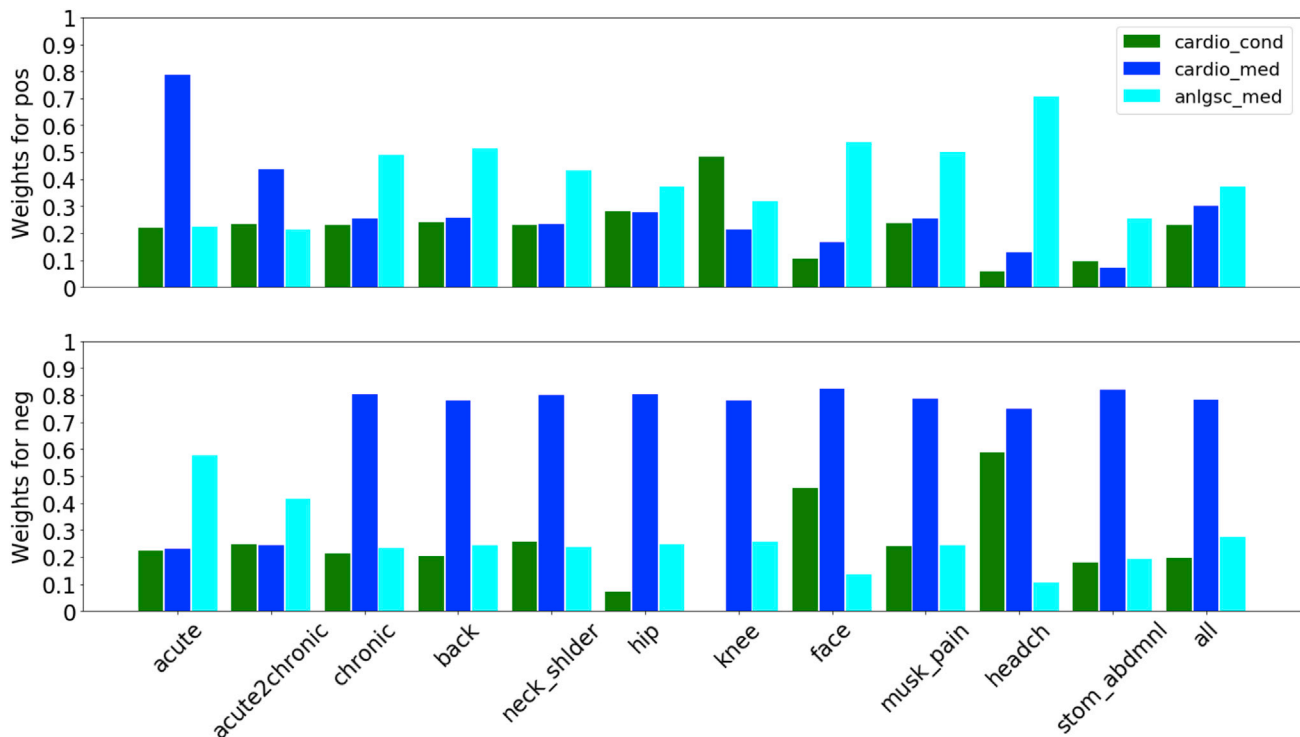


Figure 9. Contributions of cardiovascular-related conditions and medications to 12 pain types

The importance weights of cardiovascular conditions and medications in predicting the 12 pain types were calculated by the sum of their topic probabilities weighted by the linear regression coefficients across all topics (STAR Methods). Descriptions of the names of these pain types are in Table S3. The results for the cardiovascular category were based on the filtered sets, where the known pain-related conditions and medications based on a physician curated list were removed. The results of analgesic category (anlgsc_med) were based on the non-filtered set, as only full set contains analgesic medications. All of the importance weights from other filtering schemes were displayed in Figure S5.

respectively (Figure 6). In addition, the most predictive CMK pain topics contained top conditions and medications that were strongly associated with CMK pain (Figure 7), implying their potential usage as phenotyping markers. In addition, the GETM topic mixtures are robust to the removal of clinically obvious conditions and medications (Figure 5). This is reflected in the larger relative improvements over the LR classifier that operates on the raw (filtered) features. This implies its utility in predicting predisposition for CMK pain for those individuals with no reported pain symptoms. We observed consistent quantitative performance when extending to other pain types (Figure 8).

Several existing methods utilize topic models to find meaningful latent topics from electronic health record (EHR) data using structured administrative data such as the ICD codes (Li et al., 2020; Song et al., 2021). In contrast, we demonstrated the utility of GETM on the less structured and more sparse self-reported questionnaire information from the UKB including 443 conditions and 802 medications. We expect that GETM would work equally well if not better on more structured EHR data and leave that to future exploration. Indeed, although we used UKB data and focused on pain phenotypes as a case study, our approach is a generalizable and highly efficient method that can be used to characterize other phenotypes in the UKB or from other similarly scaled biobanks.

We now discuss the epidemiological implications of our results on the pain analysis in the context of existing studies. Our findings identified intriguing links between cardiovascular medications and their protective effect for chronic (but not acute) pain that stood out as particularly predictive compared with other medication or condition categories (Figure 9). In addition, cardiovascular conditions were negatively predictive for headache and face pain particularly. Recent meta-analyses have attempted to quantify the relationship between cardiovascular conditions and chronic pain (Fayaz et al., 2016b; Oliveira et al., 2019), but their relevance to our findings is low. Neither of these nor the studies that

were included in the meta-analyses explicitly quantified the effects of taking medications. In addition, the direction of causality focused on chronic pain as a risk factor for cardiovascular outcomes, including mortality, and did not consider reverse causality; authors cited diversity in outcomes and in chronic pain taxonomy making meta-analysis results generally inconclusive (Fayaz et al., 2016b). The more recent meta-analysis was more limited in scope, estimating that people with CMK pain were 1.91 times more likely to report having a cardiovascular disease compared with those without CMK pain with statistical significance (Oliveira et al., 2019).

The dominance of the predictive effect of cardiovascular medications across chronic pain phenotypes leads to several avenues for further exploration. The specifics of what medications and what conditions are found in individually highly predictive topics would need to be explored in further detail. In addition, further analyses are needed to identify whether more specific subcategories of cardiovascular medications drove the strong cardiovascular medications predictive component. The extent of similarity across specific medications within and across topics, such as common mechanisms of action, would be of particular interest. Several medications united by a common mechanism would have a higher potential for generalizability beyond the limited population subsets represented by individual items or topics. This was illustrated by our focused analysis on CMK pain, where the predictive medications identified in the top protective topic (Figures 7B and 7C, Topic 73) were all ACE inhibitors. Given that the most common indication for this medication class is hypertension, perhaps one of the protective cardiovascular medication effects at play was hypertension-associated hypoalgesia, especially because the hypoalgesia effect is maintained even if the hypertension is treated (Ghione, 1996; Suarez-Roca et al., 2019). Baroreceptor sensitivity is the most studied specific mechanism (Ghione, 1996; Suarez-Roca et al., 2019). Hypertension-associated hypoalgesia manifests through a range of phenotypes from pain sensitivity to pain chronification, with higher blood pressure associated with lower pain sensitivity or lesser chronification. This has been demonstrated in both animal experimental and human observational studies (Ghione, 1996; Bruehl and Chung, 2004). Nonetheless, hypertension-associated hypoalgesia is not well recognized.

Perhaps the present report will lead to findings supporting greater clinical and public health importance for hypertension-associated hypoalgesia in preventing pain chronification than previously thought. Beyond ACE inhibitors, cardiovascular medications classified as beta blockers are used as prophylaxis for migraine headaches (Jackson et al., 2019). They have been shown to have analgesic effects for other types of chronic pain (Tchivileva et al., 2010, 2020). There may also be further subcategories of cardiovascular medications of interest.

Limitations of the study

In terms of the data used in this study, the UKB as a cross-sectional study creates challenges in interpreting the role of temporality and identifying putative causal effects. The pain questionnaire was administered at baseline, and medications reported were taken regularly, also at baseline. Conditions were considered to occur at any time in the subjects' past, through a record of age-of-onset for each condition. Thus, the impact of medications on reported conditions do not allow for direct investigations of medication efficacy. In classifying chronic pain, the pain questionnaire referred to a time including at least the prior three months in the past from baseline, but this could extend to any length in a given subject's history. By considering the subset of individuals without baseline chronic pain, but who developed chronic pain by the time of a subsequent visit, it would be possible to approximately quantify the total length of time of pain chronicity. We may also consider chronification as temporally subsequent to the appearance of conditions and taking of medications from before baseline. This is only true for a subset of UKB subjects, however. Thus, the UKB study design has limited potential to evaluate chronic pain causality. Rather, our analyses produced hypotheses that might best be tested in the context of other cohorts that are truly longitudinal in design. In terms of the biological findings in this paper, there are variable time frames between development of cardiovascular conditions, taking of cardiovascular medications and development of chronic pain in the present study. Therefore, other longitudinal studies that control these sources of variability would be needed to better understand the protective effect suggested by the present analysis.

Although GETM provides relatively the highest chronic musculoskeletal (CMK) pain prediction accuracy (Figure 5), the absolute AUROC and AUPRC are not high. This underscores the limitation of both the UKB data and the GETM model. First of all, we only used self-reported conditions and medications in

this task, which are noisy, sparse, and sometimes inaccurate. To make more accurate predictions on CMK pain, other sources of health status data such as healthcare provider-facilitated clinical notes, laboratory tests, genomic measurements such as gene expression profiles of each patient are needed. Another important missing piece for predicting chronic conditions is the longitudinal information of each patient. As mentioned previously, we only have the baseline data at the initial visit for most of the UKB subjects. As longitudinal data become increasingly available, we will extend our model to a dynamic topic model that accounts for the evolution of subjects' health status over time. In particular, we will allow for integration of longitudinal healthcare information such as age-of-onset across conditions and multiple follow-up time-points. We will also improve embedding learning by considering a more expressive graph embedding approach than node2vec. For instance, we will consider a graph convolutional neural network (GCN) (Kipf and Welling, 2016) to produce the embedding used by the ETM model. Furthermore, we will explore a multi-relational graph approach (Schlichtkrull et al., 2018) to model the relationships within and among conditions and medications. Because both the GCN and the ETM models share the same objective function (i.e., the evidence lower bound), we will be able to perform end-to-end training instead of the pipeline approach presented here. To produce competitive performance with this approach, however, careful model fine-tuning of both the GCN and ETM will be needed at the expense of computational resources.

STAR★METHODS

Detailed methods are provided in the online version of this paper and include the following:

- **KEY RESOURCES TABLE**
- **RESOURCE AVAILABILITY**
 - Lead contact
 - Material availability
 - Data and code availability
- **METHOD DETAILS**
 - UKB data processing
 - UKB pain-related phenotype labels extractions
 - Graph-embedded topic model details
 - Ablation experiments
 - Topic quality evaluation
 - Study of medication and condition relations
 - Data imputation
 - CMK pain prediction
 - Odds ratios of CMK-related conditions or medications
 - Feature filtering schemes in removing obvious CMK-related conditions and medications
 - Calculating importance of conditions and medications for CMK pain
 - Characterizing pain types from 7 body sites and pain all over the body
- **QUANTIFICATION AND STATISTICAL ANALYSIS**

SUPPLEMENTAL INFORMATION

Supplemental information can be found online at <https://doi.org/10.1016/j.isci.2022.104390>.

ACKNOWLEDGMENTS

Y.L. is supported by Natural Sciences and Engineering Research Council (NSERC) Discovery Grant (RGPIN-2019-0621), Fonds de recherche Nature et technologies (FRQNT) New Career (NC-268592), and Canada First Research Excellence Fund Healthy Brains for Healthy Life (HBHL) initiative New Investigator start-up award (G249591). AVG and YW were partly funded by the Canadian Excellence Research Chairs (CERC09). LD was supported by the Canadian Excellence Research Chairs fund (CERC09), a Pfizer Canada Professorship in Pain Research, and CIHR (SCA-145102) for Health Research's Strategy for Patient-Oriented Research (SPOR) in Chronic Pain. The current study was conducted under UK Biobank application 20802.

AUTHOR CONTRIBUTIONS

These authors jointly supervised this work: Audrey V. Grant and Yue Li. Y.L. and A.V.G. conceived the study. Y.L. and Y.W. developed the methodology. Y.W. created the computational software and ran the analyses. Y.W., Y.L., and A.V.G. analyzed the results and wrote the initial manuscript. L.D. and R.B. participated in

initial discussions and subsequently in critical writing related to pain phenotypes and medications. All of the authors wrote the final version.

DECLARATION OF INTERESTS

The authors declare no competing interests.

Received: February 25, 2022

Revised: April 8, 2022

Accepted: May 6, 2022

Published: June 17, 2022

REFERENCES

- Al-Ghamdi, S., Shubair, M.M., El-Metwally, A., Alsalamah, M., Alshahrani, S.M., Al-Khateeb, B.F., Bahkali, S., Aloudah, S.M., Al-Zahrani, J., et al. (2021). The relationship between chronic pain, prehypertension, and hypertension. A population-based cross-sectional survey in Al-Kharj, Saudi Arabia. *Postgrad. Med.* 133, 345–350. <https://doi.org/10.1080/00325481.2020.1863716>.
- Baral, P., Udit, S., and Chiu, I.M. (2019). Pain and immunity: implications for host defence. *Nat. Rev. Immunol.* 19, 433–447. <https://doi.org/10.1038/s41577-019-0147-2>.
- Blei, D.M., Ng, A.Y., and Jordan, M.I. (2003a). Latent dirichlet allocation. *J. Mach. Learn. Res.* 3, 993–1022.
- Blei, D.M., Griffiths, T.L., Jordan, M., and Tenenbaum, J.B. (2003b). Hierarchical topic models and the nested Chinese restaurant process. In *Proceedings of the 16th International Conference on Neural Information Processing Systems. NIPS'03 (MIT Press)*, pp. 17–24.
- Bordes, A., Usunier, N., Garcia-Duran, A., and Weston, J. (2013). Translating embeddings for modeling multi-relational data. *Adv. Neural Inf. Process. Syst.* 26.
- Bruehl, S., and Chung, O.Y. (2004). Interactions between the cardiovascular and pain regulatory systems: an updated review of mechanisms and possible alterations in chronic pain. *Neurosci. Biobehav. Rev.* 28, 395–414. <https://doi.org/10.1016/j.neubiorev.2004.06.004>.
- Bycroft, C., Freeman, C., Petkova, D., Band, G., Elliott, L.T., Sharp, K., Motyer, A., Vukcevic, D., Delaneau, O., O'Connell, J., et al. (2018). The UK Biobank resource with deep phenotyping and genomic data. *Nature* 562, 203–209. <https://doi.org/10.1038/s41586-018-0579-z>.
- Chen, Q., Zhang, W., Sadana, N., and Chen, X. (2021). Estrogen receptors in pain modulation: cellular signaling. *Biol. Sex Differ.* 12, 22. <https://doi.org/10.1186/s13293-021-00364-5>.
- . CoRR abs/1907.04907 Dieng, A.B., Ruiz, F.J.R., and Blei, D.M. (2019). Topic modeling in embedding spaces. Preprint at arXiv. <https://doi.org/10.48550/arXiv.1907.04907>.
- Dueñas, M., Ojeda, B., Salazar, A., Mico, J.A., and Failde, I. (2016). A review of chronic pain impact on patients, their social environment and the health care system. *J. Pain Res.* 9, 457–467. <https://doi.org/10.2147/jpr.s105892>.
- Fayaz, A., Croft, P., Langford, R.M., Donaldson, L.J., and Jones, G.T. (2016a). Prevalence of chronic pain in the UK: a systematic review and meta-analysis of population studies. *BMJ Open* 6, e010364. <https://doi.org/10.1136/bmjopen-2015-010364>.
- Fayaz, A., Ayis, S., Panesar, S.S., Langford, R.M., and Donaldson, L.J. (2016b). Assessing the relationship between chronic pain and cardiovascular disease: a systematic review and meta-analysis. *Scandinavian J. Pain* 13, 76–90. <https://doi.org/10.1016/j.sjpain.2016.06.005>.
- Ghione, S. (1996). Hypertension-associated hypalgesia. Evidence in experimental animals and humans, pathophysiological mechanisms, and potential clinical consequences. *Hypertension* 28, 494–504. <https://doi.org/10.1161/01.hyp.28.3.494>.
- Glorot, X., Antoine, B., and Bengio, Y. (2011). Deep Sparse Rectifier Neural Networks (AISTATS).
- Groen, R.N., Ryan, O., Wigman, J.T.W., Riese, H., Penninx, B.W.J.H., Giltay, E.J., Wichers, M., and Hartman, C.A. (2020). Comorbidity between depression and anxiety: assessing the role of bridge mental states in dynamic psychological networks. *BMC Med.* 18, 308. <https://doi.org/10.1186/s12916-020-01738-z>.
- . CoRR abs/1607.00653 Grover, A., and Leskovec, J. (2016). node2vec: scalable feature learning for networks. Preprint at arXiv. <https://doi.org/10.48550/arXiv.1607.00653>.
- Heran, B.S., Wong, M.M., Heran, I.K., and Wright, J.M. (2008). Blood pressure lowering efficacy of angiotensin converting enzyme (ACE) inhibitors for primary hypertension. *Cochrane Database Syst. Rev.* 4, CD003823. <https://doi.org/10.1002/14651858.CD003823.pub2>.
- . CoRR abs/1502.03167 Ioffe, S., and Szegedy, C. (2015). Batch normalization: accelerating deep network training by reducing internal covariate shift. Preprint at arXiv. <https://doi.org/10.48550/arXiv.1502.03167>.
- Jackson, J.L., Kuriyama, A., Kuwatsuka, Y., Nickoloff, S., Storch, D., Jackson, W., Zhang, Z.J., and Hayashino, Y. (2019). Beta-blockers for the prevention of headache in adults, a systematic review and meta-analysis. *PLoS One* 14, e0212785. <https://doi.org/10.1371/journal.pone.0212785>.
- Jensen, P.B., Jensen, L.J., and Brunak, S. (2012). Mining electronic health records: towards better research applications and clinical care. *Nat. Rev. Genet.* 13, 395–405. <https://doi.org/10.1038/nrg3208>.
- Kingma, D., and Welling, M. (2014). Auto-encoding variational Bayes. Preprint at arXiv. <https://doi.org/10.48550/arXiv.1312.6114>.
- Kipf, T., and Welling, M. (2016). Semi-supervised classification with graph convolutional networks. Preprint at arXiv. <https://doi.org/10.48550/arXiv.1609.02907>.
- Li, Y., Nair, P., Lu, X.H., Wen, Z., Wang, Y., Dehaghi, A.A.K., Miao, Y., Liu, W., Ordog, T., Biernacka, J.M., et al. (2020). Inferring multimodal latent topics from electronic health records. *Nat. Commun.* 11, 2536. <https://doi.org/10.1038/s41467-020-16378-3>.
- Makovac, E., Porciello, G., Palomba, D., Basile, B., and Ottaviani, C. (2020). Blood pressure-related hypoalgesia: a systematic review and meta-analysis. *J. Hypertens.* 38, 1420–1435. <https://doi.org/10.1097/HJH.0000000000002427>.
- McCoy, T.H., Castro, V.M., Snapper, L.A., Hart, K.L., and Perlis, R.H. (2017). Efficient genome-wide association in biobanks using topic modeling identifies multiple novel disease loci. *Mol. Med.* 23, 285–294. <https://doi.org/10.2119/molmed.2017.00100>.
- Nguyen, D.Q., Billingsley, R., Du, L., and Johnson, M. (2015). Improving topic models with latent feature word representations. *Trans. Assoc. Comput. Linguist.* 3, 299–313. https://doi.org/10.1162/tacl_a_00140.
- Oliveira, C.B., Maher, C.G., Franco, M.R., Kamper, S.J., Williams, C.M., Silva, F.G., and Pinto, R.Z. (2019). Co-occurrence of chronic musculoskeletal pain and cardiovascular diseases: a systematic review with meta-analysis. *Pain Med.* 21, 1106–1121. <https://doi.org/10.1093/pm/pnz217>.
- Oliveira, C.B., Maher, C.G., Franco, M.R., Kamper, S.J., Williams, C.M., Silva, F.G., and Pinto, R.Z. (2020). Co-occurrence of chronic musculoskeletal pain and cardiovascular diseases: a systematic review with meta-analysis. *Pain Med.* 21, 1106–1121. <https://doi.org/10.1093/pm/pnz217>.
- Schlichtkrull, M., Kipf, T.N., Bloem, P., van den Berg, R., Titov, I., and Welling, M. (2018). Modeling relational data with graph

convolutional networks. In *European semantic web conference* (Springer), pp. 593–607.

Singh, U., Wangia-Anderson, V., and Bernstein, J.A. (2019). Chronic rhinitis is a high-risk comorbidity for 30-day hospital readmission of patients with asthma and chronic obstructive pulmonary disease. *J. Allergy Clin. Immunol. Pract.* 7, 279–285.e6. <https://doi.org/10.1016/j.jaip.2018.06.029>.

Song, J., Zeng, M., Wang, H., Qin, C., Hou, H., Sun, Z., Xu, S., Wang, G., Guo, C., Deng, Y., et al. (2020). Distinct effects of asthma and COPD comorbidity on disease expression and outcome in patients with COVID-19. *Allergy* 76, 483–496. <https://doi.org/10.1111/all.14517>.

Song, Z., Toral, X.S., Xu, Y., Liu, A., Guo, L., Powell, G., Verma, A., Buckeridge, D., Marelli, A., and Li, Y. (2021). Supervised multi-specialist topic model with applications on large-scale electronic health record data. In *Proceedings of the 12th ACM Conference on Bioinformatics, Computational Biology, and Health Informatics*.

BCB '21 (Association for Computing Machinery). <https://doi.org/10.1145/3459930.3469543>.

Suarez-Roca, H., Klinger, R.Y., Podgoreanu, M.V., Ji, R.R., Sigurdsson, M.I., Waldron, N., Mathew, J.P., and Maixner, W. (2019). Contribution of Baroreceptor function to pain perception and perioperative outcomes. *Anesthesiology* 130, 634–650. <https://doi.org/10.1097/ALN.0000000000002510>.

Tchivileva, I.E., Lim, P.F., Smith, S.B., Slade, G.D., Diatchenko, L., McLean, S.A., and Maixner, W. (2010). Effect of catechol-O-methyltransferase polymorphism on response to propranolol therapy in chronic musculoskeletal pain: a randomized, double-blind, placebo-controlled, crossover pilot study. *Pharmacogenetics Genom.* 20, 239–248. <https://doi.org/10.1097/FPC.0b013e328337f9ab>.

Tchivileva, I.E., Hadgraft, H., Lim, P.F., Di Giosia, M., Ribeiro-Dasilva, M., Campbell, J.H., Willis, J., James, R., Herman-Giddens, M., Fillingim, R.B., et al. (2020). Efficacy and safety of propranolol for treatment of temporomandibular disorder pain: a

randomized, placebo-controlled clinical trial. *Pain* 161, 1755–1767. <https://doi.org/10.1097/j.pain.0000000000001882>.

Turnbull, C., Scott, R.H., Thomas, E., Jones, L., Murugaesu, N., Pretty, F.B., Haili, D., Baple, E., Craig, C., Hamblin, A., et al. (2018). The 100 000 Genomes Project: bringing whole genome sequencing to the NHS. *BMJ* 361, k1687. <https://doi.org/10.1136/bmj.k1687>.

Wróbel, A., Serefko, A., Wozniak, A., Kociszewski, J., Szopa, A., Wisniewski, R., and Poleszak, E. (2020). Duloxetine reverses the symptoms of overactive bladder coexisting with depression via the central pathways. *Pharmacol. Biochem. Behav.* 189, 172842. <https://doi.org/10.1016/j.pbb.2019.172842>.

Yao, L., Zhang, Y., Wei, B., Jin, Z., Zhang, R., Zhang, Y., and Chen, Q. (2017). Incorporating knowledge graph embeddings into topic modeling. In *Thirty-first AAAI conference on artificial intelligence*, 31 (Palo Alto, California USA: AAAI Press), pp. 3119–3126.

STAR★METHODS

KEY RESOURCES TABLE

REAGENT or RESOURCE	SOURCE	IDENTIFIER
Deposited data		
UK Biobank dataset	https://biobank.ctsu.ox.ac.uk	The UK Biobank data access has been approved by McGill IRB under the project title "A replication study of pain interactions with comorbidities". The approval number is A03-M20-21B.
Software and algorithms		
Graph-embedded topic model	This paper	https://github.com/li-lab-mcgill/getm

RESOURCE AVAILABILITY

Lead contact

Requests for data and requests for additional information should be directed to the lead contact, Yue Li (yueli@cs.mcgill.ca).

Material availability

This study did not generate new unique reagents.

Data and code availability

- The UK Biobank data access has been approved by McGill IRB under the project title "A replication study of pain interactions with comorbidities". The approval number is A03-M20-21B.
- All code associated with this paper can be freely accessed and downloaded via <https://github.com/li-lab-mcgill/getm>.
- Any additional information required to reanalyze the data reported in this paper is available from the [lead contact](#) upon request.

METHOD DETAILS

UKB data processing

For conditions data, we used UKB datafield 20002, which records self-reported non-cancer diseases for each subject. This was collected by questionnaire during participant interviews. The participants were asked whether or not they had been diagnosed with certain conditions (heart attack, angina, stroke, high blood pressure, blood clot in leg, blood clot in lung, emphysema/chronic bronchitis, asthma or diabetes), were asked to add any other conditions and to provide a date of diagnosis for each when possible <https://biobank.ctsu.ox.ac.uk/crystal/field.cgi?id=20002>. Medication usage data was similarly collected during participant interviews prompted by a question on regular prescription medication use, resulting in datafield 20003 which contains treatment/medication codes <https://biobank.ctsu.ox.ac.uk/crystal/field.cgi?id=20003>.

We kept 457,461 individuals of European descent to reduce potential confounding by ethnic group. In total, 802 active ingredients were kept as medications and 443 conditions were extracted. Here we encoded the medications and conditions as binary variables. Only baseline visit data was included.

UKB pain-related phenotype labels extractions

We sought to associate pain phenotypes with comorbid conditions and medications. To this end, we used the pain-related phenotypes as the label data and patient-dependent topic mixture derived from the 443 conditions and 802 medications as the input features in a post-hoc supervised topic analysis as detailed in [CMK pain prediction](#).

The pain phenotypes were collected through questionnaire. We drew on datafield 6159; participants were asked “In the last month have you experienced any of the following that interfered with your usual activities? (You can select more than one answer).” If they answer “yes” to pain at any site, for example, back pain, they were further asked “Have you had back pains for more than 3 months?” (datafield 3571), and so on. In total, we collected 9 pain-related labels with data field identifiers listed as follows:

1. pain type(s) experienced in last month: 6,159
2. headaches for 3+ months: 3,799
3. facial pains for 3+ months: 4,067
4. neck/shoulder pain for 3+ months: 3,404
5. back pain for 3+ months: 3,571
6. stomach/abdominal pain for 3+ months: 3,741
7. hip pains for 3+ months: 3,414
8. knee pains for 3+ months: 3,773
9. general pain for 3+ months: 2,956

We can also query the UKB website to obtain relevant information for each field by replacing x in <https://biobank.ndph.ox.ac.uk/ukb/field.cgi?id=x> with the above data field identifiers.

Graph-embedded topic model details

Graph-embedded topic model (GETM) generative process

We formulated the problem of modelling discrete patient healthcare data into a topic modeling problem. In particular, we treat each of the D patient health records as a document and each observed feature in the record as a word sampled from a defined vocabulary. In our case, we have two vocabularies covering $M = 802$ medications and $C = 443$ conditions, respectively. With this analogy, each patient record is represented as a mixture of latent topics. In the original latent Dirichlet allocation (LDA) model (Blei et al., 2003b), a topic distribution over a vocabulary is defined as an independent Dirichlet prior $\beta_k \sim \text{Dirichlet}(\tau_\beta)$, where $\sum_{w=1}^V \beta_{k,w} = 1$ over V words. Inspired by a more recent work called embedded topic model (ETM) (Dieng et al., 2019), we decomposed the K topic distributions over medications $\beta^{(med)}$ into a medication-defined topic embedding $\alpha^{(med)} \in \mathbb{R}^{K \times L_1}$, and a medication embedding $\rho^{(med)} \in \mathbb{R}^{L_1 \times M}$, where L_1 denotes the medication embedding dimension and M denotes the number of unique medications. Similarly, the topic distribution over conditions $\beta^{(cond)}$ is proportional to the inner product of the condition-defined topic embedding $\alpha^{(cond)} \in \mathbb{R}^{K \times L_2}$, and condition embedding $\rho^{(cond)} \in \mathbb{R}^{L_2 \times C}$, where L_2 denotes the condition embedding dimension and C denotes the number of unique conditions.

For the d^{th} patient record ($d \in \{1, \dots, D\}$), the generative process starts by drawing the topic mixture θ_d from a logistic normal distribution $\theta_d \sim \mathcal{LN}(0, \mathbf{I})$:

$$\eta_d \sim \mathcal{N}(0, \mathbf{I}); \quad \theta_d = \text{softmax}(\eta_d) = \frac{\exp(\eta_d)}{\sum_{k=1}^K \exp(\eta_{d,k})} \quad (\text{Equation 1})$$

We then draw the i^{th} medication token or the j^{th} condition token from two respective categorical distributions:

$$w_{i,d}^{(med)} = \prod_m r_{m,d}^{[w_{i,d}^{(med)} = m]}, \quad w_{j,d}^{(cond)} = \prod_c r_{c,d}^{[w_{j,d}^{(cond)} = c]} \quad (\text{Equation 2})$$

where

$$r_{m,d} = \sum_k \theta_{d,k} \beta_{k,m}^{(med)} = \theta_{d,\cdot} \beta_{\cdot,m}, \quad r_{c,d} = \sum_k \theta_{d,k} \beta_{k,c}^{(cond)} = \theta_{d,\cdot} \beta_{\cdot,c} \quad (\text{Equation 3})$$

Here we use notation $\theta_{d,\cdot}$ to denote a $1 \times K$ row vector for the d^{th} patient record and $\beta_{\cdot,m}$ (and $\beta_{\cdot,c}$) to denote a $K \times 1$ vector for the m^{th} medication (and the c^{th} condition).

The k^{th} topic distribution for the m^{th} medication (or the c^{th} condition) is defined as the product of the corresponding topic embedding and medication (or condition) embedding followed by softmax normalization over the corresponding vocabularies:

$$\beta_{k,m}^{(med)} = \frac{\exp(\alpha_{k..}^{(med)} \rho_{..m}^{(med)})}{\sum_{m'} \exp(\alpha_{k..}^{(med)} \rho_{m'}^{(med)})}, \quad \beta_{k,c}^{(cond)} = \frac{\exp(\alpha_{k..}^{(cond)} \rho_{..c}^{(cond)})}{\sum_{c'} \exp(\alpha_{k..}^{(cond)} \rho_{c'}^{(cond)})} \quad (\text{Equation 4})$$

Here we treat both the topic embeddings $\alpha^{(med)}$ and $\alpha^{(cond)}$ and condition and medication embeddings $\rho^{(med)}$ and $\rho^{(cond)}$ as point estimates without imposing any prior distribution.

For the ease of mathematical expression below, we denote $x_{m,d}^{(med)}$ (and $x_{c,d}^{(cond)}$) as the frequency of the m^{th} medication (and the c^{th} condition) for the d^{th} patient record. Formally, modeling the likelihood of the count frequency simply requires reformulating the above multinomial likelihood:

$$p(x_{m,d}^{(med)} | \theta_d, \Theta) = \prod_i [r_{m,d}^{(med)}]^{w_{i,d}^{(med)} = m} = [r_{m,d}^{(med)}]^{\sum_i [w_{i,d}^{(med)} = m]} = [r_{m,d}^{(med)}]^{x_{m,d}^{(med)}} \quad (\text{Equation 5})$$

$$p(x_{c,d}^{(cond)} | \theta_d, \Theta) = \prod_j [r_{c,d}^{(cond)}]^{w_{j,d}^{(cond)} = c} = [r_{c,d}^{(cond)}]^{\sum_j [w_{j,d}^{(cond)} = c]} = [r_{c,d}^{(cond)}]^{x_{c,d}^{(cond)}} \quad (\text{Equation 6})$$

where θ_d is the patient topic mixture of patient d and $\Theta = \{\alpha^{(med)}, \rho^{(med)}, \alpha^{(cond)}, \rho^{(cond)}\}$.

The vector $\mathbf{x}_{..d}^{(med)}$ and $\mathbf{x}_{..d}^{(cond)}$ denote the frequency of all medications and all conditions for the d^{th} patient record, respectively. The entire data can then be represented as a $M \times D$ matrix $\mathbf{X}^{(med)}$ and a $C \times D$ matrix $\mathbf{X}^{(cond)}$ and modelled by Equations (5) and (6):

$$p(\mathbf{X}^{(med)} | \theta, \Theta) = \prod_d \prod_m p(x_{m,d}^{(med)} | \theta_d, \Theta) = \prod_d \prod_m [r_{m,d}^{(med)}]^{x_{m,d}^{(med)}} \quad (\text{Equation 7})$$

$$p(\mathbf{X}^{(cond)} | \theta, \Theta) = \prod_d \prod_c p(x_{c,d}^{(cond)} | \theta_d, \Theta) = \prod_d \prod_c [r_{c,d}^{(cond)}]^{x_{c,d}^{(cond)}} \quad (\text{Equation 8})$$

Inference of θ_d 's and learning of Θ are described next.

Model inference and estimation

To train GETM, we want to maximize the marginal likelihood of the individuals with respect to $\Theta = \{\alpha^{(med)}, \rho^{(med)}, \alpha^{(cond)}, \rho^{(cond)}\}$:

$$\alpha^{(med)}, \rho^{(med)}, \alpha^{(cond)}, \rho^{(cond)} \leftarrow \max_{\Theta} \mathcal{L}(\Theta) = \max_{\Theta} \sum_{d=1}^D \log p(\mathbf{x}_d^{(.)} | \Theta) \quad (\text{Equation 9})$$

where $\mathbf{x}_d^{(.)}$ is word frequency vector (i.e., the bag of words of medications and conditions for individual d). This marginal likelihood is intractable to compute since it involves integrating out the logistic normal topic mixture variable θ_d , $\forall d \in \{1, \dots, D\}$:

$$\begin{aligned} \log p(\mathbf{x}_d^{(.)} | \Theta) &= \log \int p(\eta_d) \prod_{m=1}^M p(x_{d,m}^{(med)} | \eta_d, \alpha^{(med)}, \rho^{(med)}) \prod_{c=1}^C p(x_{d,c}^{(cond)} | \eta_d, \alpha^{(cond)}, \rho^{(cond)}) d\eta_d \\ &= \log \int p(\eta_d) \prod_{m=1}^M (f(\eta_{d..}) f(\alpha^{(med)} \rho_{..m}^{(med)}))^{x_{m,d}^{(med)}} \end{aligned} \quad (\text{Equation 10})$$

$$\prod_{c=1}^C (f(\eta_{d..}) f(\alpha^{(cond)} \rho_{..c}^{(cond)}))^{x_{c,d}^{(cond)}} d\eta_d \quad (\text{Equation 11})$$

$$\equiv \log \int p(\eta_d) \prod_{m=1}^M (\theta_{d..} \beta_{..m}^{(med)})^{x_{m,d}^{(med)}} \prod_{c=1}^C (\theta_{d..} \beta_{..c}^{(cond)})^{x_{c,d}^{(cond)}} d\eta_d \quad (\text{Equation 12})$$

where $f(\cdot)$ denotes the softmax function and $\theta_{d..} \beta_{..m}^{(med)}$ is the inner matrix product of the $1 \times K$ vector $\theta_{d..}$ and the $K \times 1$ vector $\beta_{..m}^{(med)}$. Same for conditions.

Therefore, we took an *Amortized Variational Inference* (AVI) approach to approximate the above intractable integral. This is quite similar to the one described by Kingma and Welling (2014) (Kingma and Welling, 2014) and by Dieng et al., (2019) (Dieng et al., 2019). For the sake of completeness, we describe it below.

To approximate true posterior $p(\boldsymbol{\eta}_d | \mathbf{x}_d)$, we define the proposed distribution as a Gaussian with mean and variance produced by a feedforward neural network with parameters \mathbf{W}_θ :

$$q(\boldsymbol{\eta}_d | \mathbf{x}_d, \mathbf{W}_\theta) = \mathcal{N}(\boldsymbol{\mu}_d, \text{diag}(\boldsymbol{\sigma}_d) \mathbf{I}), \quad [\boldsymbol{\mu}_d, \log \boldsymbol{\sigma}_d^2] = \text{NNET}(\mathbf{x}_d; \mathbf{W}_\theta) \quad (\text{Equation 13})$$

The evidence lower bound (ELBO) of the above log likelihood is approximated as follows:

$$\log p(\mathbf{x}_d^{(\cdot)} | \Theta) = \log \int q(\boldsymbol{\eta}_d) \frac{p(\boldsymbol{\eta}_d)}{q(\boldsymbol{\eta}_d)} \prod_{m=1}^M (\boldsymbol{\theta}_{d,\cdot} \boldsymbol{\beta}_{\cdot,m}^{(med)})^{x_{m,d}^{(med)}} \prod_{c=1}^C (\boldsymbol{\theta}_{d,\cdot} \boldsymbol{\beta}_{\cdot,c}^{(cond)})^{x_{m,d}^{(cond)}} d\boldsymbol{\eta}_d \quad (\text{Equation 14})$$

$$\geq \int q(\boldsymbol{\eta}_d) \log p(\boldsymbol{\eta}_d) d\boldsymbol{\eta}_d \quad (\text{Equation 15})$$

$$+ \int q(\boldsymbol{\eta}_d) \sum_{m=1}^M x_{m,d}^{(med)} \log(\boldsymbol{\theta}_{d,\cdot} \boldsymbol{\beta}_{\cdot,m}^{(med)}) d\boldsymbol{\eta}_d \quad (\text{Equation 16})$$

$$+ \int q(\boldsymbol{\eta}_d) \sum_{c=1}^C x_{m,d}^{(cond)} \log(\boldsymbol{\theta}_{d,\cdot} \boldsymbol{\beta}_{\cdot,c}^{(cond)}) d\boldsymbol{\eta}_d - \int q(\boldsymbol{\eta}_d) \log q(\boldsymbol{\eta}_d) d\boldsymbol{\eta}_d$$

$$= \mathbb{E}_{q(\boldsymbol{\eta})} [\log p(\boldsymbol{\eta}_d)]$$

$$+ \mathbb{E}_{q(\boldsymbol{\eta})} [\log p(\mathbf{x}_d^{(med)} | \boldsymbol{\eta}_d, \boldsymbol{\beta}^{(med)})]$$

$$+ \mathbb{E}_{q(\boldsymbol{\eta})} [\log p(\mathbf{x}_d^{(cond)} | \boldsymbol{\eta}_d, \boldsymbol{\beta}^{(cond)})]$$

$$- \mathbb{E}_{q(\boldsymbol{\eta}_d)} [\log q(\boldsymbol{\eta}_d)] \quad (\text{Equation 17})$$

$$= \sum_{t \in \{\text{med}, \text{cond}\}} \mathbb{E}_{q(\boldsymbol{\eta})} [\log p(\mathbf{x}_d^{(t)} | \boldsymbol{\eta}_d, \boldsymbol{\beta}^{(t)})] - \text{KL}[q(\boldsymbol{\eta}_d) \| p(\boldsymbol{\eta}_d)] \quad (\text{ELBO}) \quad (\text{Equation 18})$$

$$\approx \sum_{t \in \{\text{med}, \text{cond}\}} \log p(\mathbf{x}_d^{(t)} | \hat{\boldsymbol{\eta}}_d, \boldsymbol{\beta}^{(t)}) - \text{KL}[q(\hat{\boldsymbol{\eta}}_d) \| p(\hat{\boldsymbol{\eta}}_d)] \quad (\text{Equation 19})$$

where

$$\log p(\mathbf{x}_d^{(t)} | \hat{\boldsymbol{\eta}}_d, \boldsymbol{\beta}^{(t)}) = \sum_{m=1}^M x_{m,d}^{(med)} \log(\hat{\boldsymbol{\theta}}_{d,\cdot} \boldsymbol{\beta}_{\cdot,m}^{(med)}) + \sum_{c=1}^C x_{c,d}^{(cond)} \log(\hat{\boldsymbol{\theta}}_{d,\cdot} \boldsymbol{\beta}_{\cdot,c}^{(cond)})$$

$$\mathbb{E}_{q(\boldsymbol{\eta}_d | \mathbf{x}_d)} [\log p(\hat{\boldsymbol{\eta}}_d)] = \mathbb{E}_{q(\boldsymbol{\eta}_d | \mathbf{x}_d)} [\log \mathcal{N}(\hat{\boldsymbol{\eta}}_d; \mathbf{0}, \mathbf{1})] = -\frac{1}{2} (\log(2\pi) + \boldsymbol{\mu}_d^2 + \boldsymbol{\sigma}_d^2)$$

$$\mathbb{E}_{q(\boldsymbol{\eta}_d | \mathbf{x}_d)} [\log q(\hat{\boldsymbol{\eta}}_d)] = -H(\hat{\boldsymbol{\eta}}) = -\frac{1}{2} (\log \boldsymbol{\sigma}_d^2 + 1 + \log(2\pi))$$

$$\mathbb{E}_{q(\boldsymbol{\eta}_d | \mathbf{x}_d)} [\log q(\hat{\boldsymbol{\eta}}_d) - \log p(\hat{\boldsymbol{\eta}}_d)] \equiv \text{KL}[q(\hat{\boldsymbol{\eta}}_d) \| p(\hat{\boldsymbol{\eta}}_d)] = -\frac{1}{2} (\log \boldsymbol{\sigma}_d^2 + 1 - \boldsymbol{\mu}_d^2 - \boldsymbol{\sigma}_d^2)$$

Here Equation (19) uses the sampled $\hat{\boldsymbol{\eta}}_d$ from the feedforward network encoder as a surrogate to approximate the intractable variational expectation w.r.t. $\boldsymbol{\eta}_d$. We re-parameterized the Gaussian distribution in Equation (13) to enable stochastic gradient calculation of the ELBO w.r.t. \mathbf{W}_θ and AVI training of the encoder network *NNET*:

$$q(\boldsymbol{\eta}_d | \mathbf{x}_d, \mathbf{W}_\theta) = \boldsymbol{\mu}_d + \text{diag}(\boldsymbol{\sigma}_d \mathcal{N}(0, \mathbf{1})), \quad [\boldsymbol{\mu}_d, \log \boldsymbol{\sigma}_d^2] = \text{NNET}(\mathbf{x}_d; \mathbf{W}_\theta) \quad (\text{Equation 20})$$

$$\frac{\partial \text{ELBO}}{\partial \mathbf{W}_\theta} = \frac{\text{NNET}(\mathbf{x}_d; \mathbf{W}_\theta)}{\partial \mathbf{W}_\theta} \quad (\text{Equation 21})$$

Model specifications

The encoders for GETM, partial GETMs and ETMs are 2-layer neural networks with hidden sizes of 128, ReLU activations (Glorot et al., 2011) and 1D batch normalization (Ioffe and Szegedy, 2015). The medication embedding dimension and condition embedding dimension were both 128 and they were fixed during training of the models with KG-informed pre-trained embeddings. The models were optimized with Adam optimizer at 0.01 learning rate. We trained each model with batch size of 100. For inferring individual-topic mixtures used by prediction tasks, the topic number was set to be 128.

Knowledge graph construction for learning the embedding of UKB conditions and medications

In the original ETM (Dieng et al., 2019), the word embedding \mathbf{p} can be either learned from the documents or pre-trained on a separate large data corpus such as Wikipedia using the word2vec approach. The latter approach allows leveraging the contextual data to improve the topic modeling. As one of our main contributions, we sought to develop a simple framework that exploits ETM's ability to incorporate a pre-trained medication embedding $\mathbf{p}^{(med)}$ and condition embedding $\mathbf{p}^{(cond)}$ from their taxonomic knowledge graphs (KG). Leveraging the structural KG information and the internal relational information among medications and conditions in a principled way is beneficial to modeling the UKB phenotype data and other EHR data in general, which are often sparse, noisy, and bias.

Among the graph representational learning methods, we chose node2vec (Grover and Leskovec, 2016) as it is an unsupervised approach that can directly learn from the KG without linking to a prediction task as in other methods. Specifically, we applied node2vec to separately learn the node embedding of the 443 conditions or 802 medications from their hierarchical trees (Figure 1).

For the KG of the 443 conditions, we constructed a tree graph using coding taxonomy designed by the UKB team (i.e., datafield 20002). The tree describes the topology of the conditions with 473 nodes and 4 hierarchical levels.

For the KG of the 802 medications, we constructed a KG based on the Anatomical Therapeutic Chemical (ATC) classification system <https://www.whocc.no/>. The entire tree is composed of 5 levels. We first kept the top 4 levels of ATC, of which the first level contains main anatomical or pharmacological groups; the second level includes pharmacological or therapeutic subgroups; and in the third and fourth levels are chemical, pharmacological or therapeutic subgroups. To harmonize the UKB medications with the ATC graph, we mapped the names of active ingredients from UKB datafield 20003 to the fifth level codes of ATC, which are chemical substances. Some UKB medication codes were mapped to multiple ATC fifth level codes as they could belong to different subgroups depending on their usages. In that case, we replaced all of the mapped ATC fifth level codes with the same corresponding UKB medication code of the active ingredient in the second step. The final medication graph contains 2561 nodes in total.

Both the condition graph and the medication graph were treated as undirected graphs as input to the node2vec model (Figure 1).

Algorithm 1. Topic modelling with GETM

Input : UKB phenotype matrices $\mathbf{X}^{(cond)}$ and $\mathbf{X}^{(med)}$, node2vec-pretrained embedding for conditions and medications $\mathbf{p}^{(cond)}$ and $\mathbf{p}^{(med)}$

Initialize :condition-defined topic embedding $\alpha^{(cond)}$, medication-defined topic embedding $\alpha^{(med)}$, encoder neural network W_θ

Output :Learned topic embeddings $\hat{\alpha}^{(cond)}$, $\hat{\alpha}^{(med)}$ and encoder network parameters \hat{W}_θ for patient topic mixtures

1 for epoch $i = 1, 2, \dots$ do

2 for each batch B do

3 Get normalized bag-of-words $\tilde{\mathbf{X}}_B = [\tilde{\mathbf{X}}_B^{(cond)}, \tilde{\mathbf{X}}_B^{(med)}]$

4 Sample $\hat{\theta}_B$ using $NNET(\tilde{\mathbf{X}}_B; W_\theta)$ by (Equation 21)

5 Compute approximate ELBO($\hat{\theta}_B$, $\alpha^{(med)}$, $\alpha^{(cond)}$) using (Equation 19)

6 Update parameters $\alpha^{(cond)}$, $\alpha^{(med)}$, W_θ by backpropagation to maximize ELBO

7 end for

8 end for

Summary of the GETM learning algorithm

GETM learning algorithm is summarized in Algorithm 1.

Ablation experiments

To gain a good understanding of how each component of our method contribute to the overall improvements, we performed ablation experiments. In particular, we compared GETM with 8 ablated baseline models (Table S1):

Topic quality evaluation

For medication, the topic coherence was calculated as:

$$TC_{med} = \frac{1}{K} \sum_{k=1}^K \frac{m_k}{n} \quad (\text{Equation 22})$$

where n is the number of top medications, $m_k \leq n$ is the maximum number of medications that are from the same category across all categories, and K is the number of topics. We set $n = 5$ because of the softmax-normalization of each topic, where only the top 5 conditions or medications under each topic has non-negligible probabilities.

To avoid overestimating the topic quality, the categories we used to evaluate the topic coherence were not processed from the ATC graph since it was involved in the pre-trained model. Instead, we employed 59 categories which are physician-curated and pain-focused (Table S7).

For conditions, we did not have any external gold-standard reference to compare against. Therefore, we calculated the topic coherence for conditions based on the co-occurrence of the top conditions under the same topic observed from the same patient record. Specifically, for any top two conditions under the same topic, we calculated the average pointwise mutual information:

$$TC_{cond} = \frac{1}{K} \sum_{k=1}^K \frac{1}{n(n-1)} \sum_{i=1}^n \sum_{j=i+1}^n f(c_i^k, c_j^k) \quad (\text{Equation 23})$$

where c_i^k is the i^{th} top most likely condition in topic k , and $f(\cdot, \cdot)$ is normalized pointwise mutual information:

$$f(c_i, c_j) = \frac{\log \frac{P(c_i, c_j)}{P(c_i)P(c_j)}}{-\log P(c_i, c_j)} \quad (\text{Equation 24})$$

where $P(c_i, c_j)$ is the probability of condition i and condition j co-occurring in one individual and $P(c_i)$ is the marginal probability of condition i . The probabilities were approximated by empirical counts.

Topic diversity is the percentage of total unique features of the top 5 features of all topics. We first chose top 5 features from K topics. Among $5 \times K$ features there are U unique words. The topic diversity is calculated as $TQ = \frac{U}{5 \times K}$. Therefore, the higher topic diversity the more diverse the topics are.

Study of medication and condition relations

In GETM, since the encoder takes both medication and condition information as input (Figure 1), the top medications and conditions from the same topic (i.e. the same index) are related. To quantitatively measure the ability of the model to capture known condition-medication relations, we first obtained all of the combinations from top 3 conditions and top 3 medications for each topic. We then counted the number of condition-medication pairs which are known to be related. The reference of known pairs was extracted from CTD <http://ctdbase.org/> and DrugBank <https://go.drugbank.com/>. We eventually mapped 222 conditions and 529 medications from UKB to these two databases. As a result, we had 2444 positive pairs of which the medication has treatment effects on the condition and 3231 negative pairs of which the condition belongs to the adverse effects of the medication. The number of known pairs discovered by our model and baselines were compared using different number of topic numbers (Table 1).

Data imputation

Using the generative model, we can (1) impute missing entries including conditions or medications and (2) impute medications based on only clinical conditions. Accordingly, two sets of experiments were performed to evaluate how well our model could complete these two tasks. We split the UKB data into 80%

training and 20% test. To simulate missing entries, we randomly masked 50% of medications and conditions for test individuals. Then we calculated reconstruction error as the negative log-likelihood of the held-out data:

$$-\log p(\mathbf{X}|\boldsymbol{\theta}_d, \boldsymbol{\beta}) = \frac{1}{D'} \sum_{d=1}^{D'} -\mathbf{x}_d \log(\widehat{\mathbf{x}}_d^\top) \quad (\text{Equation 25})$$

where D' is the number of test individuals and $\widehat{\mathbf{x}}_d = \boldsymbol{\theta}_d \boldsymbol{\beta}$

For medication imputation experiment, we masked the entire medication data of the test individuals and then reconstructing the medication matrix. This experiment mimics the scenario that we recommend relevant medications based only on individuals' condition history.

In addition to reconstruction error, we also calculated recall and precision rates at the top 5 predictions: recall @5 and precision @5. Specifically, we sorted the predicted probabilities of medications and chose top 5 medications, of which we calculated recall and precision with respect to medications that the individual took as true labels. The recall and precision are calculated as: $\text{Recall} = \frac{\text{TruePositives}}{\text{TruePositives} + \text{FalseNegatives}}$ and $\text{Precision} = \frac{\text{TruePositives}}{\text{TruePositives} + \text{FalsePositives}}$, respectively.

CMK pain prediction

CMK pain was defined as musculoskeletal pain lasting at least three months localized to any of five body sites: knee, back, neck/shoulder, face or hip. The label information was obtained from the UKB pain questionnaire as described in [STAR Methods](#).

We used individual topic mixture $\boldsymbol{\theta}_d$'s to predict whether a individual d has chronic musculoskeletal pain at the time of the visit as the binary outcome y_d . We split data to 80% of training set and 20% of testing set. GETM was first trained with training set to get $\boldsymbol{\theta}^{(\text{train})}$. Then we fit logistic regression model $\text{logit}(y_d) \sim \boldsymbol{\theta}_d^{(\text{train})} \boldsymbol{\omega}$, where $\boldsymbol{\omega}$ is a $K \times 1$ vector of the linear coefficients corresponding to the predictive weights of the K topics.

To evaluate model performance, we first obtained $\boldsymbol{\theta}^{(\text{test})}$ with testing set using the trained GETM (i.e., the encoder network parameters \mathbf{W}_θ) and then predicted CMK pain status using the trained logistic regression coefficients $\hat{y}^{(\text{test})} = \boldsymbol{\theta}^{(\text{test})} \hat{\boldsymbol{\omega}}$.

Odds ratios of CMK-related conditions or medications

We calculated odds ratios of each condition and medication with respect to CMK pain as the outcome ([Figure S4](#)). The odds ratios were used for two purposes: (1) obtaining conditions and medications that are obvious for chronic musculoskeletal pain prediction and to be removed for some experiments; (2) constructing a baseline reference to construct a curated list of pain-related conditions and medication by the physician in our team. For each feature (condition or medication), we formed a contingency table as below. The odds ratios were calculated as $\text{OR} = \frac{ad}{bc}$.

Feature	Pain	
	Yes	No
Yes	a	b
No	c	d

Feature filtering schemes in removing obvious CMK-related conditions and medications

Let \mathcal{C} and \mathcal{M} be the complete sets of 443 conditions and 802 medications, respectively. Based on the odds ratios, we identified 50 conditions ($\mathcal{C}1$) and 150 medications ($\mathcal{M}1$) that are significantly associated with the CMK pain ([Figure S4](#)). A physician also provides a general pain related condition list ($\mathcal{C}2$), a musculoskeletal pain related condition list ($\mathcal{C}3$) and a general pain related medication set ($\mathcal{M}2$). We created three filtered condition sets: (1) $\mathcal{C} - \mathcal{C}2$; (2) $\mathcal{C} - \mathcal{C}3$; (3) $\mathcal{C} - (\mathcal{C}1 \cup \mathcal{C}2)$, where using the set notations $\mathcal{C} - \mathcal{C}1$ denotes condition set without $\mathcal{C}1$ and $\mathcal{C}1 \cup \mathcal{C}2$ denotes the union of $\mathcal{C}1$ and $\mathcal{C}2$.

We also created two filtered medication sets: (1) $\mathcal{M} - \mathcal{M}_1$; (2) $\mathcal{M} - (\mathcal{M}_1 \cup \mathcal{M}_2)$. These filtered condition sets and filtered medication sets formed six experiment settings plus the full set (Table S2), which we used to explore CMK prediction performance as a function of reduced feature sets (Figure 5).

Calculating importance of conditions and medications for CMK pain

Using the $K \times 1$ coefficients from logistic regression $\hat{\omega}$ and the topic embeddings, we calculated the relevance of medications and conditions to CMK pain as follows (Figures 6A and 6B).

$$\mathbf{v} = \hat{\omega}^T \hat{\boldsymbol{\beta}} \quad (\text{Equation 26})$$

We selected the top $N \in \{10, 30, 50\}$ medications and conditions from $\mathbf{v}^{(med)}$ and $\mathbf{v}^{(cond)}$. We then calculated the proportions of these top N medications or conditions overlapping with the pain-related lists created by a physician. The results are visualized as barplots in Figures 6C and 6D.

Characterizing pain types from 7 body sites and pain all over the body

We used similar approaches as described for the CMK pain predictions to characterize other pain types. For all experiments, we used three different filtered feature sets of conditions and medications, which are (1) all: non-filtered conditions and medications, (2) general: physician-curated general pain-related conditions and medications were filtered, and (3) least: physician-curated / top odds ratio calculation conditions and medications were filtered.

For each pain type p , we calculated the importance weight that a medication or condition positively or negatively related to the pain using logistic regression coefficient for predicting that pain type $\hat{\omega}_p$ and GETM topic-feature mixture $\boldsymbol{\beta}$ of certain feature filtering regime (Figure 6B). Specifically, the way to calculate importance weight $w_{c,p}$ for a positive relation between a condition c and the pain type p is:

$$\hat{\omega}_{c,p}^{(+)} = \begin{cases} 0 & \text{if } \hat{\omega}_{c,p} < 0 \\ \hat{\omega}_{c,p} & \text{otherwise} \end{cases}, \quad w_{c,p}^{(+)} = \frac{[\hat{\omega}_p^{(+)}]^T \boldsymbol{\beta}_{c..}^{(cond)}}{\sum_c \hat{\omega}_{c,p}^{(+)}} \quad (\text{Equation 27})$$

Similarly, we calculated importance weight for negative relationship between a condition c and pain type p :

$$\hat{\omega}_{c,p}^{(-)} = \begin{cases} 0 & \text{if } \hat{\omega}_{c,p} > 0 \\ \hat{\omega}_{c,p} & \text{otherwise} \end{cases}, \quad w_{c,p}^{(-)} = \frac{[\hat{\omega}_p^{(-)}]^T \boldsymbol{\beta}_{c..}^{(cond)}}{\sum_c \hat{\omega}_{c,p}^{(-)}} \quad (\text{Equation 28})$$

Calculation for the positive and negative relations between every medication and every pain type is the same.

QUANTIFICATION AND STATISTICAL ANALYSIS

Topic model were evaluated using topic quality scores as described in Section 8. The topic quality scores were computed for five repeated experiments with random initialization for each method (Figure 3 and Supplemental Figure S2). One-sided t-test were performed to compare the topic quality scores from the proposed Graph-embedded topic model (GETM) with the baseline methods. Prediction of chronic musculoskeletal pain phenotypes were evaluated using the receiver operating characteristic (ROC) curves and area under the ROC curve (AUROC) as well as precision-recall curves (PRC) and area under the PRC (AUPRC). Feature selection were performed using Fisher exact test as detailed above.