# VisAGE: Integrating external knowledge into electronic medical record visualization

Edward W Huang, Sheng Wang, and ChengXiang Zhai[†]

*Department of Computer Science,*
*University of Illinois at Urbana-Champaign,*
*Urbana, IL, USA*
[†]*Email: czhai@illinois.edu*

In this paper, we present VisAGE, a method that visualizes electronic medical records (EMRs) in a low-dimensional space. Effective visualization of new patients allows doctors to view similar, previously treated patients and to identify the new patients' disease subtypes, reducing the chance of misdiagnosis. However, EMRs are typically incomplete or fragmented, resulting in patients who are missing many available features being placed near unrelated patients in the visualized space. VisAGE integrates several external data sources to enrich EMR databases to solve this issue. We evaluated VisAGE on a dataset of Parkinson's disease patients. We qualitatively and quantitatively show that VisAGE can more effectively cluster patients, which allows doctors to better discover patient subtypes and thus improve patient care.

*Keywords*: Electronic medical records; Data integration; Knowledge graphs; Visualization.

## 1. Introduction

In modern healthcare settings, doctors record the details of patient visits in electronic medical records (EMRs), which are then collected in databases. At first, EMR systems were poorly implemented; initial studies reported that they reduced physician productivity and lacked data sharing capabilities.[1] However, recent advances have improved recordkeeping and decision support. For example, EMR systems have enhanced productivity in physician workloads[2] and have increased the delivery of health behavior counseling.[3]

Despite these advances, EMR systems can still benefit from human interpretation, which allows for exploratory analysis and more control over decision-making.[4] However, human interaction with EMR systems has been hindered. 37% of participants in a previous study reported that interacting with their EMR databases was too time consuming.[5] Another study showed that when using EMRs, nurses face challenges that can threaten quality and safety of care.[6] Both shortcomings can be addressed with information visualization, which can aid doctors in processing and understanding complex, high-dimensional EMR data.

In particular, EMR visualization in a two-dimensional space is useful for observing disease subtypes in patient clusters. Coherent clusters may elucidate a patient's most significant characteristics by visualizing his or her proximity to successfully diagnosed patients.[7] For example, thoracic aortic dissections are commonly misdiagnosed as acute myocardial infarctions (MIs).[8] Misdiagnosis in these cases is extremely harmful, as patients with aortic dissections treated for MIs have mortality rates similar to that of untreated patients. Despite this risk, patients with aortic dissections have a misdiagnosis rate of 39%.[9] Fortunately, the most telltale signs

of aortic dissection (age, onset of pain, and syncope) are readily available in EMRs.[10,11] An effective visualization would utilize these features to place an undiagnosed aortic dissection patient near similar patients, reducing the chance of misdiagnosis.

Unfortunately, designing an effective visualization system is a complicated task, as EMRs are high-dimensional sources of data that consist of thousands of features. Because there are many potential medical tests that a patient can take, but only a handful are relevant to each patient's condition, EMRs also tend to be sparse. Additionally, EMRs are typically fragmented or incomplete due to human error. These reasons make it challenging to correctly group together similar patients, leading to poor or even misleading visualizations.

To address these challenges, we present **Vis**ualization **A**ssisted by Knowledge **G**raph **E**nrichment (VisAGE), a method that enriches patient records with a knowledge graph built from external databases. These databases include protein-protein interactions, genomic data, and drug-chemical associations. Performing network embedding on the knowledge graph allows us to infer associations among different types of data to accommodate inexact matchings of related features, which can alleviate data sparsity in EMRs. A major novelty of VisAGE is that it is the first to use all of these data sources in EMR visualization. In the rest of the paper, we describe our dataset, the details of VisAGE, and our evaluation process.

## 2. Data Description

While VisAGE is a general method that can be applied to any set of EMRs, we chose the Parkinson's Progression Markers Initiative (PPMI) dataset[12] for the evaluation process. The PPMI dataset contains a mix of Parkinson's disease (PD) patients and control patients suffering from other diseases. We chose this dataset for two reasons: (1) it contains many feature types, and (2) Parkinson's disease is a complicated disorder the causes of which have been attributed to complex combinations of genetic and environmental factors. We only considered the 1,579 patients with baseline visits. This dataset includes 6,013 biospecimen, genetic, drug, symptom, diagnosis, medical test, and demographic features. Feature types include binary, numerical, and categorical features. We binarized the categorical features. On average, each patient only has 261 of the 6,013 available features, which supports our previous assertion that EMRs are typically sparse.

## 3. Patient Profile Matrix

In general, each EMR can be formally represented as a high-dimensional feature vector. Features can be any attributes of interest in an EMR. We first generated a $p \times n$ patient profile matrix, $M$, from the data, where $p$ is the number of patients and $n$ is the number of features. Thus, each row corresponds to a patient record and each column corresponds to a feature. In our dataset, $p = 1,579$ and $n = 6,013$. Existing visualization methods use $M$ directly as input (see **Related Work**). However, as previously stated, $M$ is typically sparse and thus suboptimal for visualization. The main idea of VisAGE is to enrich $M$ before visualization by leveraging associations inferred from a knowledge graph. The enriched profile matrix, $M'$, then replaces $M$ as the input to any visualization method. We later show that $M'$ gives better visualizations in several applications on our dataset.

## 4. Patient Profile Matrix Enrichment

Our proposed method for enriching a profile matrix consists of three steps (Figure 1). The first constructs a knowledge graph with external data sources and EMRs. The second performs embedding on the constructed graph to learn a similarity matrix. Lastly, the third step multiplies $M$ by the similarity matrix to obtain $M'$. We now describe each step in more detail.
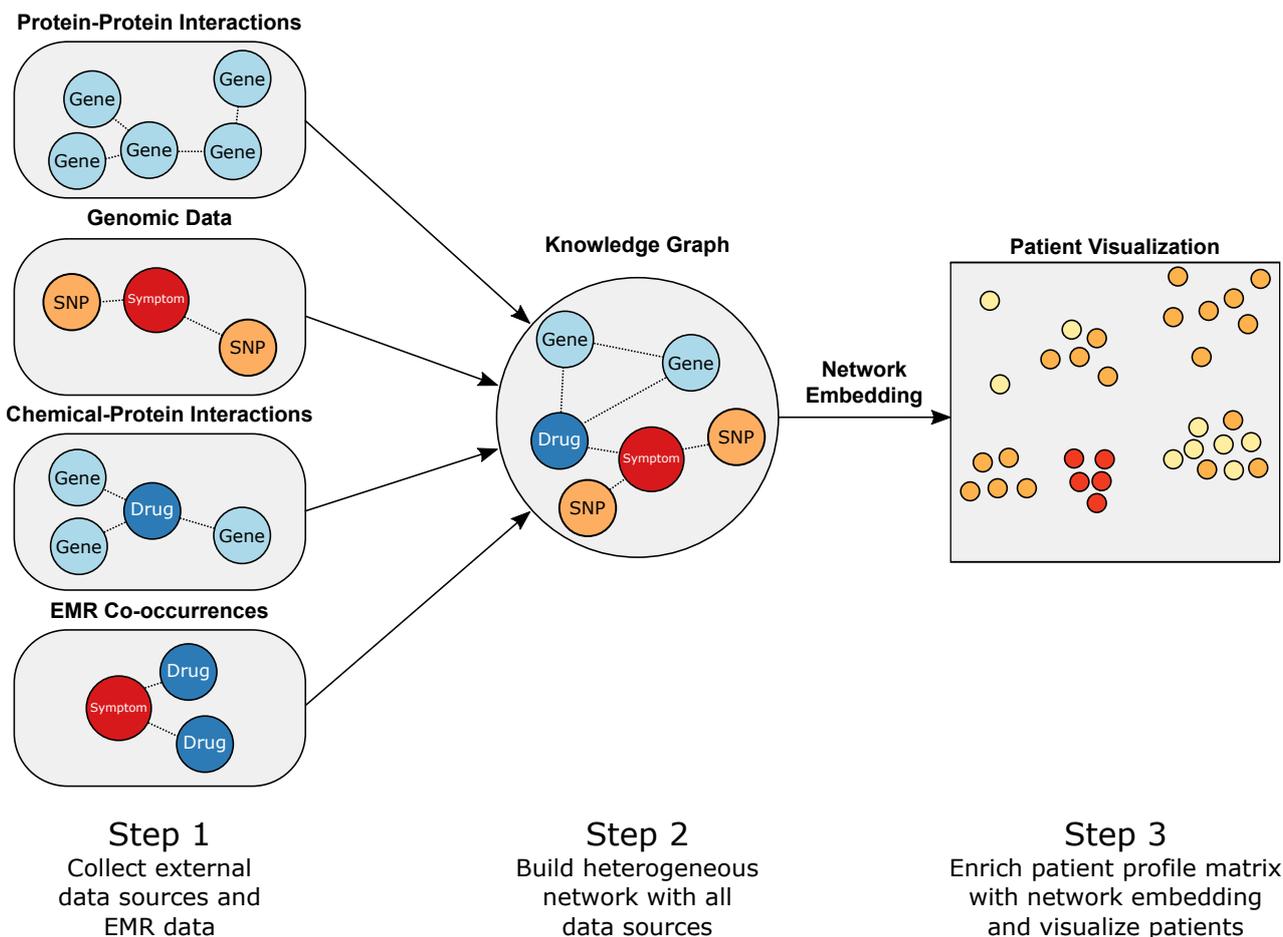


Fig. 1: The VisAGE pipeline. We first create a knowledge graph from multiple data sources. We then perform network embedding on this knowledge graph, enrich the patient profile matrix, and then visualize each patient in a two-dimensional space.

### 4.0.1. *Knowledge Graph Construction*

The knowledge graph is a heterogeneous network containing edges from four data sources.

(1) **Protein-protein interaction network.** We used the inBioMap database[13] of protein-protein interaction (PPI) edges. For a functional linkage between two proteins $p_1$ and $p_2$, we created a node for $p_1$, a node for $p_2$, and an undirected edge $\{p_1, p_2\}$ in the network. There were 17,327 proteins and 606,194 edges in this network.

(2) **Single-nucleotide polymorphism enrichment.** We integrated genomic data in the form of single-nucleotide polymorphisms (SNPs), which are single variations in the human genome. We identified SNPs that are highly enriched in PD patients in the dataset by using a one-sided Fisher's exact test.[14] Overall, we found 3,900 SNPs with $p$-values $< 0.05$. We then selected the nonsynonymous SNPs and determined if specific symptoms were enriched in SNPs with another one-sided Fisher's exact test. For each PD-enriched SNP $g$, we created a node for $g$ and an edge $\{g, s\}$ if $s$ was significantly enriched in $g$ with a $p$-value $< 0.01$. There were 34,324 SNP-symptom edges.

(3) **Chemical-protein interaction network.** We used STITCH, a database of known and predicted interactions between chemicals and proteins.[15] STITCH includes computationally predicted associations in addition to those aggregated from other databases. For each drug $d$ in the EMR data, if $d$'s active ingredient interacts with a protein $p$ in the STITCH database, then we created a node for $d$, a node for $p$, and an undirected edge $\{d, p\}$. There were 7,218 drug-protein edges in this network.

(4) **Electronic medical records.** We directly added co-occurrence edges from each medical record. For example, if a patient was diagnosed with symptom $s$ and prescribed a drug $d$, then we created a node $s$, a node $d$, and an undirected edge $\{s, d\}$. We repeated this for all elements in each patient's medical record.

The resulting network contained 23,886 nodes and 17,108,116 edges. The knowledge graph's purpose is to utilize the "guilt by association" rule: although many related medical concepts may not directly co-occur in any medical records, they may indirectly share neighbors in the knowledge graph through the protein-protein, SNP-symptom, and drug-protein edges.

## 4.1. *Similarity Matrix Learning*

We used the recently proposed method ProSNet to infer the relationships among the entities in the knowledge graph.[16] ProSNet performs a dimensionality reduction algorithm on heterogeneous networks to optimize a low-dimensional vector representation for each node. The vectors of two nodes will be co-localized in the low-dimensional space if the nodes are near each other in the heterogeneous network. After generating these vectors, we enriched the patient profile matrix. The similarities between the node vectors capture the latent relationships among the external data sources and the information in the EMRs. VisAGE constructs an $n \times n$ similarity matrix, $S$, where $n$ is the number of features and $S_{ij}$ is the cosine similarity between feature $i$'s low-dimensional vector and feature $j$'s low-dimensional vector.

### 4.1.1. *Enriching the Profile Matrix*

After obtaining $S$, we generated the enriched patient profile, $M'$, with the following operation:

$$M' = M \times S \tag{1}$$

This multiplication allows features that are related but not semantically identical to partially match through similarity scores. As stated before, each patient in the original patient profile matrix, $M$, only had an average of 261 nonzero features. On the other hand, each patient in $M'$ had 1,732 nonzero features.

## 5. Evaluation

We wished to determine whether $M'$ would lead to better visualization results than the original patient profile matrix, $M$. Thus, we compared the results between using $M'$ versus using $M$ in meaningful downstream visualization applications. Specifically, following previous studies,[17,18] we used t-distributed stochastic neighbor embedding (t-SNE),[19] an algorithm that can efficiently model high-dimensional objects as two-dimensional points, which makes it especially well-suited for visualizing our dataset. We generated our visualization by running t-SNE with default settings on the patient profile matrix $M$ for the baseline and $M'$ for VisAGE. This created a new $p \times 2$ matrix, so that each patient was finally reduced to two dimensions. We plotted this matrix as a set of points. We now discuss our results when visualizing $M'$ and $M$ in various applications.

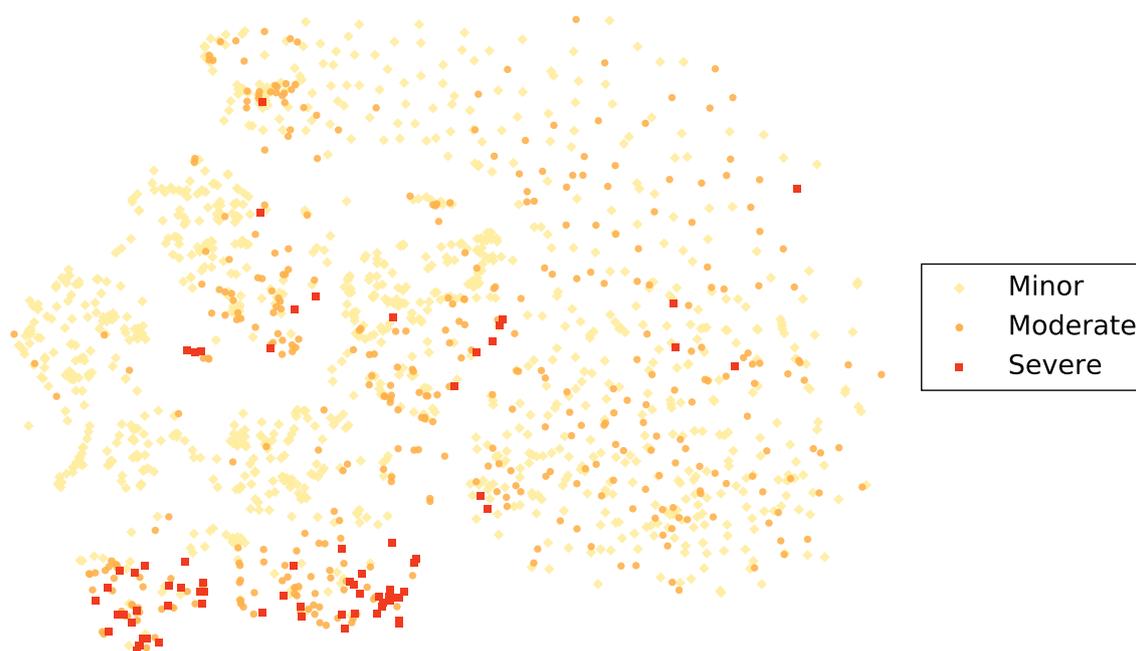### 5.1. *Two-Dimensional Visualization with UPDRS*

Using the two-dimensional representations of patient records, we labeled each record according to its unified Parkinson's disease rating scale (UPDRS) scores. The UPDRS consists of six sections, each containing survey questions that evaluate a patient's physical and mental condition.[20] The questions deal with topics ranging from anxiety to sleeping problems, with scores scaled from 0 to 4. A higher score indicates more severe impairment or disability. As in previous work, we labeled each patient with the sum of his or her UPDRS scores.[21]

The main difference between VisAGE and the baseline is that in Figure 2, moderately impaired patients (orange circles) on the right side of the plots were clustered more distinctly in the VisAGE visualization, while the same patients were less structured in the baseline visualization. The most severe Parkinson's disease patients are marked by red squares, and were clustered more tightly together in the VisAGE visualization than in the baseline. The baseline's worse performance can be attributed to the data sparsity of the EMRs.

### 5.2. *Qualitative Evaluation: Drug and Symptom Enrichment*

We qualitatively evaluated the visualization results by computing drug and symptom enrichments for each cluster. We used symptoms and drugs because they are strongly connected to patient statuses and diagnoses. Thus, if a cluster is highly enriched in a symptom or drug, then doctors will have a general idea of the cluster's disease subtype. We first clustered the two-dimensional patient representations with DBSCAN,[22] which is robust to outliers and does not need to specify the number of clusters. Because the PPMI dataset contains control patients to simulate noise, DBSCAN's robustness to outliers is especially desirable. Additionally, not having to specify the number of clusters *a priori* is useful for our application, as we do not know the exact number of patient subtypes beforehand.

For the DBSCAN parameters, we set $\epsilon = 1$ and $minPts = 10$. In the baseline method, patients were placed into 10 clusters. With VisAGE, patients were placed into 18 clusters. For each cluster $c$ and each symptom or drug $b$, we computed Fisher's exact test to determine if $c$ was significantly enriched in $b$. We only used symptoms and drugs with binary values to avoid medical tests for which all patients had non-zero values (e.g., the Epworth Sleepiness Scale[23]).

(a) Baseline visualization



(b) VisAGE visualization

Fig. 2: The two-dimensional representations of patient records, plotted with color labels determined by each record's UPDRS scores. VisAGE's visualization identifies more clusters for moderately impaired patients, and more tightly groups severely impaired patients.

### 5.3. *VisAGE Discovers More Patient Subtypes*

We show the two-dimensional plot of patients in Figures 3 and 4, with each color-shape combination corresponding to a unique cluster generated by DBSCAN. We also show the two most enriched symptoms for each cluster in the legends. With the baseline, many of the points in the upper right quadrant of the plot were determined to be noise (black circles). As a result, no patients can be deemed to be similar to these noise points. On the other hand, VisAGE was able to properly classify many of these patients into distinct clusters.

We saw that both methods mostly grouped together the patients with the highest UPDRS scores (Figure 2). The corresponding DBSCAN clusters that overlapped with these high-UPDRS patients were most enriched in parkinsonism and Parkinson's disease, as expected.



Menopausal depression
Hypothyroidism

REM sleep disorder
Enlarged prostate

Parkinsonism
Hypertension

Corrective lens user
Hypertension

Pregnancy
Polycystic ovaries

Hormone replacement therapy
Caesarean section

Parkinsonism
Parkinson's disease

Hypertension
Hypercholesterolaemia

Refraction disorder
House dust allergy
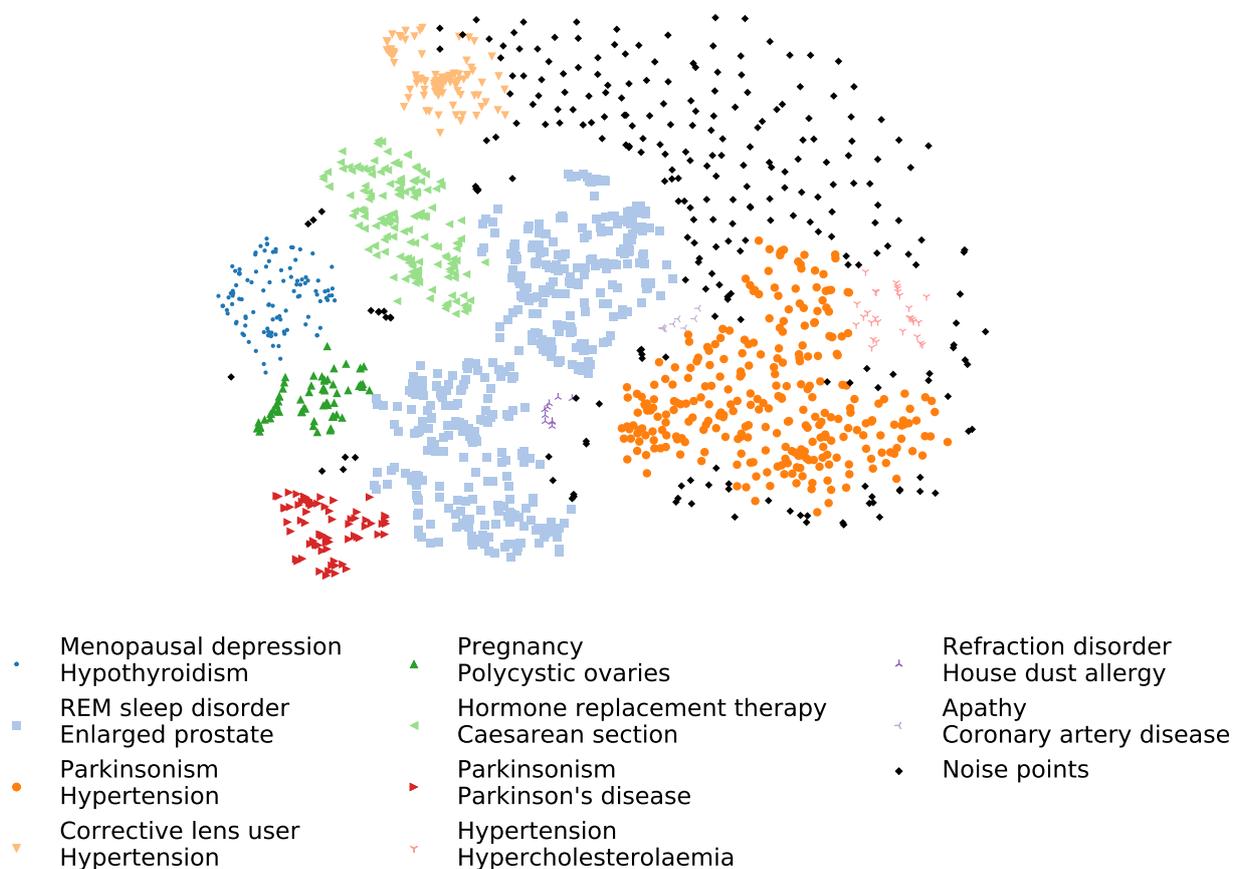
Apathy
Coronary artery disease

Noise points

Fig. 3: The baseline's two-dimensional representation of patient records, with colors determined by the DBSCAN clustering.

Both methods identified a cluster of patients enriched in parkinsonism and hypertension (orange circles in the baseline and dark green triangles pointing up in VisAGE). Indeed, hypertension is commonly known to be prevalent in PD patients.[24] However, VisAGE identified four additional clusters that were significantly enriched in PD/parkinsonism and another informative symptom. On the other hand, the baseline method combined these clusters into

Fig. 4: VisAGE's two-dimensional representations of patient records, with colors determined by the DBSCAN clustering.

larger ones, losing information in the process. We interpreted these additional clusters as PD patient subtypes that required special treatment. We now discuss these four clusters.

(1) **Parkinsonism and head injury.** The cluster of dark orange circles contained 16 patients, and was enriched in parkinsonism and head injury with $p$-values of $3.110 \times 10^{-4}$ and 0.01013, respectively. This is consistent with previous work, as head trauma is one of the most common candidates for PD causes.[25] This cluster was highly enriched in entacapone, levodopa, and carbidopa with $p$-values of $1.085 \times 10^{-25}$, $9.030 \times 10^{-10}$, and $7.099 \times 10^{-5}$, respectively. While levodopa/carbidopa (LC) is the most common drug prescribed to PD patients, entacapone is often prescribed as a supplementary drug to improve the efficacy

of LC.[26] As expected, these patients are also labeled as "Severe" in Figure 2, which would explain the need for this supplement. Furthermore, entacapone has been proposed as a possible treatment for traumatic brain injury.[27] In the baseline, this group of patients was incorrectly combined with the cluster most enriched in parkinsonism and hypertension.

(2) **REM sleep disorders and Parkinson's disease**. The cluster of light orange, down-pointing triangles contained 292 patients, and was enriched in rapid eye movement (REM) sleep behavior disorder, which is most often associated with PD ($p$-values of $1.477 \times 10^{-5}$ and $7.246 \times 10^{-3}$, respectively).[28] In addition to the standard levodopa prescription ($p$-value $= 9.787 \times 10^{-23}$), the cluster was also highly enriched in clonazepam ($p$-value $= 0.004377$). Clonazepam administered with levodopa at bedtime has been shown to reduce REM sleep disorder symptoms.[29] In the baseline, the corresponding cluster contained nearly twice as many patients (458), and was not highly enriched in Parkinson's disease.

(3) **Parkinsonism and bradykinesia.** In VisAGE's visualization, the cluster of light green, left-pointing triangles contained 159 patients, and was enriched in parkinsonism and bradykinesia with $p$-values of $1.526 \times 10^{-8}$ and $3.974 \times 10^{-8}$, respectively. As expected, bradykinesia is a key symptom of parkinsonism.[30] Additionally, this cluster was highly enriched in ropinirole with a $p$-value of $1.246 \times 10^{-7}$. Ropinirole stimulates mesolimbic $D_3$ receptors, which alleviates bradykinesia.[31] In the baseline, this group of patients was mixed with patients exhibiting parkinsonism and hypertension.

(4) **Parkinsonism and back injury**. The cluster of light brown circles contained 17 patients, and was enriched in parkinsonism and back injury with $p$-values of $1.320 \times 10^{-5}$ and $1.092 \times 10^{-4}$, respectively. A previous study showed that spinal cord injuries are associated with increased risk of PD.[32] In addition to the standard levodopa/carbidopa prescription, this cluster was significantly enriched in amantadine ($p$-value $= 6.09 \times 10^{-5}$). Amantadine is not only an antiparkinsonian agent, but has also been shown to act as a non-competitive $N$-Methyl-D-aspartate (NMDA) receptor antagonist.[33] NMDA receptor antagonists have been shown to treat acute spinal cord injuries.[34] Like in the cluster that was enriched in parkinsonism and head injury, this cluster also contained many patients with severe UPDRS scores. In the baseline, these patients were again mixed with the cluster enriched in parkinsonism and hypertension.

### 5.4. *Quantitative Evaluation: False Discovery Rate*

For each method, we compared the number of clusters highly enriched in drugs and symptoms. To this end, we excluded drugs and symptoms from both $M$ and $M'$. Additionally, we excluded these features from VisAGE's knowledge graph in order to limit data leakage. We then re-computed the enrichments for drugs and symptoms, taking the drug or symptom with the lowest $p$-value to represent each cluster. With these $p$-values, we counted the number of clusters that were significantly enriched in at least one drug or symptom.

To create a fair comparison, we used the Benjamini-Hochberg procedure[35] to control the false discovery rate at different levels of $\alpha$ (Figure 5). We see that VisAGE identified more enriched clusters than the baseline at every level of $\alpha$, which is consistent with our earlier observation that the baseline method is incapable of distinguishing among patients with less

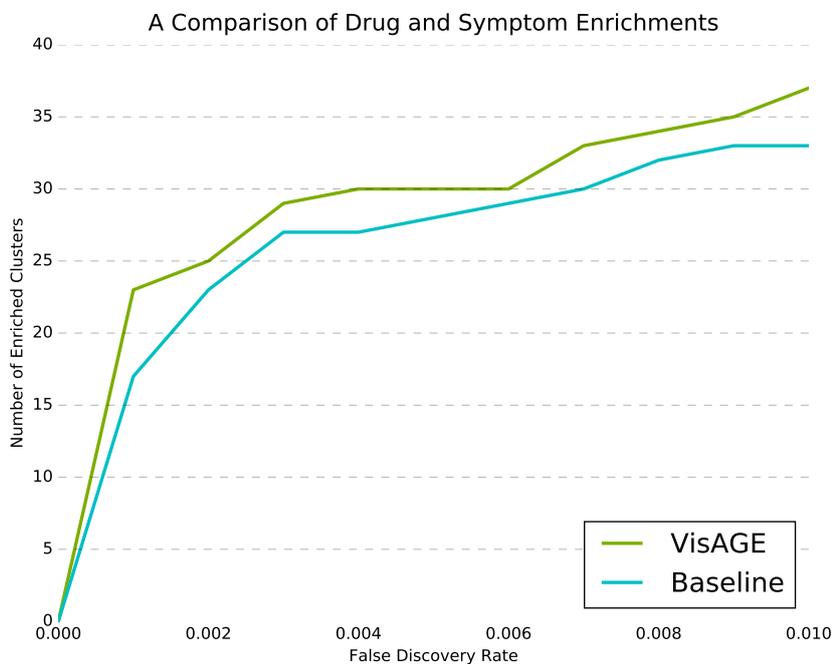severe symptoms. Thus, we conclude that VisAGE also performs better quantitatively.



Fig. 5: A plot comparing the baseline and VisAGE. VisAGE dominates the baseline in the number of clusters enriched with at least one drug or symptom at each level of $\alpha$.

## 6. Related Work

A previous study visualized high-dimensional data with a technique called LargeVis. However, it builds a k-NN network directly from the data, and then reduces the network to two dimensions without using external information.[18] Another study built upon LargeVis to visualize single cells, but still also directly computed embeddings from a k-NN network without utilizing external data.[17] Marlin *et al.* visualized a pattern discovery model's clustering parameters in the context of EMR analysis.[36] However, they focused on longitudinal data and predicting mortality outcomes rather than patient subtyping. Gotz *et al.* performed interactive visualization of EMR data, but worked with time series data to analyze patterns over time.[37] The Dynamic Icons (DICON) system clusters EMRs that are similar to a given patient, visualizing the clusters. However, it does not utilize molecular interaction networks or genomic data to compute similarities between EMRs.[38] Lastly, Perer *et al.* developed Care Pathway Explorer to visualize EMR data to investigate correlations with patient outcomes.[39] However, they use sequential pattern mining, which relies on historical EMR data to extract patterns.

## 7. Conclusions and Future Work

In this paper, we presented VisAGE, a method of improving EMR visualization by enriching EMRs with external knowledge sources. Evaluations on a PD patient dataset showed that VisAGE can generate visualizations such that similar patients are clustered together more tightly than in a baseline that does not alter the original database. We also evaluated our visualization with enrichments of drugs and symptoms, and showed that VisAGE can produce a higher quantity of fine-grained partitions of PD patients.

One limitation of our work is that the evaluation is done on only one dataset, which is mainly due to the necessity of expensive patient annotations. In the future, it is important to further evaluate the proposed enrichment method on more datasets as they become available. We also plan to build software that can implement our visualization in real application environments. Since VisAGE is a general method, the software would serve as a framework for an interactive component that can enrich any EMR database. For example, in a clinical setting, previously treated patients can serve as guidelines for doctors treating new patients. Doctors can identify these similar, previously treated patients in the two-dimensional space using the visualization tool and optimize treatment for the current patient.

## 8. Acknowledgments

## References

1. S. T. Mennemeyer, N. Menachemi, S. Rahurkar and E. W. Ford, *Journal of the American Medical Informatics Association* **23**, 375 (2016).
2. J. Bae and W. E. Encinosa, *BMC health services research* **16**, p. 172 (2016).
3. J. Bae, J. M. Hockenberry, K. J. Rask and E. R. Becker, *Health care management review* **42**, 258 (2017).
4. A. Rind, P. Federico, T. Gschwandtner, W. Aigner, J. Doppler and M. Wagner, Visual analytics of electronic health records with a focus on time, in *New Perspectives in Medical Records*, (Springer, 2017) pp. 65–77.
5. A. Rind, T. D. Wang, W. Aigner, S. Miksch, K. Wongsuphasawat, C. Plaisant, B. Shneiderman *et al.*, *Foundations and Trends® in Human–Computer Interaction* **5**, 207 (2013).
6. M. Ozkaynak, B. Reeder, L. Hoffecker, M. B. Makic and K. Sousa, *CIN: Computers, Informatics, Nursing* **35**, 465 (2017).
7. S. Ebadollahi, J. Sun, D. Gotz, J. Hu, D. Sow and C. Neti, Predicting patients trajectory of physiological data using temporal trends in similar patients: A system for near-term prognostics, in *AMIA annual symposium proceedings*, 2010.
8. G. Hals and F. Lovecchio, *Emergency Medicine Reports* **30**, 145 (2009).
9. M. S. Hansen, G. J. Nogareda and S. J. Hutchison, *The American journal of cardiology* **99**, 852 (2007).
10. P. G. Hagan, C. A. Nienaber, E. M. Isselbacher, D. Bruckman, D. J. Karavite, P. L. Russman, A. Evangelista, R. Fattori, T. Suzuki, J. K. Oh *et al.*, *Jama* **283**, 897 (2000).

11. B. K. Nallamothu, R. H. Mehta, S. Saint, A. Llovet, E. Bossone, J. V. Cooper, U. Sechtem, E. M. Isselbacher, C. A. Nienaber, K. A. Eagle *et al.*, *The American journal of medicine* **113**, 468 (2002).

12. K. Marek, D. Jennings, S. Lasch, A. Siderowf, C. Tanner, T. Simuni, C. Coffey, K. Kieburtz, E. Flagg, S. Chowdhury *et al.*, *Progress in neurobiology* **95**, 629 (2011).

13. T. Li, R. Wernersson, R. B. Hansen, H. Horn, J. M. Mercer, G. Slodkowicz, C. Workman, O. Regina, K. Rapacki, H.-H. Staerfeldt *et al.*, *bioRxiv* , p. 064535 (2016).

14. R. A. Fisher, *Journal of the Royal Statistical Society* **85**, 87 (1922).

15. D. Szklarczyk, A. Santos, C. von Mering, L. J. Jensen, P. Bork and M. Kuhn, *Nucleic acids research* **44**, D380 (2015).

16. S. Wang, M. Qu and J. Peng, Prosnet: Integrating homology with molecular networks for protein function prediction, in *Pacific Symposium on Biocomputing. Pacific Symposium on Biocomputing*, 2016.

17. J. Kim, N. Russell and J. Peng, Scalable visualization for high-dimensional single-cell data, in *PACIFIC SYMPOSIUM ON BIOCOMPUTING 2017*, 2017.

18. J. Tang, J. Liu, M. Zhang and Q. Mei, Visualizing large-scale and high-dimensional data, in *Proceedings of the 25th International Conference on World Wide Web*, 2016.

19. L. v. d. Maaten and G. Hinton, *Journal of Machine Learning Research* **9**, 2579 (2008).

20. C. Ramaker, J. Marinus, A. M. Stiggelbout and B. J. Van Hilten, *Movement Disorders* **17**, 867 (2002).

21. C. Shi, Z. Zheng, Q. Wang, C. Wang, D. Zhang, M. Zhang, P. Chan and X. Wang, *PloS one* **11**, p. e0155758 (2016).

22. M. Ester, H.-P. Kriegel, J. Sander, X. Xu *et al.*, A density-based algorithm for discovering clusters in large spatial databases with noise., in *Kdd*, (34)1996.

23. M. W. Johns, *sleep* **14**, 540 (1991).

24. A. A. Ejaz, I. S. Sekhon and S. Munjal, *European journal of internal medicine* **17**, 417 (2006).

25. S. M. Goldman, C. M. Tanner, D. Oakes, G. S. Bhudhikanok, A. Gupta and J. W. Langston, *Annals of neurology* **60**, 65 (2006).

26. R. A. Hauser, M. Panisset, G. Abbruzzese, L. Mancione, N. Dronamraju and A. Kakarieka, *Movement Disorders* **24**, 541 (2009).

27. R. D. Zafonte, J. Lexell and N. Cullen, *The Journal of head trauma rehabilitation* **16**, 112 (2001).

28. J. J. Gugger and M. L. Wagner, *Annals of Pharmacotherapy* **41**, 1833 (2007).

29. M. Stacy, *Drugs & aging* **19**, 733 (2002).

30. M. Hallett and S. Khoshbin, *Brain* **103**, 301 (1980).

31. W. H. Jost and D. Angersbach, *CNS drug reviews* **11**, 253 (2005).

32. T. Yeh, Y. Huang, H. Wang and S. Pan, *Spinal cord* **54**, 1215 (2016).

33. J. Kornhuber, G. Quack, W. Danysz, K. Jellinger, W. Danielczyk, W. Gsell and P. Riederer, *Neuropharmacology* **34**, 713 (1995).

34. A. I. Faden, J. Ellison and L. Noble, *European journal of pharmacology* **175**, 165 (1990).

35. Y. Benjamini and Y. Hochberg, *Journal of the royal statistical society. Series B (Methodological)* , 289 (1995).

36. B. M. Marlin, D. C. Kale, R. G. Khemani and R. C. Wetzel, Unsupervised pattern discovery in electronic health care data using probabilistic clustering models, in *Proceedings of the 2nd ACM SIGHIT International Health Informatics Symposium*, 2012.

37. D. Gotz, J. Sun, N. Cao and S. Ebadollahi, Visual cluster analysis in support of clinical decision intelligence, in *AMIA Annual Symposium Proceedings*, 2011.

38. N. Cao, D. Gotz, J. Sun and H. Qu, *IEEE transactions on visualization and computer graphics* **17**, 2581 (2011).

39. A. Perer, F. Wang and J. Hu, *Journal of biomedical informatics* **56**, 369 (2015).