# ICD-10 Code Representation in Embedding Vector Spaces

By InuitionLabs.ai • 6/7/2025 • 60 min read

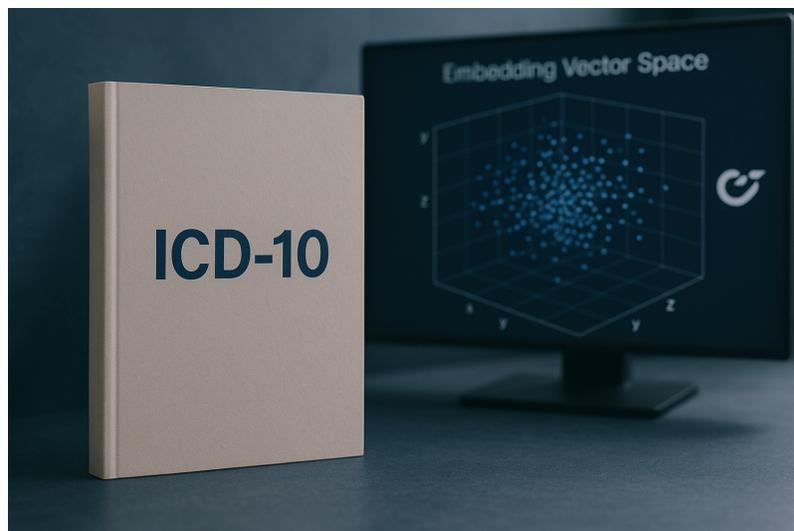icd-10    embedding vector spaces    healthcare data    machine learning    ehr    medical coding

data science    clinical informatics

# Encoding ICD-10 Codes in Embedding Vector Spaces

## Introduction

In modern healthcare, vast amounts of patient data are recorded using standardized codes. **ICD-10 codes** (International Classification of Diseases, 10th Revision) are a global standard for encoding diagnoses, conditions, and external causes of illness practicefusion.com bmcbioinformatics.biomedcentral.com. These codes capture critical clinical information (e.g. diseases, symptoms, findings) in electronic health records (EHRs) and are essential for consistency in documentation, billing, and epidemiological research. There are tens of thousands of ICD-10 codes (the U.S. ICD-10-CM modification includes over 70,000 codes) covering a comprehensive spectrum of health conditions practicefusion.com bmcbioinformatics.biomedcentral.com. Each patient's medical history can be represented as a sequence or set of ICD-10 codes, which in aggregate characterize their health profile ai.jmir.org.

Machine learning and data science techniques require numerical input, so these categorical medical codes must be converted into a numeric form to be used in algorithms. **Embedding vector spaces** offer a powerful approach: an *embedding* is a representation of a discrete item (like a word or a code) as a point in a continuous vector space geeksforgeeks.org. In other words, each code is mapped to a vector of numbers. The key idea is that the geometry of the vector space can encode meaningful relationships – codes with similar meaning or context end up with vectors that are close together. This concept has revolutionized natural language processing, where *word embeddings* capture semantic similarity (e.g. vectors for "dog" and "cat" are near each other). Likewise, for clinical data, *ICD-10 code embeddings* can capture medical similarities (e.g. different codes for closely related forms of diabetes might have similar vectors).

This report provides a comprehensive overview of encoding ICD-10 codes into embedding vector spaces. We will discuss what ICD-10 codes are and why they matter, explain embedding vector spaces in an intuitive way, and explore the benefits of embedding ICD codes versus naive encoding. We then delve into various methods for creating embeddings of categorical codes – from simple one-hot encodings to more advanced approaches like Word2Vec, FastText, GloVe, transformer-based models, and graph-based embeddings. Real-world use cases from research and industry are highlighted, such as using embeddings for predictive modeling, patient similarity clustering, clinical decision support, and phenotyping of diseases. We also address the challenges, limitations, and ethical considerations when embedding medical codes, and finally survey potential future research directions in this rapidly evolving area. Throughout, citations and examples are provided to ground the discussion in current knowledge and applications.

## What are ICD-10 Codes and Why Are They Important?

**ICD-10** stands for the 10th revision of the *International Statistical Classification of Diseases and Related Health Problems*. It is a standardized medical classification system developed by the World Health Organization (WHO) and used worldwide practicefusion.com bmcbioinformatics.biomedcentral.com. In ICD-10, diseases, injuries, and health conditions are assigned alphanumeric codes. For example, *E10.52* might denote "Type 1 diabetes mellitus with diabetic peripheral angiopathy with gangrene," whereas *E13.52* denotes "Other specified diabetes mellitus with diabetic peripheral angiopathy with gangrene" – two distinct codes that describe very similar conditions web.stanford.edu. ICD-10 codes are organized hierarchically by chapters and categories; the code structure itself often encodes information (the first character indicates a broad category like cardiology, endocrinology, etc., and additional digits add specificity).

**Relevance in healthcare:** ICD-10 codes serve numerous critical functions in health systems:

- **Consistent documentation:** They provide a *consistent and standardized method* for coding diagnoses across different providers and hospitals bmcbioinformatics.biomedcentral.com. This ensures that when a condition is recorded (e.g. "diabetes mellitus"), it's coded the same way everywhere, enabling uniform data collection and exchange.

- **Billing and reimbursement:** In many countries, including the U.S., ICD codes are used for billing insurance and payers. Healthcare providers must assign ICD-10 codes to every patient encounter to justify treatments and obtain reimbursement practicefusion.com. Thus, these codes have financial and administrative importance.

- **Clinical research and epidemiology:** Because ICD-10 is standardized, it allows aggregation of data for research. Public health agencies and researchers use ICD-coded data to track disease prevalence, identify trends, and conduct epidemiological surveillance bmcbioinformatics.biomedcentral.com. For example, one can query databases for all patients with ICD code E66 (obesity) to study obesity rates or outcomes.

- **Quality measurement:** Hospitals and regulators use ICD codes to monitor outcomes and quality of care. Certain codes (like those indicating hospital-acquired infections or complications) can flag quality issues, and overall coding patterns can be used to compute risk scores (e.g. the Elixhauser comorbidity index uses ICD codes to quantify patient risk).

- **Health policy and planning:** Aggregated ICD-10 data informs policy by showing disease burden. Policymakers can allocate resources (vaccines, screening programs) by analyzing which ICD-coded conditions are rising in frequency bmcbioinformatics.biomedcentral.com bmcbioinformatics.biomedcentral.com.

In summary, ICD-10 codes are the *language of diagnosis* in healthcare. They condense complex medical information into standard labels. This standardization is crucial for interoperable health records, large-scale analytics, and day-to-day operations in healthcare. However, the richness and complexity of this coding system (with **over 70,000 codes in the ICD-10-CM U.S.**

**extension** and ~14,000 base codes in WHO's ICD-10 practicefusion.com bmcbioinformatics.biomedcentral.com) pose challenges when we try to use them in computational models – motivating the need for better representations like embeddings.

## Embedding Vector Spaces in Machine Learning

**Embedding vector spaces** are a core concept in machine learning for representing categorical or discrete items in a numerical form that captures their meaning or context. An *embedding* is essentially a mapping of a discrete variable to a point in a continuous, multi-dimensional space. Instead of representing a concept like an ICD code as a single arbitrary index or a very sparse vector, we represent it as a dense vector of real numbers. Importantly, this vector is chosen such that it encodes useful information about the concept's relationships to others.

In formal terms, *an embedding is a mathematical representation of discrete objects or values as dense vectors in a continuous space* geeksforgeeks.org. These vectors typically have far fewer dimensions (e.g. 50, 100, or 300 dimensions) compared to the number of possible categories. The embedding space is learned so that similarity in the space reflects semantic or functional similarity in the original domain. For words in natural language, embedding algorithms learn to place synonyms or related words near each other in the vector space.

For example, in afirst character indicates a broad category like cardiology, endocrinology, etc., and additional digits add specificity).

**Relevance in healthcare:** ICD-10 codes serve numerous critical functions in health systems:

- **Consistent documentation:** They provide a *consistent and standardized method* for coding diagnoses across different providers and hospitals bmcbioinformatics.biomedcentral.com. This ensures that when a condition is recorded (e.g. "diabetes mellitus"), it's coded the same way everywhere, enabling uniform data collection and exchange.

- **Billing and reimbursement:** In many countries, including the U.S., ICD codes are used for billing insurance and payers. Healthcare providers must assign ICD-10 codes to every patient encounter to justify treatments and obtain reimbursement practicefusion.com. Thus, these codes have financial and administrative importance.

- **Clinical research and epidemiology:** Because ICD-10 is standardized, it allows aggregation of data for research. Public health agencies and researchers use ICD-coded data to track disease prevalence, identify trends, and conduct epidemiological surveillance bmcbioinformatics.biomedcentral.com. For example, one can query databases for all patients with ICD code E66 (obesity) to study obesity rates or outcomes.

- **Quality measurement:** Hospitals and regulators use ICD codes to monitor outcomes and quality of care. Certain codes (like those indicating hospital-acquired infections or complications) can flag quality issues, and overall coding patterns can be used to compute risk scores (e.g. the Elixhauser comorbidity index uses ICD codes to quantify patient risk).

- **Health policy and planning:** Aggregated ICD-10 data informs policy by showing disease burden. Policymakers can allocate resources (vaccines, screening programs) by analyzing which ICD-coded conditions are rising in frequency bmcbioinformatics.biomedcentral.com bmcbioinformatics.biomedcentral.com.

In summary, ICD-10 codes are the *language of diagnosis* in healthcare. They condense complex medical information into standard labels. This standardization is crucial for interoperable health records, large-scale analytics, and day-to-day operations in healthcare. However, the richness and complexity of this coding system (with **over 70,000 codes in the ICD-10-CM U.S. extension** and ~14,000 base codes in WHO's ICD-10 practicefusion.com bmcbioinformatics.biomedcentral.com) pose challenges when we try to use them in computational models – motivating the need for better representations like embeddings.

## Embedding Vector Spaces in Machine Learning

**Embedding vector spaces** are a core concept in machine learning for representing categorical or discrete items in a numerical form that captures their meaning or context. An *embedding* is essentially a mapping of a discrete variable to a point in a continuous, multi-dimensional space. Instead of representing a concept like an ICD code as a single arbitrary index or a very sparse vector, we represent it as a dense vector of real numbers. Importantly, this vector is chosen such that it encodes useful information about the concept's relationships to others.

In formal terms, *an embedding is a mathematical representation of discrete objects or values as dense vectors in a continuous space* geeksforgeeks.org. These vectors typically have far fewer dimensions (e.g. 50, 100, or 300 dimensions) compared to the number of possible categories. The embedding space is learned so that similarity in the space reflects semantic or functional similarity in the original domain. For words in natural language, embedding algorithms learn to place synonyms or related words near each other in the vector space.

For example, in a word embedding space, *"lion"*, *"tiger"*, and *"cat"* might cluster together while being far from vectors for unrelated terms geeksforgeeks.org geeksforgeeks.org. The distances and directions in the space have meaning: the famous example from word embeddings is that **Vector("King") - Vector("Man") + Vector("Woman")** is closest to **Vector("Queen")**, illustrating how relationships (male→female, king→queen) are captured by linear operations in the vector space.

Translating this to ICD-10 codes: an embedding space for medical codes would position *E10.52 (Type 1 diabetic angiopathy with gangrene)* near *E13.52 (Other diabetic angiopathy with gangrene)*, because these diagnoses are clinically similar web.stanford.edu. Similarly, codes for related conditions (like different types of diabetes, or various hypertensive diseases) would ideally form clusters or continuous gradations in the vector space. This property can be extremely useful — it means a machine learning model can *infer* that two codes are related without being explicitly told, simply by virtue of their vector representations.

To get these embeddings, we typically use data-driven algorithms that analyze the contexts in which codes appear (in analogy to words in sentences). The resulting vectors are *dense* (most entries are non-zero, distributing information across dimensions) and *continuous*. They allow mathematical operations: one can compute similarity (e.g. cosine similarity) between code vectors to quantify how related two diagnoses are, or use the vectors as features in predictive models.

In summary, embedding vector spaces enable us to go from the realm of discrete symbols (ICD codes) to a geometric, algebraically manipulable representation. This is a foundational step in applying modern machine learning and deep learning to health records, as algorithms can then take these vectors as input and detect patterns that would be invisible in the original discrete form.

## Why Encode ICD-10 Codes as Embeddings? (Benefits and Rationale)

Representing ICD-10 codes as learned vectors offers significant advantages over naive encoding schemes. The traditional approach to use diagnosis codes in models is **one-hot encoding** (or similar binary encoding), which has severe limitations. By understanding these limitations, we can appreciate why learned embeddings are beneficial:

- **High Dimensionality and Sparsity:** There are thousands of possible ICD-10 codes, which would translate to a very high-dimensional one-hot vector (mostly zeros, with a 1 for each code the patient has). For instance, using all ICD-10-CM codes would require a feature vector of length ~72,000 for each patient bmcbioinformatics.biomedcentral.com. This is computationally unwieldy and statistically challenging – models on such data risk overfitting and often require enormous sample sizes to estimate so many parameters. In practice, researchers often must *limit the codes* to a smaller subset (e.g. the top 100 or 1000 codes) ai.jmir.org ai.jmir.org or aggregate codes into broader categories, which sacrifices information. Embeddings, by contrast, condense this information into a dense vector of manageable size (say 100 dimensions), massively reducing dimensionality while preserving most of the important signal ai.jmir.org.

- **Ignoring Code Similarities:** One-hot encoding treats each code as an independent, unrelated feature. This is *not* true in reality – many codes are related or even nearly synonymous. For example, if a patient has *"unspecified atrial fibrillation"* vs *"atrial fibrillation with rapid ventricular response"*, those are two different ICD-10 codes, but clinically they describe the same fundamental condition (one with a qualifier). One-hot encoding would represent them as orthogonal vectors (having zero overlap), failing to recognize the similarity. In practice, *different ICD codes may represent similar or even identical underlying issues*, and one-hot "considers medical diagnoses as distinctive and unrelated features" ai.jmir.org. This loss of information limits predictive power. Learned embeddings can overcome this by placing similar codes closer together in the vector space, effectively *encoding clinical relationships*. Indeed, an ideal embedding might place E10.52 and E13.52 extremely near each other, reflecting that both describe diabetic angiopathy with gangrene web.stanford.edu. In a one-hot scheme, these overlap 0%; in an embedding, they might have a cosine similarity close to 1, indicating redundancy or closeness.

- **Hierarchical Structure and Context:** ICD-10 codes have a built-in hierarchy (chapters, blocks, categories) indicating relationships (for example, all codes starting with "I50" relate to heart failure). Standard encoding does not exploit this ontology – in fact, it breaks it by flattening codes into unrelated dummy variables. This means models might need to relearn relationships that are already known medically (like that I50.1 and I50.2 are both heart failure codes). Embedding methods can incorporate hierarchy explicitly or implicitly. By analyzing data, embeddings often *learn* that codes sharing prefixes or co-occurring frequently are related. Some advanced embedding methods even incorporate the ICD hierarchy during training to enforce that parent-child codes end up close (we'll discuss those later). Overall, encoding strategies that account for the *taxonomy of codes have been shown to better capture the medical code structure* than those that do not researchgate.net.

- **Improved Model Performance:** Using learned embeddings has been empirically shown to improve performance in healthcare machine learning tasks compared to one-hot encoding. By compressing information, embeddings can make models more generalizable and robust. For example, a recent study on **patient risk prediction** found that a model using patient vectors (embeddings of their ICD diagnosis history) outperformed a baseline model using binary indicators of diagnoses, especially in lower-dimensional representations ai.jmir.org. The embedding-based model achieved better predictive accuracy and was *more robust to missing data*, since even if some codes were missing, the vector still captured the patient's health profile in a smooth way ai.jmir.org. Other research has similarly noted that word embedding techniques applied to medical codes allow machine learning algorithms (like neural networks or gradient boosting machines) to reach their full potential with claims/EHR data ai.jmir.org – essentially because the input features are more informative and noise-tolerant than raw one-hot codes.

- **Compression of Complex Information:** An embedding can be seen as a form of *nonlinear compression*. It tries to pack the information content of a patient's myriad diagnoses into a small vector. The Pat2Vec framework, for instance, embedded entire patient diagnosis profiles into a 100-dimensional vector and demonstrated that this vector could still reproduce most of the relevant information from 10,000+ possible codes ai.jmir.org. In tests, the compressed representation retained enough detail to accurately predict outcomes like future healthcare costs, while being far more tractable to model with. The results showed "large performance gains, particularly in lower dimensions, demonstrating the embedding model's compression of nonlinear information" ai.jmir.org.

- **Generalization to Unseen Cases:** Another benefit is that embedding spaces can enable generalization to codes or combinations not explicitly seen in training. Because similar codes have similar vectors, if a model has learned to handle one code, it can naturally handle a related code it wasn't trained on, by virtue of the vectors being close. For example, if a model learned that code E10.9 (Type 1 diabetes, without complications) is an important risk factor for some outcome, it might automatically apply nearly the same logic to E11.9 (Type 2 diabetes without complications) if their embeddings are close – even if E11.9 was rare or absent in training data. This is especially useful given the **long-tailed distribution** of ICD codes where some codes are very rare. Traditional models struggle with rare categories, but an embedding can borrow strength from similar codes. In fact, research has found that embedding-based models like GRAM significantly improve prediction for *rare conditions* by leveraging ontology information – GRAM achieved 10% higher accuracy in predicting seldom-seen diseases compared to a baseline RNN, by generalizing from related codes in the hierarchy pmc.ncbi.nlm.nih.gov pmc.ncbi.nlm.nih.gov.

In short, encoding ICD-10 codes into a learned vector space allows us to handle the high dimensional, inter-related, and often noisy nature of diagnosis data in a much more effective way. It aligns the data representation with clinical reality (similar diagnoses represented similarly), which in turn lets machine learning models make better inferences. This leads to **better predictive performance, better clustering of patients, and potentially new insights** (since the embedding space might reveal relationships between diagnoses that weren't obvious from raw codes alone). The next sections will survey *how* we can generate such embeddings for ICD codes, and later we will see concrete examples of what we can do with them.

## Methods for Embedding ICD-10 Codes into Vector Spaces

There are several approaches to encode categorical medical codes like ICD-10 into vectors. These range from simple to sophisticated, each with its own strengths. We will overview the major methods, including examples of how they have been applied:

- **One-Hot Encoding (Baseline):** The simplest "embedding" is not really a learned embedding at all, but it's worth mentioning as a baseline. In one-hot encoding, each ICD-10 code is represented as a binary vector with a 1 in the position corresponding to that code and 0s elsewhere. For a patient with multiple codes, their record might be a multi-hot vector (1s for each code they have). While trivial to construct, this representation is extremely sparse and high-dimensional. It treats codes as unrelated features, failing to capture any similarities ai.jmir.org. One-hot encoding also doesn't scale well: using the full ICD-10 vocabulary leads to tens of thousands of features bmcbioinformatics.biomedcentral.com, so, as noted, practitioners often restrict to a subset (e.g. only coding the presence/absence of the 1000 most frequent codes) ai.jmir.org. This results in loss of information and can bias analyses toward common conditions. **In summary, one-hot is easy but suffers from the lack of generalization and relationship information** – hence the need for learned embeddings.

- **Distributional Embeddings (Word2Vec & Doc2Vec):** Borrowing techniques from natural language processing, we can treat each medical code like a "word" and each sequence of codes (e.g. codes in a patient visit, or all codes in a patient's history) like a "sentence" or document. **Word2Vec** is an algorithm that learns vector embeddings for words by scanning a corpus of sentences and learning which words co-occur with which others bmcbioinformatics.biomedcentral.com. Applying this to ICD codes, we feed in sequences of codes that occur together (for example, all diagnoses a patient had in a hospitalization, or the sequence of diagnoses across visits). The algorithm will then assign similar vectors to codes that tend to appear in similar contexts (meaning in similar patients or alongside similar other diagnoses). Researchers have indeed applied Word2Vec to medical coding data: for instance, the *Med2Vec* model trained embeddings on a large dataset of patients, yielding vectors for each diagnosis code bmcbioinformatics.biomedcentral.com. These embeddings captured latent groupings of diseases and were used to predict patient outcomes (Med2Vec improved mortality prediction by providing better features) bmcbioinformatics.biomedcentral.com. Another extension is **Doc2Vec**, which can learn a vector for an entire document (here, an entire patient) alongside word vectors – the Pat2Vec approach effectively uses Doc2Vec to get a patient-level embedding summarizing all their ICD-10 codes ai.jmir.org. In practice, Word2Vec-based methods (using either the continuous bag-of-words or skip-gram training strategies) have been popular for embedding medical concepts; they are *unsupervised* (learn from unlabeled co-occurrence data) and computationally efficient. Many studies have shown that basic Word2Vec on claims or EHR data yields clinically meaningful code embeddings bmcbioinformatics.biomedcentral.com bmcbioinformatics.biomedcentral.com. These vectors can then be averaged or pooled to represent patients, or used in sequence models (like feeding the code embeddings into an LSTM to predict future diagnoses web.stanford.edu web.stanford.edu). **In summary, Word2Vec learns embeddings by predicting codes from their neighbors (or vice versa) in patient records, producing vectors where co-occurring diagnoses lie close together.** Variants like *FastText* and *GloVe* (discussed next) build on similar distributional principles.

- **FastText Embeddings:** *FastText* is an extension of Word2Vec developed by Facebook AI that enriches the embedding process by considering subword information (usually character n-grams) for each token. In the context of natural language, FastText can compose the representation of a word from representations of its character pieces, which helps handle rare words or typos (e.g. it would understand "organiz(s)ation" variants). For ICD-10 codes, FastText could be applied by breaking codes into meaningful fragments – for example, the code "E11.9" might be segmented into "E11" and ".9" or even character trigrams like "E11", "11.", "1.9". The idea is that codes sharing substrings (which often means they share a category or subcategory) would have some shared representation. This can help with the *out-of-vocabulary* issue: if a certain ICD code was very rare in training, a pure Word2Vec might not learn a good vector for it, but FastText can still construct an embedding from its pieces (which it learned from other codes) medrxiv.org sciencedirect.com. For instance, suppose "A99" (unspecified viral hemorrhagic fever) is rare, but the model has seen many codes starting with "A9" (other infectious diseases). FastText will utilize the "A9" substring in embedding A99, giving it a sensible location near related infectious disease codes. Some researchers have indeed included FastText in evaluations of medical code embeddings; one study found that FastText embeddings achieved good performance in capturing ICD code similarities, especially for less frequent codes medrxiv.org. In summary, **FastText helps incorporate the internal structure of ICD codes into the embedding, leveraging the fact that code prefixes and patterns have meaning**. This typically yields more robust embeddings for rare or new codes.

- **GloVe (Global Vectors):** GloVe is another widely-used word embedding technique, which takes a more global statistical approach. Instead of scanning through contexts with a neural network like Word2Vec, GloVe builds a large matrix of co-occurrence counts (how often each code appears with each other code in patients or visits) and then factorizes this matrix to obtain embeddings sessionize.com. The result is similarly a vector for each code, but GloVe explicitly tries to make the dot product of two code vectors approximate the log of their co-occurrence probability. Using GloVe for ICD-10 would involve computing which diagnoses frequently appear together in the data (for example, diabetes appears often with hypertension; heart failure co-occurs with chronic kidney disease, etc.) and then deriving vectors that reflect these patterns. The advantage of GloVe is that it uses global corpus statistics, which can capture broad relationships (not just local context windows). For medical codes, this could reveal groupings like metabolic syndrome cluster (diabetes, hypertension, hyperlipidemia) through the embedding space because those diagnoses co-occur over many patients. Some works have explored GloVe for clinical concept embedding and found it effective in certain cases medrxiv.org – for instance, one might achieve slightly different embeddings than Word2Vec, but ones that can sometimes better reflect long-range associations in data. However, GloVe also shares some limitations with Word2Vec (it produces one static vector per code). **In practice, Word2Vec and GloVe (and FastText) all belong to the family of *distributional embeddings*, leveraging co-occurrence in data to learn representations.** They are unsupervised and relatively "shallow" two-layer models bmcbioinformatics.biomedcentral.com. Many pre-trained embeddings of medical codes available today were created with these methods (some even available open-source, e.g. pretrained vectors from claims data).

- **Transformer-Based and Contextual Embeddings:** Newer approaches use *transformer models* and large language model techniques to embed medical codes, either by treating sequences of codes like a language or by using the textual descriptions of codes. Transformers (like BERT, GPT, etc.) are deep models with self-attention mechanisms that can generate *contextualized* embeddings bmcbioinformatics.biomedcentral.com. There are two main ways they've been applied to ICD coding:
**(a) Sequence models on codes:** In a model like *BEHRT* (BERT for EHR), patients' records (sequence of codes, possibly with time steps) are input into a transformer encoder. The model produces embeddings for each code position (taking into account the surrounding codes) as well as an overall patient embedding. These code embeddings are *context-dependent* – e.g. the representation of code "I21" (acute myocardial infarction) might be slightly different if it appears in a patient who also has diabetes and hypertension versus in a patient who has none of those. Contextual models can capture nuanced differences (similar to how the word "bank" has different embeddings depending on whether the context is finance or river). In practice, researchers have shown that transformer models can be trained on large EHR datasets to learn powerful embeddings; one 2020 study introduced BEHRT and found it improved prediction of outcomes by capturing the temporal patterns of code sequences (it essentially treats the sequence of codes as a sentence and each visit as a segment) sciencedirect.com sciencedirect.com. Transformers can also naturally handle longer sequences and incorporate other information like patient age or time between visits as additional embeddings.
**(b) Embeddings from code descriptions using language models:** Another fruitful approach is to utilize the fact that each ICD-10 code comes with a textual description (the medical diagnosis name). We can feed these descriptions into a pretrained medical language model (like ClinicalBERT, BioBERT, or BioGPT) to obtain an embedding for the description, and use that as the code's vector. For example, the code "J45.909" has description "Unspecified asthma, uncomplicated". If we input that sentence into BioBERT, the model will output a sentence embedding capturing its meaning (likely positioned near other respiratory disease descriptions). A recent work by Kane et al. (2023) did exactly this: they generated *ICD-10-CM code embeddings by inputting the code descriptions into a large language model (BioGPT)* and then applied an autoencoder for dimensionality reduction bmcbioinformatics.biomedcentral.com. The result was a publicly available set of vector representations for all ICD-10-CM codes, which preserved clinical relationships – for instance, categories and subcategories of diseases were clustered together in the embedding space bmcbioinformatics.biomedcentral.com. They validated that these embeddings could be compressed down to as few as 10 dimensions while still maintaining the ability to reconstruct the original vectors (showing the redundancy in the original 768-dimensional BioGPT embeddings) bmcbioinformatics.biomedcentral.com. Such description-based embeddings leverage rich external knowledge (the language model's understanding of medical terminology) rather than raw co-occurrence in a specific dataset. This can be very useful for codes that are rare in data but whose description clearly places them in a known clinical context. It's a form of *transfer learning* from text to structured codes.
In general, transformer-based methods produce **contextual or highly informative embeddings**. They tend to be higher-dimensional and require more computation, but they address some limitations of static Word2Vec-style embeddings. Notably, they can produce different embeddings for a code depending on context (for example, the transformer might understand that "diabetes" as a *comorbidity* vs as a *primary diagnosis* could have different implications, adjusting the patient embedding accordingly). They also naturally handle sequential predictions (like predicting the next code or outcome from previous ones). As transformer models continue to advance, we expect to see even more sophisticated uses – e.g. GPT-like models generating patient embeddings from mixed

data (notes + codes). Already, the trend is towards using **Large Language Models (LLMs) in healthcare to get rich representations**, and ICD codes can be part of that input or output space bmcbioinformatics.biomedcentral.com bmcbioinformatics.biomedcentral.com.

- **Graph-Based and Ontology-Grounded Embeddings:** Another category of methods treats medical codes as nodes in a graph (network) and uses graph embedding techniques or knowledge of the ontology to derive vectors. The intuition here is that ICD-10 has a built-in graph structure (hierarchy), and also codes can be related through co-occurrence networks or through mapping to clinical ontologies like SNOMED CT or the UMLS metathesaurus. **Graph embeddings** like Node2Vec or GraphSAGE can learn node vectors such that connected nodes (or those with many shared neighbors) end up close in the vector space. If we create a graph where each ICD-10 code is connected to other codes that tend to co-occur in patients, then embedding this graph would yield vectors that group frequently co-occurring diagnoses. This is similar to Word2Vec's result but explicitly uses a graph formulation. More powerfully, we might incorporate the *ICD-10 hierarchy*: for example, connect each code to its parent category and to its sibling codes. By doing so, we inject domain knowledge that, say, J45 (asthma) is under category J40-J47 (chronic lower respiratory diseases). Embedding the ontology graph could ensure asthma's vector is near other chronic respiratory disease vectors even if data co-occurrence is sparse. Researchers Finch et al. (2021) tackled this by **explicitly incorporating hierarchical information into Word2Vec training** for ICD-10 codes researchgate.net. They added not only codes from patient data but also additional "context" tokens representing the higher-level categories (using the Clinical Classifications Software hierarchy) in the training corpus. The resulting embeddings ("hierarchy-informed Word2Vec") were better at reflecting the medical taxonomy: the learned vectors clustered according to known disease categories more often than standard Word2Vec researchgate.net researchgate.net. In clustering tests, including hierarchical data improved normalized mutual information with true groupings (61.4% vs 57.5%) researchgate.net. Moreover, using these embeddings in downstream prediction tasks modestly improved accuracy in most cases researchgate.net. This demonstrates that **blending data-driven methods with ontological structure can yield embeddings that are both data-efficient and clinically intuitive**. Another notable example is **GRAM (Graph-based Attention Model)** by Choi et al. (2017). GRAM isn't just an embedding method but a predictive model that *represents each medical concept (diagnosis) as a weighted combination of its ancestors in a knowledge graph (like the ICD ontology)* pmc.ncbi.nlm.nih.gov pmc.ncbi.nlm.nih.gov. The base embeddings are learned for the leaf codes, but when representing a patient or making a prediction, GRAM can generalize to higher-level concepts via attention. It achieved better performance on predicting conditions (especially those with limited data) by using the hierarchy to inform the representation – e.g., if a specific rare code wasn't seen enough, the model falls back on the parent category representation pmc.ncbi.nlm.nih.gov. GRAM's learned code representations were aligned with the medical ontology and helped interpretability (you could see which ancestor concepts were weighted) pmc.ncbi.nlm.nih.gov pmc.ncbi.nlm.nih.gov.
  Beyond purely hierarchical relations, **knowledge graph embeddings** can integrate multiple types of medical entities. For instance, *cui2vec* embedded UMLS concepts (which include diagnoses, medications, lab findings, etc.) into one space by mining a variety of clinical data bmcbioinformatics.biomedcentral.com. This means an ICD-10 code mapped to a UMLS Concept Unique Identifier (CUI) would get an embedding influenced by not just diagnosis co-occurrences but also by how that concept appears in clinical narratives and alongside medications. Such multi-modal embeddings can reveal connections like a diagnosis being closely related to certain treatments. Graph-based techniques can also utilize **graph neural networks (GNNs)**: e.g., a Graph Convolutional Network could be applied to the ICD hierarchy to compute embeddings that take into account neighbors and neighbors-of-neighbors. In fact, recent research on *ICD coding models* leveraged GCNs to embed label graphs, noting that incorporating the relations between ICD codes (as a label tree) improved text-to-code assignment nature.com.

In summary, graph-based embeddings use **structured relationships** (whether from the official ICD ontology or empirically derived graphs) to place codes in a vector space. They are particularly useful for ensuring that the embedding reflects known medical knowledge (e.g., grouping by disease category), and for improving the representation of rare codes by linking them with more common relatives. Many cutting-edge solutions combine graph approaches with language models (for example, use both a text encoder and a graph encoder and concatenate them). The recent **ICD2Vec (2023)** work exemplifies a hybrid approach: they fine-tuned a transformer model (GatorTron, a clinical language model) on a combination of structured and unstructured data about diseases, effectively marrying text-based embedding with knowledge of symptom-disease associations. The authors report that ICD2Vec embeddings captured meaningful semantic relationships – for example, in their embedding space, the diseases most similar to *COVID-19* were *common cold (J00)*, *unspecified viral hemorrhagic fever (A99)*, and *smallpox (B03)*, which arguably makes sense in terms of virology and symptom overlap github.com. They also introduced a way to derive a risk score (IRIS) from the embeddings, illustrating the clinical utility of such representations github.com github.com.

Each of these methods can be used in isolation or in combination. One might start with a Word2Vec embedding and then refine it using a GNN with hierarchy, or use a language model to initialize embeddings which are then fine-tuned on EHR data (a form of transfer learning). The field is moving toward **integrative approaches**, where we leverage all available signals – data co-occurrence, code descriptions, and ontology – to get the best embeddings. The good news is that many pre-trained embeddings and tools are available open-source. For example, the ICD2Vec project has released their code and embeddings on GitHub for others to use github.com, and the Pat2Vec study made their trained model publicly available for application on new data ai.jmir.org. Additionally, large healthcare datasets like *MIMIC-III* (an ICU database) have been extensively used to train concept embeddings and are open to researchers, providing a rich resource to experiment with encoding medical codes.

## Use Cases and Applications of ICD-10 Code Embeddings

Embedding ICD-10 codes into vector form is not just a theoretical exercise – it enables and enhances many practical applications in healthcare AI and data analysis. Below are several key use cases, along with examples from research and real-world implementations:

- **Predictive Modeling and Risk Stratification:** One of the primary uses of code embeddings is as features in predictive models. Instead of feeding thousands of sparse indicator variables for diagnoses into a model, clinicians and data scientists can feed a concise embedding (or a set of embeddings) that summarize a patient's condition. This has been shown to improve predictions of outcomes such as mortality, readmission, disease onset, and healthcare costs. For example, the *Pat2Vec* study created patient embeddings from ICD-10 codes and found that using these vectors led to significantly better predictions of future drug prescription costs compared to a baseline one-hot model ai.jmir.org ai.jmir.org. The embeddings captured comorbidities and health status in a way that linear models and even tree-based models could easily leverage. In another case, *Med2Vec* embeddings were used to predict inpatient mortality; the richer representation of patients' diagnoses improved the accuracy over traditional approaches bmcbioinformatics.biomedcentral.com. Embeddings have also been applied in predicting *heart failure* or other chronic conditions onset: a sequence model (like an LSTM or Transformer) that ingests a time series of embedded diagnoses can forecast the likelihood of a patient developing a certain disease. In those tasks, researchers observed that incorporating embeddings (especially those informed by medical ontology) allowed models like RNNs to perform better with less data, particularly for rare outcomes pmc.ncbi.nlm.nih.gov. In short, **embedding vectors serve as powerful features for machine learning models**, enabling more accurate risk scores and early identification of high-risk patients (for example, predicting which patients are at risk of hospital readmission or complications based on their encoded history).

- **Patient Similarity and Cohort Matching:** By representing patients as vectors (for instance, by averaging or otherwise aggregating their diagnosis code embeddings), we can measure patient-patient similarity in a quantitative way. This is very useful for finding cohorts of patients that have similar disease profiles or for case-based reasoning. Clustering algorithms can be applied to patient embeddings to discover **subpopulations or phenotypes** in the data. In the Pat2Vec study, the authors performed density-based clustering on patient vectors from 10 million patients and were able to identify meaningful subgroups in the cohort ai.jmir.org. When they visualized these patient embeddings in 2D (using UMAP for dimensionality reduction), distinct clusters emerged that corresponded to recognizable sets of diagnoses (subcohorts). Such clusters could represent, for example, a group of patients who have a combination of metabolic syndrome disorders, or a group who have both cardiovascular and renal issues, etc., even if the exact ICD codes vary individual to individual. This kind of *data-driven patient segmentation* can be used to tailor interventions – e.g. design specialized care programs for the cluster of multimorbid elderly patients with diabetes, hypertension, and kidney disease. Beyond clustering, nearest-neighbor search in the embedding space allows **patient similarity matching**: given a new patient, find past patients with the most similar diagnostic history. This can support clinical decision support (by examining what outcomes or treatments those similar patients had) and research (by finding matched controls for a study). Some healthcare applications have started to integrate such patient similarity engines, especially in personalized medicine contexts. The key enabler is the embedding: a good patient-level embedding will take into account the numerous diagnoses (and possibly other factors) in a way that two patients with the same overall condition end up near each other even if not every code is identical. As an example, if one patient has Type 1 diabetes and hypothyroidism, and another has Type 2 diabetes and rheumatoid arthritis, they might seem different if you look at raw ICD codes; but a learned embedding might place them closer because it recognizes both have endocrine/metabolic disorders plus an autoimmune condition. Indeed, a well-trained model can encode latent dimensions that correspond to clinical axes (like "autoimmune burden" or "metabolic health"), and patients can be compared on those. **In summary, embeddings make the concept of patient similarity computationally tractable and clinically meaningful**, facilitating cohort discovery and personalized analytics.

- **Clinical Decision Support & Recommendations:** Embeddings of ICD-10 codes can also directly assist in clinical decision support systems (CDSS). One use is **suggesting related diagnoses or alerts**. If a patient has been coded with certain diagnoses, an embedding model might suggest other diagnoses to consider, based on what usually occurs together. For instance, suppose a patient has an ICD-10 code for *rheumatoid arthritis* and *interstitial lung disease*; an embedding might be able to highlight the connection to *rheumatoid lung* (a complication) or suggest checking for *pulmonary fibrosis*. This could manifest as a tool that, given a set of known diagnoses, retrieves the nearest neighboring codes in the embedding space as "you might want to rule out or consider these". Such a system acts like a smart lookup that goes beyond simple hierarchy. Another application is **auto-completion or assistance in coding**: when clinicians document diagnoses, an embedding-based system could rank likely ICD-10 codes based on partial information. In a study of automated coding, researchers compared similarity search vs. machine learning for matching diagnosis descriptions to ICD-10 codes medrxiv.org. Embedding the descriptions and codes in the same vector space (via something like BioBERT) can enable finding the closest code vector to a given clinical text vector, thus suggesting the best code. This approach can improve coding efficiency and accuracy in hospital billing departments.

  Furthermore, embeddings can support *treatment recommendations* indirectly. If diagnoses are encoded well, they can be combined with other embeddings (say, embeddings for medications or procedures) in a unified space or model. For example, one could train a model where a certain medical procedure code vector is close to the diagnosis vectors for which that procedure is indicated. Then, given a patient's diagnosis vector, a CDSS could recommend considering the procedures or interventions whose vectors align with the patient's vector. This is analogous to how recommendation systems work in other domains (like suggesting products based on user embedding), but here it's suggesting clinical actions based on patient embedding. Some preliminary research has looked at *embedding entire EHRs including diagnoses, medications, lab tests* into one space, so that reasoning like this becomes possible. One can also use the similarity of code vectors to assist in **differential diagnosis**: e.g., if a physician enters a preliminary diagnosis, the system can show other diagnoses that are similar (so the doctor can ensure those aren't a better fit). An example from ICD2Vec demonstrated that their learned vectors correlated with actual disease biology github.com – meaning if one were to use their embedding to find "similar diagnoses" for a given disease, one often finds diseases that share pathophysiological features. This kind of knowledge could be turned into a decision support tool for clinicians (for instance, reminding them that certain infections have similar presentations). Overall, while still an emerging area, **CDSS can leverage embeddings to provide smarter, context-aware suggestions and checks**, thereby potentially improving diagnosis and management.

- **Phenotyping and Disease Characterization:** *Computational phenotyping* refers to the process of identifying meaningful disease subtypes or patient groups from data, often without a priori definitions. Embedding vectors have become a valuable ingredient in phenotyping algorithms. One approach, as done in **Phe2Vec**, is to define a phenotype by a "seed" ICD code and then use the embedding space to find other codes that cluster around it, effectively capturing the full spectrum of that phenotype medrxiv.org. For example, to phenotype *heart failure*, one might take the code for heart failure as a seed, then find the nearest neighbors in the embedding space (which might include codes for edema, shortness of breath, diuretics use, etc.). These neighboring concepts collectively represent the phenotype's manifestations. Phe2Vec demonstrated that using unsupervised embeddings of medical concepts, they could automatically derive phenotypes for diseases that matched or even outperformed expert-defined criteria (like the PheKB algorithms) medrxiv.org medrxiv.org. Essentially, the embedding space provided a way to *cluster clinical concepts around diseases*, and by aggregating those concepts per patient, they identified patients with the disease. This approach was validated by manual chart review, showing similar accuracy to the labor-intensive rule-based phenotypes medrxiv.org. The advantage is clear: once you have a high-quality embedding space, you can prototype a phenotype by just picking a point in that space and grabbing everything in its vicinity. This is much faster and more flexible than manually enumerating all criteria. Another use is discovering **new sub-phenotypes**. Researchers can apply clustering to the embeddings of patients or even embeddings of diagnosis combinations. For instance, by clustering patient vectors of those all labeled with a broad condition like "diabetes", one might find distinct subgroups – e.g. one cluster of diabetic patients also has obesity and hypertension (metabolic syndrome cluster), another has frequent hospitalizations for renal issues (maybe a diabetic nephropathy cluster), etc. These data-driven phenotypes can generate hypotheses for precision medicine: perhaps each cluster responds differently to treatments. In one study, unsupervised clustering of embeddings identified a subgroup of oncology patients that had significantly different outcomes, effectively finding a hidden phenotype that wasn't obvious just by a single ICD code ascopubs.org. **Phenotyping with embeddings is a powerful way to leverage the information densification that embeddings provide**, enabling recognition of patterns that involve combinations of diagnoses. This can lead to better disease definitions and improved inclusion criteria for clinical trials or observational studies.

- **Healthcare Operations and Resource Allocation:** On a more system level, aggregated embeddings can support health system management. If patient populations are embedded, a hospital could, for example, project all their patients into an embedding space and see how they are distributed. This could inform resource allocation by identifying how many patients fall into certain health categories. The Pat2Vec authors mentioned *health care resource planning based on subcohorts* identified by their embedding framework [ai.jmir.org](ai.jmir.org). For example, suppose an embedding cluster corresponds to patients with complex neurological disorders – the administration can ensure they have enough neurologists or clinic slots for that subgroup. Another cluster might correspond to frail elderly with multiple needs – indicating a need for integrated care programs. Because embeddings allow integrating many data types, operations research can combine them with cost data, length-of-stay data, etc., to make predictions about utilization. For instance, one can train a model to predict future *healthcare costs or hospital admissions* using patient embeddings as inputs [ai.jmir.org](ai.jmir.org). Those predictions help in budgeting and managing care management programs (like flagging patients for care coordination). Additionally, public health surveillance can benefit: by monitoring shifts in the distribution of patient embeddings over time, one could potentially detect emerging health trends (e.g. a new syndrome might cause a cluster of patients to start appearing as outliers in the embedding space, cueing further investigation).

These use cases illustrate the versatility of embedding ICD-10 codes. In practice, many projects combine multiple uses – for example, a population health management platform might use embeddings to segment patients (similarity use case) and then apply a predictive model to each segment (predictive modeling use case) to identify at-risk individuals, and finally suggest interventions (decision support use case). What makes all this feasible is that embeddings turn the messy, high-dimensional diagnosis data into something machines can more easily learn from and we can reason about quantitatively. It's worth noting that many of these applications have been validated in research settings (with papers reporting improved metrics), and they are gradually translating into real healthcare tools. As an example, some electronic health record systems now incorporate risk scores that implicitly rely on embedded representations of patient history, even if under the hood. We anticipate even broader adoption as the healthcare industry seeks to make sense of big data and move toward personalized, data-driven care.

## Challenges and Limitations

While encoding ICD-10 codes as embeddings offers numerous advantages, it is not without challenges. Both the embedding process and its use in practice come with limitations and considerations, especially in the sensitive domain of healthcare. Here we outline several key challenges, along with ethical considerations:

- **Data Quality and Coding Practices:** ICD-10 codes are only as good as the documentation and coding process. If the underlying data is noisy or biased (e.g. physicians have varying habits in how they choose codes, or some conditions tend to be under-coded/over-coded), the embeddings will learn those patterns – which may not reflect true patient state. For instance, socioeconomic factors can lead to differences in coding (a condition might be coded more often in some hospitals than others due to billing focus). So, embeddings might inadvertently capture hospital-specific or practice-specific patterns rather than universal clinical truth. This is a garbage-in, garbage-out issue: *if coding data is incomplete or inconsistent, the embedding will embed those flaws*. An example is mental health codes which might be underdiagnosed; their embeddings might be less reliable or far from related conditions because they seldom appear. It's crucial to be aware of these data artifacts.

- **Rare Codes and Unseen Codes:** The long tail of ICD-10 means many codes are infrequent. Standard embedding algorithms (especially Word2Vec/GloVe) may not learn meaningful vectors for extremely rare codes due to insufficient data bmcbioinformatics.biomedcentral.com. While methods like FastText or ontology-informed approaches help, rare codes can still end up with noisy embeddings. Similarly, new codes (ICD-10 gets updated regularly, and ICD-11 is on the horizon) will not be in the embedding vocabulary if the model was trained on older data. Updating embeddings to accommodate new codes is non-trivial; one might have to retrain or at least fine-tune with new data. A related challenge is **out-of-vocabulary** handling: if an embedding model encounters a code in a patient that it never saw in training, it can't directly produce an embedding (unless using subword tricks or a fallback mechanism). Some approaches map unknown codes to a generic "unknown" vector, but that loses the specific meaning. Ongoing maintenance of the embedding model is needed in a live system to handle code updates.

- **Hierarchical and Contextual Limitations:** Not all embedding techniques perfectly capture the ICD hierarchy or context. We noted that plain Word2Vec struggles with hierarchical relations bmcbioinformatics.biomedcentral.com – it might put some codes near others that often co-occur, even if they are from very different categories (for example, if two unrelated conditions often happen in the same patients due to a comorbidity pattern, Word2Vec might conflate them). This could potentially confuse interpretations. Also, static embeddings do not capture context by themselves; "contextual embeddings" require more complex models. Without context, the same code has one vector regardless of situation, which might be a limitation. For example, the code for "Chest pain" (R07.9) could mean very different things (heart attack vs musculoskeletal pain) in different patients – a single embedding can't capture all those nuances. Some advanced models address this by incorporating additional information (age, sex, other codes) when generating patient embeddings, but as a limitation, *most off-the-shelf embeddings are context-agnostic*. This means one should be careful not to over-interpret the embedding of a code without considering the patient's full context.

- **Bias and Fairness Issues:** Embeddings can encode biases present in the data. If certain diagnoses are more frequently recorded in one demographic group due to structural biases, the embedding space might reflect that in undesirable ways. For instance, if historically a condition was underdiagnosed in women, the "female" patients' vectors might cluster in a way that models mispredict or underemphasize that condition for them. Similarly, racial biases in healthcare (like unequal access leading to different co-morbidity patterns) can become baked into embeddings. An AI model might then perpetuate these biases, leading to unfair outcomes foreseemed.com foreseemed.com. For example, an algorithm might learn an embedding dimension that partly correlates with race or socioeconomic status (because those factors influence disease patterns), and then inadvertently use that to allocate less resources to a group, thinking they are lower risk when actually it's underdiagnosis. It is challenging to detect and debias embeddings in healthcare – one has to analyze the embedding space and outcomes carefully. Techniques from NLP to debias word embeddings (e.g. remove gender bias) could potentially be applied to medical embeddings (for removing bias related to race or insurance status), but this is an active area of research. The **ethical implication** is that using biased embeddings could exacerbate health disparities. We must ensure training data is representative and possibly apply corrections. Stakeholders increasingly emphasize testing these models for bias: e.g. does the risk model using embeddings perform equally well for all ethnic groups? Does the patient similarity measure inadvertently cluster patients by socioeconomic status more than by clinical features? **Transparency and fairness audits** are needed when deploying such models.

- **Interpretability and Trust:** Embeddings are inherently abstract. Each dimension of an embedding vector doesn't have a clear meaning that a clinician can understand (they are not like specific risk factors, but rather latent combinations of features). This can make it hard to explain *why* a model made a certain prediction using those embeddings. Clinicians might be wary of a "black-box" representation. For example, if a model flags a patient as high-risk because of their embedding, the doctor might ask: which diagnoses or factors contributed? With one-hot features, one could at least see which codes had weight. With an embedding, it's more opaque – the information is distributed across many dimensions. Some work like GRAM aims to improve interpretability by tying embeddings to ontology concepts (so you can say "this patient's vector is similar to chronic lung disease category"). But in general, **the loss of explicitness is a trade-off**. It's important to develop tools to interpret embeddings: one could, for instance, find the nearest code vectors to a patient's embedding to give a sense of what it represents ("this patient's embedding is closest to those of patients with COPD and heart failure"), which can be communicated to clinicians. Ensuring trust will likely involve such post-hoc explanations.

- **Privacy and Security:** On the surface, using embeddings might seem neutral for privacy since they abstract away direct identifiers. Indeed, if we only use ICD codes (which are not direct identifiers) to create an embedding, one might assume privacy is not a big concern. However, patient embeddings encapsulate a lot of personal health information. If someone obtained a patient's embedding vector, could they infer sensitive details? Potentially yes, if they had reference points. For example, an embedding could reveal that a patient likely has a certain disease cluster. There's also the risk of re-identification: if an embedding is unique enough (patients with rare combinations of diagnoses might have distinctive vectors), and an adversary has some knowledge of a person's diagnoses, they might pick them out of an embedding database. This is somewhat speculative, but privacy researchers have shown concerns about models memorizing data. A deep model that produces embeddings might inadvertently encode a rare code in a way that leaks information. From a compliance perspective (e.g. HIPAA), one needs to treat derived data carefully. Usually, aggregating and de-identifying should make embeddings safe, but it's worth noting that these are not just random numbers – they are structured encodings of health states. Thus, organizations often keep control of such models and don't share raw patient embeddings publicly. There's interest in **federated learning** approaches for embeddings, where different institutions train embeddings on their data and then share or combine them without exposing individual-level data.

- **Need for Re-training and Generalizability:** An embedding trained in one context might not generalize perfectly to another context. For example, an embedding learned from German claims data (ICD-10-GM) ai.jmir.org might encode certain patterns specific to that healthcare system. If you apply it in a U.S. hospital, it may work differently (some codes may be used differently, or patient populations differ). Likewise, an embedding trained on 2018 data might not reflect medical knowledge in 2025 if new treatments or coding guidelines changed practice. Thus, embeddings may need periodic retraining or adaptation. If a model was deployed and never updated, over years its performance could drift as the data distribution shifts (this is a general problem of dataset shift). Specifically for ICD, the transition to **ICD-11** in the future will be a big change – embeddings will have to be learned anew for ICD-11 codes, or cross-walked from ICD-10. Ensuring continuity (so that predictive models remain accurate) is a challenge. Cross-mapping embeddings (learning a transformation from ICD-10 space to ICD-11 space) could be one research direction to handle the transition.

- **Evaluation and Validation:** How do we know an embedding is "good"? There's no single right answer – an embedding that works well for one task might not for another. It's important to validate embeddings on multiple tasks (clustering quality, prediction improvements, clinician assessment of nearest neighbors etc.) researchgate.net researchgate.net. This validation can be time-consuming and may require expert feedback (for example, asking physicians if the nearest neighbors of code X in the embedding space make clinical sense as related concepts). Without careful validation, one could deploy an embedding that has odd quirks or fails for certain subsets of codes. As a simple example, perhaps the embedding didn't cluster mental health codes well because they rarely co-occur with physical health codes in the data used – a mental health professional might notice that and adjust the approach (maybe by adding some domain constraints or data). **In critical applications, regulatory requirements might also demand validation** – for instance, if an embedding-driven model is used to influence patient care, regulators might treat it like a medical device which needs evidence of safety and effectiveness.

In summary, while **embedding ICD-10 codes is a powerful technique, it must be approached with caution and domain awareness**. The limitations span technical (data sparsity, updating models), interpretative (making sense of the vectors), and ethical (bias, privacy) domains. Many of these challenges are active areas of research. For example, researchers are developing methods to create *time-aware embeddings* that evolve as data evolves medrxiv.org bmcbioinformatics.biomedcentral.com (addressing static vs dynamic issues), and methods to inject fairness constraints into training to mitigate bias. Collaboration between data scientists, clinicians, and ethicists is important to navigate these challenges, ensuring that embedding-based tools truly benefit patients across the board and do not inadvertently harm or discriminate. As with any powerful technology, careful validation and oversight are key in the healthcare context.

## Future Research Directions

The field of medical code embeddings is rapidly advancing, and there are many exciting avenues for future exploration. Below are some potential future research directions and improvements that could further enhance the encoding of ICD-10 (and other clinical codes) in vector spaces:

- **Integrating Multimodal Data:** Future embeddings are likely to go beyond diagnosis codes alone and incorporate multiple data types – procedures, medications, lab results, clinical notes, imaging findings, genetic data, etc. A truly comprehensive patient embedding might combine ICD-10 codes with other coding systems (like CPT procedure codes, ATC drug codes) and even unstructured text from clinical notes. Research in this direction includes joint embedding spaces where, say, a disease code and a relevant lab test abnormality vector are close because they indicate the same condition. One example is cui2vec which already blended data from clinical text, claims, and knowledge bases for UMLS concepts bmcbioinformatics.biomedcentral.com. In the future, we may see *fusion models* that use transformers to read clinical text and graph networks to process code ontologies, merging their outputs. This could yield embeddings that capture a "360-degree view" of a patient. Such rich embeddings could be used for holistic predictions (like forecasting overall health outcomes, not just single disease risks). They also open doors for transfer learning between modalities – for instance, using patient note information to inform code embeddings and vice versa.

- **Contextual and Dynamic Embeddings:** As healthcare data are sequential and time-dependent, an important research direction is creating embeddings that are context-sensitive and that can evolve. Rather than each ICD code having one fixed vector, we might have embeddings that are a function of patient context or time. Techniques from NLP like ELMo or GPT (which give contextual word embeddings depending on the sentence) could be adapted so that the embedding of "diabetes" in a patient with only mild disease is different from "diabetes" in a patient with many complications. *Dynamic embeddings* could also account for changes over time – e.g., the meaning of certain codes might shift if guidelines change, or a patient's state changes. One could imagine *time-aware embeddings* where the vector representation can drift year by year, or which include an age dimension (so a code vector is modulated by patient age, capturing that certain diagnoses have different implications in pediatric vs geriatric populations). Early work has attempted to incorporate time by training separate embeddings for different time windows or by adding time as an input to the model bmcmedinformdecismak.biomedcentral.com. With the rise of sequence models in EHR, this trend will continue: models like Transformers may directly output contextualized embeddings for codes as they appear in sequence. Research could also explore **hyperdynamic embeddings** using recurrent networks that update a patient's embedding with each new event, effectively learning how the vector should move in space as conditions accumulate or resolve.

- **Improved Handling of Hierarchies and Ontologies:** While some methods incorporate hierarchies, there is room for more advanced approaches, especially with the advent of *hyperbolic embeddings*. Hyperbolic geometry has been shown beneficial for embedding tree-structured data (because it can naturally represent hierarchies with less distortion). Future models might embed the ICD tree in a hyperbolic space to better respect parent-child distances (there is already some work on hyperbolic embeddings for medical ontologies dl.acm.org). Additionally, combining multiple ontologies is a challenge – e.g., linking ICD-10 with SNOMED, with phenotype ontologies, etc. A unified embedding of multi-ontology knowledge graphs could be valuable. We might also see more use of **graph neural networks** (GNNs) beyond simple GCN. For instance, a *Graph Attention Network* could learn which relationships in the ontology are more important and weight them. The work on Hierarchical Attention (like Hierarchical Attention Propagation) dl.acm.org is along these lines, and future research could extend that. As ICD-11 is implemented, which has an even more complex hierarchical and cross-referenced structure, techniques to embed such complex ontologies will be crucial. There may be efforts to seamlessly map ICD-10 embeddings to ICD-11 via the known mapping, or to train joint embeddings that span both ICD-10 and ICD-11 for compatibility during the transition.

- **Personalized and Subgroup-Specific Embeddings:** Another direction is recognizing that one size may not fit all. Perhaps separate embeddings trained for subpopulations (e.g., pediatric vs adult patients, or for different ethnic groups) might capture nuances better than a single embedding for everyone. This is tricky because one doesn't want to bake bias, but there might be legitimate physiological differences or disease prevalence differences that a single embedding space can't represent well. Future research might explore training embeddings that adapt to patient metadata. For example, a conditional embedding model that takes patient demographics as input and produces slightly adjusted code vectors tailored to that profile. This could improve accuracy of models in personalized medicine. Some evidence of this need is seen in models that had to adjust for sex-specific differences (e.g., heart attack symptoms differ in men vs women; a contextual embedding might handle that).

- **Embedding Explainability and Alignment with Clinical Knowledge:** To address interpretability concerns, research could focus on making embeddings more transparent. One idea is to align embedding dimensions with known clinical factors. This might involve *post hoc* rotation of the embedding space to maximize correlation with known variables (like one dimension might align with "chronic disease burden"). Or using sparse embeddings where each dimension is encouraged to correspond to a cluster of codes that is interpretable (some works use autoencoders with regularization to achieve something like this). Another approach is to build **explainable embeddings** by design – e.g., an embedding method that forces the vector of a code to be a weighted combination of vectors of a small number of medical concepts (maybe leveraging an ontology or expert-defined features). This would be slower and more supervised, but could yield dimensions that are easier to label (like a dimension that clearly separates cancer vs non-cancer). Also, developing visualization tools for embedding spaces can help – researchers will likely work on interactive tools where clinicians can navigate the code embedding space, see clusters, and attach labels to regions of the space (like "this region is all cardiovascular-related codes"). This can build trust and also feed back to refining embeddings.

- **Benchmarking and Standardization:** As many groups develop their own embeddings using various data sources (claims, EHRs, international datasets), it would be useful to have benchmarks to compare them. Future work might establish benchmark tasks, such as a standard set of analogies or clustering tasks for medical codes, similar to how NLP has analogy tasks (e.g., "Hypertension is to HeartDisease as Asthma is to ____?"). Some initial work shows analogical reasoning with ICD embeddings is possible [sciencedirect.com](sciencedirect.com). A community-driven evaluation platform (like **MD2Vec Challenge** etc.) could emerge. Additionally, with many embeddings floating around, there might be efforts to **standardize representations** or at least catalog them (for example, a repository where one can download pre-trained embeddings for ICD-10 trained on MIMIC-III, on CMS claims, on UK Biobank, etc., and some guidance on which to use). This would parallel the NLP field where one can choose between GloVe, word2vec Google News, etc. based on needs.

- **Adapting to ICD-11 and Beyond:** ICD-11 has been released by WHO with a much more complex structure (allowing multiple axes, post-coordination, etc.). Encoding ICD-11 will be a new challenge and opportunity. Research can explore if the techniques used for ICD-10 embeddings hold or if new methods are needed. Possibly, ICD-11 being closer to a graph with multiple parents might benefit even more from graph embedding techniques. Also, the transition period where both ICD-10 and ICD-11 are in use will require embedding spaces that can talk to each other. Future research might involve *training a model to automatically convert an ICD-10 code embedding into an ICD-11 code embedding* (leveraging the official mapping or machine learning). This could help in comparative studies and in migrating models to the new coding system with minimal retraining.

- **Real-time and Deployment Considerations:** From a practical standpoint, there is scope to research how to efficiently deploy embedding models in clinical settings. For example, can embeddings be computed in real-time as data streams in from monitors or new diagnoses get added? This leans into systems research: perhaps simplified embedding models that can run on edge devices or within EHR software. Also, methods for incremental learning: updating embeddings with new data without full retraining. Some online learning algorithms or embedding fine-tuning methods could be applied, so that an installed model keeps improving as it sees more cases (while avoiding catastrophic forgetting of past knowledge). Ensuring stability (so that small updates don't drastically change the space and confuse downstream systems) is part of that research.

- **Ethical and Bias Mitigation Research:** Finally, given the concerns about bias and fairness, future work will likely include developing **fairness-aware embeddings**. This could mean techniques that remove or neutralize the impact of protected attributes (like ensuring the embedding does not implicitly encode race unless absolutely medically relevant). Adapting methods from word embeddings (where they, for instance, subtract "gender" direction) might help foreseemed.com foreseemed.com. Another angle is privacy-preserving embeddings: using differential privacy during training so that the embeddings do not leak information about any single individual. This is an active area in general machine learning and can be applied here (though it often comes at a performance cost). With regulations like GDPR and HIPAA, demonstrating that patient embeddings cannot be reverse-engineered to identify individuals could become necessary if these are used in shared systems.

In conclusion, the future of encoding medical codes is likely to involve **richer models, more data sources, and a greater emphasis on trustworthiness**. The trajectory mirrors what we see in NLP and recommender systems, but with domain-specific twists. We expect that in the next few years, medical code embeddings will become a standard part of the healthcare AI toolkit – much like word embeddings are in NLP – and ongoing research will make them more powerful, interpretable, and equitable. The ultimate goal is that these representations enable health professionals and systems to glean insights and make predictions that improve patient care, while minimizing unintended consequences. Each research advance – whether it's a new algorithm that better captures medical knowledge, or a new way to validate embeddings against clinical benchmarks – will bring us closer to embedding technologies that are robust, fair, and beneficial across the healthcare landscape.

**References:** The information in this report was gathered from a wide range of academic studies and sources that demonstrate the concepts discussed. Key references include a 2023 JMIR study proposing the Pat2Vec embedding framework and showing its advantages over one-hot encoding ai.jmir.org ai.jmir.org, the description of ICD-10 and its uses by WHO and practice guidelines practicefusion.com bmcbioinformatics.biomedcentral.com, research on Word2Vec/Doc2Vec applications to clinical codes ai.jmir.org bmcbioinformatics.biomedcentral.com, the ICD2Vec project highlighting semantic relationships among disease embeddings github.com, the Phe2Vec framework for phenotyping medrxiv.org, and many others. These citations (noted in the text like $X^+$Ly-Lz) provide detailed evidence and examples, and interested readers are encouraged to consult them for deeper technical details or specific results.

## DISCLAIMER

The information contained in this document is provided for educational and informational purposes only. We make no representations or warranties of any kind, express or implied, about the completeness, accuracy, reliability, suitability, or availability of the information contained herein.

Any reliance you place on such information is strictly at your own risk. In no event will IntuitionLabs.ai or its representatives be liable for any loss or damage including without limitation, indirect or consequential loss or damage, or any loss or damage whatsoever arising from the use of information presented in this document.

This document may contain content generated with the assistance of artificial intelligence technologies. AI-generated content may contain errors, omissions, or inaccuracies. Readers are advised to independently verify any critical information before acting upon it.

All product names, logos, brands, trademarks, and registered trademarks mentioned in this document are the property of their respective owners. All company, product, and service names used in this document are for identification purposes only. Use of these names, logos, trademarks, and brands does not imply endorsement by the respective trademark holders.

IntuitionLabs.ai is an AI software development company specializing in helping life-science companies implement and leverage artificial intelligence solutions. Founded in 2023 by Adrien Laurent and based in San Jose, California.

This document does not constitute professional or legal advice. For specific guidance related to your business needs, please consult with appropriate qualified professionals.