

---

# Phenotypical Ontology Driven Framework for Multi-Task Learning

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1       The increasing availability of rich longitudinal Electronic Health Records (EHR)  
2       has driven the development of AI methods aimed at generating valuable clinical  
3       insights. patients in EHRs, cohorts extracted to derive outcome predictive models  
4       for specific phenotypes often result in small and imbalanced data sets that are  
5       of limited practical use for AI modeling. In this paper, we propose an Ontology  
6       driven Multi-Task Learning framework called OMTL that is able to cope with  
7       these data limitations for EHR modeling by leveraging external knowledge. The  
8       model is especially designed for situations where we want to study patients for  
9       a set of related disease categorizations, and interested in potentially different  
10      outcomes corresponding to different categories. More specifically, OMTL leverages  
11      phenotype similarity according to a predefined ontology to learn the outcome  
12      prediction tasks. The model structure allows for an increased sharing of data  
13      representations across phenotype profiles such that the information sharing is  
14      regulated by the ontological and patient level similarities. We evaluate the benefit  
15      of leveraging an ontology to relate tasks in a multi-task learning settings and  
16      demonstrate the efficacy of OMTL on several real outcome prediction problems  
17      on patients from the MIMIC III database over state-of-the-art multi-task learning  
18      schemes.

## 19   1 Introduction

20   The increasing availability of rich longitudinal patient records from large EHR has enabled novel  
21   data-driven analysis on such observational data towards clinical hypothesis generation and testing.  
22   EHR patient data can cover several aspects of patient history such as diagnosis, medication, and lab  
23   tests resulting in a diverse set of patient attributes.

24   Despite the large number of patients covered in these data sets, formulating a query to extract a  
25   cohort for the outcome of interest often results in a very modest number of data points being returned,  
26   hindering the effectiveness and reproducibility of deep learning techniques on these data sets. Let us  
27   consider a scernario where we want to study patients who have experienced Rheumatic disease of  
28   hearth valve and Mitral valve stenosis on the well-known MIMIC III data set [10]. Such a query only  
29   returns 210 ICU stays, down from the original 42000 stays present in the complete data set. In another  
30   case, on an EHR database spanning more than 55 million patients, querying for patients within  
31   the age group 21 – 40 and having experienced concussions that led to mild cognitive impairments  
32   returns fewer than 1100 patients. Further characterization of this cohort to focus only on patients that  
33   eventually experienced dementia takes that number below 200.

34   Furthermore, such histories are often sparse and difficult to be directly utilized for simple analysis.  
35   Indeed, being observational in nature, EHR data sets often exhibit widely varying regularity of  
36   encounters as well as diseases coverage. This data imbalance can manifest itself as a ‘vanishing data’

37 problem during cohort analysis of specific diseases of interest with low coverage. Any analysis on  
38 such data sets can be impractical and suffer from significant lack of robustness.

39 Interestingly, a significant amount of domain knowledge relating medical concepts is readily avail-  
40 able in several medical ontology. Phenotypes corresponding to observable physical properties of  
41 patients [1], including diseases have been organized in ontologies and expressed within EHR systems.  
42 We postulate that leveraging phenotypical relations within EHR data sets to tie and jointly learn  
43 multiple learning outcomes in a multi-task learning setup could provide solutions to this vanishing  
44 data problem and lead to more robust prediction models.

45 In this paper, we propose OMTL, a novel phenotypical ontology-driven framework for multi-task  
46 learning. The main tenant of our approach is to ground the structure of the proposed multi-task  
47 framework on phenotypical medical concepts relations between patients from different medical cohorts.  
48 Using OMTL allows us to cope with the lack of data for training specific outcomes by leveraging  
49 existing knowledge encapsulated in ontologies. Learning outcomes for different phenotypes may bor-  
50 row representation expressiveness during training for the prediction of several outcomes. This sharing  
51 allows OMTL to generate prediction models for patient phenotypical cohorts even for phenotypes with  
52 little amount of patient observations by leveraging data from other phenotypical cohorts.

53 Moreover, OMTL enables sharing information among phenotypes based on the ontology which allows  
54 the model to decouple the outcome learning from the representation learning of similar diseases.  
55 In this case, the model allows studying patients with related diseases with different outcomes. For  
56 example, it can learn about mortality prediction for diabetic patients and readmission prediction for  
57 patients with hypertension by sharing information between diabetes and hypertension.

58 The contributions of this paper are 3 folds:

- 59 1. We proposed a systematic way to tie different phenotypical outcome using existing domain  
60 knowledge in the form of ontologies of medical concepts (e.g. SNOMED CT [4]).
- 61 2. We propose a novel hierarchical deep learning architecture using gating mechanisms mirror-  
62 ing the structure of a domain ontology.
- 63 3. We demonstrate the efficacy of the proposed approach on the prediction of several outcomes  
64 on patients from the MIMIC III database.

65 In the remainder of the paper, we start with a review of related work before exposing in details the  
66 proposed method OMTL. We then present the experimental results geared towards the evaluation of  
67 the gain achieved by leveraging domain ontologies to link together multiple tasks in OMTL.

## 68 2 Related Work

69 **Multi-task learning:** Multi-task learning has received a great deal of attention in the last few years.  
70 More notably, [12] adapts Mixture-of-Experts (MoE) structure to jointly learn multiple outcomes  
71 by sharing a single expert, backed by gating networks to specialize for each outcome and learn the  
72 relationships jointly. In [8, 2], authors focus on using hierarchical relationships to learn multiple  
73 outcomes. Specifically, a two-level hierarchy has been used to regularize a learned model to share  
74 information among related outcomes. [21] extended these approaches with multi-level hierarchies.  
75 They proposed a cascade approach to progressively learn outcomes at prior levels to boost learning  
76 performance on outcomes in immediately subsequent levels.

77 The model proposed in [14] learns intermediate representations for predictive tasks at multiple  
78 scales for the same input via a cascade of connected layers. In [15], the authors proposed a way  
79 to incorporate a set of *carefully* selected semantic NLP tasks into the multi-task model. The tasks  
80 that are “easy” to learn are supervised at the lower levels of the architecture while more complex  
81 interactions are kept at deeper layers. A soft-ordering of the shared layers was proposed in [13] to  
82 learn both the layers and the order of layers for each task.

83 OMTL is motivated by such efforts to learn phenotype specific predictive models in a multi-task  
84 setting to share information across outcomes. It is also important to note that OMTL uses an ontology  
85 based hierarchy, but more crucially, decouples tasks from hierarchies such that representations learned  
86 for different phenotypes can be further specialized w.r.t. different tasks such as readmission and

87 mortality prediction. In most of the cases, the shared layers are integrated in parallel orders, where  
 88 the outcome of one layer is fed into the next layer.

89 **Knowledge graph:** Structured domain knowledge, often in the form of graphs, have been applied in  
 90 various applications to achieve more robust inferences over only data-driven models. [18] presents a  
 91 framework to enable the use of various kinds of external knowledge bases to retrieve answers relevant  
 92 to a given question. It provides a graph based model which maps the text phrases to concepts in the  
 93 external knowledge graph with several strategies: concepts only, one-hop, and two-hop. Finally, the  
 94 extracted graph is used as an input to a neural network model for entailment.

95 **Outcome prediction from EHR:** Over the years, many researchers have developed predictive  
 96 models on EHR for the prediction of various adverse events such as early identification of heart  
 97 failure [16], readmission prediction [19], and recently acute kidney injury prediction [17]. Other  
 98 notable works include [20] where the authors presented a multi-level attention mechanism to derive  
 99 patient specific representations. Specifically, the authors split the trajectory of patients into sub-  
 100 sequences and use within-subsequence and between-subsequence attention levels sequentially to  
 101 arrive at patient representations in order to handle various data issues such as data sparsity and noise.  
 102 Of note is also the work by [7] where the authors present four clinical prediction benchmarks using  
 103 data derived from the publicly available Medical Information Mart for Intensive Care (MIMIC-III)  
 104 database, including In-hospital mortality, Decompensation, Length of stay, and Acute care phenotype  
 105 classification. In addition, it proposes a deep learning model called multi-task Recurrent Neural  
 106 Networks to empirically validate the four prediction benchmarking tasks. This benchmark has been  
 107 used as the basis for evaluating the effectiveness of OMTL.

### 108 3 Methods

109 To tackle the vanishing data problem, OMTL uses augmented patients having *similar* phenotypical  
 110 profiles to the original cohort being studied. Specifically, OMTL builds on the recently published  
 111 work of [12] using a multi-task mixture-of-experts model to share information across different  
 112 learning tasks. This model essentially learns different shared patient data representations (called  
 113 experts) fed as inputs for each individual learning outcomes. OMTL takes this multi-task learning  
 114 approach further by leveraging the relationship and ontological distance among phenotypes to share  
 115 information between learning tasks.

116 We assume the existence of an ontology for the domain, such as SNOMED [4] and CSS [6], that  
 117 represents the semantic phenotypical relationships among the medical concepts describing the data  
 118 fields of interest (e.g. diagnosis codes). We represent this ontology with a directed acyclic graph over  
 119 of a set of nodes (e.g. diagnosis codes) and edges expressing parent-child *semantic* relationships  
 120 between the nodes. Let  $X = \{x_i \in \mathbb{R}^d; \forall i\}$  be the EHR data where  $x_i$  is the data record  $i$  (e.g. an  
 121 intensive care unit visit or an admission for a patient), and  $d$  is the number of features<sup>1</sup>.

122 For each data record, we assume the presence of all phenotypical medical concepts tied to the record.  
 123 Using that information, we construct a subgraph  $\mathcal{G}$  of the ontology with only the nodes corresponding  
 124 to medical concepts that are present in  $X$ . Within  $\mathcal{G}$ , the set of user-specified phenotypes of interest  
 125 corresponding to the cohort under study (collection of medical concepts) are represented by a subset  
 126 of nodes, which we call core nodes  $\mathcal{C}$  (blue nodes in  $\mathcal{G}$  depicted in Figure 1). The process for obtaining  
 127 such a graph is described in details in the supplementary materials.

128 Any outcome prediction on  $X$  can be formulated as a mapping from  $X$  into an outcome set  $Y$ . To  
 129 enable sharing across outcomes, we break down this mapping as follows. We first define an encoding  
 130 function  $\mathcal{E}$  generating representations as:

$$\mathcal{E} : X \rightarrow Z \quad (1)$$

131 where  $Z = \{z_i \in \mathbb{R}^{d'}\}$ ,  $d'$  is the dimension of the representation, and  $z_i$  represents the representation  
 132 for record  $i$ . For this representation to be useful, we enforce an unsupervised reconstruction process  
 133 by learning a reconstruction function  $\mathcal{R}$ :

$$\mathcal{R} : Z \rightarrow X \quad (2)$$

---

<sup>1</sup>For simplicity we assume that  $x_i \in \mathbb{R}^d$  but it is straightforward to apply the method to temporal data where  
 $x_i \in \mathbb{R}^{T \times d}$

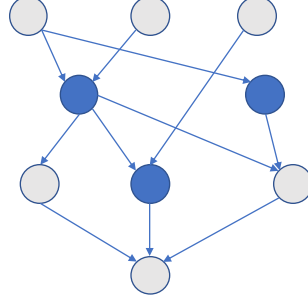


Figure 1: A directed acyclic subgraph that encompasses nodes represented in the dataset. The nodes in blue (i.e. core nodes) are of specific interest to the user.

The generated representation can be used as an input to a standard supervised machine learning algorithm to predict the outcome of interest  $Y = \{y_i; \forall i\}$ . This can be achieved by learning an outcome prediction function  $\mathcal{E}_S$ :

$$\mathcal{E}_S : Z \longrightarrow Y \quad (3)$$

Learning the parameters for the mappings  $\mathcal{E}$ ,  $\mathcal{R}$  and  $\mathcal{E}_S$  can be achieved by minimizing the following objective function:

$$\mathcal{L} = \lambda \underbrace{\sum_i \|\mathcal{R}(\mathcal{E}(x_i)) - x_i\|_2}_{\mathcal{L}_1} + \underbrace{\sum_i \ell(\mathcal{E}_S(\mathcal{E}(x_i)), y_i)}_{\mathcal{L}_2} \quad (4)$$

where,  $\mathcal{L}_1$  is the reconstruction loss function and  $\mathcal{L}_2$  is the total outcome prediction loss that aggregates the losses  $\ell$  (e.g. cross entropy for classification and mean square error for regression) at the node levels corresponding to the potentially different outcomes.  $\lambda \in \mathbb{R}^+$  is a hyper-parameter that weighs the importance of the two losses.

Learning a representation and a predictive model for one task is relatively straightforward by optimizing Equation (4). However, learning such representations for multiple patient cohorts with different phenotypical profiles and multiple predictive tasks simultaneously is a challenging problem. OMTL tackles the problem by using a structured hierarchical network that mimics  $\mathcal{G}$  and account for dissimilar coverage of patients across nodes. We formalize OMTL by first describing a multi-task learner and, subsequently, complete the formalization by describing how we include the hierarchical ontological structure as part of the model.

### 3.1 Multi-gate Mixture of Experts for Multi-task

Inspired by [12] a multi-gate mixture of experts approach can be used to learn  $\mathcal{E}$  in a multi-task settings as in Figure 2 (a). Input data flows from the top into a mixture of  $E$  expert nodes responsible for clustering the input data subspace to facilitate representation computations. We denote these expert nodes as  $E_e(\cdot)$ ,  $1 \leq e \leq E$ . These expert nodes are neural networks producing a representation  $E_e(x_i)$  given input  $x_i$ .

The outputs of these expert nodes are consumed by an array of  $P$  representation nodes denoted  $P_p(\cdot)$ ,  $1 \leq p \leq P$ . Each  $P_p$  maps directly into one medical concept (i.e. the number of representation nodes is equal to the number of nodes in the subgraph  $\mathcal{G}$  defined in the previous section). Each of these  $P_p$  nodes ingest experts outputs via a gating function  $G_p$  to learn how to combine information from all expert nodes for the generation of representations for the current medical concept. The gating node  $G_p$  controls the combination of expert nodes to produce the input  $\mathcal{M}_p$  to the representation node  $P_p$ . These are typically very simple networks that produce combination weights for the experts at each representation node via following Equation:

$$G_p(x_i) = \text{Softmax}(x_i \cdot W_p + b_p) \in \mathbb{R}^E$$

$$\mathcal{M}_p(x_i) = \sum_{e=1}^E G_p^e(x_i) E_e(x_i) \in \mathbb{R}^{d_e} \quad (5)$$

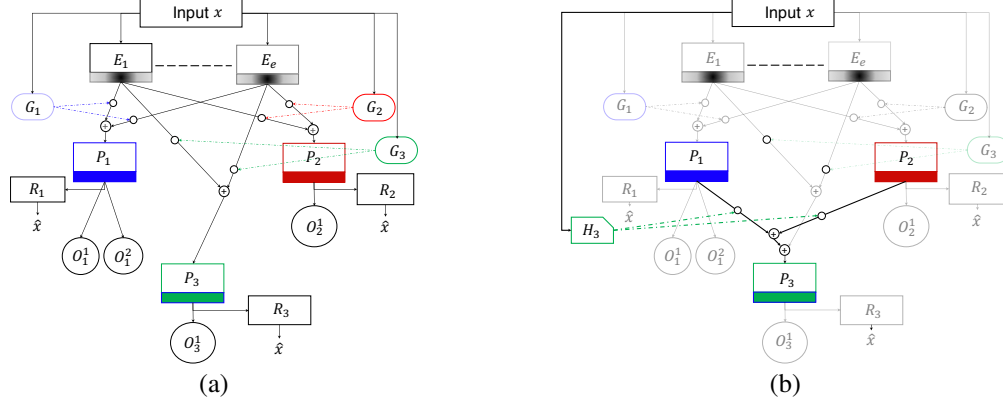


Figure 2: (a) Multi-gate mixture of experts architecture. (b) OMTL: The model is explained using only 3 nodes  $P_1, P_2, P_3$  (medical concepts) to minimize confusion. The modules that are the same as in the baseline model are faded out for clarity, while the new additional modules are plotted in bold.

where  $W_p \in \mathbb{R}^E$ ,  $b_p \in \mathbb{R}$  are parameters to be learned by the gate network during training, and  $G_p^e(x_i)$  is the entry  $e$  in the vector  $G_p(x_i)$ . The input  $\mathcal{M}_p(x_i)$  is then passed to the representation node  $P_p$  to produce representations

$$\mathcal{E}_p(x_i) = P_p(\mathcal{M}_p(x_i)) \quad (6)$$

for the current medical concept  $p$ .

The representations are then passed as input to the *outcome nodes*. The outcome nodes leverage the representations to predict these outcomes. It is quite important to note that the model can be optimized for multiple outcomes of interest *and* each representation node can have different and multiple outcomes, e.g. we can associate mortality prediction for certain medical concepts and readmission prediction for another (possibly overlapping) subset of medical concepts. For example, in Figure 2 (a), the medical concept  $P_1$  is associated with two outcomes  $O_1^1$  and  $O_1^2$  while the medical concept  $P_3$  is associated with the outcome  $O_3^1$ .

This model is trained in a multi-task setting allowing expert nodes and representation nodes to share learned representations and support outcome nodes in their prediction tasks. The loss function to be optimized is

$$\mathcal{L} = \underbrace{\lambda \frac{1}{N.P} \sum_i \sum_{\substack{p \\ x_i \text{ has } p}} (\mathcal{R}_p(\mathcal{E}_p(x_i)) - x_i)^2}_{\mathcal{L}_1} + \underbrace{\frac{1}{N.P} \sum_i \sum_{\substack{p \in \mathcal{C} \\ x_i \text{ has } p}} \sum_{\substack{o \\ o \text{ assoc. with } p}} \ell(O_p^o(\mathcal{E}_p(x_i)), y_i^o)}_{\mathcal{L}_2} \quad (7)$$

where  $\mathcal{R}_p$  is the reconstruction network that reconstructs the input  $x_i$  from the representation  $\mathcal{E}_p(x_i)$ ,  $\ell$  is the prediction loss function (e.g. cross entropy for classification) and  $y_i^o$  is the label for the outcome  $o$  for the visit  $i$ . The second sum in  $\mathcal{L}_1$  runs over all representation nodes (medical concepts) that  $x_i$  expresses forcing the representation from all nodes to be useful for the reconstruction of the original input. While the second sum in  $\mathcal{L}_2$  runs over all *core* nodes  $\mathcal{C}$  that  $x_i$  expresses, based on the user interest in specific phenotypes, the third sum runs over all outcome nodes  $o$  that are associated with the representation nodes  $p$ .

### 3.2 Hierarchical Multi-task Approach with OMTL

Optimizing Equation (7) can be used to solve the multi-task problem outlined above. However, depending on the choice of the phenotypes of interest, there can be a significant variation in the number of patients spanned by the representation nodes. For example, nodes corresponding to a higher level concept will typically span a much greater number of patients than the ones spanned by concepts at a leaf level of the ontology where concepts are very specialized. Interestingly, such nodes are semantically related via the medical concept ontology. Therefore, the generated representations across different medical concepts should share similar semantic relationships. OMTL aims to increase the information sharing among nodes using the ontology.

194 The learning architecture for OMTL is depicted in Figure 2 (b). Compared to the architecture shown  
 195 in Figure 2 (a), the representation nodes are now connected by mirroring the sub-graph  $\mathcal{G}$ . The  
 196 representation nodes without ancestors in  $\mathcal{G}$  can be computed using the same gating principle outlined  
 197 in Equation (5). However, for a node with ancestors (such as  $P_3$  in Figure 2 (b)), the model shares  
 198 both the representation from the expert nodes and its parents. Assuming that  $P_j$  has  $n$  parents,  
 199 the resultant representation  $\mathcal{E}_j$  at  $P_j$  is an aggregation of two vectors: (i)  $\mathcal{M}_j$  computed using the  
 200 combination of experts as given by Equation (5) and, (ii)  $\mathcal{P}_j$  computed using the combination of the  
 201 representations from the parents as:

$$\begin{aligned} H_j(x_i) &= \text{Softmax}(x_i \cdot W_j + b_j) \\ \mathcal{P}_j(x_i) &= \sum_{k=1}^n H_j^k(x_i) \mathcal{E}_k(x_i) \end{aligned} \quad (8)$$

202 The input to the node  $P_j$  will be the aggregation of  $\mathcal{M}_j$  and  $\mathcal{P}_j$ . The outcomes  $\mathcal{E}_j(x_i)$  can be  
 203 computed using the same procedure outlined in Equation (6). OMTL encourages information sharing  
 204 among the nodes using the graphical structure informed by the ontology. More specifically, the  
 205 ancestors of a representation node will be less affected by the ‘vanishing data’ problem and the  
 206 structure thereby acts as a forcing mechanism to help the network to learn the representation nodes  
 207 and outcome with higher efficacy. During the forward and backward pass of the training, data  $x_i$  is  
 208 only routed to representation nodes corresponding to phenotypes expressed in  $x_i$ .

## 209 4 Experiments

210 In this section, we present an experimental evaluation of OMTL focusing on the impact provided by  
 211 leveraging an ontology in a multi-task framework. We first describe the data used in the experiments  
 212 including the construction of patient cohorts. We then present the experimental setup before presenting  
 213 experimental results.

### 214 4.1 Data Description

215 We train and evaluate all models on data extracted from the publicly available MIMIC-III (Medical  
 216 Information Mart for Intensive Care) data set [10], a large, single-center database of patients admitted  
 217 to critical care units at a large tertiary hospital. While MIMIC-III includes several facets of patient  
 218 medical records, we focus our experiments on the vital signs, observations and diagnosis codes from  
 219 the cohort construction proposed in [7].

220 Since MIMIC-III tracks diagnosis codes in the ICD-9 format, only preserving a shallow two-level  
 221 hierarchical structure of medical concepts, we map all the observed ICD-9 diagnosis codes to  
 222 SNOMED CT codes, and construct a graph with a sufficiently deep hierarchy of medical concepts.  
 223 As a result, we obtain a complex knowledge graph where each node is associated with a set of ICU  
 224 stays. Next, we exploit a data augmentation algorithm<sup>2</sup> to extract cohorts for training our models as  
 225 follows. Given the knowledge graph and a few core nodes (e.g., specific SNOMED nodes identified  
 226 by domain experts for a certain phenotype), we aim to include more SNOMED nodes to augment  
 227 the ICU stays for analysis. The data augmentation algorithm starts from this set of nodes and grows  
 228 the set by choosing a series of predecessors using an MCMC sampler. We apply this node growth  
 229 for a predefined number of iterations to obtain connected graphs with sufficient nodes and stays for  
 230 training. In this manner, we obtain 3 data sets for different cohorts: liver disorder, heart disorder, and  
 231 kidney disorder for which our predictive outcomes are phenotypes “Disorders of lipid metabolism”,  
 232 “Acute myocardial infarction”, and “Chronic kidney disease”, respectively, and in-hospital mortality  
 233 prediction on both cohorts. The extracted graphs for all data sets are shown in Figure 4.

234 More statistical details on the selected cohorts can be found in Table 1. The corresponding data are  
 235 multivariate time series for each ICU stay. As typical with real-world data, these data contain missing  
 236 values. To address this issue, we replace these missing values with an average of carry-forward,  
 237 carry-backward and mean imputation methods. We then flatten the time series data to a single time  
 238 point using the mean values of each feature and each ICU stay. For each visit, we extracted 41  
 239 features for analysis.

<sup>2</sup>The details for the data augmentation algorithm used are described in the supplementary material.

Table 1: Statistics of datasets.

| Dataset          | # Core nodes | # Nodes | # Visits (+) |
|------------------|--------------|---------|--------------|
| Liver Mortality  | 5            | 26      | 4767 (837)   |
| Liver Pheno      | 5            | 26      | 9073 (1474)  |
| Heart Mortality  | 3            | 21      | 13067 (1867) |
| Heart Pheno      | 3            | 21      | 24135 (4039) |
| Kidney Mortality | 1            | 78      | 20316 (2749) |
| Kidney Pheno     | 1            | 78      | 39108 (5588) |

## 4.2 Experimental Setup

As mentioned in the Methods section, our goal is to share information among outcomes of interest at a few core nodes of interest  $\mathcal{C}$ . To realize the benefits of the augmented nodes without restricting the outcomes to be homegenous across the nodes, we train OMTL by optimizing for outcome loss at the core nodes along with reconstruction loss at all nodes. Under this setup, our primary objective is to test the hypothesis that adding the structure information from the ontology helps the model to produce more accurate predictions. To test this hypothesis, we propose to compare OMTL against the multi-task mixture of expert models [12] trained using the loss structure outline in Equation 7. This model referred as MMOE in the rest of this section, uses mixture of experts [9] to share information among outcomes without using any prior knowledge about the semantic outcome-relationship. In addition, we compare the proposed method to mixture of experts (MOE) [5] where mulitple experts were utilized but no gates, and the most commonly used shared-bottom (SB) multi-task deep learning structure [3] where there is only one expert.

OMTL augments MMOE with a gating mechanism  $\{H_j; \forall j\}$  enforcing the ontological structural relationship across representation nodes as described in the Methods section. To train OMTL, we first train it without considering the hierarchy to learn the experts and their corresponding gates  $G$ . We freeze these components and subsequently fine-tune for the outcomes and representations at the nodes by admitting the hierarchical structure.

**Implementation details:** Each module (expert, representation, reconstruction, outcome) is implemented as a one layer feed-forward neural network. We used 3 experts for MOE, MMOE, and OMTL. The size of the expert, reconstruction, representation modules are  $41 \times 5$ ,  $5 \times 41$ ,  $5 \times 5$ , respectively. The activation functions are Leaky ReLU for experts, Softplus for the representation layers, and ReLU for the reconstruction. Each expert has a dropout layer for regularization. The size of the gate  $G_p$  is  $41 \times 3$  and the size of the gate  $H_j$  is  $41 \times n_j$  where  $n_j$  is the number of parents for the representation node  $j$ . We performed hyper-parameter tuning and found good performance and stability across all models by setting the batch size to 64, learning rate to 0.001, dropout to 0.5,  $\lambda$  to 0.0001, and using Adam as the optimizer [11]. We evaluated the models using stratified 5-fold cross validation. the cross validation was done so as to ensure we have the same distribution of positive and negative examples in each node. All models were evaluated on exactly the same folds for fair comparison.

## 4.3 Multi-task Outcome Results

We compare all baselines with OMTL using AUC. The results are shown in Table 2 and Figure 3.

Table 2: AUC under ROC for all models.

| Dataset                  | SB           | MOE   | MMOE         | OMTL         |
|--------------------------|--------------|-------|--------------|--------------|
| Liver Disease Mortality  | 0.715        | 0.742 | 0.715        | <b>0.750</b> |
| Liver Disease Phenotype  | 0.572        | 0.592 | 0.576        | <b>0.626</b> |
| Heart Disease Mortality  | 0.723        | 0.700 | 0.713        | <b>0.770</b> |
| Heart Disease Phenotype  | 0.535        | 0.537 | 0.561        | <b>0.593</b> |
| Kidney Disease Mortality | <b>0.728</b> | 0.674 | 0.697        | <b>0.728</b> |
| Kidney Disease Phenotype | 0.743        | 0.736 | <b>0.749</b> | <b>0.749</b> |

Table 3: AUC under PRC for all models.

| Dataset                  | SB           | MOE   | MMOE         | OMTL         |
|--------------------------|--------------|-------|--------------|--------------|
| Liver Disease Mortality  | 0.414        | 0.448 | 0.429        | <b>0.478</b> |
| Liver Disease Phenotype  | 0.086        | 0.106 | 0.084        | <b>0.109</b> |
| Heart Disease Mortality  | 0.211        | 0.221 | 0.242        | <b>0.262</b> |
| Heart Disease Phenotype  | 0.180        | 0.182 | 0.202        | <b>0.214</b> |
| Kidney Disease Mortality | 0.252        | 0.219 | 0.233        | <b>0.267</b> |
| Kidney Disease Phenotype | <b>0.806</b> | 0.798 | <b>0.806</b> | <b>0.806</b> |

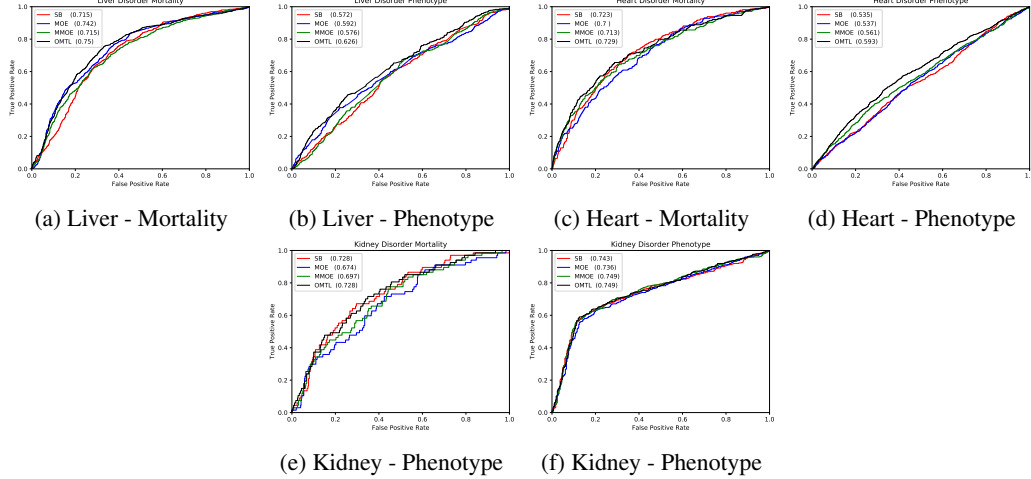


Figure 3: ROC for all experiments.

In general, we can see from Table 2 that the AUC gain that the structure is adding to MMOE ranges from 5%-9% for the liver disorder data with ontological graph shown in Figure 4 (a). It is worth noting that the core nodes in this graph are well-connected to each other. We hypothesize that the benefit of using the structural information among such well-connected core nodes helps OMTL to attain a significant boost in sharing information and hence provides more accurate predictions than MMOE.

The gain in the heart disorder graph, which is shown in Figure 4 (b), ranges from 2% to 6%. The AUC gain is not as significant as in the liver disorder case. A possible explanation for this is that the core nodes are less connected to each other. In such situation, the training may benefit by leveraging more outcome nodes. Similar conclusion are seen in the kidney graph. In the next section, we investigate this hypothesis.

Since the datasets are highly imbalanced, we also show the AUC under the precision recall curve (PRC) as shown in Figure 5 and Table 3. The OMTL model outperforms MMOE in the precision which indicates that it is able to predict the rare positive class more accurately than MMOE.

## 5 Conclusion

We proposed OMTL, a novel multi-task learning approach to analyze EHR data. OMTL models outcomes by incorporating existing ontology of phenotypical information in the multi-task setup. We have demonstrated the utility of this approach by comparing against state-of-the-art knowledge agnostic multi-task framework over a set of datasets. Our ongoing efforts are aimed at deploying the framework for large scale EHR analysis.



## References

- [1] Definition phenotype. <https://www.nature.com/scitable/definition/phenotype-phenotypes-35>, 2019. Accessed: 2019-05-22.
- [2] A. Ahmed, A. Das, and A. J. Smola. Scalable hierarchical multitask learning algorithms for conversion optimization in display advertising. In *Proceedings of the 7th ACM international conference on web search and data mining*, pages 153–162. ACM, 2014.
- [3] R. Caruana. Multitask learning. *Machine learning*, 28(1):41–75, 1997.
- [4] K. Donnelly. SNOMED-CT: The advanced terminology and coding system for e-health. *Studies in health technology and informatics*, 121:279, 2006.
- [5] D. Eigen, M. Ranzato, and I. Sutskever. Learning factored representations in a deep mixture of experts. *arXiv preprint arXiv:1312.4314*, 2013.
- [6] A. Elixhauser, C. Steiner, and L. Palmer. Clinical classifications software (CCS). *Book Clinical Classifications Software (CCS)(Editor ed^eds)*, 2014.
- [7] H. Harutyunyan, H. Khachatrian, D. C. Kale, and A. Galstyan. Multitask learning and benchmarking with clinical time series data. In *MLAH: Machine Learning for Health workshop, held at Advances in Neural Information Processing Systems*, 2017.
- [8] J. He and Y. Zhu. Hierarchical multi-task learning with application to wafer quality prediction. In *2012 IEEE 12th International Conference on Data Mining*, pages 290–298. IEEE, 2012.
- [9] R. A. Jacobs, M. I. Jordan, S. J. Nowlan, G. E. Hinton, et al. Adaptive mixtures of local experts. *Neural computation*, 3(1):79–87, 1991.
- [10] A. E. Johnson, T. J. Pollard, L. Shen, H. L. Li-wei, M. Feng, M. Ghassemi, B. Moody, P. Szolovits, L. A. Celi, and R. G. Mark. MIMIC-III, a freely accessible critical care database. *Scientific data*, 3:160035, 2016.
- [11] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- [12] J. Ma, Z. Zhao, X. Yi, J. Chen, L. Hong, and E. H. Chi. Modeling task relationships in multi-task learning with multi-gate mixture-of-experts. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pages 1930–1939. ACM, 2018.
- [13] E. Meyerson and R. Mäkelä. Beyond shared hierarchies: Deep multitask learning through soft layer ordering. In *6th International Conference on Learning Representations (ICLR)*, 2018.
- [14] R. Sanabria and F. Metze. Hierarchical multi task learning with CTC. In *Workshop On Spoken Language Technology, held at IEEE SLT*. IEEE, 2018.
- [15] V. Sanh, T. Wolf, and S. Ruder. A hierarchical multi-task approach for learning embeddings from semantic tasks. In *AAAI*, 2019.
- [16] J. Sun et. al. Combining knowledge and data driven insights for identifying risk factors using electronic health records. In *AMIA Annual Symposium Proceedings*, volume 2012, page 901. American Medical Informatics Association, 2012.
- [17] N. Tomašev, X. Glorot, J. W. Rae, M. Zielinski, H. Askham, A. Saraiva, A. Mottram, C. Meyer, S. Ravuri, I. Protsyuk, et al. A clinically applicable approach to continuous prediction of future acute kidney injury. *Nature*, 572(7767):116, 2019.
- [18] X. Wang, P. Kapanipathi, R. Musa, M. Yu, K. Talamadupula, I. Abdelaziz, M. Chang, A. Fokoue, B. Makni, N. Mattei, et al. Improving natural language inference using external knowledge in the science questions domain. *arXiv preprint arXiv:1809.05724*, 2018.
- [19] C. Xiao, T. Ma, A. B. Dieng, D. M. Blei, and F. Wang. Readmission prediction via deep contextual embedding of clinical concepts. *PloS one*, 13(4):e0195024, 2018.

- 338 [20] J. Zhang, K. Kowsari, J. H. Harrison, J. M. Lobo, and L. E. Barnes. Patient2vec: A personalized  
339 interpretable deep representation of the longitudinal electronic health record. *IEEE Access*, 6:  
340 65333–65346, 2018. ISSN 2169-3536. doi: 10.1109/ACCESS.2018.2875677.
- 341 [21] A. Zweig and D. Weinshall. Hierarchical multi-task learning: a cascade approach based on  
342 the notion of task relatedness. In *Theoretically Grounded Transfer Learning workshop, held at*  
343 *International Conference on Machine Learning (ICML)*, 2013.

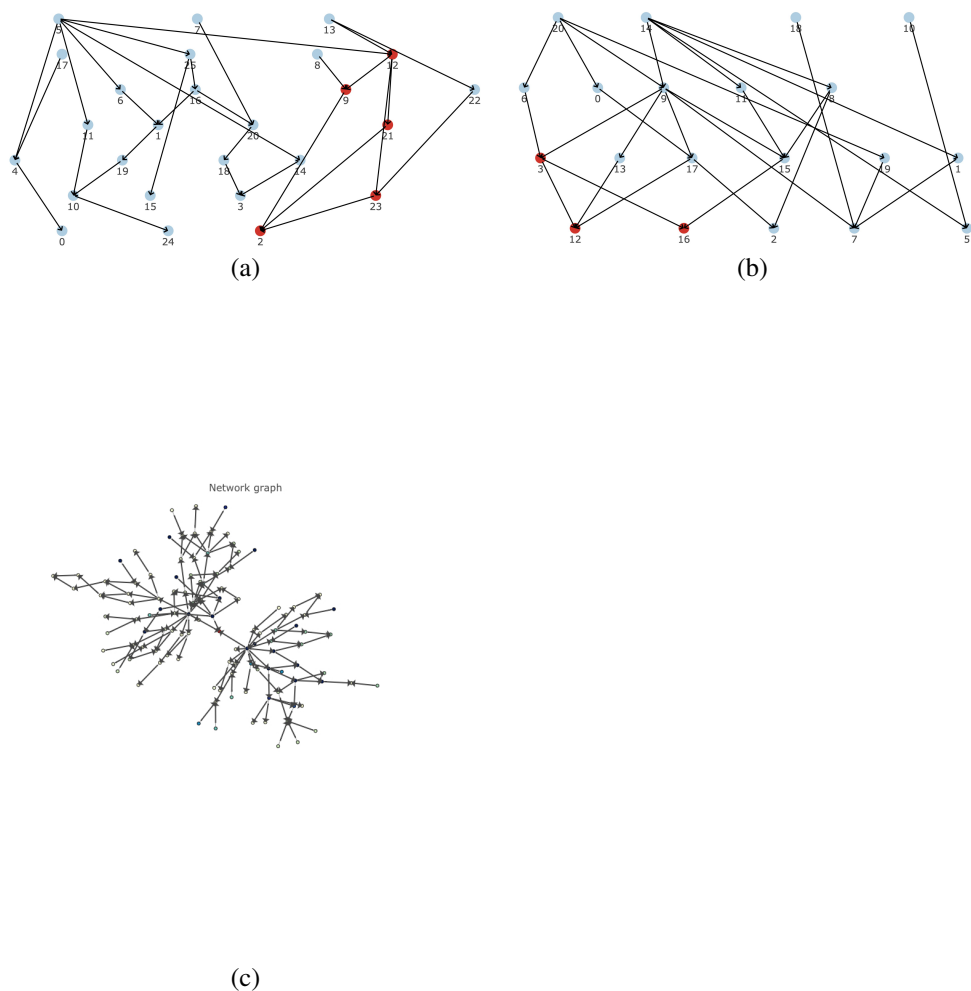


Figure 4: (a) Liver disorder graph. (b) Heart disorder graph. (c) Kidney disorder graph. Core nodes are shown in red, and augmented nodes are shown in blue. The labeled indices are used to describe the details of each node in the supplementary materials.

## 344 A Data

## 345 B Precision-recall curve

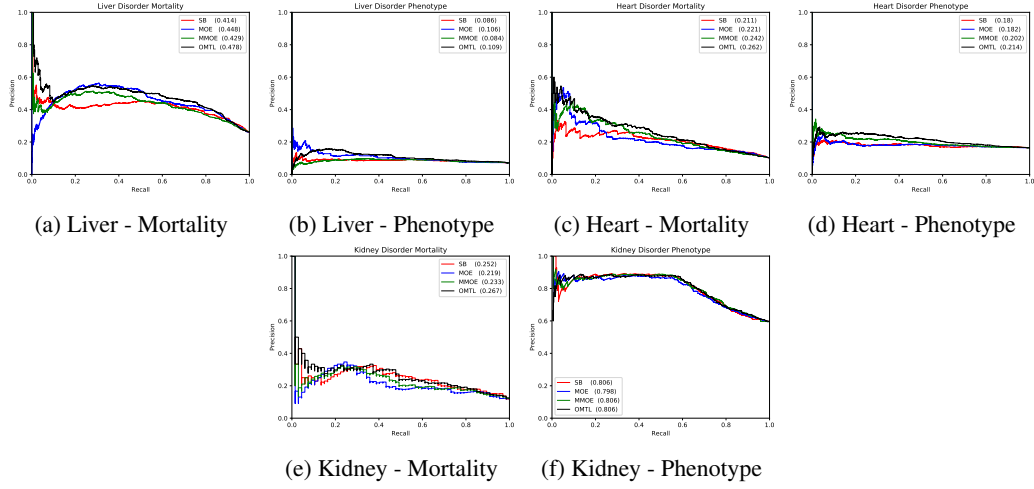


Figure 5: PRC for all experiments.