

Gene expression

# ECMarker: interpretable machine learning model identifies gene expression biomarkers predicting clinical outcomes and reveals molecular mechanisms of human disease in early stages

Ting Jin <sup>1,†</sup>, Nam D. Nguyen<sup>2,†</sup>, Flaminia Talos<sup>3,4</sup> and Daifeng Wang <sup>1,5,\*</sup>

<sup>1</sup>Department of Biostatistics and Medical Informatics, University of Wisconsin – Madison, Madison, WI 53706, USA, <sup>2</sup>Department of Computer Science, Stony Brook University, Stony Brook, NY 11794, USA, <sup>3</sup>Departments of Pathology and Urology, <sup>4</sup>Stony Brook Cancer Center, Stony Brook Medicine, Stony Brook, NY 11794, USA and <sup>5</sup>Waisman Center, University of Wisconsin – Madison, Madison, WI 53705, USA

\*To whom correspondence should be addressed.

<sup>†</sup>The authors wish it to be known that, in their opinion, the first two authors should be regarded as Joint First Authors.

Associate Editor: Pier Luigi Martelli

Received on November 12, 2019; revised on September 27, 2020; editorial decision on October 21, 2020; accepted on October 22, 2020

## Abstract

**Motivation:** Gene expression and regulation, a key molecular mechanism driving human disease development, remains elusive, especially at early stages. Integrating the increasing amount of population-level genomic data and understanding gene regulatory mechanisms in disease development are still challenging. Machine learning has emerged to solve this, but many machine learning methods were typically limited to building an accurate prediction model as a ‘black box’, barely providing biological and clinical interpretability from the box.

**Results:** To address these challenges, we developed an interpretable and scalable machine learning model, ECMarker, to predict gene expression biomarkers for disease phenotypes and simultaneously reveal underlying regulatory mechanisms. Particularly, ECMarker is built on the integration of semi- and discriminative-restricted Boltzmann machines, a neural network model for classification allowing lateral connections at the input gene layer. This interpretable model is scalable without needing any prior feature selection and enables directly modeling and prioritizing genes and revealing potential gene networks (from lateral connections) for the phenotypes. With application to the gene expression data of non-small-cell lung cancer patients, we found that ECMarker not only achieved a relatively high accuracy for predicting cancer stages but also identified the biomarker genes and gene networks implying the regulatory mechanisms in the lung cancer development. In addition, ECMarker demonstrates clinical interpretability as its prioritized biomarker genes can predict survival rates of early lung cancer patients ( $P$ -value < 0.005). Finally, we identified a number of drugs currently in clinical use for late stages or other cancers with effects on these early lung cancer biomarkers, suggesting potential novel candidates on early cancer medicine.

**Availability and implementation:** ECMarker is open source as a general-purpose tool at <https://github.com/daifengwanglab/ECMarker>.

**Contact:** daifeng.wang@wisc.edu

**Supplementary information:** [Supplementary data](#) are available at *Bioinformatics* online.

## 1 Introduction

Human disease development such as cancer is a complex, dynamic process that is fundamentally driven by abnormal molecular mechanisms. However, understanding the cancer mechanisms is still a challenging task, especially during the early cancer development

(Herbst *et al.*, 2018; Koeffler *et al.*, 1991). To this end, the tumor/node/metastasis (TNM) system has been widely used to characterize and classify the cancer development into various stages (Ludwig and Weinstein, 2005). The TNM stages were further associated with a number of individual molecular biomarkers and clinical outcomes such as survival rates (Ludwig and Weinstein, 2005). However, with

the advancement of whole genome sequencing of numerous human tumors, it became apparent that the molecular profiles of various tumor stages might not necessarily be reflected by the TNM system. Therefore, using systems biology approaches to identify the biomarkers that drive cancers from early to late stages could allow for better understanding of cancer mechanisms and offer new venues for the development of new preventive and therapeutic strategies.

However, there is a gap in understanding of the molecular biomarkers of early cancer and their underlying mechanisms at a system level. For example, the lung cancer, causing 27% of cancer-related deaths in the USA alone (Siegel *et al.*, 2018), is localized to the lung; neither lymph nodes nor other organs are believed to be affected at the early stage. As the cancer progresses to a more advanced stage, nearby lymph nodes and other organs may be affected (Frost *et al.*, 1984; Hu *et al.*, 2008). This pathological difference suggests that the underlying molecular mechanisms of the early and late stages are different. Also, if cancer is diagnosed at an early stage, such as a localized stage, the five-year survival rate is approximately 50%; this is mostly due to surgical interventions involving lung removal. After the localized stages, survival rates decrease rapidly as cases involving the lymph nodes or other metastatic sites necessitate elaborate treatment strategies (Hu *et al.*, 2008). Nearly 70% of patients with lung cancer present with locally advanced or metastatic disease at the time of diagnosis (Molina *et al.*, 2008). Thus, although a number of studies have indicated that early localized stages are easier to treat and have better survival rates, the underlying molecular mechanisms remain elusive. Here, we hypothesized that early cancers have different molecular wiring at a system level and that understanding this wiring could reveal new biomarkers and mechanisms of early cancer development. Thus, it is essential to identify the specific biomarkers of early cancer to understand the molecular mechanisms driving cancer development; this would enhance early cancer diagnosis and therefore improve survival rates.

Detecting early cancer biomarkers, however, involves the inherent challenges of relating the complex, multi-dimensional molecular processing that occurs in organs and tissues during early-stage cancer to observable clinical phenotypes in human patients. In particular, differential, temporal and spatial gene expression during early cancer result from disruptions in the complex, dynamic and multi-gene process that tightly regulates and controls the developmental integrity of organs and tissues. These temporal and spatial gene expression dynamics are fundamentally controlled by a variety of molecules called gene regulatory factors, including transcription factors (TFs) and non-coding RNAs. These factors cooperate in a gene regulatory network (GRN) to carry out correct developmental functions on a genome scale (Iyer *et al.*, 2017). The nodes of a GRN are genes, and the edges of a GRN connect regulatory factors to their target genes. Disruption of the cooperation between genes and regulatory factors in a GRN can give rise to abnormal gene expression, such as that which is present in diseases such as cancer. Therefore, a fundamental challenge for uncovering early cancer mechanisms is that of understanding the gene regulatory mechanisms, especially GRNs, controlling the changes in gene expression across cancer stages.

The collection of next-generation sequencing (NGS) data from large cohorts such as TCGA (Liu *et al.*, 2018a, b) provides measurements across multi-omics, including transcriptomics and epigenomics. This allows for studies of temporal dynamics in gene expression and regulation during cancer development and also for the systematic identification of stage-specific cancer biomarkers. Progress has been made in identification of some stage-specific molecular biomarkers of lung cancer, but systematic genome-wide analyses for identification of all potential early-stage biomarkers with predictive value for disease outcome are limited. For example, dysregulations in the epidermal growth factor receptor EGFR, associated with sensitivity of lung cancers to the tyrosine kinase inhibitor gefitinib (Iressa) (Pao *et al.*, 2004), echinoderm microtubule-associated protein-like 4 (EML4) and anaplastic lymphoma kinase (ALK) are frequently involved in oncogenic transformation (Lindeman *et al.*, 2013). In addition, v-raf murine sarcoma viral oncogene homologue B1 (BRAF) is a driver mutation gene in lung adenocarcinoma (Paik *et al.*, 2011). Although a challenging task, finding novel ways to

integrate the large-scale data provided by human tumors would enable the discovery of genome-wide early cancer biomarkers and underlying GRNs.

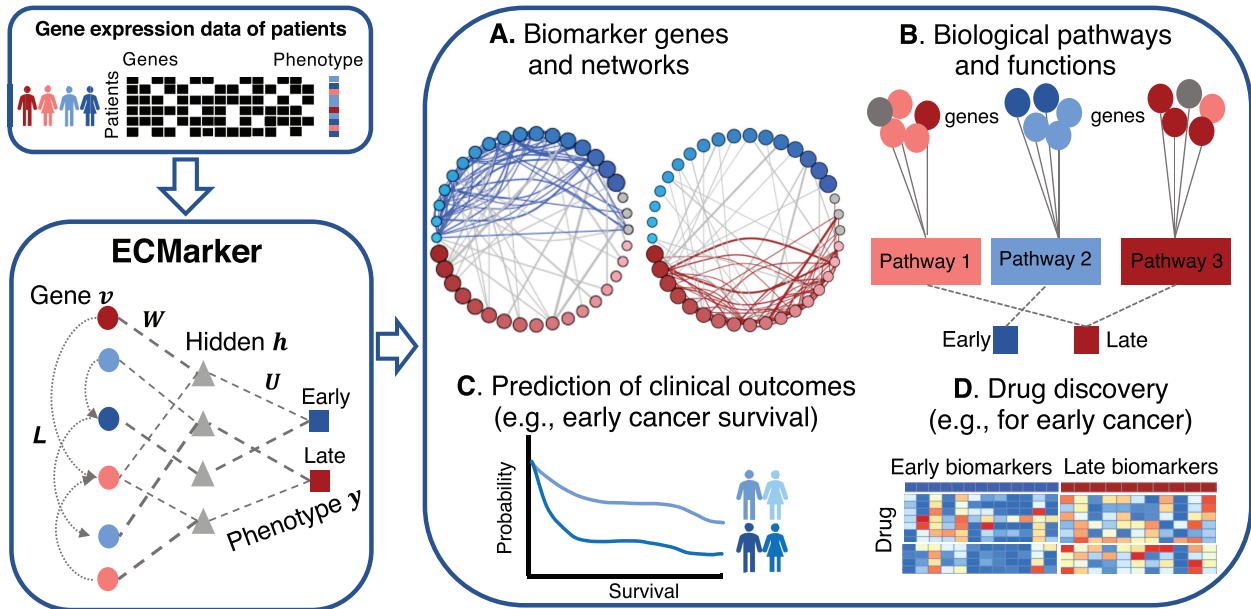
Traditionally, correlation-based models have been used to select biomarker genes involved in cancer development; e.g. 62 genes were uncovered in this way to distinguish between the early and late stages of clear cell renal cell carcinoma (ccRCC) (Jagga and Gupta, 2014; Rahimi and Gönen, 2018). However, correlation-based models only reveal linear relationships, whereas cancer development is a complex, non-linear process. Thus, machine learning has emerged as a powerful tool to predict biomarkers for various cancer features related to clinical presentation and staging; this tool has been found to be of great help in the diagnosis and treatment of various diseases (Libbrecht and Noble, 2015). For example, Statnikov *et al.*, (2008) applied random forests (RFs) and support vector machines (SVMs) to microarray data in order to aid in cancer diagnosis. Xiao *et al.*, (2018) constructed a multi-model ensemble approach to predict cancer in both normal conditions and tumor conditions. However, none of these studies revealed novel cancer mechanistic insights; these studies were limited to building an accurate classification model as a ‘black box’ but lacked any biological or clinical interpretability from the box. In addition, the biological datasets especially for genomics have the challenging of ‘curse of dimensionality’ (Clarke *et al.*, 2008); e.g. variables (e.g. genes) are much more than samples. To solve this, many machine learning methods applied prior feature selections to reduce the dimensionality, which however likely miss potentially important information at the system level.

To address these challenges, we designed a novel, interpretable machine learning approach, ECMarker, that can be used to discover gene expression biomarkers for the early disease stages, and simultaneously unravel the underlying molecular mechanisms in the ‘black box’ such as gene regulatory networks (GRNs). In particular, ECMarker is built on a neural network model, semi-restricted Boltzmann machine (SRBM) allowing lateral connections at the input gene layer for classifying disease phenotypes using population-level gene expression data. The SRBM model (Osindero and Hinton, 2007) has been used in non-biological contexts (e.g. computer vision, image classification) enabling modeling intra-connections among input variables (e.g. image patches). Based on the neural network connectivity, ECMarker further enables prioritizing genes and revealing the underlying gene network (from lateral connections) for predicting phenotypes. Compared to other methods, ECMarker is an interpretable model aiming to reveal underlying molecular mechanisms (e.g. gene regulation) while predicting phenotype. Many existing machine learning models still aim to learn a ‘black box’ with high accuracy, which is not straightforward to provide any biological insights as described above. ECMarker, instead, was designed to achieve all ‘interpretability’, ‘accuracy’ and ‘scalability’ via (i) using the lateral connections at the visible layer (i.e. genes) to reveal gene networks, (ii) simultaneously trying to achieve relatively high accuracy of classifying disease stages and (iii) inputting all genes and prioritizing genes by implicit feature selection. With applications to cancer genomic data, the prioritized genes and networks for early/late cancer stages revealed potential cancer stage-specific gene biomarkers and GRNs. Furthermore, we found the drugs that have significant effects on the ECMarker biomarkers for uncovering novel genomic medicine to early cancer development.

## 2 Materials and methods

### 2.1 ECMarker, an interpretable machine learning model to identify gene expression biomarkers and reveal underlying molecular mechanisms for disease phenotypes and clinical outcomes

ECMarker consists of three major components (Fig. 1) including (i) a neural network model, integrating the semi-restricted Boltzmann machine (SRBM) (Osindero and Hinton, 2007) with the Discriminative restricted Boltzmann machine (DRBM) (Larochelle and Bengio, 2008) for classifying disease phenotypes from the gene



**Fig. 1.** ECMarker, an interpretable machine learning framework for the identification of gene expression biomarkers of cancer stages and the prediction of clinical outcomes. ECMarker is a hierarchical neural network approach integrating semi- and discriminative-restricted Boltzmann machine models, to input the gene expression data of patients for predicting their disease phenotypes; e.g. early and late cancer stages. In particular, the ECMarker classification model consists of three layers: (1) the input gene layer  $v$ , (2) the hidden layer  $h$  and (3) the output phenotype layer  $y$ ; e.g. early versus late cancer stages. The lateral connections at the input gene layer enable identifying a gene network providing potential mechanistic insights for disease phenotypes. Thus, in addition to the phenotype prediction, ECMarker is also biologically and clinically interpretable for (A) identifying the gene expression biomarkers and gene networks for phenotypes (e.g. early and late stages); (B) revealing the associated biological functions and pathways for each phenotype; (C) predicting clinical outcomes such as survival rates, especially for early cancer patients; and (D) discovering novel drugs potentially affecting early cancer. Red and blue represent early and late stages, respectively

expression data at the population level; (ii) the prioritization of gene expression biomarkers for each phenotype using the integrated gradient method based on the neural network connectivity, and identification of a gene network using the lateral connections at the input layer; (iii) the functional and survival analyses of biomarker genes and networks for revealing underlying molecular mechanisms in the disease phenotypes (biological interpretability) and predicting clinical outcomes (clinical interpretability). We elaborated each component as follows.

## 2.2 The ECMarker classification model

The standard restricted Boltzmann machine (RBM) is an energy-based model that uses a layer of  $n$  hidden units to model a distribution over  $m$  visible units in the other layer (Hinton and Salakhutdinov, 2006). The connections between the two layers (i.e. visible layer to hidden layer) and all visible and hidden units form a bipartite graph. Thus, the connections within a layer (e.g. a visible unit to another visible unit) are prohibited. Also, the standard RBM typically takes binary values; i.e.  $v \in \{0, 1\}^m$  and  $h \in \{0, 1\}^n$ , and is often trained by the input distributions only. In the ECMarker, we have extended the RBM, based on the semi-restricted Boltzmann machine (SRBM) (Osindero and Hinton, 2007) and the Discriminative RBM (DRBM) (Larochelle and Bengio, 2008) for enabling (i) classification, (ii) inputting continuous values of visible units (e.g. gene expression) and (iii) modeling the gene relationships (i.e. a network) as follows. First, ECMarker inputs the expression profiles of  $m$  genes as  $m$  visible nodes  $v \in \mathbb{R}^m$ . Second, the hidden layer in the ECMarker consists of the binary variables  $h \in \{0, 1\}^n$ , where  $n$  is number of hidden nodes. Finally, the output layer  $y \in \{0, 1\}^K$  consists of all  $K$  phenotypes to predict. The improvements of the ECMarker classification model include:

1. To deal with the real-valued gene expression data, we replaced the binary visible units in the RBM by linear units with independent Gaussian noises. To simplify calculation, we used the

Gaussian noise  $\mathcal{N}(0, 1)$  and transformed the input data before training; i.e. standardizing features by removing the mean and scaling to unit variance.

2. We added an output layer  $y \in \{0, 1\}^K$  with discretized values modeling  $K$  phenotypes (e.g.  $K=2$ , early versus late disease stages) on the top of the hidden layer, and then used a joint distribution  $\prod_{(v,y) \in D_{\text{train}}} p(v, y)$  over the training dataset  $D_{\text{train}}$  of the input  $v \in \mathbb{R}^m$  (e.g.  $m$  gene expression values) and associated phenotype  $y$  for classification.
3. We allowed lateral connections among the visible units as SRBM did for modeling a network linking genes while predicting phenotypes.

In particular, the probability distribution represented by the ECMarker classification model with parameters  $\Theta$  is  $p(v, y, h | \Theta) \propto e^{-E(v, y, h; \Theta)}$ , where  $E(v, y, h; \Theta)$  is the energy function defined by

$$E(v, y, h; \Theta) = -b^T W v - a^T v - b^T h - c^T y - h^T U y - v^T L v \quad (1)$$

with that  $\Theta = \{W, a, b, c, U, L\}$  represents the model parameters. Note that the first three terms in the energy function are the same in the standard RBMs in which  $W \in \mathbb{R}^{n \times m}$  models the weight connections between  $m$  visible units and  $n$  hidden nodes,  $a \in \mathbb{R}^m$  is the bias of visible input units, and  $b \in \mathbb{R}^n$  is the bias of hidden nodes. The fourth and fifth terms model the contributions from the phenotype  $y$  in which  $U \in \mathbb{R}^{n \times K}$  models the weight connections between target and hidden layers, and  $c \in \mathbb{R}^K$  is the bias of target. The last term, a quadratic term of visible units  $v$  models two aspects contributing to the energy: (i) the lateral connections among genes where  $L \in \mathbb{R}^{m \times m}$  encodes the gene-gene relationships (i.e. adjacency matrix of a gene network) and (ii) the Gaussian units of gene expression inputs. We combined these two aspects in one term because they did not affect the calculation of log-likelihood gradient in training.

The conditional probability distribution of  $i$ th visible unit,  $v_i$  with real continuous value given the hidden units and other visible units (according to lateral connections among genes), is given by:

$$p(v_i|b, v_i; \Theta) = \mathcal{N}(v_i; a_i + b^T \mathbf{W}_{:i} + v^T \mathbf{L}_{:i}, 1), \quad (2)$$

where  $\mathbf{L}$  is a hollow matrix whose diagonal elements are all equal to zero, and  $\mathbf{W}_{:i}$ ,  $\mathbf{L}_{:i}$  are the  $i$ th columns of matrices  $\mathbf{W}$  and  $\mathbf{L}$  respectively.

The conditional probability distribution of the output units (of binary values) given the hidden units is as follows:

$$p(y|b; \Theta) = \sigma(c + b^T \mathbf{U}), \quad (3)$$

where  $\sigma(x) = (1 + e^{-x})^{-1}$  is the logistic sigmoid function.

The hidden units capture predictive information from both the visible inputs and the output classes. Thus, the conditional probability distribution of the hidden units (of binary values) given the visible inputs and output classes has the following form:

$$p(h|v, y; \Theta) = \sigma(b + \mathbf{W}v + \mathbf{U}y). \quad (4)$$

In training, because only  $v$  and  $y$  are observed, we calculated the marginal distribution represented by the model:

$$p(v, y; \Theta) \propto \sum_b e^{-E(v, y, b; \Theta)} = e^{-F(v, y; \Theta)},$$

where  $F(v, y; \Theta)$  is the free energy function defined by  $F(v, y; \Theta) = -\log \sum_b e^{-E(v, y, b; \Theta)}$ .

With the energy function aforementioned, the free energy can be further derived as follows:

$$F(v, y; \Theta) = -a^T v - c^T y - v^T \mathbf{L} v - \sum_{i=1}^n \log(1 + e^{b_i + \mathbf{W}_i v + U_i y}), \quad (5)$$

where  $\mathbf{W}_i$  and  $U_i$  are the  $i$ th rows of matrices  $\mathbf{W}$  and  $\mathbf{L}$  respectively. Furthermore, to address the ‘curse of dimensionality’ problem (i.e. features/genes are much more than samples) in genomic datasets, we trained our classification model with an  $\ell_1$  regularization. The  $\ell_1$  term,  $\mathbf{W}_1$ , which acts on the weight connections between the visible input layer and hidden layer, and also serves as an implicit feature selection method to automatically selecting prominent genes responsible for hidden units which models the distribution over visible input units (i.e. genes). This also enables us to add more hidden units to increase the learning capacity without being overfitted. Thus, we introduced the following loss function minimization during the model training:

$$\Theta^* \in \arg \min_{\Theta} -\log p(v, y; \Theta) + \mathbf{W}_1. \quad (6)$$

The data negative log-likelihood gradient is then:

$$-\frac{\partial}{\partial \Theta} \log p(v, y) = \frac{\partial}{\partial \Theta} F(v, y) - \sum_{v^-, y^-} p_{\Theta}(v^-, y^-) \frac{\partial}{\partial \Theta} F(v^-, y^-),$$

where  $v_j^-, y_j^- \sim p_{\Theta}(v^-, y^-)$  are generated examples from the current model’s distribution  $p_{\Theta}(v^-, y^-)$  and  $j$  is the patient index. These generated examples can be obtained by running a Markov chain to convergence using Gibbs sampling. However, in practice, the sampling process does not wait for convergence. Instead, the samples are obtained after 5-step Gibbs sampling in our case.

Finally, the learning procedure in the ECMarker is as follows (Algorithm 1). First, initialize  $\Theta_0$ ,  $y_0 \leftarrow y$  and  $v_0 \leftarrow v$ . At the learning iteration  $t$ , let  $\Theta_t$  be the model parameters. We generate  $v^-, y^- \sim q_{\Theta}$  using the Gibbs sampling. Then we update  $\Theta_{t+1} = \Theta_t + \eta_t \Delta(\Theta_t)$ , where  $\eta_t$  is the learning rate and >

$$\Delta(\Theta) \approx \sum_j \frac{\partial F(v_j, y_j)}{\partial \Theta} - \sum_j \frac{\partial F(v_j^-, y_j^-)}{\partial \Theta} + \frac{\partial \mathbf{W}_1}{\partial \Theta}. \quad (7)$$

The convergence of the algorithm toward the local optimum (since the loss function is non-convex) depends on the optimizer. We

---

#### Algorithm 1: ECMarker learning algorithm

---

**input** : training pairs  $(v, y)$   
**params**: Gibbs sampling step  $k$ , training steps  $T$ , learning rate  $\eta$   
**output** : parameters  $\Theta_{T+1}$  where  $\Theta = \{\mathbf{W}, \mathbf{a}, \mathbf{b}, \mathbf{c}, \mathbf{U}, \mathbf{L}\}$   
 initialize  $\Theta_0$ ;  $y_0 \leftarrow y$ ,  $v_0 \leftarrow v$ ;  
**for**  $t = 0 : k$  **do**  
    $h_k \sim p(h|y_k, v_k, \Theta_k)$  from equation (4);  
    $y_{k+1} \sim p(y|h_k, \Theta_k)$  from equation (3),  $v_{k+1} \sim p(v|h_k, v_k, \Theta_k)$  from equation (2);  
**end**  
 $\Theta_0 \leftarrow \Theta_k$ ;  
**for**  $t = 0 : T$  **do**  
   Update  $\Theta_{t+1} = \Theta_t + g(\Delta(\Theta_t), \eta, t)$  from equation (7) where  $g$  is an optimizer, e.g., SGD  
**end**

---

used the stochastic gradient descent (SGD) (Bottou, 2010) in ECMarker, which converges to a local minimum with an explicit convergence rate of  $O(1/T)$ , where  $T$  is the number of iterations. After training the model, the updated matrix  $L_{t+1} \in \Theta_{t+1}$  can be used as an adjacency matrix to construct a gene network, revealing the gene-gene relationships for predicting phenotypes (e.g. disease stages) and providing potential novel mechanistic insights.

### 2.3 Prioritization of the biomarker genes in ECMarker for phenotypes

Once the ECMarker classification model is trained, we further used a derivative-based method called integrated gradient for prioritizing input features (e.g. genes) (Sundararajan et al., 2017). In particular, we computed the gradient of model’s prediction with respect to each individual gene to show how the output response value (i.e. early versus late stages) changes with respect to a small change of input gene expression value. Hence, calculating these gradients for given input genes provide potential clues about which genes attribute the stage outcomes. This can be also interpreted to see which features are not selected due to  $\ell_1$  regularization since the gradients for these input genes are zeros. The output response value can be computed as the posterior class probability distribution given input  $v$  and has the following closed form:

$$p(y = y_k | v) = \frac{e^{-F(v, y_k)}}{e^{-F(v, y)}},$$

where  $F(v, y)$  is the free energy over all phenotypes in the output layer as in Equation (5), and  $F(v, y_k)$  is the free energy with regard to phenotype  $y = y_k$  ( $k = 1, \dots, K$ ), calculated as  $F(v, y_k) =$

$-a^T v - c_k y_k - v^T \mathbf{L} v - \sum_{i=1}^n \log(1 + e^{b_i + \mathbf{W}_i v + U_i y_k})$ . The exact gradient of this probability distribution can be calculated using the *autograd* package in PyTorch (Paszke et al., 2017). Furthermore, we define an importance score of each gene for the phenotype as the gradient of the gene to the phenotype. The higher positive scores that genes have, the more likely they contribute to predict the corresponded phenotype. Finally, given a phenotype, ECMarker prioritizes the genes for the phenotype via ranking its gene importance scores using the *Captum* package in PyTorch (Kokhlikyan, 2020).

### 2.4 Cancer gene expression and clinical datasets for building and testing ECMarker

We built and tested the ECMarker with the following publicly accessible gene expression datasets in lung cancer. The Gentles2015 dataset (Gentles et al., 2015) includes the log2-transformed gene expression data of 1103 non-small cell lung cancer (NSCLC) patients who had not received pre-biopsy treatment. We further imputed the missing values of the Gentles2015 dataset using the R package *impute* (Trevor Hastie, 2020), and then standardized the data per sample; e.g. a mean of zero and a standard deviation of 1. Also, we grouped the patients based on their TMN stages, with (I + IA + IB)

as the early stage ( $N=766$ ) and II, III and IV as the late stages ( $N=337$ ) and divided the dataset into balanced training and testing datasets via oversampling using the R package *ROSE* (Lunardon *et al.*, 2014). The lung adenocarcinoma (LUAD) and squamous cell carcinoma (LUSC) are two of the most common subtypes of NSCLC (Herbst *et al.*, 2018). To demonstrate the utility of ECMarker in distinct subtypes of NSCLC, we downloaded the RNA-seq gene expression datasets (FPKM values) for the LUAD and LUSC patients in TCGA (The Cancer Genome Atlas Research Network *et al.*, 2013), resulting in independent validation datasets TCGA-LUAD and TCGA-LUSC. Overall, we included 741 TCGA-LUAD patients ( $N=409$  and  $N=332$  in the early and late stages, respectively) and 758 TCGA-LUSC patients ( $N=380$  and  $N=378$  in the early and late stages, respectively) with (I + IA + IB) as the early stage and the rest as the late stage. A summary on the patient numbers of various stages is in Supplementary Table S1.

## 2.5 Identification of pathways and functions associated with early and late cancer stages via enrichment analysis

We performed the enrichment analyses of the genes in the ECMarker for revealing underlying molecular mechanisms from genes to disease stages. In particular, given a phenotype (e.g. early stage), we applied the Gene Set Enrichment Analysis (GSEA) (Subramanian *et al.*, 2005) by a R package, *fgsea* (Korotkevich *et al.*, 2019) to the ranking list of all input genes based on gene importance scores for the phenotype, and found the enriched terms including pathways, functions and oncogenic signatures from all eight gene sets of the Molecular Signatures Database (MSigDB) in the GSEA (Liberzon *et al.*, 2011, 2015). Using this enrichment analysis, we identified the enriched terms for both early and late lung cancer stages, providing mechanistic insights in lung cancer progress.

## 2.6 Survival analysis using ECMarker biomarker genes

We used the R function, *kmeans* to partition the early cancer patients into two groups using the gene expression data of top early biomarker genes in ECMarker ( $N=14$ ). The survival analyses and Kaplan–Meier plots were implemented using the R package, *survival* (Therneau, 2020).

## 2.7 Gene network analysis in ECMarker

The lateral connection weights (i.e.  $L$  matrix) from the ECMarker model how the gene-gene pairs (rather than individual genes) contribute to predict phenotypes, providing potentially additional mechanistic insights in terms of gene-gene relationships. Thus, using the  $L$  matrix as adjacency matrix, we further constructed a gene network from the ECMarker model for revealing potential gene regulatory relationships, especially on the transcription factors (TFs) to target genes (TGs). In addition, we compared the ECMarker gene network with the existing widely used methods such as GENIE3 (Huynh-Thu *et al.*, 2010) that only predict gene regulatory networks (TFs to TGs) from gene expression data, without simultaneously predicting phenotypes like ECMarker. In particular, we calculated the pairwise cosine distances between same genes; i.e. a distance of  $i$ th row vectors of  $L$  and  $G$  for Gene  $i$ ,

$$d_i(L_i, G_i) = 1 - \frac{L_i \cdot G_i}{\|L_i\| \|G_i\|}, \quad i = 1, \dots, m,$$

where  $G$  is the GENIE3's adjacency matrix. The cosine distance ranges from 0 meaning exactly the same to 2 meaning exactly opposite. Actually, the cosine distances have been widely used to measure the similarity of vectors and matrices, revealing the structural equivalence of the same vertices (i.e. same gene) between two networks via calculating weighted common neighbors divided by the geometric mean of their degrees (Salton, 1988). Furthermore, the cosine distance does not depend on the magnitude of the vector, which likely varies across different methods (e.g. ECMarker versus GENIE3) and thus is incomparable to each other. In contrast, the

other distance metrics such as the Euclidean distance or correlation take in account for the vector/matrix magnitudes, so are not chosen here for evaluating comparison. We used the Gentles2015 dataset for comparing ECMarker and GENIE3.

In addition to looking at the ECMarker network nodes and links, we also inferred the TF-TG relationships from the ECMarker network structures. In particular, we clustered the ECMarker gene network into a set of gene modules via hierarchical clustering (10 modules for the Gentles2015 dataset). The genes clustered together into a same module have strong connections for phenotype prediction, implying potential similar mechanisms such as co-regulation. Thus, we further identified the TFs with enriched target genes in each module (via TF binding sites on target gene regulatory regions) by *g:Profiler* (Raudvere *et al.*, 2019) and linked them to the modular target genes. Finally, we also calculated the centralities of the ECMarker gene network using *igraph* (Csárdi and Nepusz, 2006) and found the hub genes with high centrality (e.g. degree). Also, we applied ECMarker to the DREAM5 challenge data (Supporting Information). We found that ECMarker outperformed all other methods, including GENIE3, to infer the gene regulatory network in *Saccharomyces cerevisiae*. Since the *S.cerevisiae* network is the most complex in this DREAM5 challenge, this result demonstrates the ECMarker's high performance for predicting gene regulatory networks in complex biological systems.

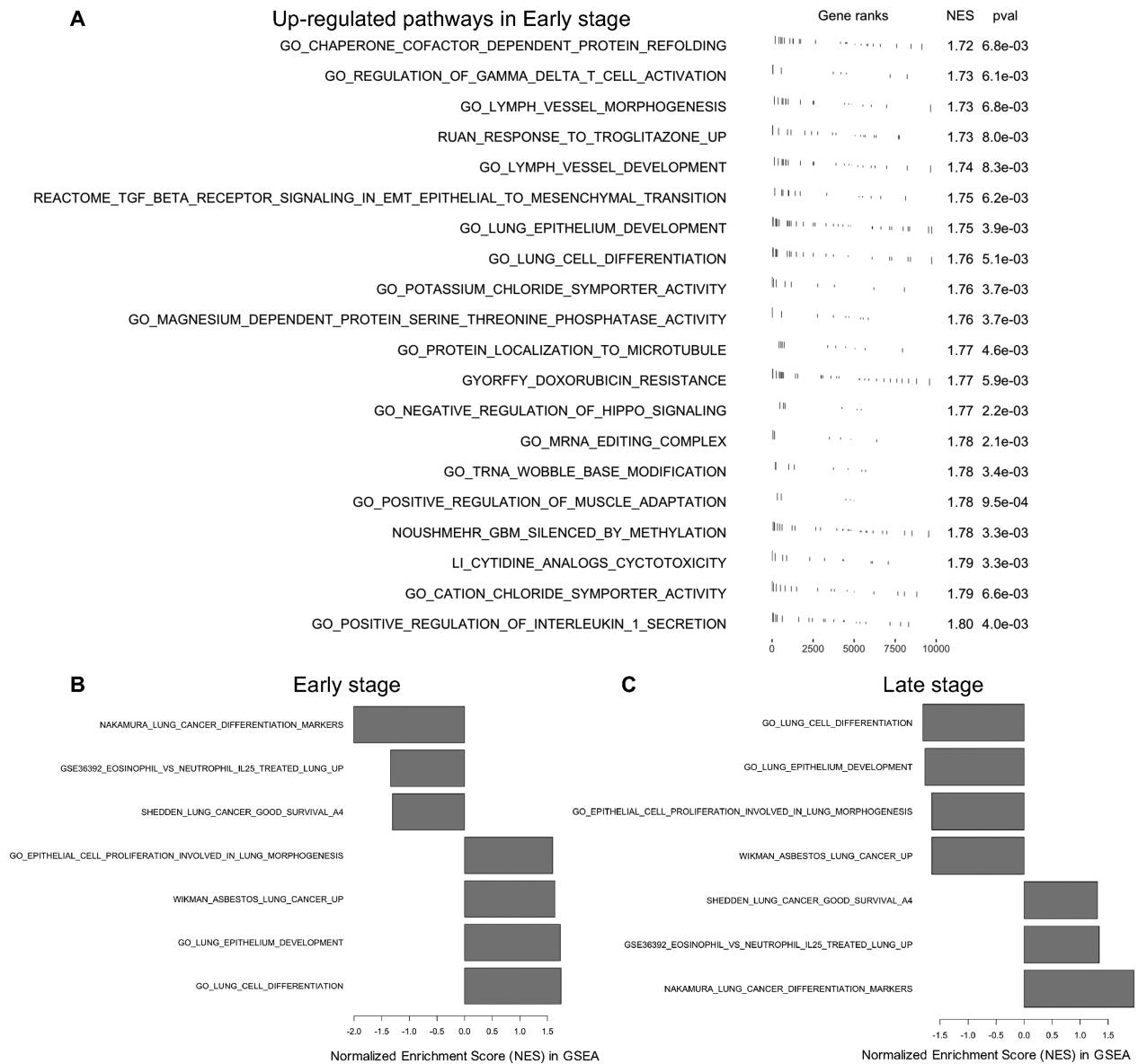
## 2.8 Identification of drugs targeting biomarker genes predicted by ECMarker

For discovering potential novel genome medicine using ECMarker, we identified the drugs targeting ECMarker biomarker genes of early and late cancer stages. In particular, we looked at a drug-gene database, GSCALite (Liu *et al.*, 2018a, b) that has calculated and summarized the  $z$ -scores of drug-gene pairs using the method in Rees *et al.* (2016) for revealing the mechanisms of action (MoA) of drugs to the target genes. The drugs with high  $z$ -scores imply potential causal mechanistic effects such as activation mechanisms and direct protein targets to the genes (Rees *et al.*, 2016). The  $z$ -scores were the Fisher's  $z$ -transformed correlation coefficients between gene expression and drug sensitivity (IC50 value) across all possible cancer cell lines in the GDSC database (Yang *et al.*, 2012), and thus removed potential biased effects from specific tissue types corresponding to the cell lines. Given a phenotype (e.g. early cancer stage), we found a number of drugs for its ECMarker biomarker genes with high  $z$ -scores (FDR < 0.05), suggesting their potential effects to the phenotype (e.g. early cancer drugs).

## 3 Results

### 3.1 Lung cancer stage biomarker genes by ECMarker reveal functions on cancer development and progress

We applied ECMarker to the Gentles2015 dataset (Section 2) for predicting the biomarker genes for lung cancer development and outcomes, especially for early-stage patients. In particular, we input the expression data of 10 102 genes from 766 early and 766 late patients (after balancing data) in the Gentles2015 dataset. After tuning hyperparameters in this ECMarker classification model, we had: (i) the input layer containing 10 102 genes; (ii) the hidden layer containing 9 hidden units; (iii) the output layer predicting early or late stage by a probability. Other hyperparameters were optimized as follows: train batch size = 50; learning rate = 0.1 and weight decay = 0.9 with the SGD method (Bottou, 2010);  $\ell_1$  penalty parameter: 0.1; number of training epoch: 1. We also performed  $k=10$  cross-validation and found that the model has the consistent relatively high balanced accuracy values with Mean = 0.74 and Variance = 0.001 compared to a baseline of 0.5 (for two phenotypes). Also, we applied another RBM-based model, elastic restricted Boltzmann machines (eRBMs) (Zhang *et al.*, 2017) that does not model lateral connections at the input layer, and found that its accuracy for predicting early and late stages is just around baseline of 0.5.



**Fig. 2.** Cancer-stage biomarker genes of ECMarker reveal the biological functions and pathways associated with lung cancer and cancer development. The enrichment analyses were accomplished by the GSEA for the ranked genes by gene importance scores for early or late lung cancer stages. (A) The upregulated functions and pathways significantly enriched in the early stage biomarker genes with  $P < 0.01$ . (B) and (C) Select functions and pathways associated with lung cancer and cancer development are up- or down-regulated at different stages

After training and cross-validating models, we used the average predictive model to further calculate the gene importance scores for both early and late lung cancer stages (Section 2), and prioritized the stage biomarker genes (i.e. high importance scores) in the lung cancer (Supplementary File S1). As shown on Figure 2A, a number of known lung, immunity and cancer related pathways, especially on cancer development are significantly enriched among top early biomarker genes after gene set enrichment analyses (Section 2); e.g. the epithelial mesenchymal transition (EMT,  $P < 6.2e-3$ ) (Lu and Kang, 2019), the gamma delta T cell activation ( $P < 6.1e-3$ ) (Pauza et al., 2018) and Interleukin-1 regulation ( $P < 4.0e-3$ ) (Lewis et al., 2006). Furthermore, top early and late genes are enriched with different upregulated pathways relating to lung cancer ( $P < 0.001$ , Fig. 2B and C); e.g. lung cell differentiation and epithelium development are upregulated for the early stage, but lung cancer survival and differential markers are upregulated for the late stage. All enriched terms for early lung cancer stage are available in Supplementary File S2.

Lung cancer is also heterogeneous; e.g. non-small cell lung cancer has two major subtypes: adenocarcinoma (LUAD) and squamous cell carcinoma (LUSC) (Lucchetta et al., 2019). To test ECMarker for predicting the cancer stages in the lung cancer subtypes, we also applied ECMarker to the TCGA-LUAD and TCGA-LUSC gene expression datasets (Section 2). To obtain an appropriate sample size for training ECMarker as the Gentles2015 dataset, we combined TCGA-LUAD and TCGA-LUSC together and trained one ECMarker model for classifying four phenotypes: the LUAD early and late stages and the LUSC early and late stages, aiming to reveal the specific early cancer mechanisms to lung cancer subtypes. After training and testing, the model achieved a high classification accuracy of 0.48 compared to a baseline of 0.27 (for four phenotypes). Furthermore, we found that top early-stage biomarker genes for LUAD and LUSC have significantly anti-correlated importance scores (Supplementary Fig. S1), suggesting potential distinct early cancer mechanisms across the lung cancer subtypes.

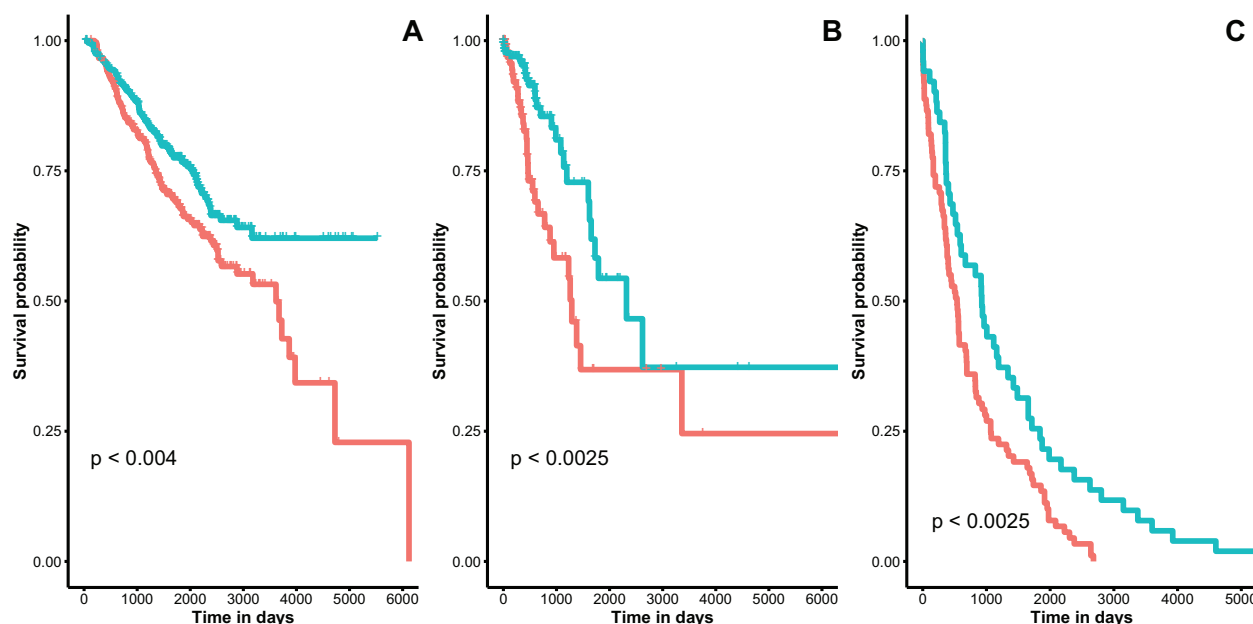


Fig. 3. ECMarker stage biomarker genes predict early cancer survival rates. Early lung cancer patients were clustered into two groups (represented by two curves on each panel) based on the top 14 ECMarker biomarker genes for early stage lung cancer. These biomarker genes were found using the Gentles2015 dataset (Gentles *et al.*, 2015). A Kaplan–Meier analysis showed that the early patient groups had significantly different survival rates ( $P < 0.005$ ) as shown in (A). In addition, application of these biomarker genes to an independent lung cancer cohorts, TCGA-LUAD and TCGA-LUSC (The Cancer Genome Atlas Research Network *et al.*, 2013), showed that the early-stage patients also had significantly different survival rates, as shown in (B) (TCGA-LUAD) and (C) (TCGA-LUSC) with  $P < 0.0025$ .

### 3.2 ECMarker biomarker genes predict clinical outcomes for early lung cancer

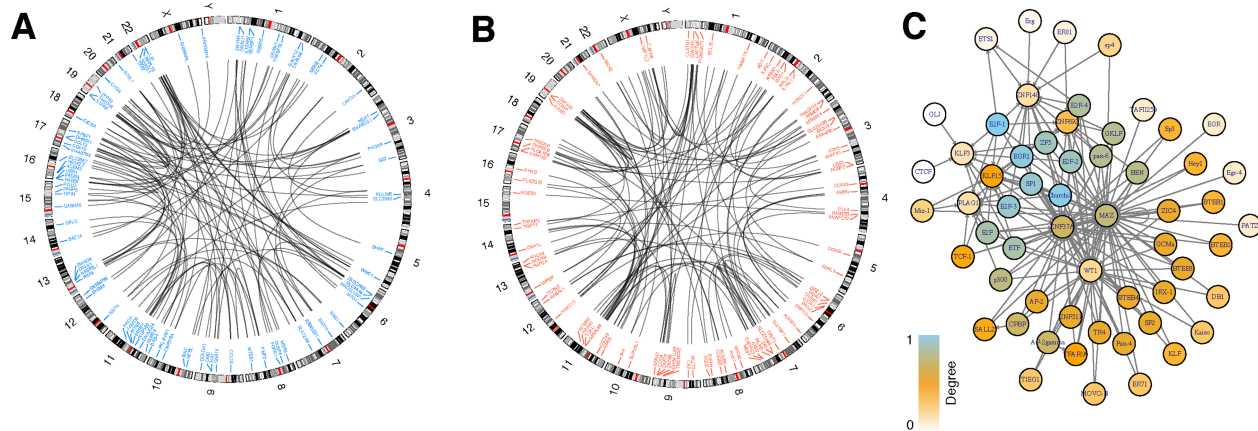
In addition to the genomic functions associated with cancer stages discovered by ECMarker, we also explored the relationships between stage biomarker genes and clinical outcomes of cancer patients. For example, we found that three lung cancer prognostic biomarker genes, CX3CR1, SLC15A2 and TFRC found by a recent multi-omics study (Haghjoo *et al.*, 2020) also have very high ECMarker importance scores for early stage (>80% genes). In addition, we found that the early lung cancer genes, SLC15A2 and TFRC are also hub genes (degree centrality in 1 and 10%) in the gene network revealed by ECMarker. To test the capability of ECMarker for predicting clinical outcomes for early cancer, we used top early-stage biomarkers learned by ECMarker (i.e. highest importance scores for early stage, Section 2) to partition the early cancer patients of Gentles2015 into two groups. We then found that two groups have significantly differential survival rates, suggesting that our ECMarker early biomarkers are able to predict early cancer survival rates ( $P < 0.004$ , Fig. 3A). Furthermore, we validated the top ECMarker early biomarker genes using early patients in two independent cohorts, TCGA-LUAD and TCGA-LUSC, and found that the early patients groups clustered by these biomarkers also have significantly differential survival rates ( $P < 0.0025$ , Fig. 3B and C). This demonstrates that our early biomarker genes have potential to predict survivals at the early cancer stage, suggesting the clinical interpretability of the ECMarker model.

### 3.3 Gene network in the ECMarker uncovers gene regulatory mechanisms in lung cancer

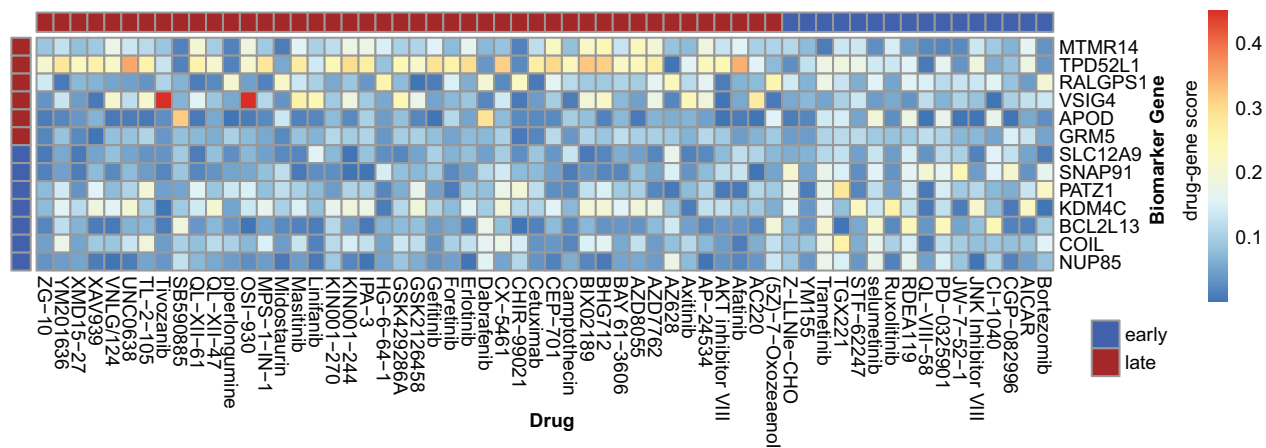
In addition to individual biomarker genes, we pursued the elucidation of the molecular mechanisms that drive the functional connectivity, especially in terms of gene regulation. Gene expression, though complex for phenotypes, is controlled by gene regulatory mechanisms. In particular, various regulatory factors, such as transcription factors (TFs), control the expression of biomarker genes to coordinate cancer phenotypes so they do not behave randomly (i.e. forming a GRN). ECMarker has the capability to model gene-gene interactions so that we could extract the weight values of lateral connections to describe the relationship of any pair of genes among

the input genes. A larger weight value indicated a stronger connection between genes. According to this standard, we are able to extract the learned GRN from any well-trained ECMarker classification model (Section 2). Specifically using the ECMarker gene network learnt from the Gentles2015 dataset for predicting lung cancer stages, we found the subnetworks linking top early and late biomarker genes are different (Fig. 4A and B), suggesting potential developmental regulatory mechanisms in the lung cancer progress.

We found that a number of known TFs related to cancer development and oncogenes involve in the ECMarker gene network (from top 1% links), including 46 epithelial-to-mesenchymal transition (EMT) signature genes (Byers *et al.*, 2013), lung cancer mutation genes (e.g. KRAS, BRAF, ALK, PIK3CA, AKT1, NRAS, EGFR, RET, ROS1) (Khoo *et al.*, 2015) and the genes in the frequency of alterations and signaling pathways of the lung cancer [53 out of 83 such genes (Li *et al.*, 2017)]. Moreover, a number of gene pairs from the top ECMarker links were also previously predicted to relate to the lung cancer gene regulation. For example, there are 2676 ECMarker top pairs (117 TF, 2348 TGs) also presented in a lung cancer regulatory network predicted by the one-class support vector machine (OC-SVM) model (Zhang *et al.*, 2018a, b). Also, a recent study has identified 10 oncogenic TFs and 11 tumor suppressing genes potentially required for NSCLC cell proliferation (Zhang *et al.*, 2018a, b). We found that all 10 such TFs and 10 out of 11 suppressing genes in our top ECMarker links (involved in 1017 pairs). In addition, a previous correlation-based analysis focusing on microRNA targets in lung cancer also predicted a set of TFs and target genes for the NSCLC (Mitra *et al.*, 2014), and was found to have 22 TFs, 101 target genes and 13 TF-TGs presenting in the top ECMarker links. Finally, we systematically compared the lung cancer networks of ECMarker and other computational methods such as GENIE3 that only use gene expression data to predict gene regulatory networks without integrating any phenotypic information (Section 2), especially on human diseases. Although ECMarker predicted a gene network specifically for predicting lung cancer development (i.e. stages), rather than generally for the lung cancer, there are still a variety of genes with high similarity between two networks ( $N = 1369$  genes, 13.6% with cosine distance  $< 0.35$ ). This shows a consistency between ECMarker and GENIE3 but also implies that



**Fig. 4.** The ECMarker gene network reveals the stage-specific gene connections and regulatory networks in the lung cancer. (A) The ECMarker gene network for early-stage biomarker genes (gene importance score  $> 0.001$  and top 100 links). (B) The ECMarker gene network for late-stage biomarker genes (gene importance score  $> 0.001$  and top 100 links). (C) A regulatory network among TFs inferred from the ECMarker gene network. The nodes are TF genes. The directed edges link TFs to the target genes (which are also TFs in this case). The node color corresponds to the degree centrality of TF in the network; e.g. skyblue represents high degree, and white means low degree. Note that all TFs shown here are predicted to be associated with the lung cancer development from the ECMarker model



**Fig. 5.** ECMarker biomarker genes discover potentially novel effective drugs for early lung cancer. The heatmap shows the effective scores of drugs to genes in terms of the MoAs of drugs to genes (Section 2) (Rees et al., 2016). The columns are the drugs with high MoAs to at least one of top 10 ECMarker stage biomarker genes. The rows are the genes from top 10 biomarker genes targeted by the drugs (blue: early stage, red: late stage)

rest of the genes without high network similarity potentially relate to the specific developmental functions and pathways in lung cancer such as immune response, cell proliferation and differentiation (Supplementary Fig. S2).

In addition to the ECMarker network nodes and links, using the network structures (e.g. gene modules), we also identified a list of TFs and TF-TG pairs for the lung cancer and cancer development (Section 2). In particular, we found a number of oncogenic TFs in our list, such as E2F genes, the cell-cycle TFs relating to tumor progression (Johnson and Schneider-Broussard, 1998), and EGR genes, the TFs regulating multiple tumor suppressors (Baron et al., 2006). Furthermore, we identified a number of potential master regulators in lung cancer development using our network. For example, E2F-1, a well-known transcription factor promoting the tumor progression for many cancer types including lung cancer (Engelmann and Putzer, 2012; Zhang et al., 2018a, b), plays a hub role (i.e. high degree) in the gene regulatory subnetwork in which target genes are TFs as well (Fig. 4C). Also, SP1, another known TF regulating lung cancer progression (Hsu et al., 2012) is in our network and also a hub gene. These hub genes in the ECMarker network imply that they regulate a number of lung cancer TFs as potential master regulators in lung cancer development. In addition, we found potential novel TFs for lung cancer development which were previously found to associate with other cancer types, such as WT1 for leukemia,

kidney and prostate cancers (Hastie, 2017) and MAZ, a MYC-associated zinc finger protein for pancreatic cancer (Maity et al., 2018). In addition, the target genes of some TFs are found to have significantly higher stage-specific importance scores than non-target genes (t-test  $P < 0.05$ ), suggesting that the cancer stage associated effects of the TFs; e.g. SP1 and AP-2 for late stage, TCF-1 and ER81 for early stage.

### 3.4 ECMarker biomarker genes link to potential novel drugs for early lung cancer

We further identified a number of drugs directly affecting ECMarker stage biomarker genes, aiming to provide potential novel candidates for early cancer medicine. Using the mechanisms of actions (MoAs) of drugs to genes (Section 2) (Rees et al., 2016), we identified a list of drugs for top 10 ECMarker early and late stage biomarker genes (Supplementary File S3). As shown on Figure 5, the drugs and stage biomarker genes can be in general clustered into early and late groups, suggesting the stage-specific drug effects on lung cancer development. Our analyses revealed that several known drugs for lung cancer also have high effects to our stage biomarker genes; e.g. the Type I RAF inhibitor—Dabrafenib and the Type II RAF inhibitor—AZ628 for the treatment of non-V600 BRAF mutant lung cancer (Noeparast et al., 2018) in the late stage group, and



YM155 for delaying the growth of NSCLC tumor xenografts is in the early stage group (Iwasa *et al.*, 2008).

Furthermore, several known drugs that were not originally used for lung cancer were predicted to have significant effects on our early-stage biomarkers; e.g. Bortezomib (Jones *et al.*, 2010), a proteasome inhibitor ameliorating breast cancer osteolytic disease, and AICAR inhibiting the cell growth in prostate cancer cells (Digregorio *et al.*, 2019). In addition, TAK1 inhibitor 5Z-7-oxozeanol for the treatment of cervical cancer is found to have potential effects on our late-stage lung cancer biomarkers (Guan *et al.*, 2017). In addition, we observed that a few drugs previously used for multi-cancer or in clinical use for the late cancer stages are in the early stage group, suggesting their potential effects to the early lung cancer; e.g. CI-1040 and PD-0325901 for advanced non-small cell lung, breast, colon or pancreatic cancers (Rinehart *et al.*, 2004). Another example is Ruxolitinib, a drug used during the Phase II study in the breast cancer (Stover *et al.*, 2018) and also possibly for NSCLC patients in all-stage to enhance oncolytic virotherapy (Patel *et al.*, 2019). Therefore, these drugs could potentially have effects on early cancer stages.

## 4 Discussion

ECMarker is an interpretable machine learning approach, built on the SRBM and DRBM for identifying gene expression biomarkers for disease phenotypes such as cancer stages. Beyond that a variety of machine learning methods typically pursuing the high prediction accuracy from genes to phenotypes (Supplementary Table S2), we demonstrated that the ECMarker model also has biological and clinically interpretabilities, in addition to high accuracy; e.g. it revealed the underlying regulatory mechanisms during lung cancer development and the stage biomarker genes predicted the survival rates of early cancer patients. Also, we showed that the drugs targeting the ECMarker biomarker genes are potential novel candidates for early cancer medicine. These biomarker genes comprise novel molecular candidates for early cancer diagnosis and detection, and the gene networks could potentially guide future experimental validations for early cancer mechanisms and treatments. Furthermore, ECMarker is scalable for inputting all possible genes and implicitly selecting biomarker genes for phenotypes via neural network regularization, and thus does not need any prior feature selections. Although this study applied ECMarker to lung cancer data specifically, ECMarker is a general-purpose method and can therefore be applied for other cancer types (Bailey *et al.*, 2018) and disease types such as neurodevelopmental and neurodegenerative diseases (De Jager *et al.*, 2018; Li *et al.*, 2018).

This study demonstrated that we are able to build the machine learning models that are biologically interpretable. This was our first round of attempts to address the lack of interpretability and translation of machine learning applications in biology and biomedicine. Given that cancer phenotypes are driven by a variety of multi-omic mechanisms (Bailey *et al.*, 2018), including transcriptomics, epigenomics, metabolomics, etc., multi-omic data integration and analyses for understanding cancer biology have been emerging (Rappoport and Shamir, 2019). Thus, we expect to develop advanced machine learning approaches that can reveal interactions across multi-omics relating to disease phenotypes in the near future; e.g. via multiview learning approaches (Nguyen and Wang, 2020). In particular, by integrating genotyping data, ECMarker can be extended to a deep hierarchical model, similar to deep neural network models (Wang *et al.*, 2018), to predict genotype-phenotype relationships and use intermediate biological connectivity and structures inside the model to reveal possible molecular mechanisms from genotype to phenotype.

The present study focused on gene expression data at the individual tissue level. However, the cancer tissues consist of different cell types with various fitness and mutational profiles (Saadatpour *et al.*, 2015). The continuous development of single-cell genomic and transcriptomic analyses for cancer research will enable us to explore how single cells contribute to cancer tissue expression and eventually affect phenotypes; e.g. single-cell deconvolution to estimate cell-

type fractions (Baron *et al.*, 2016; Wang *et al.*, 2019). Integrating single-cell data into the interpretable machine learning modeling and drug association analysis might uncover novel biological mechanisms and targetable key regulators at the cellular resolution for the advancement of precision cancer medicine.

## Acknowledgements

The authors thank Mr. Mufang Ying for helping Figure 4.

## Funding

This work was supported by the grants of National Institutes of Health, R01AG067025, R21CA237955 and U01MH116492 to Daifeng Wang and U54HD090256 to Waisman Center at UW-Madison.

*Conflict of Interest:* none declared.

## References

- Bailey, M.H. *et al.* (2018) Comprehensive characterization of cancer driver genes and mutations. *Cell*, **174**, 1034–1035.
- Baron, M. *et al.* (2016) A single-cell transcriptomic map of the human and mouse pancreas reveals inter- and intra-cell population structure. *Cell Syst.*, **3**, 346–360.e344.
- Baron, V. *et al.* (2006) The transcription factor Egr1 is a direct regulator of multiple tumor suppressors including TGFβ1, PTEN, p53, and fibronectin. *Cancer Gene Ther.*, **13**, 115–124.
- Bottou, L. (2010) *Large-Scale Machine Learning with Stochastic Gradient Descent*. Physica-Verlag HD, Heidelberg, pp. 177–186.
- Byers, L.A. *et al.* (2013) An epithelial-mesenchymal transition gene signature predicts resistance to EGFR and PI3K inhibitors and identifies Axl as a therapeutic target for overcoming EGFR inhibitor resistance. *Clin. Cancer Res.*, **19**, 279–290.
- Clarke, R. *et al.* (2008) The properties of high-dimensional data spaces: implications for exploring gene and protein expression data. *Nat. Rev. Cancer*, **8**, 37–49.
- Csárdi, G. and Nepusz, T. (2006) The igraph software package for complex network research. *InterJournal, Complex Systems*, **1695**, 1–9.
- De Jager, P.L. *et al.* (2018) A multi-omic atlas of the human frontal cortex for aging and Alzheimer's disease research. *Sci. Data*, **5**, 180142.
- Digregorio, M. *et al.* (2019) Relevance of translation initiation in diffuse glioma biology and its therapeutic potential. *Cells*, **8**, 1542.
- Engelmann, D. and Putzer, B.M. (2012) The dark side of E2F1: in transit beyond apoptosis. *Cancer Res.*, **72**, 571–575.
- Frost, J.K. *et al.* (1984) Early lung cancer detection: results of the initial (prevalence) radiologic and cytologic screening in the Johns Hopkins study. *Am. Rev. Respir. Dis.*, **130**, 549–554.
- Gentles, A.J. *et al.* (2015) Integrating tumor and stromal gene expression signatures with clinical indices for survival stratification of early-stage non-small cell lung cancer. *JNCI J. Natl. Cancer Inst.*, **107**, djv211.
- Guan, S. *et al.* (2017) TAK1 inhibitor 5Z-7-oxozeanol sensitizes cervical cancer to doxorubicin-induced apoptosis. *Oncotarget*, **8**, 33666–33675.
- Haghjoo, N. *et al.* (2020) Introducing a panel for early detection of lung adenocarcinoma by using data integration of genomics, epigenomics, transcriptomics and proteomics. *Exp. Mol. Pathol.*, **112**, 104360.
- Hastie, N.D. (2017) Wilms' tumour 1 (WT1) in development, homeostasis and disease. *Development*, **144**, 2862–2872.
- Herbst, R.S. *et al.* (2018) The biology and management of non-small cell lung cancer. *Nature*, **553**, 446–454.
- Hinton, G.E. and Salakhutdinov, R.R. (2006) Reducing the dimensionality of data with neural networks. *Science*, **313**, 504–507.
- Hsu, T.I. *et al.* (2012) Sp1 expression regulates lung tumor progression. *Oncogene*, **31**, 3973–3988.
- Hu, Z. *et al.* (2008) Genetic variants of miRNA sequences and non-small cell lung cancer survival. *J. Clin. Invest.*, **118**, 2600–2608.
- Huynh-Thu, V.A. *et al.* (2010) Inferring regulatory networks from expression data using tree-based methods. *PLoS One*, **5**, e12776.
- Iwasa, T. *et al.* (2008) Radiosensitizing effect of YM155, a novel small-molecule survivin suppressant, in non-small cell lung cancer cell lines. *Clin. Cancer Res.*, **14**, 6496–6504.
- Iyer, A.S. *et al.* (2017) Computational methods to dissect gene regulatory networks in cancer. *Curr. Opin. Syst. Biol.*, **2**, 115–122.

- Jagga,Z. and Gupta,D. (2014) Classification models for clear cell renal carcinoma stage progression, based on tumor RNAseq expression trained supervised machine learning algorithms. *BMC Proc.*, **8**, S2–S2.
- Johnson,D.G. and Schneider-Brossard,R. (1998) Role of E2F in cell cycle control and cancer. *Front. Biosci.*, **3**, d447–d448.
- Jones,M.D. et al. (2010) A proteasome inhibitor, bortezomib, inhibits breast cancer growth and reduces osteolysis by downregulating metastatic genes. *Clin. Cancer Res.*, **16**, 4978–4989.
- Khoo,C. et al. (2015) Molecular methods for somatic mutation testing in lung adenocarcinoma: EGFR and beyond. *Transl. Lung Cancer Res.*, **4**, 126–141.
- Koeffler,H.P. et al. (1991) Molecular mechanisms of cancer. *West. J. Med.*, **155**, 505–514.
- Kokhlikyan,N. et al. (2020) Captum: a unified and generic model interpretability library for pytorch. arXiv Preprint arXiv:2009.07896.
- Korotkevich,G. et al. (2019) Fast gene set enrichment analysis. bioRxiv. doi: 10.1101/060012.
- Larochelle,H. and Bengio,Y. (2008) Classification using discriminative restricted Boltzmann machines. In *Proceedings of the 25th International Conference on Machine Learning*. Association for Computing Machinery, Helsinki, Finland. pp. 536–543.
- Lewis,A.M. et al. (2006) Interleukin-1 and cancer progression: the emerging role of interleukin-1 receptor antagonist as a novel therapeutic agent in cancer treatment. *J. Transl. Med.*, **4**, 48.
- Li,M. et al.; BrainSpan Consortium. (2018) Integrative functional genomic analysis of human brain development and neuropsychiatric risks. *Science*, **362**, eaat7615.
- Li,Z. et al. (2017) The OncoPPI network of cancer-focused protein-protein interactions to inform biological insights and therapeutic strategies. *Nat. Commun.*, **8**, 14356.
- Libbrecht,M.W. and Noble,W.S. (2015) Machine learning applications in genetics and genomics. *Nat. Rev. Genet.*, **16**, 321–332.
- Liberzon,A. et al. (2015) The Molecular Signatures Database (MSigDB) hallmark gene set collection. *Cell Syst.*, **1**, 417–425.
- Liberzon,A. et al. (2011) Molecular signatures database (MSigDB) 3.0. *Bioinformatics*, **27**, 1739–1740.
- Lindeman,N.I. et al. (2013) Molecular testing guideline for selection of lung cancer patients for EGFR and ALK tyrosine kinase inhibitors: guideline from the College of American Pathologists, International Association for the Study of Lung Cancer, and Association for Molecular Pathology. *J. Thorac. Oncol.*, **8**, 823–859.
- Liu,C.J. et al. (2018a) GSCALite: a web server for gene set cancer analysis. *Bioinformatics*, **34**, 3771–3772.
- Liu,J. et al. (2018b) An integrated TCGA pan-cancer clinical data resource to drive high-quality survival outcome analytics. *Cell*, **173**, 400–416.e411.
- Lu,W. and Kang,Y. (2019) Epithelial-mesenchymal plasticity in cancer progression and metastasis. *Dev. Cell*, **49**, 361–374.
- Lucchetta,M. et al. (2019) Distinct signatures of lung cancer types: aberrant mucin O-glycosylation and compromised immune response. *BMC Cancer*, **19**, 824.
- Ludwig,J.A. and Weinstein,J.N. (2005) Biomarkers in cancer staging, prognosis and treatment selection. *Nat. Rev. Cancer*, **5**, 845–856.
- Lunardon,N. et al. (2014) ROSE: a package for binary imbalanced learning. *R. J.*, **6**, 79–89.
- Maity,G. et al. (2018) The MAZ transcription factor is a downstream target of the oncoprotein Cyt6/CCN1 and promotes pancreatic cancer cell invasion via CRAF-ERK signaling. *J. Biol. Chem.*, **293**, 4334–4349.
- Mitra,R. et al. (2014) Reproducible combinatorial regulatory networks elucidate novel oncogenic microRNAs in non-small cell lung cancer. *RNA*, **20**, 1356–1368.
- Molina,J.R. et al. (2008) Non-small cell lung cancer: epidemiology, risk factors, treatment, and survivorship. *Mayo Clin. Proc.*, **83**, 584–594.
- Nguyen,N.D. and Wang,D. (2020) Multiview learning for understanding functional multiomics. *PLoS Comput. Biol.*, **16**, e1007677.
- Noeparast,A. et al. (2018) Type II RAF inhibitor causes superior ERK pathway suppression compared to type I RAF inhibitor in cells expressing different BRAF mutant types recurrently found in lung cancer. *Oncotarget*, **9**, 16110–16123.
- Osindero,S. and Hinton,G. (2007) Modeling image patches with a directed hierarchy of Markov random fields. In *Proceedings of the 20th International Conference on Neural Information Processing Systems*. Curran Associates Inc., Vancouver, British Columbia, Canada. pp. 1121–1128.
- Paik,P.K. et al. (2011) Clinical characteristics of patients with lung adenocarcinomas harboring BRAF mutations. *J. Clin. Oncol.*, **29**, 2046–2051.
- Pao,W. et al. (2004) EGF receptor gene mutations are common in lung cancers from “never smokers” and are associated with sensitivity of tumors to gefitinib and erlotinib. *Proc. Natl. Acad. Sci. USA*, **101**, 13306–13311.
- Paszke,A. et al. (2017) Automatic Differentiation in PyTorch. 31st Conference on Neural Information Processing Systems (NIPS 2017), Long Beach, CA, USA
- Patel,M.R. et al. (2019) JAK/STAT inhibition with ruxolitinib enhances oncolytic virotherapy in non-small cell lung cancer models. *Cancer Gene Ther.*, **26**, 411–418.
- Paauw,C.D. et al. (2018) Gamma delta T cell therapy for cancer: it is good to be local. *Front Immunol*, **9**, 1305.
- Rahimi,A. and Gönen,M. (2018) Discriminating early- and late-stage cancers using multiple kernel learning on gene sets. *Bioinformatics*, **34**, i412–i421.
- Rappoport,N. and Shamir,R. (2019) Multi-omic and multi-view clustering algorithms: review and cancer benchmark. *Nucleic Acids Res.*, **47**, 1044–1044.
- Raudvere,U. et al. (2019) g:Profiler: a web server for functional enrichment analysis and conversions of gene lists (2019 update). *Nucleic Acids Res.*, **47**, W191–W198.
- Rees,M.G. et al. (2016) Correlating chemical sensitivity and basal gene expression reveals mechanism of action. *Nat. Chem. Biol.*, **12**, 109–116.
- Rinehart,J. et al. (2004) Multicenter phase II study of the oral MEK inhibitor, CI-1040, in patients with advanced non-small-cell lung, breast, colon, and pancreatic cancer. *J. Clin. Oncol.*, **22**, 4456–4462.
- Saadatpour,A. et al. (2015) Single-cell analysis in cancer genomics. *Trends Genet.*, **31**, 576–586.
- Salton,G. (1988) *Automatic Text Processing: The Transformation, Analysis, and Retrieval of Information by Computer*. Addison-Wesley, Reading, Mass.
- Siegel,R.L. et al. (2018) Cancer statistics, 2018. *CA Cancer J. Clin.*, **68**, 7–30.
- Statnikov,A. et al. (2008) A comprehensive comparison of random forests and support vector machines for microarray-based cancer classification. *BMC Bioinformatics*, **9**, 319.
- Stover,D.G. et al. (2018) Phase II study of ruxolitinib, a selective JAK1/2 inhibitor, in patients with metastatic triple-negative breast cancer. *NPJ Breast Cancer*, **4**, 10.
- Subramanian,A. et al. (2005) Gene set enrichment analysis: a knowledge-based approach for interpreting genome-wide expression profiles. *Proc. Natl. Acad. Sci. USA*, **102**, 15545–15550.
- Sundararajan,M. et al. (2017) Axiomatic attribution for deep networks. In *Proceedings of the 34th International Conference on Machine Learning – Volume 70*. JMLR.org, Sydney, NSW, Australia. pp. 3319–3328.
- The Cancer Genome Atlas Research Network et al. (2013) The Cancer Genome Atlas Pan-Cancer analysis project. *Nat. Genet.*, **45**, 1113–1120.
- Therneau,T.M. (2020) A Package for Survival Analysis in R. R package version 3.2-7. <https://CRAN.R-project.org/package=survival>, (28 September 2020, date last accessed).
- Trevor Hastie,R.T. et al. (2020) impute: impute: Imputation for microarray data. Bioconductor. R package version 1.62.0.
- Wang,D. et al.; PsychENCODE Consortium. (2018) Comprehensive functional genomic resource and integrative model for the human brain. *Science*, **362**, eaat8464.
- Wang,X. et al. (2019) Bulk tissue cell type deconvolution with multi-subject single-cell expression reference. *Nat. Commun.*, **10**, 380.
- Xiao,Y. et al. (2018) A deep learning-based multi-model ensemble method for cancer prediction. *Comput. Methods Programs Biomed.*, **153**, 1–9.
- Yang,W. et al. (2012) Genomics of Drug Sensitivity in Cancer (GDSC): a resource for therapeutic biomarker discovery in cancer cells. *Nucleic Acids Res.*, **41**, D955–D961.
- Zhang,D.L. et al. (2018a) Genome-wide identification of transcription factors that are critical to non-small cell lung cancer. *Cancer Lett.*, **434**, 132–143.
- Zhang,S. et al. (2018b) Landscape of transcriptional deregulation in lung cancer. *BMC Genomics*, **19**, 435.
- Zhang,S. et al. (2017) Elastic restricted Boltzmann machines for cancer data analysis. *Quant. Biol.*, **5**, 159–172.