1    Full title:

2    **Predicting cancer origins with a DNA methylation-based deep neural network**

3    **model**

4    Short title:

5    **DNN model for cancer origin prediction**

6

7    Authors:

8    Chunlei Zheng[1] and Rong Xu[1]*

9    [1]Department of Population and Quantitative Health Sciences, School of Medicine, Case

10   Western Reserve University, Cleveland, OH, USA

11

12   *Corresponding to: Rong Xu

13   Email: rxx@case.edu (RX)

14

15

16

17

18

19

20

# Abstract

Cancer origin determination combined with site-specific treatment of metastatic cancer patients is critical to improve patient outcomes. Existing pathology and gene expression-based techniques often have limited performance. In this study, we developed a deep neural network (DNN)-based classifier for cancer origin prediction using DNA methylation data of 7,339 patients of 18 different cancer origins from The Cancer Genome Atlas (TCGA). This DNN model was evaluated using four strategies: (1) when evaluated by 10-fold cross-validation, it achieved an overall specificity of 99.72% (95% CI 99.69%-99.75%) and sensitivity of 92.59% (95% CI 91.87%-93.30%); (2) when tested on hold-out testing data of 1,468 patients, the model had an overall specificity of 99.83% and sensitivity of 95.95%; (3) when tested on 143 metastasized cancer patients (12 cancer origins), the model achieved an overall specificity of 99.47% and sensitivity of 95.95%; and (4) when tested on an independent dataset of 581 samples (10 cancer origins), the model achieved overall specificity of 99.91% and sensitivity of 93.43%. Compared to existing pathology and gene expression-based techniques, the DNA methylation-based DNN classifier showed higher performance and had the unique advantage of easy implementation in clinical settings.

2

# Introduction

46    Identification of cancer origins is routinely performed in clinical practice as site-specific

47    treatments improve patient outcomes [1-4]. While some cancer origins are easy to be determined, others

48    are difficult, especially for metastatic and un-differentiated cancer. Cancer origin determination is

49    typically carried out with immunohistochemistry panels on the tumor specimen and imaging tests, which

50    need considerable resources, time, and expense. In addition, pathologic-based procedures have limited

51    accuracy (66-88%) in determining the origins of metastatic cancer [5-8].

52    Several gene expression- or microRNA-based molecular classifiers have been developed to

53    identify cancer origin. A k-nearest neighbor classifier based on 92 genes showed an accuracy of 84% in

54    identifying primary site of metastatic cancer via cross-validation [9]. Pathwork, a commercially available

55    platform based on similarity score of 1,550 genes between cancer tissue and reference tissue, achieved an

56    overall sensitivity of 88%, an overall specificity of 99% and an accuracy of 89% in identifying tissue of

57    origin [10, 11]. A decision-tree classifier based on 48 microRNA showed an accuracy of 85-89% in

58    identification of cancer primary sites [12, 13], and an updated version, the 64-microRNA based assay,

59    exhibited an overall sensitivity of 85% [14, 15]. A recent support vector machine-based classifier that

60    integrated gene expression and histopathology showed an accuracy of 88% in known origins of cancer

61    samples [16]. All these molecular platforms have shown better performance in identifying tissue of origin

62    as compared to pathology-based methods. However, gene expression- or microRNA-bases classifiers

63    need to handle RNA that is unstable and less convenient in clinic settings. In addition, these classifiers

64    have performance of <90% accuracy, which may further limit their wide adoption in clinical settings.

65    Hence, it is desirable to develop higher performance prediction tools for cancer origin determination,

66    which can also be easily implemented in clinical settings.

67    DNA methylation is a process by which methyl groups are added to the DNA molecule and 70-

68    80% of human genome is methylated [17]. It has been shown that DNA methylation is established in

69    tissue specific manner during development [18, 19]. Though the genomes of cancer patients exhibit

70    overall demethylation, tissue specific DNA methylation markers might be conserved [19]. Indeed, a

71    random forest-based cancer origin classifier using DNA methylation was reported to achieve a

72    performance with 88.6% precision and 97.7% recall in the validation set [20], which demonstrated the

73    usefulness of methylation data in cancer origin prediction. Recently, deep learning technologies have

74    rapidly applied to the biomedical field, including protein structure prediction, gene expression regulation,

75    behavior prediction, disease diagnosis and drug development [21, 22]. Studies show that deep learning-

76    based models often achieved higher performance than traditional machine learning methods (e.g. random

77    forest and support vector machine, etc.) in many settings, such as gene expression inference [23],

78    transcript factor binding prediction [24], protein-protein interaction prediction [25], detection of rare

79    disease-associated cell subsets [26], variant calling [27], clinic trial outcome prediction [28], among

80    others. In this study, we trained and robustly evaluated a high-performance cancer origin predictive model

81    by leveraging the large amount of DNA methylation data available in The Cancer Genome Atlas

82    (TCGA) and the recent developments in deep neural network learning techniques. We demonstrated that

83    our model performed better than traditional pathology- or gene expression-based models as well as

84    methylation-based random forest prediction model.

85

# Materials and methods

## Datasets

88    DNA methylation data (Illumina human methylation 450k BeadChip) and clinical information of

89    8,118 patients across 24 tissue types were obtained from in GDC data portal [29] using TCGAbiolink

90    (Bioconductor package, version 2.5.12) [30]. We excluded six tissue types with less than 100 cases in

91    TCGA to build robust cancer origin classifier. The final data include DNA methylation data and clinical

92    information from 7,339 patients of 18 cancer origins. TCGA data were used for both cancer origin

93    classifier training and evaluation, which were randomly and stratified split into training set (n=4,403),

94    development set (n=1,468) and test set (n=1,468) (Fig 1).

4

95 **Fig 1. Distribution of cancer samples in TCGA by tissue of origin.** A total of 7339 patients were

96 randomly and stratified split into train, dev and test sets according to 60:20:20.

97

98 In order to evaluate the classifier trained on TCGA dataset using independent data, we obtained

99 11 DNA methylation datasets (Illumina 450k platform) from Gene Expression Omnibus (GEO) [31]

100 using GEOquery (Bioconductor package, version 2.42.0) [32]. A total of 581 cancer patients covering 10

101 cancer origins were obtained and the information for each dataset was described in Table 1.

102

103 **Table 1. Characteristics of GEO datasets**

| GEO ID | Disease | Cancer origin | Cancer type | Num. of patients |
|---|---|---|---|---|
| GSE77871 | Adrenocortical carcinomas | Adrenal gland | Primary | 18 |
| GSE78751 | Triple negative breast cancer | Breast | Primary, metastatic | 23 12 |
| GSE101764 | Colorectal cancer | Colorectal | Primary | 112 |
| GSE38268 | Head and Neck Squamous Cell Carcinoma | Head and neck | Primary | 6 |
| GSE89852 | hepatocellular carcinomas | Liver | Primary | 37 |
| GSE49149 | Pancreatic cancer | Pancreas | Primary | 167 |
| GSE112047 | Prostate cancer | Prostate | Primary | 31 |
| GSE38240 | Prostate cancer | Prostate | Primary, metastatic | 2 6 |
| GSE73549 | Prostate cancer | Prostate | Metastatic | 18 |
| GSE86961 | Papillary thyroid cancer | Thyroid | Primary | 82 |
| GSE52955 | Urology cancer | Kidney, Bladder, prostate | Primary | 17, 25, 25 |

104

# Feature selection

106 Only the training data (n=4,403) from TCGA were used for feature selection. Currently, Illumina

107 450K and 27K are two commonly used platforms for genome wide analysis of DNA methylation, which

108 measure DNA methylation of around 450K and 27K CpG sites respectively. DNA methylation level of

5

109　CpG site is expressed as beta value using the ratio of intensities between methylated and unmethylated

110　alleles. Beta value is between 0 and 1 with 0 being unmethylated and 1 fully methylated. To make the

111　model with good compatibility and also reduce the dimensionality, we firstly reduced CpG sites to 27K

112　for 450K derived samples. To further remove the noise in the data, we used one-way analysis of variance

113　(one-way ANOVA) to filter the CpG sites whose beta values are not significantly different ($p > 0.01$)

114　among different tissues. Then we used the Tukey honest test to remove the CpG sites that maximal

115　differences of their beta values are less than 0.15.  The input features used for model building consisted of

116　DNA methylation from 10,360 CpG sites.

## Training a deep neural network (DNN) model for cancer origin

## classification

119　　　　We used DNA methylation data from training set (n=4,403) to build a DNN model to predict

120　cancer origins. Tensorflow [33], an open source framework to facilitate deep learning model training, was

121　used for this purpose. Four well-established techniques were used to optimize the training process,

122　including weight initialization by Xaiver method [34], Adam optimization [35], learning rate decay and

123　mini-batch training. Xaiver method can efficiently avoid gradient disappearance/explosion that random

124　initialization may bring. Adam, a combination of Stochastic Gradient Descent with momentum

125　descendent [36] and RMSprop [37], makes training process faster. Exponential learning decay (decay

126　every 1,000 steps with a base of 0.96) was used to improve model performance. Training was performed

127　in 128 mini-batch of 30 epochs to efficiently use the data.  In addition, three hyperparameters (learning

128　rate, number of hidden layer and hidden layer unit) were optimized to obtain best performance according

129　to development set performance (1,468 patients with the same distribution of cancer origins as training

130　set).

## Validating and testing DNN-based cancer origin prediction model

132　　　　We used four strategies to evaluate the performance of the DNN cancer origin classifier: (1)

6

133  evaluation in the 10-fold cross-validation in training dataset to obtain overall specificity, sensitivity, PPV

134  and NPV as well as corresponding confidence intervals of this model; (2) evaluation in the hold-out

135  testing dataset to obtain both the overall model performance and tissue-wise performance; (3) evaluation

136  in the subset of metastatic cancer samples nested in testing dataset to assess the performance of the model

137  in predicting the primary sites of metastatic cancer, which are often more difficult to be identified in

138  clinical practice and more clinically relevant; (4) evaluation in independent datasets from GEO to test the

139  robustness and generalizability of this DNN model. Metrics including specificity, sensitivity, positive

140  predictive value (PPV) and negative predictive value (NPV) were reported. Receiver Operating

141  Characteristic curve (ROC curve) was also calculated for each test data performance.

## Source code, data availability, and reproducibility

143  Source code used in this study is publicly available in a Github repository

144  (https://github.com/thunder001/Cancer_origin_prediction). We also shared a Jupyter Notebook to

145  replicate all the machine learning experiments from data processing, model building and optimization to

146  model evaluation. To execute this notebook, the environment needs to be firstly created according to a

147  YAML file available in Github. In addition, we also created a Docker image available in Docker hub

148  (https://hub.docker.com/r/thunder001/cancer_origin_prediction), where you can download it and run the

149  container directly on your computer.

150

# Results

## The overall performance of the DNN-based cancer origin classifier

## in 10-fold cross-validation setting

154  We used DNA methylation data of 7,339 patients from TCGA across 18 primary tissues to train

155  and test a DNN-based cancer origin classifier. The sample distribution in different cancer origins were

156   shown in Fig 1. The final DNN architecture consists of one input layer (10,360 neurons), two hidden

157   layers (64 neurons each layer) and one output layer (18 neurons) that represents 18 cancer origins (Fig 2).

158

159   **Figure 2.  Schematic representation of DNN architecture of cancer origin classifier.**

160

161         Evaluated in a 10-fold cross-validation setting, the model achieved an overall precision (positive

162   predictive value, PPV) of 0.9503 (95% CI:0.9373-0.9633) and recall (sensitivity) of 0.9259 (95%

163   CI:0.9187-0.9330) respectively. In addition, this model also achieved a high specificity of 0.9972 (95%

164   CI:0.9969-0.9975) (Table 2).

165

166   **Table 2.  DNN model performance using 10-fold cross validation of training data.**

| | Mean | SD | CI (95%) |
|---|---|---|---|
| **Specificity** | 0.9972 | 0.0001 | 0.9969, 0.9975 |
| **Sensitivity (Recall)** | 0.9259 | 0.0032 | 0.9187, 0.9330 |
| **PPV (Precision)** | 0.9503 | 0.0057 | 0.9373, 0.9633 |
| **NPV** | 0.9973 | 0.0001 | 0.9970, 0.9976 |

167
168         Note:  PPV: positive predictive value; NPV: negative predictive value.

169

# DNN-based cancer origin classifier shows high performance in

# testing dataset

172         We tested the classifier using test dataset, which includes 1,468 samples with similar distribution

173   with training set (Fig 1). Cancer origin classification and a confusion matrix for all samples were shown

174   in S1 and S2 Tables respectively. Model performance metrics were shown on Table 3. The specificity and

175   negative predictive value (NPV) in individual cancer origin prediction were consistently higher than 0.99.

176   The overall precision (PPV) and recall (sensitivity) reached 0.9608 and 0.9595 respectively. For many

177   cancer tissue origin predictions, including brain, colorectal, prostate, skin, testis, thymus and thyroid, this

8

178    DNN model achieved a precision of 100% (Table 3) and an average AUC of 0.99 (Fig 3).

179    **Table 3.  DNN model performance in test set.**

| CANER ORIGIN | SPECIFICITY | SENSITIVITY (RECALL) | PPV (PRECISION) | NPV |
|---|---|---|---|---|
| **AG** | 0.9993 | 0.9787 | 0.9787 | 0.9993 |
| **BLADDER** | 0.9986 | 0.9878 | 0.9759 | 0.9993 |
| **BRAIN** | 1.0000 | 1.0000 | 1.0000 | 1.0000 |
| **BREAST** | 0.9977 | 1.0000 | 0.9810 | 1.0000 |
| **COLORECTAL** | 1.0000 | 0.9861 | 1.0000 | 0.9993 |
| **ESOPHAGUS** | 0.9909 | 0.7410 | 0.7579 | 0.9902 |
| **HN** | 0.9971 | 0.9099 | 0.9619 | 0.9927 |
| **KIDNEY** | 0.9993 | 1.0000 | 0.9925 | 1.0000 |
| **LIVER** | 0.9993 | 0.9851 | 0.9851 | 0.9993 |
| **LUNG** | 0.9984 | 0.9740 | 0.9894 | 0.9961 |
| **PANCREAS** | 0.9979 | 1.0000 | 0.9167 | 1.0000 |
| **PROSTATE** | 1.0000 | 1.0000 | 1.0000 | 1.0000 |
| **SKIN** | 1.0000 | 1.0000 | 1.0000 | 1.0000 |
| **SOFT TISSUE** | 0.9993 | 0.9825 | 0.9825 | 0.9993 |
| **STOMACH** | 0.9921 | 0.9375 | 0.8721 | 0.9964 |
| **TESTIS** | 1.0000 | 1.0000 | 1.0000 | 1.0000 |
| **THYMUS** | 1.0000 | 0.8889 | 1.0000 | 0.9979 |
| **THYROID** | 1.0000 | 1.0000 | 1.0000 | 1.0000 |
| **OVERALL** | 0.9983 | 0.9595 | 0.9608 | 0.9983 |

180    Note:  PPV: positive predictive value; NPV: negative predictive value; AG: Adrenal Gland; HN: Head
181    and Neck
182

183    **Fig 3.  AUCs for individual cancer origin prediction in TCGA test set.**

184

185         There are some variations in precision and recall in different cancer origin predictions.  The

186    lowest performance occurred in esophagus origin prediction with a precision of 0.7579 and a recall of

187    0.7410. A total of 10 of 39 esophagus origins were incorrectly predicted as stomach origins (S1 and S2

188    Tables). Given that esophagus is a broad area, if a tumor is located at the border of stomach and

189    esophagus, it might be difficult for the classifier to distinguish these two tissues.  In addition, tissues from

190    adjacent regions may have similar methylation profiles so that the methylation-based prediction model

191    has difficulty in differentiating cancers with adjacent origins (e.g., esophagus vs stomach).

192

**193** **DNN-based cancer tissue classifier shows high performance in**

**194** **determining the origins of metastasized cancers**

**195** We evaluated the performance of the classifier in determining the origins of metastatic cancers

**196** that nested in our test data. Our data contained 701 samples from distantly metastasized cancers and 558

**197** of them have been used for model development. We then used remaining 143 samples from 12 cancer

**198** origins with various sample sizes for evaluation (Fig 4A). Cancer origin predictions and corresponding

**199** confusion matrix were shown in S3 and S4 Tables. Model performance metrics and ROC curves were

**200** shown in Table 4 and Fig 4B. Consistently, DNN model showed robust high performance in predicting

**201** metastatic cancer origins.

**202**

**203** **Fig 4. Performance of the DNN-based cancer origin classifier in metastatic cancer samples from**

**204** **TCGA test set.** (A) Distribution of metastatic cancer samples by tissue of origin. (B) AUCs for

**205** individual cancer origin prediction

**206** **Table 4. DNN model performance in metastatic cancer samples.**

| CANER ORIGIN | SPECIFICITY | SENSITIVITY (RECALL) | PPV (PRECISION) | NPV |
|---|---|---|---|---|
| ADRENAL GLAND | 1.0000 | 1.0000 | 1.0000 | 1.0000 |
| BLADDER | 1.0000 | 0.9643 | 1.0000 | 0.9914 |
| BREAST | 0.9929 | 1.0000 | 0.7500 | 1.0000 |
| COLORECTAL | 1.0000 | 1.0000 | 1.0000 | 1.0000 |
| ESOPHAGUS | 0.9504 | 1.0000 | 0.2222 | 1.0000 |
| HEAD AND NECK | 1.0000 | 0.8833 | 1.0000 | 0.9222 |
| KIDNEY | 1.0000 | 1.0000 | 1.0000 | 1.0000 |
| LIVER | 0.9929 | 1.0000 | 0.6667 | 1.0000 |
| LUNG | 1.0000 | 0.6667 | 1.0000 | 0.9929 |
| PANCREAS | 1.0000 | 1.0000 | 1.0000 | 1.0000 |
| STOMACH | 1.0000 | 1.0000 | 1.0000 | 1.0000 |
| THYROID | 1.0000 | 1.0000 | 1.0000 | 1.0000 |
| OVERALL | 0.9947 | 0.9595 | 0.8866 | 0.9922 |

**207** Note: PPV: positive predictive value; NPV: negative predictive value.

1

208      We noticed that performance metrics in several cancer origin predictions were poor:  a precision

209      of 0.22 for esophagus origin prediction, a precision of 0.67 for liver origin prediction and a recall of 0.67

210      for lung prediction.  The poor performance in these three cancer origin predictions may be due to small

211      sample size. As mentioned above, metastatic cancer samples comprise only a small subset of test dataset

212      in TCGA, the majority of which are primary tumors.  Only 2, 2 and 3 metastatic cancer samples from

213      esophagus, liver and lung origin respectively were included in test dataset (Fig 4A).  The classifier mis-

214      classified 6 out of 60 head and neck cancers as esophagus origin and 1 of 3 of lung cancers as liver

215      cancers (S4 Table).  Due to small sample sizes for esophagus, liver and lung cancers, a few mis-

216      classifications had significant impacts on the precision metrics.

217

## 218  DNN-based cancer tissue classifier shows high performance in

## 219  independent testing datasets

220      The DNN model was trained using DNA methylation data from TCGA. We then tested it in

221      independent datasets of 11 data series consisting of 581 tumor samples covering 10 tissue origins

222      downloaded from Gene Expression Omnibus (GEO). The sample distribution was shown in Fig 5A and

223      cancer origin predictions were listed in S5 Table. Evaluated using these independent datasets, the DNN

224      model achieved high performance with an overall precision and recall of 98.69% and 93.43% respectively

225      (Table 5). High performance was also achieved in individual cancer origin predictions (Table 5) with an

226      average AUC of 0.99 (Fig 5B). Importantly, the model achieved 100% accuracy in predicting the origins

227      of metastatic cancers in these datasets, including 24 prostate cancer that metastasized to bone, lymph node

228      or soft tissue and 12 breast cancer that metastasized to lymph node (see Table 1 for these samples).

229

230      **Fig 5.  Performance of the DNN-based cancer origin classifier in GEO dataset.**  (A) Distribution of

231      cancer samples obtained from GEO by tissue of origin. (B) AUCs for individual cancer origin prediction

232

**233**      **Table 5. DNN model performance using independent cancer samples (GEO)**

| CANER ORIGIN | SPECIFICITY | SENSITIVITY (RECALL) | PPV (PRECISION) | NPV |
|---|---|---|---|---|
| ADRENAL GLAND | 1.0000 | 0.7778 | 1.0000 | 0.9929 |
| BLADDER | 1.0000 | 1.0000 | 1.0000 | 1.0000 |
| BREAST | 0.9963 | 0.9714 | 0.9444 | 0.9982 |
| COLORECTAL | 1.0000 | 0.9643 | 1.0000 | 0.9915 |
| HEAD AND NECK | 1.0000 | 0.8333 | 1.0000 | 0.9983 |
| KIDNEY | 1.0000 | 1.0000 | 1.0000 | 1.0000 |
| LIVER | 0.9945 | 1.0000 | 0.9250 | 1.0000 |
| PANCREAS | 1.0000 | 0.8084 | 1.0000 | 0.9283 |
| PROSTATE | 1.0000 | 1.0000 | 1.0000 | 1.0000 |
| THYROID | 1.0000 | 0.9878 | 1.0000 | 0.9980 |
| OVERALL | 0.9991 | 0.9343 | 0.9869 | 0.9907 |

**234**      Note:  PPV: positive predictive value; NPV: negative predictive value.

**235**

# Discussion

**237**      We developed a deep neural network model to predict the cancer origins based on large amount

**238**   of DNA methylation data from 7,339 patients of 18 different cancer origins.  By combining DNA

**239**   methylation data with deep learning algorithm, our caner origin classifier achieved high performance as

**240**   demonstrated in four different evaluation settings.  Compared with Pathwork, a commercially available

**241**   cancer origin classifier based on gene expressions [10], our DNN model showed higher precision (95.03%

**242**   vs 89.4%) and recall (92.3% vs 87.8%) and comparable specificity (99.7% vs 99.4%). Compared with

**243**   DNA methylation-based random forest model, our DNN model achieved higher PPV (precision) (95.03%

**244**   in cross validation and 96.08% in test vs 88.6%) and comparable specificity, sensitivity and NPV.  In

**245**   addition, we showed that our DNN model is highly robust and generalizable as evaluated in an

**246**   independent testing dataset of 581 samples (10 cancer origins), with overall specificity of 99.91% and

**247**   sensitivity of 93.43%. Therefore, high performance both in primary and metastatic cancer origin

**248**   prediction and the potential for easy implementation in clinical setting make the methylation-based DNN

**249**   model a promising tool in determining cancer origins.

1

250    DNA methylation is established in tissue specific manner and conserved during cancer

251    development [19], which makes DNA methylation profile a very useful feature in cancer origin

252    prediction. Deep neural networks (DNNs) excels in capturing hierarchical features inherent in many

253    complicated biological mechanisms. Our study indicates that the trained DNN model may be able to

254    capture hierarchical patterns of cancer origins from the DNA methylation data.  While Interpretation of

255    deep learning-based models is a rapidly developing field and we expect that our model can be explained

256    in a meaningful way in the future.

257    Our DNN model has potential in predicting origins of Cancer of Unknown Primary origin (CUP).

258    CUP is a sub-group of heterogenous metastatic cancer with illusive primary site even after standard

259    pathological examination [38].  It is estimated that 3-5% metastatic cancers are CUP and the majority of

260    CUP patients (80%) have poor prognosis with overall survival of 6 -10 months [38].  Identifying primary

261    site of CUP poses challenges for treatment decisions in clinical practice. Currently, intensive pathologic

262    examination still leaves 30% of them unidentified [39, 40]. High performance of our DNA methylation-

263    based DNN model may provide an opportunity in this scenario when pathology-based approach fails.

264    However, due to the limited CUP data in both TCGA and GEO, we currently are unable to test the DNN

265    models in predicting the origins of CUP. Our future direction is to collaborate with hospital to collect

266    DNA methylation data from CUP patients to test our model. One challenge is to obtain the true primary

267    sites for these patients. Due to unknown property of CUP, true primary sites may be established in later

268    cancer development [20]. Another is through the post-mortem examination of patients since 75% of

269    primary sites of CUP were found in autopsy [41].

270    One limitation of this study is that small sizes of metastatic cancers in our data. Two resources of

271    metastatic cancer were used in this study: TCGA and GEO. TCGA has 701 metastatic cancer samples (12

272    tissues) with available methylation data from Illumina Human Methylation 450K platform. While the

273    model achieved an overall specificity of 99.47% and sensitivity of 95.95% in cross-validation using

274    TCGA data, we were unable to robustly test it using independent dataset since methylation data of

275    metastatic cancers is limited in GEO.  Further independent validation of our DNN-based model in

1

276 predicting origins of metastatic cancers, especially poorly differentiated or undifferentiated metastatic

277 cancer samples, is needed.

278

## References

280 1. Hainsworth JD, Rubin MS, Spigel DR, Boccia RV, Raby S, Quinn R, et al. Molecular gene

281 expression profiling to predict the tissue of origin and direct site-specific therapy in patients with

282 carcinoma of unknown primary site: a prospective trial of the Sarah Cannon research institute. J Clin

283 Oncol 2013. 10;31:217-23.

284 2. Varadhachary GR, Raber MN, Matamoros A, Abbruzzese JL. Carcinoma of unknown primary with a

285 colon-cancer profile-changing paradigm and emerging definitions. Lancet Oncol. 2008;9:596–9.

286 3. Varadhachary GR, Spector Y, Abbruzzese JL, Rosenwald S, Wang H, Aharonov R, et al. Prospective

287 gene signature study using microRNA to identify the tissue of origin in patients with carcinoma of

288 unknown primary. Clin Cancer Res. 2011;17:4063-70.

289 4. Varadhachary GR, Karanth S, Qiao W, Carlson HR, Raber MN, Hainsworth JD, et al. Carcinoma of

290 unknown primary with gastrointestinal profile: immunohistochemistry and survival data for this

291 favorable subset. Int J Clin Oncol. 2014;19:479-84.

292 5. Brown RW, Campagna LB, Dunn JK, Cagle PT. Immunohistochemical identification of tumor

293 markers in metastatic adenocarcinoma. A diagnostic adjunct in the determination of primary site. Am

294 J Clin Pathol. 1997;107:12-9.

295 6. DeYoung BR, Wick MR. Immunohistologic evaluation of metastatic carcinomas of unknown origin:

296 an algorithmic approach. Semin Diagn Pathol. 2000;17:184-93.

297 7. Dennis JL, Hvidsten TR, Wit EC, Komorowski J, Bell AK, Downie I, et al. Markers of

298 adenocarcinoma characteristic of the site of origin: development of a diagnostic algorithm. Clin

299 Cancer Res. 2005;11:3766-72.

300 8. Park SY, Kim BH, Kim JH, Lee S, Kang GH. Panels of immunohistochemical markers help

301    determine primary sites of metastatic adenocarcinoma. Arch Pathol Lab Med. 2007;131:1561-7

302    9.  Ma XJ, Patel R, Wang X, Salunga R, Murage J, Desai R, et al. Molecular classification of human

303        cancers using a 92-gene real-time quantitative polymerase chain reaction assay. Arch Pathol Lab

304        Med. 2006;130:465–73.

305    10. Monzon FA, Lyons-Weiler M, Buturovic LJ, Rigl CT, Henner WD, Sciulli C, et al. Multicenter

306        validation of a 1,550-gene expression profile for identification of tumor tissue of origin. J Clin Oncol.

307        2009;27:2503–8.

308    11. Pillai R, Deeter R, Rigl CT, Nystrom JS, Miller MH, Buturovic L, et al. Validation and

309        reproducibility of a microarray-based gene expression test for tumor identification in formalin-fixed,

310        paraffin-embedded specimens. J Mol Diagn. 2011;13:48-56.

311    12. Rosenfeld N, Aharonov R, Meiri E, Rosenwald S, Spector Y, Zepeniuk M, et al. MicroRNAs

312        accurately identify cancer tissue origin. Nat Biotechnol. 2008;26:462–9.

313    13. Rosenwald S, Gilad S, Benjamin S, Lebanony D, Dromi N, Faerman A, et al. Validation of a

314        microRNA-based qRT-PCR test for accurate identification of tumor tissue origin. Mod Pathol

315        2010;23:814–23.

316    14. Meiri E, Mueller WC, Rosenwald S, Zepeniuk M, Klinke E, Edmonston TB, et al. A second-

317        generation microRNA-based assay for diagnosing tumor tissue origin. Oncologist. 2012;17:801–12

318    15. Pentheroudakis G, Pavlidis N, Fountzilas G, Krikelis D, Goussia A, Stoyianni A, et al. Novel

319        microRNA-based assay demonstrates 92% agreement with diagnosis based on clinicopathologic and

320        management data in a cohort of patients with carcinoma of unknown primary. Mol Cancer.

321        2013;12:57.

322    16. Tothill RW, Shi F, Paiman L, Bedo J, Kowalczyk A, Mileshkin L, et al. Development and validation

323        of a gene expression tumour classifier for cancer of unknown primary. Pathology. 2015;47:7–12.

324    17. Kulis M, Esteller M. DNA methylation and cancer. Adv Genet. 2010;70:27–56.

325    18. Ohgane J, Yagi S, Shiota K. Epigenetics: the DNA methylation profile of tissue-dependent and

326        differentially methylated regions in cells. Placenta. 2008;29 Suppl A:S29–35.

327    19. Fernandez AF, Assenov Y, Martin-Subero JI, Balint B, Siebert R, Taniguchi H, et al. A DNA

328       methylation fingerprint of 1628 human samples. Genome Res. A DNA methylation fingerprint of

329       1628 human samples. Genome Res. 2012;22:407–19.

330    20. Moran S, Martínez-Cardús A, Sayols S, Musulén E, Balañá C, Estival-Gonzalez A, et al. Epigenetic

331       profiling to classify cancer of unknown primary: a multicentre, retrospective analysis. Lancet Oncol.

332       2016;17:1386–1395.

333    21. Min S, Lee B, Yoon S. Deep learning in bioinformatics. Brief Bioinform. 2017;18:851–869

334    22. Ching T, Himmelstein DS, Beaulieu-Jones BK, Kalinin AA, Do BT, Way GP, et al. Opportunities

335       and obstacles for deep learning in biology and medicine. J R Soc Interface. 2018;15(141). doi:

336       10.1098/rsif.2017.0387

337    23. Chen Y, Li Y, Narayan R, Subramanian A, Xie X. Gene expression inference with deep learning.

338       Bioinformatics. 2016;32(12):1832-9.

339    24. Alipanahi B, Delong A, Weirauch MT, Frey BJ. Predicting the sequence specificities of DNA- and

340       RNA-binding proteins by deep learning. Nat Biotechnol. 2015;33(8):831-8

341    25. Du T, Liao L, Wu CH, Sun B. Prediction of residue-residue contact matrix for protein-protein

342       interaction with Fisher score features and deep learning. Methods. 2016;110:97-105

343    26. Arvaniti E, Claassen M. Sensitive detection of rare disease-associated cell subsets via representation

344       learning. Nat Commun. 2017;8:14825. doi: 10.1038/ncomms14825.

345    27. Poplin R, Chang PC, Alexander D, Schwartz S, Colthurst T, Ku A, et al. A universal SNP and small-

346       indel variant caller using deep neural networks. Nat Biotechnol. 2018;36(10):983-987

347    28. Artemov AV, Putin E, Vanhaelen Q, Aliper A, Ozerov IV, Zhavoronkov A, et al. Integrated deep

348       learned transcriptomic and structure-based predictor of clinical trials outcomes. BioRxiv [Preprint].

349       2016. (doi:10.1101/095653)

350    29. GDC data portal. https://portal.gdc.cancer.gov. Accessed 7 August 2019

351    30. Colaprico A, Silva TC, Olsen C, Garofano L, Cava C, Garolini D, et al. TCGAbiolinks: a

352       R/Bioconductor package for integrative analysis of TCGA data. Nucleic Acids Res. 2016;44:e71

1

353    31. Gene Expression Omnibus. https://www.ncbi.nlm.nih.gov/geo/. Accessed 7 August 2019

354    32. Davis S, Meltzer PS. GEOquery: a bridge between the Gene Expression Omnibus (GEO) and

355        BioConductor. Bioinformatics. 2007;23:1846–7
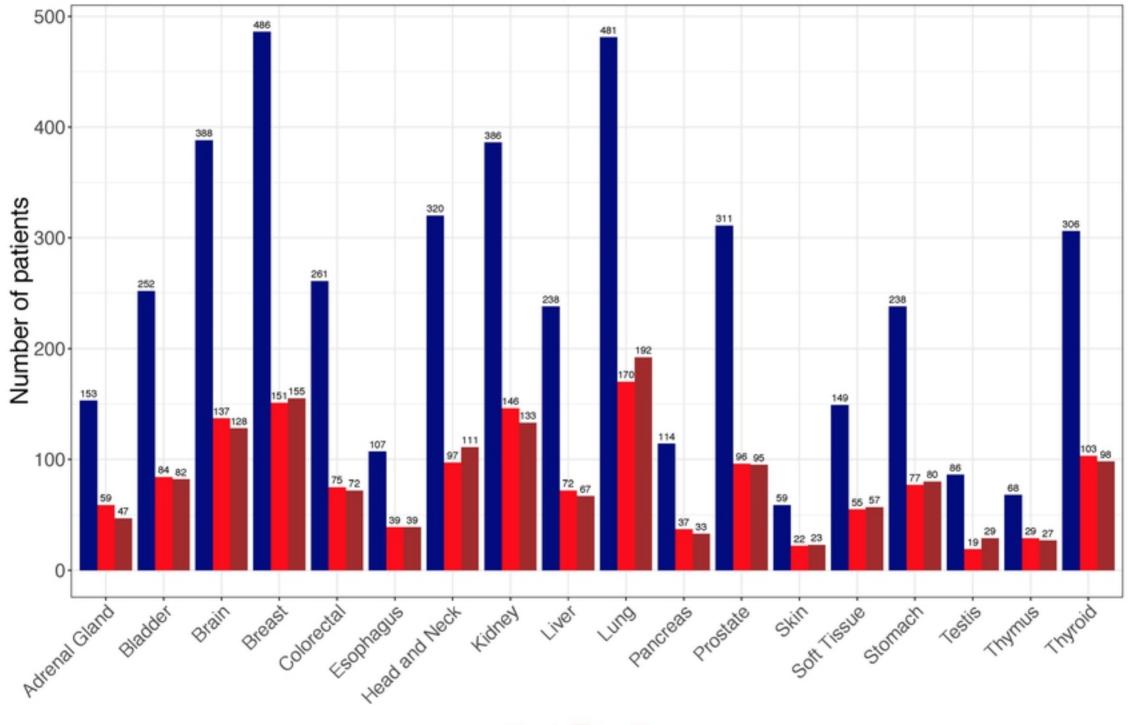
356    33. Abadi M, Agarwal A, Barham P, Brevdo E, Chen Z, Citro C, et al. TensorFlow: Large-scale machine

357        learning on heterogeneous systems.  In: OSDI'16 Proceedings of the 12th USENIX conference on

358        Operating Systems Design and Implementation. 2016;265-283

359    34. Glorot X, Bengio Y. Understanding the difficulty of training deep feedforward neural networks. In:

360        Proceedings of the International Conference on Artificial Intelligence and Statistics. 2010;249-256

361    35. Diederik P. Kingma and Jimmy Lei Ba. Adam. A method for stochastic optimization. arXiv.

362        2014;1412.6980v9

363    36. Qian N. On the momentum term in gradient descent learning algorithms. Neural Netw. 1999;12:145-

364        151.

365    37. Mcmahan HB and Streeter M. Delay-Tolerant Algorithms for Asynchronous Distributed Online

366        Learning. Advances in Neural Information Processing Systems (Proceedings of NIPS). 2014;1–9.

367    38. Varadhachary GR, Raber MN. Cancer of unknown primary site. N Engl J Med. 2014;371:757–65

368    39. Krämer A, Hübner G, Schneeweiss A, Folprecht G, Neben K. Carcinoma of Unknown Primary - an

369        Orphan Disease? Breast Care (Basel). 2008;3:164-170.

370    40. Ettinger DS, Agulnik M, Cates JM, Cristea M, Denlinger CS, Eaton KD, et al. NCCN Clinical

371        Practice Guidelines Occult primary. J Natl Compr Canc Netw. 2011;9:1358-95.

372    41. Pentheroudakis G, Golfinopoulos V, Pavlidis N. Switching benchmarks in cancer of unknown

373        primary: from autopsy to microarray. Eur J Cancer. 2007;43:2026–36

374

375

376

1

# Supporting information

**S1 Table.** **Cancer origin predictions for 1468 patient samples from TCGA.**
**(DOCX)**


**S2 Table.** **Confusion matrix for TCGA test set predictions.**
**(CSV)**


**S3 Table.** **Cancer tissue origin predictions for 143 metastatic cancer samples.**
**(DOCX)**


**S4 Table.** **Confusion matrix for metastatic cancer samples in TCGA test set.**
**(CSV)**


**S5 Table.** **Cancer origin predictions for 581 samples from GEO datasets.**
**(DOCX)**


**S6 Table.** **Confusion matrix for GEO sample predictions.**
**(CSV)**

1

Figure 1

# Figure 2



Note: HL: Hidden layer

# Figure 3

# Figure 4

# Figure 5