# An Explanatory Deep Learning Model for the Prediction of Biologically Relevant Gene Expression in Non-small Cell Lung Tumors and their Microenvironment

Vibha Rajesh Rao, Dartmouth Health; Adrienne A. Workman; Liang Lu; Xiaoying Liu, MD; Shrey S. Sukhadia, PhD

## Introduction/Background

Lung adenocarcinoma (LUAD), the most common subtype of non-small cell lung cancer, exhibits pronounced histologic and molecular heterogeneity, hindering accurate prognosis and effective treatment. Current methods to characterize this heterogeneity, such as histopathology and molecular assays, are constrained by cost, turnaround time, and tissue availability. To address these limitations, we developed XpressO-Lung, an explainable deep learning model that infers gene expression directly from H&E-stained whole slide images (WSIs), enabling scalable morpho-genomic analysis.

## Methods/Intervention

WSIs and matched RNA-seq data from 200 LUAD patients were obtained from TCGA. Data were processed using our deep learning pipeline, XpressO, which segments tumor regions of interest (ROIs) on WSIs and extracts high-dimensional patch-level feature embeddings. These features were paired with binarized RNA-seq expression data, categorized as "high" or "low" for 12 LUAD-relevant genes using median thresholds, to train gene-specific classifiers. Models were trained using a weakly supervised attention-based multiple instance learning framework, with k-fold cross-validation. Attention heatmaps enabled spatial interpretation of model predictions across WSIs.

## Results/Outcome

XpressO-Lung demonstrated strong predictive performance across all 12 biomarkers (AUC ≥ 0.84 for NAPSA, TP53I3, SLC47A1). In representative test slides, NAPSA-bright regions aligned with glandular structures composed of columnar cells, indicative of alveolar differentiation, while CDKN2A-silent regions corresponded to solid tumor nests with nuclear atypia, suggesting unchecked proliferation. This morpho-genomic pattern of high NAPSA with low CDKN2A recurred across multiple tissues and is characteristic of LUAD clones bearing 9p21 deletions. Clinically, such tumors exhibit poor response to PD-(L)1 monotherapy but are vulnerable to CDK4/6 and PRMT5-targeted agents. In the tumor microenvironment, CD8A and KRT7 co-localized in immune-inflamed stroma infiltrated with lymphocytes, while high KRT7 and low CDKN2A marked invasive tumor nests surrounded by desmoplastic stroma respectively. These spatially consistent and biologically informative patterns underscore the model's ability to localize clinically relevant morpho-genomic signals directly from histology, highlighting its potential utility for precision oncology in LUAD.

## Conclusion

XpressO-Lung accurately predicts biomarker expression and spatially localizes expression patterns from routine WSIs.

## Statement of Impact

XpressO-Lung offers a cost-effective, interpretable alternative to molecular assays, supporting precision

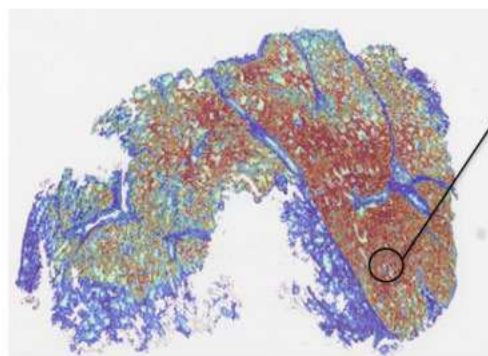oncology directly from H&E slides in routine lung cancer care.
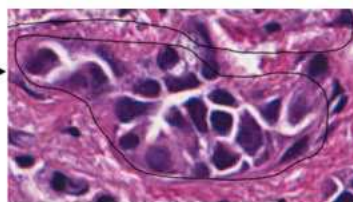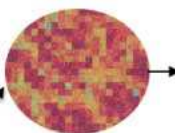


**Fig 1a.1:** NAPSA heatmap

expression p_0 = 0.90
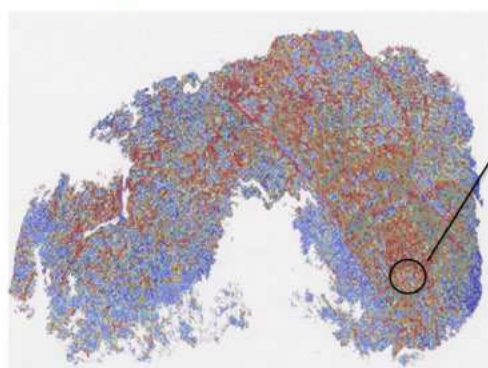
**Fig 1a.2:** NAPSA H&E
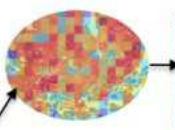
expression p_0 = 0.75

**Fig 1b.1:** CDKN2A heatmap

**Fig 1b.2:** CDKN2A H&E

Figure 1: Attention heatmaps and corresponding H&E slides for WSI TCGA-97-A4M1 showing model-predicted expressions of NAPSA (Fig 1a.1 and Fig 1a.2) and CDKN2A (Fig 1b.1 and Fig 1b.2). Predictions were generated using the top-k most highly attended patches. The model predicted high NAPSA expression with probability 0.90, and low CDKN2A expression with probability 0.75. Red regions indicate strong attention weights, highlighting spatial features driving each gene's expression prediction.
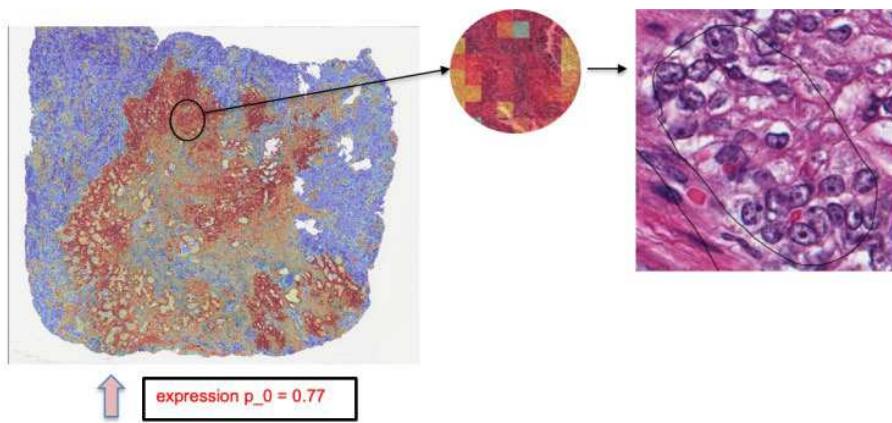
Figure 2: Attention heatmaps and corresponding H&E slides for WSI TCGA-97-A4M5 showing model-predicted expressions of KRT7 (Fig 2a.1 and Fig 2a.2) and CDKN2A (Fig 2b.1 and Fig 2b.2). Predictions were generated using the top-k most highly attended patches. The model predicted high KRT7 expression with probability 0.77, and low CDKN2A expression with probability 0.66. Red regions indicate strong attention weights, highlighting spatial features driving each gene's expression prediction.

| Biomarker | Best performing fold | AUC | Precision | Recall | F1 |
|---|---|---|---|---|---|
| NAPSA | 14th | 0.92 | 0.85 | 0.85 | 0.84 |
| SLC47A1 | 9th | 0.84 | 0.85 | 0.85 | 0.85 |
| TP53I3 | 14th | 0.84 | 0.85 | 0.85 | 0.85 |
| KLRB1 | 11th | 0.83 | 0.76 | 0.75 | 0.75 |
| FAM189A1 | 12th | 0.8 | 0.75 | 0.73 | 0.73 |
| TICAM1 | 13th | 0.78 | 0.76 | 0.73 | 0.73 |
| CD8A | 14th | 0.77 | 0.7 | 0.7 | 0.69 |
| CXCL13 | 9th | 0.76 | 0.67 | 0.65 | 0.64 |
| TTF1 | 12th | 0.72 | 0.84 | 0.7 | 0.65 |
| CDH3 | 9th | 0.66 | 0.68 | 0.65 | 0.61 |
| KRT7 | 8th | 0.65 | 0.68 | 0.65 | 0.62 |
| CDKN2A | 11th | 0.64 | 0.65 | 0.63 | 0.63 |

Table 1: Performance metrics for the predicted expression of twelve genes in the test set as a function of the best performing kth fold for each gene.

## Keywords