

Expanding the Possibilities for Cancer Patients by Predicting the Stages of Cancer

Prof. Trapti Soni¹, Shaikh Nabeel², Pankaj Jaiswal³, Atharva Chavan⁴, Aditya Tiwari⁵

^{1,2,3,4,5} Dept. Information Technology Pillai HOC College of Engineering and Technology (Mumbai University) Rasayani, India

Abstract - The "Expanding the Possibilities of Cancer Patients by Predicting the Stages of Cancer" project develops a systematic three-part system to enhance cancer diagnosis and patient care delivery. Users and administrators can securely access the system through login during the Registration Phase while providing hospital details and organizational information. The management system allows administrators to effectively organize patient records through separate cancer groups including Lung Cancer, Oral Cancer and Breast Cancer as well as demographic data involving habits and family history. During the Processing Phase medical institutions create imaging data through advanced methods including MRI alongside other scanning technologies that need special hardware to display images. The system developers created an application-specific interface which improves image examination by allowing detailed high-resolution views for clinical diagnosis. From these images the system generates a cancer stage prediction that identifies either Stage 1 through 3. The Report Generation Phase produces an extensive document that presents patient information together with diagnosis stage and hospital data as well as survival rate projections through simple visual charts.

Keywords - Artificial Intelligence, Machine Learning, Cancer Patient Journey, Image Classification, CNN.

I. INTRODUCTION

The project "Expanding the Possibilities of Cancer Patients by Predicting the Stages of Cancer" addresses global healthcare challenges by improving cancer diagnosis and data management using AI-ML techniques. It operates through three main phases: Registration, Processing, and Report Generation. In the Registration Phase, users log in securely and submit hospital details. Admins manage patient records by grouping individuals based on cancer types like Lung Cancer, Oral Cancer, and Breast Cancer, while also recording demographic information such as habits, family history, and blood groups. [1]

The Processing Phase involves analysing complex

MRI or scan images using a custom-built application that enhances image clarity for better visualization. The system predicts cancer stages (Stage 1, Stage 2, or Stage 3) using the Logistic Regression algorithm, ensuring accurate diagnosis. Finally, the Report Generation Phase compiles patient data, cancer stage details, and hospital information while estimating survival rates. The system employs the YOLO algorithm to present survival predictions in clear bar graphs for improved decision-making. [2]

II. LITERATURE REVIEW

Pathology images now receive analysis through multimodal machine learning models that combine them with other kinds of data points like gene expression profiles for improved clinical and biological research. Several current multimodal data fusion frameworks struggle to effectively model intricate relationships which exist between different data types. Researcher utilized an attention-based fusion architecture which integrates pathology image data through undirected graphs together with gene expression information for advanced analytical purposes. The model integrates modalities through attention mechanisms to unite their embedding elements and use this combination for predictive analysis of individual patient survival outcomes.

Both the pathology images and their spatial relationships between tissue regions are converted into graph structures for analysis. The graph-based embedding's receive additions from gene expression embedding's so that the model becomes capable of generating tumor-based survival predictions. The framework establishes new performance standards for survival prediction in non-small cell lung cancer patients by surpassing the present multimodal approaches in the field.

III. EXISTING SYSTEM

Multiple weaknesses exist within the present structures used to detect cancer and handle patient medical information. Cancer diagnosis together with patient data management heavily depends on human-operated entry techniques which lead to slow operation times and frequent errors in input. Several healthcare facilities store their patient records across separate systems which prevents healthcare providers from accurately tracing medical backgrounds and analyzing treatment development. MRI scans together with other medical imaging data need specialized hardware to view and analyze them adequately which leads to restricted accessibility for healthcare providers. The existing diagnostic approaches produce uncertain cancer stage predictions which leads medical practitioners to make delayed or improper disease stage recognition. The current manual approach to report generation means patient data and survival rate and critical information become difficult to display in an easily comprehensible structure. The present situation underlines the demand for an updated system which integrates AI-ML techniques for optimizing data administration together with picture analysis and diagnostic precision.

IV.METHODOLOGY

1) Registration Phase:

When users register for the system they provide protected information which forms the foundation of system operations through proper data organization. The system allows users to access the platform only when they provide their email credentials along with their password. Users access the system by providing organizational details at the start including hospital

name together with email and address information. The system gives administrators full authority over patient record management through its categorization system of Lung Cancer and Oral Cancer together with Breast Cancer. Under the oversight of admins patients' process status can be recorded as “In Testing Process.”

A patient information recording system within the program type includes the Clinical Form for collecting medical history then the Baseline Form alongside the Pathology Form for lab results while the Radiology Form accepts imaging data. Users need to provide additional demographic information including habits along with blood group and family background data in order to help identify risk factors which can increase prediction accuracy.

2) Processing Phase:

The Processing Phase functions as the vital operational phase because it applies Machine Learning techniques on medical images to determine cancer stage. Medical staff undergo various scan procedures including MRI that produces complex image data. These images currently do not work well with standard smartphones or laptops or tablets because they need unique specialized display hardware to display correctly. The system developed a customized application which enabled easy viewing of medical images to resolve this restriction. Medical professionals have access to a custom application which allows them to zoom through images in detail for analytical purposes. The feature allows precise observation of tumor characteristics in order to achieve accurate diagnosis.

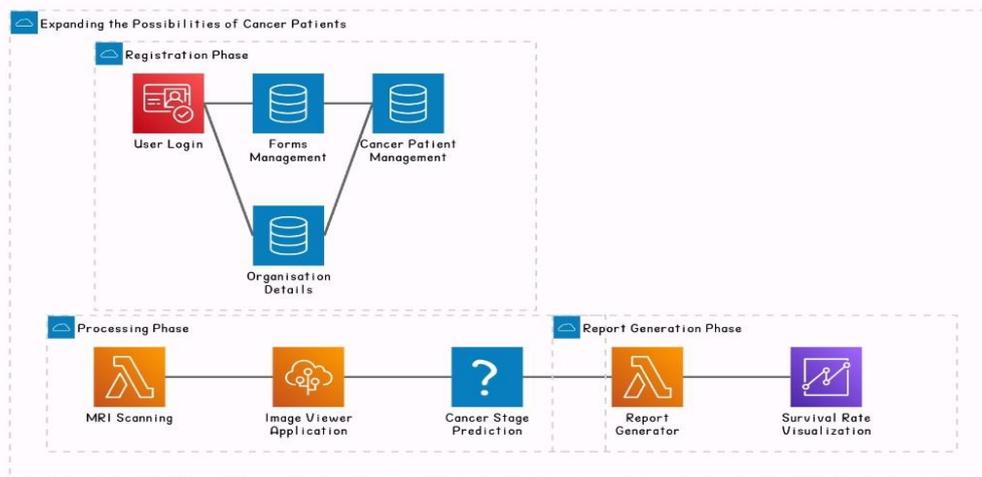


Fig. 1

The pre-processed images undergo analysis involving the powerful logistic regression model to detect cancer stages in classification-based medical tasks. The system extracts important identification features of tumor size and both shape characteristics and intensity measurements to perform cancer stage classification. Within the model the sigmoid function provides output prediction between 0 and 1 to determine cancer stages up to Stage 3 while Logistic Regression serves as the selected algorithm for its accurate medical dataset analysis.

3) Report Generation Phase:

The system finalizes its process by creating one consolidated document with patient files and medical diagnosis records plus anticipated survival rate projections. The report contains essential patient information which includes their identification along with residential details and cancer diagnosis details such as type and stage as well as hospital application data. Through its clear bar graph presentation of survival rate predictions the YOLO (You Only Look Once) algorithm provides healthcare providers with effective tools to assess patient conditions for making informed decisions.

V.MATERIALS AND METHODS

A) Study Population

We retrieved whole-slide images and bulk gene expression records and clinical as well as demographic characteristics including overall

survival data from TCGA [3] and CPTAC and NLST.

TCGA represents a historic cancer genomics program that analyzed the molecular changes found in numerous pairs of primary cancer and matched normal tissue samples across multiple cancer types. A national consortium called CPTAC [4] uses large-scale proteome and genome analysis to speed up cancer research for explaining molecular cancer basis. NLST [5] was a randomized controlled study examined whether low-dose helical computed tomography helped decrease mortality rates in high-risk individuals seeking lung cancer screening.

The use of low-dose helical computed tomography screening produces better lung cancer mortality outcomes when employed for high-risk patients in comparison to chest radiography testing. Multiple research investigations discovered genetic signatures which relate to lung cancer survival outcomes. Researchers have demonstrated increases of B cell signatures throughout LUAD and LUSC tumors but opposite tumor prognosis results happen when B cells infiltrate cancer tissue only in LUAD. The analysis included five B cell population-specific gene signatures retrieved from single-cell RNA sequencing studies of lung cancer and normal adjacent tissue: Sinjab (Plasma) with name sig#1 and Sinjab (B Cell) with name sig#2 and Sinjab (B: 1) under sig#3 as well as Sinjab (B: 0) under sig#4 [6] and Travaglini (B) with sig#5 denoted as sig#5 [7].

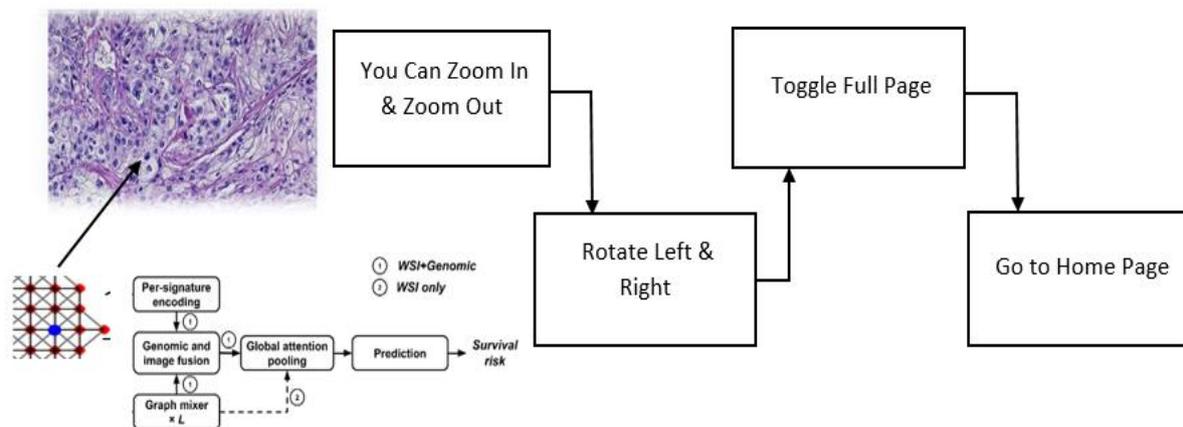


Fig. 2

Fig. 2: Graph attention-based fusion framework: The mixer framework utilizes both graph node embedding's together with gene expression signature embedding's to discover a spatial fingerprint of WSI-transcriptomic association by means of an attention-based method that enables survival prediction. Graph mixer (right) includes repeated node-mixture layers

then channel-mixture layers to identify corresponding (blue and purple) nodes and extract more meaningful features (blue to green) on the graph. Each part of the architecture contains fully-connected layers which form the per-node encoding module as well as the per-signature encoding module and the prediction module. This paper discusses

various details about the graph mixer while also describing genomic and image fusion modules along with global attention pooling modules as described in the previous Section.

signatures of areas which strongly relate to patient survival. The research developed two survival models including (a) ISM denoting the imaging survival model based on WSIs and (b) FSM serving as the Fusion Survival Model merging WSI and genomic data.

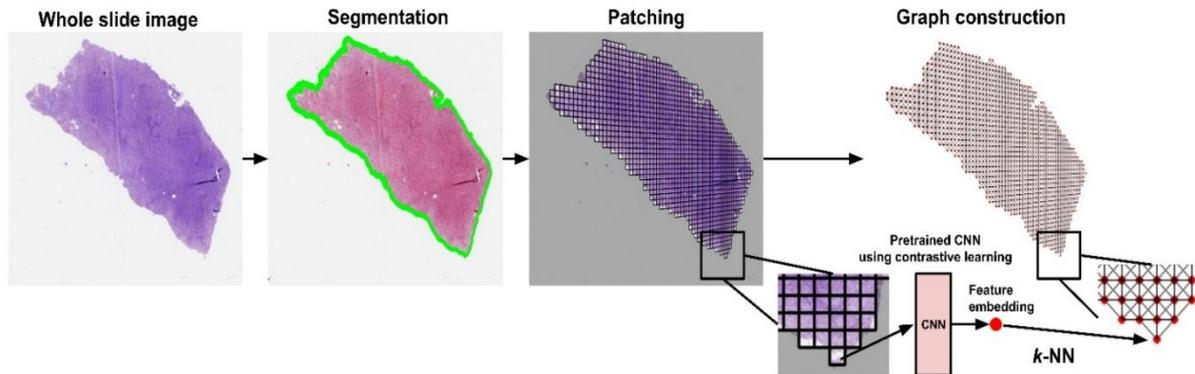


Fig. 3

B) Modeling Framework

The framework uses joint learning to decode WSIs and genomic data while making tumor survival predictions along with identifying image-genomic (v_i, v_j) belongs to E and A_{ij} equals 0 when (v_i, v_j) is not in E . Each image patch in the dataset requires minimum one neighboring relationship while maximum possible connections with eight adjacent patches thus resulting in A matrix elements between 1 and 8. A graph exhibits a node feature matrix H that contains the feature vectors of each node which has a dimension of C and ranges in size of N vertices. This matrix appears in the set $IRN \times C$. Each image patch corresponds to a feature vector that is used for computing node information with a total number of $N = |V|$ nodes.

Fig. 3: Whole slide image (WSI) processing and graph construction:

The WSI pipeline operated according to a three-stage process where foreground-background segmentation followed tessellation of images into patches before building an undirected graph. The generated node features through contrastive learning operated as graph attributes in the structure.

VI. SYSTEM FLOW

The process starts by initializing the system after

Whole slide image processing and graph construction:

The analysis pipeline consists of two sections for processing whole slide images while building the final networks. An undirected graph $G = (V, E)$ contains V which stands for the nodes associated with image patches of the WSI as well as E which describes the edges connecting adjacent nodes in V (Fig. 2). The adjacency matrix for G has its entries defined by $A = [A_{ij}]$ where A_{ij} takes value 1 when

which user input acceptance takes place. At this moment the workflow starts its operation.

The request for login credentials starts when users have to enter their username created from an email ID along with a chosen password. Only authorized users consisting of healthcare staff and administrators and system operators maintain access to the platform.

The system checks the entered credentials to verify their accuracy. The system analyzes whether the entered username and password exist exactly as recorded data in the database.

Users who supply valid login credentials to the system will receive access permissions. The system collects necessary patient information during this process which gets recorded within its database. The system collects essential patient-specific data consisting of name and age in combination with address and medical history elements and additional necessary information. The administrator runs a process which places patients into defined groups following a determination of their cancer classification (Lung Cancer, Oral Cancer, or Breast Cancer).

Following the addition of patient information the system moves forward to run the image verification

tests. The specialized Viewer application of the system accepts medical images uploaded from MRI or other scan sources. An application has been developed with improved image analysis functions that enable precise medical image evaluations for healthcare providers.

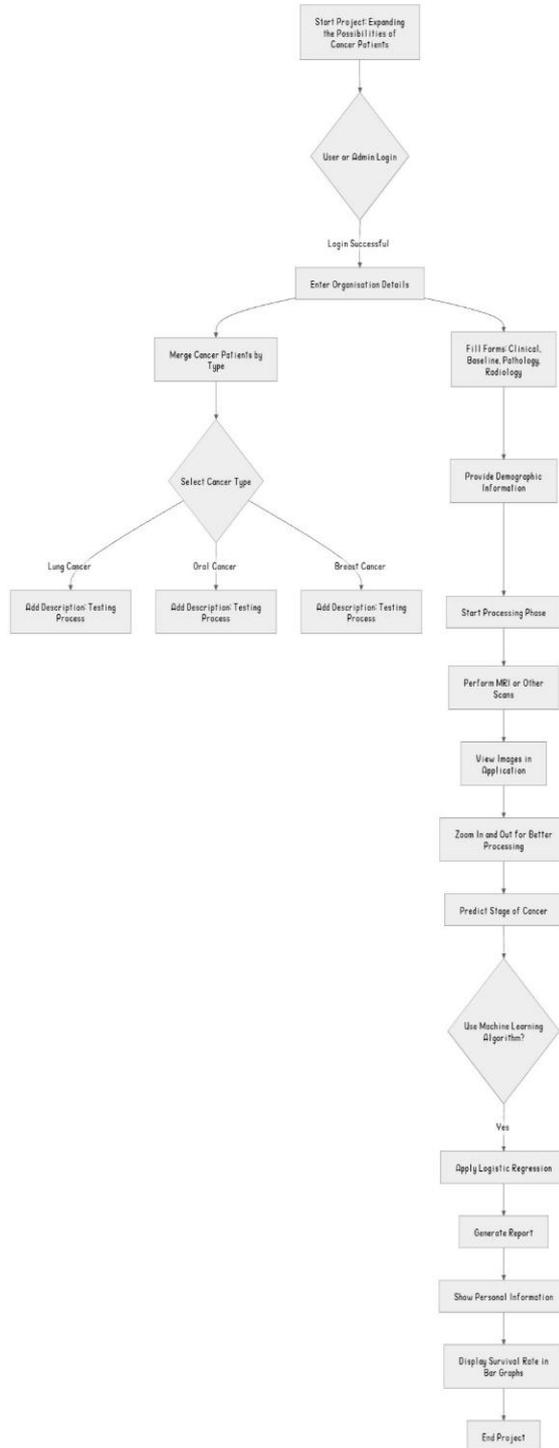


Fig. 4

The analyzed medical images serve as input for system identification of patient cancer type. Medical professionals will receive clear diagnostic information because the system performs a three-way

categorization between Lung Cancer and Oral Cancer and Breast Cancer.

A patient must submit a form which exclusively targets the cancer type diagnosed during their assessment. The system uses different types of forms such as Clinical, Baseline, Pathology and Radiology formats. Accurate diagnosis and risk assessment together with future treatment preparation depend on the information obtained.

The system execution concludes after step-by-step completion of login verification and patient data entry and image analysis and form submission. After data entry the system holds collected data and diagnostic information which enables healthcare professionals to review it.

VII. RESULT

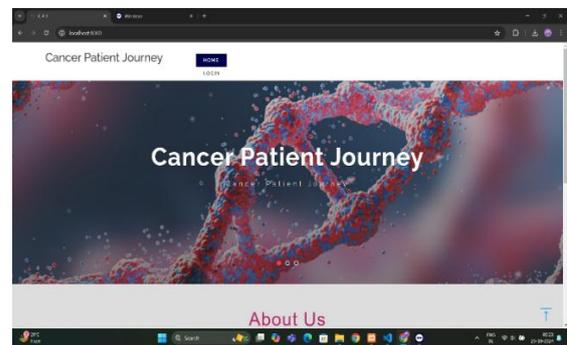


Fig. 5

Fig 5 shows the log in page of our registration phase.

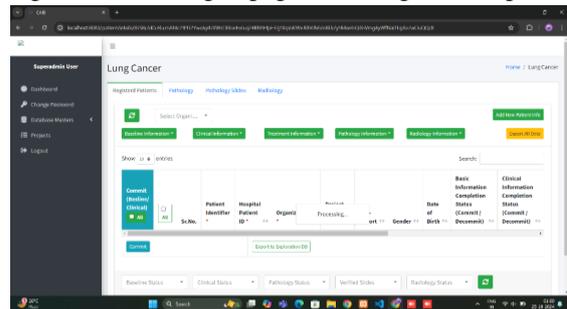


Fig. 6

In the above fig 7, it gives all details about the patient Such as Hospital name, Patient ID, Gender etc.

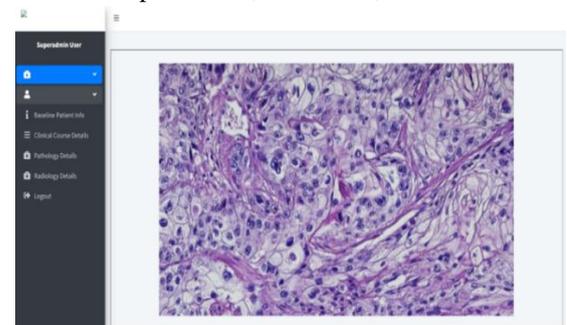


Fig. 7

This fig 7, is the Whole Slide Image(WSI) Page, where you can View Cancer Patients Pathology Images.

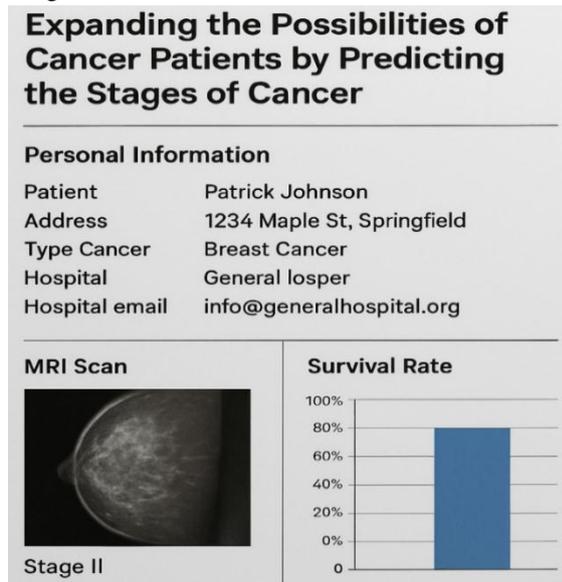


Fig. 8

Fig 8 is the Final Report of the Patients, where it shows all the details about the Patient.

VIII. CONCLUSION

The "Expanding the Possibilities of Cancer Patients by Predicting the Stages of Cancer" initiative employs AI and ML technologies to boost cancer assessment and patient treatment quality. Logistic Regression enables stage prediction in the system while the YOLO algorithm produces reports which leads to enhanced data analysis precision and efficiency. The designed application provides better medical image visualization by dealing with hardware barriers which results in improved diagnosis effectiveness. The system facilitates structured data registration with report generation that includes survival rate analysis for patients. Healthcare professionals will benefit from this integrated system because it enhances their scientific support while improving diagnostic effectiveness. The combination of precision capabilities delivers better cancer patient treatment results.

IX. FUTURE SCOPE

This project holds major potential growth and development opportunities. Additional machine learning approaches such as Random Forest or XGBoost enhance the prediction accuracy process due to their ability to detect sophisticated patterns in medical information. The following versions of the

application will implement deep learning technologies including Convolutional Neural Networks (CNNs) to enhance image investigations which will lead to better early-stage cancer identification. By adding new cancer varieties to the system and linking it to wearables in patient monitoring it enables real-time disease surveillance. A cloud-based platform developed by the team would advance doctor and global researcher access to data thereby spurring collaborative work in cancer research.

Healthcare personnel benefit from simpler report analysis through specialized visualization tools combined with interactive dashboards to design better treatment approaches and deliver improved patient care.

X. REFERENCES

- [1] B. He, L. Bergenstrahle, L. Stenbeck, A. Abid, A. Andersson, A. Borg, J. Maaskola, J. Lundeberg, and J. Zou, "Integrating spatial gene expression and breast tumour morphology via deep learning," *Nat Biomed Eng*, vol. 4, no. 8, pp. 827–834, 2020.
- [2] X. Tan, A. Su, M. Tran, and Q. Nguyen, "SpaCell: integrating tissue morphology and spatial gene expression to predict disease cells," *Bioinformatics*, vol. 36, no. 7, pp. 2293–2294, 2020.
- [3] TCGA Research Network, "The Cancer Genome Atlas Program," Available: <https://portal.gdc.cancer.gov/>.
- [4] N. J. Edwards, M. Oberti, R. R. Thangudu, S. Cai, P. B. McGarvey, S. Jacob, S. Madhavan, and K. A. Ketchum, "The CPTAC Data Portal: A Resource for Cancer Proteomics Research," *Journal of Proteome Research*, vol. 14, no. 6, pp. 2707–2713, 2015.
- [5] The National Lung Screening Trial Research Team, "Reduced Lung-Cancer Mortality with Low-Dose Computed Tomographic Screening," *New England Journal of Medicine*, vol. 365, no. 5, pp. 395–409, 2011.
- [6] A. Sinjab, G. Han, W. Treekitkarnmongkol, K. Hara, P. M. Brennan, M. Dang, D. Hao, R. Wang, E. Dai, H. Dejima, J. Zhang, E. Bogatenkova, B. Sanchez-Espiridion, K. Chang, D. R. Little, S. Bazzi, L. M. Tran, K. Krysan, C. Behrens, D. Y. Duose, E. R. Parra, M. G. Raso, L. M. Solis, J. Fukuoka, J. Zhang, B. Sepesi, T. Cascone, L. A. Byers, D. L. Gibbons, J. Chen, S. J. Moghaddam,

E. J. Ostrin, D. Rosen, J. V. Heymach, P. Scheet, S. M. Dubinett, J. Fujimoto, I. I. Wistuba, C. S. Stevenson, A. Spira, L. Wang, and H. Kadara, “Resolving the Spatial and Cellular Architecture of Lung Adenocarcinoma by Multiregion Single-Cell Sequencing,” *Cancer Discov*, vol. 11, no. 10, pp. 2506–2523, May 2021.

- [7] K. J. Travaglini, A. N. Nabhan, L. Penland, R. Sinha, A. Gillich, R. V. Sit, S. Chang, S. D. Conley, Y. Mori, J. Seita, G. J. Berry, J. B. Shrager, R. J. Metzger, C. S. Kuo, N. Neff, I. L. Weissman, S. R. Quake, and M. A. Krasnow, “A molecular cell atlas of the human lung from single-cell RNA sequencing,” *Nature*, vol. 587, no. 7835, pp. 619–625, Nov 2020.