

Molecular subtyping reveals immune alterations associated with progression of bronchial premalignant lesions

Beane et al.

Supplementary Materials and Methods

N-nitrosotris-(2-chloroethyl)urea (NTCU) mouse sample collection and library preparation

We have previously collected and banked RNA from 40 fresh frozen whole lung sections (curls) and laser microdissected (LCM) tissue isolated with an Acurtus Pixcell II, from SWR/J and A/J mice treated with NTCU. Mice had been treated topically with 15 or 25 μmol NTCU (25 μl of 40 mM NTCU for 15 or 25 weeks) as part of a study performed in accordance with IACUC approved protocol at Roswell Park Comprehensive Cancer Center (Roswell). Samples include examples of: normal (SWR/J $n=3$ LCM & 3 curls & A/J $n=2$ LCM & 1 curl), metaplasia/mild dysplasia (SWR/J $n=5$ LCM & 2 curls), moderate dysplasia (SWR/J $n=7$ LCM & 4 curls & A/J $n=2$ LCM & 1 curls), and severe dysplasia (SWR/J $n=3$ LCM & 2 curls), and *carcinoma in situ*/LUSC (A/J $n=2$ LCM & 2 curls). Samples were extracted using the Qiagen mi-RNAeasy kit according to manufacturer's protocol. Sequencing libraries were prepared from total RNA samples using Illumina® TruSeq® RNA Sample Preparation Kit v2. Each sample was sequenced five per lane on the Illumina® HiSeq 2500 to generate single-end 50-nucleotide reads.

Histological Classification of the NTCU Mouse Samples that underwent RNA Sequencing (n = collected/ n =passed QC after sequencing)							
Mouse Stains	Sample Type	Normal	Mild Dysplasia	Moderate/ Severe Dysplasia	Severe Dysplasia	CIS/SCC Tumor	Total
A/J	LCM	2/2	-	2/2	-	2/1	6/5
A/J	Curls	1/1	-	1/1	-	2/1	4/3
SWR/J	LCM	4/1	5/3	7/3	3/1		19/8
SWR/J	Curls	3/3	2/2	4/3	2/1		11/9
						Total	40/25
Mean RIN values (SD)		4.0(1.8)	3.8(0.5)	3.3(0.6)	2.55(0.1)	3.4(1.2)	

NTCU mouse data processing

De-multiplexing and creation of FASTQ files were performed using Illumina CASAVA 1.8.2. Trimmomatic was used to trim adapter sequences as well as to trim reads of poor quality using the following parameters: ILLUMINACLIP:TruSeq3-SE.fa:2:30:10, LEADING:20, TRAILING:20, SLIDINGWINDOW:4:20, and MINLEN:20. After trimming, greater than 99% of reads were retained in all samples. Samples were subsequently aligned using mm9 and 2-pass STAR¹ alignment. Gene and transcript level counts were calculated using RSEM² and Ensembl annotation. Quality metrics were calculated by STAR and RSeQC³. Initially, 15 samples were removed based on percent of uniquely aligned reads (compared to total reads) less than 15%. Subsequent sample and gene filtering was conducted separately on each set as follows: First, EdgeR⁴ was used to compute normalized data (library sizes normalized using TMM, trimmed mean of M-values, and log2 counts per million computed) and genes were excluded that either had an interquartile range equal to zero or a sum across samples equal or less than 1. Samples were excluded based on values greater than 2 standard deviations from the mean for 1) mean Pearson correlation with all other samples calculated across all filtered genes 2) the 1st or 2nd principal components calculated using the filtered gene expression matrix 3) transcript

integrity number (TIN, computed by RSeQC). After sample filtering, gene filtering was recomputed as described above on the final set of high-quality samples. The data are available from NCBI's Gene Expression Omnibus using the accession GSE111091 [<http://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE111091>].

Immunofluorescent quantification of cell type and proliferative markers

Basal and ciliated cell type markers (KRT5 and TUB1A1) and the proliferative marker (KI67) were manually enumerated for all epithelium within a biopsy in reference to DAPI staining, with a minimum of 500 cells counted per biopsy. The enumeration was conducted on different regions (independent areas of tissue) present on a slide (1-4 regions/biopsy) for each biopsy. A percent of positively stained cells was calculated for each marker in each region enumerated. A binomial mixed effects model via the lme4 R package was used to assess differences in the percentages of cells staining positive for a given protein in each region between the molecular subtypes using the total cells stained in each region as weights and adjusting for patient as a random effect.

TCGA SCC tumors data processing

Log2 transcript per million data across 20,500 genes from 476 LUSC tumors was obtained from Campbell⁵ *et al.* Genes were excluded that either had an interquartile range equal to zero or a sum across samples equal or less than 1. Samples were excluded based on values greater than 2 standard deviations from the mean for more than one of the following criteria: 1) mean Pearson correlation with all other samples calculated across all filtered genes 2) the 1st or 2nd principal components calculated using the filtered gene expression matrix 3) transcript integrity number (TIN, computed by RSeQC). After sample filtering, gene filtering was recomputed as described above (n=17,887 genes) on the final set of high-quality samples (n=471 tumors).

Software Packages

Custom scripts: R-3.3.2, GSVA v1.22.4, Limma v3.30.13, ggplot2 v3.0.0, SummarizedExperiment v1.4.0, edgeR v3.16.5, ConsensusClusterPlus v1.38.0, biomaRt v2.30.0, estimate v1.0.13, heatmap3 v1.2.2, pamr v1.55, lme4 v1.1-13

RNA-seq processing pipeline software for human data (https://github.com/joshua-d-campbell/nf-RNA_Seq_Preprocess v1.0): nextflow v0.24.4, star v2.5.2b, rsem v1.3.0, FastQC v0.11.3, Picard tools v2.8.0, GATK v3.5, rseqc v2.6.4, multiqc v0.9, samtools v1.4

RNA-seq processing software for mouse data: star v2.4.2a, samtools v1.2, picard v1.138, trimmomatic v0.33, rseqc v2.6.1, rsem v1.2.23

Supplementary Tables

	Discovery Cohort	Validation Cohort
Genomic smoking status over time	Number of Subjects	Number of Subjects
Current	9	9
Former	10	5
Current->Former	7	4
Former->Current	3	2
Current->Former->Current	1	0

Supplementary Table 1. Genomic smoking status over time by subject. The smoking status of each subject at each time point was computed based on a previously published smoking-associated gene signature⁶ (see methods for details). The rows indicate the smoking status across all time points sampled for each patient. The -> symbol indicates changes in smoking status over time. There is not a statistical difference between the distribution of subjects in the smoking status categories between the discovery and validation cohorts by a two-sided Fisher's exact Test ($p=0.90$). Source data are provided as a Source Data file.

Variable	Discovery Cohort		Validation Cohort		P-value	
Sample Type	Biopsies	Brushes	Biopsies	Brushes	Biopsies	Brushes
Batch/Illumina Flow Cell Assignment					<2e-16	<2e-16
1	19/190 (10)	12/89 (13.5)	0/105 (0)	0/48 (0)		
2	18/190 (9.5)	13/89 (14.6)	0/105 (0)	0/48 (0)		
3	22/190 (11.6)	9/89 (10.1)	0/105 (0)	0/48 (0)		
4	19/190 (10)	10/89 (11.2)	0/105 (0)	0/48 (0)		
5	29/190 (15.3)	2/89 (2.2)	0/105 (0)	0/48 (0)		
6	24/190 (12.6)	8/89 (9.0)	0/105 (0)	0/48 (0)		
7	20/190 (10.5)	11/89 (12.4)	0/105 (0)	0/48 (0)		
8	17/190 (8.9)	14/89 (15.7)	0/105 (0)	0/48 (0)		
9	22/190 (11.6)	10/89 (11.2)	0/105 (0)	0/48 (0)		
10	0/190 (0)	0/89 (0)	22/105 (10.7)	9/48 (18.8)		
11	0/190 (0)	0/89 (0)	19/105 (9.3)	10/48 (20.8)		
12	0/190 (0)	0/89 (0)	21/105 (10.2)	10/48 (20.8)		
13	0/190 (0)	0/89 (0)	20/105 (9.8)	12/48 (25)		
14	0/190 (0)	0/89 (0)	23/105 (11.2)	7/48 (14.6)		
Total Reads	45.5e+6 (7.2e+6)	45.3+6 (7.9e+6)	42.9e+6 (6.3e+6)	42.6e+6 (4.9e+6)	1.50E-03	0.014
Median Transcript Integrity Number (TIN)	78.4 (1.9)	72.6 (3.4)	76.3 (2.0)	72.3 (2.8)	2.08E-10	0.59
Percent Uniquely Mapped	90.1 (2.9)	89.0 (5.9)	83.9 (9.6)	87.6 (4.9)	2.15E-09	0.15

Supplementary Table 2. Batch information and alignment statistics on samples in the Discovery and Validation cohorts. Statistical tests between the discovery and validation cohorts were performed using two-sided Fisher's exact tests for categorical variables and two-sided Student's T-tests for continuous variables. Percentages are reported for categorical variables and mean and standard deviations are reported for continuous variables. Source data are provided as a Source Data file.

Module Number	Number of Genes	Biological Pathways Associated with Module Genes	Key Genes	FDR for Difference between Molecular Subtypes
1	514	Extracellular Matrix / Cell Adhesion	Collagens, Lamins, TGFb	2.7E-36
2	939	mRNA processing and splicing	RBMs & SRSF	7.2E-05
3	20	Transcriptional regulation in response to stimuli - (AP-1)		
4	64	Immediate-early response genes	JUN & FOS	1.9E-01
5	209	OXPPOS / ETC / TCA	COXs & NDUFs	3.3E-07
6	1295	PCNA, TOP2A, CDC, AURK, RAD, XRCC		2.0E-31
7	180	Cell Cycle / DNA replication / DNA repair	FOXJ1, DYNC	6.6E-57
8	603	Cilium organization and assembly	RPLs & RPSs	1.9E-13
9	112	Ribosomal Proteins/ Translation	CD8A, CD86, GATA, STAT, IL1B, CD163, CD68	3.3E-07
		Immune Activation and Inflammatory Response (leukocyte/lymphocyte regulation)	SP100, HLAs, STAT1	1.3E-02
		Interferon signaling and Antigen Processing and Presentation		

Supplementary Table 3. Summary of biological characteristics of the gene modules. For each gene module, the following characteristics are listed: the module number, the number of genes in the module, the biological pathways enriched in each gene module, select genes from the module, and an FDR value for the difference in GSVA scores for each module between the molecular subtypes (within the discovery cohort biopsies) are reported. The FDR value for the difference between molecular subtypes was calculated using a linear mixed model with molecular subtype as the main effect and patient as a random effect. Source data are provided as a Source Data file.

Variable	DC P-value	VC P-value
Genomic Smoking Status	2.71E-09	2.72E-04
Subject	9.66E-05	5.87E-03
Subject/Time	6.96E-04	1.40E-02
Histology	6.75E-03	9.99E-08
Location	2.57E-02	6.69E-01
Subject/Location	6.01E-02	1.95E-01
Asbestos Exposure	1.23E-01	7.47E-02
Lung Cancer History	1.32E-01	9.92E-01
Progression Status	1.60E-01	1.67E-05
High-risk Job	4.31E-01	8.30E-01
Sex	5.62E-01	8.90E-01
LUSC Tumor Subtype	9.99E-08	1.80E-06
COPD Status	1.62E-01	9.38E-03

Supplementary Table 5. Molecular subtype associations with clinical and biological characteristics within the discovery and validation cohort biopsies. Statistical tests within the discovery and validation cohorts were performed using two-sided Fisher's exact tests. DC = Discovery Cohort and VC = Validation Cohort. Source data are provided as a Source Data file.

Variable	Discovery Cohort Biopsies (n=190)				Validation Cohort Biopsies (n=105)			
	No LC History	LC History - LUSC	LC History - Other	P-Value	No LC History	LC History - LUSC	LC History - Other	P-Value
Molecular Subtype								
<i>Proliferative</i>	14	5	33		12	9	7	
<i>Inflammatory</i>	10	6	21		12	4	14	
<i>Secretory</i>	26	8	27		14	13	7	
<i>Normal-like</i>	9	3	28	p=0.19	6	1	6	p=0.10

Supplementary Table 6. Molecular Subtype associations with previous history of lung cancer. Previous history of lung cancer (LC) was categorized as follows: no history (No LC History), a previous history of LC that include a lung squamous cell carcinoma (LC History – LUSC), and a previous history of LC that does not include a lung squamous cell carcinoma (LC History – Other). Statistical tests within the discovery and validation cohorts were performed using two-sided Fisher's exact tests. Source data are provided as a Source Data file.

			Samples Used For Scoring Each Panel							
IF Panel	All Samples		K5/KI67/Ac-alpha-Tubulin		CD68/CD163		CD4		CD8	
Variable	Discovery Cohort	Validation Cohort	Discovery Cohort	Validation Cohort	Discovery Cohort	Validation Cohort	Discovery Cohort	Validation Cohort	Discovery Cohort	Validation Cohort
Number of Subjects	17	12	7	2	17	12	17	11	17	11
Number of Samples	27	20	8	2	25	18	27	19	26	18
Subtype										
Normal-like	2/27 (7)	1/20 (5)	1/8 (13)	0/2 (0)	2/25 (8)	1/18 (6)	2/27 (7)	1/19 (5)	2/26 (8)	1/18 (6)
Secretory	7/27 (26)	5/20 (25)	1/8 (13)	0/2 (0)	7/25 (28)	4/18 (22)	7/27 (26)	5/19 (26)	7/26 (27)	5/18 (28)
Inflammatory	8/27 (30)	3/20 (15)	2/8 (25)	1/2 (50)	7/25 (28)	2/18 (11)	8/27 (30)	3/19 (16)	8/26 (31)	3/18 (17)
Proliferative	10/27 (37)	11/20 (55)	4/8 (50)	1/2 (50)	9/25 (36)	11/18 (61)	10/27 (37)	10/19 (53)	9/26 (35)	9/18 (50)
Histology										
Normal/Hyperplasia	9/27 (33)	2/20 (10)	3/8 (38)	0/2 (0)	9/25 (36)	1/18 (6)	9/27 (33)	2/19 (11)	9/26 (35)	2/18 (11)
Squamous Metaplasia	3/27 (11)	3/20 (15)	0/8 (0)	1/2 (50)	3/25 (12)	3/18 (17)	3/27 (11)	3/19 (16)	3/26 (12)	3/18 (17)
Mild Dysplasia	1/27 (4)	4/20 (20)	0/8 (0)	0/2 (0)	1/25 (4)	4/18 (22)	1/27 (4)	4/19 (21)	1/26 (4)	4/18 (22)
Moderate Dysplasia	8/27 (30)	9/20 (45)	3/8 (38)	1/2 (50)	7/25 (28)	8/18 (44)	8/27 (30)	8/19 (42)	7/26 (27)	7/18 (39)
Severe Dysplasia/CIS	6/27 (22)	2/20 (10)	2/8 (25)	0/2 (0)	5/25 (20)	2/18 (11)	6/27 (22)	2/19 (11)	6/26 (23)	2/18 (11)
Lesion State										
Progressive/Persistent	7/27 (26)	10/20 (50)	1/8 (13)	1/2 (50)	7/25 (28)	9/18 (50)	7/27 (26)	10/19 (53)	6/26 (23)	9/18 (50)
Regressive	6/27 (22)	5/20 (25)	3/8 (38)	0/2 (0)	5/25 (20)	5/18 (28)	6/27 (22)	4/19 (21)	6/26 (23)	4/18 (22)
Unknown or Normal/Stable	14/27 (52)	5/20 (25)	4/8 (50)	1/2 (50)	13/25 (52)	4/18 (22)	14/27 (52)	5/19 (26)	14/26 (54)	5/18 (28)
Smoking Status										
Current	15/27 (56)	12/20 (60)	6/8 (75)	1/2 (50)	13/25 (52)	12/18 (67)	15/27 (56)	11/19 (56)	14/26 (54)	10/18 (56)
Former/Never	12/27 (44)	8/20 (40)	2/8 (25)	1/2 (50)	12/25 (48)	6/18 (33)	12/27 (44)	8/19 (42)	12/26 (46)	8/18 (44)

Supplementary Table 7. Clinical and biological characteristics of the samples used for immunofluorescence studies. For each characteristic, percentages are reported in parenthesis.

Molecular Subtype	Normal	Normal	Secretory	Secretory	Inflammatory	Inflammatory	Proliferative	Proliferative
Cohort	DC	VC	DC	VC	DC	VC	DC	VC
Number of Progressive/Persistent Lesions	5	1	17	7	7	5	15	7
Number of Regressive Lesions	3	3	8	1	4	1	15	13
Module Number								
1	ns	N/A	ns	N/A	ns	N/A	ns	ns
2	ns	N/A	ns	N/A	ns	N/A	ns	ns
3	ns	N/A	ns	N/A	ns	N/A	0.047	ns
4	0.026	N/A	ns	N/A	ns	N/A	ns	ns
5	ns	N/A	ns	N/A	ns	N/A	ns	ns
6	ns	N/A	ns	N/A	ns	N/A	ns	ns
7	ns	N/A	ns	N/A	ns	N/A	ns	ns
8	0.027	N/A	ns	N/A	0.005	N/A	ns	ns
9	ns	N/A	ns	N/A	ns	N/A	0.0017	0.03

Supplementary Table 8. Statistical associations between progression/persistence versus regression within each molecular subtype and cohort for each gene module. The P-values were calculated based on a linear model (implemented in limma) with GSVA scores for each module as the dependent variable and progression/regression status as the independent variable and patient as a random effect. P-values less than 0.05 are reported. ns= not significant and N/A= not enough samples in each group to conduct the analysis. DC = Discovery Cohort and VC = Validation Cohort. Source data are provided as a Source Data file.

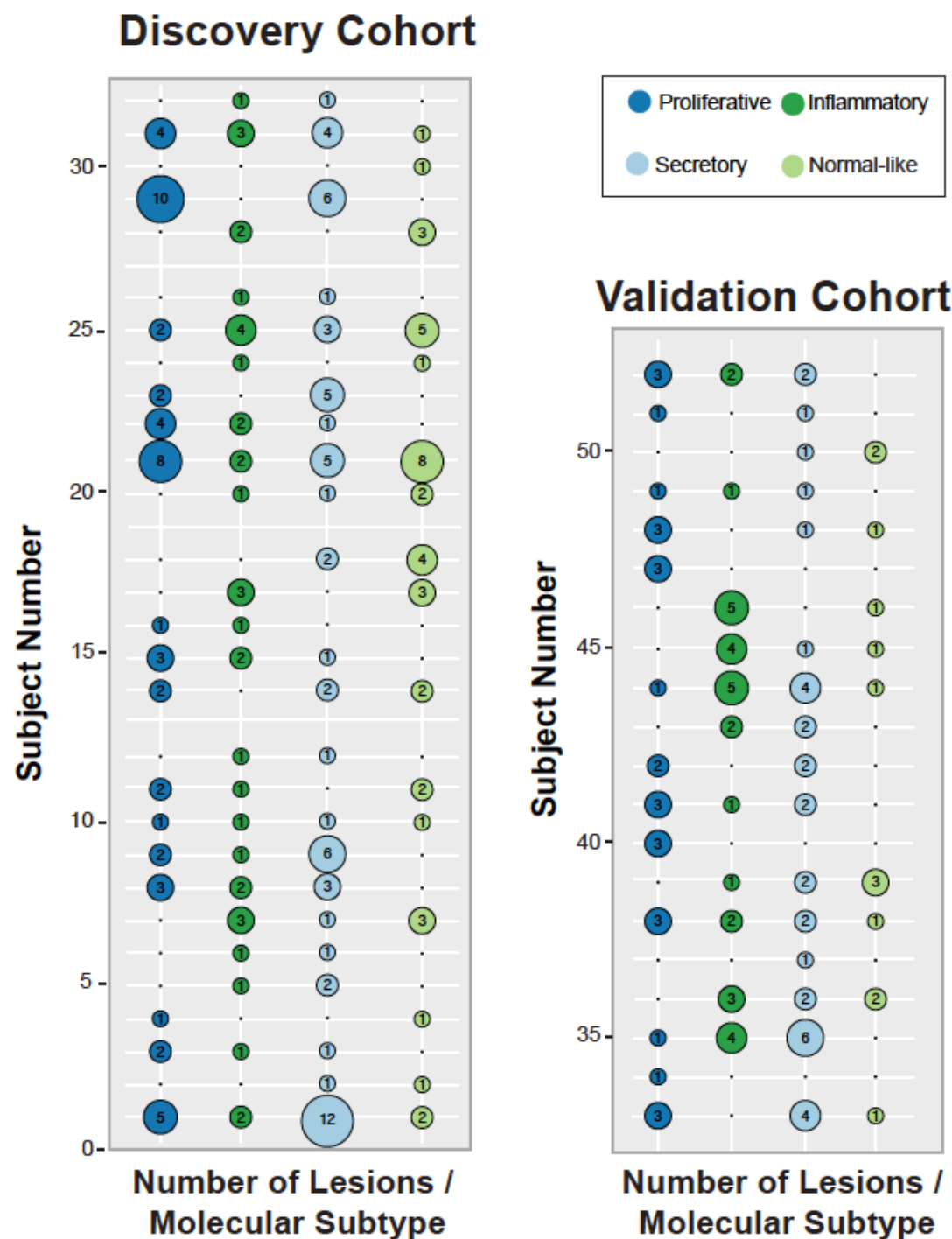
ID	Name	Description
096	VC	True Vocal Cords, Neck
051	Mouth	Floor of Mouth
007	EPIG	Epiglottis
005	ART	Arytenoids
008	FVC	False Vocal Cords
095	TR	Trachea
050	MC	Main Carina, Carina NOS
086	RMB	Right Main Bronchus, incl Secondary Carina right
091	RUL	Right Upper Lobe
093	RULO	Right Upper Lobe Orifice or opening
094	RULS	Right Upper Lobe Stump
092	RULB	Right Upper Lobe Bronchus
087	RML	Right Middle Lobe
089	RMLO	Right Middle Lobe Orifice or opening
090	RMLS	Right Middle Lobe Stump
088	RMLB	Right Middle Lobe Bronchus
082	RLL	Right Lower Lobe
084	RLLO	Right Lower Lobe Orifice
085	RLLS	Right Lower Lobe Stump
083	RLLB	Right Lower Lobe Bronchus
006	BI	Bronchus Intermedius
052	RB1	RUL Apical Segment (AS)
060	RB2	RUL Posterior Segment (PS)
063	RB3	RUL Anterior Segment (ANTS)
053	RB1/2	RUL Carina between RB1 and RB2
054	RB1/3	RUL Carina between RB1 and RB3
061	RB2/3	RUL Carina between RB2 and RB3
059	RB1A/B	RUL AS Carina between RB1 A and B
062	RB2A/B	RUL PS Carina between RB2 A and B
064	RB3A/B	RUL ANTS Carina between RB3 A and B
065	RB4	RML Lateral Segment (LS)
068	RB5	RML Medial Segment (MS)
066	RB4/5	RML LS Carina between RB4 and RB5
067	RB4A/B	RML LS Carina between RB4 A and B
069	RB5A/B	RML MS Carina between RB5 A and B
070	RB6	RLL Superior Basal Segment (SBS)
071	RB6A/B	RLL SBS Carina between RB6A and B
072	RB6A/C	RLL SBS Carina between RB6A and C
073	RB6B/C	RLL SBS Carina between RB6B and C
074	RB7	RLL Medial Basal Segment (MBS)
075	RB7A/B	RLL MBS Carina between RB7A and B
076	RB8	RLL Anterior Basal Seg (ABS)
077	RB8/9	RLL ABS Carina between RB8 and RB9
078	RB8A/B	RLL ABS Carina between RB8A and B
079	RB9	RLL Lateral Basal Segment (LBS)
080	RB9/10	RLL LBS Carina between RB9 and RB10
081	RB9A/B	RLL LBS Carina between RB9A and B
055	RB10	RLL Posterior Basal Segment (PBS)
056	RB10A/B	RLL PBS Carina between RB10A and B
057	RB10A/C	RLL PBS Carina between RB10A and C
058	RB10B/C	RLL PBS Carina between RB10B and C
001	666	Location was surgically altered or removed
002	777	Abstractor needs clinician help to code
003	888	Location code is unknown, illegible
004	999	Location code is blank, not noted
043	LMB	Left Main Bronchus, incl Secondary Carina left
044	LMBD	Left Main Bronchus, Distal
046	LUL	Left Upper Lobe
048	LULO	Left Upper Lobe Orifice or opening
049	LULS	Left Upper Lobe Stump
035	LGL	Lingula
037	LGLO	Lingula Orifice or opening
038	LGLS	Lingula Stump
047	LULB	Left Upper Lobe Bronchus
045	LUDB	Left Upper Division Bronchus
036	LGLDB	Lingular Division Bronchus, lingular bronchus
039	LLL	Left Lower Lobe
041	LLLO	Left Lower Lobe Orifice or opening
042	LLLS	Left Lower Lobe Stump
040	LLLB	Left Lower Lobe Bronchus
009	LB1+2	LUL Apical-Posterior Segment (APS)
018	LB3	LUL Anterior Segment
011	LB1/2	LUL APS Carina between LB1 and LB2
010	LB1+2/3	LUL APS Carina between LB1+2 and LB3
016	LB2A/C	LUL APS Carina between LB2 A and C
017	LB2B/C	LUL APS Carina between LB2B and C
019	LB3A/B	LUL ANTS Carina between LB3A and B
020	LB4	LUL Superior Lingular Segment (SLS)
023	LB5	LUL Inferior Lingular Segment (ILS)
021	LB4/5	LUL SLS Carina between LB4 and LB5
022	LB4A/B	LUL SLS Carina between LB4A and B
024	LB5A/B	LUL ILS Carina between LB5A and B
025	LB6	LLL Superior Segment (SS)
026	LB6A/B	LLL SS Carina between LB6A and B
027	LB6A/C	LLL SS Carina between LB6A and C
028	LB6B/C	LLL SS Carina between LB6B and C
029	LB8	LLL Antero Medial Basal Segment (AMBS)
030	LB8/9	LLL AMBS Carina between LB8 and LB9
031	LB8A/B	LLL AMBS Carina between LB8A and B
032	LB9	LLL Lateral Basal Segment (LBS)
033	LB9/10	LLL LBS Carina between LB9 and LB10
034	LB9A/B	LLL LBS Carina between LB9A and B
012	LB10	LLL Posterior Basal Segment (PBS)
013	LB10A/B	LLL PBS Carina between LB10A and B
014	LB10A/C	LLL PBS Carina between LB10A and C
015	LB10B/C	LLL PBS Carina between LB10B and C

Supplementary Table 9. Lung sites where endobronchial biopsies were obtained. The site code, name, and description are reported for each site.

Antibody	Company	Catalog	Dilution	Antigen retrieval	Species
Immune cell type markers					
CD68	Dako	m0876	1- 250	AR6	mouse
CD163	Cell Marque	163m-16	1-100	AR9	mouse
CD4	Thermo Fisher	ms1528S	1-100	AR9	mouse
CD8	Dako	M7103	1-100	AR9	mouse
Epithelial cell type and proliferation markers					
Ac- α -Tub	Sigma	T6793	1-100	citrate	mouse
KRT5	BioLegend	905-901	1-100	citrate	chicken
KI67	Abcam	ab16667	1-100	citrate	rabbit

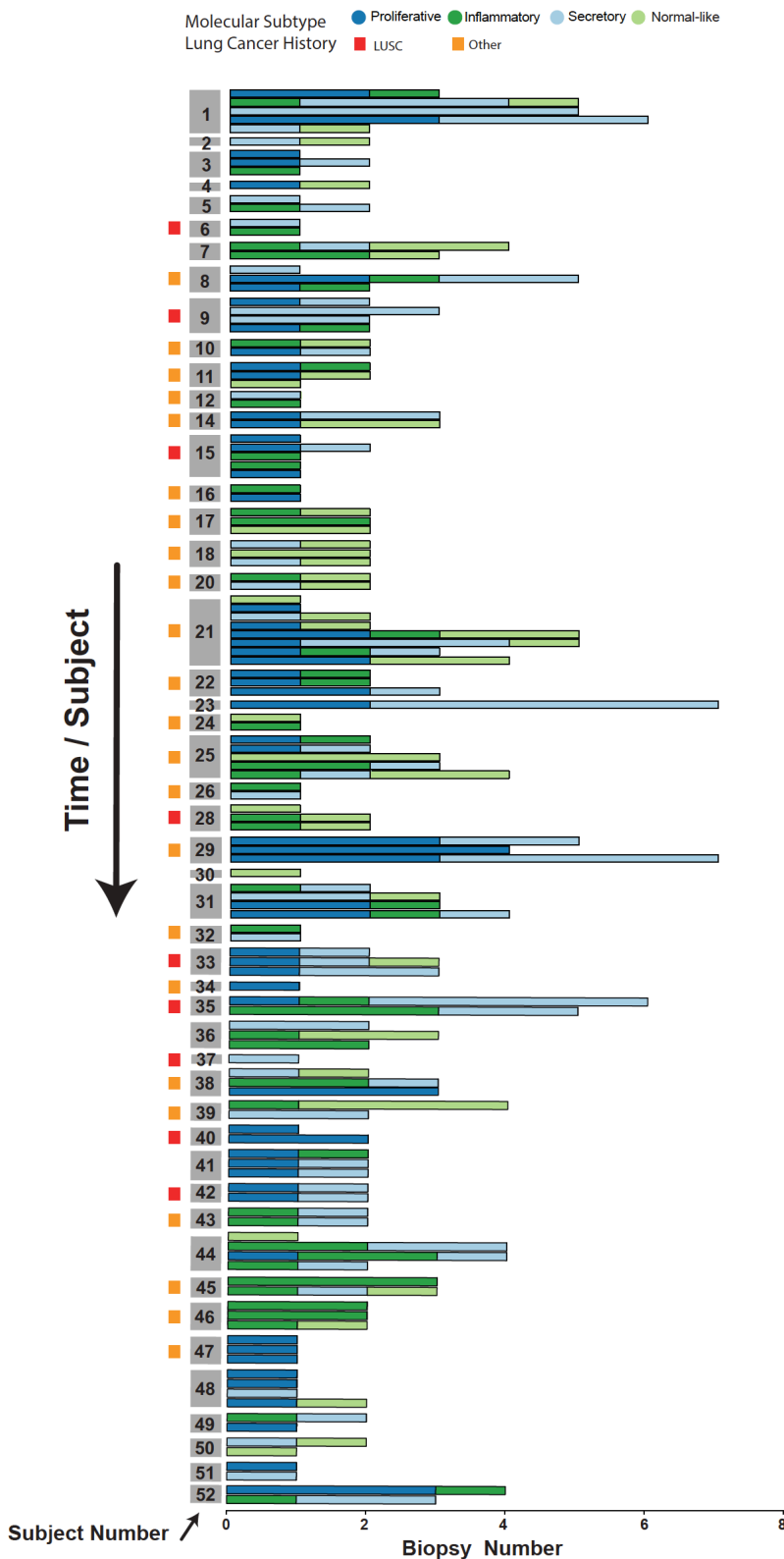
Supplementary Table 10. Antibodies used in the immunofluorescence studies.

Supplementary Figures

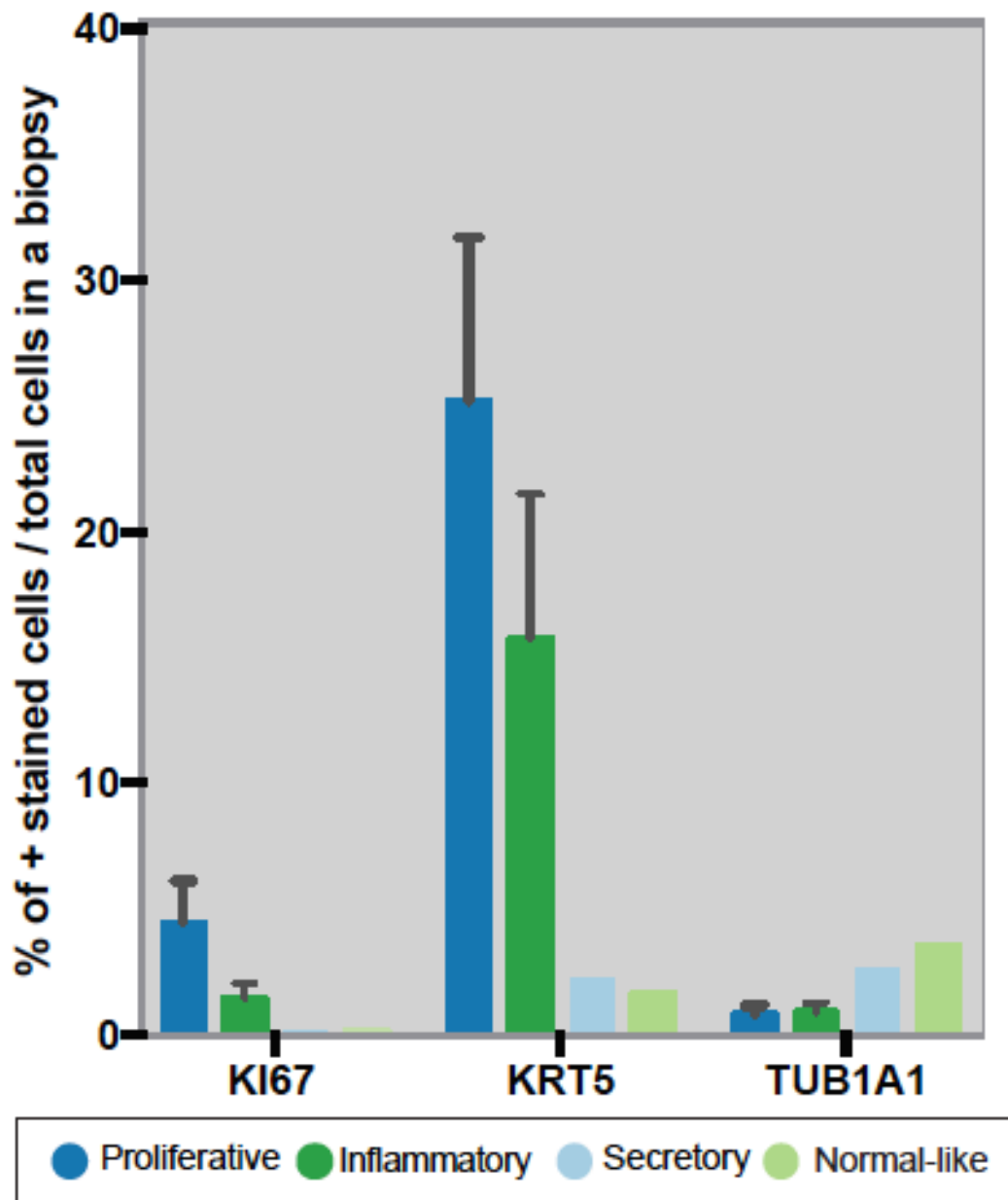


Supplementary Figure 1. Distribution of molecular subtypes by subject across the biopsies. The columns represent the 4 molecular subtypes (Proliferative, dark blue; Inflammatory, dark green; Secretory, light blue; and Normal-like, light green) and the radius of the circle is proportional to the number of samples within each subtype. The discovery cohort samples are shown on the right and the validation cohort samples are shown on the left. Source data are provided as a Source Data file.

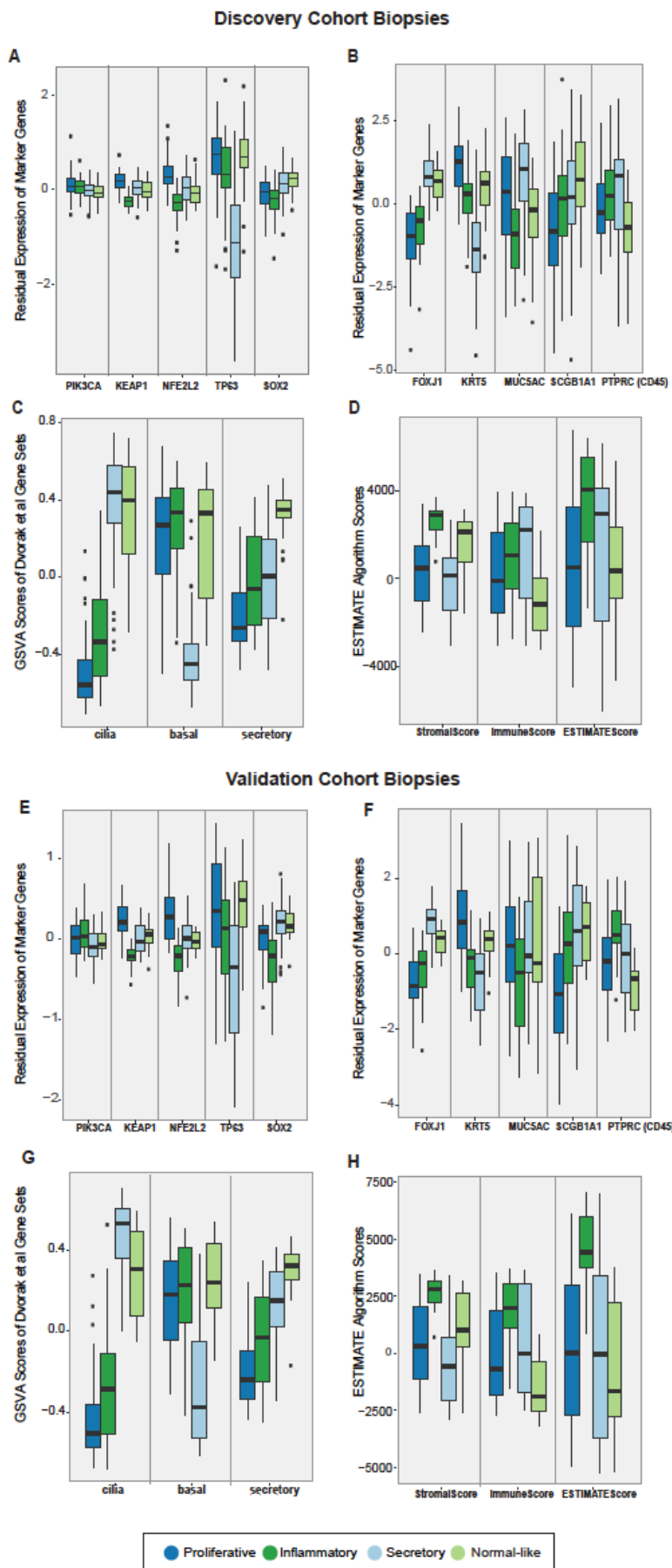
Discovery and Validation Cohort Biopsies



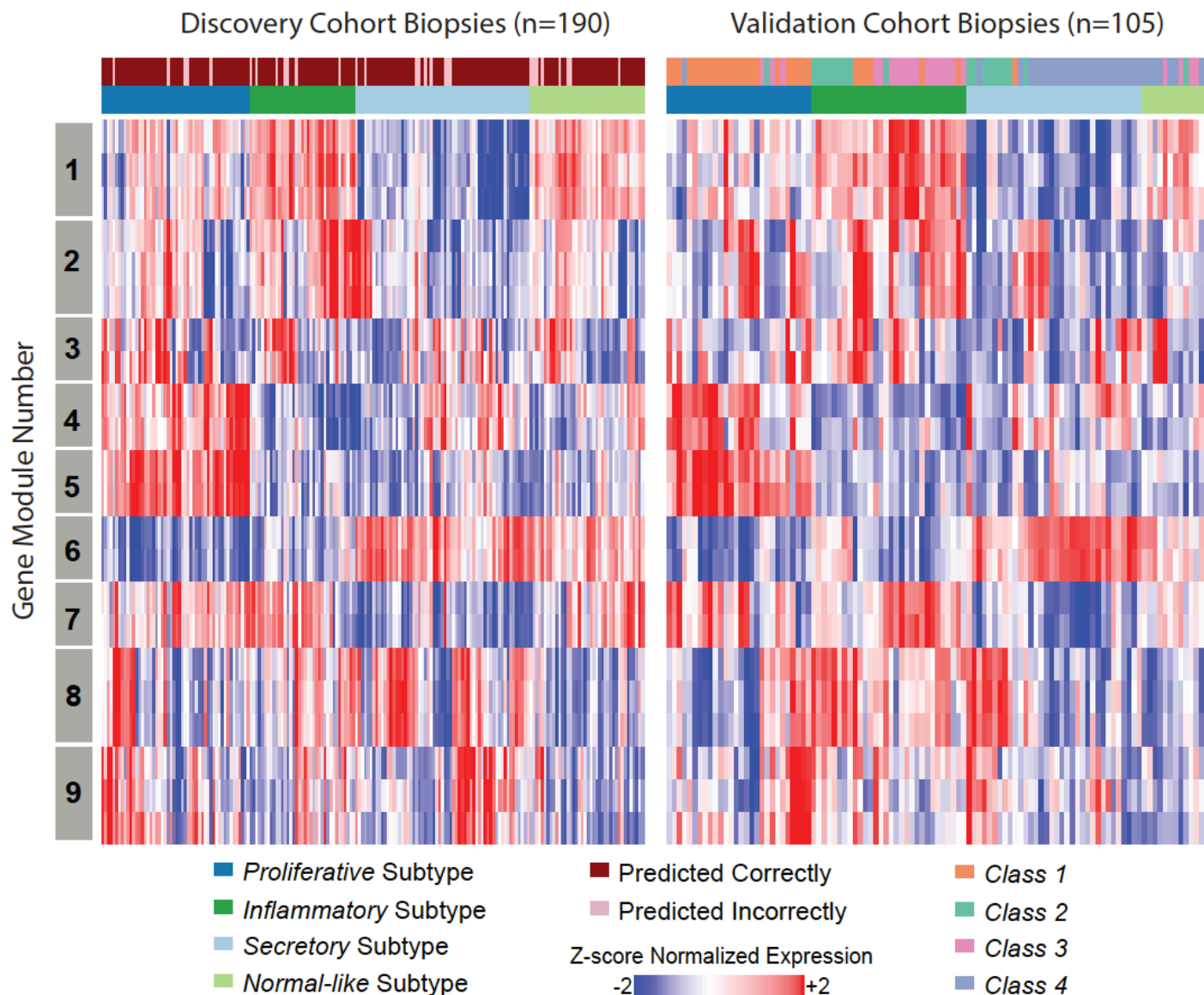
Supplementary Figure 3. Molecular subtype distribution for each subject across bronchoscopy procedures. The barplot shows for each subject and each bronchoscopy procedure the number of biopsies sampled and their corresponding molecular subtype. The y-axis indicates the subject number and whether or not that subject had a prior history of either lung squamous cell carcinoma (LUSC, red) or another type of lung cancer (Other, yellow). The discovery cohort includes subjects 1 through 32 and the validation cohort includes subjects 33 through 52. We did not detect a difference in the diversity of subtype classifications within a subject based on prior history of lung cancer (mean Shannon entropy of subtype classifications within patients with a history of lung cancer = 1.12, n=32 vs. patients without a history of lung cancer = 1.25, n = 17; Wilcoxon Rank Sum test p-value = 0.43). Source data are provided as a Source Data file.



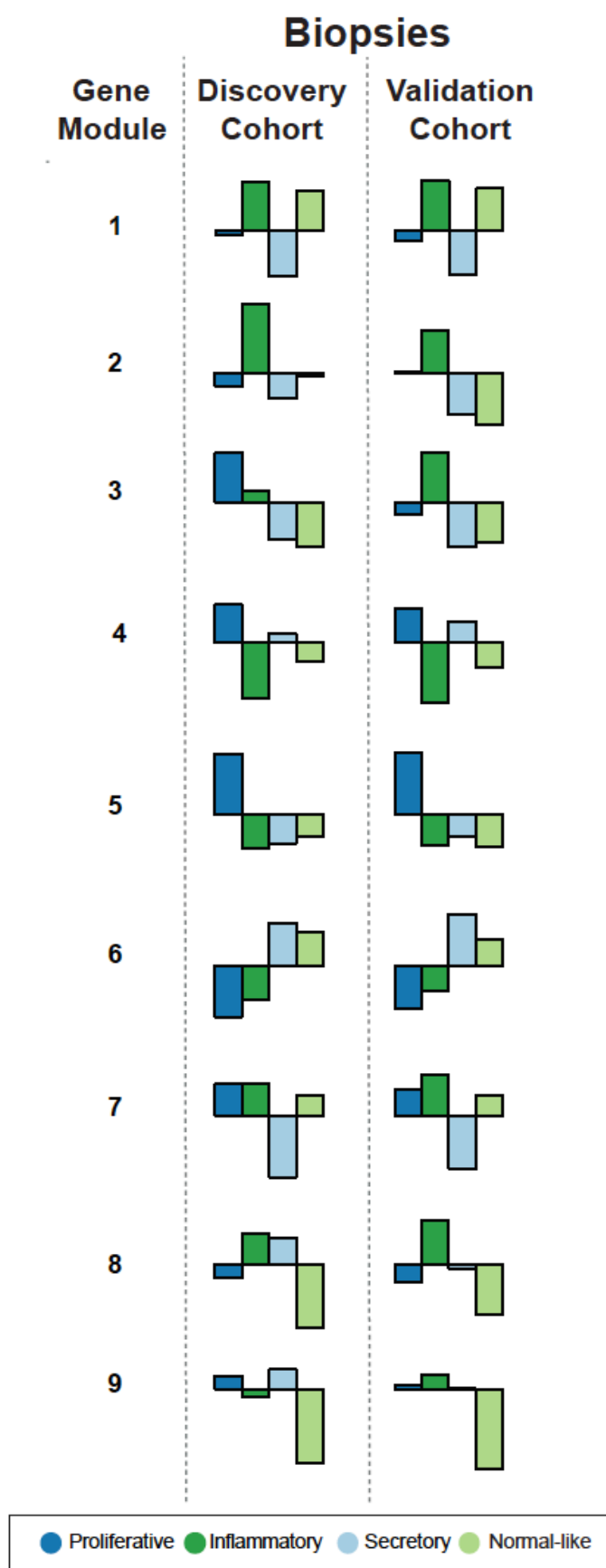
Supplementary Figure 4. Immunofluorescent staining quantitation of proliferation, basal cell, and ciliated cell markers across the molecular subtypes. Boxplot of immunofluorescent staining quantitation of KI67 (proliferation), KRT5 (basal cell) and TUB1A1 (ciliated cell) across representative samples from each molecular subtype (Proliferative n= 4, Inflammatory n=3, Secretory n=1, Normal-like n=1). KI67 and KRT5 staining are significantly higher in samples in the Proliferative subtype ($p=0.02$ and $p=0.01$ via linear models, respectively, for sample differences between the Proliferative subtype and other subtypes). TUB1A1 was lower in samples in the Proliferative and Inflammatory subtypes but did not reach statistical significance ($p=0.06$, linear model, for sample differences between Proliferative and Inflammatory subtypes versus Inflammatory and Secretory subtypes). The error bars represent the standard deviation. Source data are provided as a Source Data file.



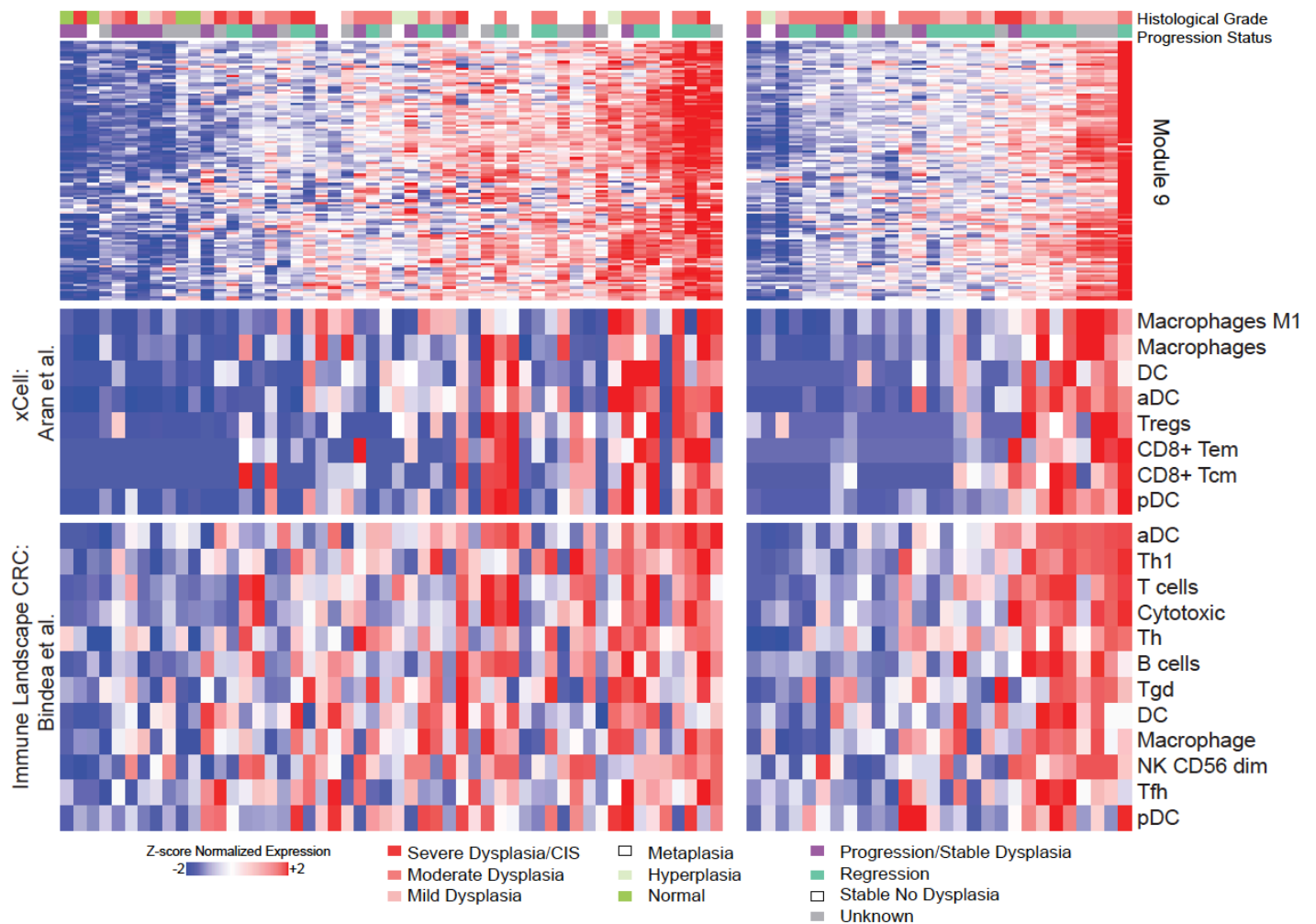
Supplementary Figure 5. Boxplots of select genes and cell type deconvolution results across the discovery and validation cohorts by molecular subtype. (A-D) Discovery cohort biopsies. **(E-H)** Validation cohort biopsies. **(A)** and **(E)** show boxplots of gene expression levels of LUSC driver genes identified by TCGA across the molecular subtypes. **(B)** and **(F)** show boxplots of gene expression levels of cell type marker genes across the molecular subtypes. **(C)** and **(G)** show boxplots of GSVA scores calculated using Dvorak *et al.* gene sets across the molecular subtypes. **(D)** and **(H)** show boxplots of ESTIMATE algorithm scores across the molecular subtypes. The ESTIMATE algorithm estimates the stromal (StromalScore), immune (ImmuneScore), and epithelial (ESTIMATEScore) cell fractions in each sample. High immune and stromal scores indicate a high fraction of stromal and immune cells while low epithelial scores indicate a high fraction of epithelial cells. In the boxplots, the upper and lower hinges correspond to the first and third quartile, center line represents the median, and whiskers extend from the hinge to the largest or smallest value at most 1.5 times the distance between the quartiles. Source data are provided as a Source Data file.



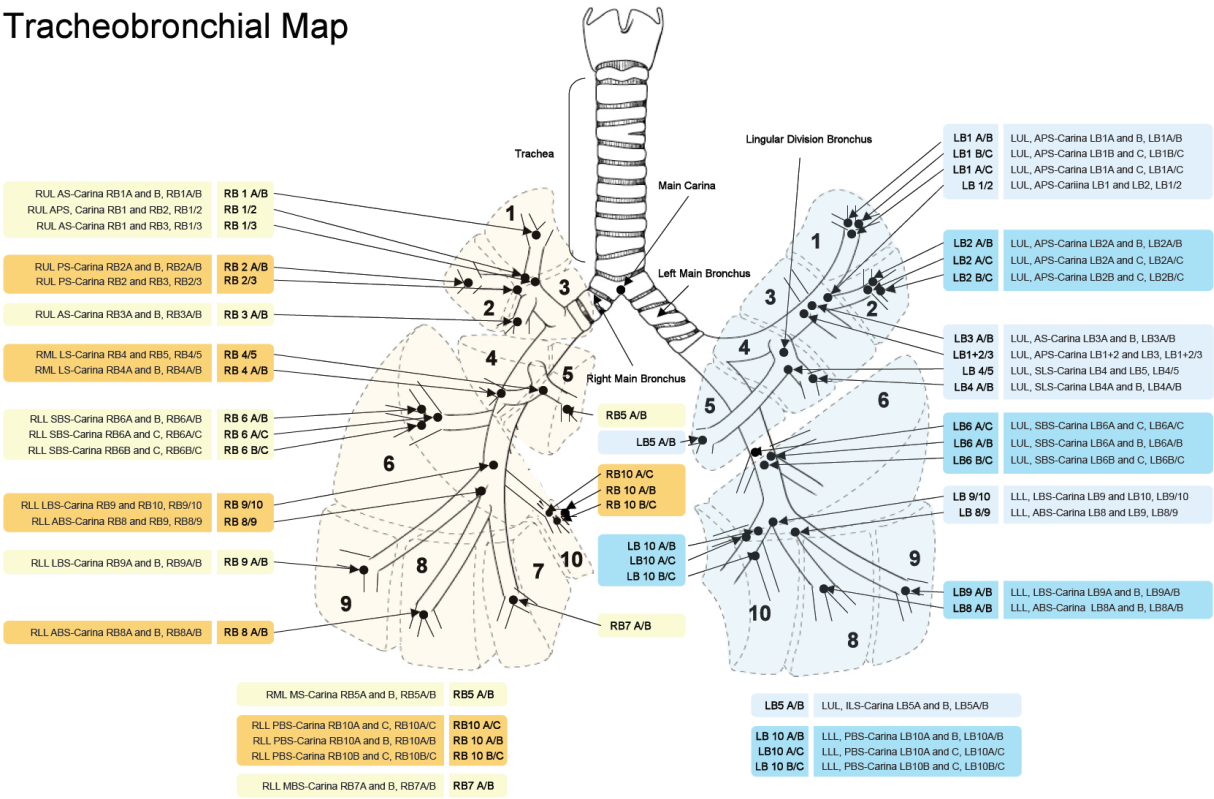
Supplementary Figure 6. Heatmap of the 22-gene molecular subtype classifier in the discovery and validation cohort biopsies. Semi-supervised hierarchal clustering of z-score normalized residual gene expression across the 22 classifier genes and 190 discovery cohort biopsies training samples (left) and the 105 Validation cohort biopsies (right). The rows of the heatmap show the gene module membership. The first column color bar shows molecular subtype membership in the discovery cohort and the 22-gene predicted molecular subtype membership in the validation cohort. The second column color bar depicts correct (dark red) and incorrect (pink) predictions in the discovery cohort using the 22-gene classifier and molecular subtypes (orange, Class 1; turquoise, Class 2; pink, Class 3; and purple/blue, Class 4) derived by performing consensus clustering across the validation cohort using n=3,936 genes. Source data are provided as a Source Data file.



Supplementary Figure 7. Gene module behavior across the molecular subtypes in the discovery and validation cohort biopsies. The mean GSVA score for each module is plotted for each molecular subtype. Source data are provided as a Source Data file.



Tracheobronchial Map



Supplementary Figure 9. Tracheobronchial map. The locations of the sites sampled by endobronchial biopsy.

Supplementary References

1. Dobin, A. *et al.* STAR: ultrafast universal RNA-seq aligner. *Bioinformatics* **29**, 15–21 (2013).
2. Li, B. & Dewey, C. N. RSEM: accurate transcript quantification from RNA-Seq data with or without a reference genome. *BMC Bioinformatics* **12**, 323 (2011).
3. Wang, L., Wang, S. & Li, W. RSeQC: quality control of RNA-seq experiments. *Bioinformatics* **28**, 2184–2185 (2012).
4. Robinson, M. D., McCarthy, D. J. & Smyth, G. K. edgeR: a Bioconductor package for differential expression analysis of digital gene expression data. *Bioinformatics* **26**, 139–140 (2010).
5. Campbell, J. D. *et al.* Distinct patterns of somatic genome alterations in lung adenocarcinomas and squamous cell carcinomas. *Nature Publishing Group* **48**, 607–616 (2016).
6. Beane, J. *et al.* Reversible and permanent effects of tobacco smoke exposure on airway epithelial gene expression. *Genome Biol.* **8**, R201 (2007).