



Published in final edited form as:

Nat Commun. ; 5: 5893. doi:10.1038/ncomms6893.

## Transcriptome Meta-Analysis of Lung Cancer Reveals Recurrent Aberrations in NRG1 and Hippo Pathway Genes

Saravana M. Dhanasekaran<sup>1,3,\*</sup>, O. Alejandro Balbin<sup>1,\*</sup>, Guoan Chen<sup>2,\*</sup>, Ernest Nadal<sup>2</sup>, Shanker Kalyana-Sundaram<sup>1</sup>, Jincheng Pan<sup>1</sup>, Brendan Veeneman<sup>1</sup>, Xuhong Cao<sup>1</sup>, Rohit Malik<sup>1</sup>, Pankaj Vats<sup>1</sup>, Rui Wang<sup>1</sup>, Stephanie Huang<sup>1</sup>, Jinjie Zhong<sup>7</sup>, Xiaojun Jing<sup>1</sup>, Matthew Iyer<sup>1</sup>, Yi-Mi Wu<sup>1</sup>, Paul W. Harms<sup>1,3,4</sup>, Jules Lin<sup>2</sup>, Rishindra Reddy<sup>2</sup>, Christine Brennan<sup>1</sup>, Nallasivam Palanisamy<sup>1</sup>, Andrew C. Chang<sup>2</sup>, Anna Truini<sup>8</sup>, Mauro Truini<sup>9</sup>, Dan R. Robinson<sup>1</sup>, David G. Beer<sup>2,#</sup>, and Arul M. Chinnaiyan<sup>1,3,5,6,#</sup>

<sup>1</sup>Michigan Center for Translational Pathology, University of Michigan Medical School, Ann Arbor, MI 48109, USA

<sup>2</sup>Thoracic Surgery, Department of Surgery, University of Michigan Medical School, Ann Arbor, Michigan

<sup>3</sup>Department of Pathology, University of Michigan Medical School, Ann Arbor, MI 48109, USA

<sup>4</sup>Department of Dermatology, University of Michigan Medical School, Ann Arbor, MI 48109, USA

<sup>5</sup>Howard Hughes Medical Institute, University of Michigan Medical School, Ann Arbor, MI 48109, USA

<sup>6</sup>Comprehensive Cancer Center, University of Michigan Medical School, Ann Arbor, MI 48109, USA

<sup>7</sup>Xinjiang Medical University, Xinjiang, 830011, China

<sup>8</sup>Lung Cancer Unit, IRCCS AOU San Martino-IST National Institute for Cancer Research, Genoa, Italy

<sup>9</sup>Department of Pathology, IRCCS AOU San Martino-IST National Institute for Cancer Research, Genoa, Italy

### Abstract

Lung cancer is emerging as a paradigm for disease molecular subtyping, facilitating targeted therapy based on driving somatic alterations. Here, we perform transcriptome analysis of 153 samples representing lung adenocarcinomas, squamous cell carcinomas, large cell lung cancer, adenoid cystic carcinomas and cell lines. By integrating our data with The Cancer Genome Atlas

Users may view, print, copy, and download text and data-mine the content in such documents, for the purposes of academic research, subject always to the full Conditions of use:[http://www.nature.com/authors/editorial\\_policies/license.html#terms](http://www.nature.com/authors/editorial_policies/license.html#terms)

#Corresponding Authors: Arul M. Chinnaiyan, M.D., Ph.D., Director, Michigan Center for Translational Pathology, Investigator, Howard Hughes Medical Institute, S. P. Hicks Endowed Professor of Pathology, American Cancer Society Professor, Professor of Urology, University of Michigan Medical School, 1400 E. Medical Center Dr. 5316 CCGC, Ann Arbor, MI 48109-5940, arul@umich.edu. David G. Beer, Ph.D., John and Carla Klein Professor of Thoracic Surgery, Department of Surgery, Professor of Radiation Oncology, University of Michigan Medical School, 1400 E. Medical Center Drive, 6304 UMCCC, Ann Arbor, MI 48109-5942, Phone: 1-734-763-0325; Fax: 734-615-6033; dgbeer@umich.edu.

\*These authors contributed equally to this article.

and published sources, we analyze 753 lung cancer samples for gene fusions and other transcriptomic alterations. We show that higher numbers of gene fusions is an independent prognostic factor for poor survival in lung cancer. Our analysis confirms the recently reported CD74-*NRG1* fusion and suggests that *NRG1*, *NF1* and Hippo pathway fusions may play important roles in tumors without known driver mutations. In addition, we observe exon skipping events in c-MET, which are attributable to splice site mutations. These classes of genetic aberrations may play a significant role in the genesis of lung cancers lacking known driver mutations.

---

Lung cancer is the leading cause of cancer-related deaths<sup>1, 2</sup> and is histologically classified as either non-small cell lung cancer (NSCLC) or small cell lung cancer (SCLC). NSCLC accounts for 80% of all lung cancers with lung adenocarcinoma (LUAD) and squamous cell carcinoma (LUSC) representing the major subtypes and large cell lung cancer (LCLC) and lung adenoid cystic carcinomas (LACC) the minor subtypes. LUAD are increasing in incidence worldwide<sup>3</sup>. Lung cancers poor overall 5-year survival rate (~15%) is primarily attributable to late diagnosis when curative surgery is no longer an option<sup>2</sup>.

Genomic analyses of LUAD have revealed mutations in many known oncogenes and tumor suppressor genes including *KRAS*, *EGFR*, *TP53*, *CDKN2A*, and *STK11*<sup>4</sup>. These tumors also harbor low frequency copy number alterations including *ERBB2* amplification which is targetable with herceptin<sup>5</sup>. Alterations in oncogenes such as *KRAS*, *EGFR*, *ALK* and *MET* influence tumor formation and maintenance and are considered “drivers” in a subset of NSCLCs yet in a substantial patient population the driver aberrations are yet to be identified (i.e., “driver mutation unknown”)<sup>6</sup>. Recent analyses by The Cancer Genome Atlas (TCGA) of both LUSC<sup>7</sup> and LUAD<sup>8</sup> revealed recurrent mutations and recurrent copy number alterations in genes that present in both subtypes and also specific to each. The histologic and molecular heterogeneity observed in lung cancer underscores the difficulties in developing effective therapies for patients.

Patients with *EGFR* mutations show responsiveness to EGFR inhibitors which are often not durable<sup>9</sup>. In addition to driver somatic gene mutations, oncogenic gene fusions including the *EML4-ALK* fusion gene have been identified in approximately 4% of LUAD<sup>10</sup>. This fusion protein links the N-terminal portion of echinoderm microtubule-associated protein-like 4 (*EML4*) with the intracellular signaling portion of a receptor tyrosine kinase, the anaplastic lymphoma kinase (*ALK*). The *EML4-ALK* translocation is mutually exclusive with *EGFR* and *KRAS* mutations, an indicator of therapeutic responsiveness to *ALK* inhibitors,<sup>10</sup> and tumors with this translocation also have fewer *TP53* gene mutations<sup>11</sup>. Additional gene fusion events have now been identified in LUAD including *KIF5B-ALK*<sup>12</sup>, *ROS1*<sup>13</sup> and *RET*<sup>14, 15</sup> gene fusions. *KIF5B-ALK* fusion-positive lung cancers may respond to *ALK* inhibitors, whereas *RET* fusions may be treated using drugs that target this kinase<sup>16</sup>. We previously identified *NFE2* and *FGFR3* gene fusions in a subset of lung cancers<sup>17, 18</sup>.

In this study, we perform transcriptome meta-analysis on a data compendium assembled by combining 153 primary NSCLCs that we sequenced with 521 NSCLCs from the TCGA and 79 samples from a published report<sup>19</sup>. The highly heterogeneous lung cancer gene fusion landscape is dominated by low recurrence and private fusions. We demonstrate that the

number of fusions in a sample is an independent prognostic factor for poor survival. We found gene fusions affecting core members of the Hippo pathway, Neurofibromatosis 1 (NF1), and Neuregulin 1 (NRG1) genes, besides the recently reported CD74-NRG1 fusion variant<sup>20, 21, 22</sup> and c-MET exon skipping event<sup>23</sup>. Upon integrating fusion, mutation and outlier expression data these events collectively account for ~16% of driver negative lung cancer samples.

## Results

### Analysis Work Flow and Mutation Landscape of NSCLC Subtypes

We sequenced mRNA from 153 samples representing major (LUAD, LUSC) and minor (LULC, LACC) subtypes of NSCLC using strand-specific, RNA paired-end sequencing (RNASeq). Our “UMICH cohort”, samples included 67 LUAD, 36 LUSC (64 stage I, 17 stage II and 22 stage III patients), 9 LCLC, 11 LACC, 24 lung cancer cell lines and 6 matched nonmalignant lung samples. Eighty-two patients were heavy smokers (>20 pack years), 13 were light-smokers (defined by <20 pack years) and smoking status of 15 patients was unknown (Supplementary Table 1). The median smoking pack years was 45 (range, 2 – 300). The average follow-up was 5.05 years. Sample acquisition details are provided in the **Methods** section. To increase the power of our analysis and to discover recurrent fusions, we included two publically available NSCLC datasets from TCGA and Korean LUAD (SEOUL cohort) studies<sup>19</sup> and assembled a RNASeq cohort that totaled 753 patient tumors. The TCGA cohort included 305 LUAD and 216 LUSC samples (250 stage I, 112 stage II, 101 stage III, and 19 stage IV cases and 39 with unknown stage).

The combined cohort included 451 LUAD, 251 LUSC, 9 LCLC, 11 LACC and 24 NSCLC cell lines, making this the most comprehensive RNA-sequencing cohort of lung cancers assembled to date. A description of the cohort assembly and sample clinical-pathological information is presented in the **Methods** section and summarized in Supplementary Table 1. The available clinical information including smoking history is presented in Supplementary Data 1.

We developed the analysis pipeline, depicted in Supplementary Fig. 1 thus assessing gene fusions among all 753 patients in the combined cohort and for integration with mutation and clinical information (see **Methods** for details). For each sample we determined the mutation status of oncogenes and tumor suppressors known to play a role in lung cancer<sup>6</sup> and reflected the previously reported mutational landscape of LUAD and LUSC (Fig. 1)<sup>4, 5, 7</sup>. *KRAS* was mutated in 30.1% and 1.6% of LUAD and LUSC respectively; *EGFR* in 13% and 1.6% of LUAD and LUSC; *BRAF* in 8% and 3.2% of LUAD and LUSC and *PIK3CA* in 7.6% and 13.5% of LUAD and LUSC respectively. As previously reported<sup>4, 5, 7</sup>, *TP53* mutations are common in both LUAD and LUSC patients, 50.3% and 65.7% respectively (Fig. 1). The mutations identified among select genes in the characterized cell lines are summarized in Supplementary Fig. 2.

In addition to the major NSCLC subtypes, we profiled 9 LCLC and 11 LACC, also called lung colloid carcinoma, a rare subtype. In LCLC we found 1 sample with *KRAS* activating mutation, 3 with *TP53* missense mutations and 4 without mutations in known lung cancer

genes (Supplementary Table 2). The pattern observed in LCLC is consistent with a recent report<sup>24</sup> supporting their reclassification into either LUAD or LUSC based on shared genetic aberrations.

In LACC, despite the small sample size, we observed a higher frequency of *RAS/RAF* pathway mutations (72%, 8/11) compared to the major NSCLC subtypes (Supplementary Table 2). The mutations were mutually exclusive, where five samples with *KRAS* mutations had *KRAS*<sup>G12C, G12V, G13C, G12D, Q61H</sup> variants respectively, while *BRAF*<sup>V600E</sup>, *HRAS*<sup>Q61L</sup>, and *NRAS*<sup>Q61R</sup> were observed in three independent samples. Interestingly, the samples with *NRAS*<sup>Q61R</sup> and *KRAS*<sup>G13C</sup> also had mutations in *TP53*<sup>R141C, R141L</sup> and *KIT*<sup>M537L</sup>. *MET*<sup>T1010I</sup> variation was also observed in the *NRAS* mutated sample. While 2 LACC samples had no mutations reported in COSMIC, one sample harbored an *IDH1*<sup>V178I</sup> variant. Interestingly, *MYB-NF1* gene fusions were absent in the lung ACC, unlike the salivary gland-ACC where it occurs in 57% of cases<sup>25</sup>. Likewise, *KRAS* mutations were common in LACC, but none were detected in the 60 salivary-ACCs sequenced recently<sup>25</sup>. However, in the salivary-ACC cohort a potential driver *HRAS* non-synonymous mutation was noted to be mutually exclusive with *MYB* gene fusion. Another recent study identified activating *BRAF* and *HRAS* mutations in breast adenoid cystic carcinoma (BACC) samples that were a distinct subset of triple negative breast cancers<sup>26, 27</sup>. Hence the previous report on BACC and our results here on LACC have identified distinct ACC subsets that harbor activating RAS/RAF mutations but lack MYB fusions that are primarily found in head and neck ACC, revealing differences in underlying molecular events despite histological similarities. Due to the small cohort size and lack of significant fusion events in LCLC and LACC samples, these cohorts were excluded from the fusion analysis presented below.

### NSCLC Fusion Landscape

In order to generate comparable results across samples from different cohorts, we developed a consistent data-driven gene fusion prediction pipeline and analysis workflow shown in Supplementary Fig. 1 (also see **Methods**). We detected 6,348 unique fusions among the 733 samples for an average of 13 fusions per tumor sample (range: 0–67). Although both LUAD and LUSC had a comparably high single nucleotide mutation rate of 8.1 mutations/Mb<sup>5, 7</sup> they differed in the average number of fusions per sample with 11 fusions in LUAD, and 17 in LUSC (Student's t-test  $P < 2.2e-16$ ). We did not observe a statistically significant difference in average number of fusions among heavy and light smokers (LUAD Student's t-test  $P = 0.06$ ; LUSC Student's t-test  $P = 0.59$ ) among different clinical stages and regardless of the tumor type (Supplementary Table 3). Tumors with missense or nonsense mutations in *TP53* showed greater average number of fusions compared to samples with wild-type *TP53* (Supplementary Fig. 3A and 3B,  $P = 0.001$ ). As most LUSC have somatic mutations in *TP53*<sup>7</sup>, this difference is consistent with the average number of fusions between LUAD and LUSC samples. In LUAD, we observed a significant correlation between the presence of oncogenic mutations (e.g., *KRAS* activating mutations) and *TP53* deleterious mutations (stop codon or splice site mutations) and the number of fusions (Fisher's exact test  $P = 0.008$ ). We could not determine whether a similar correlation exists in LUSC due to the low incidence of mutations in *KRAS*, *EGFR* or other oncogenes in the samples.

## Number of Fusions is Associated with Prognosis

We investigated the relationship between the number of fusions present in a tumor and patient prognosis. Patients in our combined cohort were first classified into three fusion categories based on distribution percentiles as low (0–7), intermediate (8–17), or high (18) and then a 10-year Kaplan-Meier survival analysis was performed. Patients with high number of fusions had significantly shorter median overall survival (35.6, 95% confidence interval (CI) 27.2–43.9) compared to patients with intermediate (49.5, 95% CI 23.9–75.1) or low number of fusions (62.3, 95% CI 44.6–80.1; Likelihood ratio test  $P=0.008$  Fig. 2). We observed similar results both for LUAD and LUSC when analyzed independently (Supplementary Fig. 4A and 4B). Statistically significant clinical covariates in the univariate Cox model (Supplementary Table 4) were used in the multivariate analysis to examine the prognostic value of fusion number. Strikingly, a high fusion incidence was independently associated with worse overall survival (HR=1.56, 95% CI 1.13–2.15,  $P=0.007$ , Supplementary Table 5) after adjusting for gender and disease stage. When *TP53*, *KRAS* and *EGFR* mutation status or smoking status was included in the multivariate analysis, the number of fusions remained independently associated with poor outcome (Supplementary Table 6).

## Private or Low Recurrence Fusions in Lung Cancer

In order to filter the fusion data and prioritize fusion candidates, we developed a random forest fusion classifier (**See Methods**). This classifier uses structural and functional annotation features of each fusion in order to prioritize gene fusion candidates involving exonic regions. Remarkably, our classifier had a true positive recovery rate greater than 90% in two independent validation datasets and automatically recapitulated the intuitive knowledge about the important structural properties defining *bona fide* fusions (Supplementary Data 2 and 3). In our fusion dataset, the top five features contributing to the fusion classifier were, in decreasing order of importance: fusion type (inter-chromosomal, intra-chromosomal, tandem-duplication), sum of the median alignment quality of reads supporting the fusions, number of spanning and encompassing reads across the fusion junction and the cohort normalized expression value for the 3'-partner gene (Supplementary Fig. 5).

Using this classifier, 422 fusions were shortlisted from the entire cohort (Supplementary Data 4). Sixty-four out of 422 fusions (15%) involved kinases (either as 3' or 5'-partner) including the known *ROS1*, *RET* and *ALK* fusions: 52 fusions involved oncogenes and 63 involved tumor suppressors (Supplementary Data 4). Moreover, of the fusions involving “informative genes”, we found 61 productive in-frame fusions, 63 out-of-frame fusions and 6 promoter fusions.

In the *KRAS* mutant population, a large NSCLC molecular subtype where chemotherapy is the only approved treatment we identified additional private fusions (Supplementary Data 4). For example, sample pt\_lung\_A25 contains a driver *KRAS*<sup>G12C</sup> and *TP53*<sup>P72R</sup> mutations in addition to fusions with abundant read support. While the TRAF interacting protein-Inositol Hexakisphosphate Kinase (*TRAIIP-IP6K1*) fusion results in loss of function of both partners, the *SLC12A7-TERT* fusion produces an in-frame ORF where, the telomerase

domain of *TERT* is retained and could serve as a potential combinatorial drug target. In another sample, pt\_lung\_A63, that harbored *KRAS*<sup>G12D</sup>, *TP53*<sup>P72R</sup> and *ATM*<sup>E2423K</sup> mutations, in addition has a *TSC1-SMARCA4* fusion. Further pt\_lung\_C028 with *TP53*<sup>R248L</sup> and *SMARCA4*<sup>E1056stop</sup> mutations also harbored a *WASF2-FGR* fusion where the kinase domain of FGR is retained. These three cases are representative examples of private fusions and the additional events that coexist in NSCLC tumors.

As our cohort was large enough, we estimated the recurrence of different gene fusions that we classified as molecular, functional or family recurrence. Molecular recurrence defined as the same 5' and 3' partners observed in different samples such as *SLC34A2-ROS1*; functional recurrence refers to when either 5' or 3' partner is the same (*CCDC6-RET* and *KIF5B-RET*); and gene family recurrence correspond to gene fusions in which 5' or 3' partners belongs to the same gene family such as *FGFR* (*FGFR3-TACC3*, *FGFR2-CCDC6*, *BAG4-FGFR1*). Functionally-recurrent kinase fusions *ROS1*, *RET* and *ALK* were found in 0.86%, 0.29%, and 0.14% across the combined cohort (Supplementary Data 5).

Interestingly, in tumors with known driver fusions the number of “classified” fusions is lower than those without driver fusions (Student t-test  $t = 2.7588$ ,  $df = 5.023$ ,  $P = 0.01985$ ) suggesting their functional importance. Similarly, *BCAS3-MAP3K3*, *MRC2-MAP3K3* is another example of family recurrence. We observed “pathway fusion recurrence” in which multiple genes in the same signaling pathway are involved in fusions. Interestingly, 10 out of 33 members of the Hippo<sup>28</sup> pathway were identified as fusion partners (Supplementary Data 6).

### Perturbation of the HIPPO Pathway in Lung Cancer

The Hippo signaling pathway is highly conserved across species and plays a major role in cell polarity, cell-cell adhesion and contact inhibition<sup>29</sup>. The mammalian homologues of the *Drosophila* Hippo and Warts core serine threonine kinases are STE20-like protein kinase (*MST1/2*) and large tumor suppressor homolog kinase (*LATS1/2*) respectively. The core kinases regulate the activity and stability of the transcriptional co-activators yes-associated protein 1 (*YAP1*) and WW domain-containing transcription regulator 1 (*WWTR1*) through phosphorylation. Un-phosphorylated *YAP/WWTR1* binds to TEA domain family (*TEAD*) transcription factors in the nucleus to regulate gene expression (Fig. 3A). Accessory members of the Hippo pathway such as *KIBRA* (*WWCI*), scribbled planar cell polarity (*SCRIB*) and Neurofibromin 2 (*NF2*) have been shown to activate the core kinases. An increasing number of studies have investigated the Hippo pathway in lung, colorectal, ovarian and liver cancers<sup>29</sup>. While animal model experiments support the Hippo pathway in tumorigenesis, no evidence for non-synonymous mutations in this pathway has been found in lung cancer. Few somatic or germ-line mutations discovered in Hippo pathway genes are found in common human cancers, with *NF2* being the only gene known to be inactivated by mutation<sup>29</sup>. We observed novel recurrent *NF2* fusions, where retention of only the first exon of *NF2*, in both *NF2-OSBP2* and *NF2-MORC2* fusions result in loss of function of this tumor suppressor gene (Fig. 3B and 3C) and several fusions involving core members of the Hippo pathway such as *LATS1*, *YAP1*, and *WWTR1* (previously known as *TAZ*) (Fig. 3C). We also identified fusions in associate members of the Hippo pathway including, *HIPK2*, *TAOK1*, *TAOK3*, *FAT1*, *DCHS2* and *PTPN14* (Fig. 3B and 3C). Detailed inspection of the

fusions revealed two intriguing aspects of these aberrations. Gene fusions in Hippo pathway tumor suppressor members such as *LATS1*, *DCHS2*, *FAT1*, *TAOK1*, *TAOK3*, *PTPN14* and *NF2* (Fig. 3B) abrogated their function by generating truncated proteins. However, fusions involving oncogenic proteins in the Hippo pathway such as *WWTR1*, *YAP1* and *HIPK2* retained their crucial functional domains Fig. 3B). Furthermore, we investigated the presence of additional genetic aberration in the index fusion samples and noticed that the vast majority lack known driver mutations (10 out of 14) (Supplementary Data 6). Using cBioportal (<http://www.cbioportal.org>) we discovered copy number loss and associated low mRNA expression of *FAT1* in the index fusion sample (Supplementary Fig S6A) and copy gain and elevated expression of *YAP1* in the sample harboring *YAP1* fusion (Supplementary Fig S6B). These observations suggest that gene fusions are a novel mechanism of altering Hippo pathway genes potentially promoting a transforming phenotype. Taken together, the fusion landscape in lung cancer is highly heterogeneous and characterized by low recurrence and private fusions (Supplementary Data 5). Despite this heterogeneity, gene fusions could still be functionally relevant in lung cancers by affecting several members of common pathways such as those of the Hippo signaling cascade we observed here.

### Inactivating fusions of *NF1* in lung cancer

Next, our integrative analysis combining fusion and mutation status revealed a total of 33 samples with aberrations in *NF1* gene such as truncating fusions: *GOSR1-NF1*, *NLK-NF1* and *NF1-PSMD11* or deleterious mutations: non-sense, frame shift or splice site (Fig. 1, Fig. 4A and Supplementary Table 7). The fusions and mutations were observed in both LUAD and LUSC predominantly in driver negative samples (27 out of 33). Loss of *NF1* promotes cell proliferation by de-repressing the mTOR pathway in a RAS-, PI3K-dependent fashion<sup>30, 31</sup>. The fusion architecture renders the tumor suppressor *NF1* inactive by either truncating ORFs (*GOSR1-NF1*, *NLK-NF1*) or by destroying its functional domains (*NF1-PSMD11*), (Fig. 4A and 4B) indicating an alternate mechanism for *NF1* inactivation in lung cancers besides somatic mutations<sup>4</sup>. To assess additional *NF1* destructive fusions in lung cancer we did a comprehensive analysis assessing fusion junctions involving either exons or introns and found two additional events of *NF1-DRG2\_Antisense* and *NF1-MYO15A\_Antisense* present in the LS2 sample (Fig. 4A and 4B). The read evidence suggests genomic deletion as the mechanism for the *NF1* fusions except in sample LS2 where centromeric inversion may be the underlying aberration (Fig. 4B). Importantly, 20 out of 29 mutated *NF1* samples and all *NF1* truncating fusions were observed in samples without known drivers, accounting for 6.2% (24/386) of this subpopulation. Interestingly, two samples had fusions accompanying somatic mutations in *NF1* potentially altering both the alleles of this tumor suppressor gene (Supplementary Table 7).

### Exon skipping and coincident splice site mutations in *c-MET*

Recently, a significant percent of driver unknown lung cancer samples have been shown to harbor fusions involving *ALK*, *ROS1*, *RET*<sup>19, 23</sup> kinases and an activating exon skipping in the *c-MET* oncogene<sup>23</sup>. Our analysis revealed that 1.3, 0.52 and 0.26 percent fusions involving *ROS1*, *RET* and *ALK* respectively among LUAD and LUSC with unknown driver. We detected *c-MET* exon-14 skipping in 15 samples, 14 of which occurred in driver unknown samples, a 3.6% (14/386) recurrence rate in this subpopulation (Fig. 5).

Importantly, in 5 out of 15 samples the skipping of *c-MET* exon-14 is likely caused by a mutation affecting the splice donor site adjacent to the amino acid position D1010 as previously described<sup>32</sup>. Our RNASeq data also validated the reported *c-MET* exon skipping event in the H596 cell line<sup>23</sup>.

### Outlier kinase expression in lung cancer

Next, integrative analysis combining the mutation, fusion and gene expression data revealed outlier expression information in the context of fusions and mutations per sample. Focusing on kinase genes for example *ROS1*, we noticed 6 samples across the combined cohort with outlier *ROS1* expression that lacked any evidence for *ROS1* fusions. A similar phenomenon was also observed in cases with *FGFR3* outlier expression. Intriguingly, tumors showing outlier expression of *ROS1* and *FGFR3* are almost exclusively driver unknown samples without evidence of fusions (Fisher exact test  $P=0.004$  and  $P=0.086$  respectively, Fig. 1). Fluorescence *in-situ* hybridization analysis of *ROS1* (n=1) and *RET* (n=3) outlier index cases did not detect any gene rearrangements. Hence while the mechanism of overexpression remains to be determined, the outlier kinase expression may act as oncogenic drivers and be potentially actionable.

### Recurrent *NRG1* rearrangements with novel fusion partners in lung cancer

Remarkably, we noted functionally recurrent gene fusion where the common 3' gene neuregulin 1 (*NRG1*) was fused to various 5' partners (Fig. 6A and Supplementary Table 8) *CD74-NRG1*, *RBPMS-NRG1*, *WRN-NRG1* and *SDC4-NRG1*, in both LUAD and LUSC samples. Importantly, *CD74-NRG1* fusion variant was recently identified by three independent groups<sup>20, 21, 22</sup>. While *CD74-NRG1*, *SDC4-NRG1* and *RBPMS-NRG1* fusion events resulted in the production of chimeric proteins, the *WRN-NRG1* fusion results in the overexpression of full length *NRG1* regulated by the *WRN* gene promoter. As a member of EGF ligand family, *NRG1* transduces its signal through the HER/ErbB family receptor tyrosine kinases<sup>33, 34</sup>. *NRG1* functional domains include kringle-like, immunoglobulin-like domain and the EGF domain located in the C-terminal region<sup>33</sup>. Notably the EGF domain is essential for receptor interaction<sup>35</sup> and preserved in all the *NRG1* fusions identified (Fig. 6A). All *NRG1* fusion index samples were found in samples without known driver mutations and displayed *NRG1* outlier expression in the tumor but not matching normal tissue (Fig. 6B and 6C). Strikingly similar to the pattern described above for the known receptor kinases fusions, we noticed *NRG1* outlier expression in both index fusion samples (n=4) and an independent set of known driver aberration negative cases (n=10) (Supplementary Table 8). Among the lung cancer cell line RNASeq data, H1793 exhibited the highest *NRG1* transcript expression (Fig. 6D and Supplementary Fig. 2). At 70% knock down with two independent *NRG1* siRNAs (Fig. 6E) H1793 cell proliferation rate was affected as assessed using cell growth assays (Fig. 6F). Conversely upon stable overexpression of the *CD74-NRG1* fusion protein in normal lung BEAS-2B cells we observed significant increase in cell proliferation, migration (Fig. 6G and Supplementary Fig. 7A) and an altered morphology relative to LacZ controls (Supplementary Fig. 7B and 7C). *CD74-NRG1* overexpression induces epithelial to mesenchymal transition (EMT) as evidenced by increased VIM and SNAIL protein expression and decreased CDH1 level by Western blot analysis (Supplementary Fig. 7D and 9). We next performed gene expression profiling of *CD74-NRG1* and LacZ control cells to



identify affected biological pathways. Significant analysis of microarrays (SAM) showed overexpression of several EMT markers such as *VIM*, *ZEB1*, *ZEB2*, *FZD7*, *TWIST1*, *VCAN*, and *CHD2* and under-expression of *RGS2* and *CDH1* among others further supporting the EMT phenotype in *CD74-NRG1*-positive cells (Supplemental Data 7). Vimentin, *ZEB1* and *ZEB2* were overexpressed more than 4-fold, while *CDH1* and *RGS2* were among the most under-expressed genes (Supplementary Fig. 7E and 8A). Gene set enrichment analysis (GSEA) identified down-regulation of cell adhesion (Supplementary Fig. 8B) and up-regulation of the SRC and ERBB pathways (Supplementary Fig. 8C and 8D) in *CD74-NRG1* cells. We examined both total and phosphorylated ERBB3, a receptor known to bind NRG1, and observed a substantial decrease in total ERBB3 upon overexpression of *CD74-NRG1* which was also reflected in its phosphorylated form as compared to LacZ control (Supplementary Fig. 8E and 9). Despite the observed decrease in total ERBB3 in the fusion expressing cells, phospho-ERBB3 was still detectable (Supplementary Fig. 8E and 9). Total ERBB3 decrease upon exposure to NRG1 has been previously demonstrated in MCF-7<sup>36</sup> and also in H568 lung cells upon *CD74-NRG1* overexpression<sup>20</sup>. In addition we observed increased levels of phospho-ERK (1.95-fold) and phosphoJNK1 (5.5-fold) relative to LacZ control (Supplementary Fig. 8E and 9) potentially promoting the oncogenic phenotype in NRG1 fusion overexpressing cells. Finally, we examined other cancer types for *NRG1* fusions and discovered one additional *RAB21L1-NRG1* fusion in the TCGA ovarian cancer RNASeq data. As observed in lung cancer, the functional EGF domain is retained in *RAB21L1-NRG1* and the fusion index case exhibited outlier *NRG1* expression (Fig. 6A and 6B). Altogether, *NRG1* is perturbed (*NRG1* fusions and/or outlier expression) in 3.9% (15/386) of driver unknown samples, supporting a causal role for NRG1 in this lung cancer patient subpopulation.

## Discussion

Increased understanding of lung cancer has resulted in the identification of therapeutic molecular targets and development of relevant targeted therapies. For example, *EGFR* activating mutations in exons 18, 19 and 21 are now routinely assessed in tumor biopsies prior to treatment with gefitinib or erlotinib; the response rate is nearly 70% in mutation-positive advanced NSCLC<sup>37</sup>. Further, fusions involving *ROS1*, *ALK* and *RET*<sup>15, 16, 38</sup> tyrosine kinases are identified primarily in younger patients with LUAD and without known driver mutations or significant smoking history. Despite the low fusion frequency, clinical trials for ALK-positive lung cancer patients have shown higher response rate and longer progression-free survival when treated with crizotinib, a drug targeting ALK, relative to chemotherapy<sup>39, 40</sup>. These results support targeting specific molecular aberrations in patients' tumors.

In this study, RNA sequencing was used to characterize the fusion landscape of NSCLC in an unbiased fashion. We find the fusion landscape highly heterogeneous dominated by private and low recurrence fusions; with a greater number of fusions per sample detected in LUSC than LUAD on average (Student t-test,  $P < 2.2e-16$ ). No statistically significant difference, with respect to any other clinical characteristics such as smoking history or disease stage, was observed (Supplementary Tables 3, 4 and 5). Importantly, a higher number of fusions were independently associated with poor overall survival (Fig. 2,

Supplementary Table 5), after adjusting for histological subtype, age, gender, disease stage and *TP53*, *KRAS* and *EGFR* mutation status (Supplementary Table 6). As RNA sequencing becomes widely adopted for profiling transcript expression and gene fusion detection, our results suggest that the number of fusions could also be used as an independent prognostic marker in lung cancers.

Our analysis of functionally recurrent fusions identified aberrations in multiple members of the Hippo pathway. This evolutionarily-conserved pathway regulates tissue growth and cell fate and has been thought to play an important role in cancer<sup>28</sup>. Functional studies conducted in mouse models showed that knock down of tumor suppressor or overexpression of oncogene members of the pathway induced tumor formation<sup>29</sup>. Furthermore, two recent reports identified recurrent fusions involving *WWTR1*, an oncogene member of the Hippo pathway and *CAMTA1* in epithelioid hemangioendothelioma<sup>41, 42</sup>. The previously reported *WWTR1* fusion and the one in our study (*WWTR1-SLC9A9*) (Fig. 3) have identical *WWTR1* gene breakpoints, whereby the functional WW domain of WWTR1 is retained in both fusion events. We also observed fusions involving 3 out of 13 core members and 7 out of 20 associate members of the Hippo pathway (Fig. 3). A recent study has demonstrated the role of STK11 (also called LKB1) in regulating the core Hippo kinases through Scribble<sup>43</sup>. The tumor suppressor STK11 is frequently inactivated in lung cancer (Fig. 1) which is associated with YAP activation. This discovery now vastly expands the incidence of Hippo pathway aberration in lung cancers. Interestingly, gene fusions in Hippo pathway tumor suppressor members appear to abrogate their function by generating truncated proteins, while fusions involving oncogenic proteins in the Hippo pathway retain their crucial functional domains (Fig. 3). Taken together, our data now present novel evidence for the involvement of the Hippo pathway in lung cancer.

The recurrent tyrosine kinase fusions mentioned earlier were found almost exclusively in LUAD not harboring known fusions and have not been previously identified in the LUSC subtype. Here, we observed a recurrent fusion with *NRG1* as 3' partner (*CD74-NRG1*, *RBPMS-NRG1* and *WRN-NRG1*) in both LUAD and LUSC (Fig. 6). *NRG1*, a growth factor that interacts with the HER/ErbB receptor tyrosine kinases, is expressed in a subset of cancers, including breast, lung and others<sup>44</sup>. *CD74* is a known 5' fusion partner of *ROS1* kinase in lung cancer. While *CD74-NRG1* and *WRN-NRG1* fusions contain the signal peptide and type II transmembrane domain required for *NRG1* localization to the plasma membrane, cellular location of *RAB21LI-NRG1* and *RBPMS-NRG1* fusion proteins is uncertain. However, of the 20 *NRG* transcript variants (transcribed from *NRG1-4*) reported, several lack the N-terminal signal sequence required for membrane localization and transport to the extracellular space. In these instances an internal hydrophobic amino acid stretch is speculated to substitute for the N-terminal signal sequence<sup>33, 35</sup>. Additionally, we identified a novel *SDC4-NRG1* fusion in two samples added to the TCGA cohort after our data freeze. The *SDC4-NRG1* fusion produces a secretory *NRG1* protein due to the signal peptide contributed by *SDC4* protein. This observation suggests that incidence of *NRG1* aberrations in lung cancer is likely to increase as more samples are characterized.

Remarkably, *NRG1* fusions are present in tumors without known driver events (Fig. 1 and Supplementary Table 8) and the index samples display outlier *NRG1* expression (Fig. 6),

similar to oncogenic fusions such as *ROS1*. Moreover, we found additional cases of *NRG1* outlier expression in samples without known driver mutations, suggesting a potential role for *NRG1* in those samples. We demonstrated that abrogating *NRG1* expression affects cell proliferation (Fig. 6) and more importantly we showed that human bronchial cells stably expressing *CD74-NRG1* promoted proliferation and migration (Fig. 6). Three independent studies have very recently associated *CD74-NRG1* fusions with mucinous LUAD subtype<sup>20, 21, 22</sup>. We further examined our samples and discovered that *HNF4A*, a recently characterized biomarker for mucinous LUAD<sup>45</sup>, showed highest expression in our *CD74-NRG1* index case, providing independent support for association of *NRG1* gene fusions with mucinous LUAD. Interestingly, the *SDC4-NRG1* index sample with the highest *NRG1* outlier expression (Fig 6B, *NRG1* expression: 380FPKM, higher than the cell line H1793) did not show high *HNF4A* expression, suggesting that *NRG1* fusions with partners other than *CD74* are perhaps more prevalent in non-mucinous LUAD. *NRG1* rearrangements have also been detected using FISH in breast cancer cell lines<sup>46</sup>. Moreover, *NRG1* over-expression was recently demonstrated in a subset of breast clinical tumor samples, and was mutually exclusive with *HER2* mutations<sup>47</sup>. These observations together with our results from lung and ovarian cancers suggest that *NRG1* rearrangements are recurrent and likely drivers of various cancers types.

The therapeutic targeting of NRG1-ERBB autocrine loop was previously suggested<sup>48</sup> and recently blocking NRG1 and other ligand-mediated HER4 signaling shown to enhance the magnitude and duration of the chemotherapeutic response of NSCLC<sup>49</sup>. Therefore, the characterization of all *NRG1* fusions presented in this study, as well as the common signaling pathways activated in both fusion and outlier expression index samples could further elucidate NRG1 mechanism of action and reveal further therapeutic opportunities.

Our integrative analysis combining mutation and fusion status extended previous observations of *c-MET* exon skipping and *NF1* truncating mutations. We detect novel truncating fusions involving several tumor suppressor genes such as *NF1*, *NF2*, *TP53* (data not shown), *LATS1*, *DCHS2*, *FAT1*, *SMARCA4*, *TAOK1* and *TAOK3* among others. These results highlight gene fusions as potentially common and a previously underappreciated mechanism for loss of function of many tumor suppressor genes. In summary, Hippo pathway fusions (2.6%), *NRG1* fusion/outlier expression (3.9%), *NF1* truncating mutations/fusions (6.2%) and *c-MET* exon skipping (3.6%) account for ~16% of driver-unknown lung cancer cases, and expanding the repertoire of lung cancer molecular subtypes. The previously documented success of targeted therapies against low recurrence oncogenic fusions and the heterogeneity of the fusion landscape, demonstrated in this study, reinforce the demand for personalized molecularly targeted drug therapies in lung cancer.

## Methods

### Sample Acquisition and Total RNA isolation

We collected tumor samples from 67 LUAD, 36 LUSC, 9 LCLC patients along with 6 matched normal lung tissues samples following surgery at the University of Michigan. The recruitment of subjects and informed consent were reviewed and approved by our Institutional Review Board. The publically available dataset from TCGA was downloaded

using the TCGA portal and the Seoul data from dbGAP. Formalin-fixed, paraffin-embedded (FFPE) sections from 11 adenoid cystic carcinoma samples were from IRCCS AOU San Martino-IST, Genova, Italy. The 24 lung cell lines were purchased from American Type Culture Collection (ATCC) and cultured following their media and growth conditions. Total RNA from frozen tissues or cell lines were isolated using miRNeasy mini kit (Qiagen Valencia, CA) while RNA was isolated from FFPE sections using FFPE RNAeasy kit (Qiagen). Only high quality RNA from frozen sections and cell lines with RNA integrity number (RIN) >8.0, upon 2100 Bioanalyzer analysis (Agilent Santa Clara, CA) were subjected to RNA sequencing (Supplementary Methods).

### Preparation of RNASeq Libraries and Sequencing

Transcriptome libraries were prepared following a previously described protocol for generating strand-specific RNASeq libraries with slight modifications<sup>50</sup> (Supplementary Methods). Libraries were next size selected in the range of 350 bps after resolving in a 3% Nusieve 3:1 (Lonza, Basel, Switzerland) agarose gel and DNA recovered using QIAEX II gel extraction reagent (Qiagen). Libraries were barcoded during the 14-cycle PCR amplification with Phusion DNA polymerase (New England Biolabs, Ipswich, MA) and purified using AMPure XP beads (Beckman Coulter, Brea, CA). Library quality was estimated with Agilent 2100 Bioanalyzer for size and concentration. The paired end libraries were sequenced with Illumina HiSeq 2000 (2×100bases, read length). Reads that passed the filters on Illumina BaseCall software were used for further analysis. The data have been deposited to Sequence Read Archive (SRA) under the SRA accession number SRP048484.

### Cloning and Expression of *CD74-NRG1* Fusion, Cell Proliferation and Migration Assays

*CD74-NRG1* fusion transcript was amplified from the index lung cancer sample tissue cDNA with forward 5'CACCATGCACAGGAGGAGAAGCAGGAGCTGT3' and reverse primers 5'TTCAGGCAGAGACAGAAAGGGAGTGGGA3' using Hi-fidelity polymerase (Qiagen). The PCR product was gel purified and cloned into pLenti-TOPO cloning vector (Invitrogen, Carlsbad, CA) and Sanger sequencing verified. The control LacZ or C-terminal V5 tagged *CD74-NRG1* constructs were transfected into the normal lung epithelial BEAS-2B cells. The stable cells were generated following selection in BEBM media (Lonza, Basel, Switzerland) containing 3 micrograms of blasticidin (Invitrogen, Carlsbad, CA). For proliferation assays, 50,000 cells were plated in 12-well plates and grown in regular media. Cells were harvested by trypsinization and counted manually at indicated time points. All assays were performed in quadruplicates. For migration assays, stable cells were re-suspended in medium without growth factors, then seeded at 50,000 cells per well into Boyden chambers (8 µm pore size, BD Biosciences) and were incubated for 24 hours in a humidified incubator at 37°C, 5% CO<sub>2</sub> atmosphere. The bottom chamber contained medium with growth factors as chemo-attractant. The top non-migrating cells were removed with a cotton swab moistened with medium and the lower surface of the membrane was stained with Diff-Quick Stain Set (Siemens). The number of cells migrating to the basal side of the membrane was visualized with an Olympus microscope at 20x magnification. Pictures of five random fields from 4 wells were obtained and the number of stained cells was quantified.

## Sequence Alignment, Fusion and Mutation Calling

Sequence alignment was performed using the Tuxedo pipeline: Bowtie2 (Bowtie2/2.0.2) and Tophat2 (TopHat/2.0.6)<sup>51</sup>. Fusion calling was performed with TopHat-fusion<sup>51</sup> (THF) on the UMICH, TCGA and Seoul cohorts. Additional details and parameter values used for sequence alignment and fusion calling are provided in the Supplementary Methods.

## Fusion Annotation and Lung Cancers Fusions Database

A database of fusions in lung cancers was developed, and for each fusion structural and functional annotation was recorded. The structural information corresponds to characteristics such as fusion type (inter-chromosomal, intra-chromosomal, tandem-duplication), number of spanning and encompassing reads, median alignment quality of reads that support 3' and 5' gene, among others (see Supplementary Methods). The functional annotation corresponds to features such as kinase status, oncogene status and tumor suppressor status among others. Moreover, the gene expression of the 5' and 3' partner genes was calculated in fragments per-kilo base per million (FPKM) using Cufflinks<sup>52</sup> and stored in the database. Furthermore, the outlier sum score<sup>53</sup> was independently calculated for the expression of both 5' and 3' partners in order to identify fusion cases for which the 3' gene partner was highly expressed relative to its median expression in the cohort. Overexpression of the 3'-partner as consequence of gene fusions has been observed in well-known fusions such as *TPRSS2-ERG* and others<sup>54</sup>. Finally, we also recorded the mutation status for each patient, allowing us to classify each patient as “driver positive” or “driver negative” according to mutation status of well-known cancer related genes (Supplementary methods).

## Fusions Classifier

All fusion-calling algorithms produce a significant number of false positive fusions when applied on RNASeq data. Many of these spurious fusions are due to diverse and difficult to model bioinformatics, sequencing and biological factors such as: template switching and chimeric events associated with amplicon regions among others<sup>55, 56, 57</sup>. Therefore, we developed a classifier in order to prioritize fusions for follow up based on the structural and functional features collected for each fusion, which were described above and stored in our fusion's database.

THF called 31304 fusions across the combined cohort, making the task of separating false positive fusions from potentially true ones far from trivial. We first reason that functional fusion proteins have open reading frames (ORFs); therefore fusions in which the exon of one gene is fused to the intron of another or two introns are fused together would not produce fusion products with ORFs. This first level filtering reduced to 6465 the number of fusions to classify. Then, we reason that fusions found in normal samples, fusions involving pseudogenes, lincRNAs, or antisense transcripts and fusions for which the median alignment quality of reads supporting any of the gene partners was equal to zero (indicating multi-mapping) are potentially false positives, and there were excluded from downstream analysis. This second level filtering reduced to 4990 the number of fusions called by THF. Because assessing the quality of each one of those fusions manually is impractical, we built a random forest classifier to prioritize what fusions to follow up out those 4990 gene fusions.

For the classification step, we train a random forest classifier with 10000 trees using the structural, functional and expression features described above (Supplementary Methods). True positives examples were selected from the TCGA, Seoul and UMICH cohorts. On one hand, the examples chosen from the TCGA and Seoul cohorts correspond to well-known fusions involving ALK, RET and ROS1 kinases. On the other, the examples chosen from the UMICH cohort correspond to fusions called by at least two independent algorithms, carefully curated manually and validated by PCR (Supplementary Data 4). False positive examples were selected representing different types of spurious fusions: e.g. overlapping genes, fusions involving highly expressed genes such as ribosomal proteins among others. After applying the classifier, we obtained 422 high quality gene fusions. Taken together our approach allowed us to efficiently prioritize the initial set of 31,304 fusions reported by THF, filtering out potential false positives. Finally, open reading frame prediction and protein domain retention analysis were performed in recurrent fusions or biologically interesting fusions found in this final set of 422 fusions.

An additional advantage of using a classifier to determine the potential true fusions, as opposed to hard filters defined a priori, is that we can learn those features or rules from the data itself. In our dataset, the top five features that contributed the most for the random forest classifier were, in decreasing order of importance, fusion type (Inter-chromosomal, Intra-chromosomal, Tandem-duplication), sum of the median alignment quality of both gene partners, number of reads spanning and encompassing reads across the fusion junction and the cohort normalized expression value of the 3' gene (Supplementary Fig. 5).

Two additional sets of true fusions were left out of the training dataset to calculate the recovery rate. First, a set of 11 fusions called in the Seoul cohort<sup>19</sup> and validated by PCR by the same authors, and a second set of 15 fusions called in the UMICH cohort by THF, and validated by PCR. In the first of these datasets, our classifier recovered 10 out 11 true fusions for a 90.1% recovery rate (Supplementary Data 2). In the second set, the classifier recovered 14 out 15 validated fusions for a 93.3% recovery rate (Supplementary Data 3).

## Supplementary Material

Refer to Web version on PubMed Central for supplementary material.

## Acknowledgments

We would like to thank Daniel Miller, Terrence Barrette, Marcin Cieslik for NGS data processing pipeline and analysis, Jyoti Athanikar and Karen Giles for critically reading the manuscript and submission, Xia Jia and John Prensner for experimental assistance.

### Grant Support

This research is supported in part by the National Institutes of Health through grant R01CA154365 to (DGB and AMC) and through the University of Michigan's Cancer Center Support Grant (5 P30 CA46592). O.A.B is supported by the F31 NIH Ruth L. Kirschstein National Research Service Awards for Individual Pre-doctoral Fellowships to Promote Diversity in Health-Related Research (F31-CA-165866) and by T32 Proteome Informatics of Cancer Training Program at the University of Michigan (T32-CA-140044). P.H. is supported by Dermatology Foundation, Dermatopathology Research Career Development Award. E. N. was supported by a Spanish Society of Medical Oncology Fellowship. J.P. is supported by the China Scholarship Council Award (201206380049). B.V. is supported by T32 Proteome Informatics of Cancer Training Program at the University of Michigan (T32-CA-140044), and by the National Science Foundation under Grant No. 0903629.

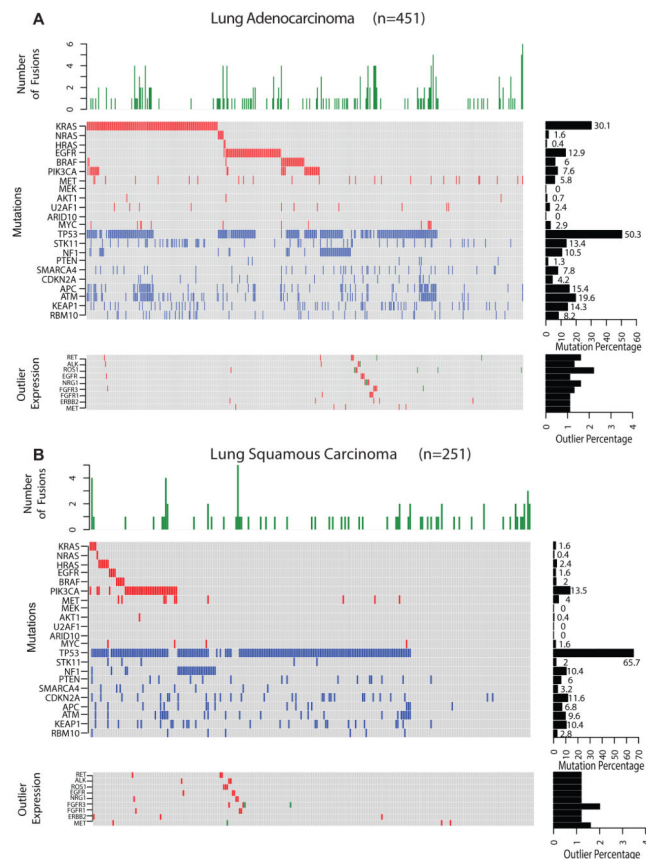
## References

1. Ferlay J, Shin HR, Bray F, Forman D, Mathers C, Parkin DM. Estimates of worldwide burden of cancer in 2008: GLOBOCAN 2008. *International journal of cancer Journal international du cancer*. 2010; 127:2893–2917. [PubMed: 21351269]
2. Siegel R, Naishadham D, Jemal A. Cancer statistics, 2013. *CA: a cancer journal for clinicians*. 2013; 63:11–30. [PubMed: 23335087]
3. Nakamura H, Saji H. A worldwide trend of increasing primary adenocarcinoma of the lung. *Surg Today*. 2013
4. Ding L, et al. Somatic mutations affect key pathways in lung adenocarcinoma. *Nature*. 2008; 455:1069–1075. [PubMed: 18948947]
5. Weir BA, et al. Characterizing the cancer genome in lung adenocarcinoma. *Nature*. 2007; 450:893–898. [PubMed: 17982442]
6. Pao W, Girard N. New driver mutations in non-small-cell lung cancer. *Lancet Oncol*. 2011; 12:175–180. [PubMed: 21277552]
7. Cancer Genome Atlas Research N. Comprehensive genomic characterization of squamous cell lung cancers. *Nature*. 2012; 489:519–525. [PubMed: 22960745]
8. Cancer Genome Atlas Research N. Comprehensive molecular profiling of lung adenocarcinoma. *Nature*. 2014; 511:543–550. [PubMed: 25079552]
9. Paez JG, et al. EGFR mutations in lung cancer: correlation with clinical response to gefitinib therapy. *Science*. 2004; 304:1497–1500. [PubMed: 15118125]
10. Soda M, et al. Identification of the transforming EML4-ALK fusion gene in non-small-cell lung cancer. *Nature*. 2007; 448:561–566. [PubMed: 17625570]
11. Inamura K, et al. EML4-ALK lung cancers are characterized by rare other mutations, a TTF-1 cell lineage, an acinar histology, and young onset. *Modern pathology: an official journal of the United States and Canadian Academy of Pathology, Inc*. 2009; 22:508–515.
12. Takeuchi K, et al. KIF5B-ALK, a novel fusion oncokinin identified by an immunohistochemistry-based diagnostic system for ALK-positive lung cancer. *Clinical cancer research: an official journal of the American Association for Cancer Research*. 2009; 15:3143–3149. [PubMed: 19383809]
13. Rikova K, et al. Global survey of phosphotyrosine signaling identifies oncogenic kinases in lung cancer. *Cell*. 2007; 131:1190–1203. [PubMed: 18083107]
14. Ju YS, et al. A transforming KIF5B and RET gene fusion in lung adenocarcinoma revealed from whole-genome and transcriptome sequencing. *Genome Research*. 2012; 22:436–445. [PubMed: 22194472]
15. Takeuchi K, et al. RET, ROS1 and ALK fusions in lung cancer. *Nature Medicine*. 2012; 18:378–381.
16. Drilon A, et al. Response to Cabozantinib in patients with RET fusion-positive lung adenocarcinomas. *Cancer discovery*. 2013; 3:630–635. [PubMed: 23533264]
17. Wang XS, et al. An integrative approach to reveal driver gene fusions from paired-end sequencing data in cancer. *Nature Biotechnology*. 2009; 27:1005–1011.
18. Wu YM, et al. Identification of targetable FGFR gene fusions in diverse cancers. *Cancer discovery*. 2013; 3:636–647. [PubMed: 23558953]
19. Seo JS, et al. The transcriptional landscape and mutational profile of lung adenocarcinoma. *Genome Research*. 2012; 22:2109–2119. [PubMed: 22975805]
20. Fernandez-Cuesta L, et al. CD74-NRG1 Fusions in Lung Adenocarcinoma. *Cancer Discov*. 2014; 4:415–422. [PubMed: 24469108]
21. Gow CH, Wu SG, Chang YL, Shih JY. Multidriver mutation analysis in pulmonary mucinous adenocarcinoma in Taiwan: identification of a rare CD74-NRG1 translocation case. *Med Oncol*. 2014; 31:34. [PubMed: 24913807]
22. Nakaoku T, et al. Druggable oncogene fusions in invasive mucinous lung adenocarcinoma. *Clin Cancer Res*. 2014; 20:3087–3093. [PubMed: 24727320]
23. Kong-Beltran M, et al. Somatic mutations lead to an oncogenic deletion of met in lung cancer. *Cancer Res*. 2006; 66:283–289. [PubMed: 16397241]

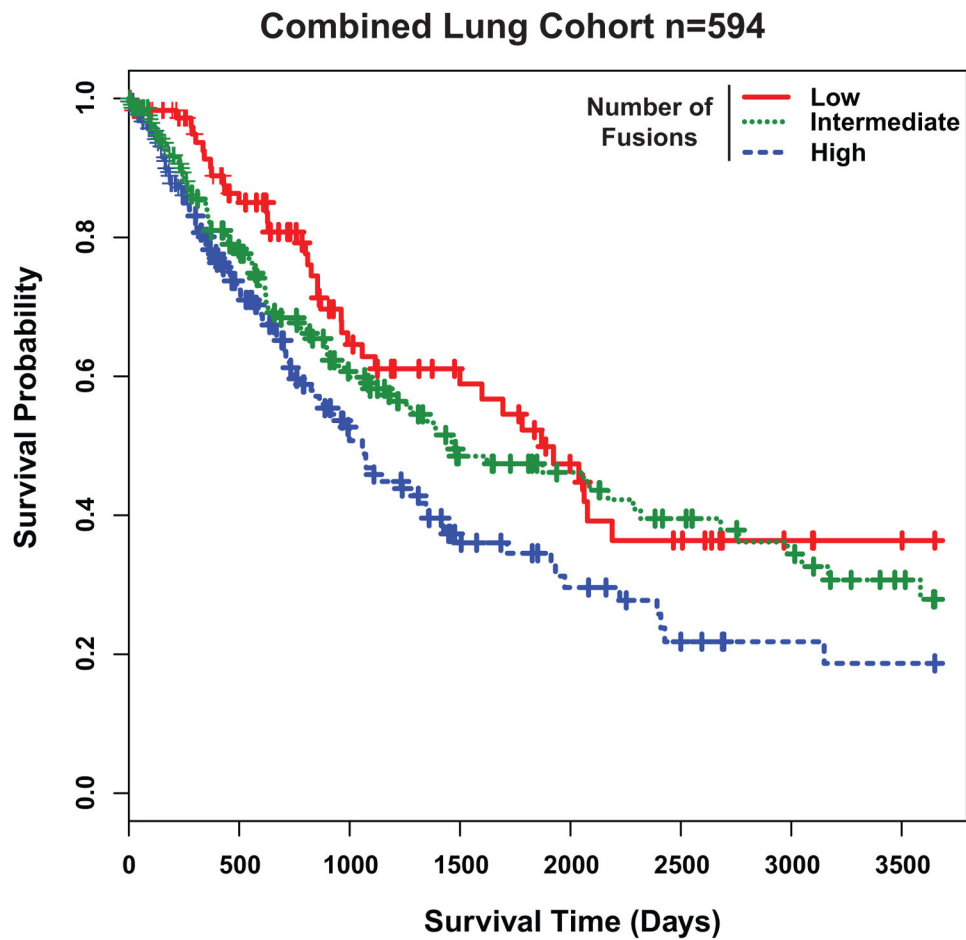
24. A genomics-based classification of human lung tumors. *Sci Transl Med.* 2013; 5:209ra153.
25. Ho AS, et al. The mutational landscape of adenoid cystic carcinoma. *Nat Genet.* 2013; 45:791–798. [PubMed: 23685749]
26. Wetterskog D, et al. Mutation profiling of adenoid cystic carcinomas from multiple anatomical sites identifies mutations in the RAS pathway, but no KIT mutations. *Histopathology.* 2013; 62:543–550. [PubMed: 23398044]
27. Wetterskog D, et al. Adenoid cystic carcinomas constitute a genomically distinct subgroup of triple-negative and basal-like breast cancers. *J Pathol.* 2012; 226:84–96. [PubMed: 22015727]
28. Zhao B, Li L, Lei Q, Guan KL. The Hippo-YAP pathway in organ size control and tumorigenesis: an updated version. *Genes Dev.* 2010; 24:862–874. [PubMed: 20439427]
29. Harvey KF, Zhang X, Thomas DM. The Hippo pathway and human cancer. *Nat Rev Cancer.* 2013; 13:246–257. [PubMed: 23467301]
30. Bollag G, et al. Loss of NF1 results in activation of the Ras signaling pathway and leads to aberrant growth in haematopoietic cells. *Nat Genet.* 1996; 12:144–148. [PubMed: 8563751]
31. Sandsmark DK, Zhang H, Hegedus B, Pelletier CL, Weber JD, Gutmann DH. Nucleophosmin mediates mammalian target of rapamycin-dependent actin cytoskeleton dynamics and proliferation in neurofibromin-deficient astrocytes. *Cancer Res.* 2007; 67:4790–4799. [PubMed: 17510408]
32. Onozato R, Kosaka T, Kuwano H, Sekido Y, Yatabe Y, Mitsudomi T. Activation of MET by gene amplification or by splice mutations deleting the juxtamembrane domain in primary resected lung cancers. *J Thorac Oncol.* 2009; 4:5–11. [PubMed: 19096300]
33. Falls DL. Neuregulins: functions, forms, and signaling strategies. *Exp Cell Res.* 2003; 284:14–30. [PubMed: 12648463]
34. Holmes WE, et al. Identification of heregulin, a specific activator of p185erbB2. *Science.* 1992; 256:1205–1210. [PubMed: 1350381]
35. Wen D, et al. Structural and functional aspects of the multiplicity of Neu differentiation factors. *Mol Cell Biol.* 1994; 14:1909–1919. [PubMed: 7509448]
36. Cao Z, Wu X, Yen L, Sweeney C, Carraway KL 3rd. Neuregulin-induced ErbB3 downregulation is mediated by a protein stability cascade involving the E3 ubiquitin ligase Nrdp1. *Mol Cell Biol.* 2007; 27:2180–2188. [PubMed: 17210635]
37. Sholl LM, et al. EGFR mutation is a better predictor of response to tyrosine kinase inhibitors in non-small cell lung carcinoma than FISH, CISH, and immunohistochemistry. *American journal of clinical pathology.* 2010; 133:922–934. [PubMed: 20472851]
38. Lipson D, et al. Identification of new ALK and RET gene fusions from colorectal and lung cancer biopsies. *Nature Medicine.* 2012; 18:382–384.
39. Koivunen JP, et al. EML4-ALK fusion gene and efficacy of an ALK kinase inhibitor in lung cancer. *Clinical cancer research: an official journal of the American Association for Cancer Research.* 2008; 14:4275–4283. [PubMed: 18594010]
40. Shaw AT, et al. Effect of crizotinib on overall survival in patients with advanced non-small-cell lung cancer harbouring ALK gene rearrangement: a retrospective analysis. *The lancet oncology.* 2011; 12:1004–1012. [PubMed: 21933749]
41. Tanas MR, et al. Identification of a disease-defining gene fusion in epithelioid hemangioendothelioma. *Sci Transl Med.* 2011; 3:98ra82.
42. Errani C, et al. A novel WWTR1-CAMTA1 gene fusion is a consistent abnormality in epithelioid hemangioendothelioma of different anatomic sites. *Genes Chromosomes Cancer.* 2011; 50:644–653. [PubMed: 21584898]
43. Mohseni M, et al. A genetic screen identifies an LKB1-MARK signalling axis controlling the Hippo-YAP pathway. *Nat Cell Biol.* 2014; 16:108–117. [PubMed: 24362629]
44. Montero JC, Rodriguez-Barrueco R, Ocana A, Diaz-Rodriguez E, Esparis-Ogando A, Pandiella A. Neuregulins and cancer. *Clin Cancer Res.* 2008; 14:3237–3241. [PubMed: 18519747]
45. Sugano M, et al. HNF4alpha as a marker for invasive mucinous adenocarcinoma of the lung. *Am J Surg Pathol.* 2013; 37:211–218. [PubMed: 23108025]



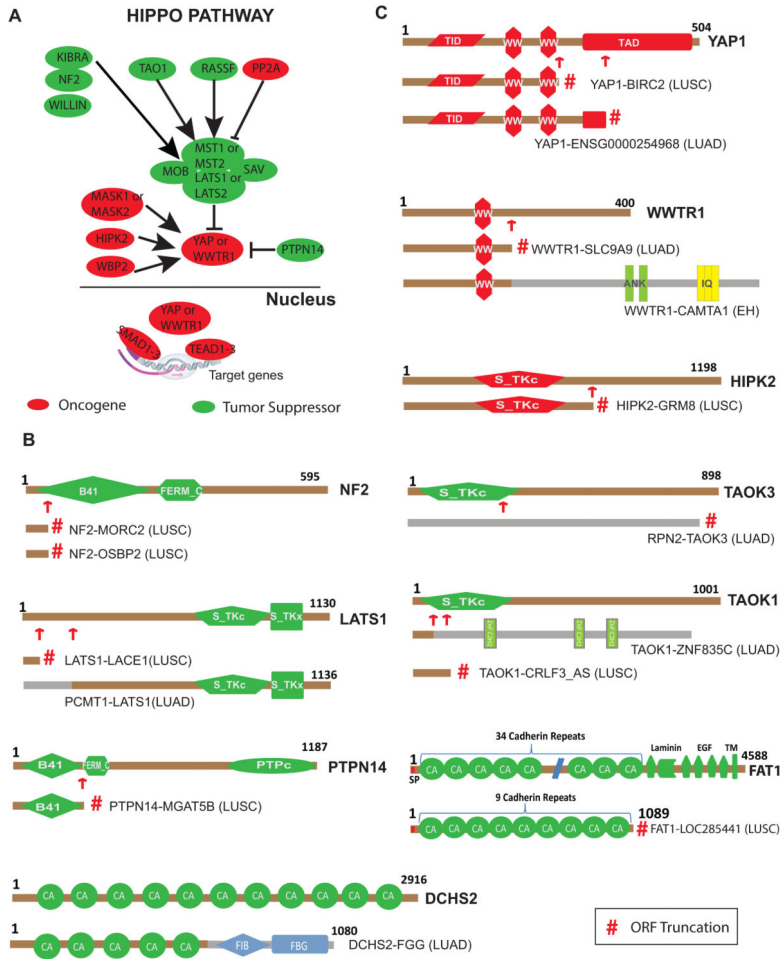
46. Adelaide J, et al. A recurrent chromosome translocation breakpoint in breast and pancreatic cancer cell lines targets the neuregulin/NRG1 gene. *Genes, chromosomes & cancer*. 2003; 37:333–345. [PubMed: 12800145]
47. Prentice LM, et al. NRG1 gene rearrangements in clinical breast cancer: identification of an adjacent novel amplicon associated with poor prognosis. *Oncogene*. 2005; 24:7281–7289. [PubMed: 16027731]
48. Gollamudi M, Nethery D, Liu J, Kern JA. Autocrine activation of ErbB2/ErbB3 receptor complex by NRG-1 in non-small cell lung cancer cell lines. *Lung cancer*. 2004; 43:135–143. [PubMed: 14739033]
49. Hegde GV, et al. Blocking NRG1 and other ligand-mediated Her4 signaling enhances the magnitude and duration of the chemotherapeutic response of non-small cell lung cancer. *Sci Transl Med*. 2013; 5:171ra118.
50. Levin JZ, et al. Comprehensive comparative analysis of strand-specific RNA sequencing methods. *Nature Methods*. 2010; 7:709–715. [PubMed: 20711195]
51. Kim D, Pertea G, Trapnell C, Pimentel H, Kelley R, Salzberg SL. TopHat2: accurate alignment of transcriptomes in the presence of insertions, deletions and gene fusions. *Genome Biology*. 2013; 14:R36. [PubMed: 23618408]
52. Trapnell C, et al. Differential gene and transcript expression analysis of RNA-seq experiments with TopHat and Cufflinks. *Nature Protocols*. 2012; 7:562–578. [PubMed: 22383036]
53. Tibshirani R, Hastie T. Outlier sums for differential gene expression analysis. *Biostatistics*. 2007; 8:2–8. [PubMed: 16702229]
54. Tomlins SA, et al. Recurrent fusion of TMPRSS2 and ETS transcription factor genes in prostate cancer. *Science*. 2005; 310:644–648. [PubMed: 16254181]
55. Carrara M, et al. State of art fusion-finder algorithms are suitable to detect transcription-induced chimeras in normal tissues? *BMC bioinformatics*. 2013; 14 (Suppl 7):S2. [PubMed: 23815381]
56. Carrara M, et al. State-of-the-art fusion-finder algorithms sensitivity and specificity. *BioMed research international*. 2013; 2013:340620. [PubMed: 23555082]
57. Ozsolak F, Milos PM. RNA sequencing: advances, challenges and opportunities. *Nature reviews Genetics*. 2011; 12:87–98.



**Fig. 1.** The gene fusion and mutational landscape of lung cancers. **A**, Lung adenocarcinoma (LUAD, n=451). **B**, Lung squamous carcinoma (LUSC, n=251). Top panels represent histograms depicting the number of high quality gene fusions identified in each sample. Central panels denote the presence or absence of activating mutations in known oncogenes (red), deleterious mutations in tumor suppressors (blue) no aberration (gray). Samples are represented in columns and genes in rows. Right middle panel are bar plot summarizing the number of samples harboring activating or deleterious mutations for each gene. Bottom panels indicate samples harboring both known and novel gene fusions (in green) involving either receptor kinase genes or *NRG1*. Samples in red indicate outlier expression pattern observed in the respective genes. Cohorts of additional non-small cell lung cancers including lung adenoid cystic carcinoma (ACC) (n=11) and large cell carcinomas (n=9) were also analyzed included in Supplementary Table 2.



**Fig. 2.** Gene fusion numbers correlates with lung cancer prognosis. **A**, Kaplan-Meier analysis for the combined cohort of lung cancer samples (n=594) with low (0–7) (n=124), intermediate (8–16) (n=237), or high (17) (n=233) number of fusions (Likelihood Ratio Test  $P=0.008$ ). Samples with high number of fusions have worse prognosis (Cox survival analysis  $P=0.005$ ). Individual Kaplan-Meier analyses with LUAD and LUSC samples are found in Supplementary Fig. 4A and 4B respectively.



**Fig. 3.** Gene fusions among the Hippo pathway genes in lung cancer. **A**, Schematic representation of core and associate members of the Hippo pathway adapted from *Harvey et al*<sup>29</sup>. Potential tumor suppressors are represented in green, while potential oncogenes are indicated in red. Phosphorylation of YAP or TAZ by LATS retains them in the cytoplasm and hinders their transcriptional regulation. **B**, Fusions in putative oncogenes of the Hippo pathway. **C**, Fusions in putative tumor suppressors of the Hippo pathway. For all fusion schematics represented, the wild-type Hippo pathway protein domain structure is presented first, numbers indicate total amino acids and domain names are abbreviated. Red arrows show the fusion junctions and red # symbol indicate protein truncation due to out-of-frame ORFs from fusion transcript analysis. The schematic of the previously reported TAZ-CAMTA1 fusion in epithelioid hemangioendothelioma (EH)<sup>42</sup> is also displayed. Protein abbreviations: MST1/2-STE20-like protein kinase; LATS1/2-Large Tumor Suppressor Homolog Kinase; YAP1-Yes-associated Protein 1; WWTR1-ww-Domain Containing Transcription Regulator 1; TEAD-TEA-Domain Family; HIPK2 Homeodomain Interacting Protein Kinase 2; TAOK1/3-TAO Kinase; FAT1-FAT Atypical Cadherin 1; DCHS2-Dachsous Cadherin-related 2; PTPN14-Protein Tyrosine Phosphatase, Non-receptor Type 14. Domain abbreviations: B4-Band 4.1 homologues; FERM\_C-FERM C-terminal PH-like Domain;

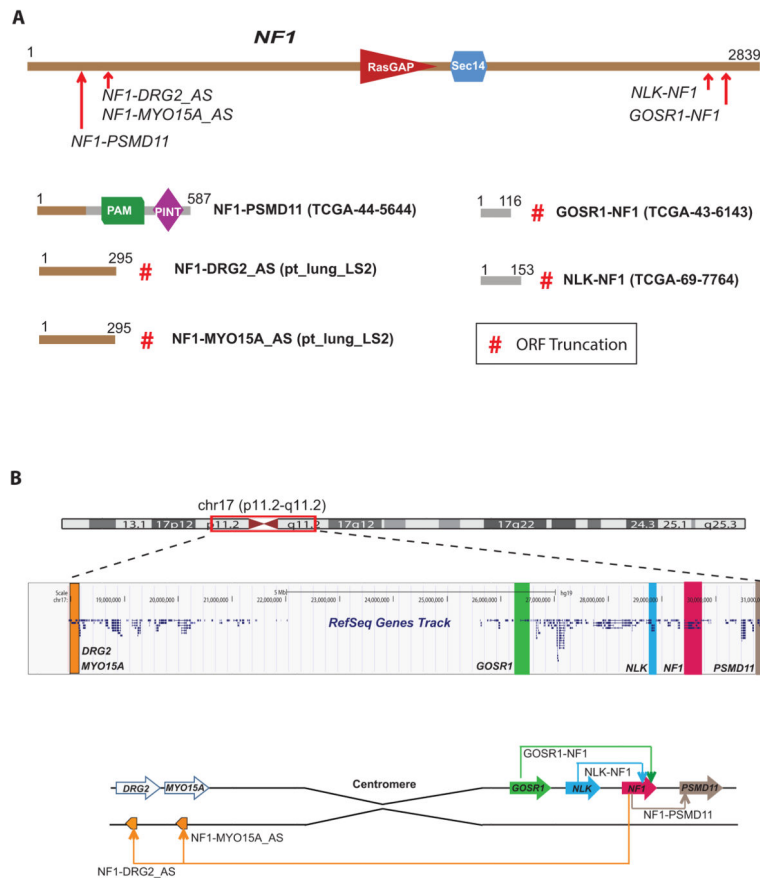
S\_TKc-Serine/Threonine Protein Kinases, Catalytic Domain; PTPc-Protein Tyrosine Phosphatase, Catalytic Domain; CA-Cadherin Repeats; FIB-Fibrinogen; FBG-Fibrinogen-related Domains; WW-Domain with 2 conserved Trp (W) residues; TID-TEAD Interacting Domain; TAD-Transactivation Domain; ANK-Ankyrin Repeats; IQ-Short Calmodulin-binding Motif; EGF-Epidermal Growth Factor-like Domain; ZnFC2H2-Zinc Finger; TM-Transmembrane Domain.

Author Manuscript

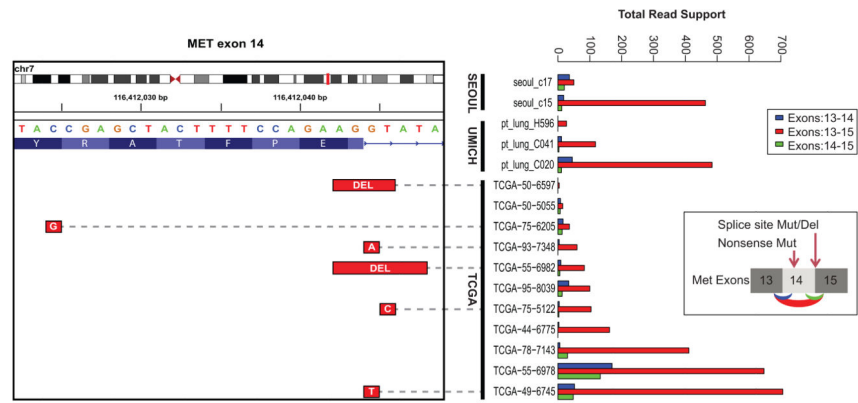
Author Manuscript

Author Manuscript

Author Manuscript

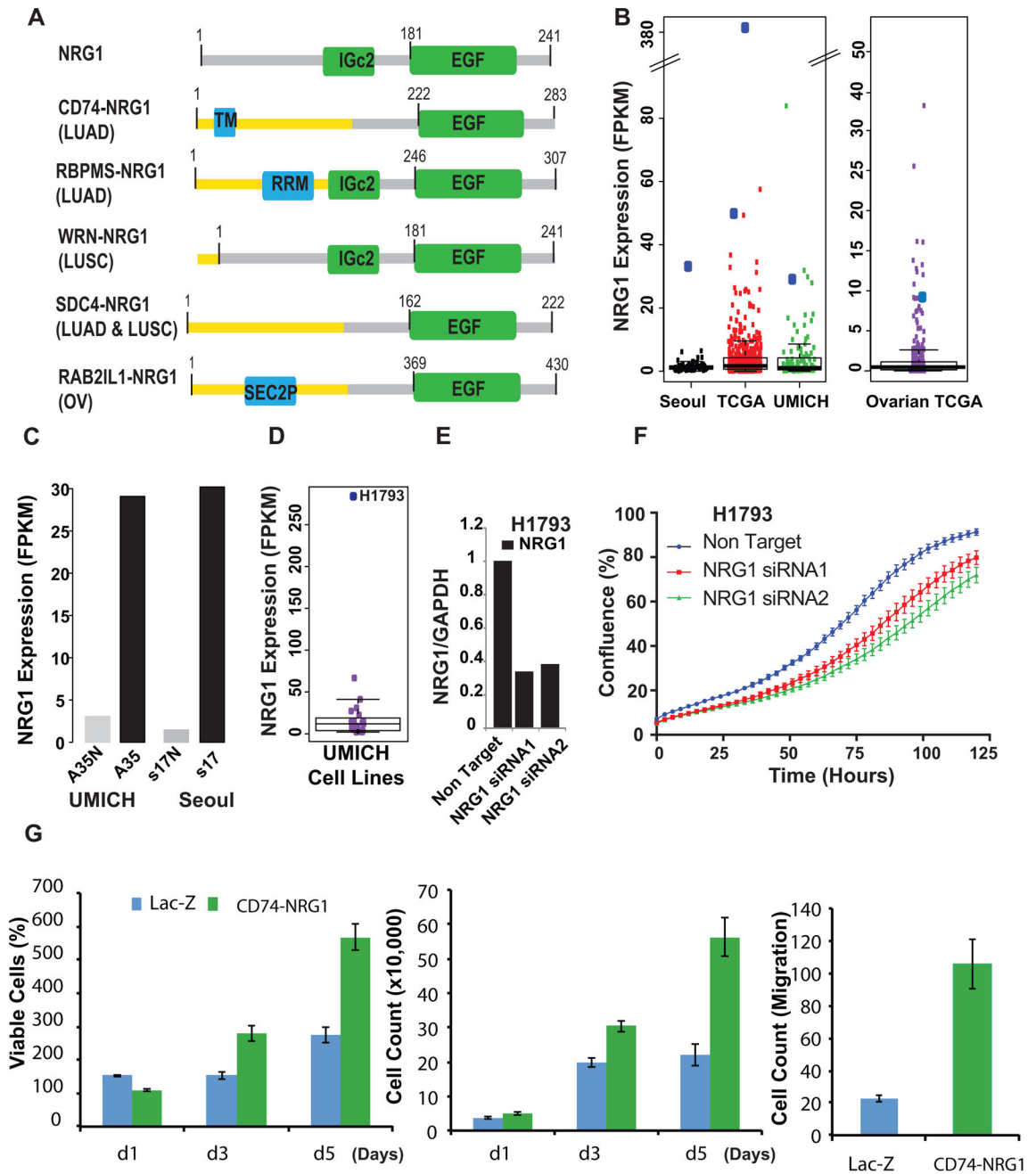


**Fig. 4.** Inactivating gene fusions of *NF1* in lung cancer. **A**, *NF1* protein schematic and the observed fusion breaks (red arrows) in the index cases are displayed on top. Recurrent *NF1* fusions with partners (*GOSR1*, *PSMD11*, *NLK*, *DRG2* and *MYO15A antisense*) resulted in loss of the *NF1* gene as illustrated by the corresponding fusion protein structure below. Index samples are indicated in parenthesis and the numbers over the protein schematic indicate total amino acids. Red # symbol indicate protein truncation due to out-of-frame ORFs from fusion transcript analysis. **B**, UCSC browser view of genomic location of *NF1* gene and its fusion partners (Top). Schematic representation of various *NF1* rearrangements on chromosome 17 identified in lung cancer (Bottom).



**Fig. 5.**

Recurrent activating *MET* exon skipping events. Right Panel: An activating *MET* exon-14 skipping event was observed in a total of 15 samples across all three cohorts. The total reads support each splice variant exon13–14 (blue), exon13–15 (red) and exon14–15 (green) are represented in the bar plot on the right. In 5 out of 11 TCGA samples where DNA mutation data was available, skipping of *MET* exon-14 was accompanied by a mutation affecting the splice donor site adjacent to position D1010 (illustrated inset on the right). Additionally one sample harbored a non-sense mutation g.chr7:116412024C>Gp.Y1003\*, which accompanied exon-14 skipping. Left Panel: IGV browser view of splice site deletions/mutations in the corresponding samples.



**Fig. 6.** Recurrent *NRG1* rearrangements in lung cancer. **A**, Recurrent fusions involving *NRG1* as a 3' partner were detected in lung adenocarcinoma and lung squamous carcinoma in the three cohorts included in this study. Schematic representation of functional domains present in the *NRG1* fusion proteins namely *CD74-NRG1*; *RBPMS-NRG1* (LUAD); *WRN-NRG1* (LUSC); *SDC4-NRG1* (LUSC) and *RAB2IL1-NRG1* (ovarian cancer from TCGA) compared to the wild-type *NRG1* (Top). The receptor binding EGF domain is preserved in all fusions. TM, transmembrane domain; RRM- domain; IGc2- domain; SEC2P-domain. **B**, Analysis of



RNASeq expression values revealed outlier *NRG1* mRNA expression in all index cases (large blue dots) within each cohort. **C**, High *NRG1* mRNA expression driven by the fusion event in the index tumor tissue compared to matched normal, in both an LUAD patient in the University of Michigan and Seoul cohorts. **D**, Boxplot showing outlier expression of *NRG1* in H1793 in the University of Michigan lung cell line cohort. **E**, Two independent siRNAs mediated knockdown of *NRG1* in H1793 cells as assessed by Q-PCR. **F**, Knock-down of *NRG1* decreased cell proliferation as monitored by IncuCyte confluence analysis. **G**, Overexpression of *NRG1* induces cell proliferation and migration. Cell proliferation by WST-1 assay (left panel) and cell counting (middle panel) on BEAS-2B cells stably transfected with Lac-Z or *CD74-NRG1* fusion. Both assays demonstrated that cells expressing the *CD74-NRG1* fusion had significantly higher proliferation rate at day 3 and 5 (Student's t-test  $P < 0.001$  for both time-points) as compared to Lac-Z. The right panel represents a cell migration assay after 24 hours. BEAS-2B cells expressing *CD74-NRG1* fusion showed a higher migration rate as compared to Lac-Z (Student's t-test  $P = 0.0014$ ).