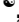
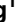




RESEARCH ARTICLE

Functionally distinct BMP1 isoforms show an opposite pattern of abundance in plasma from non-small cell lung cancer subjects and controls

Margaret K. R. Donovan¹ , Yingxiang Huang¹ , John E. Blume¹, Jian Wang¹, Daniel Hornburg¹, Shadi Ferdosi¹ , Iman Mohtashemi¹, Sangtae Kim¹, Marwin Ko¹, Ryan W. Benz¹, Theodore L. Platt¹, Serafim Batzoglou¹, Luis A. Diaz², Omid C. Farokhzad¹, Asim Siddiqui¹ *

1 Seer, Inc., Redwood City, CA, United States of America, **2** The Ludwig Center and The Howard Hughes Medical Institute at Johns Hopkins Kimmel Cancer Center, Baltimore, MD, United States of America

 These authors contributed equally to this work.

* asiddiqui@seer.bio



OPEN ACCESS

Citation: Donovan MKR, Huang Y, Blume JE, Wang J, Hornburg D, Ferdosi S, et al. (2023) Functionally distinct BMP1 isoforms show an opposite pattern of abundance in plasma from non-small cell lung cancer subjects and controls. *PLoS ONE* 18(3): e0282821. <https://doi.org/10.1371/journal.pone.0282821>

Editor: Qi Wu, Aarhus University, DENMARK

Received: November 23, 2022

Accepted: February 23, 2023

Published: March 29, 2023

Copyright: © 2023 Donovan et al. This is an open access article distributed under the terms of the [Creative Commons Attribution License](https://creativecommons.org/licenses/by/4.0/), which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.

Data Availability Statement: The data used in this project is available in ProteomeXchange Consortium (<http://proteomecentral.proteomexchange.org>) via the PRIDE partner repository with the dataset PXD017052.

Funding: This work was funded by Seer, Inc. The funders had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript.

Competing interests: OCF. has financial interest in Selecta Biosciences, Tarveda Therapeutics, and

Abstract

Advancements in deep plasma proteomics are enabling high-resolution measurement of plasma proteoforms, which may reveal a rich source of novel biomarkers previously concealed by aggregated protein methods. Here, we analyze 188 plasma proteomes from non-small cell lung cancer subjects (NSCLC) and controls to identify NSCLC-associated protein isoforms by examining differentially abundant peptides as a proxy for isoform-specific exon usage. We find four proteins comprised of peptides with opposite patterns of abundance between cancer and control subjects. One of these proteins, BMP1, has known isoforms that can explain this differential pattern, for which the abundance of the NSCLC-associated isoform increases with stage of NSCLC progression. The presence of cancer and control-associated isoforms suggests differential regulation of BMP1 isoforms. The identified BMP1 isoforms have known functional differences, which may reveal insights into mechanisms impacting NSCLC disease progression.

Introduction

Multiple isoforms of a single protein, or proteoforms, can arise due to alternative splicing (i.e., protein isoforms), allelic variation, and post translational modifications [1]. Proteoforms play key and distinct roles in biological mechanisms, including impacting complex traits [2] and disease [3]. For example, protein isoforms may differ in domain composition, where consequently each isoform may have substantially different functions and influence disease predisposition or progression. Advances in characterizing the proteomic landscape of lung cancers such as non-small cell lung cancer (NSCLC) and squamous cell lung cancer have enabled identification of important protein biomarkers [4–6], however, few proteoforms relevant to lung cancer have been identified [7], as these studies are limited to only single or few protein [8–10]

Seer. MKRD, YH, JEB, JW, DH, SF, IM, SK, MK, RWB, TLP, SB, OCF, and AS have financial interest in Seer. LAD is a member of Seer's Scientific Advisor Board and is financially compensated for that role. Only Seer, and no other companies mentioned here, was involved in the study design, data collection and analysis, and manuscript writing/editing. This does not alter our adherence to PLOS ONE policies on sharing data and materials.

or proteoforms arising from different genes [11]. Unbiased readout technologies, such as high-resolution quantitative mass spectrometry (MS), can be employed to infer and quantify peptides and proteins with high confidence (e.g., < 1% false discovery rate (FDR)). However, large-scale LC-MS/MS-based proteomics studies have historically been impractical due to cumbersome and lengthy workflows required to achieve unbiased, deep, and rapid sampling of clinically relevant biospecimens with large dynamic ranges of protein abundances, such as blood plasma [12–14].

Here, we analyze data from a previous study [15] of independent acquisition (DIA)-based MS data generated from 188 subjects (80 healthy control subjects and 108 subjects identified as having NSCLC) using the Proteograph™ workflow which uses nanoparticles (NPs) to enable high-resolution, unbiased, and deep assessment of the plasma proteome. We used a discordant peptide intensity search (Fig 1A) to infer four proteins with differentially abundant protein isoforms, including BMP1, for which we show has differential abundance of two isoforms (long and short), both of which have higher magnitude of differential abundance at later stages of NSCLC. BMP1 plays a role in collagen processing and the short isoform lacks the domains enabling its secretion, potentially impacting collagen's protective role in cancer consistent with the higher abundance of the short isoform in cancer subjects observed in this paper. Hence, BMP1 isoforms may constitute a novel biomarker previously concealed when assessing the aggregated BMP1 protein abundance.

Results

Peptide-level analyses provides unique biological insight versus protein-level

Starting from the previously derived analysis [15], we searched for proteins and peptides that are differentially abundant (DA). First, to reduce potential noise introduced by rare peptides, proteins were filtered to those present in at least 50% of subjects from either 80 healthy or 61 early NSCLC (stages 1, 2 and 3) subjects, retaining 10,280 peptides and 1,565 proteins across 141 subjects (Fig 1B). Next, as each protein may have been detected by more than one NP (each NP can be thought of as generating a separate MS fraction), we use MaxLFQ [16] to quantify a single abundance (hereto referred to as *collapsed abundances*) between healthy and early NSCLC subjects. We evaluated differential protein abundance observing 243 significantly regulated proteins (adjusted $p < 0.05$; Wilcoxon Test) (Fig 1C and 1D). To investigate NPs capacity to capture biological signal beyond abundance levels (e.g., proteoform information, or NP specific protein complexes), we treated each NP:protein feature pair as a separate observation comparing healthy and early NSCLC subjects. We identified 877 NP:protein feature pairs (Fig 1E), corresponding to a 3.6-fold increase from examining differences at the aggregated level alone. This highlights the capacity of NPs coupled with LC-MS/MS to interrogate the proteome at a finer biological resolution (i.e. protein variants and complexes) than that captured by conventional DA analysis at the aggregated protein level. In addition, we performed DA analysis using peptide abundances across all NPs (i.e., not collapsed abundances) between healthy and early NSCLC subjects and identified 5,181 DA peptides (Fig 1C and 1E), corresponding to a 6.5-fold increase from examining differences at the protein-level. Further, we identified known hallmark cancer and inflammatory biomarkers which were differentially regulated in the peptide data (S1 File, Fig 1D and 1E). Overall, this increased number of observed significant differences between proteins, protein across NPs, and peptides across NPs, verified by the presence of known cancer biomarkers, indicates substantial opportunity to increase biological insight and the potential to identify proteoforms using peptide-level, high resolution proteomics.

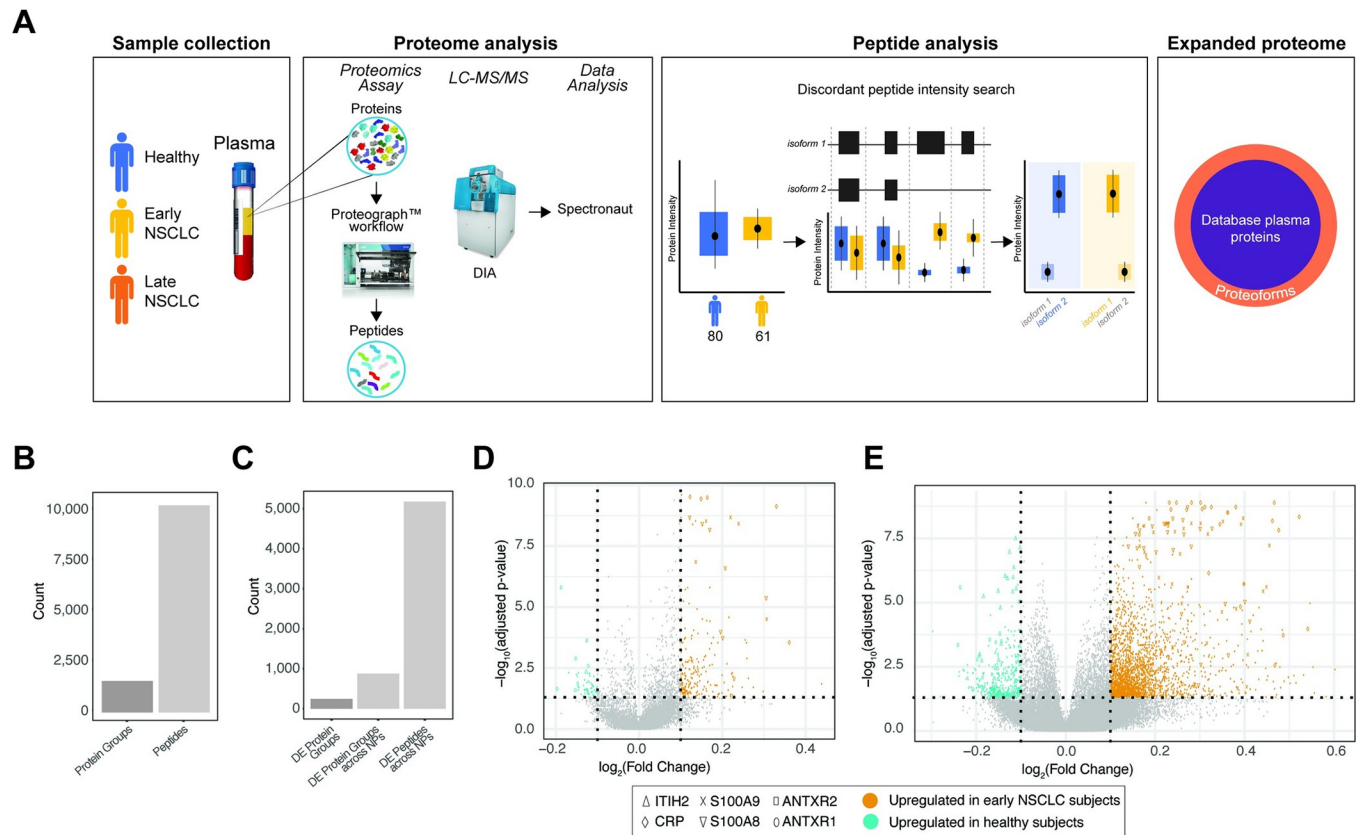


Fig 1. Proteome analysis of healthy and NSCLC subjects using a 5 NP plasma workflow. A. Overview of this proof-of-concept proteoform identification study. Plasma samples were collected from healthy (blue), early non-small cell lung cancer (NSCLC; yellow), late NSCLC (orange), and co-morbid (green) subjects (*Sample Collection*). The plasma proteomes were analyzed for each of these subjects, which included protein extraction, protein discovery using the NP-based Proteograph platform, then DIA protein/peptide identification and quantification using LC-MS/MS and search algorithms (*Proteome Analysis*). Proteoforms were then identified using a discordant peptide intensity search, which included examining peptide mappings to known protein coding isoforms and using differential abundance to discover protein isoforms. Together, these identified proteoforms represent an expanded plasma proteome database not captured in standard MS-based or targeted proteomic studies (*Expanded proteome*). B. Barplots showing the number of peptides and proteins groups retained after filtering to those present in at least 50% of subjects from either healthy or early NSCLC. C. Barplots showing the number of differentially abundant (DA): 1) protein groups, with collapsed abundances using MaxLFQ; 2) protein groups across NPs (i.e., DA independently across NPs); and 3) peptides across NPs. D. Volcano plot showing the significance (adjusted p-value; y-axis) and fold change (x-axis) from calculating the differential abundance of protein groups across NPs between healthy and early NSCLC subjects. Protein groups with a $\log_2(\text{Fold Change})$ greater or less than 1.0 and adjusted p-value < 0.05 are highlighted, where protein groups with increased abundance in early NSCLC subjects are shown in orange and protein groups with increased abundance in healthy subjects are shown in teal. Proteins with known roles in cancer and immune response (ITIH2, CRP, S100A9, S100A8, ANTXR2, and ANTXR1) are highlighted with various shapes. E. Volcano plot showing the significance (adjusted p-value; y-axis) and fold change (x-axis) from calculating the differential abundance of peptides across NPs between healthy and early NSCLC subjects. Peptides with a $\log_2(\text{Fold Change})$ greater or less than 1.0 and adjusted p-value < 0.05 are highlighted, where peptides with increased abundance in early NSCLC subjects are shown in orange and peptides with increased abundance in healthy subjects are shown in teal. Peptides mapping to proteins with known roles in cancer and immune response (ITIH2, CRP, S100A9, S100A8, ANTXR2, and ANTXR1) are highlighted with various shapes.

<https://doi.org/10.1371/journal.pone.0282821.g001>

Identification of four NSCLC-associated proteoforms using peptide-level discordant peptide search

Next, we explored whether we could use DA peptides in contrast to the average protein-level information to help resolve proteoforms. Specifically, we extracted DA peptides and retained proteins with at least one peptide over-expressed in healthy subjects and at least one peptide over-expressed in early NSCLC subjects (Fig 2A). Then, by mapping the DA peptides to genomic space, we inferred potential exon usage and proteoforms. We performed this discordant peptide intensity analysis and identified four proteins for which we potentially captured

multiple protein isoforms with significant differential behavior in early NSCLC when compared to healthy controls: BMP1, C4A, C1R, and LDHB (Fig 2B). We examined the Open Target Score [17] (Release 21.09), which is an association score of known and potential drug targets with diseases using integrated genome-wide data from a broad range of data sources, to assess the association of the four proteins with lung carcinoma targets. We found modest to low scores (Fig 2B), suggesting a mix of novel and known lung cancer-associated proteins. These proteins have all been previously identified in plasma and range from highly abundant (C4A, C1R, LDHB) to moderately abundant (BMP1) [18] (Fig 2C). BMP1, the least abundant of the four proteins, is not identified in depleted plasma published in this study, indicating this approach identified protein isoforms inaccessible with conventional depleted plasma proteomics workflows. These results indicate that, using a MS-based peptide discordant intensity search, we can infer proteoforms with possible relevance to NSCLC.

To interrogate the extent to which isoforms information adds to disease insight, we examined differences in abundances between healthy and early NSCLC subjects for BMP1 (Fig 2D–2G), C4A (S1 Fig), C1R (S2 Fig), and LDHB (S3 Fig) at the collapsed protein-level, NP:protein-level, and peptide-level. Examining BMP1, at the collapsed protein (Fig 2D) and NP:protein (Fig 2E) level, we do not observe a difference in BMP1 abundance, as a result of an averaging of peptide abundances occurring at the protein-level. However, at the peptide-level (Fig 2F), there are three significantly differential peptides: 1) peptide 1, which is significantly upregulated in early NSCLC subjects (adjusted $p = 5.29 \times 10^{-4}$; Wilcoxon Test); 2) peptide 3, which is significantly upregulated in healthy subjects (adjusted $p = 1.21 \times 10^{-2}$; Wilcoxon Test); and 3) peptide 7, which is significantly upregulated in healthy subjects (adjusted $p = 4.99 \times 10^{-2}$; Wilcoxon Test). We also observe a trend in direction of abundance differences, where the first two peptides are upregulated in early NSCLC subjects and the last five peptides are upregulated in healthy subjects (Fig 2F). Interestingly, this pattern was only observed in two NPs evaluated in this study, potentially suggesting NP-specific isoform selection.

To assess whether these two groups of peptides belong to different proteoforms, we further compared their abundance similarities across the 141 subjects. We expect peptides that belong to the same proteoform to have correlated abundances across a cohort of individuals since they belong to the same molecular entity while peptides belonging to different proteoforms should have less-correlated abundances across the same cohort of individuals. We thus performed pairwise Pearson correlation and hierarchical clustering analysis, which showed two distinct clusters driven by a high degree of correlation in peptide 1 and 2 (cluster 1) and peptides 3–7 (cluster 2) (Fig 2G). We next mapped the peptides to their genomic sequence, including four protein coding isoform transcripts (ENST00000397814, ENST00000354870, ENST00000306349, and ENST00000306385), and ordered them according to exon order (Fig 2H). We observed two distinct segments of corresponding direction of BMP1 peptide differential abundance. Specifically, peptides 1 and 2 were both upregulated in early NSCLC subjects (segment 1) and peptide 3–7 were all upregulated in healthy subjects (segment 2) (Fig 2F). Peptides mapping to segment 1 exclusively map to exons present in the short isoform (ENST00000397814), whereas peptides mapping to segment 2 exclusively map to exons present in the three longer isoforms (ENST00000354870, ENST00000306349, and ENST00000306385) (Fig 2H). The opposite pattern of abundance of the long and short isoforms in early NSCLC subjects versus healthy subjects suggest that BMP1 isoforms may play a role in cancer and may serve as a novel biomarker. This pattern is exaggerated when examining long and short isoforms in late-stage NSCLC subjects, where we observe a trend of increasing upregulation of segment 1 peptides (short BMP1 isoform) and a trend of decreasing upregulation of segment 2 peptides (long BMP1 isoform) between healthy subjects, Stage 1 and 2 NSCLC subjects, and Stage 3 and 4 NSCLC subjects (S4 Fig).

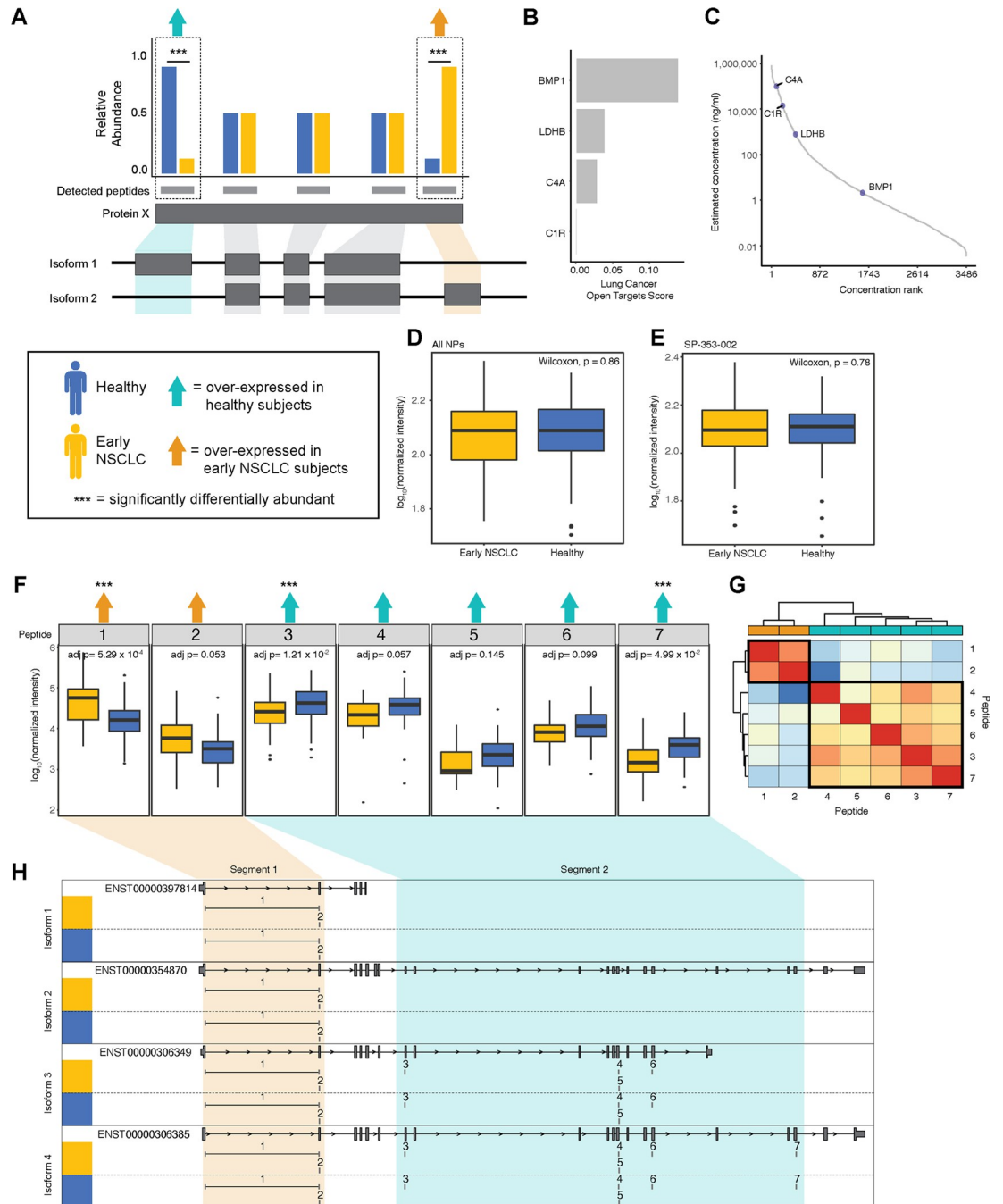


Fig 2. Identification of four proteoforms, including BMP1, in 141 healthy and early NSCLC subjects using a discordant peptide intensity search. A. Cartoon describing the discordant peptide intensity search strategy. We calculated DA across peptides between healthy (blue) and early NSCLC (yellow). Protein groups with at least one peptide significantly over-expressed (triple asterisks) in healthy subjects (teal arrow) and at least one peptide over-expressed in early NSCLC subjects (orange arrow) were identified as having putative proteoforms. Mapping the peptides to the gene structure, we inferred potential exon usage and segments suggesting the detection of more than one protein isoform. B. Barplot showing four proteins in which we potentially captured multiple protein isoforms: BMP1, C4A, C1R, and LDHB and their associated Open Target Score for lung carcinoma. C. Plot showing the four proteins with putative proteoforms matched to a reference database (HPPP) plotted as a distribution by the rank order of published concentrations (x-axis) and by the log₁₀ published concentration (ng/ml; y-axis). D. Box plot showing the log₁₀ median normalized intensities of BMP1 in early NSCLC subjects (yellow) and in healthy subjects (blue) with collapsed abundances across NPs. P-values, calculated using a Wilcoxon test, are shown. E. Box plot showing the log₁₀ median normalized intensities of BMP1 in early NSCLC subjects (yellow) and in healthy subjects (blue) in NP, SP-353-002. P-values, calculated using a Wilcoxon test, are shown. F. Series of boxplots showing the log₁₀ median normalized intensities of seven peptides mapping

BMP1 in early NSCLC (yellow) and healthy subjects (blue). Peptides that are over-expressed in healthy subjects are indicated with a teal arrow and in early NSCLC are indicated with an orange arrow. Peptides that are significantly DA are indicated with a triple asterisk. P-values, calculated using a Wilcoxon test and adjusted, are shown. G. Heatmap showing the Pearson correlation of the seven BMP1 peptide abundances, where low correlation is indicated in shades of blue and high correlation is indicated in shades of red. Correlation values were clustered using hierarchical clustering. Peptides are annotated by the direction of DA, including over-expressed in healthy subjects are highlighted in teal and early NSCLC are highlighted in orange. H. Gene structure plots of four known BMP1 protein coding transcripts (i.e., isoforms) with the seven BMP1 peptides mapped to genomic region. Peptides spanning intronic regions are indicated with a horizontal line. Peptides 1 and 2, corresponding to being over-expressed early NSCLC, are boxed in orange, creating one segment. Peptides 37, corresponding to being over-expressed healthy, are boxed in teal, creating a second segment. Segment 1 appears to correspond to the shorter isoform 1, whereas segment 2 appears to correspond to the longer isoforms 2–4.

<https://doi.org/10.1371/journal.pone.0282821.g002>

Discussion

Existing technologies, including an unbiased bottom-up NP-based methodology upstream of LC-MS/MS-based workflows and targeted methodologies, have enabled protein-centric analyses that have revealed new insights into human disease. While protein-centric bottom-up analyses have made substantial strides in our understanding of human biology, aggregating peptide level quantifications to the protein level may conceal biologically critical features, such as proteoforms arising from alternative splicing (protein isoform), allelic variation (protein variants), or post-translational modifications, which may provide mechanistic insights and novel biomarkers underlying complex traits and disease. Importantly, unbiased LC-MS/MS-based proteomic data can be re-mined to enable peptide-centric analyses that may reveal new information about proteoforms. In this study, we use peptide-level information derived from LC-MS/MS data to enable proteoform identification using discordant peptide abundance and apply that to a NSCLC cohort. Typically, protein inference engines use peptide-level data to detect the presence or absence of peptides to identify protein isoforms. However, here we show the utility of incorporating quantitative profiles of peptides mapping to known isoforms in potentially increasing the sensitivity of the underlying proteoform detection. Thus, we hypothesized that previously generated LC-MS/MS plasma proteomic data can be reanalyzed at the peptide-level using quantitative profiles to infer protein isoforms [15], potentially yielding deeper insights into disease mechanisms and we demonstrated that such a reanalysis revealed known and putative, novel disease-relevant proteoforms.

We performed peptide analysis using DIA data derived from healthy and early NSCLC subjects by conducting a discordant peptide intensity search to identify protein isoforms. We identified four proteins with DA peptides and putative isoforms, including BMP1, C4A, C1R, and LDHB. Importantly, none of these proteins showed a difference in abundance at the protein-level. For BMP1 and C1R, using peptide abundance as a proxy for functionally relevant proteins we identified potential NSCLC-related isoforms. We showed BMP1 has differential abundance of two isoforms (long and short), both of which have higher magnitude of differential abundance at later stages of NSCLC. BMP1 plays a role in collagen processing and the short isoform lacks the domains enabling its secretion, potentially impacting collagen's protective role in cancer consistent with the higher abundance of the short isoform in cancer subjects observed in this paper. Additionally, C4A showed distinct peptide abundance discordance in one segment of the protein, which did not correspond to any known protein coding isoforms, suggesting peptide-centric proteoform identification may result in novel disease-associated isoforms.

The method we used to search for protein isoforms through discordant peptide intensity is stringent in terms of the number of protein isoform candidates we can find, but easily interpretable. Similar approaches such as COPF [19] and PeCorA [20] use quantitative disagreements between peptides mapped to the same protein or peptide correlation within the same protein to detect protein isoforms and suggest proteoforms. However, as shown with our

examples where 2 of the 4 isoform candidates (C1R and LDHB) met the discordant peptide intensity criteria but failed to be readily explained by known isoforms or biological conjecture, evaluation of the validity of the isoform candidate is needed but is outside the scope of this study. It is also possible that proteoforms derived from post translational modifications rather than protein isoforms are responsible for the observed differences. In this paper, our validation was mapping the peptides back to the genomic sequence and known isoform transcripts. Manual validation (e.g., isoform specific enrichment with isoform specific antibodies) can confirm the presence of novel isoforms. This might be possible for the limited candidates arising from our stringent isoform detection process, however, other processes such COPF and PeCorA could yield dramatically more candidates.

It is possible that the finding of only four protein isoforms in 188 subjects has been impacted by limited sample sizes reducing our power to identify proteoforms. Similarly, it is also possible that other undiscovered proteoforms are not functional in plasma and may only be identified in other biofluids or tissues. While our study shows the utility of using NP-based methodology upstream LC-MS/MS-based workflows to identify proteoforms, it is possible that expanding the sample size and diversity in sample type may yield further insights into disease-associated proteoforms. In addition, LC-MS/MS enables quantifying and identifying tens of thousands of peptides with post-translational modifications precisely defined by their intact mass and fragmentation pattern.

The identification of proteoforms (protein isoforms) highlights important considerations for current approaches characterizing the impact of genetic variation on molecular phenotypes, like protein abundance, by conducting protein quantitative trait analyses (pQTLs). Specifically, recent pQTL analyses using large cohorts [21] are performed at the protein-level and largely miss or misattribute peptide-level proteoform effects. Furthermore, these studies utilize aptamer and antibody-based methodologies that, as been recently shown [21], can lead to false discoveries and uncertain identification error rates because of conceptual limitations (e.g., the presence of a non-synonymous SNP inducing an amino acid change that disrupts the binding of the aptamer or antibody). Interrogating protein abundances with this high resolution approach provides deeper insight into the molecular mechanisms underlying human biology and opens a possible new avenue for biomarker identification and therapeutic development.

Materials and methods

Identification of protein isoforms

As previously reported, plasma from healthy subjects and from subjects diagnosed with NSCLC at stage 1, 2, 3, and 4 was collected and processed with the Proteograph™ workflow [15] and DIA data was generated and processed (S1 File). From the 1,565 proteins present after filtering, we searched for peptides that had differential abundance between controls and cancer ($p < 0.05$; where p-values were adjusted using Benjamini-Hochberg correction). Discordant pairs are defined as peptides from the same protein where at least one peptide was identified with significantly higher, and another peptide was identified with significantly lower plasma abundance in healthy controls vs. early NSCLC.

This work generated no additional data from new or existing patient samples and used raw data already deposited in the public database ProteomeXchange Consortium (<http://proteomecentral.proteomexchange.org>) with the dataset PXD017052.

Protein quantification across multiple samples

Within each nanoparticle, standard MaxLFQ was used to quantify abundance at the protein level. For each peptide, the intensity ratios between every pair of samples were first computed.

The pairwise protein ratio is then defined as the median of the peptide ratios from all peptides map to the same protein. With all the pairwise protein ratios between any two samples, we can perform a least-squares analysis to reconstruct the abundance profile optimally satisfying all the protein ratios. Then the whole profile is rescaled to the cumulative intensity across samples for the final protein abundance [22]. A modified MaxLFQ was used to quantify abundance across samples and nanoparticles. For each protein, all peptides' intensities belonging to a protein from all samples and NP were employed to calculate peptide ratios and subsequent calculation steps resulting in abundance across all samples and NP.

Supporting information

S1 Fig. Identification of C4A proteoform in 141 healthy and early NSCLC subjects using a discordant peptide intensity search. A. Box plot showing the \log_{10} median normalized intensities of C4A in early NSCLC subjects (yellow) and in healthy subjects (blue) with collapsed abundances across NPs. P-values, calculated using a Wilcoxon test, are shown. B. Box plot showing the \log_{10} median normalized intensities of C4A in early NSCLC subjects (yellow) and in healthy subjects (blue) in NP, SP-006-030. P-values, calculated using a Wilcoxon test, are shown. C. Heatmap showing the Pearson correlation of the 64 C4A peptide abundances, where low correlation is indicated in shades of blue and high correlation is indicated in shades of red. Correlation values were clustered using hierarchical clustering. Peptides are annotated by the direction of DE, including over-expressed in healthy subjects are highlighted in teal and early NSCLC are highlighted in orange. D. Series of boxplots showing the \log_{10} median normalized intensities of 64 peptides mapping C4A in early NSCLC (yellow) and healthy subjects (blue). Peptides that are over-expressed in healthy subjects are indicated with a teal arrow and in early NSCLC are indicated with an orange arrow. Peptides that are significantly DE are indicated with a triple asterisk. P-values, calculated using a Wilcoxon test and adjusted, are shown. E. Gene structure plots of 2 known C4A protein coding transcripts (i.e., isoforms) with the 64 C4A peptides mapped to genomic region. Peptides spanning intronic regions are indicated with a horizontal line. Peptides 1–39 (except peptide 32), corresponding to being over-expressed in healthy subjects, are boxed in teal, creating one segment. Peptides 40–53, corresponding to being over-expressed early NSCLC, are boxed in orange, creating a second segment. Peptides 54–63 (except peptide 64), corresponding to being over-expressed in healthy subjects, are boxed in teal, creating a third segment. Segment patterns do not appear to correspond to any known protein isoforms, potentially indicating a novel isoform. (TIF)

S2 Fig. Identification of C1R proteoform in 141 healthy and early NSCLC subjects using a discordant peptide intensity search. A. Box plot showing the \log_{10} median normalized intensities of C1R in early NSCLC subjects (yellow) and in healthy subjects (blue) with collapsed abundances across NPs. P-values, calculated using a Wilcoxon test, are shown. B. Box plot showing the \log_{10} median normalized intensities of C1R in early NSCLC subjects (yellow) and in healthy subjects (blue) in NP, SP-353-002. P-values, calculated using a Wilcoxon test, are shown. C. Heatmap showing the Pearson correlation of the 17 C1R peptide abundances, where low correlation is indicated in shades of blue and high correlation is indicated in shades of red. Correlation values were clustered using hierarchical clustering. Peptides are annotated by the direction of DE, including over-expressed in healthy subjects are highlighted in teal and early NSCLC are highlighted in orange. D. Series of boxplots showing the \log_{10} median normalized intensities of 17 peptides mapping C1R in early NSCLC (yellow) and healthy subjects (blue). Peptides that are over-expressed in healthy subjects are indicated with

a teal arrow and in early NSCLC are indicated with an orange arrow. Peptides that are significantly DE are indicated with a triple asterisk. P-values, calculated using a Wilcoxon test and adjusted, are shown. Gene structure plots of 6 known CIR protein coding transcripts (i.e., isoforms) with the 17 CIR peptides mapped to genomic region. Peptides spanning intronic regions are indicated with a horizontal line. Peptides corresponding to being over-expressed in healthy subjects are boxed in teal. Peptides corresponding to being over-expressed early NSCLC are boxed in orange.

(TIF)

S3 Fig. Identification of LDHB proteoform in in 141 healthy and early NSCLC subjects using a discordant peptide intensity search. A. Box plot showing the \log_{10} median normalized intensities of LDHB in early NSCLC subjects (yellow) and in healthy subjects (blue) with collapsed abundances across NPs. P-values, calculated using a Wilcoxon test, are shown. B. Box plot showing the \log_{10} median normalized intensities of LDHB in early NSCLC subjects (yellow) and in healthy subjects (blue) in NP, SP-006-030. P-values, calculated using a Wilcoxon test, are shown. C. Heatmap showing the Pearson correlation of the 12 LDHB peptide abundances, where low correlation is indicated in shades of blue and high correlation is indicated in shades of red. Correlation values were clustered using hierarchical clustering. Peptides are annotated by the direction of DE, including over-expressed in healthy subjects are highlighted in teal and early NSCLC are highlighted in orange. D. Series of boxplots showing the \log_{10} median normalized intensities of 12 peptides mapping LDHB in early NSCLC (yellow) and healthy subjects (blue). Peptides that are over-expressed in healthy subjects are indicated with a teal arrow and in early NSCLC are indicated with an orange arrow. Peptides that are significantly DE are indicated with a triple asterisk. P-values, calculated using a Wilcoxon test and adjusted, are shown. E. Gene structure plots of 6 known LDHB protein coding transcripts (i.e., isoforms) with the 12 LDHB peptides mapped to genomic region. Peptides spanning intronic regions are indicated with a horizontal line. Peptides corresponding to being over-expressed in healthy subjects are boxed in teal. Peptides corresponding to being over-expressed early NSCLC are boxed in orange.

(TIF)

S4 Fig. Association of NSCLC stage (1, 2, 3, and 4) and the presence of BMP1 proteoforms using a discordant peptide intensity search. Series of boxplots showing the \log_{10} median normalized intensities of 7 peptides mapping BMP1 in healthy subjects (blue, stages 1 and 2 NSCLC (yellow), and stages 3 and 4 (red). *Peptides that are significantly DE (Wilcoxon test and adjusted).

(TIF)

S1 File. Supplementary information. File containing the supplementary results and supplementary methods, with the associated supplementary references.

(PDF)

Acknowledgments

We thank the individuals who took part in this study for their participation. We also thank Steven A. Carr and Robert Langer for their review of this manuscript.

Author Contributions

Conceptualization: Margaret K. R. Donovan, Yingxiang Huang, John E. Blume, Daniel Hornburg, Serafim Batzoglou, Luis A. Diaz, Omid C. Farokhzad, Asim Siddiqui.

Formal analysis: Margaret K. R. Donovan, Yingxiang Huang, John E. Blume, Jian Wang, Daniel Hornburg, Shadi Ferdosi, Iman Mohtashemi, Sangtae Kim, Marwin Ko, Ryan W. Benz, Theodore L. Platt, Asim Siddiqui.

Writing – original draft: Margaret K. R. Donovan, Yingxiang Huang.

Writing – review & editing: Margaret K. R. Donovan, Yingxiang Huang, John E. Blume, Jian Wang, Daniel Hornburg, Shadi Ferdosi, Iman Mohtashemi, Marwin Ko, Ryan W. Benz, Theodore L. Platt, Serafim Batzoglou, Luis A. Diaz, Asim Siddiqui.

References

1. Smith LM, Kelleher NL, Proteomics TC for TD. Proteoform: a single term describing protein complexity. *Nat Methods*. 2014; 10(3):186–7.
2. Li YI, Van De Geijn B, Raj A, Knowles DA, Petti AA, Golan D, et al. RNA splicing is a primary link between genetic variation and disease. *Science*. 2016 Apr 29; 352(6285):600–4. <https://doi.org/10.1126/science.aad9417> PMID: 27126046
3. Lisitsa A, Moshkovskii S, Chernobrovkin A, Ponomarenko E, Archakov A. Profiling proteoforms: promising follow-up of proteomics for biomarker discovery. *Expert Rev Proteomics*. 2014 Feb; 11(1):121–9. <https://doi.org/10.1586/14789450.2014.878652> PMID: 24437377
4. Satpathy S, Krug K, Jean Beltran PM, Savage SR, Petralia F, Kumar-Sinha C, et al. A proteogenomic portrait of lung squamous cell carcinoma. *Cell*. 2021; 184(16):4348–4371.e40. <https://doi.org/10.1016/j.cell.2021.07.016> PMID: 34358469
5. Kisluk J, Ciborowski M, Niemira M, Kretowski A, Niklinski J. Proteomics biomarkers for non-small cell lung cancer. *J Pharm Biomed Anal*. 2014 Dec; 101:40–9. <https://doi.org/10.1016/j.jpba.2014.07.038> PMID: 25175018
6. Nishimura T, Végvári Á, Nakamura H, Kato H, Saji H. Mutant Proteomics of Lung Adenocarcinomas Harboring Different EGFR Mutations. *Front Oncol*. 2020; 10(August):1–20. <https://doi.org/10.3389/fonc.2020.01494> PMID: 32983988
7. Nishimura T, Nakamura H, Végvári Á, Marko-Varga G, Furuya N, Saji H. Current status of clinical proteogenomics in lung cancer. *Expert Rev Proteomics*. 2019 Sep; 16(9):761–72. <https://doi.org/10.1080/14789450.2019.1654861> PMID: 31402712
8. Franc V, Zhu J, Heck AJR. Comprehensive Proteoform Characterization of Plasma Complement Component C8αβγ by Hybrid Mass Spectrometry Approaches. *J Am Soc Mass Spectrom*. 2018 Jun 1; 29(6):1099.
9. Gao J, McCann A, Laupsa-Borge J, Nygård O, Ueland PM, Meyer K. Within-person reproducibility of proteoforms related to inflammation and renal dysfunction. *Sci Rep*. 2022 Dec 1; 12(1). <https://doi.org/10.1038/s41598-022-11520-1> PMID: 35523986
10. Koska J, Furtado J, Hu Y, Sinari S, Budoff MJ, Billheimer D, et al. Plasma proteoforms of apolipoproteins C-I and C-II are associated with plasma lipids in the Multi-Ethnic Study of Atherosclerosis. *J Lipid Res*. 2022 Sep; 63(9):100263. <https://doi.org/10.1016/j.jlr.2022.100263> PMID: 35952903
11. Wählén K, Ernberg M, Kosek E, Mannerkorpi K, Gerdle B, Ghafouri B. Significant correlation between plasma proteome profile and pain intensity, sensitivity, and psychological distress in women with fibromyalgia. *Sci Rep*. 2020 Dec 27; 10(1):12508. <https://doi.org/10.1038/s41598-020-69422-z> PMID: 32719459
12. Anderson NL. The clinical plasma proteome: a survey of clinical assays for proteins in plasma and serum. *Clin Chem*. 2010 Feb; 56(2):177–85. <https://doi.org/10.1373/clinchem.2009.126706> PMID: 19884488
13. Geyer PE, Holdt LM, Teupser D, Mann M. Revisiting biomarker discovery by plasma proteomics. *Mol Syst Biol*. 2017; 13(9):942.
14. Geyer PE, Kulak NA, Pichler G, Holdt LM, Teupser D, Mann M. Plasma Proteome Profiling to Assess Human Health and Disease. *Cell Syst*. 2016; 2(3):185–95. <https://doi.org/10.1016/j.cels.2016.02.015> PMID: 27135364
15. Blume JE, Manning WC, Troiano G, Hornburg D, Figa M, Hesterberg L, et al. Rapid, deep and precise profiling of the plasma proteome with multi-nanoparticle protein corona. *Nat Commun*. 2020; 11(1):1–14.
16. Zhu Y, Orre LM, Tran YZ, Mermelekas G, Johansson HJ, Malyutina A, et al. DEqMS: A method for accurate variance estimation in differential protein expression analysis. *Mol Cell Proteomics*. 2020; 19(6):1047–57. <https://doi.org/10.1074/mcp.TIR119.001646> PMID: 32205417

17. Koscielny G, An P, Carvalho-Silva D, Cham JA, Fumis L, Gasparyan R, et al. Open Targets: A platform for therapeutic target identification and Validation. *Nucleic Acids Res.* 2017; 45(D1):D985–94. <https://doi.org/10.1093/nar/gkw1055> PMID: 27899665
18. Schwenk JM, Omenn GS, Sun Z, Campbell DS, Baker MS, Overall CM, et al. The Human Plasma Proteome Draft of 2017: Building on the Human Plasma PeptideAtlas from Mass Spectrometry and Complementary Assays. *J Proteom Res.* 2017; 16(12):4299–310. <https://doi.org/10.1021/acs.jproteome.7b00467> PMID: 28938075
19. Bludau I, Frank M, Dörig C, Cai Y, Heusel M, Rosenberger G, et al. Systematic detection of functional proteoform groups from bottom-up proteomic datasets. *Nat Commun.* 2021 Dec 1; 12(1). <https://doi.org/10.1038/s41467-021-24030-x> PMID: 34155216
20. Dermit M, Peters-Clarke TM, Shishkova E, Meyer JG. Peptide Correlation Analysis (PeCorA) Reveals Differential Proteoform Regulation. *J Proteome Res.* 2021; 20(4):1972–80. <https://doi.org/10.1021/acs.jproteome.0c00602> PMID: 33325715
21. Pietzner M, Wheeler E, Carrasco-Zanini J, Kerrison ND, Oerton E, Koprulu M, et al. Synergistic insights into human health from aptamer- and antibody-based proteomic profiling. *Nat Commun.* 2021 Dec 1; 12(1). <https://doi.org/10.1038/s41467-021-27164-0> PMID: 34819519
22. Cox J, Hein MY, Lubner CA, Paron I, Nagaraj N, Mann M. Accurate proteome-wide label-free quantification by delayed normalization and maximal peptide ratio extraction, termed MaxLFQ. *Mol Cell Proteomics.* 2014; 13(9):2513–26. <https://doi.org/10.1074/mcp.M113.031591> PMID: 24942700