# Supplemental Information

## EXTENDED EXPERIMENTAL PROCEDURES

### Specimen Acquisition

The tumor samples used in this manuscript are from The Cancer Genome Atlas (TCGA), and were previously characterized from a genomic perspective (TCGA Research Network, 2011). Biospecimens were collected from newly diagnosed patients with ovarian serous adenocarcinoma who were undergoing surgical resection and had received no prior treatment for their disease, including chemotherapy or radiotherapy. All cases had to be of serous histology but were collected regardless of surgical stage or histologic grade. Cases were staged according to the 1988 International Federation of Gynecology and Obstetrics (FIGO) staging system. Each tumor specimen was approximately 1 cm$^3$ in size and weighed between 100 mg and 200 mg, in general. The specimens were first used for integrated genomic analysis as part of the TCGA program. Residual material from the same specimens were analyzed as part of this Clinical Proteomic Tumor Analysis Consortium (CPTAC) project. Each specimen was embedded in optimal cutting temperature (OCT) medium and histologic sections were obtained from top and bottom portions for review. Each case was reviewed through TCGA by a board-certified pathologist to confirm that the frozen section was histologically consistent with ovarian serous adenocarcinoma. The most recent pathological re-review of the tumor samples used in this report reclassified five tumors ("TCGA-09-2056", "TCGA-24-1544", "TCGA-24-1565", "TCGA-25-1316" and "TCGA-61-2095") as other than HGSC (Vang et al., 2016), and these tumors were removed from the subtyping, survival, chromosomal instability (CIN) and differential dependency network (DNN) analyses. The top and bottom sections had to contain an average of 70% tumor cell nuclei with less than 20% necrosis. Specimens were shipped overnight from one of 15 tissue source sites using a cryoport that maintained an average temperature of less than -180°C. Use of the ESTIMATE algorithm (Yoshihara et al., 2013) to calculate a Tumor Purity score for the paired samples analyzed at JHU and PNNL indicated a statistically significant correlation (r = 0.83; $p$ value = 3.8e-9) between the two sites. When the tumor purity score was compared for the same tumor analyzed for transcriptomics and proteomics, there correlation was statistically significant (r = 0.74, $p$ value = 2.2e-16).

### Clinical Data Annotation

Clinical data were obtained from TCGA tissue source site through data collection forms. Data forms were entered electronically at the TCGA biospecimen core resource and XML files were generated. Clinical data can be accessed and downloaded from the TCGA Data Portal at http://cancergenome.nih.gov/. Demographics, histopathologic information, treatment details including chemotherapy drugs, doses and routes of administration, and outcome parameters were collected. Supplemental clinical data were collected directly from the tissue source sites as part of the CPTAC program and updated clinical data are provided in Table S1.

The definitions of most clinical variables were implicit and select variables were defined as follows: TUMORRESIDUALDISEASE was defined as the size of residual disease at the conclusion of the primary surgical procedure. This field was used to define surgical cytoreduction as optimal or suboptimal. Optimal was defined as no residual disease greater than 1cm and included the variable categories of no macroscopic disease (*i.e.* microscopic residual disease) and 1 to 10 mm. Suboptimal was defined as residual disease greater than 1cm and included the variable categories of 11 to 20 mm and greater than 20mm. PRIMARYTHERAPYOUTCOMESUCCESS was defined as the response to treatment determined after primary surgery and subsequent adjuvant chemotherapy. PERSONNEOPLASMCANCERSTATUS was defined as the last known status of disease. For the purpose of these analyses, the date of surgery was used as a surrogate for the date of initial diagnosis, since treatment planning and intervention for these cases undergoing initial surgical resection began at that time point. Overall survival was defined as the interval from the date of initial surgical resection to the date of last known contact or death. Progression free survival was defined as the interval from the date of initial surgical resection to the date of progression, date of recurrence, or date of last known contact if the patient was alive and has not recurred. For the purpose of these analyses, any patient who had died without a date of progression or recurrence was excluded from analyses of progression free survival.

Chemotherapy treatment details were reviewed to identify dates of adjuvant therapy. The date of last primary platinum treatment was also determined from the available chemotherapy details and included adjuvant therapy and consolidation treatment when given consecutively following adjuvant therapy. Some treatment data were directly obtained from the Tissue Specimen Sites to supplement this CPTAC project. The platinum free interval was defined as the interval from the date of last primary platinum treatment to

## Supplemental Information

the date of progression, date of recurrence, or date of last known contact if the patient is alive and has not recurred. Platinum status was defined as resistant if the platinum free interval was less than six months and the patient had progressed or recurred. Platinum status was defined as sensitive if the platinum free interval was six months or greater, there was no evidence of progression or recurrence, and the follow-up interval was at least six months from the date of last primary platinum treatment. Patients who have not progressed or recurred and who have been followed for less than six months from the date of last primary platinum treatment were excluded from analyses regarding platinum status and their platinum status is indicated as 'unable to determine'.

Tumors were selected for putative homologous recombination deficiency (HRD), defined by the presence of germline or somatic *BRCA1* or *BRCA2* mutations, *BRCA1* promoter methylation, or homozygous deletion of *PTEN* (Woodbine et al., 2014). Among the 122 samples selected, 67 were classified as HRD whereas 55 did not meet the above criteria (non-HRD). Of the 55 non-HRD cases, two samples had a mutation in other genes associated with homologous recombination deficiency; one each in *ATM* and *PALB2*; out of the 11 common HRD genes surveyed (*ATM, BARD1, BRIP1, CHEK1, CHEK2, FAM175A, MRE11A, NBN, PALB2, RAD51C, and RAD51D*) (Pennington et al., 2014). We have compared results with and without the samples containing mutations in *ATM* and *PALB2*, and observed no change in the final results.

Standard statistical tests were used to analyze the clinical data including, but not limited to Student's t-test, Fisher's exact test, log-rank test, and Cox proportional hazard analysis, as appropriate. Descriptive statistics were also included. All statistical tests were two-sided, corrected for multiple hypothesis testing as appropriate, and statistical significance was considered when $p$ value $< 0.05$. Analyses of clinical data were performed in R (version 3.0.1) (R Development Core Team, 2008). Clinical data were available for the 174 patients included in this report. Clinical variables are provided in Table S1.

### Additional Clinical Data Results

As shown in Table S1, the characteristics of the TCGA ovarian cases reflect the general population of women with advanced ovarian cancer.

The average age at diagnosis was 60.5 years, all cases were of serous histology, 71% of the cases were FIGO stage IIIC, 16% were FIGO stage IV, 83% were grade 3, and 64% had optimal surgical cytoreduction. Univariate correlations between select clinical variables and progression-free survival (PFS) or overall survival (OS) are shown in Table S1. Age at diagnosis and platinum status were associated with both PFS and OS. The median PFS and OS was 15.8 and 42.6 months for FIGO stage III patients and 13.1 and 37.1 months for FIGO stage IV patients, respectively ($p$ value 0.12 for PFS and $p$ value 0.48 for OS). The median PFS and OS were 14.9 and 42.6 months for optimally debulked patients and 15.4 and 39.9 months for suboptimally debulked patients, respectively ($p$ value 0.49 for PFS and $p$ value 0.81 for OS). The median OS was 56.8 months for platinum sensitive patients and 25.8 months for platinum resistant patients ($p$ value 4.8e-11).

An association between surgical cytoreduction and platinum sensitivity has been previously reported (Eisenhauer et al., 2008). We explored the relationship between surgical cytoreduction and platinum sensitivity in the TCGA ovarian cases. We found no association between platinum status and surgical cytoreduction when defined traditionally as optimal or suboptimal. However, when considering microscopic residual disease separately from other optimally or suboptimally debulked patients, there was an association between surgical outcome and platinum status. Patients with microscopic residual disease were more likely to be platinum sensitive than patients with more than microscopic residual disease ($p$ value = 0.0028; Table S1). Odds ratios could not be calculated since all patients with microscopic residual disease were platinum sensitive. These data suggest that surgical cytoreduction may have a direct impact on platinum status.

### Sample Preprocessing

All samples for the current study were obtained through the TCGA Biospecimen Core Resource as described above, and processed for mass spectrometric (MS) analysis at either of the two CPTAC proteomic characterization centers participating in this study: Pacific Northwest National Laboratory (PNNL) or the Johns Hopkins Medical Institutions of Johns Hopkins University (JHU). The samples used for this project were obtained from the same anatomic site as the samples used by TCGA. For 168 (96.6%) of the 174 samples, the two specimens were either directly adjacent or no more than 5 mm apart; the 32 samples analyzed in parallel at both JHU and PNNL meet these criteria, and are identified in Table S1. For

# Supplemental Information

six (3.4%) samples the distance between specimens may have been as great as 10 mm. The PNNL ovarian cancer tumor samples were first processed by cryopulverization. Briefly, tumor pieces were transferred into pre-cooled Covaris Tissue-Tube 1 Extra bags (Covaris) and processed in a Covaris CP02 Cryoprep device using an impact setting of 3 as all tumor tissue wet weights were less than 100 mg. The tissue powder was then transferred into precooled cryovials (Corning). All procedures were carried out on dry ice and liquid nitrogen to maintain tissue in a powdered, frozen state. JHU samples were extracted and digested directly without cryopulverization.

## Protein Extraction and Tryptic Digestion

At PNNL, approximately 50 mg of each of the pulverized TCGA ovarian tumor tissues were homogenized separately in 600 µL of lysis buffer (8 M urea, 100 mM $NH_4HCO_3$, pH 7.8, 0.1% NP-40, 0.5% sodium deoxycholate, 10 mM NaF, phosphatase inhibitor cocktails 2 and 3, 20 µM PUGNAc). Protein concentrations of tissue lysates were determined by BCA assay (Pierce). Proteins were reduced with 5 mM dithiothreitol for 1 h at 37 °C, and subsequently alkylated with 10 mM iodoacetamide for 1 hour at RT in the dark. Samples were diluted 1:2 with nanopure water, 1 mM $CaCl_2$ and digested with sequencing grade modified trypsin (Promega) at 1:50 enzyme-to-substrate ratio. After 4 h of digestion at 37 °C, samples were diluted 1:4 with the same buffers and another aliquot of the same amount of trypsin was added to the samples and further incubated at room temperature (RT) overnight (~16 h). The digested samples were then acidified with 10% trifluoroacetic acid to ~pH 3. Tryptic peptides were desalted on strong cation exchange (SCX) solid-phase extraction (SPE) (SUPELCO, Discovery-SCX) and reversed-phase C18 SPE columns (SUPELCO Discovery) and dried using Speed-Vac.

At JHU, approximately 50 mg of each of the sectioned TCGA ovarian tumor tissues were sonicated separately in 1.5 mL of 8 M urea, 0.8 M $NH_4HCO_3$, pH 8.0. Proteins were reduced with 10 mM tris (2-carboxyethyl) phosphine (TCEP) for 1 h at 37 °C, and subsequently alkylated with 12 mM iodoacetamide for 1 h at RT in the dark. Samples were diluted 1:4 with deionized water and digested with sequencing grade modified trypsin at 1:50 enzyme-to-protein ratio. After 12 h of digestion at 37 °C, another aliquot of the same amount of trypsin was added to the samples and further incubated at 37 °C overnight. The digested samples were then acidified, cleaned up (SCX and C18) and dried as described above.

## iTRAQ Labeling of Peptides

Quantitative proteomic comparisons were based on iTRAQ reporter ion intensities (Ross et al., 2004). Currently, the CPTAC consortium uses two major LC-MS/MS based methods for quantitative proteomics, label-free quantification (Zhang et al., 2014) and isobaric stable isotope labeling approaches such as iTRAQ (Mertins et al., 2014). Compared to label-free quantification, iTRAQ is widely used to reduce the variation in quantification attributable to differences in instrument performances and to provide simultaneous, integrated measurements of protein and PTM levels.  In addition, the multiplexed analysis of multiple independent samples possible with iTRAQ results in deeper quantification of particularly the phosphoproteome.   Desalted peptides were labeled with 4-plex iTRAQ (isobaric tag for relative and absolute quantitation) reagents according to the manufacturer's instructions (AB Sciex). At PNNL, peptides (500 µg) from each of the TCGA tumors were dissolved in 150 µL of 0.5 M triethylammonium bicarbonate (TEAB), pH 8.5 solution, and mixed with 5 units of iTRAQ reagent that was dissolved freshly in 350 µL of ethanol. A reference sample was created by pooling an aliquot from each individual PNNL tumor sample, and Channel 117 was used for labeling the pooled reference sample throughout the TCGA sample analysis. After 1 h incubation at RT, 1.5 mL of water was added and incubated for 30 min at RT to stop the reaction and hydrolyze the unreacted iTRAQ reagents. Peptides labeled by different iTRAQ reagents were then mixed and concentrated to ~500 µL, and were desalted on C18 SPE columns.

The JHU samples were labeled similarly, except that Channel 114 was used for labeling the reference sample created by pooling an aliquot from each individual JHU tumor sample. A single JHU ovarian tumor sample (independent of the TCGA sample collection) was also prepared and repeatedly analyzed (10 aliquots were co-randomized with the 122 TCGA samples for iTRAQ labeling) the same way as the TCGA samples as an internal quality control (QC).

## Peptide Fractionation by Basic Reversed-phase Liquid Chromatography (bRPLC)

Both sites (PNNL and JHU) used extensive fractionation by basic reversed phase liquid chromatography to reduce sample complexity and thus reduce the likelihood of peptides being co-isolated and co-fragmented;

# Supplemental Information

this approach has been well-documented to reduce iTRAQ reporter ion ratio distortion effects (Bantscheff et al., 2008; Ow et al., 2011). The percentage of statistically significantly regulated proteins and phosphopeptides detected by this iTRAQ approach is better than obtained from the application of alternative label-free method (Tabb et al., 2016), and comparable to another similar approach using mTRAQ (Mertins et al., 2012).

At PNNL, the 4-plex iTRAQ labeled sample was separated on a Waters reversed phase XBridge C18 column (250 mm × 4.6 mm column containing 5-µm particles, and a 4.6 mm × 20 mm guard column) using Agilent 1200 HPLC System as previously reported (Wang et al., 2011). After the sample loading, the C18 column was washed for 35 min with solvent A (10 mM TEAB, pH 7.5), before applying a 102-min liquid chromatography (LC) gradient with solvent B (10 mM TEAB, pH 7.5, 90% acetonitrile). The LC gradient started with a linear increase of solvent A to 10% B in 6 min, then linearly increased to 30% B in 86 min, 10 min to 42.5% B, 5 min to 55% B and another 5 min to 100% solvent B. The flow rate was 0.5 mL/min. A total of 96 fractions were collected into a 96 well plate throughout the LC gradient. These fractions were concatenated into 24 fractions by combining 4 fractions that are 24 fractions apart (i.e., combining fractions #1, #25, #49, and #73; #2, #26, #50, and #74; and so on). JHU samples were also fractionated using bRPLC with slightly modified conditions. Approximately 100 µg of 4-plex iTRAQ labeled sample was separated on a reversed phase Zorbax extend-C-18 column (4.6 x 100 mm column containing 1.8-µm particles; Agilent) using an Agilent 1220 Infinity HPLC System. The solvent consisted of 10 mM ammonium formate (pH 10) as mobile phase A and 10 mM ammonium formate and 90% acetonitrile (pH 10) as mobile phase B. The separation gradient was set as follows: 2% B for 10 min, from 2 to 8% B for 5 min, from 8 to 35% B for 85 min, from 35 to 95% B for 5 min, and 95% B for 25 min. Ninety-six fractions were collected across the LC separation and were concatenated into 24 fractions by combining fractions 1, 25, 49, 73; 2, 26, 50, 74; and so on. The samples were dried in a Speed-Vac and stored at −80°C until analysis using LC coupled to tandem mass spectrometry (LC-MS/MS) analysis.

For proteome analysis at PNNL, 10% of each concatenated fraction was dried down and re-suspended in 0.1% trifluoroacetic acid to a peptide concentration of 0.15 µg/µL for LC-MS/MS analysis. The rest of the concatenated fractions (90%) were further concatenated into 12 fractions by combining two well separated fractions (i.e., combining concatenated fractions #1 and #13; #2 and #14; and so on), concentrated down, and subjected to immobilized metal affinity chromatography (IMAC) for phosphopeptide enrichment at PNNL.

## Phosphopeptide Enrichment Using IMAC

Magnetic $Fe^{3+}$-NTA-agarose beads were freshly prepared using Ni-NTA-agarose beads (QIAGEN) for phosphopeptide enrichment. For each of the 12 fractions, peptides were reconstituted in 200 µL IMAC binding/wash buffer (80% acetonitrile, 0.1% formic acid) and incubated with 50 µL of the 5% bead suspension for 30 min at RT. After incubation, the beads were washed 3 times each with 200 µL of wash buffer. Phosphopeptides were eluted from the beads using 75 µL of 1:1 (acetonitrile:5% ammonia in 5 mM phosphate buffer, pH 8) mixed buffer, pH 10 after incubating for 5 min at RT. Samples were acidified to ~pH 3.5 and concentrated to 5-10 µL, and were reconstituted to 30 µL with 0.1% trifluoroacetic acid for LC-MS/MS analysis.

## LC-MS/MS for Global Proteome Analysis

At PNNL, the LC system was custom built using two Agilent 1200 nanoflow pumps, various Valco valves (Valco Instruments Co.), and a PAL autosampler (Leap Technologies). Full automation was made possible by custom software that allows for parallel event coordination and faster MS duty cycle through use of two analytical columns. Reversed-phase columns were prepared in-house by slurry packing 3-µm Jupiter $C_{18}$ (Phenomenex) into 35-cm x 360 µm o.d. x 75 µm i.d fused silica (Polymicro Technologies Inc.) using a 1-cm sol-gel frit for media retention. Mobile phase flow rate was 300 nL/min and consisted of 0.1% formic acid in water (A) and 0.1% formic acid in acetonitrile (B) with a gradient profile as follows (min:%B); 0:5, 1:10, 85:28, 93:60, 98:75, 100:75. Sample injection (5 µL) occurred 40 min prior to beginning the gradient while MS data acquisition lagged gradient start and end times by 10 min to account for column dead volume. Aside from the improved duty cycle gained in using two columns, it also allowed each column to be 'washed' (using multiple fast gradients) after each run and re-generated off-line without any cost to duty cycle.

MS analysis was performed using a Thermo Scientific LTQ Orbitrap Velos mass spectrometer outfitted with a custom electrospray ionization interface. Electrospray emitters were custom made using

## Supplemental Information

150 um o.d. x 20 μm i.d. chemically etched fused silica. The ion transfer tube temperature and spray voltage were 300 ºC and 1.8 kV, respectively. Orbitrap spectra (AGC 3x10$^6$) were collected from 300-1800 m/z at a resolution of 30,000 followed by data-dependent higher-energy C-trap dissociation (HCD) MS/MS (centroid mode, at a resolution of 7,500, collision energy 45%, activation time 0.1 ms, AGC 5x10$^4$) of the ten most abundant ions using an isolation width of 2.5 Da. Charge state screening was enabled to reject unassigned and singly charged ions. A dynamic exclusion time of 30 sec was used to discriminate against previously selected ions (within -0.55 Da to 2.55 Da).

At JHU, peptides were separated on a Dionex Ultimate 3000 RSLCnano system (Thermo Scientific) with a 75 μm x 15 cm Acclaim PepMap100 separating column (Thermo Scientific) protected by a 2 cm guarding column (Thermo Scientific). Mobile phase flow rate was 300 nL/min and consisted of 0.1% formic acid in water (A) and 0.1% formic acid 95% acetonitrile (B). The gradient profile was set as following: 2-22% B for 70 min, 22-29% B for 8 min, 29-95% B for 4 min, 95% B for 8 min. MS analysis was performed using an Orbitrap Velos Pro mass spectrometer (Thermo Scientific). The spray voltage was set at 2.2 kV. Orbitrap spectra (AGC 1x10$^6$) were collected from 400-1800 m/z at a resolution of 30,000 followed by data-dependent HCD MS/MS (at a resolution of 7,500, collision energy 35%, activation time 0.1 ms) of the ten most abundant ions using an isolation width of 2.0 Da. Charge state screening was enabled to reject unassigned and singly charged ions. A dynamic exclusion time of 40 sec was used to discriminate against previously selected ions.

### LC-MS/MS for Phosphoproteome Analysis

The multi-capillary metal-free LC system used an Eksigent nanoLC ultra-2D system and one Isco pump (Isco, Inc.) for delivering mobile phases. Samples were loaded onto the sample loop through a 6-ports peek valve (Valco Inc.) by a PAL auto-sampler and transferred to the micro SPE using mobile phase A by the Isco pump. The concentrated sample was transferred from the SPE column to the 50-μm i.d. LC column through a 10-port peek valve (Valco Inc.). A "back-flush" arrangement was used on the system. When the column labeled as SPE1 was in use, SPE2 was being re-equilibrated by mobile phase A to prepare the column for loading, concentrating, desalting and separating the next sample. A Pt wire (1 cm) was used to apply high voltage for electrospray ionization. The ratio of column flow rate (90 nL/min) and pump flow rate (900 nL/min) was controlled by the flow adjustor. Full automation was made possible by in-house software "LC-MSNet" that allows for parallel event coordination and faster MS duty cycle through use of two SPE columns and two analytical columns. All columns were manufactured in-house by slurry packing media into fused silica using a 1-cm sol-gel frit for media retention. SPE columns: 5-μm Jupiter C$_{18}$ (Phenomenex), 4-cm x 360 μm o.d. x 150 μm i.d. Reversed-phase analytical columns: 3-μm Jupiter C$_{18}$ (Phenomenex), 60-cm x 360 μm o.d. x 50 μm i.d. Mobile phase consisted of 0.1% formic acid in water (A) and 0.1% formic acid in acetonitrile (B) with a gradient profile as follows (min:%B); 0:5, 2:7, 120:25; 125:68, 129:80, 130:5. Sample injection (20 μL) and SPE (5 μL/min) occurred for 20-30 min before the SPE column was then placed in-line with the analytical column and the analytical gradient started. MS data acquisition lagged gradient start and end times by 50 min to account for column dead volume.

MS analysis was performed using either a LTQ Orbitrap Velos or a Q Exactive Hybrid Quadrupole-Orbitrap mass spectrometer (Thermo Scientific) outfitted with a custom electrospray ionization interface. Electrospray emitters were custom made using 360 μm o.d. x 20 μm i.d. chemically etched fused silica capillary. Analysis of the phosphoproteome sample fractions on the Velos mass spectrometer applied similar conditions as used in the global proteome sample fraction analysis, except that the spray voltage was 2.2 kV. On the Q Exactive instrument, the ion transfer tube temperature and spray voltage were 250 ºC and 2.2 kV, respectively. Mass spectra were collected from 300-1800 m/z at a resolution of 70,000 and AGC set at 3x10$^6$ followed by data-dependent HCD MS/MS (profile mode, at a resolution of 17,500, collision energy 28%, AGC 5x10$^4$) of the twelve most abundant ions using an isolation width of 2.5 m/z. Charge state screening was enabled to reject unassigned and 1+, 7+, 8+, and >8+ ions. A dynamic exclusion time of 30 sec was used.

### LC-MS/MS for Targeted Proteome Analysis

To verify the finding of histone H4 acetylation on K12 and K16, 106 samples were analyzed in a targeted method using Sequential Windowed data independent Acquisition of the Total High-resolution Mass Spectra (SWATH-MS). SWATH acquires a complete and permanent digital record for all the detectable spectra of a sample using data-independent acquisition (DIA) (Collins et al., 2013; Liu et al., 2014; Liu et al., 2013). The digital record as SWATH-MS maps, once generated, can be used for iterative analyses of

# Supplemental Information

candidate proteins discovered from the initial discovery phase or when additional candidates are identified with bioinformatics analyses.

SWATH-MS measurements were conducted using an AB Sciex 5600+ TripleTOF mass spectrometer interfaced with an Eksigent ekspert nanoLC 425 cHiPLC system. One microgram of peptides was loaded onto a 6 mm x 200 µm ChromXP C18-CL 3 µm, 120Å trap followed by separation on a 75 µm x 15 cm ChromXP C18-CL 3 µm, 120Å Nano cHiPLC column using a 120-min method (90-min gradient) at a flow rate of 300 nL/min. SWATH data were acquired using a variable window strategy wherein the sizes of the precursor ion selection windows were directly related to the m/z density for improved specificity. The average window width for precursor ion selection was 12 m/z with a range of 6 – 25 m/z. The collision energy was optimized for each window according to the calculation for a charge 2+ ion centered upon the window with a spread of 5 eV. The MS accumulation time was 250 ms, and the MS/MS accumulation time for fragment ions accumulated in high sensitivity mode was 50ms, resulting in a total duty cycle of approximately 3.5 s.

## Quantification of Global Proteomics Data

LC-MS/MS analysis of the iTRAQ-labeled, bRPLC fractionated samples generated a total of 1,922 global proteomics data files from both JHU and PNNL. All files were jointly processed at PNNL using the following pipeline. Thermo RAW files were processed with DeconMSn (v2.3.0) to determine accurate charge and m/z values of the precursor ions and create a concatenated DTA (CDTA) file. Next, DTARefinery (Petyuk et al., 2010) (v1.2, in conjunction with X!Tandem, v2007.07.01.2) processed the CDTA file to characterize and correct for any instrument calibration errors. The calibrated CDTA file was then processed with MS-GF+ (Kim and Pevzner, 2014; Kim et al., 2008) (v9881), matching against the RefSeq human protein sequence database, release version 37 (32,799 proteins), combined with 15 trypsin and keratin contaminants. The partially tryptic search used a +/-10 ppm parent ion tolerance, 0.5 m/z fragment ion tolerance, allowed for isotopic error in precursor ion selection, and searched a decoy database composed of the forward and reversed protein sequences. MSGF+ considered static carbamidomethylation (+57.0215 Da) on Cys residues and 4-plex iTRAQ modification (+144.1021 Da) on the peptide N-terminus and Lys residues, and dynamic oxidation (+15.9949 Da) on Met residues for searching the global proteome data. A comparison of tryptic and semi-tryptic filters produced no new peptide identifications, and a 3% reduction in the total number of peptides identified; quantification of protein abundance correlated at r = 0.97 between the tryptic and semi-tryptic filters.

All spectral identification files from the MSGF+ search engine were converted to IDPicker 3 index files (idpXML) and used for protein assembly using IDPicker 3 (Ma et al., 2009). Peptide identification stringency was set at a maximum false discovery rate (FDR) of 1% at peptide level; a minimum of 3 unique peptides was required to identify a given protein within the full data set, resulting in the identification of a total of 9,600 protein groups among the 174 samples (Table S2) with protein-level FDR <1%.

The intensities of all four iTRAQ reporter ions were extracted using MASIC software (Monroe et al., 2008). Next, PSMs passing the confidence threshold described above were linked to the extracted reporter ion intensities by scan number. The reporter ion intensities from different scans and different bRPLC fractions corresponding to the same protein were summed. Relative protein abundance was calculated as the ratio of sample abundance to reference abundance using the summed reporter ion intensities. The site-specific pooled reference sample was labeled with iTRAQ 117 reagent at PNNL and 114 at JHU, respectively, allowing comparison of relative protein abundances across the entire sample set at each site. The procedures for merging and correcting for batch-effects between PNNL and JHU sample sets are described below. The relative abundances were log2 transformed to obtain final, relative expression values.

The achieved dynamic range of the global proteome measurements was estimated to be >4 orders of magnitude, from low abundance proteins such as transcription factors to major structural proteins such as actins and tubulins, using a "proteomic ruler" method that normalizes the MS signal of individual proteins over histones (proportional to the amount of DNA in the sample) for estimation of the copy numbers of individual proteins per cell (Wisniewski et al., 2014), without the need for additional spike-in experiments.

## Data Processing Prior to Combining the JHU and PNNL Datasets

# Supplemental Information

Small differences in laboratory conditions and sample handling can result in systematic, sample-specific bias in the quantification of protein levels. In order to mitigate these effects while minimizing the risk of overcorrecting potentially important biological differences between individuals, we computed the median, log2 relative protein abundance over 4,476 proteins present in every sample and used re-centering to achieve a common median of 0.

The 32 overlapping samples analyzed by both JHU and PNNL allowed us to correct for laboratory-related differences in the log2 relative abundances at individual protein levels between the two sites prior to merging them into a single dataset for analysis. Specifically, potential batch effects between the two sites (PNNL and JHU), were corrected by shifting the PNNL data at individual protein levels so that the median abundances of each protein estimated over the 32 overlapping samples at PNNL and JHU were equalized (the batch effect corrected, unmerged data are provided in Table S2). For the 32 overlapping samples, the normalized PNNL and JHU measurements were averaged and used as their protein abundances in the merged dataset.

## Assessment of Analytical Precision of Proteomic Measurement

In addition to the 32 samples with overlapping measurements by both JHU and PNNL, a QC sample (a single ovarian cancer sample independent of the TCGA collection) was repeatedly analyzed 10 times spanning the entire duration and co-randomized as a part of the JHU sample proteomic analysis process. Data from the 32 pairs of duplicate measurements (4,682 proteins with no missing values) and the 10 replicate analyses of the QC samples (6,436 proteins with no missing values) provided the basis for the following assessment of analytical precision for the current proteomic dataset.

1) *Within-site coefficient of variation.* Coefficient of variation (CV), defined as the standard deviation divided by mean of replicated measurements, is a standardized measure of precision of an analytical assay(Everitt, 1998). The estimation of the within-site CV used the relative abundance data from the 10 replicate analyses of the QC sample instead of the log2 transformed data since the log2 ratio data tended to have means close to 0. The histogram of the estimated within-site CVs among the 10 replicates was plotted in Figure S1B. Among these proteins 85% had a CV < 15% (Min. 0.009, 1st Qu. 0.057, Median 0.082, Mean 0.099, 3rd Qu. 0.122, Max. 0.916).

2) *Approximated between-site coefficient of variation.* To assess the remaining between-site analytical variability after correction for batch effect, we used the 32 pairs of duplicate JHU and PNNL relative abundance data (again, without log2 transformation) to approximate the between-site CV for each protein by calculating the absolute value of the difference between two paired measurements divided by their mean, then averaged across the 32 samples. In Figure S1C, we plotted the histogram of the approximated between-site CVs of proteins with no missing values. The majority of the proteins (n = 3,426, 73%) had a CV < 20% (Min. 0.049, 1st Qu. 0.127, Median 0.161, Mean 0.176, 3rd Qu. 0.204, Max. 0.695).

3) *Correlation between samples analyzed at PNNL and JHU.* In order to understand the effect of between-site analytical variability on correlation estimation, we have in addition performed between-site correlation analyses using the 32 samples that had been analyzed by both PNNL and JHU, at the sample level across proteins and more relevantly at individual gene/protein level across the 32 paired measurements. In Figure 1 (upper panel), we show that the paired PNNL and JHU measurements of proteins had significantly higher correlations (median and mean: 0.69 and 0.67, respectively) than those between mRNA and proteins (median and mean: 0.38 and 0.45, respectively) (two-sided Wilcoxon rank-sum test, $p$ value < 2.2e-16). We also compared the concordance between protein abundances measured for the same samples at JHU and PNNL using co-clustering and principle component analysis (PCA) (Figure S1A). Using the software package MBatch (http://bioinformatics.mdanderson.org/main/TCGABatchEffectsV2:MBatch), no significant batch effect was observed between the two sites in the normalized/corrected data and further improvement from additional MBatch correction algorithms (Empirical Bayes, Median Polish and ANOVA) correction procedures was negligible.

## Final Protein Abundance Data

Further analysis and proteome-transcriptome associations used only proteins observed in all 169 HGSC samples and whose total variance (which included biological variance) among samples exceeded technical variance (3,586) in the merged data (Table S2). Proteins with missing data were excluded from the analysis to avoid problems associated with the imputation of missing values (Karpievitch et al., 2012; Webb-

# Supplemental Information

Robertson et al., 2013; Webb-Robertson et al., 2010). The data are available at the CPTAC data portal (https://cptac-data-portal.georgetown.edu/cptacPublic/).

## Quantification of Phosphopeptides

In addition to global proteomics, phosphoproteomics data were also acquired at PNNL for the 69 tumors with sufficient sample (see II A and II C). Phosphopeptide identification for these data were performed as described in **Quantification of Global Proteomics Data** (e.g., peptide level FDR <1%), with an additional dynamic phosphorylation (+79.9663 Da) on Ser, Thr or Tyr residues using the fully tryptic search. The phosphoproteome data was further processed by Ascore algorithm (Beausoleil et al., 2006) for phosphorylation site localization, and the top-scoring sequences were reported. For phosphoproteomic datasets, the iTRAQ quantitative data was not summarized by protein, but left at phosphopeptide level (Table S2).

All the peptides (phosphopeptides and global peptides) were labeled with iTRAQ reagent simultaneously. Separation into phospho- and non-phosphopeptides using IMAC was performed after the labeling. Thus, all the biases upstream of labeling are assumed to be identical between global and phosphoproteomics datasets. Therefore, to account for sample-specific biases in the phosphoproteome analysis we applied the correction factors derived from mean-centering the PNNL global proteomic dataset.

Phosphorylation site occupancy ratios (i.e., the relative extent of phosphorylation) were calculated by subtracted the log2 relative abundance of the parent protein (reflecting the overall abundance) from the phosphopeptide. The resulting phosphopeptide values more accurately reflect the changes in relative phosphorylation rates of the protein and are not confounded by changes in expression of the protein itself.

One phosphoproteome iTRAQ run of 3 samples (TCGA-24-1556, TCGA-24-2267 and TCGA-24-2260) was excluded from the final analysis because the overall number of missing values was significantly higher than that of other samples. The missing value count for all other runs fell within one standard deviation of the mean. Two samples that do not have somatic *TP53* mutation (TCGA-25-1316 and TCGA-61-2095) were also removed from the functional analyses.

Variation in warm ischemia time during the surgical sample acquisition process is a significant source of unwanted variation in the measurement of phosphopeptide abundance (Mertins et al., 2014). Although we believe that any error associated with warm ischemia should be independent of the clinical and pathological variable of interest, we conservatively filtered out all phosphopeptides mapping to any protein found to be affected by warm ischemia in a previous analysis (Mertins et al., 2014), removing 16,788 phosphopeptides mapping to 2,096 proteins (the remaining data used for downstream functional analyses are provided in Table S2). It is worth noting that the bulk phosphorylation differences observed between long and short surviving patients did not change due to this removal.

## Quantification of Acetylated Peptides

In order to identify specific sites of acetylation, we searched the global proteomics data as previously described with the addition of a dynamic acetylation to lysines (+42 Da). The PSMs were filtered by (1) FDR < 1%, (2) number of missed cleavage sites no greater than 2, (3) at least 2 PSMs for an identified peptide, and (4) at least 2 peptides for an identified protein. Each identified PSM was quantified using iTRAQ reporter ion intensities.

Additional validation from targeted proteomics experiments (SWATH) was used to orthogonally quantify the acetylated peptide abundances. To quantify the SWATH data, a synthetic peptide of interest (Histone H4; aa 9-17: GLGK(Ac)GGAK(Ac)R) was used to generate MS/MS spectrum for subsequent data processing. The MS/MS spectrum of the synthetic acetylated Histone H4 peptide served as the spectral library for the targeted data analysis of the SWATH maps to quantify K12K16 acetylation using Skyline (MacLean et al., 2010) (v. 2.6). The 5 transitions that were used for peak area integration were selected based on their rank and relative intensity in the spectrum from the synthetic peptide. Peak integrations were reviewed manually. Peak areas represent the sum of all 5 transitions. Quantitative data were exported from Skyline for statistical analysis.

## Identification of Novel Peptides

In an effort to identify novel peptides predicted from observed genomic alterations, we used a tiered searching strategy for all global proteomics data against four specialized databases; this approach was used only for identification of novel proteogenomic events. All identifications used in other analyses were performed as described above (**Quantification of Global Proteomics Data**). Spectra confidently identified

## Supplemental Information

in one search were excluded from subsequent searches. The search parameters were as follows: parent mass tolerance 20 ppm, semi-tryptic, static carbamidomethylation on Cys residues and iTRAQ modifications on the peptide N-terminus and Lys residues, and dynamic oxidation on Met residues. Each search was separately filtered for FDR control (see below).

The first search identified known (*i.e.*, not novel) peptides using the Ensembl protein sequence database (http://www.ensembl.org version GRCh37.70). This database is 59.5 MB in size. After searching all spectra against this database, we filtered the results to a 1% peptide-level false discovery rate (see below). Spectra for confidently identified peptides were removed from subsequent searches; spectra which did not have a confident identification were used in the next round of database search.

The second search was designed to identify variant peptides, corresponding to a nucleotide polymorphism identified by the TCGA sequencing data. The sequence polymorphisms which altered protein sequence were identified and used to create a database composed of candidate novel peptides using the software tool SpliceDB (Woo et al., 2014a; Woo et al., 2014b). Using 2.0 TB of RNA-seq files corresponding to the analyzed tumors and downloaded from TCGA consortium (https://cghub.ucsc.edu/software/downloads.html), a variant database of 852.8 MB was created. This variant peptide database was used in the second round of search using MSGF+ with the same parameters as described above. Because all spectra which correspond to peptides in the known Ensembl database were excluded from this search, confident identifications in this second round are exclusively variant peptides. We applied a 1% peptide-level FDR to filter the search results and report only confident variant peptides. Spectra which identified these variant peptides were removed from the next round of searching.

The third search was designed to identify peptides produced by novel splicing events. The SpliceDB tool used the 2 TB of RNA-Seq data from matched TCGA tumors to create a protein sequence database of splice junctions seen in the transcript data which are novel (*i.e.*, not part of Ensembl). This created a 641.7 MB protein sequence database. This splice-junction database was used in the third round of search using MSGF+ with the same parameters as described above. Because all confidently identified spectra from rounds one and two were excluded from this search, confident identifications in this third round are exclusively peptides which span these novel junctions and exon extensions. We applied a 1% peptide-level FDR to filter the search results and report only confident peptides. Spectra which identified these peptides were removed from the next round of searching.

The fourth and final search was against a 6-frame translation of the human genome, and was 3.1 GB in size. Because all confidently identified spectra from previous rounds were excluded from this search, confident identifications in this final round are exclusively peptides which were outside known exons or those predicted with RNA-seq. We applied a 1% peptide-level FDR as described below to filter the search results and report only confident peptides.

For each search, a matching decoy database of the same size was created for FDR control, by reversing the sequences. The FDR for a given score threshold ($\tau$) is calculated as

$$FDR(\tau) = \frac{Number\ of\ distinct\ decoy\ peptides\ with\ score \geq \tau}{Number\ of\ distinct\ target\ peptides\ with\ score \geq \tau}$$

As described above in the search overview, we applied FDR filters after every stage of searching. This strategy is an intentionally conservative approach to avoid false positives. After filtering, there are 3,863 novel peptides discovered at 1% FDR.

Peptides mapping to more than three genomic locations were removed from consideration at the event level. The remaining peptides were grouped based on genomic location, where multiple peptides supported a single novel event. After the grouping, we calculate the probability that at least one peptide in a group is true (or 1-probability all peptides are false).

$$p(True\ Event) = \prod_{i}(1 - \frac{FDR(pep_i)}{Number\ of\ locations})$$

Events with probability > 0.5 were classified into the following types by Enosi (Woo et al., 2014b): alternative-splice, novel-splice, fusion-gene, novel-gene, deletion, insertion, mutation, novel exon, exon boundary, translated UTR, non CDS gene, pseudo gene, frame shift, reverse strand, and IG gene expression. Interestingly, we also observed a large set of immunoglobulin gene peptides, suggesting a response from infiltrating B-lymphocytes.

To validate the identification of some of the novel proteogenomic peptides, we randomly selected 20 candidate novel peptides from the full list of proteogenomic peptide identifications (Table S2), and ordered the corresponding synthetic peptides from New England Peptide for LC-MS/MS analysis.

# Supplemental Information

Comparing histograms of the database searching scores (best MSGF score "SpecEValue" for each peptide) for: 1) all global "normal" peptide identifications, 2) all proteogenomic variants, and 3) the 20 tested proteogenomic variants shows that the proteogenomic variants selected for validation represent the full range of spectral quality observed in these experiments. The peptide mixture containing the 20 synthetic peptides was analyzed under the same LC-MS/MS conditions as described in **Quantification of Global Proteomics Data**. All 20 proteogenomic variants were validated using the synthetic peptides.

## Comparison of mRNA and protein subtypes

To model the protein abundance data as a mixture of subtypes we employed a model-based clustering approach (Fraley and Raftery, 2007). Where more than one protein was mapped to a gene, we selected a representative (minimum RefSeq ID) protein thus reducing the number of proteins from 3,586 down to 3,326. The 50% most variable of these proteins (1,663) were used for clustering, and individual log2 ratios were scaled by the standard deviations of the corresponding proteins to produce z-scores.

Three different approaches were considered to define the optimal number of protein clusters, including Bayesian information criteria (Fraley and Raftery, 2002), statistical resampling (consensus clustering) (Fraley and Raftery, 2007), and VISDA-based sub-phenotype clustering (Wang et al., 2007). All three approaches provided similar results, with a stable optimization on five clusters, and the grouping shown in Figure S2A was derived using the *mclust* R package (Fraley and Raftery, 2002; Fraley et al., 2012), based on Bayesian information criteria with default settings to model clusters.

Each proteomic cluster was compared to each CLOVAR subtypes using Fisher's exact test applied to separate 2x2 tables describing subtype membership (Figure S2E). Proteins over- and under-abundant in the newly identified proteomic subtypes were discovered using linear model analysis, using the *limma* package from Bioconductor (Ritchie et al., 2015). Preliminary statistical analysis of the protein subtypes for possible bias associated with specific Tumor Sample Sites or tumor stage indicated a random distribution of TSS and tumor stage across the protein subtypes.

To infer genes or gene networks that drive subtyping into 5 clusters we performed WGCNA analysis (Langfelder and Horvath, 2008) followed by correlation with subtype as a trait. Use of module discovery with default WGCNA settings updated to allow additional mergers resulted in seven co-abundant protein clusters. To relate the protein clusters to biological knowledge we performed ontology term enrichment analysis using Bioconductor "clusterProfiler" package (Yu et al., 2012) (Figure 2). The identified modules were named according to enriched ontology terms from KEGG and Reactome pathway databases (Figures S2C and S2D). Given the module-trait correlation matrix, the subject subtypes were most statistically significantly associated with a combination of modules rather than by a single module (Figures 2 and S2B).

Consensus clustering was used to assess the clustering stability (Monti et al., 2003). Specifically, 90% of the original sample pool were randomly subsampled without replacement and partitioned into 5 clusters using model-based clustering (Fraley and Raftery, 2007). This was repeated 200 times, and pairwise consensus values calculated, describing how frequently two samples are assigned to the same cluster. Stable partitions of the data will show much higher average consensus values for within-cluster sample pairs than for between-cluster pairs (Figure S2A).

## Calculation and Functional Enrichment of mRNA-Protein Correlation

The mRNA expression values for the 169 HGSC tumors analyzed in this study was obtained from FIREHOSE (https://confluence.broadinstitute.org/display/GDAC/Home) as microarray analysis. Expression values were transformed to log2 fold change versus mean expression for each gene. Analysis was restricted to 3,202 proteins with complete data for which matched mRNA data was available. Spearman correlation was calculated for each mRNA-protein pair across all 169 tumors. Spearman correlation is robust to the potential differences in distribution of the values from mRNA profiling and protein abundance, however, Pearson correlation yielded very similar results. Separately, Spearman correlation between mRNA and protein was calculated for those samples run by PNNL (82 tumors) and those run by JHU (119 tumors) giving correlations of 0.28 and 0.35, respectively. Using semi-partial correlation analysis to include site as a covariate to adjust for possible batch effect between PNNL and JHU showed negligible difference in correlation estimation. Such adjustment hence was not necessary.

To determine if functional groups (pathways or complexes) were non-randomly distributed in terms of mRNA-protein correlation, we used the non-parametric Kolmogorov–Smirnov test, and considered a $p$ value of 0.05 or less significant after multiple hypothesis correction with the Benjamini-

# Supplemental Information

Hochberg method. Functional classes were obtained from the MSIGDB (http://www.broadinstitute.org/gsea/msigdb/index.jsp) and highlighted functions in Figure 1 were chosen from the most significant non-redundant functions that were biologically informative. The individual proteins associated with pathways highlighted in the text as being significantly differentially present in more or less correlated protein-mRNA pairs are presented as Table S3.

## Correlation of CNA with mRNA Expression and Protein Abundance

Gene expression data for the samples in the cohort was obtained from FIREHOSE (https://confluence.broadinstitute.org/display/GDAC/Home) from microarray analysis and copy number data from the same source as the GISTIC 2.0 processed CNA data that has been summarized by gene. The Biomart portal (http://www.biomart.org/) was used to convert the protein Refseq ID to the respective gene names. Analysis was restricted to 3,202 proteins for which matched mRNA data was available. CNA data for one case (TCGA-13-2066) among the 169 cases was missing, so the analysis included 168 cases across the mRNA, protein and CNA data. Spearman correlation coefficients and corresponding $p$ values for 3,202 protein/transcript X 29,393 CNA loci were calculated using the R function cor.test, and the $p$ values were corrected using the Benjamini-Hochberg procedure to estimate FDR.

## Functional Enrichment Analysis of *Trans*-Regulated Proteins

Gene set analysis was performed to infer functions for groups of proteins that are significantly correlated with a given CNA locus. Gene sets were taken from MSIGDB v4.0 (http://www.broadinstitute.org/gsea/msigdb/index.jsp), and evaluated by hypergeometric test, with Benjamini-Hochberg correction of $p$ values. The $p$ values after correction were then used to filter out enriched KEGG, NCI pathway information database (PID), and Reactome pathways. We provide both the complete set of trans-affected proteins associated with each CNA and protein membership of the associated pathways in Table S3.

## Functional Analysis of CNA Loci with Strong *Trans* Effects

The most perturbed region on chromosome 2 at 55MB (Figure 3) affects 363 proteins, which are significantly enriched in the KEGG pathways for regulation of the actin cytoskeleton (hsa04810) and leukocyte transendothelial migration (hsa04670). Proteins in these pathways are likely to convey a survival advantage on malignant cancer cells, facilitating the ability to invade adjacent tissues and the vasculature. Interestingly, this enrichment for actin cytoskeleton and cell migration pathways was not observed in mRNA and CNA correlations. Enrichment of cell invasion and migration proteins was also observed for CNA loci on seven other chromosomes (2, 7, 11, 12, 15, 17, 20 and 22, Figure S3B).

Another prominent functional enrichment was for proteins related to immune function. Chromosome 20 has a distinct affected region ranging from 30-37 MB (Figure S3B) containing 146 loci. The most perturbed locus affects 473 proteins, which are significantly enriched in the KEGG pathways associated with immune function (hsa04062, hsa04612, hsa04672) and related diseases (hsa05310, hsa05320, hsa05330, hsa05332, hsa05140). This recurring theme is also conserved across five other chromosomes (chromosomes 5, 6, 12, 17 and 22). This observation underscores the importance of immune regulation as a target of selective pressure.

## Chromosomal Instability Index and Associated Proteins

A CIN index (TCGA Research Network, 2012) was calculated for each sample as the mean absolute values of copy number measurements at the 29,393 selected loci. Bootstrap resampling was used to select proteins correlated with *CIN index* at high confidence. Specifically, within each of the 100 bootstrap iterations, Spearman correlation tests were performed between sample CIN indices and proteins across the current bootstrap sample of patients. The selected proteins were those exhibiting a Benjamini-Hochberg adjusted $p$ value < 1e-4 in at least 50% of the bootstrap iterations. Within the same bootstrap procedure, the abundances of proteins were also correlated directly to CNA data at 19,898 selected loci across the current bootstrap sample of patients. For each protein, the number of loci (outside of the protein's own coding chromosome) that had a significant Spearman correlation (Benjamini-Hochberg adjusted $p$ value < 0.01) was recorded. The 19,898 loci were selected due to their relatively large CNA, defined as having an absolute value of the sum of CNAs greater than 10, a cutoff determined to capture the majority of loci with large CNAs.

# Supplemental Information

Results were summarized graphically in a heatmap showing relative abundances for the selected proteins (Figure S5A), and a dot plot showing the number of significantly affected loci for each of the selected proteins, organized by genomic position (Figure 5). An interactome network of the significant proteins (Figure S5B) was extracted from multiple curated databases (BIND, BioGRID, DIP, HPRD, IntAct, MINT, MIPS, PDZBase, Reactome) built from Yu et al. (Yu et al., 2012) using Cytoscape (Shannon et al., 2003). Functional annotation of the proteins in the interactome network was obtained through gene set enrichment analysis using MSigDB database (Subramanian et al., 2005). The phosphorylation levels of selected phosphopeptides were plotted as well (Figure S5C).

## Analysis of Homologous Repair Deficiency (HRD)

Acetylation analysis of the 119 TCGA tumors initially analyzed at JHU identified 3,527 unique acetylated peptides from 1,987 protein groups. Of these, 26 acetylated peptides were observed in all 119 samples (the 100% filter value); a total of 399 acetylated peptides were identified and quantified in at least 50% of the JHU samples. It is worth noting that the 100% filter applied for the global proteomics analysis was at the protein group level (but a protein group can be identified by different peptides), and the 50% filter for the acetylated peptides was applied at the acetylated peptide level, for specific analysis of site-specific acetylation. Setting the filter to 50% provided a large dataset for statistical analysis. Moreover, the peptides identified in this analysis were used as a candidate list for further validation with targeted proteomics methods.

Acetylation levels were compared between HRD and non-HRD cases by t-test with a $p$ value $< 0.05$ considered statistically significant, FDR $< 0.05\%$ was estimated by bootstrap/permutation test. A histone H4 acetylation score for each sample was calculated as the average value of H4 peptide acetylation at K12 and K16. The median histone H4 acetylation score was used to partition samples into high- and low-acetylation groups (n=61 per group) and compared by Kaplan-Meier analysis and log rank test.

BRCA1/BRCA2 related protein signatures were curated from the literature (TCGA Research Network, 2011; Konstantinopoulos et al., 2010) and the cBio portal (Cerami et al., 2012), with 171 proteins selected for the subsequent analysis (Table S7). We applied differential dependency network (DDN) analysis (Zhang et al., 2009; Zhang et al., 2011; Tian et al., 2014) to these proteins, comparing 52 non-HRD cases with 67 HRD cases to identify HRD-associated changes in protein interaction. The DDN analysis and network visualization were performed using CytoDDN, a Cytoscape plugin for the DDN analysis. Random permutations ($k$=5000) were used to assess the significance of the network rewiring and the $p$ value cutoff was set to 0.05. Isolated node pairs that were not connected to the rest of the network with a probability better than 0.05 were trimmed from the network and its representation in Figure 5. Those nodes and connections remaining at a $p$ value cut-off of 0.01 are also shown on Figure 5 as bold lines.

## Analysis of *Trans*-Affected Proteins for Association with Survival

We used *trans*-affected protein data to build predictive models of overall survival. We calculated correlations between CNAs and proteins across the PNNL and JHU datasets separately, then considered proteins that were predicted to be trans-affected by a CNA in each dataset as input into modeling (see below). We filtered the CNAs to those with more than 300 shared trans-affected proteins (Benjamini-Hochberg adjusted $p$ value $< 0.05$) leaving 469 sets of trans-affected proteins, which largely fell into the most influential regions (on chromosomes 2, 7, 20, and 22; see Figure 3). In each case, models were trained on the proteomics data from PNNL (82 tumors) and tested in the data from JHU (87 tumors, excluding 32 overlapping cases).

Because each CNA affects many proteins, we used a regression approach that identifies parsimonious Cox proportional hazards models with maximal predictive ability from the list of significantly correlated proteins for that CNA. Specifically, we employed the regularization path following (RPF) LASSO model previously described (Duverle et al., 2013). It infers a minimal set of components with associated weights whose expression explains the dependent variable (*i.e.*, survival) with minimal error. The RPF method incorporates an itemset mining approach to determine possible interactions between variables that are incorporated into the model. The resulting model provides a predicted, relative hazard for each sample, describing the risk of death under a proportional hazards model.

Results were visualized in Kaplan-Meier plots comparing high risk cases having predicted hazards (model scores) above the median for the cohort (50% of cases), to low risk cases with predicted hazards below the median (Figure S4A). Details of the models are provided in Table S6.

# Supplemental Information

A consensus prediction was produced by each of the four models voting on whether each patient would be in a high or low risk group, where ties (two models predicting high and two predicting low) were considered to be in the low group. The performance of this signature at predicting overall survival is shown in a Kaplan-Meier plot (Figure 4) and the performance on predicting progression-free survival is also shown (Figure S4B).

We identified 'minimal' models (Table S6) from these by considering the model for each of the four chosen CNAs that had a minimal number of components, but was still significantly predictive of survival for the training and validation sets (as opposed to those models that were most predictive for the 'full' models described above). These models can serve as a resource for further investigation of protein panels that can be used to predict patient survival.

A previously described protein signature for ovarian cancer survival, Provar (Yang et al., 2013), was applied to our MS proteomics. The five unmodified proteins in Provar (AR, BID, EEF2, HSP70, STAT5A) were all identified in our proteomics data. Of the four phosphosites in the signature (EGFR Y1173, MEK1 S217+S221, PRKCA S657, and TAZ S89) we quantified EGFR Y1173, though we did quantify PRKCA S226s, this did not match the RPPA-based signature. For pMEK1 S217+S221, pPRKCA S657, and pTAZ S89 we used the RPPA quantitation for the same tumor (Available from the TCGA data portal) and paired it with our MS-derived abundance values for the remainder of the signature. Our patient cohort was more limited than that used to derive the Provar signature, and we found that the Provar signature was only weakly predictive of survival ($p$ value 0.11 as opposed to $p$ value < 0.001 reported in the original paper; Figure S4C). Since Provar was trained on RPPA data from a large number of patients it is possible that the smaller number of patients and MS abundance data are incompatible with the signature. Combining the Provar in our consensus signature approach described above did not improve the results.

## Pathway Analysis of Short and Long Surviving Patients

We implemented a simple pathway analysis to identify functional links between peptide phosphorylation and overall survival. Within the set of 69 tumors with phosphopeptide data, 17 short survivors selected to have overall survival time less than three years (1,118 days) were compared to 19 long survivors who survived more than 5 years (1,850 days). See Table S1.

Phosphopeptides were mapped to signaling pathways in the NCI PID (http://pid.nci.nih.gov) using the gene names. For each signaling pathway in PID, relative abundances for all phosphopeptides mapping to any pathway component were identified and separated into short and long survivor groups. The difference in distributions between the set of pathway-specific peptides associated with short survivors and the set of pathway-specific peptides associated with long survivors was then assessed using a two-tailed t test. Results with a Benjamini-Hochberg adjusted $p$ value <0.05 were considered statistically significant. We also performed similar enrichment analyses using protein abundance, mRNA abundance, and CNA. Figure 7A compares pathway enrichment results between short and long surviving patient samples and Figure 7B illustrates the components of the PDGFR-beta pathway identified through mRNA and MS-based proteomic and phosphoproteomic measurements on the same samples. Components (proteins or phosphoproteins) for those pathways highlighted in Figure 7A are presented as Table S3.

To assess the robustness of the pathway results, we repeated the analysis 1,000 times, each time randomly subsampling 50% of the samples without replacement, and calculating the proportions of runs in which the pathway had a Benjamini-Hochberg adjusted $p$ value < 0.05. The analysis was repeated using different thresholds for short and long survival to ensure that the result was robust to that parameter.

# Supplemental Information

## SUPPLEMENTAL REFERENCES

Bantscheff, M., Boesche, M., Eberhard, D., Matthieson, T., Sweetman, G., and Kuster, B. (2008). Robust and sensitive iTRAQ quantification on an LTQ Orbitrap mass spectrometer. Mol. Cell. Proteomics *7*, 1702-1713.

Beausoleil, S.A., Villen, J., Gerber, S.A., Rush, J., and Gygi, S.P. (2006). A probability-based approach for high-throughput protein phosphorylation analysis and site localization. Nat. Biotechnol. *24*, 1285-1292.

Cerami, E., Gao, J., Dogrusoz, U., Gross, B.E., Sumer, S.O., Aksoy, B.A., Jacobsen, A., Byrne, C.J., Heuer, M.L., Larsson, E.*, et al.* (2012). The cBio cancer genomics portal: an open platform for exploring multidimensional cancer genomics data. Cancer Discov. *2*, 401-404.

Duverle, D.A., Takeuchi, I., Murakami-Tonami, Y., Kadomatsu, K., and Tsuda, K. (2013). Discovering combinatorial interactions in survival data. Bioinformatics *29*, 3053-3059.

Eisenhauer, E.L., Abu-Rustum, N.R., Sonoda, Y., Aghajanian, C., Barakat, R.R., and Chi, D.S. (2008). The effect of maximal surgical cytoreduction on sensitivity to platinum-taxane chemotherapy and subsequent survival in patients with advanced ovarian cancer. Gynecol. Oncol. *108*, 276-281.

Everitt, B.S. (1998). The Cambridge Dictionary of Statistics (Cambridge: Cambridge University Press).

Fraley, C., and Raftery, A.E. (2002). Model-based clustering, discriminant analysis, and density estimation. J. Am. Stat. Assoc. *97*, 611-631.

Fraley, C., and Raftery, A.E. (2007). Model-based methods of classification: using the mclust software in chemometrics. J. Stat. Softw. *18*, 1-13.

Fraley, C., Raftery, A.E., Murphy, T.B., and Scrucca, L. (2012). mclust Version 4 for R: normal mixture modeling for model-based clustering, classification (and Density Estimation Technical Report No).

Karpievitch, Y.V., Dabney, A.R., and Smith, R.D. (2012). Normalization and missing value imputation for label-free LC-MS analysis. BMC Bioinformatics *13 Suppl 16*, S5.

Kim, S., Gupta, N., and Pevzner, P.A. (2008). Spectral probabilities and generating functions of tandem mass spectra: a strike against decoy databases. J. Proteome Res. *7*, 3354-3363.

Kim, S., and Pevzner, P.A. (2014). MS-GF+ makes progress towards a universal database search tool for proteomics. Nat. Commun. *5*, 5277.

Konstantinopoulos, P.A., Spentzos, D., Karlan, B.Y., Taniguchi, T., Fountzilas, E., Francoeur, N., Levine, D.A., and Cannistra, S.A. (2010). Gene expression profile of BRCAness that correlates with responsiveness to chemotherapy and with outcome in patients with epithelial ovarian cancer. J. Clin. Oncol. *28*, 3555-3561.

Liu, Y., Chen, J., Sethi, A., Li, Q.K., Chen, L., Collins, B., Gillet, L.C., Wollscheid, B., Zhang, H., and Aebersold, R. (2014). Glycoproteomic analysis of prostate cancer tissues by SWATH mass spectrometry discovers N-acylethanolamine acid amidase and protein tyrosine kinase 7 as signatures for tumor aggressiveness. Mol.Cell.Proteomics *13*, 1753-1768.

Liu, Y., Huttenhain, R., Surinova, S., Gillet, L.C., Mouritsen, J., Brunner, R., Navarro, P., and Aebersold, R. (2013). Quantitative measurements of N-linked glycoproteins in human plasma by SWATH-MS. Proteomics *13*, 1247-1256.

Ma, Z.Q., Dasari, S., Chambers, M.C., Litton, M.D., Sobecki, S.M., Zimmerman, L.J., Halvey, P.J., Schilling, B., Drake, P.M., Gibson, B.W.*, et al.* (2009). IDPicker 2.0: Improved protein assembly with high discrimination peptide identification filtering. J. Proteome Res. *8*, 3872-3881.

# Supplemental Information

MacLean, B., Tomazela, D.M., Shulman, N., Chambers, M., Finney, G.L., Frewen, B., Kern, R., Tabb, D.L., Liebler, D.C., and MacCoss, M.J. (2010). Skyline: an open source document editor for creating and analyzing targeted proteomics experiments. Bioinformatics *26*, 966-968.

Mertins, P., Udeshi, N.D., Clauser, K.R., Mani, D.R., Patel, J., Ong, S.E., Jaffe, J.D., and Carr, S.A. (2012). iTRAQ labeling is superior to mTRAQ for quantitative global proteomics and phosphoproteomics. Mol. Cell. Proteomics *11*, M111.014423.

Monroe, M.E., Shaw, J.L., Daly, D.S., Adkins, J.N., and Smith, R.D. (2008). MASIC: a software program for fast quantitation and flexible visualization of chromatographic profiles from detected LC-MS(/MS) features. Comput. Biol. Chem. *32*, 215-217.

Monti, S., Tamayo, P., Mesirov, J., and Golub, T. (2003). Consensus clustering: A resampling-based method for class discovery and visualization of gene expression microarray data. Mach. Lear. *52*, 91-118.

Ow, S.Y., Salim, M., Noirel, J., Evans, C., and Wright, P.C. (2011). Minimising iTRAQ ratio compression through understanding LC-MS elution dependence and high-resolution HILIC fractionation. Proteomics *11*, 2341-2346.

Pennington, K.P., Walsh, T., Harrell, M.I., Lee, M.K., Pennil, C.C., Rendi, M.H., Thornton, A., Norquist, B.M., Casadei, S., Nord, A.S.*, et al.* (2014). Germline and somatic mutations in homologous recombination genes predict platinum response and survival in ovarian, fallopian tube, and peritoneal carcinomas. Clin. Cancer Res. *20*, 764-775.

Petyuk, V.A., Mayampurath, A.M., Monroe, M.E., Polpitiya, A.D., Purvine, S.O., Anderson, G.A., Camp, D.G., 2nd, and Smith, R.D. (2010). DtaRefinery, a software tool for elimination of systematic errors from parent ion mass measurements in tandem mass spectra data sets. Mol.Cell.Proteomics *9*, 486-496.

R Development Core Team. (2008). R: A Language and Environment for Statistical Computing. R Foundation for Statistical Computing. Vienna, Austria (ISBN 3-900051-07-0, URL http://www.Rproject.org).

Ritchie, M.E., Phipson, B., Wu, D., Hu, Y., Law, C.W., Shi, W., and Smyth, G.K. (2015). limma powers differential expression analyses for RNA-sequencing and microarray studies. Nucleic Acids Res. gkv007.

Shannon, P., Markiel, A., Ozier, O., Baliga, N.S., Wang, J.T., Ramage, D., Amin, N., Schwikowski, B., and Ideker, T. (2003). Cytoscape: a software environment for integrated models of biomolecular interaction networks. Genome Res. *13*, 2498-2504.

Subramanian, A., Tamayo, P., Mootha, V.K., Mukherjee, S., Ebert, B.L., Gillette, M.A., Paulovich, A., Pomeroy, S.L., Golub, T.R., Lander, E.S.*, et al.* (2005). Gene set enrichment analysis: a knowledge-based approach for interpreting genome-wide expression profiles. Proc. Natl. Acad. Sci. U.S.A. *102*, 15545-15550.

Tabb, D.L., Wang, X., Carr, S.A., Clauser, K.R., Mertins, P., Chambers, M.C., Holman, J.D., Wang, J., Zhang, B., Zimmerman, L.J., *et al*. (2016). Reproducibility of Differential Proteomic Technologies in CPTAC Fractionated Xenografts. J. Proteome Res. *15*, 691-706.

Tian, Y., Zhang, B., Hoffman, E.P., Clarke, R., Zhang, Z., Shih, I.-M., Xuan, J., Herrington, D.M., and Wang, Y. (2014). Knowledge-fused differential dependency network models for detecting significant rewiring in biological networks. BMC Systems Biol. *8*, 87.

Vang, R., Levine, D.A., Soslow, R.A., Zaloudek, C., Shih, Ie.M., and Kurman, R.J. (2016). Molecular Alterations of TP53 are a Defining Feature of Ovarian High-Grade Serous Carcinoma: A Rereview of Cases Lacking TP53 Mutations in The Cancer Genome Atlas Ovarian Study. Int. J. Gynecol. Pathol. *35*, 48-55.

# Supplemental Information

Wang, J., Li, H., Zhu, Y., Yousef, M., Nebozhyn, M., Showe, M., Showe, L., Xuan, J., Clarke, R., and Wang, Y. (2007). VISDA: an open-source caBIG analytical tool for data clustering and beyond. Bioinformatics *23*, 2024-2027.

Wang, Y., Yang, F., Gritsenko, M.A., Wang, Y., Clauss, T., Liu, T., Shen, Y., Monroe, M.E., Lopez-Ferrer, D., Reno, T.*, et al.* (2011). Reversed-phase chromatography with multiple fraction concatenation strategy for proteome profiling of human MCF10A cells. Proteomics *11*, 2019-2026.

Wax, M., and Kailath, T. (1985). Detection of Signals by Information Theoretic Criteria. IEEE T. Acoust. Speech *33*, 387-392.

Webb-Robertson, B.J., Matzke, M.M., Metz, T.O., McDermott, J.E., Walker, H., Rodland, K.D., Pounds, J.G., and Waters, K.M. (2013). Sequential projection pursuit principal component analysis--dealing with missing data associated with new -omics technologies. BioTechniques *54*, 165-168.

Webb-Robertson, B.J., McCue, L.A., Waters, K.M., Matzke, M.M., Jacobs, J.M., Metz, T.O., Varnum, S.M., and Pounds, J.G. (2010). Combined statistical analyses of peptide intensities and peptide occurrences improves identification of significant peptides from MS-based proteomics data. J. Proteome Res. *9*, 5748-5756.

Wiśniewski, J.R., Hein, M.Y., Cox, J., and Mann, M. (2014). A "proteomic ruler" for protein copy number and concentration estimation without spike-in standards. Mol. Cell. Proteomics *13*, 3497-3506.

Woo, S., Cha, S.W., Merrihew, G., He, Y., Castellana, N., Guest, C., MacCoss, M., and Bafna, V. (2014b). Proteogenomic database construction driven from large scale RNA-seq data. J. Proteome Res. *13*, 21-28.

Yoshihara, K., Shahmoradgoli, M., Martínez, E., Vegesna, R., Kim, H., Torres-Garcia, W., Treviño, V., Shen, H., Laird, P.W., Levine, D.A., et al.,(2013). Inferring tumour purity and stromal and immune cell admixture from expression data. Nat. Commun. *4*, 2612.

Yu, G., Wang, L.G., Han, Y., and He, Q.Y. (2012). clusterProfiler: an R package for comparing biological themes among gene clusters. OMICS. *16*, 284-287.

Yu, X., Wallqvist, A., and Reifman, J. (2012). Inferring high-confidence human protein-protein interactions. BMC Bioinformatics *13*, 79.

Zhang, B., Tian, Y., Jin, L., Li, H., Shih Ie, M., Madhavan, S., Clarke, R., Hoffman, E.P., Xuan, J., Hilakivi-Clarke, L.*, et al.* (2011). DDN: a caBIG(R) analytical tool for differential network analysis. Bioinformatics *27*, 1036-1038.

## EXTENDED AUTHOR LIST

### National Cancer Institute Clinical Proteomics Tumor Analysis Consortium (NCI CPTAC)

Steven A. Carr[1], Michael A. Gillette[1], Karl R. Klauser[1], Eric Kuhn[1], D. R. Mani[1], Philipp Mertins[1], Karen A. Ketchum[2], Ratna Thangudu[2], Shuang Cai[2], Mauricio Oberti[2], Amanda G. Paulovich[3], Jeffrey R. Whiteaker[3], Nathan J. Edwards[4], Peter B. McGarvey[4], Subha Madhavan[5], Pei Wang[6], Daniel Chan[7], Akhilesh Pandey[7], Ie-Ming Shih[7], Hui Zhang[7], Zhen Zhang[7], Heng Zhu[8], Leslie Cope[9], Gordon A. Whiteley[10], Steven J. Skates[11], Forest M. White[12], Douglas A. Levine[13], Emily S. Boja[14], Christopher R. Kinsinger[14], Tara Hiltke[14], Mehdi Mesri[14], Robert C. Rivers[14], Henry Rodriguez[14], Kenna M. Shaw[14], Stephen E. Stein[15], David Fenyo[16], Tao Liu[17], Jason E. McDermott[17], Samuel H. Payne[17], Karin D. Rodland[17], Richard D. Smith[17], Paul Rudnick[18], Michael Snyder[19], Yingming Zhao[20], Xian Chen[21], David F. Ransohoff[21], Andrew N. Hoofnagle[22], Daniel C. Liebler[23], Melinda E. Sanders[23], Zhiao Shi[23], Robbert J. C. Slebos[23],

# Supplemental Information

David L. Tabb[23], Bing Zhang[23], Lisa J. Zimmerman[23], Yue Wang[24], Sherri R. Davies[25], Li Ding[25], Matthew J. C. Ellis[26] & R. Reid Townsend[25]

[1]The Eli and Edythe L. Broad Institute of Massachusetts Institute of Technology and Harvard University Cambridge, Massachusetts 02142, USA. [2]Enterprise Science and Computing, Inc., 155 Gibbs St, Suite 420, Rockville, Maryland 20850, USA. [3]Clinical Research Division, Fred Hutchinson Cancer Research Center, 1100 Eastlake Avenue East, Seattle, Washington 98109, USA. [4]Department of Biochemistry and Molecular & Cellular Biology, Georgetown University Medical Center, 3900 Reservoir Rd NW, Washington, DC 20057, USA. [5]Innovation Center for Biomedical Informatics, Georgetown University Medical Center, 2115 Wisconsin Ave NW, Suite 110, Washington, DC 20057, USA. [6]Icahn Institute and Department of Genetics and Genomic Sciences, Icahn School of Medicine at Mount Sinai, Hess CSM Building, Room S8-102, 1470 Madison Avenue, New York, New York 10029, USA. [7]Department of Pathology, The Johns Hopkins University, 600 North Wolfe Street, Baltimore, Maryland 21287, USA. [8]Department of Pharmacology and Molecular Science, the Johns Hopkins University, 733 N. Broadway, Baltimore, Maryland 21287, USA. [9]Department of Oncology, the Johns Hopkins University, 733 N. Broadway, Baltimore, Maryland 21287, USA. [10]Antibody Characterization Laboratory, Advanced Technology Program, Leidos, Inc., 1050 Boyles Street, Frederick, Maryland 21701, USA. [11]Biostatistics Center, Massachusetts General Hospital Cancer Center, 55 Fruit Street, Boston, Massachusetts 02114, USA. [12]Department of Biological Engineering, Massachusetts Institute of Technology, 77 Massachusetts Avenue, Cambridge, Massachusetts 02139, USA.[13]Gynecology Service/Department of Surgery, Memorial Sloan Kettering Cancer Center, 1275 York Avenue, New York, New York 10065, USA. [14]Office of Cancer Clinical Proteomics Research, National Cancer Institute, 31 Center Drive, MS 2580 Bethesda, Maryland 20892, USA. [15]Biomolecular Measurement Division, Material Measurement Laboratory, National Institute of Standards and Technology, 100 Bureau Drive, M/S 8300, Gaithersburg, Maryland 20899, USA. [16]Department of Biochemistry and Molecular Pharmacology, Smilow Research Building, Room 201, 522 First Avenue, New York University Langone Medical Center, New York, New York 10016, USA. [17]Biological Sciences Division, Pacific Northwest National Laboratory, 902 Battelle Boulevard, Richland, Washington 99352, USA. [18]Spectragen-Informatics, Rockville, Maryland 20850, USA. [19]Department of Genetics, Stanford University, Stanford, California 94305, USA. [20]The Ben May Department for Cancer Research, University of Chicago, 929 East 57th Street, W421 Chicago, Illinois 60637, USA. [21]University of North Carolina at Chapel Hill, 130 Mason Farm Road, Chapel Hill, North Carolina 27599, USA. [22]Department of Lab Medicine, University of Washington, Campus Box 357110, Seattle, Washington 98195, USA. [23]Vanderbilt University School of Medicine, 1161 21st Avenue South, Nashville, Tennessee 37232, USA. [24]Bradley Department of Electrical and Computer Engineering, Virginia Tech, 900 N. Glebe Road, Arlington, Virginia 22203, USA. [25]Department of Medicine, Washington University in St. Louis, 660 S. Euclid Avenue, St. Louis, Missouri 063110, USA.[26]Departments of Medicine and Cellular and Molecular Biology, Baylor College of Medicine, One Baylor Plaza, Houston, Texas 77030, USA