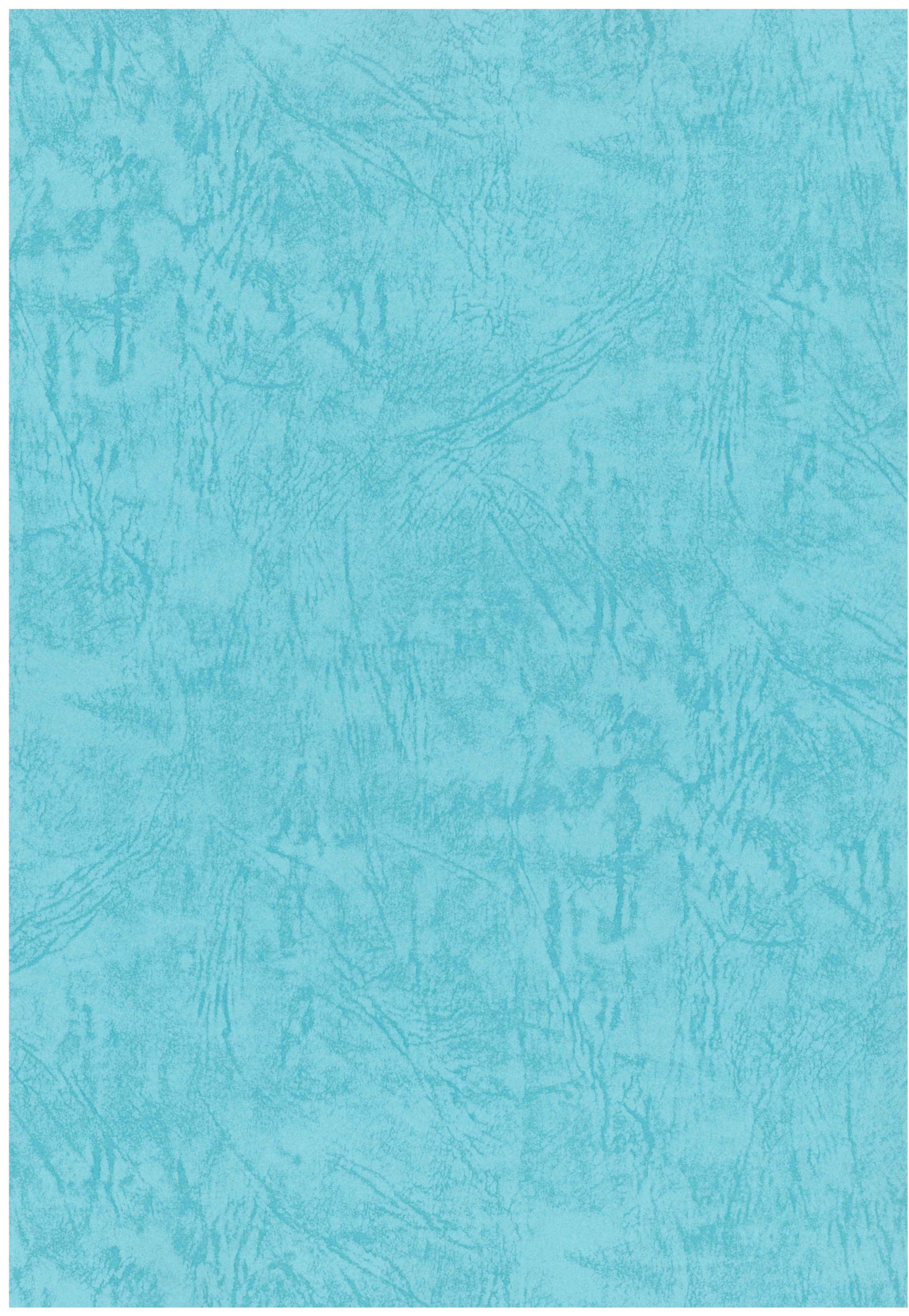


中村和行 (Nakamura, K.)	Identification of up- and down-regulated proteins in gemcitabine-resistant pancreatic cancer cells using two-dimensional gel electrophoresis and mass spectrometry.	Anticancer Res.	30	3367-72	2010
	Heat-shock protein 27 is phosphorylated in gemcitabine-resistant pancreatic cancer cells.	Anticancer Res.	30	2539-43	2010
同 上	Proteomic analysis for nuclear proteins related to tumour malignant progression: a comparative proteomic study between malignant progressive cells and regressive cells.	Anticancer Res.	30	2093-9	2010
同 上	Blue SeePico™ stain after Flamingo™ fluorescent gel stain is useful for cancer proteomic analysis by means of two-dimensional gel electrophoresis.	Anticancer Res.	30	4001-5	2010
同 上	Identification of up- and down-regulated proteins in gemcitabine-resistant pancreatic cancer cells using two-dimensional gel electrophoresis and mass spectrometry.	Anticancer Res.	30	3367-72	2010

中村和行 (Nakamura, K.)	Heat-shock protein 27 is phosphorylated in gemcitabine-resistant pancreatic cancer cells.	Anticancer Res.	30	2539-43	2010
同 上	Proteomic analysis for nuclear proteins related to tumour malignant progression: a comparative proteomic study between malignant progressive cells and regressive cells.	Anticancer Res.	30	2093-9	2010



201207005B(2/3)

厚生労働科学研究費補助金

創薬基盤推進研究事業

疾患関連創薬バイオマーカー探索研究

平成20年度～24年度 総合研究報告書

(分冊 2/2冊)

研究代表者 山西 弘一

平成 25 (2013) 年5 月

厚生労働科学研究費補助金

創薬基盤推進研究事業

疾患関連創薬バイオマーカー探索研究

平成20年度～24年度 総合研究報告書

(分冊 2/2冊)

研究代表者 山西 弘一

平成25 (2013) 年5 月

目 次

I. 総括研究報告	
疾患関連創薬バイオマーカー探索研究 -----	1
山西 弘一	
II. 分担研究報告	
1. 次世代プロテオミクス解析技術による大規模なバイオマーカーの探索と検証 -----	86
朝長 毅	
2. 疾患関連タンパク質の解析基盤の研究 -----	124
角田 慎一	
3. プロテオミクス手法による癌の創薬標的分子探索 -----	139
仲 哲治	
4. ターゲットプロテオミクスを用いた網羅的タンパク質 解析技術の開発とバイオマーカー探索への応用 -----	173
中山 敬一	
5. 創薬バイオマーカー探索研究基盤の確立とその活用 -----	180
平野 久	
6. 2DICAL 法による微量タンパク質解析技術の研究 -----	191
尾野 雅哉	
7. 循環器疾患に関連する微量タンパク質解析技術の研究 -----	199
寒川 賢治、南野 直人	
8. 精神・神経疾患に関連する微量タンパク質解析技術の研究 -----	207
高坂 新一	
9. 新規糖鎖腫瘍マーカーおよび血液中腫瘍由来 DNA の研究 -----	219
加藤 菊也、宮本泰豪	
10. 血清・血漿の前処理法に関する微量タンパク質解析技術の研究： 血清・血漿を用いたプロテオーム解析の臨床検査応用 -----	226
野村 文夫	
11. 脳腫瘍に関連する微量タンパク質解析技術の研究： 統合プロテオミクスによるバイオマーカー／治療ターゲットとなる 脳神経系腫瘍組織細胞内シグナル分子群の解析 -----	240
荒木 令江	
12. 自己抗体を活用した難治性がんのバイオマーカー探索研究 -----	258
中村 和行	
III. 研究成果の刊行に関する一覧表 -----	272
IV. 研究成果の刊行物・別刷 -----	300

Identification of Missing Proteins in the neXtProt Database and Unregistered Phosphopeptides in the PhosphoSitePlus Database As Part of the Chromosome-Centric Human Proteome Project

Takashi Shiromizu,[†] Jun Adachi,[†] Shio Watanabe,[†] Tatsuo Murakami,[†] Takahisa Kuga,[†] Satoshi Muraoka,[†] and Takeshi Tomonaga^{*,†,‡}

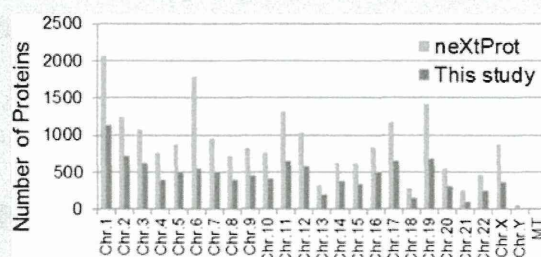
[†]Laboratory of Proteome Research, National Institute of Biomedical Innovation, Ibaraki, Osaka, Japan

[‡]Clinical Proteomics Research Center, Chiba University Hospital, Chiba, Japan

S Supporting Information

ABSTRACT: The Chromosome-Centric Human Proteome Project (C-HPP) is an international effort for creating an annotated proteomic catalog for each chromosome. The first step of the C-HPP project is to find evidence of expression of all proteins encoded on each chromosome. C-HPP also prioritizes particular protein subsets, such as those with post-translational modifications (PTMs) and those found in low abundance. As participants in C-HPP, we integrated proteomic and phosphoproteomic analysis results from chromosome-independent biomarker discovery research to create a chromosome-based list of proteins and phosphorylation sites. Data were integrated from five independent colorectal cancer (CRC) samples (three types of clinical tissue and two types of cell lines) and lead to the identification of 11,278 proteins, including 8,305 phosphoproteins and 28,205 phosphorylation sites; all of these were categorized on a chromosome-by-chromosome basis. In total, 3,033 “missing proteins”, i.e., proteins that currently lack evidence by mass spectrometry, in the neXtProt database and 12,852 unknown phosphorylation sites not registered in the PhosphoSitePlus database were identified. Our in-depth phosphoproteomic study represents a significant contribution to C-HPP. The mass spectrometry proteomics data have been deposited to the ProteomeXchange Consortium with the data set identifier PXD000089

KEYWORDS: Chromosome-Centric Human Proteome Project, missing protein, phosphopeptide, IMAC, colorectal cancer, FASP, neXtProt, PhosphoSitePlus



■ INTRODUCTION

The Chromosome-Centric Human Proteome Project (C-HPP) is a worldwide effort by proteomics researchers to create expression profiles of the approximately 20,000 genes encoded on all human chromosomes and build a database.¹ Protein expression patterns are closely associated with the location of the gene on a chromosome and are correlated with diseases associated with chromosomal abnormalities. Therefore, a comprehensive understanding of the protein expression profile of each chromosome is critical for biological studies and clinical research. The initial aim of C-HPP was to identify at least one protein isoform for every gene encoded by the human genome. Proteins not detected by antibody or proteomic analysis using mass spectrometry are called “missing proteins”.² Currently, there are about 6,000 missing proteins among all of the proteins in the neXtProt database.³ One reason why missing proteins are undetectable is that protein expression differs significantly between tissue and cell types. Although the number of proteins that can be identified in a single analysis has greatly increased due to recent advances in mass spectrometric techniques, complete expression profiles of all proteins will require the integration and analysis of data from a wide variety of samples.

C-HPP also aims to map specific protein variations such as post-translational modifications (PTMs), alternative splicing, and protease-processed variants.² Protein phosphorylation is a key regulator of cellular signal transduction processes, and its deregulation is involved in the onset and progression of various human diseases such as cancer and inflammatory and metabolic disorders.^{4–7} Recent advances in proteomics, especially phosphopeptide enrichment strategies such as immobilized metal ion affinity chromatography (IMAC) and TiO₂ affinity chromatography,⁸ have enabled the identification of up to several thousands of site-specific phosphorylation events within one large-scale analysis.^{9–19}

As participants in C-HPP, we have integrated proteomic and phosphoproteomic analysis data from human colorectal cancer tissue and cell lines and created a chromosome-based list of identified proteins. Newly detected proteins and phosphorylated peptides were identified from the neXtProt and PhosphoSitePlus databases.

Special Issue: Chromosome-centric Human Proteome Project

Received: August 30, 2012



MATERIALS AND METHODS

Tissue and Cell Culture Samples

Colorectal cancer tissue and tumor-adjacent normal tissue samples were obtained from 44 patients at Chiba University School of Medicine. Tissue samples were frozen in liquid nitrogen and stored at -80°C until analysis. Written informed consent was obtained from each patient before surgery, and the protocol was approved by the ethics committees of the Proteome Research Center, National Institute of Biomedical Innovation, and the Chiba University School of Medicine. Cell cultures used were HCT116, SW480, and SW620. HCT116, a colorectal cancer cell line, was grown in RPMI 1640 medium with 10% fetal bovine serum (Invitrogen, Carlsbad, CA, USA) and penicillin/streptomycin (Invitrogen). Cells were maintained at 37°C in an incubator supplemented with 5% CO_2 until they grew to 80% confluence. SW480 and SW620, colon cancer cell lines, were grown at 37°C and 5% CO_2 for at least five passages in SILAC media (R1780-RPMI-1640 without arginine, lysine, leucine (Sigma–Aldrich Corp., St. Louis, MO, USA) with 10% dialyzed fetal bovine serum (Invitrogen) and 100 U/mL penicillin/streptomycin (Invitrogen)) containing 84 mg/L L-arginine (Arg0) and 40 mg/L L-lysine (Lys0) (light), or $^{13}\text{C}_6$ $^{15}\text{N}_4$ -L-arginine (Arg10) and $^{13}\text{C}_6$ -L-lysine (Lys6) (heavy) and 50 mg/L L-leucine.

Protein Extraction and Digestion

Protein extraction and proteolytic digestion were performed using a filter-assisted sample preparation (FASP) protocol.²⁰ Tissue samples or pellets of cultured cells were homogenized by sonication in FASP buffer [1% SDS, 0.1 M DTT, in 0.1 M Tris/HCl, pH 7.6 and PhosSTOP phosphatase inhibitor cocktail (Roche, Mannheim, Germany)]. Protein concentration was determined using a DC protein assay kit (Bio-Rad, Richmond, CA, USA). A total of 10 mg (for phosphoproteomic analysis) or 100 μg (for proteomic analysis) of extracted proteins was digested using 1:100 (w/w) trypsin (proteomics grade; Roche) for 12 h at 37°C . Digested peptides were concentrated and purified using a C18 Sep-PAK cartridge (Waters, Milford, MA, USA).

Phosphopeptide Enrichment

Phosphopeptide enrichment was performed using immobilized Fe(III) affinity chromatography (Fe-IMAC) as described previously.²¹ The Fe-IMAC resin was prepared from Probond (Nickel-Chelating Resin; Invitrogen) by substituting Ni^{2+} on the resin with Fe^{3+} . Ni^{2+} was released from Probond upon treatment with 50 mM EDTA-2Na, and then Fe^{3+} was chelated to the ion-free resin upon incubation with 100 mM FeCl_3 in 0.1% acetic acid. The Fe-IMAC resin was packed into an open column for large-scale enrichment. Following equilibration of the resin with loading solution (60% acetonitrile/0.1% TFA), the peptide mixture was loaded onto the IMAC column. After washing with loading solution (9 times the volume of the IMAC resin) and 0.1% TFA (3 times the volume of the IMAC resin), phosphopeptides were eluted using 1% phosphoric acid (2 times the volume of the IMAC resin).

iTRAQ Labeling

Enriched peptides were labeled with isobaric tags for relative and absolute quantification reagents (iTRAQ 4 plex; Applied Biosystems, Foster City, CA, USA) according to the manufacturer's instructions. Peptide mixtures desalted with C18 Stage-Tips were incubated in the iTRAQ reagents for 1 h.

iTRAQ 115, 116, and 117 were used for labeling individual samples, and iTRAQ 114 was used as the reference sample, a mixture of aliquots of all samples. The reaction was terminated by the addition of an equal volume of distilled water. The labeled samples were combined, acidified with trifluoroacetic acid, and desalted with C18 Stage-Tips.

Strong Cation Exchange Chromatography (SCX)

The peptides were fractionated using a HPLC system (Shimadzu Prominence UFLC) fitted with an SCX column (50 mm \times 2.1 mm, 5 μm , 300 Å, ZORBAX 300SCX; Agilent Technology). The mobile phases consisted of buffer A [25% acetonitrile and 10 mM KH_2PO_4 (pH 3)] and B [25% acetonitrile, 10 mM KH_2PO_4 (pH 3), and 1 M KCl]. The labeled peptides were dissolved in 200 μL of buffer A and separated at a flow rate of 200 $\mu\text{L}/\text{min}$ using a four-step linear gradient: 0% B for 30 min, 0% to 10% B in 15 min, 10% to 25% B in 10 min, 25% to 40% B in 5 min, 40% to 100% B in 5 min, and 100% B for 10 min. Fractions were collected and desalted using C18-Stage Tips (number of fractions: CRC tissues_1 peptides, 30 fractions; CRC tissues_1 phosphopeptides, 25 fractions; HCT116 peptides, 34 fractions; HCT116 phosphopeptides, 32 fractions; SW480 + SW620 peptides, 25 fractions; SW480 + SW620 phosphopeptides, 30 fractions; CRC tissues_2 non-tumor phosphopeptides, 30 fractions; CRC tissues_2 tumor phosphopeptides, 30 fractions).

LC–MS/MS Analysis

Fractionated peptides were analyzed using an LTQ-Orbitrap Velos mass spectrometer (Thermo Fisher Scientific, Bremen, Germany) equipped with a nanoLC interface (AMR, Tokyo, Japan), a nanoHPLC system (Michrom Paradigm MS2) and an HTC-PAL autosampler (CTC, Analytics, Zwingen, Switzerland). The analytical column was made in-house by packing L-column2 C18 particles (Chemical Evaluation and Research Institute (CERI), Tokyo, Japan), into a self-pulled needle (200 mm length \times 100 μm inner diameter). The mobile phases consisted of buffer A (0.1% formic acid and 2% acetonitrile) and B (0.1% formic acid and 90% acetonitrile). Samples dissolved in buffer A were loaded onto a trap column (0.3 \times 5 mm, L-column ODS; CERI). The nanoLC gradient was delivered at 500 nL/min and consisted of a linear gradient of buffer B developed from 5% to 30% B in 180 min. A spray voltage of 2000 V was applied.

Full MS scans were performed using an orbitrap mass analyzer (scan range m/z 350–1500, with 30 K fwhm resolution at m/z 400). The 10 most intense precursor ions were selected for the MS/MS scans, which were performed using collision-induced dissociation (CID) and higher energy collision-induced dissociation (HCD, 7500 fwhm resolution at m/z 400) for each precursor ion. The dynamic exclusion option was implemented with a repeat count of 1 and exclusion duration of 60 s. Automated gain control (AGC) was set to $1.00\text{e} + 06$ for full MS, $1.00\text{e} + 04$ for CID MS/MS, and $5.00\text{e} + 04$ for HCD MS/MS. The normalized collision energy values were set to 35% for CID and 50% for HCD.

The CID and HCD raw spectra were extracted and searched separately against UniProtKB/Swiss-Prot (release-2010_05), which contains 20,295 sequences (the forward and reverse-decoy) of *Homo sapiens*, using Proteome Discoverer 1.3 (Thermo Fisher Scientific) and Mascot v2.3. The precursor mass tolerance was set to 7 ppm, and fragment ion mass tolerance was set to 0.5 Da for CID and 0.01 Da for HCD. The search parameters allowed two missed cleavage for trypsin,

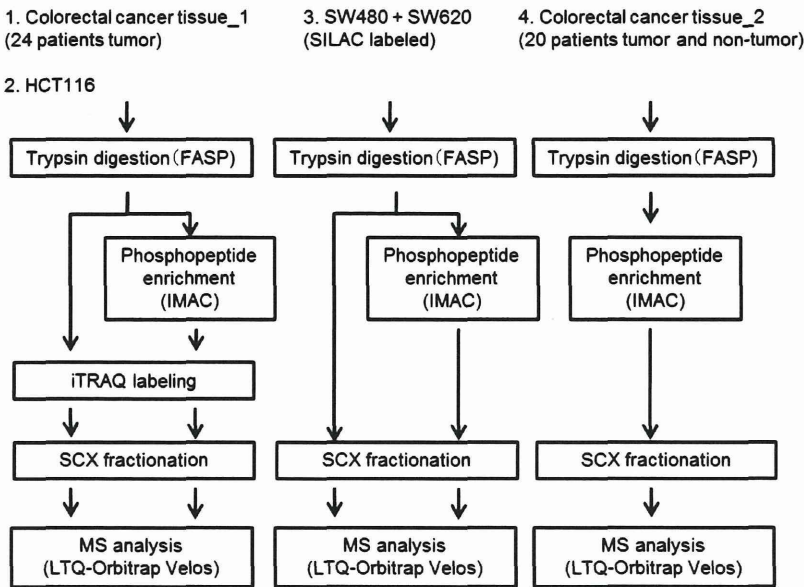


Figure 1. Schematic representation of the experimental work flow for the proteomic and phosphoproteomic analyses of the four experiments. SW480 + SW620: a mixture of protein extracts obtained from SW480 and SW620 cells. After trypsin digestion, each sample was separated for proteomic (100 μ g) or phosphoproteomic (10 mg) analysis. Digested samples were separated by using an SCX column. LC–MS/MS, requiring 3-h runs, was performed using an LTQ-Orbitrap Velos.

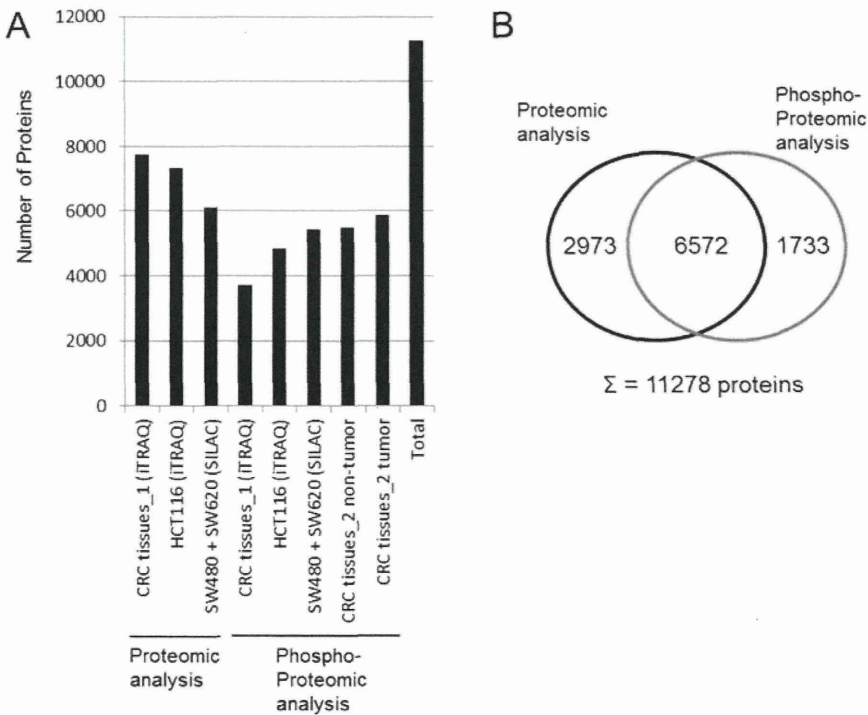


Figure 2. Number and overlap of identified proteins from the proteomic and phosphoproteomic analyses. (A) Number of identified proteins from the proteomic and phosphoproteomic analyses of the 8 data sets. (B) Proportion of proteins identified in each analysis and overlap between proteins identified by the proteomic and phosphoproteomic analyses.

fixed modifications (carbamidomethylation at cysteine), and variable modifications (oxidation at methionine). Fixed modifications were set for CRC tissue and HCT116 (iTRAQ labeling at lysine and the N-terminal residue) and SW480 + SW620 (SILAC labeling 13C(6) 15N(4) Arg, 13C(6) Lys). Variable modifications were added for phosphoproteomic analysis (phosphorylation at serine, threonine, and tyrosine). In the workflow of Proteome Discoverer 1.3, following the Mascot search, the phosphorylated sites on the identified

peptides were assigned again using the PhosphoRS algorithm, which calculated the possibility of the phosphorylated site from the spectra matched to the identified peptides.²² The score threshold for peptide identification was set at 1% false-discovery rate (FDR) and 75% phosphoRS site probability. FDR was calculated using the Percolator algorithm for peptide sequence analysis. Percolator uses >30 features of a peptide spectral match (PSM) to distinguish true positives from random matches.

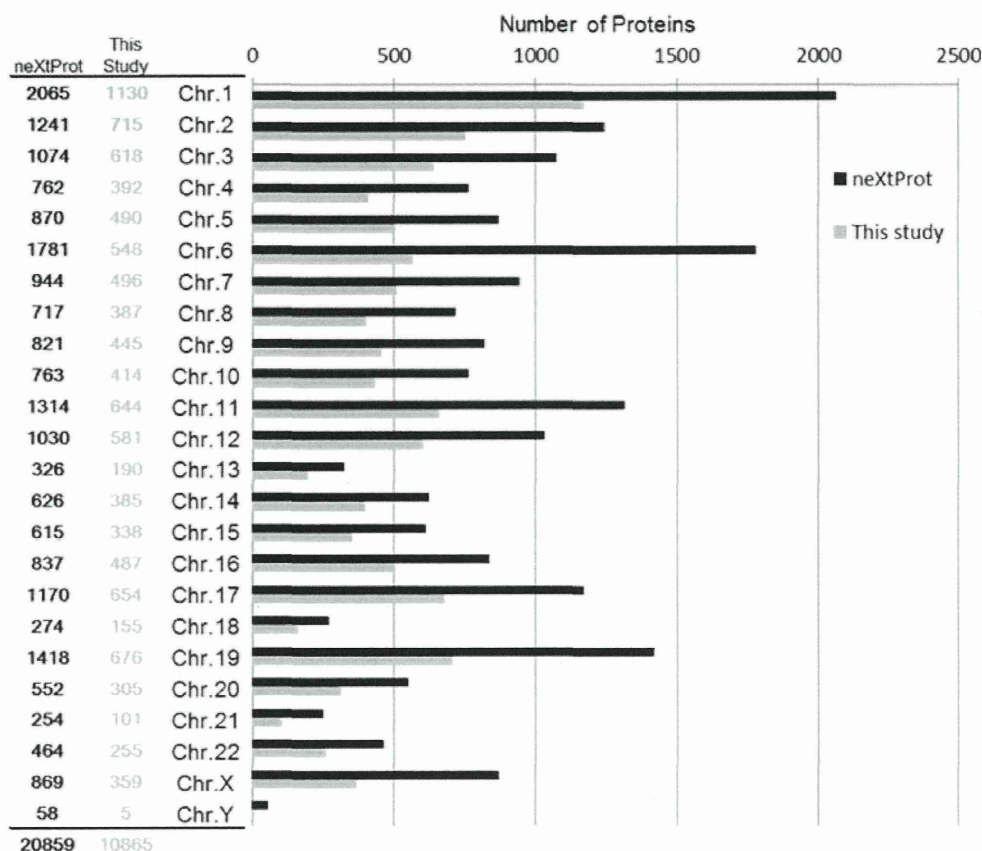


Figure 3. Chromosomal distribution of the identified proteins (gray) in relation to total proteins (black) registered in the neXtProt database.

Bioinformatics Analysis

Chromosomal locations and missing protein analyses of identified proteins were elucidated using the neXtProt database (<http://www.nextprot.org/db/>), and identified phosphorylation sites were elucidated using the PhosphoSitePlus database (<http://www.phosphosite.org/>). The function of identified missing proteins was elucidated by ingenuity pathway analysis software (Ingenuity Systems, Redwood City, USA).

Stable Isotope-Labeled Peptides

Stable isotope-labeled standard peptides (SIS peptides, crude grade) were synthesized (Thermo Fisher Scientific, Ulm, Germany). A single lysine was replaced by isotope-labeled lysine ($^{13}\text{C}_6$, 98%, $^{15}\text{N}_2$, 98%). The SIS peptides were dissolved in distilled water at a concentration of $1\ \mu\text{g}/\mu\text{L}$ and stored at $-80\ ^\circ\text{C}$. A mixture of these SIS peptides was added to colorectal carcinoma phosphoproteomic samples.

RESULTS

As part of the C-HPP project, we combined the eight data sets from four different experiments obtained from colorectal cancer tissue and colon cancer cells; these experiments included three quantitative analyses and one non-quantitative analysis. Colorectal cancer tissues and colon cancer cells were first solubilized and trypsin-digested using the FASP method.²⁰ Phosphopeptides were then enriched using the IMAC method. These peptides and phosphopeptides were fractionated on a Strong Cation-Exchange (SCX) column before LC-MS/MS using an LTQ-Orbitrap mass spectrometer (Figure 1). Proteome Discoverer 1.3 software was used to analyze the RAW data files, Mascot was used as the search engine, and UniProtKB/Swiss-Prot (release-2010_05) was the database. Following data

integration, 11,278 proteins were identified with Peptide FDR ≤ 1.0 containing at least one unique peptide corresponding to one protein in the database (Figure 2A, Supplementary Table 1–4). Of these, 8,305 proteins were identified as phosphorylated. Among the total identified proteins, 673 proteins were identified only with CID, and 386 proteins were identified only with HCD. Also, 4924 phosphopeptides were identified only with CID, and 3538 phosphopeptides were identified only with HCD. A total of 6,572 proteins were commonly identified in the proteomic and phosphoproteomic analyses (Figure 2B). However, a proportion of proteins were found not to overlap in the analyses. This is probably due to the abundance and complexity of the proteins and phosphoproteins in the samples, which prevent proteomics and phosphoproteomics to identify all of the proteins and phosphoproteins present.

Quantitative analyses were performed to investigate the differences between metastatic and non-metastatic cases by using clinical tissue and two types of cultured cells (a mixture of SW620 + SW480 and HCT116 cells). Clinical tissue samples of primary colorectal cancer obtained from 12 patients with or without metastasis were pooled. Cancers without metastasis were labeled with iTRAQ 114 or 116, and those with metastasis were labeled with iTRAQ 115 or 117. We also performed quantitative analyses between metastatic and non-metastatic cell lines. HCT116 metastatic clone was established by orthotopic implantation model mouse, and its protein expression was compared with that of the parent clone. SW620 cell line is a lymph node metastatic variant of SW480. HCT116 parent clone was labeled with iTRAQ 114 or 115, and metastatic clone was labeled with iTRAQ 116 or 117. SW480 and SW620 were reciprocally labeled with light and heavy

Table 1. Number of Identified Proteins in Each Chromosome

	neXtProt	total	proteomic analysis			phosphoproteomic analysis				
			CRC tissue_1	HCT116	SW480 + SW620	CRC tissue_1	HCT116	SW480 + SW620	CRC tissue_2 non-tumor	CRC tissue_2 tumor
Chr.1	2065	1171	808	767	611	356	494	554	568	598
Chr.2	1241	753	516	481	416	259	305	378	371	414
Chr.3	1074	642	445	425	360	209	275	310	321	341
Chr.4	762	412	267	251	185	126	178	176	210	224
Chr.5	870	508	339	319	285	168	210	249	255	281
Chr.6	1781	569	395	347	302	186	233	260	272	291
Chr.7	944	511	364	319	305	160	204	269	244	253
Chr.8	717	402	270	263	196	125	177	188	176	199
Chr.9	821	458	312	294	229	161	197	212	230	225
Chr.10	763	434	298	300	254	131	186	223	232	245
Chr.11	1314	660	477	447	388	227	299	340	352	360
Chr.12	1030	603	399	409	330	207	274	296	294	316
Chr.13	326	194	140	127	114	62	80	97	87	99
Chr.14	626	400	283	252	218	132	166	180	203	210
Chr.15	615	353	237	234	181	121	163	184	174	197
Chr.16	837	508	341	333	260	159	230	245	236	258
Chr.17	1170	678	437	484	399	239	330	339	347	361
Chr.18	274	161	102	112	74	50	67	64	70	77
Chr.19	1418	707	479	471	362	253	336	332	341	369
Chr.20	552	313	223	194	182	119	137	151	138	160
Chr.21	254	103	78	68	64	36	41	51	47	53
Chr.22	464	262	187	182	153	93	109	123	129	130
Chr.X	869	369	261	222	222	120	144	189	175	188
Chr.Y	58	6	4	2	1	1	0	0	1	2
NA ^a		101	73	30	17	11	11	17	18	16
total	20845	11278	7735	7333	6108	3355	4352	5427	5491	5867

^aNot applicable in neXtProt.

stable isotope amino acids (lysine and arginine). HCT116 has a mutation in codon 13 of the ras protooncogene, while SW480 and SW620 have a mutation in codon 12. Among 8305 proteins and 28,205 phosphopeptides, 472 proteins and 2547 phosphopeptides showed >2-fold differences between metastatic and non-metastatic tissues and cell lines (either upregulated or downregulated).

A total of 20,845 proteins have been registered in the neXtProt database. Proteins identified in this study were referred to the database and accounted for 53.6% (11,177/20,845) of all the proteins registered in the neXtProt database; their chromosomal locations are shown in Figure 3 and Table 1. Of the proteins registered in the neXtProt database, the expression of 14,612 proteins (70.1% of the total of 20,845 proteins) has been confirmed by mass spectrometry or antibody assay (protein level 'yes'), whereas 10,649 proteins (51.1% of the total of 20,845 proteins) have been identified only by MS analysis (proteomic level 'yes') (Table. 2). Cross-checking the 11,278 proteins identified in this study with the neXtProt database revealed 1,145 proteins currently lacking evidence of protein expression by mass spectrometry or

antibody assays, and 3,033 proteins lacking evidence by mass spectrometry. These “missing proteins (protein level = no and proteomic level = no)” are listed on a chromosome-by-chromosome basis (Figure 4).

In contrast, 28,205 phosphorylation sites were identified (Supplementary Table 5). When these phosphopeptides were cross-checked with the PhosphoSitePlus database, 15,353 registered phosphorylation sites were identified, or 12.2% of all registered sites in PhosphoSitePlus (15,353/125,433). Of these, 12,852 sites were not registered in PhosphoSitePlus (Figure 5A). The chromosomal locations of these phosphorylation sites are shown in Figure 5B. In order to verify the accuracy of the identified phosphopeptides, lysine at the C-terminus of two peptides (LYNSEESRPYTNK, SASQSLDKLDQELK) was labeled by stable isotope (¹³C₆, ¹⁵N₂). The SIS peptides were added to the extract of colorectal cancer tissue, and annotated mass spectra and extracted ion chromatogram data of SIS peptides were compared to those of nonlabeled endogenous peptides (Supplementary Figure 1).

Non-quantitative analyses were also performed using pooled colorectal carcinoma tissues and tumor-adjacent normal tissues 5–10 cm remote from the tumor. To investigate the association between phosphoproteins and biological function, gene ontology analysis was performed by using Ingenuity Pathway Analysis (IPA) software. Specifically identified phosphoproteins in normal (636 proteins) and carcinoma tissues (1020 proteins) were also analyzed by IPA. Molecular functions involved in cell cycle (normal = 9 proteins, tumor = 132 proteins; *p* < 0.01 Fisher's exact test) and DNA replication (normal = 15 proteins, tumor = 106 proteins ; *p* < 0.01)

Table 2. Number of Proteins Identified at the Protein or Proteomic Level			
evidence		neXtProt	this study
protein level	yes	14612	10032
	no	6233	1145
proteomics level	yes	10649	8144
	no	10196	3033

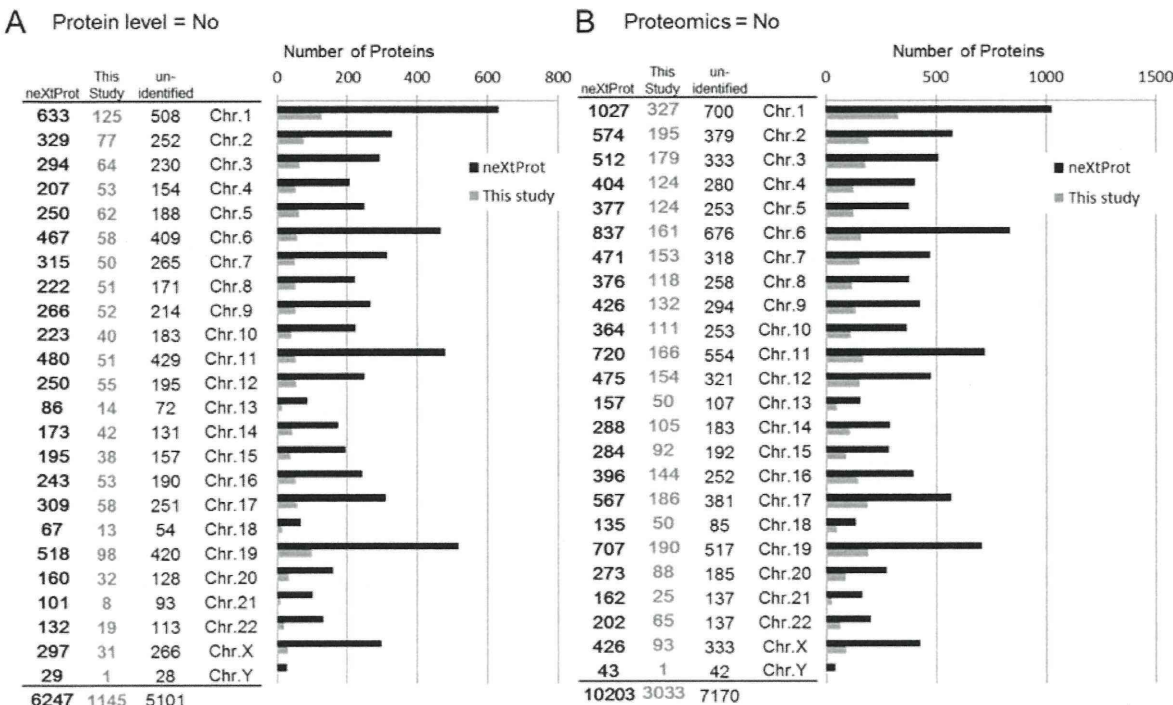


Figure 4. Chromosomal distribution of the identified proteins (gray) and total registered proteins (black) in the neXtProt database with no evidence of expression at the protein level (A) and at the proteomic level (B).

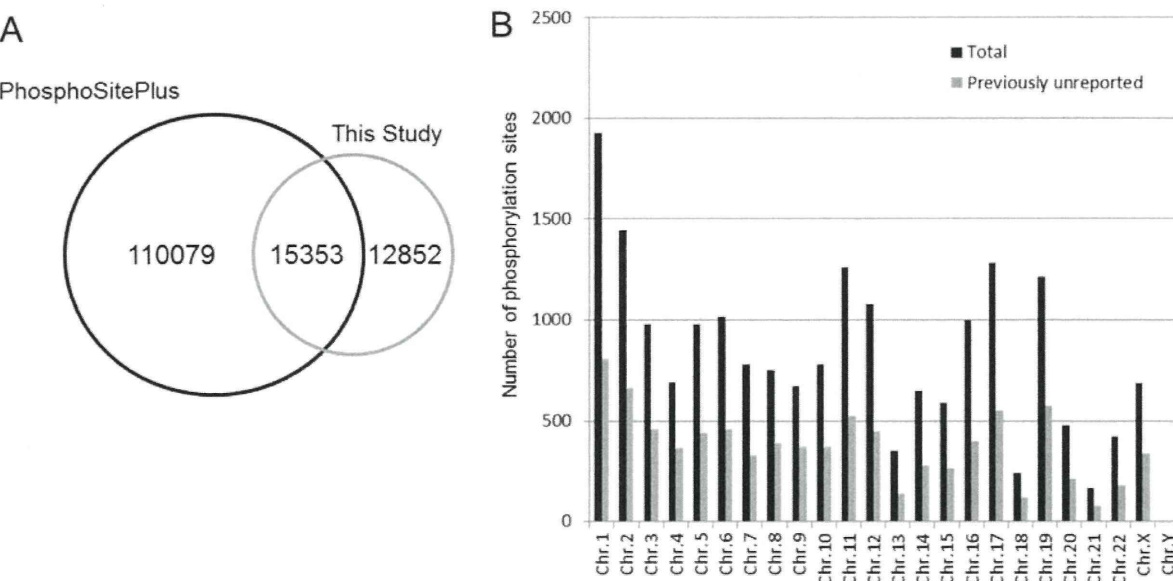


Figure 5. (A) Overlap between phosphorylation sites identified in this study and those registered in the PhosphoSitePlus database. (B) Chromosomal distribution of the identified previously unreported phosphorylation sites (gray) and total registered sites (black).

functions were more abundant in cancer tissues than in normal tissues (Supplementary Figure 2).

DISCUSSION

The objective of C-HPP is to map and annotate all protein-coding genes on each human chromosome. C-HPP also prioritizes particular protein subsets such as post-translational modifications (PTMs) and low-abundance proteins. Thus, we have integrated proteomic and phosphoproteomic data obtained from a shotgun analysis using human cancer tissue and cell lines prepared for various purposes. We have integrated quantitative and non-quantitative data; quantitative analysis was

performed for the relative quantification between metastatic and non-metastatic colorectal carcinoma samples, while non-quantitative analysis was performed to compare the tumor and normal tissues. As a result, we identified 11,278 proteins, 8,305 phosphoproteins, and 28,205 phosphorylation sites, and their chromosomal locations were defined using the neXtProt database. Furthermore, we were able to identify 3,033 missing proteins that currently lack evidence by mass spectrometry and 12,852 unknown phosphorylation sites that are not in the PhosphoSitePlus database.

Currently, the research group with the most advanced mass analysis system can identify over 10,000 proteins in a single

analysis run and identify about 50% of the proteins in their comprehensive analyses using multiple cell lines.²³ Additionally, the number of phosphorylation sites identified has exponentially increased,²⁴ largely due to improvements in phosphopeptide enrichment methods such as IMAC¹⁵ and TiO₂ affinity chromatography.⁸ A phosphoproteomic study of HeLa cells identified more than 65,000 phosphopeptides using a combination of phosphopeptide enrichment and SCX chromatography.²⁵ Several phosphoproteomic studies using tissue samples have been reported and have identified 5,195 phosphopeptides from human dorsolateral prefrontal cortex²⁶ and 5,698 phosphorylation sites from tumor tissues of melanoma model mice.²⁷ In our study, we identified 11,278 proteins and 28,205 phosphorylation sites; some had been identified in previous reports, but a number of the proteins and phosphorylation sites are not listed in the neXtProt or PhosphoSitePlus databases. Since mass analysis systems are rapidly becoming more powerful, in the future an individual research group may be able to identify all the proteins in the human genome in one analysis. However, in order to build an extensively annotated proteome database, which is one purpose of the C-HPP project, it is necessary to combine the analyzed data of various samples from many research groups.

Our analysis increased the number of identified proteins by combining the results of proteome analysis and phosphoproteome analysis on identical samples. Even using commonly studied cell lines, combining the results of post-translational modification analysis and analysis of fractionated samples will increase the number of identified missing proteins. The data presented here are based on relative quantification, and thus to confirm protein expression and examine protein abundance and localization, validation using antibodies or selected reaction monitoring (SRM) is required. Such validations will benefit from information on the identified cell line, sample preparation methods, MS analysis data, and the sequences of the identified peptides/phosphopeptides. We and other researchers, including Muraoka and colleagues,²⁸ Narumi and colleagues (unpublished data), and Kume and colleagues (unpublished data), are currently using a strategy for large-scale proteomics and SRM-based validation to discover biomarkers for various diseases and aim to obtain additional proteomics data by SRM validation and quantitation that will be integrated into the C-HPP project.

■ ASSOCIATED CONTENT

● Supporting Information

Annotated mass spectra and retention time data from liquid chromatography; results of phosphoproteomic analysis in normal and carcinoma tissues; lists of identified proteins and peptides by phosphoproteomic and proteomic analysis; list of identified phosphorylation sites. The mass spectrometry proteomics data have been deposited to the ProteomeXchange Consortium (<http://proteomecentral.proteomexchange.org>) via the PRIDE partner repository²⁹ with the data set identifier PXD000089. This material is available free of charge via the Internet at <http://pubs.acs.org>.

■ AUTHOR INFORMATION

Corresponding Author

*Tel: +81-72-641-9862. Fax: +81-72-641-9861. E-mail: tomonaga@nibio.go.jp.

Notes

The authors declare no competing financial interest.

■ ACKNOWLEDGMENTS

This work was supported by a Grant-in-Aid for Research on Biological Markers for New Drug Development (H20-0005 to T.T.) from the Ministry of Health, Labour, and Welfare of Japan. This work was also supported by Grants-in-Aid (21390354 to T.T. and 23701093 to T.S.) from the Ministry of Education, Science, Sports, and Culture of Japan. The data deposition to the ProteomeXchange Consortium was supported by PRIDE Team, EBI.

■ ABBREVIATIONS

C-HPP, Chromosome-Centric Human Proteome Project; CRC, colorectal cancer; PTMs, post-translational modifications; IMAC, immobilized metal ion affinity chromatography; FASP, filter-assisted sample preparation; SCX, strong cation-exchange; FDR, false discovery rate; SRM, selected reaction monitoring; CID, collision-induced dissociation; HCD, higher energy collision-induced dissociation; LC-MS/MS, liquid chromatography tandem mass spectrometry; CE, collision energy; LTQ, linear ion trap; fwhm, full width at half-maximum

■ REFERENCES

- (1) Hancock, W.; Omenn, G.; Legrain, P.; Paik, Y. K. Proteomics, human proteome project, and chromosomes. *J. Proteome Res.* **2011**, *10* (1), 210.
- (2) Paik, Y. K.; Jeong, S. K.; Omenn, G. S.; Uhlen, M.; Hanash, S.; Cho, S. Y.; Lee, H. J.; Na, K.; Choi, E. Y.; Yan, F.; Zhang, F.; Zhang, Y.; Snyder, M.; Cheng, Y.; Chen, R.; Marko-Varga, G.; Deutsch, E. W.; Kim, H.; Kwon, J. Y.; Aebersold, R.; Bairoch, A.; Taylor, A. D.; Kim, K. Y.; Lee, E. Y.; Hochstrasser, D.; Legrain, P.; Hancock, W. S. The Chromosome-Centric Human Proteome Project for cataloging proteins encoded in the genome. *Nat. Biotechnol.* **2012**, *30* (3), 221–223.
- (3) Lane, L.; Argoud-Puy, G.; Britan, A.; Cusin, I.; Duek, P. D.; Evalet, O.; Gateau, A.; Gaudet, P.; Gleizes, A.; Masselot, A.; Zwahlen, C.; Bairoch, A. neXtProt: a knowledge platform for human proteins. *Nucleic Acids Res.* **2012**, *40* (Database issue), D76–83.
- (4) Hanahan, D.; Weinberg, R. A. The hallmarks of cancer. *Cell* **2000**, *100* (1), 57–70.
- (5) Kaminska, B. MAPK signalling pathways as molecular targets for anti-inflammatory therapy—from molecular mechanisms to therapeutic benefits. *Biochim. Biophys. Acta* **2005**, *1754* (1–2), 253–62.
- (6) Peifer, C.; Wagner, G.; Laufer, S. New approaches to the treatment of inflammatory disorders small molecule inhibitors of p38 MAP kinase. *Curr. Top. Med. Chem.* **2006**, *6* (2), 113–49.
- (7) White, M. F. Regulating insulin signaling and beta-cell function through IRS proteins. *Can. J. Physiol. Pharmacol.* **2006**, *84* (7), 725–37.
- (8) Larsen, M. R.; Thingholm, T. E.; Jensen, O. N.; Roepstorff, P.; Jorgensen, T. J. Highly selective enrichment of phosphorylated peptides from peptide mixtures using titanium dioxide microcolumns. *Mol. Cell. Proteomics* **2005**, *4* (7), 873–86.
- (9) Collins, M. O.; Yu, L.; Coba, M. P.; Husi, H.; Campuzano, I.; Blackstock, W. P.; Choudhary, J. S.; Grant, S. G. Proteomic analysis of in vivo phosphorylated synaptic proteins. *J. Biol. Chem.* **2005**, *280* (7), 5972–82.
- (10) Molina, H.; Horn, D. M.; Tang, N.; Mathivanan, S.; Pandey, A. Global proteomic profiling of phosphopeptides using electron transfer dissociation tandem mass spectrometry. *Proc. Natl. Acad. Sci. U.S.A.* **2007**, *104* (7), 2199–204.
- (11) Wissing, J.; Jansch, L.; Nimtz, M.; Dieterich, G.; Hornberger, R.; Kerl, G.; Wehland, J.; Daub, H. Proteomics analysis of protein kinases by target class-selective prefractionation and tandem mass spectrometry. *Mol. Cell. Proteomics* **2007**, *6* (3), 537–47.

- (12) Villen, J.; Beausoleil, S. A.; Gerber, S. A.; Gygi, S. P. Large-scale phosphorylation analysis of mouse liver. *Proc. Natl. Acad. Sci. U.S.A.* **2007**, *104* (5), 1488–93.
- (13) Ballif, B. A.; Villen, J.; Beausoleil, S. A.; Schwartz, D.; Gygi, S. P. Phosphoproteomic analysis of the developing mouse brain. *Mol. Cell. Proteomics* **2004**, *3* (11), 1093–101.
- (14) Beausoleil, S. A.; Jedrychowski, M.; Schwartz, D.; Elias, J. E.; Villen, J.; Li, J.; Cohn, M. A.; Cantley, L. C.; Gygi, S. P. Large-scale characterization of HeLa cell nuclear phosphoproteins. *Proc. Natl. Acad. Sci. U.S.A.* **2004**, *101* (33), 12130–5.
- (15) Ficarro, S. B.; McClelland, M. L.; Stukenberg, P. T.; Burke, D. J.; Ross, M. M.; Shabanowitz, J.; Hunt, D. F.; White, F. M. Phosphoproteome analysis by mass spectrometry and its application to *Saccharomyces cerevisiae*. *Nat. Biotechnol.* **2002**, *20* (3), 301–5.
- (16) Lee, J.; Xu, Y.; Chen, Y.; Sprung, R.; Kim, S. C.; Xie, S.; Zhao, Y. Mitochondrial phosphoproteome revealed by an improved IMAC method and MS/MS/MS. *Mol. Cell. Proteomics* **2007**, *6* (4), 669–76.
- (17) Moser, K.; White, F. M. Phosphoproteomic analysis of rat liver by high capacity IMAC and LC–MS/MS. *J. Proteome Res.* **2006**, *5* (1), 98–104.
- (18) Trinidad, J. C.; Specht, C. G.; Thalhhammer, A.; Schoepfer, R.; Burlingame, A. L. Comprehensive identification of phosphorylation sites in postsynaptic density preparations. *Mol. Cell. Proteomics* **2006**, *5* (5), 914–22.
- (19) Li, X.; Gerber, S. A.; Rudner, A. D.; Beausoleil, S. A.; Haas, W.; Villen, J.; Elias, J. E.; Gygi, S. P. Large-scale phosphorylation analysis of alpha-factor-arrested *Saccharomyces cerevisiae*. *J. Proteome Res.* **2007**, *6* (3), 1190–7.
- (20) Wisniewski, J. R.; Zougman, A.; Nagaraj, N.; Mann, M. Universal sample preparation method for proteome analysis. *Nat. Methods* **2009**, *6* (5), 359–62.
- (21) Matsumoto, M.; Oyamada, K.; Takahashi, H.; Sato, T.; Hatakeyama, S.; Nakayama, K. I. Large-scale proteomic analysis of tyrosine-phosphorylation induced by T-cell receptor or B-cell receptor activation reveals new signaling pathways. *Proteomics* **2009**, *9* (13), 3549–63.
- (22) Taus, T.; Kocher, T.; Pichler, P.; Paschke, C.; Schmidt, A.; Henrich, C.; Mechtler, K. Universal and confident phosphorylation site localization using phosphoRS. *J. Proteome Res.* **2011**, *10* (12), 5354–62.
- (23) Geiger, T.; Wehner, A.; Schaab, C.; Cox, J.; Mann, M. Comparative proteomic analysis of eleven common cell lines reveals ubiquitous but varying expression of most proteins. *Mol. Cell. Proteomics* **2012**, *11* (3), M111 014050.
- (24) Lemeer, S.; Heck, A. J. The phosphoproteomics data explosion. *Curr. Opin. Chem. Biol.* **2009**, *13* (4), 414–20.
- (25) Dephoure, N.; Zhou, C.; Villen, J.; Beausoleil, S. A.; Bakalarski, C. E.; Elledge, S. J.; Gygi, S. P. A quantitative atlas of mitotic phosphorylation. *Proc. Natl. Acad. Sci. U.S.A.* **2008**, *105* (31), 10762–7.
- (26) Martins-de-Souza, D.; Guest, P. C.; Vanattou-Saifoudine, N.; Rahmoune, H.; Bahn, S. Phosphoproteomic differences in major depressive disorder postmortem brains indicate effects on synaptic function. *Eur. Arch. Psychiatry Clin. Neurosci.* **2012**, *262*, 657–666.
- (27) Zanivan, S.; Gnad, F.; Wickstrom, S. A.; Geiger, T.; Macek, B.; Cox, J.; Fassler, R.; Mann, M. Solid tumor proteome and phosphoproteome analysis by high resolution mass spectrometry. *J. Proteome Res.* **2008**, *7* (12), 5314–26.
- (28) Muraoka, S.; Kume, H.; Watanabe, S.; Adachi, J.; Kuwano, M.; Sato, M.; Kawasaki, N.; Kodera, Y.; Ishitobi, M.; Inaji, H.; Miyamoto, Y.; Kato, K.; Tomonaga, T. Strategy for SRM-based verification of biomarker candidates discovered by iTRAQ method in limited breast cancer tissue samples. *J. Proteome Res.* **2012**, *11* (8), 4201–10.
- (29) Vizcaino, J. A.; Cote, R.; Reisinger, F.; Barsnes, H.; Foster, J. M.; Rameseder, J.; Hermjakob, H.; Martens, L. The Proteomics Identifications database: 2010 update. *Nucleic Acids Res.* **2010**, *38* (Database issue), D736–742.

In-depth Membrane Proteomic Study of Breast Cancer Tissues for the Generation of a Chromosome-based Protein List

Satoshi Muraoka,[†] Hideaki Kume,[†] Jun Adachi,[†] Takashi Shiromizu,[†] Shio Watanabe,[†] Takeshi Masuda,[‡] Yasushi Ishihama,[§] and Takeshi Tomonaga^{*,†}

[†]Laboratory of Proteome Research, National Institute of Biomedical Innovation, Ibaraki, Osaka, Japan

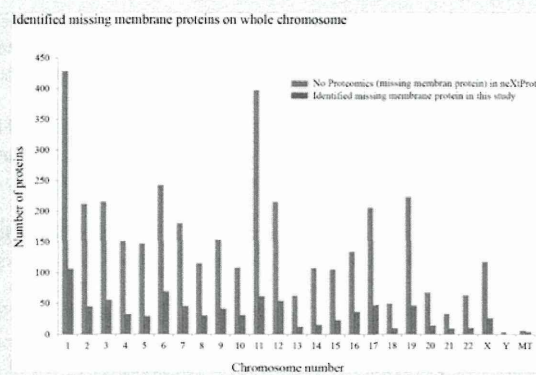
[‡]Institute for Advanced Biosciences, Keio University, Tsuruoka, Yamagata, Japan

[§]Graduate School of Pharmaceutical Sciences, Kyoto University, Sakyo-ku, Kyoto, Japan

Supporting Information

ABSTRACT: The Chromosome-centric Human Proteome Project (C-HPP) aims to define all proteins encoded in each chromosome and especially to identify proteins that currently lack evidence by mass spectrometry. The C-HPP also prioritizes particular protein subsets such as membrane proteins, post-translational modifications, and low-abundance proteins. In this study, we aimed to generate deep profiling of the membrane proteins of human breast cancer tissues on a chromosome-by-chromosome basis using shotgun proteomics. We identified 7092 unique proteins using membrane fractions isolated from pooled breast cancer tissues with high confidence. A total of 3282 proteins were annotated as membrane proteins by Gene Ontology analysis, which covered 45% of the membrane proteins predicted in 20 859 protein-coding genes. Furthermore, we were able to identify 851 membrane proteins that currently lack evidence by mass spectrometry in neXtProt. Our results will contribute to the accomplishment of the primary goal of the C-HPP in identifying so-called “missing proteins” and generating a whole protein catalog for each chromosome.

KEYWORDS: missing protein, shotgun proteomics, membrane protein, neXtProt, chromosome, C-HPP



INTRODUCTION

Completed in 2003, the Human Genome Project (HGP) was a 13-year project coordinated by the U.S. Department of Energy and the National Institutes of Health.¹ The project was to identify all of the approximately 20 000–25 000 genes in human DNA.^{1,2} Results were published as the human genome database. The age of whole-genome sequencing has made the research field of proteomics possible. In 2008, the Human Proteome Organization (HUPO) developed a strategy for the first phase of the human proteome project (HPP). The C-HPP is one component of the HPP and focuses on constructing a proteomic catalog in a chromosome-by-chromosome fashion and aims to define the full set of proteins encoded in whole-chromosomes.^{3–5} The initial goal of the C-HPP is to identify and characterize proteins that currently lack MS evidence, referred to as “missing proteins”, in neXtProt, a new human protein-centric knowledge platform.⁶ “Missing proteins” are likely due to their very low-abundance and/or absence of expression in given cells or tissues. Thus, more in-depth proteomic studies of cell lines and patient tissues are needed.

The C-HPP also underscores the mapping of particular protein subsets such as membrane proteins and/or post-translational modifications. Membrane proteins are of great interest, particularly because they could be key biomarkers for

early diagnosis, progression of diseases, and suitable drug targets; however, there have been difficulties in enrichment/solubilization and also subsequent protease digestion in membrane proteome analysis.^{7–9} Recently, several protocols have been reported to increase the solubilization and digestion of proteins, which has greatly improved membrane proteomic analysis of cells and tissues.^{10,11}

In this study, to generate a chromosome-based membrane protein list, we integrated membrane proteomic analysis data from human breast cancer tissues with previous data¹² and analyzed with Proteome Discoverer and Database for Annotation, Visualization and Integrated Discovery (DAVID) Bioinformatics Resources followed by chromosome-based categorization using the neXtProt database.

MATERIALS AND METHODS

Human Tissue Samples

Tissue samples were obtained from 18 patients with high-risk or low-risk MammaPrint breast cancer who underwent surgery at

Special Issue: Chromosome-centric Human Proteome Project

Received: August 30, 2012

Published: November 15, 2012

the Osaka Medical Center for Cancer & Cardiovascular Diseases (Supplementary Figure 1, Supporting Information). All samples were frozen by liquid nitrogen and were stored at -80°C until analysis. Written informed consent was obtained from all subjects. The Ethics Committee of our institute and the Osaka Medical Center for Cancer & Cardiovascular Diseases approved the protocol.

Enrichment of Membrane Proteins

For enrichment of membrane proteins, frozen tissue samples were homogenized in PBS containing a protease inhibitor mixture (Complete; Roche, Mannheim, Germany) using a Dounce homogenizer (WHEATON, Millville, NJ) following centrifugation ($1000\times g$) for 10 min at 4°C . The postnuclear supernatant was centrifuged at $100\,000\times g$ for 1 h at 4°C . The pellet was suspended in ice-cold $0.1\text{ M Na}_2\text{CO}_3$ solution following centrifugation ($100\,000\times g$) for 1 h at 4°C . After centrifugation, the pellet was treated using an MPX PTS reagent kit (GL sciences, Tokyo, Japan) as follows.¹⁰ Briefly, the pellet was solubilized with PTS B buffer at 95°C for 5 min followed by sonication for 5 min using a Bioruptor sonicator (Cosmo Bio, Tokyo, Japan). The solution was centrifuged at $100\,000\times g$ for 30 min at 4°C . Supernatant containing membrane proteins was stored at -80°C . Protein concentration was determined using a DC protein assay kit (Bio-Rad, USA).

In Solution Digestion and iTRAQ Labeling

Membrane proteins from pooled high-risk ($n = 9$) or low-risk ($n = 9$) breast cancer tissue samples were digested with Lys-C (Wako Pure Chemical Industries, Osaka, Japan), followed by trypsin (Proteomics grade; Roche, Swiss). Tryptic digests were treated according to the PTS protocol and desalted using C18 StageTips.¹³ Briefly, a sample of $90\text{ }\mu\text{g}$ of pooled membrane proteins was reduced with 10 mM dithiothreitol (DTT), alkylated with 20 mM iodoacetamide (IAA), and sequentially digested by 1:100 (w/w) LysC (Wako Pure Chemical Industries, Osaka, Japan) for 8 h at 37°C and 1:100 (w/w) trypsin (proteomics grade; Roche) for 12 h at 37°C . An equal volume of an organic solvent, ethyl acetate, was added to digested samples, the mixtures were acidified by 1% trifluoroacetic acid (TFA), and vortexed to transfer detergents to the organic phase. After centrifugation, the aqueous phase containing peptides was collected. BSA ($0.45\text{ }\mu\text{g}$) was spiked into membrane protein samples as a quality control for iTRAQ labeling. The tryptic digest sample was desalted using C18 stage Tips. Desalted samples were dissolved in $30\text{ }\mu\text{L}$ of dissolution buffer and labeled with two different iTRAQ reagents at room temperature for 1 h and quenched by Milli-Q water. Sample labeling was as follows: high-risk breast cancer tissue samples with 114 tag and low-risk breast cancer tissue samples with 115 tag. Labeled samples were mixed and dried by a Speed-Vac concentrator, dissolved in $100\text{ }\mu\text{L}$ of 2% acetonitrile (ACN), 0.1% formic acid (TFA), and desalted with C18 stage Tips.

Separation with Strong Cation Exchange Chromatography (SCX)

The tryptic peptide sample was fractionated using a HPLC system (Shimadzu prominence UFLC) fitted with a SCX column ($50\text{ mm} \times 2.1\text{ mm}$, $5\text{ }\mu\text{m}$, $300\text{ }\text{\AA}$, ZORBAX 300SCX, Agilent technology). The mobile phases consisted of (A); 25% ACN with $10\text{ mM KH}_2\text{PO}_4$ (pH 3.0) and (B); (A) containing 1 M KCl . The mixed sample was separated at a flow rate of $200\text{ }\mu\text{L}/\text{min}$ using a four-step linear gradient; 0% B for 30 min, 0 to 10% B

in 15 min, 10 to 25% B in 10 min, 25 to 40% B in 5 min, and 40 to 100% B in 5 min, and 100% B in 10 min.

NanoLC–MS/MS

NanoLC–MS/MS analysis was conducted by an LTQ-Orbitrap Velos mass spectrometer (Thermo Fisher Scientific, Bremen, Germany) equipped with a nanoLC interface (AMR, Tokyo, Japan), a nanoHPLC system (Michrom Paradigm MS2), and an HTC-PAL autosampler (CTC, Analytics, Zwingen, Switzerland). L-column2 C18 particles ($3\text{ }\mu\text{m}$) (Chemicals Evaluation and Research Institute (CERI), Japan) were packed into a self-pulled needle ($200\text{ mm length} \times 100\text{ }\mu\text{m}$ inner diameter) using a Nanobaume capillary column packer (Western Fluids Engineering). Mobile phases consisted of (A) 0.1% FA and 2% ACN and (B) 0.1% FA and 90% ACN. SCX-fractionated peptides dissolved in 2% ACN and 0.1% TFA were loaded onto a trap column ($0.3 \times 5\text{ mm}$, L-column ODS; CERI). The nanoLC gradient was delivered at $500\text{ nL}/\text{min}$ and consisted of a linear gradient of mobile phase B developed from 5 to 30% B in 135 min. A spray voltage of 2000 V was applied.

Data Acquisition with LTQ-Orbitrap Velos

Full MS scans were performed in the orbitrap mass analyzer of LTQ-Orbitrap Velos (scan range $350\text{--}1500\text{ m/z}$, with 30K fwhm resolution at 400 m/z). In MS scans, the ten most intense precursor ions were selected for MS/MS scans of LTQ-Orbitrap Velos respectively, in which a dynamic exclusion option was implemented with a repeat count of one and exclusion duration of 60 s. This was followed by collision-induced dissociation (CID) MS/MS scans of selected ions performed in the linear ion trap mass analyzer, and further followed by higher energy collision-induced dissociation (HCD) MS/MS scans of the same precursor ions performed in the orbitrap mass analyzer with 7500 fwhm resolution at 400 m/z . The values of automated gain control (AGC) were set to $1.00 \times 10^{+06}$ for full MS, $1.00 \times 10^{+04}$ for CID MS/MS, and $5.00 \times 10^{+04}$ for HCD MS/MS. Normalized collision energy values were set to 35% for CID and 50% for HCD. CID, also known as collision-activated dissociation, is performed in the linear ion trap. It is able to increase the number of peptide identifications, and, thus, is applied to obtain peptide sequence information. HCD is performed in the C-trap of the LTQ Orbitrap and is a useful tool for elucidating the structure of small molecules, metabolites, peptides, and PTM peptides, and for de novo sequencing of peptides. It allows quantitative information to be obtained from iTRAQ ions in the lower mass area. By analyzing the sample using a combination of CID with HCD, we are able to obtain the best conditions for both peptide sequencing and iTRAQ quantitation.

Identification and Quantification of Membrane Proteins

CID and HCD raw spectra were extracted and searched separately against UniProtKB/Swiss-Prot (release-2010_05) containing 20 295 sequences of *Homo sapiens* using Proteome Discoverer (Thermo Fisher Scientific, Beta Version 1.3) and Mascot v2.3.1. Search parameters included trypsin as the enzyme with one missed cleavage allowed; Carbamidomethylation at cysteine and iTRAQ labeling at lysine and the N-terminal residue were set as fixed modifications while oxidation at methionine and iTRAQ labeling at tyrosine were set as variable modifications. Precursor mass tolerance was set to 7 ppm and a fragment mass tolerance was set to 0.6 Da for CID and 0.01 Da for HCD. Protein identification required at least one unique peptide and quantification required at least two peptides. FDR was calculated

by enabling peptide sequence analysis using Percolator. High confidence peptide identification was obtained by setting a target FDR threshold of $\leq 1.0\%$ at the peptide level. The mass spectrometry proteomics data have been deposited to the ProteomeXchange Consortium (<http://proteomecentral.proteomexchange.org>) via the PRIDE partner repository (<http://www.ebi.ac.uk/pride/>) with the data set identifier PXD000066.

Bioinformatics Analysis

The subcellular locations of identified proteins were annotated by DAVID Bioinformatics Resources 6.7, available at <http://david.abcc.ncifcrf.gov/home.jsp>.¹⁴ The chromosomal locations and missing protein analysis of identified proteins were elucidated by neXtProt, available at <http://www.nextprot.org/db/>. The function of identified missing membrane proteins was elucidated by the Ingenuity system, available at www.ingenuity.com.

RESULTS

C-HPP is collecting protein data identified by the chromosome-independent shotgun approach and then sharing this data

Table 1. Comparison of the Number of Identified Membrane Protein with Our Result and Previously Reported Results

	Muraoka et al.	Polisetty et al. ⁷	Han et al. ¹⁷
Protein identified	7092	1834	1482
Membrane protein	3282	1027	642

according to the chromosome number to ensure a complete parts list.¹⁵ In this study, we integrated membrane proteomic analysis data from human breast cancer tissues and analyzed with Proteome Discoverer and DAVID Bioinformatics Resources and

characterized them on a chromosome-by-chromosome basis using the neXtProt database.

A total of 7092 unique proteins were identified with high confidence. A list of proteins and peptides are presented in Supplementary Tables 1 and 2, Supporting Information. Identified unique proteins were examined with respect to subcellular localization using Gene Ontology annotation analysis in DAVID Bioinformatics Resources. It revealed that 3282 (46%) were annotated to membrane proteins (Supplementary Table 3, Supporting Information), 692 (10%) proteins were extracellular space, 4030 (57%) proteins were cytoplasm proteins, and 1782 (25%) proteins were nucleus proteins by GO analysis. As shown in Table 1, this number of identified membrane proteins is much greater than previously reported.

To generate a chromosome-based membrane protein list, the identified 3282 membrane proteins were examined with respect to chromosomal location using the neXtProt database. The chromosomal distribution of protein-coding genes in the neXtProt database, identifying total proteins, and membrane proteins are shown in Figure 1. The neXtProt database annotates 7326 proteins as membrane proteins in the 20 859 protein-coding genes, and surprisingly, we identified 45% of them in this study (Figure 2). These results support the effectiveness of the method to solubilize and digest integral membrane proteins, allowing large-scale detection and identification of this protein class with no bias against membrane proteins.

A primary goal of the C-HPP is to identify and characterize proteins that currently lack MS evidence and are referred to as "missing proteins". Thus, we examined how many missing proteins were identified in this study. We compared our membrane protein list with a list of no proteomic proteins (missing proteins) in the neXtProt database. Surprisingly, 851 membrane missing proteins (22.7%) were identified in this study (Figure 3 and Supplementary Table 3, Supporting Information).

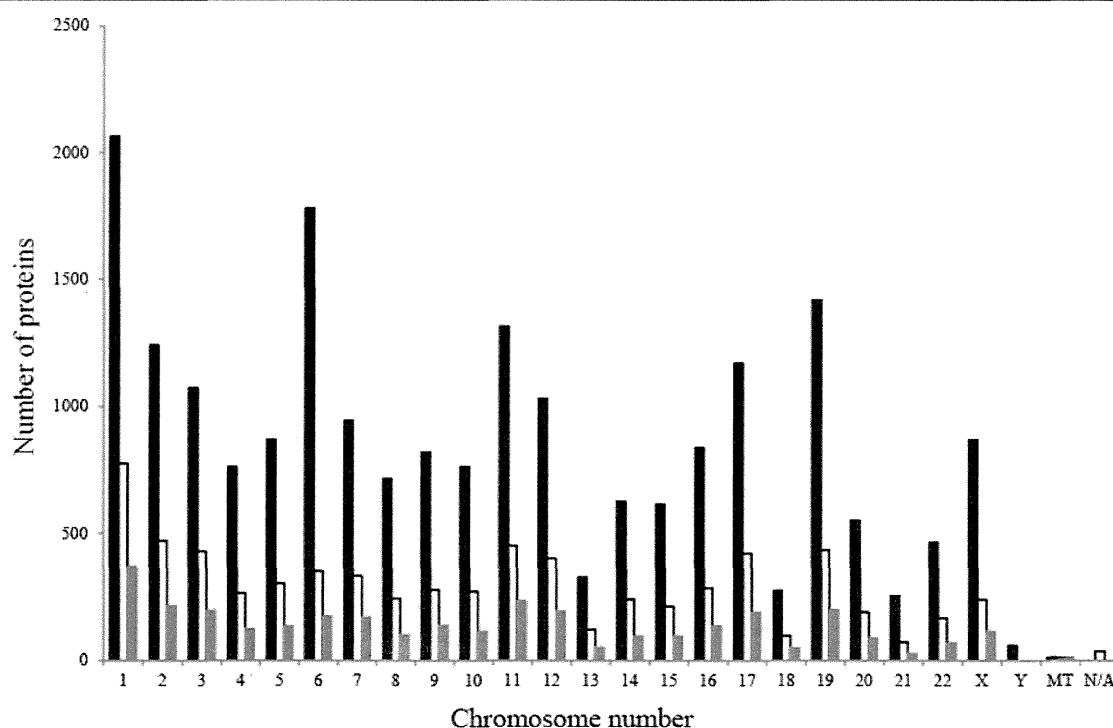


Figure 1. Distribution of identified total and membrane proteins on a whole-chromosomal location. Black bar, neXtProt database proteins; white bar, identified total proteins; gray bar, identified membrane proteins; N/A, no protein in the neXtProt database; MT, mitochondria.

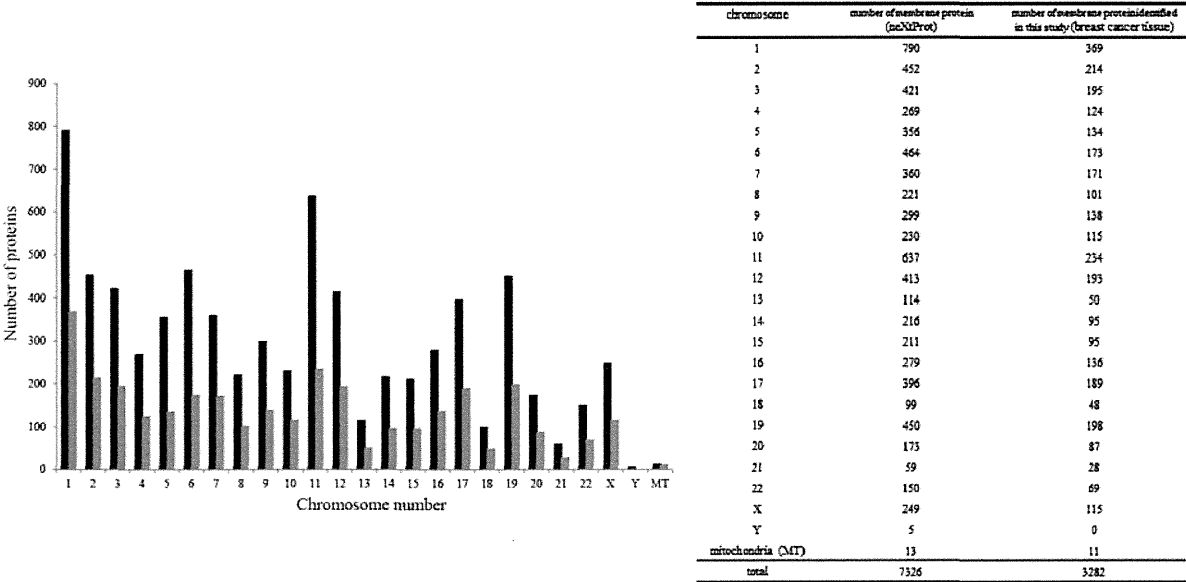


Figure 2. Comparison of chromosome-based membrane proteins annotated by the neXtProt database with identified membrane proteins in this study. (Left) Black bar, neXtProt database membrane proteins; gray bar, identified membrane proteins. (Right) Number of identified and neXtProt database membrane proteins on a whole-chromosomal location.

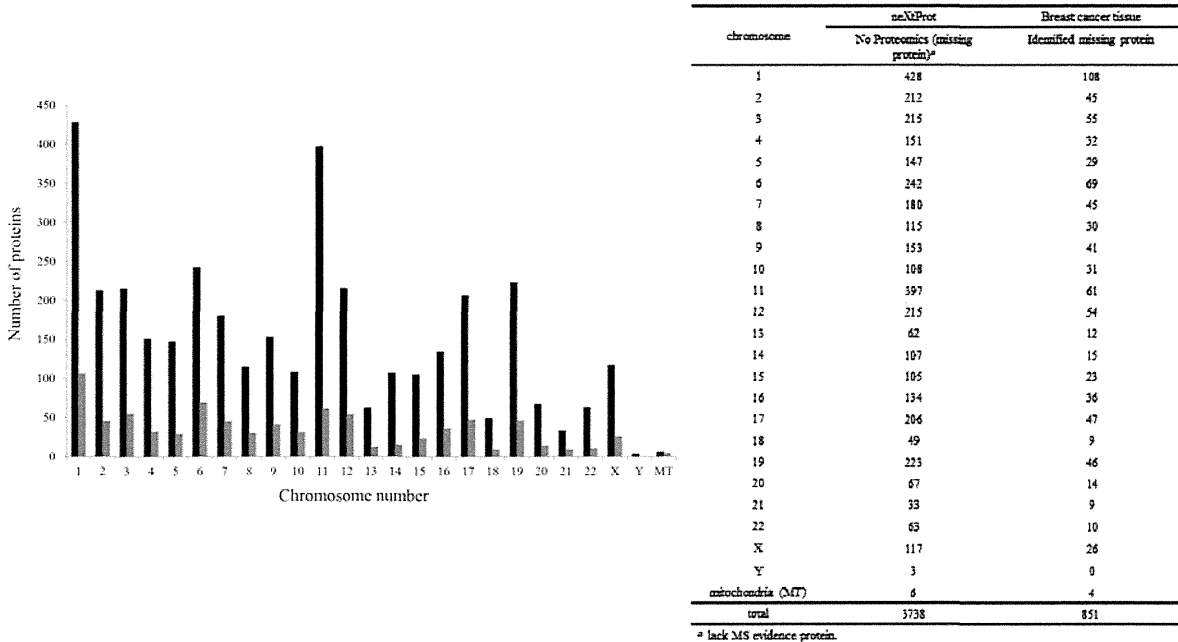


Figure 3. Identified and missing membrane proteins on a whole-chromosomal location. (Left) Black bar, no proteomics (missing membrane proteins) in neXtProt; gray bar, identified missing membrane proteins in this study. (Right) Number of identified and missing membrane proteins on a whole chromosome.

Lipid metabolism, small molecule biochemistry, cell-to-cell signaling and interaction, hematological system development and function, and immune cell trafficking were the major molecular and cellular processes identified by IPA (Figure 4). These results indicate that our in-depth membrane proteomic study of breast cancer tissue samples was able to identify and characterize a number of low-abundance missing proteins.

DISCUSSION

The objective of C-HPP is to map and annotate all protein-coding genes on each human chromosome, especially so-called “missing proteins”, which only have transcriptomic evidence and

a predicted sequence. To accomplish this, deep profiling for low-abundance proteins and subcellular proteins such as membrane proteins is needed. In this study, we performed an in-depth membrane proteomic study of breast cancer tissues. A total of 7092 proteins were identified, of which 3282 proteins were annotated as membrane proteins by Gene Ontology analysis. Furthermore, we could identify not only nearly 50% of the membrane proteins mapped on the whole chromosome but also 851 proteins among the 3738 missing membrane proteins. Several previously published reports have described membrane proteome analysis.^{8,16,17} Polisetty and co-workers recently performed a large-scale proteomic study utilizing shotgun

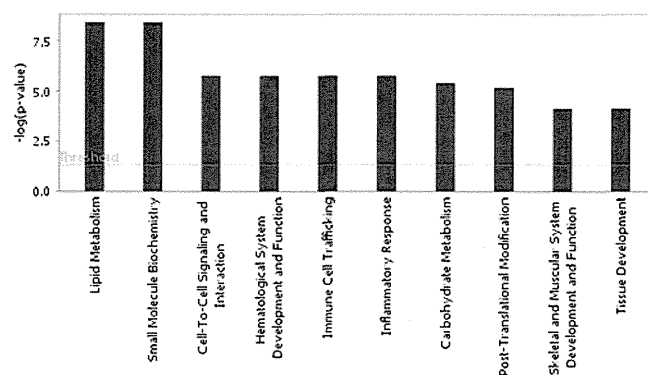


Figure 4. Bar chart indicating the cellular function of proteins found with missing proteins determined using Ingenuity software.

technology and identified 1834 distinct proteins from membrane fractions of glioblastoma multiforme patient specimens, with 56% of them (1027) being annotated as membrane proteins.⁷ In this study, we identified a total of 7092 proteins in the membrane fraction; with 46% of them (3282) being known membrane proteins associated with major cellular processes. This number of membrane proteins is much greater than those previously reported. Moreover, we were able to identify a number of missing proteins that currently lack MS evidence. This is probably due to utilization of the PTS method-based isolation of membrane proteins and SCX fractionation before LC–MS/MS analysis. Efficient isolation and solubilization of membrane proteins can be achieved with PTS by allowing the use of a high detergent concentration while avoiding interference with tryptic digestion before LC–MS/MS analysis.¹⁸ SCX prefractionation is also able to improve the number of proteins identified by reducing the complexity of clinical samples and consequently avoiding ion suppression. We have succeeded in large scale identification of membrane proteins and phosphoproteins using the above technique.^{12,19}

In conclusion, a subcellular fractionation of membrane proteins would improve low-abundance proteome coverage for identification of missing proteins. Our in-depth membrane proteomic studies of human cancer tissue will greatly contribute to the progression of the C-HPP.

■ ASSOCIATED CONTENT

⑤ Supporting Information

Supplementary figure and tables. This material is available free of charge via the Internet at <http://pubs.acs.org>.

■ AUTHOR INFORMATION

Corresponding Author

*Laboratory of Proteome Research, National Institute of Biomedical Innovation 7-6-8 Saito-Asagi, Ibaraki City, Osaka 567-0085, Japan. Tel.: +81-72-641-9862. Fax: +81-72-641-9861. E-mail: tomonaga@nibio.go.jp.

Notes

The authors declare no competing financial interest.

■ ACKNOWLEDGMENTS

This work was supported by Grants-in-Aid, Research on Biological Markers for New Drug Development H20-0005 to T.T from the Ministry of Health, Labour, and Welfare of Japan. This work was supported by Grants-in-Aid 21390354 to T.T and

22800095 to S.M from the Ministry of Education, Science, Sports, and Culture of Japan.

■ ABBREVIATIONS

C-HPP, The Chromosome-Centric Human Proteome Project; PTS, phase-transfer surfactants; CID, collision-induced dissociation; HCD, higher energy collision-induced dissociation; LC–MS/MS, Liquid chromatography–tandem mass spectrometry; LTQ, linear ion trap; fwhm, Full Width at Half Maximum; FDR, false discovery rate

■ REFERENCES

- (1) Collins, F. S.; Green, E. D.; Guttmacher, A. E.; Guyer, M. S. A vision for the future of genomics research. *Nature* **2003**, *422* (6934), 835–47.
- (2) Lander, E. S.; Linton, L. M.; Birren, B.; Nusbaum, C.; Zody, M. C.; Baldwin, J.; Devon, K.; Dewar, K.; Doyle, M.; FitzHugh, W.; Funke, R.; Gage, D.; Harris, K.; Heaford, A.; Howland, J.; Kann, L.; Lehoczy, J.; LeVine, R.; McEwan, P.; McKernan, K.; Meldrum, J.; Mesirov, J. P.; Miranda, C.; Morris, W.; Naylor, J.; Raymond, C.; Rosetti, M.; Santos, R.; Sheridan, A.; Sougnez, C.; Stange-Thomann, N.; Stojanovic, N.; Subramanian, A.; Wyman, D.; Rogers, J.; Sulston, J.; Ainscough, R.; Beck, S.; Bentley, D.; Burton, J.; Clee, C.; Carter, N.; Coulson, A.; Deadman, R.; Deloukas, P.; Dunham, A.; Dunham, I.; Durbin, R.; French, L.; Grafham, D.; Gregory, S.; Hubbard, T.; Humphray, S.; Hunt, A.; Jones, M.; Lloyd, C.; McMurray, A.; Matthews, L.; Mercer, S.; Milne, S.; Mullikin, J. C.; Mungall, A.; Plumb, R.; Ross, M.; Shownkeen, R.; Sims, S.; Waterston, R. H.; Wilson, R. K.; Hillier, L. W.; McPherson, J. D.; Marra, M. A.; Mardis, E. R.; Fulton, L. A.; Chinwalla, A. T.; Pepin, K. H.; Gish, W. R.; Chissoe, S. L.; Wendl, M. C.; Delehaunty, K. D.; Miner, T. L.; Delehaunty, A.; Kramer, J. B.; Cook, L. L.; Fulton, R. S.; Johnson, D. L.; Minx, P. J.; Clifton, S. W.; Branscomb, E.; Predki, P.; Richardson, P.; Wenning, S.; Slezak, T.; Doggett, N.; Cheng, J. F.; Olsen, A.; Lucas, S.; Elkin, C.; Uberbacher, E.; Frazier, M.; Gibbs, R. A.; Muzny, D. M.; Scherer, S. E.; Bouck, J. B.; Sodergren, E. J.; Worley, K. C.; Rives, C. M.; Gorrell, J. H.; Metzker, M. L.; Naylor, S. L.; Kucherlapati, R. S.; Nelson, D. L.; Weinstock, G. M.; Sakaki, Y.; Fujiyama, A.; Hattori, M.; Yada, T.; Toyoda, A.; Itoh, T.; Kawagoe, C.; Watanabe, H.; Totoki, Y.; Taylor, T.; Weissenbach, J.; Heilig, R.; Saurin, W.; Artiguenave, F.; Brottier, P.; Bruls, T.; Pelletier, E.; Robert, C.; Wincker, P.; Smith, D. R.; Doucette-Stamm, L.; Rubinfeld, M.; Weinstock, K.; Lee, H. M.; Dubois, J.; Rosenthal, A.; Platzer, M.; Nyakatura, G.; Taudien, S.; Rump, A.; Yang, H.; Yu, J.; Wang, J.; Huang, G.; Gu, J.; Hood, L.; Rowen, L.; Madan, A.; Qin, S.; Davis, R. W.; Federspiel, N. A.; Abola, A. P.; Proctor, M. J.; Myers, R. M.; Schmutz, J.; Dickson, M.; Grimwood, J.; Cox, D. R.; Olson, M. V.; Kaul, R.; Shimizu, N.; Kawasaki, K.; Minoshima, S.; Evans, G. A.; Athanasiou, M.; Schultz, R.; Roe, B. A.; Chen, F.; Pan, H.; Ramser, J.; Lehrach, H.; Reinhardt, R.; McCombie, W. R.; de la Bastide, M.; Dedhia, N.; Blocker, H.; Hornischer, K.; Nordsiek, G.; Agarwala, R.; Aravind, L.; Bailey, J. A.; Bateman, A.; Batzoglou, S.; Birney, E.; Bork, P.; Brown, D. G.; Burge, C. B.; Cerutti, L.; Chen, H. C.; Church, D.; Clamp, M.; Copley, R. R.; Doerks, T.; Eddy, S. R.; Eichler, E. E.; Furey, T. S.; Galagan, J.; Gilbert, J. G.; Harmon, C.; Hayashizaki, Y.; Haussler, D.; Hermjakob, H.; Hokamp, K.; Jang, W.; Johnson, L. S.; Jones, T. A.; Kasif, S.; Kasprzyk, A.; Kennedy, S.; Kent, W. J.; Kitts, P.; Koonin, E. V.; Korf, I.; Kulp, D.; Lancet, D.; Lowe, T. M.; McLysaght, A.; Mikkelsen, T.; Moran, J. V.; Mulder, N.; Pollara, V. J.; Ponting, C. P.; Schuler, G.; Schultz, J.; Slater, G.; Smit, A. F.; Stupka, E.; Szustakowski, J.; Thierry-Mieg, D.; Thierry-Mieg, J.; Wagner, L.; Wallis, J.; Wheeler, R.; Williams, A.; Wolf, Y. I.; Wolfe, K. H.; Yang, S. P.; Yeh, R. F.; Collins, F.; Guyer, M. S.; Peterson, J.; Felsenfeld, A.; Wetterstrand, K. A.; Patrino, A.; Morgan, M. J.; de Jong, P.; Catanese, J. J.; Osoegawa, K.; Shizuya, H.; Choi, S.; Chen, Y. J. Initial sequencing and analysis of the human genome. *Nature* **2001**, *409* (6822), 860–921.
- (3) Legrain, P.; Aebersold, R.; Archakov, A.; Bairoch, A.; Bala, K.; Beretta, L.; Bergeron, J.; Borchers, C. H.; Corthals, G. L.; Costello, C. E.; Deutsch, E. W.; Domon, B.; Hancock, W.; He, F.; Hochstrasser, D.; Marko-Varga, G.; Salekdeh, G. H.; Sechi, S.; Snyder, M.; Srivastava, S.; Uhlen, M.; Wu, C. H.; Yamamoto, T.; Paik, Y. K.; Omenn, G. S. The

human proteome project: current state and future direction. *Mol. Cell. Proteomics* **2011**, *10* (7), M111 009993.

(4) Hancock, W.; Omenn, G.; Legrain, P.; Paik, Y. K. Proteomics, human proteome project, and chromosomes. *J. Proteome Res.* **2011**, *10* (1), 210.

(5) Paik, Y. K.; Jeong, S. K.; Omenn, G. S.; Uhlen, M.; Hanash, S.; Cho, S. Y.; Lee, H. J.; Na, K.; Choi, E. Y.; Yan, F.; Zhang, F.; Zhang, Y.; Snyder, M.; Cheng, Y.; Chen, R.; Marko-Varga, G.; Deutsch, E. W.; Kim, H.; Kwon, J. Y.; Aebersold, R.; Bairoch, A.; Taylor, A. D.; Kim, K. Y.; Lee, E. Y.; Hochstrasser, D.; Legrain, P.; Hancock, W. S. The Chromosome-Centric Human Proteome Project for cataloging proteins encoded in the genome. *Nat. Biotechnol.* **2012**, *30* (3), 221–3.

(6) Lane, L.; Argoud-Puy, G.; Britan, A.; Cusin, I.; Duek, P. D.; Evalet, O.; Gateau, A.; Gaudet, P.; Gleizes, A.; Masselot, A.; Zwahlen, C.; Bairoch, A. neXtProt: a knowledge platform for human proteins. *Nucleic Acids Res.* **2012**, *40* (Database issue), D76–83.

(7) Polisetty, R. V.; Gautam, P.; Sharma, R.; Harsha, H. C.; Nair, S. C.; Gupta, M. K.; Uppin, M. S.; Challa, S.; Puligopu, A. K.; Ankathi, P.; Purohit, A. K.; Chandak, G. R.; Pandey, A.; Sirdeshmukh, R. LC-MS/MS analysis of differentially expressed glioblastoma membrane proteome reveals altered calcium signaling and other protein groups of regulatory functions. *Mol. Cell. Proteomics* **2012**, *11* (6), M111 013565.

(8) Josic, D.; Clifton, J. G. Mammalian plasma membrane proteomics. *Proteomics* **2007**, *7* (16), 3010–29.

(9) Russell, W. K.; Park, Z. Y.; Russell, D. H. Proteolysis in mixed organic-aqueous solvent systems: applications for peptide mass mapping using mass spectrometry. *Anal. Chem.* **2001**, *73* (11), 2682–5.

(10) Masuda, T.; Tomita, M.; Ishihama, Y. Phase transfer surfactant-aided trypsin digestion for membrane proteome analysis. *J. Proteome Res.* **2008**, *7* (2), 731–40.

(11) Wisniewski, J. R.; Zougman, A.; Nagaraj, N.; Mann, M. Universal sample preparation method for proteome analysis. *Nat. Methods* **2009**, *6* (5), 359–62.

(12) Muraoka, S.; Kume, H.; Watanabe, S.; Adachi, J.; Kuwano, M.; Sato, M.; Kawasaki, N.; Kodera, Y.; Ishitobi, M.; Inaji, H.; Miyamoto, Y.; Kato, K.; Tomonaga, T. Strategy for SRM-based Verification of Biomarker Candidates Discovered by iTRAQ Method in Limited Breast Cancer Tissue Samples. *J. Proteome Res.* **2012**, *11* (8), 4201–10.

(13) Rappsilber, J.; Mann, M.; Ishihama, Y. Protocol for micro-purification, enrichment, pre-fractionation and storage of peptides for proteomics using StageTips. *Nat. Protoc.* **2007**, *2* (8), 1896–906.

(14) Huang da, W.; Sherman, B. T.; Lempicki, R. A. Systematic and integrative analysis of large gene lists using DAVID bioinformatics resources. *Nat. Protoc.* **2009**, *4* (1), 44–57.

(15) Paik, Y. K.; Omenn, G. S.; Uhlen, M.; Hanash, S.; Marko-Varga, G.; Aebersold, R.; Bairoch, A.; Yamamoto, T.; Legrain, P.; Lee, H. J.; Na, K.; Jeong, S. K.; He, F.; Binz, P. A.; Nishimura, T.; Keown, P.; Baker, M. S.; Yoo, J. S.; Garin, J.; Archakov, A.; Bergeron, J.; Salekdeh, G. H.; Hancock, W. S. Standard guidelines for the chromosome-centric human proteome project. *J. Proteome Res.* **2012**, *11* (4), 2005–13.

(16) Chen, J. S.; Chen, K. T.; Fan, C. W.; Han, C. L.; Chen, Y. J.; Yu, J. S.; Chang, Y. S.; Chien, C. W.; Wu, C. P.; Hung, R. P.; Chan, E. C. Comparison of membrane fraction proteomic profiles of normal and cancerous human colorectal tissues with gel-assisted digestion and iTRAQ labeling mass spectrometry. *FEBS J.* **2010**, *277* (14), 3028–38.

(17) Han, C. L.; Chen, J. S.; Chan, E. C.; Wu, C. P.; Yu, K. H.; Chen, K. T.; Tsou, C. C.; Tsai, C. F.; Chien, C. W.; Kuo, Y. B.; Lin, P. Y.; Yu, J. S.; Hsueh, C.; Chen, M. C.; Chan, C. C.; Chang, Y. S.; Chen, Y. J. An informatics-assisted label-free approach for personalized tissue membrane proteomics: case study on colorectal cancer. *Mol. Cell. Proteomics* **2011**, *10* (4), M110 003087.

(18) Iwasaki, M.; Masuda, T.; Tomita, M.; Ishihama, Y. Chemical cleavage-assisted tryptic digestion for membrane proteome analysis. *J. Proteome Res.* **2009**, *8* (6), 3169–75.

(19) Narumi, R.; Murakami, T.; Kuga, T.; Adachi, J.; Shiromizu, T.; Muraoka, S.; Kume, H.; Kodera, Y.; Matsumoto, M.; Nakayama, K.; Miyamoto, Y.; Ishitobi, M.; Inaji, H.; Kato, K.; Tomonaga, T. A strategy

for large-scale phosphoproteomics and SRM-based validation of human breast cancer tissue samples. *J. Proteome Res.* **2012**, *11* (11), 5311–22.