

Global Survey of Human T Leukemic Cells by Integrating Proteomics and Transcriptomics Profiling*[§]

Linfeng Wu‡, Sun-Il Hwang‡§, Karim Rezaul‡§, Long J. Lu§¶||, Viveka Mayya‡, Mark Gerstein¶, Jimmy K. Eng**, Deborah H. Lundgren‡, and David K. Han‡ ‡‡

A global protein survey is needed to gain systems-level insights into mammalian cell signaling and information flow. Human Jurkat T leukemic cells are one of the most important model systems for T cell signaling study, but no comprehensive proteomics survey has been carried out in this cell type. In the present study we combined subcellular fractionation, multiple protein enrichment methods, and replicate tandem mass spectrometry analyses to determine the protein expression pattern in a single Jurkat cell type. The proteome dataset was evaluated by comparison with the genome-wide mRNA expression pattern in the same cell type. A total of 5381 proteins were identified by mass spectrometry with high confidence. Rigorous comparison of RNA and protein expression afforded removal of the false positive identifications and redundant entries but rescued the proteins identified by a single high scoring peptide, resulting in the final identification of 6471 unique gene products among which 98% of the corresponding transcripts were detected with high probability. Using hierarchical clustering of the protein expression patterns in five subcellular fractions (cytosol, light membrane, heavy membrane, mitochondria, and nuclei), the primary subcellular localization of 2241 proteins was assigned with high confidence including 792 previously uncharacterized proteins. This proteome landscape can serve as a useful platform for systems-level understanding of organelle composition and cellular functions in human T cells. *Molecular & Cellular Proteomics* 6: 1343–1353, 2007.

An important goal in functional genomics is to globally profile protein expression and localization in biological systems. Many studies have utilized genome-wide cDNA or oligonucleotide microarrays to measure mRNA expression level,

deducing the corresponding protein expression (1, 2). However, despite the obvious dependence of protein synthesis on mRNA, many studies have reported that more than half of the total transcripts are non-coding RNA (3–5). In addition, quantitative measurements show only moderate or even poor correlation between protein and mRNA expression level due to different translational efficiency and post-translational turnover (6–8). Furthermore subcellular localization of proteins cannot be accurately predicted based on mRNA expression. Therefore, biological systems ultimately need to be explained at the level of proteins.

The completed draft sequences of the human genome (9, 10) and several other organisms (11, 12) combined with mass spectrometry have made large scale proteomics feasible (13, 14). However, due to the huge diversity and dynamic range of expressed proteins, especially in human cells, identification of all or most of the expressed proteins in cells has remained one of the greatest challenges (15). Although many proteome-scale studies on different cells, tissues, and subcellular organelles have been reported, no comprehensive analysis of a single human cell type has been carried out to date. In this study, we performed a comprehensive survey of a human Jurkat T leukemic cell line by combining proteomics and transcriptomics profiling. Human Jurkat T leukemic cells are one of the most popular model systems for studying signal transduction because many key advances in the field of T cell receptor signaling were made using Jurkat T cells (16). Moreover this cell type is also used for studying other biological phenomena such as apoptosis and cell engulfment (17, 18). Therefore, a global survey of human Jurkat T cells can serve as a platform for many in-depth characterizations of cellular function and signaling transduction. Moreover in contrast to the recent survey of organ and organelle protein expression in mouse (8, 19), our study was carried out in a single cell type, making it more suitable for studying protein network and signaling flow within cells.

EXPERIMENTAL PROCEDURES

Whole Cell Lysate Preparation—Human Jurkat A3 T leukemic cells from American Type Culture Collection (Manassas, VA) were used for this study. Jurkat T cells were grown to a maximal density of $0.5\text{--}0.8 \times 10^6$ cells/ml and then collected by centrifugation at $400 \times g$ for 10 min at 4 °C. The cell pellets were washed twice with ice-cold PBS. To obtain whole cell lysates, cells were resuspended in lysis buffer (50

From the ‡Department of Cell Biology and Center for Vascular Biology, School of Medicine, University of Connecticut, Farmington, Connecticut 06030, ¶Department of Molecular Biophysics and Biochemistry, Yale University, New Haven, Connecticut 06520, ||Division of Biomedical Informatics, Cincinnati Children's Hospital Medical Center, Cincinnati, Ohio 45229, and **Fred Hutchinson Cancer Research Center, Seattle, Washington 98195

Received, January 19, 2007

Published, MCP Papers in Press, May 21, 2007, DOI 10.1074/mcp.M700017-MCP200

mM HEPES, pH 7.5, 100 mM NaCl, 1 mM EDTA, 1% Tween 20, and a mixture of protease inhibitors (Roche Diagnosis GmbH) on ice for 30 min and then centrifuged at $12,000 \times g$ for 20 min at 4 °C. The supernatant (designated whole cell lysate) was collected and stored at -80 °C for later analysis.

Subcellular Fractionation—Subcellular fractionation was carried out as described below. One volume of Jurkat cell pellet was incubated in 5 volumes of hypotonic buffer (buffer A: 20 mM HEPES, pH 7.5, 10 mM KCl, 1.5 mM MgCl₂, 1 mM EDTA, 1 mM EGTA, 1 mM DTT, and a mixture of protease inhibitors) for 2 min and then mixed with an equal volume of buffer A supplemented with 0.5 M sucrose resulting in a 0.25 M sucrose isolation buffer. After 10 min of incubation on ice, cells were homogenized with a glass Dounce homogenizer until ~50% of the cells became trypan blue-positive. The homogenates were centrifuged at $650 \times g$ for 10 min at 4 °C. The postnuclear supernatants were incubated on ice for other subcellular fraction preparation. The pellets were further homogenized in buffer A until ~95% of the cells became trypan blue-positive. The nuclear pellets were isolated by centrifugation at $650 \times g$ for 10 min at 4 °C and then rinsed with isolation buffer once. The nuclear pellets were gently resuspended in 2.5 ml of buffer A supplemented with 1.28 M sucrose, layered over 5 ml of buffer A supplemented with 2.3 M sucrose, and then centrifuged at $60,000 \times g$ for 90 min at 4 °C. The purified nuclear pellets were rinsed with buffer A, centrifuged at $12,000 \times g$ for 10 min, and stored at -80 °C for subsequent experiments. The postnuclear supernatants were used to isolate other subcellular fractions including cytosol, heavy membrane, light membrane, and mitochondria as described previously (20).

To isolate plasma membrane fraction, 1 volume of cells ($\sim 2 \times 10^9$) was resuspended in 5 volumes of hypotonic buffer A and homogenized with a glass Dounce homogenizer 15 times. The homogenates were centrifuged at $650 \times g$ for 10 min at 4 °C, and the pellet was removed. The supernatant was further centrifuged at $100,000 \times g$ for 1 h at 4 °C. The pellet was resuspended in buffer A supplemented with 0.25 M sucrose, homogenized with a Dounce homogenizer 10–15 times, layered on 6.5 ml of sucrose buffer (buffer A supplemented with 38% sucrose), and then centrifuged at $200,000 \times g$ for 2 h. The membranes were collected from the phase between the 0.25 M and 38% sucrose, diluted in 10 ml of 0.25 M sucrose buffer, and centrifuged at $100,000 \times g$ for 1 h. The pellet was designated as plasma membrane.

To isolate the lipid raft fraction, the PBS-washed cell pellets were resuspended in 1 ml of lysis buffer (25 mM MES, pH 6.5, 150 mM NaCl, 0.1% Triton X-100, 1 mM sodium vanadate, 5 mM EDTA, and a mixture of protease inhibitors). Cells were homogenized with a Dounce homogenizer 20–30 times and then mixed with 1 ml of 80% sucrose in lysis buffer without Triton X-100. The lysates were placed in the bottom of a 14 × 89-mm clear centrifuge tube (Beckman), gently overlaid with 6.5 ml of 30% sucrose and 3.5 ml of 5% sucrose in lysis buffer without Triton X-100, and then centrifuged at $200,000 \times g$ at 4 °C for 18 h (the machine was set at “no brake” condition). The low density membrane raft in the 5% sucrose fraction was collected and designated as lipid raft. The details of cell culture and other enrichment approaches are shown in the supplemental materials.

In-gel Digestion and Nano-LC-MS/MS Analysis—Proteins were digested with trypsin and analyzed as described previously (20) with minor modifications as outlined below. Digested proteins were analyzed using a linear ion trap mass spectrometer (Finnigan LTQ, Thermo Finnigan, San Jose, CA). Samples were loaded onto a 10-cm × 100- μ m capillary C₁₈ reversed-phase column by a microautosampler (Famos, Dionex, Sunnyvale, CA) followed by LC-MS/MS analysis on the LTQ. For stable isotope-free peptide samples, each full MS scan was followed by five MS/MS scans of the five most intense peaks in the MS spectrum with dynamic exclusion enabled. The *m/z* scan range was either 300–1700 or 400–1700 for full mass

range. For stable isotope-labeled peptide samples, which mainly come from subcellular fractions and phosphoprotein enrichment, each full MS scan was followed by one MS/MS scan of the most intense peak in the MS spectrum with dynamic exclusion enabled.

Database Searching and Data Processing—All the mass spectrometry raw files were converted to .dat files using Xcalibur software (version 1.4 SR1) that were then converted to mzXML using the conversion software dat2xml (Institute for Systems Biology). Peak lists were generated automatically without smoothing and deisotoping, and charge states were assigned based on the MS and MS/MS scans as described previously (21). The minimum signal count for full MS is 1000. All the mzXML files were searched against a local copy of the non-redundant human protein database (56,709 entries, November 30, 2004 release version) from the NCI, National Institutes of Health, Advanced Biomedical Computing Center using the SEQUEST algorithm (SEQUEST-PVM version 27 (revision 0)) (21). SEQUEST parameters were as follows: all the filtering thresholds were off; mass tolerance of 1.0 Da for precursor ions and 0.5 Da for fragment ions; full tryptic constraint allowing one missed cleavage; and allowing oxidization (+16 Da) of methionine. If the peptides contained heavy isotope-labeled amino acids, the corresponding amino acid modification was also allowed. The detailed description of labeled amino acids used in each experiment is shown in Supplemental Table S1. The database search results were processed using the INTERACT program (22) and filtered with the following criteria: Xcorr cutoff values of 1.9, 2.2, and 3.7 for 1+, 2+, and 3+ peptides, respectively; Δ Cn cutoff value of 0.1; and partially isotope-labeled peptides were excluded. Proteins identified by at least two distinct peptides in the same experimental fraction were collected for in-depth analysis. To estimate the false positive rate, the datasets were searched against a forward and reversed concatenated human protein database (23). To address the redundancy issue in the list of identified proteins, the peptides filtered by the above Xcorr and Δ Cn values were used to compute peptide probability and protein probability using PeptideProphet (version 1.0) (24) and ProteinProphet (version 2.0) softwares (25), which combine the redundant proteins into a unique protein group and indicate whether the peptide sequences are unique to the corresponding protein group. The searching parameter was set as minimum probability 0.0 to include all the results for the in-depth analysis. In addition, a genome-wide transcript profiling of the same human Jurkat cell type was performed and compared with our proteomics dataset.

Transcript Profiling—The Sentrix Human-6 Expression BeadChip (Illumina, San Diego, CA) that contains 50-mer gene-specific oligonucleotide probes corresponding to >46,000 human transcript variants was used in this study. There are on average 30× redundancy for each transcript per array. Total RNA was isolated from Jurkat cells at log phase using the guanidine thiocyanate method. All the solutions and materials were RNase-free if necessary. Cells (1×10^8) were washed twice with PBS, lysed with 4 ml of GTC (4 M guanidine thiocyanate, 30 mM sodium acetate, 1% 2-mercaptoethanol, pH 7.0), and then homogenized using a 20-gauge syringe needle 20 times. The cell lysate was gently layered on the top of 3 ml of CsCl solution (5.7 M CsCl, 30 mM sodium acetate). Then the sample was centrifuged for 20 h at 27,000 rpm in an SW 41 rotor. The supernatant was removed carefully, and the pellet was dissolved in 200 μ l of H₂O. Total RNA was further purified with phenol:chloroform:isoamyl alcohol (25:24:1, v/v/v) and precipitate with 3 M sodium acetate and ethanol. The final total RNA pellet was lyophilized in a speed vacuum and stored at -80 °C. Poly(A)-enriched mRNA was isolated from total RNA using the Absolute mRNA Purification kit (Stratagene). The qualities of total RNA and poly(A)-enriched mRNA were examined by electrophoresis on a formaldehyde, 1% agarose gel and Northern blot hybridization (Supplemental Fig. S1).

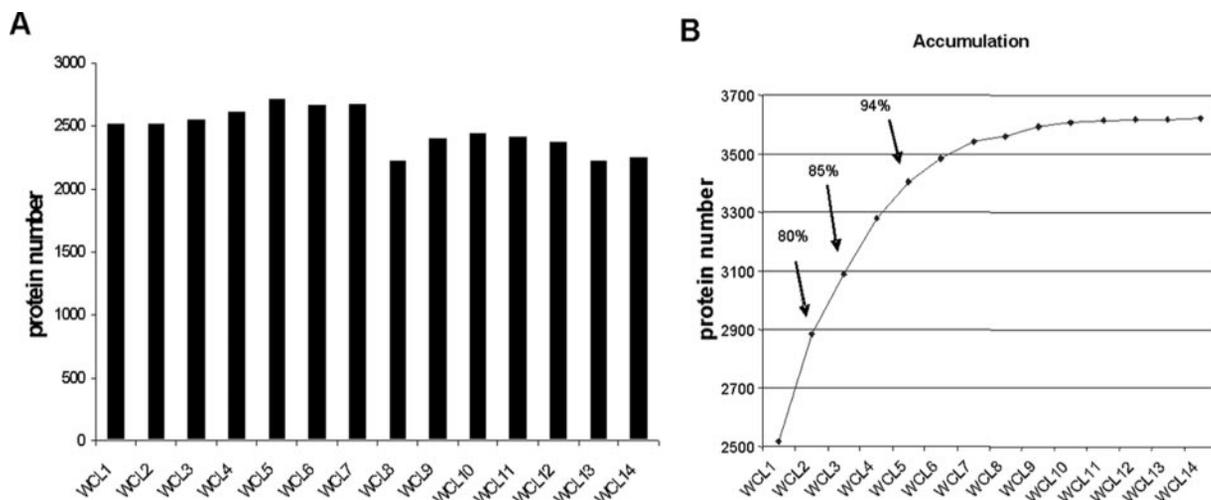


FIG. 1. **Saturated protein identification by replicate analyses.** A, the number of unique proteins identified in each replicate analysis of the human Jurkat whole cell lysates (WCL) is shown. B, the cumulative curve of unique proteins identified by replicate analyses of the whole cell lysates from human Jurkat T leukemic cells is shown. The percentages of proteome coverage after two, three, and five repeats are indicated.

The RNA sample was amplified using the TotalPrep RNA amplification kit (Ambion, Foster City, CA) followed by hybridization, labeling, and scanning of the chips according to the Illumina protocol. The data were extracted, normalized, and analyzed using the Illumina BeadStudio software.

Quantitative Analysis—The semiquantitation of protein abundance was calculated by normalizing the spectral counts of each protein in one fraction relative to the total spectral counts in the corresponding fraction. The normalized profiles were hierarchically clustered based on uncentered correlation with centroid linkage using Cluster 3.0 (26) and visualized using Java TreeView (27).

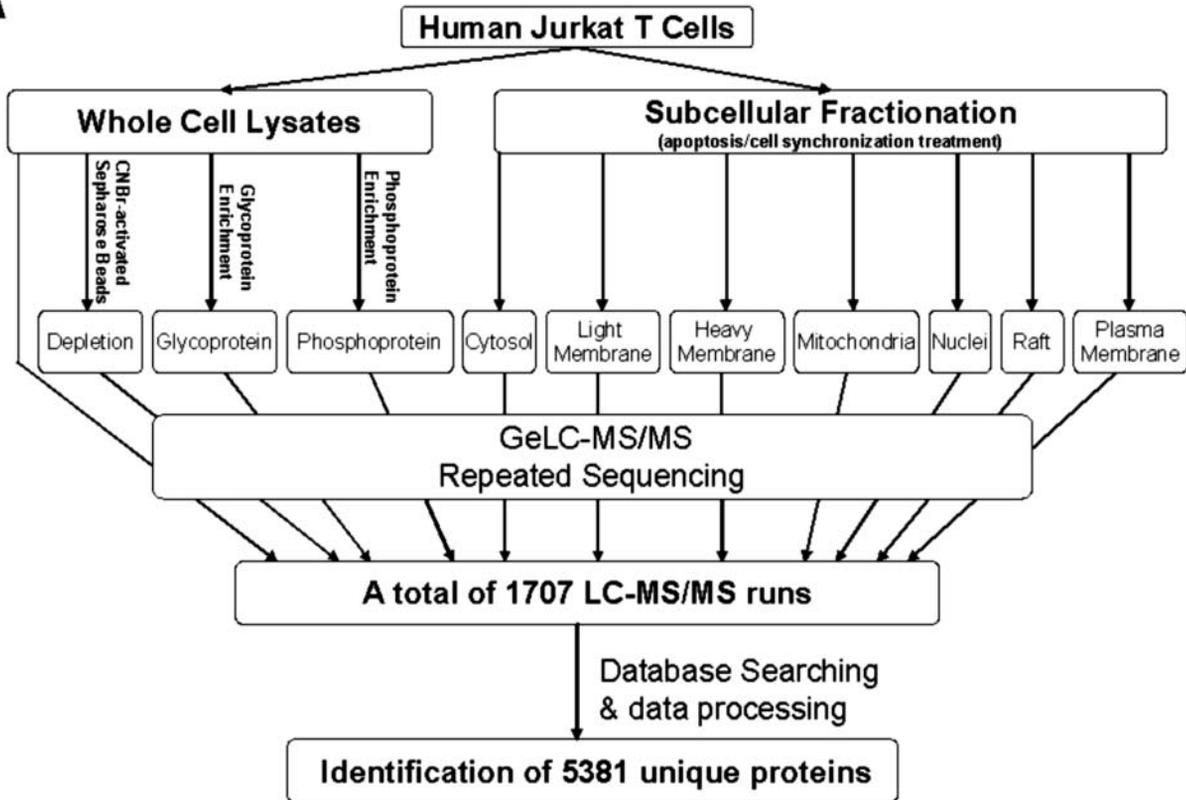
RESULTS

Saturated Protein Identification—To begin to identify the total proteome of human Jurkat T cells, we separated 33 μg of whole cell lysates by one-dimensional gel electrophoresis and cut the whole gel lane into 18 gel slices. The proteins contained in the gel slices were trypsinized, and peptides were extracted as described previously (20). The peptides were then analyzed by LC-MS/MS using the LTQ ion trap mass spectrometer. This process was repeated 14 times using the same amount (33 μg) of proteins from the same whole cell lysate sample. Each replicate analysis identified 2200–2700 proteins (Fig. 1A and Supplemental Table S1). The overlap between each two replicates is about 82% with only 1.8% standard deviation. Differences in protein identification among replicates can be attributed to the complexity of the peptide sample and random sampling during data acquisition by LC-MS/MS. When the data for cumulative total number of unique proteins were analyzed, we found that 3620 proteins were identified from the 14 replicates. We observed that the proteome coverage was enhanced by about 40% when the sample was sequenced 14 times compared with a single analysis (Fig. 1B). Moreover the first five analyses reached about 94% of the total proteins identified by 14 replicate analyses. The protein accumulation curve approaches a slope

of 0 after nine repeats, suggesting that 14 replicate analyses are more than enough to achieve saturation.

Global Proteome Survey of Human Jurkat T Leukemic Cells—Although we saturate protein identification in whole cell lysate by replicate analyses, it is obvious that we did not identify all the proteins. Therefore, we next attempted to detect proteins that were refractory to replicate analysis. We aimed to increase the proteome coverage by reducing the sample complexity using well established and validated fractionation methods (20, 28). First, we subfractionated the cell into seven fractions, including cytosolic, light membrane, heavy membrane, mitochondrial, nuclear, lipid raft, and plasma membrane fractions. Western blot assessments confirmed the appropriate partitioning of several organellar markers across these subcellular fractions, providing a basic confirmation of fraction purity (Supplemental Fig. S2). Second, whole cell lysates from Jurkat cells were used to enrich phosphorylated proteins and validated by Western blotting (Supplemental Fig. S3). Third, we enriched glycosylated proteins by using the lectin wheat germ agglutinin, which preferentially binds *N*-acetylglucosamine (GlcNAc), terminal GlcNAc structures, and sialic acid. The enrichment of the glycosylated proteins from the whole cell lysates was validated by Western analysis (Supplemental Fig. S4). Fourth, we tried to detect previously masked subsets of proteins by differential protein depletion. Because proteins have different binding rates to a particular medium, we hypothesized that the abundance level of proteins in a complex mixture might shift during the process of binding, resulting in enrichment of previously masked proteins. To investigate this hypothesis, whole cell lysates of Jurkat cells were incubated with CNBr-activated Sepharose beads, which covalently bind to free amine groups. Uncoupled proteins were collected at different time points. We found that the protein abundance pattern shifted after depletion (Supplemental Fig. S5).

A



B

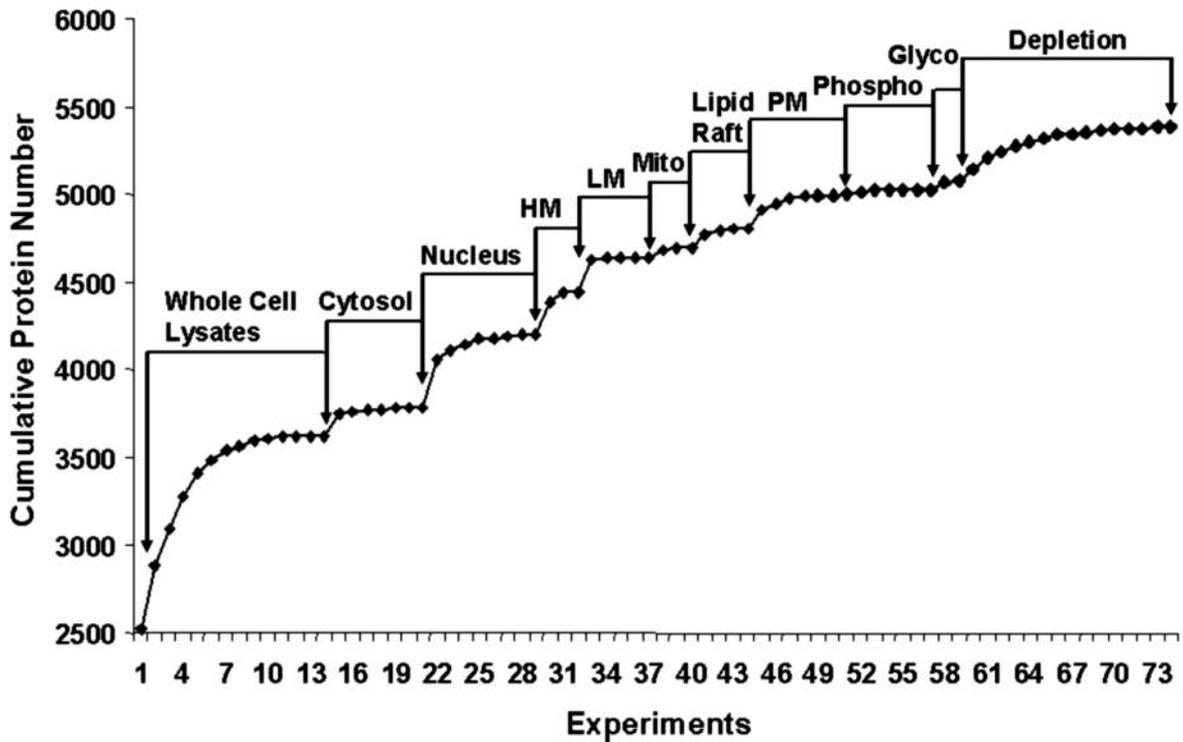


TABLE I
Summary of the proteomics data

The number of high confidence protein identifications (more than one high scoring peptide) and their associated unique peptide counts, spectral counts, and peptide false positive rates in each experimental fraction are shown.

Experimental fraction	Proteins	Unique peptides	Spectra	Peptide false positive rate
				%
Whole cell lysates	3,620	26,865	256,375	0.43
Cytosol	2,012	13,618	72,893	0.64
Heavy membrane	1,599	9,073	34,671	0.39
Light membrane	2,154	13,617	58,001	0.42
Mitochondrion	1,154	5,971	21,037	0.46
Nucleus	1,750	10,338	44,673	0.37
Raft	1,112	5,462	20,645	0.28
Plasma membrane	1,529	8,239	49,381	0.28
Glycoprotein	936	6,175	23,668	0.08
Phosphoprotein	1,033	4,964	21,693	0.50
Depletion	3,035	20,156	195,167	0.51
Total	5,381	43,693	798,204	

Next we repeatedly analyzed the proteomes from all of the above subfractions using one-dimensional gel electrophoresis combined with LC-MS/MS, and a combined total of 1707 LC-MS/MS runs were performed. The flow diagram of our experimental strategy is shown in Fig. 2A. All the mass spectrometry data were then searched against a non-redundant human protein database using the SEQUEST algorithm (21) followed by stringent filtering, resulting in the identification of 9611 unique proteins. Further selection of proteins identified by at least two high scoring unique peptides and exclusion of several apparent contaminants introduced during sample handling (e.g. trypsin and keratins) led to a total identification of 5381 proteins with high confidence with ~1000 to ~3600 proteins identified in each fraction (Fig. 2B, Table I, and Supplemental Table S2). The peptide false positive rate in each fraction was lower than 0.7% when the entire dataset was searched against the concatenated forward and reversed database (Table I). These data indicated that the criteria we used to filter our spectra are stringent, and the final protein list contained very few false positive identifications.

One essential issue in shotgun proteomics is that peptide sequences can be present in multiple protein entries due to closely related proteins (e.g. splice variants, homologs, paralogs, orthologs, or redundant entries in the protein database), leading to an overestimation of protein identification number. To address this issue, we computed protein probability and

observed that the vast majority of our identified proteins (80%) were assigned unambiguously, that is each unique protein group has only one representative protein in our identified protein list and was distinguished by at least one high scoring (peptide probability ≥ 0.9) unambiguous peptide. The remaining proteins could not be identified unambiguously because the same peptide sequences mapped to more than one protein. In these cases, it is not reliable to claim which proteins are truly expressed in the cells simply based on bioinformatics prediction. Therefore a straight bioinformatics method based on mass spectrometry data is inadequate to completely address the redundancy issues, indicating that an alternative method is needed to support large scale shotgun proteomics data.

Proteomics and Transcriptomics Profiles Comparison—To further address the redundancy issue, we next performed a genome-wide mRNA analysis in Jurkat cells as an independent approach to support our proteome profiling. Both total RNA and mRNA were prepared from Jurkat cells and then examined with a commercialized human oligonucleotide microarray that contains >46,000 transcript-specific probe sequences per array. Duplicate and triplicate arrays were analyzed in parallel using total RNA and poly(A)-enriched mRNA, respectively. Genes with detection *p* value less than 0.05 were regarded as positive identifications. A total of 15,592 and 15,286 unique gene targets were detected in total RNA and mRNA, respectively. 13,973 gene targets were jointly detected in both total RNA and mRNA (Fig. 3A).

Bridging gene symbols with protein accession numbers resulted in about 7000 gene/protein pairs for the expression comparison study. According to their different detections using mass spectrometry and microarray tools, we categorized these gene/protein pairs into four groups (A, B, C, and D) (Table II and Supplemental Table S3). Group A includes gene products detected by at least two unique peptides and high confidence mRNA expression. The majority of group A genes (4270 of 4522 genes) were matched to a single protein accession number, indicating that there is no redundancy among these protein identifications. A small set of group A genes (252 genes of 540 proteins) were matched to multiple protein accession numbers; this might be due to two reasons. First, there are redundant protein entries in our list. Second, some of the oligonucleotide probes on the chip may not be splice isoform-specific as claimed. To address this issue, we investigated the ProteinProphet results and observed that at least 332 proteins were identified in this subgroup based on at

Fig. 2. **Proteome survey of human Jurkat T leukemic cells.** A, the flow chart of overall experimental strategy for protein identification. Human Jurkat T leukemic cells were fractionated into 10 fractions including cytosolic, light membrane, heavy membrane, mitochondrial, nuclear, lipid raft, plasma membrane, phosphoprotein, glycoprotein, and depletion. Proteins extracted from the Jurkat whole cell lysates and the above subfractions were repeatedly analyzed by one-dimensional gel electrophoresis combined with LC-MS/MS (GeLC-MS/MS). A total of 1707 LC-MS/MS runs were performed resulting in the identification of 5381 proteins. B, the cumulative curve of total identified proteins following multiple enrichment methods is shown. *HM*, heavy membrane; *LM*, light membrane; *Mito*, mitochondria; *Phospho*, phosphoproteins; *Glyco*, glycoproteins; *PM*, plasma membrane; *Depletion*, proteins uncoupled from CNBr-activated Sepharose beads.

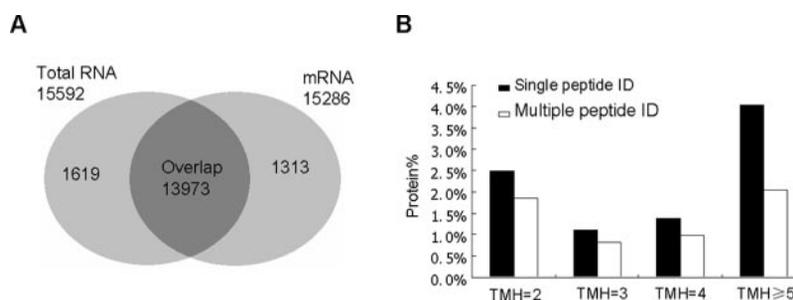


FIG. 3. **Proteomics and transcriptomics profile integration.** A, high confidence transcript detection. Using $p < 0.05$ as a cutoff, 15,592 and 15,286 unique gene targets were detected in purified total RNA and mRNA samples, respectively. 13,973 gene targets were jointly detected in both total RNA and mRNA. B, comparison between the distribution of membrane-associated proteins identified by single and multiple peptides. The accepted proteins identified (ID) by multiple peptides and a single peptide were applied for TMH prediction by TMHMM Server 2.0. The membrane protein distributions in these two categories are shown.

TABLE II

Summary of proteomics and transcriptomics profile comparison

Yes, the corresponding transcripts were detected with $p < 0.05$; No, the corresponding transcripts were not detected with $p < 0.05$. The numbers outside parentheses represent the matched protein number; the numbers inside parentheses represent the corresponding target gene number. ID, identification; N/A, not applicable.

Protein ID	Microarray target	Unique match no.	Redundant match no.	Total RNA	mRNA	ID no.	Group	
5381 proteins identified by more than one peptide	With target	4610	N/A	Yes	Yes	4180	A	
				Yes	No	32		
				No	Yes	58		
				No	No	340	B	
				N/A	Yes	Yes	536 (250)	A
				Yes	No	4 (2)		
				No	Yes	0		
Proteins identified by a single peptide	No target	N/A	N/A	N/A	N/A	211	C	
	With target	1733	0	Yes	Yes	1733	D	
				No	No	20 (10)	B	

least one unambiguous peptide sequence. Therefore a total of 4602 proteins (4270 + 332) were accepted as unambiguous identifications in group A. Group B includes 360 proteins (350 genes) identified with at least two peptides, but their transcripts were not detected with high confidence. This group might be due to either false positive identifications or mRNA levels that were too low or even degraded after translation. Among them 74 proteins were identified with more than four unique peptides including at least one unambiguous peptide sequence; these were accepted as positive identification. The remaining proteins were considered false positive identifications. In group C, the identified proteins have no corresponding gene target on the microarray. In these cases, we accepted 62 proteins that were identified with more than four unique peptides including at least one unambiguous peptide sequence. Group D comprises proteins identified with a single peptide, and their mRNA expressions were jointly detected in both purified total RNA and mRNA with high confidence. In this study, we identified more than 4000 proteins with a single peptide that have a relatively high false positive rate. However, transcript expression confirmation effectively reduced the number to 1733 proteins that we accepted in the final protein count; these proteins are listed separately in Supple-

mental Table S4. Thus, using these stringent criteria, we were able to accept a total of 6471 unique proteins (4602 + 74 + 62 + 1733).

It is known that membrane-associated proteins are more difficult to be detected with multiple peptides. Therefore we compared the distribution of integral membrane proteins identified with a single peptide (group D) and multiple peptides (combined groups A, B, and C). As expected, we observed that proteins in group D have a higher coverage of integral membrane proteins than those identified by multiple peptides (Fig. 3B).

In total, excluding the redundant protein entries and potential false positive identifications, we identified 6471 unique gene products by mass spectrometry among which 98% were verified by high confidence transcript expression (Supplemental Tables S2, S3, and S4). We also investigated the presence of membrane-associated proteins as a measure of proteome detection coverage. Using TMHMM Server 2.0 (29) to predict protein transmembrane helix (TMH),¹ we found that a total of 998 proteins in our accepted proteome dataset had

¹ The abbreviations used are: TMH, transmembrane helix; GO, Gene Ontology.

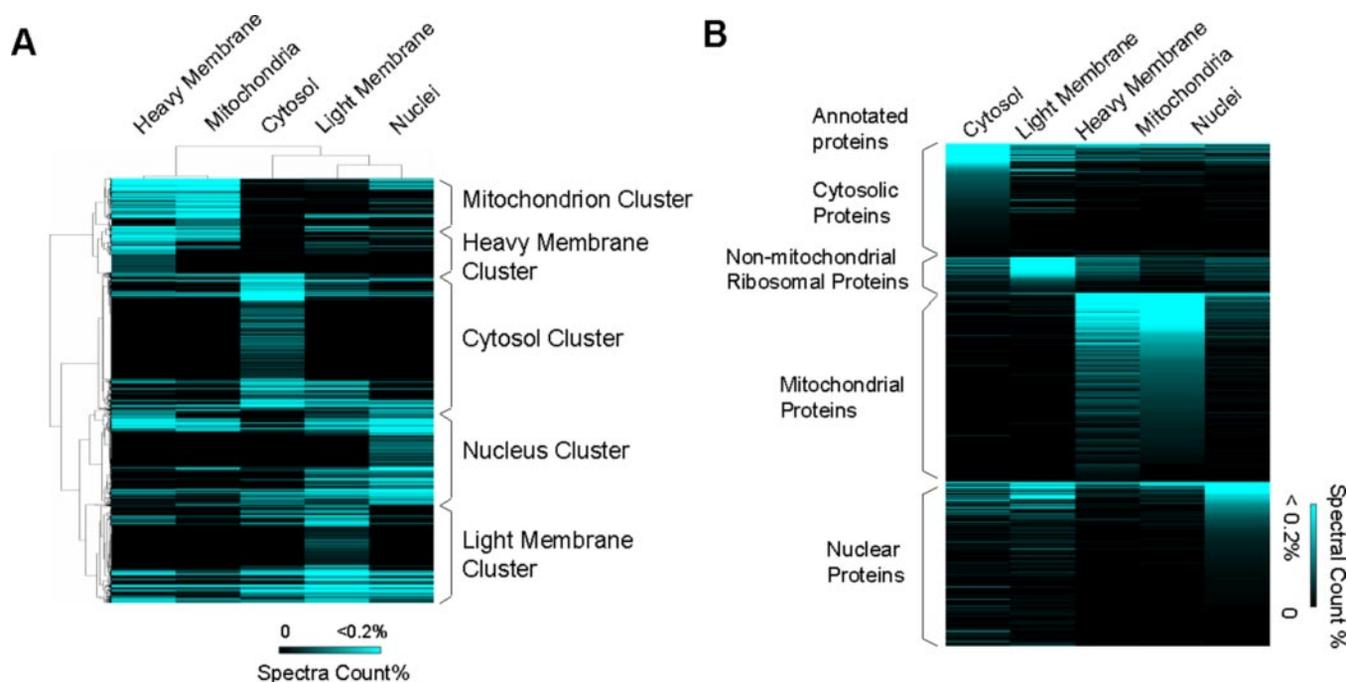


FIG. 4. Protein subcellular localization. *A*, hierarchical clustering of protein expression pattern obtained from cytosolic, light membrane, heavy membrane, mitochondrial, and nuclear fractions is shown. The protein expression level was measured by the normalized spectral count, *i.e.* the spectral count of each protein in one fraction divided by the total spectral counts of the same fraction. The five protein clusters are indicated. This pattern allowed the assignment of primary localization of each of the identified proteins except in some proteins with multicompartment distributions. Note the compartment-specific protein distribution pattern as well as the multicompartmental distribution patterns. *B*, the expression pattern of the gold standard proteins in cytosolic, light membrane, heavy membrane, mitochondrial, and nuclear fractions is shown. The gold standard proteins were selected from the identified proteins in different subcellular compartments using the GO terms.

at least one putative TMH, whereas 492 proteins had two or more TMHs (Supplemental Table S2). These results are better than a recent study where a similar scale proteome was characterized from the mouse organs and organelles (8). Therefore, we concluded that our dataset provides high coverage of membrane-associated proteins.

Subcellular Localization—One major advantage of proteomics measurements over mRNA profiling is the ability to deduce protein subcellular localization, providing insights into the organelle functions. However, one limitation associated with proteomics study is the difficulty to isolate pure subcellular organelles. Although we have attempted to purify sub-fractions efficiently and optimized the purification methods (20, 28), it may be inaccurate to assign protein localization solely based on the protein identification in the enriched fractions due to the possible cross-contamination during sample preparation. The specificity of protein assignment to particular organelles requires appropriate comparison and data analysis.

Because it has been reported that the spectral count of proteins is a semiquantitative measure of protein abundance (30), the proteomics data reported here might be useful for protein subcellular localization prediction. In this study, we focused on analyzing subcellular localization of proteins identified by at least two unique peptides but not having redun-

dant matches to the same gene target (cytosol, light membrane, heavy membrane, mitochondrion, and nucleus fractions; Supplemental Table S5) based on their spectral count. In addition to the subcellular localization analysis performed in these five fractions, we also provide a detailed protein identification list for all of the fractions that were analyzed (Supplemental Table S2). This information can be used to predict the global localizations of all of the identified proteins.

First we compared the normalized spectral counts obtained in this study with Western blotting results (Supplemental Fig. S2). It was observed that the normalized spectral counts of several biomarker proteins agree with their distribution in different organelles. We therefore assigned a primary localization to each of the proteins based on the normalized spectral count followed by hierarchical clustering (Fig. 4A and Supplemental Table S5). The clustering result shows significant pattern differences between the five subcellular fractions, confirming that our samples were enriched with distinct proteins belonging to different functional categories.

Next we used a “gold standard” protein list to assess the accuracy of our subcellular assignment. This gold standard list is constructed based on the Gene Ontology (GO) terms and comprises four specific subcellular compartments, *i.e.* cytosol, non-mitochondrial ribosome, mitochondria, and nu-

TABLE III

Validation of subcellular localization prediction of proteins based on the normalized spectral count followed by hierarchical clustering
N/A, not applicable.

Protein group	Cluster	Total protein no.	Known correct assignment no.	Enrichment score
616 gold standard proteins	Cytosol	165	121	0.73
	Nuclei	88	82	0.93
	Mitochondria	208	207	1.00
	Light membrane	130	42	0.32
	Heavy membrane	25	N/A	N/A
3370 identified proteins	Cytosol	1100	670	0.61
	Nuclei	768	520	0.68
	Mitochondria	373	259	0.69
	Light membrane	745	N/A	N/A
	Heavy membrane	384	N/A	N/A

clei. The accuracy of the GO terms assignments were further confirmed by comparison with the Human Protein Reference Database, which only contains information manually extracted from the literature by expert biologists (31). Only those proteins in agreement with the primary localization annotation in the Human Protein Reference Database were accepted as the gold standard proteins. Our list and the gold standard list have 616 proteins in common. By plotting the distribution of these gold standard proteins, we observed that cytosolic proteins are mainly enriched in the cytosolic fraction, non-mitochondrial ribosomal proteins are mainly enriched in the light membrane fraction, mitochondrial proteins are mainly enriched in the mitochondrial and heavy membrane fraction, and nuclear proteins are mainly enriched in the nuclear fraction (Fig. 4B). We hierarchically clustered these gold standard proteins as described above to assign a primary localization (Supplemental Table S6). If our assignment to a protein agreed with the annotations, we considered it a correct assignment. We quantified the degree of correct assignment in the four clusters that we were able to evaluate (cytosol, nuclei, mitochondria, and light membrane) using an enrichment score, which is defined as the ratio of the number of correct assignments to the number of total assignments in this cluster (Table III). A high enrichment score was observed for the nucleus and mitochondrion (0.93 and 1.00, respectively). The cytosol cluster had a moderate (0.73) enrichment score. We analyzed the potentially incorrect assignments in the cytosol, most of which (42 of 44) are nuclear proteins based on the annotation. This may be due to protein shuttling between organelles or incorrect protein annotation in the literature. For example, proliferating cell nuclear antigen is known to be present in the cytosol, but it also shuttles to the nucleus during cell proliferation (32). Therefore, incomplete assignment of many proteins in the literature penalized the actual enrichment score in the cytosol. For the light membrane cluster, although most non-mitochondrial ribosomal proteins were assigned to this cluster, many proteins from other organelles were also included, resulting in a relatively low enrichment score (0.32). This is because the light membrane fraction

comprises multiple membrane compartments such as lysosome, smooth/rough endoplasmic reticulum, Golgi apparatus, and even the debris of the nucleus. As for the heavy membrane fraction, also called crude mitochondrial fractions in some studies, most proteins (22 of 25) clustered in this fraction are mitochondrial proteins, which indicates that some proteins in this fraction localize in the mitochondrion. However, based on the distribution of several biomarkers, heavy membrane also enriched with other membrane structures such as endoplasmic reticulum and plasma membrane. Therefore, it is also difficult to specifically assign a primary localization to proteins clustered in heavy membrane. Therefore we concluded that the overall protein assignments in the three clusters (cytosol, nuclei, and mitochondria) are highly reliable, and caution should be paid to the light membrane and heavy membrane assignments.

We also verified the specificity of the clustering results using our protein list. We collected the annotated or predicted cytosolic, mitochondrial, and nuclear proteins based on the GO terms and PSORT II prediction (33). We observed that the majority (259 of 373) of proteins in the mitochondrion cluster are assigned to the mitochondrion, 670 of the 1100 proteins in the cytosol cluster are assigned to the cytosol, and 520 of the 768 proteins in the nucleus cluster are assigned to the nucleus. Because we obtained a higher enrichment score in these three clusters using the gold standard list, there is a high probability that the rest of the unassigned proteins in these three clusters are predicted correctly. Therefore, our proteome data provide a primary subcellular prediction for 2241 proteins with high confidence, including 792 proteins that are unannotated by the GO terms and unassigned by PSORT II. For the rest of the 1129 proteins in the light membrane and heavy membrane clusters, more dedicated subcellular prediction computation and assessment are required.

DISCUSSION

Eukaryotic cells segregate and organize functionally related proteins into discrete compartments that have distinct structures and functions. Previous organelle proteomics studies

have mainly focused on one compartment, providing insights into the biology and functions of these structures. Recently two groups performed magnificent proteomics studies on multiple organelles in mouse organ by combining subcellular fractionation and mass spectrometry technologies (8, 19). However, no comprehensive characterization of a single human cell type has been carried out to date. In this study, we combined replicate proteomics analyses and extensive subcellular fractionation/enrichment methods in Jurkat cells, identifying 5381 proteins of which 80% were assigned with at least one unambiguous peptide sequence. Based on comparison between proteomics and transcriptomics profiling in Jurkat cells, we were able to specifically exclude redundant entries and potential false positive identifications, resulting in 4738 protein identifications. Among them, more than 98% were confirmed by high confidence mRNA expression. Because we used multiple stringent criteria to filter and confirm our proteome dataset, the protein false positive rate was estimated to be closed to zero.

This proteome/transcriptome coverage is much higher than previous proteomics and transcriptomics comparison studies in mammalian cells (8, 34). It may be because previous studies either did not analyze the global expression of proteome or transcriptome comprehensively or compared proteomics and transcriptomics data generated from different biological systems (e.g. different mouse strains or cell type).

Although we performed a comprehensive proteome analysis of a single cell type resulting in the identification of a huge number of proteins with high confidence, many lower abundance and membrane-associated proteins are still refractory to rigorous identification by mass spectrometry techniques. This incomplete proteome coverage likely arose from the intrinsic limitations in instrument sensitivity and bias of data-dependent acquisition toward high abundance proteins (30). Additionally it may also be due to an overly stringent filtering of the mass spectrum search results. Consistent with this, over 4000 proteins were identified by one high scoring peptide. By simply accepting proteins identified by a single high scoring peptide if the corresponding transcript was jointly detected in both total RNA and mRNA samples, the number of protein identifications could be boosted to 6471. Moreover the accepted proteins identified with a single peptide have a higher coverage of membrane-associated proteins than the proteins identified with multiple peptides. This result indicates that proteomics and transcriptomics integration is a powerful tool to rescue false negative protein identification. Another benefit of proteomics and transcriptomics integration is to investigate how the mRNA levels reflect protein abundances and the biological mechanisms of the discordance between protein and mRNA expression. This ongoing investigation being conducted in our laboratory will be reported in the near future.

In this study, we used CNBr-activated Sepharose beads to covalently couple proteins and identified uncoupled proteins

during a time course. Using this approach, we were able to detect several hundred new unique proteins after previous extensive profiling. We found that this method helps purify specific classes of proteins including some hydrophobic proteins after a long incubation. But the reaction between proteins and the beads is not simply based on protein hydrophobicity (data not shown). The reaction rate between proteins and the beads is likely due to the combinational outcome of multiple protein properties such as hydrophobicity and the number of accessible free amine groups.

One unique advantage of proteomics profiling over transcriptomics profiling is the ability to provide information on protein post-translational modifications. In this study, we enriched proteins based on two types of post-translational modifications, phosphorylation and glycosylation. Therefore, proteins identified in these two fractions are likely phosphorylated or glycosylated. However, note that we did not specifically detect the modification sites, and the proteins were enriched based on affinity purification; some unspecific binding proteins and proteins associated with those phosphoproteins/glycoproteins may also be detected in these two fractions. Therefore large scale phosphopeptide and glycopeptide identification is needed to complement our dataset.

Another advantage of proteomics profiling is to deduce protein subcellular localization, providing insights into the biological functions of gene products and organelles. Given the difficulty of isolating completely pure organelles, we opted to combined differential protein expression in multiple subcellular fractions with hierarchical clustering to more accurately predict the primary subcellular localization of proteins. Using annotated proteins to assess the prediction accuracy, we were able to provide high confidence assignment by this method in at least three compartments (the cytosol, nuclei, and mitochondria). The primary subcellular localization assignment of 2241 proteins reported here, including 792 previously unassigned proteins by the GO term or PSORT II predictions, adds more information to the proteome composition of these organelles in this widely used human cell type. As for proteins clustered in the other two compartments (light membrane and heavy membrane), we chose to be more cautious and did not assign a specific primary localization to each of them because these two compartments comprise multiple subcellular structures. More defined subcellular fractionation approaches are required to further separate these fractions. However, the information here still provides a clue for protein localization. Together with the proteins detected in lipid raft and plasma membrane fractions, one can deduce much valuable information on those unassigned proteins in this study.

One of the main challenges confronting protein subcellular localization prediction is that many proteins likely shuttle between compartments, having multiple subcellular localizations. In this study we only assigned a primary localization to each protein, whereas most proteins reported here were detected in more than one subcellular fraction. Although some of

these cases may stem from cross-contamination during sample preparation, we cannot exclude that they may indeed reflect the real protein localization patterns. Moreover 1129 proteins clustered in the light and heavy membrane fractions were assigned a primary localization with low confidence due to the multiple organelle composition of these two fractions. Therefore, it is possible that by applying more advanced machine learning methods on the proteomics data reported here more accurate subcellular localization assignment can be expected.

Despite the caveats in the identification of post-translationally modified protein and the assignment of protein subcellular localization, the proteomics profile reported here provides a global landscape in a single human cell type, precluding the differences among tissues and more suitable for many in-depth characterizations of human systems at a cellular level such as comparison between protein and mRNA expression and integration of protein expression pattern with protein-protein interactions and biological phenotypes. In addition, some of the data in this study were obtained from cells after specific perturbations. More thorough analysis on the dynamic regulation of proteins in these cells has been (28) or will be reported separately. Because many mutant cell lines have been derived from human Jurkat cells and widely used by the biological community (16), the comprehensive proteome survey of this cell type can serve as a useful platform for more extensive experimental characterization and integration studies. Complete summaries of the peptides and proteins identified in this study are accessible in the supplemental data. In addition, all of the raw data generated in this study are being deposited with the *Molecular & Cellular Proteomics* data repository. Investigators are encouraged to utilize this rich proteomics resource.

SUPPLEMENTAL DATA

Supplemental data include experimental procedures, five figures, six summary tables, and the detailed information for each identified peptide (peptide atlas file) and their sharing results among different protein entries (ProteinProphet file).

* This work was supported by National Institutes of Health Grants RO1 HL 67569, PO1 HL 70694, and RR 13186. The costs of publication of this article were defrayed in part by the payment of page charges. This article must therefore be hereby marked "advertisement" in accordance with 18 U.S.C. Section 1734 solely to indicate this fact.

□ The on-line version of this article (available at <http://www.mcponline.org>) contains supplemental material.

§ These authors contributed equally to this work.

‡ To whom correspondence should be addressed: Center for Vascular Biology, Dept. of Cell Biology, University of Connecticut Health Center, 263 Farmington Ave., Farmington, CT 06030. Tel.: 860-679-2444; Fax: 860-679-1201; E-mail: han@nso.uconn.edu.

REFERENCES

1. Pan, Q., Shai, O., Misquitta, C., Zhang, W., Saltzman, A. L., Mohammad, N., Babak, T., Siu, H., Hughes, T. R., Morris, Q. D., Frey, B. J., and Blen-

cowe, B. J. (2004) Revealing global regulatory features of mammalian alternative splicing using a quantitative microarray platform. *Mol. Cell* **16**, 929–941

2. Su, A. I., Wiltshire, T., Batalov, S., Lapp, H., Ching, K. A., Block, D., Zhang, J., Soden, R., Hayakawa, M., Kreiman, G., Cooke, M. P., Walker, J. R., and Hogenesch, J. B. (2004) A gene atlas of the mouse and human protein-encoding transcriptomes. *Proc. Natl. Acad. Sci. U. S. A.* **101**, 6062–6067

3. Okazaki, Y., Furuno, M., Kasukawa, T., Adachi, J., Bono, H., Kondo, S., Nikaido, I., Osato, N., Saito, R., Suzuki, H., *et al.* (2002) Analysis of the mouse transcriptome based on functional annotation of 60,770 full-length cDNAs. *Nature* **420**, 563–573

4. Carninci, P., Kasukawa, T., Katayama, S., Gough, J., Frith, M. C., Maeda, N., Oyama, R., Ravasi, T., Lenhard, B., Wells, C., *et al.* (2005) The transcriptional landscape of the mammalian genome. *Science* **309**, 1559–1563

5. Cheng, J., Kapranov, P., Drenkow, J., Dike, S., Brubaker, S., Patel, S., Long, J., Stern, D., Tammana, H., Helt, G., Sementchenko, V., Piccolboni, A., Bekiranov, S., Bailey, D. K., Ganesh, M., Ghosh, S., Bell, I., Gerhard, D. S., and Gingeras, T. R. (2005) Transcriptional maps of 10 human chromosomes at 5-nucleotide resolution. *Science* **308**, 1149–1154

6. Gygi, S. P., Rochon, Y., Franza, B. R., and Aebersold, R. (1999) Correlation between protein and mRNA abundance in yeast. *Mol. Cell. Biol.* **19**, 1720–1730

7. Chen, G., Gharib, T. G., Huang, C. C., Taylor, J. M., Mizek, D. E., Kardia, S. L., Giordano, T. J., Iannettoni, M. D., Orringer, M. B., Hanash, S. M., and Beer, D. G. (2002) Discordant protein and mRNA expression in lung adenocarcinomas. *Mol. Cell. Proteomics* **1**, 304–313

8. Kislinger, T., Cox, B., Kannan, A., Chung, C., Hu, P., Ignatchenko, A., Scott, M. S., Gramolini, A. O., Morris, Q., Hallett, M. T., Rossant, J., Hughes, T. R., Frey, B., and Emili, A. (2006) Global survey of organ and organelle protein expression in mouse: combined proteomic and transcriptomic profiling. *Cell* **125**, 173–186

9. Lander, E. S., Linton, L. M., Birren, B., Nusbaum, C., Zody, M. C., Baldwin, J., Devon, K., Dewar, K., Doyle, M., FitzHugh, W., *et al.* (2001) Initial sequencing and analysis of the human genome. *Nature* **409**, 860–921

10. Venter, J. C., Adams, M. D., Myers, E. W., Li, P. W., Mural, R. J., Sutton, G. G., Smith, H. O., Yandell, M., Evans, C. A., Holt, R. A., *et al.* (2001) The sequence of the human genome. *Science* **291**, 1304–1351

11. Kunst, F., Ogasawara, N., Moszer, I., Albertini, A. M., Alloni, G., Azevedo, V., Bertero, M. G., Bessieres, P., Bolotin, A., Borchert, S., *et al.* (1997) The complete genome sequence of the gram-positive bacterium *Bacillus subtilis*. *Nature* **390**, 249–256

12. Adams, M. D., Celniker, S. E., Holt, R. A., Evans, C. A., Gocayne, J. D., Amanatides, P. G., Scherer, S. E., Li, P. W., Hoskins, R. A., Galle, R. F., *et al.* (2000) The genome sequence of *Drosophila melanogaster*. *Science* **287**, 2185–2195

13. Peng, J., Schwartz, D., Elias, J. E., Thoreen, C. C., Cheng, D., Marsischky, G., Roelofs, J., Finley, D., and Gygi, S. P. (2003) A proteomics approach to understanding protein ubiquitination. *Nat. Biotechnol.* **21**, 921–926

14. Ptacek, J., Devgan, G., Michaud, G., Zhu, H., Zhu, X., Fasolo, J., Guo, H., Jona, G., Breitkreutz, A., Sopko, R., McCartney, R. R., Schmidt, M. C., Rachidi, N., Lee, S. J., Mah, A. S., Meng, L., Stark, M. J., Stern, D. F., De Virgilio, C., Tyers, M., Andrews, B., Gerstein, M., Schweitzer, B., Predki, P. F., and Snyder, M. (2005) Global analysis of protein phosphorylation in yeast. *Nature* **438**, 679–684

15. Wu, L., and Han, D. K. (2006) Overcoming the dynamic range problem in mass spectrometry-based shotgun proteomics. *Expert Rev. Proteomics* **3**, 611–619

16. Abraham, R. T., and Weiss, A. (2004) Jurkat T cells and development of the T-cell receptor signalling paradigm. *Nat. Rev. Immunol.* **4**, 301–308

17. Arur, S., Uche, U. E., Rezaul, K., Fong, M., Scranton, V., Cowan, A. E., Mohler, W., and Han, D. K. (2003) Annexin I is an endogenous ligand that mediates apoptotic cell engulfment. *Dev. Cell* **4**, 587–598

18. Wang, P., Song, J. H., Song, D. K., Zhang, J., and Hao, C. (2006) Role of death receptor and mitochondrial pathways in conventional chemotherapy drug induction of apoptosis. *Cell. Signal.* **18**, 1528–1535

19. Foster, L. J., de Hoog, C. L., Zhang, Y., Xie, X., Mootha, V. K., and Mann, M. (2006) A mammalian organelle map by protein correlation profiling. *Cell* **125**, 187–199

20. Rezaul, K., Wu, L., Mayya, V., Hwang, S. I., and Han, D. (2005) A systematic characterization of mitochondrial proteome from human T leukemia cells. *Mol. Cell. Proteomics* **4**, 169–181
21. Eng, J., McCormack, A. L., and Yates, J. R., III (1994) An approach to correlate tandem mass spectral data of peptides with amino acid sequences in a protein database. *J. Am. Soc. Mass Spectrom.* **5**, 976–989
22. Han, D. K., Eng, J., Zhou, H., and Aebersold, R. (2001) Quantitative profiling of differentiation-induced microsomal proteins using isotope-coded affinity tags and mass spectrometry. *Nat. Biotechnol.* **19**, 946–951
23. Peng, J., Elias, J. E., Thoreen, C. C., Licklider, L. J., and Gygi, S. P. (2003) Evaluation of multidimensional chromatography coupled with tandem mass spectrometry (LC/LC-MS/MS) for large-scale protein analysis: the yeast proteome. *J. Proteome Res.* **2**, 43–50
24. Keller, A., Nesvizhskii, A. I., Kolker, E., and Aebersold, R. (2002) Empirical statistical model to estimate the accuracy of peptide identifications made by MS/MS and database search. *Anal. Chem.* **74**, 5383–5392
25. Nesvizhskii, A. I., Keller, A., Kolker, E., and Aebersold, R. (2003) A statistical model for identifying proteins by tandem mass spectrometry. *Anal. Chem.* **75**, 4646–4658
26. Eisen, M. B., Spellman, P. T., Brown, P. O., and Botstein, D. (1998) Cluster analysis and display of genome-wide expression patterns. *Proc. Natl. Acad. Sci. U. S. A.* **95**, 14863–14868
27. Saldanha, A. J. (2004) Java Treeview—extensible visualization of microarray data. *Bioinformatics* **20**, 3246–3248
28. Hwang, S. I., Lundgren, D. H., Mayya, V., Rezaul, K., Cowan, A. E., Eng, J. K., and Han, D. K. (2006) Systematic characterization of nuclear proteome during apoptosis: a quantitative proteomic study by differential extraction and stable isotope labeling. *Mol. Cell. Proteomics* **5**, 1131–1145
29. Krogh, A., Larsson, B., von Heijne, G., and Sonnhammer, E. L. (2001) Predicting transmembrane protein topology with a hidden Markov model: application to complete genomes. *J. Mol. Biol.* **305**, 567–580
30. Liu, H., Sadygov, R. G., and Yates, J. R., III (2004) A model for random sampling and estimation of relative protein abundance in shotgun proteomics. *Anal. Chem.* **76**, 4193–4201
31. Peri, S., Navarro, J. D., Amanchy, R., Kristiansen, T. Z., Jonnalagadda, C. K., Surendranath, V., Niranjana, V., Muthusamy, B., Gandhi, T. K., Gronborg, M., Ibarrola, N., Deshpande, N., Shanker, K., Shivashankar, H. N., Rashmi, B. P., Ramya, M. A., Zhao, Z., Chandrika, K. N., Padma, N., Harsha, H. C., Yatish, A. J., Kavitha, M. P., Menezes, M., Choudhury, D. R., Suresh, S., Ghosh, N., Saravana, R., Chandran, S., Krishna, S., Joy, M., Anand, S. K., Madavan, V., Joseph, A., Wong, G. W., Schiemann, W. P., Constantinescu, S. N., Huang, L., Khosravi-Far, R., Steen, H., Tewari, M., Ghaffari, S., Blobel, G. C., Dang, C. V., Garcia, J. G., Pevsner, J., Jensen, O. N., Roepstorff, P., Deshpande, K. S., Chinnaiyan, A. M., Hamosh, A., Chakravarti, A., and Pandey, A. (2003) Development of human protein reference database as an initial platform for approaching systems biology in humans. *Genome Res.* **13**, 2363–2371
32. Vrız, S., Lemaitre, J. M., Leibovici, M., Thierry, N., and Mechali, M. (1992) Comparative analysis of the intracellular localization of c-Myc, c-Fos, and replicative proteins during cell cycle progression. *Mol. Cell. Biol.* **12**, 3548–3555
33. Nakai, K., and Horton, P. (1999) PSORT: a program for detecting sorting signals in proteins and predicting their subcellular localization. *Trends Biochem. Sci.* **24**, 34–36
34. Hu, S., Li, Y., Wang, J., Xie, Y., Tjon, K., Wolinsky, L., Loo, R. R., Loo, J. A., and Wong, D. T. (2006) Human saliva proteome and transcriptome. *J. Dent. Res.* **85**, 1129–1133