# DAnTE: a statistical tool for quantitative analysis of -omics data

**Ashoka D. Polpitiya**, **Wei-Jun Qian**, **Navdeep Jaitly**, **Vladislav A. Petyuk**, **Joshua N. Adkins**, **David G. Camp II**, **Gordon A. Anderson**, and **Richard D. Smith**[*]
Pacific Northwest National Laboratory, Richland, WA 99352, USA.

## Summary

DAnTE (Data Analysis Tool Extension) is a statistical tool designed to address challenges associated with quantitative bottom-up, shotgun proteomics data. This tool has also been demonstrated for microarray data and can easily be extended to other high-throughput data types. DAnTE features selected normalization methods, missing value imputation algorithms, peptide to protein rollup methods, an extensive array of plotting functions, and a comprehensive hypothesis testing scheme that can handle unbalanced data and random effects. The Graphical User Interface (GUI) is designed to be very intuitive and user friendly.

## 1 INTRODUCTION

Although a number of tools are available for high-throughput microarray data processing (Saeed, Sharov et al. 2003; Gentleman, Carey et al. 2004), the data from LC-MS based quantitative bottom-up proteomics measurements (i.e., label-free approaches, stable isotope labeling methods, peptide/spectrum counting approaches, and the Accurate Mass and Time Tag method) pose different challenges than what these tools are designed to address. One of the major issues associated with proteomics data is often the extent of missing values that is largely due to limited dynamic range and leads to unbalanced datasets. In addition, proteomics data involves another level of grouping or "rollup" information to map peptides to proteins. Peptide abundances are often used to infer the corresponding protein abundances.

Developed to address the issues common to proteomics data, DAnTE is readily extendable. Though the target application is high-throughput proteomics, DAnTE has also been successfully demonstrated for microarray data analysis and can readily be applied to other forms of high-throughput "omics" data that bears similar characteristics (e.g., metabolomics data). A screenshot of the DAnTE user interface is illustrated in Figure 1.

## 2 DESCRIPTION

### 2.1 Dependencies

The graphical user interface (GUI) of DAnTE is implemented using the C# language, and the core algorithms are implemented in the open source R statistical environment (Ihaka and Gentleman 1996). DAnTE runs on a Microsoft Windows XP platform within a .NET 2.0 framework. The connectivity between R and the C#/.NET environment is achieved by using

the open source R(D)COM server application (Baier and Neuwirth 2007). This unique choice of environments makes DAnTE a very user friendly software tool, even though it cannot integrate into the popular Bioconductor (Gentleman, Carey et al. 2004) project.

## 2.2 Application features

**2.2.1 Data loading—**The input data to DAnTE can be any file that stores tabular data, including flat files (either CSV or tab-delimited text files) and Microsoft Excel files. A unique feature of the data loading mechanism is that it preserves peptide-to-protein mapping information for use later in plotting peptides that belong to a particular protein, as well as in the peptides-to-protein rollup methods. In addition, DAnTE can also process SEQUEST (Eng, McCormack et al. 1994) results and create spectral count tables.

**2.2.2 Factor Definitions—**Factors are used to capture the fixed and random effects in experimental design. For example, the biological condition is a fixed effect factor, while a list of liquid chromatography (LC) columns used to separate the samples can be treated as a random effect. This information is vital in normalization, imputation, and hypothesis testing methods in DAnTE. Factors can either be declared once the data is loaded or be loaded from a flat file.

**2.2.3 Investigative plots—**Various statistical plots, including histograms, box plots, correlation diagrams, and MA (or R-I: ratio-intensity) plots can be plotted in DAnTE. These plots help the user evaluate reproducibility within the study set and single out problematic datasets so that they can be excluded from further analysis.

**2.2.4 Data normalization—**As normalization is arguably the most important step in downstream data analysis, DAnTE employs several normalization methods that have been successfully tested for both proteomics data (Callister, Barry et al. 2006) and microarray genomics data (Quackenbush 2002; Smyth, Yang et al. 2003). Among them are a robust linear regression method, lowess method, and a quantile normalization method. In addition, global intensity adjustment based on median absolute deviation (MAD) and central tendency adjustment methods are also available.

**2.2.5 Missing value imputation—**Incomplete datasets due to missing values are common with high-throughput proteomics. As imputing these values is a much debated topic (Troyanskaya, Cantor et al. 2001), DAnTE offers several simple methods, as well as some advanced algorithms to chose from. The simple methods allow the user to fill in missing values with either the dataset mean/median or with a pre-chosen constant. Advanced methods include filling in with a row mean based on a user defined factor, K-nearest neighbor imputation (KNNimpute), and singular value decomposition based imputation (SVDimpute).

**2.2.6 Peptide to Protein rollup—**In most proteomics methods, peptide measurements are rolled up to corresponding protein abundances. Ideally, all peptides from a single protein should have similar abundances that manifest as similar signal intensities; however, in reality many factors, such as digestion efficiency, electro-spray ionization efficiency, etc., can affect the identifications and abundances or signal intensities of peptides. In the RRollup method available in DAnTE, peptides that originate from the same protein are first scaled on the basis of a chosen reference peptide in order to bring all peptide profiles across biological conditions to the same level and then averaged to obtain the protein abundance. During scaling, the peptide with the most observations is chosen as the reference peptide and its total abundance across datasets is used as a tie breaker. In the ZRollup method, a scaling method similar to $z$-scores (except that medians instead of means from peptide profiles across biological conditions are used) is applied first to peptides that originate from a single protein and then the scaled pepe-tides are averaged to obtain relative protein abundance. In both RRollup and Zrollup methods,

outlying peptide values are excluded from protein abundance calculations, using a Grubb's outlier test (Grubbs 1969). In the third QRollup method, peptides are selected on the basis of a user selected abundance cutoff value, and protein abundance is calculated as the average of these selected peptides.

**2.2.7 Analytical algorithms**—DAnTE offers several well characterized algorithms to further explore patterns in the data. Traditional principal component analysis (Jolliffe 2002) and associated scores and loadings plots can be useful as an unsupervised way of finding the principal variation in the data. In contrast, the partial least squares method (Wold, Albano et al. 1984) available in DAnTE can be used as a discrimination procedure whereby the grouping information is assigned using factors. Hierarchical and k-means clustering methods on features/ samples are also available as part of the heat map plotting function.

**2.2.8 Hypothesis testing**—A comprehensive ANOVA scheme for unbalanced studies that uses marginal sums of squares (Fox 1997) and mixed models (Pinheiro and Bates 2000) is included in DAnTE. The user can also test for interactions among factors in a multi-way ANOVA. The q-values are also calculated along with the p-values in order to control the false discovery rate (FDR) in multiple testing (Storey 2002). In addition, DAnTE can check whether the data follows a normal distribution by employing the Shapiro-Wilks test and features two non-parametric hypothesis tests (Wilcoxon rank sum test and Kruskal-Walis test) when the normality assumption fails to hold.

## 3 SUMMARY

DAnTE is designed as a complete downstream analysis tool that incorporates a host of algorithms for large-scale bottom-up proteomics data. This tool features an interactive GUI interface and harnesses the power of R statistical environment; its uniqueness lies in its ability to handle incomplete data and to roll peptides up to proteins. Though designed specifically for analyzing proteomics data, DAnTE performs equally well on genomics microarray data.
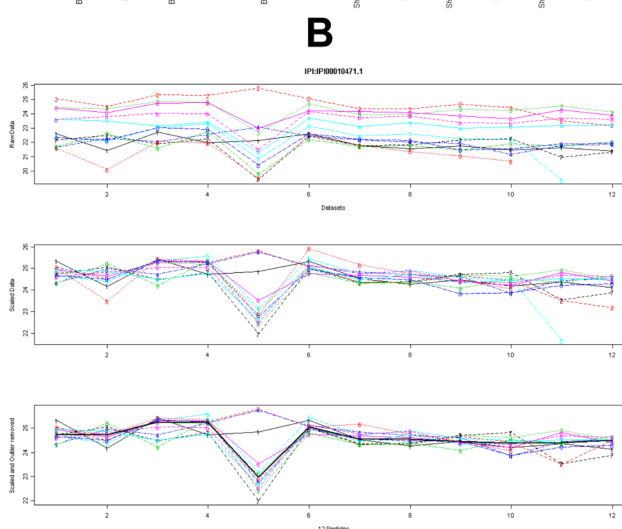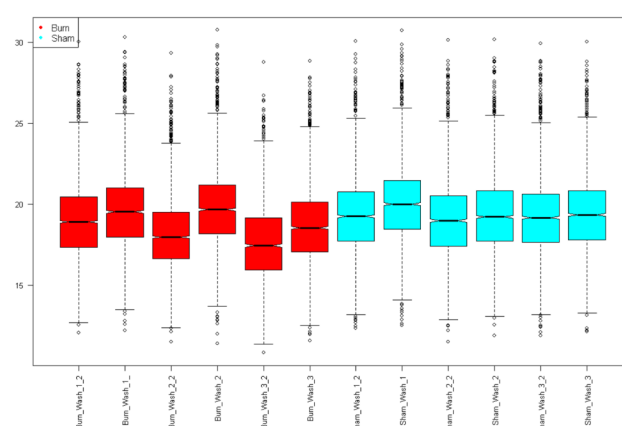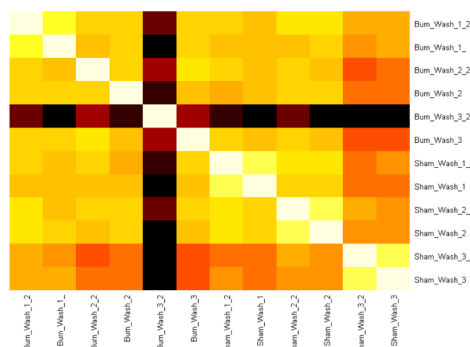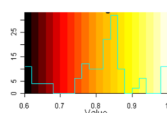
## Acknowledgments

## REFERENCES

Baier, T.; Neuwirth, E. R (D)COM Server V2.01. 2007. from http://sunsite.univie.ac.at/rcom/.

Callister SJ, Barry RC, et al. Normalization approaches for removing systematic biases associated with mass spectrometry and label-free proteomics. J Proteome Res 2006;5(2):277–286. [PubMed: 16457593]

Eng JK, McCormack AL, et al. An Approach to Correlate Tandem Mass-Spectral Data of Peptides with Amino-Acid-Sequences in a Protein Database. Journal of the American Society for Mass Spectrometry 1994;5(11):976–989.

Fox, J. Applied Regression Analysis, Linear Models, and Related Methods. Thousand Oaks, CA: Sage Publications; 1997.

Gentleman RC, Carey VJ, et al. Bioconductor: open software development for computational biology and bioinformatics. Genome Biol 2004;5(10):R80. [PubMed: 15461798]

Grubbs F. Procedures for Detecting Outlying Observations in Samples. Technometrics 1969;11(1):1–21.

Ihaka R, Gentleman R. R: A Language for Data Analysis and Graphics. Journal of Computational and Graphical Statistics 1996;5(3):299–314.

Jolliffe, IT. Principal Component Analysis. New York: Springer; 2002.

Pinheiro, JC.; Bates, DM. Mixed-Effects Models in S and S-PLUS. Springer; 2000.

Quackenbush J. Microarray data normalization and transformation. Nat Genet 2002;32:496–501. [PubMed: 12454644]

Saeed AI, Sharov V, et al. TM4: a free, open-source system for microarray data management and analysis. Biotechniques 2003;34(2):374–378. [PubMed: 12613259]

Smyth GK, Yang YH, et al. Statistical issues in cDNA microarray data analysis. Methods Mol Biol 2003;224:111–136. [PubMed: 12710670]

Storey JD. A direct approach to false discovery rates. Journal of the Royal Statistical Society Series B-Statistical Methodology 2002;64:479–498.

Troyanskaya O, Cantor M, et al. Missing value estimation methods for DNA microarrays. Bioinformatics 2001;17(6):520–525. [PubMed: 11395428]

Wold S, Albano C, et al. Modeling Data Tables by Principal Components and Pls - Class Patterns and Quantitative Predictive Relations. Analusis 1984;12(10):477–485.

**Fig.1.**
Representative screen shots from DAnTE. (A) Data grid and the navigation panel on the left; (B) Box plot of log transformed data; (C) A correlation heatmap of a set of data showing a possible outlier dataset; (D) Peptides to protein rollup results from RRollup method (panels from top to bottom: raw data; scaled data; median profile shown as a thick black line after outliers removed).