

John M. Abowd

School of Industrial and Labor Relations

Cornell University

259 Ives Hall

Ithaca, NY 14853-3901

Phone: (607) 255-4801 (CISER)

Fax: (607) 255-4496 (CISER)

E-mail: John.Abowd@cornell.edu

John Haltiwanger

Department of Economics

University of Maryland

College Park, MD 20742

Phone: (301) 405-3504

Fax: (301) 405-3542

Email: haltiwang@econ.umd.edu

Julia Lane

The Urban Institute

2100 M St NW

Washington DC 20037

Phone: (301) 675-3020

Fax: (202) 463-8522

Email: jlane@ui.urban.org

Session Title: **New Data and New Questions in Personnel Economics**

Word-processing package: **Word 2002**

Integrated Longitudinal Employee-Employer Data for the United States

by John M. Abowd, John Haltiwanger and Julia Lane*

The development of a database infrastructure that captures the complex interactions among households and businesses at the micro level and characterizes the dynamics of the modern economy is critical for the social sciences. The creation of such an infrastructure has posed a major challenge to national statistical institutes. Since most institutes collect, store and disseminate data on the engines of economic growth—businesses and households—in twin data silos, proposals to integrate the two face technical, monetary, legal, and policy obstacles that go far beyond the norm of data collection activities. Recent efforts at the Longitudinal Employer-Household Dynamics (LEHD) Program at the U.S. Census Bureau have finally made this critical data infrastructure achievable and accessible.

The potential uses of longitudinal integrated employer-employee data are far-reaching. A partial list includes: the effect of technological and structural change on earnings, employment and productivity; the analysis of the firm-specific contribution to pay; the effect of firm wage setting and turnover policies on productivity; the impact of firm policies on different groups of workers (*e.g.*, welfare recipients); the effect of firm expansion, exit and relocation decisions on neighborhood demographic composition, the analysis of worker commuting patterns and mobility; and a full accounting of the returns to firms and workers of investments in enterprise training, research and development. These data are also essential for unraveling some of the striking findings in the industrial organization and productivity literatures regarding the nature of business dynamics and the sources of micro and aggregate productivity growth.

Fields other than industrial organization and labor economics will be similarly enriched. For example, environmentalists can examine the impact of different firms' pollution levels on

worker and firm outcomes. Health specialists can examine, with large samples and without expensive clinical studies, the effects of earnings, employment history, firm personnel practices, and health benefit availability on death rates. Demographers can examine whether workers from different countries sort into firms that have hired their countrymen, what kinds of firms employ immigrants, the dynamics of immigrant mobility across firms, and a host of other issues. Finally, the burgeoning availability of integrated employer-employee data from advanced economies (many European countries as well as New Zealand) as well as developing and transition economies (*e.g.*, Colombia and Slovenia) will permit, for the first time, detailed international comparisons of the dynamic interrelationships among firms and workers—thus extending the research possibilities to international as well as national scientists.

In this paper we describe the new database infrastructure and the current status of research. We also summarize the proposed access protocols. We close with a description of our future agenda.

The Structure of the LEHD Program Data

The LEHD database infrastructure is complex. Figure 1 provides a visual summary. The core integration records are state Unemployment Insurance wage records (which are described in detail elsewhere). The integration of the business and demographic data by means of these records takes place under strict confidentiality protection protocols.¹ The UI records, from 22 partner states representing about 60% of U.S. employment, are reports filed by employers every quarter for each individual in covered employment. Using these records LEHD creates a database that provides longitudinal information on workers, firms, and the match between the two. Coverage is approximately 96% of private non-farm wage and salary employment; the coverage of agricultural and federal government employment is less comprehensive. Self-

employed individuals and independent contractors are also not covered (David Stevens, 2002). Although the identifiers in the administrative records are subject to some error, researchers have invested substantial resources in editing the identifiers and making them internally consistent (Abowd and Lars Vilhuber, forthcoming). This permits the use of edited identifiers as the primary linking protocol across all of the arrows shown in Figure 1. The secondary linking protocol is probabilistic record matching. The tertiary linking protocol is based on combinations of identifiers, geocodes, and other entity characteristics. The researcher then chooses the level of aggregation and the appropriate linked entities (Abowd, Haltiwanger, Ron Jarmin, Lane, Paul Lengermann, Kristin McCue, Kevin McKinney, and Kristin Sandusky, forthcoming).

Basic demographic information (date of birth, place of birth, sex, and a crude measure of race and ethnicity) is integrated via the person identifier link for almost all workers in the data—the non-match rate is about 4%. Other demographic survey data are integrated if their use is permitted under Title 13 of the U.S. Code. The Census Business Register is the core integration file for business data using the federal Employer Identification Number as the primary linking entity. Other economic censuses and surveys are also integrated (again if their use is permitted under Title 13 of U.S.C). Residential and establishment addresses are geocoded to the rooftop.

The sheer volume of data, while posing substantial computing challenges, is also a source of much of the analytical usefulness. The LEHD Program maintains universe files for the integration record, individual characteristics, and employer characteristics. Researchers can accurately change analysis frames from households/individuals to jobs to employing entities without adjusting or re-weight the analysis sample because of differential linking probabilities. The LEHD infrastructure universe is updated quarterly. Currently, about 80 million individual records, 5 million business records, and all of their associated wage records are processed every

quarter, which provides an unparalleled level of geographic and industrial detail. For additional technical information, see <http://lehd.dsd.census.gov> following the link to “Documentation.”

Current Status of Research and Data Development

The research program operates under a number of policy and financial constraints. All work that is undertaken using Census Bureau data must have, as its predominant purpose, the improvement of economic and demographic censuses, surveys and inter-censal population estimates. In addition, under the terms of the Memoranda of Understanding with the states, research that is not specifically identified in the MOU must also be approved by the state custodians. Finally, the very real financial exigencies of the program have meant that any applied research project must also contribute at least 50% to the development of the database infrastructure. These constraints, combined with the comparative advantage of the Principal Investigators, have meant that the focus of the research undertaken so far have been oriented towards the analysis of micro data and economic modeling.

i) Quarterly Workforce Indicators. The flagship research product of the program is the Quarterly Workforce Indicators (a subset of which is online at <http://lehd.dsd.census.gov>). These 29 indicators provide information on employment, job creation and destruction, accessions (hires and recalls), separations (exits and layoffs) at the county, metro and Workforce Investment Area, by eight age and two sex categories, and by detailed industry, for all quarters for which data are available for each partner state. The methods are based on Steven Davis, Haltiwanger and Scott Schuh, 1996, and Simon Burgess, Lane and Stevens, 2000.

The new indicators can be used in a variety of arenas. They highlight the dynamism of the U.S. economy—an illustrative example of which is provided in Table 1. The table demonstrates features of the local labor market that cannot be learned from other statistical

sources. For the illustrated period in this particular state the employment picture was quite negative. But the poor outlook varied markedly across age groups—the 19-21 year olds lost about 0.7% of jobs; the 35-44 year olds lost 0.9%, and the oldest group lost over 2%. Although jobs were lost on net, there were still job creations. This creation rate varied dramatically across age groups again: 18% of jobs for 19-21 year-olds were newly created compared with about 5% for the older cohort. Hiring continued even during the slowdown in economic activity—over 43% of the youngest cohort, 12% of 35-44 year-olds, and 7% of 55-64 year-olds were in new jobs in the next year.

The Quarterly Workforce Indicators can also be used as local labor market controls in regression analysis; to identify long term trends; to provide local context in performance evaluations, and a host of other applications. Additional indicators will be developed as the program expands. Measures of individual earnings dynamics across consecutive quarters are already in the internal system but have not been released. Indicators of cross-state flows and inter-industry mobility are in development.

ii) Human capital/productivity. Economists have long recognized that standard measures of workforce quality—typically years of education and experience—and firm characteristics—typically industry and firm size—fall woefully short of measuring the heterogeneity of workers and firms. A major research focus of the LEHD Program has been to calculate the human capital embodied in each individual in the dataset as well as the firm pay premium (Abowd, Francis Kramarz and David N. Margolis, 1999; Abowd, Lengermann and McKinney, 2002). The market value of the portable part of an individual's skill and has two components: a person effect, which does not vary over time, and a component based on labor market experience. The firm effect, which also does not vary over time and is estimated simultaneously with the other effects,

captures the average premium or discount that a given firm pays its workers, controlling for the other effects. The decomposition thus enables researchers to quantify the impact of firm quality (historical and current) on worker outcomes as well as the impact of workforce quality and turnover on firm outcomes. It is worth emphasizing how powerful these new measures are. Traditional surveys of workers that measure the “kitchen sink” of demographic characteristics—such as education, occupation, age, sex, marital status and even include some firm characteristics such as firm size and industry—are typically able to account for roughly 30% of earnings variation. Analysis including these new measures of worker and firm quality account for closer to 90% of earnings variation.

iii) Detailed industry studies/Sloan. One of the major results to come out of the analysis of micro-data on businesses has been the enormous within-industry heterogeneity in the ways in which firms organize themselves and produce their output. The integrated establishment universe maintained by LEHD in this new database infrastructure permits the in-depth examination of very narrowly defined industries using both empirical and case study approaches. A project funded by the Alfred P. Sloan foundation combines the rich industry level expertise of five Sloan Industry Centers—software, retail food, finance, trucking and semiconductors—with the new measures of workforce quality and workforce turnover at the firm level to examine the effect of how firms choose their workforce on economic growth, productivity and earnings outcomes.

iv) Transportation. Although most economists think of the market as being the most interesting aspect of the interaction between workers and firms, an important component of the LEHD Program data is the ability to incorporate the spatial aspects of the interaction. Because the data have information on both the place of work and place of residence of workers—and how these change over time—the value of the data infrastructure for researchers in transportation,

economic development and regional planning is unsurpassed. Initial research, in cooperation with the states of Florida, Illinois and Minnesota and the Bureau of Transportation Statistics, has enabled the production of a variety of releasable statistics: block level origin to destination flows of employee numbers from household to place of employment; information on the characteristics of workers by block residence (the number of workers living on each block; the proportion of workers earning low, medium, or high annual wages and mean annual wages) and information on the characteristics of businesses by block (mean quarterly pay per worker and the industry classification of firms operating on each block).

v) Low Wage Work and Welfare Recipients. The combination of skill-biased technological change and the reform of welfare laws have resulted in increasing policy concern about the employment and earnings outcomes of low-wage workers and welfare recipients. The LEHD data can contribute substantially because they permit the study of the longitudinal paths taken by workers and the direct measurement of the separate contributions of firms and workers to the outcomes of affected workers. The impact of the employer on the ability of workers to transition out of low-wage work was documented in a substantial recent project funded by the Rockefeller and Sage Foundations. (See Harry Holzer, Lane and Vilhuber, forthcoming and Fredrik Andersson, Holzer, and Lane, forthcoming.)

vi) Aging. The aging of the workforce raises an extensive range of political, economic, and social issues for the nation. The LEHD data address these questions: what types of firms employ older workers; how does the likelihood of employing older workers vary by industry and firm characteristics; what is the persistence and heterogeneity in employers' workforce composition; how heterogeneous are firms in their adjustment of workforce composition—who is hiring and firing older workers; how are the earnings outcomes of older workers related to firm

characteristics; and how does the changing nature of the firm affect older workers. New information about the turnover and earnings of workers in industries that cater to the elderly, such as the nursing home industry has been produced.

Next Steps

i) The evolution of industries and the changing nature of firm. A confluence of many factors has changed the nature of the firm—the movement towards the service sector, the ubiquitous use of computers, and the Internet in the workplace—have transformed both the way businesses do business and the way workers interact with coworkers, suppliers and customers. These changes increasingly imply that the key input into the activity of a business is the skill and knowledge of its workforce. Re-thinking the theory of the firm leads to new insights into critical questions. What constitutes a business? Why are businesses organized in the way that they are? What constitutes the value created by a business activity? Although these issues appear abstract, they profoundly impact the way we collect and process data on businesses and workers. Moreover, for many businesses (especially in the service sector), the answers to virtually all of these questions require integrating firm and worker characteristics and outcomes.

ii) Heterogeneous agent macroeconomics (measurement and analysis). The development of the LEHD data offers an unprecedented ability to build aggregate statistics from micro (household and firm) to macro in an integrated and consistent fashion. The current practice, which integrates data from disparate sources at the industry, regional or national level, misses much of the churning of workers and firms that recent empirical studies show are important for understanding aggregate fluctuations. These recent studies suggest that a key factor determining the success of an economy is how well its market structures and institutions handle this constant churn in the face of both frictions and imperfections. This heterogeneity in outcomes raises

questions about the ability of workers to insure against idiosyncratic consumption risk, and related government interventions in the marketplace. Until the development of integrated employer-employee data it has been very difficult to study the complex interactions involving the tradeoffs between economic efficiency and insurance.

iii) The Virtual Research Data Center (RDC). It goes without saying that these data are extremely sensitive in nature. Extraordinary precautions have been taken during their construction. The LEHD Program works on special-purpose dedicated computers with their own security plan and access protocols. All Privacy Act identifiers are removed from the files before they are integrated; a special individual identifier has been adopted for internal Census use. Only authorized researchers working from authorized Census-controlled areas have worked with the LEHD micro data. However, a major effort to make the data available to external researchers began in January, 2004. Important enhancements to the Census Business Register in the form of establishment-level data on workforce composition, turnover, and earnings have been made available to Census RDC researchers for approved Title 13 projects via the standard Census review process, which is managed by the Center for Economics Studies at <http://www.ces.census.gov>.

While there has been tremendous progress in terms of micro data development and substantially increased access to the micro data in the U.S. statistical community via the RDC network, there are still many limitations. There are only eight RDCs (including CES in Washington). The steps required to obtain access are non-trivial in terms of time and resources for both the researcher and the Census Bureau. Beyond the obvious cost of traveling to an RDC, there are many other barriers to access, not least of which is the steep learning curve associated with understanding the complex data infrastructure.

Reaping the full benefits of the enormous investment in the database infrastructure requires that researchers at many different academic institutions from a wide variety of disciplines collaborate on basic research, link information from a variety of sources, disseminate their results and replicate each others research. This cannot occur without the simultaneous development of a protocol to enable such activities to occur. We are developing a multi-layered access protocol that builds on recent data infrastructure developments and the access modalities as they currently exist. Key components of this multi-layered access protocol are the development of inference-valid public use synthetic micro data, access to richer synthetic micro data at a *virtual RDC*, and, in turn, limited access to the *gold standard* micro data in the Census/NSF RDC network.

Conclusion

We began this paper by arguing that the development of longitudinally integrated employer-employee data was critical to many social science researchers. The initial research represents but the beginning of a very broad agenda—not simply for the LEHD Program but for the wider research community. We hope that the general data integration approach that we have painstakingly developed can be used in many diverse research areas such as health services and geographic information systems. There are many steps to be taken. Perhaps the most important is to make the integrated data more accessible to the research community. We think that the creation of the *virtual RDC* that will permit public-use access to inference-valid synthetic data in parallel with the inherently more limited access to the internal *gold standard* data has tremendous promise for the social science research community.

References

Abowd, John M. and Vilhuber, Lars. “The Sensitivity of Economic Statistics to Coding Errors in Personal Identifiers.” Journal of Business and Economics Statistics, forthcoming.

Abowd, John M. Haltiwanger, John; Jarmin, Ron; Lane, Julia; Lengermann, Paul; McCue, Kristin; McKinney, Kevin and Kristin Sandusky. “The Relation among Human Capital, Productivity and Market Value: Building Up from Micro Evidence,” in C. Corrado, J. Haltiwanger, and D. Sichel, eds., Measuring Capital in the New Economy, Chicago: University of Chicago Press for the NBER.

Abowd, John M. Lengermann, Paul and McKinney Kevin. “The Measurement of Human capital in the U.S. Economy” LEHD Technical Paper No. TP-2002-09, 2002.

Abowd, John M. Kramarz, Francis and Margolis, David N. “High Wage Workers and High Wage Firms.” Econometrica, March 1999, 67(2), pp. 251-333.

Andersson Fredrik; Holzer, Harry and Lane, Julia. Moving Up or Moving On: Workers, Firms and Advancement in the Low-Wage Labor Market, New York: Russell Sage Press, forthcoming.

Burgess, Simon; Lane, Julia and Stevens, David. “Job Flows, Worker Flows and Churning.” Journal of Labor Economics, July 2000, 18(3), pp. 473-502.

Davis, Steven; Haltiwanger, John and Schuh, Scott. Job creation and destruction, Cambridge, MA: MIT Press, 1996.

Holzer, Harry, Lane, Julia and Vilhuber, Lars. “Escaping Poverty for Low-Wage Workers: The Role of Employer Characteristics and Changes.” Industrial and Labor Relations Review, forthcoming.

Stevens, David. “Employment that is not covered by State Unemployment.” LEHD Technical Paper no. TP-2002-16, 2002.

Table 1: 2001 Quarterly Workforce Indicators (non farm, private sector employment) for Pennsylvania by age group

	<i>19-21</i>	<i>35-44</i>	<i>55-64</i>
Total Employment	277,894	1,274,474	509,417
Net Job Change	-1,988	-12,004	-11,183
Jobs Created	49,184	81,250	27,730
New Hires	119,070	155,869	36,132

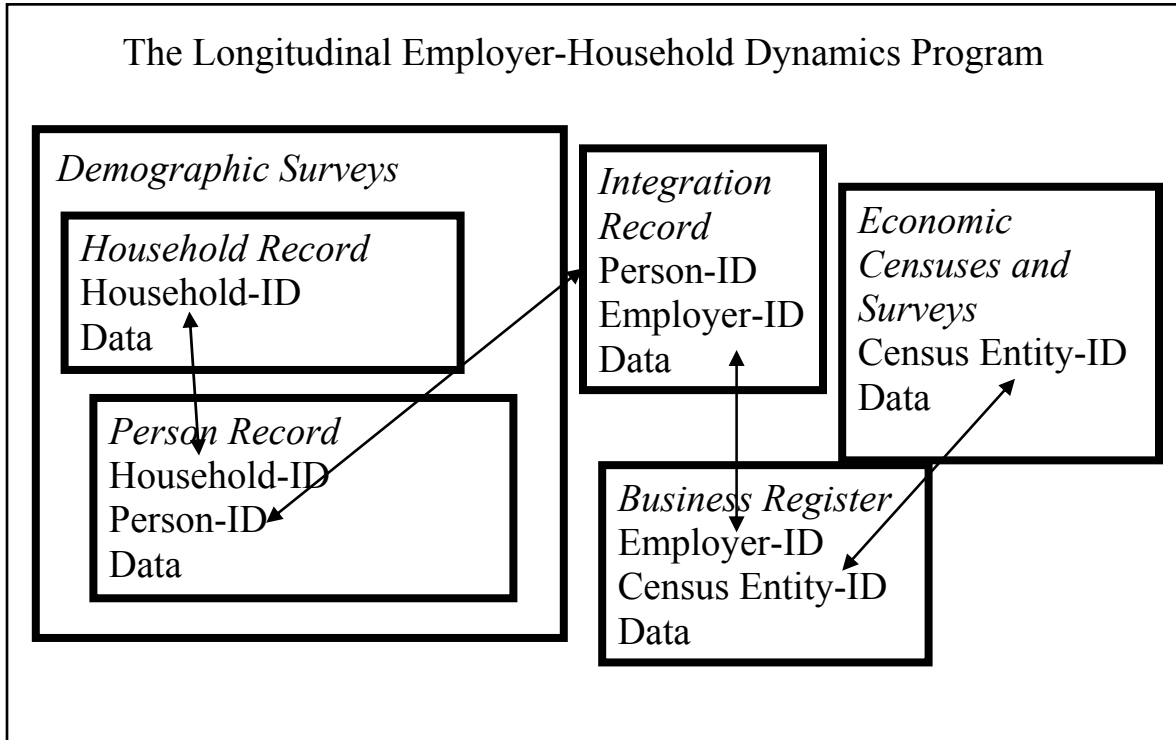


Figure 1

* School of Industrial and Labor Relations, Cornell University; Department of Economics, University of Maryland and the Urban Institute, respectively. All three authors are Senior Research Fellows at the United States Census Bureau. This research is a part of the Census Bureau's LEHD Program, which is partially supported by the National Science Foundation Grant SES-9978093 to Cornell University, the National Institute on Aging (R01-AG18854-01), and the Alfred P. Sloan Foundation. The views expressed herein are attributable only to the authors and do not represent the views of the U.S. Census Bureau, its program sponsors or data providers.

¹ The data are anonymized before use. They may be used for statistical purposes for approved projects by Census Bureau employees (including Special Sworn Status employees). The data are protected by Title 13 of the U.S. Code: employees who disclose the identity of an individual or business are subject to a penalty of five years in jail, a \$250,000 fine, or both.