

# Programming and Use of Blaise to SAS Transformation Tool for Streamlining Processing in the Panel Study of Income Dynamics

*April Beaulé and Mohammad Mushtaq,  
Survey Research Center, University of Michigan*

## 1 Introduction

The Panel Study of Income Dynamics (PSID) is a nationally representative longitudinal study of approximately 9000 U.S. families. Since 2001, the PSID has used Blaise as its main software for the data collection instrument. Due to size and complexity of the instrument, the PSID staff developed a tool to assist with extraction and transformation of initial files. This tool, referred to as Data Extraction, Transformation and Loading (DETL), streamlines first step processing of Blaise files.

This paper will discuss the development and use of the DETL tool in the processing cycle, specifically providing an overview of the development and how we improve efficiency and save time setting up first step files for further processing.

## 2 Background

In the early years of the study, the data collection tool was a paper and pencil questionnaire. In 1992, the PSID automated the survey instrument using SurveyCraft software. By 2001, the Survey Research Center at the University of Michigan began using Blaise and the PSID was re-programmed at that time in the new software. Just as the sample continues to grow, so has the instrument. In 1968, the first year of the PSID the interview took approximately 20 minutes to administer and the staff released 447 variables. By 2009, the interview length increased to an average of 89 minutes and we have released over 5012 variables in this last wave.

In order to make the processing of first step files more efficient, the DETL tool was developed. This tool allows the user to take advantage of the Blaise meta information including question text, code frames, loop levels and frequencies in order to bundle sets of variables together in a logical way.

The end product from DETL are SAS programs that the user produces by grouping together sets of variables from the basic output tables (extracted at the block level), and sub-setting for certain types of individuals when necessary.

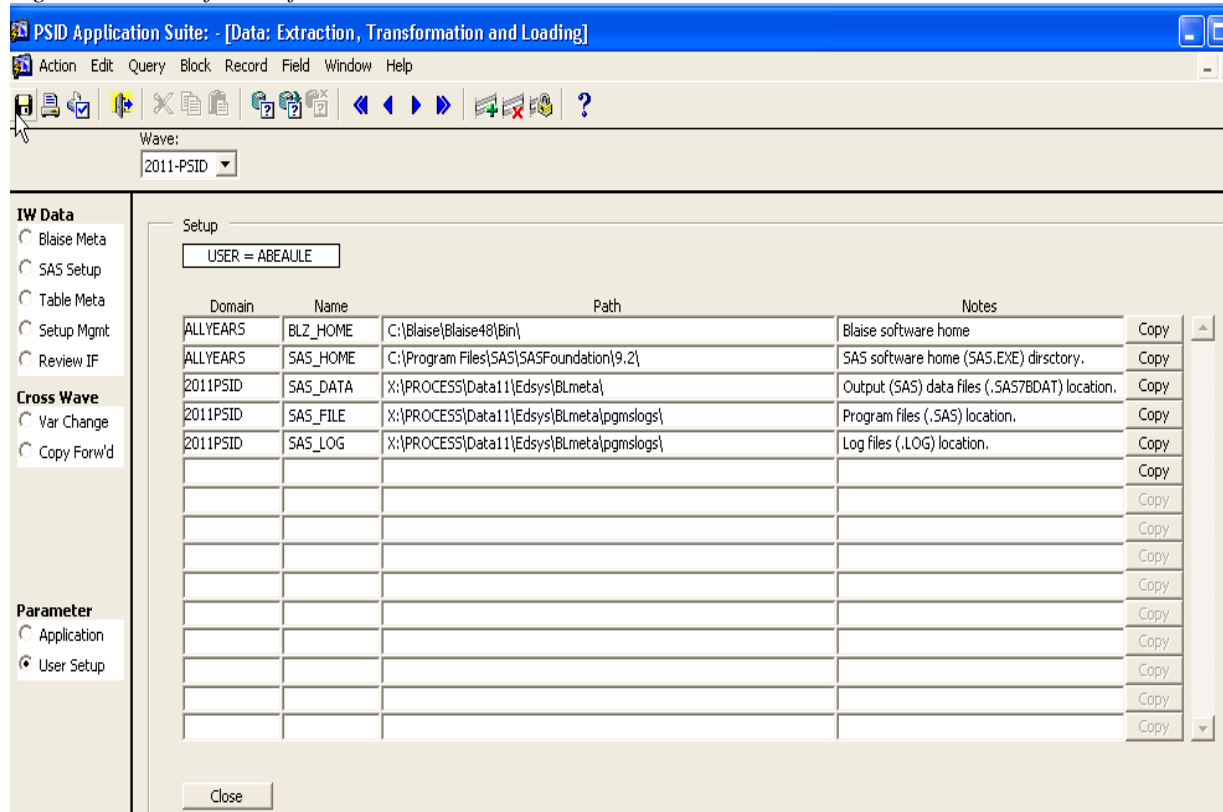
## 3 Display of Meta Information

The main purpose of first step files is to aggregate variables together at the content level and at the person level. The PSID sample is at the family level but within the instrument we ask more questions about certain people in the family, we label those 'special' people Head and Wife (if there is one). We ask fewer sets of questions about the rest of the family members and we label those individuals Other Family Unit Members (OFUM's). For Blaise programming those questions which are the same for Head and Wife and sometimes OFUM's are looped. However, for processing we have to break them apart. Using DETL we can subset the types of records we want to focus on and label the SAS setups and output resulting from the sub-setting using 'Head' or 'Wife' labels.

The basic extraction from Blaise produces tables at the block level. For Blaise programming purposes often to maximize run time speed, questions are grouped together in blocks. For processing and editing however, we need to group variables in alternative ways.

The first step for the user is to define the directories for a given year. The user defines the location where he or she wants to copy the output data as well as where the .SAS programs and .LOG files are to be written (Figure 1).

Figure 1: User Definition for DETL



Once the user has defined the parameters, they begin by reviewing the meta information about the variables and begin selecting variables for setup. The user clicks the radio button to the left under IW Data, 'Blaise Meta' (Figure 2.1). Displayed under Blaise Meta is a complete list of all the SAS files extracted by block for this BMI. The name of the output file is the name of the block as defined by the Blaise programmer. For PSID in 2009 for the main BMI, there were one hundred and twenty-three tables extracted. Variables from all these tables need to be merged together and sub-set in a coherent format for processing.

To assist the user with choosing the correct variables for each setup, the Blaise Meta screen allows the user to click on a block in the upper portion of the window (Figure 2.1). Down in the bottom pane, all the variables in this block are displayed along with their meta information. Helpful information such as the type of variable (character or numeric), the width, the number of loops and whether 'Don't Know' and 'Refused' were allowed for this variable (Figure 2.2). In the check boxes in the center of the lower pane, there are options to view the question text (Figure 2.1), the code frame, frequencies and the universe for every variable in the block selected.

Figure 2.1: Blaise Meta Display – Question Text

The screenshot shows the Blaise Meta Display interface. On the left, there are navigation options like 'IW Data', 'Blaise Meta', 'SAS Setup', 'Table Meta', 'Setup Mgmt', 'Review IF', 'Crd Wave', 'Var Change', and 'Copy Forw'd'. Below these are 'Parameter' options: 'Application' and 'User Setup'. The main area displays a list of variables with columns for varId, Name, Label, Type, Width, Dec, Loops, Min-Max, DK, RF, and Frame Universe. Variable 003, G33a\_1, is highlighted. A callout points to the 'Social Security Questions, G33 - G35' row. Below the variable list is a 'Variables' table with columns: varId, Name, Label, Type, Width, Dec, Loops, Min-Max, DK, RF, Qtext, and Freqs. Variable 010, G34, is highlighted. A callout points to the 'G34 [qText]' tab, which shows the question text: '[F1] - HELP How much was the total amount from Social Security? <<G34> w ENTER amount here, then ENTER unit of time on next screen (Month, Year)'. Another callout points to the 'Meta information - question text appears here' area.

Figure 2.2: Blaise Meta Display – Frequencies

The screenshot shows the Blaise Meta Display interface with the 'G34 [Freqs]' window open. The 'Variables' table is visible on the left, with variable 010, G34, highlighted. A callout points to the 'G34 [qText]' tab, which shows the question text. Another callout points to the 'G34 [Freqs]' tab, which displays a table of frequencies. The table has columns: Value, Count, and Percent. The first row shows a value of 1 - 150000 with a count of 2430 and a percent of 93.21. The second row shows a value of 999998 with a count of 90 and a percent of 3.45. The third row shows a value of 999999 with a count of 87 and a percent of 3.34. A callout points to the 'User clicks here to see the freqs' area.

## 4 SAS Setups

Once the user has decided on the tables and variables they need based on the meta information display, they turn their focus to defining the SAS steps. There are three basic types of SAS setups available: (a) Long, (b) Wide, and (c) Merge. The Long format setup takes a file that is collected in a wide format and flips it vertically (Figure 3). We end up with a more compact file especially when the data we collect is formatted using our standard 24 person array but we only ever have a couple of mentions.

The Wide format set up transforms in the opposite way. It takes data collected in a looped format and widens it out so that there is one record per case. For example, we use this format for mortgages. We collect a maximum of two, but we release it flattened with mortgage one variables followed by mortgage two variables. The SAS syntax that is required to produce these two transformations is lengthy. DETL is able to produce SAS code quickly, accurately with just a few clicks from the user.

Figure 3: Example of SAS Table Transformation – Wide to Long

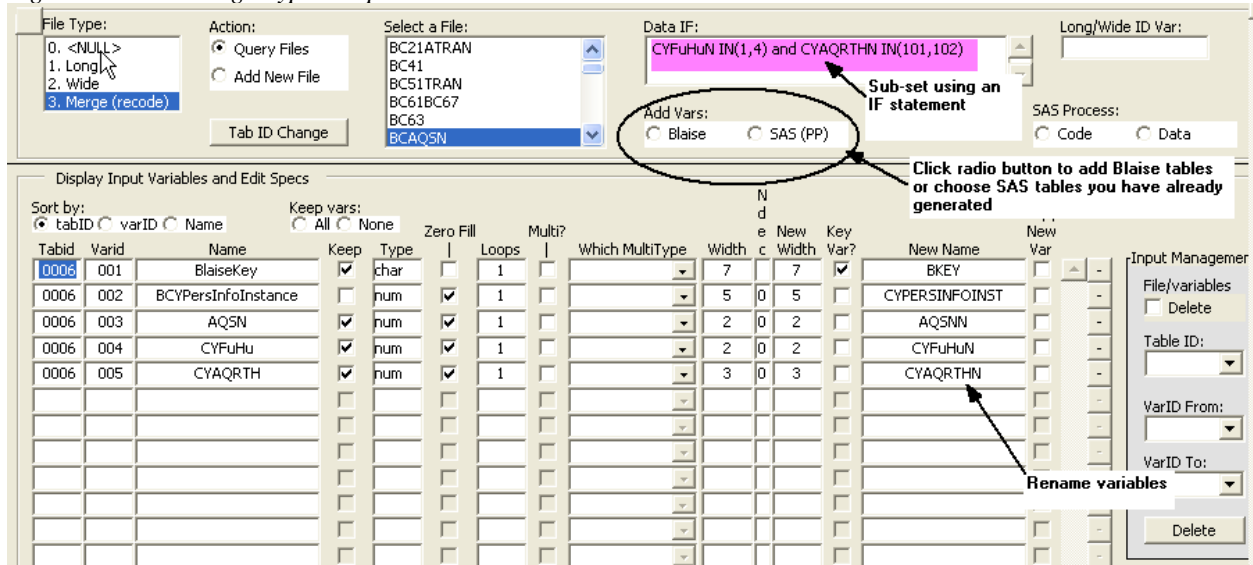
VIEWTABLE: Blio.Bdeceased								
	BlaiseKey	BDECEASEDInstance	D1	D2_01	D2_02	D2_03	D2_04	D2_05
1	0004211		1	97	1	.	.	.
2	0041134		1	1	1	.	.	.
3	0079001		1	97	1	.	.	.
4	0088003		1	97	1	2	.	.
5	0117001		1	97	1	.	.	.
6	0129004		1	1	1	.	.	.
7	0165002		1	97	1	.	.	.
8	0173002		1	97	1	.	.	.
9	0297131		1	1	1	.	.	.

VIEWTABLE: B111.D1d3					
	BKEY	DECINST	D1	i1	AQSNN
1	0004211		1	97	1
25	0041134		1	1	1
49	0079001		1	97	1
73	0088003		1	97	1
74	0088003		1	97	2
97	0117001		1	97	1
121	0129004		1	1	1
145	0165002		1	97	1
169	0173002		1	97	1
193	0297131		1	1	1
217	0301003		1	2	1
241	0393001		1	1	1
265	0418001		1	97	1
289	0419002		1	1	1

The most commonly used set up is the Merge type. The user is selecting tables extracted at the block level and merging them with other similar content variables. The Merge type set up is flexible. The user can merge together raw tables directly from Blaise or newly created DETL tables can be merged with Blaise tables (Figure 4). The user must define the merge keys between the tables. He or she can also rename variables, prefill system missing data with zeros, and subset for specific record types. One of the examples in the PSID using a merge setup is the creation of the housing variables. We ask about mortgages and foreclosures up to a maximum of two in a looped format. First we widen the loop and then we merge this wide table back with the anchor table so that we end up with one record for every case with the housing variables even if the family had neither a mortgage nor a foreclosure.

Figure 4: DETL Merge Type Setup Screen



## 5 Handling Multi-Mention Variables

The PSID has many multi-mention type variables. DETL allows them to be transformed into the output type needed quickly and easily. The first type output we may need is the standard mention order type. The question asks ‘How is your home heated – with gas, electricity, oil or what?’

Figure 5: Multi-mention question

How is your home heated--with gas, electricity, oil or what?

- ENTER all that apply
- PROBE: Any others?
- For multiple response, use space bar or dash to separate responses

<input type="checkbox"/> 1. Gas	<input type="checkbox"/> 6. Solar
<input checked="" type="checkbox"/> 2. Electricity	<input checked="" type="checkbox"/> 10. Bottled gas; propane
<input type="checkbox"/> 3. Oil	<input type="checkbox"/> 11. Kerosene
<input type="checkbox"/> 4. Wood	<input type="checkbox"/> 97. Other - specify
<input type="checkbox"/> 5. Coal	

Heating Method

In DETL, if we like to have the output as mentions, we simply mark the multi-mention type to ‘OK’ meaning ‘as is’ in mention order. But more frequently analysts do not want to see mentions displayed in order of how they were chosen but rather, they want to see if the respondent said yes to any of the options. In DETL, there is 2<sup>nd</sup> option for multi-mentions ‘Dummy’. If one chooses to transform to a dummy indicator, the output for each mention is displayed as a True/False type (Figure 6).



Figure 7: Multi-mention question – Month String

During which months of 2010 did you get this income?

- ◆ ENTER all that apply
- ◆ For multiple response, use space bar or dash to separate responses

<input type="checkbox"/> 1. Jan	<input type="checkbox"/> 7. July	<input checked="" type="checkbox"/> 13. All
<input type="checkbox"/> 2. Feb	<input type="checkbox"/> 8. Aug	
<input type="checkbox"/> 3. March	<input type="checkbox"/> 9. Sept	
<input type="checkbox"/> 4. April	<input type="checkbox"/> 10. Oct	
<input type="checkbox"/> 5. May	<input type="checkbox"/> 11. Nov	
<input type="checkbox"/> 6. June	<input type="checkbox"/> 12. Dec	

Months in 2010 Rec'd SS Ben

Here is what the SAS data looks like – transformation to a month string.

G35_M
111111111111
111111111111
111111111111
111111111111
111111111111
111111111111
111111111111
000001100000
111111111111
111111111111
000011111111
111111111111
111111111111
111111111111
111111111111
111111111111
111111111111
111111111111
000111111111
111111111111

## 6 Producing SAS Code, Log and Tables

Once the data manager is satisfied with their setups, they simply press a radio button to either generate the SAS program (so that it may be run from SAS itself) or they can choose to both generate the SAS program, run it producing the SAS table and save the log all at the same time. The program, output table and logs are all saved to the location determined in the parameter setup.

## 7 Pulling Setups Forward for the Next Wave

As a panel study, the PSID continues to have a need to be able to take advantage of all the data preparatory work in the prior wave. Because a great deal of effort goes into the DETL setups, our programmer was able to bring setups forward into the next wave; the caveat being that variables that were dropped disappear from the setups for the current wave and variables that are new need to be added to the appropriate existing setups.

DETL provides the user with critical between wave reports. The first report indicates which variables have been dropped and which variables have been added. The change report can be generated at the study, bmi or block level. A total of 580 variables were changed from main instrument (PSID\_Main.bmi) between 2009 and 2011 waves of data collection, and 687 variables were changed from all bmi's. Figure 8 below lists variables added and dropped from the foreclosure (BA37FOR) module/block. By using this detail report, the setups for the current wave can be carefully edited (Figure 8).

Similarly a second report alerts the user to width changes so that adjustments can be made in the current year setups where the variable name has remained the same but the width has been changed.

Figure 8: Variable change report – between 2009 and 2011

PSID Blaise Variable Name Change Report (PY = 2009, CY = 2011) Count = 580				
Table	Bmi	Year	Variable	CY/PY Status
BA37FOR	PSID09	2009	A37FOR2Mo	Dropped in PY.
			A37FOR2Yr	Dropped in PY.
			A37FOR3	Dropped in PY.
			A37FOR4	Dropped in PY.
			A37FOR5	Dropped in PY.
			A37FOR5SPEC	Dropped in PY.
			A37FOR6	Dropped in PY.
BA37FOR	PSID11	2011	A37B	New in CY.
			A37C	New in CY.
			A37D	New in CY.
			A37E	New in CY.
			A37F	New in CY.
			A37FSPEC	New in CY.
			A37G	New in CY.

## 8 Table Number Changes between Pretest and Production

The use of DETL starts immediately when data are delivered from test or pretest stages of data collection instrument or questionnaire. After uploading metadata information into the Oracle database, SAS setups and data files are created to analyze data from various development cycles. Based on this analysis, the Blaise data model changes significantly during an instrument development cycle: for example to add new blocks or variables or update a code frame. These structural adjustments modify the table order and thus void the prior SAS setups created by the data user (Figure 9).

Figure 9: Table number changes between stages of instrument development

If the setup is based on a single table and it has few variables, then it can be updated manually. However if the setup consists of several variables from many tables, manual editing is time consuming and is prone to error. The “Tab ID Change” function of DETL is an elegant way of making such changes. As shown in Figure 9, the “Tab ID Change” function pops up a small window

with three columns. Using Tabid to New Tabid map created from old and new bmi's, the table numbers can be changed by providing values to the New Tabid column. After applying changes, users can create new SAS setups and data files.

## 9 Miscellaneous Features

DETL collects metadata information about the setup created by the users. As mentioned earlier, users can create setups only or setups and data. The metadata information is displayed in a master-detail relationship as shown in Figure 10 below.

Figure 10: Metadata of SAS setup files

The screenshot shows the 'SAS Setup File Info' window. At the top, 'Wave:' is set to '2011-PSID'. Below this, there are radio buttons for 'File Type' with options 'Long', 'Wide', and 'Merge' (selected). The main area is split into two panels. The left panel, 'SAS Setup File Info:', contains a table with columns 'Table name', 'Nvar', 'Nobs', and 'Select'. The right panel, 'Variable Info:', shows 'Table Name' as 'A37\_B' and a table with columns 'Varid', 'Name', 'Type', 'Width', 'Ndec', and 'Key var'.

Table name	Nvar	Nobs	Select
A37_B	15	8941	<input type="checkbox"/>
BC14A14F	18	8941	<input type="checkbox"/>
BC14CTRAN	2	262	<input type="checkbox"/>
BC14FTRAN	2	262	<input type="checkbox"/>
BC18JOBTYPE	5	11487	<input type="checkbox"/>
BC21ATRAN	3	186	<input type="checkbox"/>
BC41	5	5404	<input type="checkbox"/>
BC51TRAN	3	100	<input type="checkbox"/>
BC61BC67	18	8941	<input type="checkbox"/>
BC63	6	8941	<input type="checkbox"/>

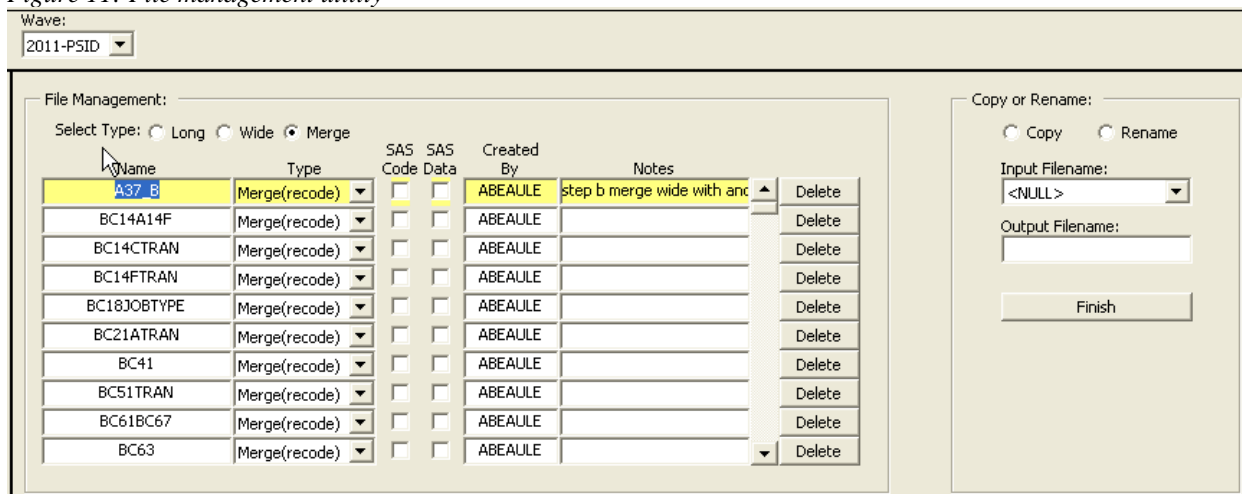
Varid	Name	Type	Width	Ndec	Key var
1	BKEY	C	7		<input checked="" type="checkbox"/>
2	A37BMO_1N	N	2	0	<input type="checkbox"/>
3	A37CYR_1N	N	4	0	<input type="checkbox"/>
4	A37D_1N	N	1	0	<input type="checkbox"/>
5	A37E_1N	N	1	0	<input type="checkbox"/>
6	A37F_1N	N	1	0	<input type="checkbox"/>
7	A37FSPEC1	C	255		<input type="checkbox"/>
8	A37G_1N	N	7	0	<input type="checkbox"/>
9	A37BMO_2N	N	2	0	<input type="checkbox"/>
10	A37CYR_2N	N	4	0	<input type="checkbox"/>

For PSID 2011, there are 203 files created by DETL, out of which 19 are long, 7 are wide and 177 are merge type/format. The number of variables (nvar) and number of observations (nobs) are displayed in the left hand side block of figure 10. The variable information, (order, name, type, width, ndec, etc.), from the selected file is displayed in the right hand side panel. This information is helpful to users as they are able to see the structure of output data tables from the same interface and make changes to the setups if needed.

DETL also offers file management utility, where file setup can be renamed or copied to a new name. User can add processing notes to a file or can delete files that are no longer needed. For further details about this interface. See Figure 11 on the next page.

Another useful feature of DETL is the collection of “data if statements” in one tab together where users can review them with proper context (i.e., along with other files) and make changes if needed. When satisfied, then one can proceed to creating SAS setup and data files again. A modified ‘if statement’ may change the number of observations. The process will update file metadata and users should check meta information before moving to the next steps.

Figure 11: File management utility



## 10 Data Model and Interface Development (Programming Details)

DETL was developed using Oracle Forms/Reports 6i and Oracle Database 9.2. The development process started in 2003 and the application went into production in early 2004. Later on, it was customized for other studies and additional features were added as needed by the data processing team.

To develop this tool, the first task is to extract metadata from Blaise data model (bmi) and merge it with metadata extracted from the dataout tables (SAS tables). During this process, the looped variables need special attention because there is a one-to-many relationship between the bmi variables and the dataout variables. Then, the next task is to develop relational data model or schema that includes tables, views, functions and procedures.

Figure 12.1 describes the relationship among Blaise-to-SAS tables. A new row is inserted into blz\_domain for each wave of data collection. Then blz\_tables, blz\_variables and blz\_fmt\_map are populated. Blz\_freqs table is updated during data extraction process and information from this table is used to create SAS setups as discussed earlier in Sections 3 and 4. Blz\_questions and blz\_universe are part of the relationship. The information from these tables is joined with blz\_variables using a one-to-many relationship between PVNAME and VNAME. Similarly, SAS format names are extracted from SAS data tables and then information is linked back to Blaise TCodes using blz\_fmt\_map and blz\_formats. It should be noted that the data model is not normalized. Social science data are problematic to model into a normalized relational database and we hope to achieve this goal during next upgrade.

Figure 12.2 describes the relationship among SAS setup tables, specified in Section 4 and further explained in Sections 5 and 6 of this paper. The SAS setup module is using an Oracle database for storage and interview data are in SAS tables. The next step is to integrate the user interface with SAS, so that users are able to send tasks to SAS and get the results back to the calling environment. The bridge between SAS and Oracle is created by a SAS/Access interface to Oracle.

As explained in previous sections, this tool creates a SAS program per user specifications and it is submitted to SAS for data processing. The application uses SAS/Macros in the background to generate SAS code and calls SAS to execute the program.

Figure 12.1: ER diagram of Blaise and dataout tables

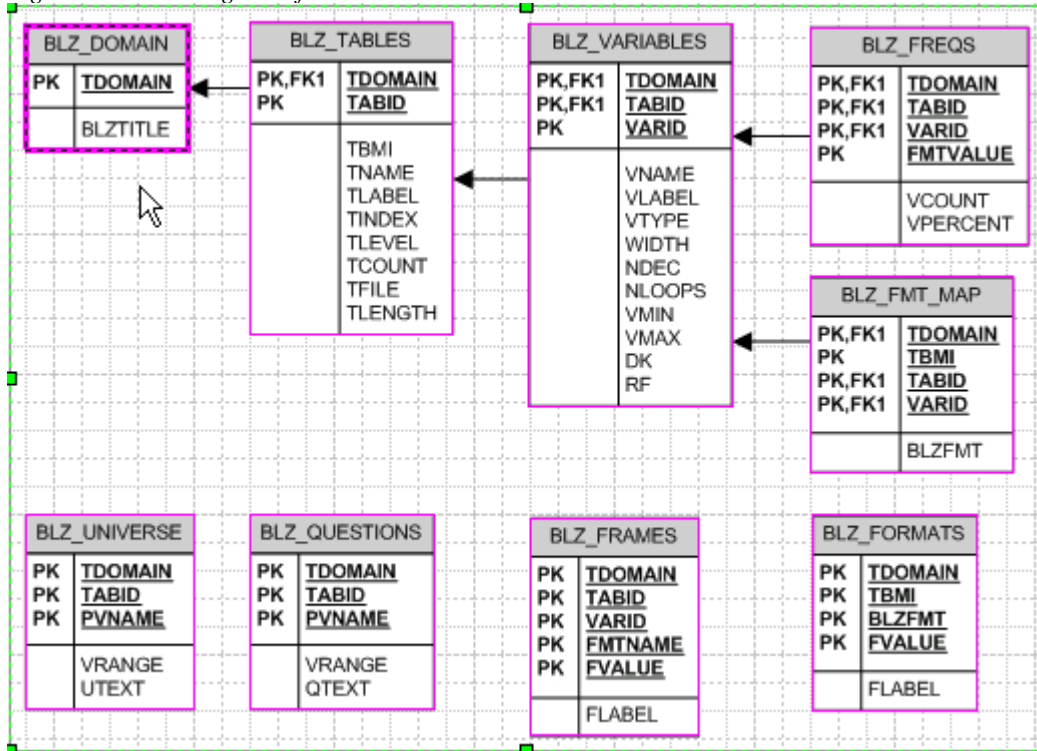
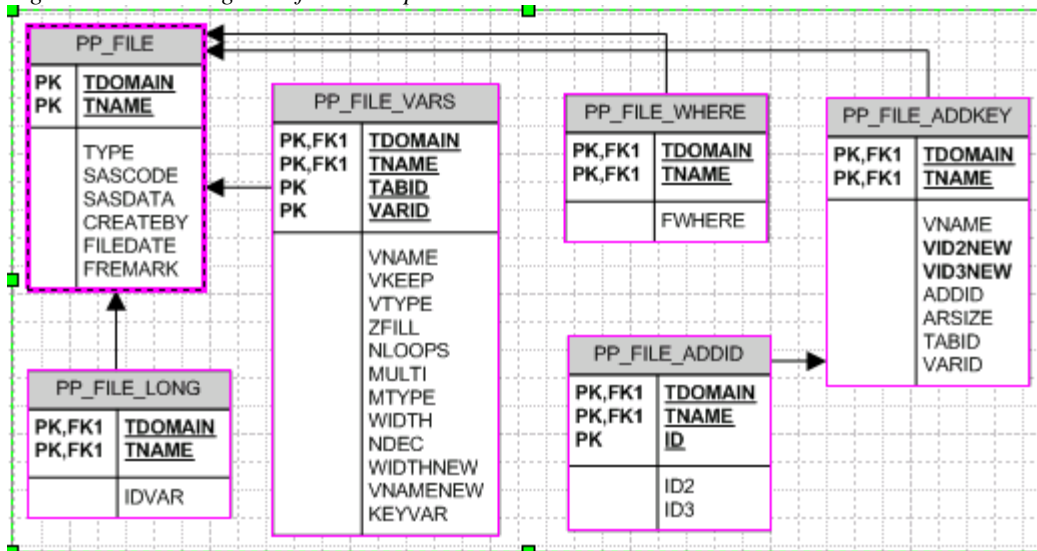


Figure 12.2: ER diagram of SAS setup tables



The SAS data, log and output are written to the folder specified in Section 3. DETL also has a function to check for errors in the log file and report them back to the calling interface.

## **11 Conclusion and Summary**

The DETL tool developed by the PSID programming staff has enabled the processing staff to methodically see all the changes that were made from one wave to the next and update first pass setups accordingly. DETL informs the user about the meta information about each variable including the loop level, the question text, the code frame and even the raw frequencies. Using this information, the user can create first step files and do basic transformations and variable renaming without having to write out the long SAS syntax to complete these.

The SAS programs generated from DETL can be run while also saving a copy of the SAS program and the log file with one click. DETL is an important tool for the processing staff that has saved time and reduced errors and confusion by providing all the Blaise meta information to the user to make informed decisions about variables, where they should be placed, and at what level.