

# Interlocking and Mutually Reinforcing System - From Shells to TLGs Automatically Generation

Zhiping Yan, Dizal Pharma

## ABSTRACT

Drawing TLG (Table, Listing and Graph) mock shell relies heavily on manual work. After spending a lot of time on generation, modification, and finalization the mock shell, statistical programmers still need to manually write repetitive code for creation of TLG.

To reduce manual work and improve efficiency, we designed a metadata driven process to automate generation of mock shell as well as TLG. By adopting Python GUI (Graphical User Interface) toolkit and NLP (Natural Language Processing), this process was implemented with two applications. With these applications, statistical programmers can create shell and final output within short time by just clicking on several buttons.

## INTRODUCTION

Scalability and reusability of metadata-driven solution used in statistical programming can help avoid inconsistencies and issues with data analysis and reporting across studies as well as reduce redundancy and unnecessary use of space. More importantly, having the metadata centralized can make programming task less time-consuming.

The working efficiency of statistical programmers can be further improved by leveraging application — Graphical User Interface (GUI) — developed by Python. The operations — like clicking on buttons, selecting combo box item — over applications can seamlessly connect the metadata system and whole TLG automation process.

## TLG AUTOMATION PROCESS

With TOC provided by statistician and global metadata files derived from global standard shell, study level XML metadata and excel file containing title and footnote can be generated. These two kinds of material can be used to generate study level mock shell. With help of NLP (Nature Language Processing) like NER (Name Entity Recognition), key words and phrases can be extracted from TLG title and mock shell. The extracted key words/phrases and SAS macro mapping metadata can be linked together with fuzzy match or machine learning to generate TLG macro parameter values. The generated macro parameter values can be filled into SAS program template and thus create study level TLG SAS program file which can generate final TLG.

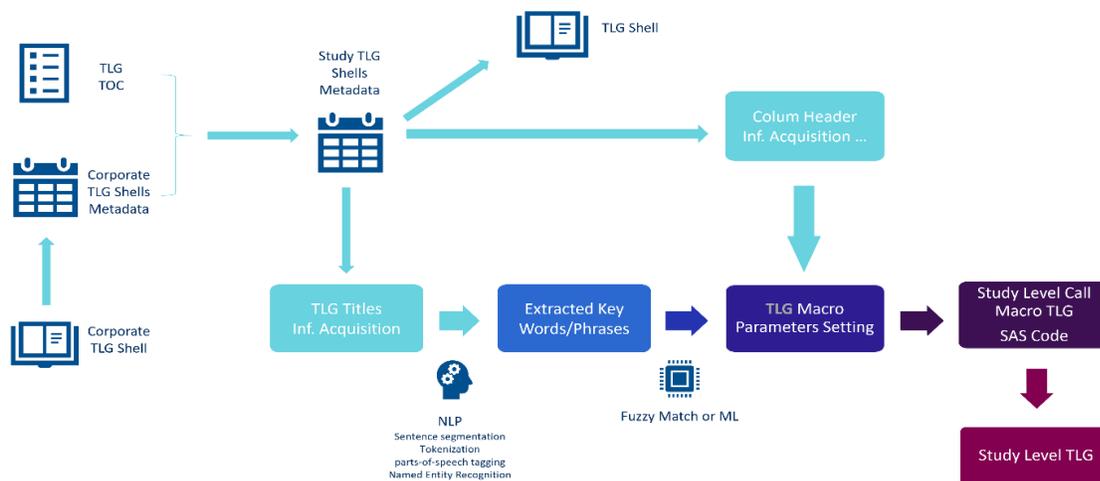


Figure 1. TLG Automation Process







- **Step 3:** By clicking buttons 'Row xml2tree' and 'Column xml2tree', overall xml files will be displayed as 'row tree' and 'column tree' in right part of GUI. Right-click menu feature as well as Delete undo/redo is provided for customization of shell. Buttons 'Row tree2xml' and 'Column tree2xml' can allow us to convert customized tree into xml files.
- **Step 4:** Initial title and footnote excel file will be generated based on global footnote metadata once 'Generate Title and Footnote Excel' is clicked. A button was also created to open this excel file for customization.
- **Step 5:** Row xml file, column xml file, title and footnote metadata file will be used to create RTF file for each TLG by clicking on 'XML2RTF'. Button 'Combine RTF' enable us to bundle individual rtf files into final TLG shell.

Repeat step 2 or 3, step 4 (no need to generate initial title and footnote metadata) and step 5 to update TLG shell once get comments from Biostatistician, Medical and other functions.

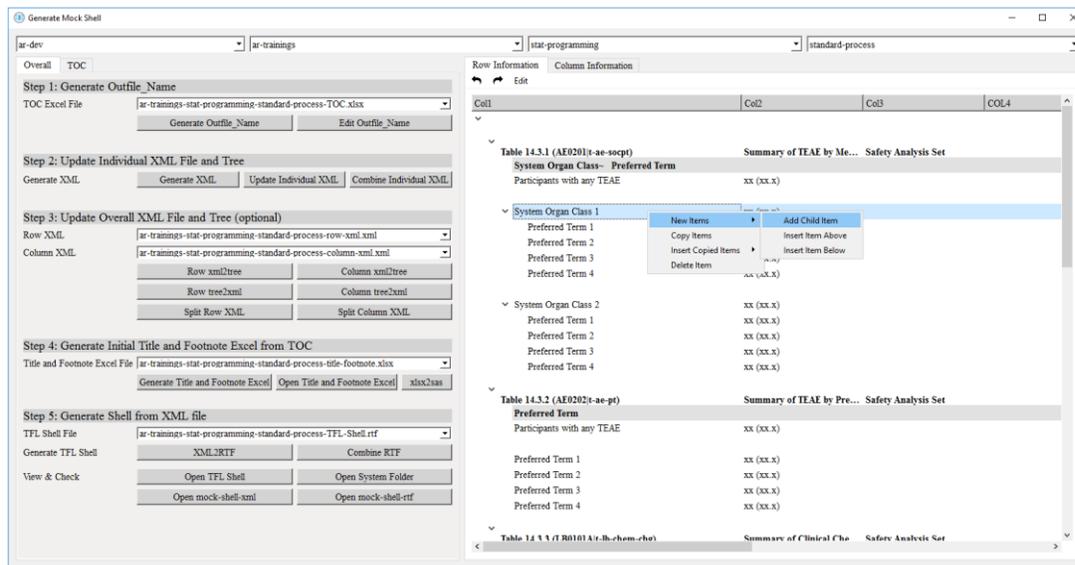


Figure 9. Generate Mock Shell Application

To make it easier to modify xml files for just one TLG, a TOC tab (see Figure 10) was also generated to contain TLG information organizing in the same way as it is stored in TLG TOC file. Right click menu on specific Outfile\_Name cell can help generate initial xml file, convert xml files into row/column tree, generate RTF and even delete xml files.

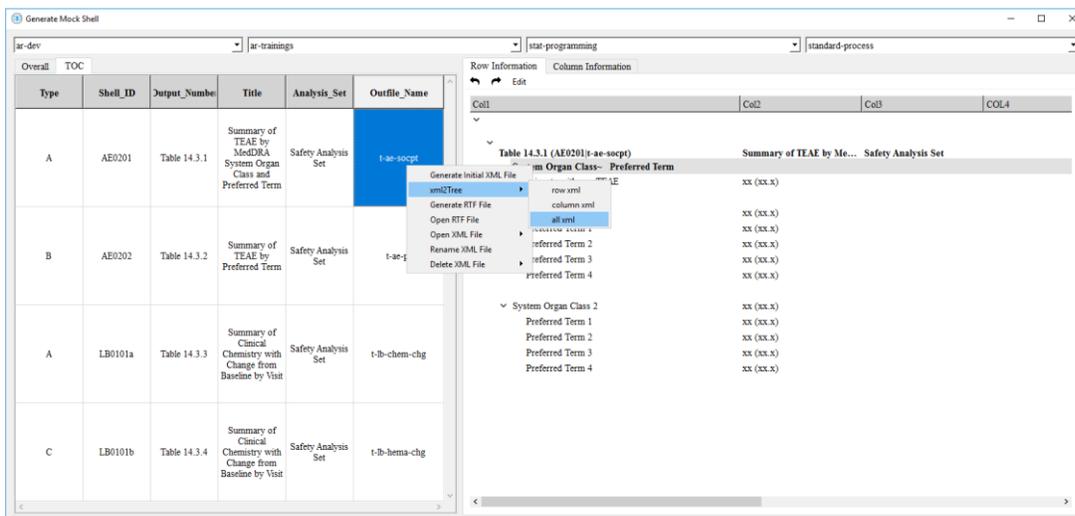


Figure 10. TOC tab of Generate Mock Shell Application

## GENERATE TLG OUTPUT

There are two types of TLG, TLG requiring many derivations and those require no or few derivations. Automatic generation of first type relies heavily on SAS reporting macros while creation of second type heavily depends on metadata file linking columns and variables together.

### TLG REQUIRING MANY DERIVATIONS

#### SAS Reporting Macro System

Besides xml metadata and footnote metadata, SAS reporting macros were also generated based on global standard TLG mockup and Shell ID (Figure 11). For example, SAS macro `m_ar_ae_socpt` was intended to develop AE summary table having shell ID like AE0201, AE0202, AE0301 and AE0302.

These reporting SAS macros have same prefix 'm\_ar\_' and developed based on same unit macros (e.g. `u_lbl2header`), analysis macros (e.g. `m_ar_count` and `m_ar_mean`) and report macro `m_ar_reports`. For efficacy TLG, SAS procedures instead of `m_ar_count` and `m_ar_mean` will be used to increase readability and understandability for reviewers.

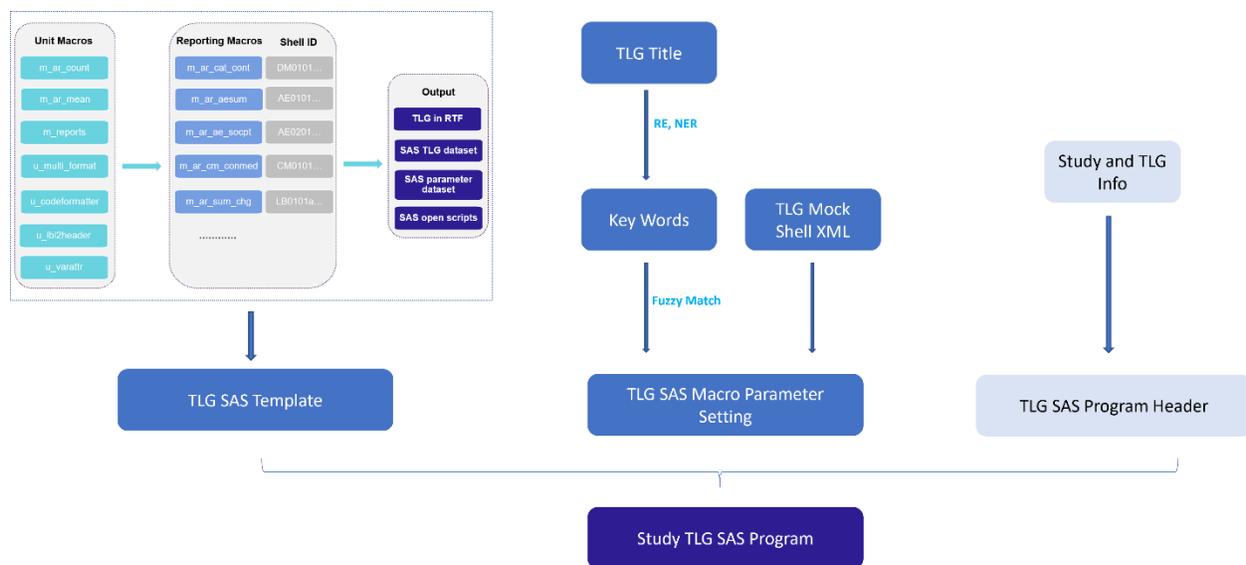


Figure 11. SAS Reporting Macro System

Outputs generated by these reporting macros include RTF file, SAS TLG output dataset, SAS parameter dataset and SAS open code. Main function of SAS open scripts is to be submitted to regulatory authorities. It can generate same TLG as parent SAS program file does. Moreover, only few sub macros like `m_ar_count`, `m_ar_mean` and `m_reports` can be found in these open scripts. Comments regarding sub macro feature and parameter usage will be added right before the code calling these sub macros. With good formatting and comments, it is very easy for reviewer to understand the code.

Multi-language is also handled by these macros. By switching value of parameter `LANGUAGE`, users can produce both English and Chinese TLG. Fixed wording will be translated directly by macros. Table content coming from ADaM or SDTM dataset should be translated at ADaM or SDTM level.

#### TLG SAS TEMPLATE

Given the fact that there is a one-to-one or one-to-many relationship between SAS reporting macro and shell ID, it is very easy to generate a SAS program template file for each TLG. Figure 12 presents a template file for table with shell ID AE0201. It consists of program header part and SAS macro part.

Information required by program header and macro parameter values can vary from table to table. To make it easy for automation, they are broken into small units enclosed in brackets. These small units can be replaced with real code or text customized for a specific TLG so that SAS program file to create this TLG output can be generated.

```

/*-----*/
Study:          [Protocol]
Program:        [program_name].sas
Keywords:      TEAE, SOC, PT
Function:       [title]
Author:        [username]
Date Completed: [date]
SAS Version:   9.4
Platform:      SUSE Linux Enterprise Server 15
Input Data:    ADSL, ADAE
Macros Called: m_ar_ae_socpt
Program Output: [outfile_name].rtf
Program log:   [outfile_name].log
Macro Parameters: None

Program Flow:
  Step 01:
  Step 02:

Limitations/Cautions:

Revisions:

/*-----*/

%macro m_ar_ae_socpt (
  ,input_dataset          = lptad.adae
  ,input_dataset_where   = saffl = 'Y' [teaeCode] [tesaeCode] [saeCode] [relatedCode] [grd3Code] [leadingCode]
  ,description_text      = Subjects with Any[relatedText] [teaeText] [tesaeText] [saeText] [grd3Text] [leadingText]
  ,trt_var               = trta
  ,level1_var            = [lv11Code]
  ,level2_var            = [lv12Code]
  ,level3_var            = [lv13Code]
  ,population_dataset    = lptad.adsl
  ,population_dataset_where= saffl = 'Y'
  ,therapy_var           =
  ,display_total         =
  ,display_trt           =
  ,grp_var               = [grpCode]
  ,display_grp           = [grpText]
  ,level1_sortby         =
  ,level2_sortby         =
  ,sortby_colnum         =
  ,bign_type             =
  ,cutoff                = [pctCode]
  ,report_yn             =
  ,page_vars             =
  ,rpt_span_headers      = [rpt_span_headers]
  ,column_width          =
  ,coll_header           =
  ,line_num_per_page     =
  ,outfile_name          = [outfile_name]
  ,title_footnote_template = lptmtsy.tlft
  ,tfl_title             =
  ,tfl_footnote          =
  ,footnote_on_last_page =
  ,output_path           =
  ,output_pgm_path       = [output_pgm_path]
  ,output_datlib         =
  ,output_dataset        =
  ,layout_orientation    =
  ,font_size             =
  ,language               = [language]
  ,destination           = [destination]
  ,help                  =
  ,debug                 =
):

```

Figure 12. SAS Program Template

### EXTRACT INFORMATION FOR TLG TITLE AND XML FILES

It is not hard to get information like protocol number, program name, TLG title, program author and substitute them for small units in program template file. But for other small units like '[teaeCode]', '[tesaeCode]' and '[lv11Code]', the only way to get related information is to check TLG title and mock shell (xml files in our case). At Dizal, we applied regular expression and rule-based matching of Spacy package to extract Named-entity recognition (NER) which can be used to derive small units. Below shows the code regarding how to perform NER for AE summary table title <sup>[2]</sup>.

The extracted group of words by Spacy package were assigned labels like byGroup, leadingGroup, relatedGroup, pctGroup, cycGroup, teaeGroup and so on. For title 'Summary of TEAE by MedDRA System Organ Class and Preferred Term', we can get three NERs which are 'TEAE', 'MedDRA System Organ Class' and 'Preferred Term'. Their entity names are **teaeGroup**, **lv11Group** and **lv12Group**.

```

title = 'Summary of TEAE by MedDRA System Organ Class and Preferred Term'

ae_expression_dict = {
'relatedGroup': r"(\w+-related|\w+-induced)",
'pctGroup': r"(((\^|\\(\*ESC)\D+\d{4,})\D+|\W+)\s*\d+\s*(percent|pct|%))",
'grd3Group': r"(grade 3 or higher|grade\s+\W+\s*\d)",
'cycGroup': r"(Cycle\s+\d+\s+or\s+Cycle\s+\d+\s+Day\s+\d+|Cycle\s+\d+\s+Day\s+\d+)(Day\s+\d+)",
'teaeGroup': r"(Treatment Emergent Adverse Event|Treatment-Emergent Adverse Event|Treatment Emergent
AE|Treatment-Emergent AE|TEAE|TEAEs)",
'tesaeGroup': r"(Treatment Emergent Serious Adverse Event|Treatment Emergent Serious AE|TESAE|TESAEs)",
'saeGroup': r"(Serious Adverse Event|SAE)",
'aeGroup': r"(adverse event|AE)"
}

ner_text = []
ner_name = []

# Regex
title1 = title
for key in ae_expression_dict:
    match_lst = re.findall(ae_expression_dict[key], title1, flags=re.IGNORECASE)
    if match_lst is not None:
        for match in match_lst:
            if isinstance(match, str) == True:
                match = match
            else:
                match = match[0]
            ner_text.append(match.strip())
            ner_name.append(key)
            title1 = title1.replace(match, "")
            title1 = re.sub(' +', ' ', title1)

nlp = spacy.load("en_core_web_sm", disable=['ner'])
# Entity for by Group
bygroup_pattern = [
{"LOWER": "by"},
{"POS": "{}", "OP": "*"},
{"POS": {'NOT_IN': ['ADP', 'PUNCT']}, "OP": "*"},
{"TEXT": {'REGEX': '(?i)(by|,|of)', "OP": "?"},
{"POS": {'NOT_IN': ['ADP', 'PUNCT']}, "OP": "*"},
{"TEXT": {'REGEX': '(?i)(by|,)', "OP": "?"},
{"POS": {'NOT_IN': ['ADP', 'PUNCT']}, "OP": "*"},
{"POS": {'NOT_IN': ['ADP', 'VERB', 'PART', 'CCONJ', 'PUNCT']}]
]
patterns_bygroup = [{"label": "byGroup", "pattern": bygroup_pattern}]
ruler_bygroup = EntityRuler(nlp, overwrite_ents=True)
ruler_bygroup.add_patterns(patterns_bygroup)
ruler_bygroup.name = 'ruler_bygroup'

# Entity for Leading to Group
leadingGroup_pattern = [
{"LOWER": "leading"},
{"LOWER": "to"},
{"TEXT": {'REGEX': '(?i)(death)', "OP": "?"},
{"POS": {'NOT_IN': ['ADP']}, "OP": "?"},
{"TEXT": {'REGEX': '(?i)(of|from)', "OP": "?"},
{"POS": {'NOT_IN': ['ADP']}, "OP": "*"},
{"POS": {'NOT_IN': ['ADP', 'VERB', 'PART', 'CCONJ']}]
]
patterns_leading = [{"label": "leadingGroup", "pattern": leadingGroup_pattern}]
ruler_leading = EntityRuler(nlp, overwrite_ents=True)
ruler_leading.add_patterns(patterns_leading)
ruler_leading.name = 'ruler_leading'

.....

nlp.add_pipe(ruler_bygroup)
nlp.add_pipe(ruler_leading)

# Entity Ruler
doc = nlp(title1)
for ent in doc.ents:
    ner_text.append(ent.text.strip())
    ner_name.append(ent.label_)

```

Figure 13. Code to Extract Key Words from TLG Title

## Macro Mapping Specification

With the macro mapping specification showed in Figure 14 and fuzzy match of fuzzywuzzy package, the extracted text and labels can be derived as small units. For example, label 'teaeGroup' can derive two small units - '[teaeCode]' and '[teaeText]'. Their respective values are **and strip(uppercase(trim))** = 'Y' and **TEAE**. Label 'lvl1Group' and text 'MedDRA System Organ Class' can get small unit '[lvl1Code]' and corresponding value is '**aebodsys**'. While label 'lvl2Group' and text 'Preferred Term' can get value of '[lvl2Code]' assigned as 'aedecod'.



Shell_ID	Output_Type	Output_Number	Title1	Population	Outfile_Name	SEQ	Col_Header	Sub_SEQ	Sub_Col_Header	Source_Dataset	Main_YN	Source_Variable	Data_Typ	Format	Codelist	Sort_YN	Group_YN	Line_Brk_YN	Column_Width
AE9301	L	16.2.7.1	All Adverse Events	Safety Analysis Set	I-ae	1	Participant ID	1	Participant ID	LPTAD ADAE(where = (raffl = "Y"))		SUBJID	Char			Y	Y	Y	9
AE9301	L	16.2.7.1	All Adverse Events	Safety Analysis Set	I-ae	2	Sex/Age*(years) Y/Race	1	Sex			SEX	Char				Y		9
AE9301	L	16.2.7.1	All Adverse Events	Safety Analysis Set	I-ae	2	Sex/Age*(years) Y/Race	2	Age (years)			AGE	Integer						9
AE9301	L	16.2.7.1	All Adverse Events	Safety Analysis Set	I-ae	2	Sex/Age*(years) Y/Race	3	Race			RACE	Char						9
AE9301	L	16.2.7.1	All Adverse Events	Safety Analysis Set	I-ae	3	Term/Preferred Term/System Qrgan/Class	1	AE Term			AETERM	Char						9
AE9301	L	16.2.7.1	All Adverse Events	Safety Analysis Set	I-ae	3	Term/Preferred Term/System Qrgan/Class	2	Preferred Term			AEDECOD	Char						9
AE9301	L	16.2.7.1	All Adverse Events	Safety Analysis Set	I-ae	3	Term/Preferred Term/System Qrgan/Class	3	System Organ Class			AEBOOSY S	Char						9
AE9301	L	16.2.7.1	All Adverse Events	Safety Analysis Set	I-ae	4	(Day)/End Date (Day) ('ESC')super	1	Start Date (Day)			AESTDTC. ADY	Char	YYMMDD01 0.		Y			9
AE9301	L	16.2.7.1	All Adverse Events	Safety Analysis Set	I-ae	4	(Day)/End Date (Day) ('ESC')super	2	End Date (Day)			AEENDTC. ADY	Char	YYMMDD01 0.		Y			9

Figure 16. Example of Listing Metadata

For listings not covered by corporate level mock shell, source dataset and variables as well as datatype, format can be populated by finding most similar label from ADaM and SDTM spec comparing against column header using NLP. Below presents code about how to map source variables.

```
import spacy
import re
import numpy as np
import pandas as pd

spec_df .....
lst_df .....

nlp = spacy.load('en_core_web_lg')
label_lst = spec_df['LABEL'].tolist()

for index, row in lst_df.iterrows():
    col_header = str(row["SUB_COL_HEADER"])
    col_header = re.sub('{.+}', '', col_header)
    col_header = re.sub('\(.+\)', '', col_header)
    search_txt = nlp(col_header)

    sims= []
    for lbl in nlp.pipe(label_lst):
        sims.append(search_txt.similarity(lbl))

    id_max = np.argmax(sims)
    matched_label = label_lst[id_max]

    if sims[id_max] > 0.6:
        matched_row = spec_df.loc[spec_df['LABEL'] == matched_label,]
        matched_variable = matched_row['VARIABLE'].values[0]
        matched_datatype = matched_row['DATATYPE'].values[0]
        matched_format = matched_row['FORMAT'].values[0]

        lst_df.loc[index, 'SOURCE_VARIABLE'] = matched_variable
        lst_df.loc[index, 'DATA_TYPE'] = matched_datatype
        lst_df.loc[index, 'FORMAT'] = matched_format

        label_lst.remove(matched_label)
    else:
        lst_df.loc[index, 'SOURCE_VARIABLE'] = ''
        lst_df.loc[index, 'DATA_TYPE'] = ''
        lst_df.loc[index, 'FORMAT'] = ''
```

Figure 17. Code to Map Source Variable with Column Header

## Convert Listing Metadata into Program File

After manually check and update the metadata, listing metadata file can be converted into SAS program file by calling SAS macro `u_shell_meta2code`. Below shows example metadata converted program file. The code is well structured. Even if required source information is not filled, similar code can also be generated. Users only need to update Process input dataset section and thus still saves a lot of time for statistical programmers.

```
/*-----
Study:          ar-trainings-stat-programming-standard-process
Program:        l-ae.sas
Keywords:       l-ae
Function:       This program is used to generate Listing 16.2.7.1 All Adverse Events (core)
Author:
-----
Limitations/Cautions:
Revisions:
-----*/
*** Read input dataset;
data indat1;
  set lptad.adae(where=(saffl = "Y"));
run;
-----
*** Process input dataset;
data indat2;
  length col1 col10 col11 col12 col2 col6 col7 col8 col9 col31 col32 col33 col41 col42 col43 col51 col52 col53 $200 col3 $600. col4 $600. col5 $600. ;
  set indat1;
  ***col1: [Cohort/~Treatment]***;
  col1="";

  ***col2: Participant-ID***;
  col2 = strip(subjid);

  ***col3: Sex/Age-(years)/Race***;
  col31 = strip(sex);
  if missing(age) = 0 then col32 = strip(put(age, best.));
  else col32 = ' ';
  col33 = strip(race);
  col3= catx('\line ',strip(col31), strip(col32), strip(col33));
run;
-----
*** Derive final dataset;
data final;
  set indat2;
run;
-----
proc sort data=final; by aestdct aeendtc; run;
-----
*** TFL title and footnote;
data tlaft;
  SHELL_ID = "AE9901S" ;
  OUTPUT_TYPE = "L" ;
  SECTION_TITLE1 = "" ;
  SECTION_TITLE2 = "" ;
  OUTPUT_NUMBER = "16.2.7.1" ;
  TITLE1 = "All Adverse Events (core)" ;
run;
-----
*** Macro purpose: Construct PROC REPORT and generate RTF file
-----
%macro reports(
  ,input_dataset          = final
  ,rpt_vars               = col1 col2 col3 col4 aestdct aeendtc col5 col6 col7 col8 col9 col10 col11 col12
  ,rpt_headers            = %nrquote([ [Cohort/~Treatment] | Participant-ID| Sex/Age-(years)/Race] AE Term/~Preferred Term/~System Organ-Class (*ESC*)(super a)|
  ,column_width           = 5 5 9 9 5 5 5 5 9 9 5 5
  ,column_alignment       = 1 1 1 1 1 1 1 1 1 1 1 1
  ,header_alignment       =
  ,sort_vars              = aestdct aeendtc
  ,noprint_vars           = aestdct aeendtc
  -----
  ,outfile_name            = l-ae-pa
  ,title_footnote_template = tlaft
  -----
  ,language                = EN
  ,destination             = A
  ,help                    = N
  ,debug                   = N
);
```

Figure 18. Example Listing Program File Created from Metadata

## CDARS APPLICATION

To improve working efficiency, we integrated the process of reviewing mockup, generating program file, reviewing program file, submitting program and opening RTF output into one application – CDARS.

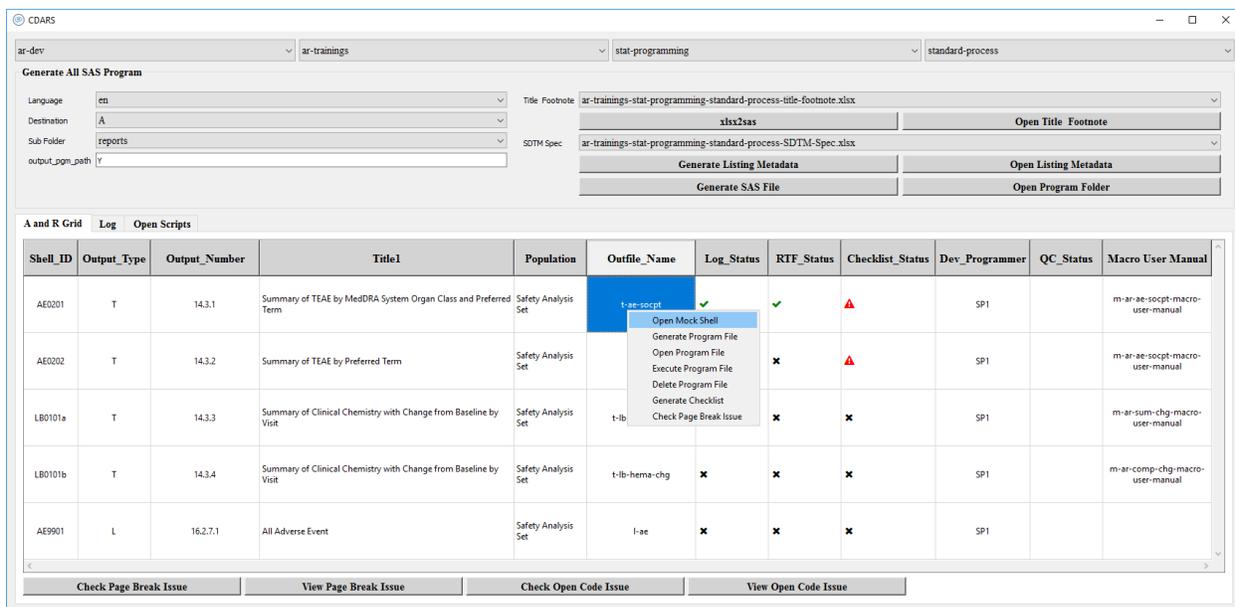


Figure 19. CDARS Application

By selecting title and footnote excel, Shell\_ID, Output\_Number, Table Title, Outfile\_Name and other information will be extracted and displayed in table as showed in Figure 19. Log status, RTF status and Checklist status can also be displayed. Reviewing QC status can be completed by pressing right click on column header of 'QC\_Status' column.

By clicking on 'Generate SAS Program File', initial SAS program containing program header and macro code will be generated for each TFL based on Shell\_ID. Already existed program file will not be over-written. By clicking on cells under 'Log\_Status', 'RTF\_Status' and 'Checklist\_Status' columns, corresponding files can be opened for review.

Another two feature about page break issue and open code issue is also integrated.

## CONCLUSION

This paper introduced how to achieve the task of TLG Auto Generation at Dizal. This process has been used by several studies and has been proved to have a great ability to increase efficiency. With more and more data recycled from studies, we believe that the productivity can be further improved.

## REFERENCES

1. Chenggeng Tian. "Decomposition and Reconstruction of TLF Shells - A Simple, Fast and Accurate Shell Designer" PharmaSUG China 2018.
2. Sumesh Kalappurakal, Srikanth Ramakrishnan, Shobha Rani, Harry Chen. "Metadata-based Auto-Programming Process – Part 3: Virtual AI assistant to automate TLG generation" PhUSE US Connect 2020

## ACKNOWLEDGMENTS

I would like to thank Huadan Li and all statistical programming team members at Dizal for their advice and feedback so that the process can be improved continuously.

## CONTACT INFORMATION

Your comments and questions are valued and encouraged. Contact the author at:

Zhiping Yan  
 Dizal Pharma  
 +86-13691423469  
 zhiping.yan@dizalpharma.com