# Metadata-based Auto-Programming Process – Part 2: Map human language to SAS® code by Natural Language Processing (NLP)

Sumesh Kalappurakal, Janssen R&D, New Jersey, US;

Qian Zhao, Harry Chen, Janssen R&D, Shanghai, China;

## ABSTRACT

Traditionally the process for programming ADaM datasets is cumbersome and relies heavily on manual work. Per regulatory requirements clinical programing algorithms should be clearly defined in the analysis specification documents in natural language(human-readable). Programmers spend most of the time developing or updating SAS® code according to specification documents. By adopting Machine Learning and leveraging the power of NLP we could analyze human-readable text from the specification documents, train the machine to convert defined algorithms to metadata and map them to the core pieces of SAS® code.

This paper is part 2 of Metadata-based Auto-Programming Process[1], and it shares an approach to automatically generate SAS® code to create ADaM datasets from source SDTM datasets via metadata and NLP methodology.

The strategy would be to extract key information from defined algorithms written in human language and existing code, populate metadata and then utilize the metadata to generate code.
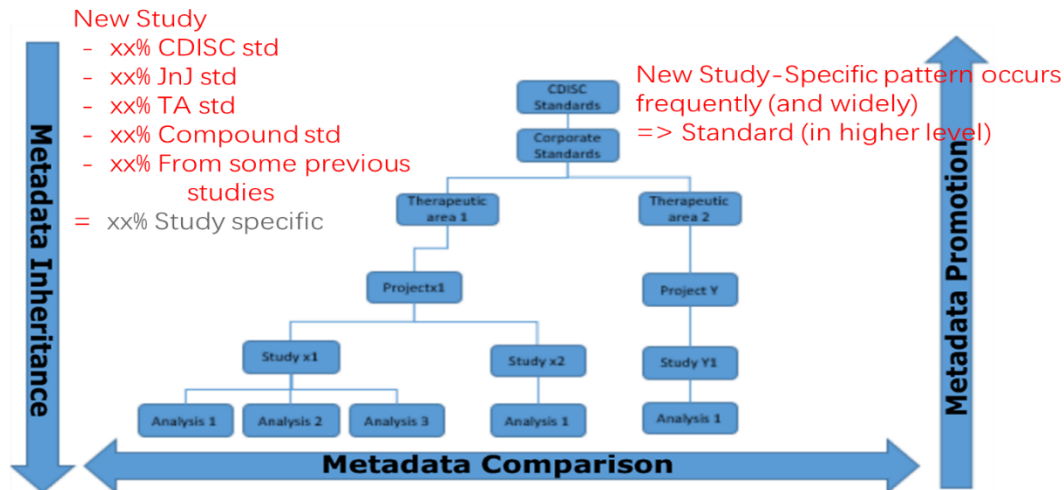
## INTRODUCTION

According to Gartner, "by 2020, natural-language generation and artificial intelligence will be a standard feature of 90% of modern business intelligence platforms"[3]. Following the trend, adopting Machine Learning and leveraging the power of NLP we can analyze human language, train the machine to convert defined algorithms to metadata and map them to the core pieces of code (e.g., SAS® or R) to automate the programing process.

We started an innovative project called Autocode in programming team. Autocode project's main objective is to automate generation of SAS® code for analysis datasets and analysis reports. The vision is to apply new technology to analyze the data structure, digitize the specification document and SAS® code and store them in the form of metadata into a database. Advantages of this model would be more automation, standardization and re-usability study after study. This project is still in the prototype stage and we will share our design, progress and experience in this paper.

## FIRST MAIN TOPIC

### METADATA DRIVEN APPROACH

With increased enforcement from regulatory agencies towards standardized data and the industry moving towards standardizing analysis requirements, we look to achieve more automation by reusing standard and previously used study algorithms stored in metadata to generate code for analysis.
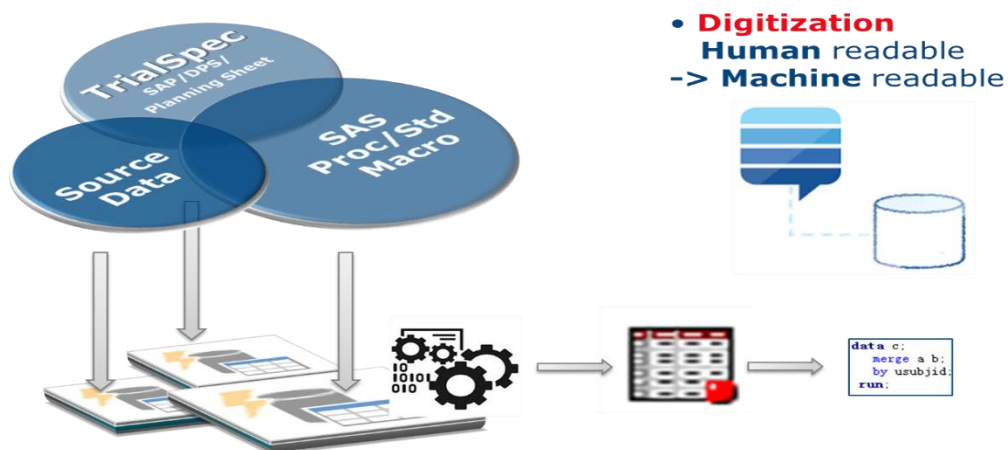
Source: PhUSE 2017 Paper SI02 Clinical Metadata – Metadata Management with a CDISC Mindset

**Figure 1. Metadata management**

## MODULES

Our model is comprised of three different modules to analyze data, SAS® code and document specifications. We use SAS® to calculate the metadata summary for source data and output data. We are using regular expressions and exploring the powerful parser generator ANTLR[10] (Another Tool for Language Recognition) to parse the SAS® code automatically and NLP methodology to analyze specification documents to calculate document similarity and extract key pieces to find the algorithm patterns.
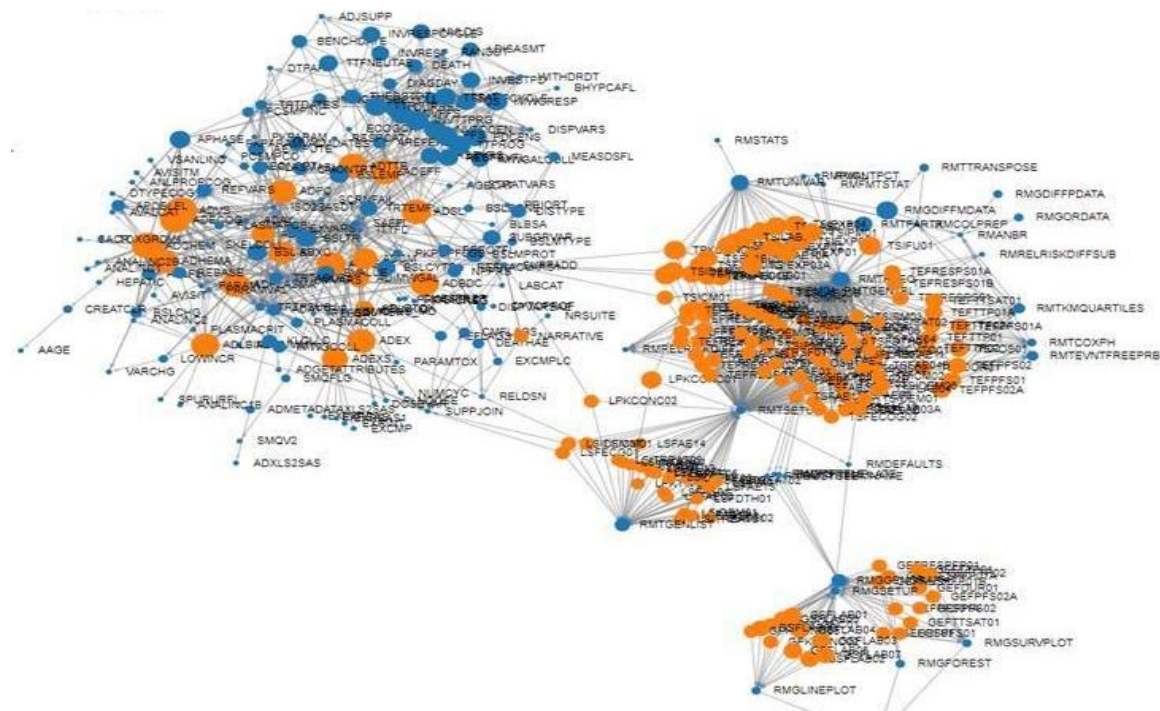


**Figure 2. Modules to analyze data, SAS® code and specifications**

The goal of the three modules is to extract key information from different sources for data analysis and make them machine readable in structured data, from which we can generate SAS code for analysis. Here are the 3 modules:

1. **PDE (PowerDataExplorer)[2]**: This module is used to generate metadata for SDTM & ADaM datasets, including data elements, structure/hierarchy, relationship, changes/differences across

multiple datasets



**Figure 3. PDE Module**

2. **PME (PowerMacroExplorer):** This module is used to create a dataset which contains program/macro dependency, macro parameters and macro calls in historical use cases by scanning and parsing codes across multiple study folders

**Figure 4. PME Module**

3. **ADaM Planning Sheet analyzed by NLP:** The Excel specification includes key information such as SAS macro name (SASMCR) and text descriptions of the derivation algorithm (ORIGCOM) necessary to guide programmer to create SAS code for the variable derivations. Generally, ORIGCOM columns may be any combination of plain English text, a formula, pseudocode, etc. It describes the source of the data or the derivation clearly and it serves as the submission-ready text for the source/derivation/comments in define.xml. This paper focus on derivation and we will introduce how to analyze the algorithm text in ORIGCOM using NLP to map it to macro metadata and generate SAS code. At this stage we focus on the planning sheet, but a similar approach can be applied to other document specifications like protocol and Statistical Analysis Plan (SAP).

| DATASET | VARNAME | VARLABEL | ORIGCOM | DERVCOPY | SASMCR |
|---|---|---|---|---|---|
| ADAE | AENTMF | Analysis End Time Imputation Flag | Set to 'M' if minute is missing and imputed. Set to 'H' if hour and minute are missing and imputed. | Derived | %ADAEDT |
| ADAE | ASTDY | Analysis Start Relative Day | ASTDY = ASTDT – ARFSTDT + 1 if ASTDT >= ARFSTDT, or ASTDT – ARFSTDT if ARFSTDT is after ASTDT. If the dates are missing ASTDY will be missing. | Derived | %ADAEDT |
| ADAE | AENDY | Analysis End Relative Day | AENDT–ARFSTDT+1 if on or after ARFSTDT, otherwise, AENDT–ARFSTDT. | Derived | %ADAEDT |
| ADAE | ADURN | AE Duration (N) | ADURN is the duration of the AE in days from onset date to AE stop date. AE. AEENDY – AE. AESTDY + 1 . If using the numeric dates AENDT – ASTDT + 1, then check that the | Derived | %ADADURN |
| ADAE | TRTEMFL | Treatment Emergent Analysis Flag | If ADSL. ARFSTDT=<ASTDT=<ADSL. TRTEDT + [xx] days then TRTEMFL= 'Y'. See supplementaldefinitions.pdf for derivation of the flag in case of missing dates. | Derived | %ADTRTEMFL |
| ADAE | AESEVN | Severity/Intensity (N) | Code AE. AESEV to numeric value. Low intensity should correspond to low value. | Derived | %ADSEV |

Instructions | ADM Dataset | ADAE | ADM Codelist | Release History
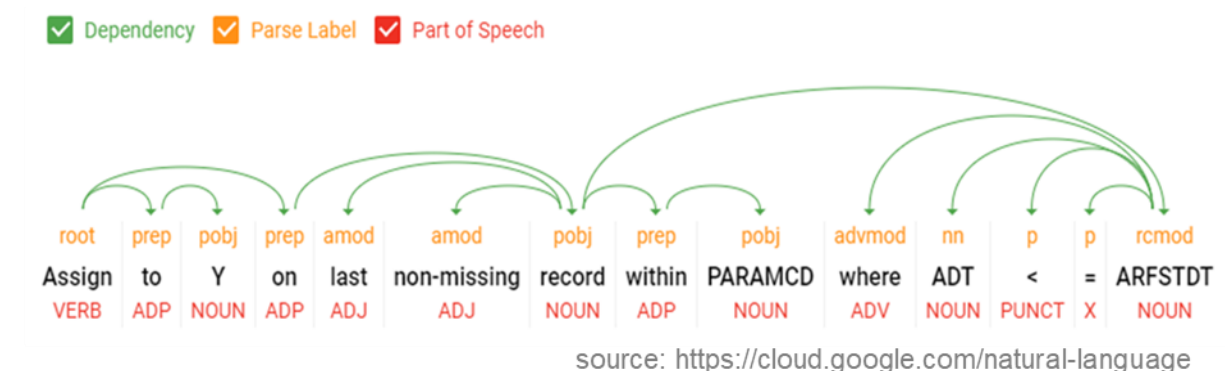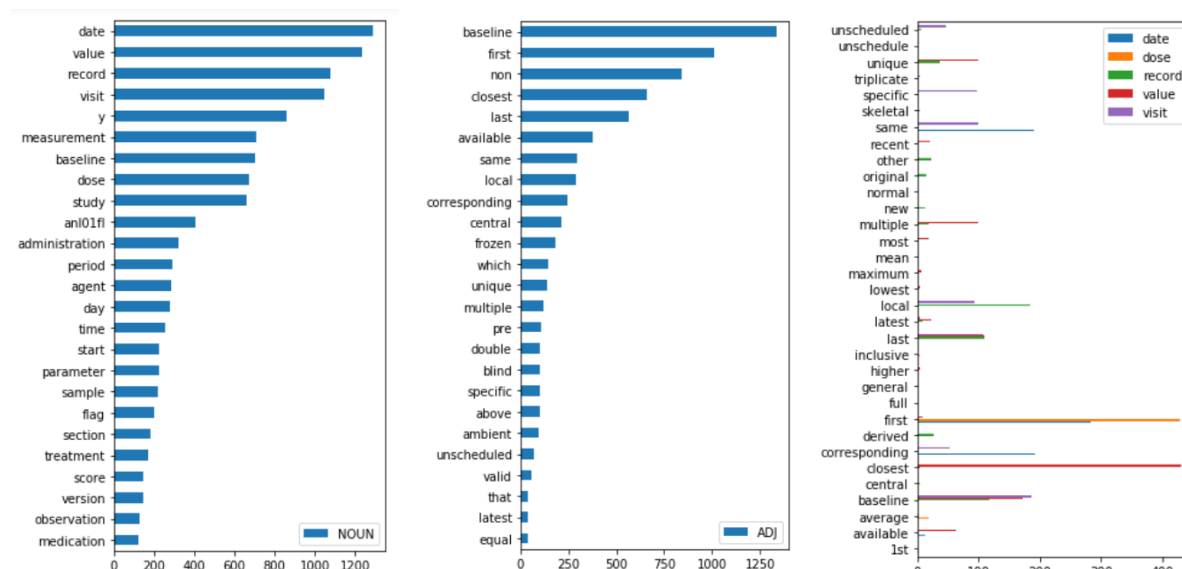
**Figure 5. ADaM Planning Sheet example**

## Natural Language Processing (NLP)

Natural language processing (NLP) is a branch of artificial intelligence that helps computers understand, interpret and manipulate human language. NLP draws from many disciplines, including computer science and computational linguistics, in its pursuit to fill the gap between human communication and computer understanding[4]. NLP is important because it helps analyze large volumes of textual data efficiently and resolve ambiguity in language to structure a highly unstructured textual data source. Below is a simple example of syntax analysis to discover structure of the algorithm text.



source: https://cloud.google.com/natural-language

**Figure 6. Syntax analysis from NLP**

We have loaded 600,000+ records of variable derivation algorithms into our prototype Metadata Repository (MDR) system collected from the planning sheets across 60+ compounds and 1000+ study folders. Taking the variable ABLFL (Baseline Record Flag) derivation for example, there are 3000+ records of algorithms and 380 unique algorithm descriptions. After some pre-processing steps to normalize (lemmatization and removing some stop words) the texts, we did basic syntax analysis to get some idea of the key nouns, adjectives and their relationships in ABLFL derivation, like the words 'first', 'last', 'closest', 'multiple', 'average' etc. that play important role to describe the algorithms to select the right value from data.



**Figure 7. Syntax analysis for variable ABLFL derivation**

There are two approaches developed to map the human-language defined algorithms to code. The <u>unsupervised approach</u> is to calculate the document similarity of algorithm sentences to get the corresponding standard/historical codes. The <u>supervised approach</u> is to split algorithm sentences to phrases and map phrases to code pieces in a training database and then incorporate the RNN (Recurrent Neural Networks)[11] model to learn the pattern and generate more flexible code.

.

- **Unsupervised approach:** For new data analysis algorithms we search the database based on text similarity score to get the most similar standard/historical algorithms and codes. We applied two approaches to calculate the text similarity, Fuzzy string-matching algorithm based on edit distance, and Word Mover's Distance (WMD) [12] algorithm based on semantic distance, and we also created the algorithm clusters from the similarity matrix.

  1. Fuzzy String Matching, also called Approximate String Matching, is the process of finding strings that approximatively match a given pattern. The degree of closeness between two strings is measured using Levenshtein Distance[14], also known as edit distance which basically is based on counting the number of primitive operations (insertion, deletion, and substitution) necessary to convert the string into an exact match. We use the fuzzywuzzy[6] package to do the calculation:

```python
from fuzzywuzzy import fuzz

textA = "'Y' if index(ATRTRF, 'BEFORE/DURING') or ASTDT=AP01SDT"
textB = "if index(atrtrf,'BEFORE/DURING') or astdt =ap01sdt then ablfl='Y'"

fuzz.ratio(textA.upper(),textB.upper())/100
```
```
0.82
```

**Figure 8. Text similarity score using fuzzywuzzy**

  2. Word Mover's Distance (WMD) algorithm is a method that allows us to assess the "distance" between two documents in a meaningful way, even when they have no words in common, via vector embeddings of words. Word embeddings (word vectors) are numeric representations of words, usually generated via dimensionality reduction on a word cooccurrence matrix for a large corpus. These vectors are used to calculate semantic similarity between words and documents:



**Figure 9. WMD algorithm via word embedding**

We use the library textacy[7] which is built on top of spaCy[8] package to compute the word mover's distance and get the text similarity score

6

```
import textacy.similarity
import spacy
nlp = spacy.load('en_vectors_web_lg')

textA = "Defined as Closest Non-Missing Value prior to the first treatment date"
textB = "Defined as nearest non-Missing value before the first dosing date"

textacy.similarity.word_movers(nlp(textA), nlp(textB), metric=u'cosine')
```
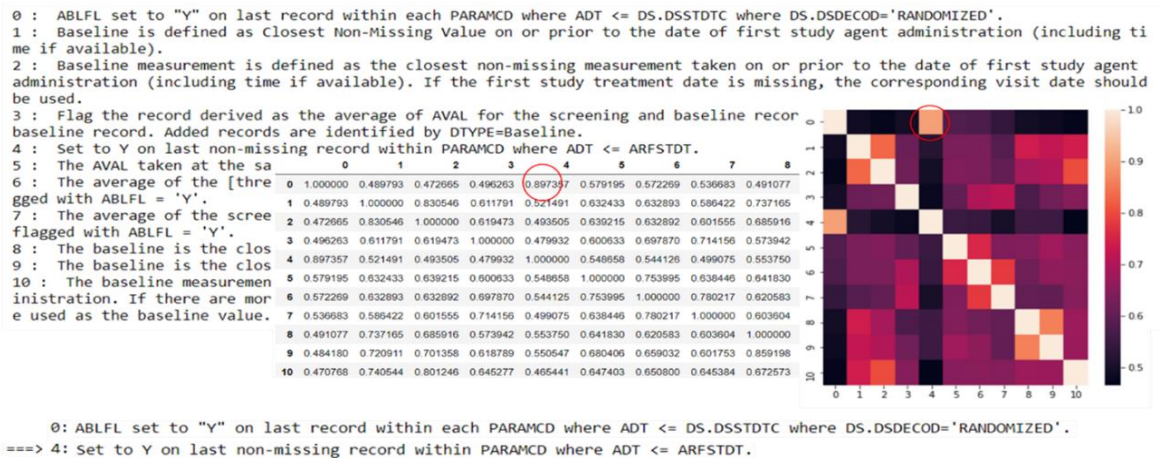
0.8109184230589997

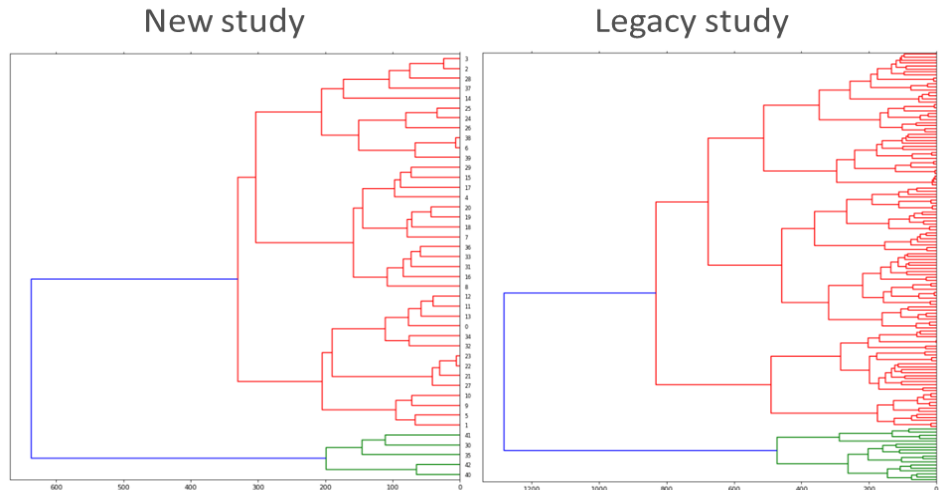**Figure 10. Text semantic similarity score using WMD**

We loaded pre-trained 'en_vectors_web_lg'[13] model which provides 300-dimensional vectors for over 1 million terms of English, and we can also create our own word embeddings to cover more domain specific vocabularies in our industry if we have more training data.

3. For a group of algorithms, we calculate similarity score for each pair of algorithm texts and generate the similarity matrix, based on which we are able to find the most similar pairs (or top 3 most similar ones), also hierarchical clustering dendrogram plots are created to show algorithm clusters.



**Figure 11. Text similarity matrix**

We have noticed that the algorithm descriptions in new studies are becoming more and more standardized than that in the legacy studies, as shown below, which lays great foundation to more automation over time.

7

**Figure 12. Comparison between new study and legacy study**

- **Supervised approach:** Small pieces of algorithms are much easier to be matched and reused than long sentences and we can map the algorithm to code based on pieces. The NLP software NLTK[9] and spaCy are used to split the algorithm documentation to meaningful pieces, and the powerful parser generator ANTLR is used to parse the SAS® code (both macro calls and data steps), and then the language pieces and code pieces would be mapped and combined to build the training database prepared for the RNN (Recurrent Neural Networks) deep learning model, based on which machine can learn the algorithm pattern and so map the new algorithm to code automatically and more flexibly. Below are the 5 steps and steps 3-5 are in a proof of concept stage.

1. Split algorithm texts into pieces using NLP tools (NLTK and spaCy) with manually reviews:

```
# Extract Keywords using noun chunks
text = str(Origcom)
doc = nlp(text)
keywords = Counter()
for chunk in doc.noun_chunks:
        keywords[chunk.lemma_] += 1

keywords.most_common(50)

[('the average', 3),
 ('reference date', 3),
 ("ablfl = ' y", 3),
 ('the baseline', 3),
 ('y', 2),
 ('the date', 2),
 ('first study agent administration', 2),
 ('time', 2),
 ('the screening', 2),
 ('the close non - miss value', 2),
 ('last record', 1),
 ('each paramcd', 1),
 ('where adt < = ds.dsstdtc', 1),
```

```
grammar = r"""
  NP:
    {<.*>+}              # Chunk everything
    }((<WP|WDT|WP$|WRB|IN>)+){      # Chink sequence
  """
```

#WDT wh-determiner which

#WP wh-pronoun who, what

#WP$ possessive wh-pronoun whose

#WRB wh-abverb where, when

#IN (Prepositions and Subordinating conjunctions)

Set to Y on last non-missing record within PARAMCD where ADT <= ARFSTDT.
=> ['set to y', 'last non-missing record', 'paramcd', 'adt < = arfstdt .']

**Display 1. Split texts into pieces with NLTK and spaCy**
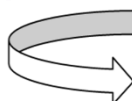
2. Parse codes using ANTLR and extract key parameters:

8

```
proc sort data=_adablflbl;
    by usubjid paramcd prev adt adtm;
run;

data _outdsn (drop=prev);
    set _adablflbl;
    by usubjid paramcd prev adt adtm;
    if last.prev and prev=1 then ablfl='Y';
run;
```

**Display 2. Code parsing with ANTLR**

3. Map algorithm text pieces to macro name and parameters in the training data



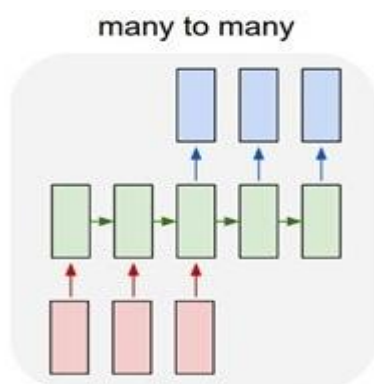**Display 3. Map into the macro name and parameters in the training data**

4. Train an RNN model



**Display 4. Training of RNN model**

5. Complex algorithm decomposed into pattern algorithms

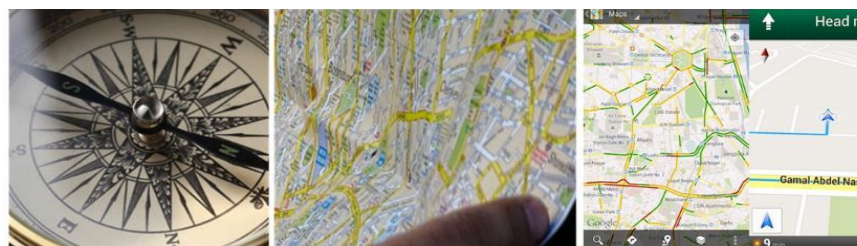**Display 5. Decompose complex algorithm into pattern algorithm**

For the deep learning model, we need to prepare large amount of training data (language pieces with corresponding code pieces). This is challenging because we need high-quality data, and this involves lots of texts and codes preprocessing and manual review, compared to the unsupervised learning process.

## CONCLUSION

We have designed and explored Metadata-based Auto-Programming Process with NLP to convert three different parts (data, SAS® code and document specifications) into machine readable metadata and connect them to try to achieve as much automation as possible in code generation for data analysis. From our experience, the analysis standard, high quality training data and technology are key factors for the success. Technological progress follows a pattern of exponential growth and there will be more and more powerful Machine Learning and NLP technology available to our daily work over time, so the top priority for us is to focus on the business side and collect 'big' standardized data, or Machine Learning cannot help as "garbage in, garbage out".

The benefits of Metadata-based Auto-Programming Process are:

- Top-Down Standard end-to-end metadata management and project management

- Metadata-driven SAS code generation to reduce manual work and potential human error which inevitably leads to wasted time in both coding and debugging

- Significant reduction in time to submission and increase in consistency and quality

- Accumulate and maintain internal knowledge in database and system

- Leave programmers more time for more complex tasks, in depth disease area understanding, more time for broadening skillsets



Exploration / Map / Auto-Navigation
Open-Code / Standard/ Auto-Code
Experience / Documentation/ System automation

In a fast-changing world, technology is changing how we live. Like Auto-Navigation brings convenience in our daily life, compared to the map and compass age. Humans needs keep-up with the fast pace technology advancement.

We are only just scratching the surface when it comes to the uses of AI and Machine Learning in the Pharmaceutical industry. However, even at this early stage, the technologies are proving to be tremendously promising when it comes to giving new mechanistic insights to disease and thereby helping to identify promising targets.

The technology will also help in terms of the industry's selection of patients for clinical trials and enable companies to identify any issues with compounds much earlier when it comes to efficacy and safety. So, the industry has much to gain by adopting AI and Machine Learning approaches. It can be used to good effect to build a strong, sustainable pipeline of new medicines[5].

## REFERENCES

1. Metadata-based Auto-Programming Process https://www.lexjansen.com/phuse-us/2018/si/SI10.pdf
2. Power Data Explorer (PDE) - Data Exploration in an All-In-One Dynamic Report Using SAS & EXCEL https://www.lexjansen.com/phuse-us/2018/dv/DV13.pdf

3. Gartner, Augmented Analytics Is the Future of Data and Analytics, Rita L. Sallam, Cindi Howson, Carlie J. Idoine, 27 July 2017
4. https://www.sas.com/en_us/insights/analytics/what-is-natural-language-processing-nlp.html
5. https://www.drugtargetreview.com/article/15400/artificial-intelligence-drug-discovery
6. https://github.com/seatgeek/fuzzywuzzy
7. https://github.com/chartbeat-labs/textacy
8. https://spacy.io
9. https://www.nltk.org
10. http://www.antlr.org
11. https://en.wikipedia.org/wiki/Recurrent_neural_network
12. http://www.cs.cornell.edu/~kilian/papers/wmd_metric.pdf
13. https://spacy.io/models/en#en_vectors_web_lg
14. https://en.wikipedia.org/wiki/Levenshtein_distance