# SDTM Mapping Based on Natural Language and Machine Learning Models

CJUG 2019

Sam Tomioka, Director of Clinical Data Programming Data Science SUNOVION Pharmaceuticals 02.07.2019



# **Brief Introduction**

- Sam Tomioka (When I am outside of Japan)
- Director, Clinical Data Programming Data Science
- Clarity Pharma -> Dainippon -> Dainippon Sumitomo -> Sunovion



# **AGENDA**



Problem to Solve



SDTM Mapping with "Machine Learning"



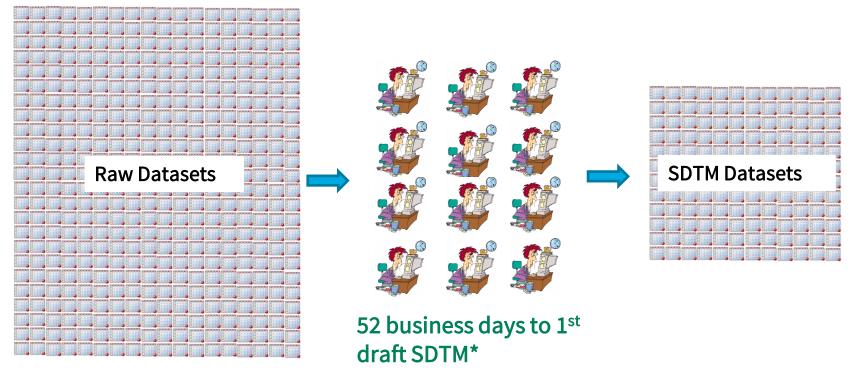
Tools Used



Thought

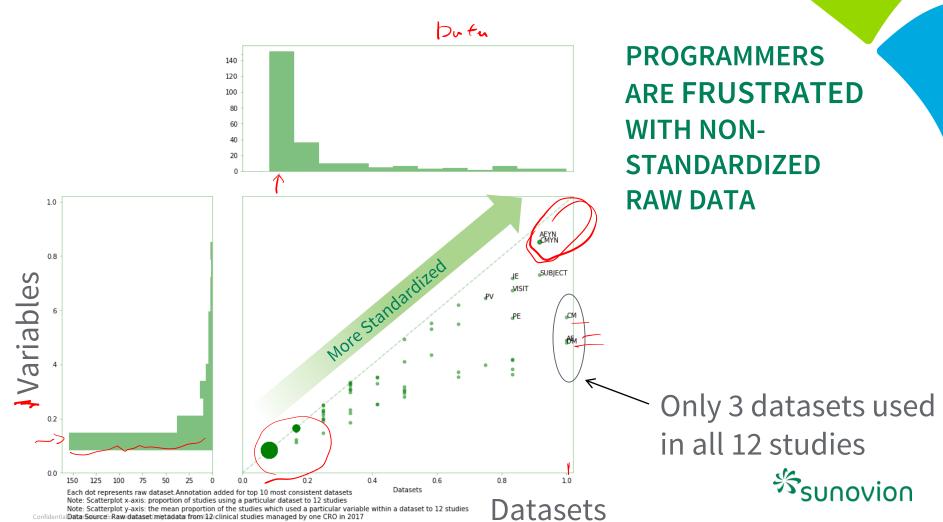


#### A PROBLEM WE WANT TO SOLVE



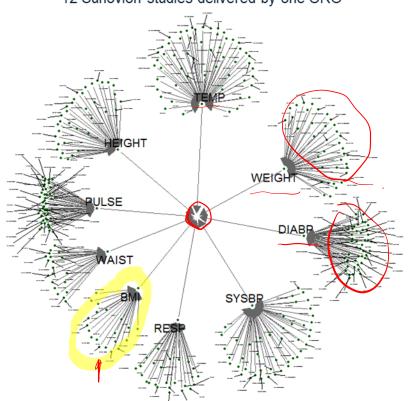


Reuse of datasets and variables across 12 studies



#### **ENDLESS MAPPING.**

SDTM.VS.VSORRES mapping for 12 Sunovion studies delivered by one CRO



#### 20 sources for BMI

"VS.BMI"

"VS.BMI\_RAW"

"VS.BMI\_Z"

"VS.BMI\_Z\_RAW"

"VS.BMIS"

"VS.BMIS\_RAW"
"VS1.BMI"

. v 21. DMT

"VS1.BMI\_RAW"

"VS1.BMI\_Z"

"VS1.BMI\_Z\_RAW"

"VS1.VS1BMI"

"VS2.BMI"

"VS2.BMI\_RAW"

"VS2.D\_BMI"

"VS2.D\_BMI\_RAW"

"VS2.VS2BMI"

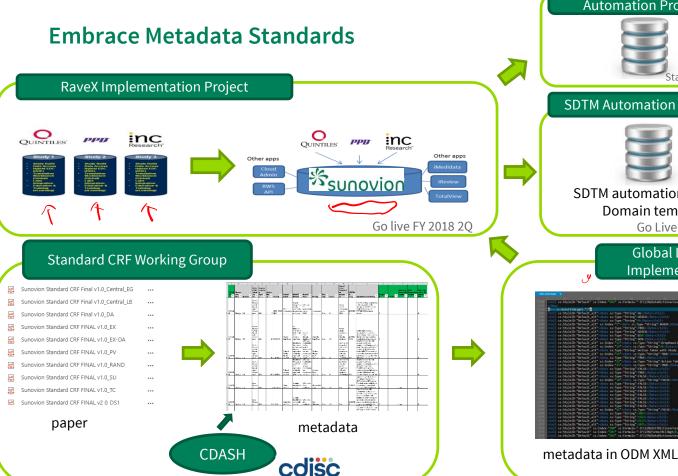
"VSMSTR.BMI"

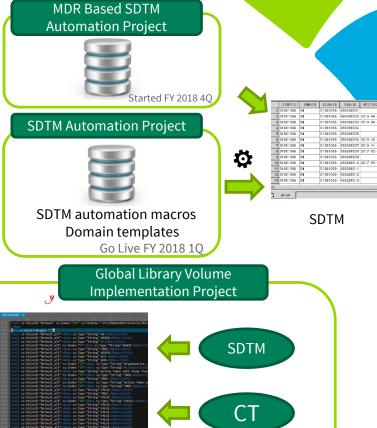
"VSMSTR.BMI\_RAW"

"VSMSTR.D\_BMI"

"VSMSTR.D\_BMI\_RAW"







cdisc

Started FY2018 10

# WHAT NEXT?

Can I use **natural language** model and **machine learning** algorithms to map raw data variables to SDTM variables?



# **AGENDA**



Problem to Solve



SDTM Mapping with "Machine Learning"



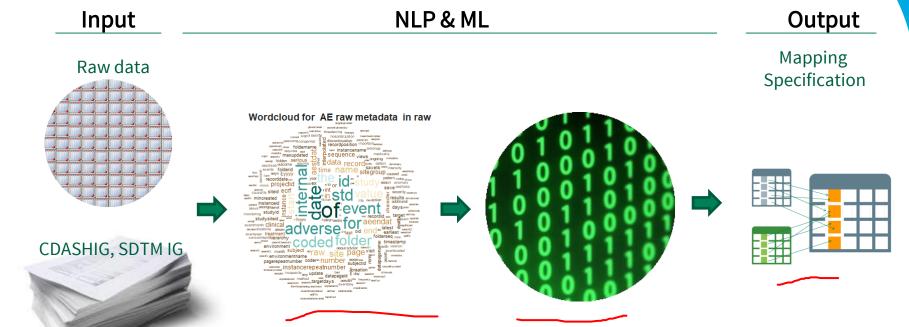
Tools Used



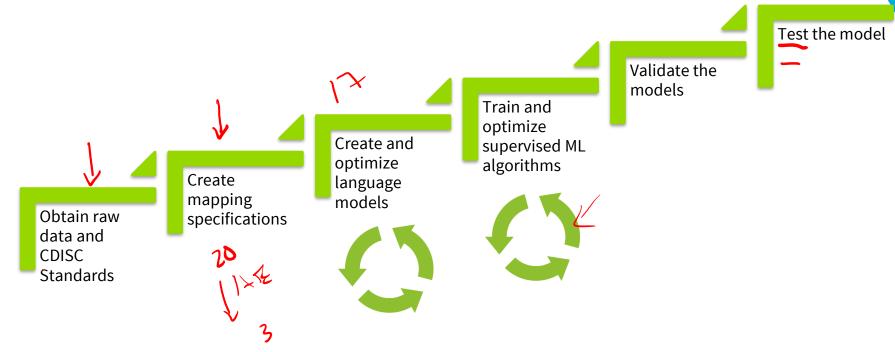
Thought



# ML based SDTM mapping for <u>fast</u>, <u>accurate</u>, <u>consistent</u> SDTM generation



#### **Steps**



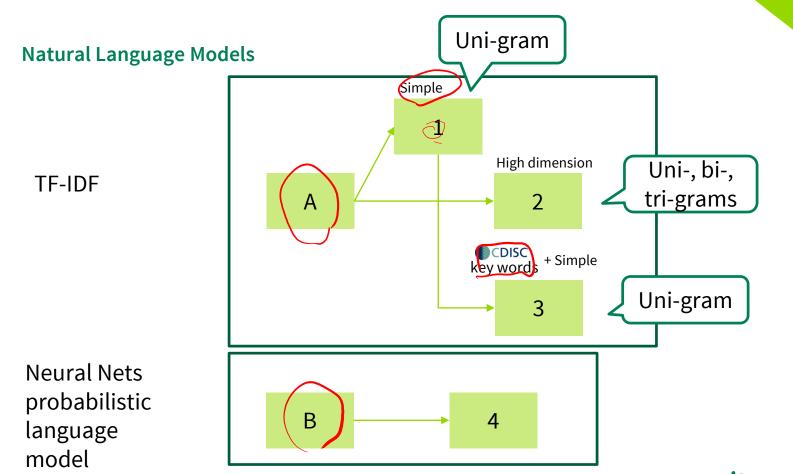


#### **Create mapping specifications (with Human Intelligence)**

Raw Variable	SDTM Variable
PT	AEDECOD •
SOC	AEBODSYS
PTNAME	AEDECOD
SOCNAME	AEBODSYS ,
•••	

<sup>\*</sup>illustration purpose only







# Natural Language Model A

TF-IDF algorithm: Weighing terms

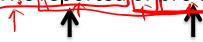
- Words occur nearby frequently are important than words that only appear once or twice Frequency (TF)
- Words that are too frequent are not important Inverse Document Frequency (DF)
- Weight

$$idf_t = \log_{10} \left( \frac{N}{16} \right)$$

$$w_{t,d} = \mathbf{tf}_{t,d} \times \mathbf{idf}_t$$

AETERM definition in CDASH IG

The reported or pre-specified name pfthe adverse event.



0.6331503 0.6105753

0.6247964 0.6611213 Logistic regression, accuracy=0.64058



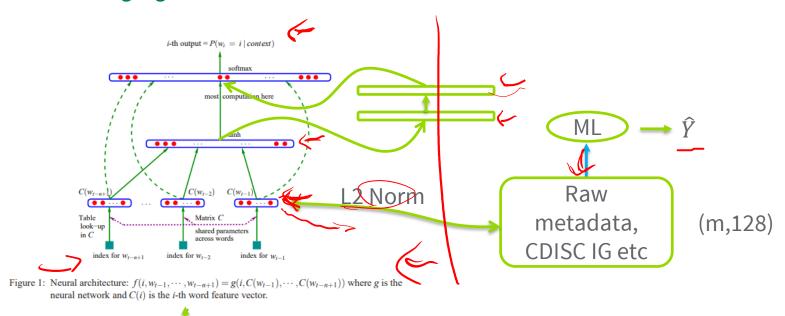


if count(t,d) > 0

otherwise

Document term matrix of 17 studies data and IG

#### Natural Language Model B



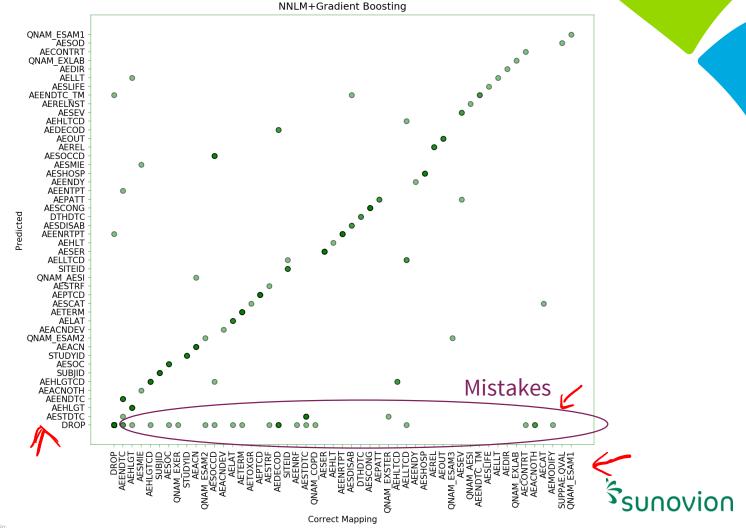
200,000,000,000 English Google News corpus

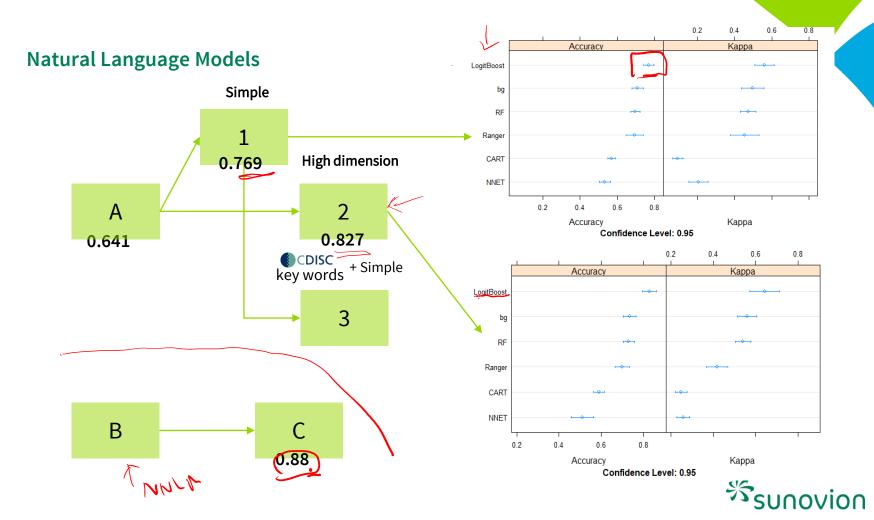
15

#### Natural Language Model B

```
The reported or pre-specified name of the adverse event. = [[ 0.4712333
     -0.08311215 0.51585186 0.1620522
      0.08530827 0.01127403 -0.09444343
                                      0.06110086
      0.04609384 -0.01650539 0.04840343 0.07840265 -0.04859402 -0.12693526
      0.04102192 -0.08049966
                           0.00896453 -0.04968296 -0.04114896
                            0.03986294
     -0.07853805 -0.01386436 0.1349262
      0.11551834 -0.06626088
                           0.01692365 0.0140045
      0.00335585 -0.01931612 0.08324809
                                      0.01950045
               0.09413652 -0.0340781
                                      0.07892839 -0.11236215 -0.17992431
      0.04012968 0.0132846
                            0.09931402 -0.09699554
     -0.11555596 -0.10244549 -0.24568063 -0.14376119 0.01996353
     -0.00676261 -0.01049919 0.09959473 -0.21953513 -0.06734764
     -0.14200029 -0.10907198 -0.06779519 0.04904051 -0.15821819 0.14623989
      0.01035561 -0.01564043 0.00844536 0.06118653 -0.12851241 -0.09161683
     -0.16958985 -0.10710518 -0.01085521 -0.1233172 -0.08650532 -0.05034445
      0.04382536 -0.12378055 0.17017077 -0.04361667 -0.00717157 0.16402934
      0.16050048 -0.19754027 0.0992441 0.00550133 -0.0244873
      0.05600908 -0.0494738 ]]
```

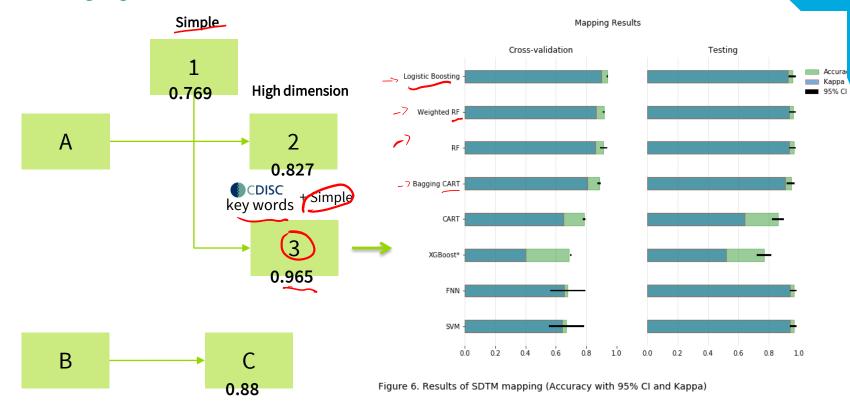
This model yield <u>0.88</u> accuracy with a gradient boosting from 10 fold cross validations



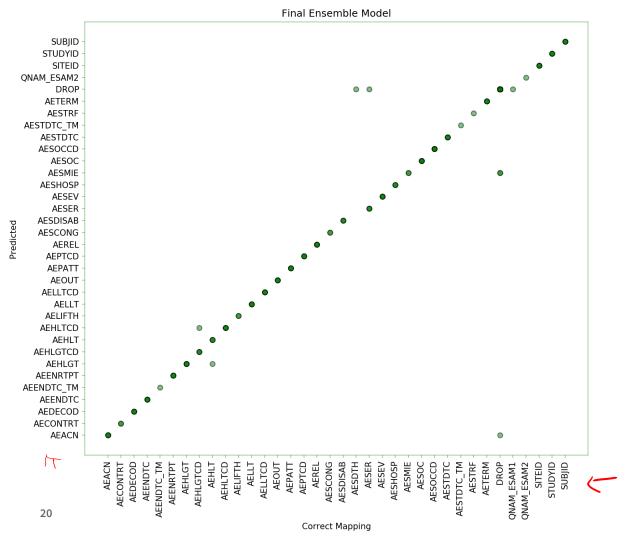


#### **Machine Learning Models**

#### **Natural Language Models**

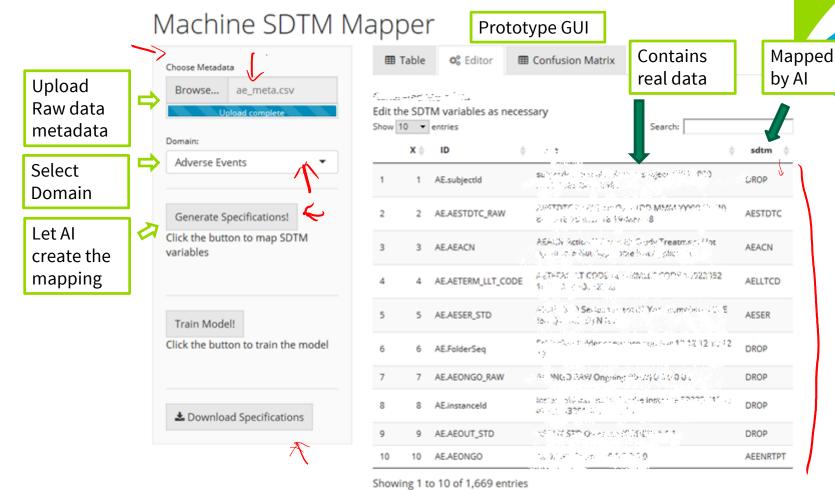






Mapping Accuracy on data from 3 new studies using Final Ensemble Model

0.97484





# **AGENDA**



Problem to Solve



SDTM Mapping with "Machine Learning"



**Tools Used** 



Thought



#### Tools used for development

Programming	IDE	ML Framework	Purpose
<b>S</b> .sas.	NA <del>C</del>	NA	Metadata extraction from sas7bdat
R	R Studio Jupyter	caret	NLP ML Visualizations
<b>?</b> python™	Jupyter	lecun	Transfer learning ML Visualizations



# **AGENDA**



Problem to Solve



SDTM Mapping with "Machine Learning"



Tools Used



Thought





**Summary** 

- This proof of concept demonstrated that machine learning along with a natural language model can produce a pretty accurate SDTM mapping specification document.
- Light weight classic approach works.
- Transfer training has potential for reducing the amount of time to prepare training data



#### **Current/Future Works**



Optimization of NNLM based algorithm

End to end machine learning to step wise approach Transfer learning from other state of the art word embeddings

Development of language model for CDISC documentations



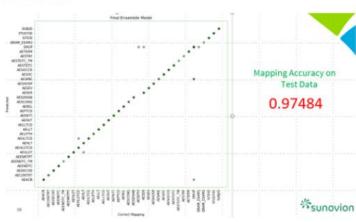
# **Machine Learning Project Sub-team**

Computational Science Working Groups



# **SDTM Mapping** $\angle$

- Project initiated based on the positive outcome of the POC presented at CDISC US Interchange 2018
  - weighting scheme for word-variable matrices
  - neural network language model
- Purpose
  - To explore and develop the SDTM mapping algorithms using NLP, ML, and NN.
  - To publish whitepaper on SDTM mapping





# Thank you



www.linkedin.com/in/STomioka <



Sam.Tomioka@Sunovion.com



Back up



#### **Unsupervised Machine Leaning Approach**

