# DOCUMENT CLASSIFICATION
## TEAM 2

TEAM MEMBERS:

ARTUR AVAGYAN

NUNE TADEVOSYAN

SUPERVISED BY:

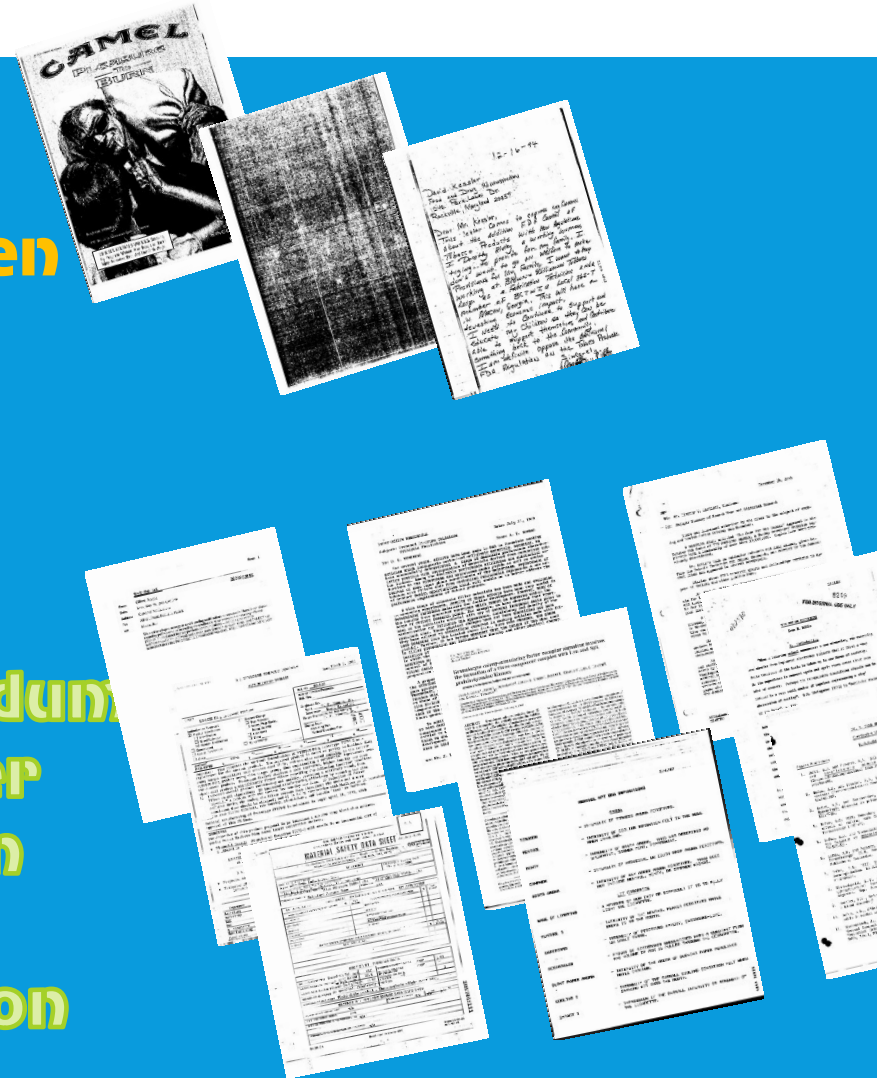ARSEN YEGHIAZARYAN

# SPLIT DATA WITH IMAGE AND TEXT LABELS

- AD
- CV
- E-mail
- File
- Handwritten
- Invoice
- Letter
- Memorandum
- Newspaper
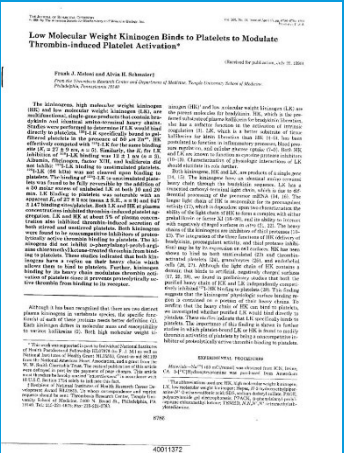- Publication
- Report
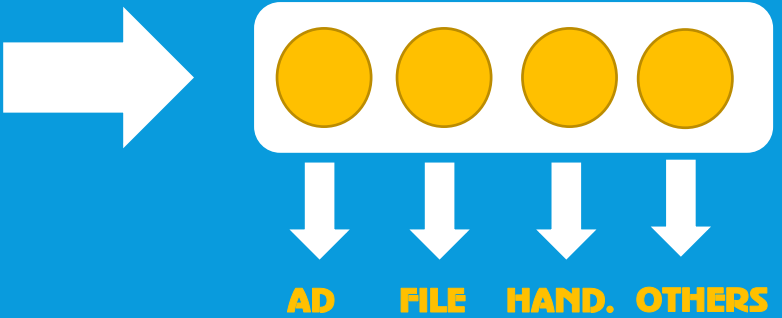- Specification
- Survey

**IMAGE**

- AD
- File
- Handwritten

**TEXT**

- CV
- E-mail
- Invoice
- Letter
- Memorandum
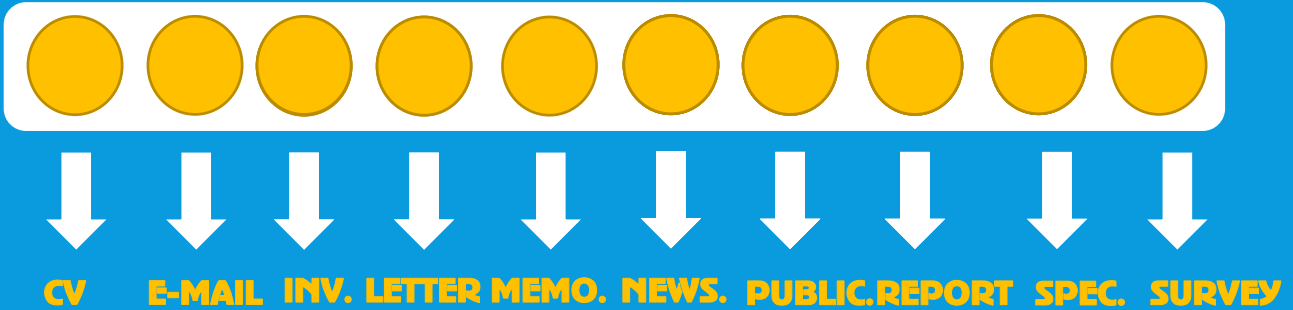- Newspaper
- Publication
- Report
- Specification
- Survey

# 4 LABEL IMAGE CLASSIFICATION

## VGG 16

**Data:** ~ 29 k
**Model size:** ~134 mln
**Number of epochs:** 20
**Device:** CPU
**Train acc:** 88.7%
**Val acc:** 88.1%

# ACCURACY METRICS FROM IMAGE CLASSIFICATION



|              | precision | recall | f1-score | support |
|--------------|-----------|--------|----------|---------|
| ad           | 0.74      | 0.78   | 0.76     | 18      |
| file         | 0.92      | 0.92   | 0.92     | 25      |
| handwritten  | 0.92      | 0.92   | 0.92     | 24      |
| others       | 0.9       | 0.86   | 0.88     | 21      |
|              |           |        |          |         |
| accuracy     |           |        | 0.88     | 88      |
| macro avg    | 0.87      | 0.87   | 0.87     | 88      |
| weighted avg | 0.88      | 0.88   | 0.88     | 88      |

# EXTRACT TEXT FROM IMAGE

## WITH PYTESSERACT

# 10 LABEL TEXT CLASSIFICATION

TRAIN DATA → 70%
VALIDATION DATA → 12%
TEST DATA → 18%

# DISTRIBUTIONS OF NUMBERS OF WORDS IN EACH CLASS

# WORDCLOUDS AFTER TEXT PRE-PROCESSING

# TEXT CLASSIFICATION MODELS
## (VALIDATION ACCURACY)

| vectorization | clf | ngram | train_accuracy | val_accuracy |
|---|---|---|---|---|
| Count_Vect | Naive Bayes | 1 | 0.835 | 0.777 |
| Count_Vect | Naive Bayes | 2 | 0.877 | 0.803 |
| Count_Vect | Naive Bayes | 3 | 0.876 | 0.801 |
| TF IDF | Naive Bayes | 1 | 0.834 | 0.777 |
| TF IDF | Naive Bayes | 2 | 0.860 | 0.792 |
| TF IDF | Naive Bayes | 3 | 0.859 | 0.792 |
| Count_Vect | Logistic Regression | 1 | 0.989 | 0.841 |
| Count_Vect | Logistic Regression | 2 | 0.992 | 0.853 |
| Count_Vect | Logistic Regression | 3 | 0.992 | 0.852 |
| TF IDF | Logistic Regression | 1 | 0.859 | 0.792 |
| TF IDF | Logistic Regression | 2 | 0.859 | 0.792 |
| TF IDF | Logistic Regression | 3 | 0.859 | 0.792 |
| Fasttext_word | Fasttext | 1 | 0.999 | 0.839 |
| Fasttext_word | Fasttext | 2 | 0.998 | 0.849 |
| Fasttext_word | Fasttext | 3 | 0.996 | 0.841 |
| Fasttext_char | Fasttext | 3 | 0.955 | 0.844 |
| Fasttext_char | Fasttext | 4 | 0.979 | 0.851 |
| Fasttext_char | Fasttext | 5 | 0.986 | 0.855 |

# TEST ACCURACY FROM BEST CLASSIFICATION MODELS

## FastText

Chargram = 5
Train accuracy = 98.6%
Val accuracy = 85.5%
Test accuracy = 86.6%

## NAIVE BAYES

Count vectorizer
Ngram = 2
Train accuracy = 87.7%
Val accuracy = 80.3%
Test accuracy = 80%

## Logistic Regression

Count vectorizer
Ngram = 2
Train accuracy = 99.2%
Val accuracy = 85.3%
Test accuracy = 79.5%

# ACCURACY METRICS FROM FASTTEXT CLASSIFIER



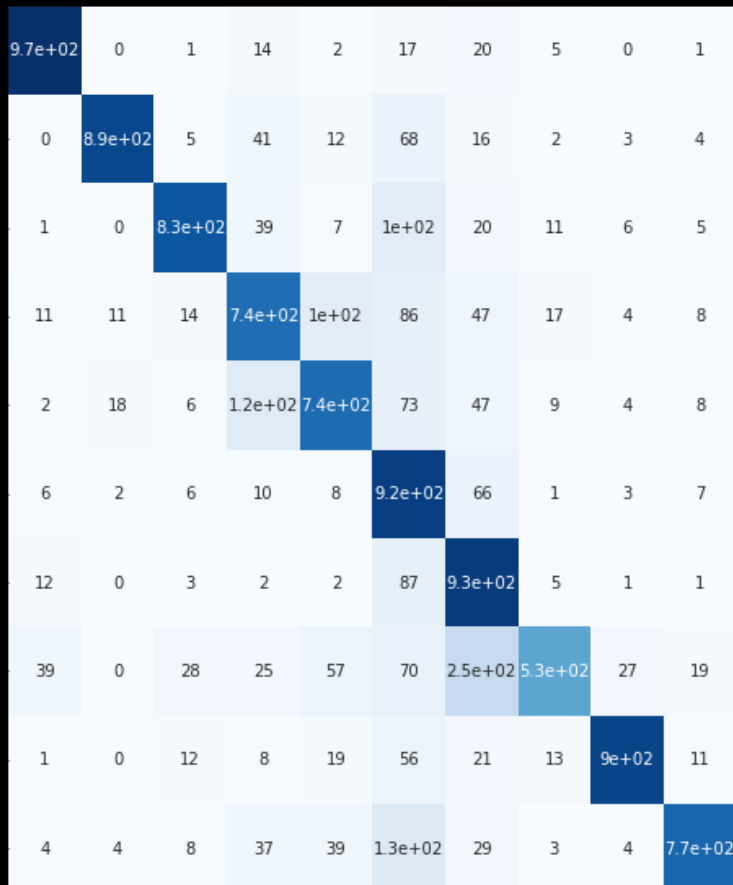|  | precision | recall | f1-score | support |
|---|---|---|---|---|
| cv | 0.97 | 0.97 | 0.97 | 1033 |
| email | 0.93 | 0.91 | 0.92 | 1038 |
| invoice | 0.91 | 0.90 | 0.90 | 1023 |
| letter | 0.82 | 0.80 | 0.81 | 1041 |
| memorandum | 0.83 | 0.81 | 0.82 | 1029 |
| newspaper | 0.77 | 0.86 | 0.81 | 1028 |
| publication | 0.80 | 0.87 | 0.83 | 1043 |
| report | 0.82 | 0.75 | 0.79 | 1037 |
| specification | 0.92 | 0.92 | 0.92 | 1039 |
| survey | 0.91 | 0.87 | 0.89 | 1026 |
| | | | | |
| accuracy | | | 0.87 | 10337 |
| macro avg | 0.87 | 0.87 | 0.87 | 10337 |
| weighted avg | 0.87 | 0.87 | 0.87 | 10337 |

# MODEL ACCURACY

**4 LABEL IMAGE CLASSIFICATION**

**10 LABEL TEXT CLASSIFICATION**

88.1%

86.6%

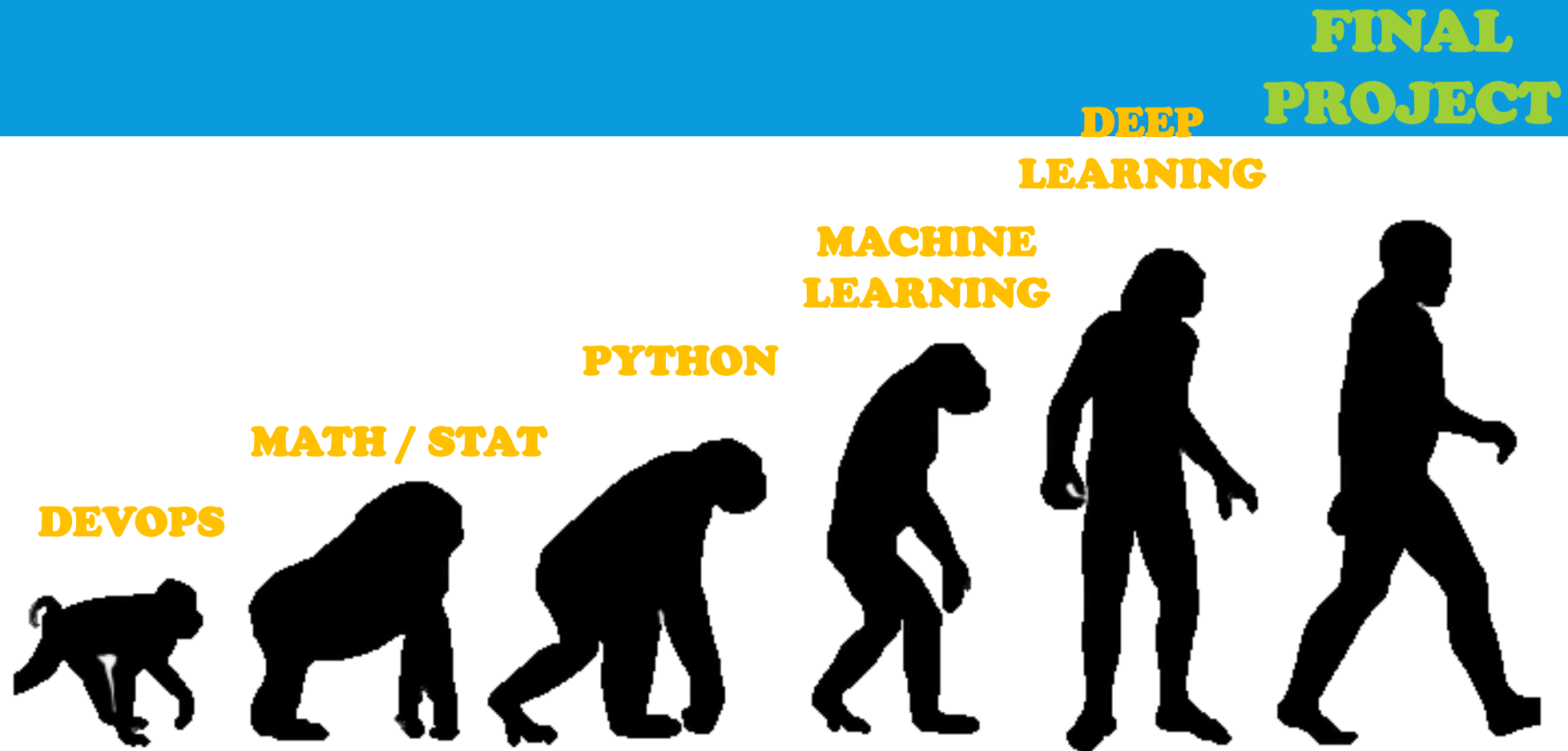**13 LABEL DOCUMENT CLASSIFICATION**

88.7%

# EVALUATION OUR MODEL PIPELINE ON SAMPLE DATA



| | precision | recall | f1-score | support |
|---|---|---|---|---|
| ad | 0.64 | 0.90 | 0.75 | 1000 |
| cv | 0.99 | 0.98 | 0.98 | 1000 |
| email | 0.97 | 0.95 | 0.96 | 1000 |
| file | 0.79 | 0.90 | 0.86 | 1000 |
| handwritten | 0.88 | 0.84 | 0.89 | 1000 |
| invoice | 0.96 | 0.90 | 0.89 | 1000 |
| letter | 0.91 | 0.90 | 0.90 | 1000 |
| memorandum | 0.92 | 0.76 | 0.91 | 1000 |
| newspaper | 0.89 | 0.89 | 0.82 | 1000 |
| publication | 0.89 | 0.89 | 0.89 | 1000 |
| report | 0.94 | 0.82 | 0.87 | 1000 |
| specification | 0.97 | 0.92 | 0.94 | 1000 |
| survey | 0.95 | 0.84 | 0.89 | 1000 |
| | | | | |
| accuracy | | | 0.89 | 13000 |
| macro avg | 0.90 | 0.89 | 0.89 | 13000 |
| weighted avg | 0.90 | 0.88 | 0.89 | 13000 |

# COURSE EVALUATION



FINAL PROJECT

DEEP LEARNING

MACHINE LEARNING

PYTHON

MATH / STAT

DEVOPS

THANKS