

HANDWRITING RECOGNITION

in the Kazakh Language

Alar Akilbekov (Math, Code)

Baktyiar Toksanbay (Code)

2024

DESCRIPTION

Handwritten Text Recognition (HTR).

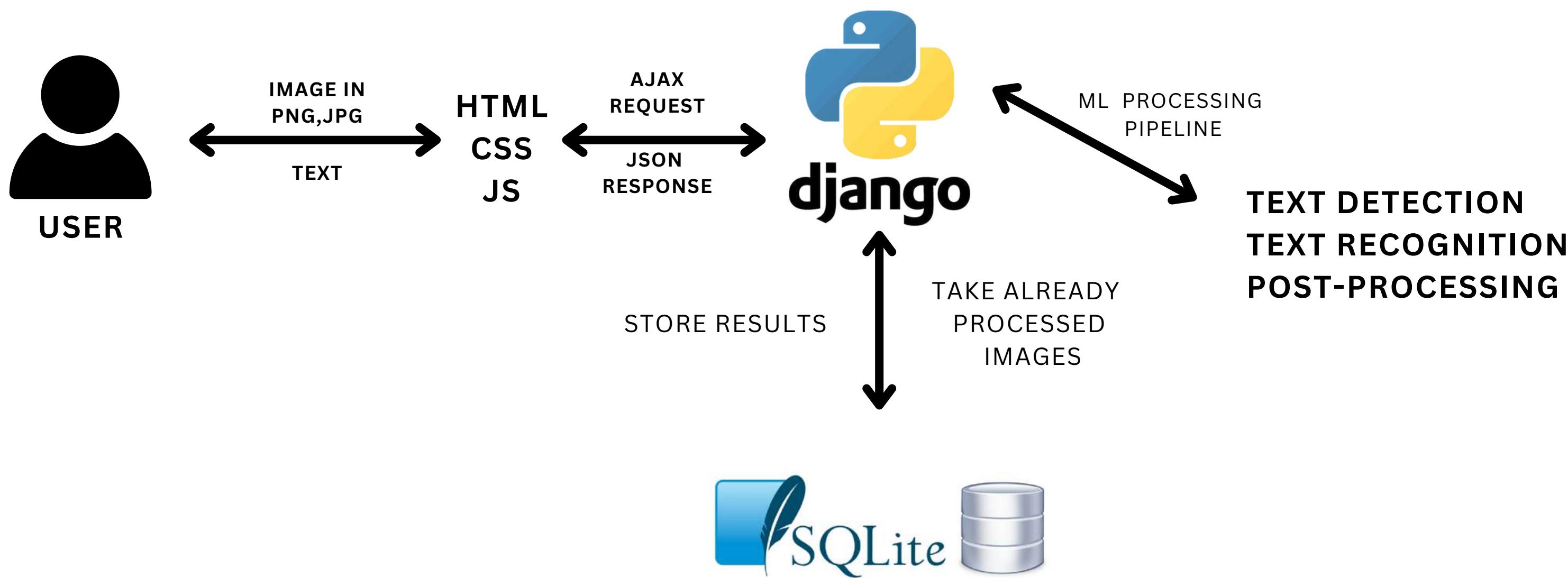
Ability of a computer to receive and interpret intelligible handwritten input from sources such as paper documents, photographs, touch-screens and other devices.



PROJECT STACK

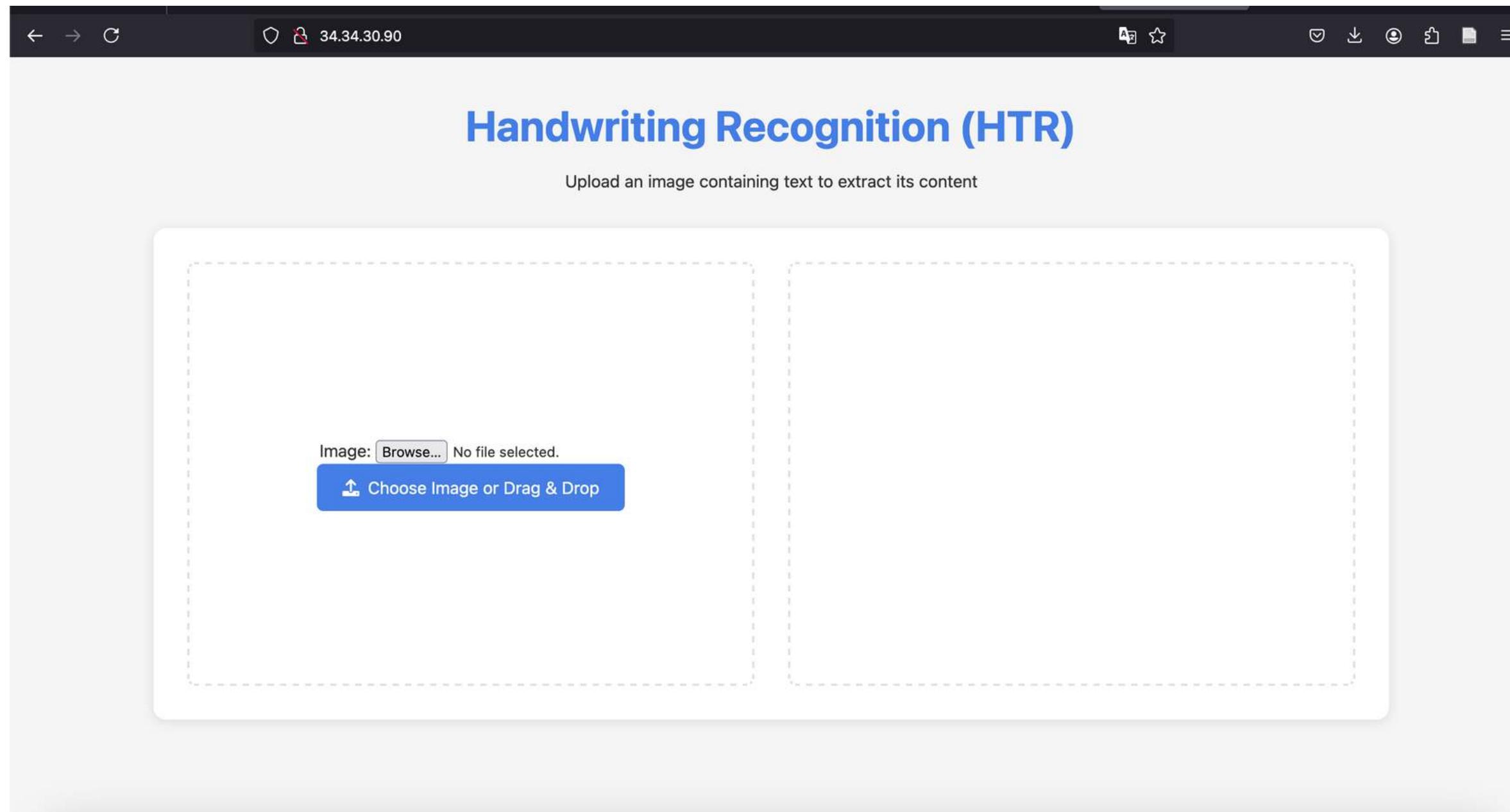


PROJECT ARCHITECTURE



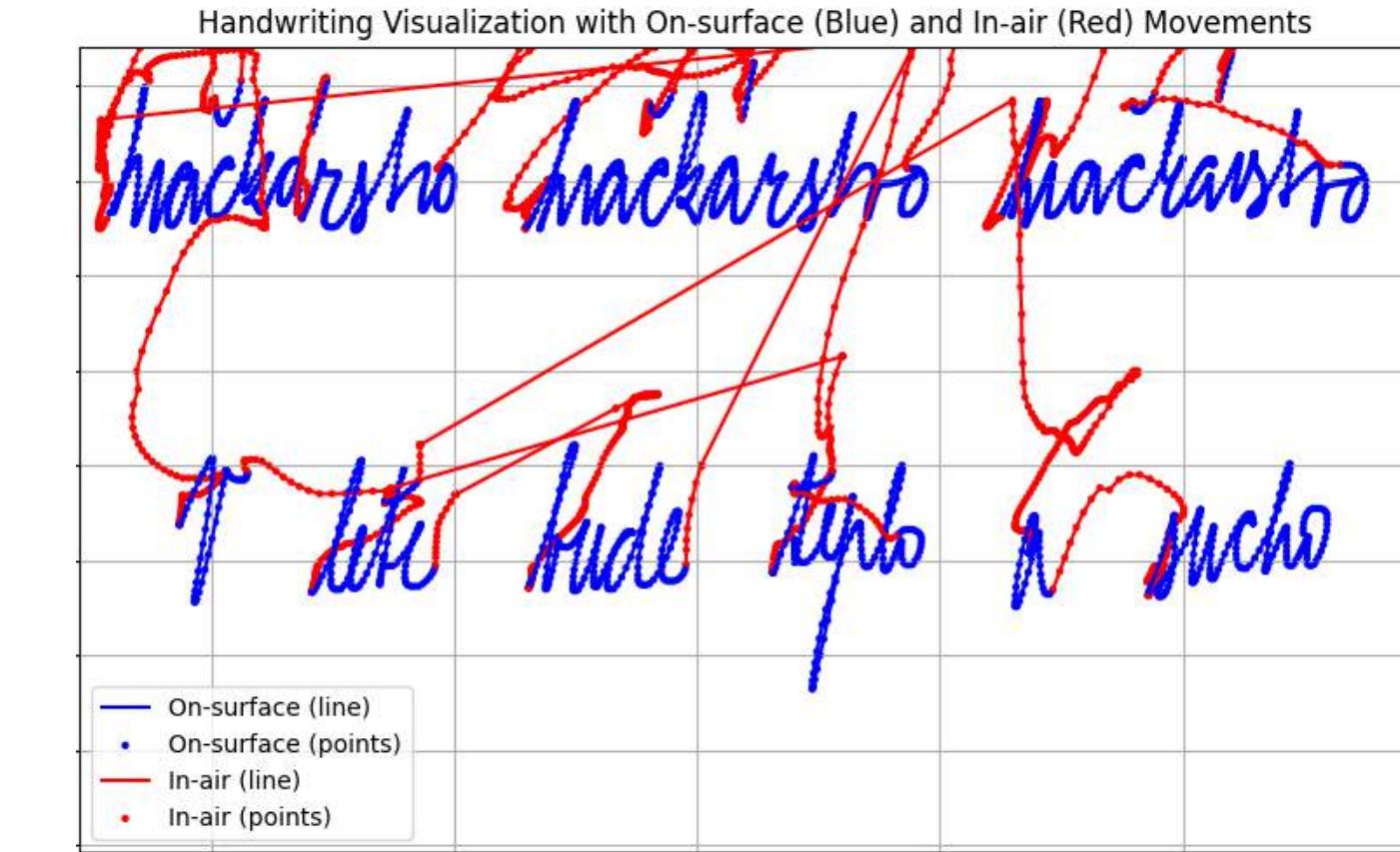
TRY IT OUT

34.34.30.90



On-line:

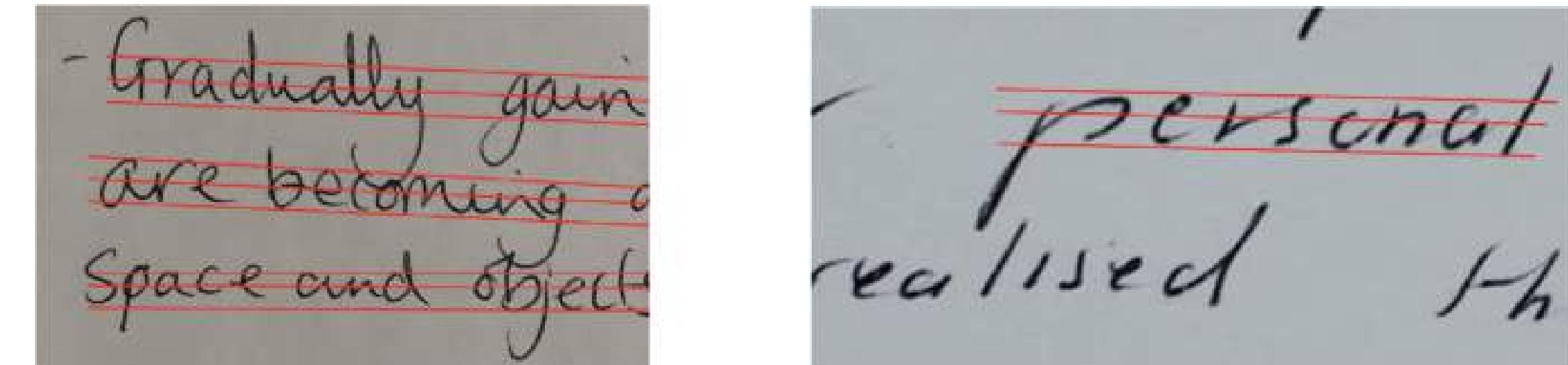
- pen-based computer screen surface
- pen-up and pen-down switching
- pen pressure
- velocity/changes of writing direction
- specifically



(<https://github.com/alarxx/Dysgraphia-Diagnosis>).

Off-line:

- piece of paper
- image



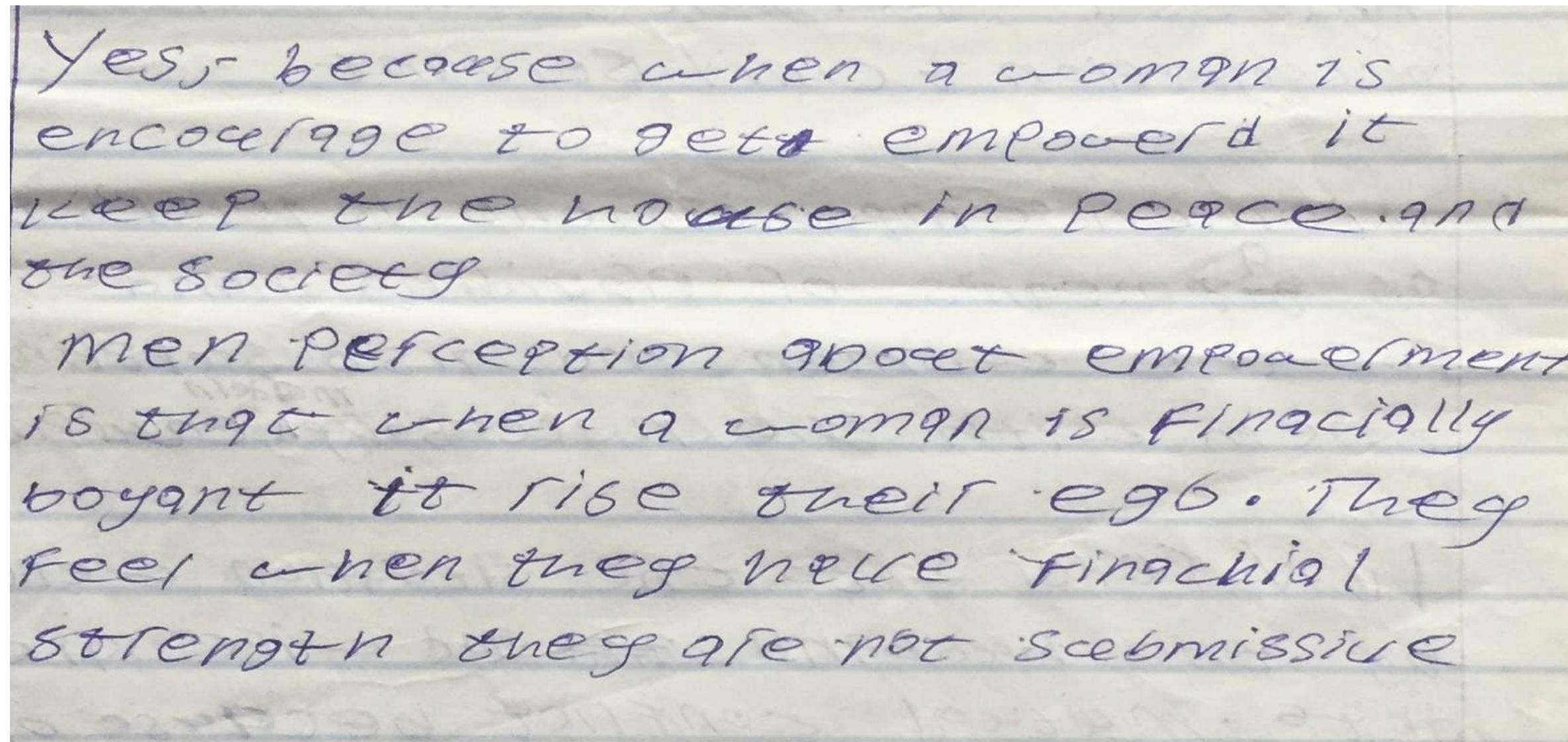
(GNHK: A Dataset for English Handwriting in the Wild)

https://dl.acm.org/doi/10.1007/978-3-030-86337-1_27

<https://github.com/GoodNotes/GNHK-dataset>

PROBLEM

Input data:



(GNHK: A Dataset for English Handwriting in the Wild)

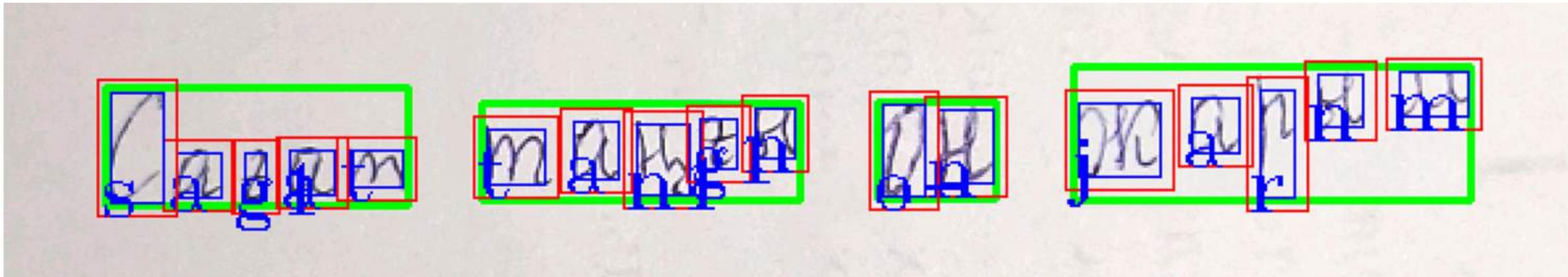
Output data: written text on the image

- **Text Detection**
- **Classification**

SOLUTION APPROACHES

- Optical Character Recognition (OCR)
- Intelligent Character Recognition (ICR)
- Intelligent Word Recognition (IWR)

OCR engines are primarily focused on character-by-character machine printed text recognition from a scanned document and **ICR** for different fonts or even handwritten text:

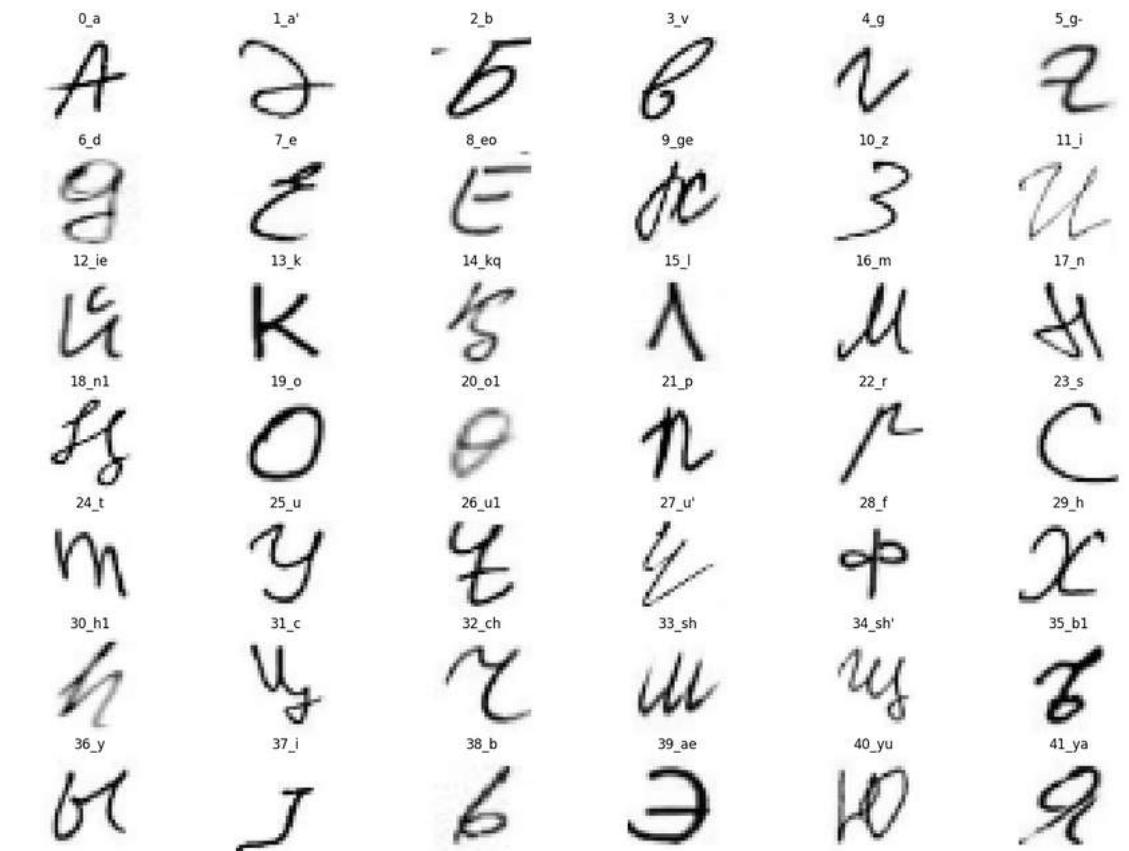


Intelligent Word Recognition (IWR) is the recognition of unconstrained handwritten words. IWR recognizes entire handwritten words or phrases instead of character-by-character, like its predecessors.

CLASSIFICATION DATASETS

Kazakh Alphabet Letters (CMNIST)

<https://github.com/bolattleubayev/cmnist>



Cyrillic Handwriting Dataset

<https://www.kaggle.com/datasets/constantinwerner/cyrillic-handwriting-dataset/data>



DATASET

KOHTD: Kazakh offline handwritten text dataset

Toiganbayeva, N., Kasem, M., Abdimanap, G., Bostanbekov, K., Abdallah, A., Alimova, A., & Nurseitov, D. (2022). Kohtd: Kazakh offline handwritten text dataset. *Signal Processing: Image Communication*, 108, 116827.

Chicago

Github: <https://github.com/abdoelsayed2016/KOHTD>

2) Шартсыз ортасынан туруу есептөөнүү
f(x) функциясының дарыткы төмөрлөккөк
минимумын және максимумын издеүтүү
есептөөнүү.

Шартсыз ортасынан туруу есептөөнүү
 $f(x)$ -шарттың ~~функция~~ максимумун табаңын
— x айналыштарын даскаршамасын аныктайып
шешиндер. DB маралыгын берген,
жоле даскаршамасын айналышташып
B мөнөттөргө индекстөөн көз-кеңзен
үйнегер сиизен ортасынан туруу
есептөөнүү маралыгын шешиндер жөнөтүштөөнүү.

Шартсыз: даскаршама нүхемдөрүнөн
ураганымын сиизтөөнүү, Ол учила
бірнеше дербес мөнөттөргөтөн мадашуу:

$$\frac{df}{dx} = -2x_1 + 6; \frac{dE}{dx} = -8x_2 + 32$$

$$\nabla f(x) = \begin{pmatrix} -2x_1 + 6 \\ -8x_2 + 32 \end{pmatrix}, \nabla F(x_0) = \begin{pmatrix} -2 \cdot 7 + 6 \\ -8 \cdot 4 + 32 \end{pmatrix} = \begin{pmatrix} -2 \\ 0 \end{pmatrix}$$

Маралы x₀ нүхемдөрүнүү:

$$x_1 = x_0 + \lambda_1 \nabla f(x_0), \nabla F(x_0) = \begin{pmatrix} 4 \\ 4 \end{pmatrix} + \lambda_1 \begin{pmatrix} -2 \\ 0 \end{pmatrix} = \begin{pmatrix} 4 - 2\lambda_1 \\ 4 \end{pmatrix}$$

Маралы нүхемдөрүнүү
мадашуу:

$$\nabla F(x_1) = \begin{pmatrix} -2 \cdot (4 - 2\lambda_1) + 6 \\ -8 \cdot 4 + 32 \end{pmatrix} = \begin{pmatrix} 10\lambda_1 - 8 \\ 0 \end{pmatrix}$$

$$\frac{B \nabla F}{\lambda_1} = \nabla F(x_0) \cdot \nabla F(x_1) = \begin{pmatrix} -2 \\ 0 \end{pmatrix} \cdot \begin{pmatrix} 10\lambda_1 - 8 \\ 0 \end{pmatrix} = 0$$



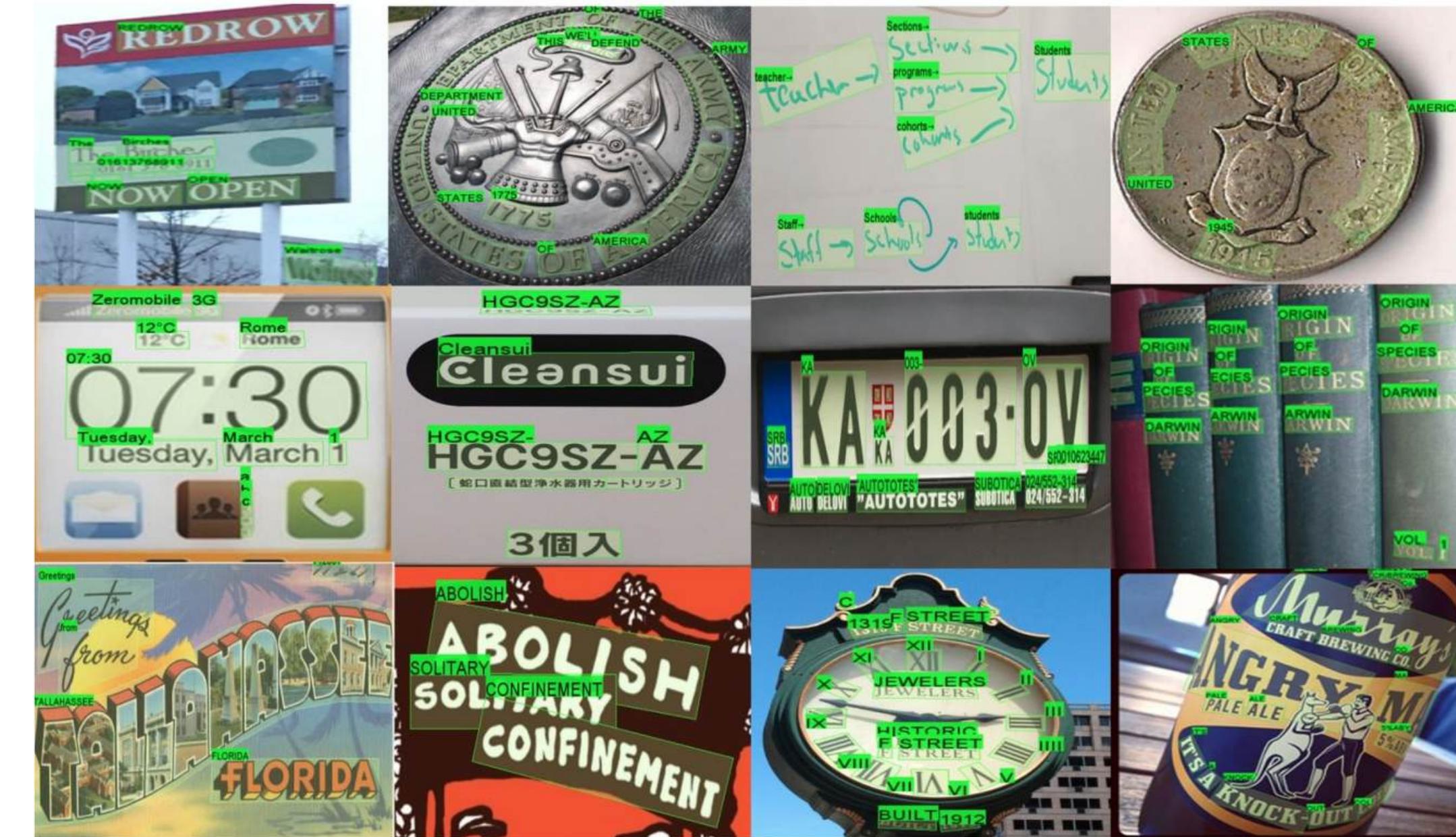
TEXT DETECTION DATASETS

MSRA Text Detection 500 Database (MSRA-TD500)

[http://www.iapr-tc11.org/mediawiki/index.php?
title=MSRA_Text_Detection_500_Database_\(MSRA-TD500\)](http://www.iapr-tc11.org/mediawiki/index.php?title=MSRA_Text_Detection_500_Database_(MSRA-TD500))

TextOCR

<https://paperswithcode.com/dataset/textocr>

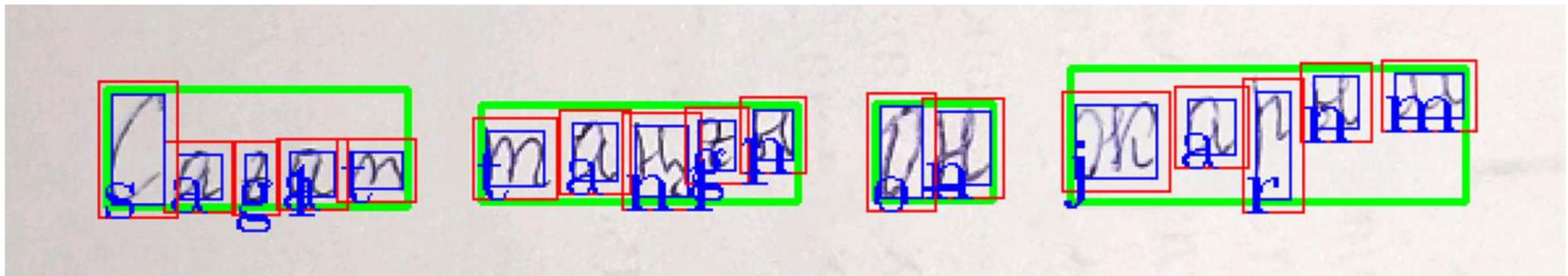
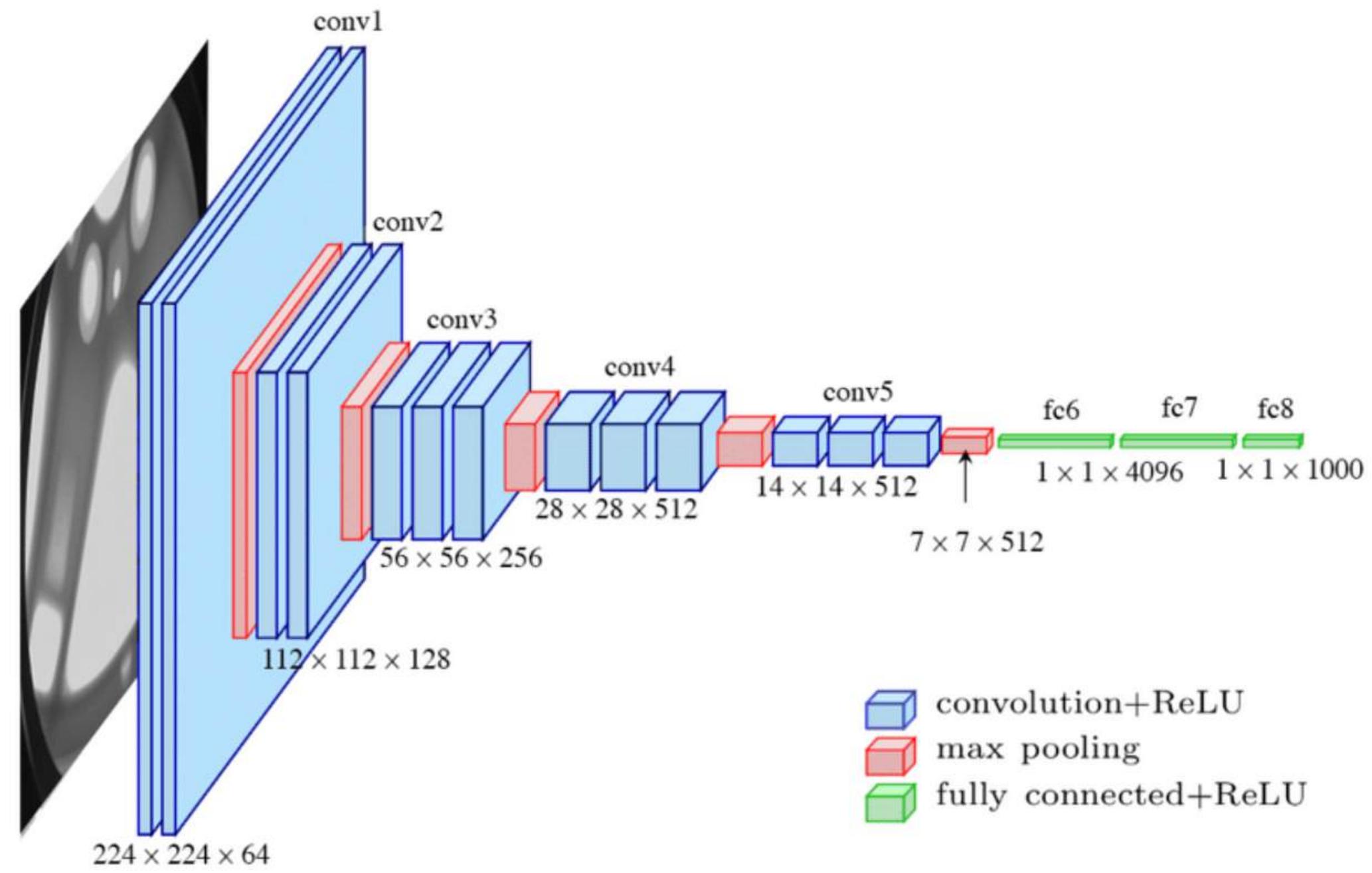


INTELLIGENT WORD RECOGNITION [IWR]

Application of deep learning: convolutional neural networks (CNNs) and recurrent neural networks (RNNs), including Long Short-Term Memory (LSTM) and Gated Recurrent Units (GRUs), to handwriting image processing.

Recently, transformers, especially those that use attention mechanisms, have become popular for sequence processing, including handwriting recognition. They allow the model to focus on important parts of the input data when generating each character.

GANs can be used to generate realistic handwritten text, which helps to increase the available data for training recognition models.



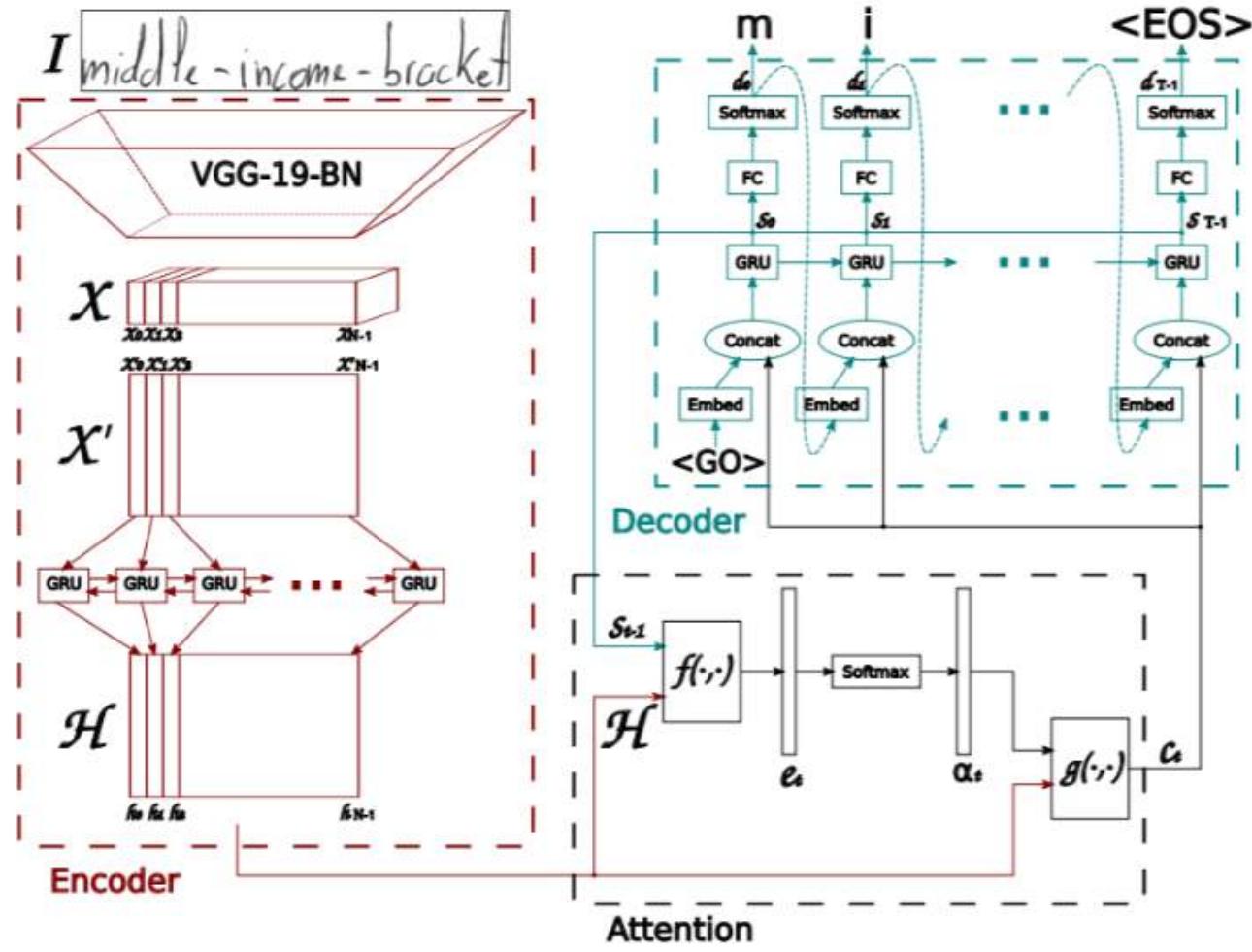


Fig. 1: Architecture of the seq2seq model with attention mechanism.

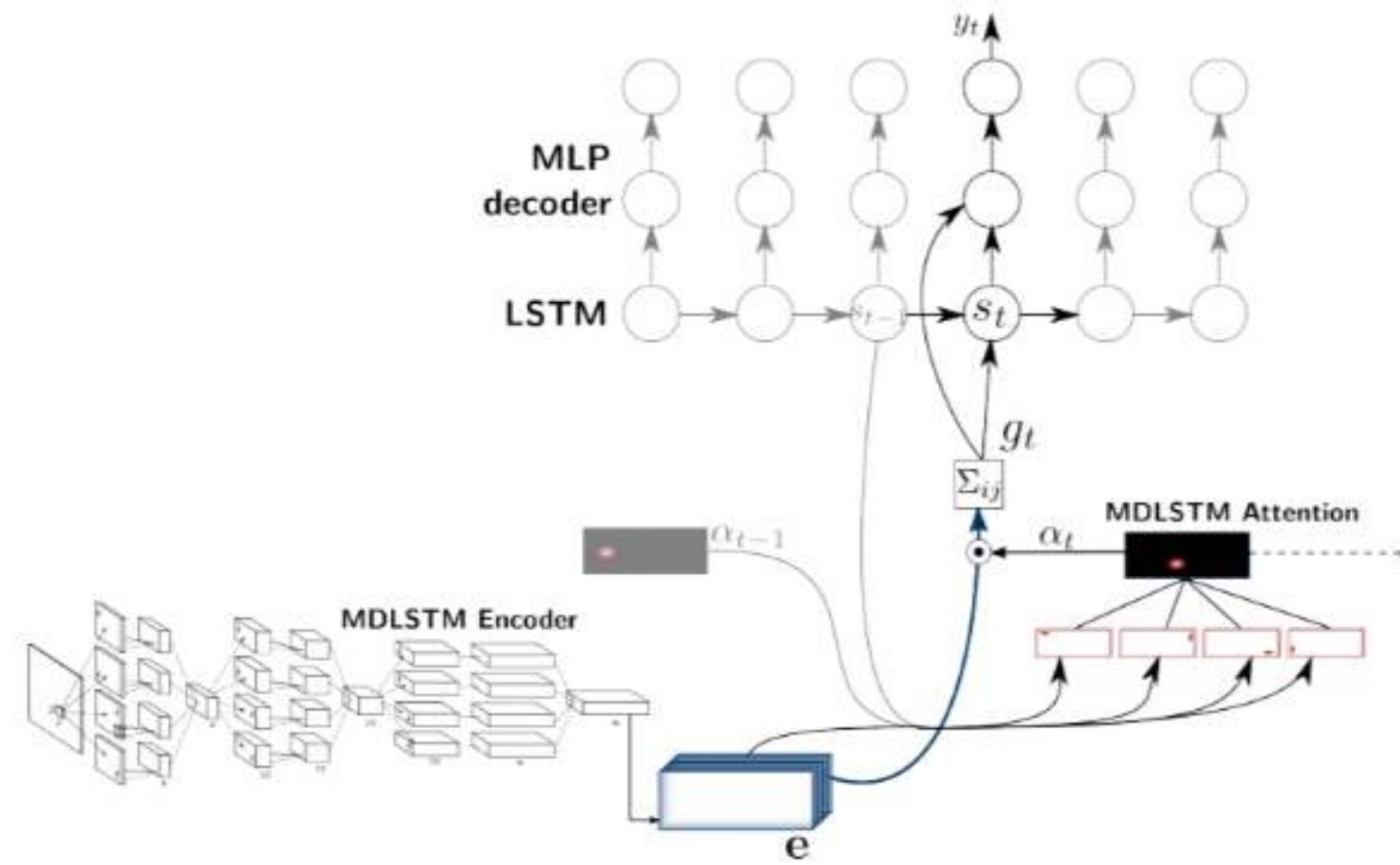


Figure 2: Proposed architecture. The encoder network has the same architecture as the standard network of Figure 1, except for the collapse and softmax layers. At each timestep, the feature maps, along with the previous attention map and state features are fed to an MDLSTM network which outputs new attention weights at each position. The weighted sum of the encoder features is computed and given to the state LSTM, and to the decoder. The decoder also considers the new state features and outputs character probabilities.

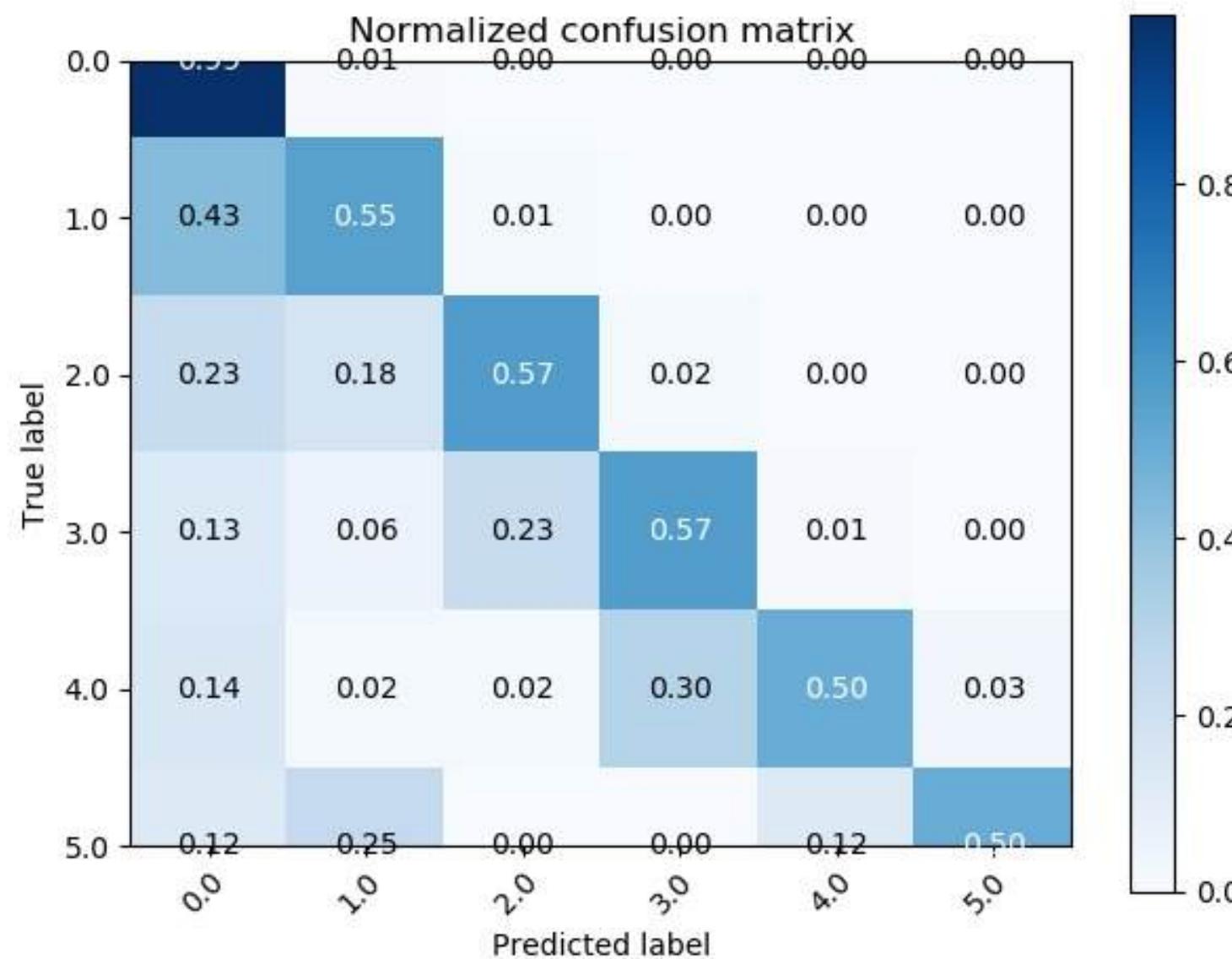
POST-EVALUATION

	h	e	I	I	o
0	1	2	3	4	5
k	1	1	2	3	4
e	2	2			
I	3				
m	4				

Levenshtein distance

Analysing the probability of words in a sentence

CONFUSION MATRIX FOR LETTERS





Tesseract OCR

Output

W-4

Form W-4
Department of the Treasury
Internal Revenue Service

Employee's Withholding Certificate

► Complete Form W-4 so that your employer can withhold the correct federal income tax from your pay.
► Give Form W-4 to your employer.
► Your withholding is subject to review by the IRS.

Step 1:
Enter Personal Information

Name and middle initial
Adrian

Address
PO Box 17598 #17900
PO Box 17598 #17900

City, State and Zip Code
Baltimore, MD 21297-1598

Baltimore, MD 21297-1598

(c) Single or Married filing separately
 Married filing jointly (or Qualifying widow(er))
 Head of household (Check only if you're unmarried and pay more than half the costs of keeping up a home for your dependents)

TESSERACT TEXT DETECTION

Tesseract Text Detection

МО дегеніміз не?
МО дегеніміз не?

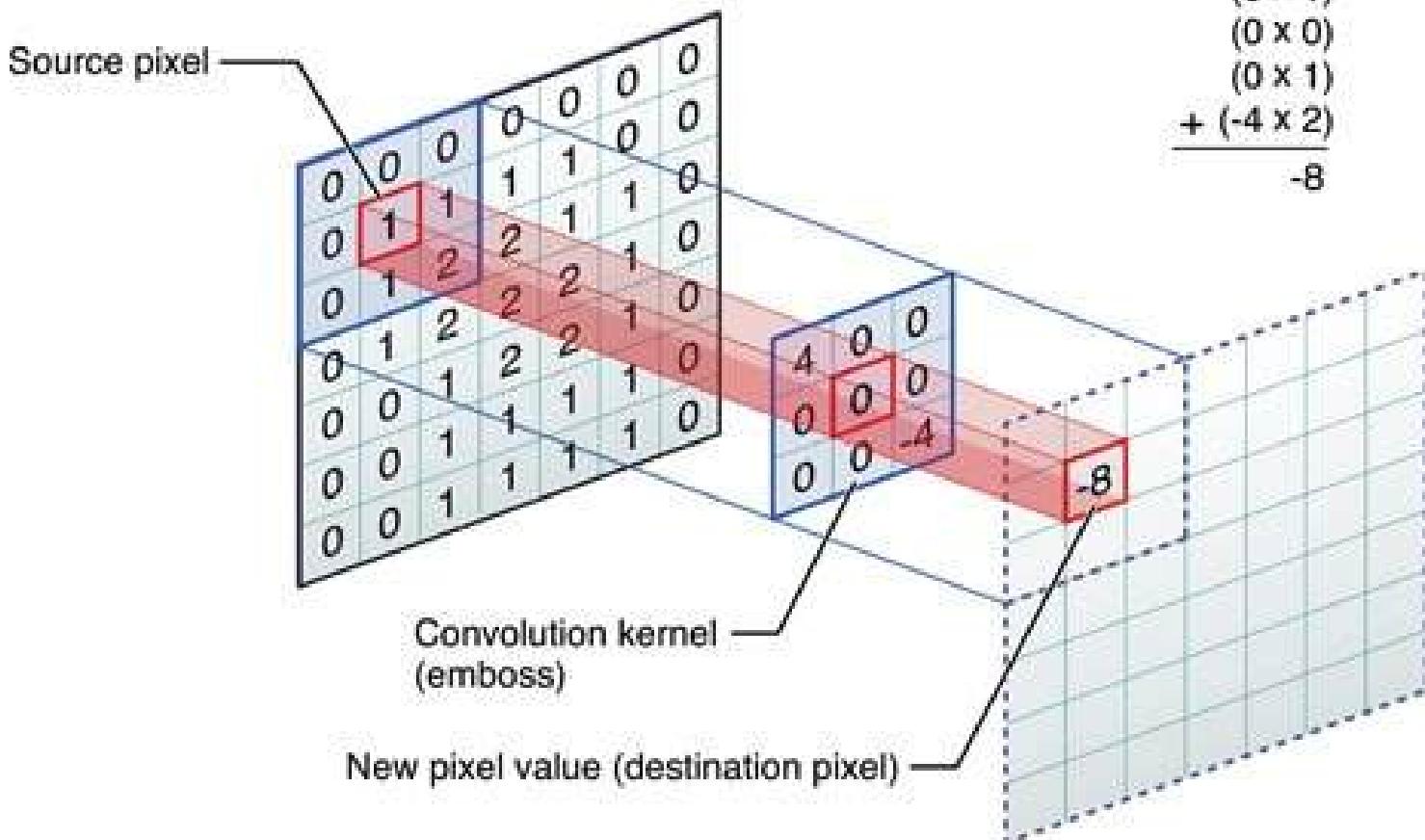
- “ Машиналық оқыту - бул адамның араласуынсыз өздігінен
▪ Машиналық оқыту турде бул адамның араласуынсыз өздігінен
автоматты турде талдайтын және шешім қабылдайтын
автоматты турде талдайтын және шешім қабылдайтын
алгоритмдер мен деректердің қамтитын жасанды интеллект
алгоритмдер мен деректердің қамтитын жасанды интеллект
қосымшасы. қосымшасы.
- “ Ол алдыңғы тәжірибе негізінде компьютердің тапсырмаларды
▪ Ол алдыңғы тәжірибе негізінде компьютердің тапсырмаларды
қалай дербес орындаитының сипаттайды.
қалай дербес орындаитының сипаттайды.
- “ Сондықтан, машиналық тілде жасанды интеллект тәжірибе
▪ Сондықтан, машиналық тілде жасанды интеллект тәжірибе
негізінде жасалады деп айта аламыз.
негізінде жасалады деп айта аламыз.

Машындың орнында 8-нада ажырылған арасалуынан
Озінің автомобилін түрде таңдалған жаңа
шешім қабылдайтын алгоритмдер мен деңгектерді
тұжырымдаудың иштемелекті қосынчыласы.

TEXT DETECTION

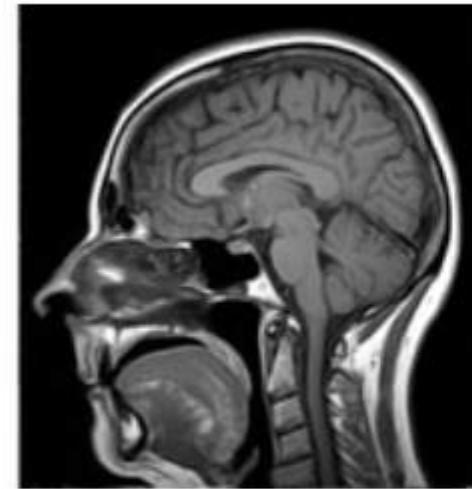
Edge Detector - searches for object contours by looking for a brightness difference

Center element of the kernel is placed over the source pixel. The source pixel is then replaced with a weighted sum of itself and nearby pixels.



$$[f * g](t) = \int_0^t f(x) g(t - x) dx$$

0	0	0
0	1	0
0	0	0



Identity Kernel: Just the starting image itself.

1	1	1
1	1	1
1	1	1



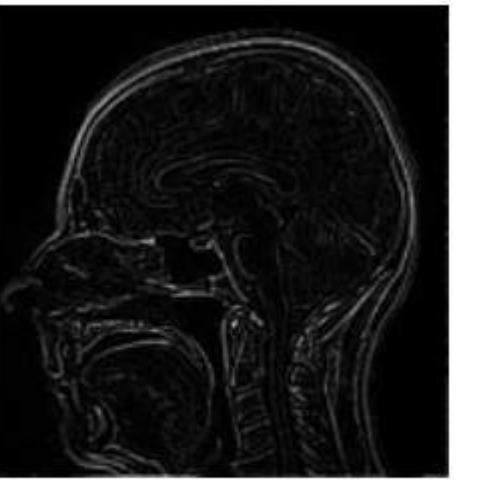
Box Blur Kernel: Each pixel is averaged equally with its 8 neighbors.

0	-1	0
-1	5	-1
0	-1	0



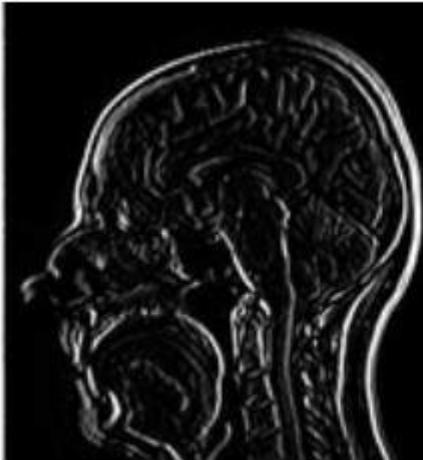
Sharpen Kernel: Differences emphasized with its adjacent pixel values.

-1	-1	-1
-1	8	-1
-1	-1	-1



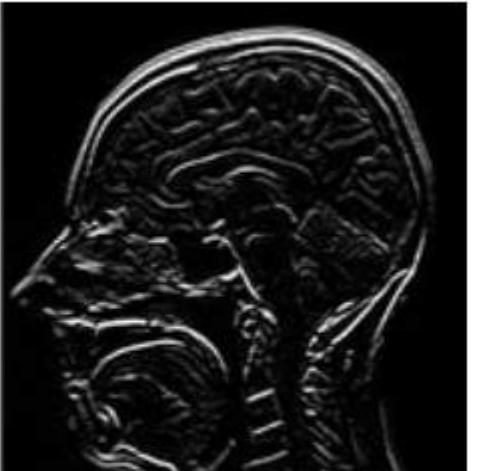
Outline ('Edge') Kernel: Highlights neighboring pixels with large differences of intensity.

1	0	-1
2	0	-2
1	0	-1



Left-to-right Sobel Kernel: Used to accentuate differences between pixels along the horizontal axis

-1	-2	-1
0	0	0
1	2	1



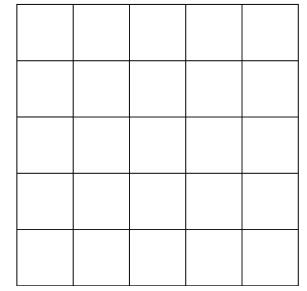
Superior-to-inferior Sobel Kernel: Used to accentuate differences between pixels along the vertical axis.

Morphological Transformations

Image



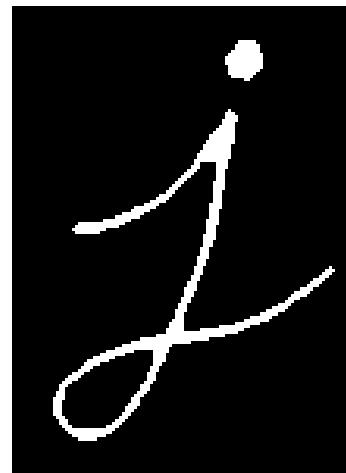
+



A pixel in the resulting image will be considered 1

- 1) only if all the pixels
 - 2) or at least one
- under the kernel is 1

1) Erosion



2) Dilation



3) Opening



4) Closing



1. The Dangerous and Thrilling Documentation Chronicles

Kismet Chameleon

This journey begins on a bleary Monday morning. Our intrepid team is in desperate need of double shot mochas, but the milk expired eight days ago. A trip to the dairy was out of the question. On Friday night, a mutant, script-injecting warlock had infected the Shetland cattle herd with a ravenous craving for tags and annotations. The security wolves were at a trust building retreat in Katchanga, and no one in the village could locate their defensive operations manual.

Weak daylight trickled across the stripped pasture, chased by distant bovine screams...

Cavern Glow

The river rages through the cavern, rattling its content...

2. The Dangerous and Thrilling Documentation Chronicles

Kismet Chameleon

This journey begins on a bleary Monday morning. Our intrepid team is in desperate need of double shot mochas, but the milk expired eight days ago. A trip to the dairy was out of the question. On Friday night, a mutant, script-injecting warlock had infected the Shetland cattle herd with a ravenous craving for tags and annotations. The security wolves were at a trust building retreat in Katchanga, and no one in the village could locate their defensive operations manual.

Weak daylight trickled across the stripped pasture, chased by distant bovine screams...

Cavern Glow

The river rages through the cavern, rattling its content...

3. [REDACTED]

[REDACTED]

4. The Dangerous and Thrilling Documentation Chronicles

Kismet Chameleon

This journey begins on a bleary Monday morning. Our intrepid team is in desperate need of double shot mochas, but the milk expired eight days ago. A trip to the dairy was out of the question. On Friday night, a mutant, script-injecting warlock had infected the Shetland cattle herd with a ravenous craving for tags and annotations. The security wolves were at a trust building retreat in Katchanga, and no one in the village could locate their defensive operations manual.

Weak daylight trickled across the stripped pasture, chased by distant bovine screams...

Cavern Glow

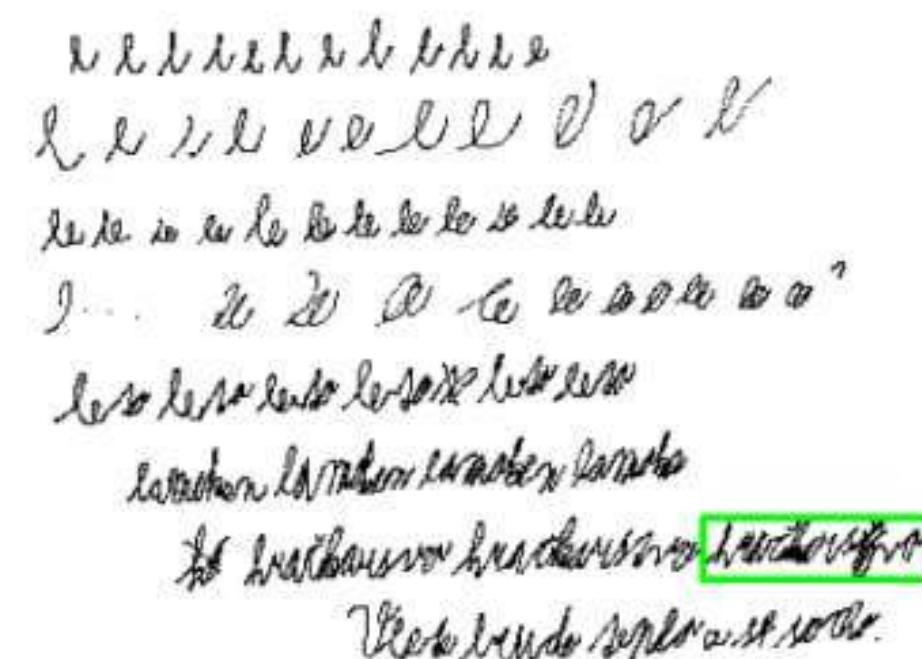
The river rages through the cavern, rattling its content...

Creating a one-word picture:

After Morphological Closing

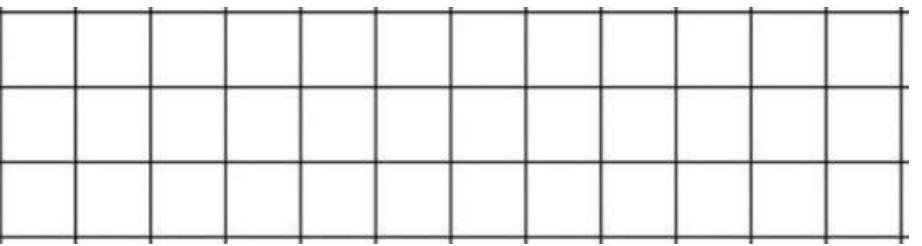
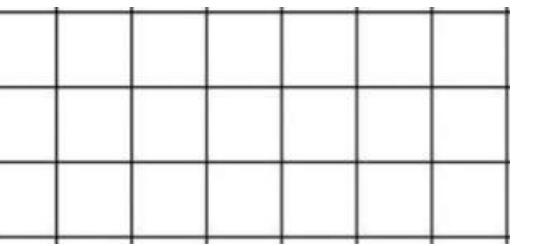
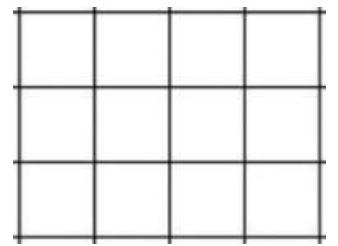


Original Image with Boundary of Word #1



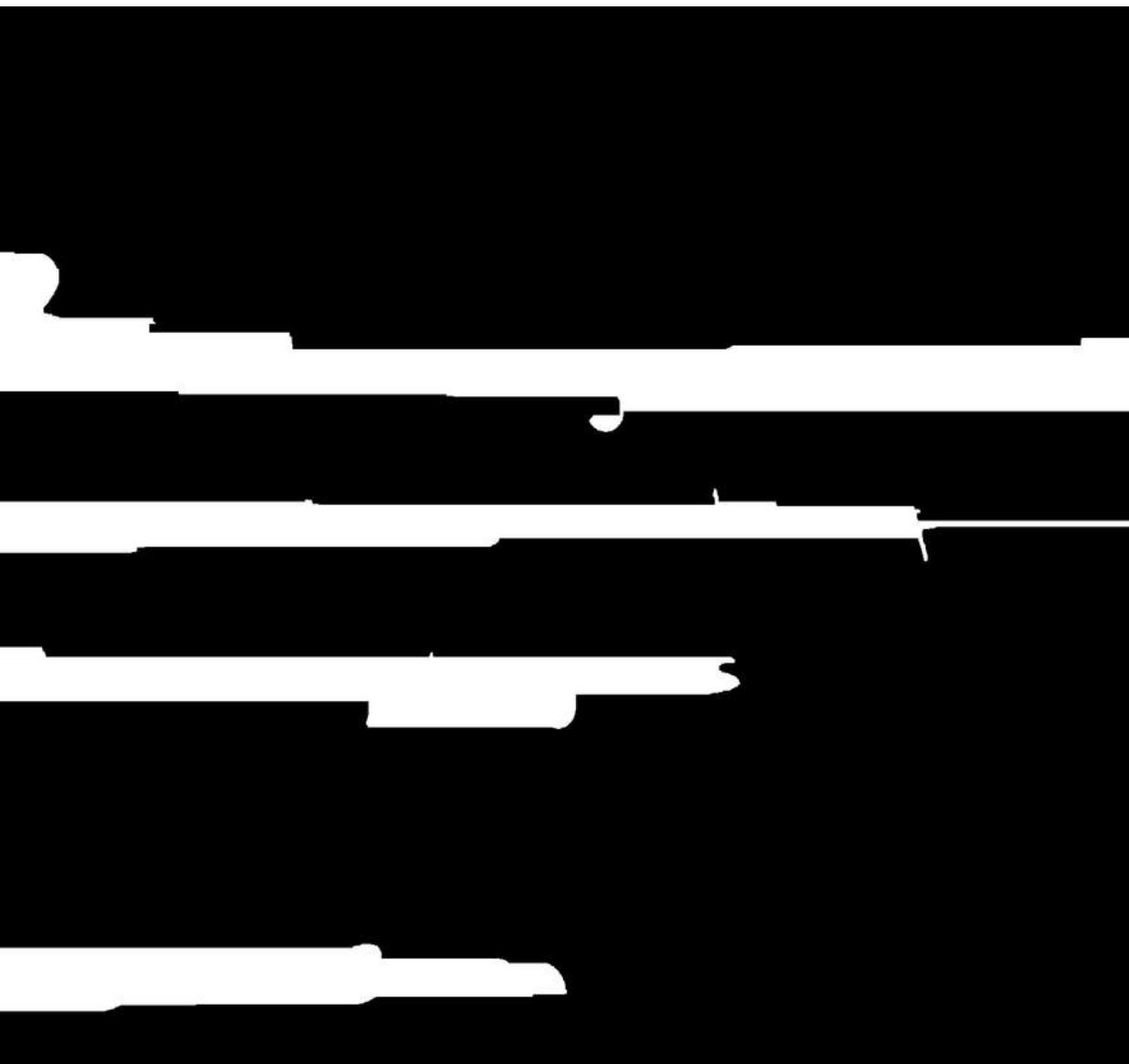
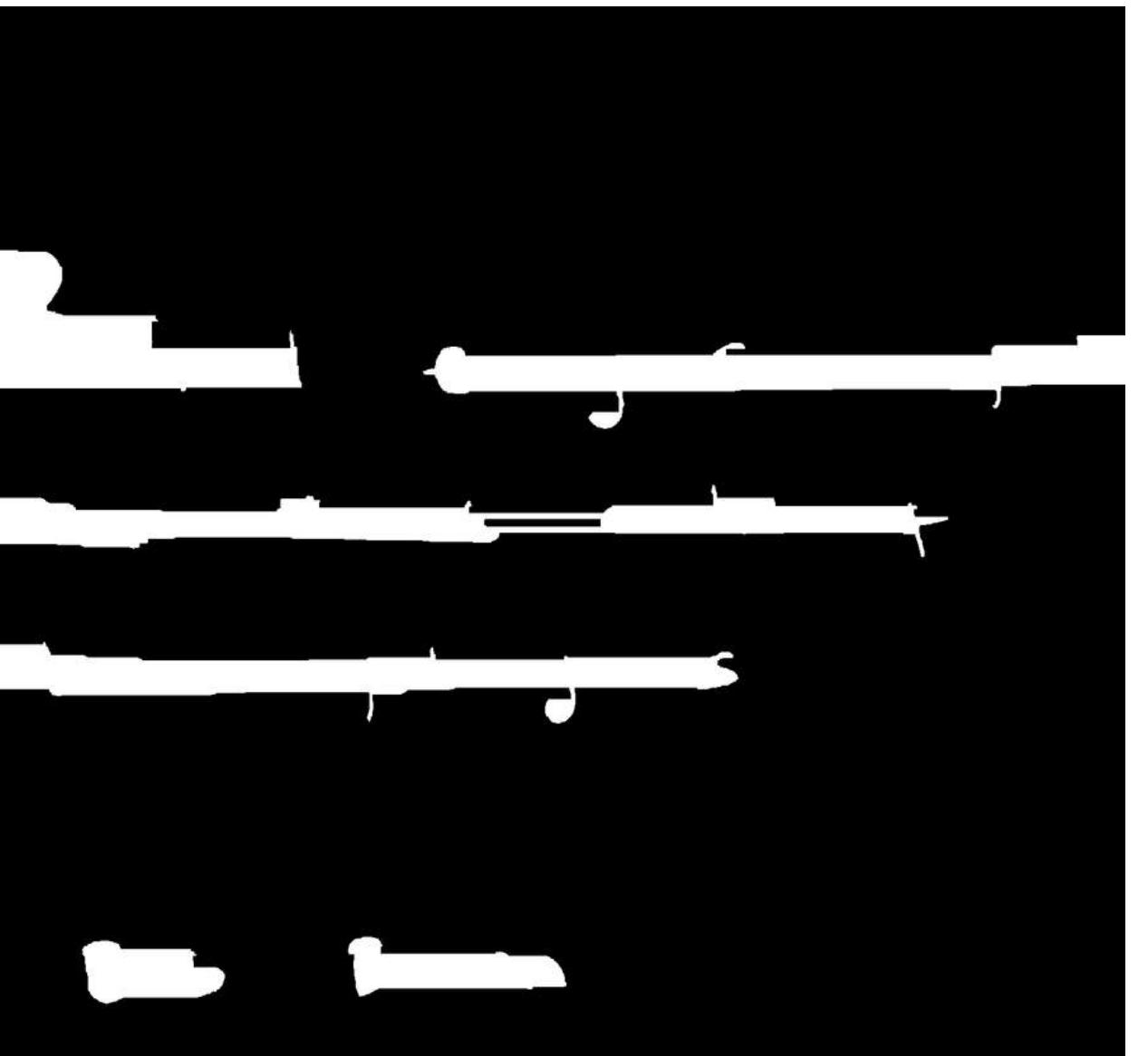
224x224 Image with Word #1





Behind every successful
man there's a lot of
unsuccessful years

Bob Brown



Behind every successful

man there's a lot of

unsuccessful years

Bob Brown

Behind every successful

man there's a lot of

unsuccessful years

Bob Brown

EAST: An Efficient and Accurate Scene Text Detector

Xinyu Zhou, Cong Yao, He Wen, Yuzhi Wang, Shuchang Zhou, Weiran He, and Jiajun Liang

Megvii Technology Inc., Beijing, China

{zxy, yaocong, wenhe, wangyuzhi, zsc, hwr, liangjiajun}@megvii.com

Abstract

Previous approaches for scene text detection have already achieved promising performances across various benchmarks. However, they usually fall short when dealing with challenging scenarios, even when equipped with deep neural network models, because the overall performance is determined by the interplay of multiple stages and components in the pipelines. In this work, we propose a simple yet powerful pipeline that yields fast and accurate text detection in natural scenes. The pipeline directly predicts words or text lines of arbitrary orientations and quadrilateral shapes in full images, eliminating unnecessary intermediate steps (e.g., candidate aggregation and word partitioning), with a single neural network. The simplicity of our pipeline allows concentrating efforts on designing loss functions and neural network architecture. Experiments on standard datasets including ICDAR 2015, COCO-Text and MSRA-TD500 demonstrate that the proposed algorithm significantly outperforms state-of-the-art methods in terms of both accuracy and efficiency. On the ICDAR 2015 dataset, the proposed algorithm achieves an F-score of 0.7820 at 13.2fps at 720p resolution.

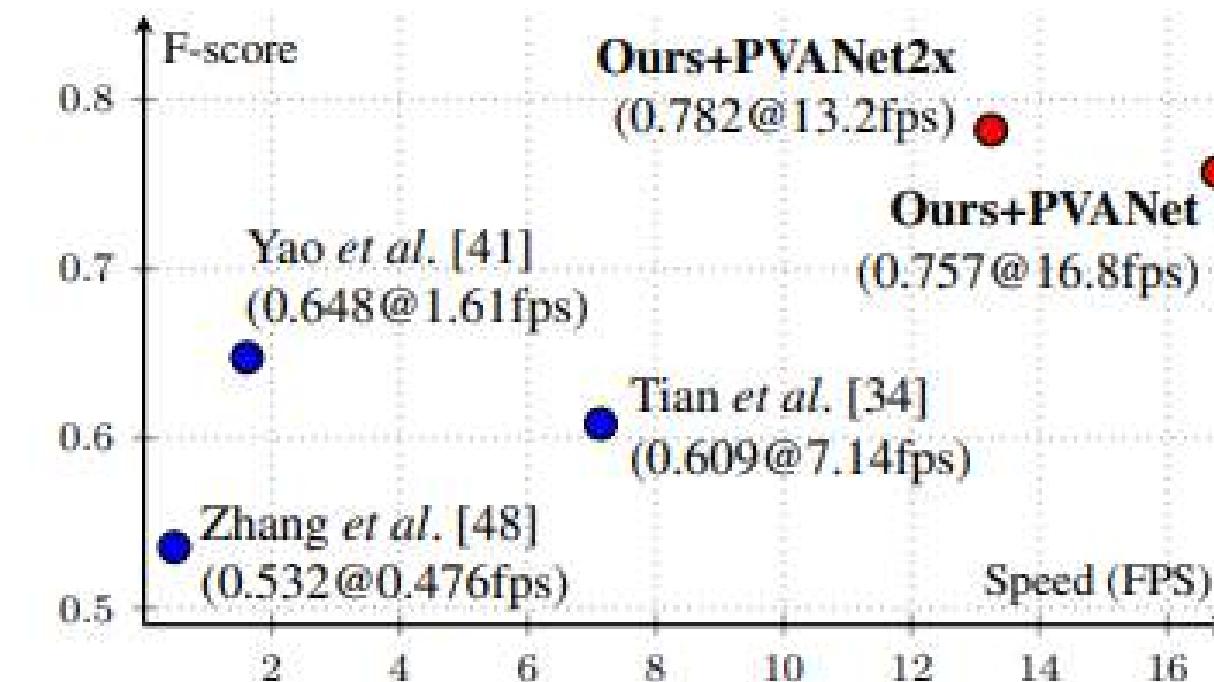


Figure 1. Performance versus speed on ICDAR 2015 [15] text localization challenge. As can be seen, our algorithm significantly surpasses competitors in accuracy, whilst running very fast. The specifications of hardware used are listed in Tab. 6.

features are manually designed [5, 25, 40, 10, 26, 45] to capture the properties of scene text, while in deep learning based methods [3, 13, 11, 12, 7, 48] effective features are directly learned from training data.

However, existing methods, either conventional or deep neural network based, mostly consist of several stages and components, which are probably sub-optimal and time-

EAST



(a)

(b)

(c)

Figure 5. Qualitative results of the proposed algorithm. (a) ICDAR 2015. (b) MSRA-TD500. (c) COCO-Text.



EAST

Gradient Boosting

In subject area: Computer Science

Gradient boosting is a type of ensemble supervised machine learning algorithm that combines multiple weak learners to create a final model. It sequentially trains these models by placing more weights on instances with erroneous predictions, gradually minimizing a loss function. The predictions of the weak learners are compared with actual value, and the difference represents the error rate of the model. This error rate is used to calculate the gradient, which is used to find the direction for model parameter adjustment in the next round of training. Unlike a neural network model, where a single model is used, gradient boosting combines the predictions of multiple models to minimize the error.

AI generated definition based on: Machine Learning Guide for Oil and Gas Using Python, 2021

[About this page](#)[Add to Mendeley](#)[Set alert](#)[Discover other topics >](#)

On this page

[Chapters and Articles](#)[Related terms](#)[Recommended publications](#)[Featured Authors](#)

Chapters and Articles

You might find these chapters and articles relevant to this topic

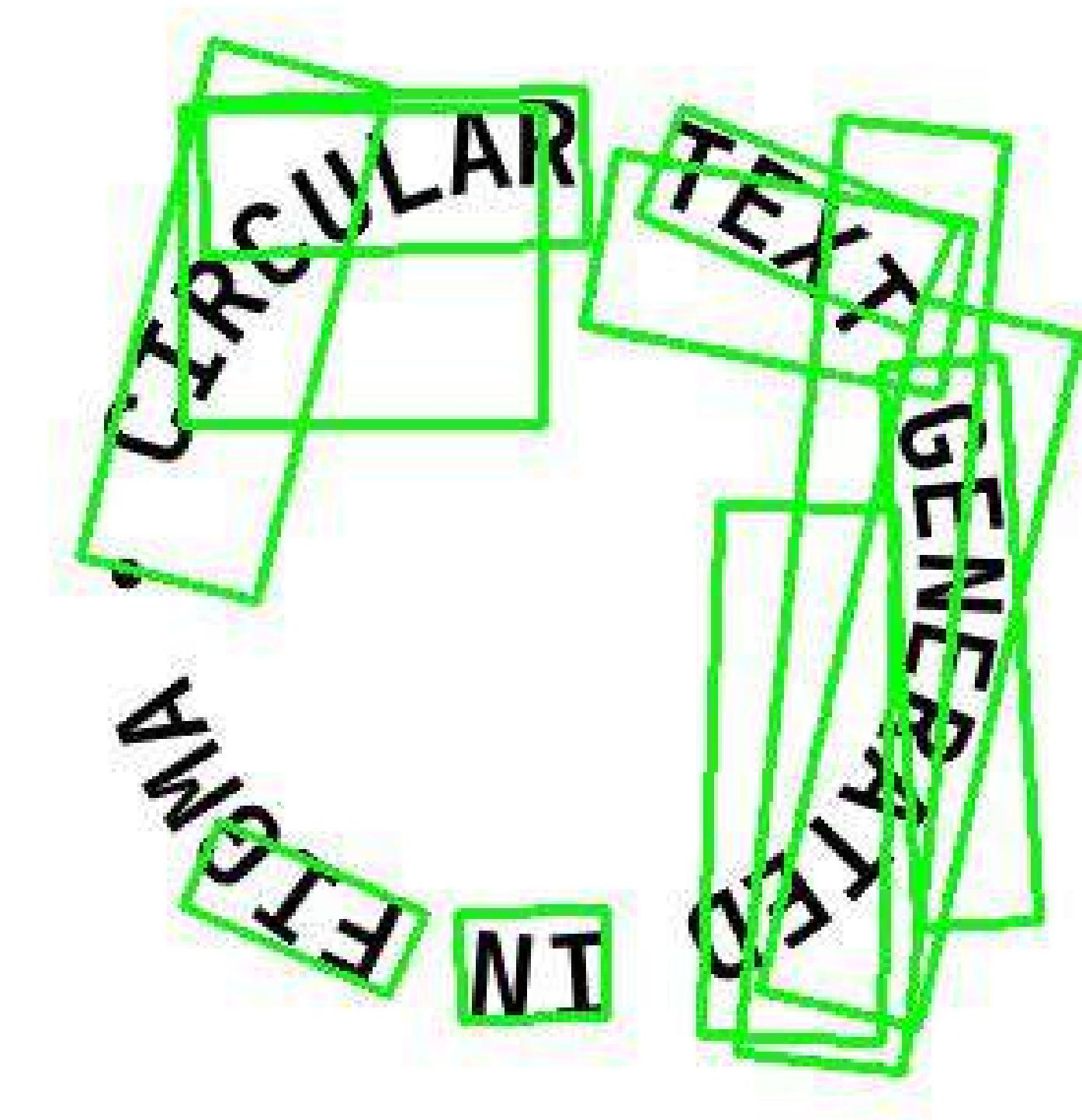
Supervised learning

Hosseini Belvadi, Alireza Haghhighat, in *Machine Learning Guide for Oil and Gas Using Python*, 2021

Gradient boosting

Gradient boosting is another type of ensemble supervised ML algorithm that can be used for both classification and regression problems. The main reason why algorithms such as random forest, extra trees, gradient boosting, etc., are called ensemble's is because a final model is generated based on many individual models. As described under the bagging/boosting comparison section, gradient boosting will train many models sequentially by placing more weights on instances with erroneous predictions. Therefore, challenging cases are the focus

EAST



EAST

PROTEIN

- for growth and development

High protein
Diet (Animals & Plant source)

ALTERED NUTRITION:

- less than body requirements UNDERWEIGHT PATIENTS
- more than body requirements OVERWEIGHT PATIENTS

- TISSUE REPAIR

BUILDING BLOCKS (hormones, enzymes, vit & minerals)

- COMPLEX STRUCTURE / LARGE MOLECULE

Breakdown into AMINO ACIDS → small molecule
(22 usually) simple substance



Hepatitis

- increase protein

Alcoholic liver disease

Liver cancer

- decrease protein

Undergo metabolism:

✓ catabolism (breakdown/digestion)

✓ anabolism (building up)

WASTE PRODUCT OF PROTEIN

AMMONIA

goes to

LIVER

→ UREA

Problem: Blood Chem

↑ ammonia (Hepatic Encephalopathy)

KIDNEYS

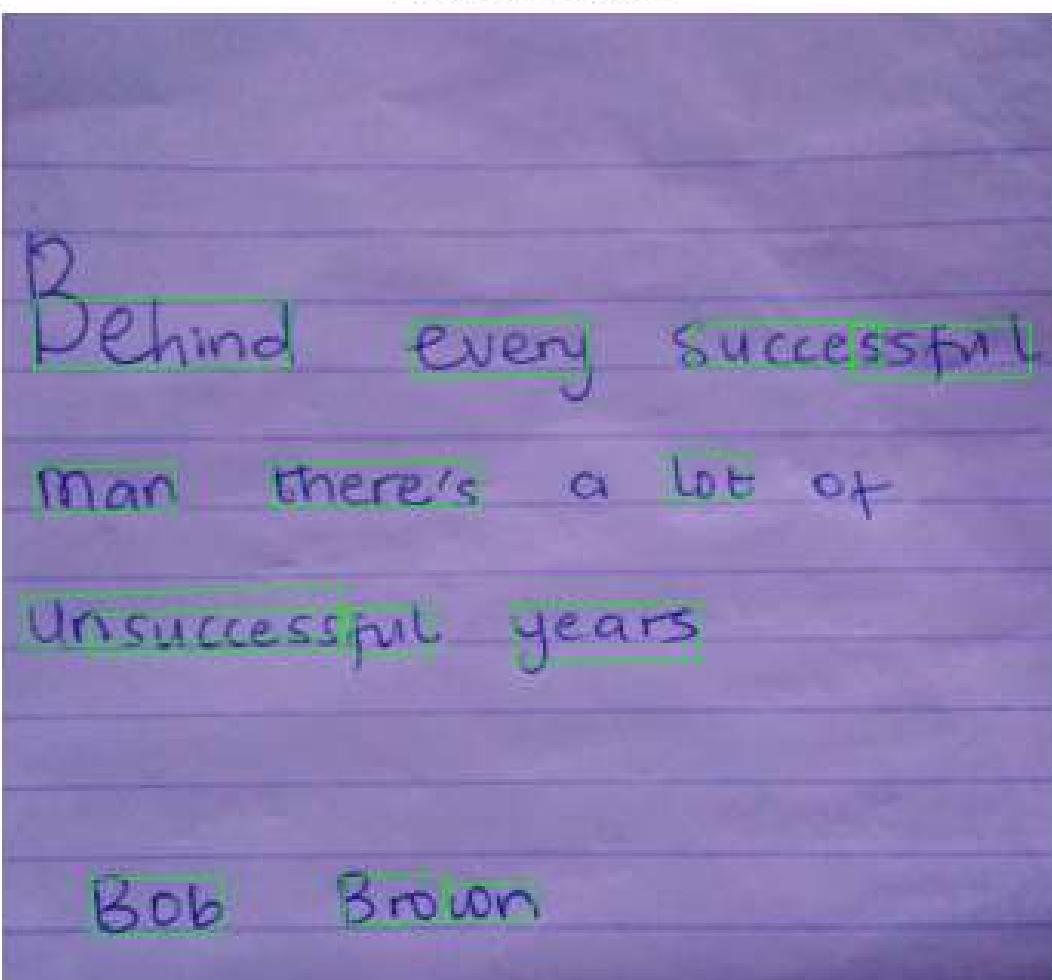
Behind every successful man there's a lot of unsuccessful years

Bob Brown

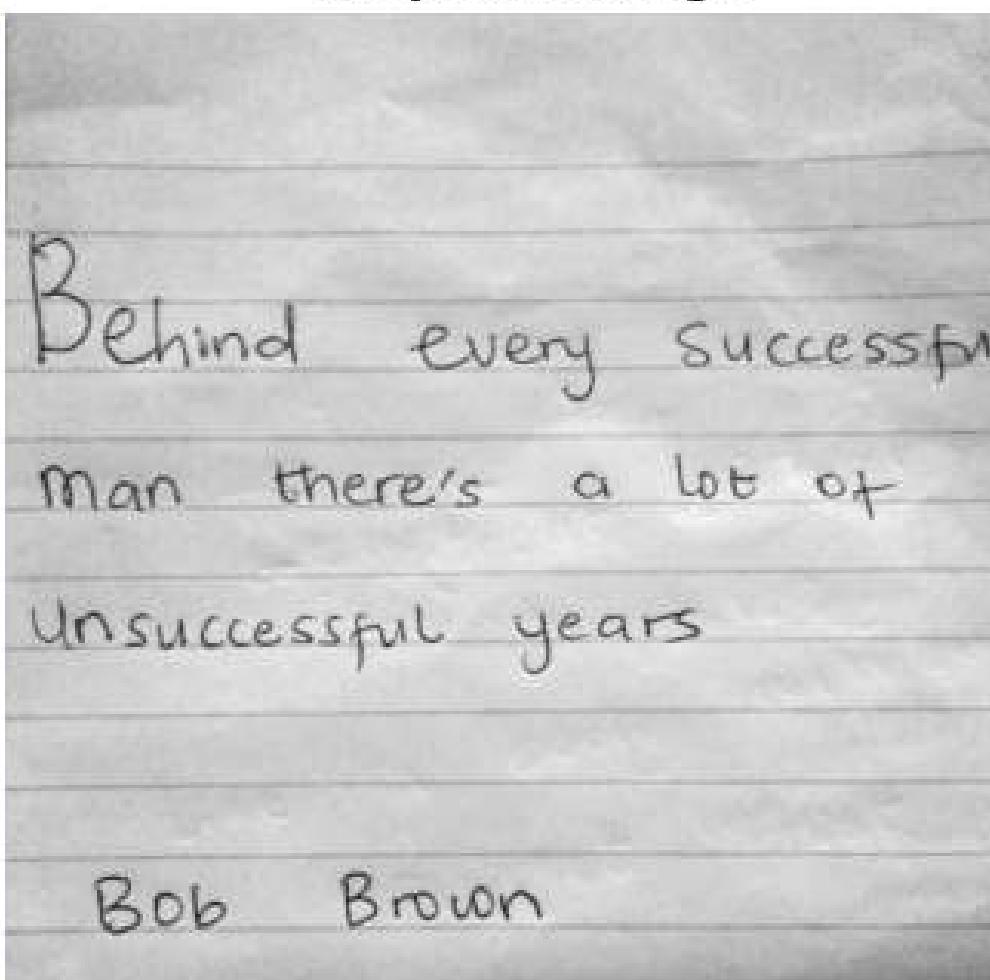
EAST: 69.0

Grayscale Image

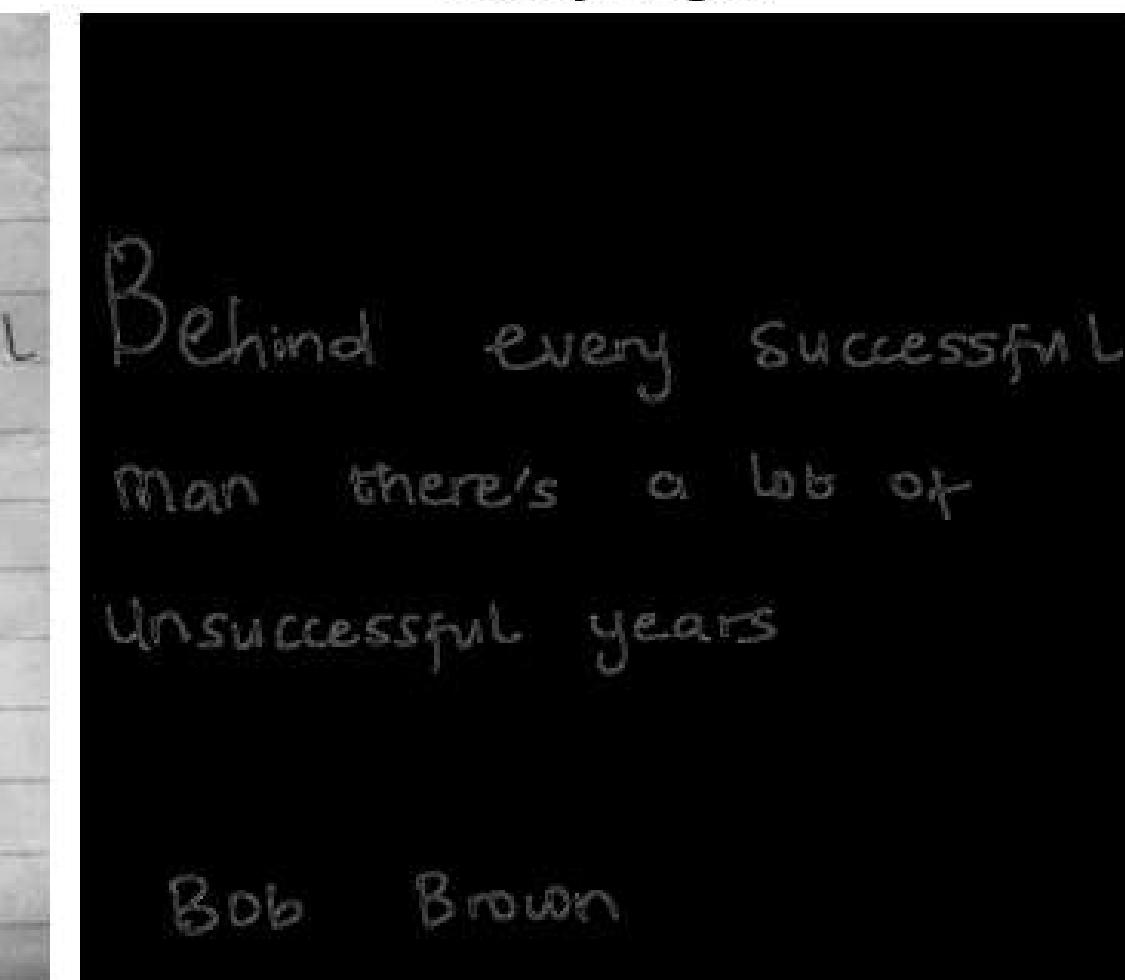
Canny Edges



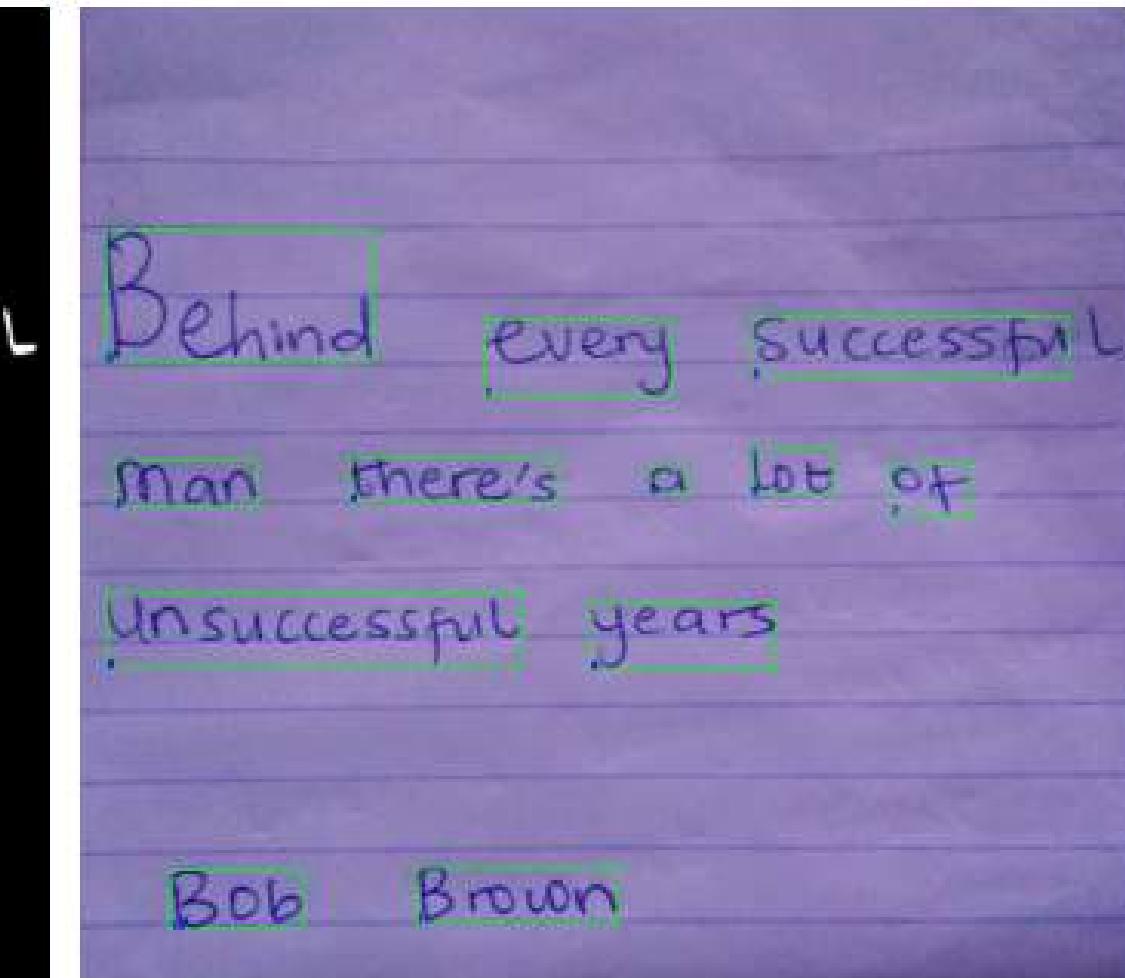
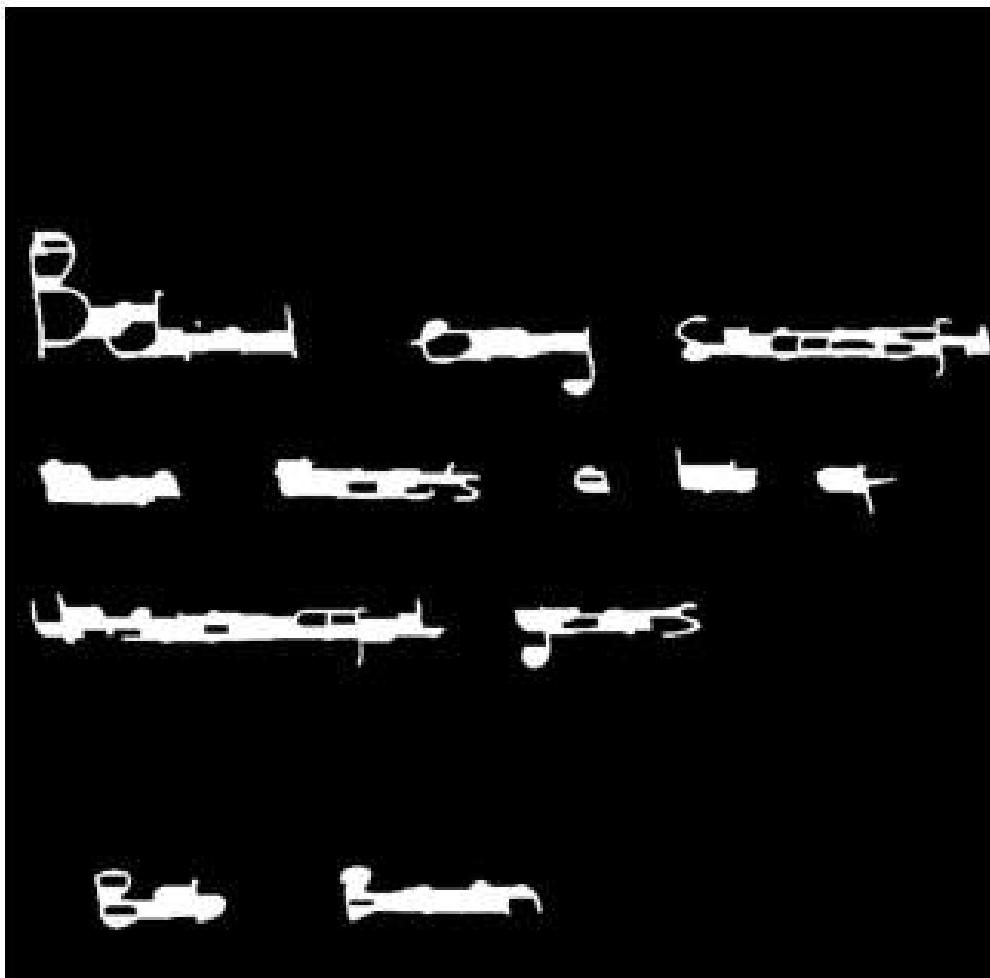
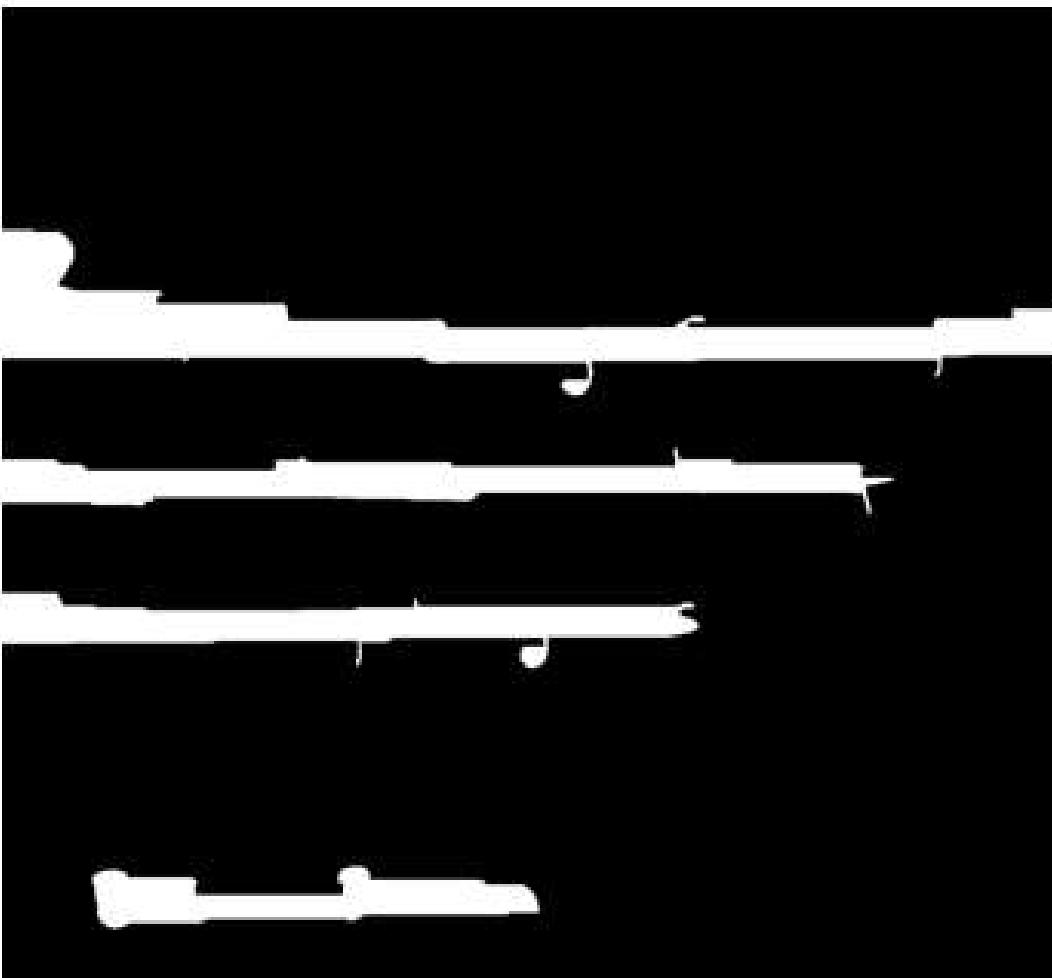
Morphological Line Detection



Morphological Word Detection



Result



General Algorithm of the applied Text Detection

1. Estimation of average text height using EAST
2. Morphological transformation finds text regions.
3. Visualization and cropping out text regions for further analysis.

The approach combines the advantages of deep model (EAST) with classical image processing methods.

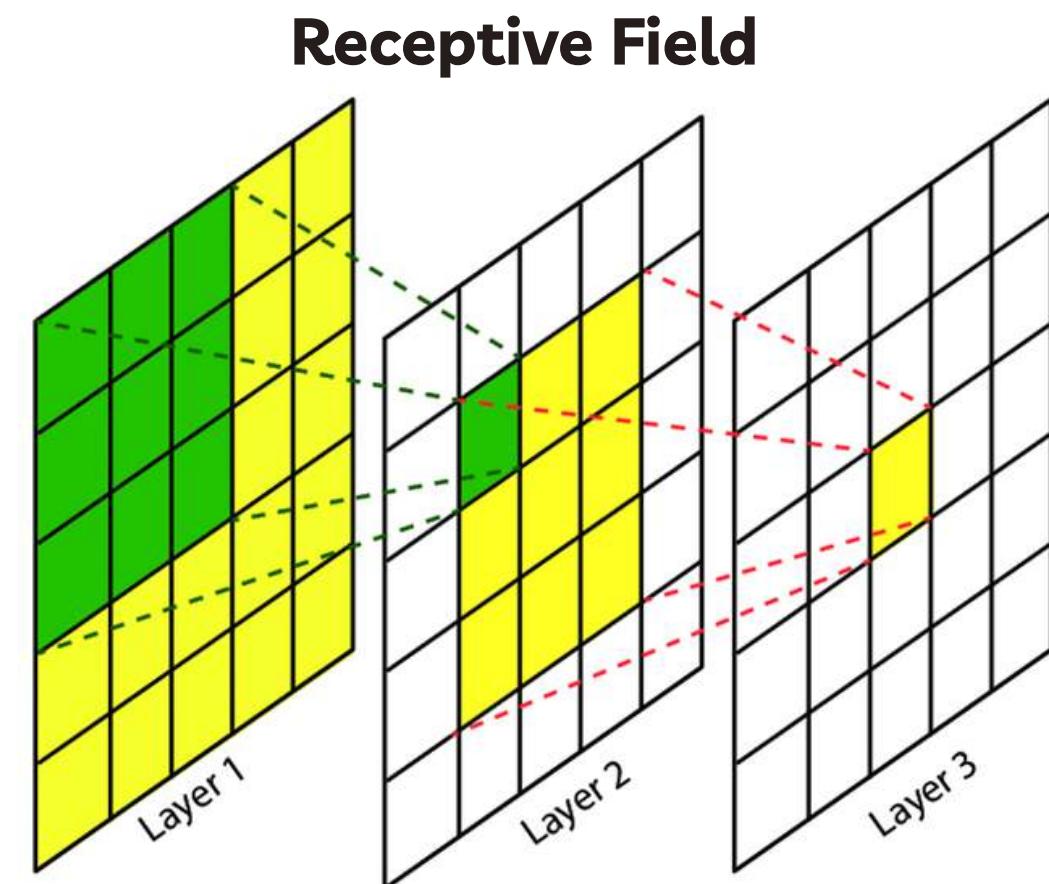
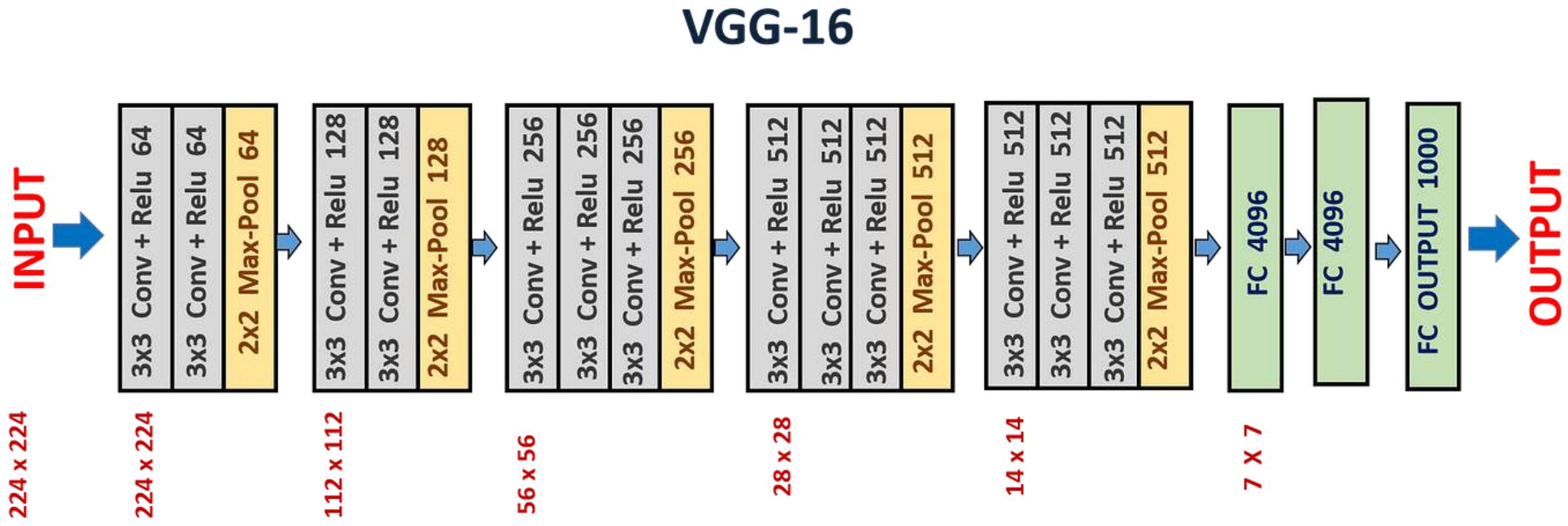
Limitations

The method is designed to detect individual words on white paper, where the lines of text are written horizontally. This is a limitation, because handwritten text can be written in different structures, in a circle, in tables, and there can be different formulas and this requires a different research.



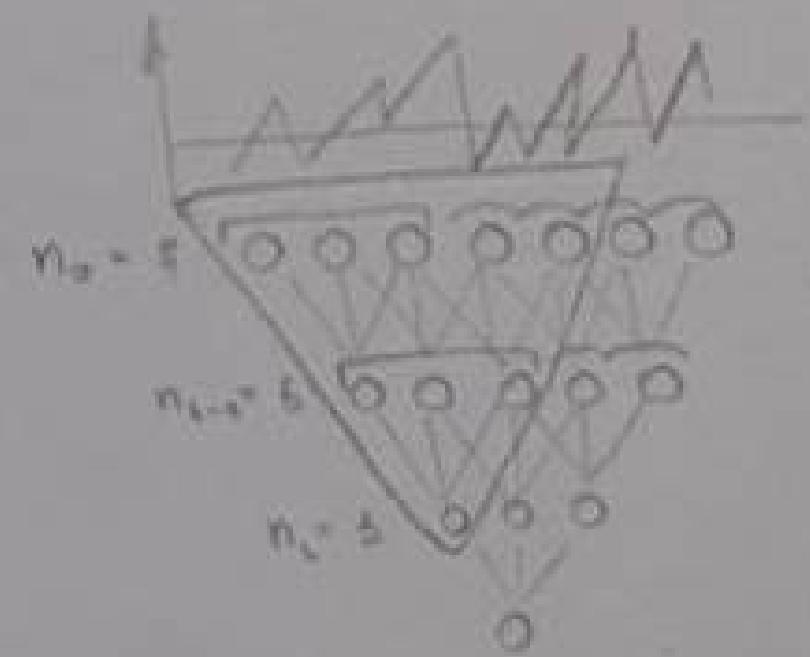
Kazakh Alphabet Letters (CMNIST)

<https://github.com/bolattleubayev/cmnist>

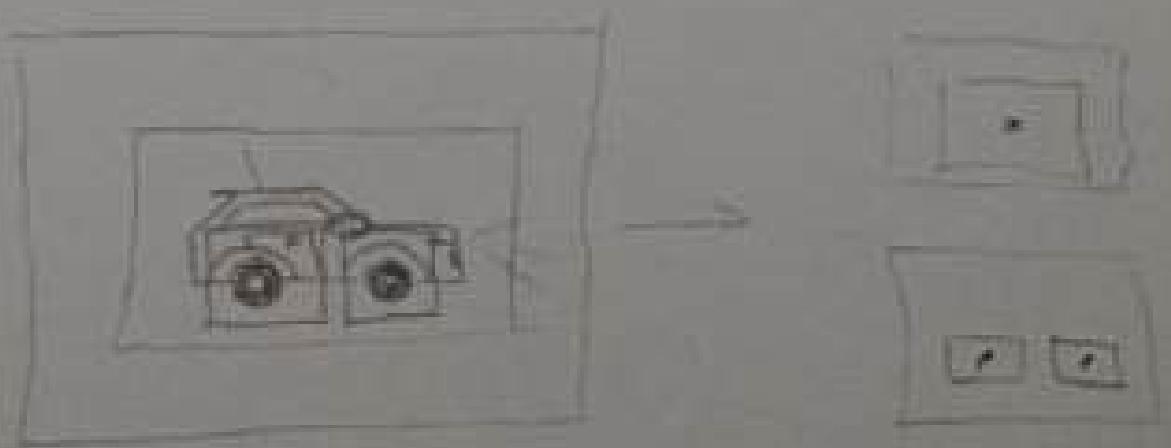


- 1) - 8

входные изображ.

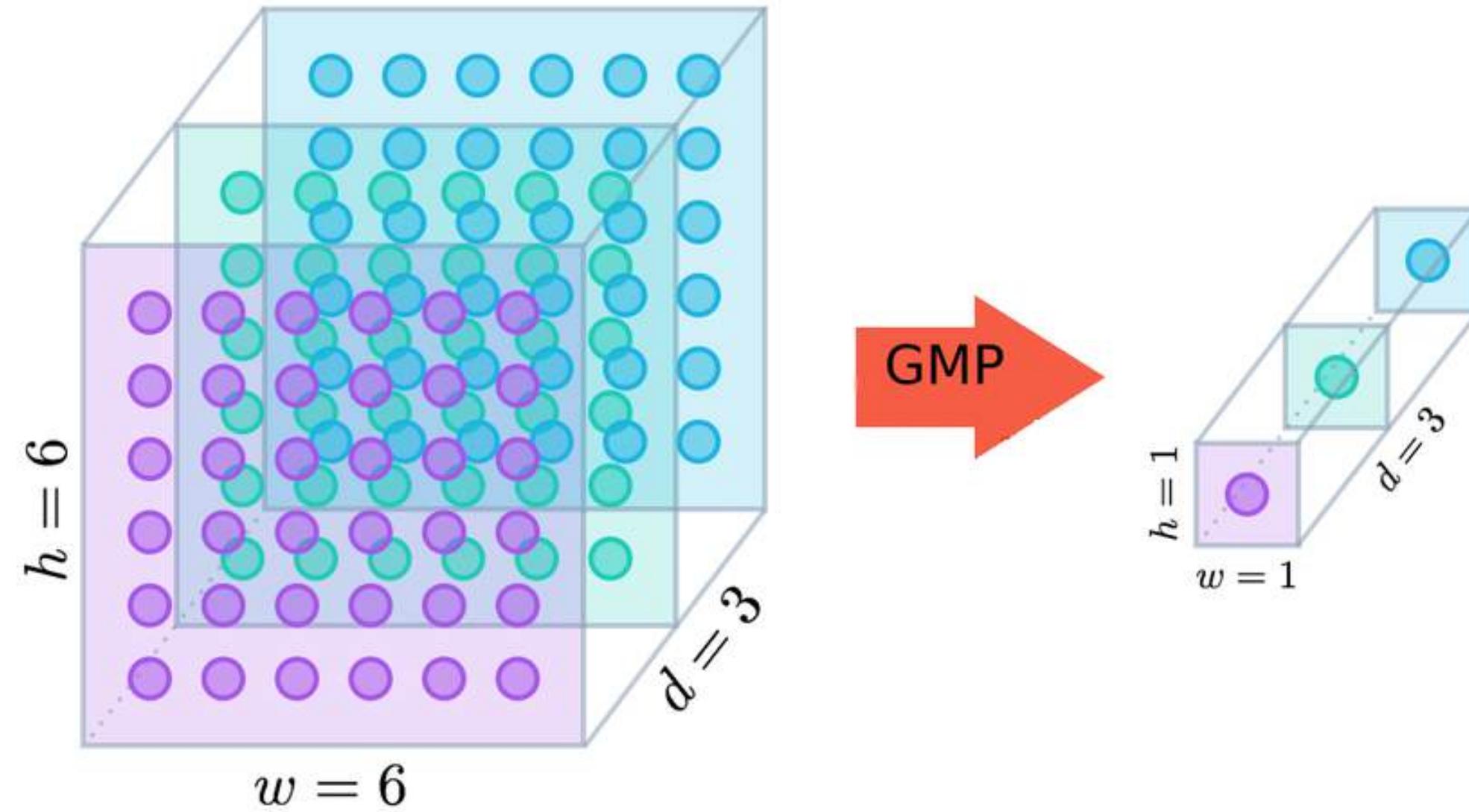


Receptive field нейронов можно рассматривать
изображения, или относительное расположение пикселей.



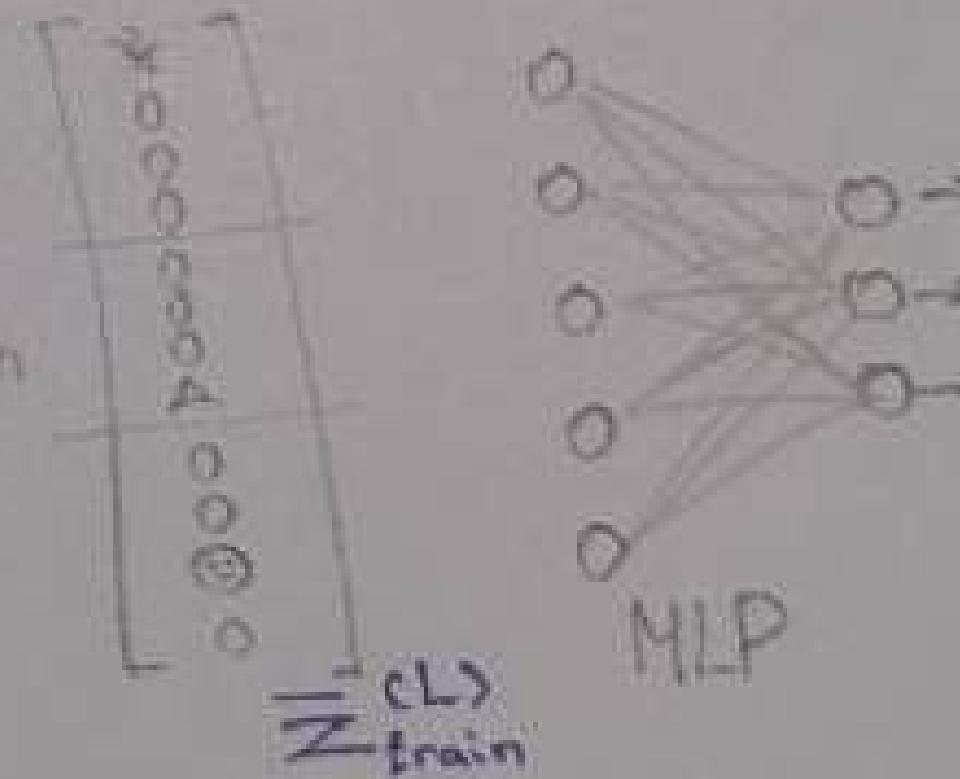
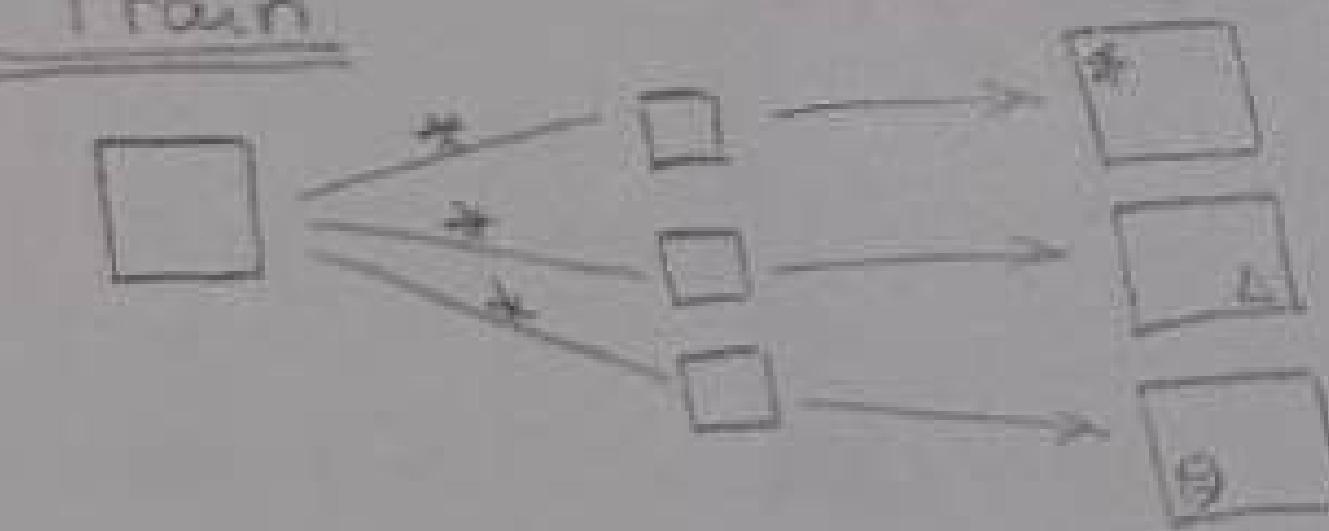
Receptive Field

3	3 × 64 × 64
5	64 × 64 × 64
10	
12	64 × 32 × 32
14	128 × 32 × 32
18	
30	128 × 16 × 16
32	256 × 16 × 16
64	
66	256 × 8 × 8
68	512 × 8 × 8

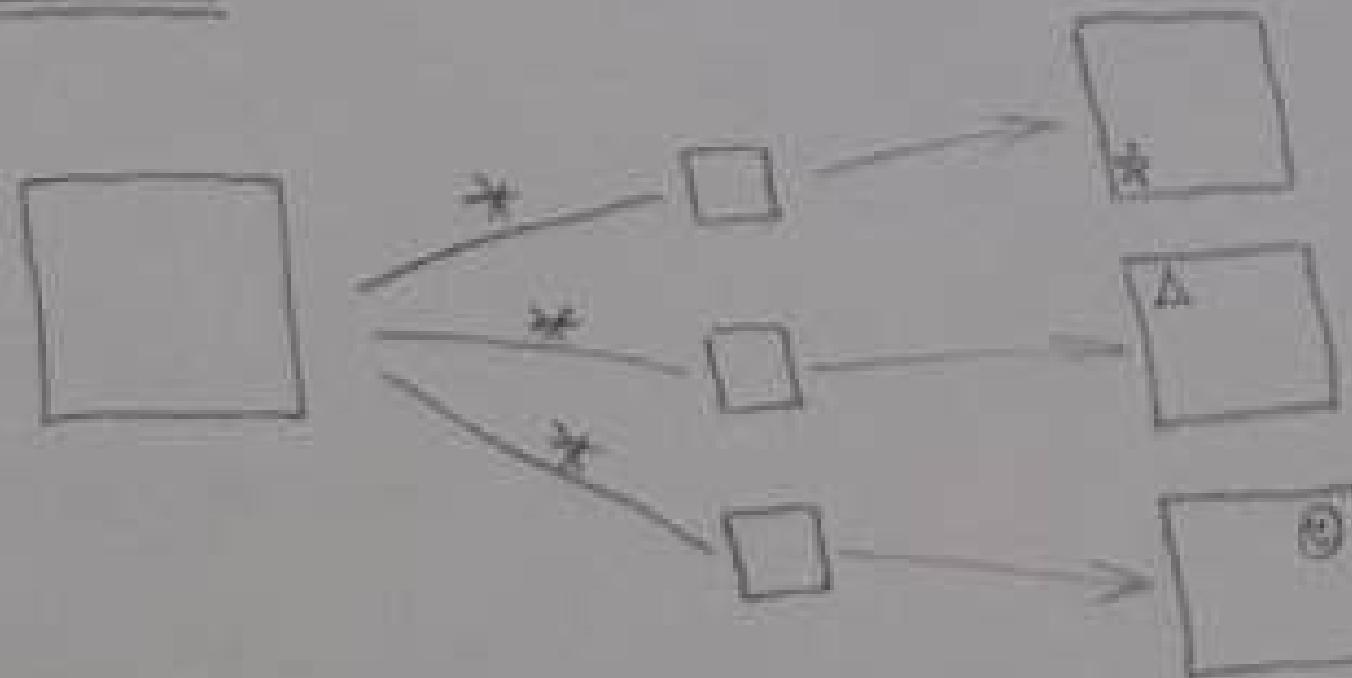


Why to use GMP

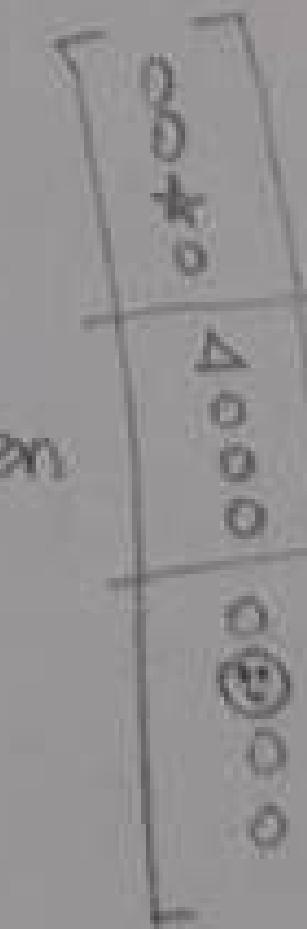
T rain



10

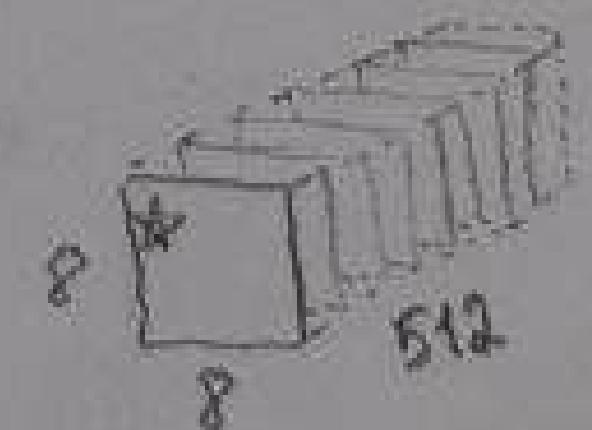


860



11

$$\bar{z}_{\text{train}}^{(l)} \neq \bar{z}_{\text{test}}^{(l)}$$

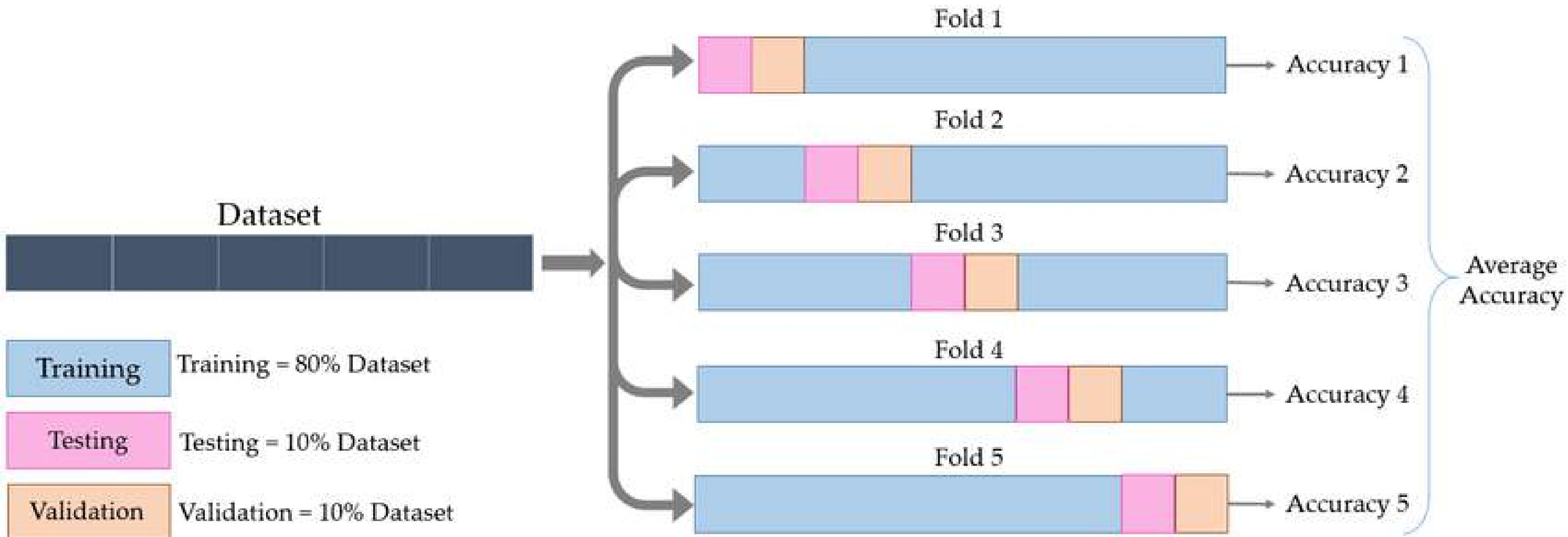


1) GMP (4,4) 00000512

2) GMP (1)

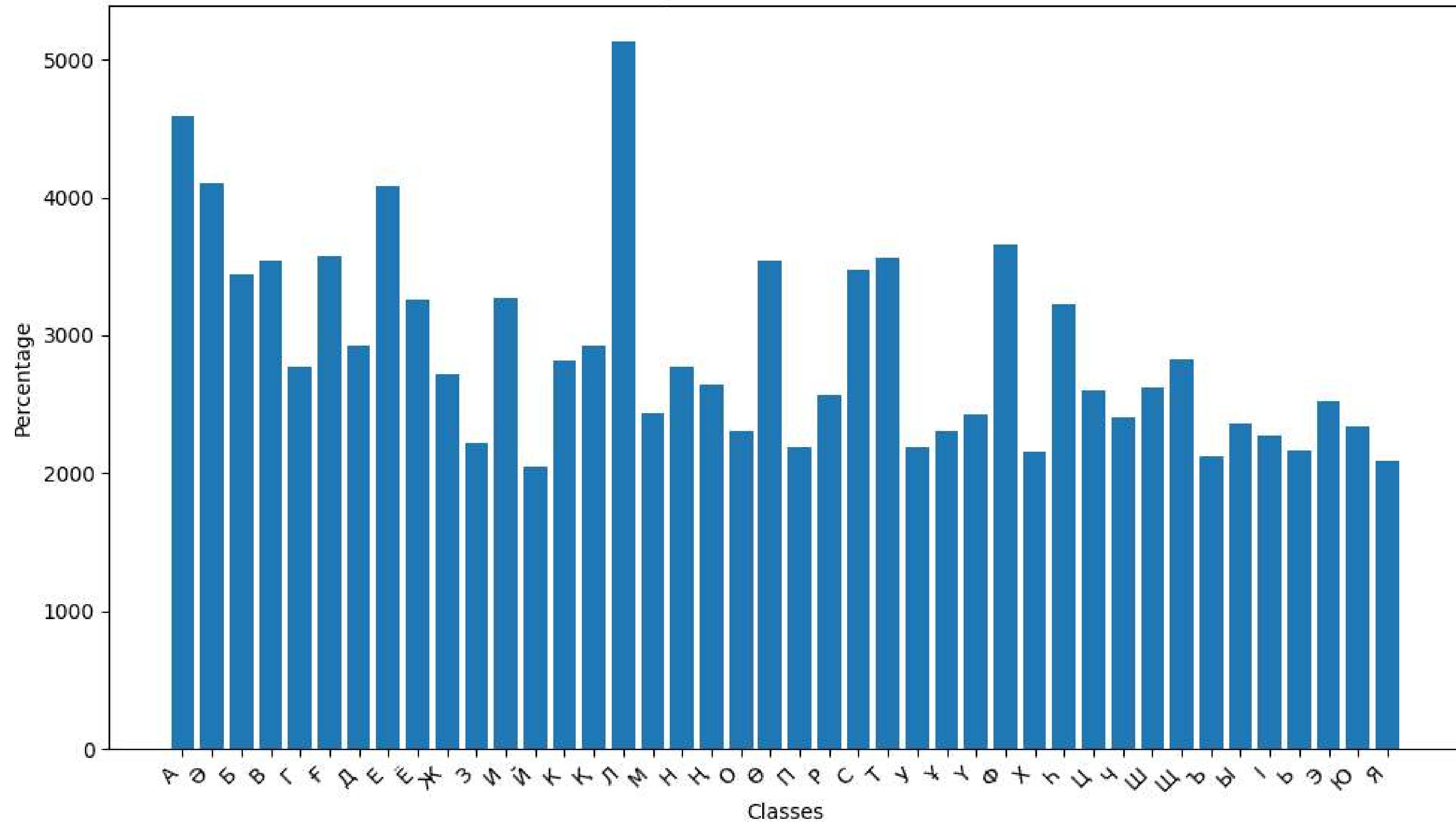


A black and white photograph of a simple metal cage. On the left side of the cage, the number "64" is handwritten in dark ink. The cage has a rectangular frame with vertical bars.

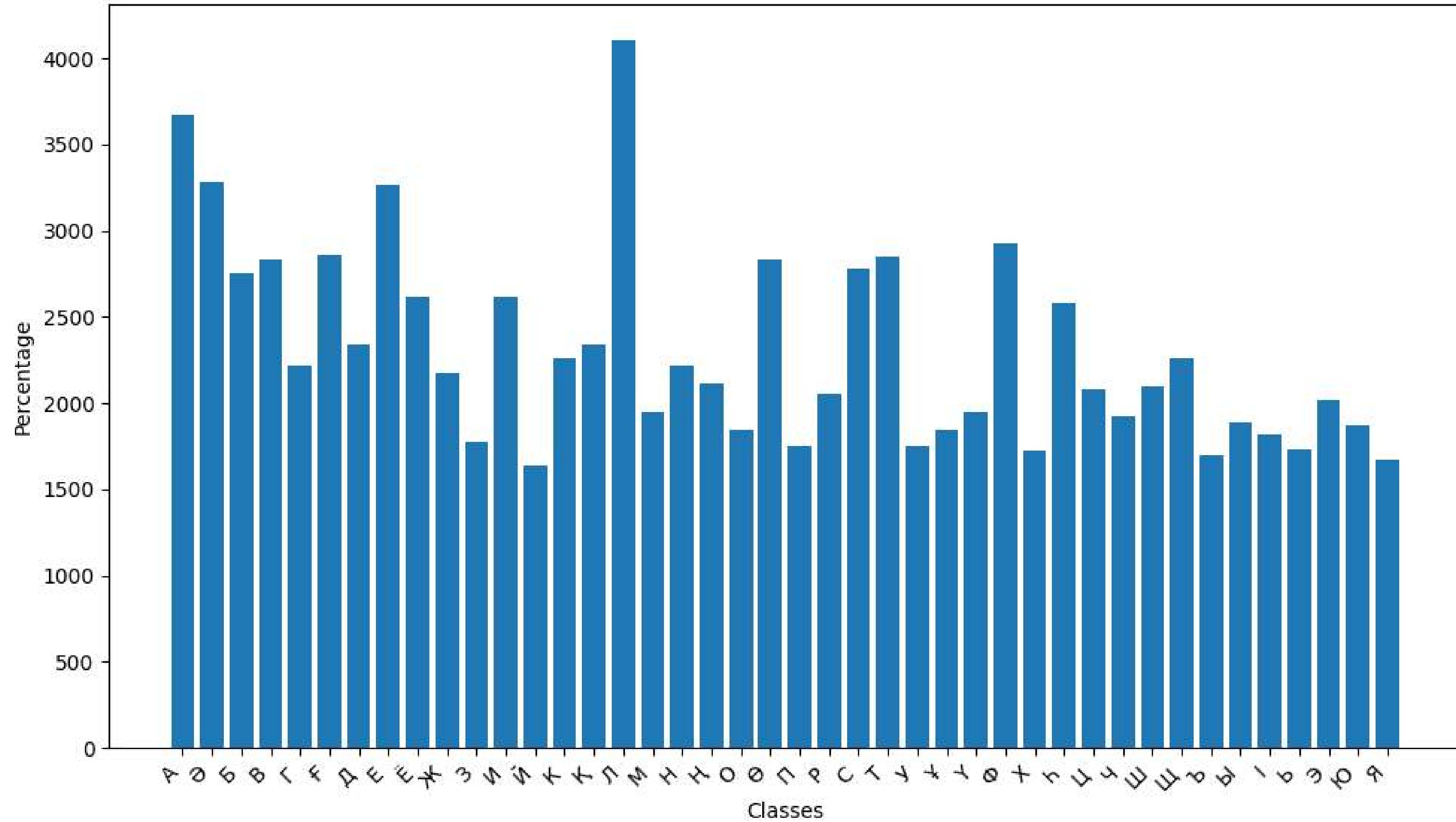


https://www.researchgate.net/figure/sual-representation-of-the-training-test-and-validation-split-using-cross-validation_fig3_338676025

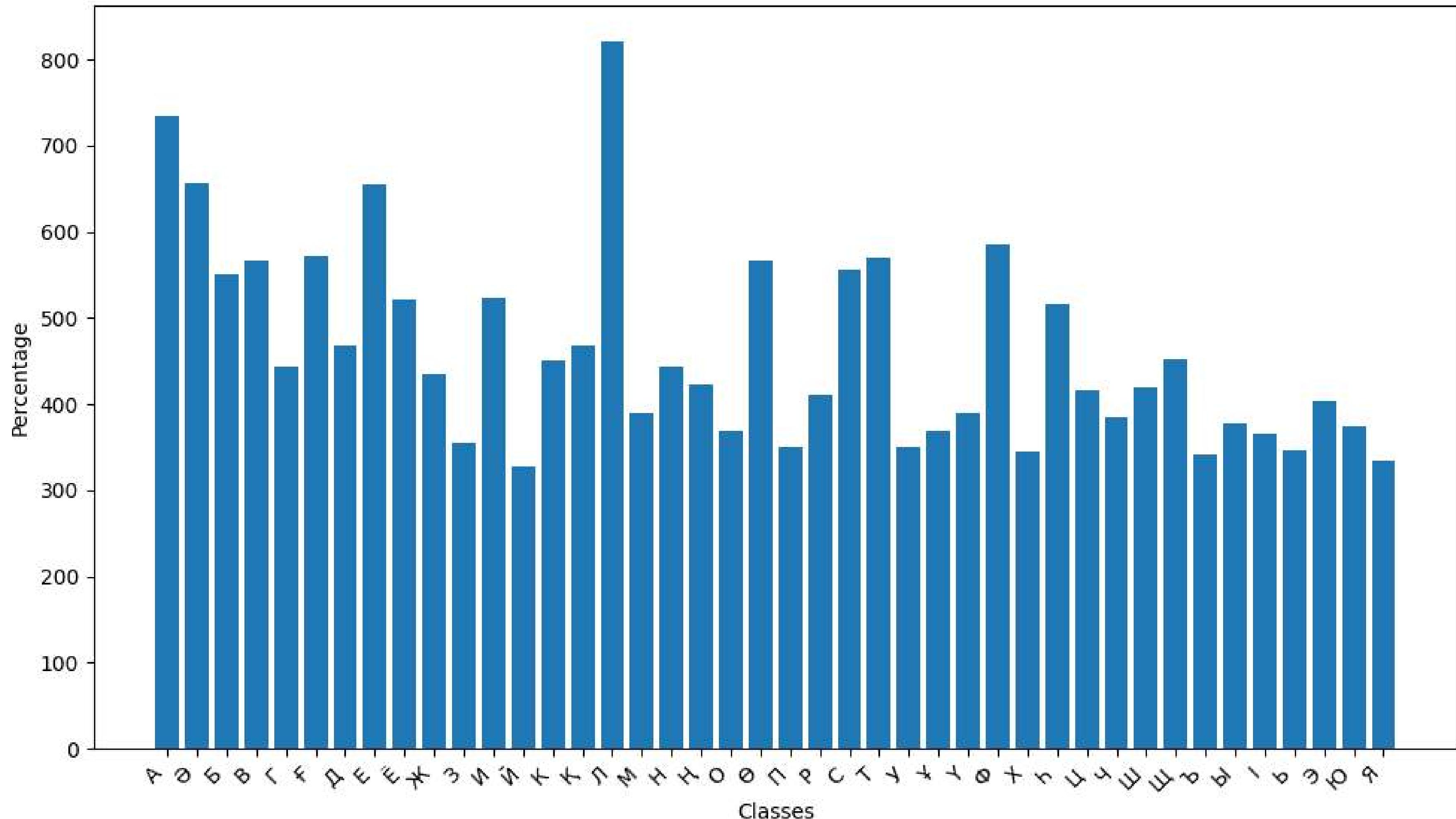
Original Class Distribution



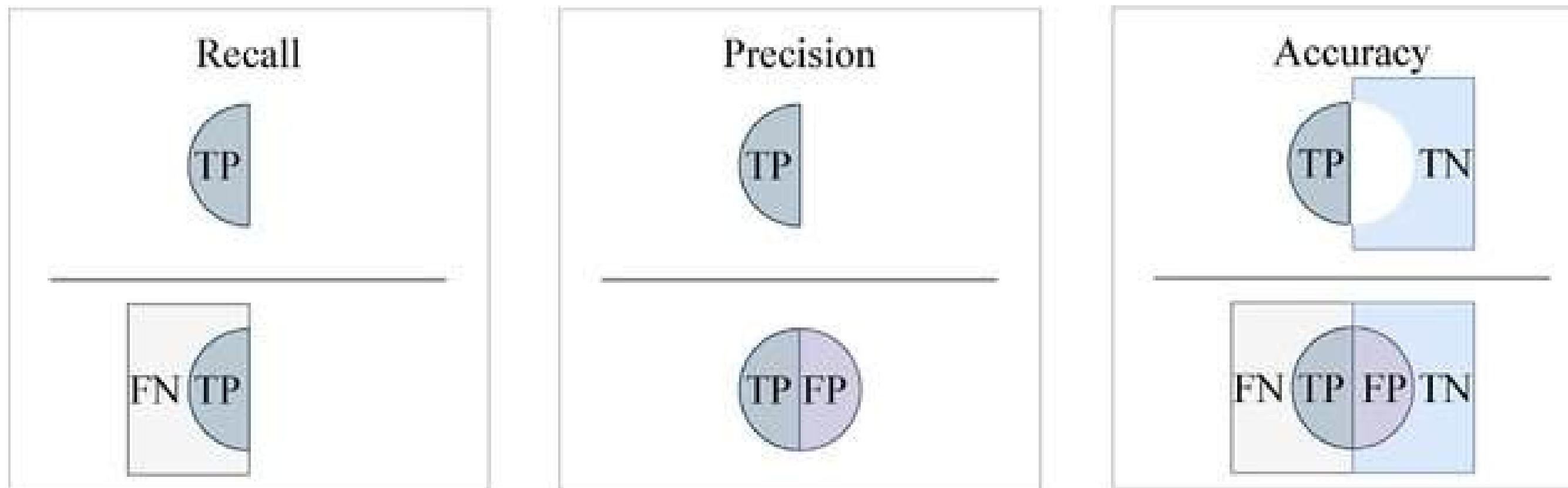
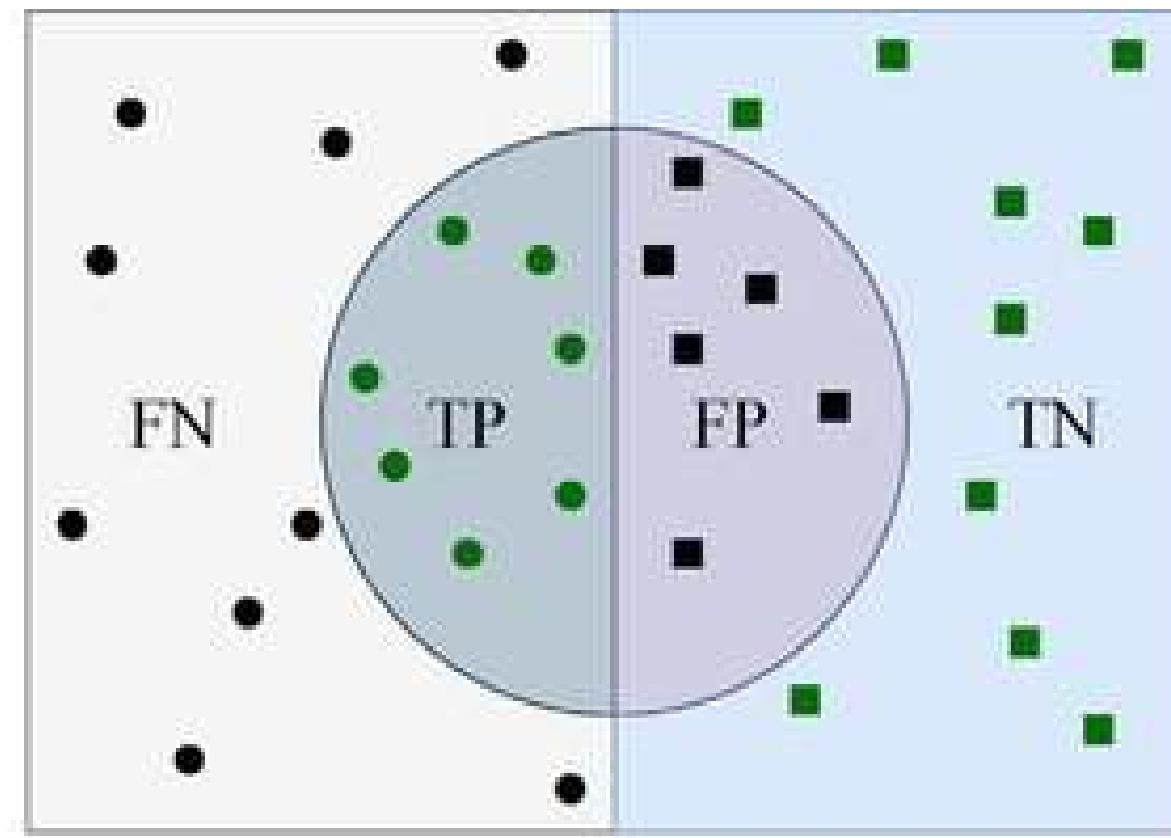
Train Class Distribution



KFold Validation Class Distribution



PRECISION & RECALL



F1-SCORE

Harmonic mean of the **precision** and **recall**

$$F_1 = \frac{2}{\text{recall}^{-1} + \text{precision}^{-1}} = 2 \frac{\text{precision} \cdot \text{recall}}{\text{precision} + \text{recall}} = \frac{2\text{TP}}{2\text{TP} + \text{FP} + \text{FN}}$$

EVALUATION METRICS

Character Error Rate (CER):

$$CER = \frac{\text{Number of incorrect characters}}{\text{Total number of characters in the reference text}} \times 100\%$$

Word Error Rate (WER):

$$WER = \frac{\text{Number of incorrect words}}{\text{Total number of words in the reference text}} \times 100\%$$

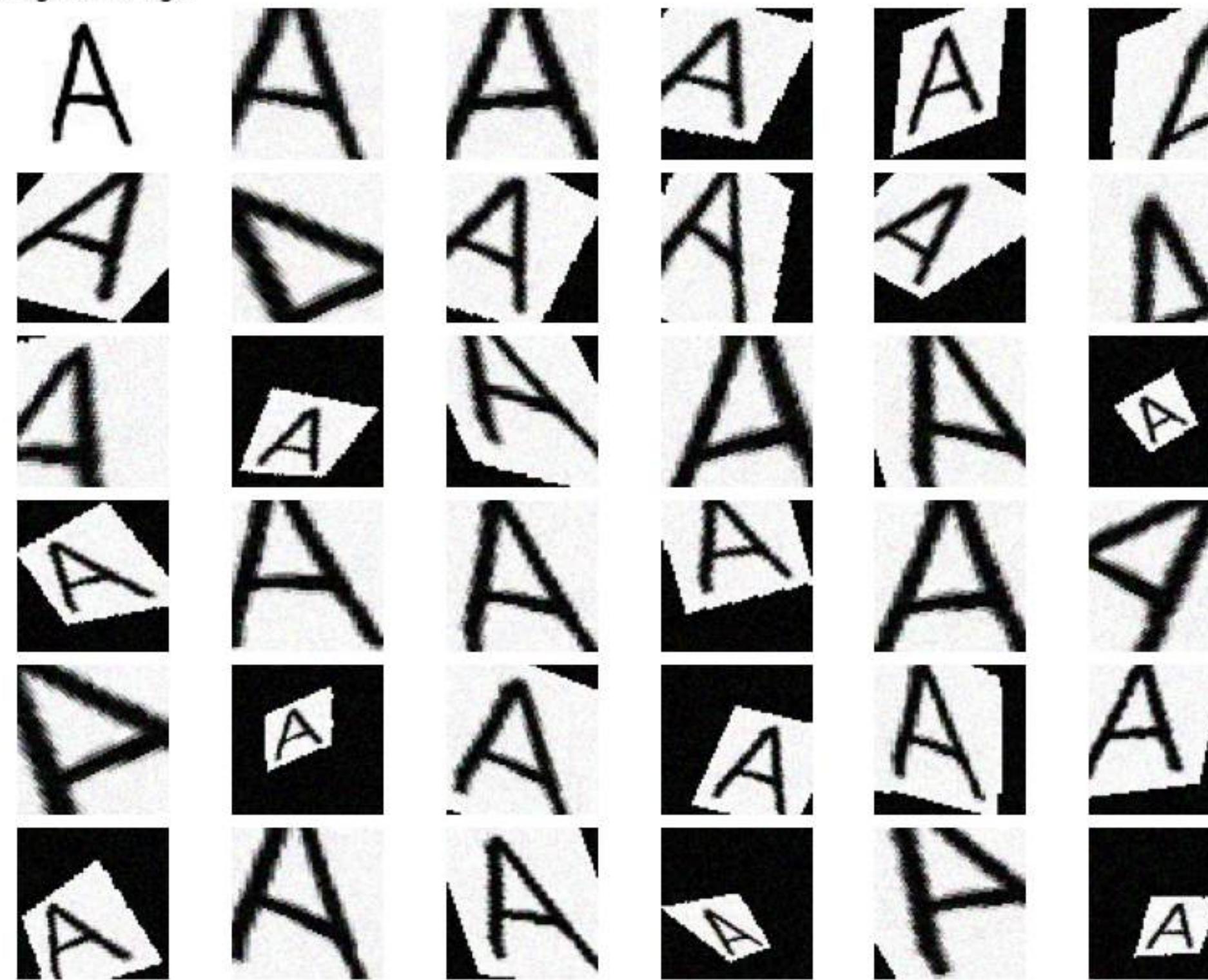
2.1. Label Error Rate

In this paper, we are interested in the following error measure: given a test set $S' \subset \mathcal{D}_{\mathcal{X} \times \mathcal{Z}}$ disjoint from S , define the *label error rate* (LER) of a temporal classifier h as the normalised edit distance between its classifications and the targets on S' , i.e.

$$LER(h, S') = \frac{1}{Z} \sum_{(\mathbf{x}, \mathbf{z}) \in S'} ED(h(\mathbf{x})) \quad (1)$$

where Z is the total number of target labels in S' , and $ED(\mathbf{p}, \mathbf{q})$ is the edit distance between the two sequences \mathbf{p} and \mathbf{q} — i.e. the minimum number of insertions, substitutions and deletions required to change \mathbf{p} into \mathbf{q} .

Original Image



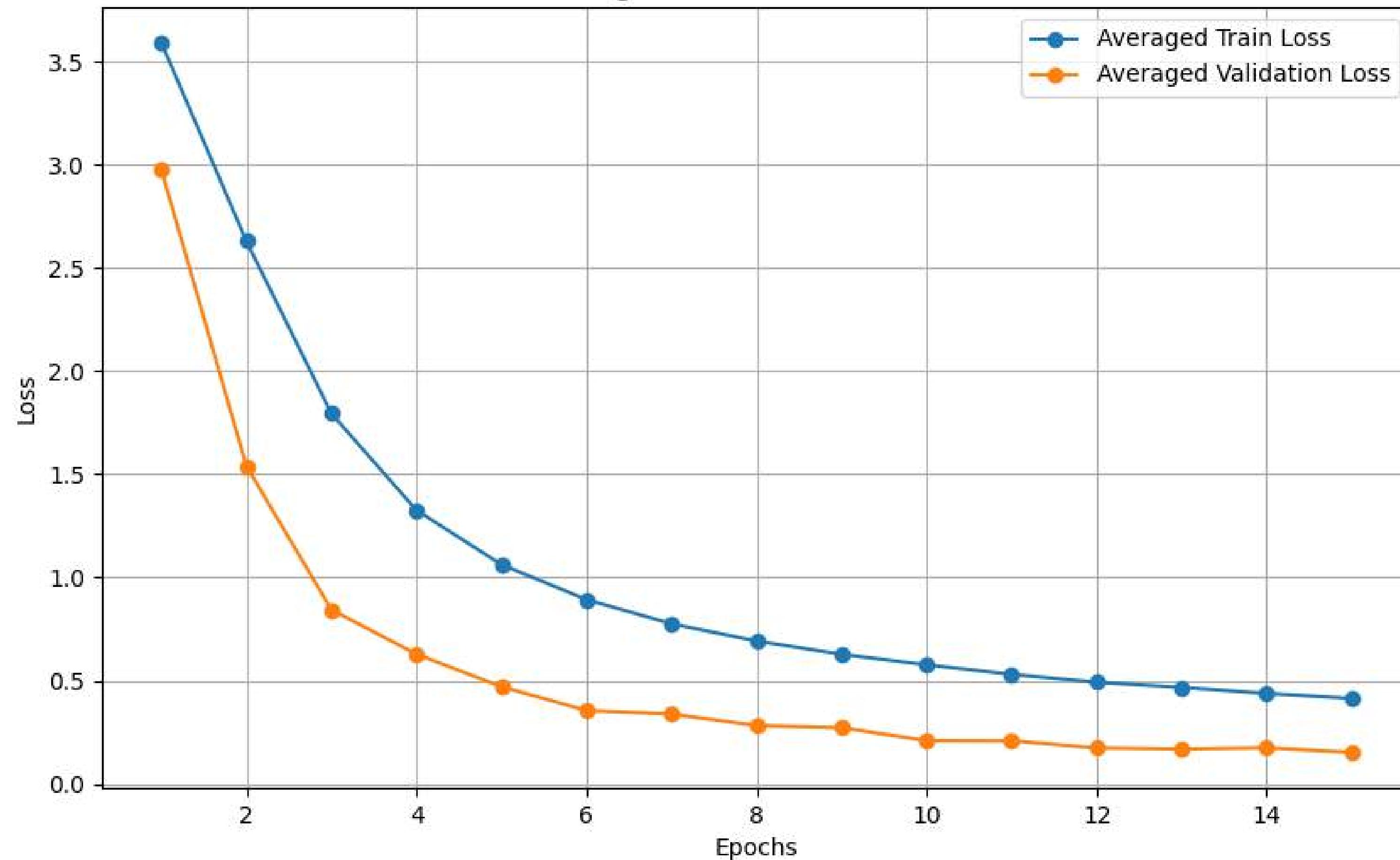
Data Augmentation

RandomAffine(rotate, translate, scale, shear)

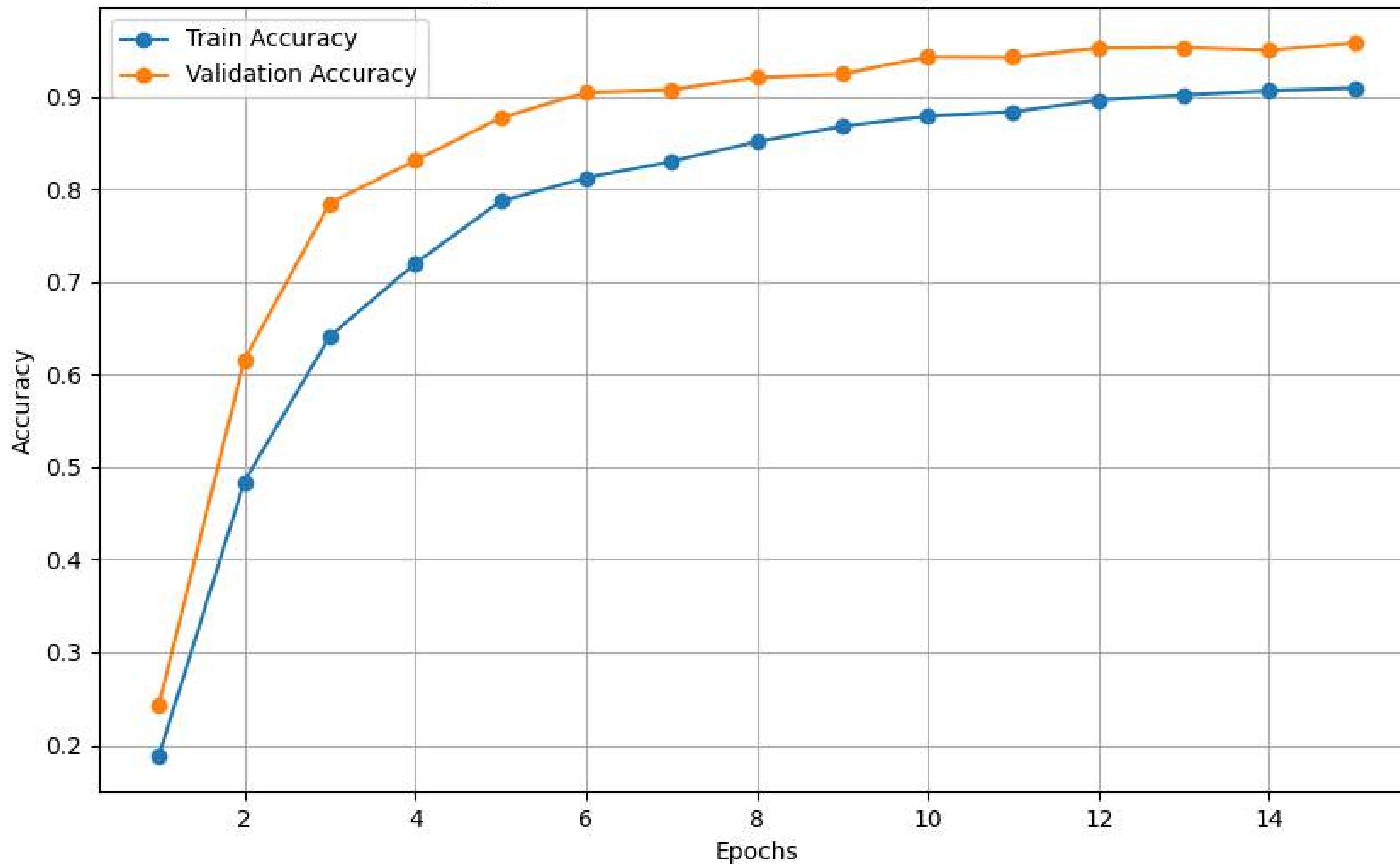
RandomPerspective

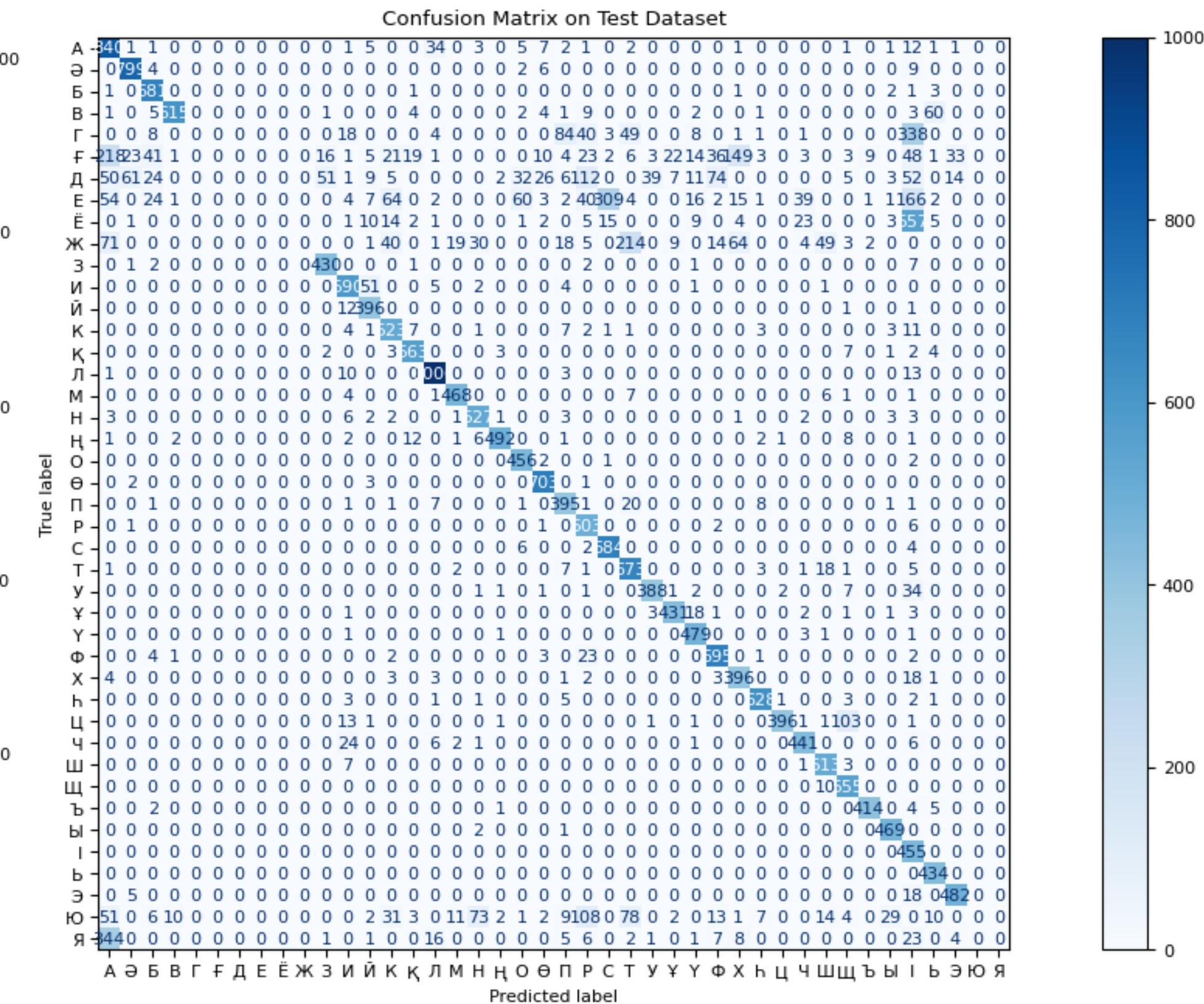
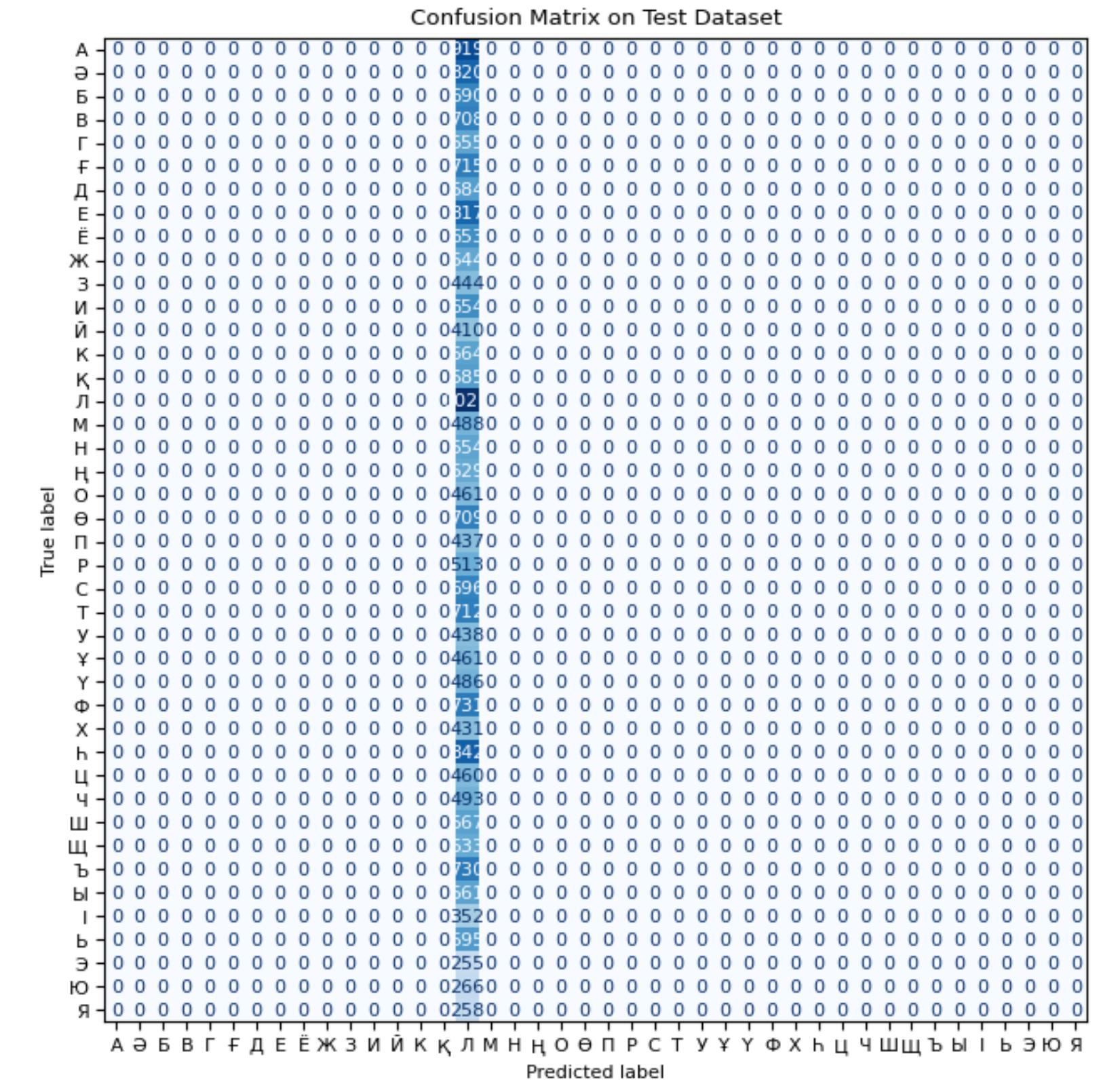
RandomNoise

Averaged Losses Across 5 Folds



Average Train and Validation Accuracy Across Folds





ImageFolder Nuance: 0, 1, 2, 3, 4 ... -> 0, 1, 10, 11, 2

Learning Duration: 4 hours 18

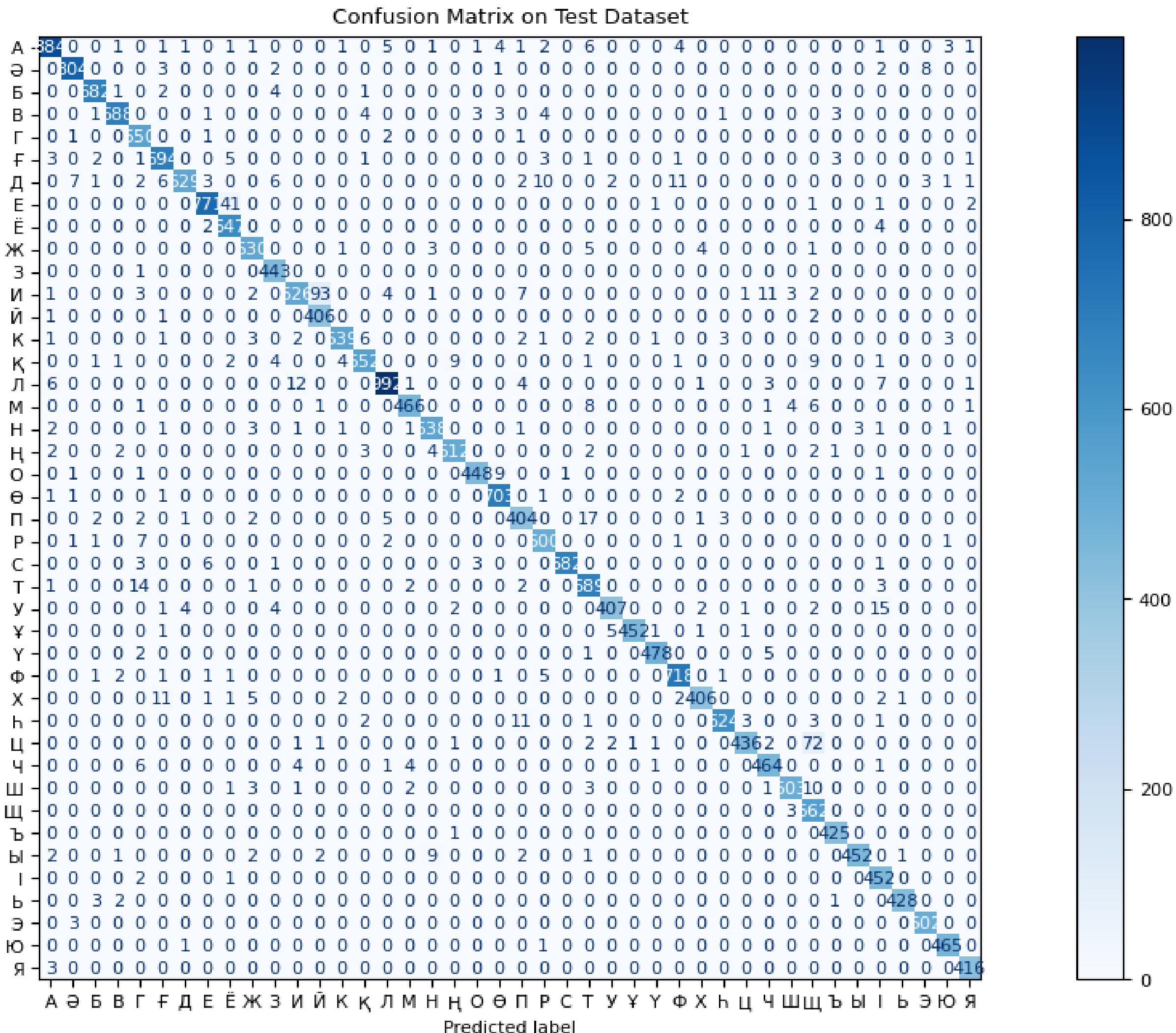
Best model is from fold 4 with

Validation F-Score: 0.9639

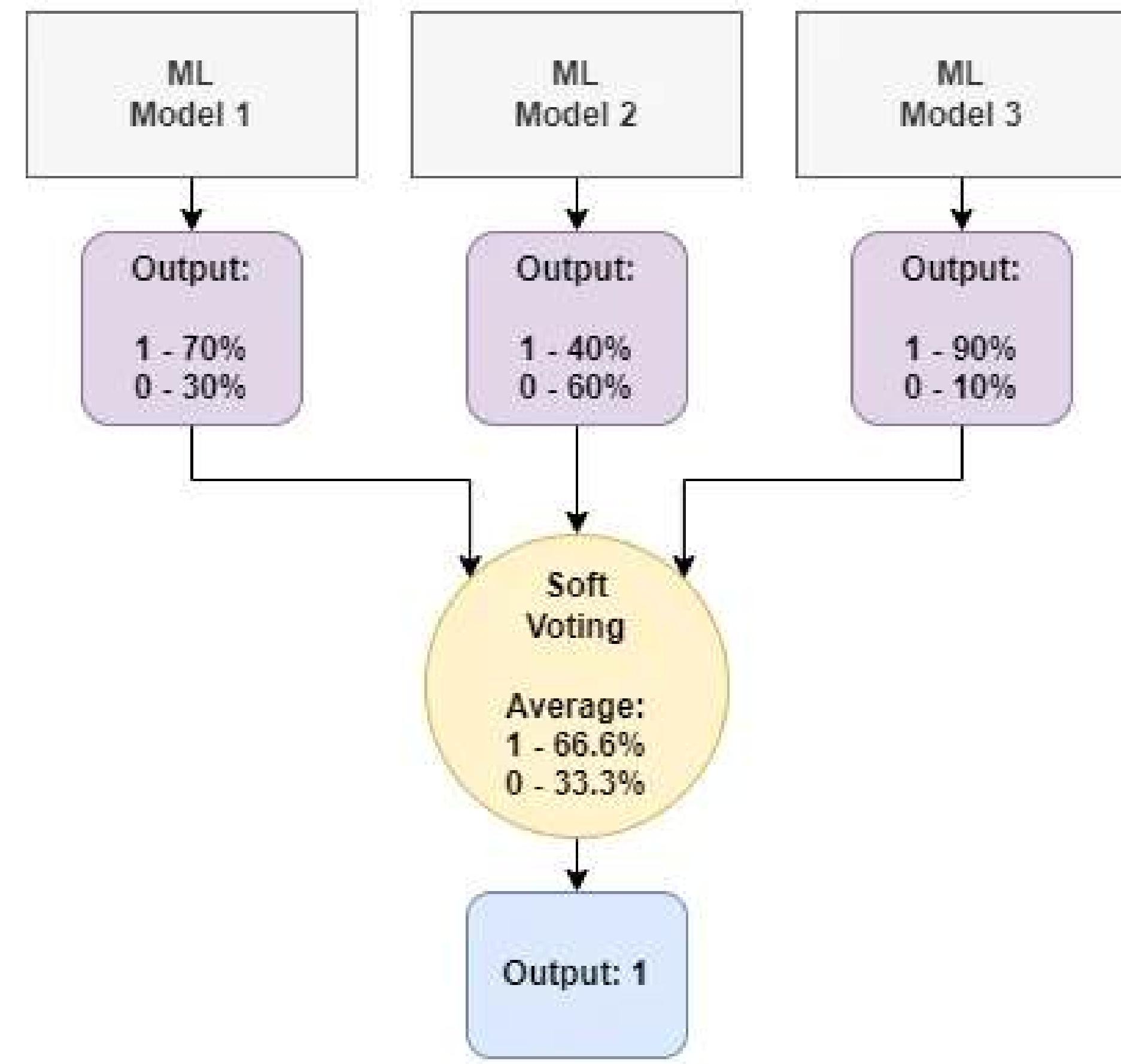
Test dataset: 24247

Test Accuracy: 96.38%

Test F1-Score: 96.33%



Voting



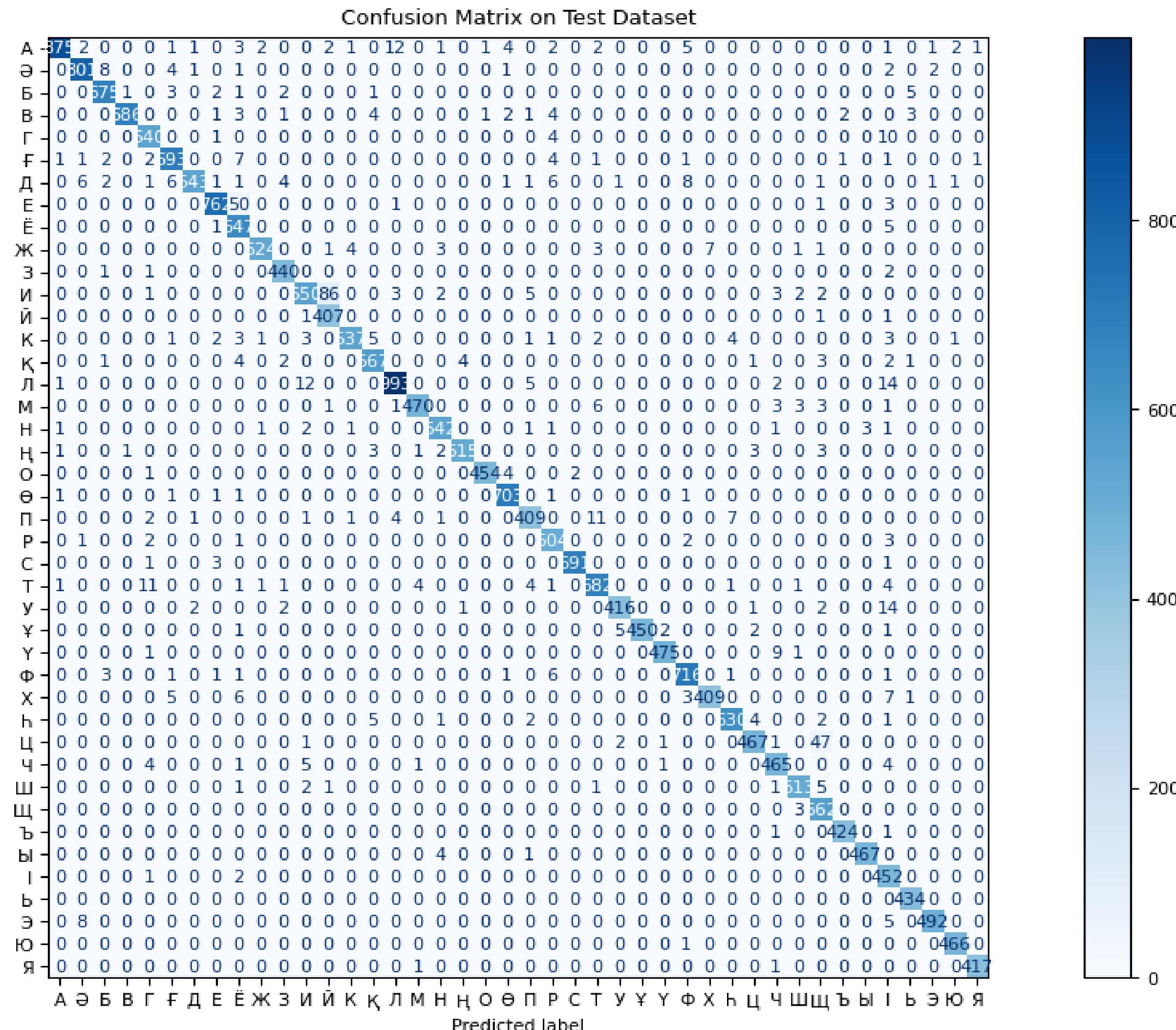
Fusion

Test dataset: 24247

Test Accuracy: 96.38% -> 96.77%

F1-Score: 96.33% → 96.78%

$$0.0039 * 24247 = 94 \text{ samples}$$



Fold 4

Fusion

KOHTD: Kazakh offline handwritten text dataset

Toiganbayeva, N., Kasem, M., Abdimanap, G., Bostanbekov, K., Abdallah, A., Alimova, A., & Nurseitov, D. (2022). Kohtd: Kazakh offline handwritten text dataset. *Signal Processing: Image Communication*, 108, 116827.

Chicago

Github: <https://github.com/abdoelsayed2016/KOHTD>

2) Шартсыз орталайсаныңу жағдайында
функцияның жарык мөндерекік
максимумы және максимиумын іздеүтін
жөнди.

Шартсыз орталайсаныңу жағдайында
 $f(x)$ -шарттың ~~функция~~ максатын жөнди
— x айналыштарындаңдардан анықташылған, DB жағдайында,
және басқарушының айналыштарындаңдардан
В мөндердегі максатын көз-кеңзен
үзүндең миссиян орталайсаныңу
жөндиндең жағдайын максатын жөнди
амалдайтын.

Миссияның: бастапқы нұхтеге ~~фрасы~~
қарашемен сабактайдын, Ол үшін
бірнеше деңгес міндеттіліктің мадасын.

$$\frac{df}{dx} = -2x_1 + 6; \frac{dE}{dx} = -8x_2 + 32$$

$$\nabla f(x) = \begin{pmatrix} -2x_1 + 6 \\ -8x_2 + 32 \end{pmatrix}, \nabla F(x_0) = \begin{pmatrix} -2 \cdot 7 + 6 \\ -8 \cdot 4 + 32 \end{pmatrix} = \begin{pmatrix} -2 \\ 0 \end{pmatrix}$$

Матса x_1 нұхтегін максатынаныз:

$$x_1 = x_0 + \lambda_1 \nabla f(x_0), \nabla f(x_0) = \begin{pmatrix} -2 \\ 0 \end{pmatrix} + \lambda_1 \begin{pmatrix} -2 \\ 0 \end{pmatrix} = \begin{pmatrix} -2 - 2\lambda_1 \\ 0 \end{pmatrix}$$

Матса нұхтеге қарашемни
мадасынаныз:

$$\nabla F(x_1) = \begin{pmatrix} -2 \cdot (-2 - 2\lambda_1) + 6 \\ -8 \cdot 0 + 32 \end{pmatrix} = \begin{pmatrix} 16\lambda_1 - 8 \\ 32 \end{pmatrix}$$

$$\frac{B \nabla F}{\lambda_1} = \nabla F(x_0) \cdot \nabla F(x_1) = \begin{pmatrix} -2 \\ 0 \end{pmatrix} \cdot \begin{pmatrix} 16\lambda_1 - 8 \\ 32 \end{pmatrix} = 0$$

мұндағын миссиянаныз.



Connectionist Temporal Classification: Labelling Unsegmented Sequence Data with Recurrent Neural Networks

Alex Graves¹

Santiago Fernández¹

Faustino Gomez¹

Jürgen Schmidhuber^{1,2}

ALEX@IDSIA.CH

SANTIAGO@IDSIA.CH

TINO@IDSIA.CH

JUERGEN@IDSIA.CH

¹ Istituto Dalle Molle di Studi sull’Intelligenza Artificiale (IDSIA), Galleria 2, 6928 Manno-Lugano, Switzerland

² Technische Universität München (TUM), Boltzmannstr. 3, 85748 Garching, Munich, Germany

Abstract

Many real-world sequence learning tasks require the prediction of sequences of labels from noisy, unsegmented input data. In speech recognition, for example, an acoustic signal is transcribed into words or sub-word units. Recurrent neural networks (RNNs) are powerful sequence learners that would seem well suited to such tasks. However, because they require pre-segmented training data, and post-processing to transform their outputs into label sequences, their applicability has so far been limited. This paper presents a novel method for training RNNs to label unsegmented sequences directly, thereby solving both problems. An experiment on the TIMIT speech corpus demonstrates its advantages over both a baseline HMM and a hybrid HMM-RNN.

1. Introduction

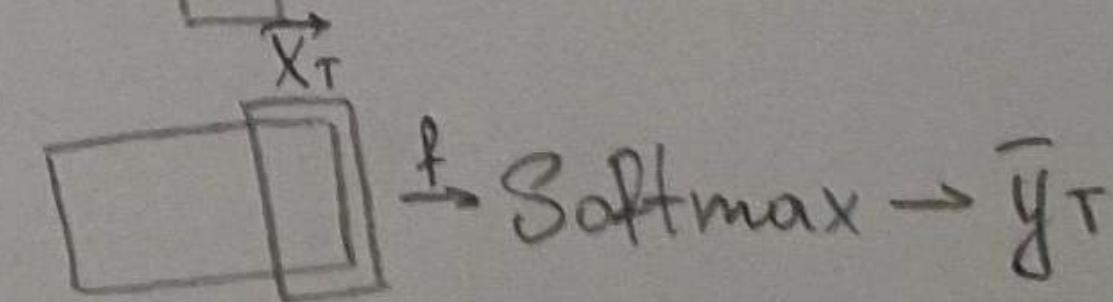
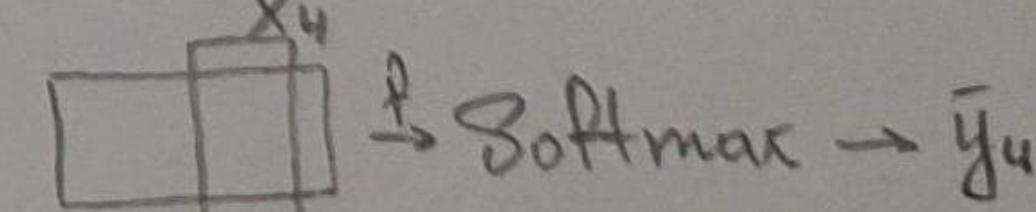
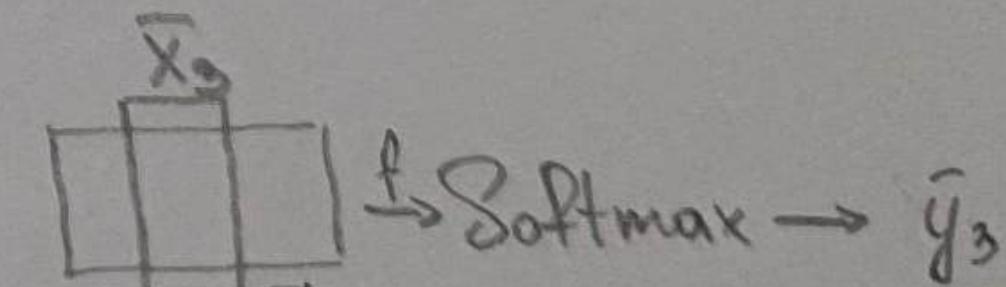
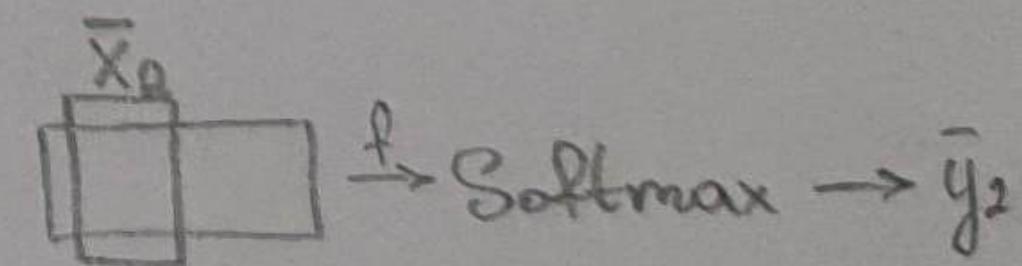
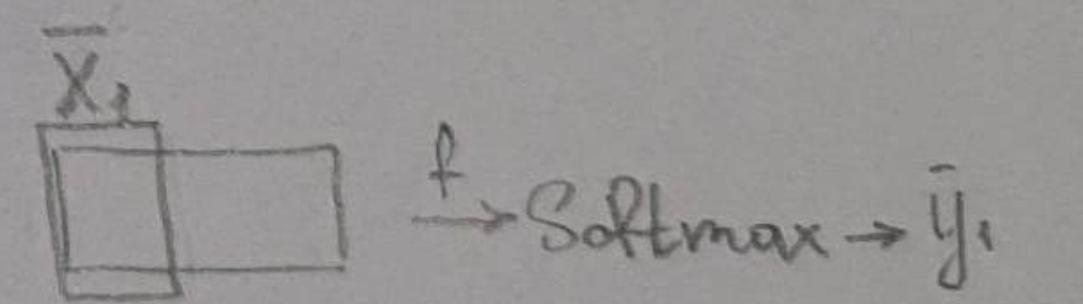
Labelling unsegmented sequence data is a ubiquitous problem in real-world sequence learning. It is particularly common in perceptual tasks (e.g. handwriting recognition, speech recognition, gesture recognition) where noisy, real-valued input streams are annotated

bellng. While these approaches have proved successful for many problems, they have several drawbacks: (1) they usually require a significant amount of task specific knowledge, e.g. to design the state models for HMMs, or choose the input features for CRFs; (2) they require explicit (and often questionable) dependency assumptions to make inference tractable, e.g. the assumption that observations are independent for HMMs; (3) for standard HMMs, training is generative, even though sequence labelling is discriminative.

Recurrent neural networks (RNNs), on the other hand, require no prior knowledge of the data, beyond the choice of input and output representation. They can be trained discriminatively, and their internal state provides a powerful, general mechanism for modelling time series. In addition, they tend to be robust to temporal and spatial noise.

So far, however, it has not been possible to apply RNNs directly to sequence labelling. The problem is that the standard neural network objective functions are defined separately for each point in the training sequence; in other words, RNNs can only be trained to make a series of independent label classifications. This means that the training data must be pre-segmented, and that the network outputs must be post-processed to give the final label sequence.

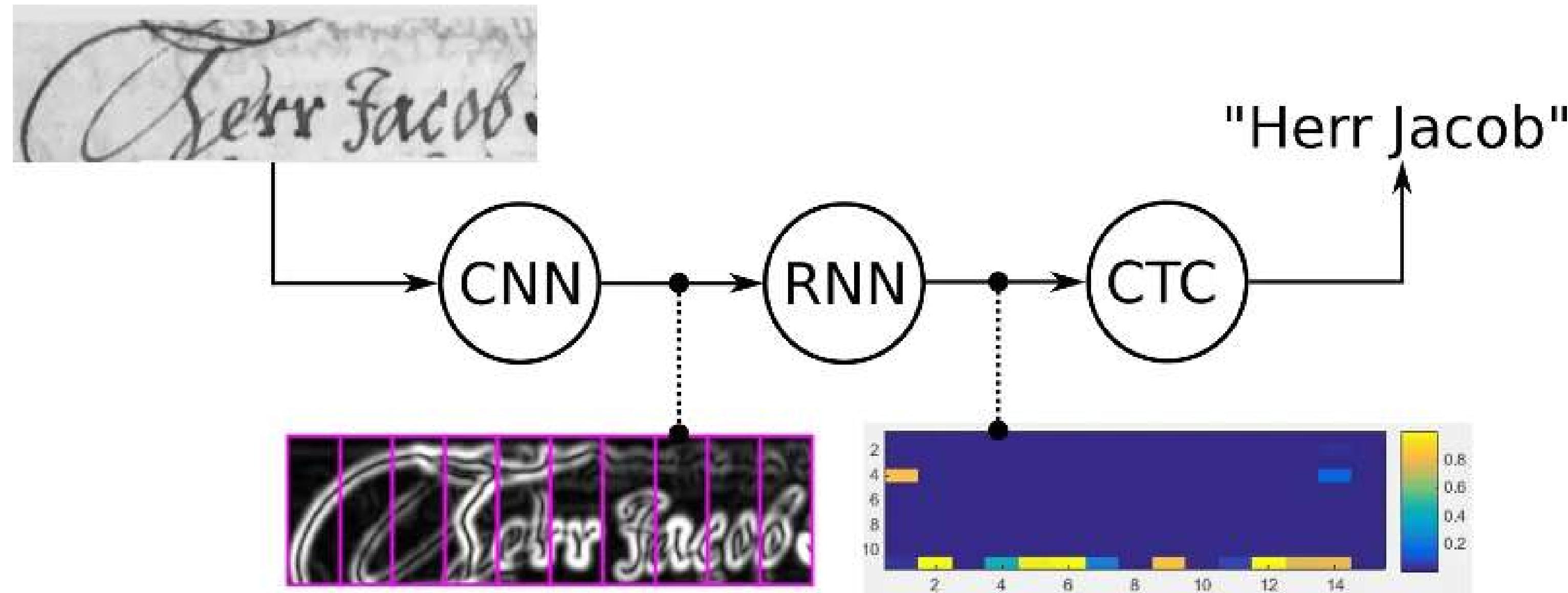
At present, the most effective use of RNNs for se-



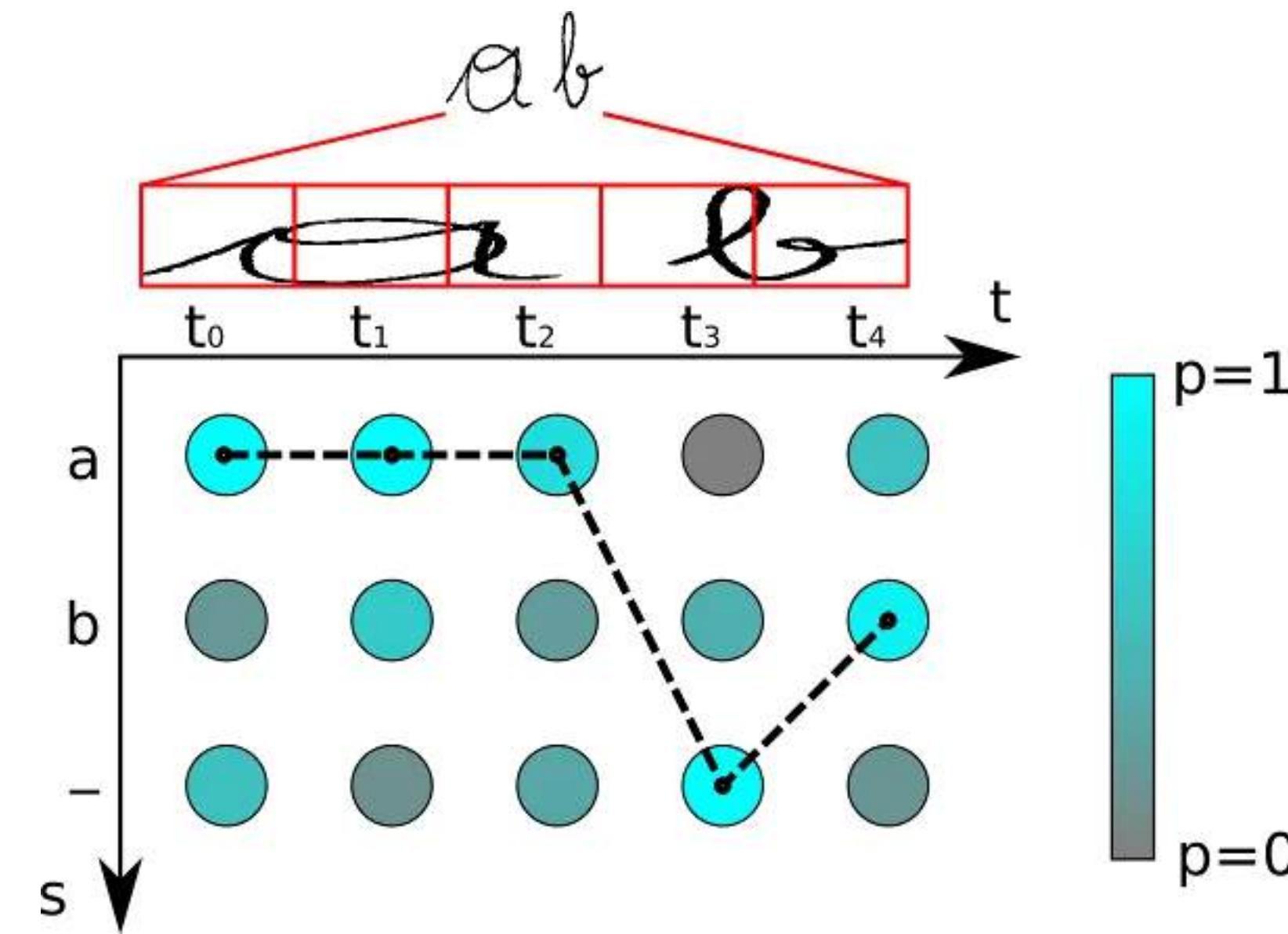
T - КОЛ-ВО ШАГОВ

-WW--000-R-DDDD-

Z = WORD



<https://towardsdatascience.com/intuitively-understanding-connectionist-temporal-classification-3797e43a86c>



$$p(\mathbf{l}|\mathbf{x}) = \sum_{\pi \in \mathcal{B}^{-1}(\mathbf{l})} p(\pi|\mathbf{x}).$$

$$h(\mathbf{x}) = \arg \max_{\mathbf{l} \in L^{\leq T}} p(\mathbf{l}|\mathbf{x}).$$

Image from: <https://towardsdatascience.com/intuitively-understanding-connectionist-temporal-classification-3797e43a86c>

Loss function

$$P(S, \mathcal{M}_\omega) = \prod_{(x,z) \in S} p(z|x) \rightarrow \max$$

$$\mathcal{O}^{\text{ML}}(S, \mathcal{M}_\omega) = - \sum_{(x,z) \in S} \ln(p(z|x))$$

$$-\ln(P(S, \mathcal{M}_\omega)) = -\sum_{(x,z) \in S} \ln(p(z|x)) \rightarrow \min$$

\mathcal{O} - objective function

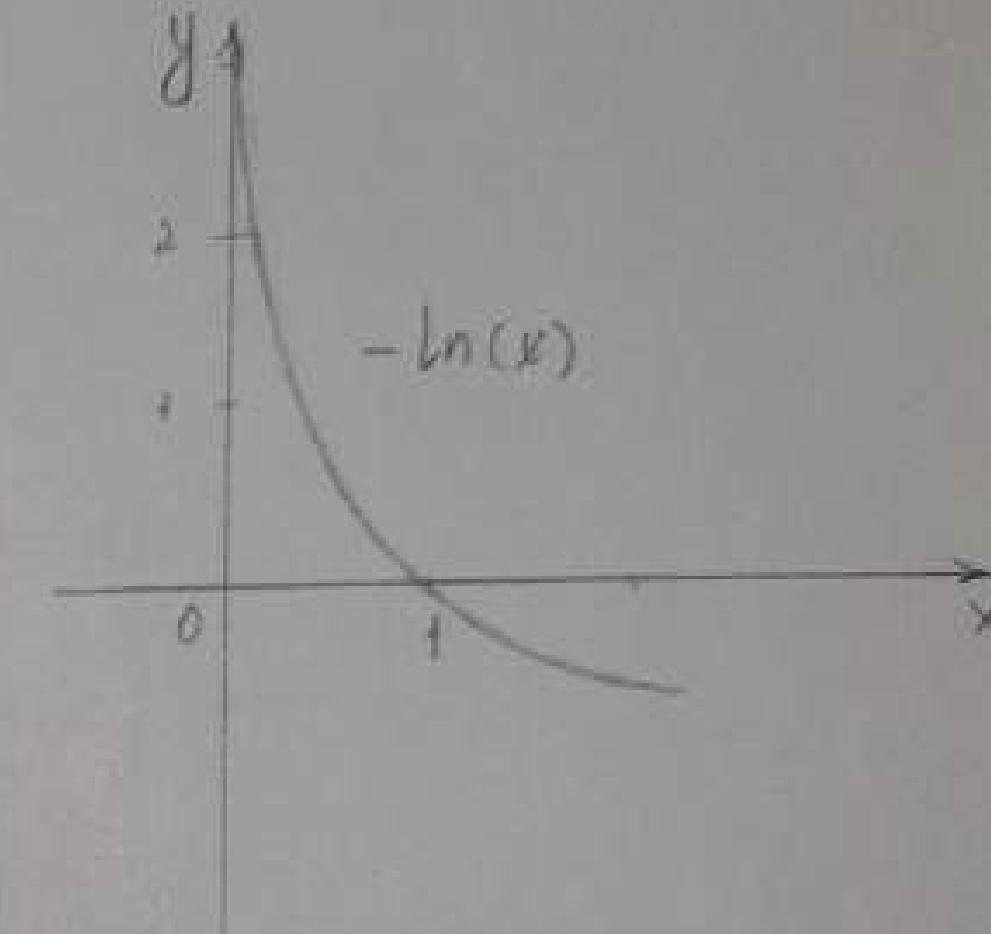
ML - maximum likelihood

S - dataset

\mathcal{M}_ω - model

x - input sequence.

z - target sequence.



$$\frac{d\mathcal{O}^{\text{ML}}(\{(x,z)\}, \mathcal{M}_\omega)}{dy_k^t} = -\frac{d\ln(p(z|x))}{dy_k^t} \quad \text{partial derivative by some } y_k^t$$

Forward Backward Algorithm

$$\alpha_t(s) \stackrel{\text{def}}{=} \sum_{\substack{\pi \in N^T : \\ \mathcal{B}(\pi_{1:t}) = \mathbf{l}_{1:s}}} \prod_{t'=1}^t y_{\pi_{t'}}^{t'}.$$

$$\beta_t(s) \stackrel{\text{def}}{=} \sum_{\substack{\pi \in N^T : \\ \mathcal{B}(\pi_{t:T}) = \mathbf{l}_{s:|1|}}} \prod_{t'=t}^T y_{\pi_{t'}}^{t'}$$

$$\alpha_1(1) = y_b^1$$

$$\beta_T(|\mathbf{l}'|) = y_b^T$$

$$\alpha_1(2) = y_{\mathbf{l}_1}^1$$

$$\beta_T(|\mathbf{l}'|-1) = y_{\mathbf{l}_{|\mathbf{l}|}}^T$$

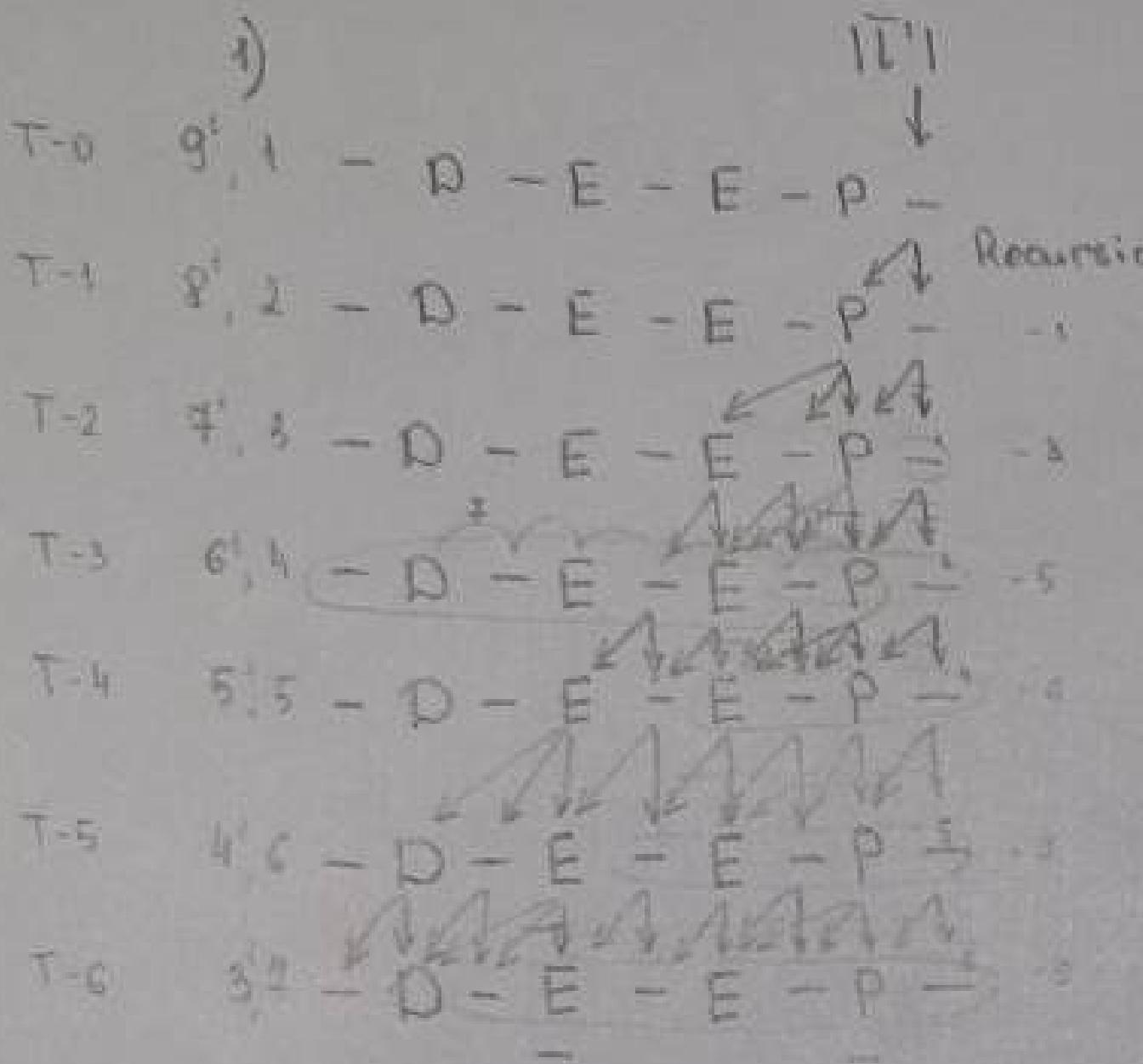
$$\alpha_1(s) = 0, \quad \forall s > 2$$

$$\beta_T(s) = 0, \quad \forall s < |\mathbf{l}'|-1$$

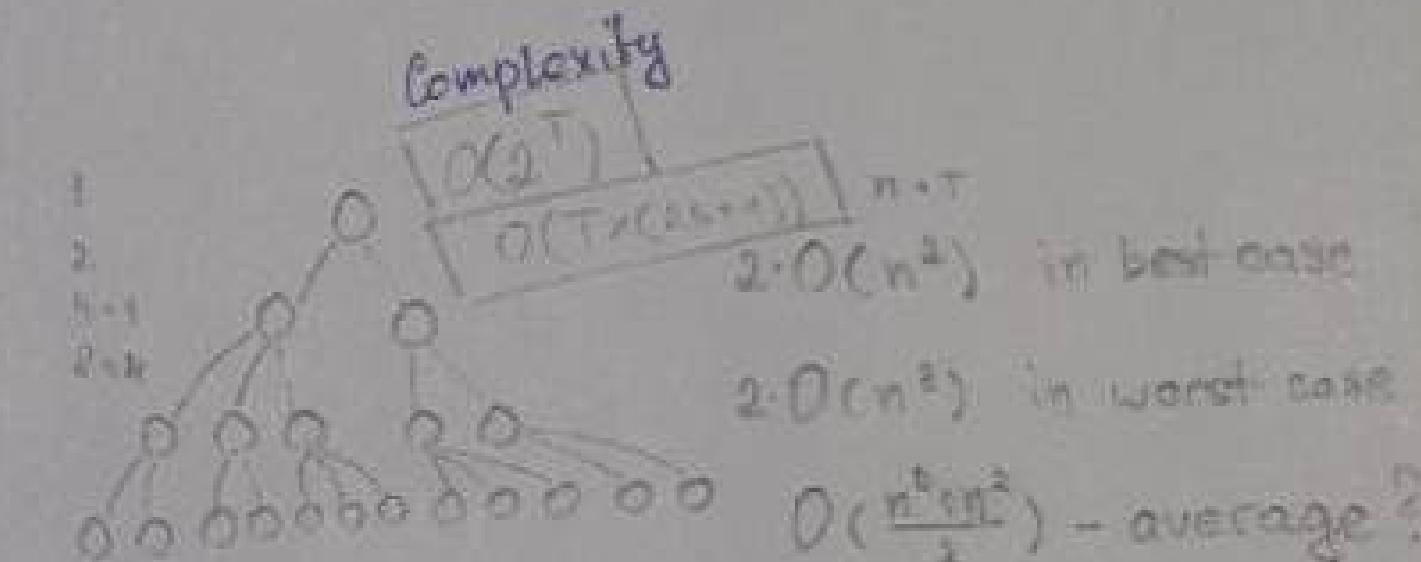
Forward - Backward Algorithm

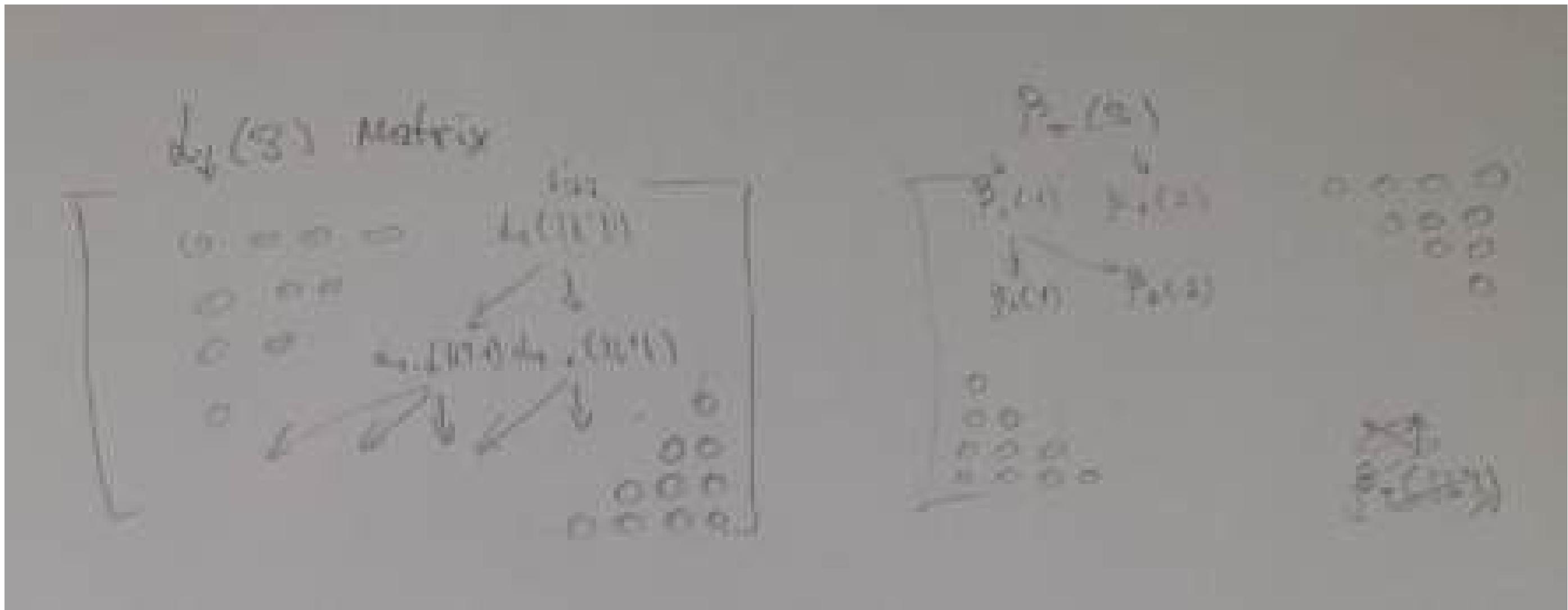
$$P(\bar{L} \mid x) = d_T(|\bar{L}'|) + d_T^{(2)}(|\bar{L}'|-1)$$

$$d_t(s) = \begin{cases} y_{l_s}^t \cdot (d_{t-1}(s) + d_{t-1}(s-1)) & \text{if } l_s = b \text{ or } l_s = l_{s-1} \\ y_{l_s}^t \cdot (d_{t-1}(s) + d_{t-1}(s-1) + d_{t-1}(s-2)) & \text{otherwise} \end{cases}$$



$$\begin{aligned} d_T(|\bar{L}'|) &= y_b^T \cdot (d_{T-1}(|\bar{L}'|) + d_{T-1}(|\bar{L}'|-1)) \\ 1) d_{T-1}(|\bar{L}'|) &= y_{l_{T-1}}^{T-1} (d_{T-2}(|\bar{L}'|) + d_{T-2}(|\bar{L}'|-1)) \\ 2) d_{T-1}(|\bar{L}'|-1) &= y_{l_{T-1}}^{T-1} (d_{T-2}(|\bar{L}'|-1) + d_{T-2}(|\bar{L}'|-2)) + \\ &\quad + d_{T-2}(|\bar{L}'|-3)) \end{aligned}$$

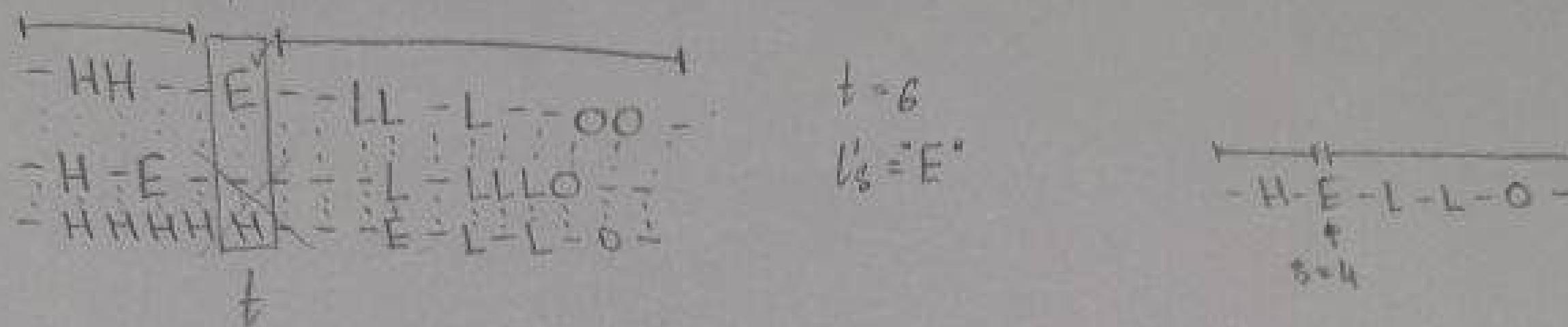




Forward - Backward Algorithm

$$\alpha_t(s) \cdot \beta_t(s) = \sum_{\pi \in B^{-1}(L)} y_{i_s}^t \prod_{t'=1}^T y_{\pi_{t'}}^t$$

$\pi_t = L_s$



$\alpha_t(s)$ - беремуси жи t үрагчын $g(s)$, ма сони $l'_{1:s}$

$\beta_t(s)$ - беремуси жи t үрагчын $g(s)$ га түрүүлж, ма сони $l'_{s:T}$

$$d_t(s) \cdot \beta_t(s) = \sum_{\pi \in B^{-1}(l)} \underbrace{y_{l_s}^t}_{\pi_t = l_s} \cdot \underbrace{\prod_{t=1}^T y_{\pi_t}^t}_{p(\pi|x)}, \quad \begin{array}{l} \text{const} \\ \downarrow \\ \pi \in B^{-1}(l) \end{array}, \quad \begin{array}{l} t = \text{const} \\ s = \text{const} \end{array}, \quad \begin{array}{l} d_1(3) = 0 \\ d_2(5) = 0 \\ \dots \end{array}$$

$$p(\pi|x) = \prod_{t=1}^T y_{\pi_t}^t, \forall \pi \in L^{*T}, L^{*T} \in L^{*}$$

$$d_t(s) \cdot \beta_t(s) \cdot \frac{1}{y_{l_s}^t} = \sum_{\pi \in B^{-1}(l)} p(\pi|x)$$

$\pi_t = l_s$

$$1) P(l|x) = \sum_{\pi \in B^{-1}(l)} p(\pi|x) = \text{no margin + here u have Sigma no braces t}$$

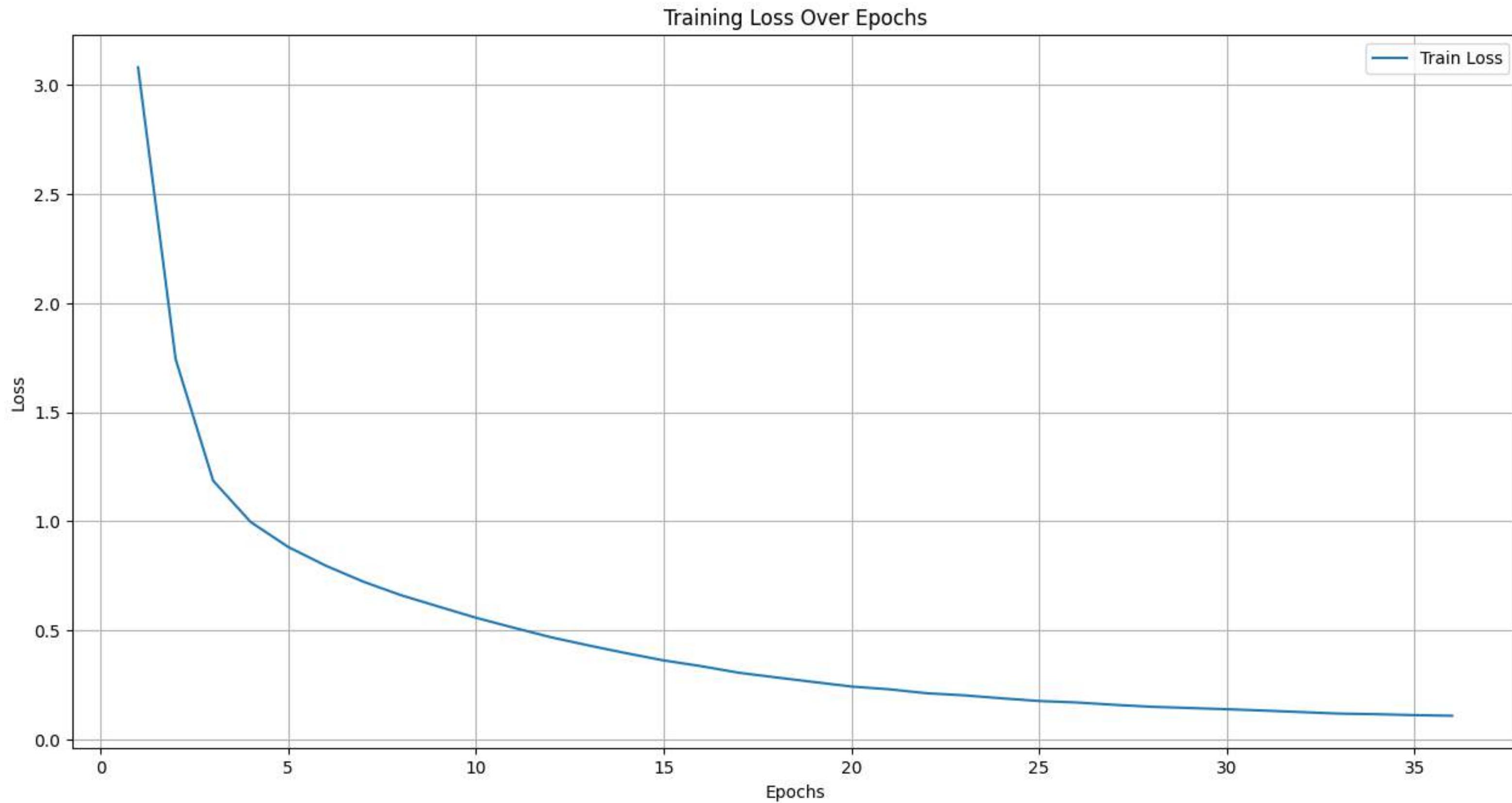
$$2) P(l|x)_t = \sum_{s=1}^{T+1} \frac{1}{y_{l_s}^t} \cdot d_t(s) \cdot \beta_t(s) \quad \begin{array}{l} \text{nonhomogeneous } p(l|x) \\ \text{no margin } t = \text{const}, \text{ max remo } 1) \neq 2) \end{array}$$

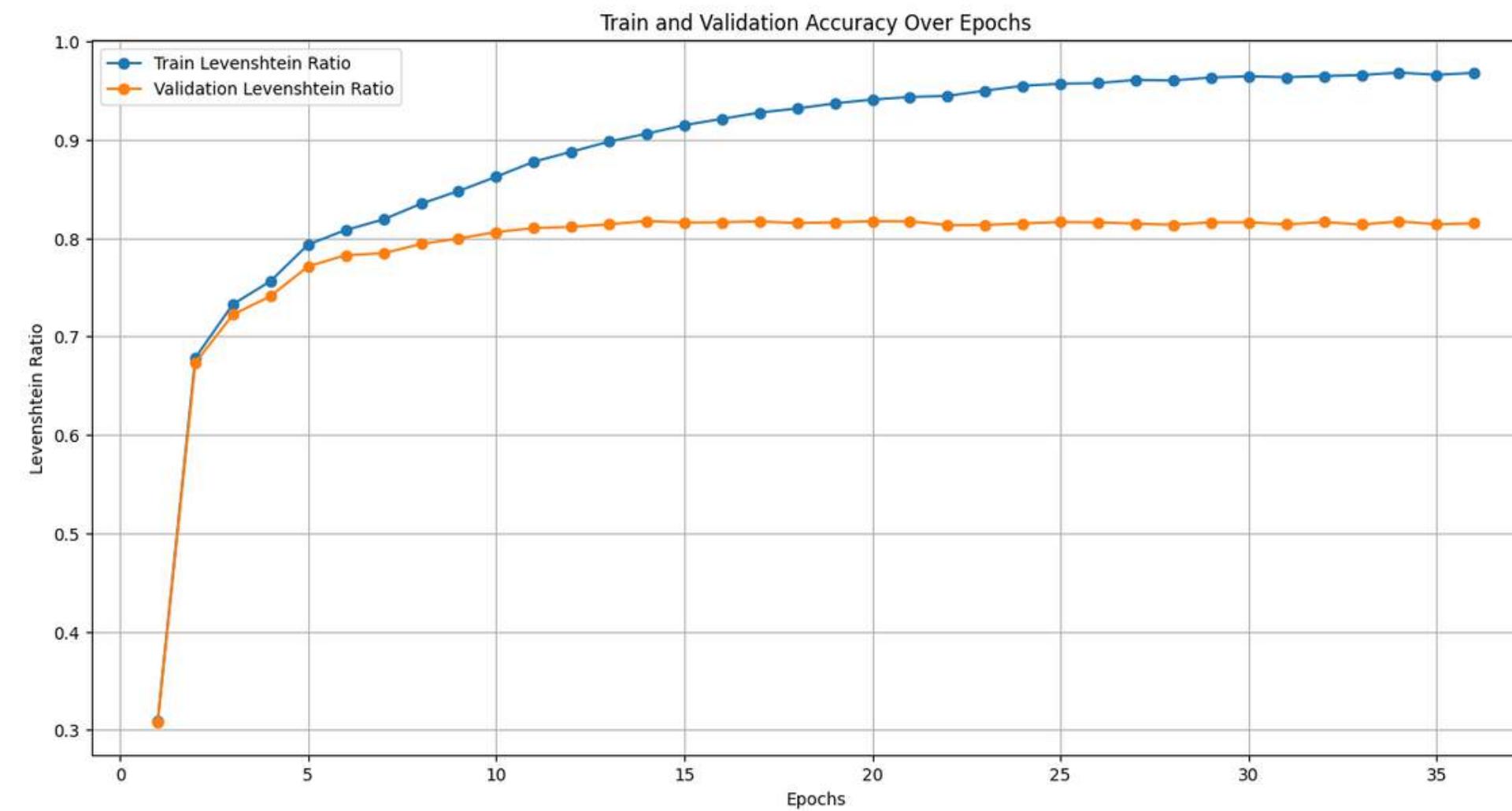
$$\sum_{s=1}^{T+1} \sum_{\pi \in B^{-1}(l)} p(\pi|x), t$$

$\pi_t = l_s$

$$E = O^{ML}(S, \Phi_\omega) = - \sum_{\substack{(x, z) \in S \\ l \in Z}} \ln(p(l|x))$$

и $\pi_t = l_s$ т.е. равнозначно или не симметрично





АЛМАЙДЫ

аңшоиды

ҚҰРАЛДАРЫНА

құралдароста

ӨНДЕУ

өңдеу

ҒАНЫҢ

ғаның

ТУРАҚТ

мұнисијат

АҚЫ

ақы

ҚОЛДАНЫС

қолданыс

КАРТОЧКАЛАРДЫ

карточкаларды

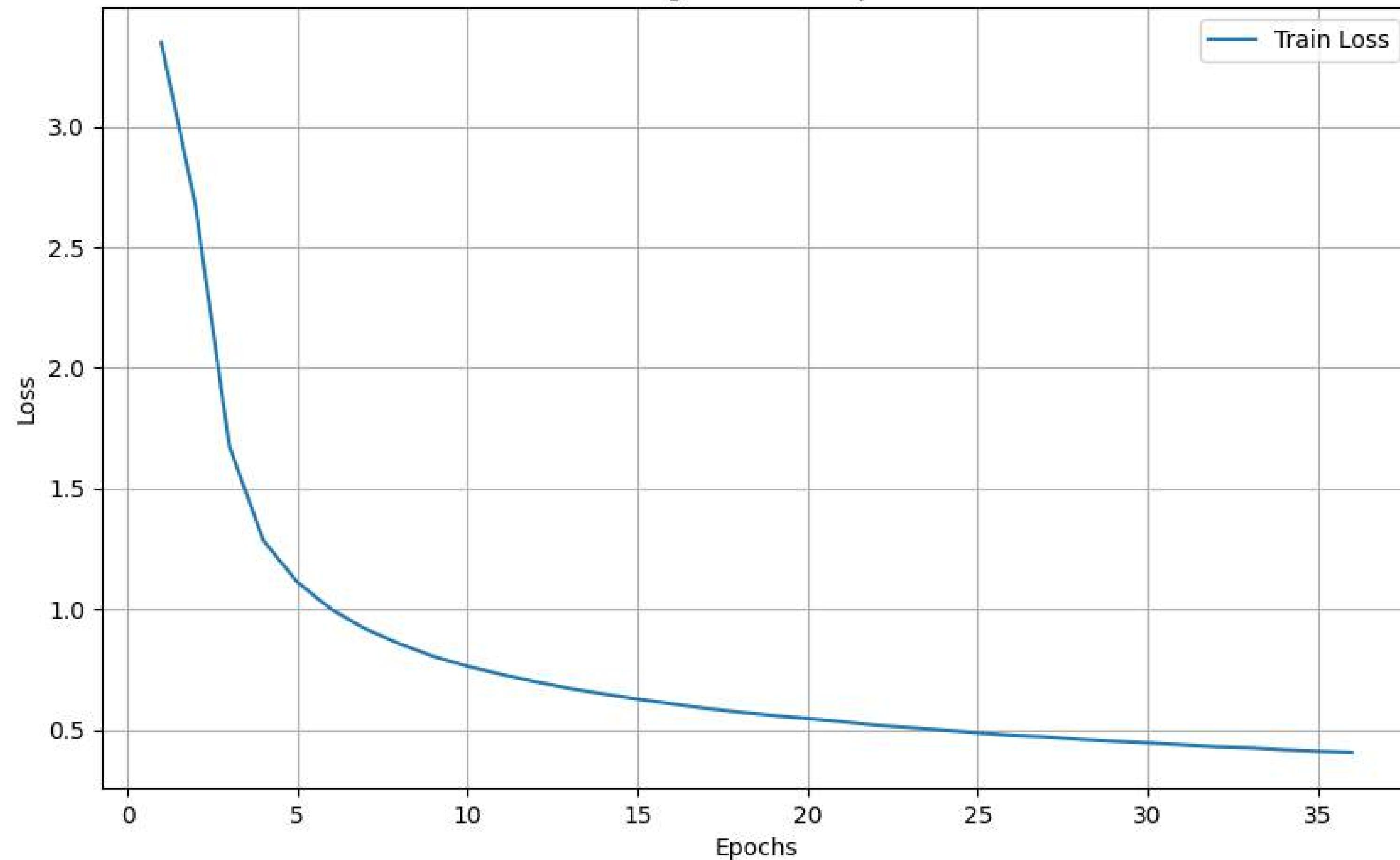
ЖАТКАН

жаткан

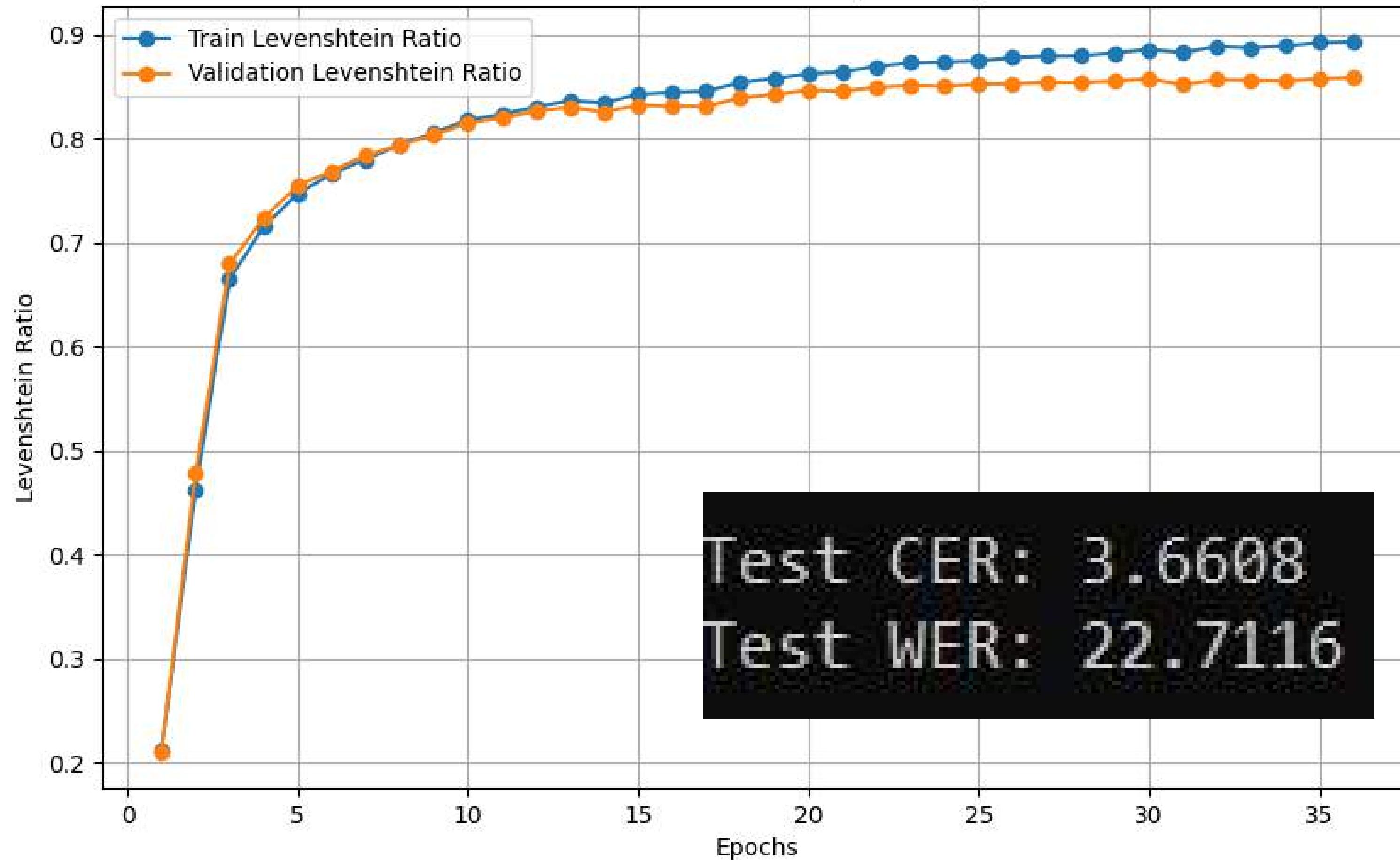
ӨЗ

өз

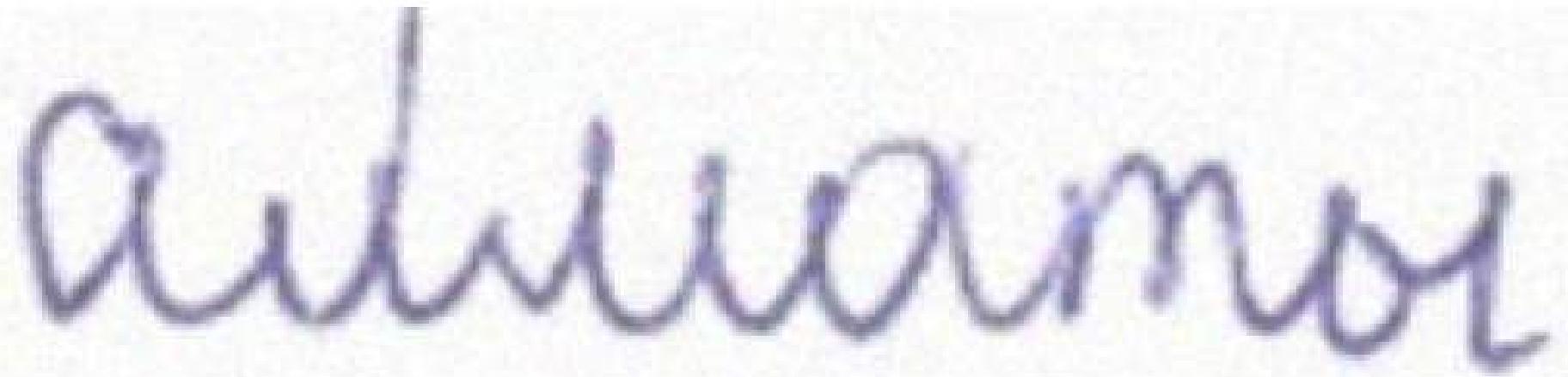
Training Loss Over Epochs



Train and Validation Accuracy Over Epochs

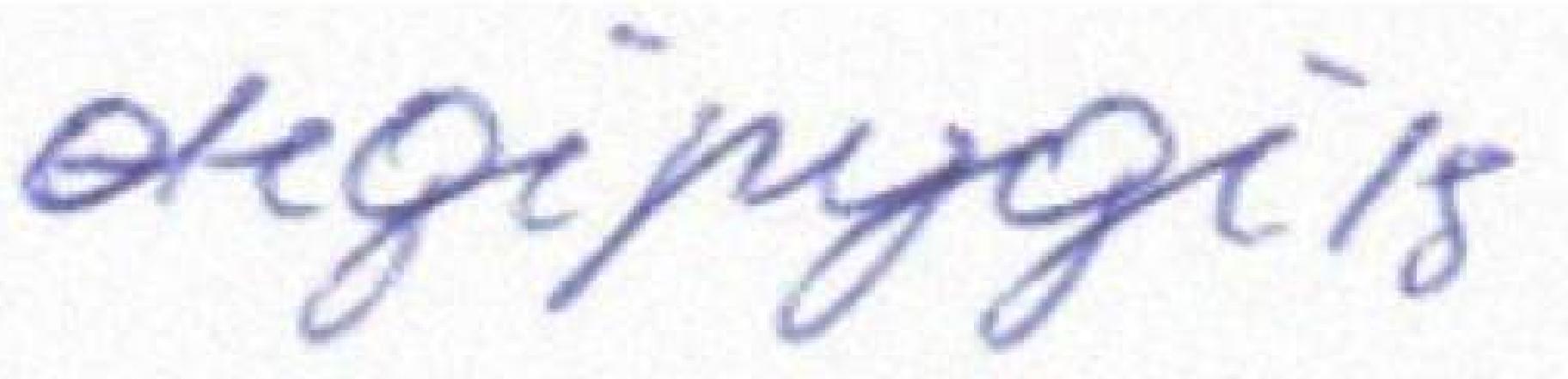


GT: АЛМАТЫ
Predicted: АЛМАТЫ
Avg Conf: 0.9713



Алматы

GT: ӨНДІРУДІҢ
Predicted: ӨНДІРУДІҢ
Avg Conf: 0.9969



өндірудің

GT: НЕМЕСЕ
Predicted: НЕМСЕ
Avg Conf: 0.9251



Немесе

GT: ОСЫ
Predicted: ОСЫ
Avg Conf: 0.9167



осы

GT: АДАМДАР
Predicted: АДАМДАР
Avg Conf: 0.9403



адамдар

GT: ЖҮЙ
Predicted: ЖҮБ
Avg Conf: 0.9958



жүб

GT: ОРТА
Predicted: ДРТА
Avg Conf: 0.8590

Жамса

GT: ЕРЕКШЕЛІГІ
Predicted: ЕРЕКШШЕІГІ
Avg Conf: 0.9248

ерекшешілігі

GT: БОЛУ
Predicted: ыАЛУ
Avg Conf: 0.7767

боялы

Recognized Text Overlaid

КӘСІПРЫН

Кәсіп отолы

ЖФЕ

ЖФЕ

ТОЛЫҚ

толық

АНБЕХКАРЫНЫ

Анбекжаконыс,

3

ЕССПОУ

ЕССПОУ

4

FUTURE RESEARCH

- Circular text recognition
- Table detection and structure recognition
- Sequential patterns of texts for recognition or post-processing
- End-to-end model for text detection and recognition
- Speech Recognition

REFERENCES

- Lee, A. W., Chung, J., & Lee, M. (2021). GNHK: a dataset for English handwriting in the wild. In Document Analysis and Recognition–ICDAR 2021: 16th International Conference, Lausanne, Switzerland, September 5–10, 2021, Proceedings, Part IV 16 (pp. 399-412). Springer International Publishing.
- C. Yao, X. Bai, W. Liu, Y. Ma and Z. Tu. Detecting Texts of Arbitrary Orientations in Natural Images. CVPR 2012 (PDF).
- Singh, A., Pang, G., Toh, M., Huang, J., Galuba, W., & Hassner, T. (2021). Textocr: Towards large-scale end-to-end reasoning for arbitrary-shaped scene text. In Proceedings of the IEEE/CVF conference on computer vision and pattern recognition (pp. 8802-8812).
- Matcha, A. C. N. (n.d.). Handwriting Recognition with ML (An In-Depth Guide). Nanonets Blog | AI Document Processing & Workflow Automation. <https://nanonets.com/blog/handwritten-character-recognition/>
- Toiganbayeva, N., Kasem, M., Abdimanap, G., Bostanbekov, K., Abdallah, A., Alimova, A., & Nurseitov, D. (2022). Kohtd: Kazakh offline handwritten text dataset. Signal Processing: Image Communication, 108, 116827.
- Kumar, A., & Pati, P. B. (2023). Offline HWR Accuracy Enhancement with Image Enhancement and Deep Learning Techniques. Procedia Computer Science, 218, 35-44.
- Nurseitov, D., Bostanbekov, K., Kurmankhojayev, D., Alimova, A., Abdallah, A., & Tolegenov, R. (2021). Handwritten Kazakh and Russian (HKR) database for text recognition. Multimedia Tools and Applications, 80(21), 33075-33097.
- Toiganbayeva, N., Kasem, M., Abdimanap, G., Bostanbekov, K., Abdallah, A., Alimova, A., & Nurseitov, D. (2022). Kohtd: Kazakh offline handwritten text dataset. Signal Processing: Image Communication, 108, 116827.
- Patterson, J., & Gibson, A. (2017). Deep learning: A practitioner's approach (First edition). O'Reilly.
- Goodfellow, I., Bengio, Y., & Courville, A. (2016). Deep learning. The MIT press.
- Raschka, S., Liu, Y. H., & Mirjalili, V. (2022). Machine Learning with PyTorch and Scikit-Learn: Develop machine learning and deep learning models with Python. Packt Publishing Ltd.

**ANY
QUESTIONS?**

THANK YOU!