

# Распознавание букв и элементов письма в русской каллиграфии

Постановка задачи и  
подход к решению (work in progress)

**Местецкий Л.М.**

МГУ имени М.В. Ломоносова  
Факультет вычислительной математики и кибернетики

# Challenge

- Все, кто учился в школе, научились читать печатные тексты и научились писать рукописные тексты – есть отработанные **методики обучения**
- Дети не могут прочесть рукописные тексты, а взрослые могут – **как они научились?**
- Считается, что тот, кто умеет писать, тот может и читать. Методики обучения чтению рукописного текста на русском языке **не существует!**
- Задача автоматического распознавания рукописного текста – **обучение робота навыкам, которым людей обучать не умеем!**

# Каллиграфия

- Каллиграфия — одна из отраслей изобразительного искусства, это искусство красивого письма.
- Мы употребляем это слово как набор навыков письма всех русскоязычных грамотных людей, полученных в процессе обучения в школе на уроках чистописания, с использованием эталонных прописей. Эти навыки мало менялись за последние сто лет, как мало менялись сами прописи.
- Мы предполагаем, что все пишущие на русском языке делают это по тому стандарту, который представлен в прописях и на основе тех навыков, которые получили в школе, а отличия образуются вследствие индивидуальных свойств людей.

# Русская каллиграфия

Аа Бб Вв Гг  
Ддд Ее Жж Зз  
Ии Кк Лл Мм  
Нн Оо Пп Рр Сс  
Ттт Уу Фф Хх  
Цц Чч Шш Щщ  
Ээ Юю Яя ьы ъ  
: ; „ “ ? ! § №  
1 2 3 4 5 6 7 8 9 0

А. Б. В. Г. Д. Е. Ж. З. И  
К. Л. М. Н. О. П. Р. С.  
Т. У. Ф. Х. Ц. Ч. Ш. Щ. Ъ  
Э. Ю. Я. ъ  
абвгдежзигклмнопрстуфхцг  
шщрытьяюяей 1. 2. 3. 4. 5. 6. 7. 8. 9. 0.

# Каллиграфическое письмо

На вершине горы лежит клад.  
Рядом живёт рыжий дракон. Он  
сторожит клад. У дракона большие  
лапы и пушистые уши и хвост.  
Принц решил одолеть дракона. При-  
шёл с оружием, а дракон убежал.  
Писал Саттуранов Тимур 1 кл.

Уди вперёд!

3.

Уди вперёд! Борись со судьбой!  
Мужайся духом и не падай!  
Пусть надежда путь живой  
Во мраке жизни предь тобой  
Сияет яркою лампадой!  
Уди вперёд! Не унывай!  
Но палкой гордаго терпенья,  
Малое трудися и страдай,  
Чтоб людей, как клад, скривай  
Свои незрелая мученя!  
Уди вперёд! Сильн будь со мной!

# Обычное письмо

Обещание

Часть 3. Очередь

В очереди люди были все серые, все глядели  
внутрь себя, все шепстели еле слышно в  
ожидании приговора. Кому год, кому три...

В очереди люди были все серые все глядели  
внутрь себя все шепстели еле слышно в ожида-  
нии приговора. Кому год кому три... Наверняка

В очереди люди были все серые, все гляде-  
ли внутрь себя, все шепстели еле слышно в  
ожидании приговора. Кому-год, +  
кому-три, ... Наверняка все, как и Там, про-

В очереди люди были все серые, все глядели  
внутрь себя, все шепстели еле слышно в  
ожидании приговора. Кому год, кому три...  
Наверняка, все, как и Там, просматривали

# Задача

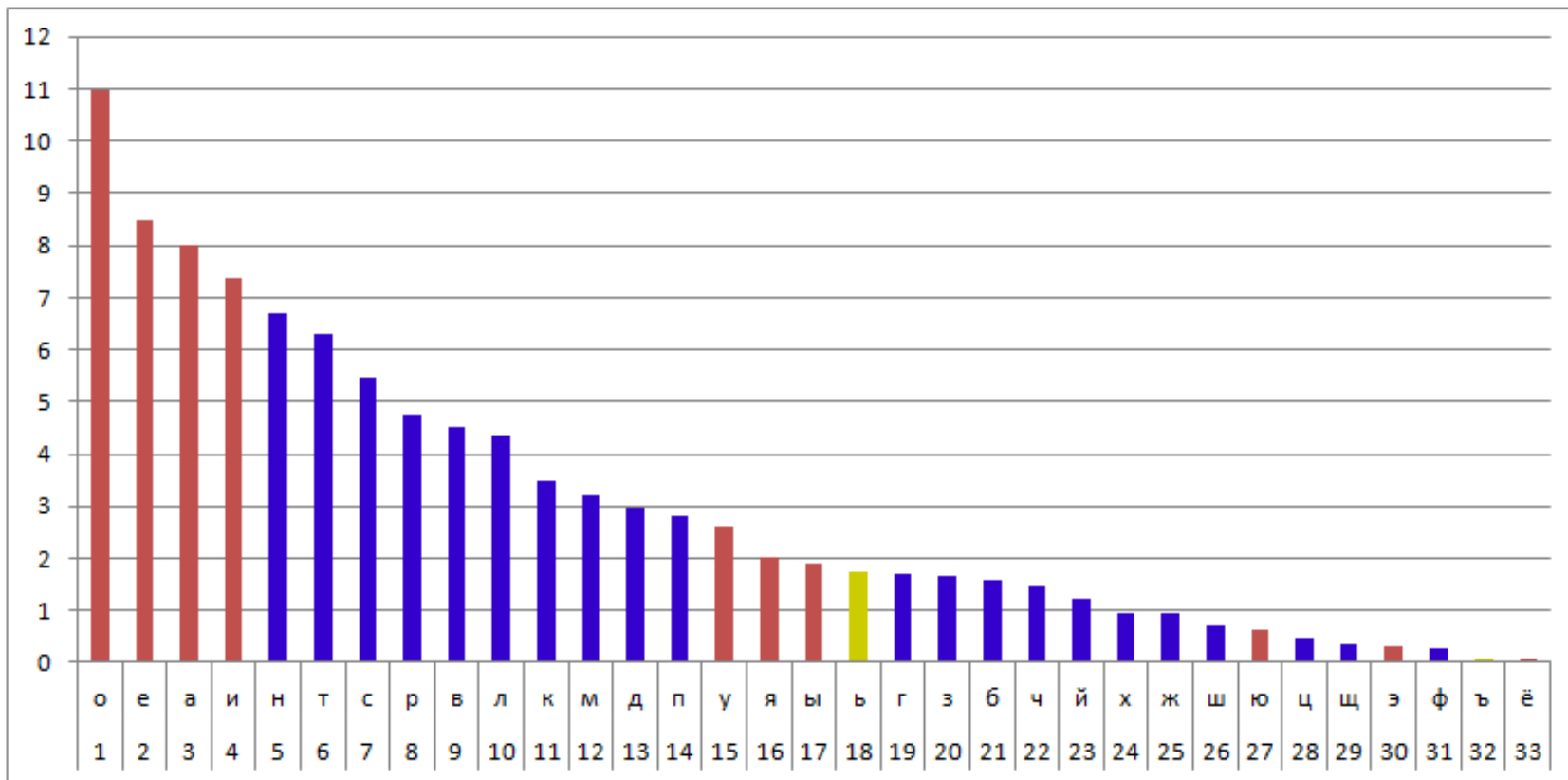
- Автоматизация оцифровки рукописных текстовых документов
- Поиск ключевых слов в рукописных документах

# Подход

- Штриховое медиальное представление цифровых изображений текстовых документов
- Распознавание базовых рукописных штрихов – овалов, крючков, палочек и их сочетаний
- Машинное обучение на основе штрихового представления

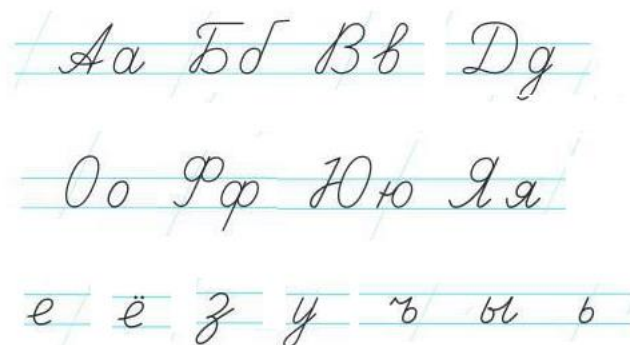
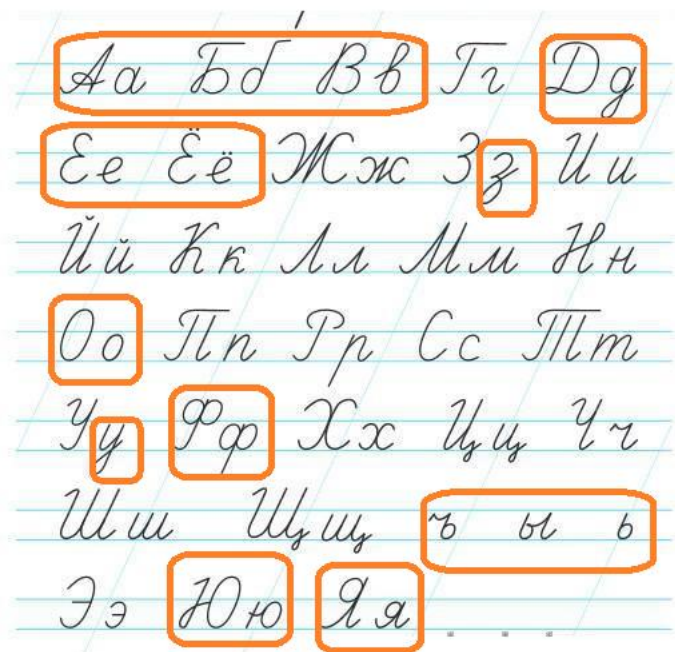


# Частотное распределение букв



место	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	33
буква	о	е	а	и	н	т	с	р	в	л	к	м	д	п	у	я	ы	ь	г	з	б	ч	й	х	ж	ш	ю	ц	щ	э	ф	ъ	ё
%	10,98	8,48	8,00	7,37	6,70	6,32	5,47	4,75	4,53	4,34	3,49	3,20	2,98	2,80	2,62	2,00	1,90	1,74	1,69	1,64	1,59	1,45	1,21	0,97	0,94	0,72	0,64	0,49	0,36	0,33	0,27	0,04	0,01

# Буквы с овальными элементами



Овалы – 52.17%

# Буквы с другими элементами «палочка», крючок, полуовал

Кк

Нн

Ии

Йй

Пп Рр

Тт

Шш Щщ

Палочка – 19.31%

Крючок – 9.66 %

Жж Зз

Сс

Хх

Ээ

Полуовал – 7.71%

# Распределение букв по каллиграфическим элементам

место	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	33
буква	о	е	а	и	н	т	с	р	в	л	к	м	д	п	у	я	ы	ь	г	з	б	ч	й	х	ж	ш	ю	ц	щ	э	ф	ъ	ё
%	10,98	8,48	8,00	7,37	6,70	6,32	5,47	4,75	4,53	4,34	3,49	3,20	2,98	2,80	2,62	2,00	1,90	1,74	1,69	1,64	1,59	1,45	1,21	0,97	0,94	0,72	0,64	0,49	0,36	0,33	0,27	0,04	0,01

о	е	а	и	н	т	с	р	в	л	к	м	д	п	у	я	ы	ь	г	з	б	ч	й	х	ж	ш	ю	ц	щ	э	ф	ъ	ё			
1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	33			
10,98	8,48	3,00		6,70	6,32		4,75	4,53		3,49		2,98		2,62	2,00	1,90	1,74		1,64	1,59					0,64				0,27	0,04	0,01				
				7,37																															

1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	33			
о	е	а	и	н	т	с	р	в	л	к	м	д	п	у	я	ы	ь	г	з	б	ч	й	х	ж	ш	ю	ц	щ	э	ф	ъ	ё			

Частота букв с овалами 52.17%

Частота букв с палочками 19.31% (всего 71.48%)

Частота букв с крючками 9.66% (всего 81.14%)

Частота букв с полуовалами 7.71% (всего 88.85%)

# Восстановление текста по буквам

Буря \_\_\_ою ебо рое\_,  
В р\_\_\_е\_ые ру\_я;  
\_о, \_а\_ зверь, о\_а завое\_,  
\_о за\_а\_е\_, а\_ д\_я,  
\_о\_о\_ров\_е обве\_а\_ой  
Вдру\_о\_о\_ой за\_у\_\_\_,  
\_о, \_а\_ у\_\_\_ за\_озда\_ый,  
\_\_а\_ в о\_о\_о за\_у\_\_\_.  
\_а\_а\_ве\_ая\_а\_у\_а  
\_\_е\_а\_ь\_а\_\_е\_а.  
\_\_о\_е\_ы, \_оя\_ару\_а,  
\_р\_у\_о\_\_\_а\_у\_о\_а?  
\_\_\_ бур\_завыва\_ье\_  
\_ы, \_ой дру\_, у\_о\_е\_а,  
\_\_\_ дре\_е\_ь\_од\_у\_а\_ье\_  
\_вое\_о\_вере\_е\_а?  
Вы\_ье\_, добрая\_одру\_а  
Бед\_ой\_ю\_о\_\_\_о\_ей,  
Вы\_ье\_ \_\_\_о\_ря; де\_е\_ру\_а?  
\_ерд\_у\_буде\_ве\_е\_ей.  
\_\_ой\_\_е\_е\_ю, \_а\_ \_\_\_а  
\_\_о\_за\_оре\_ \_\_\_а;  
\_\_ой\_\_е\_е\_ю, \_а\_ дев\_а  
За\_водой\_оу\_ру\_а.

Буквы с овалами

Буря \_\_\_ою небо кроет,  
В р\_\_\_не\_ные крутя;  
То, как зверь, она завоюет,  
То зап\_а\_ет, как д\_тя,  
То по кров\_е обвет\_а\_ой  
Вдру\_о\_о\_ой за\_у\_т,  
То, как путн\_к запозда\_ый,  
К на\_в\_око\_ко за\_ту\_т.  
На\_а\_вет\_ая\_а\_у\_ка  
\_\_пе\_а\_ь\_на\_\_те\_на.  
\_\_то\_е\_ты, \_оя\_тару\_ка,  
Пр\_у\_о\_к\_а\_у\_окна?  
\_\_\_ бур\_завыванье\_  
Ты, \_ой дру\_, у\_то\_ена,  
\_\_\_ дре\_е\_ь\_под\_у\_анье\_  
\_вое\_о\_веретена?  
Выпье\_, добрая подру\_ка  
Бедной юно\_т\_о\_ей,  
Выпье\_ \_\_\_о\_ря; де\_е\_кру\_ка?  
\_ердцу\_будет\_ве\_е\_ей.  
\_\_пой\_не\_пе\_ню, как \_\_н\_ца  
Т\_о\_за\_оре\_ \_\_\_а;  
\_\_пой\_не\_пе\_ню, как дев\_ца  
За\_водой\_поутру\_а.

Буквы с овалами  
и палочками

Буря \_\_\_ою небо кроет,  
Вихри снежные крутя;  
То, как зверь, она завоюет,  
То зап\_а\_ет, как д\_тя,  
То по кров\_е обветша\_ой  
Вдру\_со\_о\_ой зашу\_ит,  
То, как путник запозда\_ый,  
К на\_в\_окошко засту\_ит.  
Наша ветхая\_а\_ужка  
И пе\_а\_ь\_на\_и\_те\_на.  
\_\_то\_же\_ты, \_оя\_старушка,  
Приу\_о\_к\_а\_у\_окна?  
И\_и\_бури\_завыванье\_  
Ты, \_ой дру\_, у\_то\_ена,  
И\_и\_дре\_е\_шь\_под\_жужжанье\_  
Свое\_о\_веретена?  
Выпье\_, добрая подружка  
Бедной юности\_о\_ей,  
Выпье\_ с\_о\_ря; де\_же\_кружка?  
Сердцу\_будет\_весе\_ей.  
Спой\_не\_песню, как синица  
Тихо\_за\_оре\_жи\_а;  
Спой\_не\_песню, как девица  
За\_водой\_поутру\_ш\_а.

Буквы с овалами,  
палочками и  
полуовалами

# Исходный документ

Обещание

Часть 3. Очередь

В очереди люди были все серые, все внутри себя, все шепстели еле слышим ожидания приговора. Кому год, кому 10, Улаверика все, как и Тамя, просят каска свою 'жизнь, как и Тамя, что-себе обещали: заново начать, всё изменить заново сегодня так, как откладывал, завтра, только бы доктор дал ещё в/у как будто приговор уже не был запис

Обещание

Часть 3. Очередь

В очереди люди были все серые, все внутри себя, все шепстели еле слышим ожидания приговора. Кому год, кому 10, Улаверика все, как и Тамя, просят каска свою 'жизнь, как и Тамя, что-себе обещали: заново начать, всё изменить заново сегодня так, как откладывал и завтра, только бы доктор дал ещё в/у

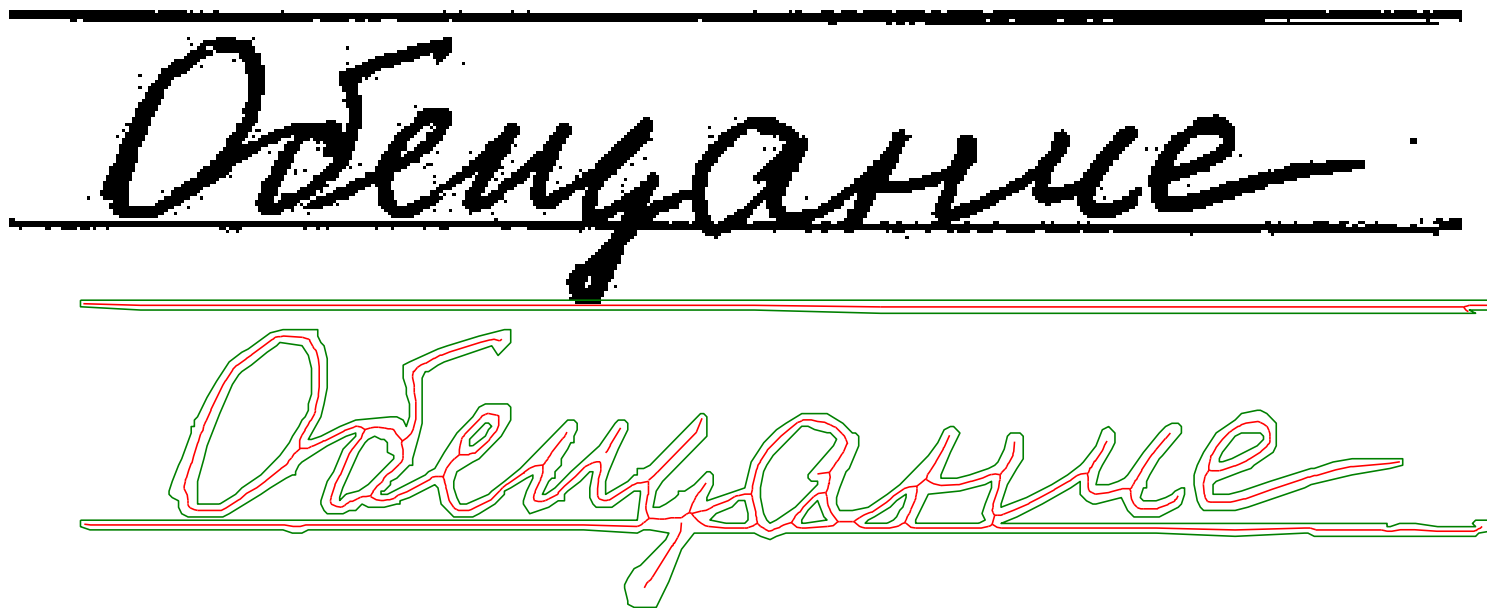
# Сегментация - бинаризация

Обещание

Часть 3. Очередь

В очереди люди били все серые, все  
внутри себя, все шелестели еле слышным  
ожиданием приговора. Кому год - кому т,  
Наверняка все как и Памя, просидит  
каспер свою жизнь, как и Памя, что-то  
себе обещали: заново начать, всё изменить  
занесть сегодня так, как складывал и  
завтра, только вы доктор да еще в

# Гранично-скелетное представление



Граница – многоугольная фигура, аппроксимирующая растровое изображение

Скелет – срединная ось состоит из точек-центров вписанных в фигуру окружностей



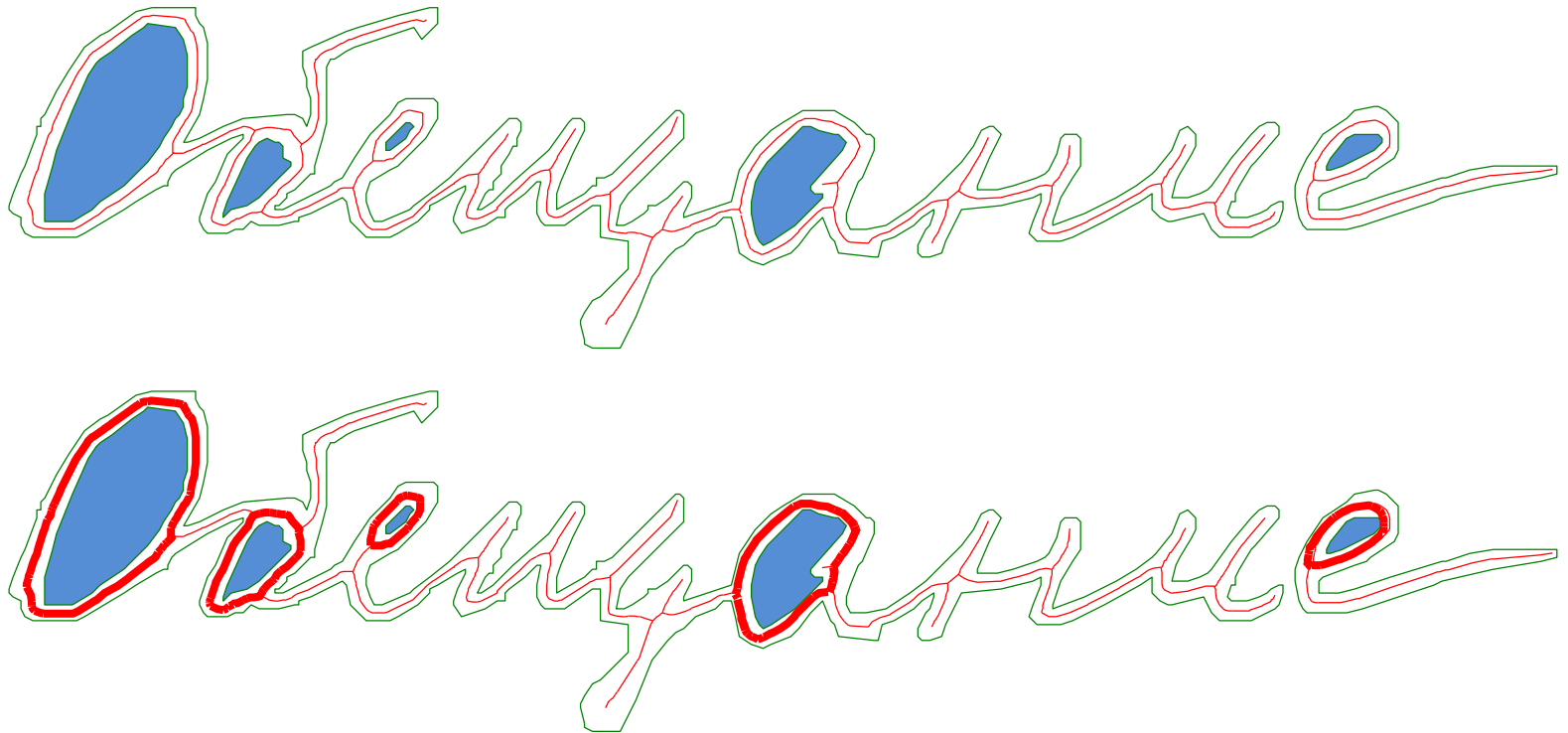
# Выделение текста



Прямолинейная разметка удаляется



# Выделение овалов



Каждый овал определяется связной компонентой разбиения изображения границей многоугольной фигуры (дыра внутри фигуры)

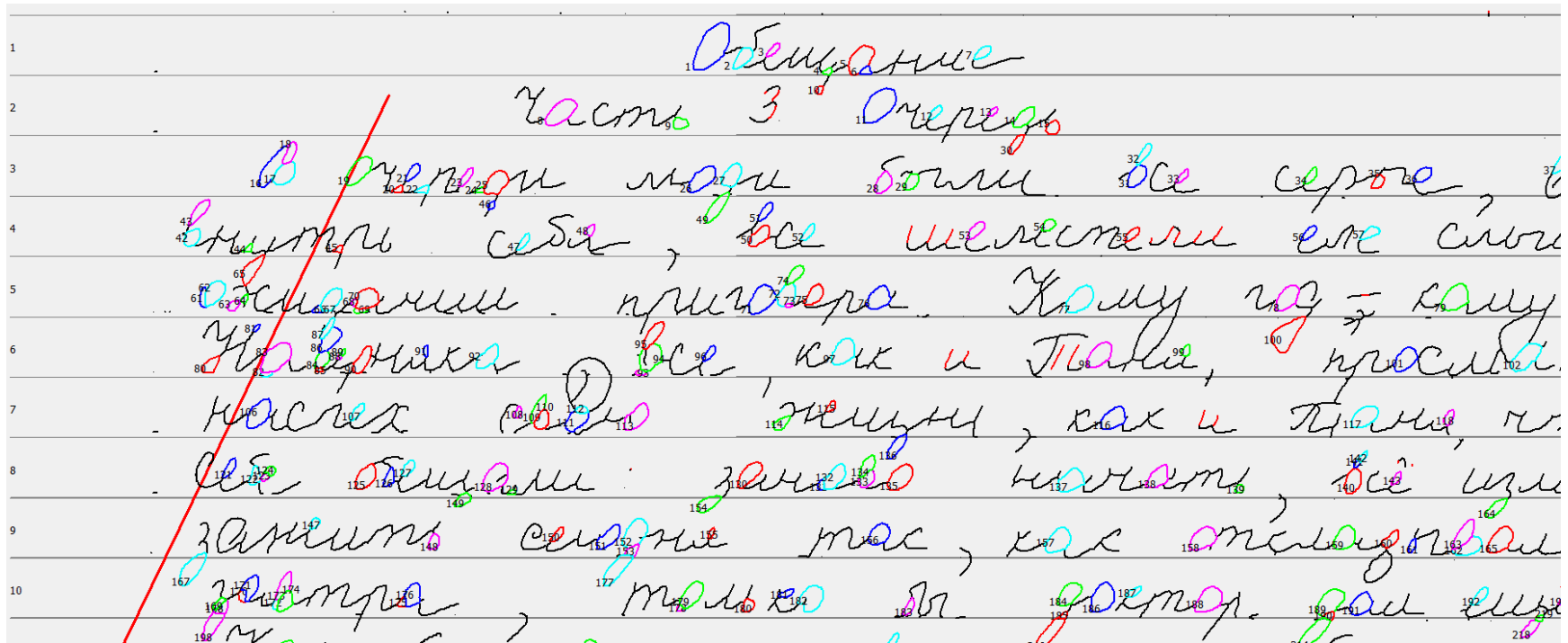
# Построение скелета

Обещание

Часть 3. Очередь

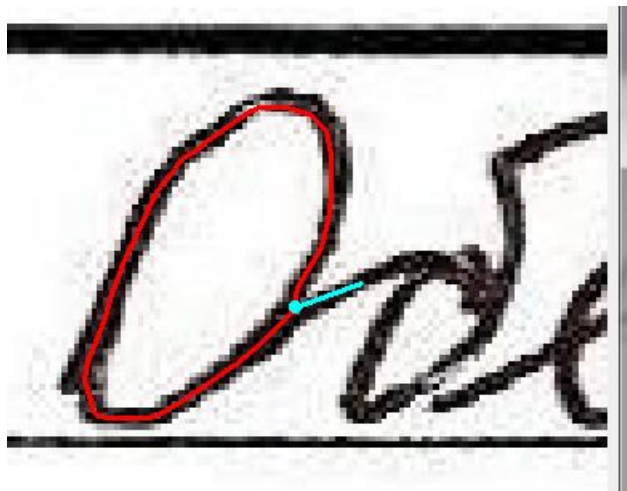
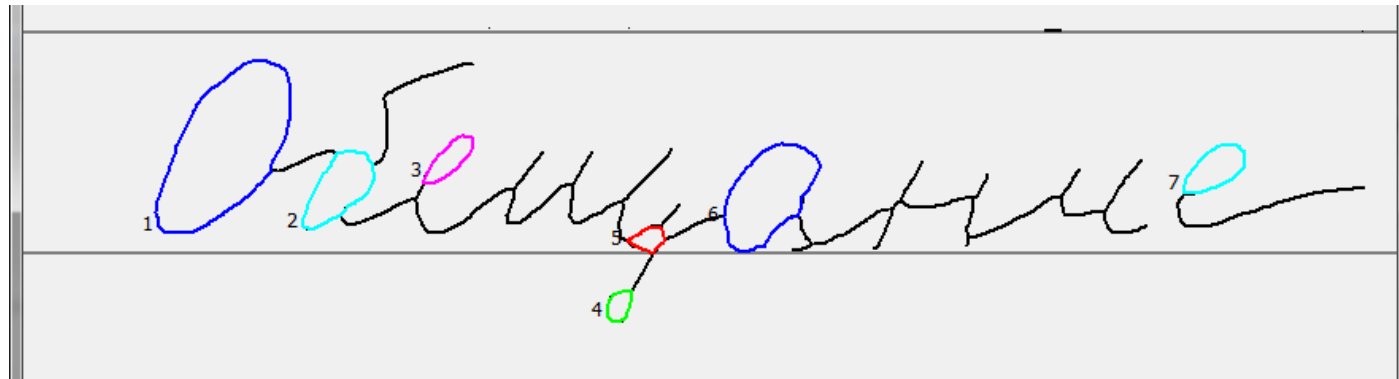
В очереди люди были все серые, все в  
внутри себя, все шелестели еле слышно  
ожиданием приговора. Кому год — кому ту  
Наверняка все, как и Памя, просидели  
наследство свою жизнь, как и Памя, что-то  
себе обещали: заново начать, всё изменить  
занести сегодня так, как откладывал на  
завтра только я. Вот она еще в

# Выделение овалов



Скелетное представление позволяет вычислить наклон письма (красная линия)

# Штрихи-ростки на овалах



Набор ростков является индивидуальным для каждой буквы и используется для классификации

# Задача обучения

## Распознавание букв

- Датасет – 35 диктантов на 2 страницах каждый
- 30 строк на странице по 30-35 букв в строке – это 1000 букв на странице
- 600 овалов на странице
- Требуемое распознавание – ранжированный список возможных вариантов для каждой буквы

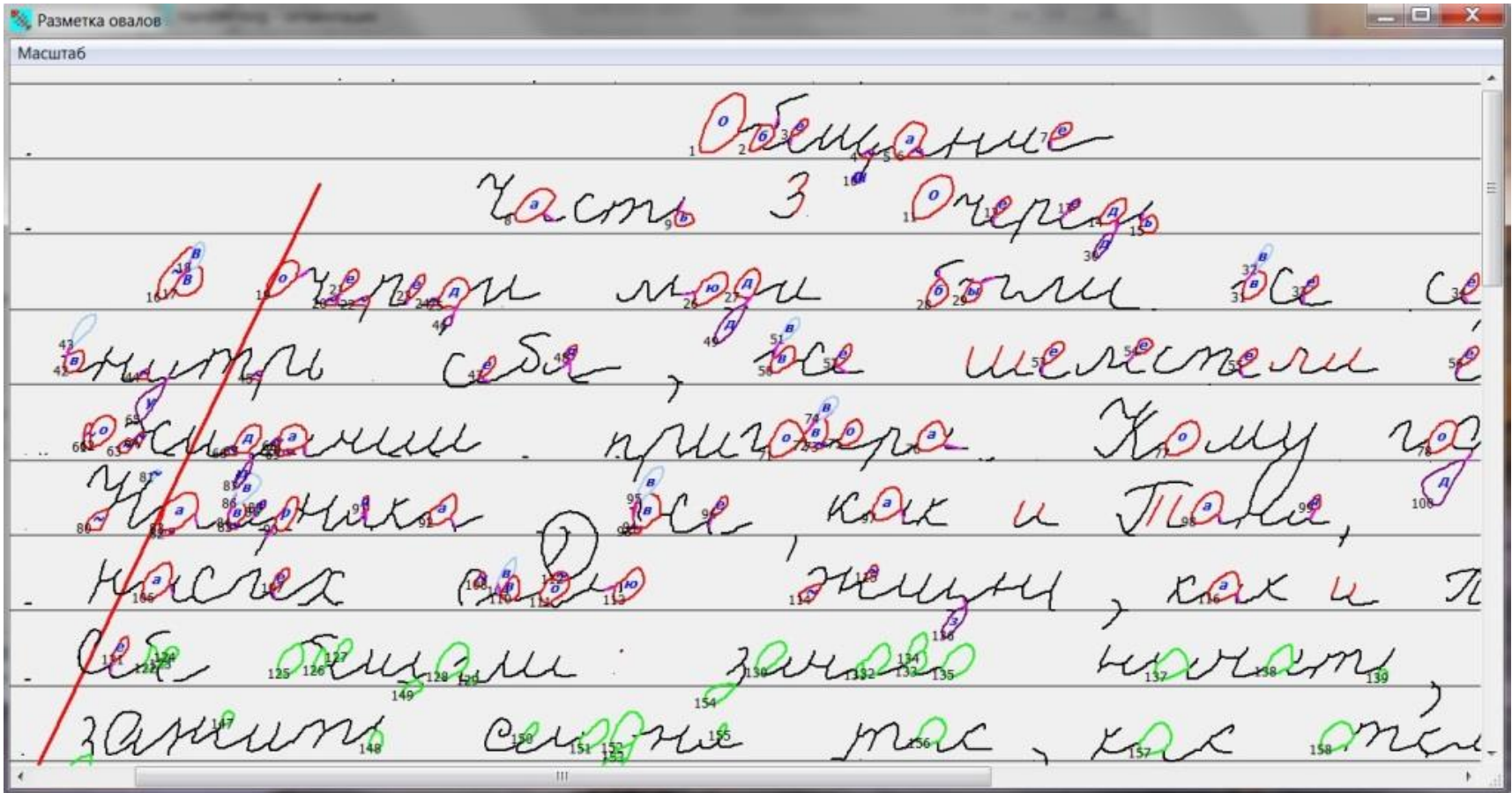
## Распознавание слов

- Синтаксические правила, работа со словарём

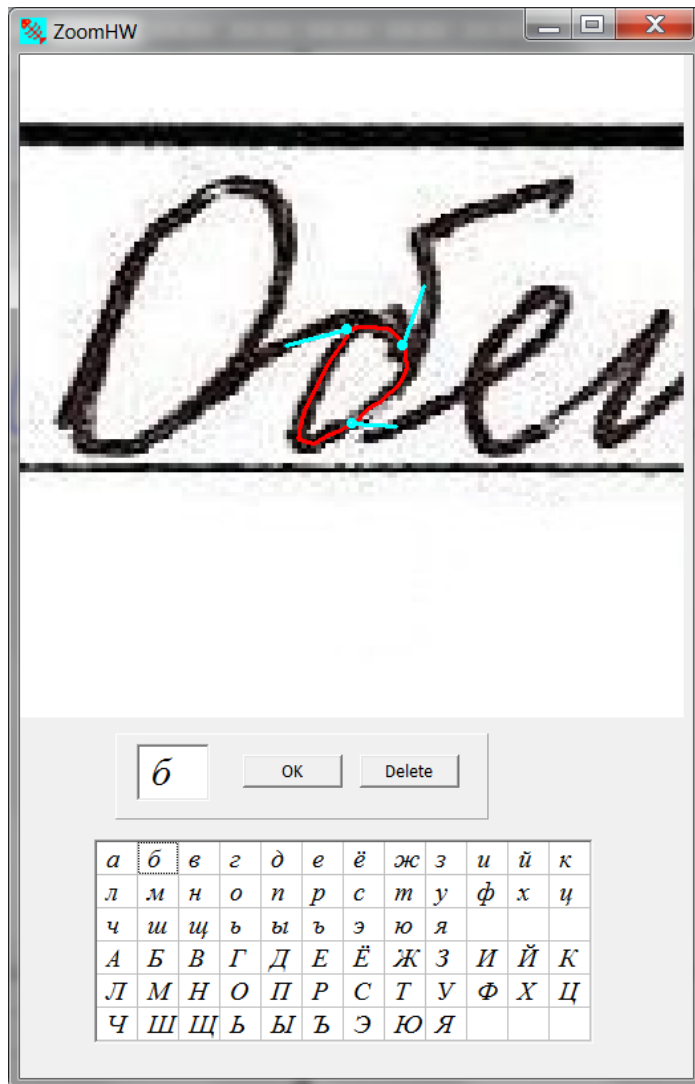
## Распознавание текста

- Семантический анализ , тематические модели

# Разметка рукописных текстов



# Разметка овалов и ростков



## Овал:

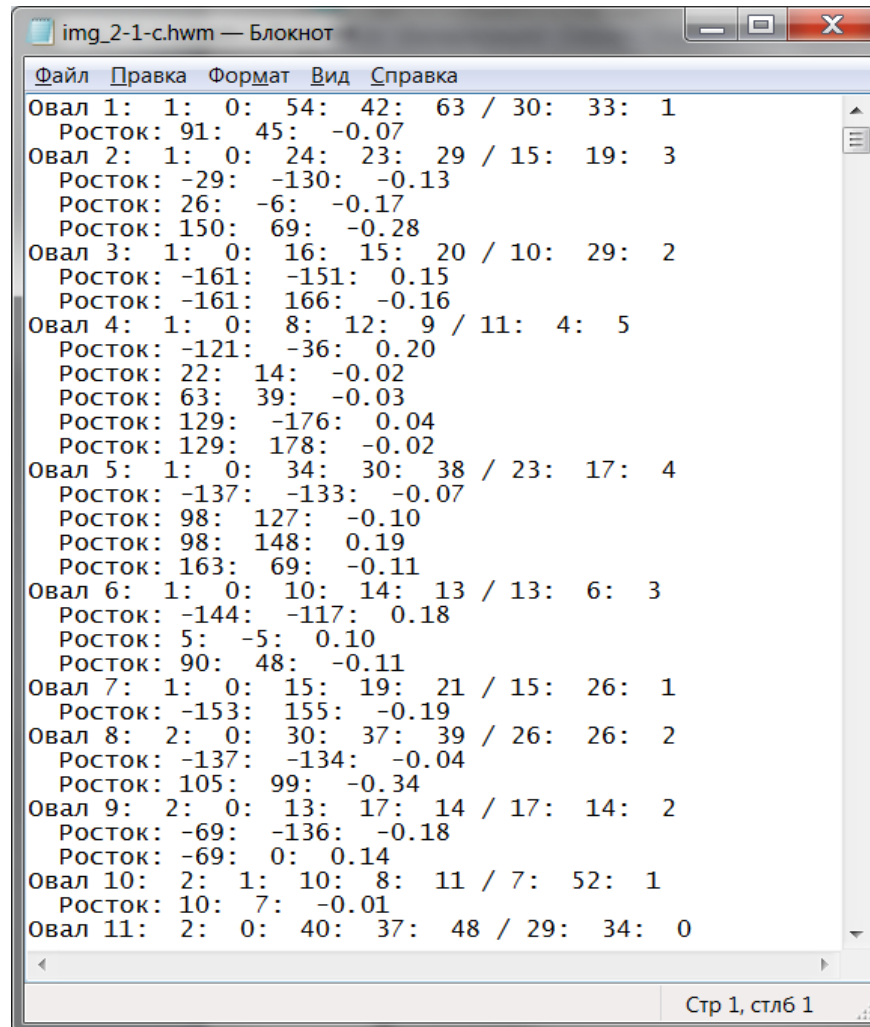
- номер овала
- буква
- рамка (высота, ширина)
- координаты центра над базовой линией

## Росток:

- азимут точки начала,
- направление
- кривизна,
- значимость



# Признаковое описание овалов



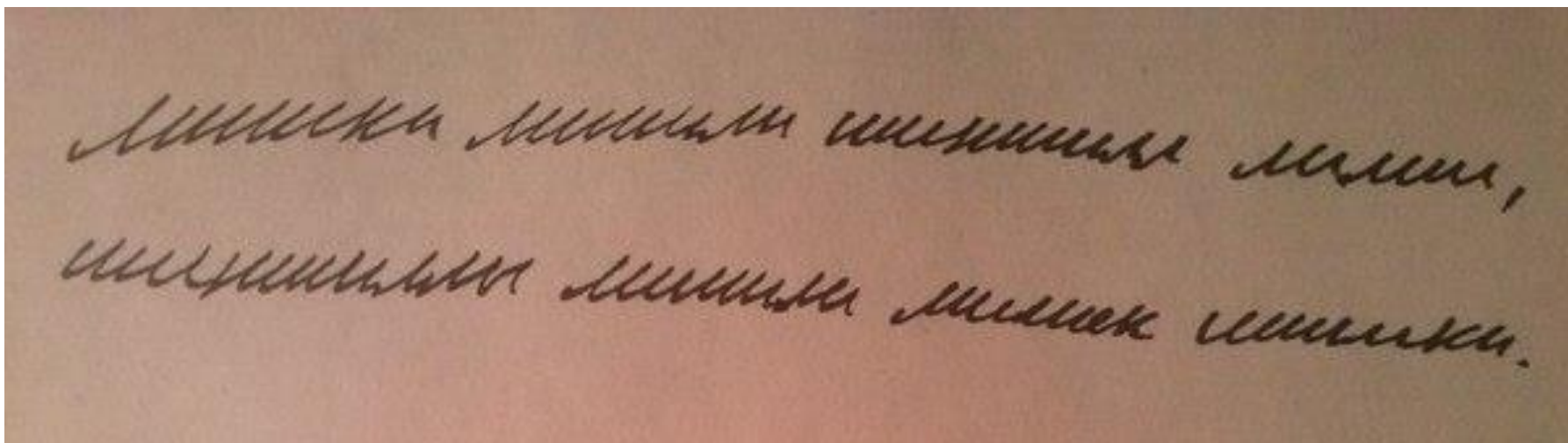
```
img_2-1-c.hwm — Блокнот
Файл Правка Формат Вид Справка
Овал 1: 1: 0: 54: 42: 63 / 30: 33: 1
Росток: 91: 45: -0.07
Овал 2: 1: 0: 24: 23: 29 / 15: 19: 3
Росток: -29: -130: -0.13
Росток: 26: -6: -0.17
Росток: 150: 69: -0.28
Овал 3: 1: 0: 16: 15: 20 / 10: 29: 2
Росток: -161: -151: 0.15
Росток: -161: 166: -0.16
Овал 4: 1: 0: 8: 12: 9 / 11: 4: 5
Росток: -121: -36: 0.20
Росток: 22: 14: -0.02
Росток: 63: 39: -0.03
Росток: 129: -176: 0.04
Росток: 129: 178: -0.02
Овал 5: 1: 0: 34: 30: 38 / 23: 17: 4
Росток: -137: -133: -0.07
Росток: 98: 127: -0.10
Росток: 98: 148: 0.19
Росток: 163: 69: -0.11
Овал 6: 1: 0: 10: 14: 13 / 13: 6: 3
Росток: -144: -117: 0.18
Росток: 5: -5: 0.10
Росток: 90: 48: -0.11
Овал 7: 1: 0: 15: 19: 21 / 15: 26: 1
Росток: -153: 155: -0.19
Овал 8: 2: 0: 30: 37: 39 / 26: 26: 2
Росток: -137: -134: -0.04
Росток: 105: 99: -0.34
Овал 9: 2: 0: 13: 17: 14 / 17: 14: 2
Росток: -69: -136: -0.18
Росток: -69: 0: 0.14
Овал 10: 2: 1: 10: 8: 11 / 7: 52: 1
Росток: 10: 7: -0.01
Овал 11: 2: 0: 40: 37: 48 / 29: 34: 0
```

Стр 1, стлб 1

# Выводы

- Предлагается подход к распознаванию рукописного текста на основании выделения и распознавания каллиграфических элементов письма и букв
- Выделение элементов осуществляется на основе скелетного представления изображения текста и геометрического анализа срединных осей
- Разработан метод выделения и классификации овальных каллиграфических элементов письма
- Классификация осуществляется на основе машинного обучения по размеченным образцам текста
- Основой эффект от предлагаемого подхода ожидается за счет высокой информативности и малой размерности признакового описания геометрических характеристик базовых каллиграфических элементов письма

# Спасибо за внимание!



Мишки лишили шиншиллы лилии, шиншиллы лишили мишек шишки