



## ТЕОРЕТИКО-ИНФОРМАЦИОННЫЕ ГРАНИЦЫ И ЗАДАЧИ ЗАЩИТЫ ИНФОРМАЦИИ, СВЯЗАННЫЕ С ОБРАБОТКОЙ ЕСТЕСТВЕННОГО ЯЗЫКА

Мещеряков Р.В., д.т.н., проф., профессор РАН (ИПУ РАН).

**Мельников С.Ю.**, д.ф.-м.н. (РУДН, ООО «Линфо»)

*Работа поддержана грантом РНФ 24-11-00340 Исследование и разработка методов обработки слабоструктурированной информации на естественных языках в условиях сильных шумов для решения задач безопасности*

16.05.24, г. Санкт-Петербург



# ПЛАН ДОКЛАДА

## 1. Зашумление (искажение) текста

1.1. О способах зашумления

1.2. Задача восстановления зашумленного текста как оптимизационная задача

1.3. Теоретико-информационные (шенноновские) границы уровня искажений, в которых возможно/невозможно восстановление зашумленного текста

## 2. Что можно понять из зашумленного текста, если его теоретически нельзя восстановить?

2.1. Подход на основе локальных флуктуаций энтропии текста

2.2. Подход на основе статистической оценки признаков текста, таких как язык, тематика, стиль, авторство и пр.

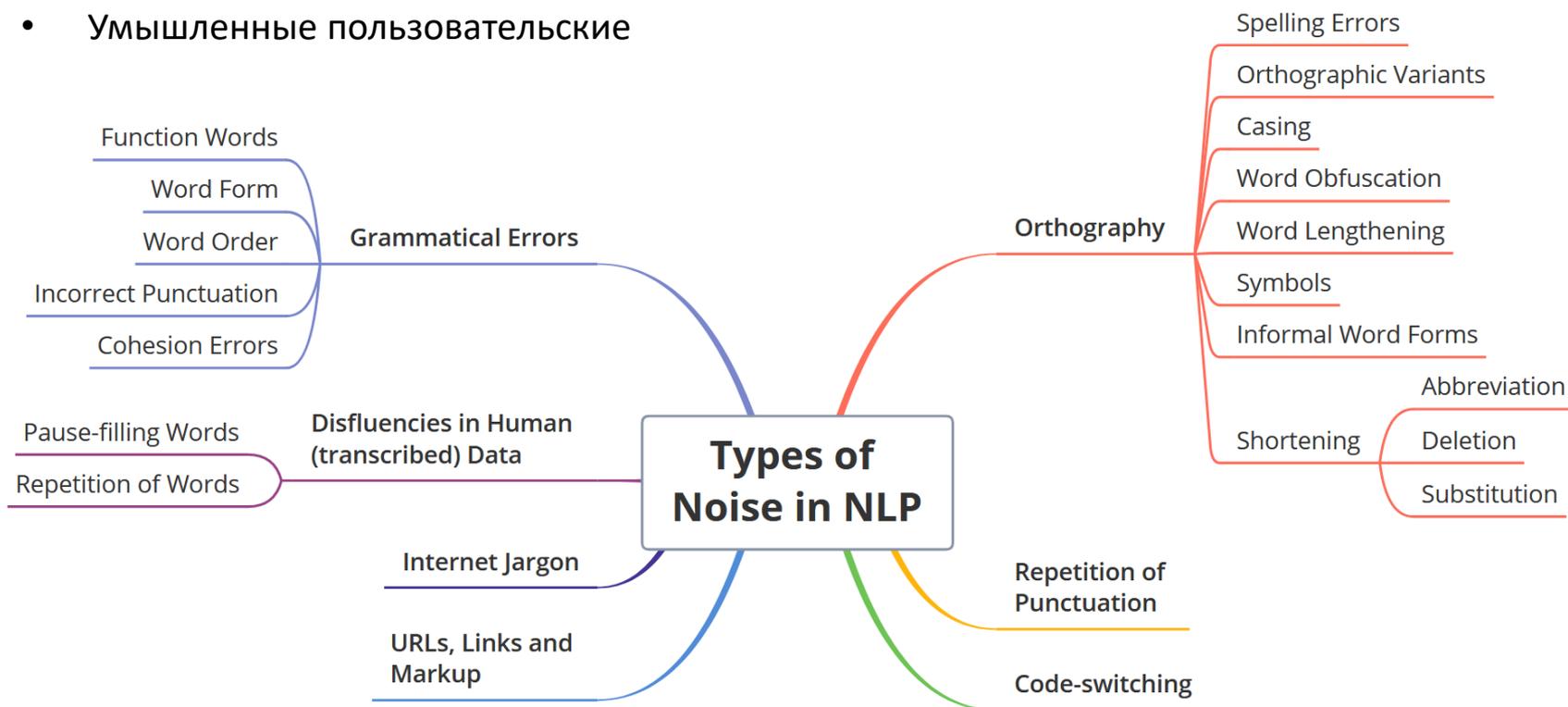
## 3. Опыты по пониманию речи в сильных шумах

## 4. Выводы

# «ЕСТЕСТВЕННЫЕ» ИСКАЖЕНИЯ ТЕКСТА

Искажения в тексте:

- Случайные
- В результате распознавания
- Умышленные пользовательские

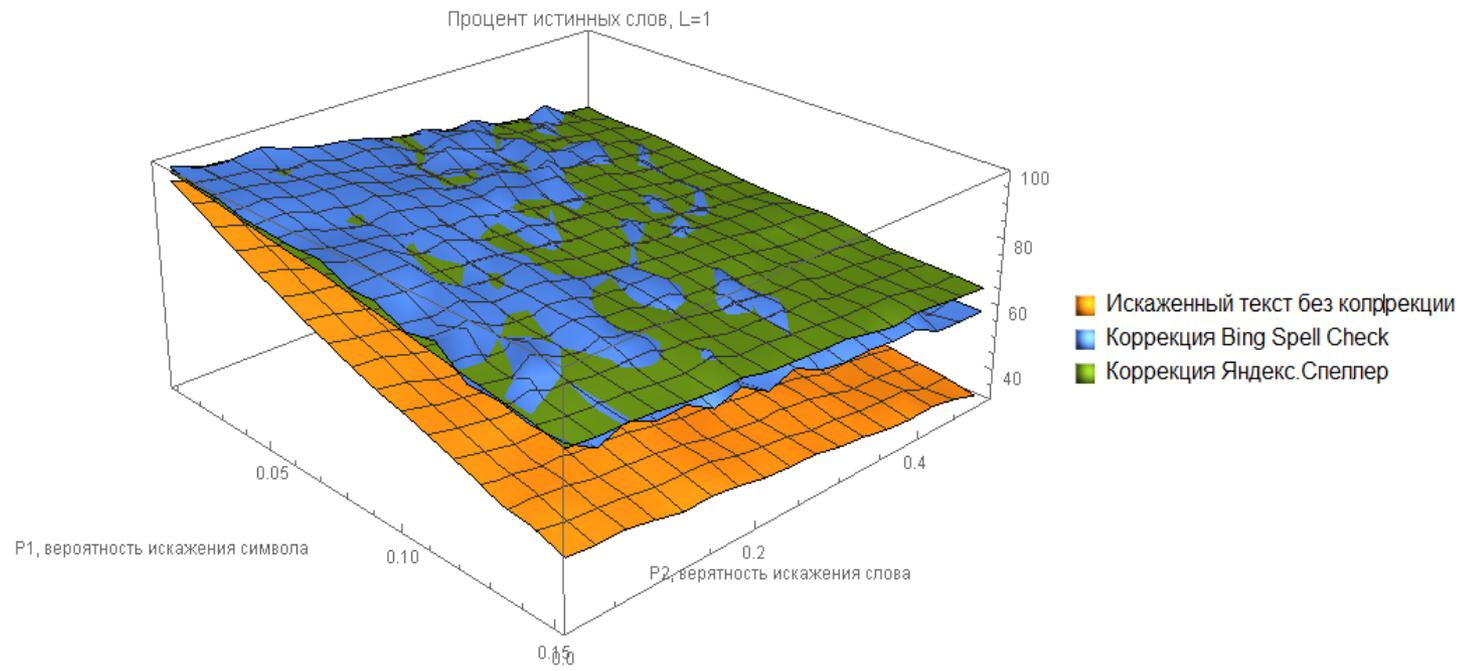




# ДВУХПАРАМЕТРИЧЕСКАЯ СХЕМА ИСКАЖЕНИЙ. РЕЗУЛЬТАТЫ КОРРЕКЦИИ ОТ BING И ЯНДЕКС

$P1$  – вероятность замены слов на близкие по Левенштейну,  $P2$  – вероятность посимвольных искажения,  $L$  – расстояние по Левенштейну

Область параметров модели искажений:  $0.01 < P2 < 0.9$ ;  $0.01 < P1 < 0.5$ ;  $L=1,2$ .



Поверхности  $Noisy(P1,P2)$ ,  $Yand(P1,P2,1)$ ,  $Bing(P1,P2,1)$ .

Д.А. Бирин, С. Ю. Мельников, В. А. Пересыпкин Об эффективности средств коррекции искаженных текстов в зависимости от характера искажений, 2018

# ПРИМЕРЫ ЗАШУМЛЕННЫХ ТЕКСТОВ. СИЛЬНЫЕ, СРЕДНИЕ И СЛАБЫЕ ИСКАЖЕНИЯ



8eb681f801374166a5ce8306b0a43ea4\_4.txt — Блокнот

Файл Правка Формат Вид Справка

Мскве. 4 бая. INTERFAX.RU - Итатльясре ВМС првлозят Еврпнеельмик опжерацієюно дставканищпушу йптримхеурною 5,7 ты. неьягаълom, спщахпнымйуаза япоследеедвое оуырщк а Сревиземно чора, сообщаемц агентствБЕFE.p в честнст, фегапте "Берсальее"у ужай щюосйвлв сорг Рляжо-ди-Клабрия778кмийратто. Ещек870й жчеловека хрифыхлх ншаз Сюилия на сборуькораблае берхговый охрнн"Аско-29". Кроге кяою, ы итальянсвий поёрть посмавлены тжщбе энискоколко ммене мноочиблеунгьём фиупепа нелалм. ов 2015ъмводу наблюдаемвся ролу поточки нлеральые жбдгрантош, пыающийёсн вохбратаьья пи морвчз Северн Афмикэ к Еровпы. Блолфшинство из цш отикльваюизэ лкиоиув нкправлениял Иалиц.

8eb681f801374166a5ce8306b0a43ea4\_4.txt — Блокнот

Файл Правка Формат Вид Справка

Москва. 4 мая. INTERFAX.RU - Итааляньские ВМС прооди в понедельникоперацію по доставкеца суму примерна 5,7 тыс. ьелегалов, спасенных за пследние двое скток в Средиземном мре,сообщает агетство EFE. В частноти, фрегат"Берсальекр" уже доставит во порт Реджо-ди-Калабрия 778 мигрантов. Еще 870 человеку прибыли на Сицилию нибоярту корабля хлересголой охтаншу "Ассо-29". Кроме того, в итальянские порты доставлены еще несколько менте многочисаленнх гхупп неёегалов. В 2015 году нлблюдается роста потока нелкегалных мигрантом, пытающихся дбрасья по молю из Северной Африки в Европу. Большинство и тих отплыва виз Лидии в напленияИталии.

8eb681f801374166a5ce8306b0a43ea4\_4.txt — Блокнот

Файл Правка Формат Вид Справка

Москва. 4 мая. INTERFAX.RU - итальянские ВМС ыроводетйв понмдельник операцію по доставке на сушу примено 5,7 тыб. нелегалов, спасенных за последние двое суток в Средиземном море, ёсообщает агентстве EFE. В частности, фрегат "Берсальер" уже доставил в порт Реджо-ди-Калабрия 778мигрантов. Еще 870 человек прибыли на Сшцилию на борту корабля берегаовой охранф "Ассо-29". Кроье того, в итальянские порты доставлены еще несколько менее многочисленных групп нелегалов. йв 2015 году наблюдается рост потока нелегальных мигрантов, пытающихся добаться по морю уз Северной Афярики вЕвропу. Большинство виз них отпйиват из ливии в направлении Италии.

Даже слабые искажения делают невозможным автоматическую обработку текста.

Kumar A., Makhija P., Gupta A. Noisy Text Data: Achilles' Heel of BERT, 2020.

# ТЕОРЕТИКО-ИНФОРМАЦИОННЫЙ ПОДХОД. ВОССТАНОВЛЕНИЕ ТЕКСТА КАК ОПТИМИЗАЦИОННАЯ ЗАДАЧА



Пусть канал связи характеризуется набором условных вероятностей  $P(S/X)$  появления последовательности наблюдений  $S = (s_1, s_2, \dots, s_n)$  при условии, что истинная последовательность есть  $X = (x_1, x_2, \dots, x_n)$ ,  $n = 1, 2, \dots$ . Задача коррекции состоит в нахождении такой последовательности  $\hat{X} = (\hat{x}_1, \hat{x}_2, \dots, \hat{x}_n)$  символов из алфавита, для которой достигается максимум условной вероятности  $P(X/S)$ :

$$\hat{X} = \arg \max_x P(X/S).$$

По формуле Байеса:

$$\arg \max_x P(X/S) = \arg \max_x P(X) \times P(S/X).$$

Таким образом, для коррекции текста необходимо решить оптимизационную задачу

$$\text{Модель языка} \quad P(X) \times P(S/X) \rightarrow \max. \quad \text{Модель искажений}$$
A diagram showing the equation  $P(X) \times P(S/X) \rightarrow \max.$  with two blue curved arrows pointing towards it. One arrow starts from the text 'Модель языка' (Language Model) on the left and points to  $P(X)$ . The other arrow starts from the text 'Модель искажений' (Distortion Model) on the right and points to  $P(S/X)$ .

Величина  $P(X)$  определяется моделью языка и обычно представляется в виде

$$P(X) = \prod_{i=1}^n p(x_i / x_{1,i-1}),$$

где  $p(x_i / x_{1,i-1})$  – вероятность появления символа  $x_i$  после фрагмента  $x_1, x_2, \dots, x_{i-1}$ ,  $i = 1, 2, \dots, n$ .

Величина  $P(S/X)$  определяется способом искажения текста, ее вид зависит от предметной области, в которой решается задача коррекции.



# МЕТОДЫ ВОССТАНОВЛЕНИЯ В ЗАВИСИМОСТИ ОТ УРОВНЯ ИСКАЖЕНИЙ

## Методы

Коррекция отдельных ошибок по словарю

Восстановление текста методами поиска с использованием моделей языка

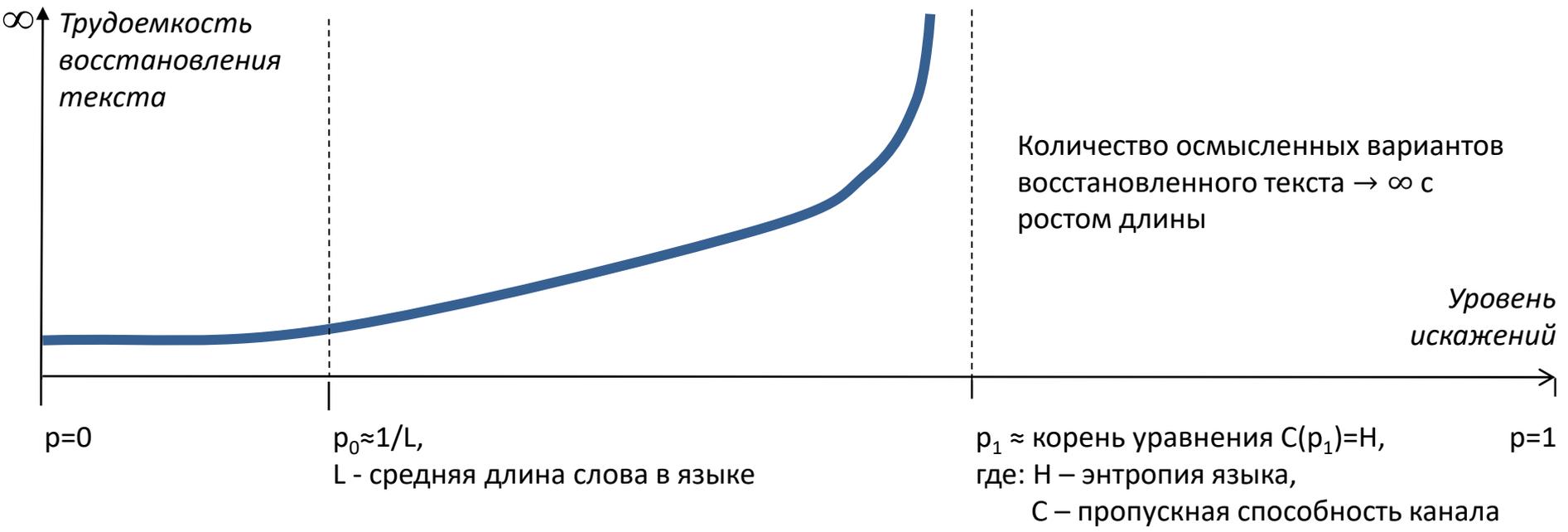
Восстановление текста невозможно

Распространенные спелл-чекеры

Специализированные системы

???

## Трудоемкость





# ЛОКАЛЬНЫЕ ФЛУКТУАЦИИ ЭНТРОПИИ ТЕКСТА СДВИГАЮТ ГРАНИЦУ «НЕВОЗМОЖНОСТИ ВОССТАНОВЛЕНИЯ»

В текстах, в зависимости от их жанра и стиля, могут встречаться низкоэнтропийные участки.

Примеры:

*«В целях исключения открытого опубликования материалов, содержащих сведения, составляющих государственную тайну, и сведений, подпадающих под действие контрольных списков по экспортному контролю, в соответствии с федеральными законами от 21.07.1993 № 5485-1 «О государственной тайне», от 18.07.1999 № 183-ФЗ «Об экспортном контроле» и во исполнение Рекомендаций по проведению экспертизы материалов, предназначенных к открытому опубликованию, одобренных Межведомственной комиссией по защите государственной тайны от 24.01.2012 № 225»... (Приказ Тольяттинского государственного университета № 4409 от 10.12.14г.)».*

*«... и лично генеральный секретарь Центрального комитета Коммунистической партии Советского Союза Леонид Ильич Брежнев»*

Наличие подобных участков приводит к локальному снижению энтропии текста и может значительно сдвинуть границу «невозможности восстановления».

Life.ru  
life.ru › Wow › Интересное

**10 образов Брежнева, после которых ни у кого не...**  
Генеральный секретарь Центрального комитета Коммунистической партии Советского Союза Леонид Ильич Брежнев и президент США Ричард Никсон (слева направо)... Читать ещё

Sports.ru  
sports.ru › Брежнев, каким мы его знали, в день 110-летнего «юбилея» б...

**Брежнев, каким мы его знали, в день 110-летнего...**  
14 июня 1977 года генеральный секретарь Центрального Комитета Коммунистической Партии Советского Союза Леонид Ильич Брежнев вместе с Первым секретарём ЦК КП Грузинской Советской Социалистической республики... Читать ещё

Picturehistory.Livejournal.com  
picturehistory.livejournal.com

**Брежнев, Ульбрихт и медведь )))**: picturehistory — ЖЖ  
Всем известно, что генеральный секретарь Центрального комитета Коммунистической партии Советского Союза Леонид Ильич Брежнев, помимо многочисленных заслуг перед народом и Отечеством, был еще и заядлым охотником. Читать ещё

Vk.com  
vk.com › wall-211186281\_70669

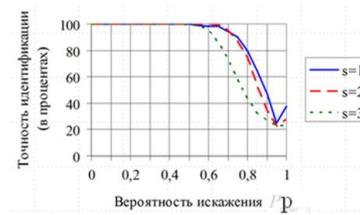
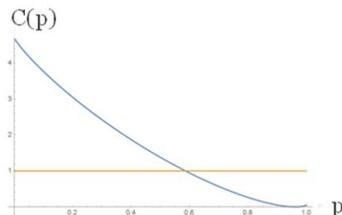
**Как Леонид Брежнев в Донецкую область приезжал...**  
И в этот раз на активе должен был присутствовать высокий гость даже не из ЦК КПУ, а куда выше — Генеральный секретарь Центрального комитета Коммунистической партии Советского Союза Леонид Ильич Брежнев. Читать ещё



# ИСПОЛЬЗОВАНИЕ СТАТИСТИЧЕСКИХ ОЦЕНОК ПРИЗНАКОВ ТЕКСТА В ЗОНЕ «НЕВОЗМОЖНОСТИ ВОССТАНОВЛЕНИЯ»

Такие признаки текста, как язык, тематика, стиль, жанр, авторская принадлежность, и др. могут быть установлены статистически на основе распределения символьных N-грамм.

Тексты на 4-х европейских языках приводились к латинскому алфавиту без пробелов и знаков препинания. Задавалось значение  $p \in (0,1)$  вероятности замены символа. Замены осуществлялись независимо друг от друга. Если решение о замене символа принималось, то он заменялся на произвольный другой символ латинского алфавита по равновероятной схеме. Построен дискретный канал без памяти, который моделирует описанные искажения. Вычислена его пропускная способность  $C(p)$ .



Построена процедура распознавания языка рассматриваемых текстов, использующая  $s$ -граммные статистики на символах,  $s=1,2,3$ . Согласно основному теоретико-информационному неравенству, исходный текст можно восстановить при условии  $H < C$ , где  $H \approx 1$  - энтропия языка исходного текста, а  $C$  - пропускная способность. Для построенного канала при  $p < 0.59$  можно восстанавливать исходный текст по искаженному, а при  $p > 0.59$  однозначное восстановление текста невозможно. Результаты экспериментов показывают, что язык текста можно определять с высокой степенью надежности даже при  $p > 0.59$ .

Поэтому, даже в случае, когда восстановить текст нельзя, его содержание можно как-то охарактеризовать.



# ОПЫТЫ ПО ПОНИМАНИЮ РЕЧИ В СИЛЬНЫХ ШУМАХ

В период 2010-2019 гг в МГЛУ проводилась серия исследований по восприятию речи в условиях сильных шумов.

*Ю. В. Абрамов, Р. К. Потапова, М. В. Хитина. Специфика смыслового восприятия спонтанного звучащего текста (в условиях эскалации уровня шума) (2011)*

*Хитина М. В. Экспериментальное исследование оценки качества смыслового восприятия спонтанного монолога на фоне аддитивных шумов (2012)*

*Хитина М. В. Перцептивно-слуховое восприятие материалов разного вида: экспериментальные исследования (2017)*

*Хитина М. В. Особенности восприятия читаемых текстов в "белом шуме" (2019) и др.*

В исследованиях принимало большое количество экспертов (студентов-филологов).

Исследовался, в том числе, вопрос, может ли эксперт воспринять основную и побочные темы произнесенного речевого фрагмента.

Одним из результатов этих исследований был вывод о сильном влиянии квалификации, опыта и индивидуальных качеств воспринимающего речь эксперта и, как следствие, значительном разбросе получаемых оценок.



# ВЫВОДЫ

1. Искажения в текстах и зашумление речевого сигнала существенно затрудняют их восприятие человеком и делают практически невозможной автоматическую обработку.
2. Избыточность естественного языка позволяет восстанавливать искаженные тексты, вплоть до уровня шумов, соответствующего шенноновской теоретико-информационной границе. Если уровень искажений превышает эту границу, однозначное восстановление текста невозможно.
3. Для случая, когда уровень искажений в тексте превышает теоретико-информационную границу, предложено два подхода к получению информации из искаженного текста. Первый связан с поиском низкоэнтропийных участков в тексте и их восстановлением, второй связан с возможностью получения интегральной информации о тексте (язык, тематика, жанр и др.) по статистике встречаемости символьных N-грамм.
4. Разработка этих подходов позволит уточнить требования и методики по защите информации для ряда каналов связи.



СПАСИБО ЗА ВНИМАНИЕ

Спасибо за внимание!

[melnikov@linfotech.ru](mailto:melnikov@linfotech.ru)