

# MISIS AI LAB

ALPHA

Автоматизация работы с документами  
посредством введения нейронных сетей



## Наша команда



**Савельев Ярослав**  
Robotics engineer -  
Sber Robotics Lab



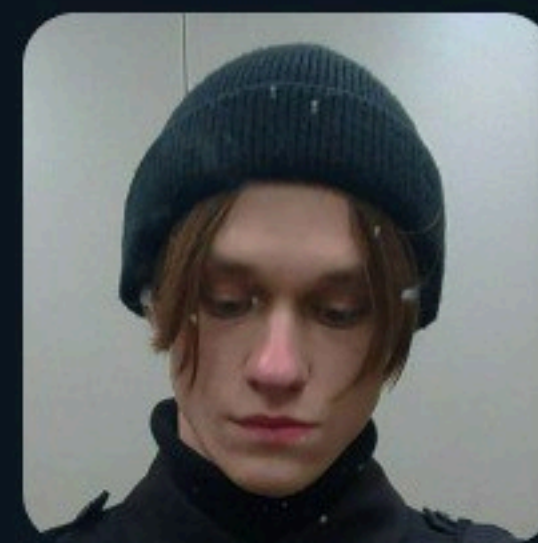
**Комлев Данила**  
Frontend/Backend - INVITRO



**Иванов Арсений**  
Исследователь данных -  
Sber AmazMe



**Волков Даниил**  
MLOps engineer -  
ГосУслуги



**Максим**  
ML engineer -  
Точка Банк



# ПОСТРОЕНИЕ ИНТЕРПРЕТИРУЕМОСТИ

**BOW/TF-IDF НА N-ГРАММАХ ПО  
СЛОВАМ + ДЕРЕВЬЯ/ЛОГРЕГ**

->

**ПЛЮСЫ:** ЛЕГКО ИНТЕРПРЕТИРОВАТЬ И  
ЗАХВАТЫВАТЬ ВСЬ ДОКУМЕНТ  
**МИНУСЫ:** КЛЮЧЕВЫЕ СЛОВА НЕПОНЯТНЫ  
В КОНТЕКСТЕ

**BERT ПО ПЕРВЫМ 512 ТОКЕНАМ  
ДОКУМЕНТА + ГРАДИЕНТНЫЕ/  
ATTENTION МЕТОДЫ ИНТЕРПРЕТАЦИИ**

->

**ПЛЮСЫ:** БЫСТРО РЕАЛИЗОВАТЬ  
**МИНУСЫ:** НЕУДОБНО ЗАНОВО  
ПРОСМАТРИВАТЬ ДОКУМЕНТ, УЧИТЫВАЕМ НЕ  
ВСЬ ТЕКСТ ДОКУМЕНТА, ТЯЖЕЛЫЙ ИНФЕРЕНС

**BERT ПО ФРАГМЕНТАМ  
ДОКУМЕНТА + УСРЕДНЕНИЕ  
СКОРА ФРАГМЕНТОВ**

->

**ПЛЮСЫ:** УДОБНО ИНТЕРПРЕТИРОВАТЬ,  
УЧИТЫВАЕМ ВСЬ ДОКУМЕНТ  
**МИНУСЫ:** ТРУДНО РЕАЛИЗОВАТЬ, ТЯЖЕЛЫЙ  
ИНФЕРЕНС

# МЕТРИКИ

		1 class	2 class	3 class	4 class	5 class	macro
	recall						0.97
BERT	f1-score						0.98
	precision						0.98
	recall						0.96
TFIDF + CatBoost	f1-score						0.97
	precision						0.97

	accuracy		
BERT	0.98		
TFIDF + CatBoost	0.97		

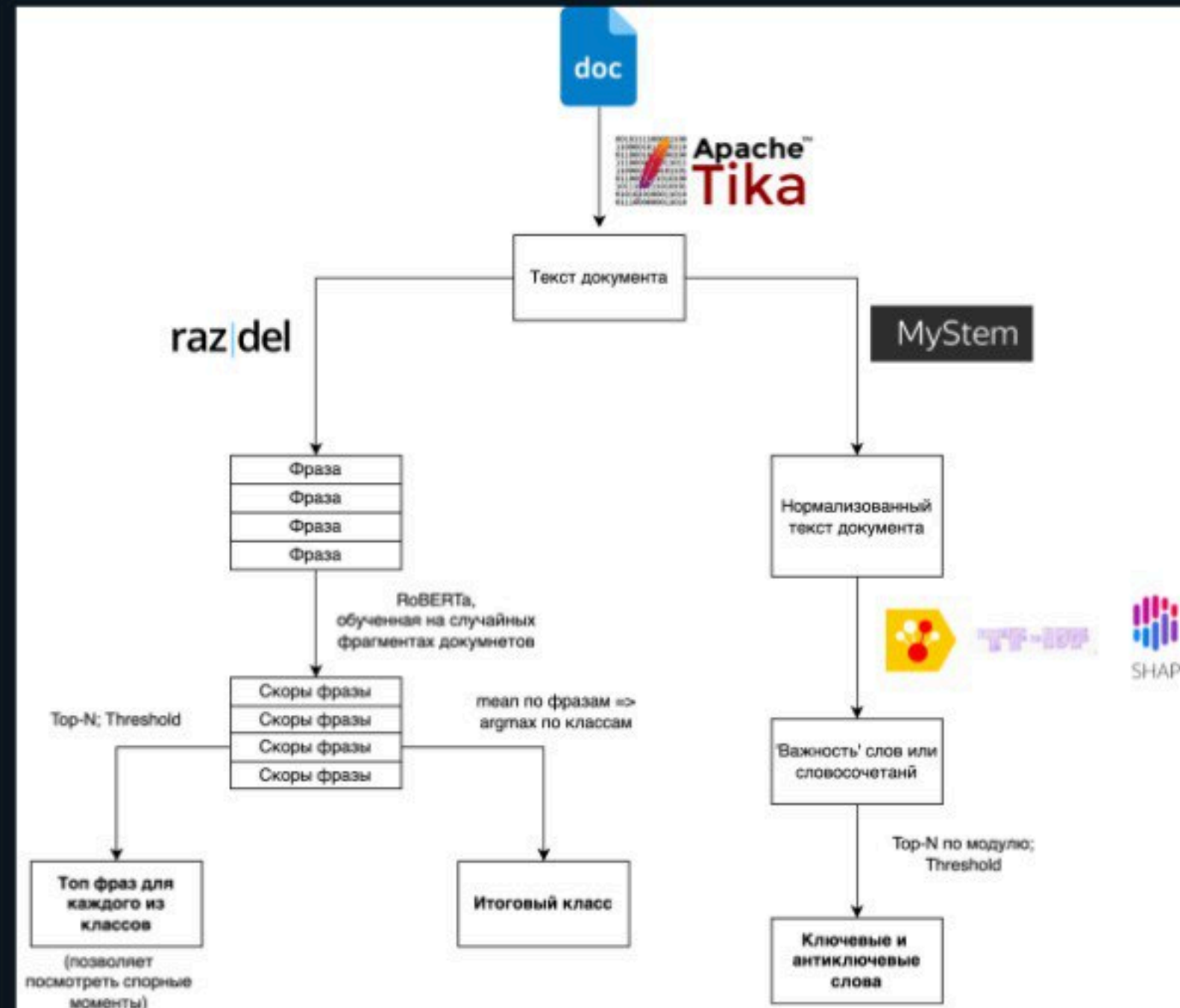
$$F1 = 2 * \frac{Precision * Recall}{Precision + Recall}$$

$$Precision = \frac{\sum TruePositive}{\sum TruePositive + FalsePositive}$$

$$Recall = \frac{\sum TruePositive}{\sum TruePositive + FalseNegative}$$



# ИНФЕРЕНС



# МАСШТАБИРОВАНИЕ МОДЕЛИ

Published as a conference paper at ICLR 2020

## DIVIDEMIX: LEARNING WITH NOISY LABELS AS SEMI-SUPERVISED LEARNING

Junnan Li, Richard Socher, Steven C.H. Hoi

Salesforce Research

{junnan.li, rsocher, shoi}@salesforce.com

### ABSTRACT

Deep neural networks are known to be annotation-hungry. Numerous efforts have been devoted to reducing the annotation cost when learning with deep networks. Two prominent directions include learning with noisy labels and semi-supervised learning by exploiting unlabeled data. In this work, we propose DivideMix, a novel framework for learning with noisy labels by leveraging semi-supervised learning techniques. In particular, DivideMix models the per-sample loss distribution with a mixture model to dynamically divide the training data into a labeled set with clean samples and an unlabeled set with noisy samples, and trains the model on both the labeled and unlabeled data in a semi-supervised manner. To avoid confirmation bias, we simultaneously train two diverged networks where each network uses the dataset division from the other network. During the semi-supervised training phase, we improve the MixMatch strategy by performing label co-refinement and label co-guessing on labeled and unlabeled samples, respectively. Experiments on multiple benchmark datasets demonstrate substantial improvements over state-of-the-art methods. Code is available at <https://github.com/LiJunnan1992/DivideMix>.

### 1 INTRODUCTION

The remarkable success in training deep neural networks (DNNs) is largely attributed to the collection of large datasets with human annotated labels. However, it is extremely expensive and time-consuming to label extensive data with high-quality annotations. On the other hand, there exist alternative and

## Beyond 512 Tokens: Siamese Multi-depth Transformer-based Hierarchical Encoder for Long-Form Document Matching

Liu Yang Mingyang Zhang Cheng Li Michael Bendersky Marc Najork

Google Research, Mountain View, CA, USA

{yangliuy, mingyang, chgli, bemike, najork}@google.com

### CT

al language processing and information retrieval problem formalized as the task of semantic matching. Existing area has been largely focused on matching between (e.g., question answering), or between a short and a (e.g., ad-hoc retrieval). Semantic matching between long-form documents, which has many important applications like news recommendation, related article recommendation and document clustering, is relatively less explored and needs more research effort. Recently, self-attention based models like Transformers [30] [5] have achieved state-of-the-art performance in the matching. These models, however, are still limited to process a few sentences or one paragraph due to the quadratic complexity of self-attention with respect to length. In this paper, we address the issue by proposing Multi-depth Transformer-based Hierarchical (SMITH) long-form document matching. Our model contains variations to adapt self-attention models for longer text. We propose a transformer based hierarchical encoder to capture document structure information. In order to better capture local semantic relations within a document, we pre-train with a novel masked sentence block language modeling on top of the masked word language modeling task used in experimental results on several benchmark datasets. Our document matching show that our proposed SMITH outperforms the previous state-of-the-art models including multi-head attention [34], multi-depth attention-based hierarchical neural network [14], and BERT. Comparing to BERT, our model is able to increase maximum input text length from 512 to 2048. We will open source a Wikipedia based dataset, code and a pre-trained checkpoint to accelerate research on long-form document matching.<sup>1</sup>

### 1.1.1.1

nce Format: Mingyang Zhang Cheng Li Michael Bendersky Marc

Hierarchical Encoder for Long-Form Document Matching. In *Proceedings of the 29th ACM International Conference on Information and Knowledge Management (CIKM '20)*, October 19–23, 2020, Virtual Event, Ireland. ACM, New York, NY, USA, 10 pages. <https://doi.org/10.1145/3340531.3411908>

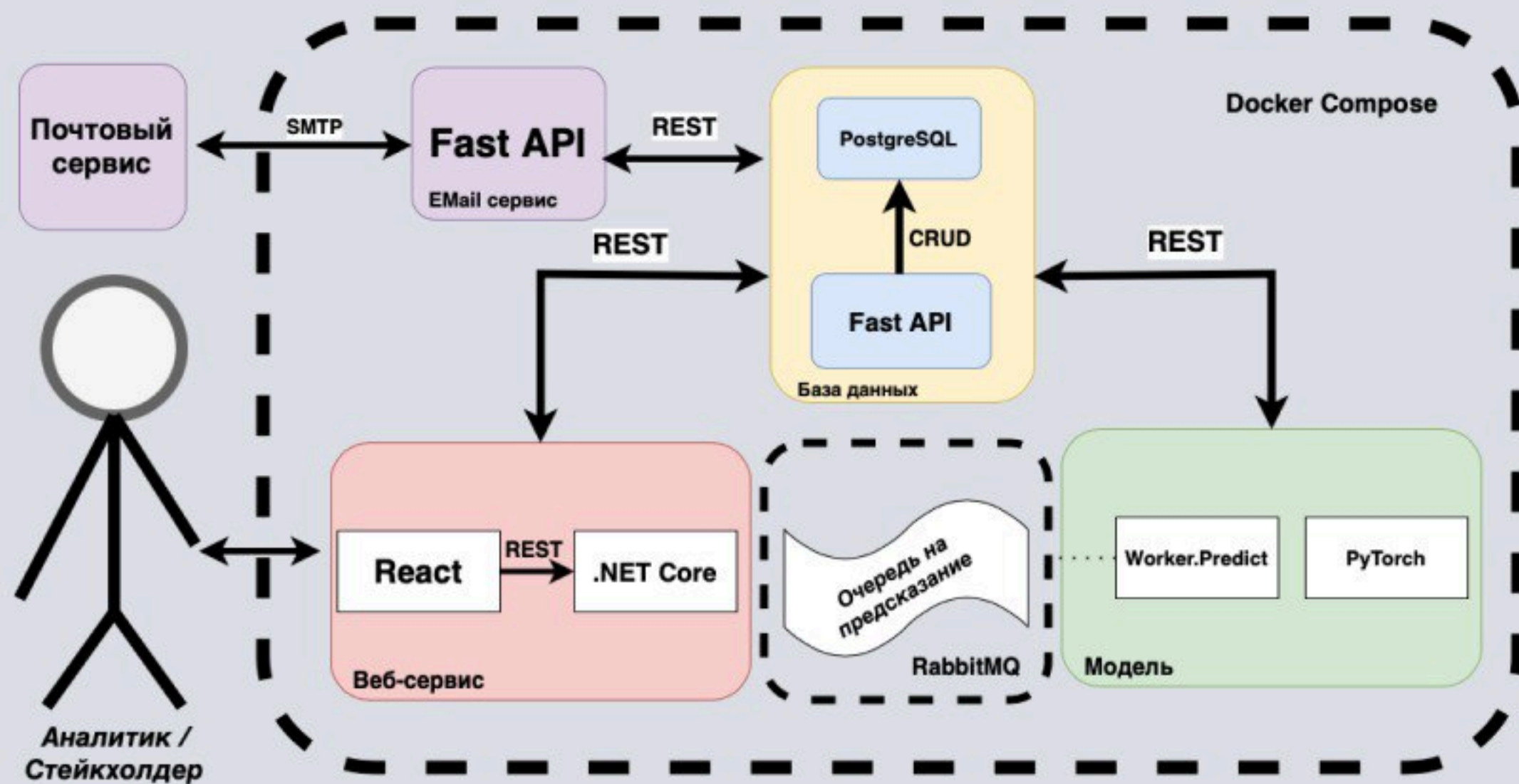
### 1 INTRODUCTION

Semantic matching is an essential task for many natural language processing (NLP) and information retrieval (IR) problems. Research on semantic matching can potentially benefit a large family of applications including ad-hoc retrieval, question answering and recommender systems [17]. Semantic matching problems can be classified into four different categories according to text length, including short-to-short matching, short-to-long matching, long-to-short matching and long-to-long matching. Table 1 shows a classification of different semantic matching tasks with example datasets. Semantic matching between short text pairs is relatively well studied in previous research on paraphrase identification [38], natural language inference [1], answer sentence selection [35], etc. Short-to-long semantic matching like relevance modeling between query/document pairs has also been a popular research topic in IR and NLP communities [3]. For long-to-short semantic matching, there are also a variety of research on tasks like conversation response ranking, which is to match a conversation context with response candidates [18]. To the best of our knowledge, semantic matching between long document pairs, which has many important applications like news recommendation, related article recommendation and document clustering, is less explored and needs more research effort. Table 2 shows an example of semantic matching between document pairs from Wikipedia. These documents have thousands of words organized in sections, passages and sentences.

Compared to semantic matching between short texts, or between short and long texts, semantic matching between long texts is a more challenging task due to a few reasons: 1) When both texts are long, matching them requires a more thorough understanding of semantic relations including matching pattern between text



# BACKEND



# FEATURES

В базе данных хранится сопоставление типов документов и емэйлов. После определения типа документа, данный файл будет отправлен на все емейлы, которые привязаны к этому типу документа.

## Документ

Волков Даниил Сегодня, 12:06  
Кому: вам



1 файл Скачать (297 КБ) Сохранить в Облако

Добрый день, уважаемые коллеги!

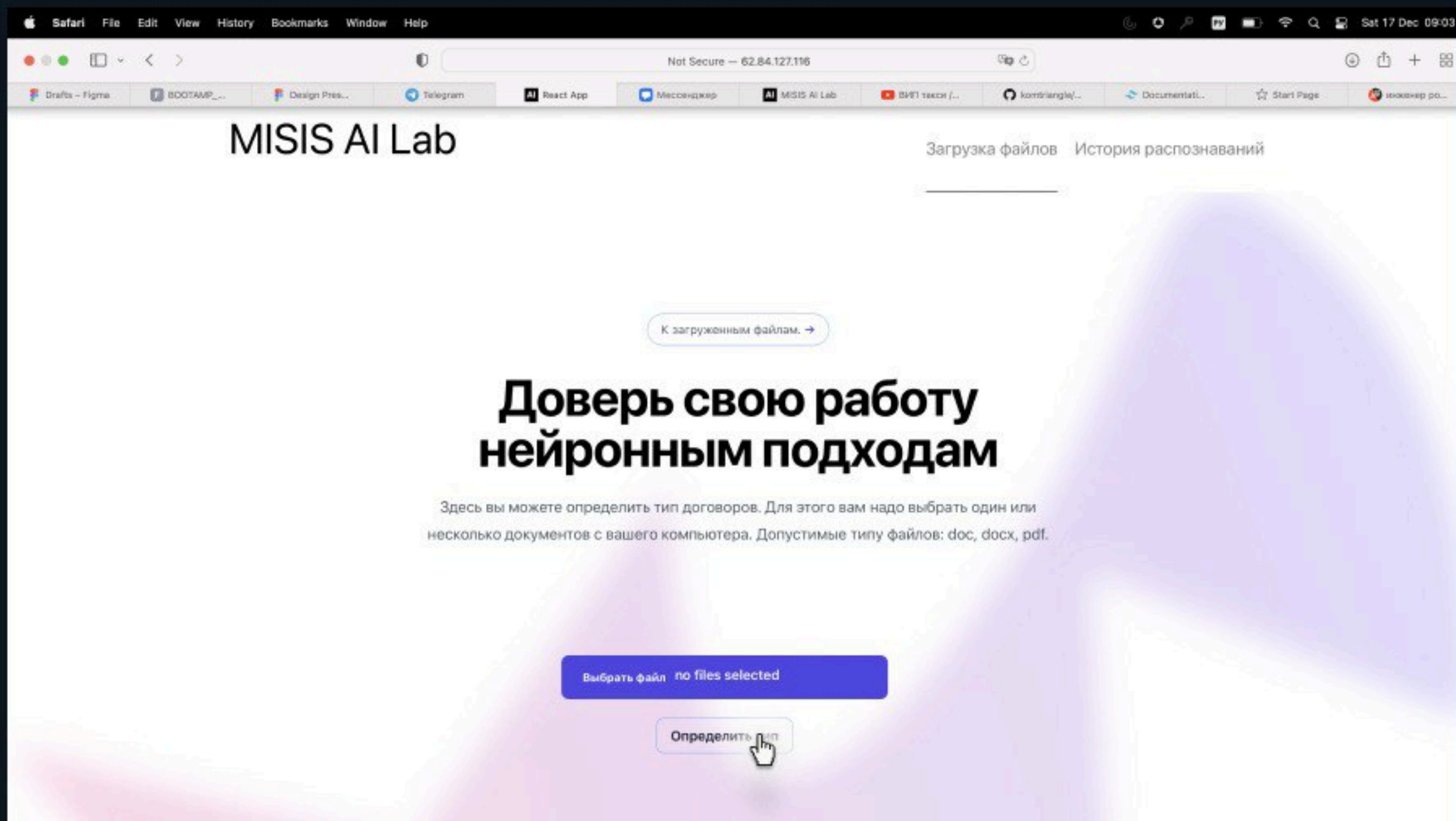
Отправляем вам необходимые документы: 1dogovorarendy\_kvartiry.pdf

Если документ получен ошибочно, перейдите по ссылке [source].  
С уважением, команда MISIS AI Lab!





Redux



# THE END

