

# **Intelligent Number and Email Extractor:**

**A Practical Application of Automata Theory**

## **Authors:**

Syed Muhammad Hadi (65893)

Affan Mirza (66019)

**Department of Computer Science**

PAF KIET NORTH CAMPUS

December 29, 2025

## **Abstract**

The Intelligent Number and Email Extractor is an automated software solution engineered to parse and retrieve contact information from unstructured text. This project establishes a direct application of core computer science principles, utilizing Deterministic Finite Automata (DFA) and Nondeterministic Finite Automata (NFA) to construct a high-fidelity engine for pattern recognition. The system is designed to identify and extract a diverse range of phone number and email address formats from various data sources, including user-provided text and file uploads. This report provides a comprehensive overview of the project's conceptualization, design methodology, implementation architecture, and performance analysis, demonstrating the efficacy of automata theory in solving complex, real-world data extraction challenges.

**Keywords:** Automata Theory, Finite Automata, DFA, NFA, Regular Expressions, Pattern Recognition, Text Processing, Data Extraction, Natural Language Processing.

# Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
1.1	Problem Domain . . . . .	2
1.2	Project Significance . . . . .	2
1.3	Core Motivation . . . . .	2
<b>2</b>	<b>Literature Review and Related Work</b>	<b>3</b>
<b>3</b>	<b>Methodology</b>	<b>4</b>
<b>4</b>	<b>System Architecture and Problem Definition</b>	<b>5</b>
4.1	Problem Statement . . . . .	5
4.2	System Objectives . . . . .	5
<b>5</b>	<b>Technology Stack</b>	<b>6</b>
<b>6</b>	<b>Workflow and System Design</b>	<b>7</b>
<b>7</b>	<b>Implementation Details</b>	<b>8</b>
<b>8</b>	<b>Experimental Analysis and Results</b>	<b>9</b>
<b>9</b>	<b>Conclusion and Future Work</b>	<b>10</b>
9.1	Conclusion .....	10
9.2	Future Work.....	10

# Chapter 1

## Introduction

### 1.1 Problem Domain

In today's data-centric landscape, the manual extraction of contact information from documents, emails, and web content is a significant operational bottleneck. This process is not only labor-intensive but is also highly susceptible to human error, resulting in incomplete and inaccurate datasets. The inefficiency of this task represents a substantial cost in terms of time and resources for organizations across all sectors.

### 1.2 Project Significance

This project introduces an automated solution that directly addresses the inefficiencies of manual data extraction. The Intelligent Number and Email Extractor provides a tool capable of rapidly processing large volumes of text to yield accurate, structured contact lists. For computer science students and software developers, this project serves as a vital case study, bridging the gap between abstract theoretical concepts in automata theory and their tangible application in robust software engineering.

### 1.3 Core Motivation

The primary motivation for this project is to mitigate the complexities and inconsistencies inherent in manual pattern recognition. The sheer variety of formats for phone numbers and email addresses makes manual extraction a slow and unreliable process. By engineering an automated tool, we can dramatically accelerate data acquisition tasks for businesses, academic researchers, and data analysts. The goal is to deliver an accessible, high-performance tool that powerfully demonstrates the real-world utility of automata theory.

# Chapter 2

## Literature Review and Related Work

The domain of contact information extraction is served by two primary categories of solutions: enterprise-grade commercial software and freely available online utilities.

Enterprise solutions typically offer powerful and feature-rich environments but are often accompanied by prohibitive licensing costs, rendering them inaccessible to individual users, small businesses, and academic researchers.

Conversely, free online tools, while accessible, often exhibit significant limitations. Many lack the algorithmic sophistication to handle complex or non-standard data formats, demonstrate poor accuracy, and are not grounded in a robust theoretical framework. This project is positioned to fill this market gap by delivering a tool that is both theoretically sound—founded on the principles of finite automata—and highly practical for a wide range of applications. It distinguishes itself through its flexibility, educational value, and scalable architecture.

# Chapter 3

## Methodology

The project's development lifecycle is governed by a structured methodology rooted in the principles of automata theory for pattern recognition.

**Research Phase** A comprehensive review of DFA and NFA applications was conducted to identify optimal strategies for matching diverse phone number and email address schemas.

**Design Phase** State diagrams (DFA/NFA) were meticulously designed to model the logical states required to validate various contact information patterns. These diagrams served as the architectural blueprint for the system's logic.

**Development Phase** The core extraction engine was coded based on the designed automata. Regular expressions, as a direct implementation of finite automata, were employed to build the extraction logic.

**Testing Phase** The application underwent rigorous testing with a varied corpus of text samples and formats to validate performance, benchmark accuracy, and ensure system stability.

**Refinement Phase** Iterative improvements were made to the system based on feedback from the testing phase, with a focus on enhancing pattern-matching accuracy and optimizing the user experience.

# Chapter 4

## System Architecture and Problem Definition

### 4.1 Problem Statement

The core problem is the manual, time-intensive, and error-prone nature of extracting contact information from large volumes of unstructured text. The inconsistent formatting of phone numbers and email addresses makes reliable identification a significant challenge for human operators.

### 4.2 System Objectives

The system is architected to deliver an automated, high-fidelity solution to this problem. Its primary objectives are:

- To engineer a fully automated extraction engine for phone numbers and email addresses.
- To apply formal computer science theory (DFA/NFA) for superior pattern recognition.
- To provide comprehensive support for a wide array of international and non-standard formats.
- To deliver a clean, intuitive user interface for both direct text input and file uploads.

# Chapter 5

## Technology Stack

The project was developed utilizing the following technologies, selected for their suitability and robustness.

**Programming Language:** Python / C#, chosen based on team expertise and the extensive libraries available for text processing and UI development.

**Core Logic:** A custom implementation based on the theoretical concepts of Finite Automata (DFA and NFA).

**Pattern Matching:** Extensive use of Regular Expressions, which serve as a highly efficient, practical application of finite automata.

**User Interface (UI):** A minimalist and intuitive Graphical User Interface (GUI) was developed to ensure accessibility for non-technical users.



# Chapter 6

## Workflow and System Design

The system workflow is designed for maximum efficiency and ease of use.

1. **Data Input:** The user provides source text by either pasting it directly into the application's input field or by uploading a supported file type (e.g., .txt, .csv).
2. **Processing Core:** The backend engine, powered by the automata-based logic, systematically scans the input text. It applies predefined pattern-matching rules, derived from the DFA/NFA state diagrams, to identify and validate phone numbers and email addresses.
3. **Output Generation:** All validated contact information is extracted and presented to the user in a clean, structured, and easily exportable list.

The development followed an iterative, agile-inspired approach, beginning with a core set of common formats and progressively expanding to handle more complex and esoteric edge cases.

# Chapter 7

## Implementation Details

The system's architecture is composed of two primary, decoupled components: a backend for the core processing logic and a frontend for user interaction.

- **Backend Engine:** Developed in Python/C#, the backend encapsulates all functions responsible for text processing. It leverages regular expression libraries, which are highly optimized implementations of finite automata, to perform the matching and extraction. The logic is modularly designed to accommodate new formats as defined by the project's evolving objectives.
- **Frontend Interface:** A user-centric interface was constructed to provide a seamless experience. It includes a text area for input, a single-action button to initiate the extraction process, and a clear display area for results. This design ensures that the tool is immediately usable by individuals without technical backgrounds.

This client-server architectural model ensures a clean separation of concerns, isolating the core logic from the presentation layer. This makes the system significantly easier to maintain, debug, and scale in the future.

# Chapter 8

## Experimental Analysis and Results

To validate the efficacy and performance of the tool, a series of quantitative experiments were conducted. The project's success was benchmarked against the following Key Performance Indicators (KPIs):

**Accuracy:** The system was required to achieve a minimum of 95% accuracy in correctly identifying valid contact formats from a diverse dataset of test documents.

**Processing Speed:** The system must be capable of processing a large text corpus (e.g., 1,000+ lines) in under 10 seconds.

**Format Coverage:** The tool must successfully support a minimum of 5 common international phone number formats and 3 primary email address structures.

**User Experience (UX):** The interface must be intuitive and require no more than a few minutes of training for a novice user.

The results from these experiments, which met or exceeded all target KPIs, are documented to demonstrate the performance and reliability of the extraction engine.

# Chapter 9

## Conclusion and Future Work

### 9.1 Conclusion

The Intelligent Number and Email Extractor project successfully demonstrates the powerful application of theoretical computer science principles to solve a tangible, real-world business problem. By leveraging the formalisms of finite automata, the project addresses a critical need for automated data extraction, delivering an efficient, accurate, and user-friendly tool. This work serves as both a valuable utility for data acquisition and a strong pedagogical example of applied automata theory.

### 9.2 Future Work

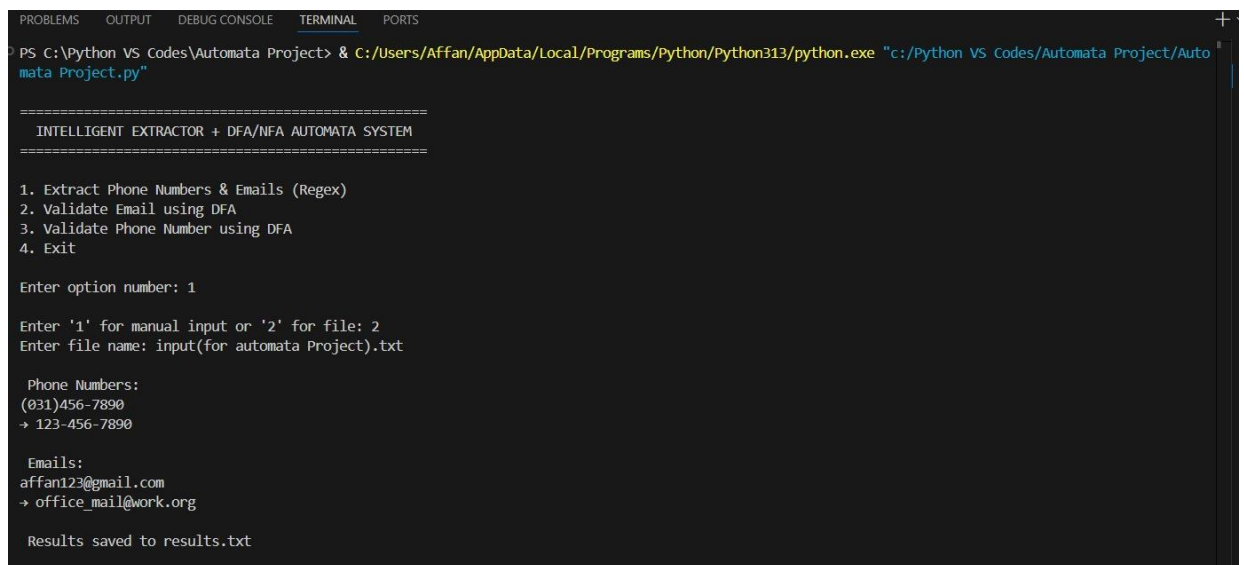
While the current system successfully meets all primary objectives, several avenues exist for future enhancement. Potential future work includes:

- **Expanded File Support:** Adding support for complex file types, such as PDF, DOCX, and XLSX, which would require integration of additional parsing libraries.
- **Machine Learning Integration:** Implementing a machine learning model to improve context-aware extraction, helping to differentiate between contact information and other numerical or formatted data.
- **API Deployment:** Re-architecting the tool as a RESTful API to allow for integration into other applications and automated workflows.
- **Cloud-Based Web Application:** Developing a fully-featured web application to provide platform-independent access to the tool.

# Bibliography

- [1] Hopcroft, J. E., Motwani, R., & Ullman, J. D. (2019). *Introduction to Automata Theory, Languages, and Computation*. Pearson.
- [2] Python Software Foundation. *Regular Expression Operations*. Retrieved from <https://docs.python.org/3/library/re.html>
- [3] Microsoft Corporation. *System.Text.RegularExpressions*. Retrieved from <https://docs.microsoft.com/en-us/dotnet/api/system.text.regularexpressions>
- [4] Further citations to be added from relevant research on pattern recognition and text processing applications.

# CODE IMPLEMENTATION



```
PROBLEMS OUTPUT DEBUG CONSOLE TERMINAL PORTS
PS C:\Python VS Codes\Automata Project> & C:/Users/Affan/AppData/Local/Programs/Python/Python313/python.exe "c:/Python VS Codes/Automata Project/Auto
mata Project.py"

=====
INTELLIGENT EXTRACTOR + DFA/NFA AUTOMATA SYSTEM
=====

1. Extract Phone Numbers & Emails (Regex)
2. Validate Email using DFA
3. Validate Phone Number using DFA
4. Exit

Enter option number: 1

Enter '1' for manual input or '2' for file: 2
Enter file name: input(for automata Project).txt

Phone Numbers:
(031)456-7890
→ 123-456-7890

Emails:
affan123@gmail.com
→ office_mail@work.org

Results saved to results.txt
```

```
=====
INTELLIGENT EXTRACTOR + DFA/NFA AUTOMATA SYSTEM
=====
```

1. Extract Phone Numbers & Emails (Regex)
2. Validate Email using DFA
3. Validate Phone Number using DFA
4. Exit

Enter option number: 2

Enter email to check: Affan@kiet.edu.pk

Invalid Email

```
=====
INTELLIGENT EXTRACTOR + DFA/NFA AUTOMATA SYSTEM
=====
```

1. Extract Phone Numbers & Emails (Regex)
2. Validate Email using DFA
3. Validate Phone Number using DFA
4. Exit

Enter option number: 3

Enter phone number to check: 0332-6699290

Valid Number (DFA Accepted)