# Web Scraping Project Report

**Title:** Extracting Book Details: A Comprehensive Book Data Collection Approach

**Objective:** To develop an automated web scraping solution for collecting structured book information from online platforms, focusing on extracting critical book metadata.

**Methodology:**

1. Web Scraping Technique: Custom Python script using BeautifulSoup and Requests libraries

2. Data Extraction Strategy: Systematic collection of book details from multiple online sources

3. Data Storage: Multiple format outputs (TXT, CSV, JSON)

**Websites Investigated:**

1. **Goodreads([https://www.goodreads.com/list/show/1.Best_Books_Ever](https://www.goodreads.com/list/show/1.Best_Books_Ever))**

   ➢ Platform: Comprehensive book listing website

   ➢ Data Potential: Book titles, authors, ratings

   ➢ Scraping Complexity: Moderate

2. **Library Thing(https://www.librarything.com/)**

   ➢ Platform: Advanced book cataloguing website

   ➢ Data Potential: Diverse book metadata

   ➢ Scraping Complexity: Advanced

**Technical Implementation:**

1. **Programming Language:** Python

2. **Key Libraries:**

   ➢ Requests (HTTP requests)

   ➢ BeautifulSoup (HTML parsing)

   ➢ JSON (Data serialization)

   ➢ CSV (Structured data storage)

**Data Extraction Capabilities:**

1. **Book Titles**

2. **Author Names**

3. **Potential for expanding to:**

   ➢ Publication Years

   ➢ Ratings

   ➢ Genre Information

**Challenges Addressed:**

- Website structure variations
- Anti-scraping mechanisms
- Data consistency
- Ethical scraping practices

**Potential Applications:**

- Literary research
- Book recommendation systems
- Reading trend analysis
- Academic and commercial book studies

**Team Composition:**

1. Vansh Shrivastava (2301EC34)
2. Riya Singh (2301PH25)
3. Shivani (2301PH28)
4. Utpal Raj Ambastha (2301PH24)

**Conclusion:** A versatile web scraping solution capable of extracting structured book information, with scalability for future enhancements.