

# 1 Задание

Зотов Антон 11-106

- 1. Скачать минимум 100 текстовых страниц с помощью краулера.
- 2. Записать каждую страницу (без html-кода) в отдельный текстовый файл.
- 3. Создать файл index.txt, в котором хранится номер документа и ссылка на страницу.

```
PS C:\Users\1etmehear\RiderProjects\NAZVANIE1\NAZVANIE1> dotnet run -- https://render.ru/xen/ https://www.gameprogrammingpatterns.com/contents.html
1665
Сохранена страница 1: https://render.ru/xen/
53
Страница https://www.gameprogrammingpatterns.com/contents.html содержит менее 1000 слов, пропускаем
1252
Сохранена страница 2: http://render.ru/
1567
Сохранена страница 3: https://render.ru/xen/whats-new/posts/
540
Страница https://render.ru/xen/search/?type=post содержит менее 1000 слов, пропускаем
290
Страница https://render.ru/xen/whats-new/ содержит менее 1000 слов, пропускаем
36
Страница https://render.ru/xen/whats-new/profile-posts/ содержит менее 1000 слов, пропускаем
400
Страница https://render.ru/xen/whats-new/latest-activity содержит менее 1000 слов, пропускаем
4
Страница https://render.ru/ru/artist/about содержит менее 1000 слов, пропускаем
```

Рис.1 – Вызов программы в терминале

|  |              |                 |                    |       |
|--|--------------|-----------------|--------------------|-------|
|  | page_94.txt  | 29.03.2025 0:21 | Текстовый докум... | 16 КБ |
|  | page_95.txt  | 29.03.2025 0:21 | Текстовый докум... | 14 КБ |
|  | page_96.txt  | 29.03.2025 0:21 | Текстовый докум... | 17 КБ |
|  | page_97.txt  | 29.03.2025 0:21 | Текстовый докум... | 16 КБ |
|  | page_98.txt  | 29.03.2025 0:21 | Текстовый докум... | 16 КБ |
|  | page_99.txt  | 29.03.2025 0:22 | Текстовый докум... | 37 КБ |
|  | page_100.txt | 29.03.2025 0:23 | Текстовый докум... | 69 КБ |

Рис. 2 – Сохраненные страницы

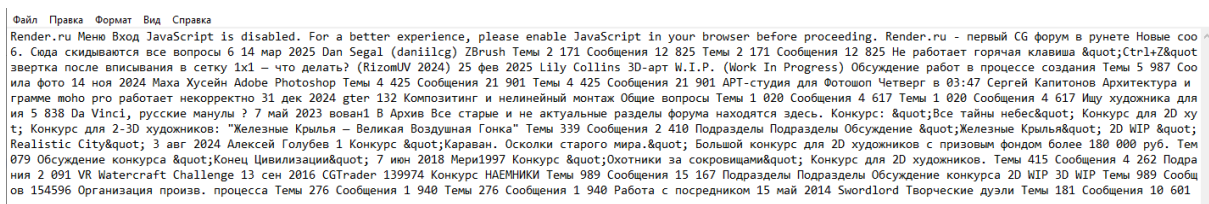


Рис. 3 – Пример страницы

```
88 https://render.ru/xen/forums/3d-wip.104/
89 https://render.ru/xen/forums/wip-konkursa.151/
90 https://render.ru/xen/forums/organizacija-proizv-processa.75/
91 https://render.ru/xen/posts/1111135/
92 https://render.ru/xen/forums/tvorcheskie-duehli.95/
93 https://render.ru/xen/forums/animacionnye-konkursy.83/
94 https://render.ru/xen/forums/wip-animacija.84/
95 https://render.ru/xen/forums/skulpting-na-skorost.86/
96 https://render.ru/xen/#top
97 http://render.ru/section
98 http://render.ru/section/8
99 https://render.ru/ru/news/post/26556
100 https://render.ru/ru/LessiDance/post/26566
```

Стр 101, стр 6 1

Рис. 4 – index.txt

```
1 usage
class WebCrawler
{
    private readonly List<string> _startUrls;
    private readonly int _minPages;
    private readonly int _minWords;
    private readonly HashSet<string> _visitedUrls = new ();
    private int _pagesDownloaded = 0;
    private const string IndexFile = "index.txt";
    private const string OutputDir = "pages";

2 usage
    public WebCrawler(List<string> startUrls, int minPages = 100, int minWords = 1000)
    {
        _startUrls = startUrls;
        _minPages = minPages;
        _minWords = minWords;

        Directory.CreateDirectory(path: OutputDir);

        File.WriteAllText(IndexFile, contents: string.Empty);
    }

3 usage
    private bool IsValidUrl(string url)
    {
        // Проверяем, что URL валиден и не является файлом (например, .pdf)
        if (!Uri.TryCreate(url, UriKind.Absolute, out var uri))
            return false;

        if (uri.Scheme != Uri.UriSchemeHttp && uri.Scheme != Uri.UriSchemeHttps)
            return false;

        var invalidExtensions = new[] { ".pdf", ".jpg", ".png", ".gif", ".zip", ".doc", ".docx" };
        return !invalidExtensions.Any(ext: string => uri.AbsolutePath.EndsWith(ext, StringComparison.OrdinalIgnoreCase));
    }
}
```

Рис. 5 – Конструктор краулера + проверка валидных ссылок

```

1 usage
private async Task<HashSet<string>> GetLinksAsync(string url, HttpClient client)
{
    try
    {
        var response = await client.GetAsync(url);
        response.EnsureSuccessStatusCode();

        var html:string = await response.Content.ReadAsStringAsync();
        var doc = new HtmlDocument();
        doc.LoadHtml(html);

        var links = new HashSet<string>();

        foreach (var link:HtmlNode in doc.DocumentNode.SelectNodes(xpath: "//a[@href]") ?? Enumerable.Empty<HtmlNode>())
        {
            var href:string = link.GetAttributeValue(name: "href", def: string.Empty);
            var absoluteUrl:string = new Uri(new Uri(url), href).AbsoluteUri;

            if (IsValidUrl(absoluteUrl))
                links.Add(absoluteUrl);
        }

        return links;
    }
    catch (Exception e)
    {
        Console.WriteLine($"Ошибка при получении ссылок с {url}: {e.Message}");
        return new HashSet<string>();
    }
}

```

Рис. 6 – Получение ссылок из страницы. Извлекает текстовое содержимое из HTML. Удаляет ненужные элементы (скрипты, стили, навигацию и т.д.). Очищает текст от лишних пробелов и переносов строк

```

private string ExtractText(string html)
{
    var doc = new HtmlDocument();
    doc.LoadHtml(html);

    var nodesToRemove :HtmlNodeCollection = doc.DocumentNode.SelectNodes(xpath: "//*[script|style|nav|footer|head|iframe]");

    foreach (var node in nodesToRemove)
        node.Remove();

    var text = doc.DocumentNode.InnerText;

    // Удаляем лишние пробелы и переносы строк
    text = Regex.Replace(input: text, pattern: @"\s+", replacement: " ").Trim();

    return text;
}

1 usage
private bool SavePage(string url, string text)
{
    var words :string[] = text.Split(new[] { ' ', '\t', '\n' }, StringSplitOptions.RemoveEmptyEntries);

    // foreach(var word in words)
    Console.WriteLine(words.Length);

    var wordCount = words.Length;
    if (wordCount < _minWords)
    {
        Console.WriteLine($"Страница {url} содержит менее {_minWords} слов");
        return false;
    }

    _pagesDownloaded++;
    var filename = Path.Combine(OutputDir, $"page_{_pagesDownloaded}.txt");

    File.WriteAllText(filename, contents: text, Encoding.UTF8);

    File.AppendAllText(path: IndexFile, contents: $"{_pagesDownloaded}\t{url}{Environment.NewLine}", Encoding.UTF8);

    Console.WriteLine($"Страница сохранена {_pagesDownloaded}: {url}");
    return true;
}

```

Рис. 7 – Метод для парсинга страницы + сохранение страницы в txt файл. Проверяет количество слов на странице. Если слов достаточно, сохраняет текст в файл

```
1 usage
private async Task<bool> DownloadPageAsync(string url, HttpClient client)
{
    if (_visitedUrls.Contains(url) || _pagesDownloaded >= _minPages)
        return false;

    _visitedUrls.Add(url);

    try
    {
        var response = await client.GetAsync(url);
        response.EnsureSuccessStatusCode();

        var html:string = await response.Content.ReadAsStringAsync();
        var text = ExtractText(html);

        if (string.IsNullOrEmpty(text))
            return false;

        return SavePage(url, text);
    }
    catch (Exception e)
    {
        Console.WriteLine($"Ошибка при скачивании {url}: {e.Message}");
        return false;
    }
}
```

Рис. 8 – Метод для скачивания страницы. Проверяет, не была ли страница уже загружена. Асинхронно скачивает страницу. Извлекает текст и сохраняет его.

```

1 usage
public async Task CrawlAsync()
{
    var handler = new HttpClientHandler
    {
        AutomaticDecompression = DecompressionMethods.GZip | DecompressionMethods.Deflate
    };

    using var client = new HttpClient(handler);
    client.DefaultRequestHeaders.UserAgent.ParseAdd(input: "Mozilla/5.0");
    client.Timeout = TimeSpan.FromSeconds(10);

    var currentLevelUrls = new HashSet<string>(_startUrls);
    var nextLevelUrls = new HashSet<string>();

    while (_pagesDownloaded < _minPages && currentLevelUrls.Count > 0)
    {
        foreach (var url:string in currentLevelUrls)
        {
            if (_pagesDownloaded >= _minPages)
                break;

            if (await DownloadPageAsync(url, client))
            {
                var links:HashSet<string> = await GetLinksAsync(url, client);
                foreach (var link:string in links)
                    nextLevelUrls.Add(link);
            }

            // Небольшая задержка чтобы не перегружать сервер
            await Task.Delay(1000);
        }

        currentLevelUrls = new HashSet<string>(collection: nextLevelUrls.Except(_visitedUrls));
        nextLevelUrls.Clear();
    }

    Console.WriteLine($"Завершено. Скачано {_pagesDownloaded} страниц.");
}

```

Рис. 9 – Основной метод краулера. Имитирует работу браузера, настраивает распаковку сжатых ответов, в цикле с самого начала проходится по стартовым ссылкам. Если скачивание успешно, извлекает все ссылки, добавляет их в основной хэшсет