



Senior Project, 2018, Fall Web Page Archiving and Content Analysis

Student: Adam Tahoun, Florida International University
Mentor: Dr. Mark Finlayson, Florida International University
Professor: Dr. Masoud Sadjadi, Florida International University



Problem

- Temporary web content
- Filtering and parsing is time consuming and difficult
- Content may need translation
- Content not available for offline access and manipulation
- No solution exists

Solution

- Custom file format
- Improved GUI for viewing and interfacing with custom file format
- Introduced deep learning translation API
- Introduced machine learning approach to content extraction
- Cross-platform

Implementation

- Python2.7 and wget for the downloading CLI
- Electron, HTML5, CSS3, and JS for the GUI
- Dragnet for the content extraction
- PyTorch and OpenNMT for the translation

Screenshots

The screenshots show the application's user interface. The first screenshot, 'Download Screen', displays a web form for downloading content with fields for 'Input URL', 'Output Directory', and 'Number of threads'. The second screenshot, 'Original Article', shows a news article in Russian with a title 'Кладбище' and a sub-header 'Стерляда в колледже в Керчи: опубликовано видео с камер наблюдения'. The third screenshot, 'Translated Article', shows the same article translated into English, with a title 'Cemetery' and a sub-header 'Shooting at College in Kerch: published a video with surveillance camera'.

Current System

- Previous projects were incomplete
- Existing partial solutions have been deprecated

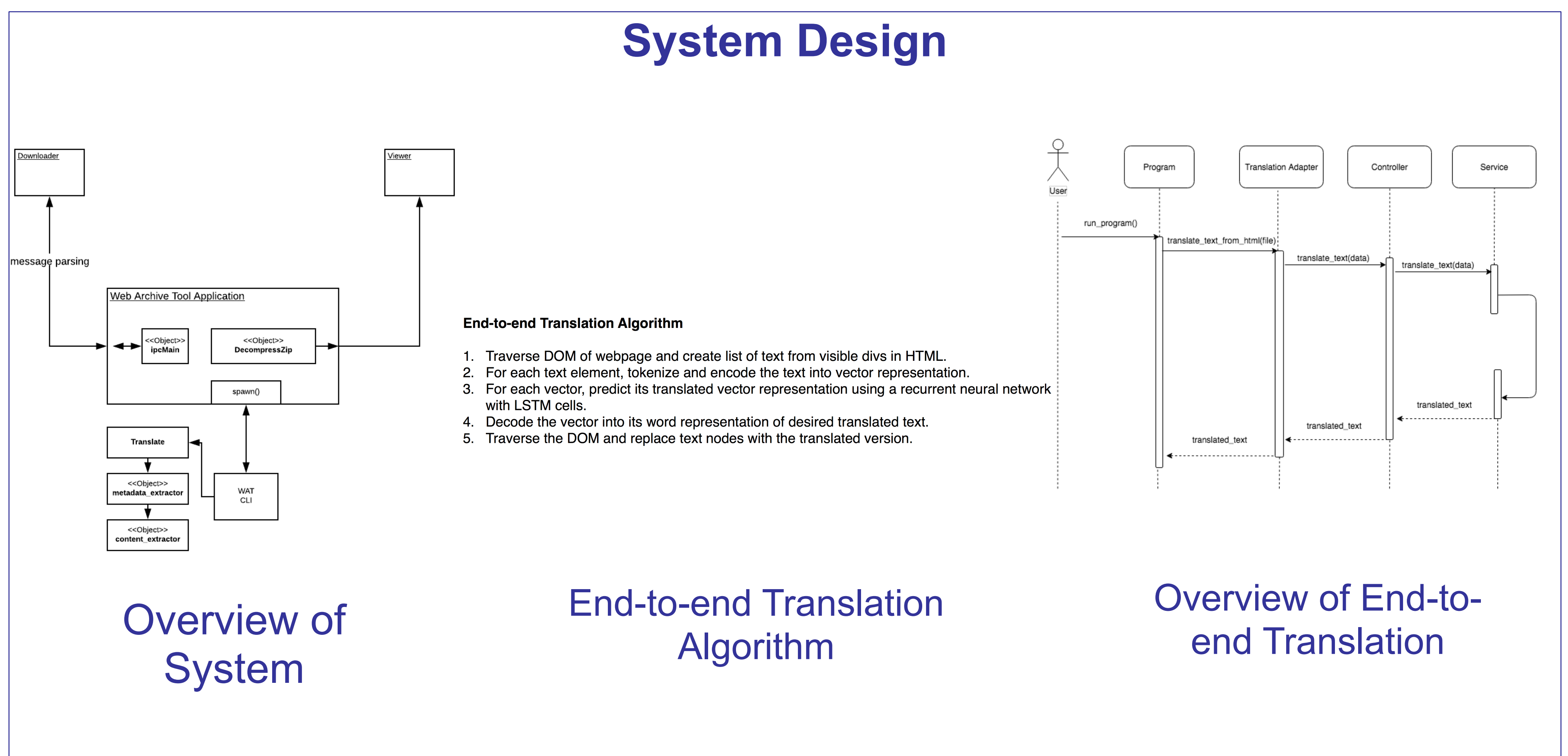
Requirements

- Batch download large sets of pages
- Encapsulate sets of pages in a session
- Continue session if download was interrupted
- Archive web page in format that allows easy programmatic access
- Encapsulate this archive in a single file
- Ability to easily view web archive
- Identify, translate, and extract the main textual content of the article
- Identify and extract the metadata of the article
- Cross-platform

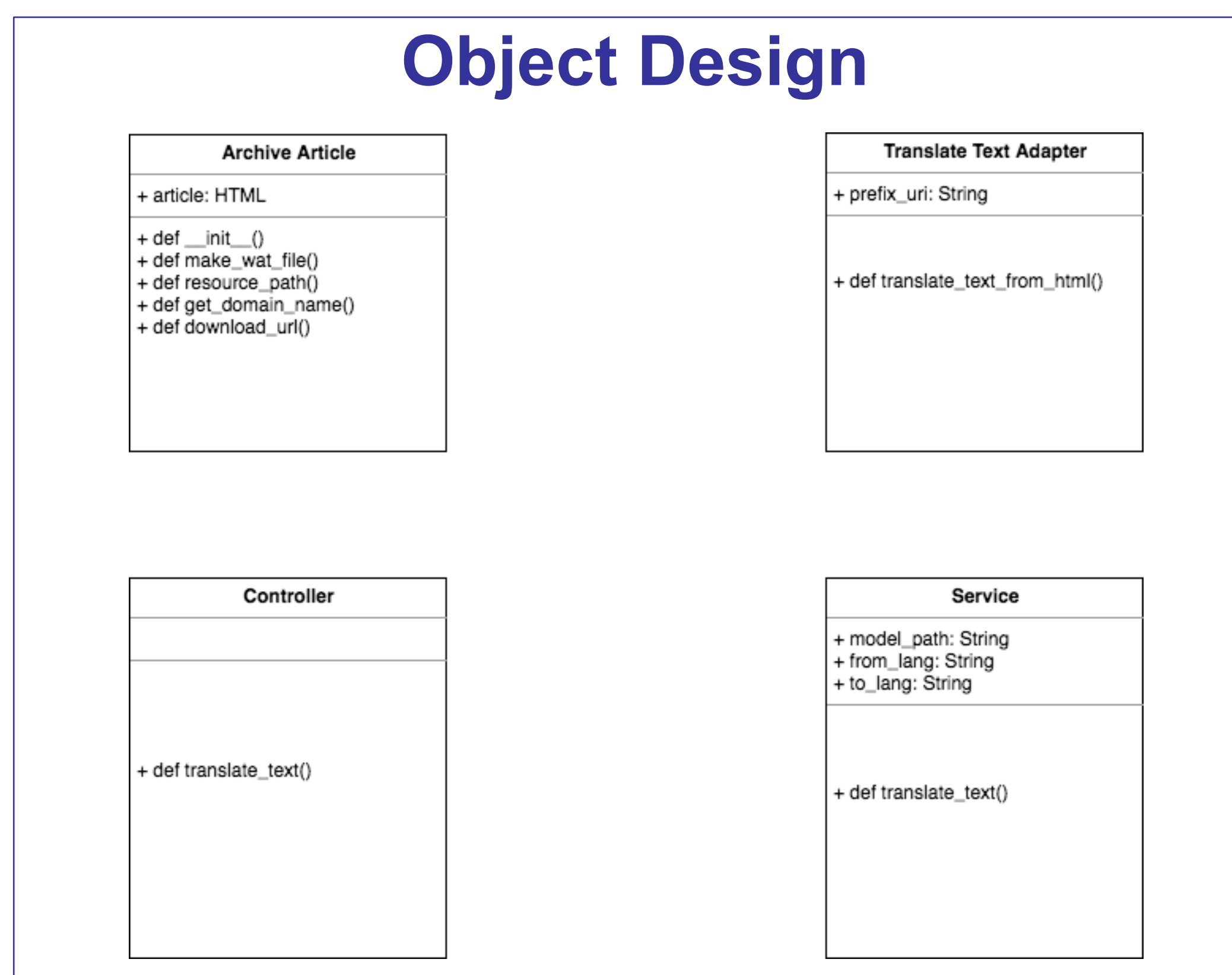
Verification

- Unit and integration testing
- Tested on 5GB of articles
- Tested on macOS, Windows, and Linux
- Content Extraction: 90% accuracy
- In-house Translation: 70% accuracy

System Design



Object Design



Summary

- Our work allows users to:
- Download and archive batches of web pages
 - Translate and extract important information from the web pages
 - View these web pages on any platform

Acknowledgment

The material presented in this poster is based upon the work supported by Mark Finlayson, Andres Cremisini, and Pedro Torres-Carrasquillo. I am thankful for the help that I received from my group member, Oscar Martinez.