# Chinese Word Segmentation

- Author: Otis B.C. Chung

## Feature

- This project aims to segment Chinese words.

- The program will download dictionary automatically.

- The default sentence is "這裡是MI2S實驗室,位於資訊工程學系的65802室", you can use the command-line arguments to pass your own sentence.
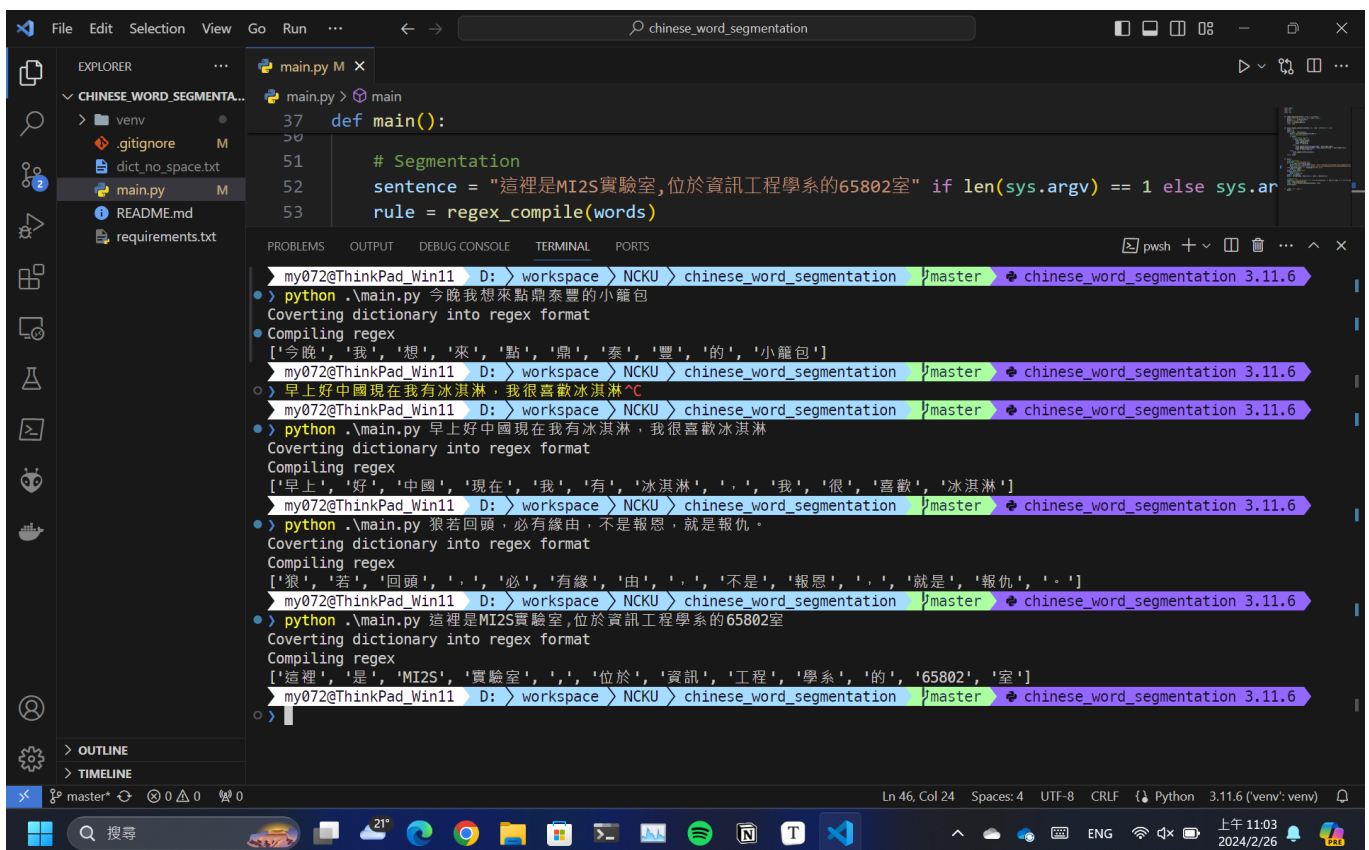
## Usage

By default, you can use

```
python main.py
```

You can pass your sentence into program using command-line arguments

```
python main.py "狼若回頭，必有緣由，不是報恩，就是報仇。"
```

## Result