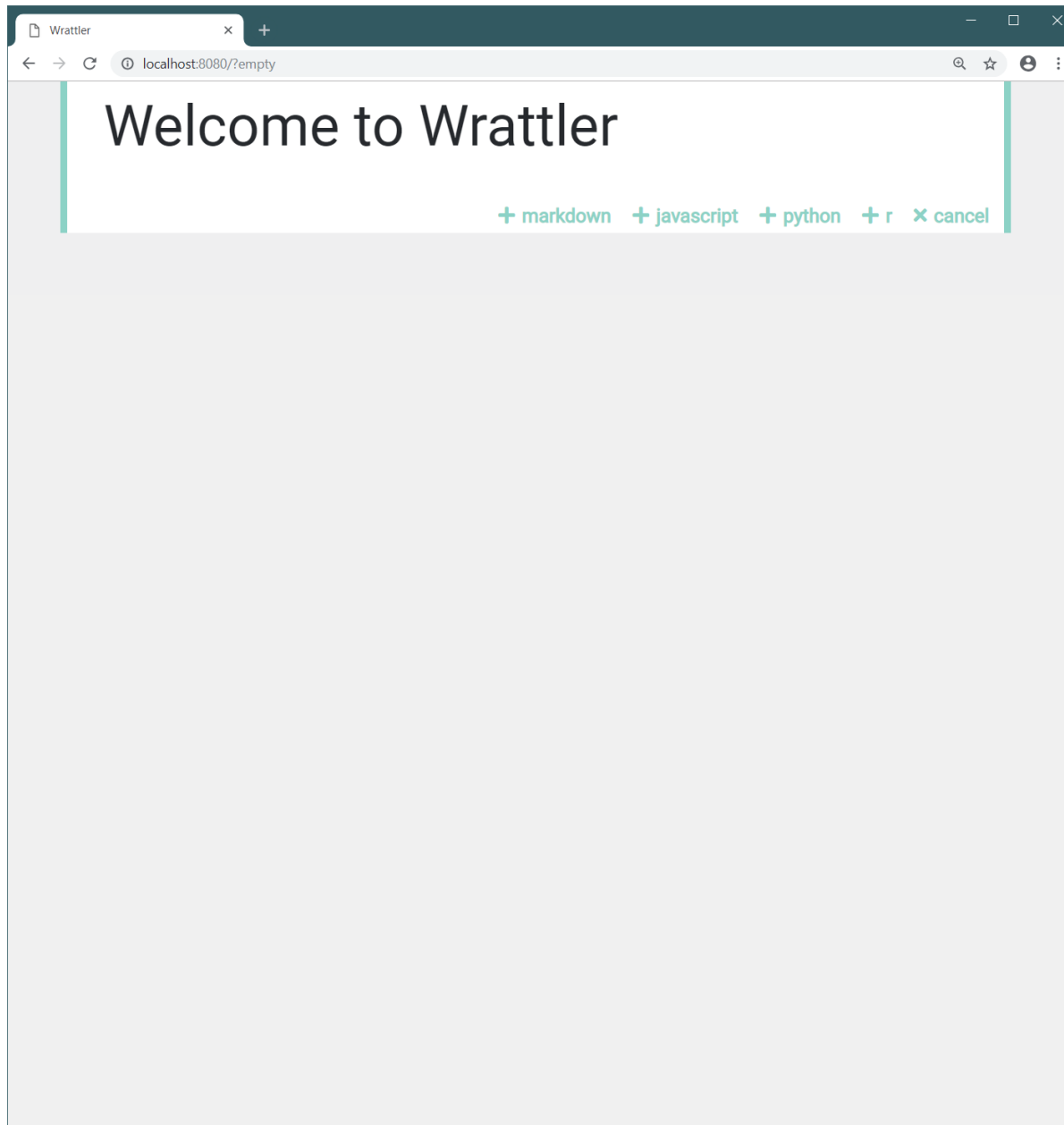# **Wrattler:** Polyglot, reproducible and smart notebooks
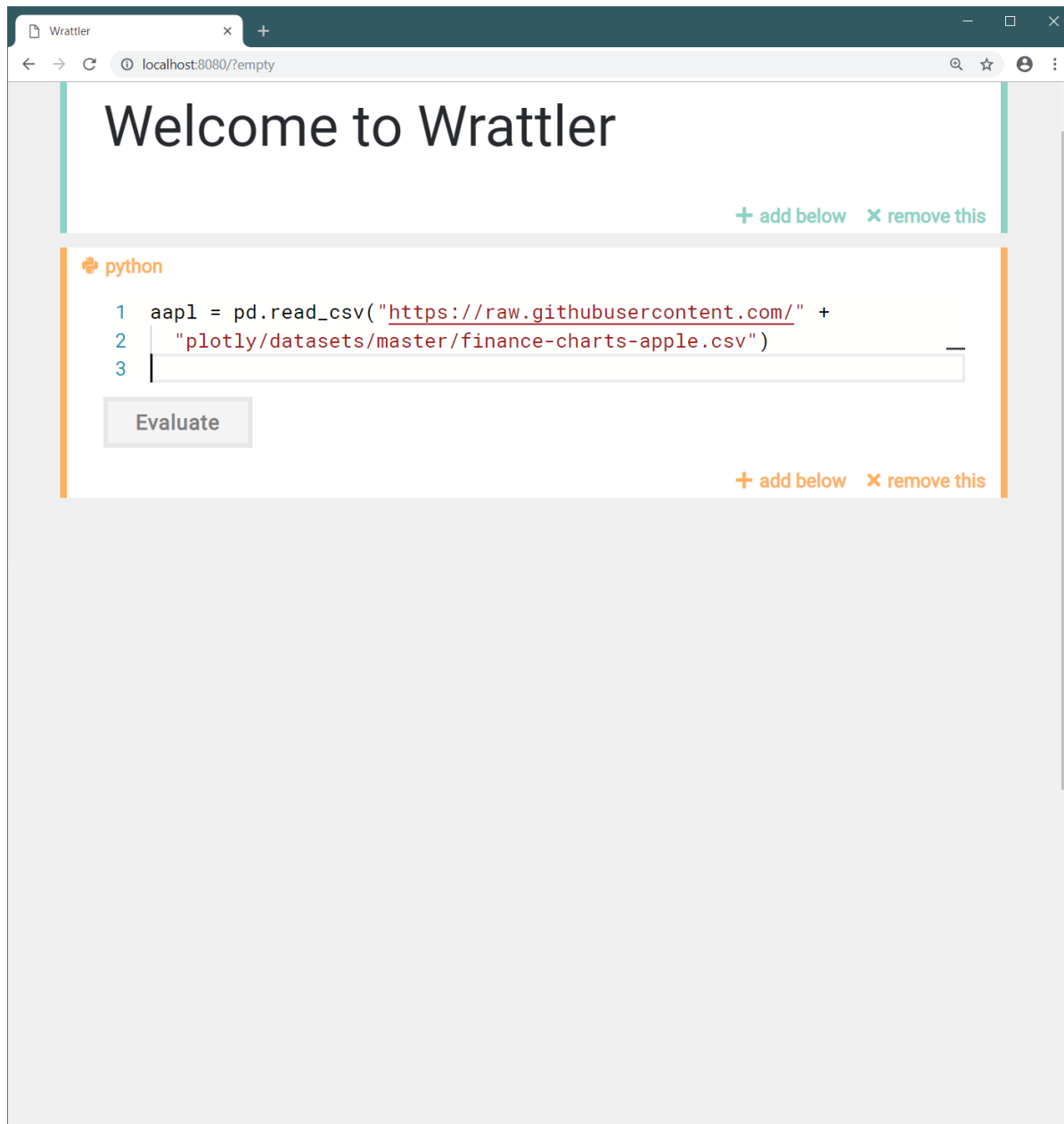
Demonstration of the developed prototype system (GDI 1.6b)

**Polyglot**
Passing data from Python to JavaScript

We start with an empty notebook. Wrattler allows us to add with **Markdown** comments and cells with **JavaScript, R** and **Python** code.

We use the Python **pandas** library to easily load CSV file from an online source.

After updating the cell and clicking the **Evaluate** button, we see a preview of the downloaded data.

Next, we add a **JavaScript** cell with source code that uses the **Plotly** library to build a visualization.

When we evaluate the added code, Wrattler runs our **JavaScript** directly **in the web browser** and shows the visualization.

# Reproducible
Recomputation using provenance

Consider a sample notebook with two cells that define data frames **one** and **two** and a third cell that concatenates the two and prints the result.

We can evaluate cells one-by-one by clicking on **Evaluate.** Here, we evaluated just the first (Python) cell.

But you **do not have to** run cells one-by-one. If we ask to evaluate the last cell, Wrattler automatically **runs all dependencies.**

Here, we evaluated the last cell and the second was evaluated automatically.

If we modify a cell, Wrattler updates the **dependency graph** it maintains. Results of all cells that depend on a modified cell are removed and need to be recomputed.

Re-evaluating the last cell also evaluates all cells that it depends on, using the new **dependency graph.**

```
python
1  one = pd.DataFrame({"name":["Jim"], "age":[51]})
```

one

| age | name |
| --- | --- |
| 51 | Jim |

**+ add below   ✕ remove this**

```
r
1  two <- data.frame(name=c("Jane"), age=c(54))
```

two

| age | name |
| --- | --- |
| 54 | Jane |

**+ add below   ✕ remove this**

```
javascript
1  addOutput(function(id) {
2    var items = one.concat(two).map(function(row) {
3      return "<li><b>" + row.name + "</b> (" + row.age + ")</li>" });
4    var html = "<ul style='margin:20px'>" + items.join("") + "</ul>";
5    document.getElementById(id).innerHTML = html;
6  });
```

output0

- **Jim** (51)
- **Jane** (54)

**+ add below   ✕ remove this**

If we revert a change back, we do not have to re-evaluate and we see **earlier outputs immediately.**

Wrattler caches past dependency graph nodes and reuses a past node that has already been evaluated.

**Smart**
Support for data cleaning tools

Wrattler supports other data wrangling and cleaning tools built as part of the **AI for Data Analytics** project such as **datadiff**. As an example, we look at the UK broadband quality data set.

```r
library(broadband)

bb2013 <- broadband2013
bb2014 <- broadband2014

stripnn <- function(x) { gsub("[^0-9\\.]", "", x) }

bb2013[["ID"]][bb2013[["ID"]] == "FTTC"] <- NA
bb2013$ID <- as.integer(bb2013$ID)
bb2013[["Headline.speed"]] <-
  as.integer(stripnn(bb2013[["Headline.speed"]]))
bb2013[["Packet.loss....24.hour"]] <-
  as.numeric(stripnn(bb2013[["Packet.loss....24.hour"]]))
bb2013[["Packet.loss....8.10pm.weekday"]] <-
  as.numeric(stripnn(bb2013[["Packet.loss....8.10pm.weekday"]]))
```

**bb2013**   bb2014

| DNS.failure....24.hour | DNS.failure....8.10pm.weekday | DNS.resolution..ms.24.hour | DNS.res |
|---|---|---|---|
| 0 | 0 | 19.139 | 18.578 |
| 0.001 | 0 | 28.55 | 28.672 |
| 0.001 | 0 | 22.078 | 21.901 |
| 0 | 0 | 17.869 | 17.895 |
| 0 | 0 | 12.948 | 13.227 |
| 0.002 | 0.002 | 20.45 | 20.641 |
| 0 | 0 | 30.22 | 32.43 |
| 0 | 0 | 12.964 | 13.29 |
| 0.002 | 0.002 | 33.325 | 34.119 |
| 0 | 0 | 27.005 | 29.667 |
| 0 | 0 | 21.693 | 21.05 |

We load two data sets, **bb2013** and **bb2014** which represent same data for two years, but with some differences in the file structure.

To do any data analysis, we need to reconcile the file structure. For this, we can run **datadiff.** by adding an R cell to our notebook.

Datadiff takes two datasets and generates **a list of patches** that can be applied to transform the structure of the first dataset into the structure of the second dataset.

We first print the inferred patches.

We can switch the tab to **bb2014nice** to see the new transformed dataset, which is now compatible with the **bb2013** dataset.

## Sample data analyses

This section contains larger data analyses created using Wrattler that use the unique features provided by Wrattler such as the ability to create polyglot notebooks.

- Analysing broadband quality in the UK - In this notebook, we analyse the difference between internet quality in rural and urban areas between 2014 and 2015.

## Data Wrangling tools

Wrattler is a part of project that aims to build a variety of tools to simplify tedious data wrangling and cleaning tasks. The following notebooks show some of the other tools that we build, integrated into Wrattler:

- Joining data using datadiff - Datadiff identifies structural differences between pairs of (related) tabular data sets and returns an executable patch.
- Inferring column types using ptype - ptype is a probabilistic type inference method that can robustly annotate a column of data.
- Better CSV parsing with CleverCSV - This notebook shows the use of the CleverCSV tool, which is a replacement for the standard Python csv module.

## Empty notebooks

- Open an empty notebook - This opens an empty Wrattler notebook where you can add JavaScript, Python, R and Markdown cells as you wish.

+ add below    ✕ remove this

In addition to **datadiff,** Wrattler also comes with examples showing **ptype** for inferring types of columns and **CleverCSV** for smart CSV parsing.

**Comprehensive**
Simplifying the data analytics process

Wrattler helps with **all stages** of the data analytics process.

We look at a larger notebook, analysing the UK broadband quality data.

We can mix **R and Python** for loading data with **JavaScript** for quickly visualizing and exploring data.

We have access to **datadiff** and other data wrangling tools, which help us make our data ready for interesting analytical tasks.

Wrattler makes it easy to analyse data using the **wide range of libraries** available for R and Python.