# AI-based Structured Web Data Extraction

Jan Joneš

June 2022

Charles University

# Web data extraction



**Verticals and Attributes** — **One Labeled Seed Site** — **Many Unseen Sites**

**Books**
- Title
- Author
- Publisher
- Publish Date

**Restaurants**
- Name
- Cuisine
- Address
- Phone

**Autos**
- Model
- Price
- Engine
- Fuel Economy

Q. Hao et al. "From One Tree to a Forest: a Unified Solution for Structured Web Data Extraction". 2011

## Motivation

- Scraping use-cases:
  - price comparison,
  - analysis,
  - creating datasets

## Motivation

- Scraping use-cases:
    - price comparison,
    - analysis,
    - creating datasets
- Manual scrapers need manual labor

**Approach**

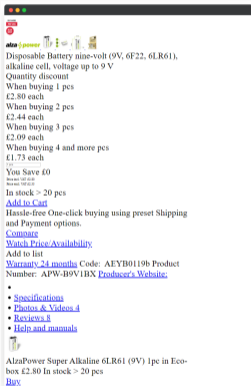Text, DOM tree, visuals → Deep learning model

# Datasets

## SWDE (2011)

| website | pages | price values | price nodes | economy values | economy nodes | engine values | engine nodes | model values | model nodes |
|---|---|---|---|---|---|---|---|---|---|
| aol | 2,000 | 2,000 | 2,001 | 1,845 | 1,845 | 0 | 0 | 2,000 | 2,009 |
| autobytel | 2,000 | 1,994 | 1,994 | 1,816 | 1,816 | 2,000 | 2,000 | 3,998 | 4,023 |
| automotive | 1,999 | 1,999 | 1,999 | 1,999 | 1,999 | 1,999 | 1,999 | 1,999 | 1,999 |
| autoweb | 2,000 | 2,000 | 2,001 | 4,000 | 4,000 | 1,998 | 1,998 | 4,000 | 4,000 |
| carquotes | 2,000 | 2,000 | 2,000 | 2,000 | 2,000 | 2,000 | 2,000 | 2,000 | 2,000 |
| cars | 657 | 647 | 741 | 607 | 607 | 1,168 | 1,168 | 657 | 683 |
| kbb | 2,000 | 2,000 | 4,732 | 2,000 | 2,000 | 2,000 | 2,938 | 4,000 | 4,000 |
| motortrend | 1,267 | 1,267 | 1,267 | 2,534 | 2,534 | 1,267 | 1,267 | 1,267 | 1,267 |
| msn | 2,000 | 3,901 | 4,007 | 3,623 | 4,081 | 2,000 | 2,000 | 2,000 | 3,177 |
| yahoo | 2,000 | 2,000 | 2,073 | 2,000 | 2,000 | 2,000 | 2,000 | 2,000 | 6,094 |
| total | 17,923 | 19,808 | 22,815 | 22,424 | 22,882 | 16,432 | 17,370 | 23,921 | 29,252 |

## Apify e-commerce (2022)

| website | pages | name | price | cat | images | short | long | spec |
|---|---|---|---|---|---|---|---|---|
| alza | 2,493 | 2,493 | 2,478 | 2,493 | 2,493 | 2,493 | 2,493 | 2,477 |
| asos | 499 | 499 | 499 | 499 | 499 | 0 | 499 | 0 |
| bestbuy | 1,000 | 1,000 | 986 | 998 | 1,000 | 0 | 1,000 | 1 |
| bloomingdales | 448 | 462 | 447 | 448 | 448 | 0 | 1,087 | 0 |
| conrad | 1,500 | 1,500 | 1,499 | 1,500 | 1,485 | 1,500 | 1,500 | 1,495 |
| etsy | 250 | 250 | 250 | 250 | 250 | 0 | 1,000 | 0 |
| ikea | 1,972 | 1,972 | 1,972 | 1,959 | 1,972 | 2,261 | 1,917 | 0 |
| notino | 808 | 808 | 702 | 808 | 808 | 808 | 785 | 783 |
| radioshack | 499 | 499 | 499 | 499 | 498 | 0 | 6,204 | 4 |
| tesco | 1,500 | 1,499 | 1,465 | 1,499 | 1,499 | 0 | 17,205 | 21,526 |
| total | 10,969 | 10,982 | 10,797 | 10,953 | 10,952 | 7,062 | 33,690 | 26,286 |

# Implementation

# Visuals

Page without CSS or JavaScript evaluated is difficult to understand.

# Visuals

Page without CSS or JavaScript evaluated is difficult to understand.
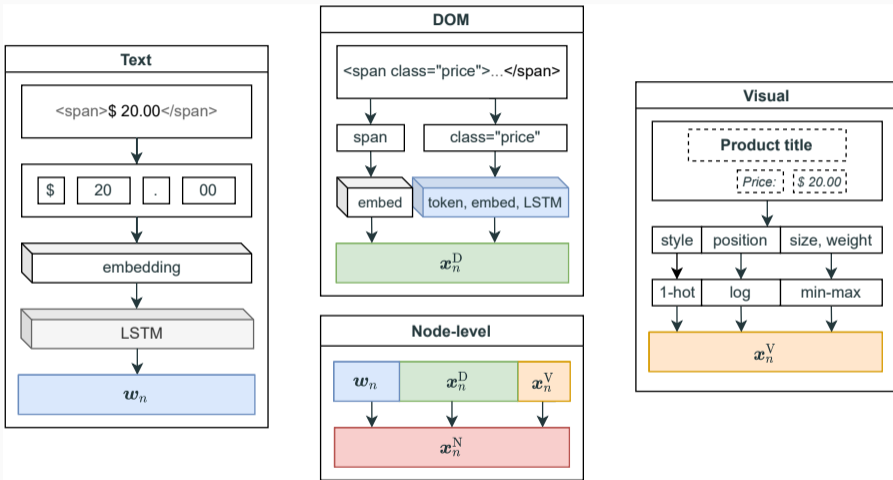




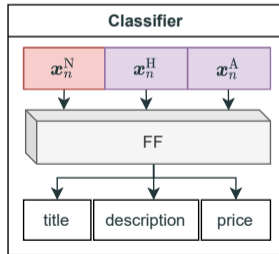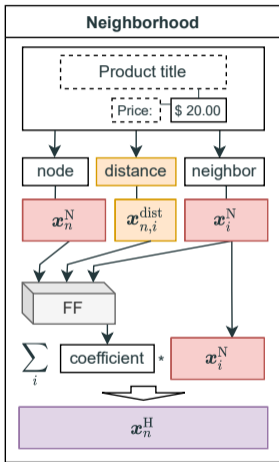From www.alza.co.uk, archived at https://archive.ph/7xXoG

**Web Page**   **DOM-tree**

Q. Hao et al. "From One Tree to a Forest: a Unified Solution for Structured Web Data Extraction". 2011

**Text**

<span>$ 20.00</span>

$ | 20 | . | 00

embedding

LSTM

$\boldsymbol{w}_n$

**DOM**

<span class="price">...</span>

span | class="price"

embed | token, embed, LSTM

$\boldsymbol{x}_n^{\mathrm{D}}$

**Node-level**

$\boldsymbol{w}_n$ | $\boldsymbol{x}_n^{\mathrm{D}}$ | $\boldsymbol{x}_n^{\mathrm{V}}$

$\boldsymbol{x}_n^{\mathrm{N}}$

**Visual**

**Product title**

*Price:* | *$ 20.00*

style | position | size, weight

1-hot | log | min-max

$\boldsymbol{x}_n^{\mathrm{V}}$

## Baseline (2021)

- No visuals
- Old dataset
- Non-separated test set

e-commerce pages [F1]

# Demo



Generated by live demo at bit.ly/awedemo

## Summary

- Node.js extractor of visual features

- Node.js extractor of visual features
- Python deep learning model

## Summary

- Node.js extractor of visual features
- Python deep learning model
- Evaluated on a modern dataset

- Node.js extractor of visual features
- Python deep learning model
- Evaluated on a modern dataset
- Docker image for reproducible training

# Summary

- Node.js extractor of visual features
- Python deep learning model
- Evaluated on a modern dataset
- Docker image for reproducible training
- Demo app for live inference

# Thank you for your attention

bit.ly/awedemo
github.com/jjonescz/awe