

Intelligent Web Crawling

Jakob Torben, Ricardo Mokhtari and Yousef Nami

Advanced Data Science Team, Imperial College London

The Team



Dr. Ovidiu Serban

Supervisor

Research Fellow, DSI
Imperial College London



Jakob Torben

Imperial College London
BSc Physics
MSc Applied Computational
Science



Ricardo Mokhtari

Imperial College London
MEng Molecular Bioengineering

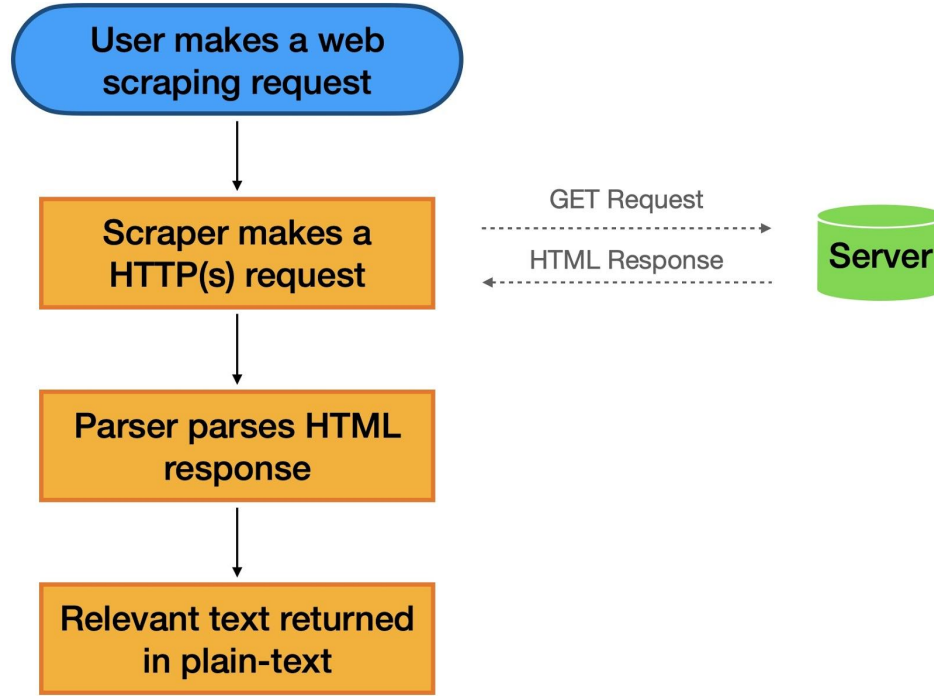


Yousef Nami

Imperial College London
MEng Mechanical Engineering

Problem Definition

Problem Definition



Temporality:

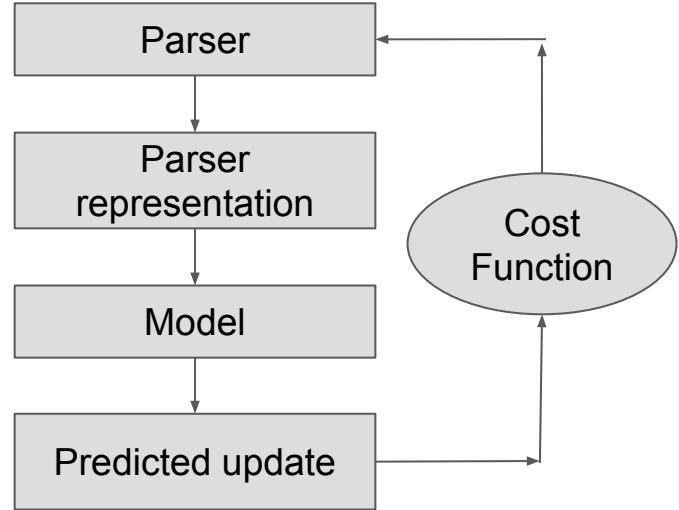
- websites change over time
- Changes in the HTML structure
- Changes in domain names

Scale:

- 100 companies
- All HTML structures different
- Different subdomains for leadership pages
- Duplicate names

Problem Definition

- Rule based methods not very effective due to very high variability
- Data driven approach could work
- Requires a 'Cost Function' to determine the similarity/difference between any two parsers



Case Study

Case Study: GSK

Corporate Executive Team

The Chief Executive Officer is responsible for the management of the business and is assisted by the Corporate Executive Team (CET). The CET manages our activities, and each member is responsible for a specific part of the business.

To view all our images in high resolution, visit our [GSK Flickr](#) page.



Dame Emma Walmsley

Chief Executive Officer

[Read more](#)



Iain Mackay

Chief Financial Officer

[Read more](#)

```

▼<li class="grid-listing__item">
  ::marker
  ▼<a class="grid-listing__link" href="/en-gb/at
    ▶<div class="grid-listing__img">...</div>
    ▼<div class="grid-listing__content">
      <h2>Dame Emma Walmsley</h2>
      <p>Chief Executive Officer</p>
      <span>Read more</span>
    </div>
  </a>
</li>
▶<li class="grid-listing__item">...</li>
▶<li class="grid-listing__item">...</li>
▶<li class="grid-listing__item">...</li>
▶<li class="grid-listing__item">...</li>

```

Case Study: GSK Parser

```

▼ <li class="grid-listing__item">
  ::marker
  ▼ <a class="grid-listing__link" href="/en-gb/at
    ▶ <div class="grid-listing__img">... </div>
    ▼ <div class="grid-listing__content">
      <h2>Dame Emma Walmsley</h2>
      <p>Chief Executive Officer</p>
      <span>Read more</span>
    </div>
  </a>
</li>
▶ <li class="grid-listing__item">... </li>
▶ <li class="grid-listing__item">... </li>
▶ <li class="grid-listing__item">... </li>
▶ <li class="grid-listing__item">... </li>

```

```

# parse the current GSK board
def parse_current(self, response):
    # define selector that contains all items
    all_people = response.css("li.grid-listing__item")
    print(all_people)
    # iterate through items
    for person in all_people:
        # manually parse name
        name = person.css("a>div>h2::text").get()
        # manually parse title
        title = person.css("a>div>p::text").get()

        now = datetime.datetime.now()
        year = now.year

        # Return item
        yield self.create_board(name, title, year)

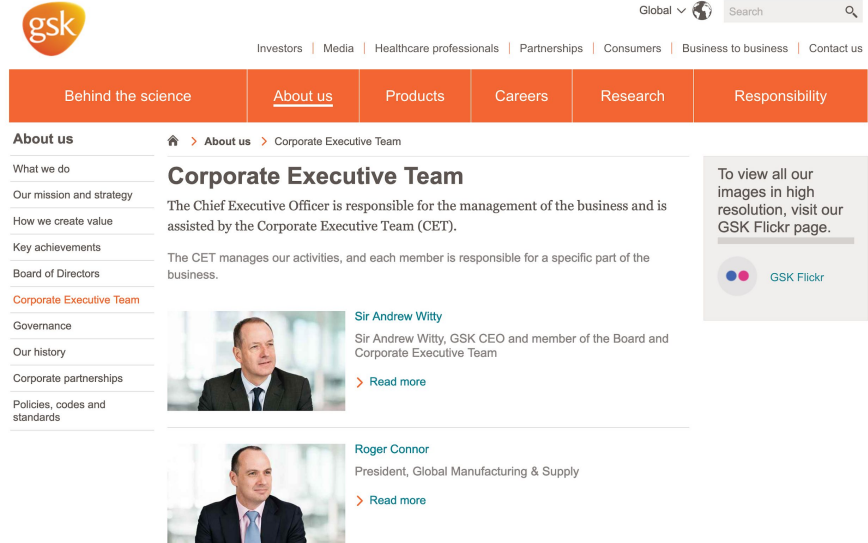
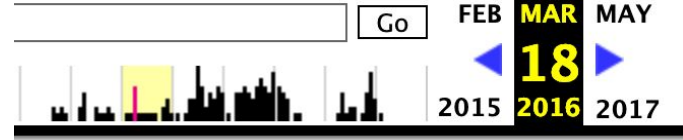
```


Case Study: Historical Information

3 Different parsers:

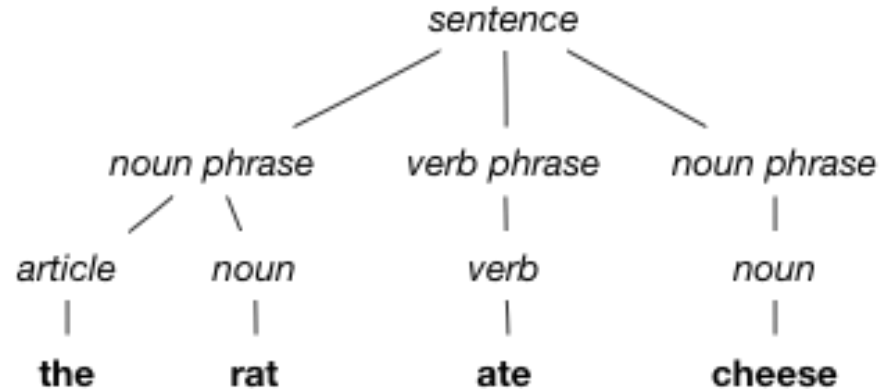
- Current web page
- 2 archived versions

How can we compute the similarity between these parsers?



Case Study: What are ASTs?

- Recognize code elements
- Understand code syntax





Case Study: Similarity of ASTs

- Plagiarism detection
- ASTs better captures intent of code
- Tree edit distance
 - Weighted number of edit operations (insert, delete, and modify) to transform one tree to another
- Removing node normalisation

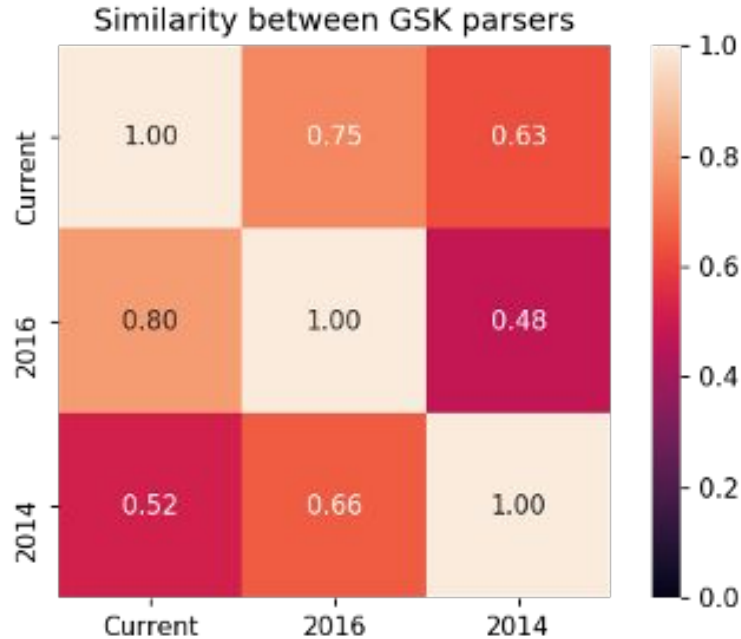
GSK: 2016

```
all_people = response.css("article.listing-item.with-image")
```

GSK: 2014

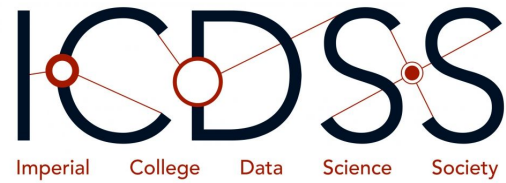
```
all_people = response.css("a.titleLink.textDecorateNone")
```

Case Study: Results



Limitations

- Historical sites might have completely different structure
- Only works for relatively similar parsers
- Websites lacking structure
- Example: Shell



EXECUTIVE COMMITTEE

The Royal Dutch Shell plc Executive Committee operates under the direction of the Chief Executive Officer and is responsible for Shell's overall business and affairs.

The Chief Executive Officer has final authority in all matters of management that are not within the duties and authorities of the Board or of the shareholders' general meeting. The Executive Committee supports the Chief Executive Officer and implements all Board resolutions and supervises all management levels in Shell.



Ben van Beurden

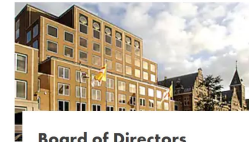
Chief Executive Officer.



Jessica Uhl

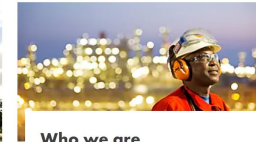
Chief Financial Officer.

MORE IN ABOUT US



Board of Directors

The Board of Directors meet to discuss reviews and reports on the business and plans of Royal Dutch Shell plc.



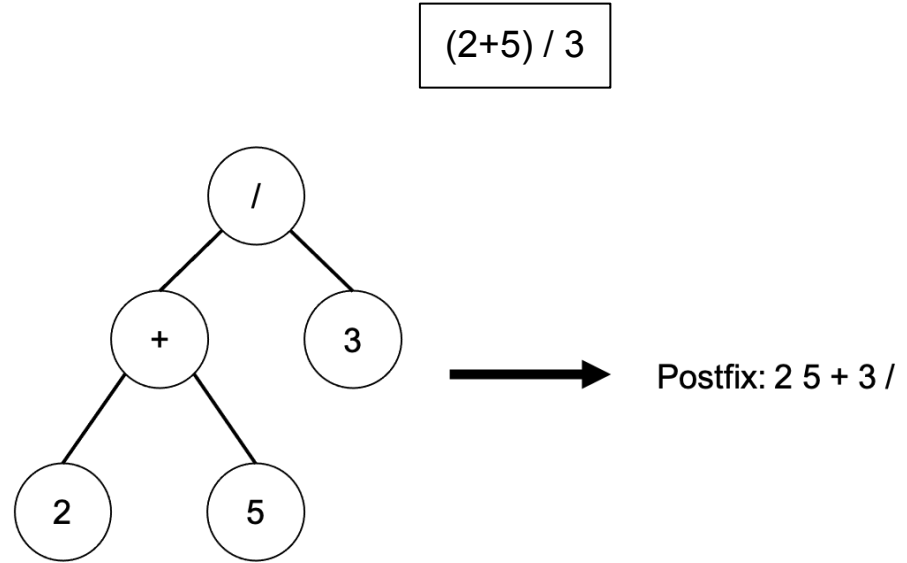
Who we are

Find out about our business, people and how we are working to power progress together with more and cleaner energy solutions.

Alternative Method

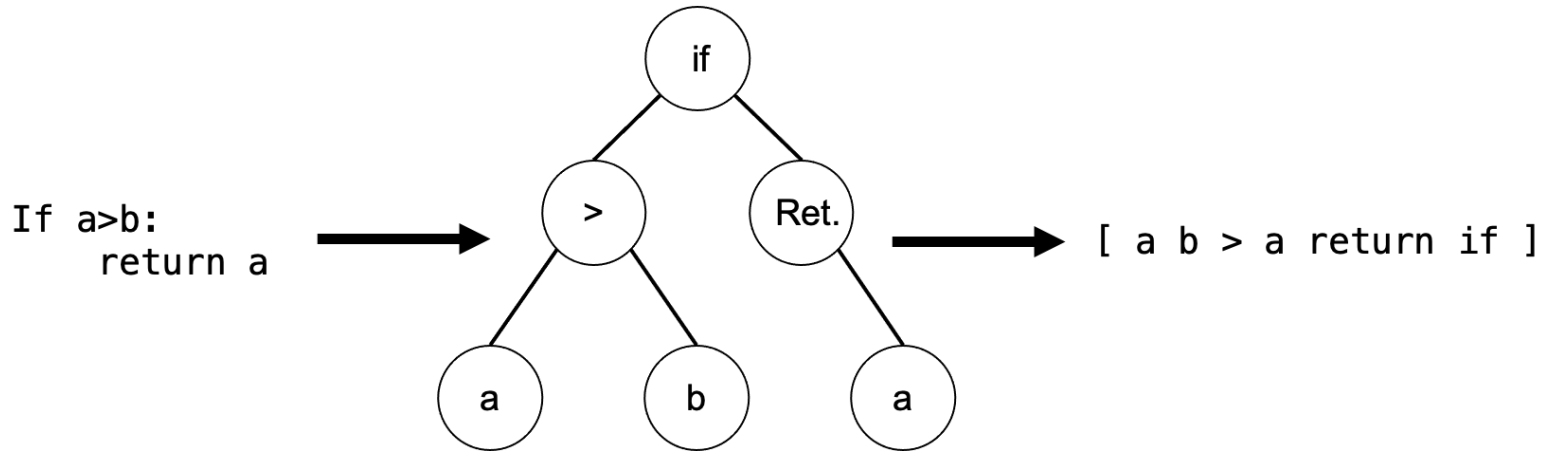
Vectorial Representation of ASTs

- Computing vector similarity is straightforward
- Standard way of representing expression trees as vectors:
Postfix representation



Vectorial Representation of ASTs

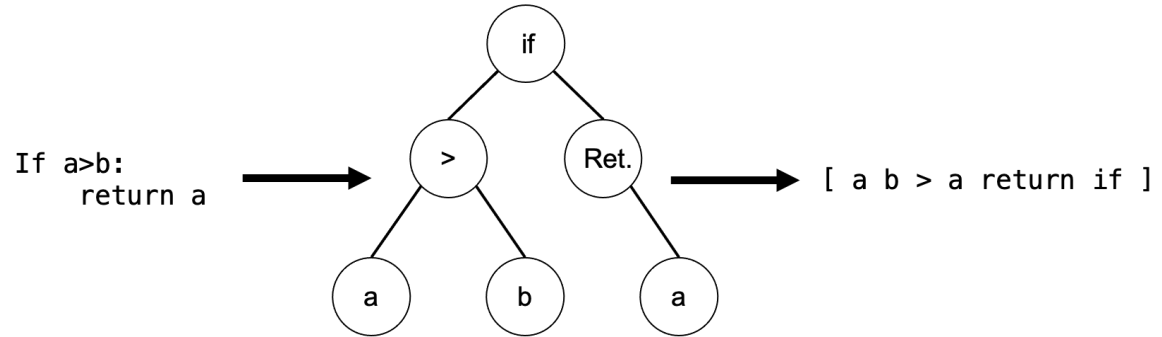
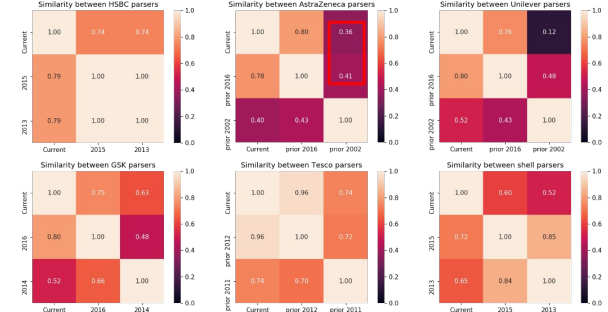
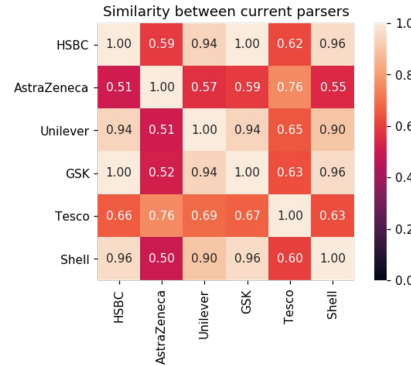
- Convert code to AST, traverse AST to produce vector representation
- To compute similarity: Jaccard, Levenshtein, Dice



Key Takeaways & Future Work

Key Takeaways

- 18 parsers created in total
- Explored Cost Function based on Abstract Syntax Tree representation
- Post-fix representation of ASTs in vector format



Future work

- Analysis of post-fix represented parsers using common vector comparison metrics (i.e. Levenshtein, Jaccard)
- Re-visit direct AST comparison by reproducing examples found in the literature
- Identify models that can automate the parser update (RL, triplet loss)
- Opportunities:
 - Reduce labour costs
 - Reduce time to update/create new parsers

Thank you for listening!

Appendix

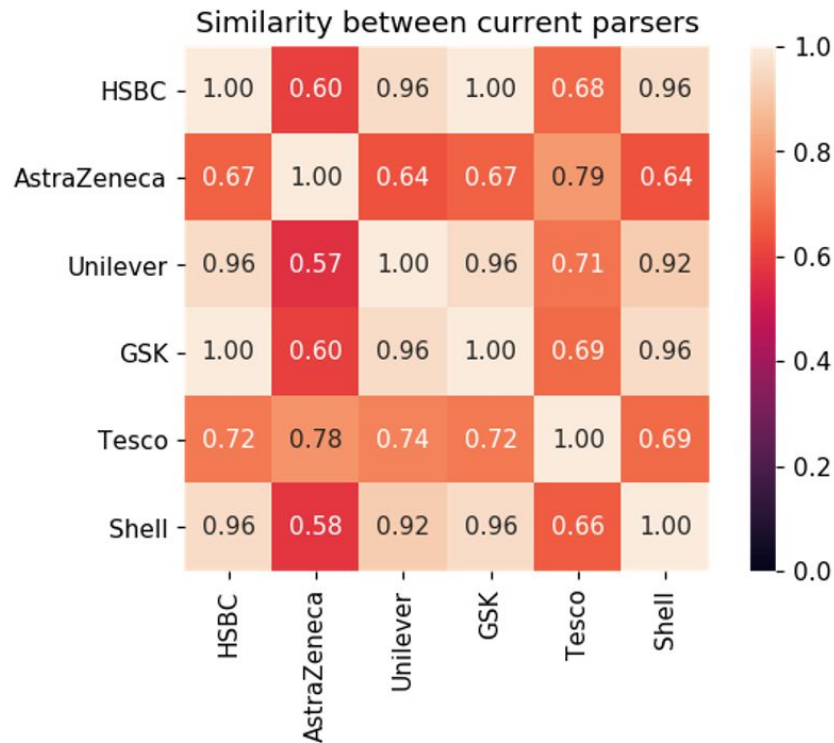
TreeDiff() - unsymmetric

Similarity between current parsers

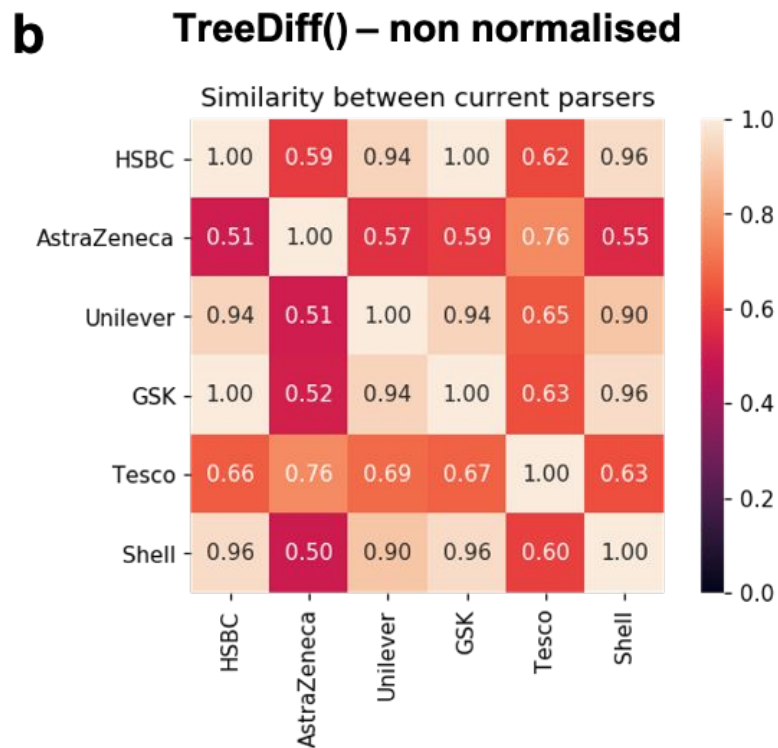
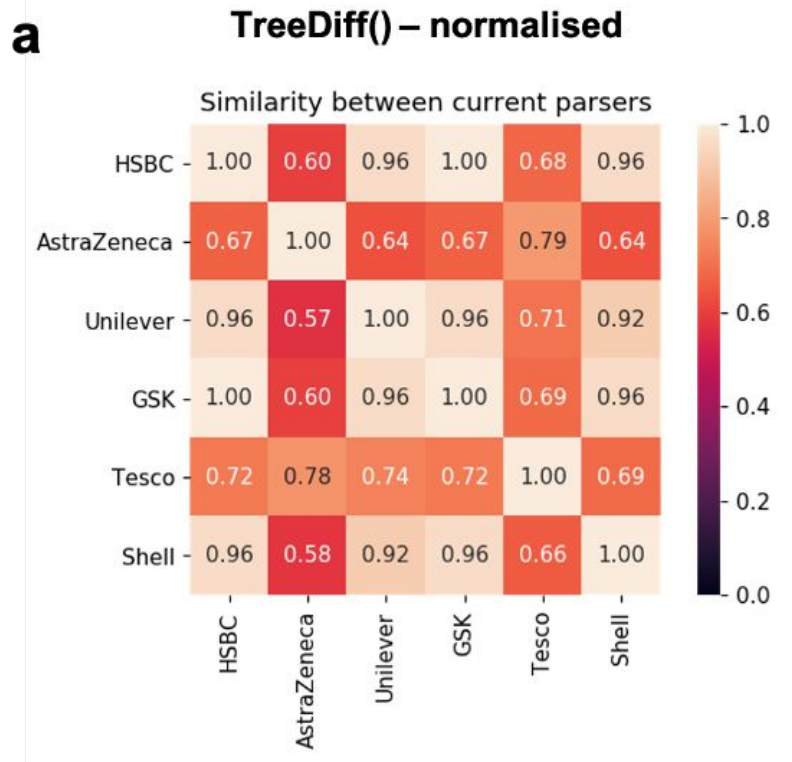


Appendix

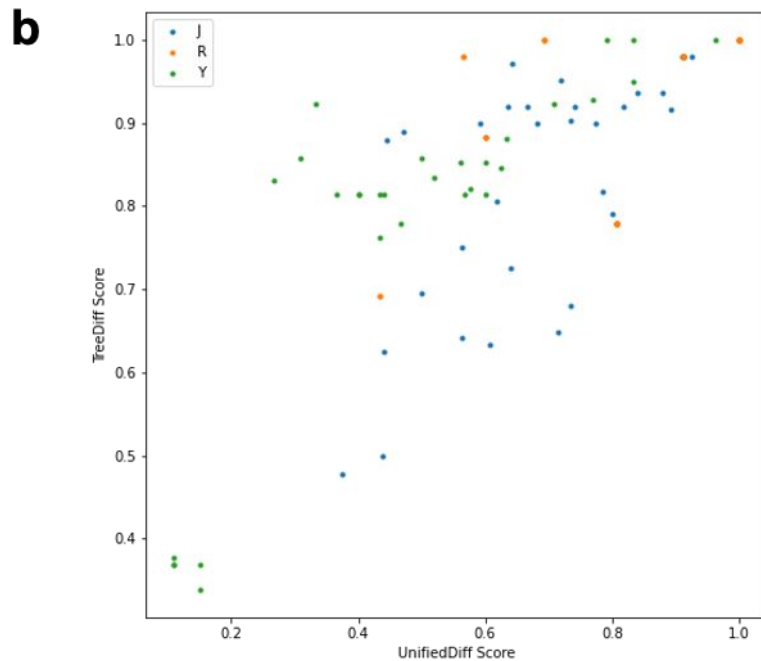
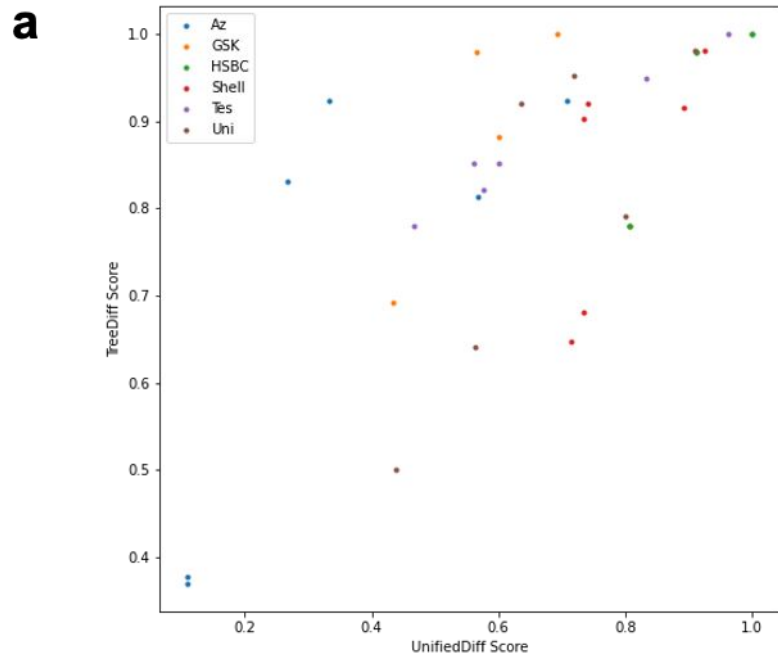
TreeDiff() – normalised



Appendix



Appendix



Appendix

