

# Modular Annotation for Natural Language Processing



Diego Montes | Bengisu Cuneyit | Jisoo Kim | Shreyas Chickerur | Sahithi Tummala | Saimonish Tunguturu

## Data Annotation for Contemporary ML Challenges

### Inspiration

- Data is the backbone of machine learning models, yet real-world data is messy: data annotations format this data in a way that a model can learn from.
- In the past, we had used and attempted to extend an open-source annotation tool: cdQA-annotator; however, there were a number of pitfalls:
  - No centralized storage: annotators had to manually download annotations after each annotation session and upload them to Sharepoint.
  - Difficult to implement new features: the tool was no longer maintained, and the code base's infrastructure did not support pluggable features.
  - A lack of a project structure and workflow.
- As such, our project's specifications revolved around these three missing features.

Im Vergleich zu Säuglingen und jüngeren Kleinkindern wurde bei Kindern ,die älter als 24 Monate waren , eine höhere Rate lokaler Reaktionen beobachtet .

Compared to infants and younger children , a higher rate of local reactions has been observed in children older than 24 months of age .

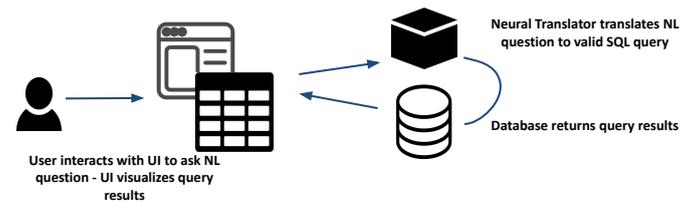
SOURCE: [pubmed.ncbi.nlm.nih.gov/17204475/](https://pubmed.ncbi.nlm.nih.gov/17204475/)

### An Example Translation Data Annotation:

- The top portion displays the text to-be-translated.
- The bottom portion shows the translated text, annotated for errors.

### Why Another Annotation Tool?

- Distinct from many open-source annotators, a priority was making our annotator pluggable for future tools and software contributors.
- Similarly, despite a number of paid annotation tool services that offer the features we were looking for, we wanted to create a tool that could annotate data for the newest ML/AI challenges:
  - One such example is NL2SQL: no annotation service offers a tool for natural language to SQL annotations (depicted in the diagram below)



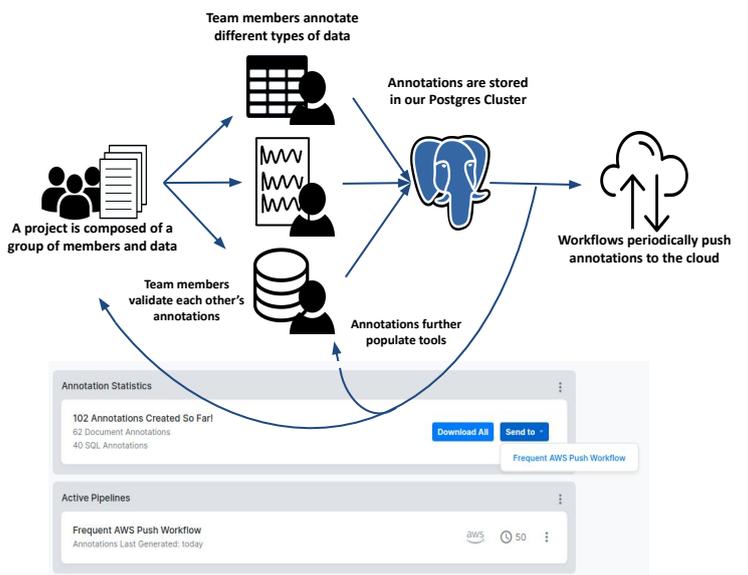
## Annotating Workflow

### Project Management

- Organization of data annotations becomes increasingly important as annotation projects get larger, this calls for:
  - Store and retrieve functionality for data and annotations.
  - Project and task structure to distribute annotation workload.
  - Workflows to push annotations to other storages automatically.

### Annotation Pipeline

- A user creates a project, inviting other members if necessary.
- Any team member uploads data belonging to the supported formats.
  - Tasks are automatically created and divided among team members.
- As annotations are created, project statistics are updated and at any point all annotations for a project can be downloaded.
- Optionally, annotations can be periodically sent to an AWS or GCP storage.



## Question-Answer Systems, Classification, and SQL Annotators

```
the value from the environment by enforcing os.environ["DJANGO_SETTINGS_MODULE"] = "mysite.settings" in your wsgi.py.
```

```
Applying WSGI middleware¶
```

```
To apply WSGI middleware you can wrap the application object. For instance you could add these lines at the bottom of wsgi.py:
```

```
from helloworld.wsgi import HelloWorldApplication
```

```
application = HelloWorldApplication(application)
```

A Sample Annotation from the Question-Answer Annotator

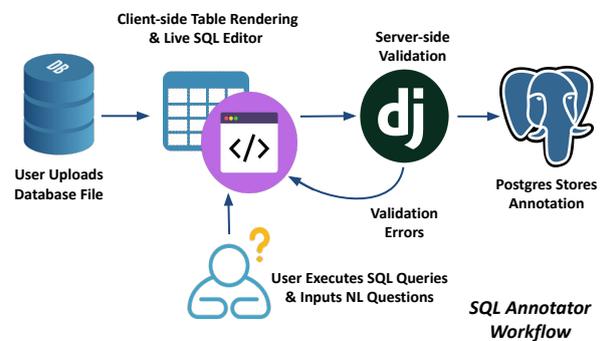
Classification Annotator Table View

### Generic Annotation Workflow

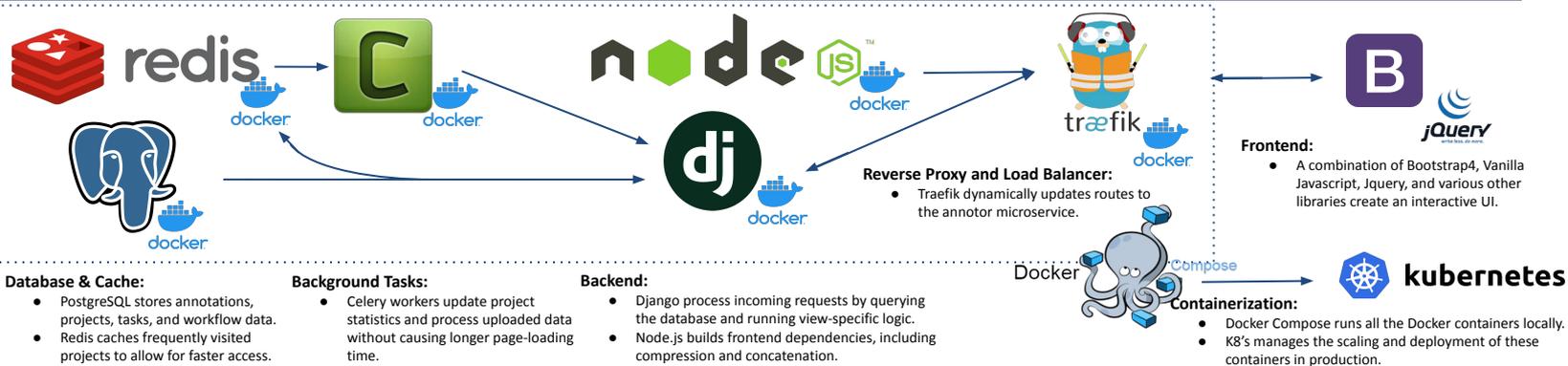
- The end-user uploads a supported file type (.db, .csv, .txt, .md) directly to the annotator or pulls a previously uploaded file from a project.
- The user creates an annotation, specific to the annotator being used:
  - For QAS, the user writes a natural-language question and selects an answer span.
  - For SQL Annotations, the user writes a NL question and a corresponding SQL query.
  - For Classification, the user drags labels onto table rows, classifying them into categories.
- The annotation is sent to Django for additional validation and bleaching.
- The annotation is saved to Postgres and project statistics are updated.

### Classification Tool Motivation

- Classification of natural language (NL) is not always on a document-by-document basis: text data can be found in tables as well
  - The tool should include an annotator that can classify rows of tabular data into a distinct categorical column.
  - Users should be able to make and search for class labels.
  - Annotating a given row should be as easy as possible.



## Technology Stack



## Conclusions & Future Goals

### Conclusions:

- The learning curve for a full stack of technologies was larger than initially anticipated:
  - The MAT has distinct 8 microservices in its stack.
- Development time for seemingly small features increases drastically as they need to be integrated with existing features.
- Necessity of a separation between frontend and backend.

### Future Technical Goals:

- Create better collaboration between the frontend design and the backend logic by using a REST or GraphQL Framework.
  - Similarly, improve frontend design development by using React or another Javascript framework versus Vanilla Javascript.
- Add more annotation tools to the MAT suite
  - Multiple choice image labeling
  - Audio transcriptions
  - Named entity labeling
- Implement integrations with other data services, such as Pure Storage.

## References & Acknowledgments

### References:

- Django (Version 1.5) [Computer Software]. (2013). Retrieved from <https://djangoproject.com>.
- Merkel, D. (2014). Docker: lightweight linux containers for consistent development and deployment. Linux Journal, 2014(239), 2.
- Bootstrap (Version 4.6) [Computer Software]. (2021). Retrieved from <https://getbootstrap.com/>

### Acknowledgments:

- We would like to thank our:
  - Corporate Partner Mentors: Sarah Rodenbeck and Justin Gould
  - Student Mentor: Rishabh Rajesh
  - And the entire Ford team for giving us this opportunity, letting us access their internal data and systems, and guiding us as we created and refined this project.
- We would also like to thank: Dr. Mark Daniel Ward, Ellen Gundlach, and Maggie Ann Betz for supporting us during the Academic Year and providing us with the resources for the successful completion of our project.