

# Intermediate Presentation – Classification of Text ML-Approach

Master Internship: Approaching Information System Challenges with Natural Language  
Processing (IN2106, IN2130)

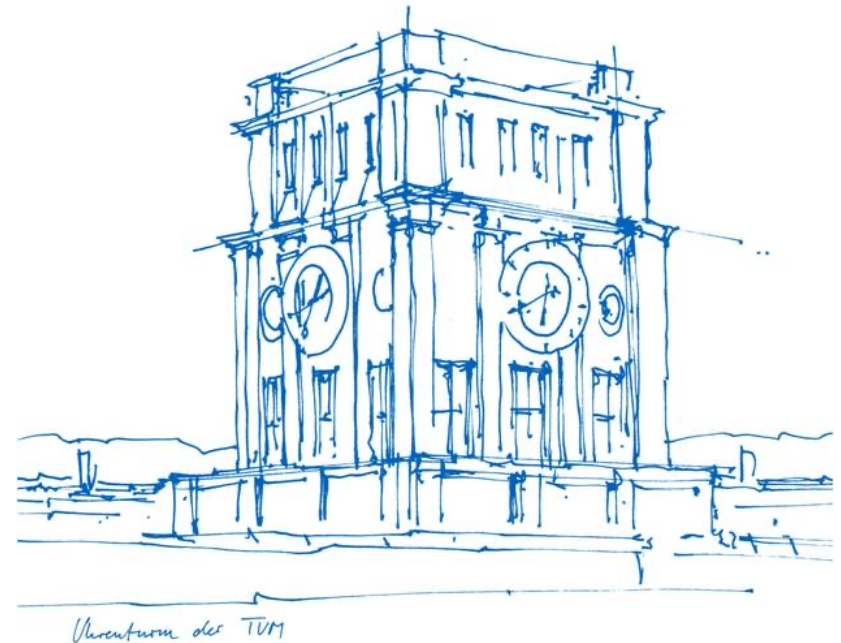
Patrick Ahrend

Garching, 5. December 2023

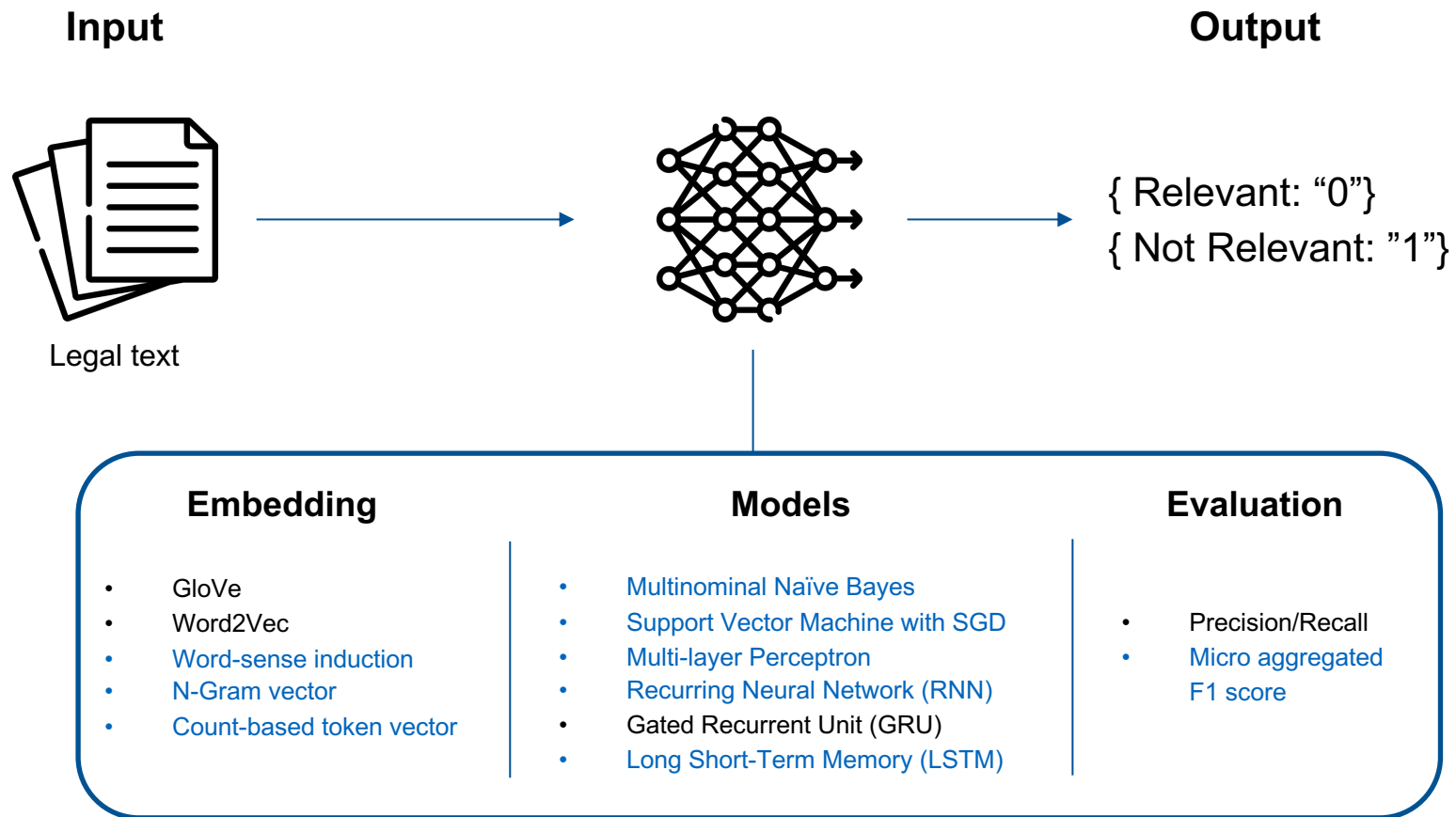


# Agenda

- 1) Project Goal
- 2) Project Plan
- 3) Data Gathering
- 4) Repository Structure



# Process Discovery Machine Learning Pipeline



Source: Automated Business Process Discovery from Unstructured Natural-Language Documents 2020, [https://link.springer.com/chapter/10.1007/978-3-030-66498-5\\_18](https://link.springer.com/chapter/10.1007/978-3-030-66498-5_18)

# Currently still in the Data Pre-Processing Step

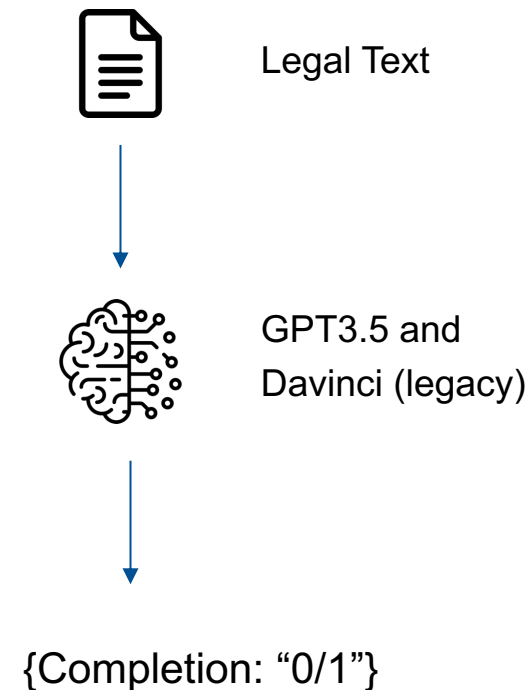


# Australian Data + GDPR for the Data

	ratio 1:10 - input approaches				
	10% (group A)		45% (group B)	45% (group C)	100%
Processes	compliance relevant text paragraphs	informative relevant text paragraphs	non-relevant text paragraphs from relevant documents	non-relevant text paragraphs from non-relevant documents	total number of text paragraphs
1: travel insurance claim	21	28	220	220	489
2: know your customer	24	7	140	140	311
3: hire employee		9	40	41	90
4: GDPR 1-Data breach		8	36	36	80
5: GDPR 2-Consent to use the data		16	72	72	160
6: GDPR 3-Right to Access		11	50	49	110
7: GDPR 4-Right to Portability		4	18	18	40
8: GDPR 5-Right to Withdraw		4	18	18	40
9: GDPR 6-Right to Rectify		2	9	9	20
10: GDPR 7-Right to be forgotten		13	59	58	130
					1470

# Data Labelling with LLM to Gather More Data

- **Davinci Fine-Tuning** : {  
  **prompt**: "Process: <Process\_Description>  
  /n/n Text: <Legal\_Passage> /n/n Relevant:"  
  **completion**: "<0/1>###"}
- **GPT 3.5 Fine-Tuning** :  
  { **system\_message**: "Determine if the text is  
  relevant to the process described" ,  
  **user\_message**: "Process Description :  
  <Process\_Description>/n/n Text to classify:  
  <Legal\_Passage>/n",  
  **system\_message**: "<0/1>###" }



# Cookie Cutter for the Repository Structure

```

├── Makefile          <- Makefile with commands like `make data` or `make train`
├── README.md         <- The top-level README for documentation and instruction on how to run the code.
├── data
│   ├── external      <- Data generated from the fine-tuned models for labeling.
│   ├── interim       <- Intermediate data that has been transformed.
│   ├── processed     <- The final data sets for modeling.
│   └── raw           <- The original, immutable data from the other repository which I decided to use.
├── models            <- Trained and serialized models, model predictions, or model summaries
├── notebooks         <- Jupyter notebooks. Contains for instance the fine-tuning of GPT for labeling notebooks.
├── references        <- Data dictionaries, manuals, and all other explanatory materials to understand the data bet
├── requirements.txt  <- The requirements file for reproducing the environment
├── setup.py          <- makes project pip installable (pip install -e .) so src can be imported
├── src               <- Source code for use in this project.
│   ├── __init__.py   <- Makes src a Python module
│   ├── data          <- Scripts to download or generate data
│   │   └── make_dataset.py
│   ├── features       <- Scripts to turn raw data into features for modeling and create word embeddings
│   │   ├── build_features.py
│   │   └── build_word_embeddings.py
│   ├── models         <- Scripts to train models and then use trained models to make
│   │                   predictions
│   │   ├── predict_model.py
│   │   └── train_model.py
│   └── visualization  <- Scripts to create exploratory and results-oriented visualizations as well as word embeddin
│       └── visualize.py

```

Source: <https://drivendata.github.io/cookiecutter-data-science/>

Thanks!  
Questions/Feedback?