

NLP – Group 23

Named Entity Recognition for German Official Documents

16. Januar 2026

Tim Greß (12412672)
Viktor Hoffmann (12433741)
Lars Böhmer (12436447)
Jan Tölken (12432831)

Dataset

Labelling

Baselines

RuleChef

LLM API

Performance & Overview

Conclusion & Process-Pipeline

Dataset

Name

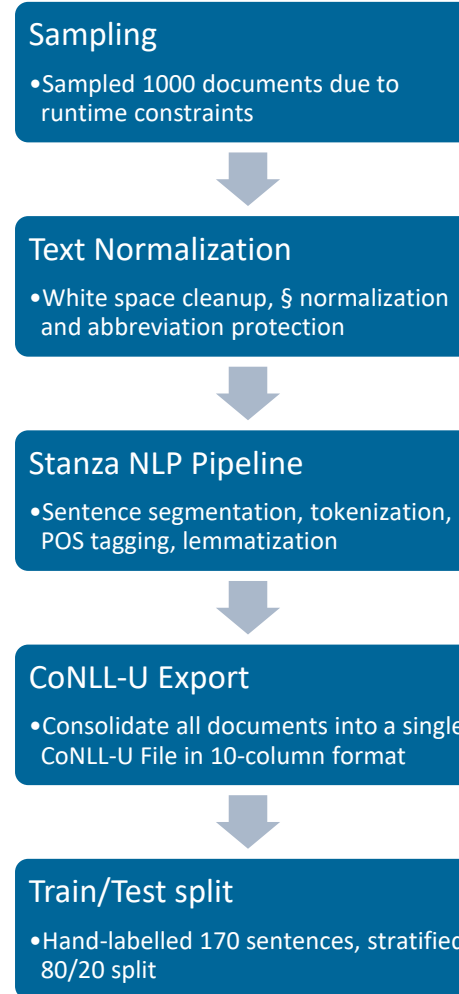
FinCorpus-DE10k

Description

German financial/legal corpus containing 10,000 PDF documents with extracted text

Source

'anhaltai/fincorpus-de-10k' on HuggingFace



Labelling Guidelines

BIO Tagging Scheme

- Begin / Inside / Outside

Entities

- Organizations (ORG)
- Monetary (MON)
- Legal References (LEG)

Organizations – ORG

Include	Example
Full company names	Landesbank Baden-Württemberg
Company abbreviations	LBBW
Legal forms	KR Labs GmbH
Financial institutions	EZB
Banks	Sparkasse
Exclude	Example
Generic organizational terms	Eine Sparkasse
All other names	Max Mustermann

Organization Example

Die **Landesbank Baden-Württemberg** ist eine deutsche Bank, die mit der **Deutschen Bank AG** und der **EZB** zusammenarbeitet und jährlich **2.000.000.000 Euro** umsetzt.

B-ORG

I-ORG

B-MON

I-MON

O

Labelling Pitfalls

- Even manually it is often hard to distinguish between abbreviations for organizations and other names
- What to include as monetary values?
 - *Return of 20%*
- Very time consuming – labelled 170 sentences

Labelling Tool with GUI

CoNLL-U Manual NER Annotation (Windows)

Sentence 4/150 (NOT CONFIRMED) | Token 12/39 | Output: test.conllu | Autosave: every 25 labels

Autosave now Save now

Sentence (current token is highlighted)

Die jeweilige Beteiligungsquote (in Prozent) des Aktionärs bezieht sich auf das gezeichnete Kapital der DZ BANK in Höhe von EUR 4.899.938.940,00 abzüglich der von der DZ BANK gehaltenen 93.247.143 eigenen Stückaktien mit einem rechnerischen Gesamtbetrag von EUR

Current token

auf

id=12 | current tag=O | context: sich [auf] das

O B-ORG I-ORG B-MON I-MON B-LEG I-LEG

Rest des Satzes auf O Review/Confirm this sentence

<< Prev token Next token >> | << Prev sentence Next sentence >> | Go to first untagged ()

Baselines

EnhancedRuleNER

- Hand-crafted regex patterns per entity type
- Validation heuristics using the POS tags (e.g. span contains at least one NOUN or PROPN)

TokenNB

- Token-level Naïve Bayes classifier
- Uniform class priors to compensate heavy class imbalance (~93% O-tags)

CRF

- Linear-chain CRF for BIO tagging
- Rich token features: casing/shape, prefixes/suffixes,
- Domain features: currency symbols (€), paragraph signs (§), company suffixes (GmbH, AG)

SpaCy + Heuristics

- Hybrid labelling combining:
 - spaCy German NER model
 - Key-word based ORG spans
 - Monetary Detection
 - Legal reference detection

German BERT

- Fine-tunes bert-base-german-cased
- Aligns labels to WordPiece tokens
- Token-level + span overlap evaluation

RuleChef

Idea

Use an LLM to *learn*
explicit extraction rules



Result

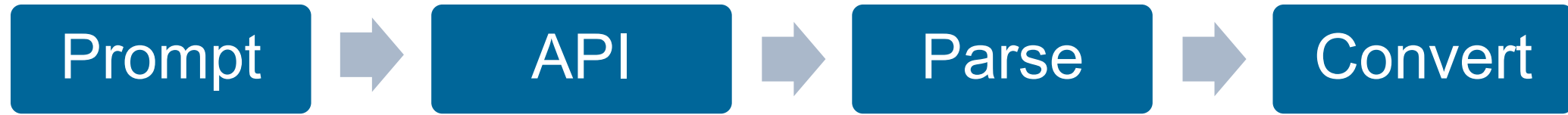
interpretable ruleset that can
be applied like a classic
rule-based NER system

Nach §15 Abs.1
ist die
vollständige....



```

{
  "name": "Extract § + number
+ Abs.",
  "format": "regex",
  "pattern":
  "$\\s*\\d+\\s*Abs\\.\\.?\\s*\\d+"
}
  
```



Input: "Die DZ BANK hat EUR 4.899.938 nach § 271 HGB ausgewiesen."

API Returns:

```

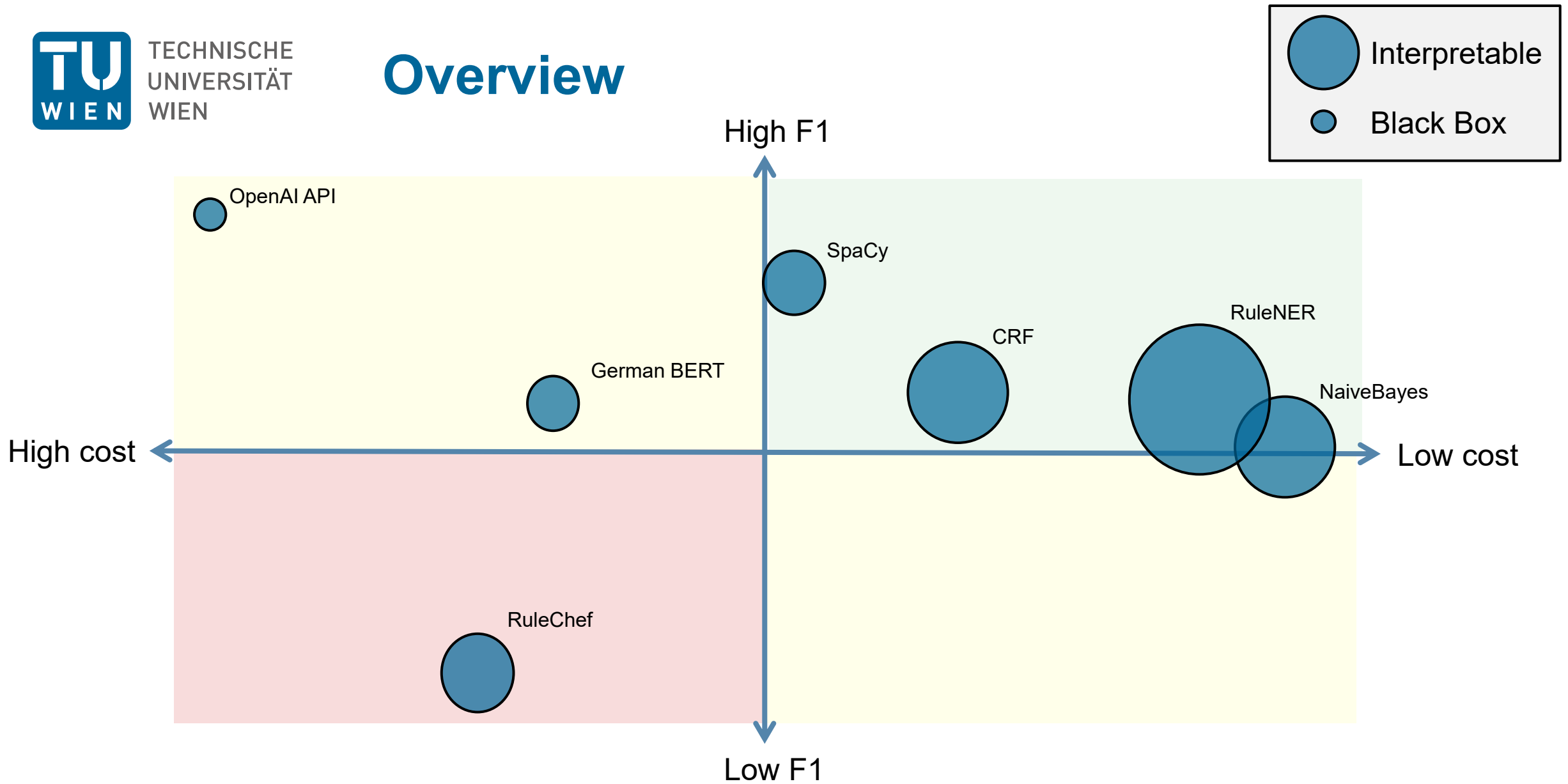
"entities": [
  {"label": "ORG", "text": "DZ BANK"},
  {"label": "MON", "text": "EUR 4.899.938"},
  {"label": "LEG", "text": "§ 271 HGB"}
]
```

Output: O B-ORG I-ORG O B-MON I-MON O B-LEG I-LEG I-LEG O

Performance

Method	LEG F1	MON F1	ORG F1	Macro F1		Run Time (s)
EnhancedRuleNER	0,59	0,67	0,20	0,61		0,01
TokenNB	0,07	0,90	0,20	0,51		0,01
CRF	0,67	0,97	0,06	0,67		0,76
SpaCy	0,86	0,82	0,44	0,77		0,24
German BERT (2 epochs)	0,43	0,69	0,62	0,68		99,52
RuleChef (3 loops)	0,07	0,10	0,01	0,22		370,75
OpenAI API (10 threads)	1,00	1,00	0,59	0,89		167,90

Overview



RuleChef

+



-



- Interpretable results as code or regex
- Fast and deterministic inference as soon as rules are created
- Strong domain control, can adapt to certain keywords, formats etc.

- Refinement often overgeneralizes
- Formatting issues with malformed rules or invalid JSON
- Struggles with unstructured classes like ORG
- Very bad performance in our case
- Inconsistent across runs

Conclusion and Insights

- NER for German official documents is a trade-off problem, no „free lunch“
- Entity structure is very important - consistent underperformance for ORG
- Neural Models improve ORG detection but are computationally expensive
- RuleChef was unstable and underperformed in our setup due to implementation and robustness issues

Suggested Pipeline

