

1. Context & Motivation

- In digital libraries, named entities are among the most relevant information to index documents and the main entry point to their retrieval. Most digitized documents are indexed through their OCRed version and OCR errors hinder their accessibility.
- Prior research shows that named entities are the first point of entry of DL users.

Goal: Quantitatively estimate the impact of OCR quality on NER performance. **<u>Problem</u>**: No text with (original OCR) and (OCR groundtruth) and (annotated NER). Methodology:

- Simulating OCR outputs from existing clean NER corpora with various levels and types of OCR noise.
- Testing state-of-the-art NER systems over clean and noisy data.
- Studying the correlation between NER results and OCR.

2. Overview on NER approaches

« The German racing driver Michael Schumacher won with Mercedes the last race in **China**. »

PER

LOC

ORG

Three main approaches:

Rule-based: rules are based on dictionaries, triggers and linguistic descriptors

- Machine-learning: training, hand crafted features
- Deep-learning: jointly training and learning effective features
- Rule-based systems are clearly disoriented to process noisy data: perpetual updates \rightarrow of rules for noisy NEs and costly manual efforts
- Machine-learning systems can be updated and generalized.
- Deep-learning systems outperform other machine-learning NER systems. \rightarrow

4. Results & Evaluation

Neural-network based system: LSTM-CRF

GT:	Pierre	Van	Н
gold_annot:	PER	PER	Р
OCR lev-0:	Dierre	Van	Н
pred_annot:	PER	PER	0
OCR blur:	Pierre	n	Н
pred_annot:	PER	0	0

Ν	
Re	
90	
90	
84	
78	
69	
64	

An Analysis of the Performance of Named Entity Recognition over OCRed Docume

Ahmed Hamdi, Axel Jean-Caurant, Nicolas Sidère, Mickaël Coustaty, Antoine Doucet

firstname.lastname@univ-lr.fr



Simulation of many OCR outputs similar to documents stored on digital libraries. Evaluation of the evolution of NER accuracy depending on the level of noise in the text. NER Accuracy drops from 90% to 60% with variable CER and WER from 1% to 7% and from 8% to 20% respectively.

Adding weights to OCR outputs at the character level and the word level. NE-focused OCR post-correction => high impact on information access in digital libraries.

n	tc

La Rochelle



Bleeding effect: simulates the verso ink that seeps through the recto side

Japan began the defence of their Asian Cup title with a lucky 2-1 win against Syria in a Group C championship But China saw their luck desert them in the second match of the group , crashing to a surprise 2-0 defeat to

AL-AIN, United Arab Emirates 1996-12-06 Japan began the defence of their Asian Cup title with a lucky 2-1 win against Syria in a Group C championship But China saw their luck desert them in the second match of the group, crashing to a surprise 2-0 defeat to

Japan began the defence of their Asian Cup title with a lucky 2-1 win against Syria in a Group C championship But China saw their luck desert them in the second match of the group , crashing to a surprise 2-0 defeat to

