

Key Information Extraction from documents: Kleister NDA/Charity challenges

Filip Graliński
[filip.gralinski@applica.ai](mailto:fيلip.gralinski@applica.ai)

Aplica.ai / Adam Mickiewicz University

March, 15th

Kleister: A novel task for Information Extraction involving Long Documents with Complex Layout

Filip Graliński^{1,2}, Tomasz Stanisławek^{1,4}, Anna Wróblewska^{1,4},
Dawid Lipiński¹, Agnieszka Kaliska^{1,2}, Paulina Rosalska^{1,3},
Bartosz Topolski¹, Przemysław Biecek^{4,5}

¹Applica.ai, 15 Zajęcza, Warsaw, 00351, Poland, firstname.lastname@applica.ai

²Adam Mickiewicz University, 1 Wieniawskiego, Poznań, 61712, Poland, firstname.lastname@amu.edu.pl

³Nicolaus Copernicus University, 11 Gagarina, Toruń, 87100, Poland, firstname.lastname@umk.pl

⁴Warsaw University of Technology, Koszykowa 75, Warsaw, Poland, firstname.lastname@pw.edu.pl

⁵Samsung R&D Institute Poland, Plac Europejski 1, Warsaw, Poland, firstname.lastname@samsung.com

Abstract

State-of-the-art solutions for Natural Language Processing (NLP) are able to capture a broad range of contexts, like the sentence-level context or document-level context for short documents. But these solutions are still struggling when it comes to longer, real-world documents with the information encoded in the spatial structure of the document, such as page elements like tables, forms, headers, openings or footers; complex page layout or presence of multiple pages.

To encourage progress on deeper and more complex Information Extraction (IE) we introduce a new task (named *Kleister*) with two new datasets. Utilizing both textual and structural layout features, an NLP system must find the most important information, about various types of entities, in long formal documents. We propose *Pipeline* method as a text-only baseline with different Named Entity Recognition architectures (Flair, BERT, RoBERTa). Moreover, we checked the most popular PDF processing tools for text extraction (pdf2djvu, Tesseract and Textract) in order to analyze behavior of IE system in presence of errors introduced by these tools.

1 Introduction

Information Extraction (IE) requires quick but careful skimming through the whole document. We often have to not only search for pieces of information, but also to generate final output for specific entity type (e.g. aggregate multiple occurrences of organization names into one). In practice, this means that the results should be presented in an appropriate form (e.g. data points such as addresses normalized to



Information Extraction Challenges

- ▶ Applica Kleister
 - ▶ Kleister NDA
 - ▶ Kleister Charity
- ▶ PolEval 2020 Annual Reports

SPRAWOZDANIE Z DZIAŁALNOŚCI GRUPY KAPITAŁOWEJ ZM „ROPCZYCE” S.A. za I półrocze 2012 r.

1. INFORMACJE O SPÓŁKACH WCHODZĄCYCH W SKŁAD GRUPY KAPITAŁOWEJ

JEDNOSTKA DOMINUJĄCA – PREZENTACJA SPÓŁKI



Zakłady Magnezytowe „ROPCZYCE” S.A. (ZMR S.A.)

Siedziba: Ropczyce, woj. podkarpackie

Adres: ul. Przemysłowa 1, 39-100 Ropczyce

Regon: 690026060

NIP: 818-00-02-127

www.ropczyce.com.pl

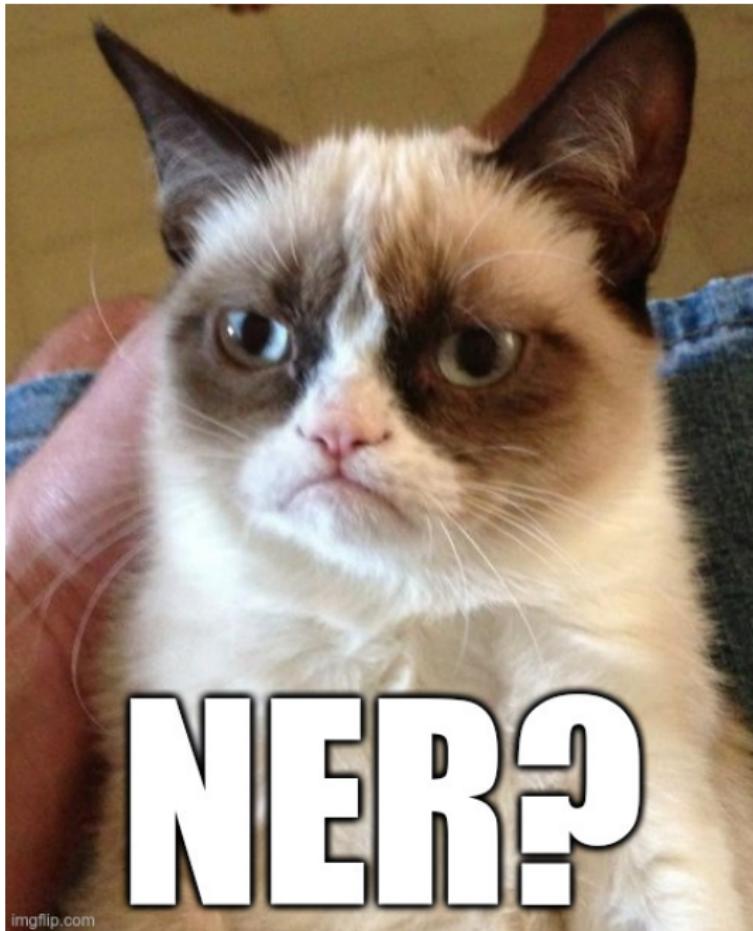
PRZEDMIOT DZIAŁALNOŚCI

Przedmiot działalności ZMR S.A. obejmuje produkcję i sprzedaż zasadowych wyrobów o których są niezbędnym elementem konstrukcji wyłożen pieców i urządzeń cieplnych przy wysokich temperaturach, głównie w hutnictwie żelaza i stali, hutnictwie metali nieżelaznych przemysle cementowo-wapienniczym, odlewniczym.

Spółka świadczy także usługi w zakresie nawęglania i ulepszania cieplnego wyrobów oraz prowadzenia badań badawczo-rozwojowe w dziedzinie związanej z przedmiotem jej działalności.

period_from 2012-01-01
period_to 2012-06-30
postal_code 39-100
city Ropczyce
...

company, drawing_date,
period_from, period_to,
postal_code, city, street,
street_no, people



imgflip.com

... no!

This is a **Key Information Extraction** task, not NER*.

- ▶ we are interested in the information not where it is
- ▶ not just any person, but CEO, etc.

* But of course you could use NER as a part of the pipeline

TRUSTEES' REPORT (INCLUDING DIRECTORS' REPORT AND STRATEGIC REPORT)

The Trustees, who are directors of the Charity, present their annual report on the affairs of the Charity and the Group, together with the financial statements and auditor's report for the year ending 31 March 2018. The Trustees have adopted the provisions of the statement of recommended practice (SORP) "Accounting and Reporting by Charities" (FRS 102) in preparing the annual report and financial statements of the Charity.

NOTES

COMPANY NUMBER

4311665

DIRECTOR'S STATEMENT

DIRECTOR'S STATEMENT

Director(s) appointed: 01/01/2018

AUDITORS

Kirklees Active Leisure Limited

37 Victoria Place

Huddersfield, HD1 2AR

BANKERS

Barclays Bank plc

17 Victoria Place

Huddersfield, HD1 2AB

STRUCTURE, GOVERNANCE AND MANAGEMENT

General information
 Kirklees Active Leisure ("KAL") was formed as a Company Limited by Guarantee, not having share capital and having charitable status, on 29 November 2001. The Charity operates community recreation facilities and part of the Kirklees Leisure Group which wholly owns KAL. Kirklees Active Leisure Ltd ("KALT"). The Charity is required to comply with both the Companies Act 2006 and the Charities Act 2011. Accounting and Reporting by Charities (FRS 102) and other more general Charity Commission regulations.

The Memorandum and Articles of Association are the Charity's constitution.

CHAIRPERSON

D. Stephenson

TRUSTEES

J.A. Fringe

M.T. Steele

S.S. Khan

D.C. Hulme

C.R. W.J. Doddle (resigned 25/7/18)

A.N. Fletcher

C.R. Solal

S.M. Sopta

E.C. Morris

J.S. Foster

D. Morley

C.R. M.S. Thompson (Appointed 23/3/18)

KIRKLEES ACTIVE LEISURE
CONSOLIDATED STATEMENT OF FINANCIAL ACTIVITIES (including Income and Expenditure Account)

for the year ended

31 MARCH 2018

Year	Unrestricted funds £	2018		2017	
		Revised funds £	Total funds £	Unrestricted funds £	Total funds £
Income and endowments from:					
2. Charitable activities	12,953,044	-	12,953,044	13,256,540	-
3. Other trading activities	102,119	-	102,119	213	-
4. Investments	5,043	-	5,043	6,467	-
5. Other	1,983,137	72,442	1,983,579	2,312,285	-
Total income	14,623,197	137,843	14,760,040	15,812,535	-
Expenditure etc:					
6. Repairs and maintenance	122,119	-	122,119	168,517	-
7. Charitable activities	16,16,204	122,508	16,16,712	15,899,468	-
Total expenditure	16,18,343	122,508	16,44,851	16,069,445	-
Net income/(expenditure)	(1,66,346)	35,434	(16,44,851)	(16,069,445)	-
24. Other recognised gains/(losses) on defined benefit pension scheme	1,299,808	-	1,299,808	(2,569,000)	-
Net movement in funds before	(39,348)	35,434	(160,112)	(3,765,957)	-
Corporation tax charge on voluntary	(354)	-	(354)	(356)	-
25. Net movement in funds after tax	(36,696)	35,434	(160,466)	(3,766,316)	-
26. Total funds brought forward	(6,779,297)	391,859	(6,487,738)	(2,703,637)	-
27. Total funds brought forward	(7,215,400)	437,293	(5,860,494)	(6,887,748)	-

All income and expenditure derive from continuing activities.

The statement of financial activities includes all gains and losses recognised during the year.



Trustees' Annual Report for the period

From 01/04/2017 To 31/03/2018

Section A Reference and administration details

CHARITY NAME	The Al Perry Charitable Foundation
Other names charity is known by	
REGISTERED CHARITY NUMBER (if any)	90073
CHARITY'S PRINCIPAL ADDRESS	Guthrie Bank Plc, Treaders Dispensary, 8th Floor, 111 Greyfriars, London, EC2R 8AF
Postcode	EC2R 0FF

Names of the charity trustees who manage the charity

Trustee name	Office (if any)	Dates acted if not for whole year	Name of person (if being replaced) to whom replaced (if any)
1. NatWest Bank Plc			
2.			
3.			
4.			
5.			
6.			
7.			
8.			
9.			
10.			
11.			
12.			
13.			
14.			
15.			
16.			
17.			
18.			
19.			
20.			

Names of the trustees for the charity, if any, (for example, any custodian trustees)

Name	Dates acted if not for whole year

TAR

April 2008



Change lives. For good.
action:aid

	The Al Perry Charitable Foundation	90073	CC16e
Receipts and payments accounts			
For the period	01/04/2017	To	31/03/2018

Section A Receipts and payments

	Unrestricted funds in the account	Restricted funds in the account	Endowment funds in the account	Total funds in the account	Last year
All Receipts					
Income from donations	21,000	-	-	21,000	20,000
Taxes imposed	-	-	-	-	11
Interest	-	-	-	-	11
Sub total	21,000	-	-	21,000	20,000
A2 Asset and investment sales, etc.					
	-	-	-	-	200,280
Total receipts	21,000	-	-	21,000	200,280
A3 Payments					
Charitable expenses	24,000	-	-	24,000	17,000
Administrative expenses	7,000	-	-	7,000	7,000
Professional fees	5,000	-	-	5,000	5,000
Independent Examiner Fee	400	-	-	400	375
Sub total	32,400	-	-	32,400	27,375
A4 Asset and investment purchases, etc.					
	-	-	-	-	100,000
Total payments	32,400	-	-	32,400	100,000
Net of receipts/(payments)					
AS Transfers between funds	9,000	-	-	9,000	1,000
AS Change funds last year	0	-	-	0	1,000
AS Current year	9,000	-	-	9,000	1,000
Sub total	9,000	-	-	9,000	1,000
AN Asset and investment purchases, etc.					
	-	-	-	-	200,280
Total payments	32,400	-	-	32,400	200,280
Net of receipts/(payments)					
AS Transfers between funds	11,000	-	-	11,000	1,000
AS Change funds last year	0	-	-	0	1,000
AS Current year	11,000	-	-	11,000	1,000
Sub total	11,000	-	-	11,000	1,000
Net assets at the beginning of the year					
	10,000	-	-	10,000	10,000
Total assets carried forward					
	9,000	9,007	10,000	9,007	10,000
Net assets at the end of the year					
	9,000	9,007	10,000	9,007	10,000

Non	Restricted funds	Investment funds	Charitable funds	Administrative funds	Unrestricted funds	Total funds
Income:						
- Donations and legacies	2a	14,000	23,207	48,217	10,700	27,360
- Trading	2a	-	-	-	11	11
- Investment income	2a	-	902	902	-	79
Income from charitable activities:						
- Grants	2b	11,022	616	10,100	10,518	10,135
-raising income from charitable activities	2b	-	2,007	2,007	-	11
Total income	28,443	29,004	61,037	61,037	20,308	59,436
Expenses:						
-Repairs and maintenance	3a	14,805	12,020	9,143	20,006	49,936
-Running funds	4	229	8,005	8,123	77	8,240
-Funding/trading costs: costs of goods sold and other costs	4	-	33	33	-	5
Charitable activities	6	30,458	18,203	11,143	33,006	64,802
Total expenditure	56,079	57,235	40,143	38,073	21,047	59,436
Net movement in funds						
Net movement in funds	(27,636)	(7,231)	(10,000)	(18,036)	(1,025)	(18,036)
Net assets at the beginning of the year						
	10,000	10,000	10,000	10,000	10,000	10,000
Total assets carried forward						
	9,000	9,007	10,000	9,007	10,000	10,000
Net assets at the end of the year						
	9,000	9,007	10,000	9,007	10,000	10,000

Notes on pages 40–48 form part of these financial statements. There are no recognised gains and losses.



Dataset name	CoNLL 2003	WikiReading	FUNSD	SROIE	Kleister NDA	Kleister Char- ity
Source	Reuters news Documents Pages Entities	Wikipedia 1,393 — 35,089	forms 4.7M — 18M	receipts 199 199 9,743	EDGAR 540 3,229 2,160	UK Com. 2,778 61,643 21,612
train docs	946	16.03M	149	626	254	1,729
dev docs	216	1.89M	—	—	83	440
test docs	231	0.95M	50	347	203	609
Input/Output on token level	✓	✗	✓	✓	✗	✗
Long Document	✗	✓	✗	✗	✓	✓
Complex layout	✗	✗	✓	✓	✓	✓
OCR	✗	✗	✓	✓	✗	✓

Table: Summary of the existing English datasets and the Kleister sets.

Entities	Description
<i>NDA dataset</i>	
party	parties appearing in the agreement (each of them is treated as a separate entity)
jurisdiction	state or country whose law governs the agreement
effective_date	date on which the contract becomes legally binding
term	duration of the agreement
<i>Charity dataset</i>	
address__post_town	post town (part of a charity address)
address__postcode	postcode (part of a charity address)
address__street_line	street with the house number (part of a charity address)
charity_name	name of the charitable organization
charity_number	identification number in the charity register
report_date	date of reporting
income_annually	annual income in British pounds (GBP)
spending_annually	annual spending in British pounds (GBP)

Evaluation metric

F1-score will be used as the evaluation metric

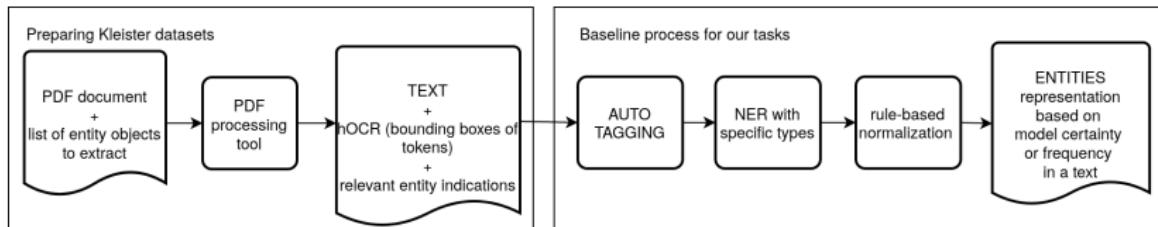
--metric MultiLabel-F1

Possible approaches

- ▶ handcrafted rule, e.g. regexps (as a baseline)
- ▶ standard NER (for general entities) + role classification
- ▶ specialized NER (but you need to **autotag** entities)
- ▶ end-to-end (generative models)
- ▶ ensembles

Łukasz Garncarek, Rafał Powalski, Tomasz Stanisławek, Bartosz Topolski, Piotr Halama, Filip Graliński, *LAMBERT: Layout-Aware (Language) Modeling using BERT for information extraction*, <https://arxiv.org/abs/2002.08087>

Autotagging approach



Kleister-NDA dataset (pdf2djvu)

Entity name	Flair	BERT	RoBERTa	LayoutLM	LAMBERT	Autotag.	Human
effective_date	79.37	80.20	81.50	82.08	85.27	79.00	100 %
party	70.13	71.60	80.83	75.28	78.70	33.15	98 %
jurisdiction	93.87	95.00	92.87	94.40	96.50	54.10	100 %
term	60.33	45.73	52.27	48.34	55.03	74.10	95 %
ALL	77.83	78.20	81.00	78.68	81.77	60.09	97.86 %

Kleister-Charity dataset (Azure CV)

post_town	83.30	77.03	77.70	79.97	81.03	66.04	98 %
postcode	82.63	87.10	88.40	81.06	82.97	87.60	100 %
street_line	68.17	62.23	72.03	70.92	75.33	75.02	96 %
charity_name	72.40	75.93	78.03	78.82	79.10	67.00	99 %
charity_number	96.73	96.67	95.37	95.76	96.57	98.60	98 %
income	70.93	64.43	69.73	72.86	76.90	69.00	97 %
report_date	95.67	96.60	96.77	95.42	95.80	89.00	100 %
spending	61.67	67.30	68.60	71.20	74.33	73.00	92 %
ALL	80.10	78.33	81.50	80.74	82.97	78.16	97.45 %

Comparing OCR engines

Kleister NDA dataset (born-digital PDF files)					
PDF tool	Flair	BERT	RoBERTa	LayoutLM	LAMBERT
Azure CV	78.03 \pm 0.12	77.67 \pm 0.18	79.33 \pm 0.68	77.43 \pm 0.29	80.57 \pm 0.25
pdf2djvu	77.83 \pm 0.26	78.20 \pm 0.17	81.00 \pm 0.05	78.47 \pm 0.76	81.77\pm0.09
Tesseract	76.57 \pm 0.49	76.60 \pm 0.30	77.81 \pm 0.97	77.70 \pm 0.48	81.03 \pm 0.23
Texttract	77.37 \pm 0.08	74.83 \pm 0.45	79.49 \pm 0.32	77.40 \pm 0.40	77.37 \pm 0.08

Kleister Charity dataset (mixture of born-digital and scanned PDFs) (*)					
PDF tool	Flair	BERT	RoBERTa	LayoutLM	LAMBERT
Azure CV	81.17 \pm 0.12	78.33 \pm 0.08	81.50 \pm 0.23	81.53 \pm 0.23	83.57\pm0.29
Tesseract	72.87 \pm 0.81	71.37 \pm 1.25	76.23 \pm 0.15	77.53 \pm 0.20	81.50 \pm 0.07
Texttract	78.03 \pm 0.12	73.30 \pm 0.43	80.08 \pm 0.15	80.23 \pm 0.41	82.97 \pm 0.21

Table: F_1 -scores for different PDF processing tools and models checked on Kleister challenges test sets over 3 runs with standard deviation. (*) pdf2djvu does not work on scans. We used the Base version of the BERT, RoBERTa, LayoutLM and LAMBERT models.

Financial Review

The table below summarises the group's financial performance:

	2017/18 £'000	2016/17 £'000
INCOME		
From customers	12,861	13,229
From customers of subsidiary	183	213
	13,044	13,442
Investment income	5	7
Funding by Kirklees MC	1,656	2,352
Other funding	2	21
	14,722	15,852
EXPENDITURE		
Staff costs, including self-employed instructors	9,241	8,925
Other costs incurred by KAL	6,024	6,414
Costs incurred by KALT	54	90
	15,319	15,429
EXCESS OF EXPENDITURE OVER INCOME, BEFORE THE EFFECTS OF THE DEFINED BENEFIT PENSION SCHEME	(517)	421

2017/18 was an especially challenging year financially with a huge reduction in funding support from the local authority of £750,000, while significant local competitors to the fitness market also opened new facilities. Given that the local economy remains difficult and wage growth relatively flat, the Charity faced considerable financial challenge, resulting in an outcome position below the planned budget position.

The group's income for the year was £14,722,729 (2016/17 £15,852,553), of which £12,938,544 (2016/17 £13,280,540) was generated through charitable training activities. The trading subsidiary contributed a further £182,573 (2016/17 £213,321) from its activities.

Total resources expended amounted to £16,440,851 (2016/17 £16,069,445) including costs incurred by the trading subsidiary of £53,697 (2016/17 £89,861).

Before the actuarial effects of the defined benefit pension scheme on the group, net resources expended amounted to £1,659,112 (2016/17 £16,892).

Performance as a function of document's length

BERT Flair LAMBERT LayoutLM RoBERTa

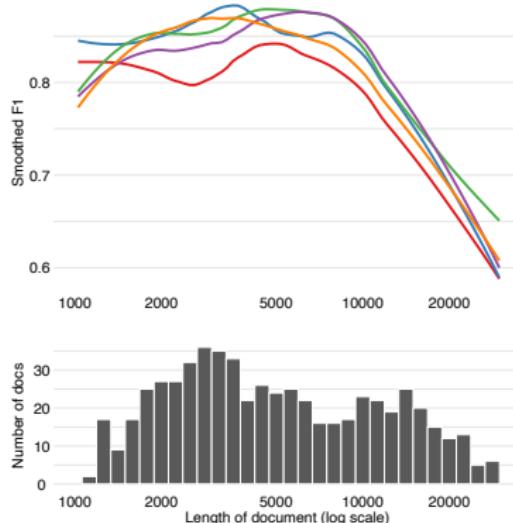


Figure: Normalization issues for an income entity (amount in the table should be multiplied by 1000).

Figure: Relationship between F_1 -scores and document length in the Kleister Charity test set for the Azure CV OCR.

Aplica Kleister challenges

<https://gitlab.com/kleister-challenge-2021/kleister-nda.git>

<https://gitlab.com/kleister-challenge-2021/kleister-charity.git>

Thank you!