

# The Financial Document Structure Extraction Shared task: FinToc 2020

**Najah-Imane Bentabet**

Fortia Financial Solutions  
Paris, France

najah-imane.bentabet@fortia.fr

**Rémi Juge**

Fortia Financial Solutions  
Paris, France

remi.juge@fortia.fr

**Ismail El Maarouf**

Fortia Financial Solutions  
Paris, France

ismail.elmaarouf@fortia.fr

**Virginie Mouilleron**

Fortia Financial Solutions  
Paris, France

virginie.mouilleron@fortia.fr

**Dialekti Valsamou-Stanislowski**

Fortia Financial Solutions  
Paris, France

dialekti.valsamou@fortia.fr

**Mahmoud El-Haj**

Lancaster University  
Lancaster, UK

m.el-haj@lancaster.ac.uk

## Abstract

This paper presents the FinTOC-2020 Shared Task on structure extraction from financial documents, its participants results and their findings. This shared task was organized as part of The 1st Joint Workshop on Financial Narrative Processing and MultiLing Financial Summarisation (FNP-FNS 2020), held at The 28th International Conference on Computational Linguistics (COLING'2020). This shared task aimed to stimulate research in systems for extracting table-of-contents (TOC) from investment documents (such as financial prospectuses) by detecting the document titles and organizing them hierarchically into a TOC. For the second edition of this shared task, two subtasks were presented to the participants: one with English documents and the other one with French documents.

## 1 Introduction

The use of PDF electronic documents is recurrent in the financial domain. They are used to share and broadcast information concerning investment strategies, policy and regulation. Even with a great layout, long documents can be hard to navigate, hence, the presence of a table-of-contents (TOC) can provide a valuable assistance for potential investors or regulators by increasing readability and facilitating navigation.

In this shared task, we focus on extracting the TOC of financial prospectuses. In these official documents, investment funds accurately depict their characteristics and investment modalities. Depending on their country of origin, they might be edited with or without a TOC, and they might follow a template as well. But even though their format is regulated, the choice of the text format, the layout, the graphics and tabular presentation of the data is in the hand of the editor. Thus, the TOC is of fundamental importance to tackle sophisticated NLP tasks such as information extraction or question answering on long documents.

In this paper, we report the results and findings of the FinTOC-2020 shared task.<sup>1</sup> The Shared Task was organized as part of The 1st Joint Workshop on Financial Narrative Processing and MultiLing Financial Summarisation (FNP-FNS 2020), to be held at The 28th International Conference on Computational Linguistics (COLING'2020).

A total of 5 teams submitted runs and contributed 5 system description papers. All system description papers are included in the FNP-FNS 2020 workshop proceedings and cited in this report.

<sup>1</sup><http://wp.lancs.ac.uk/cfie/fintoc2020/>

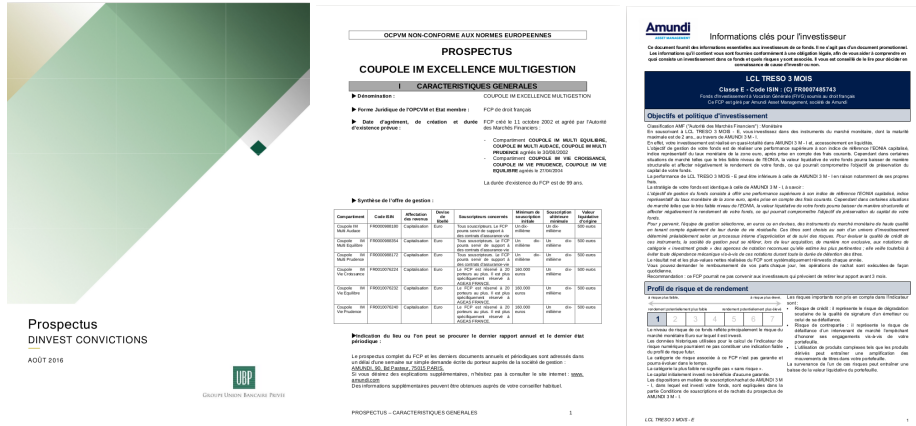


Figure 1: Random pages from the shared task datasets. We observe a strong variability of complex layouts.

## 2 Previous Work on TOC extraction

There are mainly two concepts in the literature to approach TOC extraction. The first one parses the hierarchical structure of sections and subsections from the TOC pages embedded in the document. This area of research was mostly motivated by the INEX (Drešević et al., 2009) and ICDAR competitions (Doucet et al., 2013; Beckers et al., 2010; Nguyen et al., 2018) which aim at extracting the TOC of old and lengthy OCR-ised books. The documents we target in this shared task are very different: they contain graphical elements, and the text is not displayed to respect a linear reading direction but is optimized to condense information and catch the eye of the reader. Apart from these competitions, we find the methods proposed by El-Haj et al (El Haj et al., 2014; El Haj et al., 2019b; El-Haj et al., 2019a), also based on the parsing of the TOC page.

In the second category of approaches, we find algorithms that detect the titles of the document using learning methods based on layout and text features. The set of titles is then hierarchically ordered according to a predefined rule-based function (Doucet et al., 2013; Liu et al., 2011; Gopinath et al., 2018).

Lately, we find systems that address the hierarchical ordering of the titles as a sequence labelling task, using neural networks models such as Recurrent Neural Networks and LSTM networks (Bentabet et al., 2019).

## 3 Task Description

As part of the FNP-FNS Workshop, we present a shared task on Financial Document Structure Extraction.

Participants to this shared task were given two sets of financial prospectuses with a wide variety of document structure and length. Their systems had to automatically process the documents to extract their document structure, or TOC. In fact, the two sets were specific to two different subtasks:

- **TOC extraction from French documents:** The set of French documents is rather homogeneous in terms of structure, due to the existence of a common template. However, the words and phrasing can differ from one prospectus to another. Also, French prospectuses never include a TOC page that could be parsed.
- **TOC extraction from English documents:** English prospectuses are characterized by a wide variety of structures as there is no template to constrain their format. Contrary to the French documents, there is always a TOC page but the latter is usually highly incomplete as only the higher level section titles are displayed.

For both sets, we observe that:

- some documents contain specific titles that do not appear in any other document
- the same title in two different documents can have a different position in the hierarchy

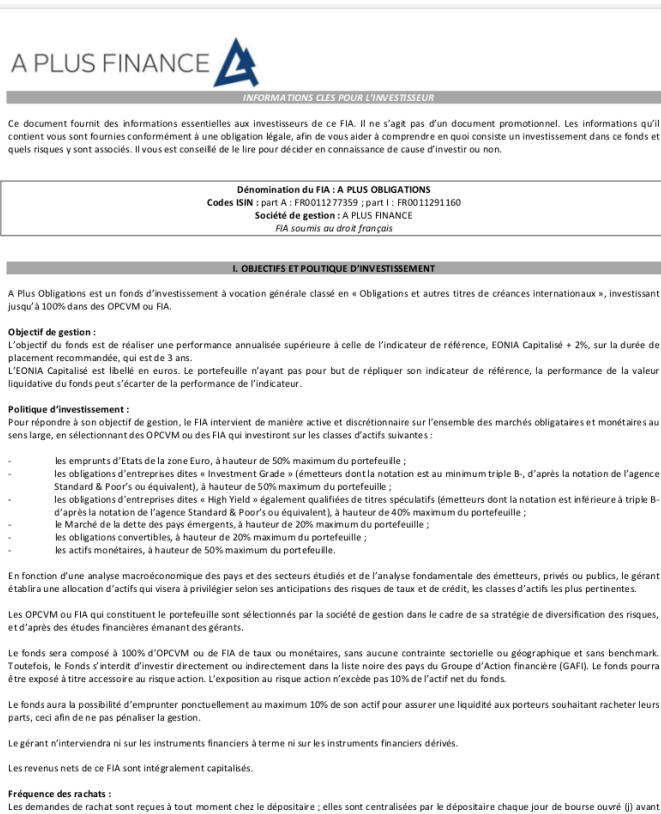


Figure 2: A French prospectus with its JSON annotation file.

- two titles that follow each other can have the same layout but a different position in the TOC
- the font size of a higher-level title can be smaller than the font size of a lower-level one
- and a title can have the exact same layout as its associated paragraph.

For each subtask, all participating teams were provided with a training dataset which included the original PDFs alongside their corresponding JSON file representing the TOC of the document. This JSON represented the TOC by giving the titles, their pages, their depths and their IDs, as shown in Fig. 2. A private test set was used to evaluate the TOCs generated by the participants systems. As stated in Section 2, most of the previous research on TOC generation has focused on short papers such as research publications (*Arxiv* database), or weakly graphical material such as digitalized books. However, the task of extracting the TOC of commercial documents with a complex layout structure in the domain of finance is not much explored in the literature.

## 4 Shared Task Data

In this section, we discuss the corpus of documents used for the TOC extraction subtasks.

### 4.1 Corpus annotation

Investment documents can be accessed online in PDF format, and are also made available from asset managers. We compiled a list of 71 French documents, and 72 English documents from Luxembourg, to create the datasets of each subtask. We chose documents with a wide variety of layouts and styles. We provided annotators with the original PDFs and a software that was developed internally to manually annotate the TOC of any PDF document. Once the annotator finishes their annotation task, the software produces a file containing the TOC-entries (title, page number, depth, and id) in a hierarchically structured format.

## Techniques et instruments utilisés

### Actifs (hors dérivés intégrés)

#### - Actions

En sa qualité de SICAV éligible au PEA, le portefeuille est investi au minimum à 75 % en titres de sociétés et en parts ou actions d'OPC éligibles au PEA. Les titres de société éligibles au PEA sont ceux dont le siège social est établi dans un État membre de l'Union européenne ou dans un autre État partie à l'accord sur l'Espace économique européen (EEE) non membre de l'Union européenne ayant conclu avec la France une convention fiscale contenant une clause administrative en vue de lutter contre la fraude ou l'évasion fiscale.

Les titres sont sélectionnés selon les critères présentés dans la stratégie d'investissement.

La gestion est orientée sur le marché français. Sur opportunité, des investissements peuvent être réalisés sur des valeurs d'autres zones géographiques présentant des perspectives particulièrement attractives.

La sélection des titres s'effectue sans a priori sur la taille des sociétés. La gestion ne s'intéresse pas seulement aux principales capitalisations, même si les grandes capitalisations demeurent majoritaires dans le portefeuille. Le poids accordé aux grandes capitalisations par rapport aux capitalisations plus petites n'est pas figé, il varie en fonction des opportunités de marché et des valorisations relatives entre les différents titres.

#### - Actions ou parts d'autres placements collectifs de droit français ou d'autres OPCVM, FIA ou fonds d'investissement de droit étranger

La SICAV peut investir jusqu'à 50 % de son actif en parts ou actions d'OPCVM français ou européens ou de fonds d'investissement à vocation générale de droit français, dans des actions ou parts de fonds de capital investissement, d'OPC investissant plus de 10 % en parts ou actions d'un autre véhicule de gestion collective, d'OPC nourriciers, de fonds professionnels à vocation générale, de fonds professionnels spécialisés, de fonds d'investissement constitués sur le fondement d'un droit étranger répondant aux critères prévus aux articles R214-32-42 ou R214-13 du Code monétaire et financier ou de l'article 422-95 du Règlement général de l'AMF, ainsi que des parts ou actions de fonds de fonds alternatifs.

Les OPC sont sélectionnés afin de respecter la politique de gestion ci-dessus présentée.

Afin d'augmenter l'exposition actions ou taux, la SICAV se réserve également la possibilité d'investir dans des OPC indiciels cotés (ETF ou trackers).

La SICAV se réserve la possibilité d'acquérir des parts ou actions d'OPC érés par LA BANQUE POSTALE ASSET MANAGEMENT ou une société liée.

La sélection d'OPCVM et de fonds d'investissement non gérés par LA BANQUE POSTALE ASSET MANAGEMENT ou une société liée repose sur une analyse quantitative des performances passées ainsi que sur une analyse qualitative de leurs processus d'investissement.

Figure 3: In this example, we can see that the titles tagged in green have the same style as the plain text of their paragraphs. Only the indentation is insightful to detect them.

Each annotator was asked to:

1. Identify the title: Locate a title inside the PDF document.
2. Associate the entry level in the TOC: Every title is tagged with an integer representing the depth of the title in the TOC tree. The depth ranges from 1 to 10.
3. Tag the next title.

Each document was annotated independently by two people and a third person would review the annotations to resolve possible conflicts. For each dataset, the agreement scores between annotators are depicted in Table 1 and Table 2. We can observe high agreement scores, allowing us to be confident enough about the quality of our datasets.

	Xerox F1	Inex08 F1
tagger 1 & tagger 2	89.8%	77.0%
tagger 1 & reviewer	92.1%	82.8%
tagger 2 & reviewer	90.1%	79.6%

Table 1: Agreement scores between different annotators of the French investment document dataset (71 documents).

**Annotation Challenge: Title identification** Investment prospectuses are commercial documents whose complex layout is optimized to highlight specific information such that a potential investor can identify it

	Xerox F1	Inex08 F1
tagger 1 & tagger 2	87.7%	82.4%
tagger 1 & reviewer	95.6%	91.6%
tagger 2 & reviewer	91.8%	90.0 %

Table 2: Agreement scores between different annotators on a validation set of 62 documents from the English investment document dataset (79 documents).

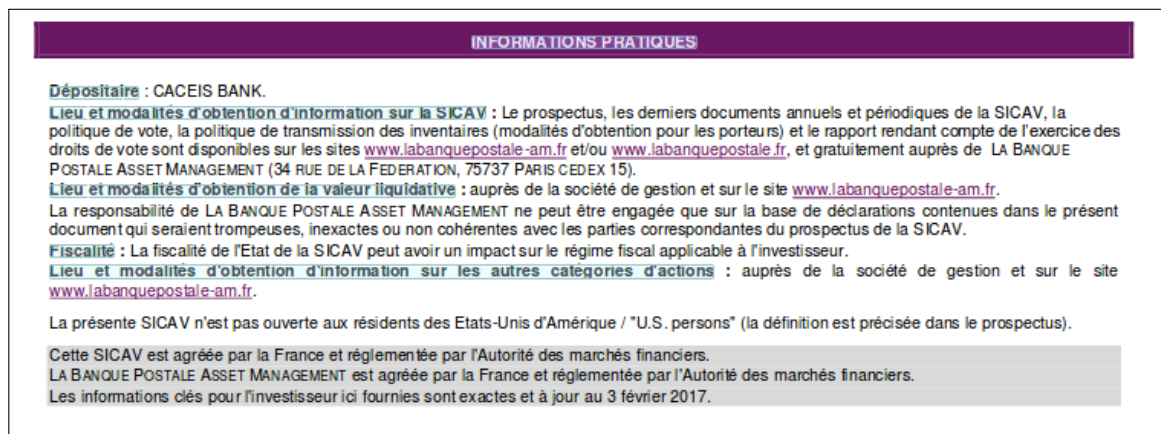


Figure 4: In this second example, the identification of titles tagged in light blue is not evident because they might be followed by plain text in the same line.

quickly. Hence, annotating a title and its level in the TOC hierarchy is a difficult task as one cannot rely on the visual appearance of the title to do so. Some examples can be observed in Fig. 3 and Fig. 4.

**Annotation Challenge: Tagging PDF documents.** The annotation of PDF documents is not an easy task since they are meant to be displayed. The tool we used for the annotations allows the annotators to directly tag on the PDF, however, the text selection relies on the HTML encoding of the PDF, where the text might slightly differ from what is actually displayed. For instance, it is possible that a piece of text is impossible to select if it is from an image. It is also possible that the tagged text has additional or missing characters.

## 4.2 Corpus Description

In the following, we provide an analysis of the data used for the shared task.

We simplified greatly the format of the annotation files compared to the first edition of the shared task (Juge et al., 2019). Instead of the XML format inherited from the Structure Extraction Competition (SEC) (Doucet et al., 2013) that implicitly encodes the title level, we used a simple JSON file containing a list of entries, where each entry has the following information: textual content, id, level, page number. An example of a JSON extract is provided in Fig. 2. In particular, the title level is explicitly stated. Statistics about levels on the French and English datasets are presented in Table 3.

In addition to the annotation files, the public dataset provided to the participants contained documents in PDF format. The private dataset on which participants were ranked contained documents in PDF format only. The french subtask (respectively the english subtask) had 47 (respectively 50) public documents. The rest was kept private for the final ranking.

## 5 Participants and Systems

A total of 50 teams registered in the shared task all from different institutions. Eventually, 5 teams participated and submitted a paper with the description of their method, see Table 5 for more information about their affiliation. In Table 4, we show the details on the submissions per task. All the participants that submitted a standard run, sent a paper describing their approach as well.

	French dataset	English dataset
number of documents	71	72
average number of pages	28	91
level 1 (% of titles)	2%	5%
level 2 (% of titles)	11%	21%
level 3 (% of titles)	29%	30%
level 4 (% of titles)	24%	25%
level 5 (% of titles)	21%	11%
level 6 (% of titles)	13%	4%
level 7 (% of titles)	0%	2%
level 8 (% of titles)	0%	1%
level 9 (% of titles)	0%	1%
level 10 (% of titles)	0%	0%

Table 3: Statistics on the subtasks datasets.

	# teams	# std runs
French subtask	4	6
English subtask	5	7
papers	5	-

Table 4: Statistics on the participation on French and English subtasks

Participating teams explored and implemented a wide variety of techniques and features. In this section, we give a brief description of each system. More details could be found in the description papers published in the proceedings of the FNP-FNS 2020 Workshop.

**AMEX-AI Labs (Premi et al., 2020):** This team participated in the English subtask only. Several pre-processing steps, including headers and footers removal, are performed to segment the textual content of the documents into elements that are classified later as titles and non-titles. They separately trained two title detectors, one on an external dataset, and the other on the English dataset provided by the FinTOC2020 shared task. Then, they concatenated both models to form their final title detector.

**Daniel (Giguet et al., 2020):** Unlike the other teams, this team focused on removing tables to clean the textual content of the documents. They detected titles by looking for numbered lines, leveraging stylistic properties, and checking the existence of a line in the list of training titles. For the hierarchy part, they clustered the titles previously detected to infer their hierarchical level using stylistic features. It is interesting to see that this approach ranked well on the English dataset but not so well on the French dataset.

**DNLP (Kosmajac et al., 2020):** The DNLP team participated in both subtasks. They used *tesseract* an open-source OCR tool to extract the text regions. Then, they defined a set of features to use with

Team	Affiliation	Tasks
AMEX-AI Labs (Premi et al., 2020)	American Express AI Labs, Bangalore	E
Daniel (Giguet et al., 2020)	STIH, Sorbonne Univeristy	F and E
DNLP (Kosmajac et al., 2020)	Dalhousie University	F and E
Taxy.io (Haase and Kirchhoff, 2020)	Taxy.io	F and E
UWB (Hercig and Král, 2020)	University of West Bohemia	F and E

Table 5: List of the 5 teams that participated in Subtasks of the FinTOC2020 Shared Task. "F" refers to the French subtask and "E" refers to the English subtask



three different algorithms: linear regression, random forest and SVM. For both title detection and TOC extraction steps, their best performing models are random forest models.

**Taxy.io (Haase and Kirchhoff, 2020):** The Taxy.io team participated in both subtasks, with a multilingual pipeline. They used an unsupervised learning approach to tackle text block detection. They first run a DBSCAN clustering on pages characters to extract features and then run a second DBSCAN clustering to identify the text blocks. Finally, they classify each text block, represented by features from the previous step plus text features extracted with a multilingual BERT model.

**UWB (Hercig and Král, 2020):** UWB team participated in both subtasks but contributed to the title detection part only. They state title detection as a binary classification on text segments, for which they use a Maximum entropy classifier, on top of a diverse set of features including orthographic characters and character n-grams.

## 6 Results and Discussion

**Evaluation Metric** Since both subtasks tackle the same problem but on different corpora, we used the same evaluation metric.

For the TOC generation part, we adapted the metrics proposed by the Structure Extraction Competition (SEC) held at ICDAR 2013 (Doucet et al., 2013): we adapted the script, replaced the customized Levenshtein distance specifically designed for SEC by a standard Levenshtein distance whose edit cost is 1 in all cases, and removed the constraint on first and last 5 characters.

The final ranking is based on the harmonic mean between *Inex F1 score* and *Inex level accuracy*. In the calculation of the *Inex F1 score*, correct entries in the predicted TOC are those which match the title of an entry in the groundtruth TOC *and* have the same page number as this entry. The *Inex level accuracy* evaluates the hierarchy of the predicted TOC. If we denote by  $E_{ok}$  an entry in the predicted TOC with a correct page number, and by  $E'_{ok}$  an entry in the predicted TOC with a correct page number *and* a correct hierarchical level, then the Inex level accuracy is:

$$\frac{\sum E'_{ok}}{\sum E_{ok}}$$

We also provided scores for the title detection part separately: we used the F1 score, and considered as correct entries the predicted entries which match the titles of groundtruth entries according to the standard Levenshtein distance.

For both parts, the threshold on the Levenshtein *score* was set to 0.85<sup>2</sup>. Moreover, the Inex scores and title F1 score are calculated for each document and then averaged over the documents of the private set to produce two performance figures per team submission: one for TOC extraction, and another for title detection (TD).

**Baseline** For comparison purposes, we implemented a simple baseline TOC extractor consisting of:

- extracting textual content from the PDF documents using `pdftohtml` utility from Poppler library<sup>3</sup>
- assigning groundtruth labels (title or non-title) to text segments by fuzzy string matching with the annotations
- vectorizing text segments into one-dimensional vectors of length 3 encoding the following features: `is_bold`, `is_italic`, `is_all_capitalized`
- training a SVM on the obtained dataset
- assigning to a predicted title the most frequent hierarchy level found in the training set

Table 6 (respectively Table 7) reports the results obtained by the participants and the baseline on TOC extraction from French documents (respectively English documents).

<sup>2</sup>The script implementing these metrics can be found here: <https://drive.google.com/file/d/1TzsS2F79af8U5F5ivDEsc9ezUTX97aeW/view?usp=sharing>

<sup>3</sup>see <https://poppler.freedesktop.org/>

Team	TD	Team	TOC
<b>UWB</b>	<b>0.81</b>	<b>DNLP</b>	<b>0.37</b>
<b>Taxy.io</b>	0.69	<b>taxy.io</b>	0.32
<b>Daniel 1</b>	0.66	<b>Baseline</b>	0.32
<b>DNLP</b>	0.64	<b>Daniel 1</b>	0.22
<b>Daniel 2</b>	0.64	<b>Daniel 2</b>	0.22
<b>Daniel 3</b>	0.64	<b>Daniel 3</b>	0.20
<b>Baseline</b>	0.57		

Table 6: Results obtained by the participants for the first FinTOC2020 subtask : TOC extraction from French documents. The title detection (TD) ranking is based on F1-score, while the Table-Of-Content (TOC) ranking is based on the harmonic mean between Inex F1 score and Inex level accuracy

Team	TD	Team	TOC
<b>Amex 1</b>	<b>0.79</b>	<b>DNLP</b>	<b>0.34</b>
<b>UWB</b>	0.77	<b>Daniel 3</b>	0.28
<b>Daniel 1</b>	0.69	<b>Daniel 2</b>	0.28
<b>Daniel 3</b>	0.63	<b>Daniel 1</b>	0.26
<b>Daniel 2</b>	0.62	<b>taxy.io</b>	0.24
<b>DNLP</b>	0.59	<b>Amex 1</b>	0.23
<b>Taxy.io</b>	0.55	<b>Amex 2</b>	0.23
<b>Baseline</b>	0.19	<b>Baseline</b>	0.18

Table 7: Results obtained by the participants for the second FinTOC2020 subtask : TOC extraction from English documents. The title detection (TD) ranking is based on F1-score, while the Table-Of-Content (TOC) ranking is based on the harmonic mean between Inex F1 score and Inex level accuracy

**Discussion.** Title detection is the easiest problem encountered in this competition. All the submitted models show a high increase of performance from the baseline. In addition, the numbers show that it is slightly easier to detect titles from French investment documents than it is from English investment documents. Clearly, supervised methods from UWB, Taxy.io, and AMEX-AI Labs perform better than heuristic methods such as the one proposed by team Daniel. Nevertheless, the supervised multi-lingual model from team Taxy.io performed well on the French documents only.

Concerning TOC extraction on French documents, we observe that the baseline, which naively assigns the most frequent label found in the training set, performs as well as the BERT model used by team Taxy.io, and that the unsupervised approach from team Daniel scores worse than the baseline. This probably indicates that the number of training documents provided is not enough for the diversity encountered among these documents, and that TOC extraction problem on this data is hard. TOC extraction on the English documents is an even harder task as can be inferred from the figures in Table 7. However, team DNLP stands out from the rest of the participants with a 6% to 11% increase in performance.

## 7 Conclusions

In this paper we presented the setup and results for the Financial Document Structure Extraction task (FinToc) 2020, organized as the 1st Joint Workshop on Financial Narrative Processing and MultiLing Financial Summarisation (FNP-FNS 2020) collocated with the 28th International Conference on Computational Linguistics(COLING’2020). A total of 50 teams registered and 5 teams participated in the shared task with a wide variety of techniques. All participating teams contributed with a paper describing their system.

This edition introduced the community to a new dataset, composed of French investment documents, and annotated for the TOC extraction problem. This dataset supplements previously released datasets



for English (Juge et al., 2019). TOC extraction for PDF documents is a realistic problem in everyday applications which explain the interest from and participation of both public universities and profit organizations.

## Acknowledgments

We would like to thank our dedicated annotators who contributed to the building of the corpora used in this Shared Task: Anais Koptient, Aouataf Djillani, and Lidia Duarte, and Fortia’s DLA team.

## References

- Thomas Beckers, Patrice Bellot, Gianluca Demartini, Ludovic Denoyer, Christopher M. De Vries, Antoine Doucet, Khairun Nisa Fachry, Norbert Fuhr, Patrick Gallinari, Shlomo Geva, Wei-Che Huang, Tereza Iofciu, Jaap Kamps, Gabriella Kazai, Marijn Koolen, Sangeetha Kutty, Monica Landoni, Miro Lehtonen, Véronique Moriceau, Richi Nayak, Ragnar Nordlie, Nils Pharo, Eric Sanjuan, Ralf Schenkel, Xavier Tannier, Martin Theobald, James A. Thom, Andrew Trotman, and Arjen P. De Vries. 2010. Report on INEX 2009. *Sigir Forum*, 44(1):38–57, June. Article disponible en ligne : <http://www.cs.otago.ac.nz/homepages/andrew/papers/2010-4.pdf>.
- Najah-Imane Bentabet, Rémi Juge, and Sira Ferradans. 2019. Table-of-contents generation on contemporary documents. In *Proceedings of ICDAR 2019*.
- Antoine Doucet, Gabriella Kazai, Sebastian Colutto, and Günter Mühlberger. 2013. Icdar 2013 competition on book structure extraction. In *Document Analysis and Recognition (ICDAR), 2013 12th International Conference on*, pages 1438–1443. IEEE.
- Bodin Dresevic, Aleksandar Uzelac, Bogdan Radakovic, and Nikola Todic. 2009. Book layout analysis: Toc structure extraction engine. In Shlomo Geva, Jaap Kamps, and Andrew Trotman, editors, *Advances in Focused Retrieval*, pages 164–171, Berlin, Heidelberg. Springer Berlin Heidelberg.
- Mahmoud El Haj, Paul Rayson, Steven Young, and Martin Walker, 2014. *Detecting document structure in a very large corpus of UK financial reports*. LREC’14 Ninth International Conference on Language Resources and Evaluation. Proceedings of the Ninth International Conference on Language Resources and Evaluation (LREC-2014) . European Language Resources Association (ELRA), Reykjavik, Iceland, pp. 1335-1338.
- Mahmoud El-Haj, Paulo Alves, Paul Rayson, Martin Walker, and Steven Young. 2019a. Retrieving, classifying and analysing narrative commentary in unstructured (glossy) annual reports published as pdf files. *Accounting and Business Research*, pages 1–29.
- Mahmoud El Haj, Paul Edward Rayson, Steven Eric Young, Paulo Alves, and Carlos Herrero Zorita, 2019b. *Multilingual Financial Narrative Processing: Analysing Annual Reports in English, Spanish and Portuguese*. World Scientific Publishing, 2.
- Emmanuel Giguet, Gael Lejeune, , and Jean-Baptiste Tanguy. 2020. Daniel@fintoc’2 shared task: Title detection and structure extraction. In *The 1st Joint Workshop on Financial Narrative Processing and MultiLing Financial Summarisation of COLING 2020*.
- Abhijith Athreya Mysore Gopinath, Shomir Wilson, and Norman Sadeh. 2018. Supervised and unsupervised methods for robust separation of section titles and prose text in web documents. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 850–855.
- Frederic Haase and Steffen Kirchhoff. 2020. Taxy.io@fintoc’2: Multilingual document structure extraction using transfer learning. In *The 1st Joint Workshop on Financial Narrative Processing and MultiLing Financial Summarisation of COLING 2020*.
- Tomas Hercig and Pavel Král. 2020. Uwb@fintoc-2020 shared task: Financial document title detection. In *The 1st Joint Workshop on Financial Narrative Processing and MultiLing Financial Summarisation of COLING 2020*.
- Remi Juge, Imane Bentabet, and Sira Ferradans. 2019. The FinTOC-2019 shared task: Financial document structure extraction. In *Proceedings of the Second Financial Narrative Processing Workshop (FNP 2019)*, pages 51–57, Turku, Finland, September. Linköping University Electronic Press.
- Dijana Kosmajac, Stacey Taylor, and Mozghan Saeidi. 2020. Table of contents detection in financial documents. In *The 1st Joint Workshop on Financial Narrative Processing and MultiLing Financial Summarisation of COLING 2020*.

- Caihua Liu, Jiajun Chen, Xiaofeng Zhang, Jie Liu, and Yalou Huang. 2011. Toc structure extraction from ocr-ed books. In *International Workshop of the Initiative for the Evaluation of XML Retrieval*, pages 98–108. Springer.
- Thi Tuyet Hai Nguyen, Antoine Doucet, and Mickael Coustaty. 2018. Enhancing table of contents extraction by system aggregation. In *Proceedings of the International Conference on Document Analysis and Recognition, ICDAR*.
- Dhruv Premi, Amogh Badugu, and Himanshu Sharad Bhatt. 2020. Amex-ai-labs: Investigating transfer learning for title detection in table of contents generation. In *The 1st Joint Workshop on Financial Narrative Processing and MultiLing Financial Summarisation of COLING 2020*.