

SORBONNE UNIVERSITÉ UNIVERSITÉ

ECOLE DOCTORALE INFORMATIQUE, TÉLÉCOMMUNICATIONS ET ÉLECTRONIQUE - ED130

INRIA DE PARIS / ÉQUIPE ALMANACH

THÈSE DE DOCTORAT

Discipline : Informatique

Présentée par

Pedro Javier ORTIZ SUÁREZ

Dirigée par

Laurent ROMARY et Benoît SAGOT

Pour obtenir le grade universitaire de

DOCTEUR de SORBONNE UNIVERSITÉ

On Modern NLP Methods for Document Enriching

Présentée et soutenue publiquement le 32 octobre 2021 devant le jury composé de :

Alonzo CHURCH	Princeton University	Examineur
Margaret HAMILTON	University of Michigan	Rapporteur
Emmy NOETHER	Georg-August-Universität Göttingen	Examineur
Laurent ROMARY	Inria - ALMANACH	Directeur
Benoît SAGOT	Inria - ALMANACH	Co-directeur
Claude SHANNON	MIT	Examineur
Alan TURING	Princeton University	Rapporteur

ABSTRACT

Scientific documents often use \LaTeX for typesetting. While numerous packages and templates exist, it makes sense to create a new one. Just because.

CONTENTS

1	INTRODUCTION	1
1.1	Why?	1
1.2	How?	1
1.3	Features	2
1.3.1	Typesetting mathematics	2
1.3.2	Typesetting text	3
1.4	Changing things	3
I	DATA	5
2	OSCAR	7
2.1	goclassy	7
3	MODERN FRENCH DATA	9
3.1	LEM17	9
3.2	presto max	9
3.3	presto gold	9
4	ANCIENT/MEDIEVAL FRENCH DATA	11
4.1	BERTrade Corpus	11
5	OTHER DATA	13
II	MODELS	15
6	CAMEMBERT	17
7	FRELMo	19
8	D’ALEMBERT	21

Contents

9	BERTRADE	23
III DOWNSTREAM TASKS		25
10	PARSING	27
11	POS TAGGING	29
12	NAMED-ENTITY RECOGNITION	31
13	TEXT NORMALIZATION	33
14	DOCUMENT STRUCTURATION	35
IV REAL WORLD APPLICATION		37
15	LE PETIT LAROUSSE	39
16	BASNUM	41
17	SCIENTIFIC PAPERS	43
18	MODERN FRENCH TREEBANK	45
19	NAMED-ENTITY RECOGNITION CORPORA	47
BIBLIOGRAPHY		51

1 INTRODUCTION

In which the reasons for doing this Ph.D. are laid bare for the whole world to see and we encounter some answers to questions in which, frankly, only an extremely small number of people were interested in the first place.

This package contains a minimal, modern template for writing your thesis. While originally meant to be used for a Ph. D. thesis, you can equally well use it for your honour thesis, bachelor thesis, and so on—some adjustments may be necessary, though.

1.1 WHY?

I was not satisfied with the available templates for \LaTeX and wanted to heed the style advice given by people such as Robert Bringhurst [1] or Edward R. Tufte [2, 3]. While there *are* some packages out there that attempt to emulate these styles, I found them to be either too bloated, too playful, or too constraining. This template attempts to produce a beautiful look without having to resort to any sort of hacks. I hope you like it.

1.2 HOW?

The package tries to be easy to use. If you are satisfied with the default settings, just add

```
\documentclass{mimosis}
```

at the beginning of your document. This is sufficient to use the class. It is possible to build your document using either \LaTeX , \XeTeX , or \LuaTeX . I personally prefer one of the latter two because they make it easier to select proper fonts.

Package	Purpose
<code>amsmath</code>	Basic mathematical typography
<code>amsthm</code>	Basic mathematical environments for proofs etc.
<code>booktabs</code>	Typographically light rules for tables
<code>bookmarks</code>	Bookmarks in the resulting PDF
<code>dsfont</code>	Double-stroke font for mathematical concepts
<code>graphicx</code>	Graphics
<code>hyperref</code>	Hyperlinks
<code>multirow</code>	Permits table content to span multiple rows or columns
<code>paralist</code>	Paragraph ('in-line') lists and compact enumerations
<code>scrlayer-scrpage</code>	Page headings
<code>setspace</code>	Line spacing
<code>siunitx</code>	Proper typesetting of units
<code>subcaption</code>	Proper sub-captions for figures

Table 1.1: A list of the most relevant packages required (and automatically imported) by this template.

1.3 FEATURES

The template automatically imports numerous convenience packages that aid in your typesetting process. [Table 1.1](#) lists the most important ones. Let's briefly discuss some examples below. Please refer to the source code for more demonstrations.

1.3.1 TYPESETTING MATHEMATICS

This template uses `amsmath` and `amssymb`, which are the de-facto standard for typesetting mathematics. Use numbered equations using the `equation` environment. If you want to show multiple equations and align them, use the `align` environment:

$$V := \{1, 2, \dots\} \tag{1.1}$$

$$E := \{(u, v) \mid \text{dist}(p_u, p_v) \leq \epsilon\} \tag{1.2}$$

Define new mathematical operators using `\DeclareMathOperator`. Some operators are already pre-defined by the template, such as the distance between two objects. Please see the template for some examples. Moreover, this template contains a correct differential operator. Use `\diff` to typeset the differential of integrals:

$$f(u) := \int_{v \in \mathbb{D}} \text{dist}(u, v) \, \mathrm{d}v \tag{1.3}$$

You can see that, as a courtesy towards most mathematicians, this template gives you the possibility to refer to the real numbers \mathbb{R} and the domain \mathbb{D} of some function. Take a look at the source for more examples. By the way, the template comes with spacing fixes for the automated placement of brackets.

1.3.2 TYPESETTING TEXT

Along with the standard environments, this template offers `para\list` for lists within paragraphs. Here's a quick example: The American constitution speaks, among others, of (i) life (ii) liberty (iii) the pursuit of happiness. These should be added in equal measure to your own conduct. To typeset units correctly, use the `siunitx` package. For example, you might want to restrict your daily intake of liberty to 750 mg.

Likewise, as a small pet peeve of mine, I offer specific operators for *ordinals*. Use `\th` to typeset things like July 4th correctly. Or, if you are referring to the 2nd edition of a book, please use `\nd`. Likewise, if you came in 3rd in a marathon, use `\rd`. This is my 1st rule.

1.4 CHANGING THINGS

Since this class heavily relies on the `scrbook` class, you can use *their* styling commands in order to change the look of things. For example, if you want to change the text in sections to **bold** you can just use

```
\setkomafont{sectioning}{\normalfont\bfseries}
```

at the end of the document preamble—you don't have to modify the class file for this. Please consult the source code for more information.

PART I

DATA

2 OSCAR

2.1 GOCLASSY

3 MODERN FRENCH DATA

3.1 LEM17

3.2 PRESTO MAX

3.3 PRESTO GOLD

4 ANCIENT/MEDIEVAL FRENCH DATA

4.1 BERTRADE CORPUS

5 OTHER DATA

PART II

MODELS

6 CAMEMBERT

7 FrELMo

8 D'ALEMBERT

9 BERTRADE

PART III

DOWNSTREAM TASKS

10 PARSING

11 POS TAGGING

12 NAMED-ENTITY RECOGNITION

13 TEXT NORMALIZATION

14 DOCUMENT STRUCTURATION

PART IV

REAL WORLD APPLICATION

15 LE PETIT LAROUSSE

16 BASNUM

17 SCIENTIFIC PAPERS

18 MODERN FRENCH TREEBANK

19

NAMED-ENTITY RECOGNITION CORPORA

BIBLIOGRAPHY

1. R. Bringhurst. *The Elements of Typographic Style*. 4th ed. Hartley & Marks Publishers, Vancouver, British Columbia, Canada, 2012.
2. E. R. Tufte. *Envisioning information*. Graphics Press, Cheshire, CT, USA, 1990.
3. E. R. Tufte. *The visual display of quantitative information*. 2nd ed. Graphics Press, Cheshire, CT, USA, 2001.