

Le circuit des métadonnées à travers les plateformes de mathématiques

Thierry Bouche



*Dynamiques de l'édition scientifique,
de l'industrie de l'information,
de la documentation*

Meudon, mardi 4 novembre 2014

Sommaire

- 1 Le CEDRAM
- 2 Le circuit idéal d'un article
- 3 L'environnement de production du CEDRAM
- 4 EuDML

Le CEDRAM

Introduction

- Un projet CNRS (INSMI) et UJF de soutien aux revues académiques de recherche en mathématiques (principalement)
- Réalisation : cellule Mathdoc (Grenoble)
- « Académiques » : indépendants (labos. . .) ou sociétés savantes
- Un ensemble d'outils de production
- Un portail de diffusion www.cedram.org
- En projet : élargir le périmètre et le nombre de services

Le CEDRAM

Principes

- Les articles sont des travaux de recherche originaux validés scientifiquement sous la responsabilité d'un comité éditorial
- Les actes de séminaires ou conférences suivent une politique plus légère (un responsable scientifique)
- Une fois publiés, les textes ne sont jamais modifiés
- La navigation des collections est libre, les métadonnées sont ouvertes largement, elles peuvent être révisées périodiquement, au fil des enrichissements (liens...)
- L'accès aux articles est libre à l'issue d'une barrière mobile (BM) raisonnable (À ce jour : 4 revues et 5 séries d'actes en accès libre, 1 revue BM = 2 ans, 2 revues BM = 5 ans)
- Les articles sont versés dans une bibliothèque numérique de référence (libre accès à terme, voie « orange »)

Le CEDRAM

Interopérabilité

- Les sources des articles sont archivées (pour leur préservation)
- Les articles sont versés dans NUMDAM 1-2 ans après parution (pour l'accès sur le long terme)
- NUMDAM est partenaire de EuDML et bientôt GDML (pour une plus grande visibilité et la construction d'une bibliothèque de référence couvrant toute la discipline)
- Plus généralement : OAI-PMH (oai_dc, mini-dml, NLM/EuDML)

Le circuit idéal d'un article

Eurêka \LaTeX auteur

Prépublication arXiv/HAL

Soumission Choix de la revue, rapport de referee

Production Corrections, mise aux normes

Publication Texte et métadonnées définitifs

Indexation MathScinet, zbMath, Crossref

Archivage indépendant, pérenne & distribué
Projet : Private CLOCKSS européen

DML Bibliothèque de référence en libre accès à terme

Partagée, mondiale, pérenne, sans but lucratif

Enrichissements globaux : consolidation des métadonnées, liens, similarité, *math mining*, etc.

L'environnement de production du CEDRAM

Objectifs

- Un chaîne de production *maîtrisable* par de petites structures
 - Uniquement des formats habituels en entrée : \LaTeX , Bib \TeX . . .
- Produisant une édition de *qualité*
 - Une structure assez universelle pour permettre des maquettes Web et papier très variables mais des fonctions de base constantes
 - Éditions papier & électronique identiques
 - Métadonnées exactes (déduites des articles)
 - (Méta)données précises et universelles (Unicode, PDF, XML, MathML. . .)
 - Métadonnées compatibles avec NUMDAM, OAI-PMH, faciles à transformer et exporter (dégrader. . .)

L'environnement de production du CEDRAM

Principes

- 1 Les articles et leurs métadonnées sont produits à partir des mêmes sources \LaTeX
- 2 Une métadonnée est déclarée *une fois* au plus, dans un fichier source *ad hoc*
- 3 Toute métadonnée qui n'est pas pertinente dans un fichier source ne doit pas s'y trouver (exemple : tomaisou ou pagination dans un article)
- 4 Tout ce qui peut être calculé *doit* l'être
Pas de duplication de métadonnées pas d'heuristiques ou de copier-coller (*un processus*) \implies pas d'erreurs de dernière minute !
- 5 L'environnement doit être d'appropriation facile

L'environnement de production du CEDRAM

Les niveaux

Revue

- Les constantes de la revue : titre, ISSN, maquette. . .
- Les quasi-constantes : pages de titre, ours, comité de rédaction, instructions aux auteurs. . .

Volume

- Données bibliographiques (année, mois, tome, fascicule, lieu. . .)
- Le folio du premier article
- La liste ordonnée des articles
- Variables globales et extras (éditorial, pubs. . .)

Article

- Langue, titre, DOI, auteurs, résumés, texte, biblio, etc.

Pages

- Tous les numéros de pages sont calculés à la volée

Diversité graphique

Même source, présentation variable (1).



Annales de l'institut Fourier



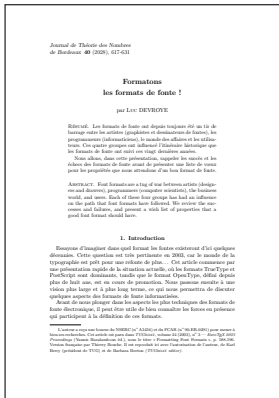
Annales de la faculté des sciences de Toulouse (mathématiques)

Diversité graphique

Même source, présentation variable (2).



*Annales mathématiques
Blaise-Pascal*



*Journal de théorie des nombres
de Bordeaux*

Diversité graphique

Même source, présentation variable (3).

Mathematics in Action
Vol. 4, 1-11 (2012)

**Formations
les formats de fonte !**

LUC DEWEGE*

*MATH Université, Université ILLIAC (France).
E-mail address: luc@illiac.org

Résumé

Les formats de fonte ont depuis toujours été un jeu de langage entre les auteurs (graphistes et documentalistes de livres), les programmeurs (algorithmiciens), le monde des éditeurs et les utilisateurs. Ces quatre groupes sont désormais (partiellement) identifiés par les formats de fonte qui ont ainsi un usage beaucoup moins "bien défini" dans cette profession. Après les avoir vu à travers les formats de fonte avant de présenter une liste de notes pour les graphistes qui sont intéressés par ces formats de fonte.

Abstract

Font formats are a play of language between authors (designers and documentalists of books), the programmers (algorithmists), the business world, and users. Each of these four groups has had an influence on the path that font formats have followed. We review the common and different, and present a wish list of properties that a great font format should have.

1. Introduction

Essayer d'imaginer dans quel format les fontes existaient d'il y a quelques décennies. Cette question est très pertinente en 2012, car le monde de la typographie est peut-être plus en volée de gloire... Cet article commence par une présentation rapide de la situation actuelle, où les formats TrueType et PostScript sont dominants, tandis que le format OpenType, utilisé depuis plus de huit ans, est en cours de promotion. Nous pouvons ensuite à une vision plus large et à plus long terme, ce qui nous permettrait de discuter quelques aspects des formats de fonte informatiques. Avant de nous éloigner dans les aspects les plus techniques des formats de fonte informatiques, il peut être utile de faire connaître les formes ou parties qui participent à la création de ces formats.

D'abord et avant tout, les utilisateurs administrateurs disposent de formats simples et utiles, faciles à manipuler et modifier. Ils veulent à la fois sentir les notes sur les notes d'un auteur par de grands typographes et les raffinement techniques fournis par les experts de fonte informatique. En outre, les utilisateurs professionnels veulent un certain degré de flexibilité des fontes, de façon à pouvoir à la fois les modifier.

Les auteurs et typographes veulent une solution compatible avec les formats de fonte à leur disposition. Les premiers typographes étaient à peu près tous des artisans. Au vingtième siècle, diverses machines technologiques furent utilisées pour des conceptions comme Linotype et Compugraphic. Une autre est venue dans l'histoire, celle de 1982, et à la fin du 20e siècle, les programmes (Quark, Metafont, etc.), sont venus à l'honneur (PostScript et TrueType). Une autre technologie est venue à l'honneur et a été l'élément de l'avenir de Real Time Graphics de TPC et de Helvetica (TrueType) elle-même.

1

Mathematics in Action

Ministère de l'éducation et de la formation
Québec
Volume 36 (2012) (10-11)

**FORMATONS
LES FORMATS DE FONTE !**

Luc Dewege

Résumé

Les formats de fonte ont depuis toujours été un jeu de langage entre les auteurs (graphistes et documentalistes de livres), les programmeurs (algorithmiciens), le monde des éditeurs et les utilisateurs. Ces quatre groupes sont désormais (partiellement) identifiés par les formats de fonte qui ont ainsi un usage beaucoup moins "bien défini" dans cette profession. Après les avoir vu à travers les formats de fonte avant de présenter une liste de notes pour les graphistes qui sont intéressés par ces formats de fonte.

Abstract

Font formats are a play of language between authors (designers and documentalists of books), the programmers (algorithmists), the business world, and users. Each of these four groups has had an influence on the path that font formats have followed. We review the common and different, and present a wish list of properties that a great font format should have.

1. Introduction

Essayer d'imaginer dans quel format les fontes existaient d'il y a quelques décennies. Cette question est très pertinente en 2012, car le monde de la typographie est peut-être plus en volée de gloire... Cet article commence par une présentation rapide de la situation actuelle, où les formats TrueType et PostScript sont dominants, tandis que le format OpenType, utilisé depuis plus de huit ans, est en cours de promotion. Nous pouvons ensuite à une vision plus large et à plus long terme, ce qui nous permettrait de discuter quelques aspects des formats de fonte informatiques.

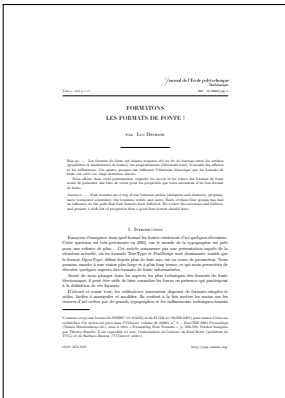
D'abord et avant tout, les utilisateurs administrateurs disposent de formats simples et utiles, faciles à manipuler et modifier. Ils veulent à la fois sentir les notes sur les notes d'un auteur par de grands typographes et les raffinement techniques fournis par les experts de fonte informatique. En outre, les utilisateurs professionnels veulent un certain degré de flexibilité des fontes, de façon à pouvoir à la fois les modifier.

Les auteurs et typographes veulent une solution compatible avec les formats de fonte à leur disposition. Les premiers typographes étaient à peu près tous des artisans. Au vingtième siècle, diverses machines technologiques furent utilisées pour des conceptions comme Linotype et Compugraphic. Une autre est venue dans l'histoire, celle de 1982, et à la fin du 20e siècle, les programmes (Quark, Metafont, etc.), sont venus à l'honneur (PostScript et TrueType). Une autre technologie est venue à l'honneur et a été l'élément de l'avenir de Real Time Graphics de TPC et de Helvetica (TrueType) elle-même.

*Séminaire de Théorie spectrale
et géométrie*

Diversité graphique

Même source, présentation variable (5).



Journal de l'École polytechnique

Cahiers GUTenberg

Diversité graphique

Mêmes formats, présentation Web variable.

Journal de l'École polytechnique
Mathématiques

Table des matières de ce fascicule | Table of contents

Revue d'abstracts

Table des matières de ce fascicule | Table of contents

Revue d'abstracts

Table des matières de ce fascicule | Table of contents

Revue d'abstracts

Table des matières de ce fascicule | Table of contents

Journal de l'École polytechnique

Confluentes Mathematiqi

Table des matières de ce fascicule | Table of contents

Revue d'abstracts

Table des matières de ce fascicule | Table of contents

Revue d'abstracts

Table des matières de ce fascicule | Table of contents

Confluentes Mathematiqi

ANNALES de l'INSTITUT FOURIER

Table des matières de ce fascicule | Table of contents

Revue d'abstracts

Table des matières de ce fascicule | Table of contents

Revue d'abstracts

Table des matières de ce fascicule | Table of contents

Annales de l'Institut Fourier

ANNALES FACULTÉ DES SCIENCES TOULOUSE

Table des matières de ce fascicule | Table of contents

Revue d'abstracts

Table des matières de ce fascicule | Table of contents

Revue d'abstracts

Table des matières de ce fascicule | Table of contents

Annales Fac. Sci. Toulouse

Les enrichissements

- Les métadonnées comportent les biblios
- $\text{T}_{\text{E}}\text{X} \mapsto \text{XML}/\text{MathML}$ (Tralics)
- Outils de *matching* (articles et biblios) : MathSciNet, Crossref (services); Jahrbuch, zbMATH, mini-DML, EuDML (outils maison)
- Export NLM + plein texte vers EuDML

La vision EuDML



La bibliothèque numérique de mathématiques devrait s'efforcer de réunir un corpus mathématique **aussi vaste que possible** pour

- le rendre **disponible en ligne**
 - en accès **libre à terme**,
 - sous la forme d'une collection **de référence**,
 - **alimentée** en continu par les nouveautés des éditeurs,
 - **valorisée** par des outils de recherche et référencement sophistiqués,
 - développée et entretenue par un réseau d'**institutions**
- + et le **préserver** à très long terme

EuDML

Enrichissements après import

- Intégration de métadonnées (notamment en anglais) en provenance du zbMath
- Conversion des formules en MathML (OCR Infty, Maxtract → recherche de formules)
- Calcul de similitude des articles utilisant le texte intégral, les formules, les métadonnées
- Expérimental : textes accessibles

GDML

Projet en cours de définition

- Corpus mondial interoperable
 - Infrastructure basse rendant possible la construction de services « sémantiques » et « computationnels »
 - Format de description des structures mathématiques (objets, énoncés, preuves, définitions. . .)
 - Outils de traitement automatisé de cette langue surnaturelle
- ⇒ Activités à Washington et Séoul organisées par le CEIC
- ⇒ Rapport du Conseil national de la recherche (NRC, Nat. Sc. Acad. USA)
Developing a 21st Century Global Library for Mathematics Research, The National Academies Press, 2014
- ⇒ Groupe de travail de l'union mathématique internationale

Merci !

Thierry Bouche

Directeur cellule Mathdoc (UMS 5638 UJF/CNRS)

Président EuDML initiative

Membre CEIC (IMU)

Membre GDML-WG (IMU)