



PROJET TUTORÉ
M1 TECHNOLOGIES DE L'INTERNET

**Conception et développement d'une
application d'annotation thématique dans
l'environnement Gate**

Auteurs :

Roland BARY
Charles FOLLET

Tuteurs :

Marie-Noëlle BESSAGNET
Annig LACAYRELLE
Albert ROYER
Christian SALLABERRY

Remerciements

Nous tenons à remercier nos tuteurs pour leur pédagogie et leur encadrement. Monsieur Royer pour sa précision et sa connaissance pointue du domaine. Madame Lacayrelle pour son soutien et sa clarté. Monsieur Sallaberry pour nous avoir remis sur de bonnes pistes quand nous nous égarions. Et enfin, madame Bessagnet pour avoir assuré la coordination et le suivi de ce projet.

Table des matières

I	Introduction	3
II	Cahier des charges	5
I	Contexte	6
II	Description	7
III	Diagramme de Gantt prévisionnel	7
III	Cadre d'analyse	8
IV	Définition de concepts	9
1	Gazetier	9
2	Entité nommée	9
3	Expression régulière	9
V	Outil	9
1	L'environnement Gate	9
IV	Développement	10
VI	Prise en main de l'environnement	11
VII	Définition des dimensions d'annotation et leur contenu	11
VIII	Première recherche d'entités nommées avec les gazetiers . . .	11
IX	Deuxième recherche d'entités nommées avec les règles JAPE .	11
V	Conclusion	12

Première partie

Introduction

Avec l'évolution de manière significative des volumes d'informations sur internet, on peut observer une évolution du web vers une approche dans laquelle chaque donnée acquiert un sens afin de rendre possible une interprétation du contenu des pages web par des machines. Cette extension constitue le web sémantique. L'une des principales motivations du web sémantique est la recherche d'information sémantique.

C'est donc dans ce cadre que nous sommes intervenus pour répondre à l'appel d'offre de nos encadrants. L'objectif est l'annotation sémantique d'un document texte spécifique, qui constitue effectivement la première étape dans un processus d'indexation et de recherche d'information sémantique.

Au regard de ce qui a été exprimé en amont, se pose les problématiques suivantes :

- Existe-t-il des outils qui se prêtent aisément à l'annotation sémantique ?
- Quelle approche de conception peut nous permettre de réaliser cette étape d'annotation sur un document texte non-structuré ?

La résolution de ces différentes problématiques, nous à donc amené à organiser ce document comme suit : /*Il faut caller notre plan ici */ :) Une première partie dans laquelle nous présenterons le cahier des charges. Ensuite une seconde partie décrira quelques connaissances existantes sur le sujet avec les technologies utilisées au sein du projet.

Deuxième partie

Cahier des charges

I Contexte

A partir des travaux de Georges DEPEYROT sur les monnaies carolingiennes, nous avons travaillé pour une équipe parisienne de numismates sur l'annotation du Numéraire Carolingien ¹.

Sur celui-ci, l'équipe a besoin d'effectuer des recherches :

Temporelles : Quelles étaient les pièces en circulation de l'an 859 à l'an 865 ?

Spatiales : Dans quels ateliers, les pièces de type Obole de Charlemagne ont été produites ?

Thématiques : Combien d'exemplaires de la monnaie d'or de Charles le Chauve ont été étudiés ?

Répondre à cette demande implique de définir puis d'explorer les dimensions temporelles, spatiales et thématiques de l'ouvrage.

Pour cela, il est nécessaire de connaître le domaine et l'ouvrage afin de savoir quelle information correspond à quelle dimension.

Une fois cet apprentissage fait, nous pouvons construire des règles dans une chaîne de traitement permettant d'annoter chaque information en fonction de sa dimension.

Les monnaies carolingiennes sont le domaine central pour la réalisation du projet. Les ressources nécessaires à l'annotation (ici sous forme de gazetiers) ont été construites à partir des données de l'ouvrage.

Fort de son expérience dans le domaine, la maîtrise d'ouvrage nous a demandé d'utiliser la boîte à outils logicielle GATE qui sera utile pour le traitement du langage naturel.

En résumé, les caractéristiques du projet sont :

- l'apprentissage et la compréhension du domaine considéré,
- l'étude des principes d'annotation de documents,
- le développement d'une chaîne d'annotation dans GATE,
- la mise en place d'une visualisation des résultats.

II Description

III Diagramme de Gantt prévisionnel

1. <http://www.cgb.fr/le-numeraire-carolingien-moneta-77-3e-edition-depeyrot-georges, Ln71,a.html>

Troisième partie

Cadre d'analyse

Introduction

IV Définition de concepts

- 1 Gazetier
- 2 Entité nommée
- 3 Expression régulière

V Outil

- 1 L'environnement Gate

Quatrième partie

Développement

- VI Prise en main de l'environnement
- VII Définition des dimensions d'annotation et leur contenu
- VIII Première recherche d'entités nommées avec les gazetiers
- IX Deuxième recherche d'entités nommées avec les règles JAPE

Cinquième partie

Conclusion