

# Recherche d'Information (RI)

**Mathias Géry**

Mathias.Gery@univ-st-etienne.fr

Laboratoire Hubert Curien, UMR CNRS 5516

Université Jean Monnet Saint-Étienne

# Projet RI

- Sujet : Moteur de recherche de documents structurés.
  - Développer un Système de RIS.
  - Tâche INEX : « Relevant In Context » (RIC).
  - Collection : sous-ensemble de la collection INEX 2009.
  - Ressources téléchargeables sur cours en ligne.
  - Étapes intermédiaires, à chaque fois :
    - Envoi de « runs » RIC INEX,
    - Vérification → évaluation → résultats.
    - **Principe 1** : j'évalue seulement les runs reçus **sur Claroline** en temps et en heure.
    - **Principe 2** : j'évalue seulement les runs valides.
  - Rendu final : runs + rapport.

- Questions (précises SVP) : [Mathias.Gery@univ-st-etienne.fr](mailto:Mathias.Gery@univ-st-etienne.fr)

Mathias.Gery@pers.univ-st-etienne.fr

# Projet RI : calendrier

- 14/11 : équipes de 4 étudiants (par mail) :
  - Nom d'équipe : concaténation des 4 prénoms, sans espace ni accent.
  - Ex : JoeJackWilliamAverell
- 14/11 : envoi des premiers runs INEX (étapes 1 et 2).
  - ➔ 18/11 : vérification et évaluation ➔ 1<sup>ers</sup> résultats (en cours).
- J-2 de chaque cours : envoi des runs suivants (étapes 3-10).
  - ➔ Chaque cours : résultats des runs (06/12, 13/12 et 5/1/2017).
- 11/1 : envoi des runs définitifs.
  - ➔ 17/1 : évaluation finale des runs.
- 25/1 : deadline rendu rapport.

# Rappel : exercice avec Michel Beigbeder

- Write a program in the language of your choice (perl, sh, lex, C, java, C++, python, etc.) which index documents contained in a file.
- (...)
- Several collection files of increasing size are available on:

<https://ead-sciences.univ-st-etienne.fr/claroline/course/index.php?cid=RI2016>

- 55k 01-Text\_Only-Ascii-Coll-1-10-NoSem.gz
- 52k 02-Text\_Only-Ascii-Coll-11-20-NoSem.gz
- 103k 03-Text\_Only-Ascii-Coll-21-50-NoSem.gz
- 96k 04-Text\_Only-Ascii-Coll-51-100-NoSem.gz
- 357k 05-Text\_Only-Ascii-Coll-101-200-NoSem.gz
- 559k 06-Text\_Only-Ascii-Coll-201-500-NoSem.gz
- 747k 07-Text\_Only-Ascii-Coll-501-1000-NoSem.gz
- 1.2M 08-Text\_Only-Ascii-Coll-1001-2000-NoSem.gz
- 4.1M 09-Text\_Only-Ascii-Coll-2001-5000-NoSem.gz

# Projet RI : collection

- Collection M2DSC7Q\_2016-17 :
  - 7 requêtes d'INEX 2009.
  - 9804 articles XML (sur 2,6 millions).
  - Version texte (même format que celle utilisée avec Michel Beigbeder),
    - 76,4Mo
    - Text\_Only-Ascii-Coll-MWI-NoSem.gz
  - Version XML (pour RIS !).
    - 237Mo
    - XML-Coll-MWI-withSem.tar.gz

## 7 Requetes :

2009011 : olive oil health benefit  
2009036 : notting hill film actors  
2009067 : probabilistic models in information retrieval  
2009073 : web link network analysis  
2009074 : web ranking scoring algorithm  
2009078 : supervised machine learning algorithm  
2009085 : operating system +mutual +exclusion

# Projet RI : étape 1

- Étape 1/10 (premiers runs « texte ») :
  - Construire 5 runs (maxi).
  - Avec la version texte des documents.
  - Avec la fonction de pondération *ltn* ou *ltc*.
  - Sans structure ➔ un run = 1500 articles complets **par requête**.
  - Objectif : runs valides syntaxiquement, sans recouvrement, etc. Peu importe la performance.
  - Remarque : anti-dico, lemmatisation, etc., sont optionnels.
- Étape 2/10 (premiers runs « XML-articles ») :
  - Idem, 5 runs (maxi), mais avec la collection au format XML.
- Runs à déposer sur cours en ligne au plus tard le 14/11 (pour que je puisse les vérifier/évaluer pour le cours du 18/11).

Conseil : commencez dès aujourd'hui !

# Projet RI : étape 1

- Nom du fichier :
  - NomEquipe\_n°etape\_n°run\_pondération\_granularité\_paramètres.txt
    - N° étape sur 2 chiffres,
    - N° run unique sur 2 chiffres,
    - Pondération = un code court, ex. ltn, ltc, bm25, bm25f, etc. ou “perso”.
    - Granularité  $\in \{\text{elements, articles}\}$
    - Paramètres : description synthétique des paramètres
  - Exemple : JoeJackWilliamAverell\_01\_08\_bm25\_articles\_k1.2b0.8.txt
  - Dépôt de plusieurs runs : un seul fichier compressé sans sous-répertoire.
- Vérification et évaluation des runs valides déposés dans les temps.
- Collection au format XML → utiliser un Parser XML :
  - Java : par exemple SAXParserFactory,
  - Perl : XML::LibXML : un exemple est fourni.