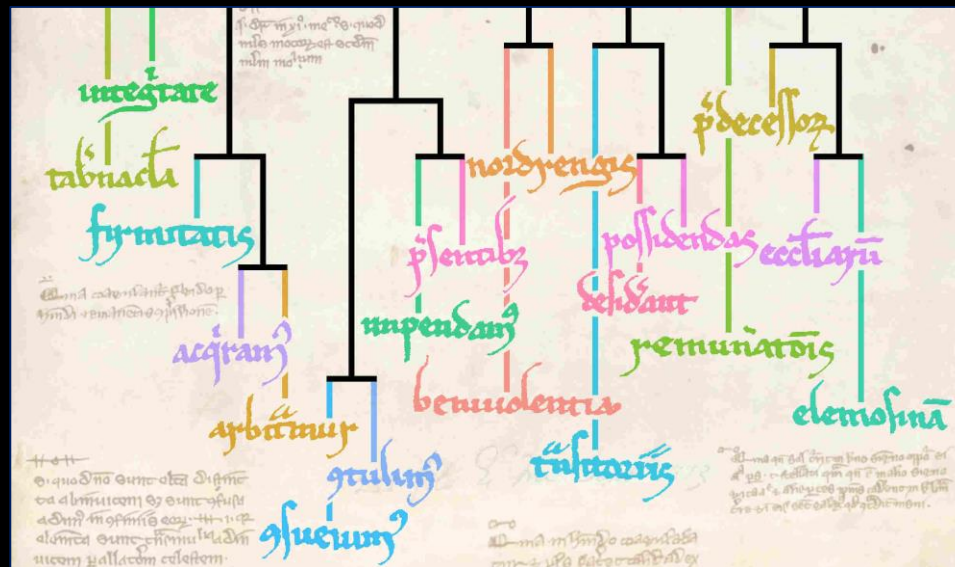


# Du parchemin à la fouille de données

Nouveaux outils pour la création, la formalisation  
et l'analyse des corpus médiévaux



IRAT

LAMoP  
UMR 8509  
CNRS  
UNIVERSITÉ PARIS 1 PANTHÉON-SORBONNE

HN

Huma-Num

UCL

Université  
catholique  
de Louvain



UNIVERSITÉ  
DE NAMUR

ICHEC

BRUSSELS MANAGEMENT SCHOOL

Paul Bertrand  
Étienne Cuvelier  
Sébastien de Valeriola  
Nicolas Perreaux  
Nicolas Ruffini-Ronzani



# 1. Fonctions R à utiliser

Tâche à effectuer	Fonction (du paquet <code>stringr</code> )
Détecter un motif	<code>str_detect()</code> <code>str_which()</code>
Extraire un motif	<code>str_extract()</code> , <code>str_extract_all()</code> <code>str_subset()</code>
Localiser un motif dans une chaîne	<code>str_locate()</code> <code>str_locate_all()</code>
Remplacer un motif par un autre	<code>str_replace()</code> <code>str_replace_all()</code>
Diviser une chaîne à l'emplacement d'un motif	<code>str_split()</code>
Supprimer les espaces blancs au début et à la fin d'une chaîne	<code>str_trim()</code>



## 2. Séquences d'échappement

<code>\n</code>	<code>\r</code>	<code>\t</code>	<code>\'</code>	<code>\"</code>
nouvelle ligne	retour de chariot	tabulation	apostrophe	guillemet



### 3. Opérateurs

.	[...]	[^...]	$[x_1 - x_2]$	
n'importe quel caractère	un des caractères écrits entre les crochets	un caractères qui n'est pas écrit entre les crochets	intervalle des caractères entre $x_1$ et $x_2$	OU logique

## 4. Quantificateurs

*	+	?	{ $n$ }	{ $n, \quad$ }	{ $n, m$ }
au moins 0 fois	au moins 1 fois	au plus 1 fois	exactement $n$ fois	au moins $n$ fois	entre $n$ et $m$ fois

## 5. Comportement des quantificateurs

comportement	explication	comment l'obtenir ?
gourmandise ( <i>greediness</i> )	répète le motif autant que possible	comportement par défaut
paresse ( <i>laziness</i> )	répète le motif aussi peu que possible	ajouter ? après le quantificateur



## 6. Classes de caractères

<code>[ :digit: ]</code> ou <code>\d</code>	<code>\D</code>	<code>[ :lower: ]</code>	<code>[ :upper: ]</code>	<code>[ :alpha: ]</code>	<code>[ :alnum: ]</code>
chiffres ( $\equiv [0-9]$ )	non-chiffres ( $\equiv [^0-9]$ )	lettres minuscules ( $\equiv [a-z]$ )	lettres majuscules ( $\equiv [A-Z]$ )	caractères alphabétiques ( $\equiv [A-z]$ )	caractères alphanumériques ( $\equiv [A-z0-9]$ )
<code>\w</code>	<code>\W</code>	<code>[ :xdigit: ]</code>	<code>[ :blank: ]</code>	<code>[ :space: ]</code>	<code>\s</code>
caractères "mots" ( $\equiv [A-z0-9_]$ )	caractères "non-mots" ( $\equiv [^A-z0-9_]$ )	chiffres hexadécimaux ( $\equiv [0-9A-Fa-f]$ )	caractères blancs (espaces et tabulations)	caractères d'espacement (tabulation, nouvelle ligne, etc.)	espace
<code>\S</code>	<code>[ :punct: ]</code>	<code>[ :graph: ]</code>	<code>[ :print: ]</code>	<code>[ :cntrl: ]</code>	
non-espace	ponctuation ( <code>! " # \$ % &amp; ' ( )</code> <code>* + , - . / : ; &lt; =</code> <code>&gt; ? @ [ ] ^ _ ` {</code> <code>  } ~</code> )	caractères graphiques ( $\equiv [[:alnum:]]$ <code>[ :punct: ]</code> )	caractères visibles ( $\equiv [[:alnum:]]$ <code>[ :punct: ] \s</code> )	caractères de contrôle ( <code>\n</code> , <code>\r</code> , etc.)	



## 7. Position du motif dans la chaîne

^	\$	\b	\B
début de la chaîne	fin de la chaîne	caractère vide désignant le début d'un mot	caractère vide désignant la fin d'un mot





## 8. Retro-références (*backreferences*)

( . . . )	( ? : . . . )
groupe capturant	groupe non-capturant