

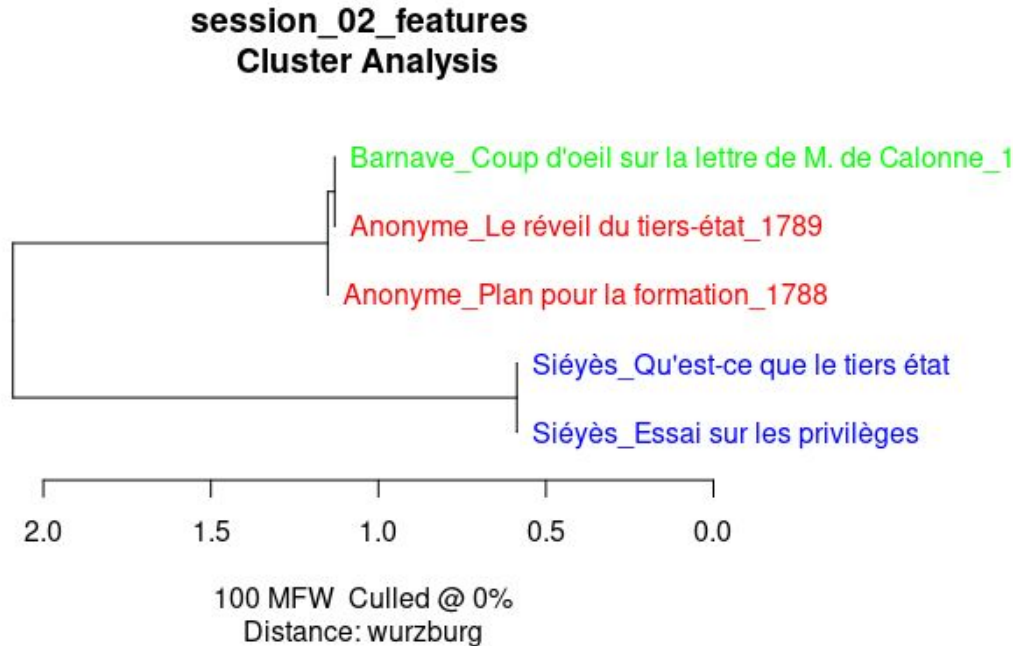
Stylometry with R

- Part 3. Words behind trees

Joanna Byszuk and
Artjoms Šeļa

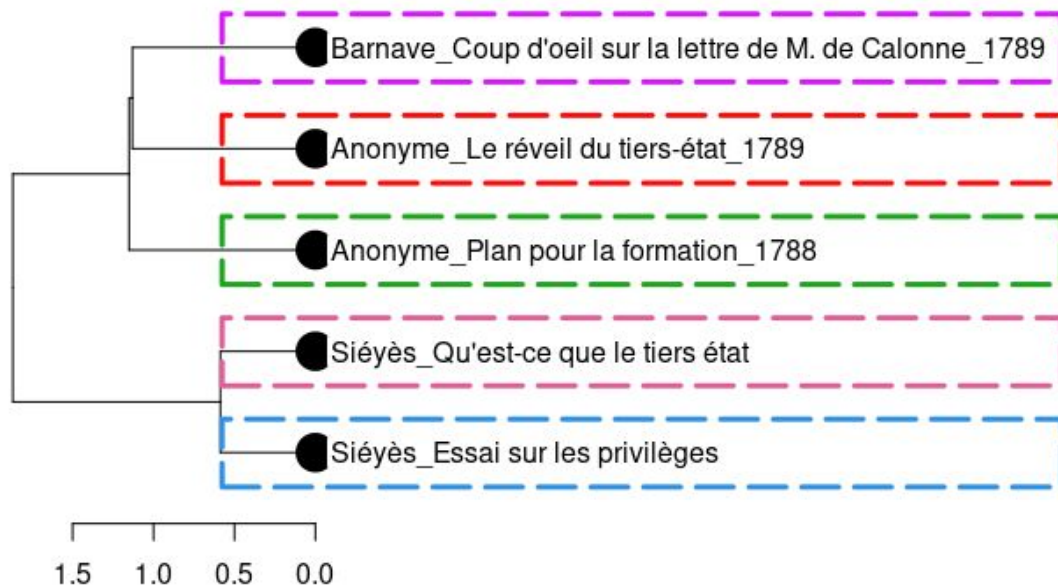


Cluster analysis: grouping suggestion



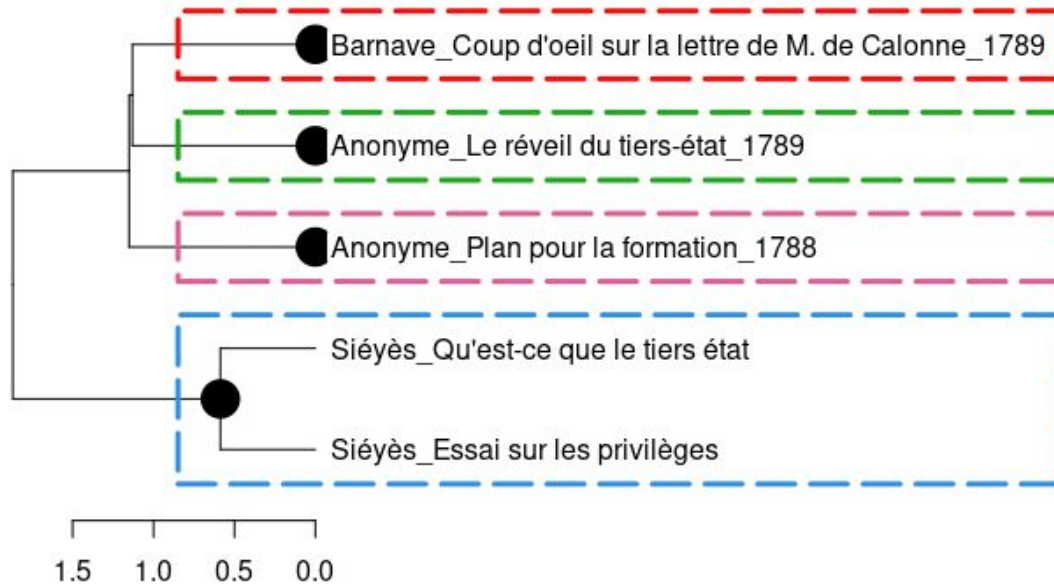
Cutting trees by groups!

Hierarchical clustering, cut at k= 5



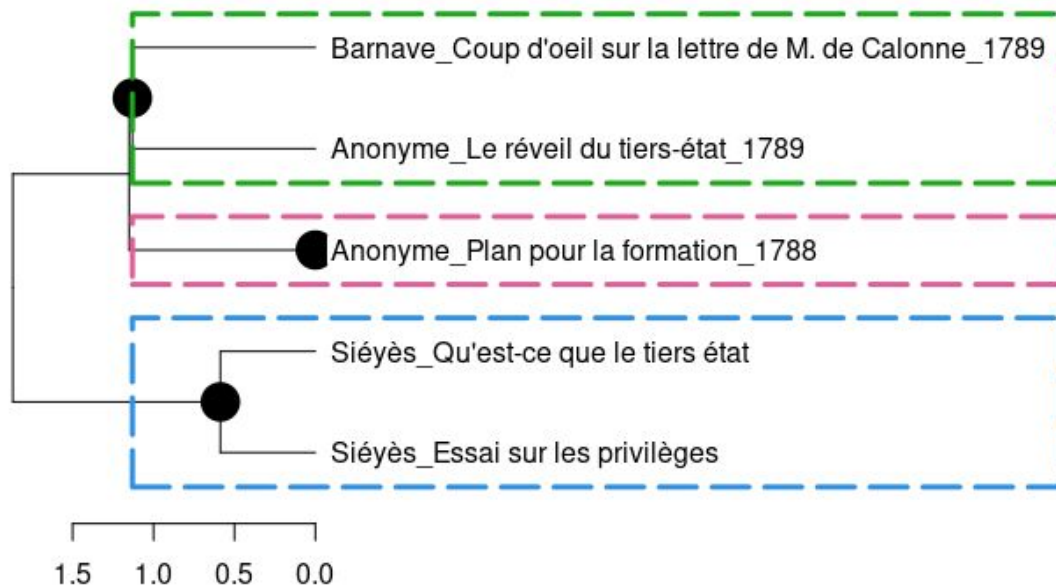
Cutting trees by groups!

Hierarchical clustering, cut at k= 4



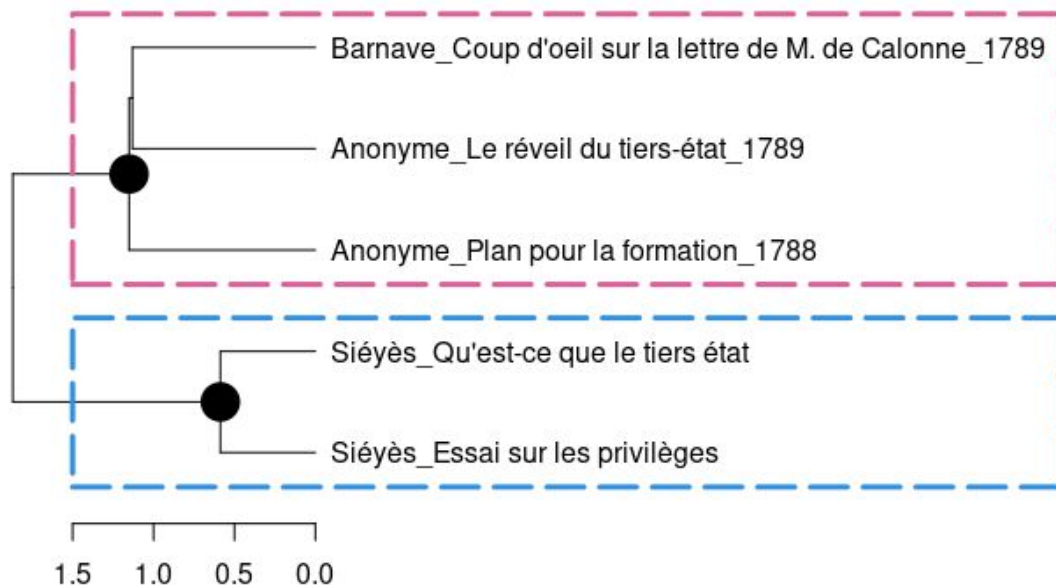
Cutting trees by groups!

Hierarchical clustering, cut at k= 3



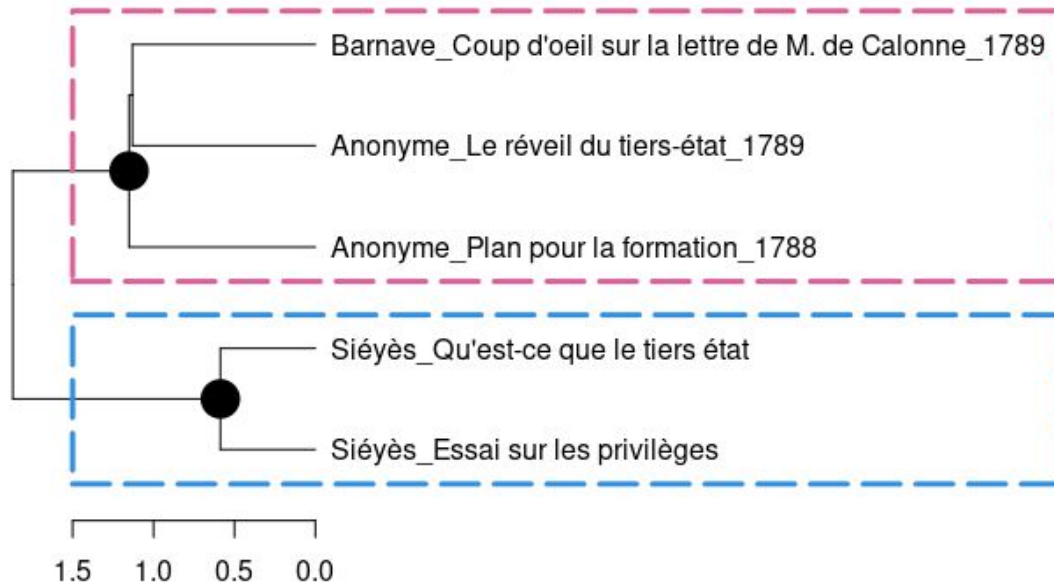
Cutting trees by groups!

Hierarchical clustering, cut at k= 2



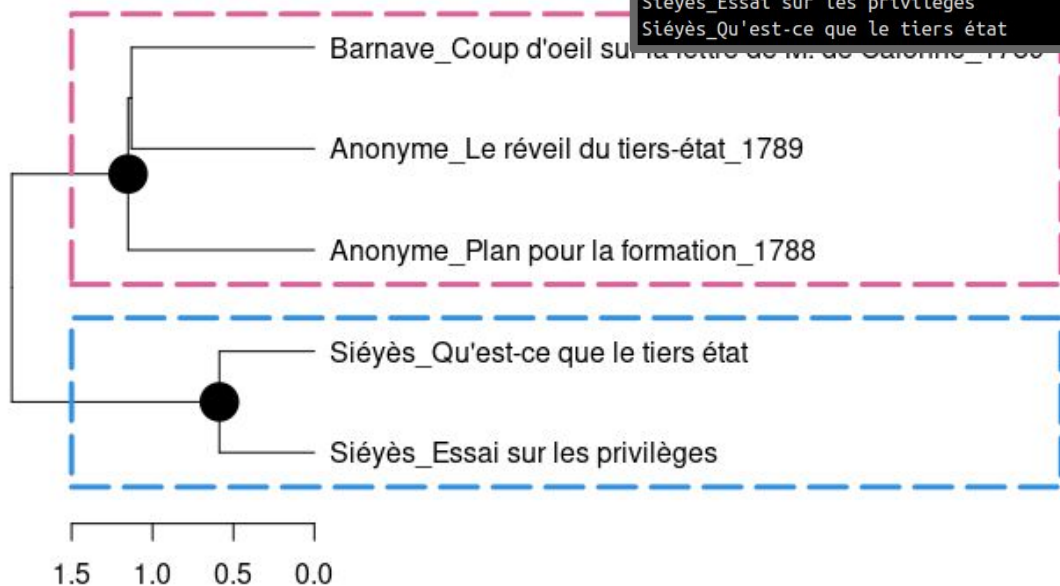
Think of new 'pink' and 'blue' classes

Hierarchical clustering, cut at k= 2



Classes driven by word frequencies!

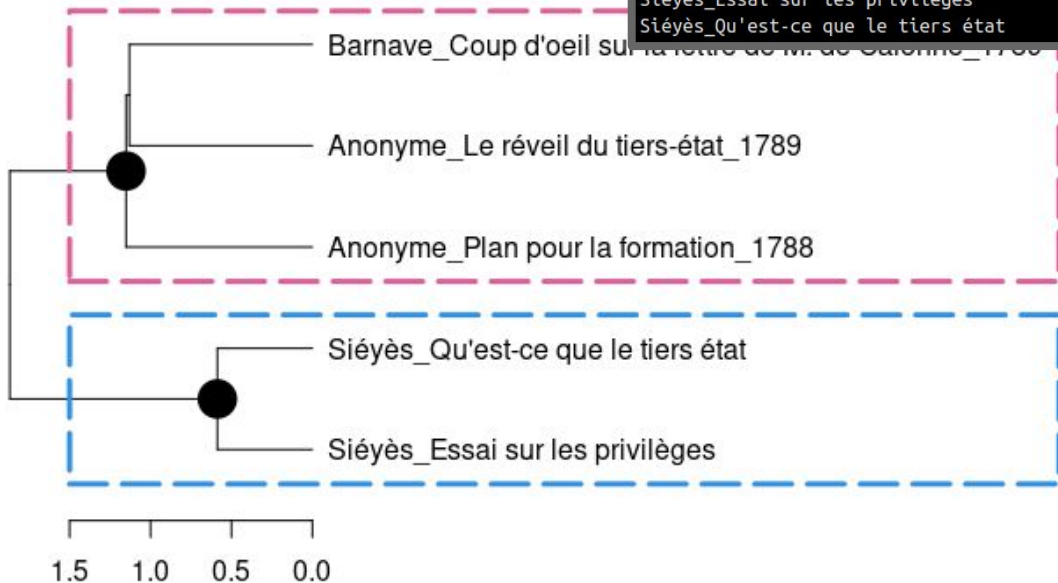
Hierarchical clustering, cut at k= 2



	de	la	les	l	à	le
Anonyme_Le réveil du tiers-état_1789	4.9919255	3.1490453	2.9068111	2.0993635	1.8333808	1.9663722
Anonyme_Plan pour la formation_1788	4.3847242	3.1117397	4.6676096	1.6030174	1.3672796	1.7444602
Barnave_Coup d'oeil sur la lettre de M. de Calonne_1789	4.2275472	3.0196766	2.8443405	2.2209234	1.1494253	2.1235145
Siéyès_Essai sur les privilèges	4.3275072	2.729227	2.4114403	2.0095336	2.3460136	1.8412936
Siéyès_Qu'est-ce que le tiers état	3.9608393	2.9340656	2.3968003	1.874459	2.0416082	1.8983375

Feature ~ cluster association

Hierarchical clustering, cut at k= 2



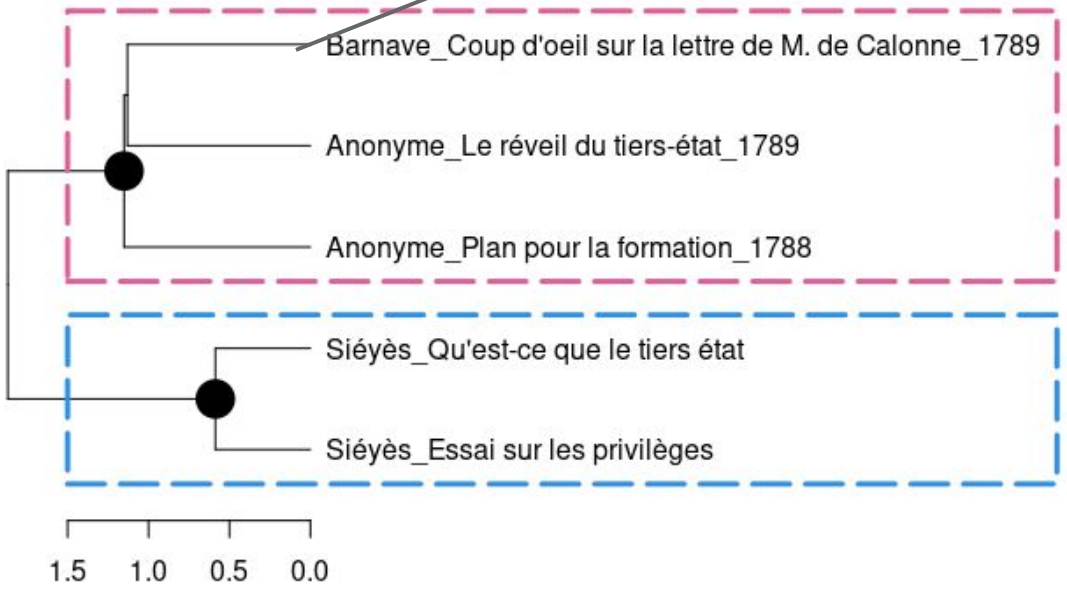
	de	la	les	l	à	le
Anonyme_Le réveil du tiers-état_1789	4.9919255	3.1490453	2.9068111	2.0993635	1.8333808	1.9663722
Anonyme_Plan pour la formation_1788	4.3847242	3.1117397	4.6676096	1.6030174	1.3672796	1.7444602
Barnave_Coup d'oeil sur la lettre de M. de Calonne_1789	4.2275472	3.0196766	2.8443405	2.2209234	1.1494253	2.1235145
Siyès_Essai sur les privilèges	4.3275072	2.729227	2.4114403	2.0095336	2.3460136	1.8412936
Siyès_Qu'est-ce que le tiers état	3.9608393	2.9340656	2.3968003	1.874459	2.0416082	1.8983375

You can use 'emerging' class information from the tree to define corpora and check **which features differ** across clusters.

Think about it as **keyword** problem.

	de	la	les	l	à	le
Anonyme_Le réveil du tiers-état_1789	4.9919255	3.1490453	2.9068111	2.0993635	1.8333808	1.9663722
Anonyme_Plan pour la formation_1788	4.3847242	3.1117397	4.6676096	1.6030174	1.3672796	1.7444602
Barnave_Coup d'oeil sur la lettre de M. de Calonne_1789	4.2275472	3.0196766	2.8443405	2.2209234	1.1494253	2.1235145
Siéyès_Essai sur les privilèges	4.3275072	2.729227	2.4114403	2.0095336	2.3460136	1.8412936
Siéyès_Qu'est-ce que le tiers état	3.9608393	2.9340656	2.3968003	1.874459	2.0416082	1.8983375

Hierarchical clustering, cut at k= 2



Cluster 2

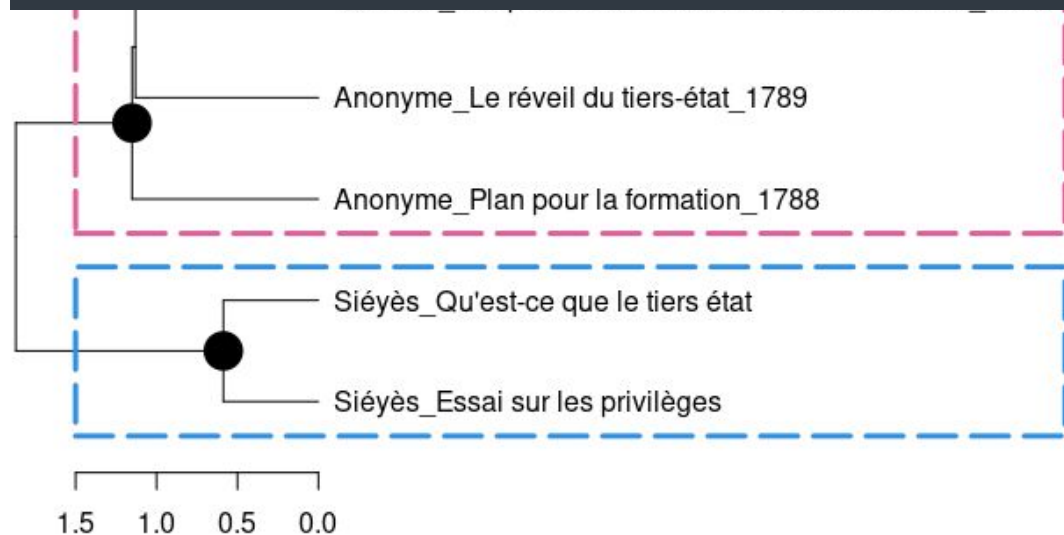
est
sont
sans
sa
sur
et
ses

Cluster 1

fe



pute, lequel fe l'Assemblée de Témion ou arrondissement. Ces Députés ne pourront être élus que parmi les propriétaires domiciliés ou parmi les forains qui auront des propriétés dans le lieu payant cinquante livres de charges réelles & pour être élu, il ne fera pas nécessaire d'être présent à l'Assemblée. Les Etats indiqueront les chefs-lieux des arrondissements ailleurs que dans les villes qui ont des Députés particuliers & pour la première convocation, les Députés de l'Assemblée de Grenoble fe réuniront à Villeceux de l'Election de Vienne, à Bouroin; ceux de l'Election de Romans à Beaupaire; ceux de l'Election de Valence, Chabeuil ceux de l'Election de Gap, aux Charges ceux de l'Election de Montebellard les Députés qui devront représenter le tiers état; aux Etats les Procès-verbaux fera envoyer au Secrétaire le nom des



Cluster 2

est
sont
sans
sa
sur
et
ses

Cluster 1

fe

An early alpha function from me

```
library(stylo)
source("src/view_tree.R")
|
res1 <- stylo(corpus.dir="../session_01_sample",
              gui=F,
              distance.measure="wurzburg")

view_tree(res1, ## saved results variable from
            k=2, ## to how many groups you want
            p=0.05, ## p-value threshold for view
            color_leaves=F ## should it color the
            )
```

```
1 CLUSTER 1
2 =====
3 TEXTS
4 Anonyme_Plan pour la formation 1788
5 Anonyme_Le réveil du tiers-état 1789
6 Barnave_Coup d'oeil sur la lettre de M. de Calonne 1789
7 =====
8 FEATURES associated (p<0.05)
9
10 fe
11
12
13
14
15 CLUSTER 2
16 =====
17 TEXTS
18 Siéyès_Essai sur les privilèges
19 Siéyès_Qu'est-ce que le tiers état
20 =====
21 FEATURES associated (p<0.05)
22
23 est sont sans sa sur et ses
24
25
```


An early alpha function from me

Some tips:

- Bad idea to use when a group includes one 'leaf'
- Bad idea to use when **$k > 5$**
- For now, does not work with consensus trees
- Untested, might throw horrible errors 🧛‍♂️ 🧛‍♀️ 🧛‍♂️

Assembling a corpus based on **$k=2$** cut and doing