

TP1

Map Reduce simplifié

Objectif : Compter des éléments dans un dictionnaire à partir d'un texte à analyser.

1. Composer le dictionnaire de manière à ce que chaque mot différent du texte constitue une entrée.
2. Chaque entrée permet de dénombrer le nombre d'occurrences de chaque mot.
3. Trier le dictionnaire avec les entrées les plus nombreuses en haut.
4. Proposer une visualisation des données avec un histogramme en barres horizontales.

Texte de base :

Le petit poucet c'est perdu dans la campagne ! Il a peur ... A vous de l'aider. Pas facile ? Le petit poucet compte sur vous !

Consignes :

Chaque occurrence sera représentée par une étoile :

Non trié :

UNSORT	
<campagne>	*
<l'aider.>	*
<!>	**
<?>	*
<>	*
<de>	*
<facile>	*
<Il>	*
<petit>	**
<vous>	**
<a>	*
<peur>	*
<poucet>	**
<c'est>	*
<perdu>	*
<Le>	**
<dans>	*
<la>	*
<A>	*
<Pas>	*
<compte>	*
<sur>	*

Après tri sur les clés :

```
    SORT by Keyes
<!>          **
<?>          *
<A>           *
<Il>          *
<Le>          **
<Pas>         *
<a>           *
<campagne>    *
<compte>      *
<c'est>       *
<dans>        *
<de>          *
<facile>      *
<la>          *
<l'aider.>    *
<perdu>       *
<petit>       **
<peur>        *
<poucet>      **
<sur>         *
<vous>        **
<...>        *
```

Après tri sur les valeurs :

```

SORT by Values
<petit>          **
<Le>             **
<!>             **
<poucet>         **
<vous>           **
<campagne>       *
<Il>             *
<facile>         *
<de>            *
<a>             *
<peur>          *
<?>            *
<c'est>         *
<perdu>         *
<l'aider.>      *
<dans>          *
<la>            *
<A>             *
<Pas>           *
<compte>        *
<~>            *
<sur>           *
```

Consignes autour du projet : Le projet sera réalisé dans une console pour faire simple. La version de base doit pouvoir supporter au moins le jeu d'essai sur une œuvre d'Harry Potter.

Vous pouvez vérifier votre projet avec l'intégrale de « Harry Potter à l'école des sorciers ! »

Evidemment, pour vérifier la régularité du traitement, il va falloir lire l'œuvre intégrale :)

Le mot « de » apparaît 3226 fois `<1><de><3226>` et le mot « Harry » 652 fois `<17><Harry><652>` ! Le traitement génère 15021 entrées dans le dictionnaire.
`<15021><courageux...><1> *`

Attention, le mot Harry apparaît également dans d'autres formats, « Harry. » « Harry, » « ? Harry » ... Il s'agit donc de valeurs brutes à retravailler pour obtenir une vraie tendance, sur les répétitions. Mais cela n'était pas l'objectif de cette version.

[illegible]