

## What has been done?

- Setup of new CPU machine → ready and running
- Started Research on ANN and DNN methods to use Word Embeddings
- started looking at alternative corpora datasets
- Currently building Word2Vec without lemmatization
- Playing with similarities
- Building Dictionary is longer on CM? exactly 1h?

## short term

- playing with - ANN chatbot
  - parallel algorithm
  - protocol to evaluate proactive
- Literature and test of existing chatbot solution
- Combine with ANN
- Document in report Gensim life
- Compare with : fastText, word2Vec, Glove
- Build alternative Word2Vec with different corpora
- Play with Word2Vec
  - geometries
  - analogies (capital of science)
- Compare unlemmatized vs lemmatized
- Compare Premade W2VC(google) vs mine

## overall progress

- fixed the memory problem with gensim by splitting data sets by pages and use a generator by line
- Tried to use cloud based machines. Expensive and not worth it for CPU computation, same perf as on hera-gr.
- Full wikipedia EN trained model : 35 et 20h on hera-gr.  
→ with lemmatization
- Word2Vec operations working
- Testing on a CPU dedicated machine to train

## Questions:

- meeting next week? A partir mercredi apres-midi

Faire :

- - stopword
- Generate similar sentences
- Use Word2Vec as it is
- "drunk guy in a pub"

Avancer sur la doc

Écrire tout ce qui a été fait jusqu'à présent  
→ DUP

- 1) <sup>Versions</sup> Dans quel mesure peut-on exploiter le W2V, sans intelligence
- 2) Seq2Seq, article + Implementation, avec quelque chose tout fait  
→ Comparer les modules

Explorer la génération de phrases  
→ puis greffer

2 semaines de rédaction

→ mettre à jour le plan

—