# what has been done?

- full wikipedia english training : 3 j 20 h on heia-fr
- lemmatization vs non-lemmatization
  - without → performance → less shades
- playing with operations by vectors
  - words
  - still working on sentences
- tested analogies operations: Athens is to Greece
  as Baghdad it to Iraq
- translation random
- Plus proche voisin → le chat est un mamifaire
  - vecteur appliquer au reste de la phrase
  - vecteur de translation

# short term

- starting new sprint on - 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 : Fast Text, word2Vec, Glove

- translation dirigé → shift

# overall progress

- fixed the memory problem with gensim by
  splitting data sets by pages and use a generator by line
- Tryied to use cloud based machines. Expensive and
  not worth it for CPU computation, same perf as on heia-fr.
- Full wikipedia EN trained model : 3 j et 20 h on heia-fr
  → with lemmitization
- Word2Vec operations working

Prendre note des examples
->

est-um
sur base d'exemple

Exemple 10ème 20ème de Capital

    -> Average of dimension
    -> Vecteur de différence
    -> Capital [ ]

        -> Capital of science

    -> Est le marie de..

  -> Est-ce que les relations W2V
      sont les mêmes que les relations des
                          entités

  -> la force des similarités
      -> Cosine (direction)
      -> Normalizé
      -> Non normalisé

-> Articles avec des trucs Bizzare qu'on
      a vu
-> Why people are not using full lemized vocab
-> Thinking a metric to compare via similar
      -> mine VS google

methodd (wv_a, wv_b, list_ de_mot, top-n)
    &#8627; ratio du nombre de mots commun
    ↖ jolie contribution

Idée :
    → Evolution de la femme dans
         le word2vec
        → avec les backups de wikipedia