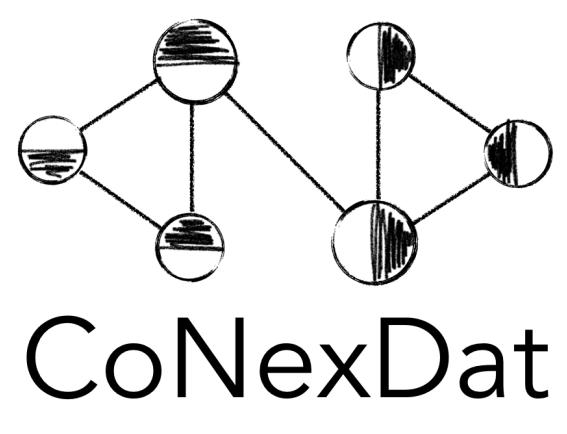


Longitudinal word2vec



CoNexDat



Carlos Selmo^(1,4)
Julián F. Martínez^(1,2,3)
Mariano G. Beiró^(1,2,3)
J. Ignacio Alvarez-Hamelin^(1,2,3)

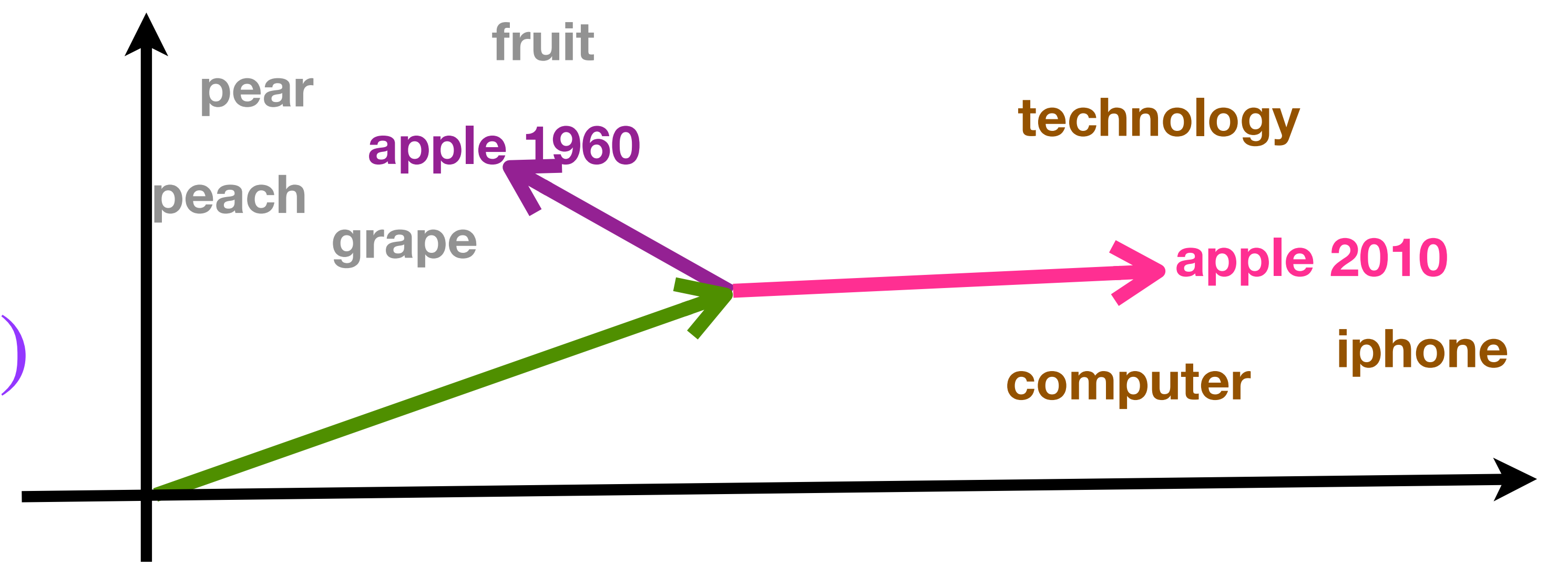
OpLaDyn
project



Introduction

How to make a **word embedding** where a word has a relative position according to its meaning on every **year**?

$$\vec{word}_{year} = \vec{common}(word) + \vec{\delta}_{year}(word)$$



Dataset

The **New York Times** newspaper has fully digitalized articles from 1970 to nowadays. We selected articles from 1990 to 2016.



Model

A **Neural Network** is trained by maximising the pseudo-likelihood:

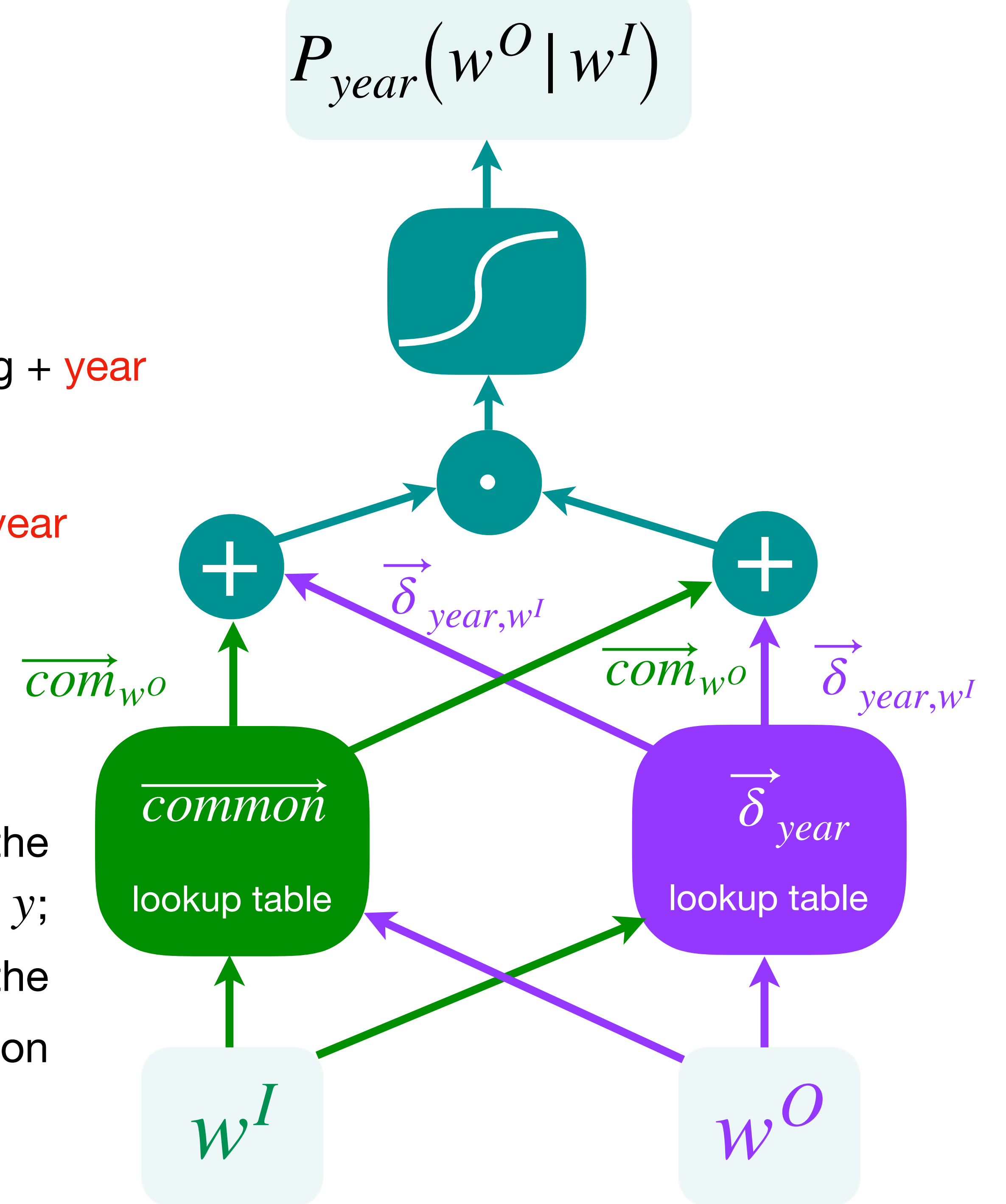
$$\mathcal{L} = \sum_{y=1}^Y [\mathcal{L}_{pos}^y + \mathcal{L}_{neg}^y + \mathcal{L}_{reg}^y]$$

$$\mathcal{L}_{pos}^y = \frac{1}{T_y} \sum_{t=1}^{T_y} \sum_{-c \leq j \leq c, j \neq 0} \log \sigma(\vec{v}'_{y,w_y,t+j} \cdot \vec{v}_{y,w_y,t}), \text{ word2vec's positive sampling + year}$$

$$\mathcal{L}_{neg}^y = \sum_{k=1}^K \mathbb{E}_{P_n(\cdot|y)} \left[\log \sigma(-\vec{v}'_{y,W_0^k} \cdot \vec{v}'_{y,W_1^k}) \right], \text{ word2vec's negative sampling + year}$$

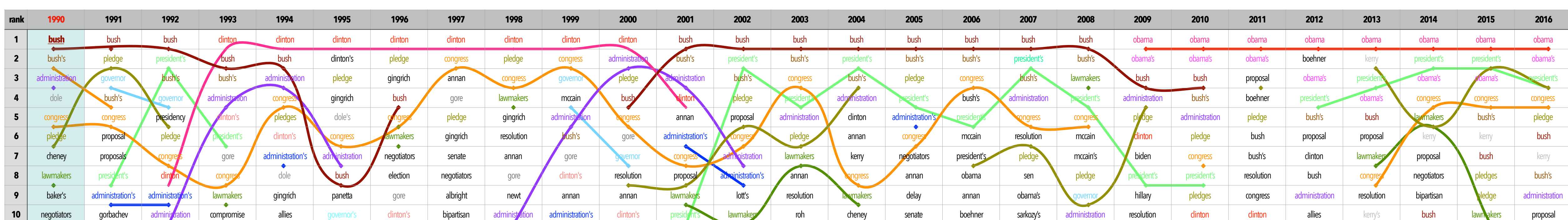
$$\mathcal{L}_{reg}^y = - \sum_{w \in V_y} \lambda_{y,w} \|\delta_{y,w}\|^2, \text{ our regularization term}$$

where: v' is the vector representing a word w ; y is the **year**; $-c \leq j \leq c$ is the word j of the context window with size c ; T_y is the size of vocabulary at year y ; W_0^k and W_1^k are independent random words sampled from $P_n(\cdot|y)$, the distribution obtained by $P_n(w|y) \propto U_y(w)^{\frac{3}{4}}$ with $U_y(w)$ the unigram distribution of the **year** y ; x^T is the vector x transposed.

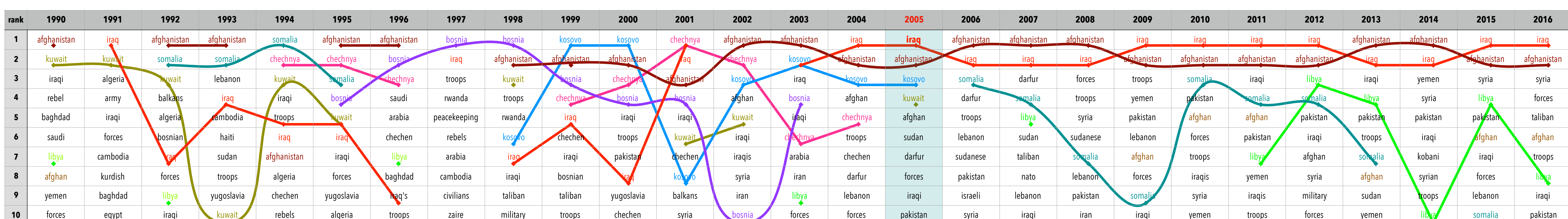


Results

A top-10 words close to **Bush 1990** along years



A top-10 words close to **Iraq 2005** along years



Git-hub
repository
for NYT
1990-2016
embedding

