SCRIBE - news summarizer MVP

For this project, I am planning to build a news summarizer/ generate the most representative sentence from a news article, after selecting similar stories from different news stories in the dataset based on:
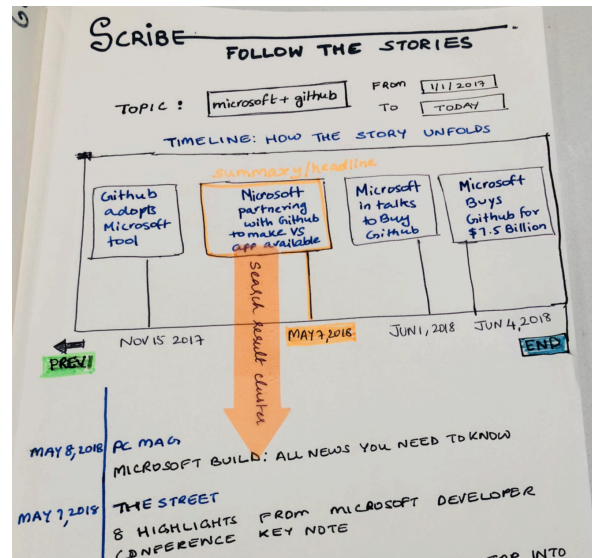1. query combination
2. date

Data
For the MVP, I have two large datasets that I am trying to examine.
Both contain:



1. Date
2. url
3. headline
4. complete text
5. summary
6. author
7. source

Known unknowns:

1. LDA
2. choosing between summarizing and extracting most representative sentence - based on my specific use case
3. summarizing technique: extraction vs abstraction
4. have specific doubts on whether to use jupyter notebooks or set up GPU cluster on aws