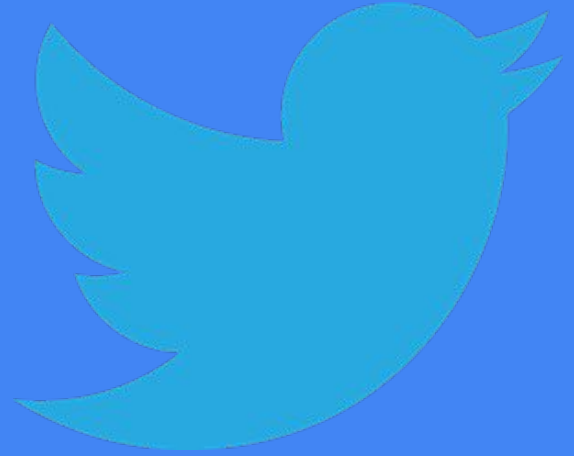
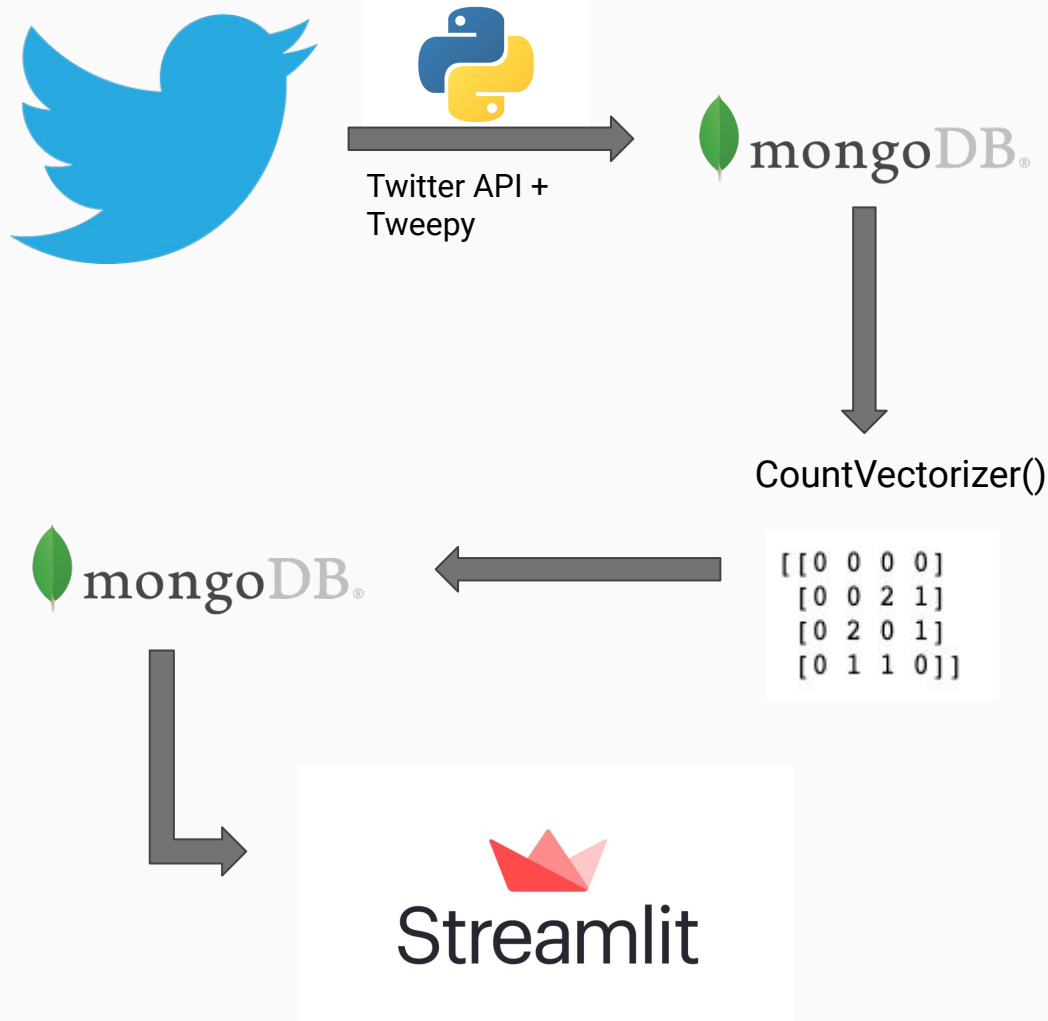


Twitter Keyword Search

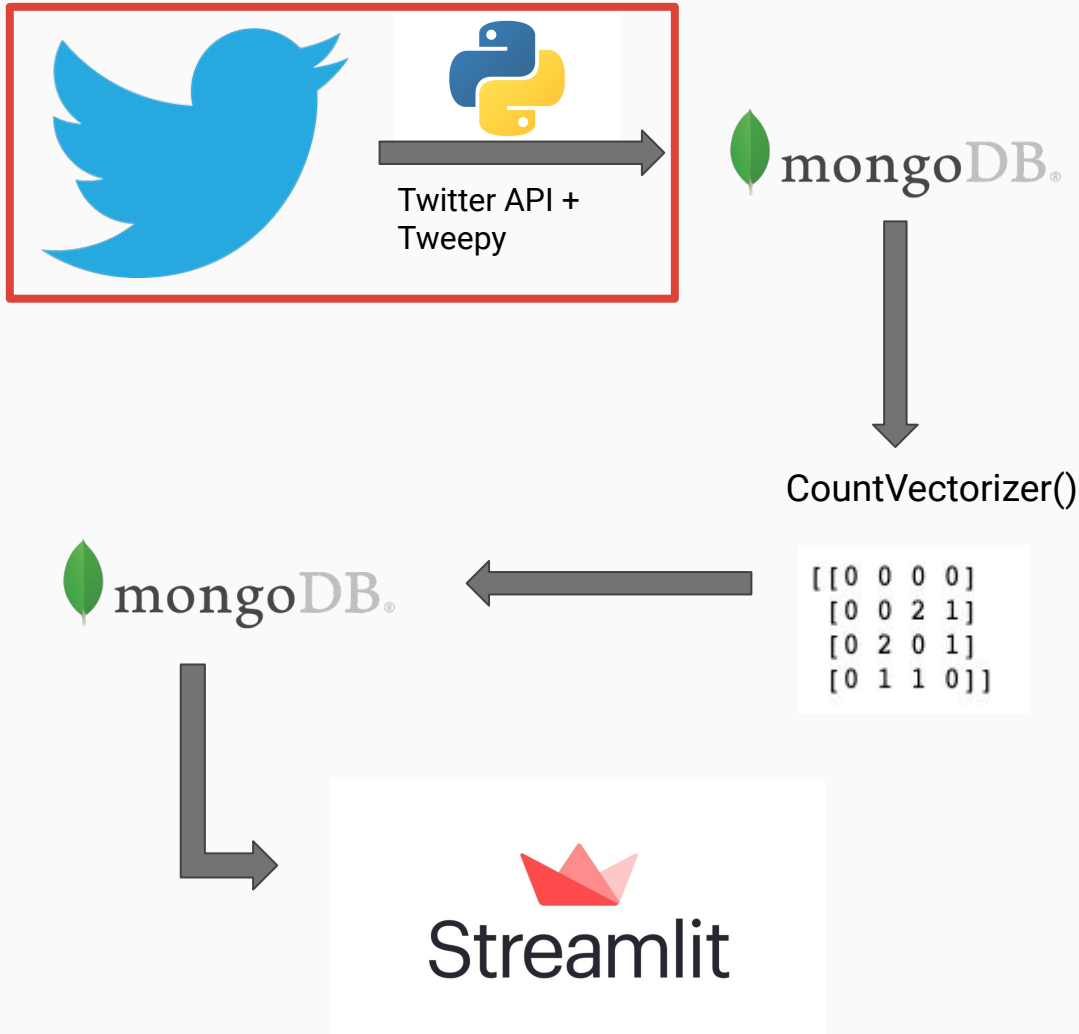


Pipeline Roadmap

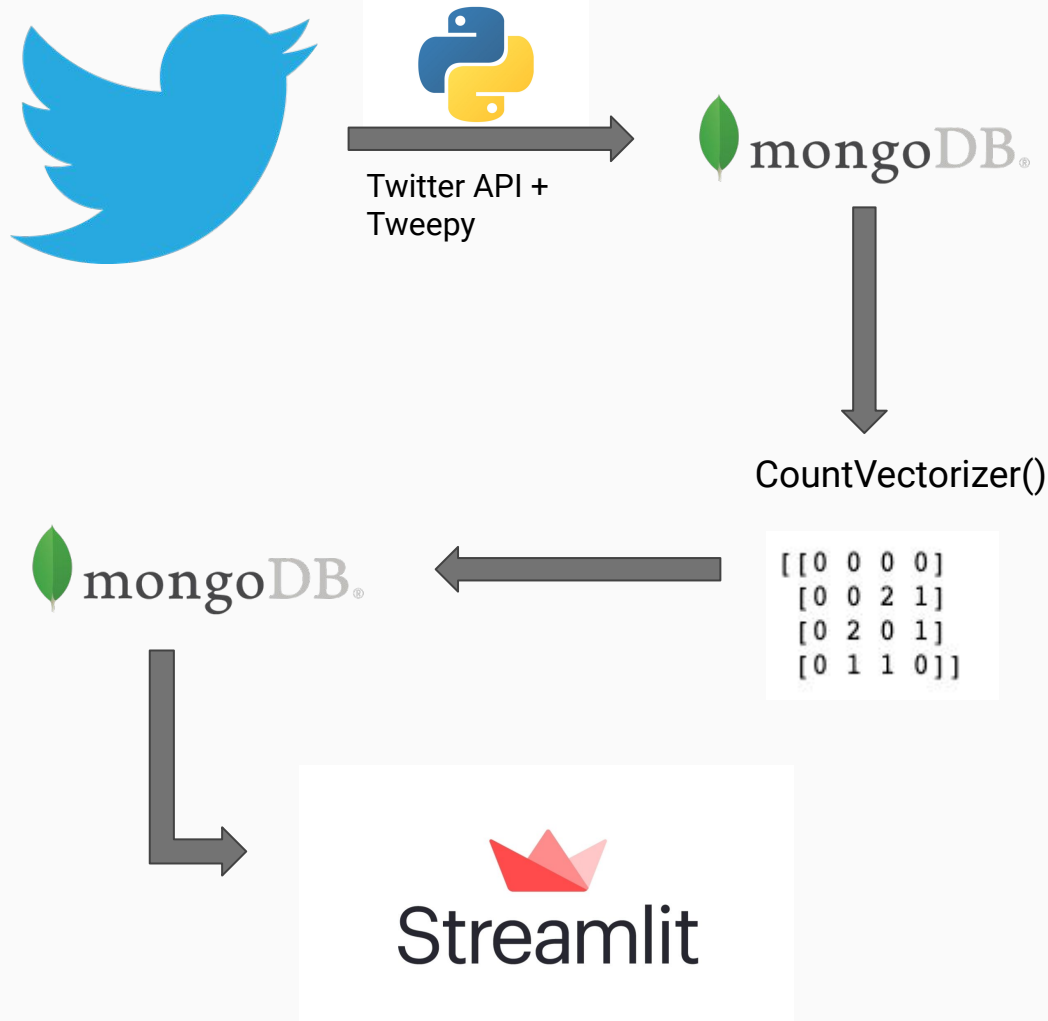


Twitter API + Tweepy

- Searched on tags: a, the, I, you, u
 - This is to avoid biasing the scraped tweets.
- Also filtered down to English

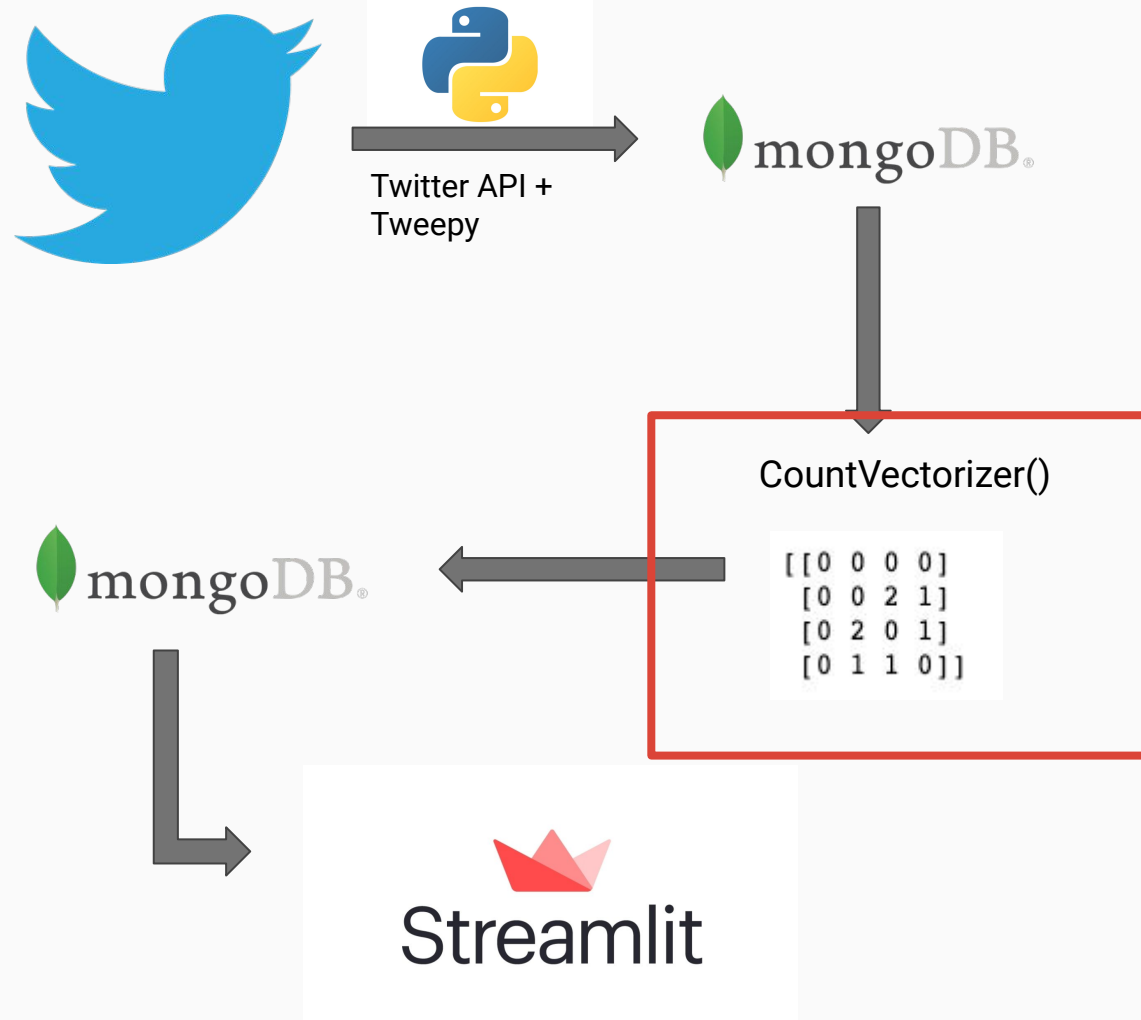


Pipeline Roadmap



CountVectorizer()

- Counts 'intra-documental' relations and stores them cumulatively in a matrix
- Better to vectorize in smaller chunks due to the sparsity of the resultant matrix
 - A lot of words are useless and just add to the run-time
- Used MongoDB to contain the equivalent of the matrix

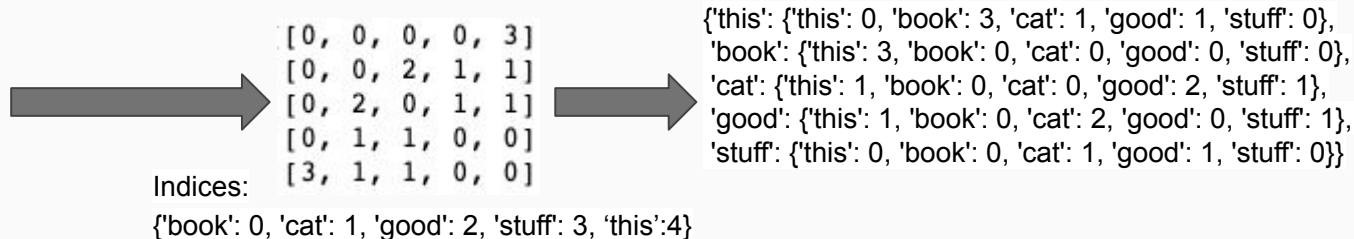


CountVectorizer()

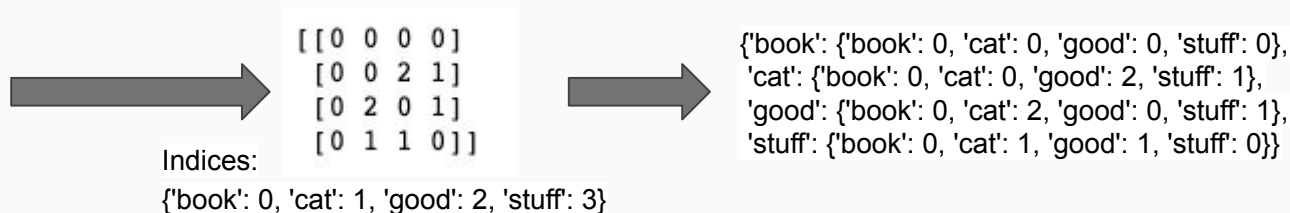
Assume:

docs = ['this this this book',
'this cat good',
'cat good stuff']

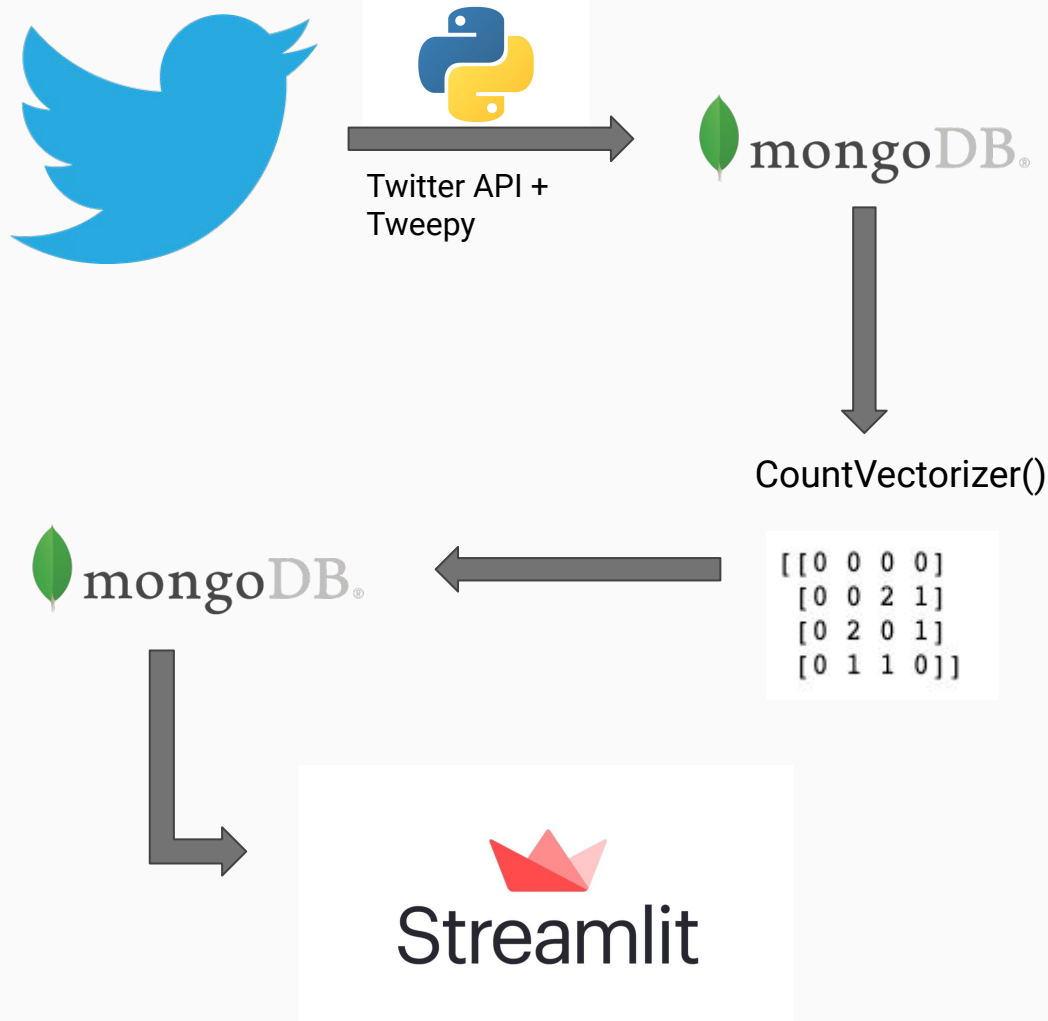
No stop words vocabulary:
[book, cat, good, stuff, this]



With stop words vocabulary:
[book, cat, good, stuff]

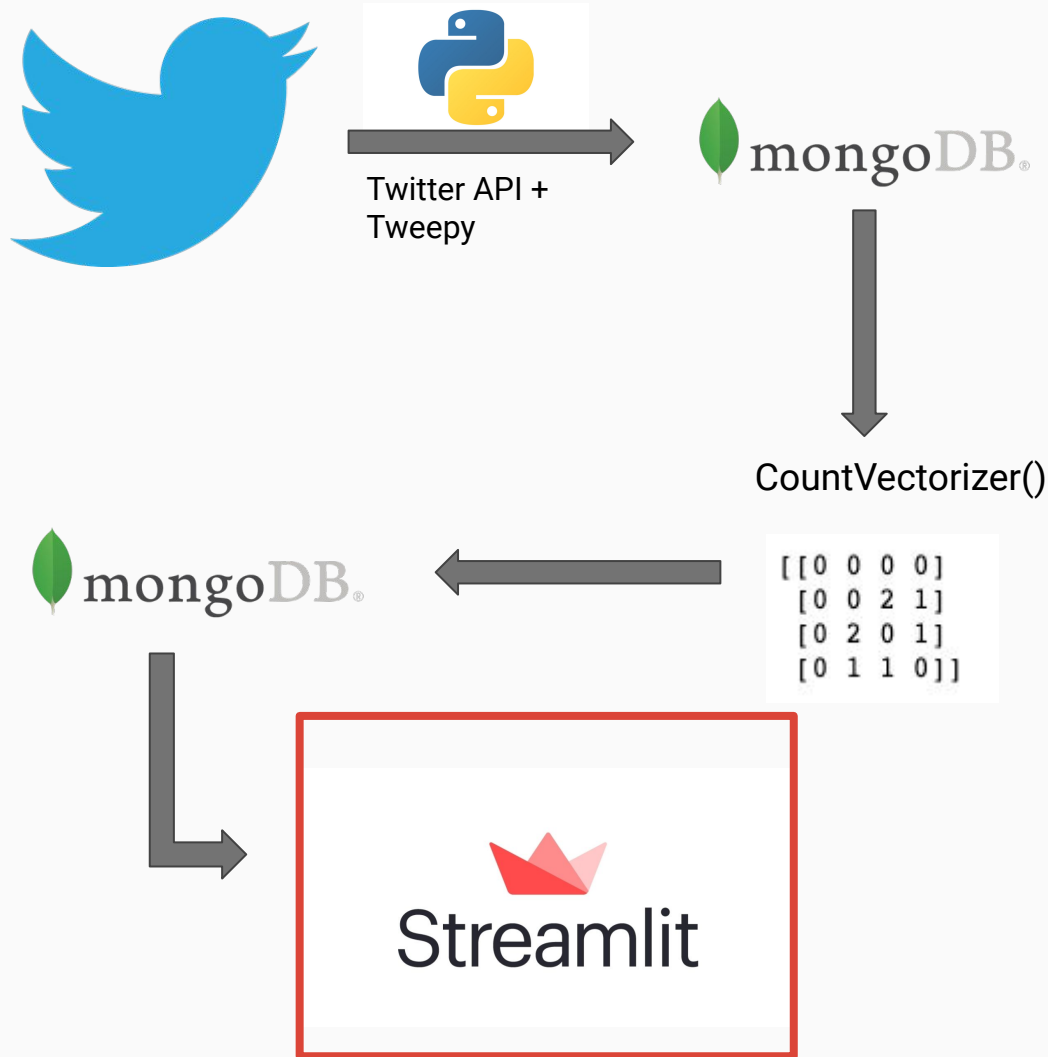


Pipeline Roadmap



Querying through Streamlit

- Vectorized strings are transformed into Mongo-readable format and stored.
- Front-end Streamlit app queries directly from the database to retrieve keywords for an inputted search term.



Video demonstration

Some Takeaways:

- Pipeline construction matters!
 - CountVectorizer() runtime vs Storage runtime
- Not everything tabular should be stored tabularly
- Twitter is a scary place.

Questions?