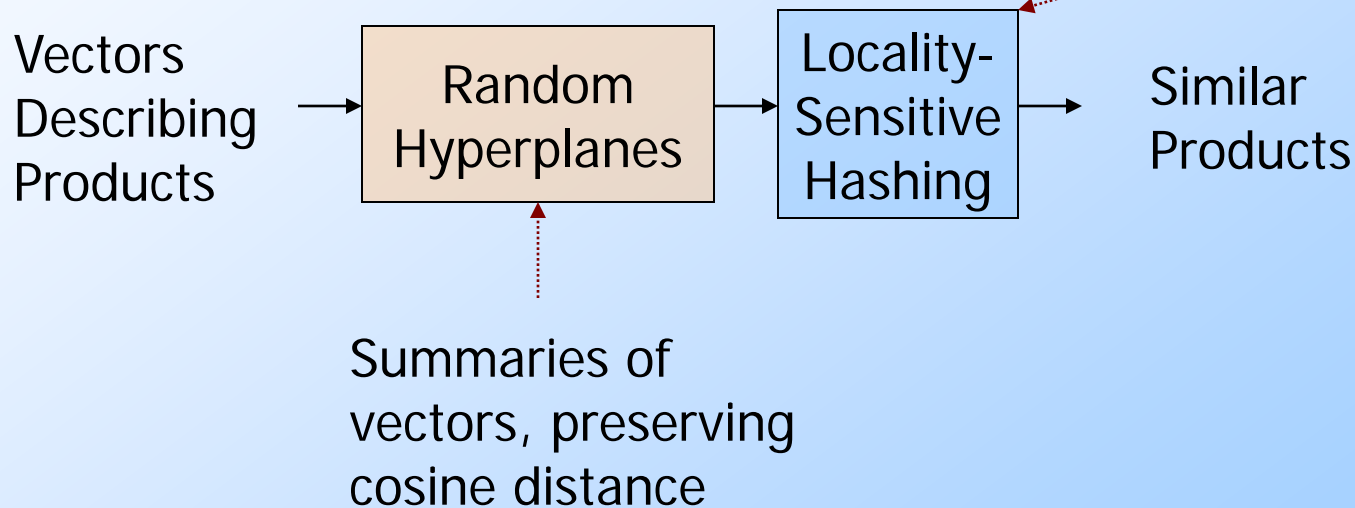
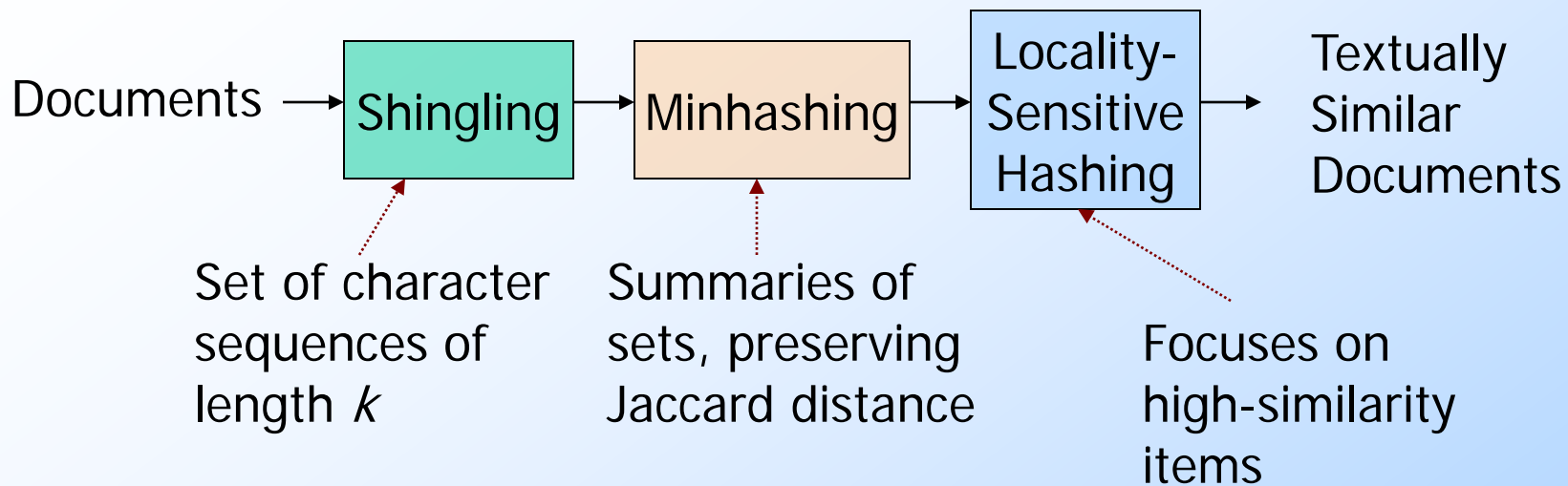


# Near-Neighbor Search

Finding Variants of a News Article

Finding Similar Restaurants



# Organizing News Articles

- ◆ Applications such as Google News see thousands of on-line news articles.
- ◆ Many come from the same story, with modifications by the publisher.
- ◆ Need to cluster by underlying story.
- ◆ Perfect for shingling – minhashing – locality-sensitive hashing.

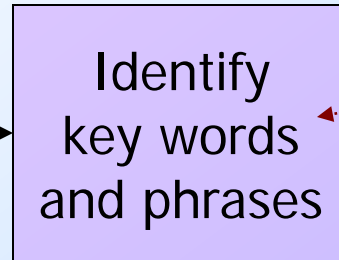
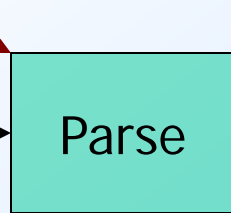
# Restaurant Advisor

- ◆ Celixis is a Stanford startup trying to create “advisors.”
  - ◆ **First application:** restaurant advisor.
- ◆ Uses data from restaurant reviews plus one or more restaurants you select as examples of what you like.

# Celixis Architecture

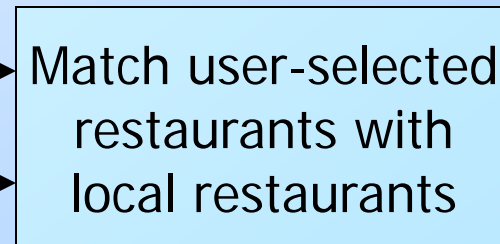
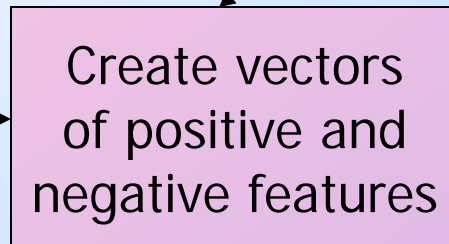
Chris Manning's NLP software

Restaurant reviews



TF.IDF-like strategy

Sentiment information



Cosine distance used for matches

Easier than general case; you don't care if "big portions" are good or bad, as long as "small portions" is the opposite.

# Further Reading

- ◆ A new book by Anand Rajaraman and Jeff Ullman, titled *Mining of Massive Datasets*, covers these topics.
- ◆ See it on-line at [infolab.stanford.edu/~ullman/pub/book.pdf](http://infolab.stanford.edu/~ullman/pub/book.pdf)