Finding near-duplicate documents



































Experiments (1996) by Broder et. al.

- 30 million HTML and text docs (150GB) from Web crawl
- 10-word shingles
- 600 million shingles (3GB)
- 40-bit shingle "fingerprints"
- Sketch using 4% shingles (variation of alg. we've seen)
- Used count of shingles for similarity
- Using threshold t = 50%, found
 - 3.6 million clusters of 12.3 million docs
 - 2.1 million clusters of identical docs 5.3 million docs
 - remaining 1.5 million clusters mixture:
 "exact duplicates and similar"

19