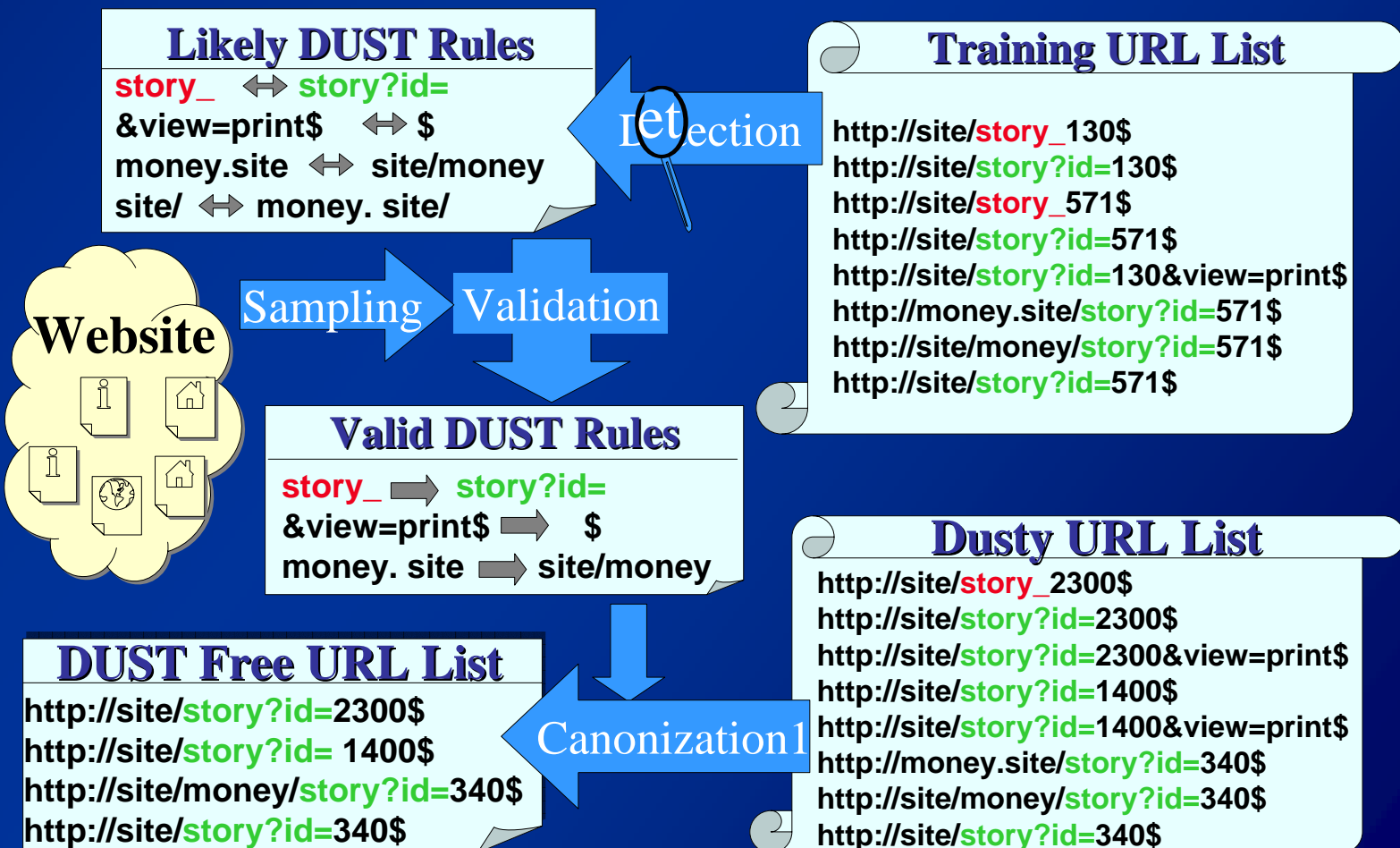


Do Not Crawl In The DUST: Different URLs, Similar Text

Uri Schonfeld, Ziv Bar-Yossef, and Idit Keidar (Technion)



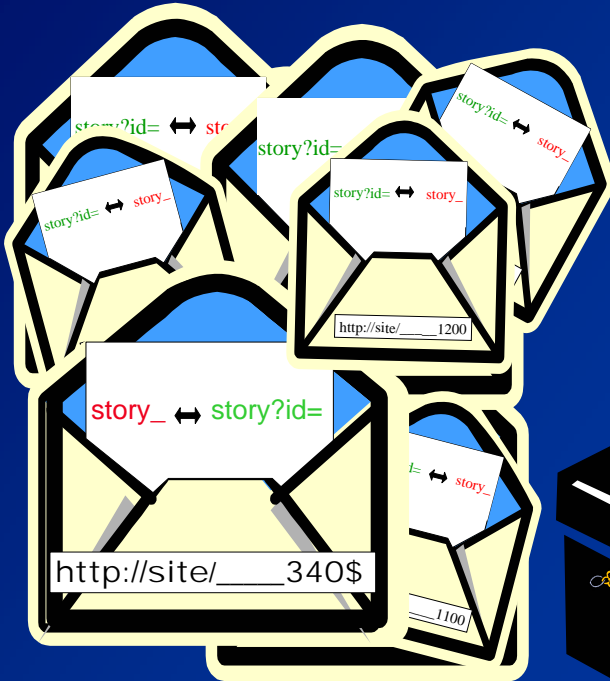
DustBuster Flow



Why Does DustBuster Work?

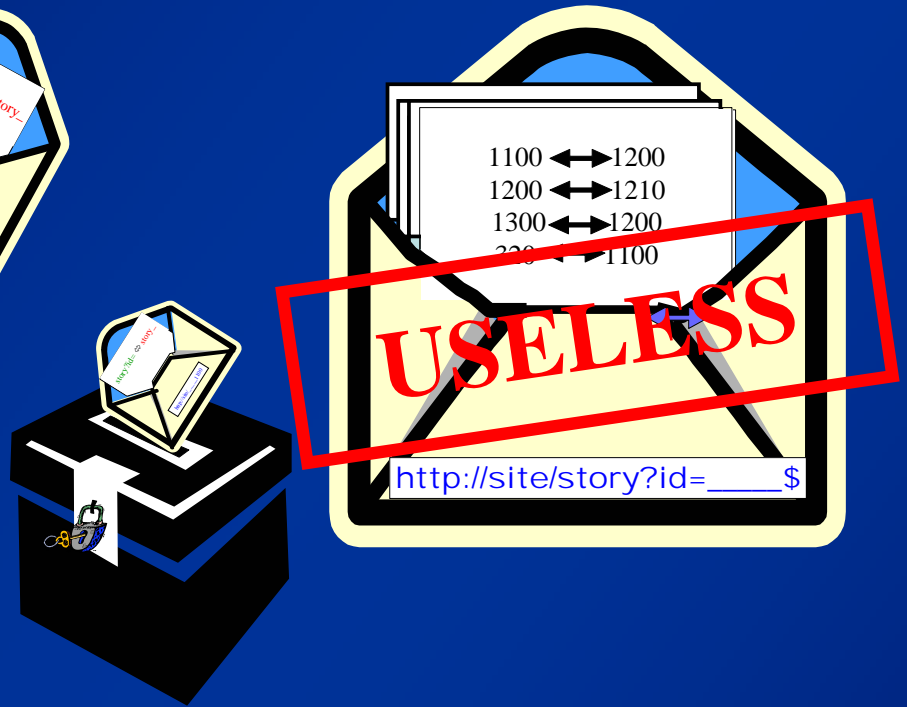
Large Support Principle

Valid Rules Get Many Votes



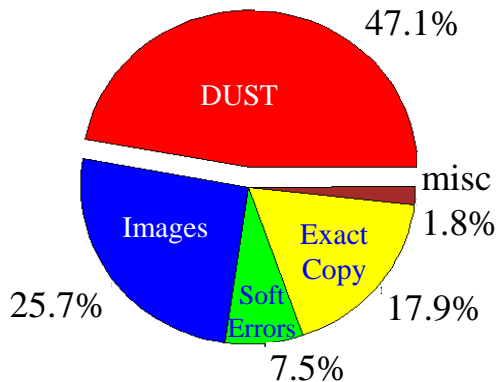
Fat Envelopes Principle

Fat Envelopes Are Useless

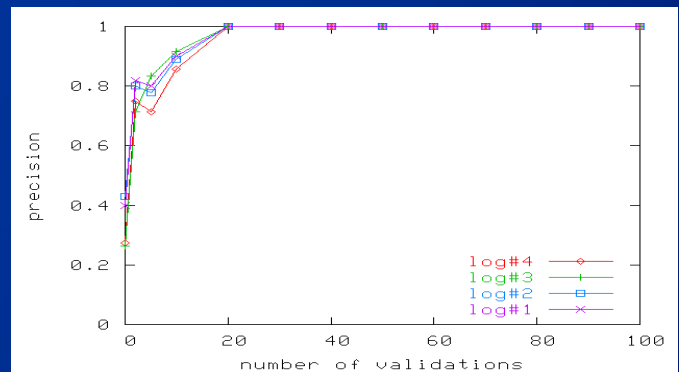


Results From Real Web Sites

DUST Distribution



Precision vs. Validations



Precision@k with and without size

