















Difference from classical IR

- ✤ Retrieve <u>high quality</u> pages:
 - Static documents: text, audio
 Dynamic documents: generated by queries on data bases via the Web (e.g., Amazon books)
- Huge collections of data (Google is one of the largest indexes, but it collects only 10% of the Web pages)
- * No persistency (data is volatile)
- Heterogeneity (type of documents, languages, no typographical control, etc.)
- * Duplication (nearly same documents, different URLs)
- * Removing non relevant info (e.g., banners)



Web and IR tools to help searching

- * Hierarchical directories (e.g., Yahoo)
- * Topic-specialized engines
- $\boldsymbol{\diamond}$ Search by example (from a set of URLs)
- Meta-information (e.g.,compare search engines)
- Clustering (from IR)
- * Categorization
- * Summarization
- * Latent semanting indexing































This document was created with Win2PDF available at http://www.daneprairie.com. The unregistered version of Win2PDF is for evaluation or non-commercial use only.