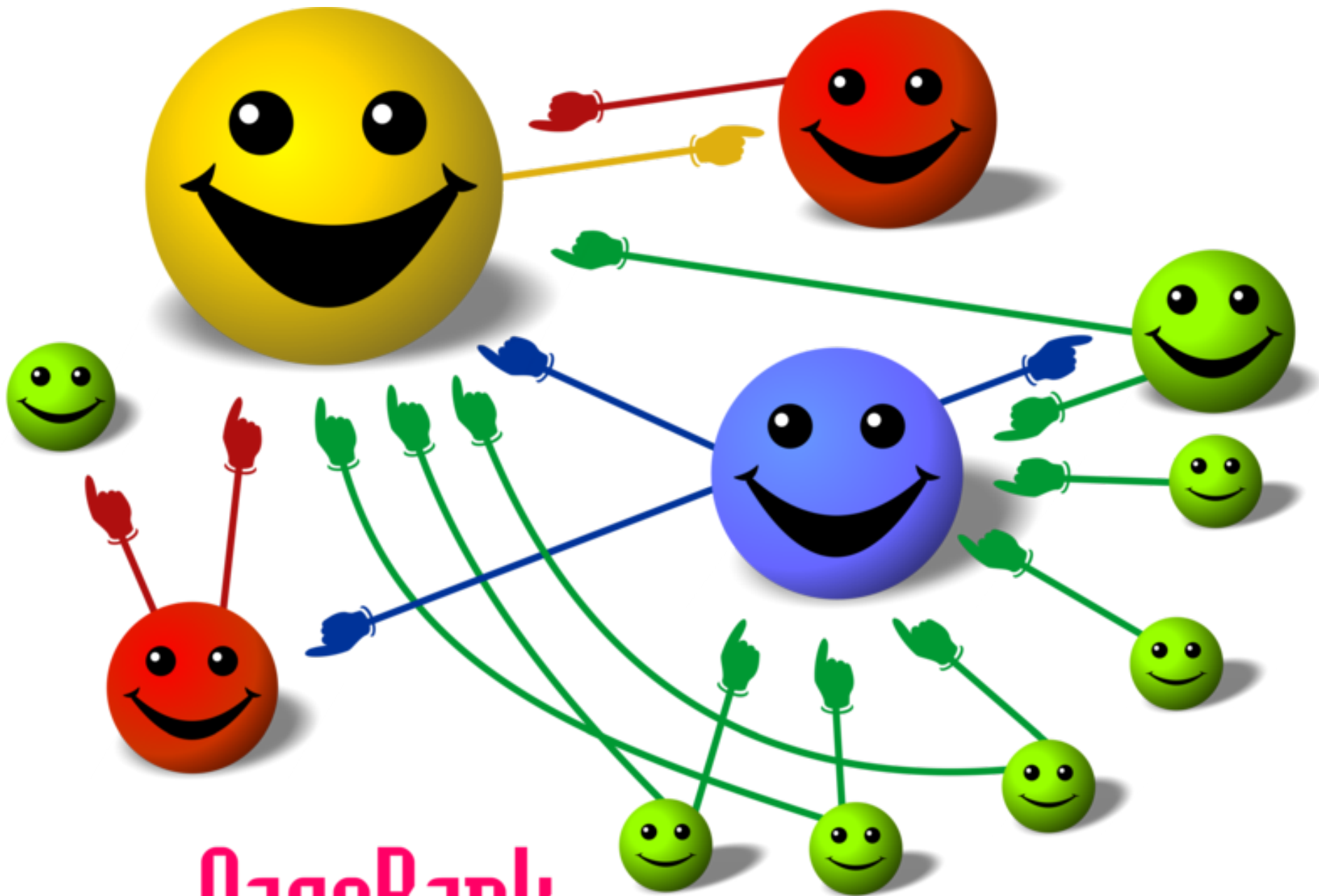# Web Search
# Information Retrieval

# "Web Data Analytics"

- Get the digital data (from web or from scanning)
  - Need to crawl web (? Solved "engineering" problem)
- Preprocess data to get searchable things (words positions)
- Form **Inverted Index** mapping words to documents
- Typically use **TF-IDF** (term frequency, Inverse Document frequency) to quantify importance of word match
- Rank relevance of documents: **PageRank**
- Lots of technology for advertising, "reverse engineering" "preventing reverse engineering"
- Clustering of documents into topics (as in Google News)

PageRank

Size of face proportional to PageRank

# Modern Recommendation Systems (from Yahoo)

- Goal (Function to Optimize – Long Term dollars)
  - Serve the right item to a user in a given context to optimize long-term business objectives

- A scientific discipline that involves
  - Large scale Machine Learning & Statistics
    - Offline Models (capture global & stable characteristics)
    - Online Models (incorporates dynamic components)
    - Explore/Exploit (active and adaptive experimentation)
  - Multi-Objective Optimization
    - Click-rates (CTR), Engagement, advertising revenue, diversity, etc
  - Inferring user interest
    - Constructing User Profiles
  - Natural Language Processing to understand content
    - Topics, "aboutness", entities, follow-up of something, breaking news,…

Recommend search queries

Recommend packages:
Image
Title, summary
Links to other pages

Pick 4 out of a pool of *K*
*K* = 20 ~ 40
Dynamic

Routes traffic other pages

Recommend applications

Recommend news article

# Some examples from content optimization

- ## Simple version
  - I have a content module on my page, content inventory is obtained from a third party source which is further refined through editorial oversight. Can I algorithmically recommend content on this module? I want to improve overall click-rate (CTR) on this module

- ## More advanced
  - I got X% lift in CTR. But I have additional information on other downstream utilities (e.g. advertising revenue). Can I increase downstream utility without losing too many clicks?

- ## Highly advanced
  - There are multiple modules running on my webpage. How do I perform a simultaneous optimization?