

Mining Search Engine Query Logs via Suggestion Sampling

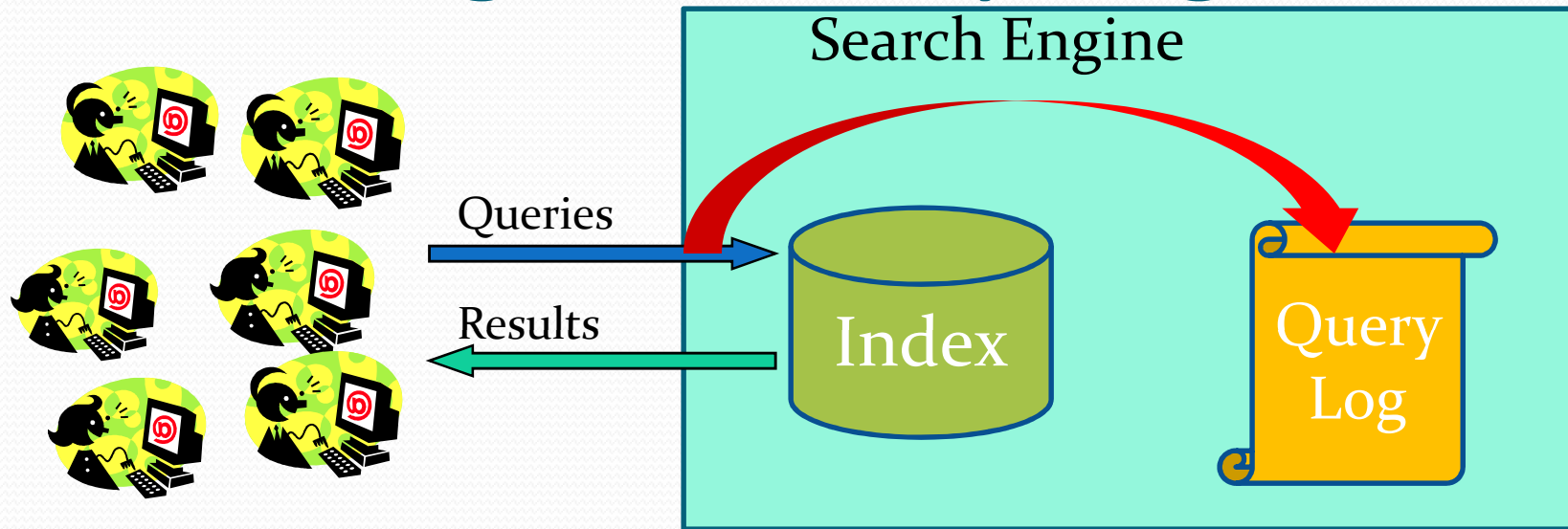
Ziv Bar-Yossef

Technion and Google

Maxim Gurevich

Technion

Search Engine Query Logs



- Used by search engines to improve search results
- Contains private information
 - Of users **and** of the search engine itself
- Not disclosed by search engines



Applications of Query Log Analysis

- Keyword based advertising
- Search quality evaluation
- User modeling

Keyword Based Advertising

- Compare keyword popularity
- Track keyword popularity over time
- Find related keywords



Web Results 1 - 10 of about 591,000 for vldb. (0.56 seconds)

[VLDB Endowment Inc.](#)

Very Large Data Base Endowment Inc.: a non-profit organisation for promoting and exchanging scholarly work in databases and related fields.

www.vldb.org/ - 15k - [Cached](#) - [Similar pages](#)

[VLDB Conference Hints](#)

Organizing a **VLDB** conference is a demanding task, and anyone underestimating this task will find himself in exasperation, frequently causing irritation and ...

www.vldb.org/hints.html - 31k - [Cached](#) - [Similar pages](#)

[VLDB 08 - VLDB 08](#)

34th International Conference on Very Large Data Bases. Auckland, New Zealand, 24-30, August, 2008.

www.vldb2008.auckland.ac.nz/ - 13k - [Cached](#) - [Similar pages](#)

Sponsored Links

[The New Fastest VLDB](#)

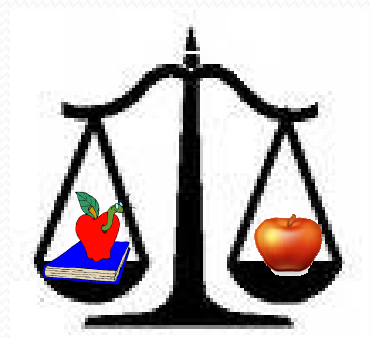
Vertica: A Faster Analytic Database
Read our Winter Corp. Review Now!
www.vertica.com

[Very Large Database](#)

Scales to 30TB with Infobright
Data Warehouse. Free Whitepaper
www.InfoBright.com

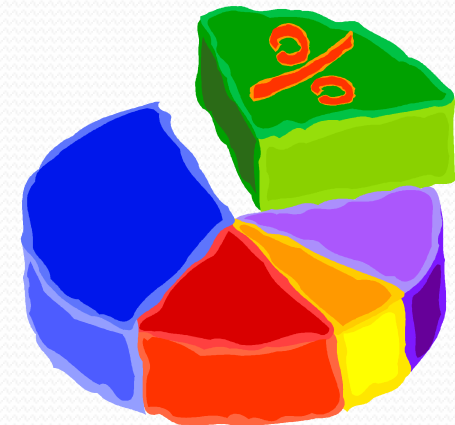
Search Quality Evaluation

- Estimate the amount of undesirable content sent to users
 - Spam
 - Stale results
 - Non-existent results
 - Pornography
 - Hate materials
 - Virus contaminated pages
 - ...
- Estimate search engine bias towards
 - Authoritative sources
 - Certain domains
 - Certain languages
 - ...

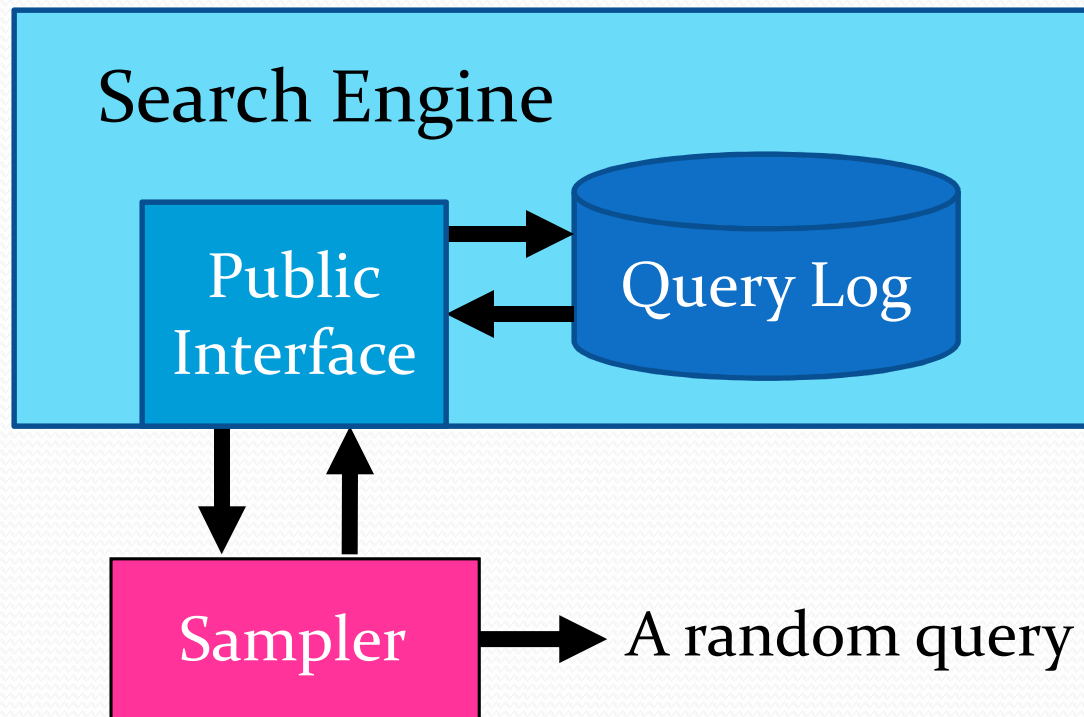


User modeling

- Distribution of query types [Broder 02]
 - Navigational
 - Informational
 - Transactional
- Density of commercial queries
- Fraction of geographical queries
- ...

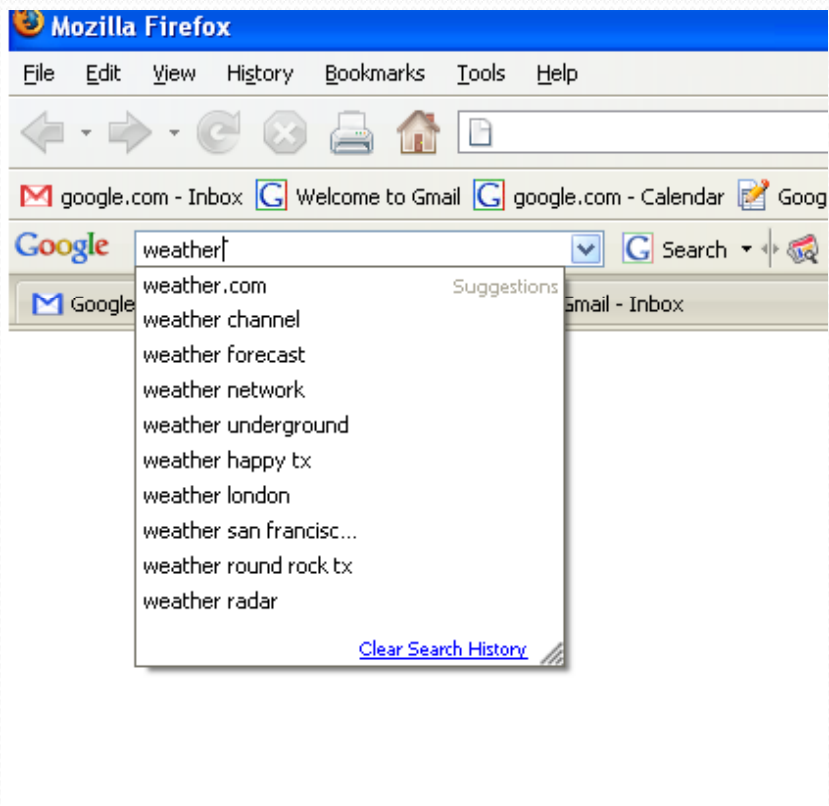


External Query Log Mining



- Sampling (uniform or by popularity)
- Computing aggregate (privacy preserving) functions

Suggest: Trapdoor to Query Logs



- Query auto-completion
 - Suggests query completions
 - More popular first
 - Offered by major search engines
- Backed by a hidden underlying database
 - Assumption: Derived from query logs

Our Contribution

- Algorithm for sampling queries **uniformly** from the query log using the suggestion service
 - Practical (few suggestion requests)
 - Unbiased
- Algorithm for sampling queries from the query log **proportionally to their popularity** using the suggestion service
 - Practical (few suggestion requests)
 - Slightly biased



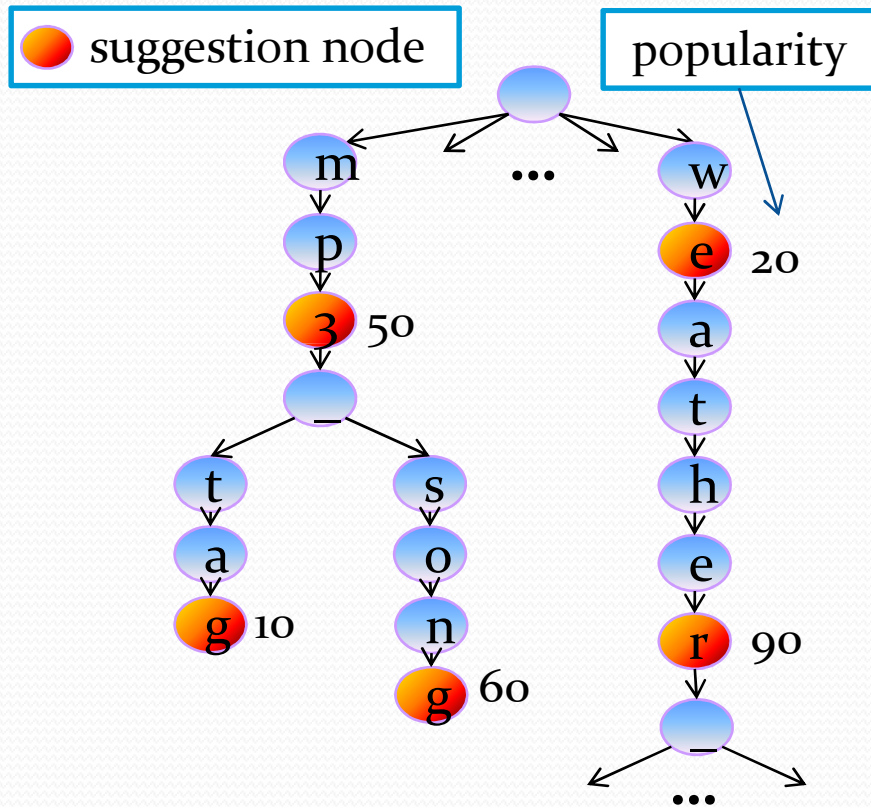
Focus of
this talk

Related Work

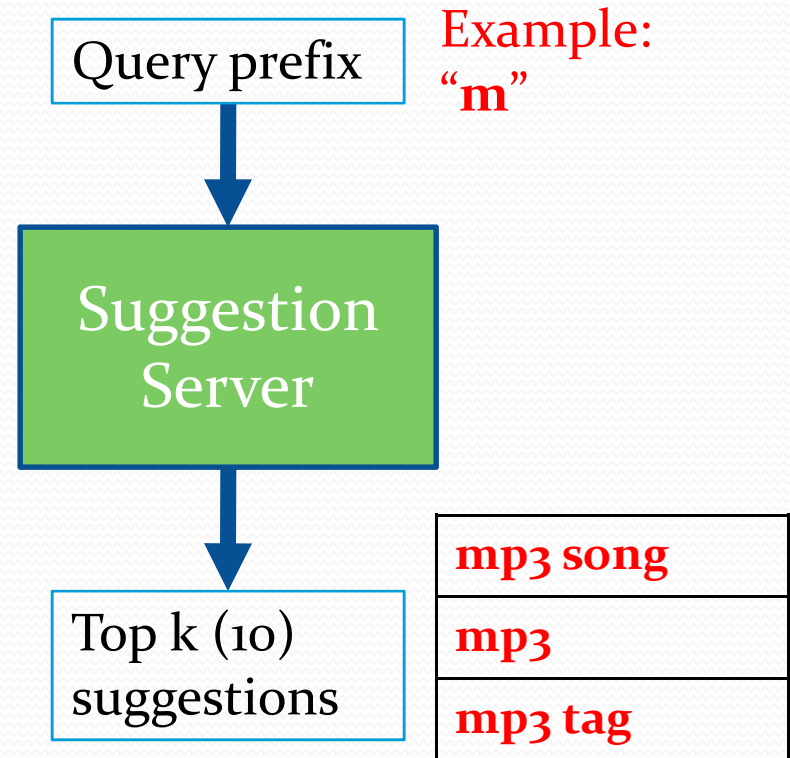
- Sampling documents from search engine index
[BarYossef et al 06,07, Broder et al 07, ...]
 - Different problem and setting
- Sampling from B-trees
[Wong et al 80, Olken et al 89,95]
 - B-tree specific assumptions
 - Inefficient for query log mining
- Sampling from databases behind web forms
[Dasgupta et al 07]
 - Different setting, inefficient for query log mining
- Uniform sampling of combinatorial structures
[Jerrum et al. 86]
 - Theoretical, the basis of our sampling algorithm

Uniform Suggestion Sampling

Suggestion TRIE



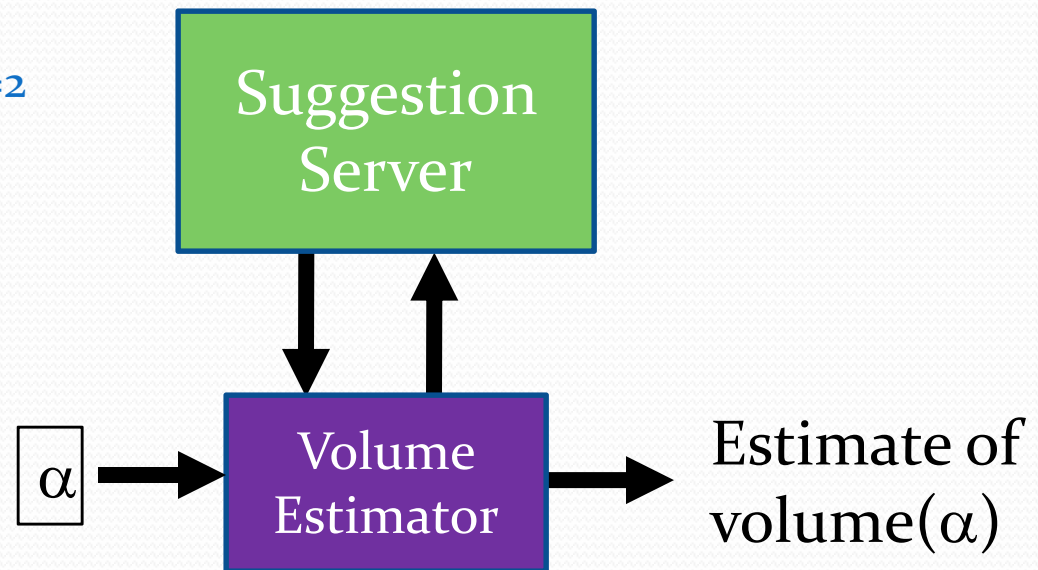
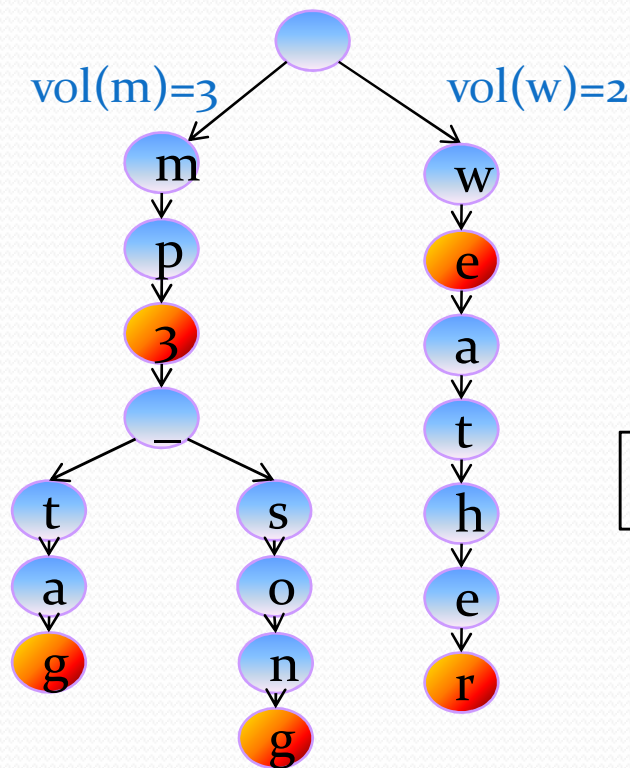
Suggestion Server



- Goal: Sample suggestion nodes uniformly

Volume Estimators

- Define: $\text{volume}(\alpha)$ = # of suggestions starting with α

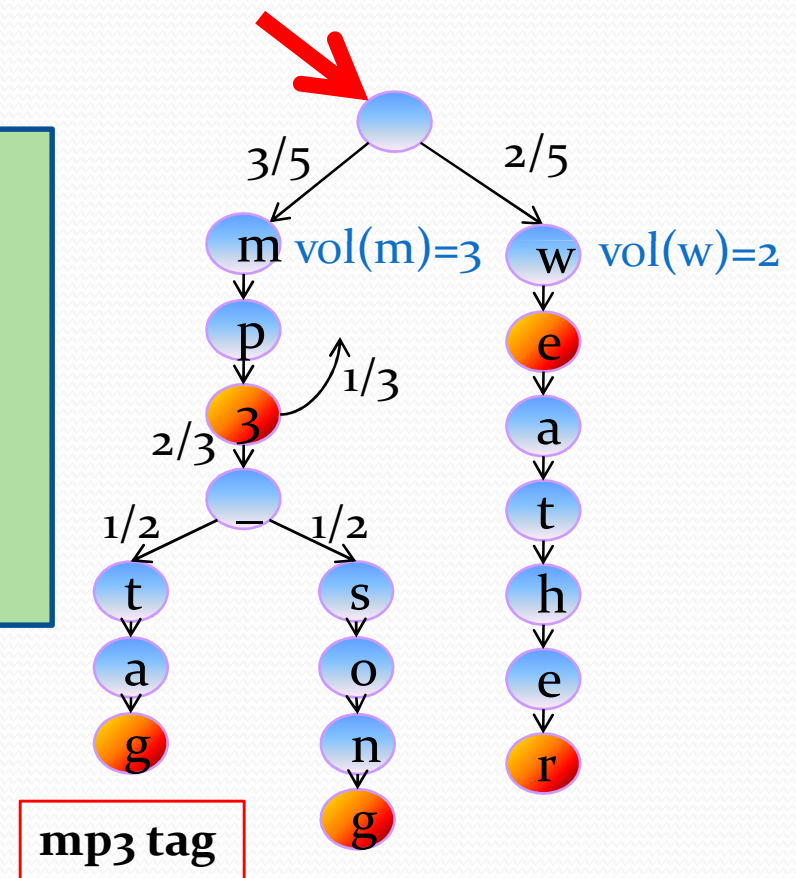


Random Walk Tree Sampling

- **Assumption:** we have a perfect volume estimator

- **current** = root
- while true
 - If current is a suggestion node
 - Return **current** with probability $1/\text{volume}(\text{current})$
 - Go to child x with probability $\propto \text{volume}(x)$

- **Theorem:** If volumes are accurate, the samples are uniform



How To Estimate Volumes

- Input: Prefix string α
- Output: \approx # of suggestions starting with α
- **Naïve estimator**: $\text{volume}(\alpha) \approx$ # of suggestions the server returns on α
- **Popularity based estimator**: $\text{volume}(\alpha) \approx$ popularity(most popular suggestion for α)
 - Rationale: Power Law distribution of popularity
(procedure for popularity estimation is included in the paper)
- **Sample based estimator**: $\text{volume}(\alpha) \approx$ normalized number of suggestions for α in a previously available query log
- Final estimator aggregates all the three results

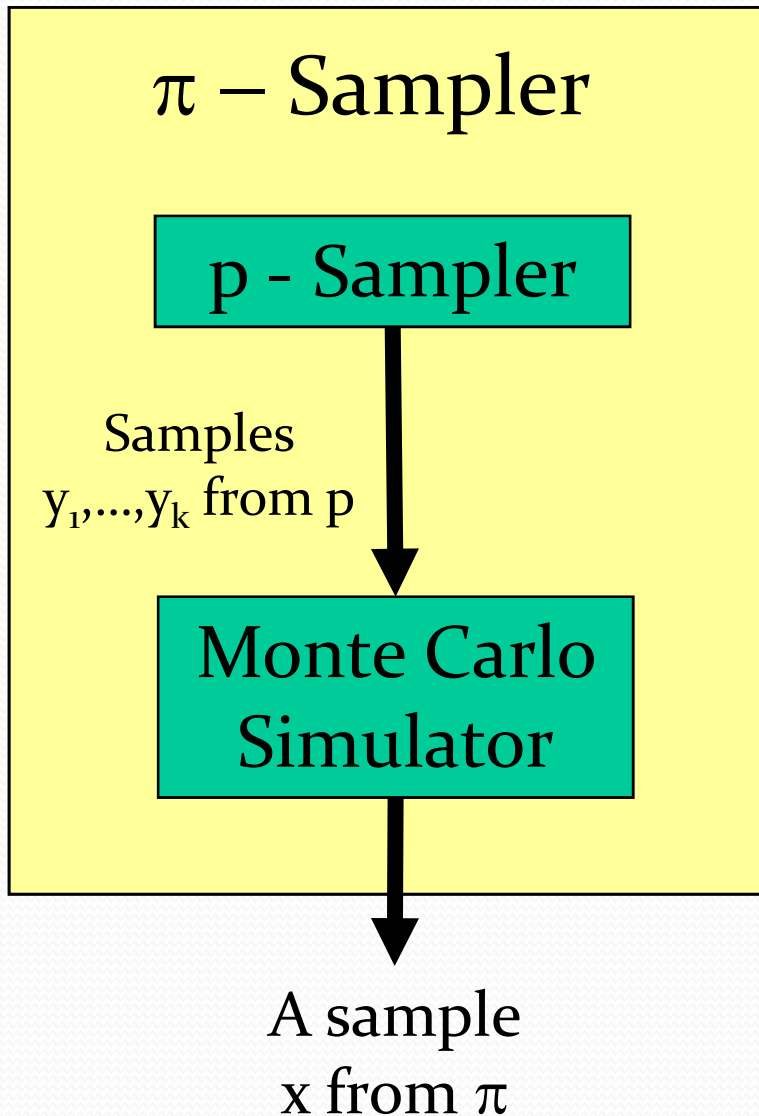
Caveat

- Random Walk Tree Sampler assumed volumes are known perfectly
- We can only approximate volumes (heuristically)



- Suggestion samples are not uniform

Monte Carlo Stochastic Simulation



- π = **target distribution** = uniform distribution on suggestions
 - Can deal with other target distributions as well
- p = **trial distribution** on S
 - Can compute $p(x)$ for each x
 - $p \neq \pi$ but $\text{support}(p)$ should contain $\text{support}(\pi)$
 - p should be easy-to-sample-from

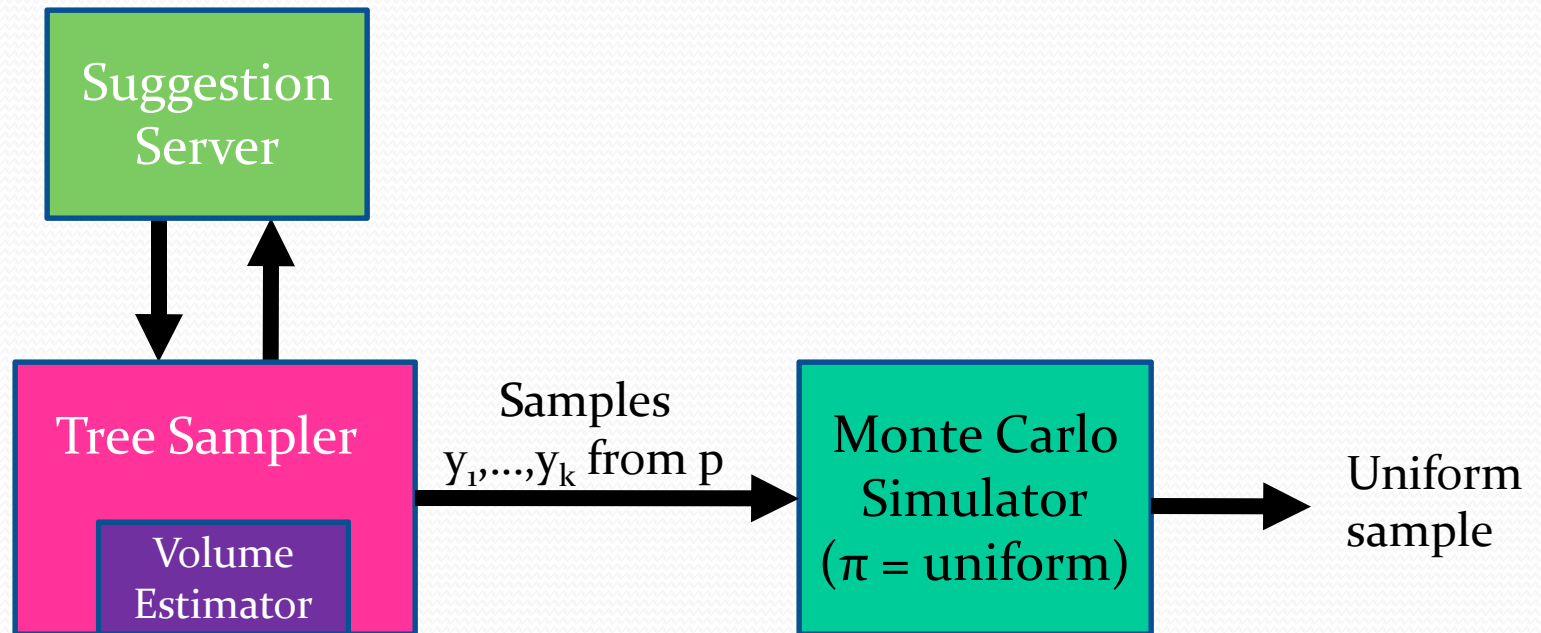
Rejection Sampling

[von Neumann 63]

- `accepted` := false
- while (not `accepted`)
 - Sample suggestion q from p
 - Calculate $p(q)$ and $\pi(q)$
 - Toss a coin whose heads probability is $\frac{\pi(q)}{C \cdot p(q)}$
 - if coin comes up heads, `accepted` := true
- return q

- $\Pr(q \text{ accepted}) = p(q) \cdot \frac{\pi(q)}{C \cdot p(q)} \propto \pi(q)$

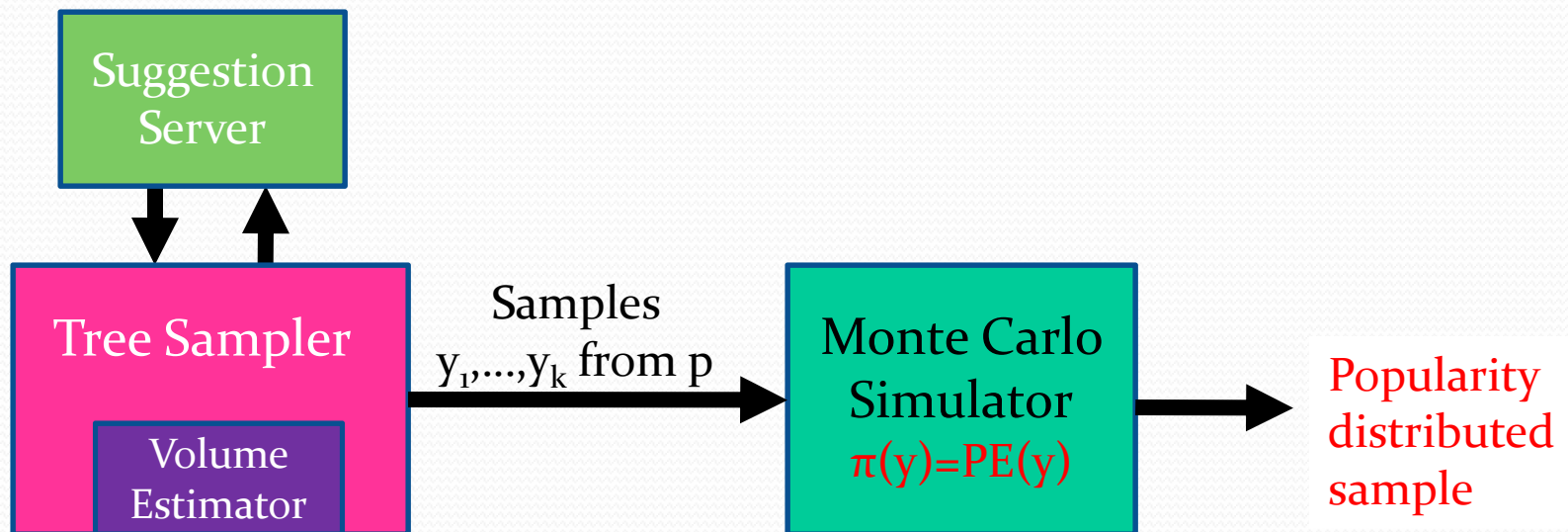
Recap – Uniform Sampling



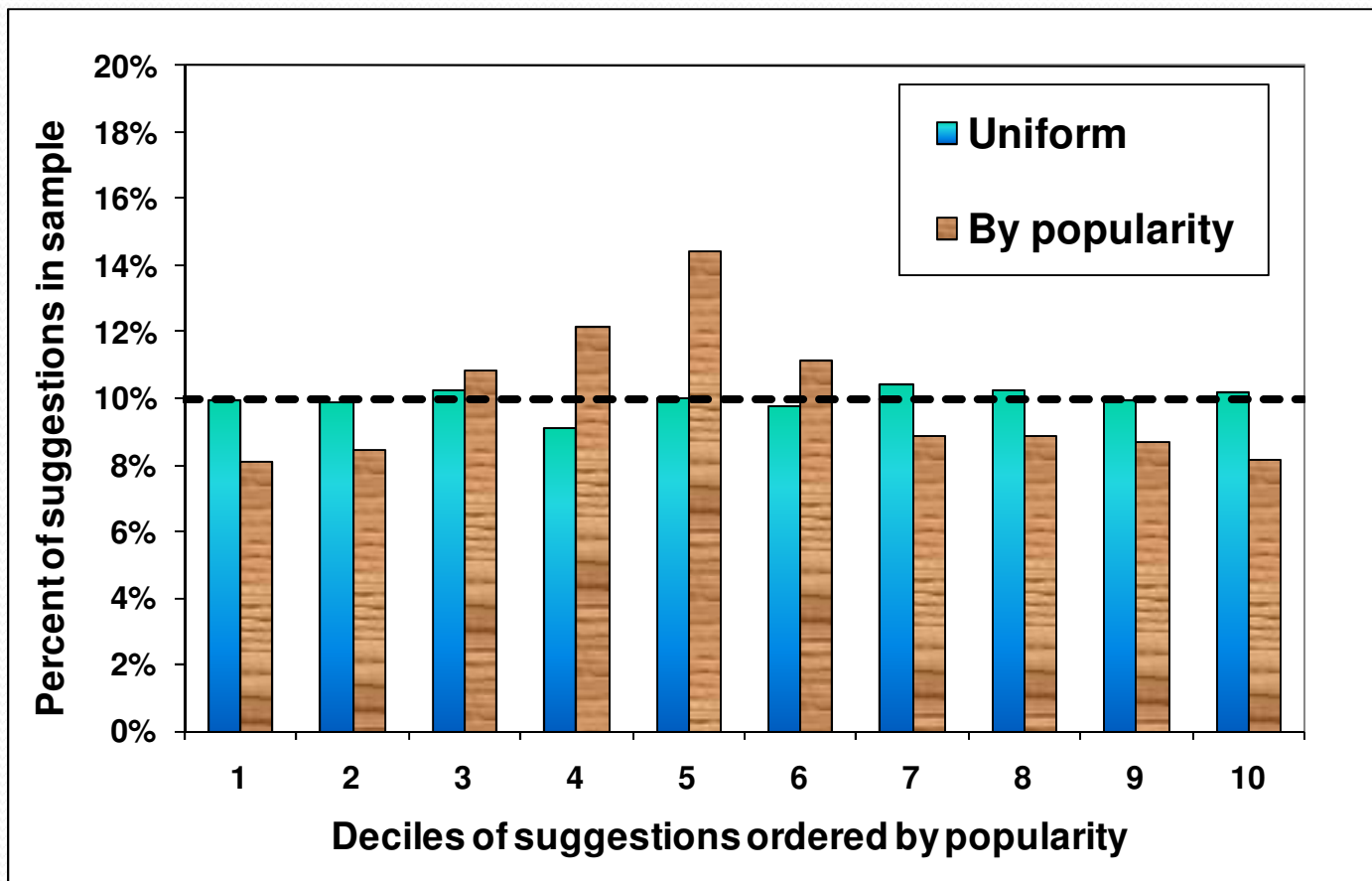
- ~6000 suggestion server requests per uniform sample

Popularity based suggestion sampling

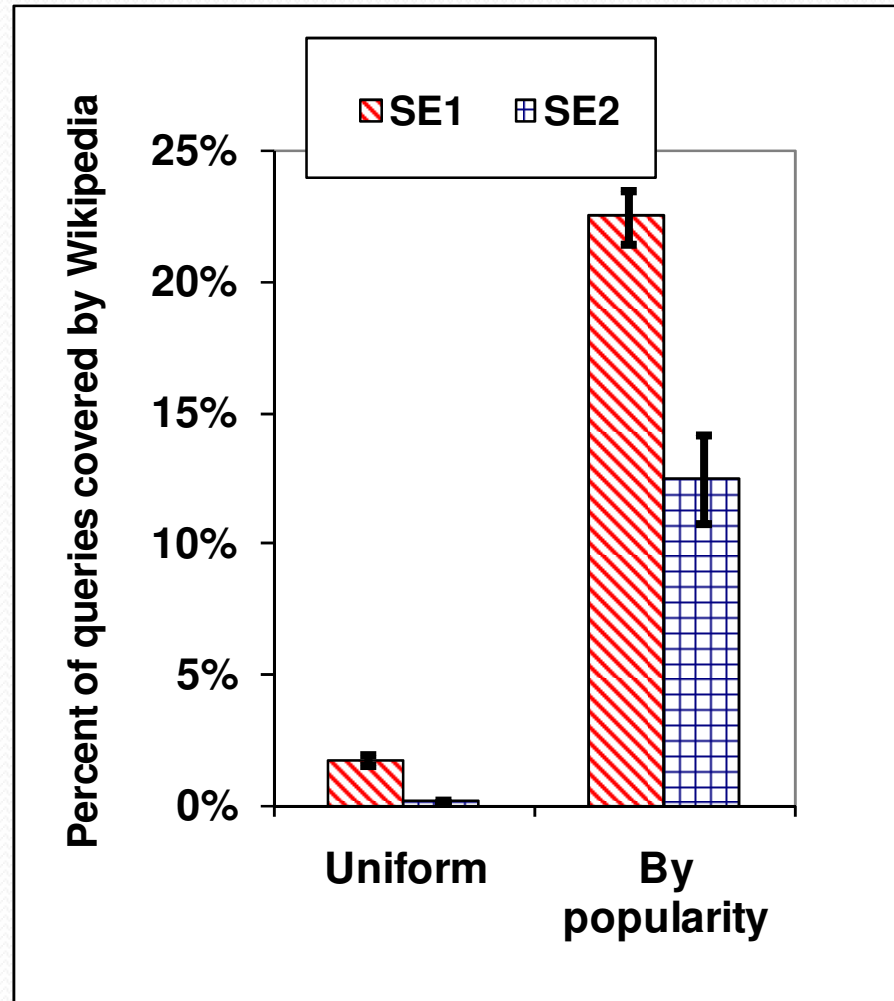
- Assumptions
 - Popularity is distributed according to Power Law
 - The Power Law exponent is known a priori
- Basic building block: Popularity Estimator (PE)



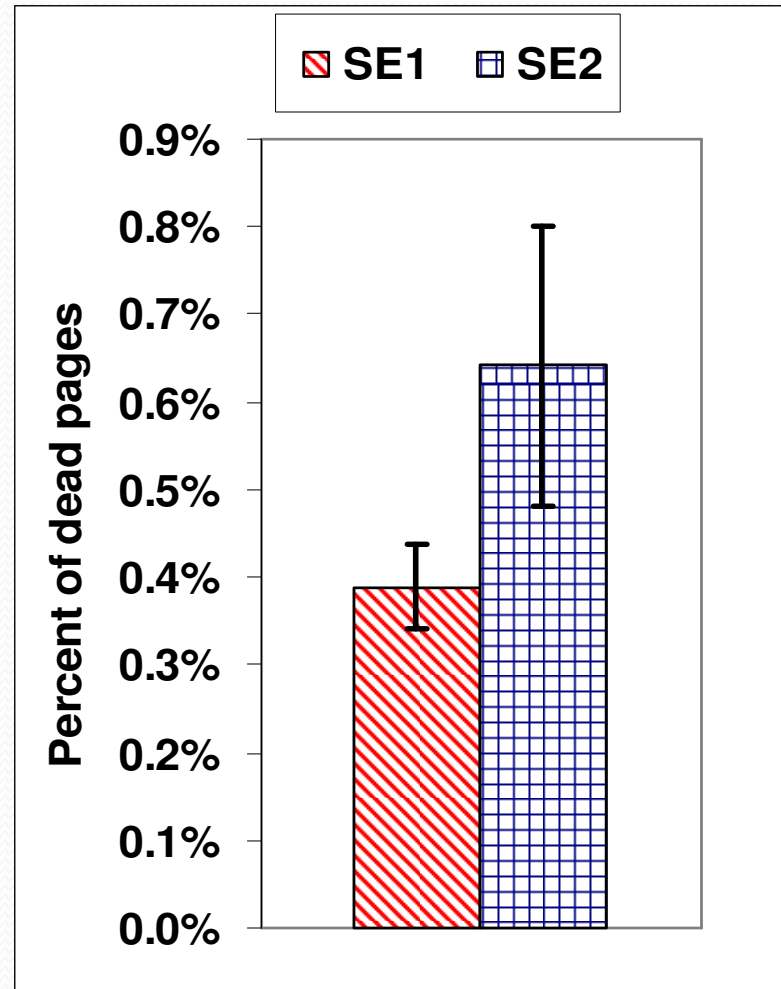
Sampling Bias



Coverage of sampled queries by Wikipedia



Percent of non-existent search results





Conclusions

- Algorithms for sampling queries randomly from a search engine query log
 - Uniformly or by popularity
 - Useful for keyword based advertising, search engine evaluation, user behavior studies
 - Via public suggestion interface only
 - Practical (accurate and efficient)



Thank You