# Experiences with Selecting Search Engines Using Metasearch

**Daniel Dreilinger**
MIT Media Laboratory

**Adele E. Howe**
Colorado State University

1997.10.23

---

## Outline

λ **Background**
λ **Main Idea**
λ **Evaluation**
λ **Conclusion**

# Background

λ **Motivation**

- the advent of many search engines on the Internet
- no single search engine is likely to return more than 45% relevant results

λ **Goal**

- develop the metasearch engine that can <u>automatically</u>, <u>carefully</u>, and <u>simultaneously</u> query several Internet search engines
  - ⇨ **minimize resource consumption**
  - ⇨ **maximize search quality**

# Background

λ **Architecture of metasearch engines**

# Background

λ **SavvySearch query form**

# Background

λ **Search plan**

# Background

λ **Key features**
- – <u>search plan</u> approach: searching advice
- – <u>metaindex dispatch</u> approach: search engine selection

λ **Issues**
- – the Web is indexed by the other search engines
- – both <u>general</u> and <u>specific</u> search engines are involved
- – the capabilities of search engines change regularly
- – to be a good citizen of the Web, <u>resource consumption</u> must be balanced against <u>results quality</u>

# Main Idea

λ **Solutions**
- – a metaindex tracks prior query experiences
- – rank the search engines for each query
- – control the degree of parallelism

λ **Metaindex**
- – a <u>t X n</u> matrix
- – value in a cell is a signed number
  - ⇨**positive: good performance**
  - ⇨**negative: bad performance**
- – accumulate user feedback passively
  - ⇨**Visit, No Result**

# Main Idea

λ **Ranking**
- based on information in metaindex
  - ⇨ **for IR documents:**
  - ⇨ **for search engines:**

- based on the recent performance of the search engines
  - ⇨ **penalty of hits:**
  - ⇨ **penalty of reponse time:**
- overall rank for <u>search engine s</u> and <u>query q</u> is

# Main Idea

λ **Concurrency**
- expected network load: <u>Web server log</u>
- local CPU load: <u>UNIX uptime</u>
- discrimitive values: <u>specific/general measure</u>

*5*

# Evaluation

λ **Experiment I.**

– Approach A: group order, selection order
  ⇨ **Visit: 2, Self-report: 72%**

– Approach B: random group order, selection order
  ⇨ **Visit: 1.76, Self-report: 60%**

– Approach C: group order, random selection order
  ⇨ **Visit: 1.89, Self-report: 65%**

– Approach D: both random
  ⇨ **Visit: 1.55, Self-report: 60%**

– quality is significantly improved by the search engine ranking: A>C, B>D

# Evaluation

# Evaluation

λ **Experiment II.**
- Were the additional knowledge meant better performance?
  ⇨ **No Result: short-term**
  ⇨ **Visit: long-term**
- If some of the search engines are truly <u>comprehensive</u>, metasearch might be unnecessary

# Evaluation

# Evaluation

λ **Experiment III.**
- comparison with <u>preprogrammed design</u>
  - ⇨ **Visit: 46%, No Result: 12%**
  - ⇨ **Visit: 40%, No Result: 14%**
    - ➔ Visit: 100~200, No Result: 4~5
- require considerable experience with a word to surpass programmed approach on the <u>Visit</u> measure, but only a few on the <u>No Result</u> measure

# Conclusion

λ **Critiques**
- <u>general framework</u> of the metasearch engines is not consistent with text
- terms in <u>mutiword queries</u> may have different contributions with a <u>Visit event</u>
- the penalty of reponse time may be related to the <u>queries</u> expect for the <u>search engines</u>

λ **Difficulties**
- remote <u>search engines</u> vary in numerous ways
- <u>user population</u> vary, too