

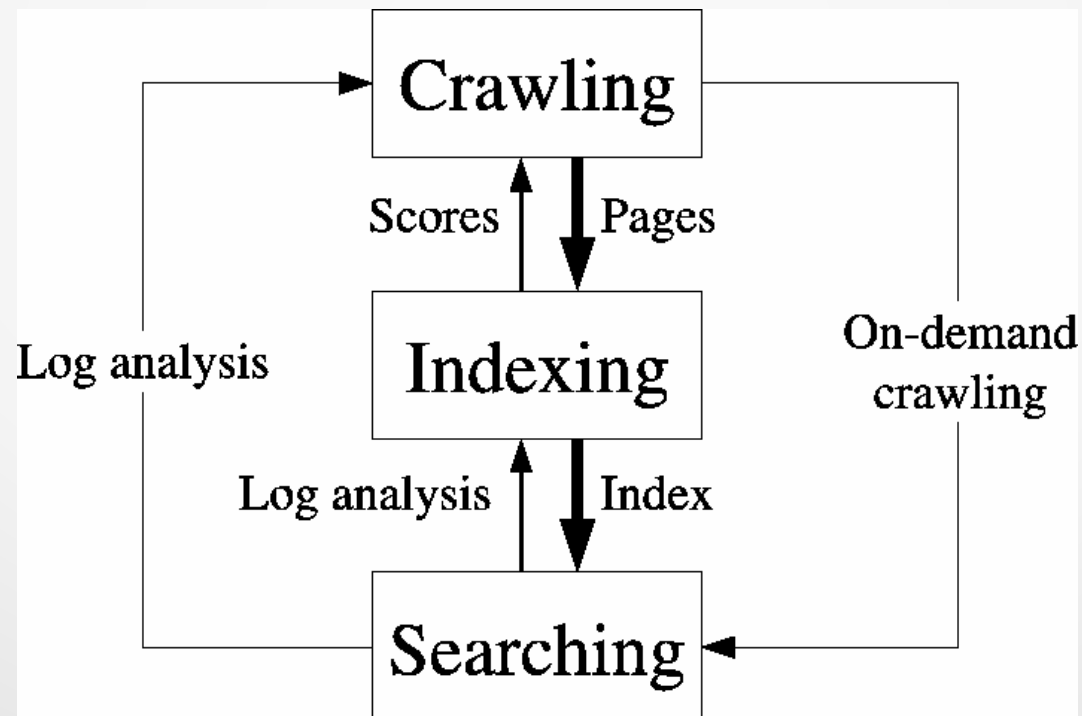
Uma Heurística para Focused Web Crawlers

Alberto Hideki Ueda
ueda@dcc.ufmg.br



O Problema

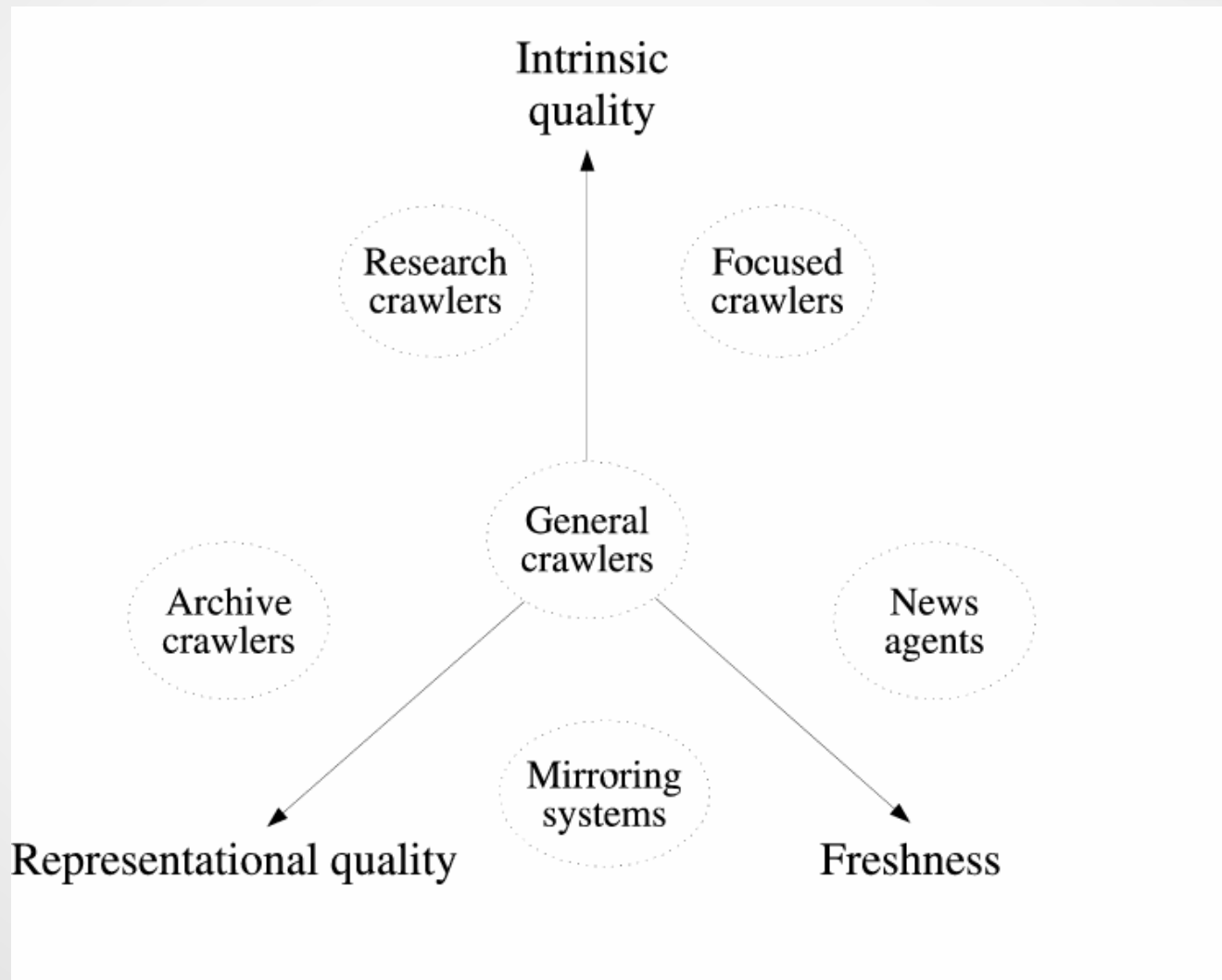
- Primeiro passo para a implementação de máquinas de busca modernas
- O que é um Coletor (*Crawler*)? [1]



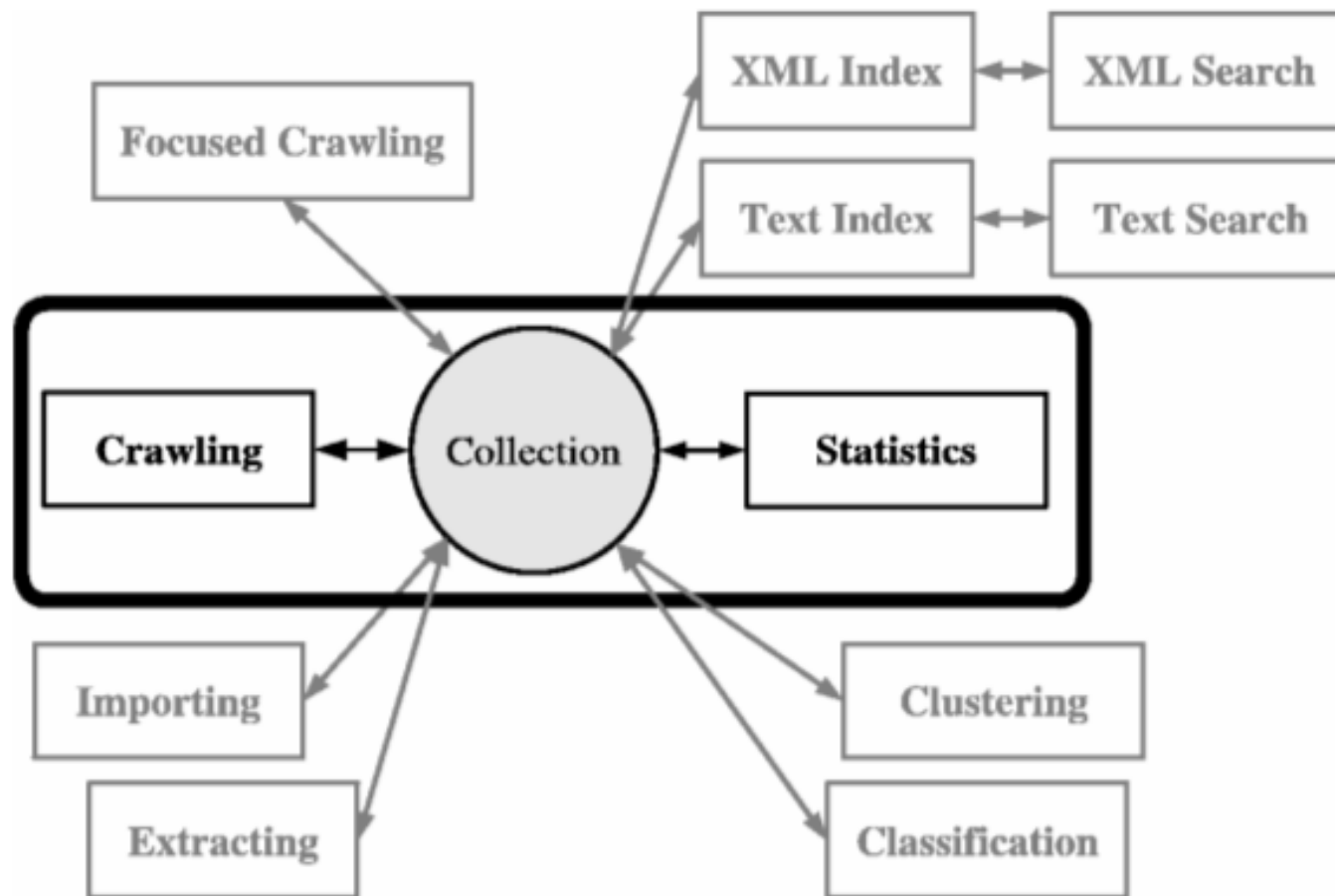
O Problema

- Primeiro passo para a implementação de máquinas de busca modernas
- O que é um Coletor (*Crawler*)? [1]
- Em 2005, foi demonstrado que o nível de cobertura das principais máquinas de buscas existentes está entre 58% e 76% da Web [7].
- Em 2004, tal custo foi estimado em US\$ 1,5 milhão para coletores de larga escala [6].
- **Inviabilidade da Solução Exata**

O Problema

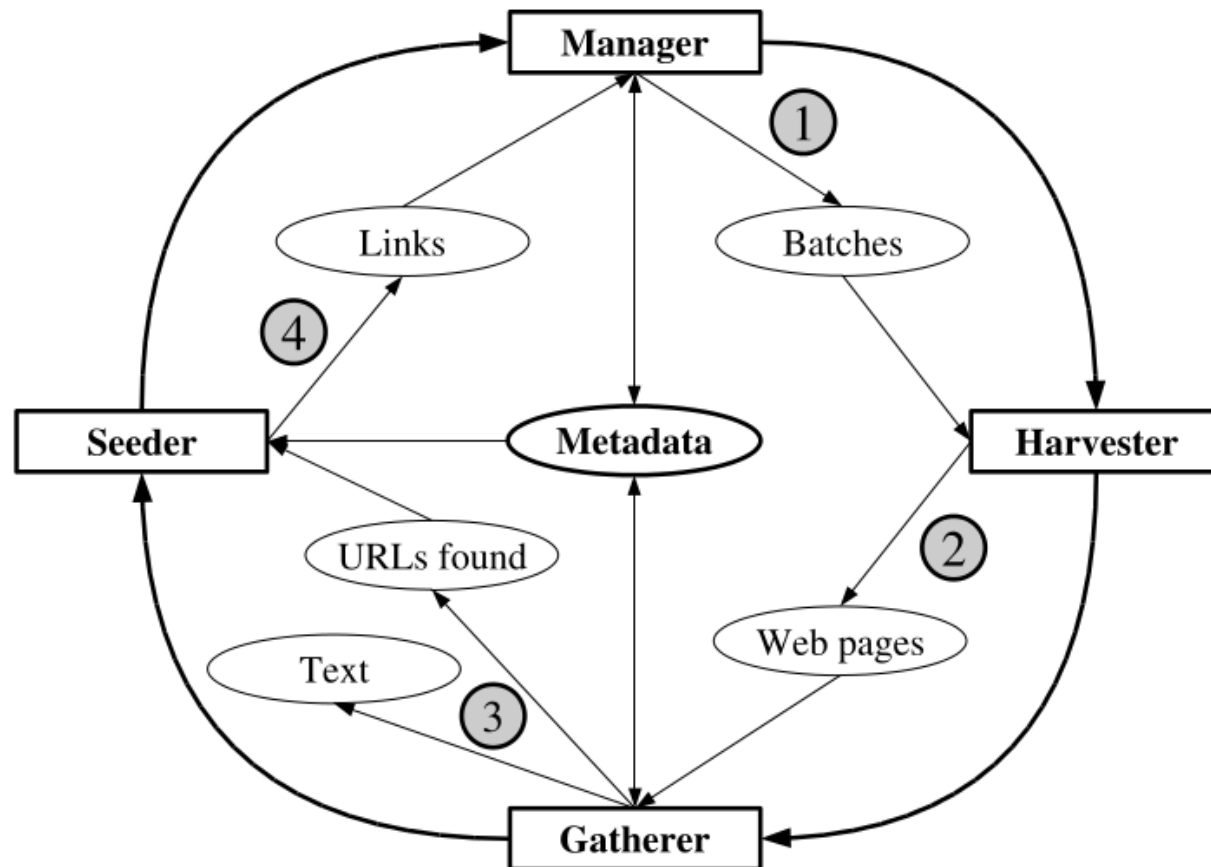


O Problema

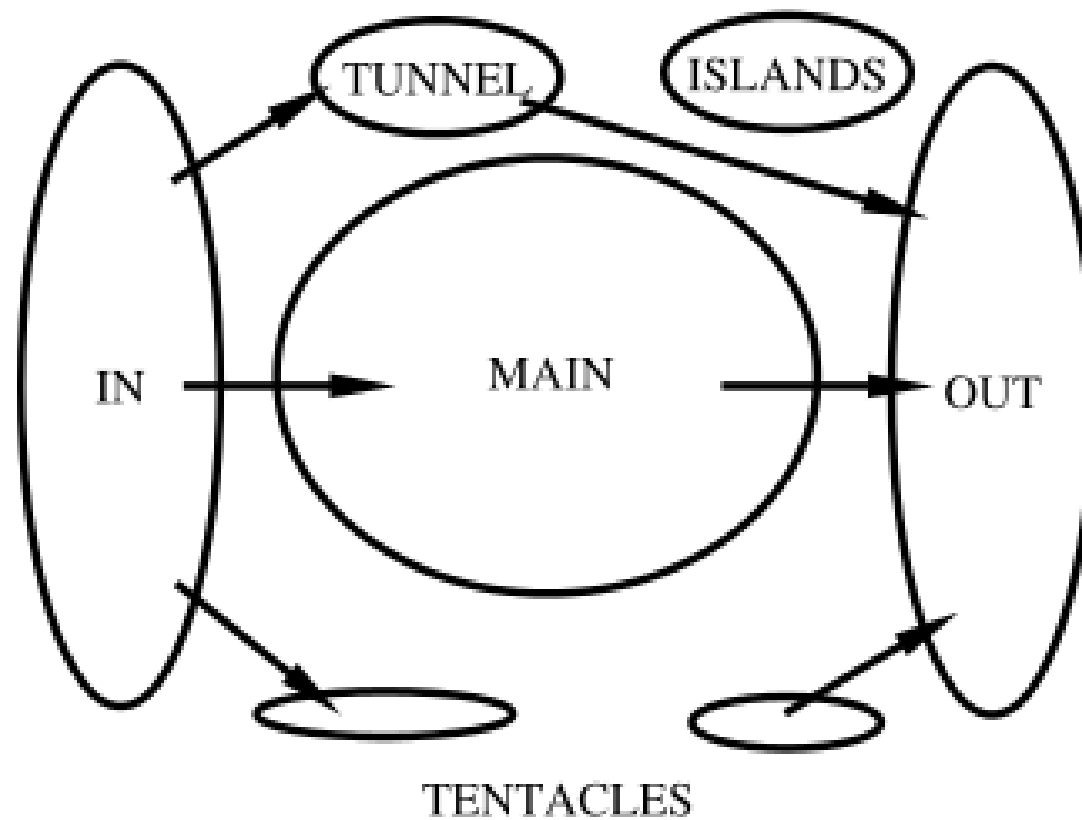


Baseline

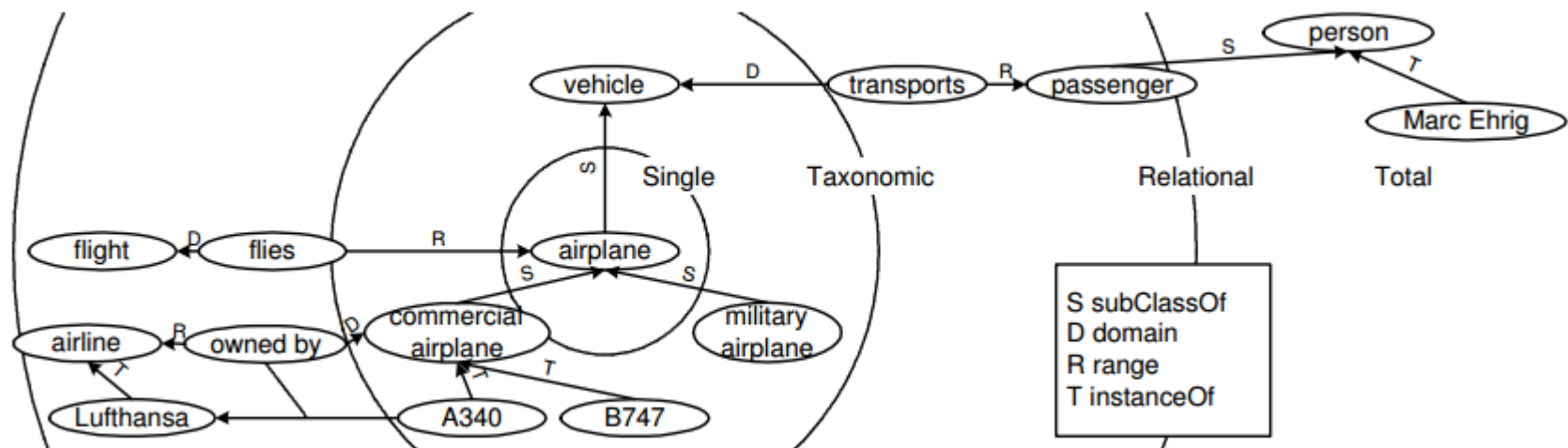
- **WIRE [4]**



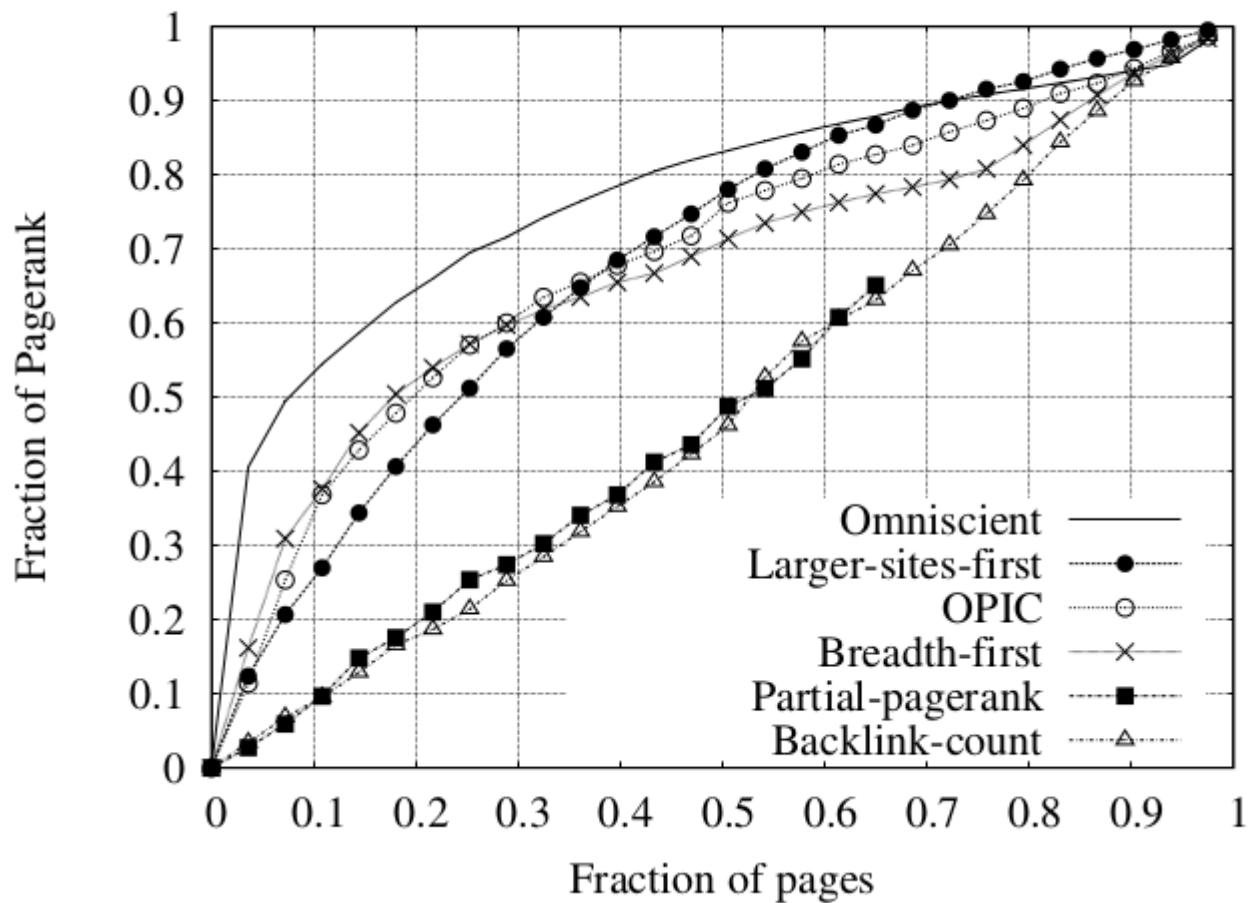
Modelagem em Grafo



Heurística Utilizada



Experimentos e Resultados



Referências

- Ricardo Baeza-Yates and Berthier Ribeiro-Neto. Modern Information Retrieval: The concepts and technology behind search. Pearson Education Limited, 2nd edition, 2011.
- Albert-Laszlo Barabasi. Linked the new science of networks, 2002.