# CONCEPTUALISATIONS OF WEB SEARCH
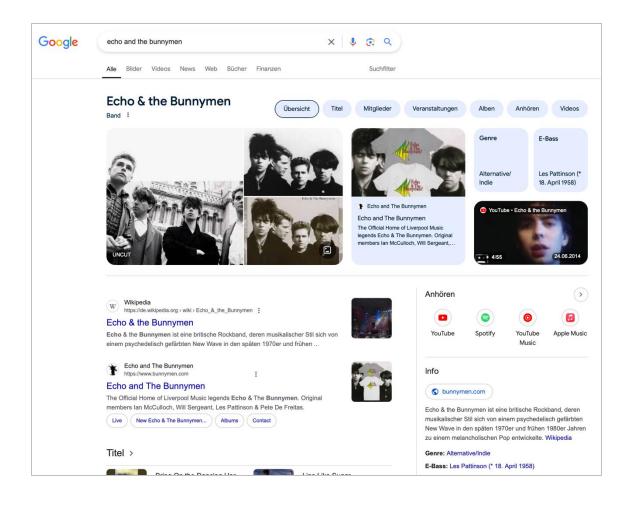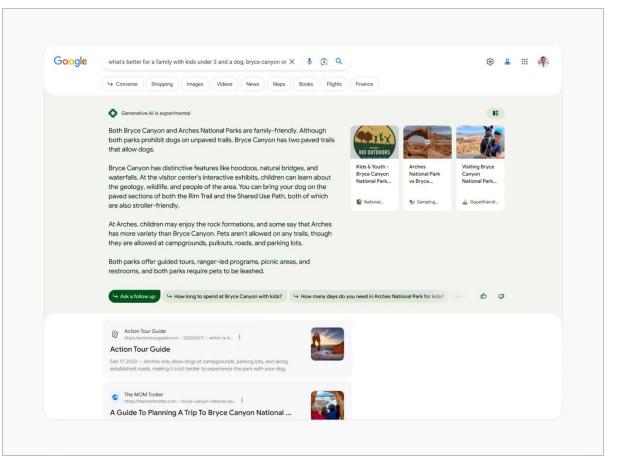
**Dirk Lewandowski**

Hamburg University of Applied Sciences, Hamburg, Germany

University of Duisburg-Essen, Duisburg, Germany

dirk.lewandowski@haw-hamburg.de

HAW HAMBURG

# SEARCH ENGINE RESULT PAGES

# TWO TYPES OF RESPONSES

**Documents: Information objects shown to help a user answer their question**

- Generated from an index of documents

**Answers: Text that directly answers a user's (anticipated) question**

- Key sentence(s) from documents (through information extraction)
- Facts (from specific data sources like Wikidata)
- AI-generated answers (longer, generated from a pre-processed index; LLMs)

HAW
HAMBURG

# TWO QUESTIONS

**1. What's the appropriate terminology?**

**2. What is a search engine?**

HAW
HAMBURG

# 1. WHAT'S THE APPROPRIATE TERMINOLOGY?

**Response**

- All result/answer types

**Result**

- Document result

**Answer**

- Text generated by the search engine

**Search Engine Result Page (SERP)**

- The page presenting the results: SERP

It makes sense to distinguish between results and answers because they are different and lead to different problems with evaluation, judgment, and information literacy. We see a mix of results and answers in search systems, and this will continue.

HAW
HAMBURG

# 2. WHAT IS A SEARCH ENGINE?

**Definition** (from "Understanding Search Engines")

"A search engine (also: Web search engine; universal search engine) is a computer system that captures distributed content from the World Wide Web via crawling and makes it searchable through a user interface, listing the results in a presentation ordered according to relevance assumed by the system."

**Elements not fulfilled in chat-based systems generating answers**

- Distributed content from the web? Crawling?
- "Presentation ordered according to relevance assumed by the system" – does a SE have to show a ranked list?

**Can LLM-generated answers be considered "search results?**

- No documents / Answers not directly generated from documents (like direct answers)
- As soon as sources are shown, very similar to "traditional" search results (the answer is just a summary of documents)
- Retrieval-Augmented Generation (RAG) uses a document collection (index) *plus* an LLM to generate the answers

**No problem with "mixed systems" (search engines extended through LLM-generated answers) but with purely chat-based systems**

- Information seeking / question answering is only one use case among many
- Is it realistic to have purely chat-based systems vs. "traditional" search engines?

HAW
HAMBURG