# The Impact of Web Search Result Quality on Decision-Making

September 12, 2024

Jan Heinrich
Merker

Lena
Merker

Alexander
Bondarenko

Friedrich-Schiller-Universität Jena
Martin-Luther-Universität Halle-Wittenberg

https://webis.de

# The Impact of Web Search Result Quality on Decision-Making

## Comparative questions and decision-making



Sources: https://smartpastamaker.com/pizza-vs-pasta-which-is-healthier/, https://petsoid.com/cats-vs-dogs/

What is healthier, pizza or pasta?

# The Impact of Web Search Result Quality on Decision-Making

## Comparative questions and decision-making



Sources: https://smartpastamaker.com/pizza-vs-pasta-which-is-healthier/, https://petsoid.com/cats-vs-dogs/

Should I adopt a dog or a cat?

# The Impact of Web Search Result Quality on Decision-Making

## Search engines and decision-making



**More reasons pro dog**



→ **How do you make a decision?**
... unless it's an obvious choice 😊

Comparative argument retrieval

- ❏ Goal: Retrieve relevant, high-quality arguments

- ❏ Comparative questions used as search queries
  [Bondarenko et al., WSDM'20]

- ❏ Examples: args.me or ArgumenText
  [Wachsmuth et al., EMNLP'17; Stab et al., NAACL-HLT'18]

- → Yet, many use "normal" search engines, like Google

  - – Known to be biased [Azzopardi, CHIIR'21]

  - – Impact on decision-making unclear!

# The Impact of Web Search Result Quality on Decision-Making

## Search engines and decision-making



**More reasons pro dog**



➔ How do you make a decision?

... unless it's an obvious choice 😊

## Comparative argument retrieval

❑ **Goal:** Retrieve relevant, high-quality arguments

❑ Comparative questions used as search queries
[Bondarenko et al., WSDM'20]

❑ Examples: args.me or ArgumenText
[Wachsmuth et al., EMNLP'17; Stab et al., NAACL-HLT'18]

➔ Yet, many use "normal" search engines, like Google

   – Known to be biased [Azzopardi, CHIIR'21]

   – Impact on decision-making unclear!

# The Impact of Web Search Result Quality on Decision-Making
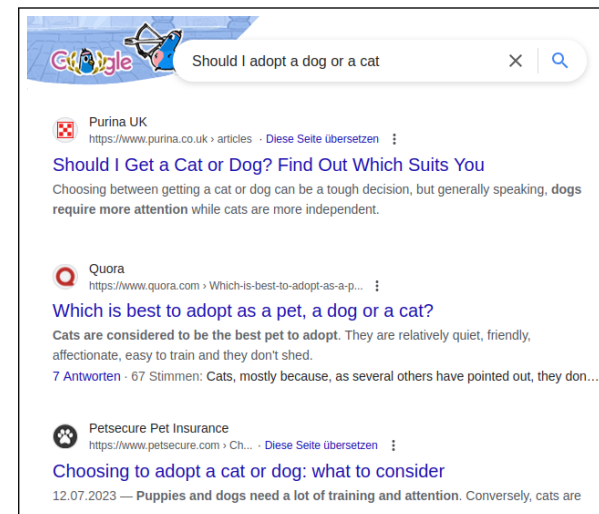
## Search engines and decision-making



More reasons pro dog



➔ How do you make a decision?
... unless it's an obvious choice 😊

## Comparative argument retrieval

❑ **Goal:** Retrieve relevant, high-quality arguments

❑ Comparative questions used as search queries
[Bondarenko et al., WSDM'20]

❑ Examples: args.me or ArgumenText
[Wachsmuth et al., EMNLP'17; Stab et al., NAACL-HLT'18]

➔ Yet, many use "normal" search engines, like Google

– Known to be biased [Azzopardi, CHIIR'21]

– Impact on decision-making unclear!

# The Impact of Web Search Result Quality on Decision-Making
## Hypotheses

1. Subjective comparisons lead to less confident decisions than factual.
   Intuition: Factual comparative topics often "better" answered by search engines than subjective comparative topics  [Bondarenko et al., WSDM'20]

2. Low-quality results lead to less confident decisions than high-quality results.
   Intuition: Desire to make best decision based on available information  [Peterson, '17]

3. The higher a document's quality, the more likely it influences the decision.
   Intuition: Same as for 2

4. More confident users are less influenced by low-quality documents.
   Intuition: Confident users rely more on own knowledge than on ad hoc information  [Peterson, '17]

5. Documents that take a stance have a higher impact on the decision.
   Intuition: Relevant documents often expected to take a stance  [Bondarenko et al., WSDM'22]

# The Impact of Web Search Result Quality on Decision-Making

## Hypotheses (and spoilers)

1. ~~Subjective comparisons lead to less confident decisions than factual.~~

   Intuition: Factual comparative topics often "better" answered by search engines than subjective comparative topics  [Bondarenko et al., WSDM'20]

2. ~~Low-quality results lead to less confident decisions than high-quality results.~~

   Intuition: Desire to make best decision based on available information  [Peterson, '17]

3. The higher a document's quality, the more likely it influences the decision.

   Intuition: Same as for 2

4. More confident users are less influenced by low-quality documents.

   Intuition: Confident users rely more on own knowledge than on ad hoc information  [Peterson, '17]

5. Documents that take a stance have a higher impact on the decision.

   Intuition: Relevant documents often expected to take a stance  [Bondarenko et al., WSDM'22]

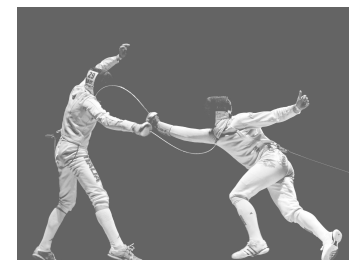# The Impact of Web Search Result Quality on Decision-Making
## Methodology

- ❑ Develop document quality criteria for comparative topics

- ❑ Assess quality, relevance, and stance of top-4 Google results for 30 topics

- ❑ Conduct user study on the decision-making

    - – Decision and confidence before/after seeing results
    - – Influence of retrieved documents

Data

Touché shared tasks

- ❑ Comparative argument retrieval task in 2020–2022
  100 topics comparing two or more options (e.g., dog vs. cat)

- → 30 topics used for quality assessment
  (comparing 2 options, easy to understand)

Google search engine

- ❑ Most popular search engine in Europe

- → Top-4 results used (after excl. ads / media-only results)

Source: https://touche.webis.de/

Source: https://unsplash.com/

# The Impact of Web Search Result Quality on Decision-Making

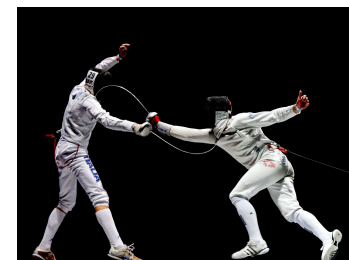## Methodology

- ❑ Develop document quality criteria for comparative topics

- ❑ Assess quality, relevance, and stance of top-4 Google results for 30 topics

- ❑ Conduct user study on the decision-making

  - – Decision and confidence before/after seeing results
  - – Influence of retrieved documents

## Data

Touché shared tasks

- ❑ Comparative argument retrieval task in 2020–2022
  100 topics comparing two or more options (e.g., dog vs. cat)

- → 30 topics used for quality assessment
  (comparing 2 options, easy to understand)

Source: https://touche.webis.de/

Google search engine

- ❑ Most popular search engine in Europe

- → Top-4 results used (after excl. ads / media-only results)

Source: https://unsplash.com/

# The Impact of Web Search Result Quality on Decision-Making
## Quality criteria

*Quality* [Lewandowski et al., '08]

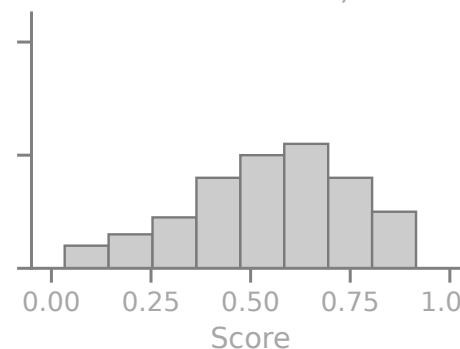| | | | | *Other criteria* | |
|---|---|---|---|---|---|
| **Content** | **Usability** | **Credibility** | **Up-to-dateness** | Relevance | Stance |
| Completeness, scope, language | Media types, structure | Source, author, truthfulness, verifiability | Date, updates | Topical relevance | Referral, emphasis, direction, magnitude |

- ❏ Based on prior quality assessment frameworks: WebQual, 2QCV3Q, AIMQ, Touché

- ❏ Relevance and quality also assessed for comparison purposes

## Quality assessment

- ❏ 120 documents assessed (Google's top-4 of 30 topics)

- ❏ 10 volunteer assessors (media / computer science stud.)

- ❏ Agreement measured based on randomly selected topic (Fleiss' $\kappa$; 6 aspects with insufficient agreement excluded)

- ❏ Calculate aggregated quality score per document and topic

Quality
(mean: 0.55, std: 0.21
median: 0.57)



Score

# The Impact of Web Search Result Quality on Decision-Making
## Quality criteria

*Quality* [Lewandowski et al., '08]                                    *Other criteria*

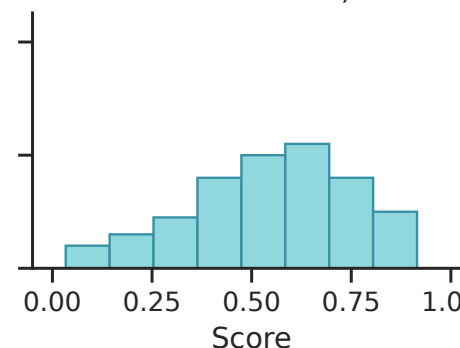| 📄 | ♿ | 🛡 | 📅 | 🧩 | ⚖ |
|---|---|---|---|---|---|
| **Content** | **Usability** | **Credibility** | **Up-to-dateness** | **Relevance** | **Stance** |
| Completeness, scope, language | Media types, structure | Source, author, truthfulness, verifiability | Date, updates | Topical relevance | Referral, emphasis, direction, magnitude |

- ❑ Based on prior quality assessment frameworks: WebQual, 2QCV3Q, AIMQ, Touché

- ❑ Relevance and quality also assessed for comparison purposes

### Quality assessment

- ❑ 120 documents assessed (Google's top-4 of 30 topics)

- ❑ 10 volunteer assessors (media / computer science stud.)

- ❑ Agreement measured based on randomly selected topic (Fleiss' $\kappa$; 6 aspects with insufficient agreement excluded)

- ❑ Calculate aggregated quality score per document and topic

Quality
(mean: 0.55, std: 0.21
median: 0.57)

Score

# The Impact of Web Search Result Quality on Decision-Making

## Quality criteria

*Quality* [Lewandowski et al., '08]                                            *Other criteria*

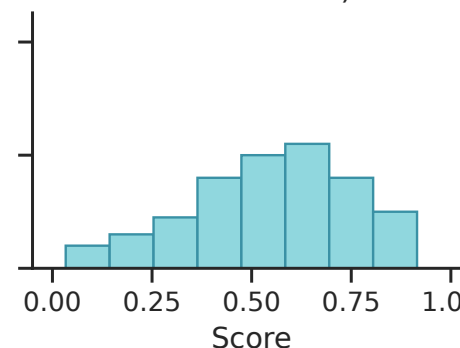| 📄 | ♿ | 🛡 | 📅 | 🧩 | ⚖ |
|---|---|---|---|---|---|
| **Content** | **Usability** | **Credibility** | **Up-to-dateness** | **Relevance** | **Stance** |
| Completeness, scope, language | Media types, structure | Source, author, truthfulness, verifiability | Date, updates | Topical relevance | Referral, emphasis, direction, magnitude |

❑ Based on prior quality assessment frameworks: WebQual, 2QCV3Q, AIMQ, Touché

❑ Relevance and quality also assessed for comparison purposes

## Quality assessment

❑ 120 documents assessed (Google's top-4 of 30 topics)

❑ 10 volunteer assessors (media / computer science stud.)

❑ Agreement measured based on randomly selected topic (Fleiss' $\kappa$; 6 aspects with insufficient agreement excluded)

❑ Calculate aggregated quality score per document and topic

Quality
(mean: 0.55, std: 0.21
median: 0.57)

Score

## Quality criteria

*Quality* [Lewandowski et al., '08]                                      *Other criteria*

| Content | Usability | Credibility | Up-to-dateness | Relevance | Stance |
|---|---|---|---|---|---|
| Completeness, scope, language | Media types, structure | Source, author, truthfulness, verifiability | Date, updates | Topical relevance | Referral, emphasis, direction / magnitude |

- ❑ Based on prior quality assessment frameworks: WebQual, 2QCV3Q, AIMQ, Touché

- ❑ Relevance and quality also assessed for comparison purposes

## Quality assessment

- ❑ 120 documents assessed (Google's top-4 of 30 topics)
- ❑ 10 volunteer assessors (media / computer science stud.)
- ❑ Agreement measured based on randomly selected topic (Fleiss' $\kappa$; 6 aspects with insufficient agreement excluded)
- ❑ Calculate aggregated quality score per document and topic

Quality
(mean: 0.55, std: 0.21
median: 0.57)



Score

# The Impact of Web Search Result Quality on Decision-Making
## User study

- ❏ Select 8 topics and screenshot top-4 results
  - – Exclude topics with missing quality judgments
  - – Cover wide range of topic-wise avg. quality

- ❏ Survey layout:
  - – Introduction and topic description
  - – Prior knowledge assessment
  - – Decision/confidence before seeing results
  - – Screenshots of documents
  - – Decision/confidence after seeing results
  - – Self-assessment of decision-making process (6 statements)

- ❏ 442 volunteer participants (German univ. students)

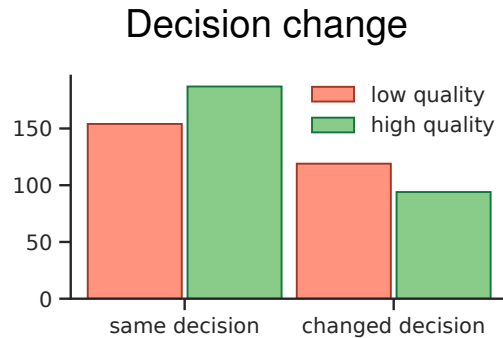- ❏ 554 study responses (1–8 topics per participant)
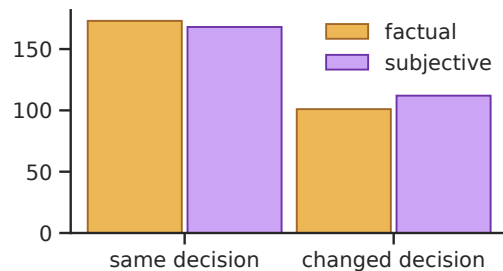


Survey view



Document screenshot

# The Impact of Web Search Result Quality on Decision-Making
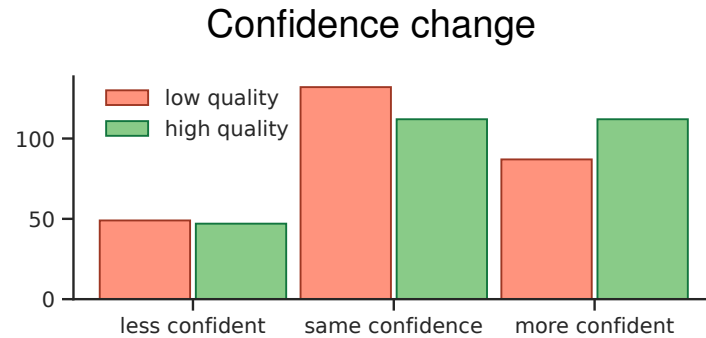
## Results: Decision and confidence change
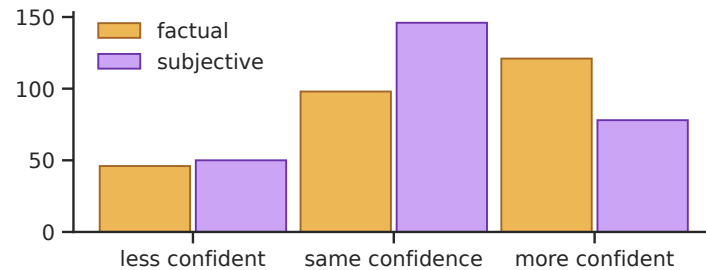


Decision change

$$\chi^2(1) = 5.59, p = 0.018$$

$$\chi^2(1) = 0.45, p = 0.502$$

Confidence change
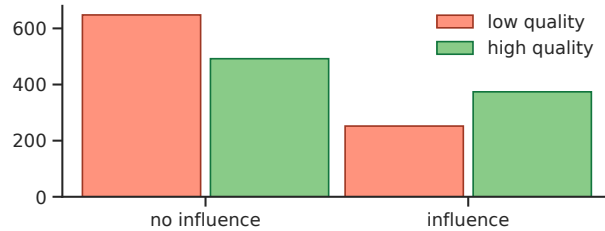
$$\chi^2(2) = 4.81, p = 0.090$$

$$\chi^2(2) = 18.76, p < 0.001$$

- ❑ Majority did not change decision, but 37% gained confidence

- ❑ Decision changed more often with overall low-quality results,
  but more confident with high-quality results

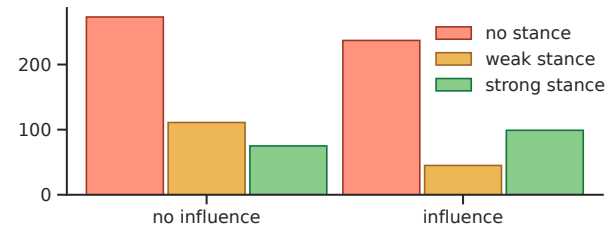- ❑ Decision confidence significantly increased more for factual topics

# The Impact of Web Search Result Quality on Decision-Making
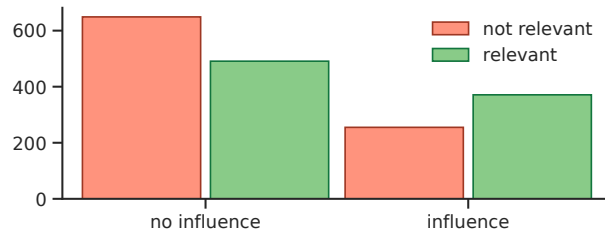
## Results: Influence of documents
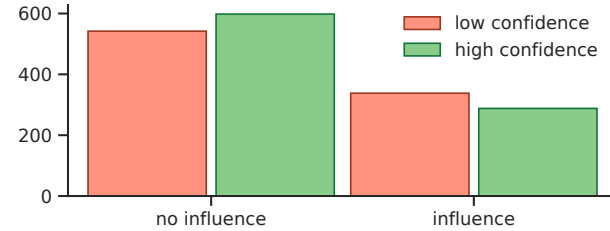
**Influenced decision-making**



$$\chi^2(4) = 44.49, p < 0.001$$

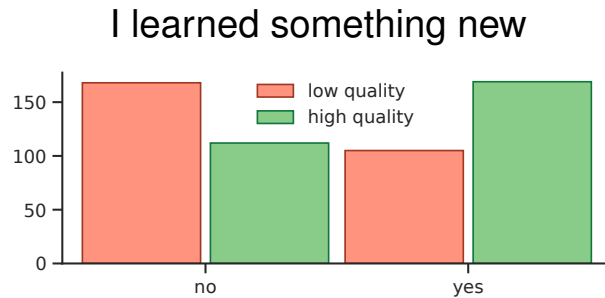$$\chi^2(2) = 26.76, p < 0.001$$

$$\chi^2(1) = 41.77, p < 0.001$$
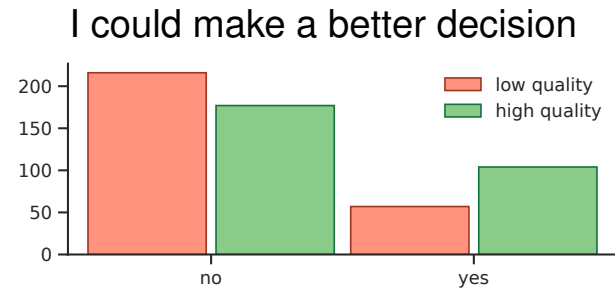
$$\chi^2(1) = 4.51, p = 0.034$$

- ❑ More likely to influence decision-making:
  - – high-quality
  - – relevance
  - – strong stance
- ❑ Confident users less likely to be influenced by low-quality documents

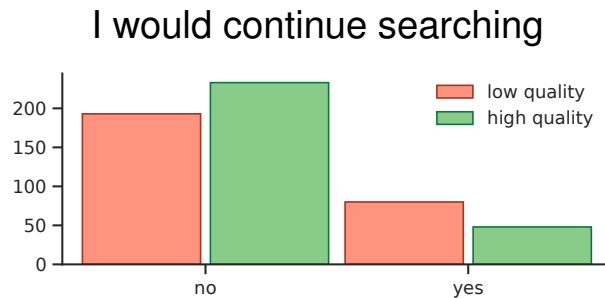# The Impact of Web Search Result Quality on Decision-Making

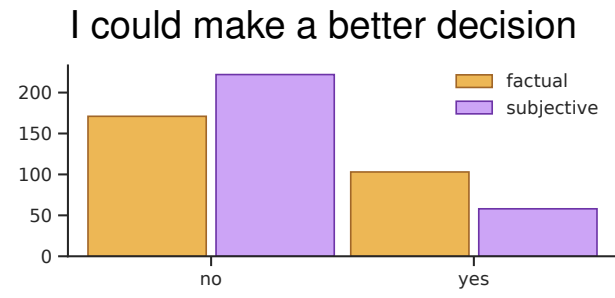## Results: Self-assessment of decision-making process



**I learned something new**

$$\chi^2(1) = 25.18, p < 0.001$$

**I could make a better decision**

$$\chi^2(1) = 18.32, p < 0.001$$

**I would continue searching**

$$\chi^2(1) = 10.96, p < 0.001$$

**I could make a better decision**

$$\chi^2(1) = 16.71, p < 0.001$$

- ❑ Participants learned sth. new (50%), but would still often continue search (25%)

- ❑ High-quality results more helpful, less likely to continue searching

- ❑ Better decisions with factual topics and high-quality results

## Results: Hypotheses

1. ~~Subjective comparisons lead to less confident decisions than factual.~~

   → Reject: Better decision for factual topics but not significantly more confident.

2. ~~Low-quality results lead to less confident decisions than high-quality results.~~

   → Reject: Slight increase in confidence with high-quality topics, better decision with high-quality topics; but overall not significant.

3. The higher a documents's quality, the more likely it influences the decision.

   → Accept: Low-quality documents influence decisions significantly less often than high-quality documents; position bias ruled out.

4. More confident users are less influenced by low-quality documents.

   → Accept: Confident users significantly less likely to be influenced by low-quality documents.

5. Documents that take a stance have a higher impact on the decision.

   → Accept: Docs. with strong stance influenced decision more often than with weak stance.

# The Impact of Web Search Result Quality on Decision-Making
## Summary

- First step for quality assessment of comparative queries

- Quality has significant impact on decision-making process

- Potential ranking factors: quality, stance (especially for subjective topics)

- Limitations: only German student participants, single search engine

- Future work: larger study (e.g., more participants / topics / search engines)

## Code and data

⌕ github.com/webis-de/CLEF-24

📄 doi.org/978-3-031-71736-9_5

# The Impact of Web Search Result Quality on Decision-Making
## Summary

- First step for quality assessment of comparative queries

- Quality has significant impact on decision-making process

- Potential ranking factors: quality, stance (especially for subjective topics)

- Limitations: only German student participants, single search engine

- Future work: larger study (e.g., more participants / topics / search engines)

## Code and data

- github.com/webis-de/CLEF-24

- doi.org/978-3-031-71736-9_5

*Thank you & merci!*