# Contents

## 0.1 Database Indexing and File Organization Review

**Short Answer Questions:**

1. **Explain how data is organized on a hard disk drive in terms of tracks, sectors, and blocks.**
2. **What is the main idea behind the fixed-length record storage approach, and what modification is made to avoid extra block accesses?**
3. **Describe two different approaches to handling record deletion in a file.**
4. **Differentiate between heap and sequential file organization in terms of record placement.**
5. **Explain why sequential file organization might be inefficient for queries involving the natural join of two tables.**
6. **What is the purpose of an index in a database, and what are the trade-offs of using one?**
7. **Differentiate between a dense index and a sparse index.**
8. **Why are secondary indexes always dense?**

9. **What is a multilevel index, and why is it used?**
10. **Describe the benefits and drawbacks of using a B+ tree index.**

**Answer Key:**

1. Data on a hard disk is organized into concentric circles called **tracks**. Each track is divided into **sectors**, which are the smallest units of data that can be read or written. **Blocks**, consisting of a contiguous sequence of sectors, are used to transfer data between the disk and main memory.
2. Fixed-length record storage aims to store record $i$ starting at byte $n(i$ -1), where $n$ is the fixed record size. However, this can lead to records spanning multiple blocks. To avoid this, the approach is modified to ensure each record fits within a single block.
3. Two approaches to handling record deletion are: **(1) Record shifting:** Move subsequent records to fill the gap created by the deleted record. **(2) Linked lists:** Instead of moving records, link free spaces using pointers, starting with a header pointing to the first free record and each free record pointing to the next.
4. **Heap file organization** allows records to be placed anywhere in the file where there is space, while **sequential file organization** stores records in a specific order based on a search key.
5. Sequential file organization requires traversing the entire file in the specified order. When performing a natural join, this means potentially reading both tables multiple times to find matching records, leading to inefficiency.
6. An **index** in a database is a data structure that improves the speed of data retrieval. The trade-off is that indexes require additional storage space and need to be updated whenever the data they index is modified.
7. A **dense index** contains an entry for every search key value in the data file, while a **sparse index** contains entries for only a subset of the search key values.
8. Secondary indexes are always dense because they don't determine the physical order of the data file. Therefore, to locate any record based on the secondary key, an index entry is required for every possible value.
9. A **multilevel index** is essentially an index to another index. It is used when the primary index is too large to fit in memory, improving search performance by reducing the number of disk accesses.
10. **B+ tree indexes** offer fast search, insertion, and deletion operations due to their balanced tree structure. They eliminate the need for periodic reorganization. However, they can be less efficient than other structures for certain types of queries, like range queries on non-indexed attributes.

**Essay Questions:**

1. Discuss the advantages and disadvantages of storing multiple tables in a single file (file grouping) versus storing each table in a separate file. Consider factors such as query performance, storage efficiency, and complexity of data management.

2. Explain the differences between primary indexes and secondary indexes. Describe scenarios where each type of index would be most beneficial.
3. Detail the steps involved in inserting a new record into a table that has a B+ tree index. Explain how the B+ tree structure ensures efficient search operations even after multiple insertions.
4. Compare and contrast the use of fixed-length records and variable-length records in database file organization. Discuss the implications of each approach on storage space, performance, and complexity of implementation.
5. Explain the concept of prefix compression in the context of indexing strings. How does prefix compression help improve the efficiency of B+ tree indexes on string attributes?

**Glossary of Key Terms:**

- **Track:** A circular path on a disk platter where data is magnetically recorded.
- **Sector:** A subdivision of a track, representing the smallest unit of data that can be read or written from a disk.
- **Block:** A contiguous sequence of sectors, forming the basic unit of data transfer between disk and memory.
- **Heap file organization:** A file organization method where records are stored in no particular order.
- **Sequential file organization:** A file organization method where records are stored in sequence based on a search key.
- **Index:** A data structure used to speed up data retrieval by providing a quick lookup mechanism for finding records based on specific attributes.
- **Dense index:** An index that contains an entry for every search key value in the data file.
- **Sparse index:** An index that contains entries for only a subset of search key values, typically used when the data file is sequentially ordered on the indexed attribute.
- **Multilevel index:** An index to another index, used to improve search performance for large indexes by reducing disk accesses.
- **B+ tree index:** A balanced tree data structure commonly used for indexing in databases, offering efficient search, insertion, and deletion operations.
- **Prefix Compression:** A technique used to shorten index keys by storing only a prefix of the key, reducing space requirements and improving efficiency.

# 1 Relational Algebra Study Guide

## 1.1 Short-Answer Questions

**Instructions:** Answer the following questions in 2-3 sentences each.

1. What is a relational algebra operator?
2. How can a SQL query be translated into an expression in relational algebra?
3. What are the two types of relational algebra operators? Give an example of each.
4. Explain the concept of "cost" in the context of relational algebra operations. How is it typically measured?
5. What does the notation "R = (A1::T1, ..., An::Tn)" represent in the context of relational schemas?
6. Define the generalized projection operation in relational algebra. Provide an example.
7. What is the purpose of the selection operation in relational algebra? How is it denoted?
8. Describe the linear search algorithm for the selection operation and state its estimated cost.
9. What is the purpose of a selectivity factor in estimating the size of results from relational algebra operations?
10. Explain the concept of a "natural join" in relational algebra. Provide an example.

## 1.2 Short-Answer Key

1. A relational algebra operator is a function that takes one or more relations (tables) as input and produces a new relation as output. These operators allow for the manipulation and retrieval of data within a relational database.
2. A SQL query can be translated into an expression in relational algebra by breaking down the query into smaller components, each representing a specific relational algebra operation. These operations are then combined to form an expression that represents the original query.
3. The two types of relational algebra operators are logical operators and physical operators. A logical operator defines the operation to be performed (e.g., selection, projection). A physical operator defines how the logical operation is implemented (e.g., linear search, index scan).
4. "Cost" in relational algebra refers to the resources used when performing an operation. It is typically measured in terms of the number of disk block transfers, as disk I/O is often the most time-consuming aspect. This measurement helps compare the efficiency of different algorithms.
5. This notation represents a relational schema for a relation named "R". It consists of attributes A1 to An, each with its corresponding data type T1 to Tn. For example, "Name::string" indicates an attribute "Name" with the data type "string".
6. Generalized projection extends the standard projection by allowing computations on attributes. It is denoted as $\Pi$ f1,...,fn(R), where f1 to fn are functions applied to

each tuple of relation R. For example, to get the annual salaries of professors from a table "Professor(Name, Salary)", we can use Π Name, Salary*12(Professor).

7. The selection operation retrieves tuples from a relation that satisfy a given condition (predicate). It is denoted as $\sigma$ P(R), where P is the predicate applied to each tuple of relation R. For example, to select professors with a salary greater than 60000, we would use $\sigma$ Salary > 60000(Professor).

8. The linear search algorithm scans each block of the relation and checks if each tuple satisfies the selection condition. Its estimated cost is br block transfers + 1 block access, where br is the number of blocks containing tuples in relation R.

9. A selectivity factor estimates the fraction of tuples in a relation that will satisfy a given condition. It is used to estimate the size of intermediate and final results after operations like selection or join. This helps in optimizing query execution plans.

10. A natural join combines tuples from two relations based on the equality of all attributes with the same name in both relations. For example, if we have relations R(A, B, C) and S(B, C, D), the natural join R ⋈ S will return tuples where R.B = S.B and R.C = S.C.

## 1.3 Essay Questions

**Instructions:** Answer the following questions in an essay format.

1. Discuss the differences between logical and physical operators in relational algebra. Explain how these concepts are related to the process of query optimization.

2. Explain the concept of a B+ tree index and describe how it can be used to optimize the selection operation in relational algebra. Compare the performance of using a B+ tree index to that of a linear scan for different selection predicates.

3. Compare and contrast the nested-loop join and merge-join algorithms for implementing the join operation in relational algebra. Discuss the advantages and disadvantages of each algorithm and the scenarios in which one might be preferred over the other.

4. Explain how the selectivity factor is used in estimating the size of results for relational algebra operations. Discuss the assumptions involved in calculating selectivity factors and their impact on the accuracy of the estimations.

5. Discuss the different strategies for implementing the duplicate elimination operation in relational algebra. Compare the costs and benefits of each approach and explain how they relate to the concept of external sorting.

## 1.4 Glossary of Key Terms

**TermDefinition**RelationA table with rows (tuples) and columns (attributes), representing a set of related data.SchemaDefines the structure of a relation, specifying the

name of each attribute and its data type.TupleA row in a relation, representing a single instance of the data described by the schema.AttributeA named column in a relation, representing a specific characteristic of the data.DomainThe set of possible values that an attribute can take.Relational AlgebraA formal query language that provides a set of operators for manipulating relations.Logical OperatorA relational algebra operator that defines the operation to be performed, such as selection, projection, or join.Physical OperatorA concrete implementation of a logical operator, specifying the algorithms and data structures used to perform the operation.Query OptimizationThe process of finding the most efficient way to execute a given query, typically by considering different execution plans and choosing the one with the lowest estimated cost.Disk Block TransferThe movement of a block of data between disk and memory. The number of block transfers is often used as a measure of the cost of a relational algebra operation.B+ Tree IndexA tree-like data structure that allows for efficient searching, insertion, and deletion of data. B+ tree indexes can be used to optimize relational algebra operations by providing fast access to tuples based on the values of indexed attributes.Selectivity FactorA measure of the fraction of tuples in a relation that satisfy a given condition. It is used to estimate the size of results from relational algebra operations.Nested-Loop JoinA join algorithm that iterates over each tuple in the outer relation and for each such tuple, scans the entire inner relation to find matching tuples.Merge-JoinA join algorithm that sorts both relations on the join attribute(s) and then scans the sorted relations once to find matching tuples.Duplicate EliminationThe process of removing duplicate tuples from a relation.External SortingA sorting algorithm that is designed to handle data sets that are too large to fit in memory. External sorting algorithms typically involve dividing the data into smaller chunks, sorting the chunks in memory, and then merging the sorted chunks to produce the final sorted result.Concatenation (Union All)A relational algebra operation that combines the tuples from two relations without removing duplicates.IntersectionA relational algebra operation that returns only the tuples that are present in both input relations.Difference (Minus)A relational algebra operation that returns the tuples that are present in the first relation but not in the second relation.

## 2 Relational Algebra and Query Optimization Study Guide

### 2.1 Short-Answer Questions

1. **What is a query evaluation plan, and how does it relate to query optimization?**
2. **Explain the difference between materialization and pipelining in query evaluation.**
3. **What is the purpose of selectivity factor in query optimization, and how is it calculated for a selection operation?**

4. **Describe the concept of "pushing selections" in query optimization and explain its benefits.**
5. **Explain why the order of natural joins in a query can significantly impact performance.**
6. **What is a left-deep join tree, and how does it differ from other join trees?**
7. **Why are heuristics often used in query optimization instead of always striving for the absolute optimal plan?**
8. **Give an example of a heuristic rule used for join order selection and explain its rationale.**
9. **What are the advantages of using dynamic programming for join order optimization compared to a brute-force approach?**
10. **Briefly describe how hybrid approaches combine heuristics and cost-based optimization in query optimization.**

## 2.2 Short-Answer Key

1. A query evaluation plan outlines the specific steps and algorithms used to execute a query. Query optimization aims to find the most efficient plan by considering different logical equivalencies and physical operators.
2. Materialization stores intermediate results on disk as temporary tables, while pipelining directly passes results between operators without storing them. Materialization uses less memory but more disk space, while pipelining conserves disk space but requires more memory.
3. Selectivity factor estimates the fraction of tuples that satisfy a selection condition. For selection $\sigma$ (P, r), it's calculated as the number of tuples satisfying P divided by the total number of tuples in r.
4. "Pushing selections" applies selection operations as early as possible in the query tree. This reduces the size of intermediate relations, improving overall query performance.
5. Different join orders produce intermediate relations of varying sizes. Optimizing the join order minimizes the size of intermediate results, leading to faster query execution.
6. A left-deep join tree has the property that the right operand of each join is always a base relation, not the result of an intermediate join. This structure is often preferred for pipelined evaluation.
7. Finding the absolute optimal query plan can be computationally expensive, especially for complex queries. Heuristics offer good solutions quickly, even if they might not be the absolute best.
8. One heuristic is to join the smallest relations first. This minimizes the initial intermediate relation size, potentially leading to faster subsequent joins.
9. Dynamic programming breaks down the problem of finding the optimal join order into smaller subproblems and stores their solutions. This avoids redundant

computations and significantly reduces the search space compared to a brute-force approach.

10. Hybrid approaches apply heuristics for parts of the query and use cost-based optimization for others, balancing performance and optimization cost. For example, they might use heuristics for join order within specific subqueries while employing a cost-based approach for optimizing between subqueries.

## 2.3 Essay Questions

1. Discuss the trade-offs between materialization and pipelining in query evaluation. When is one approach preferred over the other? Consider factors such as memory usage, disk space, and query complexity.
2. Explain the concept of relational algebra equivalences and their significance in query optimization. Provide examples of at least three equivalence rules and illustrate how they can be used to transform a query into a more efficient form.
3. Compare and contrast heuristic-based and cost-based query optimization techniques. Discuss the advantages and disadvantages of each approach. Explain scenarios where one approach might be more suitable than the other.
4. Describe the dynamic programming approach to join order optimization. Explain the steps involved in building the dynamic programming table and how the optimal join order is determined. What are the time and space complexities of this approach?
5. Explain the concept of hybrid query optimization. Discuss how hybrid approaches combine heuristics and cost-based optimization. Provide examples of how different aspects of a query might be optimized using different techniques within a hybrid approach.

## 2.4 Glossary of Key Terms

TermDefinitionQuery Evaluation PlanA sequence of steps and algorithms used by a database management system to execute a query. It outlines the order of operations, access methods, and algorithms for each operator.Query OptimizationThe process of finding the most efficient way to execute a database query, considering factors like execution time, resource utilization, and overall cost.MaterializationAn evaluation strategy where intermediate results of a query are stored on disk as temporary tables. This approach reduces memory usage but may increase disk I/O.PipeliningAn evaluation strategy where the results of one operator are directly passed as input to the next operator without storing them as temporary relations. This approach minimizes disk I/O but requires more memory for buffering intermediate results.Selectivity FactorA statistical measure used in query optimization to estimate the fraction of tuples in a relation that will satisfy a given selection predicate."Pushing Selections"A query optimization technique that applies selections as early as possible in the query tree to reduce the size of intermediate relations and improve performance.Join OrderThe order in which relations are joined in a query

involving multiple tables. Choosing an efficient join order can significantly impact query performance.Left-Deep Join TreeA query tree representing joins where the right operand of each join is always a base relation, not an intermediate result. This structure is often preferred for pipelined evaluation.HeuristicA rule of thumb or guideline used in query optimization to make quick decisions and find good, but not necessarily optimal, execution plans. Heuristics are often used to reduce the search space for optimization.Cost-Based OptimizationA query optimization approach that assigns costs to different execution plans based on factors like disk I/O, CPU time, and communication costs. The optimizer chooses the plan with the lowest estimated cost.Dynamic ProgrammingA technique used in query optimization to find the optimal join order for a query by breaking down the problem into smaller, overlapping subproblems and storing their solutions to avoid redundant calculations.Hybrid Query OptimizationA combination of heuristic-based and cost-based optimization techniques. Hybrid approaches aim to balance optimization efficiency with the quality of the generated execution plans.

# 3 Information Retrieval Study Guide

## 3.1 Short Answer Questions

**Instructions:** Answer the following questions in 2-3 sentences each.

1. What is information retrieval (IR), and what kind of data does it typically deal with?
2. Explain three key differences between relational databases and IR systems.
3. Describe two different types of query languages used in IR systems.
4. What is the purpose of using proximity operators in IR queries? Provide an example.
5. How does Google Search handle natural language queries, and what are some operators it uses?
6. What are the main factors that determine the relevance of a document to a query?
7. Briefly describe the Boolean model of information retrieval. What are its limitations?
8. How does the vector space model represent documents and queries? How is relevance measured in this model?
9. What is TF-IDF, and how is it used in the vector space model? Explain the idea behind it.
10. Explain the purpose and process of stemming and synonym handling in the context of IR.

## 3.2 Short Answer Key

1. Information retrieval (IR) is the process of finding documents from a collection that are relevant to a user's query. IR systems typically deal with unstructured, natural language data.
2. Three key differences: a) Relational databases deal with structured data in tables, while IR systems handle unstructured data like text. b) Relational databases use structured query languages (SQL), while IR systems often employ keyword-based or natural language queries. c) Relational databases return exact matches, while IR systems rank results based on relevance.
3. Two types of query languages: a) Boolean queries, where users combine keywords with operators like AND, OR, and NOT. b) Natural language queries, where users phrase their information need in everyday language, like asking a question.
4. Proximity operators specify how close certain terms should be to each other in a document. For example, "science NEAR technology" would retrieve documents where these terms appear close together, suggesting a stronger relationship between the concepts.
5. Google Search uses a sophisticated algorithm to understand the intent and meaning behind natural language queries. It also utilizes operators like quotes for exact phrase matching, OR for alternative terms, "-" for exclusion, "site:" for limiting results to a specific website, and "intitle:" for searching within page titles.
6. Relevance is determined by factors such as: a) Term frequency: how often query terms appear in a document. b) Inverse document frequency: how rare a term is across the entire document collection (rarer terms are considered more informative). c) Links to a document: the number and quality of links pointing to a page can indicate its importance and authority.
7. The Boolean model represents documents as sets of terms and uses Boolean operators (AND, OR, NOT) to construct queries. Results are either a match or not, with no ranking. Limitations include difficulty in expressing complex information needs and no notion of partial matches or relevance ranking.
8. The vector space model represents both documents and queries as vectors in a multi-dimensional space, where each dimension corresponds to a term. Relevance is measured by the cosine similarity between the query vector and document vectors, with higher similarity indicating greater relevance.
9. TF-IDF stands for Term Frequency-Inverse Document Frequency. It's a weighting scheme used in the vector space model to assess the importance of a term within a document and across the entire collection. The idea is that terms appearing frequently in a document but rarely in others are more informative and discriminative.
10. Stemming reduces words to their root form (e.g., "running," "runs" become "run"), grouping similar words together. Synonym handling involves using a thesaurus or ontology to consider alternative terms that convey the same meaning, improving retrieval accuracy and recall.

### 3.3 Essay Questions

1. Discuss the advantages and disadvantages of using natural language queries compared to Boolean queries in information retrieval systems.
2. Compare and contrast the Boolean model and the vector space model of information retrieval. Discuss their strengths, weaknesses, and the types of applications where each might be more suitable.
3. Explain the concept of relevance in information retrieval. Discuss the various factors that contribute to document relevance and how they can be measured and utilized by IR systems.
4. Describe the structure and purpose of an inverted index in information retrieval. Explain how it is used to efficiently process queries and retrieve relevant documents.
5. Discuss the challenges and techniques associated with evaluating the performance of information retrieval systems. Describe commonly used metrics like precision, recall, and F-score, and explain how they can be interpreted and used to compare different systems or approaches.

### 3.4 Glossary of Key Terms

**Boolean Model:** An IR model that represents documents and queries as sets of terms, using Boolean operators (AND, OR, NOT) for retrieval.

**Cosine Similarity:** A measure of similarity between two vectors, often used in the vector space model to determine the relevance of a document to a query.

**Information Retrieval (IR):** The process of finding relevant information from a collection of documents in response to a user's query.

**Inverted Index:** A data structure used in IR systems that maps terms to the documents containing them, enabling efficient retrieval.

**Natural Language Query:** A query expressed in everyday language, allowing users to interact with IR systems in a more intuitive way.

**Precision:** A metric in IR that measures the proportion of retrieved documents that are relevant to the query.

**Proximity Operator:** A search operator used to specify how close certain terms should be to each other in a retrieved document (e.g., NEAR, WITHIN).

**Recall:** A metric in IR that measures the proportion of relevant documents that are successfully retrieved.

**Relevance:** The degree to which a document satisfies the information need expressed in a user's query.

**Stemming:** The process of reducing words to their root form to improve retrieval efficiency and effectiveness.

**Stop Words:** Common words (e.g., "the," "a," "is") that are often removed from documents and queries as they carry little informational value.

**Synonym Handling:** The use of a thesaurus or ontology to consider alternative terms with similar meanings during retrieval, enhancing recall.

**TF-IDF (Term Frequency-Inverse Document Frequency):** A weighting scheme used in IR to assess the importance of a term within a document and across the entire collection.

**Vector Space Model:** An IR model that represents documents and queries as vectors in a multi-dimensional space, using cosine similarity to measure relevance.

## 3.5 Web Information Retrieval Study Guide

### 3.5.1 Glossary of Key Terms

### 3.5.2 TermDefinitionWeb CrawlerA program that systematically browses the World Wide Web, typically for the purpose of web indexing.IndexingThe process of creating a data structure that allows for efficient searching of a collection of documents.Inverted IndexA data structure that maps words to the documents they appear in, facilitating fast keyword searches.TF-IDFTerm Frequency-Inverse Document Frequency: a numerical statistic reflecting how important a word is to a document in a collection. Used for information retrieval and text mining.PageRankA Google algorithm that assigns a numerical weighting to each element of a hyperlinked set of documents, such as the World Wide Web, with the purpose of "measuring" its relative importance within the set.Anchor TextThe visible, clickable text in a hyperlink.Hits AlgorithmAn algorithm that measures the importance of a webpage based on the number of other important pages that link to it.SnippetA short extract of text from a web page, displayed on a search engine results page (SERP) to give the user a preview of the content of the page.BreadcrumbA navigational aid that shows the user the path they have taken to reach their current location on a website.Knowledge PanelA box that appears on the right-hand side of the Google search results page, providing a summary of information about a particular topic, person, place, or thing.Short Answer Questions

**Instructions:** Answer the following questions in 2-3 sentences each.

1. What is the fundamental difference between a web search engine and a traditional search engine?

2. Explain how web crawlers utilize hyperlinks to discover new web pages.
3. Why is it crucial for web search engines to have robust indexing systems?
4. Describe the challenges posed by spam pages to web search engines and how they can be addressed.
5. Why did solely relying on TF-IDF prove insufficient for ranking web pages effectively?
6. Explain the concept of "site popularity" and how it contributes to search result ranking.
7. How does PageRank utilize a random walk model to calculate webpage importance?
8. What is the significance of anchor text in evaluating webpage relevance for specific topics?
9. How does Google utilize user search behavior to enhance the accuracy and relevance of its search results?
10. Describe the different types of information presented in a typical Google search result, beyond just a link to the webpage.

### 3.5.3 Short Answer Key

1. Traditional search engines operate on a predefined collection of documents, while web search engines constantly discover and index new pages on the ever-evolving World Wide Web.
2. Web crawlers start with a set of seed pages and follow the hyperlinks present on those pages to find new documents. This process continues recursively, expanding the collection of indexed pages.
3. Indexing systems organize and categorize web pages based on their content, enabling quick and efficient retrieval of relevant pages in response to user queries.
4. Spam pages manipulate search engine algorithms by stuffing keywords or creating artificial backlinks. Search engines combat this through algorithms that identify and penalize spam practices, and by incorporating user feedback.
5. The vastness of the web and the existence of spam pages meant that TF-IDF alone could not accurately reflect the true relevance or authority of a web page.
6. Site popularity refers to metrics like the number of visits a website receives or how many other websites link to it. It serves as a proxy for the quality and trustworthiness of a site, influencing its ranking in search results.
7. PageRank simulates a random web surfer who clicks links randomly. The probability of landing on a particular page reflects its PageRank, with pages receiving more links from high-ranking pages having a higher score.
8. Anchor text provides contextual clues about the content of the linked page. Search engines use this information to assess the page's relevance for queries containing those keywords.
9. Google tracks user clicks, dwell time, and other interactions with search results. This data helps refine ranking algorithms, promoting pages that users find genuinely helpful and demoting those that fail to engage users.

10. A Google search result typically includes the page title, URL, a snippet of relevant text, breadcrumbs, and sometimes additional information like the date, sitelinks, or rich snippets (e.g., images, reviews).

### 3.5.4 Essay Questions

1. Discuss the challenges involved in building and maintaining a large-scale web search engine, considering aspects like crawling, indexing, and handling the dynamic nature of the web.
2. Explain how the concept of "link analysis" is used to determine the importance and authority of web pages. Discuss different link analysis algorithms and their strengths and weaknesses.
3. Evaluate the effectiveness of combining traditional relevance measures like TF-IDF with link-based popularity metrics like PageRank for ranking web search results. Discuss the potential limitations of this approach.
4. Analyze the evolution of web search engines beyond simple keyword matching. Discuss how modern search engines strive to understand user intent and provide more sophisticated search results, including rich snippets, knowledge panels, and personalized recommendations.
5. Discuss the ethical implications of web search algorithms. Consider issues like filter bubbles, algorithmic bias, and the impact of search engine rankings on information access and societal discourse.