

搜索引擎实现技术方案

1. 功能描述

这经过精心设计的搜索引擎能够将用户输入的文本内容分解成关键词，并基于这些关键词实现智能搜索。搜索结果会根据相关性智能排序后呈现给用户。该搜索引擎的架构分为三个关键层次：初始数据层、倒排索引层和二级索引层。

初始数据层：

我们存储了数据源文件，总计 900+ 个文件，其中包含百万行句子。每行代表唯一的 ID，后跟 ID 对应的句子。句子内的单词以空格分隔，每行可以看作一个搜索链接。

倒排索引层：

与传统的按文件记录方式不同，这一层按独立词汇进行记录。每个词汇都记录了其出现在哪些文件中以及对这些文件的相关性信息。

二级索引层：在这一层，我们记录了每个词汇在倒排索引中的位置信息，以加速检索过程。

这个架构的优势在于其高效处理大量文本数据的能力，能够将用户查询分解为关键词，快速检索相关文本，并根据相关性进行排序，从而提供更智能和相关性更高的搜索结果。

2. 技术选型

A. 项目概况

设计和开发一个离线搜索引擎，该引擎需要处理庞大的原始数据集并实现高效的检索功能，同时生成便于搜索的多项要素。这个项目不仅需要深刻理解大数据技术，还需要熟练应用这些技术以应对复杂的数据处理和搜索需求。

B. 需求分析

- 原始数据的分割与预处理

处理庞大的原始数据集时，初步处理是不可或缺的。我们将大量原始数据按照指定的句子数量或文件大小进行分割，形成一个文件集合。这样做的目的是为了能够并行处理大量数据，提高后续工作的效率。

- 数据的智能排序和倒排索引构建

搜索引擎的核心功能是根据搜索查询中的具体单词返回相关结果。为了实现快速且精确的查询，我们需要对数据进行智能排序，以满足搜索引擎对单词敏感的特性。这就需要构建倒排索引，将数据按单词进行排序，并记录每个单词的相关信息，如单词频率、包含该单词的文件列表等。

- 分区策略的应用

- 倒排索引文件的大小通常接近原始数据的大小。如果我们直接在倒排索引上进行检索，将导致搜索效率低下，同时也失去了初步处理原始数据以提高后续工作效率的优势。因此，我们需要对倒排索引进行分区。分区后，搜索工作只需根据查询的单词首字母就能够定位到相应的分区，从而缩小了检索范围。例如，查询单词为“goat”，在倒排索引未分区时需要遍历整个索引文件才能找到单词信息。而在分区后，只需根据首字母“g”即可定位到相应的分区，大幅提高了检索效率。

- 建立单词地址簿

为了进一步优化查询效率，我们可以在一个文件中记录每个单词的单词地址。查询时，根据查询所对应的单词地址，可以直接定位到倒排索引中的单词信息，加速了查询过程。

- 依赖关系的理解

搜索引擎的建立是一个复杂的过程，从原始数据到倒排索引再到二级索引的构建。然而，搜索查询的实际工作是从二级索引开始，然后定位到倒排索引，最终找到与查询相关的原始数据。因此，二级索引功能是依赖于倒排索引的，倒排索引是搜索引擎的核心组成部分。

C. 选型依据

根据功能需求分析，我们需要一个解决方案，它能够高效地存储和管理大量文件，同时具备数据并行处理的能力。以下是性能要求和选型依据的优化和扩展：

功能需求：

我们需要解决存储和管理大量文件以及数据并行处理两个主要要求。

性能要求：

- 并发性： 由于原始数据量较大，我们需要解决方案能够实现并发处理多个任务，以提高效率。
- 响应时间： 响应时间需要尽可能灵活，根据数据量和任务复杂性在几秒到数分钟之内完成响应。尽管不要求实现即时搜索引擎，但在几秒到数分钟之内的响应仍然是可接受的范围。
- 吞吐量： 搜索引擎只需要在合理时间内处理给定数据集的搜索工作。由于原始数据大小约为 4.0GB，吞吐量（MBps）是单位时间内倒排索引程序每秒处理的数据。考虑到计算机性能，我们可以接受较低的吞吐量。
- 成本： 我们希望解决方案成本较低，最好是开源的，并且能够使用廉价机器组成服务器集群来处理数据。
- 生态系统： 我们需要一个具备强大生态系统的解决方案，包括各种工具、库和组件，以解决各种问题，减少整合和开发的复杂性，从而提高项目的生产力。
- 社区支持： 一个活跃的开源社区对于知识共享和技术支持非常重要。
- 文档和资源： 丰富的文档和教程等资源可以帮助团队成员更快地学习和掌握技术，降低学习曲线。

综合考虑以上要求，我们发现 Hadoop 生态圈是一个能够满足所有项目需求的强大解决方案。它包括 Hadoop 分布式文件系统（HDFS）用于存储和管理大量文件 MapReduce 等分布式计算框架用于数据并行处理，满足并发性需求。同时，Hadoop 具备灵活的响

应时间，并且可以在廉价机器上构建服务器集群，符合低成本要求。此外，Hadoop 拥有一个活跃的开源社区，提供丰富的文档和资源，帮助团队更好地应对挑战。因此，Hadoop 生态圈技术被认为是一个高度匹配的离线搜索引擎解决方案。

D. Hadoop 进行开发搜索引擎的优势与不足

优势：

- 分布式架构： Hadoop 具备分布式架构，将大数据处理引擎靠近数据存储，提高数据本地性。这减少了数据传输和网络开销，从而提高了数据处理效率。
- MapReduce 框架： MapReduce 框架支持将单个任务分解并并行处理，然后将结果加载到数据仓库中。这种并行计算能力对于复杂的计算任务，如构建倒排索引，有助于提高处理速度和效率。
- 容错性： Hadoop 自动保存数据的多个副本，并能够自动重新分配失败的任务，保障了系统的容错性。
- 可扩展性： Hadoop 能够在可用的计算机集群之间分配数据并完成计算任务，使系统具备了良好的可扩展性。
- 生态系统： Hadoop 拥有完整的生态系统，包括 MapReduce 框架、YARN 资源管理系统、HDFS 分布式文件系统、HBase 分布式数据库、Zookeeper 分布式协作服务等，提供了实现离线搜索引擎所需的所有服务和工具。

不足：

- 资源消耗： Hadoop 需要大规模的硬件资源来运行，这包括硬件成本、电力消耗和维护成本。在小型搜索引擎项目或资源有限的环境中，部署和维护 Hadoop 集群可能导致较高的资源消耗，不太适合在资源有限的情况下选择 Hadoop 作为搜索引擎的基础技术。
- 实时性： Hadoop 最初设计为批处理和离线数据处理，MapReduce 框架的任务调度和执行需要一定时间，无法提供低延迟的查询响

应。因此，在实时性方面存在不足，不适合快速响应用户的查询并即时更新索引。

- 数据一致性： Hadoop 的 HDFS 使用数据副本来提高容错性，这可能导致数据更新问题，因为数据的多个版本存在于不同节点上。如果用户同时查询不同的节点，不一致的数据可能导致错误的搜索结果，引发数据一致性问题。

3. 功能实现

首先，我们开始处理老师提供的数据集 “sentence.txt”。我们编写了自己的 Python 程序，根据每一行的 ID，将数据划分成多个文件，总计生成了 940+ 个文件，以获取原始数据。值得注意的是，通常情况下，一个真正的搜索引擎会使用爬虫技术来收集数据。

然而，由于本次实验的目标是实现离线搜索引擎，因此我们仅使用了老师提供的数据集。接下来，我们将这些原始数据导入倒排索引程序进行处理。一旦倒排索引构建完成，单词将被按照首字母 0-9 和 a-z 进行分区，共生成了 36 个文件。每个文件的格式是每行一个单词，其后是包含该单词的文件列表，根据单词在文件中的频率进行排序，最终构建了倒排文件。

考虑到倒排索引文件的庞大，为了提高搜索效率，我们编写了二级索引程序，生成了二级索引文件。这些文件记录了单词在倒排文件中的偏移量，也就是实际位置，以便更快地定位单词的位置，从而显著提高了搜索效率。最后，我们将倒排文件存储在 HBase 数据库中的一个表中，以便进行检索。在这个过程中，我们采用了倒排内的分区和二级索引两种方法，以进一步提高搜索效率。当用户输入搜索内容后，搜索引擎会将搜索内容分割成单词，并在二级索引文件中查找它们在倒排索引分区文件中的位置。接着，我们提取这些单词所在文件的信息，并计算它们的相关性之和。

最后，将这些文件按照相关性从高到低的顺序呈现给用户，以提供高度相关的搜索结果。这样，我们已经完成了数据处理的所有工作，用户可以获得高效的搜索体验。

4. 工作计划

序号	工作安排	时间
1	学习 Hadoop 相关内容，包括 Hadoop 环境搭建、HDFS、MapReduce、Yarn、HBase、Zookeeper 等	4/9 - 10/9
2	搭建 Hadoop 分布式环境	6/9-7/9
3	运用 MapReduce 框架实现倒排索引	8/9-6/10
4	倒排索引存入 HBase	1/10-5/10
5	第一次优化搜索引擎，运用 MapReduce 框架在倒排索引的基础上实现二级索引	23/9-28/9
6	第二次优化搜索引擎，partition	27/9
7	整合实现搜索引擎	4/10-5/10
8	编写文档	5/10-8/10

5. 组织结构

组长：洪子翔

组员：何秉翰、陈晓如、杨珍奇、金泰显

分工	分工内容	人员以及具体分工内容	
算法小组	负责 MapReduce 框架	何秉翰	倒排索引， 文件代码优化， 倒排优化
		洪子翔	文件分割， 倒排索引分区
		陈晓如	二级索引
Hbase、 Zookeeper 小组	负责 HBase 框架	洪子翔	HBase 使用
		何秉翰	HBase 使用
		杨珍奇	HBase 使用
		金泰显	HBase 使用
环境搭建组	负责搭配 Hadoop 集群 ， Zookeeper， HBase 和 IDEA	洪子翔	hadoop 集群搭建， Zookeeper 搭建， HBase 搭建， IDEA 搭建
		何秉翰	hadoop 集群搭建， ZooKeeper 搭建， HBase 搭建 IDEA 搭建
文档组	编写搜索引擎技术实现 方案和实验报告	杨珍奇；全体成员	主要文档工程负责 人，全部队员

6. 软件质量保证、非功能保证性

A. 软件质量保证

- 质量目标

在合理时间范围内完成原始数据的处理，生成倒排索引和二级索引文件，实现搜索引擎的离线工作过程。

- 测试策略

模拟源文件生成测试用例，模仿程序输入的源文件内容，生成一个简单的版本来节省测试程序的时间。例如倒排索引程序需要输入原始数据文件，则模拟原始数据的样式来生成一个小文件，用于试运行倒排索引程序，减少时间代价。

B. 非功能性保证

- 性能保证

在资源管理系统中获取程序的响应时间，按测试结果，响应时间范围一般在几秒到数分钟以内。

- 可用性保证

Hadoop 的特点之一是自动保存数据的多个副本，因此数据丢失是可自动恢复的。

- 可维护性保证

将软件版本控制在与搜索引擎相同的版本就可以避免兼容问题，便于维护。