

海量数据，选择批处理？还是流处理？

嗨，你好，我是**猿java**

在大数据时代，我们常常需要处理各种量级以及不同场景的数据，通常有`批处理 (Batch Processing)` 和`流处理 (Stream Processing)` 两种方式，那么它们是如何工作的？两者之间有存在什么区别？我们又该如何选择？这篇文章我们将一一解答这些问题。

批处理

什么是批处理？

批处理 (Batch Processing) 是一种传统的数据处理方法，它在一段时间内（如小时、天、甚至数周）收集和存储数据，然后在预定的时间间隔内对这些数据进行批量处理。如下示意图：

![batch-processing.png](https://p6-xtjj-sign.byteimg.com/tos-cn-i-730wjymdk6/665089f2e6d54aaa8fdec93f599bc6ad~tplv-730wjymdk6-watermark.image?rk3s=f64ab15b&x-expires=1722045715&x-signature=M4u8k6s17T9JlyVbXTqMPODkNOM%3D)

适合的场景

批处理这种方法非常适合以下任务：

- * **高吞吐量：**批处理在处理大量数据时表现出色，适用于数据仓库、数据集成和数据清洗等任务。
- * **复杂处理：**批处理适用于需要复杂算法、数据转换和数据汇总的任务。
- * **离线处理：**批处理通常在离线模式下进行，数据以批次处理，结果存储供以后使用。

比如：

- * 数据仓库和ETL（抽取、转换、加载）过程
- * 数据备份和归档
- * 科学模拟和数据分析

批处理的优点

- * 效率：批处理通过大批量处理数据来优化资源利用，减少频繁加载和处理数据的开销。
- * 简便性：由于处理静止数据，不需要复杂的事件处理系统，因此批处理通常较容易实现和管理。
- * 可扩展性：批处理系统可以通过在多台机器上分配工作负载进行水平扩展，从而处理海量数据集。
- * 容错性：批处理框架（如Apache Hadoop和Apache Spark）提供内置的容错机制，确保数据完整性并允许恢复失败的作业。
- * 成本效益：与实时处理相比，批处理在非高峰时段高效利用计算资源，可能更具成本效益。

批处理的缺点

- * 延迟：批处理在数据收集和处理之间引入了延迟，因为数据按计划的时间间隔批量处理。这种延迟可能不适合需要实时或接近实时结果的应用程序。
- * 缺乏实时洞察：由于存在延迟，批处理可能无法提供最新洞察或实现实时决策。
- * 对突变响应不佳：批处理系统可能无法立即适应意外的数据激增或模式变化。
- * 数据陈旧：在处理完成时，批处理的数据可能已经陈旧，尤其是在快速变化的环境中。
- * 资源密集：大批量处理可能需要大量计算资源，如果管理不当，可能影响系统性能。

使用框架

批处理框架在处理大量数据时非常有用，以下是一些流行的批处理处理框架：

Apache Hadoop

- * 特点：分布式存储和处理大数据。
- * 组成部分：HDFS（分布式文件系统）和MapReduce（分布式计算框架）。
- * 优点：可靠性、可扩展性、高容错性。

Apache Spark

- * 特点：处理速度快，支持批处理和实时处理。
- * 组成部分：核心引擎和丰富的库，如Spark SQL、MLlib、GraphX和Spark Streaming。
- * 优点：内存计算，大规模数据处理，支持多种语言（如Java、Scala、Python、R）。

Apache Flink

- * 特点：高性能，低延迟的数据流处理。
- * 组成部分：流处理和批处理引擎。
- * 优点：统一的批处理和流处理，事件时间处理，容错性好。

Apache Beam

- * 特点：统一的编程模型，可在多种执行引擎上运行。
- * 组成部分：编程模型和SDK。
- * 优点：可移植性，支持批处理和流处理，兼容多种执行引擎（如Apache Flink、Apache Spark和Google Cloud Dataflow）。

流处理？

=====

什么是流处理？

流处理（Stream Processing）是一种数据处理方法，它在数据到达时立即对其进行处理，而不是像批处理那样等待一段时间后再进行处理。如下示意图：

![stream-processing.png](https://p6-xtjj-sign.byteimg.com/tos-cn-i-730wjmld6/430fe4dd8ad341c1b9ef19d3ab02091b~tplv-730wjmld6-watermark.image?rk3s=f64ab15b&x-expires=1722045715&x-)

适用场景

流处理特别适用于以下任务：

- * **实时洞察**：适合需要即时洞察和响应的任务，如实时监控、欺诈检测和实时推荐系统。
- * **低延迟任务**：适合需要快速处理和低延迟的应用，如物联网传感器数据处理和社交媒体监控。
- * **持续处理任务**：适合需要持续不断处理数据的任务，如金融市场交易监控和网络安全监控。

比如：

- * **实时欺诈检测和预防**：在金融交易中实时检测和防止欺诈行为。
- * **社交媒体监控和情感分析**：实时监控社交媒体平台上的内容，分析用户情感和趋势。
- * **物联网传感器数据处理和分析**：实时收集和分析物联网设备生成的数据，如环境传感器数据和健康监测设备数据。

流处理的优点

- * **实时处理**：流处理允许立即处理数据，这对于依赖及时洞察的应用程序至关重要，如欺诈检测系统或实时监控工具。
- * **可扩展性**：现代流处理框架设计为可以横向扩展，意味着可以通过增加更多资源来处理越来越多的数据量。
- * **灵活性**：流处理器可以适应数据类型和处理逻辑的变化，而不会导致显著的停机时间。

流处理的缺点

- * **复杂性**：管理流处理系统可能很复杂，因为需要连续处理数据流并具备故障容错机制。
- * **更高的成本**：根据基础设施的不同，流处理可能会由于需要更强大的系统和持续运行而产生更高的运营成本。

* **数据丢失的潜在风险**: 在系统故障的情况下, 存在丢失实时数据的风险, 除非有先进的故障容错措施。

使用框架

常见的流处理框架包括:

- * **Apache Kafka**: 一个分布式流处理平台, 用于构建实时数据流应用和数据管道。
- * **Apache Flink**: 一个用于分布式流处理和批处理的开源框架, 支持低延迟和高吞吐量的数据处理。
- * **Apache Storm**: 一个分布式实时计算系统, 适用于处理高速和大规模的数据流。

批处理和流处理的区别

批处理是指按计划的时间间隔处理大量数据, 而流处理则是指实时或近乎实时地处理数据。两者的主要区别主要有以下几点:

- * **数据处理的定义和性质**
- * **延迟和处理时间**
- * **使用案例和应用**
- * **容错性和可靠性**
- * **可扩展性和性能**
- * **复杂性和设置**
- * **工具和平台示例**

数据处理的定义和性质

* **批处理**: 涉及处理大块数据, 或在一定时期内收集数据后进行批量处理, 数据被存储起来, 一旦有足够的数据, 或者在经过一定时间后, 它就会立即被处理

◦ **流处理**: 流处理旨在实时或近乎实时地处理数据, 一旦数据到达, 它就会被处理, 并不需要等待一批数据累积。

延迟和处理时间

* 批处理：通常具有更高的延迟，因为数据不会立即处理，它等待批处理完成或等待特定计划触发处理。

* 流处理：提供更低的延迟，因为数据在流入系统时会立即进行处理，这使得它更适合实时分析或需要即时洞察的任务。

用例和应用

* 批处理：在不需要立即处理数据的情况下很常见。比如月度工资单处理、日终报告生成和大规模数据分析。

* 流处理：用于需要根据传入数据立即采取行动的情况，例如银行业中的欺诈检测、电子商务中的实时推荐或实时仪表板更新。

容错性和可靠性

* 批处理：如果批处理作业失败，可以从中断的位置重新启动，也可以重新处理整个批处理。

* 流处理：需要更复杂的容错机制，如果数据流中断，系统需要处理中断并确保数据不会丢失的方法。

可扩展性和性能

* 批处理：由于一次处理大量数据，因此系统通常针对吞吐量进行了优化。根据用例，它们可能会垂直扩展（更强大的计算机）或水平扩展（更多计算机）。

* 流处理：系统需要针对高吞吐量和低延迟进行设计。它们通常水平缩放以处理不同的数据速度。

复杂性和设置

* 批处理：可能具有更简单的设置和设计，因为它并不总是需要考虑实时处理的复杂性。

* 流处理：通常需要更复杂的设置，尤其是在确保容错、管理状态和处理无序数据事件时。

工具和平台示例

- * 批处理：常见于 Hadoop MapReduce、Apache Hive 和 Apache Spark 等框架。
- * 流处理：常见于 Apache Kafka Streams、Apache Flink 和 Apache Storm 等框架。

如何选择？

当我们在批处理和流处理两者之间进行选择时，需要考虑以下因素：

- * 数据量：如果处理大量数据，批处理可能是更好的选择。
- * 实时需求：如果应用需要根据输入数据立即提供洞察或采取行动，流处理是更适合的方法。
- * 复杂性：如果处理任务需要复杂算法和数据转换，批处理可能更合适。
- * 数据性质：您的数据是有限且可预测大小，还是一个无界、持续的数据流？批处理更适合前者，流处理更适合后者。

总结

批处理和流处理是数据管理中的两种截然不同的范式，批处理更提供高吞吐量，它通常具有更高的延迟。

流处理可以连续、实时地管理数据，非常适合实时分析和监控，由于其始终在线的性质，它需要有弹性的系统。

批处理可能占用大量资源且灵活性较差，而流可能面临一致性问题，它们之间的选择取决于所讨论的数据任务的特定需求。

学习交流

如果你觉得文章有帮助，请帮忙点个赞呗，或者公众号：猿java，持续输出硬核文章。

原文链接: <https://juejin.cn/post/7390588322768683044>