



程序代写代做 CS 编程辅导

MONASH
INFORMATION
TECHNOLOGY



FIT5202 – Data Processing for Big Data

WeChat: esutores

Stream Data Processing

Prajwol Sangat

Edited by Ting Chee Ming (7 May 2021)

QQ: 749389476

<https://tutorcs.com>



Last Lecture

程序代写代做 CS编程辅导

- Collaborative Filtering



WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

This Week

程序代写代做 CS编程辅导

- **Reminder:** Student Evaluation of Teaching and Units (SETU)
- Stream data processing
- Streaming processing technology



WeChat: cstutorcs

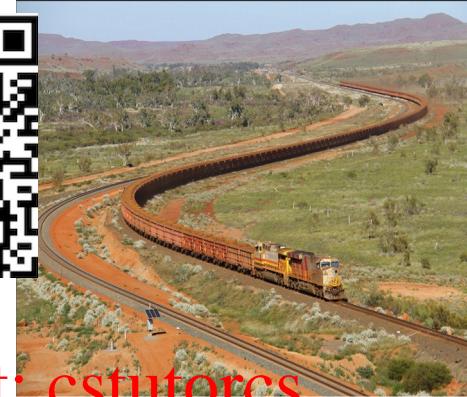
Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

Data Stream Applications



#trending

WeChat: cstutorcs

Track monitoring



Email: tutorcs@163.com

QQ: 749389476

#Stock market

<https://tutorcs.com>



#Security

Characteristics of the application

- Unbounded (infinite) data: Data comes in without boundary
- Real time query at any given time
 - Get response quickly
 - Differ from batch query applied to entire table/data block
- Data comes continuously
 - Uniform rate (e.g., data comes every 1 second)
 - Bursty (e.g., data records come in one single reading)
- Interested in an “event” or trends across time.



WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

Data stream

程序代写代做 CS编程辅导

A data stream is a real-time sequence of data items, ordered (implicitly by arrival time or explicitly by timestamp) set.



WeChat: cstutorcs

Assignment Project Exam Help

t_n

Email: tutorcs@163.com

8

QQ: 749389476

<https://tutorcs.com>



Data Stream Query Examples

What is the total number of attendees?

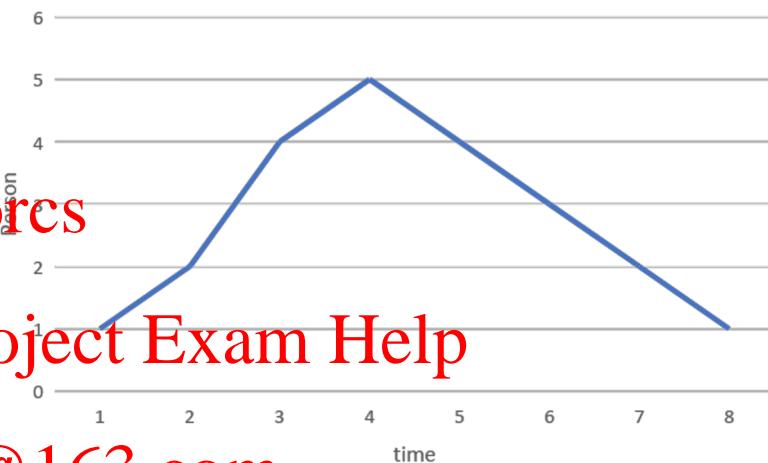


What is the average number of attendees
(at time t_x)

Will the next number be lower or higher than the current reading?

WeChat: cstutorcs

Assignment Project Exam Help



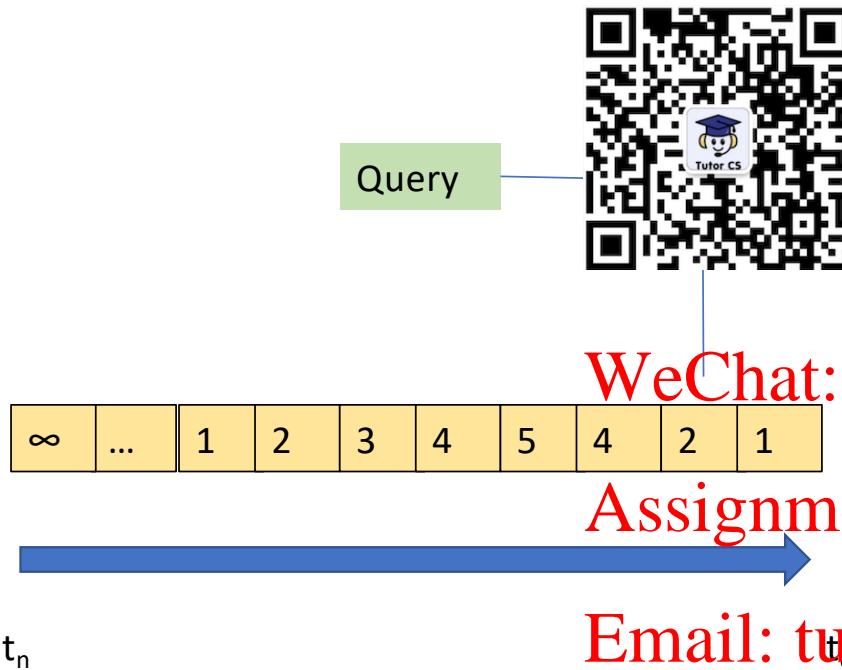
Email: tutorcs@163.com

t_n

QQ: 749389476

<https://tutorcs.com>

What is the total number of attendees up to t_x ? 程序代写代做CS编程辅导



Assume memory can only keep 5 tuples.

Keep an accumulator for the sum. There is no need to keep the tuples.

Possible issues:

- What to do with incoming tuples if the incoming tuple arrival rate > processing time?

<https://aseigneurin.github.io/2018/08/27/kafka-streams-processing-late-events.html>

Data streams: Can't keep all incoming data because data is unbounded

Database: Have persistent storage to store all data

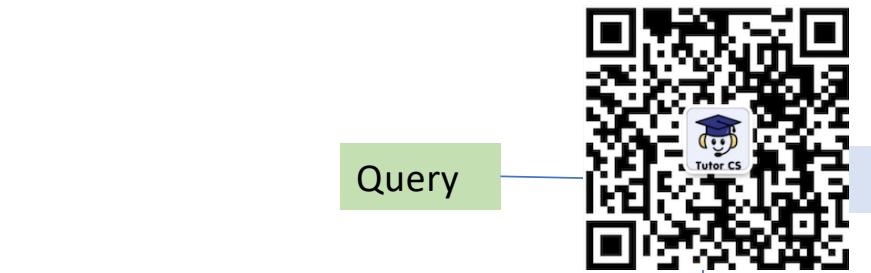
→ Rely on accumulator as a summary of the stream
<https://tutorcs.com>

What is the average number of attendees so far? (at time t_x)



Keep two accumulators for the sum and count. There is no need to keep the tuples.

Will the next number be lower or higher than the current reading?



∞	...	1	2	3	4	5	4	2	1
----------	-----	---	---	---	---	---	---	---	---

WeChat: cstutorcs

Assignment Project Exam Help

Prediction may need all/some previous tuples.
Memory can only keep 5 tuples.
What can we do?

t_n

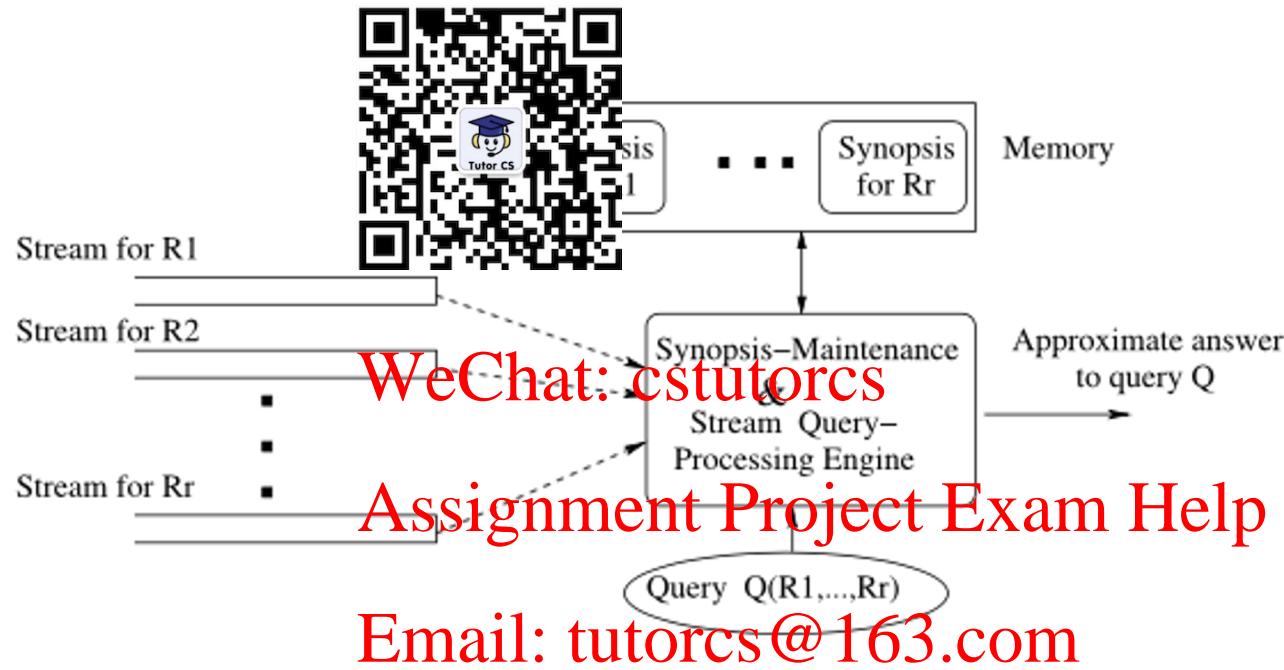
Learn from previous data (create summary/model based on previous data)

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

Stream Processing System Architecture



Garofalakis M., Gehrke J., Rastogi R. (2016) Data Stream Management: A Brave New World. In: Garofalakis M., Gehrke J., Rastogi R. (eds) Data Stream Management. Data-Centric Systems and Applications. Springer, Berlin, Heidelberg

<https://tutorcs.com>

Time Decay

程序代写代做 CS编程辅导

- Is it possible to reduce the influence of older tuples, without eliminating them?
- Popular approach: window
- Query processing is performed only on the tuples inside the window.
- Window types:
 - Time-based
 - Tuple-based (count-based)



WeChat: cstutorcs

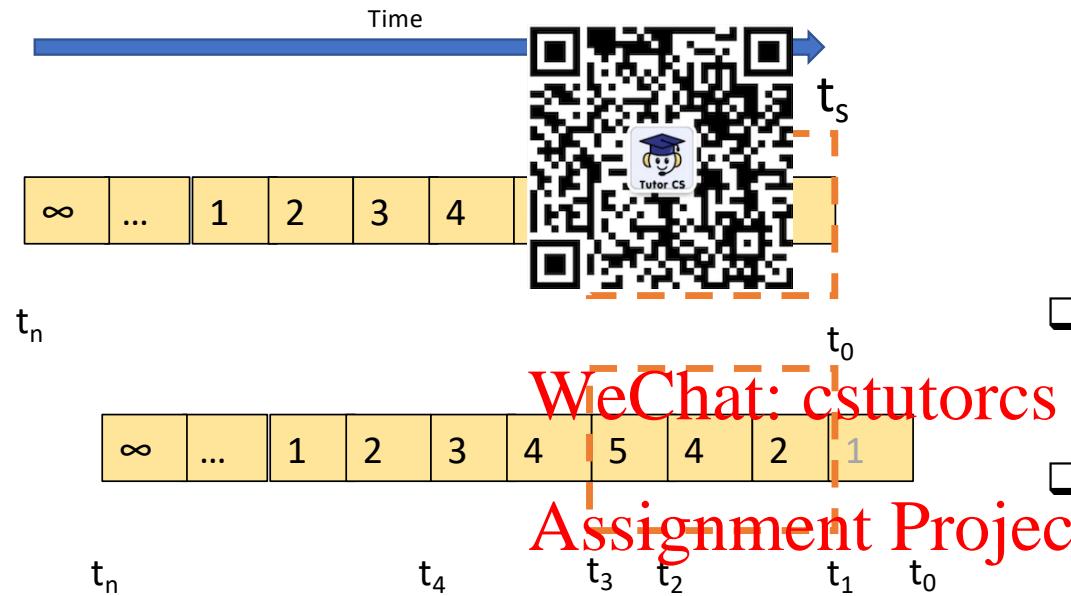
Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

Stream Window – Time Based Window

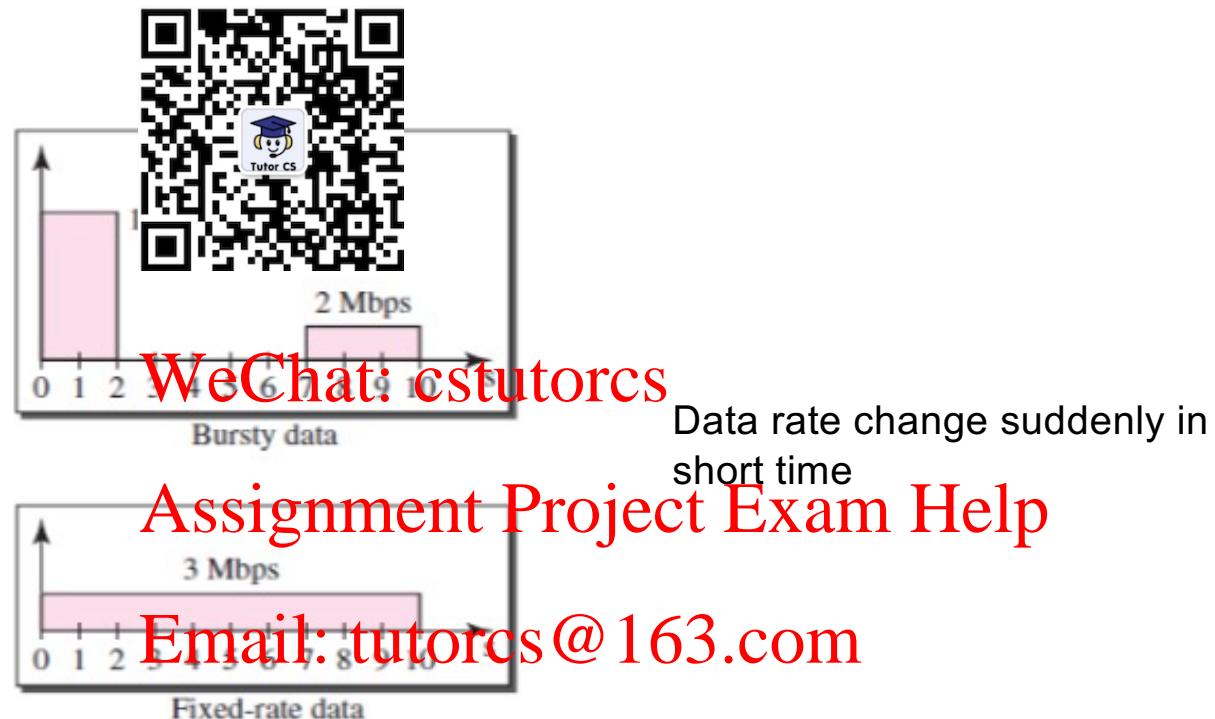
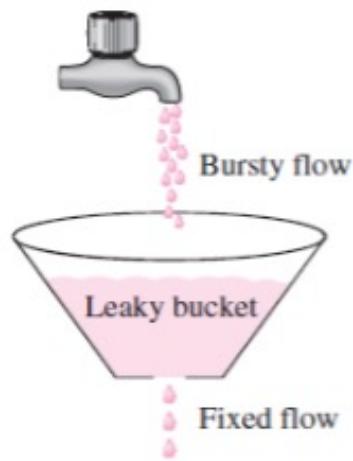


Window size 2 seconds
Slides 1 second.

- Size of window is specified by fixed time duration
- Number of tuples in windows may vary depending on how many tuples are available while processing
- It may be high due to bursty arrival of data due to delay in data transfer
- Tuples in the windows will be processed, no matter how many tuples in the windows

程序代写代做 CS编程辅导

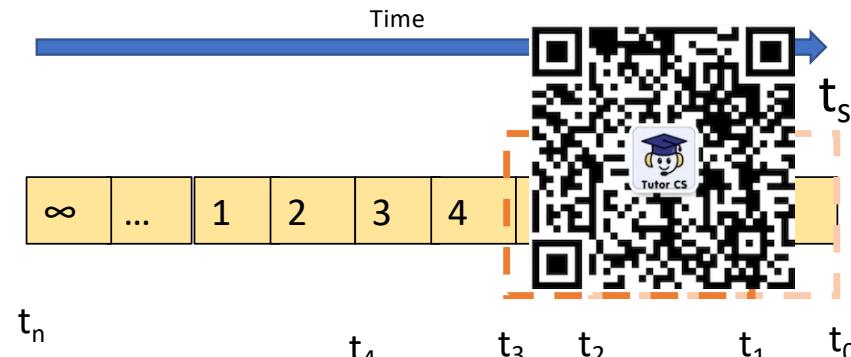
Bursty data



<https://tutorcs.com>

Image source: <https://www.quora.com/What-is-the-difference-between-token-bucket-and-leaky-bucket-algorithms>

Stream Window – Time Based Window



Window size 2 seconds
Slides 1 second.

WeChat: cstutorcs

Overlapping Sliding Window

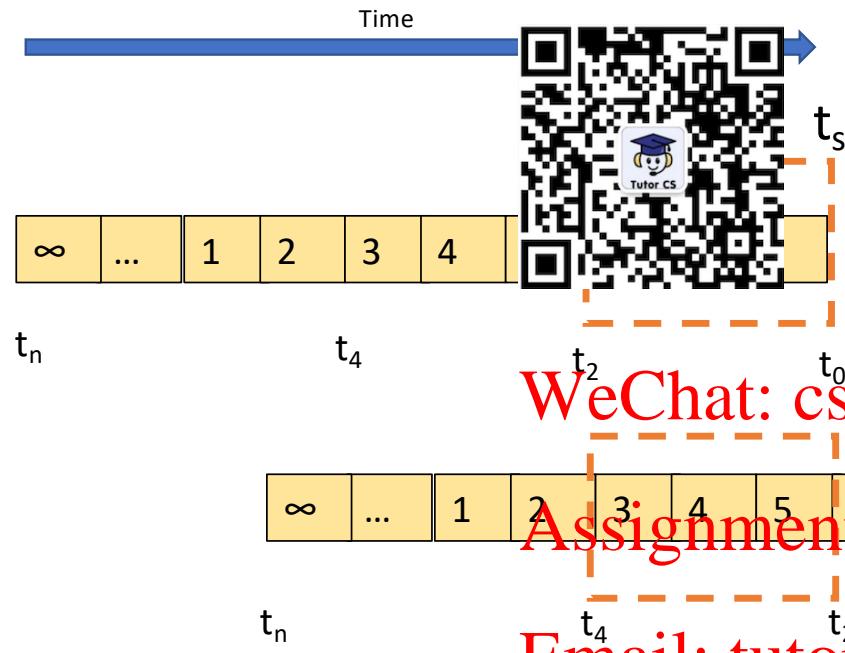
Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

Stream Window – Time Based Window



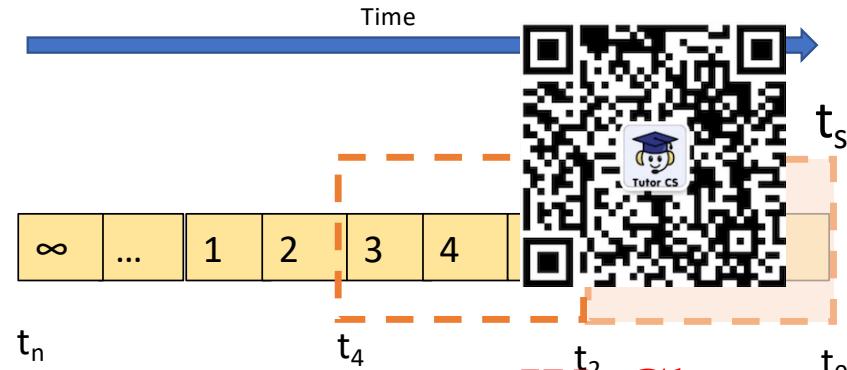
- Window size is based on time, eg 2 seconds
- Window can be advanced by:
 - $t_e - t_s$ (window size)
 - Duration less than the window size (sliding window).
- In uniform data rate, the number of tuples will be the same for each window.

Non-Overlapping Sliding Window

QQ: 749389476

- Slide of window is equal to size of window
- Two windows are disjoint
- No influence of data in previous window on current window

Stream Window – Time Based Window



- Window size is based on time, eg 2 seconds
- Window can be advanced by:
 - $t_e - t_s$ (window size)
 - Duration less than the window size (sliding window).
- In uniform data rate, the number of tuples will be the same for each window.

Non-Overlapping Sliding Window

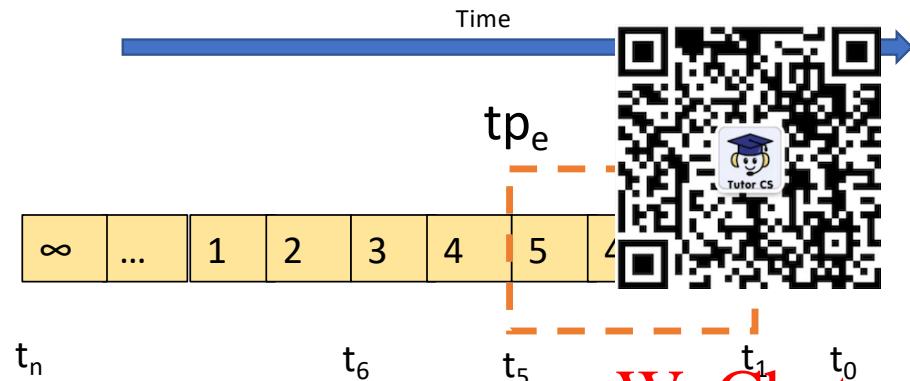
Assignment Project Exam Help

Email: tutorcs@163.com

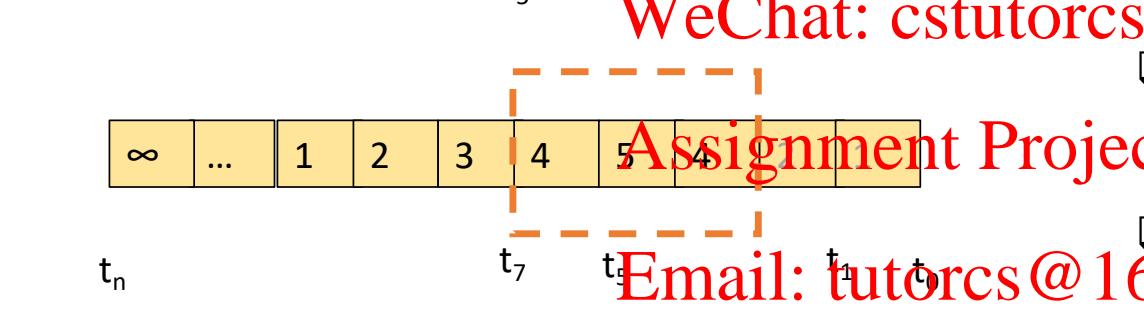
QQ: 749389476

<https://tutorcs.com>

Stream Window – Tuple Based Window



Window size is 3 tuples
Slide window by 1 tuple

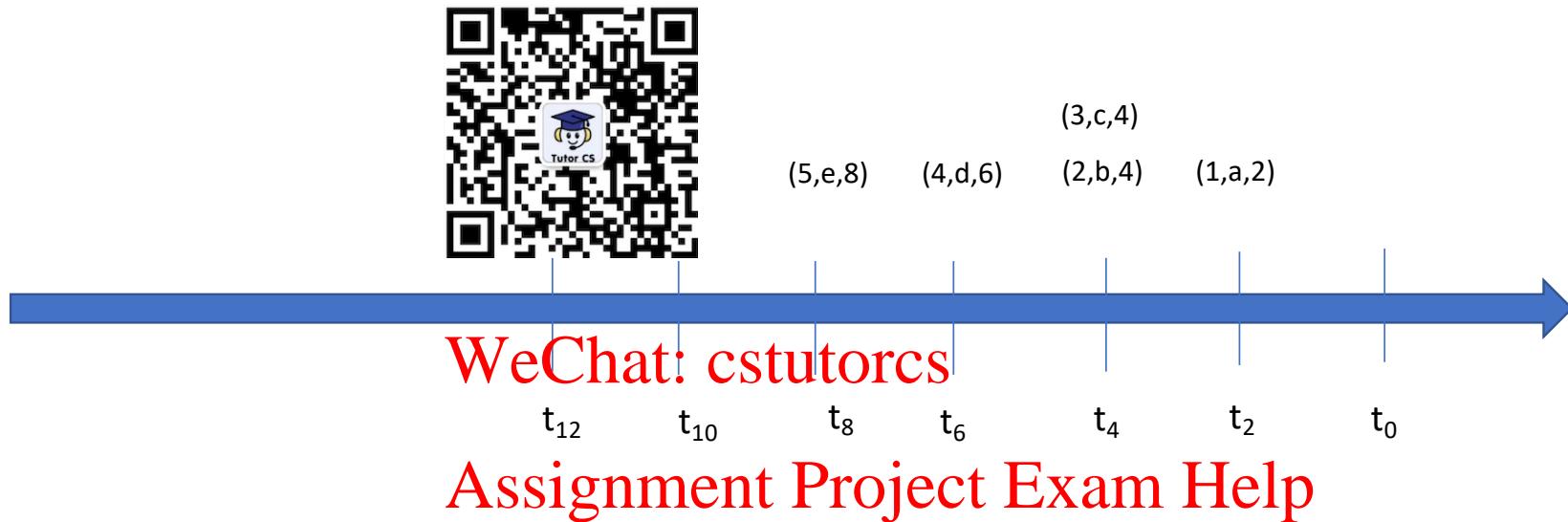


- ❑ Size of window is specified by number of tuples in the window
- ❑ The window always has a fixed number of tuples
- ❑ We count the number of tuples and create windows. If the window has required numbers of tuples, then that window will be processed
- ❑ The time duration of windows may vary

QQ: 749389476

<https://tutorcs.com>

Example of Burst Data Arrival



(eventTime, value, processingTime)

Email: tutorcs@163.com

Event Time: The time stamp when the data is generated

Processing Time: the time stamp when the data arrived at the processing.

QQ: 749389476

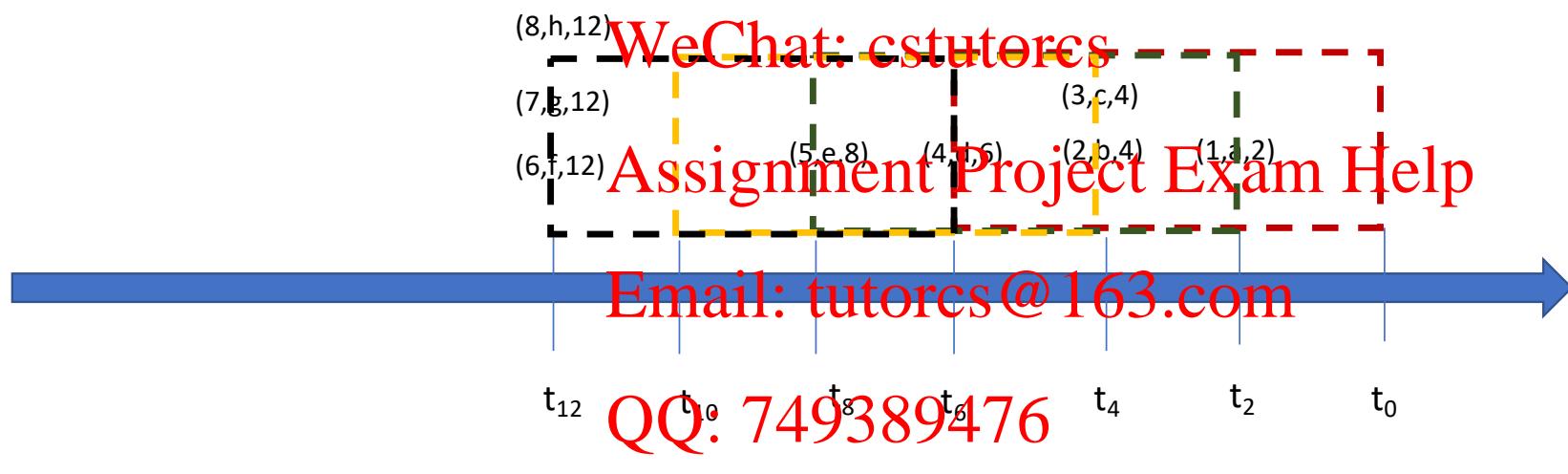
<https://tutorcs.com>

Time Based Window

- Processing time window
- Window size 6 units, slide by 1 unit

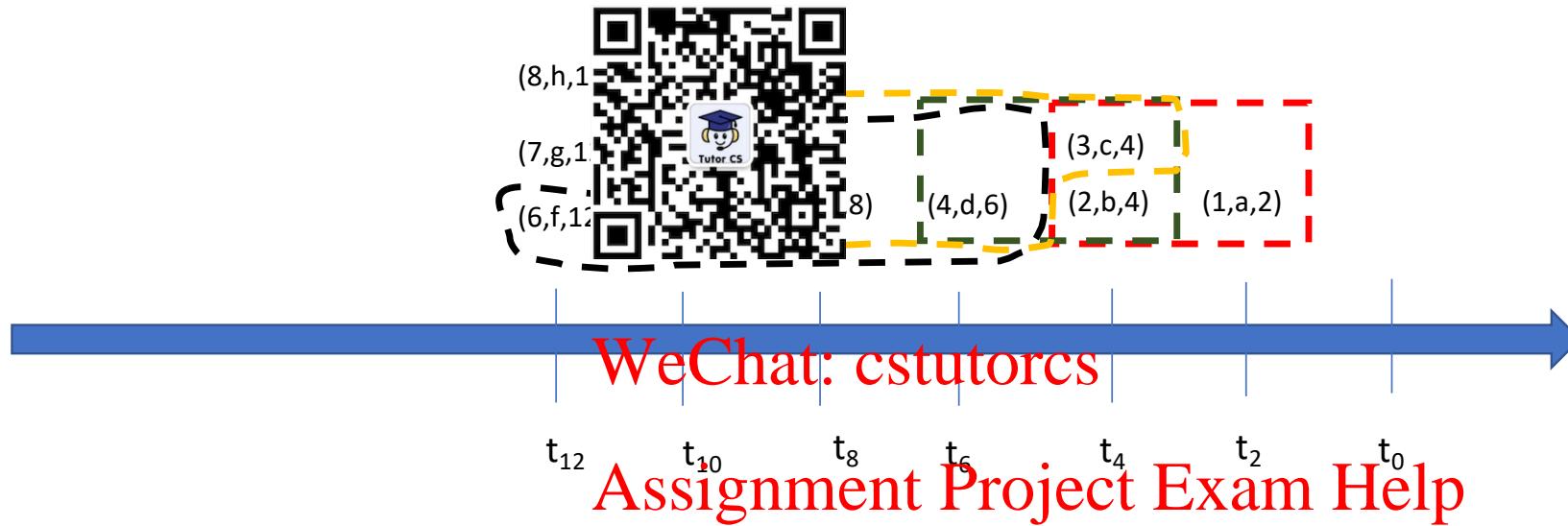


$W_1\{(1,a,2),(2,b,4),(3,c,4),(4,d,6)\}$
 $W_2\{(1,a,2),(2,b,4),(3,c,4),(4,d,6),(5,e,8)\}$
 $W_3\{(2,b,4),(3,c,4),(4,d,6),(5,e,8)\}$
 $W_4\{(4,d,6),(5,e,8),(6,f,12),(7,g,12),(8,h,12)\}$



<https://tutorcs.com>

Tuple Based Window



Window size is 3 tuples.
Slide by 1 tuple

Email: tutorcs@163.com

$W_1\{(1,a,2),(2,b,4),(3,c,4)\}$
 $W_2\{(2,b,4),(3,c,4),(4,d,6)\}$
 $W_3\{(3,c,4),(4,d,6),(5,e,8)\}$
QQ: 749389476

Need to drop some of the tuple
with processing timestamp 12

<https://tutorcs.com>

Practice Question



A data stream application processes data from a data source. The data source emits a single tuple per second (constant rate). The tuple received by the data stream application in the format of $(eventTime, value)$. Once it arrives at the application, the processing timestamp is added to the tuple. The final tuple format after the addition of the processing timestamp is $(eventTimestamp, value, processingTimestamp)$.

Examples of tuples are as followed:

{1,a,1}
{2,b,4}
{3,c,4}
{5,e,6}
{6,f,8}
{7,g,8}
{8,h,9}
{9,i,10}
{10,j,11}

WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

Based on the above data and scenario answer next exercise.

<https://tutorcs.com>

程序代写代做 CS编程辅导

Exercise 2



- Let's say we follow a timestamp windowing mechanism. The window size is three seconds and the overlap is two seconds. Assume the timing start at processing timestamp 1. What will be the content of the third window?

WeChat: cstutorcs

- A. {1,a,1}, {2,b,4}, {3,c,4}
- B. {2,b,4}, {3,c,4}, {5,e,6}
- C. {5,e,6}, {6,f,8}, {7,g,8}
- D. {3,c,4}, {5,e,6}, {6,f,8}

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

Event vs Processing time



- Event time: the time when data is produced by the source.
- Processing time: the time when data is arrived at the processing server.
- In ideal situation, event time equals processing time.
- In real world event time is earlier than the processing time due to network delay.
- The delay can be uniformed (ideal situation) or non-uniform (most of real network situation).
- Data may arrive in “burst” (bursty network).

WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

Database vs Stream Processing

- Bounded data (assessable data).
- Relatively static data
- Complex, ad-hoc query
- Possible to backtrack during processing
- Exact answer to a query
- Tuples arrival rate is low



WeChat: cstutorcs
Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

- Unbounded data (data keeps coming without end).
- Dynamic data.
- Simple, continuous query
- No backtracking, single pass operation.
- Approximate answer to a query
- Tuples arrival rates is high

Sliding Window: Produce an approximate answer to a data stream query is to evaluate the query not over the entire past history of data streams, but rather only over sliding windows of recent data from the streams.



MONASH
University

Stream Processing Technology

Streaming Data Processing
Prajwol Sangat

Updated by Ting Chee Ming

程序代写代做 CS编程辅导



WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>



Streaming Platforms



APAC
STO



WeChat: cstutorcs



Assignment Project Exam Help

Email: tutorcs@163.com

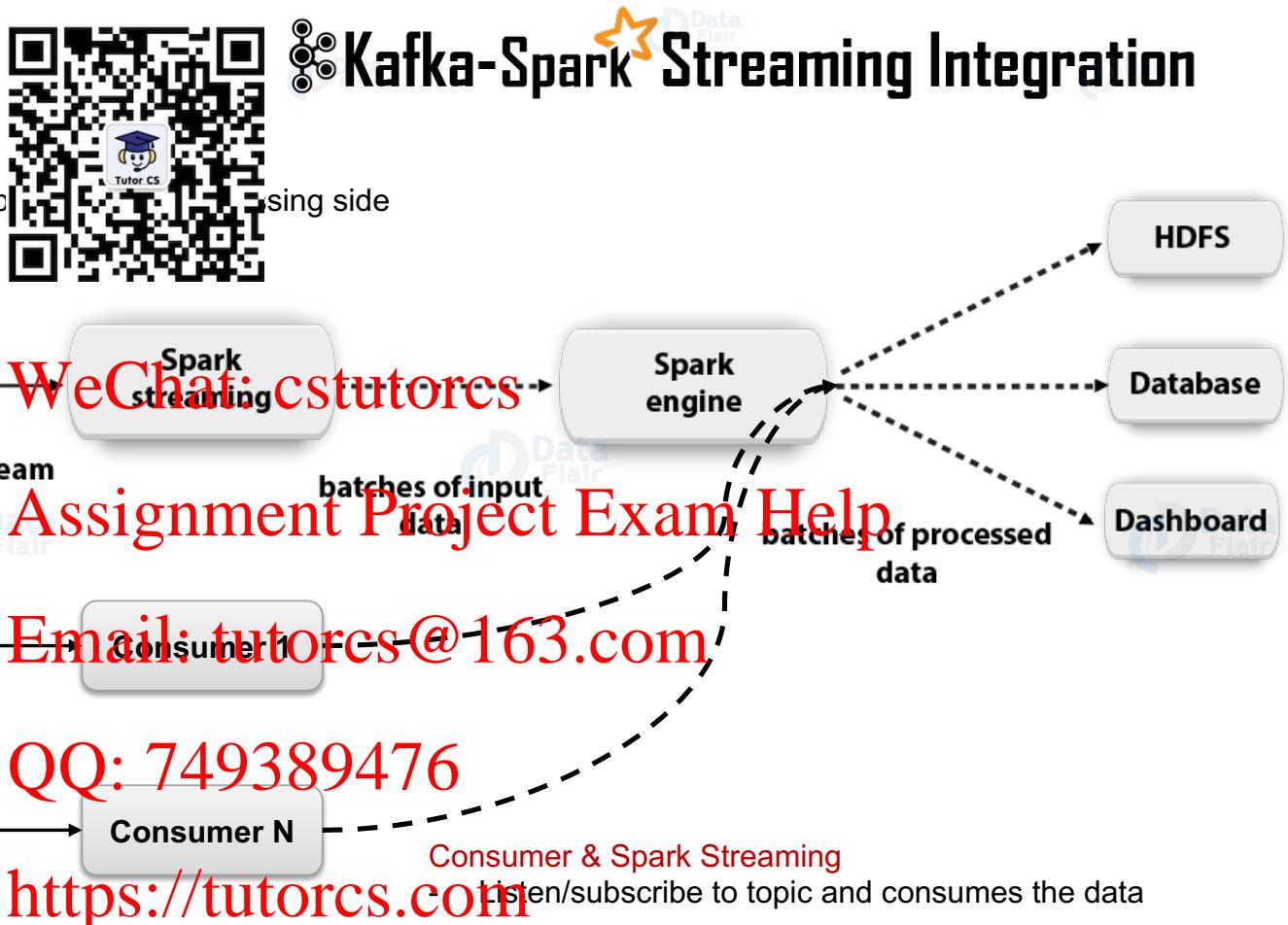
QQ: 749389476



A distributed streaming platform

<https://tutorcs.com>

Real-Time Streaming Architecture



What is Apache Kafka? 程序代写代做 CS 编程辅导

- Publish-subscribe messaging system
- Scalable, Fault-tolerant
- Enables distributed applications
- Powers web-scale Internet companies
LinkedIn, NetFlix, AirBnB, and many others.



WeChat: cstutorcs

□ Producer publish streams of data records into topics

□ Consumers subscribe to the topics and process the stream of records

□ Data in Kafka is stored in topic

- Topic is category/feed name where records are stored
- A topic is associated with a log – data structure on disk
- Each topic is indexed and stored with timestamp

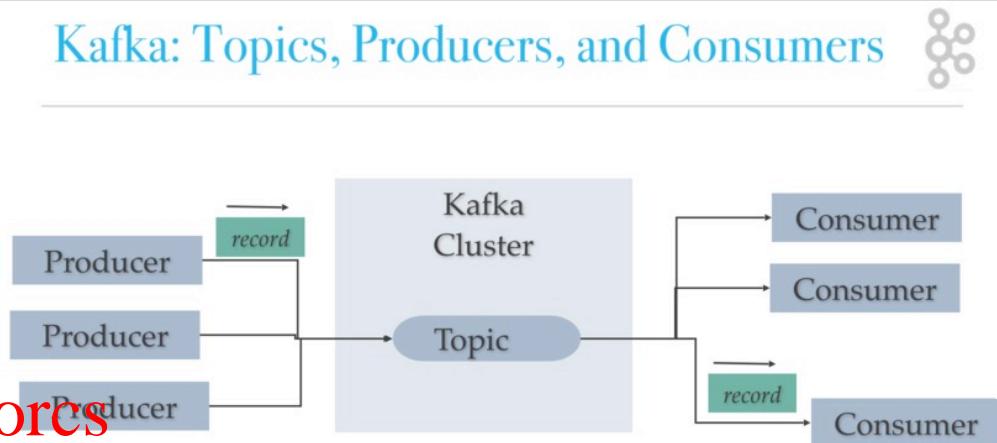
Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

Kafka: Topics, Producers, and Consumers



<https://dzone.com/articles/kafka-architecture>

□ Kafka is run as a cluster comprised of one or more servers (called brokers)

- A broker receive messages from producers and stores them on disk keyed by unique offset.
- A broker allows consumers to fetch messages by topic, partition and offset

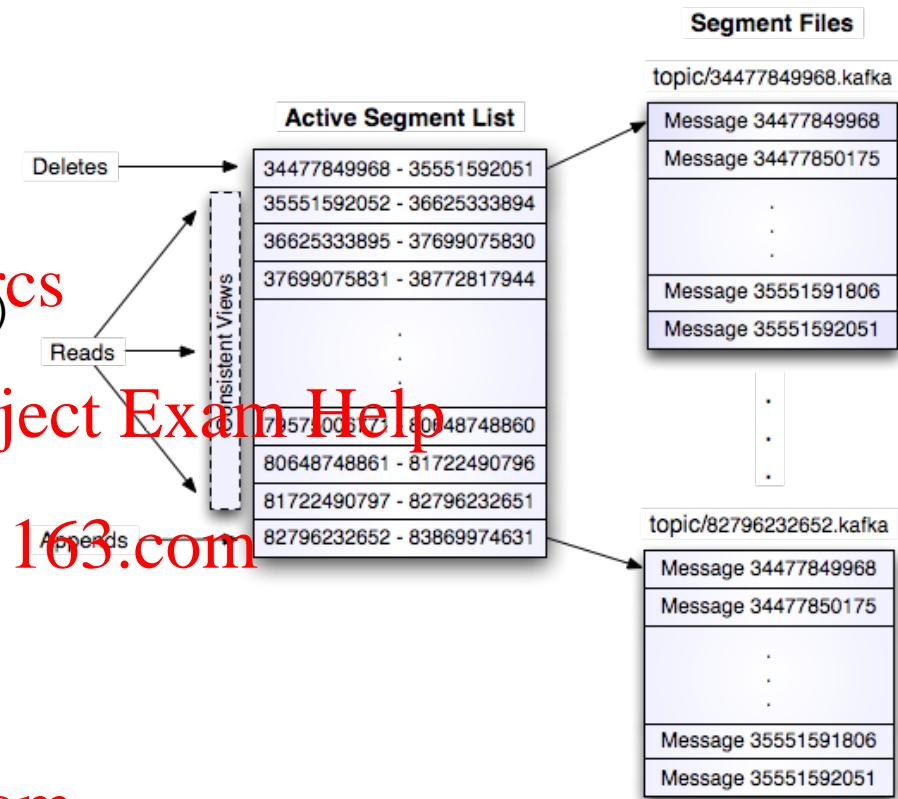
How does Kafka Work?

- Topics represent commit log data stored on disk



- Topics are divided into partitions which are further divided into segments
- Each segment has a log file to store the actual message
- Log – time-ordered, sequence of messages that is continually appended (each log entry can be array or bytes)
- Partitions are replicated and distributed over servers in Kafka cluster (brokers) for fault tolerance (when a server in the cluster fails so messages remain available)

Kafka Log Implementation



Email: tutorcs@163.com

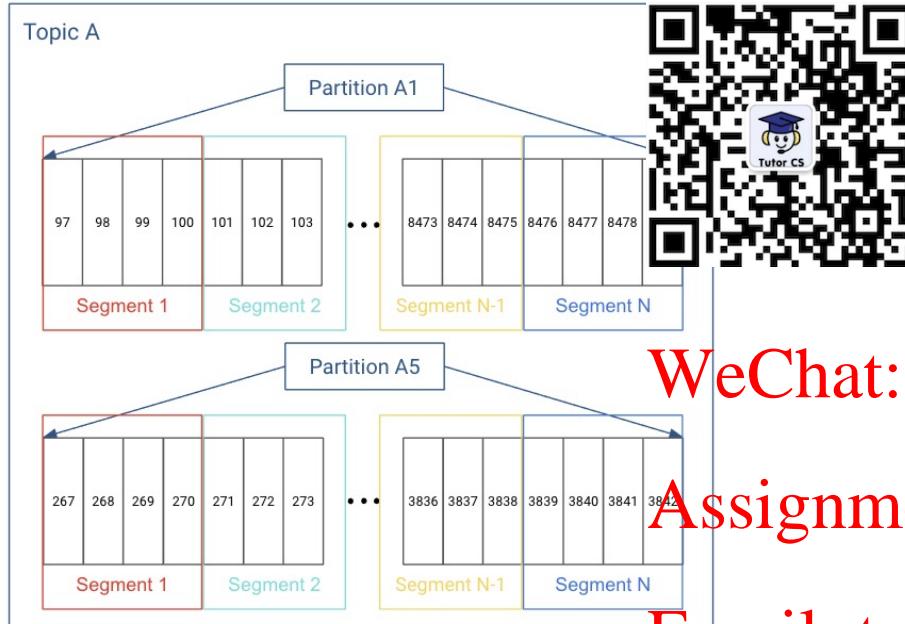
- Messages stay around for a configurable period of time (i.e., 7 days, 30 days, etc.).

QQ: 749389476

- Can recover lost messages during time out or lost connection

<https://tutorcs.com>

Kafka Storage Structure



WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

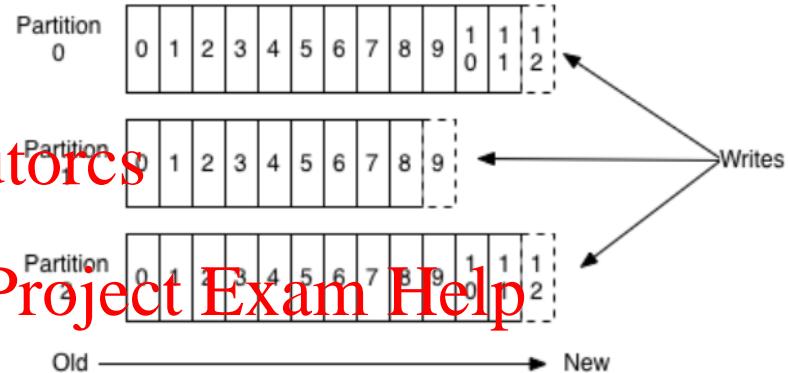
Messages are written to it in an append-only fashion

- Topic is logical grouping
- Partition is actual unit of data storage

QQ: 749389476

<https://tutorcs.com>

Anatomy of a Topic



What Kafka doesn't do? 程序代写代做 CS 编程辅导

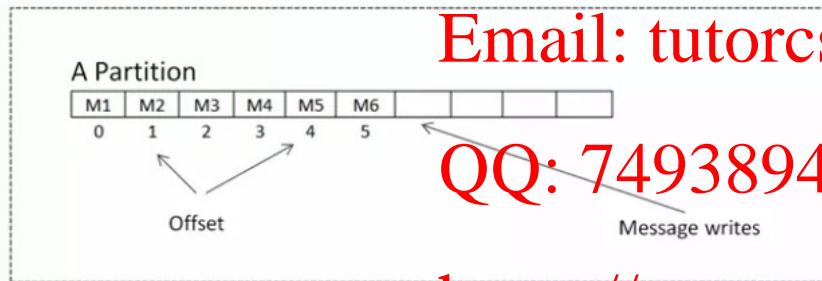
- Kafka does not have individual message IDs. Messages are simply addressed by their offset.
- Kafka also does not track consumers that a topic has or who has consumed what messages.
- There are no deletes. Kafka keeps all parts of the log for the specified time.



What is offset?

Assignment Project Exam Help

A sequence id given to messages as they arrive in a partition.



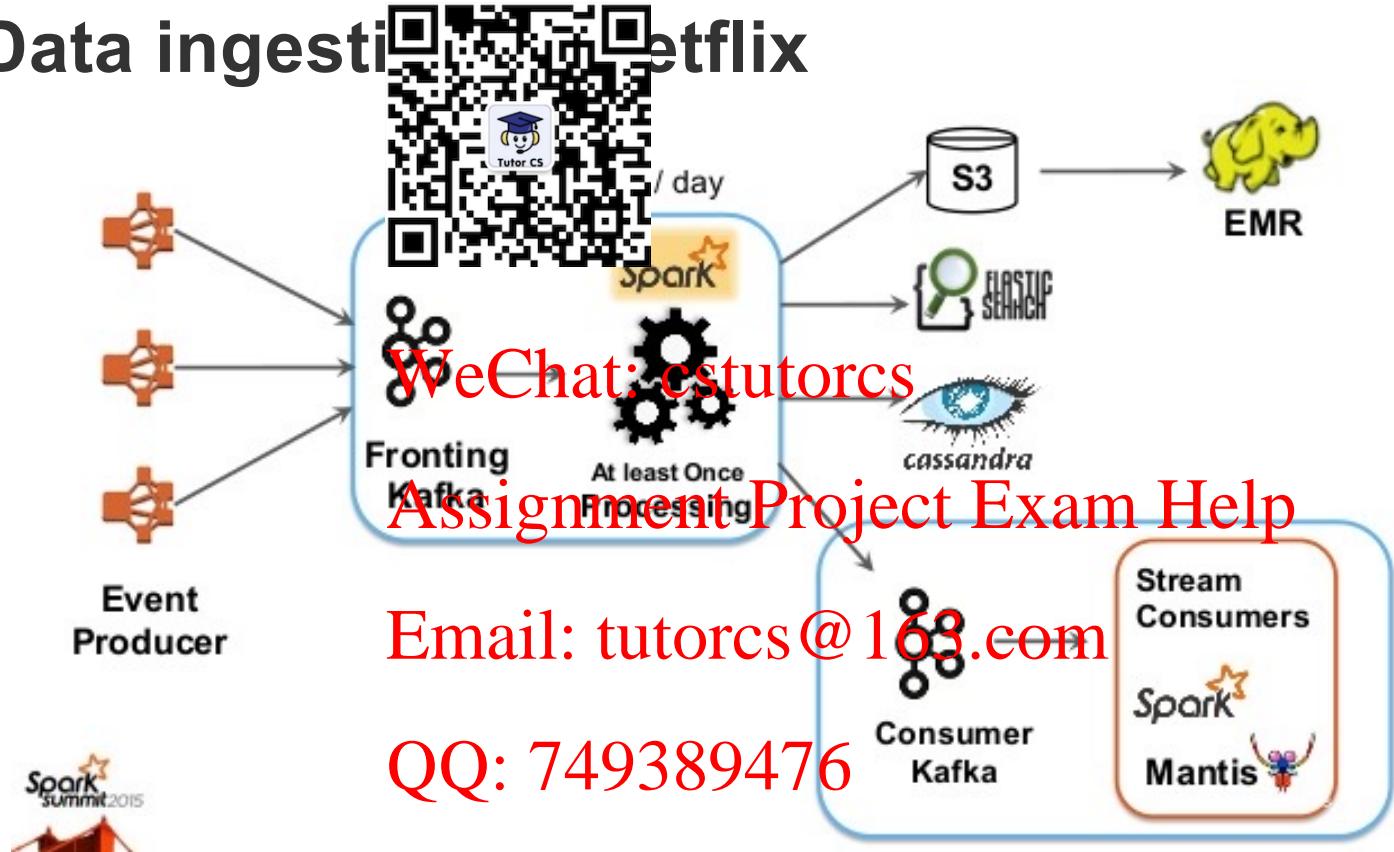
Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

Kafka and big data at web scale companies

- Big Data ingestion



<https://www.slideshare.net/SparkSummit/spark-and-spark-streaming-at-netflix-sedakar-daxini>

Kafka and big data at web scale companies

- BP OIL USE CASE



<https://www.linkedin.com/pulse/real-time-streaming-data-sirenilles-apache-kafka-spark-steven-murhula/>

Kafka and big data at web scale companies

- Transform, Store & Explore Healthcare Dataset



<https://mapr.com/blog/streaming-data-pipeline-transform-store-explore-healthcare-dataset-mapr-db/>

Should you use Apache Kafka?

- Kafka fits a class of system that a lot of web-scale companies and enterprises have, but just as the traditional message broker is not a one size fits all, neither is Kafka.
- If you're looking to build a set of resilient data services and applications, Kafka can serve as the source of truth by collecting and keeping all of the "facts" or "events" for a system.

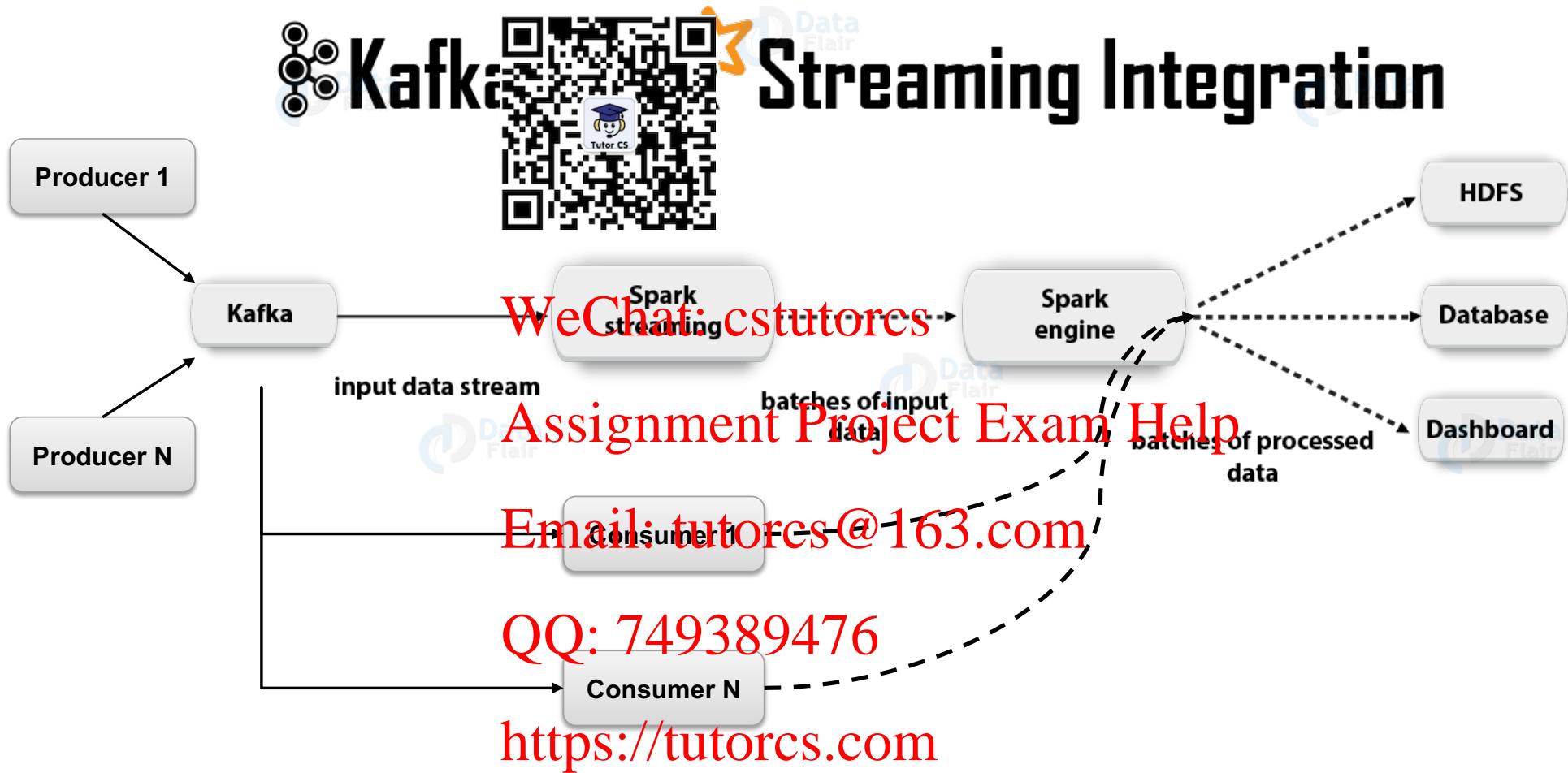


Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

Real-Time Streaming Architecture



DEMO

程序代写代做 CS编程辅导

- What is the total number of attendees up to t_x ?



Assume memory can only keep 5 tuples.

Query

ACC

WeChat: cstutores

Keep an accumulator for the sum. There is no need to keep the tuples.

Assignment Project Exam Help

∞	...	1	2	3	4	5	4	2	1
----------	-----	---	---	---	---	---	---	---	---

Email: tutorcs@163.com

t_n

QQ: 749389476 t_0

<https://tutorcs.com>

Thank You

程序代写代做 CS编程辅导

Questions?



WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>