



**BIRLA INSTITUTE OF TECHNOLOGY & SCIENCE, PILANI  
WORK INTEGRATED LEARNING PROGRAMMES**

**Digital**

**Part A: Content Design**

<b>Course Title</b>	STREAM PROCESSING AND ANALYTICS
<b>Course No(s)</b>	DSECL ZC556
<b>Credit Units</b>	5
<b>Credit Model</b>	
<b>Content Authors</b>	PRAVIN PAWAR

**Course Description**

Data is moving at very rapid space because of which necessarily of scalable systems capable of processing and analyzing this fast, streaming data has arisen. This course introduces the students with the architecture of streaming data processing systems. This course also enables students to understand the complete end-to-end solution for cost-effective analysis and visualization of streaming data with the help of various open source solutions available in this space. This course also helps students to learn the implementation and application of algorithms and data structures required for the streaming applications. Advanced streaming applications like Streaming SQL, Streaming Machine Learning will be discussed at proper length.

**Course Objectives**

<b>No</b>	
CO1	To introduce the applications of streaming data systems
CO2	To introduce the architecture of streaming data systems
CO3	To introduce the algorithmic techniques used in streaming data systems
CO4	To present survey of tools and techniques required for streaming data analytics

**Text Book(s)**

T1	Streaming Data: Understanding The Real-Time Pipeline, Andrew G.Psaltis, 2017, Manning Publications
T2	Real-Time Analytics: Techniques to Analyze and Visualize Streaming Data, Byron Ellis, 2014, Wiley

**Reference Book(s) & other resources**

R1	Big Data – Principles and best practices of scalable real-time data systems, Nathan Marz, James Warren, 2017, Manning Publications
R2	Designing Data Intensive Applications, Martin Kleppmann, O'Reilly



**Learning Outcomes:**

No	Learning Outcomes
LO1	Understand the components of streaming data systems with their capabilities and characteristics
LO2	Learn the relevant architecture and best practices for processing and analysis of streaming data
LO3	Gain knowledge about the development of system for data aggregation, delivery and storage using Open source tools
LO4	Get familiarity with the advance streaming applications like Streaming SQL, Streaming machine learning

**Part B: Learning Plan**

<b>Academic Term</b>	II Semester 2019 -2020
<b>Course Title</b>	STREAM PROCESSING AND ANALYTICS
<b>Course No</b>	DSECL ZC556
<b>Lead Instructor</b>	Prof. Maninder Singh Bawa

**Glossary of Terms**

<b>Module</b>	<b>M</b>	Module is a standalone quantum of designed content. A typical course is delivered using a string of modules. M2 means module 2.
<b>Contact Hour</b>	<b>CH</b>	Contact Hour (CH) stands for a hour long live session with students conducted either in a physical classroom or enabled through technology. In this model of instruction, instructor led sessions will be for 32 CH.
<b>Recorded Lecture</b>	<b>RL</b>	RL stands for Recorded Lecture or Recorded Lesson. It is presented to the student through an online portal. A given RL unfolds as a sequences of video segments interleaved with exercises.
<b>Lab Exercises</b>	<b>LE</b>	Lab exercises associated with various modules
<b>Self-Study</b>	<b>SS</b>	Specific content assigned for self study
<b>Homework</b>	<b>HW</b>	Specific problems/design/lab exercises assigned as homework

## Modular Structure

No.	Title of the Module
M1	Scalable Streaming Data Systems
M2	Streaming Data Systems Architecture
M3	Streaming Data Frameworks
M4	Streaming Analytics
M5	Advanced Streaming Applications

## Detailed Lecture Plan

### **M1: Scalable Streaming Data Systems**

#### **Session 1 to 3 / Contact Hour 1 - 6**

Time	Type	Description/Plan	Reference
Session 1	CH1	<ul style="list-style-type: none"> <li>Thinking about Data Systems</li> <li>Reliable, Scalable and Maintainable Data Applications</li> <li>Properties of Data</li> </ul>	R1 Ch1 R2 Ch2
	CH2	<ul style="list-style-type: none"> <li>Scaling with the traditional databases</li> <li>Big Data Systems</li> <li>Desired properties of Big Data Systems</li> </ul>	R2 Ch1
Session 2	CH3	<ul style="list-style-type: none"> <li>Data Model for Big Data</li> <li>Generalized Big Data System Architecture</li> </ul>	R2 Ch2 Class Notes
	CH4	<ul style="list-style-type: none"> <li>Real time systems</li> <li>Difference between Batch processing and Stream Processing</li> <li>Difference between real time and streaming systems</li> </ul>	T1 Ch1 Class Notes
Session 3	CH5	<ul style="list-style-type: none"> <li>Streaming Data Applications</li> <li>Databases and Streams</li> <li>Usage patterns of Streaming Data</li> </ul>	Class Notes R1 Ch11 Class Notes
	CH6	<ul style="list-style-type: none"> <li>Sources of Streaming Data</li> <li>Complex Event Processing Systems</li> </ul>	T2 Ch1 Class Notes
Post CH	SS	<ul style="list-style-type: none"> <li>Explore more on the non functional requirements of Data Intensive Applications               <ul style="list-style-type: none"> <li>✓ <a href="#">Non-functional Requirements for Real World Big Data Systems</a></li> <li>✓ <a href="#">IBM Big Data &amp; Analytics RA_VI</a></li> </ul> </li> <li>Explore more on the differences between the batch processing and streaming data applications               <ul style="list-style-type: none"> <li>✓ <a href="#">Batch vs Real time data processing</a></li> </ul> </li> </ul>	



		<ul style="list-style-type: none"> <li>Identify the use cases of Complex Event Processing Systems <ul style="list-style-type: none"> <li>✓ <a href="#">What is stream processing?</a></li> <li>✓ <a href="#">complex-event-processing</a></li> </ul> </li> </ul>
--	--	--

## M2: Streaming Data Systems Architecture

### Session 4 to 7 / Contact Hour 7 - 14

Time	Type	Description/Plan	Reference
Session 4	CH7	<ul style="list-style-type: none"> <li>Generalized Streaming Data Architecture</li> </ul>	T1 Ch 1 T1 Ch 2
	CH8	<ul style="list-style-type: none"> <li>Lambda Architecture</li> <li>Kappa Architecture</li> </ul>	Class Notes
Session 5-6	CH9	<ul style="list-style-type: none"> <li>Streaming Data system Component</li> <li>Features of Real time Architecture</li> <li>A real time architecture checklist</li> </ul>	T2 Ch2
	CH 10	<ul style="list-style-type: none"> <li>Service Configuration and Coordination Systems</li> <li>Maintaining the state</li> <li>Apache ZooKeeper</li> </ul>	T2 Ch3
	CH 11	<ul style="list-style-type: none"> <li>Data Flow Manager</li> <li>Managing distributed data flows</li> </ul>	T2 Ch4
	CH 12	<ul style="list-style-type: none"> <li>Apache Kafka</li> </ul>	T2 Ch4 Kafka Docs
Session 7-8	CH13	<ul style="list-style-type: none"> <li>Streaming Data Processor Concepts</li> <li>Timing Concepts</li> </ul>	T2 Ch 5 T1 Ch 5
	CH14	<ul style="list-style-type: none"> <li>Windowing</li> <li>Joins</li> </ul>	T1 Ch5 R1 Ch11
	CH15	<ul style="list-style-type: none"> <li>Storage for Streaming Data</li> <li>NoSQL storage Systems</li> <li>Choosing a Storage technology</li> </ul>	T2 Ch6
	CH16	<ul style="list-style-type: none"> <li>Delivery of Streaming Metrics</li> </ul>	T2 Ch7
Post CS	SS	<ul style="list-style-type: none"> <li>Explore in detail about issues with Lambda Architecture <ul style="list-style-type: none"> <li>✓ <a href="#">questioning-the-lambda-architecture</a></li> <li>✓ <a href="#">a-brief-introduction-to-two-data-processing-architectures</a></li> </ul> </li> <li>Explore the Java APIs exposed by following systems <ul style="list-style-type: none"> <li>✓ <a href="#">Apache ZooKeeper</a></li> </ul> </li> </ul>	



		<ul style="list-style-type: none"> <li>✓ <a href="#">Apache Kafka</a></li> <li>• Explore the data models of NoSQL data systems</li> <li>✓ <a href="#">MongoDB</a></li> <li>✓ <a href="#">Cassandra</a></li> </ul>	
--	--	---	--

### M3: Streaming Data Frameworks

#### Session 8 to 11 / Contact Hour 15 - 22

Time	Type	Description/Plan	Reference
Session 8	CH 15	<ul style="list-style-type: none"> <li>• Key features of Streaming Data Frameworks</li> <li>• Survey of Streaming Data Systems</li> </ul>	Class Notes
	CH 16	<ul style="list-style-type: none"> <li>• Apache Spark Streaming</li> </ul>	<a href="#">Spark Streaming Guide</a>
Session 9	CH 17	<ul style="list-style-type: none"> <li>• Apache Flink</li> <li>• Apache Samza</li> </ul>	<a href="#">Flink Docs</a> <a href="#">Samza Docs</a>
	CH 18	<ul style="list-style-type: none"> <li>• Apache Kafka Streaming</li> </ul>	<a href="#">Kafka Streaming Guide</a>
Session 10	CH 19	<ul style="list-style-type: none"> <li>• Apache Storm Architecture</li> </ul>	<a href="#">Storm Docs</a>
	CH 20	<ul style="list-style-type: none"> <li>• Apache Storm Concepts</li> <li>• Apache Storm Groupings</li> </ul>	T2 Ch 5
Session 11	CH 21	<ul style="list-style-type: none"> <li>• Apache Storm Running Example</li> </ul>	<a href="#">Storm Docs</a>
	CH 22	<ul style="list-style-type: none"> <li>• Storm – Kafka Integration Example</li> </ul>	Class Notes
Post CH	SS	<ul style="list-style-type: none"> <li>• Compare the different streaming data platforms and identify the use cases for which they are suitable</li> <li>• Implement the streaming data pipeline using the Kafka Streaming library</li> <li>• Implement a streaming data application with Spark streaming</li> </ul>	<a href="#">Kafka Streaming Guide</a>  <a href="#">Spark Streaming Guide</a>

## M4: Streaming Analytics

### Session 12 to 13 / Contact Hour 23 - 26

Time	Type	Description/Plan	Reference
Session 12	CH 23	<ul style="list-style-type: none"> <li>Exact Aggregation of Streaming Data</li> <li>Time Series Analysis</li> </ul>	T2 Ch 8
	CH 24	<ul style="list-style-type: none"> <li>Quantization Framework</li> <li>Stochastic Optimization</li> </ul>	T2 Ch8
Session 13	CH 25	<ul style="list-style-type: none"> <li>Registers and Hash Functions</li> <li>The Bloom Filter</li> </ul>	T2 Ch 10
	CH 26	<ul style="list-style-type: none"> <li>Distinct Value Sketches</li> <li>The Count-Min Sketch</li> </ul>	T2 Ch 10
Post CH	SS	<ul style="list-style-type: none"> <li>Study illustrations for Streaming data concepts</li> <li>Explore algorithms for aggregation of streaming data</li> <li>Explore more about the streaming data processing algorithms for exact results</li> </ul>	Class Notes

## M5: Advanced Streaming Applications

### Session 14 to 15 / Contact Hour 27 - 30

Time	Type	Description/Plan	Reference
Session 14	CH25	<ul style="list-style-type: none"> <li>Necessity of Streaming SQL</li> <li>Streaming SQL : Windows</li> <li>Streaming SQL : Joins</li> <li>Streaming SQL : Patterns</li> </ul>	<a href="#">Streaming SQL Blog</a>
	CH26	<ul style="list-style-type: none"> <li>Apache Storm support for Streaming SQL</li> <li>Apache Flink support for Streaming SQL</li> <li>Streaming SQL for Apache Kafka</li> </ul>	<a href="#">storm-sql</a> <a href="#">flink-stream-sql</a> <a href="#">Kafka Streaming SQL</a>
Session 15	CH27	<ul style="list-style-type: none"> <li>Models for Streaming Data - Linear models</li> <li>Models for Streaming Data - Logistic Regression models</li> </ul>	T2 Ch 11
	CH 28	<ul style="list-style-type: none"> <li>Forecasting with Models - Exponential Smoothing methods</li> <li>Forecasting with Models - Regression methods</li> </ul>	T2 Ch 11
Session 15	CH 29	<ul style="list-style-type: none"> <li>Streaming ML Frameworks I</li> </ul>	<a href="#">structured-streaming-ml</a>
	CH 30	<ul style="list-style-type: none"> <li>Streaming ML Frameworks II</li> </ul>	



Post CH	SS	<ul style="list-style-type: none"> <li>Get familiarized with Streaming SQL tools <ul style="list-style-type: none"> <li>✓ <a href="#">storm-sql</a></li> <li>✓ <a href="#">Kafka Streaming SQL</a></li> </ul> </li> <li>Build and deploy machine learning models using Spark structured streaming <ul style="list-style-type: none"> <li>✓ <a href="#">structured-streaming-ml</a></li> </ul> </li> </ul>	
---------	----	---	--

### **Session 16 / Contact Hour 31 - 32**

Time	Type	Description/Plan	Reference
Session 16	CH31	<ul style="list-style-type: none"> <li>Review of Streaming Data Systems and Architectures</li> </ul>	CH 1 to 16
	CH32	<ul style="list-style-type: none"> <li>Review of Streaming Data Techniques and Applications</li> </ul>	CH 17 to 32

### **Evaluation Scheme:**

Legend: EC = Evaluation Component; AN = After Noon Session; FN = Fore Noon Session

No	Name	Type	Duration	Weight	Day, Date, Session, Time
EC-1	Assignment-1	Take-home, Programming and use of platforms	-	10%	TBD
	Assignment-2		-	15%	TBD
	Quiz-1	Online	30 mins	5	TBD
EC-2	Mid-Semester Test	Closed Book	2 hours	30%	TBD
EC-3	Comprehensive Exam	Open Book	3 hours	40%	TBD

### **Notes:**

Syllabus for Mid-Semester Test (Closed Book): Topics in Session Nos. 1 to 8 (contact hours 1 to 16)

Syllabus for Comprehensive Exam (Open Book): All topics

### **Important links and information:**

Elearn portal: <https://elearn.bits-pilani.ac.in>

Students are expected to visit the Elearn portal on a regular basis and stay up to date with the latest announcements and deadlines.

Contact sessions: Students should attend the online lectures as per the schedule provided on the Elearn portal.

### **Evaluation Guidelines:**

- EC-1 consists of either two Assignments or three Quizzes. Students will attempt them through the course pages on the Elearn portal. Announcements will be made on the portal, in a timely manner.
- For Closed Book tests: No books or reference material of any kind will be permitted.
- For Open Book exams: Use of books and any printed / written reference material (filed or bound) is permitted. However, loose sheets of paper will not be allowed. Use of calculators is permitted in all exams. Laptops/Mobiles of any kind are not allowed. Exchange of any material is not allowed.
- If a student is unable to appear for the Regular Test/Exam due to genuine exigencies, the student should follow the procedure to apply for the Make-Up Test/Exam which will be made available on the Elearn portal. The Make-Up Test/Exam will be conducted only at selected exam centres on the dates to be announced later.

It shall be the responsibility of the individual student to be regular in maintaining the self



study schedule as given in the course handout, attend the online lectures, and take all the prescribed evaluation components such as Assignment/Quiz, Mid-Semester Test and Comprehensive Exam according to the evaluation scheme provided in the handout.