

## Review of the Paper on “Discretized Streams: Fault-Tolerant Streaming Computation at Scale”

**Sahithya Arveti Nagaraju (Person number: 50559752, UBITName: sarvetin)**

The paper introduces Discretized Streams (D-Streams), an innovative model specifically designed for scalable and fault-tolerant streaming data processing. D-Streams, as a more improved version of the Spark framework, maintain the benefits of both batch and stream processing, as they relate to each other through a representation of a stream as a successive strategy of deterministic batch computations on small time intervals. This approach allows D-Streams to maintain high throughput, one-second latency, and strong fault regulations.

### **Key Contributions**

The authors present the fact that distributed stream processing faces several major challenges such as fault tolerance, straggler handling, consistency maintenance, as well as the seamless integration of batch and stream processing. They introduce D-Streams to overcome these problems through Resilient Distributed Datasets (RDDs), which can perform lineage-based recomputation and thus be fault-tolerant without having to be replicated elsewhere or making any disk I/O. This design makes recovery from failure possible without incurring a large overhead that is typical for classical methods like copying or upstream backup.

The paper demonstrates how standard streaming operations, like stateful computation and sliding windows, can be implemented in the D-Stream model. By combining streaming, batch, and interactive queries in one framework, the authors simplify the programming model and offer a high-level functional API. Experiments prove that Spark Streaming, which is the implementation of D-Streams, is capable of processing more than 60 million records per second with sub-second latency on a 100-node cluster, hence, demonstrating the model's scalability and efficiency.

### **Strengths**

1. **Fault Tolerance:** The adoption of RDDs for lineage-based recomputation is a major milestone from the traditional methods in terms of reducing the cost of fault recovery and increasing scalability.
2. **Framework:** D-Streams integrate the batch and streaming computation thus it is no longer necessary the separate systems and reducing the programming complexity for those applications.

3. **High Performance:** The experimental results show that the system achieves better throughput and fault recovery times than existing stream processing systems like Storm and S4.

#### **Limitations:**

1. **Latency Constraints:** Millisecond latency is quite acceptable for most use cases, however, could be unsuitable for applications that require fast computations (millisecond-scale).
2. **Checkpointing Overhead:** The new habit of periodical checkpointing for stateful computations brings the complication of higher frequency of updates, for example, to update things.
3. **Applicability to Complex Workflows:** The batch model, which is non-deterministic, the consequence of it may be in the situations when you have very dynamic and unpredictable workloads.

#### **Conclusion**

The document very strongly supports that D-Streams are a solid and scalable model for distributed streaming processing.

By combining the robustness of batch systems with the low-latency requirements of streaming applications, D-Streams is a very useful tool for processing real-time big data. Despite the existence of some drawbacks, like the one-Time delay issue of formal data monitoring and the dependency on backups, the highlighted advantages of scalability, fault tolerance, and user-friendliness clearly dominate.