

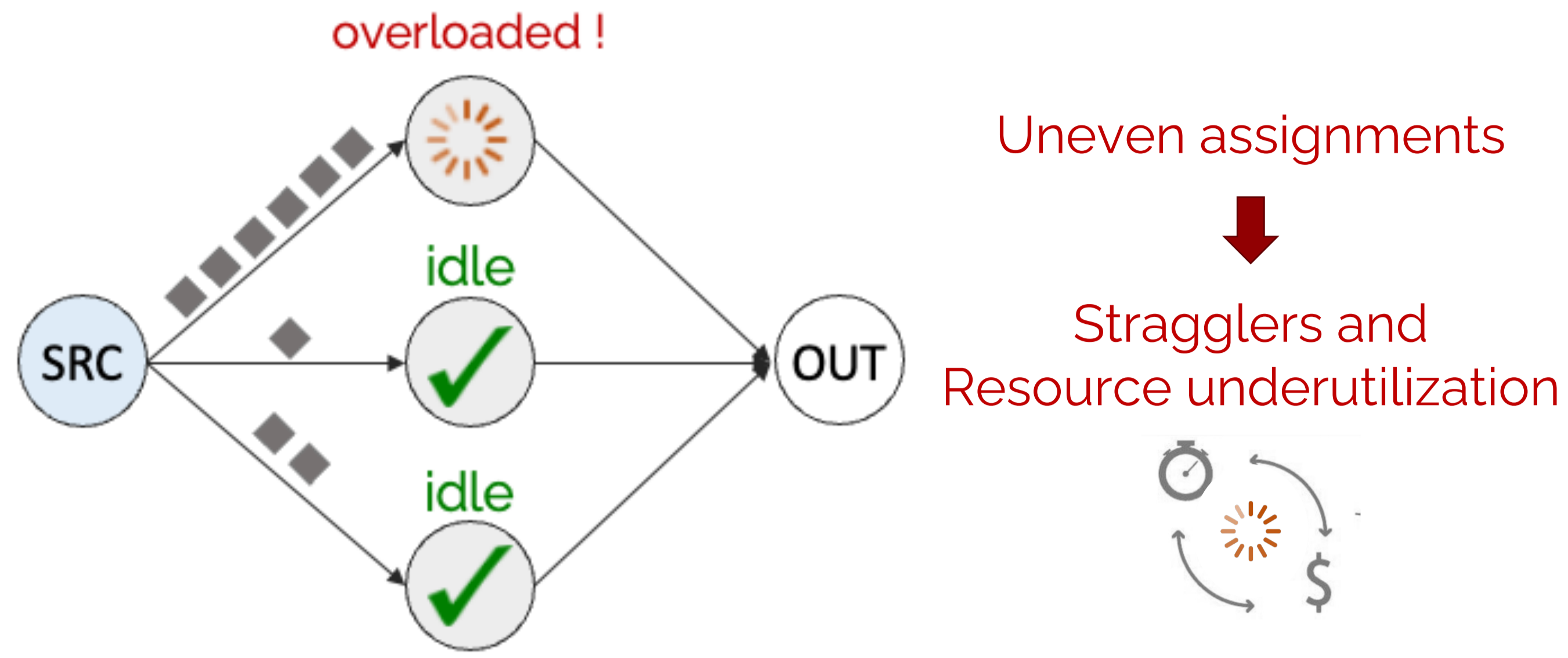


# Dalton: Learned Partitioning for Distributed Data Streams

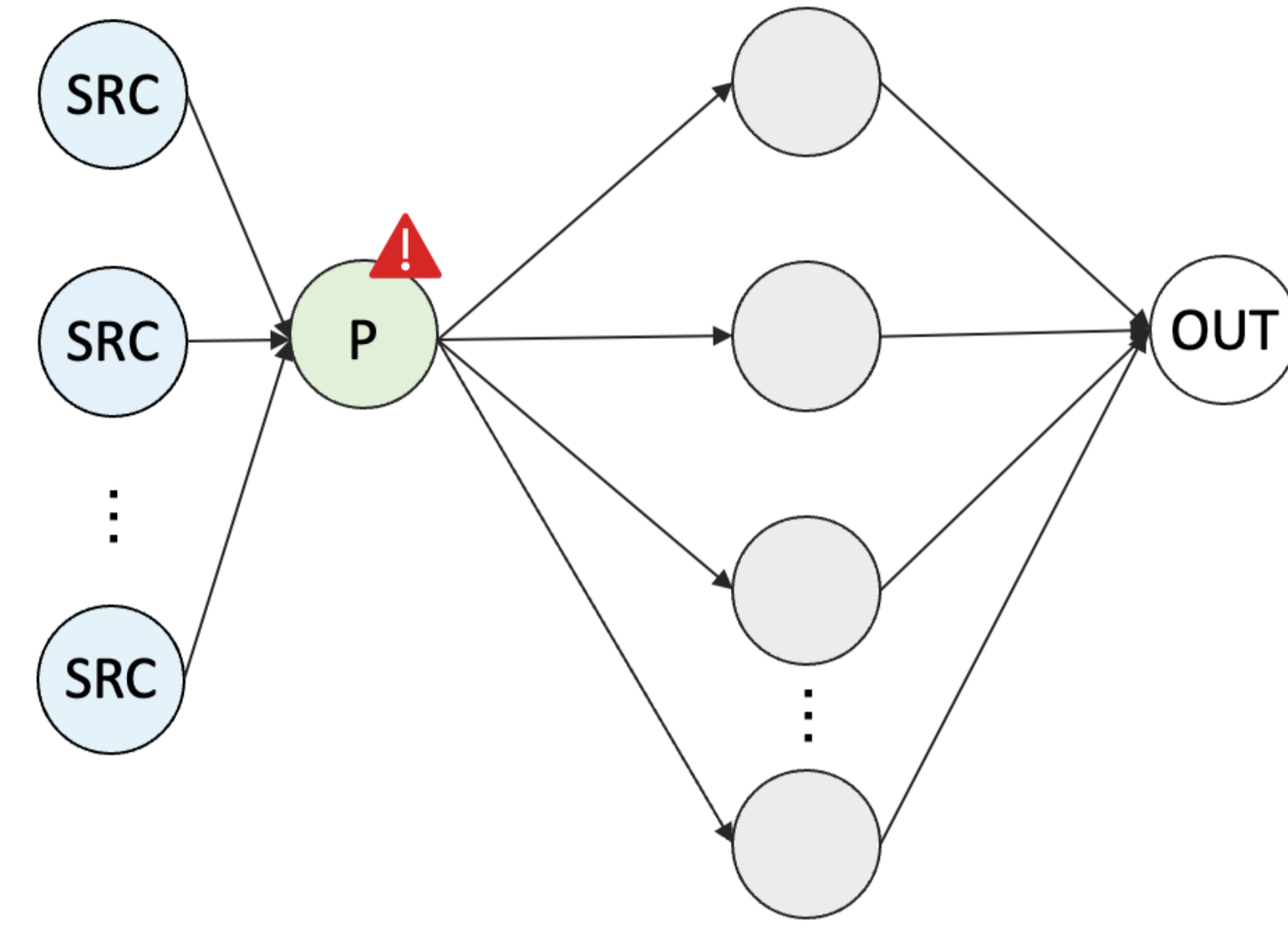
Eleni Zapridou, Ioannis Mytilinis, Anastasia Ailamaki  
firstname.lastname@epfl.ch

## 1. Partitioning must adapt to the workload

## 2. Partitioning must be scalable

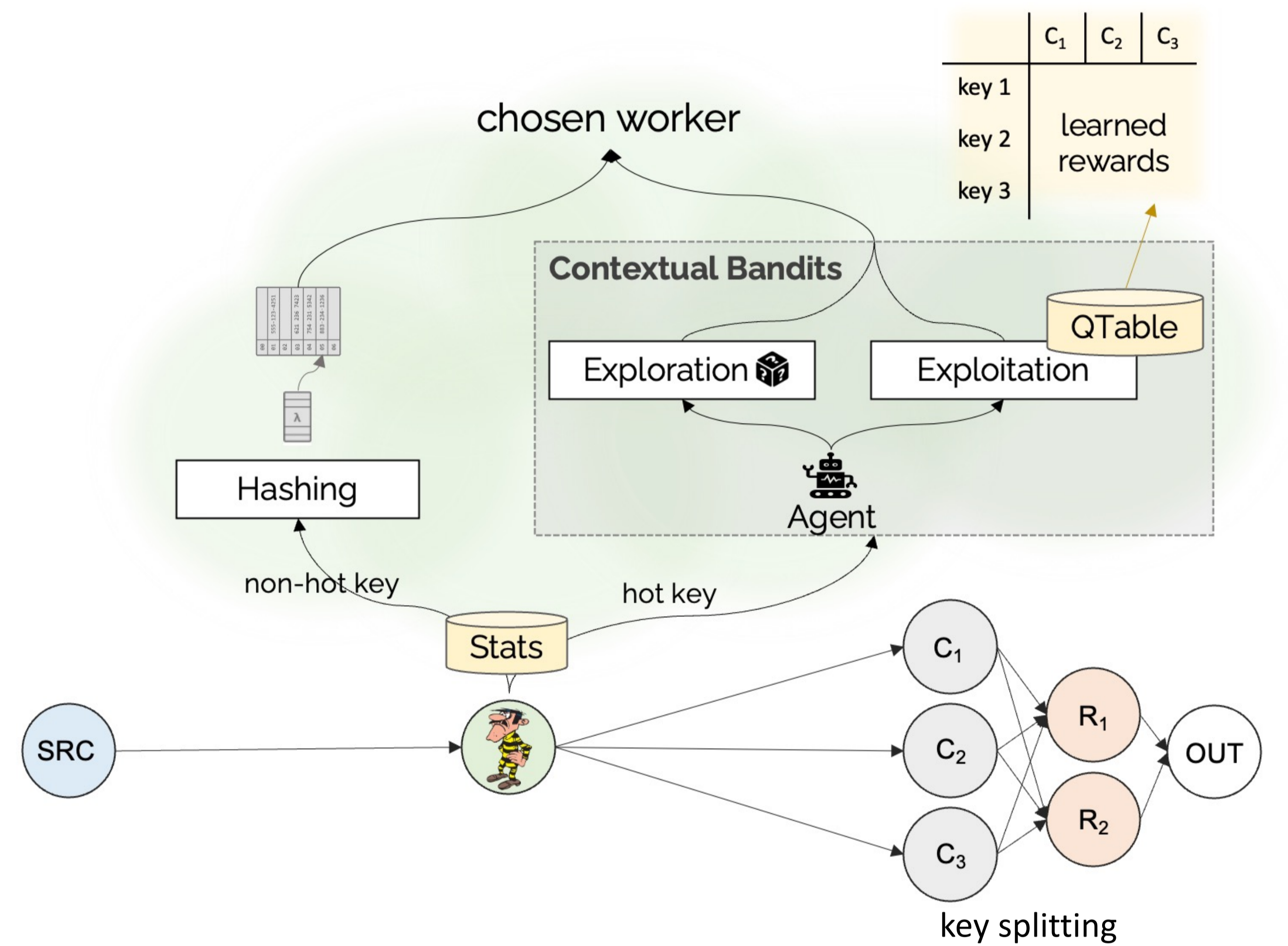
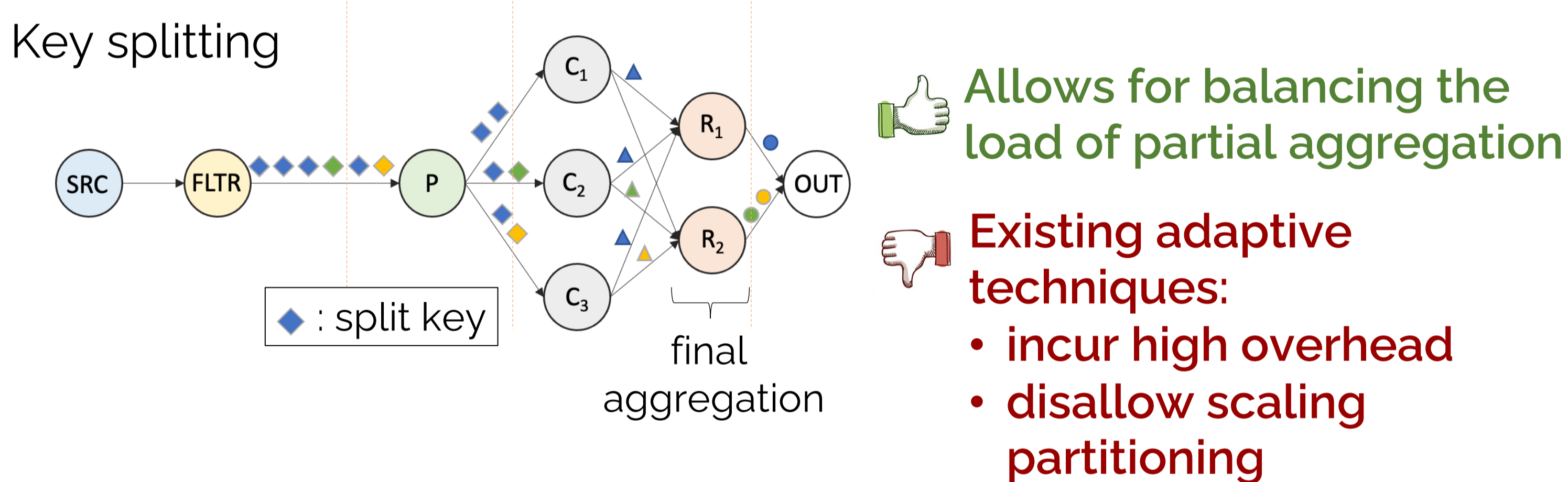
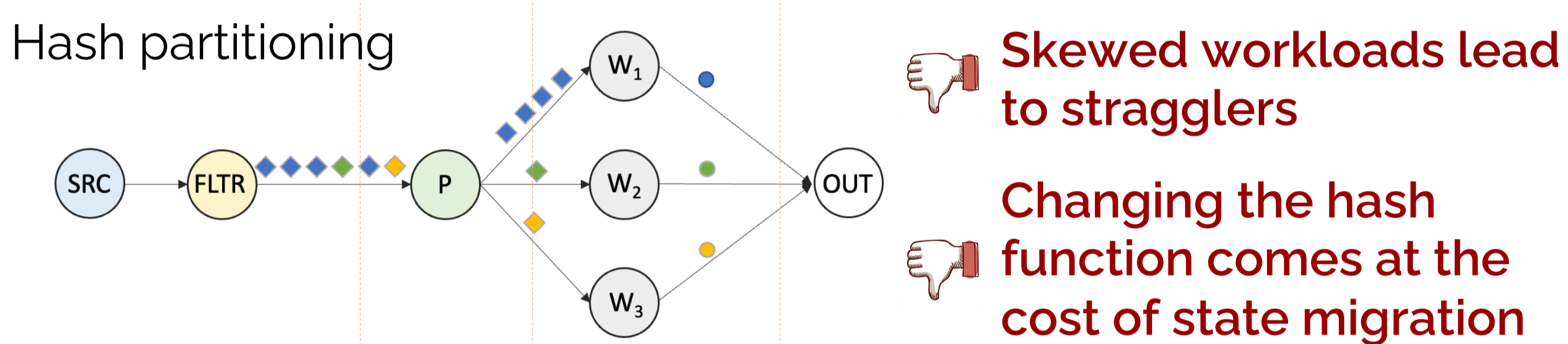
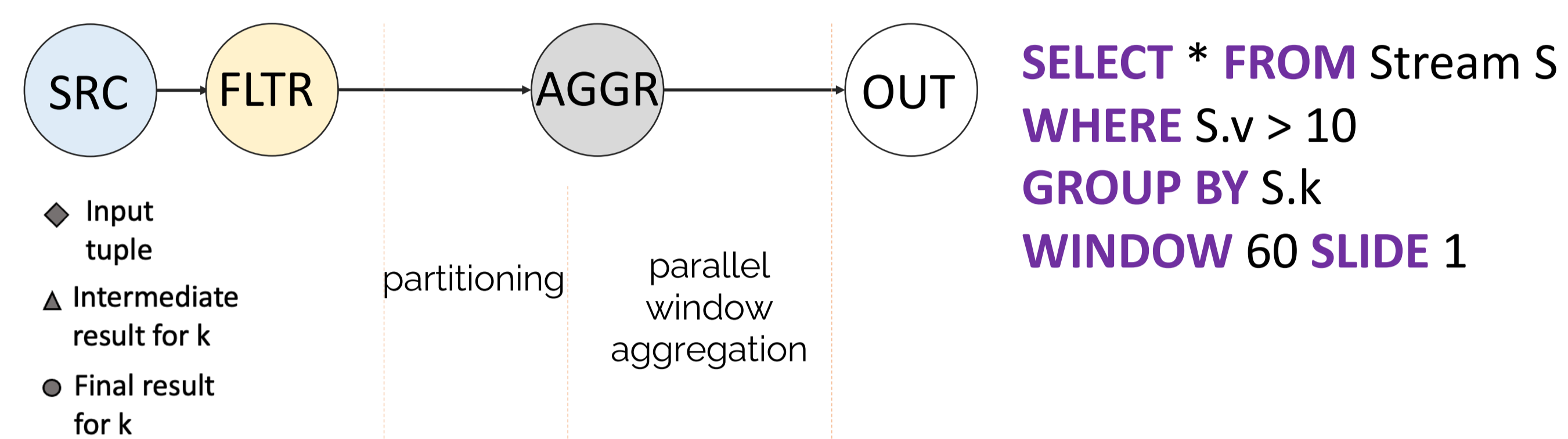


The distribution of data streams changes at runtime



## 3. Partitioning: How it is currently done

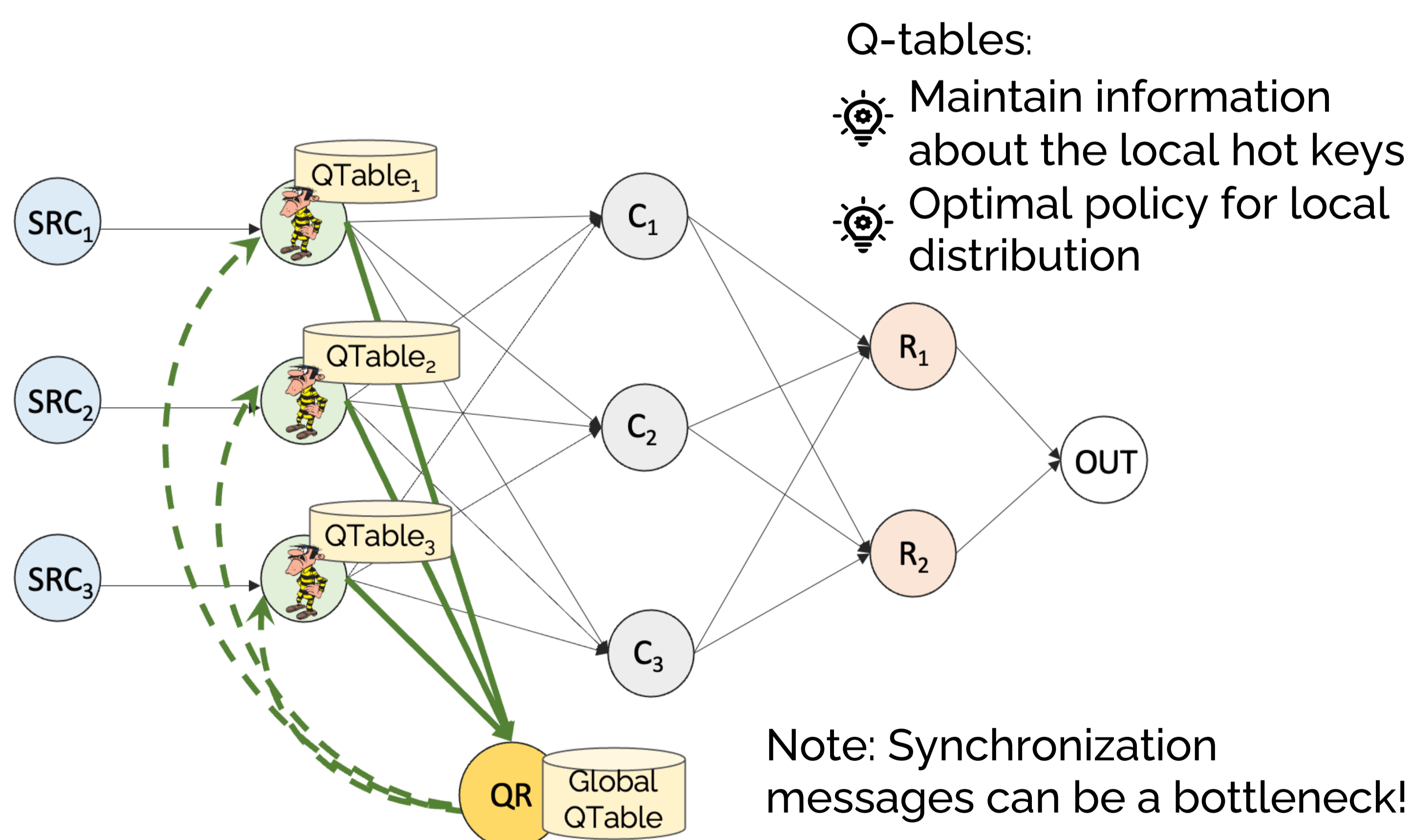
## 4. Dalton adapts partitioning at runtime



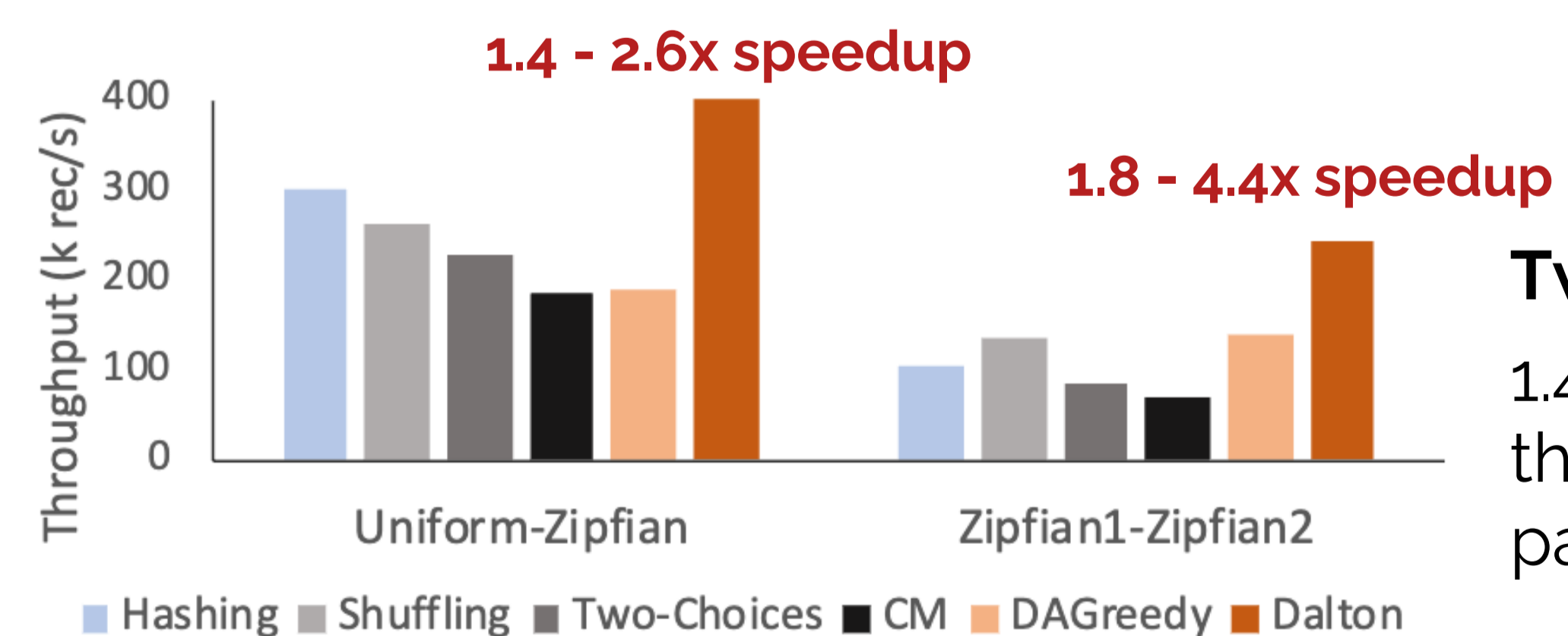
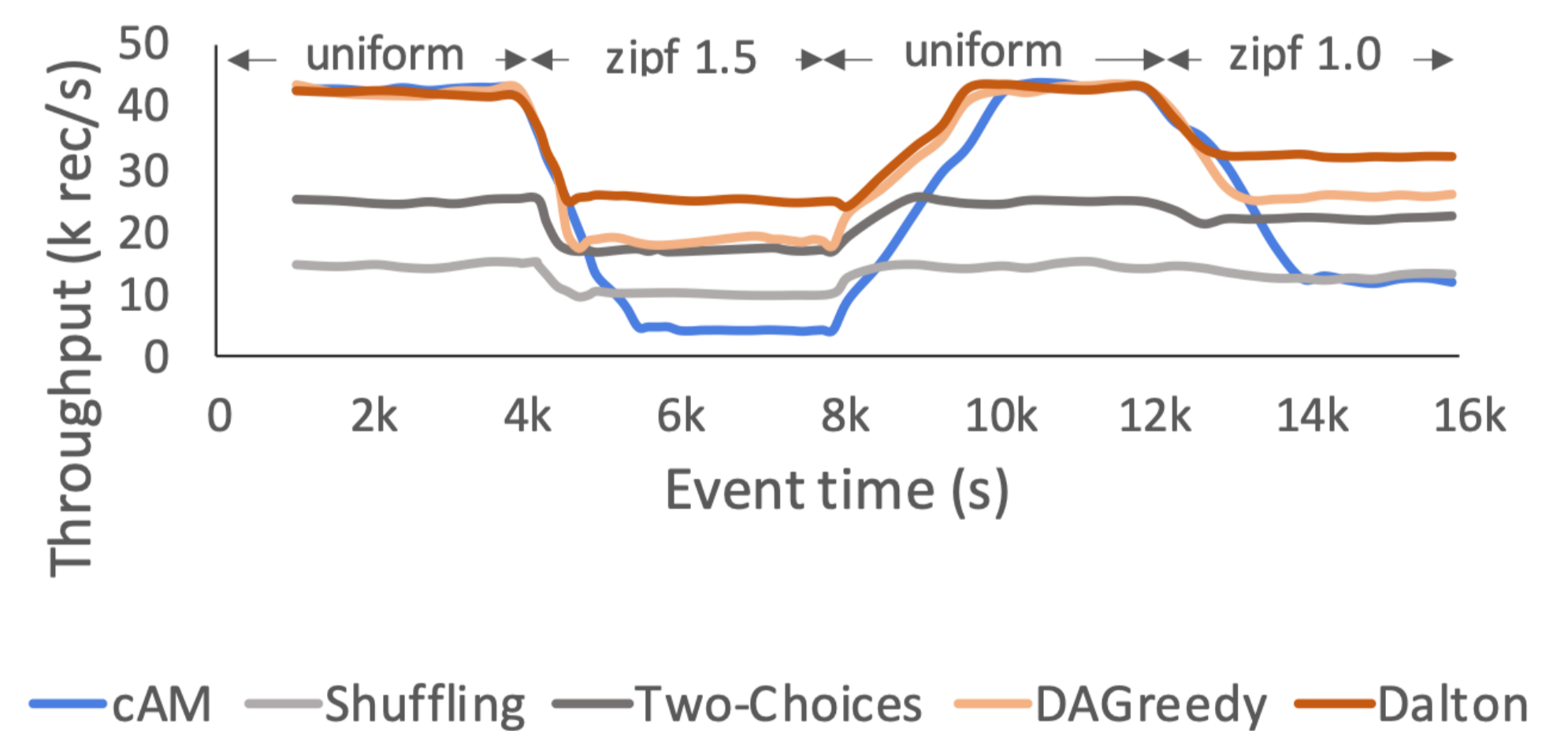
- Rewards computed by a cost model that balances partial and final aggregation
- Continuously learn rewards
- Exploitation: leverage acquired experience
- Exploration: is more splitting beneficial?

## 5. Dalton scales to many partitioners

## 6. Dalton maximizes throughput



**Dynamic workload**  
1.3-6.3x higher throughput when the data distribution is skewed



**Dalton is the only algorithm that adapts to the data distribution and scales to multiple instances**

## 7. Conclusion

- Dalton
- learns partitioning policies at runtime with minimal overhead
  - quickly adapts to the distribution and is able to scale not only the processing workers but also the partitioners
  - outperforms the state-of-the-art by a factor of 1.4-6.3x