

Move Fast and Meet Deadlines: Fine-grained Real-time Stream Processing with Cameo

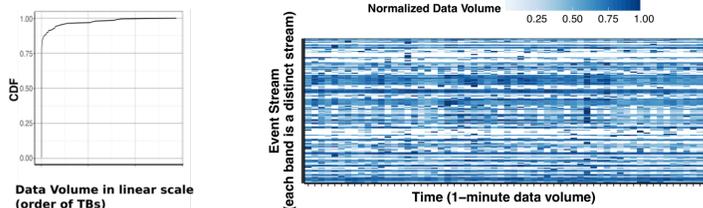
Le Xu[‡], Shivaram Venkataraman[‡], Indranil Gupta[‡], Luo Mai[¶], and Rahul Potharaju[°]

[‡]University of Illinois at Urbana Champaign, [¶]University of Wisconsin-Madison, [¶]University of Edinburgh, [°]Microsoft



Resource provisioning for data stream applications is challenging because:

- ❖ Long tail streams cause resource over-provisioning
- ❖ Temporal variation makes resource prediction difficult
- ❖ Latency requirements vary across jobs

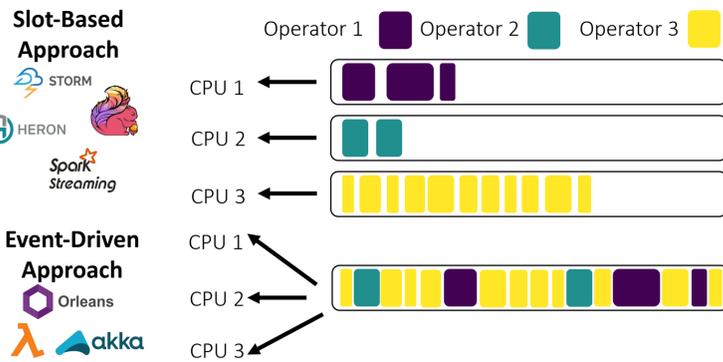


Data volume varied by streams Data volume varied by time

Cameo's Goal: Real-time, fine-grained resource multiplexing for multi-tenant data stream applications. Cameo's contributions include:

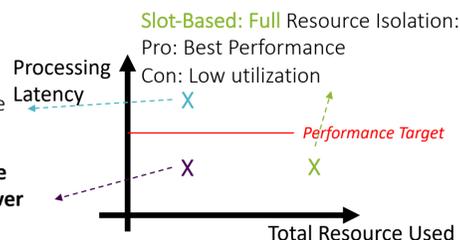
- ❖ Deriving operator priorities dynamically by analyzing pending messages.
- ❖ Scheduling mechanism and programmable interface that enables plugged in scheduling policy.
- ❖ Scheduling at message-level granularity with low overhead.

WHY FINE-GRAINED SCHEDULING?



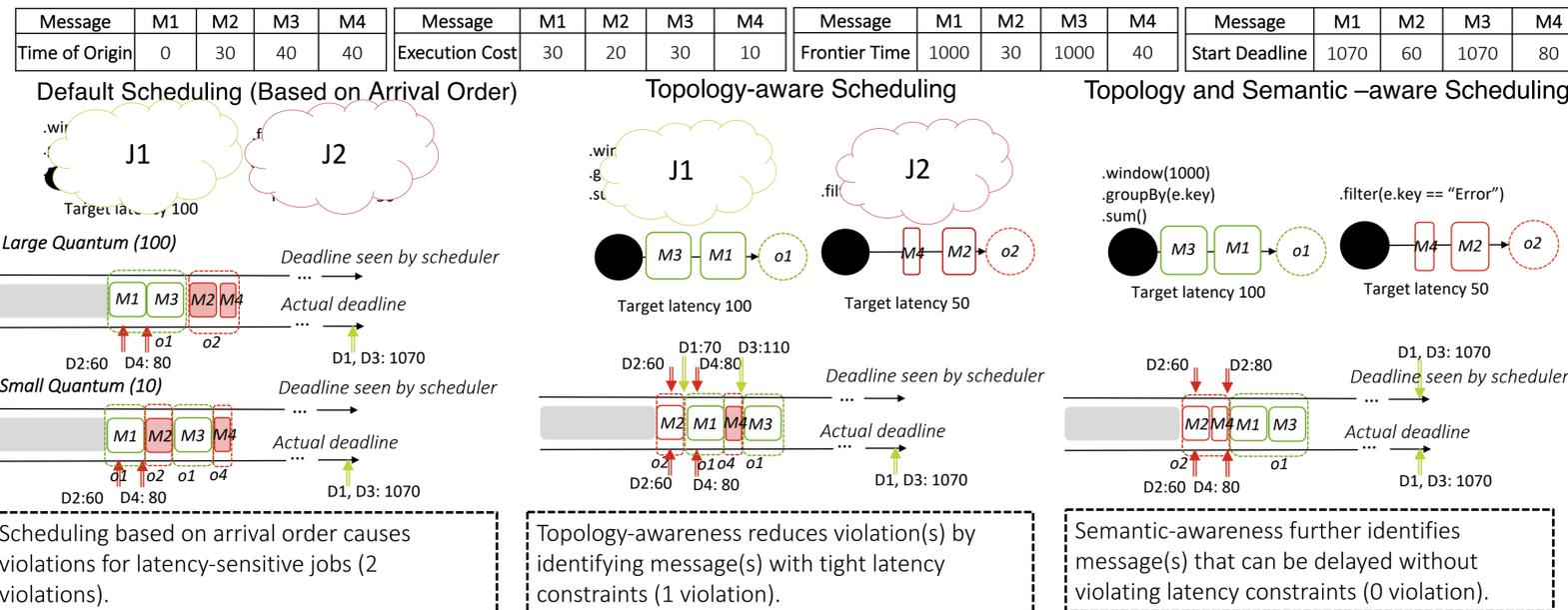
- ❖ Today's SPEs mostly adopt slot-based approach that binds operators to CPUs.
 - ❖ Actual resource requirement is difficult to predict
 - ❖ Requires to provision for peak load to avoid performance degradation
- ❖ Event-driven framework achieves high utilization by sharing work among CPUs.
 - ❖ Difficult to enforce performance isolation without application awareness

Event-Driven: No Resource Isolation
 Pro: High Utilization
 Con: Compromised Performance



Cameo: Providing performance guarantee without resource over provisioning

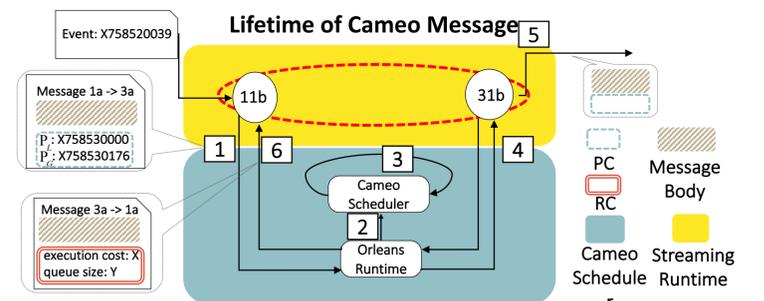
SCHEDULING WITH APPLICATION AWARENESS



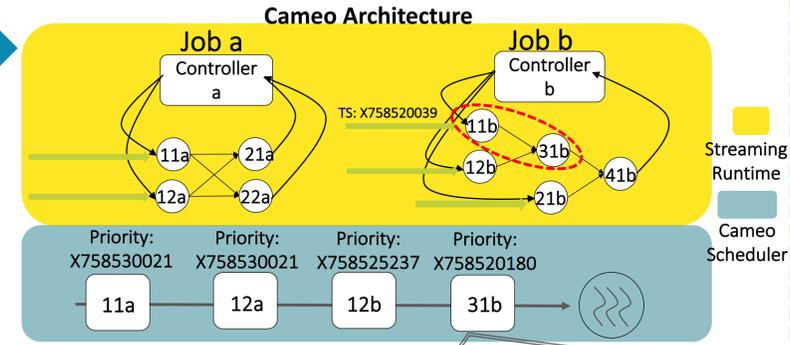
PROGRAMMABLE SCHEDULING INTERFACE

Cameo enables data-driven, plug-and-play scheduling policy by using both **Scheduling Contexts** and user-defined **Context Handlers**:

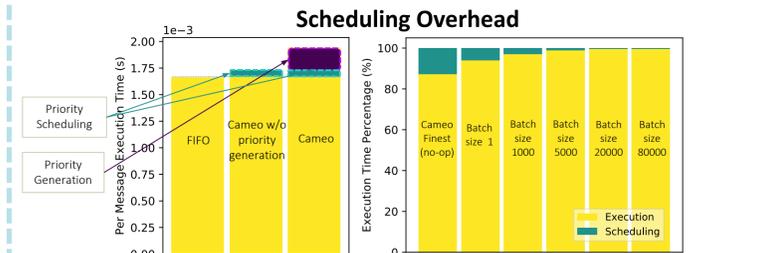
- ❖ Handler functions are invoked **before** and **after** processing every request and acknowledgement message.
- ❖ Handler functions enforces priority generation strategy by **creating**, **interpreting** and **modifying** scheduling contexts.
- ❖ Message scheduler determines message and operator priority by observing scheduling context attached to messages. Message scheduler is policy-agnostic and therefore remains stateless.



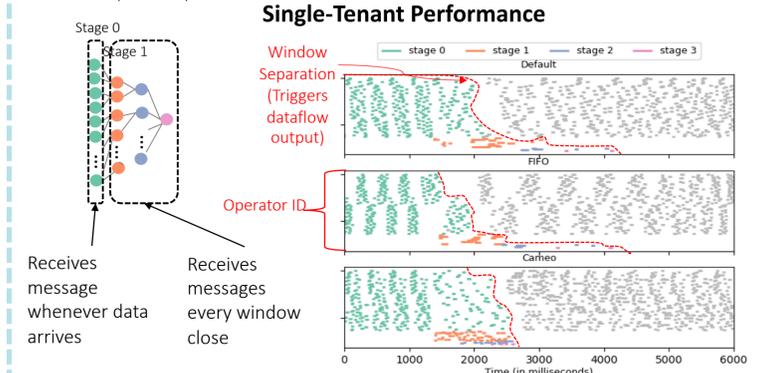
1. PC generation (**BuildCtxAtSource**)
2. Request insertion
3. Request scheduled
4. Request message execute (**BuildCtxAtOperator**)
5. Trigger target operator and reply (**PrepareReply**)
6. Reply message execute (**ProcessCtxFromReply**)



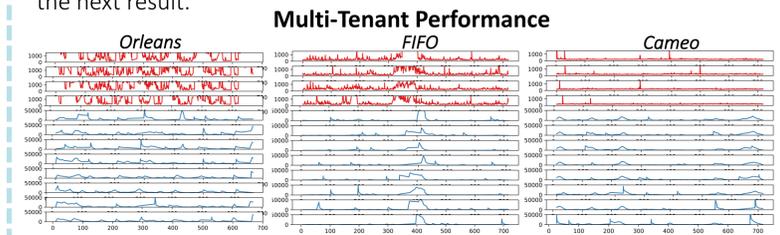
CAMEO IN ACTION



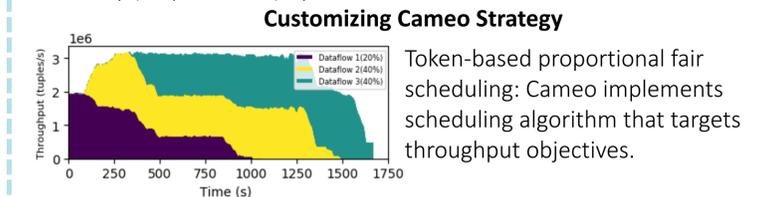
No-op message: 4% overhead from priority-based scheduling and 11% from priority generation. In practice, scheduling overhead is much smaller (<6.4%).



Cameo improves query latency (2.7 - 3.2X) by preventing messages influencing future results to be processed before messages targeting the next result.



Pareto Workload: Cameo weathers transient spikes by reducing tail latency (99 percentile) by 21 - 30X.



Token-based proportional fair scheduling: Cameo implements scheduling algorithm that targets throughput objectives.

