# StreamScope:
# Continuous Reliable Distributed Processing of Big Data Streams

*Wei Lin*, Haochuan Fan*, <u>Zhengping Qian (ZP)</u>*, Junwei Xu, Sen Yang, Jingren Zhou*, Lidong Zhou*

Microsoft
*Now with Alibaba Group

@NSDI'16

A new transaction got reflected in the output within 3s
The system processed up to 50 million events/s

State:
Completion:
Run Time:
Useful PN Hours: 14888:29:40.446
Bonus PN Hours: 43.45%

Runtime Name: scopecep_hcfan_201411
Submitted By: PHX\yuyao
Submit Time: 8/14/2014 8:18:45 PM
Compilation Time: 35 seconds
Queued Time: 1 seconds
Start Time: 8/14/2014 8:19:20 PM
End Time:
Yielded Time:

Cluster: cosmos11-prod-cy2
VC: adCenter.BICore
Priority: 800
Tokens: 482
Allocation(%): 11

Root Process Id: c3b436b7-292f-4a9f-8e5
Root Process Node: cy2sch030020747

Bytes Read: 61,336,252,895,109
Bytes Left: 249,538,700,649
Bytes Written: 60,124,233,285,838

Total Nodes: 2,876
-Completed:
-Running:
-Failed: 35

Job Diagnostics: [Diagnose]

Alert(s):

ⓘ Data Skew: 0 shallow issue(s) detected.

Investigate

Job Details:

Script | Algebra | VertexDef | Code
| Resources | Debug Stream |

Display: Progress  ▣ Succeeded  ▣ Failed  ▣ Running  ▣ Waiting

**High complexity**
48 stages
18 joins of 5 different types
21.3 TB in-memory state

**Massive scalability**
Reads 61TB + Write 61TB
7 billions of input events
6 billions of output events
3000+ long-running tasks

**Fault tolerance**
Handles both planned failures and
unplanned outages automatically

Play

# Streaming dataflow



Input streams

Streaming events

R

R

X

X

X

Vertices

Channels

M

Output streams

# Streaming dataflow



Input streams

Replay of upstream events

R

R

Rebuild the state

X

X

X

Vertices

Missing or duplicate events

M

Output streams

# rStream



- *Properties*
  - *There is a unique value associated with each sequence number*
  - *A read returns only after a successful write, for the same* `seq`
  - *If a write of* `(seq,e)` *succeeds, then for the following reads that reach position* `seq`, *they eventually return* `(seq,e)`

# Execution of a vertex

3,4,5,6,7

$X: t_1$

1

$s_1 = <\{2\}, \{1\}, t_1>$

# rVertex



$s_1 = \langle \{2\}, \{1\}, t_1 \rangle$    $s_2 = \langle \{3\}, \{2\}, t_2 \rangle$    $s_3 = \langle \{4\}, \{4\}, t_3 \rangle$

Restart from a snapshot

# Failure recovery



$s_1 = <\{2\}, \{1\}, t_1>$         $s_2 = <\{3\}, \{2\}, t_2>$         $s_3 = <\{4\}, \{4\}, t_3>$

# Optimization

- *Naïve implementation of rStream: writing events to reliable store*
  - *Synchronous writes introduce significant latencies*
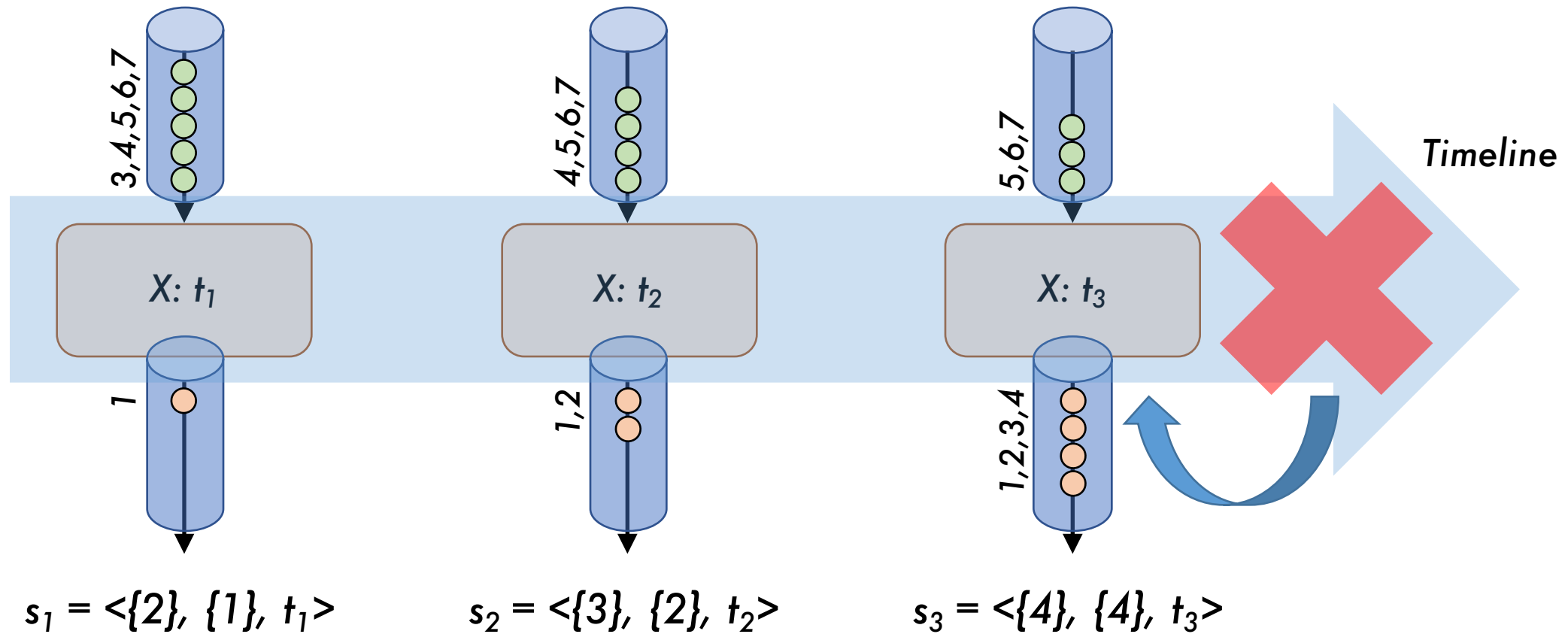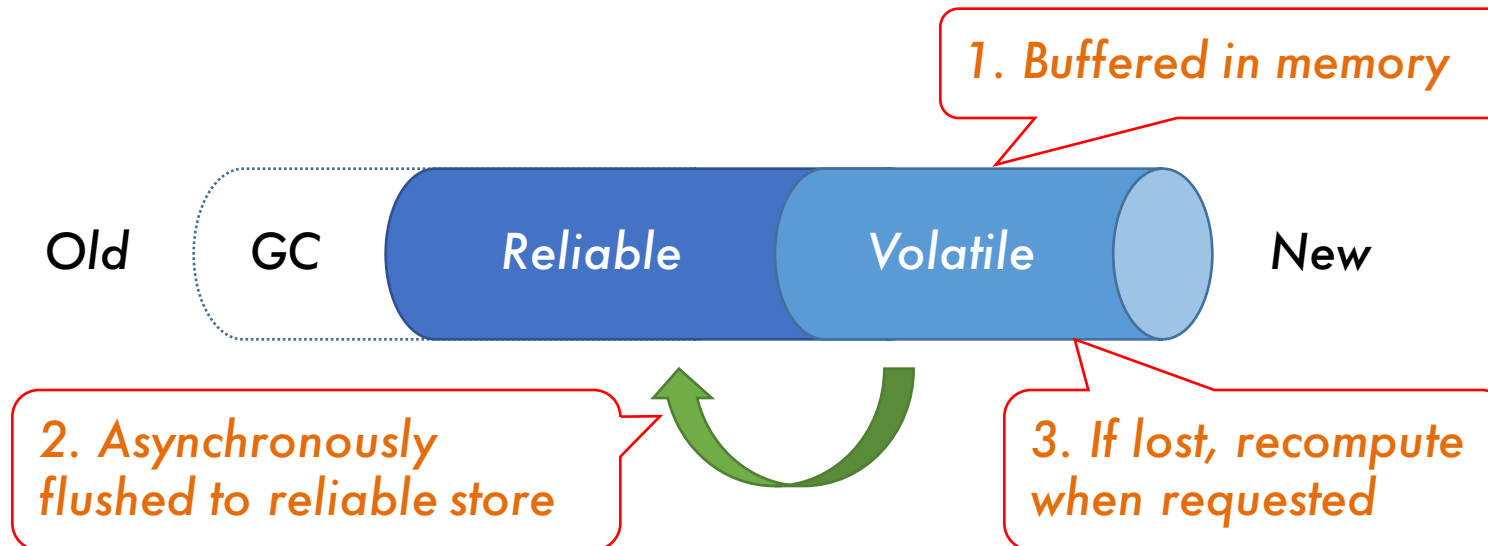
- *Uses a hybrid scheme that moves writes out of the critical path while providing the illusion of reliable channels*

*1. Buffered in memory*

Old    GC    Reliable    Volatile    New

*2. Asynchronously flushed to reliable store*

*3. If lost, recompute when requested*

# Different failure recovery strategies

- Recomputation using dependency tracking at runtime
- Checkpoint/log replay
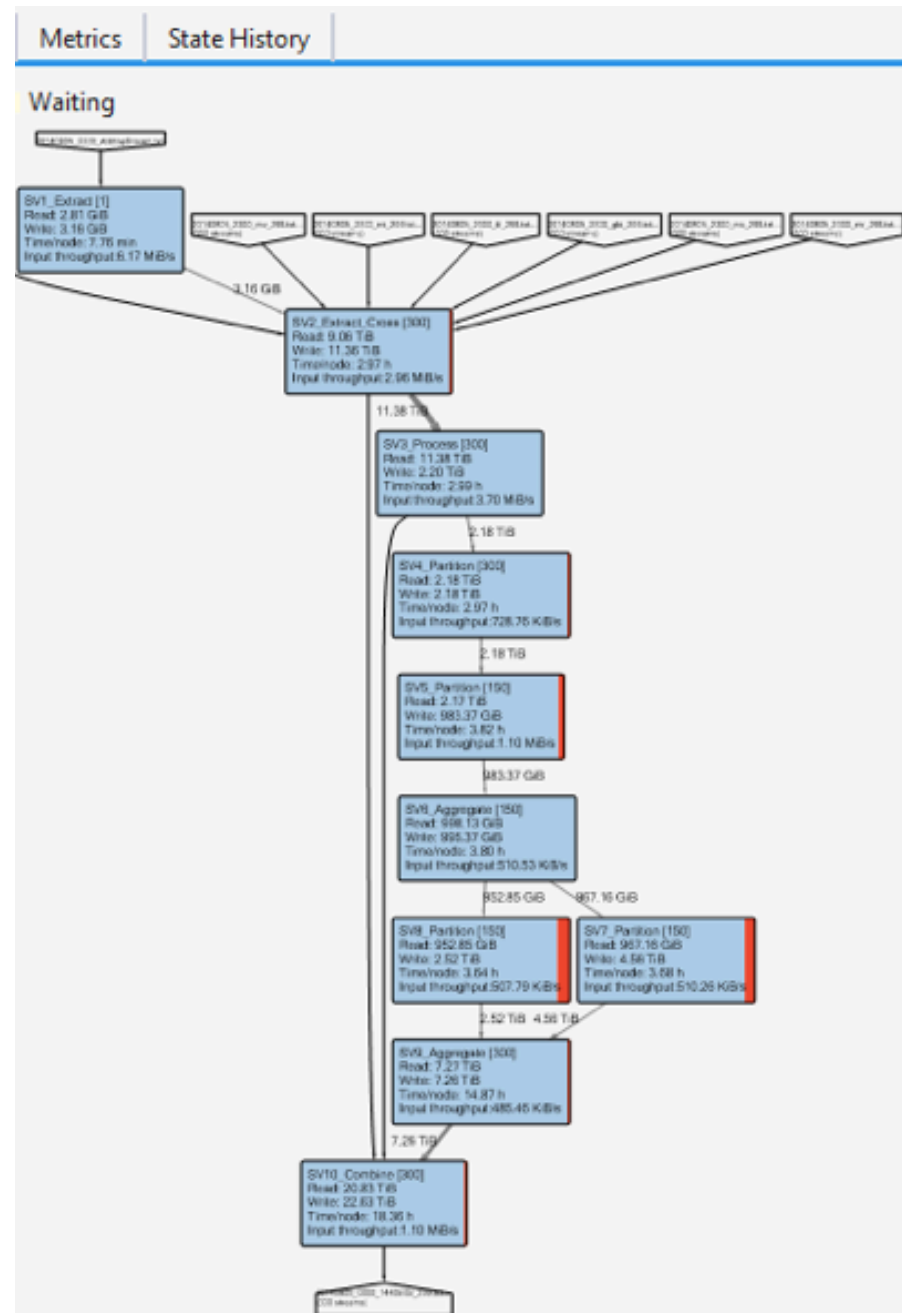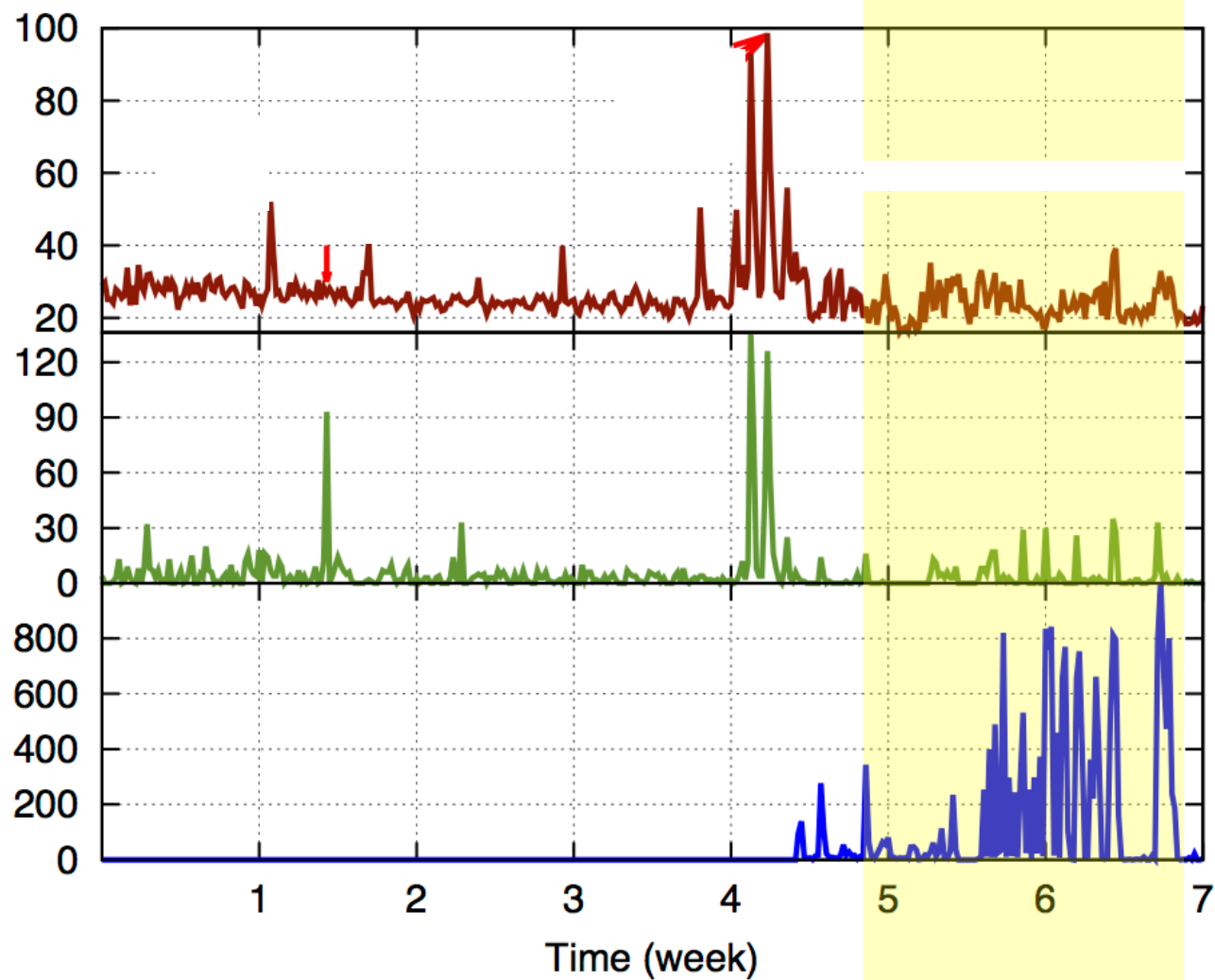- Persistent state/streams
- Hybrid

# Development/debugging

- *Greatly leveraged and tightly integrated with existing system*
  - *Integrated language, optimizer, scheduling, etc.*

- *Distributed streaming made easy*
  - *Off-line mode: starting with finite inputs with minimum resources to validate/debug a streaming application*
  - *Later switched to on-line, live execution transparently*
  - *Greatly improves developer productivity in lifecycle of an application*
    - *E.g., Can even debug/profile a vertex without impacting the running job*

# Deployment

- *Re-examination of segments of execution in the past for auditing*

- *Dynamic scaling and robustness to load fluctuation*

- *Continuous operation during system maintenance*

- *Straggler handling*

- *Dynamic reconfiguration/patching to resolve data anomalies*

# Conclusion

- *Cloud-scale stream computation is challenging due to the complexity of dependencies*

- *StreamScope introduces two new abstractions, rVertex and rStream, to manage the complexity through decoupling*

- *The abstractions separate system properties from the actual implementation to,*
  - *Enable powerful optimizations*
  - *Develop different failure recovery strategies*
  - *Better support the lifecycle of streaming applications in production*