# Malleable Flow for Time-Bounded Replica Consistency Control
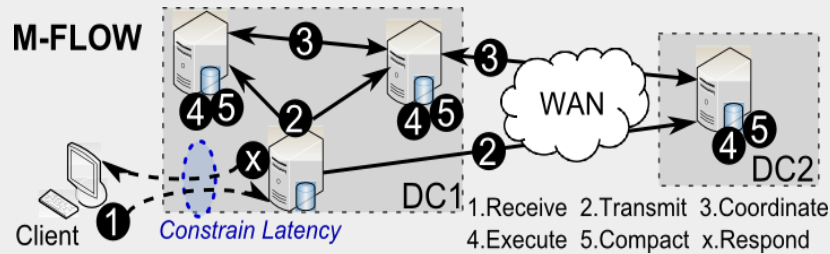
Yuqing Zhu[#1], Jianmin Wang[#2], Philip S. Yu[*#3]

[#]Tsinghua University, School of Software

[*]University of Illinois at Chicago, Dept. of Computer Science

[1]zhuyuqing@gmail.com, [2]jimwang@tsinghua.edu.cn, [3]psyu@uic.edu
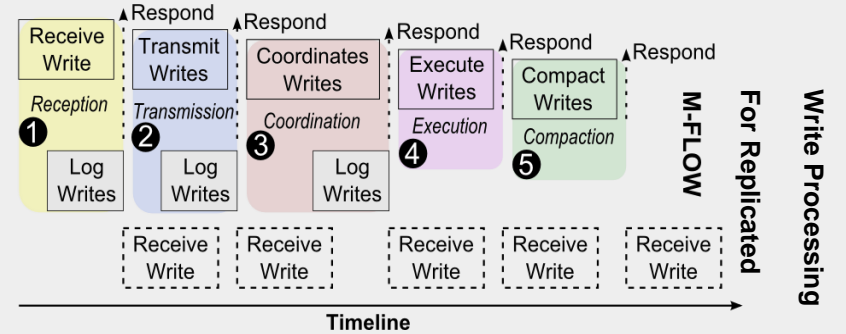
THSS-ISE Lab     Tsinghua University

## Big Picture





A system with malleable flow (M-flow) supports latency-bounded operations that maximize replica consistency within the given time. A better consistency status indicates a shorter latency for a following consistent read, or a more recent value returned by a following read within the same latency.

The M-flow replication process starts from when the request is received till when the response is sent. It is decomposed into six stages, i.e. reception, transmission, coordination, execution, compaction and acquisition.

The M-flow replication strategy allows update-anywhere, eager synchronous and lazy asynchronous replication simultaneously. The decomposed replication process enables this flexibility and the control of replica consistency (and latency) by reforming a suitable execution process with carefully selected stages and writes.
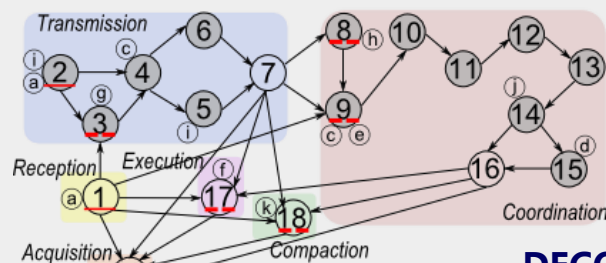
## Motivation

- Replication is important for availability, scalability and reliability in large-scale systems
- Trade-off between replica consistency and availability, latency exists
- Guaranteeing best replica consistency within a given latency
  o Not provided yet
  o But highly DESIRABLE

## Key Idea

- To decompose the replication process into stoppable stages and then into a directed graph of ordered minute steps
- To reform an execution flow by finding a path of steps in the graph that
  o Meets the time requirement
  o Maximizes replica consistency status
  o Guarantees durability and fault tolerance

## Bounded Time & Traded Consistency



Bounded Latencies (the given latency bounds versus the measured 99 percentile latencies):

- M-flow can control the general trend on achieving the latency versus consistency trade-offs.



Consistency versus latency under different cross-DC bandwidths (the number of returned values versus the given read latency constraints):

- A larger latency bound for instantaneous reads following writes leads to a larger number of returned values.

## The Malleable Flow (M-FLOW)



The decomposed steps and their execution orders form a directed graph, consisting of six partitions corresponding to the six stages.
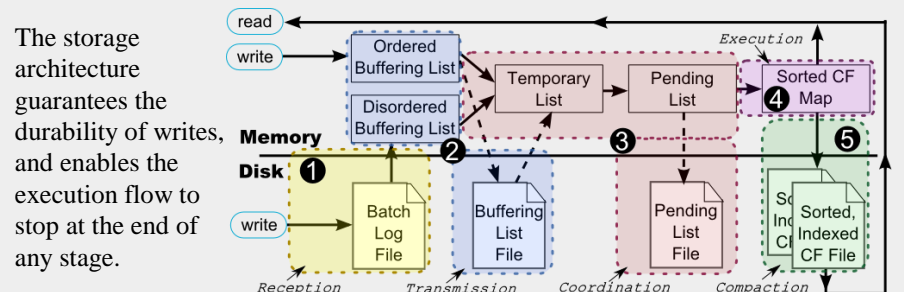
### DECOMPOSE / REFORM

Principles for reformation:

- The path either covers or bypasses a stage

- The step subset is maximized with execution time bounded

- The set of writes processed within a step is maximized

  o Computing from Stage execution, coordination, transmission, to Stage compaction

  o Increasing executed and executable writes to improve the replica consistency status

- The path for a write covers Stage reception; that for a read covers Stage acquisition

**Algorithm : dCON Algorithm(finding the legal path with maximized consistency)**

```
Input: t_P, S_P
Output: A path of nodes S_P
1  Initiate all i_s to 0;
2  If ① then                    // Executable writes
3      Compute S_P's path cost C_P;
4      Estimate the number m_e of executable writes within t_P − C_P;
5      if m_e > 0 then
6          Set i_12 ← MIN(m_e, p-list write number);
7          Add nodes 17 to S_P;              Execution
8  If ⓒ ∧ ⓔ ∧ ① then            // Coordinatable writes
9      Compute S_P's path cost C_P;
10     Estimate the number m_c of writes coordinatable and executable
       within t_P − C_P;
11     if m_c > 0 then
12         Set i_7 ← MIN(1/2 m_c, b-list size+b-file write number);
13         Set i_12 ← i_12 + m_c;
14         if i_7 > b-list size then Add node 8 to S_P;
15         Add nodes 9-14,16 to S_P;         Coordination
16 If ⓒ ∧ ⓖ then               // Transmittable writes
17     Compute S_P's path cost C_P;
18     Estimate the number m_t of writes transmittable within t_P − C_P;
19     if m_t > 0 then
20         Set i_4 ← i_3 ← MIN(m_t, bl-file write number);
21         Add nodes 3-5, 7 to S_P;          Transmission
22 If ⓚ then                   // Compactable writes
23     Compute S_P's path cost C_P;
24     Estimate the number m_z of compactions within t_P − C_P;
25     if m_z > 0 then
26         Set j_1, k_1 based on m_z;
27         Add nodes 18 to S_P;              Compaction
```
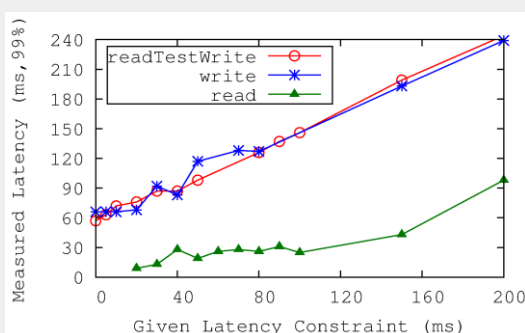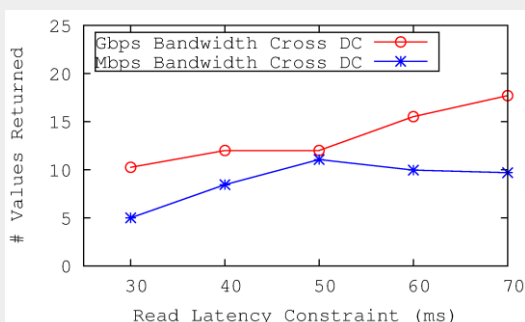
## Implementation over Cassandra

The storage architecture guarantees the durability of writes, and enables the execution flow to stop at the end of any stage.