



Low End-to-End Latency atop a Speculative Shared Log with Fix-Ante Ordering



Shreesha G.

Bhat



Tony Hong









Ram

Aishwarya Ganesan Alagappan

University of Illinois Urbana-Champaign

Xuhao Luo

Jiyu Hu

Shared logs are widely used by today's real-time, data-driven applications



Shared logs are widely used by today's real-time, data-driven applications

Well studied: Corfu, Scalog, Boki, LazyLog...



Shared logs are widely used by today's real-time, data-driven applications

Well studied: Corfu, Scalog, Boki, LazyLog...

State-of-the-art shared logs incur high delivery latencies



Shared logs are widely used by today's real-time, data-driven applications

Well studied: Corfu, Scalog, Boki, LazyLog...

State-of-the-art shared logs incur high delivery latencies

high delivery latency \rightarrow high end-to-end latency



Shared logs are widely used by today's real-time, data-driven applications

Well studied: Corfu, Scalog, Boki, LazyLog...

State-of-the-art shared logs incur high delivery latencies

high delivery latency \rightarrow high end-to-end latency

Cause: expensive ordering before record delivery



We note an opportunity for speculative execution

We note an opportunity for speculative execution

SpecLog Abstraction

Allow speculative consumption of records by predicting order Overlap ordering and application compute Confirm the order after compute

We note an opportunity for speculative execution

SpecLog Abstraction

Allow speculative consumption of records by predicting order Overlap ordering and application compute Confirm the order after compute

Fix-Ante Ordering

Predetermine the global order for easy prediction Make system adhere to that order

High speculation success!

Contributions

- SpecLog A new speculative shared log abstraction, uses fix-ante ordering to enable near perfect speculation
- Belfast Implementation of SpecLog abstraction and fix-ante ordering
 - Addresses practical challenges that arise in implementing fix-ante ordering
 - Enables low e2e latency while retaining benefits and guarantees of today's shared logs
- E2E latency benefits in end applications like fraud detection, intrusion detection and high-frequency trading

Outline

- Introduction
- Motivation
- SpecLog Abstraction and Interface
- Belfast An Implementation of SpecLog
- Evaluation

Expose simple interface



Expose simple interface

pos_t append(record_t r); record_t read(pos_t pos); // streaming records stream<record_t> subscribe(pos_t start);

Durable & linearizably ordered records

• Respect real-time order of appends

Expose simple interface

Durable & linearizably ordered records

• Respect real-time order of appends



Expose simple interface

Durable & linearizably ordered records

• Respect real-time order of appends



Expose simple interface

Durable & linearizably ordered records

• Respect real-time order of appends



Expose simple interface

Durable & linearizably ordered records

• Respect real-time order of appends

Store records over multiple storage shards



Expose simple interface

Durable & linearizably ordered records

• Respect real-time order of appends

Store records over multiple storage shards

• Total order of records across shards



Durability first – appended records first made durable



Durability first – appended records first made durable



Durability first – appended records first made durable Batched ordering – amortize cost of ordering through batching



Durability first – appended records first made durable Batched ordering – amortize cost of ordering through batching



Durability first – appended records first made durable Batched ordering – amortize cost of ordering through batching



Durability first – appended records first made durable Batched ordering – amortize cost of ordering through batching



Flexible Data Placement Seamlessly Reconfigurable Scalablity

Durability first – appended records first made durable Batched ordering – amortize cost of ordering through batching



Flexible Data PlacementSeamlessly ReconfigurableScalablitySimilar designs adopted by Boki [SOSP'21], FlexLog [HPDC'23]









Problem: High Delivery Latency

Note:

Shard reports <L> imply a local log length of L at the shard

Sequencing Layer



Problem: High Delivery Latency

Note:

Shard reports <L> imply a local log length of L at the shard

Sequencing Layer






















Fraud detection: consume txn \rightarrow match w/ history (compute) \rightarrow flag txn

Latency Demands of Real-Time Applications

Latency Demands of Real-Time Applications

Many applications care about e2e latency

• Real-time analytics, high-frequency trading, fraud and intrusion detection systems

These apps demand a low delivery latency from the shared log

Latency Demands of Real-Time Applications

Many applications care about e2e latency

• Real-time analytics, high-frequency trading, fraud and intrusion detection systems

These apps demand a low delivery latency from the shared log

2023 streaming report by RedPanda

 \sim 35% practitioners primarily care about delivery latency of their streaming system



Outline

- Introduction
- Motivation
- SpecLog Abstraction and Interface
- Belfast An Implementation of SpecLog
- Evaluation



SpecLog identifies opportunity for speculative execution





SpecLog identifies opportunity for speculative execution



SpecLog identifies opportunity for speculative execution



SpecLog identifies opportunity for speculative execution



SpecLog identifies opportunity for speculative execution



SpecLog identifies opportunity for speculative execution



SpecLog identifies opportunity for speculative execution



Abstraction

Interface

Abstraction

SpecLog appends return after global ordering

Interface

// append to shard; return position of record in log
pos_t append(record_t r, sid_t shard);

Abstraction

SpecLog appends return after global ordering

SpecLog allows speculative consumption of records

Interface

// append to shard; return position of record in log
pos_t append(record_t r, sid_t shard);

// speculatively consume records after start
stream<record_t> subscribe(pos_t start);

Abstraction

SpecLog appends return after global ordering

SpecLog allows speculative consumption of records

SpecLog confirms or fails speculation

Interface

// append to shard; return position of record in log
pos_t append(record_t r, sid_t shard);

// speculatively consume records after start
stream<record_t> subscribe(pos_t start);

// callbacks for confirmations and mis-spec void confirm_spec(pos_t upto); void fail_spec(pos_t after);

Impact to Apps

Interface

// append to shard; return position of record in log
pos_t append(record_t r, sid_t shard);

// speculatively consume records after start
stream<record_t> subscribe(pos_t start);

// callbacks for confirmations and mis-spec void confirm_spec(pos_t upto); void fail_spec(pos_t after);

Impact to Apps

Apps listen to confirmations and mis-speculations

Interface

// append to shard; return position of record in log
pos_t append(record_t r, sid_t shard);

// speculatively consume records after start
stream<record_t> subscribe(pos_t start);

// callbacks for confirmations and mis-spec void confirm_spec(pos_t upto); void fail_spec(pos_t after);

Impact to Apps

Apps listen to confirmations and mis-speculations

Upon confirm_spec(k) → expose compute results till k

Interface

// append to shard; return position of record in log
pos_t append(record_t r, sid_t shard);

// speculatively consume records after start
stream<record_t> subscribe(pos_t start);

// callbacks for confirmations and mis-spec void confirm_spec(pos_t upto); void fail_spec(pos_t after);

Impact to Apps

Apps listen to confirmations and mis-speculations

Upon confirm_spec(k) → expose compute results till k

Upon fai1_spec(k) → rollback state until k and recompute

Interface

// append to shard; return position of record in log
pos_t append(record_t r, sid_t shard);

// speculatively consume records after start
stream<record_t> subscribe(pos_t start);

// callbacks for confirmations and mis-spec void confirm_spec(pos_t upto); void fail_spec(pos_t after);
















Challenge – Prediction is Hard

Hard to predict position in global order



predetermined cuts – fix beforehand global cuts system is expected to produce

predetermined cuts – fix beforehand global cuts system is expected to produce quota – fixed number of records in each shard report

predetermined cuts – fix beforehand global cuts system is expected to produce
quota – fixed number of records in each shard report
Sequencing layer waits for predetermined cut before sending global cut

predetermined cuts – fix beforehand global cuts system is expected to produce quota – fixed number of records in each shard report Sequencing layer waits for predetermined cut before sending global cut



primary backup

Shard I

ו	
J	IJ

primary backup

Shard 2

primary backup

Shard 3

predetermined cuts – fix beforehand global cuts system is expected to produce quota – fixed number of records in each shard report Sequencing layer waits for predetermined cut before sending global cut



predetermined cuts – fix beforehand global cuts system is expected to produce quota – fixed number of records in each shard report Sequencing layer waits for predetermined cut before sending global cut



Predetermined cuts - <2, I, I > <4, 2, 2 > <6, 3, 3 > ... <2n, n, n >

Each shard knows the sequence of global cuts the system will produce Any shard can predict positions of its records easily!





primary backup Shard 2

Quota: I		
	D	
primary	backup	



0 1 2 3

Each shard knows the sequence of global cuts the system will produce Any shard can predict positions of its records easily!



Each shard knows the sequence of global cuts the system will produce

Any shard can predict positions of its records easily!



What if a shard has more

0 1 2 3

Each shard knows the sequence of global cuts the system will produce

Any shard can predict positions of its records easily!



Each shard knows the sequence of global cuts the system will produce

Any shard can predict positions of its records easily!



What if a shard has more

Each shard knows the sequence of global cuts the system will produce

Any shard can predict positions of its records easily!



If shards meet quota, predetermined order same as actual global order

If shards meet quota, predetermined order same as actual global order In this case \rightarrow speculation succeeds

If shards meet quota, predetermined order same as actual global order In this case \rightarrow speculation succeeds If not \rightarrow speculation fails

If shards meet quota, predetermined order same as actual global order In this case \rightarrow speculation succeeds If not \rightarrow speculation fails

Cases:

If shards meet quota, predetermined order same as actual global order In this case \rightarrow speculation succeeds If not \rightarrow speculation fails Cases: \longrightarrow shard reports



If shards meet quota, predetermined order same as actual global order In this case \rightarrow speculation succeeds If not \rightarrow speculation fails Cases: \longrightarrow shard reports

If shards meet quota, predetermined order same as actual global order In this case \rightarrow speculation succeeds If not \rightarrow speculation fails Cases: \longrightarrow shard reports

shard internal failure

If shards meet quota, predetermined order same as actual global order In this case \rightarrow speculation succeeds If not \rightarrow speculation fails Concerned where \rightarrow shard reports

Cases:



If shards meet quota, predetermined order same as actual global order In this case \rightarrow speculation succeeds If not \rightarrow speculation fails Concerned where \rightarrow shard reports

Cases:



No mis-speculation Shard can internally mask failures!

If shards meet quota, predetermined order same as actual global order In this case \rightarrow speculation succeeds If not \rightarrow speculation fails Concerned where \rightarrow shard reports

Cases:



Shard can internally mask failures!



primary backup

If shards meet quota, predetermined order same as actual global order In this case \rightarrow speculation succeeds If not \rightarrow speculation fails Concerned where \rightarrow shard reports

Cases:



No mis-speculation Shard can internally mask failures!



If shards meet quota, predetermined order same as actual global order In this case \rightarrow speculation succeeds If not \rightarrow speculation fails Concerned where \rightarrow shard reports

Cases:



No mis-speculation Shard can internally mask failures! whole shard failure



primary backup

If shards meet quota, predetermined order same as actual global order In this case \rightarrow speculation succeeds If not \rightarrow speculation fails Concerned where \rightarrow shard reports

Cases:



No mis-speculation Shard can internally mask failures! whole shard failure



primary backup

Mis-speculations occur Shard cannot meet quota

If shards meet quota, predetermined order same as actual global order In this case \rightarrow speculation succeeds If not \rightarrow speculation fails Cases: \longrightarrow shard reports Rare!


















Real-Time App atop SpecLog











Outline

- Introduction
- Motivation
- SpecLog Abstraction and Interface
- Belfast Design and Implementation
- Evaluation

Fix-ante ordering and SpecLog provide a general framework for speculation

Fix-ante ordering and SpecLog provide a general framework for speculation

Belfast is an implementation, solves practical challenges

Fix-ante ordering and SpecLog provide a general framework for speculation

Belfast is an implementation, solves practical challenges

• Right-sizing quotas – rate-based quotas

Fix-ante ordering and SpecLog provide a general framework for speculation

Belfast is an implementation, solves practical challenges

- Right-sizing quotas rate-based quotas
- Quickly absorbing bursts lag-fix mechanism

Fix-ante ordering and SpecLog provide a general framework for speculation

Belfast is an implementation, solves practical challenges

- Right-sizing quotas rate-based quotas
- Quickly absorbing bursts lag-fix mechanism
- Dealing with long-term rate changes at shards

Fix-ante ordering and SpecLog provide a general framework for speculation

Belfast is an implementation, solves practical challenges

- Right-sizing quotas rate-based quotas
- Quickly absorbing bursts lag-fix mechanism
- Dealing with long-term rate changes at shards
- Retaining seamless reconfigurability

- speculation lease window

Fix-ante ordering and SpecLog provide a general framework for speculation

Belfast is an implementation, solves practical challenges

- Right-sizing quotas rate-based quotas
- Quickly absorbing bursts lag-fix mechanism
- Dealing with long-term rate changes at shards
- Retaining seamless reconfigurability
- Mis-speculations and failure handling view-change protocol

· speculation lease window



- What are the end-to-end latency benefits of Belfast?
- Does Belfast benefit end applications?

End-to-End Latency Benefits

Workload: 4KB records, downstream computation of 1.5ms per batch of consumed records

End-to-End Latency Benefits

Workload: 4KB records, downstream computation of 1.5ms per batch of consumed records



Benefits in Applications

Build 3 applications – intrusion detection, fraud monitoring and high-frequency trading



Benefits in e2e latency for realworld applications

More details and experiments in the paper

Many more experiments in the paper

- App evaluation under failure scenarios and mis-speculations
- Evaluation under bursts, rate changes
- End-to-end latency at scale

The paper covers many more discussions about Belfast



- Today's shared logs suffer high delivery latencies due to expensive ordering
- SpecLog, a new abstraction, allows speculative delivery by predicting global order
- SpecLog uses Fix-Ante ordering, enables high speculation success
- Belfast, an implementation of SpecLog, enables low e2e latencies for end applications





Available on GitHub: https://github.com/dassl-uiuc/speclog-artifact