DeltaINT: Toward General In-band Network Telemetry with Extremely Low Bandwidth Overhead

Siyuan Sheng^{1,3}, Qun Huang², and Patrick P. C. Lee³,

¹University of Chinese Academy and Sciences

²Peking University

³The Chinese University of Hong Kong

In-band Network Telemetry (INT)

Source pushes control information and device-internal states

- > Transit pushes states according to control information
- Sink extracts INT information and reports an event



Limitations of INT

Significant bandwidth overhead

- Linearly grow with the length of forwarding path
- Reduce effective bandwidth for network applications
- Increase likelihood of IP-level fragmentation
- ➤ Example
 - 5-node fat-tree topology in data center
 - Trace device ID, ingress port, and egress port, of 4B each
 - 12B per-node states and 8B INT control information
 - 68B in total \rightarrow at least 4.53% of 1,500B MTU in Ethernet

Existing Studies

Sampling-based methods

- Embed INT information to only a subset of sampled packets
- Reduce bandwidth overhead yet with slow convergence
- Cannot retrieve INT information unless collecting sufficient packets
- Other methods
 - Designed for specific telemetry tasks
- > All existing methods suffer from **low generality**
 - Cannot support all families of common applications

Our Contributions

DeltaINT, a general INT framework

- Extremely low bandwidth overhead
- High generality and convergence
- > Theoretical analysis on bandwidth mitigation guarantees
- Software simulation for various applications
 - For example, reducing up to 93% bandwidth cost in gray failure detection
- P4-based hardware implementation
- Open-source DeltaINT prototype

Four Families of Applications

Per-packet-per-node monitoring

- Collect per-node states for each packet (e.g., gray failure detection)
- Per-packet aggregation
 - Aggregate per-node states for each packet (e.g., congestion control)
- Static per-flow aggregation
 - Collect static per-node states for each flow (e.g., path tracing)
- Dynamic per-flow aggregation
 - Aggregate per-node states for each flow (e.g., latency measurement)

Our Solution

Key observation

- Delta, the change between current state and embedded state
- Delta is often **negligible** at most time in typical applications
 - For example, relatively stable hop latency and static device IDs
- Motivating example



Per-node Architecture in DeltaINT

Per-node architecture

- Calculate the delta between current states and embedded states
- Only if the delta exceeds a threshold, we insert current states into a packet and update the embedded state



> How to maintain embedded states efficiently in data plane?

Sketching in DeltaINT

Sketch-based technique

- Store approximate information with limited memory and computations
- Track embedded states in the data plane with limited resources
- Per-node sketch data structure
 - Each bucket stores a flowkey and the embedded states
 - Each entry of a packet includes a bitmap and the states being embedded



Primitives in DeltaINT

Four primitives to form DeltaINT workflow

- StateLoad
 - Hash flowkey and load embedded states from the first bucket matching flowkey
- DeltaCalc
 - Calculate the delta and compare with the predefined threshold
- StateUpdate
 - Update flowkey and relevant embedded states in the hashed buckets
- MetadataInsert
 - Insert a bitmap and the states with non-negligible deltas into the packet

Fit DeltaINT into applications with slight changes to primitives

Update Example



Evaluation

Methodology

- For software simulation, we use both bmv2 and NS3
- For hardware implementation, we compile P4 in Barefoot Tofino switch
- For sketch in the data plane, we keep 1MB memory and 1 hash function

➤ Experiments

- Gray failure detection
- Congestion control
- Path tracing
- Latency measurement
- Hardware resource usage

Gray Failure Detection

Tracked states

• 8-bit device ID, 8-bit ingress port, 8-bit egress port, and 32-bit latency

Bandwidth usage

- DeltaINT (8.1 bits) mitigates 93% bandwidth usage of INT-Path (112 bits)
- Reason: DeltaINT only embeds critical states with non-negligible deltas





Congestion Control

Tracked state: 8-bit link utilization

- Bandwidth usage
 - DeltaINT (\approx 1 bit) is better than PINT (8 bits)
 - Reason: DeltaINT only needs a 1-bit bitmap for negligible delta such that controller can be aware of the stable link utilization





Path Tracing

Tracked state: 8-bit device ID

- Bandwidth usage
 - DeltaINT (≈1 bit) is better than PINT (8 bits)
 - Reason: DeltaINT only needs a 1-bit bitmap for non-first packets of each flow due to static device ID with negligible delta



Path Tracing

Convergence

- Average number of required packets: DeltaINT (1) vs. PINT (120)
- Tail (99th percentile) number of required packets: DeltaINT (1) vs. PINT (350)
- Reason
 - DeltaINT only embeds per-node device ID in the first packet of each flow
 - PINT needs sufficient sampled packets to retrieve per-flow device IDs



Latency Measurement

- Tracked state: 8-bit latency
- Bandwidth usage
 - Web search workload: DeltaINT (2.6 bits) is better than PINT (10.3 bits)
 - Hadoop workload: DeltaINT (2.4 bits) is better than PINT (9.9 bits)
 - Reason: DeltaINT only embeds critical latency with non-negligible delta





Hardware Resource Usage

Hardware resource usage

- Percentages in brackets are fractions of total resource usage
- DeltaINT incurs slightly more SRAM, stages, and stateful ALUs
 - DeltaINT needs to track embedded states in the data plane
- INT incurs more PHV sizes and actions
 - INT has larger bandwidth overhead and hence more information to process and transmit

SRAM (KB) No. stages	No. actions No. ALUs	PHV size (bytes)
F1 224 KB (1.46%) 4 (33%)	9 (nil) 2 (4.17%)	115 (15%)
F2 224 KB (1.46%) 6 (50%)	10 (nil) 3 (6.25%)	111 (15%)
F3 304 KB (1.98%) 2 (16%)	4 (nil) 1 (2.08%)	104 (14%)
F4 224 KB (1.46%) 4 (33%)	9 (nil) 2 (4.17%)	115 (15%)
INT 176 KB (1.15%) 4 (33.%)	42 (nil) 0 (0%)	231 (30%)

Conclusion

- DeltaINT, a novel INT framework to achieve extremely low bandwidth overhead
 - Generality
 - Convergence
- Evaluation on various applications
 - DeltaINT incurs less bandwidth usage than state-of-the-art methods
- Source code:
 - http://adslab.cse.cuhk.edu.hk/software/deltaint

Thank You! Q & A