# FlowRadar:
# A Better NetFlow For Data Centers

## Yuliang Li

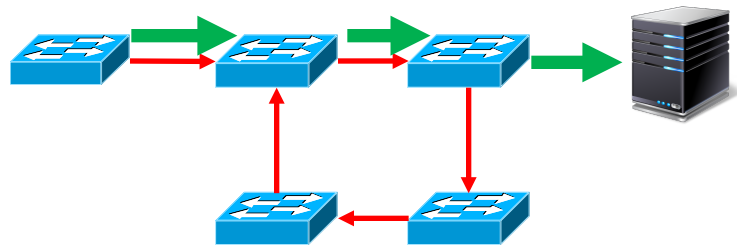Rui Miao          Changhoon Kim          Minlan Yu
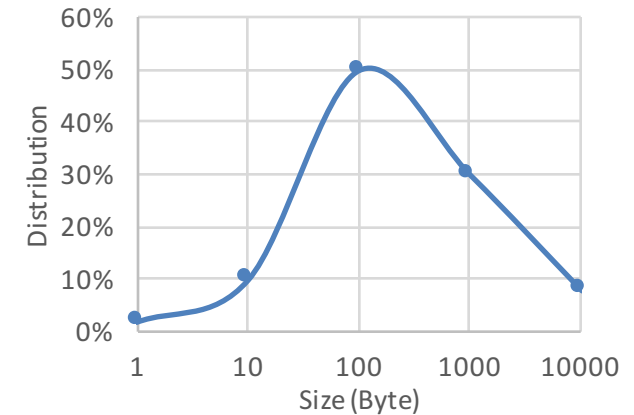
# Flow coverage in data centers

- ## Flow coverage
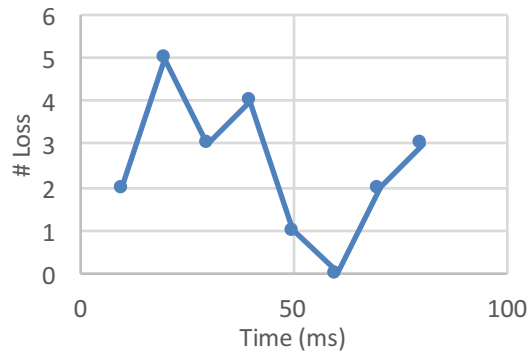  - Traffic monitoring needs to cover all the flows



Transient loop/blackhole



Fine-grained traffic analysis

# Temporal coverage in data centers

- **Temporal coverage**
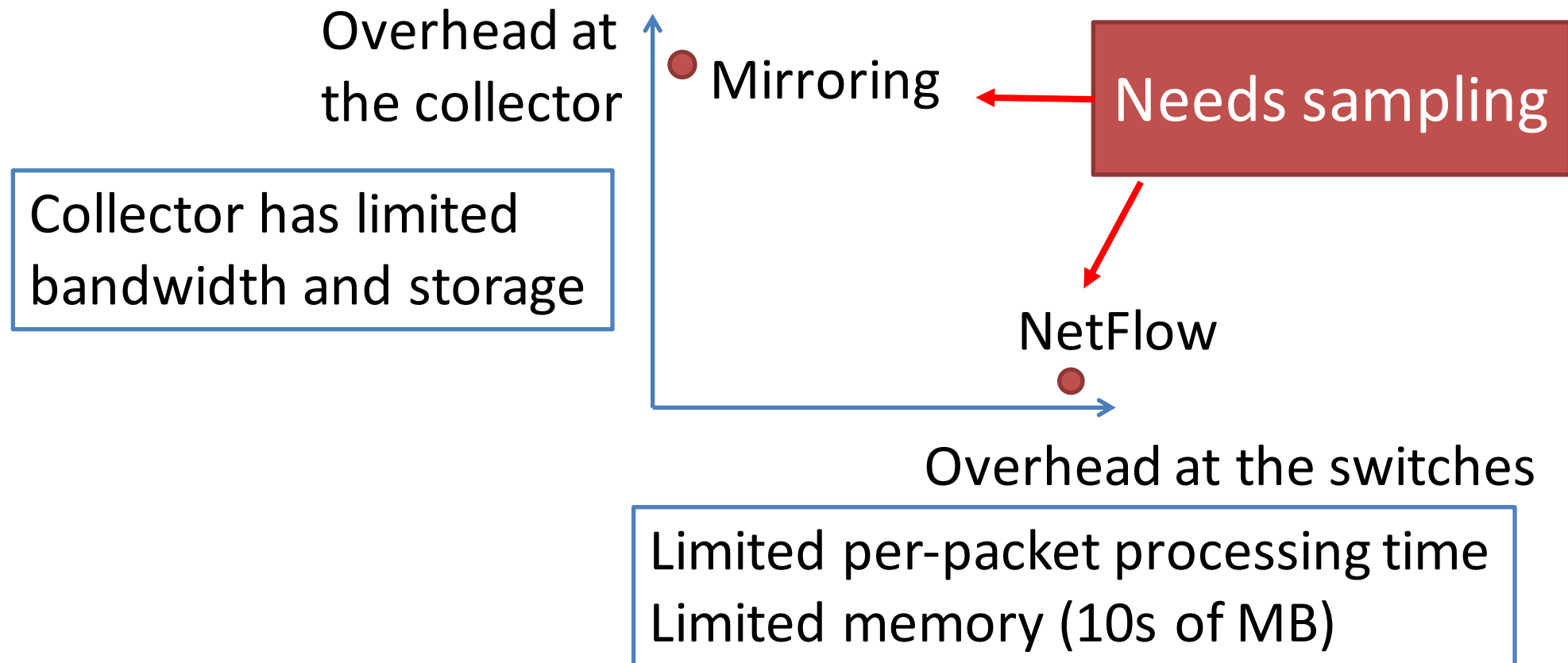  - Traffic monitoring needs millisecond-level flow information



Short-time scale Loss rate



Timely attack detection

# Key insight: division of labor

- Goal: report counters for all flows in fine-grained time granularity

Overhead at the collector

Mirroring

Needs sampling

Collector has limited bandwidth and storage

NetFlow

Overhead at the switches

Limited per-packet processing time
Limited memory (10s of MB)

# Key insight: division of labor

- Goal: report counters for all flows in fine-grained time granularity

Overhead at the collector

Mirroring

**FlowRadar**

Keep the memory usage small

NetFlow

Overhead at the switches

Use fixed operations per-pkt in switch

# FlowRadar architecture



Analyze

Flow&counter

Collector

Correlate network-wide info to extract per-flow counter

Network

Periodic report

Each switch maintains a fast and efficient data structure for half-baked per-flow counter

# Challenge: handling collision?

- Handling hash collision is hard
  - Large hash table → high memory usage
  - Linked list/Cuckoo hashing→ multiple, non-constant memory accesses

# Switch embraces collisions!

- Handling hash collision is hard
  - Large hash table → high memory usage
  - Linked list/Cuckoo hashing→ multiple, non-constant memory accesses
- Embrace the collision
  - Less memory and constant #accesses

# Switch embraces collisions!

- Embrace the collision: xor up all the flows
  - Less memory and constant #accesses



Counting table

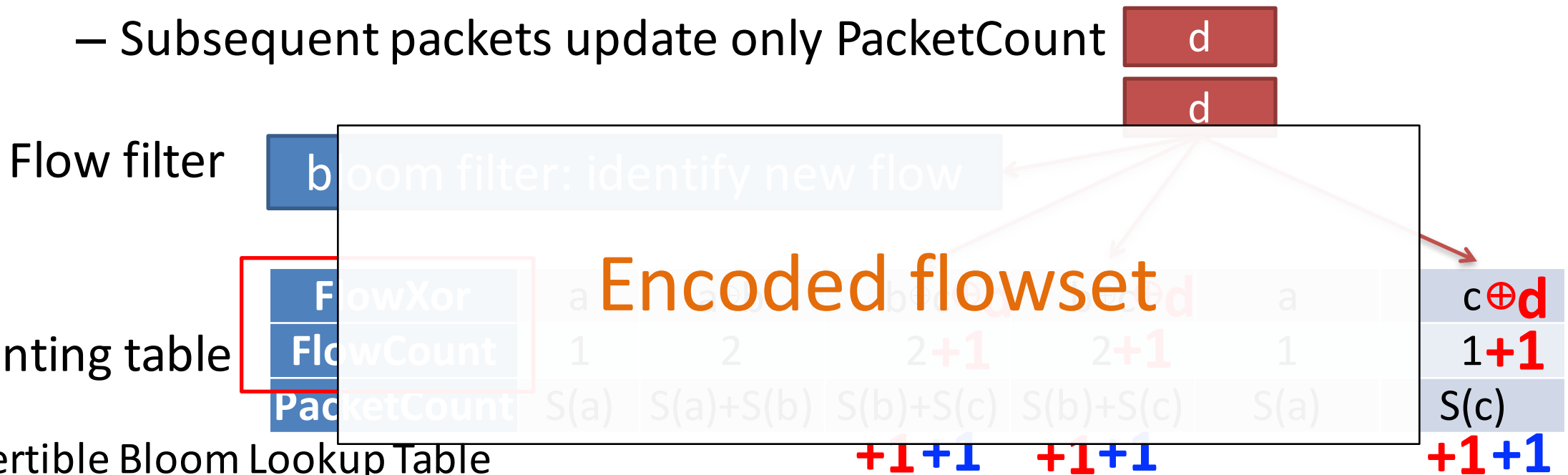| | | | | | | |
|---|---|---|---|---|---|---|
| **FlowXor** | a | a⊕b | b⊕c | b⊕c | a | c |
| **FlowCount** | 1 | 2 | 2 | 2 | 1 | 1 |
| **PacketCount** | S(a) | S(a)+S(b) | S(b)+S(c) | S(b)+S(c) | S(a) | S(c) |

S(x): #packets in x

[Invertible Bloom Lookup Table (arXiv 2011)]

# Switch embraces collisions!

- 1. Check and update the flow filter

- 2. Update counting table
  - Packet from a new flow, update all fields
  - Subsequent packets update only PacketCount

Flow filter

Counting table

[Invertible Bloom Lookup Table (arXiv 2011)]

Bloom filter: identify new flow

Encoded flowset

| | FlowXor | | | | | | $c \oplus d$ |
| --- | --- | --- | --- | --- | --- | --- | --- |
| FlowCount | | | 2+1 | 2+1 | | | 1+1 |
| PacketCount | S(a) | S(a)+S(b) | S(b)+S(c) | S(b)+S(c) | S(a) | | S(c) |

# Easy to implement in merchant silicon

- Switch data plane
  - Fixed operations in hardware
  - Small memory, 2.36MB for 100K flows

- Switch control plane
  - Control plane gets the small flowset every 10ms

- We implemented it using P4 Language.

# FlowRadar architecture

# Stage1. SingleDecode

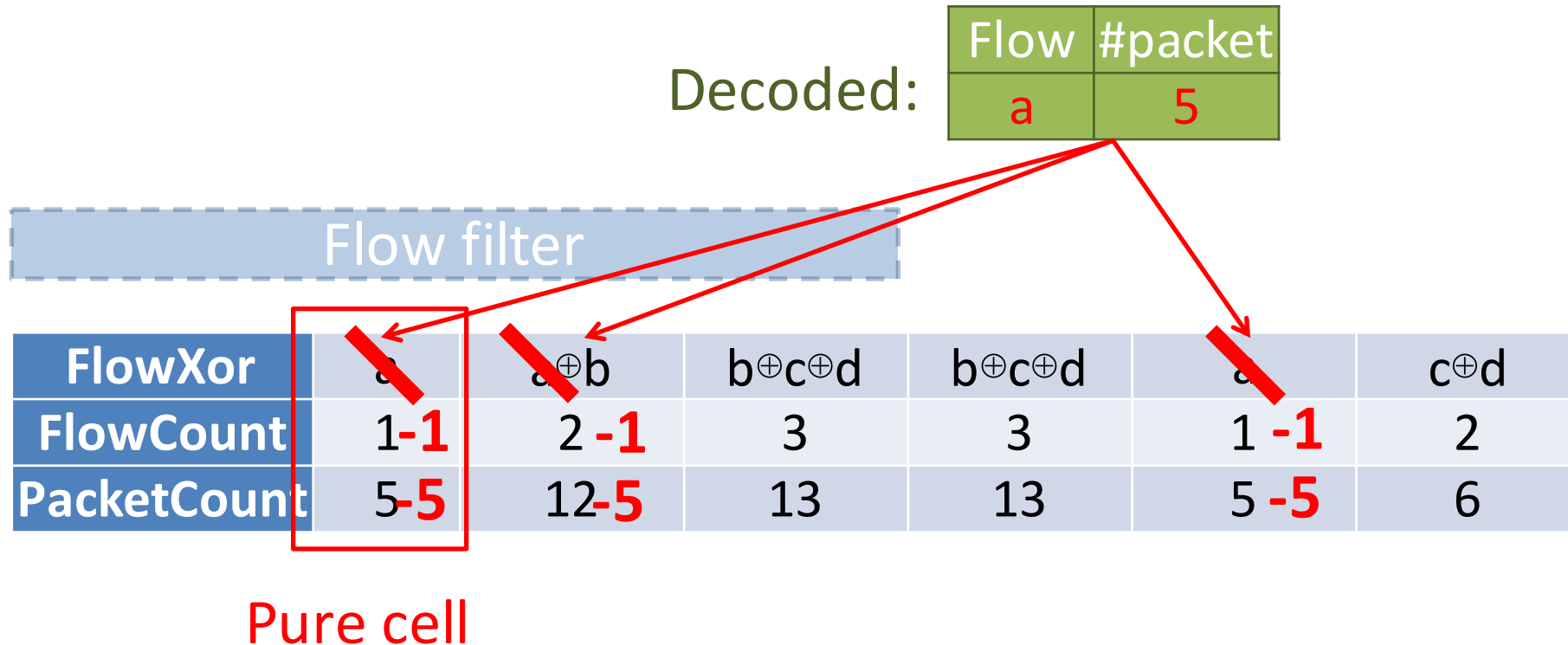Input: a single encoded flowset

Output: per-flow counters

Flow filter

| Bloom filter |
|:---:|

Counting table

| FlowXor | ... |
|:---:|:---:|
| FlowCount | ... |
| PacketCount | ... |

| Flow | #packet |
|:---:|:---:|
| a | S(a) |
| ... | ... |

# Stage1. SingleDecode

- Find a **pure cell**: a cell with one flow
- Remove the flow from all cells

Decoded:

| Flow | #packet |
|------|---------|
| a    | 5       |

Flow filter

| | | | | | | |
|---------|-----|---------|---------------|---------------|---------|---------------|
| **FlowXor** | a | a⊕b | b⊕c⊕d | b⊕c⊕d | a | c⊕d |
| **FlowCount** | 1-**1** | 2 -**1** | 3 | 3 | 1 -**1** | 2 |
| **PacketCount** | 5-**5** | 12-**5** | 13 | 13 | 5 -**5** | 6 |

Pure cell

# Stage1. SingleDecode

- Find a cell with one flow (pure cell)
- Remove the flow from all cells
  - Create more pure cells

Decoded:

| Flow | #packet |
|------|---------|
| a    | 5       |

- Iterate until no pure cells

Flow filter

| FlowXor | 0 | b | b⊕c⊕d | b⊕c⊕d | 0 | c⊕d |
|---------|---|---|-------|-------|---|-----|
| FlowCount | 0 | 1 | 3 | 3 | 0 | 2 |
| PacketCount | 0 | 7 | 13 | 13 | 0 | 6 |

# Stage1. SingleDecode

Decoded:

| Flow | #packet |
|------|---------|
| a | 5 |
| b | 7 |

**We want to leverage the network-wide info to decode more flows**

| | Flow filter | | | | | |
|---|---|---|---|---|---|---|
| **FlowXor** | 0 | 0 | c⊕d | c⊕d | 0 | c⊕d |
| **FlowCount** | 0 | 0 | 2 | 2 | 0 | 2 |
| **PacketCount** | 0 | 0 | 6 | 6 | 0 | 6 |

# FlowRadar architecture

Analyze

Flow&counter

**Collector**

Stage1.SingleDecode → Stage2. Network-wide Decode

**Network**

**Periodic report**

Encoded flowset

Encoded flowset

Encoded flowset

# Key insight: overlapping sets of flows

- The sets of flows overlap across hops
  - We can use the redundancy to decode more flows

- Use different hash functions across hops

Provision memory based on avg(#flows), not max(#flows)
- SingleDecode for normal case
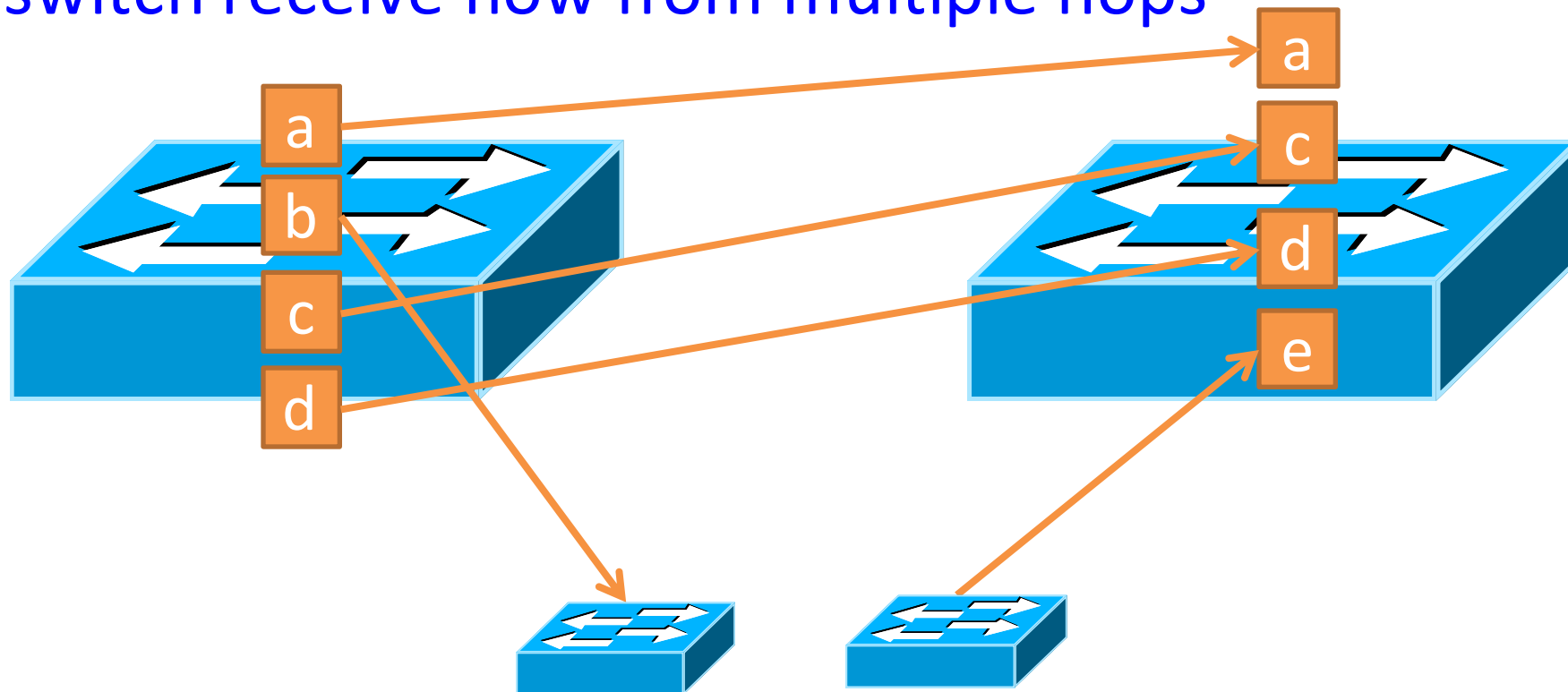- Network-wide decoding for bursts of flows

Use ~~10~~ 5 cells to decode 4 flows

Collector can leverage flowsets from all switches to decode more
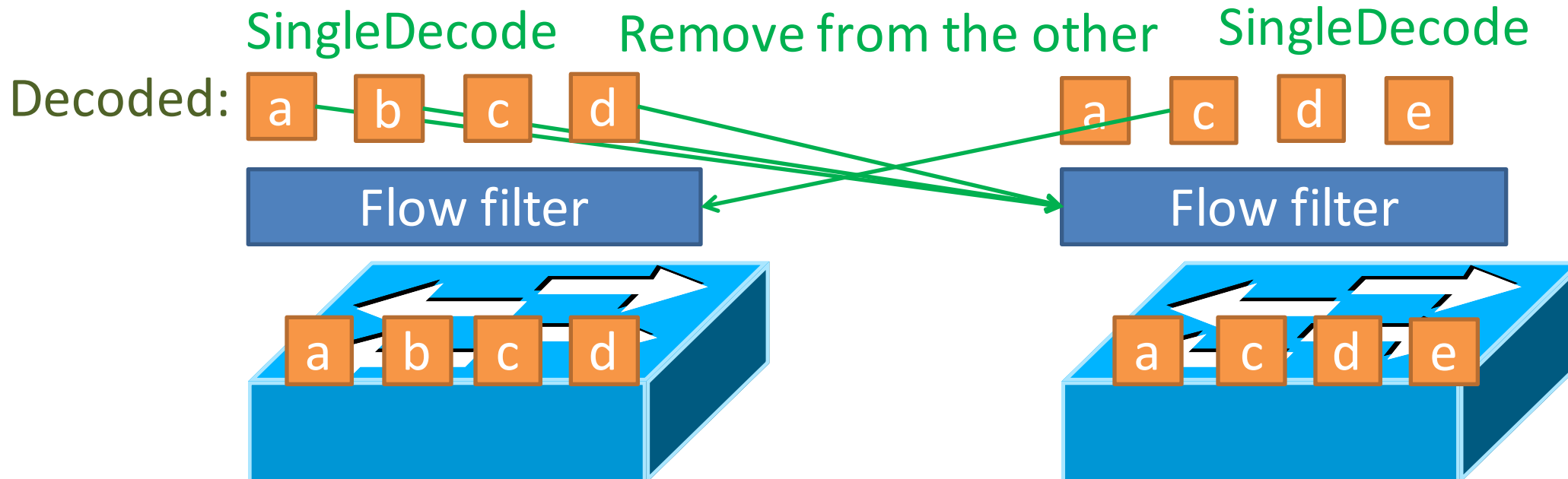
# Challenge 1: sets of flows not fully overlapped

- Flows from one switch may go to different next hops
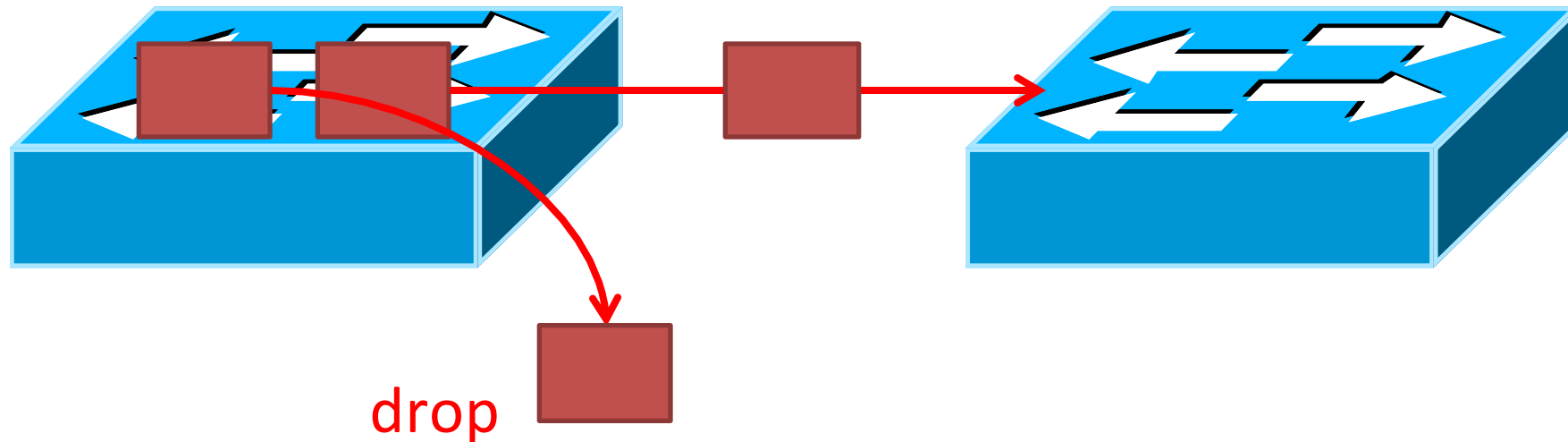- One switch receive flow from multiple hops

# Challenge 1 solution: use flow filter to check

- Generalize to network
  - No need for routing info
  - Incremental deployment

SingleDecode    Remove from the other    SingleDecode

Decoded:    a  b  c  d        a  c  d  e

Flow filter        Flow filter

a  b  c  d        a  c  d  e

# Challenge 2: counters are different across hops

- The counter of a flow may be different across hops
  - Some packets may get lost
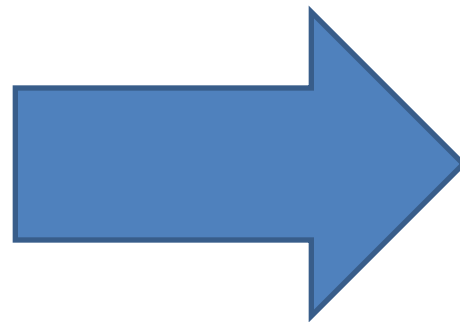  - On-the-fly packets

drop

# Challenge 2 solution: solve linear equations

- We got full list of flows
- Combine with counting table
- Construct and solve a linear equation system for each switch
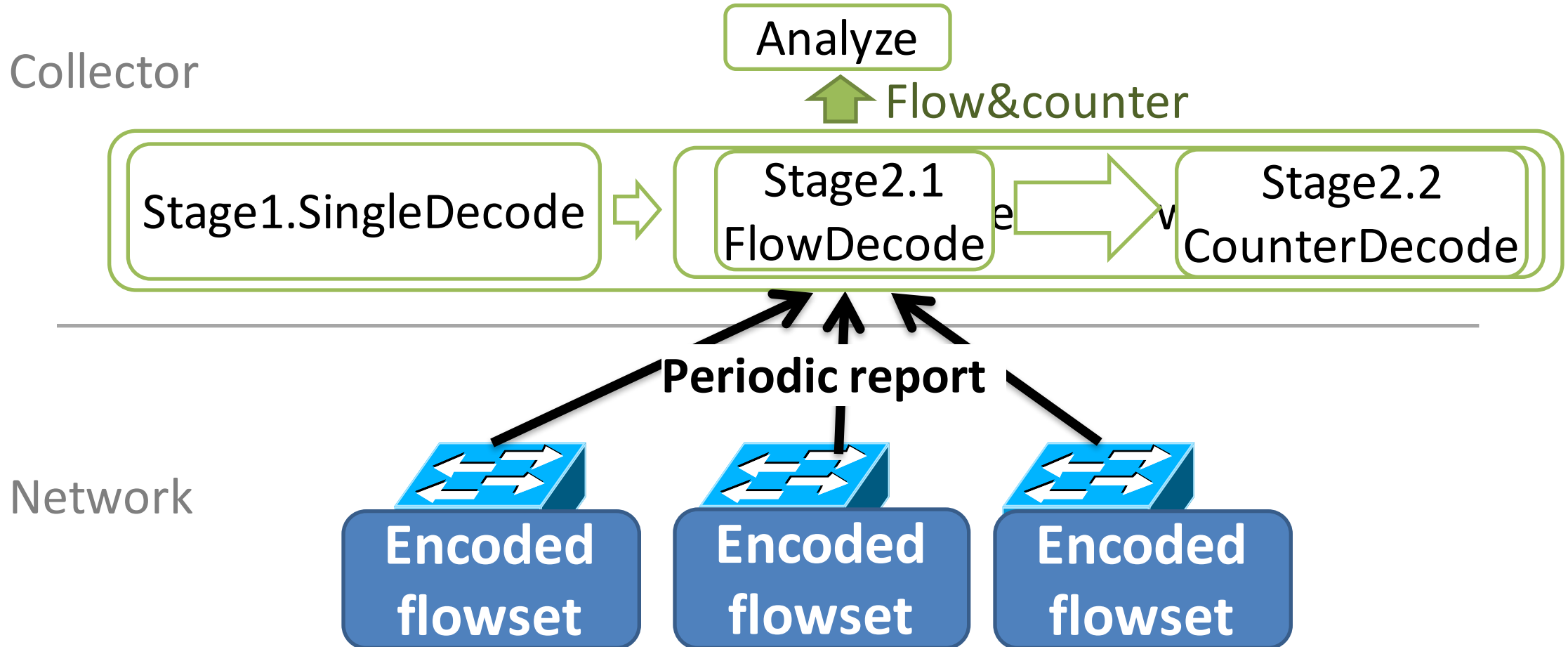- Speed up by using counter's properties to stop solver earlier



| a | b | c | d |

| FlowXor | |
|---|---|
| | ... |
| **FlowCount** | ... |
| **PktCount** | ... |

| Flow | #pkt |
|---|---|
| a | 5 |
| b | 7 |
| c | 4 |
| d | 2 |

# FlowRadar architecture

Analyze

Collector

Flow&counter

Stage1.SingleDecode → Stage2.1 FlowDecode e → v Stage2.2 CounterDecode

**Periodic report**

Network

**Encoded flowset**  **Encoded flowset**  **Encoded flowset**

# Evaluations

SingleDecode vs. Network-wide Decode

Collector

Analyze

Flow&counter

Stage1.SingleDecode ➡ Stage2.1 FlowDecode ➡ Stage2.2 CounterDecode

Bandwidth usage

**Periodic report**
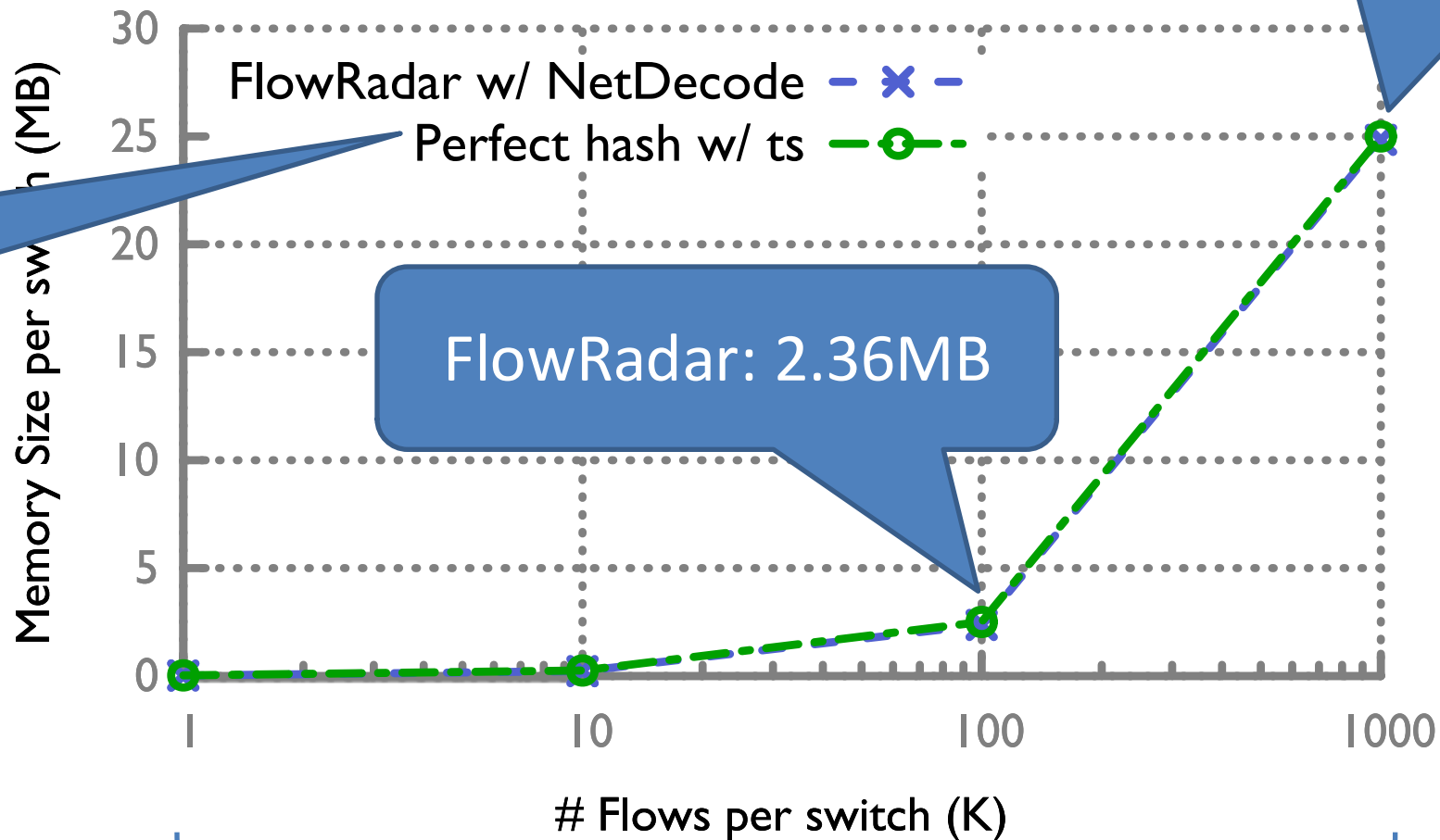
Network

**Encoded flowset**

**Encoded flowset**

**Encoded flowset**

Memory efficiency

24

# Evaluation

- Simulation of k=8 FatTree (80 switches, 128 hosts) in ns3
- Config the memory base on avg(#flow),
  - when burst of flows happens, use network-wide decode
- The worst case is all switches are pushed to max(#flow)
  - Traffic: each switch has same number of flows, and thus same memory
- Each switch reports the flowset every 10 ms.

# Memory efficiency

FlowRadar: 24.8MB

#cell=#flow
(Impractical)
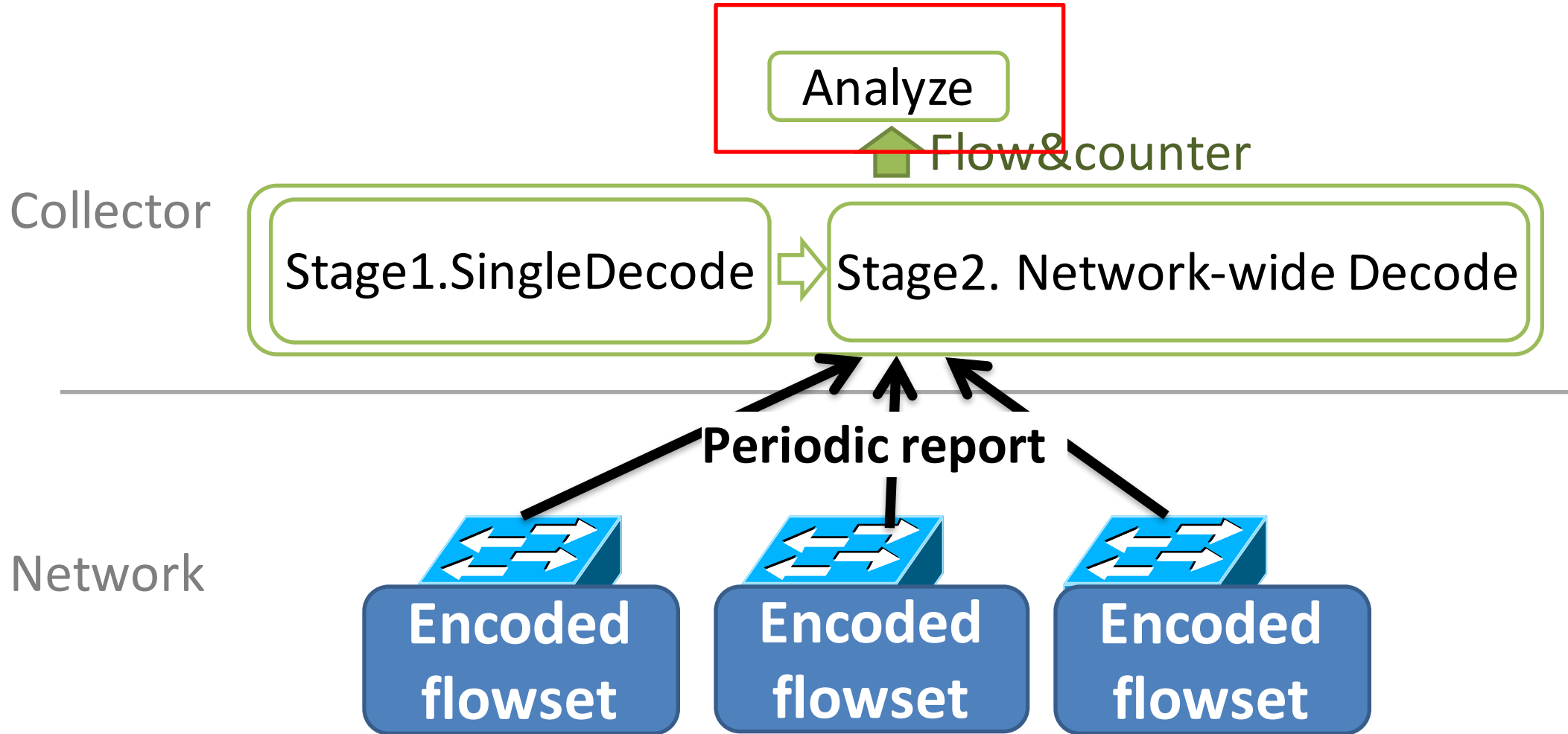
FlowRadar: 2.36MB

FlowRadar w/ NetDecode
Perfect hash w/ ts

Memory Size per switch (MB)

30
25
20
15
10
5
0

1    10    100    1000

# Flows per switch (K)

Log scale

26

# Other results

- Bandwidth usage
  - Only 0.52% based on topology and traffic of Facebook data centers (sigcomm'15)

- NetDecode improvement over SingleDecode
  - SingleDecode 100K flow, which takes 10ms
  - NetDecode 26.8% more flows with the same memory, which takes around 3 sec

# FlowRadar analyzer

Analyze

Flow&counter

Collector

Stage1.SingleDecode → Stage2. Network-wide Decode

Network

**Periodic report**

**Encoded flowset**

**Encoded flowset**
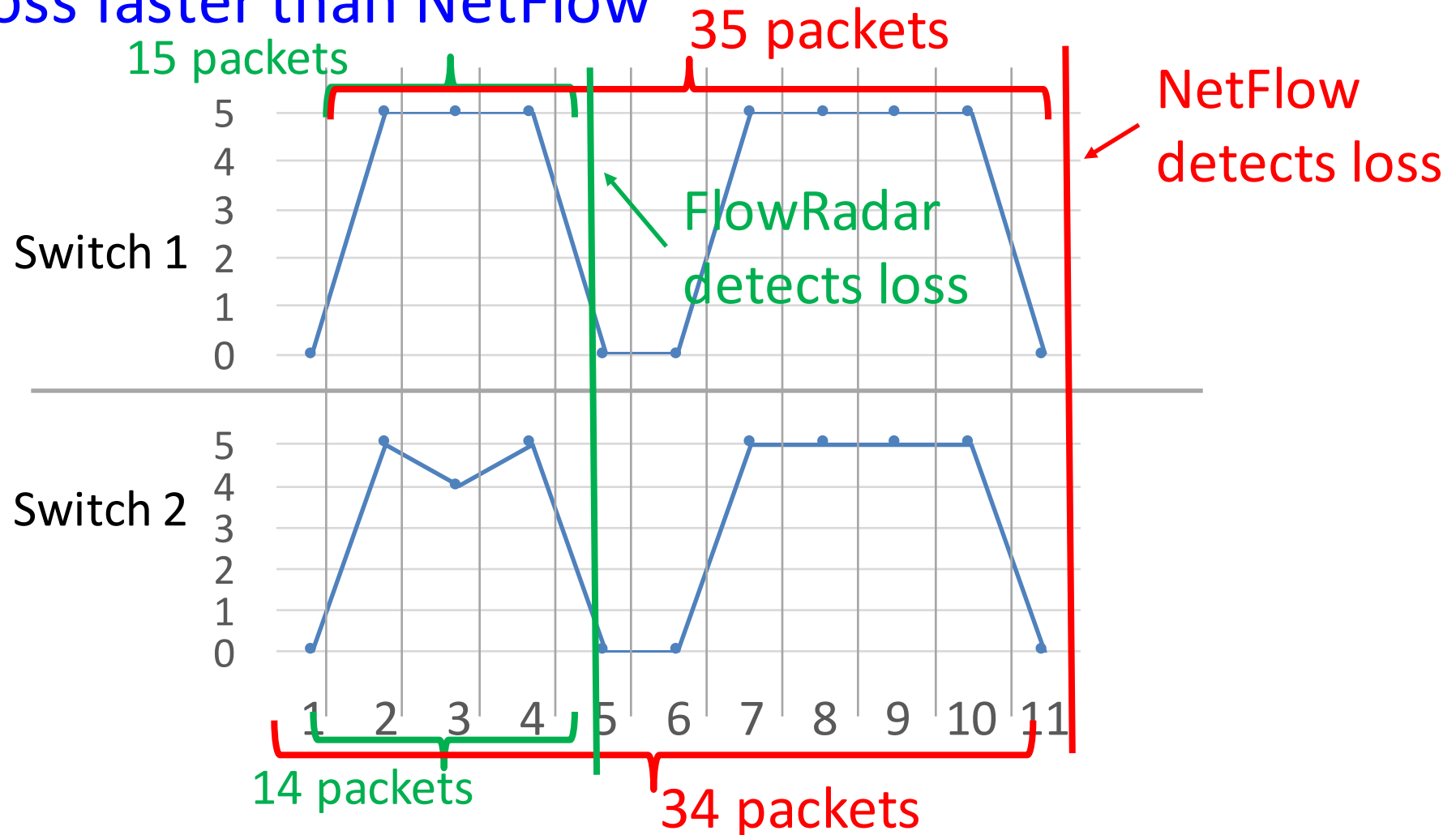
**Encoded flowset**

# Analysis applications

- Flow coverage
  - Transient loop/blackhole
  - Error in match-action table
  - Fine-grained traffic analysis

- Temporal coverage
  - **Short time-scale per-flow loss rate**
  - ECMP load imbalance
  - Timely attack detection

# Per-flow loss map: better temporal coverage

- Detect loss faster than NetFlow

# Conclusion

- Report counters for all flows in fine-grained time granularity
- Fully leverage the capability of both the switches and the collector
  - Switch: fixed per-packet processing time, memory-efficient
  - Collector: Network-wide decoding

Overhead at
the collector

Mirroring

**FlowRadar**

NetFlow

Overhead at the switches