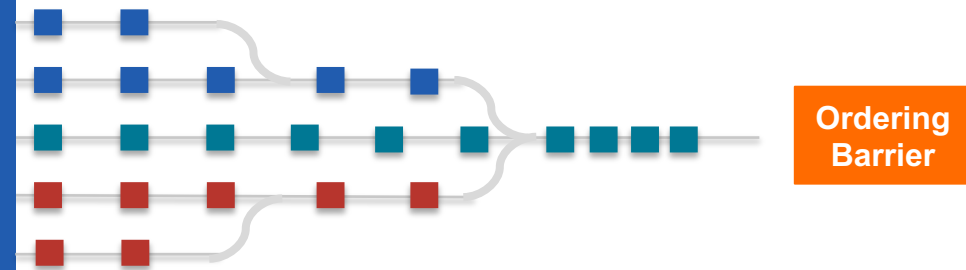


# Triggerless Data Acquisition for Online Reconstruction in High-Rate Experiments

Ring-CAM timestamp resequencing in FPGA

16-min talk · Data Acquisition and Trigger Architectures Session

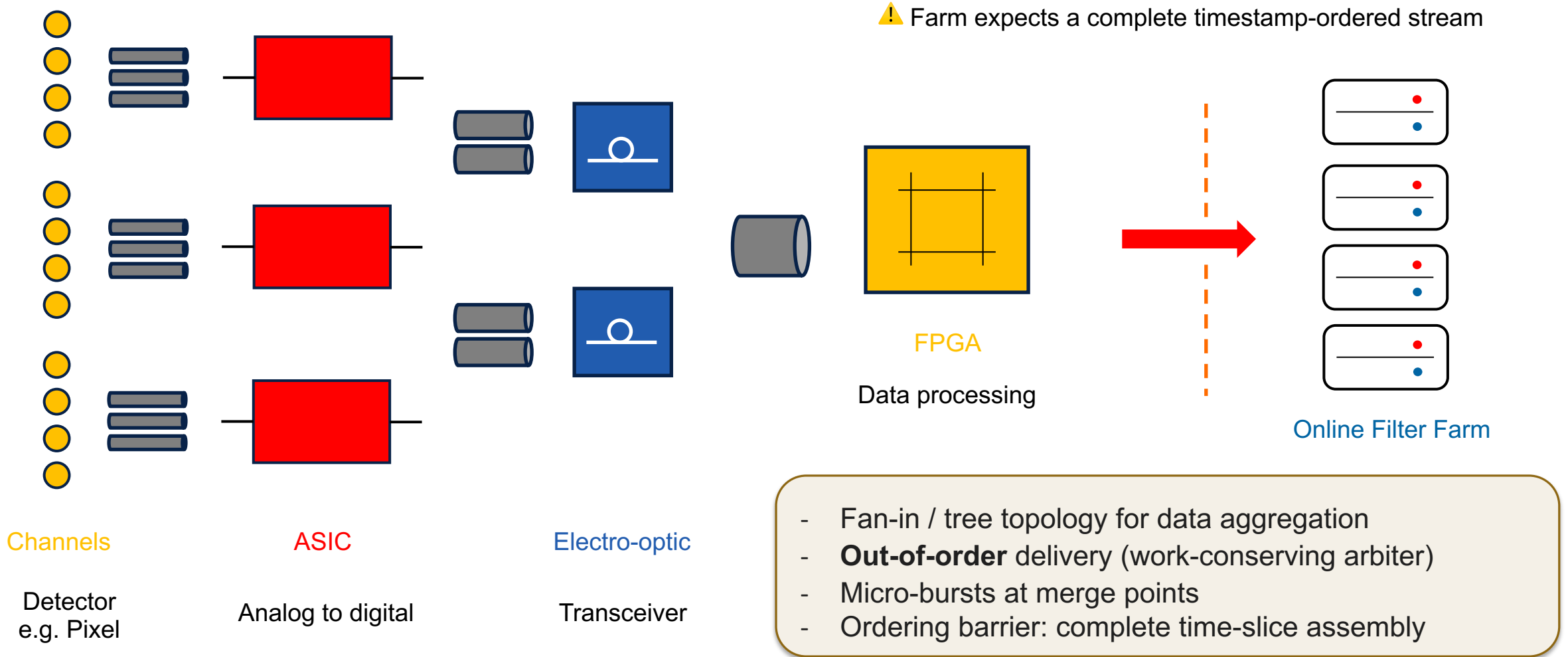
Yifeng Wang (ETH Zürich)



25<sup>th</sup> IEEE Real Time Conference  
La Biodola, Elba, Italy



## Triggerless DAQ: Fan-In Creates an Ordering Barrier



A detector time slice is distributed across many readout lanes but requires synchronization at the last stage

## HEP Bottleneck: Time-Slice Assembly, Not Average Bandwidth

### Future Physics Demands

- High data rate, low physics loss
- Smart-trigger and triggerless readout modes
- Online filtering is the prevailing strategy
- Traffic driven by physics occupancy, not a fixed trigger accept rate

### Engineering blockers

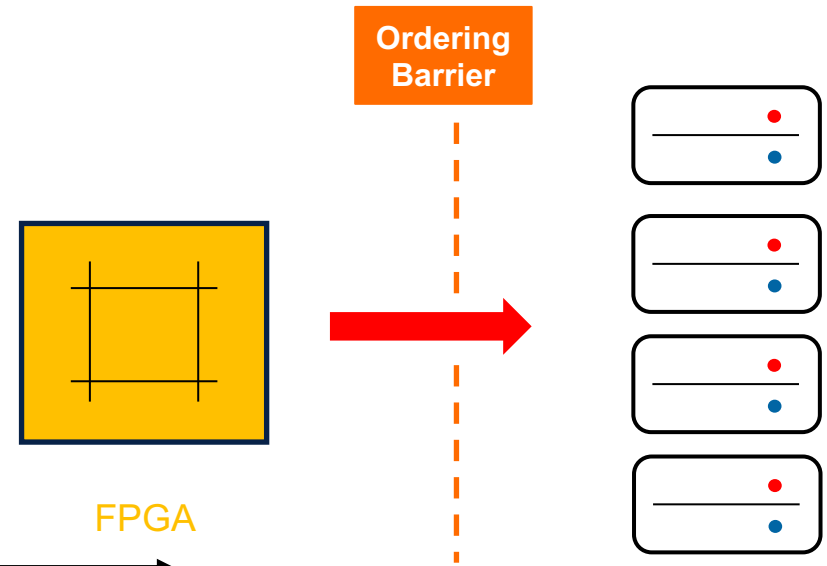
- Faster transceivers
- Better integration methods, e.g. chiplets and silicon-photonics
- Remaining blocker: deterministic time-slice (*frame*) transport across many lanes
- **frame** = *fixed-width time slice*.

At high rate, the critical path shifts from average bandwidth to bounded lane skew and complete time-slice (*frame*) delivery

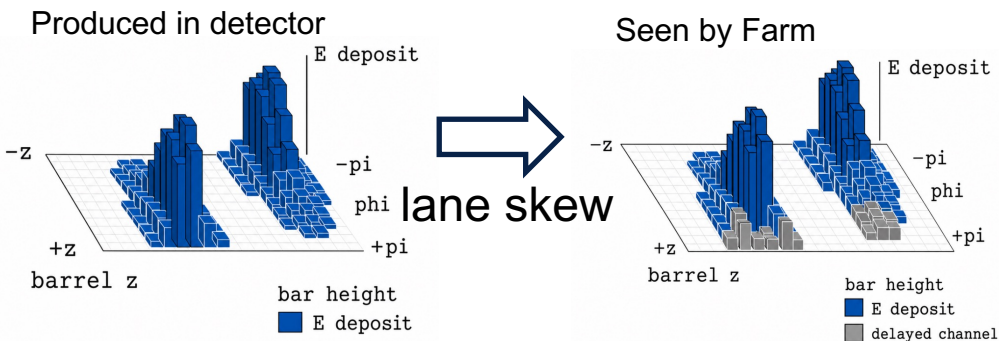
# Time-Slice Assembly Waits for the Slowest Lane

- **Hit:** a timestamped detector readout word.
- **Lane/link:** an independent readout stream carrying hits and/or progress markers.
- **Time slice:** fixed timestamp interval  $[T, T + W)$ . It is complete when all lanes have delivered hits in the interval or advanced a watermark beyond it.
- A frame needs hits or empty/watermark evidence from all lanes; lane skew makes the frame incomplete

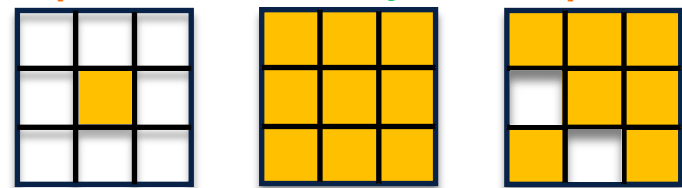
The point where downstream processing cannot safely close or process timestamp interval  $T$  until every input lane has either delivered all hits for  $T$  or advanced a watermark beyond  $T$ .



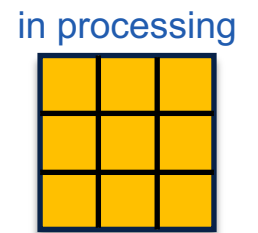
Example Event Display of Sparse/Zero-suppressed Detector a lane may have no hit



Incomplete frame ready Incomplete frame



Online Filter Farm



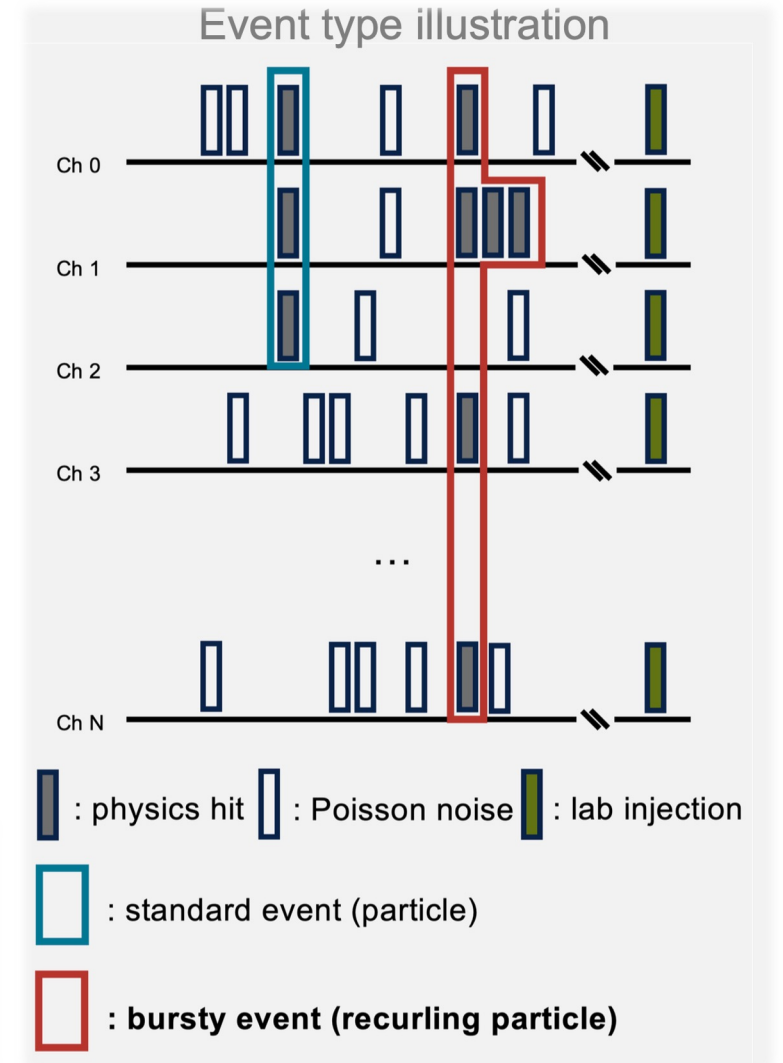
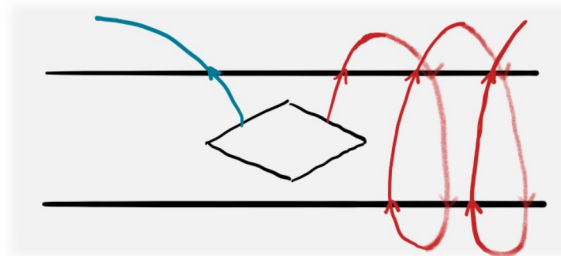
**Head-of-Line Blocking:** later timestamps may be present but cannot be released before the incomplete earlier frame

## Physics Events Create Short Bursts

- ✓ Average link bandwidth can still be sufficient
- ! Physics-rich events concentrate hits in the same time window

### Examples

- Fan-in concentrates simultaneous hits from many lanes into local bursts.
- A charged-particle track is already bursty at the source: correlated hits share nearby timestamps.

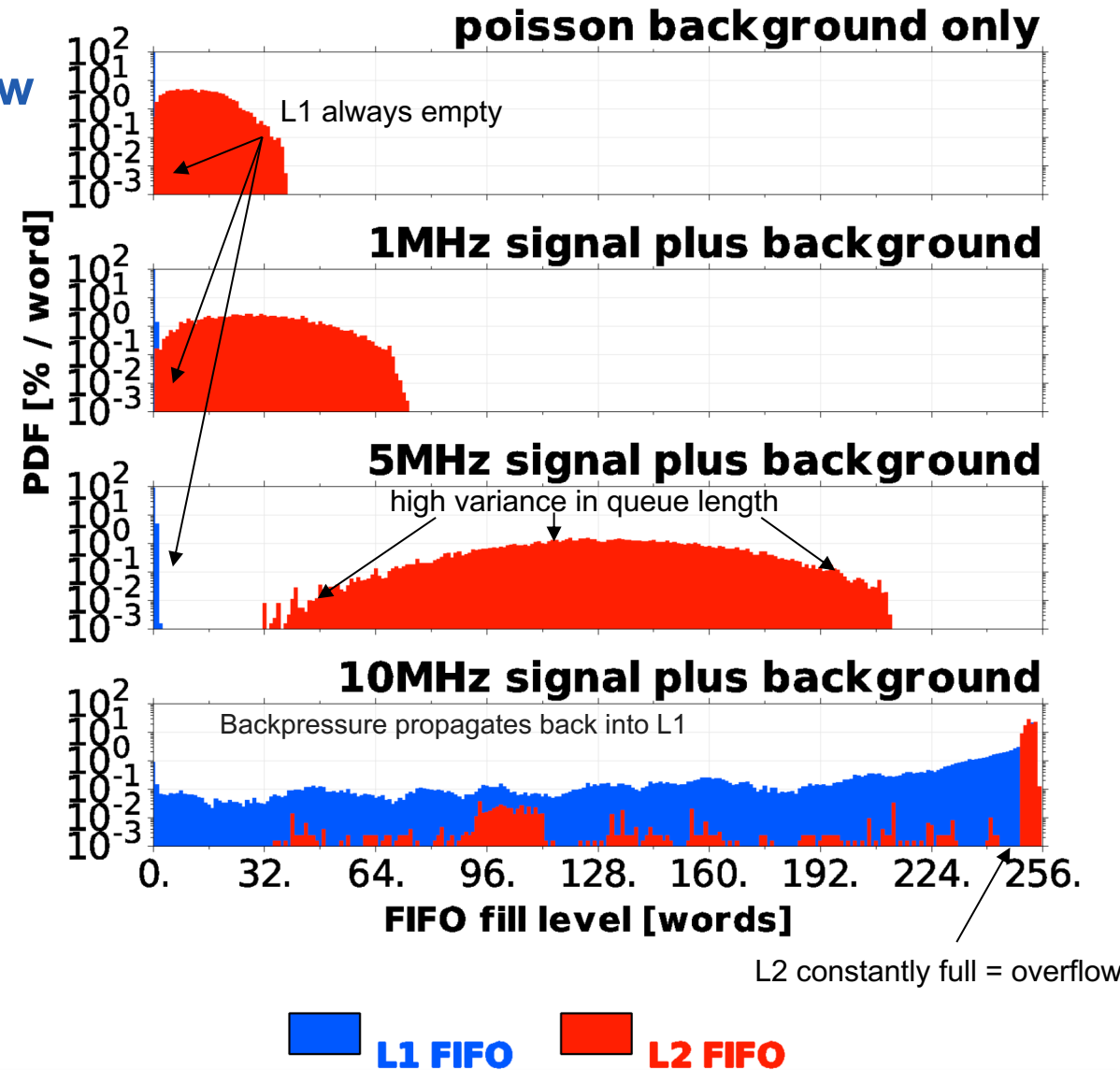
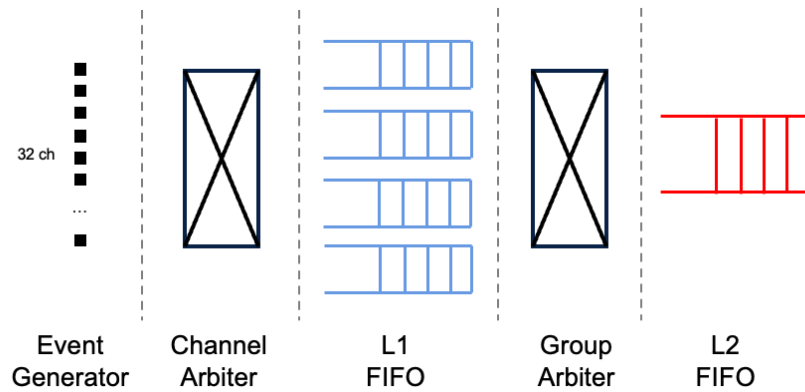


Bursts may contain high-value physics, so local loss is not an acceptable mitigation

## Local FIFO Pile-Up Creates Loss and Lane Skew

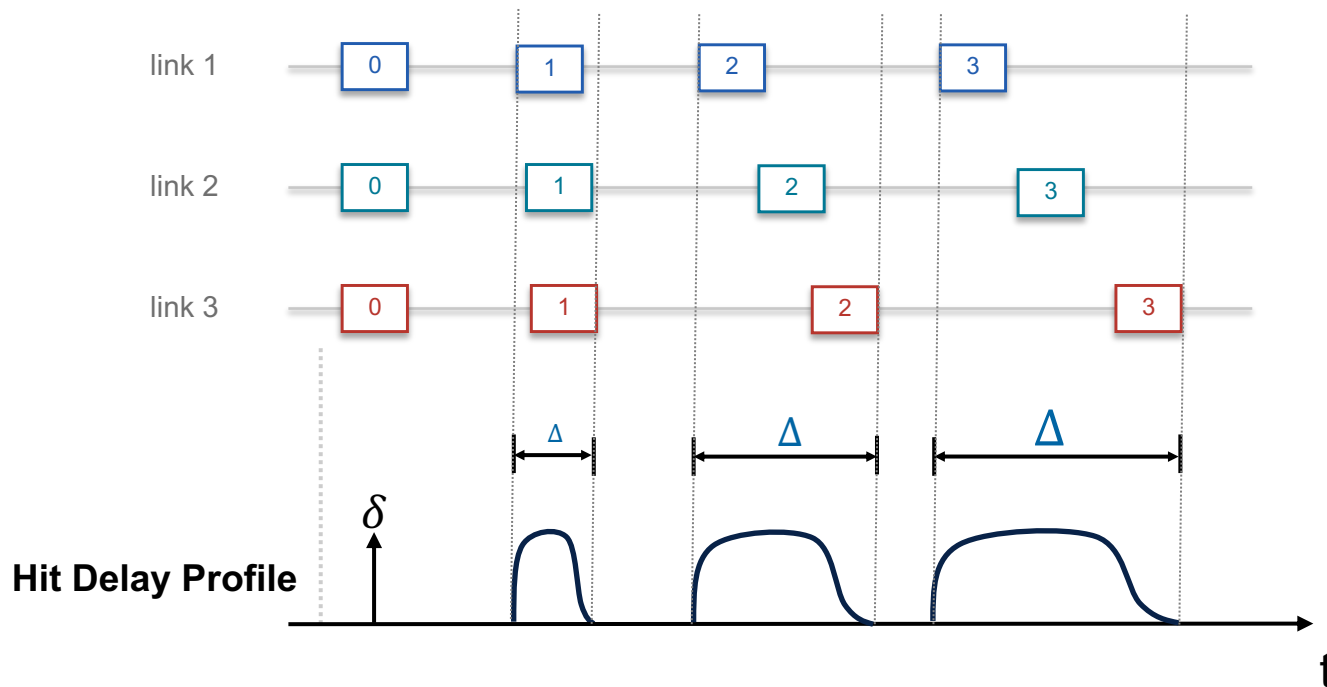
- ✔ No persistent bandwidth bottleneck required
- ! Short physics bursts fill local FIFOs
- ✘ FIFO pile-up becomes overflow and lane skew
- ✘ Too fast for trigger feedback; needs local mitigation

### Example ASIC Architecture



FIFO fill level maps directly to residence time; unequal residence time across lanes becomes lane skew

## Lane Skew Must Be Reset Locally



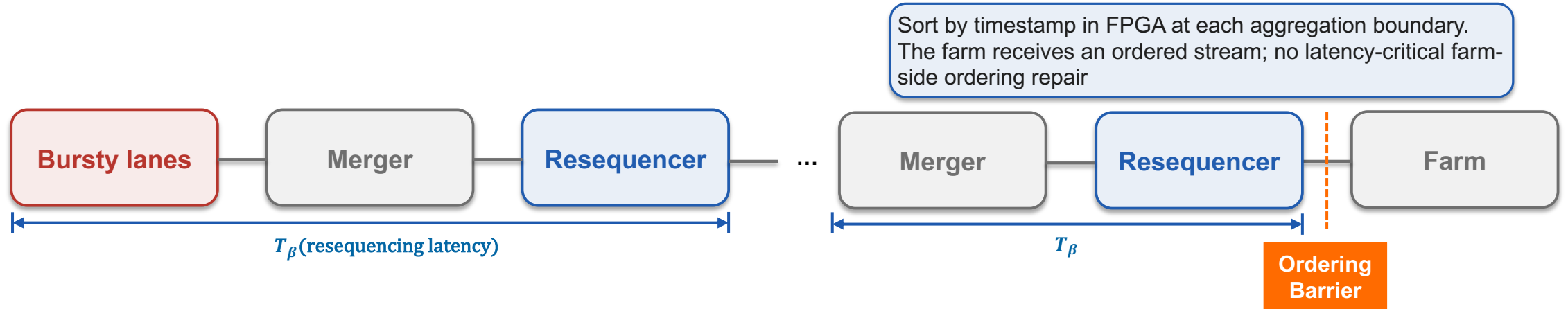
- **Hit delay:**  $d_i = t_i^{arrival} - t_i^{hit}$

- **Lane Skew:**  $\Delta = \max_i d_i - \min_i d_i$ , for hits or watermarks belonging to the same timestamp/frame. It is also known as *delay spread*.

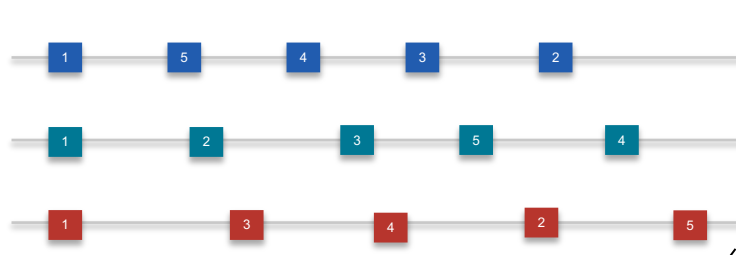
- **Skew Bound:**  $\Delta_{max}$  is the configured worst-case bound used for resequencer sizing.

Skew accumulation across stages, if no local reset

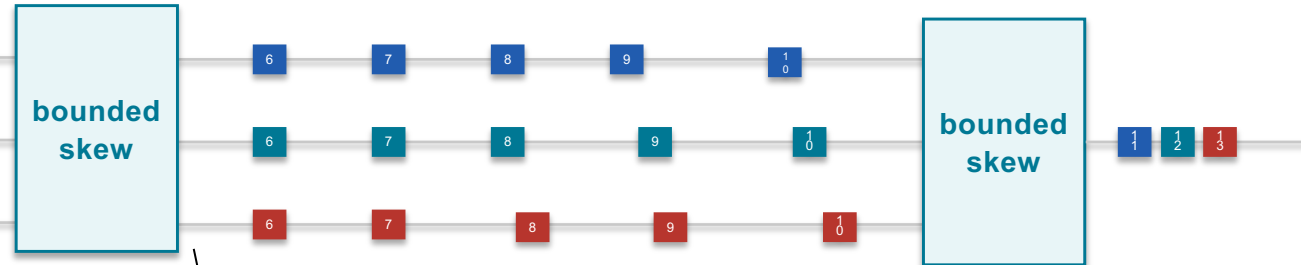
# FPGA Resequencing Resets Skew Before the Farm



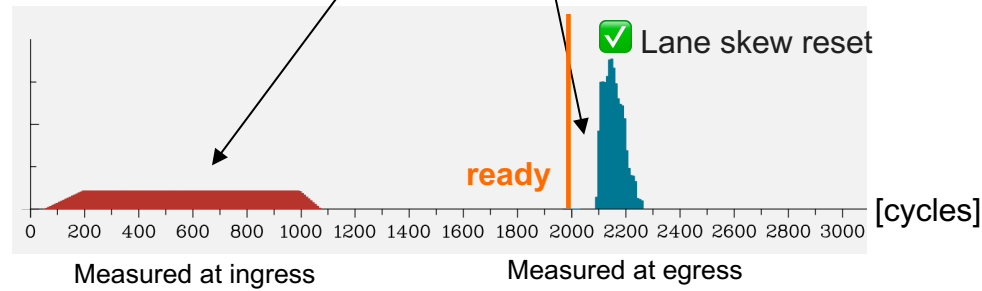
## Before resequencing



## After resequencing



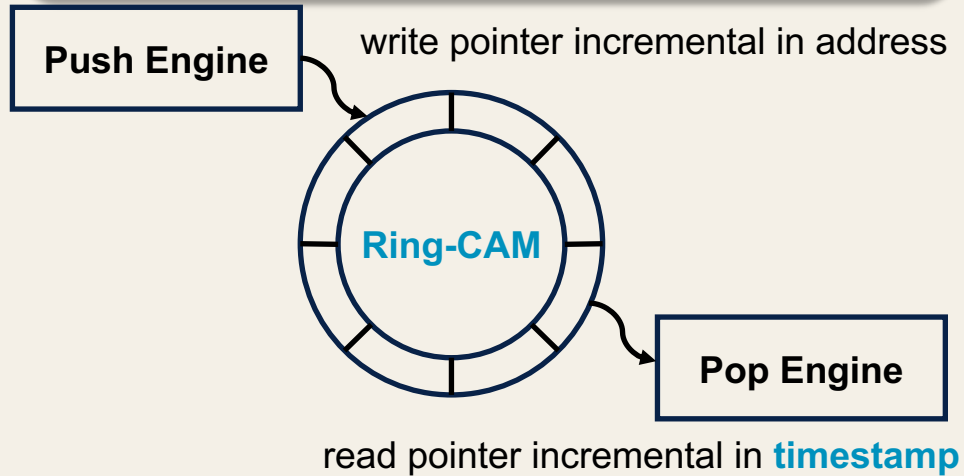
## Hits delay profile of lab measurement



## Ring-CAM: Write by Arrival, Read by Timestamp

### Ring-CAM (resequencer)

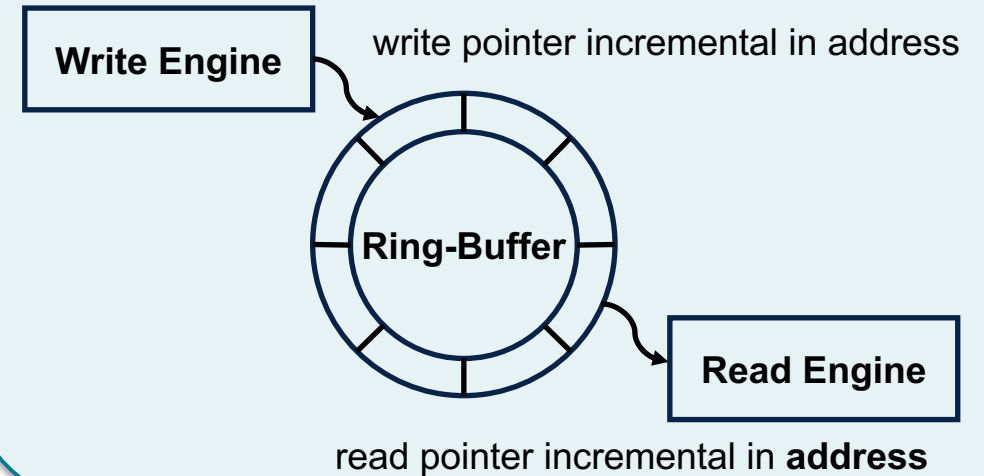
- Write hits as they arrive; read the earliest timestamp that is ready. Effective sorted in timestamp.
- **Ready:** timestamp  $\tau$  is ready when  $\tau$  is older than the current time plus delay of skew window,  $\Delta_{\max} + t$ .
- Fixed push/pop path; depth is sized by the skew window



**CAM** = Content Addressable Memory  
 Compare to **RAM**: write op is the same, read is content or address indexed

### Compare with normal Ring-Buffer

- Normal FIFO: storage order = arrival order
- Ring-CAM: storage order = arrival order, output order = timestamp order



altera™



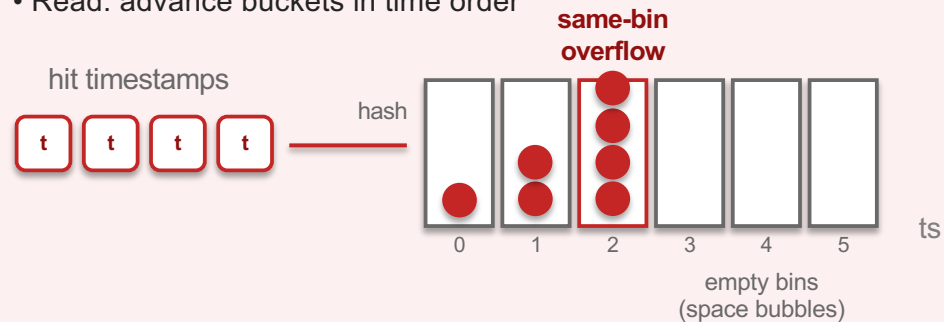
Source code available:  
 RTL Ring-CAM resequencer  
 VHDL and SystemVerilog versions

## Why Ring-CAM Instead of a Calendar Queue?

Both produce timestamp-ordered output; the difference is how FPGA storage is allocated under bursts.

### Calendar Queue — classic bucketed sorter

- Priority queue introduced for simulation event sets (Brown, CACM 1988)
- Write: bucket =  $\text{floor}(\text{timestamp} / W) \bmod N$ ;  $W$ =bucket width,  $N$ =total bins
- Read: advance buckets in time order

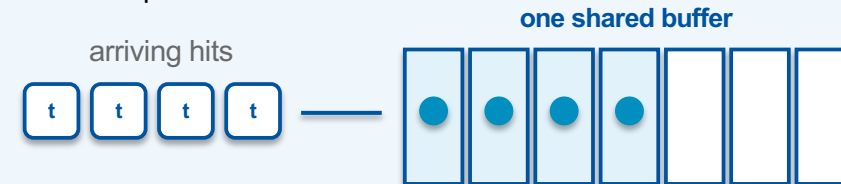


DAQ issue: a hardware calendar queue gives each timestamp bin bounded capacity. Bursty hits can overflow one bin while total memory is still available elsewhere.

Calendar Queue: Randy Brown, CACM 1988

### Ring-CAM — shared storage + timestamp lookup

- Write: next free slot in ring RAM (arrival order)
- Index: CAM records timestamp  $\rightarrow$  slot address
- **Read (CAM lookup):** timestamp  $\rightarrow$  list/bitmap of valid slots. Multiple hits with the same timestamp are emitted before advancing to the next timestamp.



All timestamps share the full depth. Size the buffer by input rate  $\times$  maximum lane skew window, plus margin.

**Takeaway: Calendar Queue sorts by placing hits into time bins; Ring-CAM sorts by looking up timestamps in shared storage.**

## Sizing Rule: Depth $\geq$ Rate $\times$ Skew Window

For lossless resequencing:

$$D \geq R_{\text{in,peak}} \times \Delta_{\text{max}} + \text{margin}$$

$$D \geq \sup_t [A(t) - A(t - \Delta_{\text{max}})] + \text{margin}$$

$D$ : required Ring-CAM hit depth

$R_{\text{in,peak}}$ : worst-case accepted aggregate rate over the skew window, not long-term average rate.

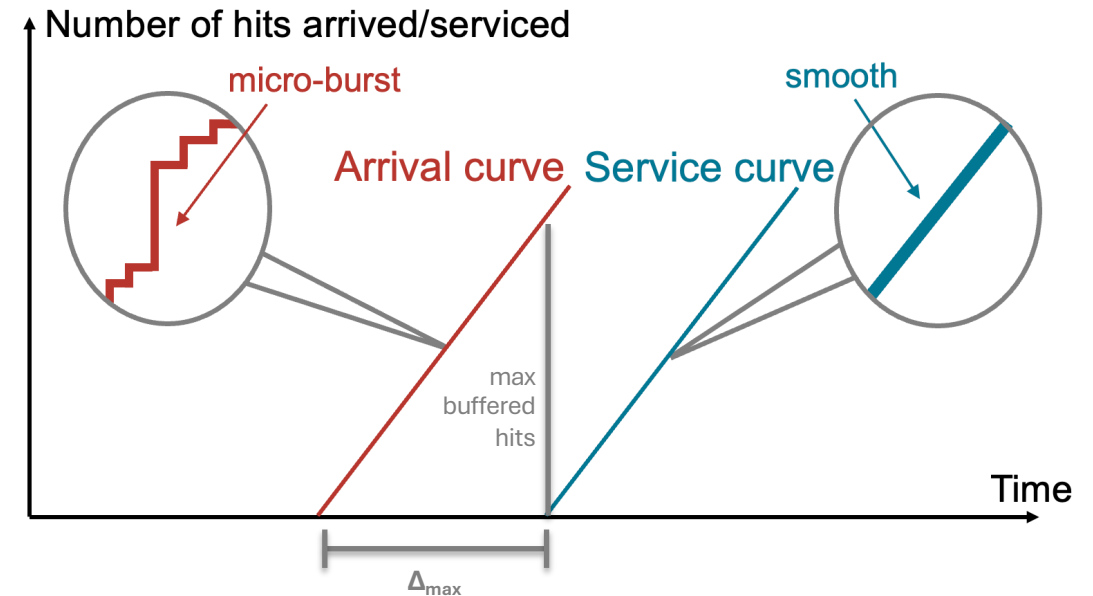
$\Delta_{\text{max}}$ : max lane-arrival skew for same timestamp

margin: clocking, burst granularity, etc

$A(t)$ : cumulative accepted input-hit arrivals.

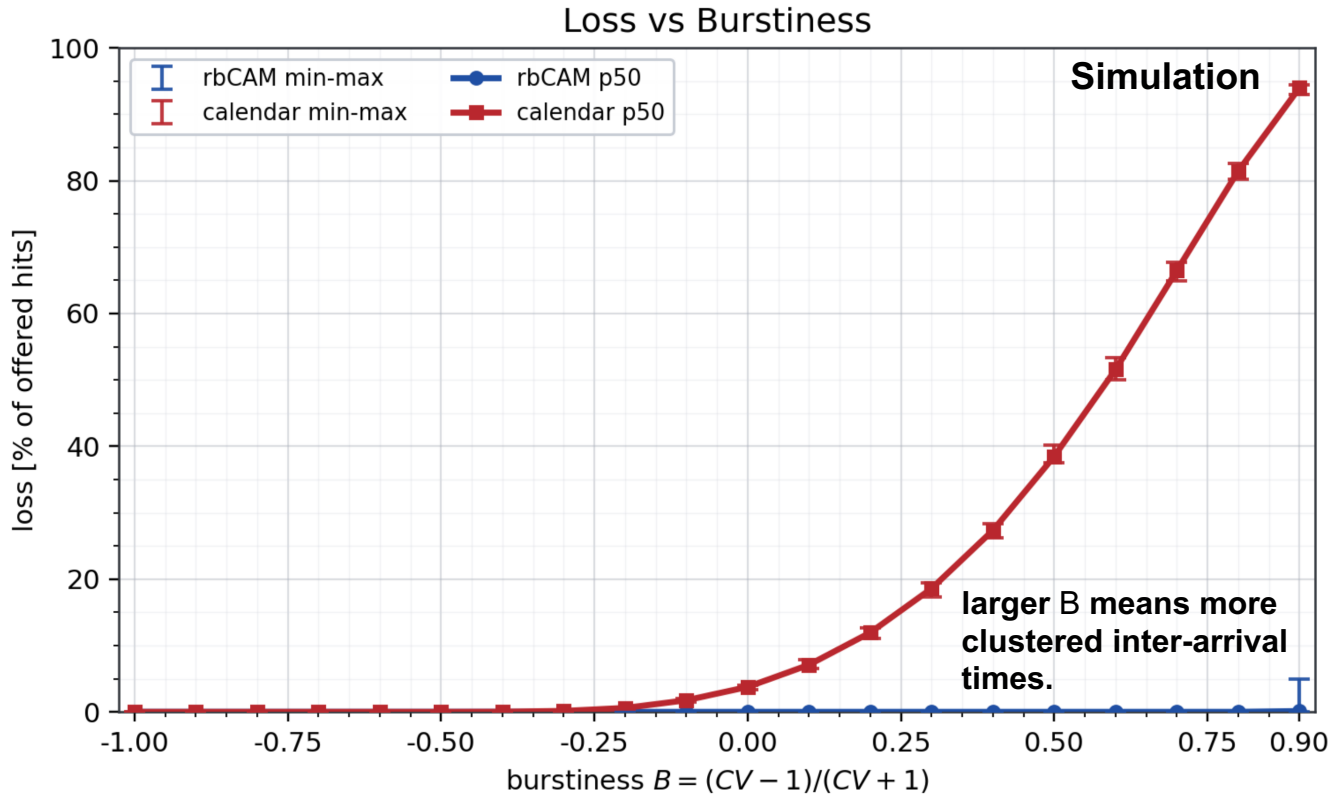
**Depth is set by the aggregate number of accepted hits that can accumulate during the skew window, not by fixed per-timestamp bucket occupancy.**

**Assumption:** output service capacity is sufficient; the depth rule is for skew-induced backlog, not persistent overload.



**Hardware interpretation: one shared buffer absorbs the skew window.**

# Benchmark 1: Ring-CAM Stays Lossless Under Bursts



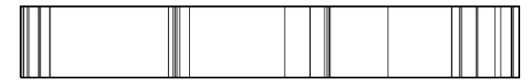
SV-ingress rbCAM vs calendar 256x4; both depth=1024; push=0.125 hit/cyc (64 KHz)

**[Burstiness Formal Definition]**  
<https://arxiv.org/abs/physics/0610233>

Calendar queue loss rises with burstiness; Ring-CAM remains at zero loss in this configured sweep.

## Example of Burstiness (B)

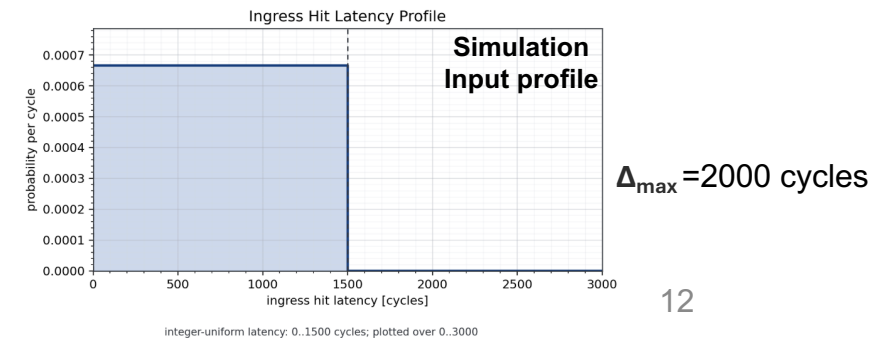
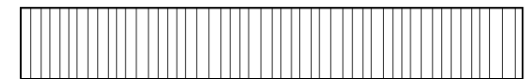
$B = 1$



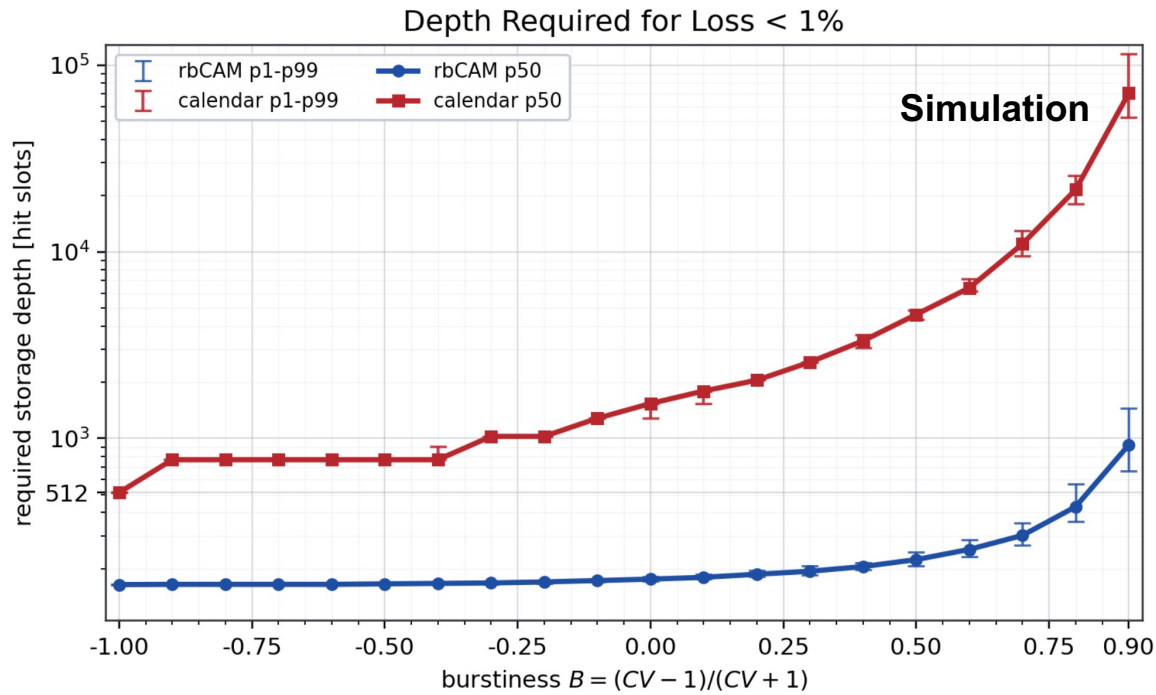
$B = 0$  (Poisson)



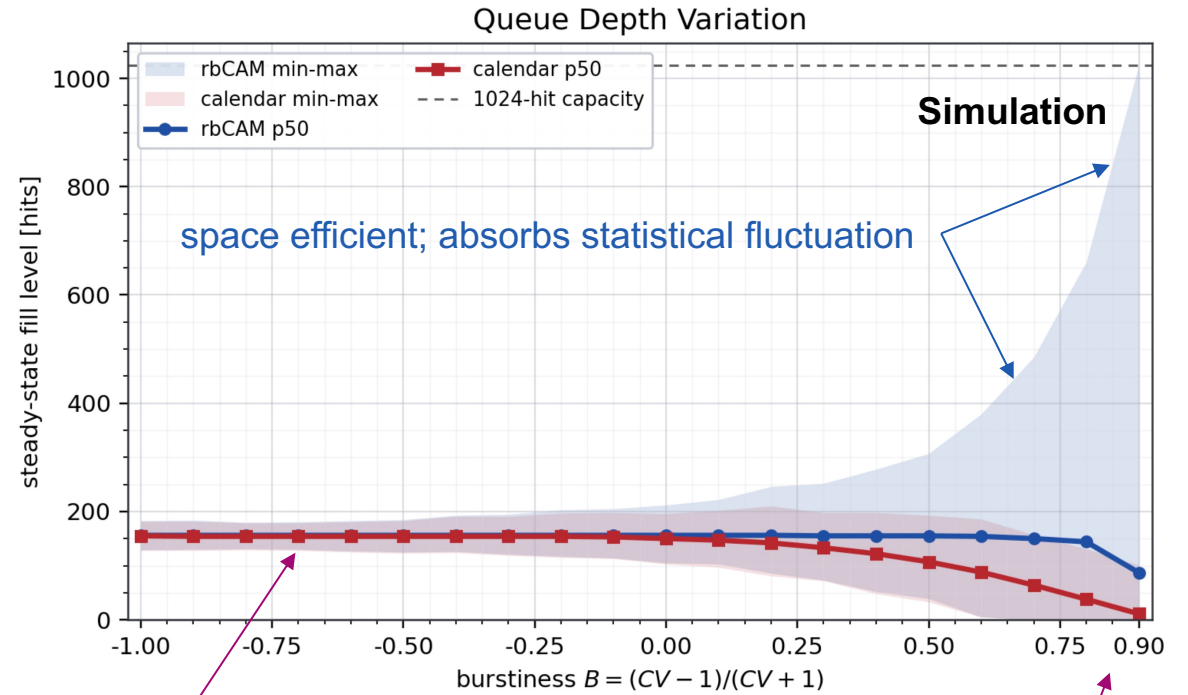
$B = -1$  (periodic)



## Benchmark 2: Ring-CAM Needs Less Depth Under Bursts



rbCAM: SV-ingress ideal TLM; calendar: 256 bins x cap; push=0.125 hit/cycle; bars=p1-p99



input=2 hits/bin, push=0.125 hit/cyc, mean residence=1250 cyc; bands=min-max

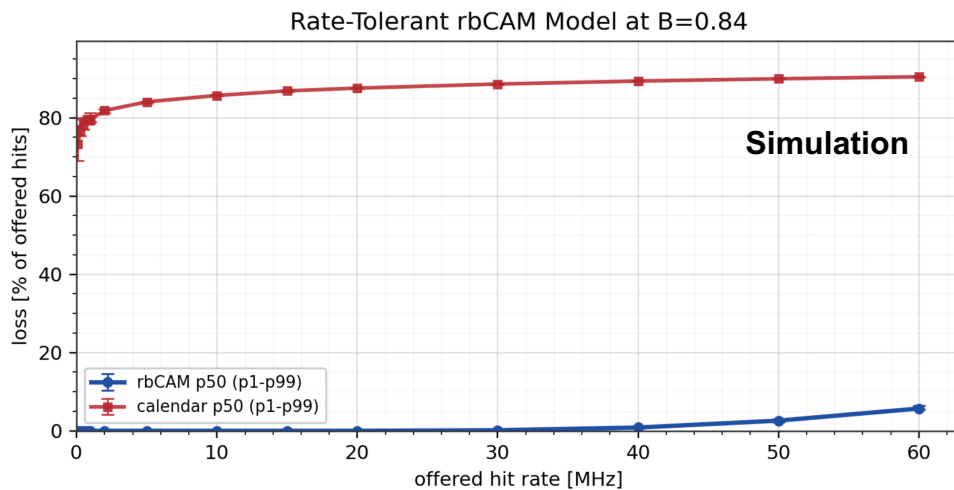
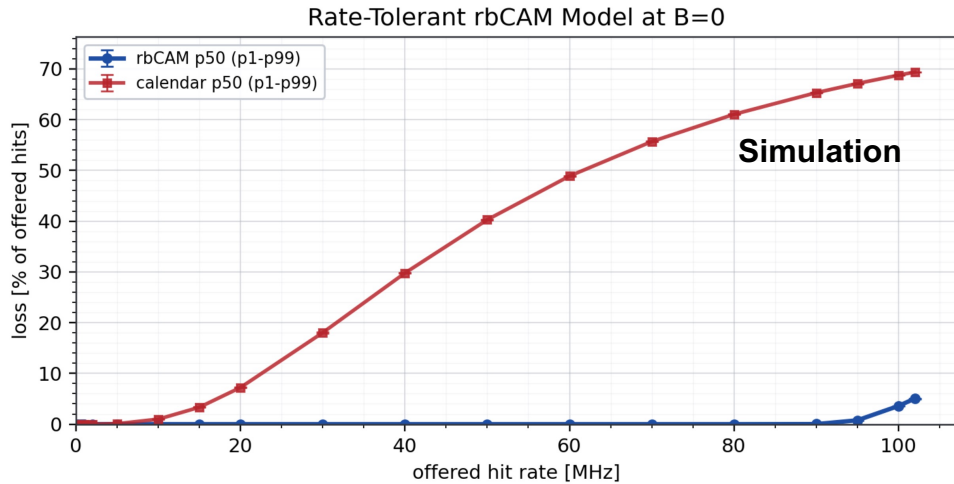
Wastes space

always overflows at high burstiness

Ring-CAM shares storage across timestamps; calendar queues strand memory in empty bins

# Benchmark 3: Full-Rate Pressure Test

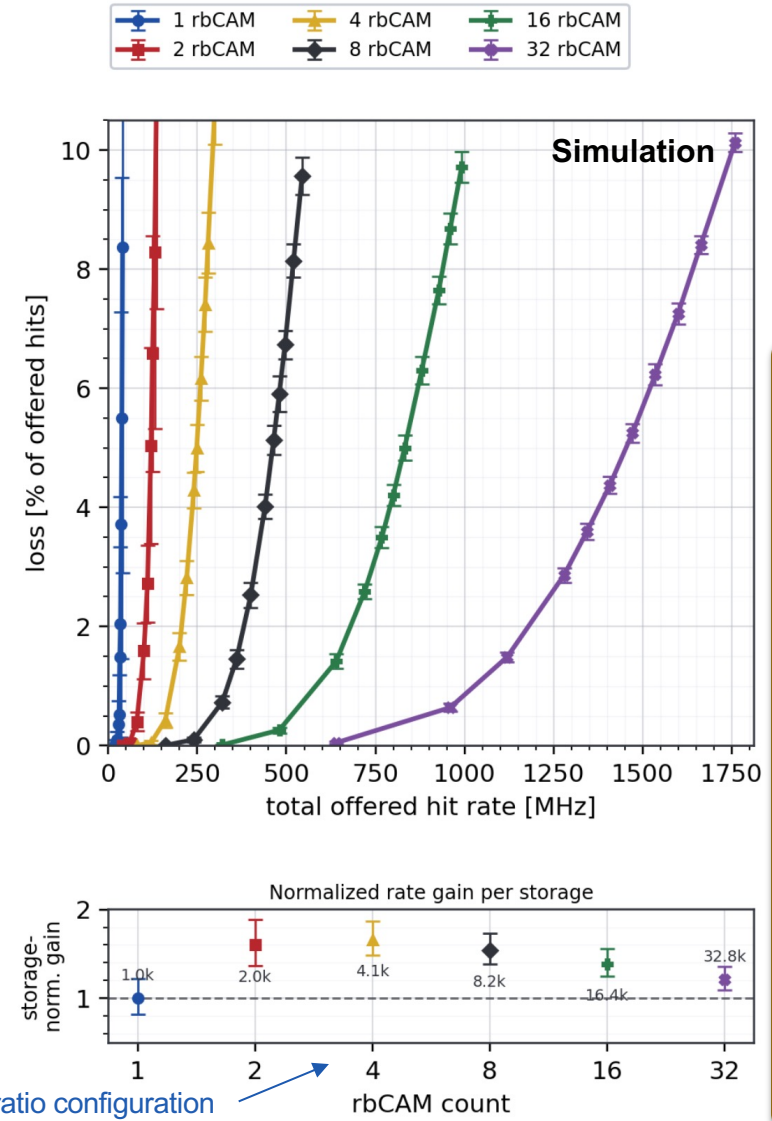
- Rate bottleneck of ring-CAM is **only** in the pop engine
- **Load balancing** with simple routing thanks to its burstiness tolerance



**Single ring-CAM:**  
near-lossless until service-rate limit

**Calendar queue:**  
bucket overflow under bursty events

rbCAM Stack at B=0.84



**Stack ring-CAM:**  
throughput scales with stack count with **even better** efficiency of space.

**Stacking:**  
K Ring-CAMs run in parallel; hits are routed by timestamp interleaving/mod; a final scheduler preserves timestamp order.

highest throughput resource ratio configuration

upper: 1024 slots/rbCAM; lower normalized by total storage N x depth Gaussian mu +/- 1 sigma over 80 seeded trials; lower bars propagated at 5%

# FPGA Implementation: Fully Synthesizable and Timing-Closed

**Sizing**

Search Key Width: 8 bits  
 Subheaders Per Frame: 128  
 Ring Depth: 1024  
 Side Data Width: 31 bits

**Derived sizing**

Global depth: 1024 entries  
 Subheaders per frame: 128  
 Partition depth: 256 entries per encoder slice  
 Entry address width: 10 bits  
 Partition address width: 8 bits  
 Interleaving bits consumed from timestamp: 2  
 Byte-aligned CAM storage: 8192 bits  
 Side-RAM storage: 40960 bits

**Delivered Footprint**

Standalone resource and timing summary  
 Exact measured profile: SV P4 depth 1024, DEBUG=0 on Arria V 5AGXBA7D4F31C5.

**Fitter resources**

ALMs: 2920  
 Registers: 3614  
 RAM blocks: 35  
 Block memory bits: 271872

**Timing at nominal 125 MHz**

Slow 85C setup WNS: 0.917 ns  
 Slow 0C setup WNS: 1.185 ns  
 Worst reported hold slack: 0.150 ns  
 Slow-corner Fmax from the 125 MHz fit: 141.18 MHz

**Resource efficient:**  
 shown configuration: depth = 1024, timestamp width = 8, stack = 1, target = 125 MHz

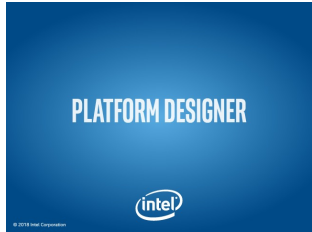
3% ALM, 3% M10K  
 Arria V, 5AGXBA7D4F31C

Auto-pipeline insertion for target clock and configuration

**Resource**  
 ALM 2,920  
 M10K 35

**Timing**  
 Fmax 141.18 MHz  
 F<sub>sign-off</sub> 125 MHz

- Synthesized Ring-CAM IP
- Configuration-aware resource and timing estimate
- Easy to integrate

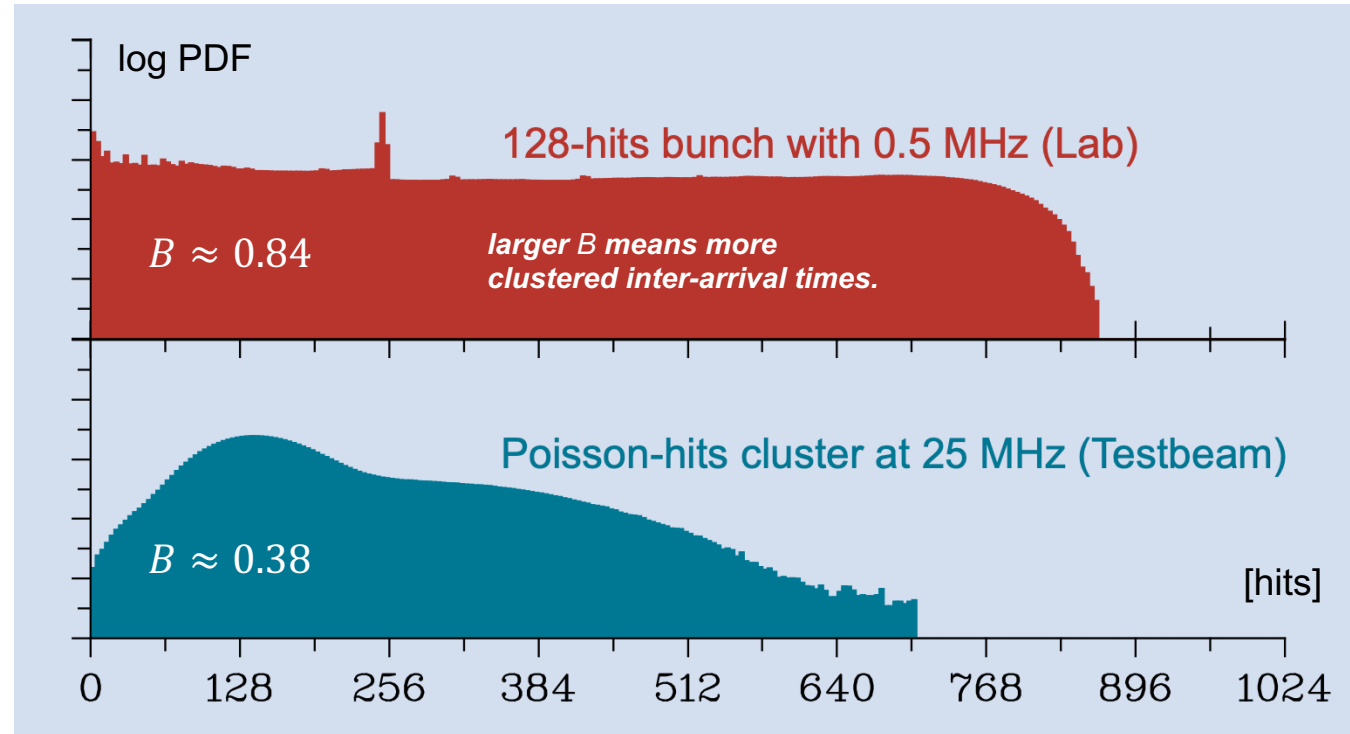


Source code available:  
 RTL Ring-CAM reseencer  
 VHDL and SystemVerilog versions

## Hardware-in-the-Loop Testbeam: 1024-Hit Depth Remains Lossless

- Ring-CAM depth is set by the aggregate accepted hits accumulated during the skew window; **shared storage avoids fixed per-timestamp bucket overflow.**
- In PSI PiM1 / Mu3e SciFi testbeam online monitoring, fill level always stays below the configured 1024-hit depth.
- Configured service capacity is 25 MHz per stack; measured fill-level distributions remain well below the 1024-hit depth, **leaving large overflow margin.**

$$D \geq R_{\text{in,peak}} \times \Delta_{\text{max}} + \text{margin} \rightarrow \text{lossless}$$



Real physics data from PSI PiM1 testbeam with Mu3e SciFi detector

## Verification Signoff: No Lost, Ghost, or Misordered Hits

### buckets of scenarios

status	bucket	catalog_planned	promoted
✓	BASIC	158	158
⚠	EDGE	140	140
⚠	PROF	137	137
⚠	ERROR	138	138

### code coverage

status	metric	pct	target
✓	stmt	96.58	95.0
⚠	branch	89.81	90.0
ℹ	cond	68.90	-
ℹ	expr	82.26	-
✓	fsm_state	100.00	95.0
⚠	fsm_trans	55.56	90.0
✓	toggle	82.11	80.0

### signoff nightly runs

status	run_id	kind	build	seq	txns	cross_pct
✓	CROSS-008	anchored_hybrid	default_p2_pipe4	nightly promoted X117 anchor: GOOD(2048) → ERROR(64) → FLUSH → GOOD(2048)	2048	88.8
✓	CROSS-009	anchored_hybrid	default_p2_pipe4	nightly promoted X118 anchor: GOOD(2048) → TERM → IDLE → RUN_PREPARE → RUN → GOOD(2048)	2048	88.8
✓	CROSS-010	anchored_hybrid	default_p2_pipe4	overwrite-pressure GOOD(pool=1, λ=1.0, 10k) → X119 anchor → recovery GOOD(random, 10k)	784	80.9
✓	CROSS-015	anchored_hybrid	default_p2_pipe4	nightly curated all-bucket mix: B005/B006, E002, P001-style random push-pop, X117-style error+flush recovery, plus overwrite-pressure windows, all separated by random idle gaps	2048	88.8
✓	CROSS-076	seed_sweep	default_p2_pipe4	nightly hotspot overwrite soak derived from P111, 131072 txn same-ts pressure	131072	88.8

- ✓ Functional and code coverage
- ✓ Nightly soak runs with real physics stimuli
- Caveats: ⚠ = waived unreachable or configuration-infeasible coverage bins.

**Signoff focus:** no ghost hits, no lost hits, timestamp order preserved.

**lost** = input hit not emitted;  
**ghost** = output hit not in input;  
**misordered** = timestamp decreases at output

## Summary: Bounded-Skew Timestamp Resequencing

**Main claim:** reset lane skew locally in FPGA so triggerless DAQ can deliver complete timestamp-ordered streams to the farm.

### Problem

Fan-in creates micro-bursts  
FIFO delay time differs by lane  
The farm waits for the slowest lane (**tail latency**)

### Ring-CAM mechanism

Write hits in arrival order  
Index timestamp -> slot  
list/bitmap  
Emit all slots for earliest ready  
timestamp

### Measured evidence

Simulated burst sweeps: lossless  
for Ring-CAM  
Resource/timing and UVM signoff  
complete  
Physics testbeam proven  
lossless and promising margin

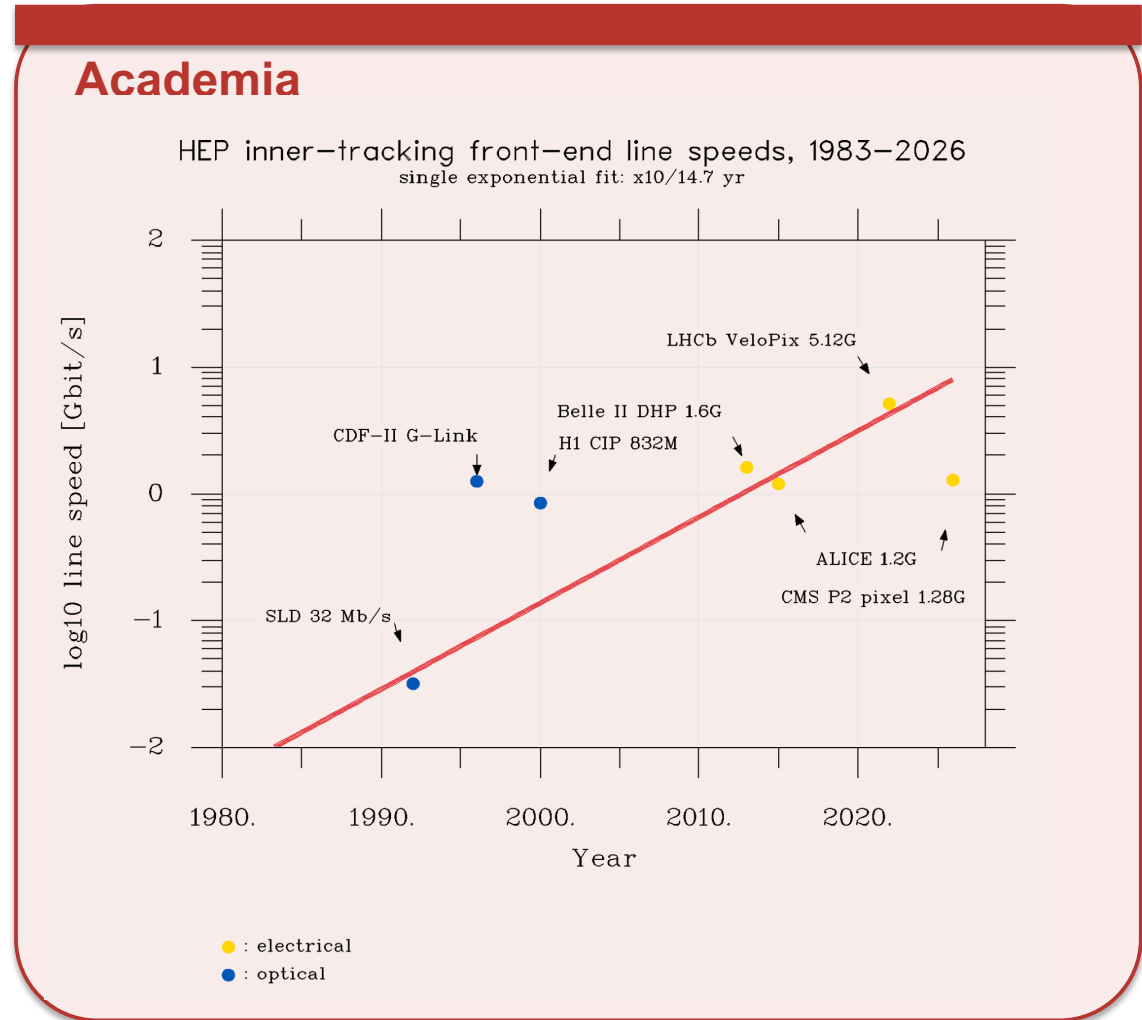
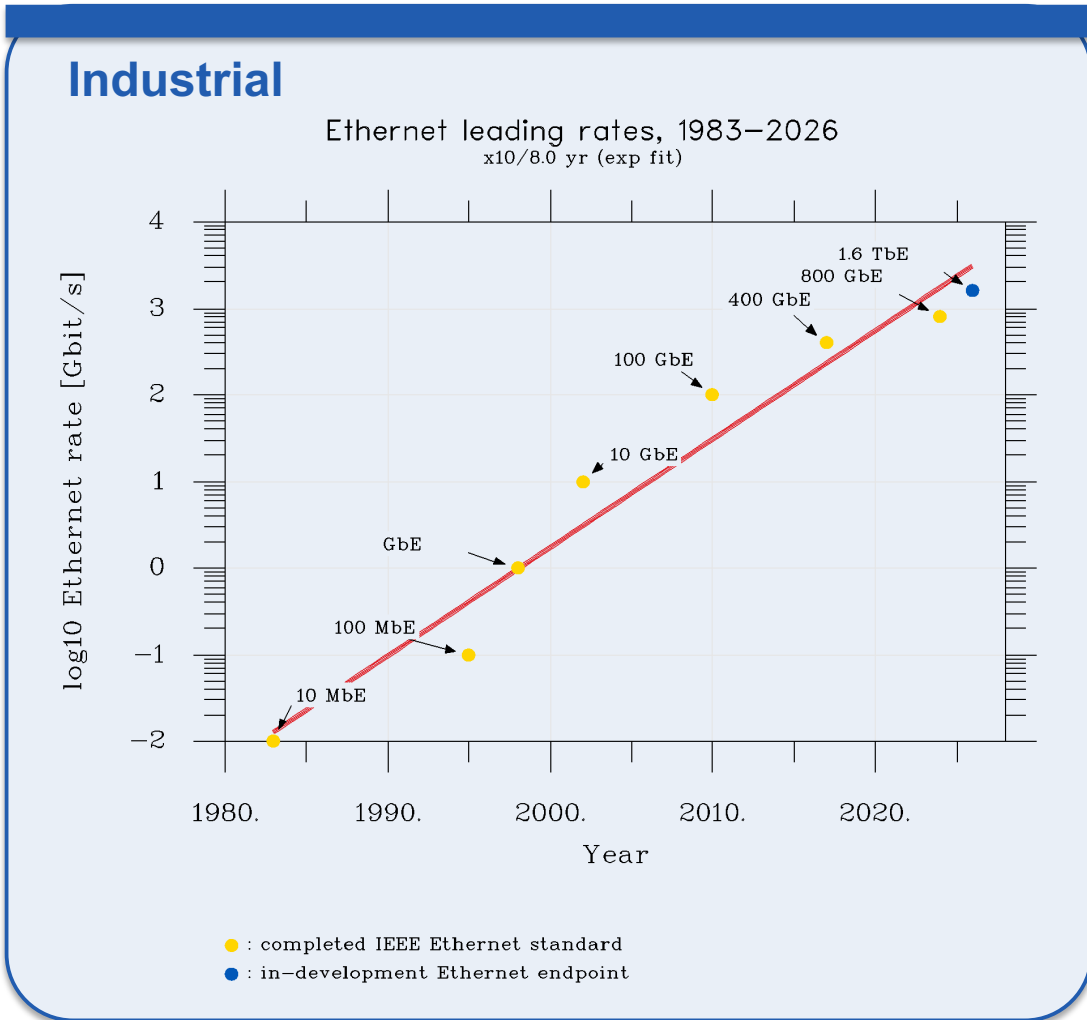
**Design rule:**  $D \geq R_{in,peak} \times \Delta_{max} + \text{margin}$

**Contribution:** a bounded-skew triggerless DAQ architecture plus validated Ring-CAM FPGA resequecer IP.

# BACKUP

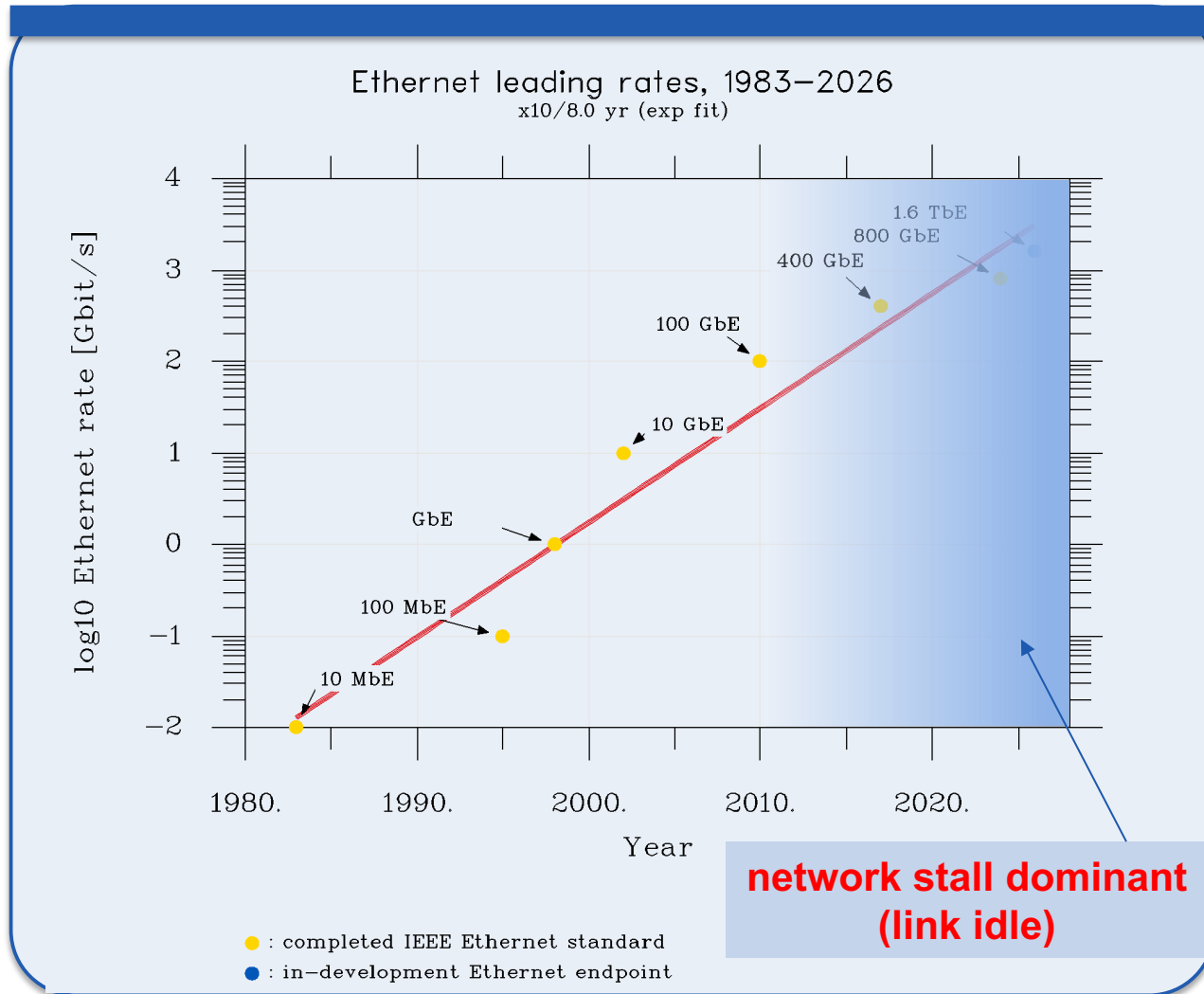
Main talk ends here. **Thank you for your attention.**  
Following slides are for technical questions.

# Readout Bandwidth Is Scaling Up



Faster links reduce the average-bandwidth problem, but do not remove timestamp synchronization and lane-skew closure

Different domain, same synchronization pattern: progress is gated by the slowest participant



External analogy: synchronization waits for the slowest transfer

The duration of each communications round is determined by the **slowest transfer... Network failures...** become expensive in lost GPU time.

[OpenAI/Microsoft]  
<https://cdn.openai.com/pdf/resilient-ai-supercomputer-networking-using-mrc-and-srv6.pdf>

All-to-All communication can account for **50%-60% of the total time** ... the primary communication bottleneck arises from the **mandatory data synchronizations (network)**.

[USENIX]  
<https://dl.acm.org/doi/abs/10.5555/3768039.3768101>

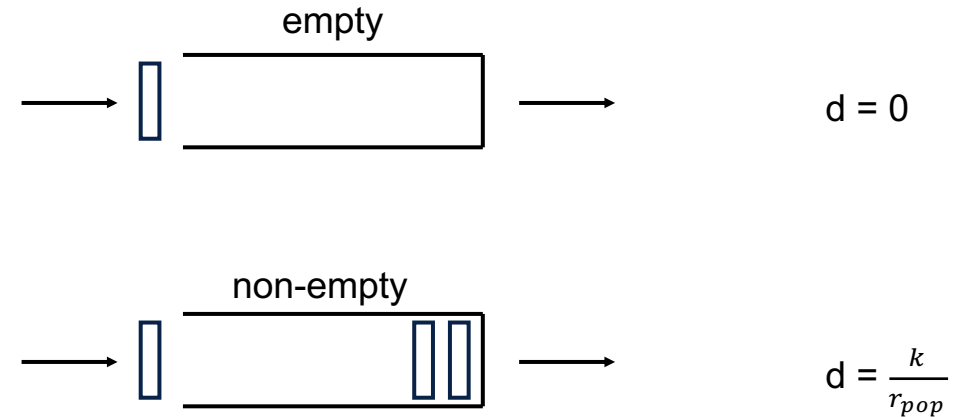
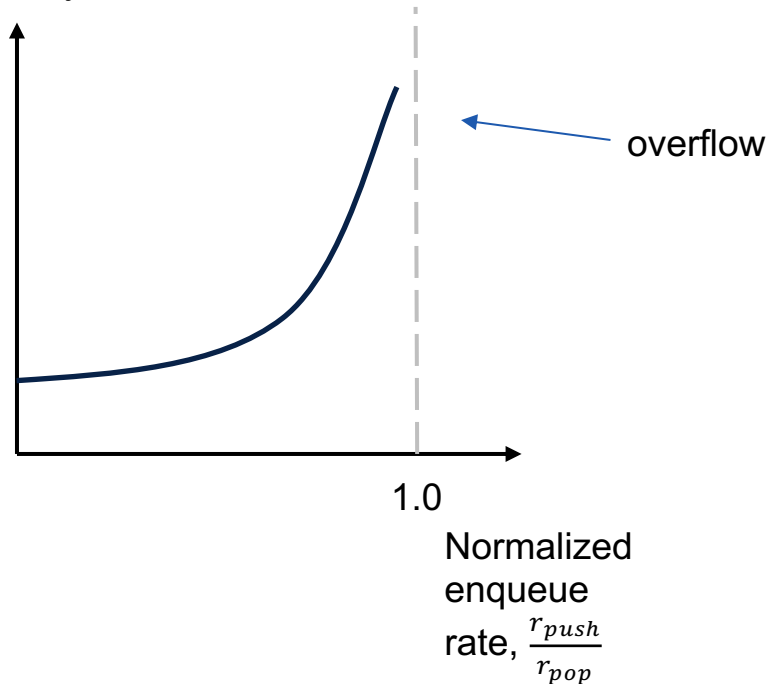
HEP analogue: fan-in waits for a complete time slice

# Backup: FIFO Residence Time Varies with Occupancy

FIFO residence-time intuition / Work-conserving FIFO model

M/M/1/∞ queue

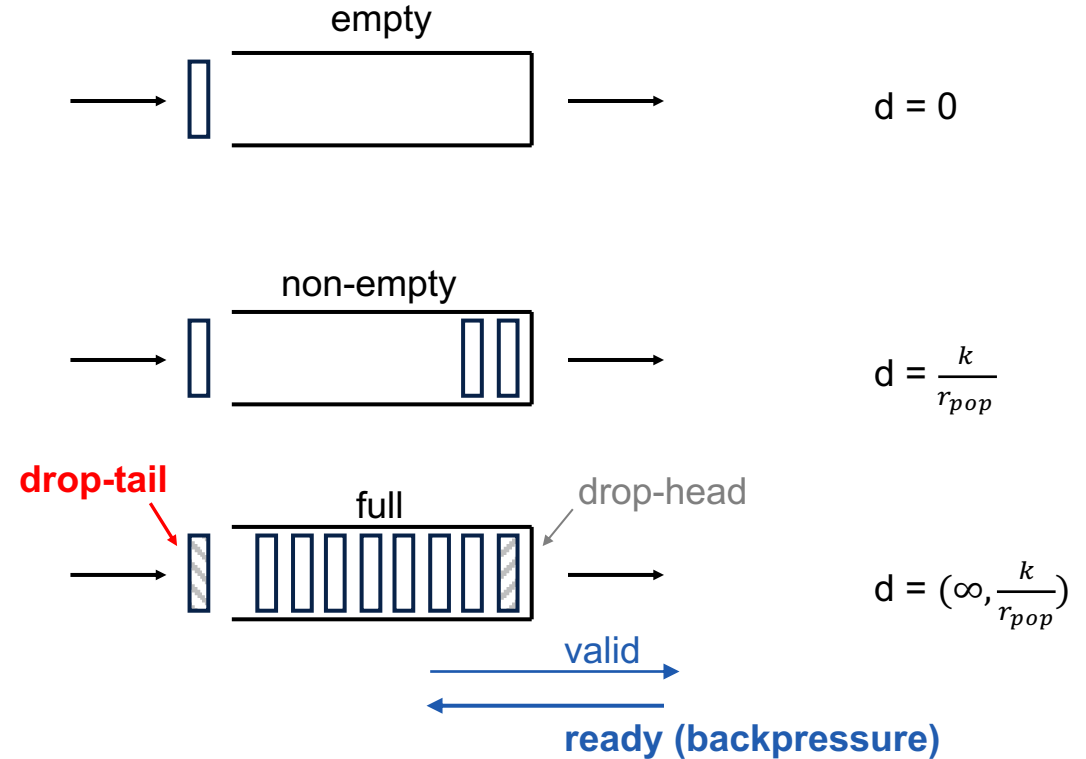
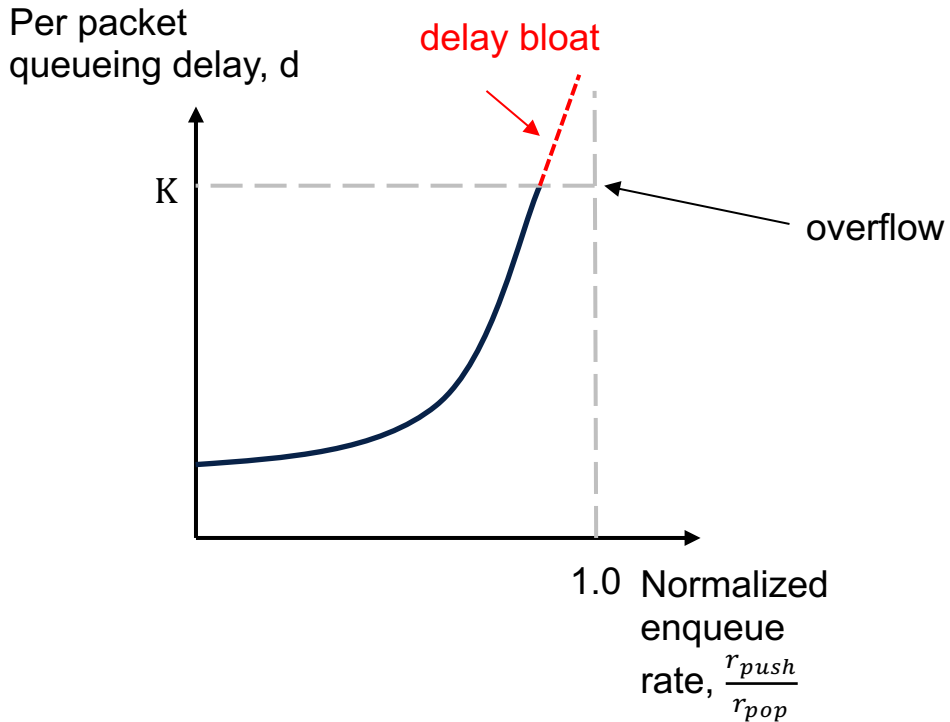
Per packet queueing delay,  $d$



- At each merge stage, FIFOs absorb short bursts
- ⚠ FIFO occupancy becomes variable residence time
- ⚠ More depth prevents immediate loss but increases worst-lane delay

# Backup: Backpressure Converts Burst Pressure into Delay

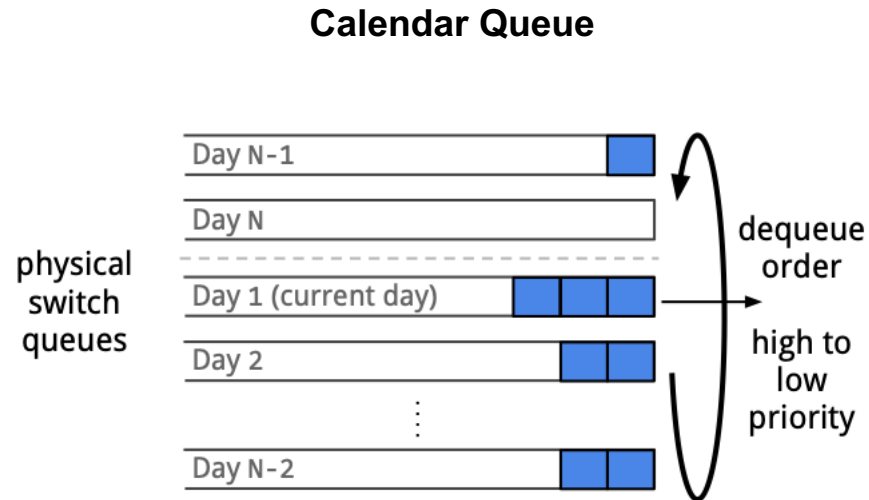
generic (with backpressure)  
 M/G/1/K queue  
 infinite depth



- ⚠ Backpressure avoids immediate loss but can grow residence time
- ⚠ Because congestion is local, lanes accumulate different delay

FIFO residence-time spread becomes lane skew

## Backup: Fixed-capacity hardware bucketed sorter



Sort hits by timestamp using fixed buckets

### Workflow

- Fixed bucket size and running admission window
- Write hits to the bucket associated with timestamp
- Read hits from the earliest bucket

### Drawback

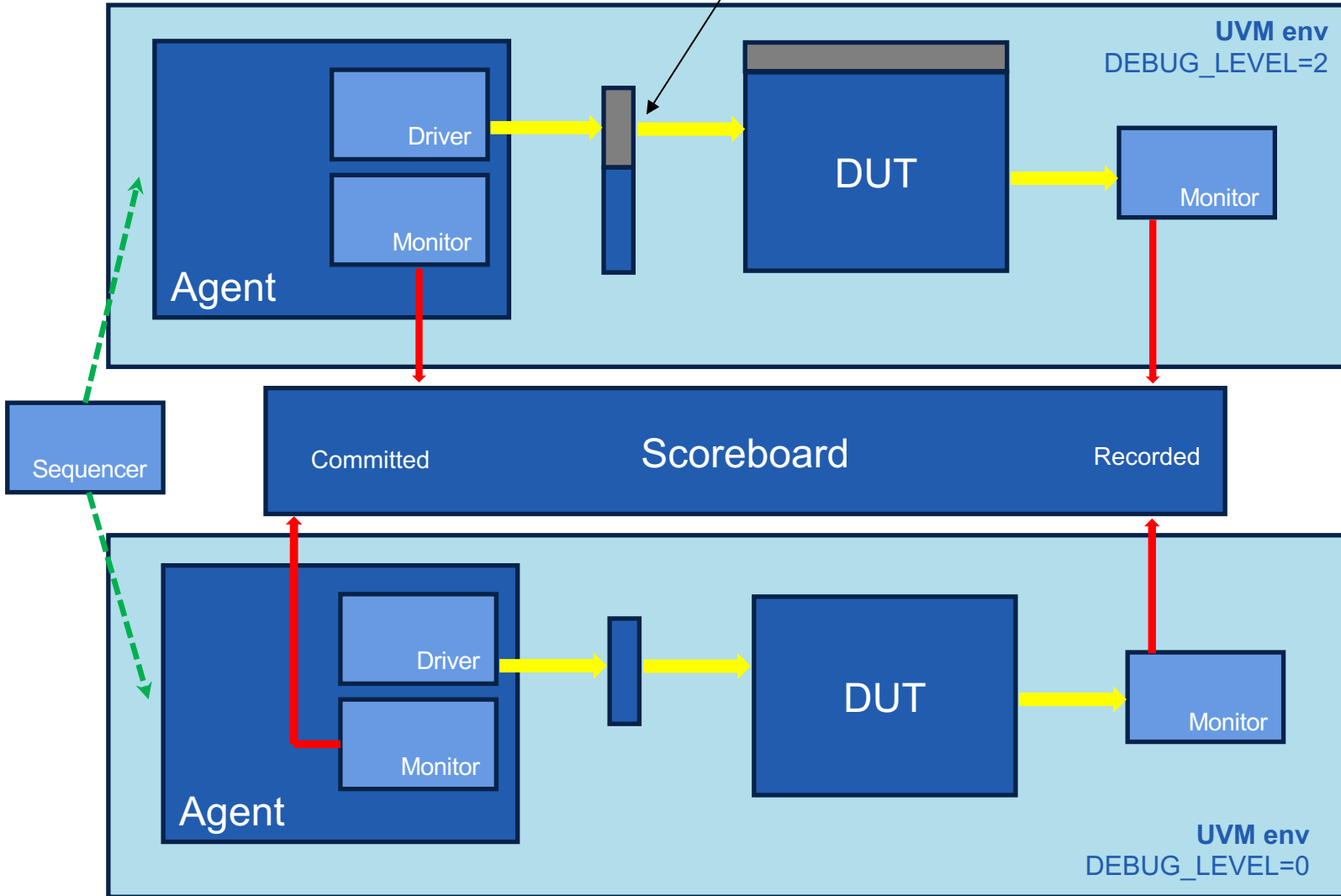
- Fixed bucket size can waste storage under bursty occupancy

[NSDI] <https://www.usenix.org/system/files/nsdi20-paper-sharma.pdf>

[NSDI] <https://www.usenix.org/conference/nsdi25/presentation/alcoz>

# Backup: Detailed UVM Regression Verification

Meta info per hit :  
 • UID  
 • True Timestamp



**Challenge:** complex gray-box DUT with memory; underflow/overflow risks.

**Mitigation:** dual UVM environments, per-hit UID tracking, directed regressions, randomized soak tests with physics stimuli.

# Resequencer IP-core

Ring-buffer CAM  
ring\_buffer\_cam Details

Configuration Identity Interfaces Register Map

**Delivered Profile**

**Catalog revision**  
This release is packaged as 26.2.13.0516.

**Packaging provenance**  
Default git stamp 0x23CE513F (23ce513f). Git describe: v26.2.6.0422-25-g23ce513f.

**Delivered interface contract**  
The packaged image keeps the established hit\_type1, hit\_type2, run\_control, and filllevel interface names so existing Platform Designer systems can be upgraded in place while picking

**Timing status**  
The packaged SystemVerilog payload currently reports 145.52 MHz slow-corner Fmax with 1523 ALMs, 2016 registers, and 19 RAM blocks, closing the required 137.5 MHz standalone

**Versioning**

**Common identity header**  
Word 0 is UID.  
Word 1 is META: write 0=VERSION, 1=DATE, 2=GIT, 3=INSTANCE\_ID.

**VERSION encoding**  
VERSION[31:24] = MAJOR, VERSION[23:16] = MINOR, VERSION[15:12] = PATCH, VERSION[11:0] = BUILD.

**Packaged defaults**  
Default VERSION\_GIT = 0x23CE513F (23ce513f). Git describe = v26.2.6.0422-25-g23ce513f.

**Editability**  
IP\_UID and INSTANCE\_ID remain integration-editable. Version, build, date, and git provenance fields are locked to the packaged image.

UID:

Version Major:

Version Minor:

Version Patch:

Build Stamp:

Version Date:

Git Stamp:

Instance ID:

**Four version-identification fields:  
release, Git hash,  
date, build**



# Resequencer IP-core

Running monitoring with control plane

counter freeze = 1

Reading counters

Reading counters

...

counter freeze = 0

Ring-buffer CAM  
ring\_buffer\_cam

Configuration Identity Interfaces Register Map

CSR Window

Word	Byte	Name	Access	Description
0x00	0x000	UID	RO	Software-visible IP identifier. Default is ASCII <b>RBCM</b> but it is integration-time overridable.
0x01	0x004	META	RW/RO	Read-multiplexed metadata word. Write 0=VERSION, 1=DATE, 2=GIT, 3=INSTANCE_ID. VERSION is packed as MAJOR[31:24]
0x02	0x008	CTRL	RW	Bit 0 <b>go</b> , bit 1 <b>soft_reset</b> , bit 4 <b>filter_iner</b> , bit 5 <b>counter_freeze</b> . Writing counter_freeze=1 snapshots all 64-bit diagnostic coun
0x03	0x00C	EXPECTED_LATENCY	RW	Read-pointer delay target in cycles. Reset default is 2000.
0x04	0x010	FILL_LEVEL	RO	Live fill-level estimate derived from push, pop, and overwrite counters.
0x05	0x014	INERR_COUNT_LO	RO	Low word of filtered ingress timestamp-error hit count.
0x06	0x018	PUSH_COUNT_LO	RO	Low word of accepted push operations.
0x07	0x01C	POP_COUNT_LO	RO	Low word of drained hits.
0x08	0x020	OVERWRITE_COUNT_LO	RO	Low word of overwrite events.
0x09	0x024	CACHE_MISS_COUNT_LO	RO	Low word of cache-miss / empty-search events.
0x0A	0x028	INERR_COUNT_HI	RO	High word of filtered ingress timestamp-error hit count.
0x0B	0x02C	PUSH_COUNT_HI	RO	High word of accepted push operations.
0x0C	0x030	POP_COUNT_HI	RO	High word of drained hits.
0x0D	0x034	OVERWRITE_COUNT_HI	RO	High word of overwrite events.
0x0E	0x038	CACHE_MISS_COUNT_HI	RO	High word of cache-miss / empty-search events.
0x0F	0x03C	DEASM_FULL_DROP_LO	RO	Low word of ingress/deassembly FIFO full-induced drop observations.
0x10	0x040	DEASM_FULL_DROP_HI	RO	High word of ingress/deassembly FIFO full-induced drop observations.
0x11	0x044	POP_CMD_FULL_DROP_LO	RO	Low word of pop-command FIFO full observations.
0x12	0x048	POP_CMD_FULL_DROP_HI	RO	High word of pop-command FIFO full observations.
0x13	0x04C	EGRESS_NOT_READY_DROP_LO	RO	Low word of egress hit words generated while the downstream sink was not ready.
0x14	0x050	EGRESS_NOT_READY_DROP_HI	RO	High word of egress hit words generated while the downstream sink was not ready.

- CSR (control and status register) read back of counters for observability



# Resequencer System Integration

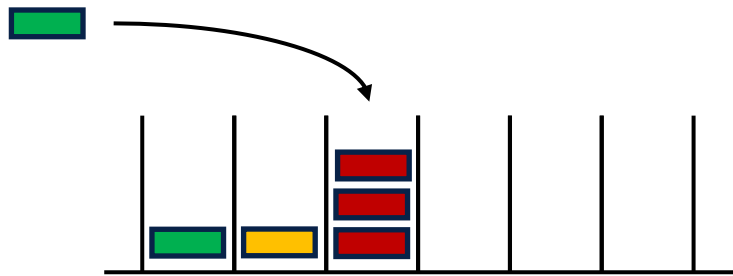


Name	Description
datapath_clock_0	Clock Source
clk_in	Clock Input
clk_in_reset	Reset Input
clk	Clock Output
clk_reset	Reset Output
ring_buffer_cam_0	Ring-buffer CAM
csr	Avalon Memory Mapped Slave
hit_type1	Avalon Streaming Sink
hit_type2	Avalon Streaming Source
clock_interface	Clock Input
reset_interface	Reset Input
run_control	Avalon Streaming Sink
ring_buffer_cam_1	Ring-buffer CAM
csr	Avalon Memory Mapped Slave
hit_type1	Avalon Streaming Sink
hit_type2	Avalon Streaming Source
clock_interface	Clock Input
reset_interface	Reset Input
run_control	Avalon Streaming Sink
ring_buffer_cam_2	Ring-buffer CAM
csr	Avalon Memory Mapped Slave
hit_type1	Avalon Streaming Sink
hit_type2	Avalon Streaming Source
clock_interface	Clock Input
reset_interface	Reset Input
run_control	Avalon Streaming Sink
ring_buffer_cam_3	Ring-buffer CAM
csr	Avalon Memory Mapped Slave
hit_type1	Avalon Streaming Sink
hit_type2	Avalon Streaming Source
clock_interface	Clock Input
reset_interface	Reset Input
run_control	Avalon Streaming Sink
data_splitter_0	Avalon-ST Splitter
clk	Clock Input
reset	Reset Input
in	Avalon Streaming Sink
out0	Avalon Streaming Source
out1	Avalon Streaming Source
out2	Avalon Streaming Source
out3	Avalon Streaming Source
feb_frame_assembly_0	FEB (Front-end Board) Frame Asse...
hit_type2_0	Avalon Streaming Sink

**Reduces integration effort from custom RTL work to GUI-level system integration**

# Sorting algo

- **Calendar queue**       $\Leftarrow$  state-of-the-art
- Ring-CAM



Workflow:

- Write to bin addressed by TS

Write:

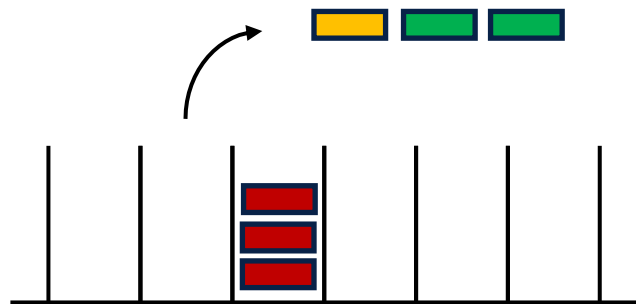
- $\text{Hash}(\text{TS}) = \text{address}$

# Sorting algo

- **Calendar queue**       $\Leftarrow$  state-of-the-art
- Ring-CAM

Workflow:

- Write to bin based on TS
- Read from bin sequentially



Read:

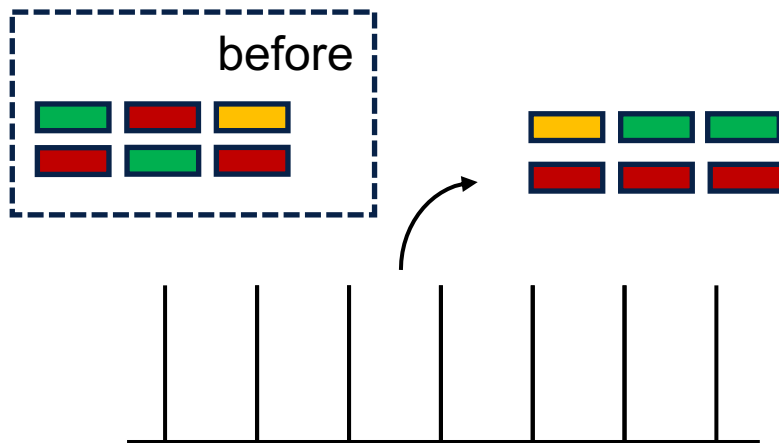
- If (empty), then address = address + 1

# Sorting algo

- **Calendar queue**       $\Leftarrow$  state-of-the-art
- Ring-CAM

Workflow:

- Write to bin based on TS
- Read from bin sequentially



Read:

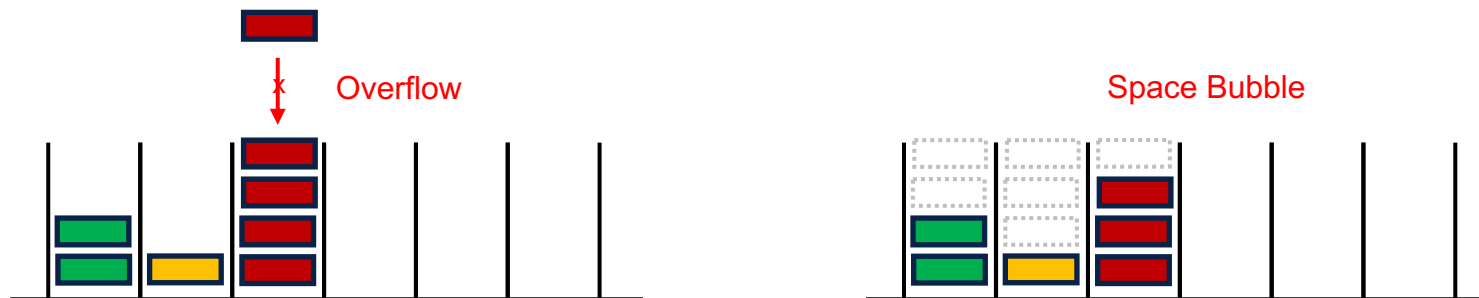
- If (empty), then address = address + 1

# Sorting algo

- **Calendar queue**       $\Leftarrow$  state-of-the-art
- Ring-CAM

## Problem of calendar queue:

- Fixed bin size, not  $O(1)$  space

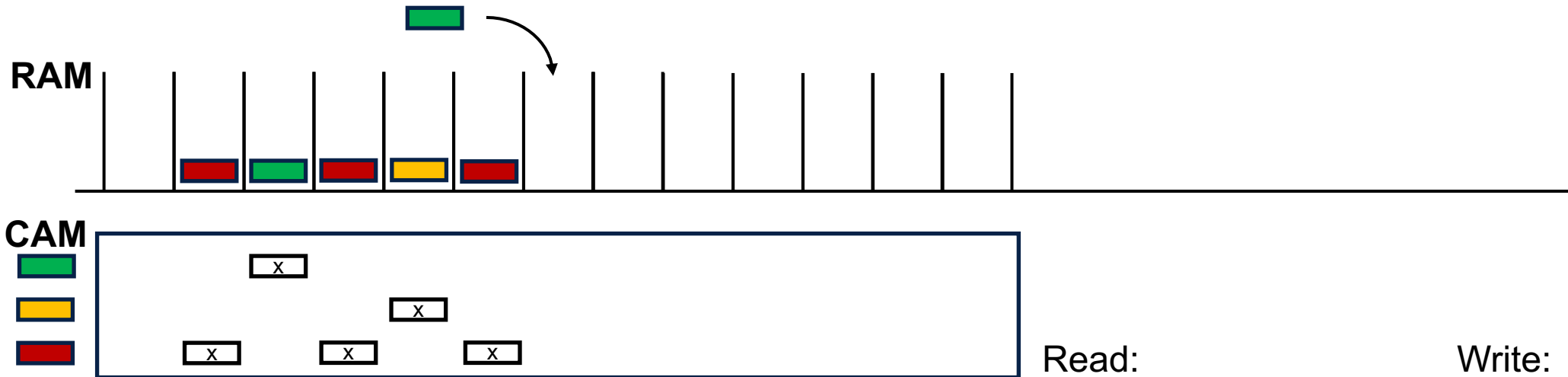


# Sorting algo

- Calendar queue  $\Leftarrow$  state-of-the-art
- **Ring-CAM**

Workflow:

- Write sequentially



Read:

- $\text{RAM}(\text{address}) = \text{data}$
- $\text{CAM}(\text{data}) = \text{address}$

Write:

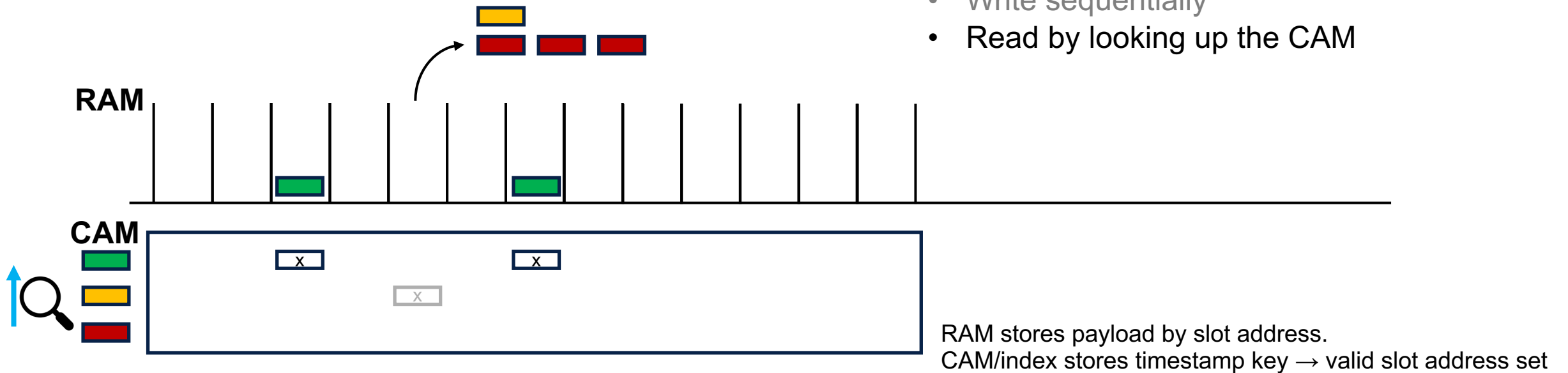
- $\text{RAM}(\text{address}) = \text{data}$
- $\text{CAM}(\text{address}) = \text{data}$

# Sorting algo

- Calendar queue  $\Leftarrow$  state-of-the-art
- **Ring-CAM**

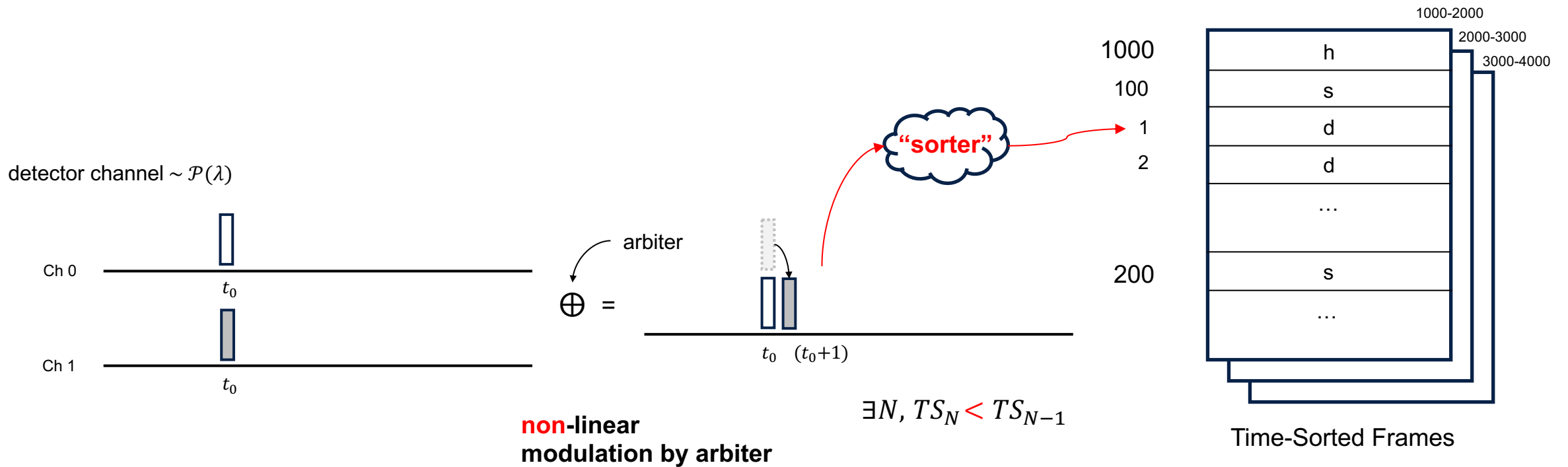
Workflow:

- Write sequentially
- Read by looking up the CAM



# Problem formulation

Q: What ruins the event **ordering**?



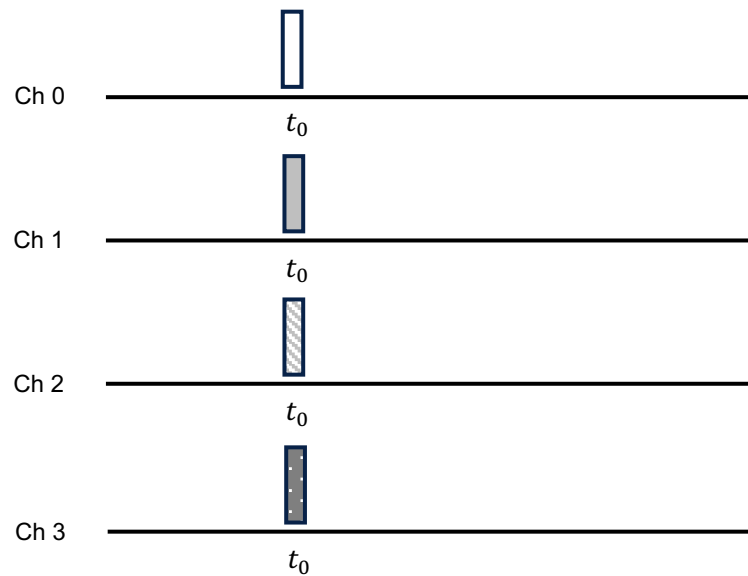
Timestamp of hit :=  $h + s + d$  (sum of offsets)

## Problem formulation (cont.)

Q: What ruins the event **ordering**? (cont.)

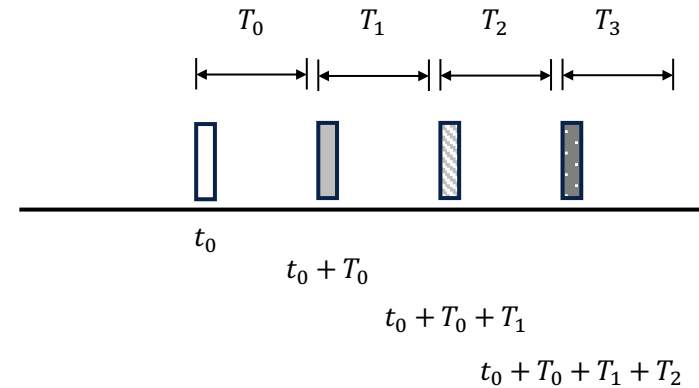
Arbiter choice 1: **time multiplexed**

detector channel  $\sim \mathcal{P}(\lambda)$



time-multiplexed arbiter

$\oplus =$



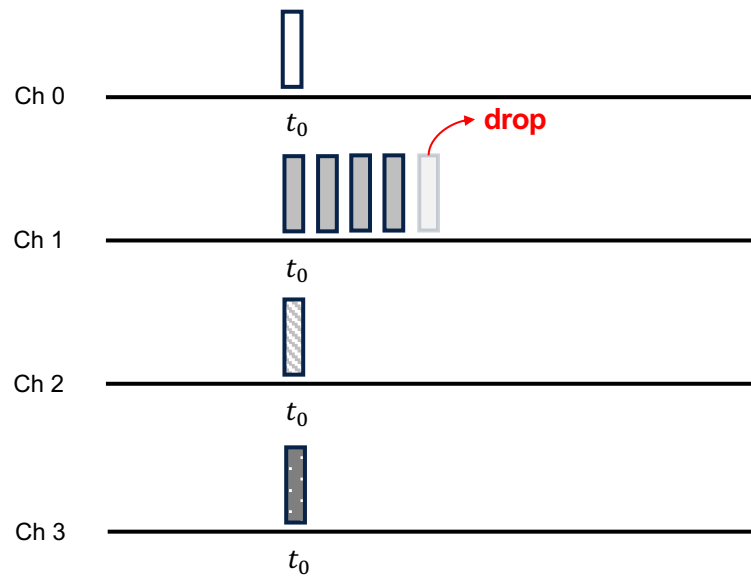
Example: JESD 204B/C protocol for ADC/DAC

# Problem formulation (cont.)

Q: What ruins the event **ordering**? (cont.)

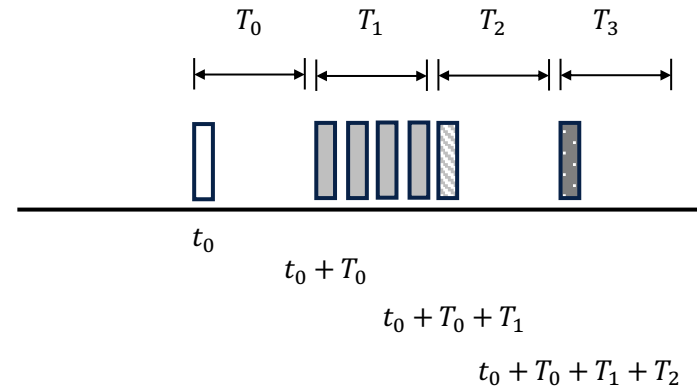
➤ Q: What is the issue with time-multiplexing? – **fixed bandwidth** and deterministic latency

detector channel  $\sim \mathcal{P}(\lambda)$



not suitable for event-driven

time-multiplexed arbiter  
 $\oplus =$



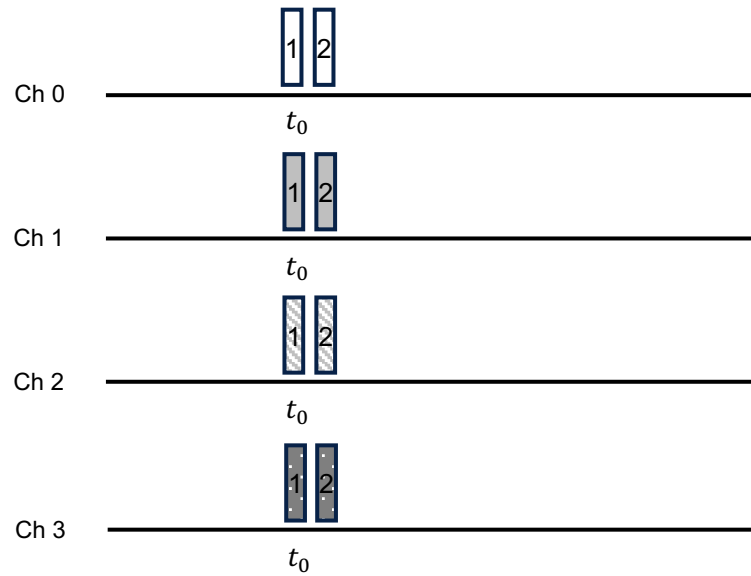
Example: JESD 204B/C protocol for ADC/DAC

# Problem formulation (cont.)

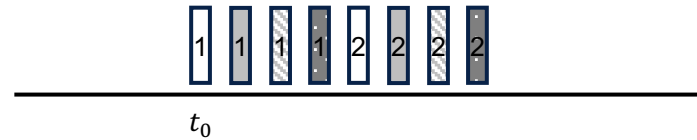
Q: What ruins the event **ordering**? (cont.)

Arbiter choice 2: **rotating priority** (Round-Robin)

detector channel  $\sim \mathcal{P}(\lambda)$



RR arbiter  
 $\oplus =$



t

t := timestamp

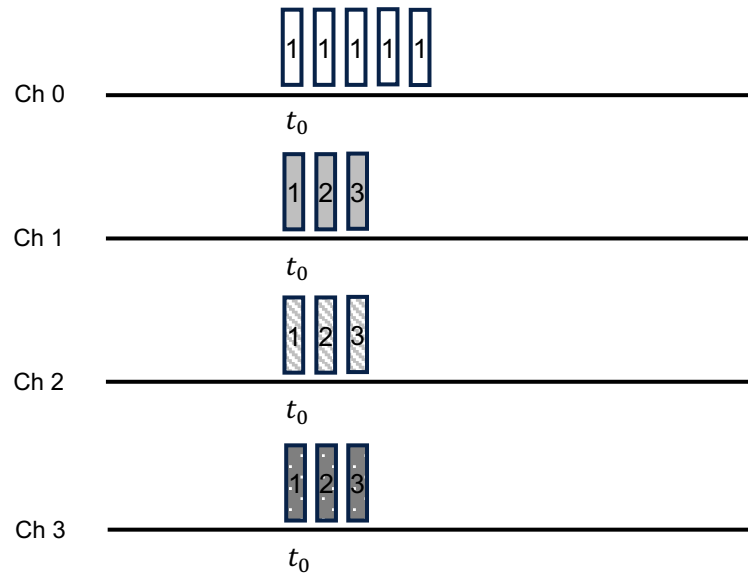
# Problem formulation (cont.)

Q: What ruins the event **ordering**? (cont.)

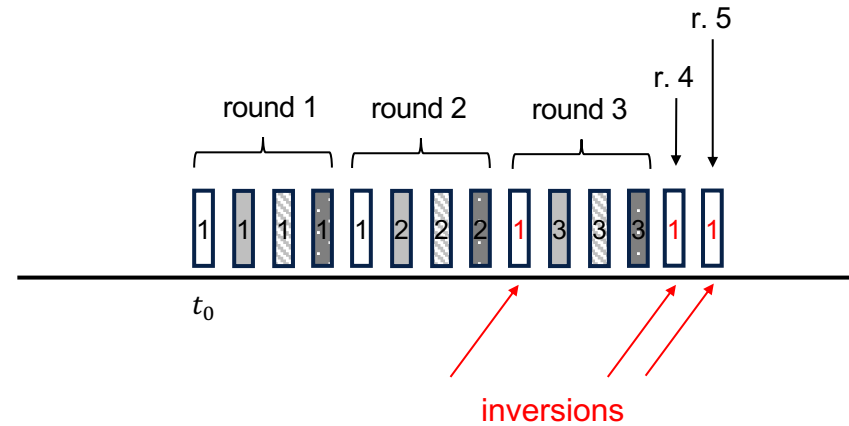
Arbiter choice 2: **rotating priority** (Round-Robin)

➤ Q: What is the issue with RR? – bandwidth “fair” sharing but **inversions**

detector channel  $\sim \mathcal{P}(\lambda)$



RR arbiter  
 $\oplus =$



t

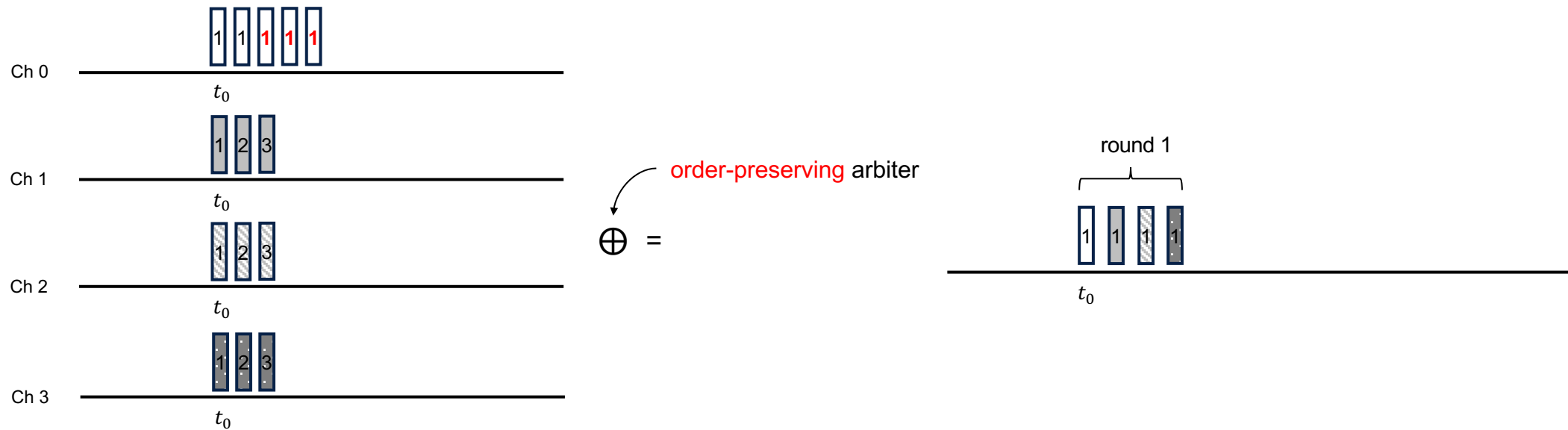
t := timestamp

# Problem formulation (cont.)

Q: What ruins the event **ordering**? (cont.)

Arbiter choice 3: **order-preserving** (granting the smallest Head-of-Line)

detector channel  $\sim \mathcal{P}(\lambda)$



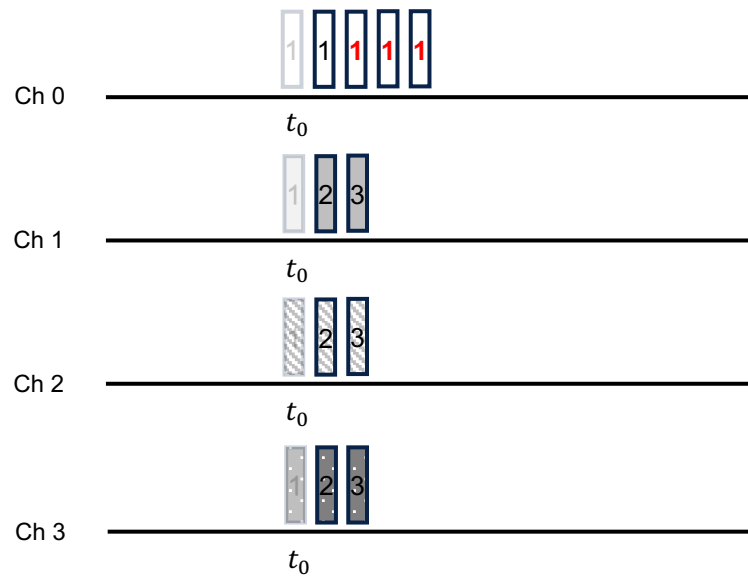
t := timestamp

# Problem formulation (cont.)

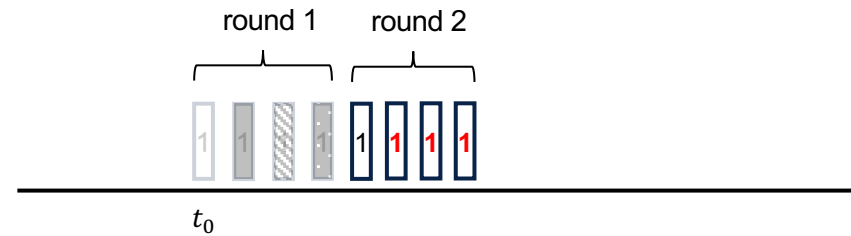
Q: What ruins the event **ordering**? (cont.)

Arbiter choice 3: **order-preserving** (granting the smallest HOL)

detector channel  $\sim \mathcal{P}(\lambda)$



$\oplus =$  **order-preserving** arbiter



t := timestamp

# Problem formulation (cont.)

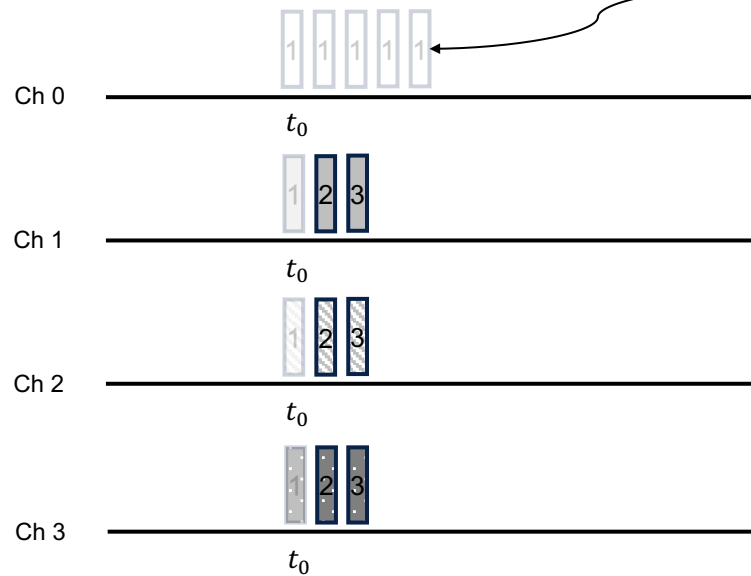
Q: What ruins the event **ordering**? (cont.)

Arbiter choice 3: **order-preserving** (granting the smallest HOL)

➤ Q: What is the issue with order-preserving? – complexity and stall (**delay jitter**)

no showahead timestamp to decide

detector channel  $\sim \mathcal{P}(\lambda)$

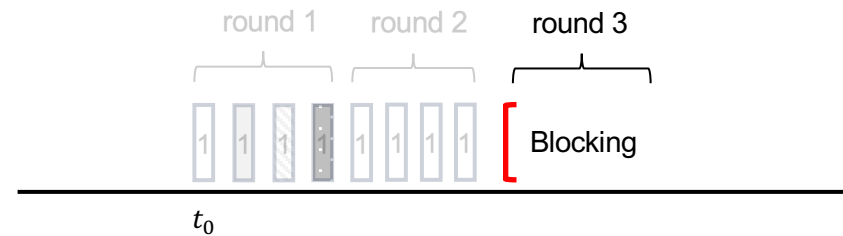


order-preserving arbiter

$\oplus =$

$t$

$t := \text{timestamp}$



Consequence of non-work-conserving:

- unbalanced channel rate induces blocking
- stall time can **accumulate** across stages
- insert beacon signal can help, but merely bound the stall time.
- if packet loss and overflow? **Deadlock** (infinite HOL block)
- **delay jitter depends on the input timestamp distribution**

**Out-of-order arrival is inevitable. No single solution of arbiter can solve the unsortiness problem. Sorter is needed.**

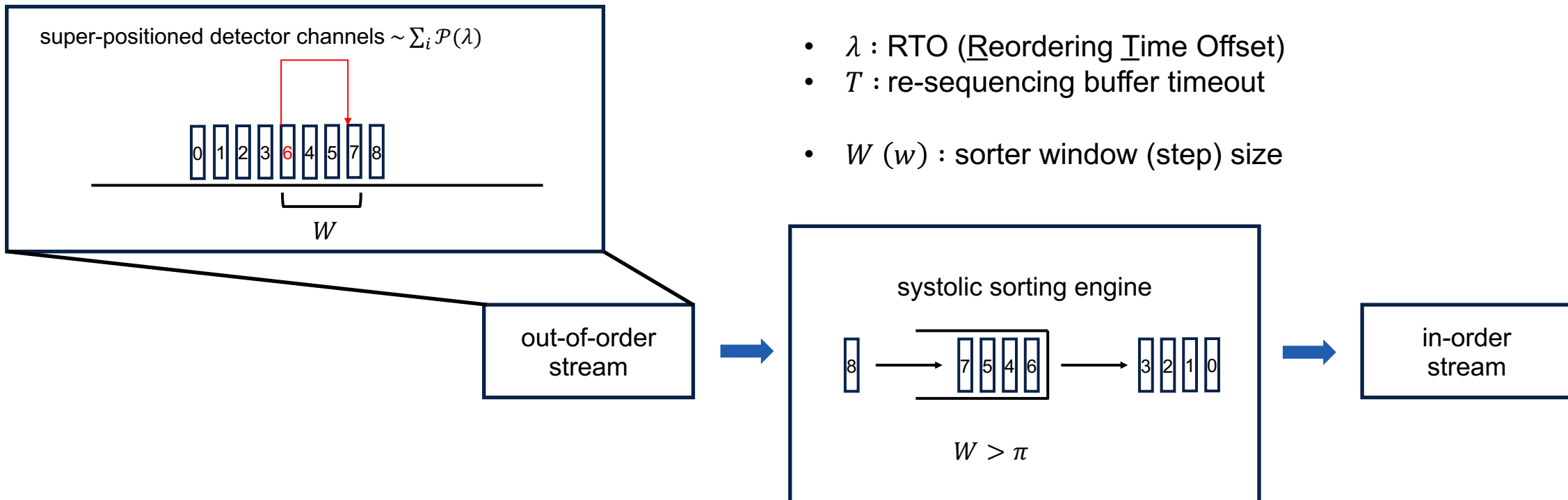
# Re-sequencing buffer (“sorter”) implementation

Q: How to sort?

Sorter choice 1: **systolic array**

How to quantify the **timestamp reordering** in terms of displacement in space and time

- $\pi$  : **RBO** (Reordering Byte Offset) ←
- $B$  : re-sequencing buffer size
- $\lambda$  : **RTO** (Reordering Time Offset)
- $T$  : re-sequencing buffer timeout
- $W$  ( $w$ ) : sorter window (step) size



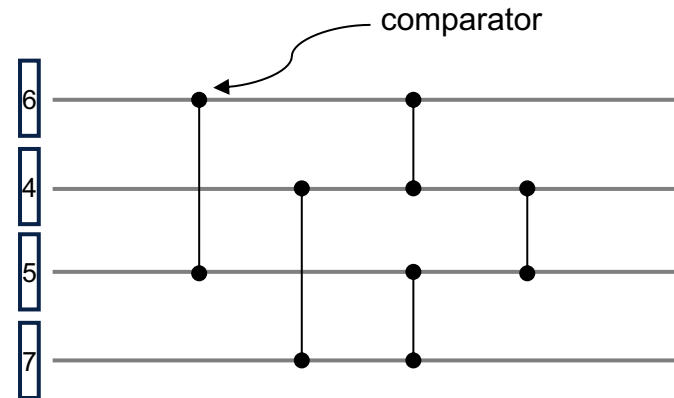
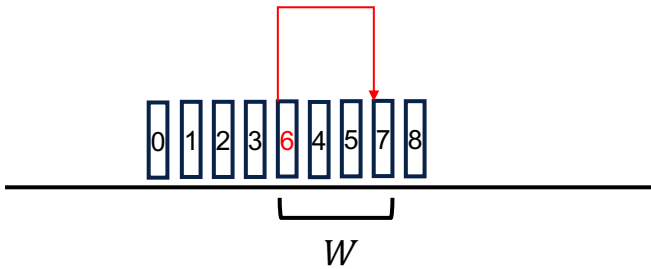
# Re-sequencing buffer (“sorter”) implementation

Q: How to sort?

Sorter choice 1: **systolic array** (cont.)

- $\pi$  : **RBO** (Reordering Byte Offset) ←
- $\lambda$  : RTO (Reordering Time Offset)
- $W (w)$  : sorter window (step) size

super-positioned detector channels  $\sim \sum_i \mathcal{P}(\lambda)$



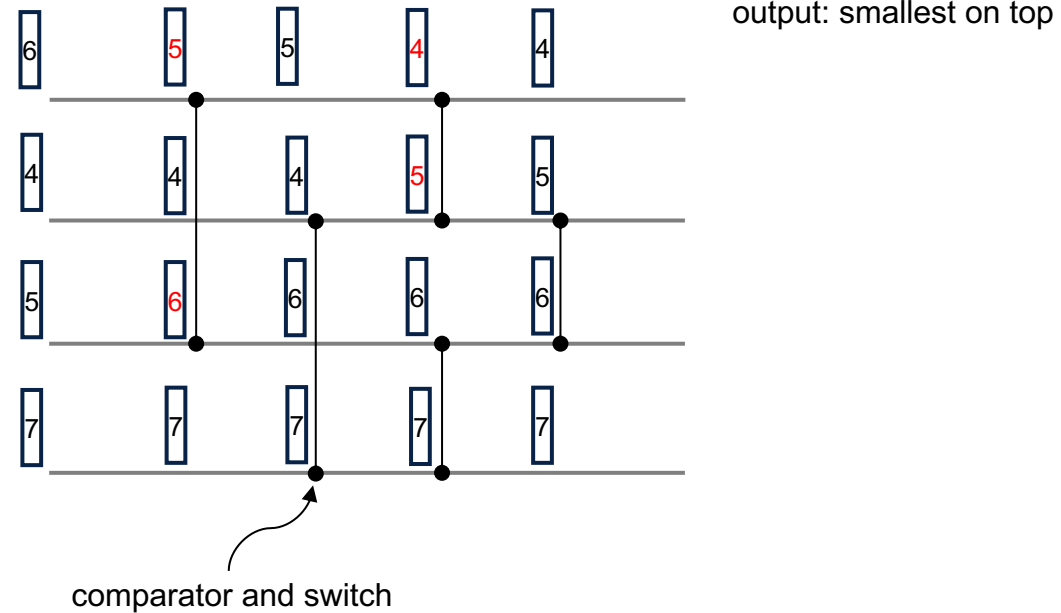
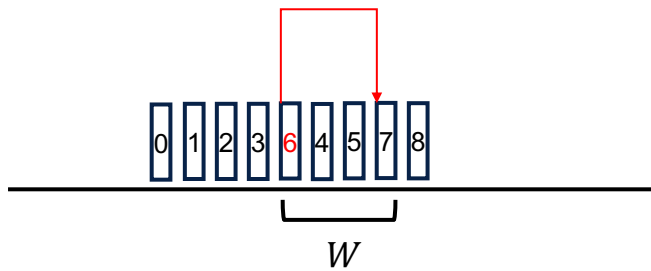
# Re-sequencing buffer (“sorter”) implementation

Q: How to sort?

Sorter choice 1: **systolic array** (cont.)

- $\pi$  : **RBO** (Reordering Byte Offset) ←
- $\lambda$  : RTO (Reordering Time Offset)
- $W(w)$  : sorter window (step) size

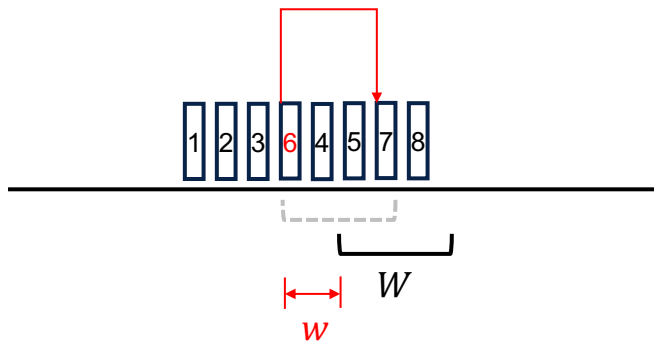
super-positioned detector channels  $\sim \sum_i \mathcal{P}(\lambda)$



# Re-sequencing buffer (“sorter”) implementation

Q: How to sort?

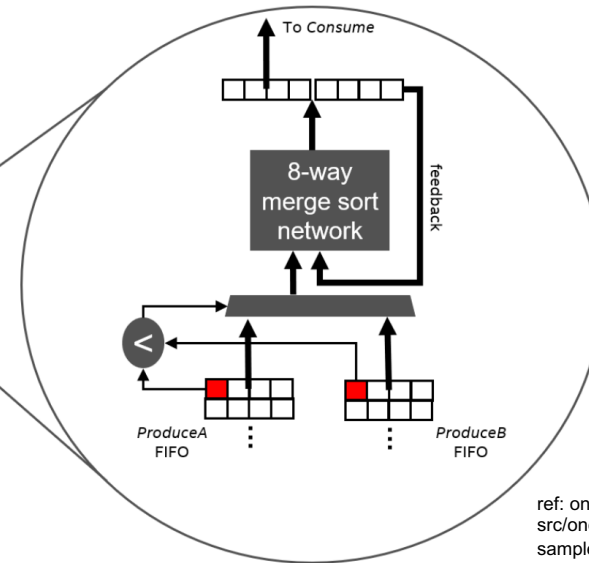
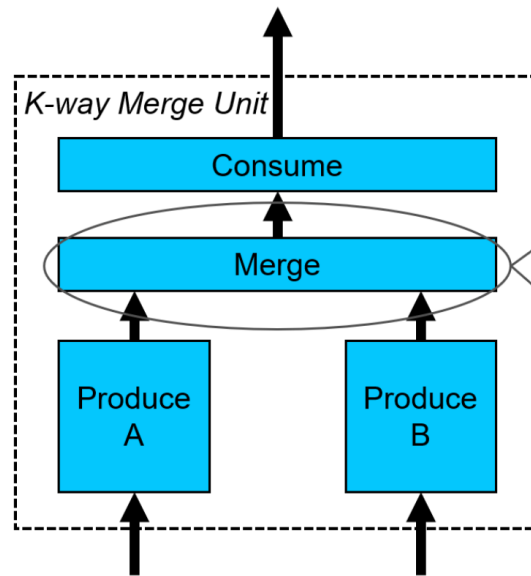
super-positioned detector channels  $\sim \sum_i \mathcal{P}(\lambda)$



Sorter choice 1: **systolic array** (cont.)

➤ Q: What is the issue with systolic array?

- $\alpha$  : arrival curve
- $V$  : delay jitter bound
- $B$  : re-sequencing buffer size
- $T$  : re-sequencing buffer timeout



ref: oneAPI ([https://github.com/oneapi-src/oneAPI-samples/tree/master/DirectProgramming/C++SYCL\\_FPGA/ReferenceDesigns/merge\\_sort](https://github.com/oneapi-src/oneAPI-samples/tree/master/DirectProgramming/C++SYCL_FPGA/ReferenceDesigns/merge_sort))

Require :

$$W > \pi$$

$$w \rightarrow 1$$

RBO bound :

$$\pi < \alpha(V) - 1 \quad \text{Packet arrival in time window of jitter}$$

Buffer size bound :

$$B > \max(\pi, \alpha(V + T)) \quad \text{Packet arrival in time window of jitter}$$

Issues :

- for RR :  $\pi$  (RBO) is **loosely** bounded
- Scalability is poor, space  $O(N^2) - O(N^3)$ ,  $W < 64$
- TS dist. limits  $\min(w)$ , time  $> O(1)$

# Re-sequencing buffer (“sorter”) implementation-2

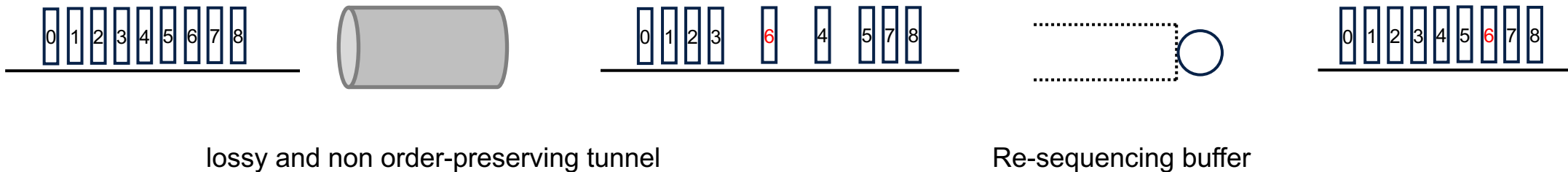
Q: How to sort?

Sorter choice 2: ???

➤ Q: How to bound the RTO?

- $\pi$  : RBO (Reordering Byte Offset)

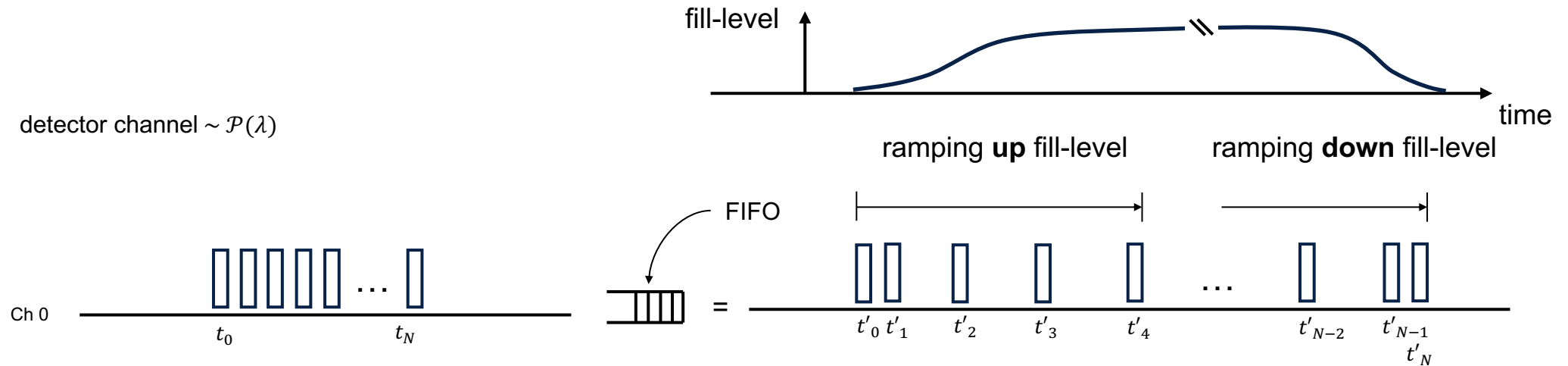
- $\lambda$  : **RTO** (Reordering Time Offset) ←



RTO : time offset of the inverted packet w.r.t its **assumed arrival time**

# Delay of a non memoryless system

**FIFO ruins events arrival-time displacement (delay spread / introduce delay jitter)**



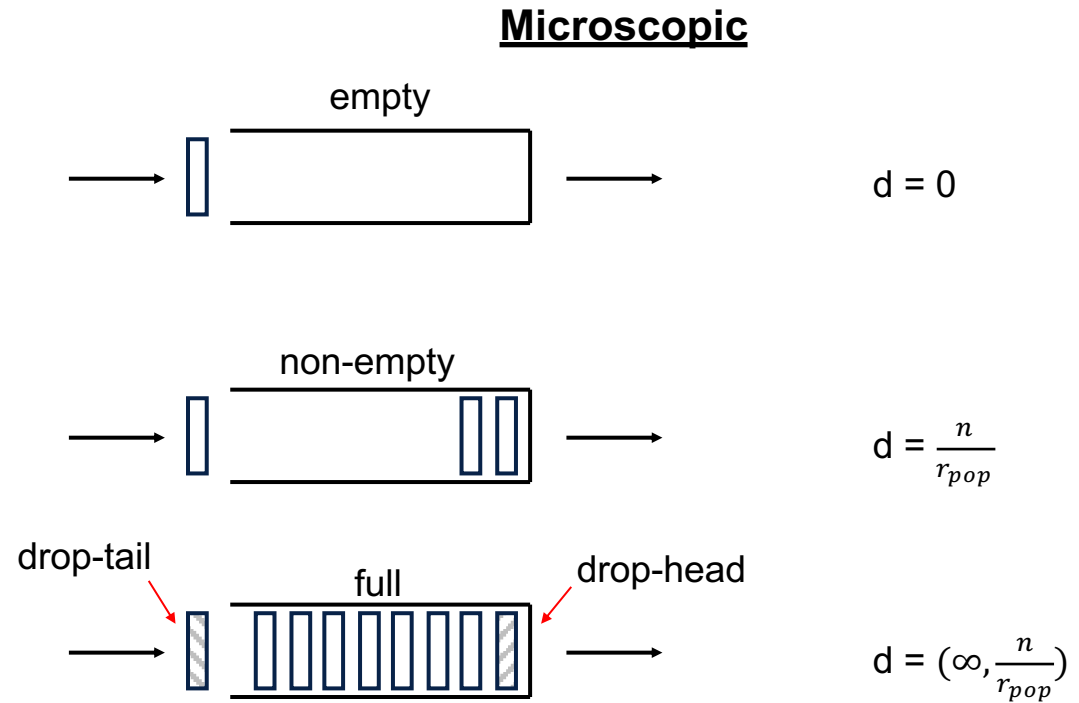
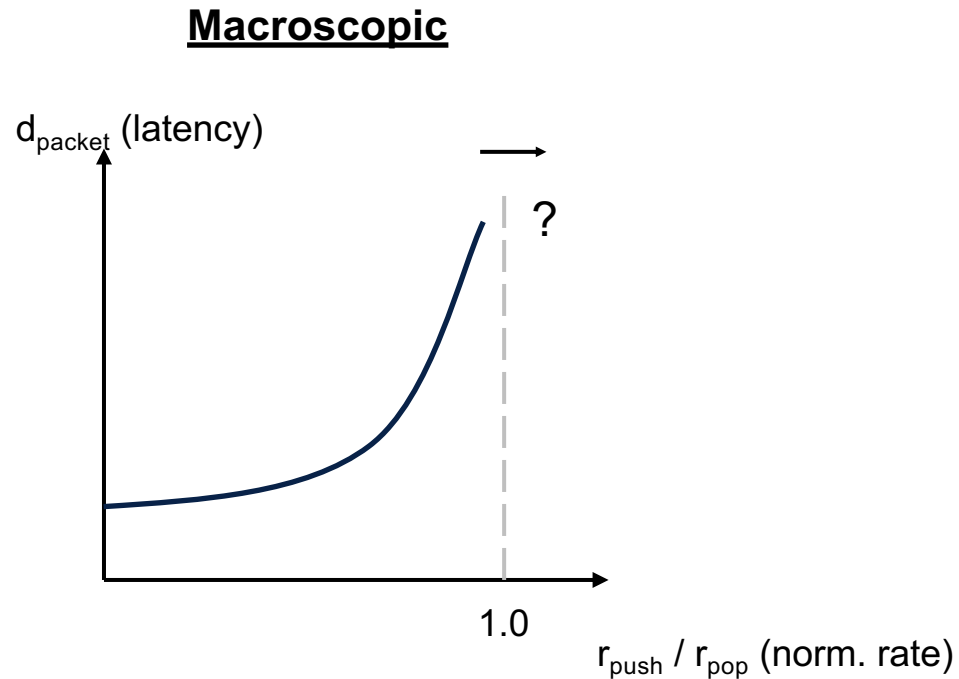
**time-variant**  
**modulation by FIFO (stateful FIFO delay)**

$$d_N : \text{Latency of packet } N := t'_N - t_N$$

Can **RTO** be bounded from above?  $\exists \lambda \quad d_N < \lambda \quad \forall N$

# Delay of a non memoryless system (cont.)

## Latency by FIFO (fill-level)

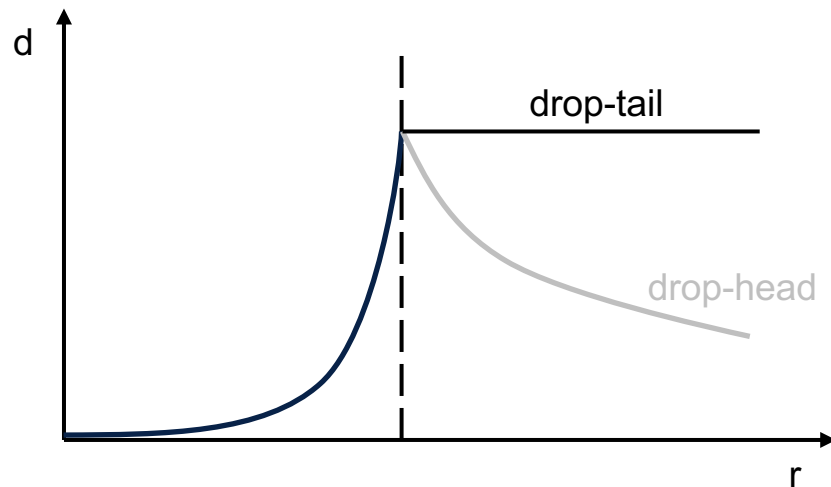


$n$  : instantaneous fill-level  
 $r$  : push/pop rate

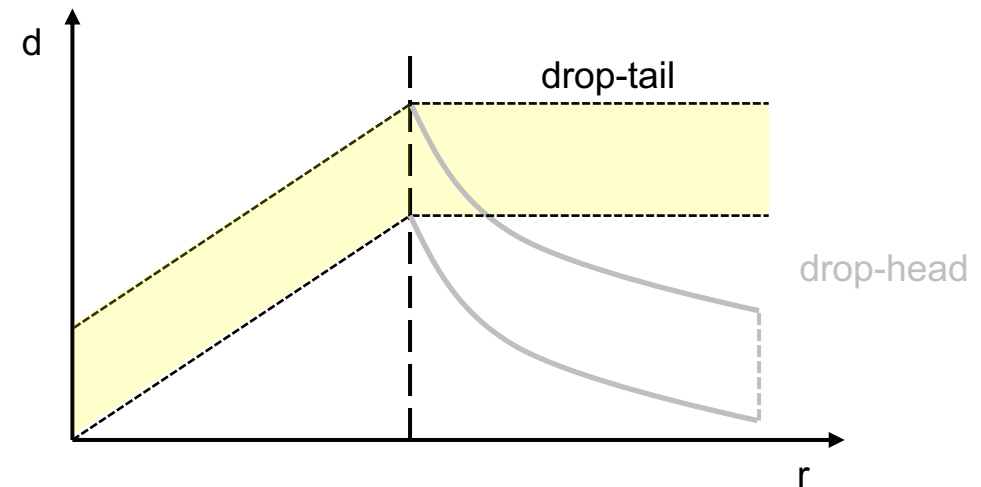
## Delay of a non memoryless system (cont.)

### stateful FIFO delay (operation mode)

Cut-through mode



Store-and-Forward mode



drop-head : possible recovery from deadlock  
 drop-tail : easy to implement

# How to bound RTO in the system?

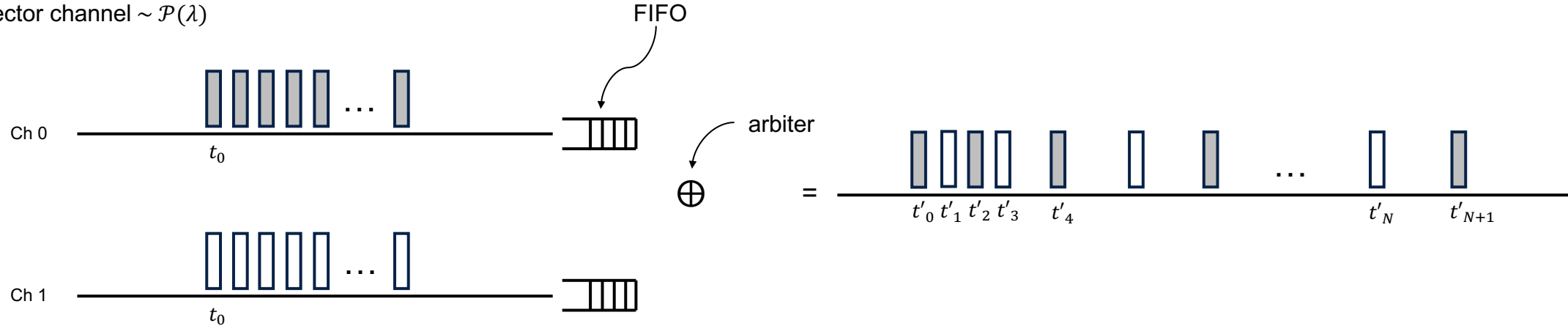
Sorter choice 2: ???

Q: How to bound the RTO?

➤ A: need to estimate the **delay jitter**

## stateful FIFO delay and lane-skew accumulation

detector channel  $\sim \mathcal{P}(\lambda)$



**time-variant** and **non-linear**  
modulation by FIFO and arbiter

$d_N$  : Latency of packet N :=  $t'_N - t_N$

**Out-of-Order:** There exist output positions  $i < j$  such that  $TS_i > TS_j$ .