

Reed Solomon FEC to correct for packet erasures

RS(8,2) GF(4) x N

- Designed for 100 Gbps software implementation
- Designed for minimal FPGA gate utilization (one dot product per restored symbol)
- Systematic RS code: Receiver only processes missing packet data
- 4 bit symbols, 8 data nibbles + 2 parity nibbles: Trivial gate count

Table lookup at decoder eliminates costly CHIEN search and Syndrome calculation

- Vector correction recognizing that for a packet, error locations are the same
- Bit level interleaving for full RS burst error protection
- Usable with any UDP streaming protocol (not EJFAT specific)

Pre-Computed H Matrix: Erasure Recovery Worked Example

1. Receive codeword — identify erasures from sequence numbers

2. Index erasure pattern to retrieve pre-computed H matrix

3. Substitute surviving parity symbols for erased positions

Only the 2 correction rows (d₃, d₄) require computation. The 3 identity rows are a pass-through — no multiply needed. For RS(10,8) with T=1, only 9 H matrices cover all possible single-packet erasure patterns.

m_corrected = [1 2 3 4 5] Voilà!

Parallel RS Operations & H Matrix Storage

Byte → Nibble Packing → 2 Parallel Codewords

Single Lost Packet → Same Erasure for All Codewords

Pre-Computed H Matrix Storage

45	2 x 8	360 B
H matrices	symbols per matrix	total storage
Choose {10,2} erasure patterns	2 correction rows x 8 cols	45 x 8 bytes (720 nibbles)

E2SAR Reassembler — Receive Path Architecture

Network / Socket: UDP from Load Balancer, Kernel SO_RCVBUF 3 MB, per-thread UDP sockets

Recv Thread xN: select() + recvfrom(), Parse RE Header, Reassemble into event buf

LOCK-FREE Event Queue: boost::lockfree::queue, capacity: 1000 events, ptr handoff — no copy

Application: getEvent(), recvEvent() — blocking, Owns buffer: delete[]

Recv Thread: per-packet receive + reassembly (N threads, each owns a subset of UDP ports)

Copies & Buffers: kernel → recvBuffer, recvBuffer → Event Buffer, Event Buffer → memcopy(), memcopy() → eventQueue.

GC Thread x1: Sweeps every 500 ms, Free timed-out events, Mutex on eventsInProgress

NO ZERO COPY: malloc/free per packet, 3 copies per event, No DPDK / io_uring / recvmmsg

SendState Thread x1: PID controller on queue, gRPC → control plane, Every 100 ms

Threads: Recv (xN) — select() + recvfrom() + reassemble, GC (x1) — timeout sweep, SendState (x1) — gRPC PID feedback, App — getEvent() / recvEvent()

E2SAR Segmenter — Send Path Architecture

Application: addToSendQueue(buf, len), sendEvent(buf, len), Owns event buffer

LOCK-FREE Lock-Free Queue: capacity: 2047 Items, Ptr only. No copy

Send Thread x1: Busy-wait poll, Dequeues events, Posts to pool

Thread Pool xN: N = # sockets

Worker Thread: _send() (N pool threads, one event each)

UDP Send: sendmsg(), sendmmsg(), batch io_uring async, 4 UDP sockets, Random src ports, 3-MB-send

On Wire: IP Hdr 20B, UDP Hdr 8B, LB Hdr 16B, RE Hdr 18B, Payload

Sync Thread x1: Sends SyncHdr every ~1s, Separate UDP socket → Control Plane

ZERO COPY: addToSendQueue() stores ptr, ptr [iov[1] → caller's buffer, No memcopy ever

Threads: App (x1), Send (x1), Workers (xN), Sync (x1)

EJFAT / E2SAR | Segmenter Data Path: Application → Queue → Workers → UDP

E2SAR Software library. Fully open source.

- Incorporated in dozens of projects
- Tested at upto 100G per compute node
- Python and C++ distributions
- Docker / Conda / Debian Packages

- Multi threaded with CPU pinning
- Optimized memory to memory handoff
- 100% Claude compatible

Opus Plan with high effort · API Usage Billing · yak.kumar@gmail.com's Organization · /private/tmp/example

Plan in mode: study how e2sar_py works by looking at the E2SAR github repository, and write a source and sink example that transmits a numpy array using ejfat. ask questions for clarification.

Interleaving Bits in Segments to FEC Code Words

Incoming Segment Perspective

Reed Solomon Symbol Perspective — Used for Parity Calculation

Complete FEC Block contains: 32 segments x 8192 bytes/segment = 256 KiB (max) = 2097152 bits (max)

Complete FEC Block contains: 65536 words x 8 symbols/word x 4 bits/symbol = 256 KiB (max) = 2097152 bits (max)

Burst Error Protection: Now we can correct up to 16 packets lost out of 48 total as long as they corrupt no more than 2 symbols.

- In production and deployed at ESnet and Jefferson Lab
 - 800 Gbps WAN capacity
 - Instant availability from any facility
- Early adoption and integration:
 - Advanced Light Source
 - Advanced Photonic Source
 - Jefferson Lab (CEBAF / EIC)
 - SLAC (LCLS)
 - FRIB
 - NSLS-II
- HPC Facilities:
 - NERSC (Perlmutter)
 - Oakridge OLCF (Defiant / Lux)
 - Jefferson lab (HPC)
 - UCSD (NRP / OSG)
 - NSF (FABRIC)

EJFAT ESnet Jefferson Lab FPGA Accelerated Transport

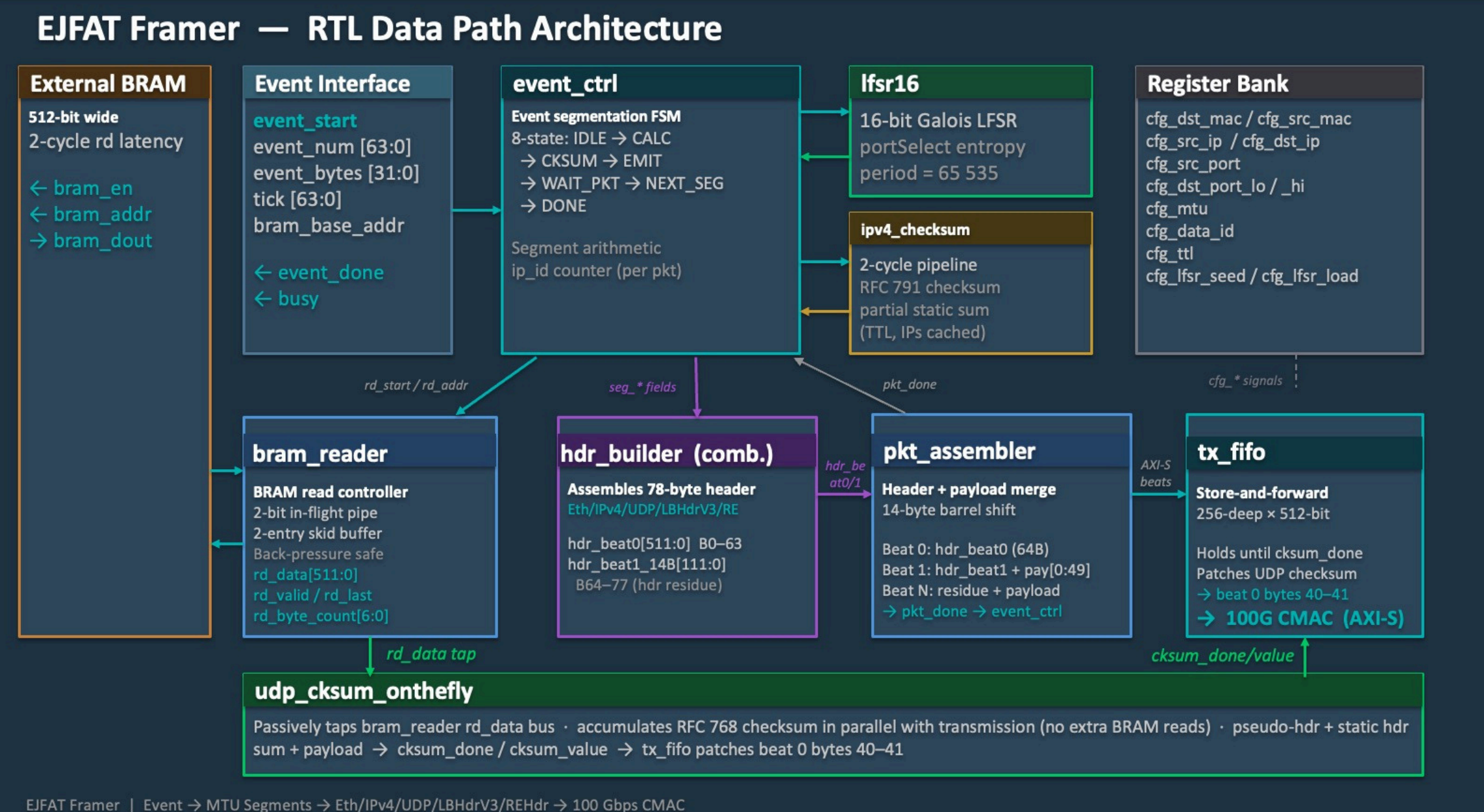
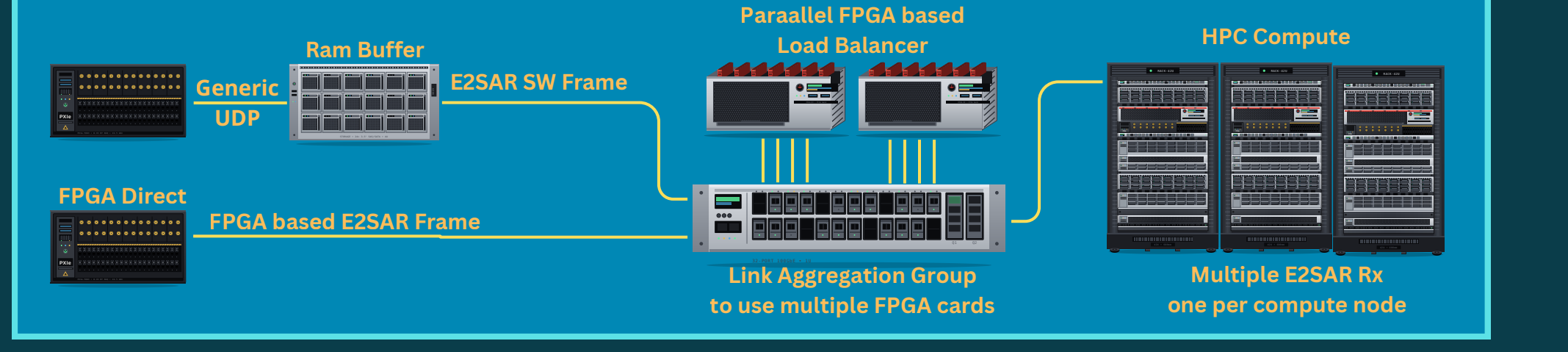
Real Time streaming to multiple supercomputer sites simultaneously!
 Dynamic load balancing, with uSecond granularity
 Dynamic node connection / removal
 Terabit scale

ESnet ENERGY SCIENCES NETWORK

Jefferson Lab

<https://github.com/JeffersonLab/E2SAR>

contact: Ilya Baldin (baldin@jlab.org)
 Yatish Kumar (yak@es.net)



LUTs REG BRAM

Framer	7000	4000	1
FEC encoder	11000	6400	64
FEC decoder	23480	14000	144

- VHDL framer developed and tested at JLab
- Verilog framer developed and tested at ESnet
- FEC Verilog and SW modules in development
- 100 Gbps design