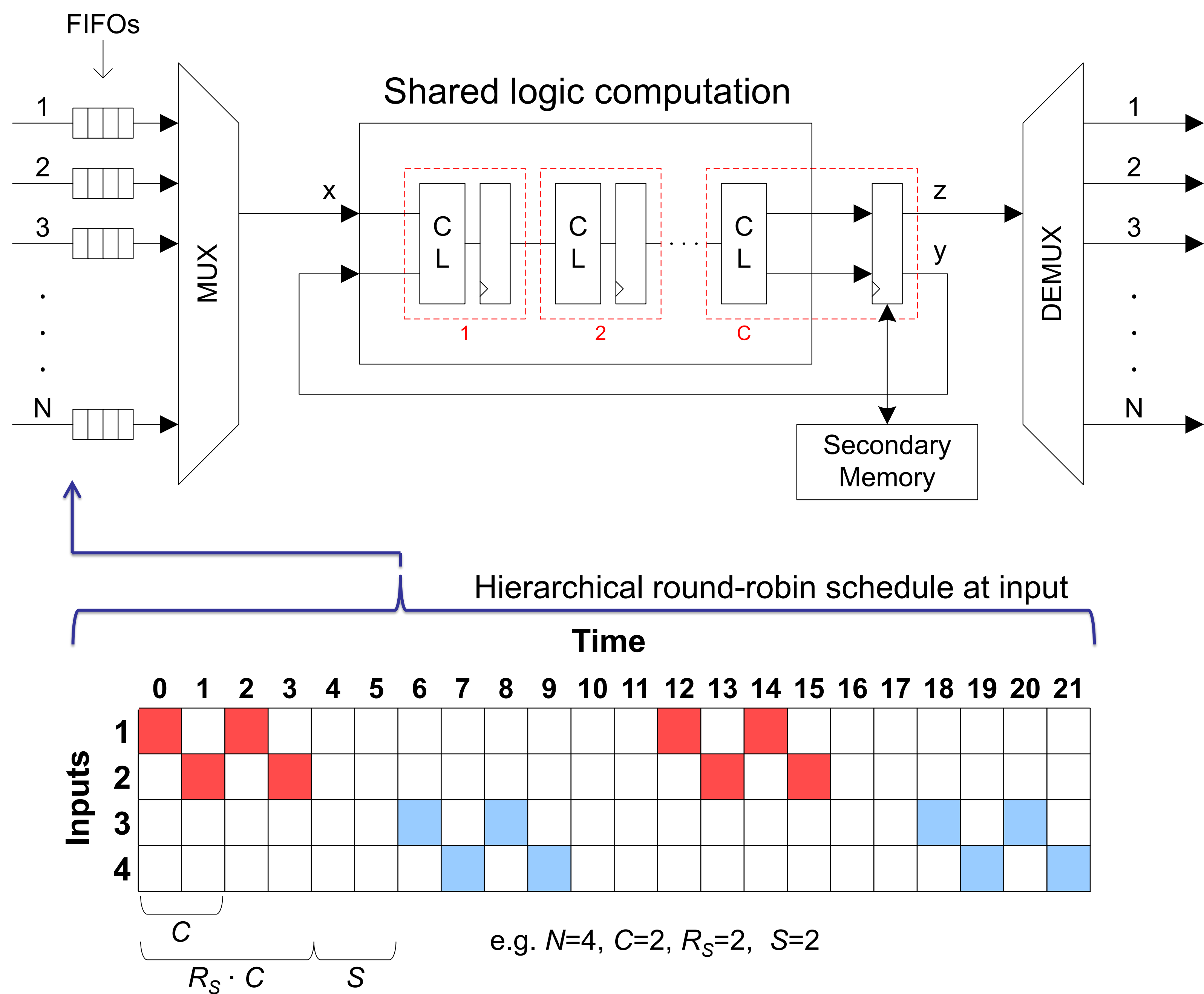
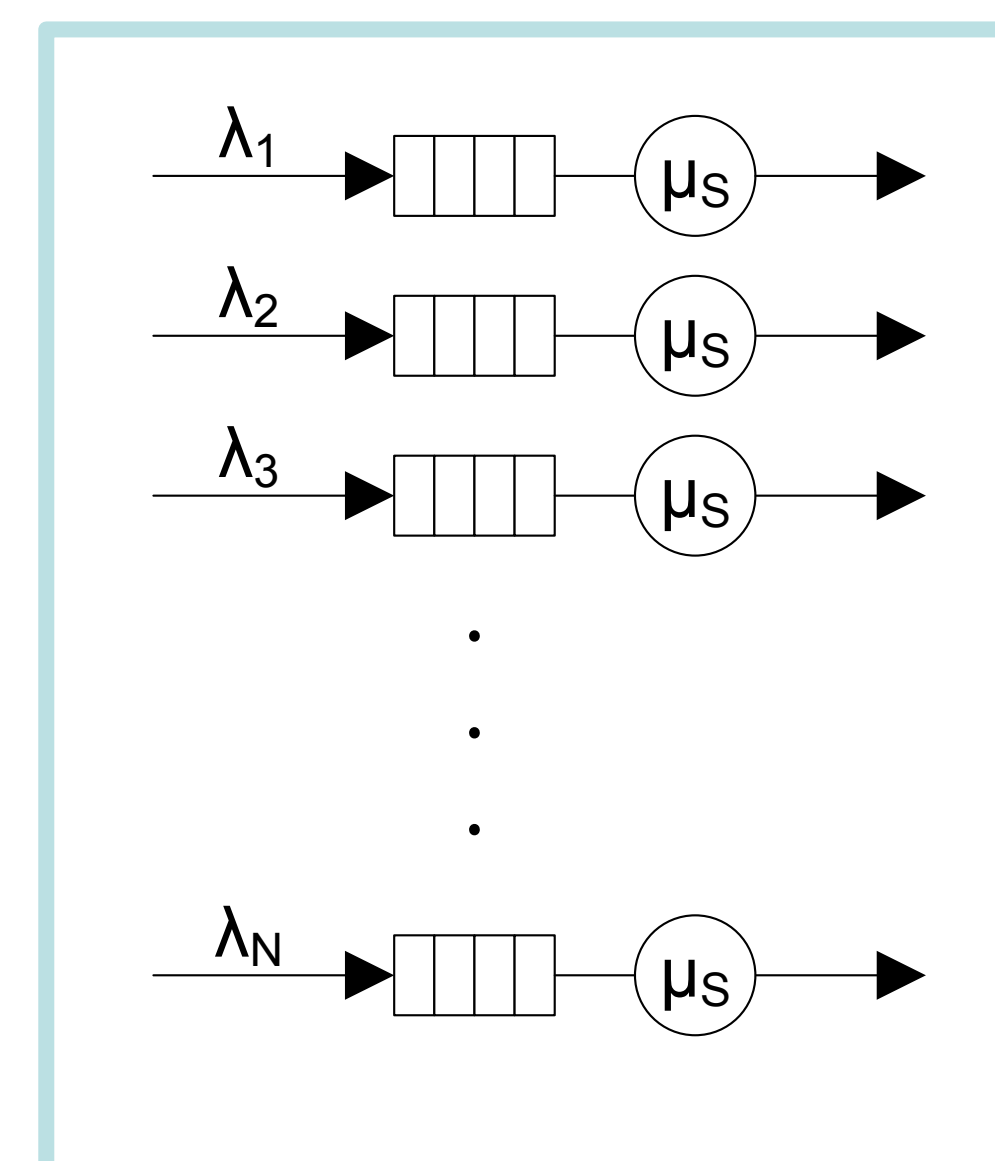


Virtualized Logic Computation



Queueing Model



Model definition:

$$T_{put}, Latency, Occupancy = f(Circuit, Tech, C, N, S, R_S, \lambda)$$

Variable	Definition
<i>Circuit</i>	Logical circuit description (e.g. written in Verilog or VHDL)
<i>Tech</i>	Target technology (e.g. FPGA or ASIC)
<i>C</i>	Pipeline depth (fine-grain contexts)
<i>N</i>	Total number of contexts (requires secondary memory if $N > C$)
<i>S</i>	Cost of a context switch (for secondary memory)
R_S	Scheduling period (number of rounds of <i>C</i> contexts that execute before context-switching to secondary memory)
λ	Arrival rate (e.g. data elements per second)

M/D/1 modeling expressions

Service rate of each virtual logic computation:

$$\mu_s = \frac{R_S}{(R_S N + S N / C) \cdot t_{CLK}} \text{ elements/s}$$

New equations in red boxes

Total achievable throughput:

$$T_{TOT} = N \cdot \mu_s = \frac{R_S}{(R_S + S / C) \cdot t_{CLK}} \text{ elements/s}$$

Waiting time expressions: (Latency)

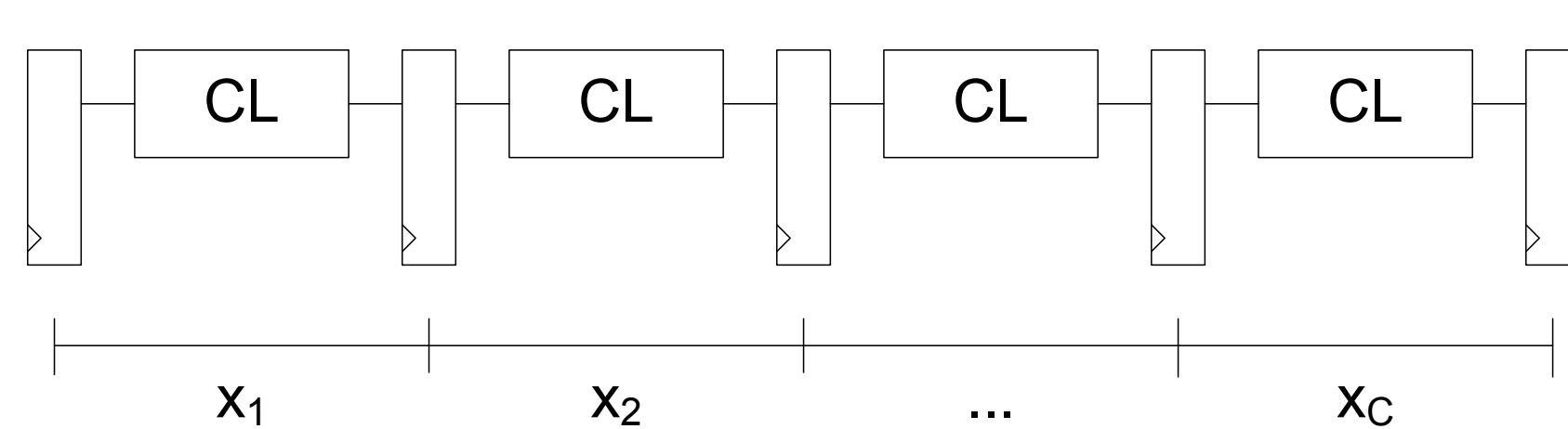
$$W_q = \frac{1}{\mu_s} \cdot \frac{\rho}{2(1-\rho)} \quad W_s = \frac{1}{\mu_s} = C \cdot t_{CLK}$$

$$W_h = \frac{(R_S(N-C) + S N / C) \cdot t_{CLK}}{2(R_S N + S N / C)} \quad W = W_q + W_h + W_s$$

Queue occupancy:

$$N_q = \lambda \cdot W_q$$

Calibration



Variable	Definition
<i>X</i>	Random variable of stage-to-stage delay for <i>C</i> random samples, x_1, x_2, \dots, x_C
<i>C</i>	Pipeline depth
t_{CL}	Total comb. logic delay
<i>k</i>	Curve fit parameter

Upper bound on t_{CLK} from order statistics:

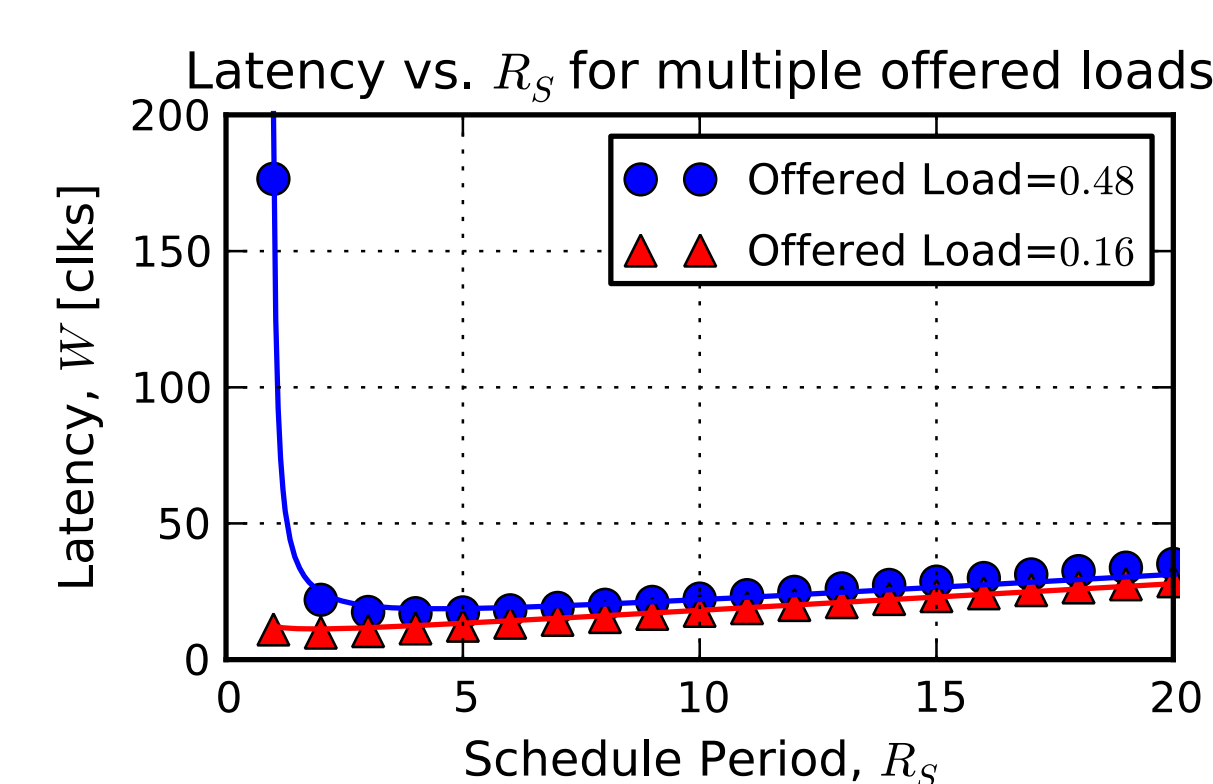
$$t_{CLK} = E[\max(X)] \leq \mu_X + \sigma_X \sqrt{C-1}$$

Clock period sub-model:

$$t_{CLK} = t_{CL} / C + k \cdot \sqrt{C-1}$$

$$T_{TOT} = \frac{1}{t_{CLK}}$$

Model Validation



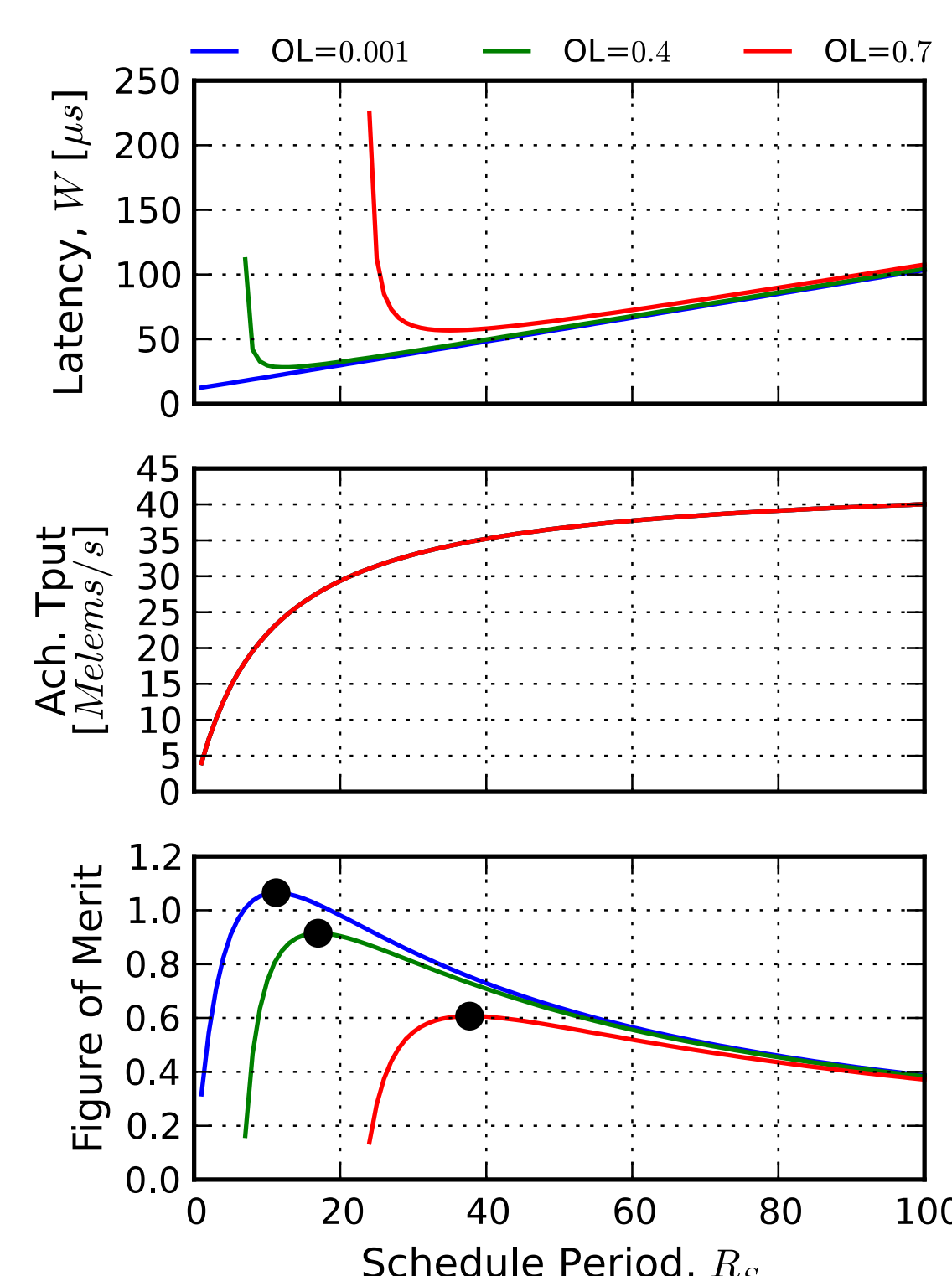
$N = 8, C = 4, S = 4, t_{CLK} = 100 \text{ MHz}$

- Offered load is the ratio of aggregate arrival rate to the peak service rate of the system and is equal to $N \cdot \lambda \cdot t_{CLK}$.
- Points are discrete event simulation
- Curves are analytic expressions

Ways to Use the Model and Results

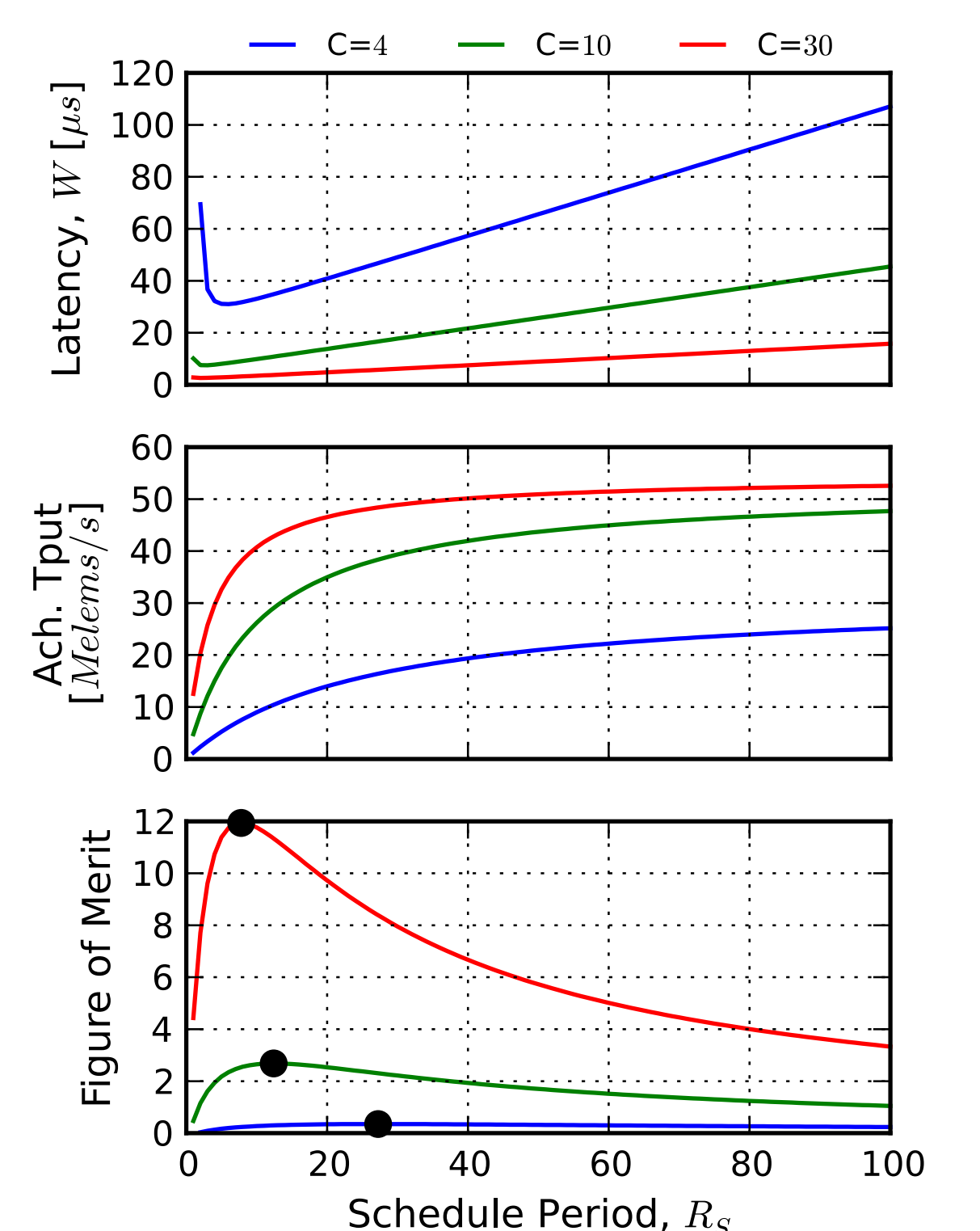
Example Design - Case 1

Given:
• *Circuit*=COS, *Tech*=FPGA,
 $N=100, C=10, S=100, OL$ varies
Design Params: R_S
Optimize: $FoM = T_{put} / Latency$

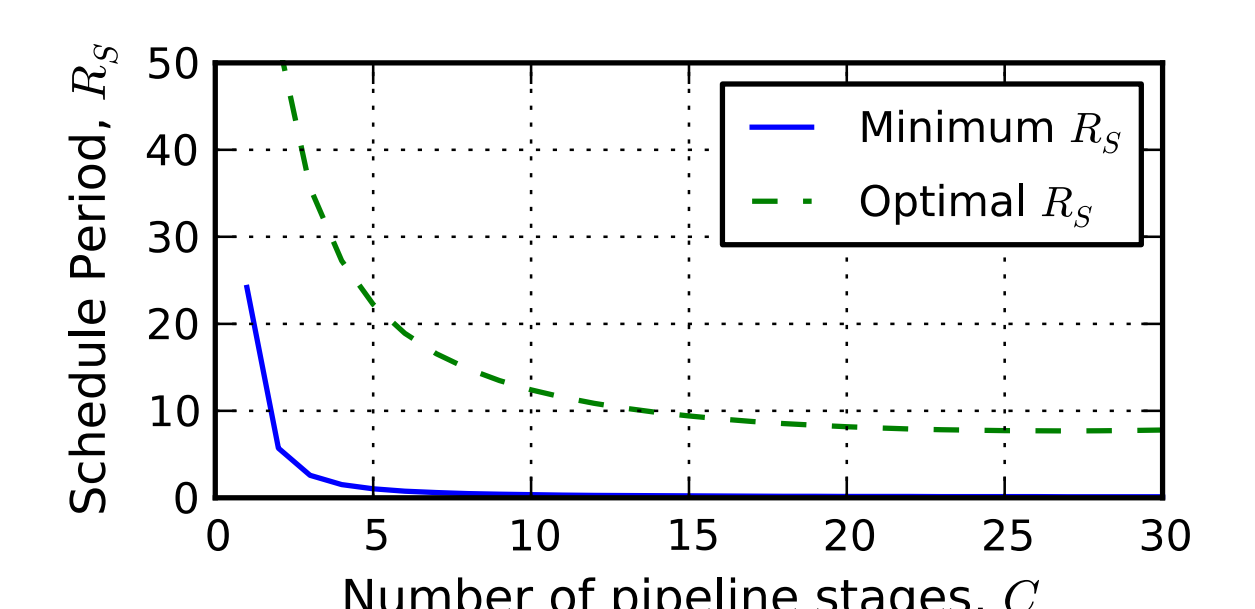
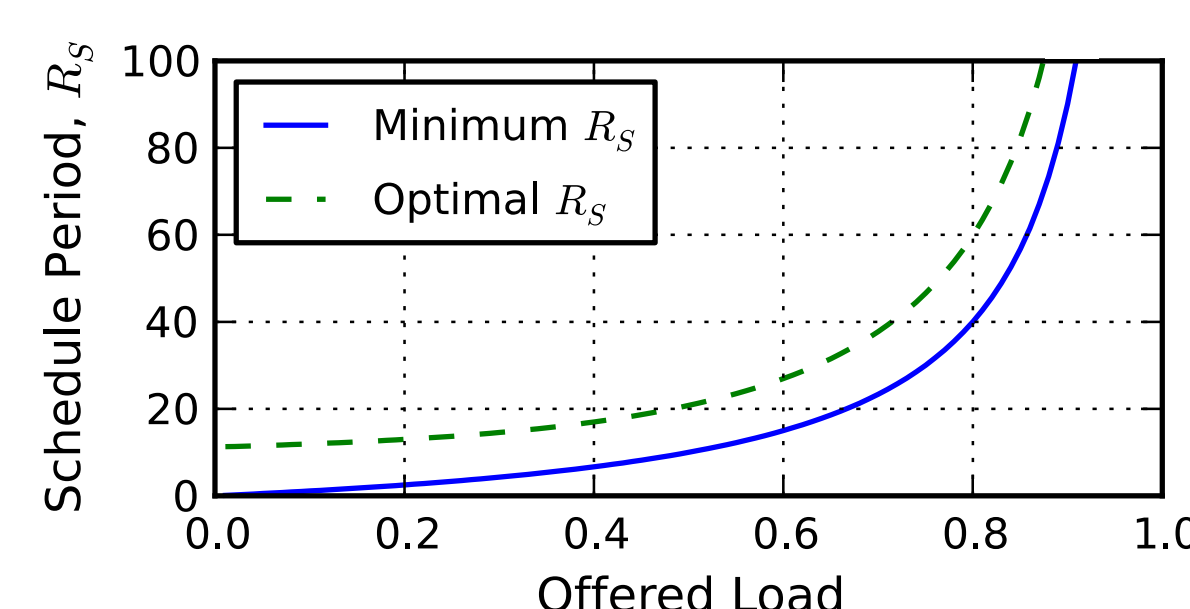


Example Design - Case 2

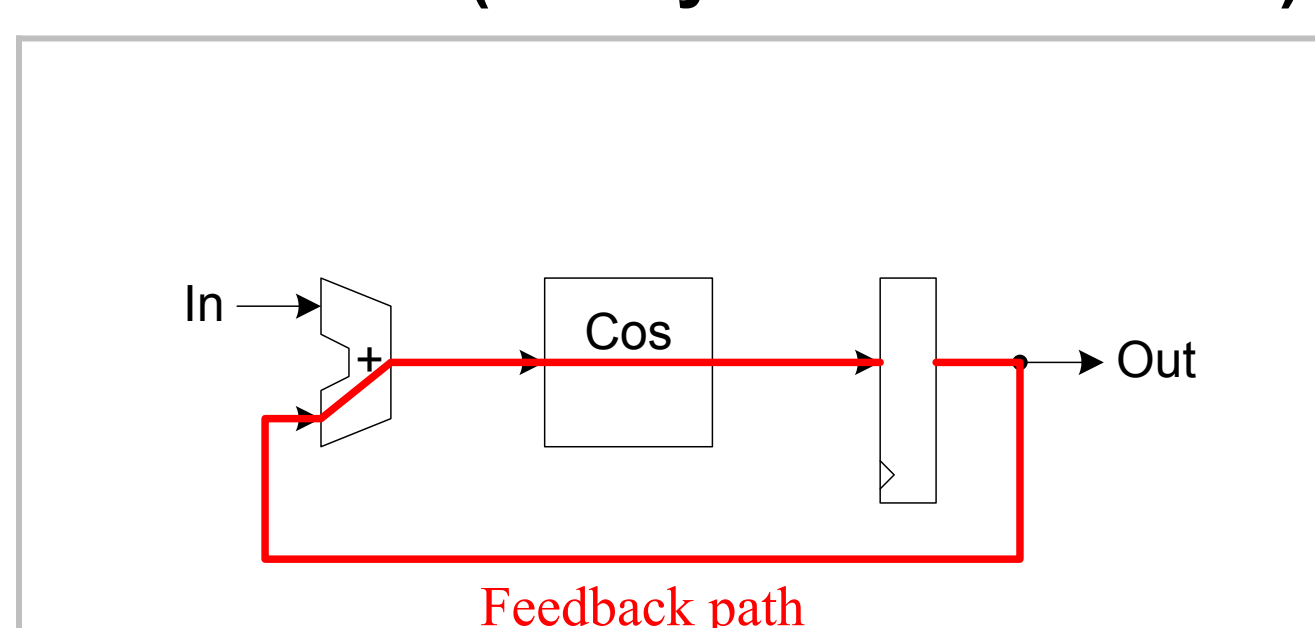
Given:
• *Circuit*=AES, *Tech*=ASIC,
 $N=60, S=100, \lambda=30$ Kelems/s
Design Params: C, R_S
Optimize: $FoM = T_{put} / Latency$



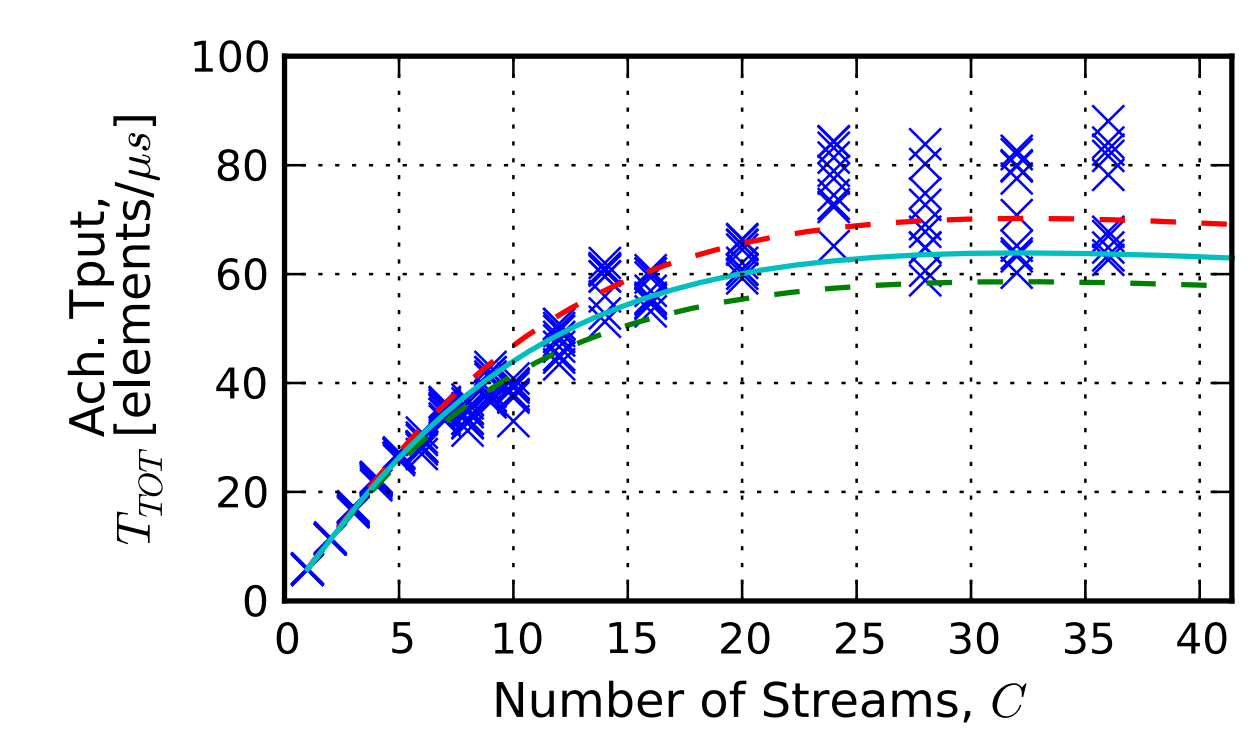
Optimization of Figure of Merit



Synthetic Cosine Application with Feedback (20 Taylor-series terms)

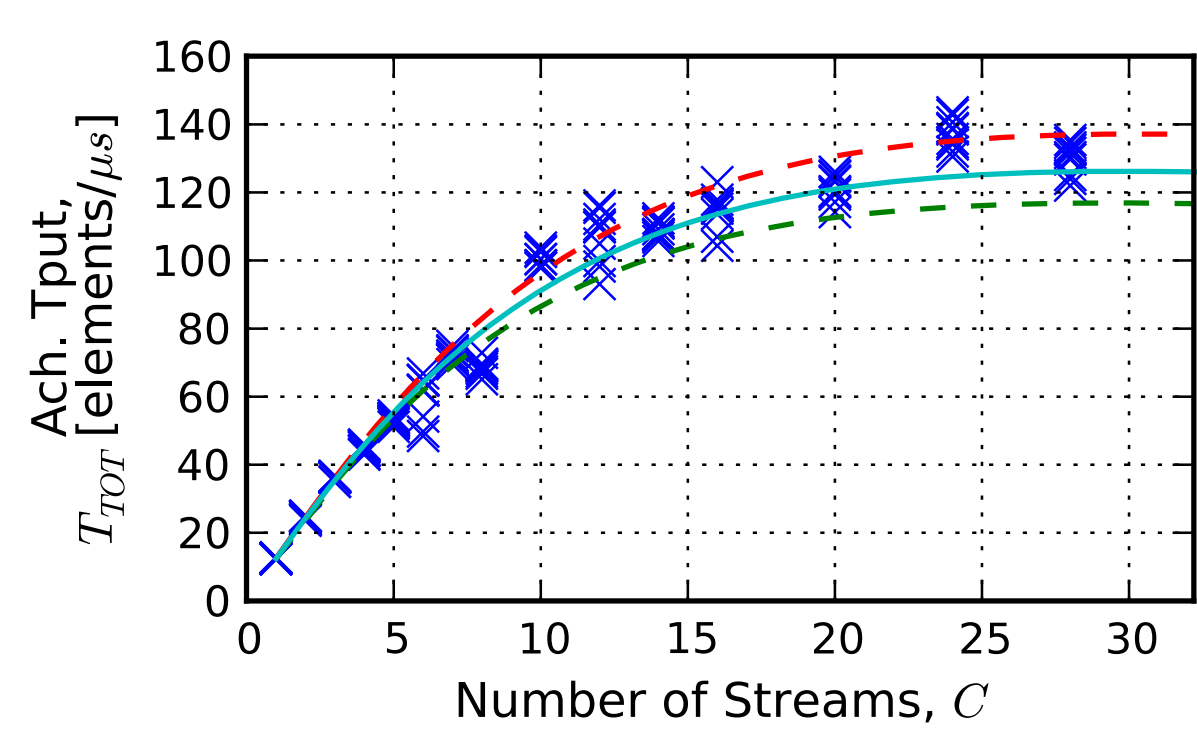
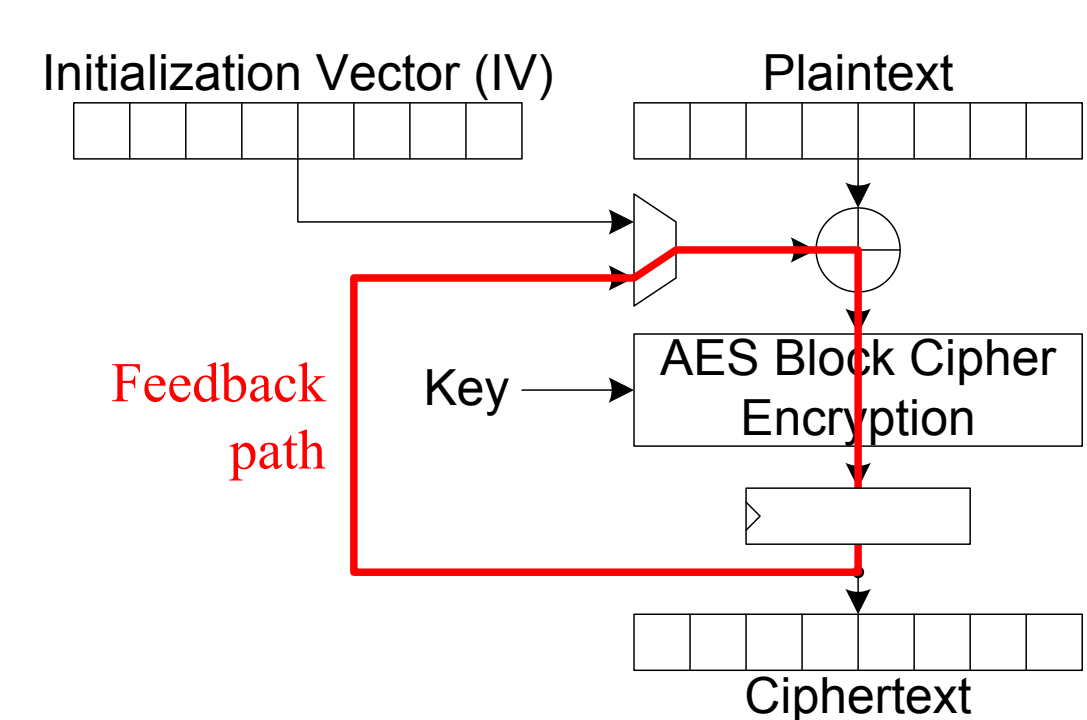


FPGA



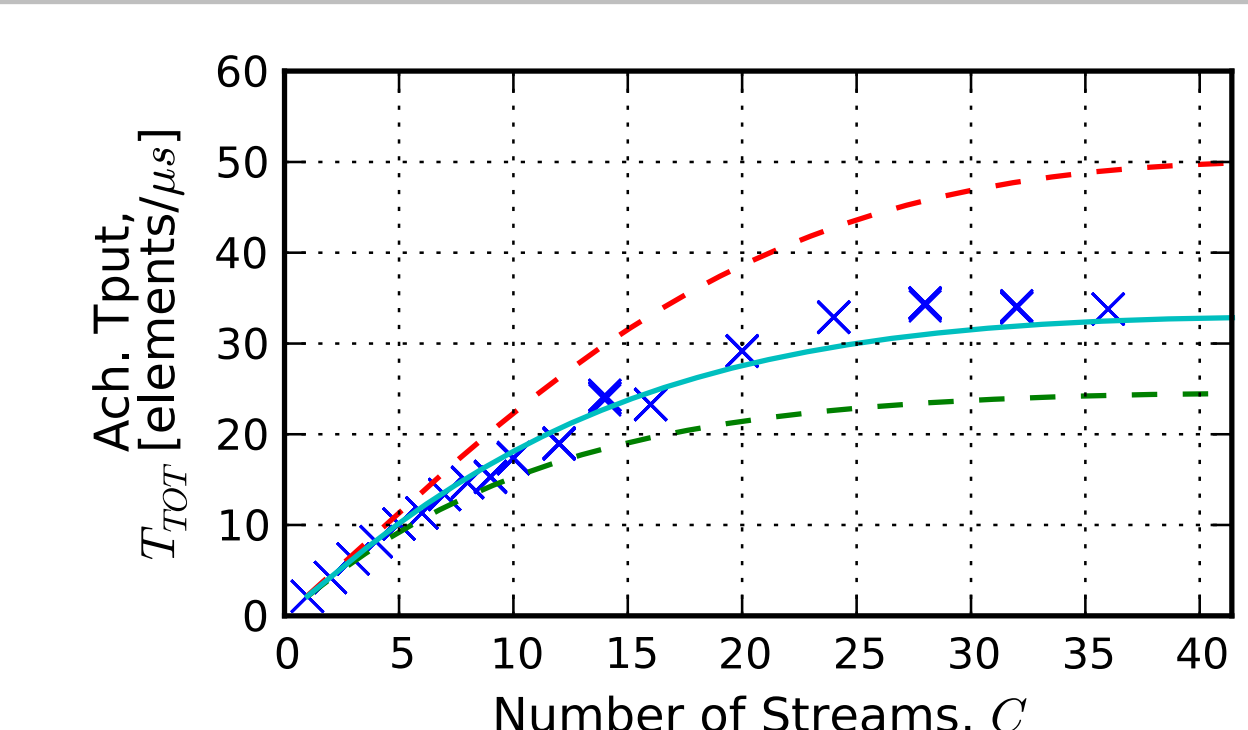
$$t_{CLK} = 8.6 \frac{ns}{term} \cdot \frac{N_t}{C} + (1.9 ns - 4.7 \frac{ps}{term} \cdot N_t) \cdot \sqrt{C-1}$$

AES Encryption Cipher in CBC Block Mode (14 fully unrolled rounds)

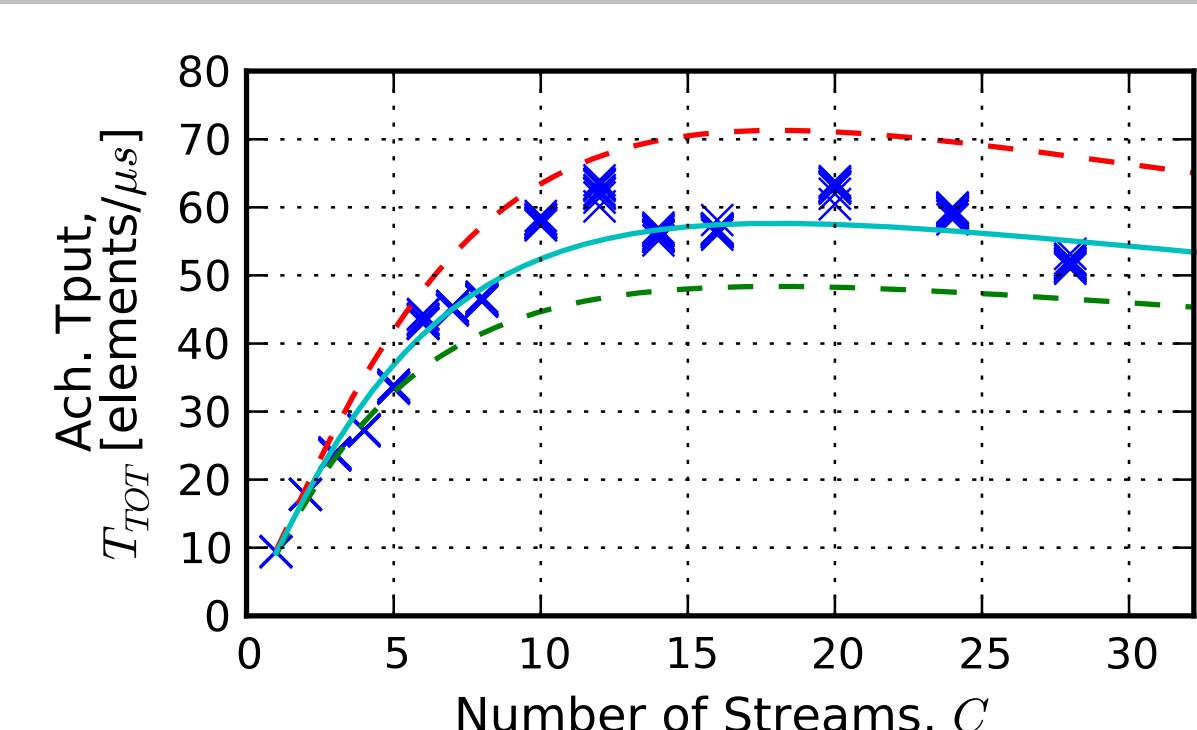


$$t_{CLK} = 5.7 \frac{ns}{rnd} \cdot \frac{N_r}{C} + (1.7 ns - 55 \frac{ps}{rnd} \cdot N_r) \cdot \sqrt{C-1}$$

ASIC



$$t_{CLK} = 2.3 \frac{ns}{term} \cdot \frac{N_t}{C} + (3.2 ns - 9.8 \frac{ps}{term} \cdot N_t) \cdot \sqrt{C-1}$$



$$t_{CLK} = 7.8 \frac{ns}{rnd} \cdot \frac{N_r}{C} + (4.4 ns - 122 \frac{ps}{rnd} \cdot N_r) \cdot \sqrt{C-1}$$