# Using M/G/1 Queueing Models with Vacations to Analyze Virtualized Logic Computations

Michael J. Hall
VelociData, Inc.
St. Louis, MO, USA

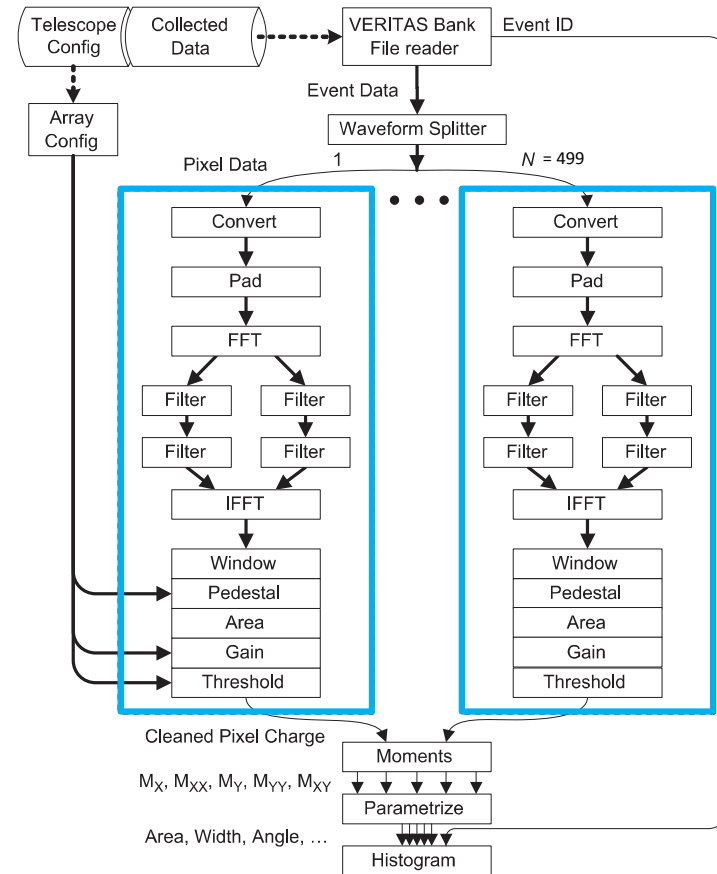Roger D. Chamberlain
Washington University in
St. Louis, MO, USA

October 19, 2015

33rd IEEE
International
Conference on
Computer Design
*ICCD 2015*

# Motivation for virtualized logic computations:
# An example big computation

- Telescope application with N=499 pixels
- Each pixel requires a channel of computation
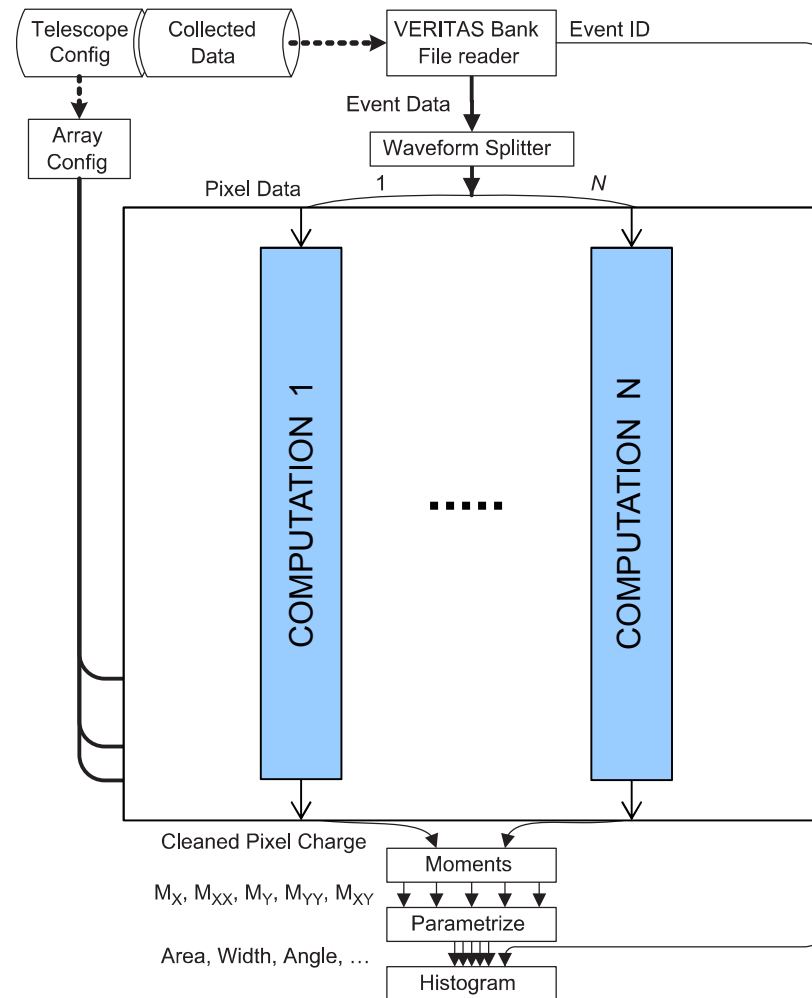- Replicating logic requires much hardware
- More hardware = more cost





Telescope Config — Collected Data ⟶ VERITAS Bank File reader — Event ID

Event Data

Array Config

Waveform Splitter

Pixel Data — 1 ... $N = 499$

Convert → Pad → FFT → Filter / Filter → Filter / Filter → IFFT → Window / Pedestal / Area / Gain / Threshold

Convert → Pad → FFT → Filter / Filter → Filter / Filter → IFFT → Window / Pedestal / Area / Gain / Threshold

Cleaned Pixel Charge

$M_X, M_{XX}, M_Y, M_{YY}, M_{XY}$ — Moments

Area, Width, Angle, ... — Parametrize

Histogram

VERITAS gamma-ray signal processing pipeline in X.

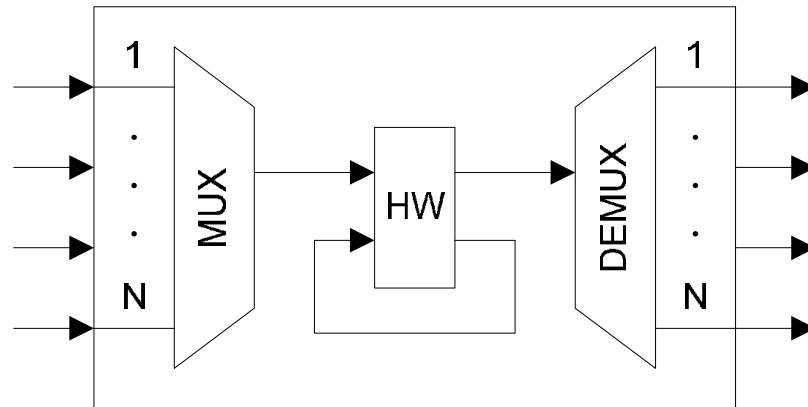[Tyson et al. 2008]

# Example big computation



VERITAS gamma-ray signal processing pipeline in X.

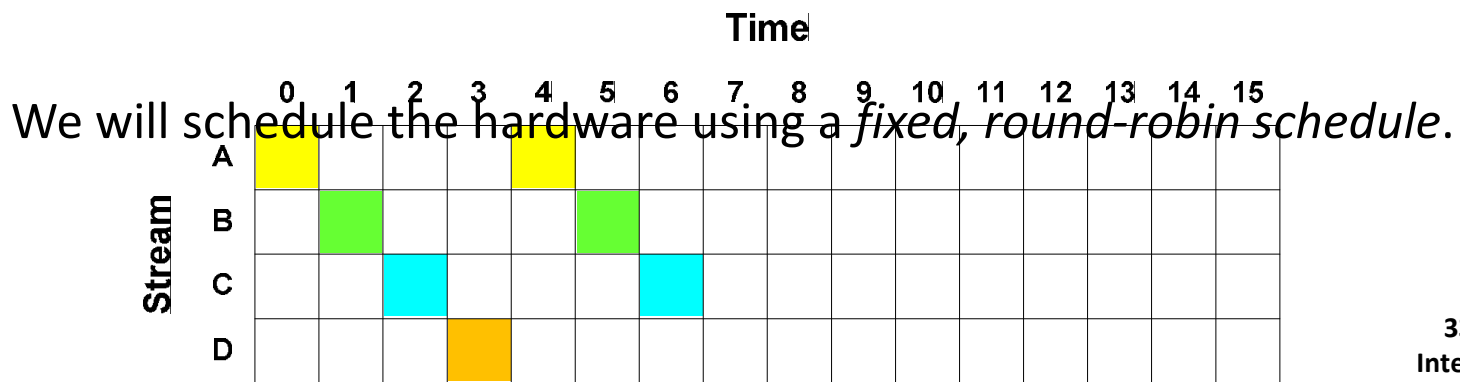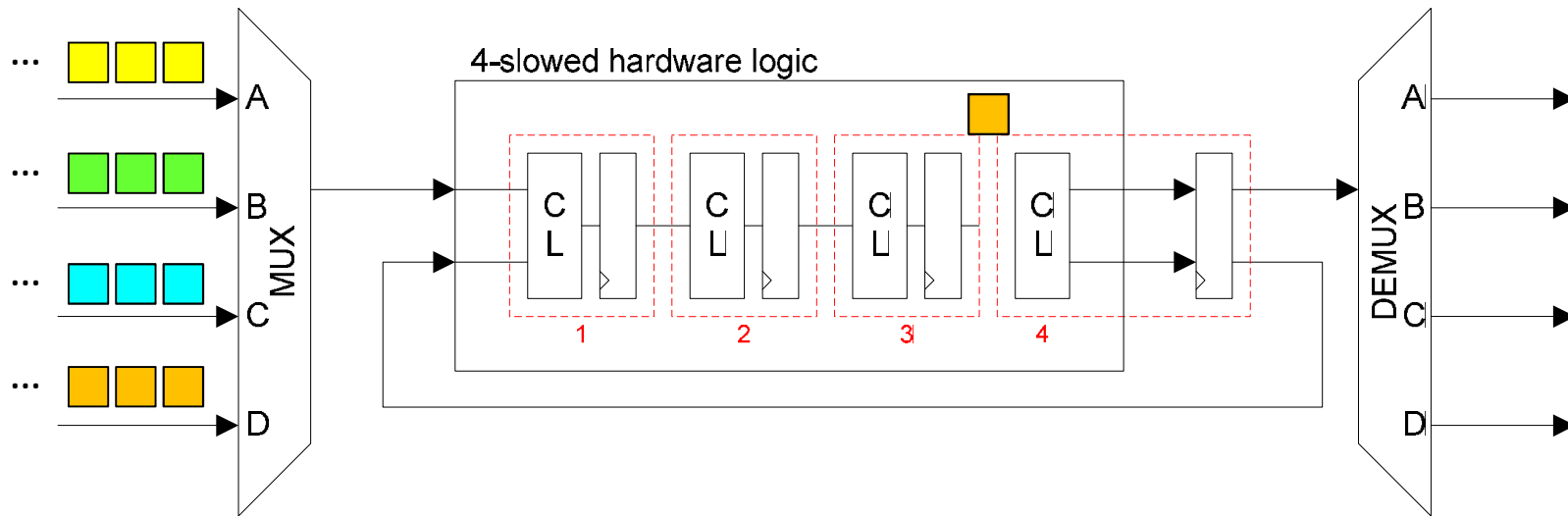# Example big computation



VERITAS gamma-ray signal processing pipeline in X.
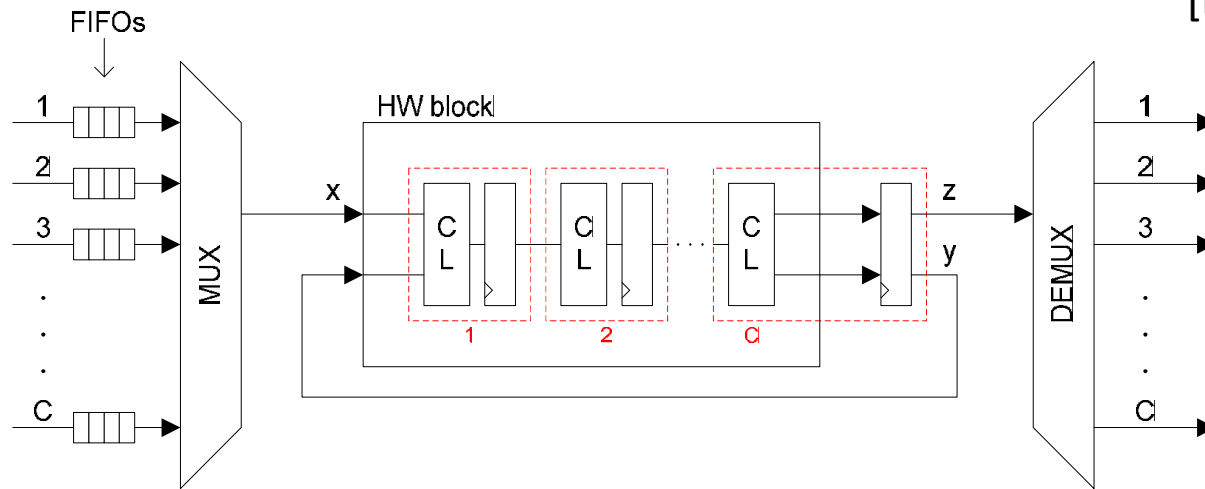
# Hardware virtualization



- Hardware virtualization for N distinct data streams that perform the same computation

- Interested in the case where there is feedback

- We will virtualize the hardware by applying a C-slow technique [Leiserson and Saxe 1991]

- We will derive queueing model equations to predict circuit performance
  - This is our main contribution
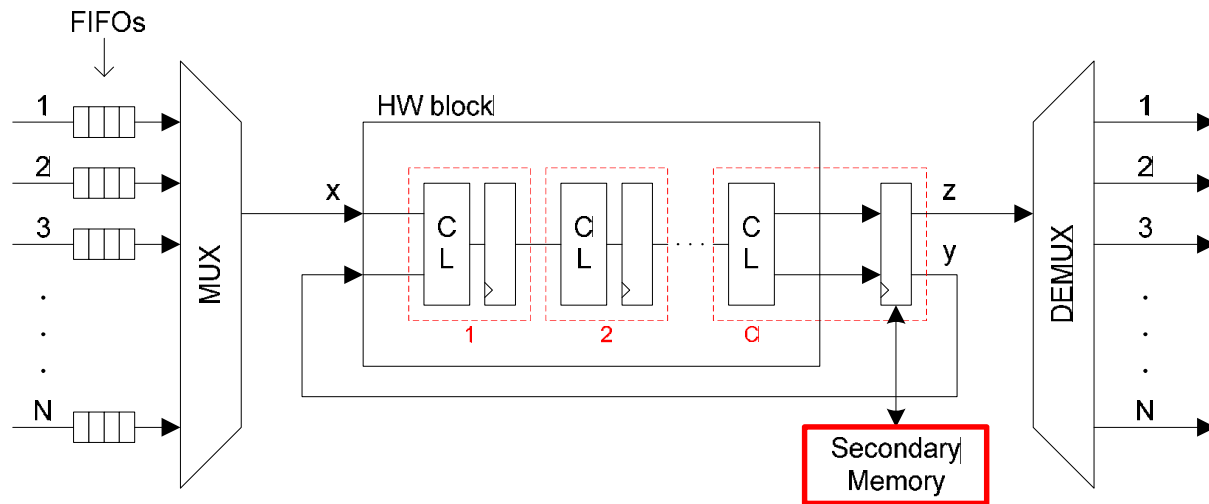
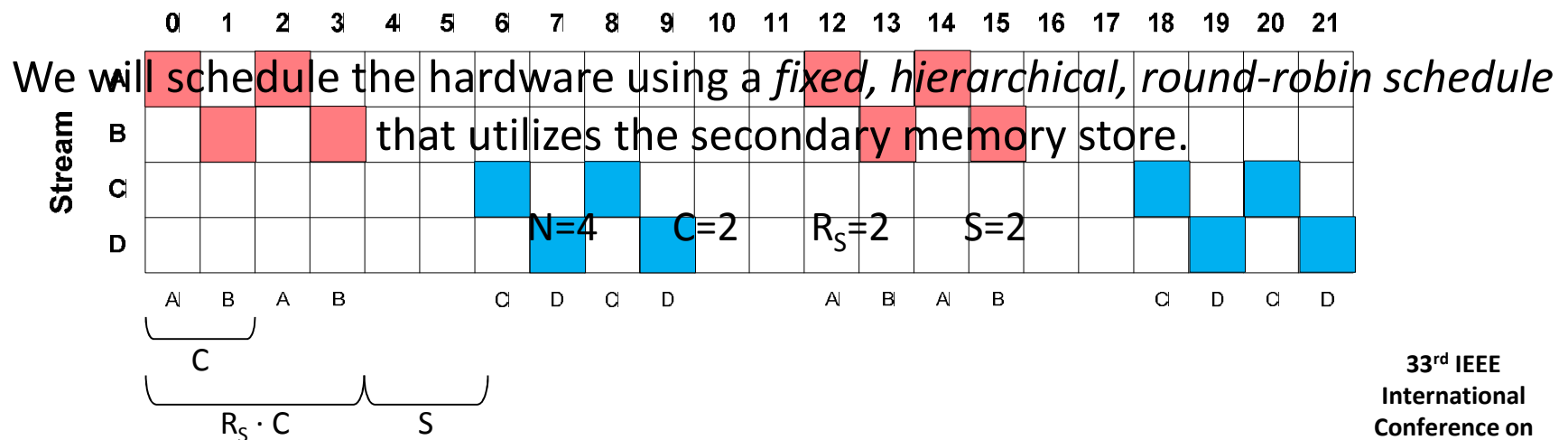# 4-slow virtualization example



We will schedule the hardware using a *fixed, round-robin schedule*.

6

# C-slow general virtualized hardware

[Leiserson and Saxe 1991]

# C-slow general virtualized hardware



We will schedule the hardware using a *fixed, hierarchical, round-robin schedule* that utilizes the secondary memory store.

$N=4$   $C=2$   $R_S=2$   $S=2$

| Var. | Definition |
|------|------------|
| N | Total contexts |
| C | Pipeline depth |
| $R_S$ | Scheduling period |
| S | Cost of context switch |

# Queueing model



Queueing model

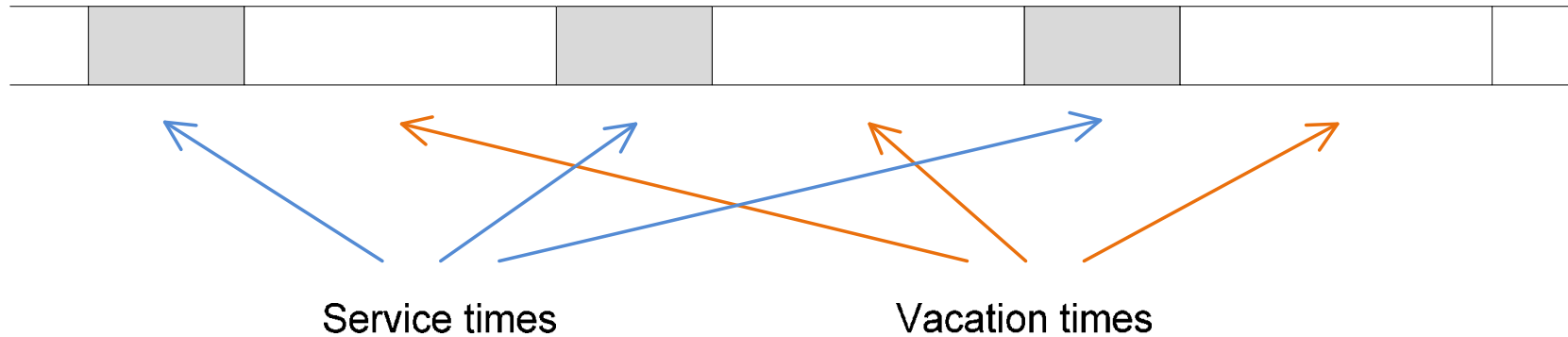- Each queueing station is modeled as an M/G/1 queueing model with vacations

- M/G/1 is *M*arkovian, or memoryless, arrival process; *G*eneral service process; and *1* server

# Service and vacation time modeling



Service times                    Vacation times

- Service and vacation times are regular and deterministic

- This is for 1 queueing station, or virtual instance of the hardware logic

# Model definition

Tput, Latency, Occupancy $= f$ (Circuit, Tech, $C$, $N$, $S$, $R_S$, $\lambda$)

| Variable | Definition |
|---|---|
| Circuit | Logical circuit description (e.g. AES-256) |
| Tech | Target technology (e.g. FPGA or ASIC) |
| C | Pipeline depth (also represents the number of fine-grain contexts) |
| N | Total number of contexts (requires secondary memory if N > C) |
| S | Cost of a context switch (to/from secondary memory) |
| $R_S$ | Scheduling period (number of rounds of C contexts that execute before doing a context switch to secondary memory) |
| $\lambda$ | Arrival rate (e.g. data elements per second) |

VelociData, Inc.

# Performance model

Total achievable throughput:

$$T_{TOT} = \frac{R_S}{(R_S + S/C) \cdot t_{CLK}}$$

Total wait time (latency):

$$W_T = \frac{\lambda \overline{X^2}}{2(1 - \rho)} + \frac{\overline{V}}{1 - \rho} + \overline{X}$$

Number in queue:

$$N_q = \frac{\lambda^2 \overline{X^2}}{2(1 - \rho)} + \frac{\lambda \overline{V}}{1 - \rho}$$

| Variable | Definition |
|----------|------------|
| C | Pipeline depth |
| N | Number of streams |
| S | Context switch cost |
| $R_S$ | Scheduling period |
| λ | Mean arrival rate |
| ρ | Utilization |
| $\overline{X}$ | Mean service time |
| $\overline{X^2}$ | Service time second moment |
| $\overline{V}$ | Mean vacation waiting |

# Ways to use the model in design

Model definition:
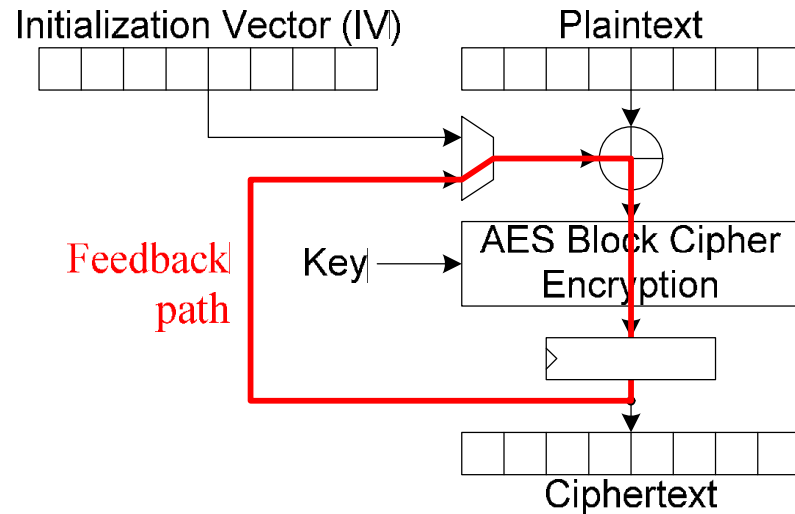
$$\text{Tput, Latency, Occupancy} = f\,(\text{Circuit, Tech, } C, N, S, R_S, \lambda)$$

- Subset of parameters are given
  - E.g., Circuit, Tech, N, S
- Remainder under control of designer
  - E.g., C, $R_S$, $\lambda$
- Design goal
  - E.g., Latency

| Variable | Definition |
|----------|------------|
| C | Pipeline depth |
| N | Number of contexts |
| S | Cost of a context switch |
| $R_S$ | Scheduling period |
| $\lambda$ | Arrival rate |

# Experimental setup

- AES-256 encryption application in CBC mode

- Fully unrolled, $N_r$ = 14 rounds

- Targeting Xilinx Virtex-4 XC4VLX100 FPGA



Calibrated $t_{CLK}$ model:

$$t_{CLK}(N_r, C) = \left[ \frac{1.8 + 5.2 \cdot N_r}{C} + (2.56 - 0.038 \cdot N_r) \cdot (\ln C)^{0.7} \right] \text{ns}$$
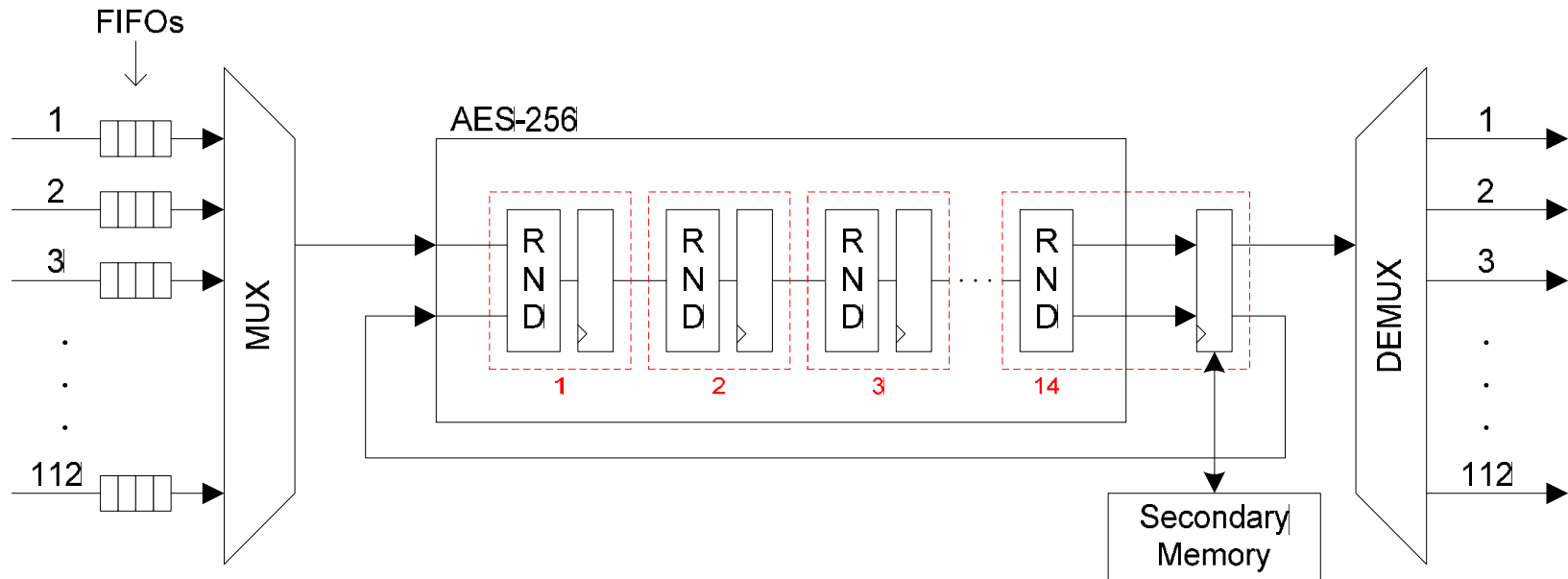
Fixed model parameters:

Pipeline depth, C = 14
Total contexts, N = 8C (or 112)
Context switch cost, S = 120 clock cycles

# 14-slow AES-256



- Each round (RND) performs a series of operations on a block of data propagating through to the output
  - Substitutions, shifts, multiplications, and logical operations
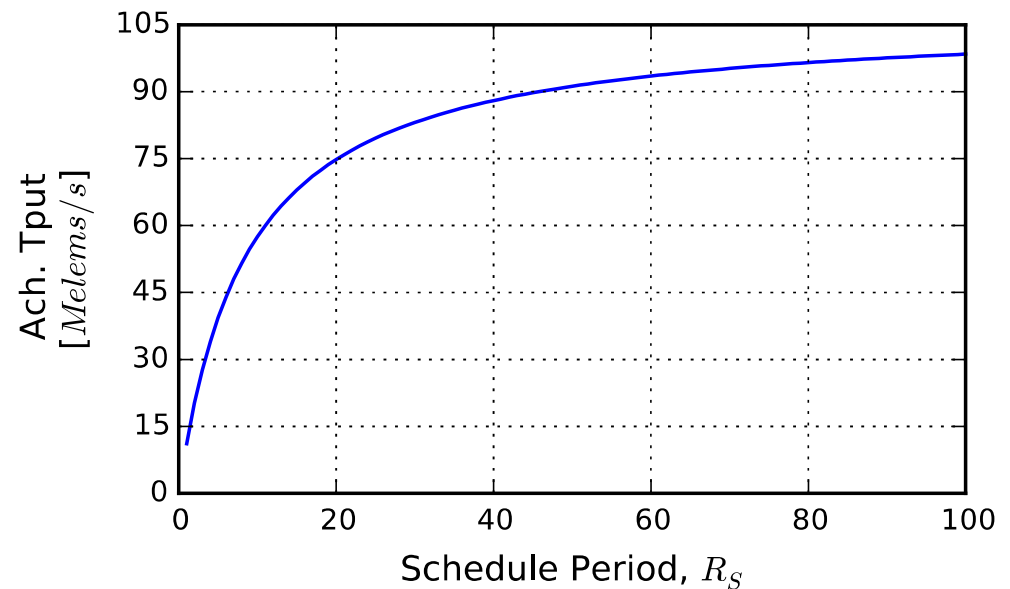- We have 112 virtual copies

# Total achievable throughput

- Model parameters:
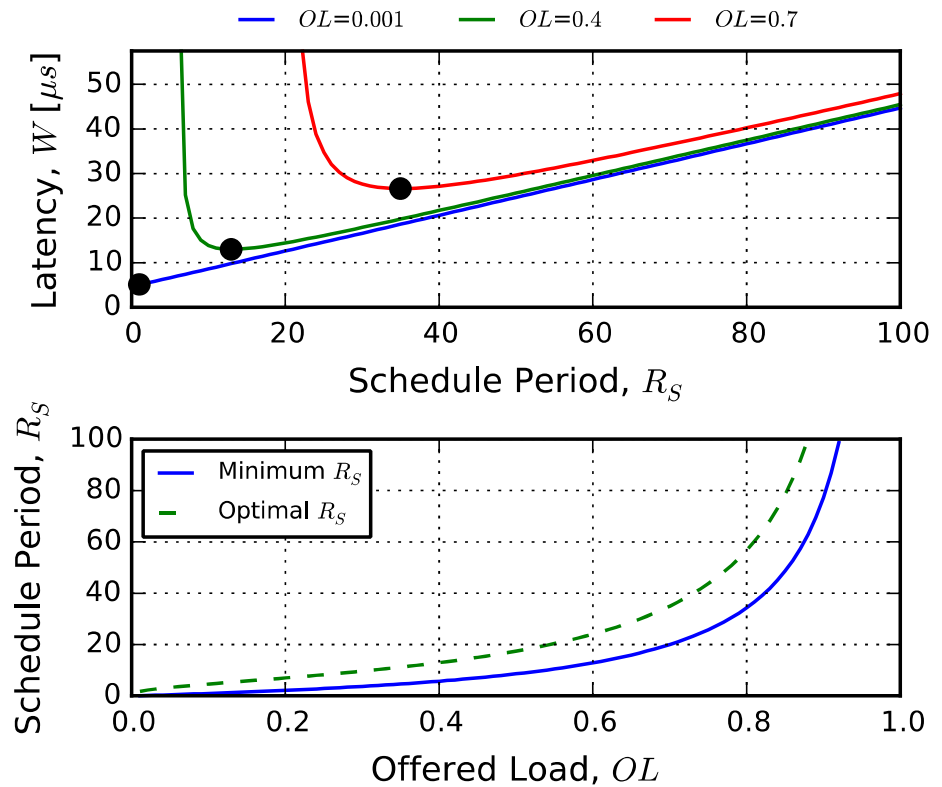
  $C = 14$

  $N = 8C$

  $S = 120$ clock cycles
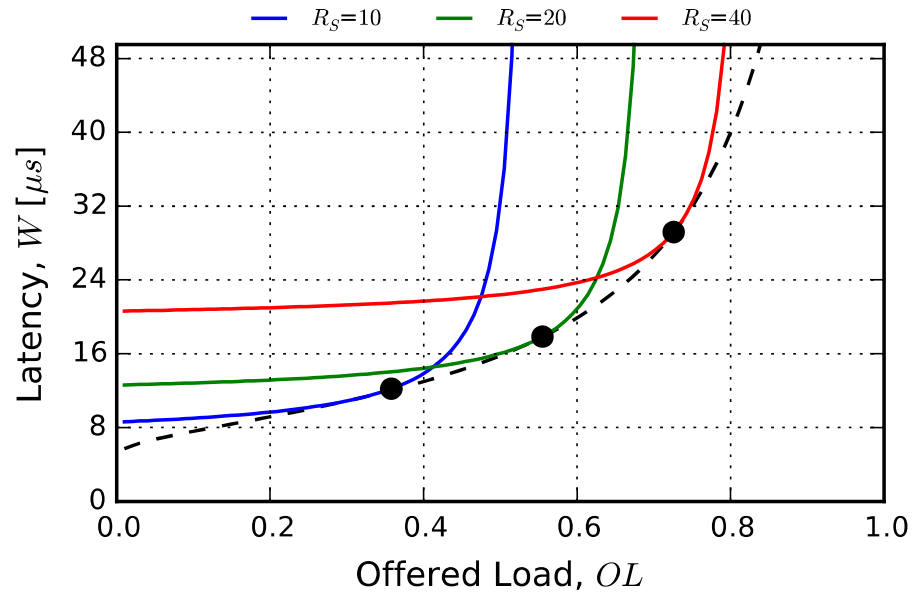
- Sweep $R_S$

# Latency prediction and optimization

- Model parameters:

  $C$ = 14
  $N$ = 8C
  $S$ = 120 clock cycles

- Sweep $R_S$

- Optimize Latency



$$OL \propto \lambda$$

**33rd IEEE International Conference on Computer Design**
*ICCD 2015*  17

# Latency prediction vs. offered load

- Model parameters:

    $C$ = 14
    $N$ = 8C
    $S$ = 120 clock cycles

- Sweep $OL$

- Optimize Latency



OL $\propto \lambda$

# Conclusion

- Developed a vacation-based M/G/1 queueing model for virtualized custom logic functions

- The model predicts throughput, latency, and queue occupancy

- Inputs to the model are the circuit, technology, clock period, pipeline depth, number of contexts, schedule period, input arrival rate, and overhead of a context switch

# Future directions

- Evaluate the assumption that the input process is Poisson; many real systems may act differently by buffering up data and sending in bursts

- Extend the model to support additional scheduling algorithms; the current schedule is not work-conserving, meaning that an empty input queue will still get scheduled

# Questions?