

**DESIGNCON<sup>®</sup> 2015**



# ***CRITICAL MEMORY PERFORMANCE METRICS FOR DDR4 SYSTEMS***

**Barbara Aichinger**

**Vice President New Business Development**

**FuturePlus Systems Corporation**



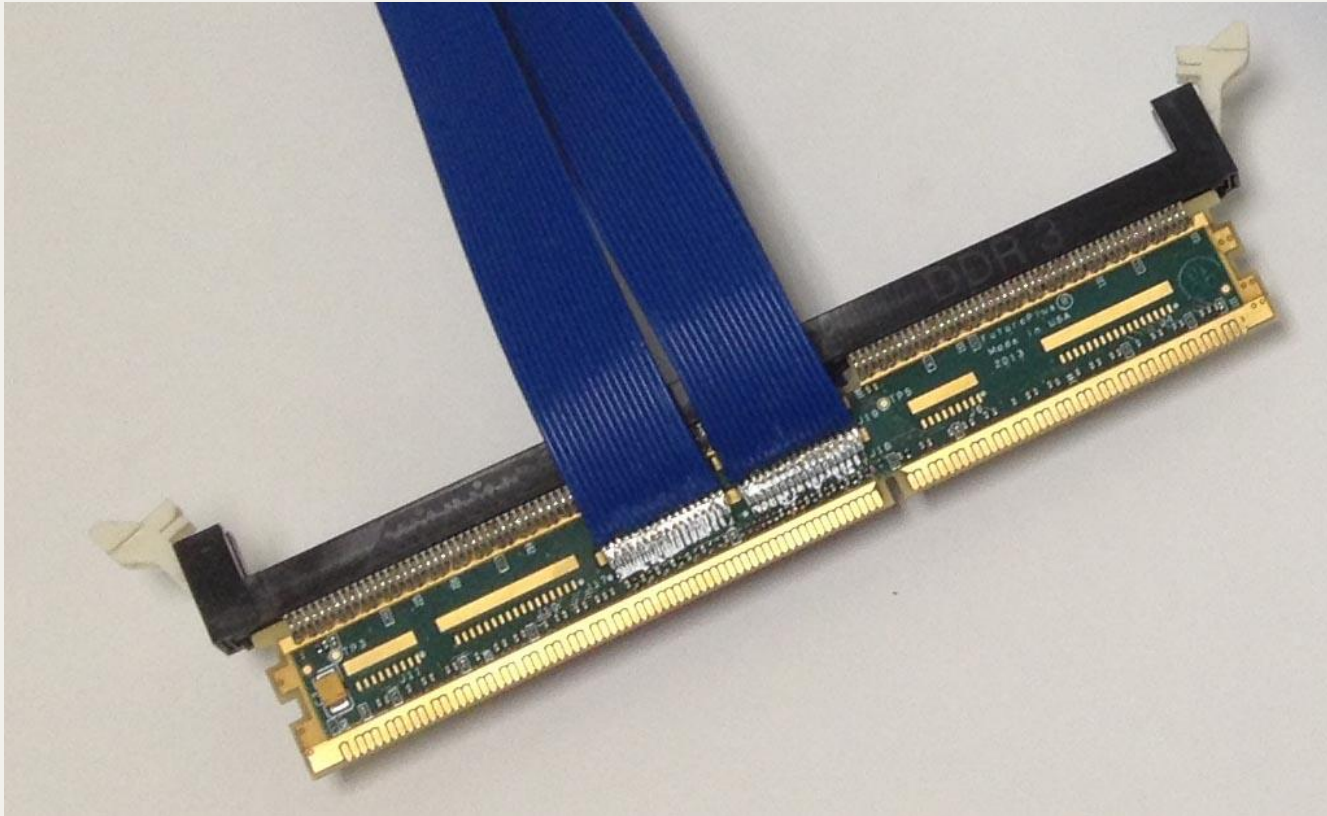
## Outline

- How can we monitor the very fast low voltage DDR4 memory without effecting the system?
- Traditional Performance Metrics
  - Bandwidth, Latency and Power Management
- New Performance Metrics
  - Bus Mode Analysis
  - Page Hit Analysis
  - Multiple Open Banks Analysis
  - Bank Utilization Analysis
  - Row Hammer (excessive Activate) Detection
- Measured DDR4 Performance in a real target

## How to Monitor the DDR4 Memory

- Use a slot interposer to ‘listen’ to the traffic between the DIMM and the Memory Controller
  - A small amount of current is ‘tapped’ off the bus
  - Only the Address, Command and Control bus needs to be monitored
  - Same probing method can be used for SODIMM and ‘memory down’ using a BGA adapter

## DDR4 DIMM Interposer

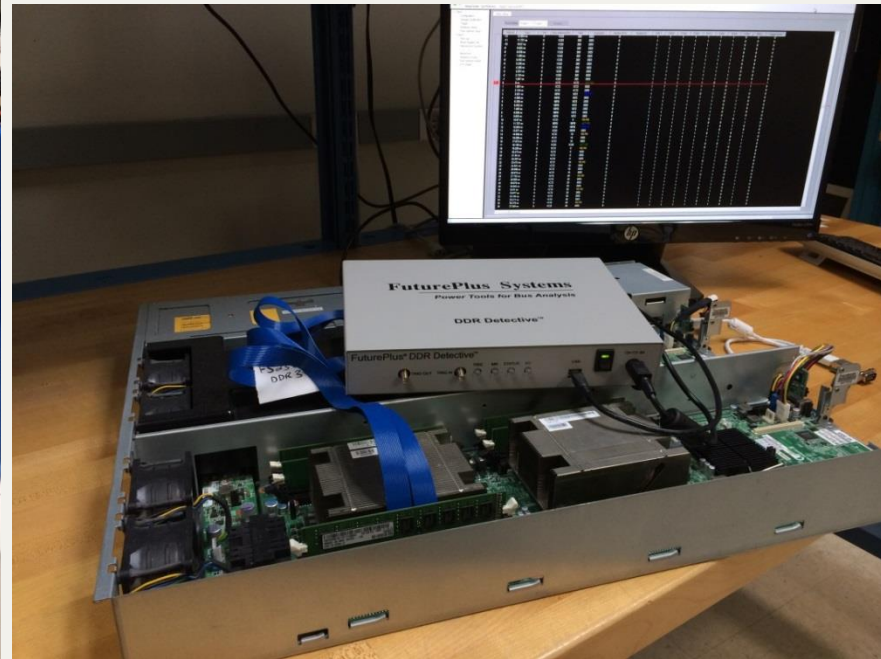
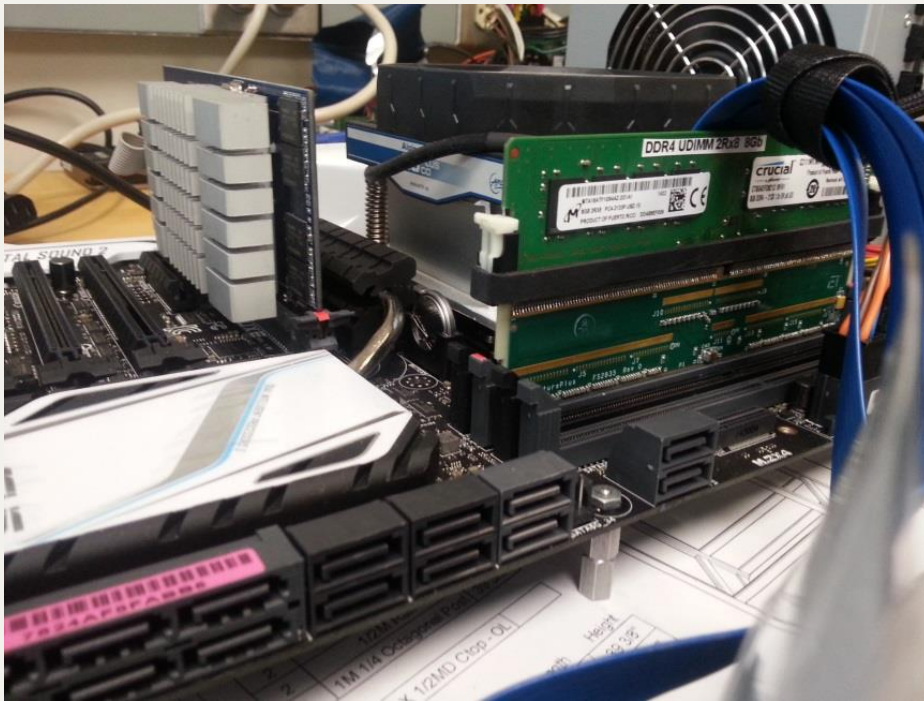


## Monitoring the DDR4 Memory

- Highly attenuated signals are then amplified and the bus is 'replicated' inside the equipment



## The system boots and runs never knowing the equipment is present



## Performance Metrics

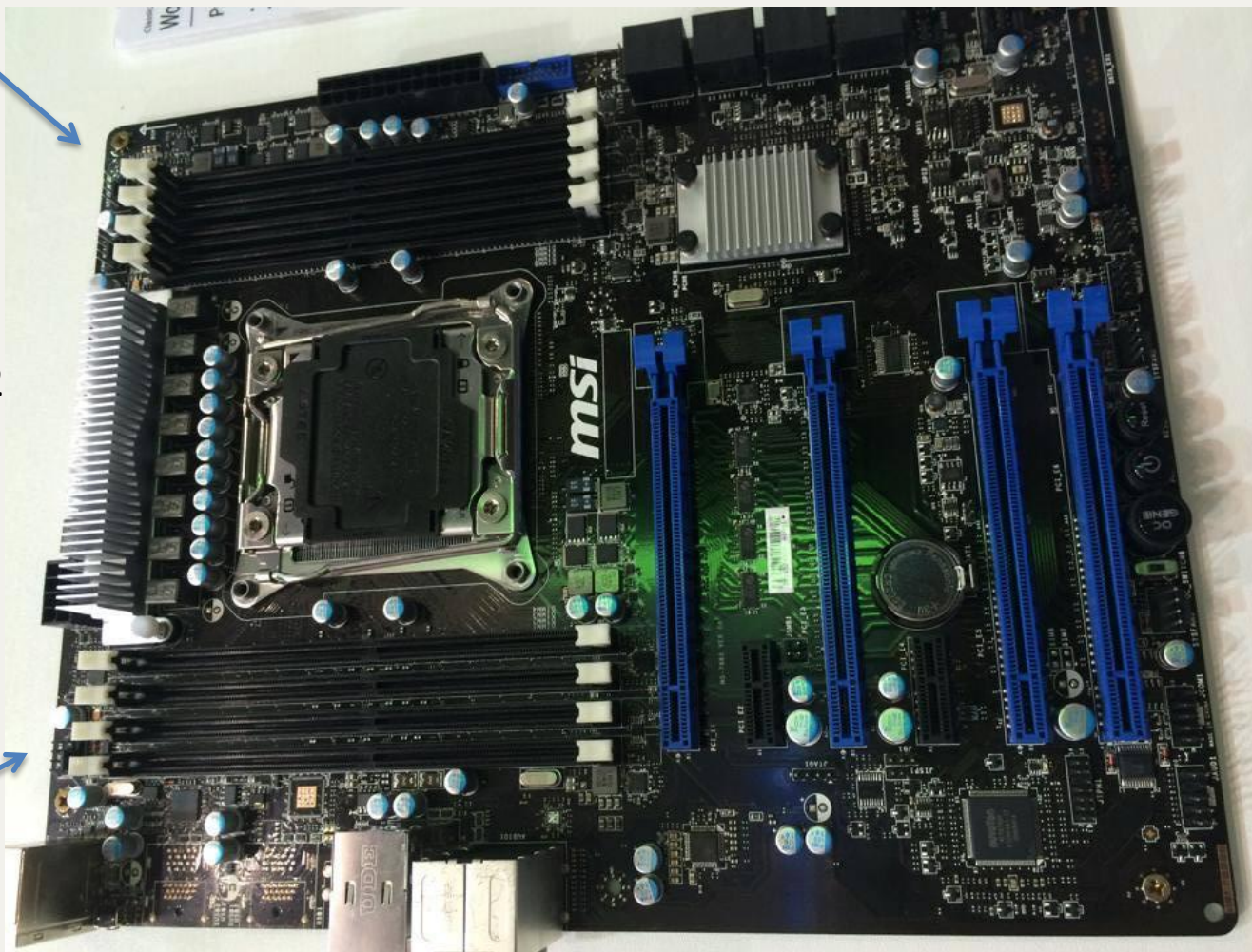
- Billions of clock cycles every second
- Old methods of measurement captures small amount of time in large trace buffers
- New method uses sophisticated event counters to:
  - Never miss a clock cycle
  - Count Performance Metrics for long periods of time REAL TIME

## Work Smarter not Harder

- For Performance metrics the DDR Detective<sup>®</sup> uses counters instead the traditional trace memory
  - To capture a second of DDR4 traffic would take 4.5Gbytes of logic analyzer/protocol analyzer trace depth \$\$\$\$!
    - 1 hour = 270 Gbytes of trace depth and then time to sift through it and post process!
  - By using large counters and counting events and the time between events we can achieve hours and days worth of metrics with no trace buffer memory and with no time consuming post processing



4 Memory Channels  
each  
channel is 2  
slots



ASUS X99  
DDR4  
Motherboard

## Traditional Measurements

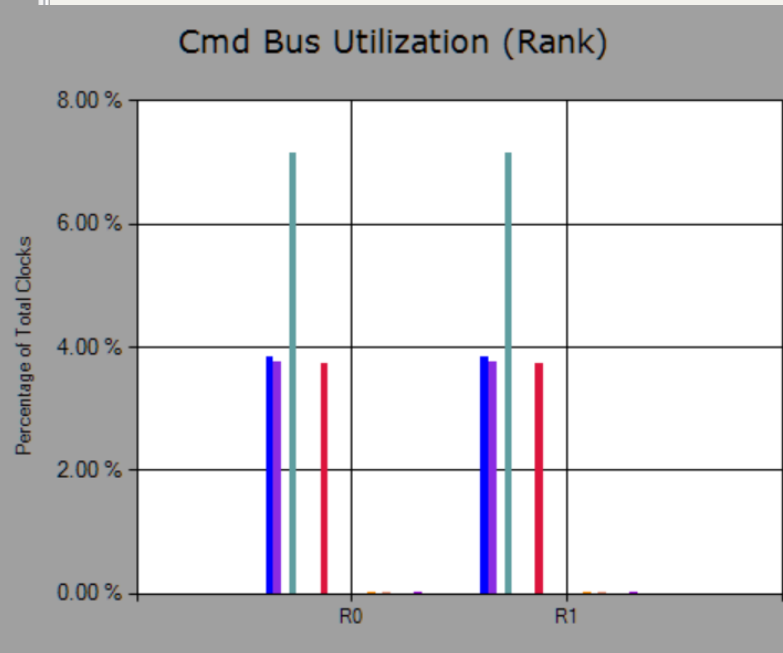
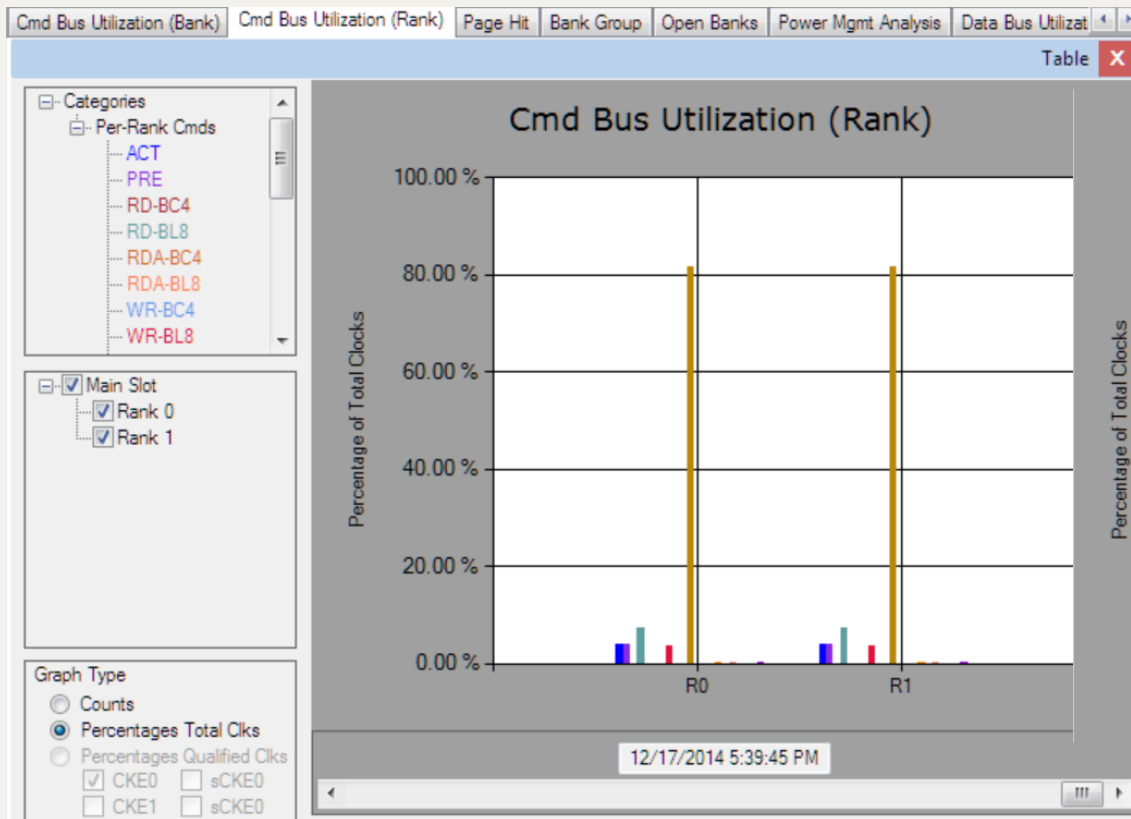
- Bandwidth
  - Command Bus Utilization
  - Data Bus Utilization
- Power Management
- Latency

## Bandwidth

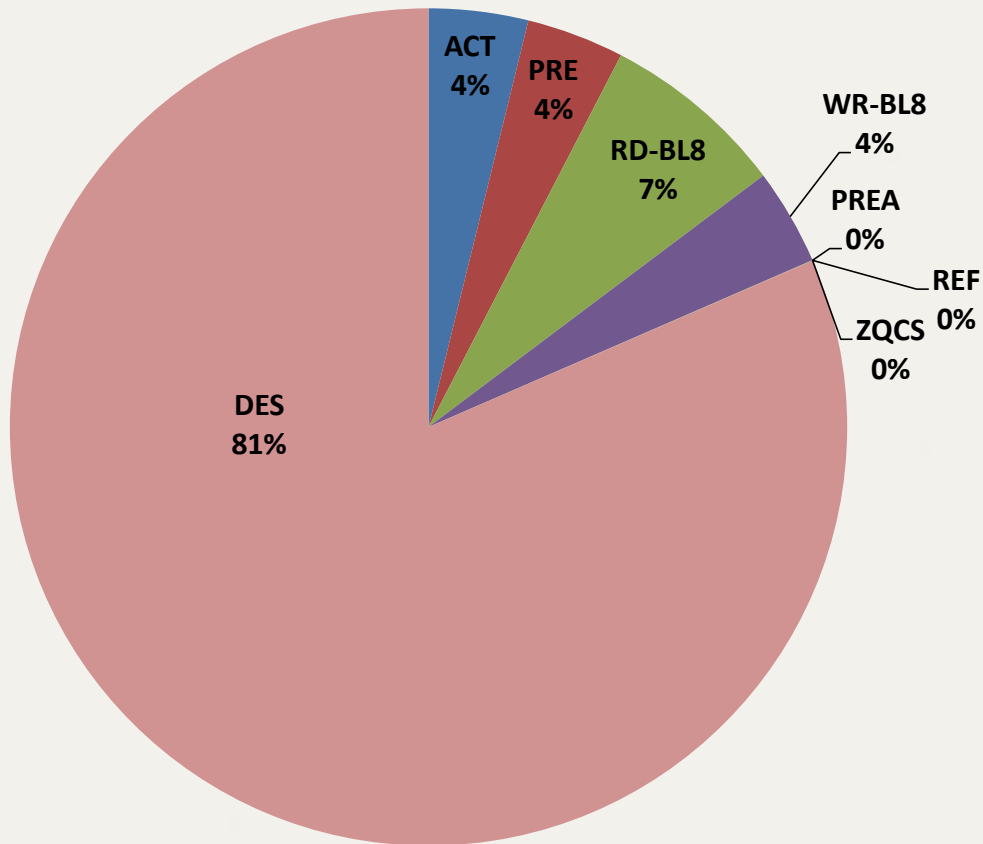
- Overhead
  - Any use of the bus other than a Read or a Write
  - Command Bus Utilization
- Data Bus
  - Utilization: the % of the time that Read or Write Data is being transferred
  - Bandwidth: the amount of data transferred per second

## Command Bus Utilization

2400MT/s

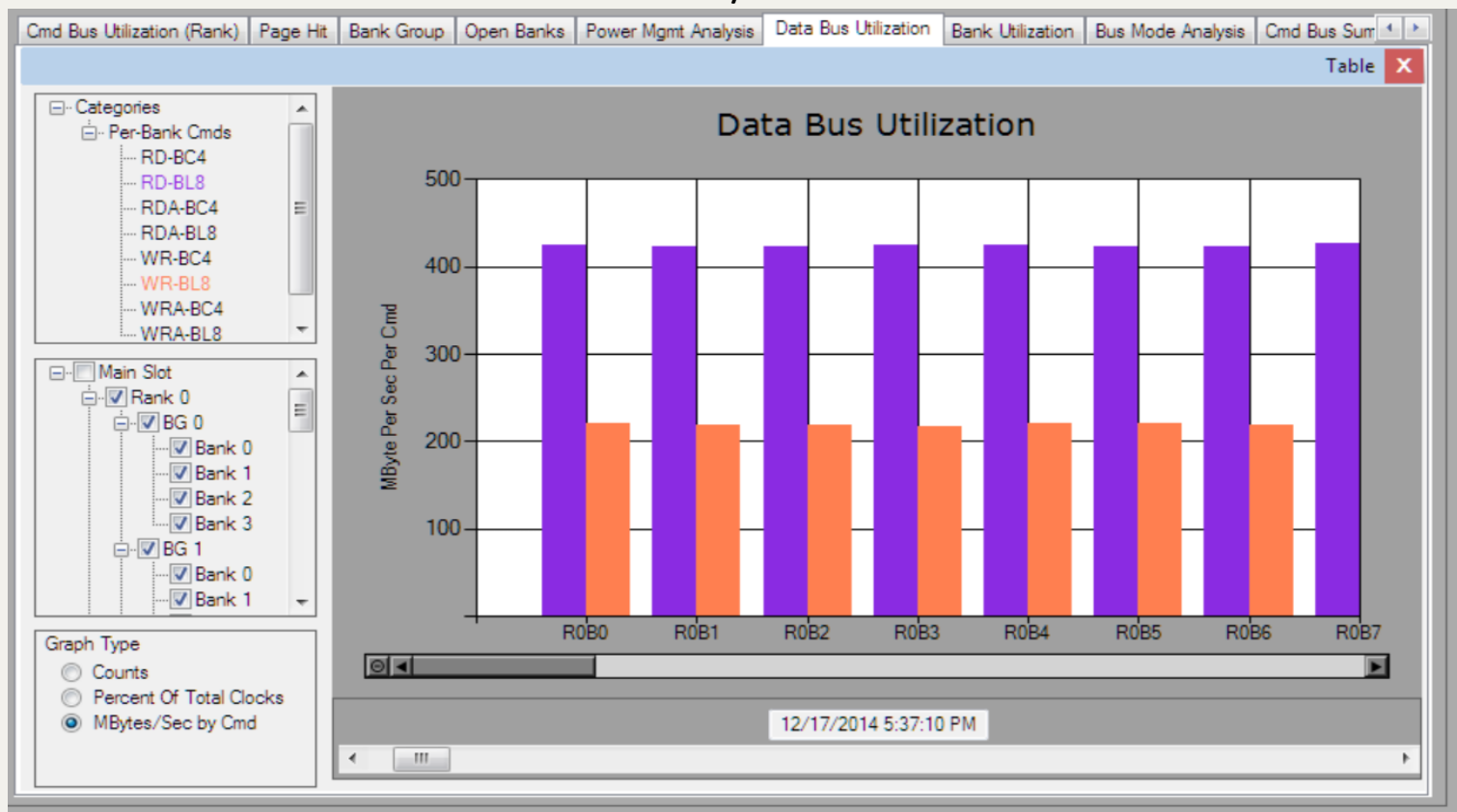


## Command Bus Utilization

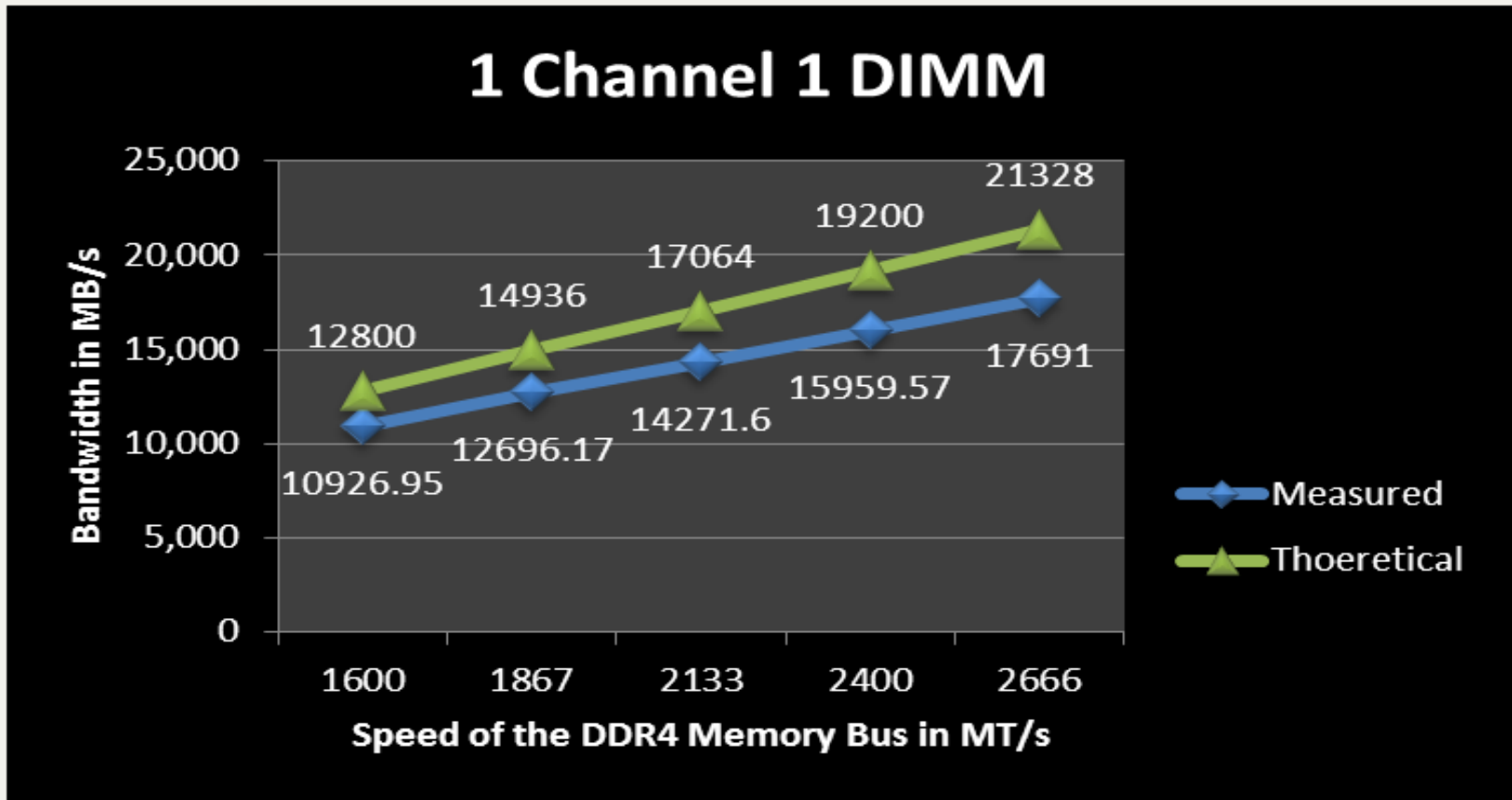


## Data Bus Utilization

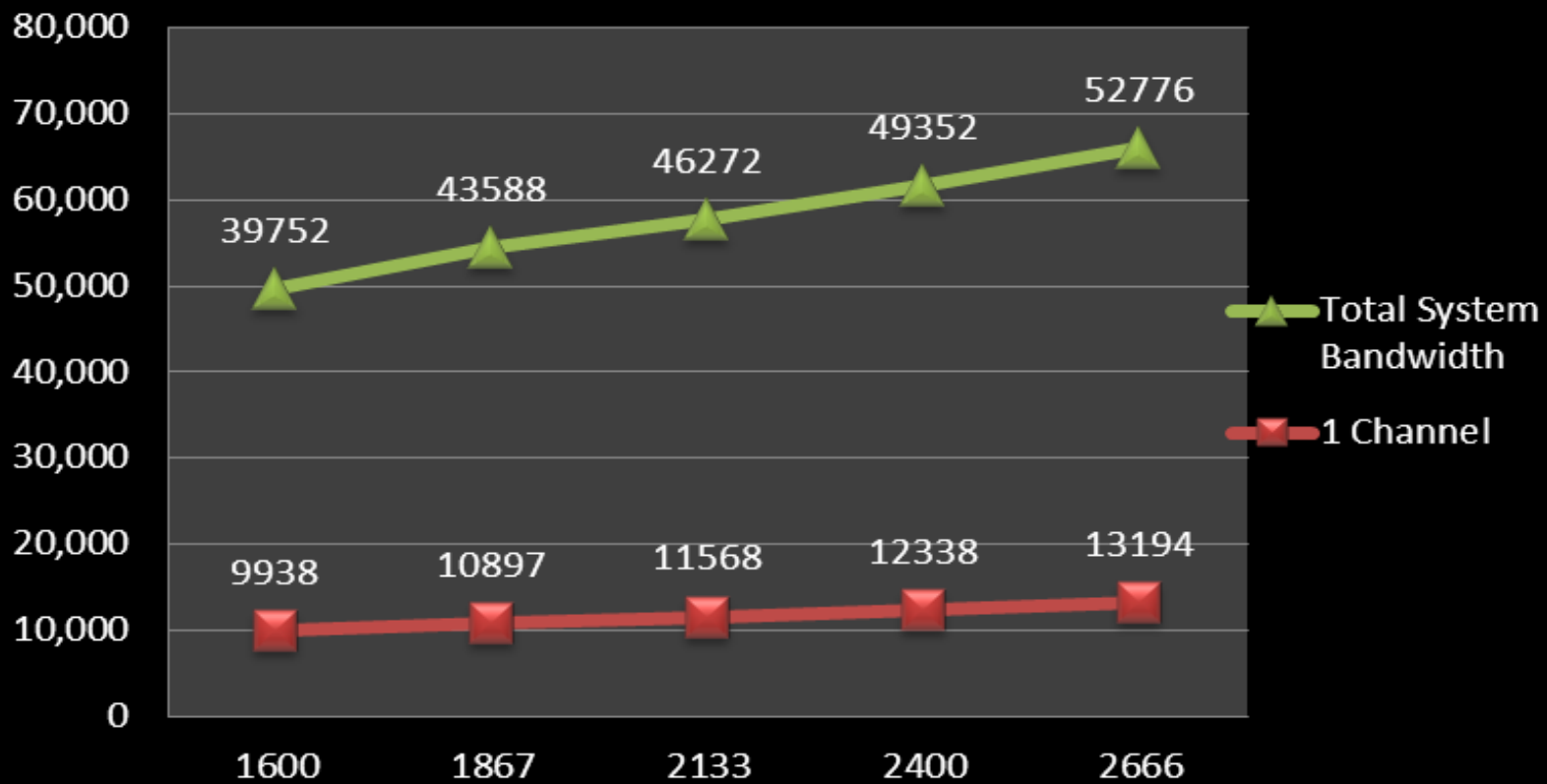
2400MT/s



## DDR4 Bandwidth

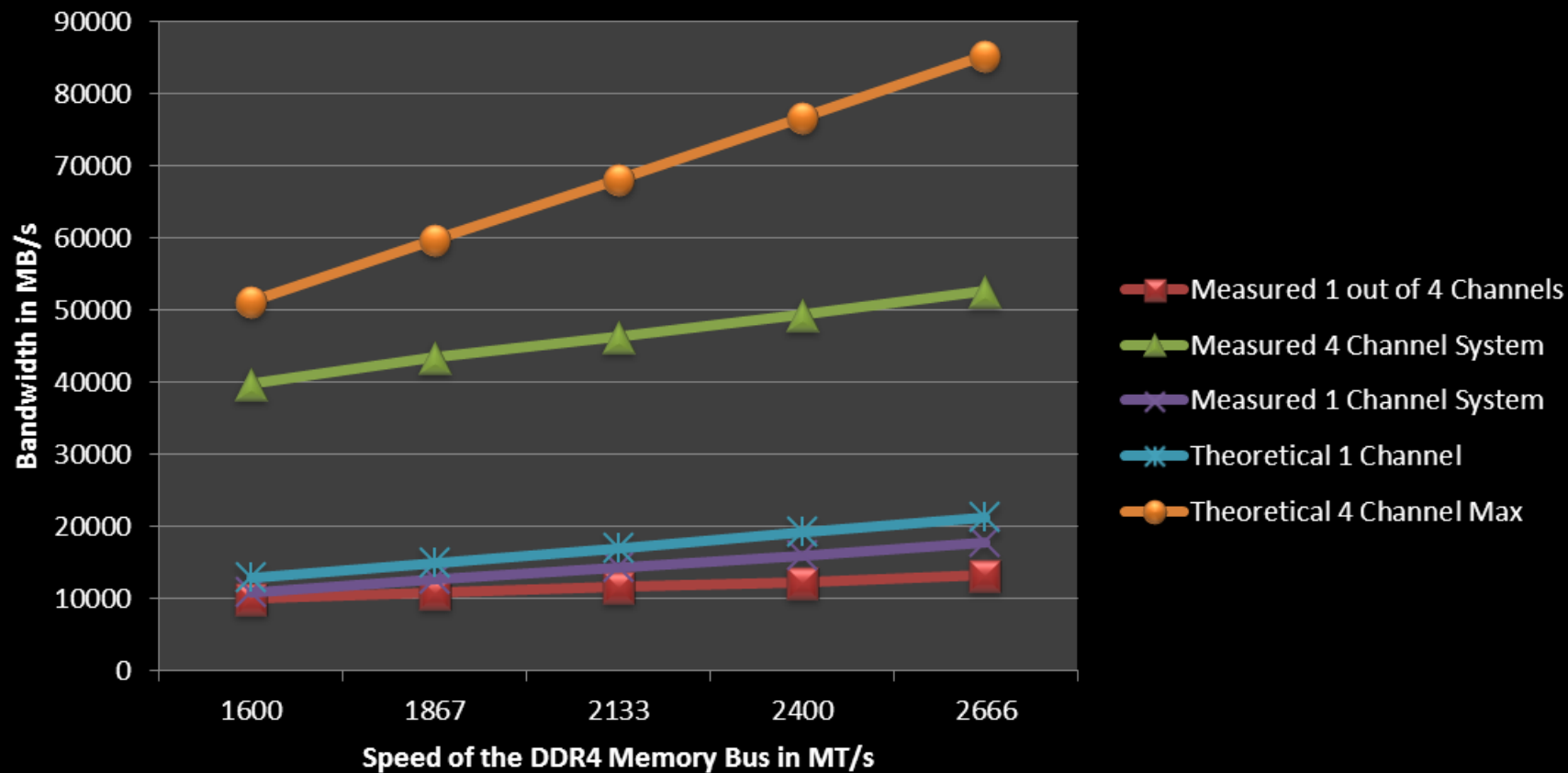


## 4 Channels 1 DIMM per Channel

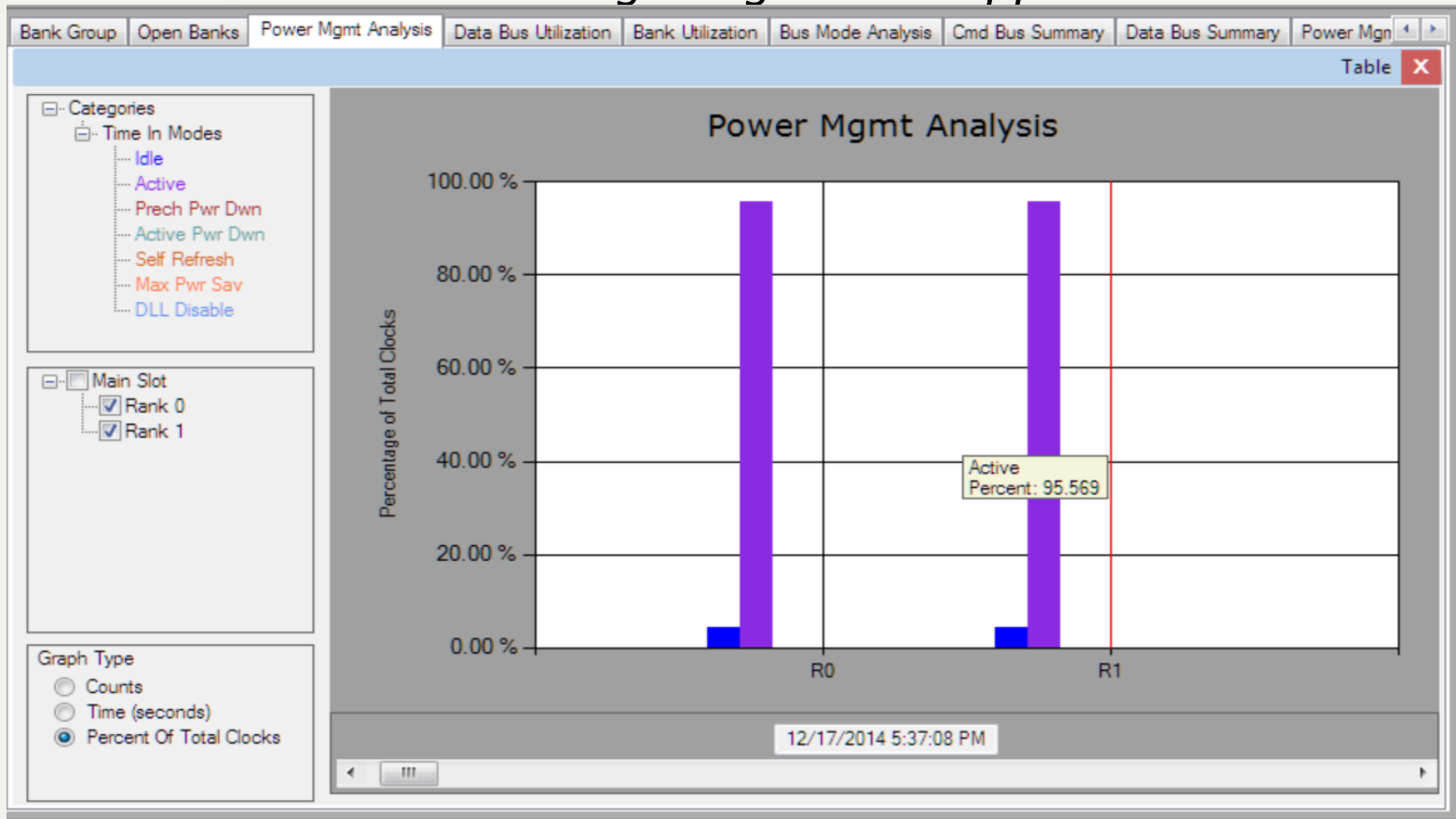




## Bandwidth 4 Channel vs 1 Channel System

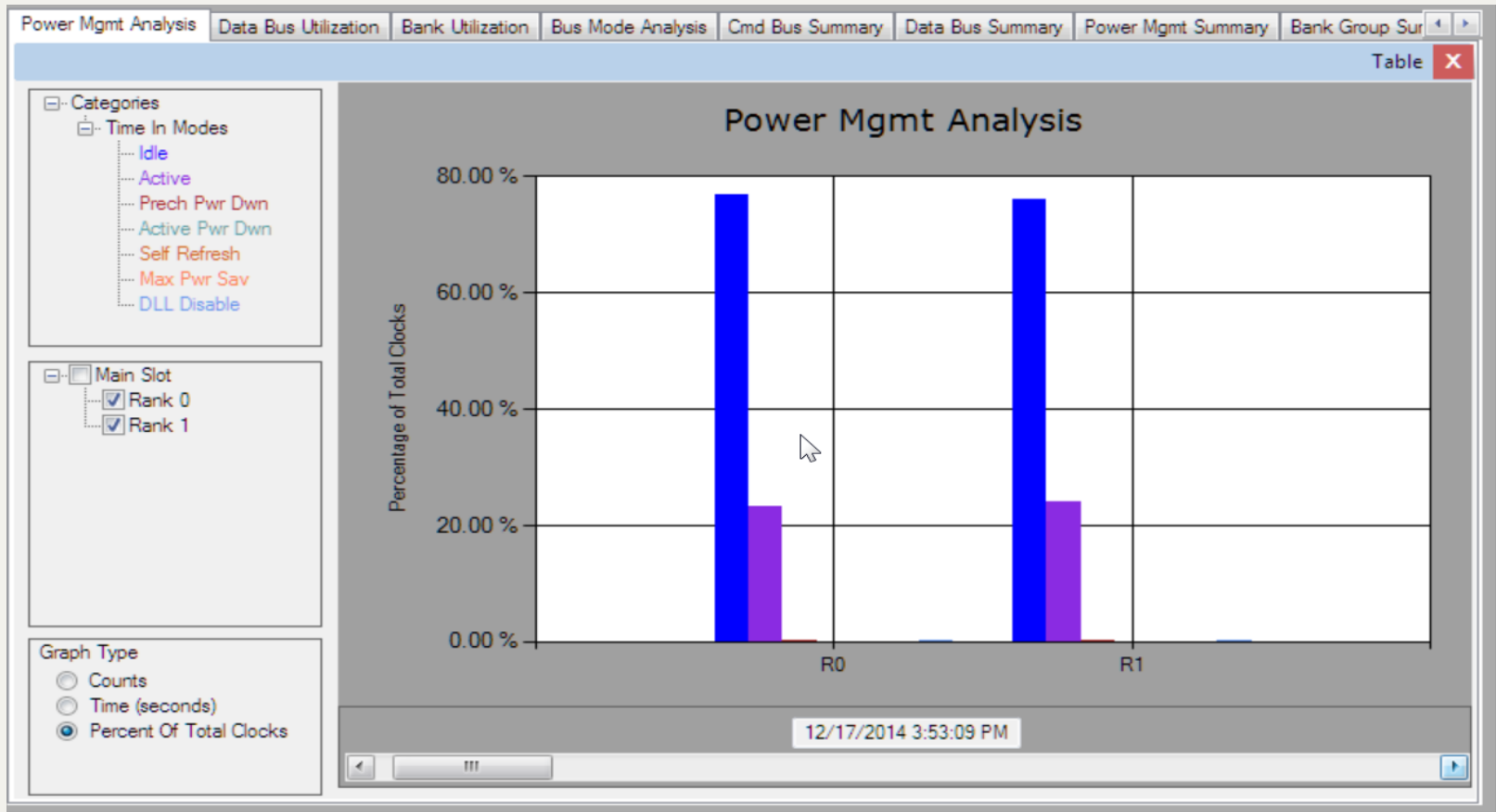


## Power Management *while running Google StressApp*



# Power Management Analysis

## Ubuntu Boot



## Power Management

- ~50M servers Servers World Wide
- Each Server averages 16-24 DIMMs
  - 800M to 1.2B DIMMs
- Even small power savings per DIMM can add up

*Every time **Facebook's** data center engineers figure out a way to **reduce** server consumption **by a single watt**, the improvement, at Facebook's scale, has the potential to add **millions of dollars** to the company's bottom line.*



Yevgeny Sverdlik  
Editor in Chief  
Data Center Knowledge

# Latency

- Several Jedec Parameters apply:
  - RD to WR same rank tSR\_RTW
  - RD to PRE/PREA same Rank tRTP
  - WR to PRE(SB) or PREA (SR) tWR
  - Read to Read different Rank tDR\_RTR
  - Read to Write different Rank DR\_RTW
  - Write to Read different Rank tDR\_WTR
  - Write to Write different Rank tDR\_WTW

D to Same Bank Group  
 R to Same Rank  
 /R Same Bank Group  
 D Different Bank Group  
 /R Different Bank Group  
 to ACT Same Bank Group  
 to ACT Different Bank Group  
 tFAWmin Same Rank  
 D Same Bank Group  
 RD Different Bank Group  
 PRE or PREA Same Rank  
 PRE(SB) or PREA (SR)  
 reset to any Command or CKE  
 ODT Enabled  
 MRS  
 Other Command or ODT High  
 ODT Enabled  
 and PAR\_IN have odd # of 1's  
 CL after Reset Low to High  
 ODT Enabled  
 1st ZQCL after Reset Low to High  
 ODT Enabled  
 om ZQCS to any Command  
 ODT Enabled  
 low to High, then CKE Low to High  
 ODT Enabled  
 Command  
 Non-Deselect  
 ZQCL or ZQCS  
 RD or CKE Low or ODT Hi  
 ODT Enabled  
 SRX  
 Non-Deselect  
 inimum Pulse Width  
 PDX is less than tPDmin  
 PDX is greater than tPDmax  
 PDE

- V31 - PRE or PREA to PDE
- V32 - RD or RDA to PDE
- V33 - WR to PDE < tWRPDEN\_cc
- V34 - WRA to PDE < tWRAPDEN\_cc
- V35 - WR to PDE < tWRPBC4EN\_cc
- V36 - WRA to PDE < tWRAPBC4EN\_cc
- V37 - REF to PDE
- V38 - MRS to PDE
- V39 - PRE, PREA to ACT, Other Commands
- V40 - ACT to PRE, PREA (Min)
- V41 - ACT to PRE or AutoPRE (Max)
- V42 - ACT to RD or WR
- V43 - ACT to ACT or REF
- V44 - REF to non-DES
- V45 - REF to REF Average Interval
- V46 - REF to REF Maximum Interval
- V47 - Read or Write to an Inactive Bank
- V48 - Refresh to an Active Bank
- V49 - Activate to an Active Bank
- V50 - MRS with an Active Bank
- V51 - Self Refresh Entry with an Active Bank
- V52 - ZQCS or ZQCL with an Active Bank
- V53 - Read to Read (Different Rank)
- V54 - Read to Read (Different DIMM)
- V55 - Read to Write (Different Rank)
- V56 - Read to Write (Different DIMM)
- V57 - Write to Read (Different Rank)
- V58 - Write to Read (Different DIMM)
- V59 - Write to Write (Different Rank)
- V60 - Write to Write (Different DIMM)
- V61 - ODT high to Low Time < ODT\_H4\_c
- V62 - ODT high to Low Time < ODT\_H8\_c
- V63 - PDE or SRE Followed by non-DES
- V64 - WRA, WRA to Command
- V65 - RDA to Command

tSR\_RTW 8 nC  
 V11 - RD to PRE or PREA Same Rank  
 tRTP 8 nC  
 V12 - WR to PRE(SB) or PREA (SR)  
 tWR 31 nC  
 V53 - Read to Read (Different Rank)  
 DR\_RTR 5 nC  
 V55 - Read to Write (Different Rank)  
 DR\_RTW 7 nC  
 V57 - Write to Read (Different Rank)  
 DR\_WTR 3 nC  
 V59 - Write to Write (Different Rank)  
 DR\_WTW 5 nC

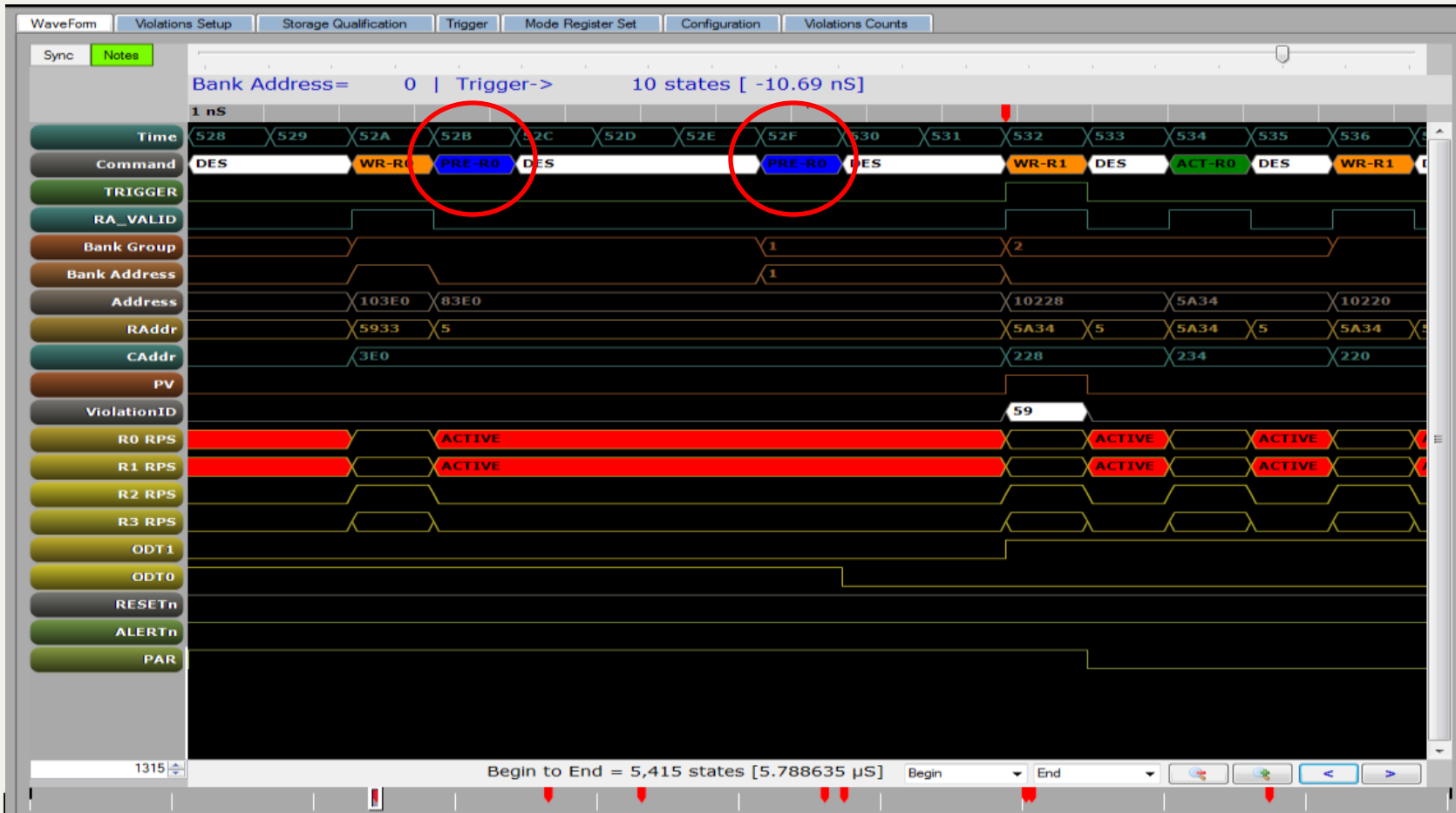
DDRx Detective - DDR4  
 File Help  
 Setup Guide Eye Detector Apply Captured MRS  
 Input: Configuration, Storage Qualification, Trigger, Violations Setup  
 Output: Run Log, Mode Register Set, State Listing, WaveForm, Violations Counts, FPS Charts  
 Violations Counts | Violations Setup | Configuration | Trigger  
 Violation 59: Write to Write (Different Rank)  
 R0: 0, R1: 0, R2: 0, R3: 0, Total: 762

# Latency Measurements

measurement made at 1867

V#	Parameter	Description	Spec	Measured
V2	tSR_RTW	RD to WR same Rank	8	10
V11	tRTP	RD to PRE same Rank	8	8
V12	tWR	WR to PRE SB or PREA SR	31	31
V53	tDR_RTR	RD to RD diff Rank	5	6
V57	tDR_WTR	WR to RD diff Rank	3	6
V59	tDR_WTW	WR to WR diff Rank	5	8

## Intervening Commands





# Latency Measurements

V#	Parameter	Description	Spec	Measured
V1	tCCD_L	RD to RD Same Bank Group	5	6
V3	tCCD_L	WR to WR Same Bank Group	5	6
V4	tCCD_S	RD to RD diff Bank Group	4	4
V5	tCCD_S	WR to WR diff Bank Group	4	4
V6	tRRD_L	ACT to ACT Same Bank Group	5	5
V7	tWTR_L	ACT to ACT diff Bank Group	4	4
V9	tWTR_L	WR to RD Same Bank Group	22	23
V10	tWTR_S	WR to RD Diff Bank Group	17	19

## Latency

- Good designs operate on the edge of the spec
- Architectural tradeoffs will occur
- Do I need margin?
  - Design for the worst case and buy quality parts

## New Performance Metrics

- Page Hit Analysis

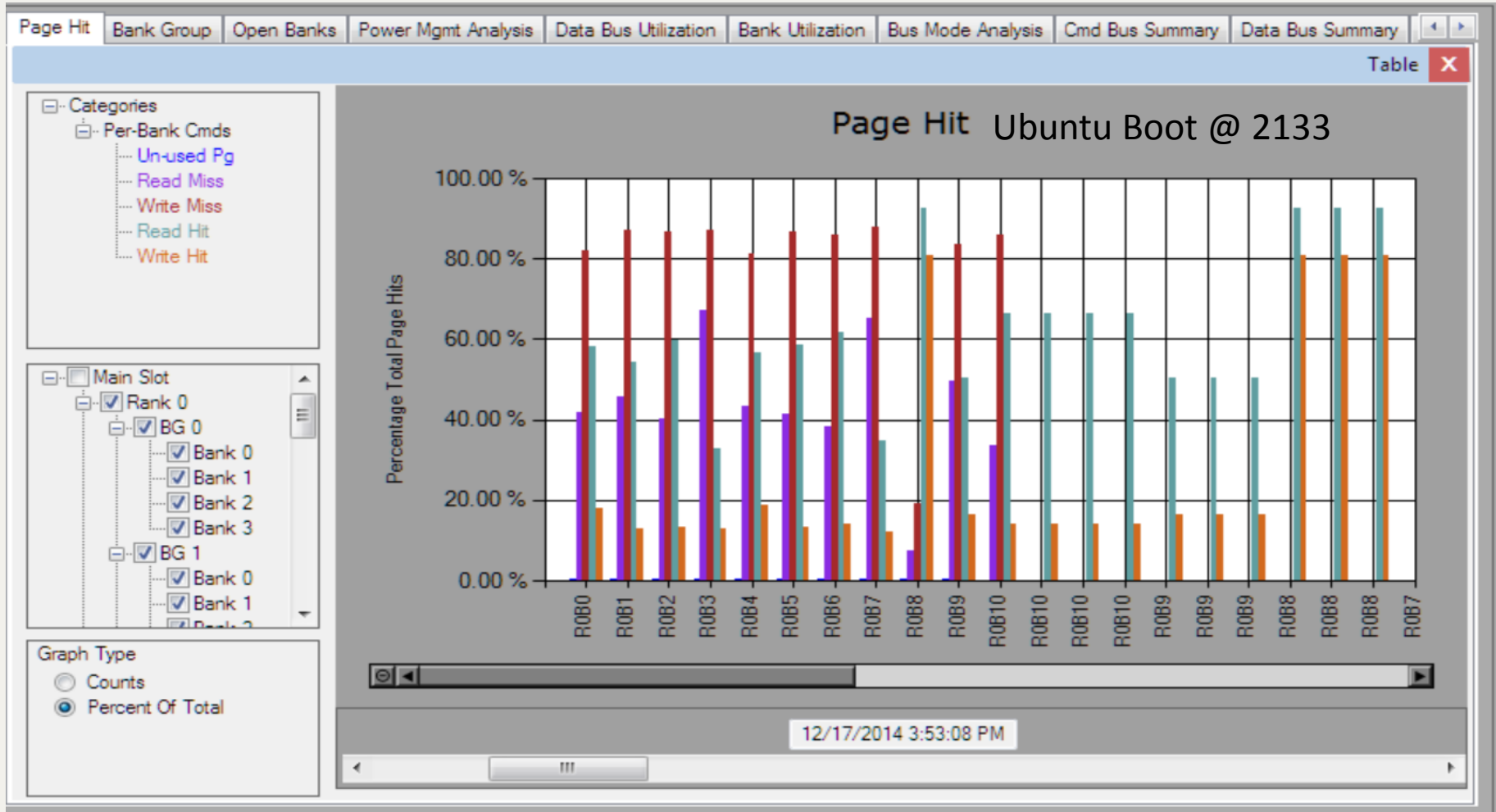
- Read Hit: Page was Open
- Read Miss : Page was not Open, Transaction was preceded by an ACT
- Write Hit: Page was Open
- Write Miss: Page was not Open, Transaction was preceded by an ACT
- Unused: Page was opened and closed and never accessed

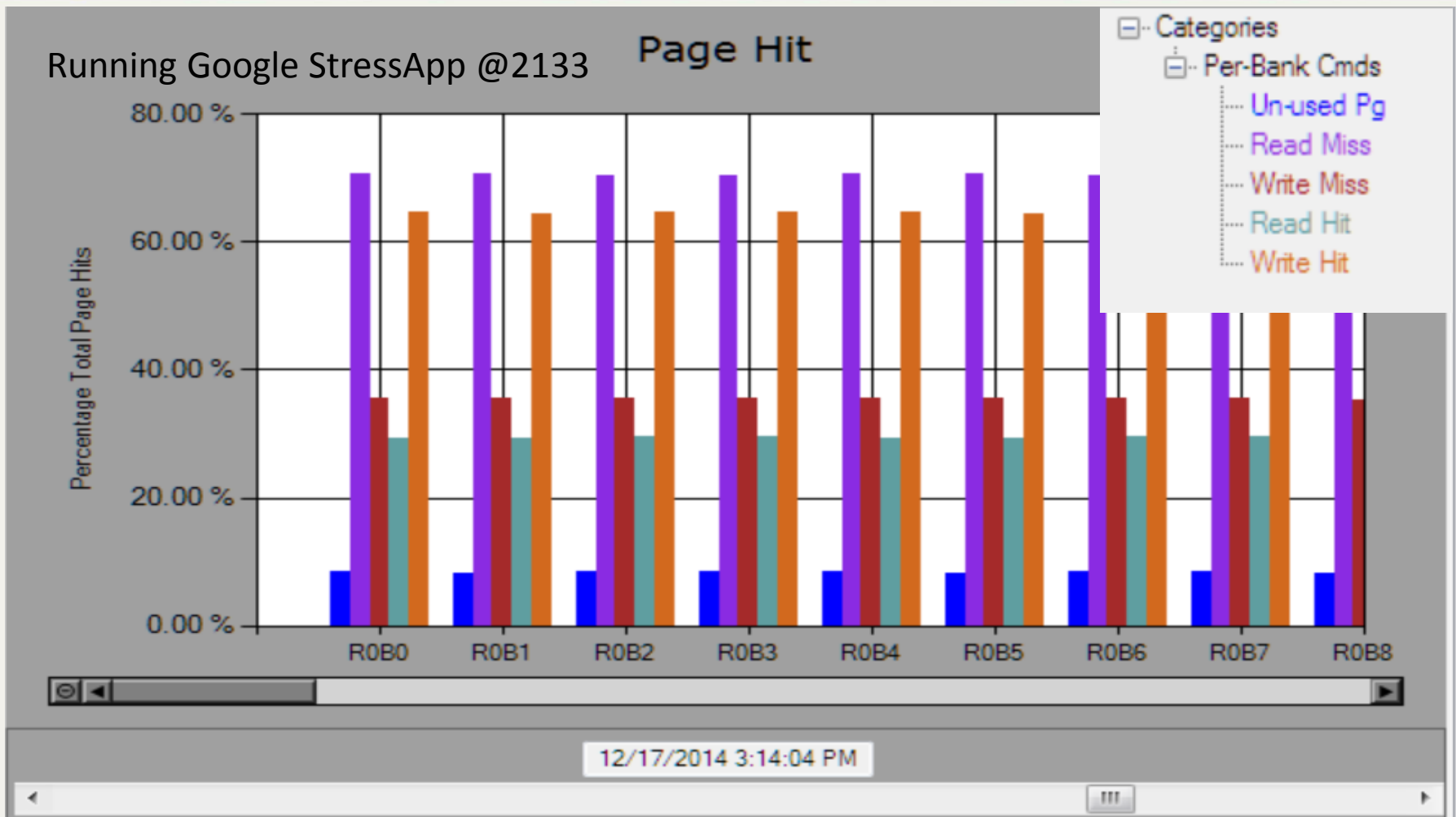
- Multiple Open Banks

- Open Banks make for faster access IF your going to that bank on the next access...performance hit if your not
- Power hit when banks are open

- Bank Group Analysis

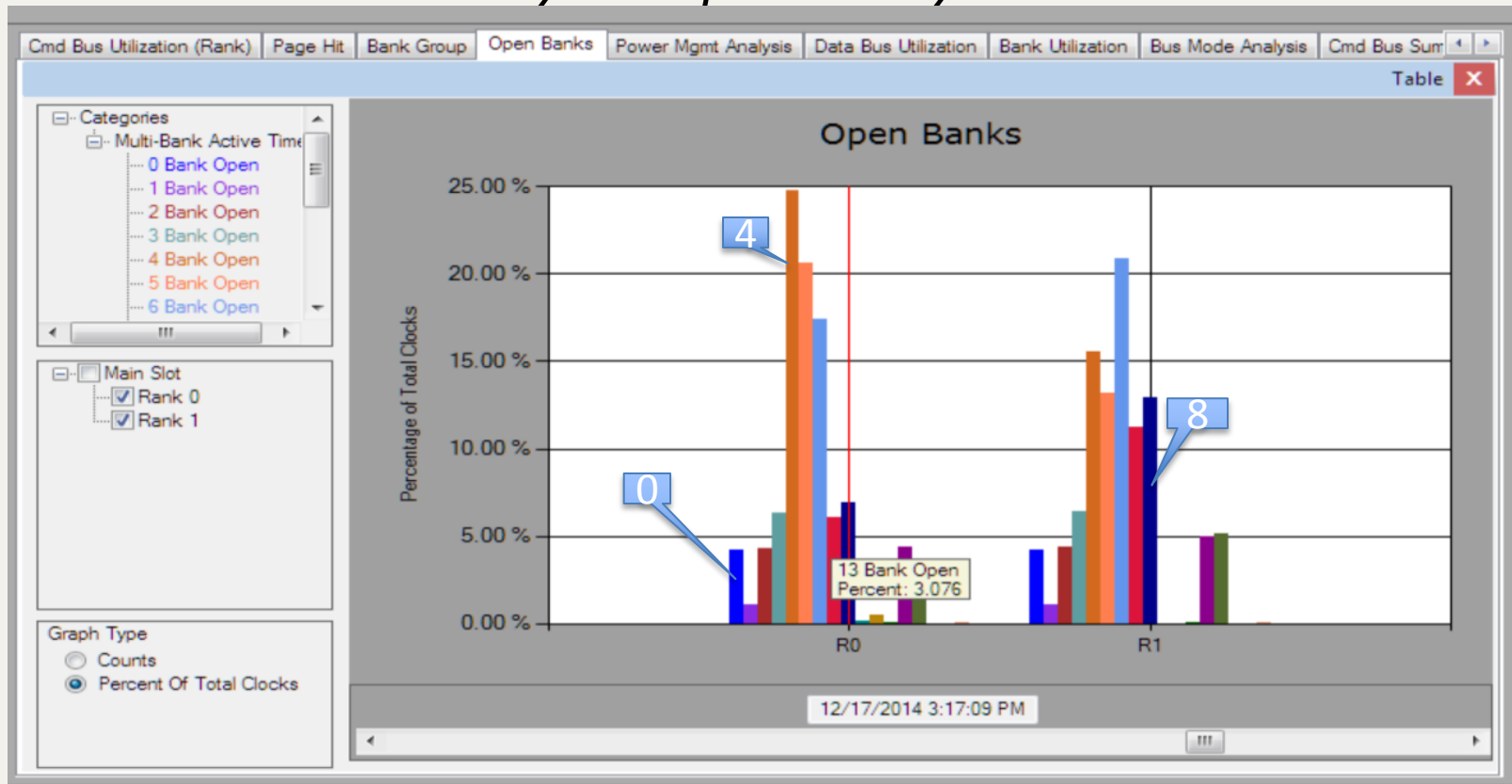
- New for DDR4: Back to back access to same bank is a performance hit
- Faster to have back to back accesses to different bank groups





## Multiple Open Banks

*How many are open at any one time*

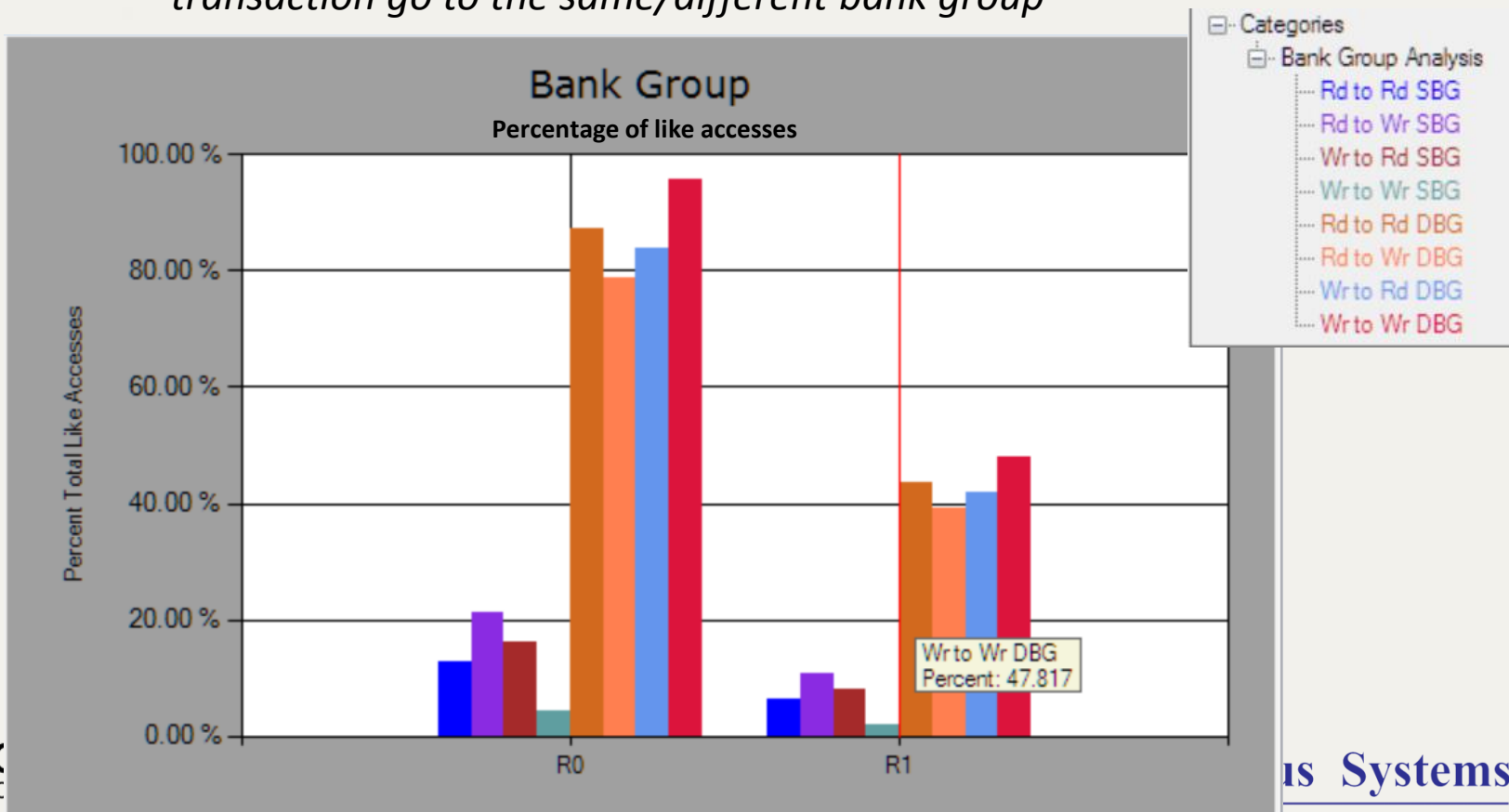


## Bank Group Access Analysis

- tCCD\_L
  - Takes longer for back to back RD/WR accesses to the same bank group
- tCCD\_S
  - Can reduce latency by going to different bank groups

## Bank Group Access Analysis

*Relative to the previous transaction how many times did the following transaction go to the same/different bank group*



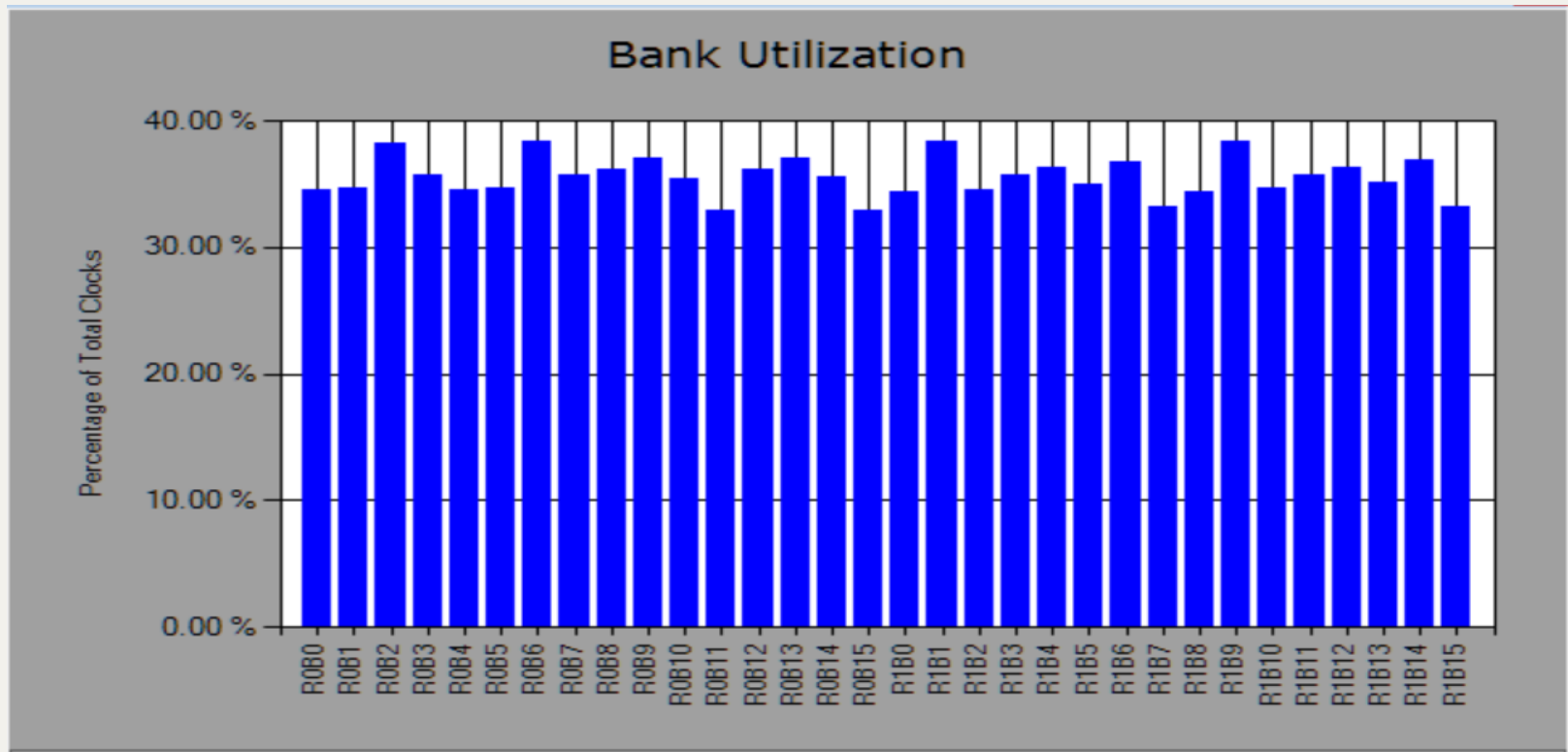


## More Performance Metrics!

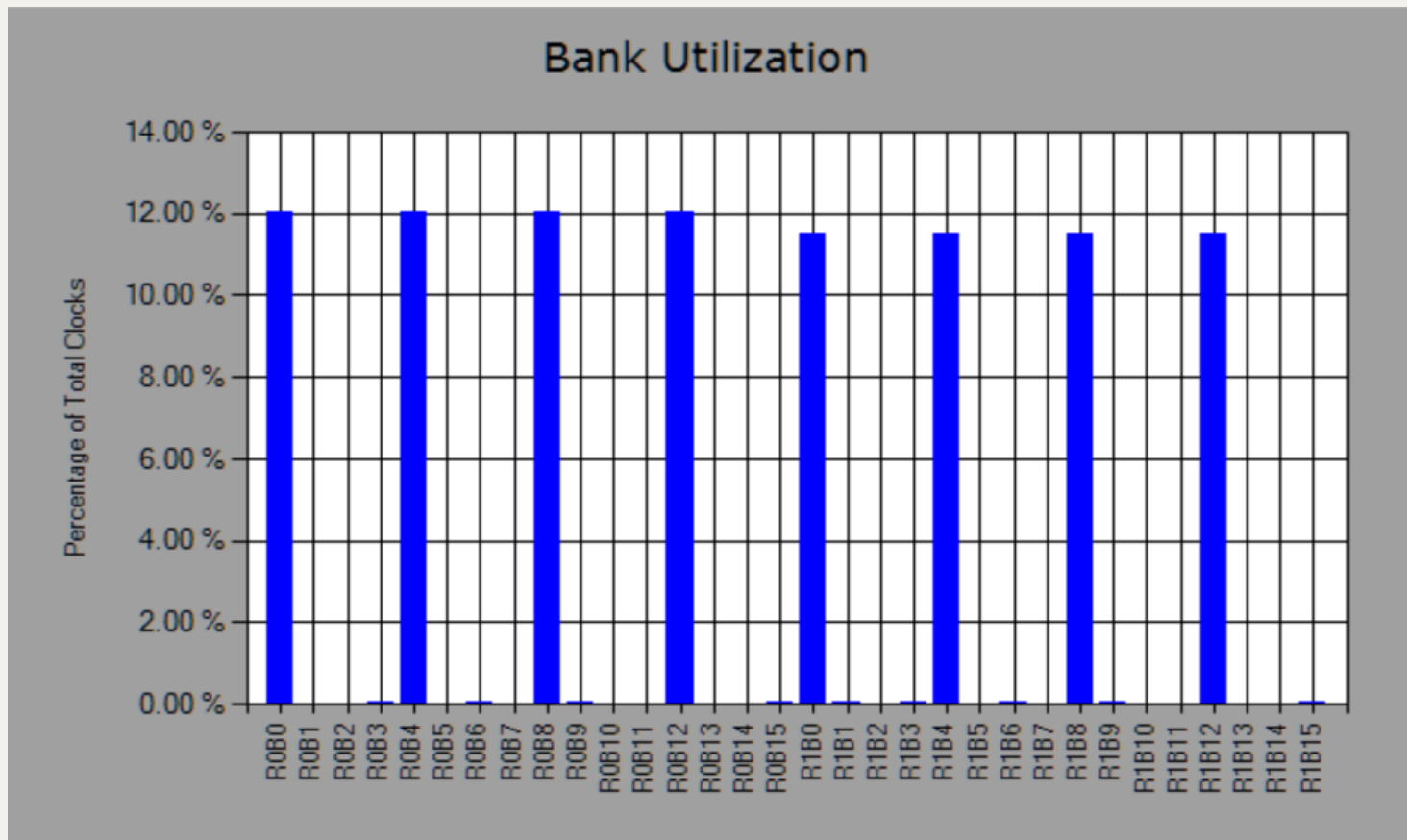
- Bank Utilization
  - What happens during a chip kill or page retirement scenario?
  - How does the traffic reallocate?
  - What are the performance implications?
- Do I have system hot spots?
  - Row Hammer (excessive Activates)
- Fast Boot
  - Why does the system take so long to boot?

## Bank Utilization

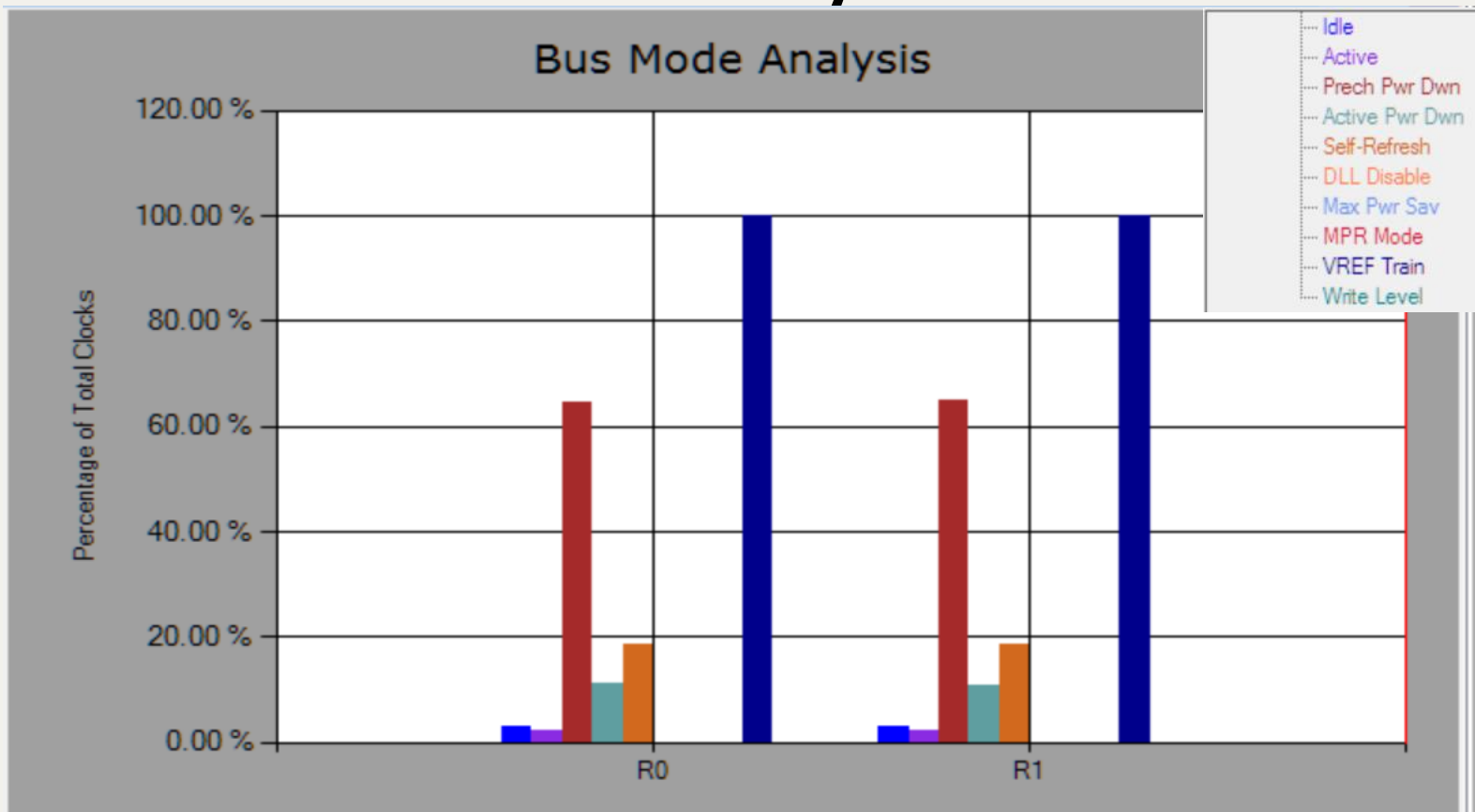
*Percentage of total cycles the banks are open (running Google Stress App)*



## Bank Utilization



## Boot Analysis



## Row Hammer

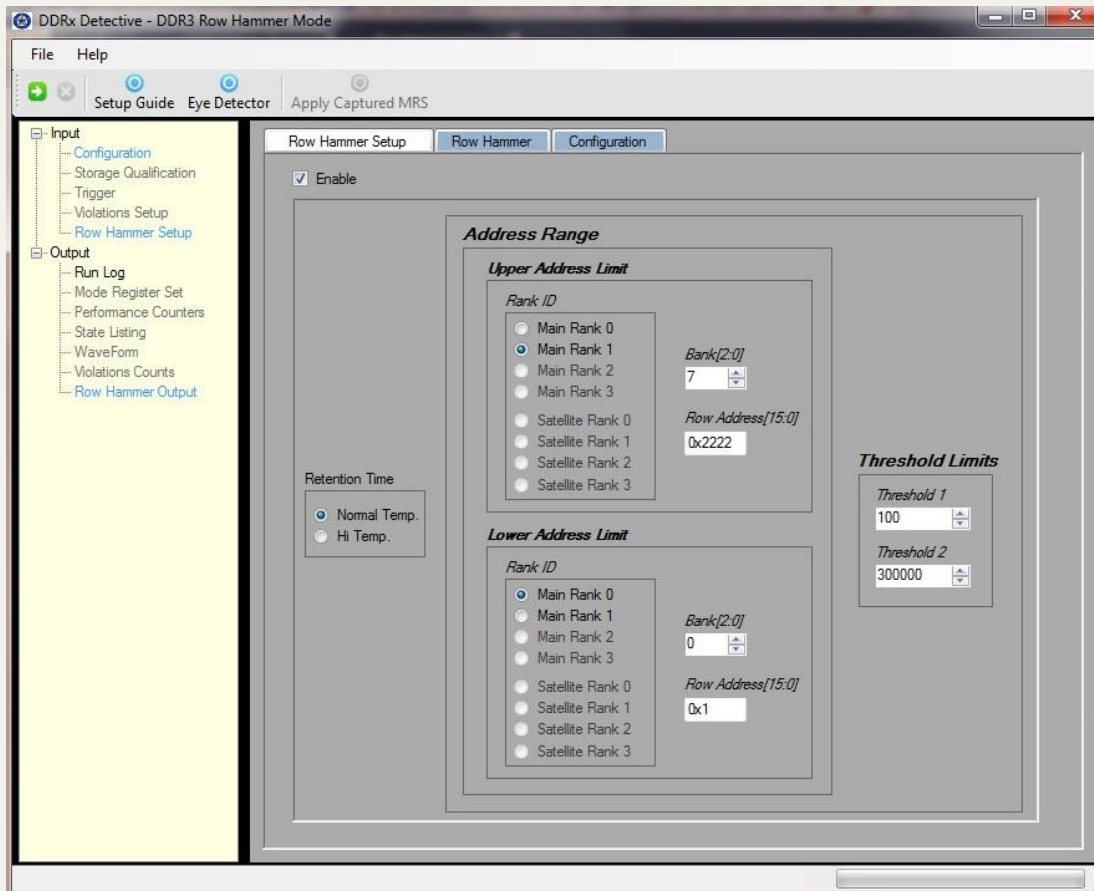
*Excessive ACTIVATE commands to a single Row*

- An identified failure mechanism in DDR3 memory
- The result of charge leakage
- Current work around is to increase the refresh rate
  - Burns power, lowers performance
  - Lowers the statistical probability of an error

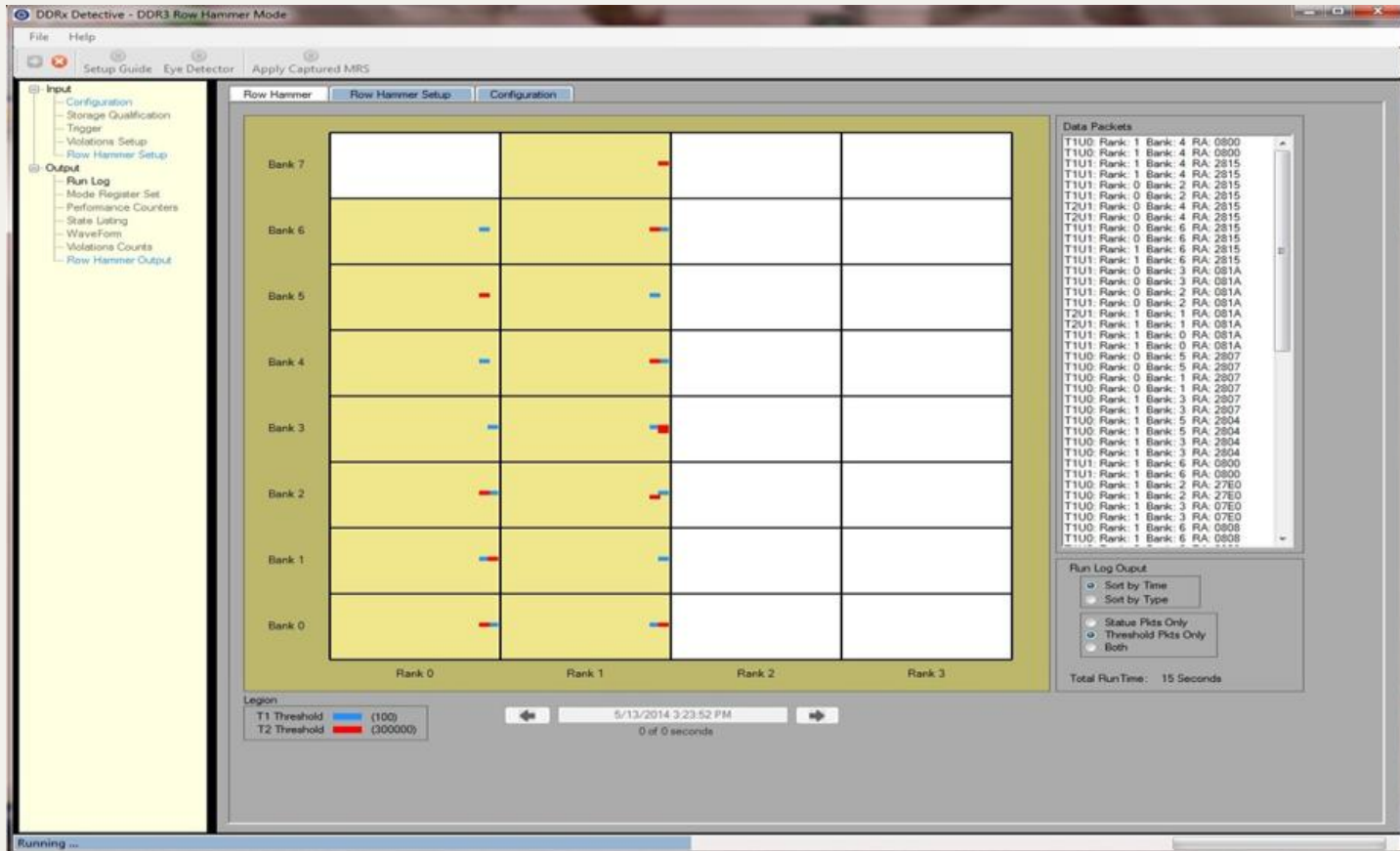
## Critical Applications cannot tolerate Row Hammer failures

- Memory Vendors are specifying how many ACTIVATES per Retention cycle to a single row the memory can tolerate
- However there is no way the memory controller or the memory can count this
- Critical Software applications should be tested to ensure that the Row Hammer Threshold is not reached

## One method to detect Row Hammer



## Detecting Excessive ACTIVATE Commands





## Critical DDR4 Performance Metrics

- Memory Controller/System Architecture
  - Can this insight lead to better designs?
  - Benchmark Servers Memory Performance
- Workload Analysis
  - Should the Memory Controller settings be based on criteria set by the workload?
  - Can compilers be made better?
  - Can critical applications be written to avoid Row Hammer?
- Do we all need a DDR5?
  - Work Smarter not Harder and understand what we have

## Summary

- Power Management, Bandwidth, Latency
- NEW Metrics:
  - Page Hit Analysis
  - Multiple Open Banks
  - Bank Group Analysis
  - Bank Utilization
  - Boot Analysis
  - Row Hammer Detection
  - New Measurements give insight into new designs and better architectures

# Contact Information

Barbara Aichinger

Vice President New Business Development

FuturePlus Systems

[Barb.Aichinger@FuturePlus.com](mailto:Barb.Aichinger@FuturePlus.com)

[www.FuturePlus.com](http://www.FuturePlus.com)

Check out our new website dedicated to

DDR Memory! [www.DDRDetective.com](http://www.DDRDetective.com)