

# Active/Idle Toggling with Low-Power Idle

January 2008  
IEEE 802.3az Task Force

Presenter: Robert Hays  
Intel Corporation

Contributors: Aviad Wertheimer, Eric Mann



IEEE 802.3az January 2008 Interim Meeting



# Supporters

- Ozdal Barkan (Marvell)
- Jim Barnette (Vitesse)
- Brad Booth (AMCC)
- Joseph Chou (Realtek)
- Dan Dove (HP ProCurve)
- Robert Hays (Intel)
- Adam Healey (LSI)
- Sanjay Kasturia (Teranetics)
- David Koenen (HP)
- David Law (3Com)
- Brian Murray (LSI)
- Gavin Parnaby (Solarflare)
- Wiren Perera (Plato Networks)
- Aviad Wertheimer (Intel)
- Bill Woodruff (Aquantia)
- George Zimmerman (Solarflare)

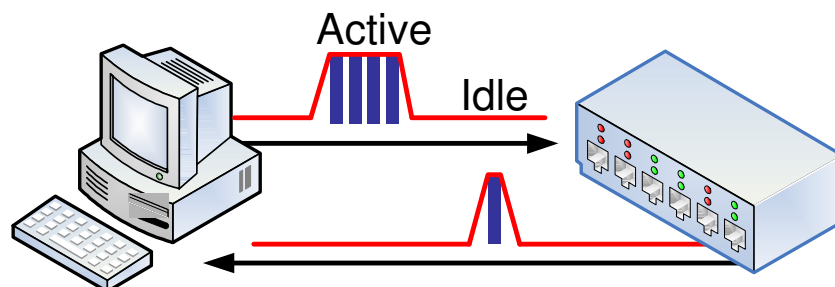
# Agenda

1. Updates from November (see `hays_01_1107`)
  - Glossary
  - Active/Idle Toggling Concept
  - Low-Power Idle Overview
  - Power Consumption
2. Elaboration on some Elements
  - Asymmetric Operation
  - Supporting Deep Sleep Levels
  - Auto-Negotiation
  - Initiating Transitions
3. Benefits of Active/Idle Toggling
4. Areas for Further Investigation

# Glossary

- **Electrical Energy Terms:**
  - **Operating Power** - (Watts) The rate at which electrical energy is delivered to a circuit or system
  - **Energy Consumption** - (Joules) Aggregate power consumed by a system over a period of time
  - **Energy Efficiency** - (Joules/bit) Energy required to complete a unit of work. E.g. energy required to transmit/receive each bit of data.
  - **Average Power** - (Watts) Energy consumed divided by period of time
- **Ethernet Operating States:**
  - **Active** - Sending packets. Higher power. Defined today for all PHYs.
  - **Normal Idle (N\_IDLE)** - Not sending packets. Same or less power than Active. Defined today as "Idle" for all PHYs.
  - **Low-Power Idle (LP\_IDLE)** - Not sending packets. Minimal power. To be defined by IEEE 802.3az.

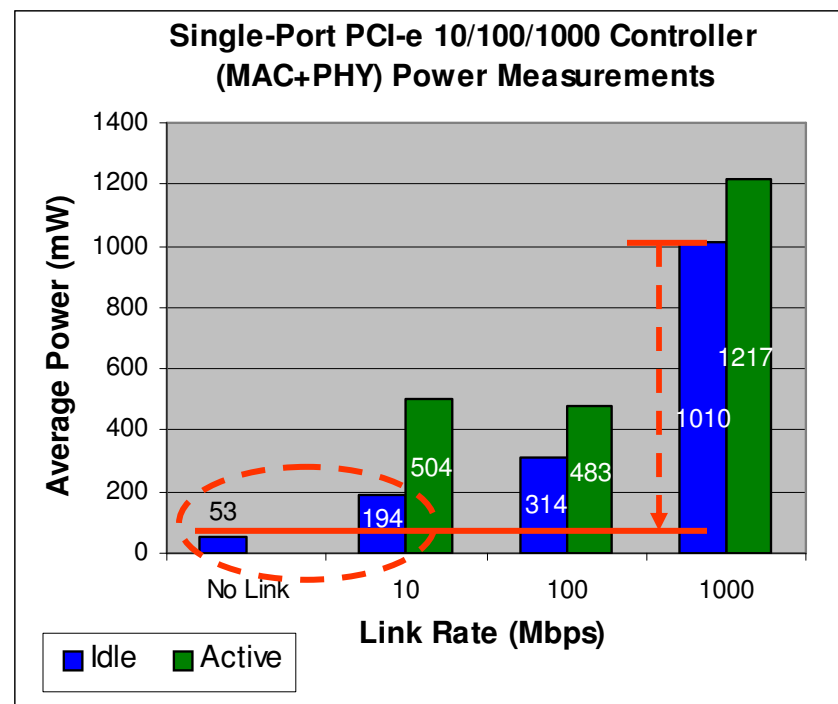
# Active/Idle Toggling Concept



- Principle: Transmit data as fast as possible, return to Low-Power Idle
  - Highest rate provides the most energy-efficient transmission (Joules/bit)
  - IP\_IDLE consumes minimal power (Watts)
- Energy savings come from cycling between Active & Low-Power Idle
  - Power is reduced by turning OFF unused circuits during LP\_IDLE (e.g. portions of PHY, MAC, interconnects, memory, CPU)
  - Energy consumption scales with bandwidth utilization
- Transmitter initiates LP\_IDLE transitions, Receiver acquiescent
  - Control policy is managed by system entity beyond IEEE 802.3 scope

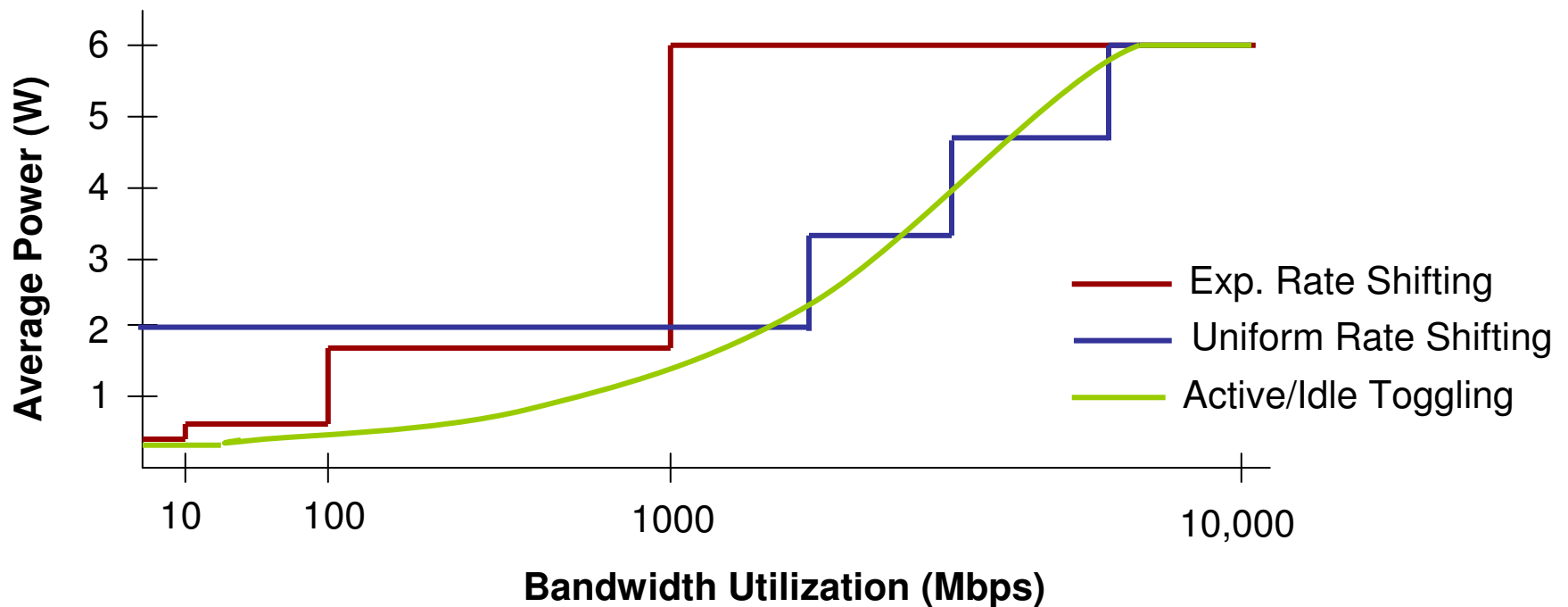
# Low-Power Idle Overview

- LP\_IDLE is a “quiet” line that consumes minimal power
  - It is used when no data is being transmitted
  - Only essential circuitry (e.g. timing recovery) must remain ON
  
- Rate-specific solutions required:
  - 100BASE-TX (see chou\_01\_0108)
  - 1000BASE-T (see healey\_01\_0108)
  - 10GBASE-T (see zimmerman\_01\_0108)
  - 10GBASE-KR (see booth\_01\_0108)
  - 10GBASE-KX4 (see booth\_01\_0108)
  
- Gigabit LP\_IDLE power estimate:
  - “No Link”  $\leq$  LP Idle  $\leq$  10Mbps Idle
  - e.g. 53mW  $\leq$  LP Idle  $\leq$  194mW
  - Should be closer to “No Link”



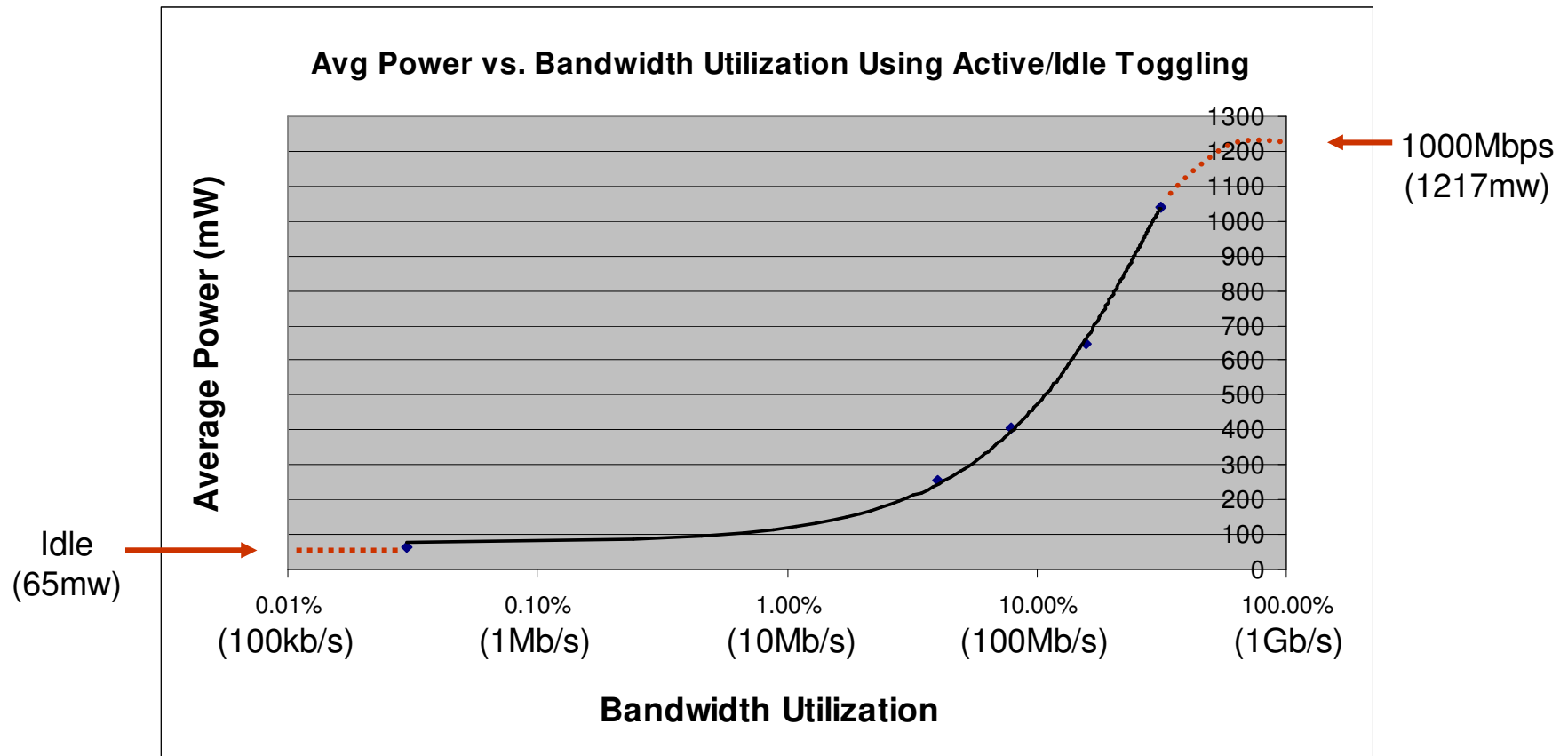
Source: Intel labs. Intel® 82573L Gigabit Ethernet Controller, 0.13μm, “Idle” = no traffic, “Active” = line-rate, bi-directional

# Conceptual Average Power vs. BW Utilization



- **Exponential Rate Shifting** offers power steps at  $1/10^{\text{th}}$ ,  $1/100^{\text{th}}$ ,  $1/1000^{\text{th}}$  rates for savings during periods of low-utilization ( $<10\%$ )
- **Uniform Rate Shifting** offers power steps on  $1/4^{\text{th}}$  rate increments for savings during periods of medium to high utilization (25%-75%)
- **Active/Idle Toggling** with Low-Power Idle allows smooth power averaging across a broad range of bandwidth utilization ( $<80\%$ ?)

# Simulated Active/Idle Toggling Avg. Power



Source: Intel labs. Simulation program source code and sample traffic pattern trace files posted on the EEE Tools web page: <http://grouper.ieee.org/groups/802/3/az/public/tools/index.html>

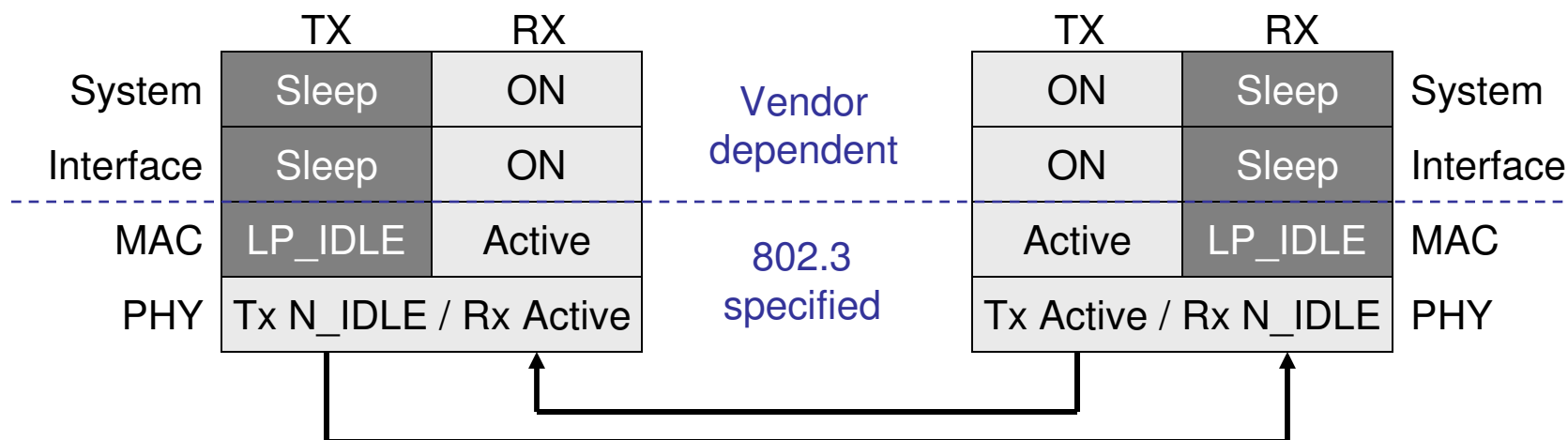
## Input Assumptions:

- Traffic Input = Trace\_VOIP\_\*.txt
- 1000Mbps Active Power = 1217mW
- LP\_IDLE Power = 65mW
- LP\_IDLE Initiation Wait = 10 $\mu$ s
- LP\_IDLE Transition Latency = 1 $\mu$ s
- Active Resume Latency = 10 $\mu$ s





# Asymmetric Operation

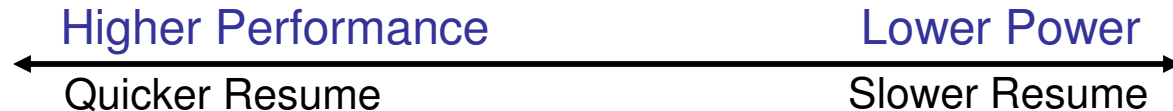


- Asymmetric operation would further improve energy efficiency
  - Independent Tx & Rx transitions into LP\_IDLE
  - End-node traffic is typically weighted toward either send or receive
- Asymmetric toggling is valuable at MAC-layer and above
  - Tx & Rx data paths already operate independently above the PHY
  - Transition initiation would need to occur between MACs
  - PHYs would only enter LP\_IDLE if both Rx & Tx are in N\_IDLE

# Supporting Deeper Sleep Levels

	Active PC Example	Quick-Resume PC Example (~10µs)	Longer-Resume PC Example (~100µs)	
Memory	M0 Active	M0 Active	M1 Standby (100µs)	Vendor dependent
PCIe	L0 Active	L0s Standby (3µs)	L1 Standby (6µs)	
MAC	Active	LP_IDLE (1µs)	LP_IDLE (1µs)	802.3 specified
PHY	Active	LP_IDLE (10 µs)	LP_IDLE (10µs)	

- Variable resume latencies allow performance vs. power optimization



- Resume predictability allows more intelligent power management
  - Greater power savings doesn't come from just longer LP\_IDLE duration, it comes from being able to safely turn OFF/ON more circuitry
  - Two ways to provide predictability:
    - Rx tells Tx how long to wait before sending data (via negotiated resume latency)
    - Tx tells Rx how long it will be in LP\_IDLE (via notification of sleep duration)

# Auto-Negotiation

- Negotiate EEE capabilities during Auto-negotiation:
  1. EEE support for each speed
    - a. 10G
    - b. 1G full-duplex
    - c. 100M full-duplex
  2. LP\_IDLE Resume Latency values
    - a. Maximum T\_RESUME (may be specified by 802.3az)
    - b. Minimum T\_RESUME (may be specified by 802.3az)
    - c. Desired T\_RESUME
  3. Possibly... LP\_IDLE Duration parameters:
    - a. Maximum T\_LP\_IDLE (PHY or system limitation)
    - b. Minimum T\_LP\_IDLE (for effective power saving)
- Updates (e.g. T\_RESUME changes) could be negotiated via MAC control frames or other means

# Initiating Transitions

- Transition control policy is managed by a system entity beyond IEEE 802.3 scope
- Transition initiated by Tx (data source), Rx acquiescent
  - 2-way negotiation or Acks are unnecessary
- Example transition to/from LP\_IDLE:
  1. When no data to transmit, Tx signals entry into LP\_IDLE
  2. Rx detects entry into LP\_IDLE and may reduce it's power
  3. PHYs may periodically wake for Link Training
    - Training may only be necessary for some PHYs, e.g. 10GBASE-T
  4. When data to transmit, Tx PHY enters N\_IDLE and MAC waits negotiated T\_RESUME before beginning data transmission

## Benefits of Active/Idle Toggling for EEE

- Reduced power during low utilization
- Energy consumption scales with bandwidth utilization
- Minimal impact to performance
- Turning circuits ON/OFF is easier than rate shifting
- Integrates well with PC & server power management
- Simple, one-way transition initiation
- May allow Asymmetric operation to save additional energy

## Areas for Further Investigation

- Low-Power Idle state for each PHY type
- Negotiating resume latencies and/or LP\_IDLE durations
- Transition signaling scheme
- MAC-PHY sync control
- Asymmetric operation

# Thank You!

- Questions?