# ORACLE®

# Wavelength Stealing: An Opportunistic Approach to Channel Sharing in Multi-chip Photonic Interconnects

Arslan Zulfiqar
University of Wisconsin – Madison
MICRO-46

# COLLABORATORS

- Oracle Labs
  - Pranay Koka
  - Herb Schwetman
  - Xuezhe Zheng
  - Ashok Krishnamoorthy

- University of Wisconsin – Madison
  - Mikko Lipasti

ORACLE

"Approved for Public Release. Distribution Unlimited"
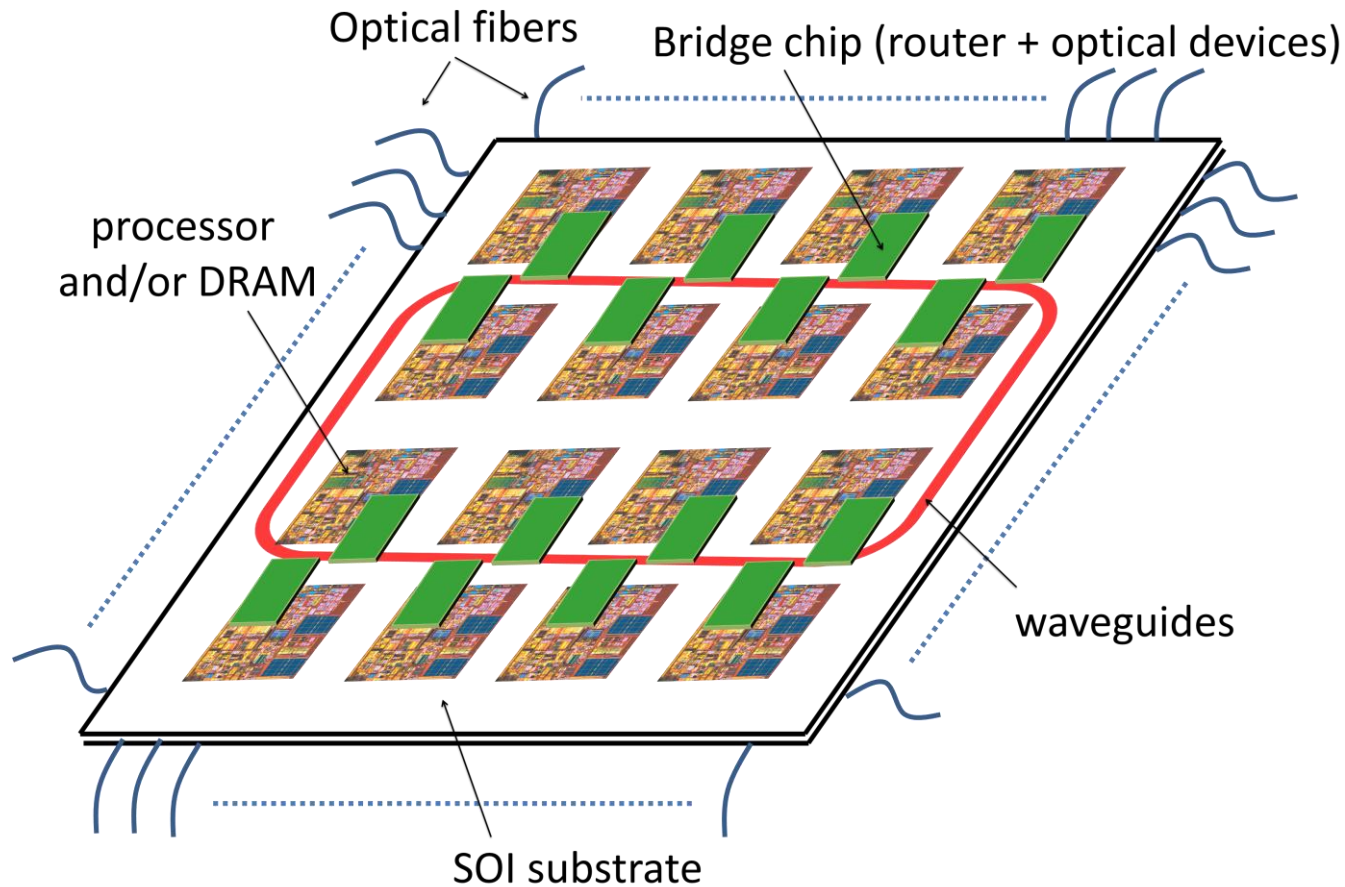
# EXECUTIVE SUMMARY

- Silicon-photonics offers integration of multiple chips
  - High sustainable performance
  - Improved process yields
  - Lower energy-per-bit consumption

- This work:
  - Focus on **channel sharing** nanophotonic designs
  - Model to determine limits and gains of channel sharing
  - Novel shared channel architecture: **Wavelength Stealing**
    - Arbitration-free
    - Strong fairness guarantees
    - Up to 28% better EDP than baseline on HPC workloads

**ORACLE**

"Approved for Public Release. Distribution Unlimited"
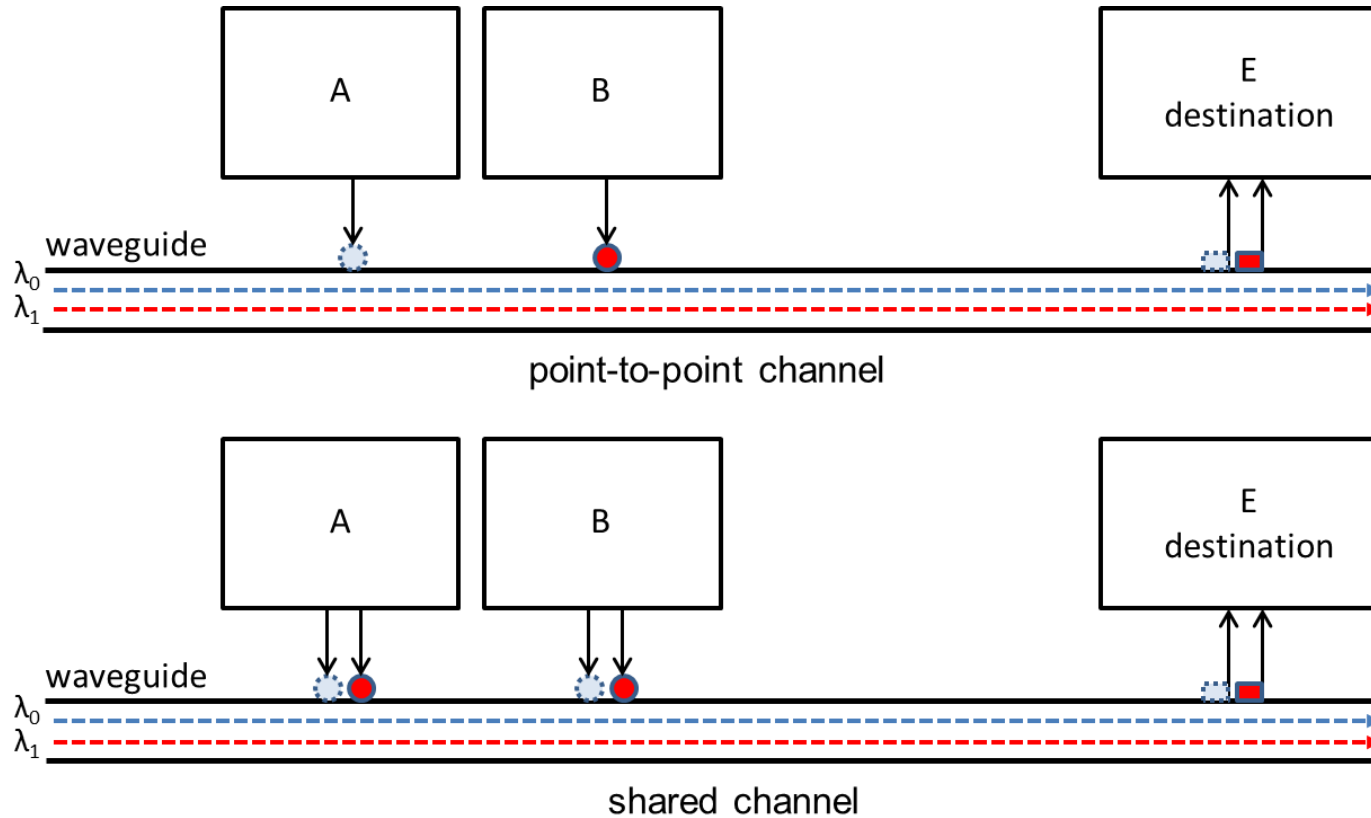
# MOTIVATION

- Technology trend: More cores

- Scaling single chip systems
  - Increasing fabrication costs
  - Low process yields
  - Power delivery limitations

- Multichip systems – require enormous off-chip communication
  - Low bandwidth densities for off-chip I/O
  - High power consumption of serial links

"Approved for Public Release. Distribution Unlimited"

# ORACLE'S "MACROCHIP" VISION

- Aggregate multiple chips with photonic communication

# WAVELENGTH SHARING LOSSES



point-to-point channel

shared channel

Extra ring-resonators on shared wavelengths increase link-loss leading to higher laser power consumption

ORACLE®

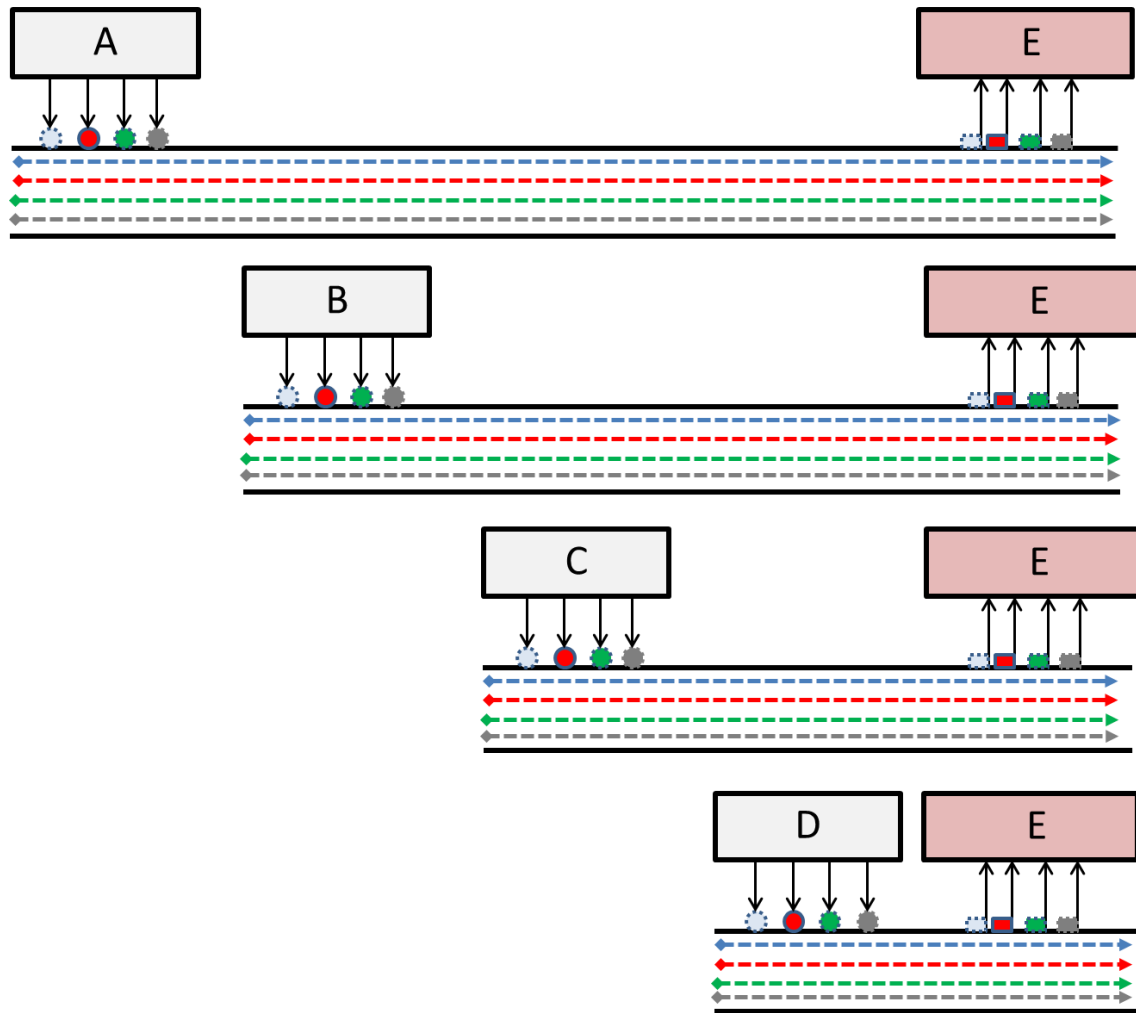"Approved for Public Release. Distribution Unlimited"

# TECHNOLOGY IMPLICATIONS

- Photonic networks are static power dominated
  - Laser power
  - Ring tuning power

- Efficiencies of commercial WDM lasers: $1 - 5\%$
  - Laser power consumption biggest contributor to static power

- Optimizing for laser power first-order design constraint

Fixed input laser-power budget for all designs

ORACLE®

## P2P (unshared)



**All-to-All traffic**
- 4 x 4b/cycle

  = 16b/cycle

**Permutation traffic**
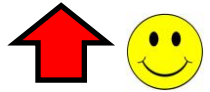- 4b/cycle

ORACLE

**2-way sharing**



All-to-All traffic
- 4 x 3b/cycle
  = 12b/cycle

Permutation traffic
- 6b/cycle

**vs**

All-to-All traffic
- 4 x 4b/cycle
  = 16b/cycle

Permutation traffic
- 4b/cycle

ORACLE®

# IMPLICATIONS OF LASER POWER BUDGET (III)

## 4-way sharing



All-to-All traffic
- 4 x 1b/cycle

  = 4b/cycle

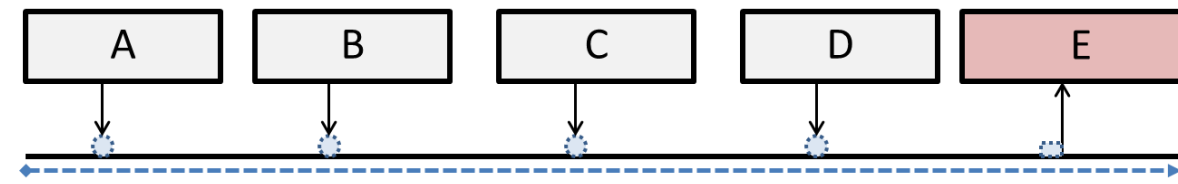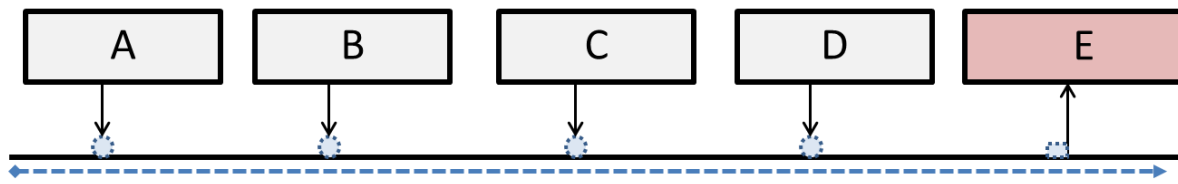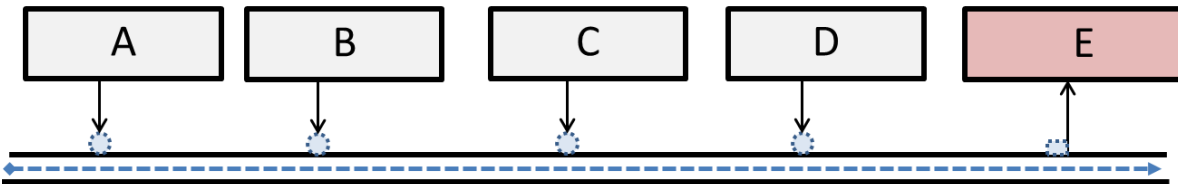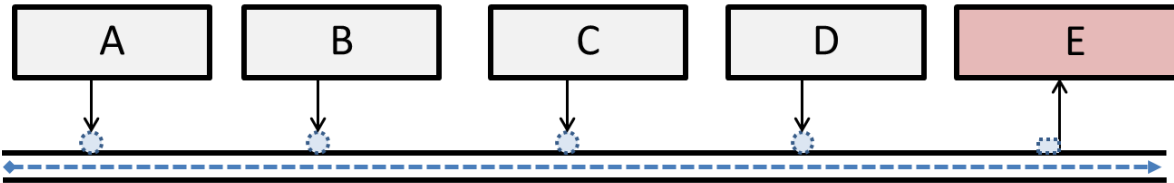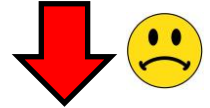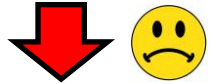Permutation traffic
- 4b/cycle

**vs**

All-to-All traffic
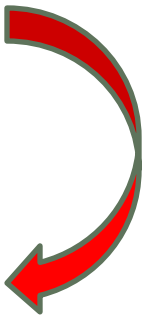- 4 x 3b/cycle

  = 12b/cycle

Permutation traffic
- 6b/cycle

# IMPLICATIONS OF LASER POWER BUDGET (IV)

- Increasing sharing degree: '$s$'
  - **Reduces effective capacity**
    - Lower performance on all-to-all (uniform random) traffic than P2P
  - **Increases peak node-node BW followed by <u>drop-off</u>**
    - <u>Potentially</u> better performance on permutation traffic than P2P

- **Can we estimate ideal sharing degree '$s_{ideal}$' and ideal node-node BW gain over P2P? <span style="color:red">Yes</span>**

"Approved for Public Release. Distribution Unlimited"

# IDEAL SHARING GAINS

Optimal sharing gain

~ 1.70x

Optimal sharing degree

$s_{ideal} = 3$



- Ignore wavelength/ time overheads of sharing
- Conservative device assumptions

# IDEAL SHARING GAINS SUMMARY

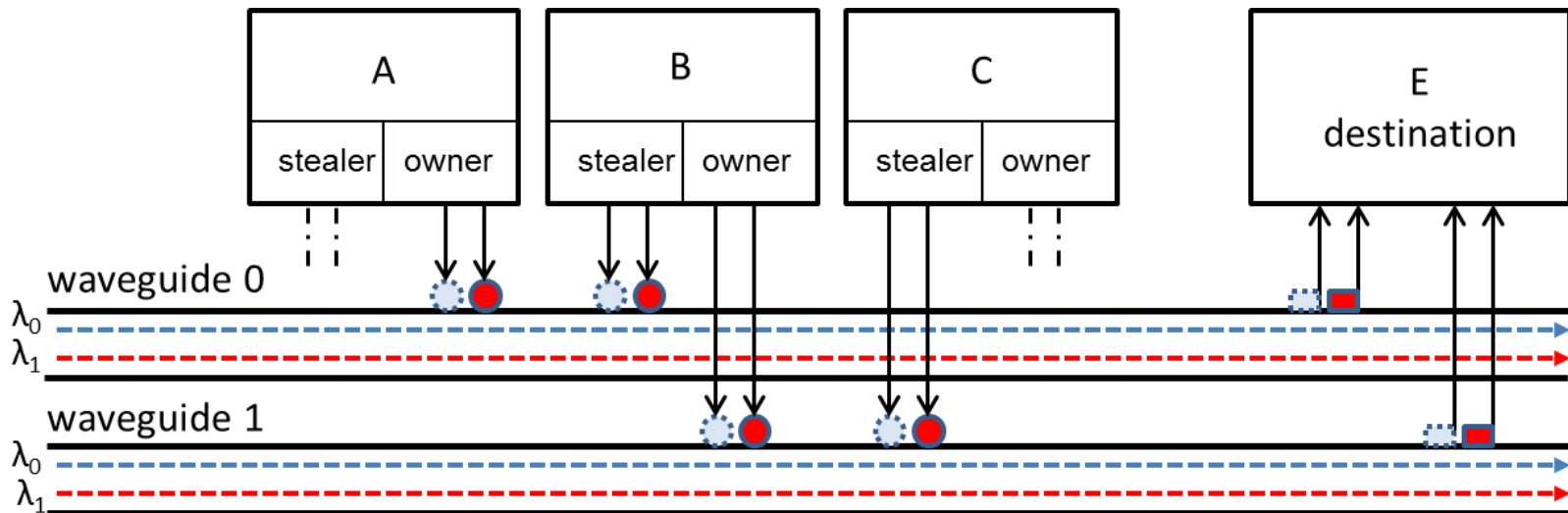| P2P | Channel Sharing |
|---|---|
| (+) High capacity | (+) High N-N BW only when $2 \leq s \leq 3$ |
|    ➢High-radix traffic |    ➢Speedup $\leq 1.70 \times$<br>   ➢Low-radix traffic |
| (-) Low N-N BW | (-) Low capacity |
|    ➢Low-radix traffic |    ➢High-radix traffic |

**ORACLE**

# WAVELENGTH STEALING

- Same topology as the P2P network: $N^2$ channels
- Every channel has **one owner** and **one or more stealers**

### 2-way stealing

"Approved for Public Release. Distribution Unlimited"

ORACLE®

# IMPLEMENTATION REQUIREMENTS

- **Owner node**
  - Guaranteed non-blocking access

- **Stealer node**
  - Arbitration-free access on an owner's channel: <u>possible packet corruption</u>
  - Notification to halt stealing when channel busy

- **Destination node**
  - Valid phit: identify sender (owner or stealer?)
  - Corrupted phit: perform correction
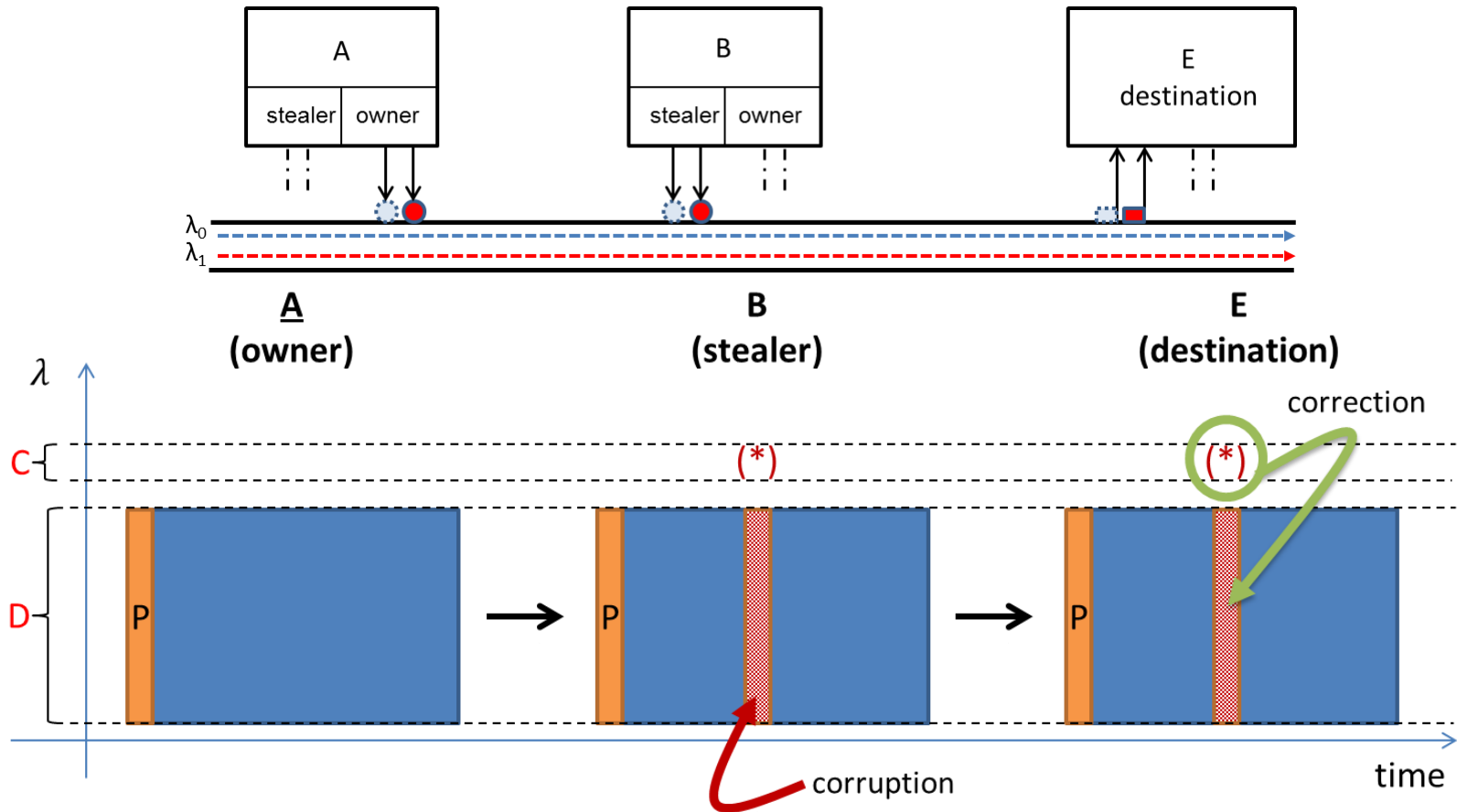
# IMPLEMENTATION MECHANISM

ERASURE CODING

**+**

CONTROL WAVELENGTHS PER CHANNEL

"Approved for Public Release. Distribution Unlimited"

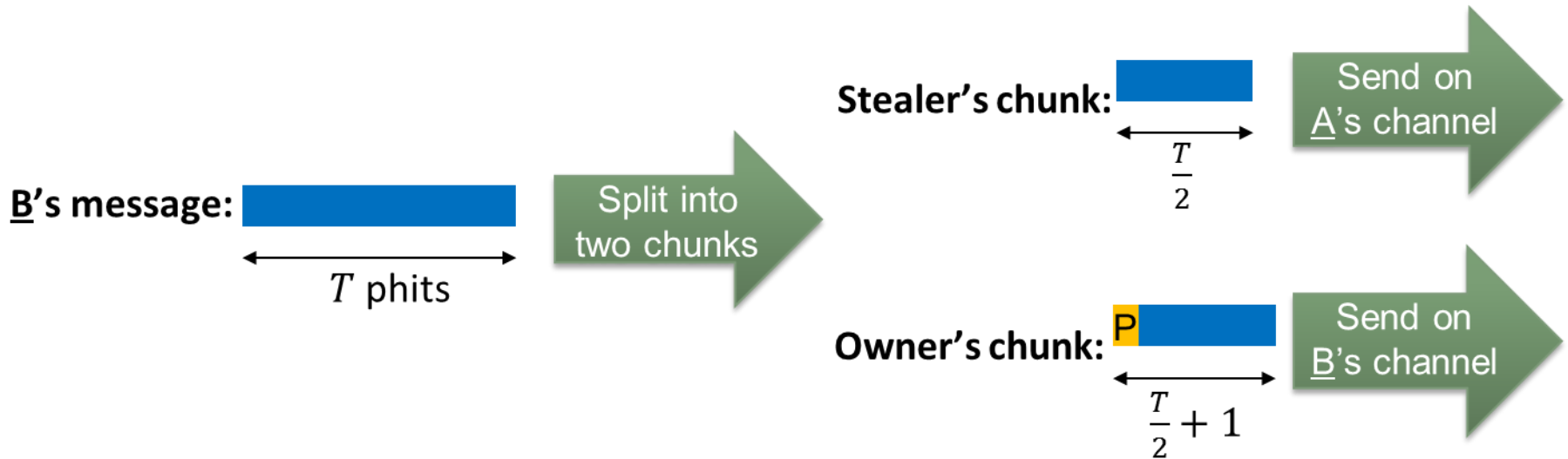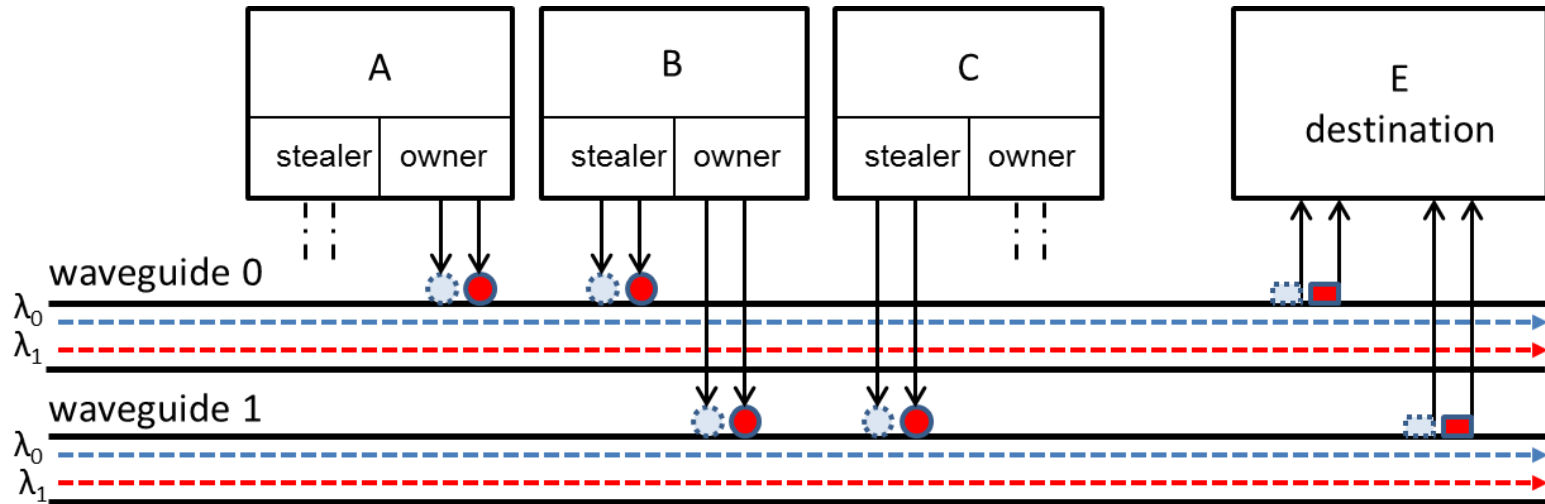ORACLE®

# ERASURE CODING

- Erasure coding is used at the destination to correct corruptions due to a collision
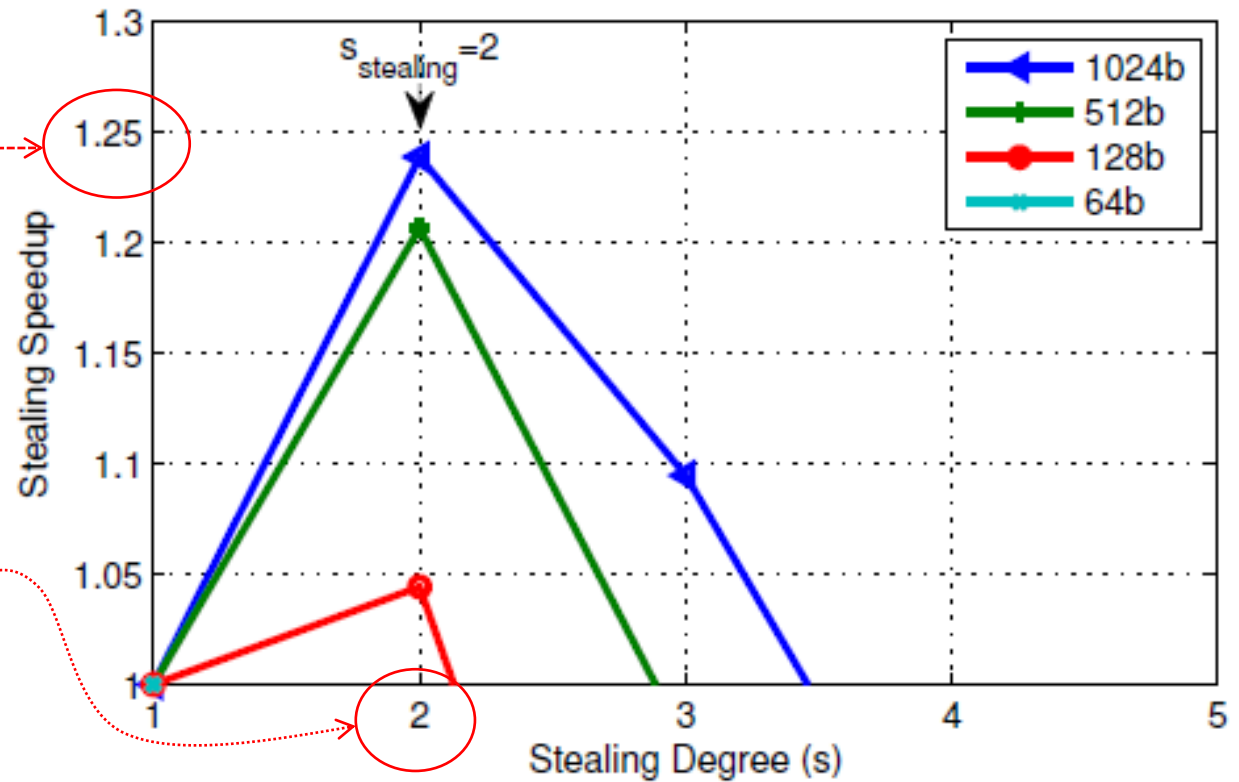
# CONTROL WAVELENGTHS

- **Functionality**
  - Mark location of corruption for erasure correction
  - Inform stealer to halt stealing when owner becomes active
  - Inform destination of the ID (owner, stealer, corrupted) of the received communication (phit)

- **Two designs – different trade-offs**
  - **Abort**
  - **Sense**

ORACLE®

# PROTOCOL OPERATION

# WAVELENGTH STEALING GAINS



Optimal stealing gain
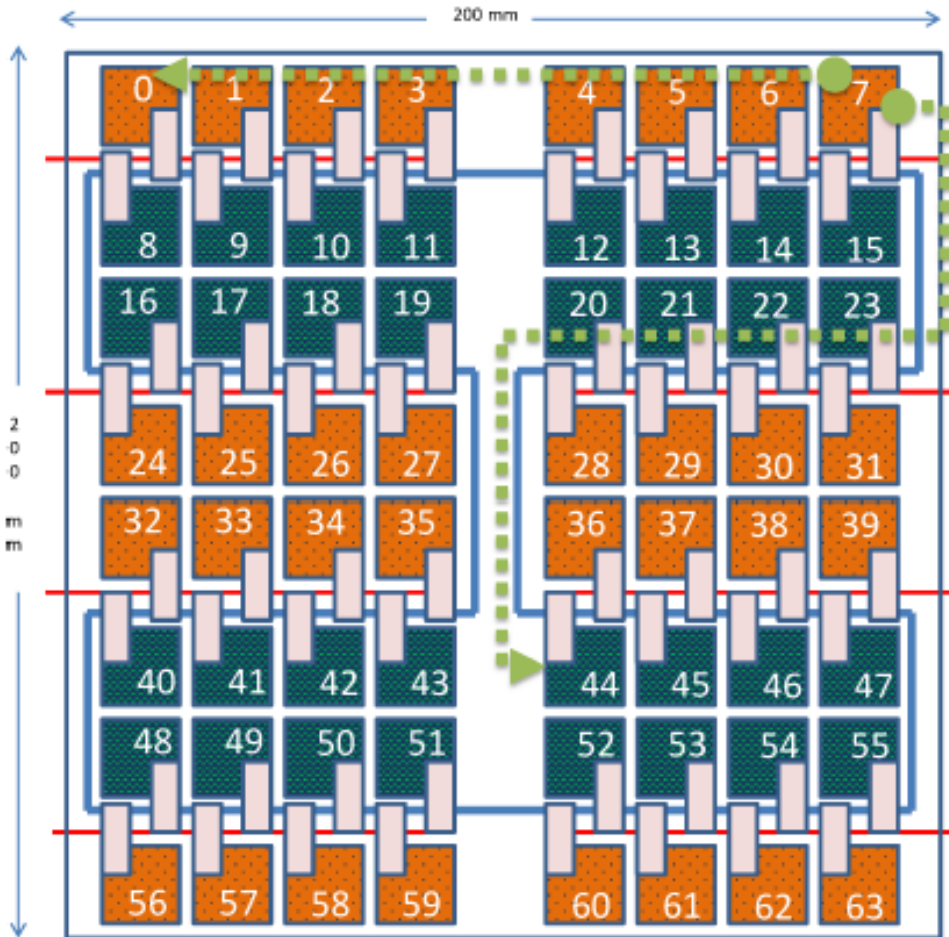
~ 1.25x

Optimal stealing degree

$s_{stealing}$ = 2

- Loss in performance due to
  - Control (wavelength) overheads $\propto \dfrac{1}{data\ wavelengths}$
  - Coding (time) overheads $\propto \dfrac{1}{message\ size}$

**ORACLE**

"Approved for Public Release. Distribution Unlimited"
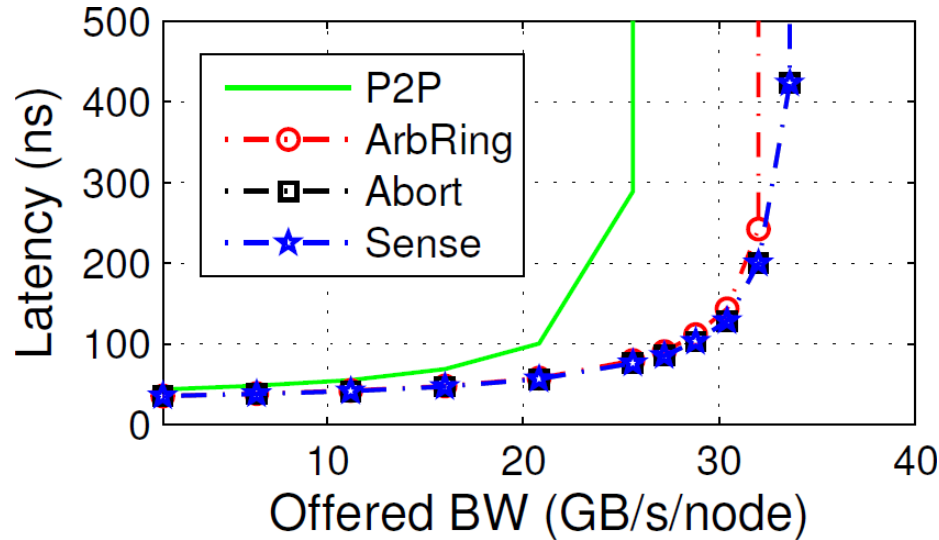
# EVALUATION – SETUP

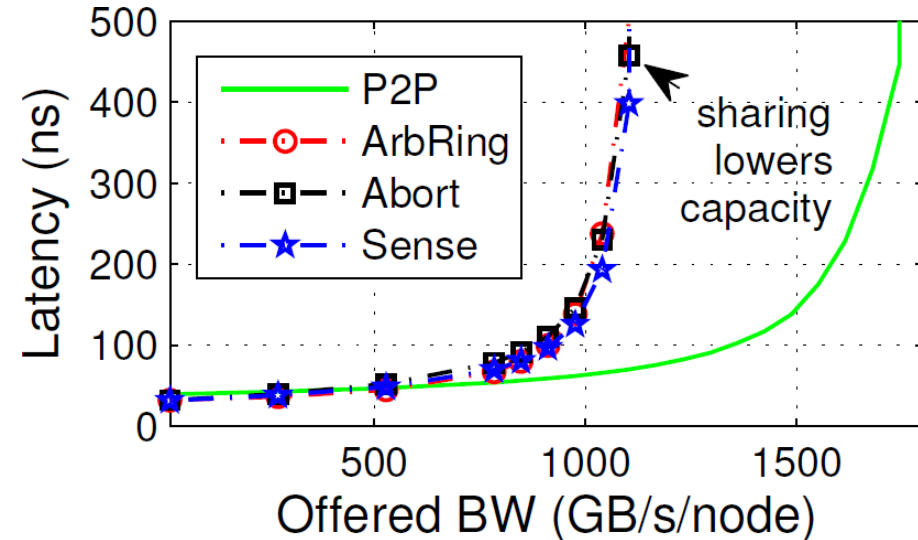**8 x 8 Macrochip System**



- Synthetic workloads
  - Uniform random
  - Permutation
  - Asymmetric

- Application workloads: NAS
  - BT: Block tri-diagonal solver
  - CG: Conjugate gradient kernel
  - DT WH: "White Hole" graph analysis
  - DT BH: "Black Hole" graph analysis
  - DT SH: "Shuffle" graph analysis

# EVALUATION – SYNTHETIC WORKLOADS (I)

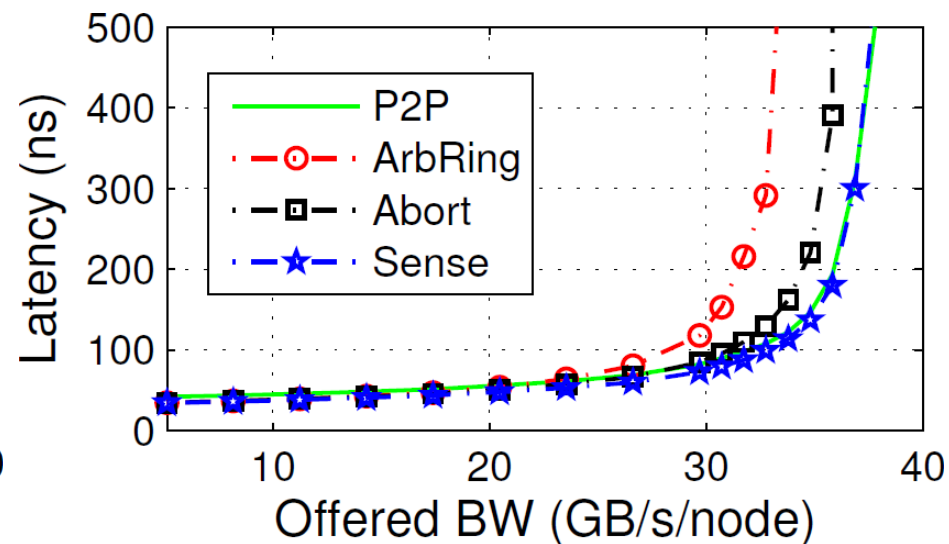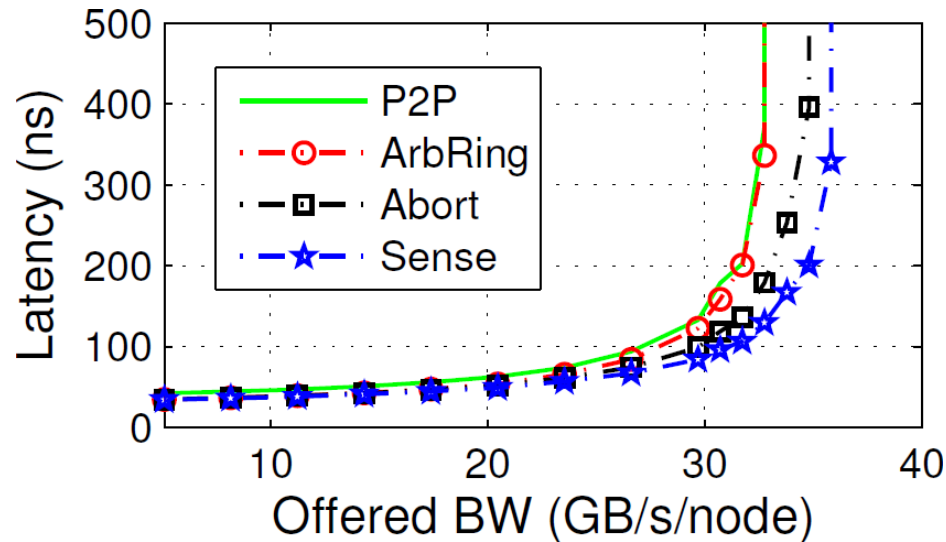**Bit-Complement (No Contention)**

**Unif. Random (Unif. Contention)**



- Sharing designs provide higher (lower) throughput than P2P in the absence (presence) of contention
- Sharing designs exhibit lower capacity

ORACLE®

# EVALUATION – SYNTHETIC WORKLOADS (II)
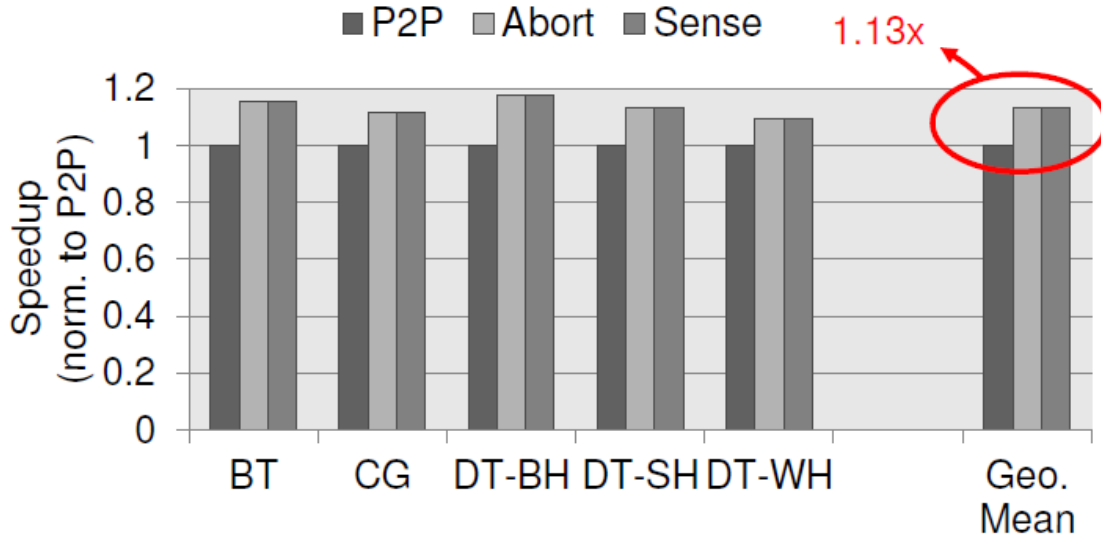
## Asymmetric K (Variable Contention)

### K = 20



### K = 30



- As contention is increased
  - Wavelength stealing provides better throughput than Token-ring design
  - Sense design outperforms abort design
  - Throughput performance of P2P improves

"Approved for Public Release. Distribution Unlimited"

- All workloads exhibit speedups
  - Max: **1.17x**
  - Average: **1.13x**

- Differences in speedups due to
  - Traffic patterns
  - Message sizes
  - Message frequencies

**ORACLE®**

"Approved for Public Release. Distribution Unlimited"

# EVALUATION – APPLICATION WORKLOADS (II)



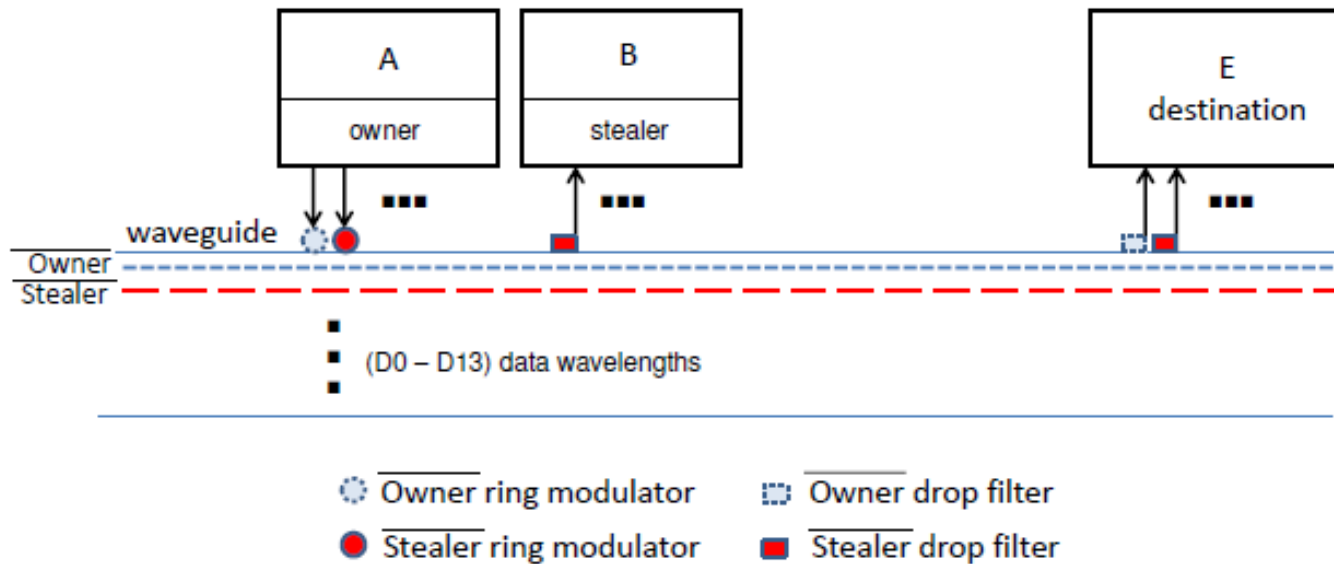- Wavelength stealing architectures achieve up to 28% lower EDP than the P2P network

- Average EDP improvement: 20% for Abort, 23% for Sense
  - Sense uses fewer ring-resonators

# CONCLUSION

- Channel sharing improves peak node-node BW compared to P2P but at the cost of reduced capacity

- Developed an analytical model to quantify limits and gains of channel sharing
  - sharing degree $\leq 3$
  - sharing gain $\leq 1.70x$

- Wavelength Stealing architecture
  - Arbitration-free accesses
  - Strong fairness guarantees
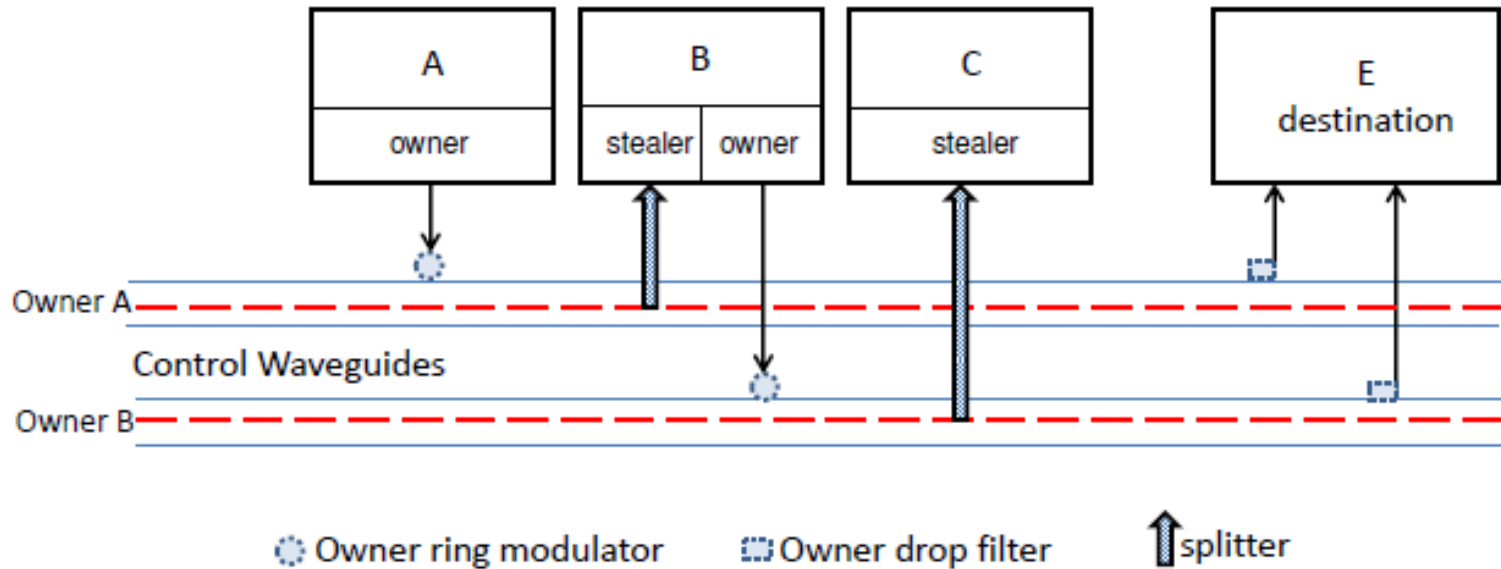  - **Guaranteed gains on VMs**

# Backup Slides

"Approved for Public Release. Distribution Unlimited"

ORACLE®

# ABORT DESIGN



| Active | A | | B | E | | Received |
| --- | --- | --- | --- | --- | --- | --- |
| Sender | $\overline{Own.}$ | $\overline{St.}$ | $\overline{St.}$ | $\overline{Own.}$ | $\overline{St.}$ | |
| $A$ | 0 | 1 | – | 0 | 1 | $A$ |
| $B$ | 1 | 0 | 0 | 1 | 0 | $B$ |
| $A, B$ | 0 | 1 | 1 | 0 | 0 | $Collision$ |
| (Invalid) | 1 | 1 | – | 1 | 1 | (Invalid) |

"Approved for Public Release. Distribution Unlimited"

# SENSE DESIGN
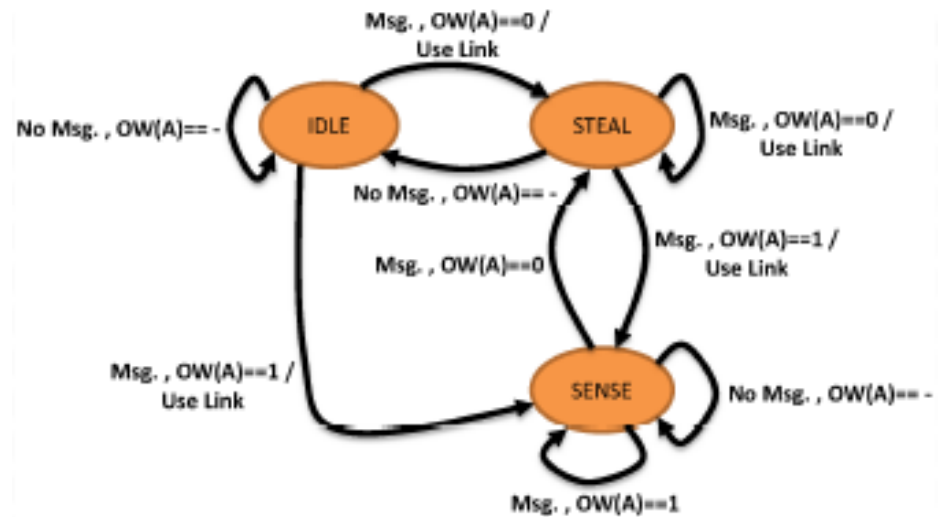


- Employs broadband splitters (non-destructive reads)
- Implemented using state machines at
  - Owner
  - Stealer
  - Destination

"Approved for Public Release. Distribution Unlimited"

# SENSE DESIGN FUNCTIONALITY (I)



Owner (A) State Machine

Stealer (B) State Machine

"Approved for Public Release. Distribution Unlimited"

**ORACLE**

# SENSE DESIGN FUNCTIONALITY (II)



Destination (E) State Machine

"Approved for Public Release. Distribution Unlimited"

ORACLE®

# ABORT VS. SENSE TRADE-OFFS

- Abort
    - (+) Fewer waveguides
    - (−) Conservative performance
    - (−) More ring-resonators
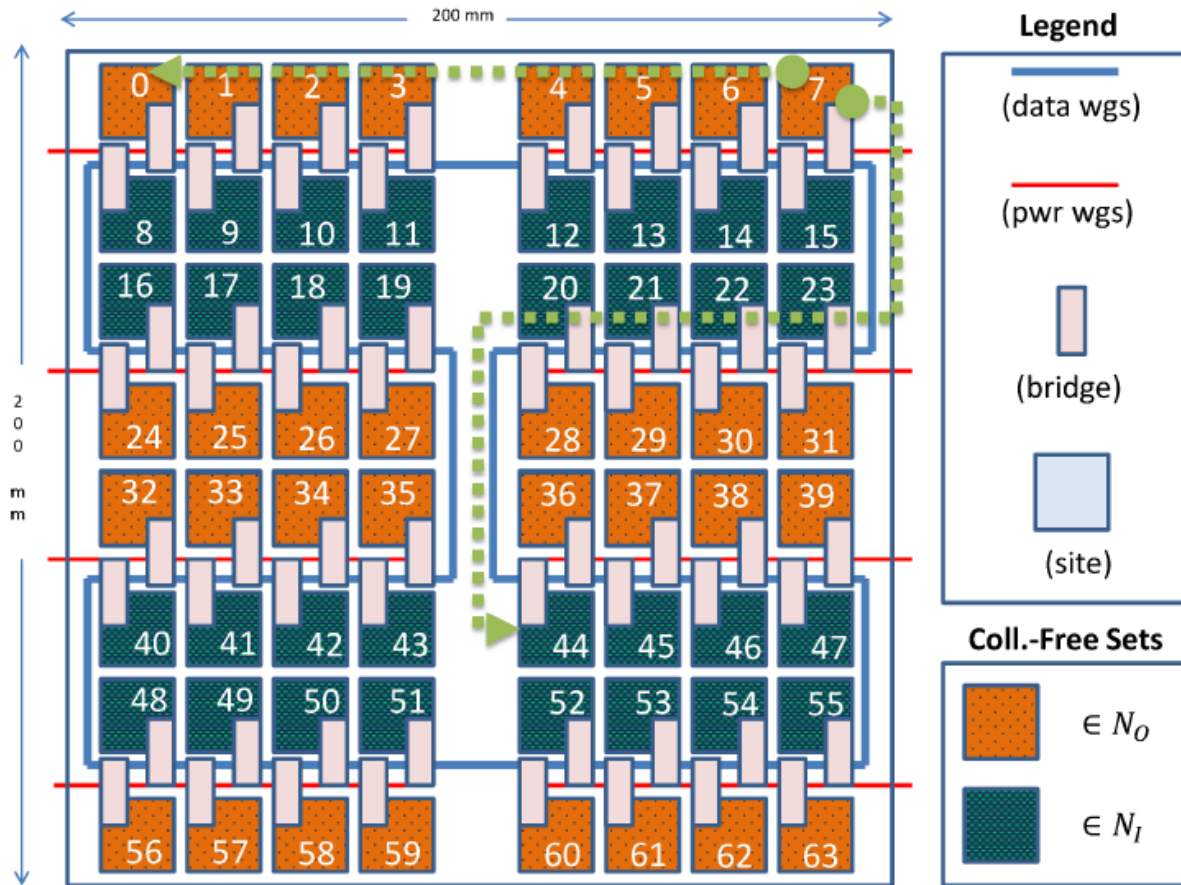
- Sense
    - (+) Aggressive performance
    - (+) Fewer ring-resonators
    - (−) More waveguides

"Approved for Public Release. Distribution Unlimited"

# OPTICAL DEVICE PARAMETERS

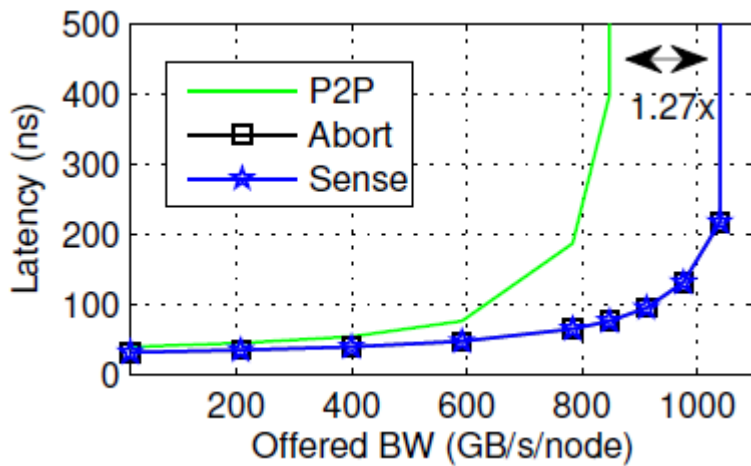| Parameter | Assumption |
|---|---|
| Mod. (Insertion) Ring Loss | $4dB$ |
| Inactive Mod. Ring Loss | $0.5dB$ |
| Active Drop-Filter Ring Loss | $1dB$ |
| Passive Ring Loss | $0.05dB$ |
| Waveguide Loss | $0.05dB/cm$ |
| Bridge Chip Waveguide Loss | $1dB$ |
| Coupler Loss | $2dB$ |
| Receiver Sensitivity Margin | $4dB$ |
| Receiver Sensitivity Level | $-21dBm$ |
| Ring Tuning Power | $0.3mW/ring$ |
| Mod. Driver | $35fJ/bit$ |
| Detector Driver | $65fJ/bit$ |
| Max. Fiber WDM-Factor | 32 |
| Max. Waveguide WDM-Factor | 16 |
| Max. Port Fibers | 2500 |
| Power per Fiber | $32mW$ |

ORACLE®

# VIRTUALIZATION GAINS (I)

- Virtualization: many VMs share the system
  - Better utilization of system resources

# VIRTUALIZATION GAINS (II)

## Domain Uniform Random



## Four 16-Node VMs

"Approved for Public Release. Distribution Unlimited"

ORACLE