# Challenges and Opportunities for the World's Largest Integrated Circuits

Mike Hutton
Office of the CTO
Altera Corporation

**Acknowledgements:**
Vaughn Betz, Stephen Brown, David Mendel, Argy Krikelis,
Mario Khalaf, Deshanand Singh, Altera Applications

VII Southern Programmable Logic Conference
Córdoba, Argentina
April 14, 2011

# Overview

- Applications / Market Drivers

- Architecture and Circuits at 28nm

- New Systems and Methodologies for 2011

Note: This version has had several of slides removed to allow for publication on the SPL conference website

# Industries that use FPGAs

**Consumer**
**Automotive**

**Test, Measurement & Medical**

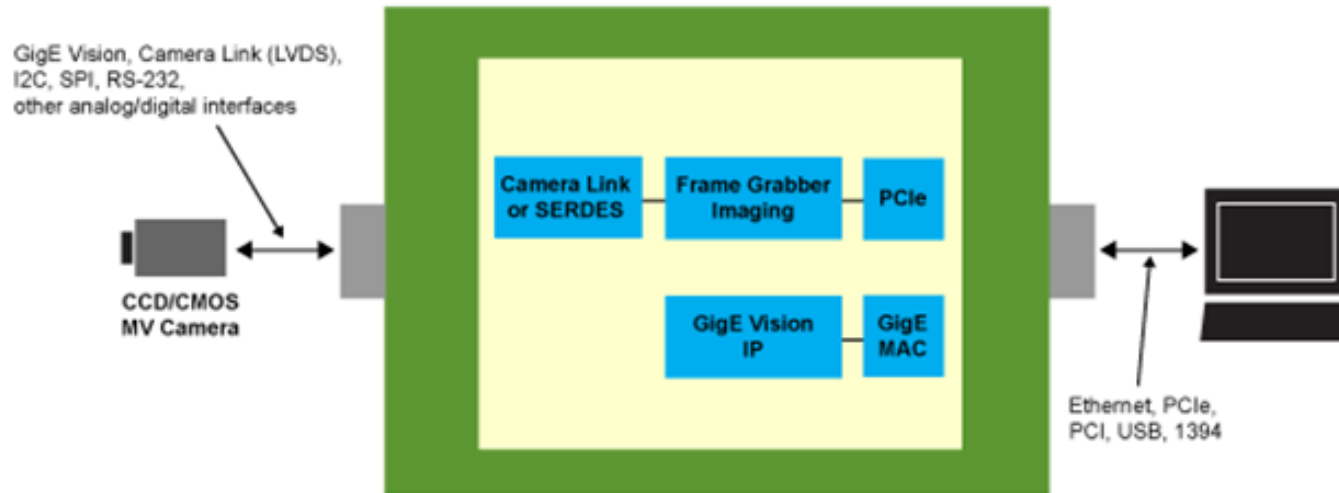**Communications Broadcast**
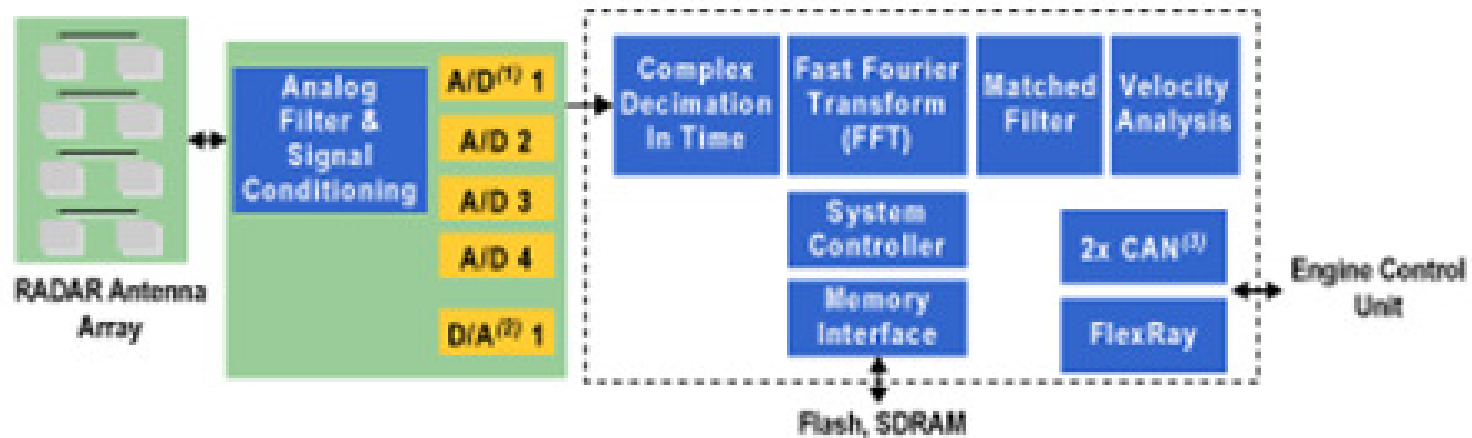
**Military & Industrial**

**Computer & Storage**

More than 50% of revenue

# E.g. Industrial Machine Vision



GigE Vision, Camera Link (LVDS), I2C, SPI, RS-232, other analog/digital interfaces

CCD/CMOS MV Camera

Camera Link or SERDES

Frame Grabber Imaging

PCIe

GigE Vision IP
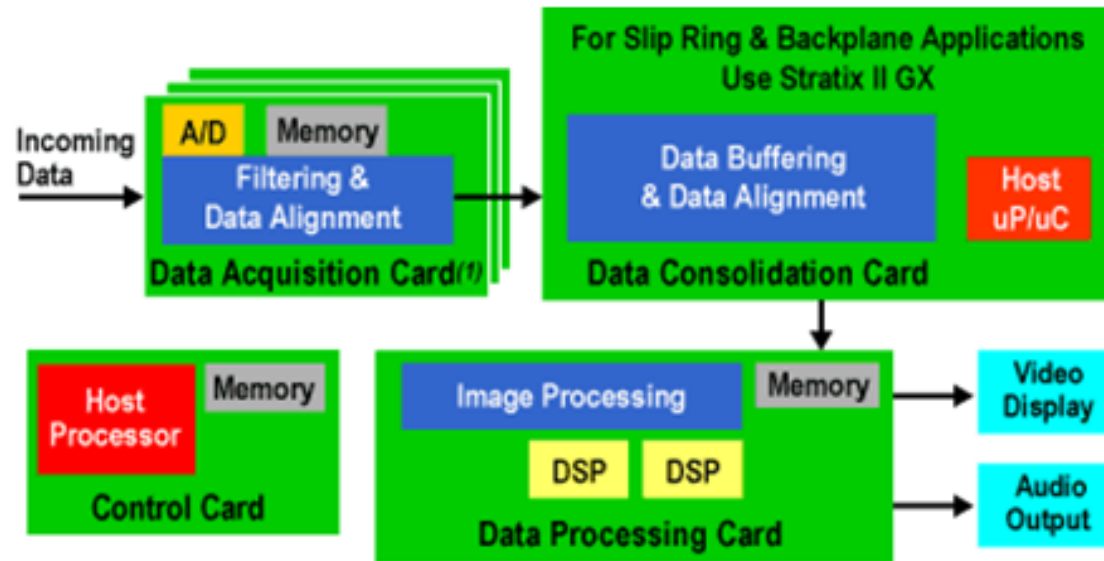
GigE MAC

Ethernet, PCIe, PCI, USB, 1394

Needs:  IP, 9x9 and 18x18 DSP, I/O interfaces

# E.g. Automotive radar / cruise control



Needs: Embedded processor, DSP, IP, I/O interfaces
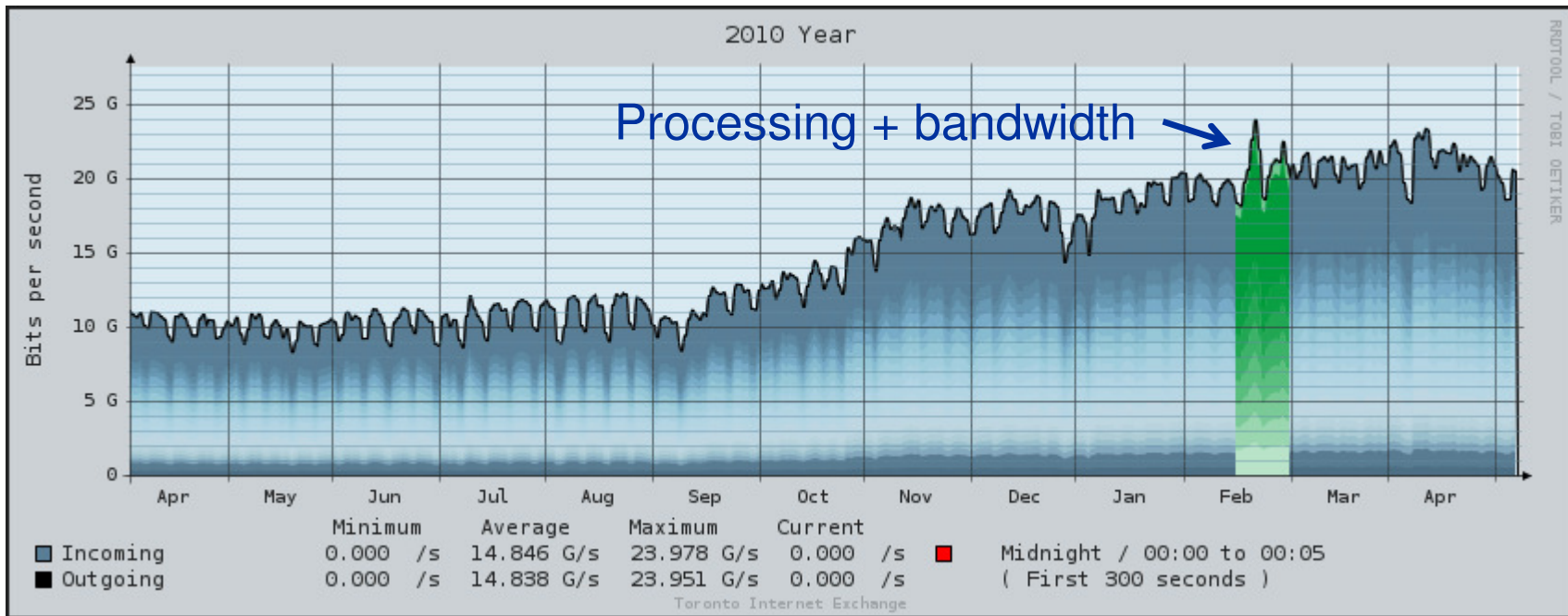
# E.g. High-end medical imaging



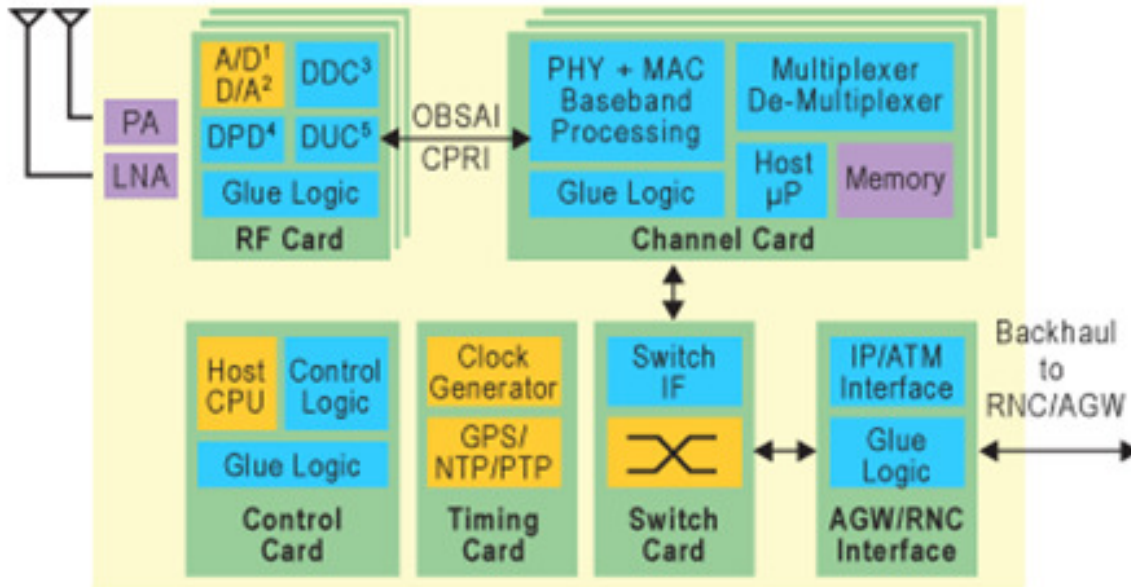Needs:    *DSP*:  18x18, 27x27, **complex**, and **FP/DP**

# Communication Processing

- More bandwidth: 25% to 131% / year across domains [Cisco]

- More data through fixed channel + more processing per packet

- Security and quality of service needs → deep packet inspection

**E.g. Toronto Internet Exchange BPS, 2009-2010**



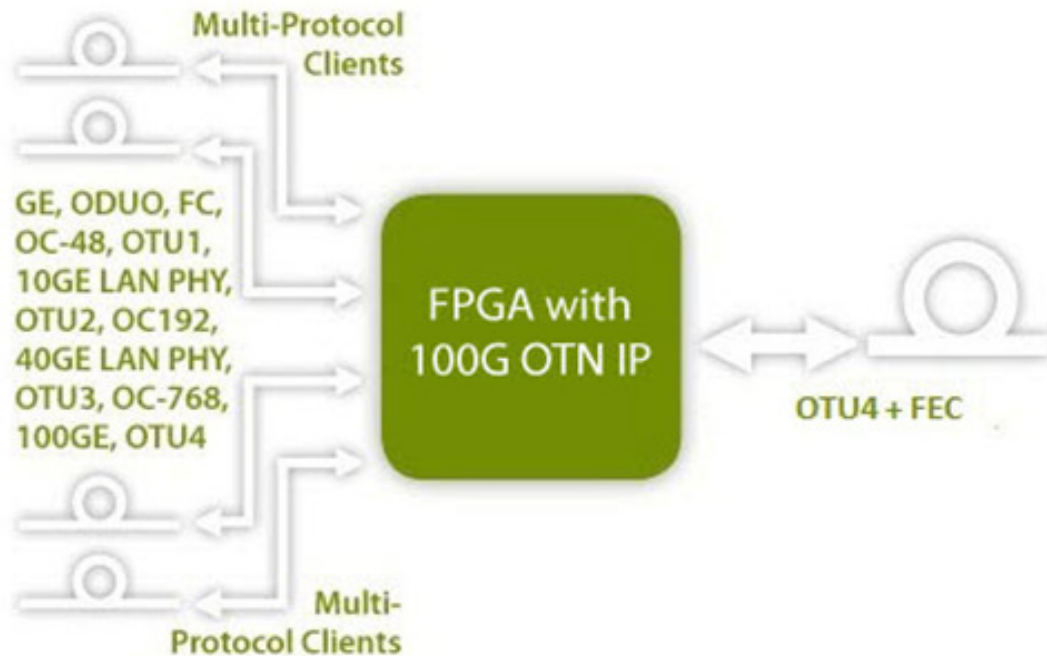**[Courtesy W. Gross, McGill]**
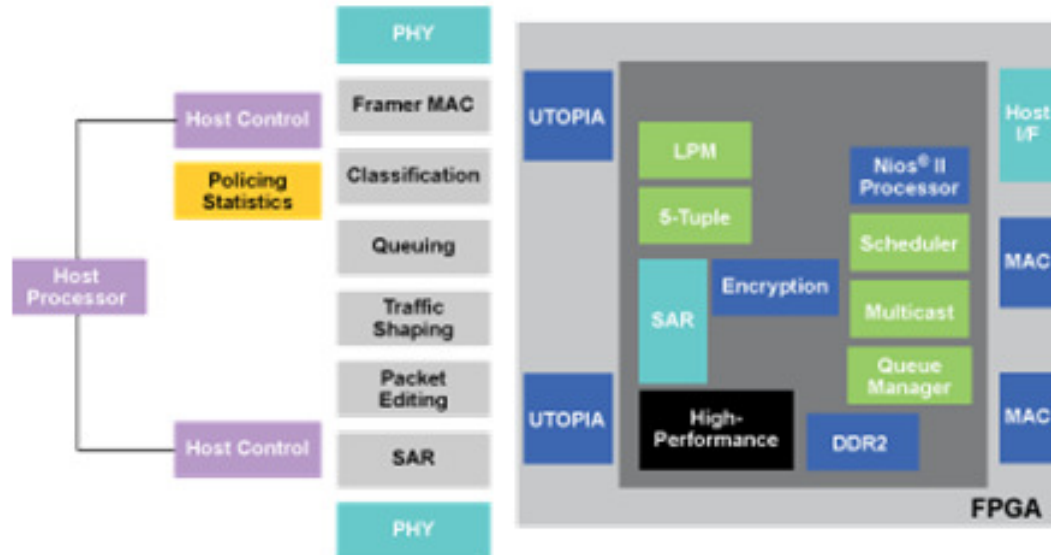
# LTE Base-station



Needs:  Integer DSP (and lots of it), IP,
Designer Productivity, Lower Power

# Wireline - Optical Transport Muxponder



Needs: I/O, memory and fabric performance, IP,
**_fast protocol changes_**, **_fractional PLLs_**,
flexible clocking

# Wireline Packet Processing



Needs: I/O, memory and fabric **performance**,
**hard protocols** (ILKN, Eth)

# Challenges

- ## Utilize the latest process generation for:

  - Greater processing power (logic, memory, I/O)
  - But, without increasing cost, power consumption or pins

- ## Requires:

  - Targeted performance improvement
  - Adding faster serial I/O – double-step from 10G to 28G
  - Disparate markets – e.g. DSP range, power vs performance
  - Improve design methodologies and CAD

# 28nm Stratix V

FPGA Fabric

- 1M LE
- 60 Mb RAM (20 Tbps)
- 4K VP DSP

HSSI

- 28 Gbps PMA
- Hard PCS
- Hard PCIe / GbE

Clocking

- High-precision PLL
- 417 clock domains

7x72b Hardened 800 MHz DDR (100 Gbps)

ALTERA®

# Families and Derivatives

**Stratix Series**
E, GS, GX, GT

Highest bandwidth and density

**Arria Series**

Balanced cost, power and performance

**Cyclone Series**

Low-power and cost

**MAX Series**

CPLD – instant-on, near-0 static power

**HardCopy Series**

Low-cost ASIC migration

ALTERA

# Embedded HardCopy Blocks
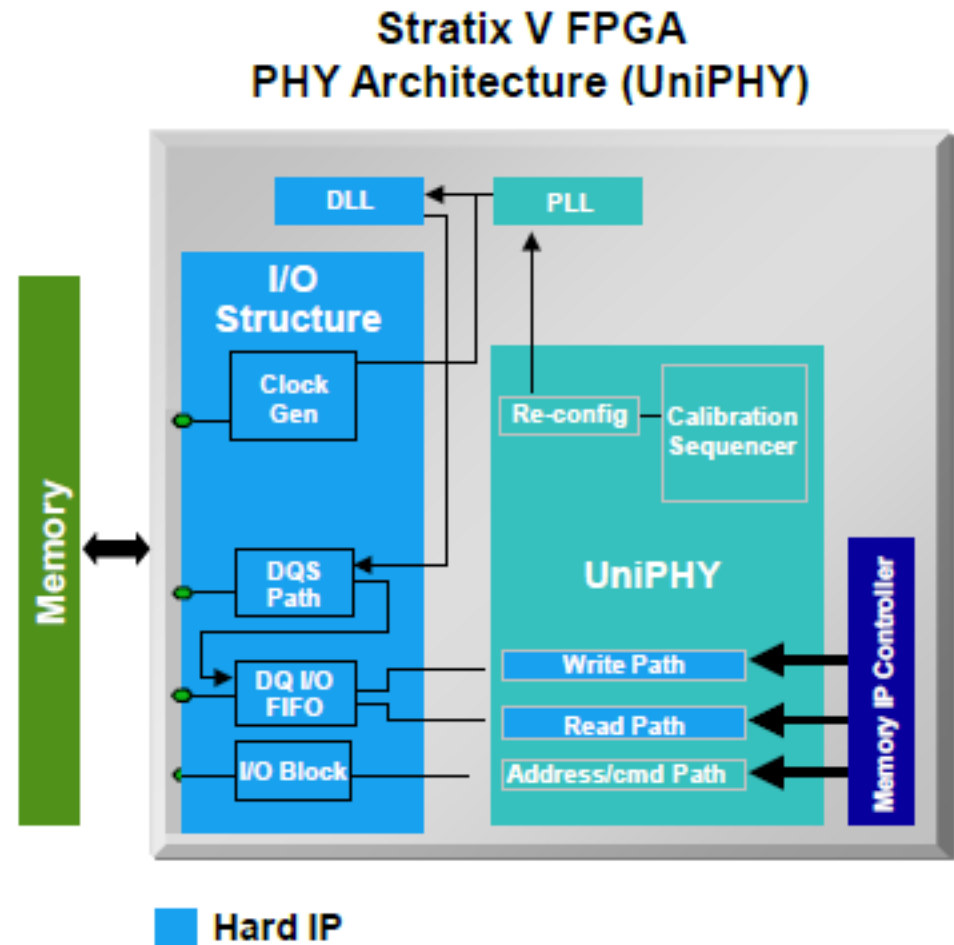
Embedded HardCopy Block

Embedded HardCopy Block

PCIe Gen3

40G/100G Ethernet

Other/Custom

- Metal-programmed ASIC blocks
  - HardCopy methodology
  - Inexpensive derivative devices
- Replaces 700K equivalent LEs
  - 5X area reduction
  - 65% AC power reduction
  - Low DC power when unused

ALTERA®

# Stratix V Memory Interfaces (DDR II/III)

- High-performance / matched delay PHY made **hard**

- PHY calibration stays **soft** for parameterization:

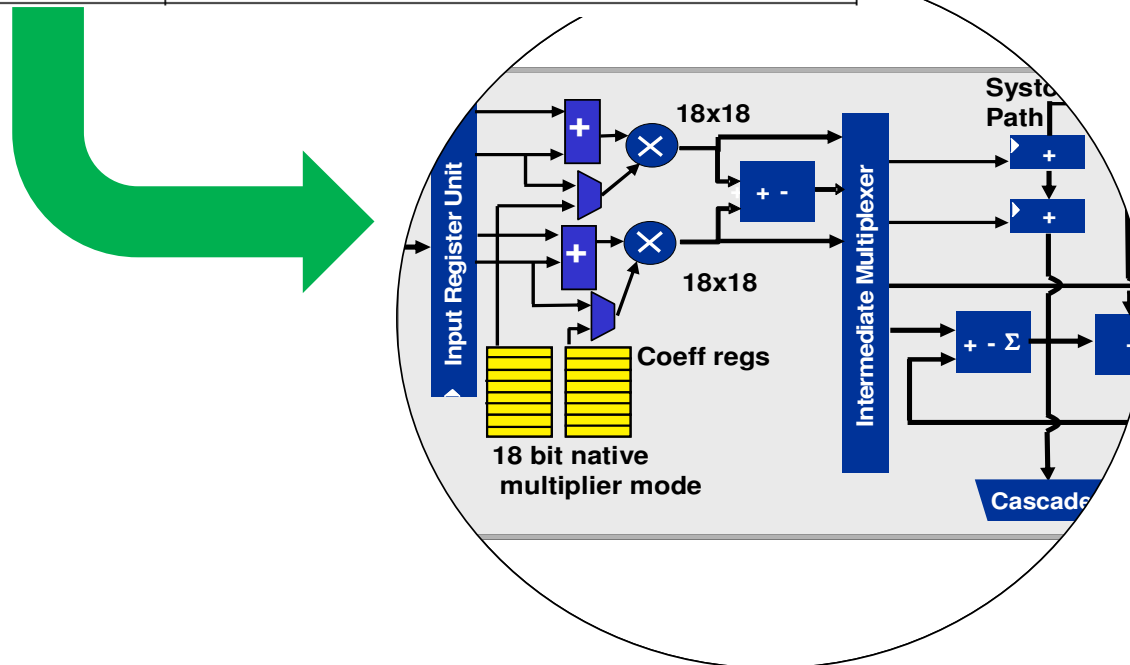- Memory controller stays **soft** for flexibility



**Stratix V FPGA PHY Architecture (UniPHY)**

# Variable-Precision DSP

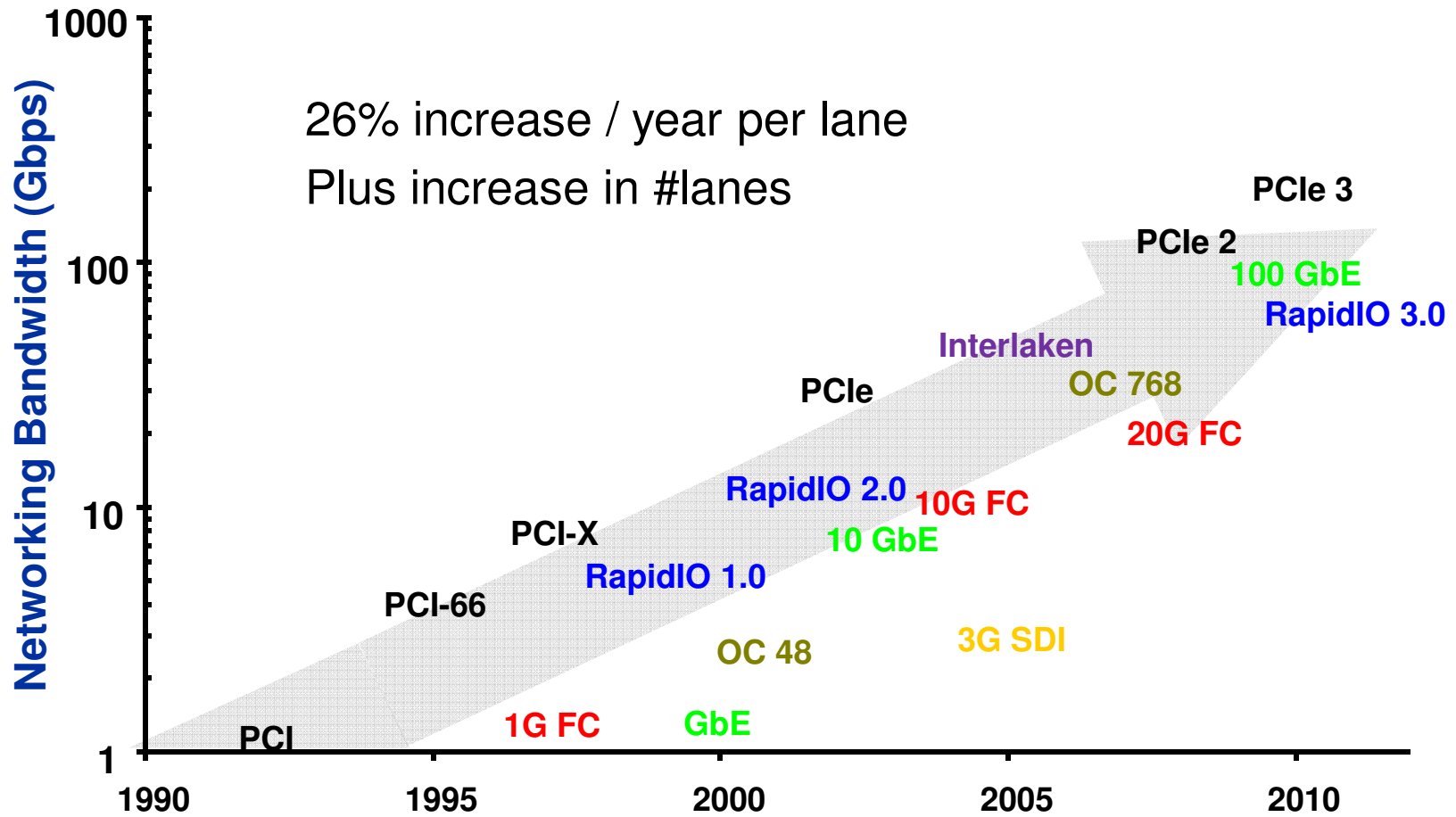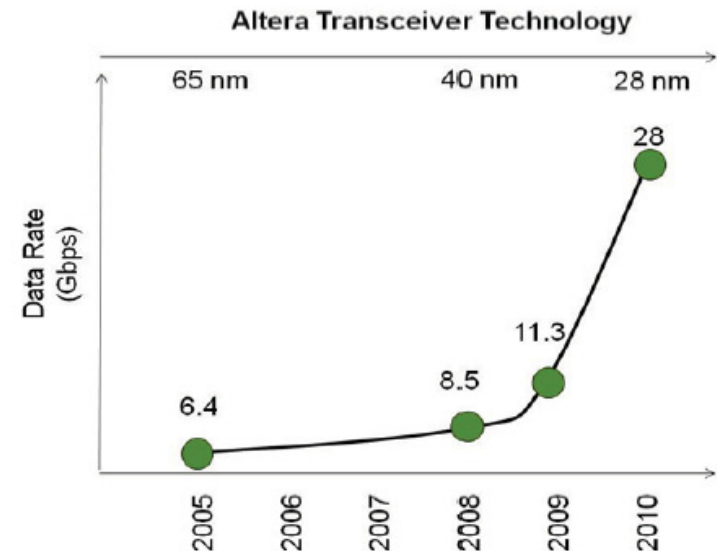| Multiplier Size | DSP Block Resources | Expected Usage |
|---|---|---|
| 9×9 bit | 1/3 of Variable Precision DSP Block | Low precision fixed point |
| 18×18 bit | 1/2 of Variable Precision DSP Block | Medium precision fixed point |
| 27×27 bit | 1 Variable Precision DSP Block | High precision fixed or single precision floating point |
| 36×36 bit | 2 Variable Precision DSP Blocks | Very high precision fixed point |
| 54×54 bit | 4 Variable Precision DSP Blocks | Double precision floating point |

Video

Wireless

Imaging/Military

# High-speed Serial Evolution



26% increase / year per lane
Plus increase in #lanes

Networking Bandwidth (Gbps)

1000 — 100 — 10 — 1

1990 — 1995 — 2000 — 2005 — 2010

PCI
PCI-66
PCI-X
1G FC
GbE
OC 48
RapidIO 1.0
RapidIO 2.0
10 GbE
10G FC
PCIe
3G SDI
Interlaken
OC 768
20G FC
PCIe 2
PCIe 3
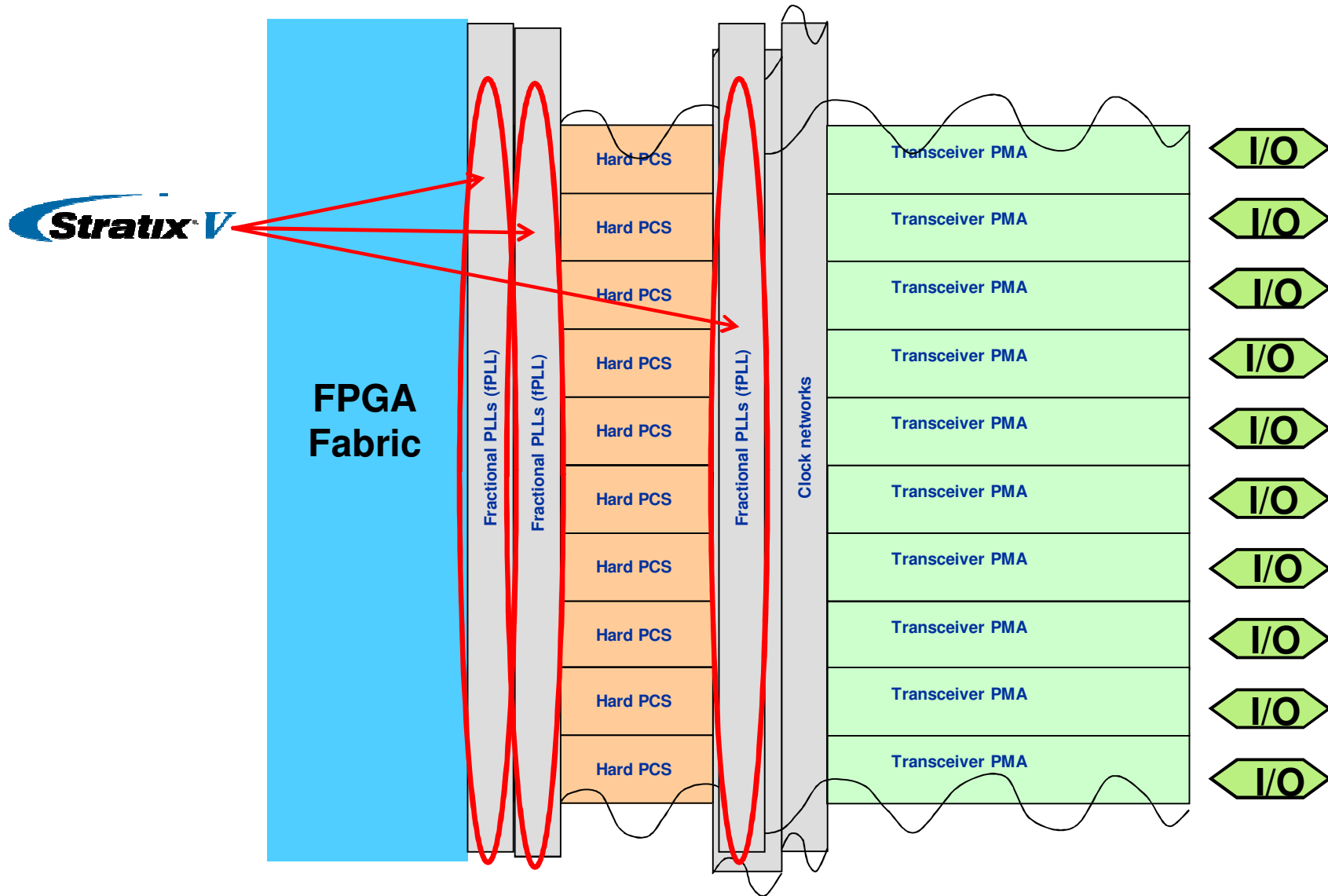100 GbE
RapidIO 3.0

17

# 28 nm Transceivers @ 28 Gbps



- **Measured on 28nm Si, shipping this quarter**

# Transcceiver Interface

# And coming soon…



FPGA with Optical Interfaces
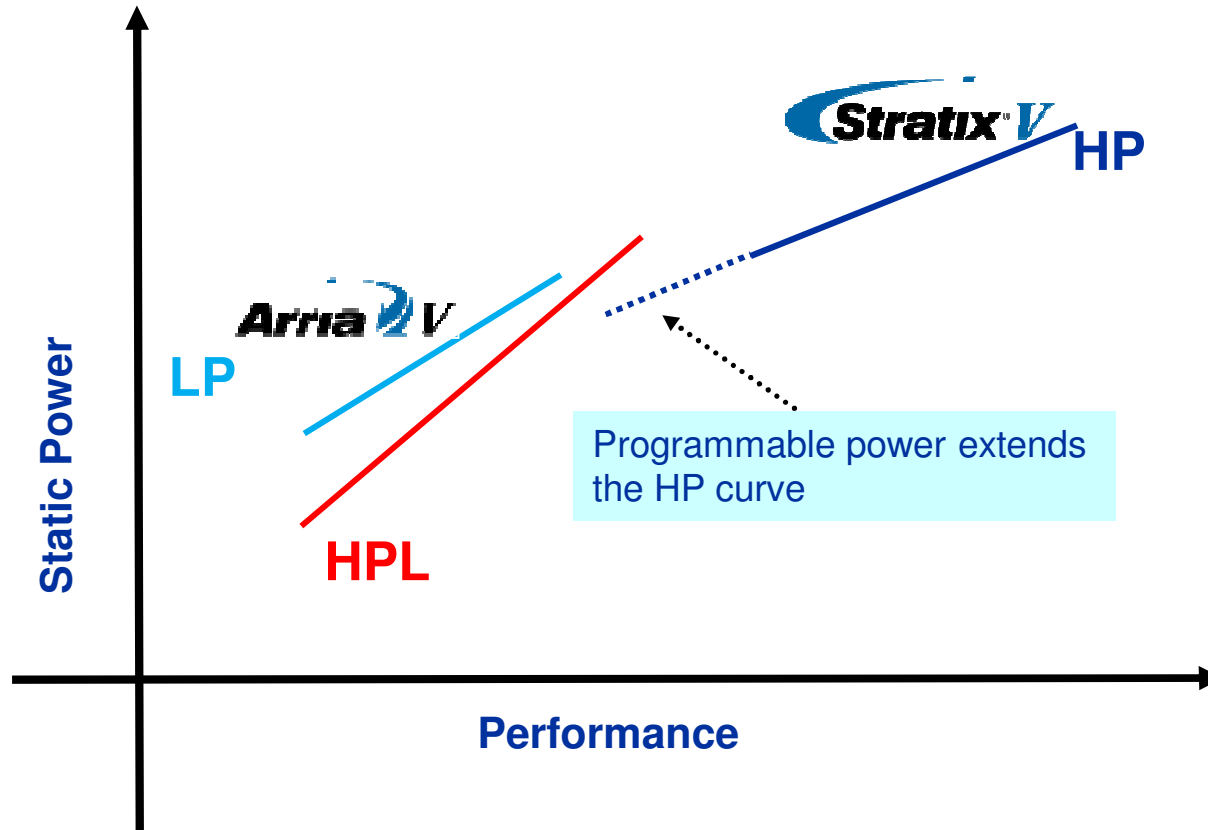
# Fabric Power and Performance

- Twice as many transistors every generation
- Power budget per device fixed
  - 10W to 20W for most higher-end applications


→ Primary application drivers will not tolerate any significant reductions in performance.

→ Must innovate to control power while improving on fmax and IO speeds.

# Choosing the Right 28-nm Process…

# Programmable Power Technology

- ## Automatic tradeoff by software
  - Performance where needed, minimizes static power elsewhere
  - Power reduction with no performance penalty

# Quartus II Software Power Optimization

**Power-Aware**

**Synthesis**

- Minimize RAM-block accesses per cycle
- Promote RE/WE to CE
- Absorb high-toggling nets

**Fitter**

- Optimize high-speed / low-power tiles
- Localize high-toggling nets for min-C
- Place circuitry to minimize clock power
- Power-efficient DSP block configurations

**Assembler**

- Program unused circuitry to minimize both AC and DC power

**15% lower dynamic power (average)**

# Fabric Architecture:

- ## Stratix V ALM architecture adds 2 additional DFF per ALM
  - Allows more pipelining, without LE-cost
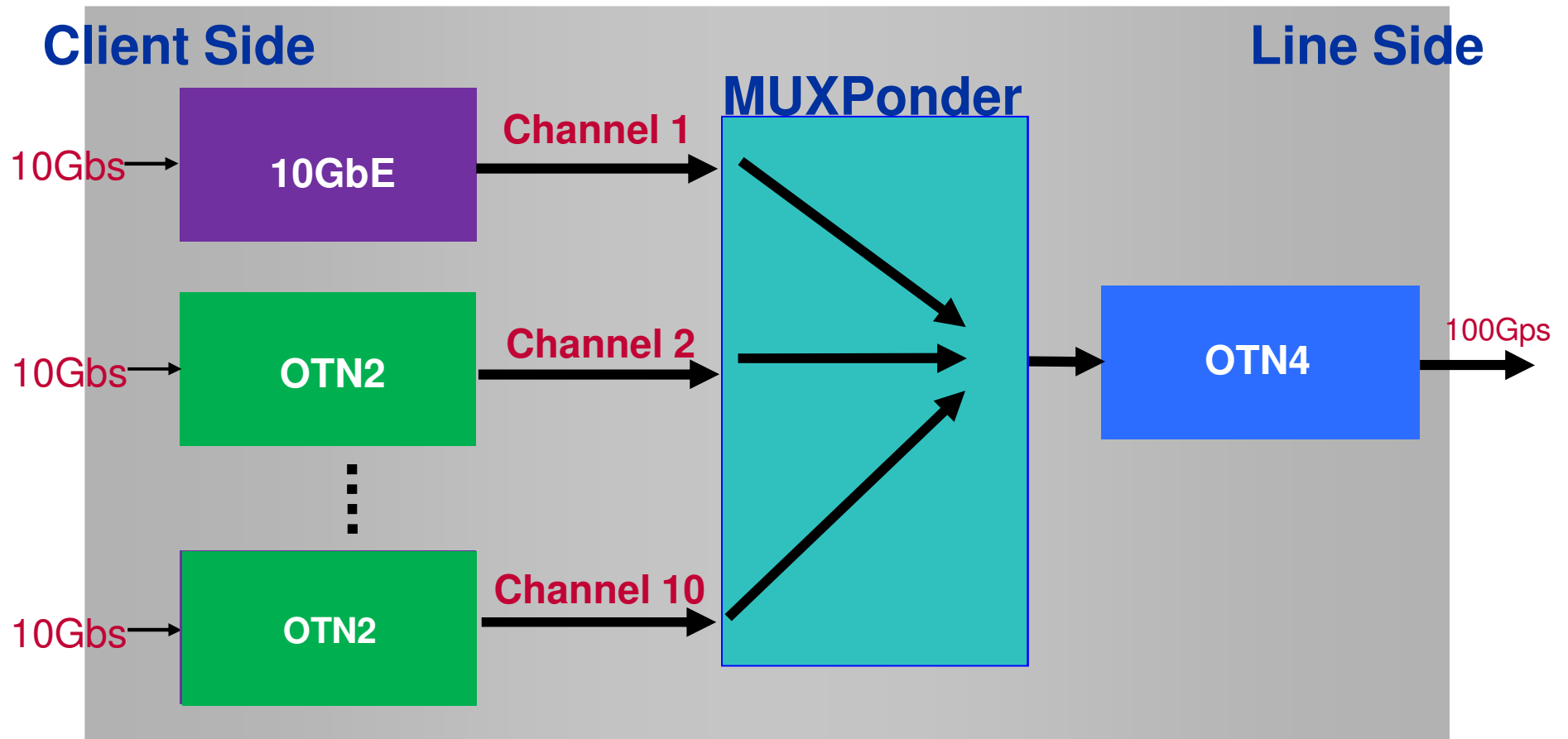  - Resulting in 6% better LE utilization over Stratix IV



- ## Re-balance of wire/metal (H3, H6, H20, V4, V12) to compensate for increased wire resistance.

# Partial Reconfiguration

# Partial Reconfiguration Flow

■ **Design Partition Based**

  – Based on Quartus Incremental Compile Flow (*partition*, *region*)

  – Enter *design intent*, automate low-level details

  – Simulation flow for operation, *including reconfiguration*

  – Each region can be configured while rest of the device is running


■ **Flow targets *multi-modal* applications**

  – Small constant number of re-configurable regions

  – Multiple *persona* (context) are compiled to a fixed location

  – An application-driven subset of the hardware capability

    ● HW can reconfigure any LAB, RAM, rmux or individual CRAM.

# Example System: 10*10Gbps→OTN4 Muxponder

# Design Entry & Simulation

- One set of HDL
- Tools to simulate during reconfig
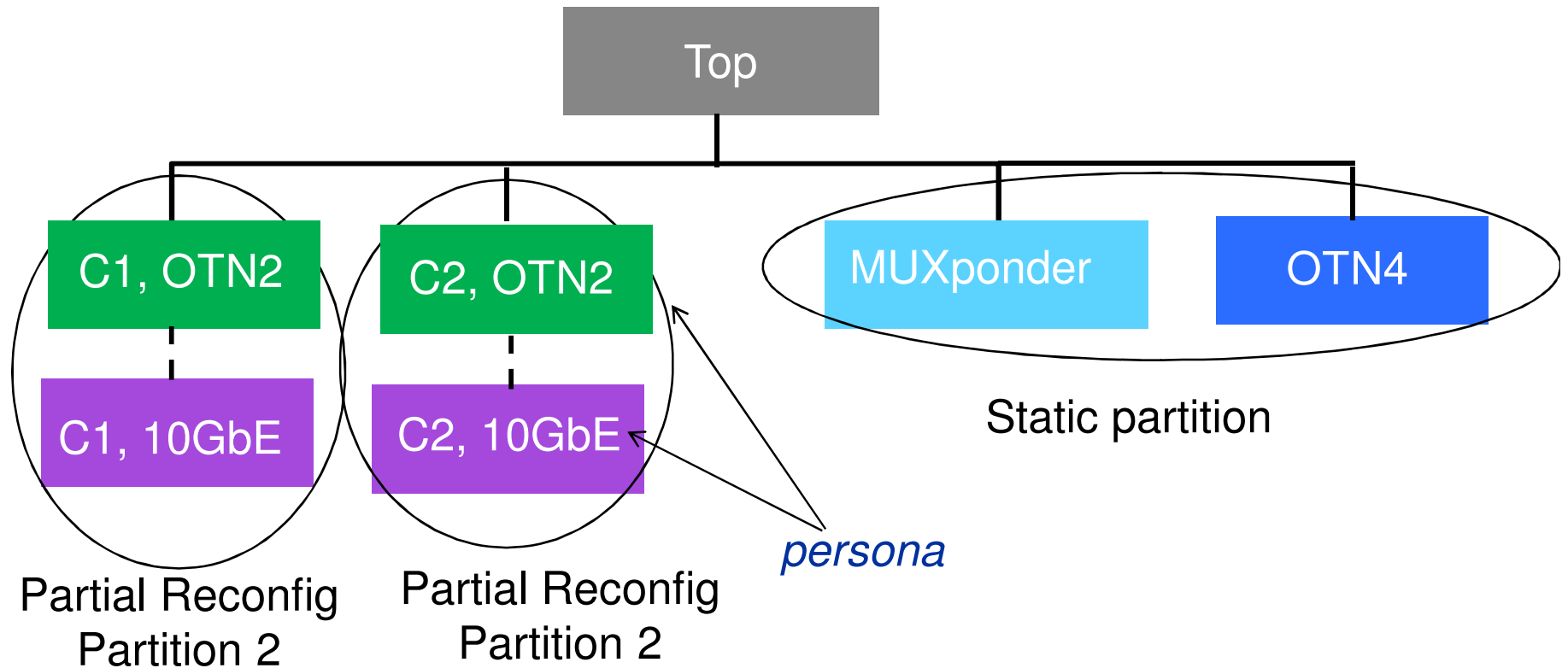
```
module reconfig_channel (clk, in, out);
input clk, in;
output [7:0] out;

parameter VER = 2;  // 1 to select 10GbE, 2 to select OTN2

generate
        case (VER)
        1: gige m_gige (.clk(clk), .in(in), .out(out));
        2: otn2 m_otn2 (.clk(clk), .in(in), .out(out));
        default: gige m_gige(.clk(clk), .in(in), .out(out));
        endcase
endgenerate

endmodule
```

# Incremental Design Flow Background

- Specify *partitions* in your design hierarchy

- Can independently recompile any partition

  - CAD optimizations across partitions prevented

  - Preserve SP&R of unchanged partitions

# *Persona* – a Partial Reconfiguration Instance



- A *revision* is a compiled subdesign of a persona
- Also, *aggregate revisions* for debug

# Partial Reconfiguration: Floorplanning

■ **Define partial reconfiguration regions**

- – Non-rectangular OK
- – Any number OK

■ **Works in conjunction with transceiver dynamic reconfiguration for dynamic protocol support**

Partial Reconfiguration for Core

# Physical:  Boundary I/O

- ## PR region I/Os must stay in same spot
  - So rest of design can communicate with any instance

- ## Solution is automatically inserted "wire LUT"
  - Tools lock this down in same spot for all instances

# Physical: Route-Throughs

- ## Routing through partially reconfigurable regions
  - Quartus II records routing reserved for top-level use
  - Prevents PR instances from using it also

# OpenCL Design Entry

# Design Flow Challenges

- HDL: low-level parallel programming language

- Timing closure
  - Device performance not increasing much in new nodes
  - Datapaths widening, device sizes growing exponentially
  - 4x28 Gbps → 336 bit datapath @ 333 MHz → need good P & R

- Compile, test, debug cycle slower than SW
  - And tools to observe HW state less mature

- Firmware development needs working HW


- Especially important for specific application domains – e.g. computing, medical imaging

ALTERA

# Altera's approach: OpenCL

- Explicitly Parallel C, not High-Level Synthesis
- Familiar to end-customer
- OpenCL programming model allows us to:
  - **Define Kernels**
    - Data-parallel computational units → can hardware accelerate
    - Including communication mechanism to kernels
  - **Describe parallelism within & between kernels**
  - **Manage Entire Systems**
    - Framework for mix of HW-accelerated and software tasks
- Multi-target

# OpenCL Structure

```
__kernel void sum { … }

__kernel void transpose {…}

float cross_product { … }
```

Program: kernels and functions

Task-level parallelism, overall framework

```
__kernel void sum
 (__global const float *a,
  __global const float *b,
  __global float *answer)
{
   int xid = get_global_id(0);
   answer[xid] = a[xid] + b[xid];
}
```

Kernels: data-level parallelism

Suitable for HW or parallel SW implementation

Specify memory hierarchy

ALTERA.

# Summary

# Challenges

- ## Explosive demand for more processing
    - Both on and off-chip bandwidth
    - Process alone not delivering the requirements

- ## FPGAs becoming SoCs / programmable ASSP
    - More hard blocks, more configurable
    - Market specialization and segmentation

- ## Designer Productivity Needs
    - Improved ease-of-use / timing closure / debug cycle
    - New paradigms to access the hardware in familiar ways

# Solutions

- ## Hardware Innovation:  28 nm & Stratix V

  - 1.5x I/O bandwidth with 28G transceivers and DDR assist

  - 1.5x – 2x more processing capability from process / hardening

  - 30% to 50% reduction in power consumption

- ## New methodologies

  - Application-driven partial reconfiguration

  - Event-driven task processors

  - SOPC and Network on Chip

  - OpenCL design initiative

- ## More needed, but 40nm/28nm have both seen big leaps in FPGA capabilities.

# Thank You!

# Technology Issues Facing the World's Largest Integrated Circuits

SPL 2011 (Cordoba, Argentina)

April 16, 2011

Mike Hutton

FPGAs are amongst the world's largest and most complex integrated circuits, and they continue to be very early adopters of the latest process technology. This talk will describe some of the driving applications and technology trends pushing FPGAs to 28 nm and smaller process nodes. We will also highlight how FPGA architecture is evolving, as exemplified by Altera's Stratix V FPGAs. Power management and silicon efficiency issues are pushing FPGAs to become somewhat more application-targeted, and to incorporate larger amounts of hard logic that makes them more complete systems-on-a-chip. In addition, the very high I/O bandwidth requirements of next-generation systems are driving innovation in both high-speed memory interface design and high-speed serial transceiver design.

Stratix V supports partial reconfiguration to increase silicon efficiency by swapping in different functionality over time. We will describe both the hardware that enables partial reconfiguration, and the software tools that will enable efficient design without becoming entangled in low-level physical details. Finally, we will discuss both software challenges and promising research efforts to create CAD tools that will help designers productively create the very large systems enabled by modern FPGAs.