# Topic 2

## SoC System Approach I
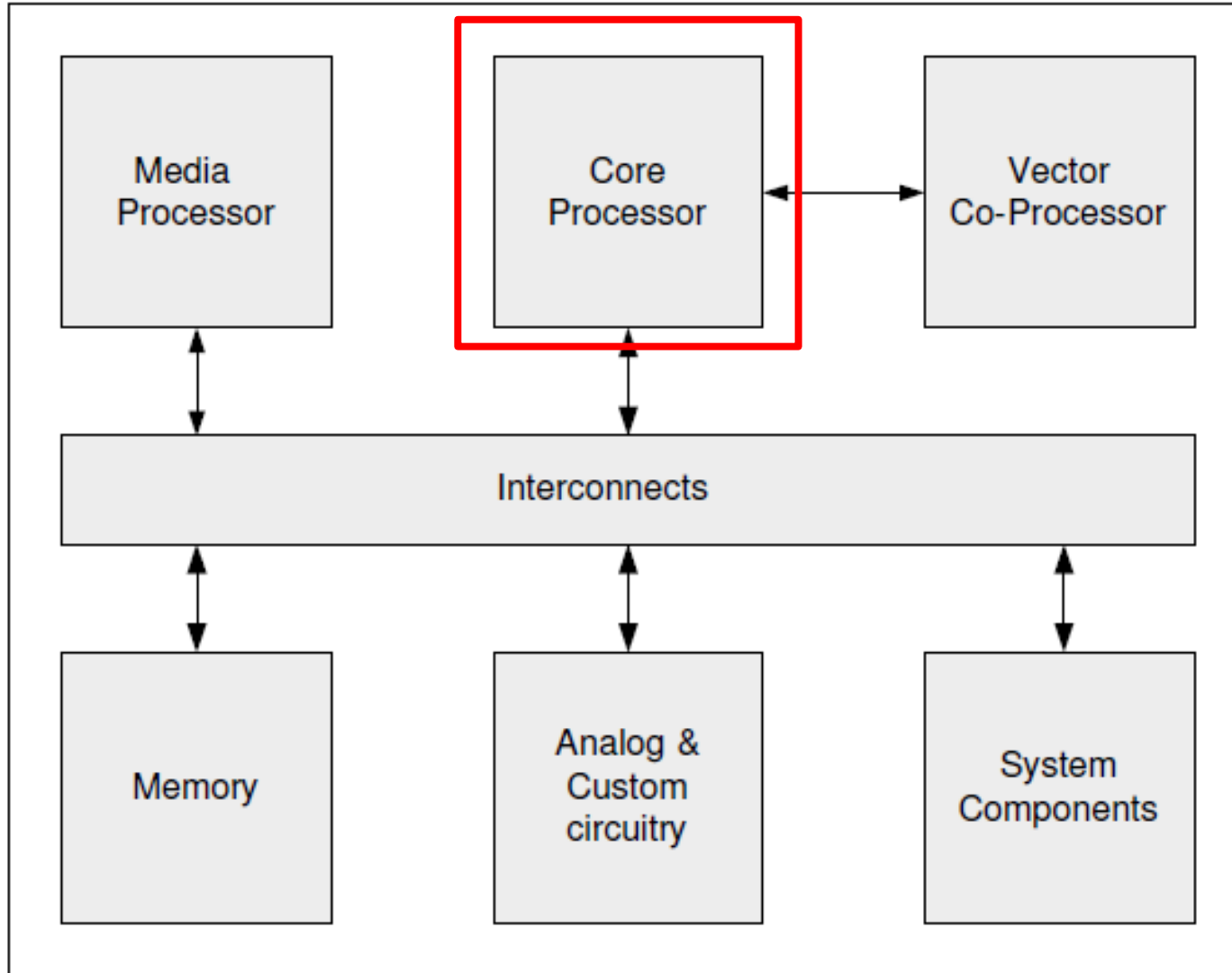
**Xinfei Guo**
**xinfei.guo@sjtu.edu.cn**

**September 25th, 2024**

# T2 learning goals

- **The philosophy of designing the SoCs…**
  - Section I
    - SoC System Design Flow Overview
    - SW/HW partition
    - NRE cost
  - Section II
    - SoC processors
    - Memory addressing
    - Design reuse

# Basic system-on-chip model

# SoC vs. processors on chip

"What distinguishes a system on a chip from the conventional general purpose computer plus memory on a board is the specific nature of the design target. The application is assumed to be <span style="color:red">known and specified, so that the elements of the system can be selected, sized and evaluated during the design process</span>." – From the textbook

# SoC vs processors on chip

- With lots of transistors, designs move in 2 ways:
  - Complete system on a chip
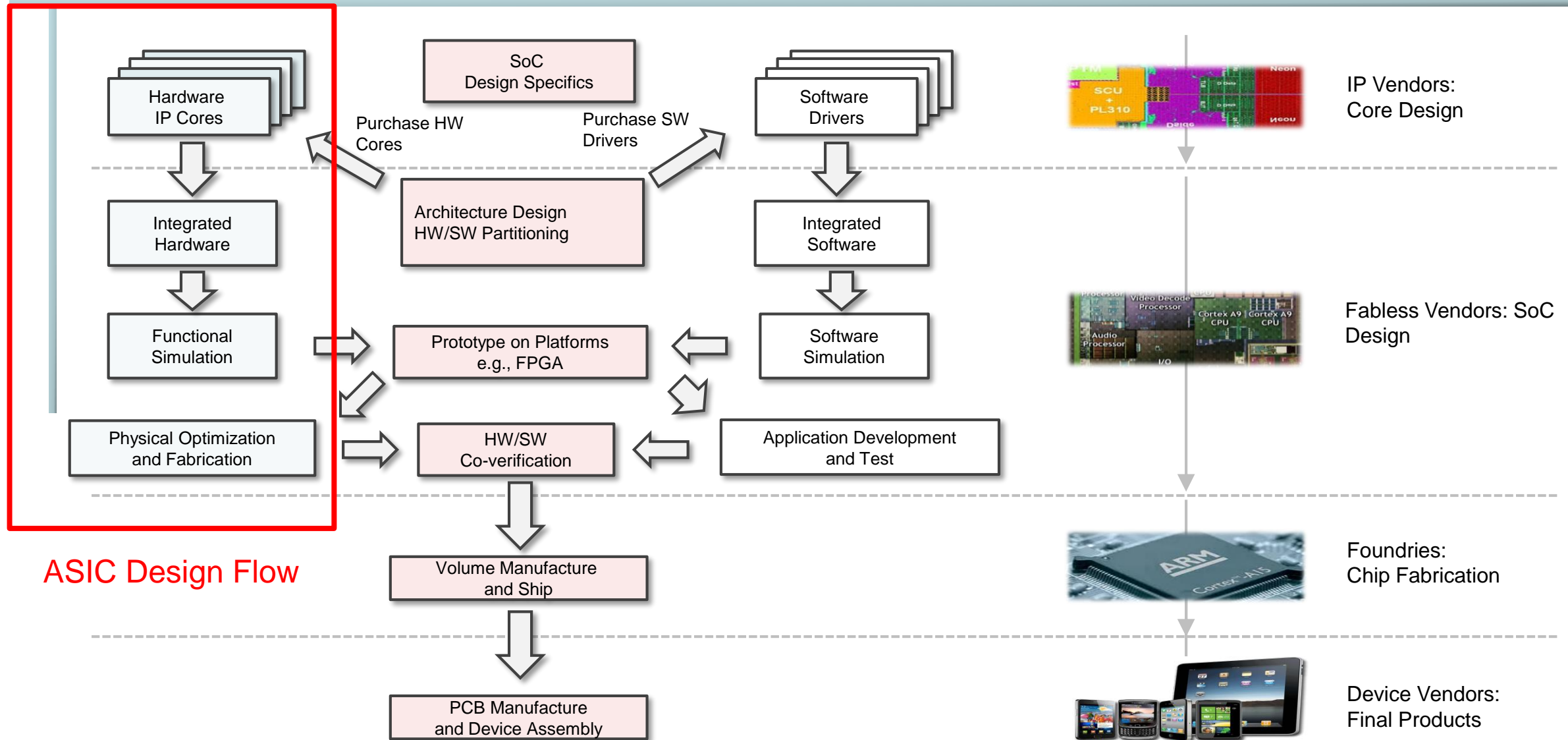  - Multi-core processors with lots of cache

|  | *System on chip* | *Processors on chip* |
|---|---|---|
| processor | multiple, simple, heterogeneous | few, complex, homogeneous |
| cache | one level, small | 2-3 levels, extensive |
| memory | embedded, on chip | very large, off chip |
| functionality | special purpose | general purpose |
| interconnect | wide, high bandwidth | often through cache |
| power, cost | both low | both high |
| operation | largely stand-alone | need other chips |

# SoC architecture and design

- System-on-chip (SOC)
  - Processors: become components in a system
- SOC covers many topics
  - Processors, cache, memory, interconnect, design tools
- Need to know
  - User view: variety of processors
  - Basic information: technology and tools
  - Processor internals: effect on performance
  - Storage: cache, embedded and external memory
  - Interconnect: buses, network-on-chip
  - Evaluation: processor, cache, memory, interconnect
  - Advanced: specialized processors, reconfiguration
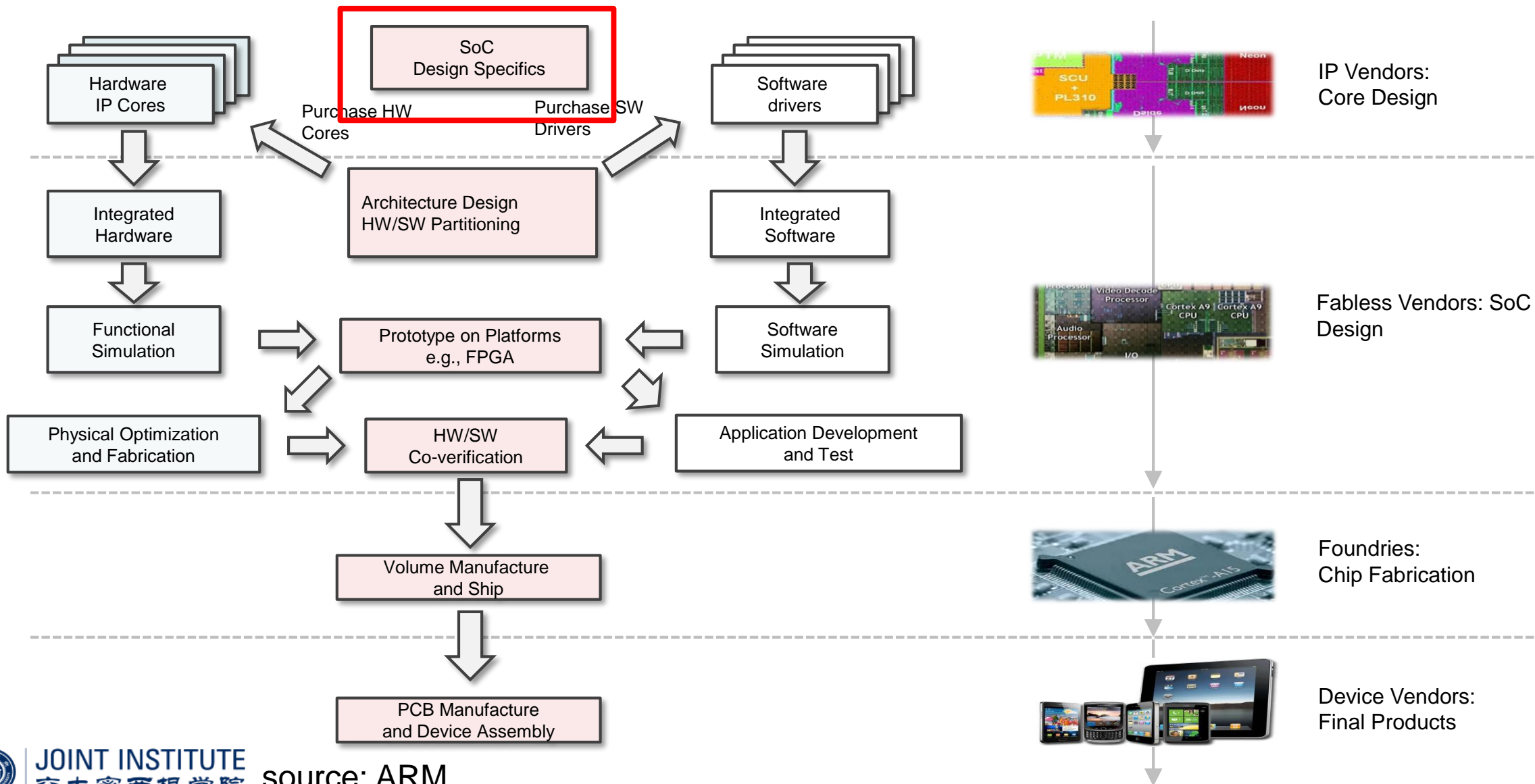  - Design productivity: system modelling, design exploration

JOINT INSTITUTE
交大密西根学院

# SOC SYSTEM DESIGN FLOW OVERVIEW
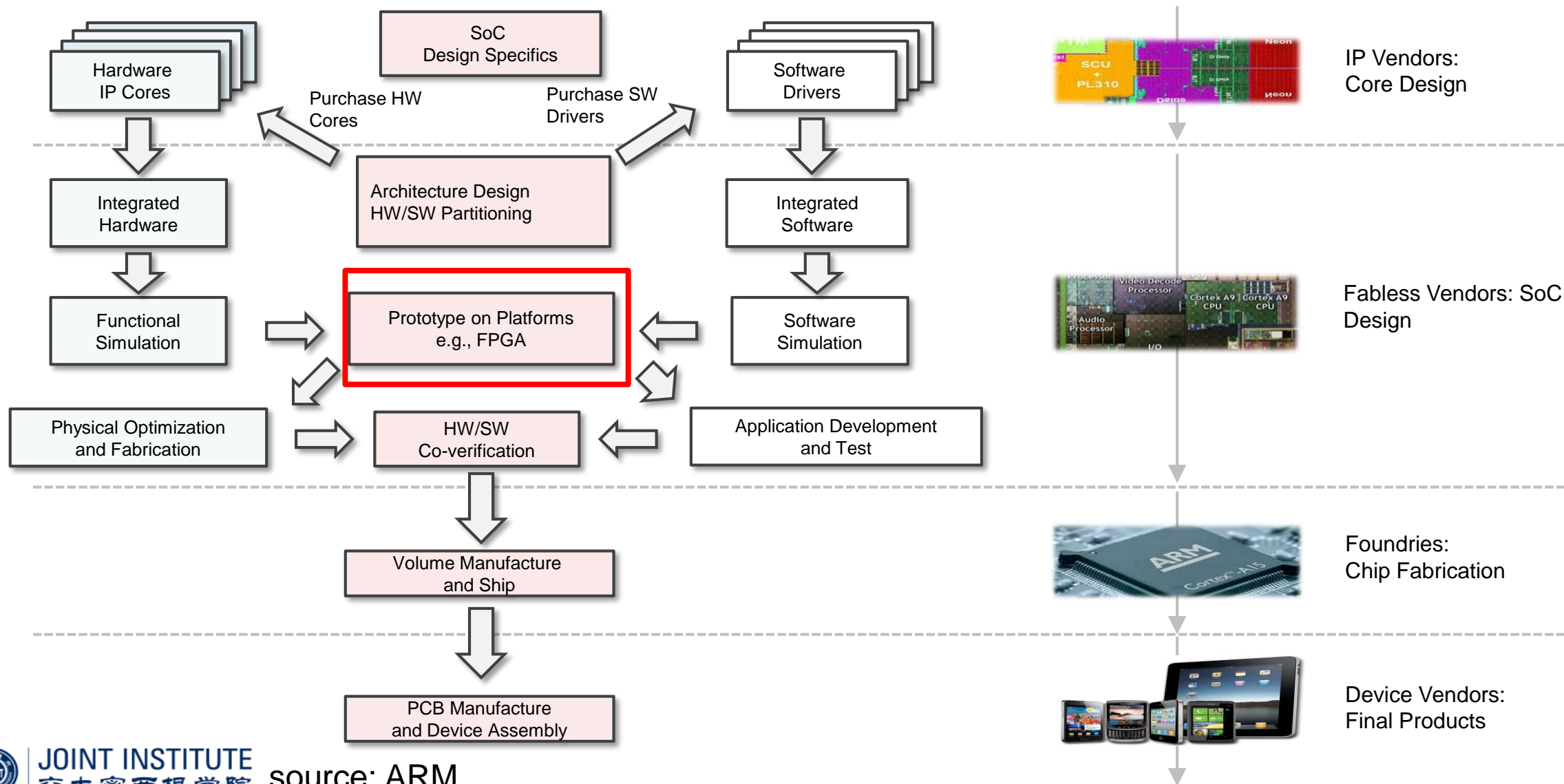
# SoC Design Flow Overview



source: ARM

# SoC Design Flow



source: ARM

# SoC Design Flow



source: ARM

# SoC Design Flow

Hardware IP Cores

SoC Design Specifics

Software Drivers

IP Vendors: Core Design

Purchase HW Cores

Purchase SW Drivers

Integrated Hardware

Architecture Design HW/SW Partitioning

Integrated Software

Fabless Vendors: SoC Design

Functional Simulation

Prototype on Platforms e.g., FPGA

Software Simulation

Physical Optimization and Fabrication

HW/SW Co-verification

Application Development and Test

Volume Manufacture and Ship

Foundries: Chip Fabrication

PCB Manufacture and Device Assembly

Device Vendors: Final Products

source: ARM

# SoC Design Flow



source: ARM

# SoC Design Flow



source: ARM

# SoC Design Flow



Hardware IP Cores

SoC Design Specifics

Software Drivers

Purchase HW Cores

Purchase SW Drivers

Integrated Hardware

Architecture Design HW/SW Partitioning

Integrated Software

Functional Simulation

Prototype on Platforms e.g., FPGA

Software Simulation

Physical Optimization and Fabrication

HW/SW Co-verification

Application Development and Test

Volume Manufacture and Ship

PCB Manufacture and Device Assembly

IP Vendors: Core Design

Fabless Vendors: SoC Design

Foundries: Chip Fabrication

Device Vendors: Final Products

source: ARM

# Quiz

Q: Which of the following statements is incorrect?

A. An SoC is typically built using cores from different vendors.

B. Testing an SoC requires hardware/software co-simulation.

C. An SoC is normally developed for a specific application.

D. The SoC hardware components can be easily replaced if faulty.

JOINT INSTITUTE
交大密西根学院

# HW/SW PARTITION

# HW vs. SW tradeoffs

■ Question: Which components in the SoC are to be implemented in HW and in SW?

| | Benefits | Drawbacks | |
|---|---|---|---|
| hardware | fast, low power consumption | inflexible, unadaptable, complex to build and test | ASIC (no fetching and decoding) |
| software | flexible, adaptable, simple to build and test | slow, high power consumption | General purpose CPU |

Implement the performance critical parts in HW and rest in SW.

# HW/SW Partitioning

- Objectives
    - Meet performance target with minimum hardware added
    - Minimize energy/power while meeting performance constraint
    - Maximum performance while keeping certain functions flexible
- Challenge: inter-dependence between the different metrics
    - Does adding application specific HW increase or decrease power?
    - How does application-specific HW impact overall system performance?



image: https://nanohub.org/courses/ECE695R/

# Programmability vs. performance

Pure HW

FPGA: faster than software, while being more flexible and having shorter development times than ASIC

- FPGA - Field-Programmable gate arrays
- CGRA - Coarse Grained Reconfigurable Architecture
- ASIP - Application-Specific Instruction Processors

Pure SW



Custom ASIC

Structured ASIC

CGRA  FPGA  ASIP  DSP

GPP

Peak performance: number of operations per watt

Low    Programmability    High

GPP: General purpose processor

# CGRA vs. FPGA

- FPGAs pay a high price for **bit-level** configurability!
  - This is often not needed by applications!
- Compilation can be easier
  - Less configurability ➔ fewer choices for compiler
  - But harder to synthesize arbitrary high-level code onto a limited array



CGRA

FPGA

# How about power?



source: https://nanohub.org/courses/ECE695R/

# Why is ASIC more efficient?

- A simple example: `add` instruction

- Q: How much of the energy consumption of an add instruction in an embedded RISC processor goes into the adder itself?

  - Answer: 1-4%

# Why is ASIC more efficient?

- Abundant Parallelism

  - SW: limited by the degree of ILP that can be exploited by the processor (superscalar), or by the number of cores (multi-core), vector lanes (vector processor), etc.

  - SW: General-purpose processors are highly pipelined, and cannot exploit specific sequences of operations that can be performed in a single clock cycle

  - HW: customizing the processing elements to be as small as needed, we can utilize orders of magnitude more "processing elements" for a given power or area budget

JOINT INSTITUTE
交大密西根学院

# Example: Google TPU

## TPUv2 Chip



- 16 GB of HBM
- 600 GB/s mem BW
- Scalar/vector units: 32b float
- MXU: 32b float accumulation but reduced precision for multipliers
- 45 TFLOPS



HBM 8 GB

core — scalar/vector units — MXU 128x128

core — scalar/vector units — MXU 128x128

HBM 8 GB

source: Google

# Example: Google TPU



**Figure 9.** Relative performance/Watt (TDP) of GPU server (blue bar) and TPU server (red bar) to CPU server, and TPU server to GPU server (orange bar). TPU' is an improved TPU (Sec. 7). The green bar shows its ratio to the CPU server and the lavender bar shows its relation to the GPU server. Total includes host server power, but incremental doesn't. GM and WM are the geometric and weighted means.

source: Google

# Any other accelerated functions?

- Multimedia (Audio / Video / Image)
- Graphics (GPUs)
- Network protocol processing
- Cryptography
- Computer vision
- Recognition / Computer Vision
- Data Mining and Search
- Scientific Computation
- AI/ML
- …

# Hardware Acceleration

■ An example: 4K HDR Video (Sample-3) Decoding



**CPU Usage: Hardware Acceleration vs Software Decoding – Computer 1**

source: https://www.5kplayer.com/video-music-player/paper-hardware-acceleration-vs-software-decoding.htm

# Hardware Acceleration

- An example: 4K HDR Video (Sample-3) Decoding



Conclusion: hardware acceleration helps decode and render videos without occupying much CPU without scarifying quality!

source: https://www.5kplayer.com/video-music-player/paper-hardware-acceleration-vs-software-decoding.htm

# A great reading

- On canvas



Hameed R, Qadeer W, Wachs M, et al. Understanding sources of inefficiency in general-purpose chips, Proceedings of the 37th annual international symposium on Computer architecture. 2010: 37-47.

# HW/SW Partitioning

- Differ primarily in how the custom HW is integrated with the processor
  - Co-processor (or HW accelerator)
  - Custom instruction

# Co-processors / HW Accelerators

- Minimal changes to the processor itself

- Accelerator is connected to the system bus or a dedicated co-processor interface

- The instruction set and pipeline structure do not change

Figure: Ravi M, Sewa A, Shashidhar T G, et al. FPGA as a hardware accelerator for computation intensive maximum likelihood expectation maximization medical image reconstruction algorithm[J]. IEEE Access, 2019, 7: 111727-111735.

# Custom Instruction Units

- Require fine-grained integration into the processor
  - Part of the processor's pipeline
- Typically need to be synthesized together with the processor
- ISA changes
- More feasible with open ISA such as RISC-V



Image: https://www.bdti.com/InsideDSP/2013/12/11/Cadence

**Cadence Xtensa processors**

# Which one to choose?

| Consideration | HW accelerator | Custom instruction |
|---|---|---|
| How is custom HW interfaced? | System bus Dedicated co-proc. I/F | Custom instruction interface Directly into processor pipeline |
| How does SW access the HW | Memory-mapped I/O. Driver abstracts the operation performed by the accelerator as API for SW | Compiler maps operations to custom instructions. Intrinsics (macros). |
| Ability to access memory hierarchy | Direct access to main memory, cannot access cache / registers | Access typically limited to register file, same view of memory as processor |
| Parallel execution with the processor | Yes, assuming SW does not block | Pipelining custom instruction allows other instructions to execute concurrently |
| Granularity of computations targeted | Coarse (100s – millions of cycles) | Fine (few – tens of cycles) |

source: https://nanohub.org/courses/ECE695R/

# How to do partitioning?

- a.k.a How to decide which part to accelerate?

# The four "C"s in SoC Design

- Is it a **C**omputational bottleneck?
- What is the **C**ommunication overhead?
- Is the HW implementation **C**ost-efficient?
- Is it a **C**ommonly used function with a Clean interface?

Amdahl's law

$$Speedup = \frac{T_{orig}}{T_{unaccel} + T_{accelerated} + T_{communication}}$$

$$\leq \frac{1}{\dfrac{T_{unaccel}}{T_{orig}} + \dfrac{T_{communication}}{T_{orig}}}$$

source: https://nanohub.org/courses/ECE695R/

# HW/SW partitioning challenges

- Specification abstraction level
- Granularity
- System-component allocation
- Metrics and estimations
- Partitioning algorithms
- Objective and closeness functions
- Flow of control and designer interaction

- …

# NRE COST

# NRE Cost

- Non-Recurring Engineering (NRE)
  - Costs spent up front on development
    - Engineering Design Time
    - Prototypes
    - Mask costs
    - …
  - Recurring Engineering - Costs to produce each chip
  - Cost(Nchips) = Cost (NRE) + Nchips x Cost(manufacturing) /chip

# Cost: product program vs engineering



*Fixed costs*

*Variable costs*

Marketing, sales, administration

Engineering

Manufacturing costs

Product cost

Chip design

Verify & test

Software

CAD support

Labor costs

Engineering costs

Mask costs

CAD programs

Capital equipment

Fixed project costs

40

# Variable (RE) Cost

- Variable costs (cost per part)
  - Wafer cost
  - Wafer processing
  - Die size (# die per wafer)
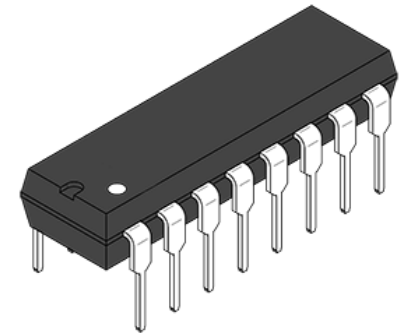    - Size of design (# gates)
    - Technology (# gates per sq. inch)
    - % utilization of die
  - Production yield = f(defect density, die size)
  - Packaging



https://www.waferworld.com/post/top-causes-of-silicon-wafer-breakage



https://www.quick-pcba.com/pcb-news/ic-packaging-types.html

# Cost is increasing!



source: https://semiengineering.com/how-much-will-that-chip-cost/

# Product volume dictates design effort

Design time
and effort

Basic
physical
tradeoffs

Balance point depends on
n, number of units

# What is the solution?

- Amortize NRE with Volume

Cost(Nchips) = Cost (NRE) + Nchips x Cost(manufacturing) /chip

Cost/chip = Cost (NRE) / N(chips) + Cost(manufacturing) /chip

# Market requirements

- Economics force fewer, more customizable chips
- Mask costs in the millions of dollars
- Custom IC design NRE 10s—100s of millions of dollars
- Number of unique chips must decrease
- Increase the programmability is one option



**Number of ASIC Design Starts**

Legend: ■ Traditional ASIC ☐ Structured ASIC

Values annotated: -1.8, -2.4, -6.4

© 2006 Gartner, Inc. All Rights Reserved.

Fig. 1

Gartner.

# SoC Profit Model

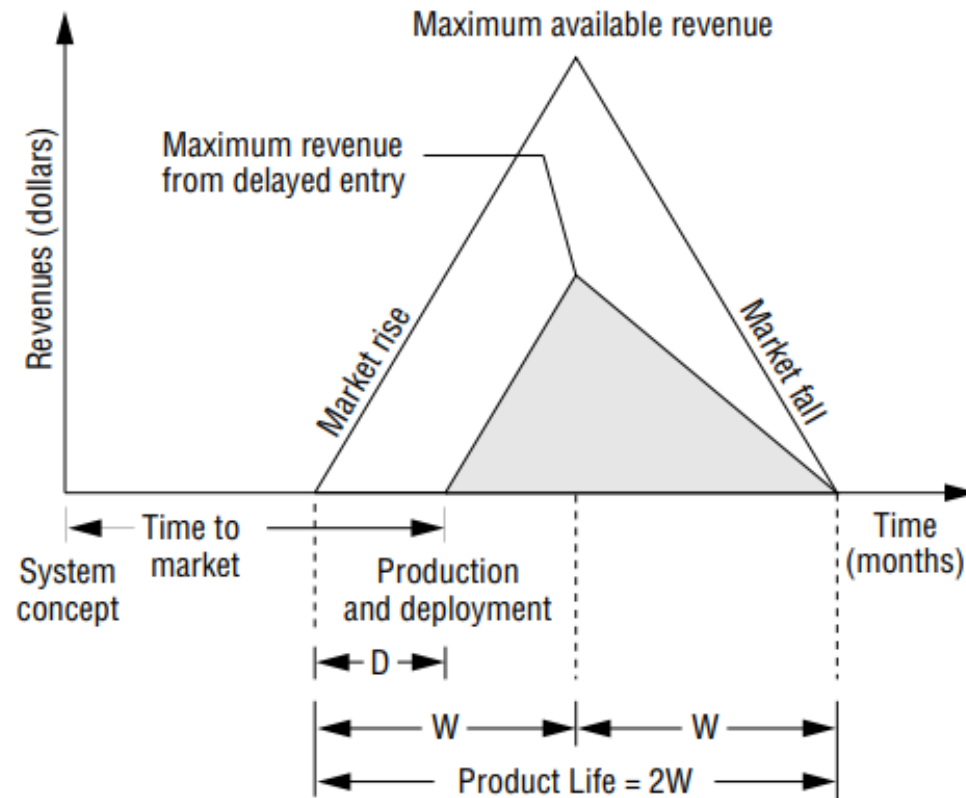- Simplified revenue model
  - Product life = 2W, peak at W
  - Time of market entry defines a triangle, representing market penetration
  - Triangle area equals revenue
- Loss
  - The difference between the on-time and delayed triangle areas

Debardelaben J A, Madisetti V K, Gadient A J. Incorporating cost modeling in embedded-system design[J]. IEEE Design & Test of Computers, 1997, 14(3): 24-35. **(Available on Canvas)**

# Percentage revenue loss



- On-time = ?
- Delayed = ?
- Percentage revenue =

(On-Time – Delayed)/On-Time*100% = ?

source: https://www.arteris.com/blog/bid/112221/what-does-it-cost-you-when-your-soc-is-late-to-market

47

# Quick Summary - SoC design: key ideas

- To design and evaluate an SoC, designers need to understand:
  - its components: processors, memory, interconnect
  - applications that it targets
- SoC economics heavily dependent on:
  - costs: initial design, marginal production
  - volume: applicability, lifetime
- Reducing design complexity
  - Intellectual Property (IP)
  - reconfigurable technology

# Where are we Heading?

- SoC System Approach II

# Action Items

- Lab #0 due by Sunday Sept. 28th 23:59

- Watch out announcements for lab logistics, lab group assignment, etc.

- Reading Materials
  - Ch. 1.2 - 1.4,1.9, Two research papers on canvas

# Acknowledgement

Slides in this topic are inspired in part by material developed and copyright by:

- Dr. Wayne Luk (Imperial College)
- Dr. Gul N. Khan (Ryerson University)
- Dr. André DeHon (UPenn)
- Dr. Xuan Zhang (WUSTL)
- Dr. Andreas Gerstlauer (UT Austin)
- Dr. Anand Raghunathan (Purdue)

JOINT INSTITUTE
交大密西根学院