# TU WIEN Informatics

# Advanced Computer Architecture

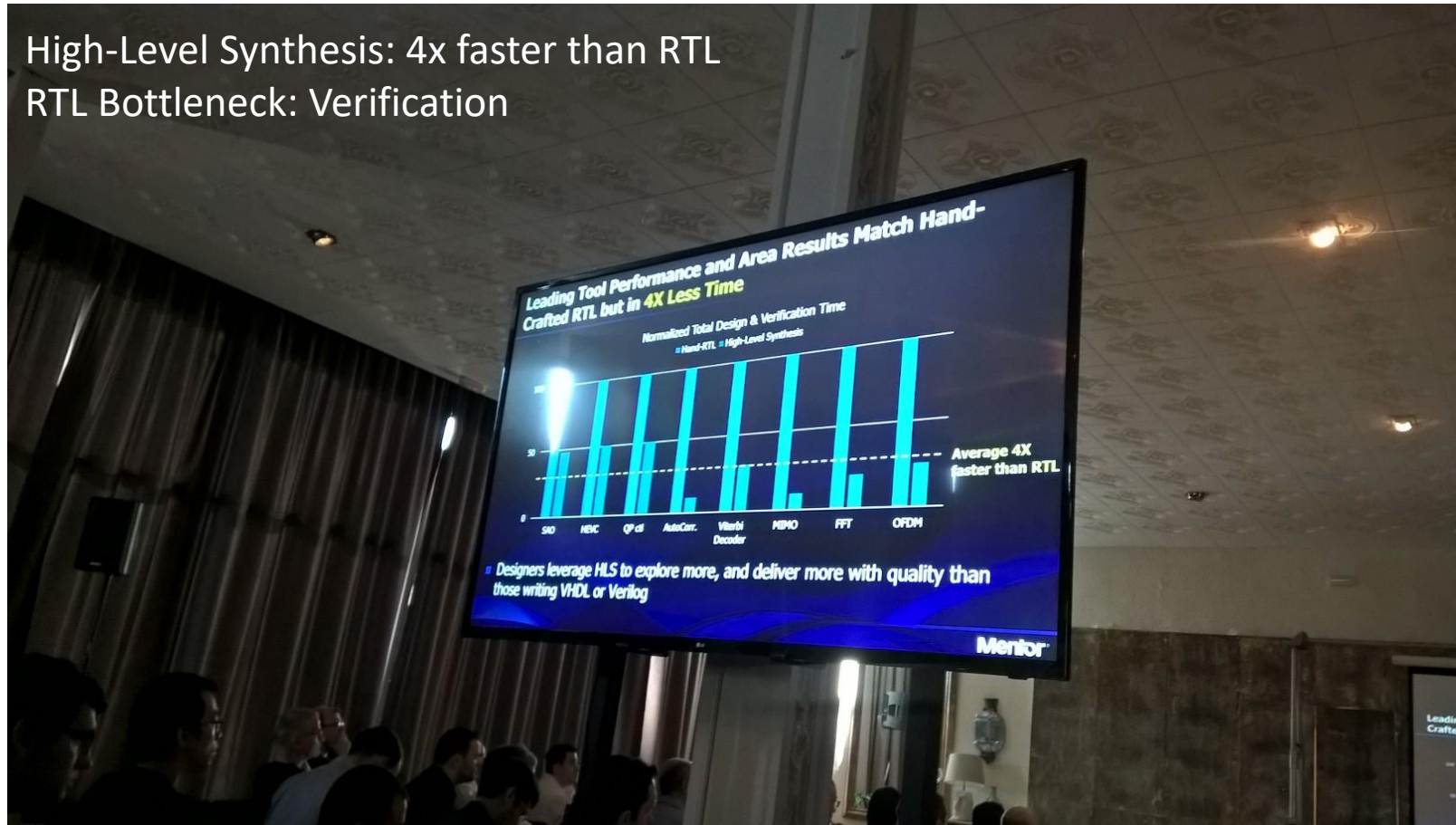D1 – Introduction to High Level Synthesis (HLS)

Daniel Mueller-Gritschneder

# Motivation for HLS

Source: WALDEN C. RHINES
President and Chief Executive Officer , Mentor, a Siemens Business
24th IEEE International Symposium on On-Line Testing and Robust System Design 2018

High-Level Synthesis: 4x faster than RTL
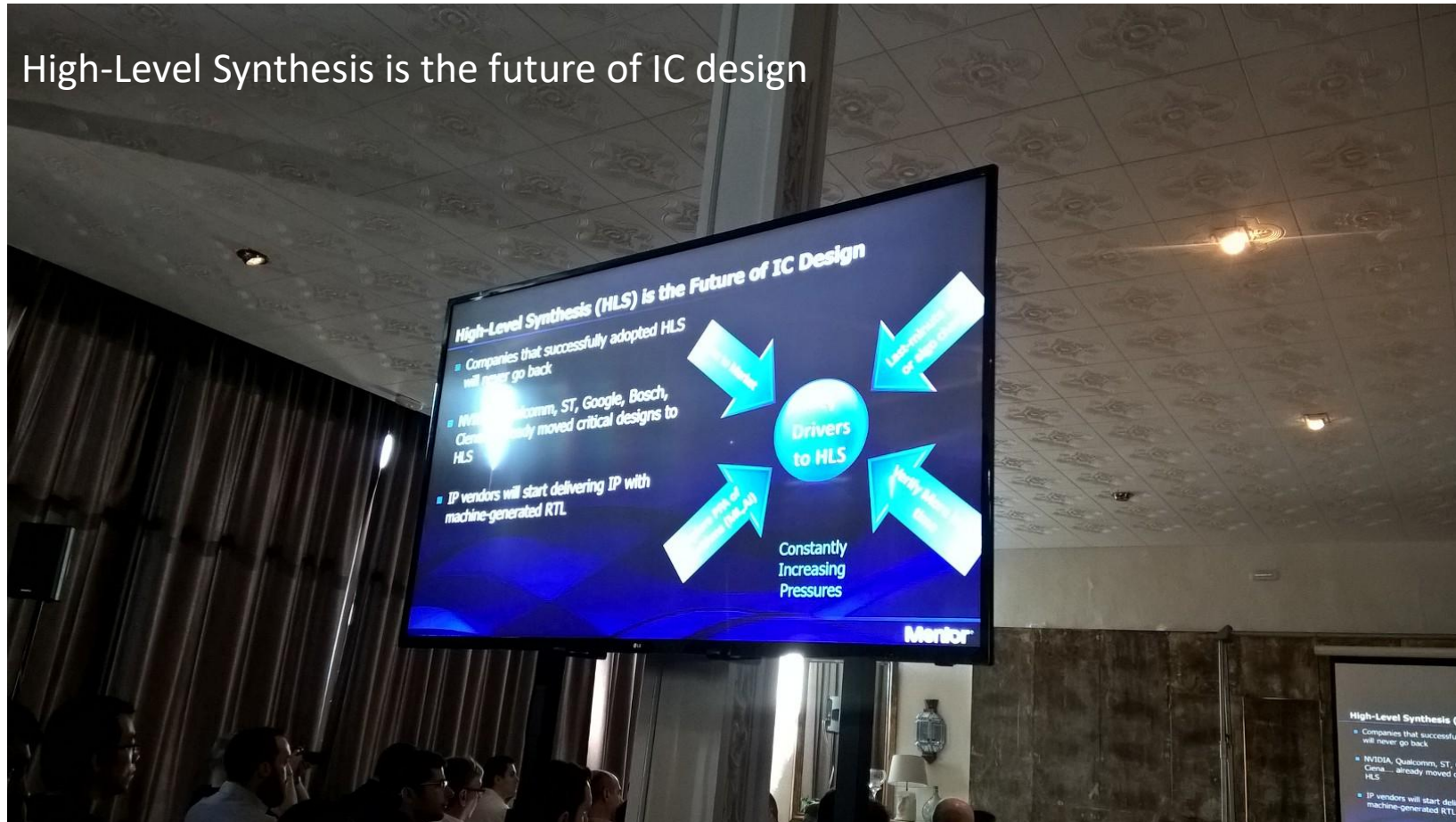
RTL Bottleneck: Verification

Source: WALDEN C. RHINES
President and Chief Executive Officer , Mentor, a Siemens Business
24th IEEE International Symposium on On-Line Testing and Robust System Design 2018



High-Level Synthesis is the future of IC design
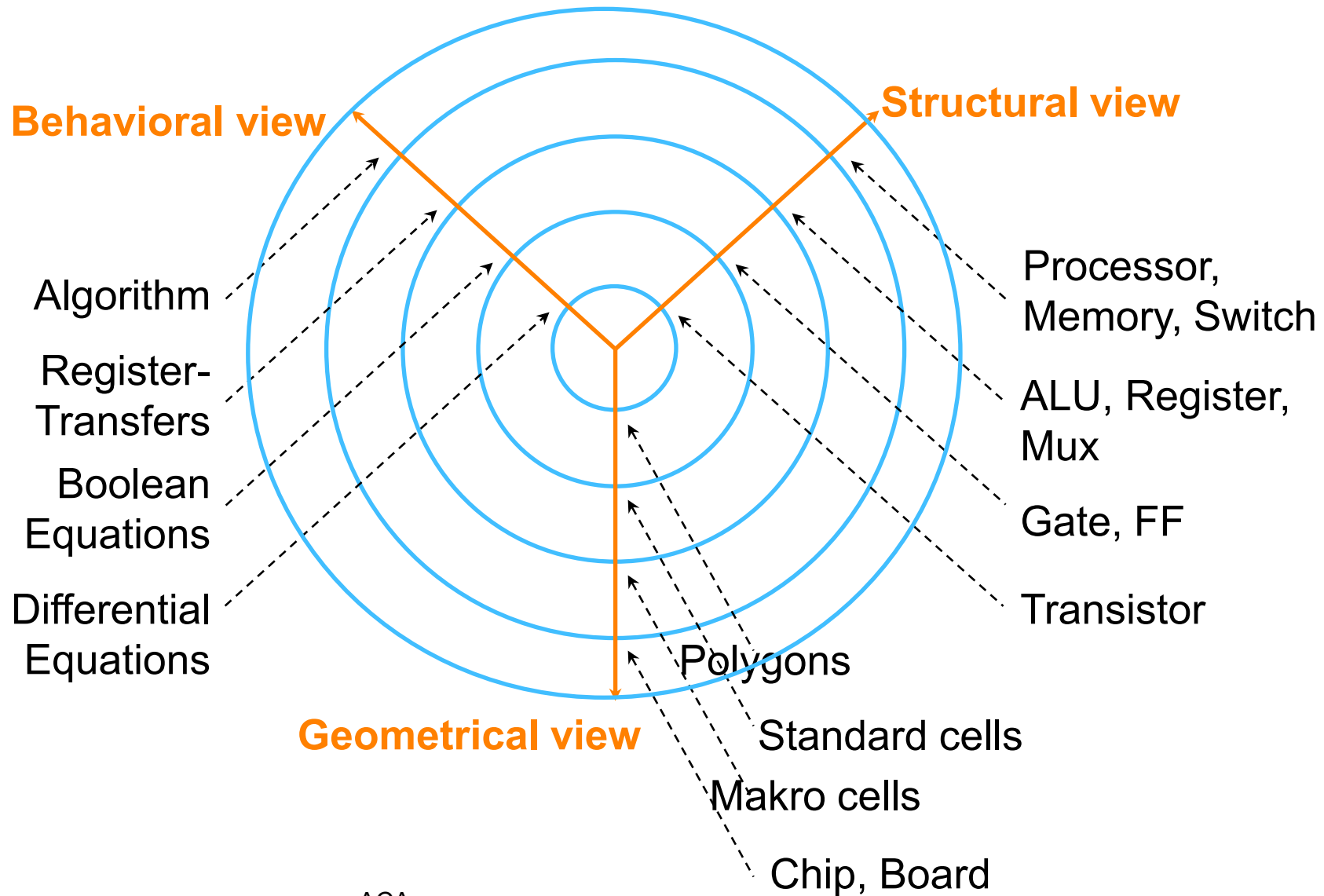
# D1-1 HW Design Flow in a Nutshell

- Literature:

- *„Specification and Design of Embedded Systems"* Daniel D. Gajski, Prentice Hall 1994

- *„Digitale Hardware/Software Systeme"*, Jürgen Teich, Springer 1997

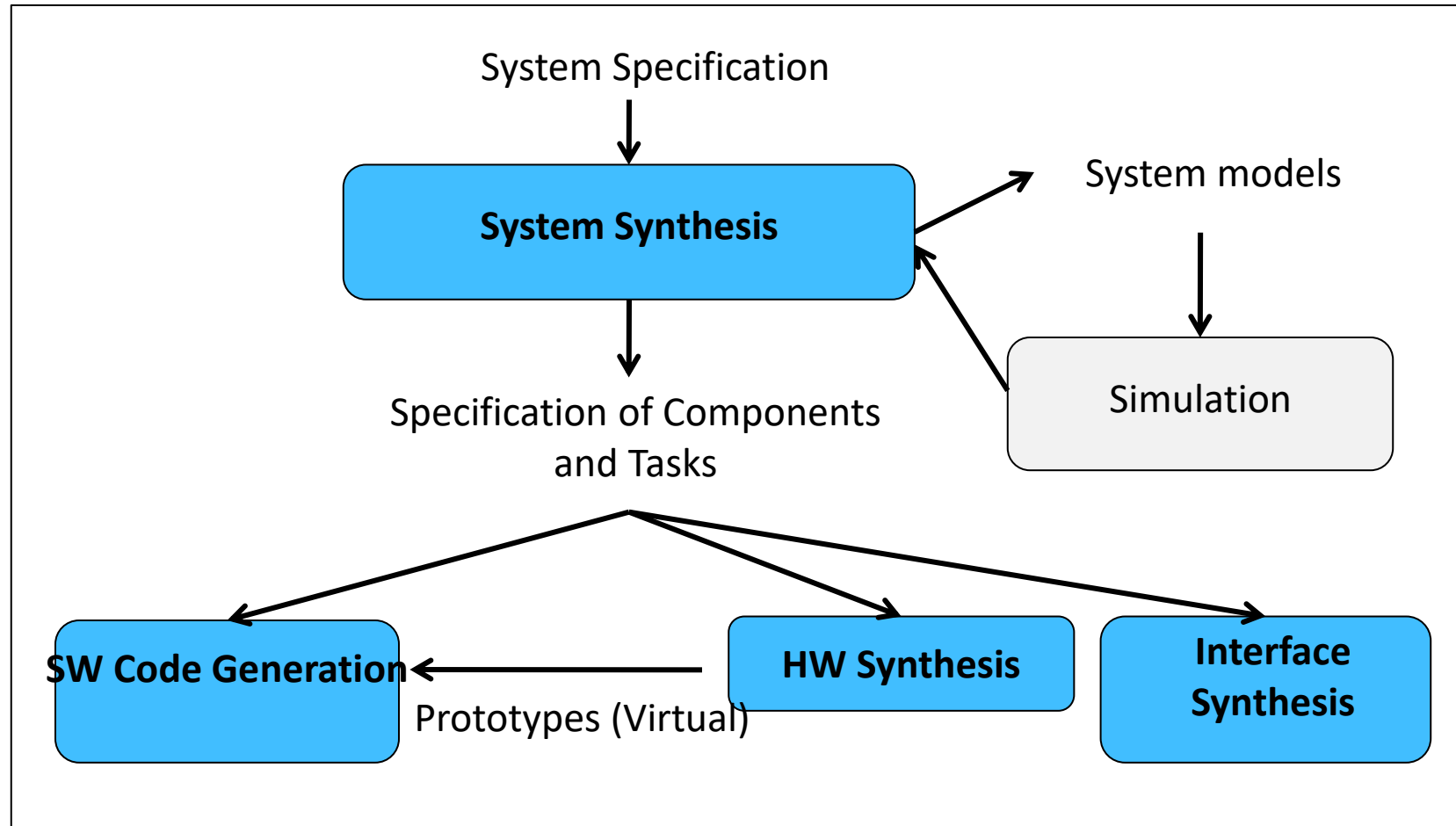- *„Embedded System Design"*, Daniel D. Gajski et.al., Springer 2009

# Abstraction Levels & Design Views

| | | Design View | | |
|---|---|---|---|---|
| | | *Behavior* | *Structure* | *Geometry* |
| *Abstraction Level* | *System* | System Specification | Connected Components | Chip, Board |
| | *Architecture* | Algorithms | CPU, Bus, HW-accelerator | Floor plan |
| | *Register Transfer* | Register Transfers / FSMs | Module netlist (ALU, Mux, Register) | Makro-cells (IP-blocks) |
| | *Logic* | Boolean Equations | Gate netlist (Gates, FlipFlops) | Standard cells, library cells |
| | *Circuit* | Differential Equations | Transistor netlist | Mask data |

# Abstraction Levels & Design Views



Behavioral view

Structural view

Geometrical view

Algorithm

Register-Transfers

Boolean Equations

Differential Equations

Processor, Memory, Switch

ALU, Register, Mux

Gate, FF

Transistor

Polygons

Standard cells

Makro cells

Chip, Board

System Specification

**System Synthesis**

System models

Simulation

Specification of Components and Tasks

**SW Code Generation**

**HW Synthesis**

**Interface Synthesis**

Prototypes (Virtual)

ACA
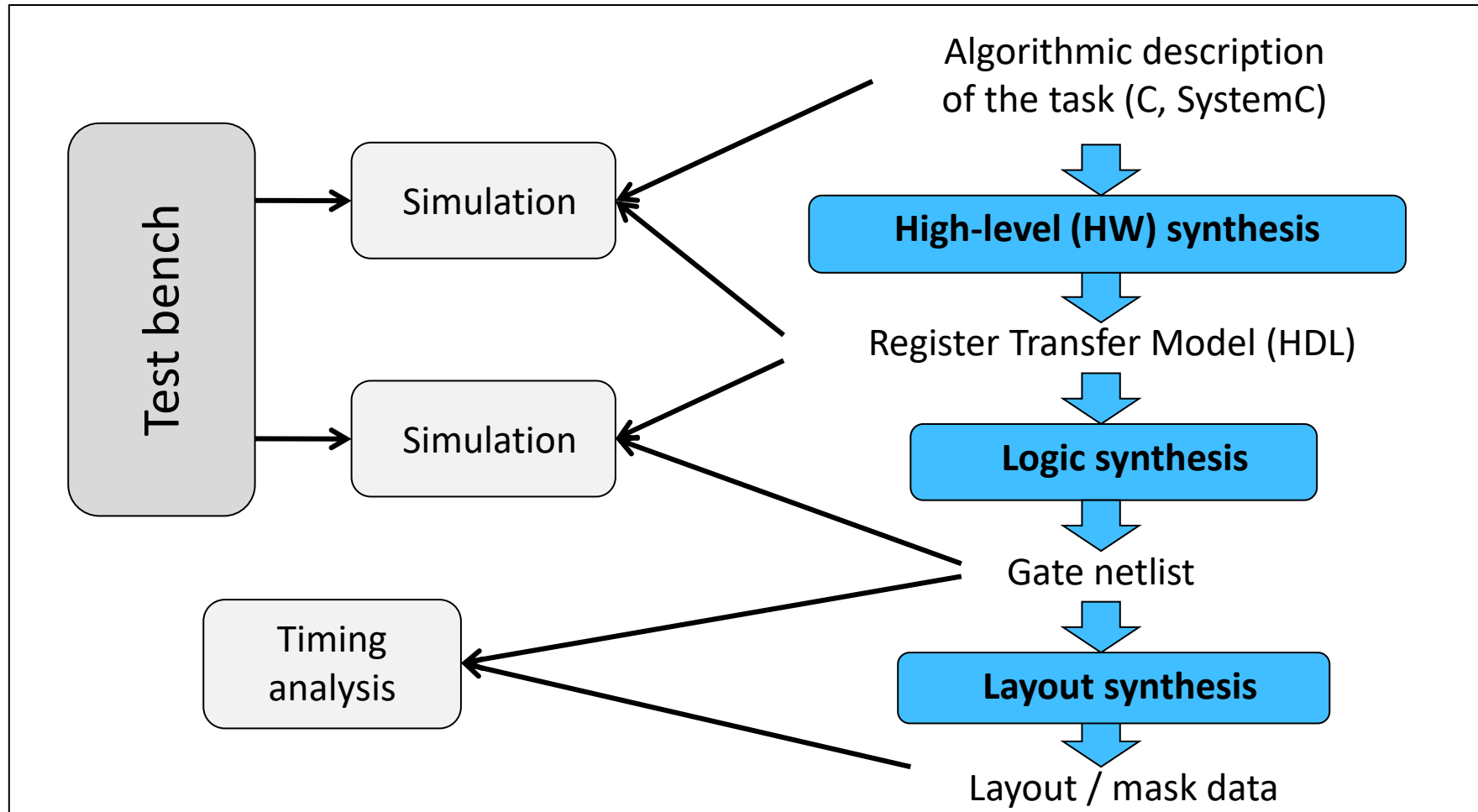
# System Synthesis

- **Inputs:**
  - Specification of the System: Description of the functionality and design constraints (Written text, specification languages)
- **Typical synthesis steps:**
  - Description of functionality as set of communicating g tasks
  - Description of behavior of tasks on algorithmic level
  - Description of task communication
  - Allocation of system components such as processors, buses, memory, …  Buses, memory,…
  - Binding of tasks and inter-task communication to system components (HW/SW Partitioning)
- **Output**
  - Specification of components, tasks and Inter-task communication that guarantees to meet system specification

ACA

# ASIC HW Synthesis Flow

ACA

# HLS Synthesis Step

- **Input**
  - Algorithmic description of a task, e.g. in C, C++, SystemC
  - Design constraints (Maximal latency, available resources, …)

- **Synthesis steps:**
  - Static code analysis and code optimization
  - Data path synthesis (Scheduling, allocation, binding)
  - Control unit synthesis (FSM implementation)

- **Output:**
  - Description of hardware module on RT level

ACA

- **Input:**
  - Description of HW module on RT level
  - Design constraints (minimal clock frequency, maximal area,...)
  - Gate library

- **Synthesis steps:**
  - Logic optimization
  - Technology mapping


- **Output:**
  - Gate netlist

ACA

- **Input**
  - Gate library
  - Design constraints
  - Layout library (P-cells)

- **Synthesis steps:**
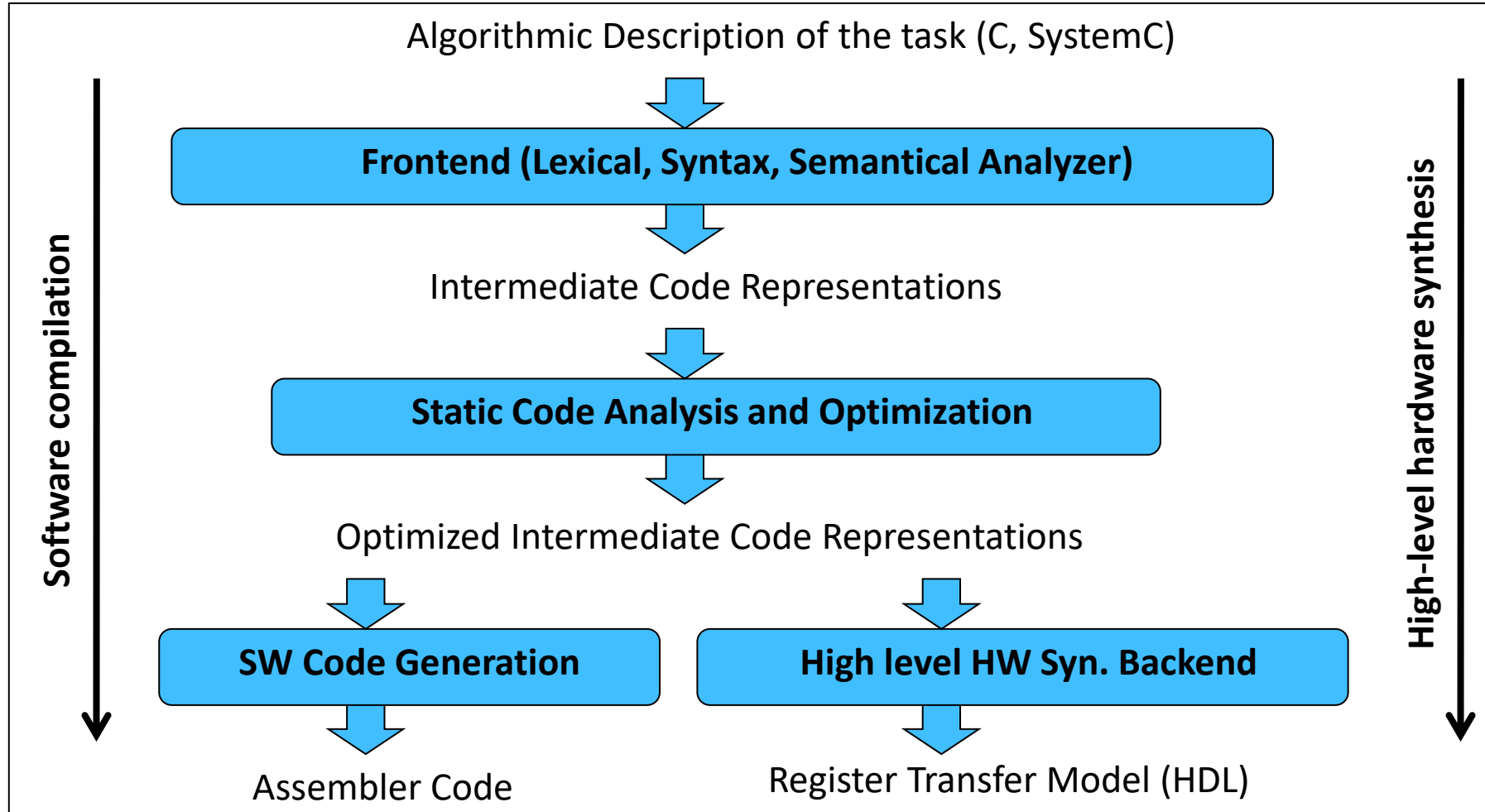  - Placement of modules
  - Routing of signal nets

- **Output**
  - Layout, mask data

ACA

# Software Compilation

- Inputs
  - Algorithmic description of task

- Synthesis steps
  - Static Code Analysis and Optimization
  - Code Generation (Instruction Selection, Register Allocation and Assignment)
  - Assembler/linker/loader

- Outputs:
  - Assembly code/machine code for target processor

ACA

# Interface Synthesis Step

- Input
  - Description of Inter-task communication
  - Design constraints (protocols, data rates, …)

- Outputs
  - Drivers, bus interfaces, …

ACA

# High level HW Synthesis (HLS) vs. SW Compilation Flow

Algorithmic Description of the task (C, SystemC)

**Frontend (Lexical, Syntax, Semantical Analyzer)**

Intermediate Code Representations

**Static Code Analysis and Optimization**

Optimized Intermediate Code Representations

**SW Code Generation**

**High level HW Syn. Backend**

Assembler Code

Register Transfer Model (HDL)

Software compilation

High-level hardware synthesis

# D1-2 The HLS Synthesis Task

- Literature:

- *„Specification and Design of Embedded Systems"* Daniel D. Gajski, Prentice Hall 1994

- *„Digitale Hardware/Software Systeme"*, Jürgen Teich, Springer 1997

- *„Embedded System Design"*, Daniel D. Gajski et.al., Springer 2009

# Basic Task

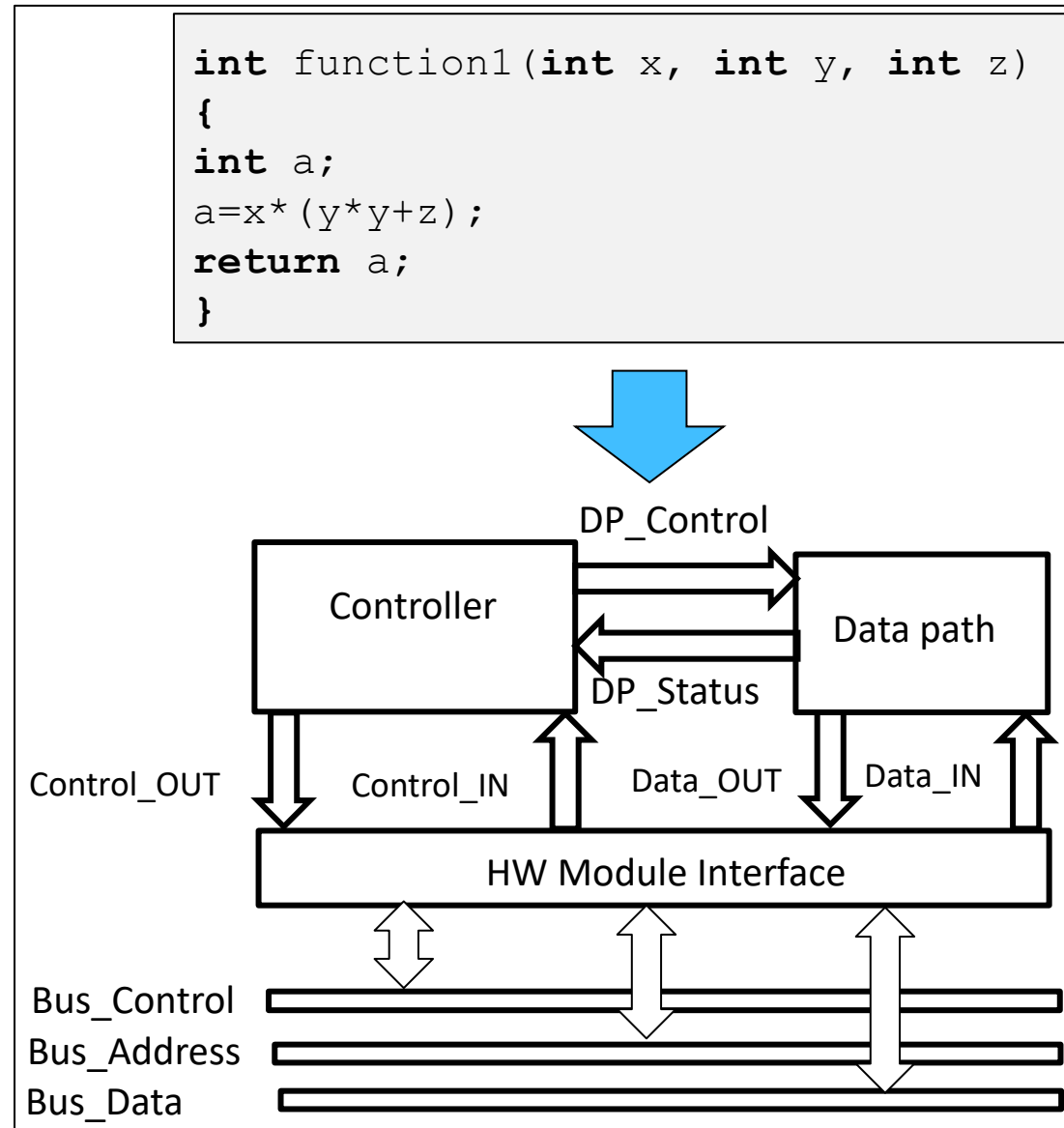Algorithmic description of the task (C, SystemC)

```
int function1(int x, int y, int z)
{
int a;
a=x*(y*y+z);
return a;
}
```

High-level HW Synthesis

RT model of Hardware module in VHDL or Verilog

This is called usually an HW accelerator or IP block

DP_Control

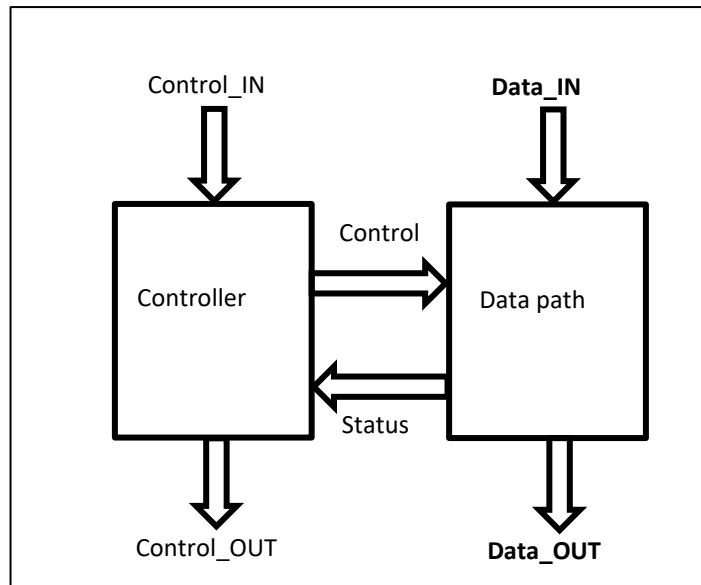Controller          Data path

DP_Status

Control_OUT     Control_IN        Data_OUT     Data_IN

HW Module Interface

Bus_Control
Bus_Address
Bus_Data

Many names:

- High-level Hardware synthesis
- **High-level synthesis (HLS)**
- Algorithmic synthesis
- Behavioral synthesis
- C Synthesis
- …

ACA

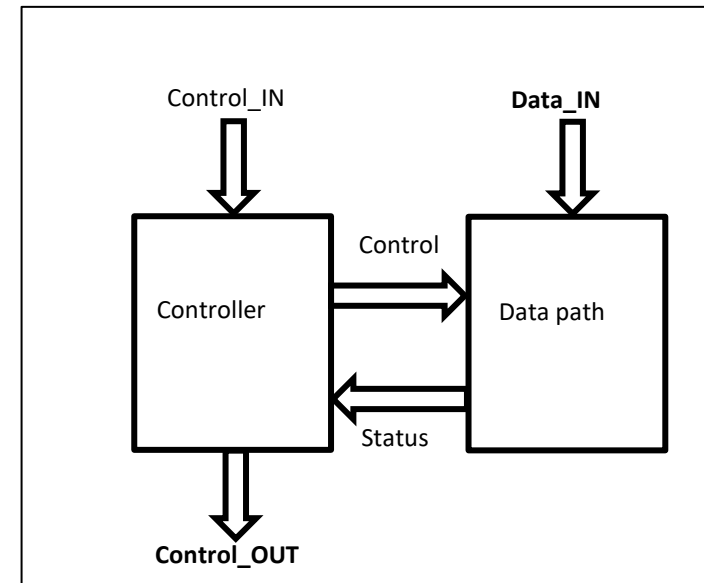# Classes of Hardware Components

- ## Data-oriented designs

HLS works better on data-oriented designs



Examples: Video signal processing, compression, encryption,…

- ## Control-oriented designs



Examples: Traffic light control, industrial machines control, …

- **Clock Cycle Time** $\Delta T$
  - Cycle duration of the driving clock of the HW module
  - The combinatorial path in the circuit with largest delay places an lower limit on the clock cycle time (critical path).

- **Latency** $\Lambda$
  - Number of clock cycles between the start of processing a block of data and the point of time at which the result is ready at the output.

- **Processing time**

$$t_{exe} = \Lambda \cdot \Delta T$$

- **Throughput** $T$
  - Number of blocks of data that can be processed in a fixed time.

ACA

- **Chip area (ASIC)**
  - Estimated via gate count.
  - Data path: Number of Hardware Operation Units such as multipliers, ALUs, registers, multiplexers,…

- **FPGA Resources**
  - Number of Luts, Number of DSP Blocks, …
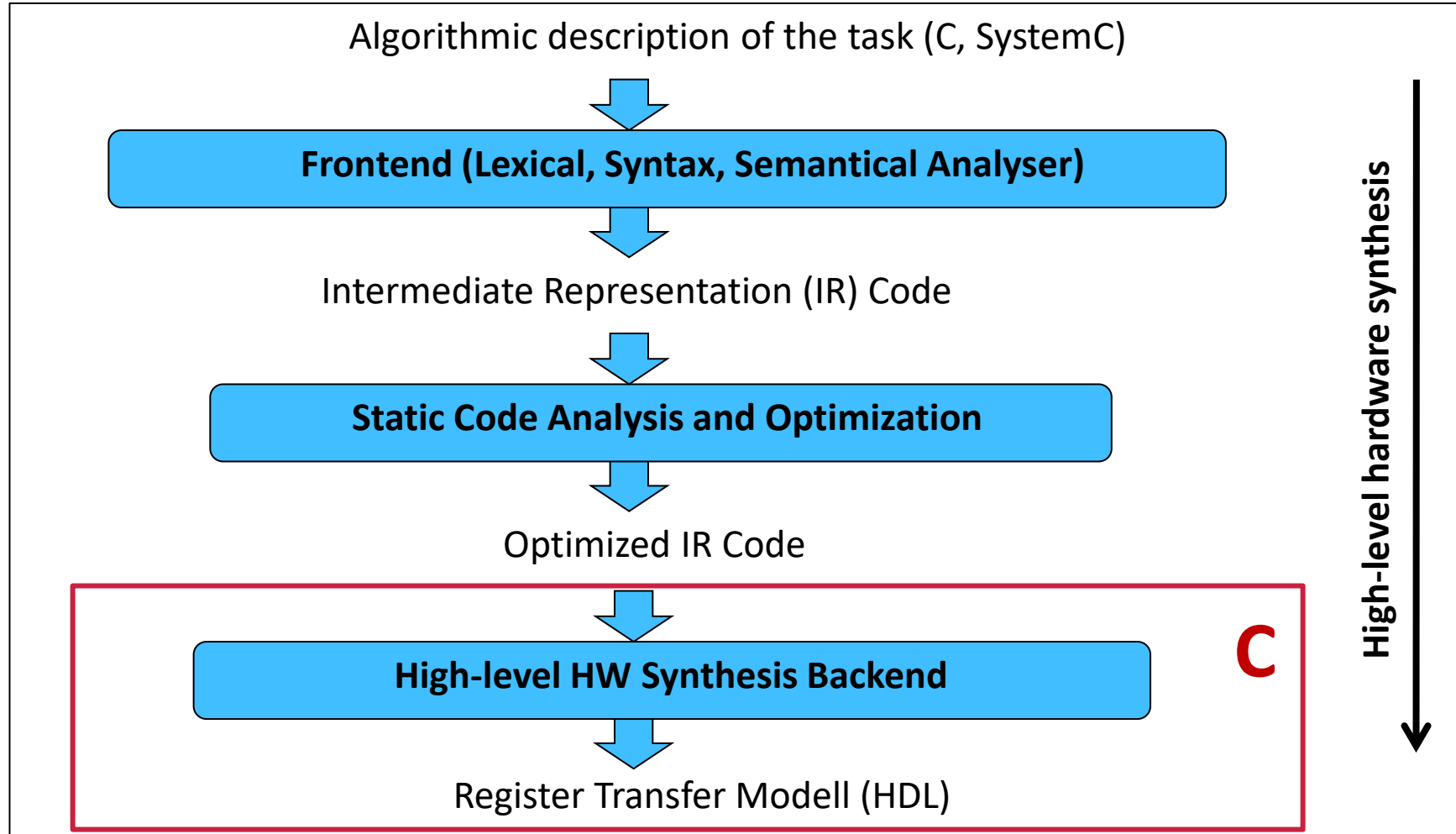
- **Power/Energy Consumption**
  - Dynamic power consumption: Power consumed by switching transistors in the circuit
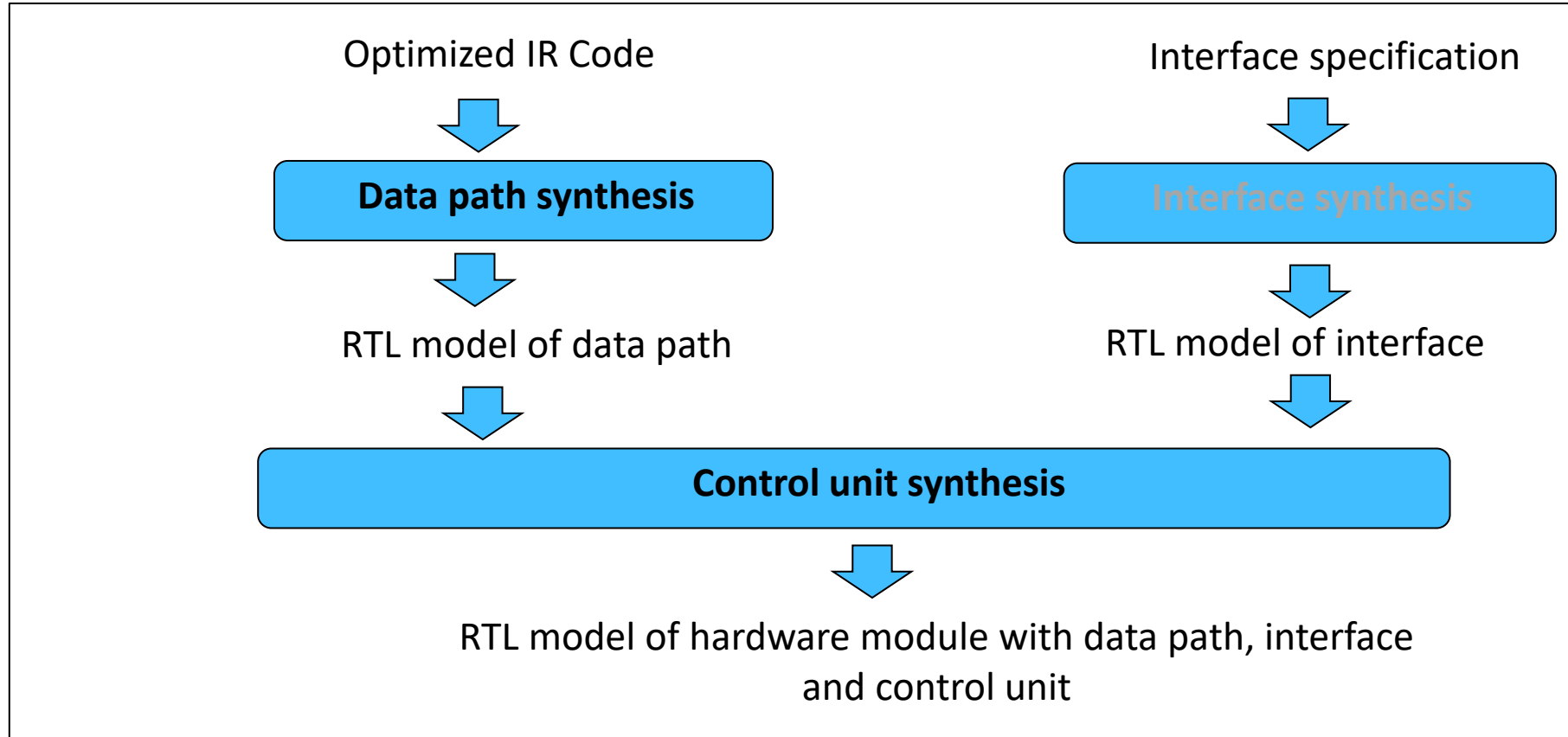  - Static power consumption: Power consumed due to leakage currents.

ACA

# Design Goals and Constraints

- Synthesis algorithms handle two typical cases:

- **Timing constrained:**
  - Constrained: Implement task such that it can compute result in maximal number of clock cycles (maximal latency).
  - Goal: Minimize number of functional units (adders, ALUs, multipliers) in data path.
  - Second goal: Minimize number of registers (register sharing), multiplexers, control unit states, …

- **Resource constrained:**
  - Constrained: Implement task with fixed maximal number of functional units (adders, ALUs, multipliers) in data path.
  - Goal: Minimize latency.
  - Second goal: Minimize number of registers (register sharing), multiplexers, control unit states, …
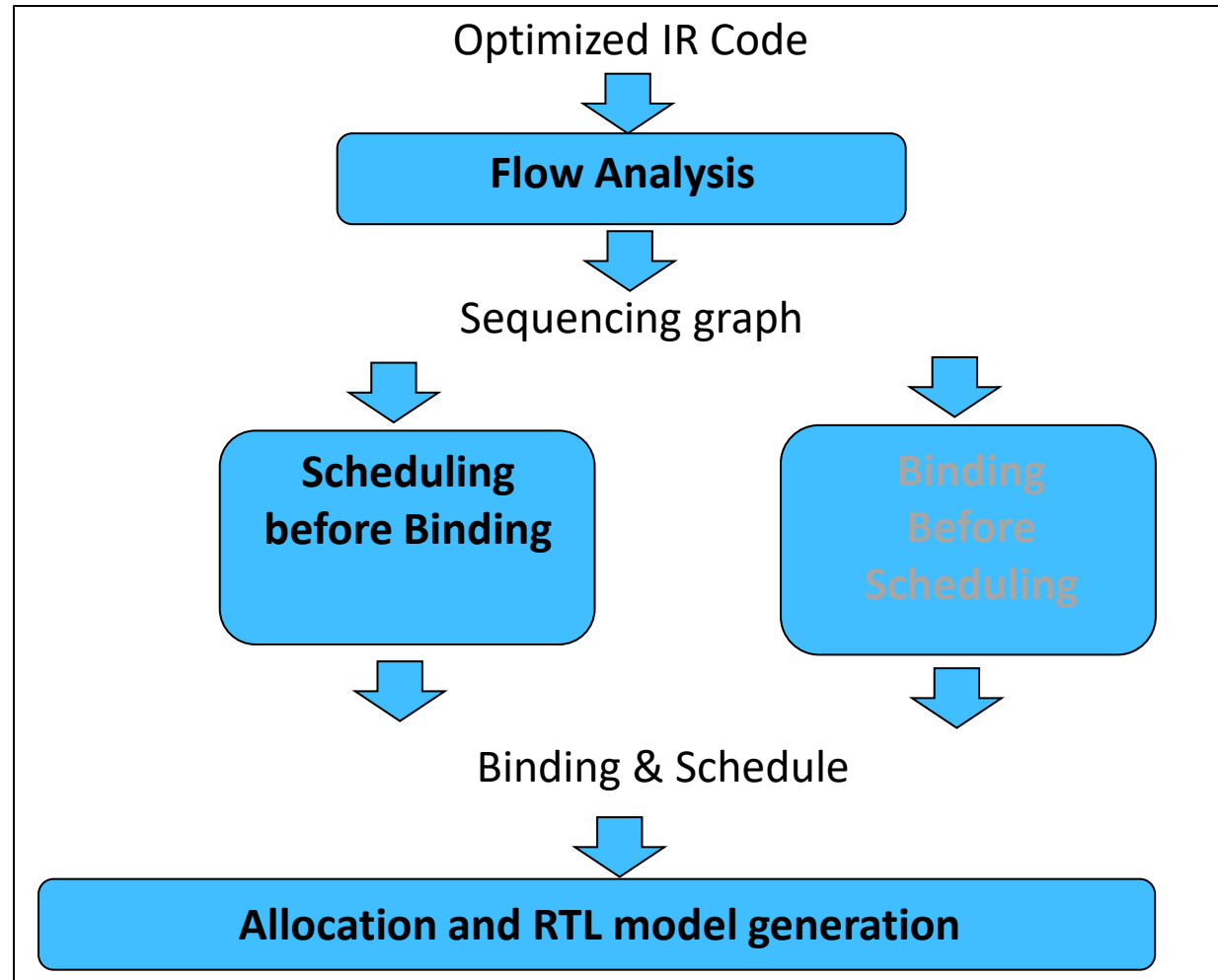
ACA

# Synchronous HW Design

- All registers in control unit and data path share same clock.

- Assumptions for simplification:
  - Functional units have a fixed and known delay such that the number of clock cycles to execute operation is assumed to be fixed and data-independent.
  - The delay of interconnect and multiplexers can be neglected.

- Real life:
  - Longest combinatorial path in the circuit will determine the maximal clock frequency.
  - Logic synthesis will try to optimize circuit dependent on target clock frequency and area.
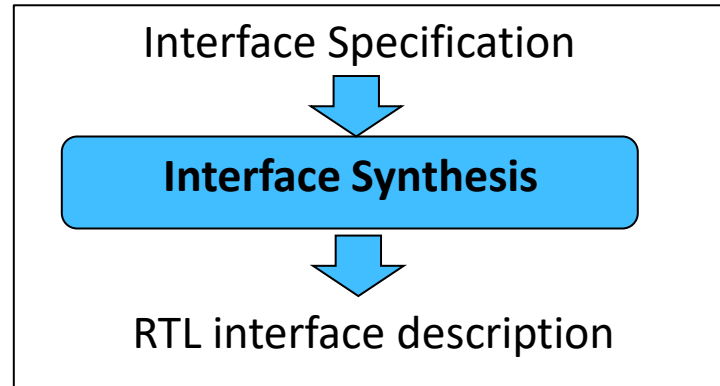
ACA

# High-level HW Synthesis Flow



Algorithmic description of the task (C, SystemC)

**Frontend (Lexical, Syntax, Semantical Analyser)**

Intermediate Representation (IR) Code

**Static Code Analysis and Optimization**

Optimized IR Code

**High-level HW Synthesis Backend**

**C**

Register Transfer Modell (HDL)

High-level hardware synthesis

# High-level HW Synthesis Backend

ACA

ACA

# Data Path Synthesis Steps

- **Scheduling**:
  - Determines the start time of each operation

- **Binding:**
  - Determine on which functional units the operation is executed.
  - Determine in which register variables are saved.

- **Allocation**:
  - Selection of resources such as functional units, registers and multiplexers.

# Interface Synthesis

Interface Specification

⬇

**Interface Synthesis**

⬇

RTL interface description

- Interfaces can differ strongly.
- Interface may consist of:
  - Memory, Registers or FIFOS as data buffers.
  - FSMs for communication (bus) protocols.
- Crossing of clock domains possible, e.g. between bus clock and HW module clock.

ACA

# Control Unit Synthesis

ACA

# D1-3 Data Path Synthesis
# HW Resources

- Literature:
- *„Specification and Design of Embedded Systems"* Daniel D. Gajski, Prentice Hall 1994
- *„Digitale Hardware/Software Systeme"*, Jürgen Teich, Springer 1997
- *„Embedded System Design"*, Daniel D. Gajski et.al., Springer 2009

- Functional units: Adder, Multiplier, ALUs,…
  - Execute operations on data, e.g., Add, Shift, AND, OR, Mult, …
  - Fixed and known delay
  - Fixed and known area demand

- Signal nets and multiplexers
  - Delay and area demand is neglected.

- Memory elements: Registers
  - Delay and area demand is neglected.

- NFU (Nonfunctional Unit):
  - Nonexistent  helper resource
  - used to execute special NOP, LOOP, BRANCH, CALL operations (more later)

ACA

- Functional units are identified by a pair

$$(k_r, z_r)$$

  - of their type:

$$k_r \in K$$

$$K = \{\text{ALU}, \text{MULT}, ...\}$$

  - and an index

$$z_r = 1, 2, ...$$

  - Example:

$$(ALU, 1), (ALU, 2), (MULT, 1)$$

- X-axis: Resources
  - List allocated operational units
  - Assign operations to operational units (Binding)

- y-axis: Time
  - Division in clock cycles.
  - Plan temporal order of the operations
  - Select start times of operations (Schedule)
  - Values must be saved in registers between clock cycles.

- Example: Goertzel Algorithm

  (Basic block B3)

```
t6= s_prev1 * s_prev1
t7= s_prev2 * s_prev2
t8= s_prev1 * s_prev2
t9= t8 * coeff
t10= t6+t7
power= t10 – t9
```
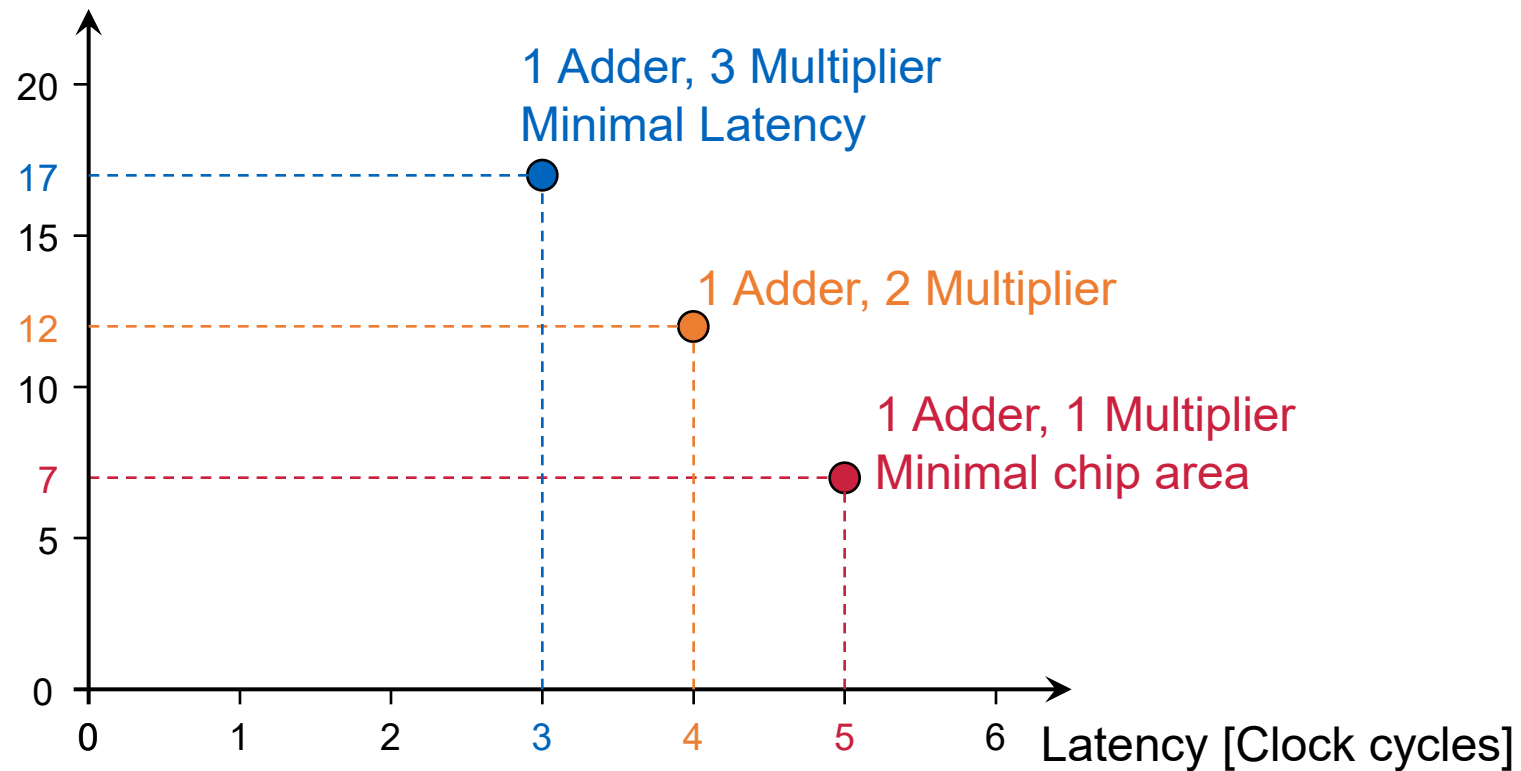
Resources (Functional units)

| | Add,1 | Mult,1 | Mult,2 |
|---|---|---|---|
| CC 1 | | t6= s_prev1* s_prev1 | t8= s_prev1* s_prev2 |
| CC 2 | | t7= s_prev2* s_prev2 | t9= t8* coeff |
| CC 3 | t10= t6+t7 | | |
| CC 4 | power= t10–t9 | | |

Time in clock cycles (CC)

ACA

- Example: Goertzel Algorithm

  (Basic block B3)

```
t6= s_prev1 * s_prev1
t7= s_prev2 * s_prev2
t8= s_prev1 * s_prev2
t9= t8 * coeff
t10= t6+t7
power= t10 – t9
```

| | Add,1 | Mult,1 |
|---|---|---|
| CC 1 | | t6= s_prev1* s_prev1 |
| CC 2 | | t7= s_prev2* s_prev2 |
| CC 3 | t10= t6+t7 | t8= s_prev1* s_prev2 |
| CC 4 | | t9= t8* coeff |
| CC 5 | power= t10–t9 | |

ACA

- Example: Goertzel Algorithm
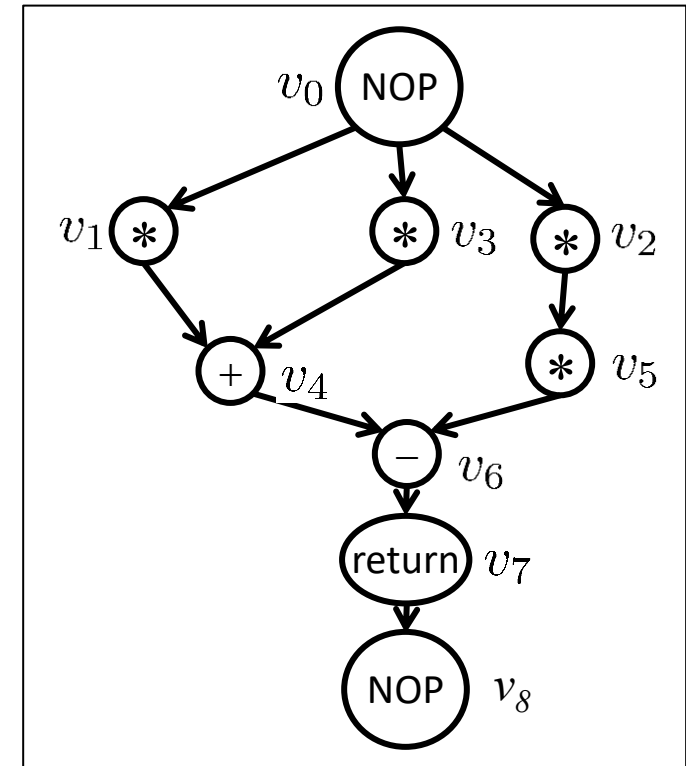
  (Basic block B3)

```
t6= s_prev1 * s_prev1
t7= s_prev2 * s_prev2
t8= s_prev1 * s_prev2
t9= t8 * coeff
t10= t6+t7
power= t10 – t9
```

| | Add,1 | Mult,1 | Mult,2 | Mult,3 |
|------|-------|--------|--------|--------|
| CC 1 | | t6= s_prev1* s_prev1 | t7= s_prev2* s_prev2 | t8= s_prev1* s_prev2 |
| CC 2 | t10= t6+t7 | t9= t8* coeff | | |
| CC 3 | power= t10–t9 | | | |
| | | | | |

ACA

- Which is the best solution?

- Solution is Pareto optimal, if there exist no solution that is better in all design performance metrics.

- Different Pareto-optimal solutions allow different trade-offs between the design performance metrics.

- Best solution is picked based on preferences on design performance metrics.

ACA

Chip Area [units]
(Demand: Adder=2 area units, Multiplier=5 area units)

1 Adder, 3 Multiplier
Minimal Latency

1 Adder, 2 Multiplier

1 Adder, 1 Multiplier
Minimal chip area

Latency [Clock cycles]

ACA

# D1-4 Sequencing Graphs

# Sequencing Graph (SG)

- Sequencing graph: $G_s(V_s, E_s)$

- Hierarchy of directed acyclic graphs (DAGs). Each DAG is called a sequencing graph unit (SGU).
- SGUs are polar: One source and one sink node is added, which is labeled No operation (NOP).

- Nodes: $V_s = \{v_i : i = 0, ..., n\}$
  - No operation (NOP)
  - Operations (+,>,<,*,…)
  - Hierarchical node (CALL, BR, LOOP)

- Edges: $E_s = \{(v_i, v_j) : i, j = 0, ..., n\}$
  - between nodes in one SGU unit: Data dependency between two operations
  - between Source/sink and hierarchical nodes: Connections between SGU on different hierarchical levels
- Paths describe concurrent operations that may possibly be executed in parallel,

ACA

- Example: SGU for basic block B3 of Goertzel algorithm

- Example: Goertzel Algorithm (Basic block B3)

```
t6= s_prev1 * s_prev1
t7= s_prev2 * s_prev2
t8= s_prev1 * s_prev2
t9= t8 * coeff
t10= t6+t7
power= t10 – t9
```

Data flow graph



Sequencing graph unit

ACA

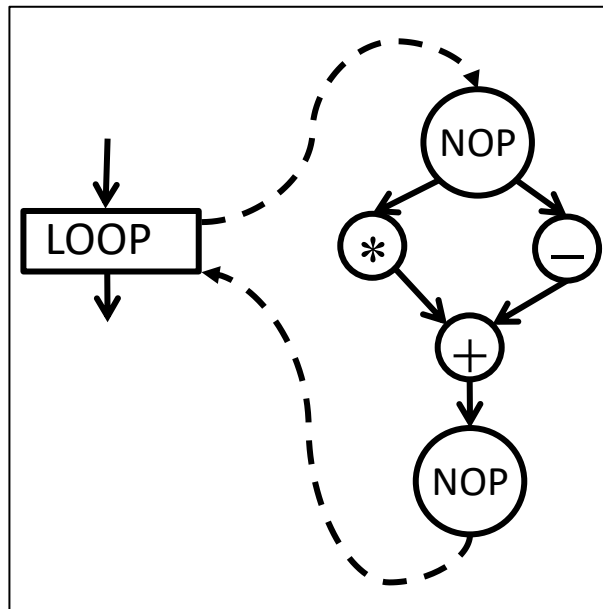- Hierarchical nodes: CALL, LOOP, BR

Call to procedure

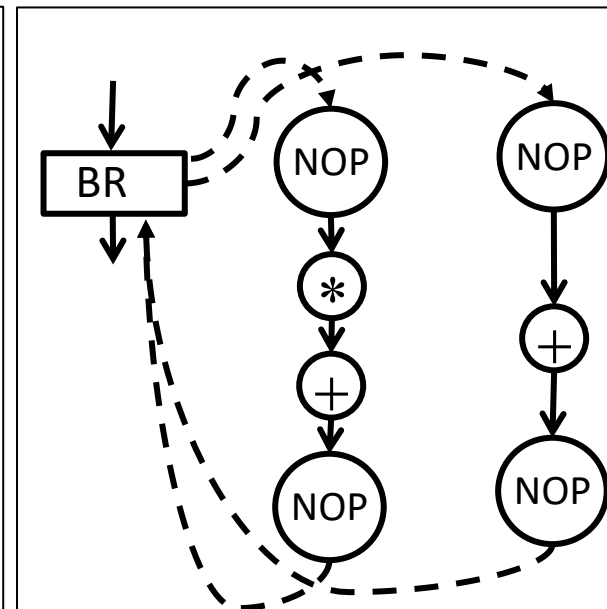Control flow loop

Control flow branch



Called SGU of one lower hierarchical level is executed once.
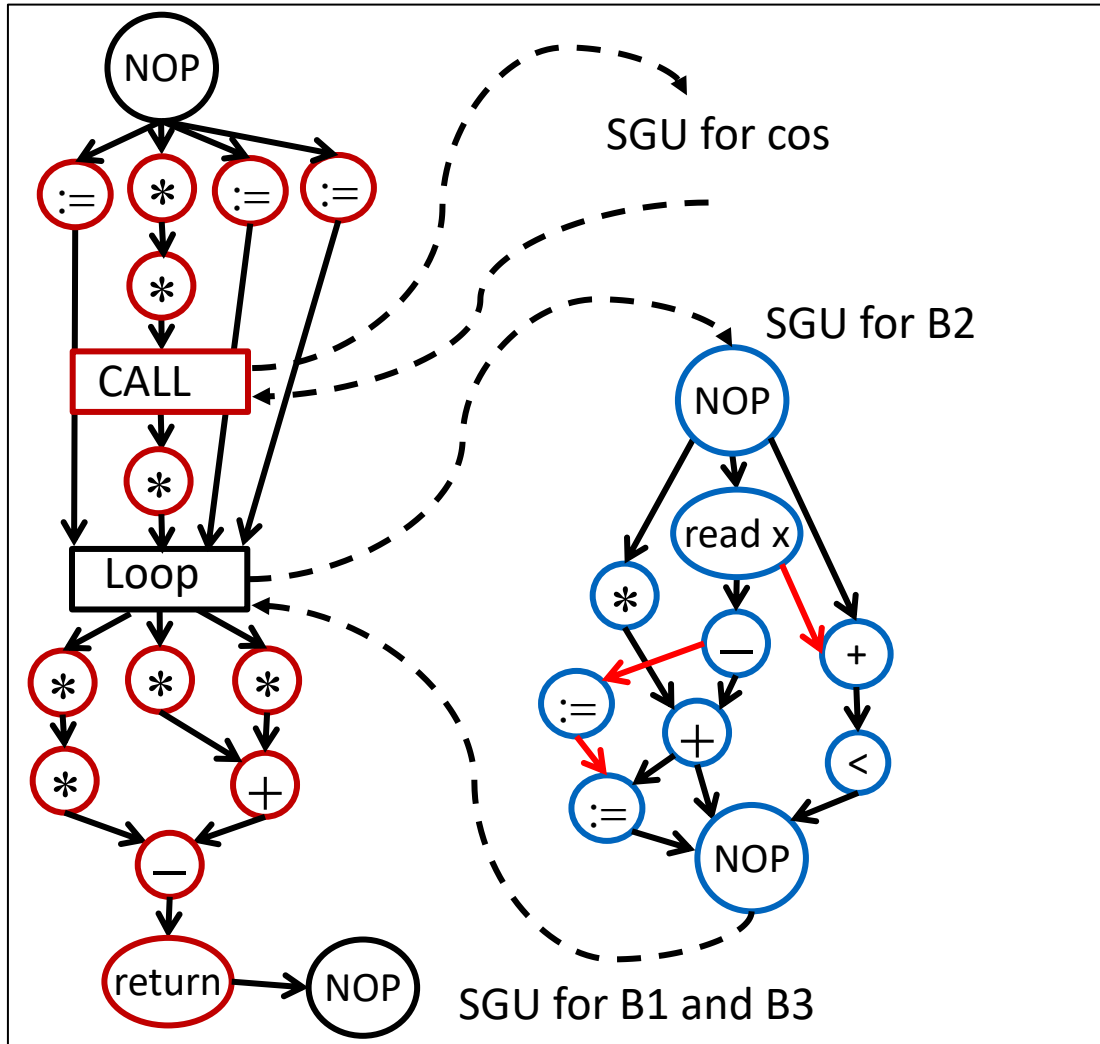
SGU of lower hierarchical level is executed 0 to N times.

Only one of the two SGU of lower hierarchical level is executed once.

ACA

- Example: Goertzel algorithm



SGU for cos

SGU for B2

SGU for B1 and B3

```
B1: s_prev1 := 0.0
    s_prev2 := 0.0
    i:=0
    t1 := 2*3.14
    f := t1 * freq
    param f
    t2 := call cos,1
    coeff:=2.0*t2
```

```
B2: t3:= coeff * s_prev1
    t4:= x[i]
    t5 := t4 - s_prev2
    s := t3 + t5
    s_prev2 := s_prev1
    s_prev1 := s
    i:=i+1
    if i < 64 goto B2
```

```
B3: t6:= s_prev1 * s_prev1
    t7:= s_prev2 * s_prev2
    t8:= s_prev1 * s_prev2
    t9:= t8 * coeff
    t10:= t6+t7
    power:= t10 - t9
    return power
```

ACA

- Example: Goertzel Algorithm

  (Basic block B3)

```
t6= s_prev1 * s_prev1
t7= s_prev2 * s_prev2
t8= s_prev1 * s_prev2
t9= t8 * coeff
t10= t6+t7
power= t10 - t9
return power
```
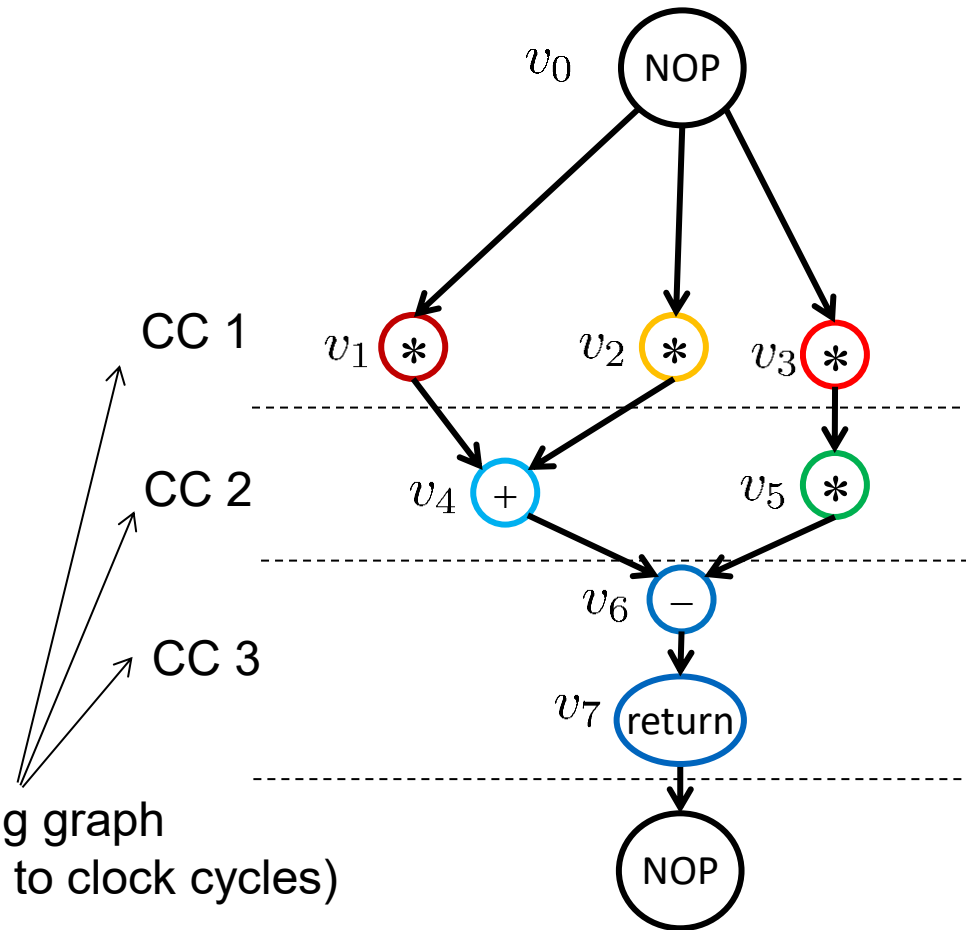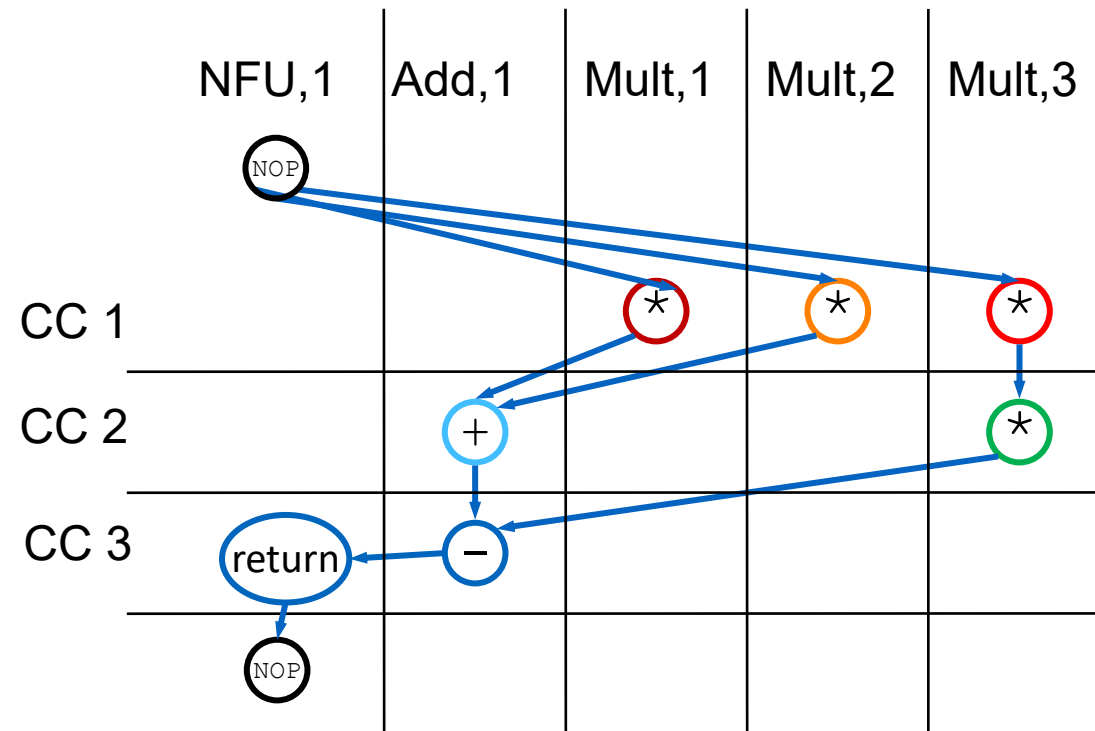
ACA

- Example: Goertzel Algorithm

(Basic block B3)

```
t6= s_prev1 * s_prev1
t7= s_prev2 * s_prev2
t8= s_prev1 * s_prev2
t9= t8 * coeff
t10= t6+t7
power= t10 - t9
return power
```

CC 1

CC 2

CC 3

$v_0$ NOP

$v_1$ *

$v_2$ *

$v_3$ *

$v_4$ +

$v_5$ *

$v_6$ −

$v_7$ return

NOP

Scheduled sequencing graph
(Operations assigned to clock cycles)

ACA
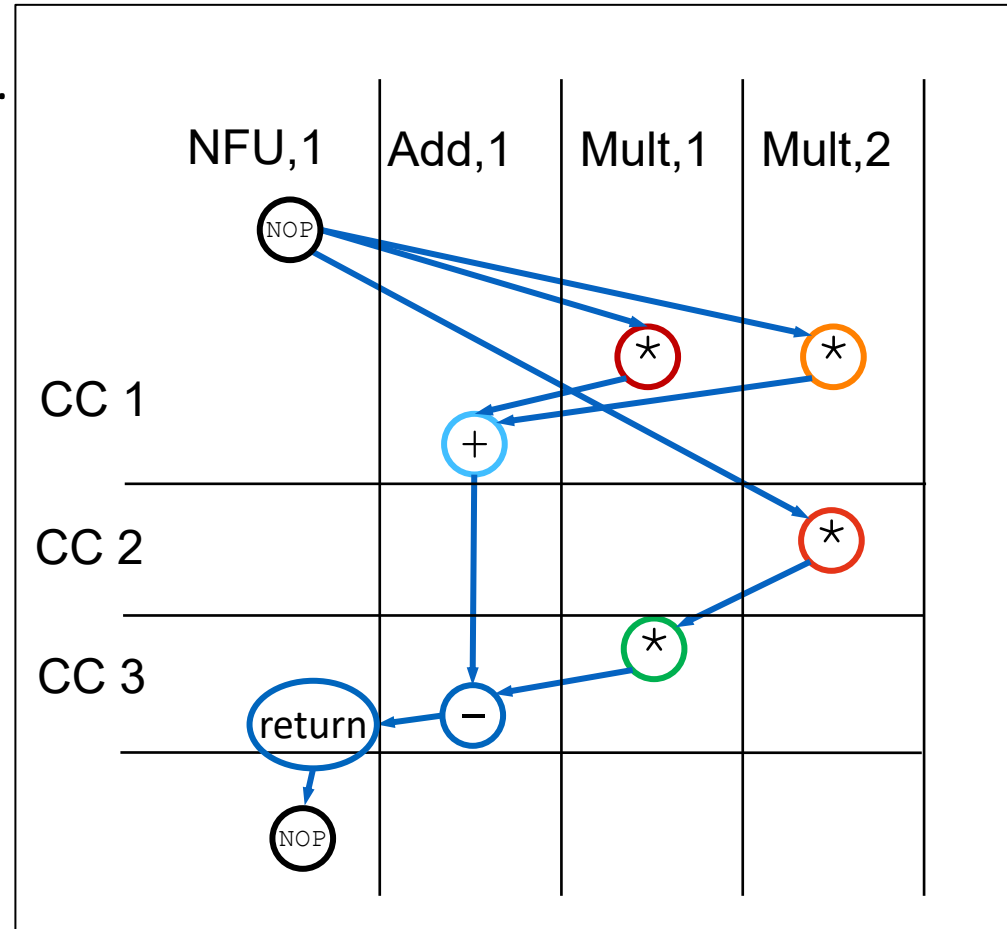
- Example: Goertzel Algorithm

  (Basic block B3)

```
t6= s_prev1 * s_prev1
t7= s_prev2 * s_prev2
t8= s_prev1 * s_prev2
t9= t8 * coeff
t10= t6+t7
power= t10 – t9
return power
```



Scheduled sequencing graph with binding
(Operations assigned to clock cycles and operational units)
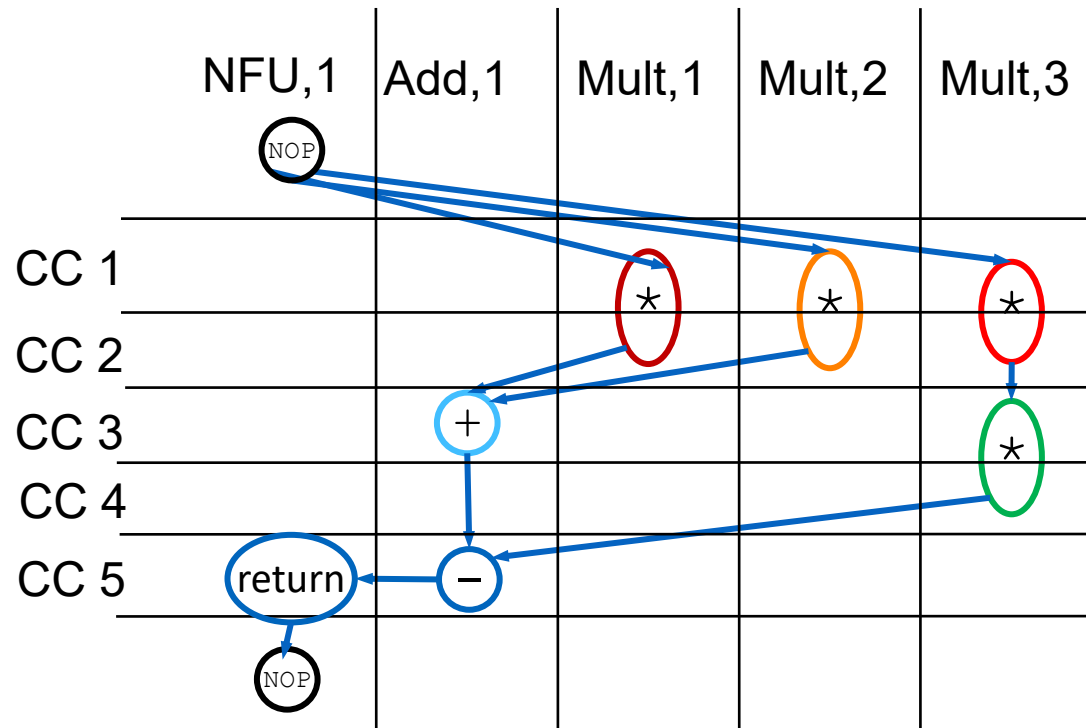
ACA

# Operation Chaining

- Delay of operational units allow two operations in one clock cycle.

- Example: Goertzel Algorithm
  - Addition and multiplication in same clock cycle possible.
  - Operational units must be switched in series (chained).
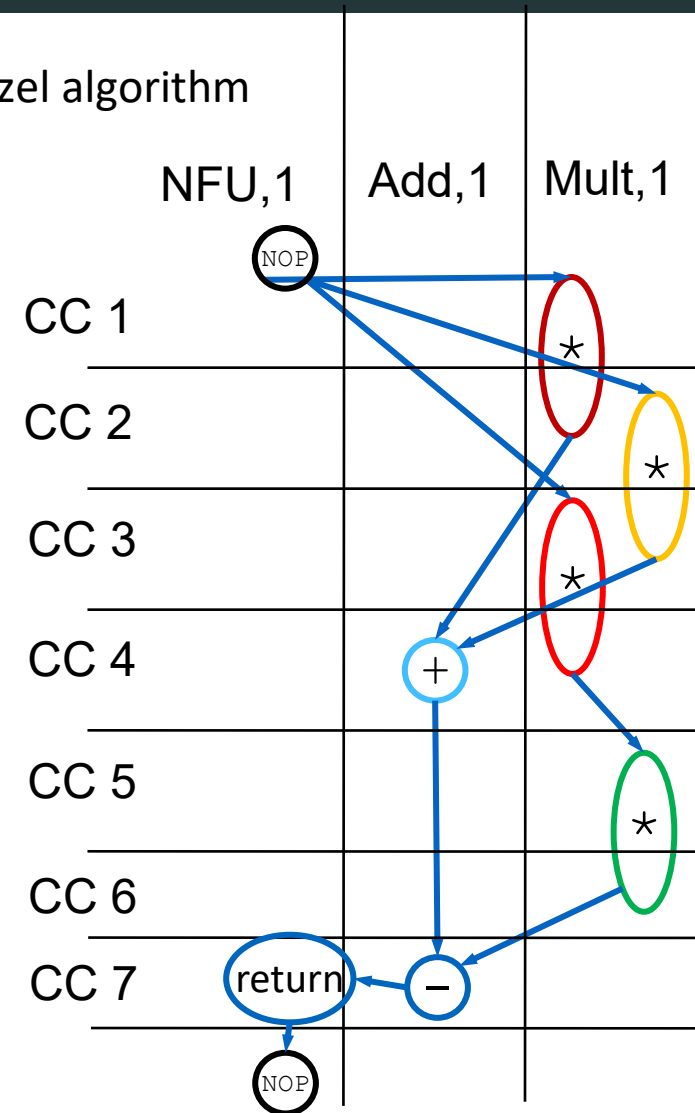
# Multi-cycle Operations

- Delay of functional elements requires several clock cycles for the execution of the operation.

- Example:
  Goertzel algorithm

# Pipelined Operational Units

- Example: Goertzel algorithm

- New operation can start before previous operation has finished.

- Number of concurrent operations is equal to pipeline depth.

- Operational units has internal registers to save intermediate values.

# Summary

- HLS is a step that uses C specification to design an HW IP block
- This HW IP block executes the algorithm specified in C