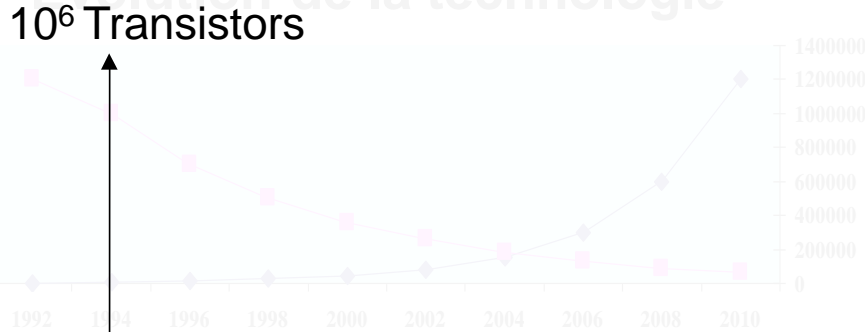# *Une introduction à la synthèse de haut-niveau*

## *(ou comment générer des architectures matérielles à partir du langage C)*

**Université de Bretagne-Sud
Lab-STICC**

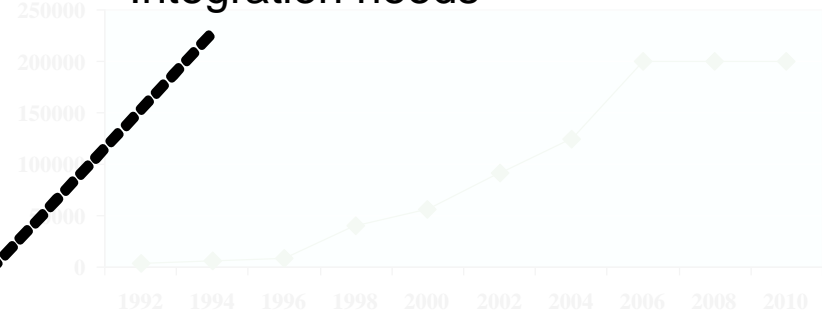**Philippe COUSSY**
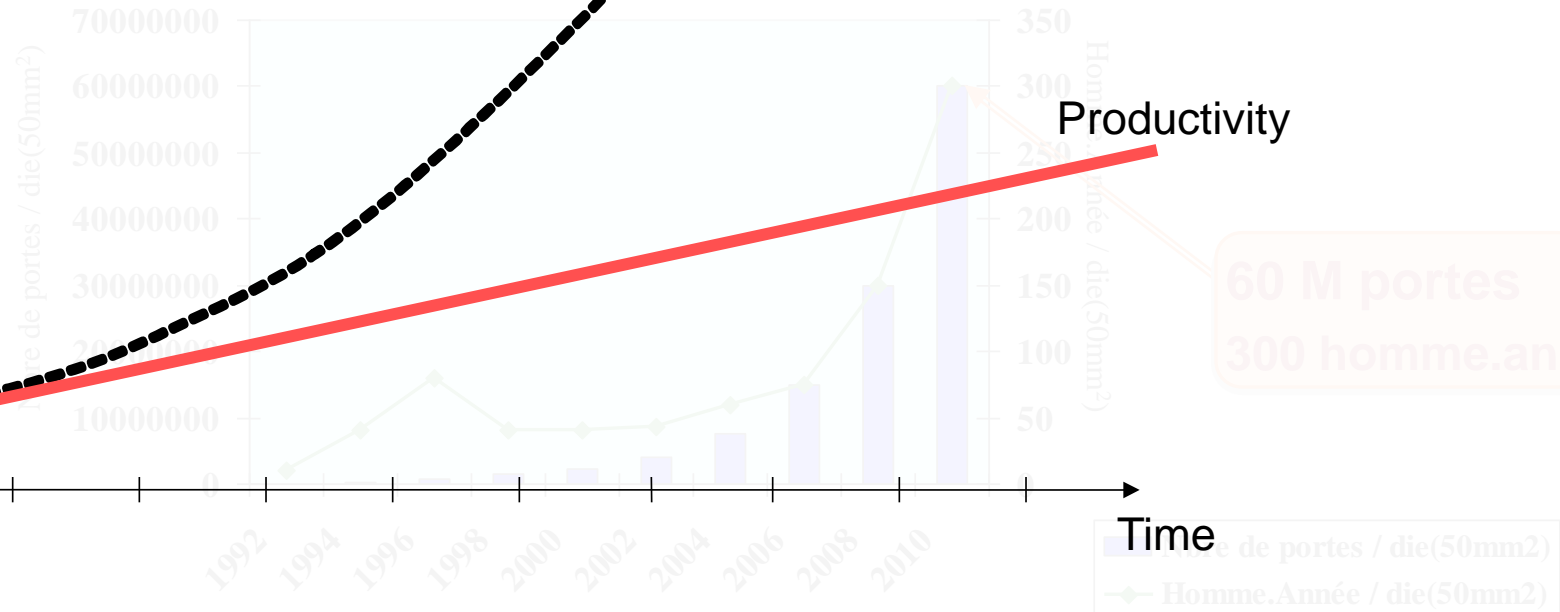**philippe.coussy@univ-ubs.fr**

# Productivity gap



$10^6$ Transistors
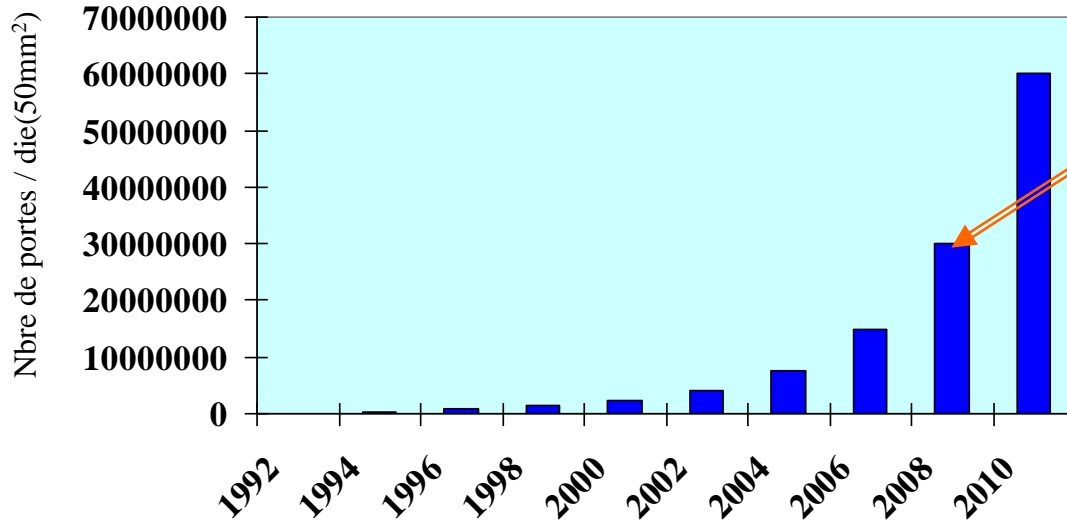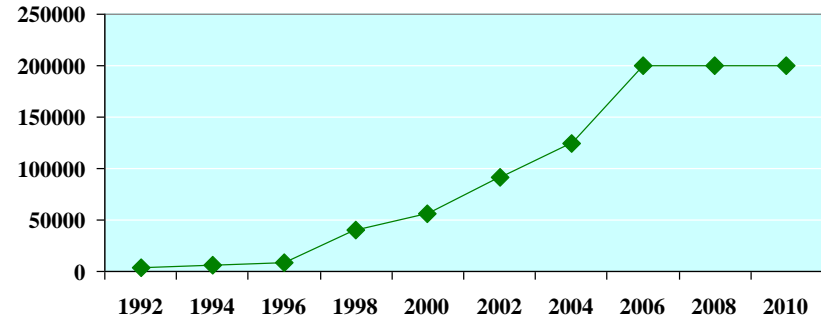
Integration needs

Productivity

Time

60 M portes
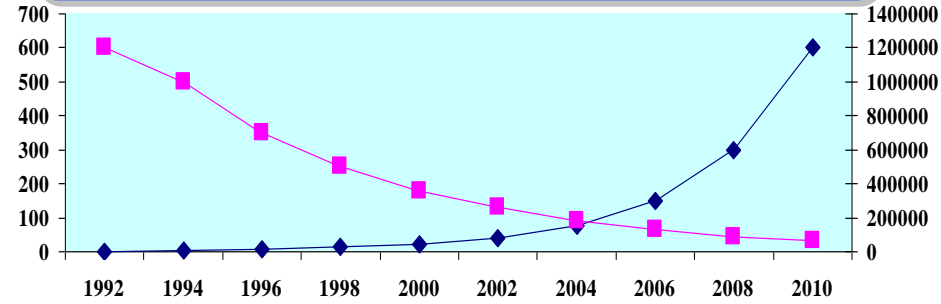300 homme.an

# Context



- Decreasing technology scale
- Increasing transistor density

- Designer productivity

- 30 M gates

Nbre de portes / die(50mm$^2$)
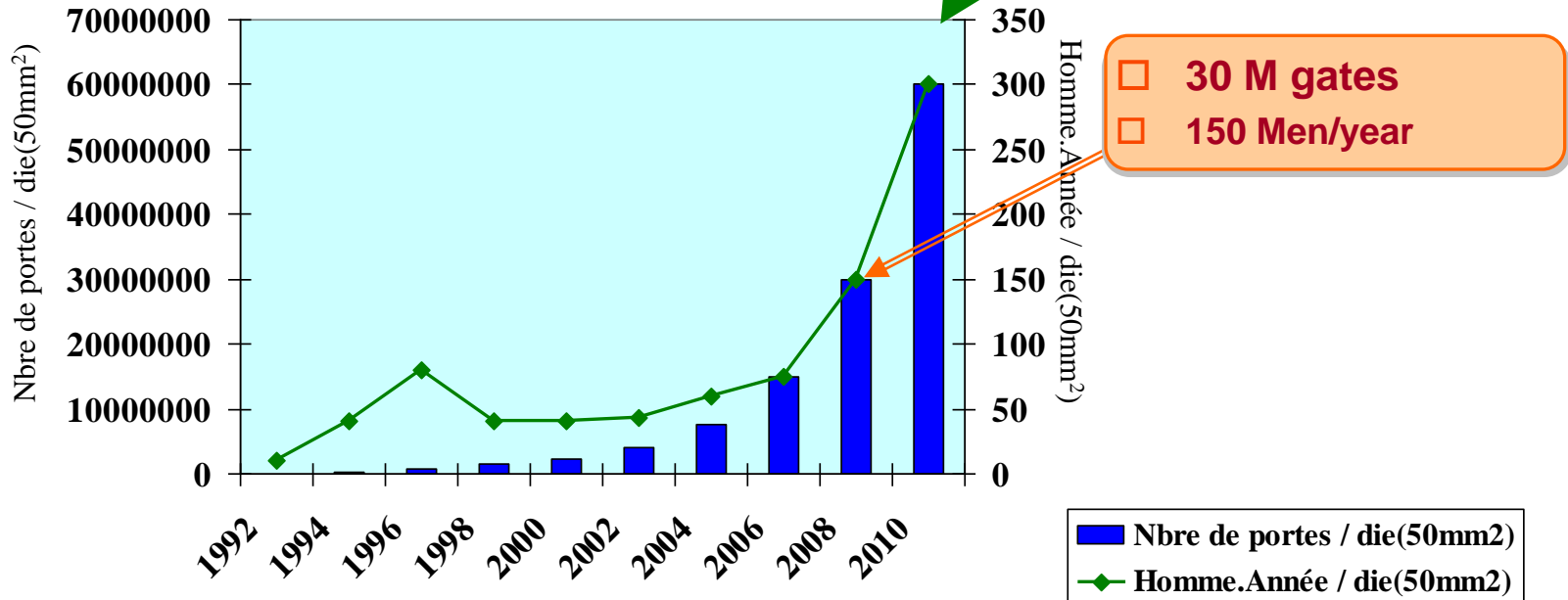
- Nbre de portes / die(50mm2)

# Context

- **Decreasing technology scale**
- **Increasing transistor density**

- **Designer productivity**

- **30 M gates**
- **150 Men/year**

Nbre de portes / die(50mm²)

Homme.Année / die(50mm²)

# Design methodologies

☐ **Synthesis and verification automation has always been key factors in the evolution of the design process**

  ■ Allow to explore the design space efficiently and rapidly
  ■ Correct by construction design

# Design methodologies

## ☐ Software domain

- Machine code (binary sequence)
- 1950s: concept of assembly language (and assembler)
  - ☐ *based on mnemonics*
  - ☐ *Maurice V. Wilkes de l'université de Cambridge*
- Later: High-level languages and compilers
  - ☐ *1951: First compiler*
    - (A-0 system) par Grace Hopper
  - ☐ *Fortran 1954-1957: First high-level language*
    - FORmula TRANslator
  - ☐ *Cobol 1959, Basic 1964, C 1972, C++ 1983…*

## ☐ High-level language

- Platform independent
- Follow the rules of human language
  - ☐ *with a grammar, a syntax and a semantic*
- Provide flexibility and portability
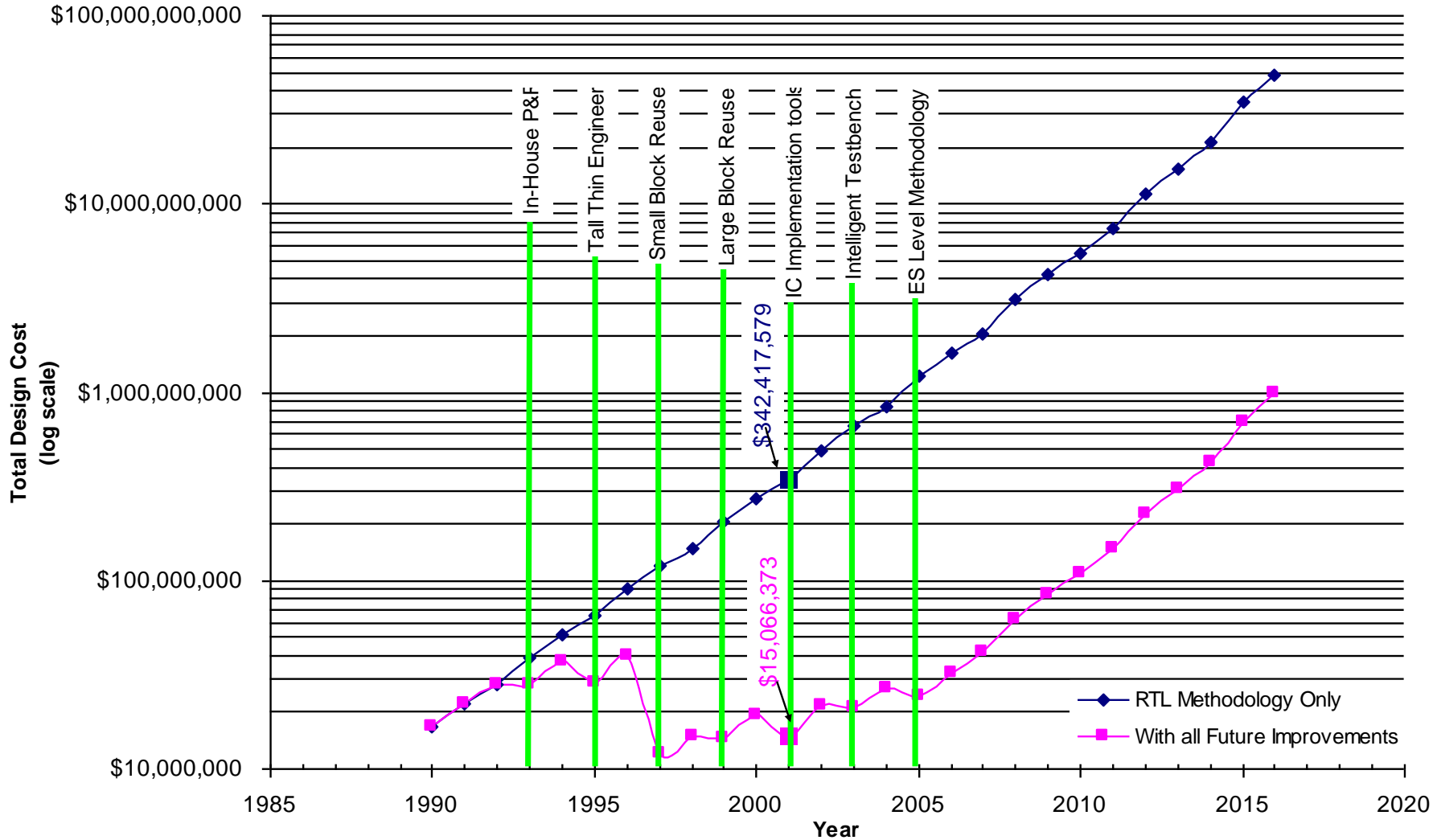  - ☐ *by hiding details of the computer architecture*

# Design methodologies

## ☐ Hardware domain

- ■ 1960: IC were done by hand
  - ❏ *designed, optimized and laid out*
- ■ 1970: Gate-level simulation
- ■ end of 70: Cycle-based simulation
- ■ 1980: Wide automation
  - ❏ *place & route, schematic circuit capture, formal verification and static timing analysis*
- ■ Mid 1980: Hardware description language
  - ❏ *1986 Verilog, 1987 VHDL*
- ■ 1990: logic synthesis
  - ❏ *VHDL and Verilog synthesizable subsets*
- ■ Mid 1990:
  - ❏ *High-level synthesis (First gen),*
  - ❏ *Co-design, IP-core reuse…*
- ■ 2000 : Electronic System Level ESL
  - ❏ *System level language*
    - ▪ SystemC, SystemVerilog…,
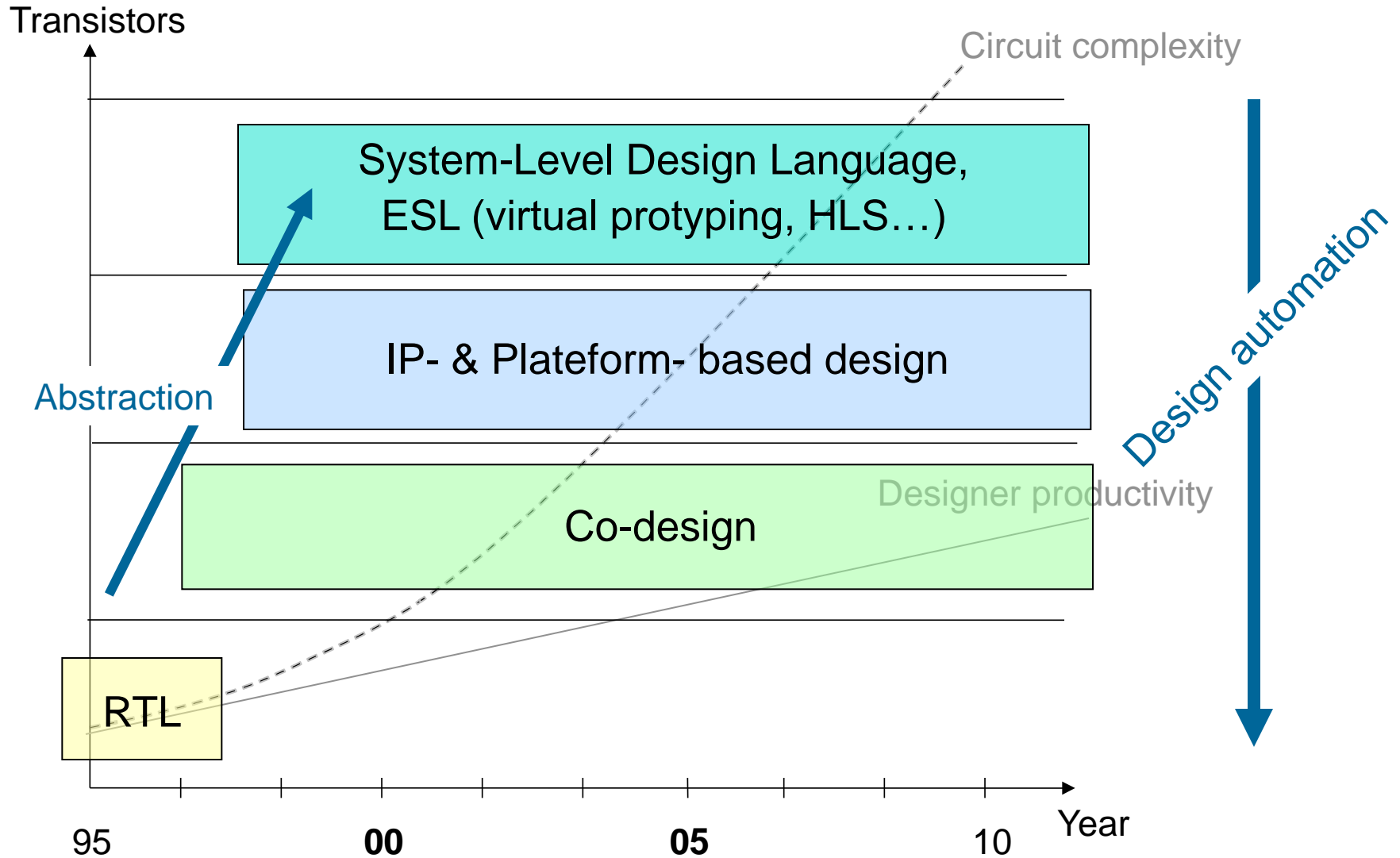    - ▪ Virtual prototyping, Transaction Level Modellin TLM ...
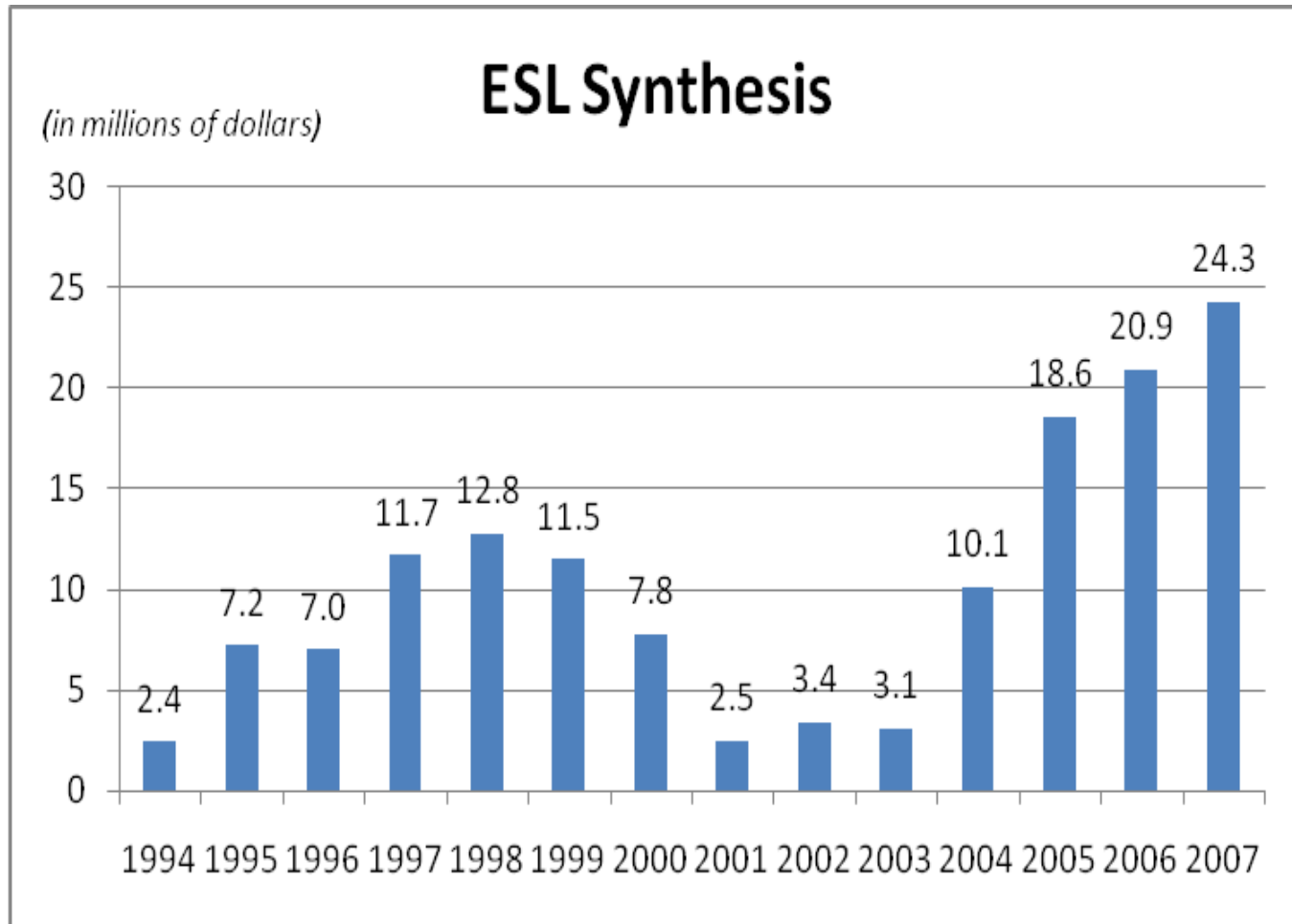
# Design gap



**SOC Design Cost Model**

(A. B. Kahng, G. Smith, "A New Design Cost Model", For the 2001 ITRS)

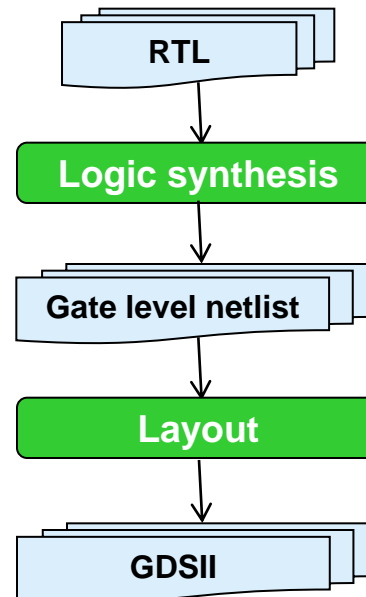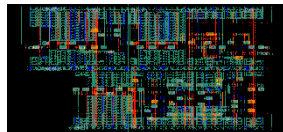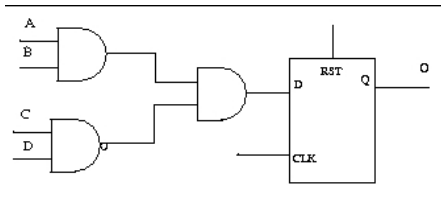# Electronic System Level Design (ESLD)

# ESL Market

# Outline

☐ **Lab-STICC**

☐ **General context**

☐ **High-Level Synthesis**

  ■ Brief introduction

  ■ "In details"

☐ **GAUT**

  ■ Overview

  ■ Results

☐ **Conclusion**

☐ **References**

# Typical HW design flow
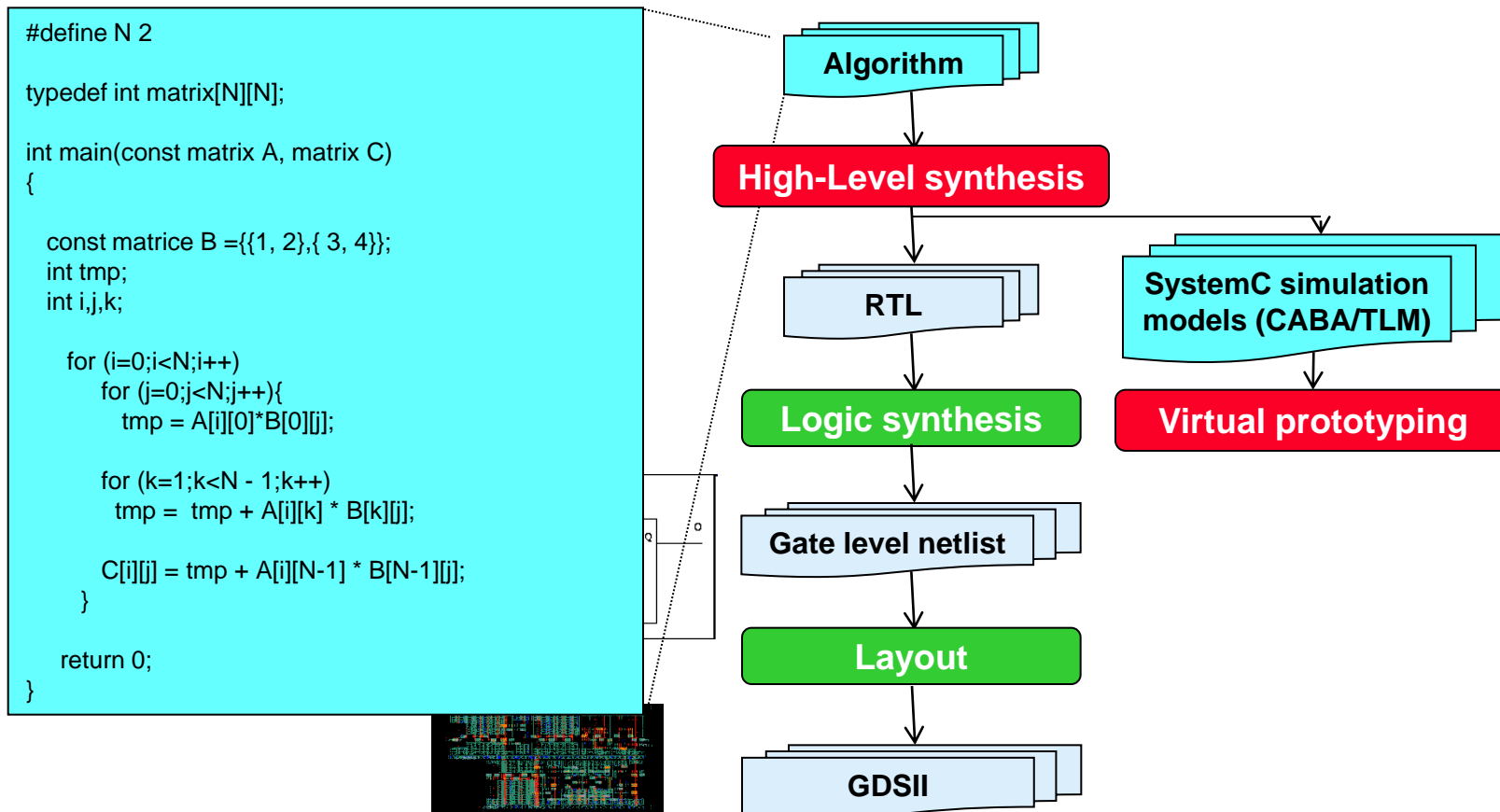
☐ **Starting from a Register Transfer Level description, generate an IC layout**



```
process( CLK, RST)

    if( RST = '1' ) then
        Q <= '0';                    .
    else if rising_edge( CLK) then
        Q <= A and B and C nand D ;
```

**RTL**

↓

**Logic synthesis**

↓

**Gate level netlist**

↓

**Layout**

↓

**GDSII**

# Typical HW design flow

☐ **Starting from a functional description, automatically generate an RTL architecture**

```
#define N 2

typedef int matrix[N][N];

int main(const matrix A, matrix C)
{

  const matrice B ={{1, 2},{ 3, 4}};
  int tmp;
  int i,j,k;

    for (i=0;i<N;i++)
       for (j=0;j<N;j++){
         tmp = A[i][0]*B[0][j];

       for (k=1;k<N - 1;k++)
         tmp =  tmp + A[i][k] * B[k][j];

         C[i][j] = tmp + A[i][N-1] * B[N-1][j];
       }

    return 0;
}
```

**Algorithm**

**High-Level synthesis**

**RTL**

**SystemC simulation models (CABA/TLM)**

**Logic synthesis**

**Virtual prototyping**

**Gate level netlist**

**Layout**

**GDSII**

# High-level synthesis

☐ **Starting from a functional description, automatically generate an RTL architecture**

☐ **Constraints**
- Timing constraints: latency and/or throughput
- Resource constraints: #Operators and/or #Registers and/or #Memory, #Slices...

☐ **Objectives**
- Minimization: area i.e. resources, latency, power consumption…
- Maximization: throughput

# Synthesis steps

☐ **Compilation**
  - ■ Generates a formal modeling of the specification

☐ **Selection**
  - ■ Chooses the architecture of the operators

☐ **Allocation**
  - ■ Defines the number of operators for each selected type

☐ **Scheduling**
  - ■ Defines the execution date of each operation

☐ **Binding (or Assignment)**
  - ■ Defines which operator will execute a given operation
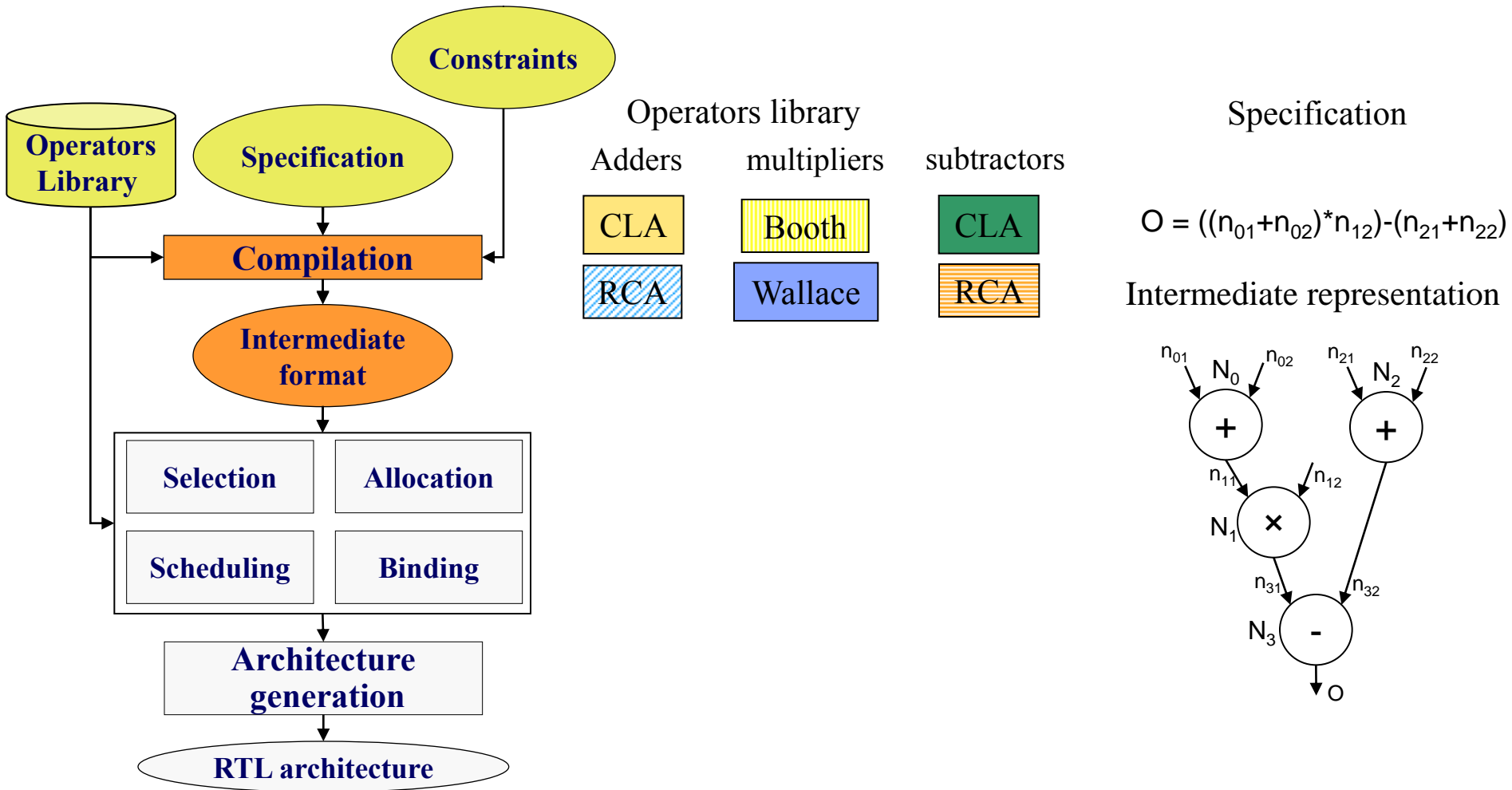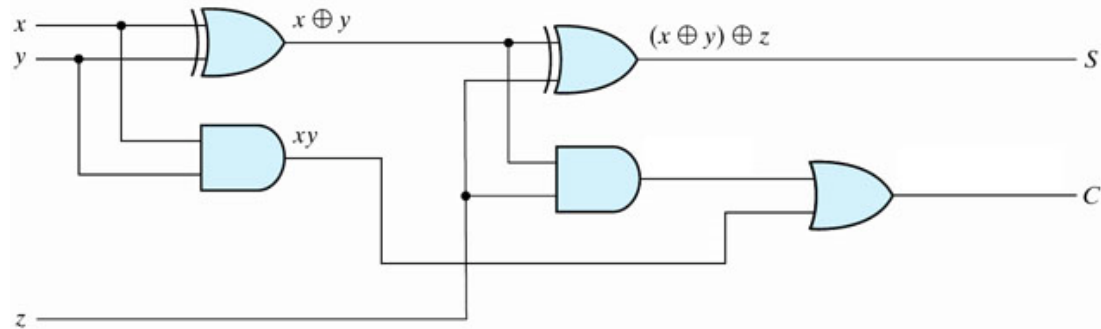  - ■ Defines which memory element will store a data

☐ **Architecture generation**
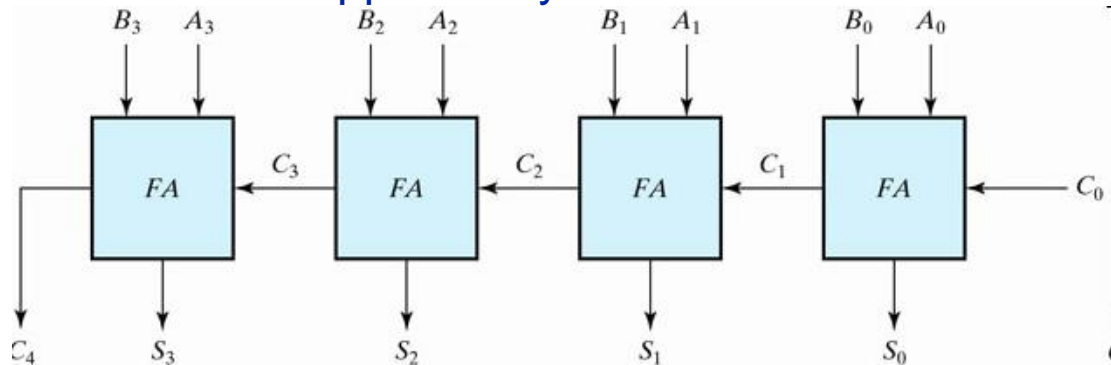  - ■ Writes out the RTL source code in the target language e.g. VHDL

# HLS steps: inputs



Operators library

| Adders | multipliers | subtractors |
|--------|-------------|-------------|
| CLA | Booth | CLA |
| RCA | Wallace | RCA |

Specification

$$O = ((n_{01}+n_{02})*n_{12})-(n_{21}+n_{22})$$

# HLS steps: Compilation

# Synthesis steps

☐ **Compilation**
  ■ Generates a formal modeling of the specification

☐ **Selection**
  ■ Chooses the architecture of the operators

☐ **Allocation**
  ■ Defines the number of operators for each selected type

☐ **Scheduling**
  ■ Defines the execution date of each operation

☐ **Binding (or Assignment)**
  ■ Defines which operator will execute a given operation
  ■ Defines which memory element will store a data

☐ **Architecture generation**
  ■ Writes out the RTL source code in the target language e.g. VHDL

# Operator architecture

## ☐ Full 1-bit adder : X + Y + Z

- ■ X, Y are the operands
- ■ Z is the input carry



## ☐ Ripple Carry Adder

- ■ Add two integers A and B
- ■ Cascade of 1-bit adders => Ripple-Carry Adder

# Operator architecture

## ☐ Carry Look-ahead adder CLA

- Uses a carry generator to compute all the carries concurrently
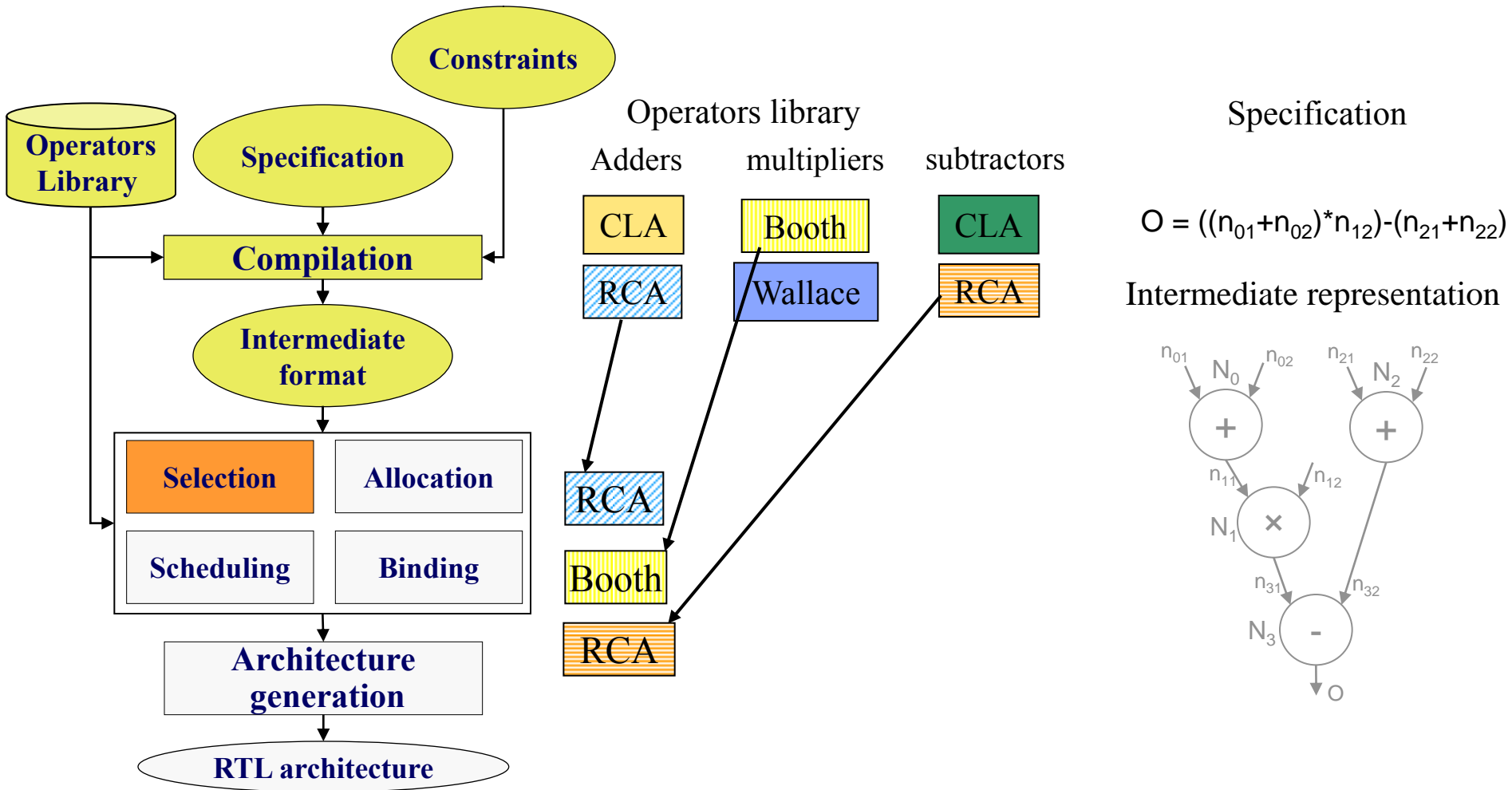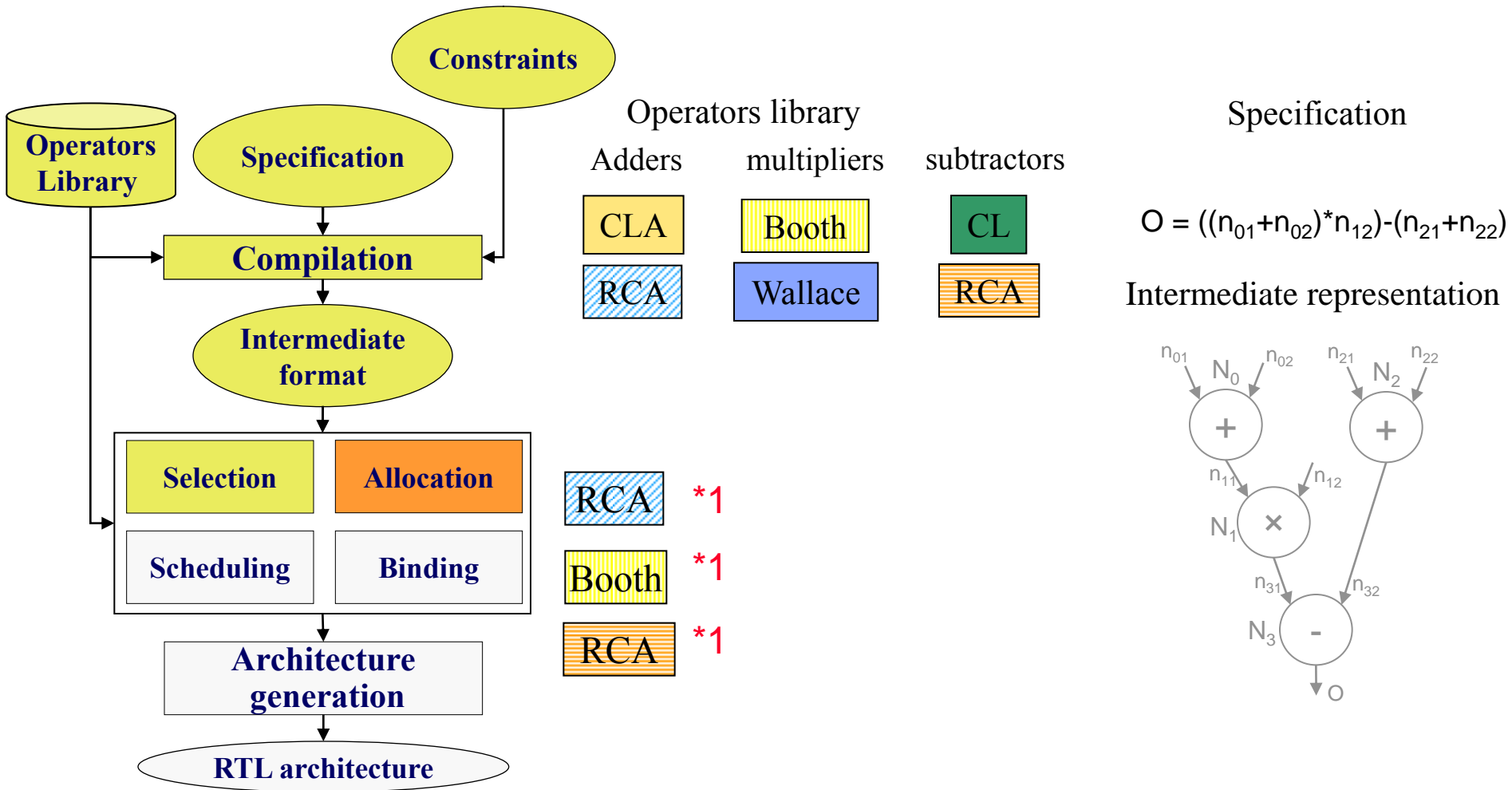  - ☐ *faster but also larger than the RCA*

# Library characterization

- ☐ **RTL architecture produced by HLS depends on the capabilities and characteristics of the operators**
- ☐ **Library processing reads the available libraries and determines the functional, timing, and area characteristics of the available parts.**

# HLS steps: Selection

# Synthesis steps

☐ **Compilation**
  - ■ Generates a formal modeling of the specification

☐ **Selection**
  - ■ Chooses the architecture of the operators

☐ **Allocation**
  - ■ Defines the number of operators for each selected type

☐ **Scheduling**
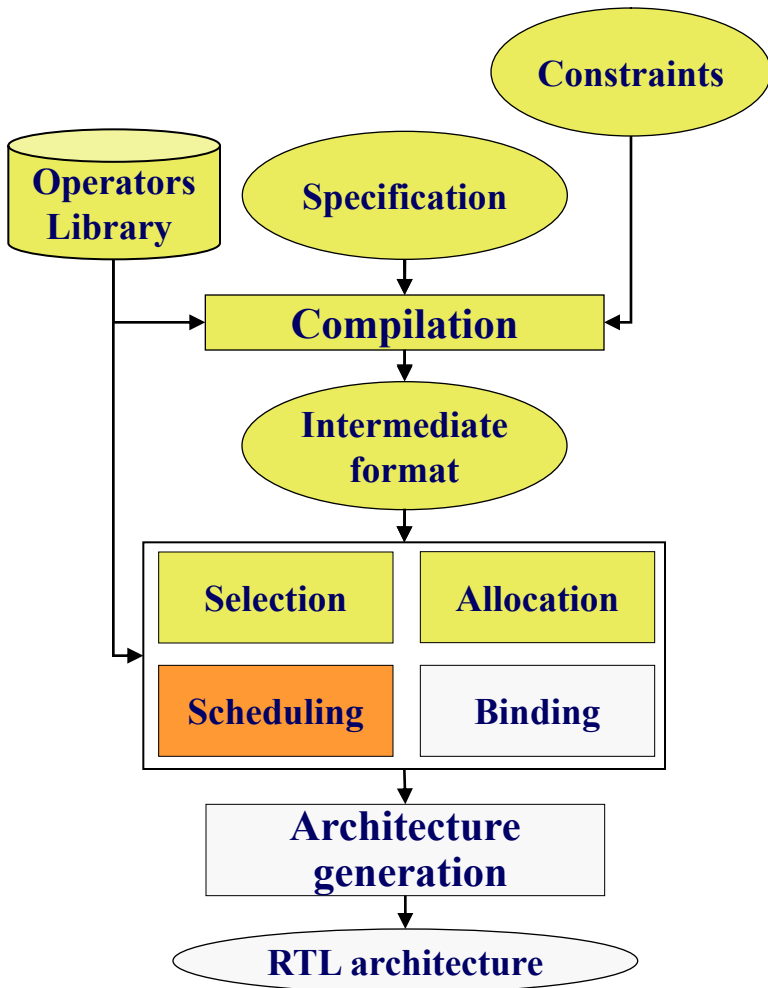  - ■ Defines the execution date of each operation

☐ **Binding (or Assignment)**
  - ■ Defines which operator will execute a given operation
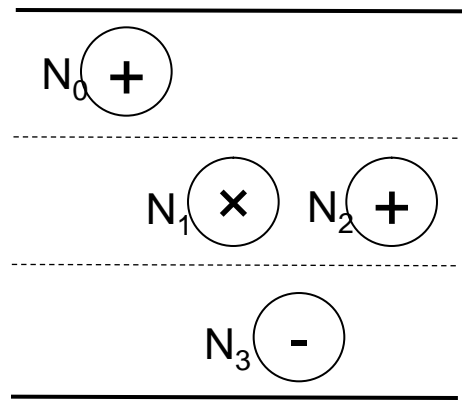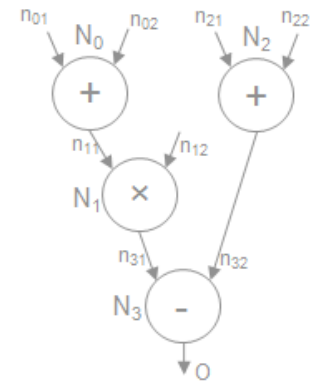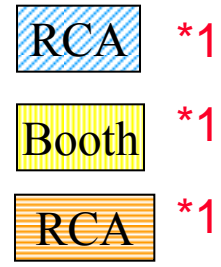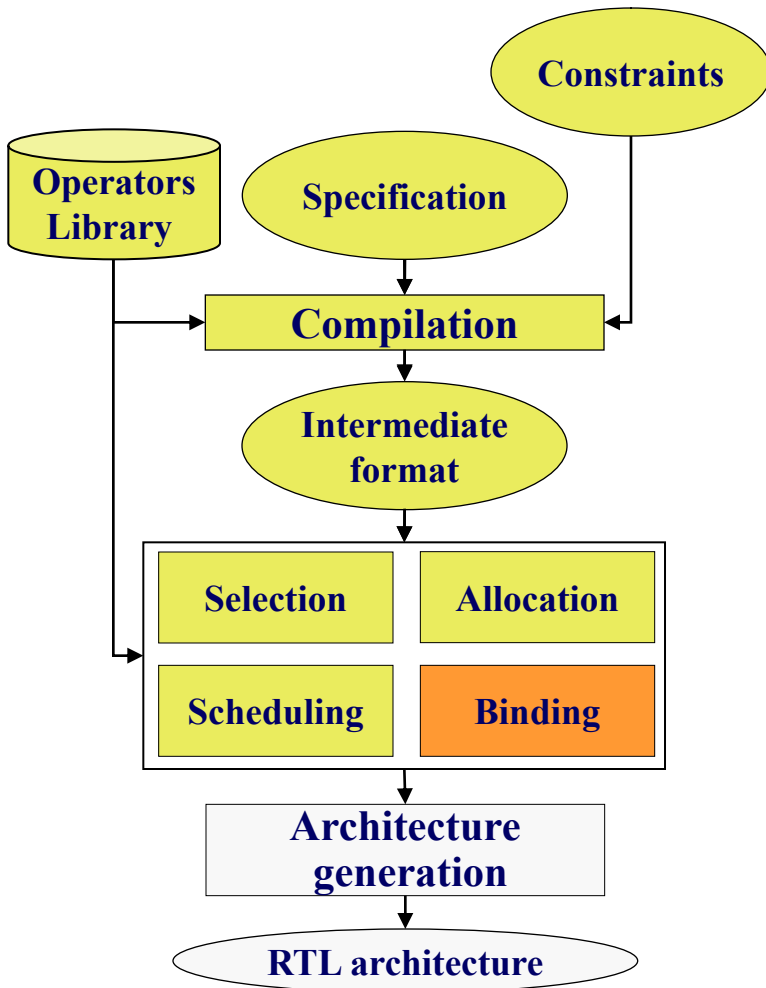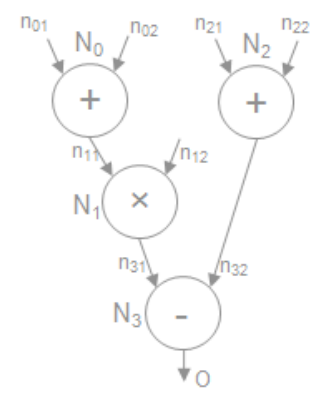  - ■ Defines which memory element will store a data

☐ **Architecture generation**
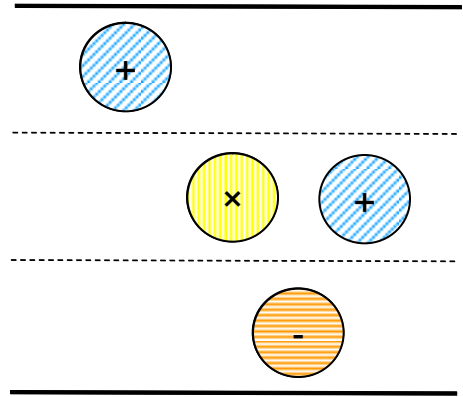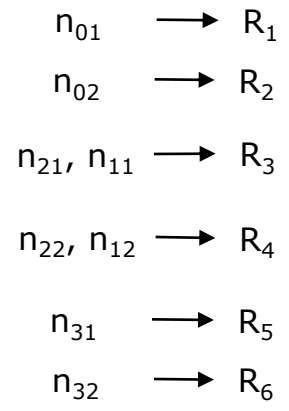  - ■ Writes out the RTL source code in the target language e.g. VHDL

# HLS steps: allocation

Constraints

Operators Library

Specification

Compilation

Intermediate format

Selection

Allocation

Scheduling

Binding

Architecture generation

RTL architecture

Operators library

| Adders | multipliers | subtractors |
|--------|-------------|-------------|
| CLA | Booth | CL |
| RCA | Wallace | RCA |

RCA  *1

Booth  *1

RCA  *1

Specification

$O = ((n_{01}+n_{02})*n_{12})-(n_{21}+n_{22})$

Intermediate representation

$n_{01}$  $N_0$  $n_{02}$   $n_{21}$  $N_2$  $n_{22}$

$+$   $+$

$n_{11}$   $n_{12}$

$N_1$  $\times$

$n_{31}$   $n_{32}$

$N_3$  $-$

$O$

# Synthesis steps

☐ **Compilation**
  - ■ Generates a formal modeling of the specification

☐ **Selection**
  - ■ Chooses the architecture of the operators

☐ **Allocation**
  - ■ Defines the number of operators for each selected type

☐ **Scheduling**
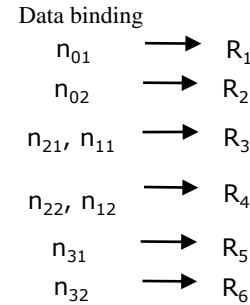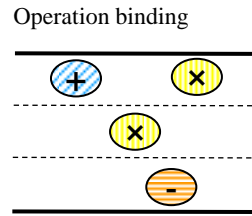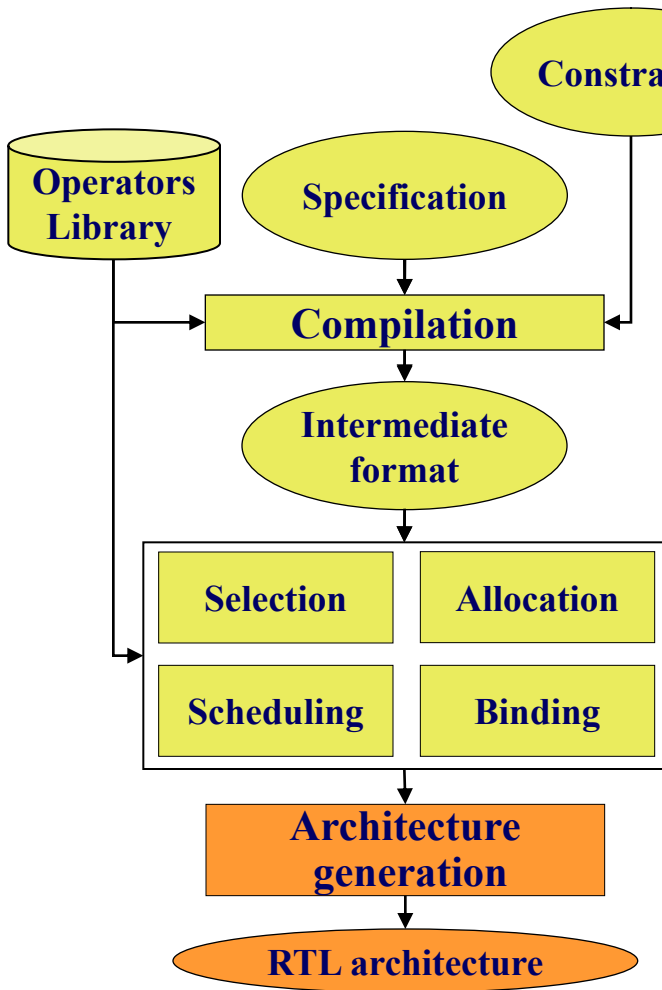  - ■ Defines the execution date of each operation

☐ **Binding (or Assignment)**
  - ■ Defines which operator will execute a given operation
  - ■ Defines which memory element will store a data

☐ **Architecture generation**
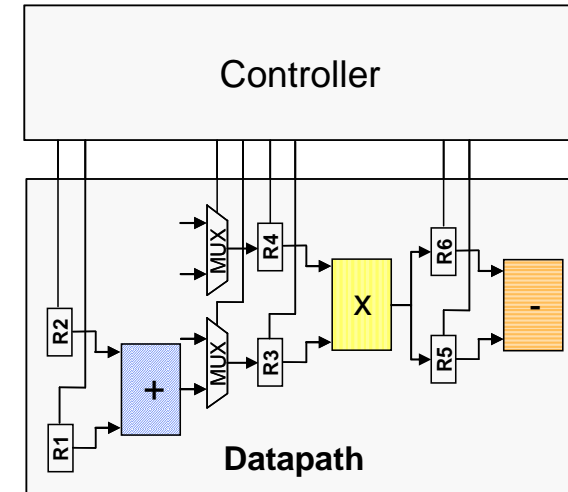  - ■ Writes out the RTL source code in the target language e.g. VHDL

# HLS steps: scheduling

# Synthesis steps

☐ **Compilation**
  - Generates a formal modeling of the specification

☐ **Selection**
  - Chooses the architecture of the operators

☐ **Allocation**
  - Defines the number of operators for each selected type

☐ **Scheduling**
  - Defines the execution date of each operation

☐ **Binding (or Assignment)**
  - Defines which operator will execute a given operation
  - Defines which memory element will store a data

☐ **Architecture generation**
  - Writes out the RTL source code in the target language e.g. VHDL

# HLS steps: binding



Constraints

Operators Library

Specification

Compilation

Intermediate format

Selection    Allocation

Scheduling    Binding

Architecture generation

RTL architecture

RCA  *1

Booth  *1

RCA  *1

Intermediate representation

$n_{01}$  $N_0$  $n_{02}$  $n_{21}$  $N_2$  $n_{22}$

$+$  $+$

$n_{11}$  $n_{12}$

$N_1$  $\times$

$n_{31}$  $n_{32}$

$N_3$  $-$

$O$

Operation binding

$+$

$\times$  $+$

$-$

Data Binding

$n_{01}$  $\longrightarrow$  $R_1$

$n_{02}$  $\longrightarrow$  $R_2$

$n_{21}, n_{11}$  $\longrightarrow$  $R_3$

$n_{22}, n_{12}$  $\longrightarrow$  $R_4$

$n_{31}$  $\longrightarrow$  $R_5$

$n_{32}$  $\longrightarrow$  $R_6$

# Synthesis steps

☐ **Compilation**
- Generates a formal modeling of the specification

☐ **Selection**
- Chooses the architecture of the operators

☐ **Allocation**
- Defines the number of operators for each selected type

☐ **Scheduling**
- Defines the execution date of each operation

☐ **Binding (or Assignment)**
- Defines which operator will execute a given operation
- Defines which memory element will store a data

☐ **Architecture generation**
- Writes out the RTL source code in the target language e.g. VHDL

# HLS steps: output



Operation binding

Data binding

$n_{01} \longrightarrow R_1$
$n_{02} \longrightarrow R_2$
$n_{21}, n_{11} \longrightarrow R_3$
$n_{22}, n_{12} \longrightarrow R_4$
$n_{31} \longrightarrow R_5$
$n_{32} \longrightarrow R_6$

Controller
- FSM controller
- Programmable controller

Datapath components
- Storage components
- Functional units
- Connection components

# RTL Architecture

□ **Controller**
- ■ FSM controller
- ■ Programmable controller

□ **Datapath components**
- ■ Storage components
- ■ Functional units
- ■ Connection components

Source :
Embedded System Design, © 2009, Gajski, Abdi, Gerstlauer, Schirner

# Example

☐ **This architecture performs the following operations:**

- store two variables coming from the port P1 in R1 and R2
- store one variable coming from the port P2 in R3
- add the variables stored in R1 and R3 and put the result in R4
- add the variables stored in R2 and R3 and put the result in R4
- connect either R1 or R2 to A1
  - ❑ *the control unit manages this connection through M1*

# RTL architecture



Source :
Embedded System Design, © 2009, Gajski, Abdi, Gerstlauer, Schirner

# Problem examples and design flow

# Resource constrained HLS

☐ **Limited number of resources**

- ■ e.g.: 2 multipliers, 3 adders
- ■ Pseudo architecture

☐ **Schedule operations according to the available operators in the current control step**

☐ **Objectives**

- ■ Minimize the latency or maximize the throughput
  - ❑ *based on operations mobility i.e. operations urgency*

# Resource constrained HLS

☐ **Limited number of resources**

  ◼ e.g.: 2 multipliers, 3 adders

  ◼ Pseudo architecture

☐ **Schedule operations according to the available operators in the current control step**

☐ **Objectives**

  ◼ Minimize the latency or maximize the throughput

    ❑ *based on operations mobility i.e. operations urgency*

⟹ **Allocation and then Scheduling**

# Time constrained HLS

□ **Latency constraint**
- e.g. 5 clock cycles to process all the data

□ **Throughput constraint**
- Cadency, initiation interval…
- e.g. process each 5 cycles a new set of input data

□ **Schedule operations by using operators as much as needed**

□ **Objective**
- Minimize the circuit area

# Time constrained HLS

☐ **Latency constraint**
- e.g. 5 clock cycles to process all the data

☐ **Throughput constraint**
- Cadency, initiation interval…
- e.g. process each 5 cycles a new set of input data

☐ **Schedule operations by using operators as much as needed**

☐ **Objective**
- Minimize the circuit area

⟹ **Scheduling and then Allocation**

# Design flows

□ **No unique design flow i.e. synthesis steps order**

- Allocation → Scheduling
  - ❑ *resource constrained*

- Scheduling → Allocation
  - ❑ *Time constrained*

- Scheduling → Binding

- Binding → Scheduling

- Scheduling & Binding

- ...

# And a lot of other problems...

- ☐ **Variable merging Storage Sharing**
- ☐ **Operation merging Operator sharing**
- ☐ **Connection merging**
  - ■ Bus sharing
- ☐ **Register merging**
  - ■ Register file...
- ☐ **Chaining**
  - ■ Several sequential operations in a cycle
- ☐ **Multi-cycling**
  - ■ One operation takes more than one clock cycle to execute
- ☐ **Pipelining**
  - ■ Pipelined Datapath, pipelined operator, pipelined controller
- ☐ **...**

# Chaining, multi-cycling



Chaining

Several sequential operations in a cycle

Multi-cycling

One operation takes more than one clock cycle to execute

# Outline

☐ **Lab-STICC**

☐ **General context**

☐ **High-Level Synthesis**

  ■ Brief introduction

  ■ "In details"

☐ **GAUT**

  ■ Overview

  ■ Results

☐ **Conclusion**

☐ **References**

# Synthesis steps

□ **Compilation**
  - Generates a formal modeling of the specification

□ **Selection**
  - Chooses the architecture of the operators

□ **Allocation**
  - Defines the number of operators for each selected type

□ **Scheduling**
  - Defines the execution date of each operation

□ **Binding (or Assignment)**
  - Defines which operator will execute a given operation
  - Defines which memory element will store a data

□ **Architecture generation**
  - Writes out the RTL source code in the target language e.g. VHDL

# High-level synthesis goal

- **Starting from a functional description, automatically generate an RTL architecture**

  - Mathematic formula
  - Matlab/Simulink
  - C/C++/SystemC
  - …

# Synthesizable models

## C for the synthesis:

- No pointer
  - *Statically unresolved*
  - *Arrays are allowed!*

- No standard function call
  - *printf, scanf, fopen, malloc…*

- Function calls are allowed
  - *Can be in-lined or not*

- Finite precision
  - *Bit accurate integers, fixed point, signed, unsigned…*
  - *Based on SystemC or Mentor Graphics data types*
    - sc_int, sc_fixed
    - ac_int, ac_fixed

# Synthesizable models

## ☐ C for the synthesis:

- Finite precision
  - *bit accurate integer, fixed point, signed, unsigned…*

S/W C: Overflow checks everywhere.

```
unsigned int    x, y, z, cy;

z = x + y;
if (0xFF..FF – x >= y)
        cy = 0;    // bit 32
else
        cy = 1;    // bit 32
```

H/W C: Check unnecessary.

```
sc_uint<32>    x, y;
sc_uint<33>    z;

z = x + y;
```

# Purely functional Example #1: a simple C code

```c
#define N 16

int main(int data_in, int *data_out)
{ static const int Coeffs [N] = {98,-39,-327,439,950,-2097,-1674,9883,9883,-1674,-2097,950,439,-327,-39,98};

  int Values[N];
  int temp;
  int sample,i,j;

  sample = data_in;
  temp = sample * Coeffs[N-1];

  for(i = 1; i<=(N-1); i++){
          temp += Values[i] * Coeffs[N-i-1];
  }

  for(j=(N-1); j>=2; j-=1 ){
     Values[j] = Values[j-1];
  }

  Values[1] = sample;
  *data_out=temp;

  return 0;
}
```

# Purely functional example #2: bit accurate C++ code

```cpp
#include "ac_fixed.h"   // From Mentor Graphics
#define PORT_SIZE ac_fixed<16, 12, true, AC_RND,AC_SAT>
// 16 bits, 12 bits after the point, quantization = rounding, overflow = saturation
#define N 16
int main(PORT_SIZE data_in, PORT_SIZE &data_out)
{
  static const PORT_SIZE Coeffs [N]={1.1, 1.5, 1.0, 1.0, 1.7, 1.8, 1.2, 1.0, 1.6, 1.0, 1.5, 1.1, 1.9, 1.3, 1.4, 1.7};
  PORT_SIZE Values[N];
  PORT_SIZE temp;
  PORT_SIZE sample;

  sample= data_in;
  temp = sample * Coeffs[N-1];
  for(int i = 1; i<=(N-1); i++){
     temp = Values [i] * Coeffs[N-i-1] + temp;
  }

  for(int j=(N-1); j>=2; j-=1 ){
     Values[j] = Values [j-1];
  }
  Values[1] = sample;


  data_out=temp;
  return 0;
}
```

# Fixed-point

## ☐ Fixed point:

| S | $b_{m-1}$ | $b_{m-1}$ | | | $b_1$ | $b_0$ | $b_{-1}$ | $b_{-2}$ | | | $b_{-n+2}$ | $b_{-n+1}$ | $b_{-n}$ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

$b = m + n + 1$ bits

format: *(b,m,n)*

Integer part: *m* bits        Fractional part: *n* bits

## ☐ Ac_fixed<W,I,S,Q,O>

- ■ W : word size (m+n+1)
- ■ I : integer part size (m+1)
- ■ S : signed or unsigned
- ■ Q : rounding mode
- ■ O : overflow mode
- ■ Equivalent to SystemC data type sc_fixed

fixed point value $n_{fx}$

$n_q$

0.11
0.10
0.01
0.00
1.11
1.10
1.01
1.00

floating point value $n_{fl}$

# Fixed point: rounding mode

SC_RND

SC_TRN

# Fixed point: overflow mode

SC_SAT

SC_WRAP

# Bit accurate operation

☐ **Sign extension before the computation**

signe

&

# Fixed point operation

☐ **Alignment before the computation**

0

signe

&

# High-level synthesis goal

- **Starting from a functional description, automatically generate an RTL architecture**

  - Algorithmic description
    - *no timing notion in the source code*

  - Behavioral description
    - *Notion of step / local timing constraints in the source code*
      - by using the wait statements of SystemC for example

  - The description can be
    - *"RTL oriented"*
    - *"Function oriented"*

# High-level synthesis

☐ **Starting from a functional description, automatically generate an RTL architecture**

■ Algorithmic description
- ❏ *No timing notion in the source code*
- ❏ *Mainly oriented toward data dominated application*
  - ▪ Highly processing algorithm like filters…
- ❏ *Initial description can be*
  - ▪ "RTL oriented"
  - ▪ "Function oriented"

■ Behavioral description
- ❏ *Notion of step / local timing constraints in the source code*
  - ▪ by using the wait statements of SystemC for example
- ❏ *Can be used for both data and control dominated application*
  - ▪ Interface controller, DMA…
  - ▪ Filters…

# High-level synthesis

☐ **Starting from a functional description, automatically generate an RTL architecture**

- ■ Algorithmic description
  - ❑ *No timing notion in the source code*
  - ❑ *Mainly oriented toward data dominated application*
    - ▦ Highly processing algorithm like filters…
  - ❑ *Initial description can be*
    - ▦ "RTL oriented"
    - ▦ "Function oriented"

- ■ Behavioral description
  - ❑ *Notion of step / local timing constraints in the source code*
    - ▦ by using the wait statements of SystemC for example
  - ❑ *Can be used for both data and control dominated application*
    - ▦ Interface controller, DMA…
    - ▦ Filters…

# Behavioral description

## Behavioral description

### ❑ Notion of step / local timing constraints in the source code

- ▪ by using the wait statements of SystemC for example

```
...
void addmul() {
    sc_signal<sc_uint<32> > tmp1;
    tmp1 = 0;
    result = 0;
    wait();
    while (1) {
        tmp1 = b * c;
        wait();
        result = a + tmp1;
        wait();
    }
}
...
```

Reset state

First state

Second state

Cycle-by-cycle FSMD
with reset state

# Function v.s. RTL description

```
01:    int OnesCounter(int Data){
02:     int Ocount = 0;
03:     int Temp, Mask = 1;
04:     while (Data > 0) {
05:      Temp = Data & Mask;
06      Ocount = Data + Temp;
07:      Data >>= 1;
08:     }
09:     return Ocount;
10:    }
```

Function-based C code

```
01:    while(1) {
02:      while (Start == 0);
03:      Done = 0;
04:      Data = Input;
05:      Ocount = 0;
06:      Mask = 1;
07:      while (Data>0) {
08:        Temp = Data & Mask;
09:        Ocount = Ocount + Temp;
10:        Data >>= 1;
11:      }
12:      Output = Ocount;
13:      Done = 1;
14:    }
```

RTL-based C code

Source :
Embedded System Design, © 2009, Gajski, Abdi, Gerstlauer, Schirner

# High-level transformations

□ **Loops**
- ■ loop unrolling
  - □ *None, partially, completely*
- ■ Loop merging
- ■ Loop tiling
- ■ …

```
for (i = 0; i<4; i++)
{
    r[i] = a[i] + b[i];
}
```



**No Unrolling**
1 Adder shared for 4 additions
Latency = 4 cycles

```
r[0] = a[0] + b[0];
r[1] = a[1] + b[1];
r[2] = a[2] + b[2];
r[3] = a[3] + b[3];
```



**Unrolling = 4 (Full)**
4 Adders in parallel
Latency = 1 cycle

```
for (i = 0; i<32; i++)
{
    a[i] = b[i] * c[i];
}
for (i = 0; i<16; i++)
{
    z[i] = a[i] + x[i];
}
```

**No Merging**
Loops execute sequentially
Latency = 48 cycles

```
for (i = 0; i<32; i++)
{
    atmp = b[i] * c[i];

    if (i<16)
        z[i] = atmp + x[i];
}
```

**Merging Enabled**
Loops execute in parallel
Latency = 32 cycles

Figures from Mentor Graphics

# High-level transformations

☐ **Loops**
- Loop pipelining,
- loop unrolling
  - ☐ *None, partially, completely*
- Loop merging
- Loop tiling
- …

☐ **Arrays**
- Arrays can be mapped on memory banks
- Arrays can be synthesized as registers
- Constant arrays can be synthesized as logic
- …

☐ **Functions**
- Function calls can be in-lined
- Function is synthesized as an operator
  - ☐ *Sequential, pipelined, functional unit…*
- Single function instantiation
- …

# Compilation

☐ **Optimization**

- Constant folding
- Dead code elimination
- Common sub-expression elimination
  - ❏ *Eliminate redundant operations*
- ...

☐ **Formal model**

- Inputs, outputs, and operations of the algorithm are identified
- Data and/or control dependencies are determined
- Intermediate representation is generated

# Control-Flow Graph CFG

☐ **Exhibits operation sequences**

- ■ Through control dependencies

☐ **The sequence of operations comes directly from the source code**

- ■ The sequence is kept unchanged
  - ❑ *This limits the parallelism which should be limited if this representation is used to model control-oriented application*

# Example

1: t = a+b;
2: u = a'-b';
3: *if* (a<b)
4:    v = t+c;
*else*
{
5:    w = u+c';
6:    v = w-d;
}
7: x = v+e;
8: y = v-e;

Source code



Graphical
representation

# Example (2)

**Filter.c**

Compiler

BB2 → If (N>0) goto BB4
Else goto BB3

BB4

BB5

BB6 → Y[j] = Y[j] + C[i]*X[N-1-i]
i++
If (i<N) goto BB6
Else goto BB7

BB7 → J++
If (j<N) goto BB5
Else goto BB3

BB3

```
Void Filtre (int N, int C[N], int X[N], int Y[N]){
    int i,j;
    for (j =0; j<N; j++){
        Y[j] = 0;
        for (i=0; i<N, i++){
            Y[j]= Y[j] + C[i]*X[N-1-i];
        }
    }
}
```

# Data Flow Graph DFG

## ☐ Exhibits the parallelism between operations

- Through data dependencies
  - *Variable node, operation node*

$$O = ((n_{01}+n_{02})*n_{12})-(n_{21}+n_{22})$$

Intermediate representation

# CDFG => DFG

## □ Exhibits the parallelism between operations
- Through data dependencies
  - □ *Variable node, operation node*

## □ Loops are completely unrolled

for i : 0 → 2
    c[i] = a[i] + b[i]

➡️

c[0] = a[0] + b[0]
c[1] = a[1] + b[1]
c[2] = a[2] + b[2]

## □ Conditional assignments are transformed
- i.e. *if/switch* constructs, are resolved by creating multiplexed values

# Example

Source code

```
1: t = a+b;
2: u = a'-b';
3: if (a<b)
4:    v = t+c;
else
{
5:    w = u+c';
6:    v = w-d;
}
7: x = v+e;
8: y = v-e;
```

# Example

1: t = a+b;
2: u = a'-b';
3: *if* (a<b)
4:    v = t+c;
*else*
{
5:    w = u+c';
6:    v = w+d;
}
7: x = v+e;
8: y = v-e;

Source code

# Data Flow Graph DFG

## ☐ Scheduling

- ■ Resource constrained
  - ❑ *Latency minimization*
    - ▪ *List-Scheduling…*
  - ❑ *Throughput maximization*
    - ▪ *Modulo scheduling (IMS, SMS…)*

- ■ Time constrained
  - ❑ *Resource minimization*
    - ▪ *Force-directed scheduling, ILP...*

## ☐ Linear FSM controller

- ■ Worst execution time for the conditional assignments

# Synthesis steps

☐ **Compilation**
  ■ Generates a formal modeling of the specification

☐ **Selection**
  ■ Chooses the architecture of the operators

☐ **Allocation**
  ■ Defines the number of operators for each selected type

☐ **Scheduling**
  ■ Defines the execution date of each operation

☐ **Binding (or Assignment)**
  ■ Defines which operator will execute a given operation
  ■ Defines which memory element will store a data

☐ **Architecture generation**
  ■ Writes out the RTL source code in the target language e.g. VHDL

Perfrom ASAP

Perfrom ALAP

Determine mobilities

Create ready list

Sort ready list by mobilities

Schedule ops from ready list

Delete scheduled ops from ready list

Add new ops to ready list

Increment state index

**RC algorithm**

no — All ops scheduled? — yes

# List-scheduling



## Constraints

- 1 adder (1 cycle)
- 1 subtractor (1 cycle)
- 1 comparing component (1 cycle)
- No chaining

# List-scheduling



**Constraints**
- 1 adder (1 cycle)
- 1 subtractor (1 cycle)
- 1 comparing component (1 cycle)
- No chaining

ASAP

ALAP

mobility

Priority = 1/Mobility

# List-scheduling



## Constraints

- 1 adder (1 cycle)
- 1 subtractor (1 cycle)
- 1 comparing component (1 cycle)
- No chaining

# List-based scheduling

□ **Scheduling under throughput constraint (cadency)**

- First operator allocation that a priori support the required parallelism
    - ❑ *In many HLS approach, an initial resource allocation is performed and subsequently modified during scheduling and/or binding => it is a lower bound*

- The average parallelism is calculated separately for each *type* of operation of the DGF

$$avr\_opr(type) = \left\lceil \frac{nb\_ops(type)}{\left\lfloor \frac{II}{T(opr)} \right\rfloor} \right\rceil$$

With *II* the Initiation Interval (cadency)
*nb_ops(type)* the number of operators of type *type*
*T(opr)* the propagation time of the operator (in cycles)

# List-based scheduling



$$\#\text{adder} = \lceil 4 * 1/3 \rceil = 2$$
$$\#\text{sub} = \lceil 3 * 1/3 \rceil = 1$$
$$\#\text{cmux} = \lceil 1 * 1/3 \rceil = 1$$

□ **Constraint**
  ■ Throughput : one iteration each 3 cycles

# Impact on the memory



**Input : X[N]**
**Constant : H[N] // in memory**

Without memory constraints
Iteration_period = 60ns
Nb_opr(*) = 2
Nb_opr(+) = 1



=> 2 memory banks
Latency(arch) = 60 ns



The memory mapping has to be done by the user

# Memory constraints



With memory constraints
=> 1 memory bank
Iteration_period = 100ns
Nb_opr(*) = 1
Nb_opr(+) = 1



Latency(arch) = 90 ns

The memory access is the bottleneck !

# Impact on the I/O interface



Input : X[N]
Constant : H[N] // in memory

Without I/O constraints
Iteration_period = 60ns
Nb_opr(*) = 2
Nb_opr(+) = 1

=> 2 input ports
Latency(arch) = 60 ns

# I/O timing constraints

With I/O constraints
4 input data in parallel
Latency = 50ns
Iteration_period = 60ns
Nb_opr(*) = 3 (and not 4)
Nb_opr(+) = 1

=> 4 input ports / 1 output port

# I/O timing constraints



With I/O timing constraints
1 data per 4 cycles
Latency (arch) = 150 ns

Iteration_period = 170ns
Nb_opr(*) = 1
Nb_opr(+) = 1

=> 1 input port / 1 output port

# Synthesis steps

☐ **Compilation**
  - ■ Generates a formal modeling of the specification

☐ **Selection**
  - ■ Chooses the architecture of the operators

☐ **Allocation**
  - ■ Defines the number of operators for each selected type

☐ **Scheduling**
  - ■ Defines the execution date of each operation

☐ **Binding (or Assignment)**
  - ■ Defines which operator will execute a given operation
  - ■ Defines which memory element will store a data

☐ **Architecture generation**
  - ■ Writes out the RTL source code in the target language e.g. VHDL

# Specification

```
Void  example (int a, int b, int g, int c, int d, int h)

{

        int e,f;
        e = a+b;
        g = a+e;
        f = c+d;
        h = f+d;

}
```

# Compilation => DFG

Void  example (int a, int b, int g, int c, int d, int h)

{

        int e,f;
        e = a+b;
        g = a+e;
        f = c+d;
        h = f+d;

}

# Scheduling

# Timing information



Data lifetimes

# Formal model for variable binding



(a) Data lifetimes

(b) Compatibility graph

# Timing information and formal model



(c) Compatibility graph

# Operation binding



Binding on $A_1$

Binding on $A_2$

# Compatibility and conflict graphs

Clique partitioning :  Binding based on a compatibility graph.
Edge exists between two data which lifetimes are not overlapping: they can share the same register.

Compatibility graph



clique (sub-graph)

Graph coloring: Binding based on a conflict graph.
 Edge exists between two data which lifetimes are overlapping:  they can not share the same register.

Incompatibility graph



Graph coloring

# (weighted) Bipartite Graph

☐ **A bipartite graph is a graph whose vertices can be divided into two disjoint sets A and B such that every edge connects a vertex in A to one in B**



Bipartite Graph

Weighted Bipartite Graph

# Example

Goal: maximize the use of existing connections between operators (Muxes optimization) while minimizing their size
weight = combination between the size and the number of connection



For each cycle (control step):

• Create a bipartite graph: free operators, operations to bind
• Compute weights

# Bipartite Weighted Matching



Maximum Weighted Bipartite  Matching : Hungarian method (munkres algorithm)

# Clique partitioning algorithm: Tseng's Algorithm



| Edge | Common neighbors |
|------|------------------|
| $e'_{1,3}$ | 1 |
| $e'_{1,4}$ | 1 |
| $e'_{2,3}$ | 0 |
| $e'_{2,5}$ | 0 |
| $e'_{3,4}$ | 1 |
| $e'_{4,5}$ | 0 |

| Edge | Common neighbors |
|------|------------------|
| $e'_{13,4}$ | 0 |
| $e'_{2,5}$ | 0 |
| $e'_{4,5}$ | 0 |

| Edge | Common neighbors |
|------|------------------|
| $e'_{2,5}$ | 0 |

Cliques:
$$s_{134} = \{v_1, v_3, v_4\}$$
$$s_{25} = \{v_2, v_5\}$$

1. Group nodes which have the greatest common neighbor number
2. Repeat until all the edges are removed
3. Each clique corresponds to a storage unit

# Data binding : the Left Edge algorithm

☐ **Data are ordered by increasing birth date**

☐ **Leftmost data are bound to distinct registers**

# Data binding : the Left Edge algorithm

☐ **Data are ordered by increasing birth date**

☐ **Leftmost data are bound to distinct registers**

# Data binding : the Left Edge algorithm



Left-edge algorithm does not take into account multiplexor cost

# Resource Binding

**Multiplexer and interconnect costs are significant.**

**Cyclic inter-dependency exists between FU binding and register binding**
To minimize interconnection, one task needs the other's result to make accurate decision



(1)A scheduling example

**(2a) FU binding + REG binding**    (2b) REG binding + FU binding

**Resource constraints: 2 FUs,  2 REGs**

The inter-dependency is far more complicated in real designs

Use « manual allocation » to change FU binding arround the best point

Use a metaHeuristic: Variable Neighborhood Search, simulated annealing or Tabu search

# Resource Binding

**Register files may be used to hide the multiplexers, which are replaced by dedicated decoders**

**Merge registers with non-overlapping access dates**



(a)   (b)   (c)   (d)   (e)

**Legend:**

MUX

FU1

FU2

REG x

# Outline

☐ **Lab-STICC**

☐ **General context**

☐ **High-Level Synthesis**

- ■ Brief introduction
- ■ "In details"

☐ **GAUT**

- ■ Overview
- ■ Results

☐ **Conclusion**

☐ **References**

# GAUT

- **An academic, free and open source HLS tool**

- **Dedicated to DSP applications**
  - Data-dominated algorithm
    - *1D, 2D Filters*
    - *Transforms (Fourrier, Hadamar, DCT…)*
    - *Channel Coding, source coding algorithms*

- **Input : bit-accurate C/C++ algorithm**
  - bit-accurate integer and fixed-point from Mentor Graphics

- **Output : RTL Architecture**
  - VHDL
  - SystemC
    - *CABA: Cycle accurate and Bit accurate*
    - *TLM: Transaction level model*
    - *Compatible with both SocLib and MPARM virtual prototyping platforms*

- **Automated Test-bench generation**

- **Automated operators characterization**

# GAUT: Constraints

**Bit accurate
Algorithm in bit-accurate C/C++**

GAUT

**Synthesis constraints**
- Initiation Interval (Data average throughput )
- Clock frequency
- FPGA/ASIC target technology

- Memory architecture and mapping
- I/O Timing diagram (scheduling + ports)
- GALS/LIS Interface (FIFO protocol)

Clock enable

Bus controller

Req(i)
Data(i)
Ack(i)

GALS/LIS interface

Controller

Data Path

Memory Unit

Specific
links &
protocols

Internal
buses

# GAUT: Design flow

Bit accu
Algorithm in bit-a

traints
erage throughput )

logy

mapping
(luling + ports)
protocol)

Specific
links &
protocols

Bus controller

Req(i
Data(i
Ack(i)

C/C++ Specification

Function library → Compilation

Characterization ← DFG → Bit-width analysis

Component library

Clustering

Constraints

- Throughput
- Clock period
- Memory mapping
- I/O timing diagram

PU synthesis

MEMU synthesis

COMU synthesis

Allocation
Scheduling
Binding
Resizing
Optimization

VHDL RTL Architecture

SystemC Simulation Model (CABA/TLM-T)

# GAUT: Compilation

# GAUT: DFG viewer

# GAUT: Operators characterization

# GAUT: Synthesis steps

# GAUT: I/O and memory constraints

# GAUT: Gantt viewer

# GAUT: Interface synthesis

# GAUT: Test-bench generation



Test-bench Generation
Modelsim Script Generation
Result File Generation

# GAUT: more than 100 downloads each year



From more than 60 countries

# Outline

☐ **Lab-STICC**

☐ **General context**

☐ **High-Level Synthesis**

  ■ Brief introduction

  ■ "In details"

☐ **GAUT**

  ■ Overview

  ■ Results

☐ **Conclusion**

☐ **References**

# Experimental results: MJPEG decoding



**Block Diagram of mjpeg baseline decoder**

| Function | Time ratio |
|---|---|
| IDCT | 43,41% |
| yuv2rgb | 15,07% |
| Entropy decoding | 8,41% |
| DeQuantization | 5,10% |
| others (each function is <5%) | 28,01% |

**Execution time ratio for software MJPEG decoding** *(by using gprof)*

# Resource estimation for IDCT



**Registers**

**LUTs**

**DSP Block**

**BRAM**

**Latency**

# Resource estimation for IDCT



**Slice** — Lower Bound Slice (chart, Slice vs Latency)

MAPPING ?

**% Error Rate** — LUT Error, Register Error (chart vs Latency)

ACCURATE

**SpeedUp = Logic Synthesis Time/HLS Synthesis Time** — SpeedUp (chart vs Latency)

FAST

# Synthesis results

| Parallelism 1 (read/write 32 bits) | | | | | | | |
|---|---|---|---|---|---|---|---|
| Operators | | | | Reg (1 bit flip flop) | Mux 2:1 | Area (slices) | Latency (cycles) | Freq (Mhz) |
| add | mult | sra | sub | | | | | |
| 2 | 3 | 1 | 1 | 2653 | 3236 | 7033 | 129 | 123,5 |
| 2 | 2 | 1 | 1 | 2818 | 3525 | 6948 | 188 | 128,1 |
| 1 | 2 | 1 | 1 | 3304 | 3905 | 6988 | 228 | 124 |
| 1 | 1 | 1 | 1 | 2876 | 3858 | 6192 | 348 | 123,7 |
| 1 | 1 | 1 | 1 | 2421 | 3938 | 6422 | 448 | 125,7 |

| Parallelism 2 (read/write 64 bits) | | | | | | | |
|---|---|---|---|---|---|---|---|
| Operator | | | | Reg (1 bit flip flop) | Mux 2:1 | Area (slices) | Latency (cycles) | Freq (Mhz) |
| add | mult | sra | sub | | | | | |
| 4 | 7 | 3 | 2 | 2904 | 2965 | 9409 | 97 | 126 |
| 2 | 3 | 1 | 1 | 2942 | 3268 | 7863 | 156 | 120,2 |
| 2 | 3 | 1 | 1 | 3112 | 3300 | 8101 | 196 | 128,9 |
| 1 | 2 | 1 | 1 | 3429 | 3969 | 7529 | 316 | 128,4 |
| 1 | 1 | 1 | 1 | 2880 | 3106 | 6498 | 416 | 121,9 |

| Parallelism 4 (read/write 128 bits) | | | | | | | |
|---|---|---|---|---|---|---|---|
| Operator | | | | Reg (1 bit flip flop) | Mux 2:1 | Area (slices) | Latency (cycles) | Freq (Mhz) |
| add | mult | sra | sub | | | | | |
| 9 | 15 | 6 | 5 | 3459 | 2694 | 12070 | 33 | 138,9 |
| 3 | 4 | 2 | 2 | 3021 | 2947 | 9282 | 92 | 132,1 |
| 2 | 3 | 1 | 1 | 2917 | 3091 | 7812 | 132 | 128,7 |
| 1 | 2 | 1 | 1 | 3462 | 4257 | 7846 | 252 | 122,5 |
| 1 | 1 | 1 | 1 | 2850 | 3314 | 6719 | 352 | 120,8 |

IDCT

YUV2RGB

| Parallelism 1 (read/write 32 bits) | | | |
|---|---|---|---|
| Reg (1 bit flip flop) | Area (slices) | Latency (cycles) | Freq (Mhz) |
| 388 | 525 | 12 | 249,18 |
| 362 | 524 | 13 | 282,11 |
| 272 | 462 | 14 | 188,96 |
| 238 | 460 | 15 | 188,96 |

# Synthesis results

**Parallelism 1 (read/write 32 bits)**

| add | mult | sra | sub | Reg (1 bit flip flop) | Mux 2:1 | Area (slices) | Latency (cycles) | Freq (Mhz) |
|---|---|---|---|---|---|---|---|---|
| 2 | 3 | 1 | 1 | 2653 | 3236 | 7033 | 129 | 123,5 |
| 2 | 2 | 1 | 1 | 2818 | 3525 | 6948 | 188 | 128,1 |
| 1 | 2 | 1 | 1 | 3304 | 3905 | 6988 | 228 | 124 |
| 1 | 1 | 1 | 1 | 2876 | 3858 | 6192 | 348 | 123,7 |
| 1 | 1 | 1 | 1 | 2421 | 3938 | 6422 | 448 | 125,7 |

**Parallelism 2 (read/write 64 bits)**

| add | mult | sra | sub | Reg (1 bit flip flop) | Mux 2:1 | Area (slices) | Latency (cycles) | Freq (Mhz) |
|---|---|---|---|---|---|---|---|---|
| 4 | 7 | 3 | 2 | 2904 | 2965 | 9409 | 97 | 126 |
| 2 | 3 | 1 | 1 | 2942 | 3268 | 7863 | 156 | 120,2 |
| 2 | 3 | 1 | 1 | 3112 | 3300 | 8101 | 196 | 128,9 |
| 1 | 2 | 1 | 1 | 3429 | 3969 | 7529 | 316 | 128,4 |
| 1 | 1 | 1 | 1 | 2880 | 3106 | 6498 | 416 | 121,9 |

**Virtual prototyping**

IDCT

**Parallelism 4 (read/write 128 bits)**

| add | mult | sra | sub | Reg (1 bit flip flop) | Mux 2:1 | Area (slices) | Latency (cycles) | Freq (Mhz) |
|---|---|---|---|---|---|---|---|---|
| 9 | 15 | 6 | 5 | 3459 | 2694 | 12070 | 33 | 138,9 |
| 3 | 4 | 2 | 2 | 3021 | 2947 | 9282 | 92 | 132,1 |
| 2 | 3 | 1 | 1 | 2917 | 3091 | 7812 | 132 | 128,7 |
| 1 | 2 | 1 | 1 | 3462 | 4257 | 7846 | 252 | 122,5 |
| 1 | 1 | 1 | 1 | 2850 | 3314 | 6719 | 352 | 120,8 |

**Hardware prototyping**

YUV2RGB

**Parallelism 1 (read/write 32 bits)**

| Reg (1 bit flip flop) | Area (slices) | Latency (cycles) | Freq (Mhz) |
|---|---|---|---|
| 388 | 525 | 12 | 249,18 |
| 362 | 524 | 13 | 282,11 |
| 272 | 462 | 14 | 188,96 |
| 238 | 460 | 15 | 188,96 |

# SoCLib: a virtual prototyping platform

☐ French National Research Project (ANR)

☐ Free and open source virtual prototyping environment
- **Library of SystemC simulation models**
- **Hardware components**
  - ☐ *CPUs, HW-ACCs, memories, busses*
  - ☐ *VCI/OCP interface protocol is used*

- **Two types of model are available for each HW component**
  - ☐ *CABA (Cycle Accurate / Bit Accurate)*
  - ☐ *TLM-DT (Transaction Level Modeling with Distributed Time)*

- **Software components**
  - ☐ *OS, API…*

- **Associated tools**
  - ☐ *Simulation, configuration, debug*
  - ☐ *Automatic generation of simulation models*

www.soclib.fr

☐ GAUT is used, to generate simulation models of HW-ACC
- **CABA and TLM-DT**

# SoCLib: Design flow

# Experiments: architecture #1



**Pure software implementation on a mono-processor architecture**

# Experiments: architecture #2



**Software implementation on a mono-processor architecture + IDCT as HW accelerator**

# Experiments: architecture #3



**Parallelized software implementation on a multiprocessor architecture**

# Experiments: architecture #4

# MJPEG Results



Execution time of the application (in cycles)
to process 50 images of 48*48 pixels

IDCT generated by GAUT reduces the application latency by 14%

Parallelization of the application on 4 CPUs reduces the latency by 21%

# MJPEG Results



The 4 HW IDCT in the multiprocessor architecture further reduce the latency by 10%

Execution time of the application (in cycles) to process 50 images of 48*48 pixels

# MJPEG Results



Execution time of the application (in ... to process 50 images of 48*48 ...

**Simulation time increase**

38%

65%

10%

Simulation time (in secondes)

# MJPEG: Hardware prototyping

☐ **Real time decoding: 24 QCIF images/sec**
  - ■ IDCT: maximum I/O bandwidth (4 parallel input ports) and the lower latency (33 cycles, Freq. 138,9Mhz)
  - ■ YUV2RGB: minimum latency (12 cycles, Freq. 249,18Mhz)

☐ **Compared to a pure SW implementation**
  - ■ 10x speed-up for the IDCT function
  - ■ 5x speed-up for the yuv2rgb function



**SoC design on a FPGA Xilinx Virtex 5 LX110 (XUPV5) board**

# Viola Jones: Hardware prototyping

■ Block Diagram of a Viola Jones Face detector



■ 7x speed-up compared to a pure sw implementation



Rgb2gray | Contrast Enhancement | Noise Reduction | Canny Edge Detector | Face Detection

# HLS for Hardware prototyping

Slope detection : acos (cordic) hwpu

Texture detection: gaussian filter and  square root hwpu

Soc Leon3 interface (AHB, Grlib)

SpeedUp >= 140

Error <= 0.00006

**SpeedUp**

**Error**

$$2^{-n} \approx 0.00001 \le error \le 2^{-(n-2)} \approx 0.00006$$

n =16, number of rotation

```
/* __ieee754_acos(x)
 * For |x|<=0.5 acos(x) = pi/2 - (x + x*x^2*R(x^2))
 * where
 * R(x^2) is a rational approximation
 *                 of (asin(x)-x)/x^3
 * For x>0.5
 * acos(x) = pi/2 - (pi/2 - 2asin(sqrt((1-x)/2)))
 */
```

# Prototyping platform

## Sundance platform

Mother board

Daughter boards
*DSP C62 C67 (Texas Instrument)*
*FPGA Virtex 1000E (Xilinx)*

Interconnection matrix
*Point to point links : Com Port (CP, up to 20 Mbytes/sec) and Sundance Digital Bus (SDB, up to 200 Mbytes/sec)*

# DVB-DSNG receiver architecture mapping

# DVB-DSNG receiver

■ Synchronization and interleaving : Sw : C62 DSP

■ Viterbi and Reed Solomon decoders : Hw : Virtex-1000E FPGA

■ 4 SDB links

■ 26 Mbps throughput (limited by the synchronization bloc…C64 for higher throughputs)

# Viterbi decoding

- functional/application parameters : state number, throughput

| State Number | 8 | 16 | 32 | 64 | 128 |
|---|---|---|---|---|---|
| Throughput (Mbps) | 44 | 39 | 35 | 26 | 22 |
| Synthesis Time (s) | 1 | 1 | 3 | 9 | 27 |
| *Number of logic elements* | 223 | 434 | 1130 | 2712 | 7051 |

- DVB-DSNG standard : throughput : 1.5 to 72 Mbps, 64 states Viterbi decoder

# Reed Solomon decoding

- functional/application parameters : number of input symbols, data symbols, throughput



- DVB-DSNG standard : 1.5 to 72 Mbps, RS (204/188) decoder

# OpenMP/HLS & manycores/HWPU



*Application with annotated function(s) to accelerate*

parser

GCC

OpenMP expansion

# OpenMP/HLS & manycores/HWPU



Application with
annotated function(s)
to accelerate

**parser** → **OpenMP expansion** →

**GCC**

outlined function(s)
to accelerate

synthesis results
(latency, area…)

**GAUT**

# OpenMP/HLS & manycores/HWPU



**GCC**

parser → OpenMP expansion →

*Application with annotated function(s) to accelerate*

*outlined function(s) to accelerate*

**GAUT**

*synthesis results (latency, area…)*

**multicore simulation**

ISS | statistics | Interconnect | memory

**HWPU**

**SystemC modules**

# OpenMP/HLS & manycores/HWPU



*Application with annotated function(s) to accelerate*

**GCC**
- parser
- OpenMP expansion
- backend

*outlined function(s) to accelerate*

*application executable*

*synthesis results (latency, area…)*

**GAUT**

**SIMULATION**
- SW
- HW

**MPARM**

**multicore simulation**
- ISS
- statistics
- Interconnect
- memory

**SystemC modules**
- HWPU

*simulator executable*

**GCC**

# MPARM Architecture



Logarithmic Interconnect (**MoT**)

ARM#0  I$  MAST PORT

ARM#M  I$  MAST PORT

HWPU#0  M/S PORT

HWPU#N  M/S PORT

L2/L3 BRIDGE

SLAVE PORT — Test-and-set semaphores

SLAVE PORT — BANK #0

SLAVE PORT — BANK #1

SLAVE PORT — BANK #N

*SHARED L1 MEM*

**Mesh Of Tree**

Cores | Routing tree | Arb tree | Mem banks

P0 P1 P2 P3

M0 M1 M2 M3 M4 M5 M6 M7

lev 1 | lev 2 | lev 3 | lev 1 | lev 2

Luca Benini, Andrea Marongiu, Paolo Burgio, *University* of *Bologna*

# Target architecture

# HWPU Integration

☐ **Interface de communication**

■ Maître / Esclave

■ Registres de configuration

☐ *Nombre d'entrées*

☐ *Nombre de sorties*

☐ *Emplacements des entrées*

☐ *Emplacements des sorties*

☐ *Mode*

☐ *Etat du HWPU*

☐ *Démarrage*

■ Les registres de configuration peuvent être doublés

☐ *Recouvrement de la configuration et du calcul*

☐ **Interface de programmation**

| Function name | Brief description |
|---|---|
| bool acc_busy () | Returns TRUE if no programming channel is available |
| void acc_reset () | Once a channel has been granted resets programming registers |
| void acc_set_input_count (int count) | Sets number of inputs |
| void acc_set_output_count (int count) | Sets number of outputs |
| void acc_set_in_addrs (int addr) | Sets current input parameter's address |
| void acc_set_out_addrs (int addr) | Sets current input parameter's address |
| void acc_trigger () | Initiates execution |
| void acc_wait () | Waits for the HWPU to complete execution |

# Example

```
void foo()
{
  int A, B, C;

  #pragma omp accelerate input(A, B) output(C)
  C = A + B;
}
```

# Example

```
void foo()
{
  int A, B, C;

  #pragma omp accelerate input(A, B) output(C)   ←
  C = A + B;
}
```

**Compilation**

```
void newfunc_acc_0 (const int * in1, const int * in2,
                    int * out1)
{ *out1 = *in1 + *in2; }
```

# Example

```
void foo()
{
  int A, B, C;

  #pragma omp accelerate input(A, B) output(C)
  C = A + B;
}
```

**Compilation**

```
void newfunc_acc_0 (const int * in1, const int * in2,
                    int * out1)
{ *out1 = *in1 + *in2; }
```

```
void foo() {
  int A, B, C;
  int *in_addrs, *out_addrs;

  in_addrs = {&A, &B};
  out_addrs = {&C};
  while (!omp_program_HWPU (2, 1,
                            in_addrs, out_addrs));

  omp_acc_wait (); }
```

# Results

# GAUT 4 (not yet available, but soon…)

- □ **An open source HLS tool**
  - ■ For both data and control-dominated algorithms (CDFG)
- □ **Input :**
  - ■ C/C++ bit-accurate integer sand fixed-points from Mentor Graphics
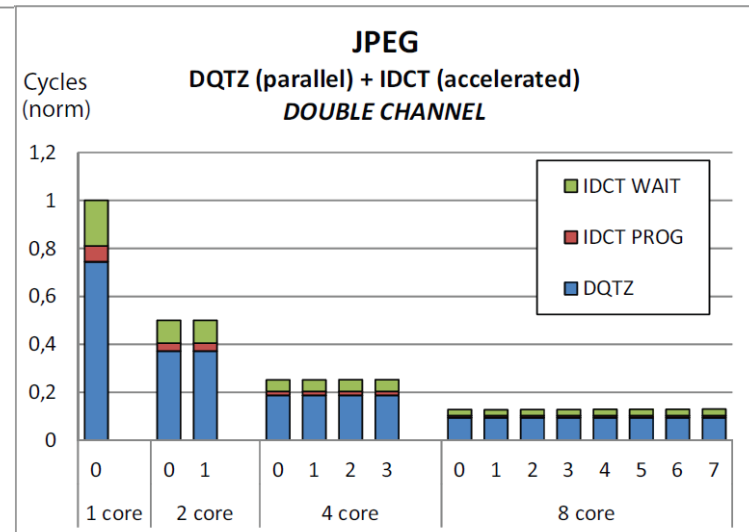  - ■ SystemC : C and C++ lack the constructs and semantics to represent design hierarchy, timing, synchronization/concurrency
  - ■ Floating point
- □ **Output : RTL Architecture**
  - ■ VHDL , Verilog
  - ■ SystemC  (CABA + TLM)
  - ■ Resource and timing estimation
- □ **Automated Test-bench generation**
- □ **Automated operators characterization**
- □ **Automated interface generation**
  - ■ AXI, AHB, FSL, …

# GAUT 4 (not yet available, but soon…)

## ☐ Constraints

- Clock, I/O protocols, loop transformations (unrolling, merging, loop pipelining with Initiation Interval), memory mapping, function inlining, resource constraints

## ☐ Objectives

- Minimization: area i.e. resources, latency, power consumption…
- Maximization: throughput

## ☐ Keys features

- *Used robust and state of the art compilation technology to extract instruction-level (Vectorization) and loop level parallelism (Polyhedral model: graphite for GCC, Polly for LLVM)*
- *Many scheduling strategies : modulo scheduling (SMS,IMS) , Force Directed List Scheduling (FDLS), System of difference constraint (SDC)…*
- *Memory analysis and optimizations: automatic partitioning of array elements to reduce conflicts and increase throughput*
- *Pattern mining for efficient resource sharing*
- *Hierarchy synthesis and function level parallelism/pipelining*
- *Design Space Exploration with directives (Loop transformation, memory partitionning) and constraints (script): one body of code, many hardware outcomes*

# Conclusion

- **HLS allows to automatically generate several RTL architectures**
  - From an algorithmic/behavioral description and a set of constraints

- **HLS allows to generate**
  - VHDL models for synthesis purpose
  - SystemC simulation models for virtual prototyping

- **HLS allows to explore the design space of**
  - Hardware accelerators
  - MPSoC architectures including HW accelerators

- **GAUT is free downloadable at**
  - http://lab-sticc.fr/www-gaut

# References

HIGH–LEVEL SYNTHESIS

Introduction to Chip and System Design

Edited by

Daniel D. Gajski

KLUWER ACADEMIC PUBLISHERS

High Level Synthesis
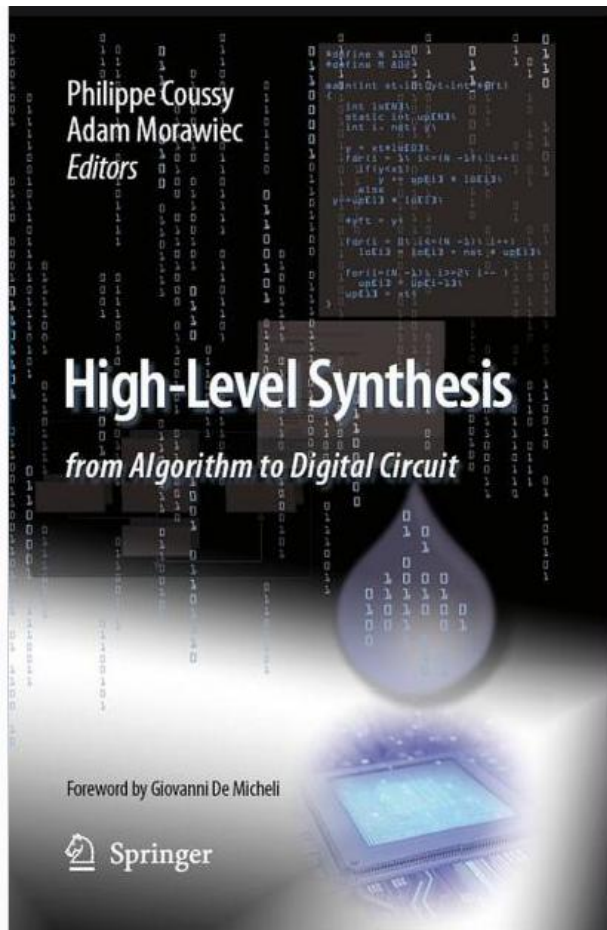of ASICs Under
Timing and
Synchronization
Constraints

David C. Ku
Giovanni De Micheli

Kluwer Academic Publishers

# References

# Academic tools

☐ **Streamroller (Univ. Mich.)**

☐ **SPARK (UCSD)**

☐ **xPilot (UCLA)**

☐ **UGH (TIMA+LIP6)**

☐ **MMALPHA (IRISA+CITI+…)**

☐ **ROCCC (UC Riverside)**

☐ **GAUT (UBS / Lab-STICC)**

☐ **…**

# Commercial tools

❑ **CatapultC (Mentor Graphics => Calypto)**

❑ **PICO (Spin-off HP => Synfora => Synopsys)**

❑ **Cynthecizer (Forte design)**

❑ **Cyber (NEC)**

❑ **AutoPilot (AutoESL => Xilinx)**

❑ **C to Silicon (Candence)**

❑ **Synphony (Synopsys)**

❑ **…**

# *Une introduction à la synthèse de haut-niveau*

## *(ou comment générer des architectures matérielles à partir du langage C)*

**Université de Bretagne-Sud**
**Lab-STICC**

**Philippe COUSSY**
**philippe.coussy@univ-ubs.fr**