



Bridge of Life  
Education

# SOC Design

## HLS Introduction

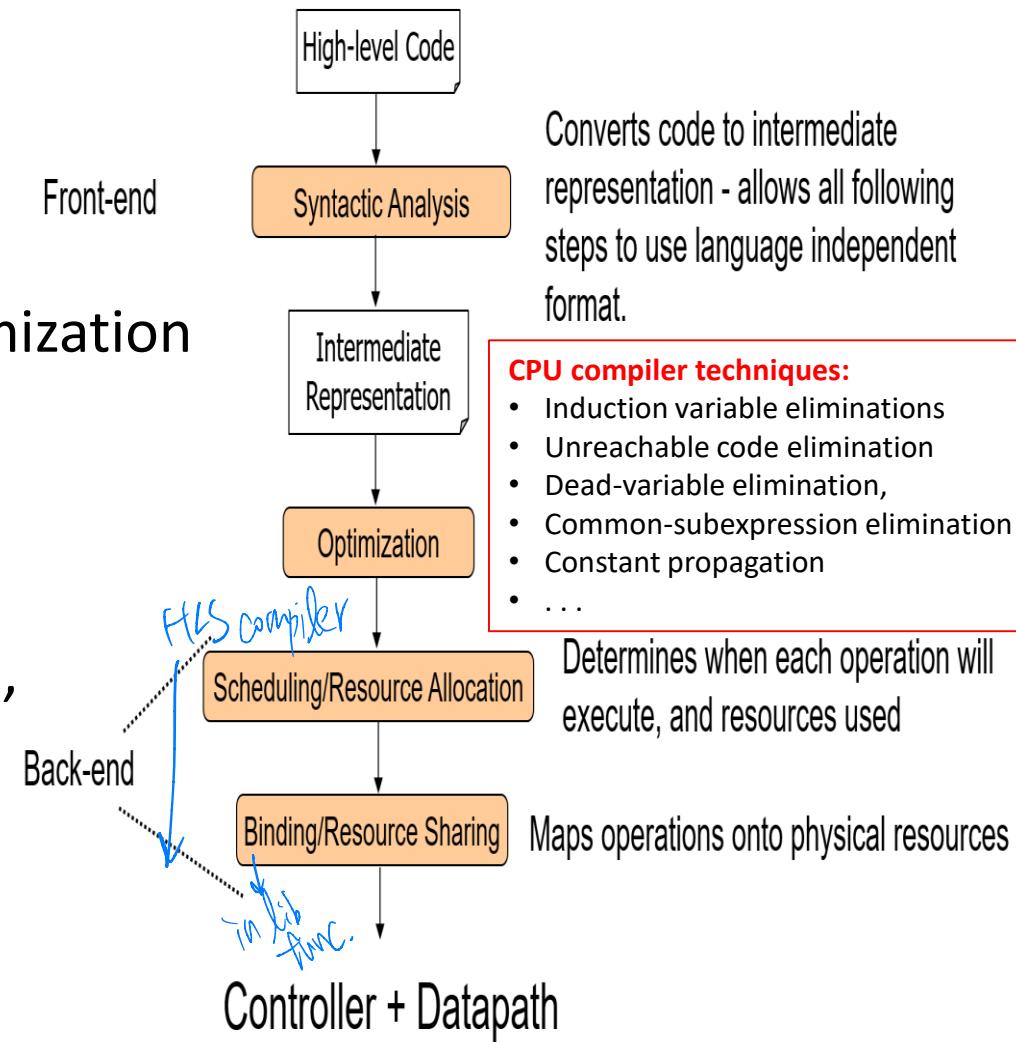
# Topics

- Unsupported C/C++ Constructs
- C to RTL Mapping
- Examples Demonstrate HLS Process
  - Expression – Datapath
  - Control Flow

# High Level Synthesis - HLS

- Convert (C/C++/OpenCL) into a RTL circuit
  - Optimize for power, performance, area, timing
  - Use Directives (**Pragma**) to direct compile/optimization process
- Vendors – FPGA vendors, IC
  - Xilinx Vitis-HLS, Intel HLS Compiler
  - Siemens/Mentor **Catapult**, Cadence Stratus HLS, (Synopsys Symphony HLS)
- Focus on the Back-end part
  - Scheduling/Resource Allocation
  - Binding/Resource Sharing

*Dynamic compiler*



# Unsupported C/C++ Constructs

# Unsupported C/C++ Constructs

- System Calls
- Dynamic Memory Usage (malloc)  
*Dynamic*
- No C++ dynamic polymorphism nor dynamic virtual function call
  - Static/Compile-time polymorphism (function/operator overloading) is ok
- Pointer Limitation  
*Mem size*  
*symmetric*  
*for loop*
- Recursive Functions

All resource must be statically allocated at compilation stage

# System Calls

① C sim  
↓  
② Call sin( RTL sim)

- e.g., printf(), malloc(), getc(), time(), sleep()
- HLS defined macro **SYNTHESIS** to exclude non-synthesized code
- **SYNTHESIS** is only defined in HLS
- Maintain the same copy of the source code for C-simulation and C/RTL co-simulation

*kernel code*

```
void hier_func4(din_t A, din_t B, dout_t *C, dout_t *D)
{
    dint_t apb, amb;

    sumsub_func(&A,&B,&apb,&amb);
#ifndef SYNTHESIS
    FILE *fp1;
    char filename[255];
    sprintf(filename,"Out_apb_%03d.dat",apb);
    fp1=fopen(filename,"w");
    fprintf(fp1, "%d \n", apb);
    fclose(fp1);
#endif
    shift_func(&apb,&amb,C,D);
}
```

# Dynamic Memory Usage

on chip  
mem buffer size

- Memory allocation: malloc(), alloc(), and free()
- User-defined macro NO\_SYNTH

```
#include "malloc_removed.h"
#include <stdlib.h>
// #define NO_SYNTH

dout_t malloc_removed(din_t din[N], dsel_t width) {
#ifdef NO_SYNTH
    long long *out_accum = malloc(sizeof(long long));
    int* array_local = malloc(64 * sizeof(int));
#else
    long long out_accum; fabie 因太小再point
    long long *out_accum = &out_accum;
    int _array_local[64];
    int* array_local = &array_local[0];
#endif
.....
}
```

# C to RTL Mapping

# Mapping of Key Attributes of C Code

Verilog

port

arbitrary  
precision

Function: design hierarchy, mapped to MODULE

Arguments : mapped to Input/output interface of the hardware

Types: All variables are of a defined type, influence the area and the performance

float IEEE 754

32 64

Loops: impact on area and performance, HLS opt with Directive Pragma

Control flow: Control logic

Arrays: impact the device area, and performance bottleneck

Expression/Operators: Function unit.  
Allocation/Scheduling (Sharing) to meet performance and area

```
46 #include "fir.h"
47
48 void fir (
49     data_t *y,
50     coer_t c[N],
51     data_t x
52 );
53
54 static data_t shift_reg[N];
55 acc_t acc;
56 data_t data;
57 int i;
58
59 acc=0;
60 Shift_Accum_Loop: for (i=N-1;i>=0;i--) {
61     if (i==0) {
62         shift_reg[0]=x;
63         data = x;
64     } else {
65         shift_reg[i]=shift_reg[i-1];
66         data = shift_reg[i];
67     }
68     acc+=data*c[i];
69     *y=acc;
70 }
```

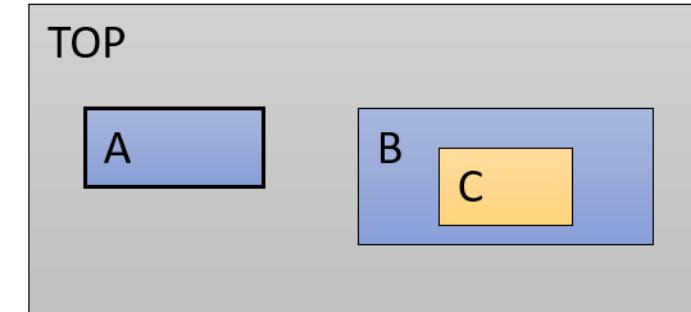
# Function Hierarchy

*↔ Verilog module Hierarchy*

*Top module*

- Top-level function becomes the top level of the RTL → (Host)
- Sub-functions are synthesized into blocks in the RTL design
- **Inlined** to dissolve the hierarchy
  - Provide greater optimization opportunity
- Vitis requires C++ kernels to be declared as extern “C” to avoid name mangling issues

```
void A { ... Body A ...}  
void C { ... Body C ...}  
void B { C; }  
void TOP() {  
    A ( ... )  
    B ( ... )  
}
```

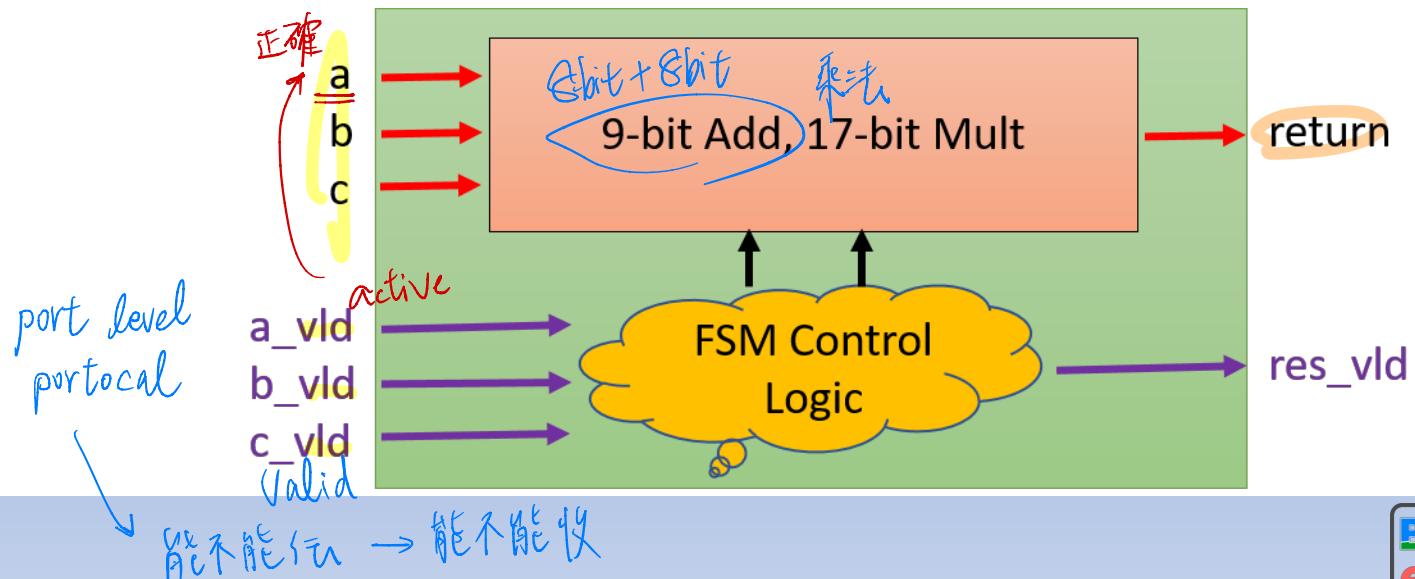


# Function Arguments

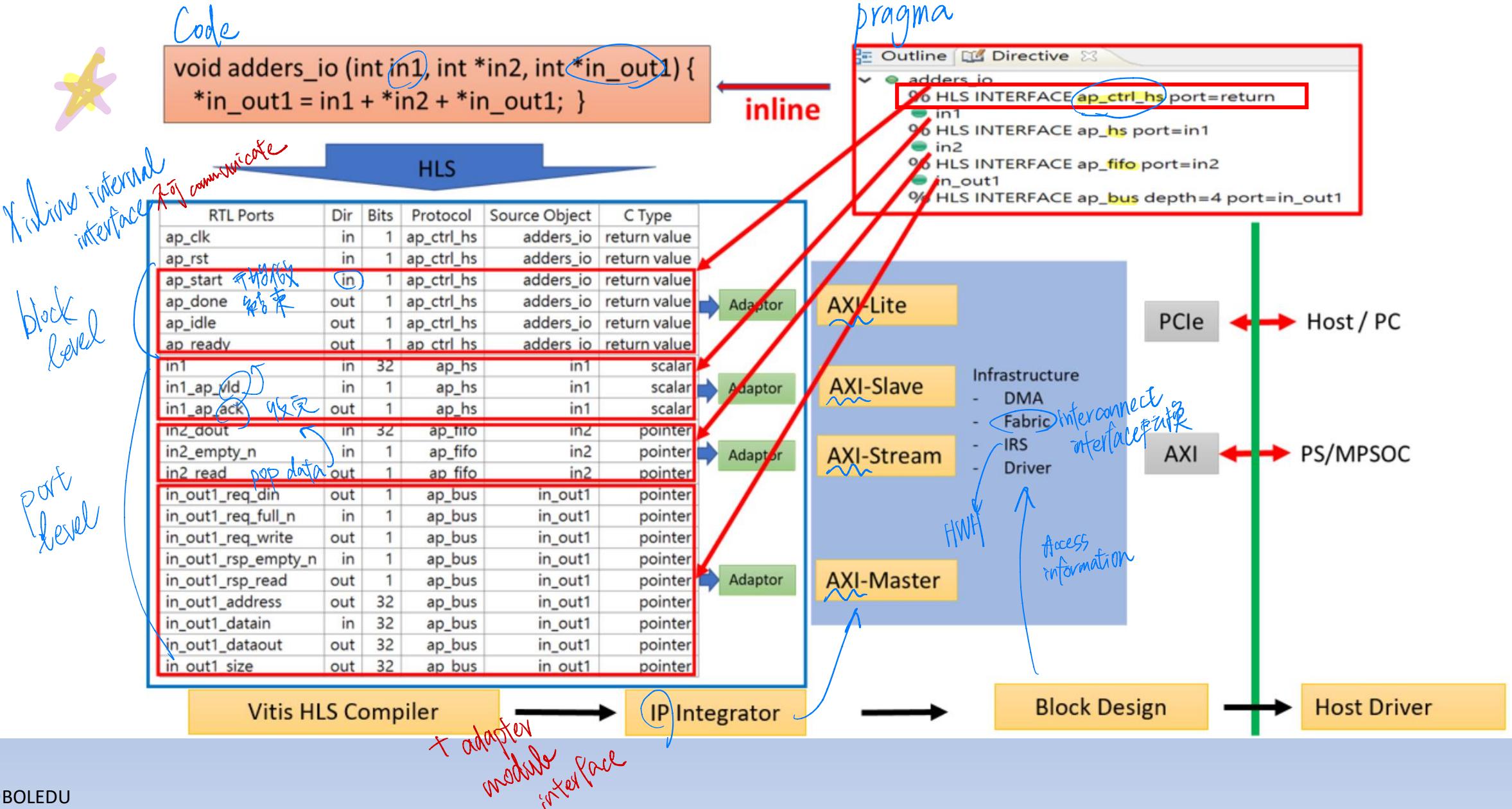
RTL port

- Function arguments mapped to **ports** on the RTL blocks
- **Global variable** if accessed only local to the function, no io port created.
- Insert control ports (**Port-level Protocol**) to automatically synchronize data exchange among blocks  
*Data 传输*
- Insert **Block-level Protocol** on Top level function to communicate with Host  
*Module on/off*    *Communicate Host*
- Arbitrary precision bit-width to reduce resource and latency

```
int17 foo_top(int8* a, int8* b, int8* c, int17* ret)
{
    int sum, multi;
    sum = *a + *b;
    multi = sum * *c;
    return multi;
}
```

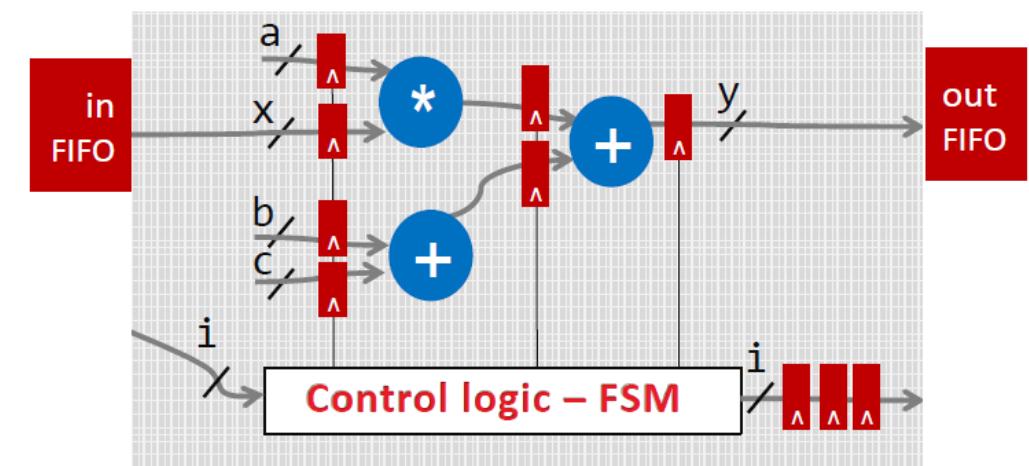
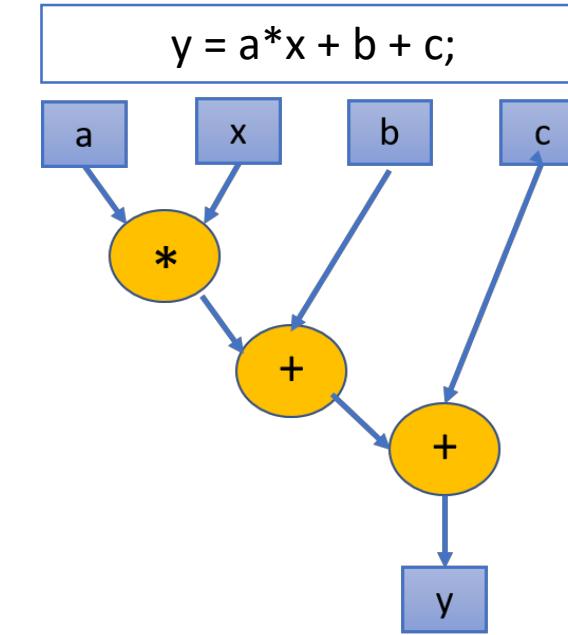


# Host/Kernel Communication



# <sup>C code</sup> Expressions – Data Flow Graph

- Start by analyzing the data dependencies between the various steps in the expression shown above. This analysis leads to a **Data Flow Graph (DFG)**
- Expression is translated to datapath and its control path (FSM)



# Resource Allocation, Scheduling, Binding

- **Resource allocation:** Each operation is mapped to a hardware resource, annotated with both timing and area information

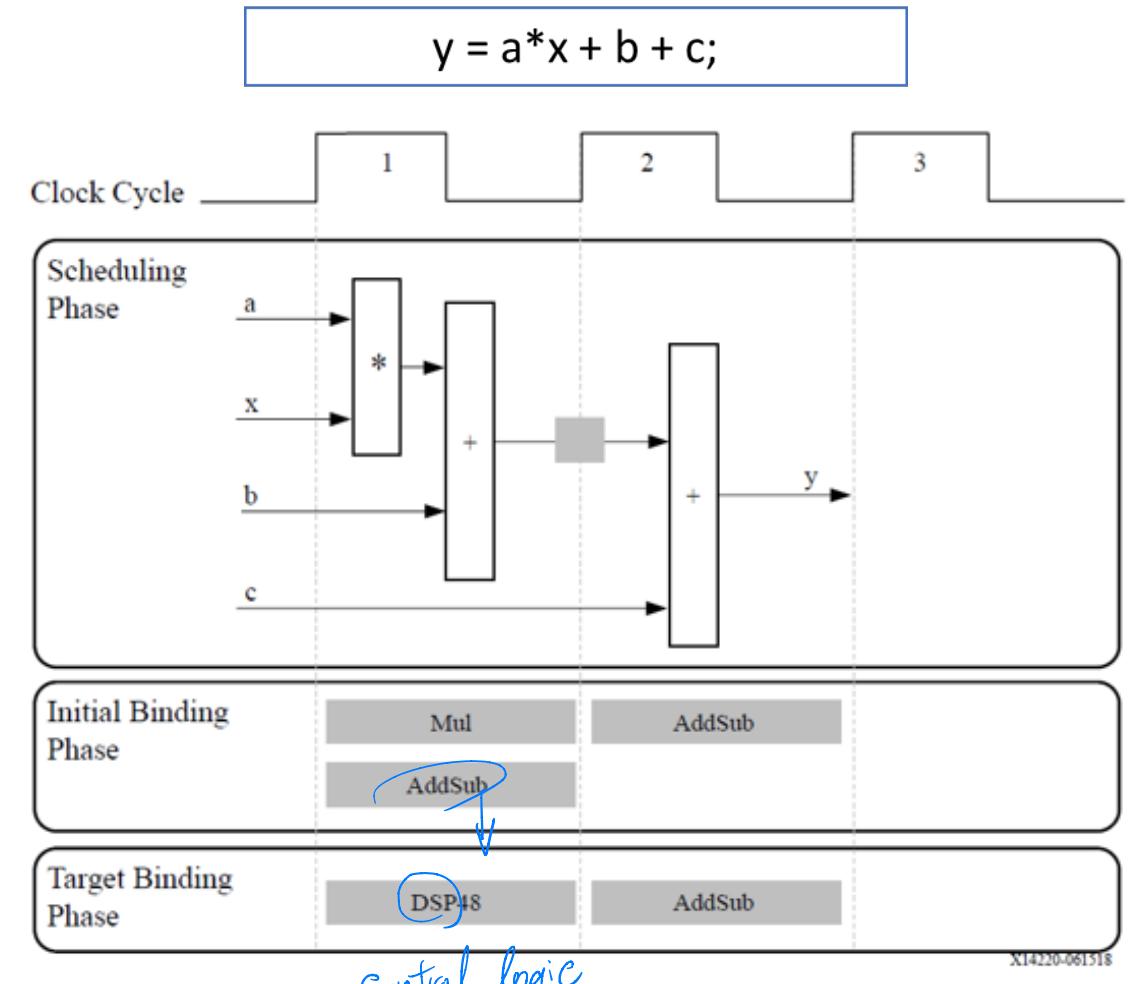
#pragma HLS allocation operation instance = add limit = 1

- **Scheduling:** decide which clock cycle to perform what operations

- **Binding:** mapped to the hardware resource.

#pragma HLS bind\_op variable=<variable> op=<type> impl=<value> latency=<int>

LUT / DSP



# Arrays

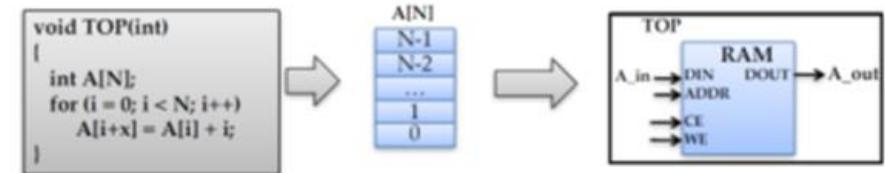
bottle necks

- Typically implemented by a memory block

- Read & write array mapped to RAM
- Constant array mapped to ROM

- Memory access is often the performance bottleneck
  - HLS default memory model assumes 2-port BRAM
  - Array can be reshaped and/or partitioned to remove bottleneck

Verilog code synthesis  
out = mem[addr]



```
void foo (...) {
...
SUM_LOOP:for(i=2;i<N;++i) {
    sum += mem[i] + mem[i-1] + mem[i-2];
}
}
```

See UG902 to get full throughput on this example  
• (Chap 3 – Array Accesses and Performance)



3CC太多 → 改寫

Example: Code implies three reads from a RAM, prevents full throughput

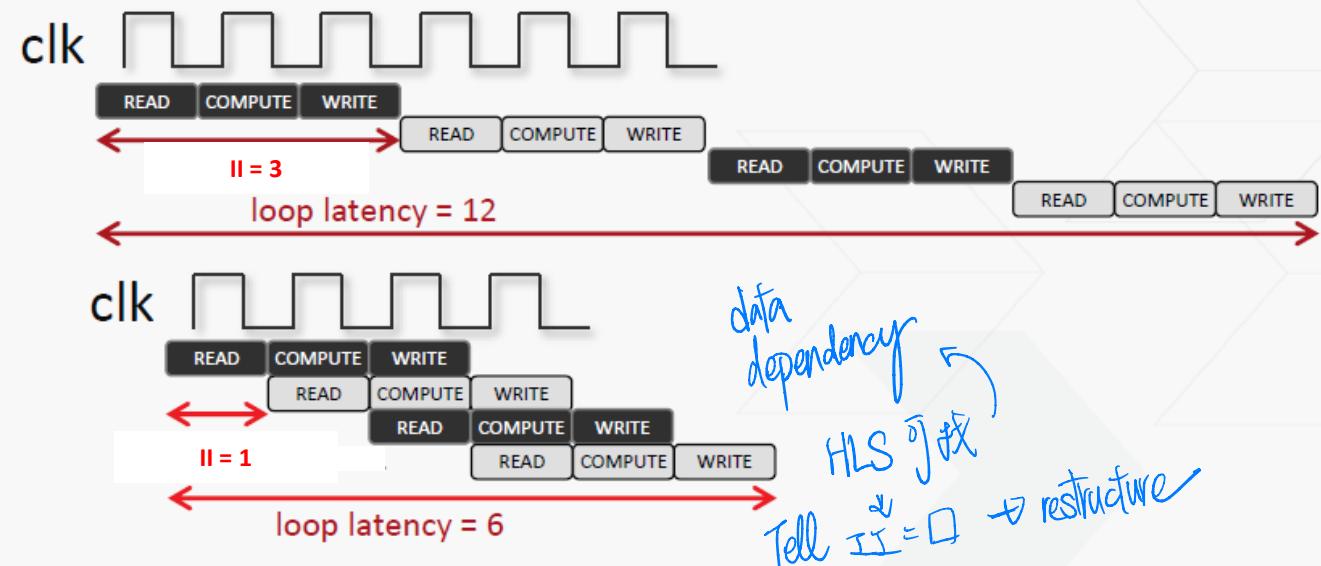
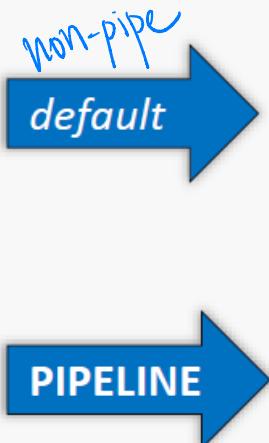
# Control Flow: Loop

- Loops are the main area of parallelism in an algorithm
- Loops can be
  - pipelined,
  - Unrolled, Partially unrolled,
  - Merged
  - Flattened
- HLS generates the datapath and control logic

# Loop - Pipeline

- One of the most important optimization
- Allow a new iteration to begin before the previous iteration is complete
- Key metric: **Initiation Interval (II)**

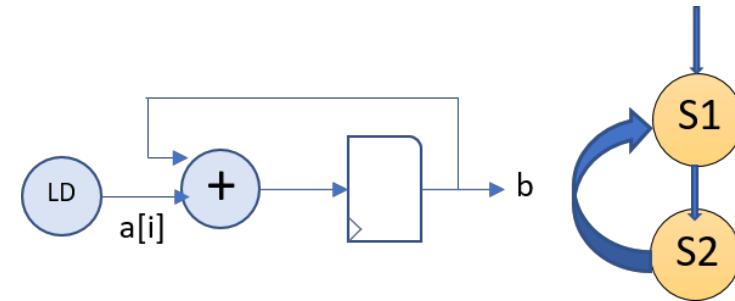
```
void F (...) {  
...  
add: for (i=0;i<=3;i++) {  
# PRAGMA HLS PIPELINE  
    op_READ;  
    op_COMPUTE;  
    op_WRITE;  
}  
...  
}
```



# Control Flow – Rolled

- By default, loops are rolled
  - Each loop iteration corresponds to a “sequence” of states (DAG)
  - The state sequence will be repeated multiple times based on the loop trip count.
  - The resource (adder) is repeatedly used in the loop iteration.
  - Efficient use the resource, but longer latency

```
void F ( . . . ) {  
    . . .  
    add: for (i=0; i <= 3; i++) {  
        b += a[i] + b;  
    . . .
```



# Loop - Unroll

手写注释：  
→ hardware  
→ memory port num

- Rolled loops can be made unrolled or partially unrolled by

**#pragma UNROLL [factor = n]**

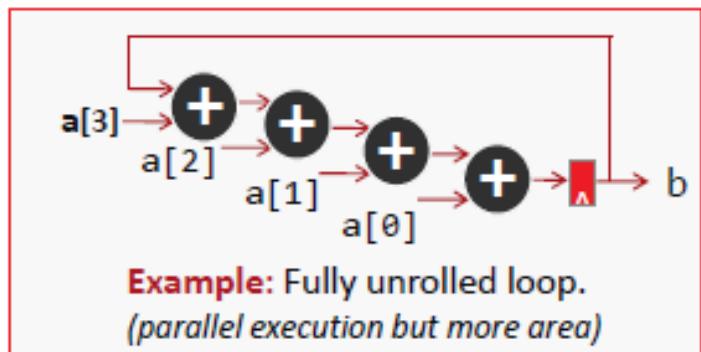
- Pros

- Decrease loop overhead
  - Increase parallelism for scheduling

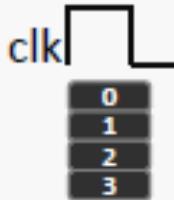
- Cons

- Increase operator count, negatively impact area, power and timing

```
void F ( . . . ) {  
    . . .  
    add: for (i=0; i <= 3; i++) {  
        #pragma UNROLL  
        b += a[i];  
    . . .
```



**Unroll: 1 cycle** →



Note: A tight timing constraint could lead to a latency different than 1 clock cycle.

# Task-Level Parallelism - Dataflow

pipeline → latency ↑ buffer ↑  
Dataflow

- > By default a C function producing data for another is fully executed first

```
// This memory can be a FIFO during optimization  
rgb_pixel inter_pix[MAX_HEIGHT][MAX_WIDTH];  
  
// Primary processing functions  
sepia_filter(in_pix, inter_pix);  
sobel_filter(inter_pix, out_pix2);
```



- > Dataflow allows Sobel to start as soon as data is ready
  - >> Functions operate concurrently and continuously
  - >> The interval (hence throughput) is improved
  - >> Channel buffer has to be filled before consumed for ping-pong

- > Dataflow creates memory channels
  - >> Created between loops or functions to store data samples
  - >> “Ping-pong” channel holds all the data
  - >> “FIFO” for sequential access, no need to store all the data

