

#### Módulo 4: Fragmentación y Predicción en SAN

#### Outline

- HLS with predictive FUs
  - FUs design
  - Centralized management
  - Distributed management
- Fragmentation in HLS
  - Dealing with heterogeneous specs
  - Techniques for area optimization
  - Techniques for time optimization
  - Pre-synthesis optimizations

#### **Prediction in HLS**

- Main requirements
  - Low percentage of mispredictions
  - Rapid recover from mispredictions
  - Low hardware overhead
  - Ideally, mispredictions should not delay other computations
  - Integrating these ideas in the HLS flow
  - Focus:
    - Designs with regularity are mainly based on interconnected cells
      - Problem: delays provoked by inter-cell communication
    - Approach:
      - Predict information from previous cells
    - Computation time depends of prediction accuracy: well suited for the asynchronous domain

#### A classical proposal<sup>1</sup>



1. W.F. Wallace, S.S. Dlay and O.R. Hinton. "Probabilistic carry state estimate for improved asynchronous adder performance", IEE Proc. Comput. Digit. Tech., 2001, 148, (6), pp. 221-226.

#### **History based prediction**

- Basic ideas:
  - Exploit data correlation to increase hit rates
  - Similar to branch prediction
  - Simple structures required: very small area penalty
- Pure prediction
  - Decide actual carry based on past information
- Hybrid prediction
  - In some cases, history is not necessary because the right carry can be easily obtained.

#### **Hybrid Prediction**



#### Simple adder design



Román Hermida Universidad Complutense

#### **Alternative adder design**



Universidad Complutense

#### **Speculative multiplier design**



## Using speculative FUs in HLS: example



#### **Common execution flow**



Román Hermida Universidad Complutense

#### **SFUs + Centralized Management**



Román Hermida Universidad Complutense

## Controller for Centralized Management



#### **SFU + Distributed Management**



Román Hermida Universidad Complutense

## Distributed Management: theoretical foundation & requirements <sup>1</sup>

- State: 1 controller per SFU
  - 1 state value per operation bound to the SFU
- Commit Signals Logic Unit. Decides:
  - If an operation can be written (commit) in its corresponding register
  - If the SFU controller can advance to the next state
    - HW to generate these signals is deduced from DFG preprocessing
- Static schedule only defines a partial order
  - Strict order is dynamically calculated
  - Actual scheduling is influenced by mispredictions and hazards
- Multicycling and chaining must be supported

Low area overhead required

1. A.A.D. Barrio, S. Ogrenci Memik, M.C. Molina, J.M. Mendías, and R. Hermida. A distributed controller for managing speculative functional units in high level synthesis. IEEE TCAD, 30:350-363, 2011.

## Controller for Distributed Management



### **Distributed Management: commit** conditions



### **Distributed Management: solving** hazards



Román Hermida Universidad Complutense

#### **Enable signals: example**

- Study of M1 controller
- Synthesis process must deduce the value for mux inputs (challenging problem)
- Remark: this is only a canonical representation. Actual implementation can be simplified.



#### **Routing and load signals: example**

- Study of register R1
- Loaded when
  - M2 controller in state 3 and state transition enabled, or
  - M1 controller in state 1 or 7, and state transition enabled
- Register input: if M2 controller in state 3, then M2; otherwise M1



# Support for multicycling and chaining

- Important requirement of any HLS strategy
- Multicycling is specially important with SFUs
  - In case of misprediction, only the last part of the operation must be done again
- HW implication of multicycling:
  - Small counter associated to each multicycled SFU
  - Loaded with: Number of computation cycles / number of recovery cycles.
- HW implications of chaining:
  - Chained operations must commit simultaneously

#### **Computation time**



- Input specs: DiffEq, 2EWF, DCT, IDCT, Lattice Filter, LMS
- Input data: from fully random to highly correlated
- Multicicycling and chained allowed



#### **Area overhead**



## Fragmentation: opportunities for HLS

#### Key idea:

• HLS golden rule: Balanced designs at the clock cycle level

• However... imbalance very common (problem-inherent causes)

#### Area optimization

- Balance of computational effort at the bit level among the different control steps: Minimizing idle FUs
- Execution time optimization
  - Balance the bit-level delays of operations in every cycle
- Approach
  - Splitting/merging operation slices than can be executed in several control steps (maybe non-consecutive)
  - Adopted in the context of heterogeneous specs ... but can also be applied to homogeneous specs



Román Hermida Universidad Complutense



- Operation scheduling with fragmentations is quite different from multicycling
  - One operation scheduled in several (consecutive or not) cycles: fragments
  - Each fragment could be implemented by a different FU
  - Every fragment narrower than original operation
  - Storage of operands not required in every cycle
  - Relaxing RAW dependences
- Also, quite different from pipelining
  - The same FU could be used to implement several fragments

- Allocation & binding
  - Some operations implemented merging previously allocated HW modules





#### **Example**



#### **Example: Scheduling**



#### **Example: Allocation**



#### **Example: Allocation**

Operation	Fragmentation	Cycle 1	Cycle 2	Cycle 3
$\mathbf{E} = \mathbf{A} \times \mathbf{B}$		⊗4×4		
$\mathbf{F} = \mathbf{C} \times \mathbf{D}$		⊗4×4		
G = E + F		<b>⊕</b> 8		
$I = H \times G$	$I1 = H \times G_{30}$		⊗8×4	
	$I2 = H_{30} \times G_{74}$		⊗4×4	
	$I3 = H_{74} \times G_{74}$		⊗4×4	
	$I4 = I1_{114} + I2$		⊕ 8	
	$I5 = (`000' \& I4_{84}) + I3$		$\oplus 8$	
	$I = I5 \& I4_{30} \& I1_{30}$			
L = J + K		<b>(D)</b> 8		
$N - L \times M$	$N1 = L \times M_{118}$	⊗8×4		
	$N2 = L \times M_{30}$			⊗8×4
	$N3 = L_{30} \times M_{74}$			⊗4×4
	$N4 = L_{74} \times M_{74}$			⊗4×4
	$N5 = N3 + N2_{114}$			⊕ 8
	$N6 = N4 + (`000' \& N5_{8.4})$			⊕ 8
	N7 = N1 + (`000' & N6)			⊕ 12
	$N = N7 \& N5_{30} \& N2_{30}$			
$\mathbf{R} = \mathbf{P} + \mathbf{Q}$	$R1 = P_{110} + Q_{110}$	<b>1</b> 2		
	$R2 = P_{2312} + Q_{2312}$		<b>1</b> 2	
	R = R2 & R1			
		<b>92 FA</b>	92 FA	92 FA

No FUs waste  $\otimes$  4×4  $\otimes$  4×4  $\oplus$  8  $\otimes$  4×8  $\oplus$  12  $\oplus$  8

Román Hermida Universidad Complutense

#### **Example: Allocation**

#### FUs vs. full circuit

	Conventional	Conventional Proposed	
Datapath FUs	⊗12×8, ⊗4×4, ⊕ 24	$\otimes 8 \times 4$ , 2 $\otimes 4 \times 4$ , $\oplus$ 12, 2 $\oplus$ 8	
FU's area	1401 equivalent gates	911 equivalent gates	35 %
Circuit area	3324 equivalent gates	2298 equivalent gates	30.86 %
Cycle length	23.1 ns	21.43 ns	7.22 %

#### Area optimization: implementation

- HLS algorithm including fragmentation during scheduling and allocation
  - Scheduling based on force-directed algorithm
    - Force considers the HW cost of every different operation or <u>fragment</u>
    - Priority: Favours ops with less mobility in cycles with small computational load
    - Bound: initially, the most uniform distribution. Relaxed as algorithm progresses.
  - Allocation & binding phase produces a datapath where all the FUs execute one operation of its same width in at least one cycle

#### **Area optimization: implementation**

Multiplications allocation (mxn). Iteratively:

- 1. Allocation and binding w/o hw waste. When size mxn is repeated in every cycle. Bind 1 op per cycle.
- 2. mxn size not exactly repeated every cycle, but waste below threshold. Bind 1 op per cycle.
- 3. <u>Fragmentation</u>: Multiplication fragments explored:
  - If a smaller fragment appears in cycles where the bigger does not
  - The allocation produces two new fragments to be treated



#### Area optimization: experiments



#### Area optimization: experiments

- Synthesis of an homogeneous specification
  - 5th order elliptic filter
    - 34 operations: 26 additions and 8 multiplications
    - Unsigned operations of 16 bits width

computational load		Comm tool	Proposed algorithm	% saved
) <b>– 9</b>	Area	7654	7398	3.34 %
<b>λ=o</b>	Clock cycle	58.63	37.27	36.43 %
2 -11	Area	7065	6438	8.87 %
λ=11	Clock cycle	51.59	41.81	18.95 %
λ=16	Area	6794	4953	27.09 %
	Clock cycle	32.27	31.13	3.53 %
Reuse fragme	of smaller ents			

#### **Execution time optimization**

- Main ideas: similar to those applied for area reductionGoal:
  - Try to break long bit-level propagation chains in the different operations
- Approach: Pre-synthesis optimization
  - Transform the original spec into a new one where a HLS commercial tool can produce better results
  - Schedule fragments in different cycles
  - Preserve the benefit of using commercial tools

#### **Example**



#### **Scheduling: small area**



Universidad Complutense

#### **Allocation: small area**



#### **Scheduling: fast execution**



#### **Allocation: fast execution**



#### **Pre-synthesis: vertical fragmentation**



#### **Alternative scheduling**



#### **Alternative allocation**







#### **Cycle length estimation**



- Calculation of every path execution time
- Operations are crossed from output to inputs
- Estimation of clock cycle length



### Partitioning

- Operations with execution times longer than CL
  - Minimum number of fragments
  - Fragments of similar size  $\rightarrow$  HW reuse
  - Compute ASAP/ALAP of every fragment
    - Fragments with identical ASAP/ALAP already scheduled
    - Fragments with different ASAP/ALAP to be scheduled by the HLS tool
  - Multiplications
    - Can be treated as rows of additions (ICCAD 2005)
    - Specific HW modules for fragmentation of multiplications (DATE 2006, TCAD 2007)
      - Better in terms of area. Similar in terms of time.

#### Some results<sup>1</sup>

- Classical HLS benchmarks
  - Performance improvement 67 % on average
  - Area increment 7 % on average
- ADPCM decoder
  - Original: Initial description + HLS tool + RTL tool
  - Optimized: transformed description + HLS tool + RTL tool
  - Identical tool in both cases

ADPCM	2	Cycle duration (ns)		Area	
module	λ	Original	Optimized	Saved	decrement
IAQ	3	6.96	2.4	66 %	2.4 %
TTD	5	9.28	3.66	61 %	6.2 %
OPFC + SCA	12	9.39	2.36	75 %	3.3 %

1. R. Ruiz-Sautua et al., Behavioural transformations to improve circuit performance in HLS, DATE 2005

#### To probe further ...

- Fragmentation and bit-lvel HLS
  - Performance-driven Read-After-Write Dependencies Softening in High-Level Synthesis, <u>ICCAD</u> 2005
  - Bitwise Scheduling to Balance the Computational Cost of Behavioural Specifications, IEEE Trans. on CAD, 2006
  - Pre-synthesis Optimization of Multiplications to Improve Circuit Performance, DATE 2006
  - Exploiting Bit-Level Delay Calculations to Soften Read-After-Write Dependences in Behavioral Synthesis, **IEEE Trans. on CAD**, 2007
  - Area optimization of multi-cycle operators in high-level synthesis, **DATE** 2007
  - Frequent-pattern Guided Multi-level Decomposition of Behavioural Specifications, <u>IEEE Trans. on CAD</u>, 2009
  - Subword Switching Activity Minimization to Optimize Dynamic Power Consumption, <u>IEEE D&T of Computers</u>, 2009
  - Power Optimization in Heterogenous Datapaths, **DATE** 2011
- Exploiting speculative FUs in HLS
  - Applying branch prediction techniques to implement functional units, <u>ICCD</u> 2008
  - Using Speculative Functional Units in high level synthesis, <u>DATE</u> 2010
  - A Distributed Controller for Managing Speculative Functional Units in High Level Synthesis, **IEEE Trans. on CAD**, 2011