



# Multimedia System-on-Chip Design

with specialization on

Application Acceleration with High-Level-Synthesis



## **賴瑾 Jiin Lai**

**Founder, CTO, VIA Technologies Inc.**

- ***Bachelor's Degree of National Taiwan University  
majoring in Electrical Engineering 1983***
- ***Master of Science of University of Texas, Austin  
majoring in Computer Engineering 1987***

**Jiin Lai was the Chief Technology Officer for VIA Technologies. He has over 30 years experience in the PC industry, and in the past 12 years in storage area. Early in his career, he is a software engineer developing EDA tools. Later he co-founded VIA technologies, developing PC chipsets, and x86 processor. He led the engineering team to develop Intel and AMD compatible chipsets, and x86-compatible processors. In the past decade, he developed SSD controller, and later, shift focus on developing distributed computational storage system. His responsibility including product and architecture development, with an eye toward to future computing architecture need. Mr. Lai holds over 50 US patents.**

# Topics

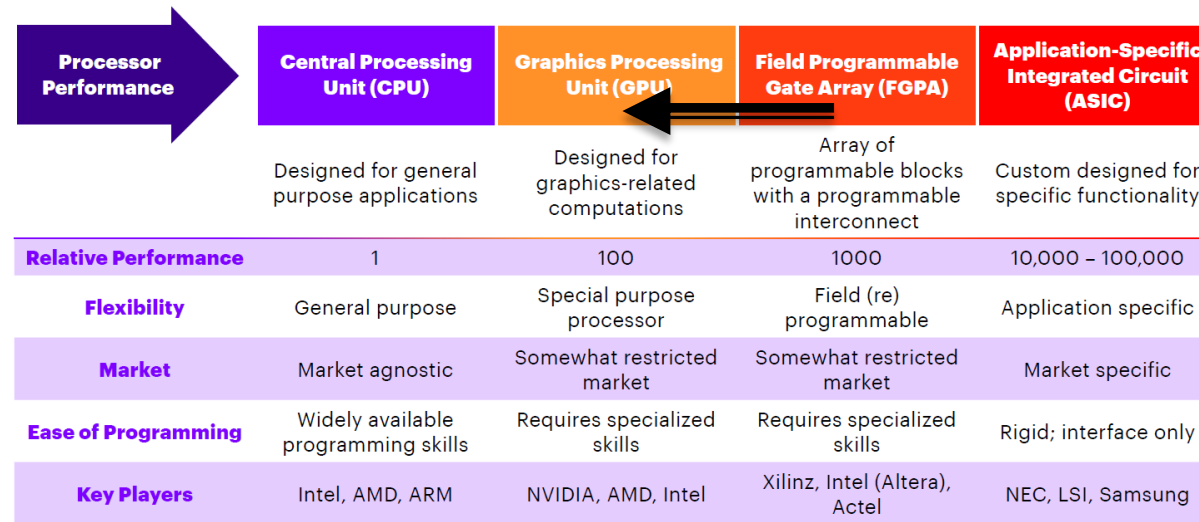
- Objectives of HLS
- HLS Course Logistics

# High-Level Statements

- Computational power hits a plateau
- Hardware accelerators to the rescue

**Figure 2:** Relative hierarchy of hardware accelerators

(Source: Accenture analysis)



- The rise of custom accelerator marketplace

Software defined Hardware – the new era of computing infrastructure

[https://www.accenture.com/\\_acnmedia/pdf-84/accenture-software-defined-hardware-v09.pdf](https://www.accenture.com/_acnmedia/pdf-84/accenture-software-defined-hardware-v09.pdf)

# FPGA Development Made Easy

- HW language are low-level and very difficult

- Know nothing about hardware design

- How software application interacts with FPGA

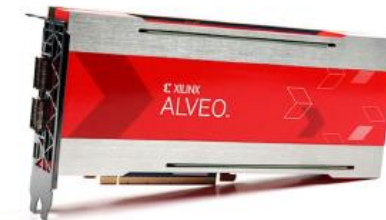
- Use C, C++, OpenCL, Python, or TensorFlow

- Parallel programming Concept apply

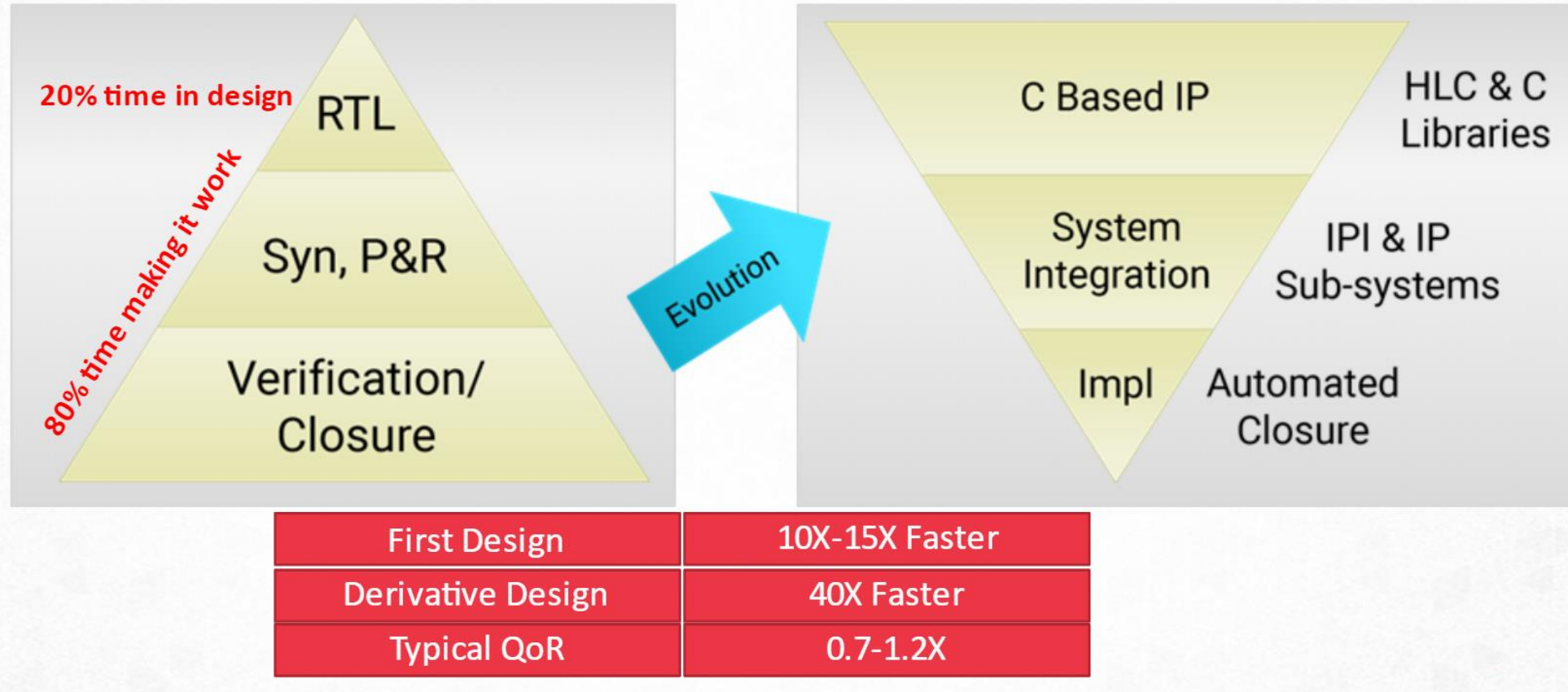
- Off-shelf Platform ready

# Use C, C++ to program FPGAs

```
// 2D Convolution (11x11)
for (int y = 0; y < height; y++) {
  for (int x = 0; x < width; x++) {
    window = in_stream.read();
    sum = 0;
    for(int row=0; row<11; row++) {
      for(int col=0; col<11; col++) {
        if (index_in_range(x+col-5, y+col-5, width, height)) {
          sum += window.val[row][col]*coeffs[row][col];
        }
      }
    }
    out_stream.write(sum);
  }
}
```

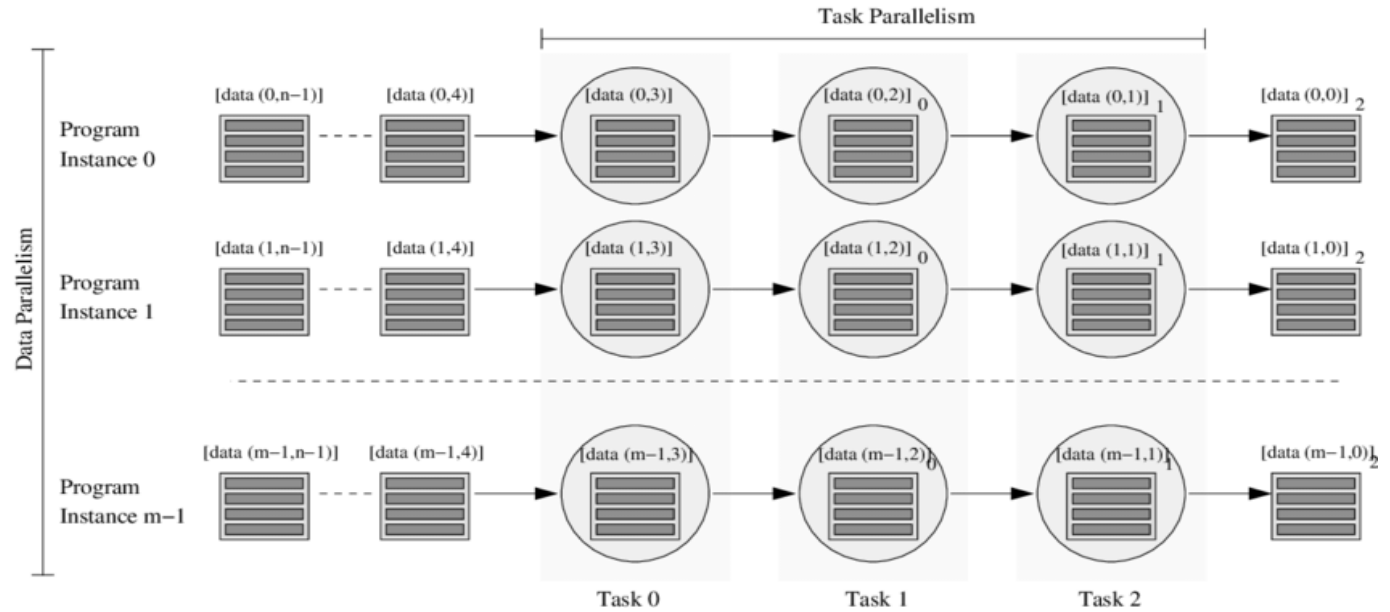


# RTL vs C-Based Design

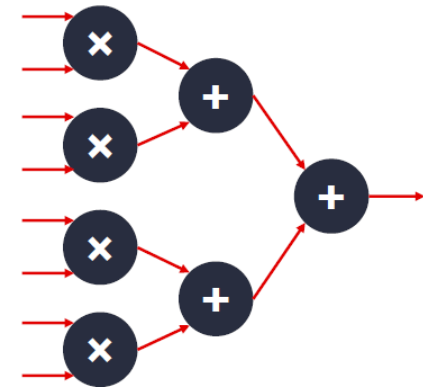


# Think “Parallel”

- Data-level Parallelism
- Task-Level Parallelism
  
- Instruction (operator) - Level Parallelism

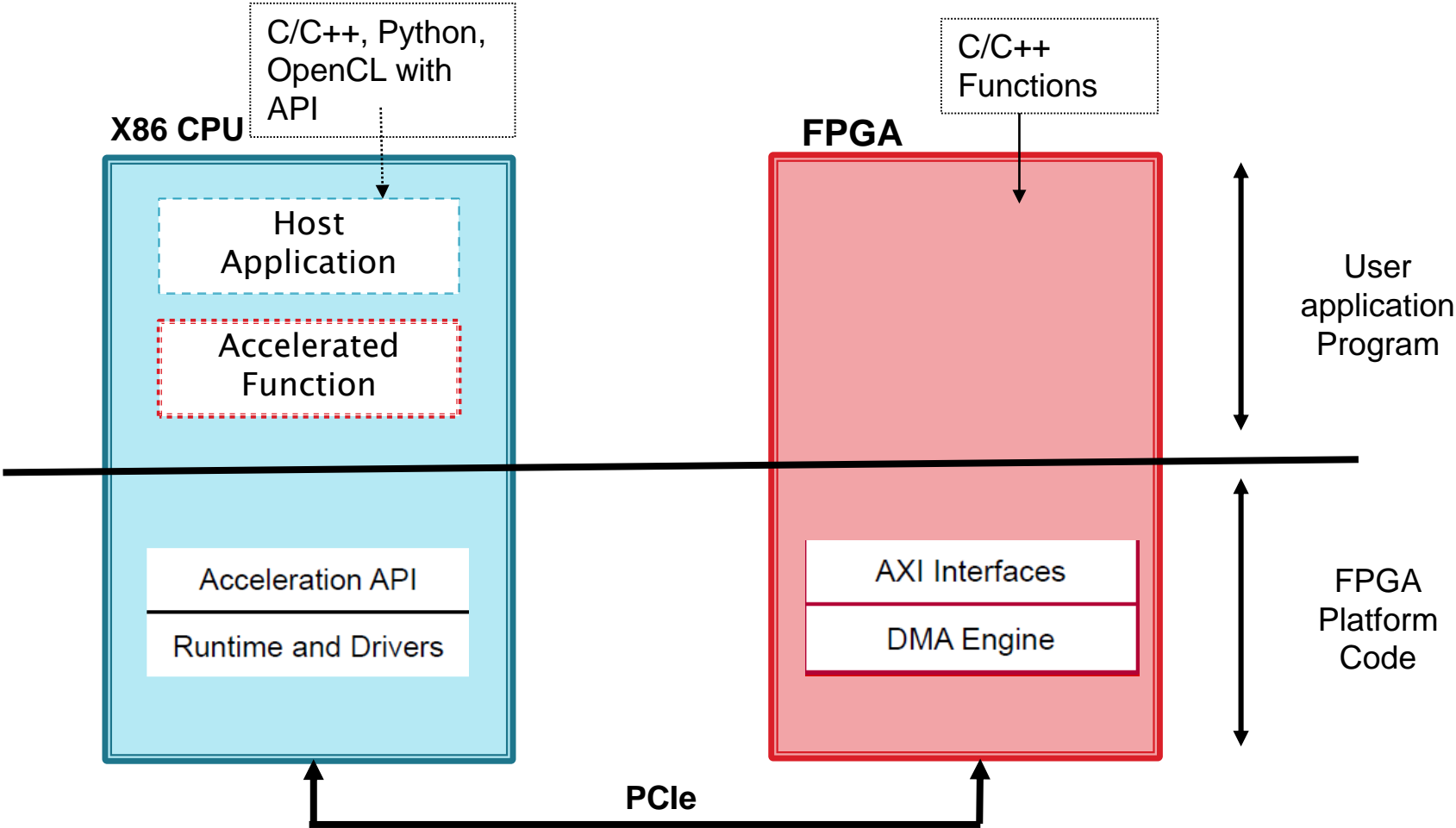


```
for (int i=0; i<N; i++)  
{  
    acc += A[i] * B[i];  
}
```





# Software Interacts with FPGA



# Speedup Development by Libraries

## Use Extensive, Open Source Libraries



### Domain-Specific Libraries



Vision & Image



Quantitative Finance



Data Analytics & Database



Data Compression

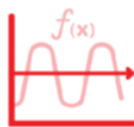


Data Security



Partner Libraries

### Common Libraries



Math



Linear Algebra



Statistics

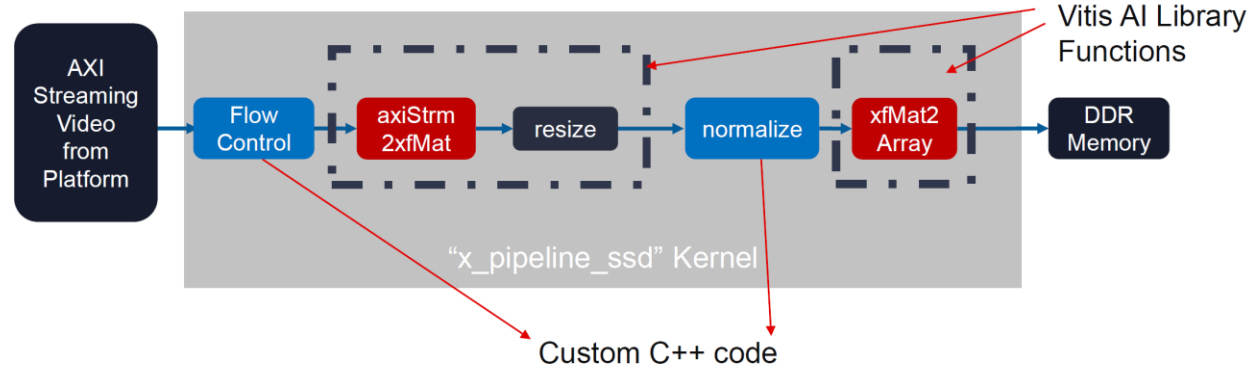


DSP



Data Management

400+ functions across multiple libraries for performance-optimized out-of-the-box acceleration



## Code Snippet

Top Level Function Definition

```
void x_pipeline_ssd(hls::stream<ap_axiu<24, 3, 1, 1>> &image_in,
                  ap_uint<AXI_WIDTH> *image_out,
                  int width_in, int height_in, int width_out,
                  int height_out, int use_mean, float scale_r,
                  float scale_g, float scale_b, unsigned char mean_r,
                  unsigned char mean_g, unsigned char mean_b)
```

Interface definitions

```
#pragma HLS INTERFACE axis port = image_in
#pragma HLS INTERFACE m_axi port = image_out offset = slave bundle = image_out_gmem depth = 131072
#pragma HLS INTERFACE s_axilite port = image_out bundle = control
#pragma HLS INTERFACE s_axilite port = width_in bundle = control
#pragma HLS INTERFACE s_axilite port = height_in bundle = control
#pragma HLS INTERFACE s_axilite port = width_out bundle = control
#pragma HLS INTERFACE s_axilite port = height_out bundle = control
#pragma HLS INTERFACE s_axilite port = use_mean bundle = control
#pragma HLS INTERFACE s_axilite port = scale_r bundle = control
#pragma HLS INTERFACE s_axilite port = scale_g bundle = control
#pragma HLS INTERFACE s_axilite port = scale_b bundle = control
#pragma HLS INTERFACE s_axilite port = mean_r bundle = control
#pragma HLS INTERFACE s_axilite port = mean_g bundle = control
#pragma HLS INTERFACE s_axilite port = mean_b bundle = control
#pragma HLS INTERFACE s_axilite port = return bundle = control
```

Dataflow the processing

Internal streaming variable declarations

```
#pragma HLS DATAFLOW
```

Synchronize to start of frame

```
hls::stream<ap_axiu<24, 0, 0, 0>> px_in_synced;
xf::cv::Mat<XF_BUC3, MAX_IN_HEIGHT, MAX_IN_WIDTH, NPC> in_mat(height_in, width_in);
#pragma HLS STREAM variable=in_mat.data depth=256 dim=1
xf::cv::Mat<XF_BUC3, MAX_IN_HEIGHT, MAX_IN_WIDTH, NPC> in_rgb(height_in, width_in);
#pragma HLS STREAM variable=in_rgb.data depth=256 dim=1
xf::cv::Mat<XF_BUC3, MAX_OUT_HEIGHT, MAX_OUT_WIDTH, NPC> out_rgb(height_out, width_out);
#pragma HLS STREAM variable=out_rgb.data depth=256 dim=1
xf::cv::Mat<XF_BUC3, MAX_OUT_HEIGHT, MAX_OUT_WIDTH, NPC> out_mat(height_out, width_out);
#pragma HLS STREAM variable=out_mat.data depth=256 dim=1
```

Convert from axi stream to xf::Mat

```
flow_control(image_in, px_in_synced, height_in, width_in);
```

```
xf::cv::axiStrm2xfMat<24, XF_BUC3, MAX_IN_HEIGHT, MAX_IN_WIDTH, NPC>(px_in_synced, in_mat);
```

Resize the image

```
xf::cv::resize<XF_INTERPOLATION_MN,
XF_BUC3,
MAX_IN_HEIGHT,
MAX_IN_WIDTH,
MAX_OUT_HEIGHT,
MAX_OUT_WIDTH,
NPC,
MAX_POWER>(in_mat, out_rgb);
```

Subtract the mean values and apply scale

```
image_normalize(out_rgb, out_mat, use_mean, scale_r, scale_g, scale_b, mean_r, mean_g, mean_b);
```

Convert from xf::Mat to memory mapped interface

```
xf::cv::xfMat2Array<AXI_WIDTH, XF_BUC3, MAX_OUT_HEIGHT, MAX_OUT_WIDTH, NPC>(out_mat, image_out);
```

© Copyright 2020 Xilinx

# Example of Oil, Gas workload

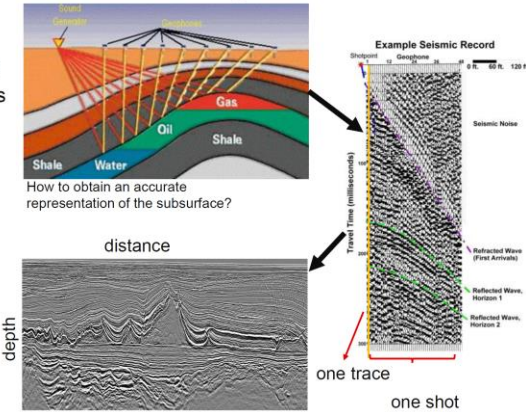
## Productivity



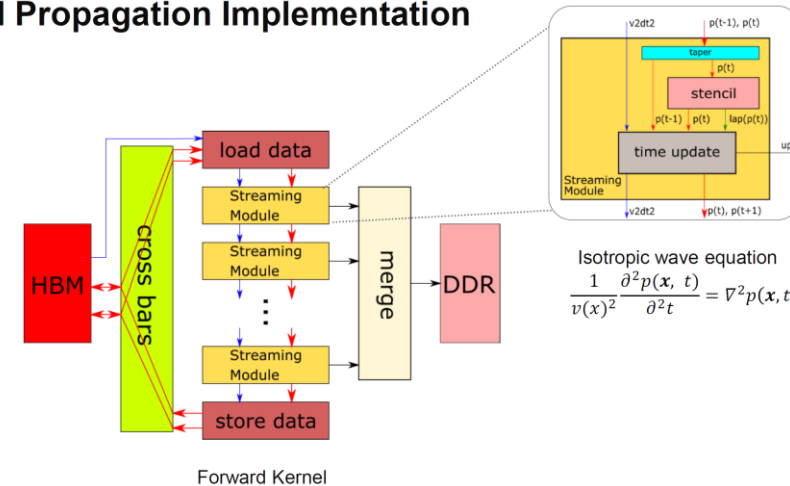
- ▶ Not the traditional programming model for FPGAs:
  - ▶ One Software Engineer, no previous O&G experience, one month to describe & implement entire RTM Algo in C++
    - No optimized library calls, completely described in C++
    - < 500 lines of code, < 50 Pragmas
  - ▶ Standard language, open source tools and libraries

## Seismic Method for Oil and Gas industry

- ▶ Seismic Imaging Technology
  - Seismic Survey: Acoustic wave sampling
  - Seismic Imaging: Mathematically process the wave traces to create an image
- ▶ RTM (Reverse Time Migration)
  - High-fidelity algorithm for imaging complex sub-surface structures
  - Cross-correlation between source wavefield and receiver wavefield
  - Wavefield reconstruction by saved boundaries



## Forward Propagation Implementation



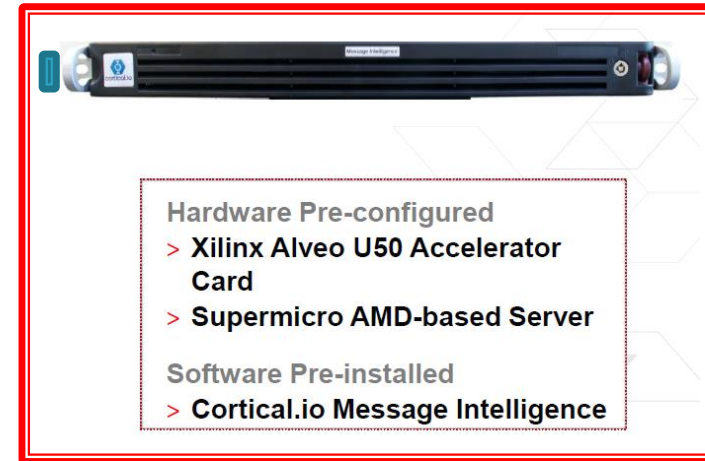
# Message Intelligence Appliance Cortical.io

- Semantic Supercomputing for NLU (Natural Language Understanding)
- Automatically classifies message based on semantics/meaning of the content

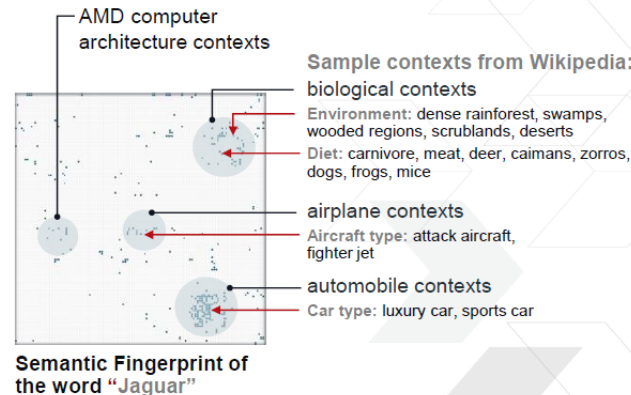
## Semantic Folding Explained

Words, sentences & paragraphs are represented by a semantic fingerprints

- > Each word is represented by **16K binary** contexts in a 2D vector
- > All operations are **binary**
- > **Minimal** source material required: reference material, textbooks, data sheets, emails, etc.
- > Creation of the semantic fingerprints is **completely unsupervised**
- > **All meanings** of a word are represented

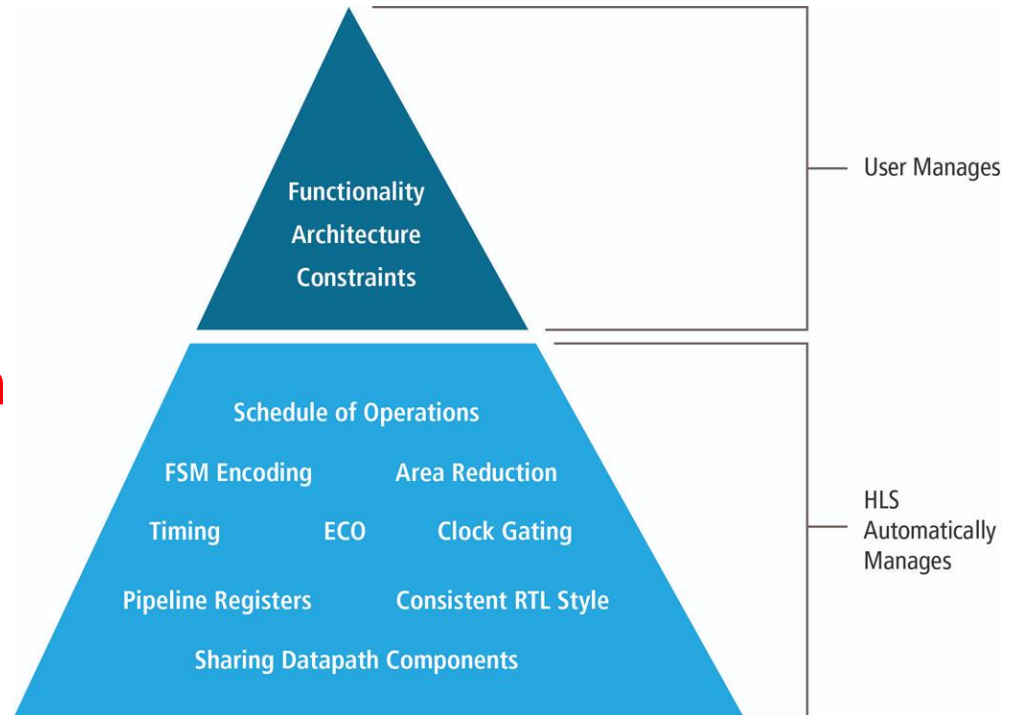


## Semantic Fingerprint



# Why HLS?

- Productivity (Design and Verification)
- IP Reuse
- Better QoR
- **End-to-end application acceleration by software designer**
- For academic, A great tool/skill for research.



# Course Objectives

Empower software designers to develop efficient application accelerator

# Course Contents

- Course Texts:
  - R. Kastner, Parallel Programming for FPGAs, arXiv, 2018
  - Xilinx ug902
- Supplementary Materials:
  - Reference Papers
  - Manual/Datasheets
- Lecture ppt & video – 16 sessions
- Labs – ~200 optional lab references
- In-class presentation – 5 sessions + final project presentation
- Final project & presentation



# Logistics

- Off-class lecture & lab/assignment
  - Lecture is self-paced
  - Lab/assignment is self-paced with lab-work submission
- In-class presentation & discussion
  - Sign-up by Google-Form (submit by Thursday 3pm)
  - Presentation selected based on available time slots, weight, and submission time.
- Refer to “HLS Course Plan.doc”

# In-class schedule and subjects

<https://cool.ntu.edu.tw/courses/3773/modules/items/110288>

session	Date	Suggest lecture title - self-paced	pdf	video	In-class discussion topics	Assignment
1	18-Sep	Course Introduction	§	§		HLS flow
		Introduction PYNQ & Lab2	§	§		ug871 labs <b>ug871-[1:7]</b>
		Vitis OpenCL XRT and Lab3	§	§		Lab1: Tool Installation
		Introduction to FPGA				Lab2: PYNQ axi-m & stream Lab3: OpenCL/XRT
2	16-Oct	Kernal IO Interface	§	§	ug871 Labs ug871-[1:7]	Xilinx Training Lab - <b>xtrain-[1:13]</b>
		Introduction to High Level Synthesis	§	§	HLS, Vivado, Vitis usage experience	Xilinx HLS Coding Style
		FPGA - CLB	§	§	Lab1, Lab2, Lab3 sharing	Xilinx HLS Design - <b>xdesign-[1:15]</b>
		FPGA - Memory	§	§		
			§	§		
3	30-Oct	System Optimization - Host	§	§	Xilinx Training Lab -xtrain-[1:13]	Vitis Tutorial - <b>vitis-[1:7]</b>
		System Optimizatin - Kernel	§	§	Xilinx HLS Design - xdesign-[1:15]	UCSD Lab <b>ucsd-[1:5]</b>
		FPGA - DSP	§	§		Cornell - ECE5775 <b>cornell-[1:4]</b>
		FPGA - Interconnect	§	§		
4	13-Nov	Kernel Optimization - Area	§	§	Vitis Tutorial - vitis-[1:7]	<b>pp4fpga-[1:8]</b>
		Kernel Optimization - Latency	§	§	UCSD Lab <b>ucsd-[1:5]</b>	Xilinx HLx Examples <b>xhls-[1:18]</b>
		Kernel Optimization - Pipeline	§	§	Cornell - ECE5775 <b>cornell-[1:4]</b>	
5	27-Nov	Design Examples	§	§	<b>pp4fpga-[1:8]</b>	
		Application Cases	§	§	Xilinx HLx Examples <b>xhls-[1:18]</b>	Xilinx Application Notes <b>xapp-[1:13]</b>
6	11-Dec				All topics	Final Project
						Refer to project resource weight=10
	22-Jan	Final Project presentation				

# Lectures – Self-Paced

## 1. Tools & Platform

- a. Introduction to PYNQ & Lab2
- b. Vitis OpenCL XRT and Lab3

## 2. FPGA (Xilinx)

- a. Introduction to FPGA
- b. FPGA – CLB
- c. FPGA – Memory
- d. FPGA – DSP
- e. FPGA – Interconnect

## 3. Concept of System Performance and Optimization

- a. Host Optimization
- b. Kernel Optimization

## 1. HLS Development

- a. Introduction to High Level Synthesis
- b. Kernel IO Interface
- c. Kernel Optimization – Area
- d. Kernel Optimization – Latency
- e. Kernel Optimization – Pipeline

## 2. Design Examples and Application

- a. Design Examples
- b. Application Cases

# Platform/Tools for Labs

- Vivado HLS - For Kernel optimization
  - Vivado HLS – C-sim, Synthesis, Co-Sim, IP-generation
  - Analyze resource, latency, timeline/scheduling, waveform
- Pynq (MPSOC - AXI) – Embedded System
  - Run HLS – C-sim, Co-sim, IP-generation
  - Vivado – IP integration, block-design, generate bit-stream
  - Download to Zedboard/PYNQ-Z2 and run Jupyter Notebook
- Vitis & AWS-F1 (FPGA-PCIe) – Cloud Application
  - Run HLS – C-sim, Co-sim, IP-generation
  - Vitis - run SW-emulation, HW-emulation, Bitstream generation
  - Upload to AWS, run application (host code) at host PC
  - Profiling and analyzing application performance

# Xilinx Tools & Exercise

- Exercises/tutorials provided to gain proficiency in design flow
  - Vivado HLS 2019.2
  - Vivado Design Suite 2019.2
  - Xilinx Vitis IDE/Makefile
  - AWS

Refer to “Xilinx Tool Flow.ppt”

# Develop Basic Skill/Tools in the first two weeks

- The following three labs in the first two weeks
  - Lab#1 - Tool installation and Implementation Flow.
  - Lab#2 - Application Acceleration for Embedded System (PYNQ-Zedboard).
  - Lab#3 - Application Acceleration for Cloud Environment (Amazon)

# Lab/Assignment & Submission Criteria

## • Lab/Project reference resources

- Refer to “HLS Lab Project Resources.xls” <https://cool.ntu.edu.tw/courses/3773/modules/items/110289>

## • Weight Categories

Weight	Description	Submission
1,2	Single item, exercise optimization pragma, coding style, setup/running/analysis effort: 30min-1hr	1. Screen dump: HLS, latency, resource, io interface, timeline
3,4	code hoist, Exercise multiple optimization, comparative analysis, effort: 2-3 hr	2. Vitis summary, HLS synthesis_report
5-9	algorithm level: code hoist, comparative analysis, effort: days	3. ppt/word: description of observation and learning
10	application level, need domain knowledge/background, effort: weeks Candidate for final team project	1. <b>ppt &amp; presentation</b> - introduce domain knowledge/theorem - optimization method - comparison of optimization merit / tradeoff 2. <b>github submission for publication</b>

## • Other Lab/Project proposal is welcomed. Weight will be assigned.

### Lab/Project References

Optimization	Ref	Topic
p204p	1	<a href="#">https://github.com/ChuanHe/CoolHLS/blob/master/README.md</a>
	2	CDCM
	3	DOT
	4	Learn Vector
	5	Matrix Multiplication
VLSI And/or Course Lab	1	<a href="#">https://github.com/ChuanHe/CoolHLS/blob/master/README.md</a>
	2	DOT
	3	codehoist
	4	dotnet_refactor
	5	project
Vitis Application Notes	1	<a href="#">https://www.xilinx.com/en/learning_centers/whitepapers/2019/07/2019-07-16-vitis-application-notes.html</a>
	2	ANIP100 - Matrix Multiplication for Neural Network
	3	ANIP101 - Sparse Matrix-Vector Multiplies on Xilinx
	4	ANIP102 - Linear Algebra: Optimal Flow Algorithm
	5	ANIP103 - Matrix Chain Order Problem and the Best Convolver
	6	ANIP104 - Digital Top Converter
	7	ANIP105 - Image Approximation Color
	8	ANIP106 - Sparse Approximation
	9	ANIP107 - Sparse Approximation
	10	ANIP108 - Sparse Approximation
Corral - EC3775	1	<a href="#">https://github.com/ChuanHe/CoolHLS</a>
	2	CDCM
	3	DOT
	4	codehoist
	5	dotnet_refactor
Vitis Benchmark	1	<a href="#">https://github.com/ChuanHe/CoolHLS</a>
	2	codehoist
	3	DOT
	4	codehoist
	5	dotnet_refactor
VLSI Vitis Decoder	1	<a href="#">https://github.com/ChuanHe/CoolHLS/blob/master/README.md</a>
	2	codehoist
	3	DOT
	4	codehoist
	5	dotnet_refactor
Vitis Lab - SCL	1	<a href="#">https://github.com/ChuanHe/CoolHLS/blob/master/README.md</a>
	2	codehoist
	3	DOT
	4	codehoist
	5	dotnet_refactor
Vitis HLS Examples	1	<a href="#">https://github.com/ChuanHe/CoolHLS/blob/master/README.md</a>
	2	Vitis - Using Vitis
	3	Vitis - Using Vitis
	4	Vitis - Using Vitis
	5	Vitis - Using Vitis
	6	Vitis - Using Vitis
	7	Vitis - Using Vitis
	8	Vitis - Using Vitis
	9	Vitis - Using Vitis
	10	Vitis - Using Vitis
	11	Vitis - Using Vitis
	12	Vitis - Using Vitis
	13	Vitis - Using Vitis
	14	Vitis - Using Vitis
	15	Vitis - Using Vitis
Vitis HLS Coding Style	1	<a href="#">https://github.com/ChuanHe/CoolHLS/blob/master/README.md</a>
	2	codehoist
	3	DOT
	4	codehoist
	5	dotnet_refactor
	6	codehoist
	7	DOT
	8	codehoist
	9	dotnet_refactor
	10	codehoist
	11	DOT
	12	codehoist
	13	dotnet_refactor
	14	codehoist
	Vitis HLS Design	1
2		codehoist
3		DOT
4		codehoist
5		dotnet_refactor
6		codehoist
7		DOT
8		codehoist
9		dotnet_refactor
10		codehoist
11		DOT
12		codehoist
13		dotnet_refactor
14		codehoist
Vitis HLS Design Tutorial		1
	2	codehoist
	3	DOT
	4	codehoist
	5	dotnet_refactor
Vitis Tutorial	1	<a href="#">https://github.com/ChuanHe/CoolHLS/blob/master/README.md</a>
	2	codehoist
	3	DOT
	4	codehoist
	5	dotnet_refactor
Vitis Training Lab	1	<a href="#">https://github.com/ChuanHe/CoolHLS/blob/master/README.md</a>
	2	codehoist
	3	DOT
	4	codehoist
	5	dotnet_refactor



# Lab/Project Resources



LabName	ID	Wt	Topic	
pp4fpga			<a href="https://github.com/KastnerRG/pp4fpgas/tree/master/examples">https://github.com/KastnerRG/pp4fpgas/tree/master/examples</a>	Kastnter pp4fpga text book
	1	3	FIR	
	2	3	CORDIC	
	3	3	DFT	
	4	3	Spare Matrix Vector	
	5	3	Matrix Multiplication	
	6	3	Prefix Sum and Histogram	
	7	3	Video System	
	8	5	Huffman Encoding	
185	h264	10	<a href="https://github.com/adsc-hls/synthesizable_h264">https://github.com/adsc-hls/synthesizable_h264</a>	H.264 Video Decoder
186	swater	10	<a href="https://github.com/necst/coursera-sdaccel-practice">https://github.com/necst/coursera-sdaccel-practice</a>	Smith-Waterman - gene sequence
187	cirrna	10	<a href="https://github.com/necst/circFAXOHW18public">https://github.com/necst/circFAXOHW18public</a>	circular RNA aligner
188	point5	10	<a href="https://bitbucket.org/necst/xohw18_5points_public/src/master/">https://bitbucket.org/necst/xohw18_5points_public/src/master/</a>	five point relative pose problem
189	sha256	5	<a href="https://github.com/dowenberghmark/FPGA-SHA256">https://github.com/dowenberghmark/FPGA-SHA256</a>	SHA256
190	beamf	5	<a href="https://developer.xilinx.com/en/articles/beamforming-acceleration.html">https://developer.xilinx.com/en/articles/beamforming-acceleration.html</a>	beamforming
191	ethash	10	<a href="https://developer.xilinx.com/en/articles/part1-introduction-to-ethash.html">https://developer.xilinx.com/en/articles/part1-introduction-to-ethash.html</a>	blockchain - hashing for Ethereum
192	mcarlo	10	<a href="https://github.com/KitAway/FinancialModels_AmazonF1/tree/master">https://github.com/KitAway/FinancialModels_AmazonF1/tree/master</a>	Monte Carlo financial models
193	profax	10	<a href="https://bitbucket.org/necst/profax-src/src/master/">https://bitbucket.org/necst/profax-src/src/master/</a>	Protein Folding Algorithm
194	graph	10	<a href="https://github.com/Xtra-Computing/ThunderGP">https://github.com/Xtra-Computing/ThunderGP</a>	Graph Processing





# Course Credits

- Earn credits from the followings:
  - Submit Lab/Assignment choose from lab/project references or propose yours
    - Credit based on the weight category
  - Class Presentation – 5 minutes presentation
    - weight category \* quality <0.7 - 1.3> ( **insight** + presentation skill)
  - Final project – 10 minutes presentation
    - Weight <5-10> \* quality <0.7 - 1.3> (**insight** + presentation skill)
- Where is the insight from
  - Fully understand the material
  - Deeper observation on the analysis report
  - Try out different optimization, make trade-off, and comparative analysis