# Contents

# 12   Heterogeneous Computing - DSLs and HLS

## 12.1   Introduction to Heterogeneous Computing

**Heterogeneous Computing** (or **Heterogeneous Processing**) refers to **systems that use multiple types of processors or accelerators to handle different workloads more efficiently**.

- In contrast to traditional homogeneous systems (which only use CPUs), heterogeneous systems combine different processing units such as CPUs, GPUs, DSPs, and FPGAs.

- The **goal is to match the right processor to the right task**, achieving higher performance and energy efficiency.

> **Example 1: Heterogeneous Processing**
>
> A self-driving car requires CPUs for decision-making, GPUs for image recognition, and FPGAs for real-time sensor fusion.

### 🎛 Energy-Efficient Computing Strategies

When designing a heterogeneous system, performance isn't the only goal; **energy efficiency is just as critical**. Given a fixed power budget, simply **increasing performance without considering power constraints is inefficient**. Specialized hardware (e.g., FPGAs, ASICs) achieves better performance per watt than general-purpose processors.

There are two main strategies for improving energy efficiency:

1. **Use Specialized Processors**. **CPUs are not energy-efficient** due to instruction decoding, branch handling, and pipeline management overhead. Specialized hardware (FPGAs, ASICs) reduces overhead, leading to more computations per joule.

$$\text{Power} = \frac{\text{Op}}{\text{second}} \times \frac{\text{Joules}}{\text{Op}}$$

2. **Minimize Data Movement**. **Memory access consumes more energy than computation!** Optimizing data locality reduces power consumption. For example, moving computation closer to memory (e.g., using tensor core inside GPUs) significantly reduces energy cost.

## 12.2   Heterogeneous parallel programming

⚠ **Challenges of Writing Portable and Efficient Parallel Code**

Writing parallel programs for heterogeneous systems is difficult due to the following reasons:

1. **Diverse Hardware Architectures**. A CPU, GPU, and FPGA all have different programming models. **Code written for one hardware type may not perform well on another**.

2. **Performance vs. Productivity Trade-offs**.

   - **Performance**: Low-level programming (e.g., CUDA, OpenCL, Verilog) allows fine-tuned optimizations but **is hard to program**.
   - **Productivity**: High-level abstractions (e.g., OpenMP, DSLs) improve productivity but **may introduce performance overhead**.

3. **Memory Management**. Different memory models (shared vs. distributed) require different optimizations. Data movement between CPU and GPU memory can be costly if not handled efficiently.

4. **Scalability Issues**. Some **programs scale well on GPUs but poorly on CPUs** due to synchronization and memory bandwidth limitations.

✅ **The Ideal Parallel Programming Language**

An **ideal parallel programming model should provide a balance of**:

- ✔ **Performance**. Optimized execution across different hardware.

- ✔ **Productivity**. Easy to use and develop.

- ✔ **Generality**. Works across different architectures.

However, **most existing languages optimize only one or two** of these factors, leading to trade-offs.

| Approach | Performance | Productivity | Generality |
|---|---|---|---|
| **CUDA/OpenCL** | ✔ High | ✘ Low | ✘ Low |
| **OpenMP (CPU)** | ✔ High | ✔ Medium | ✘ Low |
| **MPI (Distributed)** | ✔ High | ✘ Low | ✔ High |
| **FPGA/Verilog/VHDL** | ✔ Very High | ✘ Very Low | ✘ Low |
| **High-Level Synthesis** | ✔ High | ✔ Medium | ✘ Low |

❓ **Why is this important?**

If we want **portable parallel programs**, we need **new high-level abstractions** like Domain-Specific Languages (DSLs), which will be covered in the next section.

## 12.3   DSLs and Halide

❷ **What are Domain-Specific Languages (DSLs)?**

A **Domain-Specific Language (DSL)** is a **specialized programming language** designed for a **specific application domain**. The main characteristics of DSLs are:

- **Restricted expressiveness** (focused on a single domain)

- **High-level, declarative syntax** (easier than general purpose languages)

- **Optimized performance** for the target domain

- **May be standalone or embedded** in another language

| DSL Name | Target Domain | Key Benefits |
|---|---|---|
| **Halide** | Image Processing | Separates algorithm from scheduling for optimized execution. |
| **TensorFlow** | Machine Learning | Optimized computation graphs for AI workloads. |
| **SQL** | Databases | Declarative queries for efficient data retrieval. |
| **Verilog/VHDL** | Hardware Design | Describes digital circuits for synthesis. |

Table 11: Examples of DSLs.

⚖ **Embedded vs. External DSLs**

DSLs can be classified as:

- **External** DSLs:

    🔖 Have **their own** *syntax* and *compiler/interpreter*.

    ❷ **Example**: SQL, Halide, Verilog.

    ✅ **Advantages**: can be **more optimized** but require custom compilers.

- **Embedded** DSLs:

    🔖 **Built inside another general-purpose language**.

    ❷ **Example**: TensorFlow (embedded in Python).

    ✅ **Advantages**: benefit from **integration** with the host language

## ⚡ DSL Use Case: Halide for Image Processing

**Halide** is a **Domain-Specific Language (DSL) for high-performance image processing**.

### ❓ Why does image processing need DSLs?

- 🔋 Image processing is **data-intensive** and <mark>requires high performance</mark>.
- 👎 <mark>Traditional solutions</mark> (`C++`, `CUDA`, `OpenCV`) <mark>require manual optimizations</mark>.
- ☹ Optimizing code for **parallelism and memory efficiency** is **difficult**.

### ❓ Why Halide?

- ✂ **Separates "*what*" is computed from "*how*" it is executed**.
- </> **Expresses computations at a high level**, leaving optimizations to the compiler.
- 🌐 **Portable** across CPUs, GPUs, and FPGAs.

## 🔧 How Halide Works: Separating Algorithm from Schedule

In Halide, a **key feature is the separation** of **what** a program computes (*computation/algorithm*) from **how** it executes (*schedule*). This means:

- The **algorithm** specifies **what operations should be performed**. In other words, specifies **what to compute** (like a mathematical formula).

- The **schedule** defines **how those operations should be executed efficiently** on the hardware. In other words, specifies **how to execute the computation** (parallelism, memory layout, vectorization).

Now we see the difference between the traditional approach we have always used and the halide approach:

- **Traditional Approach (`C++`, `CUDA`, `OpenCV`)**. In traditional programming languages (e.g., `C++`, `OpenCV`, `CUDA`), the **algorithm and execution strategy are mixed together**. This means that if we want to **change parallelization or memory access optimizations**, we must **rewrite parts of the algorithm itself**. This makes it *hard to experiment* with different optimizations.

  ### ⚠ Problems with the Traditional Approach

  1. **If we want to optimize** for vectorization, parallel execution, or memory layout, we <mark>must modify the algorithm itself</mark>.
  2. The <mark>same code cannot easily be reused</mark> for different architectures (e.g., CPU, GPU, FPGA).

**Example 2: Problems with the Traditional Approach**

```
1  void box_blur(const Image &in, Image &out) {
2      for (int y = 1; y < in.height() - 1; y++) {
3          for (int x = 1; x < in.width() - 1; x++) {
4              out(x, y) = (
5                  in(x-1, y) + in(x, y) + in(x+1, y)
6              ) / 3;
7          }
8      }
9  }
```

The problem here is that we need to specify the "*what*", i.e. what operations should be performed, but also "*how*" these operations should be performed.

- **Halide's Approach: Separate Algorithm from Execution**. Halide splits the computation into two parts:

  1. **Computation** (Algorithm) - *What to Compute*
     - Defines the mathematical computation.
     - Remains unchanged across different hardware targets.
  2. **Schedule** - *How to Execute Efficiently*
     - Controls memory layout, parallelization, and optimization.
     - Can be changed without modifying the algorithm.

**Example 3: Box Blur in Halide**

The computation part is separated from the scheduling! So a change in the algorithm can be made without affecting the control of the execution. Therefore, we define simply:

1. Computation (Algorithm), stays the same:

```
1  Var x, y;
2  Func blurx, blury;
3
4  // First pass: horizontal blur
5  blurx(x, y) = (in(x-1, y) + in(x, y) + in(x+1, y))
       / 3;
6
7  // Second pass: vertical blur
8  blury(x, y) = (blurx(x, y-1) + blurx(x, y) + blurx(
       x, y+1)) / 3;
```

This part **only describes the math**, <u>not how</u> it should run.

2. Schedule, controls execution (can be changed easily):

```
1  blury.tile(x, y, xi, yi, 256, 32)
2      // Vectorized execution for SIMD
3      .vectorize(xi, 8)
4      // Parallel execution over y-dimension
5      .parallel(y);
```

```
6
7  blurx.compute_at(blury, x)   // Compute blurx only
       when needed by blury
8      .vectorize(x, 8);
```

This part **controls execution strategy** but does **not
modify the algorithm**. The **same algorithm** can now
**run efficiently on different hardware architectures** just
by changing the schedule.

✅ **Why DSLs Matter for Performance and Productivity: Advantages**

✔ **Performance Optimization**. A Halide program **can be better than
hand-optimization `C++` code**. Scheduling decisions affect **parallel ex-
ecution, memory locality, and vectorization**.

✔ **Productivity**. Instead of manually optimizing, **Halide allows rapid
exploration of different schedules**. Easier to **port to different ar-
chitectures** (CPU, GPU, FPGA).

In conclusion, DSLs like Halide **automate low-level optimizations**, enabling
*faster* and *more efficient* code for specialized domains.

## 12.4 Scheduling & Performance Optimization in Halide

This section focuses on **how different scheduling strategies in Halide affect performance** and how to **choose the best schedule for different hardware architectures**.

### ❷ Why Scheduling Matters

Scheduling is the **key to optimizing performance in Halide**.

- A **poorly scheduled program can be $10\times$ slower than an optimized one**.

- **Memory access, cache locality, parallelism, and vectorization** all depend on scheduling.

As we have already seen, in traditional languages (C++, CUDA, OpenMP), scheduling decisions must be hard-coded into the algorithm. In Halide, the schedule is separate and can be changed without changing the algorithm.

### ❷ How Scheduling Affects Performance

Different **scheduling strategies** impact how the computation is executed on hardware:

| Scheduling Strategy | Impact |
|---|---|
| **Serial Execution** | Simple, but slow. No parallelism. |
| **Parallel Execution** | Uses multiple CPU cores. Good for multi-core CPUs. |
| **Vectorization (SIMD)** | Uses wide registers for efficiency (e.g., AVX). |
| **Tiling** | Improves cache locality by processing data in chunks. |
| **Compute-at** | Controls when intermediate results are computed. |
| **Store-at** | Controls where intermediate results are stored. |

### ⠇☰ Scheduling Strategies in Halide

Let's take the Box Blur Algorithm as an example. The **Box Blur** is a simple **image processing technique used for smoothing or blurring an image**. It performs two main steps:

1. **Each pixel** in the output image is computed as the **average of its neighboring pixels**.

2. It applies a **moving average filter** over a **small window** (e.g., $3 \times 3$ or $5 \times 5$ pixels).

It is used to reduce noise in images, to create a smooth, blurry effect and, above all, because it is fast and efficient, since it involves only two operations: adding and dividing.

However, we will now analyze the algorithm implemented in Halide using different strategies:

⧗ **Default Serial Execution (Slow)**. Without scheduling, Halide will execute in **serial order** (one pixel at a time).

```
1  Var x, y;
2  Func blurx, blury;
3
4  blurx(x, y) = (in(x-1, y) + in(x, y) + in(x+1, y)) / 3;
5  blury(x, y) = (blurx(x, y-1) + blurx(x, y) + blurx(x, y+1)) /
       3;
```

     ✖ **Problems**

        ✖ No parallelism or vectorization.

        ✖ Poor memory access patterns.

        ✖ Slow execution.

⏱ **Parallel Execution**. We can **add parallelism** to use multiple CPU cores:

```
1  blury.parallel(y);
```

     ✔ **Advantages**

        ✔ Halide automatically splits the work across CPU cores.

        ✔ Useful for multi-threaded execution on CPUs.

⏱ **Vectorization (SIMD)**. Modern CPUs support **SIMD instructions** (e.g., AVX, NEON) to process multiple pixels at once:

```
1  blury.vectorize(x, 8);
```

     ✔ **Advantages**

        ✔ Uses **SIMD registers** for faster execution.

        ✔ Works best for **data-parallel workloads** like image processing.

⏱ **Tiling for Better Cache Performance**. Instead of processing the whole image at once, we **divide it into smaller tiles**:

```
1  blury.tile(x, y, xi, yi, 256, 32);
```

     ✔ **Advantages**

        ✔ **Each tile fits better in cache**, reducing memory latency.

✔ **Improves locality of reference** (less cache thrashing).

⚡ **Optimized Schedule: Combining Techniques**.  We can **combine multiple scheduling strategies** for maximum performance:

```
1  // Process in 256x32 tiles
2  blury.tile(x, y, xi, yi, 256, 32)
3      // Vectorized execution
4      .vectorize(xi, 8)
5      // Parallel execution across CPU cores
6      .parallel(y);
7
8  // Compute blurx only when needed
9  blurx.compute_at(blury, x)
10     .vectorize(x, 8);
```

✅ **Advantages**

✔ Breaks image into tiles for cache efficiency.

✔ Uses SIMD vectorization for fast execution.

✔ Runs in parallel on multiple CPU cores.

✔ Intermediate results (`blurx`) are computed only when needed.

⚖ **Trade-offs in Scheduling**

But *how do we choose the right scheduling?*  Well, we need to find a good trade-off. Different scheduling choices affect **performance trade-offs**:

| Scheduling Strategy | Performance Impact |
|---|---|
| **Parallel Execution** | Increases throughput, uses multiple cores. |
| **Vectorization (SIMD)** | Improves performance on CPUs/GPUs. |
| **Tiling** | Improves cache locality, reduces memory overhead. |
| **Compute-at** | Avoids redundant computations. |
| **Store-at** | Reduces memory footprint but increases recomputation. |

Table 12: Performance trade-offs in Scheduling.

✔ **Conclusion**

In conclusion, the key takeaways are:

- **Scheduling** is the **key to performance** in Halide.

- **Parallel execution**, **vectorization**, and **tiling** significantly improve performance.

- **Halide's flexibility** allows **quick** experimentation with **different schedules**.

- The **right schedule depends on hardware constraints** (CPU, GPU, FPGA).