# Support of C++ Compiler for Embedded Multicore DSP Systems

Jenq Kuen Lee      Chi-Bang Kuan

Department of Computer Science

National Tsing Hua University
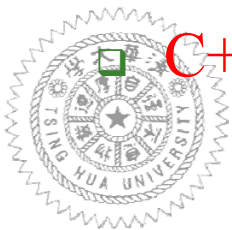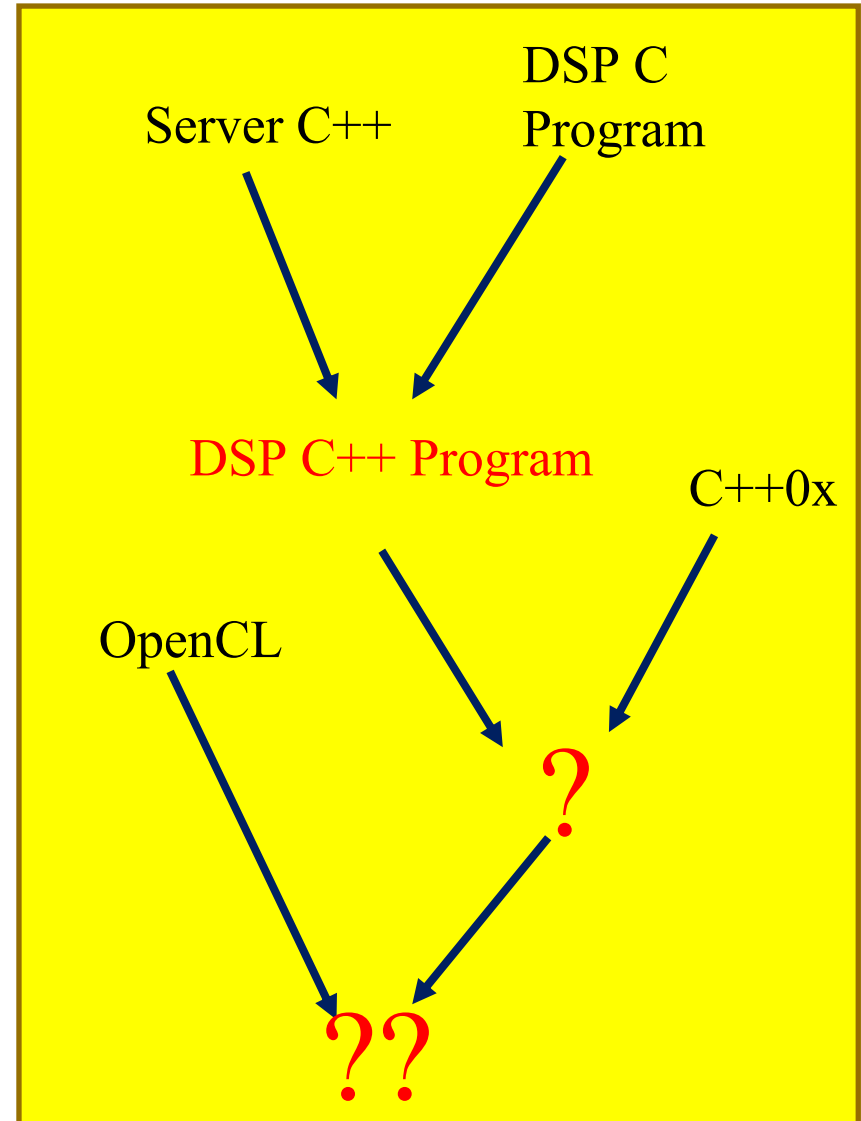
HsinChu, Taiwan

# Outline

- Motivations
- Case Study: DSP C++ compiler
  - SIMD classes to hide DSP intrinsic functions
  - Code size issues with C++ libraries for DSP systems
  - Design patterns
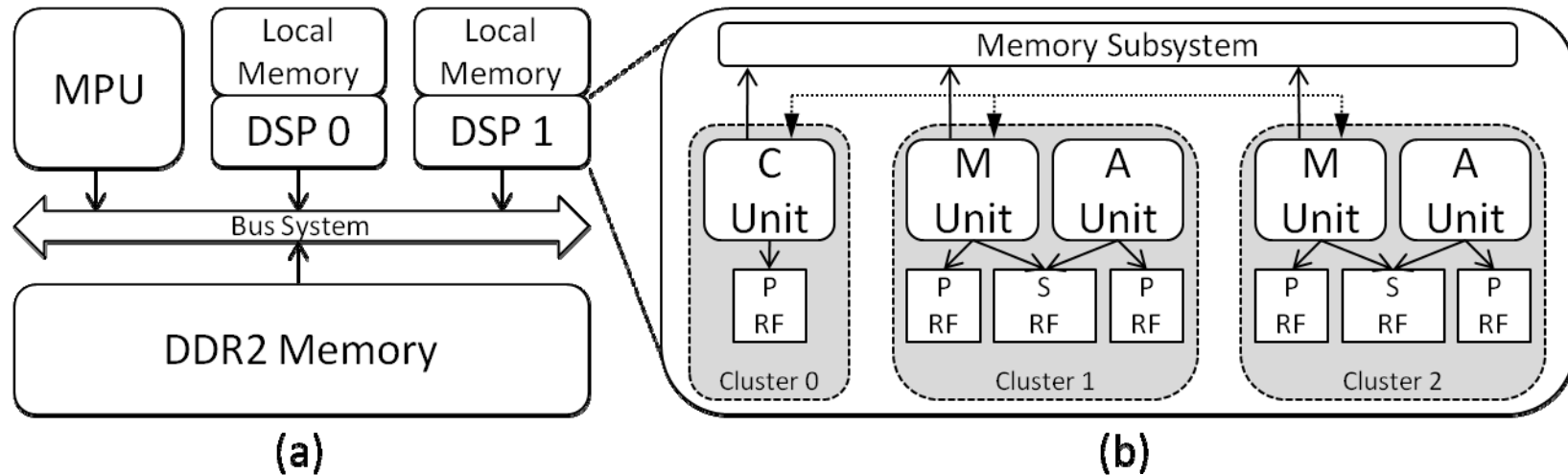- Experimental results
- Discussions

# Why Supporting C++ for Embedded DSPs?

- Increasing numbers of standards provide C++ interfaces.
  - C++ binding of OpenCV 2.0, some Android codes, C++ legacy codes
- To enhance **programming models** by providing high-level abstraction, such as
  - Overloading vectors for SIMD programming.
- To build **design patterns** for specific applications.
  - Parallel design patterns.
  - Client-server pattern/Pipe-filter pattern/ Map-N-reduce patterns
- C++ in Junction

# Case Study for DSP C++ Compiler: Platform
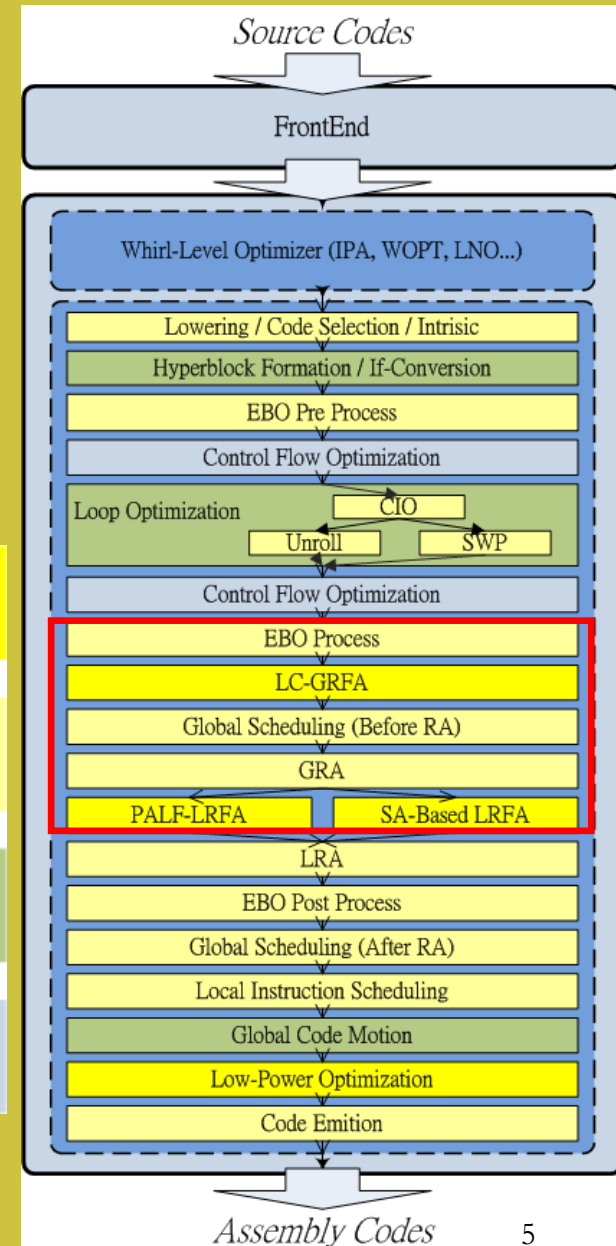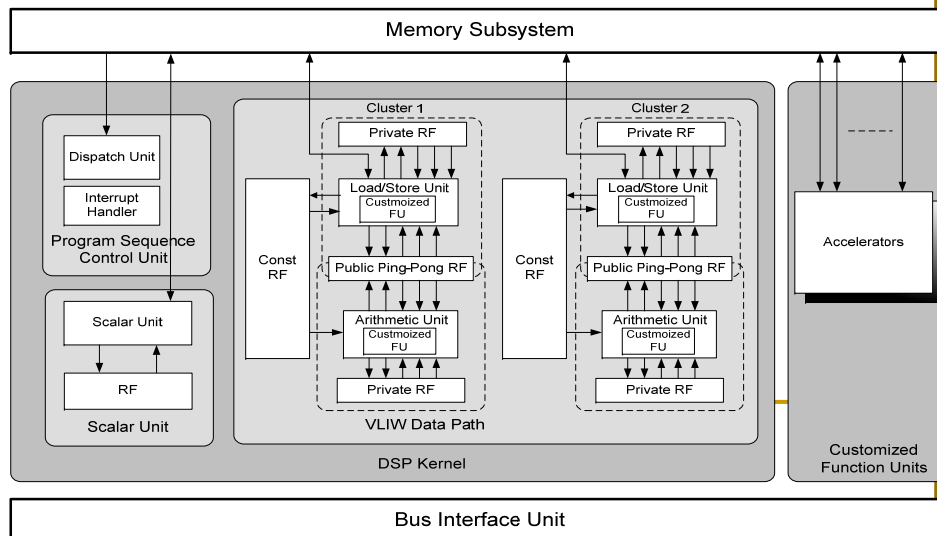


(a)  (b)

- **3 Cores SoC**
  - ARM926EJS x 1
  - PAC DSP x 2
- **Shared memory architecture**

- **ITRI PAC DSP**
  - 5-way issue VLIW
  - 3 clustered processing units
  - Distributed register file
  - SIMD instruction set

# Case Study: PACDSP Compiler

- Based on Open64 compiler
- VLIW DSP compilers for distributed register files
- PALF scheduling policies for ILP (CPC 2006)
- GRA scheme for distributed register files (CPC 2007)
- SIMD compiler optimizations + intrinsics/extrinsics
- Copy propagations for distributed register architectures (LCPC 2006)
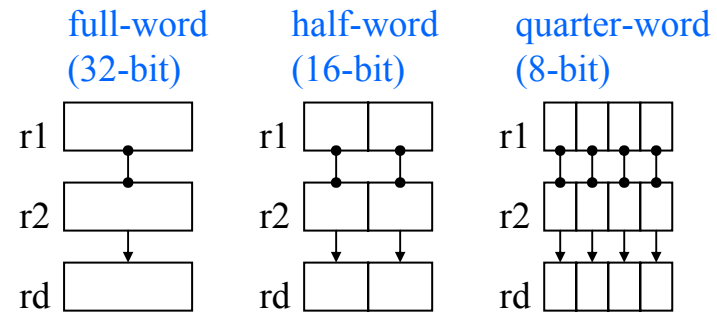- Register spills among distributed register banks (CPC 2009)



Memory Subsystem

Program Sequence Control Unit
- Dispatch Unit
- Interrupt Handler

Scalar Unit
- Scalar Unit
- RF

DSP Kernel

Cluster 1
- Private RF
- Load/Store Unit / Custmoized FU
- Public Ping-Pong RF
- Arithmetic Unit / Custmoized FU
- Private RF

Const RF

Cluster 2
- Private RF
- Load/Store Unit / Custmoized FU
- Public Ping-Pong RF
- Arithmetic Unit / Custmoized FU
- Private RF

VLIW Data Path

Accelerators

Customized Function Units

Bus Interface Unit



Source Codes

FrontEnd

Whirl-Level Optimizer (IPA, WOPT, LNO...)

Lowering / Code Selection / Intrisic
Hyperblock Formation / If-Conversion
EBO Pre Process
Control Flow Optimization
Loop Optimization — CIO — Unroll — SWP
Control Flow Optimization
EBO Process
LC-GRFA
Global Scheduling (Before RA)
GRA
PALF-LRFA — SA-Based LRFA
LRA
EBO Post Process
Global Scheduling (After RA)
Local Instruction Scheduling
Global Code Motion
Low-Power Optimization
Code Emition

New Phase for PAC DSP

Specially Tuned for PAC DSP

Ported for Target Dependency

Original Phases of OPEN64

Assembly Codes

5

# SIMD Features of PACDSP

- ## Subword instructions
  - short vector instructions that operate data in separated parts

    | | | |
    |---|---|---|
    | *add* | → | one 32-bit add |
    | *add.d* | → | two 16-bit add |
    | *add.q* | → | four  8-bit add |

full-word (32-bit)  half-word (16-bit)  quarter-word (8-bit)

r1  r2  rd

- ## Permutation instructions
  - diverse instructions for reordering data in registers

swap2 rd, r1    unpack2 rd, r1    permh2 rd, r1, r2, #1, #2

r1 | H | L |    r1 | H | L |    r1 | 0 | 1 |  r2 | 2 | 3 |

rd | L | H |    rd+1 | H |  rd | L |    rd

- ## Parallel instruction issuing
  - issue up to two-multiply or five-other operations per cycle

{ Scalar , LSU ,  ALU , LSU ,  ALU }
{ nop  , nop , mpy , nop , mpy }
{ add  , add , add , add , add }

6

# SIMD Intrinsics for PACDSP

- Summary of SIMD intrinsics

| *intrinsics* | *subdivisions* | | *instances in prototypes* |
|---|---|---|---|
| subword intrinsics | | full-word | int __builtin_add (int, int)  … |
| | | half-word | int __builtin_add_d (int, int)  … |
| | | quarter-word | int __builtin_add_q (int, int)  … |
| cluster intrinsics | | cluster 1 | int __builtin_c1_add (int, int)<br>int __builtin_c1_add_d (int, int)<br>int __builtin_c1_add_q (int, int)  … |
| | | cluster 2 | int __builtin_c2_add (int, int) … |

- All intrinsics are mapped to DSP instructions
  - subword intrinsics
    - __builtin_add_d → add.d (a dual 16-bit add)
    - __builtin_add_q → add.q (a quad 8-bit add) …
  - cluster intrinsics (*also have subword semantics*)
    - __builtin_c1_add_d → add.d (a dual 16-bit add at cluster 1)
    - __builtin_c2_add_d → add.d (a dual 16-bit add at cluster 2) …

# SIMD Intrinsic Programming: The LMS Filter in DSPstone

**1.1 identify parallel regions or parallelize programs**

```
1
2      short H[N], X[N];
3      short *p_H, *p_X, *p_X2;
4      int f, y=0, d=7, delta=1, error;
5
6      /* state update and dot-product */
7      for(f = 1; f < N; f++)
8        y += H[N-f] * (X[N-f] = X[N-f-1]);
9      y += H[0] * (X[0] = x);
10
11     /* coefficient update */
12     error = (d - y) * delta;
13     for (f = 0; f < N; f++)
14       H[f] += error * X[f];
15
```

Quick summary:
1. The data type of input arrays is 16-bit integer
2. There are two loops, one for dot-product and one for coefficient-update
3. The dot-product is not parallel, while the coefficient-update is parallel.

8

# SIMD Intrinsic Programming: The LMS Filter with SIMD intrinsics

**1. find the best way to process data**

```
1    short __attribute__ ((aligned(4))) H[N];
2    short __attribute__ ((aligned(4))) X[N];
3    int *vX = (int *)X, *vH = (int *)H;
4    int v_err=0, sum1=0, sum2=0, sum3=0, sum4=0;
5
6    for(f = 1, X[-1] = x; f < (N/2)+1; f+=4)
7    {
8      vX[(N/2)-f] = __builtin_c1_permh2(vX[(N/2)-f-1], vX[(N/2)-f], 1, 2);
9      sum1 = __builtin_c1_mac_d(sum1, vH[(N/2)-f], vX[(N/2)-f]);
10     vX[(N/2)-f-1] = __builtin_c1_permh2(vX[(N/2)-f-2], vX[(N/2)-f-1], 1, 2);
11     sum2 = __builtin_c1_mac_d(sum2, vH[(N/2)-f-1], vX[(N/2)-f-1]);
12
13     vX[(N/2)-f-2] = __builtin_c2_permh2(vX[(N/2)-f-3], vX[(N/2)-f-2], 1, 2);
14     sum3 = __builtin_c2_mac_d(sum3, vH[(N/2)-f-2], vX[(N/2)-f-2]);
15     vX[(N/2)-f-3] = __builtin_c2_permh2(vX[(N/2)-f-4], vX[(N/2)-f-3], 1, 2);
16     sum4 = __builtin_c2_mac_d(sum4, vH[(N/2)-f-3], vX[(N/2)-f-3]);
17   }
18   y = __builtin_c1_mergea(__builtin_c1_add_d(sum1, sum2)) +
19       __builtin_c2_mergea(__builtin_c2_add_d(sum3, sum4));
20
21   error = (d - y) * delta;
22   v_err= __builtin_permh2(error, error, 0, 0);
23   for (f = 0; f < (N/2); f+=2)
24   {
25     vH[f] = __builtin_c1_mac_d(vH[f], vX[f], v_err);
26     vH[f+1] = __builtin_c2_mac_d(vH[f+1], vX[f+1], v_err);
27   }
```

9

# SIMD Intrinsic Programming:
# The LMS Filter with SIMD intrinsics

**2. prepare packed data for subword processing**

```
1   short __attribute__ ((aligned(4))) H[N];
2   short __attribute__ ((aligned(4))) X[N];
3   int *vX = (int *)X, *vH = (int *)H;
4   int v_err=0, sum1=0, sum2=0, sum3=0, sum4=0;
5
6   for(f = 1, X[-1] = x; f < (N/2)+1; f+=4)
7   {
8    vX[(N/2)-f] =  __builtin_c1_permh2(vX[(N/2)-f-1], vX[(N/2)-f], 1, 2);
9    sum1 = __builtin_c1_mac_d(sum1, vH[(N/2)-f], vX[(N/2)-f]);
10   vX[(N/2)-f-1] = __builtin_c1_permh2(vX[(N/2)-f-2], vX[(N/2)-f-1], 1, 2);
11   sum2 = __builtin_c1_mac_d(sum2, vH[(N/2)-f-1], vX[(N/2)-f-1]);
12
13   vX[(N/2)-f-2] = __builtin_c2_permh2(vX[(N/2)-f-3], vX[(N/2)-f-2], 1, 2);
14   sum3 = __builtin_c2_mac_d(sum3, vH[(N/2)-f-2], vX[(N/2)-f-2]);
15   vX[(N/2)-f-3] = __builtin_c2_permh2(vX[(N/2)-f-4], vX[(N/2)-f-3], 1, 2);
16   sum4 = __builtin_c2_mac_d(sum4, vH[(N/2)-f-3], vX[(N/2)-f-3]);
17  }
18  y = __builtin_c1_mergea(__builtin_c1_add_d(sum1, sum2)) +
19      __builtin_c2_mergea(__builtin_c2_add_d(sum3, sum4));
20
21  error = (d - y) * delta;
22  v_err= __builtin_permh2(error, error, 0, 0);
23  for (f = 0; f < (N/2); f+=2)
24  {
25    vH[f] = __builtin_c1_mac_d(vH[f], vX[f], v_err);
26    vH[f+1] = __builtin_c2_mac_d(vH[f+1], vX[f+1], v_err);
27  }
```

# SIMD Intrinsic Programming:
# The LMS Filter with SIMD intrinsics

**3. utilize subword instructions and distribute computation**

```
1    short __attribute__ ((aligned(4))) H[N];
2    short __attribute__ ((aligned(4))) X[N];
3    int *vX = (int *)X, *vH = (int *)H;
4    int v_err=0, sum1=0, sum2=0, sum3=0, sum4=0;
5
6    for(f = 1, X[-1] = x; f < (N/2)+1; f+=4)
7    {
8      vX[(N/2)-f] = __builtin_c1_permh2(vX[(N/2)-f-1], vX[(N/2)-f], 1, 2);
9      sum1 = __builtin_c1_mac_d(sum1, vH[(N/2)-f], vX[(N/2)-f]);
10     vX[(N/2)-f-1] = __builtin_c1_permh2(vX[(N/2)-f-2], vX[(N/2)-f-1], 1, 2);
11     sum2 = __builtin_c1_mac_d(sum2, vH[(N/2)-f-1], vX[(N/2)-f-1]);
12
13     vX[(N/2)-f-2] = __builtin_c2_permh2(vX[(N/2)-f-3], vX[(N/2)-f-2], 1, 2);
14     sum3 = __builtin_c2_mac_d(sum3, vH[(N/2)-f-2], vX[(N/2)-f-2]);
15     vX[(N/2)-f-3] = __builtin_c2_permh2(vX[(N/2)-f-4], vX[(N/2)-f-3], 1, 2);
16     sum4 = __builtin_c2_mac_d(sum4, vH[(N/2)-f-3], vX[(N/2)-f-3]);
17   }
18   y = __builtin_c1_mergea(__builtin_c1_add_d(sum1, sum2)) +
19       __builtin_c2_mergea(__builtin_c2_add_d(sum3, sum4));
20
21   error = (d - y) * delta;
22   v_err= __builtin_permh2(error, error, 0, 0);
23   for (f = 0; f < (N/2); f+=2)
24   {
25     vH[f] = __builtin_c1_mac_d(vH[f], vX[f], v_err);
26     vH[f+1] = __builtin_c2_mac_d(vH[f+1], vX[f+1], v_err);
27   }
```

11

# SIMD Programming with C++

**Intrinsic programming is effective, though it makes code not portable and hard to read. Any solution?**

- ❑ Use C++ classes to encapsulate SIMD details in vector data types.

```
class short4
{
  int subword1;
  int subword2;

public:
  ...

};
```

```
inline short4 operator+(
    const short4& opnd1, const short4& opnd2)

{

  short4 result(
    _dsp_c1_add_d(opnd1.subword1, opnd2.subword1),
    _dsp_c2_add_d(opnd1.subword2, opnd2.subword2));

  return result;

}                              /* C++ operator overloading */
```

- ❑ Seamless integration with SIMD intrinsics to provide vector abstraction.
- ❑ Better portability and readability with similar performance (see examples and performance in next slides).

# More C++ Classes for DSPs

- **Emulated long vectors** with cluster intrinsics to encourage programmers to utilize VLIW data paths.
  - ❏ E.g. emulated **short4** types for 32-bit DSP processors
  - ❏ Implicit computation distribution over VLIW clusters.
- **Miscellaneous data types** to facilitate signal processing
  - ❏ Integers with saturation arithmetic.
  - ❏ Fractional numbers with different rounding modes.
  - ❏ Complex numbers.

|  | **Vector Types** | **Generic Vector Operations** |
|---|---|---|
| Native subword vectors | [u]char4,  [u]short2 | ■ arithmetic:      +, -, *, >>, <<, … <br> ■ reduction:      sum, dot_product, min, max, … |
| Emulated full-word vectors | [u]char8,  [u]short4, <br> [u]int2,      [u]int4, <br> cmplx,      cmplx2 | ■ load/store:    lw, sw <br> ■ access:        [ ] <br> ■ permutation:  perm, pack, unpack |

# Convolution Kernels with SIMD
## intrinsics vs. vector classes

```
/* Convolution with SIMD Intrinsics */

ushort coef [BLK_SZ];

int convolve ( uchar src[BLK_SZ*4] ) {
    int ycbcr, y, cbcr, COEF, Y, CBCR;
    __m64 tmp64;

    for( int i=0; i<BLK_SZ; i++ ) {
        ycbcr = ((int*)src)[i];

        tmp64 = _dsp_unpack4u(ycbcr);
        y     = _dsp_get64l (tmp64);
        cbcr  = _dsp_get64h(tmp64);

        COEF = _dsp_permh2(coef[i], 0, 0, 0);
        Y    = _dsp_mac_d(Y, y, COEF);
        CBCR = _dsp_mac_d(CBCR, cbcr, COEF);
    }
    return a pixel result with Y and CBCR;
```

- direct access to DSP instructions
- platform-dependent
- expose users to hardware details

- compact representation
- with similar performance

```
/* Convolution with SIMD Vectors */

ushort coef [BLK_SZ];

int convolve ( uchar src[BLK_SZ*4] ) {
    uchar4 ycbcr;
    ushort2 y, cbcr, CBCR, Y;

    for( int i=0; i<BLK_SZ; i++ ) {
        // load a char4 vector from &src[i]
        ycbcr.lw(&src[i*4]);

        // unpack a uchar4 into two ushort2
        ycbcr.unpack(y, cbcr);

        // vector operations on short2 vectors
        Y    += y * coef[i];
        CBCR += cbcr * coef[i];
    }
    return a pixel result with Y and CBCR;
}
```

- Green functions are subword intrinsics for accessing DSP instructions.
- Blue data types are SIMD vectors; red operators are vector operations.
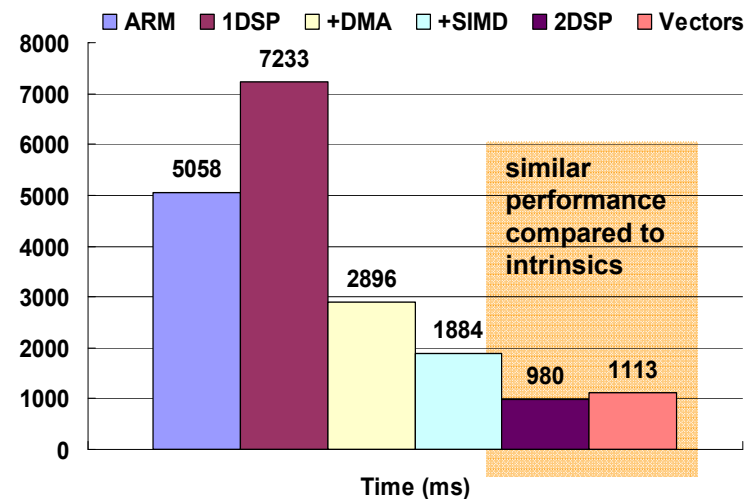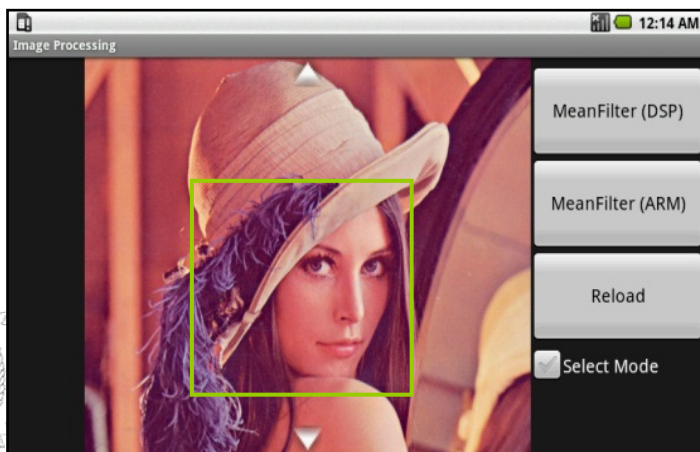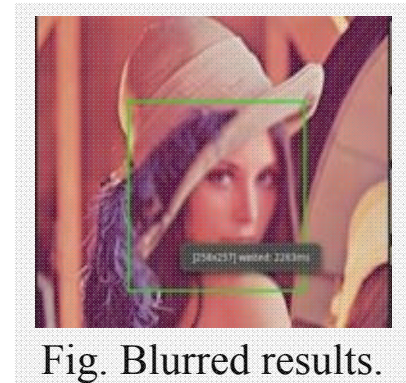
# Applications with SIMD Vectors

- **Image blurring**
  - Blurring an arbitrarily selected block and displaying results via Android UI.
- Major kernels:
  - Color space transformations (rgb2yuv, yuv2rgb).
  - 2D convolution with a 8x8 coefficient matrix (prev page).
- Performance footnotes:
  - ARM:     single-core program with ARM GCC -O3.
  - +DMA:  DMAs are used to facilitate data transmission (2D data) between shared memory and DSP locals.
  - +SIMD: DSP subword instructions are used via SIMD intrinsics.
  - Vectors: SIMD intrinsics are replaced with C++ vector operations.
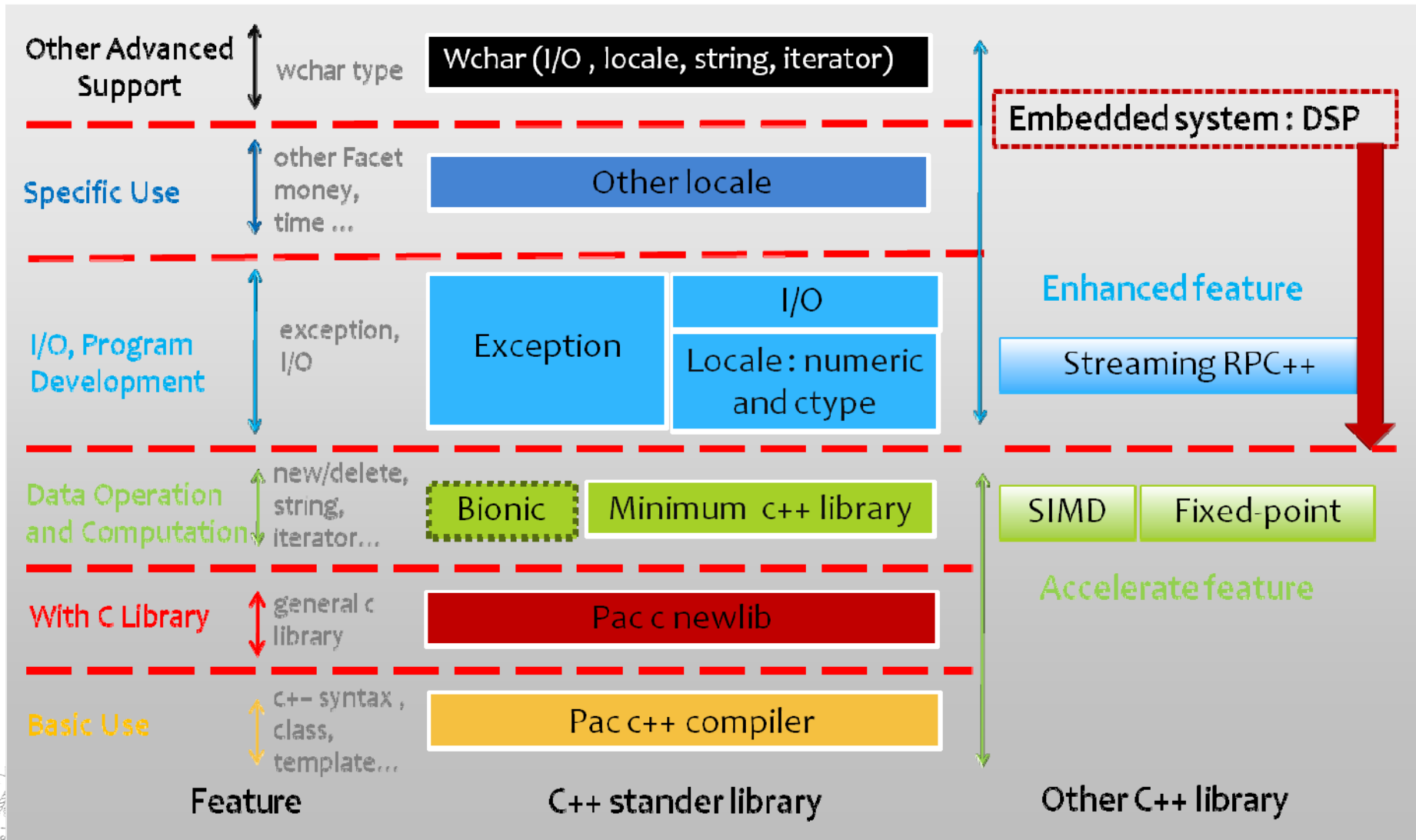


Fig. Blurred results.

# Code Size Issues of C++ Library

- Embedded processors typically have small instruction memory
  - E.g. PACDSP with 16KB~32KB i-cache in cache or SPM mode.

- Size analysis of C++ library
  - Apache C++ library compiled with -O2 and -Os.
  - Apparently, some of the C++ components are inappropriate for embedded.

- C++ also depends on partial C library for basic functions.
  - The overall lib stack is heavy.
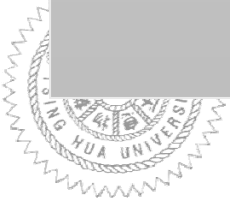  - Doesn't fit instruction cache, or cache miss penalty hurts.

| Component | Code Size | Percentage |
|---|---|---|
| Language | 4 KB | 0.3% |
| Diagnostics | 10 KB | 1.0% |
| Utilities | 3 KB | 0.3% |
| Strings | 63 KB | 6.0% |
| Localizations | 654 KB | 62.4% |
| Containers | 4 KB | 0.3% |
| Iterators | 48 KB | 4.6% |
| Algorithms | 0 KB | 0.0% |
| Numerics | 2 KB | 0.2% |
| Input/Output | 259 KB | 24.7% |
| Total | 1047 KB | 100.0% |

# DSP C/C++ Library Stack and Layers



| Feature | | C++ stander library | Other C++ library |
|---|---|---|---|
| **Other Advanced Support** | wchar type | Wchar (I/O , locale, string, iterator) | Embedded system : DSP |
| **Specific Use** | other Facet money, time … | Other locale | |
| **I/O, Program Development** | exception, I/O | Exception / I/O / Locale : numeric and ctype | Enhanced feature / Streaming RPC++ |
| **Data Operation and Computation** | new/delete, string, iterator... | Bionic / Minimum c++ library | SIMD / Fixed-point |
| **With C Library** | general c library | Pac c newlib | Accelerate feature |
| **Basic Use** | c+- syntax , class, template... | Pac c++ compiler | |

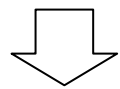**Trade-off between programmability, performance and code size.**

# Design Patterns

- Brief introduction to design patterns
  - Proposed by C. Alexander for city planning and architecture.
  - Introduced to software engineering by Beck and Cunningham.
  - Become prominent in object-oriented programming by GoF.

- Design patterns describe "good solutions" to recurring problems in a particular context.
  - Patterns for **object-oriented programming**
    - Creational patterns, Structural patterns, Behavioral patterns, etc.
  - Patterns for **limited memory systems**
    - Compression, Small data structures, Memory allocation, etc.
  - Patterns for **parallel programming**
    - Finding concurrency, Algorithm structure, Supporting structures and Implementation mechanisms.
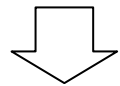
# Parallel Design Patterns

**Parallelization can be a process to transform problems
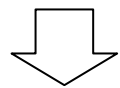to programs by selecting appropriate patterns.**

Design Space of Parallel Patterns

- **Finding Concurrency**

  ⬇ decomposition
     of problems
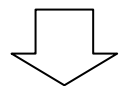
- **Algorithm Structure**

  ⬇ appropriate
     algorithms

- **Supporting Structures**

  ⬇ appropriate
     program constructs

- **Implementation Mechanisms**

  ⬇ parallelized
     programs

- **Decomposition** patterns: {data, task}
- **Dependency analysis** patterns:
  {group tasks, order tasks, data sharing}
- **Design evaluation** pattern

How the given problem is organized?
- By **tasks**: {task parallelism, divide & conquer}
- By **data decomposition**: {geometric, recursive}
- By **flow of data**: {pipeline, event-based coord.}

Software constructs to express parallel algorithms
- **Program structures**: {SPMD, master/worker, loop parallelism, fork-join, client-server, SIMD}
- **Data structures**: {shared data, shared queue, distributed array}

- UE management: {thread/process creation/destr.}
- Synchronization: {barrier, mutex, mem fence}
- Communication: {msg passing, collective comm}

These patterns are summarized from the book,
"Patterns for Parallel Programming" by Mattson et al

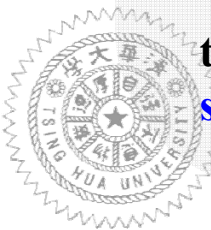# Multicore Programming Design Patterns with C++

- Client-server pattern and data stream abstraction for hiding remote task creation and channel manipulation details.
- Well-defined interface instead of APIs taking function pointers.

```
/* RPC-client program */
class task : public srpc::call_template
{
    CALL_INIT(task);
    void TX() {
        istrm << shape(len, b, s)
            << INPUT << flush;
    }
    void RX() {
        ostrm >> shape(len)
            >> RESULT >> pop;
    }
};

task task_call("task", DSP0);
srpc::srpc_call<task> call(task_call);
```

```
/* RPC-server program */
class task : public srpc::service_template
{
    SERVICE_INIT(task);
    void SERVE() {
        // compute data with istrm.
        // return results via ostrm.
    }
};

int main()
{
    // register a "task" service
    task task_service("task");
    srpc::srpc_server<task>
        server(task_service);
}
```

# Multicore Programming Design Patterns with C++

- **Client-server** pattern and data stream abstraction for hiding remote task creation and channel manipulation details. Well-defined interface in terms of API using function pointers.

> User-defined functor which implements remote call interface

> Client interface which defines essential routines & data structures for RPC

```
/* RPC-client program */
class task : public srpc::call_template
{
    CALL_INIT(task);
    void TX() {
        istrm << shape(len, b, s)
            << INPUT << flush;
    }
    void RX() {
        ostrm >> shape(len)
            >> RESULT >> pop;
    }
};
```

> Encapsulated I/O streams with operators and manipulators to accommodate various data formats

> All RPC invocation details are hidden in srpc::srpc_call

```
task task_call("task", DSP0);
srpc::srpc_call<task> call(task_call);
```
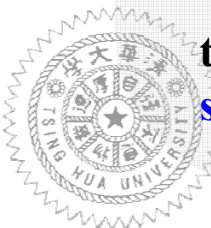
```
/* RPC-server program */
class task : public srpc::service_template
{
    SERVICE_INIT(task);
    void SERVE() {
        // compute data with istrm.
        // return results via ostrm.
    }

int main()
{
    // register a "task" service
    task task_service("task");
    srpc::srpc_server<task>
        server(task_service);
}
```

# Program Parallelization

- Belief Propagation (BP) method for Stereo Vision

```
/* Client-side Program */

void bp_cb::TX()
{
  for( t=0; t < ITER; t++)
   for( y=0; y < HEIGHT; y++)
    for( x=(y+t)%2; x<WIDTH; x+=2)
     istrm << u[x][y+1] << d[x][y-1]
           << l [x+1][y] << r[x-1][y]
           << c[x][y]    << flush;
}


void bp_cb::RX()
{
  /* get the four output vectors
     <ou, od, ol, or> via ostrm */
}
```
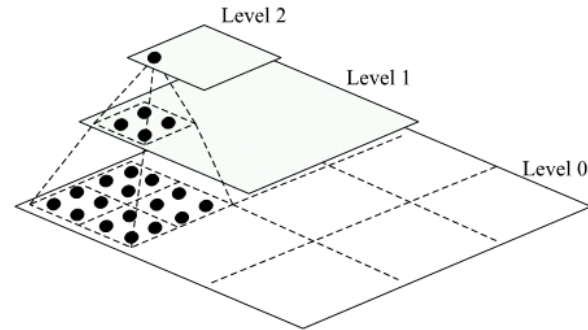
```
/* Server-side Program */

void bp_cb::SERVE()
{
  short16 u, d, l, r, c, ou, od, or, ol;
  istrm >> u >> d >> l
        >> r >> c >> pop;

  // compute disparity vectors
  // with vector arithmetic
  update(u, l, r, c, ou);
  update(d, l, r, c, od);
  update(u, d, r, c, or);
  update(u, d, l, c, ol);


  ostrm << ou << od << ol
        << or << flush;
}
```
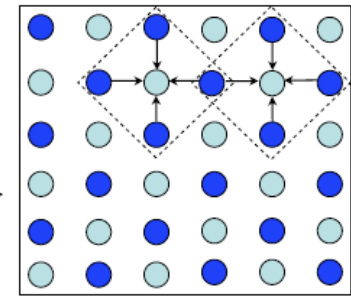


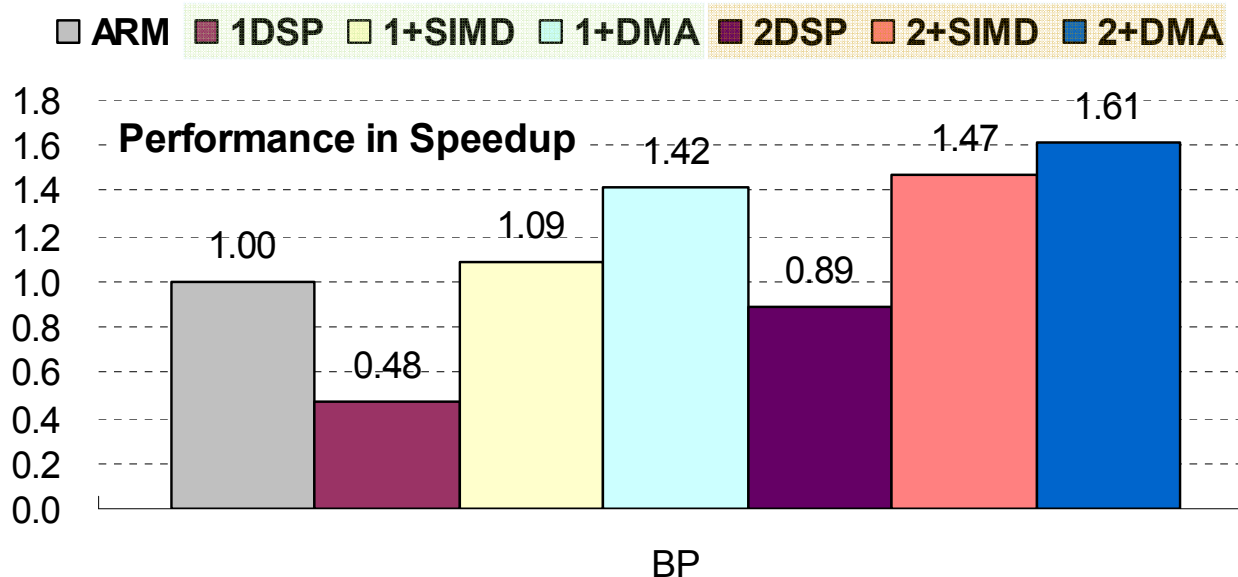belief propagation over multi-scale planes



message update by neighboring pixels

| Finding Concurrency | ■ Data decomposition<br>■ Data sharing (read-only) |
|---|---|
| Algorithm Structures | ■ Geometric decomposition (chunk: tuple<u,d,l,r,c>) |
| Supporting Structures | ■ Client-server (msg-update)<br>■ SIMD (disparity vectors) |
| Implementation Mechanisms | ■ RPC<br>■ Data streaming |

# Experimental Results on BP

- About the performance with +SIMD
  - SIMD vectors are employed to update disparity vectors.
  - Besides SIMD, few branches are eliminated by MIN intrinsics.
- About the performance with +DMA
  - DMAs on the multi-DSP system are incorporated to facilitate data streaming, implicitly in stream implementation.
- About the performance degradation (0.48 and 0.89)
  - MPU waiting for results from DSPs (especially serious for 1DSP).
  - Data communication and RPC overhead.

# OpenCL Compilation Flow for PAC Duo

**\*Work-item Coalescing**

Serialize work-items to avoid context switching.

**Cluster-aware Work-item Dispatching**
Transform each statement with **cluster intrinsic** to explicitly dispatch work-items to clusters.

```
gid = get_global_id(0);
C[gid] = A[gid] + B[gid];
```
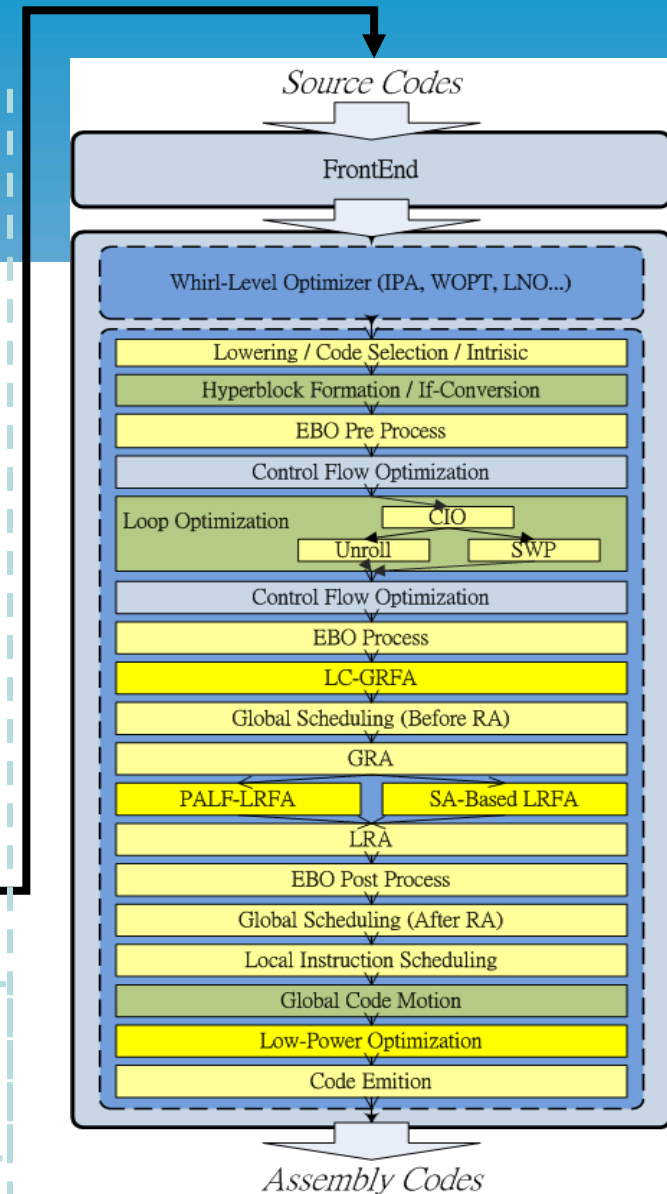
```
for (i=0; i<4; i++) {
    gid = get_global_id(0);
    C[gid] = A[gid] + B[gid];
}
```

```
for (i=0; i<4; i+=2) {
    gid  = get_global_id(0);
    gid2 = get_global_id(0) + 1;
    C[gid]  = __builtin_c1_add(A[gid],  B[gid]);
    C[gid2] = __builtin_c2_add(A[gid2], B[gid2]);
}
```

**Source-to-source Transformation:**
*Work-item Coalescing: Clang
Cluster-aware Work-item Dispatching: *Cetus

Source Codes

FrontEnd

Whirl-Level Optimizer (IPA, WOPT, LNO...)

Lowering / Code Selection / Intrisic
Hyperblock Formation / If-Conversion
EBO Pre Process
Control Flow Optimization
Loop Optimization — CIO — Unroll — SWP
Control Flow Optimization
EBO Process
LC-GRFA
Global Scheduling (Before RA)
GRA
PALF-LRFA — SA-Based LRFA
LRA
EBO Post Process
Global Scheduling (After RA)
Local Instruction Scheduling
Global Code Motion
Low-Power Optimization
Code Emition

Assembly Codes

"OpenCL Update", Yu-Te Lin, Shao-Ching Wang, Jia-Jer Li, Jenq-Kuen Lee, Open64 Developers Forum, Cupertino, USA, June 2011.

* J. Lee, et al, "An OpenCL Framework for Heterogeneous Multicores with Local Memory," PACT'10
*S. Lee, et al, "Cetus - An Extensible Compiler Infrastructure for Source-to-Source Transformation," LCPC'03

# Discussions: Emerging Issues C++0x and OpenCL/C++

- C++0x: Lambda functions allow programmers to define functions "directly" in places where their expressions are used.
  - Similar to **anonymous classes** in Java.
  - Better ways to provide a syntax sugar for parallel design patterns

```
typedef pair IntPair       /* A sorting example with lambdas */
std::vector v;

// sort v by the second element of the pair
std:sort(v.begin(), v.end(),
         [](const IntPair &a, const IntPair &b)  // call by reference
         { return a.second < b.second; }         // return value
);
```
//Source: Sorting example is referenced with Intel Software Blogs. Author: Stefanus Du Toit (Intel)

- C++0x vs. OpenCL:
  - C++0x thread model needs compilers for cares not to introduce data racing in optimizations.
  - C++0x thread model is suitable for OpenCL program with host, but might be a competing linguistic with kernel programs.

# Summary

- We present techniques in supporting C++ compiler and programming models for multi-core DSP systems.

- C++ SIMD vector class and overloading can be used to hide intrinsic functions of C programs for better abstractions.

- C++ program can help build parallel design patterns.

- Program parallelization with parallel design patterns on stereo-vision and image-processing applications are presented with a multi-DSP platform as a case study.

- Issues ahead when the embedded system meets C++0x and embedded OpenCL meets C++0x.

- Related references can be seen in http://www.cs.nthu.edu.tw/~jklee