

语义信息在领域编译中的角色

Using Semantic To Guide Compiler Optimizations

中科院计算所

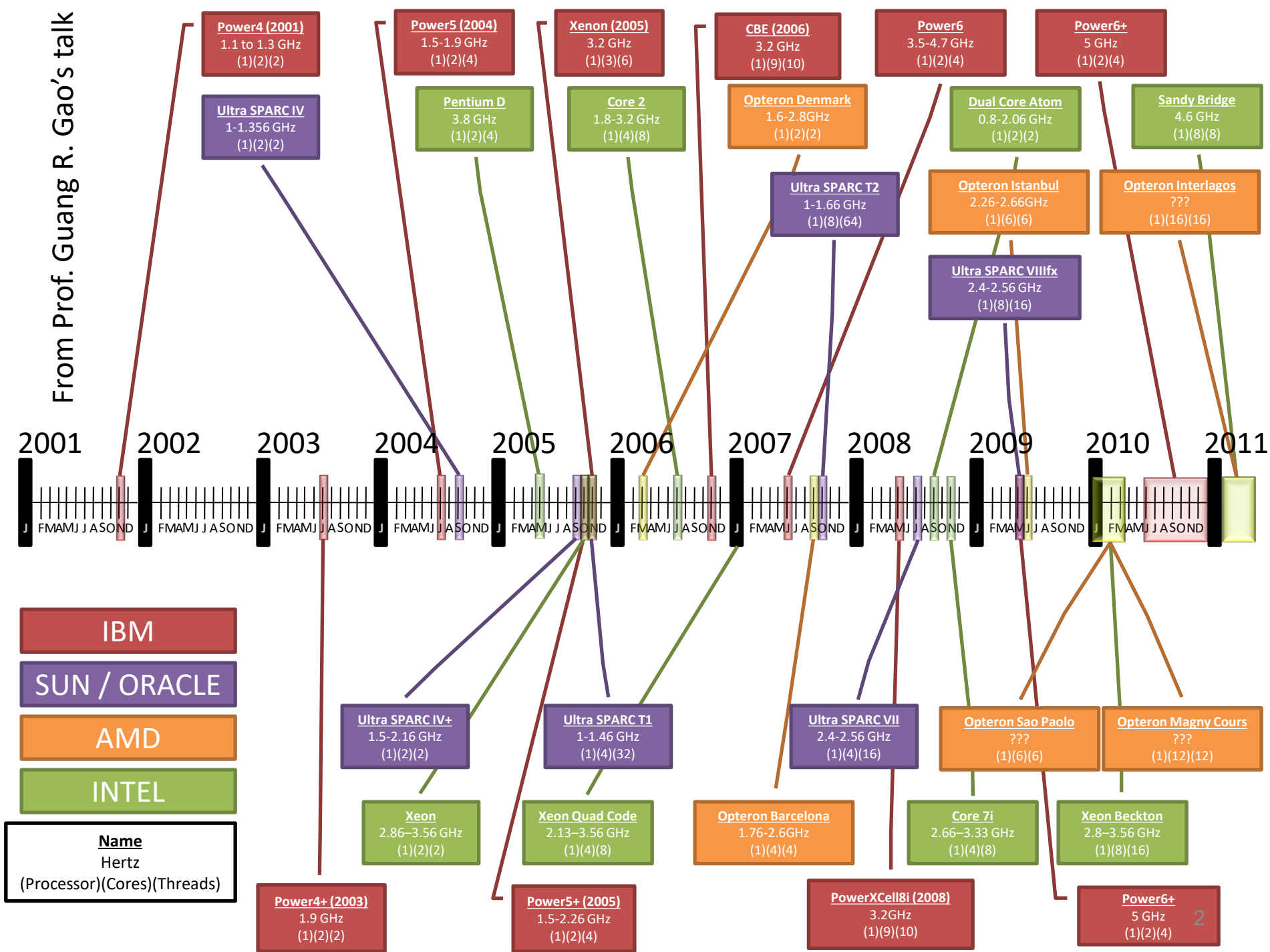
崔慧敏

2020年10月

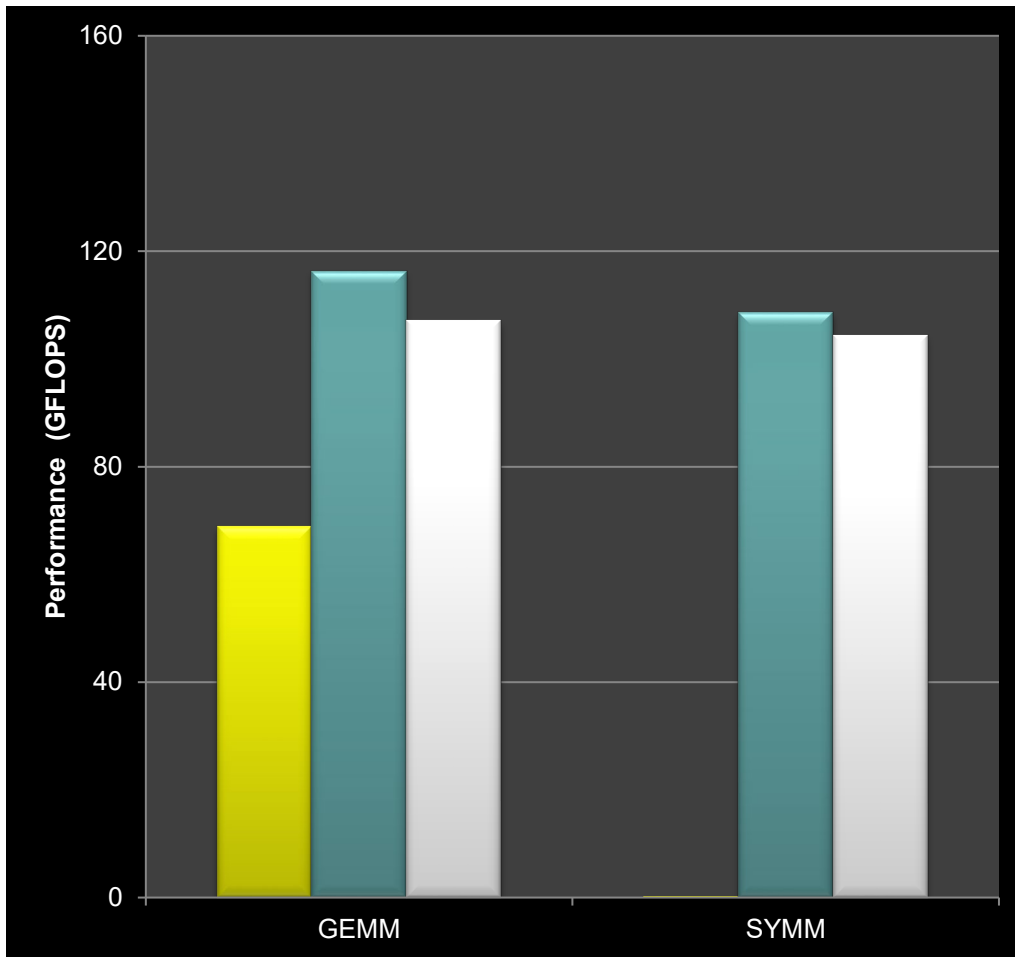


中国科学院计算技术研究所
INSTITUTE OF COMPUTING TECHNOLOGY, CHINESE ACADEMY OF SCIENCES

From Prof. Guang R. Gao's talk



Semantic Case 1: SYMM



SYMM Source Code:

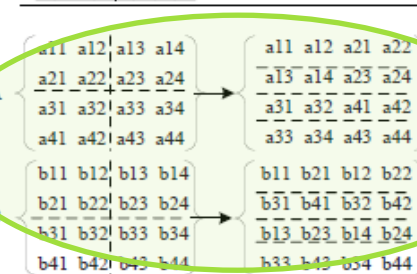
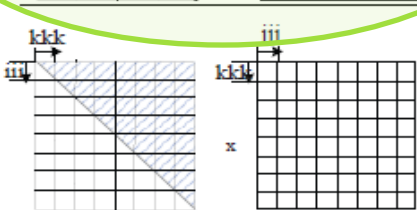
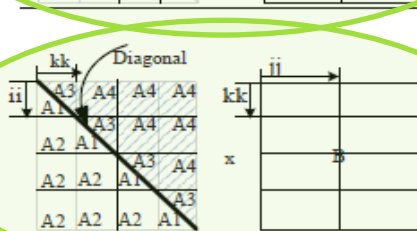
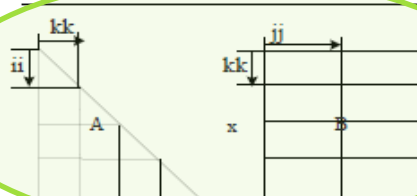
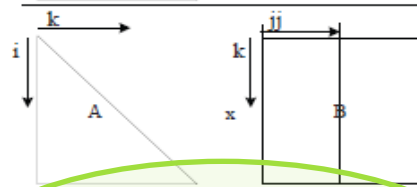
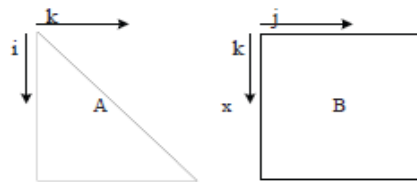
```
Li: for (i = 0; i < M; i++)  
  Lj: for (j = 0; j < N; j++)  
  {  
    Lk: for (k = 0; k < i; k++)  
    {  
      C[i][j] += A[i][k] * B[k][j];  
      C[k][j] += A[i][k] * B[i][j];  
    }  
    Ld: C[i][j] += A[i][i] * B[i][j]; //diagonal  
  }
```

SYMM Source Code:

```

Li: for (i = 0; i < M; i++)
  Lj: for (j = 0; j < N; j++)
  {
    Lk: for (k = 0; k < i; k++)
    {
      C[i][j] += A[i][k] * B[k][j];
      C[k][j] += A[i][k] * B[k][j];
    }
    Ld: C[i][j] += A[i][i] * B[i][j]; //dia
  }

```



第一步, 线程划分: 对循环迭代j进行strip-mining, 并对循环次序进行调整, 原始的循环顺序变为 **j, i, jj, k**. 并将最外层循环j分配到不同的线程。

第二步, 循环分块: 对循环迭代i和k进行循环分块 (tiling), 改善L2缓存的局部性。这时新的循环顺序为**j, i, k, ii, jj, kk**。

第三步, 循环分裂: 将循环体分裂为三个, 分别对应real area (A_1, \dots, A_2, \dots), shadow area (A_3, \dots, A_4, \dots), 和对角线元素。

第四步, 循环剥离 (peeling): 将三角形区域从矩形区域中剥离出来, 即将左图中的A1和A3分别从A2和A4中剥离。

第五步: 循环分块: 对剥离出的三角形区域A1和A3分别实施循环分块。对每个矩形区域 (其循环顺序为ii, jj, kk) 实施循环分块, 这时新的循环顺序为**j, i, k, ii, jj, kk, iii, kkk**

第六步: 数据布局重组: 将源矩阵A和B进行数据布局重组, 从行优先或者列优先转换为块优先, 块内和块间的顺序则遵从原有的行/列优先顺序。如左图所示, 这一变换是为了提高L1数据缓存的局部性。

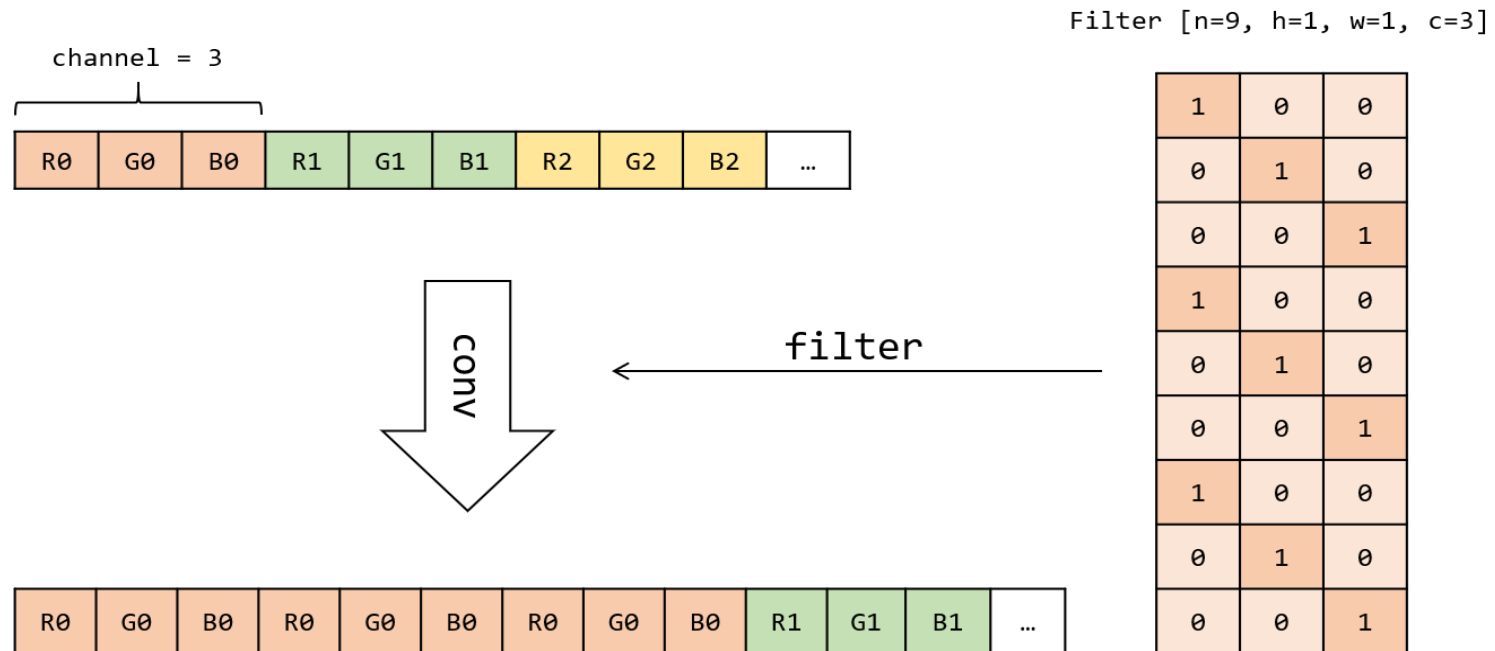
Triangular area

Symmetry

Matrix-mul

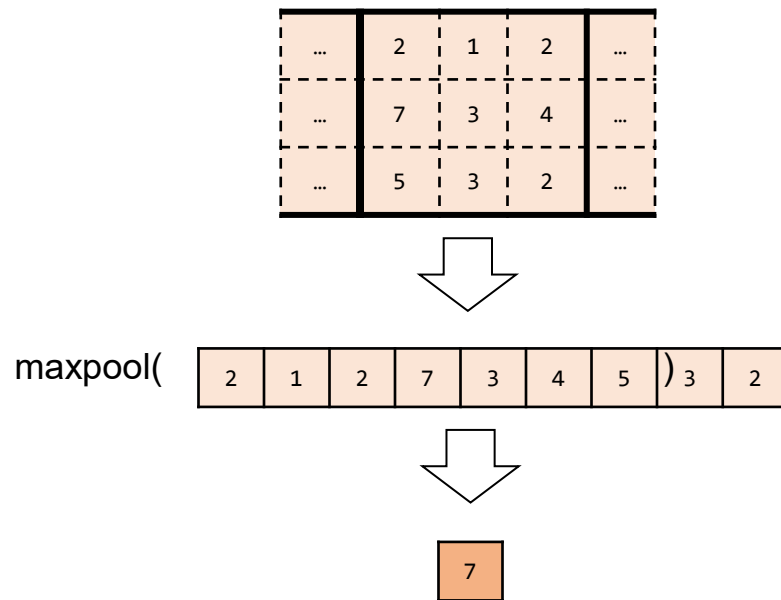
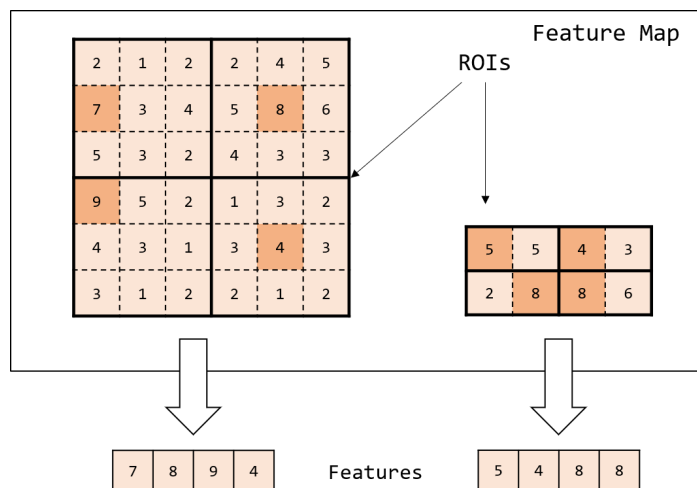
Semantic Case2 : resize

- We can use conv. to compute resize



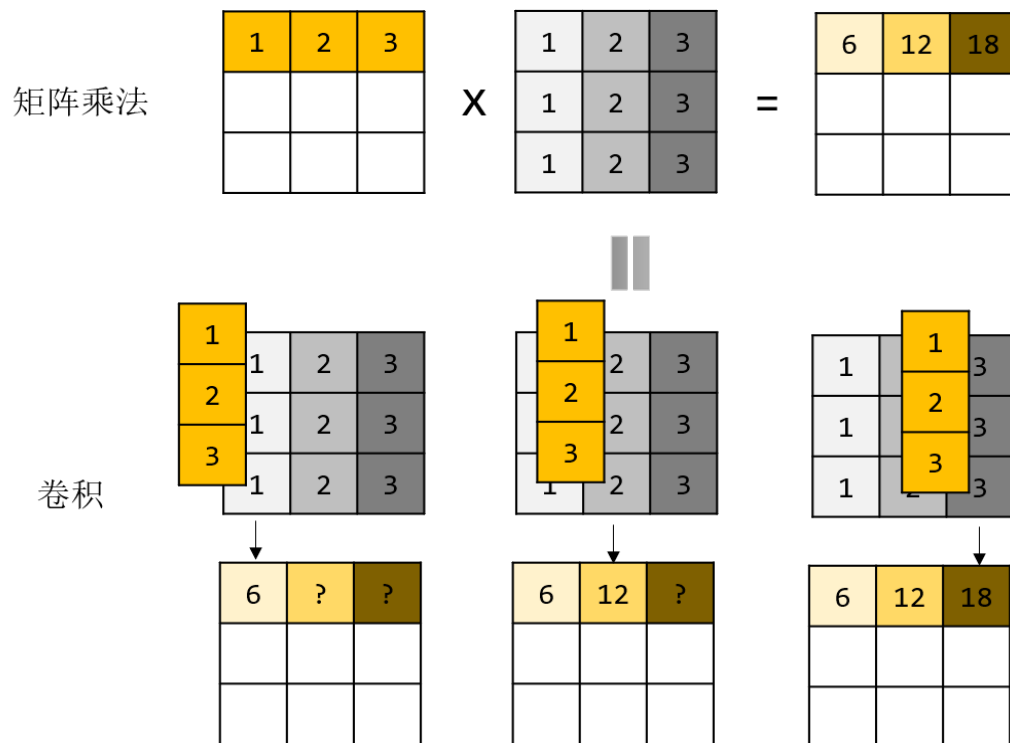
Semantic Case3: ROI Pooling

We can use Maxpooling to compute ROI pooling.

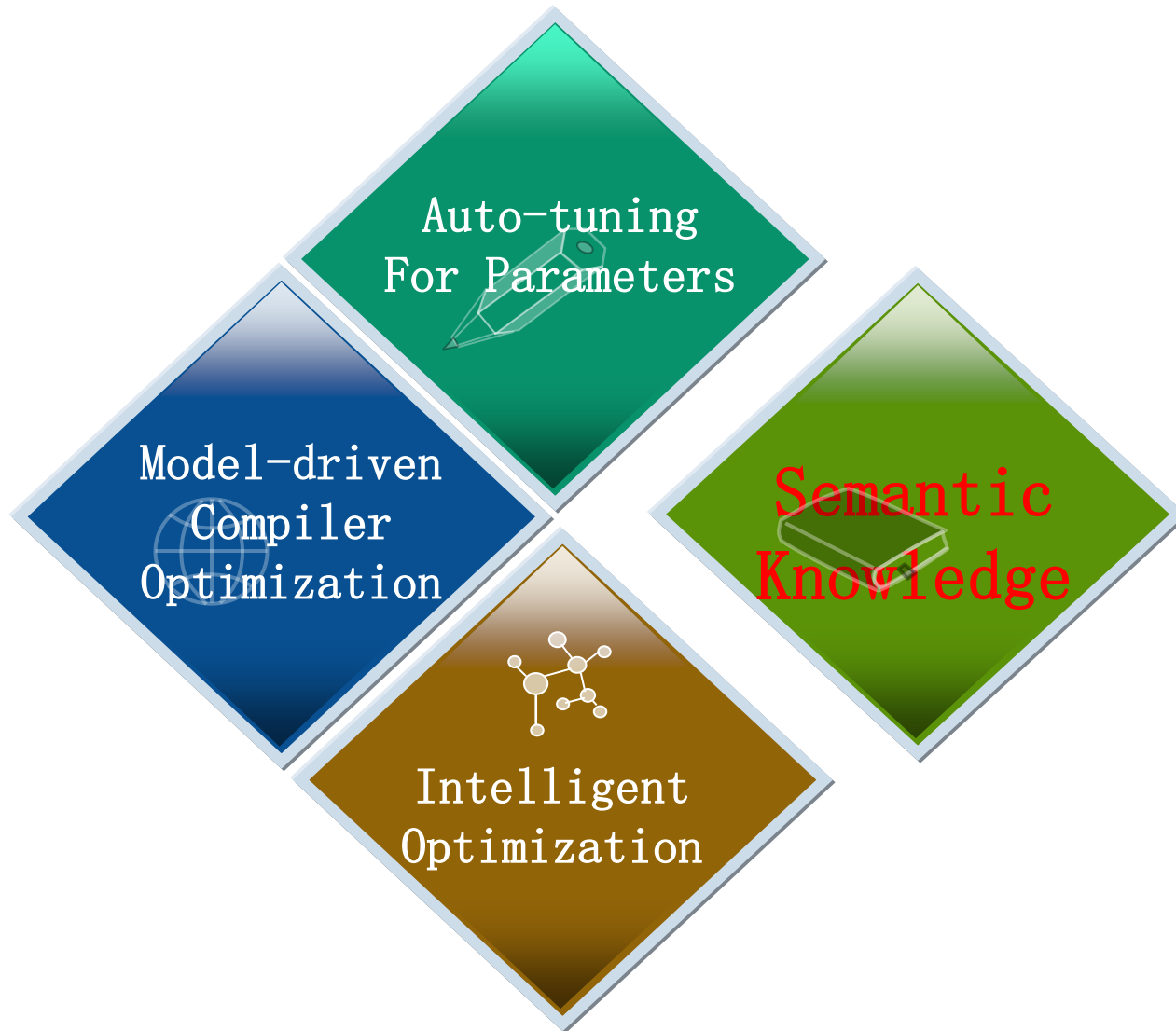


Semantic Case4: GEMM

We can use Conv to compute GEMM.



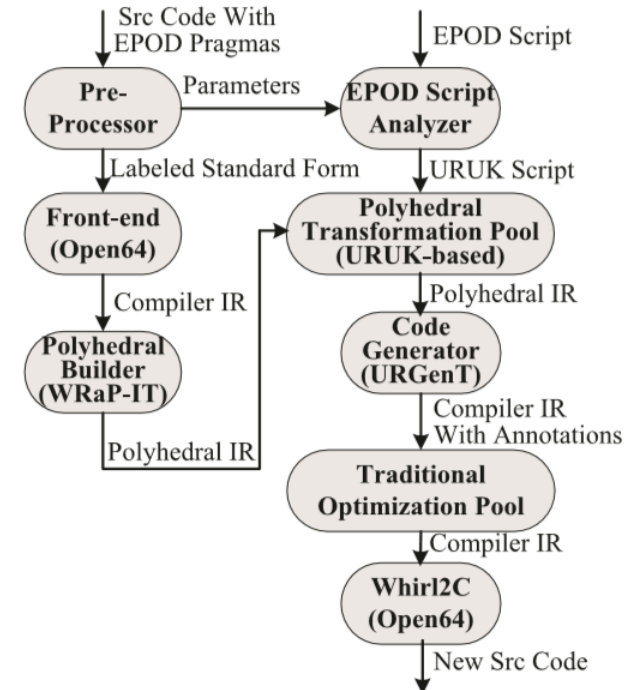
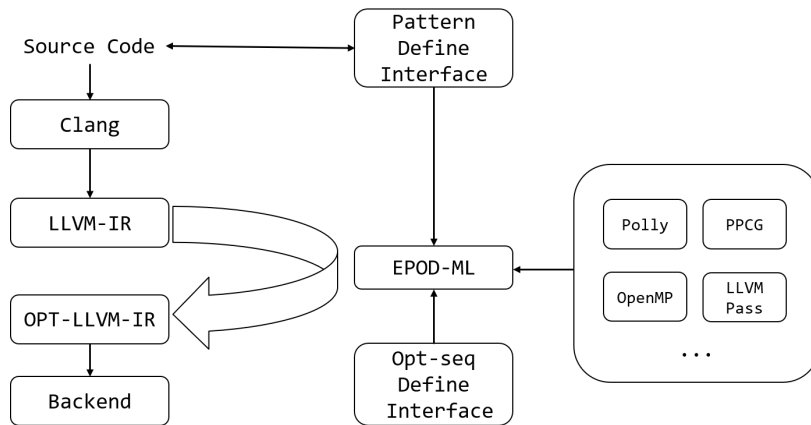
How to get Peak Performance?



Key Challenges

- **How to connect into the workflow of compilers?**
- **How to specify an optimization sequence?**
- **How to apply specific opt to specific code region?**

Extendable Pattern for Compilers



UDOL – User-defined Optimization Language

```
stmts {  
  stmt_1 {  
    domain { i: (0, 63), j:(0, 63)},  
    position {i, j, 0, 0},  
    mem_acc {  
      store_0 { mem {C}, domain {i:(0, 63), j:(0, 63)}}  
    },  
    dep { waw { stmt_1, stmt_2 } },  
    expr { st_0(const(0)) }  
  },  
  stmt_2 {  
    domain { i: (0, 63), j:(0, 63), k(0, 63)},  
    position {i, j, 1, k},  
    mem_acc {  
      load_0 { mem {C}, domain {i:(0, 63), j:(0, 63)}},  
      load_1 { mem {A}, domain {i:(0, 63), k:(0, 63)}},  
      load_2 { mem {B}, domain {k:(0, 63), j:(0, 63)}},  
      store_0 { mem {C}, domain {i:(0, 63), j:(0, 63)}}  
    },  
    dep { waw { stmt_1, stmt_2 } },  
    expr { st_0(add(ld(0), mul(ld(1), ld(2))))}  
  }  
}
```

1. Define a semantic pattern

- Similar with Polyhedral SCoP
- Each statement includes:
 - Loop domain
 - Position
 - Memory Reference
 - Dependency
 - Expression

UDOL – User-defined Optimization Language

```
def rep_conv(loop) {  
  s = loop.get_stmt();  
  d = s.get_domain(0);  
  col = d.range();  
  row = d.next().range();  
  dst = s.get_store();  
  src = s.get_load(d);  
  filter = s.get_load();  
  s.replace_with("llvm_conv(dst, src, filter,  
                           col, 1, row, 1, 1, 1, 1);");  
}
```

```
def matmul {  
  I = 16, J = 16  
  ll = polly.loop_tiling($$, [0, 1], [I, J])  
  l_0, l_1 = polly.loop_fission(ll, 3);  
  llvm.auto_vec(l_0);  
  llvm.rep_conv(l_1);  
}
```

2. Define an optimization seq.

- Polly for loop transformations.
 - Easy for extension
 - e.g. split tiling/ overlap tiling/ diamond tiling
- Backend opt. on LLVM/Open64.
 - Vectorization
 - On-chip memory usage

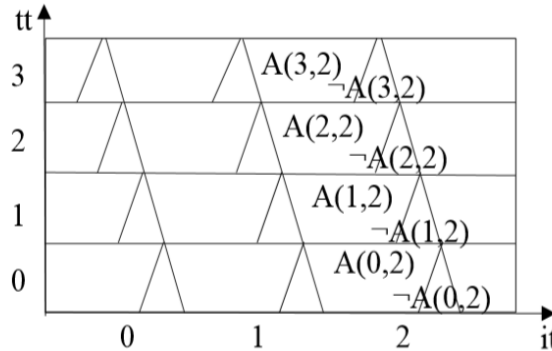
UDOL – User-defined Optimization Language

Original Loop:

```
for (t=1; t < T; t++) {
  for (i = 1; i < N; i++)
    B[i] = (A[i-1]+A[i]+A[i+1])/3;
  for (i = 1; i < N; i++)
    A[i] = B[i];
}
```

Single Statement Form:

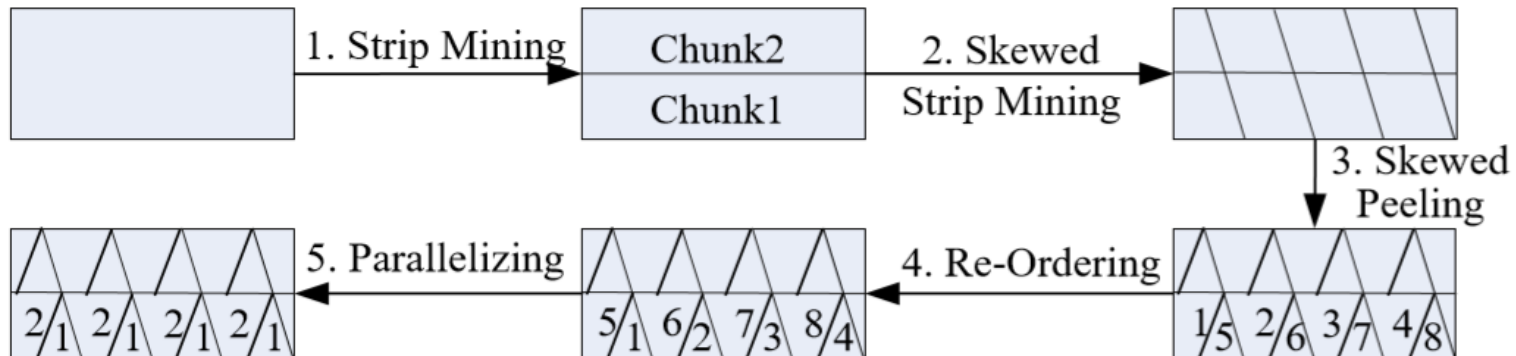
```
for (t=1; t < T; t++){
  for (i = 1; i < N; i++)
    A[t, i] = (A[t-1,i-1]+A[t-1,i]+A[t-1,i+1])/3;
}
```



(a) Tile Shape

```
for (tt=0; tt<t_tiles; tt++) {
  forall(it=0; it<i_tiles, it++)
    compute triangle  $\neg A(tt,it)$ ;
  barrier;
  forall (it=0; it<i_tiles, it++)
    compute trapezoid  $A(tt,it)$ ;
  barrier;
}
```

(b) Execution Order

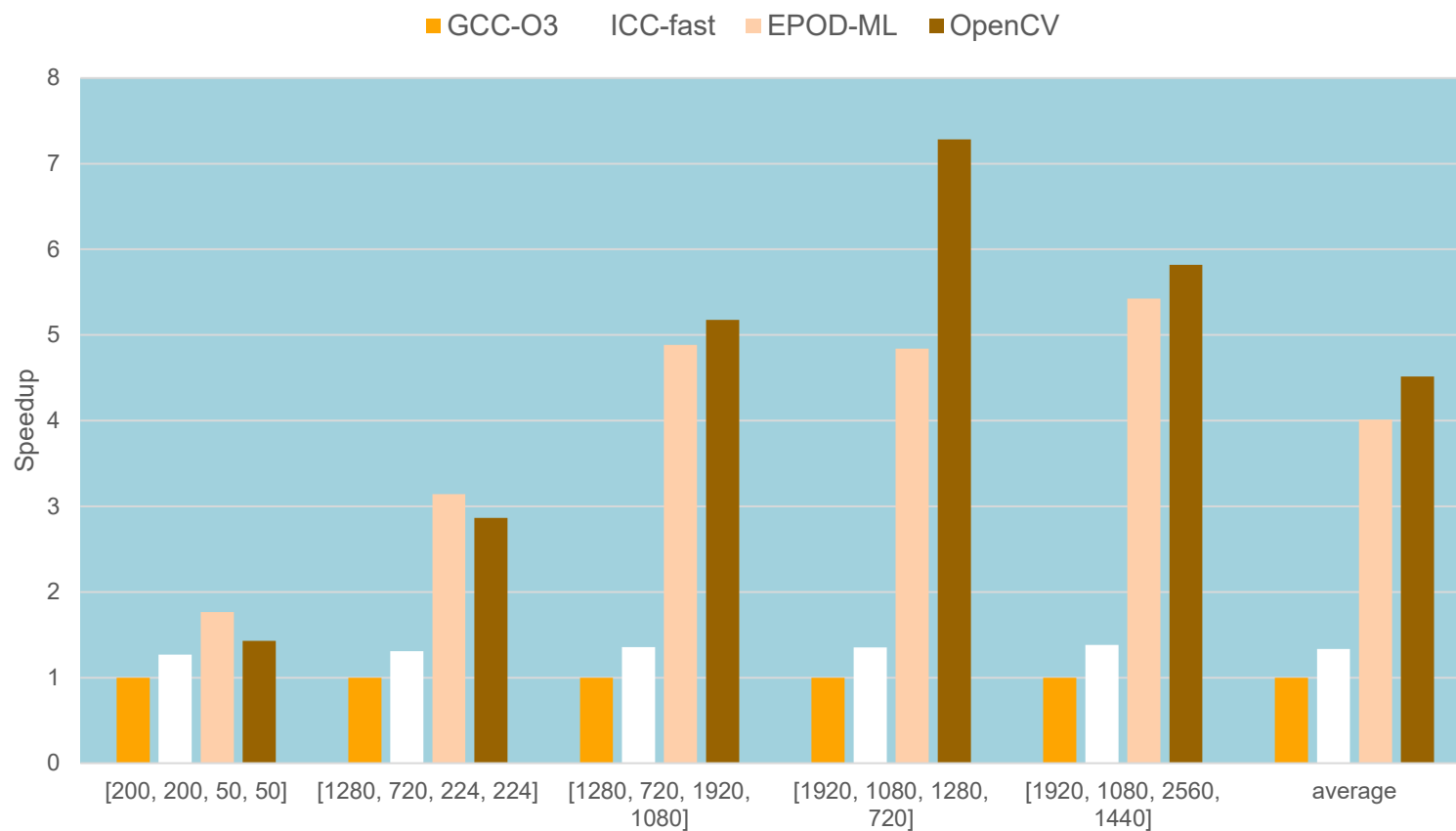


UDOL – User-defined Optimization Language

3. Specify a semantic region

```
#pragma EPOD-ML jacobi
for (t=0; t<T; t++) {
  for (i=0; i<N; i++) {
    for (j=0; j<N; j++) {
      if (i!=j) {
        sum += A[i][j]*x1[i];
      }
    }
    x[i] = (b[i]-sum)/A[i][i];
  }
  for (int i=0; i<N; i++) {
    x1[i] = x[i];
  }
}
#pragma EPOD-ML end
```

Evaluation



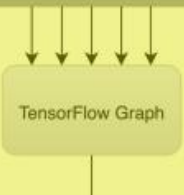
Resize加速比-CPU

人工智能编程语言与编译器

■问题1：人工智能需要通用的编程语言吗？

Machine Learning Systems are stuck in a Rut
Paul Barham and Michael Isard [HotOS2019]

- Currently much focus on optimizing 5 year old ML benchmarks
- New ideas in ML often require new primitives that are hard to compile and optimize for the zoo of modern hardware



The diagram illustrates a workflow where five separate inputs (represented by arrows) are fed into a single processing unit labeled 'TensorFlow Graph'. This visualizes the complexity of integrating multiple components into a unified system.

■问题2：芯片敏捷化 vs. 编译器敏捷化

■问题3：编译器敏捷化 vs. 峰值性能

Thanks



cuihm@ict.ac.cn