

# 编译优化：SLP矢量化

皮昊旋

- [1].[参考推文](#)
- [2].[参考论文1-Exploiting Superword Level Parallelism with Multimedia Instruction Sets](#)
- [3].[参考论文2-Loop-Aware SLP in GCC](#)

*Superword Level Parallelism*

将**相似**的**独立**指令组合成向量指令

*SIMD*指令

**单**指令**多**数据



```
1 void foo(float a1, float a2, float b1, float b2, float *A) {  
2     A[0] = a1*(a1 + b1);  
3     A[1] = a2*(a2 + b2);  
4     A[2] = a1*(a1 + b1);  
5     A[3] = a2*(a2 + b2);  
6 }
```



```
1  # %bb.0:  
  
2      addss    %xmm0, %xmm2  
  
3      mulss    %xmm0, %xmm2  
  
4      movss    %xmm2, (%rdi)  
  
5      addss    %xmm1, %xmm3  
  
6      mulss    %xmm1, %xmm3  
  
7      movss    %xmm3, 4(%rdi)  
  
8      movss    %xmm2, 8(%rdi)  
  
9      movss    %xmm3, 12(%rdi)  
  
10     retq
```

*clang++ case.cpp -O3 -fno-slp-vectorize -S*

<i>%xmm0</i>	$a_1$
<i>%xmm1</i>	$a_2$
<i>%xmm2</i>	$b_1$
<i>%xmm3</i>	$b_2$



```
1 void foo(float a1, float a2, float b1, float b2, float *A) {  
2     A[0] = a1*(a1 + b1);  
3     A[1] = a2*(a2 + b2);  
4     A[2] = a1*(a1 + b1);  
5     A[3] = a2*(a2 + b2);  
6 }
```



```
1  # %bb.0:  
  
2      addss    %xmm0, %xmm2  
  
3      mulss    %xmm0, %xmm2  
  
4      movss    %xmm2, (%rdi)  
  
5      addss    %xmm1, %xmm3  
  
6      mulss    %xmm1, %xmm3  
  
7      movss    %xmm3, 4(%rdi)  
  
8      movss    %xmm2, 8(%rdi)  
  
9      movss    %xmm3, 12(%rdi)  
  
10     retq
```

<i>%xmm0</i>	$a_1$
<i>%xmm1</i>	$a_2$
<i>%xmm2</i>	$a_1 + b_1$
<i>%xmm3</i>	$b_2$



```
1 void foo(float a1, float a2, float b1, float b2, float *A) {
2     A[0] = a1*(a1 + b1);
3     A[1] = a2*(a2 + b2);
4     A[2] = a1*(a1 + b1);
5     A[3] = a2*(a2 + b2);
6 }
```



```
1  # %bb.0:
2      addss    %xmm0, %xmm2
3      mulss    %xmm0, %xmm2
4      movss    %xmm2, (%rdi)
5      addss    %xmm1, %xmm3
6      mulss    %xmm1, %xmm3
7      movss    %xmm3, 4(%rdi)
8      movss    %xmm2, 8(%rdi)
9      movss    %xmm3, 12(%rdi)
10     retq
```

<i>%xmm0</i>	$a_1$
<i>%xmm1</i>	$a_2$
<i>%xmm2</i>	$a_1 * (a_1 + b_1)$
<i>%xmm3</i>	$b_2$



```
1 void foo(float a1, float a2, float b1, float b2, float *A) {
2     A[0] = a1*(a1 + b1);
3     A[1] = a2*(a2 + b2);
4     A[2] = a1*(a1 + b1);
5     A[3] = a2*(a2 + b2);
6 }
```



```
1  # %bb.0:
2      addss    %xmm0, %xmm2
3      mulss    %xmm0, %xmm2
4      movss    %xmm2, 4(%rdi)
5      addss    %xmm1, %xmm3
6      mulss    %xmm1, %xmm3
7      movss    %xmm3, 8(%rdi)
8      movss    %xmm2, 12(%rdi)
9      movss    %xmm3, 16(%rdi)
10     retq
```

	<i>%xmm0</i>	$a_1$
	<i>%xmm1</i>	$a_2$
$A[0] =$	<i>%xmm2</i>	$a_1 * (a_1 + b_1)$
	<i>%xmm3</i>	$b_2$



```
1 void foo(float a1, float a2, float b1, float b2, float *A) {
2     A[0] = a1*(a1 + b1);
3     A[1] = a2*(a2 + b2);
4     A[2] = a1*(a1 + b1);
5     A[3] = a2*(a2 + b2);
6 }
```



```
1  # %bb.0:
2      addss    %xmm0, %xmm2
3      mulss    %xmm0, %xmm2
4      movss    %xmm2, 4(%rdi)
5      addss    %xmm1, %xmm3
6      mulss    %xmm1, %xmm3
7      movss    %xmm3, 8(%rdi)
8      movss    %xmm2, 12(%rdi)
9      movss    %xmm3, 16(%rdi)
10     retq
```

	%xmm0	$a_1$
	%xmm1	$a_2$
$A[0] =$	%xmm2	$a_1 * (a_1 + b_1)$
$A[1] =$	%xmm3	$a_2 * (a_2 + b_2)$



```
1 void foo(float a1, float a2, float b1, float b2, float *A) {
2     A[0] = a1*(a1 + b1);
3     A[1] = a2*(a2 + b2);
4     A[2] = a1*(a1 + b1);
5     A[3] = a2*(a2 + b2);
6 }
```



```
1  # %bb.0:
2      addss    %xmm0, %xmm2
3      mulss    %xmm0, %xmm2
4      movss    %xmm2, 4(%rdi)
5      addss    %xmm1, %xmm3
6      mulss    %xmm1, %xmm3
7      movss    %xmm3, 12(%rdi)
8      movss    %xmm2, 8(%rdi)
9      movss    %xmm3, 16(%rdi)
10     retq
```

	%xmm0	$a_1$
	%xmm1	$a_2$
$A[2] = A[0] =$	%xmm2	$a_1 * (a_1 + b_1)$
$A[3] = A[1] =$	%xmm3	$a_2 * (a_2 + b_2)$





```
1 void foo(float a1, float a2, float b1, float b2, float *A) {
2     A[0] = a1*(a1 + b1);
3     A[1] = a2*(a2 + b2);
4     A[2] = a1*(a1 + b1);
5     A[3] = a2*(a2 + b2);
6 }
```

*clang++ case.cpp -O3 -S*

%xmm0	$a_1$
%xmm1	$a_2$
%xmm2	$b_1$
%xmm3	$b_2$



```
1 # %bb.0:
2     unpcklps    %xmm1, %xmm0                # xmm0 = xmm0[0],xmm1[0],xmm0[1],xmm1[1]
3     unpcklps    %xmm3, %xmm2                # xmm2 = xmm2[0],xmm3[0],xmm2[1],xmm3[1]
4     addps       %xmm0, %xmm2
5     mulps       %xmm0, %xmm2
6     movlhps     %xmm2, %xmm2                # xmm2 = xmm2[0],0]
7     movups      %xmm2, (%rdi)
8     retq
```



```
1 void foo(float a1, float a2, float b1, float b2, float *A) {
2     A[0] = a1*(a1 + b1);
3     A[1] = a2*(a2 + b2);
4     A[2] = a1*(a1 + b1);
5     A[3] = a2*(a2 + b2);
6 }
```



```
1  # %bb.0:
2      unpcklps    %xmm1, %xmm0                # xmm0 = xmm0[0],xmm1[0],xmm0[1],xmm1[1]
3      unpcklps    %xmm3, %xmm2                # xmm2 = xmm2[0],xmm3[0],xmm2[1],xmm3[1]
4      addps       %xmm0, %xmm2
5      mulps       %xmm0, %xmm2
6      movlhps     %xmm2, %xmm2                # xmm2 = xmm2[0],0]
7      movups      %xmm2, (%rdi)
8      retq
```

%xmm0	$[a_1, a_2, -, -]$
%xmm1	$a_2$
%xmm2	$[b_1, b_2, -, -]$
%xmm3	$b_2$



```
1 void foo(float a1, float a2, float b1, float b2, float *A) {
2     A[0] = a1*(a1 + b1);
3     A[1] = a2*(a2 + b2);
4     A[2] = a1*(a1 + b1);
5     A[3] = a2*(a2 + b2);
6 }
```



```
1 # %bb.0:
2     unpcklps    %xmm1, %xmm0                # xmm0 = xmm0[0],xmm1[0],xmm0[1],xmm1[1]
3     unpcklps    %xmm3, %xmm2                # xmm2 = xmm2[0],xmm3[0],xmm2[1],xmm3[1]
4     addps       %xmm0, %xmm2
5     mulps       %xmm0, %xmm2
6     movlhps     %xmm2, %xmm2                # xmm2 = xmm2[0],0]
7     movups      %xmm2, (%rdi)
8     retq
```

<i>%xmm0</i>	$[a_1, a_2, -, -]$
<i>%xmm1</i>	$a_2$
<i>%xmm2</i>	$[a_1 + b_1, a_2 + b_2, -, -]$
<i>%xmm3</i>	$b_2$



```
1 void foo(float a1, float a2, float b1, float b2, float *A) {
2     A[0] = a1*(a1 + b1);
3     A[1] = a2*(a2 + b2);
4     A[2] = a1*(a1 + b1);
5     A[3] = a2*(a2 + b2);
6 }
```

%xmm0	$[a_1, a_2, -, -]$
%xmm1	$a_2$
%xmm2	$[a_1*(a_1 + b_1), a_2*(a_2 + b_2), -, -]$
%xmm3	$b_2$



```
1 # %bb.0:
2     unpcklps    %xmm1, %xmm0                # xmm0 = xmm0[0],xmm1[0],xmm0[1],xmm1[1]
3     unpcklps    %xmm3, %xmm2                # xmm2 = xmm2[0],xmm3[0],xmm2[1],xmm3[1]
4     addps       %xmm0, %xmm2
5     mulps       %xmm0, %xmm2
6     movlhps     %xmm2, %xmm2                # xmm2 = xmm2[0],0]
7     movups      %xmm2, (%rdi)
8     retq
```



```
1 void foo(float a1, float a2, float b1, float b2, float *A) {
2     A[0] = a1*(a1 + b1);
3     A[1] = a2*(a2 + b2);
4     A[2] = a1*(a1 + b1);
5     A[3] = a2*(a2 + b2);
6 }
```

%xmm0	$[a_1, a_2, -, -]$
%xmm1	$a_2$
%xmm2	$[a_1*(a_1 + b_1), a_2*(a_2 + b_2),$ $a_1*(a_1 + b_1), a_2*(a_2 + b_2),]$
%xmm3	$b_2$



```
1 # %bb.0:
2 unpcklps    %xmm1, %xmm0           # xmm0 = xmm0[0],xmm1[0],xmm0[1],xmm1[1]
3 unpcklps    %xmm3, %xmm2           # xmm2 = xmm2[0],xmm3[0],xmm2[1],xmm3[1]
4 addps       %xmm0, %xmm2
5 mulps       %xmm0, %xmm2
6 movlhps     %xmm2, %xmm2           # xmm2 = xmm2[0],0]
7 movups      %xmm2, (%rdi)
8 retq
```

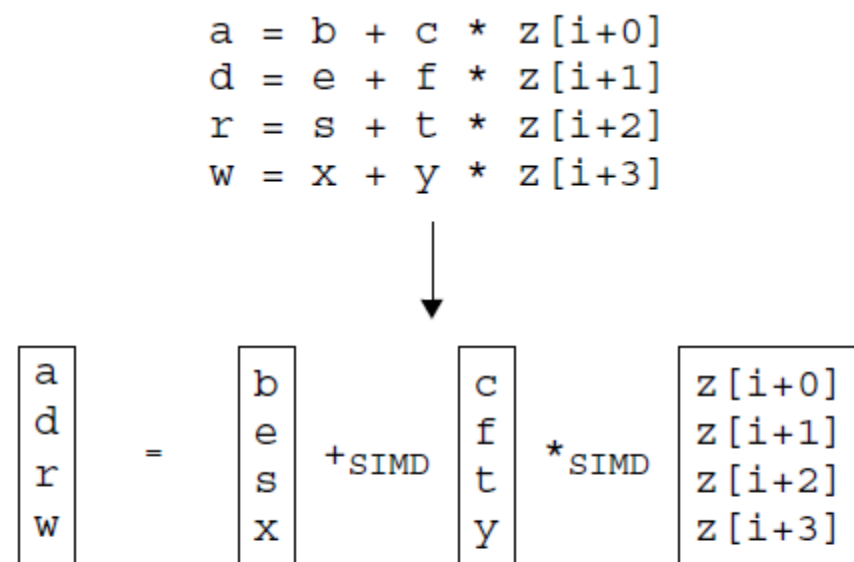


Figure 1: Isomorphic statements that can be packed and executed in parallel.

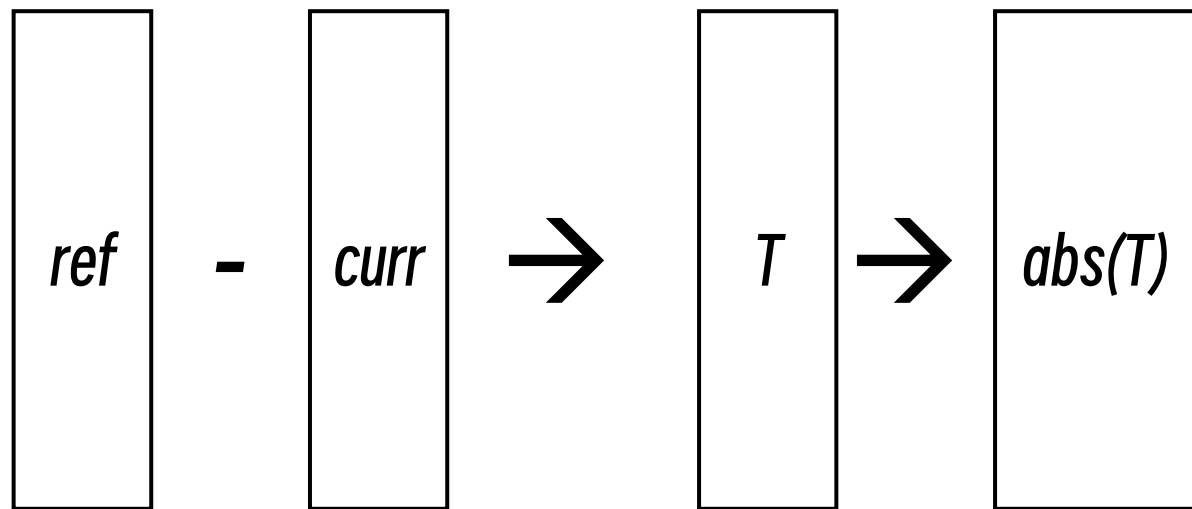
```
for (i=0; i<16; i++) {  
    localdiff = ref[i] - curr[i];  
    diff += abs(localdiff);  
}
```

体现出较强的数据依赖性  $\longrightarrow$  对向量化造成阻碍

(a) Original loop.

```
for (i=0; i<16; i++) {  
    T[i] = ref[i] - curr[i];  
}  
  
for (i=0; i<16; i++) {  
    diff += abs(T[i]);  
}
```

(b) After scalar expansion and loop fission.



减少数据依赖之后，即可使用SIMD指令，在一个时钟周期内计算多个元素的差值，从而提高循环效率。

```

for (i=0; i<16; i+=4) {
    localdiff = ref[i+0] - curr[i+0];
    diff += abs(localdiff);

    localdiff = ref[i+1] - curr[i+1];
    diff += abs(localdiff);

    localdiff = ref[i+2] - curr[i+2];
    diff += abs(localdiff);

    localdiff = ref[i+3] - curr[i+3];
    diff += abs(localdiff);
}

```

(c) Superword level parallelism exposed after unrolling.

```

for (i=0; i<16; i+=4) {
    localdiff0 = ref[i+0] - curr[i+0];
    localdiff1 = ref[i+1] - curr[i+1];
    localdiff2 = ref[i+2] - curr[i+2];
    localdiff3 = ref[i+3] - curr[i+3];

    diff += abs(localdiff0);
    diff += abs(localdiff1);
    diff += abs(localdiff2);
    diff += abs(localdiff3);
}

```

(d) Packable statements grouped together after renaming.

**手动构造类似于向量的代码，  
从而减少循环次数，提高效率  
(循环展开后暴露了并行性)**

**重命名消除数据依赖  
将同类型操作放在一起**



```
do {
    dst[0] = (src1[0] + src2[0]) >> 1;
    dst[1] = (src1[1] + src2[1]) >> 1;
    dst[2] = (src1[2] + src2[2]) >> 1;
    dst[3] = (src1[3] + src2[3]) >> 1;

    dst += 4;
    src1 += 4;
    src2 += 4;
}
while (dst != end);
```

循环向量化:

- ①先将do while 循环转换成for循环 (恢复归纳变量)
- ②将展开循环恢复成未展开的状态 (rerolling)
- ③再进行上述优化

Figure 3: An example of a hand-optimized matrix operation that proves unvectorizable.

SLP: 直接做**基本块 (basic block)** 间优化

循环向量化: 更关注**迭代间**的优化  
SLP: 更关注**迭代内**的优化



预处理

循环展开 *Loop Unrolling*

对齐分析 *Alignment analysis*

预 优 化 *Pre-Optimization*

```
SLP_extract: BasicBlock  $B \rightarrow$  BasicBlock  
  PackSet  $P \leftarrow \emptyset$   
   $P \leftarrow \text{find\_adj\_refs}(B, P)$   
   $P \leftarrow \text{extend\_packlist}(B, P)$   
   $P \leftarrow \text{combine\_packs}(P)$   
  return  $\text{schedule}(B, [], P)$ 
```

识别相邻内存引用 *find\_adj\_refs*

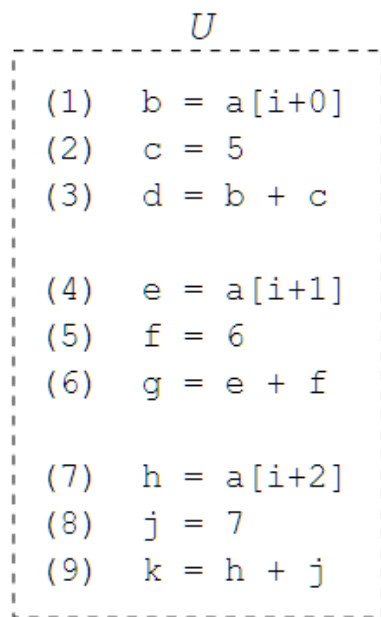
“扩张”PackSet *extend\_packlist*

组合Packs *combine\_packs*

准备调度 *schedule*

# 识别相邻内存引用 *find\_adj\_refs*

```
find_adj_refs: BasicBlock  $B \times \text{PackSet } P \rightarrow \text{PackSet}$   
  foreach Stmt  $s \in B$  do  
    foreach Stmt  $s' \in B$  where  $s \neq s'$  do  
      if has_mem_ref( $s$ )  $\wedge$  has_mem_ref( $s'$ ) then  
        if adjacent( $s, s'$ ) then  
          Int  $align \leftarrow \text{get\_alignment}(s)$   
          if stmts_can_pack( $B, P, s, s', align$ ) then  
             $P \leftarrow P \cup \{\langle s, s' \rangle\}$   
  return  $P$ 
```

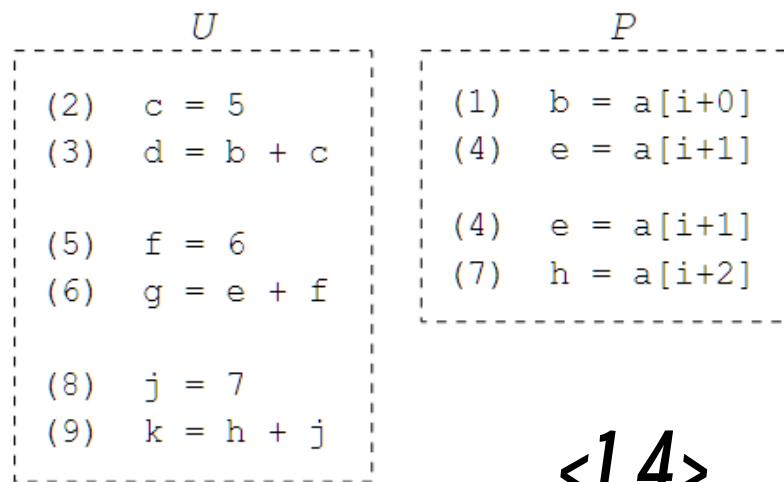


(a)

*get\_alignment(s)*

*int align*

## 对齐分析 *Alignment analysis*



(b)

$\langle 1, 4 \rangle$

$\langle 4, 7 \rangle$

```

find_adj_refs: BasicBlock  $B \times$  PackSet  $P \rightarrow$  PackSet
  foreach Stmt  $s \in B$  do
    foreach Stmt  $s' \in B$  where  $s \neq s'$  do
      if has_mem_ref( $s$ )  $\wedge$  has_mem_ref( $s'$ ) then
        if adjacent( $s, s'$ ) then
          Int  $align \leftarrow$  get_alignment( $s$ )
          if stmts_can_pack( $B, P, s, s', align$ ) then
             $P \leftarrow P \cup \{s'\}$ 
  return  $P$ 

```

stmts\_can\_pack: BasicBlock  $B \times$  PackSet  $P \times$   
 Stmt  $s \times$  Stmt  $s' \times$  Int  $align \rightarrow$  Boolean

```

  if isomorphic( $s, s'$ ) then
    if independent( $s, s'$ ) then
      if  $\forall \langle t, t' \rangle \in P. t \neq s$  then
        if  $\forall \langle t, t' \rangle \in P. t' \neq s'$  then
          Int  $align_s \leftarrow$  get_alignment( $s$ )
          Int  $align_{s'} \leftarrow$  get_alignment( $s'$ )
          if  $align_s \equiv \top \vee align_s \equiv align$  then
            if  $align_{s'} \equiv \top \vee align_{s'} \equiv align + \text{data\_size}(s')$  then
              return true
    return false

```

# “扩张”PackSet *extend\_packlist*

`extend_packlist`: BasicBlock  $B \times$  PackSet  $P \rightarrow$  PackSet

**repeat**

PackSet  $P_{prev} \leftarrow P$

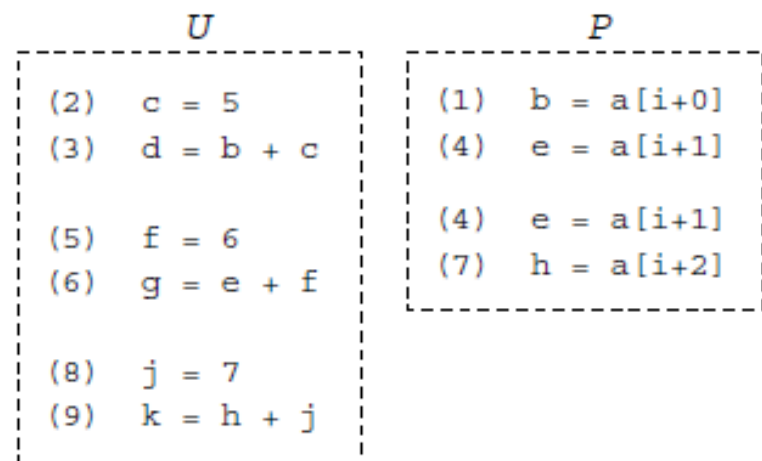
**foreach** Pack  $p \in P$  **do**

$P \leftarrow \text{follow\_use\_defs}(B, P, p)$

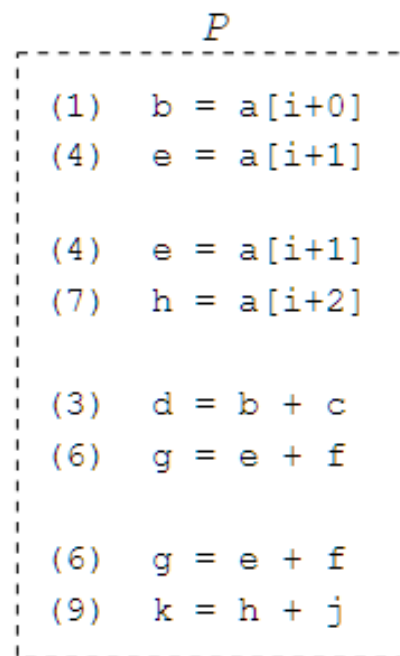
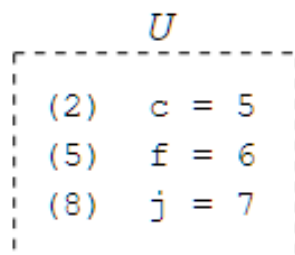
$P \leftarrow \text{follow\_def\_uses}(B, P, p)$

**until**  $P \equiv P_{prev}$

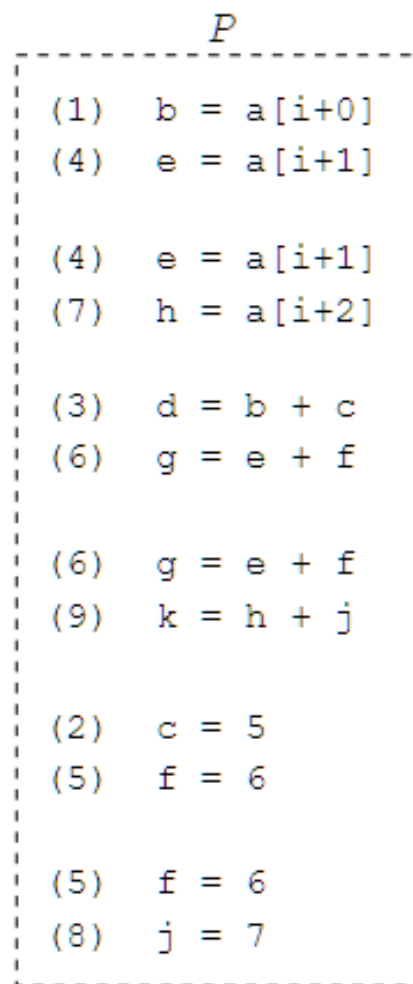
**return**  $P$



(b)



(c)



(d)

```

follow_use_defs: BasicBlock  $B \times \text{PackSet } P \times \text{Pack } p \rightarrow \text{PackSet}$ 
    where  $p = \langle s, s' \rangle$ ,  $s = [x_0 := f(x_1, \dots, x_m)]$ ,  $s' = [x'_0 := f(x'_1, \dots, x'_m)]$ 
    Int  $align \leftarrow \text{get\_alignment}(s)$ 
    for  $j \leftarrow 1$  to  $m$  do
        if  $\exists t \in B. t = [x_j := \dots] \wedge \exists t' \in B. t' = [x'_j := \dots]$  then
            if  $\text{stmts\_can\_pack}(B, P, t, t', align)$ 
                if  $\text{est\_savings}(\langle t, t' \rangle, P) \geq 0$  then
                     $P \leftarrow P \cup \{\langle t, t' \rangle\}$ 
                     $\text{set\_alignment}(s, s', align)$ 
    return  $P$ 

```

```

follow_def_uses: BasicBlock  $B \times \text{PackSet } P \times \text{Pack } p \rightarrow \text{PackSet}$ 
    where  $p = \langle s, s' \rangle$ ,  $s = [x_0 := f(x_1, \dots, x_m)]$ ,  $s' = [x'_0 := f(x'_1, \dots, x'_m)]$ 
    Int  $align \leftarrow \text{get\_alignment}(s)$ 
    Int  $savings \leftarrow -1$ 
    foreach Stmt  $t \in B$  where  $t = [\dots := g(\dots, x_0, \dots)]$  do
        foreach Stmt  $t' \in B$  where  $t \neq t' = [\dots := h(\dots, x'_0, \dots)]$  do
            if  $\text{stmts\_can\_pack}(B, P, t, t', align)$  then
                if  $\text{est\_savings}(\langle t, t' \rangle, P) > savings$  then
                     $savings \leftarrow \text{est\_savings}(\langle t, t' \rangle, P)$ 
                    Stmt  $u \leftarrow t$ 
                    Stmt  $u' \leftarrow t'$ 
    if  $savings \geq 0$  then
         $P \leftarrow P \cup \{\langle u, u' \rangle\}$ 
         $\text{set\_alignment}(u, u')$ 
    return  $P$ 

```

# 组合Packs *combine\_packs*

combine\_packs: PackSet  $P \rightarrow$  PackSet

**repeat**

PackSet  $P_{prev} \leftarrow P$

**foreach** Pack  $p = \langle s_1, \dots, s_n \rangle \in P$  **do**

**foreach** Pack  $p' = \langle s'_1, \dots, s'_m \rangle \in P$  **do**

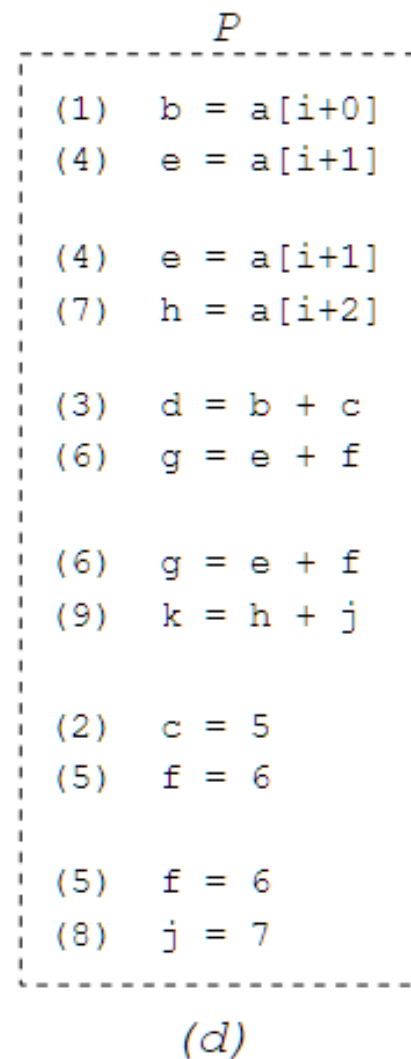
**if**  $s_n \equiv s'_1$  **then**

$P \leftarrow P - \{p, p'\} \cup \{\langle s_1, \dots, s_n, s'_2, \dots, s'_m \rangle\}$

**until**  $P \equiv P_{prev}$

**return**  $P$

去重、打包





# 组合Packs *combine\_packs*

combine\_packs: PackSet  $P \rightarrow$  PackSet

**repeat**

PackSet  $P_{prev} \leftarrow P$

**foreach** Pack  $p = \langle s_1, \dots, s_n \rangle \in P$  **do**

**foreach** Pack  $p' = \langle s'_1, \dots, s'_m \rangle \in P$  **do**

**if**  $s_n \equiv s'_1$  **then**

$P \leftarrow P - \{p, p'\} \cup \{\langle s_1, \dots, s_n, s'_2, \dots, s'_m \rangle\}$

**until**  $P \equiv P_{prev}$

**return**  $P$

(1)  $b = a[i+0]$   
(4)  $e = a[i+1]$   
(7)  $h = a[i+2]$

(3)  $d = b + c$   
(6)  $g = e + f$   
(9)  $k = h + j$

(2)  $c = 5$   
(5)  $f = 6$   
(8)  $j = 7$

(e)

$\begin{matrix} b \\ e \\ h \end{matrix} = \begin{matrix} a[i+0] \\ a[i+1] \\ a[i+2] \end{matrix}$

$\begin{matrix} c \\ f \\ j \end{matrix} = \begin{matrix} 5 \\ 6 \\ 7 \end{matrix}$

$\begin{matrix} d \\ g \\ k \end{matrix} = \begin{matrix} b \\ e \\ h \end{matrix} + \begin{matrix} c \\ f \\ j \end{matrix}$

(f)

## 去重、打包

# 准备调度 *schedule*

```
schedule: BasicBlock  $B \times$  BasicBlock  $B' \times$  PackSet  $P$   
     $\rightarrow$  BasicBlock  
    for  $i \leftarrow 0$  to  $|B|$  do  
        if  $\exists p = \langle \dots, s_i, \dots \rangle \in P$  then  
            if  $\forall s \in p. \text{deps\_scheduled}(s, B')$  then  
                foreach Stmt  $s \in p$  do  
                     $B \leftarrow B - s$   
                     $B' \leftarrow B' \cdot s$   
                return schedule( $B, B', P$ )  
            else if  $\text{deps\_scheduled}(s_i, B')$  then  
                return schedule( $B - s_i, B' \cdot s_i, P$ )  
    if  $|B| \neq 0$  then  
         $P \leftarrow P - \{p\}$  where  $p = \text{first}(B, P)$   
        return schedule( $B, B', P$ )  
    return  $B'$ 
```

$x = a[i+0] + k1$   
 $y = a[i+1] + k2$

$q = b[i+0] + y$   
 $r = b[i+1] + k3$   
 $s = b[i+2] + k4$

$z = a[i+2] + s$



$x = a[i+0] + k1$   
 $y = a[i+1] + k2$   
 $z = a[i+2] + s$

$q = b[i+0] + y$   
 $r = b[i+1] + k3$   
 $s = b[i+2] + k4$

$z$  依赖  $s$   
 $q$  依赖  $y$

分开**最早**未被调度的组

**最早**: 根据原始basic block进行判定  
即把xyz组中的 $y$ 单独拿出先调度

# 缺点

## 不适合长向量架构

*The drawback to this method is that it may not be applicable to long vector architectures. Since the unroll factor must be consistent with the vector size, unrolling may produce basic blocks that overwhelm the analysis and the code generator. As such, this method is mainly applicable to architectures with short vectors.*

### **3.1 Loop Unrolling**

Loop unrolling is performed early since it is most easily done at a high level. As discussed, it is used to transform vector parallelism into basic blocks with superword level parallelism. In order to ensure full utilization of the superword datapath in the presence of a vectorizable loop, the unroll factor must be customized to the data sizes used within the loop. For example, a vectorizable loop containing 16-bit values should be unrolled 8 times for a 128-bit datapath. Our system currently unrolls loops based on the smallest data type present.

*Thank you*