# PSLP: Padded SLP Automatic Vectorization

Vasileios Porpodas[†], Alberto Magni[‡]
and Timothy M. Jones[†]
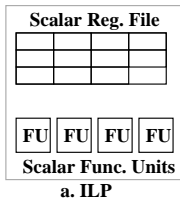
University of Cambridge[†]
University of Edinburgh[‡]

# Why SIMD Vectorization?
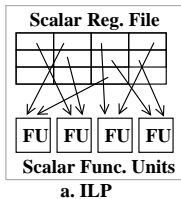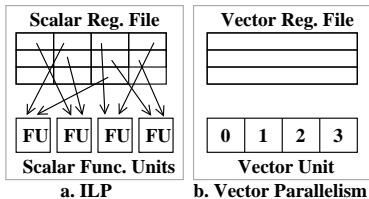
- Scalable parallelism



a. ILP

# Why SIMD Vectorization?

- Scalable parallelism
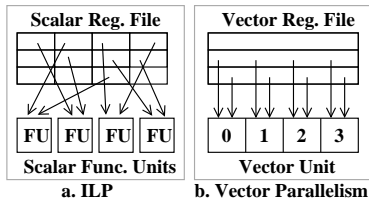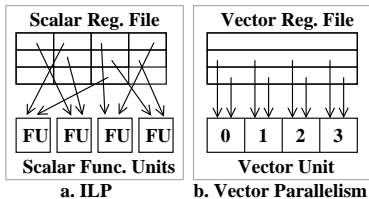


a. ILP

# Why SIMD Vectorization?

- Scalable parallelism



Scalar Reg. File | Vector Reg. File

FU FU FU FU

Scalar Func. Units

0 1 2 3

Vector Unit

a. ILP | b. Vector Parallelism

# Why SIMD Vectorization?

- Scalable parallelism



a. ILP     b. Vector Parallelism

# Why SIMD Vectorization?

- Scalable parallelism
- High Performance
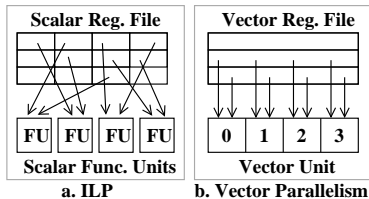


a. ILP  b. Vector Parallelism

# Why SIMD Vectorization?

- Scalable parallelism
- High Performance
- Energy efficiency



a. ILP    b. Vector Parallelism

# Why SIMD Vectorization?

- Scalable parallelism
- High Performance
- Energy efficiency
- Supported since mid 90's
- Frequent updates of vector ISAs



a. ILP  b. Vector Parallelism

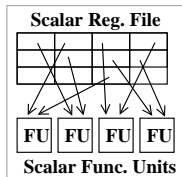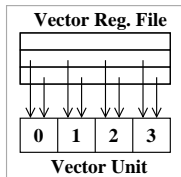# Why SIMD Vectorization?

- Scalable parallelism
- High Performance
- Energy efficiency
- Supported since mid 90's
- Frequent updates of vector ISAs
- Vector generation not done in hardware
- Low-level programming or capable compiler



a. ILP

b. Vector Parallelism

# SLP Straight-Line Code Vectorizer

- Superword Level Parallelism [Larsen PLDI'00]

# SLP Straight-Line Code Vectorizer

- Superword Level Parallelism [Larsen PLDI'00]
- State-of-the-art straight-line code vectorizer

# SLP Straight-Line Code Vectorizer

- Superword Level Parallelism [Larsen PLDI'00]
- State-of-the-art straight-line code vectorizer
- Implemented in most compilers (including GCC and LLVM)

# SLP Straight-Line Code Vectorizer

- Superword Level Parallelism [Larsen PLDI'00]
- State-of-the-art straight-line code vectorizer
- Implemented in most compilers (including GCC and LLVM)
- In theory it should be a superset of loop-vectorizer

# SLP Straight-Line Code Vectorizer

- **S**uperword **L**evel **P**arallelism [Larsen PLDI'00]
- State-of-the-art straight-line code vectorizer
- Implemented in most compilers (including GCC and LLVM)
- In theory it should be a superset of loop-vectorizer
  - Unroll loop and vectorize with SLP
  - Even if loop-vectorizer fails, SLP could partly succeed

# SLP Straight-Line Code Vectorizer

- Superword Level Parallelism [Larsen PLDI'00]
- State-of-the-art straight-line code vectorizer
- Implemented in most compilers (including GCC and LLVM)
- In theory it should be a superset of loop-vectorizer
  - Unroll loop and vectorize with SLP
  - Even if loop-vectorizer fails, SLP could partly succeed
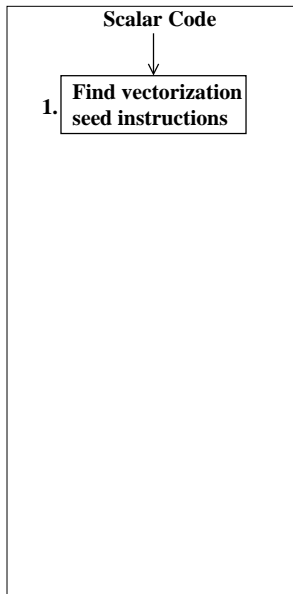- In practice it is missing features present in the Loop vectorizer (Interleaved Loads, Predication)

# SLP Vectorization Algorithm
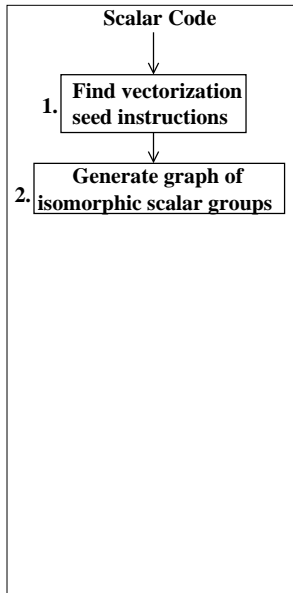
- Input is scalar IR

**Scalar Code**

# SLP Vectorization Algorithm

- Input is scalar IR
- Seed instructions are:
  1. Consecutive Stores
  2. Reductions

**Scalar Code**
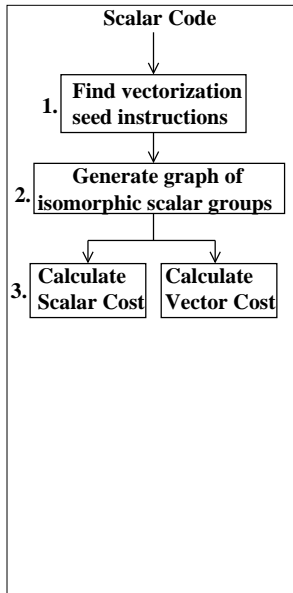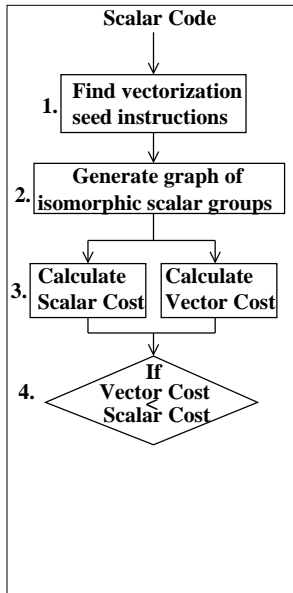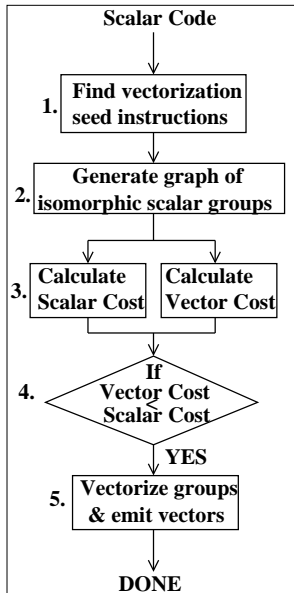
1. **Find vectorization seed instructions**

# SLP Vectorization Algorithm

- Input is scalar IR
- Seed instructions are:
  1. Consecutive Stores
  2. Reductions
- Graph contains vectorizable isomorphic instructions

**Scalar Code**

1. **Find vectorization seed instructions**
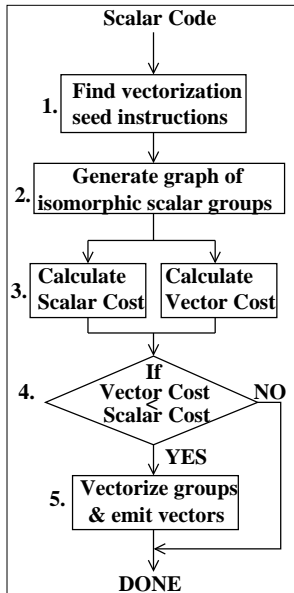
2. **Generate graph of isomorphic scalar groups**

# SLP Vectorization Algorithm

- Input is scalar IR
- Seed instructions are:
  1. Consecutive Stores
  2. Reductions
- Graph contains vectorizable isomorphic instructions
- Cost: weighted instr. count



**Scalar Code**

1. **Find vectorization seed instructions**

2. **Generate graph of isomorphic scalar groups**

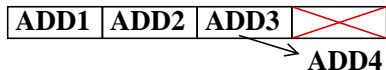3. **Calculate Scalar Cost** | **Calculate Vector Cost**

# SLP Vectorization Algorithm

- Input is scalar IR
- Seed instructions are:
  1. Consecutive Stores
  2. Reductions
- Graph contains vectorizable isomorphic instructions
- Cost: weighted instr. count
- Check vectorization profitability

# SLP Vectorization Algorithm

- Input is scalar IR
- Seed instructions are:
  1. Consecutive Stores
  2. Reductions
- Graph contains vectorizable isomorphic instructions
- Cost: weighted instr. count
- Check vectorization profitability
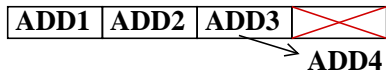- Emit vectors only if profitable

# SLP Vectorization Algorithm

- Input is scalar IR
- Seed instructions are:
  1. Consecutive Stores
  2. Reductions
- Graph contains vectorizable isomorphic instructions
- Cost: weighted instr. count
- Check vectorization profitability
- Emit vectors only if profitable



Scalar Code

1. Find vectorization seed instructions

2. Generate graph of isomorphic scalar groups

3. Calculate Scalar Cost | Calculate Vector Cost

4. If Vector Cost Scalar Cost — NO

YES

5. Vectorize groups & emit vectors

DONE

# When SLP Fails

1. Data Dependencies

# When SLP Fails

1. Data Dependencies

2. Too many gather/scatter instructions. Costs outweigh benefits.

| ADD1 | ADD2 | ADD3 | ✕ |
|------|------|------|------|

ADD4

| Original | Vectorized | | | |
|----------|------------|---|---|---|
| ADD1 | Insert1 | | | |
| ADD2 | Insert2 | | | |
| ADD3 | Insert3 | | | |
| ADD4 | Insert4 | | | |
| | ADD1 | ADD2 | ADD3 | ADD4 |
| | Extract1 | | | |
| | Extract2 | | | |
| | Extract3 | | | |
| | Extract4 | | | |

# When SLP Fails

1. Data Dependencies

2. Too many gather/scatter instructions. Costs outweigh benefits.

3. Non-isomorphism

| ADD1 | ADD2 | ADD3 | ✕ |
|------|------|------|---|

ADD4

| Original | Vectorized |
|----------|------------|
| ADD1 | Insert1 |
| ADD2 | Insert2 |
| ADD3 | Insert3 |
| ADD4 | Insert4 |
| | ADD1 ADD2 ADD3 ADD4 |
| | Extract1 |
| | Extract2 |
| | Extract3 |
| | Extract4 |

| ADD1 | ADD2 | MUL | ADD4 |
|------|------|-----|------|

# SLP Fails due to non-isomorphism

```
      ...
B[i]  = A[i]  * 7.0 + 1.0;
B[i+1]= A[i+1]      + 5.0;
      ...
```
   **a. Input C code**

| **(X)** **Instruction Node or Constant** | **→** **Data Flow Edge** |

# SLP Fails due to non-isomorphism



```
          ...
B[i]   = A[i]   * 7.0 + 1.0;
B[i+1] = A[i+1]       + 5.0;
          ...
```
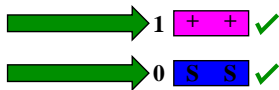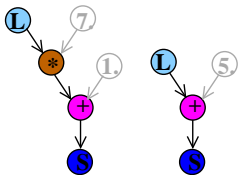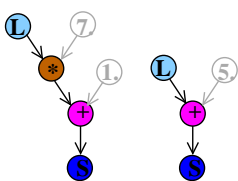**a. Input C code**



**b. DFG**

| Ⓧ Instruction Node or Constant | → Data Flow Edge |
| --- | --- |

# SLP Fails due to non-isomorphism



**a. Input C code**

```
  ...
B[i]   = A[i]   * 7.0 + 1.0;
B[i+1] = A[i+1]       + 5.0;
  ...
```

**b. DFG**    **c. SLP internal graph**    **d. SLP vectorized groups**
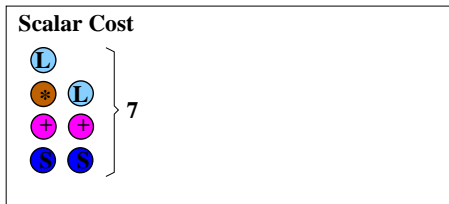
X Instruction Node or Constant → Data Flow Edge
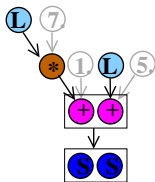
# SLP Fails due to non-isomorphism



```
      ...
B[i]   = A[i]   * 7.0 + 1.0;
B[i+1] = A[i+1]       + 5.0;
      ...
```
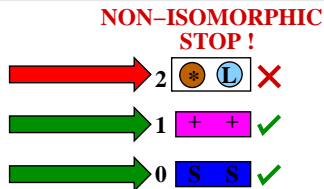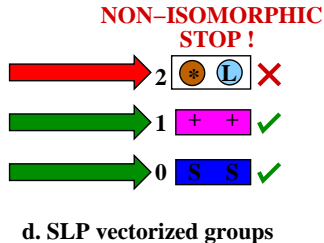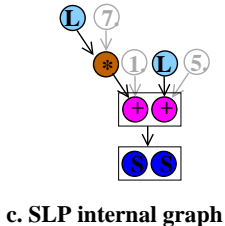**a. Input C code**

**b. DFG**

**c. SLP internal graph**
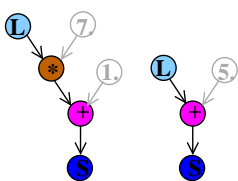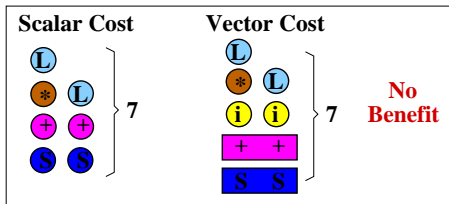
**d. SLP vectorized groups**

Ⓧ **Instruction Node or Constant** → **Data Flow Edge**

# SLP Fails due to non-isomorphism



```
      ...
B[i]   = A[i]    * 7.0 + 1.0;
B[i+1] = A[i+1]      + 5.0;
      ...
```

**a. Input C code**

**b. DFG**

**c. SLP internal graph**

**d. SLP vectorized groups**

**NON–ISOMORPHIC STOP !**

Ⓧ **Instruction Node or Constant**  →  **Data Flow Edge**

# SLP Fails due to non-isomorphism



**a. Input C code**

**Scalar Cost**

**b. DFG**

**c. SLP internal graph**

**d. SLP vectorized groups**

**NON–ISOMORPHIC STOP !**

(X) **Instruction Node or Constant** → **Data Flow Edge**

# SLP Fails due to non-isomorphism



a. Input C code

b. DFG

c. SLP internal graph

d. SLP vectorized groups

(X) Instruction Node or Constant → Data Flow Edge

# PSLP fixes Non-Isomorphism



**a. PSLP graphs**

| | |
|---|---|
| **(X)** **Instruction or Constant** | →**Data Flow Edge** |

# PSLP fixes Non-Isomorphism



**a. PSLP graphs**          **b. PSLP padded graphs**

(X) **Instruction or Constant**          →**Data Flow Edge**

# PSLP fixes Non-Isomorphism



**a. PSLP graphs**          **b. PSLP padded graphs**

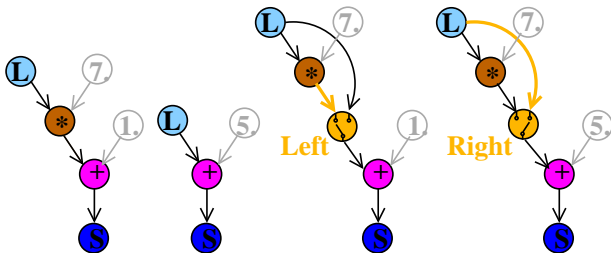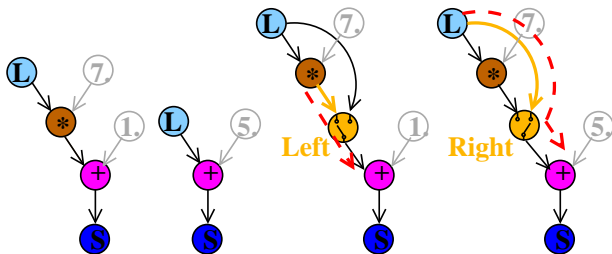| Ⓧ  Instruction or Constant | →Data Flow Edge |

# PSLP fixes Non-Isomorphism



**a. PSLP graphs**        **b. PSLP padded graphs**
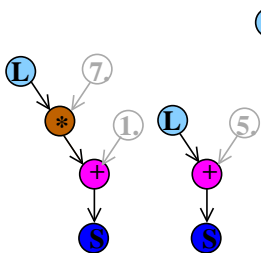
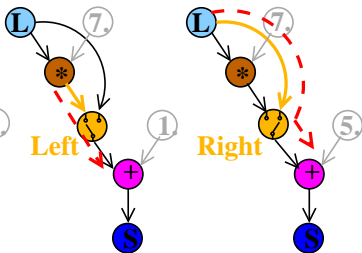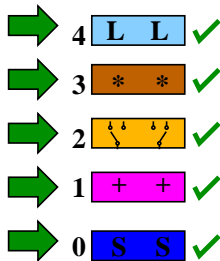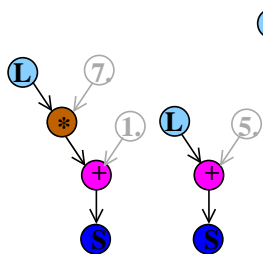| (X) Instruction or Constant | (Y) Select Instruction | →Data Flow Edge |

# PSLP fixes Non-Isomorphism



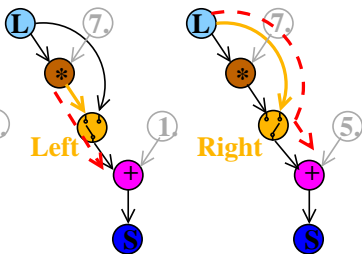a. PSLP graphs      b. PSLP padded graphs

(X) **Instruction or Constant**    (Y) **Select Instruction**    →**Data Flow Edge**

# PSLP fixes Non-Isomorphism



a. PSLP graphs          b. PSLP padded graphs          c. PSLP groups

---

$(X)$ **Instruction or Constant**     $(Y)$ **Select Instruction**     $\rightarrow$**Data Flow Edge**

# PSLP fixes Non-Isomorphism



**a. PSLP graphs**   **b. PSLP padded graphs**   **c. PSLP groups**

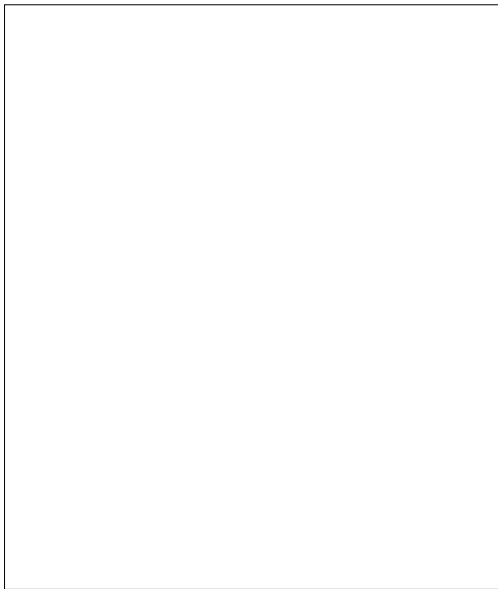Ⓧ **Instruction or Constant**   🔆 **Select Instruction**   →**Data Flow Edge**
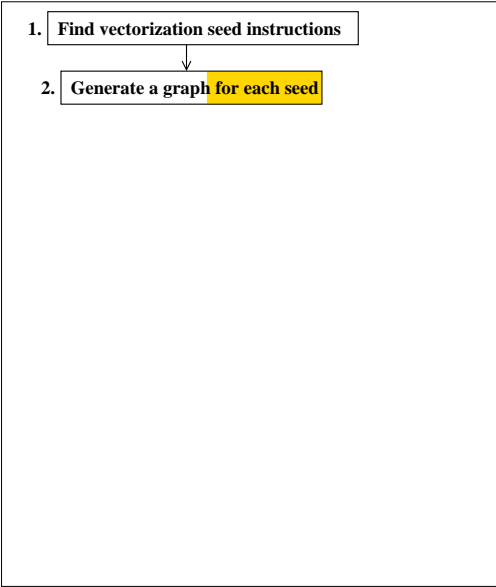
# PSLP Algorithm

- Extension to SLP

# PSLP Algorithm

- Extension to SLP

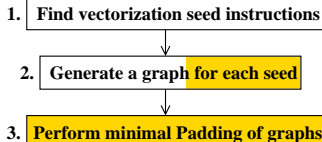1. **Find vectorization seed instructions**

# PSLP Algorithm

- Extension to SLP
- Generate multiple graphs (unlike SLP)



1. Find vectorization seed instructions

2. Generate a graph for each seed
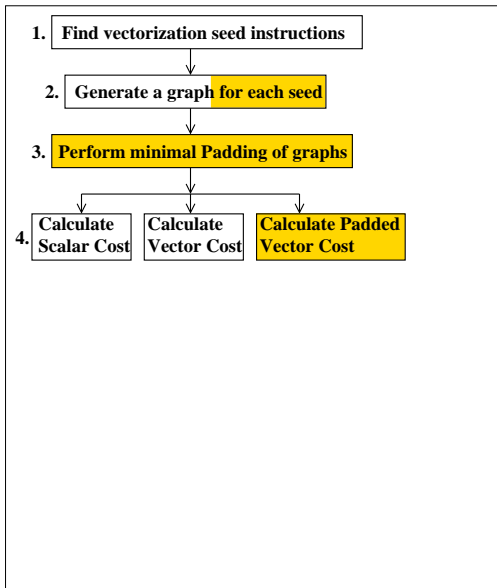
# PSLP Algorithm

- Extension to SLP
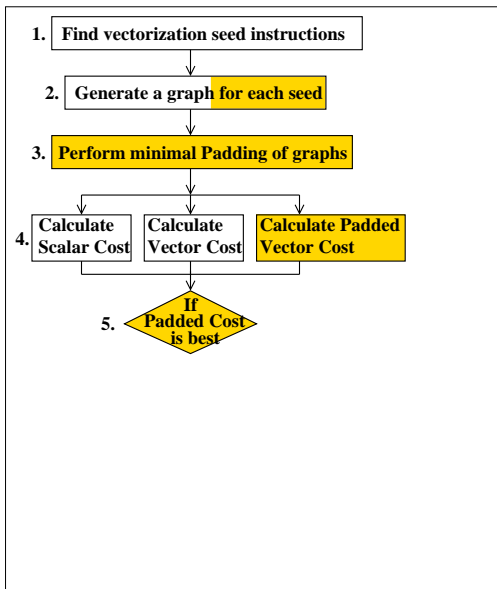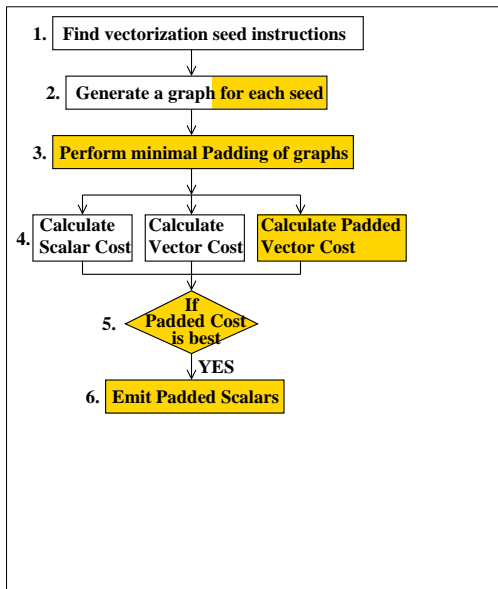- Generate multiple graphs (unlike SLP)
- Minimal Padding

# PSLP Algorithm

- Extension to SLP
- Generate multiple graphs (unlike SLP)
- Minimal Padding
- Cost estimation

# PSLP Algorithm

- Extension to SLP
- Generate multiple graphs (unlike SLP)
- Minimal Padding
- Cost estimation

# PSLP Algorithm

- Extension to SLP
- Generate multiple graphs (unlike SLP)
- Minimal Padding
- Cost estimation
- Emit redundant code to create isomorphism

# PSLP Algorithm

- Extension to SLP
- Generate multiple graphs (unlike SLP)
- Minimal Padding
- Cost estimation
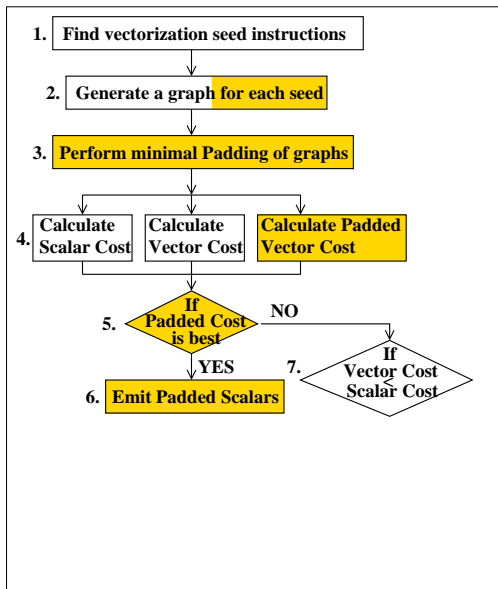- Emit redundant code to create isomorphism

# PSLP Algorithm

- Extension to SLP
- Generate multiple graphs (unlike SLP)
- Minimal Padding
- Cost estimation
- Emit redundant code to create isomorphism
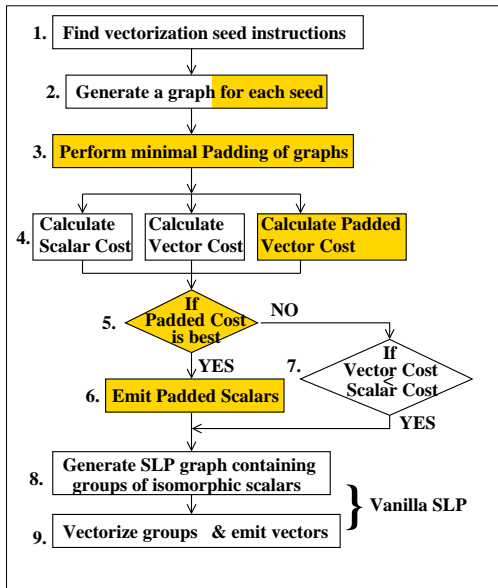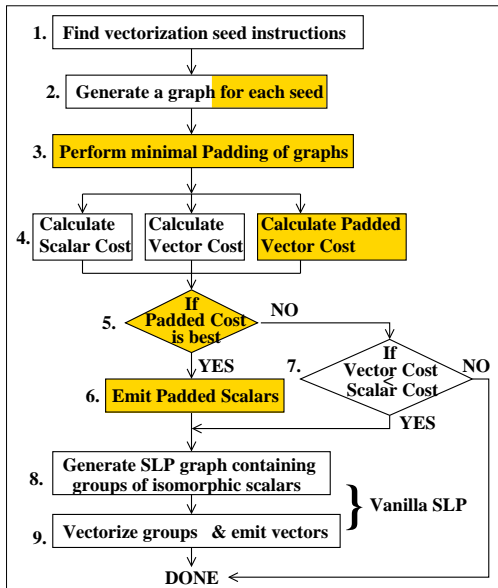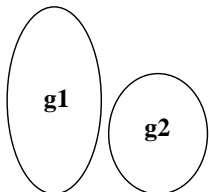- Code vectorized by original SLP
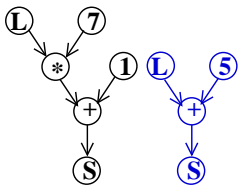
# PSLP Algorithm

- Extension to SLP
- Generate multiple graphs (unlike SLP)
- Minimal Padding
- Cost estimation
- Emit redundant code to create isomorphism
- Code vectorized by original SLP



1. Find vectorization seed instructions
2. Generate a graph for each seed
3. Perform minimal Padding of graphs
4. Calculate Scalar Cost | Calculate Vector Cost | Calculate Padded Vector Cost
5. If Padded Cost is best — NO
   YES
6. Emit Padded Scalars
7. If Vector Cost < Scalar Cost — NO
   YES
8. Generate SLP graph containing groups of isomorphic scalars
9. Vectorize groups & emit vectors
} Vanilla SLP
DONE

# Minimal Padding Algorithm

# Minimal Padding Algorithm

# Minimal Padding Algorithm

# Minimal Padding Algorithm

# Minimal Padding Algorithm

# Minimal Padding Algorithm
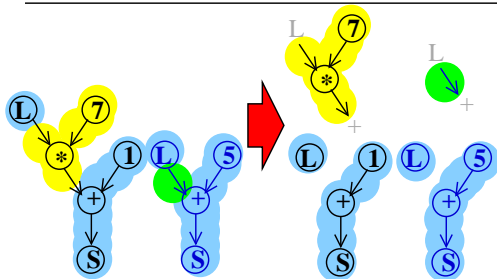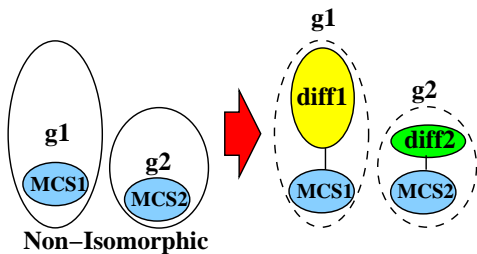
# Minimal Padding Algorithm

# Minimal Padding Algorithm

# We can do better: Remove redundant Selects



**EXAMPLE:  Instruction acting as Select**

# We can do better: Remove redundant Selects



**EXAMPLE: Instruction acting as Select**

# We can do better: Remove redundant Selects



**EXAMPLE: Instruction acting as Select**

# We can do better: Remove redundant Selects



**EXAMPLE: Instruction acting as Select**

# We can do better: Remove redundant Selects



**EXAMPLE: Instruction acting as Select**



**a. Instruction acting as Select**

# We can do better: Remove redundant Selects



**EXAMPLE: Instruction acting as Select**



**a. Instruction acting as Select**

# We can do better: Remove redundant Selects



**EXAMPLE: Instruction acting as Select**



**a. Instruction acting as Select**     **b. Select constants**

# We can do better: Remove redundant Selects



**EXAMPLE: Instruction acting as Select**

**a. Instruction acting as Select**    **b. Select constants**

# We can do better: Remove redundant Selects



**EXAMPLE: Instruction acting as Select**

**a. Instruction acting as Select**  **b. Select constants**  **c. Select same node**

# We can do better: Remove redundant Selects



**EXAMPLE: Instruction acting as Select**

**a. Instruction acting as Select**  **b. Select constants**  **c. Select same node**

# Opportunities for PSLP in real-life applications

- Non-isomorphic source code (e.g. computing conjugates in 433.milc)

```
b[0].real =   a[0].real
b[0].imag = − a[0].imag
b[1].real =   a[1].real
b[1].imag = − a[1].imag
```

| ... | **Memory** |
|---|---|
| a[0].real | |
| a[0].imag | |
| a[1].real | |
| a[1].imag | |
| ... | |

# Opportunities for PSLP in real-life applications

1. Non-isomorphic source code (e.g. computing conjugates in 433.milc)

```
b[0].real =   a[0].real
b[0].imag = ⊟ a[0].imag
b[1].real =   a[1].real
b[1].imag = ⊟ a[1].imag
```

| ... | **Memory** |
|-----------|---|
| a[0].real | |
| a[0].imag | |
| a[1].real | |
| a[1].imag | |
| ... | |

2. Isomorphic source code but non-isomorphic IR due to high-level optimizations (jdct of cjpeg)

```
tmp1 = quantval[0]*16384
tmp2 = quantval[1]*22725
tmp3 = quantval[2]*21407
tmp4 = quantval[3]*19266
```

# Opportunities for PSLP in real-life applications

1. Non-isomorphic source code (e.g. computing conjugates in 433.milc)

```
b[0].real =     a[0].real
b[0].imag = ⊟ a[0].imag
b[1].real =     a[1].real
b[1].imag = ⊟ a[1].imag
```

| ... | **Memory** |
|-----------|
| a[0].real |
| a[0].imag |
| a[1].real |
| a[1].imag |
| ... |

2. Isomorphic source code but non-isomorphic IR due to high-level optimizations (jdct of cjpeg)

```
tmp1 = quantval[0]*16384            tmp1 = quantval[0]<<14
tmp2 = quantval[1]*22725     opt    tmp2 = quantval[1]*22725
tmp3 = quantval[2]*21407            tmp3 = quantval[2]*21407
tmp4 = quantval[3]*19266            tmp4 = quantval[3]*19266
```

# Opportunities for PSLP in real-life applications

1. Non-isomorphic source code (e.g. computing conjugates in 433.milc)

```
b[0].real =   a[0].real
b[0].imag = ─ a[0].imag
b[1].real =   a[1].real
b[1].imag = ─ a[1].imag
```

| ... | **Memory** |
|-----|-----|
| a[0].real | |
| a[0].imag | |
| a[1].real | |
| a[1].imag | |
| ... | |

2. Isomorphic source code but non-isomorphic IR due to high-level optimizations (jdct of cjpeg)

```
tmp1 = quantval[0]*16384          tmp1 = quantval[0]<<14
tmp2 = quantval[1]*22725    opt   tmp2 = quantval[1]*22725
tmp3 = quantval[2]*21407          tmp3 = quantval[2]*21407
tmp4 = quantval[3]*19266          tmp4 = quantval[3]*19266
```

# Experimental Setup

- Implemented PSLP in the trunk version of the LLVM 3.6 compiler.

# Experimental Setup

- Implemented PSLP in the trunk version of the LLVM 3.6 compiler.
- Target: Intel Core i5-4570 @ 3.2Ghz

# Experimental Setup

- Implemented PSLP in the trunk version of the LLVM 3.6 compiler.
- Target: Intel Core i5-4570 @ 3.2Ghz
- Compiler flags: -O3 -allow-partial-unroll -march=core-avx2 -mtune-core-i7 -ffast-math

# Experimental Setup

- Implemented PSLP in the trunk version of the LLVM 3.6 compiler.
- Target: Intel Core i5-4570 @ 3.2Ghz
- Compiler flags: -O3 -allow-partial-unroll -march=core-avx2 -mtune-core-i7 -ffast-math
- Kernels, SPEC 2006 and Mediabench II
- We evaluated the following cases:

# Experimental Setup

- Implemented PSLP in the trunk version of the LLVM 3.6 compiler.
- Target: Intel Core i5-4570 @ 3.2Ghz
- Compiler flags: -O3 -allow-partial-unroll -march=core-avx2 -mtune-core-i7 -ffast-math
- Kernels, SPEC 2006 and Mediabench II
- We evaluated the following cases:
  1. All loop, SLP and PSLP vectorizers disabled (O3)

# Experimental Setup

- Implemented PSLP in the trunk version of the LLVM 3.6 compiler.
- Target: Intel Core i5-4570 @ 3.2Ghz
- Compiler flags: -O3 -allow-partial-unroll -march=core-avx2 -mtune-core-i7 -ffast-math
- Kernels, SPEC 2006 and Mediabench II
- We evaluated the following cases:
  1. All loop, SLP and PSLP vectorizers disabled (O3)
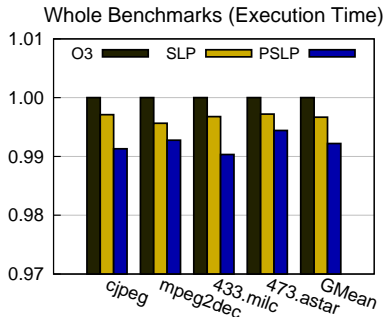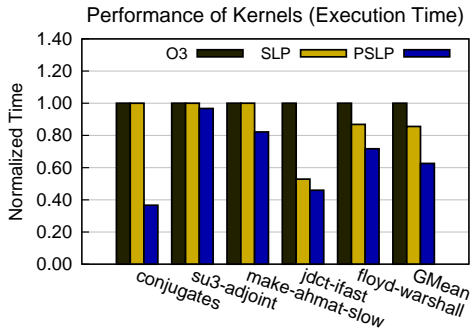  2. O3 + SLP enabled (SLP)

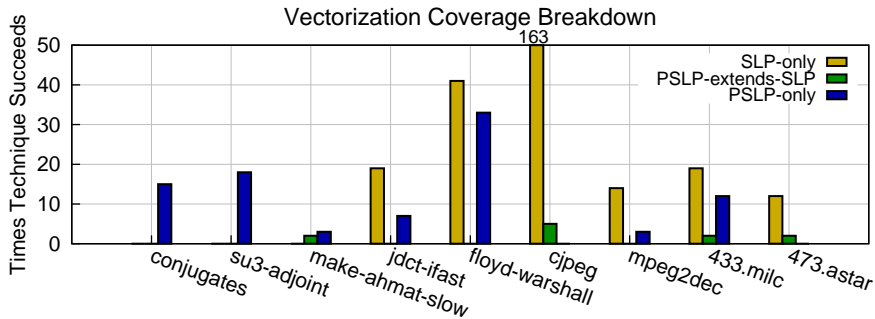# Experimental Setup

- Implemented PSLP in the trunk version of the LLVM 3.6 compiler.

- Target: Intel Core i5-4570 @ 3.2Ghz

- Compiler flags: -O3 -allow-partial-unroll -march=core-avx2 -mtune-core-i7 -ffast-math

- Kernels, SPEC 2006 and Mediabench II

- We evaluated the following cases:
  1. All loop, SLP and PSLP vectorizers disabled (O3)
  2. O3 + SLP enabled (SLP)
  3. O3 + PSLP enabled (PSLP)

# PSLP increases performance



Performance of Kernels (Execution Time)

Whole Benchmarks (Execution Time)

PSLP enables or extends vectorization

# PSLP enables or extends vectorization



Vectorization Coverage Breakdown

- SLP is adequate
- SLP stops at non-isomorphic code. PSLP extends it.

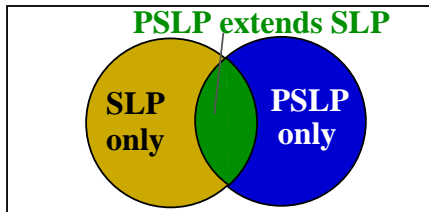# PSLP enables or extends vectorization

# Optimizing away redundant Selects

- Select-removal optimizations remove about 21% of the Selects



Percentage of Selects per region before and after Optimizations

# Conclusion

- PSLP improves vectorization coverage compared to the state-of-the-art

# Conclusion

- PSLP improves vectorization coverage compared to the state-of-the-art
- Converts non-isomorphic code into isomorphic by:

# Conclusion

- PSLP improves vectorization coverage compared to the state-of-the-art
- Converts non-isomorphic code into isomorphic by:
  - Relying on the Min Common Supergraph for minimal injection of redundant code

# Conclusion

- PSLP improves vectorization coverage compared to the state-of-the-art
- Converts non-isomorphic code into isomorphic by:
  - Relying on the Min Common Supergraph for minimal injection of redundant code
  - Emitting Select instructions to guarantee correctness

# Conclusion

- PSLP improves vectorization coverage compared to the state-of-the-art
- Converts non-isomorphic code into isomorphic by:
  - Relying on the Min Common Supergraph for minimal injection of redundant code
  - Emitting Select instructions to guarantee correctness
  - Optimizing away redundant Selects

# Conclusion

- PSLP improves vectorization coverage compared to the state-of-the-art
- Converts non-isomorphic code into isomorphic by:
  - Relying on the Min Common Supergraph for minimal injection of redundant code
  - Emitting Select instructions to guarantee correctness
  - Optimizing away redundant Selects
- PSLP performs better compared to SLP on commodity SIMD-capable hardware