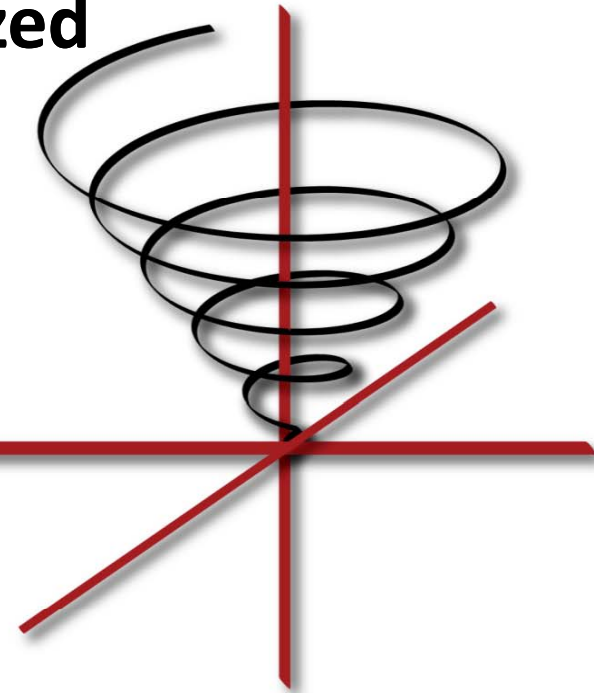


# Automatic Generation of Vectorized Fast Fourier Transform Libraries for the Larrabee and AVX Instruction Set Extension



**Daniel McFarlin**

Franz Franchetti

Markus Püschel

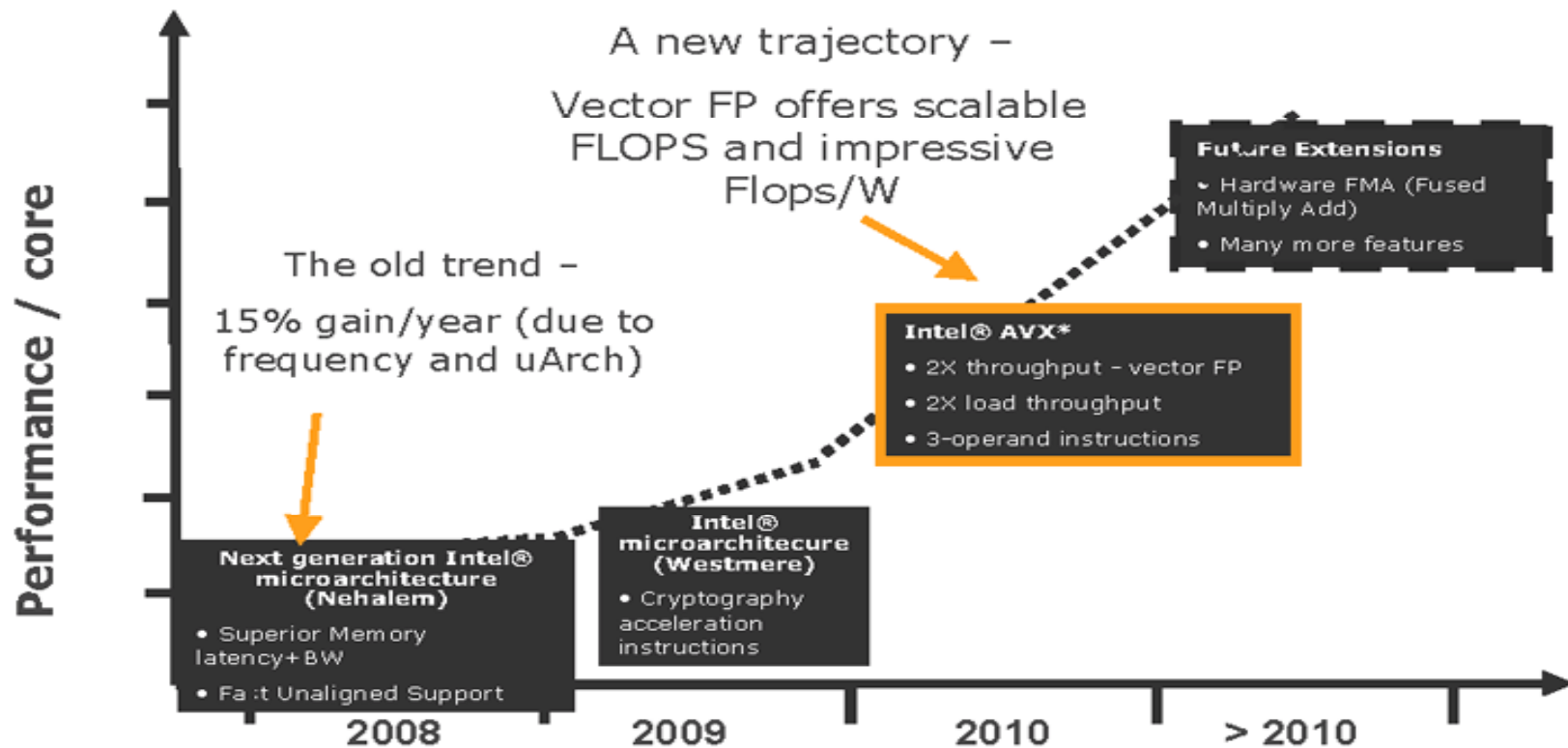
Carnegie Mellon University

HPEC, September 2009, Lexington, MA, USA

*This work was supported by DARPA DESA program, NSF-NGS/ITR, NSF-ACR, and Intel  
Daniel McFarlin was supported by NPSC and NDSEG Fellowships*

# Performance Trends

## Setting the Pace for Intel Instruction Set

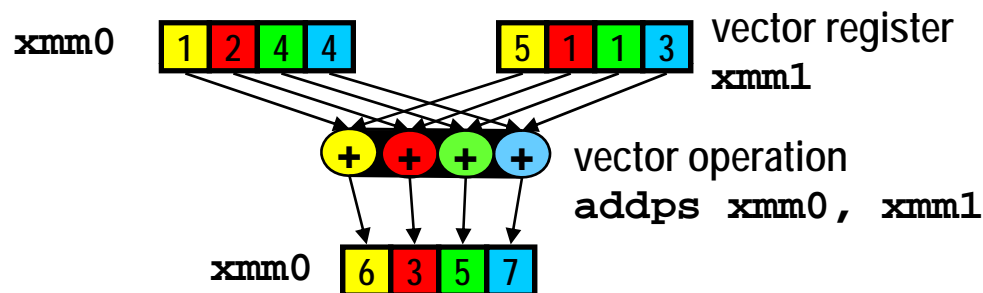


All timeframes, dates and products are subject to change without further notification

*Vector Units Will Dominate Performance on CPUs, GPUs and GPGPUs*

# Vectorization Challenges

- Must extract fine-grain data level parallelism



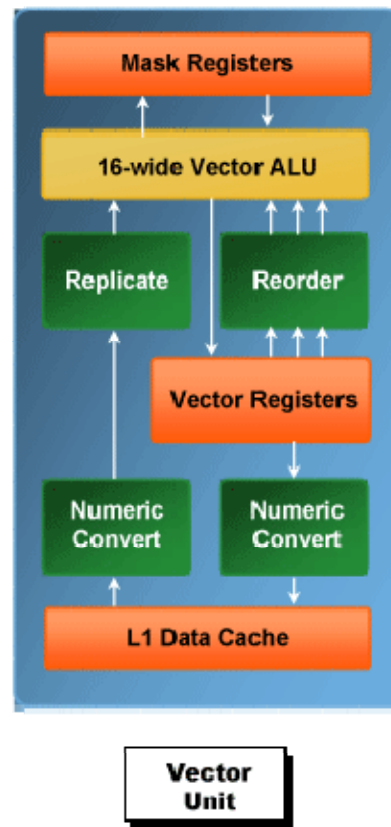
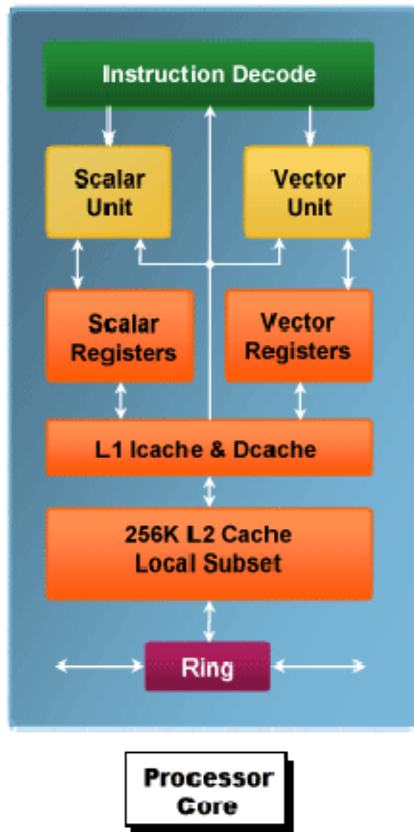
- **Problems:**

- Not standardized
- Compiler vectorization limited
- Low-level issues (data alignment,...)
- Reordering data kills runtime

- Intel MMX
- AMD 3DNow!
- Intel SSE
- AMD Enhanced 3DNow!
- Motorola AltiVec
- AMD 3DNow! Professional
- Itanium
- Intel XScale
- Intel SSE2
- AMD-64
- IBM BlueGene/L PPC440FP2
- Intel Wireless MMX
- Intel SSE3
- Intel AVX
- Intel LRBni

*One can easily slow down a program by vectorizing it*

# Larrabee



- Many enhanced in-order x86 cores
  - Simple Cores, Complex Instructions
- Completely New Vector Instructions
  - 512-bit vectors (16 x 32-bit floats)
  - Ternary, Fused multiply-add/sub
  - Writemasking/Predication
  - Load-op (third operand from mem)
  - Broadcasting
  - Swizzling (pre-defined shuffle ops)

*Code Must Be Vectorized to Attain High Performance*

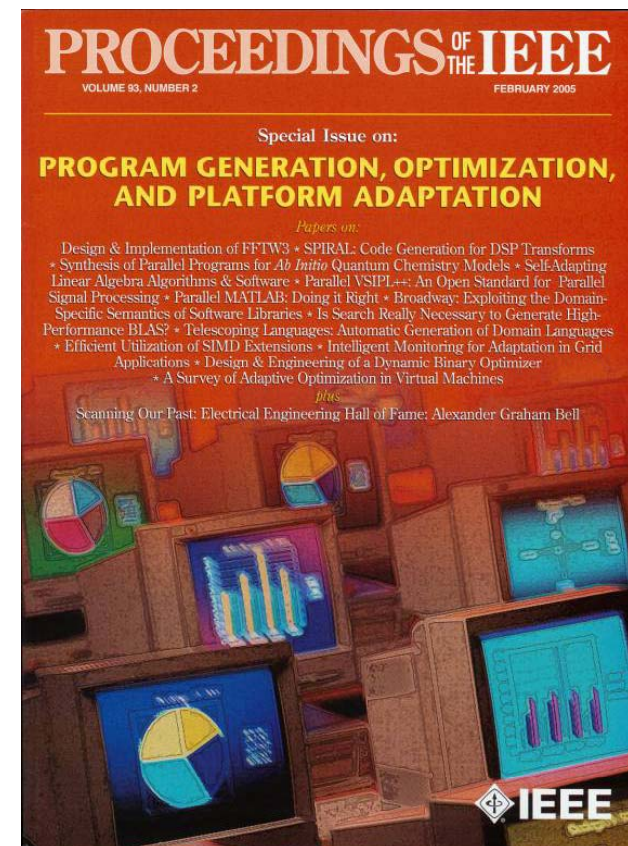
# Automatic Performance Tuning

- **Current vicious circle:** Whenever a new platform comes out, the same functionality needs to be rewritten and reoptimized

- **Automatic Performance Tuning**

- BLAS: ATLAS, PHiPAC
- Linear algebra: Sparsity/OSKI, Flame
- Sorting
- Fourier transform: FFTW
- **Linear transforms: Spiral**
- ...others
- New compiler techniques

*New challenge: ubiquitous parallelism*



Proceedings of the IEEE special issue, Feb. 2005



# Outline

- General Vectorization Challenges
- **Spiral**
- **Vectorization in Spiral**
- **Vectorization Challenges with AVX and Larrabee**
- **Results**
- **Conclusions**

# Spiral

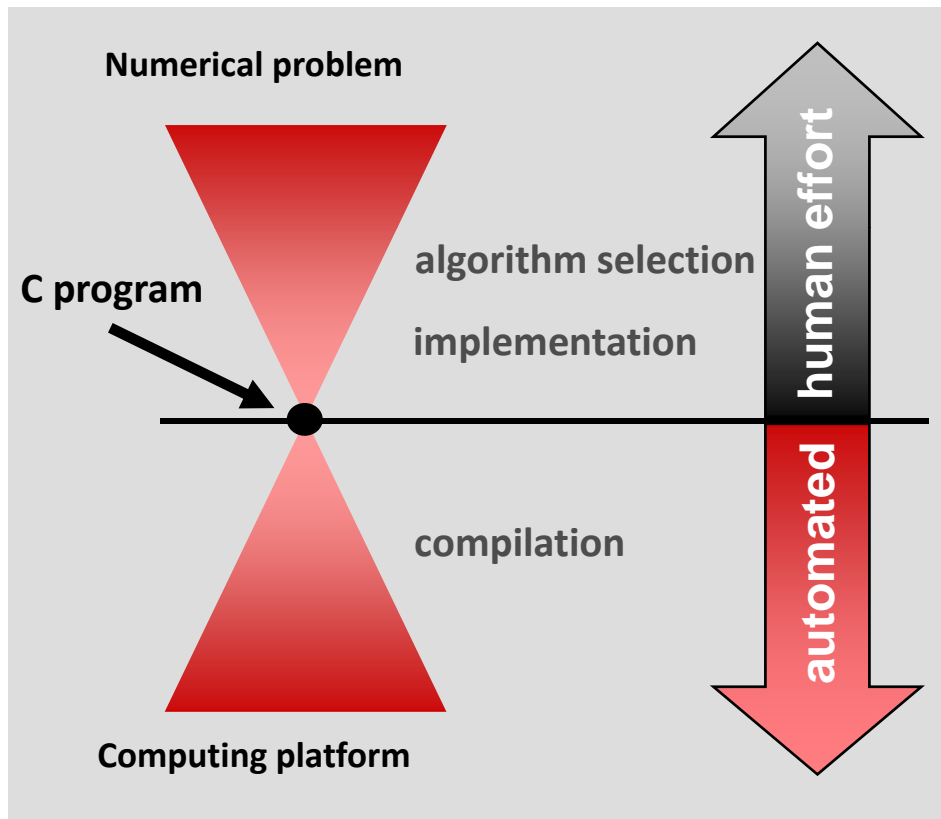
- Library generator for linear transforms (DFT, DCT, DWT, filters, ....) *and recently more ...*
- Wide range of platforms supported: scalar, fixed point, **vector, parallel, Verilog, GPU**
- **Research Goal: “Teach” computers to write fast libraries**
  - Complete automation of implementation and optimization
  - Conquer the “high” algorithm level for automation
- When a new platform comes out: **Regenerate a retuned library**
- When a new platform paradigm comes out (e.g., CPU+GPU): **Update the tool rather than rewriting the library**

*Intel uses Spiral to generate parts of their MKL and IPP libraries*



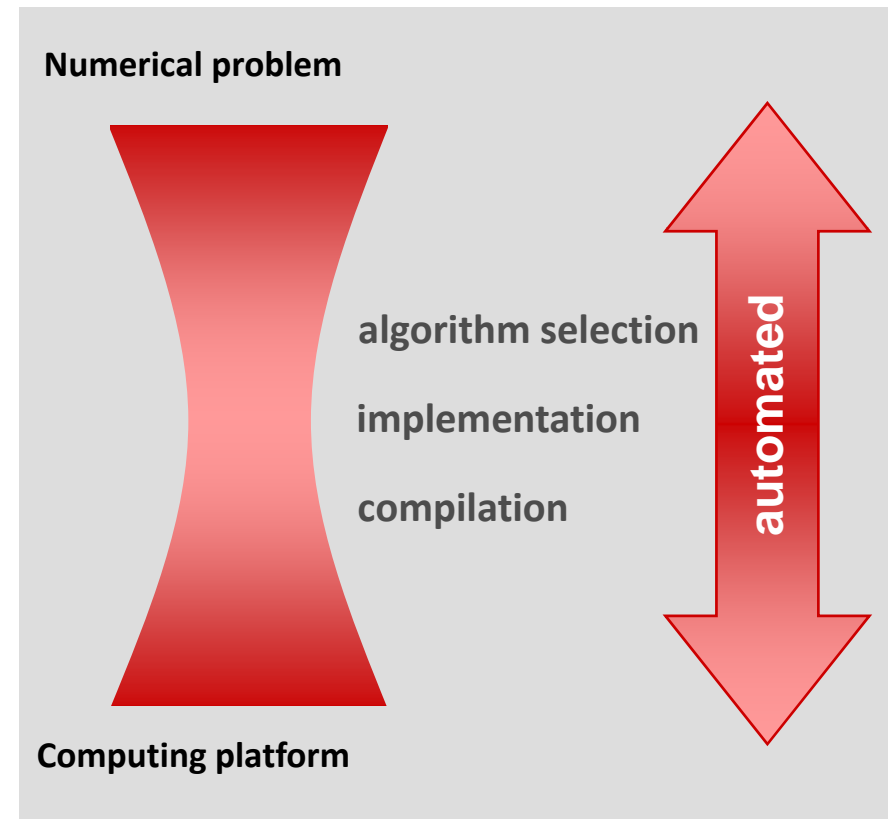
# Vision Behind Spiral

## Current



- C code a singularity: Compiler has no access to high level information

## Future



- Challenge: conquer the high abstraction level for **complete automation**

# What is a (Linear) Transform?

- Mathematically: Matrix-vector multiplication

$$x \mapsto y = T \cdot x$$

input vector (signal)  $x$     output vector (signal)  $y$     transform = matrix  $T$

- Example: Discrete Fourier transform (DFT)

$$\text{DFT}_n = [e^{-2k\ell\pi i/n}]_{0 \leq k, \ell < n}$$

# Transform Algorithms: Example 4-point FFT

Cooley/Tukey fast Fourier transform (FFT):

$$\begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & j & -1 & -j \\ 1 & -1 & 1 & -1 \\ 1 & -j & -1 & j \end{bmatrix} = \begin{bmatrix} 1 & \cdot & 1 & \cdot \\ \cdot & 1 & \cdot & 1 \\ 1 & \cdot & -1 & \cdot \\ \cdot & 1 & \cdot & -1 \end{bmatrix} \begin{bmatrix} 1 & \cdot & \cdot & \cdot \\ \cdot & 1 & \cdot & \cdot \\ \cdot & \cdot & 1 & \cdot \\ \cdot & \cdot & \cdot & j \end{bmatrix} \begin{bmatrix} 1 & 1 & \cdot & \cdot \\ 1 & -1 & \cdot & \cdot \\ \cdot & \cdot & 1 & 1 \\ \cdot & \cdot & 1 & -1 \end{bmatrix} \begin{bmatrix} 1 & \cdot & \cdot & \cdot \\ \cdot & \cdot & 1 & \cdot \\ \cdot & 1 & \cdot & \cdot \\ \cdot & \cdot & \cdot & 1 \end{bmatrix}$$

Fourier transform

Diagonal matrix (twiddles)

$$\text{DFT}_4 = (\text{DFT}_2 \otimes \text{I}_2) \text{T}_2^4 (\text{I}_2 \otimes \text{DFT}_2) \text{L}_2^4$$

Kronecker product
Identity
Permutation

- Algorithms reduce arithmetic cost  $O(n^2) \rightarrow O(n \log(n))$
- Product of structured sparse matrices
- Mathematical notation exhibits structure: **SPL (signal processing language)**

# Program Generation in Spiral (Sketched)

Transform  
user specified

DFT<sub>8</sub>



Fast algorithm  
in SPL  
many choices

$$(DFT_2 \otimes I_4) T_4^8 (I_2 \otimes ((DFT_2 \otimes I_2) \cdot T_2^4 (I_2 \otimes DFT_2) L_2^4)) L_2^8$$



Σ-SPL:

$$\sum (S_j DFT_2 G_j) \sum \left( \sum (S_{k,l} \text{diag}(t_{k,l}) DFT_2 G_l) \sum (S_m \text{diag}(t_m) DFT_2 G_{k,m}) \right)$$



C Code:

```
void sub(double *y, double *x) {
    double f0, f1, f2, f3, f4, f7, f8, f10, f11;
    f0 = x[0] - x[3];
    f1 = x[0] + x[3];
    f2 = x[1] - x[2];
    f3 = x[1] + x[2];
    f4 = f1 - f3;
    y[0] = f1 + f3;
    y[2] = 0.7071067811865476 * f4;
    f7 = 0.9238795325112867 * f0;
    f8 = 0.3826834323650898 * f2;
    y[1] = f7 + f8;
    f10 = 0.3826834323650898 * f0;
    f11 = (-0.9238795325112867) * f2;
    y[3] = f10 + f11;
}
```

Optimization at all  
abstraction levels



parallelization  
vectorization



loop  
optimizations



constant folding  
scheduling

.....

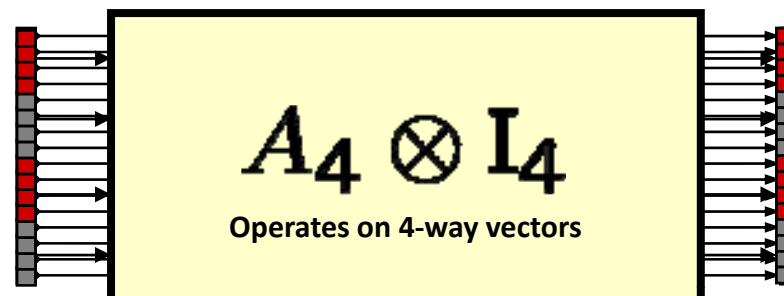
# Vectorization In Spiral

- Naturally vectorizable construct

Franchetti and Püschel (IPDPS 2002/2003)

$$y = (A \otimes I_\nu)x$$

vector length (any two-power)



- Rewriting rules to vectorize formulas

Introduces **data reorganization (permutations)**

$$I_n \otimes A^{k \times m} \rightarrow I_{n/\nu} \otimes \underline{L_\nu^{k\nu}} \left( \underline{A^{k \times m} \otimes I_\nu} \right) \underline{L_m^{m\nu}}$$

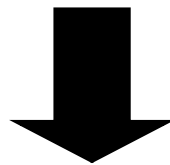
$$L_m^{m\nu} \rightarrow \left( \underline{I_{m/\nu} \otimes L_\nu^{\nu^2}} \right) \left( \underline{L_{m/\nu}^m \otimes I_\nu} \right)$$

vector construct  
further rewriting  
base permutation

# Vectorized DFT

## Standard FFT

$$\text{DFT}_{mn} \rightarrow (\text{DFT}_m \otimes \text{I}_n) \text{T}_n^{mn} (\text{I}_m \otimes \text{DFT}_n) \text{L}_m^{mn}$$



Automatic formula rewriting

$$\text{DFT}_{mn} \rightarrow \left( \text{I}_{\frac{mn}{\nu}} \otimes \underline{\text{L}_{\nu}^{2\nu}} \right) \left( \underline{\text{DFT}_m \otimes \text{I}_{\frac{n}{\nu}} \otimes \text{I}_{\nu}} \right) \overline{\text{T}_n^{mn}}$$

$$\left( \text{I}_{\frac{m}{\nu}} \otimes \left( \underline{\text{L}_{\nu}^{2n}} \otimes \text{I}_{\nu} \right) \left( \text{I}_{\frac{2n}{\nu}} \otimes \underline{\text{L}_{\nu}^{\nu^2}} \right) \left( \underline{\text{I}_{\frac{n}{\nu}} \otimes \text{L}_2^{2\nu}} \otimes \text{I}_{\nu} \right) \left( \underline{\text{DFT}_n \otimes \text{I}_{\nu}} \right) \right) \left( \text{L}_{\frac{m}{\nu}}^{\frac{mn}{\nu}} \otimes \underline{\text{L}_2^{2\nu}} \right)$$

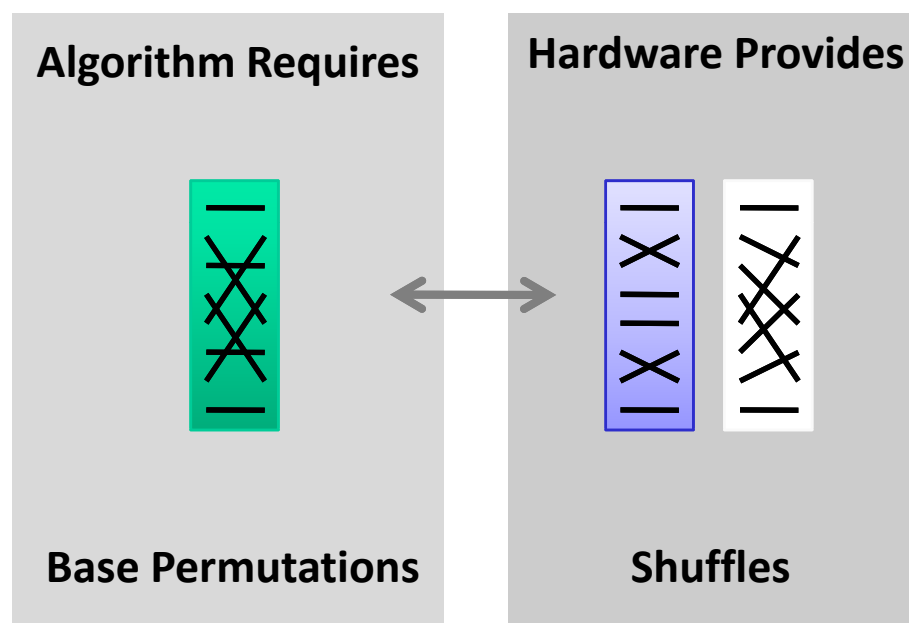
Vectorized arithmetic

Architecture specific

data reorganization

$$\text{L}_2^{2\nu}, \text{L}_{\nu}^{2\nu}, \text{L}_{\nu}^{\nu^2}$$

# Base Permutations To Code



*Goal: Find Fast Instruction Sequences That Implement Base Permutations*

# Base Permutation Example: $L_{\nu}^{2\nu}$ $\nu=4$



$$\begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 \end{bmatrix}$$

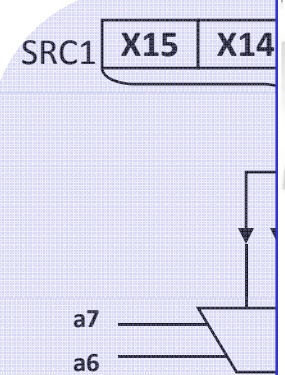
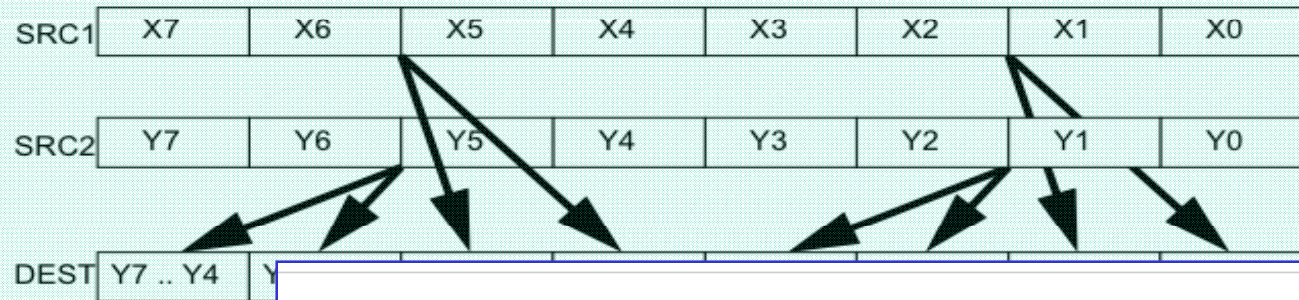
$$\begin{bmatrix} 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 \end{bmatrix}$$

```
y0 = _mm_unpacklo_ps(x0, x1);
```


```
y1 = _mm_unpackhi_ps(x0, x1);
```



# A Tale of Two Shuffles



Some Date



Anonymous

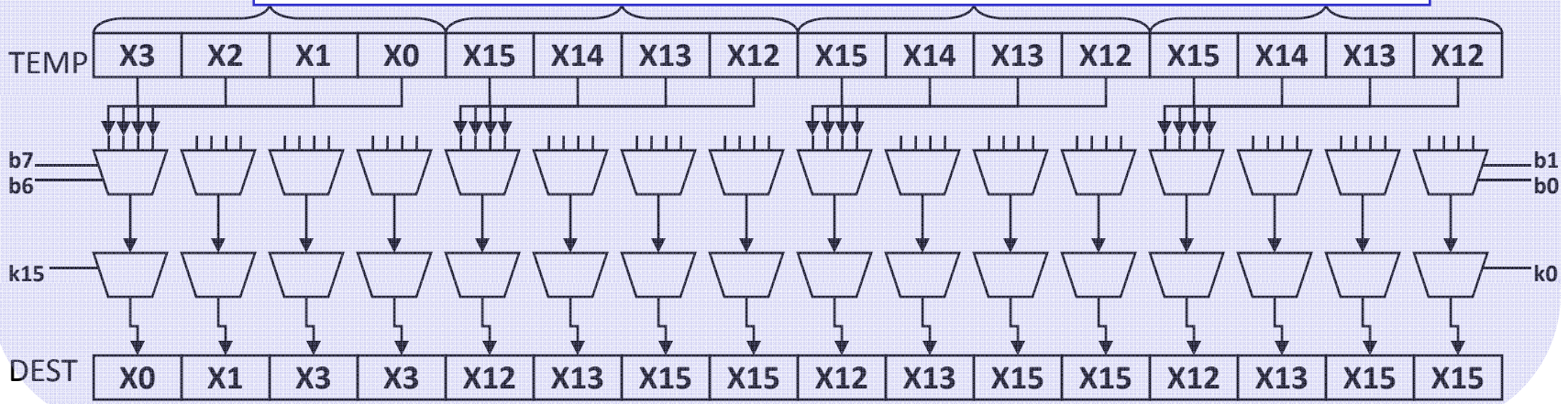
What instruction(s) can I use to achieve the following:  
`_M512 va = {a0,a1,a2,a3, a4,a5,a6,a7, a8,a9,a10,a11, a12,a13,a14,a15} ;`  
`_M512 vb = {b0,b1,b2,b3,b4,b5,b6,b7,b8,b9,b10,b11,b12,b13,b14,b15} ;`

and I want

`_M512 vc = {a0,b0, a1,b1, a2,b2, a3,b3, a4,b4, a5,b5, a6,b6, a7,b7};`

What if I want vc to have some random selection of the float values in va and vb?  
 The shuffle instruction (`_mm512_shuf128x32`) does not seem to be able to do that.

X0

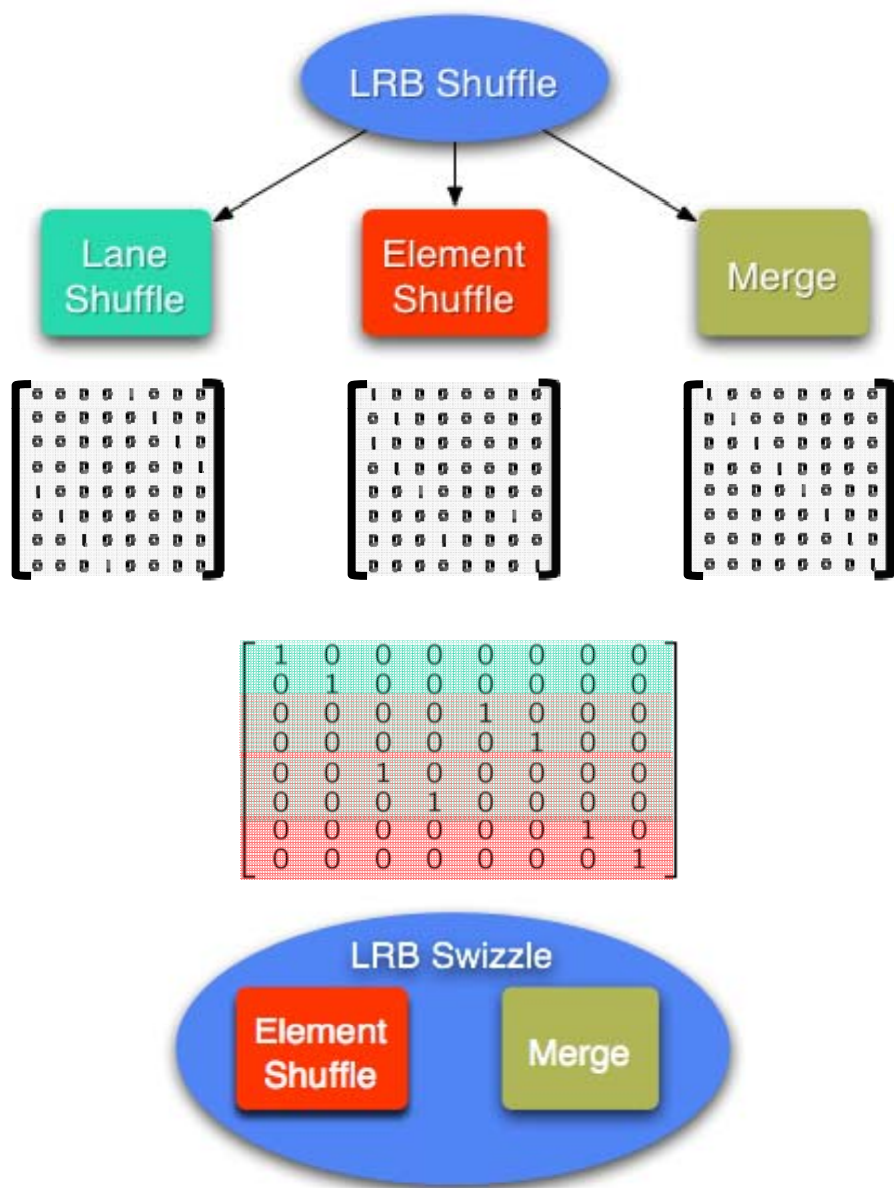


# Implementing Base Permutation Matrices

$$\begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 \end{bmatrix} = \begin{bmatrix} & & & & & & & \\ & & & & & & & \\ & & & & & & & \\ & & & & & & & \\ & & & & & & & \\ & & & & & & & \\ & & & & & & & \\ & & & & & & & \end{bmatrix} \begin{bmatrix} & & & & & & & \\ & & & & & & & \\ & & & & & & & \\ & & & & & & & \\ & & & & & & & \\ & & & & & & & \\ & & & & & & & \\ & & & & & & & \end{bmatrix}$$

- **Equivalent to Factoring A Binary Matrix**
  - **Computationally Intractable**
- **Encode instructions as matrices**
  - Generate and multiply matrix sequences
  - Minimize sequence length and cost
  - **Four Billion Variants of the LRB Shuffle Instruction!**

# Search Space Optimization

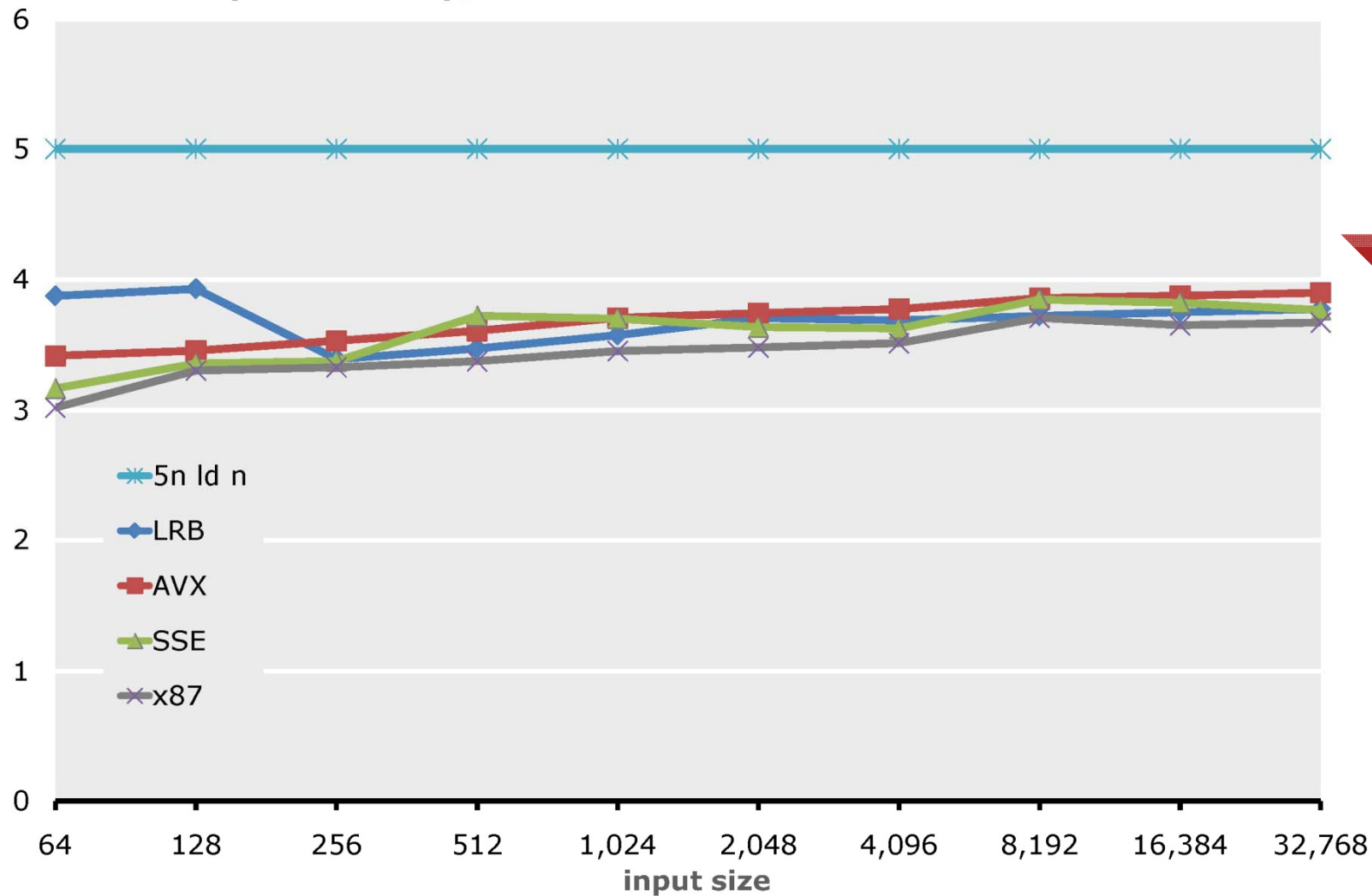


- **Micro-Op Reduction**
  - Use heuristics to filter  $\mu$ -ops
- **Fast Binary Matrix Multiply**
- **Parallelize Search Space**
- **Strength Reduction**
  - Micro-Op Fusion
  - Pattern Matching
  - Prefix Trees
    - `unpackhi/lo`

# Results: Numerical Cost

## Normalized FFT Cost (single-precision)

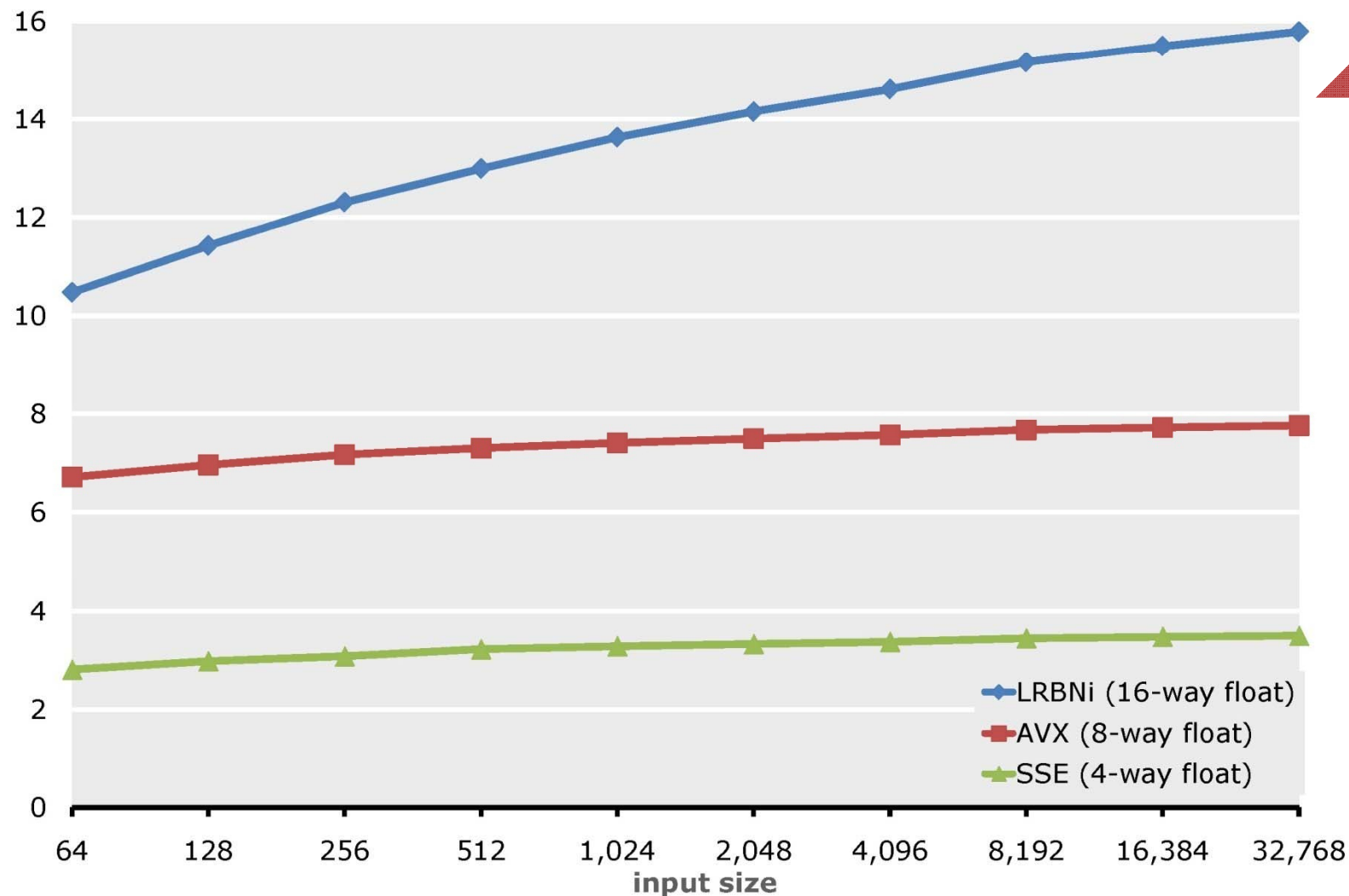
normalized opcount in flop/  $n \log_2 n$



# Results: Vectorization Efficiency

## Operations Count Reduction (single-precision)

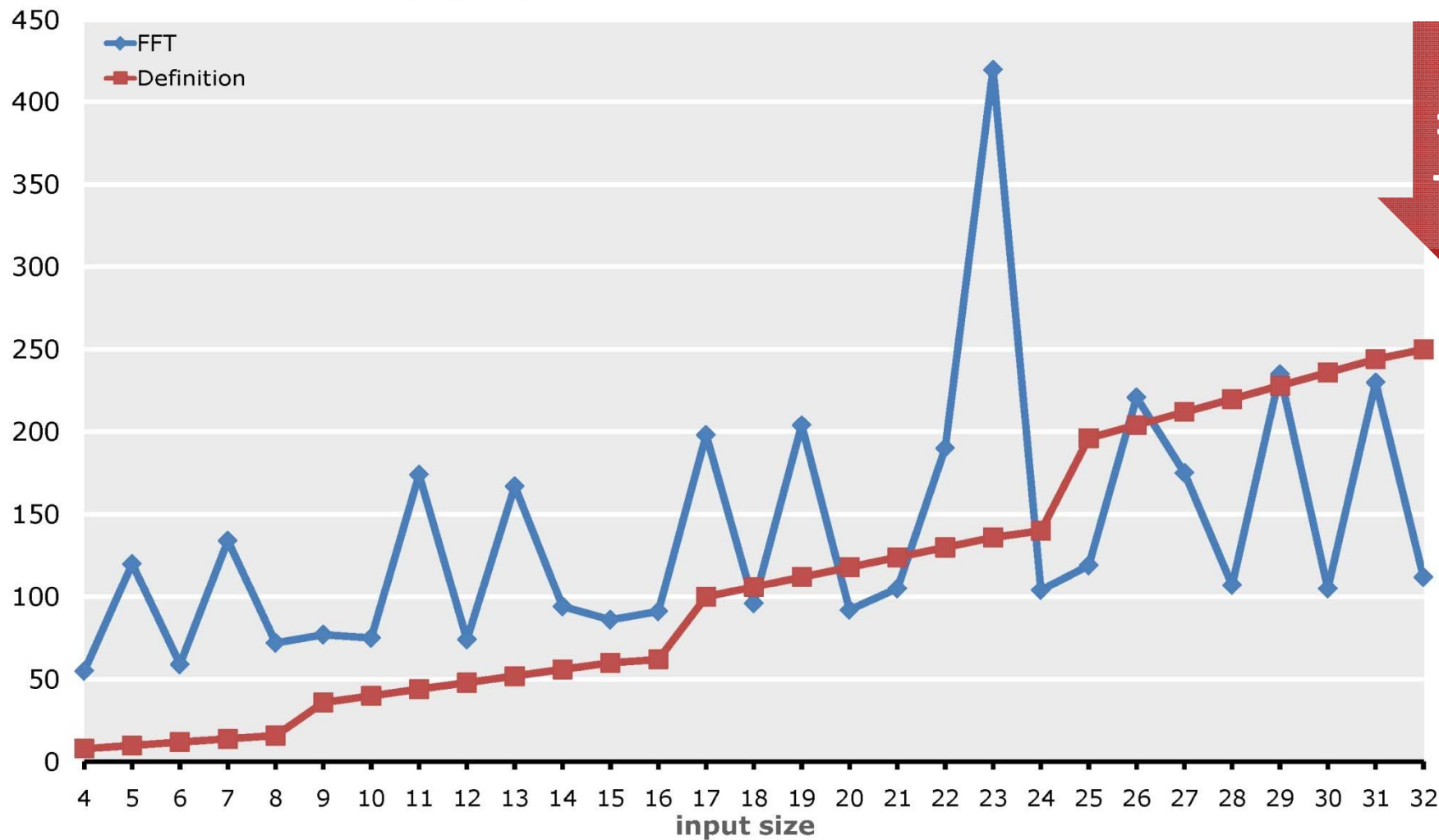
ratio of x87 to Vector Architecture opcounts



# Results: Flexible Algorithms

## DFT Vector Opcount (single precision)

arithmetic + shuffle instructions



# Conclusions

- **New processors increase vector ISA power but also ISA complexity**
- **Reordering operations dominate vector performance**
- **Minimize reordering op cost through optimized search**
- **Fully Automate Spiral Vectorization Framework**
- **Enable production of High Quality Code**