

High-Performance Code Generation for FIR Filters and the Discrete Wavelet Transform Using SPIRAL

Aca Gacic

Markus Püschel

José M. F. Moura

SPIRAL project

<http://www.spiral.net>

Electrical and Computer Engineering Department

Carnegie Mellon University



Motivation

- ✍ **FIR filters** - image enhancement, equalization, speech synthesis, biomedical signal processing, ...
- ✍ **Wavelets** - image compression, detection, denoising, signal recovery
- ✍ Numerically intensive - **significant impact on application performance**

High-Performance Implementations (How-To)

✍ Choose **Fast Algorithm**

- ✍ Arithmetic cost – only a rough estimate of performance
- ✍ **Many algorithms** with similar cost – which one?

✍ Design **Efficient Code**

- ✍ Architecture conscious – **machine dependent**
- ✍ **Obsolete** when platform is changed/upgraded

✍ **Best implementation:** Algorithm + Machine + Compiler

- ✍ Requires experts in algorithms and computer architecture
- ✍ Frequent **re-implementation**

Better way: automatic performance tuning

Existing FIR filtering and DWT software

Matlab[®], Mathematica, S+Wavelets, WaveLab, Wave++, IPP

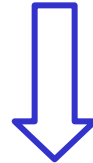
Application oriented – not machine specific, or

Hand-tuned code for specific platform

Automatic platform-adapted solutions not available !



We want to close this gap



Automatic performance tuning packages

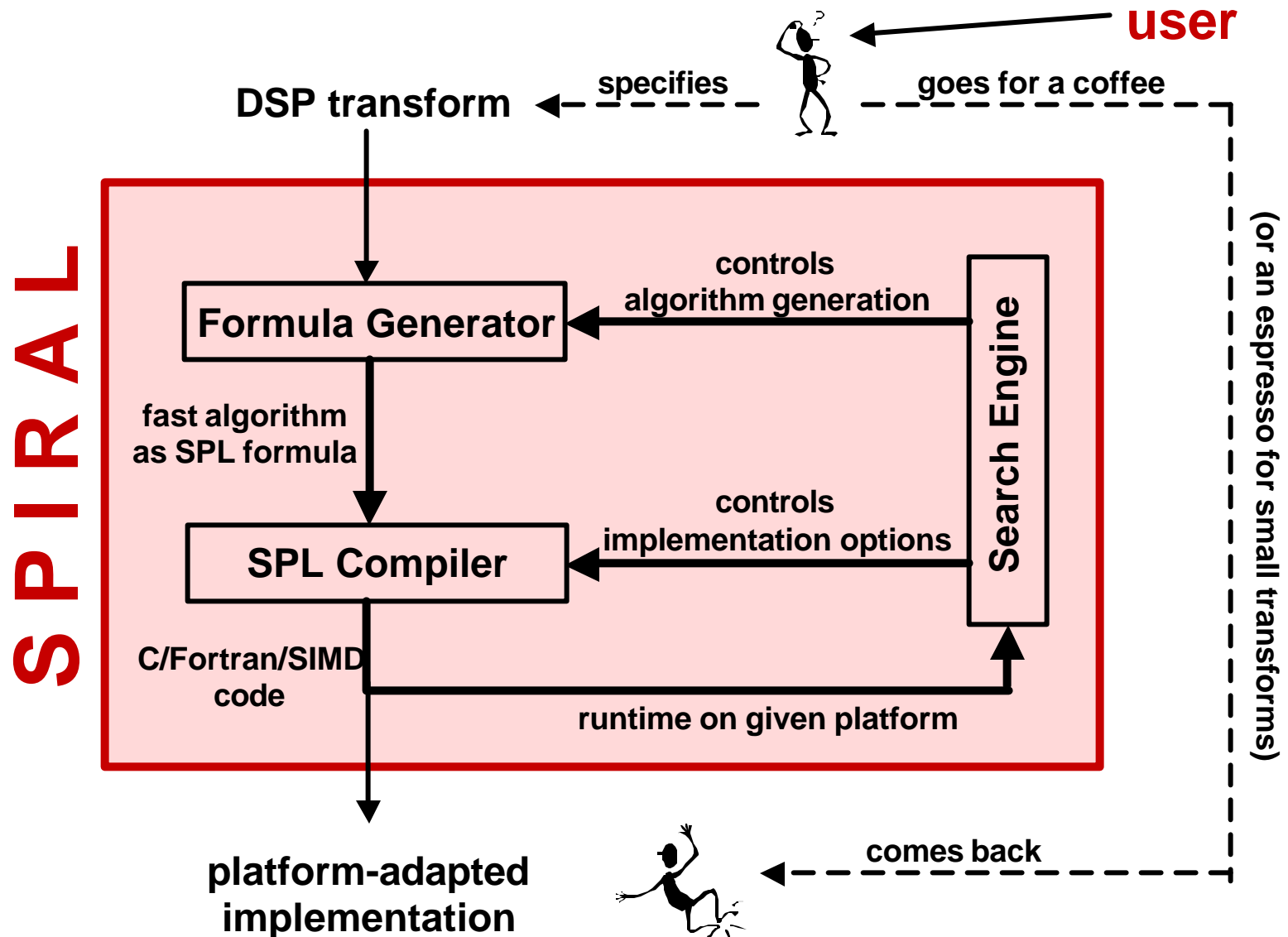
ATLAS, PhiPAC, SPARSITY, FFTW, SPIRAL

Include several basic linear algebra and DSP functions

No filtering and wavelet kernels !

SPIRAL system

Generator of optimized DSP transform implementations



SPIRAL's Four Key Concepts

Transform

$$DFT_n = \{e^{j2\pi k l / n}\}_{k,l=0,\dots,n-1}$$

parameterized matrix

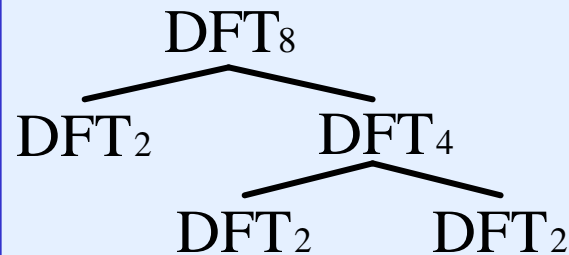
Rule

$$DFT_{nm} = DFT_n \otimes I_m \cdot T_n^{nm} \cdot I_n \cdot DFT_m \cdot L_m^{nm}$$

Diagonal matrix (twiddles)
Kronecker product
Identity
Permutation

- a breakdown strategy - product of sparse matrices
- captures important structural information

Rule Tree



- recursive application of rules
- uniquely defines an algorithm
- efficient representation
- easy manipulation

Formula

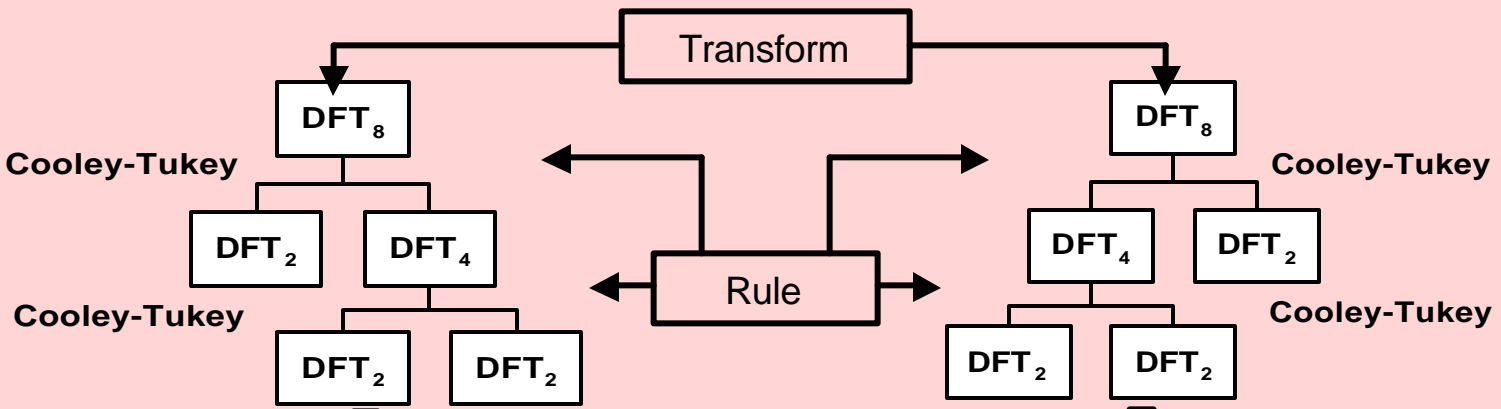
$$DFT_8 = F_2 \otimes I_4 \cdot T_4^8 \cdot I_2 \cdot I_2 \cdot F_2 \cdot L_2^8$$

- few constructs and primitives
- uniquely defines an algorithm
- can be translated into code

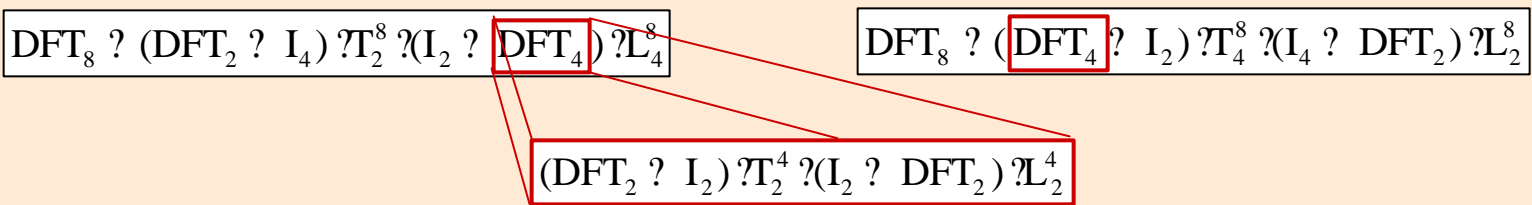


Rule Trees = Formulas = Algorithms

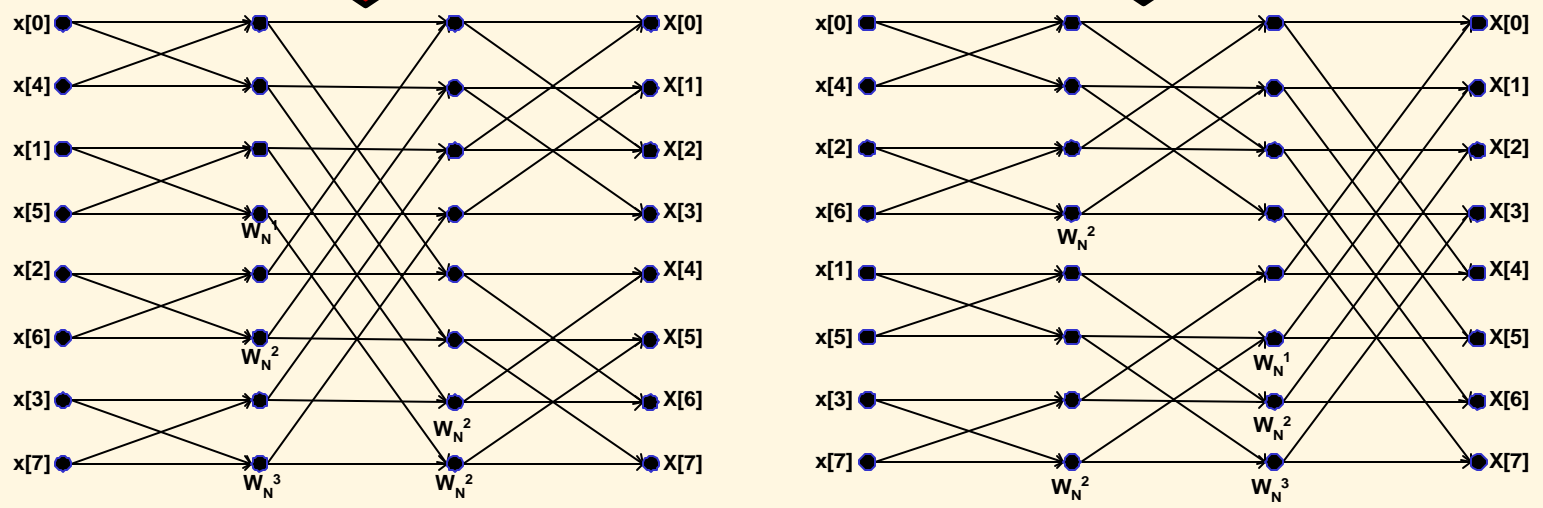
Rule Trees



Formulas



Algorithms



Increasing level of abstraction



Examples: Rules and Rule Trees

Trigonometric Rule

$$DFT_n \rightarrow \text{CosDFT}_n \rightarrow i \cdot \text{SinDFT}_n$$

DCT-Recursive Rule

$$\text{CosDFT}_n \rightarrow S_n \rightarrow \text{CosDFT}_{n/2} \rightarrow DCT_{n/4}^{(II)} \rightarrow S_{n/4} \rightarrow L_2^n$$

$$\text{SinDFT}_n \rightarrow S_n \rightarrow \text{SinDFT}_{n/2} \rightarrow DCT_{n/4}^{(II)} \rightarrow S_{n/4} \rightarrow L_2^n$$

Type-Split Rule

$$DCT_n^{(II)} \rightarrow P_n \rightarrow DCT_{n/2}^{(II)} \rightarrow DCT_{n/2}^{(IV)} \rightarrow I_{n/2} \rightarrow F_2 \rightarrow P_n$$

Type-Conversion Rule

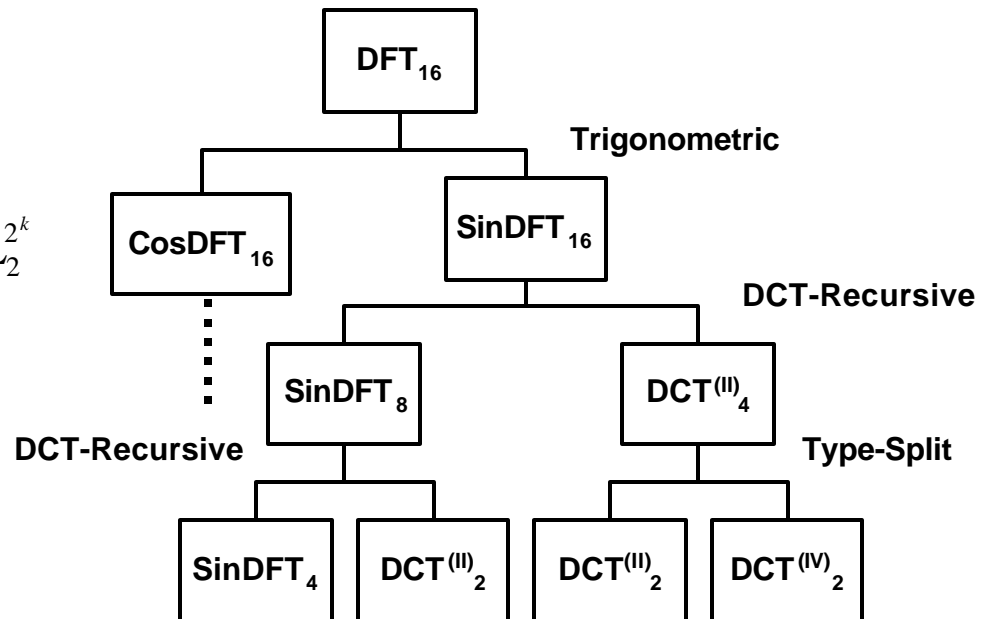
$$DCT_n^{(IV)} \rightarrow S_n \rightarrow DCT_n^{(II)} \rightarrow D_n$$

Recursive RHT Rule

$$RHT_{2^k} \rightarrow RHT_{2^{k-1}} \rightarrow I_{2^{k-1}} \rightarrow F_2 \rightarrow I_{2^{k-1}} \rightarrow L_2^{2^k}$$

Number of rule trees grows exponentially with

- Size of the transform
- Number of applicable rules



SPIRAL also includes search over implementation options, such as the **degree of loop unrolling**

FIR Filters and the DWT in SPIRAL: Challenges

- ✍ Existing filtering and wavelet algorithm representations do not capture all **important structural information**
- ✍ SPIRAL uses a **concise math framework** to represent many algorithms for several transforms (DFT, DCTs, DHT, RDFT, ...)
- ✍ Filtering and wavelet algorithms have **considerably different structure** from the current SPIRAL transforms
 - ✍ Current framework **not enough** to capture algorithms' structure
 - ✍ Existing algorithm representations need to be “spiralized” – i.e., captured concisely using the **rule formalism**

- ✍ Need changes on several levels:
 - ✍ New transforms, constructs, rules, translation templates
 - ✍ Modified search
 - ✍ Need to rephrase several concepts of the framework
- ✍ Existing framework has to be **redesigned** to accommodate new constructions and decomposition strategies

New framework has to be integrated in SPIRAL

FIR Filters: Rules and Algorithms

FIR Filter Transform

$$y_n = \sum_{i=0}^{k-1} h_i x_{n-i}$$

standard "sum" notation

$$F_n(\mathbf{h}) = \begin{bmatrix} h_0 & & & & \\ h_1 & h_0 & & & \\ \square & h_1 & \square & & \\ \square & \square & \square & \square & \\ \square & h_{k-1} & \square & \square & \\ \square & \square & \square & \square & \\ \square & \square & \square & \square & \\ \square & \square & \square & \square & \\ \square & \square & \square & \square & \\ \square & \square & \square & \square & \end{bmatrix} \begin{bmatrix} x_n \\ x_{n-1} \\ x_{n-2} \\ x_{n-3} \\ x_{n-4} \\ x_{n-5} \\ x_{n-6} \\ x_{n-7} \\ x_{n-8} \\ x_{n-9} \end{bmatrix}$$

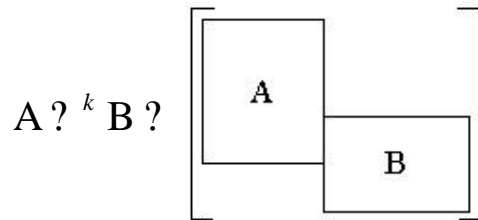
SPIRAL transform notation

Baseline Rule

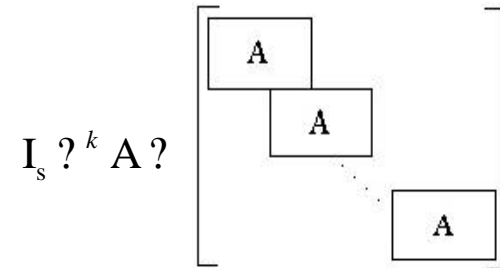
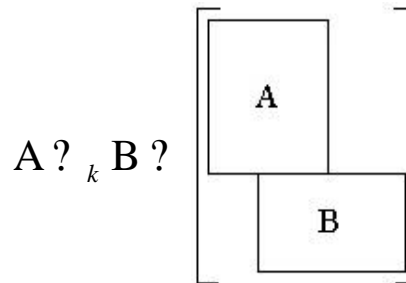
$$F_n \mathbf{h} I_n^k \mathbf{h}^T$$

SPIRAL rule
captures the structure

New constructs



column and row overlapped direct sum



overlapped tensor product

Time Domain Methods

Overlap-Add Rule

$$F_n \mathbf{h} I_{n/b}^k F_b \mathbf{h}$$

concise rule notation

- Localized access of the input data
- Input is blocked into length b segments





Overlap-Save Rule

$$F_n \mathbf{h} \quad T \mathbf{h}_L \quad I_{m/b} \quad F_b^T \mathbf{h} \quad T \mathbf{h}_R$$

|
Toeplitz matrix

- **Localized storage** of the output
- Output segments are computed independently

Blocking Rule

$$F_n \mathbf{h} \quad I_{n/b} \quad I_{b/b} \quad T \mathbf{h}_i$$

- **Multilevel blocking**
- Control over **input & output locality**

Frequency Domain Methods

Expanded Circulant Rule

$$F_n \mathbf{h} \quad C \quad E_{k,n} \mathbf{h} \quad E_{n,k} \quad E_{n,k} \quad I_n$$

|
Circulant matrix |
sparse matrix |
real DFT

Cyclic-RDFT Rule

$$C \mathbf{a} \quad \text{RDFT}_n^{-1} \quad X \quad \text{RDFT}_n \mathbf{a} \quad \text{RDFT}_n$$

|
Circulant matrix |
sparse matrix |
real DFT

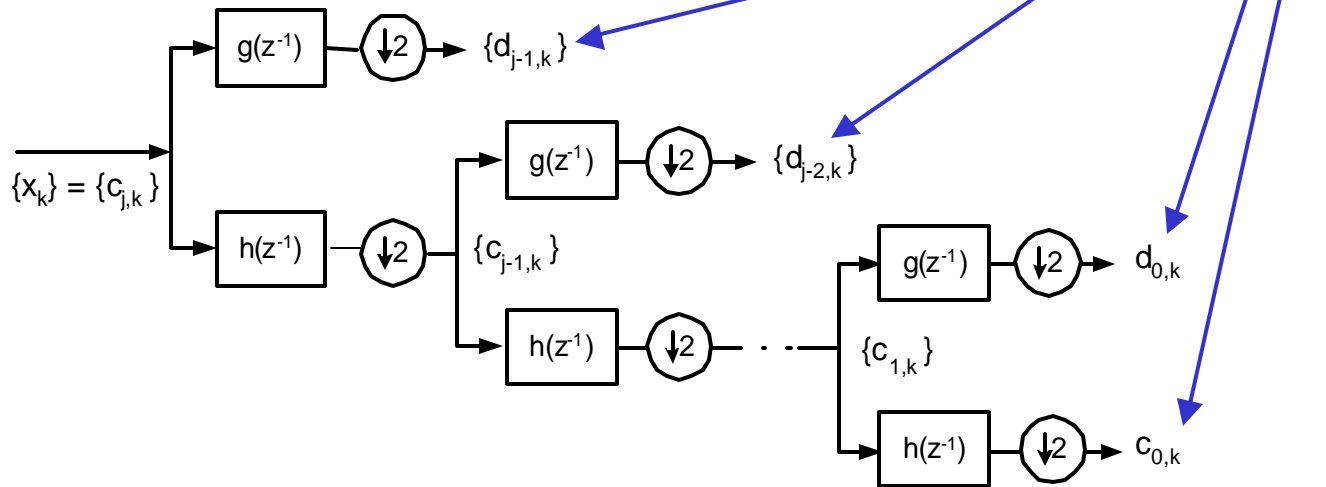
Radix-2 Cooley-Tukey RDFT Rule

$$\text{RDFT}_n \quad \text{RDFT}_{n/2} \quad I_2 \quad X(w^{kl}) \quad I_{n/2} \quad \text{RDFT}_2 \quad L_2^n$$

Time domain - better data access structure
Frequency domain - potentially reduce arithmetic cost

DWT: Rules and Algorithms

Filter Bank – Mallat's Algorithm



Mallat Rule

$$DWT_n(\mathbf{h}, \mathbf{g}) = DWT_{n/2}(\mathbf{h}, \mathbf{g}) \begin{bmatrix} \mathbf{I}_{n/2} & \mathbf{0} \\ \mathbf{0} & \mathbf{2} \end{bmatrix} F_{n/2}(\mathbf{h}) \begin{bmatrix} \mathbf{I}_{n/2} & \mathbf{0} \\ \mathbf{0} & \mathbf{2} \end{bmatrix} F_{n/2}(\mathbf{g}) \mathbf{E}_{n,l,r}^*$$

Extension matrix

Fused Mallat Rule

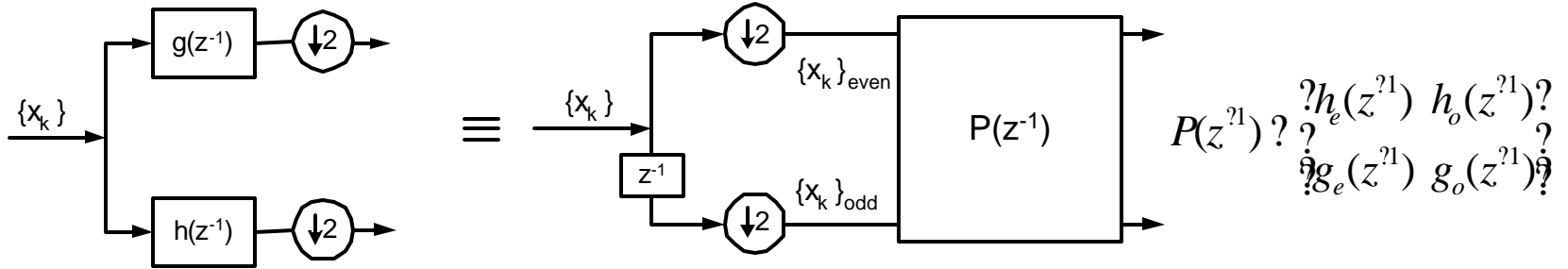
$$DWT_n(\mathbf{h}, \mathbf{g}) = DWT_{n/2}(\mathbf{h}, \mathbf{g}) \begin{bmatrix} \mathbf{I}_{n/2} & \mathbf{0} \\ \mathbf{0} & \mathbf{2} \end{bmatrix} H_n(\mathbf{h}, \mathbf{g}) \mathbf{E}_{n,l,r}^*$$

$$H_n(\mathbf{h}, \mathbf{g}) = \begin{bmatrix} \mathbf{I}_{n/2} & \mathbf{0} \\ \mathbf{0} & \mathbf{2} \end{bmatrix} \begin{bmatrix} \mathbf{h}^T \\ \mathbf{g}^T \end{bmatrix}$$

single-level DWT

- Redundant computations removed
- Filter structure lost

Polyphase Representation

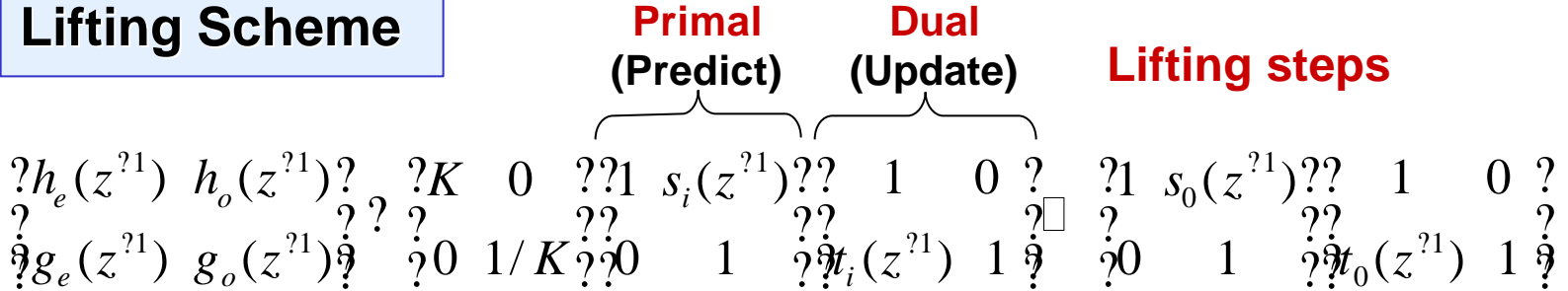


Polyphase Rule

Gateway to **frequency domain** computation

$$H_n(\mathbf{h}, \mathbf{g}) = \begin{bmatrix} F_{n/2}^T(\mathbf{h}_o) & F_{n/2}^T(\mathbf{h}_e) \\ F_{n/2}^T(\mathbf{g}_o) & F_{n/2}^T(\mathbf{g}_e) \end{bmatrix} L_2^{n \times k \times 2}$$

Lifting Scheme



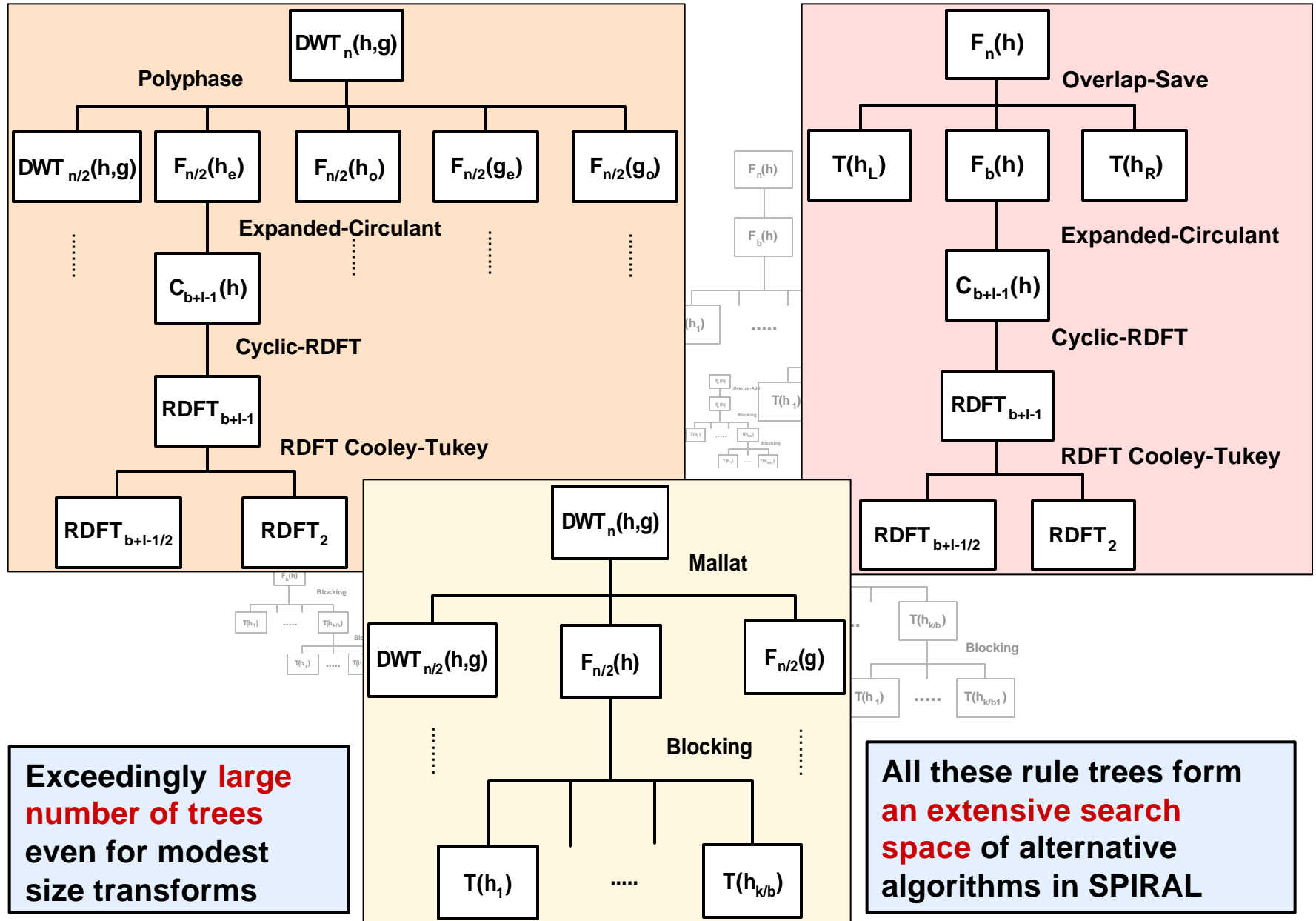
- Asymptotically **reduce cost by 50%**
- **In-place** computations

Lifting Rule

$$H_n(\mathbf{h}, \mathbf{g}) = \begin{bmatrix} K \mathbf{I}_{n/2} & \mathbf{0} \\ \mathbf{0} & 1/K \mathbf{I}_{n/2} \end{bmatrix} \begin{bmatrix} F_{n/2}^T(\mathbf{s}_i) & \mathbf{0} \\ \mathbf{0} & F_{n/2}^T(\mathbf{t}_o) \end{bmatrix} L_2^{n \times k \times 2}$$



Examples: FIR Filter and DWT Rule Trees



Exceedingly large number of trees even for modest size transforms

All these rule trees form an extensive search space of alternative algorithms in SPIRAL



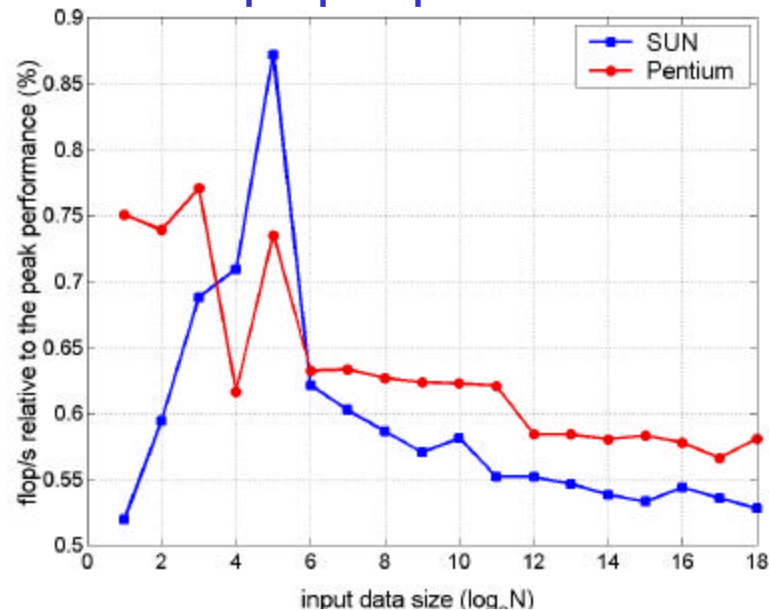
Time-Domain Methods

FIR Filter

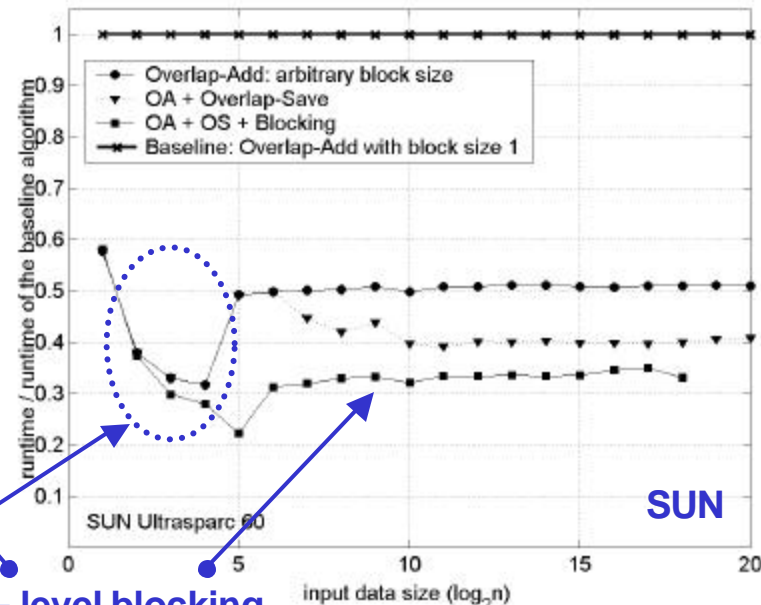
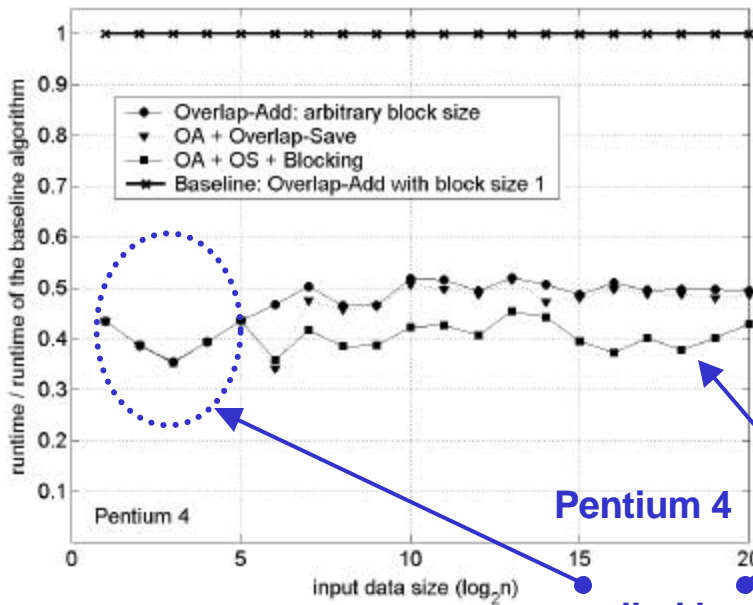
- overlap-add, overlap-save, and blocking
- all methods have the **same cost**
- **baseline** method = overlap-add with $b=1$
- sanity check – **percentage of the peak performance**

60-70% improvement over baseline
Same arithmetic cost

flops / peak performance



comparison of different time-domain methods



unrolled loops

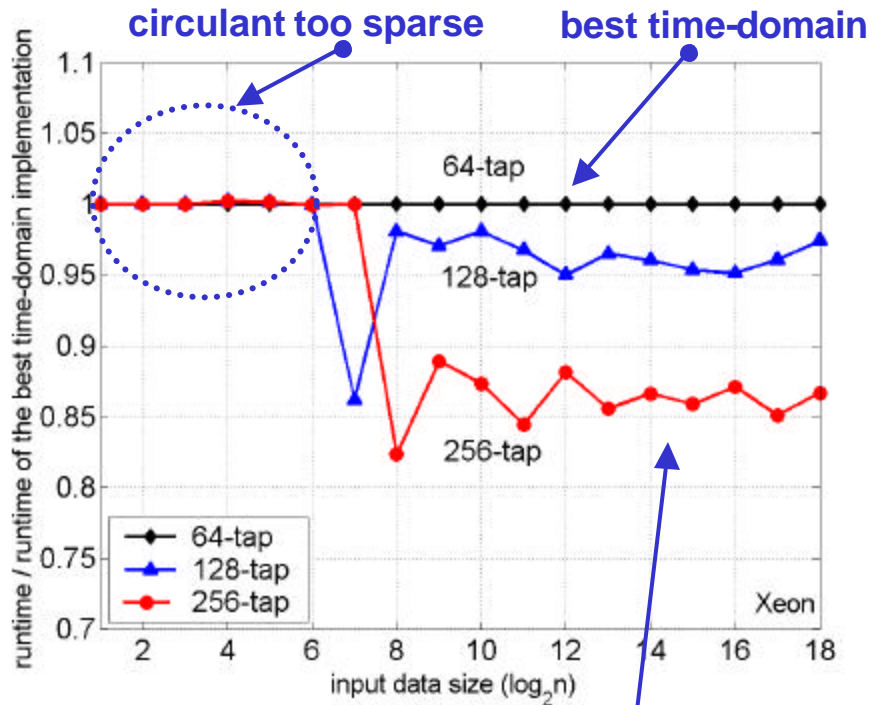
1-level blocking
blocks size 2



Frequency Domain vs. Time Domain

FIR Filter

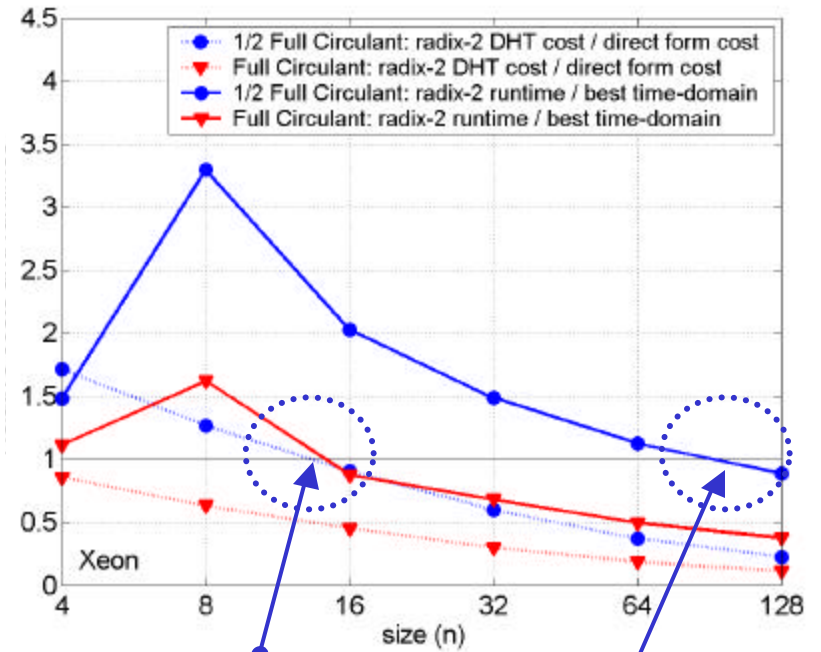
Frequency domain – potentially lower cost
 Time domain – better data access patterns



FD methods better on large “fat” circulants
 Filter lengths > 64

Circulant

“half-dense” circulants emerge in frequency domain methods



arithmetic cost cross-over

runtime cross-over

Considerable discrepancy between arithmetic costs and actual runtimes

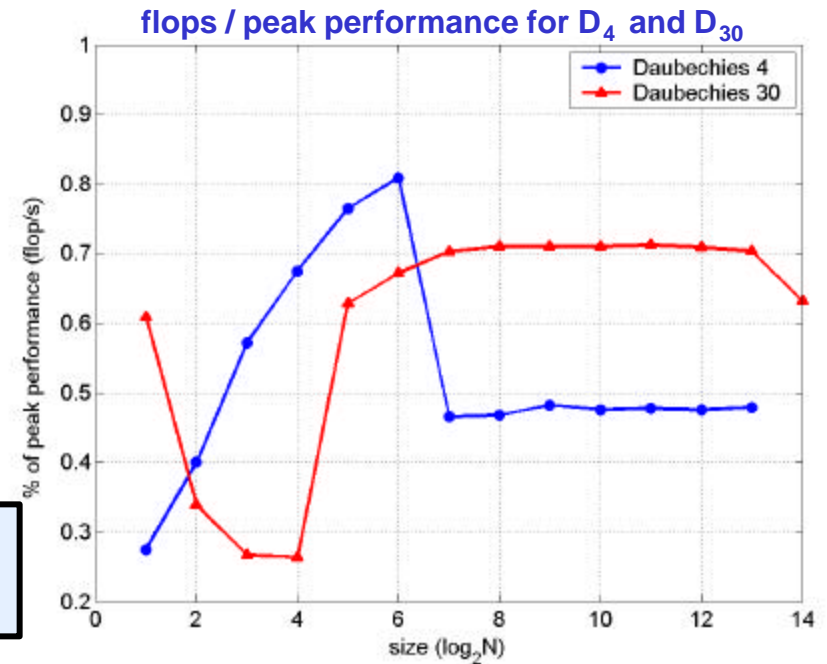


Daubechies D_4 DWT

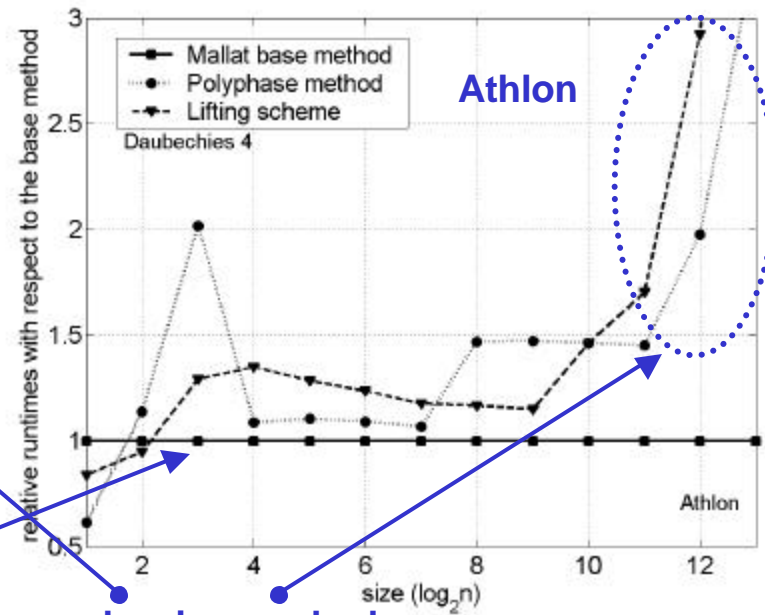
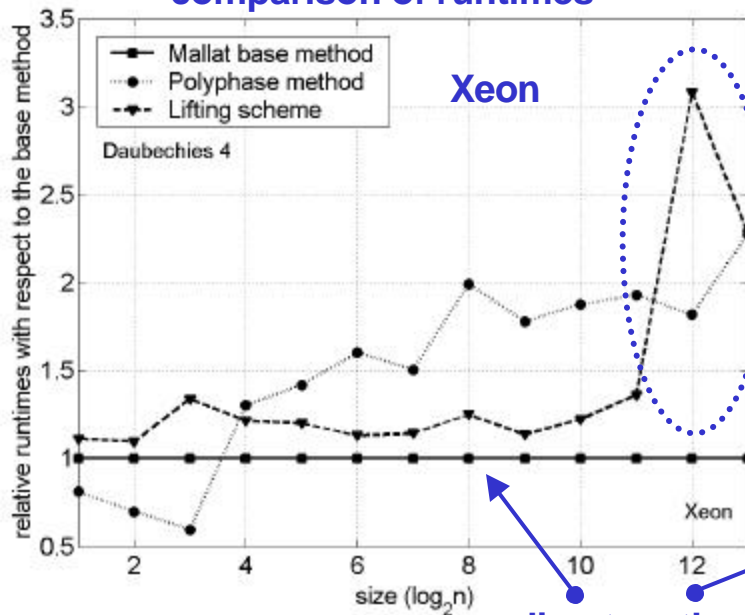
Short Wavelet

- single-level DWT using D_4
- 3 different classes of rule trees based on the initial rule
 - **Fused Mallat Rule** (direct method)
 - **Polyphase Rule** (frequency domain)
 - **Lifting Rule** (lifting scheme)

Frequency domain trees not found
Lifting scheme – 30% lower cost – slower



comparison of runtimes



direct method

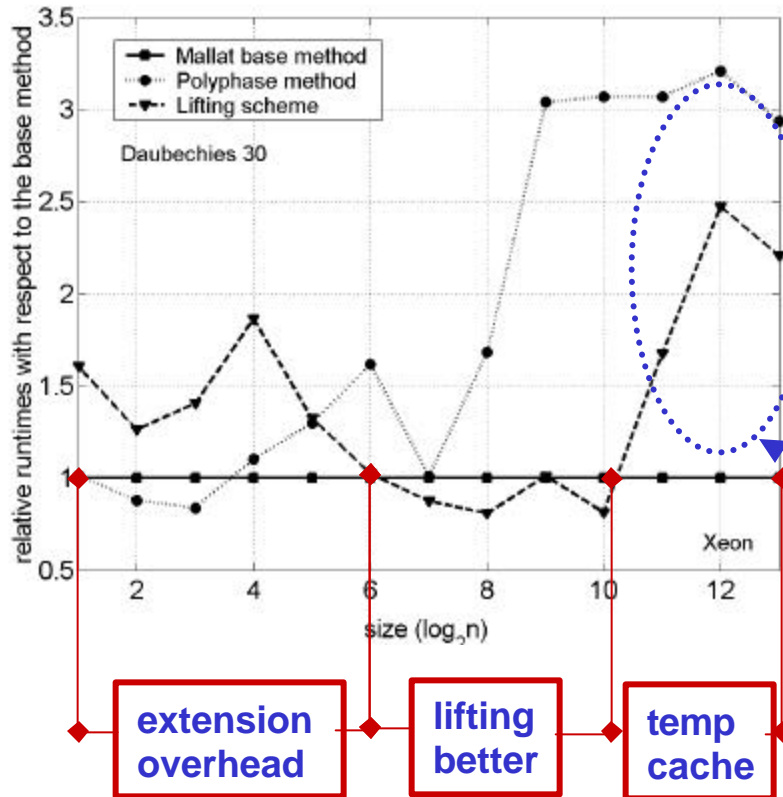
cache size reached



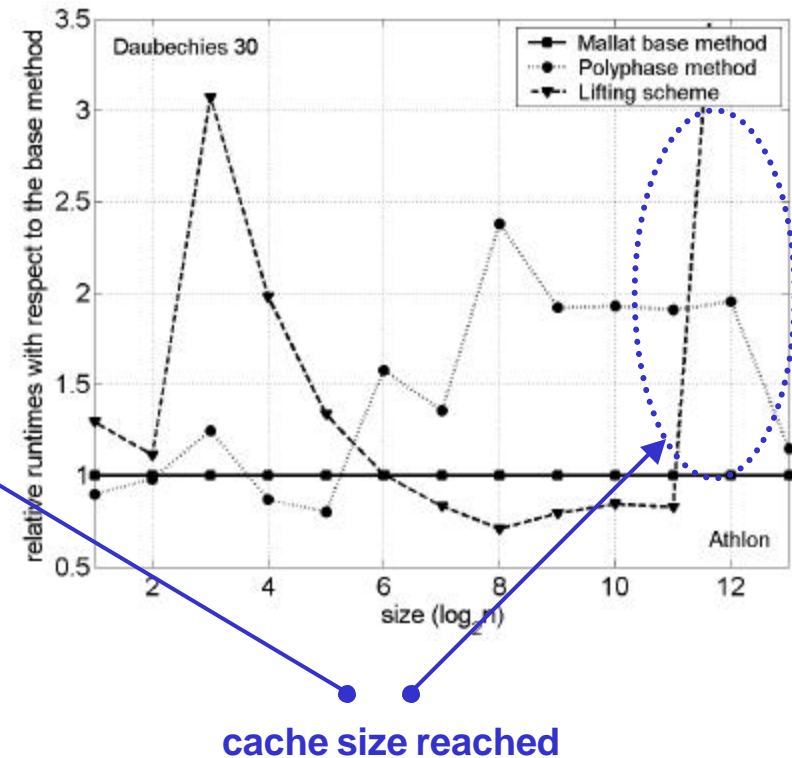
Daubechies D_{30} DWT Experiment

Long Wavelet

Xeon



Athlon



Different methods are found for different size ranges
Best implementations vary from machine to machine

Summary

- ✍ **Fast automatically generated code** for FIR filters and DWT
 - ✍ Significant contribution – **no existing solutions**
- ✍ **Preliminary results**
 - ✍ Arithmetic cost **not a reliable predictor** of performance
 - ✍ Best implementation **not obvious for a human programmer**
 - ✍ Best implementations **vary over different computer platforms**

What lies ahead?

- ✍ **Further algorithmic and implementation optimizations**
 - ✍ Fusion of the extension, exploiting in-place structure, scheduling,..
- ✍ **Entire class** of wavelets and extension methods
 - ✍ Orthogonal, biorthogonal, frames,...
 - ✍ Periodic, symmetric, polynomial extension, zero-padding
- ✍ **Extend to other wavelet transforms**
 - ✍ 2D-DWT, M-band, wavelet packet transform
- ✍ **Utilize special instruction set extensions**
 - ✍ Single Instruction Multiple Data (SIMD)
 - ✍ Fused Multiply-Add (FMA)