**SPIRAL**
www.spiral.net

# Computer Generation of Efficient Software Viterbi Decoders

**Frédéric de Mesmay, Srinivas Chellappa,**
**Franz Franchetti, Markus Püschel**

**Electrical and Computer Engineering**
**Carnegie Mellon University**

**Co-Founder**
**SpiralGen, Inc.**

# Viterbi Decoder

- **Error correction**
  - Forward Error Correction
  - Digital cellular (CDMA, GSM), modems, satellite/deep space communications, 802.11 wireless LANs
  - Software defined radio (SDR)

- **Pattern Recognition**
  - Speech recognition
  - text recognition
  - computational linguistics
  - bioinformatics

**GSM (TCH/FS)
K=5 rate=1/2**

**CDMA2000/UMTS/IS-95
K=9 rate=1/3**

**NASA Cassini
Orbiter:
K=15 rate=1/6**

*SDR requires efficient Viterbi decoder software implementations*

SPIRAL
www.spiral.net

# Software Defined Radio

**WiFi transmitter on Intel Atom Dualcore**

Run time per OFDM symbol [µs] vs. data rate [Mbit/s]



Straightforward C code
but minimizing op count

8 x

Best standard C code

*realtime*

Spiral: computer generated

6.3 x

**Parallelism:
2 threads
4-16 way SIMD**

*Compilers fail to optimize: 50x*

**SPIRAL**
www.spiral.net

# Spiral: Viterbi Software Generation

**Select convolutional code**

Select a preset code or customize parameters

○ *custom*

○ Voyager     rate    1 / [2]     code rate (?)

○ NASA-DSN     K    [7]     constraint length (?)

◉ CCSDS/NASA-GSFC     polynomials [79]

○ WiMax            [91]

○ CDMA IS-95A

○ LTE (3GPP - Long Term Evolution)

○ UWB (802.15)

○ CDMA 2000

○ Cassini

○ Mars Pathfinder & Stereo

polynomials for the code in decimal notation (?)

**Select implementation options**

frame length     [2048]     unpadded frame length in bits (?)

Vectorization level     [SSE 16-way ▾]     type of code (?)

[ Generate Code ]   [ Reset ]

**"Click": Push-button code generation**

http://www.spiral.net/software/viterbi.html

# Spiral: Generated SSE Viterbi Code

```
void viterbi_ccsds(unsigned char *Y, unsigned char *X, unsigned char *syms,
    unsigned char *dec, unsigned char *Branchtab) {
    for(int i9 = 0; i9 <= 1026; i9++) {
        unsigned char a75, a81; int a73, a92;

        ...
        a71 = ((__m128i *) X);              s18 = *(a71); a72 = (a71 + 2);
        s19 = *(a72); a73 = (4 * i9);       a74 = (syms + a73); a75 = *(a74);
        a76 = _mm_set1_epi8(a75);           a77 = ((__m128i *) Branchtab);
        a78 = *(a77);                       a79 = _mm_xor_si128(a76, a78);
        b6 = (a73 + syms);                  a80 = (b6 + 1);
        a81 = *(a80);                       a82 = _mm_set1_epi8(a81);
        a83 = (a77 + 2);                    a84 = *(a83);
        a85 = _mm_xor_si128(a82, a84);  t13 = _mm_avg_epu8(a79,a85);
        a86 = ((__m128i ) t13);             a87 = _mm_srli_epi16(a86, 2);
        a88 = ((__m128i ) a87);
        t14 = _mm_and_si128(a88, _mm_set_epi8(63, 63, 63, 63, 63, 63, 63 , 63, 63, 63,
            63, 63, 63, 63, 63 , 63));
        t15 = _mm_subs_epu8(_mm_set_epi8(63, 63, 63, 63, 63, 63, 63 , 63, 63, 63, 63,
            63, 63, 63, 63 , 63), t14); m23 = _mm_adds_epu8(s18, t14);
        m24 = _mm_adds_epu8(s19, t15);  m25 = _mm_adds_epu8(s18, t15);
        m26 = _mm_adds_epu8(s19, t14);  a89 = _mm_min_epu8(m24, m23);
        ...
    }
    ...
}
```

Generate Code    Reset

"Click": Push-button code generation

`http://www.spiral.net/software/viterbi.html`

# Organization

- **Spiral**

- **Generating software Viterbi decoders**

- **Performance results**

- **Summary**
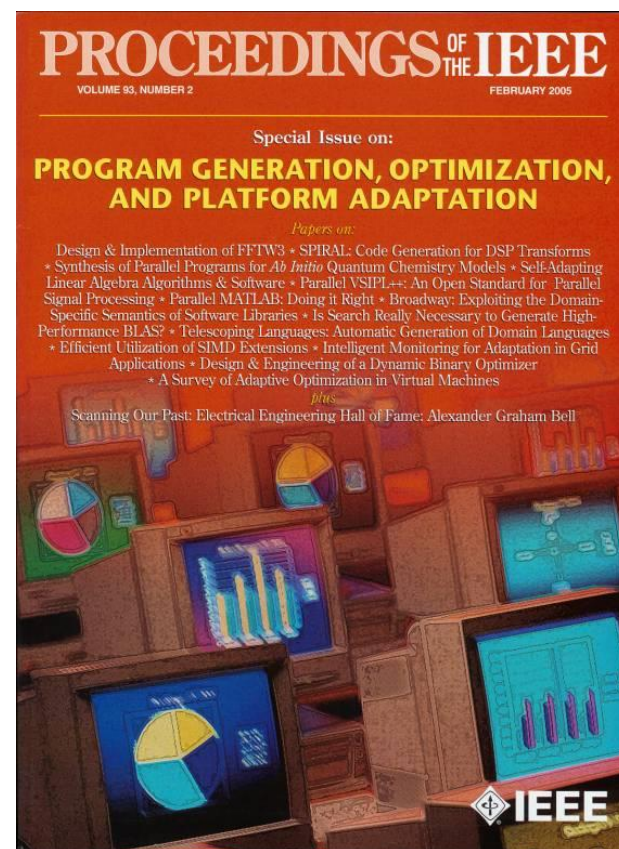
# Organization

- **Spiral**

- **Generating software Viterbi decoders**

- **Performance results**

- **Summary**

# Automatic Performance Tuning

- **Current vicious circle:** **Whenever a new platform comes out, the same functionality needs to be rewritten and reoptimized**
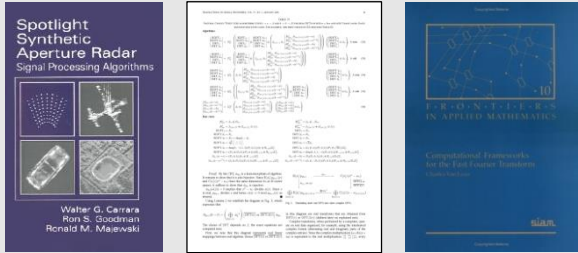
- **Automatic Performance Tuning**
  - BLAS: ATLAS, PHiPAC
  - Linear algebra: Sparsity/OSKI, Flame
  - Sorting
  - Fourier transform: FFTW
  - Linear transforms (and Viterbi): Spiral
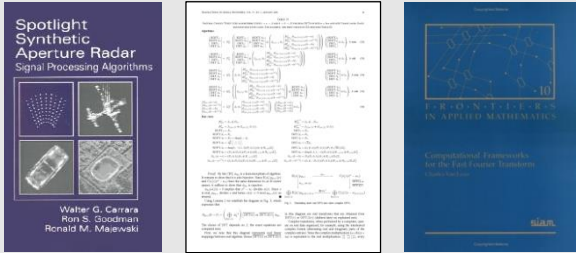  - …others

*New problem class: software Viterbi decoders*



PROCEEDINGS OF THE IEEE
VOLUME 93, NUMBER 2    FEBRUARY 2005

Special Issue on:
PROGRAM GENERATION, OPTIMIZATION, AND PLATFORM ADAPTATION

Proceedings of the IEEE special issue, Feb. 2005

# What is Spiral?

# Idea: Common Abstraction and Rewriting

**Model:** common abstraction
= spaces of matching formulas
= domain-specific language



**abstraction**

$$\underset{\text{vec}(\nu)}{\underline{A_n \otimes I_\nu}}$$

**defines**

$$\underset{\text{smp}(p,\mu)}{\underline{I_p \otimes A_n}}$$

**abstraction**

**rewriting**

**pick**

**architecture space**

**algorithm space**

**search**

$$(\text{DFT}_2 \otimes I_4)\, T_4^8 \,(I_2 \otimes (\ldots\ldots))\, L_2^8$$

Architectural parameter:
Vector length,
#processors, …

**optimization**

Kernel:
problem size,
algorithm choice

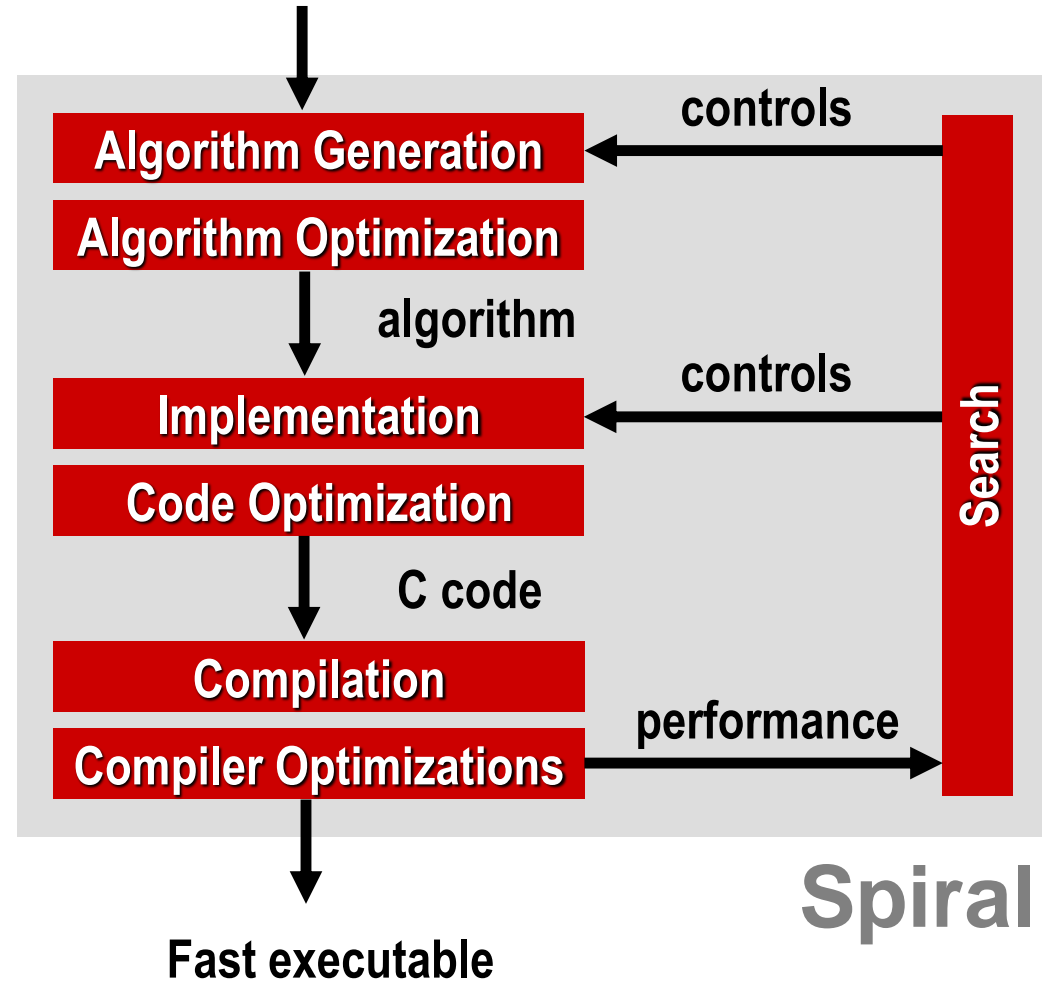# Program Generation in Spiral

**Spiral:**

**Complete automation of the implementation and optimization task**

**Basic ideas:**

**Declarative representation of algorithms**

**Rewriting systems to generate and optimize algorithms at a high level of abstraction**

**Problem specification (transform)**

**Algorithm Generation** ← controls

**Algorithm Optimization**

algorithm

**Implementation** ← controls

**Code Optimization**

C code

**Compilation**

**Compiler Optimizations** → performance

**Search**

**Fast executable**

Spiral

Markus Püschel, José M. F. Moura, Jeremy Johnson, David Padua, Manuela Veloso, Bryan Singer, Jianxin Xiong,
Franz Franchetti, Aca Gacic, Yevgen Voronenko, Kang Chen, Robert W. Johnson, and Nick Rizzolo:
**SPIRAL: Code Generation for DSP Transforms.** Special issue, Proceedings of the IEEE 93(2), 2005

# Some Kernels as Operator Formulas

## Linear Transforms

$$\mathbf{DFT}_n \;\rightarrow\; (\mathbf{DFT}_k \otimes I_m)\, T_m^n (I_k \otimes \mathbf{DFT}_m)\, L_k^n, \quad n = km$$

$$\mathbf{DFT}_n \;\rightarrow\; P_n(\mathbf{DFT}_k \otimes \mathbf{DFT}_m)Q_n, \quad n = km, \ \gcd(k,m)=1$$

$$\mathbf{DFT}_p \;\rightarrow\; R_p^T(I_1 \oplus \mathbf{DFT}_{p-1})D_p(I_1 \oplus \mathbf{DFT}_{p-1})R_p, \quad p \text{ prime}$$

$$\mathbf{DCT\text{-}3}_n \;\rightarrow\; (I_m \oplus J_m)\, L_m^n(\mathbf{DCT\text{-}3}_m(1/4) \oplus \mathbf{DCT\text{-}3}_m(3/4))$$

$$\cdot (F_2 \otimes I_m)\begin{bmatrix} I_m & 0 \oplus -J_{m-1} \\ & \frac{1}{\sqrt{2}}(I_1 \oplus 2\,I_m) \end{bmatrix}, \quad n = 2m$$

$$\mathbf{DCT\text{-}4}_n \;\rightarrow\; S_n\mathbf{DCT\text{-}2}_n \operatorname{diag}_{0 \le k < n}(1/(2\cos((2k+1)\pi/4n)))$$

$$\mathbf{IMDCT}_{2m} \;\rightarrow\; (J_m \oplus I_m \oplus I_m \oplus J_m)\left(\left(\begin{bmatrix} 1 \\ -1 \end{bmatrix} \otimes I_m\right) \oplus \left(\begin{bmatrix} -1 \\ -1 \end{bmatrix} \otimes I_m\right)\right) J_{2m}\,\mathbf{DCT\text{-}4}_{2m}$$

$$\mathbf{WHT}_{2^k} \;\rightarrow\; \prod_{i=1}^{t}(I_{2^{k_1+\cdots+k_{i-1}}} \otimes \mathbf{WHT}_{2^{k_i}} \otimes I_{2^{k_{i+1}+\cdots+k_t}}), \quad k = k_1 + \cdots + k_t$$

$$\mathbf{DFT}_2 \;\rightarrow\; F_2$$

$$\mathbf{DCT\text{-}2}_2 \;\rightarrow\; \operatorname{diag}(1, 1/\sqrt{2})\, F_2$$

$$\mathbf{DCT\text{-}4}_2 \;\rightarrow\; J_2\, R_{13\pi/8}$$

## Viterbi Decoding



$$\underbrace{\mathbf{Vit}}_{\text{vec}(v)} \;\rightarrow\; \underbrace{\left(\prod(L \times I) \circ (I \otimes C)\right) \circ Id}_{\text{vec}(v)}$$

$$\rightarrow \left(\prod \underbrace{(L \times I) \circ (I \otimes C)}_{\text{vec}(v)}\right) \circ Id$$

$$£ \quad \rightarrow \left(\prod(L \otimes I_v \times I) \circ (I \otimes C \otimes I_v) \circ (\vec{L} \times I)\right) \circ Id$$

$$\rightarrow \prod(L \otimes I_v \times I) \circ (I \otimes (B \otimes I_v)) \circ (\vec{L} \times I)$$

## Matrix-Matrix Multiplication



$$\mathbf{MMM}_{1,1,1} \rightarrow (\cdot)_1$$

$$\mathbf{MMM}_{m,n,k} \rightarrow (\otimes)_{m/m_b \times 1} \otimes \mathbf{MMM}_{m_b,n,k}$$

$$\mathbf{MMM}_{m,n,k} \rightarrow \mathbf{MMM}_{m,nb,k} \otimes (\otimes)_{1 \times n/nb}$$

$$\mathbf{MMM}_{m,n,k} \rightarrow ((\Sigma_{k/k_b} \circ (\cdot)_{k/k_b}) \otimes \mathbf{MMM}_{m,n,k_b}) \circ$$
$$((L_{k/k_b}^{mk/k_b} \otimes I_{k_b}) \times I_{kn})$$

$$\mathbf{MMM}_{m,n,k} \rightarrow (L_m^{mn/n_b} \otimes I_{n_b}) \circ$$
$$((\otimes)_{1 \times n/n_b} \otimes \mathbf{MMM}_{m,n_b,k}) \circ$$
$$(I_{km} \times (L_{n/n_b}^{kn/n_b} \otimes I_{n_b}))$$

## Synthetic Aperture Radar (SAR)



$$\mathbf{SAR}_{k \times m \rightarrow n \times n} \;\rightarrow\; \mathbf{DFT}_{n \times n} \circ \mathbf{Interp}_{k \times m \rightarrow n \times n}$$

$$\mathbf{DFT}_{n \times n} \;\rightarrow\; (\mathbf{DFT}_n \otimes I_n) \circ (I_n \otimes \mathbf{DFT}_n)$$

$$\mathbf{Interp}_{k \times m \rightarrow n \times n} \;\rightarrow\; (\mathbf{Interp}_{k \rightarrow n} \otimes_i I_n) \circ (I_k \otimes_i \mathbf{Interp}_{m \rightarrow n})$$

$$\mathbf{Interp}_{r \rightarrow s} \;\rightarrow\; \left(\bigoplus_{i=0}^{n-2} \mathbf{InterpSeg}_k\right) \oplus \mathbf{InterpSegPruned}_{k,\ell}$$

$$\mathbf{InterpSeg}_k \;\rightarrow\; G_f^{u \cdot n \rightarrow k} \circ \mathbf{iPrunedDFT}_{n \rightarrow u \cdot n} \circ \left(\frac{1}{n}\right) \circ \mathbf{DFT}_n$$

# Same Approach for Different Paradigms



## Threading:

$$\underbrace{\mathbf{DFT}_{mn}}_{\mathsf{smp}(p,\mu)} \to \underbrace{\left((\mathbf{DFT}_m \otimes \mathbf{I}_n)\mathsf{T}_n^{mn}(\mathbf{I}_m \otimes \mathbf{DFT}_n)\mathsf{L}_m^{mn}\right)}_{\mathsf{smp}(p,\mu)}$$

$$\cdots$$

$$\to \underbrace{\left(\mathbf{DFT}_m \otimes \mathbf{I}_n\right)}_{\mathsf{smp}(p,\mu)} \underbrace{\mathsf{T}_n^{mn}}_{\mathsf{smp}(p,\mu)} \underbrace{\left(\mathbf{I}_m \otimes \mathbf{DFT}_n\right)}_{\mathsf{smp}(p,\mu)} \underbrace{\mathsf{L}_m^{nm}}_{\mathsf{smp}(p,\mu)}$$

$$\cdots$$

$$\to \left((\mathsf{L}_m^{mp} \otimes \mathbf{I}_{n/p\mu}) \otimes_\mu \mathbf{I}_\mu\right)\left(\mathbf{I}_p \otimes_\parallel (\mathbf{DFT}_m \otimes \mathbf{I}_{n/p})\right)\left((\mathsf{L}_p^{mp} \otimes \mathbf{I}_{n/p\mu}) \otimes_\mu \mathbf{I}_\mu\right)$$

$$\left(\bigoplus_{i=0}^{p-1} \mathsf{T}_n^{mn,i}\right)\left(\mathbf{I}_p \otimes_\parallel (\mathbf{I}_{m/p} \otimes \mathbf{DFT}_n)\right)\left(\mathbf{I}_p \otimes_\parallel \mathsf{L}_{m/p}^{mn/p}\right)\left((\mathsf{L}_p^{pm} \otimes \mathbf{I}_{m/p\mu}) \otimes_\mu \mathbf{I}_\mu\right)$$

## Vectorization:

$$\underbrace{\left(\overline{\mathbf{DFT}_{mn}}\right)}_{\mathsf{vec}(\nu)} \to \underbrace{\left((\mathbf{DFT}_m \otimes \mathbf{I}_n)\mathsf{T}_n^{mn}(\mathbf{I}_m \otimes \mathbf{DFT}_n)\mathsf{L}_m^{mn}\right)}_{\mathsf{vec}(\nu)}$$

$$\cdots$$

$$\to \underbrace{\overline{(\mathbf{DFT}_m \otimes \mathbf{I}_n)}^\nu}_{\mathsf{vec}(\nu)} \underbrace{(\overline{\mathsf{T}_n^{mn}})^\nu}_{\mathsf{vec}(\nu)} \underbrace{\overline{(\mathbf{I}_m \otimes \mathbf{DFT}_n)\mathsf{L}_m^{mn}}^\nu}_{\mathsf{vec}(\nu)}$$

$$\cdots$$

$$\to (\mathbf{I}_{mn/\nu} \otimes \underbrace{\mathsf{L}_\nu^{2\nu}}_{\mathsf{sse}})(\overline{\mathbf{DFT}_m \otimes \mathbf{I}_{n/\nu}}\vec\otimes \mathbf{I}_\nu)\underbrace{(\overline{\mathsf{T}_n^{mn}})^\nu}_{\mathsf{sse}}$$

$$\left(\mathbf{I}_{m/\nu} \otimes (\overline{\mathsf{L}_\nu^n}\vec\otimes \mathbf{I}_\nu)(\mathbf{I}_{n/\nu} \otimes (\mathsf{L}_\nu^{2\nu}\vec\otimes \mathbf{I}_\nu)(\mathbf{I}_2 \otimes \underbrace{\mathsf{L}_\nu^{\nu^2}}_{\mathsf{sse}})(\mathsf{L}_2^{2\nu}\vec\otimes \mathbf{I}_\nu))(\overline{\mathbf{DFT}_n}\vec\otimes \mathbf{I}_\nu)\right)$$

$$\left((\mathsf{L}_m^{mn} \otimes \mathbf{I}_2)\vec\otimes \mathbf{I}_\nu\right)(\mathbf{I}_{mn/\nu} \otimes \underbrace{\mathsf{L}_2^{2\nu}}_{\mathsf{sse}})$$

## GPUs:

$$\underbrace{\left(\mathbf{DFT}_{r^k}\right)}_{\mathsf{gpu}(t,c)} \to \underbrace{\left(\prod_{i=0}^{k-1} \mathsf{L}_r^{r^k}\left(\mathbf{I}_{r^{k-1}} \otimes \mathbf{DFT}_r\right)\left(\mathsf{L}_{r^{k-i-1}}^{r^k}(\mathbf{I}_{r^i} \otimes \mathsf{T}_{r^{k-i-1}}^{r^{k-i}}\underbrace{\mathsf{L}_{r^{i+1}}^{r^k}}_{\mathsf{vec}(c)}\right)\mathsf{R}_r^{r^k}\right)}_{\mathsf{gpu}(t,c)}$$

$$\cdots$$

$$\to \left(\prod_{i=0}^{k-1}(\mathsf{L}_r^{r^n/2}\vec\otimes \mathbf{I}_2)\left(\mathbf{I}_{r^{n-1}/2} \otimes_\times \underbrace{(\mathbf{DFT}_r\vec\otimes \mathbf{I}_2)\mathsf{L}_r^{2r}}_{\mathsf{shd}(t,c)}\right)\mathsf{T}_i\right)$$

$$(\mathsf{L}_r^{r^n/2}\vec\otimes \mathbf{I}_2)(\mathbf{I}_{r^{n-1}/2} \otimes_\times \underbrace{\mathsf{L}_r^{2r}}_{\mathsf{shd}(t,c)})(\mathsf{R}_r^{r^{n-1}}\vec\otimes \mathbf{I}_r)$$

## Verilog for FPGAs:

$$\underbrace{\left(\mathbf{DFT}_{r^k}\right)}_{\mathsf{stream}(r^s)} \to \underbrace{\left[\prod_{i=0}^{k-1} \mathsf{L}_r^{r^k}\left(\mathbf{I}_{r^{k-1}} \otimes \mathbf{DFT}_r\right)\left(\mathsf{L}_{r^{k-i-1}}^{r^k}(\mathbf{I}_{r^i} \otimes \mathsf{T}_{r^{k-i-1}}^{r^{k-i}})\mathsf{L}_{r^{i+1}}^{r^k}\right)\right]\mathsf{R}_r^{r^k}}_{\mathsf{stream}(r^s)}$$

$$\cdots$$

$$\to \left[\prod_{i=0}^{k-1} \underbrace{\mathsf{L}_r^{r^k}}_{\mathsf{stream}(r^s)} \underbrace{\left(\mathbf{I}_{r^{k-1}} \otimes \mathbf{DFT}_r\right)}_{\mathsf{stream}(r^s)} \underbrace{\left(\mathsf{L}_{r^{k-i-1}}^{r^k}(\mathbf{I}_{r^i} \otimes \mathsf{T}_{r^{k-i-1}}^{r^{k-i}})\mathsf{L}_{r^{i+1}}^{r^k}\right)}_{\mathsf{stream}(r^s)}\right] \underbrace{\mathsf{R}_r^{r^k}}_{\mathsf{stream}(r^s)}$$

$$\cdots$$

$$\to \left[\prod_{i=0}^{k-1} \underbrace{\mathsf{L}_r^{r^k}}_{\mathsf{stream}(r^s)}\left(\mathbf{I}_{r^{k-s-1}} \otimes_s (\mathbf{I}_{r^{s-1}} \otimes \mathbf{DFT}_r)\right)\underbrace{\mathsf{T}_i'}_{\mathsf{stream}(r^s)}\right]\underbrace{\mathsf{R}_r^{r^k}}_{\mathsf{stream}(r^s)}$$
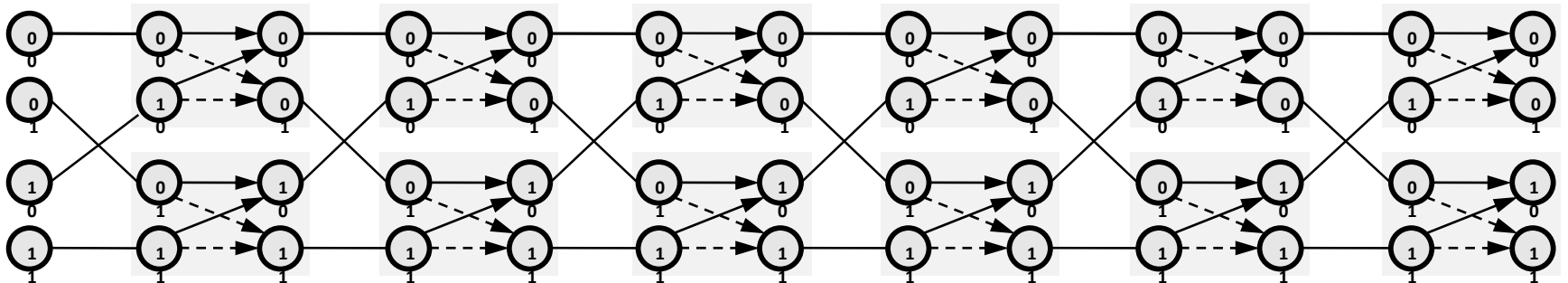
# Organization

- **Spiral**

- **Generating software Viterbi decoders**

- **Performance results**

- **Summary**

# Structure of Viterbi Decoders

## State machine



## Viterbi trellis (data flow)



*Key observation: similarity to Walsh-Hadamard transform (WHT)*

# Viterbi Language (VL)

**VL in Backus-Naur Form (BNF)**

$$
\begin{aligned}
\langle\text{op}\rangle ::= \quad & F_{K,F} & & \textit{Viterbi forward pass} \\
& |\ I_n & & \textit{identity} \\
& |\ L_n^{mn} & & \textit{stride permutation} \\
& |\ B_{i,j} & & \textit{Viterbi butterfly} \\
& |\ \langle\text{op}\rangle\ \langle\text{op}\rangle & & \textit{composition} \\
& |\ \textstyle\prod \langle\text{op}\rangle & & \textit{iterative composition} \\
& |\ \langle\text{op}\rangle \otimes \langle\text{op}\rangle & & \textit{tensor product}
\end{aligned}
$$

**Viterbi decoder forward pass in VL**

$$
F_{K,F} \rightarrow \prod_{i=1}^{F} \left( (I_{2^{K-2}} \otimes_j B_{F-i,j}) L_{2^{K-2}}^{2^{K-1}} \right)
$$

$$
B_{i,j} : \begin{cases} \pi_U = \min_{d_U}(\pi_A + \beta_{A \rightarrow U}, \pi_B + \beta_{B \rightarrow U}) \\ \pi_V = \min_{d_V}(\pi_A + \beta_{A \rightarrow V}, \pi_B + \beta_{B \rightarrow V}) \end{cases}
$$

# Compiling VL To Code

| construct | code |
|-----------|------|
| $y = (CD)x$ | `t = D(x);`<br>`y = C(t);` |
| $y = \prod_{i=0}^{l-1} C^i x$ | `y = C(l-1, x);`<br>`for (i=l-2;i>=0;i--)`<br>`   y = C(i, y);` |
| $y = (I_m \otimes_j C_n^j)x$ | `for (j=0;j<m;j++)`<br>`    y[j*n:j*n+n-1] =`<br>`        C(j, x[j*n:j*n+n-1]);` |
| $y = L_m^{mn}x$ | `for (i=0;i<m;i++)`<br>`    for (j=0;j<n;j++)`<br>`        y[i+m*j]=x[n*i+j];` |
| $y = B_{i,j}x$ | *see equation last slide* |

# Vectorization Through Rewriting

## Vectorization Rule Set

$$
\begin{aligned}
\underline{CD}\ \nu &\rightarrow\ \underline{C}\ \nu\ \underline{D}\ \nu \\
\underline{\prod C}\ \nu &\rightarrow\ \prod \underline{C}\ \nu \\
\underline{I_m \otimes_j C^j}\ \nu &\rightarrow\ I_m \otimes_j \underline{C^j}\ \nu \\
\underline{C \otimes I_\nu}\ \nu &\rightarrow\ C \bar{\otimes}\, I_\nu \\
\underline{B_{i,l\nu+j} \otimes_j I_\nu}\ \nu &\rightarrow\ \vec{B}_{i,l}^{\,\nu} \\
\underline{\mathsf{L}_\nu^{2\nu}}\ \nu &\rightarrow\ \vec{\mathsf{L}}_\nu^{2\nu}
\end{aligned}
$$

## Vectorized Viterbi Decoder

$$
\underline{\mathbf{F}_{K,F}}\ \nu \rightarrow \prod_{i=1}^{F} \left( \left( I_{2^{K-2}/\nu} \otimes_{j_1} \vec{\mathsf{L}}_\nu^{2\nu} \vec{B}_{F-i,j_1}^{\,\nu} \right) \left( \mathsf{L}_{2^{K-2}/\nu}^{2^{K-1}/\nu} \bar{\otimes}\, I_\nu \right) \right)
$$

# VL Compilation System

**VL Expression**

**VL Compiler**

*s*calar decoder

**Execution**

*metric spread overflow factors*

**Target Architecture**

**Vectorization by Rewriting**

**VL Compiler**

**Peephole Optimization**

**Vectorized Decoder**
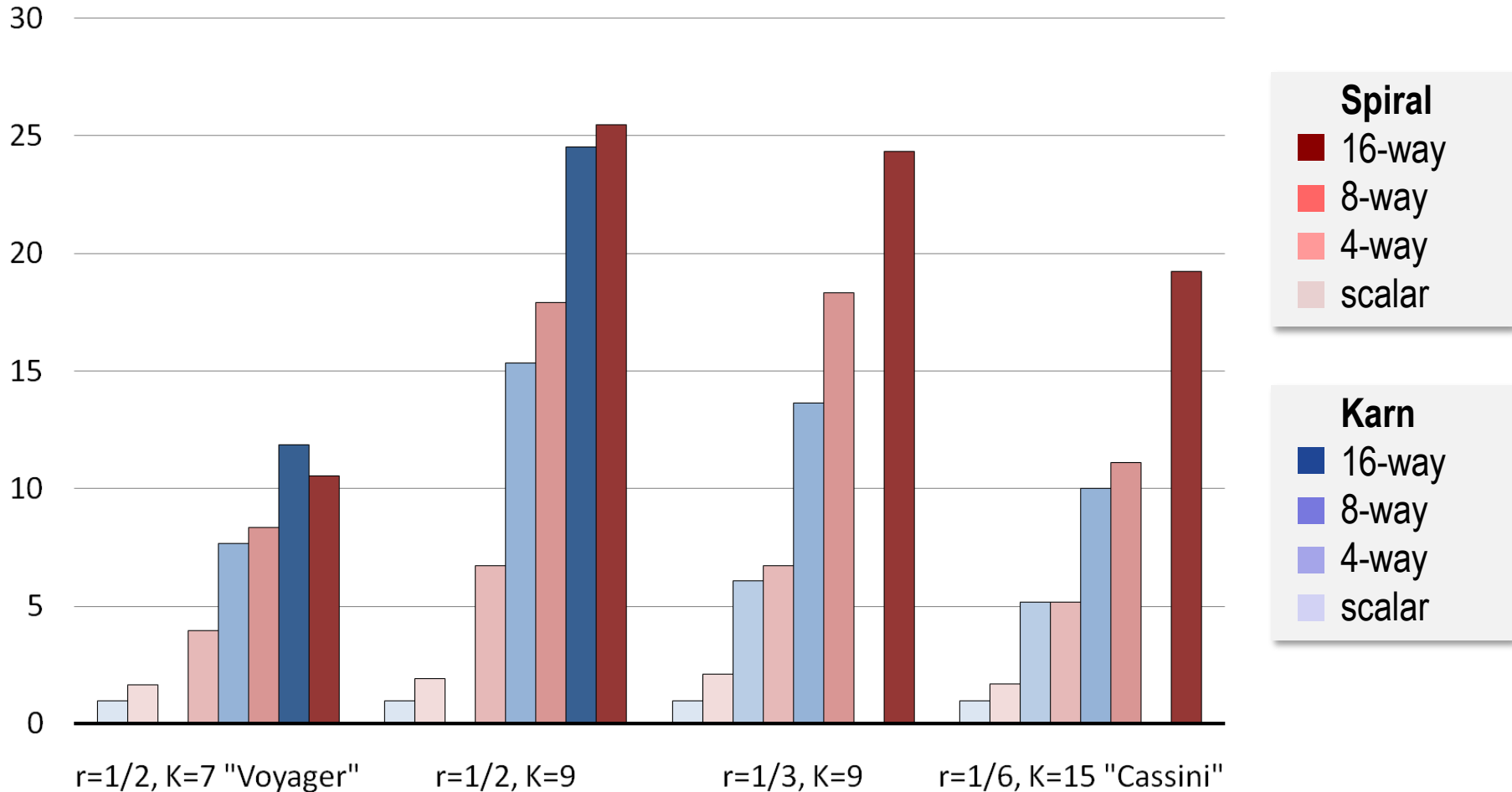
# Organization

- ■ Spiral

- ■ Generating Software Viterbi Decoders

- ■ **Performance results**

- ■ Summary

# Comparison to Hand-Tuned Code

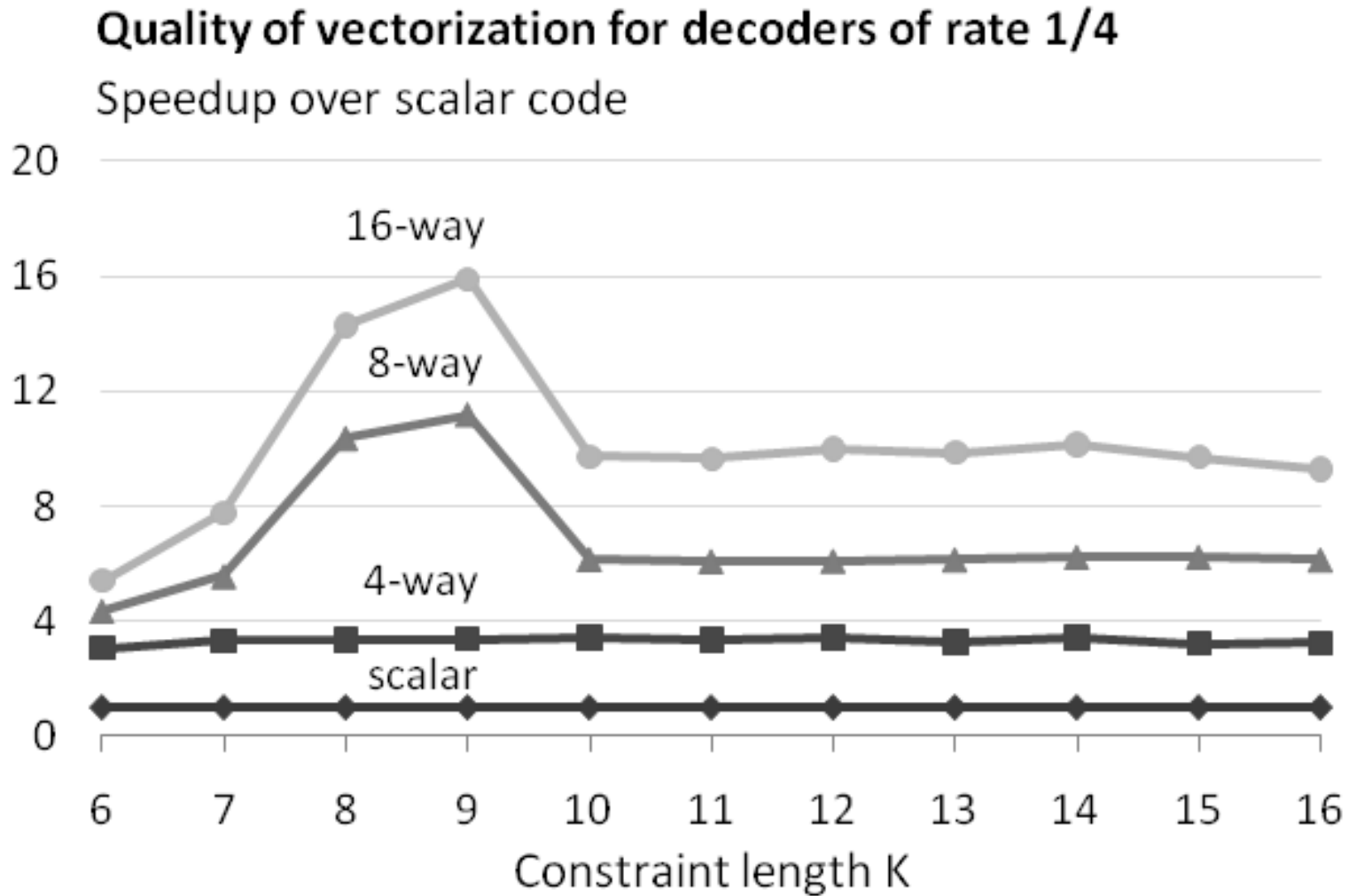**Karn's implementation:** hand-written assembly for 4 specific Viterbi codes

**Performance Gain of Various Generated Viterbi Decoders**

Speedup over Karn's C implementation



**Single core of Core2 Extreme (quad-core), 3 GHz, Intel C++ compiler 10.0**
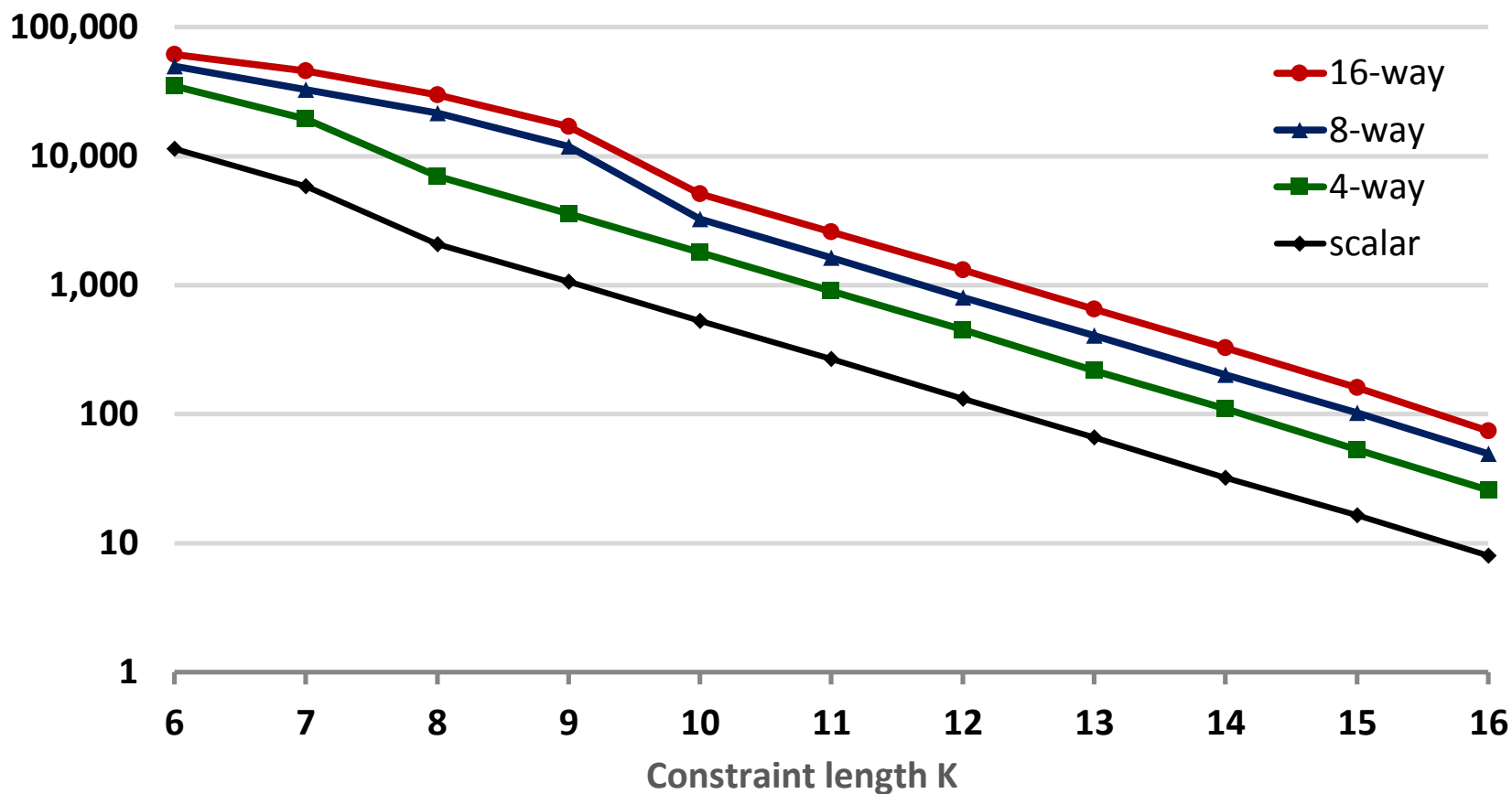
# Vectorization Speed-Up



Quality of vectorization for decoders of rate 1/4
Speedup over scalar code

**Single core of Core2 Extreme (quad-core), 3 GHz, Intel C++ compiler 10.0**

# Data Rate Results

## Decoders for rate 1/4

**Performance (kbit/s)**



**Single core of Core2 Extreme (quad-core), 3 GHz, Intel C++ compiler 10.0**

# Organization

- **Spiral**

- **Generating Software Viterbi Decoders**

- **Performance results**

- **Summary**

# Summary

- **Platforms are powerful yet complicated**
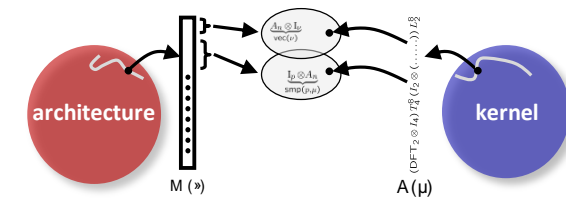  optimization will stay a hard problem
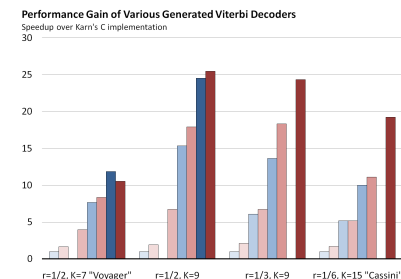

Image: Intel

- **Automatic generation of Viterbi decoder**
  from high-level specification

$$F_{K,F} \rightarrow \prod_{i=1}^{F} \left( (I_{2^{K-2}} \otimes_j B_{F-i,j}) L_{2^{K-2}}^{2^{K-1}} \right)$$
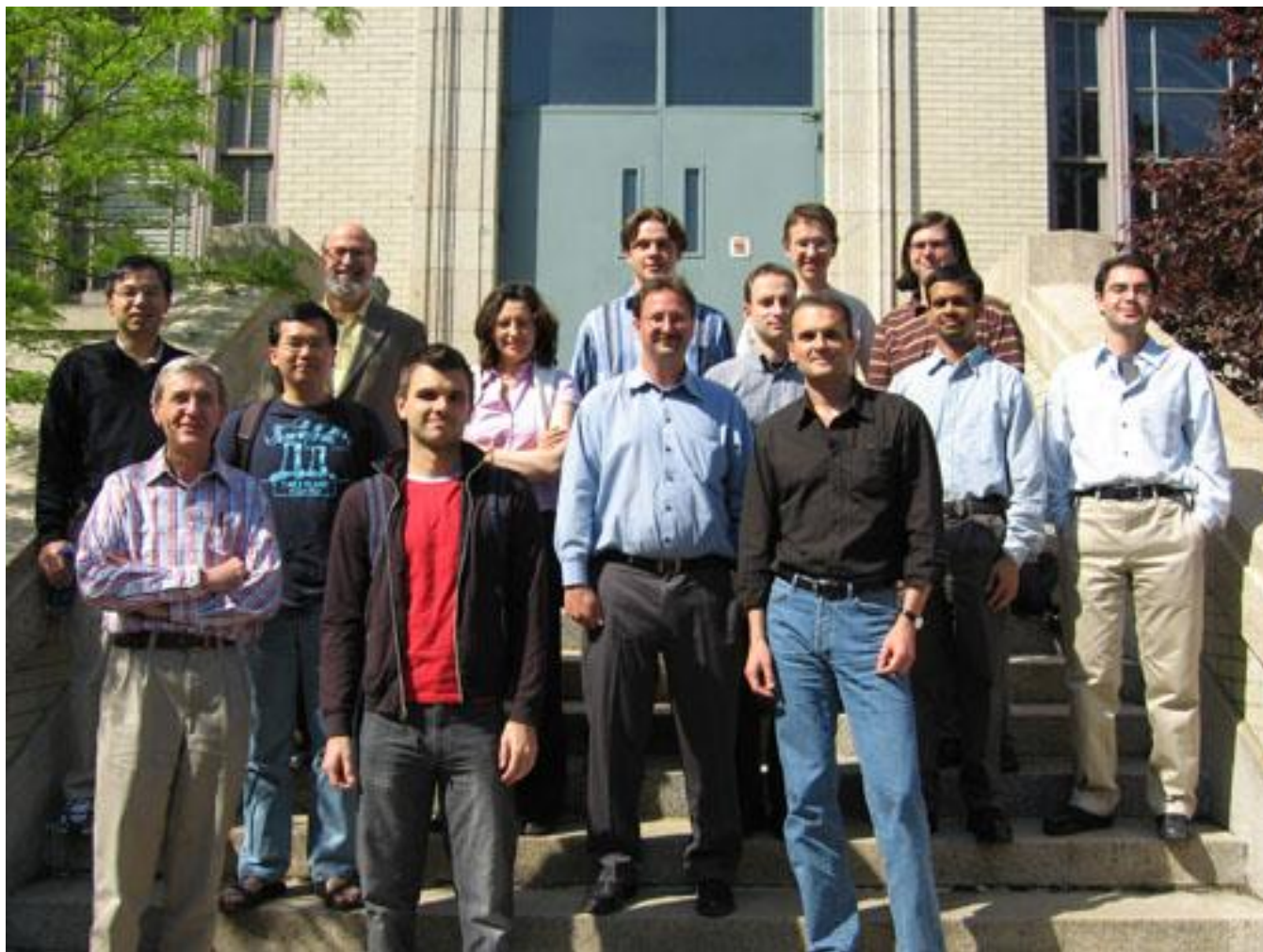
- **Spiral: program generation and autotuning**
  can provide full automation



- **Performance of Spiral's Viterbi decoders**
  is competitive with expert hand tuning

# (Part of the) Spiral Team



**www.spiral.net**