# intrinsity ™

# An Innovative High-Performance Architecture for Vector and Matrix Math Algorithms

## *Presented by: Tim Olson, Architect*

## **HPEC 2002 – September 24, 2002**

*Authors: Veeraraghavan Anantha, Ph.D.;*
*Christophe Harlé, Ph.D.; Tim Olson, George Yost, Ph.D.*

# Intrinsity FastMATH™
# Vector and Matrix Math Processor

**intrinsity™**

## Optimized for real-time and adaptive signal processing needs:
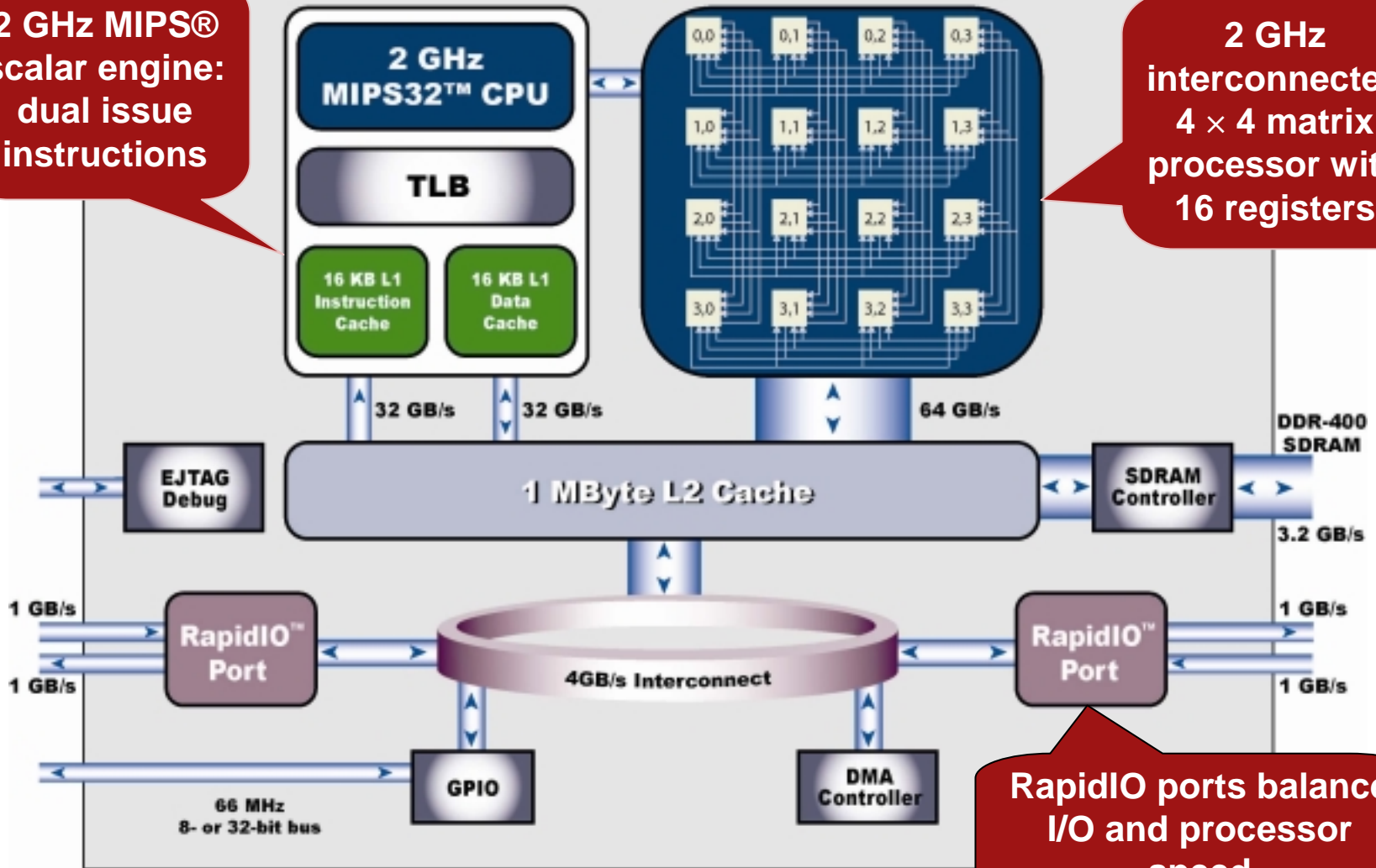
**Innovative architecture:**

- *2 GHz SIMD 4 × 4 matrix engine with multiprocessor scalability due to high bandwidth RapidIO™ interfaces*

- **Fixed-point math**

- **High-level (e.g., C) language programmable**
  - **Compiler built-in matrix intrinsics**
  - **Vector/matrix library**

- **On-chip matrix coprocessor and MIPS32™ ISA RISC core**

- **4 × 4 array of processors, each with sixteen 32-bit registers, two 40-bit MACs**

- **64 GOPS (peak)**

- **Matrix and vector math native instructions: 1-, 8-, 16-, 32-bit support; convenient complex math**

- **Descriptor-based DMA controller**

- **1 Mbyte on-chip cache-coherent L2 cache**

## Speed *plus* an architecture designed for parallel computations

# Intrinsity FastMATH Vector and Matrix Math Processor

**The matrix engine has 16 matrix registers, each with 16 32-bit values. Halfword and word arithmetic is supported.**

**Single instruction, element-wise addition of two $4 \times 4$ matrices**

**Matrix Registers**

| 0,0 | 0,1 | 0,2 | 0,3 |
| 1,0 | 1,1 | 1,2 | 1,3 |
| 2,0 | 2,1 | 2,2 | 2,3 |
| 3,0 | 3,1 | 3,2 | 3,3 |

**$M_2$**

=

| 0,0 | 0,1 | 0,2 | 0,3 |
| 1,0 | 1,1 | 1,2 | 1,3 |
| 2,0 | 2,1 | 2,2 | 2,3 |
| 3,0 | 3,1 | 3,2 | 3,3 |

**$M_1$**

+

| 0,0 | 0,1 | 0,2 | 0,3 |
| 1,0 | 1,1 | 1,2 | 1,3 |
| 2,0 | 2,1 | 2,2 | 2,3 |
| 3,0 | 3,1 | 3,2 | 3,3 |

**$M_0$**

| Re | Im |

or

| 32-bits |

**Complex data by halfwords**

**Word data**

**Matrix-multiply of two $4 \times 4$ submatrices by halfword, for example to support 16-bit complex arithmetic**

**One instruction**

- **Four cycles (2 ns @ 2 GHz)**
- **128 operations**

$$\sum_{k=0}^{3} M0^h(0,k) \times M1^h(k,0)$$

**matmulhh.m.m  M2,M0,M1**

```
for i = 0 to 3
    for j = 0 to 3
        sum = 0
        for k = 0 to 3
            sum = sum + M0^h(i,k) × M1^h(k,j);
    M2^h(i,j) = sum;
```

**High-high halfword multiply, e.g., re × re**

**Matrix Registers**

| 0,0 | 0,1 | 0,2 | 0,3 |
|-----|-----|-----|-----|
| 1,0 | 1,1 | 1,2 | 1,3 |
| 2,0 | 2,1 | 2,2 | 2,3 |
| 3,0 | 3,1 | 3,2 | 3,3 |

**$M_2$**

=

| 0,0 | 0,1 | 0,2 | 0,3 |
|-----|-----|-----|-----|
| 1,0 | 1,1 | 1,2 | 1,3 |
| 2,0 | 2,1 | 2,2 | 2,3 |
| 3,0 | 3,1 | 3,2 | 3,3 |

**$M_0$**

×

| 0,0 | 0,1 | 0,2 | 0,3 |
|-----|-----|-----|-----|
| 1,0 | 1,1 | 1,2 | 1,3 |
| 2,0 | 2,1 | 2,2 | 2,3 |
| 3,0 | 3,1 | 3,2 | 3,3 |

**$M_1$**

**Can subdivide large matrices into $4 \times 4$ parts for multiplication**

**cache**

User 1
User 2
User 3
User 4

16 elements of 1 user

m0  m1  m2  m3

4 elements of 4 users

m0  m1  m2  m3

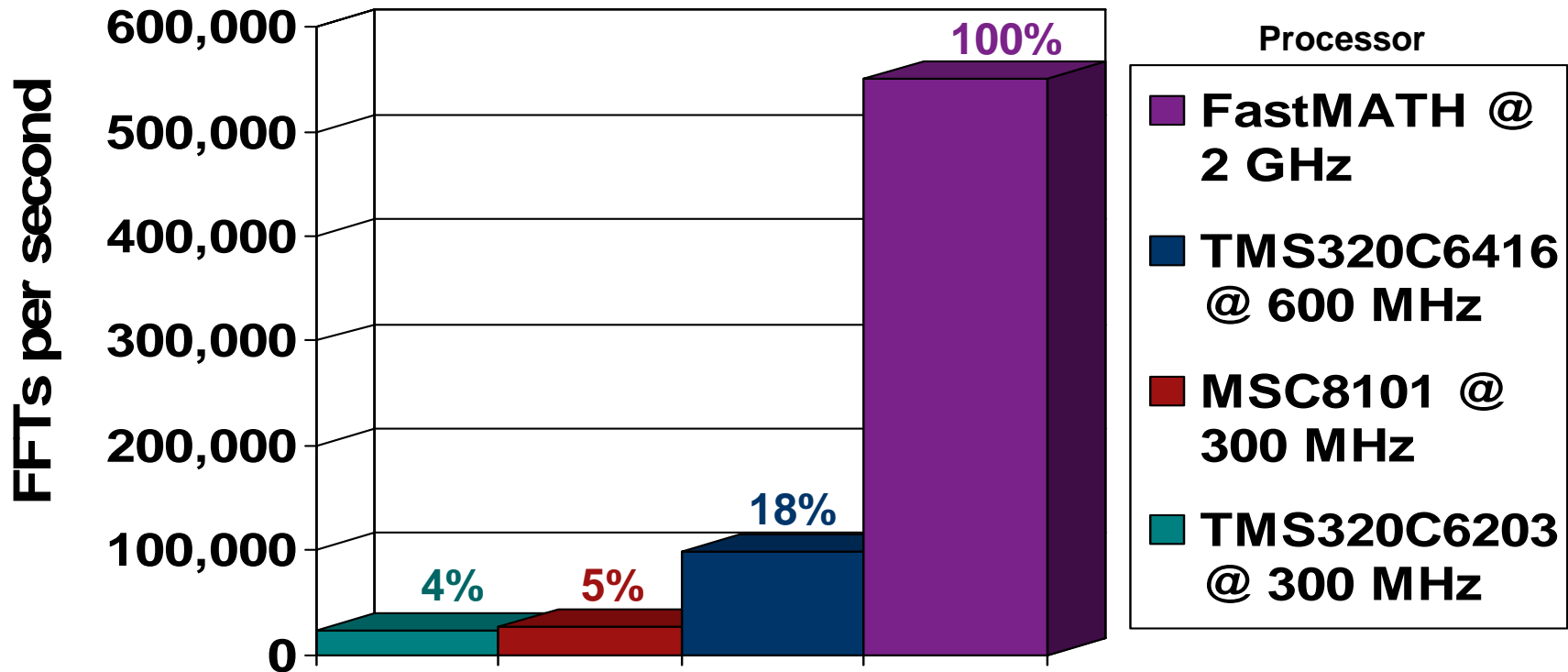1 element of 16 users

m0  m1  m2  m3

**Load 4 or 16 data streams (users) and re-block for SIMD parallel processing**

- **Original register `load` instructions**

- **`block4` (four cycles): matrix operations on four streams**

- **For SIMD operations on *16 parallel data streams*: continue rearrangement with block data movement instructions—70 cycles (35 ns) total**

# FastMATH Performance Example: FFT to Implement OFDM

**Multiple antennas**

**Front-end processing (e.g., FIR filter)**

***N* samples of 16-bit complex data**

**N-point FFT**   **N-point FFT**   **N-point FFT**

**Orthogonal Frequency-Division Multiplexing**

***N* frequency coefficients for each antenna**

**Smart antenna beamforming or symbol-rate processing**

**Example results:**
**for 8 antennas, 10 Msamples per second, 1024-pt complex FFT: requires *14.4% FastMATH processor***

# FastMATH Performance Example: Smart Antennas

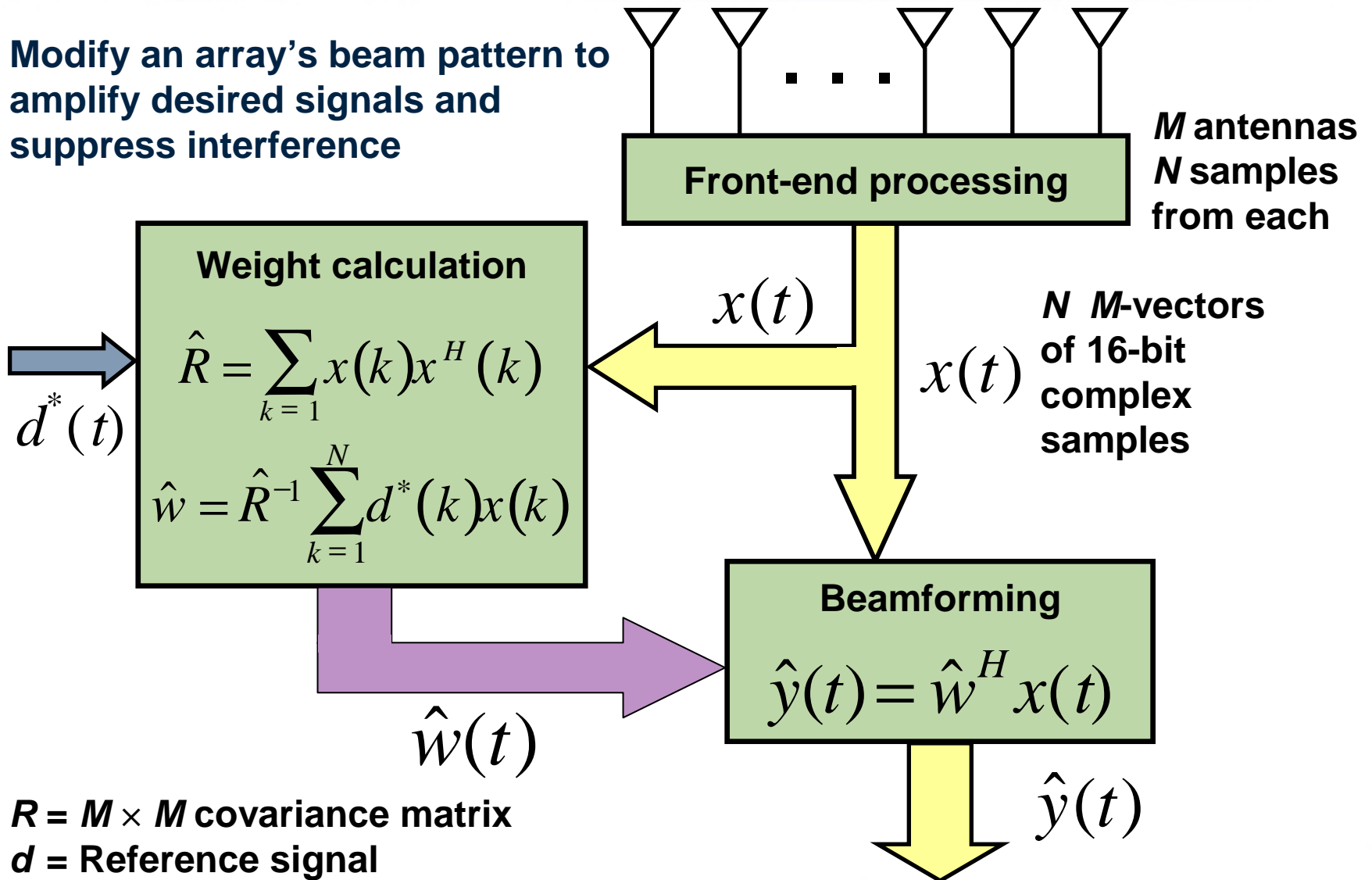**Modify an array's beam pattern to amplify desired signals and suppress interference**

**M antennas**
**N samples from each**

**Front-end processing**

**Weight calculation**

$$\hat{R} = \sum_{k=1} x(k)x^{H}(k)$$

$$\hat{w} = \hat{R}^{-1} \sum_{k=1}^{N} d^{*}(k)x(k)$$

$d^{*}(t)$

$x(t)$

$x(t)$ **N M-vectors of 16-bit complex samples**

**Beamforming**

$$\hat{y}(t) = \hat{w}^{H} x(t)$$

$\hat{w}(t)$

$\hat{y}(t)$

**R = M × M covariance matrix**
**d = Reference signal**

## Background

- **More users than antennas $\Rightarrow$ orthogonal beams not possible**
- **No a priori information about signal directions $\Rightarrow$ need real-time adaptation**
- **Input stream is 16-bit complex data**

## FastMATH Implementation

- **Covariance matrix calculated by *complex matrix-matrix multiplications on $4 \times 4$ submatrices*, then re-assembling full matrix**

- **Covariance matrix inverted by Cholesky decomposition; use *block matrix manipulation* instructions to rearrange input into blocks for SIMD parallelization**

- **Beamforming using matrix-matrix multiplications; more efficient than simple vector math**

## WCDMA Example Results

- **With 64 voice users and 16 antennas, 4 rake fingers per user, weights updated every slot: 0.73 FastMATH processors**

## Algorithms

- *Mitigate interference between users in CDMA*
- Solve for estimators for correct symbols, beginning with user-user correlation matrix *R* and user input vector *y*
- Difference equation for interference on symbol m of desired user from near-by symbols of all other users:

$$y_m = \sum_{k=-K}^{K} R_{m-k} \hat{b}_{m-k}$$

- $\hat{b}$ is desired estimator vector for symbol *m* of *N* users to be found

## Implementation

- *Jacobi iteration*: Solve for matrix *B* of *M* symbols for *N* users. Perform matrix-matrix multiplications distributed over processors
- Calculate correlation matrices *R* on chip; large capacity L2 cache reduces data transfer
- At each iteration exchange partial results over RapidIO port via DMA
- RapidIO interfaces work in background in parallel with computations – data transfer time efficiently hidden

# Scaled Multiprocessor Example: WCDMA Short Code Multi-User Detection

- **Data transfer in parallel with computation**
- **Scalable multiprocessor system distributing tasks and results over RapidIO interface via coherent L2 cache**



**Data Transfer Bandwidth Used: GB/s**

1 Chip

Add 2nd Chip

4 Chips

FastMATH processor

FastMATH processor

FastMATH processor

FastMATH processor

RapidIO-chained processor array

1.0 — RapidIO Bandwidth

*Data transfer: RapidIO ports and large L2 enable up to 134 users*

**Mitigate user-user interference in WCDMA via MUD: Jacobi algorithm**

0.5

PCI Bandwidth

48    68    134    **Users**