**tensilica**

*The Engine of SOC Design*

**DATE 2008:  Tutorial A:**
**Automatically Realising Embedded Systems**
**From High-Level Functional Models:**
*- From Platform-Independent Models to Platform-Specific Implementations for Processor-Centric Design*

Grant Martin
Chief Scientist, Tensilica
Monday 10 March 2008:  0930-1800

---

**tensilica Agenda**

- Platform-Independent Models for processor-centric design
- Options for implementation
- ASIPs
- Determining the platform specifics
  - Manual methods
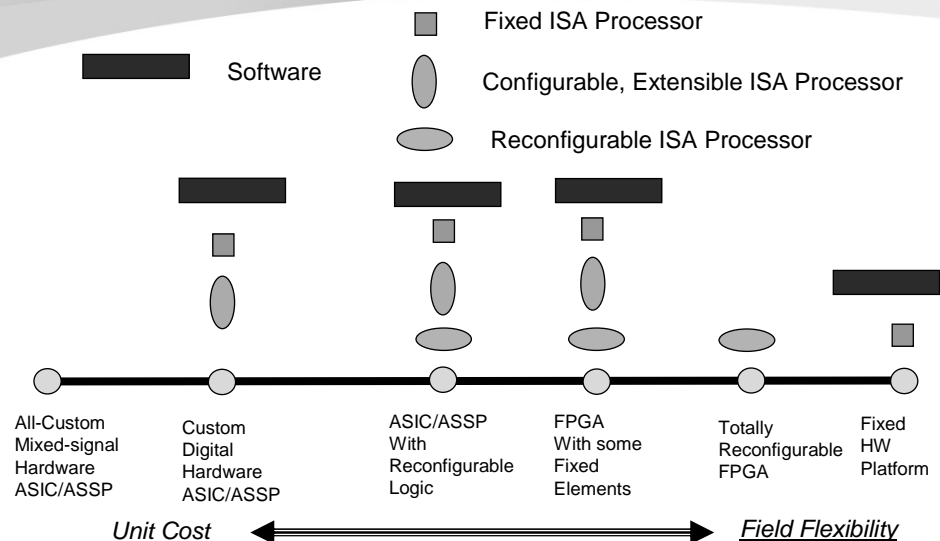  - Automated methods
- Examples and Quality of results

## Platform-Independent Models for processor-centric design

*tensilica*

- Usually algorithmic or control code in C, C++, possibly SystemC (e.g. C/C++ with SystemC interfaces) or other C dialect
- Might be generated from higher level specification: e.g. SDL, a UML profile, Mathworks Matlab/Simulink
- Such models are as "generic as they come"
  – Will compile and execute on just about any host platform
  – Usually will compile and execute (possibly poorly) on embedded processors of interest
- How does a design team decide what implementation platform to map this model to?

© 2008. Tensilica Inc.

---

## The "lumpy continuum" for functional implementation

*tensilica*



Fixed ISA Processor

Software

Configurable, Extensible ISA Processor

Reconfigurable ISA Processor

| All-Custom Mixed-signal Hardware ASIC/ASSP | Custom Digital Hardware ASIC/ASSP | ASIC/ASSP With Reconfigurable Logic | FPGA With some Fixed Elements | Totally Reconfigurable FPGA | Fixed HW Platform |

*Unit Cost* ⟵⟶ *Field Flexibility*

© 2008. Tensilica Inc.

## tensilica  Choosing implementation options

- Systems design teams will choose an implementation form based on many criteria
  - Technical:
    - Performance, power, energy consumption, size, weight, heat, programmability, packaging …..
  - Economics/Business:
    - Cost, Risk, Design effort, Verification effort, Field flexibility, Price, Volume, Life-cycle cost, Expected iterations/derivatives, ….
- Often these criteria narrow the choice window
  - For example, cost + target design time and risk might dictate an ASSP or ASIC choice
  - But even within this window, there are several choices – e.g. custom digital HW, SW on different kinds of processor
- Performance considerations will usually further narrow the choice
  - Design teams need to consider *all* the reasonable alternatives to optimise product architectures

---

## tensilica  Processor options

- Not as narrow a choice as it might seem
  - Processors may vary by 100 to 1 or more on criteria such as performance and energy consumption for a specific task
- Even for embedded applications, processors range widely in their suitability for specific applications or tasks
  - As long as they have sufficient memory, any processor should be able to run any task
    - Thus all processors have flexibility
    - But specific processor features make one more or less suited to particular tasks

## Processor features

*tensilica*

- Rated in order of application-specificity (approximate)
  - Size, number and types of local and system memories
  - Multi-operation instructions (for applications with high unstructured concurrency
- Following often found in DSPs in particular
  - SIMD operations (for applications with high structured concurrency – e.g. vectorisation)
  - Low or zero-overhead looping
  - Special function units – multipliers, MACs
  - Rounding and truncation instructions
- Usually found in ASIPs:
  - Highly specialised application-specific instructions

---

## ASIPs

*tensilica*

- Application-Specific Instruction-set Processors
- The highest form of platform-specific implementation for processor-centric design
- Might be manually designed
- However, usually designed through automated tool flows driven by:
  - GUI-based configuration parameters
  - Architecture description languages (ADL's)
- Considerable research into this
  - E.g. Masaharu Imai (Osaka U.) PEAS-I,II, III, ASIPMeister
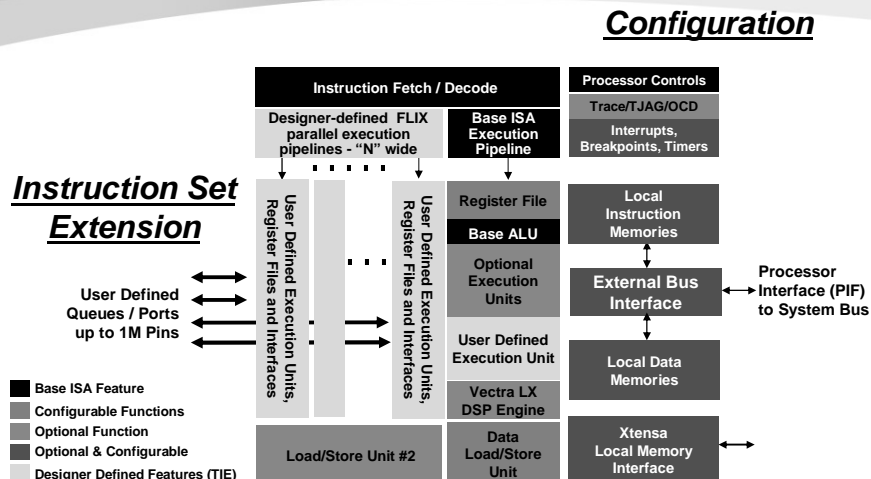- Many commercial capabilities

## Architecture-Description Language ASIPs

- Often resulted from research projects
  - E.g. nML – TU Berlin
  - Lisa – RWTH Aachen
- Research often used to generate software tools only
  - Eg Chess/Checkers – ISS and compilers– IMEC
  - Lisa – ISS
- Sometimes commercialised
  - Target compilers
  - Axys – now ARM; Lisatek – now CoWare
- Sometimes entirely commercial
  - E.g. Tensilica Instruction Extension Language (TIE)
    - Used for most of HW and SW toolchain but not for all parts of processor
    - Rest configured using GUI configurator

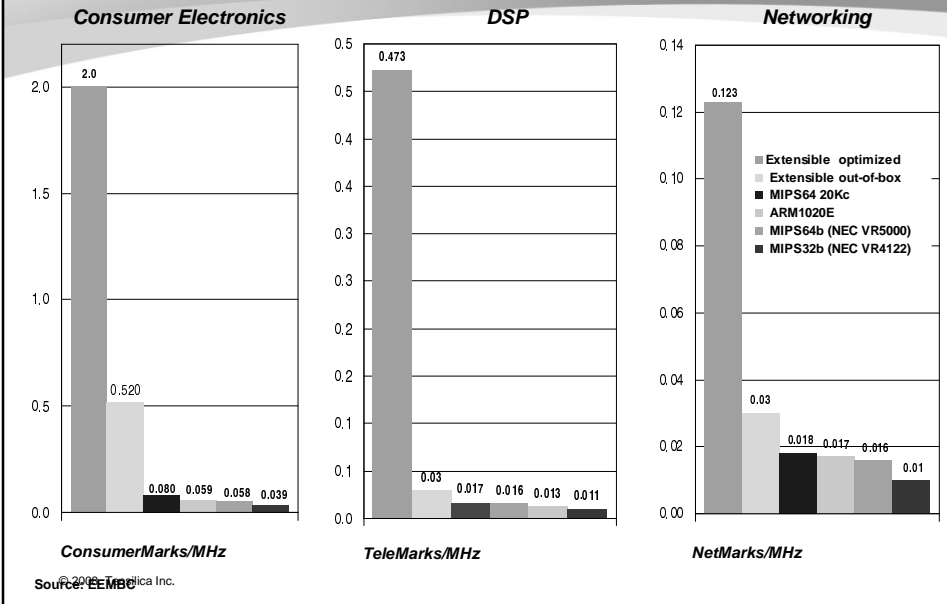## Example of configuration/extension space for ASIP

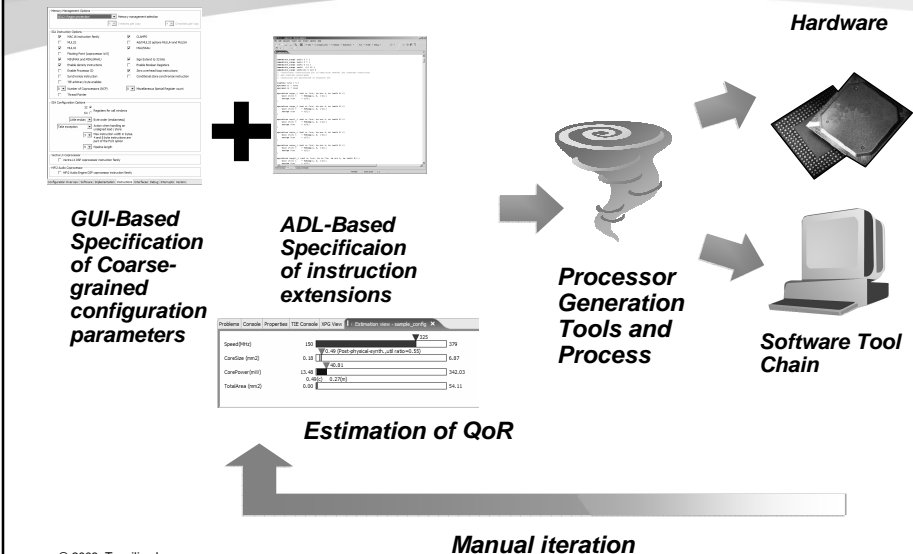## Benefits of Configurability

**Consumer Electronics**　　　**DSP**　　　**Networking**

Consumer Electronics chart:
- 2.0
- 0.520
- 0.080  0.059  0.058  0.039

*ConsumerMarks/MHz*

DSP chart:
- 0.473
- 0.03  0.017  0.016  0.013  0.011

*TeleMarks/MHz*

Networking chart:
- 0.123
- 0.03
- 0.018  0.017  0.016  0.01

Legend:
- Extensible optimized
- Extensible out-of-box
- MIPS64 20Kc
- ARM1020E
- MIPS64b (NEC VR5000)
- MIPS32b (NEC VR4122)

*NetMarks/MHz*

Source: EEMBC

© 2008. Tensilica Inc.

---

## An example ASIP generation flow

*Hardware*

***GUI-Based Specification of Coarse-grained configuration parameters***

***ADL-Based Specificaion of instruction extensions***

***Processor Generation Tools and Process***

***Software Tool Chain***

Problems | Console | Properties | TIE Console | NPG View | Estimation view - sample_config
- Speed(MHz)
- CoreSize (mm2)
- CorePower (mW)
- TotalArea (mm2)

***Estimation of QoR***

***Manual iteration***

© 2008. Tensilica Inc.

## Determining the platform specifics for ASIPs

- Manual Methods
  - Profiling source code execution on target processor
  - Finding time consuming routines
  - Finding time consuming loop nests
  - Devising new instructions
  - Optimising memory accesses
  - Modifying source code, and
    - ………….. Round the loop again
- Automated Methods
  - Seek to move from source code to ASIP configuration in a single step

© 2008. Tensilica Inc.

## Manual profiling for optimisation



**Profile Disassembly:**
**View disassembled instructions with cycle counts**

**Pipeline View:**
**See pipeline stalls due to cache misses, branch interlocks, etc**

© 2008. Tensilica Inc.

## Manual profiling for optimisation

*Round the loop comparison of different configurations*

*Exploration of alternative cache configurations*

## Energy Consumption Profiling and Optimisation

**GUI-based assistance for generating fused instructions manually**

© 2008. Tensilica Inc.

---

**Automated methods for ASIP specification**

- Long a research topic
  – For example, the work of Ienne and Pozzi at EPFL has been reported in many papers over the years
    • Concentrated on instruction fusions
- Also been available for some commercial extensible processors
  – For example, Tensilica's XPRES tool

© 2008. Tensilica Inc.

**Automated Methods of ASIP generation**

*Compile Original C/C++ Code with XCC*

```
int main()
{
 int i;
 short c[100];
 for (i=0;i<N/2;i++)
 {
```

*Run XPRES Compiler*

*Designer selects "best" configuration*

Application Performance / Hardware Cost Tradeoff

*Run Xtensa Processor Generator*

*Compile & run original, unmodified C code*

```
int main()
{
 int i;
 short c[100];
 for (i=0;i<N/2;i++)
 {
```
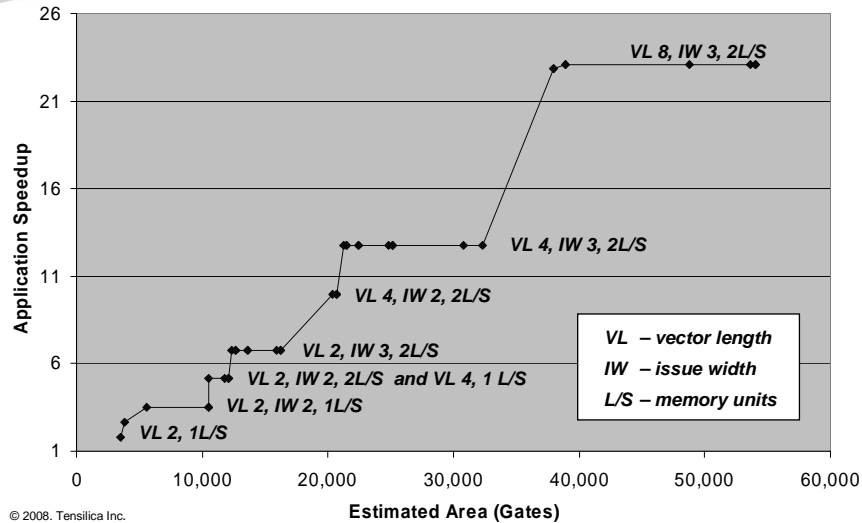
**Evaluates millions of possible extensions using SIMD/vector operations, operator fusion and parallel execution**

**Build processor, including RTL, SW, etc.**

**Build chip, system**

*Option: manually refine configuration*

© 2008. Tensilica Inc.

---

**XPRES Compiler Flow**

- Profile / Analyze application
  - Identify performance critical functions / loops
  - Collect operation mix, dependence graphs
  - Provide feedback to user – code transformations to better target meta-ISA (esp. vectorization)
- Generate sets of ISA extensions
  - Each set implements some dimension of meta-ISA
  - Evaluate each set across all functions / loops
  - Performance and cost estimates allow exploration of large design space
- Collect ISA extensions that together provide maximum performance improvement
  - Each collection of ISA extensions forms ASIP
  - Generate family of ASIPs for varying hardware cost budget

© 2008. Tensilica Inc.

## Example: Family of Architectures

**for (i=0; i<n; i++) c[i] = (a[i] * b[i]) >> 4;**



VL 8, IW 3, 2L/S

VL 4, IW 3, 2L/S

VL 4, IW 2, 2L/S

VL 2, IW 3, 2L/S

VL 2, IW 2, 2L/S  and VL 4, 1 L/S

VL 2, IW 2, 1L/S

VL 2, 1L/S

*VL  – vector length*
*IW  – issue width*
*L/S – memory units*

Application Speedup

Estimated Area (Gates)

© 2008. Tensilica Inc.

---

## TIE Techniques Employed by XPRES Compiler:  3 kinds of parallelism

- **Instruction-Level Parallelism: FLIX**
  – Bundling independent operations
- **Data-Level Parallelism:  SIMD / Vector**
  – Single instruction, multiple-data operations
- **Pipeline Parallelism: Fusion**
  – Merging of compound, dependent operations
  – E.g. a multiply-add or a multiply feeding a shift


- **All three techniques can be combined**
  – E.g. Parallel FLIX execution of Fused-SIMD operations

© 2008. Tensilica Inc.

## Instruction-Level Parallelism (ILP)

- XPRES supports ILP with multiple-issue
  - 32-bit or 64-bit instructions
  - Each instruction can contain 1-15 slots
  - Each slot can contain arbitrary mix of operations
- Implemented in Xtensa using FLIX
  - Similar to VLIW without code size increase
- Xtensa C/C++ compiler exploits FLIX automatically
  - Software pipelining for loops
  - List scheduler for other code
  - Code scheduling and packing in both compiler and assembler

---

## XPRES: Instruction Parallelism (exploited using FLIX)

### Original C Code

```
for (i=0; i<n; i++)
    c[i] = (a[i] * b[i]) >> 4
```

- **Performance**
  - **Increased IPC**
- **Hardware cost**
  - **Replicated function units**
  - **Additional register file ports**

### 64-bit 3-issue FLIX

```
format f  64 { s0, s1, s2 }
slot_opcodes s0 { L32I, S32I }
slot_opcodes s1 { SRAI, MULL }
slot_opcodes s2 { ADDI }
```

Generated Assembly
(1 iteration in 3 cycles, 2.6x)

```
loop:
 { l8ui a8,a9,0;    mul16u a12,a10,a8; addi  a9,a9,1 }
 { l8ui a10,a11,0;  nop;               addi  a11,a11,1 }
 { s8i  a13,a14,508; srai  a13,a12,4;  addi  a14,a14,1 }
```

### 64-bit 2-issue FLIX

```
format f  64 { s0, s1 }
slot_opcodes s0 { L32I, S32I }
slot_opcodes s1 { ADDI, SRAI, MULL }
```

(2 iterations in 8 cycles, 2x)

```
loop:
 { l8ui   a4,a11,12; addi    a13,a13,4 }
 { l8ui   a2,a9,12;  mul16u a3,a3,a2 }
 { l8ui   a14,a9,4;  srai    a12,a14,4}
  …
```

Results: Instruction Parallelism

© 2008. Tensilica Inc.

# Data Parallelism

- XPRES supports data parallelism with SIMD operations
  - Vectors of length 2, 4, 8, and 16 (up to 256 bits)
  - Support for unaligned vector loads and stores
  - Vector C–operators, MIN, MAX, ABS, reductions
  - Vector user-defined operations
- Xtensa C/C++ compiler exploits automatically
  - Automatic loop vectorization
  - Source attributes and compiler options specify aliasing and alignment directives to enable additional vectorization opportunities

© 2008. Tensilica Inc.

## XPRES: Data Parallelism (exploited using SIMD)

**Original C Code**

```
for (i=0; i<n; i++)
    c[i] = (a[i] * b[i]) >> 4
```

- **Performance**
  - **Increased computation bandwidth**
- **Hardware cost**
  - **Replicated function units**
  - **Vector register file(s)**

**Vector Length 8**

```
regfile vr8x8 64 4 v
regfile vr16x8 128 2 x
operation ashri16x8 { out vr16x8 a, …}
…
```

Generated Assembly
(8 iterations in 6 cycles, 10.6x)

```
loop:
   liu8x8         v0,a8,8
   ashri16x8      x0,x0,4
   liu8x8         v1,a9,8
   cvt16x8sr8x8s  v2,x0
   mpy8r16x8u     x0,v0,v1
   siu8x8         v2,a10,8
```

**Vector Length 2**

```
regfile vr8x2 16 4 v
regfile vr16x2 32 2 x
operation ashri16x2 { out vr16x2 a, … }
…
```

© 2008. Tensilica Inc.

## Results: Data Parallelism



© 2008. Tensilica Inc.

14

## Pipeline Parallelism

- XPRES supports pipeline parallelism with fused operations
  - Fused operation composed of two or more other operations, plus possibly constant values
  - Latency of fused operation usually less then combined latency of composing operations
  - Limits on input/output operands, number of composing operations, hardware cost, max latency, etc.
  - Graphical support for manual fused operation generation
- Tradeoff performance and generality
  - Performance: large fusions w/ fixed constants
  - Generality: smaller fusions, fewer fixed constants
- Xtensa C/C++ compiler exploits automatically
  - Fused operations automatically replace sequences of composing operations

---

## XPRES: Pipeline Parallelism (exploited using Fusion)

Original C Code

```
for (i=0; i<n; i++)
    c[i] = (a[i] * b[i]) >> 4
```

- **Performance**
  - **Decrease instruction count**
  - **Decrease computation latency**
- **Hardware cost**
  - **Logic to share function units**
  - **Register file ports**



Generated Assembly
(1 iteration in 5 cycles, 1.6x)

```
loop:
    l8ui    a12,a11,0
    l8ui    a13,a10,0
    addi.n  a10,a10,1
    addi.n  a11,a11,1
    fusion.mul16u.srai.s8i.addi a9,a12,a13
```

## Results: Pipeline Parallelism



Legend: idwt, autocor, rgbcmyk, fft-radix4, gsm encode, gsm decode, compress95, djpeg, m88ksim

Y-axis: Application Speedup (1 to 1.8)
X-axis: Estimated Area (Gates) (0 to 40,000)

© 2008. Tensilica Inc.

---

## Exploiting Multiple Types of Parallelism

- XPRES combines FLIX, SIMD, and Fusion

  Original C Code

  ```
  for (i=0; i<n; i++)
      c[i] = (a[i] * b[i]) >> 4
  ```

  FLIX and Vectorization  (16 iterations in 4 cycles, 32x)

  ```
  { cvt16x8sr8x8s  v2,x1;      liu8x8  v0,a8,8;    mpy8r16x8u  x2,v1,v2 }
  { si8x8          v2,a10,8;   liu8x8  v1,a9,8;    ashri16x8   v1,v0,4  }
  { siu8x8         v3,a10,16;  liu8x8  v2,a9,8;    ashri16x8   x3,x2,4  }
  { cvt16x8sr8x8s  v3,x3;      liu8x8  v1,a8,8;    mpy8r16x8   x0,v0,v1 }
  ```

  FLIX, Vectorization, and Fusion  (16 iterations in 3 cycles, 42.6x)

  ```
  { si8x8    v2,a10,8;     liu8x8  v0,a8,8;    fusion.mpy.ashri.cvtx8  v2,v0,v1 }
  { siu8x8   v5,a10,16;    liu8x8  v1,a9,8;    fusion.mpy.ashri.cvtx8  v5,v3,v4 }
  { liu8x8   v3,a8,8;      liu8x8  v4,a9,8;    nop }
  ```
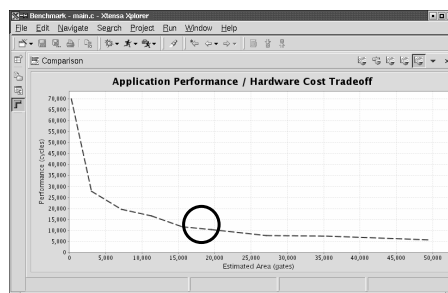
  © 2008. Tensilica Inc.

# Instruction + Data + Pipeline Parallelism

---

# Example: SAD (Sum of Absolute Differences) Automatic Solution Search

**Visually Presented**



***Wide Range of Choices of
Performance Increase -v- Hardware Cost***

***Generated Xtensa LX
Configuration Parameters***

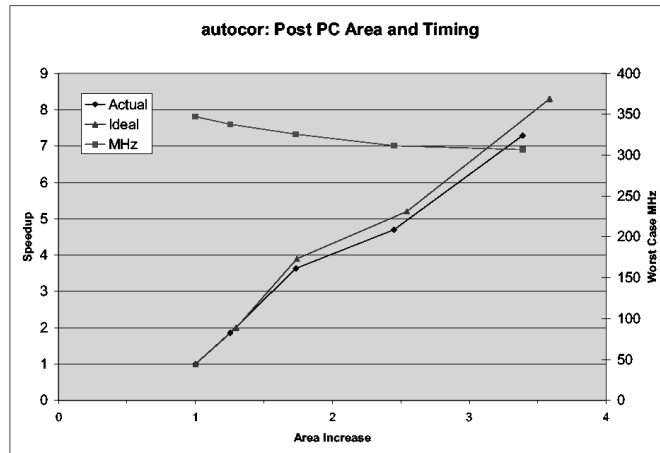| i | Speedup -v- Base processor | Gates Added | SIMD Factor | FLIX Width (Slots) | Load / Store Units | Fusion |
|---|---|---|---|---|---|---|
| 1 | 18.6x | 49,681 | 8 | 3 | 2 | Yes |
| 2 | 14.1x | 35,439 | 8 | 2 | 2 | Yes |
| 3 | 13.9x | 27,082 | 4 | 3 | 2 | Yes |
| 4 | 10.1x | 19,547 | 4 | 2 | 2 | Yes |
| 5 | 9.0x | 15,527 | 2 | 3 | 2 | Yes |
| 6 | 6.8x | 8994 | 4 | 1 | 1 | Yes |
| 7 | 5.3x | 7158 | 2 | 2 | 1 | Yes |
| 8 | 1.5x | 301 | 1 | 1 | 1 | Yes |

## Example: SAD Configuration #4 – Generated TIE



## Exploration of Estimation vs. Actuals – Accuracy of design tradeoffs – monotonicity and correlation

**Experiments using Physical Compiler**
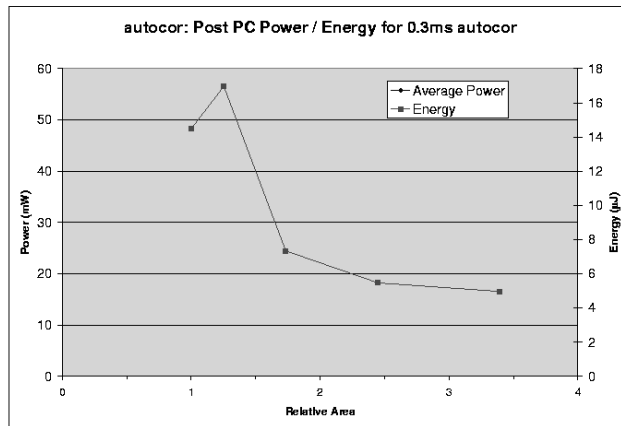
autocor: Post PC Area and Timing

© 2008.



**Experiments looking at power/energy**

autocor: Post PC Power / Energy

© 2008. Tensilica Inc.

**tensilica**

# Energy consumption for fixed performance

**autocor: Post PC Power / Energy for 0.3ms autocor**

---

**tensilica**

# Experiments with fft application

**fft: Post PC Area and Timing**

## Improving the quality of results of automated methods

- Although the ideal is to use unmodified source code, often simple transformations will yield much better results
- Maximise instruction fusion
  - Sometimes can combine manual (deep) fusions with automated techniques
- Maximise utilisation of multi-operation instructions:
  - Enable software pipelining
  - Enable loop unrolling
  - Fixing aliasing
  - Removing control flow (if not supported in automated techniques or microarchitecture)

## Maximise use of SIMD:

- Expose data-wise parallelism and enable vectorisation
- Aliasing
- Strides and gaps to allow efficient data loading and storage via contiguous data
- Outer vs. inner loop vectorisation in nest
- Function calls – inline where possible
- Data alignment
- Pointer dereferencing

## tensilica Summary

- ASIPs are a powerful mechanism to go from platform-independent models to highly efficient implementations of software on embedded processors
- Many choices of ASIP specification methods available
  - Both commercial and research
  - Manual and automated
- Use of fine-grained instruction extension and coarse-grained structural adaptation can mean up to a 100 to 1 performance improvement over a general purpose processor, up to 10 to 1 energy reduction, in similar or less area

## tensilica References

- ADLs:
  - Prabhat Mishra and Nikil Dutt, "Processor Modeling and Design Tools", Chapter 8 in Volume I: *EDA for IC System Design, Verification and Testing*, in the *Electronic Design Automation for Integrated Circuits Handbook*, edited by Louis Scheffer, Luciano Lavagno and Grant Martin, CRC/Taylor and Francis, 2006.
- ASIPs:
  - Akira Kitajima, Makiko Itoh, Jun Sato, Akichika Shiomi, Yoshinori Takeuchi, Masaharu Imai: Effectiveness of the ASIP design system PEAS-III in design of pipelined processors. *ASP-DAC 2001*: 649-654
  - Paolo Ienne and Rainer Leupers (editors), *Customizable Embedded Processors*, Elsevier Morgan Kaufmann, 2006
  - Matthias Gries and Kurt Keutzer (editors), *Building ASIPS: the MESCAL methodology*, Springer, 2005.
- Tensilica Examples
  - Chris Rowen and Steve Leibson, *Engineering the Complex SOC: Fast, Flexible Design with Configurable Processors*, Prentice-Hall PTR 2004.
  - Steve Leibson, *Designing SOCs with Configured Cores: Unleashing the Tensilica Xtensa and Diamond Cores*, Elsevier Morgan Kaufmann 2006.