David May

Cambridge December 2005

### Long term trends

since the 1940s, doubling about every 18 months Computer performance/cost has followed an exponential path

This has enabled

Increasing function/performance at constant cost

: 으

Decreasing cost at constant function/performance

December 2005

### Moore's Law

made life easy - more transistors, faster clocks ... For three decades, Moore's Law applied to microcomputers has

almost every known technique was applied to use up the Over the 1990s there was little architectural innovation, but transistors!

innovate This is changing - to keep on the long term trend, we will need to

... especially in architecture, design and software ...

## Increasing performance

transistors (Moore's Law) ... has involved increasing clock speeds and using more

pertormance. transistors does not always provide a corresponding increase in ... but increasing the processor clock speed and using more

transistors ... because the speed of the wires can't keep up with the

We are well behind Moore's law on achieved general purpose pertormance

### Decreasing cost

telephones, cameras, games, toys ...) ... has had a dramatic impact allowing computers to be embedded in an increasing variety of products (cars,

10, increase the volume by factor 100! improvements in performance/cost - reduce the cost by factor New applications will continue to emerge as a result of

than increasing performance Commercially - and socially - decreasing cost has more effect

Note: Disposable Computing started around 2000!

## Computer Architecture

There's little evidence to support architectural decisions

... we don't seem to have a good analysis of what programs do

... we are often fooled by averages

and to make matters worse

together! ... computers, compilers, languages - and applications evolve

Cambridge

## Computer Architecture

of an existing software base incremental enhancement - trying to increase the performance In general purpose architectures, we have resorted to

programming languages? Concurrency? Input and Output? What about ideas that can't be represented easily in existing

What about new applications?

innovation In embedded architectures, there has been more scope for

Cambridge December 2005

# General purpose microprocessors

Faster and faster clocks

Deeper and deeper pipelines

More and more execution units

More and more strange instructions

Longer and longer context switches Bigger and bigger branch prediction/recovery mechanisms

More and more stuff offloaded to hardware gadgets

Cambridge

# General purpose pertormance

... relies on uniform access to a random access memory

A complex memory hierarchy is an approximation to this

... some programs work well; some don't

approximation to fast sequential processing In the same way, Instruction Level Parallelism provides an

Maybe we've made these processors too fast!

Cambridge

## General purpose computing

A lot of effort is going into big multi-core chips

complex memory hierarchies Predictably these focus on symmetric multiprocessors with

Its not clear what these are for

- we won't need or want big processors in PCs we'll be using thin clients and portables
- we will need big servers but these can be built from larger numbers of slower, more power-efficient processors

David May

## Expectations vs. Reality

Symmetric multiprocessors are easy to program

... provided you're not bothered about performance!

its memory system A single general purpose processor places a heavy demand on

... you can't expect it to support several processors

algorithms - but then the architecture can be much simpler! Multiprocessors need concurrent programs and/or parallel

David May

#### David May

### The Exabyte effect

The Internet in 2010:

- Exabytes everywhere!
- Petabyte servers, petaflop supercomputers, a billion hosts, a zetabyte online
- Services to ubiquitous and nomadic clients including wearables
- Many different forms of content visuals, soundscapes ...

Cambridge

### The Exabyte effect

maximise energy efficiency We will see changes in Internet infrastructure aiming to

These will include

- more, lower-speed processors
- programmable accelerators for specialised tasks such as transcoding, cryptography, searching ...

towards Ethernet as the standard interconnect The *accelerator-in-server* market will be opened up by the trend

#### David May

## Architecture for exabytes

We're going to need

... big addresses

... arithmetic and logic on big numbers

... accurate numerics on big numbers

... efficient data representation

#### count by factor 4 Long arithmetic

instructions - if we have the right instructions! To multiply two n-word integers we have to execute  $n^2$  multiply

carryout, result := MUL op1, op2, carryin

and doubling the wordlength reduces the multiply operation

times faster than a 32 bit CPU with the wrong ones! So (eg) a 64 bit CPU with the right instructions will be 10-100

Let's think about a machine with 128-bit registers:

... big integers:  $2^{128}$  is around  $10^{39}$ 

... 64.64 format gives a range of around  $10^{20}$  at 1 in  $10^{20}$  precision

... enough addresses for every byte - or bit - on the planet

... accurate double precision floating point multiply-accumulate

powerful SIMD data representation

## Computing without power

without the need for wires - or replaceable batteries We have the potential for a new generation of sealed devices -

tracking, implants, packaging ... RFID is just the beginning be everywhere - in environmental control, industrial monitoring, A disruptive technology - sentient, communicating devices will

#### They depend on

- low-energy computing and communications
- embedded sensors and actuators
- low-cost batteries or scavenged power

December 2005

#### Wearables

We can embed technology in clothes - or wear it like jewellery!

gyros, GPS, RFID .. Technologies include audio, cameras, accelerometers and

Applications include sports, healthcare, lifestyle, leisure ...

But the big markets will probably be in fashion!

electronic devices? Will the design houses play a major role in the next wave of

#### Event driven systems including software Minimal operations including data transfers Dynamically switching off stuff when it isn't in use Low power logic design Low voltage circuits

What does low energy involve?

David May Cambridge December 2005

#### Software

compensating Moore's Law May's Law: Software efficiency halves every 18 months,

A mixture of

- shortage of skills
- adding too many features
- copy-paste programming
- massive overuse of windows and mouse-clicks
- reliance on Moore's law to solve inefficiency problems

#### Software

- ... together with an extreme reluctance to re-write software
- when it's full of bugs, is too big and complex to understand, and the authors have left (died?), it's time to re-write it!
- or at least, it should be moved from the company assets to the liabilities!

value by increasing performance and power-efficiency in This is a big opportunity - efficient software can add a lot of

- ubiquitous systems
- wearables
- high-performance systems and supercomputers

## Software and Algorithms

double the battery life In ubiquitous systems, halving the instructions executed can

algorithms: And big data sets bring big opportunities for better software and

a dramatic effect when N is large Reducing the number of operations from  $N \times N$  to  $N \times log(N)$  has

technology improvements! ... for N=30 billion this change is as good as 50 years of

Cambridge

December 2005

### Generic products

costs are escalating 'State-of-the-art' design, verification and manufacturing set-up

This suggests a move to:

- generic programmable and/or reconfigurable chips
- generic 'platforms' customisable at low cost
- :: ??

Potential opportunity for new industry ...

# The Fabless, IP-less chip company

the ASICs made possible by generic platforms Manufacturing companies will not be able to design a fraction of

programming and configuring platforms - and by selling the result - not by selling the programs and configurations So there's space for new companies which do ASICs purely by

and market expertise but will not require huge investment These companies will have specialised high-value applications

Cambridge December 2005

#### David May

## Innovations in Architecture

Concurrency, concurrency and concurrency

inside processors

in collections of processors

in systems on a chip

Expose it and exploit it - don't hide it

### Concurrency

few computers in the 1980s but have only just come into widespread use: Large collections of processors were successfully employed in a

... Google search engine - tens of thousands of processors

. digital animation (typically 1,000 processors)

... supercomputers - a few thousand processors

and we are beginning to see

... chips full of processors

### Processor arrays

... are relatively easy to design, verify and test

... can exploit local clocks for high speed

... could deliver tera-op performance within 5 years

... require innovative programming tools

They are potentially an important generic platform technology

#### David May

# General purpose processor arrays

Two key issues:

- scalable interconnect grows as  $n \times log(n)$  just like random access memory addressing!
- subroutines efficient serial re-use of parallel resources

essence of many algorithms and applications Interconnect is not an 'overhead' - communication is the

### Innovations in Design

programming languages for years We have been making hardware design tools look like

target software, re-configurable hardware - and hardware from the same source But we have failed to produce a single language which can

on design efficiency - and on the efficiency of designs! So this is still a potential opportunity. It would have a big impact

Concurrent languages are the key!

#### Robotics

... has finally come of age!

we know how to build the control systems

we know how to do the sensors - even vision

we have lightweight materials

there is a market - vacuum cleaners, lawn-mowers, pets ...

... a market for sensors, actuators and embedded intelligence

Cambridge

#### Robotics

Scalable real-time control - computers in the loop:

thousands of unnamed ones ..." "... there are about 650 named muscles in the human body and

or perhaps:

and thousands of unnamed ones ..." "... there are around 650 named actuators in a humanoid robot

and every one will need a microcontroller!

Cambridge

### New technologies

the 1940s and there are likely to be more ... There have been several jumps in computer technology since

... exotic technologies based on molecular structures

plastic devices which are flexible - and printable

... and with silicon, we can continue to reduce costs

And even battery technology is moving on - to flexible, printable devices ...

David May

#### Summary

The long-term cost-performance improvement will continue

There is plenty of scope for new ventures

There is a big opportunity for innovation

- architecture and software
- business models
- applications
- technology