

# Optimising code for apeNEXT

D. Pleiter

DESY

28. November 2006, LAP06, Zeuthen



# Agenda

Introduction

Hardware Characteristics

System Software

Optimisation

Optimisation Steps

Summary

# Agenda

Introduction

Hardware Characteristics

System Software

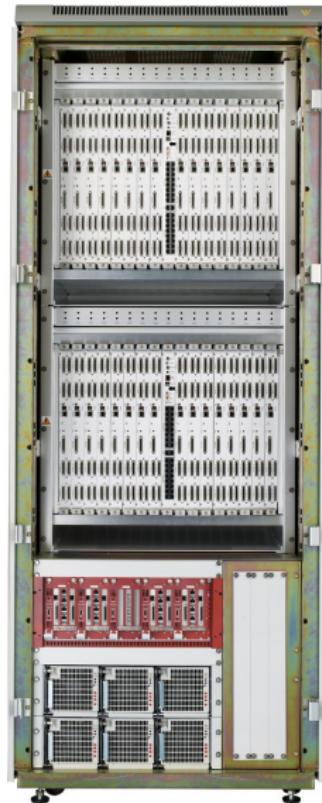
Optimisation

Optimisation Steps

Summary

# Introduction

- ▶ Joined European development project
  - ▶ DESY
  - ▶ Istituto Nazionale di Fisica Nucleare (INFN)
  - ▶ University Paris-Sud
- ▶ Fully custom designed machine
  - ▶ Custom hardware
  - ▶ Custom software



# Agenda

Introduction

Hardware Characteristics

System Software

Optimisation

Optimisation Steps

Summary

# Floating Point Unit (FPU)

- ▶ FPU performs one operation  $a \times b + c$  per clock cycle, where  $a, b, c$  complex numbers
- ▶ 64-bit IEEE floating point numbers

Peak performance

→ 8 Flops / cycle = 1.2 GFlops/sec

# Memory Hierarchy

- ▶ Memory controller
  - ▶ Supports 256 MBytes DDR-SDRAM (with ECC)
  - ▶ Maximum bandwidth:  $2 \times 64$  bit word/cycle

## Peak performance

→ 2.4 GBytes/sec or 0.5 Flop/Byte

- ▶ Very large register file:  $2 \times 256$  64-bit registers
- ▶ Data pre-fetch queues (no data cache)
- ▶ Instruction buffer
  - ▶ Burst instruction fetch
  - ▶ Different operation modes:
    - ▶ Fifo mode: read once, execute once (default)
    - ▶ Cache mode: Prefetch instructions for **multiple** execution

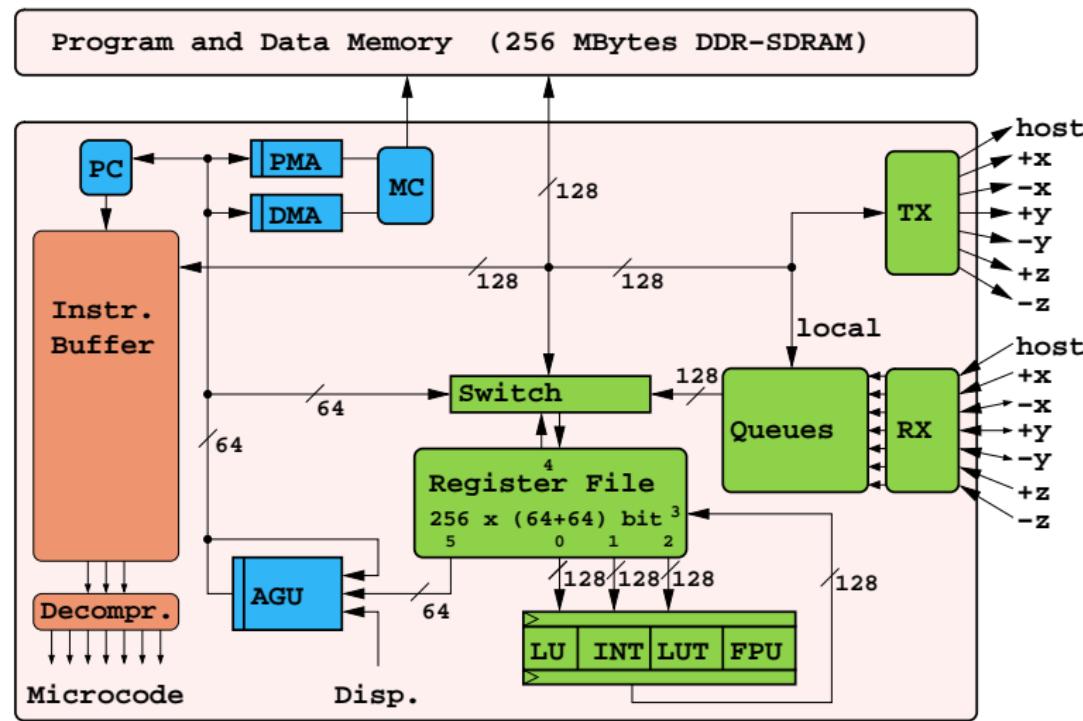
# Network

- ▶ 7 bi-directional LVDS links:  $\pm x$ ,  $\pm y$ ,  $\pm z$ , 7th
- ▶ Gross bandwidth per link and direction: 8 bit/cycle
- ▶ Concurrent send and receive and concurrent transfer along orthogonal directions
- ▶ Low latency:  $\approx 125$  ns
- ▶ Transmission by frames of 128 bit data + 16 bit CRC

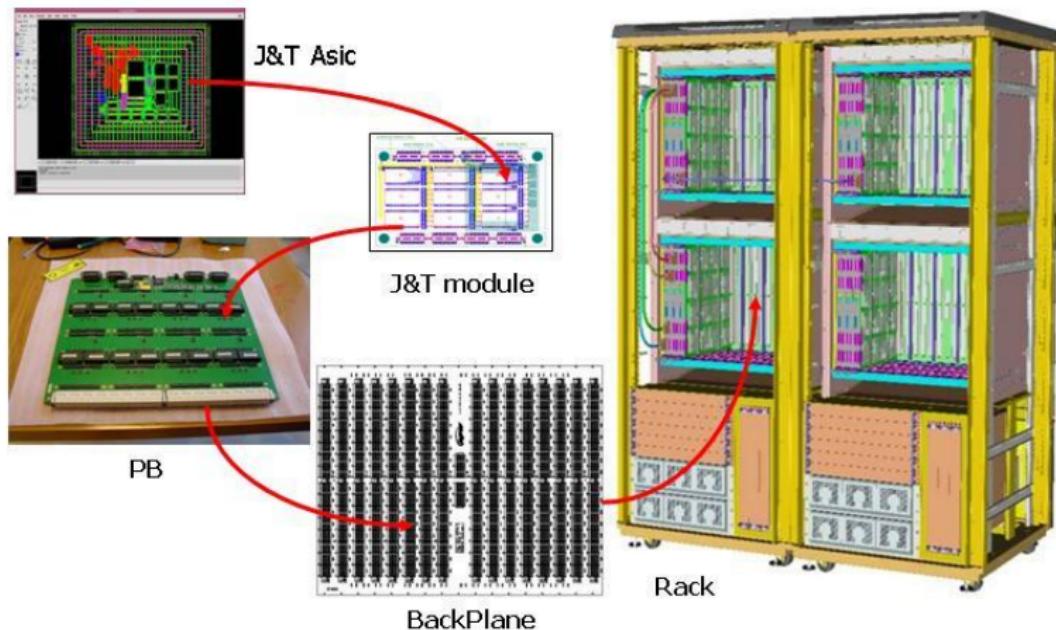
## Performance

→  $< 16 \text{ Bytes} / 18 \text{ cycles} = (1/16) \text{ memory bandwidth}$

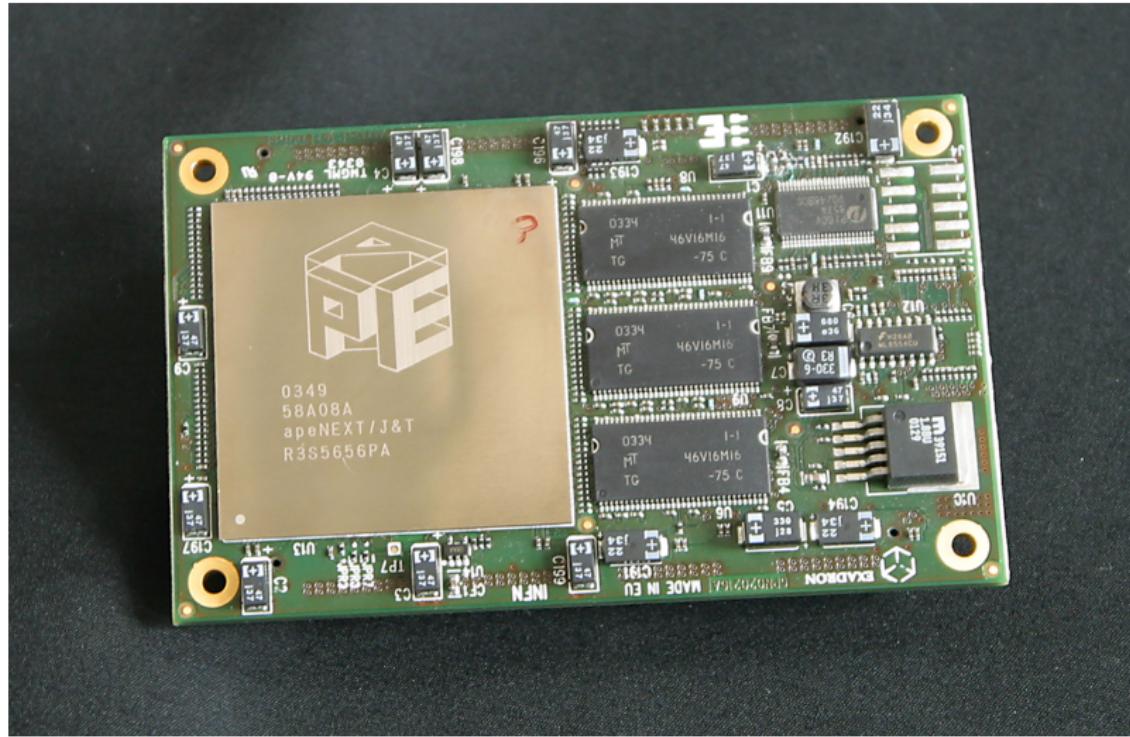
# Processor Overview



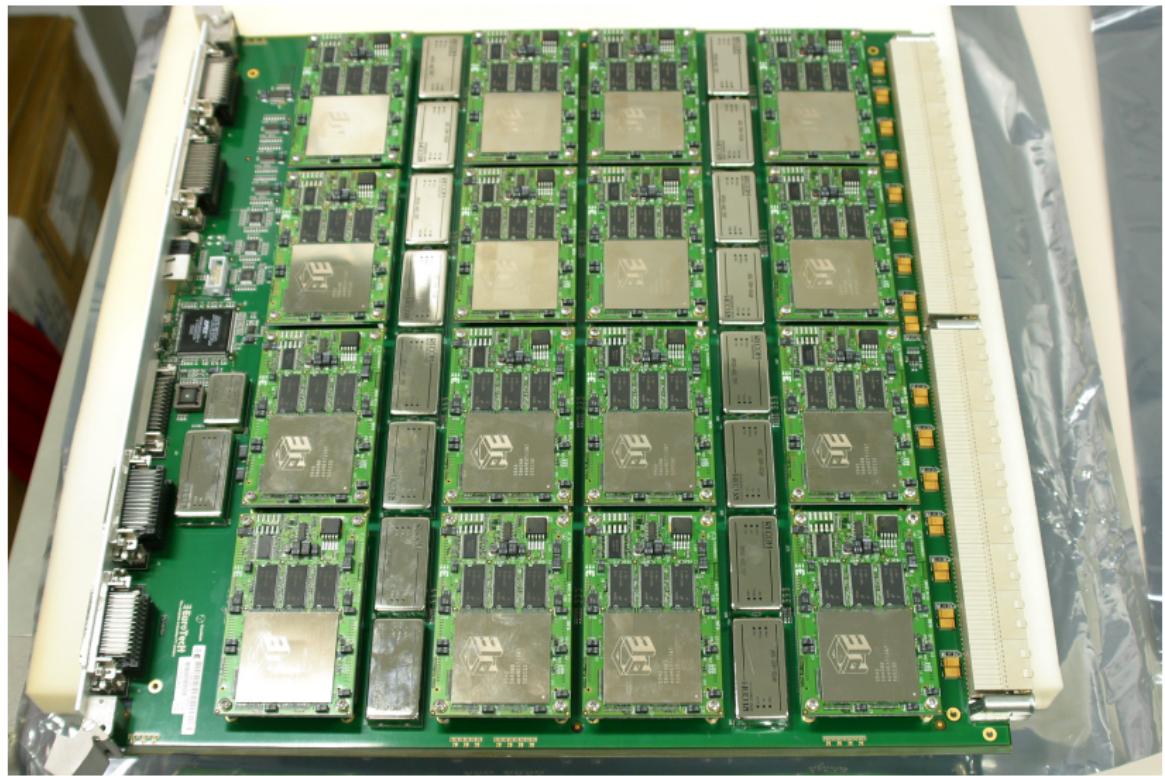
# System Overview



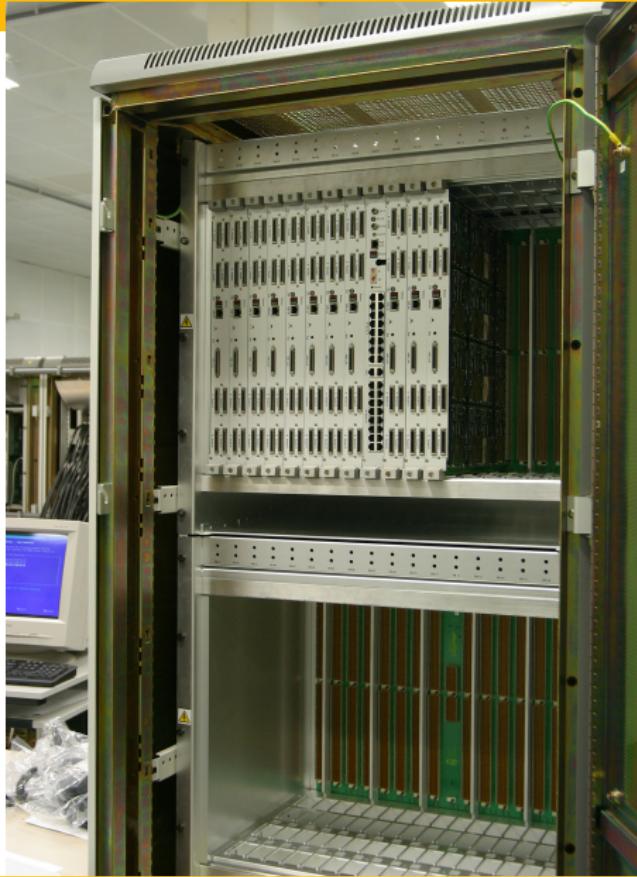
# Processing Module



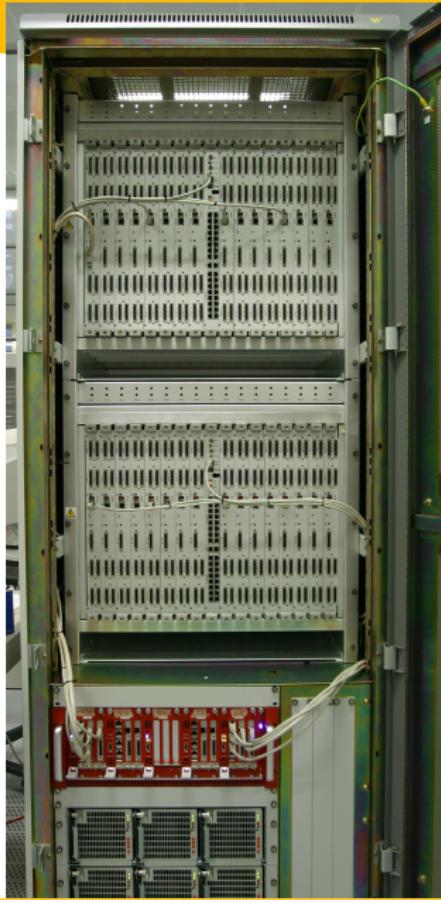
# Processing Board



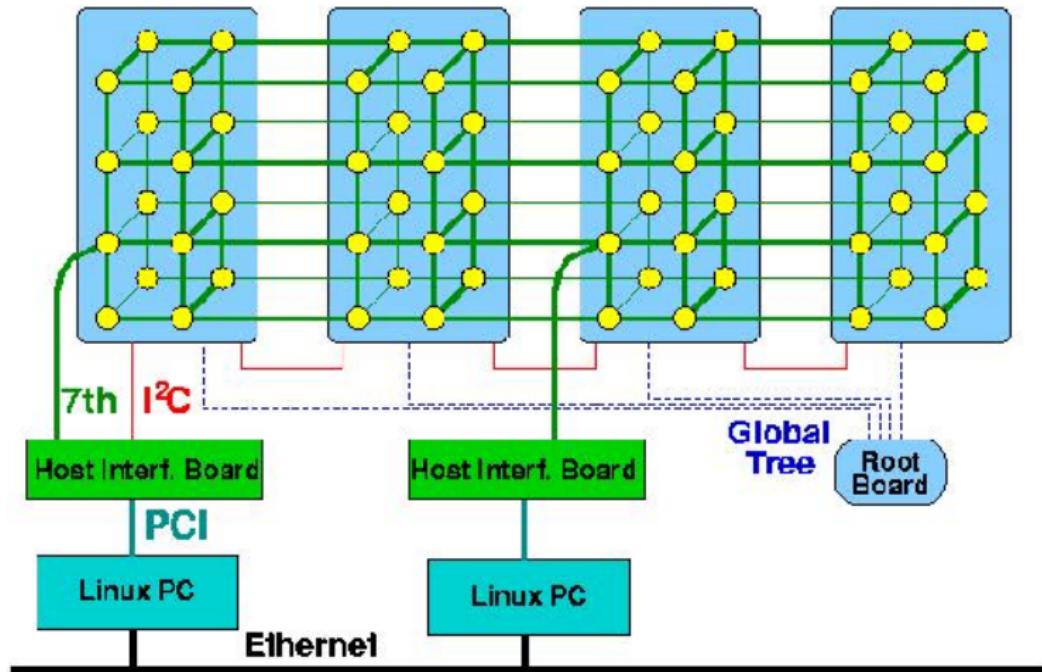
# During Installation



# Complete Rack



# Global Architecture



# Host Interface Board

- ▶ 4 I2C links
- ▶ 1 LVDS link
- ▶ PCI interface



# Agenda

Introduction

Hardware Characteristics

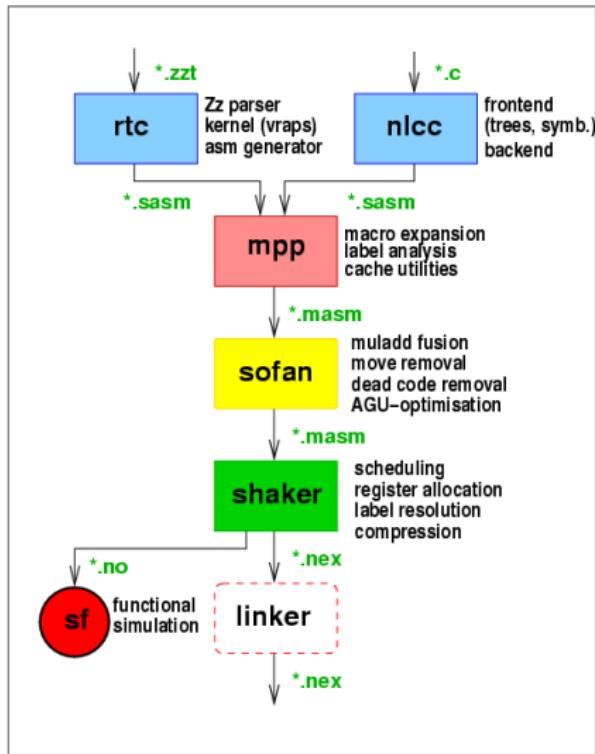
**System Software**

Optimisation

Optimisation Steps

Summary

# Compiler Chain



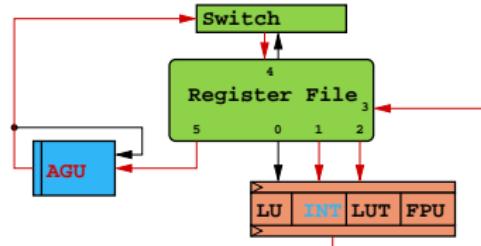
- ▶ Compile high level program:  
`rtc myprog.zzt`
- ▶ Generate low level assembly:  
`mpp -os7 myprog.sasm`
- ▶ Optimise assemble:  
`sofan myprog.masm`
- ▶ Generate microcode (VLIW):  
`npsk +a myprog.masm`

# Programming Languages

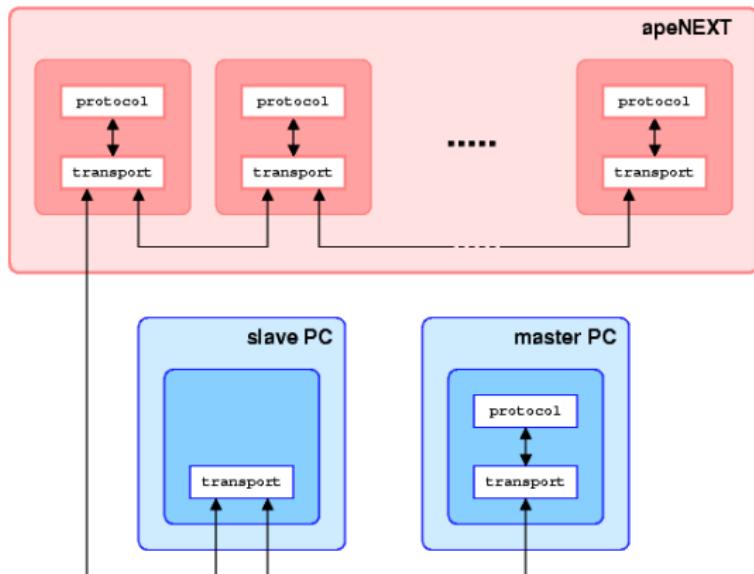
- ▶ TAO
  - ▶ Procedural programming language (like Fortran, C)
  - ▶ Dynamical grammar allows OO-style programming
- ▶ C (experimental)
  - ▶ Based on freely available lcc + custom implementation of libc
  - ▶ Most of ISO C99 standard supported
  - ▶ Few APE-specific language extensions
- ▶ SASM
  - ▶ High level assembler

# Assembly Optimiser: SOFAN

- ▶ Optimisation operating on **low-level assembly**
- ▶ Reconstruction of control- and data-flow graphs
- ▶ Optimisation steps:
  - ▶ merge multiply-add operations
  - ▶ remove dead code
  - ▶ eliminate register moves
  - ▶ optimise address generation:
  - ▶ ...



# System Services



- ▶ Essential services only:  
read/write, fopen/fclose,  
...
- ▶ Parallel I/O transport  
operations:
  - ▶ Global write
  - ▶ Slice write
  - ▶ Broadcast read
  - ▶ Multidata read
- ▶ Multi-stage transport  
→ latencies

# Agenda

Introduction

Hardware Characteristics

System Software

Optimisation

Optimisation Steps

Summary

# Optimisation Strategies

## Optimise for *instruction level parallelism*

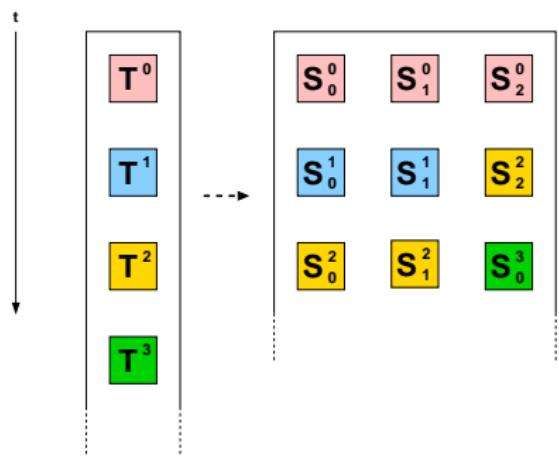
- ▶ Task = set of **independent** subtasks
- ▶ Hardware devices may allow for concurrent execution, e.g.
  - ▶ Arithmetic unit
  - ▶ Address generation unit
  - ▶ Memory controller
  - ▶ Network

## Optimise for *pipelining*

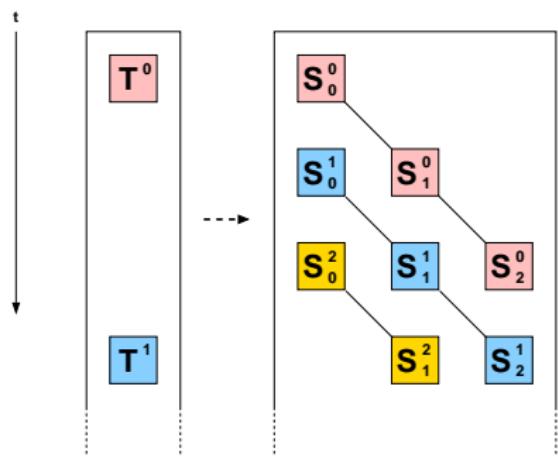
- ▶ Task = set of **dependent** subtasks
- ▶ Hardware device = pipeline
- ▶ Next subtask scheduled before previous completed

# Optimisation Strategies: Cartoons

Instruction level parallelism:



Pipelining:



# Operation Costs

- ▶ Latency  $\lambda$ :  $t(\text{first result}) - t(\text{start first operation})$
- ▶ Throughput  $\tau$  ( $=1/\text{bandwidth}$ ): resource occupation time
- ▶ Distance  $\delta$ : minimal distance between two operations

Operation	$\lambda$	$\tau$	$\delta$
complex mul-add	10	1	-
fp division	107	11	-
integer add	4	1	-
memory load	20	1	4
memory store	10	1	4
jump	$O(10-20)$	-	-

# Performance Analysis

- ▶ Static performance analysis

- ▶ Analysis at compile time
- ▶ Tools:

- ▶ `mem2ncd myprog.mem → myprog.ncd`
- ▶ `nperf myprog.mem`
- ▶ `(nprof)`

- ▶ Run-time performance analysis

- ▶ Gather run-time information from hardware counters:

<code>clock</code>	clock cycles while not in system mode
<code>run</code>	clock cycles while not stalling
<code>qwait</code>	clock cycles while waiting for data in queue

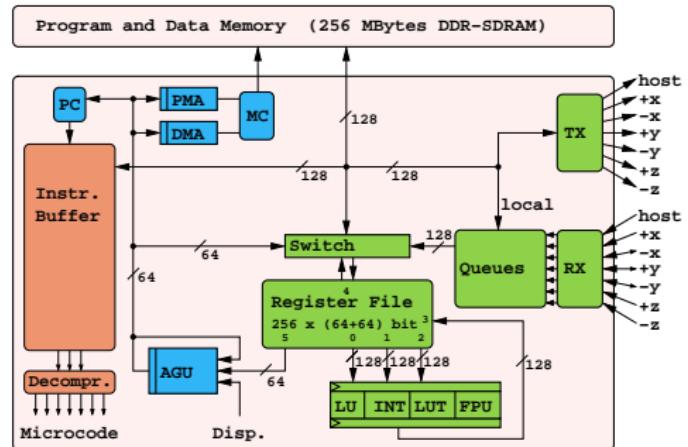
- ▶ TAO syntax: `get clockcnt in clockcnt_var`

- ▶ Example:

```
clockcnt runCounter  
...  
get clockcnt in runCounter
```

# Performance Critical Devices

- ▶ Arithmetic unit
  - ▶ Pipeline filling
- ▶ Memory
  - ▶ Effective bandwidth
- ▶ Register file
  - ▶ Input/output port occupation (P4)
- ▶ Network
  - ▶ Bandwidth



# Agenda

Introduction

Hardware Characteristics

System Software

Optimisation

Optimisation Steps

Summary

# Optimise Memory Access

- ▶ Memory bus is used for
  - ▶ Data load
  - ▶ Data store
  - ▶ Instruction load
- ▶ **Goal:** reduce memory bus occupancy

Possible solutions:

- ▶ Optimise data load/store
  - ▶ Avoid re-loading data (**data re-use**)
  - ▶ Reduce number of memory operations (**burst memory access**)
- ▶ Optimise instruction load
  - ▶ **Compress microcode** (`npsk -z ...`)
  - ▶ **Use instruction cache**

# Instruction Load Optimisation

## ► Example:

```

1 matrix complex su3.[9]
  su3 v[512], w[512]
3
begin cache
5 do i = 0, 511
    v[i] = w[i]
7 enddo
end cache

```

## ► Results for size of loop body (static performance analysis):

Compress	Cache	instructions	words
no	no	82	82
no	yes	61	61
yes	no	82	36
yes	yes	61	23

## ► Caveat:

- Only gain if there is a resource conflict
- Need time to pre-fetch instructions

# Keep Data in Physical Registers

- ▶ **virtual registers:** lifetime ends at control flow instruction
- ▶ **physical registers:** lifetime ends at subroutine call/end

```

constant N = 512
2 complex a, x[N], y[N], z[N]
4
6 do i = 0, N-1
    z[i] = a * x[i] + y[i]
8 enddo

```

```

constant N = 512
2 complex a, x[N], y[N], z[N]
physreg complex ra
4
ra = a      !! load a only once
6 do i = 0, N-1
    z[i] = ra * x[i] + y[i]
8 enddo

```

- ▶ Change of loop body size: 90 → 82 instructions
  - ▶ Remark: We assume the use of Sofan
- ▶ We are still far from optimal:
  - ▶ Memory bandwidth: 16 byte per cycle
  - ▶ FPU throughput: 1 complex add per cycle

# Burst Memory Access

```

constant N = 512
2 complex x[N], y[N]
4
6 do i = 1, N-1
    x[i] = y[i]
8 enddo

```

```

constant N = 512
2 constant B = 64
matrix complex bcomplex.[B]
4 bcomplex x[N/B], y[N/B]
6 do i = 1, (N/B)-1
    x[i] = y[i]
8 enddo

```

- ▶ Number of memory operations reduced  $64 \times$
- ▶ Change of loop body size:  $69 \rightarrow 192$  instructions  
 $(192 \ll 64 \cdot 69 = 4416)$

# Burst Memory Access (2)

- ▶ Alternative implementation:

```
constant N = 512
2 constant B = 64
matrix complex bcomplex.[B]
4 complex x[N], y[N]
register bcomplex rx, ry
6
do i = 1, N-1, B
8   extract rx from x[i]
     ry = rx
10  replace ry into y[i]
enddo
```

- ▶ Advantages:

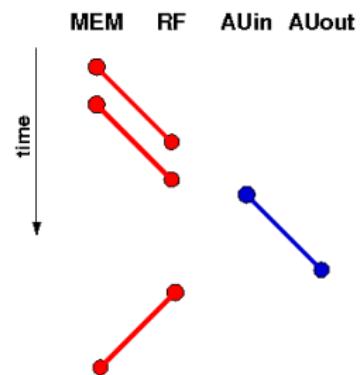
- ▶ Use “natural” objects outside application kernels
- ▶ Possible to use different burst length

- ▶ Keyword `register` defines *virtual registers*

# Pipeline Optimisation

## ► Example:

```
1 constant N = 512
  complex x[N], y[N], z[N]
3
begin cache
5 do i = 0, N-1
    x[i] = y[i] + z[i]
7 enddo
end cache
```



- Size of loop body: 49 instructions

# Pipeline Optimisation (2)

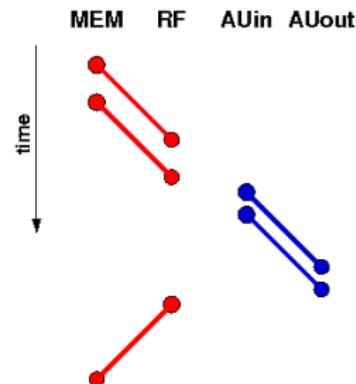
## ► Optimised Example:

```

1 constant N = 512
2 constant B = 2
3 matrix complex bcomplex.[B]
4 bcomplex x[N/B], y[N/B], z[N/B]

6 begin cache
7 do i = 0, N/B-1
8   x[i] = y[i] + z[i]
9 enddo
10 end cache

```



- Size of loop body: 56 instructions ( $\ll 2 \cdot 49$ )
- Expect optimal choice  $B \simeq \lambda_{\text{complex add}}$
- Double gain:
  - Improved filling of FPU pipeline
  - Burst memory access

# Latency Hiding by Data Prefetch

- ▶ Aim: use time needed to load data for other operations
- ▶ Example:

$$\phi_{\mu,i}^a = U_i^{ab} \psi_{\mu,i}^b$$

```

1 /include <qcd>
2
3 constant N = 512
4 su3          u[N]
5 color_spin  psi[N], phi[N]
6
7 begin cache
8 do i = 0, N-1
9   phi[i] = u[i] * psi[i]
10 enddo
11 end cache

```

```

1 ...
2 su3                      u[N]
3 physreg su3              ru
4 color_spin               psi[N], phi[N]
5 physreg color_spin      rpsi
6 register color_spin     res
7
8 rpsi = psi[0]            !! load(0)
9 ru   = u[0]
10 begin cache
11 do i = 0, N-1
12   res    = ru * rpsi    !! compute(i)
13   rpsi   = psi[i+1]     !! load(i+1)
14   ru    = u[i+1]
15   phi[i] = res         !! store(i)
16   enddo
17 end cache

```

# Explicit Use of Prefetch Queues

- ▶ Data load via queue:
  - ▶ MTQ: Load data from memory to queue
  - ▶ QTR: Load data from queue to register file
- ▶ Processor will stall at QTR if queue is empty
- ▶ TAO instructions:

```
1 begin prefetch      !! start prefetch section
  queue = x           !! load word from memory to queue
3 rx = queue          !! load word from queue to RF
  end prefetch        !! end prefetch section
```

- ▶ **begin/end pre-fetch** enforces local memory access for all implicit load operations and changes scheduling strategy
- ▶ Correct order of pre-fetch/fetch is **user responsibility!**

# Latency Hiding: Using the Queue

```
...
2 begin prefetch
queue = psi[0]          !! prefetch(0)
4 queue = u[0]

6 begin cache
do i = 0, N-1
8   rpsi    = queue      !! load(i)
   ru      = queue
10  queue   = psi[i+1]    !! prefetch(i+1)
   queue   = u[i+1]
12  phi[i] = ru * rpsi   !! compute(i)
enddo
14 end cache

16 rpsi    = queue      !! load(N)
   ru      = queue
18 end prefetch
```

# Latency Hiding: Results

- ▶ Results:

	Number of instructions
Original version:	113
With pre-fetch:	89

- ▶ Attention to pre-fetch during last loop iteration:

- ▶ Must be **suppressed or flushed**
- ▶ May need **extended array boundaries**

# Example: Optimise (Local) Vector Product

- ▶ Consider operation:  $c = \sum_i |z_i|^2$
- ▶ Naive implementation:

```
constant N = 512
2 complex c, z[N]

4 c = complex_0

6 do i = 0, N-1
    c = c + z[i] * z[i]~
8 enddo
```

# First Attempt

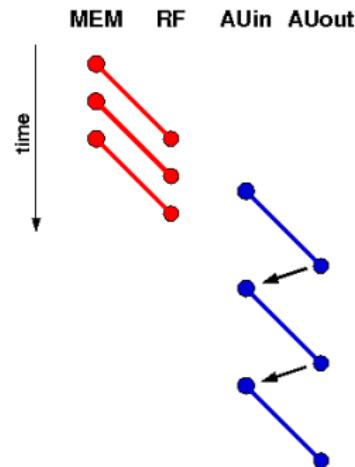
```
constant N = 512
2 constant B = 16
matrix complex bcomplex.[B]
4
complex          c, z[N]
6 register bcomplex  rz
physreg complex   rc
8
rc = complex_0
10 begin cache
do i = 0, N-1, B
12 extract rz from z[i]
 /for k = 0 to B-1 {
14     rc = rc + rz.[k] * rz.[k]^
}
16 enddo
end cache
18 c = rc
```

# First Attempt (2)

- ▶ Change of loop body size:  $93 \rightarrow 179$  instructions  
= 11 instructions per array element
- ▶ Problem **read-after-write** dependency:

```

2   extract rz from z[i]
3   /for k = 0 to B-1 {
4       rc = rc + rz.[k] * rz.[k]^
    }
```



- ▶ **Solution:** Partial sums

# Implementation of Partial Sums

```
...
2 physreg bcomplex    rc
register bcomplex    rz
4
5 /for i = 0 to B-1 {
6   rc.[i] = complex_0
7 }
8 begin cache
9   do i = 0, N-1, B
10    extract rz from z[i]
11    /for k = 0 to B-1 {
12      rc.[k] = rc.[k] + rz.[k] * rz.[k]^
13    }
14 enddo
15 end cache
16 c = rc.[0]
17 /for i = 1 to B-1 {
18   c = c + rc.[i]
19 }
```

# Implementation of Partial Sums (2)

- ▶ Number of instructions in loop body:

	total	per element
Original version	53	53
Optimisation 1	179	11
Partial sums	59	3

- ▶ Optimal number of instructions per element: 1
  - ▶ For larger B it is possible to get close to optimum

# Condition Stack

- ▶ Effect of operations may be conditioned
- ▶ Example:

```

1 integer i, j
2 real s
3
4 s = 0.
5 if (i > j)
6   s = 1.
7 endif

```

```

1 integer i, j
2 real s
3
4 s = 0.
5 where (localint(i) > localint(j))
6   s = 1.
7 endwhere

```

- ▶ Change of main program size: 102 → 66 instructions
- ▶ Performance penalties due to control flow statements:
  - ▶ Jump costs O(20) clock cycles
  - ▶ **Pipeline break**

# Agenda

Introduction

Hardware Characteristics

System Software

Optimisation

Optimisation Steps

Summary

# Relevant Optimisation Tricks

- ▶ Optimise data load/store
  - ▶ Avoid re-loading data
    - ▶ Use physical registers
  - ▶ Perform burst memory access
    - ▶ TAO: extract, replace
  - ▶ Hide latency by data pre-fetch
- ▶ Optimise instruction load
  - ▶ Compress microcode
  - ▶ Use instruction cache
- ▶ Optimise pipeline usage
  - ▶ Loop unrolling (/for)
- ▶ Eliminate instruction dependencies
- ▶ Avoid pipeline breaks