## So you want to buy a supercomputer?

James Davenport
Hebron & Medlock Professor of Information Technology

University of Bath (U.K.)
(visiting Waterloo)

15 May 2009
Many thanks to Prof. Guest (Cardiff)

## University of Bath

- Good (9th out of 117 in the U.K.: Guardian 12 May 2009)
- Heavily co-op
- Strengths in Science, Engineering, Mathematics

## University of Bath

- Good (9th out of 117 in the U.K.: Guardian 12 May 2009)
- Heavily co-op
- Strengths in Science, Engineering, Mathematics

But small — 538 Faculty

- Nationally run (EPSRC etc. $\approx$ NSERC) major supercomputers

- Nationally run (EPSRC etc. $\approx$ NSERC) major supercomputers
- HECToR (current one) 29th in TOP 500

## U.K. scene — generalities

- Nationally run (EPSRC etc. $\approx$ NSERC) major supercomputers
- HECToR (current one) 29th in TOP 500
- Time bid for on competitive grants (virtual money)

- Nationally run (EPSRC etc. $\approx$ NSERC) major supercomputers
- HECToR (current one) 29th in TOP 500
- Time bid for on competitive grants (virtual money)
- Hence you need a 'track record'

## U.K. scene — generalities

- Nationally run (EPSRC etc. $\approx$ NSERC) major supercomputers
- HECToR (current one) 29th in TOP 500
- Time bid for on competitive grants (virtual money)
- Hence you need a 'track record'

- Nationally run (EPSRC etc. $\approx$ NSERC) major supercomputers
- HECToR (current one) 29th in TOP 500
- Time bid for on competitive grants (virtual money)
- Hence you need a 'track record'

Basically, Mark 4 v 25:

# U.K. scene — generalities

- Nationally run (EPSRC etc. $\approx$ NSERC) major supercomputers
- HECToR (current one) 29th in TOP 500
- Time bid for on competitive grants (virtual money)
- Hence you need a 'track record'

Basically, Mark 4 v 25: "to him that hath shall be given".

- EPSRC etc. ($\approx$ NSERC) now allow depreciation on computing resources to be charged to grants
- (Previously, you had to *buy* your own machine

## U.K. scene — recent developments

- EPSRC etc. ($\approx$ NSERC) now allow depreciation on computing resources to be charged to grants
- (Previously, you had to *buy* your own machine and *run* it)

## U.K. scene — recent developments

- EPSRC etc. ($\approx$ NSERC) now allow depreciation on computing resources to be charged to grants
- (Previously, you had to *buy* your own machine and *run* it)
- Government announce Science Research Infrastructure Fund (£500M/year)
- (largely buildings, but equipment not excluded)

## U.K. scene — recent developments

- EPSRC etc. ($\approx$ NSERC) now allow depreciation on computing resources to be charged to grants
- (Previously, you had to *buy* your own machine and *run* it)
- Government announce Science Research Infrastructure Fund (£500M/year)
- (largely buildings, but equipment not excluded)
- Bath share about £5M/year

N.B. "year" = H.M. Treasury Year

## U.K. scene — recent developments

- EPSRC etc. ($\approx$ NSERC) now allow depreciation on computing resources to be charged to grants
- (Previously, you had to *buy* your own machine and *run* it)
- Government announce Science Research Infrastructure Fund (£500M/year)
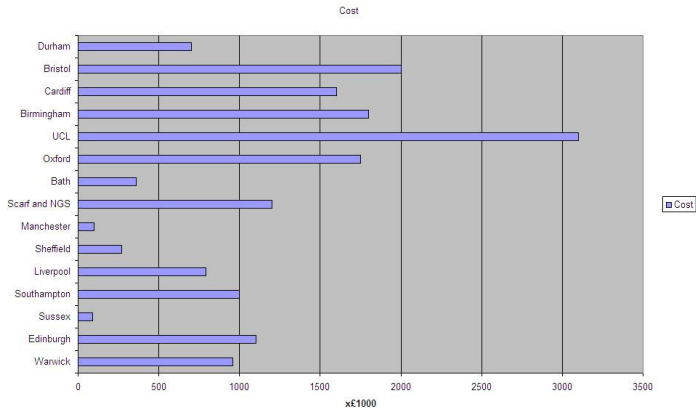- (largely buildings, but equipment not excluded)
- Bath share about £5M/year

N.B. "year" = H.M. Treasury Year

Brainwave: if I purchase a supercomputer, then I can depreciate it, and have money to buy a new one.

# Recent UK spend, excluding machine rooms etc.

Cardiff £1.6M on machine, £1.4M on converting machine room and (high-quality) air conditioning.

## Machine Rooms — a major problem

Cardiff £1.6M on machine, £1.4M on converting machine room and (high-quality) air conditioning.

Bristol £2M on machine, £2M+ on building machine room and including chilled water.

## Machine Rooms — a major problem

Cardiff £1.6M on machine, £1.4M on converting machine room and (high-quality) air conditioning.

Bristol £2M on machine, £2M+ on building machine room and including chilled water.

Imperial (Central London) £3M on $CO_2$-cooled machine room.

## Machine Rooms — a major problem

Cardiff £1.6M on machine, £1.4M on converting machine room and (high-quality) air conditioning.

Bristol £2M on machine, £2M+ on building machine room and including chilled water.

Imperial (Central London) £3M on $CO_2$-cooled machine room.

## Machine Rooms — a major problem

Cardiff £1.6M on machine, £1.4M on converting machine room and (high-quality) air conditioning.

Bristol £2M on machine, £2M+ on building machine room and including chilled water.

Imperial (Central London) £3M on $CO_2$-cooled machine room.

Bath had an old machine room from the 1970s.

+ I doubt very much Bath would have spent those sort of sums on a new machine room

## Old Machine Rooms — a mixed blessing

+ I doubt very much Bath would have spent those sort of sums on a new machine room
+ Comparative speed: I took under a year from initial decision to Phase 1 installed

## Old Machine Rooms — a mixed blessing

+ I doubt very much Bath would have spent those sort of sums on a new machine room
+ Comparative speed: I took under a year from initial decision to Phase 1 installed
− It will, just about, cope with the current smallish machine: I think in a few years we'll *need* a new machine room

## Old Machine Rooms — a mixed blessing

+ I doubt very much Bath would have spent those sort of sums on a new machine room
+ Comparative speed: I took under a year from initial decision to Phase 1 installed
− It will, just about, cope with the current smallish machine: I think in a few years we'll *need* a new machine room
− The University don't realise what a bargain they're getting

## Old Machine Rooms — a mixed blessing

+ I doubt very much Bath would have spent those sort of sums on a new machine room
+ Comparative speed: I took under a year from initial decision to Phase 1 installed
− It will, just about, cope with the current smallish machine: I think in a few years we'll *need* a new machine room
− The University don't realise what a bargain they're getting
− Despite the Estates Department's promises, the power supply *did* need upgrading

## Old Machine Rooms — a mixed blessing

+ I doubt very much Bath would have spent those sort of sums on a new machine room
+ Comparative speed: I took under a year from initial decision to Phase 1 installed
− It will, just about, cope with the current smallish machine: I think in a few years we'll *need* a new machine room
− The University don't realise what a bargain they're getting
− Despite the Estates Department's promises, the power supply *did* need upgrading
+ Contracts signed this week on a new machine room with chilled water!

1/2007 I am tasked with looking into this

1/2007 I am tasked with looking into this
5/2007 Top management buys the case

1/2007 I am tasked with looking into this
5/2007 Top management buys the case

1/2007 I am tasked with looking into this

5/2007 Top management buys the case

So what was the case?

- Researchers think they can support £450K of equipment

1/2007 I am tasked with looking into this

5/2007 Top management buys the case

So what was the case?

- Researchers think they can support £450K of equipment
- (i.e. earn that much depreciation over 3 years)

## Actual Timescale

1/2007 I am tasked with looking into this

5/2007 Top management buys the case

So what was the case?

- Researchers think they can support £450K of equipment
- (i.e. earn that much depreciation over 3 years)
- 6 year commitment with 2-year reviews/refreshes

1/2007 I am tasked with looking into this

5/2007 Top management buys the case

So what was the case?

- Researchers think they can support £450K of equipment
- (i.e. earn that much depreciation over 3 years)
- 6 year commitment with 2-year reviews/refreshes

So 4 years warning of decommitment

1/2007  I am tasked with looking into this

## Actual Timescale

1/2007 I am tasked with looking into this

5/2007 Top management buys the case: RFP for £360K

1/2007 I am tasked with looking into this

5/2007 Top management buys the case: RFP for £360K

    \* There was already a national pre-qualified list

## Actual Timescale

1/2007 I am tasked with looking into this

5/2007 Top management buys the case: RFP for £360K

    * There was already a national pre-qualified list

9/2007 "So what's your final offer?"

## Actual Timescale

1/2007 I am tasked with looking into this

5/2007 Top management buys the case: RFP for £360K

* There was already a national pre-qualified list

9/2007 "So what's your final offer?"

10/2007 Purchase decision

## Actual Timescale

1/2007 I am tasked with looking into this

5/2007 Top management buys the case: RFP for £360K

&ast; There was already a national pre-qualified list

9/2007 "So what's your final offer?"

10/2007 Purchase decision

1/2008 Phase 1 delivery

## Actual Timescale

1/2007 I am tasked with looking into this

5/2007 Top management buys the case: RFP for £360K

  * There was already a national pre-qualified list

9/2007 "So what's your final offer?"

10/2007 Purchase decision

1/2008 Phase 1 delivery

3/2008 Phase 1 acceptance

## Actual Timescale

1/2007 I am tasked with looking into this

5/2007 Top management buys the case: RFP for £360K

* There was already a national pre-qualified list

9/2007 "So what's your final offer?"

10/2007 Purchase decision

1/2008 Phase 1 delivery

3/2008 Phase 1 acceptance

- UK Treasury FY ends 5 April!

## Actual Timescale

1/2007 I am tasked with looking into this

5/2007 Top management buys the case: RFP for £360K

* There was already a national pre-qualified list

9/2007 "So what's your final offer?"

10/2007 Purchase decision

1/2008 Phase 1 delivery

3/2008 Phase 1 acceptance

- UK Treasury FY ends 5 April!

10/2008 Phase 2 decision (*not* to delay)

## Actual Timescale

1/2007 I am tasked with looking into this

5/2007 Top management buys the case: RFP for £360K

* There was already a national pre-qualified list

9/2007 "So what's your final offer?"

10/2007 Purchase decision

1/2008 Phase 1 delivery

3/2008 Phase 1 acceptance

- UK Treasury FY ends 5 April!

10/2008 Phase 2 decision (*not* to delay)

1/2009 Phase 2 delivery

## Actual Timescale

1/2007 I am tasked with looking into this

5/2007 Top management buys the case: RFP for £360K

* There was already a national pre-qualified list

9/2007 "So what's your final offer?"

10/2007 Purchase decision

1/2008 Phase 1 delivery

3/2008 Phase 1 acceptance

- UK Treasury FY ends 5 April!

10/2008 Phase 2 decision (*not* to delay)

1/2009 Phase 2 delivery

5/2009 Acceptance

## Equipment Purchased

Clustervision: a UK/Dutch firm of system integrators: the boards are Supermicro.

Clustervision: a UK/Dutch firm of system integrators: the boards are Supermicro.

- 100 nodes; $2 \times$ 4-core 2.8GHz Intel Harpertown

Clustervision: a UK/Dutch firm of system integrators: the boards are Supermicro.

- 100 nodes; $2 \times 4$-core 2.8GHz Intel Harpertown
- (3.0 gave less power/£; 2.66 pushed the power envelope)

Clustervision: a UK/Dutch firm of system integrators: the boards are Supermicro.

- 100 nodes; $2 \times 4$-core 2.8GHz Intel Harpertown
- (3.0 gave less power/£; 2.66 pushed the power envelope)
- 2 nodes/power supply

Clustervision: a UK/Dutch firm of system integrators: the boards are Supermicro.

- 100 nodes; $2 \times$ 4-core 2.8GHz Intel Harpertown
- (3.0 gave less power/£; 2.66 pushed the power envelope)
- 2 nodes/power supply
- 2GB/core main memory

## Equipment Purchased

Clustervision: a UK/Dutch firm of system integrators: the boards are Supermicro.

- 100 nodes; $2 \times 4$-core 2.8GHz Intel Harpertown
- (3.0 gave less power/£; 2.66 pushed the power envelope)
- 2 nodes/power supply
- 2GB/core main memory
- * Specified this way as 2/4 core wasn't obvious

## Equipment Purchased

Clustervision: a UK/Dutch firm of system integrators: the boards are Supermicro.

- 100 nodes; $2 \times 4$-core 2.8GHz Intel Harpertown
- (3.0 gave less power/£; 2.66 pushed the power envelope)
- 2 nodes/power supply
- 2GB/core main memory
- \* Specified this way as 2/4 core wasn't obvious
- $=$ 1.6TB main memory — it adds up!

## Equipment Purchased

Clustervision: a UK/Dutch firm of system integrators: the boards are Supermicro.

- 100 nodes; $2 \times 4$-core 2.8GHz Intel Harpertown
- (3.0 gave less power/£; 2.66 pushed the power envelope)
- 2 nodes/power supply
- 2GB/core main memory
- \* Specified this way as 2/4 core wasn't obvious
- = 1.6TB main memory — it adds up!
- Double Data Rate Infiniband

1. Phase 1: Linpack benchmark

1. Phase 1: Linpack benchmark
   - We had linear algebra compiled for the previous chip!

1. Phase 1: Linpack benchmark
   - We had linear algebra compiled for the previous chip!
2. Phase 2: a range of tests related to major users

1. Phase 1: Linpack benchmark
   - We had linear algebra compiled for the previous chip!
2. Phase 2: a range of tests related to major users
 * Very grateful to Prof. Guest for organising

## Acceptance Tests

1. Phase 1: Linpack benchmark
   - We had linear algebra compiled for the previous chip!
2. Phase 2: a range of tests related to major users
* Very grateful to Prof. Guest for organising
   - MPI defaults were badly wrong

1. Phase 1: Linpack benchmark
   - We had linear algebra compiled for the previous chip!
2. Phase 2: a range of tests related to major users
*. Very grateful to Prof. Guest for organising
   - MPI defaults were badly wrong
   - DDR Infiniband was running out of steam faster than expected

## Acceptance Tests

1. Phase 1: Linpack benchmark
   - We had linear algebra compiled for the previous chip!
2. Phase 2: a range of tests related to major users
* Very grateful to Prof. Guest for organising
   - MPI defaults were badly wrong
   - DDR Infiniband was running out of steam faster than expected
   - Several partial failures.

# Partial Failures

## Partial Failures

*Very* frustrating and hard to diagnose: typically one job would take "longer than expected".

## Partial Failures

*Very* frustrating and hard to diagnose: typically one job would take "longer than expected".

- Observe this is happening, and feel very confused

## Partial Failures

*Very* frustrating and hard to diagnose: typically one job would take "longer than expected".

- Observe this is happening, and feel very confused
- Eventually spot that it happens when node 78 is used!

*Very* frustrating and hard to diagnose: typically one job would take "longer than expected".

- Observe this is happening, and feel very confused
- Eventually spot that it happens when node 78 is used!
- Convince the manufacturer to run their tests on node 78

## Partial Failures

*Very* frustrating and hard to diagnose: typically one job would take "longer than expected".

- Observe this is happening, and feel very confused
- Eventually spot that it happens when node 78 is used!
- Convince the manufacturer to run their tests on node 78

## Partial Failures

*Very* frustrating and hard to diagnose: typically one job would take "longer than expected".

- Observe this is happening, and feel very confused
- Eventually spot that it happens when node 78 is used!
- Convince the manufacturer to run their tests on node 78

Failure modes

## Partial Failures

*Very* frustrating and hard to diagnose: typically one job would take "longer than expected".

- Observe this is happening, and feel very confused
- Eventually spot that it happens when node 78 is used!
- Convince the manufacturer to run their tests on node 78

Failure modes

1. Node 78 (and another one since) — poor Infiniband

# Partial Failures

*Very* frustrating and hard to diagnose: typically one job would take "longer than expected".

- Observe this is happening, and feel very confused
- Eventually spot that it happens when node 78 is used!
- Convince the manufacturer to run their tests on node 78

Failure modes

1. Node 78 (and another one since) — poor Infiniband
2. twice so far: a node loses 4GB of memory on a reboot

## Partial Failures

*Very* frustrating and hard to diagnose: typically one job would take "longer than expected".

- Observe this is happening, and feel very confused
- Eventually spot that it happens when node 78 is used!
- Convince the manufacturer to run their tests on node 78

Failure modes

1. Node 78 (and another one since) — poor Infiniband
2. twice so far: a node loses 4GB of memory on a reboot
3. Others?

## Partial Failures

*Very* frustrating and hard to diagnose: typically one job would take "longer than expected".

- Observe this is happening, and feel very confused
- Eventually spot that it happens when node 78 is used!
- Convince the manufacturer to run their tests on node 78

Failure modes

1. Node 78 (and another one since) — poor Infiniband
2. twice so far: a node loses 4GB of memory on a reboot
3. Others?

## Partial Failures

*Very* frustrating and hard to diagnose: typically one job would take "longer than expected".

- Observe this is happening, and feel very confused
- Eventually spot that it happens when node 78 is used!
- Convince the manufacturer to run their tests on node 78

Failure modes

1. Node 78 (and another one since) — poor Infiniband
2. twice so far: a node loses 4GB of memory on a reboot
3. Others?

"One footsore soldier can delay a regiment" — Duke of Wellington

- Get it in writing from Estates.

- Get it in writing from Estates.
- Know your (potential) users early

## Lessons I already knew

- Get it in writing from Estates.
- Know your (potential) users early
- (devise acceptance tests accordingly)

## Lessons I already knew

- Get it in writing from Estates.
- Know your (potential) users early
- (devise acceptance tests accordingly)
- It's hard to explain to management

- It's *very* hard to explain to management

- It's *very* hard to explain to management
- Acceptance tests are very important, especially

- It's *very* hard to explain to management
- Acceptance tests are very important, especially
- Car-Parrinello Molecular Dynamics (CPMD) for interconnect

- It's *very* hard to explain to management
- Acceptance tests are very important, especially
- Car-Parrinello Molecular Dynamics (CPMD) for interconnect
- Partial failure is far worse than total failure

- It's *very* hard to explain to management
- Acceptance tests are very important, especially
- Car-Parrinello Molecular Dynamics (CPMD) for interconnect
- Partial failure is far worse than total failure
- Even DDR Infiniband has trouble with 8 cores/node

## Lessons I know now

- It's *very* hard to explain to management
- Acceptance tests are very important, especially
- Car-Parrinello Molecular Dynamics (CPMD) for interconnect
- Partial failure is far worse than total failure
- Even DDR Infiniband has trouble with 8 cores/node
  (There's a good paper (*now*!) by HP)

- Good ways of detecting partial failure

- Good ways of detecting partial failure
- How to manage software licencing if you can't afford to licence every node

- Good ways of detecting partial failure
- How to manage software licencing if you can't afford to licence every node
- How to persuade management to deliver on the promised refreshes

- Good ways of detecting partial failure
- How to manage software licencing if you can't afford to licence every node
- How to persuade management to deliver on the promised refreshes
- Will the assumptions hold up:

- Good ways of detecting partial failure
- How to manage software licencing if you can't afford to licence every node
- How to persuade management to deliver on the promised refreshes
- Will the assumptions hold up:
  - Assumptions on grant-getting

## Lessons I know I still don't know

- Good ways of detecting partial failure
- How to manage software licencing if you can't afford to licence every node
- How to persuade management to deliver on the promised refreshes
- Will the assumptions hold up:
  - Assumptions on grant-getting
  - Assumptions on actual usage $\Rightarrow$ price/hour

- With the exception of a "short test" queue, allocation is based on whole nodes.

## Price per node hour: 52p≈CAN$0.9

- With the exception of a "short test" queue, allocation is based on whole nodes.
- Allocation is based on entitlements rather than retrospective billing

## Price per node hour: 52p≈CAN\$0.9

- With the exception of a "short test" queue, allocation is based on whole nodes.
- Allocation is based on entitlements rather than retrospective billing
- The Maui scheduler has (too?) many knobs in this area

## Price per node hour: 52p≈CAN$0.9

- With the exception of a "short test" queue, allocation is based on whole nodes.
- Allocation is based on entitlements rather than retrospective billing
- The Maui scheduler has (too?) many knobs in this area

    48% Equipment depreciation

## Price per node hour: 52p≈CAN$0.9

- With the exception of a "short test" queue, allocation is based on whole nodes.
- Allocation is based on entitlements rather than retrospective billing
- The Maui scheduler has (too?) many knobs in this area

    48% Equipment depreciation
    15% Equipment maintenance

## Price per node hour: 52p≈CAN\$0.9

- With the exception of a "short test" queue, allocation is based on whole nodes.
- Allocation is based on entitlements rather than retrospective billing
- The Maui scheduler has (too?) many knobs in this area

  48% Equipment depreciation

  15% Equipment maintenance

  10% Machine electricity

## Price per node hour: 52p≈CAN$0.9

- With the exception of a "short test" queue, allocation is based on whole nodes.
- Allocation is based on entitlements rather than retrospective billing
- The Maui scheduler has (too?) many knobs in this area

  - 48% Equipment depreciation
  - 15% Equipment maintenance
  - 10% Machine electricity
  - 8% Air conditioning (incl. depreciation)

## Price per node hour: 52p≈CAN$0.9

- With the exception of a "short test" queue, allocation is based on whole nodes.
- Allocation is based on entitlements rather than retrospective billing
- The Maui scheduler has (too?) many knobs in this area

  48% Equipment depreciation

  15% Equipment maintenance

  10% Machine electricity

  8% Air conditioning (incl. depreciation)

  17% 1 Programmer (1/3 of team of 3)

## Price per node hour: 52p≈CAN$0.9

- With the exception of a "short test" queue, allocation is based on whole nodes.
- Allocation is based on entitlements rather than retrospective billing
- The Maui scheduler has (too?) many knobs in this area

  48% Equipment depreciation
  15% Equipment maintenance
  10% Machine electricity
  8% Air conditioning (incl. depreciation)
  17% 1 Programmer (1/3 of team of 3)
  2% My time

# Lessons I don't know I don't know?

Any questions?