

# Scaling the Peak: Maximizing floating point performance on the Epiphany NoC

Anish Varghese, Gaurav Mitra, Robert Edwards and Alistair Rendell

Research School of Computer Science  
The Australian National University

May 30, 2015

# Outline

- 1 Introduction
- 2 Programming the Epiphany
- 3 Performance Experiments
- 4 Implementing High Performance Applications
- 5 Energy Measurement
- 6 Summary

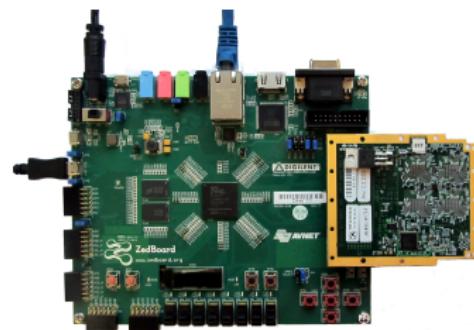
# Energy efficient HPC at ANU

Parallella-16



NVIDIA Jetson TK1 - ARM-GPU SoC

Parallella-64 Prototype



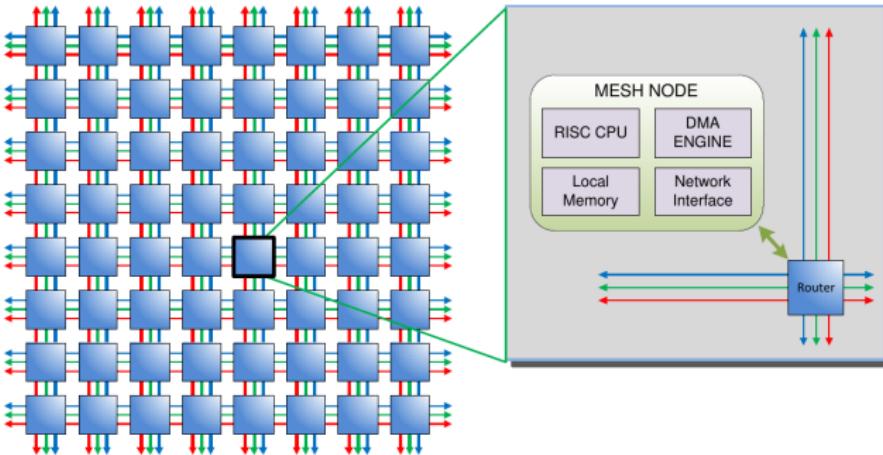
TI Keystone II - ARM-DSP SoC



# Outline

- 1 Introduction
- 2 Programming the Epiphany
  - Programming Considerations
- 3 Performance Experiments
- 4 Implementing High Performance Applications
- 5 Energy Measurement
- 6 Summary

# Programming Considerations



- Store code and data in different banks
- Interleave FMADD with load/stores
- Unroll inner loops

# Outline

- 1 Introduction
- 2 Programming the Epiphany
- 3 Performance Experiments
  - Experiment Platform
  - On-chip Communication
  - Off-chip Communication
- 4 Implementing High Performance Applications
- 5 Energy Measurement
- 6 Summary

# Experiment Platform

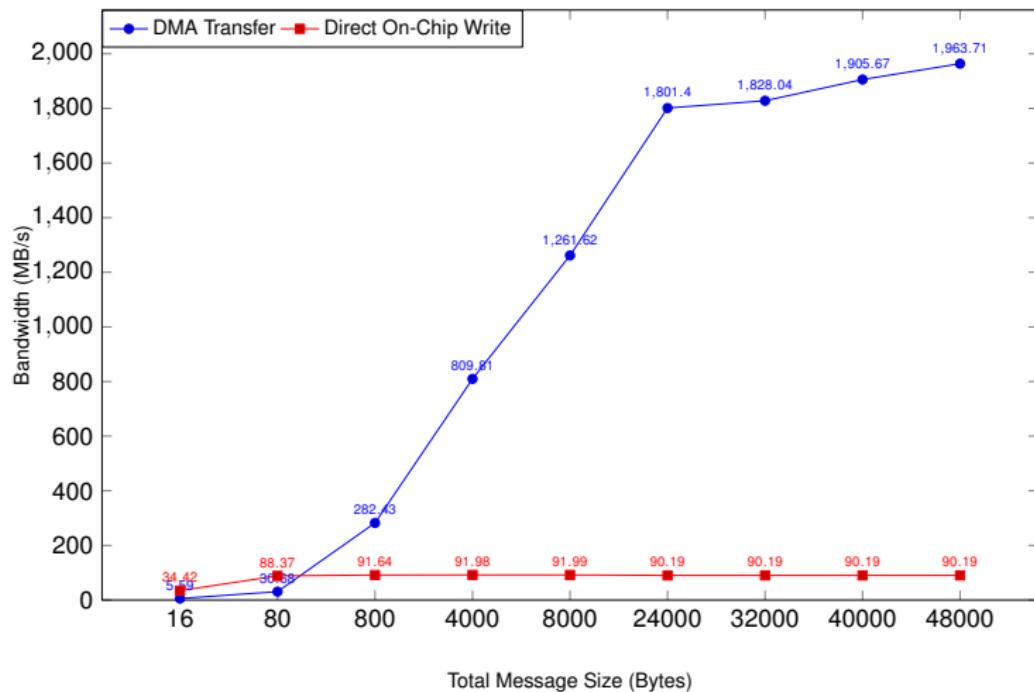
## Parallella-64 Prototype

- ZedBoard evaluation module with Zynq SoC
- Daughter card with Epiphany-IV 64-core (E64G401)
- Dual core ARM Cortex-A9 host at 667 MHz
- Epiphany eCores at 600 MHz
- 512 MB of DDR3 RAM on the host
- 32 MB shared with eCores

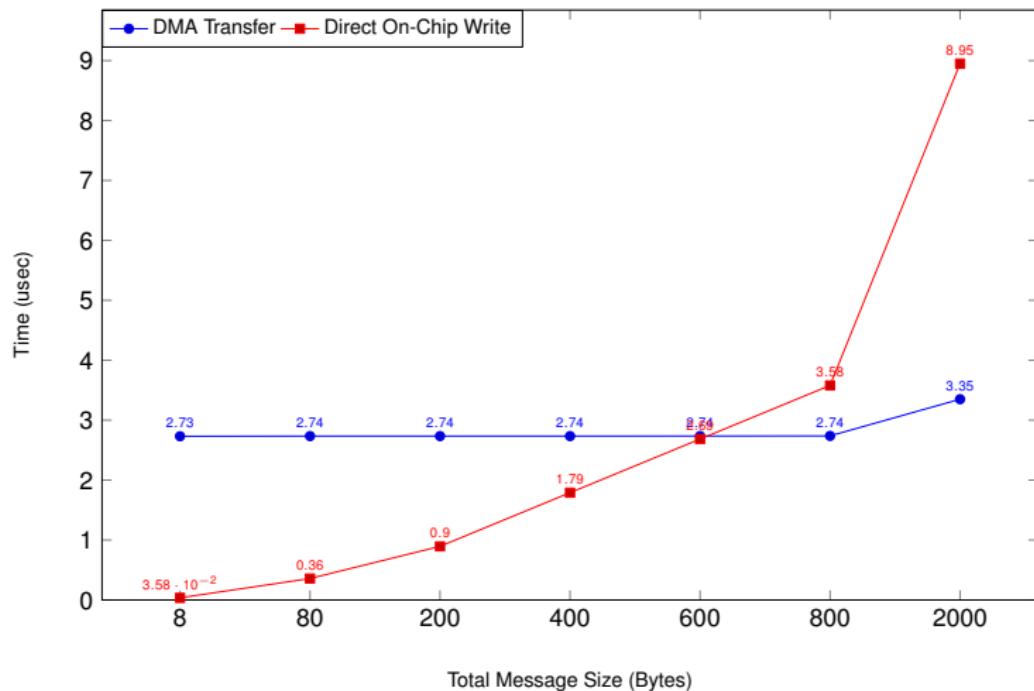
# Outline

- 1 Introduction
- 2 Programming the Epiphany
- 3 Performance Experiments
  - Experiment Platform
  - On-chip Communication
  - Off-chip Communication
- 4 Implementing High Performance Applications
- 5 Energy Measurement
- 6 Summary

# DMA vs Direct Write Bandwidth



# Message Transfer Latency

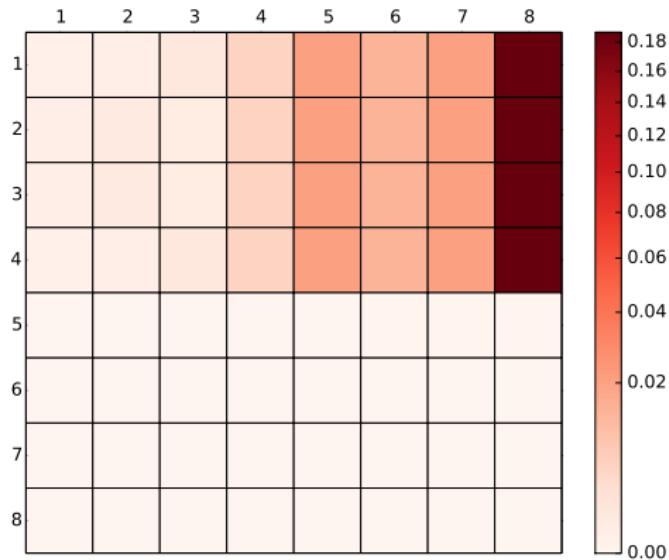


- Latency of  $\approx 7$  cycles per transfer

# Outline

- 1 Introduction
- 2 Programming the Epiphany
- 3 Performance Experiments
  - Experiment Platform
  - On-chip Communication
  - Off-chip Communication
- 4 Implementing High Performance Applications
- 5 Energy Measurement
- 6 Summary

# Core-wise utilization of external memory link



- 150 MB/sec using Direct Writes
- 400 MB/sec using DMA

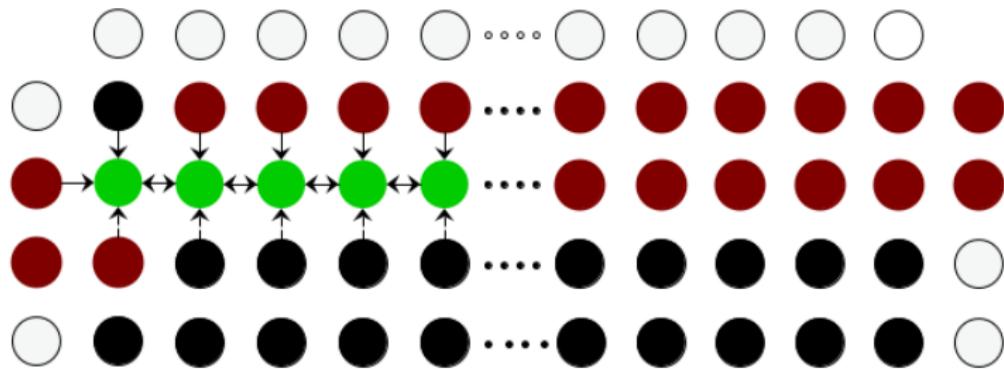
# Outline

- 1 Introduction
- 2 Programming the Epiphany
- 3 Performance Experiments
- 4 Implementing High Performance Applications
  - Heat Stencil
  - Matrix Multiplication
- 5 Energy Measurement
- 6 Summary

# Five-point star-shaped heat stencil

$$\begin{aligned} T_{new_{i,j}} = & w_1 * T_{prev_{i,j+1}} + w_2 * T_{prev_{i,j}} \\ & + w_3 * T_{prev_{i,j-1}} + w_4 * T_{prev_{i+1,j}} \\ & + w_5 * T_{prev_{i-1,j}} \end{aligned}$$

# Single-Core Stencil - 20 columns



# Diving into Assembly

```
1 .Lb:  
2   dogrid1 26,27,28,29,30,47,48,49,50,51,52,53, 1, 2, 3, 4, 5, 6, 7, 8, 9,10  
3   ldr r20,[r0,#0 + stride]  
4   dogrid0 31,32,33,34,35,52,53,54,55,56,57,58, 6, 7, 8, 9,10,11,12,13,14,15  
5   ldr r41,[r1,#0]  
6   dogrid1 36,37,38,39,40,57,58,59,60,61,62,63,11,12,13,14,15,16,17,18,19,20  
7   str r21,[r1],#stride  
8 ; 2nd row  
9   dogrid0 43,44,45,46,47,20,21,22,23,24,25,26,16,17,18,19,20,1+stride,2+stride,3+stride  
    ,4+stride,5+stride  
10  add r0,r0,#(stride * 4)  
11  dogrid1 48,49,50,51,52,25,26,27,28,29,30,31, 1, 2, 3, 4, 5, 6, 7, 8, 9,10  
12  ldr r42,[r0,#stride]  
13  dogrid0 53,54,55,56,57,30,31,32,33,34,35,36, 6, 7, 8, 9,10,11,12,13,14,15  
14  ldr r63,[r1,#0]  
15  dogrid1 58,59,60,61,62,35,36,37,38,39,40,41,11,12,13,14,15,16,17,18,19,20  
16  str r43,[r1],#stride  
17 ; 1st row  
18  dogrid0 21,22,23,24,25,42,43,44,45,46,47,48,16,17,18,19,20,1+stride,2+stride,3+stride  
    ,4+stride,5+stride  
19  add r0,r0,#(stride * 4)  
20  sub r2,r2,#1  
21  nop  
22  bne .Lb
```

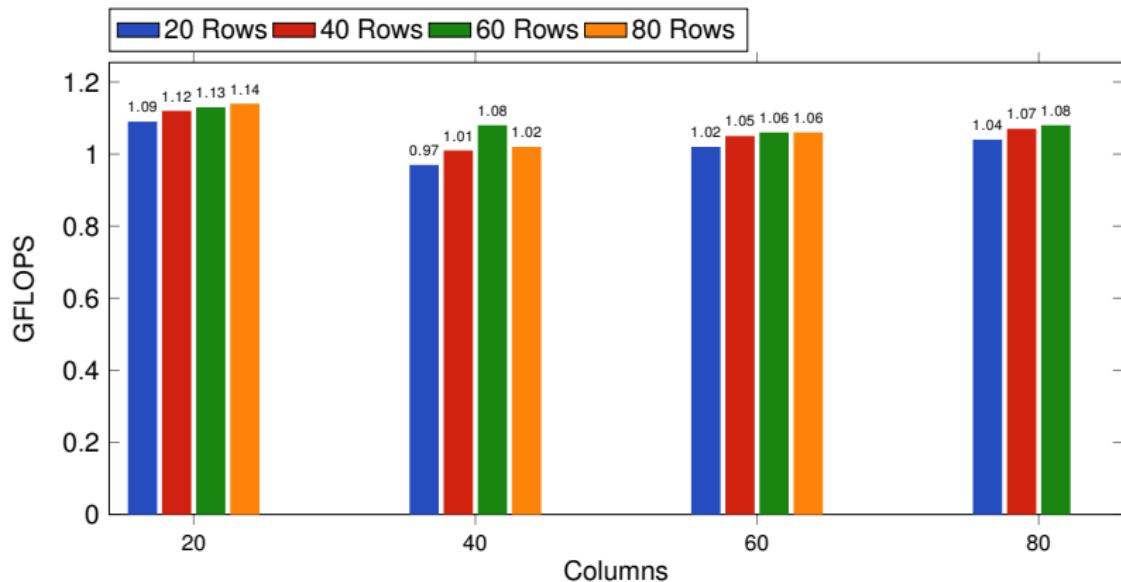
# Writing Macros

```
.macro dogrid0 a0,a1,a2,a3,a4,b0,b1,b2,b3,b4,b5,b6,o0,o1,o2,o3,o4,i0,i1,i2,i3,i4
2  fmadd r15,r\@a0,r3
   str r8,[r0,#\o0]
4  fmadd r16,r\@a1,r3
   str r9,[r0,#\o1]
6  fmadd r17,r\@a2,r3
   str r10,[r0,#\o2]
8   fmadd r18,r\@a3,r3
   str r11,[r0,#\o3]
10  fmadd r19,r\@a4,r3
    str r14,[r0,#\o4]
12  fmadd r15,r\b0,r4
    fmadd r16,r\b1,r4
14  fmadd r17,r\b2,r4
    fmadd r18,r\b3,r4
16  fmadd r19,r\b4,r4
    fmadd r15,r\b1,r5
18 .
20 .
```

Interleaved pair

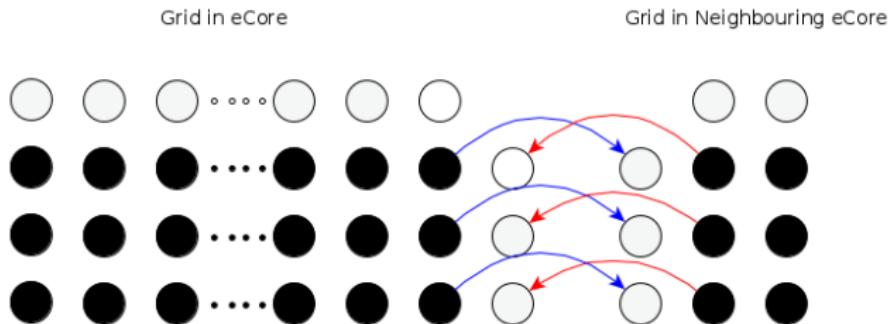
Code available at [https://github.com/parallella/parallella-examples/tree/master/heat\\_stencil](https://github.com/parallella/parallella-examples/tree/master/heat_stencil)

# Floating point performance - Single-Core



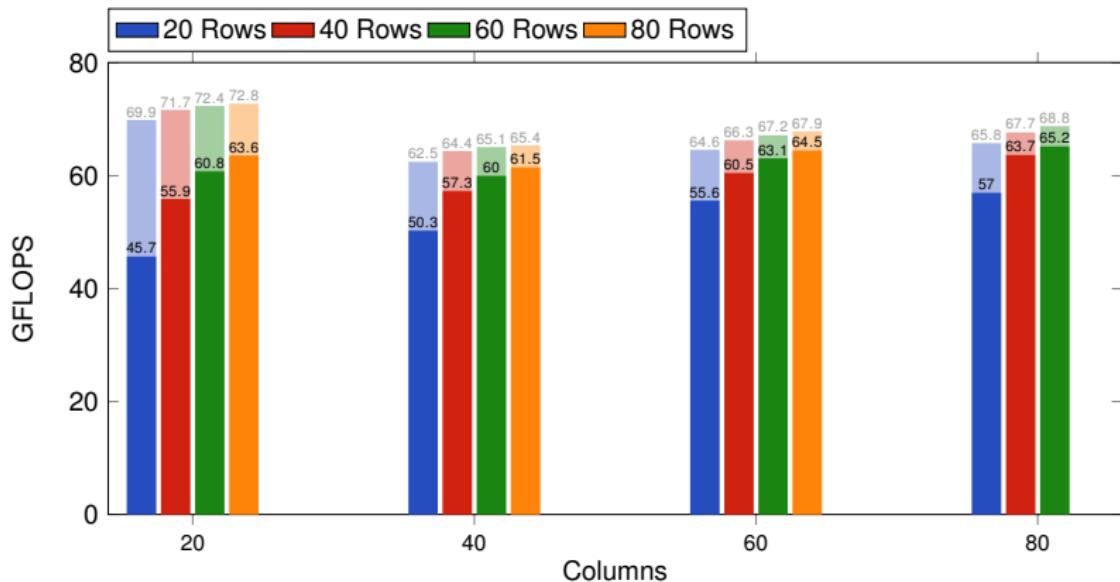
- Stencil evaluated for 50 iterations.
- 81-95% of peak performance

# Multi-Core Stencil



- Synchronization between neighbouring eCores
- DMA for boundary transfers

# Floating point performance - Multi-Core



- Lighter colors show performance without communication
- 65.2 GFLOPS: 85% of peak performance with communication

# Outline

- 1 Introduction
- 2 Programming the Epiphany
- 3 Performance Experiments
- 4 Implementing High Performance Applications
  - Heat Stencil
  - Matrix Multiplication
- 5 Energy Measurement
- 6 Summary

# Single-Core Matmul

$$C_{11} = A_{11}B_{11} + A_{12}B_{21} + A_{13}B_{31} + \dots$$

$$C_{12} = A_{11}B_{12} + A_{12}B_{22} + A_{13}B_{32} + \dots$$

⋮

$$C_{21} = A_{21}B_{11} + A_{22}B_{21} + A_{23}B_{31} + \dots$$

$$C_{22} = A_{21}B_{12} + A_{22}B_{22} + A_{23}B_{32} + \dots$$

- Operand matrices in different memory banks
- Stack moved to bottom half of bank 1
- Hand-tuned assembly code

# Diving into Assembly (again)

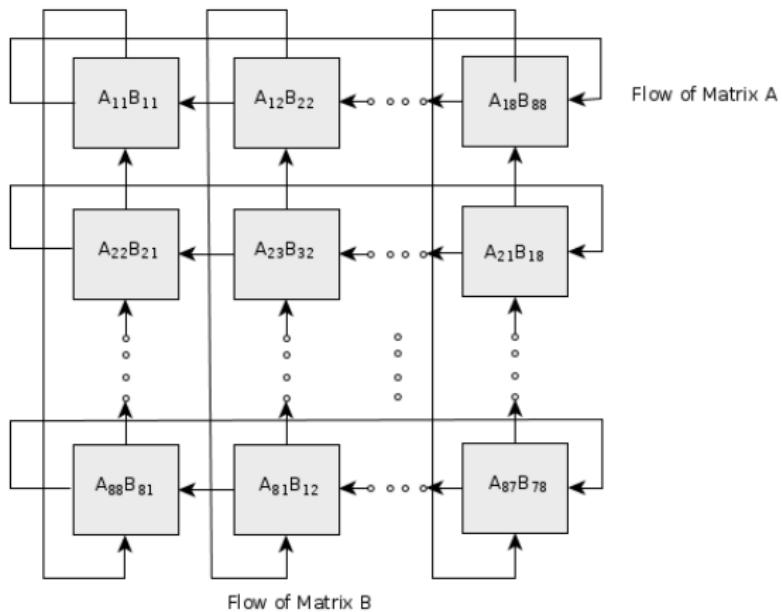
```
.macro doMult areg, index,  
2   fmadd r32,r\areg,r16  
4   ldrd r22,[r1,#\index + 3]  
    fmadd r33,r\areg,r17  
6   ldr r\aprev,[r0,#\inc]  
    fmadd r34,r\areg,r18  
8   fmadd r35,r\areg,r19  
    fmadd r36,r\areg,r20  
10  fmadd r37,r\areg,r21  
    ldrd r16,[r1,#\index + 4]  
12  fmadd r38,r\areg,r22  
    ldrd r18,[r1,#\index + 5]  
14  fmadd r39,r\areg,r23  
    ldrd r20,[r1,#\index + 6]  
16  fmadd r40,r\areg,r16  
    ldrd r22,[r1,#\index + 7]  Interleaved pair  
18 .  
19 :  
20 :
```

Code available at [https://github.com/parallellica/parallellica-examples/tree/master/matmul\\_optimized](https://github.com/parallellica/parallellica-examples/tree/master/matmul_optimized)

# Floating point performance - Single-Core

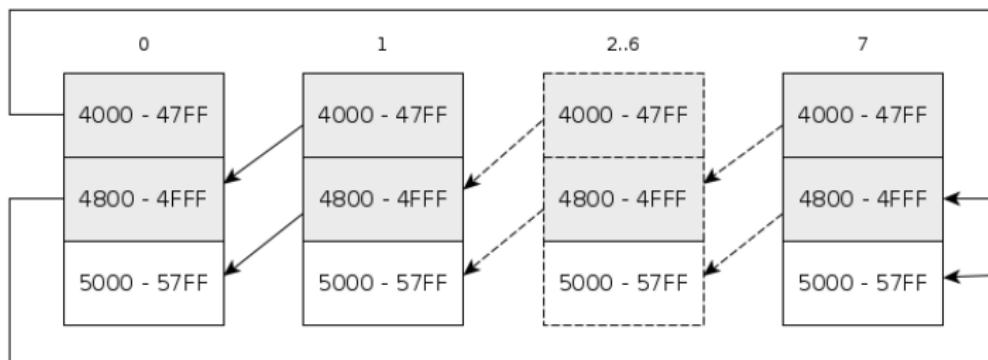
Dimensions	GFLOPS	% Peak
$8 \times 8$	0.85	70.5
$16 \times 16$	1.07	89.5
$20 \times 20$	1.11	92.5
$24 \times 24$	1.12	93.4
$32 \times 32$	1.15	95.9

# Multi-Core Matmul - Cannon's Algorithm



# Buffering Strategy

## Transfer of Matrix A and Matrix B



# Floating Point Performance - Multi-Core

Matrix C (per-core)	Num of eCores			
	4 × 4		8 × 8	
	GFLOPS	% Peak	GFLOPS	% Peak
8 × 8	5.1	26	20.3	26
16 × 16	12.8	67	51.4	67
20 × 20	14.4	75	57.6	75
24 × 24	15.4	80	62.2	81
32 × 32	16.3	85	65.3	85

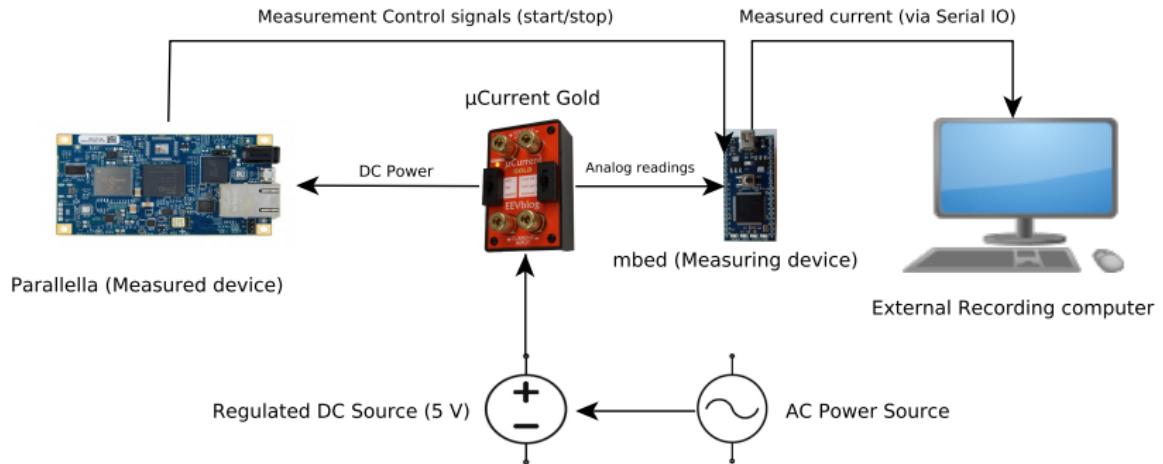
# Floating Point Performance - Off-chip

Matrix C	GFLOPS	% Peak	% Memory Transfers
512 × 512	8.3	10.8 %	87.2 %
1024 × 1024	8.6	11.1 %	86.9 %
1536 × 1536	6.3	8.2 %	89.1 %

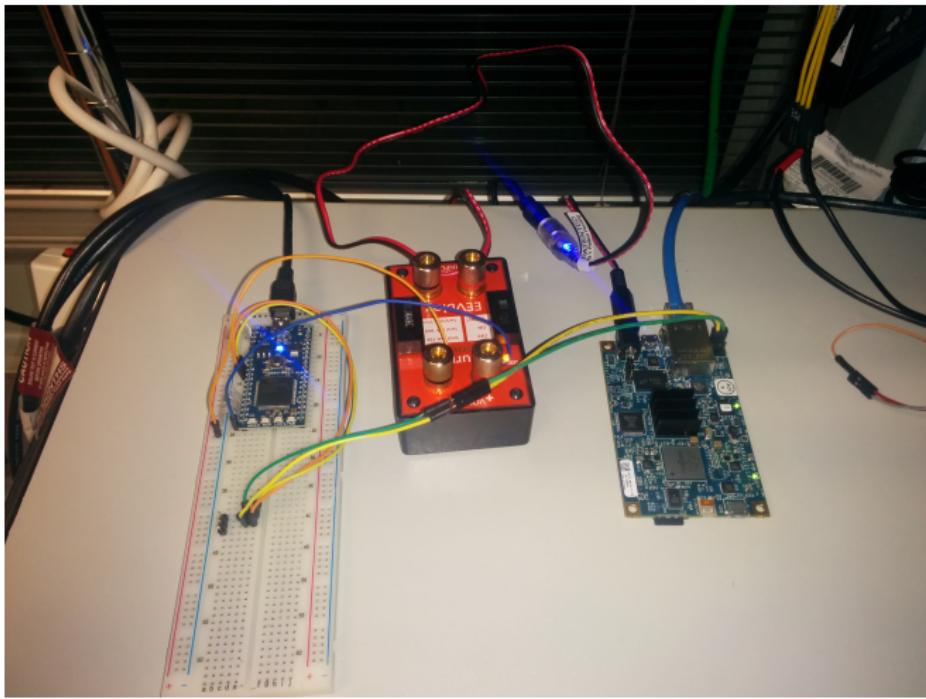
# Outline

- 1 Introduction
- 2 Programming the Epiphany
- 3 Performance Experiments
- 4 Implementing High Performance Applications
- 5 Energy Measurement
  - Setup
  - API
  - Measurements
- 6 Summary

# Measurement Setup



# Measurement Setup



# Outline

- 1 Introduction
- 2 Programming the Epiphany
- 3 Performance Experiments
- 4 Implementing High Performance Applications
- 5 Energy Measurement
  - Setup
  - API
  - Measurements
- 6 Summary

# Measurement API

```
2      /* Initialize serial port IO objects from Measured Device */
3      void power_init()
4      {
5          io = new boost::asio::io_service();
6          serial = new boost::asio::serial_port(*io, serial_port);
7          serial->set_option(boost::asio::serial_port_base::baud_rate(baud_rate));
8      }
9
10     /* Send single character 'r' to reset ADC timer and start measurement */
11     void power_start()
12     {
13         boost::asio::write(*serial, boost::asio::buffer("r",1));
14     }
15
16     /* Send single character 's' to stop ADC timer and measurement */
17     void power_stop()
18     {
19         boost::asio::write(*serial, boost::asio::buffer("s",1));
20     }
21
22     /* Perform serial IO object cleanup */
23     void power_cleanup()
24     {
25         delete io;
26         delete serial;
27     }
```

# Outline

- 1 Introduction
- 2 Programming the Epiphany
- 3 Performance Experiments
- 4 Implementing High Performance Applications
- 5 Energy Measurement
  - Setup
  - API
  - Measurements
- 6 Summary

# Energy Efficiency for Matmul - Comparisons

**Parallelia-16**

Dimensions	nJoules/Flop
$128 \times 128$	0.33
$256 \times 256$	1.82
$512 \times 512$	1.98
$1024 \times 1024$	1.83

**NVIDIA Jetson**

Dimensions	nJoules/Flop
$128 \times 128$	9.7
$256 \times 256$	1.06
$512 \times 512$	0.16
$1024 \times 1024$	0.06

- $128 \times 128$  experiment on Parallelia-16 used only on-chip SRAM
- Same numbers can be derived analytically

# Summary & Future Work

- Achieved high-performance
  - 85% of Peak for On-chip
- Introduced a method to measure energy to solution
- High programming effort was required
  - Compiler could be optimized further
- Energy efficiency
  - Observed on-chip efficiency of 330 pJoules/Flop
  - 4096 eCores => 10 pJoules/Flop
- Higher bandwidth and more SRAM per eCore needed

# Questions?