



HASKELL ARRAYS, ACCELERATED USING GPUS

Manuel M. T. Chakravarty
University of New South Wales

JOINT WORK WITH
Gabriele Keller
Sean Lee

General Purpose GPU Programming (GPGPU)



146X
 Interactive visualization of volumetric white matter connectivity¹

36X
 Ionic placement for molecular dynamics simulation on GPU²

19X
 Transcoding HD video stream to H.264 for portable video³

149X
 Financial simulation of LIBOR Model with swaptions⁴

47X
 GLAME@lab: An M-script API for Linear Algebra Operations on GPU⁷

20X
 Ultrasound medical imaging for cancer diagnostics⁸

MODERN GPUS ARE FREELY PROGRAMMABLE

But no function pointers & limited recursion

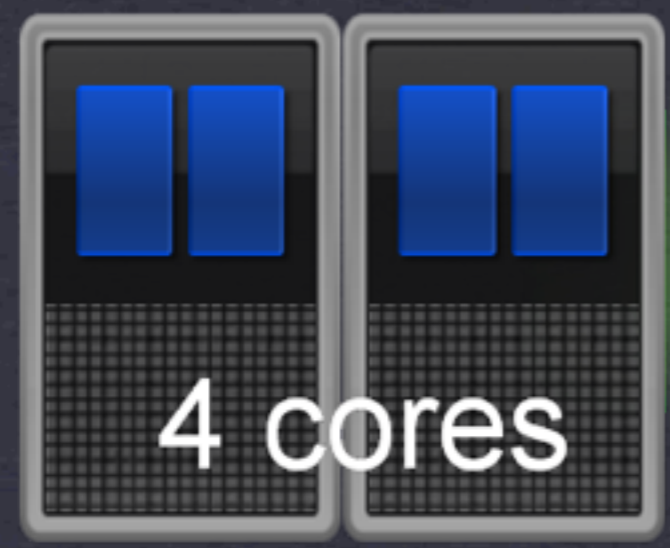
Very Different Programming Model

(Compared to multicore CPUs)



Quadcore
Xeon CPU

Tesla T10
GPU



hundreds of threads/core

MASSIVELY PARALLEL PROGRAMS NEEDED

Tens of thousands of dataparallel threads

Programming GPUs is
hard! Why bother?

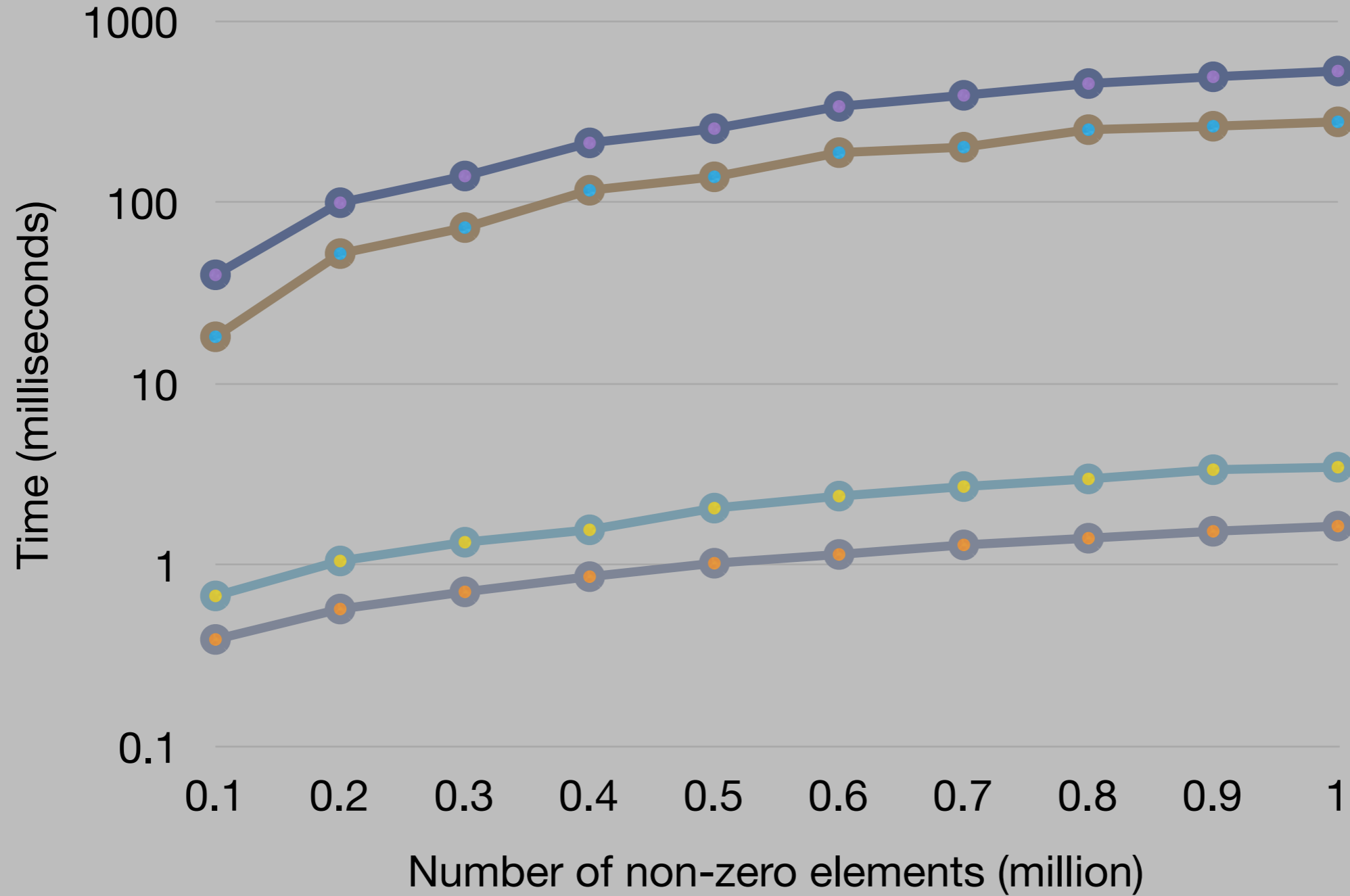




Reduce power consumption!

- * GPU achieves 20x better performance/Watt (judging by peak performance)
- * Speedups between 20x to 150x have been observed in real applications

Sparse Matrix Vector Multiplication



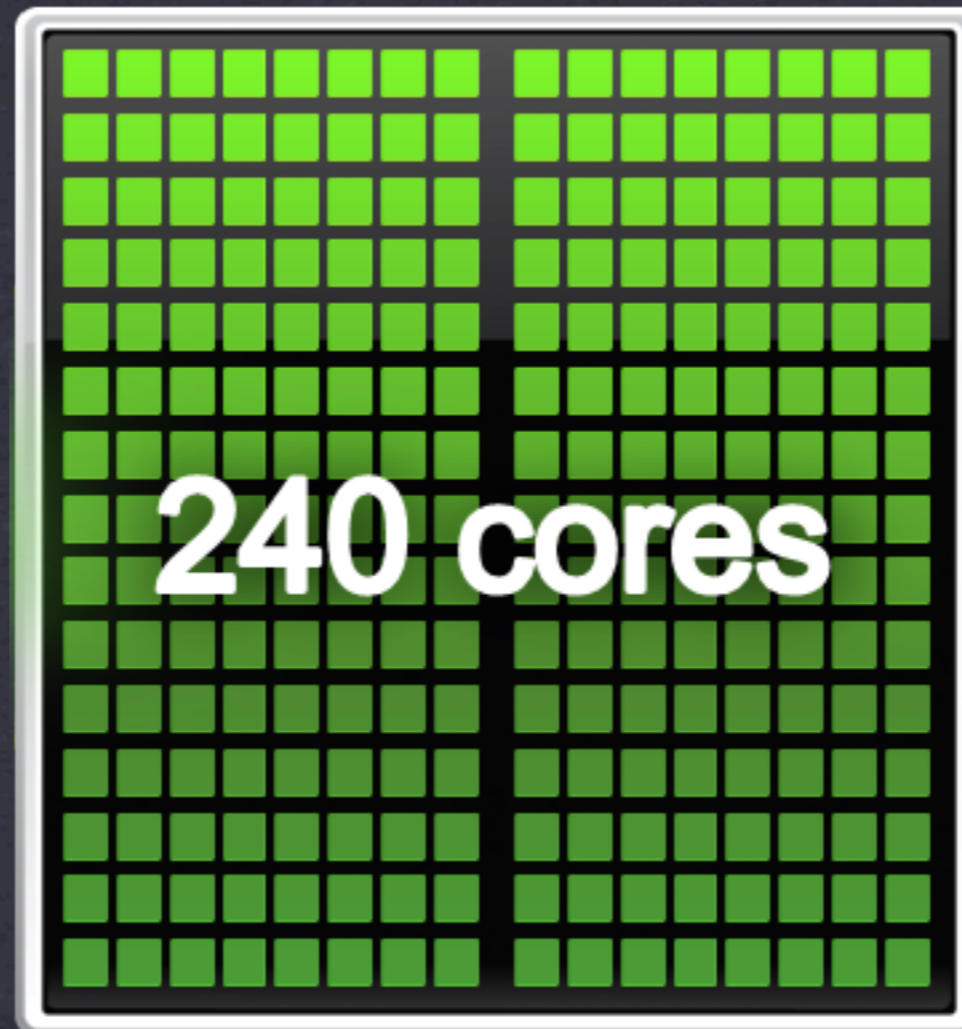
- Plain Haskell, CPU only (AMD Sempron)
- Plain Haskell, CPU only (Intel Xeon)
- Haskell EDSL on a GPU (GeForce 8800GTS)
- Haskell EDSL on a GPU (Tesla S1070 x1)

Prototype of code generator targeting GPUs
Computation only, without CPU ↔ GPU transfer

Challenges

- * Code must be **massively dataparallel**
- * Control structures are limited
 - ▶ No function pointers
 - ▶ Very limited recursion
- * Software-managed cache, memory-access patterns, etc.
- * Portability...



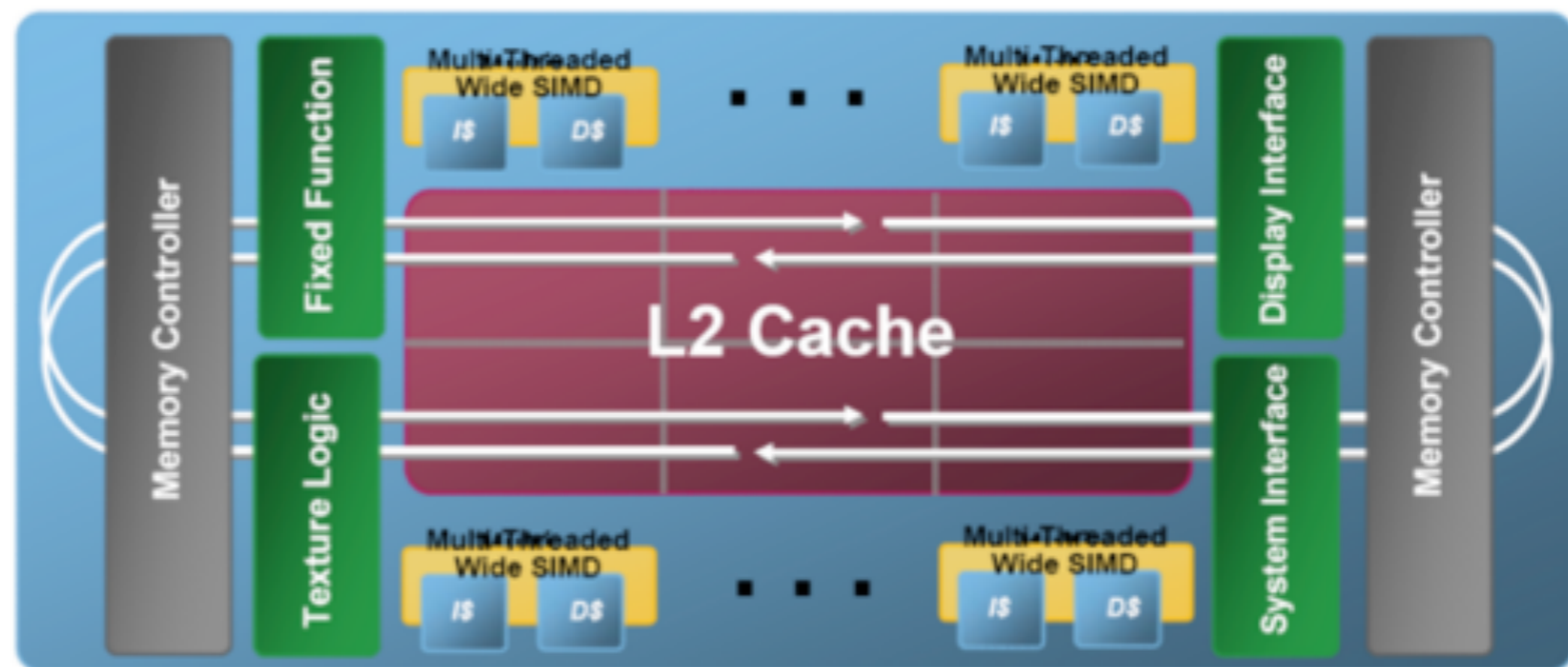


**Tesla T10
GPU**

OTHER COMPUTE ACCELERATOR ARCHITECTURES

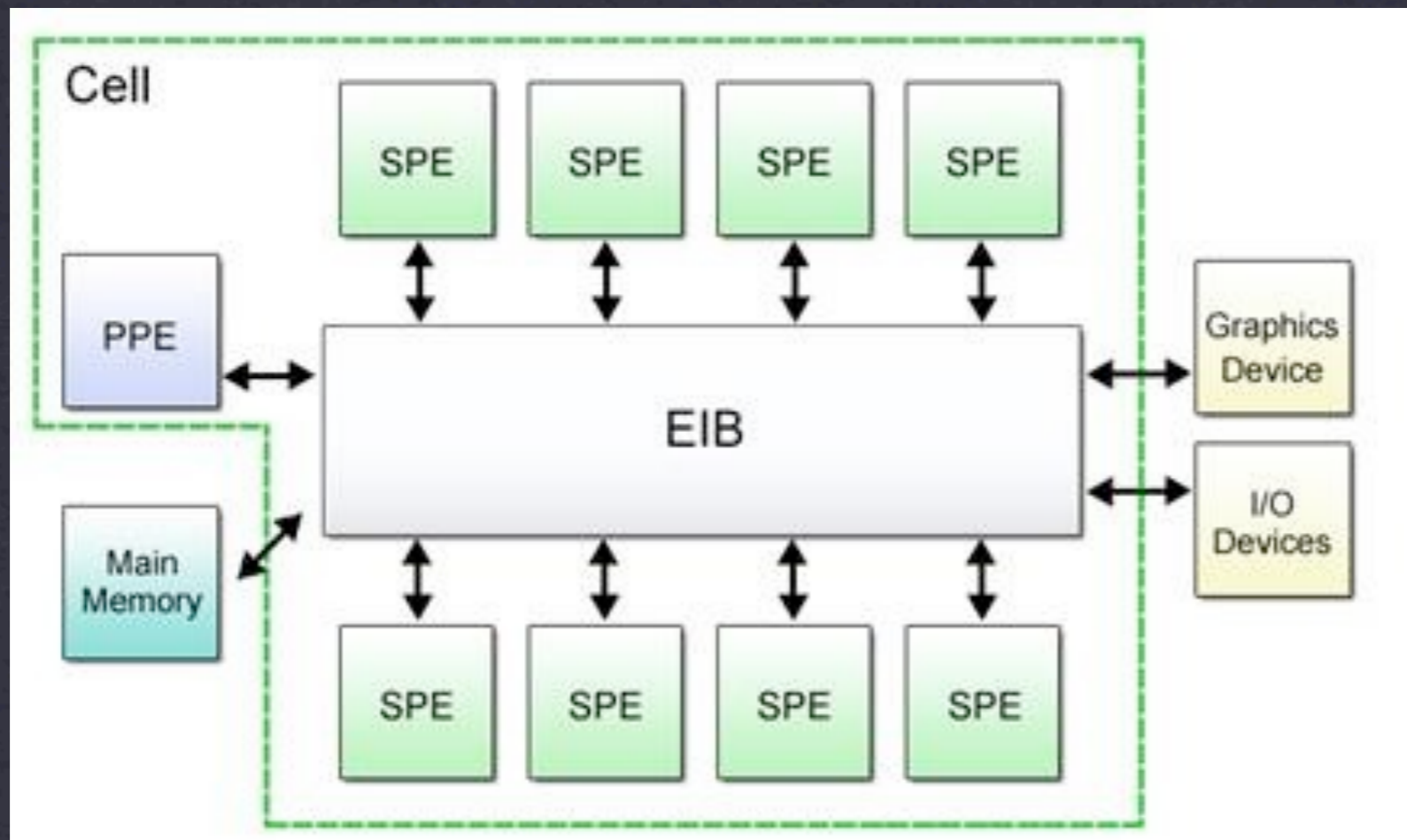
Goal: portable data parallelism

Larrabee Block Diagram



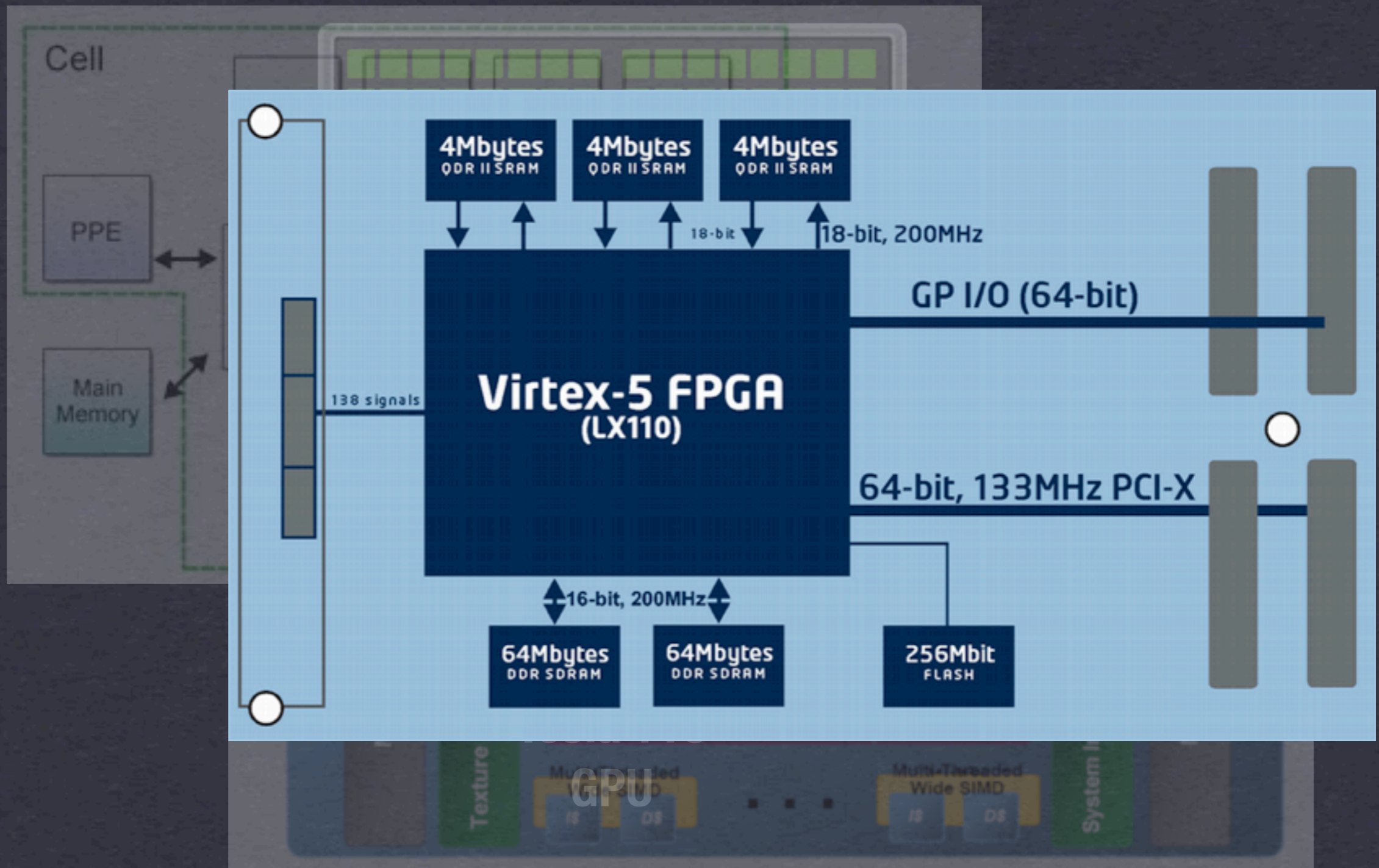
OTHER COMPUTE ACCELERATOR ARCHITECTURES

Goal: portable data parallelism



OTHER COMPUTE ACCELERATOR ARCHITECTURES

Goal: portable data parallelism



OTHER COMPUTE ACCELERATOR ARCHITECTURES

Goal: portable data parallelism

Data.Array.Accelerate

- * Collective operations on multi-dimensional regular arrays
- * Embedded DSL
 - ▶ Restricted control flow
 - ▶ First-order GPU code
- * Generative approach based on combinator templates
- * Multiple backends



Data.Array.Accelerate

- * Collective operations on multi-dimensional regular arrays
- * Embedded DSL
 - ▶ Restricted control flow
 - ▶ First-order GPU code
- * Generative approach based on combinator templates
- * Multiple backends

✓ massive data parallelism



Data.Array.Accelerate

- * Collective operations on multi-dimensional regular arrays
 - ✓ massive data parallelism
- * Embedded DSL
 - ▶ Restricted control flow
 - ✓ limited control structures
 - ▶ First-order GPU code
- * Generative approach based on combinator templates
- * Multiple backends



Data.Array.Accelerate

- * Collective operations on multi-dimensional regular arrays
 - ✓ massive data parallelism
- * Embedded DSL
 - ▶ Restricted control flow
 - ✓ limited control structures
 - ▶ First-order GPU code
 - ✓ hand-tuned access patterns
- * Generative approach based on combinator templates
- * Multiple backends



Data.Array.Accelerate

- * Collective operations on multi-dimensional regular arrays
 - ✓ massive data parallelism
- * Embedded DSL
 - ▶ Restricted control flow
 - ✓ limited control structures
 - ▶ First-order GPU code
 - ✓ hand-tuned access patterns
- * Generative approach based on combinator templates
 - ✓ portability
- * Multiple backends



import Data.Array.Accelerate

Dot product

```
dotp :: Vector Float -> Vector Float
      -> Acc (Scalar Float)
dotp xs ys
  = let
      xs' = use xs
      ys' = use ys
  in
    fold (+) 0 (zipWith (*) xs' ys')
```



import Data.Array.Accelerate

Haskell
array

Dot product

```
dotp :: Vector Float -> Vector Float
      -> Acc (Scalar Float)
dotp xs ys
  = let
      xs' = use xs
      ys' = use ys
  in
    fold (+) 0 (zipWith (*) xs' ys')
```



import Data.Array.Accelerate

Haskell
array

Dot product

```
dotp :: Vector Float -> Vector Float
      -> Acc (Scalar Float)
dotp xs ys
  = let
      xs' = use xs
      ys' = use ys
  in
    fold (+) 0 (zipWith (*) xs' ys')
```

EDSL array =
desc. of array comps



import Data.Array.Accelerate

Dot product

Haskell
array

```
dotp :: Vector Float -> Vector Float
      -> Acc (Scalar Float)
dotp xs ys
  = let
      xs' = use xs
      ys' = use ys
  in
    fold (+) 0 (zipWith (*) xs')
```

EDSL array =
desc. of array comps

Lift Haskell arrays into
EDSL — may trigger
host → device transfer



import Data.Array.Accelerate

Dot product

Haskell
array

```
dotp :: Vector Float -> Vector Float  
      -> Acc (Scalar Float)  
dotp xs ys  
  = let  
      xs' = use xs  
      ys' = use ys  
  in  
  fold (+) 0 (zipWith (*) xs' ys')
```

EDSL array =
desc. of array comps

Lift Haskell arrays into
EDSL — may trigger
host → device transfer

EDSL array
computations



import Data.Array.Accelerate

Sparse-matrix vector multiplication

```
type SparseVector a = Vector (Int, a)
type SparseMatrix a = (Segments, SparseVector a)

smvm :: Acc (SparseMatrix Float)
      -> Acc (Vector Float)
      -> Acc (Vector Float)
smvm (segd, smat) vec
  = let
      (inds, vals) = unzip smat
      vecVals      = backpermute (shape inds)
                              (\i -> inds!i) vec
      products     = zipWith (*) vecVals vals
  in
    foldSeg (+) 0 products segd
```

import Data.Array.Accelerate

$[0, 0, 6.0, 0, 7.0] \approx [(2, 6.0), (4, 7.0)]$

Sparse-matrix vector multiplication

```
type SparseVector a = Vector (Int, a)
type SparseMatrix a = (Segments, SparseVector a)

smvm :: Acc (SparseMatrix Float)
      -> Acc (Vector Float)
      -> Acc (Vector Float)
smvm (segd, smat) vec
  = let
      (inds, vals) = unzip smat
      vecVals      = backpermute (shape inds)
                              (\i -> inds!i) vec
      products     = zipWith (*) vecVals vals
  in
    foldSeg (+) 0 products segd
```

import Data.Array.Accelerate

$[0, 0, 6.0, 0, 7.0] \approx [(2, 6.0), (4, 7.0)]$

Sparse-matrix vector multiplication

```
type SparseVector a = Vector (Int, a)
type SparseMatrix a = (Segments, SparseVector a)

smvm :: Acc (SparseMatrix Float)
      -> Acc (Vector Float)
      -> Acc (Vector Float)
smvm (segd, smat) vec
  = let
      (inds, vals) = unzip smat
      vecVals      = backpermute (shape inds)
                            (\i -> inds!i) vec
      products     = zipWith (*) vecVals vals
  in
    foldSeg (+) 0 products segd
```

$[[10, 20], [], [30]] \approx ([2, 0, 1], [10, 20, 30])$

Architecture of Data.Array.Accelerate



```
map (\x -> x + 1) arr
```

```
map (\x -> x + 1) arr
```

Reify & HOAS -> de Bruijn

```
Map (Lam (Add `PrimApp`  
          (ZeroIdx, Const 1))) arr
```

```
map (\x -> x + 1) arr
```

Reify & HOAS -> de Bruijn

```
Map (Lam (Add `PrimApp`  
          (ZeroIdx, Const 1))) arr
```

Recover sharing
(CSE or Observe)

```
map (\x -> x + 1) arr
```

Reify & HOAS -> de Bruijn

```
Map (Lam (Add `PrimApp`  
         (ZeroIdx, Const 1))) arr
```

Recover sharing
(CSE or Observe)

Optimisation
(Fusion)


```
map (\x -> x + 1) arr
```

Reify & HOAS -> de Bruijn

```
Map (Lam (Add `PrimApp`  
         (ZeroIdx, Const 1))) arr
```

Recover sharing
(CSE or Observe)

Optimisation
(Fusion)

Code generation

```
__global__ void kernel (float *arr, int n)  
{...
```

```
map (\x -> x + 1) arr
```

Reify & HOAS -> de Bruijn

```
Map (Lam (Add `PrimApp`  
        (ZeroIdx, Const 1))) arr
```

Recover sharing
(CSE or Observe)

Optimisation
(Fusion)

Code generation

```
__global__ void kernel (float *arr, int n)  
{...
```

nvcc

0	1	0	
1	0	0	1
0	1	1	
1	1	0	1
1	1		
0	0	0	
	0	1	0
1	0	0	1

```
map (\x -> x + 1) arr
```

Reify & HOAS -> de Bruijn

```
Map (Lam (Add `PrimApp`  
        (ZeroIdx, Const 1))) arr
```

Recover sharing
(CSE or Observe)

Optimisation
(Fusion)

Code generation

```
__global__ void kernel (float *arr, int n)  
{...
```

nvcc

0	1	0	
1	0	0	1
0	1	1	
1	1	0	1
1	1		
0	0	0	
	0	1	0
1	0	0	1



package
plugins

The API of Data.Array.Accelerate

(The main bits)



Array types

`data Array dim e` — regular, multi-dimensional arrays

`type DIM0 = ()`

`type DIM1 = Int`

`type DIM2 = (Int, Int)`

⟨and so on⟩

`type Scalar e = Array DIM0 e`

`type Vector e = Array DIM1 e`

Array types

`data Array dim e` — regular, multi-dimensional arrays

`type DIM0 = ()`

`type DIM1 = Int`

`type DIM2 = (Int, Int)`

⟨and so on⟩

`type Scalar e = Array DIM0 e`

`type Vector e = Array DIM1 e`

EDSL forms

`data Exp e`

— scalar expression form

`data Acc a`

— array expression form

Array types

`data Array dim e` — regular, multi-dimensional arrays

`type DIM0 = ()`

`type DIM1 = Int`

`type DIM2 = (Int, Int)`

⟨and so on⟩

`type Scalar e = Array DIM0 e`

`type Vector e = Array DIM1 e`

EDSL forms

`data Exp e`

— scalar expression form

`data Acc a`

— array expression form

Classes

`class Elem e`

— scalar and tuples types

`class Elem ix => Ix ix`

— unit and integer tuples

Scalar operations

`instance Num (Exp e)` — overloaded arithmetic

`instance Integral (Exp e)`

⟨and so on⟩

`(==*)`, `(/=*)`, `(<*)`, `(<=*)`, — comparisons

`(>*)`, `(>=*)`, `min`, `max`

`(&&*)`, `(||*)`, `not` — logical operators

`(?)` :: Elem t — conditional expression

=> Exp Bool -> (Exp t, Exp t) -> Exp t

`(!)` :: (Ix dim, Elem e) — scalar indexing

=> Acc (Array dim e) -> Exp dim -> Exp e

`shape` :: Ix dim — yield array shape

=> Acc (Array dim e) -> Exp dim

Collective array operations — creation

— use an array from Haskell land

```
use :: (Ix dim, Elem e)
    => Array dim e -> Acc (Array dim e)
```

— create a singleton

```
unit :: Elem e => Exp e -> Acc (Scalar e)
```

— multi-dimensional replication

```
replicate :: (SliceIx slx, Elem e)
           => Exp slx
           -> Acc (Array (Slice slx) e)
           -> Acc (Array (SliceDim slx) e)
```

— Example: `replicate (All, 10, All) twoDimArr`

Collective array operations — slicing

— slice extraction

```
slice :: (SliceIx slx, Elem e)
      => Acc (Array (SliceDim slx) e)
      -> Exp slx
      -> Acc (Array (Slice slx) e)
```

— Example: `slice (5, All, 7) threeDimArr`



Collective array operations — mapping

```
map :: (Ix dim, Elem a, Elem b)
    => (Exp a -> Exp b)
    -> Acc (Array dim a)
    -> Acc (Array dim b)
```

```
zipWith :: (Ix dim, Elem a, Elem b, Elem c)
        => (Exp a -> Exp b -> Exp c)
        -> Acc (Array dim a)
        -> Acc (Array dim b)
        -> Acc (Array dim c)
```



Collective array operations — reductions

fold :: (Ix dim, Elem a)
=> (Exp a -> Exp a -> Exp a) — associative
-> Exp a
-> Acc (Array dim a)
-> Acc (Scalar a)

scan :: Elem a
=> (Exp a -> Exp a -> Exp a) — associative
-> Exp a
-> Acc (Vector a)
-> (Acc (Vector a), Acc (Scalar a))



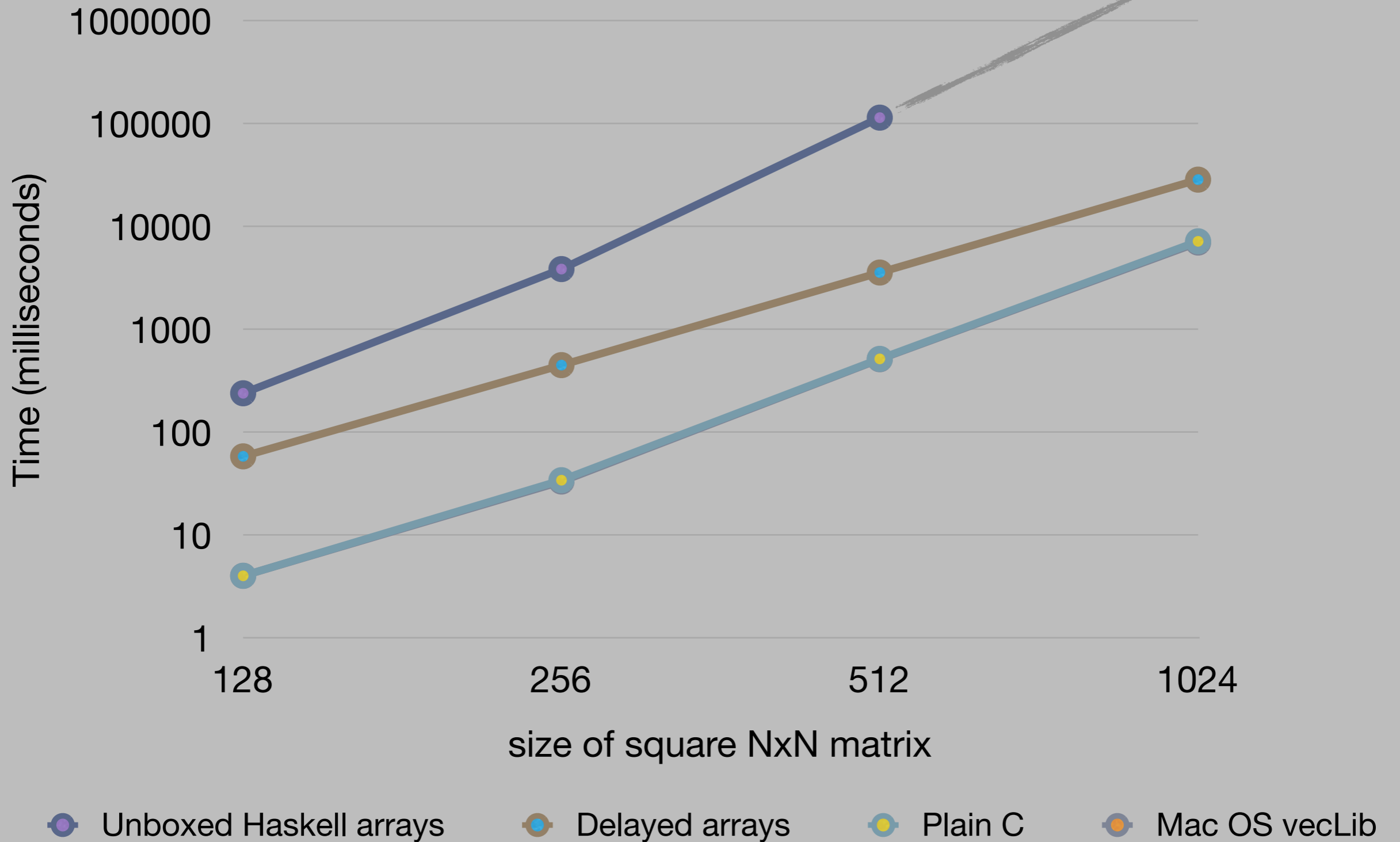
Collective array operations — permutations

```
permute :: (Ix dim, Ix dim', Elem a)
=> (Exp a -> Exp a -> Exp a)
-> Acc (Array dim' a)
-> (Exp dim -> Exp dim')
-> Acc (Array dim a)
-> Acc (Array dim' a)
```

```
backpermute :: (Ix dim, Ix dim', Elem a)
=> Exp dim'
-> (Exp dim' -> Exp dim)
-> Acc (Array dim a)
-> Acc (Array dim' a)
```



Dense Matrix-Matrix Multiplication



Regular arrays in package dph-seq

@ 3.06 GHz Core 2 Duo (with GHC 6.11)

Conclusion

- * EDSL for processing multi-dimensional arrays
- * Collective array operations (highly data parallel)
- * Support for multiple backends
- * Status:
 - ▶ Very early version on Hackage (only interpreter)
<http://hackage.haskell.org/package/accelerate>
 - ▶ Currently porting GPU backend over
 - ▶ Looking for backend contributors

