

Multiscale Dataflow Computing

(Building vertically in a horizontal world)



Oliver Pell
MuCoCoS 2013

Multiscale dataflow computing

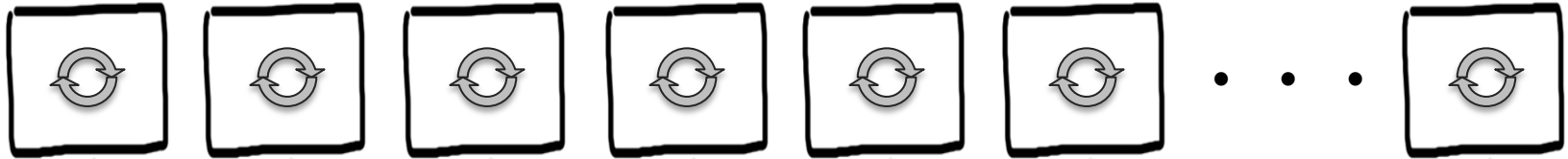
Definiton: “Multiscale”

Problems which have important features at multiple scales

Multiple scales of computing	Important features for optimization
complete system level	⇒ balance compute, storage and IO
parallel node level	⇒ maximize utilization of compute and interconnect
microarchitecture level	⇒ minimize data movement
arithmetic level	⇒ tradeoff range, precision and accuracy = discretize in time, space and value
bit level	⇒ encode and add redundancy
transistor level	=> create the illusion of ‘0’ and ‘1’

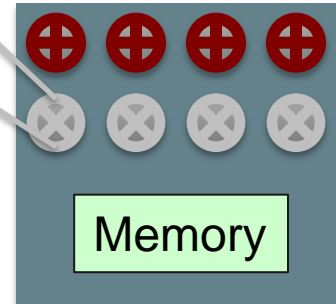
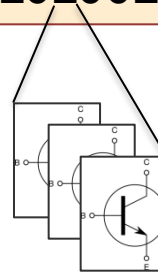
A heterogeneous system

CPU



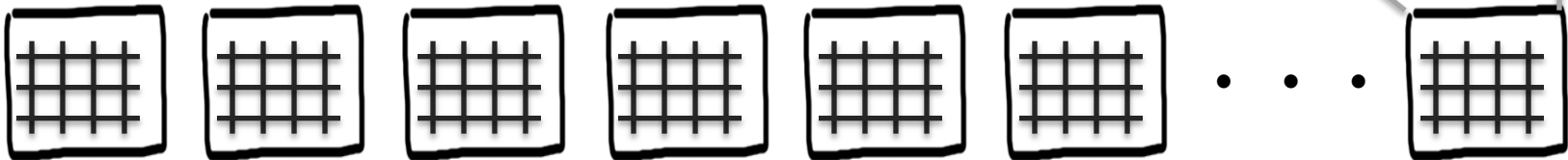
Infiniband

0110101001010010

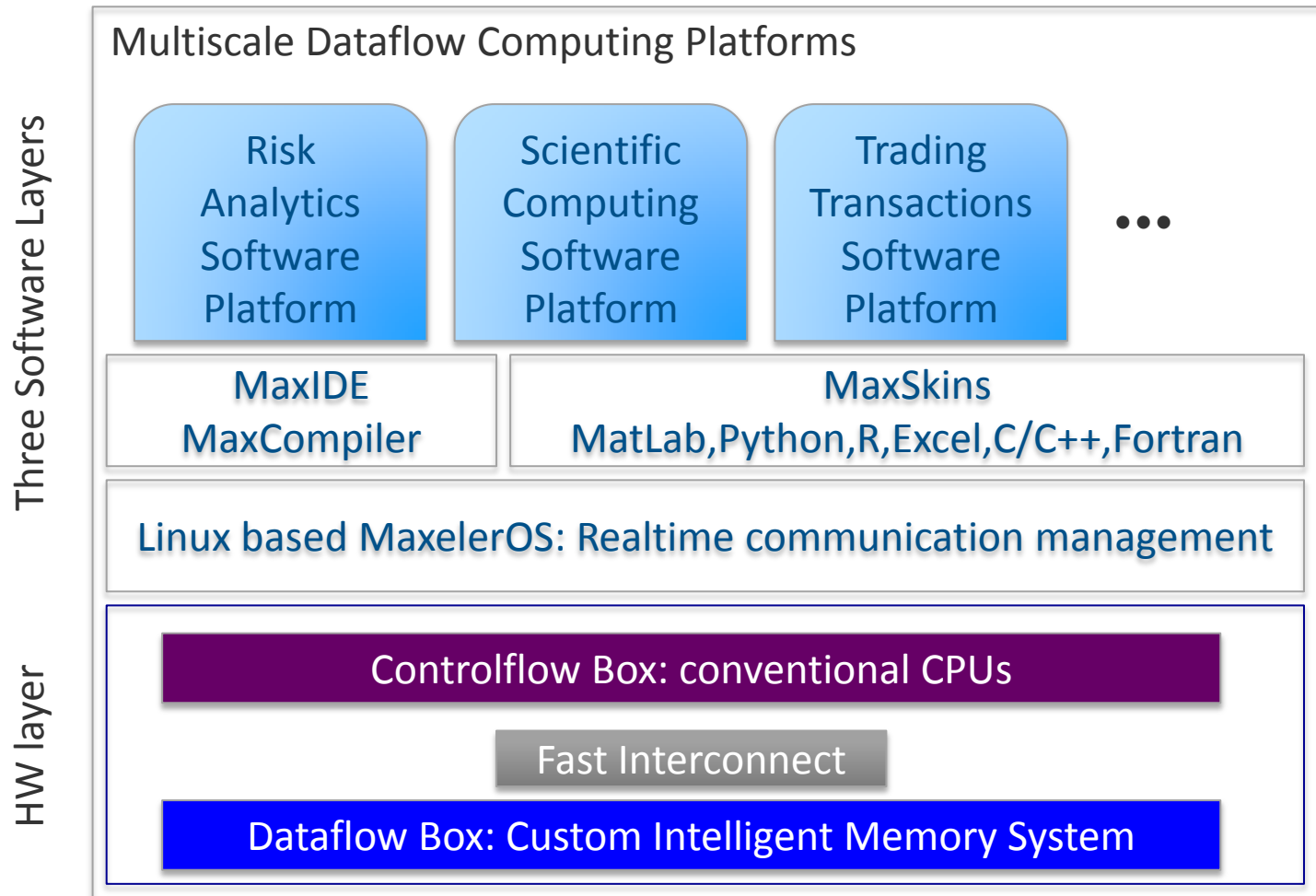


Memory

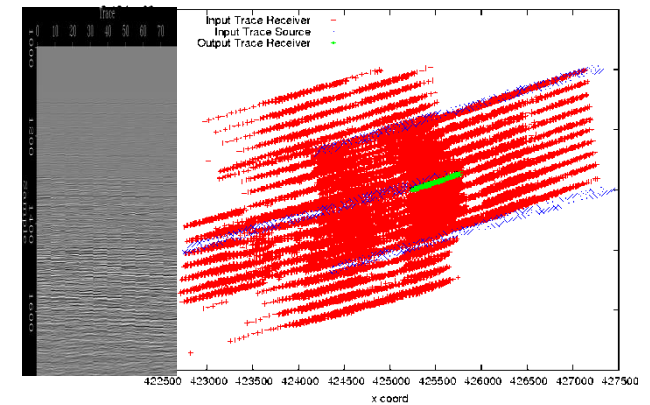
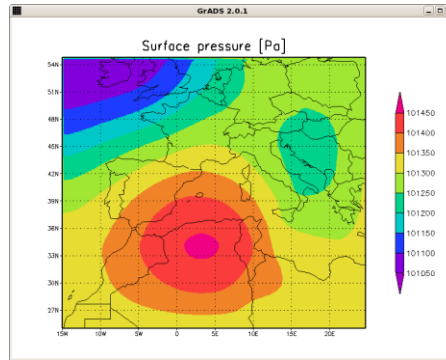
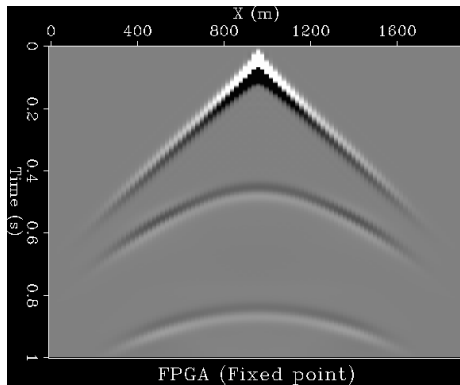
Dataflow Engines



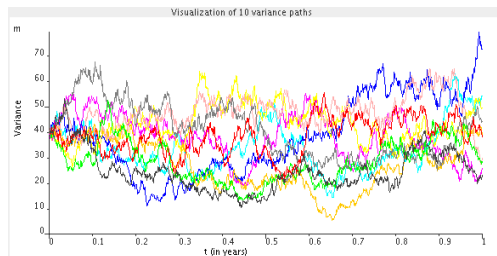
The Multiscale Dataflow Computer



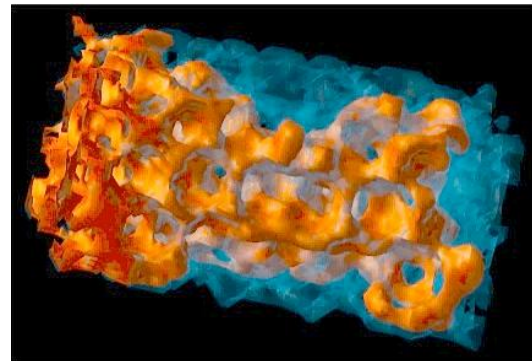
Multiscale Dataflow Advantage



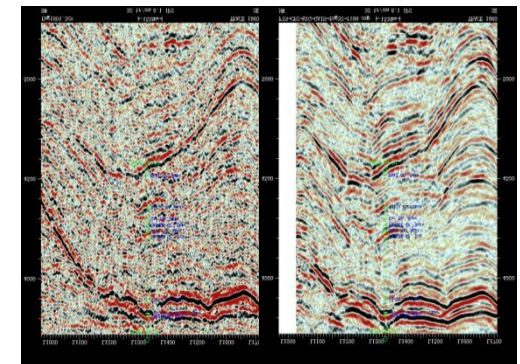
Acoustic Modelling 25x Weather Modeling 60x Trace Processing 22x



Financial Risk 32x



Fluid Flow 30x



Seismic Imaging 29x

Amdahl's Laws

- First law: Serial processors always win; too much time spent in programming the parallel processor.
- Second law: fraction of serial code, s , limits speedup to:

$$S_p = T_1 / (T_1 (s) + T_1 (1-s)/p) \text{ or}$$

$$S_p = 1 / (s + (1-s)/p)$$



Gene Amdahl

Slotnick's law (of effort)

“The parallel approach to computing does require that some original thinking be done about numerical analysis and data management in order to secure efficient use.

In an environment which has represented the absence of the need to think as the highest virtue this is a decided disadvantage.”

-Daniel Slotnick





Dataflow Application Areas

Finance, Geophysics, Chemistry
Physics, Genetics, Astronomy

Application Programming Interface

MaxCompiler: Dataflow in Space and Time
Heterogeneous dataflow+controlflow optimization

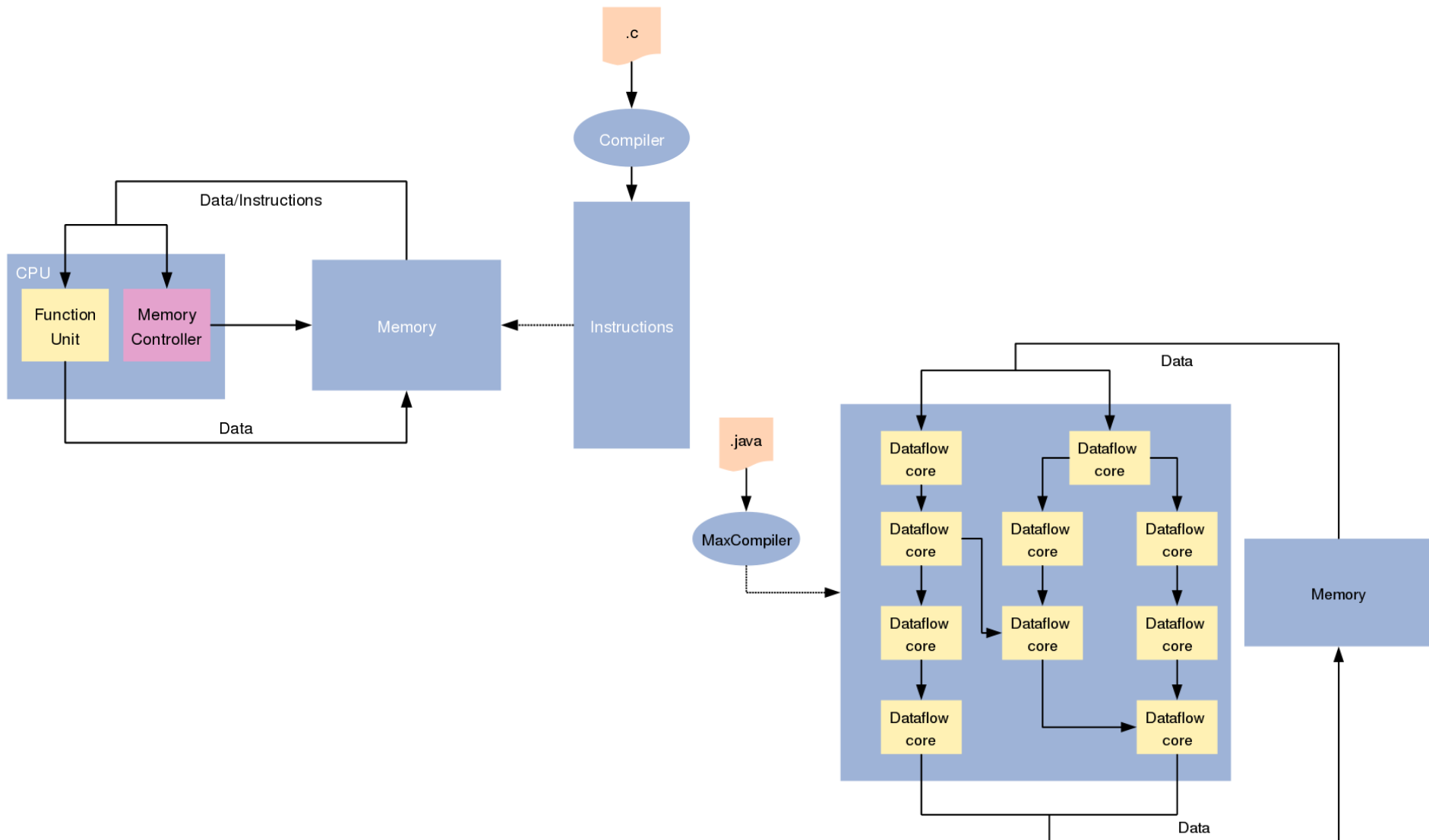
Transactions Management

MaxelerOS manages dataflow transaction
MaxelerOS keeps Dataflow-Controlflow balance

Architecture: Static Dataflow

Static Dataflow microarchitecture, cards, boxes
Predictable execution time and efficiency

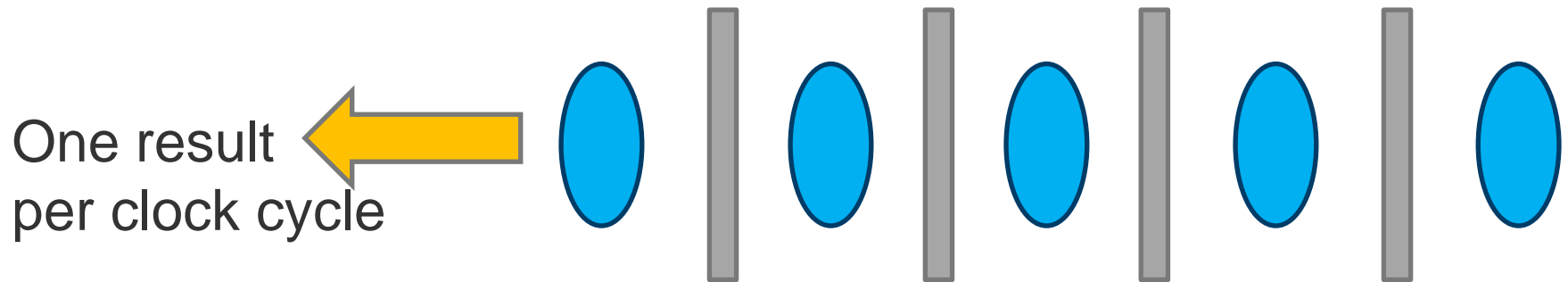
Control flow vs. Dataflow



Static Dataflow

“Systolic Arrays” without nearest neighbour interconnect restrictions

Static ultradeep (>1000 stage) computing pipelines



~~The power challenge~~

The data movement challenge

	Today	2018-20
Double precision FLOP	100pj	10pj
Moving data on-chip: 1mm	6pj	
Moving data on-chip: 20mm	120pj	
Moving data to off-chip memory	5000pj	2000pj

- Moving data on-chip will use as much energy as computing with it
- Moving data off-chip will use **200x more energy!**
 - And is much slower as well

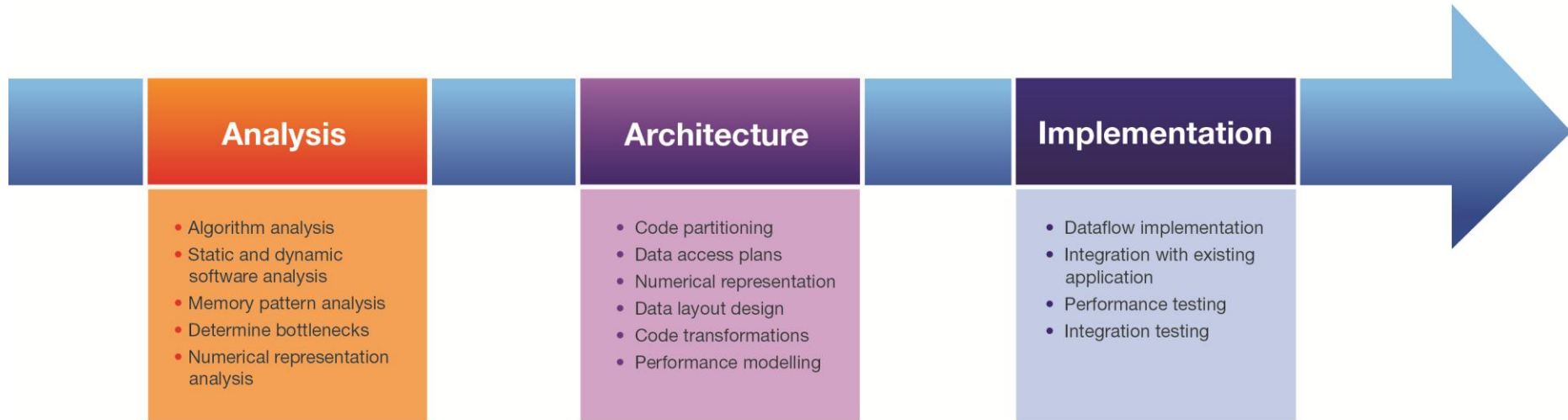
The memory hierarchy (challenge)

John von Neumann, 1946:

“We are forced to recognize the possibility of constructing a hierarchy of memories, each of which has greater capacity than the preceding, but which is less quickly accessible.”



Vertical Co-design of applications



- Deploy domain expertise to **co-design** application, algorithm and computer architecture

$$17 \times 24 = ?$$

Thinking Fast and Slow

Daniel Kahneman

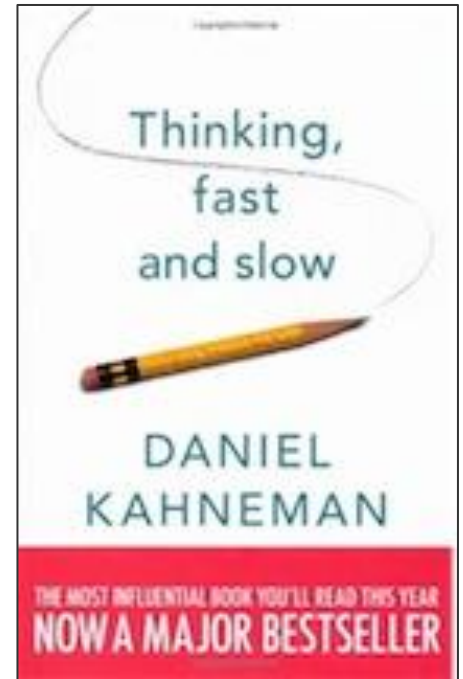
Nobel Prize in Economics, 2002

back to 17×24

Kahneman splits thinking into:

System 1: fast, hard to control ... 400

System 2: slow, easier to control ... 408



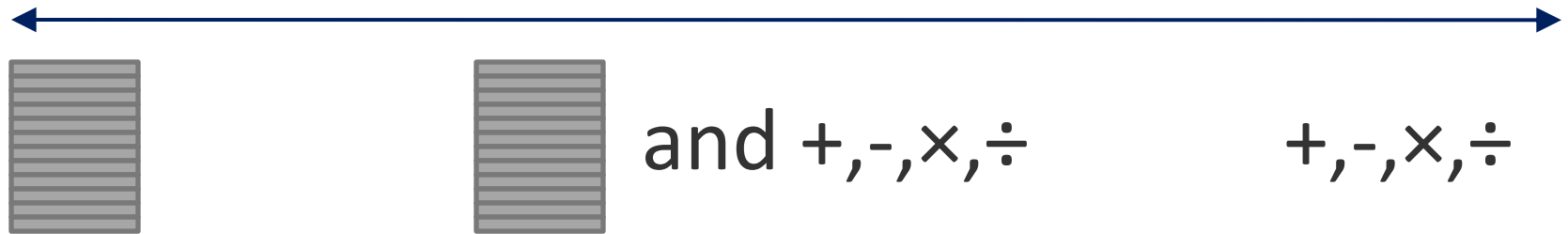
Putting it all together on the arithmetic level

Computing $f(x)$ in the range $[a,b]$ with $|E| \leq 2^{-n}$

Table

Table+Arithmetic

Arithmetic



Tradeoff: number of coefficients, number of bits per coefficient, range versus precision of result and maximal versus average error of result

Dong-U Lee, Altaf Abdul Gaffar, Oskar Mencer, Wayne Luk

Optimizing Hardware Function Evaluation

IEEE Transactions on Computers. vol. 54, no. 12, pp. 1520-1531. Dec, 2005.

Given range and precision for the result, what is the optimal table+arithmetic solution?

Minimal Latency (Optimized for Latency)

Range [bits]	Precision [bits]					
	4	8	12	16	20	24
24	sin: <i>tp2</i> , 24, 78 log: <i>po2</i> , 8, 30 sqr: <i>tp2</i> , 18, 912	sin: <i>tp2</i> , 28, 348 log: <i>tp2</i> , 12, 78 sqr: <i>tp2</i> , 24, 2400	sin: <i>tp2</i> , 32, 792 log: <i>tp2</i> , 15, 384 sqr: <i>tp2</i> , 26, 10368	sin: <i>tp2</i> , 32, 3168 log: <i>tp2</i> , 21, 1056 sqr: <i>tp2</i> , 30, 23808	sin: <i>tp2</i> , 40, 7872 log: <i>tp2</i> , 24, 4800 sqr: <i>tp3</i> , 34, 17920	sin: <i>tp3</i> , 42, 5504 log: <i>tp2</i> , 28, 11136 sqr: <i>tp4</i> , 39, 12800
20	sin: <i>po2</i> , 19, 61 log: <i>po2</i> , 7, 27 sqr: <i>tp2</i> , 18, 456	sin: <i>tp2</i> , 24, 300 log: <i>tp2</i> , 12, 78 sqr: <i>tp2</i> , 20, 2016	sin: <i>tp2</i> , 28, 696 log: <i>tp2</i> , 15, 384 sqr: <i>tp2</i> , 24, 4800	sin: <i>tp2</i> , 32, 3168 log: <i>tp2</i> , 21, 1056 sqr: <i>tp2</i> , 32, 12288	sin: <i>tp2</i> , 32, 6336 log: <i>tp2</i> , 24, 4800 sqr: <i>tp3</i> , 33, 8704	sin: <i>tp3</i> , 38, 4992 log: <i>tp2</i> , 27, 10752 sqr: <i>tp3</i> , 37, 19456
16	sin: <i>tp2</i> , 16, 54 log: <i>po2</i> , 7, 27 sqr: <i>tp2</i> , 17, 324	sin: <i>tp2</i> , 19, 240 log: <i>tp2</i> , 12, 78 sqr: <i>tp2</i> , 18, 912	sin: <i>tp2</i> , 24, 600 log: <i>tp2</i> , 15, 384 sqr: <i>tp2</i> , 24, 2400	sin: <i>tp2</i> , 28, 2784 log: <i>tp2</i> , 21, 1056 sqr: <i>tp2</i> , 26, 10368	sin: <i>tp2</i> , 32, 6336 log: <i>tp2</i> , 24, 4800 sqr: <i>tp2</i> , 30, 23808	sin: <i>tp3</i> , 32, 4224 log: <i>tp2</i> , 27, 10752 sqr: <i>tp3</i> , 34, 17920
12	sin: <i>po2</i> , 9, 31 log: <i>po2</i> , 7, 27 sqr: <i>tp2</i> , 12, 156	sin: <i>tp2</i> , 16, 204 log: <i>tp2</i> , 12, 78 sqr: <i>tp2</i> , 18, 456	sin: <i>tp2</i> , 20, 504 log: <i>tp2</i> , 15, 384 sqr: <i>tp2</i> , 20, 2016	sin: <i>tp2</i> , 24, 2400 log: <i>tp2</i> , 21, 1056 sqr: <i>tp2</i> , 24, 4800	sin: <i>tp2</i> , 28, 5568 log: <i>tp2</i> , 24, 4800 sqr: <i>tp2</i> , 31, 12288	sin: <i>tp2</i> , 32, 12672 log: <i>tp2</i> , 27, 10752 sqr: <i>tp3</i> , 33, 8704
8	sin: <i>po2</i> , 6, 22 log: <i>po2</i> , 7, 27 sqr: <i>tp2</i> , 7, 27	sin: <i>tp2</i> , 10, 132 log: <i>tp2</i> , 11, 72 sqr: <i>tp2</i> , 17, 324	sin: <i>tp2</i> , 16, 408 log: <i>tp2</i> , 15, 384 sqr: <i>tp2</i> , 18, 912	sin: <i>tp2</i> , 19, 1920 log: <i>tp2</i> , 21, 1056 sqr: <i>tp2</i> , 24, 2400	sin: <i>tp2</i> , 24, 4800 log: <i>tp2</i> , 24, 4800 sqr: <i>tp2</i> , 26, 10368	sin: <i>tp2</i> , 28, 11136 log: <i>tp2</i> , 27, 10752 sqr: <i>tp2</i> , 30, 23808
4	sin: <i>po2</i> , 6, 22 log: <i>po2</i> , 7, 27 sqr: <i>po2</i> , 7, 25	sin: <i>tp2</i> , 10, 132 log: <i>tp2</i> , 11, 72 sqr: <i>tp2</i> , 12, 156	sin: <i>tp2</i> , 15, 384 log: <i>tp2</i> , 15, 384 sqr: <i>tp2</i> , 18, 456	sin: <i>tp2</i> , 19, 1920 log: <i>tp2</i> , 17, 1056 sqr: <i>tp2</i> , 20, 2016	sin: <i>tp2</i> , 23, 4608 log: <i>tp2</i> , 24, 4800 sqr: <i>tp2</i> , 24, 4800	sin: <i>tp2</i> , 28, 11136 log: <i>tp2</i> , 27, 10752 sqr: <i>tp2</i> , 31, 12288

Architecture Level: Star versus Cube Stencil

More Computation in Less Time?

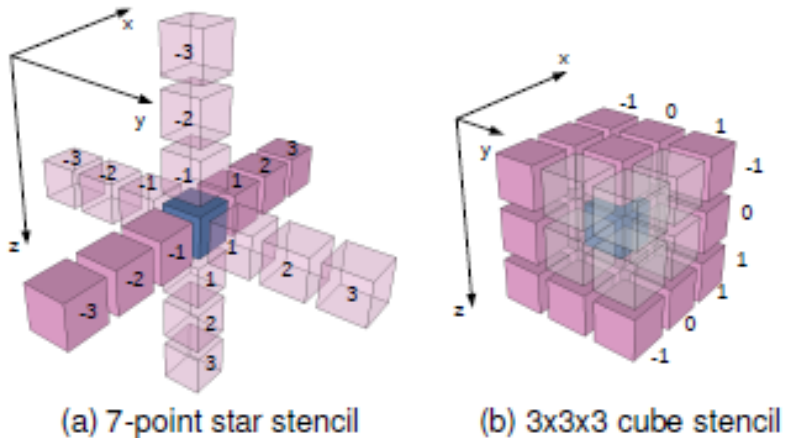
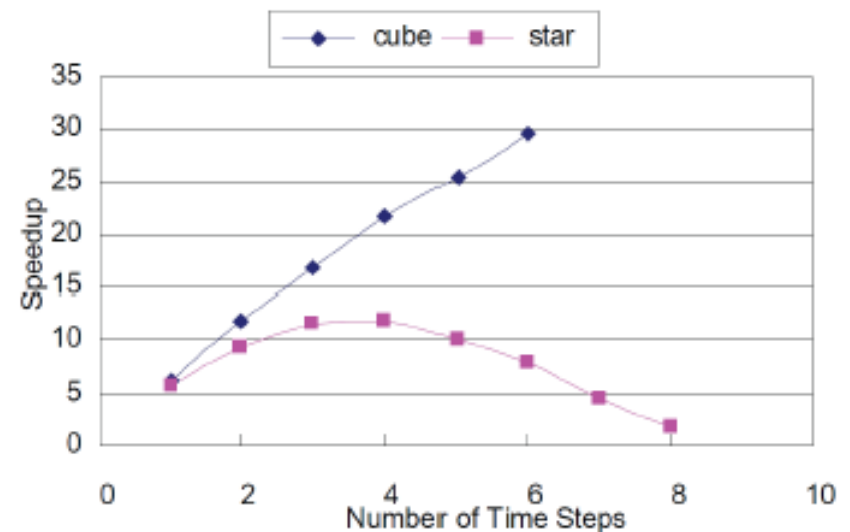
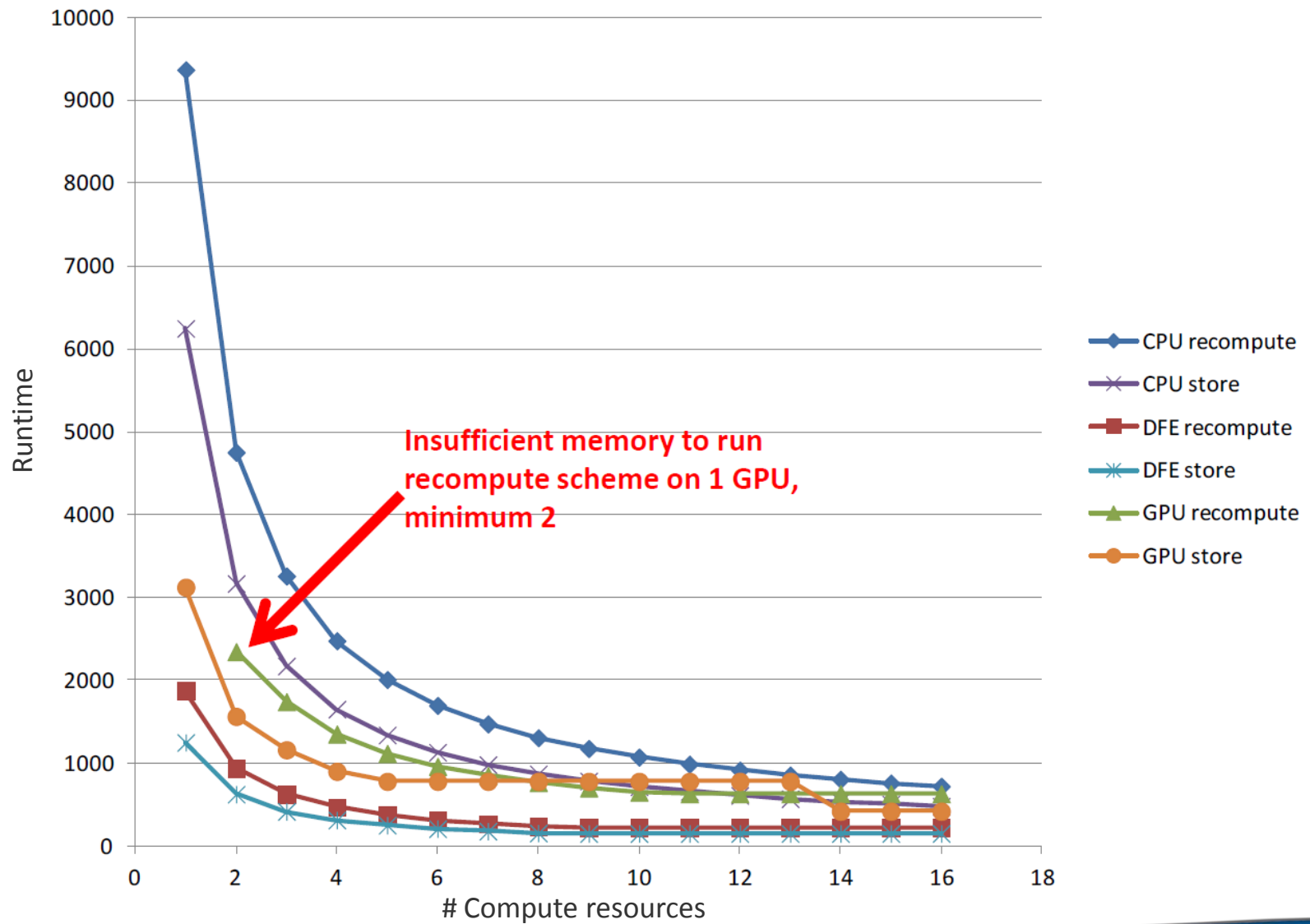


Figure 1: 2 Alternative stencil choices.

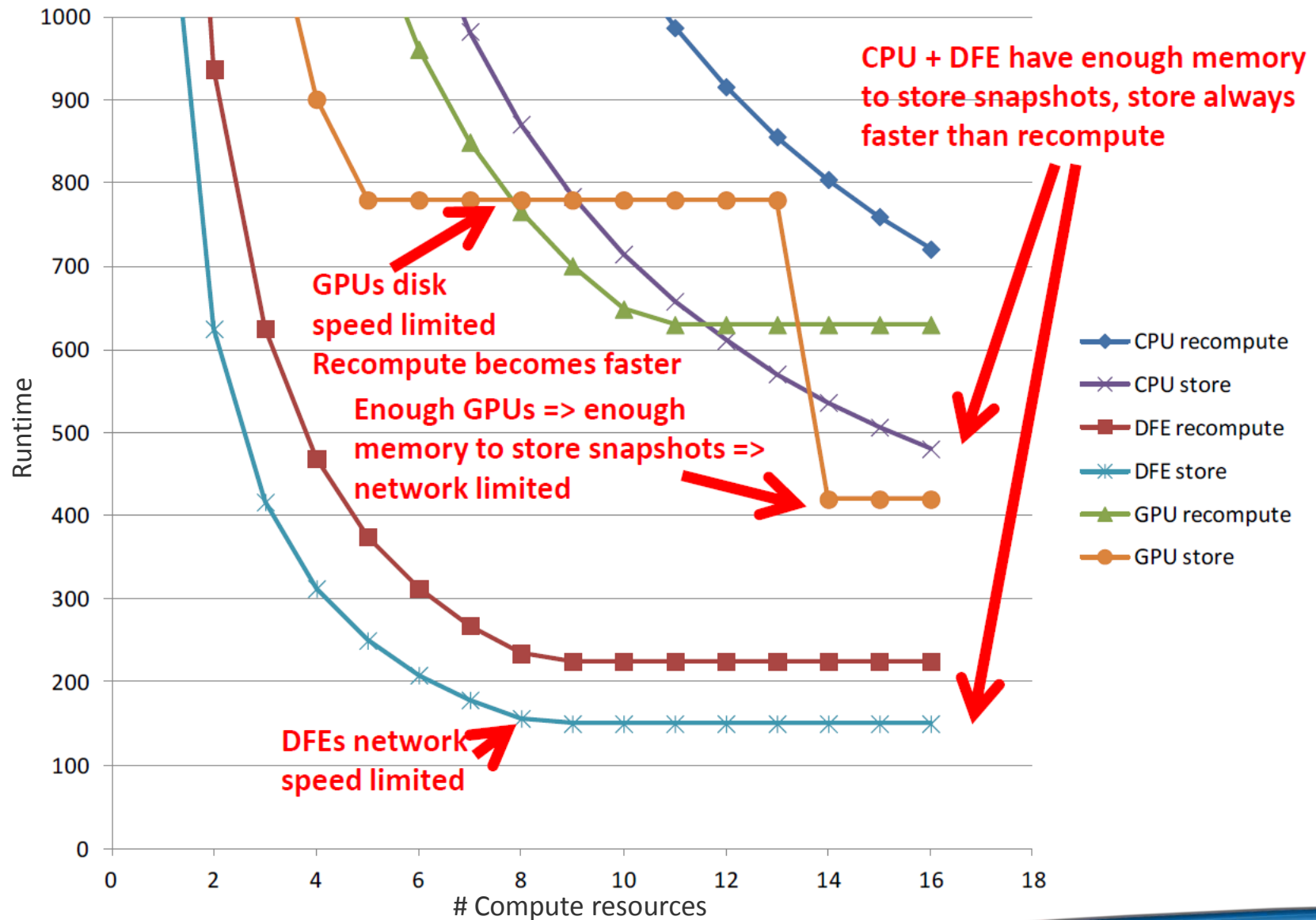
Local temporal parallelism
=> Cascading timesteps



System level: algorithm vs. resource



System level: algorithm vs. resource

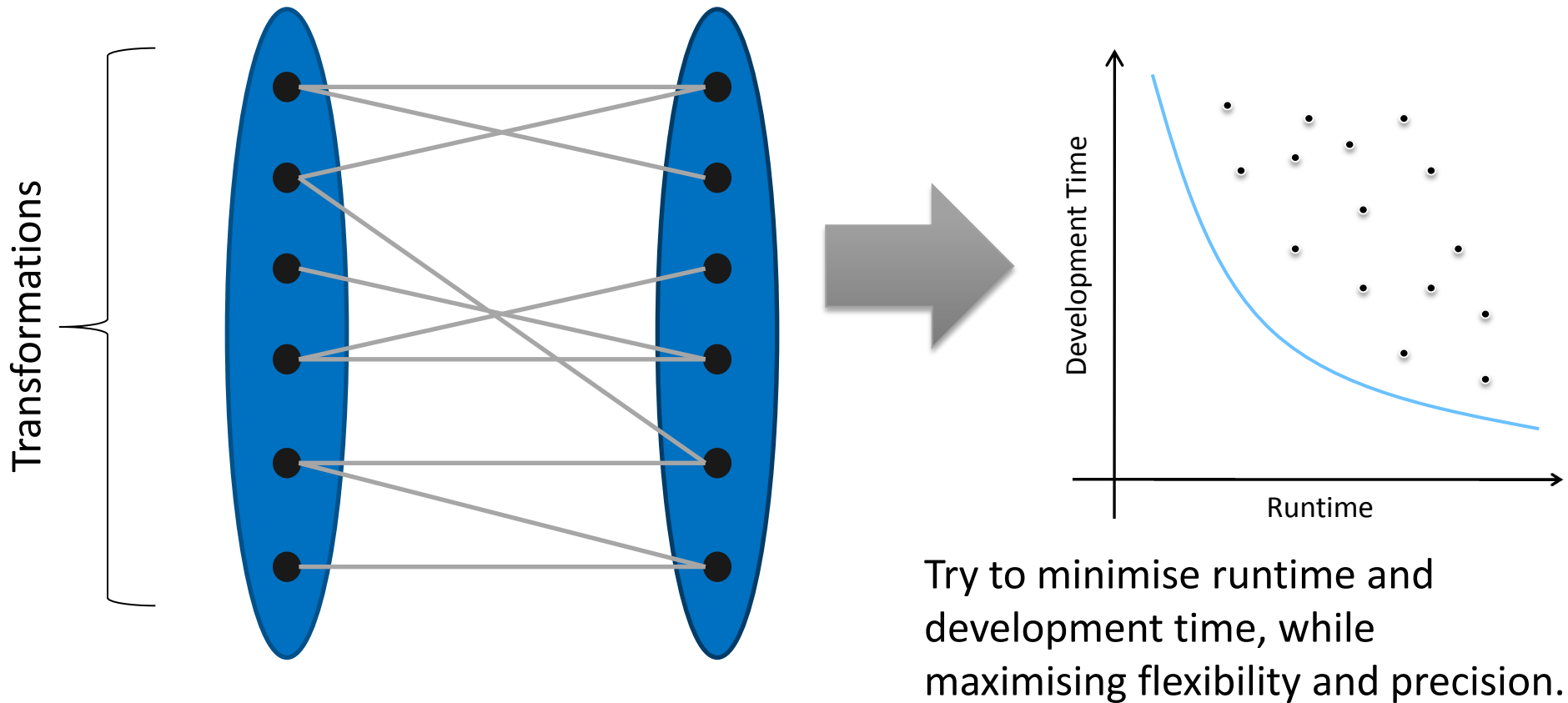


Identify and classify options

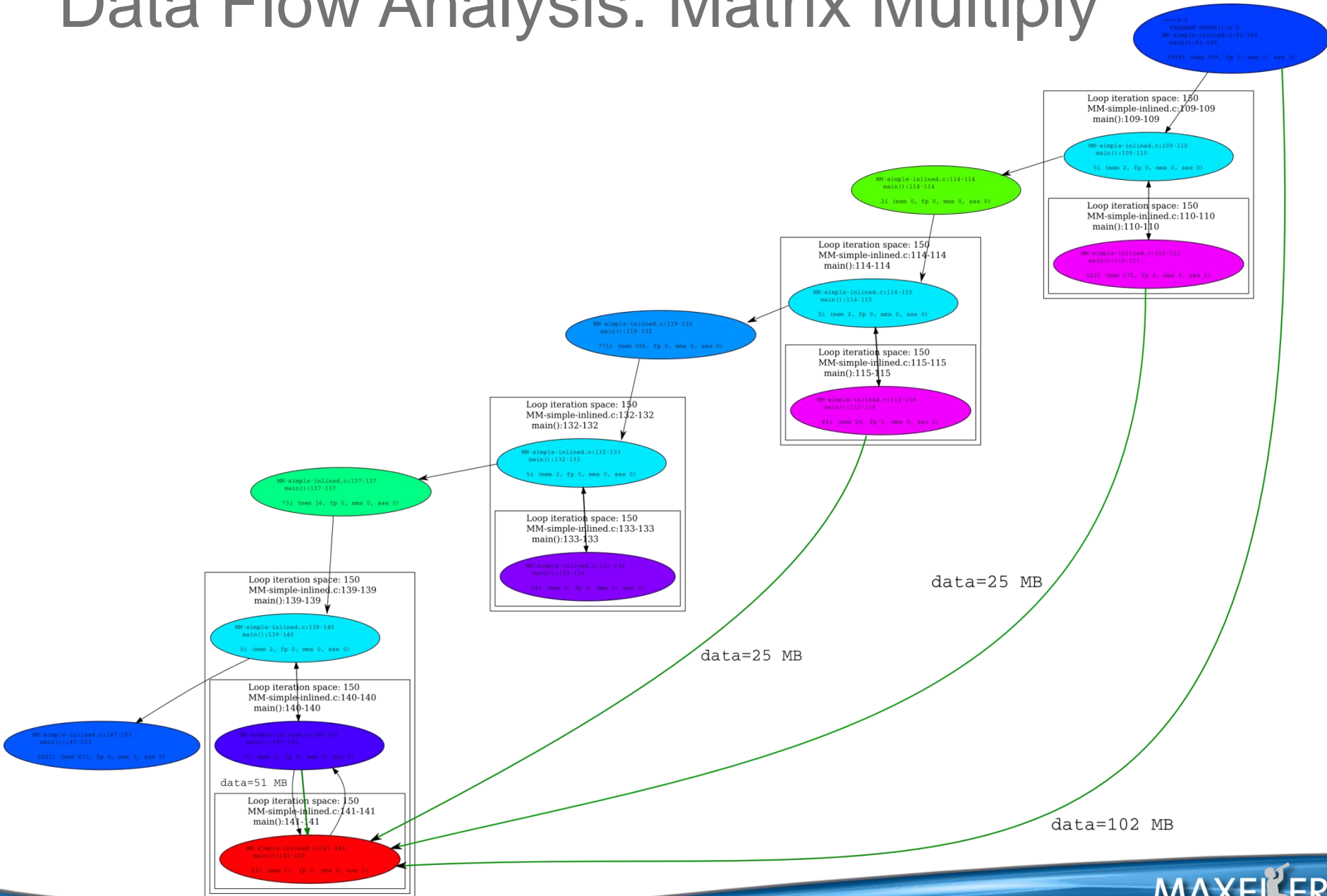
Data Access Plans

Code Partitioning

Pareto Optimal Options



Data Flow Analysis: Matrix Multiply



Maxeler Dataflow Computers



CPU's plus DFEs

Intel Xeon CPU cores and up to 4 DFEs with 192GB of RAM



DFEs shared over Infiniband

Up to 8 DFEs with 384GB of RAM and dynamic allocation of DFEs to CPU servers



Low latency connectivity

Intel Xeon CPUs and 1-4 DFEs with up to twelve 40Gbit Ethernet connections



MaxWorkstation

Desktop development system

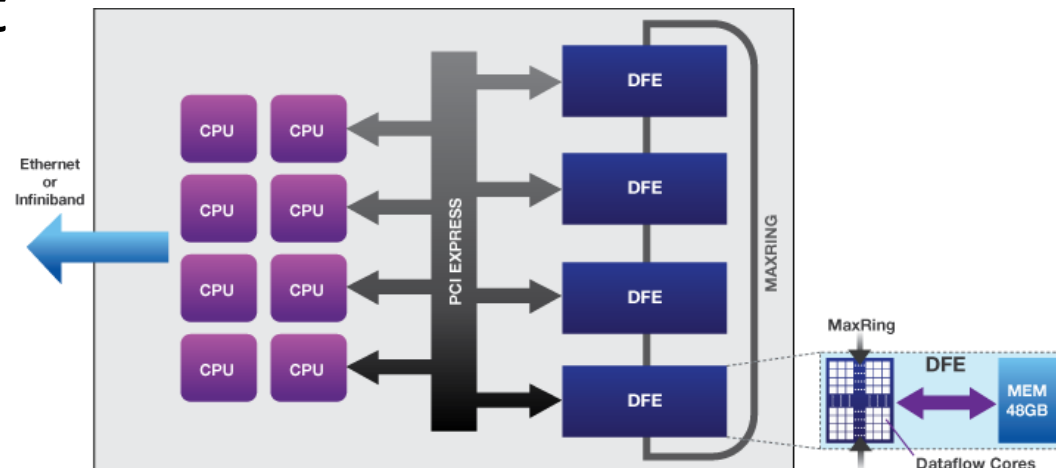
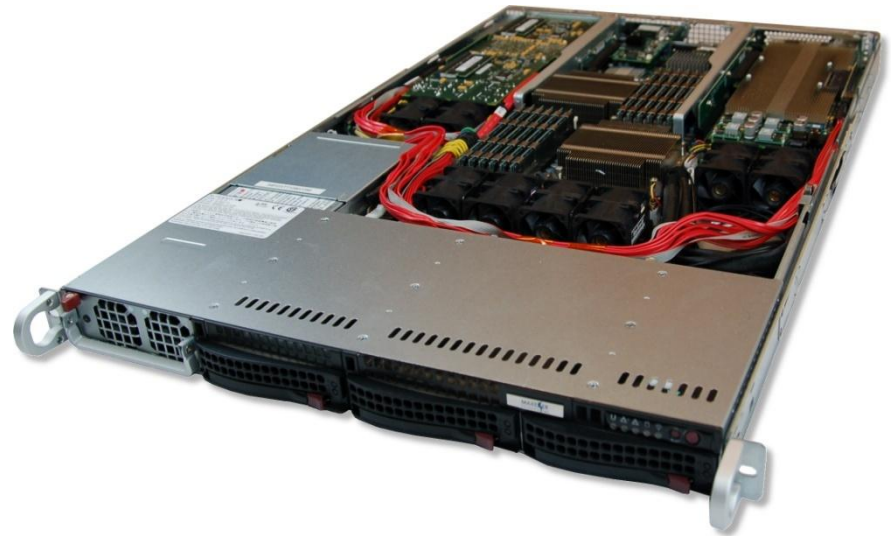


MaxCloud

On-demand scalable accelerated compute resource, hosted in London

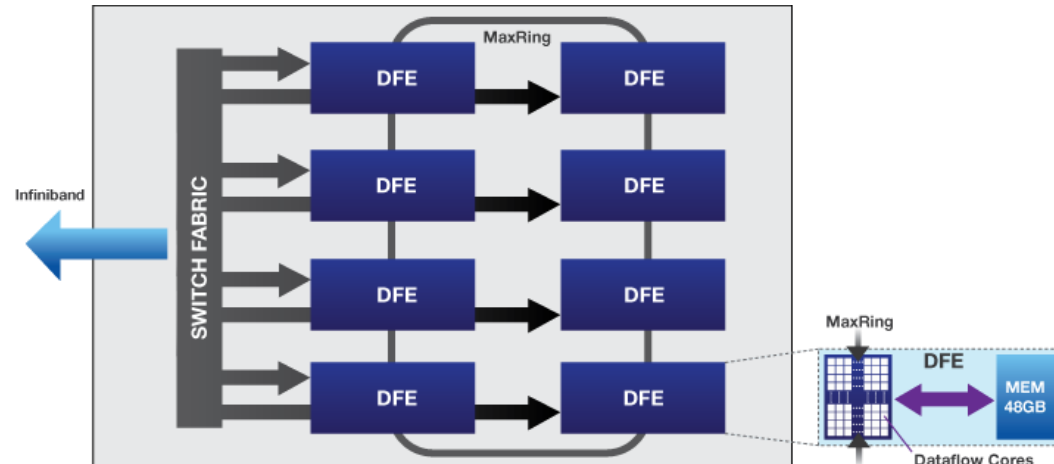
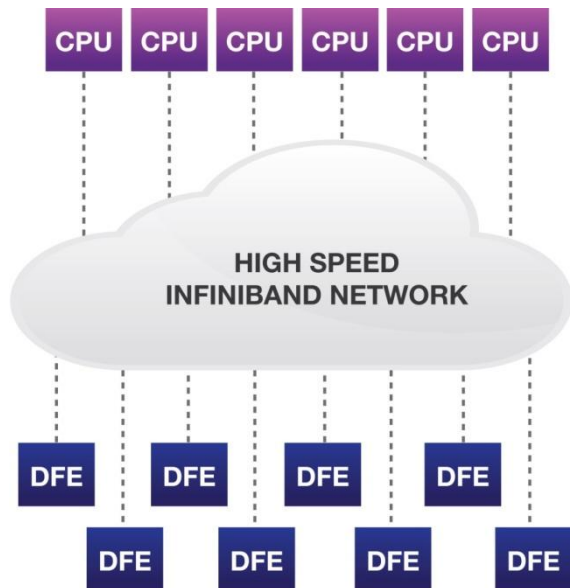
MPC-C500

- 1U Form Factor
- 4x dataflow engines
- 12 Intel Xeon cores
- 96GB DFE RAM
- Up to 192GB CPU RAM
- *MaxRing* interconnect
- 3x 3.5" hard drives
- Infiniband/10GigE



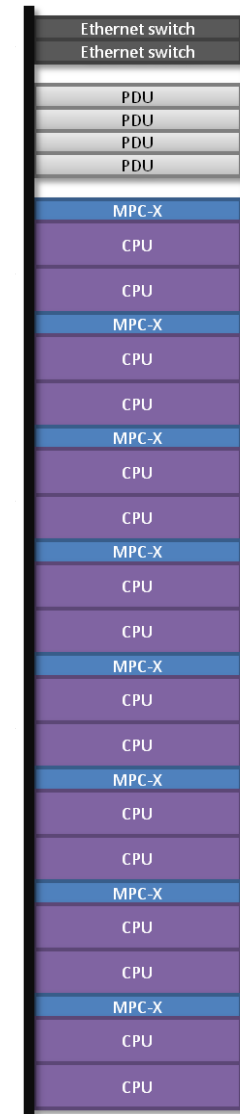
MPC-X1000

- 8 dataflow engines (384GB RAM)
- High-speed MaxRing
- Zero-copy RDMA between CPUs and DFEs over Infiniband
- Dynamic CPU/DFE balancing

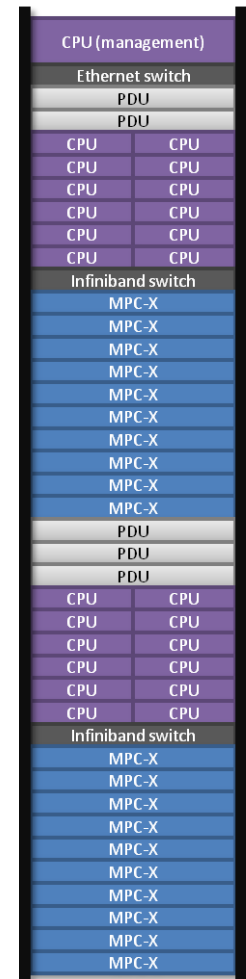


Dataflow clusters

- Optimized to balance resources for particular application challenges
- Flexible at design-time and at run-time

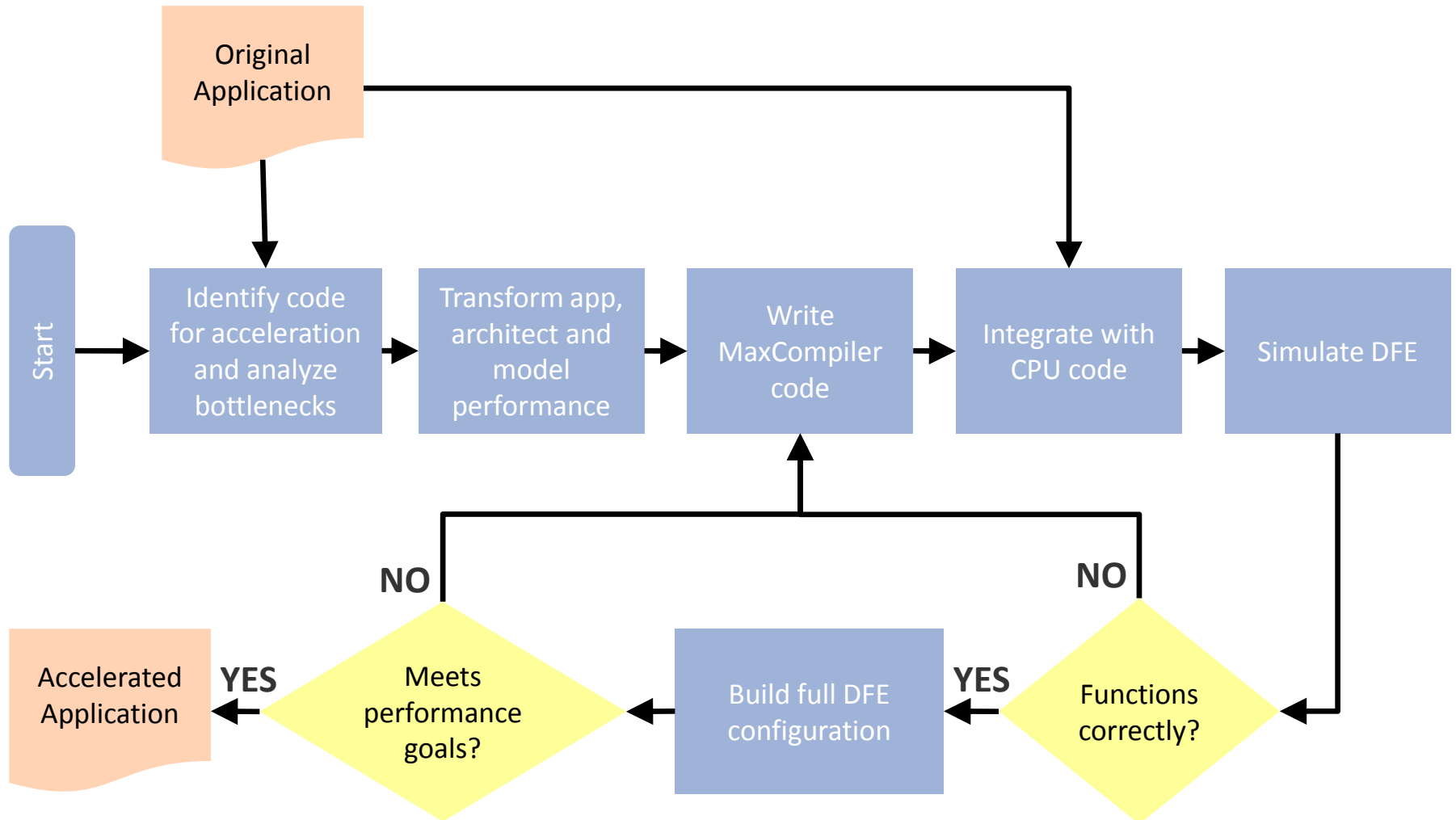


48U seismic
imaging cluster

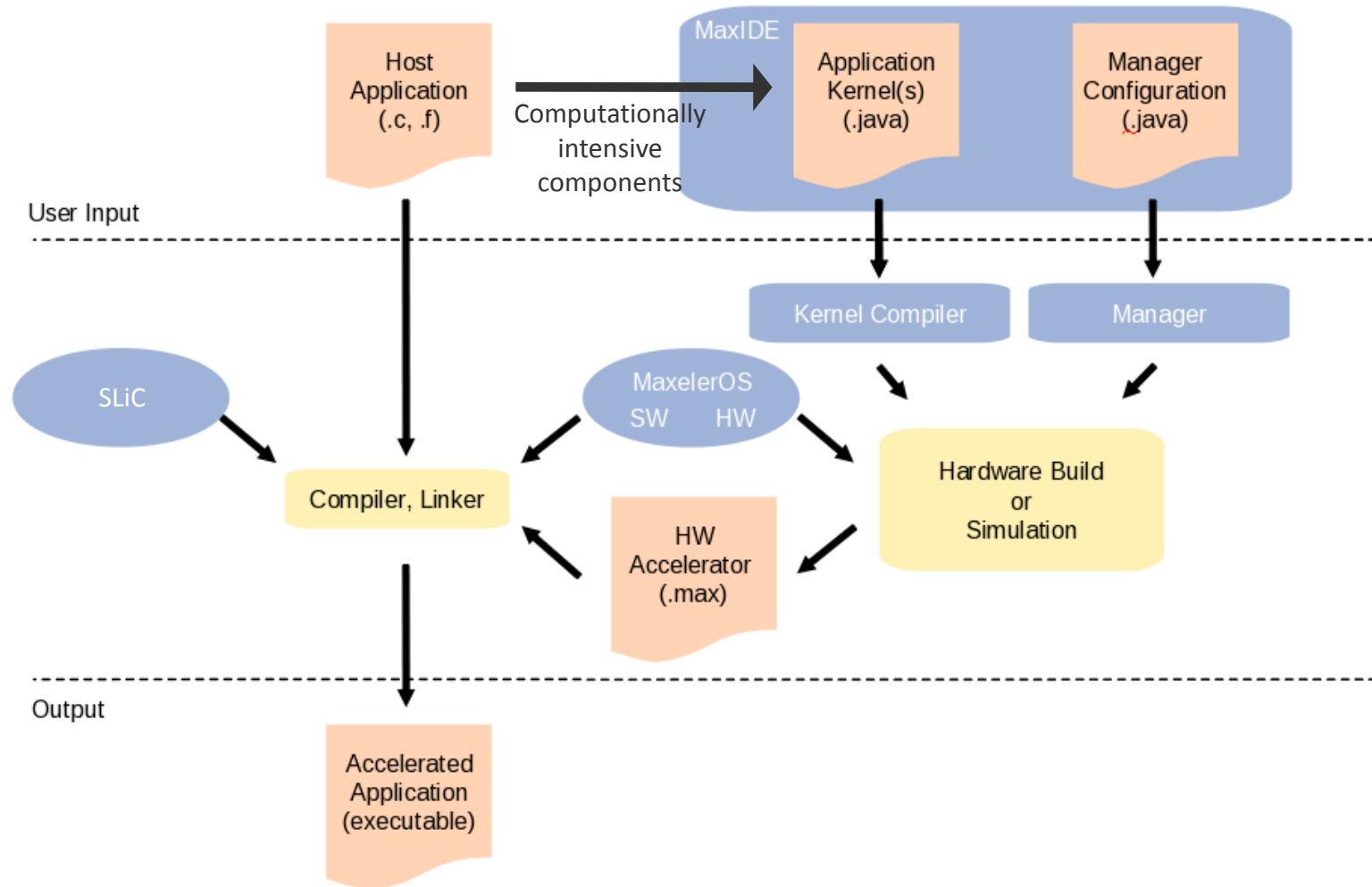


42U in-memory
analytics cluster

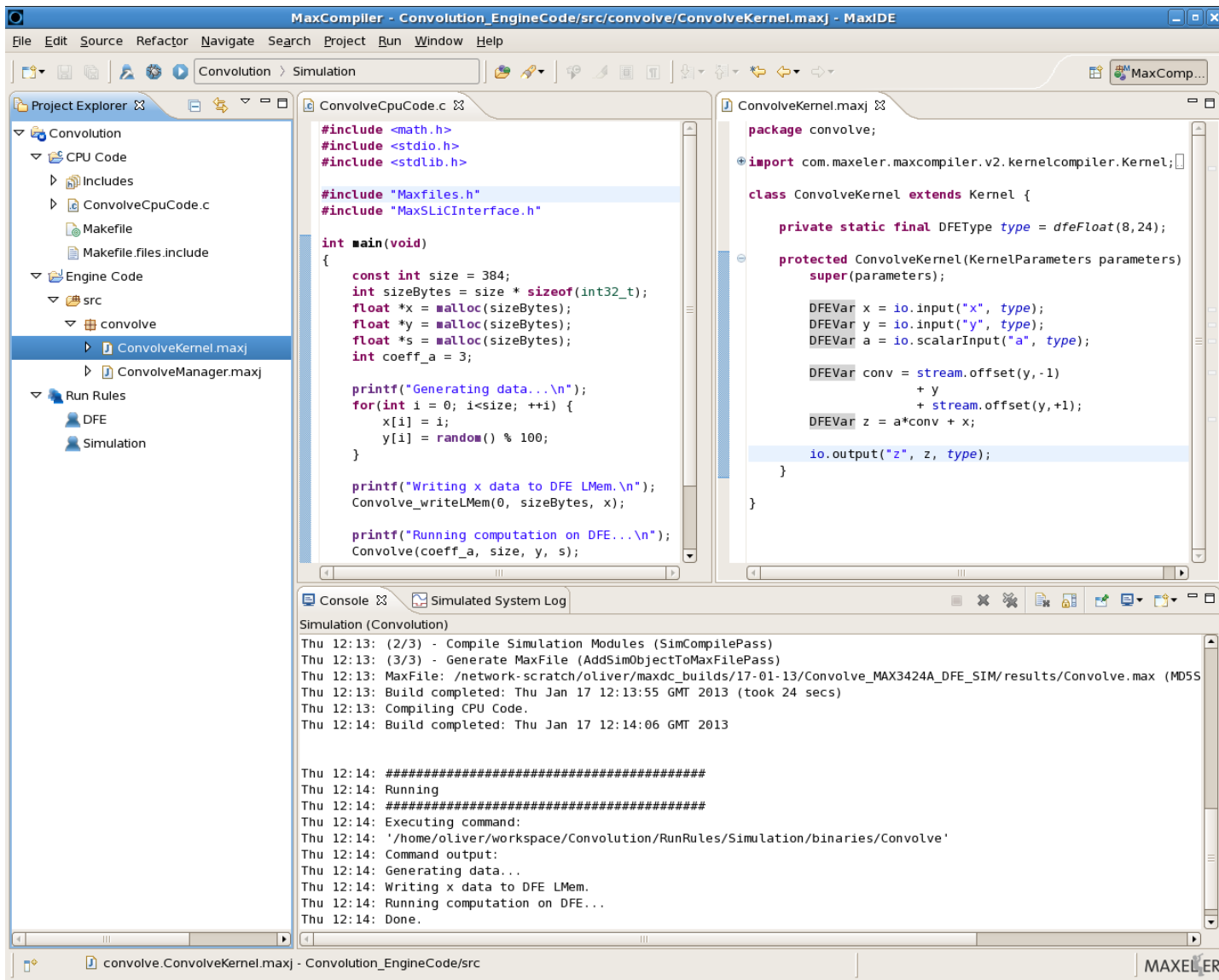
Application Programming Process



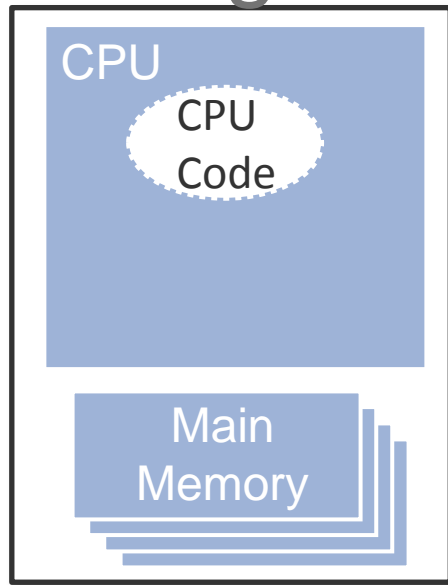
Programming with MaxCompiler



Programming with MaxCompiler



Programming with MaxCompiler



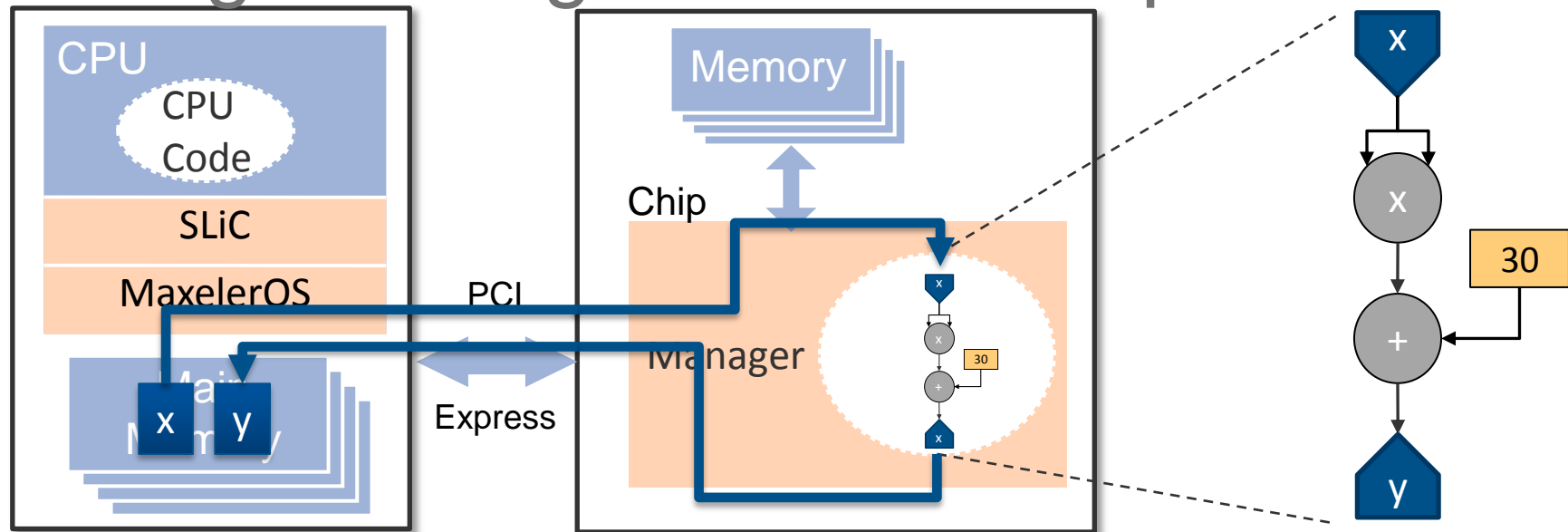
CPU Code (.c)

```
int *x, *y;  
for (int i =0; i < DATA_SIZE; i++)  
    y[i]= x[i] * x[i] + 30;
```



$$y_i = x_i \times x_i + 30$$

Programming with MaxCompiler



CPU Code (.c)

```
#include "MaxSLiCInterface.h"
#include "Calc.max"
int *x, *y;

Calc(x, y, DATA_SIZE)
```

Manager (.java)

```
Manager m = new Manager("Calc");
Kernel k =
    new MyKernel();

m.setKernel(k);
m.setIO(
    link("x", CPU),
    link("y", CPU));
m.createSLiCInterface();
m.build();
```

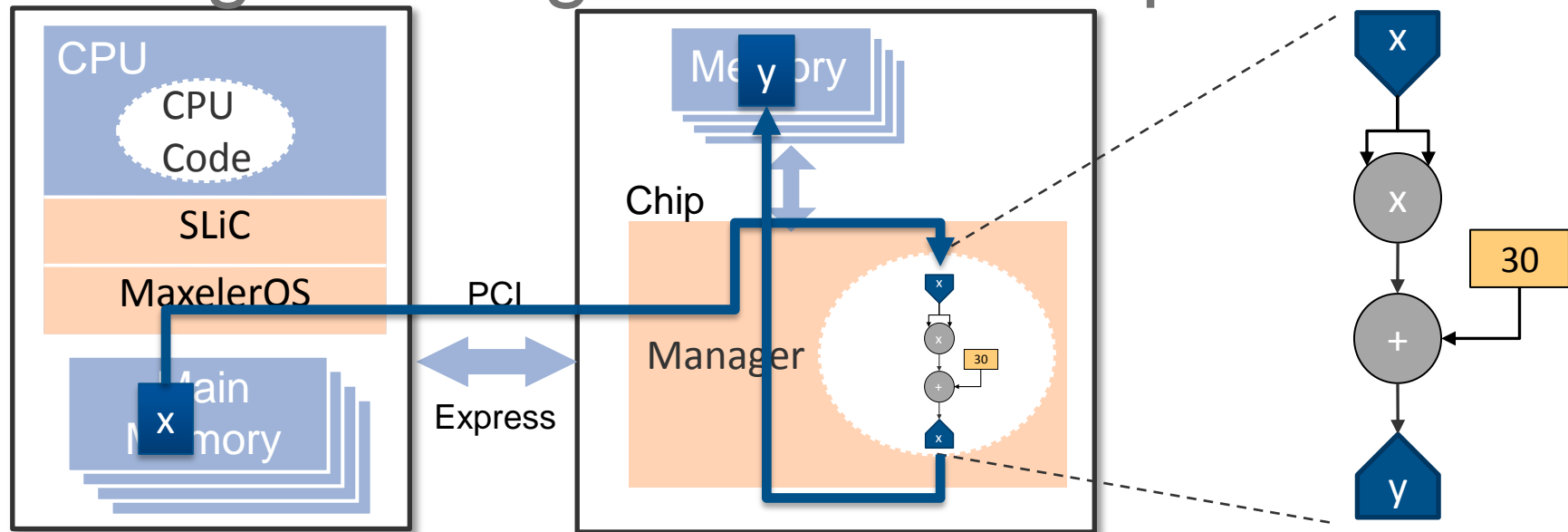
MyKernel (.java)

```
DFEVar x = io.input("x", dfeInt(32));

DFEVar result = x * x + 30;

io.output("y", result, dfeInt(32));
```

Programming with MaxCompiler



CPUCode (.c)

```
#include "MaxSLiCInterface.h"
#include "Calc.max"
int *x, *y;

Calc(x, DATA_SIZE)
```

Manager (.java)

```
Manager m = new Manager("Calc");
Kernel k =
    new MyKernel();

m.setKernel(k);
m.setIO(
    link("x", CPU),
    link("y", LMEM_LINEAR1D));
m.createSLiCInterface();
m.build();
```

MyKernel (.java)

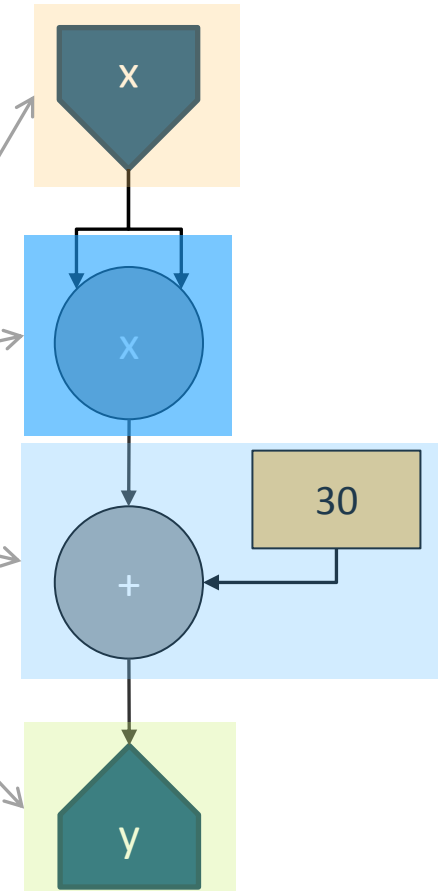
```
DfEVar x = io.input("x", dfelnt(32));

DfEVar result = x * x + 30;

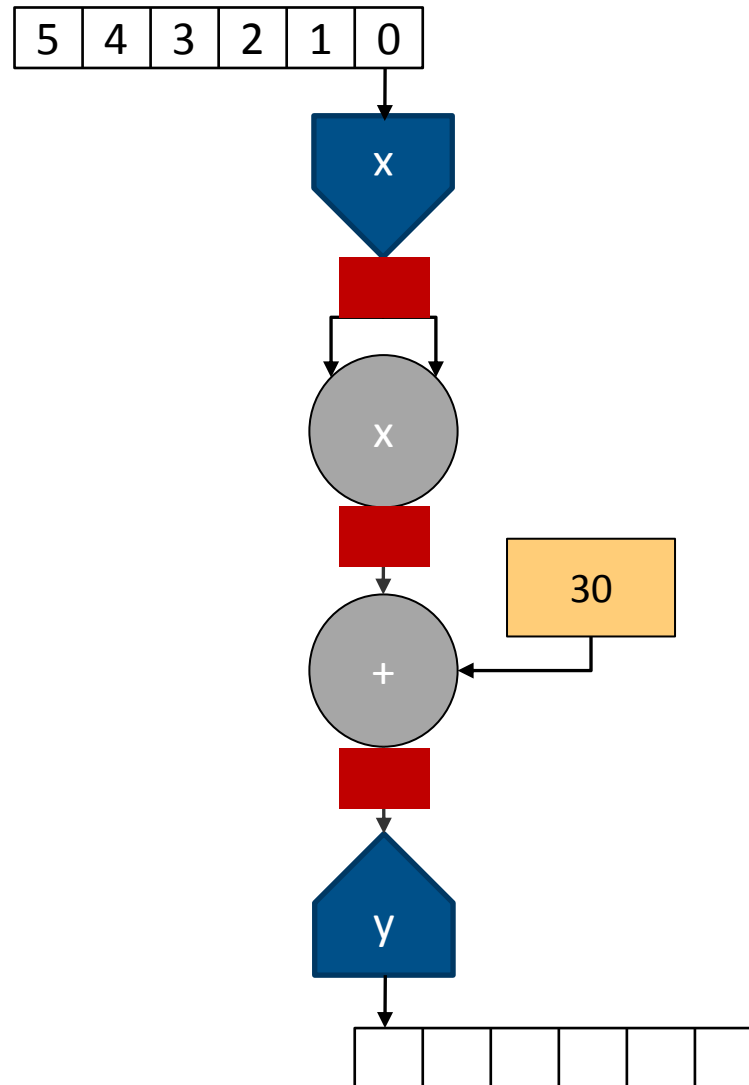
io.output("y", result, dfelnt(32));
```


The Full Kernel

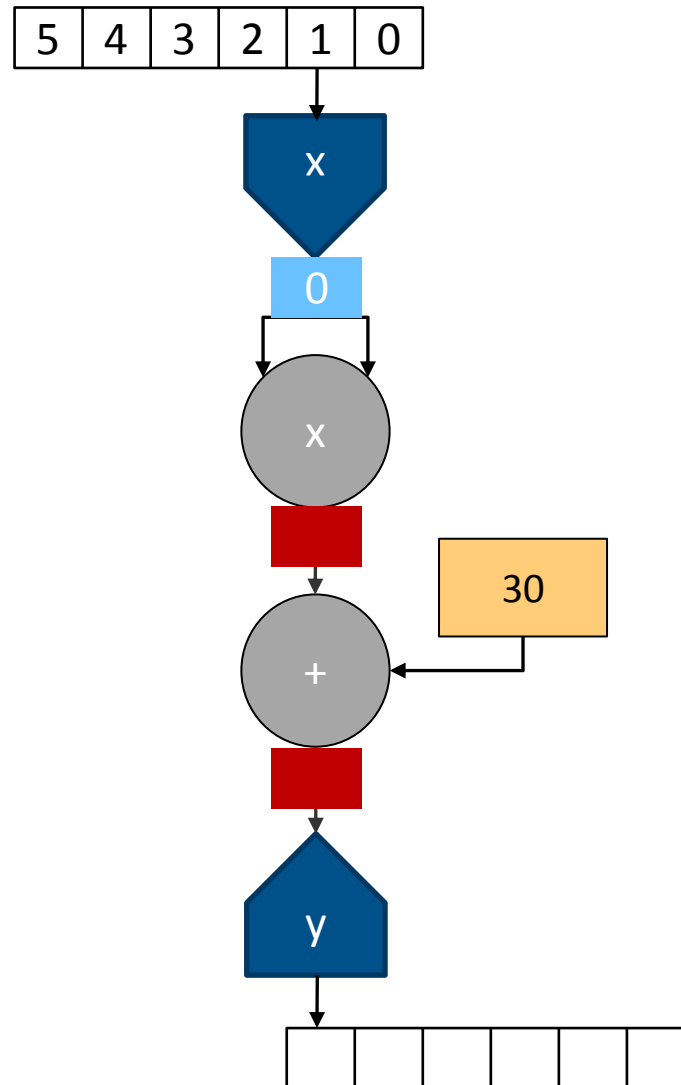
```
public class MyKernel extends Kernel {  
  
    public MyKernel (KernelParameters parameters) {  
        super(parameters);  
  
        HWVar x = io.input("x", hwInt(32));  
  
        HWVar result = x * x + 30;  
  
        io.output("y", result, hwInt(32));  
  
    }  
}
```



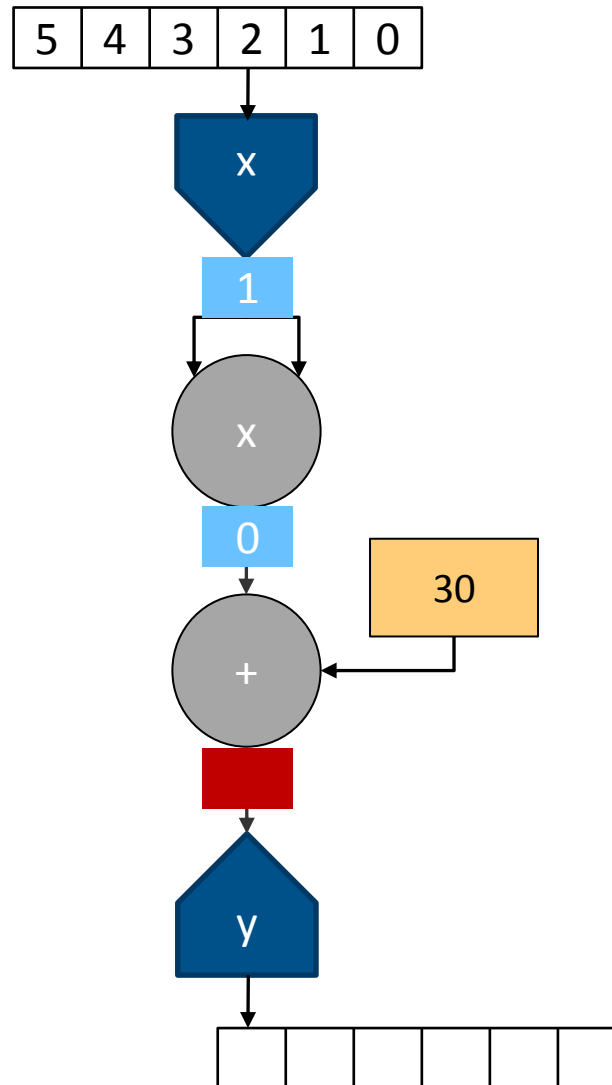
Kernel Streaming



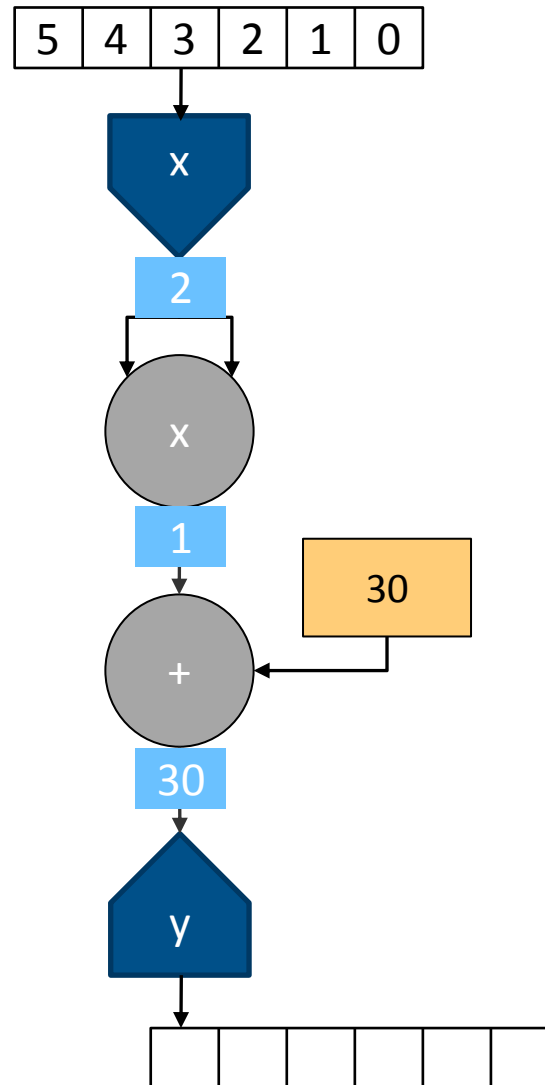
Kernel Streaming



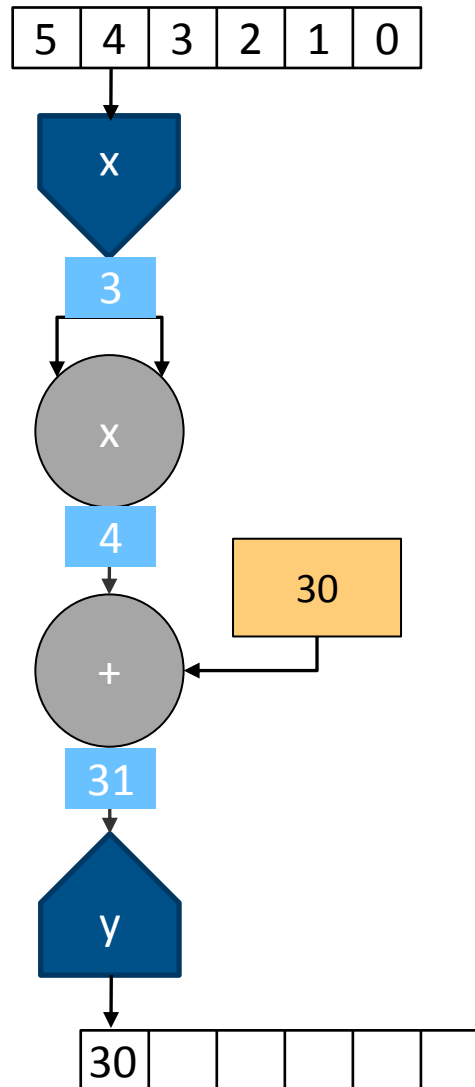
Kernel Streaming



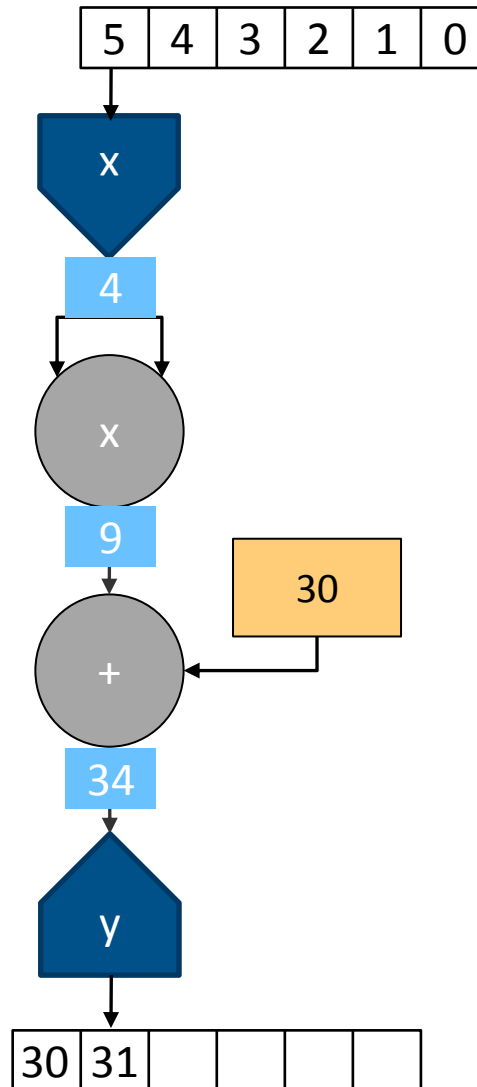
Kernel Streaming



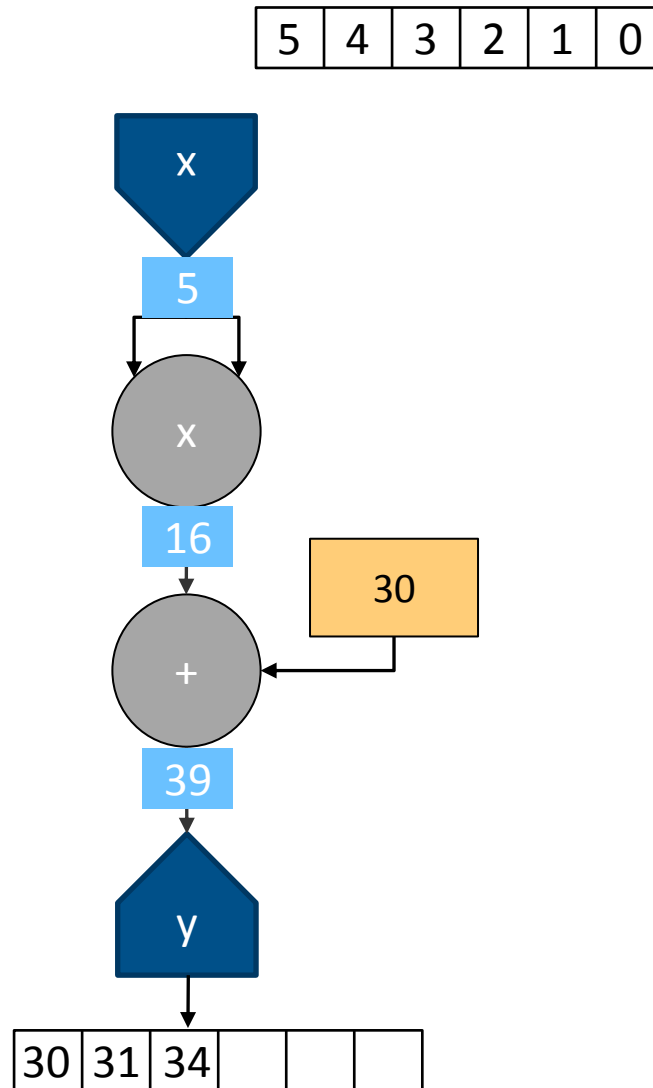
Kernel Streaming



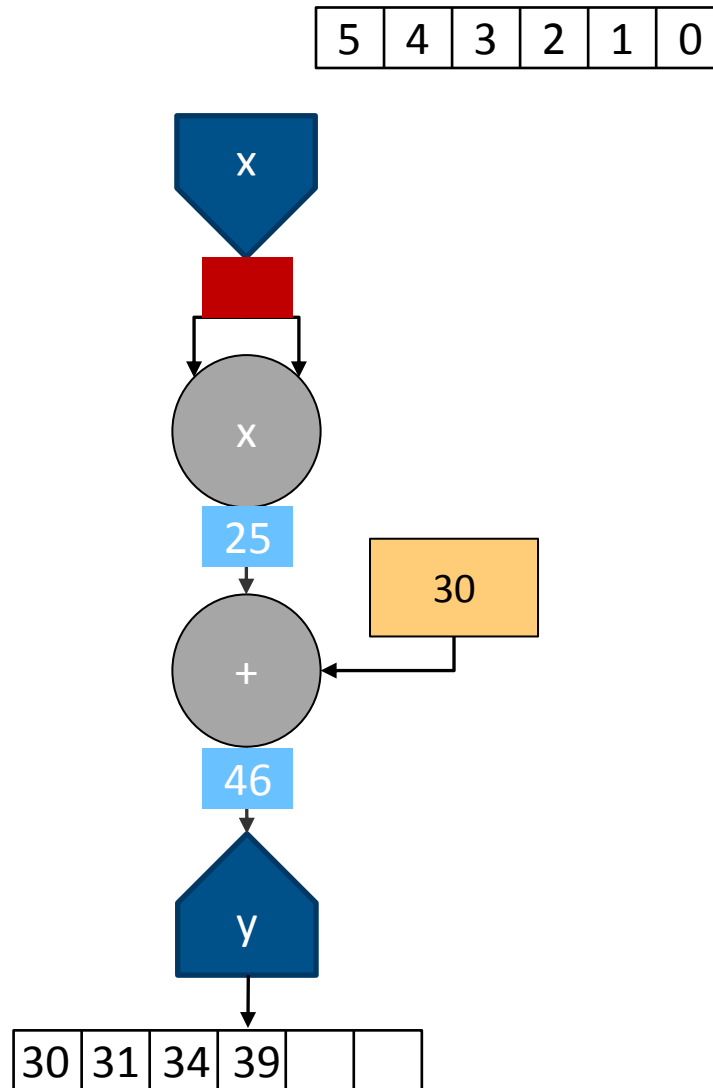
Kernel Streaming



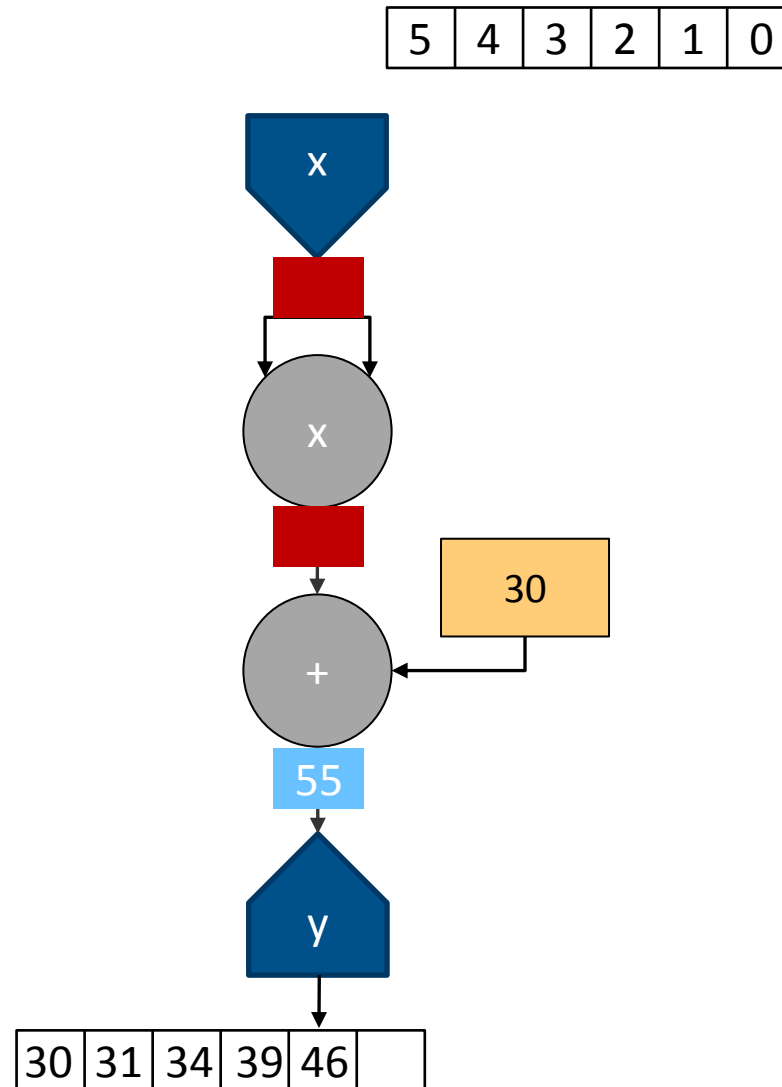
Kernel Streaming



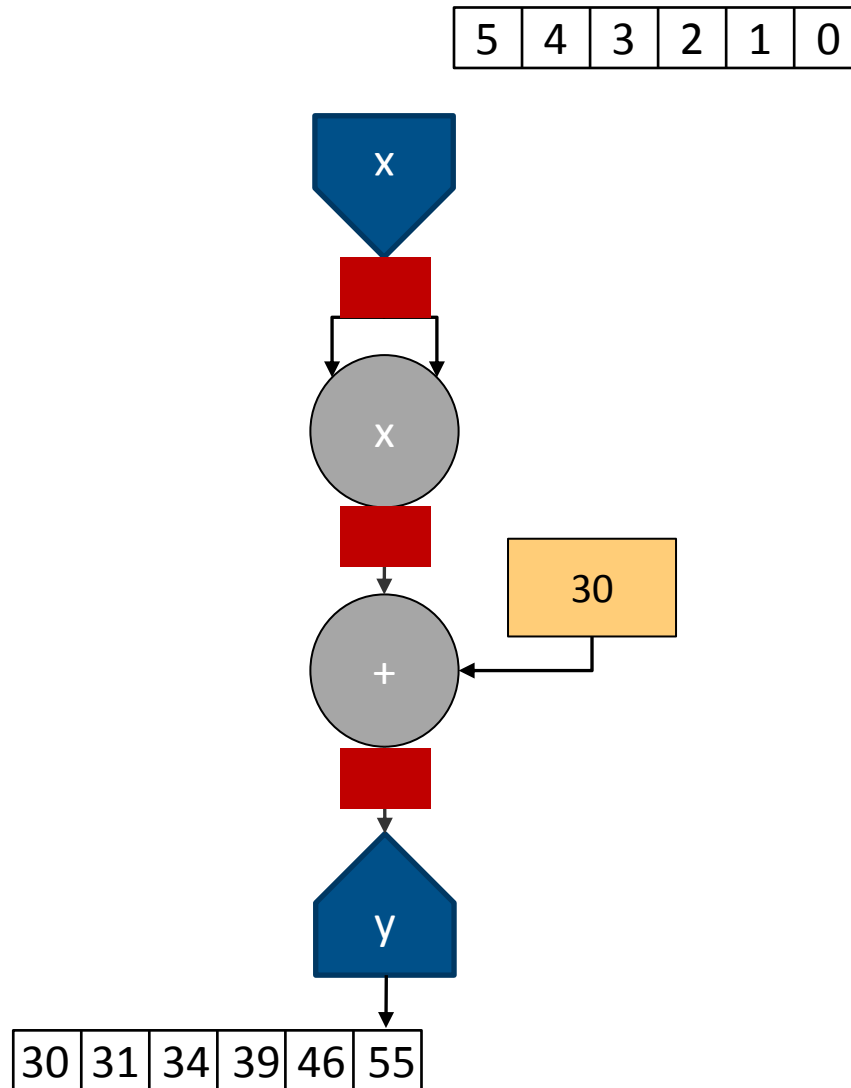
Kernel Streaming



Kernel Streaming

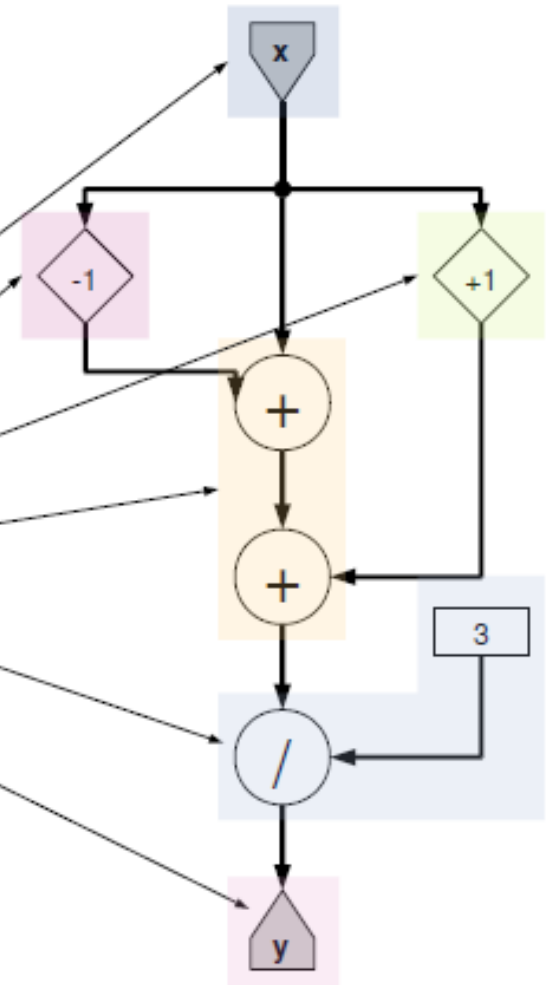


Kernel Streaming

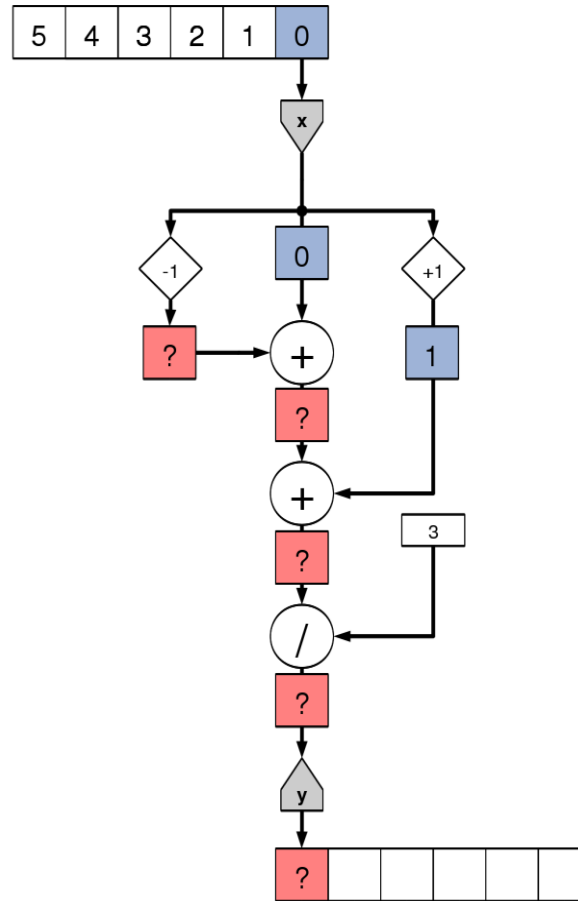


A (slightly) more complex kernel

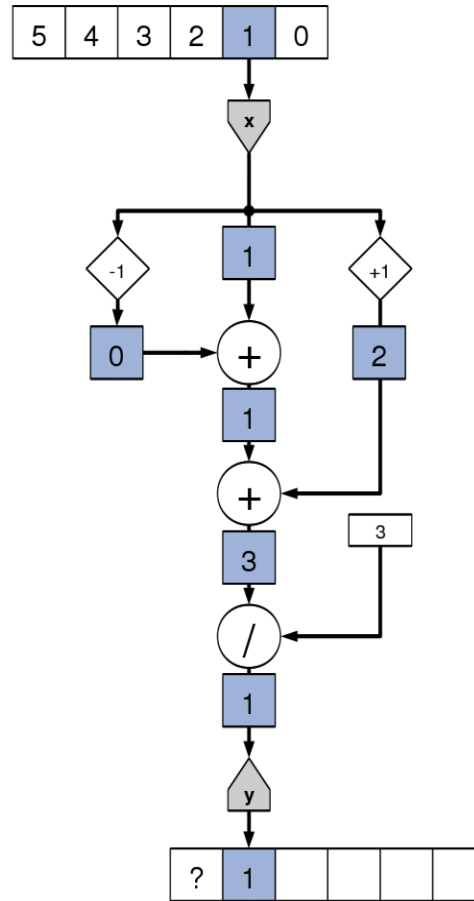
```
14 class MovingAverageSimpleKernel extends Kernel {  
15  
16     MovingAverageSimpleKernel(KernelParameters parameters) {  
17         super(parameters);  
18  
19         DFEVar x = io.input("x", dfeFloat(8, 24));  
20  
21         DFEVar prev = stream.offset(x, -1);  
22         DFEVar next = stream.offset(x, 1);  
23         DFEVar sum = prev + x + next;  
24         DFEVar result = sum / 3;  
25  
26         io.output("y", result, dfeFloat(8, 24));  
27     }  
28 }
```



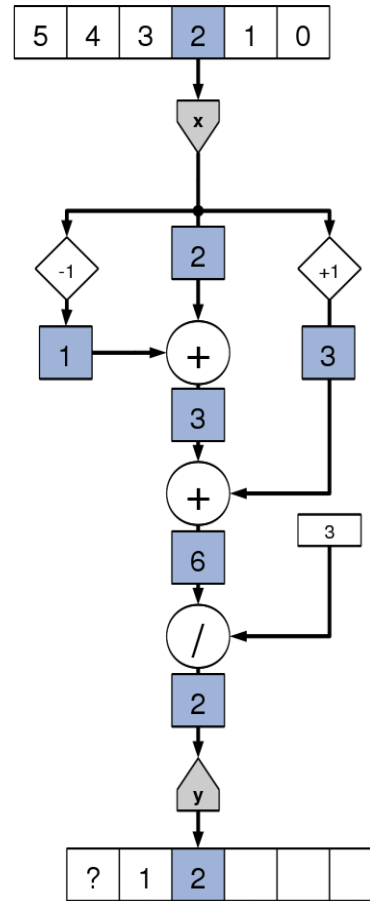
Kernel Execution



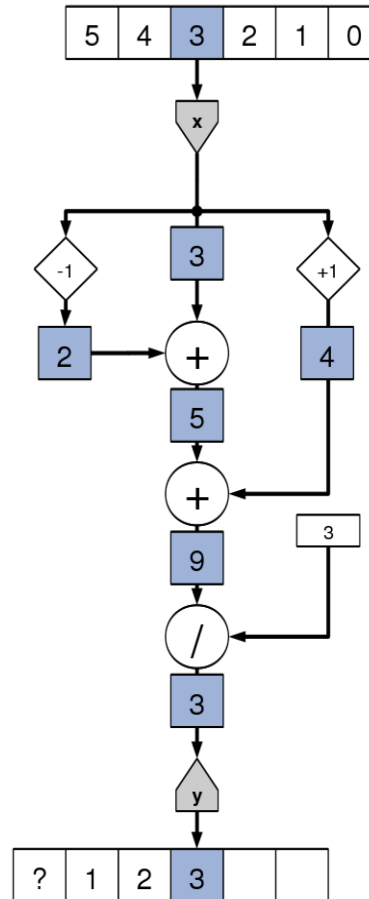
Kernel Execution



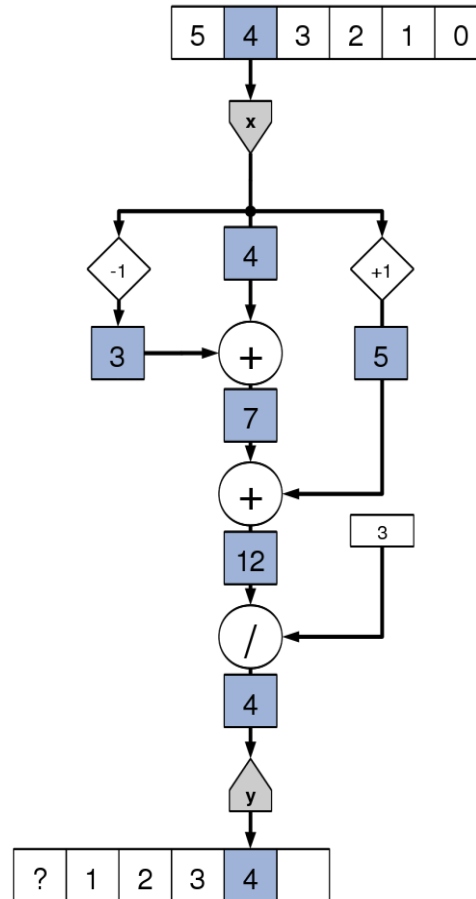
Kernel Execution



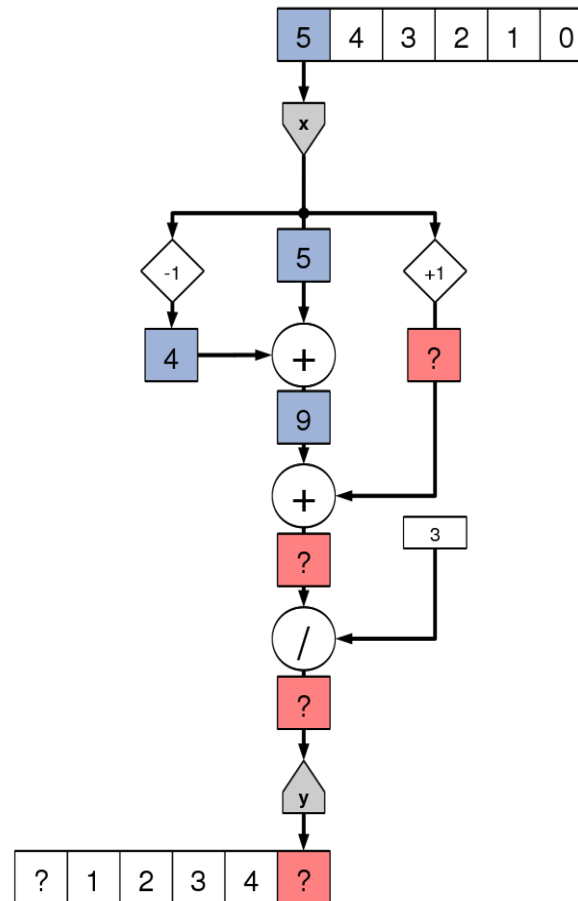
Kernel Execution



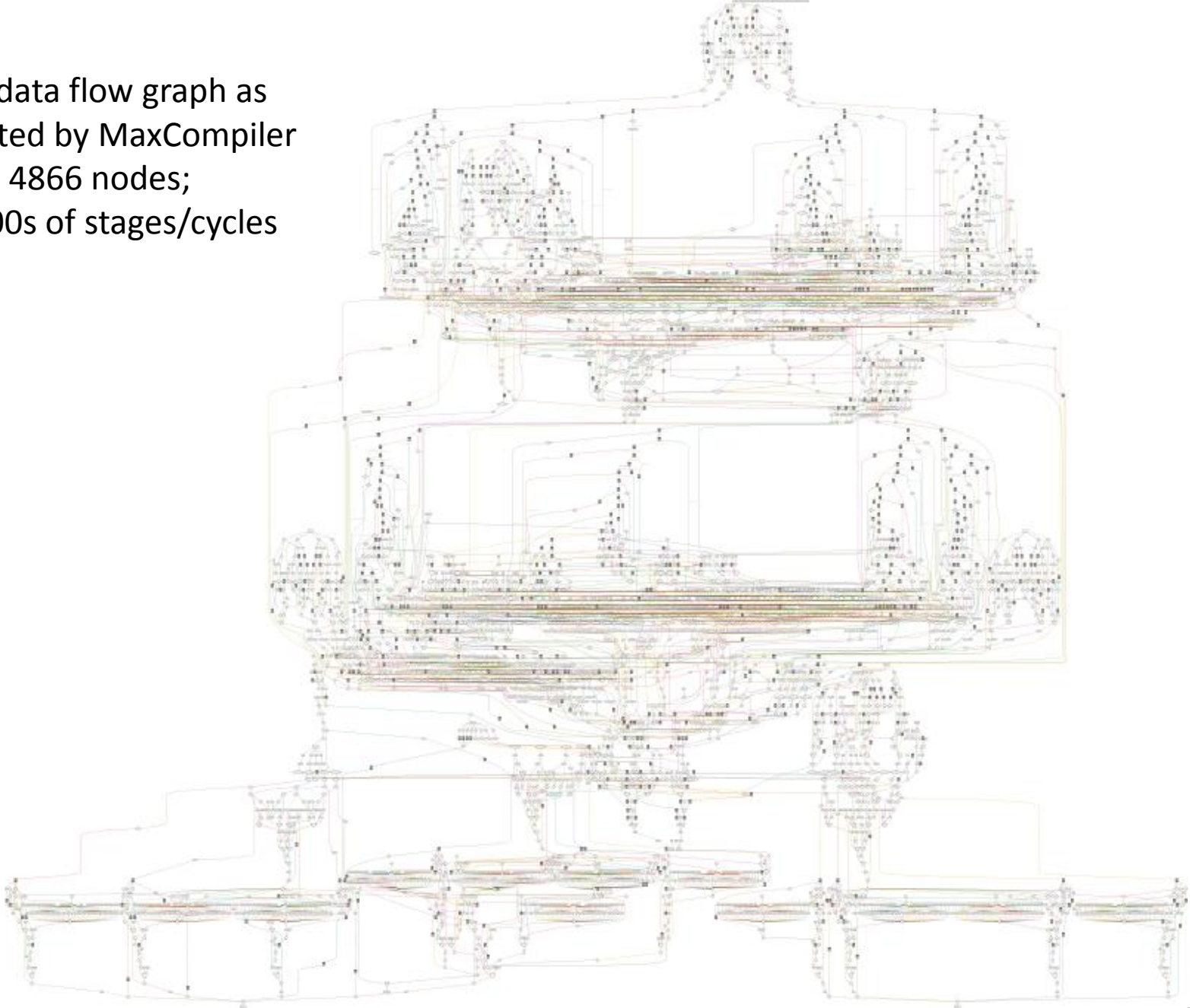
Kernel Execution



Kernel Execution



Real data flow graph as
generated by MaxCompiler
4866 nodes;
10,000s of stages/cycles



Maxeler University Program



Summary & Conclusions

- Tackling a vertical problem at multiple scales can allow you to make major jumps in capability
- Dataflow computing achieves high performance through:
 - Explicitly putting data movement at the heart of the program
 - Employing massive parallelism at low clock frequencies
 - Embodying application co-design
- Many scientific applications can benefit from this approach