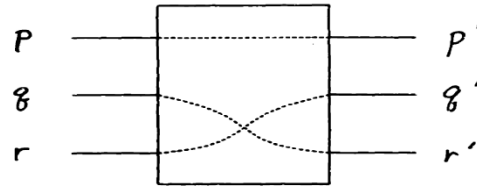
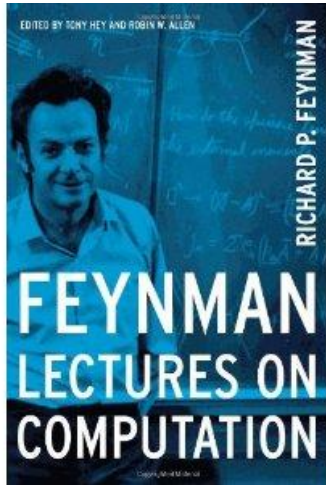


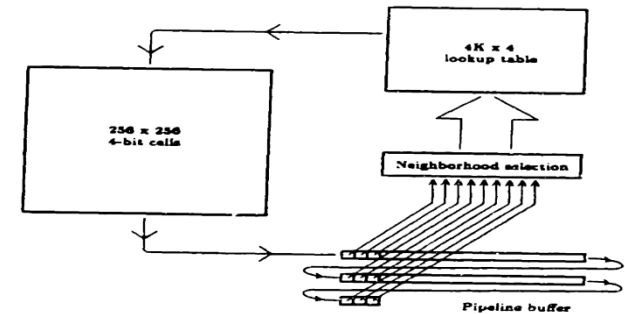
# Dataflow



# Some Dataflow History



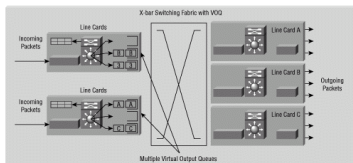
Toffoli => Fredkin Gate



CAM architecture

*PhD Thesis of Norman Margolus, MIT 1987*

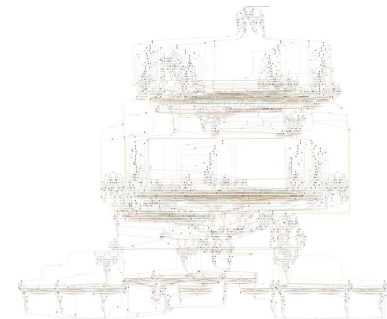
Today:



CISCO Internet Switch

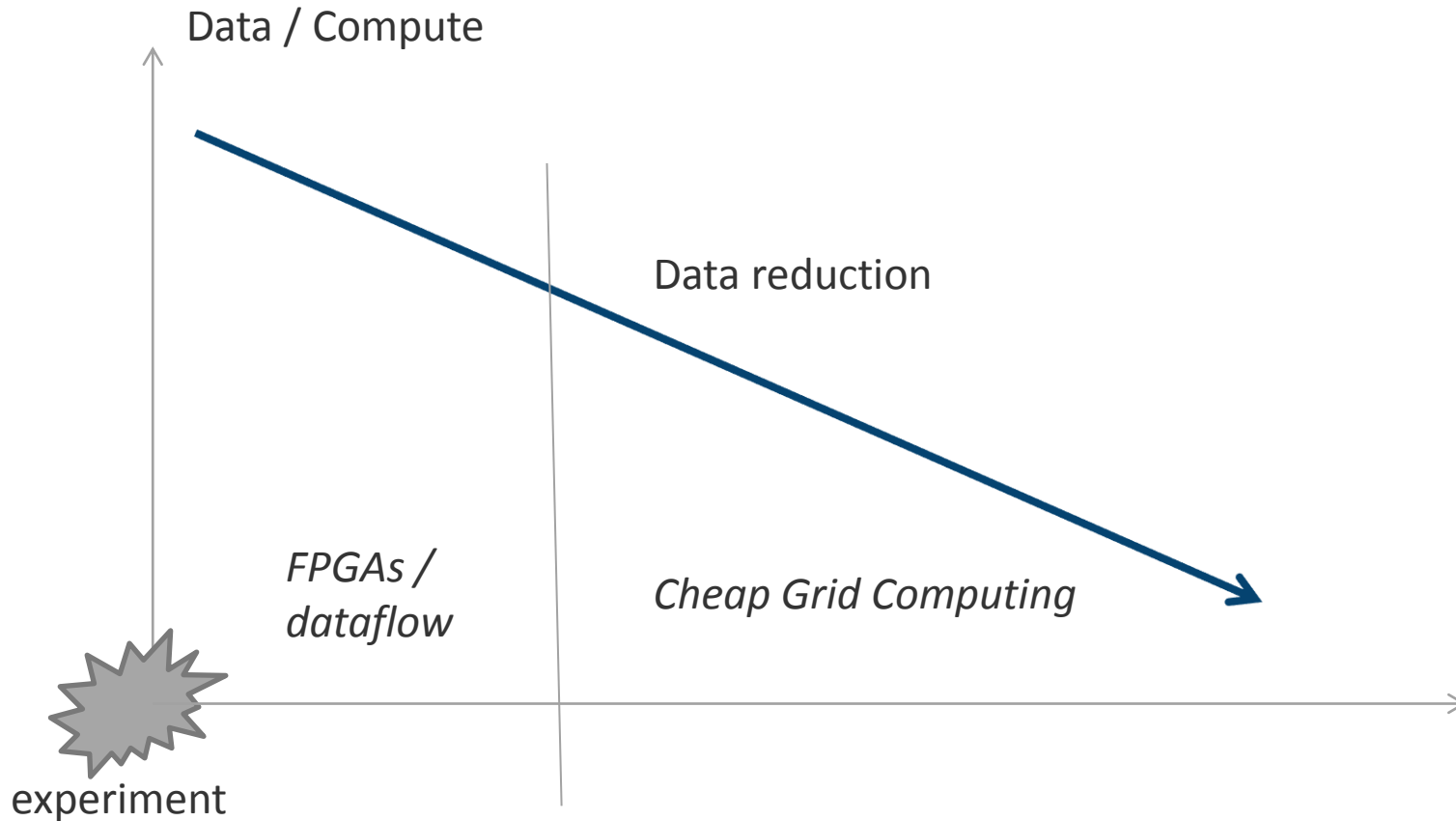


Data reduction in Physics



Maxeler Dataflow Engines

# CERN Large Hadron Collider



**“Is the LHC throwing away too much data?” New Scientist 2012.**

# Future of Medicine...

Sequencing

- Cancer: multiple genomes per tissue
- Spain: 28PB for 1 genome per cancer

Modelling

- Max Planck Institute:  $10^{18}$  core-hours for Europe
- Cancer, Diabetes, etc...

Big Data

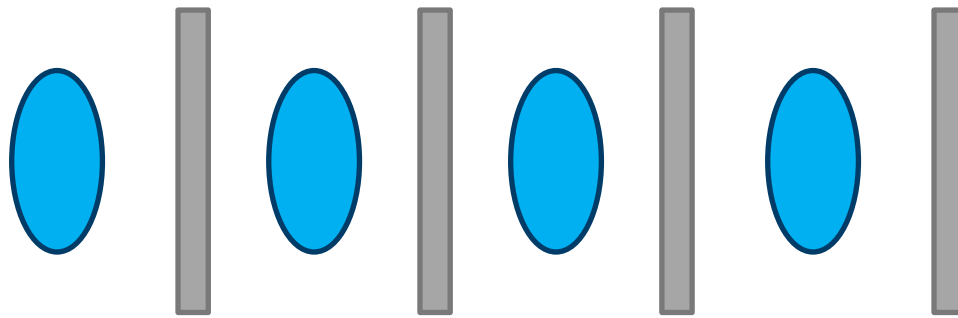
- Healthcare requires complete medical records
- Images, genomes, transcriptome, proteom, ...



# Technology

resulting from our research at Stanford, Bell Labs, and Imperial College

## MAXELER DATAFLOW COMPUTING



One result  
per clock cycle



Dynamic (switching) Power Consumption:

$$P_{avg} = C_{load} \cdot V_{DD}^2 \cdot f$$

Minimal frequency  $f$  achieves maximal performance, thus for a given power budget, we get Maximum Performance Computing (MPC)!

# Maxeler Architecture Series 2012

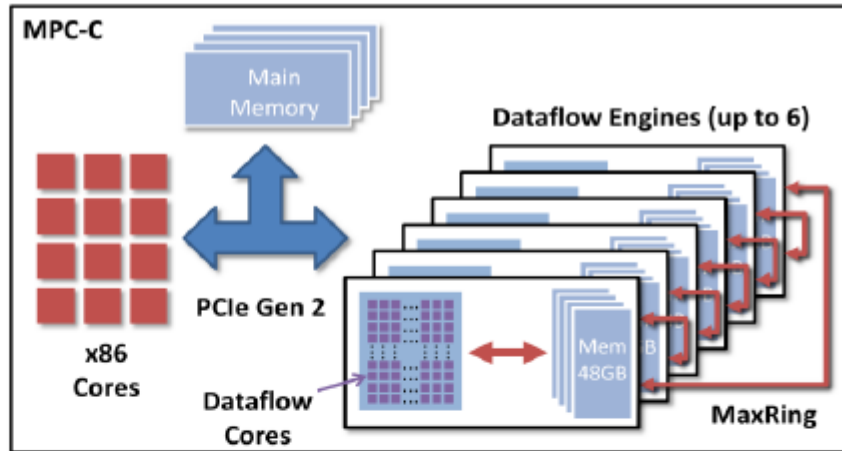


Figure 3. MPC-C Series Architecture

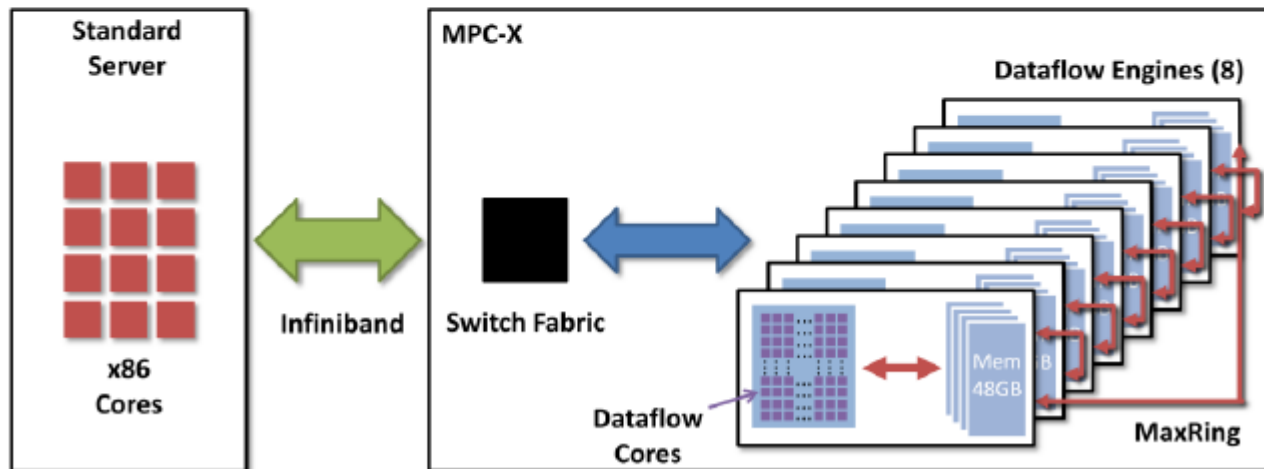
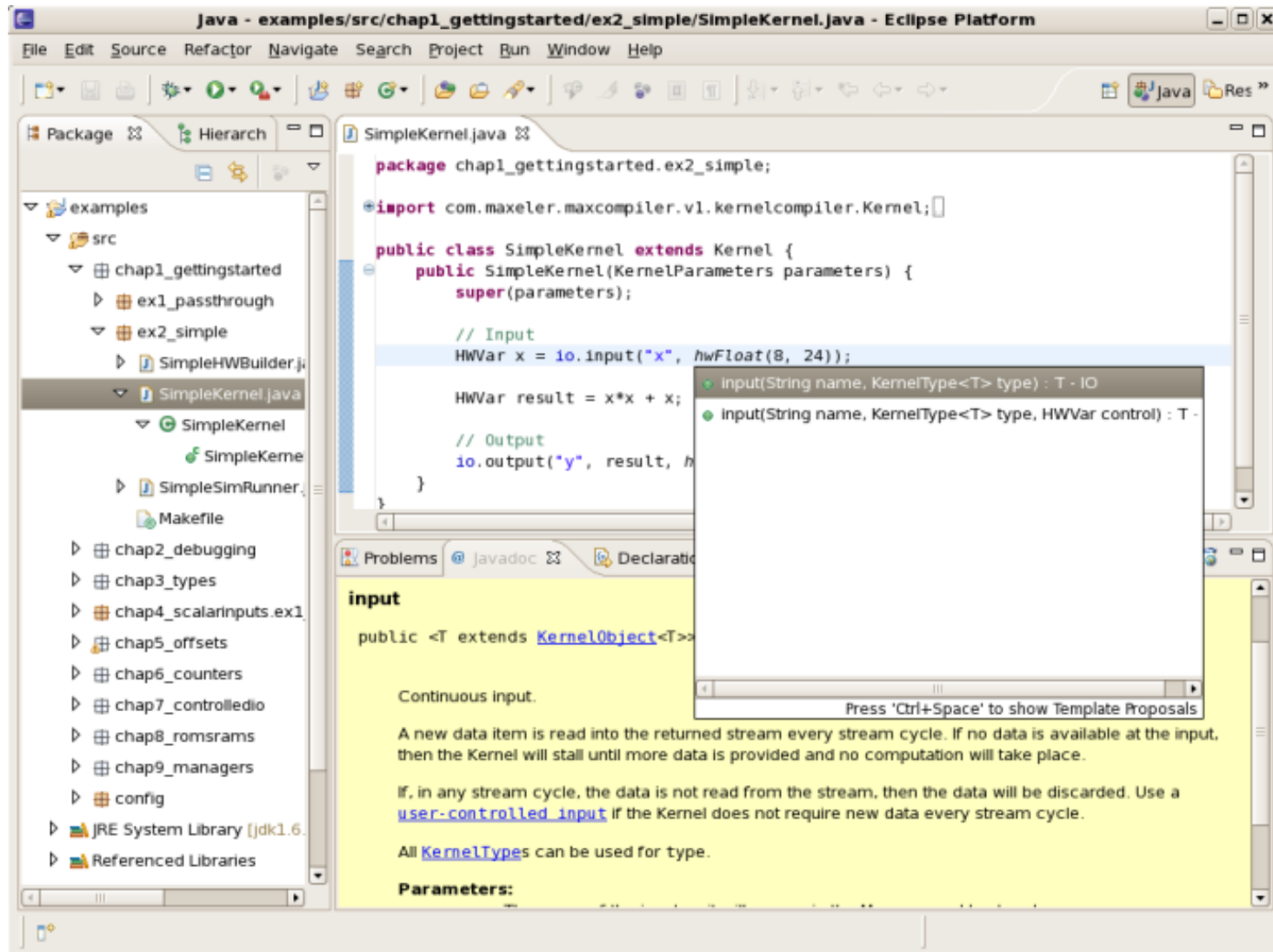


Figure 4. MPC-X Series Architecture

# Maxeler Programming Environment in Java



The screenshot shows the Eclipse IDE with the following components:

- Package Explorer:** Shows a project structure with packages like `examples`, `src`, `chap1_gettingstarted`, `ex1_passthrough`, `ex2_simple`, and `SimpleKernel.java`.
- Editor:** Displays the `SimpleKernel.java` file with the following code:

```
package chap1_gettingstarted.ex2_simple;

import com.maxeler.maxcompiler.v1.kernelcompiler.Kernel;

public class SimpleKernel extends Kernel {
    public SimpleKernel(KernelParameters parameters) {
        super(parameters);

        // Input
        HWVar x = io.input("x", hwFloat(8, 24));

        HWVar result = x*x + x;

        // Output
        io.output("y", result, h
    }
}
```
- Tooltip:** A tooltip is visible over the `io.input` call, showing the signature: `input(String name, KernelType<T> type) : T - IO` and `input(String name, KernelType<T> type, HWVar control) : T -`.
- Javadoc Window:** Shows the documentation for the `input` method:

```
input
public <T extends KernelObject<T>>
Continuous input.
A new data item is read into the returned stream every stream cycle. If no data is available at the input, then the Kernel will stall until more data is provided and no computation will take place.
If, in any stream cycle, the data is not read from the stream, then the data will be discarded. Use a user-controlled input if the Kernel does not require new data every stream cycle.
All KernelTypes can be used for type.
Parameters:
```



# Cox-Ingersol (CIR) Interest Rate Model

$$dr(t) = \alpha(b - r(t))dt + \sigma \sqrt{r(t)}dW(t)$$

$$r(t_{i+1}) = r(t_i) + \alpha(b(t_i) - r(t_i))[t_{i+1} - t_i] + \sigma \sqrt{r(t_i)} \sqrt{t_{i+1} - t_i} Z_{i+1}$$

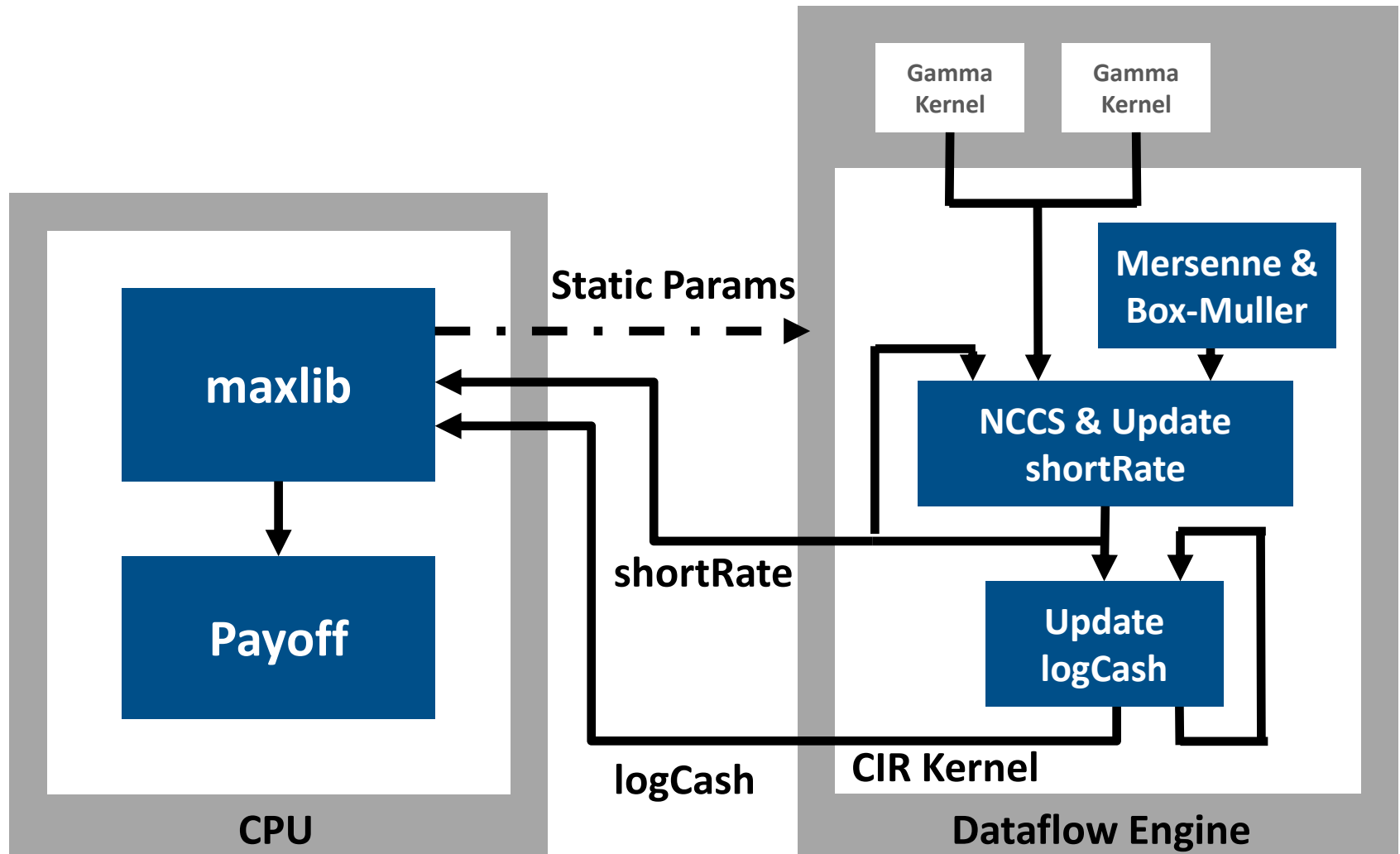
$$r(t_{i+1}) = \frac{\sigma^2(1 - e^{-\alpha(t_{i+1} - t_i)})}{4\alpha} \chi^2_d \left( \frac{4\alpha e^{-\alpha(t_{i+1} - t_i)}}{\sigma^2(1 - e^{-\alpha(t_{i+1} - t_i)})} r(t_i) \right), \quad d = 4\bar{b}\alpha/\sigma^2$$

$$B(t, T) = e^{-A(t, T)r(t) + C(t, T)}$$

$$A(t, T) = \frac{2(e^{\gamma(T-t)} - 1)}{(\gamma + \alpha)(e^{\gamma(T-t)} - 1) + 2\gamma}$$

$$C(t, T) = \frac{2\alpha b}{\sigma^2} \log \left( \frac{2\gamma e^{(\alpha + \gamma)(T-t)/2}}{(\gamma + \alpha)(e^{\gamma(T-t)} - 1) + 2\gamma} \right)$$

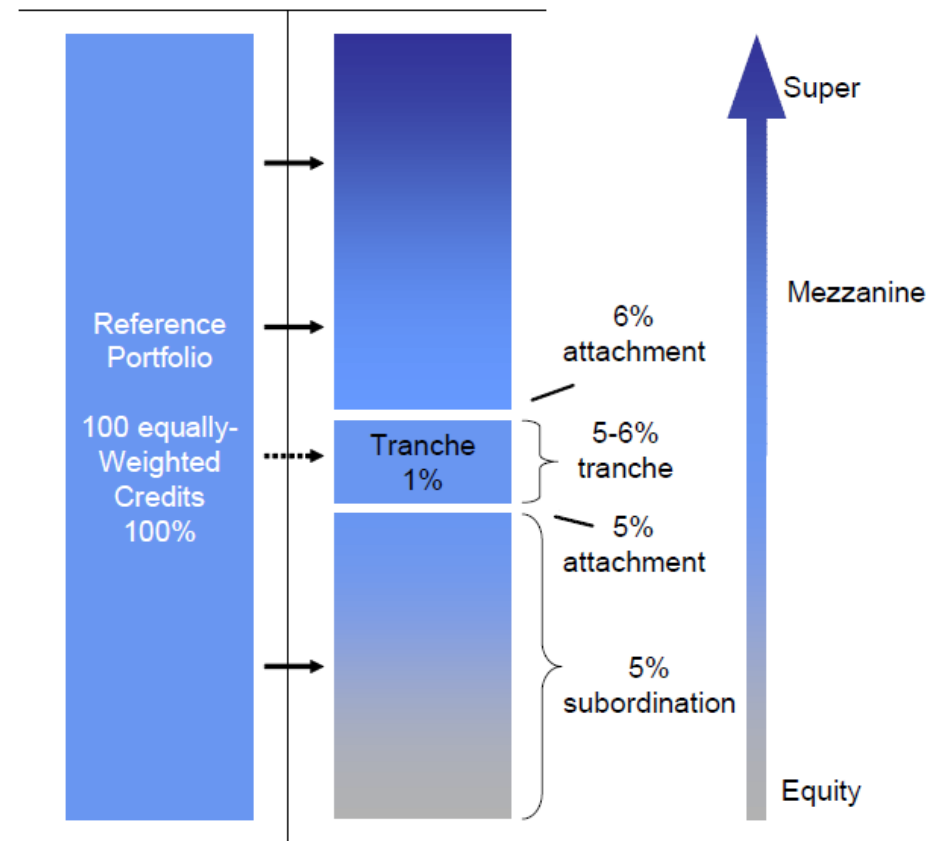
# CIR Monte Carlo Dataflow Architecture



# JP Morgan Credit Derivatives Pricing

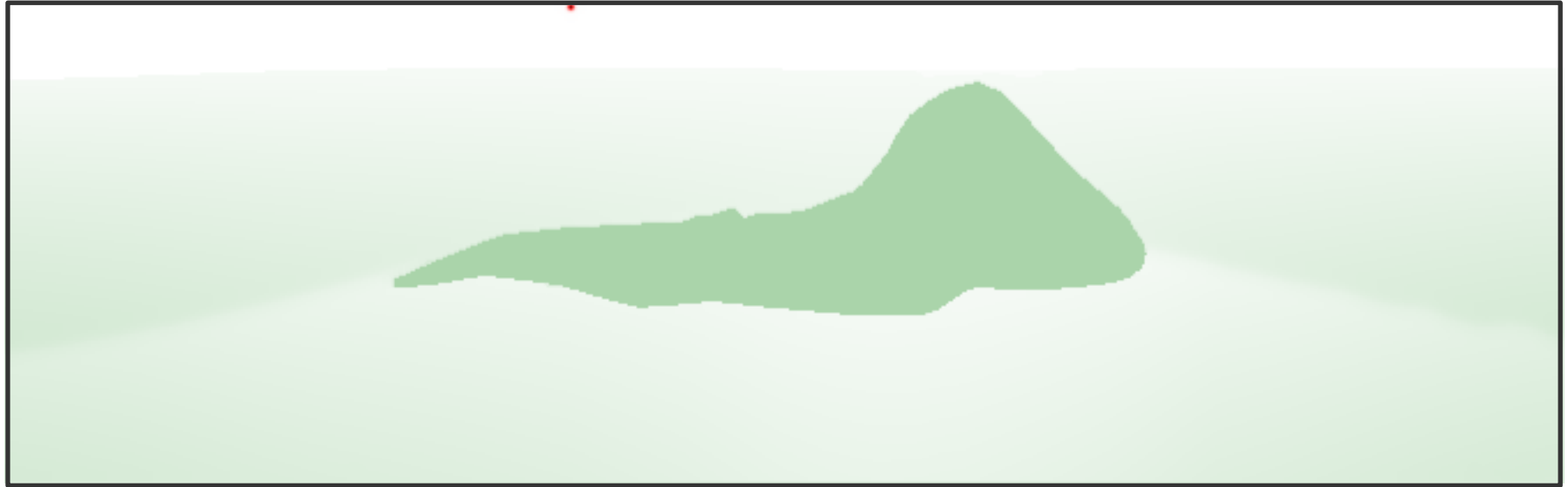
O. Mencer and S. Weston, 2010

- Compute value of complex financial derivatives (CDOs)
- Typically run overnight, but beneficial to compute in real-time
- Many independent jobs
- **Speedup: 220-270x**
- Power consumption per node drops from 250W to 235W/node



Source: JPMorgan

# Running with SEG Salt Model

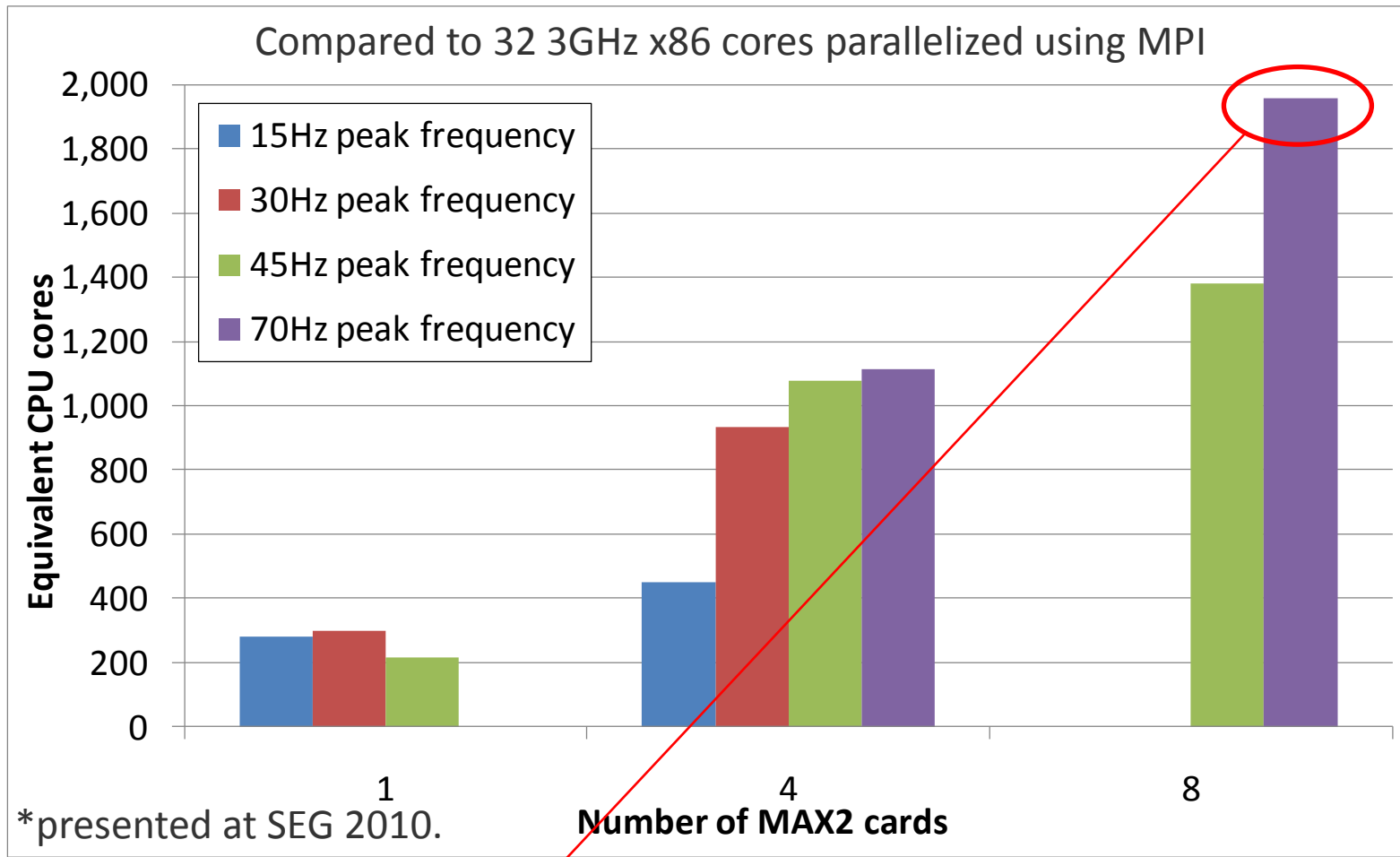


- Running on MaxNode servers
  - 8 parallel compute pipelines per chip
  - 150MHz => low power consumption!
  - 30x faster than microprocessors

**Beyond Traditional Microprocessors for Geoscience High-Performance Computing Applications** O. Lindtjorn<sup>†</sup>, R. G. Clapp<sup>†</sup>, O. Pell<sup>‡</sup>, O. Mencer<sup>‡</sup>, M. J. Flynn<sup>‡</sup> and H. Fu<sup>§</sup>.

<sup>†</sup>Stanford, <sup>‡</sup>Maxeler, <sup>§</sup>Tsinghua, IEEE Micro, March 2011.

# 3000<sup>3</sup> Modeling



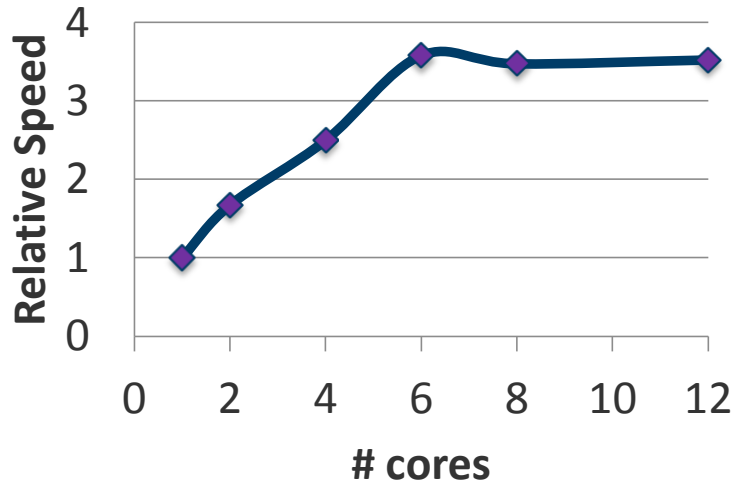
**8 Full Intel Racks ~100kWatts => 2 MaxNodes (2U) Maxeler System <1kWatt**

# Typical Scalability of Sparse Matrix

## Eclipse Benchmark

(2 node Westmere 3.06 GHz)

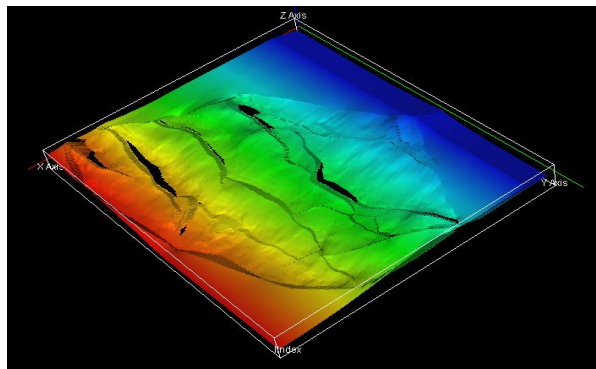
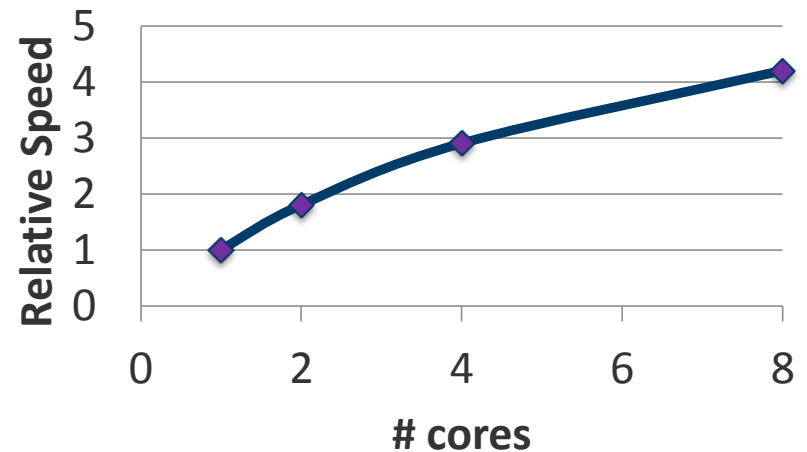
## E300 2 Mcell Benchmark



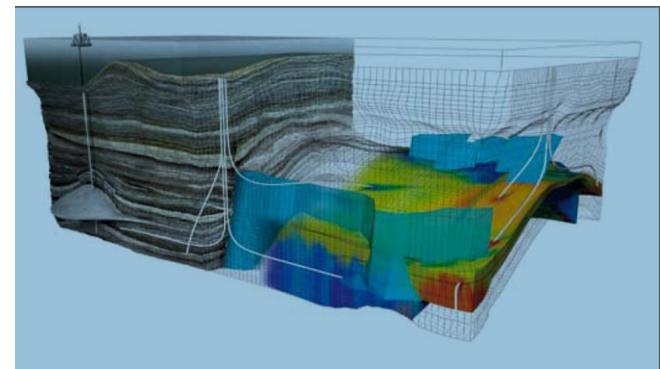
## Visage – Geomechanics

(2 node Nehalem 2.93 GHz)

## FEM Benchmark



Schlumberger



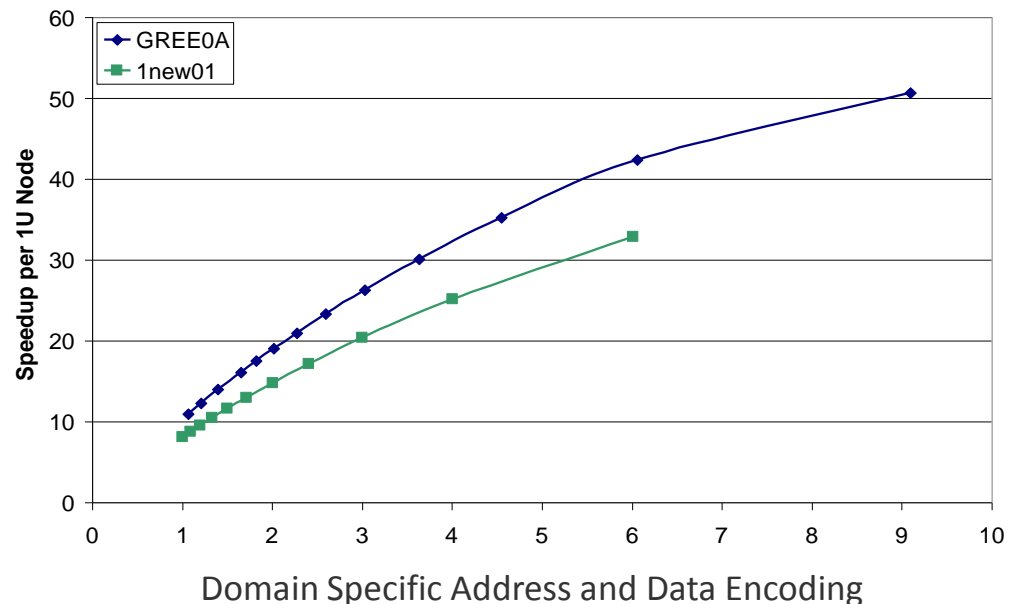
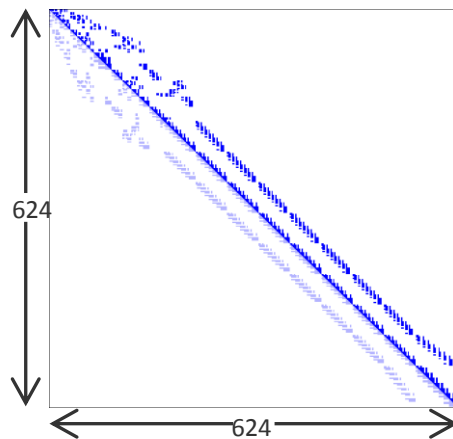
# Sparse Matrix Solving

O. Lindtjorn et al, 2010

- Given matrix  $\mathbf{A}$ , vector  $b$ , find vector  $x$  in:

$$\mathbf{Ax} = b$$

- Typically memory bound, not parallelisable.
- 1 MaxNode achieved 20-40x the performance of an Intel node.



# Economics of Computation (TCO)

## Example Solution 1

- » **50x** Speed-up per 1U server node
- » **32** Maxeler Node Solution
- » Equivalent to **1600** CPU-only Nodes
- » **\$3.2m** Operational cost savings over 3 years

## Example Solution 2

- » **30x** Speed-up per 1U server node
- » **40** Maxeler Node Solution
- » Equivalent to **1200** CPU-only Nodes
- » **\$1.8m** Operational cost savings over 3 years

## Example Solution 3

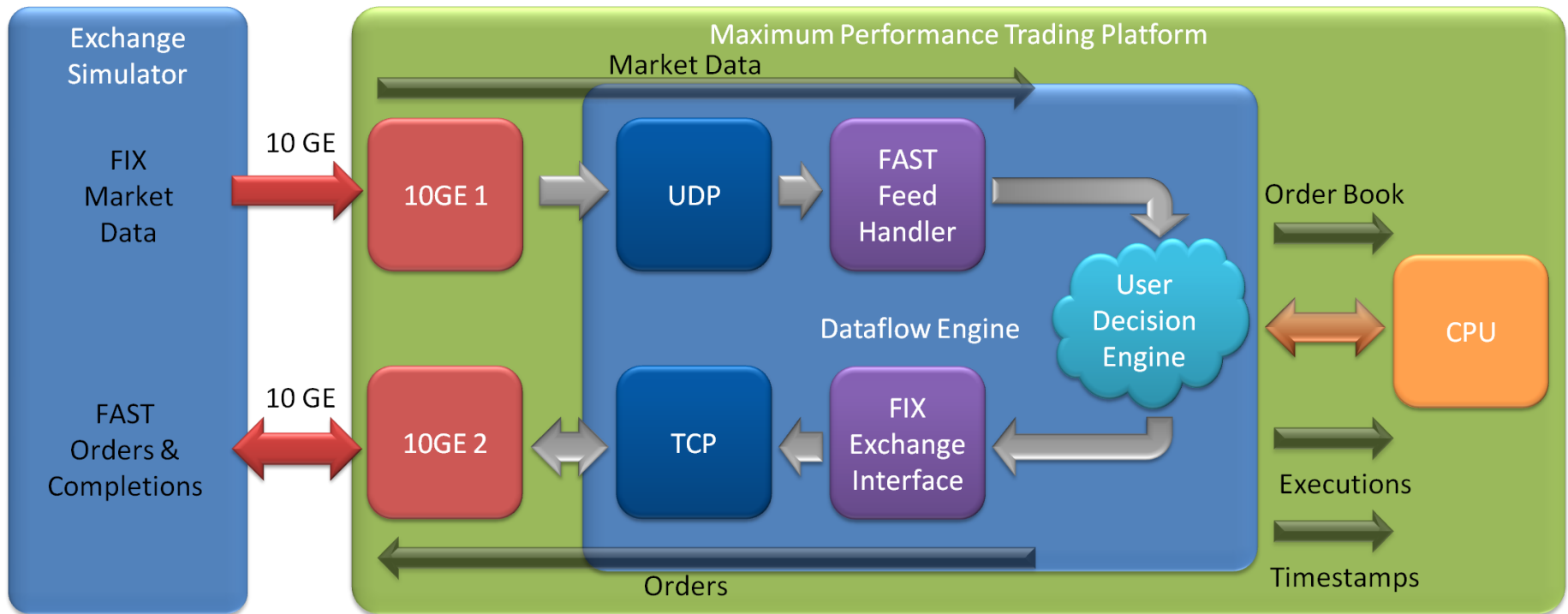
- » **20x** Speed-up per 1U server node
- » **50** Maxeler Node Solution
- » Equivalent to **1000** CPU-only Nodes
- » **\$1.7m** Operational cost savings over 3 years

## Example Solution 4

- » **40x** Speed-up per 1U server node
- » **32** Maxeler Node Solution
- » Equivalent to **1280** CPU-only Nodes
- » **\$2.6m** Operational cost savings over 3 years



# Trading: Direct Market Access and HFT



# Maxeler Hardware Solutions 2012



## CPU's plus DFEs

Intel Xeon CPU cores and up to 6 DFEs with 288GB 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 1-2 DFEs with up to six 10Gbit Ethernet connections



## MaxWorkstation

Desktop development system



## MaxCloud

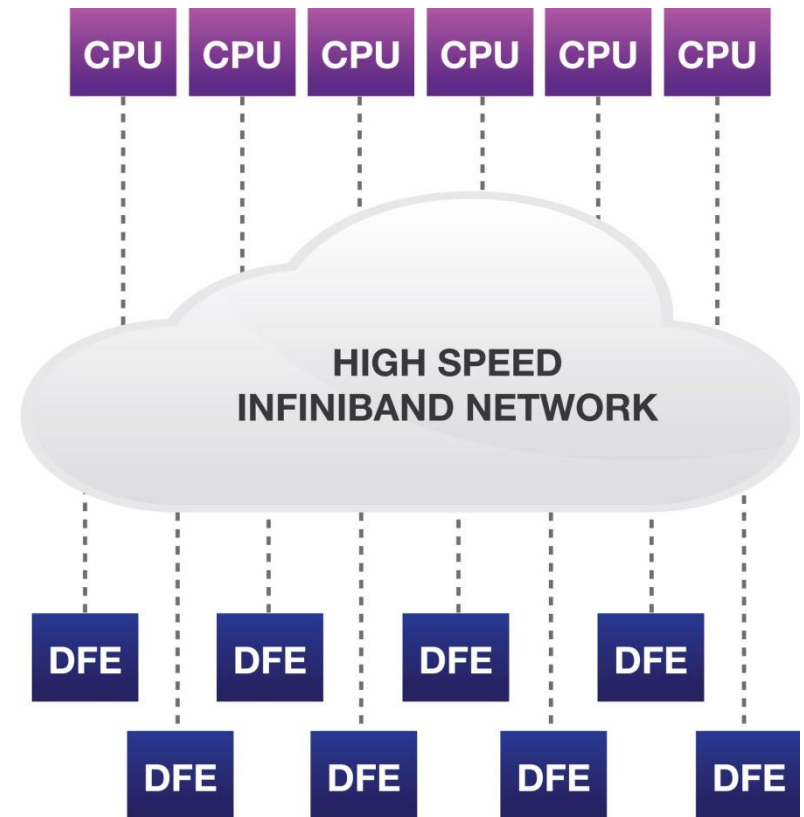
On-demand scalable accelerated compute resource, hosted in London



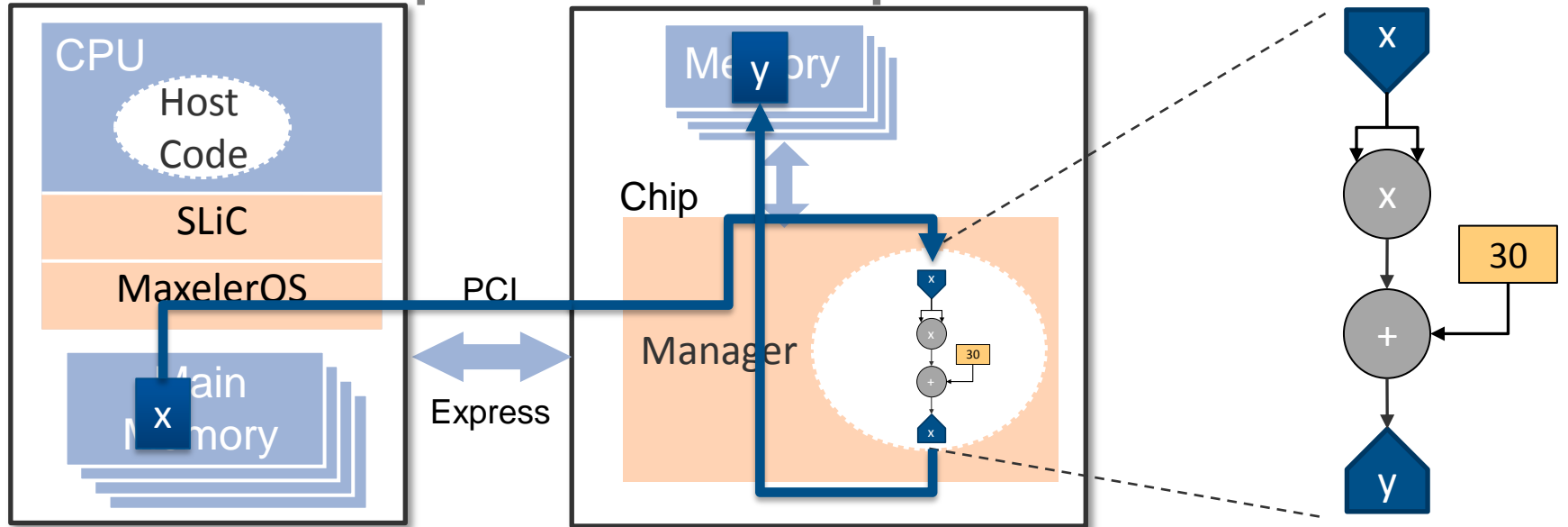
*1U dataflow cloud providing dynamically scalable compute capability over Infiniband*

## **MPC-X1000**

- 8 *vectis* dataflow engines (DFEs)
- 192GB of DFE RAM
- Dynamic allocation of DFEs to conventional CPU servers
  - Zero-copy RDMA between CPUs and DFEs over Infiniband
- Equivalent performance to 40-60 x86 servers



# MaxCompiler Development Process



CPUCode (.c)

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

Calc(x, DATA_SIZE)
```

Manager (.java)

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

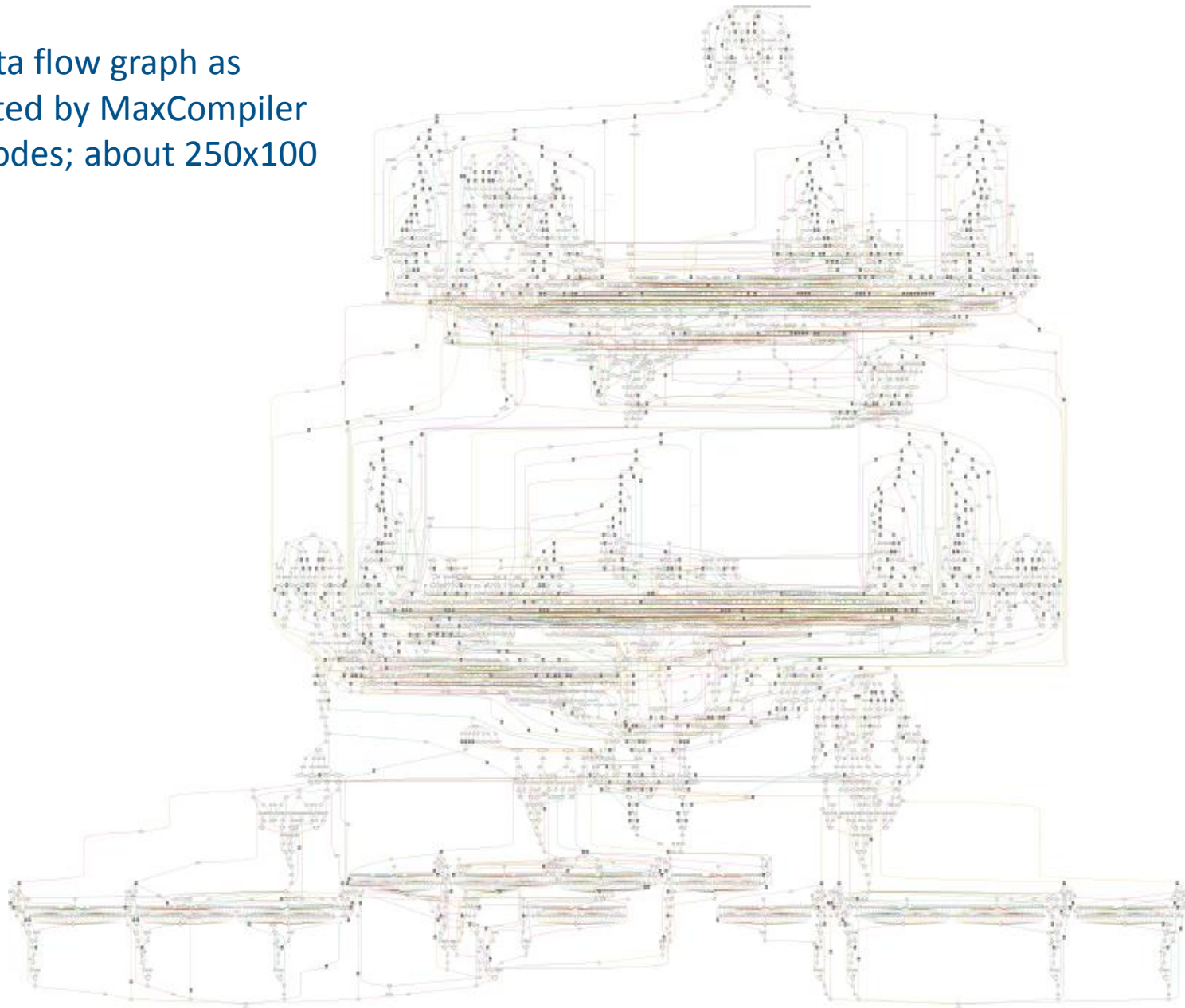
m.setKernel(k);
m.setIO(
    link("x", PCIE),
    link("y", DRAM_LINEAR1D));
m.addMode(modeDefault());
m.build();
```

MyKernel (.java)

```
HWVar x = io.input("x", hwInt(32));
HWVar result = x * x + 30;

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

Data flow graph as  
generated by MaxCompiler  
4866 nodes; about 250x100



# Chip Resource Usage

The goal is to maximize utilization of resources on the chip, and bandwidth on the memory bus.

LUTs	FFs	BRAMs	DSPs	: MyKernel.java
727	871	1.0	2	: resources used by this file
0.24%	0.15%	0.09%	0.10%	: % of available
71.41%	61.82%	100.00%	100.00%	: % of total used
94.29%	97.21%	100.00%	100.00%	: % of user resources
				:
				: public class MyKernel extends Kernel {
				: public MyKernel (KernelParameters parameters) {
				: super(parameters);
1	31	0.0	0	: HWVar p = io.input("p", hwFloat(8,24));
2	9	0.0	0	: HWVar q = io.input("q", hwUInt(8));
				: HWVar offset = io.scalarInput("offset", hwUInt(8));
8	8	0.0	0	: HWVar addr = offset + q;
18	40	1.0	0	: HWVar v = mem.romMapped("table", addr,
				: hwFloat(8,24), 256);
139	145	0.0	2	: p = p * p;
401	541	0.0	0	: p = p + v;
				: io.output("r", p, hwFloat(8,24));
				: }
				: }

# Measuring Utilization

- *Top* measures % of time CPU is running
- *Maxtop* monitors % of time the DFE is running

```
MaxTop Tool 2011.2
Found 2 Maxeler card(s) running MaxelerOS 2011.2
Card 0: MAX3A (P/N: 13424) S/N: 219270088 Mem: 24GB DFE(s): 1 /dev/maxeler0
Card 1: MAX3A (P/N: 13424) S/N: 000025559 Mem: 24GB DFE(s): 1 /dev/maxeler1

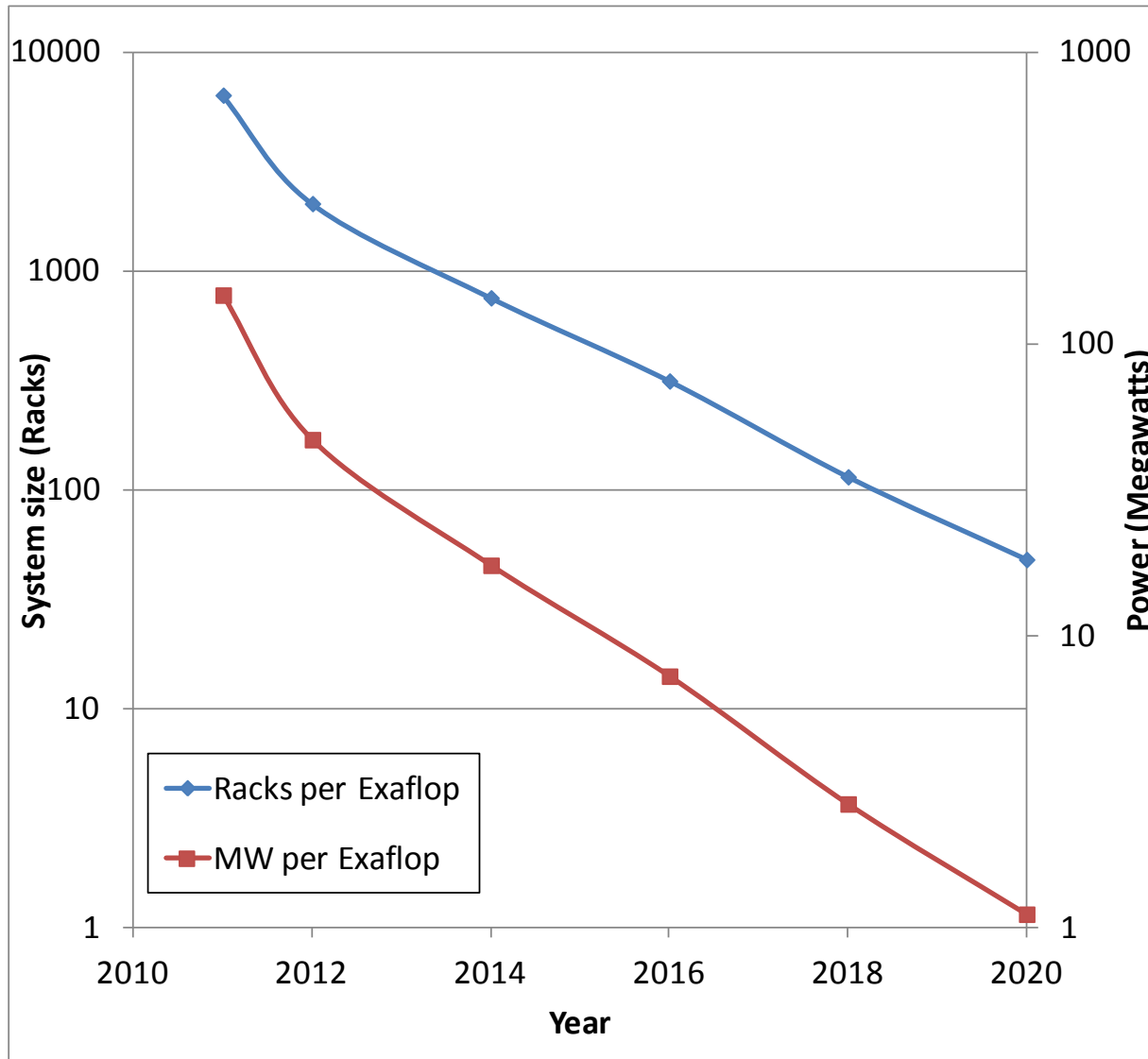
DEVICE      %DFE      TEMP      BITSTREAM      PID      USER      TIME      COMMAND
maxeler0    66.6%     57.1C     9d9de1...     12333    jspooner  00:00:39  model
maxeler1    0.0%     54.6C     9d9de1...     -        -        -        -
```

# The Exascale Supercomputer (2018)

- 1 exaflop =  $10^{18}$  FLOPS
- Using processor cores with 8FLOPS/clock at 2.5GHz
- **50M CPU cores**
- What about power?
  - Assume power envelope of 100W per chip
  - Moore's Law scaling: 6 cores today → ~100 cores/chip
  - 500k CPU chips
- **50MW (just for CPUs!) → 100MW likely**
- 'Jaguar' power consumption: 6MW for 2 Petaflops



# Exascale Dataflow Computing



- Exascale application performance by 2018
  - 100 racks
  - <3MW power
- A custom silicon solution could improve this by 5x
  - 20 racks
  - <1MW power

# Recent Scientific Publications with Clients

- **Rapid Computation of Value and Risk for Derivatives Portfolios**

S. Weston<sup>†</sup>, J. Spooner<sup>‡</sup>, S. Racanière<sup>‡</sup> and O. Mencer<sup>‡§</sup>.

<sup>†</sup>JPMorgan, <sup>‡</sup>Maxeler Technologies, <sup>§</sup>Imperial College London

Concurrency and Computation: Practice and Experience, Special Issue Paper, July 2007, doi: 10.1002/cpe.1778.

- **Beyond Traditional Microprocessors for Geoscience High-Performance Computing Applications**

O. Lindtjorn<sup>†</sup>, R. G. Clapp<sup>†</sup>, O. Pell<sup>‡</sup>, O. Mencer<sup>‡</sup>, M. J. Flynn<sup>‡</sup> and H. Fu<sup>§</sup>.

<sup>†</sup>Stanford University and Schlumberger, <sup>‡</sup>Maxeler Technologies, <sup>§</sup>Tsinghua University

IEEE Micro, vol. 31, no. 2, March/April 2011.

- **Accelerating the Computation of Portfolios of Tranched Credit Derivatives**

S. Weston<sup>†</sup>, J-T. Marin<sup>†</sup>, J. Spooner<sup>‡</sup>, O. Pell<sup>‡</sup> and O. Mencer<sup>‡</sup>.

<sup>†</sup>JPMorgan, <sup>‡</sup>Maxeler Technologies

IEEE Workshop on High Performance Computational Finance, New Orleans, USA, November 2010.

- **FD modeling beyond 70Hz with FPGA acceleration**

D. Oriato<sup>†</sup>, O. Pell<sup>†</sup>, C. Andreoletti<sup>‡</sup> and N. Bienati<sup>‡</sup>. <sup>†</sup>Maxeler Technologies, <sup>‡</sup>Eni E&P Division

SEG 2010 HPC Workshop, Denver, USA, October 2010.

- **Surviving the End of Scaling of Traditional Microprocessors in HPC**

O. Lindtjorn<sup>†‡</sup>, R. G. Clapp<sup>‡</sup>, O. Pell<sup>§</sup>, O. Mencer<sup>§</sup> and M. J. Flynn<sup>§</sup>. <sup>†</sup>Schlumberger, <sup>‡</sup>Stanford University, <sup>§</sup>Maxeler Technologies

IEEE HOT CHIPS 22, Stanford, USA, August 2010.

- **Fast 3D ZO CRS Stack - An FPGA Implementation of an Optimization Based on the Simultaneous Estimate of Eight Parameters**

P. Marchetti<sup>†</sup>, D. Oriato<sup>‡</sup>, O. Pell<sup>‡</sup>, A.M. Cristini<sup>§</sup> and D. Theis<sup>§</sup>.

<sup>†</sup>Eni E&P Division, <sup>‡</sup>Maxeler Technologies, <sup>§</sup>CRS4

72nd European Association of Geoscientists and Engineers (EAGE) Conference, Barcelona, June 2010.

- **Computational acceleration of credit and interest rate derivatives**

O. Mencer<sup>†</sup> and S. Weston<sup>‡</sup>. <sup>†</sup>Maxeler Technologies, <sup>‡</sup>JPMorgan

Global Derivatives Trading and Risk Management, Paris, May 2010.

# MAX-UP: Maxeler University Program

- Founding Members: Imperial, Tsinghua, U Tokyo and Stanford



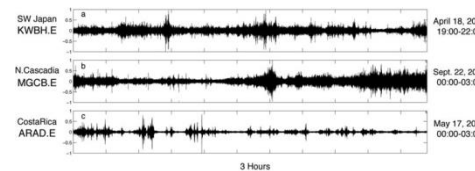
MaxWorkstation



- MaxAcademy course material for teaching

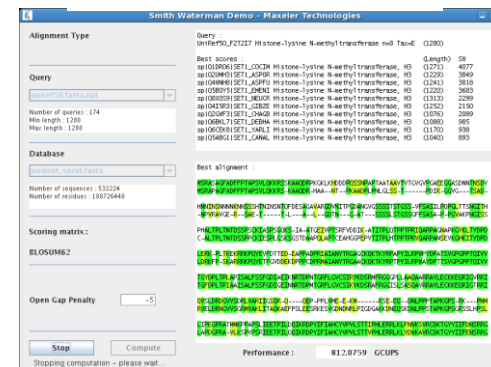
- MaxTraining at Maxeler: learn dataflow programming

## Earthquake Prediction



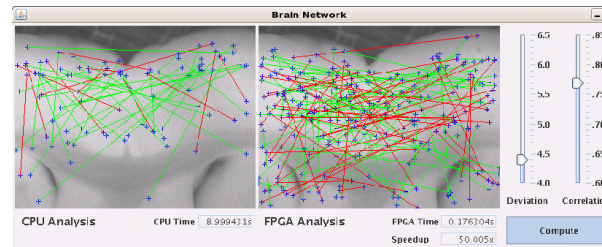
Smith-Waterman

- Internships and collaboration



- Joint research grant applications

## Brain Network



- Demos