THE UNIVERSITY
of EDINBURGH

ACCELERATING FINANCIAL MODELS

Christos Delivorias¹, Erik Vynckier¹, Peter Richtárik², Martin Takáč²¹Scottish Widows Investment Partnership, 60 Morrison Street, Edinburgh EH3 8BE, UK²School of Mathematics, University of Edinburgh, James Clerk Maxwell Building, Edinburgh EH9 3JZ, UK
Christos.Delivorias@swip.com, Erik.Vynckier@swip.com, Peter.Richtarik@ed.ac.uk, M.Takac@sms.ed.ac.uk

PERFORM

CONTRIBUTION:

Implement Heston's stochastic volatility model with Andersen's QE discretisation using Quasi-random numbers for variance reduction.**Accelerate** this model using:

CUDA GPGPU AWS HPC	Matlab PCT GPGPU	Techila Azure Cloud	FPGA
2 Tesla M2050 GPU	2 Tesla M2090 GPU	356 cores @1.6GHz	4 x Max3 cards

Evaluate Improvement in acceleration.

QE-Scheme DISCRETISATION

Leif Andersen [1] proposed a new scheme to discretise the stochastic volatility and the price of an underlying asset.

Require: The present value for the variance, $\hat{V}(t)$ **Ensure:** The value for the variance in the subsequent time-step, $\hat{V}(t + \Delta t)$

- 1: Compute $m \leftarrow \theta + (\hat{V}(t) - \theta)e^{-\kappa\Delta t}$
- 2: Compute $s^2 \leftarrow \frac{\hat{V}(t)\xi^2 e^{-\kappa\Delta t}}{\kappa}(1 - e^{-\kappa\Delta t}) + \frac{\theta\xi^2}{2\kappa}(1 - e^{-\kappa\Delta t})^2$
- 3: Compute $\psi \leftarrow \frac{m^2}{s^2}$
- 4: **if** $\psi \leq \psi_c$ **then**
- 5: Compute $b \leftarrow 2\psi^{-1} - 1 + \sqrt{2\psi^{-1} - 1}$
- 6: Compute $a \leftarrow \frac{m}{1+b^2}$
- 7: Generate Normal random variate Z_V
- 8: **return** $\hat{V}(t + \Delta t) \leftarrow a(b + Z_V)^2$
- 9: **else**
- 10: Compute $p \leftarrow \frac{\psi-1}{\psi+1} \in [0, 1)$
- 11: Compute $\beta \leftarrow \frac{1-p}{m}$
- 12: Generate Uniform random variate U_V
- 13: **if** $0 \leq U_V \leq p$ **then**
- 14: **return** $\hat{V}(t + \Delta t) \leftarrow 0$
- 15: **else**
- 16: **return** $\hat{V}(t + \Delta t) \leftarrow \beta^{-1} \ln\left(\frac{1-p}{1-U_V}\right)$
- 17: **end if**
- 18: **end if**

BIBLIOGRAPHY

- [1] L. Andersen, Efficient Simulation of the Heston Stochastic Volatility Model. Available at SSRN: <http://ssrn.com/abstract=946405>, 2007.
- [2] G. Levy, An introduction to quasi-random numbers. <http://bit.ly/13AUKVj>, 2002.
- [3] H. Niederreiter, Random number generation and Monte Carlo method. In SIAM, 1992.
- [4] C. Delivorias, <https://github.com/cmdel/mc-sim>, 2012.

Quasi-Random VARIATES

Such numbers can be sampled from the so called "low discrepancy" sequences. A sequence's discrepancy is a measure of its uniformity and is defined by the following definition [see Levy[2]].

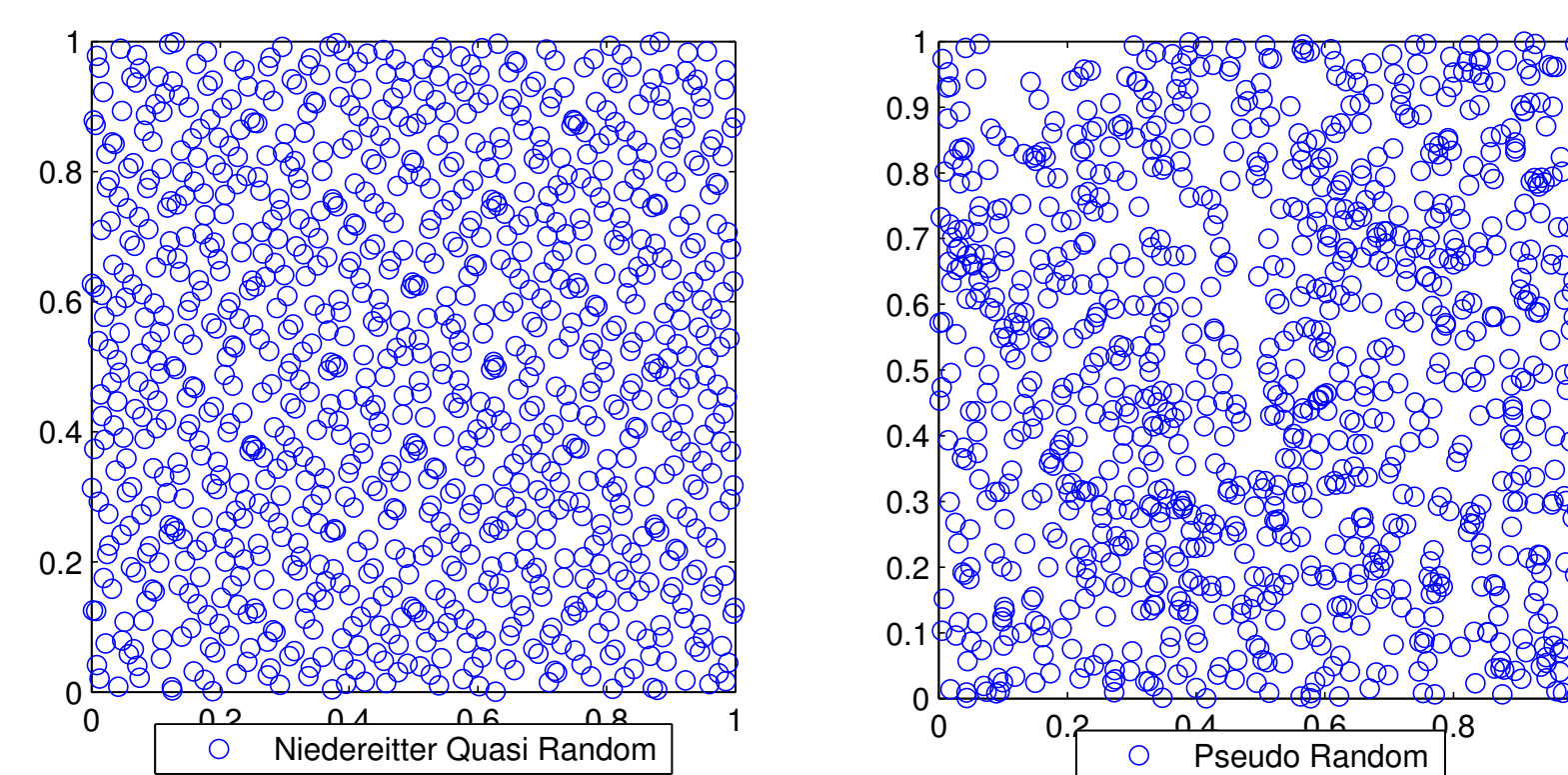
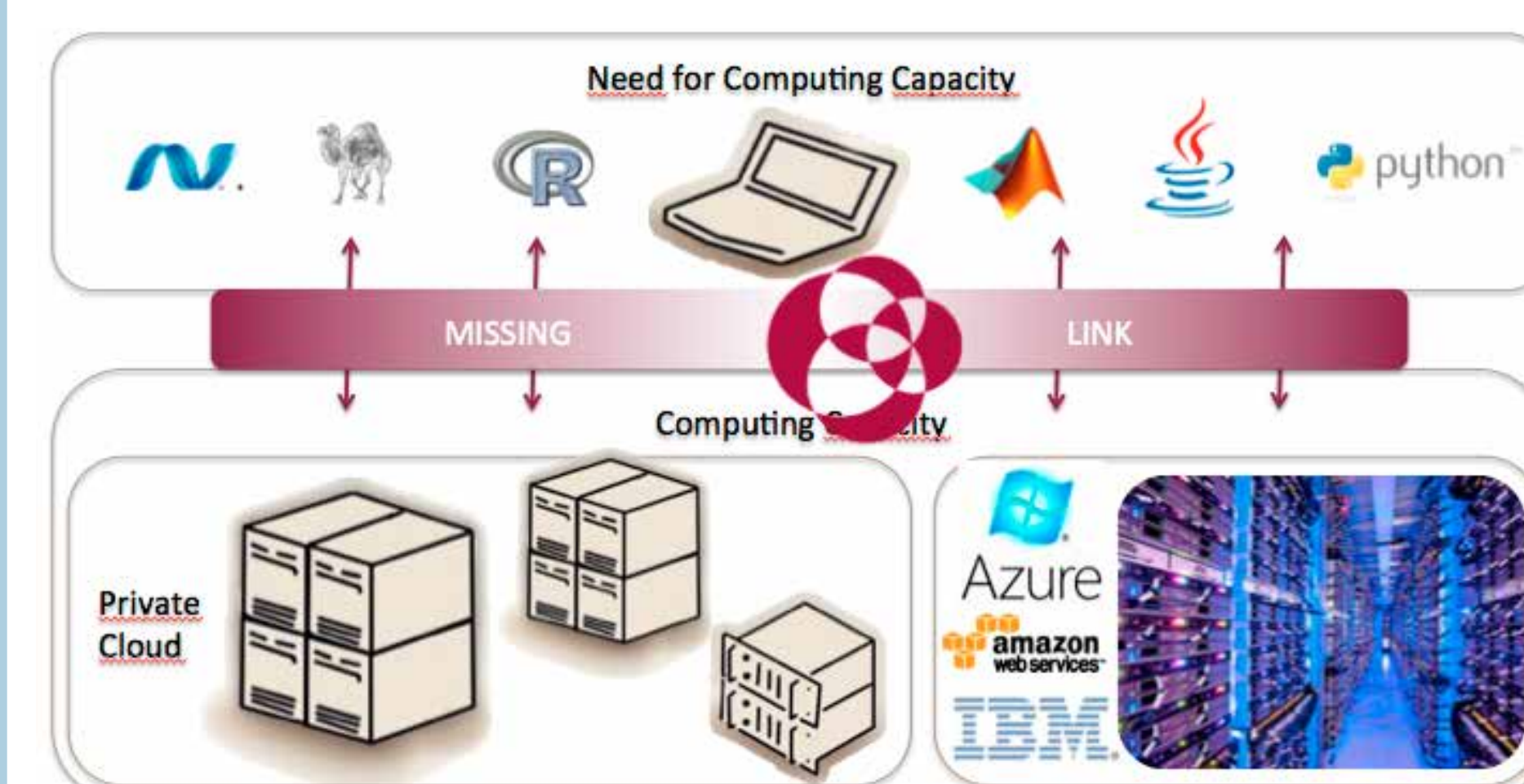


Figure 1: The Niederreiter [3] sequence is generated in MATLAB with the NAG functions g05yl and g05ym

TECHILA

Techila Technologies provided a middleware layer for parallel distribution. The grid comprised of 354 cores clocked at 1.6GHz on Microsoft's Azure cloud. The only change required in the code was the substitution of the *for* command with the *cloudfor* one.

```

1 % Main Monte Carlo loop
2 cloudfor pth = 1: paths
3 %cf:force:largedata
4 (...)
5     Payoff(pth) = max(S(pth,end)-K,0);
6     % Call option
7 (...)
8 cloudend

```

PROJECT CONTRIBUTORS



MATLAB GPGPU

The Heston model was ported to the GPU architecture by use of Matlab's Parallel Computing Toolbox. All functions that are executed on the GPGPU card are within the *gpuArray* namespace.

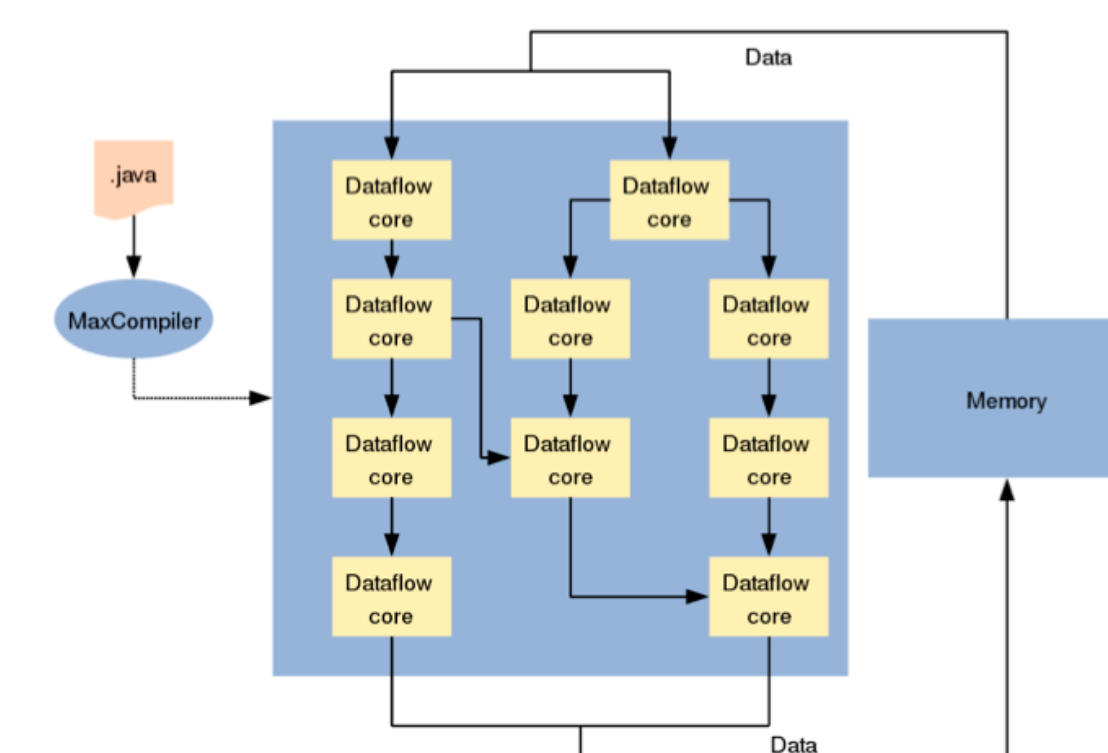
```

1 (...)
2 % Assign each worker to a different GPU
3 spmd
4     gpuDevice(labindex);
5 end
6 (...)
7 % Main Monte Carlo loop
8 for step = 1: steps
9     Zn1 = parallel.gpu.GPUArray.randn(
10         paths,1);
11     Zn2 = parallel.gpu.GPUArray.randn(
12         paths,1);
13     Uv = parallel.gpu.GPUArray.rand(paths,1);
14     [S,V]=arrayfun(@myGPUFun, Zn1,Zn2,Uv,
15         'S',V,theta,kappa,dt,xi,r,K0,
16         K1,K2,K3,K4);
17 end
18 (...)

```

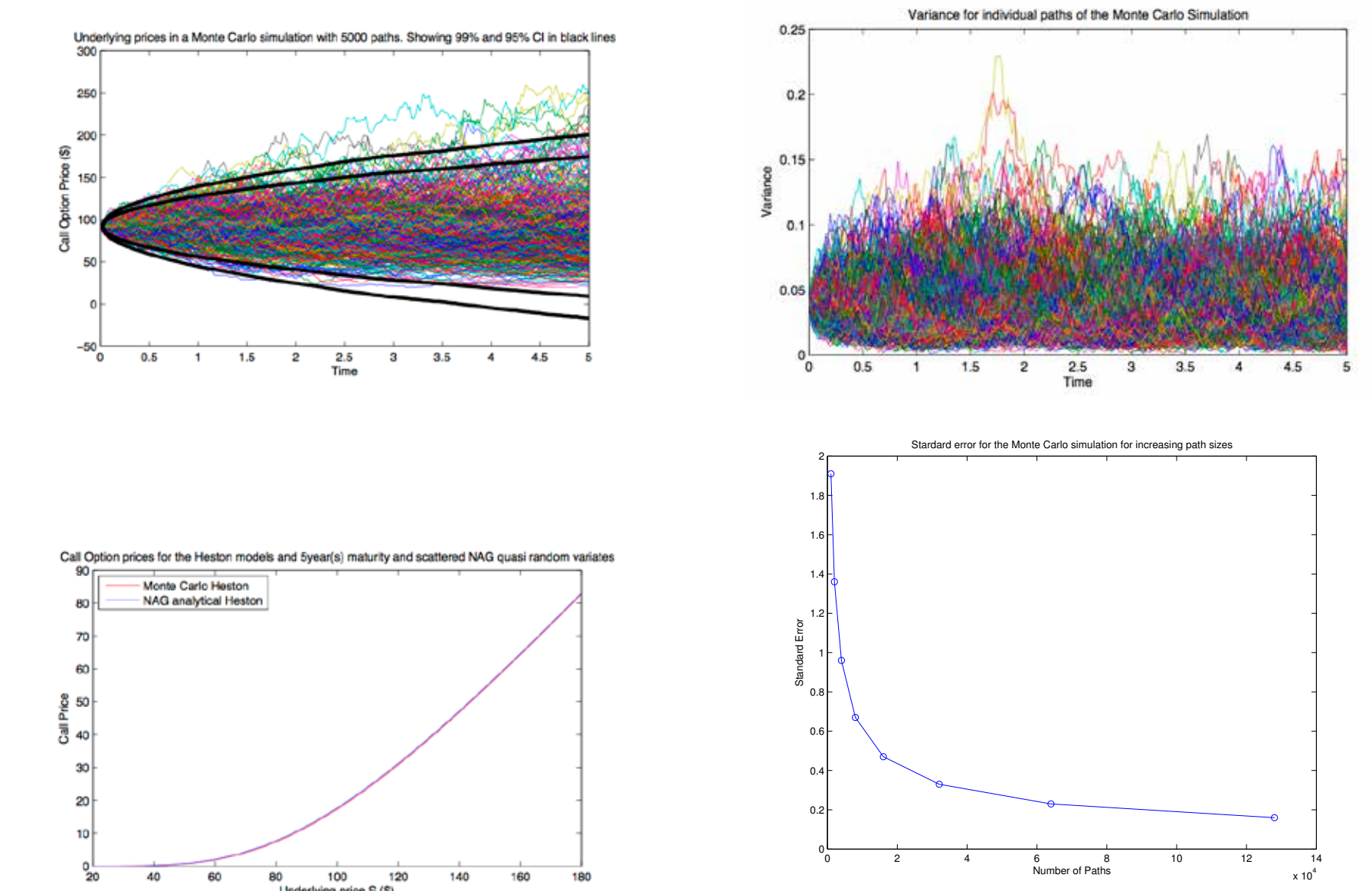
DATAFLOW ON FPGA

What Maxeler provide is a way to create workers on an FPGA that are highly specialised, and extremely quick at conducting a specific operation. The kernels are large synchronous dataflow pipelines that implement the mathematics and the control of the problem. They are asynchronously coupled to other kernels and I/O sources and sinks (DRAM, PCIe, inter-chip links,) by the manager. The overall process is as quick as the flow.



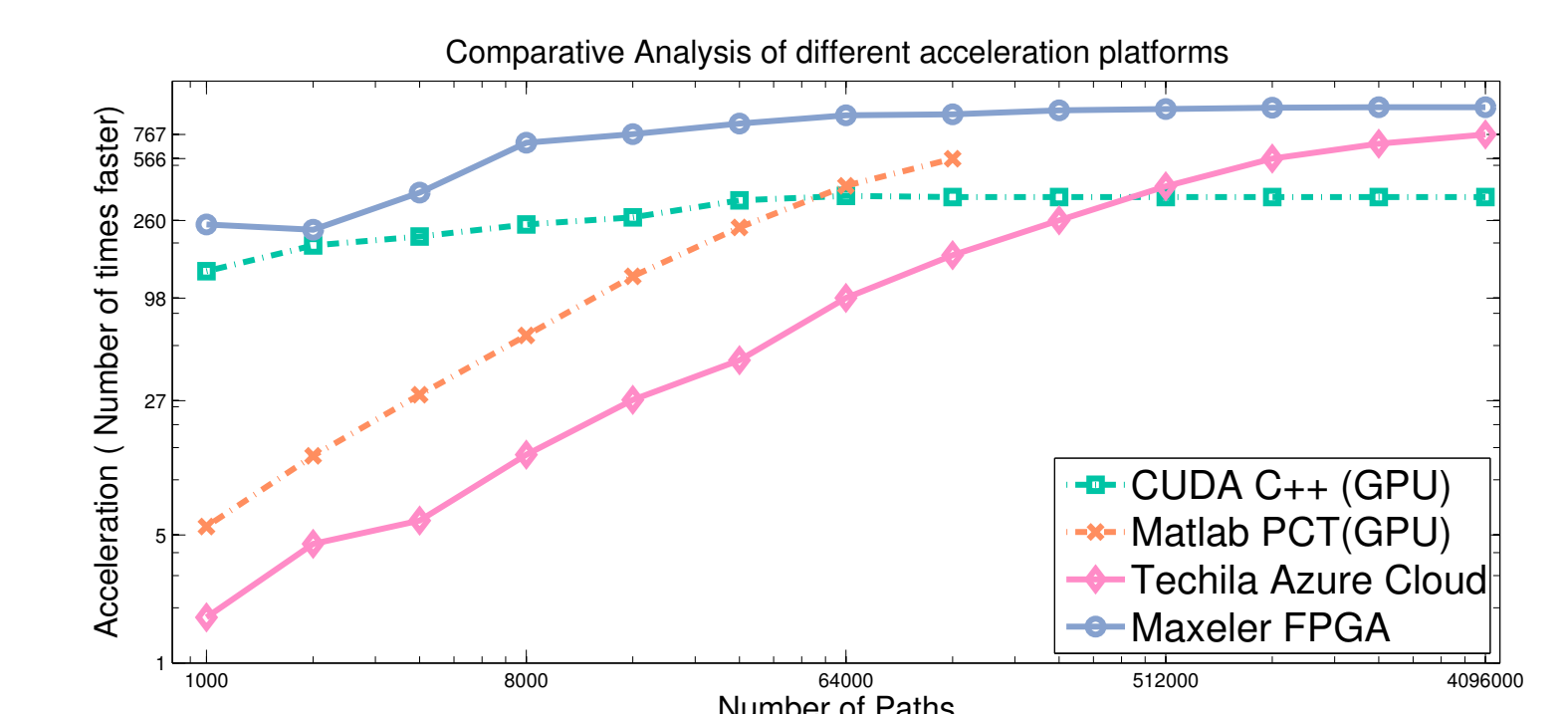
PRECISION and ACCURACY

The QE-scheme has high accuracy, compared to the analytic price calculated with the NAG function, and converges stably and quickly to a low standard error. The asset's variance is strictly positive and strongly repellant from zero.

Figure 2: Matlab figures for the Monte Carlo simulation of the Heston model. The parameters for this simulation are: $S_0 = 100$, $V_0 = 0.04$, $K = 100$, $\rho = -0.2$, $\xi = 0.2$, $\theta = 0.04$, $\kappa = 1.5$, $T = 5.0$, $NoSteps = T * 250$, $NoPaths = 1000$, $\lambda = 0$, $r = 0.0$, and $q = 0.0$. For the bottom two plots the S price is perturbed in the prices $S_0 = [20 \ 30 \ \dots \ 100 \ \dots \ 170 \ 180]$.

RESULTS

All platforms seem to converge to the same acceleration boundary. This is however with a different cost for each one of them. Total cost of ownership and ease of development should be taken into account. For code samples see [4]



CONCLUSIONS

1. FPGAs fill the pipe early on and plateau at a certain acceleration; needs to scale horizontally with additional cards
2. FPGAs had the best acceleration; high development friction
3. Easiest implementation was with Techila, followed by Matlab PCT
4. Performance via middle-ware suffers at small path number sizes, but catches up eventually