



Facultad de Informática de A Coruña  
UNIVERSIDADE DA CORUÑA

MASTER DISSERTATION

MASTER IN HIGH PERFORMANCE COMPUTING

# GPU-Accelerated Pricing of Credit Value Adjustment for Interest Rate Swap under a Gaussian HJM Model

**Student:** Michel Agustín Herrera Sanchez

**Supervisor:** Dr. Diego Andrade Canosa

A Coruña, 28 de julio de 2021.

## **Acknowledgements**

Throughout the writing of this dissertation, I have received a great deal of support and assistance.

I would first like to thank my supervisor, Professor Dr. Diego Andrade Canosa, whose expertise was invaluable in formulating the research questions and methodology. Your insightful feedback pushed me to sharpen my thinking and brought my work to a higher level.

I would also like to thank Dr. Leonidas Kosmidis from Barcelona Super Computing for their valuable guidance throughout the Barcelona Supercomputing Center GPU Hackathon. You provided me with the tools and ideas that I needed to choose the right direction and successfully complete the application benchmark.

Finally, I would like to thank my family for their patience and support.

## Abstract

The primary motivation for the current work is the acceleration of the Credit Value Adjustment (CVA) computation for an Interest Rate Swap portfolio written on forward contracts in a fixed income market. *Monte Carlo* method is used to simulate a set of *stochastic differential* equations which help to predict the prices of the portfolio at different times in the future. Numerical *Monte Carlo simulation* methods are *embarrassingly parallel* problems that turn out to be suitable for parallel acceleration and adapt easily to heterogeneous hardware constraints. These numerical methods converge to the true solution, require less memory and are easier to program than other numerical methods. It's well known that numerical integration with *Monte Carlo method* is a good candidate to execution speedup in *heterogeneous processors* like NVIDIA capable GPU. GPGPU technology has shown great potential in conventional derivatives pricing as well as in some risk-management applications. Three different implementations were devised on CPU, NVIDIA single-GPU and multi-GPU to evaluate the potential speed-up the heterogeneous processors technology offers. The risk factor *Monte Carlo simulation* speedup observed was up to 170x between single GPU vs a single CPU and of up to 600x when comparing four GPUs to a single CPU.

# Contents

1. <i>Introduction</i> .....	5
1.1    Organization.....	5
2. <i>The Credit Valuation Adjustment for an Interest Rate Swap Portfolio</i> .....	6
2.1    Credit Valuation Adjustment .....	6
2.2    Computing the Credit Valuation Adjustment .....	7
2.3    The Counterparty Credit Risk Computational Challenge.....	8
3. <i>Expected Exposure Profile Simulation for CVA</i> .....	9
3.1    Vanilla Interest Rate Swap Product Valuation.....	9
3.2    Risk Factor Simulation .....	10
3.3    Mark to Market .....	16
3.4    Expected Exposure Aggregation.....	16
4. <i>Parallelizing the Expected Exposure Profile simulation</i> .....	18
4.1    Parallelizing the Forward Rate Simulation with CUDA.....	18
4.2    Interest Rate Swap Mark to Market.....	20
4.3    Expected Exposure Aggregation.....	20
5. <i>Optimization</i> .....	21
5.1    Exploiting oversubscription .....	21
5.2    Using Shared Memory .....	22
5.3    Scaling with multiple GPUs .....	23
6. <i>Performance Results</i> .....	25
7. <i>Conclusion</i> .....	30
8. <i>References</i> .....	31
9. <i>Appendix A Gaussian random variates generation</i> .....	32
10. <i>Appendix B Simulation Grid of a Gaussian HJM Model</i> .....	34
11. <i>Appendix C Numerical Simulation Driver</i> .....	35
12. <i>Appendix D Naïve Parallelized Numerical Simulation Algorithm</i> .....	36
13. <i>Appendix E Accelerated Numerical Simulation Algorithm</i> .....	37
14. <i>Appendix F Accelerated Vanilla Interest Rate Swap Mark to Market</i> .....	38
15. <i>Appendix G Simulation kernel naïve implementation</i> .....	39
16. <i>Appendix H Simulation kernel implementation with shared memory</i> .....	41
17. <i>Appendix I Multi-gpu implementation</i> .....	43
18. <i>Appendix J System Configuration</i> .....	44
19. <i>Appendix K Full Code Implementation</i> .....	45

# Chapter 1

## Introduction

Simulation has already become an industry-approved method for estimating financial security prices for which a *simple closed-form solution* does not exist. Many problems in Mathematical Finance entail the computation of a particular integral. These integrals can be valued analytically, using numerical integrations or computed using partial differential equations. However, when the number of dimensions in the problem is large, *partial differential equations* and numerical integrals become intractable. In these cases, *Monte Carlo* simulation methods give better results at a fastest convergence rate.

Numerical *Monte Carlo simulation* methods are *embarrassingly parallel* problems that turn out to be suitable for parallel acceleration and adapt easily to heterogeneous hardware constraints. These numerical methods converge to the true solution, require less memory and are easier to program than other numerical methods. Nevertheless, it is well known that they are very time-consuming and computationally intensive. That is why numerical integration with *Monte Carlo method* is a good candidate to accelerate the implementation in *heterogeneous processors* like NVIDIA capable GPU. GPGPU (General-Purpose-Graphics-Processing-Unit) technology has shown great potential in conventional derivatives pricing as well as in some risk-management applications.

The primary motivation for the current work is the acceleration of the Credit Value Adjustment (CVA) computation for an Interest Rate Swap portfolio written on forward contracts in a fixed income market. *Monte Carlo* method is used to simulate a set of *stochastic differential* equations which help to predict the prices of the portfolio at different times in the future. Three different implementations were devised on CPU, NVIDIA single-GPU and multi-GPU to evaluate the potential speed-up the heterogeneous processors technology offers.

### 1.1 Organization

The study of the current work is based on the acceleration of *Monte Carlo simulation* of stochastic differential equations using multiple GPUs. A problem-driven approach has been used along the study where initially the problem is stated, a sequential and parallel solutions are proposed, optimization opportunities are considered for code acceleration and a performance analysis is given at the end. The *asses, parallelize, optimize and deploy* methodology developed by NVIDIA is used through the whole work [NV02].

After a short introduction to OTC market, the chapter 2 explain a step-by-step explanation on how to calculate the CVA metric for a vanilla interest rate swap portfolio. Then, the concrete problem of computational challenge that represents the calculation of counterparty credit risk for financial derivatives portfolios is introduced. Chapter 3 explain the fixed income derivative interest rate swap, valuation and pricing is detailed step by step. In chapter 3.2 the Gaussian Heath, Jarrow and Morton model for forward rate simulation using Musiela parametrization stochastic differential equation is explained. Chapter 3 provides insides on how to generate the expected exposure profile and the initial implementation that led to identify which part of the code are suitable to accelerate on GPU.

Chapter 4 looks further into a basic GPU parallelization and which optimization methods could be applied to improve the speed up and the application throughput. The performance results and discussed in Chapter 5. Finally, some possible extensions of the work are suggested, and conclusions are given in chapter 6.

# Chapter 2

## The Credit Valuation Adjustment for an Interest Rate Swap Portfolio

Financial derivatives are products that are based on the performance of another, uncertain, underlying asset. An *Over-the-counter* (OTC) market is a decentralized market in which market participants trade derivatives directly between two parties and without a central exchange or broker. The OTC derivatives market is divided by asset class:

- Interest Rate
- Foreign exchange
- Credit derivatives
- Commodity derivatives

Interest Rate derivatives is the biggest market with a notional amount standing value of 466,494 billions of US dollars, according to Bank of International Settlements. [BIS]

Interest is the money paid for credit or similar liability. Examples are bond yields, interest paid for bank loans, and return of savings. Money worth today 1 USD might not have the same value in the future. Time value of money refers to money earned today have the potential to increase its value over a given period. Money deposited into a saving account earns an interest rate. Interest is typically calculated as a percentage of the amount owed to the lender. The percentage that is paid over a certain time period is the interest rate. Interest Rate are market prices which are determined by supply and demand. The *LIBOR* rate is a living example of such rates obtained by observing the market. Forward rates are the interest rates applicable to a financial transaction that will happen in the future.

Interest Rates can be considered fixed or floating. Fixed rates as the word implies are fixed during the duration of the whole contract and can be reset with certain frequency. Floating rates in the opposite depends on how the market behaves in the future, which gives them a certain uncertainty and its trajectory in the future can be modelled by a stochastic model. [CGR][WILL][CHY]. Trading in the interest rate market is done directly with a financial counterparty. Therefore, if any of the two-counterparty involved in the contract default the value of the contract can be considered a loss.

### 2.1 Credit Valuation Adjustment

Counterparty credit risk is the risk that the counterparty to a financial contract will default prior to the expiration of the contract and will not make all the payments required by the contract. Only the contracts privately negotiated between counterparties over-the-counter derivatives and security financing transactions are subject to counterparty risk. [PYT]. Counterparty risk arises from a broad class of financial products. The most traded are interest rate derivatives hence during this case of study we will be concentrating on a portfolio containing one vanilla interest rate swap derivative contract. The risk-free price of a contract must be different from a risky (counterparty might default) contract. The price of a risky derivative could be thought as the risk-free price minus the component correcting for the counterparty risk.

$$P_{\text{risky}} = P_{\text{riskfree}} - CVA$$

The latter component is called Credit Valuation Adjustment (CVA). This counterparty risk charge should be calculated in a sophisticated way to account for all aspect that define the CVA including

- The underlying contract(s) exposure
- The default probability of the counterparty
- Netting of existing transactions with the same counterparty
- Collateralization

For the sake of simplicity, the CVA calculation in the current work will be uncollateralized and no risk mitigation like netting set will be included in the scope of the current work. Since the contract value changes unpredictable over time, as the market moves only the current exposure is known with certainty, while the future exposure is uncertain [PYK]. Due to the need to investigate the future for those possible scenarios in which the counterparty can default or there is a move in the market against the positions traded in the contract it would be required to make use of *Monte Carlo* Simulations to address the issue in more than one dimension.

## 2.2 Computing the Credit Valuation Adjustment

Computing CVA on a portfolio of interest rate swap can be separated on several steps:

- Calibration of underlying models to simple market instruments
- Generation of forward rate, discount factors and portfolio scenarios
- Computation of the exposure profile
- Discounting exposures to current time and weighting them with the corresponding default probabilities.

Given the Exposure modelled framework explained before where exposure is simulated at a fixed set of simulation dates  $\{t_k\}$   $k = 1, N$  the computation of the CVA can be approximated as:

$$CVA = LGD (1 - R) \sum_{k=0}^n DF(t_k, t) EE(t_k) PD(t_k, t_k - 1) \quad (1)$$

Where

- LGD is the loss fraction (of total value), given default
- R is the recovery rate
- DF ( $t, \tau$ ) Is the discount factor, discounting the value from time of default back to the valuation time.
- EE ( $t$ ) is the Expected Exposure at time  $t$
- PD ( $t_k, t_{k-1}$ ) is the probability of default of the counterparty at time  $t$ .

The *expected exposure curve* is calculated as an average of all *exposure curves* generated and filtered out through a Mark to Market process. The *probability of default curve* ( $PD[t_i]$ ) is the risk neutral probability of counterparty default between  $s$  and  $t$ . These probabilities were obtained from the term structure of credit default swap (CDS) spreads Sovereign Country. The JPMorgan method is generally used for the *probability of survival term structure* bootstrapping.

## 2.3 The Counterparty Credit Risk Computational Challenge

Calculating XVA is one of the largest challenges a bank faces daily due to the required volume of calculation to produce a figure. The generation of the *expected exposure profile* curve (*EE*) is the most complex and computation intensive step in the computation of the *CVA*. Which is the process where all *Monte Carlo simulations* happen, and we must aggregate and average the results to obtain the *expected exposure profile curve*. The number of *exposure profiles* generated depends on the number of scenarios and total number of assets.

To illustrate the volume and complexity of the problem consider an investment bank generally needs to process for the bigger counterparties and for all scenarios:

- 1,000,000 of trades
- 5000 pricing points
- 400 scenarios

This results in 200TB of simulations to be completed in a 24h or less. Derivative pricing are generally small problems compared with their scientific applications counterparts but requires a high volume of processing to be completed fasted based on the size in the number of assets held by an investment bank. [RUI]. This is when the application of High-Performance Computing, in specifically the data parallel programming implemented by NVIDIA CUDA and the heterogeneous computing hardware accelerates the calculation of the XVA figures in a fast way and with high degree of accuracy. The parallel processing allows to speed up the calculations by allowing the processing of multiple scenarios at the same time by taking advantage of throughput-oriented design of modern GPUs.

# Chapter 3

## Expected Exposure Profile Simulation for CVA

The Exposure Profile calculation is the most compute and memory intensive part of the CVA calculation. It has three phases applied for every trade present in the portfolio:

- Scenario Generation: Risk Factor Model Simulation
- Instrument Valuation: Mark to Market
- Aggregation: Expected Exposure Profile generation

This chapter demonstrates how the expected exposure curve calculation is implemented for an Interest Rate Swap portfolio. For each step the sequential implementation details are given with the intention to introduce later the parallelization techniques.

### 3.1 Vanilla Interest Rate Swap Product Valuation

An interest rate swap contract is an agreement between two counterparties to exchange fixed interest rate payments for floating interest rate payments, based on a predetermined notional principal, at the start of each of several successive periods. [BUK] There are two legs associated with each party. The fixed leg is associated with fixed interest rate payments and the floating leg is associated with future changes in a rate generally indexed by *LIBOR*.

Interest Rate Swap are used to manage exposure changes in interest rates. To illustrate this let's consider the following example:

Suppose a bank lends money to a customer A with a fixed interest of 8% to pay back in 30 years. As the interest is fixed a change in the market, if the interest rate grows, could potentially produce loss to the bank. In order to hedge this problem, the bank engages in a contract with another counterparty such that the bank receives 8% fixed interest rate payments and pays a floating interest rate based on a notional. That way every time the interest rate goes below 8% the bank protects itself from fluctuations in the interest rate market.

A vanilla interest rate swap is a contract agreed at time  $t$  and expires at time  $T$ , whereby a party exchanges a fixed against a floating interest rate (the LIBOR in the sequel) payment over  $n$  reset dates  $T_0, \dots, T_{n-1}$  and payment periods  $T_1 - T_0, \dots, T_n - T_{n-1}$ , with  $T \leq T_0$  on a certain notional amount. The length,  $T_n - T_0$  is known as the tenor of the interest rate swap. (Figure 1)

Let  $L(T_{i-1}, T_i)$  denote the LIBOR rate applying to the period from  $T_{i-1}$  to  $T_i$ . The payoff at  $T_i$  to the counterparty paying the fixed rate  $\delta_{i-1}$  ( $L(T_{i-1}, T_i) - K$ ) (\*), where  $K$  is the fixed interest rate and  $\delta_{i-1}$  is the length of the reset interval that follows a 30/360 convention. Therefore, the value at  $t$  of a forward starting interest rate swap is

$$PayOff(K; DF; L; T_1, \dots, T_n) = N \sum_{i=1}^n DF_i \delta_{i-1} (L(T_{i-1}, T_i) - K) \quad (2)$$

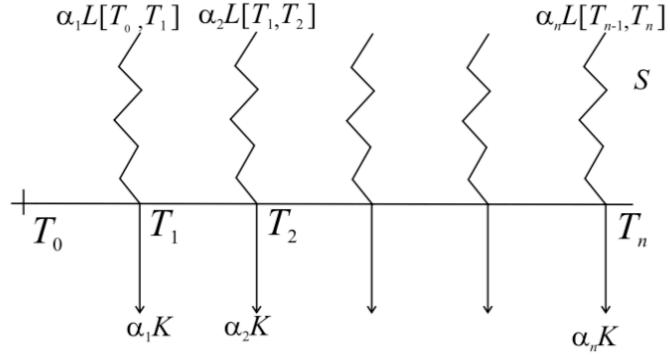


Figure 1. Interest Rate Swap Cashflows exchanges

The *Interest Rate Swap* payoff is the summation of the differences between the value of the floating leg minus the value of the fixed leg at all timepoints from 0 till expiry. Where  $N$  represents the notional principal amount,  $r_\tau$  denotes the instantaneous interest rate as of time  $\tau$ , and  $DF_i$  is the discount factor taken under the risk-neutral probability  $Q$ , and conditional on the information at time  $t$ . The future value of the  $L(T_{i-1}, T_i)$  is unknown but can be modelled using an stochastic model. The value of  $L(T_{i-1}, T_i)$  can be calculated from simulated forward rates.

### 3.2 Risk Factor Simulation

The underlying *forward rate* asset behaviour of the Interest Rate Swap derivative is modelled using the Heath, Jarrow and Morton Model. The stochastic models of the underlying asset are governed by a collection of *stochastic differential equations* which need to fit the model's parameters to provide rates according to the real market situation. As a result of the model parameters calibration and the simulation of the *stochastic differential equation* the term structure of bonds and forward rate curve is obtained to price the vanilla interest rate swap derivative. [GAT]

Since the stochastic differential equations modelled in the 3-factor gaussian HJM model does not have a known closed-form solution the simulation is needed to find the risk factors which are needed to find the future realization of the interest rate swap portfolio exposure. [BMM]

The value of the stochastic differential equation at each simulation step and time point represents the future realization of the forward rate. At each simulation step a new *stochastic differential equation* is evaluated on each tenor point. The forward rate is required to calculate the *LIBOR* rate and the discount factor on 51 points that represent all the cash flows exchanged across time. The number of simulations is calculated by dividing expiry( $T$ ) over the simulation time step  $dt$  giving it 2500 in total. To increase the number of simulations is required to decrease the value of  $dt$ .

At each timepoint the forward rate is calculated by solving a stochastic differential equation using Euler Maruyana method (JAECK)(HDE)

$$f(t + dt) = f(t) + m(t) * d(t) + \sum_{n=1}^3 (\text{vol}_{n[t]} * \text{phi}[n] * \sqrt{dt}) + \frac{dF}{dtau} * dt \quad (3)$$

Where:

- $f_j(t+dt)$  is the forward rate change over time  $dt$
- $t$  takes values between 0 and  $T$  expiry
- $dt$  is the size of the simulation step (0.01)
- $dtau$  is the difference between two consecutive tenors
- $m(t)$  the calibrated drift
- $\text{vol}_i(t)$  the calibrated volatility  $i = 0..2$
- $\text{phi}[n]$  is a gaussian variate with distribution (0, 1)
- $dF$  represents the difference between two consecutive rates. If the forward rate is not part of the right boundary (at expiry) then  $dF$  is calculated as  $f(t+dt) - f(t)$ , otherwise  $f(t) - f(t-dt)$ .

Although details about the model calibration will not be provided. It's worth to mention that to obtain the volatilities associated with the 3-factor risks model. A covariance matrix is created from a collection of Yield Curve from Bank of England expanded across 20 years. The covariance matrix is then decomposed using dimensional reduction method commonly known as *Principal Component Analysis* (PCA). As a result, the trend is resumed in three factors: shift, twist and butterfly [NOV] that define the shape and dynamics of all simulated forward rate curves. The result of the whole process are three volatility curves expanded from 0 to an expiry day that are used to obtain the drifts and the diffusion term on each *stochastic differential equation*. [GAT]

As mentioned before there is a realization of a different stochastic differential equation on each tenor and simulation point given the usage of different calibrated drift and three volatility factors. Each of the stochastic differential equations generated at each pricing and simulation point is used to produce the discount factors and forward rates to price the *Vanilla Interest Rate Swap* futures cash flows.

A whole round of simulation (forward rates, discount factors) is performed every time the interest rate swap product is priced. The smaller the size of simulation timesteps the higher the number of simulations to be performed. Figure 2 details the required steps to simulate the forward rates  $f(t,T)$  and calculate the discount factor  $Z(t, T)$ . A high number of simulations guarantee that the price for the derivative converges or a required minimum in its variance is achieved, holding the *Law of Large Numbers* for a convergence criteria. The number of simulations is calculated by dividing expiry( $T$ ) over the simulation time step  $dt$  giving it 2500 in total. To increase the number of simulations is required to decrease the value of  $dt$ .

**Simulation** Simulate an evolution of the whole risk-neutral curve for the necessary length of time, from today  $t^*$  to  $T^*$  (0 .. 51)

- For a risk neutral forward curve use GLC data (UK Gilts)
- Realizations of the entire curve  $f(t, T)$  over time steps  $dt = 0.01$  are in rows.
- Generate normally distributed  $N(0,1)$  pseudo random numbers  $\phi_0, \phi_1, \phi_2$
- $f(t, T)$  uses a simple forward Euler Maruyama scheme to compute:  $f(t) = f(t+dt) + dfbar$

where  $dfbar = m(t)*dt + \text{SUM}(\text{Vol}[i]*\phi[i] * \text{SQRT}(dt)) + dF/dtau*dt$  and having  $i = 0..2$

- The paths of forward rates for discretised tenors tau are in columns.

**Discount Factors** Obtain ZCB values for all required tenors up to  $T^*$ . Pricing a Zero Coupon Bond with any stochastic short rate  $r(t)$  under HJM.

$$Z(t, T) = \mathbb{E}^{\mathbb{Q}} \left[ \exp \left( - \int_t^T r(s) ds \right) \right]$$

Looking at the HJM output, is known that there is only an evolution of the first point on the yield curve  $f(t,t) = r(t)$ .

$$Z(t, T) = \exp \left( - \sum r_t \Delta t \right) \quad \text{under MC}$$

Monte Carlo Simulation is always discretized over time step  $dt$ . Therefore, integration becomes summation.

**Cashflows.** Calculate the value of cashflows of interest and apply discounting. Consider  $(FRA - K)^+$  payoff and discount.

The forward rate simulated at each simulation step is used as an input for the next simulation. There is a dependency across simulation timesteps and across tenors

Figure 2. Heath-Jarrow-Morton Framework with Musiela Parametrization Simulation

A simulation grid containing all forward rates (See appendix B for an example) is obtained as the result of applying the steps described in Figure 2. The top boundary points of the simulation grid are initialized with the spot rates from Bank of England and the inner points update the futures forward rates iteratively. As stated, before a different *stochastic differential equation* is evaluated at each point. Each inner point updates its *forward rate* by using the previous calculated forward rate of its neighbouring points and itself. The rate updating operation for all points in the grid needs to last long enough which means many iterations are required to get the final converged values. Figure 3 shows the calibrated input parameters volatilities and drift that are specific on each pricing point. The first row of the grid represented in figure 4 is the spot rate curve and subsequent rows represent the simulated forward rates. A different simulation grid is generated for each of the 50,000 scenarios (Figure 4). For each product pricing is required to create a grid of 51 points in the horizontal direction, representing pricing points of an interest rate swap of 25 years with cash flows exchanges every 6 months and 2500 simulations points in the vertical direction. The number of simulated rates per scenario is close to 127500 points.

Tenor $\Rightarrow$	0.00	0.5	1.0	1.5	2.0	2.5	3.0	3.5	4.0	4.5	5.0	5.5	6.0	6.5	7.0	7.5	8.0	8.5	9.0	9.5	10.0	10.5	11.0	11.5	12.0	12.5	13.0	13.5	14.0	14.5
Drift	0.000000	0.000096	0.000069	0.000099	0.000128	0.000155	0.000182	0.000208	0.000233	0.000257	0.000280	0.000303	0.000324	0.000345	0.000364	0.000381	0.000397	0.000412	0.000426	0.000438	0.000449	0.000459	0.000469	0.000478	0.000487	0.000496	0.000505	0.000515	0.000526	0.000538
Vol_1	0.006431	0.006431	0.006431	0.006431	0.006431	0.006431	0.006431	0.006431	0.006431	0.006431	0.006431	0.006431	0.006431	0.006431	0.006431	0.006431	0.006431	0.006431	0.006431	0.006431	0.006431	0.006431	0.006431	0.006431	0.006431	0.006431	0.006431	0.006431	0.006431	
Vol_2	0.003557	0.003812	0.004010	0.004155	0.004249	0.004295	0.004295	0.004252	0.004169	0.004048	0.003893	0.003705	0.003488	0.003245	0.002977	0.002587	0.002380	0.002056	0.001719	0.001371	0.001025	0.000655	0.000291	0.000072	0.000432	0.000787	0.001133	0.001468	0.001790	0.002095
Vol_3	0.004751	0.003909	0.003135	0.002427	0.001783	0.001201	0.000677	0.000211	0.000020	0.000052	0.000082	0.000073	0.001137	0.001355	0.001531	0.001764	0.001855	0.001853	0.001823	0.001766	0.001685	0.001583	0.001462	0.001324	0.001171	0.001007	0.000832	0.000651	0.000464	

Figure 3 Calibrated Drift, Volatilities

Figures 4 offer a glimpse on how the entirely forward rate evolves across timepoints and simulations. It's not difficult to appreciate the random behaviour of the forward rate and the reached maximum value of 7% and a minimum value of 2%. In Figure 5 it's possible to appreciate the shape of the simulated curves that captures the different directions or trends (*shift, twist and butterfly*) obtained by the PCA method during the calibration process for the 3-factor volatility curves.

Although convergence analysis is a mandatory step to validate the Monte Carlo simulation it's omitted on this implementation and instead a considerable high number of simulations are run assuming this can potentially reduce the standard deviation between results as the number of simulations increase.

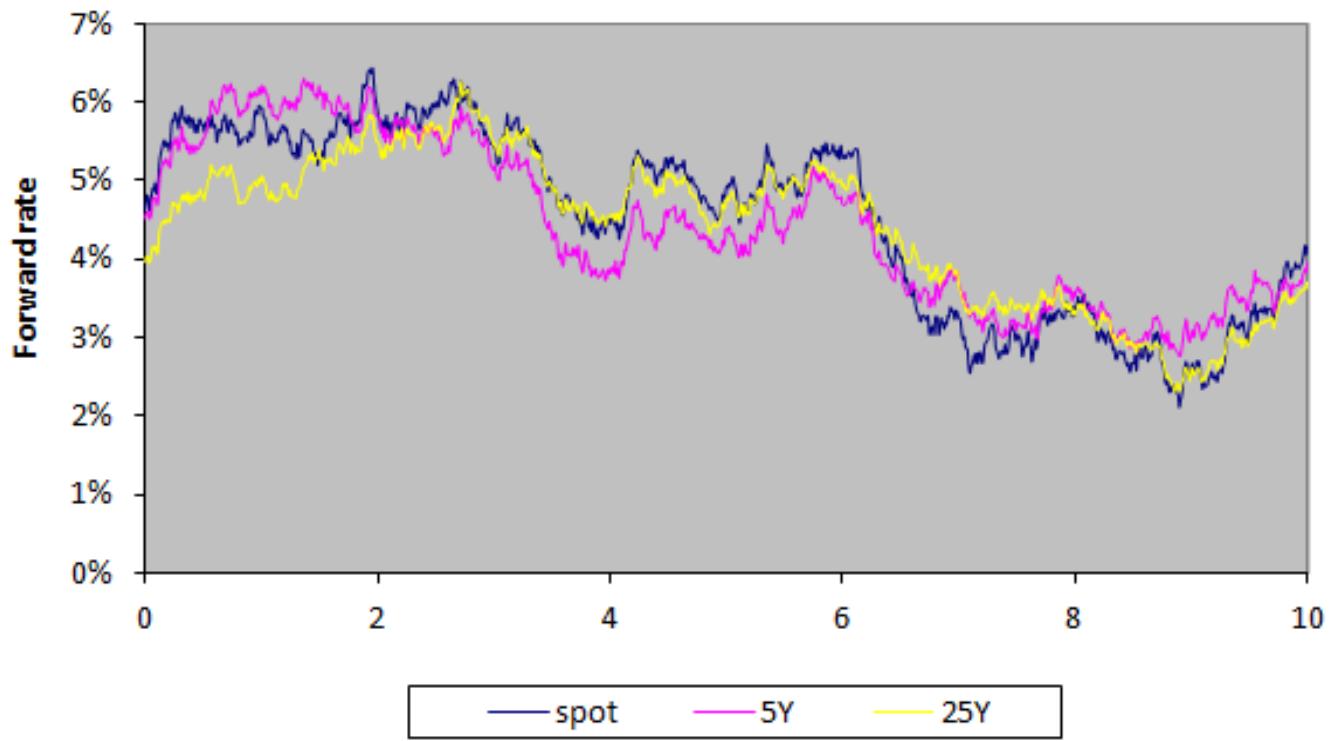


Figure 4. Forward Rates Evolution for each Tenor (Series by Row)

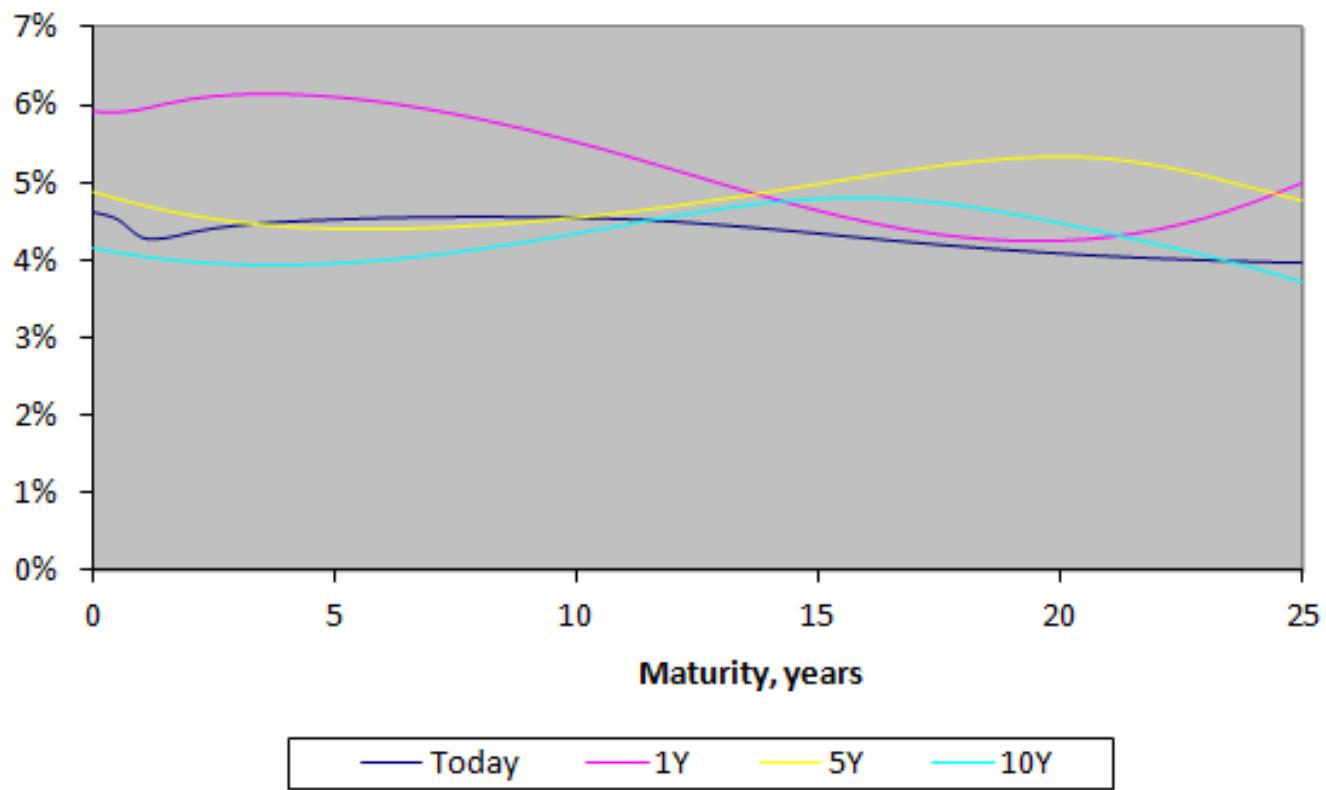


Figure 5. Forward Rates Structures – Yield Curve

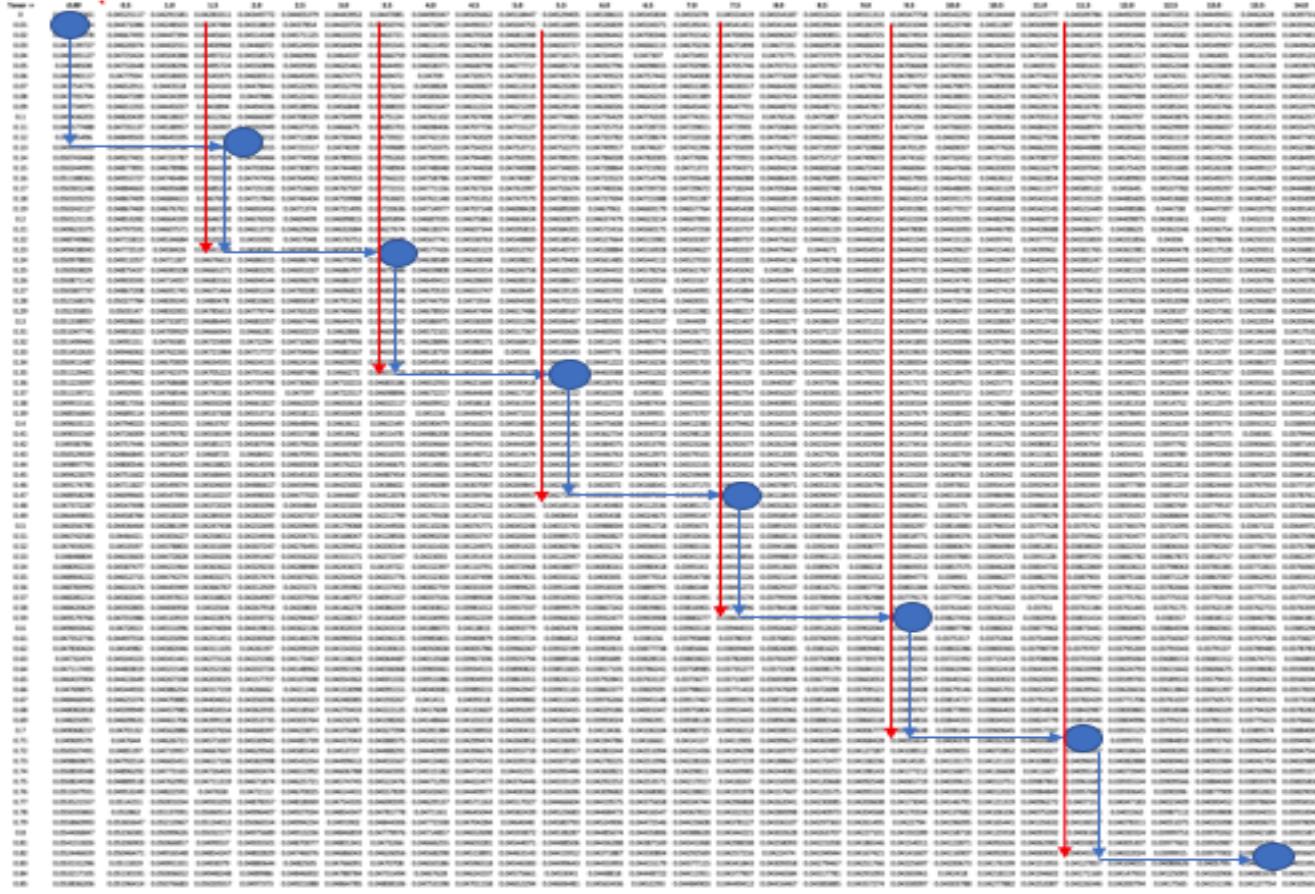


Figure 6. Forward Rate and Discount Factor 2D array Simulation Grid.

If we consider the simulation grid as a 2D array `sim_grid`. Simulations can be found in the Y dimension as increments of  $dt$  and a different *stochastic differential equation* is simulated at each pricing point across the X dimension. The Forward Rate A *forward rate* is generally represented as  $F(t; t_1, t_2)$ . The value  $r(t)$  given to calculate the  $F(t; t_1, t_1)$  can be found in the simulation grid as `sim_grid[t/dt, t_2]`, which is represented as blue dots in Figure 6. For example, the rate  $r(t)$  to calculate forward rate  $F(1Y; 0, 6M)$  can be found located at `sim_grid[100, 1]`.

Bonds are generally represented by  $B(t, T)$ , are calculated first as an exponential of the summation of all simulated rates across the whole column that represent the tenor  $t$  up to row  $\text{sim\_grid}[t/dt]$  times  $dt$ . Figure 6 is represented by red arrows where summation occurs from top to bottom direction. All previous forward rates simulations are required to be considered in order to be able to calculate subsequent *discount factors*. See Figure 2 for details on how the calculation is done.

The LIBOR rate is the ratio between two continuous bonds. Its calculated nearly equivalent to the Forward Rate calculation. LIBOR rate is given by the formula:

$$L(t, T, S) = \frac{1}{tau} * \left( \frac{B(t, T)}{B(t, S)} - 1 \right) = \frac{1}{tau} * \left( \frac{e^{-\sum_{i=0}^T r(i)dtau}}{e^{-\sum_{i=0}^S r(i)dtau}} - 1 \right) \quad (4)$$

$$L(t; T, S) = \frac{1}{\tau} * (e^{-r(t)*\tau} - 1) \text{ where } \tau = T - S \quad (5)$$

The HJM framework algorithm simulate each element  $r(t)$  of an  $N \times M$  grid of all rates. There are strong dependencies between rates calculated at futures simulations times in the vertical direction although in the horizontal direction rates can be computed independently. There is  $O(N^3)$  of computational complexity as shown in Appendix B. The main simulation driver needs to execute three nested loops to produce a cube of simulations rates. From the external path loop, passing for the intermediate scenario generation loop and ending the most inner loop that cover each of the points of the whole forward curve. Figure 7 shows the cube generated across all scenarios 51 points for each of the 2500 simulation steps. Each of the vertical slices of the cube (scenarios) can be calculated as separated sequential chunks which allows it to accommodate a high number of forward rate values per scenario. The implementation is very similar to a 2D heat conduction application that is a stencil-based code.

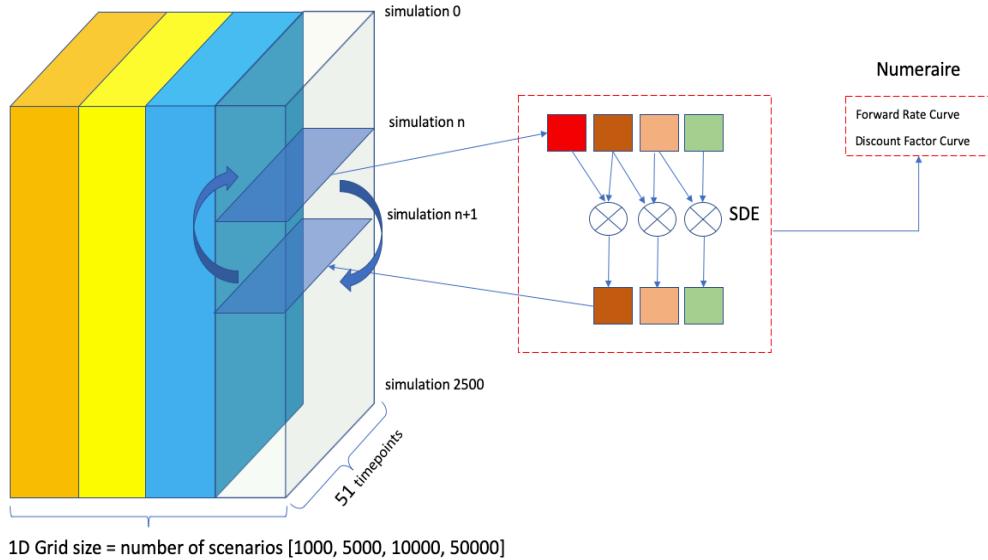


Figure 7. Simulated Risk Factor Cube / Simulation Algorithm implementation.

Figure 7 shows the implementation to generate a 2-D forward rate simulation grid with `pathN(2500)` and 51 points in the horizontal direction. Initially the forward rate curve is equal to the spot rate across all points of `variable_simulated_rates`. The kernel function `_generatedPaths_kernel_cpu` calls the *stochastic differential equation* to simulate the stochastics rates at each simulation point which are stored in `simulated_rates0` float array. After each iteration step (`dt`), `simulated_rates0` stores the new forward rates and then swap the pointer so that the next iteration the input kernel points to the current simulated new forward rates. When all scenarios are completed, the pointers need to be swapped. On each simulation step the kernel function `_generatePaths_kernel_cpu` is invoked and the forward rates and their cumulative summation are stored in the arrays `accum_rates` and `d_numeraire`.

As discussed above at each simulation step of size the function `_generatePaths_kernel_cpu` numerically simulates 51 stochastics differential equations. At each simulation point the value of each *stochastic differential equation* is calculated independently but it always depends on values of the previous simulation step. Which somehow is like the concepts of stencil on each point. The neighbour's simulated forward rate values along the volatility, drift and gaussian variates are passed to the function to numerically simulate the *stochastic differential equation* at each of the 51 simulation points. Internally one of terms of the *stochastic differential equation* depends on the neighbour's values. If the rate can be represented by `rate[i][j]` then his value depends on the neighbours values `rate[i][j-1]` and `rate[i+1][j-1]` if it's an internal simulation point, at the border the neighbours values are `rate[i-1][j-1]` and `rate[i][j-1]`.

### 3.3 Mark to Market

Once the future market scenarios are simulated for a fixed set of simulation dates using evolution models of the risk factors (PYK), the next step is to proceed to the *vanilla interest rate swap* (IRS) instrument valuation for each simulation date and for each scenario, hence simulated forward rates and discount factors.

The dynamic of the *exposure curve* is driven by the generation of multiple scenarios in the future by simulating the *discount factors* and *forward rates* stochastic models. Those numeraires are used to price the vanilla interest rate swap derivative contract on each simulation point  $t$  from zero to expiry ( $T$ ). A mark to Market process is used in place to filter out our results and focus only on positive values across pricing points because we are interested in potential loss.

A central concept in *Mark to Market* portfolio calculation is the net present value (NPV) of the vanilla interest rate swap derivative at time  $t$ . Basically the interest rate swap is valued as explained in 3.1 for every pricing point between 0 and expiry. The *cashflows* are valued and accumulated from  $[t+1, T]$  for each pricing point  $t$  for each scenario. As a result, the exposure curve is build comprising the points from  $[0, T]$ . The Exposure is given by  $NPV(t)^+ = \max \{ 0, NPV(t) \}$  as the positive *net present value* at time  $t$  of the underlying Vanilla Interest Rate Swap product contract, which is commonly known as mark to market. The result is a curve representing the portfolio exposure on each timepoint which is known as *exposure profile*.

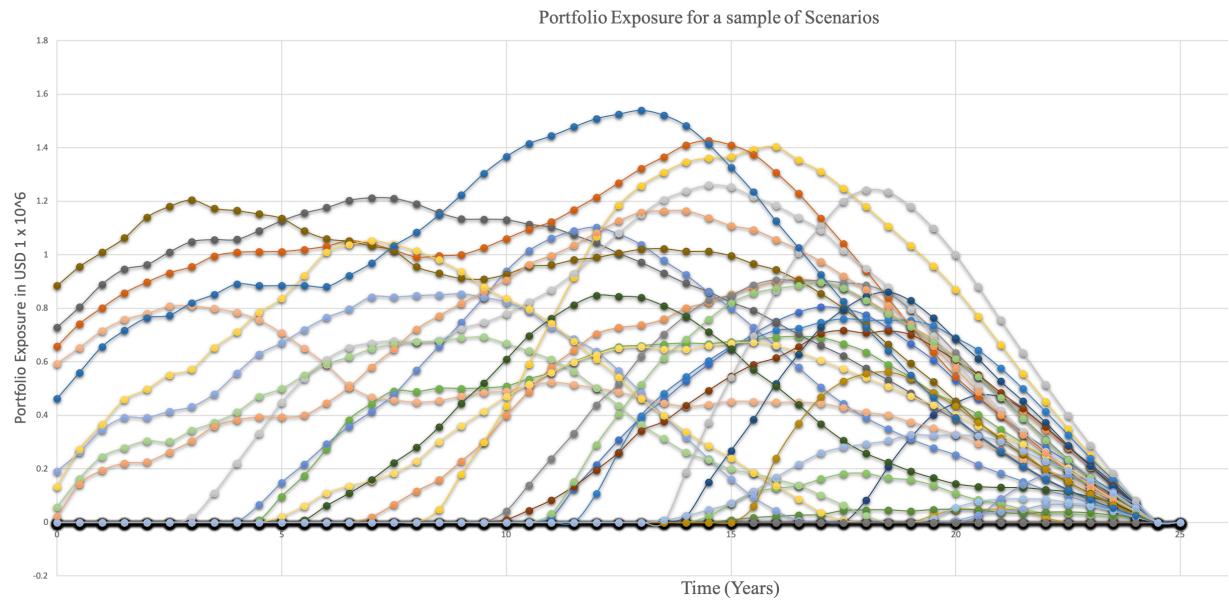


Figure 8. Portfolio Exposure Sample Scenarios Trajectories

### 3.4 Expected Exposure Aggregation

The *expected exposure profile* (EE) for a *vanilla interest rate swap* is a curve that expands in time from the initial day of the contract till the expiry day. The *expected exposure profile* is basically an average on each simulation point of all the realizations of the exposure profiles in the future. A small representation of the simulated exposure curves to be aggregated can be observed in Figure 8.

The problem to be tackled here is the reduction of matrix columns across all rows of a dense matrix  $M \times N$ . Rows in the matrix are each of the different simulated exposures profiles ( $M$ ) and  $N$  is a vector with values 1 with a dimension equals to the number of simulation points where exposure profile curve is evaluated.

It would be worth it to point out how the approach of reducing the columns of a matrix through the multiplication of a row by the same matrix can be generalized to perform the linear combination of an ensemble of vectors. The *expected exposure profile* aggregation can be implemented not by performing parallel pair-summations across all columns but instead using a linear algebra transformation that multiplies the matrix (one exposure by row) with a ones vector ( $d_x$ ) of dimension equals to the number of scenarios using BLAS `cblas_gemv` function. This is a very efficient way as it takes advantage of the software acceleration done by Intel MKL.

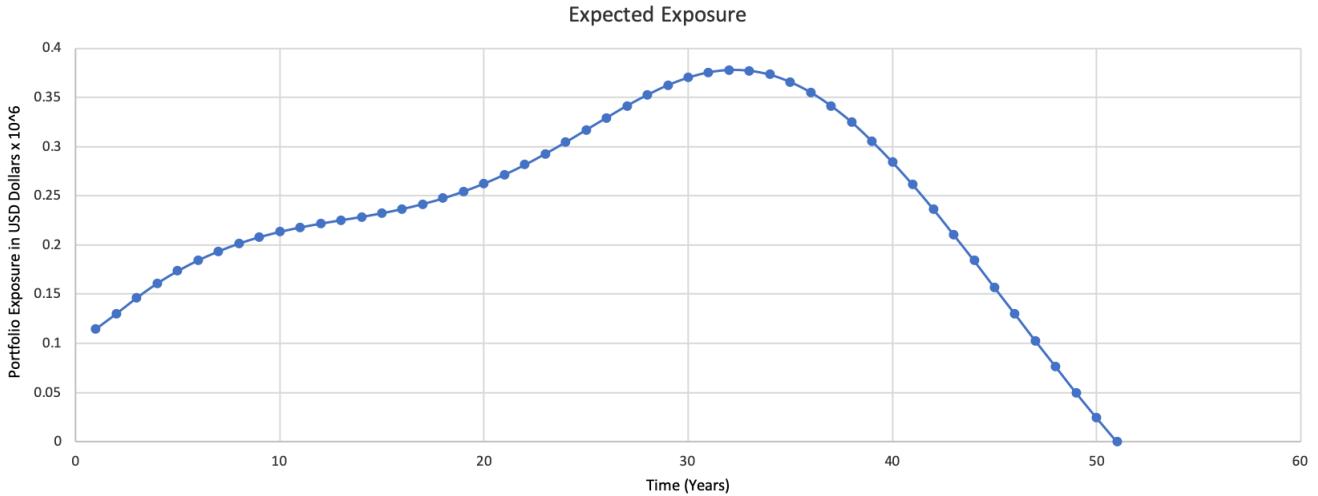


Figure 10. Expected Exposure Profile in blue colour

Observe in Figure 10 the hump shape that takes the exposure profiles for an Interest Rate Swap portfolio. The value of the exposure in the exposure profile curve increases until reaching the peak value then the value of the exposure decreases as the result of evaluating less and less cash flows so his value decreases as the product expires time is reached. The maximum loss represented by a value of 400,000 dollars from a portfolio value of 10,000,000 dollars is reached in simulation point 30, and later it decreases to 0 as the axis time approaches the expiry duration of the contract at point 51. This curve is known as the expected exposure profile.

# Chapter 4

## Parallelizing the Expected Exposure Profile simulation

This chapter explains how to leverage the CUDA programming model to parallelize the expected exposure profile simulation on GPU hardware. In section 4.1 the parallelization of the simulation process is explained. The mark to Market process is addressed in section 4.2 and the parallelization of the aggregation process is developed in section 4.3. The initial solution is not a fully optimized version but provided hints to take into consideration for high throughput implementation of the GPU hardware.

### 4.1 Parallelizing the Forward Rate Simulation with CUDA

The implementation of the Forward Rate Simulation is a stencil-based code which has been extensively implemented and accelerated on GPUs. The main flow of the application generates 50,000 scenarios, each of which requires 2500 simulation steps to build the whole *forward rate curve* of 51 points.

The external loop is used to iterate across all simulation timesteps (Appendix D line 738). All threads in the same thread group simulate the rates for a given *forward rate curve*. Basically, the same kernel `_generatePaths_kernel` function is called across all thread groups. Due to the capacity of GPU device to execute thousands of threads at the same time, there is a direct mapping between the size of the execution grid and the number of futures exposure scenarios to simulate. That is why the grid size has been set to the number of exposure scenarios. Although each SM of the GPU has a maximum number of active blocks and a maximum of 1024 threads per thread group, oversubscription has been put in place to take advantage of the latency-hiding technique embedded in the GPU hardware processing. To this end, all *forward rate curves* are generated in parallel by using a single loop parallelization. This way, if a block of threads cannot progress due to global memory access latency another group of threads could be scheduled for execution. This strategy increases the throughput as all the SMs of the GPU are occupied most of the time.

The main intention is to exploit the high-throughput capabilities of the GPU device. The current hardware setup can scale the execution of the kernel function `_generatePaths_kernel` across 40 SM. In a few words, at each simulation step represented by the variable path (Appendix D Line 740), a total number of 50,000 threads blocks are fan out or scheduled for execution across all SM. Giving it a ratio of 1250 threads blocks per SM. This oversubscription of threads allows to scale the application. This capability comes naturally for GPU devices allowing to apply loop parallelization techniques where all the independent iterations of the loop are concurrently executed across different processors. In this case, due to the *exposure profiles* being independent from each other, it was feasible to generate them in parallel across all the SM present in the GPU.

The kernel execution configuration is passed through the operator `<<<>>>` where the first parameter (the grid size) is set to the number of scenarios and the second parameter (the block size) is set to the number of tenors in the forward rate curve (64).

To reduce the data transfer time, data is kept on the GPU whenever is possible (Appendix D Line 756). A double buffer is used to avoid exchange data with the host between simulation steps. On simulation step (i-1) the buffer simulated\_`rates0` is the input and simulated\_`rates` are the output. At the end of simulation step the role of both buffers is exchanged. As it can appreciated, because the device pointer simulated\_`rates0` is only used for swapping with the pointer device simulated\_`rates`, the data exchange only happens in the memory device side, preventing data transfers between CPU and GPUs.

It could be trivial to map each of the pricing/simulation points into a 2-dimension structure of size (number of exposures profiles) x (simulation points). This design could limit the number of concurrent blocks executed per device because of GPU memory constraints. Also, it would increase memory latency as memory should be accessed every time a previous simulation result needs to be read from the previous timestep. The *CUDA* programming model can exploit *SIMD* parallelism. Thus, a group of threads can execute the same instruction at the same time if grouped on groups of 32. The internal loop used to evolve the Forward Rate across 51 points is mapped to a group of 64 threads. Each thread of a group can simulate a rate with a device function that implements the stochastic differential equation simulation. The drawback of this approach is that in a 64-thread group with only 51 threads active, 13 threads are always wasted.

As discussed before, the simulation of the Musiela Parametrization *stochastic differential equation* could be seen as a 2-dimentional grid where pricing points were represented in the horizontal direction and timestep simulations in the vertical direction. The number of simulations in the axis depends on the size of each simulation timestep. The kernel `__generatePaths_kernel` implements the numerical simulation of each of the different *stochastic differential equation* defined at each tenor. See *CUDA* implementation in Appendix G.

At each simulation step each thread looks for the data at his memory position and his neighbour of the previous simulation step to evaluate the *stochastic differential equation* value at a given tenor. The stochastic differential equation `_musiela_sde2` could be a weighted sum of the volatilities times a random factor. Every thread from the 64-thread 1-dimensional group invokes the kernel `__generatePaths_kernel` at every simulation timestep. Internally the thread scheduler splits the thread group into another 32-group size called *warp*. Inside the *warp* threads execute in a lockstep fashion, executing one instruction at a time on all threads simultaneously. In this case the execution scheduler splits the thread group into two warps and schedules them for execution in the SM. For efficiency the thread group is recommended to be multiple of warp size hence multiple of 32. Since only 51 threads out of 64 threads are used to calculate the whole forward curve, a total number of 13 threads are discarded for execution based on boundary conditions. The same kernel instructions code is executed across all threads inside the warp but on different memory locations.

Each thread has a one-to-one mapping to each tenor inside the Forward Rate Curve. The memory location where the instruction will be applied is calculated (Appendix G Line 388) and stored in variable `gindex` as a function of `blockIdx.x`, total number of tenors *forward rate curve* and `threadIdx.x`. Each thread checked based on his index information it is still capable of accessing a valid memory location (Appendix G Line 391). A memory location is considered valid if the local thread index is less than the number of tenors of the *forward rate curve* and the global thread index is less than the linearized kernel execution configuration `gridsize`. The first simulated curve takes the values from the initial spot rate curve (Appendix G Line 393). For this case, data is pulled from array `d_spot_rates` which is stored in constant memory. The *CUDA* runtime generally stores the read only constant memory in cache which substantially reduces the number of global memory accesses by all the threads in a thread group.

Otherwise, if the simulation step is intermediate or final for the simulation block ( $\text{path} > 0$ ), then all the parameters required to evaluate the stochastic differential equation at a given pricing point are pre-computed. Boundary conditions are evaluated for all tenors (Appendix G Lines 400 - 405). As all threads inside the warp need to look at *forward rates* simulated on the previous path a phenomenon known as spatial proximity is observed between threads belonging to the same warp. In that sense, instead of having each thread accessing one global memory location at a time, only one memory transaction is needed to fulfil all the memory requests for the entire warp. Access to global memory is considered a long-latency operation that could cause a whole thread group to be swapped out from the active warp execution group in order to facilitate the execution of other thread group for latency hiding.

A mapping is implemented based on the `blockId`, simulation path and number of volatilities to retrieve the pseudo random values for all threads in the same thread group running the same simulation step. Then the simulation of the *stochastic differential equation* is invoked across the whole thread group but is applied over different parameters located at different memory locations, see (Appendix G Lines 398 – 442). The simulated rate is stored in a register as a result of calling the device function `_musiela_sde2`. The usage of registers

allows the threads to use a private and faster memory location to store intermediate results. Generally, an SM has 65KB of memory capacity for registers. Once the rate is calculated it is first accumulated in the global memory variable `accum_rates` and captured in the vector `simulated_rates0` for next simulations step executions, see (Appendix G Lines 456 – 465). To finish the execution of the simulation (Appendix G Line 382), the step condition is evaluated to verify any interest of storing the accumulated rates and the simulated rate for a given tenor as described before.

## 4.2 Interest Rate Swap Mark to Market

As described in section 3.3, for each simulated scenario only positive values are kept after applying the Mark to Market process. Once all the positive exposure profiles have been generated, an average reduction across all timepoints is required to generate the expected exposure profile, which is the vector input parameters for the CVA calculation.

In general, the operations performed in the mark to market process take a binary operator  $\oplus$  and an array of elements  $[a_0, a_1, a_2, \dots, a_n]$ , and returns the array  $[(a_1 \oplus a_2 \oplus \dots \oplus a_{n-1}), (a_1 \oplus a_2 \oplus \dots \oplus a_{n-2}), \dots, a_{n-1}]$ . This transformation is a reversed prefix-sum operation which can be parallelized as explained in [MAH2]. In the case of the mark to market operation each cashflow is computed with the simulated values (forward rates ( $L$ ), discount factors( $DF$ )) generated after the scenario simulation phase is finished. As shown in formula [6] each cash flow value of the vanilla interest rate swap is computed at each timepoint:

$$\text{cashflow}(K; DF; L; T_1, \dots, T_n) = \sum_{i=1}^n D F_i \delta_{i-1} (L(T_{i-1}, T_i) - K) \quad (6)$$

$$\text{price}(i) = \sum_{k=i+1}^n \text{cashflow}[i] \quad (7)$$

The whole prefix-sum can be parallelized [MHA3]. First, a parallel map function represented at formula [6] is applied over each component of the input vector array and produces a temporal intermediate array, then a `prefix_sums` method is invoked to calculate the exposure values at each timepoint of the exposure curve, as seen in formula [7]. See the Mark to Market CUDA kernel implementation in Appendix F.

## 4.3 Expected Exposure Aggregation

The expected exposure aggregation follows the same idea as the one explained before in the cpu implementation context. The main difference relies on that the unitary vector and the matrix where each row is represented by an exposure profile are stored in the GPU device memory. In order to avoid transferring that data to host memory to perform the final computation and averaging, the CUBLAS function `cublas_gemv` is used to perform the multiplication and summation inside the GPU device memory. The `cublas_gemv` makes extensive use of Fused Multiplies and Add (FMA) instructions which help to increase throughput with the potential cost of less accuracy.

As the number of simulated exposure profiles increases, the application is exposed to floating point summations problems. When adding a huge number of floating-point numbers using the recursive summation approach the application ends up by experiencing round-off error. This problem of floating-point accuracy leads the application to provide different results after different runs or invalid floating-point results. For details on how to address floating point issues in numerical methods see [HIGH] The pair-wise summation is a better choice than the recursive summation, which results in a summation tree just like the one we use in our parallel reduction. Thus, the parallel reduction efficient on GPUs and it can improve accuracy. [PODL]

# Chapter 5

## Optimization

Different optimization acceleration strategies are explained in this chapter. Aiming to improve kernel execution throughput three techniques have been depicted:

- Modify the kernel execution configuration
- Take advantage of data locality facilities
- Scale the application across multiples GPUs.

### 5.1 Exploiting oversubscription

The parallelized simulation performance can be improved by analysing the constraints between the software and the underlying accelerator hardware. In this case a thread block in *CUDA* could be composed of up to 1024 threads. The current implementation divides kernels by flattening a 2D problem into one dimension and splitted the processing work between different SM of the same GPU device at *CUDA* thread grid level.

It would be possible to improve performance by increasing only the thread group size as a potential processing throughput increase is expected due to more jobs per thread group. The number of total simulations increases with the number of scenarios. For instance, to obtain 3000 exposure profile realizations, we need 70,000 simulations, if we jump to 10000 exposure profile realizations implies 25000000 simulations.

The execution configuration for kernel \_\_generatedPaths\_kernel is specified as following:

```
// Kernel Execution Parameters
int N = 1;
int blockSize = N* BLOCK_SIZE;
int numberofPoints = N * TIMEPOINTS;
int gridSize = scenarios_GPUs / N;
```

Rather than following the standard way of increasing the amount of work per thread or applying a loop tiling, the approach described here exploits oversubscription. Therefore, there is a map between one or multiple scenarios and a thread group. The *CUDA* Kernel Execution Configuration is defined by the number of blocks or grid size and the number of threads per block. The gridSize is equal to the number of scenarios to be simulated. The blockSize is basically the nearest multiple of 64 that is greater than the total number of points required to simulate the whole rates curve.

As the value of constant N increases, the block Size go up by N times which means that a smaller number of threads blocks per grids need to be distributed across the SM. This results in more work to be performed by thread blocks but constrained by the maximum amount of register per SM and the total number of blocks executed that can be active inside the device for the execution.

What has been benchmarked here is the sweet spot of the number of threads per block based on the code of the \_\_generatePaths kernel. The simulation of 50,000 scenarios was benchmarked using 128, 256, 512 threads per blocks. None of the increased blocks size configuration managed to improve the application throughput due to excessive global memory access. This indicates that using more threads do not necessarily imply higher throughput. [VOLK1][VOLK2]

## 5.2 Using Shared Memory

As shown before the CPU code executes a CUDA kernel function until all simulations paths are reached. Although the results are correct the executed kernel didn't make an efficient use of the memory due to access to global memory to fetch the previous data on each simulation step across all simulations. Hence, reducing the global memory traffic in `__generatePaths_kernel` would help to keep Symmetric Multiprocessors (SM) while the simulation kernel is running.

To simulate the rate with *stochastic differential equation* on each pricing point requires 18 memory access and 6 floating point operations. Calculating the compute-to-global-access-memory ratio gives a result of 0.3 which is an indication that kernel `__generatedPaths_kernel` can be considered as a memory bound.

The CUDA programming library allows developers to access faster memory located inside the SM. That's the case of *registers* and *shared memory*. *Registers* are private to the thread under execution, but *shared memory* is a kind of scratch path memory that can be shared by all threads that belong to the same thread block. Since *registers* and *shared memory* are inside the SM chip the memory access is faster in several orders of magnitudes with respect to the global memory.

The kernel function `__generatedPaths_kernel` provides the following performance sweet spots:

- Data elements are loaded multiple times for different threads inside the same thread group.
- Data is only required to be stored in global memory every 50 simulation steps.
- The size of the tile is small enough to schedule different concurrent threads blocks inside the same SM

As explained before each thread group will generate a grid of 2500 x 51 simulated rates. Rather than reserving memory for the whole grid only a portion or a tile is created and used by the thread group. Only a fraction of the memory is needed to generate the whole simulation grid for all thread groups. This type programming pattern is known as a *tiling*. Each tile belonging to one thread group basically occupies a fraction of the shared memory. The current hardware in use supports 128k shared memory size per SM. There is a trade-off with the usage of shared memory in the sense that it's faster but small when the global memory is bigger but slower. The kernel should use a tile that fits inside the shared memory and is small enough to allow other thread blocks to run in parallel. It has been proved before that shared memory is a key accelerator factor for stencil/tiled based code (WEIK).

The kernel function `__generatedPaths_kernel4` addresses the potential tuning opportunities exposed in the previously discussed before. First, the total number of times the CPU invokes the CUDA kernel is reduced by a factor of 50. In case of 50,000 scenarios are needed the number of calls to invoke the CUDA kernel drop down considerably from 125,000000 to 2,500000. Therefore, less global memory traffic was required during the simulation process.

All references to code in the rest of this section will be based on the CUDA kernel implementation `__generatedPaths_kernel4` shown in Appendix H.

The usage of shared memory is not only about speeding up memory access but also reusing data by different threads inside the same thread group. A float array of size 64 elements is instantiated (Line 493) and allocated in the shared memory. The CUDA directive `__shared__` signals the *CUDA* compiler that the variable `_ssimulated_rates[]` will be stored in shared memory. Each thread in the thread group access global memory to fetch the rate generated in the previous simulation and store the data in the shared memory location array (Lines 507 – 512). To simulate the *stochastic differential equation* each thread looks at the previous intermediate simulation value stored in the shared memory array and makes a difference of two contiguous elements. It happens that for contiguous threads in the thread group this data is reused and stored in shared memory (Line 526 – 531). Bear in mind that for all intermediate simulations with path [n, n+stride], where

stride size is equal to 50 there is no need to write global memory access on every simulation step. Observe that in line (512) all threads in the thread group executes at the same time but it's mandatory to synchronize them to avoid invalid access to memory that has not been initialized.

The top compute and memory access section of the kernel `__generatedPaths_kernel4` occurs between lines (515 – 571). The logic is almost the same as kernel `__generatedPaths_kernel` except that all memory access to simulated data is performed on shared memory. Therefore, a huge performance impact in terms of execution time is noticed as explained in the performance results section.

As elements of the same array `_ssimulated_rates` are reused between threads. The access to memory is also reduced by half. Allowing this an increase on the floating-point computation per global memory access. The `__synchronized` method to coordinate the update of array variable `_ssimulated_rates[]` which is stored in shared memory (Lines 562 and 569). The rest of the implementation is like `__generatedPaths_kernel` except that the latest simulated rate is the one to be stored in global memory. It would also be possible to improve performance by allowing each thread to process more than one rate element at a time. The section of Performance Analysis will show how the kernel `__generatedPaths_kernel4` outperform the rest of the accelerated simulation implementation.

### 5.3 Scaling with multiple GPUs

Up to this point it have been explained the *CUDA* implementation and optimization of the application in a single GPU. However, most GPU systems include multi-GPUs devices on each node. A multi-GPU implementation is the topic covered in this chapter.

The command `nvidia-smi` displayed the list of GPU devices present in the multi-GPU system (Figure 11). Observe how each GPU is labelled by number from 0 up to 3. The aggregate global memory increased from 8GB to 32 GB and the throughput processing capability is also increased if work is distributed evenly across all GPU devices. In this case multiple-GPU are connected over the same PCIe bus in a single node.

The simulation of each exposure profile curve is done independently. Therefore, to take advantage of the existence of more than one GPU-device the exposure profile simulation workload needs to be partitioned in such a way that each GPU equally shares an equal amount of work to do. The calculation is self-contained per GPU and, there is no dependency between partitions, and therefore no data exchange is required across GPUs.

To provide parallelism across multiple GPU, the set of exposure profile simulations is divided into equals the number of CUDA-capable GPU installed in the system. The problem domain size is divided by the number of GPUs, processed in parallel and therefore the application improves throughput and efficiency. Once all computations are finished and output data written back to global memory per GPU device the host code transfer the partial exposures profiles from each sibling GPU to host memory and accumulates the results as a vector summation.

The application buffers initially allocated across all exposure profiles are now divided in chunks of a quarter of the total size and created on global memory for each GPU device. Analogy, constant memory data is initialized on each GPU as well. Global memory allocation operations invoke the method `cudaSetDevice(GPUDevice)` where `GPUDevice` parameters identify the GPU among others in the multi-GPU system. NVIDIA CUDA API manages and executes kernels on multiple GPUs using an asynchronous programming model. In order to execute the same kernel across multi-GPU the host code need to create a thread per each GPU device and associate it with an individual GPU CUDA context.

```
C:\ProgramData\NVIDIA Corporation\CUDA Samples\v10.2\bin\win64\Release>nvidia-smi -L
GPU 0: GeForce RTX 2070 SUPER (UUID: GPU-49699fa2-d04c-ce2e-1036-63ede770f6e1)
GPU 1: GeForce RTX 2070 SUPER (UUID: GPU-9da01d6e-0ce2-764e-a1e5-b6bd89f5581a)
GPU 2: GeForce RTX 2070 SUPER (UUID: GPU-cbe94ae7-d531-10ad-c58c-2b85fe12ac37)
GPU 3: GeForce RTX 2070 SUPER (UUID: GPU-a09918d9-1f7b-39ac-44d1-c32af26e8c89)

C:\ProgramData\NVIDIA Corporation\CUDA Samples\v10.2\bin\win64\Release>nvidia-smi
Sun Jun 06 09:27:01 2021
+-----+
| NVIDIA-SMI 456.71      Driver Version: 456.71      CUDA Version: 11.1 |
+-----+
| GPU  Name           TCC/WDDM   Bus-Id     Disp.A  Volatile Uncorr. ECC |
| Fan  Temp  Perf  Pwr:Usage/Cap | Memory-Usage | GPU-Util  Compute M. |
|-----+-----+-----+-----+-----+-----+-----+-----+
| 0  GeForce RTX 207... WDDM    00000000:02:00.0 Off   N/A
| 48%   26C     P8    19W / 215W | 416MiB / 8192MiB | 5%       Default |
|-----+-----+-----+-----+-----+-----+-----+-----+
| 1  GeForce RTX 207... WDDM    00000000:17:00.0 Off   N/A
| 49%   26C     P8    5W / 215W | 262MiB / 8192MiB | 0%       Default |
|-----+-----+-----+-----+-----+-----+-----+-----+
| 2  GeForce RTX 207... WDDM    00000000:65:00.0 Off   N/A
| 44%   26C     P8    17W / 215W | 218MiB / 8192MiB | 0%       Default |
|-----+-----+-----+-----+-----+-----+-----+-----+
| 3  GeForce RTX 207... WDDM    00000000:B3:00.0 Off   N/A
| 48%   27C     P8    17W / 215W | 185MiB / 8192MiB | 0%       Default |
+-----+-----+-----+-----+-----+-----+-----+-----+
```

Figure 11. Multi-GPU system configuration.

The OpenMP with support of the Visual Studio C++ compiler allows the creation of multiples threads to parallelize the execution of the sequential code. Therefore, it was convenient in the sense that each cuda-context is required to be executed on a different thread. The number of threads to be used depends on the number of GPUs present in the system. An OpenMP parallel region with 4 threads equals to the number of GPUs is created (Appendix I Line 1079). Each OpenMP thread is assigned a unique identifier from 0 to 4 and associated with the GPUDevice number (Appendix I Line 1083). Once the host thread has been associated with the GPU device the context is established and therefore, with the help of the function riskFactorSim4 the kernel `_generatePaths_kernel4` is invoked on each GPU (Appendix I Lines 1094).

Gaussian variates are generated in blocks for each GPU device independently. All pseudo random variates are generated and stored in memory using the cuRand API beforehand. A different block of gaussian variates is generated for each GPU. The offset parameter is based on the GPUDevice number times the number of random number chunk size. That way the sequence of Randoms will not be repeated across all the GPU devices. (Appendix A)

The exposureCalculation kernel is executed independently across all the GPU devices as data is locally stored on each GPU device. Initially the result of the exposure profiles is scattered across the four different GPU devices global memory location. Therefore, no data exchange occurs between host nor GPUs. The host code performs a gather operation to fetch the partial accumulated exposures profiles and accumulate them, by calling the function `vsAdd`, into the host variable `partial_exposure`.

The modifications required to implement the multi-GPU version of the code only required to introduce the multithreaded code and the scatter and gather pattern. The CUDA kernel did not suffer any change but executes in parallel across different GPUs devices with less amount data. Exploiting the nature of the algorithm that allows to partition the data in smaller chunks to be processed independently for each device was a key driver for this implementation.

# Chapter 6

## Performance Results

The performance results are discussed in this chapter for each of the optimization techniques introduced before. The use of shared memory and multi-GPU is expected to yield significant performance speedup. The execution time of the top consuming CUDA kernel is improved across different runtime configurations and optimization techniques. The benchmark performed during this work is based on 4 GPU cards NVIDIA RTX 2070. More details about the hardware can be found in Appendix J. The top consuming kernel `_generatePaths_kernel` was the focus for analysis and optimization. As it was the higher contributor to computation time latency.

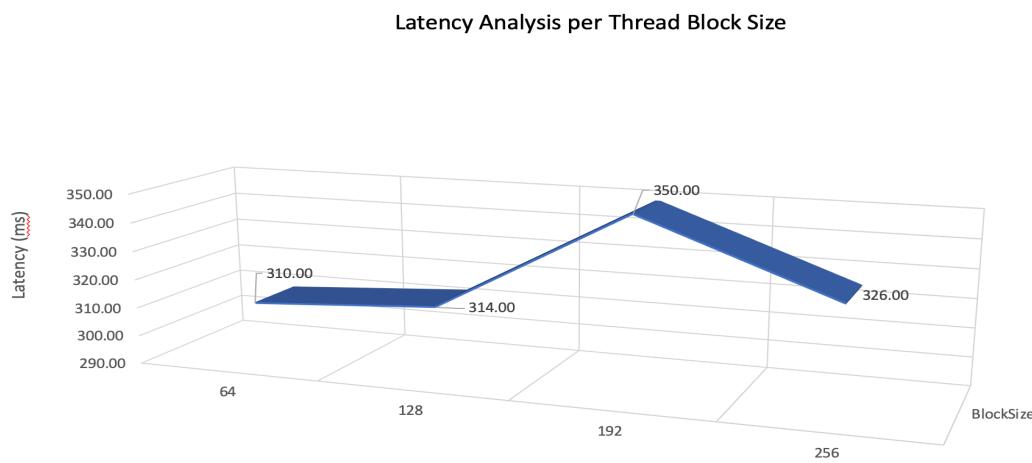


Figure 12. Latency Analysis per thread block size.

The first exploration on performance improvement was to increase the number of threads per block for the `_generatePaths_kernel` kernel execution configuration. As shown in Figure 12 four different blocks sizes 64, 128, 192, 256 were configured to run the 50,000 scenarios simulation test in single-precision mode. Due to the huge global memory traffic it was not possible to improve the performance. Even though 256 thread block size was the perfect size of thread block the best latency results were obtained with block size of 64 threads. In this case increasing the thread block size does not guarantee a performance improvement. [VOLK2]

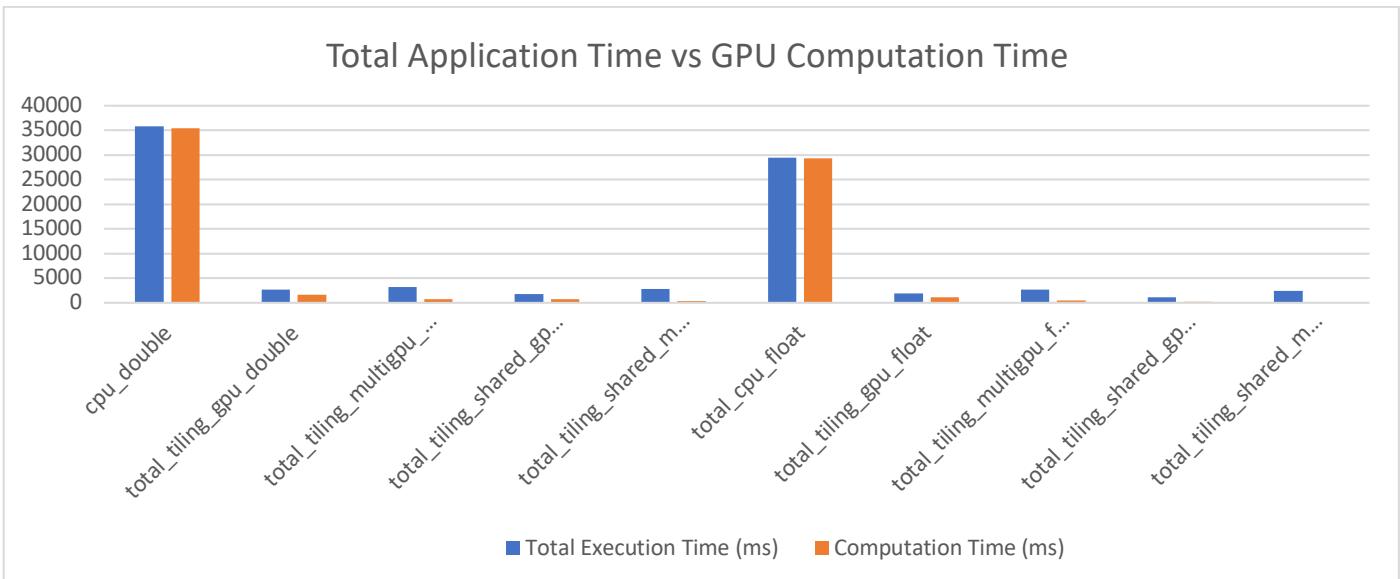


Figure 13. Latency comparison between CPU, GPU and multi-GPU solution.

Implementation	Total application execution-time (ms)	Initialization-time (ms)	Computation-time (ms)	Computation-time Speed up
cpu_double	35777		35461	
tiling_GPU_double	2642	852.49	1583	22.40
tiling_multiGPU_double	3154	1816	500	70.92
tiling_shared_GPU_double	1765	772	543	65.30
tiling_shared_multiGPU_double	2846	1824	150	236.40
cpu_float	29479		29279	
tiling_GPU_float	1916	1837	967	30.27
tiling_multigpu_float	2674	1837	342	85
tiling_shared_GPU_float	1064	751	157.486	185.91
tiling_shared_multigpu_float	2153	1808	46	636

Table 1. Latency CPU, GPU and multi-GPU solution for 50000 scenarios.

Figure 13 and Table 1 shows the total execution time of the application and the computation time measured in milliseconds. According to Table 1 all kernel names suffixed with double run in double precision and those kernels names suffixed with float run in single precision mode. The kernel `_generatePath_kernel_cpu_double` is the reference kernel to calculate the speed up for all double precision kernels. The kernel `_generatePath_cpu_float` is the reference kernel to calculate the speed up for all single precision kernels. All kernels implemented the tiling optimization technique and those using shared memory contains the shared word in their names. Single GPU and multiGPU specified if the application run the kernels across one or multiples GPUs. The speedup gained from parallelizing some of the applications kernels in GPU is the ratio of the single CPU execution time to the GPU parallel execution time ( $T$ ). Hence, the computation time Speed up =  $T(\text{cpu}) / T(\text{gpu})$ .

The computation time was the dominant kernel when running the application in CPU. Once the simulation kernel was ported to GPUs the dominant time of the application changed due to the acceleration speed up provided by the GPU solution. It's not hard to appreciate that initialization time and random generation kernels are the major latency contributors after the computation kernel was accelerated. The test was configured to simulate 50000 scenarios in single and double precision floating point accuracy. The best combination of total application execution time and computation-time speed up is provided by the implementation `tiling_shared_GPU_float` which makes extensive use of shared memory and run-in single precision. Nevertheless `tiling_shared_multigpu_float` resulted in the kernel with best computation time speed up.

Implementation	Rnd generation (ms)
cpu_double	1373
tiling_GPU_double	87
tiling_multiGPU_double	52
tiling_shared_GPU_double	167
tiling_shared_multiGPU_double	51
cpu_float	932
tiling_GPU_float	12
tiling_multiGPU_float	9
tiling_shared_GPU_float	13
tiling_shared_multiGPU_float	10

Table 2. Time spent in Gaussian variates generation

The single precision (float) implementation outperformed the double precision. Nevertheless, the float implementation lost accuracy in the tenor 50 of the expected exposure curve due to floating point round off. According to Table 1 and Table 2 all kernel names suffixed with double run in double precision and those kernels names suffixed with float run in single precision mode. The kernel \_\_generatePath\_kernel\_cpu\_double is the reference kernel to calculate the speed up for all double precision kernel. The kernel \_\_generatePath\_cpu\_float is the reference kernel to calculate the speed up for all single precision kernels.

The table 3 shows the mean square error and the accuracy achieved during the benchmark. Notice how the results in double precision using shared memory are slower but more accurate.

Implementation	Mean-square errors	Error
tiling_GPU_double	0.01377044	10^-2
tiling_multiGPU_double	0.01397286	10^-2
tiling_shared_GPU_double	0	10^-6
tiling_shared_multiGPU_double	0.00059256	10^-4
tiling_GPU_float	0.01380047	10^-2
tiling_multiGPU_float	0.0143285	10^-2
tiling_shared_GPU_float	0.000006	10^-6
tiling_shared_multiGPU_float	0.003549	10^-3

Table 3. Floating Point Accuracy and mean square errors

By doing a breakdown of the computation time it can be appreciated that the multi-GPU kernel implementation combined with the usage of shared memory and single floating-point accuracy provides an acceleration of 632x. Which somehow indicates that the kernel execution scales linearly with the number of GPU present in the system. Due to the low latency of the generate\_kernel\_path4 of nearly 30ms per GPU device running in parallel. The point here is the trade-off between speed and accuracy. Being the latter is a strong requirement in the financial domain. Table 4 shows the *Monte Carlo* simulation time for each of the CUDA kernels analysed in this study. The programming techniques used to implement access to shared memory somehow had a positive impact in the accuracy of the recursive summation results. See also Table 3 to check the accuracy for implementation kernels tiling\_shared\_GPU\_float and tiling\_shared\_GPU\_double.

MC Simulation kernel name	MC simulation execution time (ms)	MC Simulation speed up
__generatePath_kernel_cpu_double	21114	-
__generatePath_kernel_double	1470	14.36326531
__generatePath_kernel_double	509	41.48133595
__generatePath_kernel4_double	540	39.1
__generatePath_kernel4_double*	164.59	128.2823987
__generatePath_kernel_cpu_float	19335	-
__generatePath_kernel_float	1030	18.77184466
__generatePath_kernel_float	397	48.70277078
__generatePath_kernel4_float	111.863	172.8453555
__generatePath_kernel4_float*	30.5459	632.9818404

Table 4. Monte Carlo Simulation execution time. (\*) multi-GPU

The compute and memory access characterization of kernels \_\_generatePaths\_kernel and \_\_generatePaths\_kernel4 are compared using the metrics collected in the roofline model [WWP] for simple precision floating point accuracy. The ratio of compute to memory accesses for a kernel is called its arithmetic intensity. The estimated flop's performance value is also calculated for both kernels. In Figure 14 the arithmetic intensity is represented on the x-axis and the performance or flops on the y-axis. The thick magenta points plotted in Figure 14 represent both kernels. The magenta point with the red circle represents the kernel \_\_generatePaths\_kernel4 and the magenta point with blue circle represent the kernel \_\_generatePaths\_kernel. The straight lines through the graph represent the limit in terms of performance for double precision 0.5 teraflops and the other straight line at the top represents the limit in terms of performance for single precision 10 teraflops. The diagonal roof represents the limit on a kernel's peak flops imposed by memory bandwidth constraints for a single and double precision arithmetic intensity.

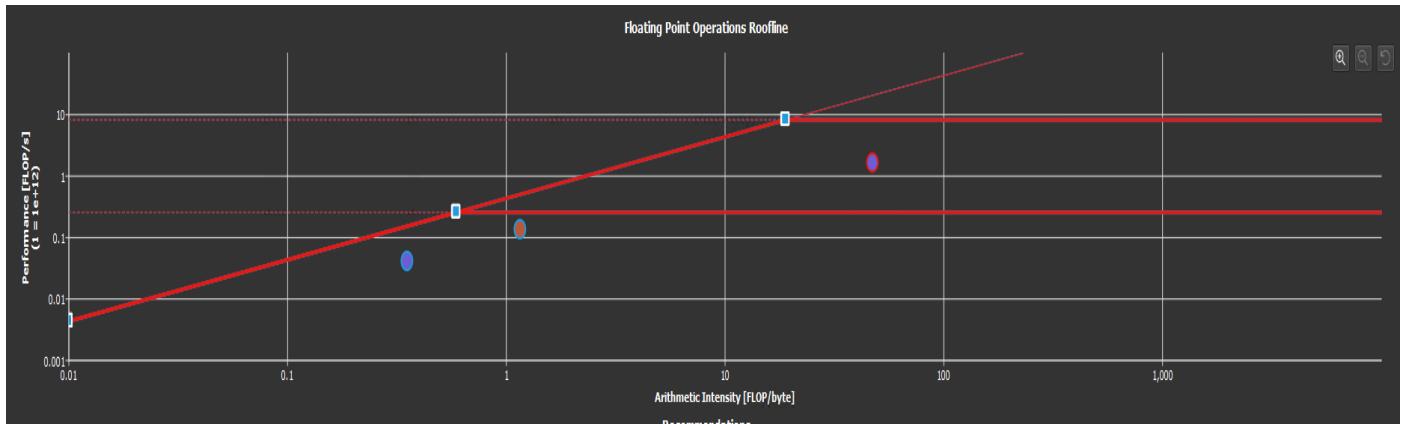


Figure 14. Roof line analysis comparison between \_\_generatePaths\_kernel and \_\_generatePaths\_kernel4

The kernel \_\_generatePaths\_kernel achieved an arithmetic intensity of 1.16 flops/bytes and the flop's performance value is 130 gigaflops. Given the fact that the \_\_generatePaths\_kernel measures were taken for single precision the observed performance is poor because the magenta dot with the blue circle represented by the \_\_generatePaths\_kernel is under the roof for double precision limits. A kernel using single floating-point precision is faster than the same version of it running in double precision. The poor performance could be related to the excessive not coalesced global memory access in all the simulation steps. In contrast, kernel \_\_generatePaths\_kernel4, represented in Figure 14 as a magenta dot with a red circle, is closer to the roof of floating-point precision limits. The improvement in arithmetic intensity for single floating-point accuracy is 47.5 FLOPs/bytes and a flops performance over 1 TB/s. The improvement in performance is in part due to the

use of shared memory to improve the application locality, data reusing between threads inside the same thread block, and partial reduces of all global memory access by using the fast on chip shared memory for intermediate simulation results.

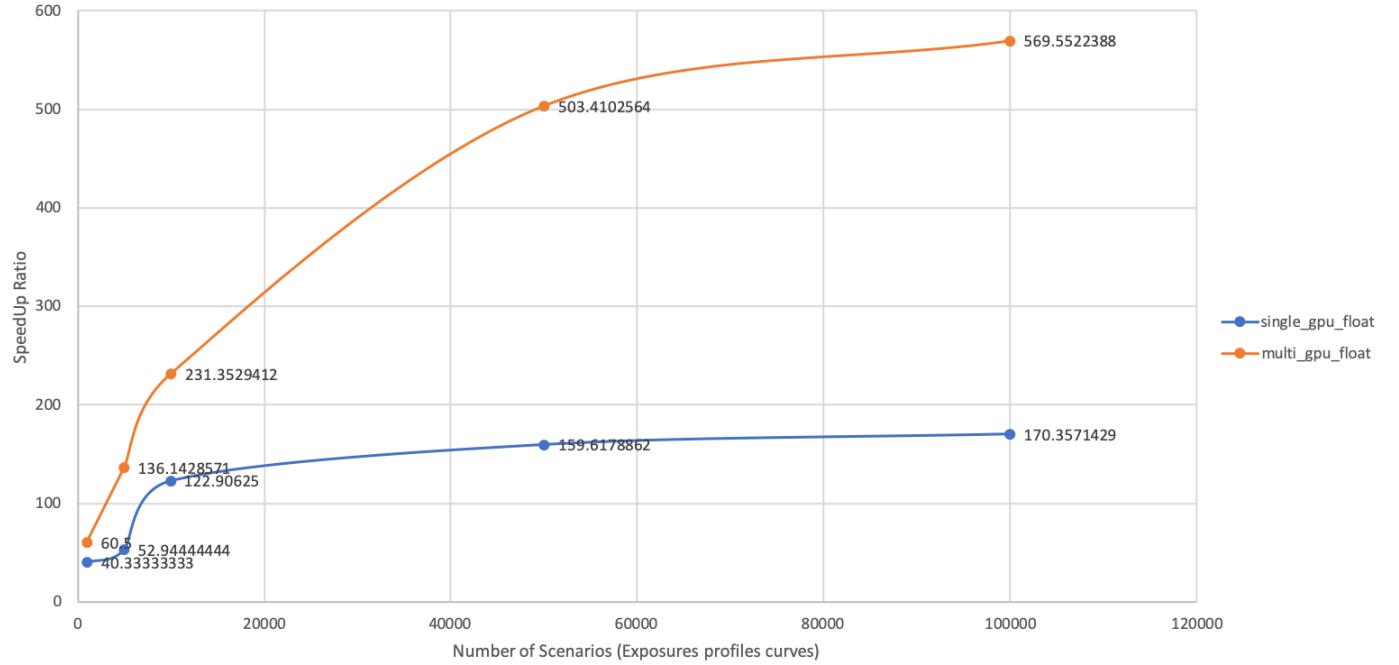


Figure 15. Speed Up ratio of kernel \_\_generatePaths\_kernel4 on NVIDIA GeForce RTX 2070 fp 32 in/out fp 32 compute

The single precision speed up ratio achieved by kernel \_\_generatePath\_kernel4 is shown in figure 15. The multi-GPU invocation of kernel \_\_generatePath\_kernel4 scales linearly as the number of scenarios (number of simulated exposures profiles) increases. The single GPU solution running in single precision flattened the scaling curve factor after reaching 50000 simulated scenarios. In contrast the multi-GPU solution increased their speed up to a value close to a ratio value of 600. This observation shows how the Monte Carlo simulation kernel scales linearly as the kernel \_\_generatePath\_kernel4 is executed across 4 different GPUs in a multi-GPU system.

# Chapter 7

## Conclusion

The CVA calculation can be considered one of the computation dwarfs and the problem has been shown that can be tackled using parallel heterogeneous programming. It has been shown how GPUs can be used to accelerate the calculation of CVA interest rate swap. The GPU accelerated implementation outperforms the CPU implementation in computing the expected exposure profile by using Monte Carlo simulation to evolve the risk factors. The usage of multi-GPU improved even better the final benchmarks results as a larger number of scenarios for an interest rate swap exposure profiles increased. Although the multi-GPU reduces the time spent in simulation time the generation of random numbers penalized the whole calculation time and shows a slight drop in accuracy with respect to the single GPU implementation. It has been illustrated the use of GPU acceleration on NVIDIA hardware, CUDA Programming can achieve a considerable speed up for embarrassing parallel problems like applying Monte Carlo to simulate the stochastic differential equation that represents the Gaussian Heath Jarrow Merton model.

There are several interesting extensions and areas for further work within the area of XVA calculations. First, solving stochastic differential equations numerically by using Monte Carlo simulation is not only about path generation and averaging but also increasing the convergence rate by using variance reduction methods. To improve the floating-point accuracy a fix-point arithmetic could be implemented or incorporate the usage of the predictor-corrector path generation as the number of simulation steps increased [GLAS][JAECK]. Although we achieved good performance results on the acceleration side it would be better if we could increase the amount of calculation work per thread, leading to process more than one exposure profile per thread group. For better utilization of the SM a reduced number of blocks could generate all the required exposure profiles by applying the grid-stride loop pattern [MHA4]. It would be equally interesting to implement a stream-based solution since it would be possible to launch more kernels per SM rather than launching only one kernel at a time. The expected exposure profile would be calculated by using parallel reduction for better accurate and stable results as a result of large floating-point summations [HIGH][MHA3]. An extension of the current work would be the instrumentation of *Monte Carlo* path generation to calculate path-wise sensitivities [SAV][GLAS].

I strongly believe that High Performance Computing is a strong technology that not only helps to obtain faster and accurate results but also will contribute with efficiency to achieve greener energy compliant investment banking.

# Chapter 8

## References

- ASHL Ashley, John. (2007). Computational Finance with GPUs : what is next  
BIS OTC derivatives outstanding. <https://www.bis.org/statistics/derstats.htm>  
BUK England, Bank of. (2009) Yield Curve Terminology and Concepts  
CGR Cornelis and Grzelak. (2019) Mathematical Modeling and Computation in Finance  
HDE Higham, Desmond. (2001). An Algorithmic Introduction to Numerical Simulation of Stochastic Differential Equations. SIAM  
HIGH Higham, Nicholas J. (2002). Accuracy and Stability of Numerical Algorithms. 2nd Edition. SIAM  
HOU Houston University, An OpenACC example code for c-based heat conduction code  
GAT Dariusz Gaterek, Przemyslaw Bachert and Robert Maksymiuk. The LIBOR Market Model in Practice  
BMM Alen Brace, Marek Musiela. A multifactor Gaus Markov implementation of Heath, Jarrow and Morton  
JAECK Jaeckel, Peter. (2000). Monte Carlo Methods in Finance.  
MHA Harris, Mark. (2000). An even easy introduction to CUDA.  
NOV Novosyolov, Arcady (2008). Global Term Structure Modeling using Principal Component Analysis.  
MHA2 Harris, Mark. (2000). Parallel Prefix Sum (Scan) with CUDA.  
MHA3 Harris, Mark. (2000). Parallel Reduction with CUDA.  
MHA4 Harris, Mark. (2008). CUDA Pro Tip: Write Flexible Kernels with Grid-Stride loops  
PODL Podlozhnyuk, Victor and Harris, Mark. (2008) Monte Carlo Option Pricing.  
PYK Pykhtin, M, and Zhu, S. (2007). A Guide to Modelling Counterparty Credit Risk  
NV01 NVIDIA, Curand User Guide.  
NV02 NVIDIA, CUDA C Programming Guide.  
NV03 NVIDIA, CUDA Best Practices Guide.  
NV05 NVIDIA, Precision and Performance Floating point IEEE754 Compliance for NVIDIA GPU  
SAV Savigne, Antoine. (2018). Modern Computational Finance, AAD and Parallel Simulations.  
RUI Ruiz, I. (2015). XVA Desks – A New Era for Risk Management Understanding, Building and Managing Counterparty, Funding and Capital Risk. Palgrave Macmillan  
THO Thomas-Collignon, Thomas. (2018). Volta Performance architecture and performance optimization.  
VOLK1 Volkov, Vasily. (2018). Understanding Latency hiding on GPU.  
VOLK2 Volkov, Vasily. (2018). Better Performance at Lower Occupancy.  
WEIK Wen Wei and Kirk. (2018). Programming Massively Parallel Processors. 3rd edition.  
WILL Willmot, Paul. Introduces Quantitative Finance.  
WWP Williams, Samuel Patterson, David. Roofline: An Insightful Visual Performance Model for Floating-Point Programs and Multicore Architectures

# Appendix A

## Gaussian random variates generation

Pseudorandom sequence: a sequence of numbers, generated by a deterministic algorithm, that has most of the properties of a truly random sequence. The pseudo-random numbers will be used in the case study example to generate gaussian variates in the sde model numerically solved by using Monte Carlo Methods. The Musiela Parametrization stochastic differential equation requires 3 different pseudo random normal distributed variates on each simulation timestep.

The NVIDIA cuRAND library provides a CUDA-based library for rapid pseudo-random number generation.

Pseudo Random Generator Options:

- RNG Algorithm: An RNG algorithm with which to generate a random sequence
- Statistical Distribution: a distribution to which the returned values will adhere
- Seed: A 64-bit integer that initialize the starting state of the pseudorandom number generator
- Offset: A parameter used to sort skip ahead in the sequence. Set offset = 100 to return the 100th number in the sequence first.

Although RNGs could be generated inside the main simulation kernel it was preferred to generate a large batch of RNGs on GPU and store them in global memory for the HJM simulations. This prevents us from having early thread divergence in code as all threads in a block share the same 3 RNGs numbers on each simulation step. The number of RNGs generated was equal to the number of 3 times HJM simulations times the number of blocks.

Figure 16 shows a sample distribution of independent pseudo random number variates for 2500 simulation time steps using the NVIDIA cuRAND library.

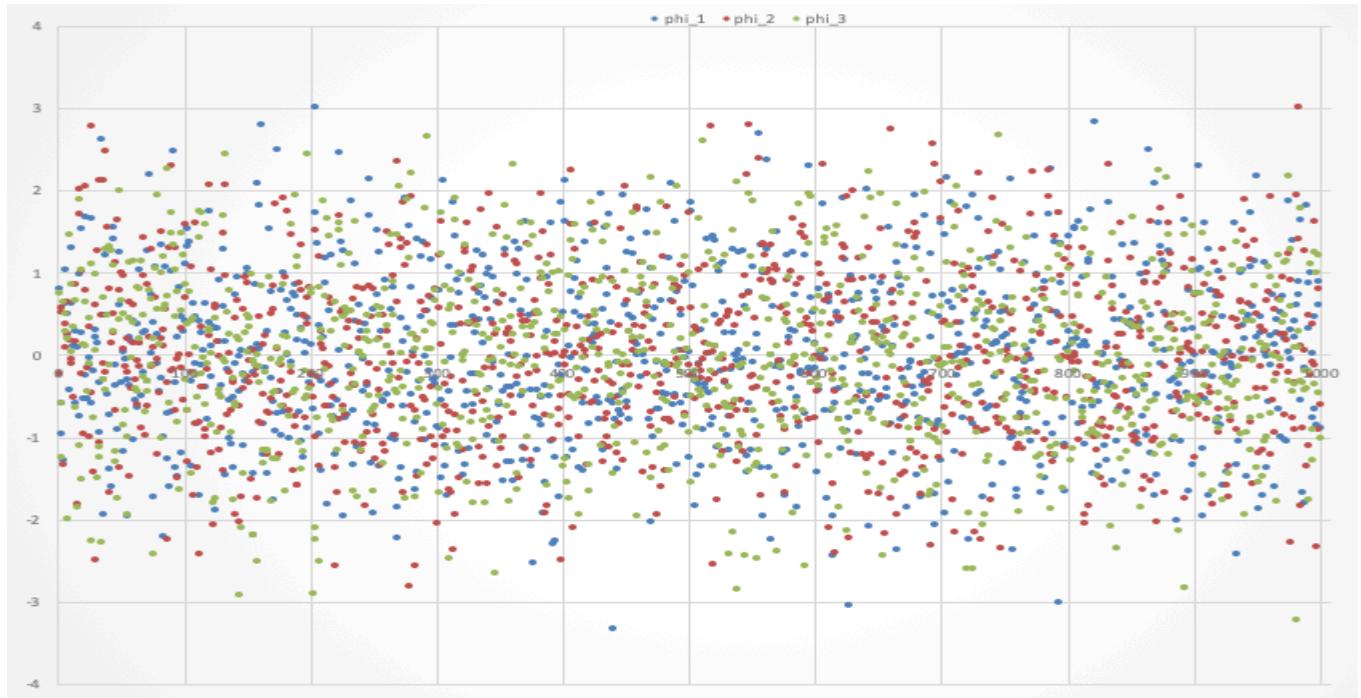


Figure 16. Independent Random Number Variates (0, 1.0)

```

100  /**
101   * * RNG init Kernel
102   */
103
104 #ifndef RNG_HOST_API
105 void initRNG2_kernel(float* rngNmVar, const unsigned int seed, int rnd_count)
106 {
107     const float mean = 0.0; const float stddev = 1.0;
108     curandGenerator_t generator;
109     CURAND_CALL( curandCreateGenerator(&generator, CURAND_RNG_PSEUDO_DEFAULT) );
110     CURAND_CALL( curandSetPseudoRandomGeneratorSeed(generator, 1234ULL) );
111     CURAND_CALL( curandGenerateNormal(generator, rngNmVar, rnd_count, mean, stddev) );
112     CUDA_RT_CALL( cudaDeviceSynchronize() );
113     CURAND_CALL( curandDestroyGenerator(generator));
114 }

```

Figure 17. Using the NVIDIA Host API for pseudo random generation.

As discussed, 3 normal distributed random variates are generated per simulation timestep. Where rnd\_count represents the total number of pseudo random numbers to generate. Given \_simN (total number of exposure profiles) equals to 50000, pathN as the result of dividing dtau/dt equals to 2500 which represents the number of paths required to simulate the Musiela SDE with 3 factors (actual number of pseudo random presents in the diffusion term of the SDE  $f(t+dt)$  [4.3.2.1] a total number of 37500000 floating RNG values are generated.

The pseudo random algorithm used by the application is the one that comes with the value CURAND\_RNG\_PSEUDO\_DEFAULT (Figure 30). Usually, financial quant libraries use other pseudo random generators of better quality as Mersenne Twister or MRG32k3a [JAECK] but for simplicity the default pseudo random generator CURAND\_RNG\_PSEUDO\_XORWOW provided by NVIDIA Developer Toolkit was used for this case. The function curandGenerateNormal is used to generate a normal distributed value with variance 1.0 and mean 0.0.

The cuRAND API allows to create and store the RNG values in device memory. That way accessing this information although in global memory is more efficient than moving data from RAM to Device Global memory. Figure 30 shows how the cuRAND API is used by the application developed for this study.

It has been observed a very low occupancy of their GPU device code as it is required to access an array of random numbers from global memory. The ideal is to find a balance between performance and quality of random numbers, which is problem dependent. The most important part of cuRAND-based projects is to understand the randomness requirements.

# Appendix B

## Simulation Grid of a Gaussian HJM Model

# Appendix C

## Numerical Simulation Driver

```
409 int blockSize = payOff.dtau / dt;
410
411 std::cout << "intermediate simulations " << blockSize << "number of paths" << pathN << std::endl;
412
413 for (int path = 0; path < pathN; path += blockSize)
414 {
415     // open mp parallel region
416     for (int s = 0; s < exposuresCount; s++)
417     {
418         // internal openmp vectorization
419         __generatePaths_kernel_cpu(
420             &numeraires[s * _TIMEPOINTS],
421             _TIMEPOINTS,
422             drift,
423             volatilities,
424             &rngNrmVar[s * pathN * 3],
425             spot_rates,
426             &simulated_rates[s * _TIMEPOINTS],
427             &simulated_rates0[s * _TIMEPOINTS],
428             &accum_rates[s * _TIMEPOINTS],
429             pathN,
430             path
431         ); //dt = 0.01, dtau = 0.5
432     }
433
434     // update simulated rates (swap pointers)
435     std::swap(simulated_rates, simulated_rates0);
436 }
437
438 t_end = std::chrono::high_resolution_clock::now();
439 elapsed_time_ms = std::chrono::duration<double, std::milli>(t_end - t_start).count();
440 std::cout << "total time taken to run all " << pathN * exposuresCount << " HJM MC simulation " << elapsed_time_ms << "(ms)" << std::endl;
```

# Appendix D

## Naïve Parallelized Numerical Simulation Algorithm

```
720  /*
721   * Risk Factor Generation naive acceleration
722   */
723
724  void riskFactorSim(
725      int gridSize,
726      int blockSize,
727      int numberofPoints,
728      double2* numeraires,
729      double* rngNrmVar,
730      double* simulated_rates,
731      double* simulated_rates0,
732      double* accum_rates,
733      const int pathN,
734      double dtau = 0.5,
735      double dt = 0.01)
736  {
737
738      for (int path = 0; path < pathN; path++)
739      {
740          __generatePaths_kernel <<< gridSize, blockSize >>> (
741              numberofPoints,
742              numeraires,
743              rngNrmVar,
744              simulated_rates,
745              simulated_rates0,
746              accum_rates,
747              pathN,
748              path,
749              dtau,
750              dt
751          );
752
753          CUDA_RT_CALL(cudaDeviceSynchronize());
754
755          // update simulated rates (swap pointers)
756          std::swap(simulated_rates, simulated_rates0);
757      }
758  }
```

# Appendix E

## Accelerated Numerical Simulation Algorithm

```
590  /*
591   * Risk Factor Generation block simulation with Shared Memory
592   */
593
594  void riskFactorSim4(
595      int gridSize,
596      int blockSize,
597      int numberofPoints,
598      double2* numeraires,
599      double* rngNrmVar,
600      double* simulated_rates,
601      double* simulated_rates0,
602      double* accum_rates,
603      const int pathN,
604      double dtau = 0.5,
605      double dt = 0.01)
606  {
607      int simBlockSize = dtau / dt;
608
609      for (int path = 0; path < pathN; path += simBlockSize)
610      {
611          __generatePaths_kernel4 <<< gridSize, blockSize >>> (
612              numberofPoints,
613              numeraires,
614              rngNrmVar,
615              simulated_rates,
616              simulated_rates0,
617              accum_rates,
618              pathN,
619              path,
620              dtau,
621              dt
622          );
623
624 #ifdef CUDA_SYNC
625         CUDA_RT_CALL(cudaDeviceSynchronize());
626 #endif
627
628         // update simulated rates (swap pointers)
629         std::swap(simulated_rates, simulated_rates0);
630     }
631 }
```

# Appendix F

## Accelerated Vanilla Interest Rate Swap Mark to Market

```

804 __global__
805 void _exposure_calc_kernel(real* exposure, real2* numeraires, real notional, real K, int simN, real dtau = 0.5f)
806 {
807     __shared__ real cash_flows[TIMEPOINTS];
808     real discount_factor;
809     real forward_rate;
810     real libor;
811     real cash_flow;
812     real sum = 0.0;
813     real m = (1.0 / dtau);
814
815     int globaltid = blockIdx.x * TIMEPOINTS + threadIdx.x;
816
817     // calculate and load the cash flow in shared memory
818     if (threadIdx.x < TIMEPOINTS) {
819         forward_rate = numeraires[globaltid].x;
820
821 #ifdef SINGLE_PRECISION
822         libor = m * (expf(forward_rate / m) - 1.0);
823 #else
824         libor = m * (exp(forward_rate / m) - 1.0);
825 #endif
826         discount_factor = numeraires[globaltid].y;
827         cash_flow = discount_factor * notional * d_accrual[threadIdx.x] * (libor - K);
828         cash_flows[threadIdx.x] = cash_flow;
829
830 #ifdef EXPOSURE_PROFILES_DEBUG
831         printf("Block %d Thread %d Forward Rate %f libor %f Discount %f CashFlow %f \n", blockIdx.x, threadIdx.x, forward_rate, libor, discount_factor, cash_flow);
832 #endif
833     }
834     __syncthreads();
835
836 #ifdef EXPOSURE_PROFILES_DEBUG
837     if (threadIdx.x == 0) {
838         for (int t = 0; t < TIMEPOINTS; t++) {
839             printf("t = %d CashFlow %f \n", t, cash_flows[t]);
840         }
841     }
842 #endif
843
844     // calculate the exposure profile
845     if (threadIdx.x < TIMEPOINTS) {
846         for (int t = threadIdx.x + 1; t < TIMEPOINTS; t++) {
847             sum += cash_flows[t];
848         }
849         sum = (sum > 0.0) ? sum : 0.0;
850         exposure[globaltid] = sum;
851
852 #ifdef EXPOSURE_PROFILES_DEBUG
853         printf("Block %d Thread %d Exposure %f \n", blockIdx.x, threadIdx.x, sum);
854 #endif
855     }
856     __syncthreads();
857 }
858 }
```

# Appendix G

## Simulation kernel naïve implementation

```

361  template <typename real, typename real2>
362  __global__
363  void __generatePaths_kernel(
364      int numberofPoints,
365      real2* numeraires,
366      real* rngNrmVar,
367      real* simulated_rates,
368      real* simulated_rates0,
369      real* accum_rates,
370      const int pathN, int path,
371      real dtau = 0.5, double dt = 0.01)
372  {
373      // calculated rate
374      real rate;
375      real sum_rate;
376
377 #ifdef RNG_HOST_API
378     real phi0;
379     real phi1;
380     real phi2;
381 #endif
382
383     // Simulation Parameters
384     int stride = dtau / dt; //
385     real sqrt_dt = sqrtf(dt);
386
387     int t = threadIdx.x % TIMEPOINTS;
388     int gindex = blockIdx.x * numberofPoints + threadIdx.x;
389
390     // Evolve the whole curve from 0 to T ( 1:1 mapping t with threadIdx.x)
391     if ((threadIdx.x < numberofPoints) && (gindex < gridDim.x * numberofPoints))
392     {
393         if (path == 0) {
394             rate = d_spot_rates[t];
395         }
396         else {
397             // Calculate dF term in Musiela Parametrization SDE
398             real dF = 0;
399
400             if (t == (TIMEPOINTS - 1)) {
401                 dF = simulated_rates[gindex] - simulated_rates[gindex - 1];
402             }
403             else {
404                 dF = simulated_rates[gindex + 1] - simulated_rates[gindex];
405             }
406
407             // Normal random variates
408 #ifdef RNG_HOST_API
409             real* rngNrms = (real*)rngNrmVar;
410             int rndIdx = blockIdx.x * pathN * VOL_DIM + path * VOL_DIM;
411             phi0 = rngNrms[rndIdx];
412             phi1 = rngNrms[rndIdx + 1];
413             phi2 = rngNrms[rndIdx + 2];
414 #else
415             if (threadIdx.x == 0) {
416                 curandStateMRG32k3a* state = (curandStateMRG32k3a*)rngNrmVar;
417                 curandStateMRG32k3a localState = state[blockIdx.x];
418                 phi0 = curand_uniform(&localState);
419                 phi1 = curand_uniform(&localState);
420                 phi2 = curand_uniform(&localState);

```

```

421     state[blockIdx.x] = localState;
422 }
423 __syncthreads();
424 #endif
425
426 // simulate the sde
427 rate = __musiela_sde2(
428     d_drifts[t],
429     d_volatilities[t],
430     d_volatilities[TIMEPOINTS + t],
431     d_volatilities[TIMEPOINTS * 2 + t],
432     phi0,
433     phi1,
434     phi2,
435     sqrt_dt,
436     dF,
437     simulated_rates[gindex],
438     dtau,
439     dt
440 );
441
442
443 #ifdef HJM_PATH_SIMULATION_DEBUG
444     printf("Path %d Block %d Thread %d index %d Forward Rate %f phi0 %f phi1 %f phi2 %f \n",
445         path, blockIdx.x, threadIdx.x, gindex, rate, phi0, phi1, phi2);
446 #endif
447
448 // accumulate rate for discount calculation
449 sum_rate = accum_rates[gindex];
450 sum_rate += rate;
451 accum_rates[gindex] = sum_rate;
452
453 // store the simulated rate
454 simulated_rates0[gindex] = rate; //
455
456 // update numeraire based on simulation block
457 if (path % stride == 0) {
458     if (t == (path / stride)) {
459         numeraires[gindex].x = rate;
460         numeraires[gindex].y = __expf(-sum_rate * dt);
461 #ifdef HJM_NUMERAIRE_DEBUG
462         printf("Path %d Block %d Thread %d index %d Forward Rate %f Discount %f\n",
463             path, blockIdx.x, threadIdx.x, gindex, rate, __expf(-sum_rate * dt));
464 #endif
465     }
466 }
467
468 */
469 /**
470 * Shared Memory & Global Access Memory optimizations & block simulation
471 */

```

# Appendix H

## Simulation kernel implementation with shared memory

```

473 template <typename real, typename real2>
474 __global__
475 void __generatePaths_kernel4(
476     int numberofPoints,
477     real2* numeraires,
478     real* rngNrmVar,
479     real* simulated_rates,
480     real* simulated_rates0,
481     real* accum_rates,
482     const int pathN,
483     int path,
484     real dtau = 0.5, real dt = 0.01)
485 {
486     // calculated rate
487     real rate;
488     real sum_rate = 0;
489     real phi0;
490     real phi1;
491     real phi2;
492
493     __shared__ real _ssimulated_rates[BLOCK_SIZE];
494
495     // Simulation Parameters
496     int stride = dtau / dt; //
497     real sqrt_dt = sqrtf(dt);
498
499     //int t = threadIdx.x % TIMEPOINTS;
500     int t = threadIdx.x;
501     int gindex = blockIdx.x * numberofPoints + threadIdx.x;
502
503     // load the accumulated rate for a given timepoint
504     sum_rate = accum_rates[gindex];
505
506     // load the latest simulated rate from global memory
507     if ((threadIdx.x < numberofPoints) && (gindex < gridDim.x * numberofPoints)) {
508         if (path > 0) {
509             _ssimulated_rates[threadIdx.x] = simulated_rates[gindex];
510         }
511     }
512     __syncthreads();
513
514     //
515     for (int s = 0; s < stride; s++) {
516
517         if ((threadIdx.x < numberofPoints) && (gindex < gridDim.x * numberofPoints))
518         {
519             if (path == 0) {
520                 rate = d_spot_rates[t];
521             }
522             else {
523                 // Calculate dF term in Musiela Parametrization SDE
524                 real dF = 0;
525
526                 if (t == (TIMEPOINTS - 1)) {
527                     dF = _ssimulated_rates[threadIdx.x] - _ssimulated_rates[threadIdx.x - 1];
528                 }
529                 else {
530                     dF = _ssimulated_rates[threadIdx.x + 1] - _ssimulated_rates[threadIdx.x];
531                 }
532
533                 // Normal random variates broadcast if access same memory location in shared memory

```

```

534     real* rngNrms = (real*)rngNrmVar;
535     int rndIdx = blockIdx.x * pathN * VOL_DIM + (path + s)* VOL_DIM;
536     phi0 = rngNrms[rndIdx];
537     phi1 = rngNrms[rndIdx + 1];
538     phi2 = rngNrms[rndIdx + 2];
539
540     // simulate the sde
541     rate = __musiela_sde2(
542         d_drifts[t],
543         d_volatilities,
544         d_volatilities[TIMEPOINTS + t],
545         d_volatilities[TIMEPOINTS * 2 + t],
546         phi0,
547         phi1,
548         phi2,
549         sqrt_dt,
550         dF,
551         _ssimulated_rates[threadIdx.x],
552         dtau,
553         dt
554     );
555 }
556
557     // accumulate rate for discount calculation
558     sum_rate += rate;
559
560 }
561
562 __syncthreads();
563
564 if ((threadIdx.x < numberOfPoints) && (gindex < gridDim.x * numberOfPoints))
565 {
566     _ssimulated_rates[threadIdx.x] = rate;
567 }
568
569 __syncthreads();
570 }
571
572 // update the rates and the rate summation for the next simulation block
573 if ((threadIdx.x < numberOfPoints) && (gindex < gridDim.x * numberOfPoints))
574 {
575     simulated_rates0[gindex] = rate;
576     accum_rates[gindex] = sum_rate;
577 }
578
579 // update numeraire based on simulation block
580 if ( t == (path + stride) / stride ) {
581     numeraires[gindex].x = rate; // forward rate
582 #ifdef double_ACC
583     numeraires[gindex].y = expf(-sum_rate * dt);
584 #else
585     numeraires[gindex].y = exp(-sum_rate * dt);
586 #endif
587 }
588 }
589 }
```

# Appendix I

## Multi-gpu implementation

```
1078 #ifdef MULTI_GPU_SIMULATION
1079 #pragma omp parallel num_threads(num_gpus)
1080 {
1081     int gpuDevice = omp_get_thread_num();
1082 #endif
1083     cudaSetDevice(gpuDevice);
1084
1085     // Kernel Execution Parameters
1086     int N = 1;
1087     int blockSize = N * BLOCK_SIZE;
1088     int numberofPoints = N * TIMEPOINTS;
1089     int gridSize = scenarios_gpus / N;
1090
1091     nvtxRangePush("Simulation");
1092
1093     TIMED_RT_CALL(
1094         riskFactorSim4(
1095             gridSize,
1096             blockSize,
1097             numberofPoints,
1098             d_numeraire[gpuDevice],
1099             rngNrmVar[gpuDevice],
1100             simulated_rates[gpuDevice],
1101             simulated_rates0[gpuDevice],
1102             accum_rates[gpuDevice],
1103             pathN,
1104             payOff.dtau,
1105             dt
1106         ),
1107         "Execution Time Partial HJM MC simulation"
1108     );
1109
1110     // TRACE main
1111     nvtxRangePop();
1112
1113     nvtxRangePush("Pricing");
1114     // Exposure Profile Calculation TODO (d_exposures + gpuDevice * TIMEPOINTS)
1115     // Apply Scan algorithm here
1116     TIMED_RT_CALL(
1117         exposureCalculation(scenarios_gpus, BLOCK_SIZE, d_exposures[gpuDevice], d_numeraire[gpuDevice], payOff.notional, payOff.K, scenarios_gpus),
1118         "exposure calculation"
1119     );
1120     nvtxRangePop();
1121
1122     //Replace all this block by a column reduction of the matrix
1123     // Partial Expected Exposure Calculation and scattered across gpus
1124     nvtxRangePushA("Exposure");
1125
1126     TIMED_RT_CALL(
1127         __expectedexposure_calc_kernel(partial_exposure[gpuDevice], d_exposures[gpuDevice], d_x[gpuDevice], d_y[gpuDevice], cublas_handle[gpuDevice], scenarios_gpus),
1128         "partial expected exposure profile"
1129     );
1130     nvtxRangePop();
1131
1132 #ifdef MULTI_GPU_SIMULATION
1133     }
1134 #endif
1135
```

# Appendix J

## System Configuration

### Software Platform

OS	Windows 10.0.1.19041 x86 64
NVIDIA SDK	NVIDIA SDK 10.0
Libraries	Intel MKL, cuBLAS, cuRAND

### CPU

Description	Intel® Xeon® W-2123 CPU @ 3.60GHz
Number of Cores	8
Memory	64GB

### GPU

Description	GeForce RTX 2070 SUPER
Number of GPU	4
GPU Chip	TU104
SM Count	40
Device Memory Size	8GB
Device Memory Theoretical B/W	427.31 GiB/s
Clock Rate	1.80GHz

### Visual Studio C++ Compiler Options

Host C/C++	/GS /GL /W3 /Zc:wchar_t /I"./" /I"C:\Program Files\NVIDIA GPU Computing Toolkit\CUDA\v10.2\include" /I"../common/inc" /I"C:\Program Files\NVIDIA GPU Computing Toolkit\CUDA\v10.2\include" /I"C:\Program Files (x86)\IntelSWTools\compilers_and_libraries\windows\mkl\include" /Zi /Gm- /O2 /Fd"x64\Release\vc142.pdb" /Zc:inline /fp:precise /D "WIN32" /D "_MBCS" /errorReport:prompt /WX- /Zc:forScope /Gd /MT /FC /Fa"x64\Release/" /EHsc /nologo /Fo"x64\Release/" /Fp"x64\Release\CVA_IRS.pch" /diagnostics:column
CUDA C/C++	C:\Program Files\NVIDIA GPU Computing Toolkit\CUDA\v10.2\bin\nvcc.exe" -gencode=arch=compute_75,code=\"sm_75,compute_75\" --use-local-env -ccbin "C:\Program Files (x86)\Microsoft Visual Studio\2019\Professional\VC\Tools\MSVC\14.27.29110\bin\HostX86\x64" -x cu -I./ -I"../common/inc" -I./ -I"C:\Program Files\NVIDIA GPU Computing Toolkit\CUDA\v10.2\include" -I"../common/inc" -I"C:\Program Files\NVIDIA GPU Computing Toolkit\CUDA\v10.2\include" -I"C:\Program Files (x86)\IntelSWTools\compilers_and_libraries\windows\mkl\include" -I"C:\Program Files\NVIDIA GPU Computing Toolkit\CUDA\v10.2\include" --keep-dir x64\Release -maxrregcount=0 --machine 64 --compile -cudart static -Xcompiler "/openmp /wd 4819" -lineinfo -lnvToolsExt -DWIN32 -DWIN32 -D_MBCS -D_MBCS -Xcompiler "/EHsc /W3 /nologo /O2 /Fdx64\Release\vc142.pdb /FS /Zi /MT" -o x64\Release\simulation_GPU.cu.obj "C:\ProgramData\NVIDIA Corporation\CUDA Samples\v10.2\7_CUDA\libraries\MC_SingleAsianOptionP\src\simulation_gpu.cu"

# **Appendix K**

## **Full Code Implementation**

The code implementation used for this work is publicly available at  
[https://github.com/mahsanchez/cva\\_hjm\\_cuda](https://github.com/mahsanchez/cva_hjm_cuda)