

# 场外衍生品定价中的蒙特卡洛模拟算法优化方法研究

羊隽婷\*      辛霄\*      陆雨<sup>†</sup>      蔡天舒<sup>†</sup>

2022 年 11 月 24 日

## 概述

蒙特卡洛方法 (Monte Carlo method), 也称统计模拟方法, 是指使用随机数 (或者伪随机数) 来解决复杂计算问题的方法, 在许多领域有着非常广泛的应用。

特别是在金融工程学中, 自 Myron Scholes 与 Fischer Black 提出 Black Scholes 模型, 以几何布朗运动 (geometric Brownian motion) 描述资产价格的随机过程, 蒙特卡洛方法就成为衍生品定价的基本方法。尽管近几十年来在 BS 模型的框架下, 学者与量化分析师们已经为数十种常见的复杂衍生品推导出了闭式解或近似解, 如各类美式障碍期权、美式香草期权、亚式期权等, 但随着场外衍生品市场的快速发展, 高度异质化、定制化的结构层出不穷, 以致目前大部分活跃交易的场外衍生品结构只能采用数值方法 (二叉树, 三叉树, 有限差分法等) 或蒙特卡洛方法进行相关计算。再者, 在 BS 模型的基础上, Heston model、CEV model、SABR model、Jump Diffusion model 等随机过程模型也被广泛应用于金融工程领域中。这些模型引入了更多的风险因子, 能更精细地刻画资产价格的随机过程, 但同时也为衍生品的定价带来极大的挑战。对于某些问题, 采用蒙特卡洛算法模拟随机过程进行求解几乎是唯一可行的方法。

近年来, 随着国内场外衍生品市场的飞速发展, 业务量和业务复杂度也相应地极速提升, 对相关从业机构与从业人员都提出了更高的要求。当前,

---

\*安信证券股份有限公司

<sup>†</sup>上海镓链科技有限公司

蒙特卡洛方法作为金融工程领域最为基础与通用的算法，虽然得到普遍运用，但运算效率上差强人意，难以满足业务、交易以及风险管理的要求。

因此，以满足国内场外衍生品市场潜在发展为目标，通过完善随机过程模型、优化设计框架，结合计算机软硬件相关先进算法，打造高效通用的蒙特卡洛模拟定价引擎，成为能否在未来场外衍生品市场竞争中处于优势地位的关键一步。

本文以场外衍生品定价中的蒙特卡洛模拟算法优化为目标，在镜像法的基础上提出了基于旋转变换的快速高斯随机数算法 (FASTNORM)，并运用单一指令多数据流 (SIMD) CPU 向量化计算技术以及基于 OpenMP 的跨平台多线程计算技术，将蒙特卡洛模拟方法的定价效率提升超 20 倍。最后，本文也研究了基于 GPU 的蒙特卡洛模拟算法，利用 CUDA 开发套件 (CUDA toolkit) 充分释放 GPU 的并行数据处理能力，实现了对蒙特卡洛模拟算法定价效率的跨越式优化。

**关键词：**蒙特卡洛模拟, FASTNORM, SIMD, OpenMP, GPU, CUDA

## 1 几何布朗运动及其离散化

### 1.1 几何布朗运动

几何布朗运动 (Geometric Brownian Motion, GBM) 是一个简单连续随机过程，对随时间演变的随机行为建模，该模型在物理和金融领域中被广泛应用。通常情况下，资产价格  $S(t)$  的时间演化模型能够通过几何布朗运动这一数学模型进行模拟，并由以下随机微分方程 (Stochastic Differential Equation, SDE) 表示：

$$dS(t) = \mu S(t)dt + \sigma S(t)dW_t \quad (1)$$

其中  $W_t$  是一个维纳过程 (Wiener process)。 $\mu$  表示漂移， $\sigma$  表示随机波动的幅度 (即波动率)，在 Black-Scholes-Merton 的模型假设下，它们都是常数。

### 1.2 随机过程的离散化

根据伊藤引理 (Itô's lemma)

$$d \ln S(t) = (\mu - \frac{1}{2}\sigma^2)dt + \sigma dW_t \quad (2)$$

46 将 (2) 式两边从  $t$  到  $t + \Delta t$  积分, 可得

$$S(t + \Delta t) = S(t) \exp((\mu - \frac{1}{2}\sigma^2)\Delta t + \sigma W_{\Delta t}) \quad (3)$$

47 其中,  $W_{\Delta t} \sim \mathcal{N}(0, \sqrt{\Delta t})$ 。

48 因此, 资产价格从  $t$  到  $t + \Delta t$  的离散化表达是

$$S(t + \Delta t) = S(t) \exp((\mu - \frac{1}{2}\sigma^2)\Delta t + \sigma\sqrt{\Delta t}\epsilon) \quad (4)$$

49 或

$$\ln S(t + \Delta t) = \ln S(t) + (\mu - \frac{1}{2}\sigma^2)\Delta t + \sigma\sqrt{\Delta t}\epsilon \quad (5)$$

50 其中,  $\epsilon \sim \mathcal{N}(0, 1)$ 。

51 上述即为 Black-Scholes 模型下的欧拉离散化 (Euler discretization),  
 52 其中 (5) 式为其对数形式的表达方法。采用对数形式计算可以避免复杂的  
 53 幂运算, 仅需进行简单的加法运算, 因此可以有效地节省计算机资源从而提  
 54 升计算效率。在后续的蒙特卡洛计算中, 均模拟资产的对数价格路径进行期  
 55 权现值的计算。

## 2 一般算法分析

56  
 57 蒙特卡洛模拟是期权定价的重要工具, 采用该方法能够灵活地为期权  
 58 进行定价。除了一些简单的期权, 例如普通欧式香草期权等, 目前市场上活  
 59 跃的大多数复杂期权都没有简单的闭式解形式。由于蒙特卡洛方法在实际  
 60 应用中的广泛性, 该方法逐渐成为目前为复杂金融期权进行定价的流行计  
 61 算工具。然而, 在大多数情况下, 使用蒙特卡洛模拟的计算效率并不高效。  
 62 为了能够同时满足使用该方法为期权定价的高精度以及高效率的要求, 对  
 63 蒙特卡洛算法的计算各步骤展开研究。

## 2.1 蒙特卡洛模拟在期权定价中的具体运用

使用蒙特卡洛模拟对期权进行定价需要生成资产价格的随机路径，并且每个路径都需有相应的赔付，将这些收益折现取平均值即可得到期权的价值。在已确定的衍生品结构之下，定义生成一次蒙特卡洛模拟路径所需的步骤如下：

- 生成随机数：根据模型假设确定对应的分布类型，例如在 Black-Scholes 模型假设下对数资产价格的随机数服从正态分布；
- 根据上述的随机数生成路径：由于 Black-Scholes 随机方程有解析解，因此只需要模拟特定时间点的标的价格。例如，对于离散观察日的障碍期权，只需要模拟每个观察日的标的价格即可；然而，对于欧式结构而言，仅需模拟到期日的价格即可；
- 根据所生成的标的对数价格路径、贴现因子以及期权定义，计算期权现值。

## 2.2 算法时效性分析

为了分析蒙特卡洛模拟在期权定价中的具体运用，使用第十二代英特尔酷睿 i7 芯片并在 Ubuntu 环境下进行算法时效性的测试研究，具体的机器配置详见附录 1。采用 C++ boost 库中的 MT19937 随机数生成器生成对应的正态分布随机数，期权的测试结构为离散观察的向上敲出看涨期权，图 1 和图 2 分别是该期权的具体测试结构条款以及相应的定价参数。

根据以上设置生成 100 万条蒙特卡洛模拟路径并计算期权现值，分析各算法步骤的具体耗时情况。通过初步实验结果发现，随机数生成这一步骤占用了绝大部分的计算资源。因此，在下文中将针对生成随机数这一步骤进行深入研究，进而在最大程度上实现蒙特卡洛模拟算法优化这一目标。各算法步骤的占比时间如下图：

## 3 关于正态分布随机数生成的优化方法

如前所述，使用伪随机数生成器生成随机数往往需要大量的时间。由正态随机变量的性质，可以对伪随机数生成器产生的随机数应用更快速的线

期权结构名称	
现货价格	90
障碍价格	100
障碍方向	向上敲出
观察类型	按日离散观察
到期时间	1 年
年交易天数	244 日
到期支付类型	香草看涨期权
到期行权价	80

表 1: 测试结构条款

无风险利率	分红率	波动率
1.00%	3.00%	20.00%

表 2: 期权定价参数

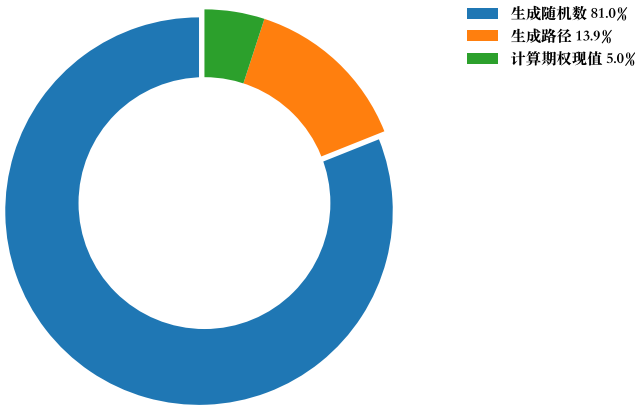


图 1: 蒙特卡洛模拟各算法步骤的计算时间占比，总耗时为 1398 毫秒。蓝色部分为随机数生成，耗时为 1133 毫秒，占总耗时约 81.0%。橙色部分为生成随机路径，耗时为 195 毫秒，占总耗时的 13.9%。绿色部分为计算期权现值，耗时为 70 毫秒，占总耗时的 5.0%。

性变换，得到更多的随机数，从而减少生成随机数的时间，提升蒙特卡洛模拟的运算速度。

基于这一点，本节将介绍并使用镜像法和 FASTNORM 【加引用】方法提升生成随机数的速度。

### 3.1 镜像法

在生成符合正态分布的随机数时，伪随机数生成器首先依均匀分布产生一个随机数，随后根据这个均匀分布的随机数来计算相应的正态随机数，这个过程往往比较耗时。如果能够根据一组已经产生的正态随机数生成更多的正态随机数，那么就可减少伪随机数生成器的使用，从而提升蒙特卡洛模拟的运算速度。针对这一点，一种简单可行的方法是将原本随机数的相反数作为另一个随机数使用。

由标准正态随机变量的性质，如果  $X \sim N(0, 1)$ ，那么  $-X \sim N(0, 1)$ 。因此，生成随机数时在每生成一个随机数  $x$  之后，取  $-x$  作为另一个随机数， $-x$  也是根据正态分布产生的随机数。这样就省去了一半的随机数生成时间。这种方法的另一个好处是生成的随机数样本均值为 0，这正是  $X$  的均值。

利用正态随机变量的性质，可以将镜像法扩展到更一般 FASTNORM 方法，通过线性变换和镜像法来成倍地产生随机数，从而减少运算的时间。

### 3.2 FASTNORM

如果  $\mathbf{X}$  是一个由  $n$  个独立标准正态随机变量组成的向量， $\mathbf{R}$  是一个  $m \times n$  矩阵，且  $\mathbf{R}\mathbf{R}' = \mathbf{I}$ ，那么  $\mathbf{Y} = \mathbf{R}\mathbf{X}$  是一个由  $m$  个独立标准正态随机变量组成的向量，这是因为随机向量  $\mathbf{Y}$  的期望和协方差矩阵分别为

$$\begin{aligned} \mathbf{E}(\mathbf{Y}) &= \mathbf{E}(\mathbf{R}\mathbf{X}) = \mathbf{R}\mathbf{E}(\mathbf{X}) = \mathbf{R}\mathbf{0} = \mathbf{0}, \\ \mathbf{V}(\mathbf{Y}) &= \mathbf{R}\mathbf{V}(\mathbf{X})\mathbf{R}' = \mathbf{R}\mathbf{R}' = \mathbf{I}. \end{aligned} \tag{6}$$

注意，如果  $\mathbf{R}\mathbf{R}' = \mathbf{I}$ ，那么  $m \leq n$ 。

令  $\mathbf{R} = -\mathbf{I}$ ，可以看出镜像法实际上是 FASTNORM 的特例。

二维旋转矩阵

$$\begin{bmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{bmatrix} \tag{7}$$

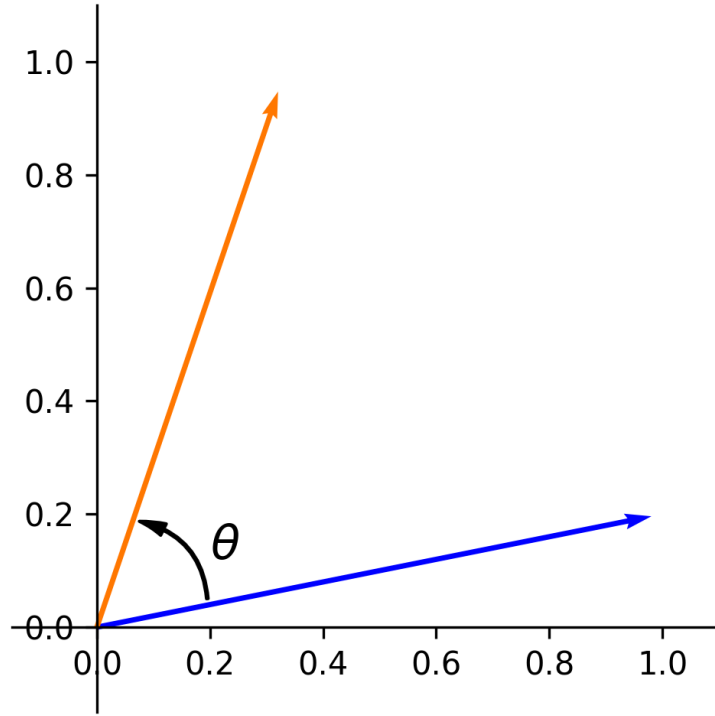


图 2: 向量旋转: 蓝色向量逆时针旋转得到橙色向量

116 就是一类符合以上条件的  $\mathbf{R}$ 。二维空间中向量的左乘旋转矩阵就能得  
117 到逆时针旋转  $\theta$  之后的向量，如下图所示：

118 应特别注意，在二维情形中， $\theta$  的取值应满足  $\theta \neq \pi + k\pi, \theta \neq \frac{\pi}{2} + k\pi, k =$   
119  $0, \pm 1, \pm 2, \dots$ ，以避免与原向量或镜像后的向量重复。

120 使用这种方法，能由  $n$  个随机数  $\mathbf{X}$  立即得到  $m$  个新的随机数  $\mathbf{Y}$ 。在  
121 此基础上再应用镜像法，能够得到  $-\mathbf{X}$  和  $-\mathbf{Y}$ 。因此，只需要生成  $m$  个伪  
122 随机数，就能得到  $2(m + n)$  个随机数。

123 尽管 FASTNORM 能大幅提升随机数的产生速度，但这种方法的在一  
124 定程度上牺牲了随机数的质量，这是由于  $\mathbf{X}$ 、 $-\mathbf{X}$ 、 $\mathbf{Y}$ ，以及  $-\mathbf{Y}$  并不是相  
125 互独立的。有

$$\begin{aligned}\text{Cov}(X_j, Y_i) &= r_{ij}, \\ \text{Cov}(X_j, -X_j) &= -1, \quad 1 \leq j \leq n, 1 \leq i \leq m.\end{aligned}\tag{8}$$

126 从第一个式子中也可以看出, 选取较大的  $m$  和  $n$  能够降低新随机数与  
127 原本随机数之间的相关性, 从而降低 FASTNORM 在相关性上产生的影响。  
128 下一小节将进一步讨论镜像法和 FASTNORM 在蒙特卡洛算法中的实  
129 现。

### 130 3.3 算法实现与分析

131 使用 FASTNORM 时, 矩阵  $\mathbf{R}$  的选取是相当重要的。理想的线性变换  
132 矩阵  $\mathbf{R}$  应满足以下条件: 1)  $\mathbf{R}\mathbf{R}' = \mathbf{I}$  (这意味着  $m \leq n$ ); 2)  $\mathbf{R}$  中的所有  
133 元素都不能为 1 或  $-1$ , 以确保生成的新随机数及其相反数不与原随机数重  
134 复。基于这两点, 可以选取矩阵

$$\mathbf{R} = \begin{bmatrix} 1/2 & 1/2 & -1/2 & 1/2 \\ 1/2 & -1/2 & 1/2 & 1/2 \\ 1/2 & -1/2 & -1/2 & -1/2 \\ -1/2 & -1/2 & -1/2 & 1/2 \end{bmatrix}.\tag{9}$$

135 这也是 Wallace 推荐的方法。容易验证  $\mathbf{R}$  是规范正交的。因此由线性  
136 变换产生的每组随机数是独立的正态变量。

137 确定了  $\mathbf{R}$  之后, FASTNORM 的实现大致可分为三步: 1) 用伪随机数  
138 生成器产生一组随机数 (以下将称为原随机数); 2) 依镜像法, 取原随机数  
139 的相反数, 得到另一组随机数; 3) 对原随机数进行线性变换, 并取线性变换  
140 的相反数, 得到另外两组随机数。具体实现方式如下:

141 假设每条路径的生成需要  $N$  个随机数。用伪随机数生成器生成 4 个随  
142 机数, 记为  $\mathbf{x}_1$ , 然后计算  $\mathbf{y}_1 = \mathbf{R}\mathbf{x}_1$ ; 然后用伪随机数再次生成 4 个随机数  
143  $\mathbf{x}_2$ , 并计算得到  $\mathbf{y}_2 \dots$  直到  $\mathbf{x}_N$  和  $\mathbf{y}_N$ 。然后, 用  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N$  的第  $i$  个元  
144 素生成一条标的价格的路径, 这里  $i = 1, 2, 3, 4$ 。对  $\mathbf{y}_1, \dots, \mathbf{y}_N$ 、 $-\mathbf{x}_1, \dots, -\mathbf{x}_N$ 、  
145  $-\mathbf{y}_1, \dots, -\mathbf{y}_N$  也用类似方式生成路径。

146 这种方法每循环一次就能得到总共 16 条路径, 而伪随机数生成器只生  
147 成了其中 4 条路径所需要的随机数, 因此这种方法能够提升随机数产生的  
148 速度。



方法名称	
原始方法	1.358816
镜像法	1.358905
FASTNORM	1.358890

表 3: 期权现值比较

	<i>t</i> 统计量	<i>P</i> 值
原始方法与镜像法	0.36897580272371977	0.7121848363482105
原始方法与 FASTNORM	0.7378756158538999	0.4606766455176965

表 4: *t* 检验：原始方法、镜像法，以及 FASTNORM 方法得到的估计量的均值及 *t* 检验结果。由 2000 次蒙特卡洛模拟得到

149        下一节将展示镜像法和 FASTNORM 带来的性能优化，以及讨论随机  
150 数的非独立性对蒙特卡洛估计的影响。

151    **3.4 分析对比**

152    **3.4.1 运算速度**

153        下图对比了运行一次原始方法、镜像法和 FASTNORM 方法的所需的  
154 时间，模拟次数为 100 万次。如图所示，运行一次原始蒙特卡洛模拟算法需  
155 要 1403 毫秒。使用镜像法之后耗时缩短至 820 毫秒，是原始算法的 58.4%；  
156 而使用 FASTNORM 之后，耗时仅为 418 毫秒，是原始算法的 29.8%。由  
157 此可见，FASTNORM 能将蒙特卡洛模拟的运算速度提升超过三倍。

158    **3.4.2 估计量的期望**

159        前文提到镜像法和 FASTNORM 方法会导致随机数之间的相关性，可  
160 能会对蒙特卡洛的估计量带来影响。镜像法和 FASTNORM 方法使用的随  
161 机数（不论是由伪随机数算法产生的，还是通过线性组合得到的）都是符合  
162 标准正态分布的，因此它们的蒙特卡洛估计量仍然是无偏的，如下表所示  
163 （这里使用了 2000 次模拟的结果）：

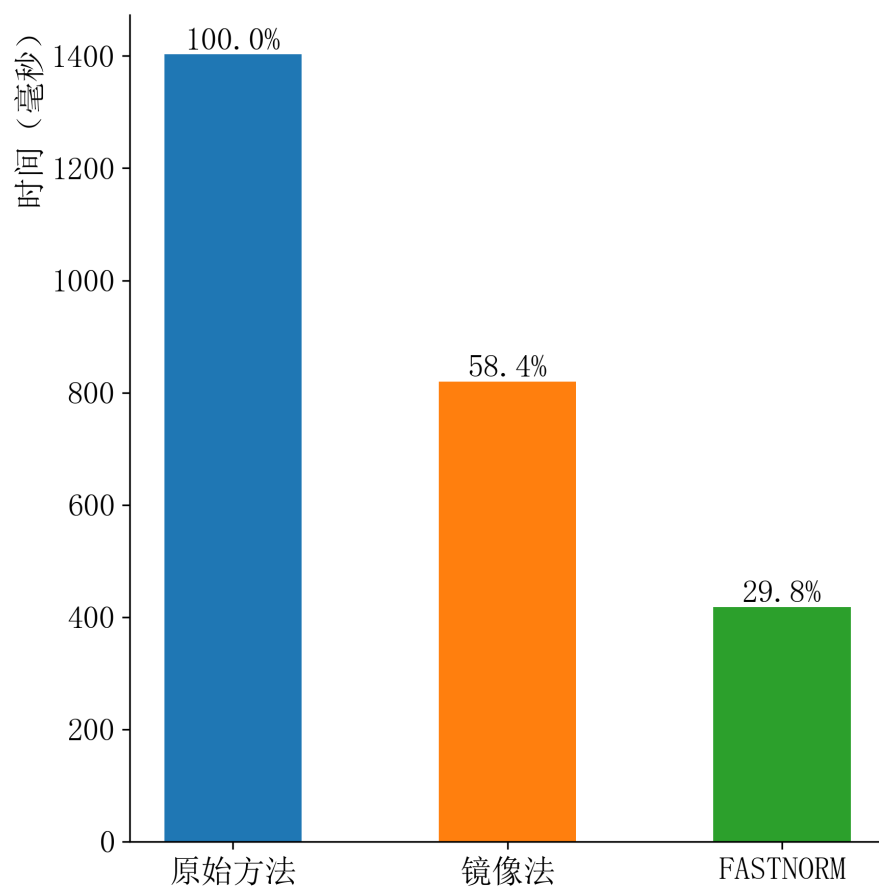


图 3: 运行 100 万次蒙特卡洛模拟的统计时间。蓝色部分为原始蒙特卡洛模拟算法，橙色部分为镜像法，绿色部分为 FASTNORM 方法

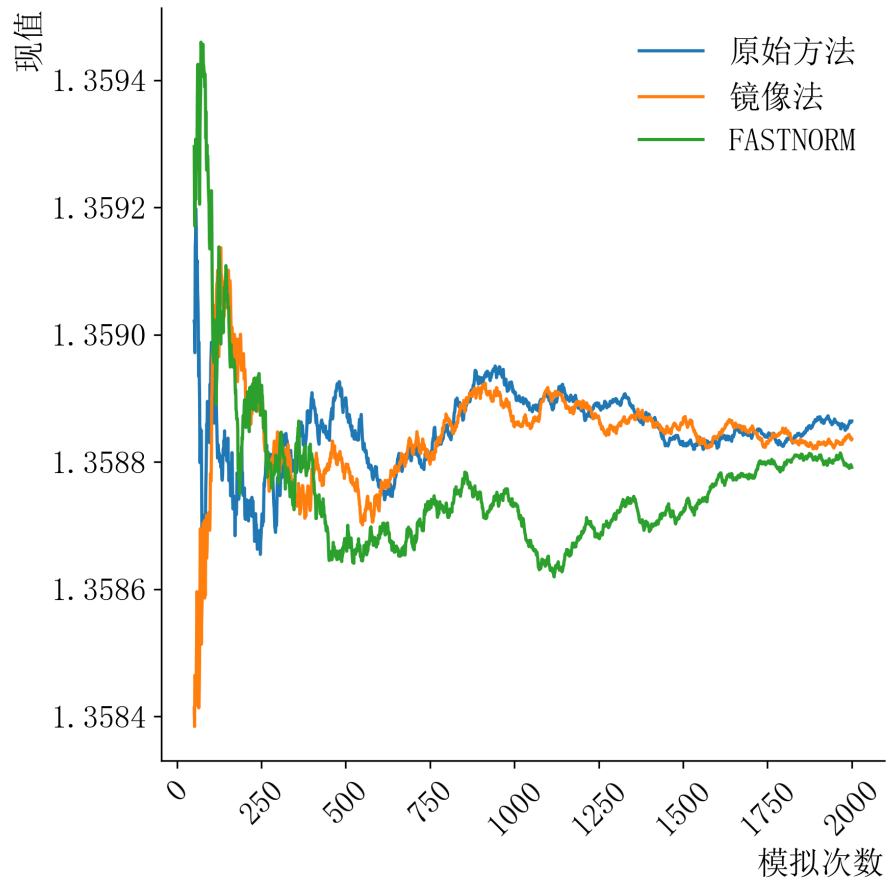


图 4: 需要修改图片，注明横轴的单位是百万次

164 上图展示了三种方法的估计随模拟路径数增加的变化。可以看出在当  
 165 路径数量增大时，三种方法的结果都向 1.3588 附近靠近。

166 估计量的无偏性确保了镜像法和 FASTNORM 的可行性。然而如前文  
 167 所述，FASTNORM 方法对蒙特卡洛估计量的影响主要体现在有效性上，即  
 168 估计量的标准差（或方差）上，这种影响可以通过比较估计量的标准差来讨  
 169 论。

方法名称	
原始方法	0.00350160089209261
镜像法	0.0034882495542142957
FASTNORM	0.003828353638954014

表 5: 各方法的标准差比较

	<i>F</i> 统计量	<i>P</i> 值
原始方法与镜像法	1.0071655989622432	0.8731991106943031
原始方法与 FASTNORM	0.836164800985319	6.419816278691966e-05

表 6: *F* 检验: 原始方法、镜像法, 以及 FASTNORM 方法得到的估计量的标准差及 *F* 检验结果, 模拟的路径数均为 100 万条.

170 **3.4.3 估计量的标准差**

171       下表展示了原始方法、镜像法, 以及 FASTNORM 方法的标准差, 用  
172       以比较估计量的有效性。由中心极限定理, 原始方法的标准差可以通过计算  
173       100 万条路径对应的贴现收益的标准差得到。而在镜像法和 FASTNORM  
174       方法中, 随机数并不是独立的, 因而贴现收益不一定独立, 所以中心极限定  
175       理关于方差的部分并不适用。出于这个原因, 镜像法和 FASTNORM 方法  
176       的估计量的标准差通过运行 2000 次模拟并计算标准差得到。

177       虽然中心极限定理关于方差的部分不适用于镜像法和 FASTNORM 方  
178       法产生的估计量, 但和原始方法一样, 这两种方法产生的估计量仍然依分  
179       布收敛于正态变量, 因此这里可以使用 *F* 检验来检验原始方法与镜像法或  
180       FASTNORM 方法的方差 (或标准差) 的差异。可以看出, 镜像法和原始方  
181       法的标准差没有显著差异, 然而对比原始方法, FASTNORM 的标准差上升  
182       了 9.3%。因此, FASTNORM 下 100 万条路径的模拟精度大约相当于原始  
183       方法下 80 万条路径的模拟精度, 然而前者的耗时却不到后者的 1/4, 对比  
184       镜像法, FASTNORM 也只需要一半的时间。该结果表明 FASTNORM 是  
185       比原始方法更高效的算法。

## 4 基于硬件平台的进一步优化

本文不仅从算法层面上对随机数生成进行了优化，也考量了现代 CPU (central processing unit, 中央处理器) 在提升蒙特卡洛模拟计算效率上的可能性。随着 CPU 的更新迭代，现代 CPU 都能够支持多核多线程，并且大多数 CPU 支持单指令多数据流 (Single Instruction Multiple Data, SIMD) 的向量化 (vectorization) 运算以及多线程并行运算。基于此基础上，下文将从 CPU 层面上对提升随机数生成速率进行相关研究。

### 4.1 基于 SIMD 的 CPU 向量化

向量化是提高现代 CPU 性能的关键工具，而 SIMD 可以实现 CPU 的向量化运算。SIMD 认为计算机具有多个处理元素，并且这些处理元素在多个数据点上能够同时执行相同的运算。向量化能够实现从一次操作多个值转换为一次操作一组值，并且现代 CPU 能够直接支持 SIMD 向量运算。例如，一个 512 位寄存器的 CPU 可以容纳 16 个 32 位单精度浮点数并进行单次运算，比单次执行一条指令快 16 倍。因此，将向量化与多核多线程 CPU 相结合，可以带来数量级的性能提升。

通常情况下，SIMD 单元中的操作包括基本的数学运算（比如加、减、乘、除）及其他常见的数学运算，包括绝对值 (abs) 和平方根 (sqrt) 的运算。比如，SIMD 接收两个向量作为输入（每个向量都有一组运算对象），对这两组运算对象（每个向量的每个操作数）执行相同的运算，并输出一个带结果的向量，如图 4 所示。

在现代 CPU 的发展中，指令集架构 (Instruction Set Architecture, ISA) 也在不断地更新，实现了更高效地进行与 SIMD 相关的数据处理的可能性。指令集架构是能够同时多个数据对象上执行相同操作并且提升性能的额外指令，其中包括 SIMD、流 SIMD 拓展 (Streaming SIMD Extensions, SSE)、高级矢量拓展 (Advanced Vector Extensions, AVX)。SSE 是对现有的 32 位 (x86) CPU 架构的 SIMD 指令集扩展，用于取代最初的、性能较差的 MMX (MultiMedia EXtensions) 架构。随着技术的更新，AVX 的出现提供了宽度更高的向量、新的拓展性语法以及更加丰富的功能，并且在 SSE 基础上实现更高效的数据管理功能以及性能提升。

目前，最新的处理器都具有多核，在日益增加的大型数据集上使用 SIMD 执行单指令，能够在提升性能的同时降低能耗。

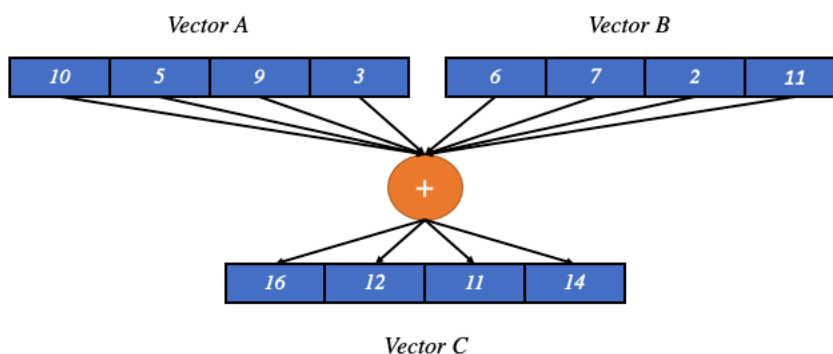


图 5: SIMD 向量加法示例。将输入的两个向量分别命名为 *Vector A* 和 *Vector B*，这两个向量均包含四个操作数，对这两组数并行执行四个相同的加法运算，得到结果向量 *Vector C*。

#### 217 4.1.1 SIMD 在蒙特卡洛模拟中的具体应用

218 在蒙特卡洛模拟中，使用基于 CPU 的 SIMD 向量化计算技术能够提升  
219 矩阵运算和随机数生成的速度。例如，FASTNORM 中矩阵的旋转变换、以  
220 及利用向量化进行快速随机数生成等。

##### 221 (1) 矩阵运算

222 在上文中提及的 FASTNORM 算法中所运用的矩阵旋转变换可以通过  
223 基于 SIMD 优化的矩阵运算实现，进而提升计算速率。

224 在实际运算中，使用 SIMD 实现常见的线性代数底层程序，即 BLAS。  
225 BLAS 规范了常见的线性代数运算的内核底层程序，是线性代数库的标准底  
226 层程序。BLAS 功能分为三个被称为级别 (levels) 的程序，分别对应于相  
227 关定义和发布时间的顺序，以及多项式次数对应的算法复杂度。第一级别的  
228 BLAS 操作通常耗时为线性的，第二级别操作为二次方时间，第三级别操作  
229 为立方时间。现代化的 BLAS 实现通常提供所有的级别。由于矩阵乘法在  
230 大部分数学应用中至关重要，使用 BLAS 能够高效、便捷、广泛地应用在  
231 数学运算中。

232 尽管 BLAS 是通用的规范，但是 BLAS 的实现通常针对特殊的机器进  
233 行优化。例如，Intel oneMKL 的 BLAS 针对 Intel 体系结构进行了优化，因  
234 此使用该数学库的 BLAS 接口进行矩阵运算能够在 Intel 芯片上实现优异

235 的性能提升。

## 236 (2) 随机数生成

237 为了充分利用现代 CPU 对 SIMD 向量化计算的优势，随机数生成器能  
238 够通过调用高度优化的基本随机数生成器和向量数学函数进行开发，从而  
239 实现大幅提升随机数生成速率的目标。例如，面向 SIMD 的快速梅森旋转  
240 (SIMD-oriented Fast Mersenne Twister, SFMT) 是一个线性反馈移位寄存  
241 器 (Linear Feedbacked Shift Register, LFSR) 生成器，可以一次生成 128  
242 位伪随机整数，其随机数生成速度是使用普通梅森生成速度的两倍。

## 243 4.2 CPU 多线程

244 除了基于 SIMD 的向量化运算之外，利用现代 CPU 的多核多线程也能  
245 够大幅提升计算效率。目前，几乎所有的现代 CPU 都支持多线程。不仅如  
246 此，多核 CPU 上的多线程能够一次处理不同的任务，从而使得计算机性能  
247 更加高效、耗能更低。

248 线程 (Thread) 是并发编程 (concurrent programming) 中的一个执行  
249 单元，而多线程则允许 CPU 同时执行一个进程中的多个任务。在计算机操  
250 作系统中，并发编程指的是在重叠时间段内执行多个任务，并且没有特定的  
251 执行顺序，而线程编程 (threaded programming) 则最常用于共享内存并行  
252 (shared memory parallel) 计算机中。线程与进程 (processes) 类似，不同  
253 之处在于线程之间可以彼此共享内存 (以及拥有私有内存)。在实际操作中，  
254 线程之间必须能够交换数据以便进行有用的并程序，也可以通过读写共  
255 享数据进行通讯。例如，线程 1 向共享变量 a 写入一个值，线程 2 可以从  
256 共享变量 a 中读取这个值。

257 多线程编程 (Multithreaded programming) 是指对多个并发执行的线  
258 程进行编程从而实现高性能的机制，在单核 CPU 或者多核 CPU 上均能够  
259 运行多线程。值得注意的是，单核 CPU 上的多线程并不是并行运行的。因  
260 为在实际运用中，单核 CPU 是通过调度算法 (scheduling algorithm)，或  
261 者根据外部输入 (或干扰) 的组合以及线程间的优先级进行切换。然而，多  
262 核 CPU 的多线程则是并行的，因为多核 CPU 是通过多个微处理器的共同  
263 协作从而提升性能。

264 并行化 (Parallelism) 是指在多核 CPU 等具有多种计算资源的硬件上  
265 同时运行多个任务或同一任务。当 CPU 内部运行的时钟频率 (clock speed)  
266 已经最大化时，并行化则是在多核 CPU 上提升性能的唯一办法。由于多核

267 CPU 的内部运行的时钟频率是单核 CPU 的多倍速，因此，采用多核 CPU  
268 不仅能在单核 CPU 的基础上有效地处理指令速度，还能够降低计算机的耗  
269 能。

#### 270 4.2.1 具体应用:OpenMP

271 CPU 与多线程之间有多种交互方式，为了计算的通用型及代码的可移  
272 植性，本文采用 OpenMP 简化多线程，进而提升蒙特卡洛模拟的计算速度。  
273 OpenMP (Open Multi-Processing) 是一个应用程序编程接口 (API)，它支  
274 持在各不同的平台、指令集架构和操作系统 (包括 Solaris、AIX、FreeBSD、  
275 HP-UX、Linux、macOS 和 Windows) 上针对 C, C++ 和 Fortran 的一组  
276 拓展，进行多平台共享内存的多线程编程。

277 OpenMP 通过指令 (directive) 和标记 (sentinel) 简化多线程。一个  
278 OpenMP 指令是指一行只针对某些编译器有意义的源代码，通常在行首由  
279 一个特殊标记进行区分。

280 OpenMP 使用并行执行的 Fork-join 模型从而实现并行化。其中，并行  
281 区域 (parallel region) 是 OpenMP 中的基本并行结构，用于定义程序的一  
282 部分。当程序开始时，在一个单独的线程 (主线程，即 master thread) 上  
283 开始执行。当遇到第一个并行区域时，主线程创建一个线程组 (fork-join 模  
284 型，如图 3 所示)。

##### 285 (1) 并行循环

286 循环是大部分应用程序中进行并行的主要来源，OpenMP 对并行循环  
287 (parallel loops) 提供了广泛的支持。当一个循环的迭代是独立的，即可以以  
288 任何顺序完成时，那么其能够在不同的线程中共享迭代，并认为这一个迭代  
289 或一组迭代是一个任务。在 OpenMP 中，仅需利用一行编译指令指示编译  
290 器生成代码，进而对多线程之间的循环迭代进行拆分，如 Algorithm 1 所示。

##### 291 (2) 降维

292 除了并行循环外，降维 (reduction) 可以实现并行执行一些重复计算的  
293 形式：从加、乘、最大值、最小值、和、或等关联操作中产生一个值。在降  
294 维中，每个线程可以积累自己的私有副本，然后这些副本被降维以得到最终  
295 结果。如果操作的数量远远大于线程的数量，那么大多数操作可以在并行中  
296 进行，如 Algorithm 2 所示。



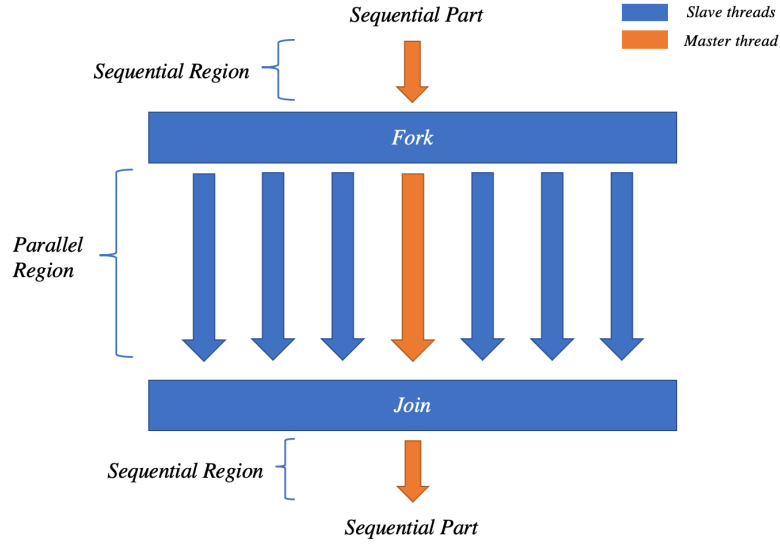


图 6: Fork-join 模型。每个线程执行并行区域内的命令，在并行区域结束后，主线程等待其他线程结束后，继续执行下一个命令。

---

**Algorithm 1** OpenMP: 并行循环

---

**Input:** Vectors  $\mathbf{A}^{(n)}$ ,  $\mathbf{B}^{(n)}$

**Output:** Element-wise addition of  $\mathbf{A}^{(n)}$ ,  $\mathbf{B}^{(n)}$

**Initialize:**

$\mathbf{A}^{(i)} \leftarrow i$

$\mathbf{B}^{(i)} \leftarrow i$

$i = 1, \dots, n$

**Parallel Loops:**

#pragma omp parallel for default(none)

shared(a, b)

**for**  $i = 1 : n$  **do**

$\mathbf{A}^{(i)} += \mathbf{B}^{(i)}$

**end for**

**return**  $\mathbf{A}^{(n)}$

---

---

**Algorithm 2** OpenMP: 降维

---

**Input:** Vector  $\mathbf{A}^{(n)}$ , double  $\mathbf{b}$

**Output:** Sum of the elements in  $\mathbf{A}^{(n)}$

**Initialize:**

$\mathbf{A}^{(i)} \leftarrow i$

$\mathbf{b} \leftarrow 0$

$i = 1, \dots, n$

**Reduction:**

#pragma omp parallel for default(none)

shared(a)

reduction(+: b)

**for**  $i = 1 : n$  **do**

$\mathbf{b} += \mathbf{A}^{(i)}$

**end for**

**return**  $\mathbf{b}$

---

297 **4.3 结果分析对比**

298 插入图片：插图，展示使用 SIMD 带来的性能提升]

299 **5 希腊值**

300 衍生品或其他资产的希腊值是该资产对不同风险因子的敏感程度，因此

301 能够描述该资产对不同风险因子的暴露程度。希腊值越大，则对应的风险因

302 子的暴露程度越大。比较常见的风险因子有 Delta、Gamma、Vega、Theta

303 和 Rho。这些希腊值分别定义为：

$$\begin{aligned} \text{Delta} &= \frac{\partial f}{\partial S}, \\ \text{Gamma} &= \frac{\partial^2 f}{\partial S^2}, \\ \text{Vega} &= \frac{\partial f}{\partial \sigma}, \\ \text{Theta} &= \frac{\partial f}{\partial \tau}, \\ \text{Rho} &= \frac{\partial f}{\partial r}. \end{aligned} \tag{10}$$

304 这里  $f$  表示衍生品价值,  $S$  表示标的价格,  $\sigma$  标的的隐含波动率,  $\tau$  表  
305 示衍生品剩余期限,  $r$  表示无风险利率。

306 由于风险因子会不断变化, 从而导致衍生品价值的变化, 因此在风险管  
307 理中, 希腊值起着非常重要的作用。希腊值不仅可以直观地描述资产的风  
308 险, 也可以帮助调整整个投资组合的风险暴露程度。在对冲中, 可以根据希  
309 腊值调整用于对冲的资产的持有数量, 从而使整个组合的对各风险因子的  
310 暴露达到理想的状态。

### 311 5.1 希腊值的计算方法

312 即使是在 Black-Scholes 模型下, 对于奇异衍生品来说, 希腊值很少有  
313 解析解。因此, 往往使用有限差分法来计算衍生品的希腊值。一阶导数的前  
314 向差分格式为

$$f'(x) \approx \frac{f(x + \Delta x) - f(x)}{\Delta x}. \quad (11)$$

315 这里  $\Delta x > 0$ 。可以分析这种格式的误差。如果  $f''$  存在, 由 Taylor 定  
316 理,

$$f(x + \Delta x) = f(x) + f'(x)\Delta x + \frac{1}{2}f''(t)(\Delta x)^2 \quad (12)$$

317 对某个  $t \in (x, x + \Delta x)$  成立。于是,

$$\left| \frac{f(x + \Delta x) - f(x)}{\Delta x} - f'(x) \right| \leq \frac{|f''(t)|}{2} \Delta x. \quad (13)$$

318 因此, 如果能够精确地 (即无误差地) 计算  $f(x + \Delta x)$  和  $f(x)$ , 则应  
319 取尽可能小的  $\Delta x$ 。

320 然而, 蒙特卡洛模拟和其他数值方法都会存在一定的误差。假定只能计  
321 算  $\hat{f}(x + \Delta x)$  和  $\hat{f}(x)$ , 且两者相对  $f(x + \Delta x)$  和  $f(x)$  的误差为  $\varepsilon > 0$ , 则  
322 用数值方法结果代替精确值进行前向有限差分的误差的一个上界为

$$\left| \frac{\hat{f}(x + \Delta x) - \hat{f}(x)}{\Delta x} - f'(x) \right| \leq \frac{2\varepsilon}{\Delta x} + \frac{|f''(t)|}{2} \Delta x. \quad (14)$$

323 如果知道了  $|f''|$  的上界, 则可以在最小化该误差上界的意义上确定一  
324 个最佳的  $\Delta x$ 。虽然对于奇异衍生品的希腊值来说, 前者往往是无法精确计

325 算的，但是该结果仍能够表明，并不是使用越小的  $\Delta x$ ，前向差分法计算的  
326 一阶导数就越精确。

327 二阶导数的离散格式采用中心差分法：

$$f''(x) \approx \frac{f(x + \Delta x) - 2f(x) + f(x - \Delta x)}{(\Delta x)^2}. \quad (15)$$

328 使用和一阶导数前向差分法类似的分析，有

$$\left| \frac{\hat{f}(x + \Delta x) - 2\hat{f}(x) + \hat{f}(x - \Delta x)}{(\Delta x)^2} - f''(x) \right| \leq \frac{4\varepsilon}{(\Delta x)^2} + \frac{|f^{(4)}(t)|}{6}(\Delta x)^2. \quad (16)$$

329 这里  $t$  是某个在  $(x - \Delta x, x + \Delta x)$  内的数。同样，在近似计算二阶导数  
330 时，更小的  $\Delta x$  也并不能保证更精确的结果。

## 331 5.2 结果展示

332 插入图片：希腊值相关展示图]

# 333 6 基于 GPU 的蒙特卡洛模拟算法

334 除了采用现代 CPU 的多核多线程优化蒙特卡洛模拟的计算速率外，下  
335 文将研究基于 GPU (graphics processing unit, 图形处理器) 的蒙特卡洛模  
336 拟算法，从而进一步提升蒙特卡洛模拟的适用性和实践价值。

## 337 6.1 GPU 计算的优势

338 无论是个人计算还是商业计算，GPU 已经成为最重要的计算技术之一。  
339 GPU 专为并行处理而设计，其高度并行结构能够比 CPU 更高效地并行处  
340 理大型数据块的算法，被广泛应用于图形和视频渲染等领域。GPU 的线程  
341 数量很多并且能够在线程之间快速切换，因此能够确保硬件一直处于繁忙  
342 状态，并且有效地隐藏内存延迟。不仅如此，将这些特点与现代 CPU 中可  
343 用的存储器带宽 (memory bandwidth) 层相结合，从而能够更高效地提升  
344 GPU 的性能。

345 GPU 是一个可编程的处理器，在 GPU 上能够并行运行数千个核心。  
346 GPU 程序也被称为内核 (kernel)，是通过单程序多数据 (Single Program

Multiple Data, SPMD) 模型进行编写而成。为了避免将 SPMD 和上文提及的 SIMD 混淆, 此处对二者进行比较说明。SPMD 允许多个线程使用它们自身的指令指针运行相应的任务, 而 SIMD 则恰好相反, SIMD 是完全相同的指令被多次执行, 不能像 SPMD 那样同时执行不同的操作。因此, SPMD 的运用更为宽泛。不仅如此, SIMD 定义的是执行模型 (execution model), 模型中有相应的专门用于并行化的程序和硬件, 但是 SPMD 实际上是一个程序描述符 (program descriptor)。

SPMD 模型在 GPU 上的工作原理是启动数千个线程用来运行同一个内核, 进而处理不同的数据。SPMD 是一种可用于实现并行的范式, 其特征是将任务细分为更小的部分, 并从单个程序中并行运行这些部分。经过不断的优化, GPU 为 SPMD 模型提供了非常高的浮点运算吞吐量 (throughput), 因此, 现代 GPU 有着强大的计算能力。

随着技术的不断更新迭代以及现代科学技术的发展需求, 所有现代 GPU 均为 GPGPU (General purpose computing on graphic processing units, 图形处理器上的通用计算)。GPGPU 方法是用于 GPU 通用目的的计算, 除了用于计算机图形计算的传统用途之外, 还能够将 GPU 广泛应用于各科学技术领域中。GPGPU 将 GPU 用于通用计算中不仅能够加速程序的一部分, 还能将剩余部分在 CPU 上继续运行, 进而增强 CPU 架构。最终 GPGPU 通过结合 CPU 和 GPU 的处理能力, 创建整体上更快、高性能的应用程序。

### 6.1.1 蒙特卡洛模拟和 GPU

上文提到, 当采用蒙特卡洛模拟对期权进行定价时, 随机数生成这一步骤占用了大部分的计算时长。

GPU 能够针对蒙特卡洛模拟中的随机数生成和对数价格路径生成进行并行加速。因为, 在标准计算机上生成百万条路径通常也十分耗时, 并且使用 cpu 也无法满足计算高精度的需求, 通过使用 GPU 的并行路径可以有效地解决该问题。由于 GPU 能够非常高效地实现并行化, 它能够将每个路径分配给单个线程, 从而并行模拟成千上万条路径, 进而大幅降低路径生成这一步骤的耗能和计算时间。

## 6.2 具体应用：CUDA

GPU 并行编程是加速密集型工作负载处理的最佳方法之一,其中 CUDA (Compute Unified Device Architecture) 编程模型提供了一个连接软件应用程序及其在 GPU 硬件上实现的一个抽象化 GPU 架构。CUDA 架构利用了不同的处理方法,采用一组流多处理器 (streaming multiprocessors, 简称为 SM) 执行相同的指令集,包括处理不同数据区域上多线程的分支条件, Nivida® 创造了单指令多线程 (single instruction multiple threads, SIMT) 这个术语用来描述这种架构。

CUDA 平台是一个软件层,它可以直接访问 GPU 的虚拟指令集和并行计算元素,以执行计算内核。

### 6.2.1 CUDA 具体应用：CuRAND

使用 cuRAND 库生成随机数。

## 6.3 分析比较

插入图片：带有 GPU 相关结果的图]

## 7 结论与致谢

本文从软件和硬件对蒙特卡洛模拟算法进行了最大程度的优化,实现了计算速率上的大幅提升。

本研究由安信证券股份有限公司与上海镓链科技有限公司联合支持。

## 8 参考文献

## 附录

(1) 测试机器参数

(2) 具体结果

操作系统	Ubuntu 22.04.1 LTS, 64-bit
处理器	12th Gen Intel(R) Core i7-12700
内存	16 GB 3200 MHz DDR4
图形处理器	Mesa Intel(R) Graphics (ADL-S GT1)

表 7: 期权定价参数