'MODEL' MARRIAGES: PAIRING IMPLEMENTATIONS & HARDWARE

Luke Markham ML Engineer - Graphcore





Optimizing model runtimes



Imagine...

- -> You're running a quant team
- -> You've validated a model
- -> It's too slow to run in production
- -> You need to optimize its runtime

(RUNTIME RELEVANT) INGREDIENTS OF A PRODUCTION MODEL







(RUNTIME RELEVANT) INGREDIENTS OF A PRODUCTION MODEL























CASE STUDY



CRR ALGORITHM – AN OVERVIEW

- Forward pass: create a tree of prices by assuming price of the option will increase/decrease by a specific factor
- Backward pass: iteratively calculate binomial value at the previous timepoint. This is the expected value discounted by the risk-free rate.

















Software Implementation





Software Implementation Flexible MIMD processor. Great ease-of-use experience. Vector vs Scalar CPU SIMD processor and high flop Floating point Scalar is the base case. ٠ rates. GPU representation Vector implementations are usually faster. ٠ MIMD processor. High Vectorisation not always possible. ٠ FLOP rates. Large on-Language & chip SRAM IPU Framework Matrix operations, Tree-based workloads, Iterative algorithms, How does it Fixed sequential steps vs batched inputs, or conditional control loops, recursive compare to FLOPS computed convergence vs other map functions, etc. algorithms. on-chip and flow. condition? per Byte read off-chip? ₿ ${f O}$ Т Termination Input-Memory Model Arithmetic Parallel Sequential dependent condition footprint intensity Sc Theory

routing

Software Implementation



Software Implementation



Software Implementation



Software Implementation



Software Implementation



Software Implementation



CASE STUDY – SOFTWARE AND HARDWARE





COX ROSS RUBINSTEIN

Software implementation

Scalar implementations are the base case.

Can be vectorized in multiple ways: across the tree, or across multiple options

```
double Binomial(int n, double Spot, double K, double r, double q, double v, doub
le T, char PutCall, char OpStyle)
 int i, j;
 vector<vector<double> > S(n + 1, vector<double>(n + 1));
 vector<vector<double> > 0p(n + 1, vector<double>(n + 1));
 double dt, u, d, p;
 dt = T / n;
u = exp(v*sqrt(dt));
 d = 1 / u;
 p = (exp((r - q)*dt) - d) / (u - d);
 // Build the binomial tree
 for (j = 0; j <= n; j++) {
  for (i = 0; i <= j; i++) {
   S[i][j] = Spot*pow(u, j - i)*pow(d, i);
 // Compute terminal payoffs
 for (i = 0; i <= n; i++) {</pre>
  if (PutCall == 'C')
   Op[i][n] = max(S[i][n] - K, 0.0);
   Op[i][n] = max(K - S[i][n], 0.0);
 // Backward recursion through the tree
 for (j = n - 1; j \ge 0; j - -) {
  for (i = 0; i \le j; i++) {
   if (0pStyle == 'E')
   Op[i][j] = exp(-r*dt)*(p*(Op[i][j + 1]) + (1 - p)*(Op[i + 1][j + 1]));
   else {
    if (PutCall == 'C')
    Op[i][j] = max(S[i][j] - K, exp(-r*dt)*(p*(0p[i][j + 1]) + (1 - p)*(0p[i + 1))
1][j + 1])));
     Op[i][j] = max(K - S[i][j], exp(-r*dt)*(p*(0p[i][j + 1]) + (1 - p)*(0p[i + 1))
1][j + 1])));
 return Op[0][0];
```

VECTORISED IMPLEMENTATIONS CAN HAVE SIDE EFFECTS





AUTOCALLABLES

Break statements inside for loops... vectorization=hard.

```
def AutoCallableNote(valuationDate, couponDates, strike, pastFixings,
    autoCallBarrier, couponBarrier, protectionBarrier, hasMemory, finalRedemptionFormula,
    coupon, notional, dayCounter, process, generator, nPaths, curve):
    # loop through all simulated paths
    for path in paths:
        payoffPV = 0.0
        unpaidCoupons = 0
        hasAutoCalled = False
        # loop through set of coupon dates and index ratios
        for date, index in zip(couponDates, (path / strike)):
            # if autocall event has been triggered, immediate exit from this path
            if(hasAutoCalled): break
            payoff = 0.0
```



MAPPING OPTIONS TO TILES





OPTION PRICING BENCHMARK

Throughput comparison between IPU and CPU



American Call Options – IPU vs CPU



THANKYOU

For:

- Updates on our GPU benchmarking
- Updates on IPU performance
- Trying out Cox-Ross-Rubinstein out on IPUs for yourself!
- Any other Graphcore related queries

Email lukem@graphcore.ai





APPENDIX





ABSTRACT

"'Model' marriages: Pairing implementations and hardware".

Box famously said that all models are wrong but some are useful. Financial firms could add a proviso: "a model is only useful if we can execute it quickly and efficiently". High performance depends on three variables: the problem size/shape, the implementation code, and the hardware platform. But financial models vary widely, hardware choices expand every few months, and implementation options abound.

Technologists can benefit significantly from a systematic approach to understanding how the three variables interact. Luke has just such an approach. With the help of a real-life study that Graphcore performed for a hedge fund (optimizing Cox-Ross-Rubenstein to price millions of options), Luke will explain a process that categorizes and grades workloads and determines the best pairing of implementation and hardware.

