

CHAPTER 12

Differentiation of the Simulation Library

12.1 ACTIVE CODE

In this chapter, we apply the AAD library from the previous chapter to differentiate our simulation library from [Part II](#). More precisely, we apply AAD to compute the differentials of prices to *model parameters*. In the next chapter, we discuss differentials of prices to *market variables*, also called market risks or *hedge coefficients*.

What will clearly appear is that an efficient AAD library is certainly necessary, but absolutely not sufficient for the efficient differentiation of complex valuation code. It takes work, art, and skill to correctly instrument calculation code. An efficient AAD library provides all the necessary pieces, but to articulate them together in an efficient instrumentation requires an intimate knowledge of the instrumented algorithms (known as *domain-specific knowledge*) and a comprehensive understanding of the inner mechanics of the AAD library. The notion that one can template calculation code, link to an efficient AAD library,

and obtain differentials with AAD speed is a myth. Instrumentation takes at least as much skill and effort as the production of the AAD library itself. A naive instrumentation does produce differentials in constant time, but it may be a very long constant time, often slower than finite differences. While the notion of an efficient instrumentation depends on the differentiated code, this chapter provides general advice and battle-tested methodologies while instrumenting the simulation code of [Part II](#) with remarkable results.

It is especially hard to instrument calculation code that was not written with AAD in mind to start with. We wrote our simulation library knowing that we would eventually instrument it for AAD. This will make its instrumentation substantially easier. In fact, we wrote templated code to start with throughout the simulation library. But not all of it. In particular, we did *not* template the code for the generation of uniform and Gaussian random numbers. All that part of the code uses doubles to represent numbers.

Why did we code it this way? We know that the final result, the average payoff across the simulated paths, depends on the random numbers. Those numbers have non-zero adjoints. However, our goal is to differentiate the final result with respect to the model parameters: spot, local volatilities, and so forth. We know that the random numbers do *not* depend on those parameters. Hence, their adjoints do not propagate all the way back to them.

Formally, with the notations of previous chapters, we differentiate:

$$V_0[(G_i)_{1 \leq i \leq N}, \mathbf{a}] = \frac{1}{N} \sum_{i=1}^N g[\mathbf{h}(G_i, \mathbf{a})]$$

with respect to the vector of model parameters \mathbf{a} . G_i is the vector of Gaussian random numbers for path i , \mathbf{h} is the path generation function, and g is the payoff function. It is clear from those notations that G_i is another argument to \mathbf{h} and has no dependency on \mathbf{a} . The derivatives we want to compute are:

$$\frac{\partial V_0[(G_i)_{1 \leq i \leq N}, a]}{\partial a} = \frac{1}{N} \sum_{i=1}^N \frac{\partial g[h(G_i, a)]}{\partial h} \frac{\partial h(G_i, a)}{\partial a}$$

The G_i s do not depend on a . In particular,

$$\frac{\partial G_i}{\partial a} = 0$$

In AAD lingo, we call the G_i s *inactive*.

We know that because we know exactly how MC simulations work: this is domain-specific knowledge in action. Because we know that, we don't need to propagate from the random numbers at all. Hence, we don't need to propagate *to* them, either. Their adjoints don't contribute to the final derivatives, so we don't need to compute these adjoints at all. To instantiate the code that produces random numbers with the Number type would result in unnecessary AAD overhead at evaluation time (to put operations and local derivatives on tape) and propagation time (conducting unnecessary propagations to and from the random numbers). Because we know that, we can easily leave the random numbers and their generation outside of the adjoint logic. All it takes is use *doubles* for those calculations in place of a templated number type.

Code that uses the templated number type is said to be *instrumented*. This code is recorded on tape and participates in the adjoint propagation at the cost of AAD overhead. Code that uses native types like *doubles* is said to be *noninstrumented*. It is not recorded or propagated and free of AAD overhead. Calculations that depend directly or indirectly on the inputs and, therefore, must participate in adjoint propagation, are called *active*. Code that is independent of inputs and, therefore, does not contribute to their adjoints, is called *inactive*. We have a first and probably most important rule for the efficient differentiation of calculation code:

- Only instrument active code.

We call this major optimization *selective instrumentation*. On the other hand, if we forget to instrument active code, we produce wrong derivatives. Active adjoints are not propagated and their contribution is lost. If we instrument inactive code, like random number generators, we produce unnecessary AAD overhead. We must instrument all the active code, but only the active code. The identification of inactive code is a key step of an efficient differentiation.

Readers are encouraged to review our simulation code and verify that we correctly applied selective instrumentation.

12.2 SERIAL CODE

We start with the serial code. We have nothing to do about random numbers, and the models and products are already instrumented (although we will need minor extensions on models). Therefore, our work focuses on the template algorithm of [Chapter 6](#), which we recall here:

```

1 // MC simulator: free function that conducts simulations
2 //      and returns a matrix (as vector of vectors) of payoffs
3 //      (0..nPath-1 , 0..nPay-1)
4 inline vector<vector<double>> mcSimul(
5     const Product<double>&           prd,
6     const Model<double>&             mdl,
7     const RNG&                      rng,
8     const size_t                     nPath)
9 {
10    // Work with copies of the model and RNG
11    //      which are modified when we set up the simulation
12    // Copies are OK at high level
13    auto cMdl = mdl.clone();
14    auto cRng = rng.clone();
15
16    // Allocate results
17    const size_t nPay = prd.payoffLabels().size();
18    vector<vector<double>> results(nPath, vector<double>(nPay));
19    // Init the simulation timeline
20    cMdl->allocate(prd.timeline(), prd.defline());
21    cMdl->init(prd.timeline(), prd.defline());
22    // Init the RNG
23    cRng->init(cMdl->simDim());
24    // Allocate Gaussian vector
25    vector<double> gaussVec(cMdl->simDim());
26    // Allocate and initialize path
27    Scenario<double> path;
28    allocatePath(prd.defline(), path);
29    initializePath(path);
30
31    // Iterate over paths
32    for (size_t i = 0; i < nPath; i++)
33    {
34        // Next Gaussian vector, dimension D
35        cRng->nextG(gaussVec);
36        // Generate path, consume Gaussian vector
37        cMdl->generatePath(gaussVec, path);
38        // Compute payoffs
39        prd.payoffs(path, results[i]);
40    }
41
42    return results; // C++11: move
43 }
```

We now differentiate this code, instantiating valuation code with the Number type in place of doubles and inserting additional logic for the back-propagation of adjoints and, more generally, the management of the tape. Following the discussion of the previous chapter, we back-propagate adjoints path-wise to avoid unreasonable tape growth.

```

1 #include "AAD.h"
2
3 // AAD instrumentation of mcSimul()
4
5 // returns the following results:
6 struct AADSsimulResults
7 {
8     AADSsimulResults(
9         const size_t nPath,
10        const size_t nPay,
11        const size_t nParam) :
12             payoffs(nPath, vector<double>(nPay)),
13             aggregated(nPath),
14             risks(nParam)
15     {}
16
17     // matrix(0..nPath - 1, 0..nPay - 1) of payoffs
18     // same as mcSimul()
19     vector<vector<double>> payoffs;
20
21     // vector(0..nPath) of aggregated payoffs
22     vector<double> aggregated;
23
24     // vector(0..nParam - 1) of risk sensitivities
25     // of aggregated payoff, averaged over paths
26     vector<double> risks;
27 };
28
29 // Default aggregator = 1st payoff = payoff[0]
30 const auto defaultAggregator = [] (const vector<Number>& v) {return v[0]; };
31
32 template<class F = decltype(defaultAggregator)>
33 inline AADSsimulResults
34 mcSimulAAD(
35     const Product<Number>& prd,
36     const Model<Number>& mdl,
37     const RNG& rng,
38     const size_t nPath,
39     const F& aggFun = defaultAggregator)
40 {
41     // Work with copies of the model and RNG
42     // which are modified when we set up the simulation
43     // Copies are OK at high level
44     auto cMdl = mdl.clone();
45     auto cRng = rng.clone();
46
47     // Allocate path and model
48     // do not initialize yet
49     scenario<Number> path;
50     allocatePath(prd.defline(), path);
51     cMdl->allocate(prd.timeline(), prd.defline());
52
53     // Dimensions
54     const size_t nPay = prd.payoffLabels().size();
55     const vector<Number*>& params = cMdl->parameters();
56     const size_t nParam = params.size();
57
58     // ...

```

The function has the same inputs as the original, except they are now instantiated with Numbers instead of doubles and an additional “aggregator.” Following the discussion of the previous chapter, we are dif-

ferentiating one result for now. In case the product has multiple payoffs, we differentiate an aggregate provided by client code as a function of the payoffs, represented by a lambda or another callable type. For example, the product may be a portfolio of transactions, and we may want to differentiate the value of the entire book, the sum of the values of the transactions weighted by the notional. In this case, we would apply an aggregator that computes the “payoff” of the portfolio as the sum of the payoffs of the products, weighted by the notional:

```
vector<double> notional;
// fill notional
// ...
const auto portfolioAggregator = [&notional] (const vector<Number>& v)
{
    return inner_product(
        notional.begin(),
        notional.end(),
        v.begin(),
        Number(0.0));
}
// call AADSimulResults() with portfolioAggregator() as aggregator
```

where *inner_product()* is a standard C++ library algorithm from header <numeric> that does exactly what its name says. Note that it is not the values, but the payoffs that are aggregated. The aggregator effectively turns multiple payoffs into one aggregated payoff and it is this aggregated payoff that is differentiated. The default aggregator simply picks the first payoff.

The return type is different: we return a matrix of path-wise payoffs as previously, but we also return a vector of path-wise aggregated payoffs and the vector of all the sensitivities of its value (path-wise average) to all the model parameters.

Like in valuation, we start by cloning the model and the RNG, modified by the subsequent code. We allocate the path and the model's working memory *but we don't initialize them yet* and pick the dimensions: number of payoffs and number of parameters. We take a refer-

ence on the vector of *pointers* to the model parameters so we easily address it thereafter under the name *params*.

What follows is new and deserves an explanation:

```
1 // ...
2 // AAD - 1
3 // Access tape
4 Tape& tape = *Number::tape;
5 // Clear and initialize tape
6 tape.clear();
7
8 // Put parameters on tape
9 // note this also initializes all adjoints
10 cMdl->putParametersOnTape();
11
12 // Init the model
13 // CAREFUL: initialization is recorded on tape
14 // Model parameters must be on tape prior to initialization
15 cMdl->init(prd.timeline(), prd.define());
16
17 // Initialize path
18 initializePath(path);
19
20 // Mark the tape straight after initialization
21 tape.mark();
22
23 // ...
24
25 // ...
26
27 // ...
28 // ...
```

We take a reference to the tape and initialize it.

The next thing we do (line 13) is put all the model parameters on tape. They are the inputs to our calculation with respect to which we compute derivatives. They must go on tape before any calculation is executed. See the discussion on page 391 for details and a simple example.

This is the modification we need in our simulation code: we must develop a method on the base Model class to put all the parameters on tape. The modification is minor, because we can implement it directly on the base class without modification to concrete models. The (small) difficulty is that the Model class is templated in its number type, but only Numbers go on tape. We must implement some simple template specialization magic so that

Model < Number >:: putParametersOnTape() puts its parameters on tape, whereas for any other type T,

Model < T >:: putParametersOnTape() does nothing.¹ Readers unfamiliar with template specialization may learn all about it in [76].

The code below completes the base Model class in mcBase.h.

```
// Model base class
// ...

// Put parameters on tape, only valid for T = Number
void putParametersOnTape()
{
    putParametersOnTapeT<T>();
}

private:

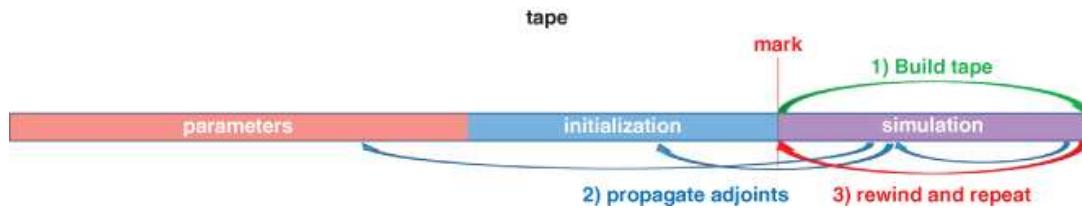
// Template specialization here

// If T is not Number : do nothing
template<class U>
void putParametersOnTapeT()
{
}

// If T is Number : put all parameters on tape
template <>
void putParametersOnTapeT<Number>()
{
    for (Number* param : parameters()) param->putOnTape();
}
```

Back to our simulation code, the next thing we do (lines 18 and 21) is initialize the model and the scenario. It is crucial to do this here and not before. Initialization must be conducted *after* the parameters are on tape. The initialization of a model typically conducts calculations with the model parameters. In Black and Scholes, we pre-calculate deterministic amounts for the subsequent generation of the paths. In Dupire, we pre-interpolate the local volatilities in time. So the parameters must be on tape, or these calculations result in nodes with dangling pointers to nonexistent parameter adjoints. We again refer to the detailed discussion on page 391.

At this point the parameters, and the initial calculations, those that are not repeated path-wise, are all on tape and correctly linked. We *mark* the tape there (line 24). Remember, we are going to compute path-wise derivatives. We will construct the tape for a path, propagate adjoints through it, and *rewind* it for the next path. We are going to rewind it *up to the mark* so that the inputs and the initialization remain on tape throughout the whole simulation, their adjoints accumulating contributions across simulations.



Next in our simulation code, we initialize the RNG, preallocate working memory for the simulations, and allocate results. This is admin code; it is not recorded and no different than valuation.

```

1 // ...
2
3 // Init the RNG
4 cRng->init(cMdl->simDim());
5
6 // Allocate workspace
7 vector<Number> nPayoffs(nPay);
8 // Gaussian vector
9 vector<double> gaussVec(cMdl->simDim());
10
11 // Results
12 AADSimulResults results(nPath, nPay, nParam);
13
14 // ...

```

The important initialization phase is complete; we now proceed with the simulations:

```

1 // ...
2
3 // Iterate through paths
4 for (size_t i = 0; i<nPath; i++)
5 {
6     // AAD - 2
7     // Rewind tape to mark
8     // parameters and initializers stay on tape
9     // but the simulation is overwritten
10    tape.rewindToMark();
11    //
12
13    // Next Gaussian vector, dimension D
14    cRng->nextG(gaussVec);
15    // Generate path, consume Gaussian vector
16    cMdl->generatePath(gaussVec, path);
17    // Compute payoffs
18    prd.payoffs(path, nPayoffs);
19    // Aggregate
20    Number result = aggFun(nPayoffs);
21
22    // AAD - 3
23    // Propagate adjoints
24    result.propagateToMark();
25    // Store results for the path
26    results.aggregated[i] = double(result);
27    convertCollection(
28        nPayoffs.begin(),
29        nPayoffs.end(),
30        results.payoffs[i].begin());
31    //
32 }
33
34 // ...

```

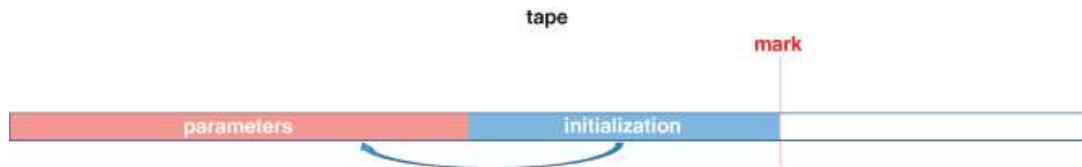
A path is processed in three stages:

1. Rewind the tape *to the mark* (line 10) so as to reuse the simulation piece of the tape across simulations and avoid growth while leaving the model parameters and initialization on tape throughout all simulations so that their adjoints correctly accumulate contributions from all the simulations.
2. Generate and evaluate the path (lines 13–20). This code is identical to the valuation code. It executes the three steps of the processing of a path: generate Gaussian numbers, simulate a path, evaluate payoffs over the path. The difference is that this code is now instantiated with Numbers. So it goes on tape when executed and builds the simulation piece of the tape, overwriting the previous simulation so that the size of the tape remains fixed to the size of one simulation; see figure on page 414.

We have an additional line 20, where the payoffs for the current path are aggregated in one result, with an aggregator provided by the client code, as explained earlier.

3. Propagate adjoints from the result (which adjoint is seeded to 1) *to the mark*. We back-propagate the adjoints of the current simulation, but only those. As a result, the adjoints on the persistent piece of tape, where the parameters and initialization were recorded, accumulate their contribution to the current simulation.

After the simulation loop completes, the adjoints on the persistent piece of tape accumulated their contribution to all the simulations. What is left is back-propagate adjoints through the persistent piece so that the adjoints of the entire Monte-Carlo valuation accumulate into the model parameters.



This takes *one* line of code:

```

1 // ...
2
3 // AAD - 4
4 // Mark = limit between pre-calculations and path-wise operations
5 // Operations above mark have been propagated and accumulated
6 // We conduct one propagation mark to start
7 Number::propagateMarkToStart();
8 //
9
10 // ...

```

It follows that we back-propagate simulation adjoints repeatedly, but initialization adjoints only once. This is another one of the strongest optimizations in our code. If we propagated adjoints all the way, repeatedly over each simulation, to differentiate a Monte-Carlo valuation would be *many times* slower. This optimization, to a large extent, is the dual of the strongest optimization we included in our simulation

code of [Chapter 6](#). We gathered as much work as possible into the initialization, conducted once, so as to minimize work during the repeated simulations. It follows that we have large a initialization tape and small simulation tapes. In the differentiation algorithm, we back-propagate repeatedly over the small tapes and once over the large tape, reaping the benefits of our optimization of [Chapter 6](#) for differentiation, too.

All that remains is pick sensitivities to parameters on the parameters adjoints, pack,² clean, and return:

```
1 // ...
2
3 // Pick sensitivities, summed over paths, and normalize
4 transform(
5     params.begin(),
6     params.end(),
7     results.risks.begin(),
8     [nPath] (const Number* p) {return p->adjoint() / nPath; });
9
10 // Clear the tape
11 tape.clear();
12
13 return results;
14 }
```

The code is found in mcBase.h.

12.3 USER INTERFACE

We must upgrade our user interface in main.h to run the differentiation algorithm. We develop two functions: *AADriskOne()*, which differentiates one selected payoff in a product, and *AADriskAggregate()*, which differentiates a portfolio of payoffs with given notional. Both functions take a model and a product in the memory store, developed in [Section 6.5](#), run the differentiation algorithm, and return:

1. The labels of the payoffs in a vector<string>
2. The corresponding values in a vector<double>

3. The value of the differentiated aggregate, as a double
4. The names of all the parameters in a vector<string>
5. The corresponding risk sensitivities of the aggregate to parameters in a vector<double>

```

1   struct
2   {
3       vector<string> payoffIds;
4       vector<double> payoffValues;
5       double riskPayoffValue;
6       vector<string> paramIds;
7       vector<double> risks;
8   } results;

```

The differentiation algorithm works with instrumented models and products, Model <Number> and Product <Number>. We upgrade the memory store in store.h so we can store and retrieve both instrumented and non-instrumented models and products. In practice, we store a pair of models or products, one instrumented and one not:

```

1 #include "mcBase.h"
2 #include "mcMdl.h"
3 #include "mcPrd.h"
4 #include <unordered_map>
5 #include <memory>
6 using namespace std;
7
8 using ModelStore =
9     unordered_map<string,
10         pair<unique_ptr<Model<double>>, unique_ptr<Model<Number>>>>;
11 using ProductStore =
12     unordered_map<string,
13         pair<unique_ptr<Product<double>>, unique_ptr<Product<Number>>>>;
14
15 ModelStore modelStore;
16 ProductStore productStore;

```

We upgrade the functions to initialize models and products in memory accordingly:

```

1 void putBlackScholes(
2     const double          spot,
3     const double          vol,
4     const bool            qSpot,
5     const double          rate,
6     const double          div,
7     const string&         store)
8 {
9     // We create 2 models, one for valuation and one for risk
10    unique_ptr<Model<double>> mdl = make_unique<BlackScholes<double>>(
11        spot, vol, qSpot, rate, div);
12    unique_ptr<Model<Number>> riskMdl = make_unique<BlackScholes<Number>>(
13        spot, vol, qSpot, rate, div);
14
15    // And move them into the map
16    modelStore[store] = make_pair(move(mdl), move(riskMdl));
17 }
18
19 // Same pattern for putDupire
20
21 void putEuropean(
22     const double          strike,
23     const Time             exerciseDate,
24     const Time             settlementDate,
25     const string&         store)
26 {
27     // We create 2 products, one for valuation and one for risk
28     unique_ptr<Product<double>> prd = make_unique<European<double>>(
29         strike, exerciseDate, settlementDate);
30     unique_ptr<Product<Number>> riskPrd = make_unique<European<Number>>(
31         strike, exerciseDate, settlementDate);
32
33     // And move them into the map
34     productStore[store] = make_pair(move(prd), move(riskPrd));
35 }
36
37 // Same pattern for other products

```

We template the functions to retrieve models and products, so client codes calls `getModel < double > ()` for pricing or `getModel < Number > ()` for risk:

```
1 template<class T>
2 const Model<T>* getModel(const string& store);
3
4 template<>
5 const Model<double>* getModel(const string& store)
6 {
7     auto it = modelStore.find(store);
8     if (it == modelStore.end()) return nullptr;
9     else return it->second.first.get();
10 }
11
12 template<>
13 const Model<Number>* getModel(const string& store)
14 {
15     auto it = modelStore.find(store);
16     if (it == modelStore.end()) return nullptr;
17     else return it->second.second.get();
18 }
19
20 // Same pattern for products
```

The entire code is found in `store.h`. With the store upgraded to deal with instrumented objects, we develop the user interface for the differentiation of simulations in `main.h` below. We write two functions, one for the differentiation of one named payoff, and one for the differentiation of a portfolio. The project in our repository exports these functions to Excel.

```

1 // Generic risk
2 inline auto AADriskOne(
3     const string&           modelId,
4     const string&           productId,
5     const NumericalParam&  num,
6     const string&           riskPayoff = "")
7 {
8     // Get instrumented model and product
9     const Model<Number>* model = getModel<Number>(modelId);
10    const Product<Number>* product = getProduct<Number> (productId);
11
12    if (!model || !product)
13    {
14        throw runtime_error(
15            "AADrisk() : Could not retrieve model and product");
16    }
17
18    // Random Number Generator
19    unique_ptr<RNG> rng;
20    if (num.useSobol) rng = make_unique<Sobol>();
21    else rng = make_unique<mrg32k3a>(num.seed1, num.seed2);
22
23    // Find the payoff to differentiate
24    size_t riskPayoffIdx = 0;
25    if (!riskPayoff.empty())
26    {
27        const vector<string>& allPayoffs = product->payoffLabels();
28        auto it = find(allPayoffs.begin(), allPayoffs.end(), riskPayoff);
29        if (it == allPayoffs.end())
30        {
31            throw runtime_error("AADriskOne() : payoff not found");
32        }
33        riskPayoffIdx = distance(allPayoffs.begin(), it);
34    }
35
36    // Simulate
37    const auto simulResults = num.parallel
38        ? mcParallelSimulAAD(*product, *model, *rng, num.numPath,
39            [riskPayoffIdx] (const vector<Number>& v)
40            {return v[riskPayoffIdx]; })
41        : mcSimulAAD(*product, *model, *rng, num.numPath,
42            [riskPayoffIdx] (const vector<Number>& v)
43            {return v[riskPayoffIdx]; });
44
45    // We return: a number and 2 vectors :
46    // - The payoff identifiers and their values
47    // - The value of the aggregate payoff
48    // - The parameter identifiers
49    // - The sensitivities of the aggregate to parameters
50    struct
51    {
52        vector<string>  payoffIds;
53        vector<double>  payoffValues;
54        double          riskPayoffValue;
55        vector<string>  paramIds;
56        vector<double>  risks;
57    } results;
58

```

```

59     const size_t nPayoffs = product->payoffLabels().size();
60     results.payoffIds = product->payoffLabels();
61     results.payoffValues.resize(nPayoffs);
62     for (size_t i = 0; i < nPayoffs; ++i)
63     {
64         results.payoffValues[i] = accumulate(
65             simulResults.payoffs.begin(),
66             simulResults.payoffs.end(),
67             0.0,
68             [i](const double acc, const vector<double>& v)
69             { return acc + v[i]; }
70         ) / num.numPath;
71     }
72     results.riskPayoffValue = accumulate(
73         simulResults.aggregated.begin(),
74         simulResults.aggregated.end(),
75         0.0) / num.numPath;
76     results.paramIds = model->parameterLabels();
77     results.risks = move(simulResults.risks);
78
79     return results;
80 }
81
82 inline auto AADriskAggregate(
83     const string&           modelId,
84     const string&           productId,
85     const map<string, double>& notionals,
86     const NumericalParam&   num)
87 {
88     // Get instrumented model and product
89     const Model<Number>* model = getModel<Number>(modelId);
90     const Product<Number>* product = getProduct<Number>(productId);
91
92     if (!model || !product)
93     {
94         throw runtime_error(
95             "AADriskAggregate() : Could not retrieve model and product");
96     }
97
98     // Random Number Generator
99     unique_ptr<RNG> rng;
100    if (num.useSobol) rng = make_unique<Sobol>();
101    else rng = make_unique<mrg32k3a>(num.seed1, num.seed2);
102
103    // Vector of notionals
104    const vector<string>& allPayoffs = product->payoffLabels();
105    vector<double> vnots(allPayoffs.size(), 0.0);
106    for (const auto& notional : notionals)
107    {
108        auto it = find(allPayoffs.begin(), allPayoffs.end(), notional.first);
109        if (it == allPayoffs.end())
110        {
111            throw runtime_error("AADriskAggregate() : payoff not found");
112        }
113        vnots[distance(allPayoffs.begin(), it)] = notional.second;
114    }
115

```

```

116     // Aggregator
117     auto aggregator = [&vnorts](const vector<Number>& payoffs)
118     {
119         return inner_product(
120             payoffs.begin(),
121             payoffs.end(),
122             vnorts.begin(),
123             Number(0.0));
124     };
125
126     // Simulate
127     const auto simulResults = num.parallel
128         ? mcParallelSimulAAD(*product, *model, *rng, num.numPath, aggregator)
129         : mcSimulAAD(*product, *model, *rng, num.numPath, aggregator);
130
131     // We return: a number and 2 vectors :
132     // - The payoff identifiers and their values
133     // - The value of the aggregate payoff
134     // - The parameter identifiers
135     // - The sensitivities of the aggregate to parameters
136     struct
137     {
138         vector<string> payoffIds;
139         vector<double> payoffValues;
140         double riskPayoffValue;
141         vector<string> paramIds;
142         vector<double> risks;
143     } results;
144
145     const size_t nPayoffs = product->payoffLabels().size();
146     results.payoffIds = product->payoffLabels();
147     results.payoffValues.resize(nPayoffs);
148     for (size_t i = 0; i < nPayoffs; ++i)
149     {
150         results.payoffValues[i] = accumulate(
151             simulResults.payoffs.begin(),
152             simulResults.payoffs.end(),
153             0.0,
154             [i](const double acc, const vector<double>& v)
155             { return acc + v[i]; }
156         ) / num.numPath;
157     }
158     results.riskPayoffValue = accumulate(
159         simulResults.aggregated.begin(),
160         simulResults.aggregated.end(),
161         0.0) / num.numPath;
162     results.paramIds = model->parameterLabels();
163     results.risks = move(simulResults.risks);
164
165     return results;
166 }

```

where `mcParallelSimulAAD()`, the parallel version of the differentiation algorithm, is developed next.

For test purposes, and to assess the speed and the accuracy of our derivatives, we also implement bump risk. Bumps are implemented in a trivial manner to produce itemized risk reports³ in a trivial manner. Bump risk is up to thousands of times slower than AAD, and the result-

ing risk sensitivities are expected to be virtually identical. The following self-explanatory function is found in main.h:

```
1 // Returns a vector of values and a matrix of risks
2 // with payoffs in columns and parameters in rows
3 // along with ids of payoffs and parameters
4
5 struct RiskReports
6 {
7     vector<string> payoffs;
8     vector<string> params;
9     vector<double> values;
10    matrix<double> risks;
11 },
12
13 inline RiskReports bumpRisk(
14     const string&             modelId,
15     const string&             productId,
16     const NumericalParam&    num)
17 {
18     auto* orig = getModel<double>(modelId);
19     const Product<double>* product = getProduct<double>(productId);
20
21     if (!orig || !product)
22     {
23         throw runtime_error(
24             "bumpRisk() : Could not retrieve model and product");
25     }
26
27     RiskReports results;
28
29     // Base values
30     auto baseRes = value(*orig, *product, num);
31     results.payoffs = baseRes.identifiers;
32     results.values = baseRes.values;
33
34     // Make copy so we don't modify the model in memory
35     auto model = orig->clone();
36
37     results.params = model->parameterLabels();
38     const vector<double*> parameters = model->parameters();
39     const size_t n = parameters.size(), m = results.payoffs.size();
40     results.risks.resize(n, m);
41
42     // Bump parameters one by one
43     for (size_t i = 0; i < n; ++i)
44     {
45         // Bump
46         *parameters[i] += 1.e-08;
47         // Reval
48         auto bumpRes = value(*model, *product, num);
49         // Unbump
50         *parameters[i] -= 1.e-08;
51
52         // Compute finite differences for all payoffs
53         for (size_t j = 0; j < m; ++j)
54         {
55             results.risks[i][j] = 1.0e+08 *
56                 (bumpRes.values[j] - baseRes.values[j]);
57         }
58     }
59
60     return results;
61 }
```

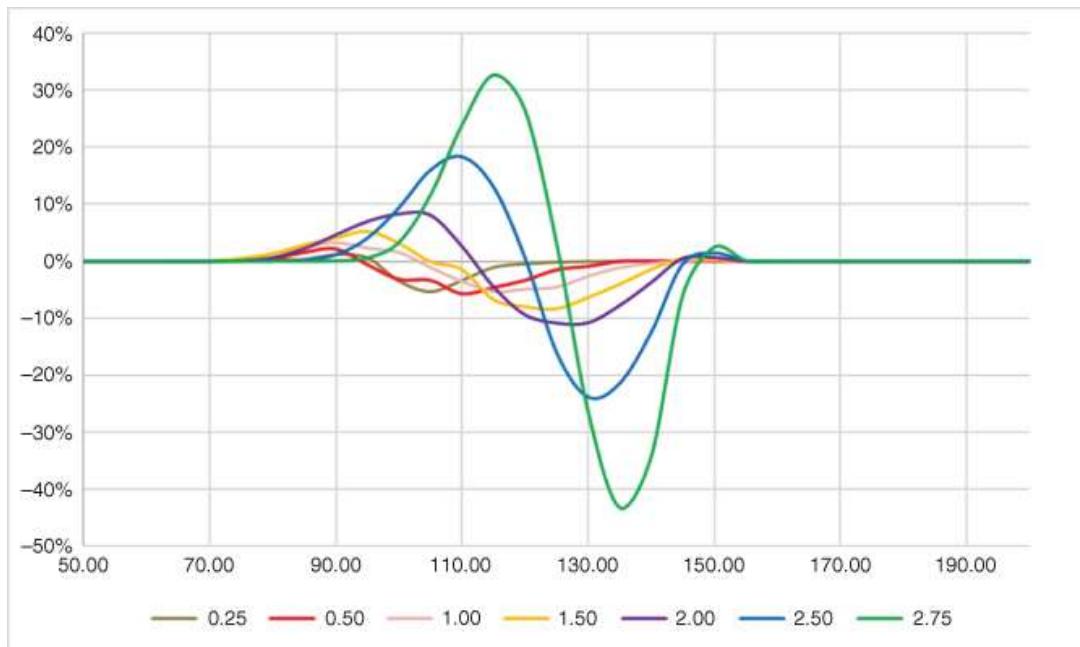
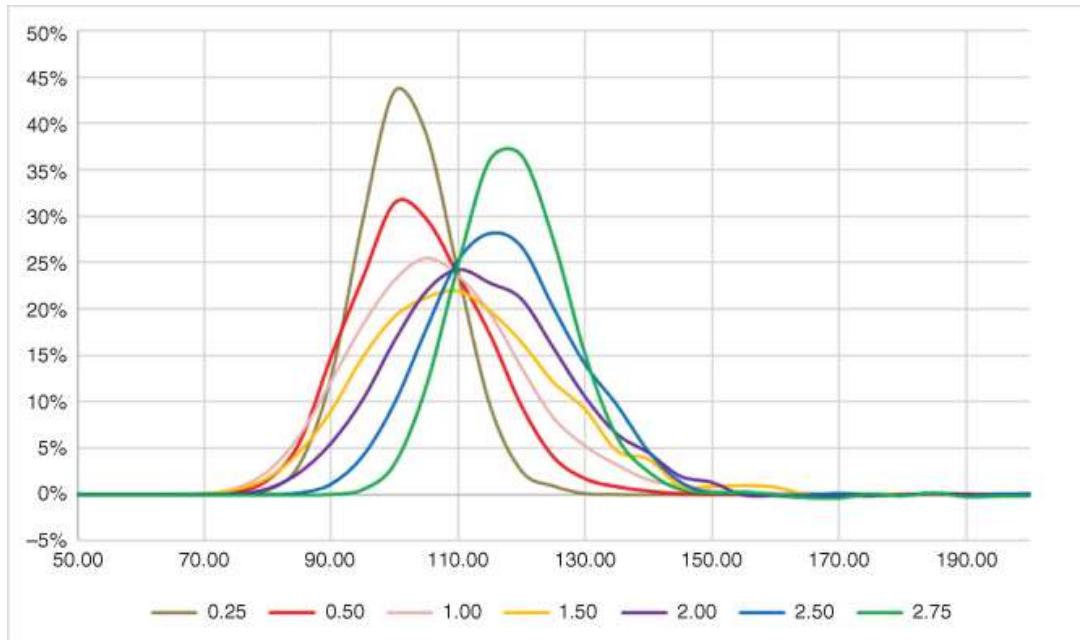
12.4 SERIAL RESULTS

We can now differentiate the results of our simulation library with respect to model parameters. For instance, we compute the sensitivities of a European or barrier option to all local volatilities. Sensitivities to local volatilities, what Dupire calls a *microbucket*, is valuable information for research purpose. For the purpose of trading, however, we really need the sensitivities to tradable market variables, not model parameters. In the context of Dupire's model, we want the sensitivities to the *market-implied* volatilities, what Dupire calls a *superbucket*, because those are risks that may be directly hedged by trading European options in the market. Local volatilities are calibrated to implied volatilities, so superbuckets can be computed from microbuckets, as discussed in detail in [Chapter 13](#). In this chapter, we focus on the production of microbuckets.

We differentiate the example of [Section 6.6](#): a 3y European call of strike 120, together with a 150 barrier, in Dupire's model with a local volatility of 15% over 60 times and 30 spots, simulating 500,000 paths over 156 weekly steps, where a pricing took around 3 seconds with Sobol. With 1,801 sensitivities to compute, we need 1,802 prices with finite differences, so the bump risk takes around 1.5 hours.⁴

Our AAD implementation produces the entire risk in around 8 seconds (8,250 ms), less than three times a single valuation. With the improvements of [Chapter 15](#), when we further accelerate AAD by factor around 2, we produce this report in less than 5 seconds (4,500 ms). With the parallel algorithm, it takes *half a second* (575 ms).

The results are displayed on the chart below, first for the European, then for the barrier, with a smoothing of 5. The charts display the differentials to local volatilities on seven different times, as functions of the spot.



The speed is remarkable and may look suspect. Didn't we establish the theoretical maximum speed for AAD at three times a valuation? With the improvements of [Chapter 15](#), our example computes in the time of 1.5 valuations, *half the theoretical minimum*.

This is mainly due to our selective code instrumentation. A lot of the valuation time is spent generating random numbers and turning them into Gaussians with the not-so-cheap inverse Gaussian, executed **$ND = 78M$** times. All those calculations are inactive. We did not instrument them. They execute without ADD overhead, in the exact

same time in differentiation or evaluation mode. Only the instrumented code bears AAD overhead. The remarkable speed is not only due to AAD's constant time, and the implementation of the AAD library, but also to the optimizations in our instrumented code, mainly the selective instrumentation and the partial back-propagation.

Risk sensitivities are almost identical to bumping. Sensitivities to local volatilities match to many significant figures. Difference in delta is around $2.5e - 05$. The reason is subtle. In the barrier payoff code, we used a smoothing factor of $x\%$ of the current spot. When the spot is bumped, the smoothing factor changes and affects the price. But this is numerical noise, not market risk. Hence, we made the smoothing factor a *double* in our code so that it does not contribute to derivatives, and the result we get is market risk alone. Our delta is not off; on the contrary, it is more accurate.

12.5 PARALLEL CODE

We now move on to the AAD instrumentation of our *parallel* code from [Chapter 7](#), combining the idioms of that section for parallelism with the techniques of this chapter for differentiation. The instrumentation of the parallel code is more complicated, but well worth the effort, since it is really the combination of parallel simulations and AAD (and expression templates, see [Chapter 15](#)) that gives us this “magic” performance. The AAD library helps. The static tape pointer on the Number type is thread local, so all the calculations executed on a thread land on that thread's own tape. But it doesn't mean our work is done. First, the tape *pointer* is thread local, but we still must declare all the tapes for all the threads and set the thread local pointer to each thread's tape *from that thread*. More importantly:

1. With AAD, the model is a *mutable* object. Even though *Model* `:: generatePath()` is marked `const`, it is not really the case any more; see the discussion on page 391. When the Numbers stored on the model (parameters and initialized amounts) are *read*,

their adjoints on tape are *written* into during back-propagation. So, we must make copies of the model for every thread like other mutable objects.⁵

2. In the single-thread version, we were careful to put the model parameters on tape and initialize the model so the initialized amounts are recorded, too. In the parallel version, we must put the parameters and initializers, not on one tape, but on every thread's tape. The only way to do this is call

Model :: putParametersOnTape() and *Model :: init()* from *each thread* before the first simulation on that thread.

3. Adjoints are accumulated on thread local tapes, so after the simulation loop, the adjoints of the parameters and initializers are scattered across different tapes. We must back-propagate over initialization on multiple tapes and aggregate all the sensitivities to parameters.

The code is essentially a combination of the parallel code of Section 7 with the differentiation code of this chapter, with a special response to these specific challenges.

To address number 2, we separate the initialization in a dedicated function:

```

1 // Put parameters on tape and init model
2 inline void initModel4ParallelAAD(
3     // Inputs
4     const Product<Number>& prd,
5     // Cloned model, must have been allocated prior
6     Model<Number>& clonedMdl,
7     // Path, also allocated prior
8     Scenario<Number>& path)
9 {
10    // Access tape
11    Tape& tape = *Number::tape;
12
13    // Rewind tape
14    tape.rewind();
15
16    // Put parameters on tape
17    // note this also initializes all adjoints
18    clonedMdl.putParametersOnTape();
19
20    // Init the model
21    // not before the parameters are on tape:
22    // initialization uses parameters as arguments
23    clonedMdl.init(prd.timeline(), prd.defline());
24
25    // Path
26    initializePath(path);
27
28    // Mark the tape straight after parameters
29    tape.mark();
30    //
31 }
```

The template algorithm has the same signature as the single-threaded one and returns results in the same format:

```

1 #include "ThreadPool.h"
2
3 // Parallel version of mcSimulaAAD()
4 template<class F = decltype(defaultAggregator)>
5 inline AADSimulResults
6 mcParallelSimulaAAD(
7     const Product<Number>& prd,
8     const Model<Number>& mdl,
9     const RNG& rng,
10    const size_t nPath,
11    const F& aggFun = defaultAggregator)
12 {
13    const size_t nPay = prd.payoffLabels().size();
14    const size_t nParam = mdl.numParams();
15
16    // Allocate results
17    AADSimulResults results(nPath, nPay, nParam);
18
19    // Clear and initialise tape
20    Number::tape->clear();
21    // ...
```

and the entry-level function in main.h is already set up to fire it.

Next, we allocate working memory (without initializing anything quite yet), one mutable object for every thread (including the main thread, so number of worker threads + 1):

```
1 // ...
2
3 // We need one of all these for each thread
4 // 0: main thread
5 // 1 to n : worker threads
6
7 ThreadPool *pool = ThreadPool::getInstance();
8 const size_t nThread = pool->numThreads();
9
10 // Allocate workspace, do not initialize anything yet
11
12 // One model clone per thread
13 // models are mutable with AAD
14 vector<unique_ptr<Model<Number>>> models(nThread + 1);
15 for (auto& model : models)
16 {
17     model = mdl.clone();
18     model->allocate(prd.timeline(), prd.defline());
19 }
20
21 // One scenario per thread
22 vector<Scenario<Number>> paths(nThread + 1);
23 for (auto& path : paths)
24 {
25     allocatePath(prd.defline(), path);
26 }
27
28 // One vector of payoffs per thread
29 vector<vector<Number>> payoffs(nThread + 1, vector<Number>(nPay));
30
31 // ~workspace
32
33 // Tapes for the worker threads
34 // The main thread has one of its own
35 vector<Tape> tapes(nThread);
36
37 // Model initialized on this thread?
38 vector<int> mdlInit(nThread + 1, false);
39
40 // ...
```

We also initialized a vector of tapes on line 35, one for each *worker* thread. The main thread already owns a tape: the static tape declared on the Number type instantiates of copy for the main thread when the application starts, as for any static object.

Finally, we have a vector of booleans on line 38, one for every thread, tracking what threads completed the initialization phase, starting with false for everyone.⁶

We must also initialize the RNGs and preallocate the Gaussian vector, but for that we need the simulation dimension D , and only an initialized model can provide this information. So the next step is initialize the main thread's model (the main thread always has index 0, the worker threads indices being 1 to nThread):

```

1 // ...
2
3 // Initialize main thread
4 initModel4ParallelAAD(prd, *models[0], paths[0]);
5
6 // Mark main thread as initialized
7 mdlInit[0] = true;
8
9 // Init the RNGs, one per thread
10 // One RNG per thread
11 vector<unique_ptr<RNG>> rngs(nThread + 1);
12 for (auto& random : rngs)
13 {
14     random = rng.clone();
15     random->init(models[0]->simDim());
16 }
17
18 // One Gaussian vector per thread
19 vector<vector<double>> gaussVecs
20 (nThread + 1, vector<double>(models[0]->simDim()));
21
22 // ...

```

Next, we dispatch the simulation loop for parallel processing to the thread pool, applying the same techniques as in [Chapter 7](#):

```

1 // ...
2
3 // Reserve memory for futures
4 vector<TaskHandle> futures;
5 futures.reserve(nPath / BATCHSIZE + 1);
6
7 // Start
8
9 size_t firstPath = 0;
10 size_t pathsLeft = nPath;
11 while (pathsLeft > 0)
12 {
13     size_t pathsInTask = min<size_t>(pathsLeft, BATCHSIZE);
14
15     futures.push_back(pool->spawnTask([&, firstPath, pathsInTask]()
16     {
17         // ...
18     }));

```

The body of the lambda executes on different threads. This is the concurrent code. Here, the tape is the executing thread's tape and `pool->threadNum()` returns the index of the executing thread. The first thing we do is note this index, and set the thread local tape accordingly:

```
1 // ...
2
3     const size_t threadNum = pool->threadNum();
4
5     // Use this thread's tape
6     // Thread local magic: each thread has its own pointer
7     // Note main thread = 0 is not reset
8     if (threadNum > 0) Number::tape = &tapes[threadNum - 1];
9
10 // ...
```

This is the core of parallel AAD right here. The tape is accessed through the *thread local* pointer `Number :: tape`. We set it *from the executing thread* to the tape we allocated for that thread in our vector of tapes. From there on, all the subsequent calculations coded in the lambda and executed on this thread are recorded on that tape. The main thread (0) already correctly refers to its global tape. It does not need resetting.

Only at this stage can we put the parameters on tape and initialize the model (other than allocation, we allocated all models earlier, on the main thread, knowing from [Chapter 3](#) that we must avoid concurrent allocations):

```
1 // ...
2
3     // Initialize once on each thread
4     if (!mdlInit[threadNum])
5     {
6         // Initialize and put on the
7         // executing thread's tape
8         initModel4ParallelAAD(
9             prd,
10            *models[threadNum],
11            paths[threadNum]);
12
13        // Mark as initialized
14        mdlInit[threadNum] = true;
15    }
16
17 // ...
```

Next, we pick this thread's RNG and skip it ahead, like in [Chapter 7](#):

```
1 // ...
2
3     // Get a RNG and position it correctly
4     auto& random = rngs[threadNum];
5     random->skipTo(firstPath);
6
7 // ...
```

Finally, the rest of the concurrent code executes the inner batch loop of [Chapter 7](#) with the same contents as the serial differentiation code: rewind tape (to mark), generate random numbers, simulate scenario, compute payoffs, aggregate them, back-propagate adjoints (to mark), store results, repeat):

```

1 // ...
2
3 // And conduct the simulations, exactly same as sequential
4 for (size_t i = 0; i < pathsInTask; i++)
5 {
6     // Rewind tape to mark
7     // Notice : this is the tape for the executing thread
8
9     Number::tape->rewindToMark();
10    // Next Gaussian vector, dimension D
11    random->nextG(gaussVecs[threadNum]);
12    // Path
13    models[threadNum]->generatePath(
14        gaussVecs[threadNum],
15        paths[threadNum]);
16    // Payoffs
17    prd.payoffs(paths[threadNum], payoffs[threadNum]);
18
19    // Propagate adjoints
20    Number result = aggFun(payoffs[threadNum]);
21    result.propagateToMark();
22    // Store results for the path
23    results.aggregated[firstPath + i] = double(result);
24    convertCollection(
25        payoffs[threadNum].begin(),
26        payoffs[threadNum].end(),
27        results.payoffs[firstPath + i].begin());
28 }
29
30 // Remember tasks must return a bool
31 return true;
32 }));
33
34 pathsLeft -= pathsInTask;
35 firstPath += pathsInTask;
36 }
37
38 // Wait and help
39 for (auto& future : futures) pool->activeWait(future);
40
41 // ...

```

Like for serial differentiation, all adjoints are accumulated on the persistent piece of the tape at this stage, and we must propagate the initializers' adjoints back to the parameters' adjoints; see figure on page 416. The difference is that the adjoints are now spread across the multiple tapes belonging to multiple threads, but we are out of the concurrent code and back on the main thread. Never mind; this is one propagation conducted one time; we propagate all the threads' tapes from the main thread:

```

1 // ...
2
3 // Mark = limit between pre-calculations and path-wise operations
4 // Operations above mark have been propagated and accumulated
5 // We conduct one propagation mark to start
6 // On the main thread's tape
7 Number::propagateMarkToStart();
8 // And on the worker thread's tapes
9 Tape* mainThreadPtr = Number::tape;
10 for (size_t i = 0; i < nThread; ++i)
11 {
12     if (mdlInit[i + 1])
13     {
14         // Set tape pointer
15         Number::tape = &tapes[i];
16         // Propagate on that tape
17         Number::propagateMarkToStart();
18     }
19 }
20 // Reset tape to main thread's
21 Number::tape = mainThreadPtr;
22
23 // ...

```

We first propagate over the main thread's tape on line 7. Then, we sequentially propagate the other n tapes, from the main thread. In order to propagate the tape number i from the main thread, we temporarily set the main thread's tape as the tape i (line 15), and propagate (line 17).

In the end, we reset the main thread's tape to its static tape on line 21.⁷

We are almost home safe. All the adjoints are correctly accumulated, but still across multiple tapes and multiple models. The tape of each thread contains the sum of the adjoints over the paths processed on that thread. The adjoints are accessible for every thread's copy of the model. We sum-up the adjoints across models, hence across tapes, normalize, pack, and return:

```

1 // ...
2
3 // Sum sensitivities over threads
4 for (size_t j = 0; j < nParam; ++j)
5 {
6     results.risks[j] = 0.0;
7     for (size_t i = 0; i < models.size(); ++i)
8     {
9         if (mdlInit[i]) results.risks[j] +=
10             models[i]->parameters() [j]->adjoint();
11     }
12     results.risks[j] /= nPath;
13 }
14
15 // Clear the main thread's tape
16 // The other tapes are cleared when the vector of tapes exists scope
17 Number::tape->clear();
18
19 return results;
20 }
```

The code is found in mcBase.h.

12.6 PARALLEL RESULTS

We reproduce the example on page 424, this time in parallel: a 3y call 120 with a 150 barrier, in Dupire's model 1,800 local volatilities (of which 1,080 active) of 15% over 60 times and 30 spots, simulating 500,000 paths over 156 weekly steps with Sobol. With the serial implementation, it took around 8 seconds to compute the microbuckets (precisely 8,250 ms). With the improvements of [Chapter 15](#), it takes less than 5 seconds (precisely 4,500 ms).

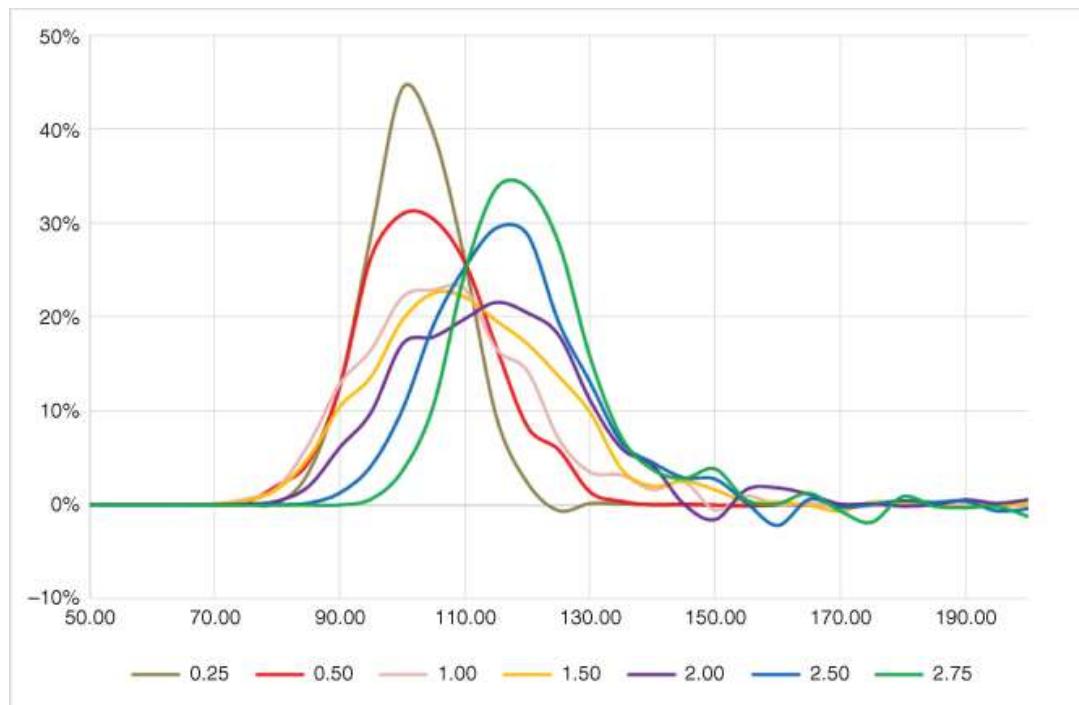
With the parallel implementation, we get the exact same results in exactly one second, with perfect parallel efficiency. With the improvements of [Chapter 15](#), we get them in half a second (precisely 575 ms).

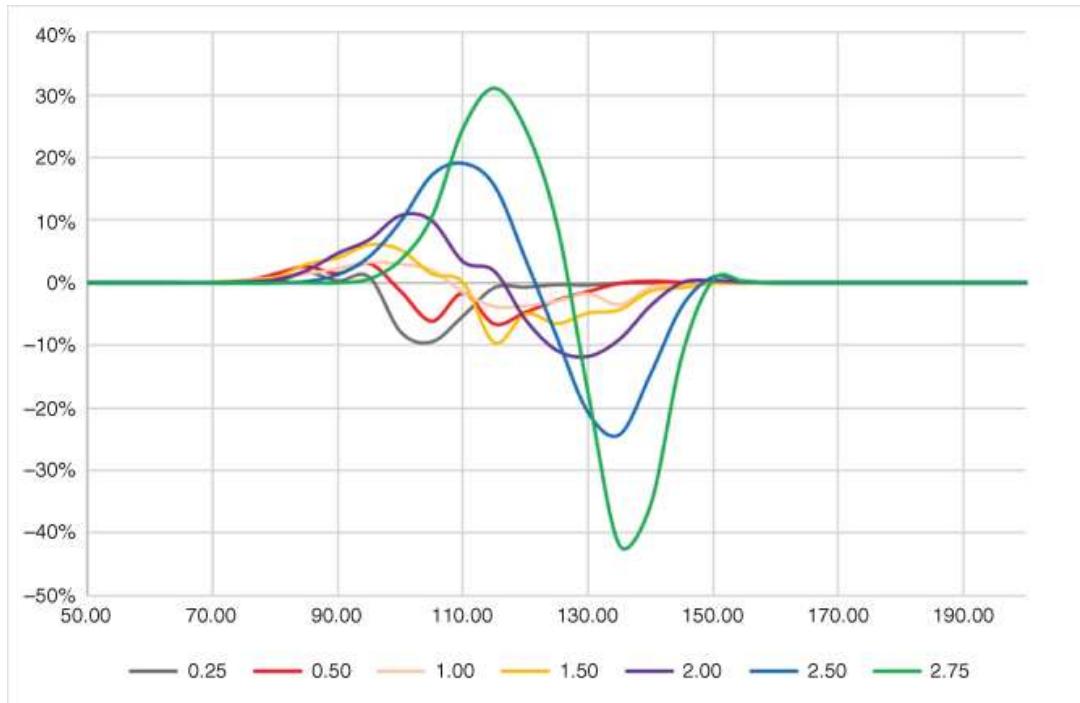
We produce more than 1,000 active sensitivities over 500,000 paths, 156 steps, in half a second.

With an efficient AAD library, coupled with an effective instrumentation of the simulation code, we computed derivatives over simulations with unthinkable speed. Not only can we produce lightning-fast risk

reports on workstations and laptops, including for xVA, we can also conduct research that would have been impossible otherwise.

For instance, a parallel price with Sobol takes 350 ms in our example, that is, more than five minutes for a finite difference risk over the 1,080 active local volatilities (which perhaps may be improved to two to three minutes with substantial effort and invasive code). This is too slow. It is therefore common practice to reduce accuracy to produce risk reports in a manageable time. With 36 monthly time steps and 100,000 paths, a parallel price takes 20 ms, so a finite difference risk report could perhaps be produced in as little as ten seconds. The results, however, are not pretty, as seen in the charts below, first for the European, then for the barrier, to be compared to the charts on page.

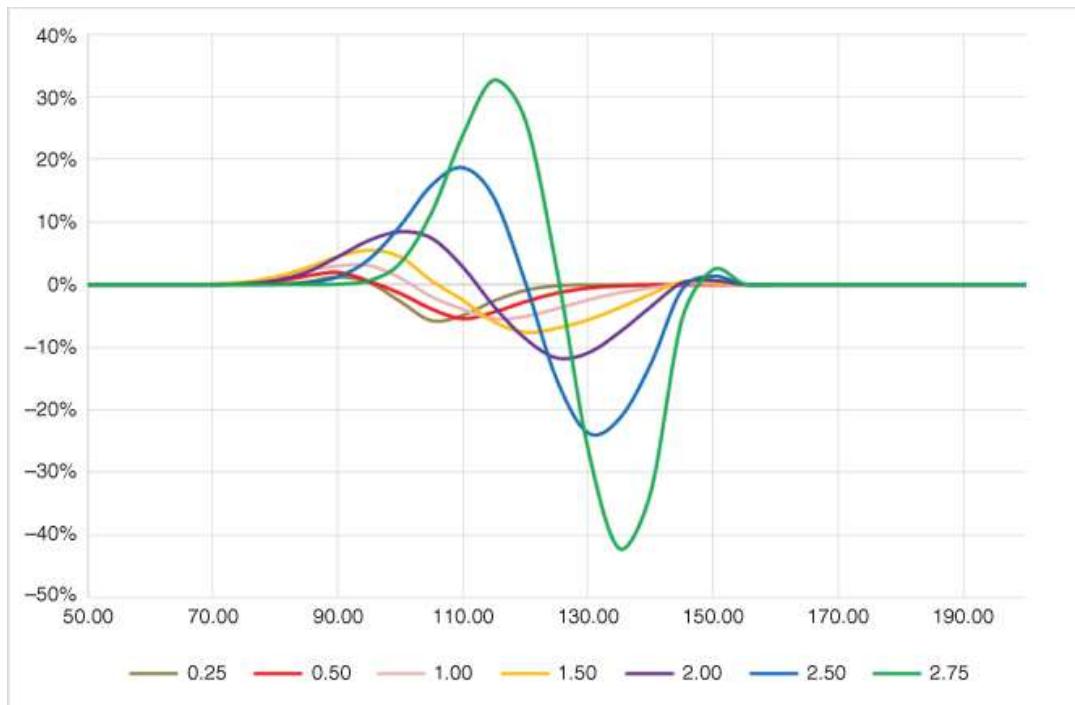
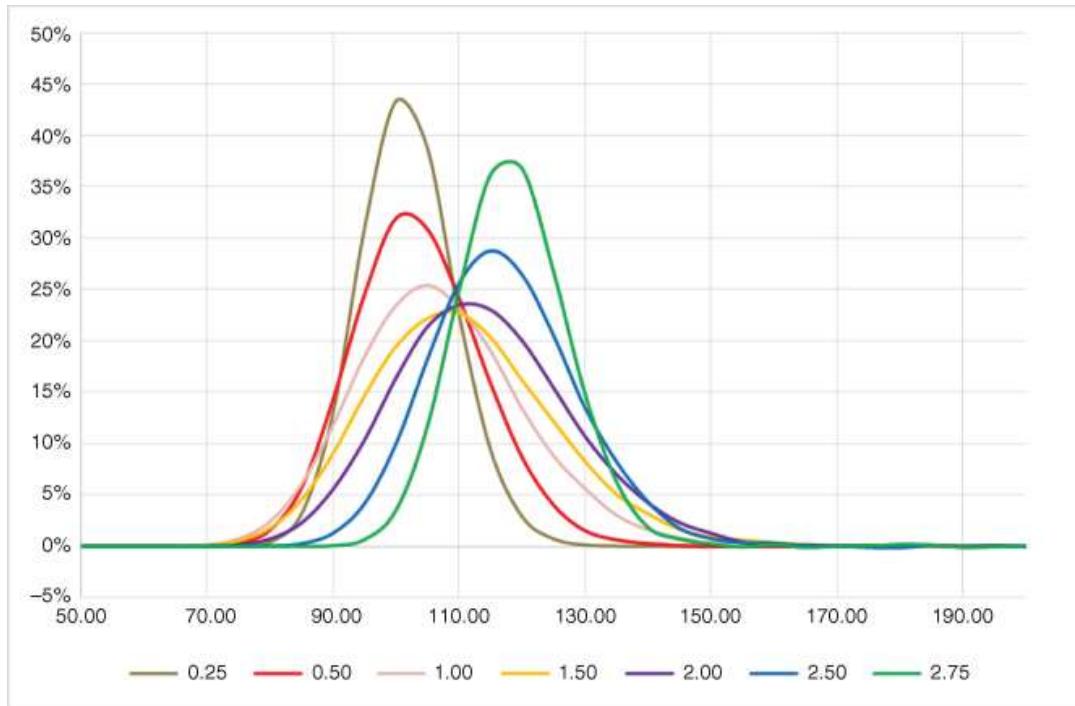




Slow financial software has the unfortunate consequence that users are encouraged to settle for approximate, inaccurate risk computed in reasonable time. With AAD, we produce the accurate charts on page 425 in half a second. We could run 1 M path to produce even better, smoother risk reports in a second.

Interestingly perhaps, one second is approximately the time it takes to produce a similar risk report with comparable accuracy, with FDM in place of Monte-Carlo, on a single core with finite differences. The combination of AAD and parallelism roughly brings FDM speed to Monte-Carlo simulations.

With 3M paths, we produce the much smoother charts below in around three seconds (that is three seconds *each*, although we produce them both in around three seconds in [Chapter 14](#), and three seconds is with the improvements of [Chapter 15](#)):



It would have taken at least half an hour with smart, multi-threaded bumping, at least four hours in serial mode. For all intents and purposes, to obtain such precise microbuckets with Monte-Carlo simulations is not feasible without AAD. AAD opens new research possibilities into financial risk management. One example is that it produces accurate microbuckets in a few seconds. For other, perhaps unex-

pected applications, see, for example, Reghai's [97] and Andreasen's [82].

We reiterate that the main purpose of derivatives systems is not to produce sensitivities to model parameters, but risks to market instruments. For this reason, this comparison between AAD and bumping is somewhat invidious. With bumping, sensitivities to hedge instruments could be produced faster by directly bumping a smaller number of market prices and repeating the whole valuation process, including the calibration of the model parameters. However, to produce market risks in this manner could lead to instabilities from differentiating through calibration. Further, the sensitivity to model parameters would not be computed, and that is valuable information in its own right. Finally, it would not offer the flexibility to compute risks with different instruments than those used to calibrate model parameters.

We have seen that AAD produces sensitivities to a large number of model parameters extremely fast. With the concern for costs of computing a large number of sensitivities out of the way, the better strategy for producing risks is to first compute sensitivities to model parameters, as explained in this chapter, and next turn them into sensitivities to hedge instruments, as explained in the next chapter, where we calibrate Dupire's model to market data and produce risks, not to local volatilities, but to implied volatilities, which directly reflect the market prices of European options.

NOTES

¹ Or we could throw an exception, or even implement a static assert that would result in a compilation error when code attempts to call the method on a `Model<T>` other than `Model <Number>`, although we implement the simple solution here.

² The parameters' adjoint accumulated the *sum* of contributions over all the paths; we must divide by N to obtain the sensitivities of values,

which are averages of payoffs across paths.

3 One risk report for each payoff, as opposed to one risk report for an aggregate—we produce itemized risk reports with AAD in [Chapter 14](#).

4 To be fair, with a 3y maturity, we “only” have 1,080 active local volatilities, so we would be done in an hour. More generally, our bump risk implementation is particularly dumb, only for reference against the AAD code. Before AAD made its way to finance, derivatives systems typically implemented a “smart bump” risk, a conscientious, invasive work on models to recalculate as little as possible in the bumped scenarios. We don't do that here because it is not the subject, but we acknowledge that a smart bump risk could complete in this example in around 20 minutes. With multi-threading over the 8 cores of our iMac Pro, it could complete in less than 3 minutes. As a comparison, AAD risk computes in around 8 seconds with the current implementation, 5 seconds with the improvements of [Chapter 15](#), and *half a second* with multi-threading. This is still at least 300 times faster than the smartest possible implementation with bumps.

5 Products are unaffected because they are not supposed to store Numbers in their state, although this is not enforced, so we must be very careful there.

6 We used a `vector<int>` in place of a `vector<bool>`. For legacy reasons, `vector<bool>` is not thread safe in standard C++ for concurrent writes into separate entries, and should be avoided.

7 This code is not exception safe. If an exception occurs here, the main thread's tape is never reset. Ideally, we would code some RAI^I logic to reset the main thread's tape no matter what.
