

How is OpenMP used

OpenMP is usually used to parallelize loops.

- Find your most time consuming loops.
- Split them up between threads.

Sequential program

```
void main()
```

```
{
```

```
int i, k, N = 1000;
double A[N], B[N], C[N];
for (i=0; i<N; i++) {
    A[i] = B[i] + k*C[i]
}
```

```
}
```

(Continued on next page)

Parallel Program

```
#include "omp.h"
```

```
void main()
```

```
{
```

```
int i, k, N = 1000;
double A[N], B[N], C[N];
#pragma omp parallel for
for (int i=0; i<N; i++) {
    A[i] = B[i] + k*C[i]
}
```

```
}
```

(Continued on next page)

Single Program Multiple Data (SPMD)

OpenMP Constructs

1) Parallel Regions

```
#pragma omp parallel
```

2) Work Sharing

```
#pragma omp for #pragma omp sections
```

3) Data Environment

```
#pragma #pragma omp barrier
```

5) Synchronization

3) Runtime functions/environment variables

```
int my_thread_id = omp_get_num_threads();
```

```
omp_set_num_threads(8);
```

4) Data Environments

```
#pragma omp parallel shared / private
```

When you enter a `#pragma omp parallel` region there is a team of threads. Variables in program can be visible to all threads (`shared`) or each thread can get its own copy (`private`). Choosing right scope is the core of correct parallel code.

- Variables declared outside a parallel region are shared by default (unless you specify otherwise)
 - Local automatic variables declared inside the parallel region are private by default
- Clauses: `shared(…)`, `private(…)`, `firstprivate(…)`, ^{Date}`lastprivate(…)`, `reduction(…)`, `decrements(…)`.

Example 1: RACE condition

```
#include <iostream>
#include <omp.h>
int main() {
    int counter = 0;           // shared by default
    #pragma omp parallel           // multiple threads execute this block
    {
        // race: multiple threads increment the same 'counter'
        counter++;               // concurrently
    }
    cout << "counter = " << counter << endl;
    return 0;
}
```

You'll often see counter < number_of_threads because of races.

Fix A (Atomic)

```
#pragma omp parallel
{
    #pragma omp atomic
    counter++;
}
```

Fix B (Reduction)

```
int N = omp_get_max_threads();
int counter = 0;
#pragma omp parallel for reduction
for (int i=0; i<N; i++) counter += counter;
// reduction (+:counter) gives each thread a partial sum and combines them safely at the end (fast & lock-free).
```

private vs firstprivate

`private(var)` → each thread has an uninitialized private copy.
`firstprivate(var)` → each thread's private copy is initialized from the original variable.

Talkpoint: `private` is used when threads should compute their own independent value; `firstprivate` when they need the original initialization.

Last private (retrieve last iteration value) → when you want the final iteration's private value copied back to original value.

`default(none) + explicit clauses`

Talkpoint: This forces you to declare the scope of each variable (good practice!) `default(none)` forces explicit thinking about each variable's scope - reduces mistakes!

Short notes on other constructs

- Reduction (t : var) - good for sums/products: much faster than locking.
- Thread private - makes a global variable private per thread. (advances, less common)
- Always prefers local private variables where possible, minimize shared writable state.

What is OpenMP

Open source Multiprocessing

Multiprocessing \rightarrow Multithreading formed by many companies like AMD, Intel, IBM, NVidia, Oracle etc.

OpenMP is available for languages C++, C & Fortran.

Released in 1997 for fortran language. In next year it released for C/C++ language.

Each process starts with one main thread. This thread is called master thread in OpenMP. For a particular block of code, we create multiple threads along with master thread are called slave threads.

OpenMP is called Fork-Join model. It is because, all slave threads after execution get joined to the master thread.

i.e. Process starts with single master thread & ends with single master thread.

For C++ we have to include "omp.h" header file.

↗ fork because creating multiple threads.
Join because all slave threads are giving result to master thread.

OpenMP Fork and Join Model.

printf("program begin\n")

serial

N = 1000;

#pragma omp parallel for

for (i=0; i < N; i++)

$$A[i] = B[i] + C[i];$$

Parallel

M = 500;

Serial

#pragma omp parallel for

for (j=0; j < M; j++)

$$P[j] = Q[j] - S[j];$$

Parallel

printf("program done.\n")

Serial

OpenMP is called as Shared Memory model as OpenMP is used to create multiple threads & these multiple threads share the memory of main process.

A simple example to create threads.

```
#pragma omp parallel
```

```
{
```

```
    printf ("Hello World")
```

By default, number of threads is equal to no. of processor cores. (In i3, i5, it there are generally 4 cores so 4 threads)

To create required no. of threads?

```
#pragma omp parallel num-threads(7)
```

```
{
```

```
    printf ("Hello World")
```

17 threads
will be created.

OpenMP program for Array Addition.

```
#include <stdio.h>
```

```
#include <omp.h>
```

```
void main()
```

```
int a[5] = {1, 2, 3, 4, 5};
```

```
int b[5] = {6, 7, 8, 9, 10};
```

```
int tid;
```

```
#pragma omp parallel num-threads (5)
```

```
{
```

```
    tid = omp_get_thread_num();
```

```
    c[tid] = a[tid] + b[tid]
```

```
    printf ("c[%d] = %d\n", tid, c[tid]);
```

```
}
```

compile cmd -

gcc -fopenmp arrayaddition.c

Outputs of array sum program

$c[1] = 9$
 $c[2] = 11$
 $c[3] = 13$
 $c[4] = 15$
 $c[0] = 7$

$c[2] = 11$
 $c[0] = 7$
 $c[4] = 15$
 $c[1] = 9$
 $c[3] = 13$

$c[1] = 9$
 $c[0] = 7$
 $c[3] = 13$
 $c[4] = 15$
 $c[2] = 11$

We are getting result in different order because all threads are running simultaneously so whichever thread completes its computation first that will get printed first.

We can create multiple for "for" block.

pragma omp parallel for num_threads (5)

```
for (i=0; i<5; i++)
```

```
{ printf("Hello World"); }
```

Allocating different work to different thread

In openMP, we can allocate different work to different thread by using "sections".

pragma omp parallel sections num_thread (3)

pragma omp section

```
printf("Hello World One")
```

pragma omp section

```
printf("Hello World Two")
```

pragma omp section

```
printf("Hello World Three")
```

Synchronizing Threads

In OpenMP, we are creating multiple threads. By default, variables are shared by all the threads. A block of code where resources (variables) are shared by multiple threads is called "critical section".

There is race condition i.e. competition among threads to use these resources (variables).

Can avoid such race condition by using processor directive.

"#pragma omp critical"

E.g.

pragma omp parallel num_threads(5)

{

pragma omp critical

{

$x = x + 1;$

}

As we are mentioning "# pragma omp critical", only one thread can do the operation of " $x = x + 1$ " at one instant.

Synchronization

A **barrier** makes all threads in the team wait until every thread reaches that barrier. After the barrier, memory updates made before the barrier are guaranteed to be visible to other threads.

(Barrier implies a flush/consistency point).

Key points

- + barrier is a collective synchronization - every thread must reach it. If some thread doesn't, your program deadlocks.
- Many OpenMP constructs insert implicit barriers (e.g. #pragma omp for has a barrier at the end, #pragma omp parallel has a barrier at the end of region unless nowait is used in some constructs.)
- Use barrier when you need an explicit phase separation e.g. Phase-1: all threads fill part of an array.
Phase-2: all threads read the full array.

Implicit barriers vs explicit barriers

Common constructs include implicit barriers:

- #pragma omp for - implicit barriers at the end (unless nowait)
- #pragma omp single - implicit barriers at end (unless nowait)
- #pragma omp sections - implicit barriers.

So you often don't need an explicit #pragma omp barrier.
Use barriers when the structure of your code requires it,

pragma omp parallel

{ if (get_omp_get_thread_num() == 0) {

pragma omp barrier // WRONG: only thread 0 hits barrier
other threads don't → deadlock.

}

All threads must reach the barrier, placing it inside a branch taken by only some threads is illegal.

pragma omp critical - serializes a block (only one thread executes at a time). Useful for small updates (printing, pushing to shared container)

pragma omp atomic - lightweight atomic update for simple operations (e.g. $x += v$)

pragma omp taskwait - waits for child tasks of the concurrent task (not the same as barrier)

flush - force memory visibility without synchronization.
(correctly needed explicitly, barriers implies flush.)

Talk point: prefer default (none), prefer reduction for aggregates, avoid unnecessary barriers - "they cost time"

Structured Block.

Most OpenMP constructs apply to structured blocks.

- Structured block: one point of entry at the top & one point of exit at bottom.

The only "branches" allowed are STOP statements if **exit()** in C/C++.

A structured block

```
#pragma omp parallel
{
```

more: do_big_job(id);

if (++count > 1) goto more;

pointf ("All done\n");

Not a structured block.

```
if (count == 1) goto more;
```

```
#pragma omp parallel
{
```

more: do_big_job (id);

if (++count > 1) goto done;

done: if (!really_done()) goto more;

In C++/C a block is a single statement or a group of stuff between brackets **{ }** .

Book Introduction to parallel Computing OpenMP

The whole point of parallel processing = faster programs (performance) and better cost-performance ratio.

Some programs are so demanding (e.g., weather simulation, AI-training, scientific modeling) that even the fastest single CPU can't finish them in reasonable time.

Solution: run parts of the program at the same time on multiple processors.

Why lost performance

Some programs are so demanding (e.g. weather) even if you can run your program on very high-end CPU, you might get the same (or better) performance by using several cheap processors together.
→ Parallel processing can give more speed with less money.

Rewards

Raw performance → finish computation faster.

Price performance → same or better using cheaper hardware.

Costs

Increased design complexity (need to think how tasks are split among processors)

Increased programming complexity (writing correct & efficient parallel code is harder than writing serial code)

(you have to avoid race conditions, handle synchronization, & manage data-sharing)

Why OpenMP

- Writing parallel programs is complex.

- OpenMP's goal: make it easier to add parallelism to existing programs
it does this by:

- Providing a standard API (work across platforms)

- Using compiler directives / pragmas so you can add parallelism without rewriting whole program.

Shared-Memory Multiprocessors

multiple CPUs / cores that all share the same main memory.

Eg: My laptop with 4 cores - all cores access the same RAM.

This is different from distributed-memory systems which usually require MPI.

(e.g. a cluster with many separate computers)

Fallpoint:

- Parallel processing is about speed & cost-efficiency, running a program across many processors instead one.
- The trade off is increased coding complexity, that's where OpenMP comes in, it makes it easier to parallelize existing code, especially for shared memory systems.

The NWS weather model share shows how OpenMP can take a computation from impossible to practical, turning days of simulation into hours, by scaling across many processors.

This speedup directly translates to timely, accurate forecasts. It's not just about speed, it's also about speed it's also about cost efficiency since multiple commodity CPUs can outperform a single expensive supercomputer.

X-axis = Number of processors used

Y-axis = Speedup

(how many times faster the program runs compared to a single program).

• Near linear scaling at first

• Doubling processors \approx doubling speed-up to a point

• As no. of processors increases:

Speedup still improves but with diminishing returns (due to overheads & performance limits - Amdahl's law)

1 CPU = baseline speed 8 CPU's = $\sim 60 \times$ faster 128 CPU's = close to $80 \times$ faster.

Even fastest single processors can't match this combined performance.

Cray YMP supercomputer \rightarrow 512 MFLOPS per CPU (fast at the Hm) still not enough for real-time forecasting.

Origin 2000 with many cheaper processors beats it in total throughput. This is more cost-effective than building a single huge super-fast processor.

CFD (Computational Fluid Dynamics)

Solving equations that model fluid flow (air, water, etc.)

(computation)-heavy Aerospace (aircraft design) Alternative/aerodynamics Turbine design

NAS Parallel Benchmarks, APPFLU

APPFLU is one of the NAS Parallel Benchmarks (NPB) used to measure how well a program scales on multiple processors.

In this benchmark, OpenMP's speedup is almost identical to MPI's up to the tested processor count.

This proves that OpenMP's speedup is almost performance comparable to MPI for shared-memory systems, without the extra complexity of message passing.

Real World point,

If your problem fits into shared memory (like a single powerful server), OpenMP can give you near-MPI performance with less development overhead.

Parallel computing enables huge performance boosts but only if,

- The problem is parallelizable (large independent chunks of work).
- Overhead (thread management, memory creation) is small compared to actual computation.

Example application - Crash simulation (full scale car crash analysis)

- Extremely expensive if done physically.
- Parallel computing makes high-detail simulations economically feasible.
- Used in auto industry & govt. safety agencies.

Performance of a leading crash code

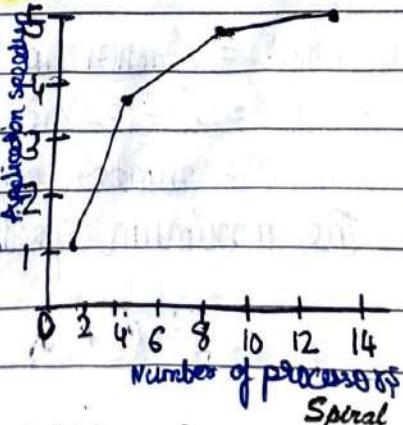
The curve rises quickly at first, but flattens as processor increase - classic diminishing returns.

Early gains: going from 1 → 4 → 8 processors gives large speedups.

Later gains: adding processors beyond that yields smaller improvement due to:

- Parallel overhead

- Non-parallelizable code sections (Amdahl's Law).



Talkpoint: Even with diminishing returns parallel computing greatly reduces simulation time for such complex engineering problems, making them feasible for regular use in product design cycles.

Benefits of OpenMP

Scalability

APPF11's near identical OpenMP/MPI scaling shows that OpenMP is not inherently "slower" than MPI in shared memory systems, the bottleneck is usually memory bandwidth or algorithmic limits.

Incremental parallelism

OpenMP lets you parallelize your code in stages, ideal when you can't dedicate months to a full rewrite.

Economic feasibility

Real crash tests or large fluids simulations are very expensive in the real world. Simulating them on parallel systems is far cheaper, faster and safer.

Amdahl's Law connection

Speedup is limited by the serial fraction of code. For ex: if 10% of the code is serial, the max speedup is 10x no matter how many processors you throw at it.

Amdahl Law Formula (additional info from internet)

f_s = fraction of time spent in serial code

$p = 1 - f_s$ = fraction that can be parallelized

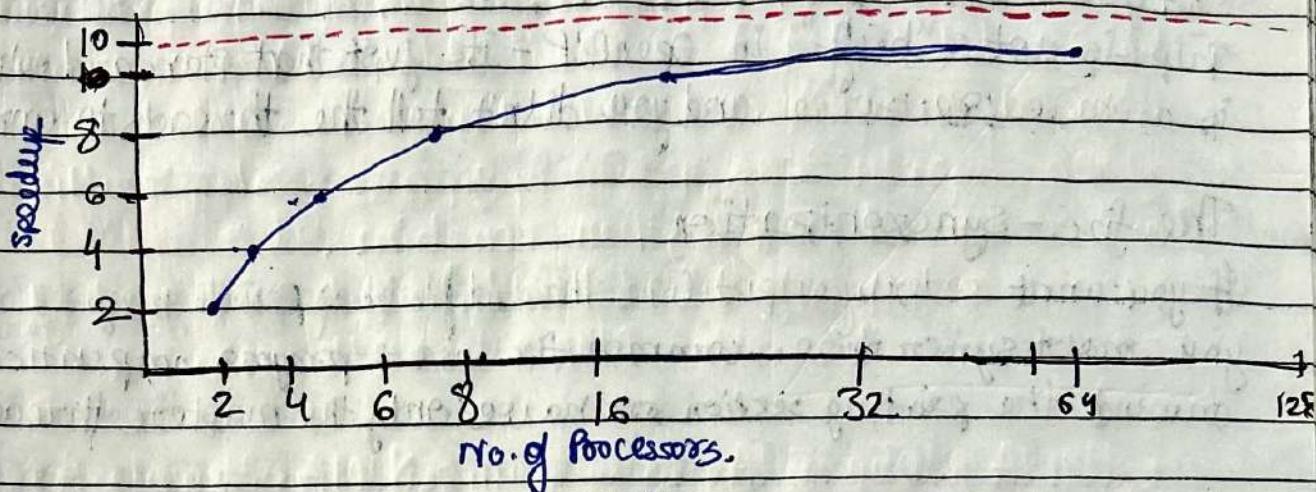
N = number of processors

The maximum speedup is:

$$S(N) = \frac{1}{f_s + \frac{p}{N}}$$

MPI = group of people in separate rooms talking via phone calls
 OpenMP = group of people at the same table sharing one notebook (message passing)
 For APPFLU, both groups work equally fast because the problem is well suited to shared memory (shared memory).

Amdahl's Law 10% serial fraction



- > With 10% serial code even infinite processors can't go beyond 9x speedup.
- > The curve rises quickly at first ($1 \rightarrow 8$ processors) but flattens - diminishing returns

OpenMP basics

- > Controlled via directives (#pragma omp ...), runtime function `omp_get_thread_num()` & environment variables (`OMP_NUM_THREADS`)
- > Use fork-join model: master thread \rightarrow forklteam \rightarrow parallel work \downarrow joinback.

Directives # pragma omp parallel { /* work */ }

Key functions `omp_get_thread_num()` \rightarrow returns thread id.

`omp_get_num_threads()` \rightarrow returns total threads in the team.

`omp_set_num_threads(n)` \rightarrow set threads programmatically.

Environment variable

`OMP_NUM_THREADS` \rightarrow sets default threads without changing code.

Ex: set `OMP_NUM_THREADS=4` (Windows)

Output can get jumbled.

- > When you run # pragma omp parallel & each thread prints its ID, all threads share the same output screen.

- > They run independently and don't wait for each other.

That means one thread might print while another immediately starts printing before the first thread finishes - so the text appears interleaved.

This is not a "bug" in OpenMP - it's just that standard output is a shared resource, and you didn't tell the threads to coordinate.

The fix - synchronization.

If you want orderly output (one thread finishes printing before another), you must synchronize. common fix: use #pragma omp critical around the printing section or have only the master thread print results.

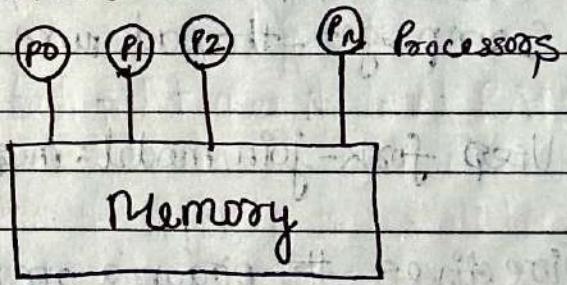
Synchronization in OpenMP is broad topic & includes:

Borders (wait until all threads arrive) Critical Sections
Atomic Updates Ordered execution.

THE OPENMP PARALLEL COMPUTER

Several CPUs / cores connected to one, globally addressable memory (all threads see the same addresses)

WHAT IT LOOKS LIKE



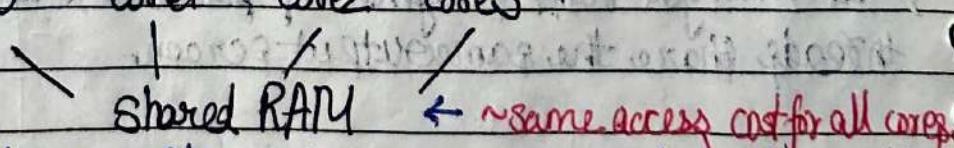
Shared Memory multiprocessors

Two common flavors

- 1) UMA (SMP)
- 2) ccNUMA

UMA (a.k.a. SMP) - uniform memory access

All cores see one shared RAM and the same latency/bandwidth to any address.



- Biggest win case for OpenMP: performance mostly depends on total memory bandwidth & algorithm.

ccNUMA - cache-coherent Non-Uniform Memory Access.

- One address space (still "shared memory") but RAM is split into per-socket banks.

Socket 0: cores + local RAM 0

Socket 1: cores + local RAM 1

Interconnect (keeps caches coherent)

Core on S0 \rightarrow RAM 0 = fast (local)

Core on S0 \rightarrow RAM 1 = slower (remote)

- Good OpenMP code must place data near threads that use it (first-touch). Otherwise many accesses become remote \rightarrow slower.

"first-touch" matters on ccNUMA

- The OS maps physical pages for a big array when they're first written.
- If thread 0 alone initializes a 20 GB array, every page lands in socket 0's memory.
- Later, when all sockets read that array, sockets 1...N do remote access \rightarrow big slowdown.

Fix: initialize in parallel so each thread "touches (writes) its own chunk first, pages locally for that thread's socket."

Program which shows bad (serial, first touch) vs good (parallel, first touch) and measures difference in the compute phase with result uploaded in later file

Talkpoint:

UMA: Any core \rightarrow any memory with roughly equal cost. OpenMP is straight forward
memory placement rarely bites you!

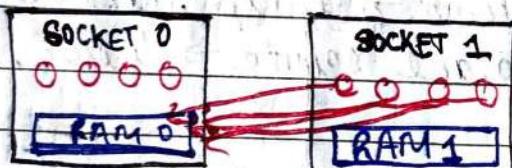
ccNUMA: Still one address space, but placement matters. Use parallel first-touch, stick to static chunking when data is block-partitioned, & consider OMP_PROC_BIND to keep threads near their data serial

Rules of Thumb

- Initialize big arrays in parallel (first-touch)
- Keep each thread working mostly on its own chunk. Avoid straying far from home.
- Use schedule (static) for predictable locality; use dynamic^(per-row) only if you truly need load balance.
- For I/O or reductions, minimize cross-socket traffic; do per-thread partials, then combine.
- If performance drops as you add threads past one socket, code count, it's often NUMA bandwidth / latency.

ccNUMA Memory Access

- Bad First Touch

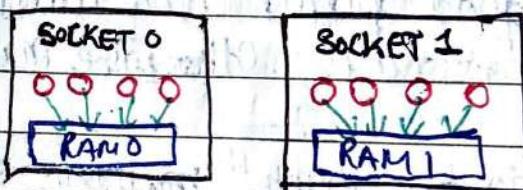


Bad first touch

Socket 1 does access RAM 0
(remote)

ccNUMA Memory Access

- Good First Touch



Good first touch

Each socket uses its local RAM

Pitfalls to know: (especially on NUMA)

- Memory bandwidth can bottleneck at high thread counts.
- NUMA placement matters: if one initializes a huge array, all its pages may land on one socket's memory → other sockets read remotely (slower).

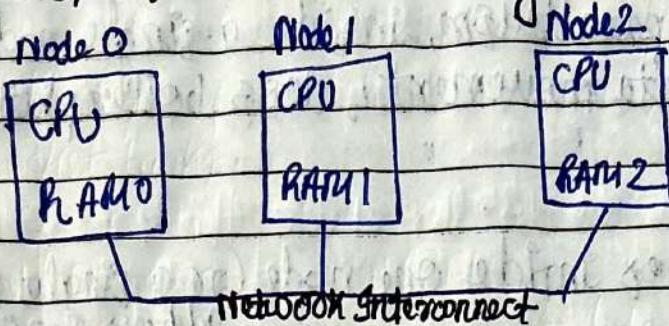
Pro tip: "first-touch" initialization

Initialize big arrays in parallel so pages are near the threads that will use them.

HPC → high Performance Computing .

Date _____

What is distributed - memory (message-passing) architecture



- Each node/process has its own private memory. There is no single globally addressable RAM.
- If process A needs data that lives in process B's, A must ask B i.e. send a message.
- Examples of systems: clusters, networked workstations, IBM SP2, modern HPC clusters.

Implication: Programmers must design where data lives & explicitly move (send/receive). This gives scalability (can use many nodes) but increases complexity.

Trade-offs: MPI (distributed) Vs OpenMP (shared)

Pros (MPI / distributed Memory)

- Scales across many nodes - can use huge aggregate memory
- Good when dataset exceeds single-node RAM.

Cons (MPI)

- You must explicitly send/receive data → more code, possible bugs.
- Network latency/bandwidth matters - communication cost can dominate if granularity is small.

For (OpenMP / shared memory)

- Easier to add parallelism inside a single node via pragmas.
- No explicit data movement, less boilerplate.

(Con_ns (OpenMP))

- Limited to cores inside one node (or a single shared address space).
- Must manage NUMA Locality (first-touch covered).

Rule:

- If program will run on a cluster, across nodes → use MPI (or hybrid)
- If you only target a single multi-core machine → OpenMP is simpler.

Hybrid MPI between nodes + OpenMP inside a Node
 (common in High Performance Computing (HPC))

Why hybrid?

- Use MPI to partition global domain across nodes.
- Within each node, use OpenMP threads to use all cores & reduce MPI ranks (less MPI communication).

Important: MPI and threads interact via MPI_THREAD_*

Simple hybrid code of get product using this on Latex file made

- MPI_THREAD_* thread lets you control thread level support. For many hybrid programs MPI_THREAD_FUNNELED or MPI_THREAD_SERIALIZED suffice. MPI_THREAD_MULTIPLE allows concurrent MPI calls from multiple threads but depends on MPI implementation performance.
- Hybrid reduces number of MPI ranks (lower memory per rank, fewer messages) but increases complexity (thread+message interactions).

WHY OPENMP

- Shared memory multiprocessor systems have become much more common, powerful, and affordable in recent years.
- But, existing programming models were often designed for distributed memory systems, not shared memory.
- This caused a gap b/w hardware capabilities.

Goal of OpenMP

Provide a standard & portable API for programming on shared memory parallel systems.

Hardware Trends:

In recent years

- Quantity of processors in shared memory systems has increased rapidly.
- Scalability (ability to handle more processes effectively) is improving.
- First multiprocessor systems had only 2 processors → quickly moved to 4 & beyond.
- Medium to high-end systems now support 2 to 32 RISC processors (examples: Compaq AlphaServer, SunEnterprise).
- Falling prices → high-end systems used in smaller organisations too.

Software Gap

- Different vendors had their own APIs for shared memory programming.
- Distributed memory models (like MPI) became common but they're harder to use for shared memory because they need explicit message passing.
- Without a standard API, developers often:
 - Used vendor-specific solutions.
 - Relied on message passing even when all processors share memory, which is inefficient.
- OpenMP emerged as the standard, making shared parallel programming more portable and easier.

Shared vs Distributed Memory Programming

Feature	Shared Memory	Distributed Memory
Ease of parallelizing an app.	Easy to do; less effort	Harder; often all or nothing work.
Scaling to many processors	Limited by hardware vendor's design (e.g., up to few dozen CPUs)	Scales to large number of processors
Extra code needed	Minimal, compiler helps	More extra code, must explicitly send/receive message
Impact on code quality	Slight complexity increase; debugging easier.	Code more complex; harder to debug.
Portability	Requires shared memory	Runs anywhere, even hardware & compiler support without shared memory.

Why Not Just Use MPI for Shared Memory?

- MPI is standardized, widely available & portable but
- Requires explicit message passing.
- Programmer must manually partition data.
- Minimal incremental parallelization - you often have to rewrite a lot of program.
- Many modern processors already share memory cache, so message passing can be unnecessary overhead.

Other Alternatives

- Threads: Good for shared memory, but
- Lower-level → more complex & awkward for large-scale scientific computing
- Portability issues b/w platforms.
- Proprietary vendor APIs: Not portable, locks you to one system.
- New languages (e.g. Java): Adoption is slow, porting existing codebases is hard.

Why OpenMP Wins

- High level, easy to learn (complex directives like #pragma omp for)
- Portable across vendors and systems.
- Allocates incremental parallelism:
You can parallelize parts of your code gradually.
- Works well with existing serial code - don't have to rewrite everything.

History of OpenMP

- OpenMP: Standard API for shared memory parallel programming (est. 1997)
- Predecessors: XJHB, HPE - failed due to hardware variety + portability issues
- Pre-openMP Problem: Hardware capable, but no portable model.
- Motivation: Developer demand for a single standard across vendors.
- Development: Led by industry group (hardware + software vendors)
- ARB: Governing body for OpenMP standard.
- First Standard: Released July 1997.

ARB → Architectural Review Board (ARB).

Chapter 2 Getting Started with OpenMP

Date _____

A parallel programming language must provide three things:

1. Specifying parallel execution.

→ How to tell computer "this program to run parallel"

2. Communication between threads.

→ Threads may need to share or pass values

• In OpenMP (shared memory), communication = just reading/writing shared variables

• OpenMP uses #pragma omp parallel for to say split this loop among threads.

• In MPI (distributed memory) communication = explicit send/receive messages.

3. Synchronization between threads

→ We must coordinate so that don't step on each other, too.

Examples in OpenMP

#pragma omp barrier (all wait here until everyone arrives)

#pragma omp critical (only one thread enters this block at a time)

Different approaches to parallelism

The book mentions 3 approaches that languages/libraries take care:

(a) New language constructs

• Example: Fortran 95 has forall keyword,

• The language itself add syntax for parallelism

(b) Directive-based (OpenMP's way)

- We use compiler directives (#pragma omp)

• Directives looks like comments, so if the compiler doesn't support OpenMP, it just ignores them → code runs as normal serial
Advantage: one code works for both serial and parallel.

(c) Library calls (MPI, threads)

- No new language syntax, but special functions calls (MPI_send, MPI_recv)
- More manual, but also more flexible.

(pthread-create)

Two Big benefits of OpenMP directive

1. Same code works everywhere.

- With OpenMP pragmas → if OpenMP is off, they're ignored, program runs serially.
- So you don't have to maintain two versions of the code.

2. Incremental parallelism:

- can start with a working sequential program.
- Then gradually add #pragma omp directives to parallelize more parts.
- Safer and easier to debug.

At the highest level, OpenMP is:

- A set of compiler directives (#pragma omp..)
- Runtime library routines (like omp_get_thread_num())
- Environment variables (like OMP_NUM_THREADS)

→ parallel languages must handle execution, communication, synchronization.

- Approaches: task constructs (Fortran forall), directive-based (OpenMP), libraries (MPI, threads)

- OpenMP is directive based, so:

1. Same code runs serially if pragmas are ignored.

2. Parallelism can be added incrementally.

In C/C++, we use #pragma omp pragmas + OpenMP runtime functions

OpenMP Directives (`#pragma omp`)

In C++, OpenMP relies on the `#pragma` keyword.
`#pragma` = tells the compiler to do something special, but if the compiler doesn't understand it, it just ignores it.

- If you compile with OpenMP (`g++ -fopenmp file.cpp`)
 - multiple threads run

- If you compile without OpenMP (`g++ file.cpp`)
 - The directive is ignored → only one thread points to the code.

Talk point: You don't need two versions of code. One program working for both serial and parallel.

Conditional Compilation with ~~Open~~ OPENMP

Sometimes you want your program to behave differently depending on whether OpenMP is enabled or not.

```
#include <iostream>
```

```
using namespace std;
```

```
int main() {
```

```
#ifdef OPENMP
```

```
cout << "Compiled with OpenMP!" << endl;
```

```
#else
```

```
cout << "Compiled without OpenMP!" << endl;
```

- If compiled with `-fopenmp` → prints the OpenMP message

- Without `-fopenmp` → prints sequential message

This is super useful for debugging, or for writing portable software that should run anywhere.

Combining Directives & Normal Code

```
#include <iostream>           // with OpenMP → loop is split across
#include <omp.h>             threads, then results are combined
using namespace std;
int main() {                  // without OpenMP → loop runs normally still
    int sum = 0;               correct.
    #pragma omp parallel for reduction(+:sum)
    for (int i = 1; i <= 100; i++) {
        sum += i;              // Each thread works on a piece of loop
    }
    cout << "Sum from 1 to 100 = " << sum << endl;
}
```

- If you're working on a large scientific simulation
 - You want it to run fast on a supercomputer (many threads)
 - But you also want to test it on your laptop (single processor)

With OpenMP directives:

- On laptop → directives are ignored → program still runs
- On supercomputer → directives activate → massive speedup.

- If OpenMP not enabled → treated as comment → safe fallback
- -OPENMP macro → lets you check if code is compiled with OpenMP
- Benefit: single codebase, portable, incremental development

Example: parallel loops, conditional messages.

Using `# pragma omp critical` Using reduction (`# total-items`)

Phase 1: `total-Iters = 0`

Thread 1:

compute depth (1)

--- Enter critical ---

`total-Iters = total-Iters + depth (1)`

--- Leaves critical ---

Thread 2:

compute depth (2)

--- Enter critical ---

`total-Iters = total-Iters + depth (2)`

--- Leaves critical ---

Thread 3: --- Leaves critical ---

Thread 1: `local-sum = 0`

`local-sum += depth (1)`

Thread 2: `local-sum = 0`

`local-sum += depth (2)`

Thread 3: ...

(Each thread keeps its own private copy, no racing!)

critical = safe, but slower (thread at the end)

total-Iters = sum of all local sum

total-Iters = depth (3) + depth (2) + depth (1)

Critical:

Only one thread at a time updates total-hits. Others must wait (causing delays).

Safe, but slower (threads wait in line.)

Reduction:

No waiting. Each thread works independently. At the end, OpenMP combines results automatically.

Safer + faster (no waiting, automatic combination).

(conditional compilation - done)

Parallel struct control structures - the "fork-join" model

Communication & the data environment (shared vs private vs reduction)

OpenMP is shared-memory: all threads can see the same global / heap variables.

But each thread has its own stack (automatic variable inside a function are private to thread unless shared by clause).

A. shared - one copy visible to all threads.

int hits = 0;

#pragma omp parallel shared (hits)

// WARNING: writing 'hits++' from many threads races!

B **private** - each thread gets its own copy (uninitialized)

```
#pragma omp parallel private(tmp)
```

```
int tmp;
```

// private space

```
tmp = compute_local();
```

// tmp is invisible to other threads

{}

C **firstprivate** - private but initialized from the original value-

```
int seed = 123;
```

```
#pragma omp parallel firstprivate(seed)
```

// each thread starts with seed = 123, they can modify its own copy

{}

D **lastprivate** - after a worksharing loop, keep the last iteration's value

```
int last_i = -1;
```

```
#pragma omp parallel for lastprivate(last_i)
```

```
for(int i=0; i < N; ++i){
```

// ...

```
last_i = i;
```

{}

// here last_i = N-1;

E **reduction** - private during the region, combined at end.

```
double sum = 0.0;
```

```
#pragma omp parallel for reduction( + : sum )
```

```
for (int j=0; j < N; ++j){
```

sum += a[j]; "each thread has a private sum _{spiral} they're added at end."

{}

Without reduction, there would be a race:

Best practice: force yourself to be explicit.

#pragma omp parallel for default(none) shared(a,N)
for (int i=0; i<N; ++i) sum += a[i];
reduction(+:sum)

default(none) makes the compiler warn if you forgot to declare scope.

Communication rule of thumb: communication is just reads/writes to shared variables but correctness/visibility is guaranteed at synchronization points (e.g., end of parallel for, barrier, critical, atomic). Don't assume another thread sees your write immediately unless there's a sync.

Black Scholes Problem.

→ An option is a contract giving the right to buy an asset at strike price K at expiry T . A European call pays max($S_T - K, 0$) at time T .

Goal: given current price S and time t , find the option price, $V(S, t)$.

Black Scholes derive a partial differential equation (PDE) for option price $V(S, t)$:

$$\frac{\partial V}{\partial t} + \frac{1}{2} \sigma^2 S^2 \frac{\partial^2 V}{\partial S^2} + \sigma S \frac{\partial V}{\partial S} - rV = 0$$

with final condition $V(S, T) = \max(S - K, 0)$.

boundary conditions

$$V(0, t) = 0$$

$$V(S, \infty) \approx S$$

S = the underlying asset (stock) price.

It's a variable that can range 0 to ∞ .

On numerical grid, this is space axis.

t = time (current time b/w now and expiry).

In the PDE form, we usually solve backward in time from the final condition $t = T$ back to $t = 0$.

$V(S, t)$ = the option price when the asset is at price S at time t .

This is the unknown we are solving for.

Ex: If today ($t=0$) the stock price is $S=70$, then $V(70, 0)$

is the fair value of option today.

At expiry ($t=T$), we already know $V(S, T) = \max(S - K, 0)$

Closed-form Black-Scholes solution:

$$C(S, t) = S \phi(d_1) - K e^{-r(T-t)} \phi(d_2)$$

where

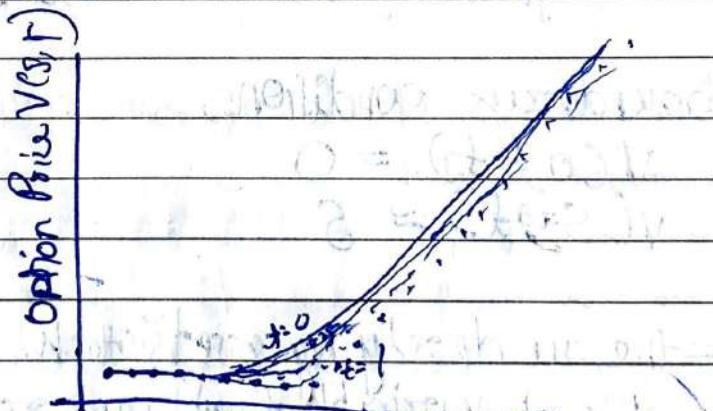
$$d_1 = \frac{\ln(S/K) + (r + \frac{1}{2}\sigma^2)(T-t)}{\sigma\sqrt{T-t}}$$

$$d_2 = d_1 - \sigma\sqrt{T-t}$$

$\sigma \rightarrow$ volatility

$r \rightarrow$ rate

$T \rightarrow$ expiry (time to expiry)



Ex:

$$K = 85$$

$$\sigma = 0.1$$

$$\sigma = 0.2$$

$$T = 1$$

Stock Price S

$$\text{Exercise } S_{\text{man}} = 3K = 255$$

$$\text{Spatial steps } \Delta M = 51 \rightarrow \Delta S = S_{\text{man}}/m = 255/51 = 5.0$$

$$\text{Time steps } N = 100 \rightarrow \Delta t = T/N = 0.01$$

For small S (far below strike) option ≈ 0

For large S option $\approx S - Ke^{-\sigma(T-t)}$

(non-linear for deep in money).

As time decreases toward maturity (in t toward 0 in our indexing), the curves move toward the payoff line.

At $t=0$ & $S=S_{\text{max}} = 255$, the analytic price $\approx 255 - Ke^{-\sigma T}$.

Compute $e^{-\sigma T} = e^{-0.1} \approx 0.9048374180$.

$$Ke^{-\sigma T} = 85 \times 0.904374180 = 76.91118053.$$

$$V(S_{\text{max}}, 0) \approx 255 - 76.91118053 = 178.0880$$

1942

Matched top range of curve,

first derivative

$$\frac{\partial V}{\partial S_j} = \frac{V_{j+1} - V_{j-1}}{2\Delta S}$$

second derivative

$$\frac{\partial^2 V}{\partial S^2} = \frac{V_{j+1} - 2V_j + V_{j-1}}{\Delta S^2}$$

$$\text{Ex: } \Delta S = 5, \Delta t = 0.01$$

$$S_{\text{center}} = 100 \cdot (j=20), t = 0.99 (n=49)$$

$$V(95) = 10.084958 \quad \frac{\partial V}{\partial S} = \frac{10}{5} = 1.0$$

$$V(100) = 15.08458$$

$$V(105) = 20.084958$$

$$\frac{\partial^2 V}{\partial S^2} = \frac{20.084958 - 2 \times 15.08458 + 10.084958}{5^2}$$

$$\Delta t \leq \frac{\Delta S^2}{\sigma^2 S_{\max}^2}$$

$$\Delta S = 5 \Rightarrow \Delta S^2 = 25.$$

$$\sigma = 0.2 \Rightarrow \sigma^2 = 0.04,$$

$$S_{\max} = 255 \Rightarrow S_{\max}^2 = 255^2 = 65025$$

$$\sigma^2 S_{\max}^2 = 0.04 \times 65025 = 2601.0$$

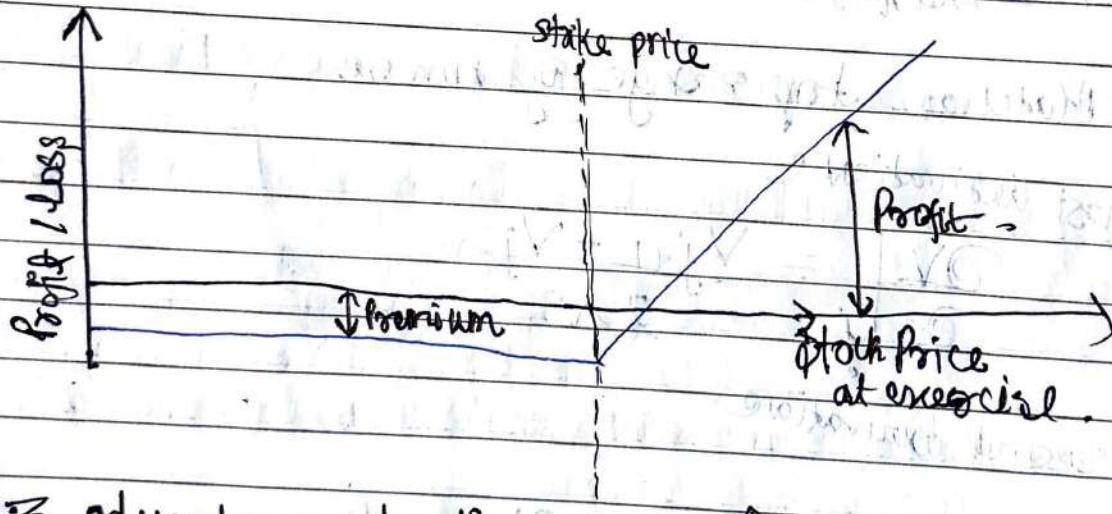
$$\Delta t \leq \frac{25}{2601.0} \approx 0.009611688$$

Compare

$$\Delta t = 0.01$$

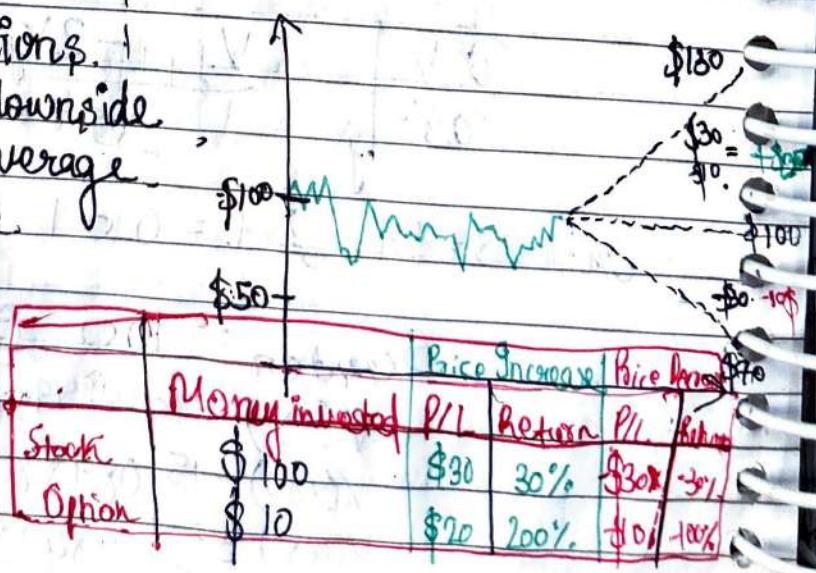
$$0.01 > 0.009611688$$

Profit and loss on an option.



3 advantages of options.

1. Limits your downside.
2. Provide leverage.
3. Use as a hedge.



~~Implicit~~

We want to solve Black Scholes PDE for $V(S, t)$

$$\frac{\partial V}{\partial t} + \frac{1}{2} \sigma^2 S^2 \frac{\partial^2 V}{\partial S^2} + rS \frac{\partial V}{\partial S} - rV = 0,$$

with final condition $V(S, T) = \max(S - K, 0)$

We discretize $S_j = j/\Delta S$ and $t^n = n\Delta t$.

Let $V_j^n \approx V(S_j, t^n)$.

spatial central difference approximations:

$$\left. \frac{\partial V}{\partial S} \right|_j \approx \frac{V_{j+1} - V_{j-1}}{2\Delta S} \quad \left. \frac{\partial^2 V}{\partial S^2} \right|_j \approx \frac{V_{j+1} - 2V_j + V_{j-1}}{\Delta S^2}$$

PDE can be written as

$$\frac{\partial V}{\partial t} = L[V] \quad \text{where} \quad L[V] = -\frac{1}{2} \sigma^2 S^2 \frac{\partial^2 V}{\partial S^2} - rS \frac{\partial V}{\partial S} + rV.$$

We moved spatial terms to RHS so we march backward in time.

Explicit (forward time)

Replace time derivative by forward in marching direction (explicit difference using known values). at the late time level:

$$\frac{V_j^n - V_j^{n+1}}{\Delta t} = L[V]^{n+1}$$

Conditionally stable - requires very small Δt , up if σ or S_{max} are large. Date _____

Or more commonly written stepping backward,

$$V_j^n = V_j^{n+1} + \Delta t [L[V]]^{n+1}.$$

(When we march from $n+1$ (new) to n (earlier time), the right hand side uses the known values V^{n+1} .)

Plugging central differences:

$$V_j^n = V_j^{n+1} + \Delta t (A_j (V_j^{n+1} - 2V_{j+1}^{n+1} + V_{j-1}^{n+1}) + B_j (V_{j+1}^{n+1} - V_{j-1}^{n+1})) + \alpha V_j^{n+1},$$

$$\text{where } A_j = \frac{1}{2} \sigma^2 S_j^2 \quad B_j = \frac{\sigma S_j}{2 \Delta t}.$$

Parallel across spatial nodes.

Implicit (backward time)

Be the spatial at the unknown new time level:

$$\frac{V_j^n - V_j^{n+1}}{\Delta t} = L[V]^n.$$

Now RHS depends on unknowns V^n at all nearby spatial points. Rearranged this yields a linear system.

$$V_j^n - \Delta t [L[V]]^n = V_j^{n+1}.$$

After assembling spatial discretization, this forms a Δt -diagonal linear system.

$$\alpha_j V_j^n + \beta_j V_j^n + \gamma_j V_{j+1}^n = \text{RHS};$$

V_j . Solving to get V^n .

- dampen high freq noise

Unconditionally stable for linear PDEs like Blackholes.
Can use large Δt .

Frank Nicolson.

Time centered avg of explicit & implicit.

$$\frac{V_j^n - V_j^{n+1}}{\Delta t} = \frac{1}{2} (L[V]^n + L[V]^{n+1}).$$

This yield a tridiagonal system too, but is second order accurate in time.

Final Black Scholes presentation.

The Black-Scholes model (1973) is a PDE Model for option pricing. Before Black-Scholes and Merton developed this approach, finding a fair price of a derivative was a "big problem". In simple terms, the model follows a geometric Brownian motion with constant volatility & interest rate, and that one can form a risk free hedge.

Under these assumptions the option price $V(s, t)$ (as a function of underlying stock price s and time t) satisfies parabolic partial differential equation. Black-Scholes derived this eq. in 1973, Merton gave parallel derivation. Black-Scholes & Merton later won the 1997 Nobel Prize for discovery.

The practical problem solved is: given the current stock price s , time to expiration T , volatility σ , and risk-free interest rate r , what should be "fair" premium (price) of an option? We consider a European call option, which gives the right but not obligation to buy the stock at a fixed price K at the expiration date T .

The Black-Scholes model provides an explicit (closed-form) answer for such options, which balances the effects of time decay, volatility and discounting.

To understand PDE,

- European call option: A contract giving its holder the right to purchase one share of underlying stock price K , but only on the expiration date (not before). The holder would exercise the option to buy cheaply at K and immediately sell for the market price, otherwise option expires worthless.
- Strike Price (K): The fixed price at which the option holder can buy (for a call) or sell (for a put) the underlying asset. For ex: if the strike price of a call is \$50, the holder can buy the stock for \$50 at expiry no matter its market value.

The strike price is a key determinant of the option intrinsic value (the immediate profit potential) and moneyness.

- Expiration Time (T): The date (often measure from "today") when the option contract expires. For a European option, exercise is allowed at $t = T$.

The time remaining until expiry is often denoted $\tau = T - t$. All else equal, a longer time to expiry typically increases options premium, because there is more time to move stock advantageously.

- Volatility (σ): A measure of how erratic the stock price movements are, typically defined as $\text{standard deviation}/\text{mean}$. High volatility means the stock price can swing widely, making the option more valuable (since large favorable moves are more likely). In the Black-Scholes model, the stock's volatility is assumed constant. Volatility is an important variable in the valuation of options.

• Risk-free interest rate: The continuously compounded rate of return of a risk-free investment (e.g. a government bond) for ex: \$1 invested today grows to $\$ \exp(rT)$ after time T . The risk-free rate enters the PDE because option pricing is done under a risk-neutral valuation. Essentially one assumes all investors accept returns on average so that expected stock growth under price measure is γ .

• European call's payoff at expiration T is $\max(S_T - K, 0)$ if the stock price S_T exceeds K , the option payoff is $S_T - K$, otherwise it pays 0. An option's premium today must reflect (under risk-neutral valuation) the expected discounted payoff under the model.

The Black-Scholes PDE

Under the standard assumptions (no dividends, long-run stock, constant γ, σ) one can show the call price $V(S, t)$ satisfies the Black-Scholes PDE:

$$\frac{\partial V}{\partial t} + \frac{1}{2} \sigma^2 S^2 \frac{\partial^2 V}{\partial S^2} + \gamma S \frac{\partial V}{\partial S} - rV = 0$$

\downarrow
represents
decay of
option value.

* Expressing fast moving
prices in terms of slower
changing factors like
implied volatility & strike
which allow for easy
comparison b/w all the different
contracts like

(Comparing Option Prices and M)
(Comparing option prices in diff stocks)

would be chaotic and meaningless as prices are influenced by so many different factors: Stock price, time to expiration, volatility & each can vary widely b/w stocks.

but when you compare implied volatilities you're speaking common
 that cut through the noise & lets you focus on what really matters
 Option Price Stock A \$10
 Option IV = 20% Stock B option price = \$3
 more expensive Option IV = 60% in terms of implied volatility.

Delta $\Delta = \frac{\partial C}{\partial S}$ how the call price changes with the change in stock price S .

Gamma $\Gamma = \frac{\partial^2 C}{\partial S^2} = \frac{\partial \Delta}{\partial S}$ rate of change of delta as stock price changes.

Theta $\Theta = \frac{\partial C}{\partial t}$ how the call price changes as time passes.

$$\Theta + \frac{1}{2} \sigma^2 S^2 \Gamma + \cancel{\rho S \Delta} - \cancel{\delta C} = 0$$

Price of stock at expiration
(Due to interest).

Discount call price for present value
(remove interest from future price)

Option Strike is \$100 & Stock Price is \$75.

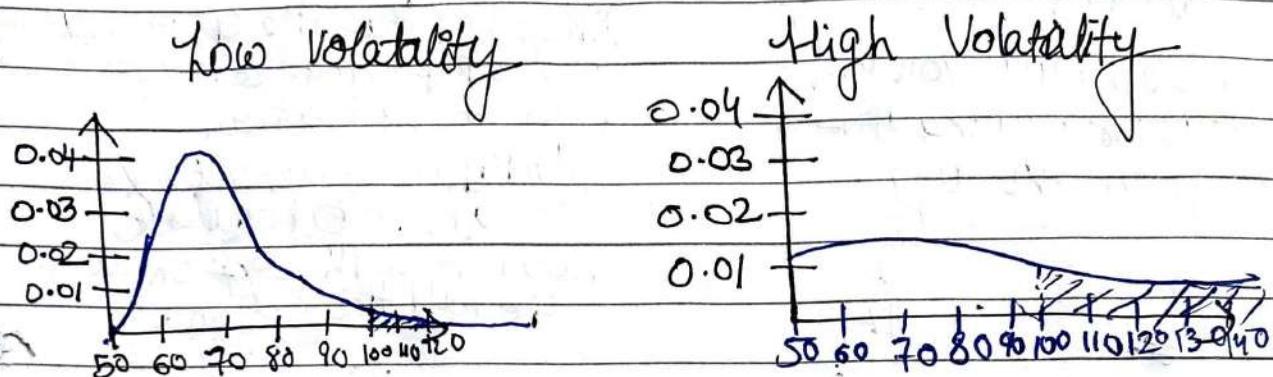
σS represent the forward price of stock this is the stock value at expiration if only interest affects its price.

negative δC discounts the call option's value to present terms meaning it removes the effects of interest from the option's future value at expiration.

Volatility (σ) is the standard deviation of the assumed log normal distribution showing how the stock prices movements affect the rate of the Delta changes volatility reflects the uncertainty or variability in the stock's price.

and the higher it is the greater the potential for options price changes this means there is a higher chance of finishing in the money (ITM) and thus higher volatility will increase options prices.

Option strike is \$100 and stock price is \$75.



We can see here that the same option with a strike of 100 and current stock price of 75 will be priced higher if volatility is higher since it has greater probability of finishing above 100.

This is represented by larger area under the curve one key assumption in the Black Scholes model is that volatility stays constant throughout options life. That's not realistic but makes maths simpler while giving a useful framework.

$$C = S N(d_1) - K e^{-rt} N(d_2)$$

$$d_1 = \frac{\ln(S/K) + (r + \sigma^2/2)t}{\sigma \sqrt{t}}$$

$$d_2 = d_1 - \sigma \sqrt{t}$$

To actually calculate these effects like Greeks and call price the equation needs to be solved. The solution has been done for us and known as the famous Black-Scholes model.

$$C = [SN(d_1)] - Ke^{-rt} N(d_2)$$

average value of
the stock above
the strike price

likelihood of paying strike price
(ie likelihood of finishing ITM)

in other
words

in other
words:

The amount you will
earn for owning this
option on avg

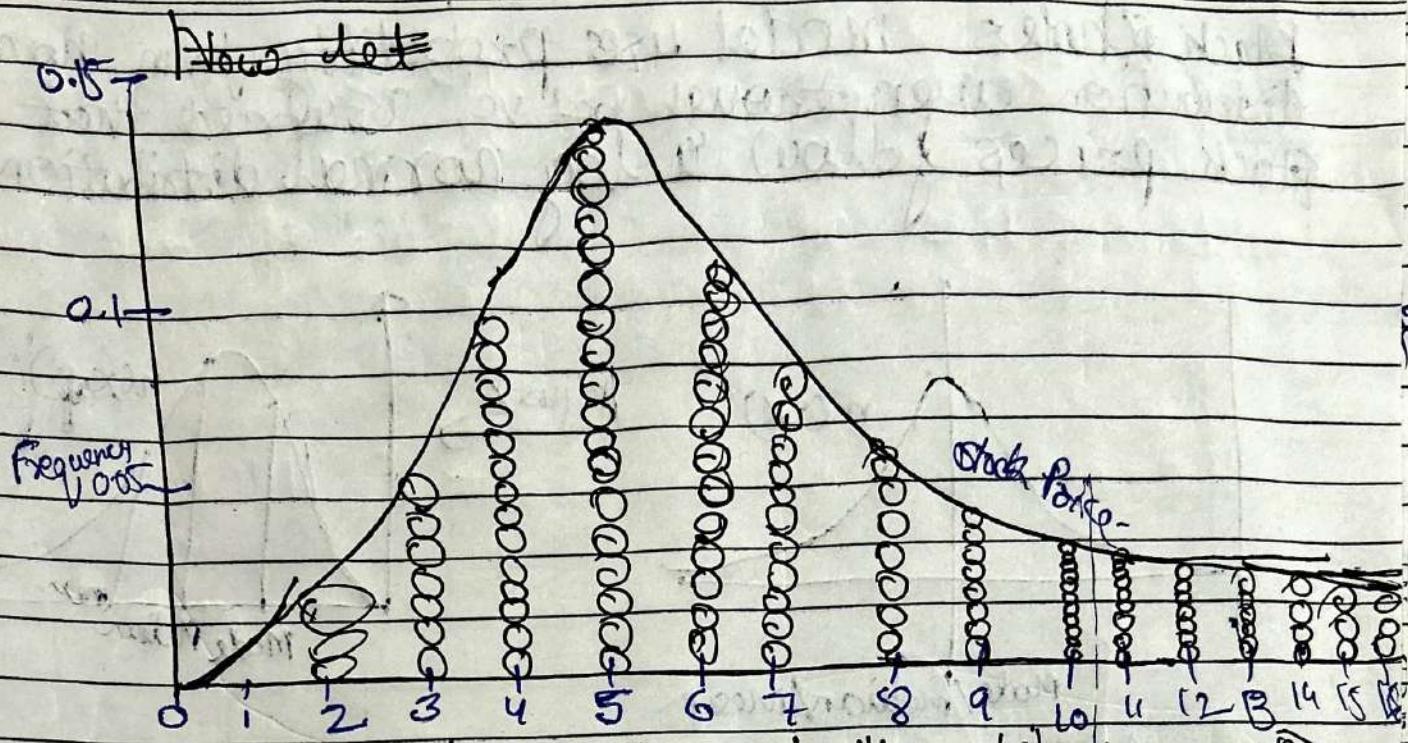
likelihood that the option holder will
end up paying the exercise
price at expiration

Multiplying by strike X
we can calculate
how it will cost on avg
to hold this option.

$$C = \text{Profit}_{\text{avg}} - \text{Cost}_{\text{avg}}$$

To explain this concept further we will imagine
log normal distribution of stock prices at
expiration where each possible price is at 50 cent
interval.

A log normal distribution is just one where
taking the log of values gives you a normal distribution
it works for stock prices because stocks ~~for~~ can't drop
below zero and tend to follow this pattern historically



Now let's evaluate a call option with a strike price of

$$\text{Profit}_{\text{avg}} = \frac{\text{sum of values above strike price}}{\text{Total # of values}} \cdot \$10.5$$

$$= \frac{\sum \$11 \times 8 + \dots + 20 \times 1}{\text{Total # of values}} = \$3.30$$

$$P(\text{ITM}) = \frac{\text{# of values above strike price}}{\text{Total # of values}} = 24.00\%$$

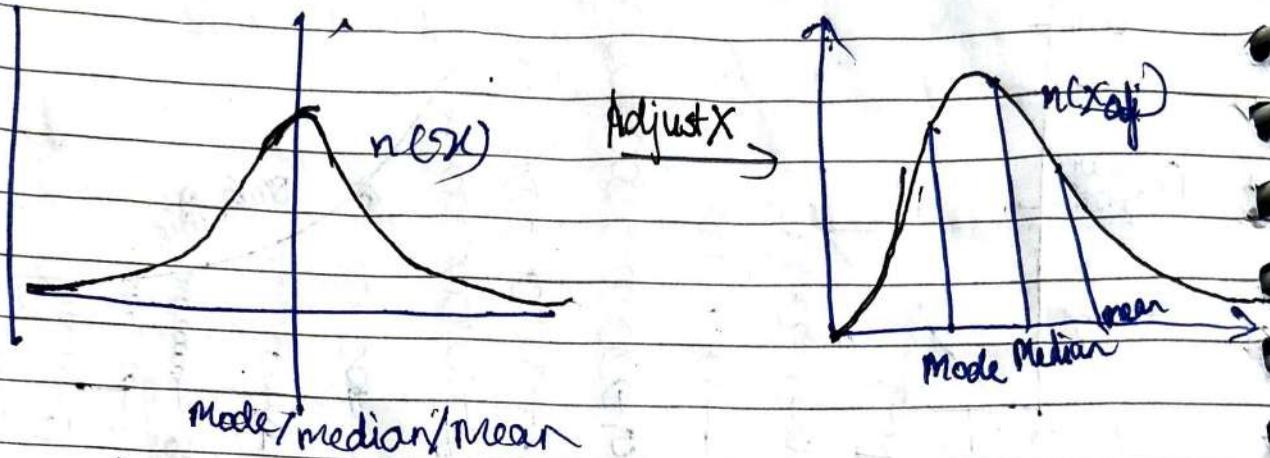
$$P(\text{ITM}) = 24\%$$

$$\text{Cost}_{\text{avg}} = P(\text{ITM}) * \text{strike} = 24\% \times 10.5 \\ = \$2.52$$

$$C = \text{Profit}_{\text{avg}} - \text{Cost}_{\text{avg}} = \$0.78$$

From this we can conclude that the current fair value of the option is ~~the~~ 78 cents.

Black-Scholes model uses probability from normal distribution even though we've assumed that stock prices follow a log-normal distribution.



In a perfect normal distribution these three align but in a log-normal distribution they're spread out understanding these points helps adjust from a log-normal to a normal framework.

In Black-Scholes the relationship b/w strike price & stock price is key for normal dist it's just the stock price S minus strike price X $S - X$

But in log-normal form we have to adjust by taking log giving $\ln(S/X)$ $\ln(S/X)$

If $S > X$ value is +ve & call is in-the-money (ITM)

If $S < X$ it's -ve & call is out-of-money (OTM).

Now let's further adjust for two key components that affect stock prices.

- interest rate.
- volatility.

Since option price valued based on stock price & stock price depend upon interest rate.
we need to account for interest rates over the life of option this adjustment ensures that pricing reflect the present value of the stock's expected price at expiration.

$$\ln\left(\frac{S}{X}\right) + \gamma t$$

Volatility

In a log normal dist stock price are skewed with a longer right tail this results the mean of dist being shifted right side to the mode.

The shift is mathematically eqvt to $\sigma^2 t / 2$.

Combining adjustments interest rate & volatility

Give num^x for d,

$$\ln\left(\frac{S}{X}\right) + \left(\gamma + \frac{\sigma^2}{2}\right)t$$

To wrap things up we need to convert this value into number of standard deviations because S.D help measure how far the strike price is from the mean in a log normal distribution.

option is in or out of the money

$$d_1 = \ln\left(\frac{S}{X}\right) + \left(\gamma + \frac{\sigma^2}{2}\right)t$$

*forward price adjustment
(from interest)*

*Adjustment for mean
lognormal distribution*

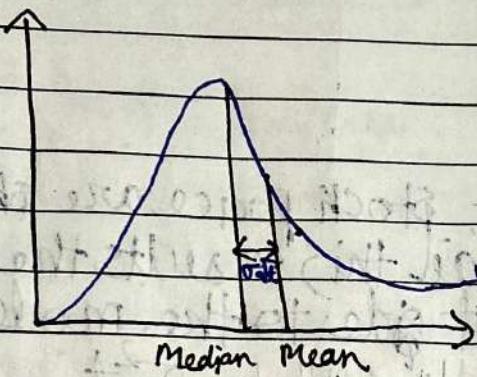
$\sigma \sqrt{T}$

*Normalisation factor (one standard deviation)
to determine the total amount of standard deviation*

Essentially d_1 tells us how many standard deviations the strike price is from the mean.

Combining this with the final price of stock gives us the expected avg profit of the option at expiration.

$$\text{Average profit at Expiration} = S e^{rt} N(d_1)$$



To figure out likelihood the option will be exercised we use the median of the log normal distribution. This is the point that splits the total occurrences in half and it usually falls to the left of the mean.

$$d_2 = d_1 - \sigma \sqrt{t}$$

The value d_2 uses the median to estimate the probability of the being in the money at expiration multiplying this probability by the exercise price gives us the avg amount we'll pay if we own the option.

$$\text{Average price paid at Expiration} = X N(d_2)$$

Combining this with the first term will give us the expected value of call option at expiration ~~but because we're paying for option today~~

$$\text{Expiration} = S e^{rt} N(d_1) - X n(d_2)$$

but because we are paying for the option today
 we have to adjust for present value multiplying
 both terms by e^{-rt} gives us the familiar Black
 Scholes model -

$$\boxed{C_{\text{today}} = S N(d_1) - X e^{-rt} N(d_2)} \quad \star\star\star$$

$$C_t + \sigma S C_S + \frac{1}{2} \sigma^2 S^2 C_{SS} = \sigma C$$

$$C(s, T) = \max(s - k, 0)$$

$$S = K e^x \Rightarrow x = \ln(S/k)$$

$$\frac{dC}{dS} = \frac{dc}{dx} \frac{1}{S} \Rightarrow S \frac{dc}{dS} = \frac{dc}{dx}$$

$$\frac{\partial^2 C}{\partial S^2} = \frac{\partial^2 c}{\partial x^2} \left(\frac{1}{S} \right)^2 = \frac{\partial^2 c}{\partial x^2} - \frac{c}{S}$$

$$\frac{\partial^2 C}{\partial t^2}$$

$$\text{Let } d = T - \frac{T(\sigma^2/2)}{2}$$

$$C_t = \frac{\partial c}{\partial t} \cdot \frac{\partial T}{\partial t} = c_T \left(\frac{\sigma^2}{2} \right)$$

$$-\frac{\sigma^2}{2} c_T + \sigma c_x + \frac{1}{2} \sigma^2 [c_{xx} - c_x] = \sigma C$$

$$-c_T + \lambda c_x + c_{xx} - c_x = \lambda C$$

$$c_T = (x-1)c_x + c_{xx} - \lambda C$$

$$C = KV \quad \boxed{V_T = (\lambda-1)V_x + V_{xx} - \lambda V}$$

Black-Scholes Equation (Explicit)

Definition

Black-Scholes Equation (BSE) is used to calculate the fair market value of a European option.

$$\text{The PDE} \rightarrow \frac{\partial V}{\partial t} + \frac{1}{2} \sigma^2 S^2 \frac{\partial^2 V}{\partial S^2} + rS \frac{\partial V}{\partial S} - rV = 0 \quad \text{--- (1)}$$

- K = strike price (pre-determined price at which holder of the option can buy the stock)
- V = price of the option \Rightarrow what we're trying to calculate
- S = current price of the stock/asset (constantly changing)
- t = time, T = expiration time of the option
- r = risk-free interest rate (like the interest rate on a Govt Bond)
- σ = Volatility of the asset
 - high σ (startup) \Rightarrow high value \Rightarrow expensive
 - low σ (utility company) \Rightarrow low " \Rightarrow cheap

$\frac{\partial V}{\partial t}$ = rate of change of option's price wrt time

$\frac{\partial V}{\partial S}$ = " " " " " asset's price

$\frac{\partial^2 V}{\partial S^2}$ = " " " " $\frac{\partial V}{\partial S}$ " " "

Derivation

In order to computerize this equation, we take approximations of the derivatives according to FDM.

We will use grid $V_{ij} \Rightarrow$ option value at time i , stock price j for $\frac{\partial^2 V}{\partial S^2}$, we take the central difference formula for second derivative

$$\text{i.e. } u''(x) = \frac{u(x+h) + u(x-h) - 2u(x)}{h^2}$$

in the grid V_{ij} , $u(x+h) \Rightarrow V_{i+1,j}$, $u(x-h) \Rightarrow V_{i-1,j}$, $u(x) \Rightarrow V_{i,j}$

$$h \rightarrow \Delta S = S_{\max} - S_{\min}$$

$$\therefore \frac{\partial^2 V}{\partial S^2} = \frac{V_{ij+1} + V_{ij-1} - 2V_{ij}}{(\Delta S)^2} \quad \text{--- (2)}$$

We take $S_{max} = 2K$ for a wide enough grid

S_{min} is taken as 0

M is the no. of steps we take in the price

$N = \dots \dots \dots \dots \dots \dots \text{time}$

$$\text{so } \Delta S = \frac{S_{max} - S_{min}}{M} = \frac{2K}{M}$$

$$\Delta t = T/N$$

for $\frac{\partial V}{\partial S}$, we ~~can~~ take central difference formula for first derivative

$$\frac{\partial V}{\partial S} \text{ i.e. } u'(x) = \frac{u(x+h) - u(x-h)}{2h}$$

$$\therefore \frac{\partial V}{\partial S} = \frac{V_{ij+1} - V_{ij-1}}{2 \Delta S} \quad \text{--- (3)}$$

for $\frac{\partial V}{\partial t}$, we take backward difference formula for first derivative

If we were to take central difference formula, it would be more accurate, but then we would have 3 different time steps
 $i-1$, i and $i+1$.

So at the cost of some accuracy, we use the backward difference formula so as to have only $i-1$ and i . s.t. $i-1$ can be taken as the only unknown variable.

$$\therefore \frac{\partial V}{\partial t} = \frac{V_{ij} - V_{i-1j}}{\Delta t} \quad \text{--- (4)}$$

from (1), (2), (3) and (4), we get

$$\frac{V_{ij} - V_{i,j+1}}{\Delta t} + \frac{1}{2} \sigma^2 S_j^2 \left(\frac{V_{ij+1} + V_{ij-1} - 2V_{ij}}{(\Delta s)^2} \right) + r S_j \left(\frac{V_{ij+1} - V_{ij-1}}{2\Delta s} \right) - r V_{ij} = 0$$

$$S_j = j \Delta s$$

Transformed a PDE to an algebraic formula

$$\Rightarrow \frac{V_{ij} - V_{i,j+1}}{\Delta t} = -\frac{1}{2} \sigma^2 (j \Delta s)^2 \left(\frac{V_{ij+1} + V_{ij-1} - 2V_{ij}}{(\Delta s)^2} \right) - r(j \Delta s) \left(\frac{V_{ij+1} - V_{ij-1}}{2 \Delta s} \right) + r V_{ij}$$

$$= -\frac{1}{2} \sigma^2 j^2 (V_{ij+1} + V_{ij-1} - 2V_{ij}) - \frac{1}{2} r j (V_{ij+1} - V_{ij-1}) + r V_{ij}$$

$$V_{i-1,j} = V_{ij} + \Delta t \left[-\frac{1}{2} \sigma^2 j^2 (V_{ij+1} + V_{ij-1} - 2V_{ij}) + \frac{1}{2} r j (V_{ij+1} - V_{ij-1}) - r V_{ij} \right]$$

$$\text{Coeff of } V_{ij+1} \rightarrow -\frac{1}{2} \sigma^2 j^2 \Delta t - \frac{1}{2} r j \Delta t = \frac{1}{2} \Delta t j (\sigma^2 j - r) = a$$

$$\text{Coeff of } V_{ij} \rightarrow -\frac{1}{2} \sigma^2 j^2 \Delta t - r \Delta t = 1 - \Delta t (\sigma^2 j^2 + r) = b$$

$$\text{Coeff of } V_{ij-1} \rightarrow \frac{1}{2} \sigma^2 j^2 \Delta t + \frac{1}{2} r j \Delta t = \frac{1}{2} \Delta t j (\sigma^2 j + r) = c$$

$$V_{i-1,j} = a_j V_{ij+1} + b_j V_{ij} + c_j V_{ij-1}$$

$$\therefore V_{i-1,j} = a_j V_{ij+1} + b_j V_{ij} + c_j V_{ij-1} \quad \text{--- (5)}$$

This is the equation we will use to code this up

Stability

This method only gives reliable solution if

$$\Delta t \leq \frac{1}{\sigma^2 M^2}$$

To get this, increase N to a suitable number.
 ↳ No. of time steps.

Boundary conditions

In order to use eq 5, we will need some boundary conditions i.e. the values at the edges of our grid, V_{ij}

1) At Expiration ($t = T$)

At the time of expiry, the value of the option is its immediate cash value. The holder will only exercise their right to buy the stock if $S > K$ (current price of stock is more than the agreed upon price) and their profit is $S - K$. If $S \leq K$, the option is worthless at value 0.

$$\therefore R(t, T, j) V_{T,j} = \max(S_j - K, 0) + t_j$$

2) At zero stock price ($S=0$)

At any time if S reaches 0, the value of the option is 0 as it is the right to buy a stock at rate K . 0 is an absorbing barrier, once a stock's S hits 0, it can't go any higher given any amount of time (Bankrupt)

$$\therefore V_{t,0} = 0 + t$$

3) At a very high stock price ($S = S_{\max}$)

If S hits very high, almost certainly the option will be exercised at expiration. \therefore holder will pay K for S .

But the payment of K will happen in future at expiration so to get its present value by discounting it back to t using risk free rate r ,

$$R(t, S_{\max}, V_{t, \max}) = S_{\max} - K e^{-r(T-t)}$$

With these three conditions in place, we can write the code

Later on, I learned that the convention is to write $V(S, T)$ instead of $V(T, S)$
So just made the changes in the code.

Black Scholes Equation (Implicit)

Reason

The explicit method is fast but suffers with conditional stability. Implicit method is used to overcome this weakness. It requires more calculation but is **unconditionally stable**. Here, no matter the choice of M and N , the solution is stable.

Derivation

$$\frac{\partial V}{\partial t} = \cancel{\frac{V_{i+1,j} - V_{i-1,j}}{\Delta t}} \quad \text{This is the same as in explicit}$$

$$\frac{\partial^2 V}{\partial S^2} = \cancel{\frac{V_{i+1,j+1} - 2V_{i,j+1} + V_{i-1,j+1}}{(\Delta S)^2}}$$

$$\frac{\partial V}{\partial S} = \frac{V_{i+1,j+1} - V_{i-1,j+1}}{2\Delta S}$$

$$\frac{\partial V}{\partial t} + \frac{1}{2} \sigma^2 S_j^2 \frac{\partial^2 V}{\partial S^2} + rS \frac{\partial V}{\partial S} - rV = 0$$

$$\Rightarrow \frac{V_{ij} - V_{i+1,j}}{\Delta t} = -\frac{1}{2} \sigma^2 S_j^2 \frac{V_{i-1,j+1} - 2V_{i,j+1} + V_{i+1,j+1}}{(\Delta S)^2} + rS \frac{V_{i-1,j+1} - V_{i-1,j+1}}{2\Delta S} + rV_{i-1,j}$$

$$S_j = j\Delta S$$

$$\Rightarrow V_{ij} - V_{i-1,j} = -\Delta t \left[\frac{1}{2} \sigma^2 (\Delta S)^2 \frac{V_{i-1,j+1} - 2V_{i,j+1} + V_{i+1,j+1}}{(\Delta S)^2} + rj\Delta S \frac{V_{i-1,j+1} - V_{i-1,j+1}}{2\Delta S} + rV_{i-1,j} \right]$$

$$V_{ij} = V_{i-1,j} + \Delta t \left[\frac{1}{2} r^2 j^2 (V_{i-1,j+1} - 2V_{i,j+1} + V_{i+1,j+1}) + \frac{1}{2} rj (V_{i-1,j+1} - V_{i-1,j+1}) - rV_{i-1,j} \right]$$

$$\text{Coeff of } V_{i-1,j+1} \rightarrow -\frac{1}{2} \sigma^2 j^2 \Delta t + \frac{1}{2} rj \Delta t = -\frac{1}{2} j \Delta t (\sigma^2 j^2 + r) = a_j$$

$$\text{Coeff of } V_{i-1,j} \rightarrow 1 + \sigma^2 j^2 \Delta t + r \Delta t = 1 + \Delta t (\sigma^2 j^2 + r) = b_j$$

$$\text{Coeff of } V_{i-1,j+1} \rightarrow -\frac{1}{2} \sigma^2 j^2 \Delta t - \frac{1}{2} rj \Delta t = -\frac{1}{2} j \Delta t (\sigma^2 j^2 + r) = c_j$$

$$V_{ij} = a_j V_{i+j+1} + b_j V_{i-j} + c_j V_{i+j+1}$$

Known Unknowns

\therefore There are 3 unknowns, we must use $\boxed{Ax=b}$

$$x \rightarrow \text{Unknowns vector} \rightarrow [V_{i+1,1}, V_{i+1,2}, \dots, V_{i+1,M+1}]^T$$

$$b \rightarrow \text{Knowns vector} \rightarrow [V_{i,1}, V_{i,2}, \dots, V_{i,M+1}]^T$$

$A \rightarrow \text{coeff matrix} \rightarrow \text{Tridiagonal} \rightarrow \begin{matrix} \text{main diagonal} \Rightarrow b_j \\ \text{lower "} \Rightarrow a_j \\ \text{upper "} \Rightarrow c_j \\ \text{All other entries are 0} \end{matrix}$

Coding logic

$V \rightarrow M+1 \times N+1$ matrix

Boundary conditions same as in explicit ($t=T$, s_0 , $s=s_{\max}$)

The change!

- loop from $i=N+1$ to 1 $\rightarrow b$ Happens outside loop
 At each step, i is the known col $j-1$ is the unknown col

- Assemble A as the tridiagonal Matrix
- Assemble b using i^{th} col of V
- solve $Ax=b$
- solved x containing the new $i-1^{\text{th}}$ col of V

Once the loop finishes, the option prices at $t=0$ are in the 0^{th} col of V