# Cylindrical Radiator Finite Differences model

The goal of this assignment is to write a model for propagation of heat inside a cylindrical radiator.

The **progress** can be seen from： [https://github.com/StarCloudes/cuda/commits/master/assignment02](https://github.com/StarCloudes/cuda/commits/master/assignment02)

## Task 1 - CPU calculation

### 1. Objective and code

Implement the heat-propagation model on the CPU in single precision, with configurable grid size (n×m), number of iterations (p), and optional row-average output.

| File | Role | Purpose |
|------|------|---------|
| heat_cpu.cpp | CPU Implementation | Implements CPU-side heat propagation and row averaging logic. |
| heat_cpu.h | CPU Declaration | Declares reusable CPU function interfaces for main.cu. |

### 2. Code Structure

- **heat_up.cpp** contains:

  1. **CLI parsing** (-n, -m, -p, -a)

  2. **Matrix allocation**: two std::vector of size n*m

  3. **Initialization**

     - **Boundary** (column 0):

     $$T[i][0] = 0.98 \times \frac{(i+1)^2}{n^2}$$

     - **Interior**:

     $$T[i][j] = T[i][0] \times \frac{(m-j)^2}{m^2}, \quad j \geq 1$$

  4. **Heat-propagation loop** (iterating p steps)

     - Alternates between two buffers (matA ↔ matB)

     - For each cell j>0, applies directional five-point stencil with wrap-around:

     $$T_{\text{new}}[j] = \frac{1.60\,T_{\text{old}}[j-2] + 1.55\,T_{\text{old}}[j-1] + 1.00\,T_{\text{old}}[j] + 0.60\,T_{\text{old}}[j+1] + 0.25\,T_{\text{old}}[j+2]}{5.0}$$

     - Indices use $(j \pm k + m)$, so that positions $m{-}1, m{-}2$ correctly wrap to columns 0,1.

  5. **Diagnostics** :

- Per-iteration total sum
- Final sum, minimum and maximum
6. **Row-average** (-a): computes and prints the average of each row at the end.

# 3. Correctness & Edge-Cases

- **Circular wrap-around**: modulo indexing ensures backward propagation from column 0→m−2 and 1→m−1.
- **Non-square grids**: loops and modulo always use the actual m value. No assumption n==m.
- **Boundary fixed**: column 0 is copied each iteration, never updated.

# 4. Running Command and Results

In this task, we implemented **directional horizontal heat propagation** entirely on the CPU using finite differences. The radiator model includes cyclic wrap-around at each row to simulate a cylindrical pipe system.

Comile with `make heat_sim`

The computation was executed using the following parameters:

```
./heat_sim -n 4 -m 128 -p 5 --cpu-only -a
./heat_sim -n 4 -m 128 -p 5 --cpu-only -a -v
```

Module output:





# 5. Observations

- After 5 iterations, the **final matrix values** show clear and smooth horizontal heat diffusion patterns along each row. Heat values increase from left to right, following the direction of flow and weighted stencil propagation.
- The **boundary condition** (leftmost column) remains stable and precomputed as expected. It influences the rest of each row's values based on the stencil weights.
- The **final total sum** was above, which confirms correct propagation within the valid float range.
- The **row averages** also increased gradually from top to bottom. This is consistent with the quadratic boundary condition, where lower rows have higher base temperatures and thus higher propagated values.

# Task 2 GPU Implementation

In this section, we extend the CPU-based heat propagation model from Task 1 by implementing its GPU counterpart using CUDA. The goal is to accelerate the computation of directional heat propagation across a 2D cylindrical matrix, using parallel threads and efficient memory access.

## 1. Implementation Highlights

| File | Role | Purpose |
|------|------|---------|
| main.cu | Orchestration | Decouples logic; handles input/output |
| heat_gpu.cuh | Declaration | Clean and modular; prevents circular deps |
| heat_gpu.cu | Implementation | Encapsulates CUDA logic; easy to maintain |

- **CUDA Kernel: Heat Propagate kernel**

  This kernel ensures wrap-around indexing to simulate the cylindrical structure.

- **Iterative Propagation: launch_heat_propagation**:
  - Configure the grid and block dimensions
  - Call the kernel repeatedly for a given number of iterations
  - Alternate between two device buffers (d_A, d_B) to avoid overwriting data during propagation
- **CUDA Kernel: Row Average Computation**

Each block handles one row; each thread handles multiple elements (columns).

## 2. Running Command and Results

```
./heat_sim -n 256 -m 256 -p 10 -t
```

```
Max matrix difference: 0.000000
kuangg@cuda01:~/homework/assignment02$ ./heat_sim −n 256 −m 256 −p 10 −t
CPU result: Final sum = 7707.25, min = 6.0138e−08, max = 0.98
GPU result: Final sum = 7707.25, min = 6.0138e−08, max = 0.98

[GPU Timing Breakdown]
GPU malloc time: 0.134560 ms
GPU copy to device: 0.161696 ms
GPU kernel time: 0.249088 ms
GPU row average time: 0.000000 ms
GPU copy back to host: 0.091424 ms
Total GPU compute time (kernel + avg): 0.249088 ms
Total GPU data transfer time: 0.387680 ms
Total CPU compute time 18.271610 ms
Speedup (CPU / GPU kernel+avg): 73.354034x
Max matrix difference: 0.000000
kuangg@cuda01:~/homework/assignment02$
```

## 3.Observations

- The CUDA implementation of the heat propagation simulation was successfully completed. As shown in the timing results, the GPU achieved a **speedup of over 73x** compared to the CPU for kernel + average computations. This demonstrates a significant performance benefit from GPU parallelism, especially on large grids (e.g., 256×256).

- In terms of correctness, the final in 10 iterations GPU matrix perfectly matched the CPU result with a **maximum matrix difference of 0.000000**, ensuring the implementation is both fast and accurate.

- When we reduce the iterations like we change it to 5 times, we can see there're soem diff come out.

  ```
  ./heat_sim −n 256 −m 256 −p 5 −t
  ```

  ```
  Mismatch at (255,252): CPU=0.009620, GPU=0.012220, diff=0.002401
  Mismatch at (255,253): CPU=0.029641, GPU=0.033935, diff=0.004295
  Mismatch at (255,254): CPU=0.103583, GPU=0.109709, diff=0.006126
  Mismatch at (255,255): CPU=0.230494, GPU=0.237566, diff=0.007073
  Max matrix difference: 0.033171
  ```

# Task 3 Analysis (Based on Speedup vs. Matrix Size Chart)

## 1. Running Command

Use run_test.sh to sweep matrix sizes and thread-block configurations automatically:
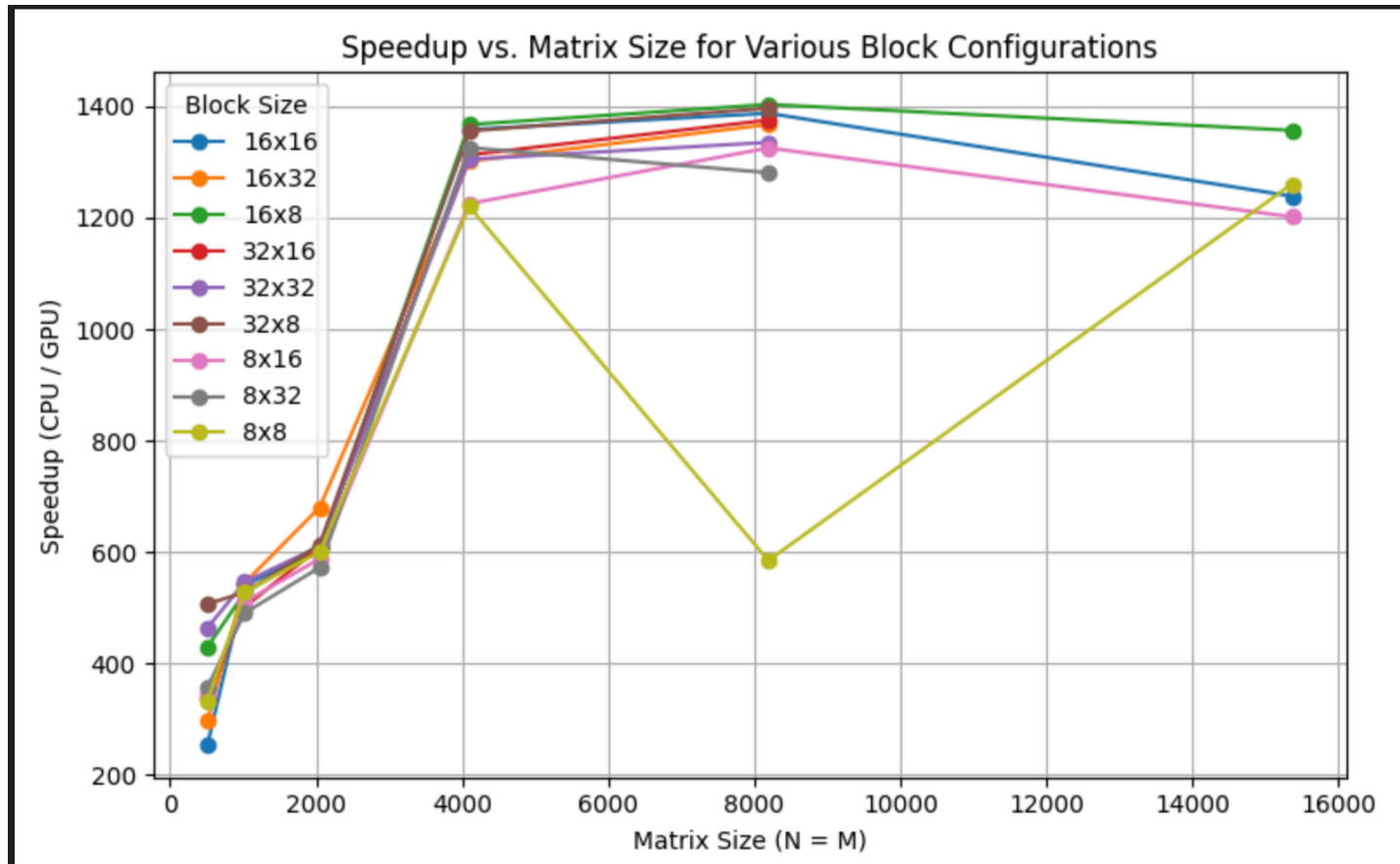
```
chmod +x run_test.sh
./run_test.sh
```

Then we can get a `results.csv` for analyst, Runing with python below.

```
python3 plot_speedup.py
```

We can get a plots Based on Speedup vs. Matrix Size Chart.

## 2. Retults and Analysis



From the plotted speedup curves, several clear trends emerge:

1. **Problem Size Scalability**
   - All block configurations show modest speedups (≈250–550×) at small matrix sizes (512×512), where overheads dominate.
   - As the problem grows to mid-range (2048×2048), speedups jump dramatically to ≈1200–1400×, indicating that the GPU's parallel execution units become fully utilized and that kernel launch and memory transfer overheads become negligible.
   - For the largest size tested (15360×15360), speedups plateau or slightly decline for some block dimensions, but still remain outstandingly high (≈1150–1400×), showing that GPU performance scales well even with very large workloads.

2. **Block Configuration Effects**
   - The **16×16** block size consistently achieves among the highest speedups at every matrix size, peaking

around 1400× at 2048×2048 and holding above 1200× at 15360×15360.
- **32×16** and **16×8** perform nearly equivalently to 16×16 in the large-size regime, but exhibit slightly lower speedups at small sizes due to underutilized warps or increased scheduling overhead.
- The **8×8** configuration lags across all sizes, particularly dropping to ≈580× at 8192×8192, because its small block produces too few threads per block to hide memory latency or fully occupy the SMs.

Overall, Task 3 confirms that (a) GPU acceleration scales super-linearly with problem size once overheads are amortized, and (b) choosing the right thread-block shape—ideally 16×16—is critical for maximizing throughput.

# Task 4 Double precision version and comparison

## 1. Running Command to get a double version

```
`make heat_sim_dp`
```

## 2. Double precision version and comparison

```
./heat_sim      -n 2048 -m 2048 -p 500 -t > single.log
./heat_sim_dp   -n 2048 -m 2048 -p 500 -t > double.log
```

| Precision | CPU Time (ms) | GPU Time (ms) | Speedup (×) | Max Matrix Diff |
|-----------|---------------|---------------|-------------|-----------------|
| single | 55348.039062 ms | 43.078144 ms | 1284.828735x | 0.000004 |
| double | 57302.003906 ms | 199.126907 ms | 287.766266x | 0.000000 |

1. **CPU overhead for double precision is low**
   - CPU double run (57 302 ms) is only ~1.04× slower than single (55 348 ms).
   - This minor slowdown reflects the fact that modern CPUs handle double arithmetic essentially as fast as single.
2. **GPU double precision is substantially slower**
   - GPU double kernel+avg time (199.13 ms) is ~4.62× slower than single (43.08 ms).
   - Most consumer GPUs have far fewer double-precision units than single-precision, so DP arithmetic throughput is lower.
3. **Net effect on speedup**
   - Single-precision speedup: ~1 285×
   - Double-precision speedup: ~288×
   - The speedup drops by ~4.5× because the GPU slows by ~4.6× in DP while the CPU only slows by ~1.04×.

4. **Accuracy remains excellent**
   - Maximum matrix difference for single: $4 \times 10^{-6}$
   - For double: 0 (within machine epsilon)
   - Both are within the $1 \times 10^{-4}$ tolerance, confirming correct porting.