

Heterogenous CPU-GPU ordinary differential equation solver

Dániel Nagy

Co-author: Dr. Ferenc Hegedűs

BME, Department of Hydrodynamic Systems

GPU Day 2024



Department of
Hydrodynamic
Systems

1 Introduction

- Motivation
- Per-thread approach

2 Methods

- CUDA Streams
- Overlapping
- Test equations

3 Results

- TEST1, Bernoulli
- TEST2, Duffing

4 Conclusion

Developing a heterogeneous CPU-GPU ordinary differential equation solver

- Using CUDA
- Utilizing both CPUs and GPUs for optimal performance
- Speeding-up the solution of parameter studies
- Possible application for high-complexity problems
 - Delay differential equations
 - Coupled differential equations

Parameter study

Systematic examination of the parameter dependent behavior of a dynamic system.

Each thread calculates the same ODE¹:

$$\dot{x} = f(x, t; p) \quad (1)$$

Can be solved numerically using the traditional 4th order Runge-Kutta method as

$$x_{n+1}^{\{i\}} = x_n^{\{i\}} + \frac{\Delta t}{6} \left(k_1^{\{i\}} + 2k_2^{\{i\}} + 2k_3^{\{i\}} + k_4^{\{i\}} \right) \quad (2)$$

$$k_1^{\{i\}} = f(x_n^{\{i\}}, t; p^{\{i\}}) \quad k_2^{\{i\}} = \dots \quad (3)$$

Notice

- i is the thread index, $i = 0 \dots N - 1$; N is the number of threads, n is the current time-step.
- The threads do not exchange data.

¹Dániel Nagy, Lambert Plavec, and Ferenc Hegedűs. "The art of solving a large number of non-stiff, low-dimensional ordinary differential equation systems on GPUs and CPUs". In: *Communications in Nonlinear Science and Numerical Simulation* 112 (2022), p. 106521.

- Thread divergence if a time-step is rejected (adaptive methods)
- Register-spills (e.g., solving the Keller-Miksis equation using RKCK45)
- Calculating delay terms (uncoalesced memory operations)
- Not enough register memory to store coupling matrices

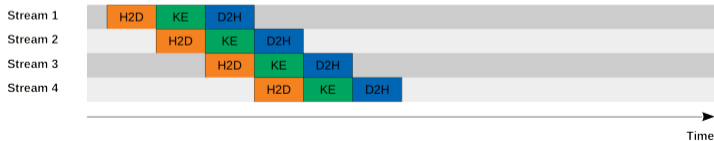
CUDA Stream

A sequence of operations that execute asynchronously on a CUDA device. It allows for concurrent execution of multiple tasks on the GPU, enabling parallel processing and overlapping of computation with data transfers.

Serial Model



Concurrent Model

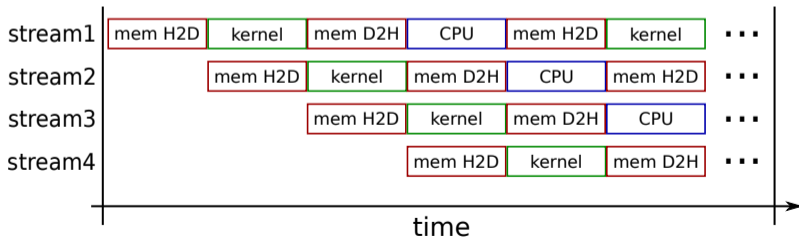


Source: Lei Mao (<https://leimao.github.io/blog/CUDA-Stream/>)

- each stream handles 1/4 of the workload

- each stream handles 1/4 of the workload
- each stream does the following algorithm
 - 1 copy x_n, t_n from CPU to GPU (mem H2D)
 - 2 kernel call: 1 Runge-Kutta step in parallel to determine x_{n+1}, t_{n+1}
 - 3 copy x_{n+1}, t_{n+1} from GPU to CPU (mem D2H)
 - 4 execute serial operations on the CPU on the data

- each stream handles 1/4 of the workload
- each stream does the following algorithm
 - 1 copy x_n, t_n from CPU to GPU (mem H2D)
 - 2 kernel call: 1 Runge-Kutta step in parallel to determine x_{n+1}, t_{n+1}
 - 3 copy x_{n+1}, t_{n+1} from GPU to CPU (mem D2H)
 - 4 execute serial operations on the CPU on the data
- each substep requires the same amount of time

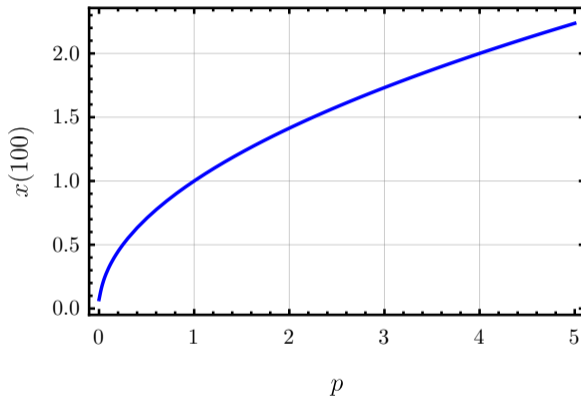


$$\dot{x}(t) = px - x^3; \quad x(0) = 1 \quad (4)$$

- Analytical solution exists, used for validating the results

$$x(t) = \frac{\sqrt{p}e^{pt}}{\sqrt{-1 + e^{2pt} + p}} \quad (5)$$

- Solution in $t \in [0,100]$ and $p \in [0,5]$
- $\Delta t = 0.1$, RK4 is used



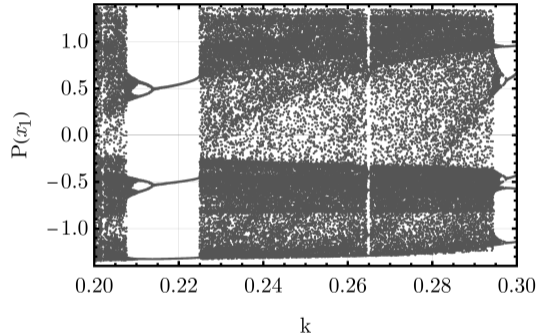
End values at $t = 100$

$$\dot{x}_1 = x_2 \quad (6)$$

$$\dot{x}_2 = x_1 - x_1^3 - kx_2 + 0.3 \cos(t) \quad (7)$$

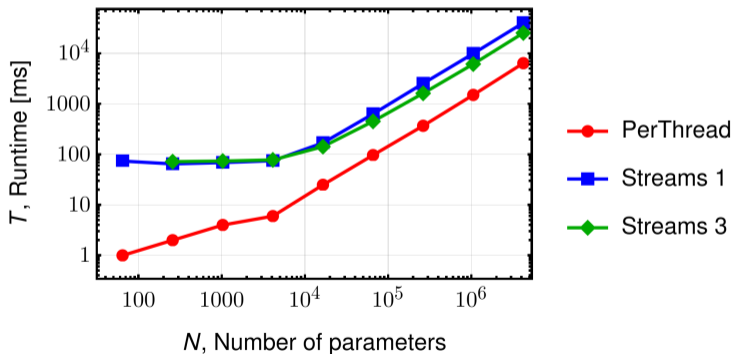
$$x_1(0) = -0.5 \quad x_2(0) = -0.1 \quad (8)$$

- $k \in [0.2, 0.3]$
- Transient simulation for $t \in [0, 1024 \cdot 2\pi]$
- 32 Poincare sections at $t = (1024 + i) \cdot 2\pi$, where $i = 1 \dots 32$
- Adaptive 5th order Runge-Kutta Cash-Karp method using abs. and rel. tolerance 10^{-6}



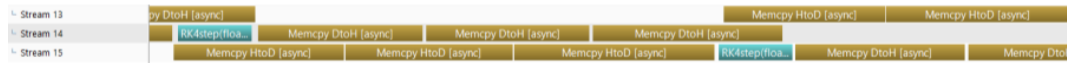
Bifurcation diagram as a result

- Pure GPU (PerThread) is the fastest
- Some overlapping is possible, using 3 streams is $1.6\times$ faster

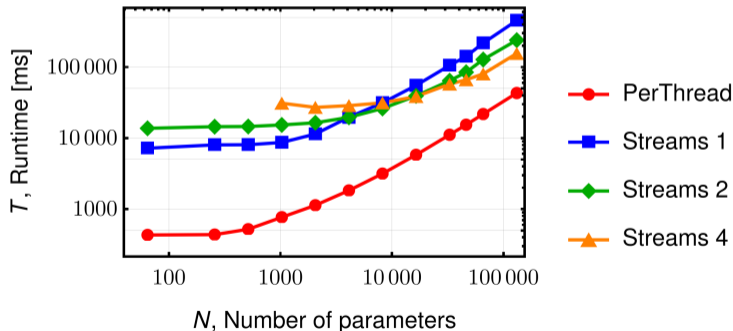


Problem:

- A single RK4 step requires much less time than the memory operations
- Streams are idle



- Pure GPU (PerThread) is the fastest
- Using 2 streams is 1.9 \times faster than 1 stream
- Using 4 streams is 2.9 \times faster than 1 stream



- CPU calculations are short (only checks if the end time is reached)
- Overlapping is possible and makes things faster



- Overlapping memory operations and CPU/GPU calculations is possible
- For the test problems slower than the per-thread approach
- Problem: **global memory operations are slow**

Future plans:

- Can be good if a lot of serial calculations are required.
- Usage for an **adaptive delay differential equation solver** where the delayed terms are found at arbitrary memory locations.

- 1 Introduction
 - Motivation
 - Per-thread approach
- 2 Methods
 - CUDA Streams
 - Overlapping
 - Test equations
- 3 Results
 - TEST1, Bernoulli
 - TEST2, Duffing
- 4 Conclusion



Nagy, Dániel, Lambert Plavec, and Ferenc Hegedűs. "The art of solving a large number of non-stiff, low-dimensional ordinary differential equation systems on GPUs and CPUs". In: *Communications in Nonlinear Science and Numerical Simulation* 112 (2022), p. 106521.