

VU Visualisierung 2 (186.833)

Kernel density estimation in accelerators

Research Paper Implementation

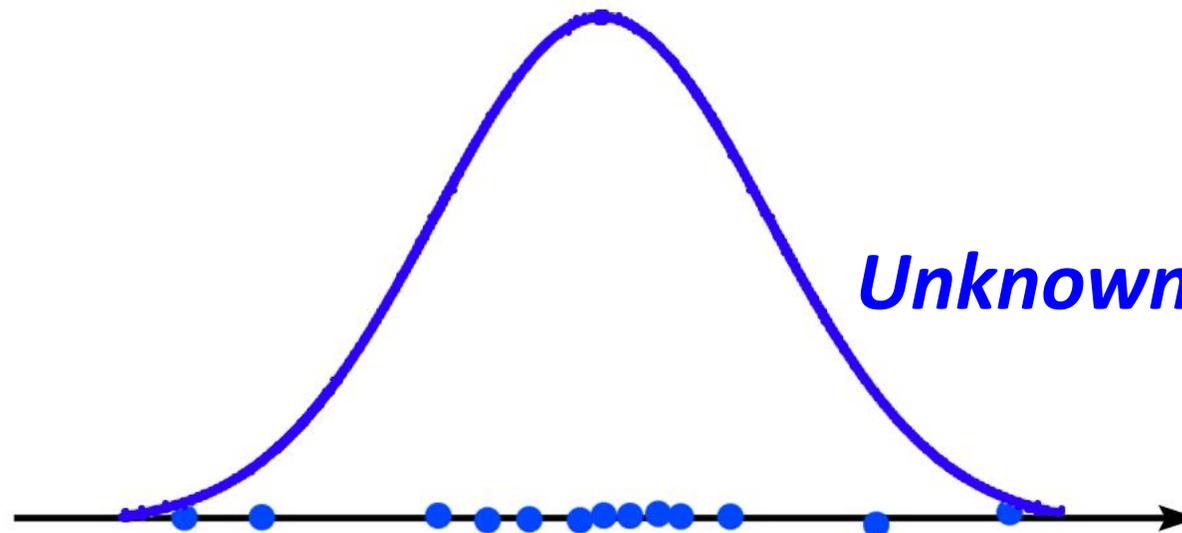
Franciszek Szewczyk (12510815) / Sam Engel (12532780)

1. You observe a *limited set of samples* from some unknown distribution



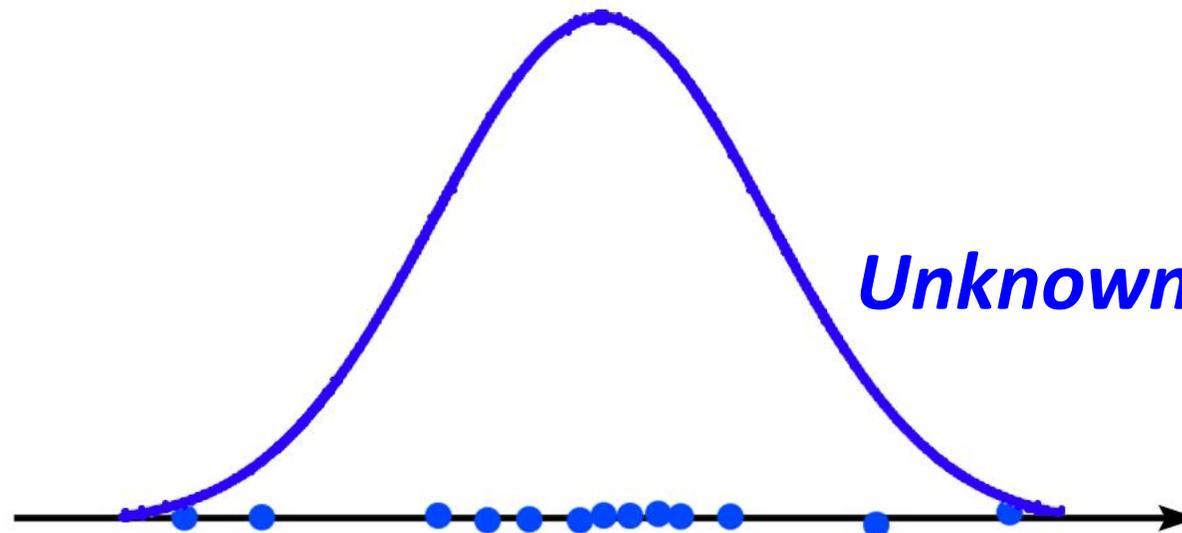
Motivation

1. You observe a *limited set of samples* from some unknown distribution
2. How do we *estimate the original distribution* it came from?



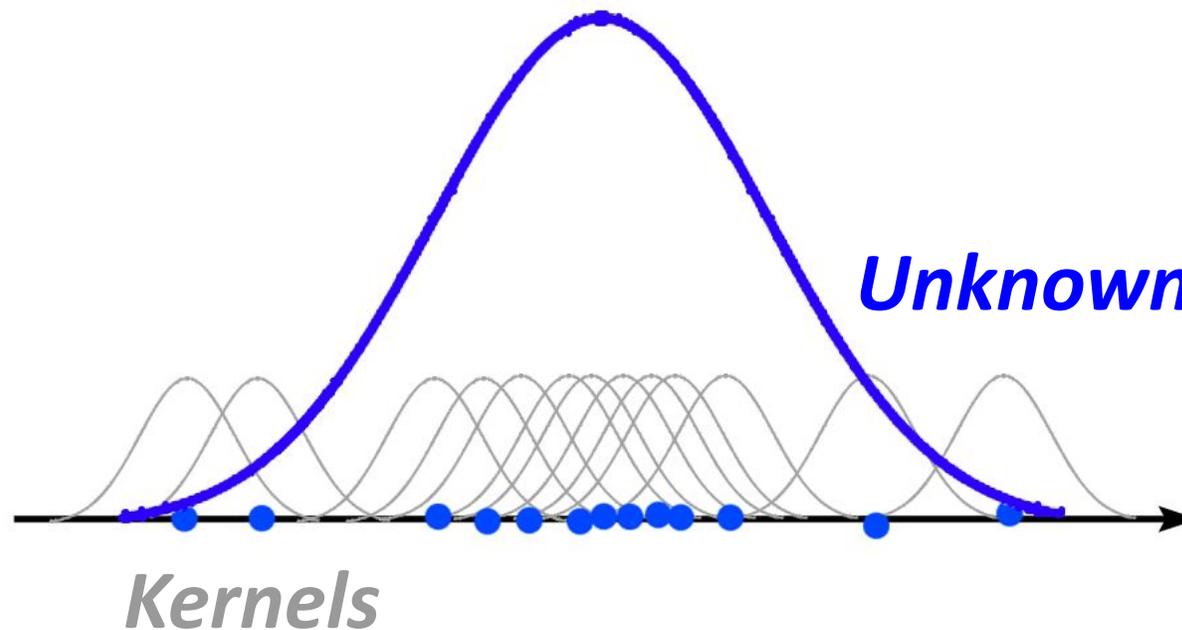
Motivation

1. You observe a **limited set of samples** from some unknown distribution
2. How do we **estimate the original distribution** it came from?
 - a. If we knew it, we could simulate more samples
 - b. Visualize the data
 - c. Understand the problem at hand better



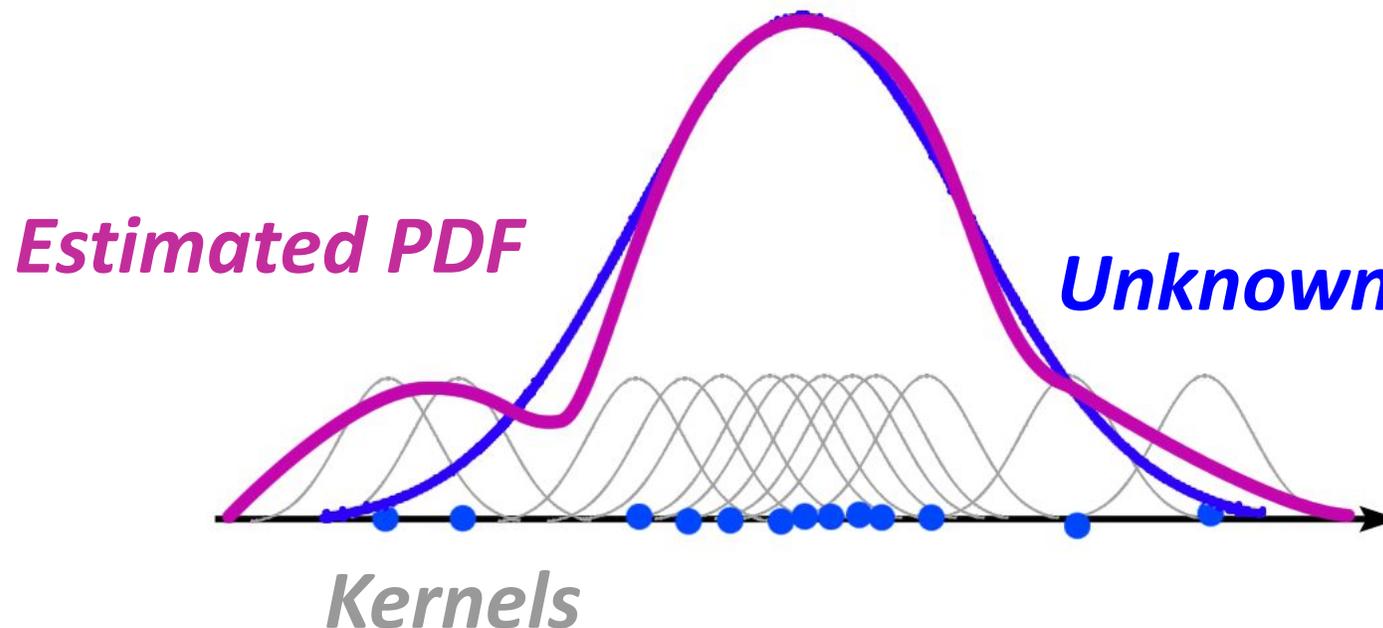
Motivation

1. You observe a **limited set of samples** from some unknown distribution
2. How do we **estimate the original distribution** it came from?
3. We can add a small **kernel** around each sample



Motivation

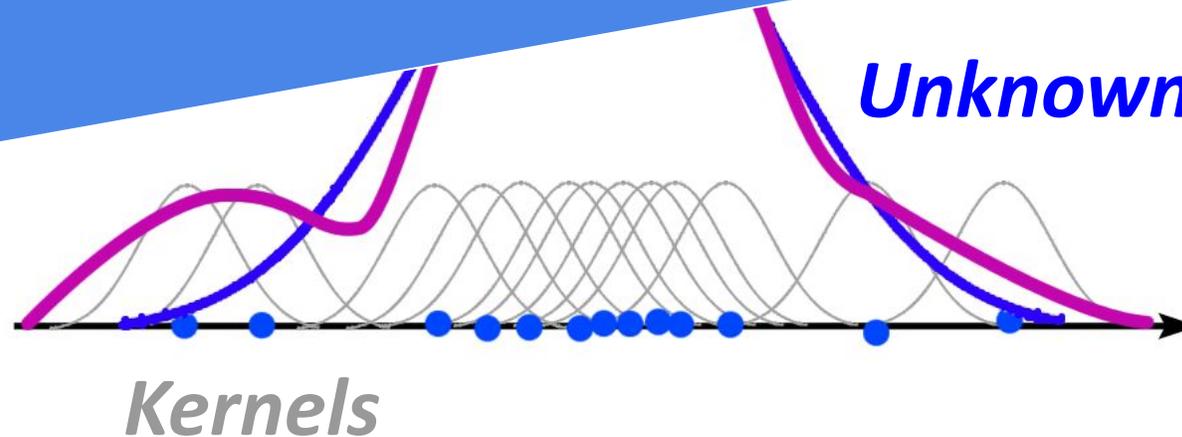
1. You observe a **limited set of samples** from some unknown distribution
2. How do we **estimate the original distribution** it came from?
3. We can add a small **kernel** around each sample
4. And **evaluate their sums** in a dense grid



Motivation

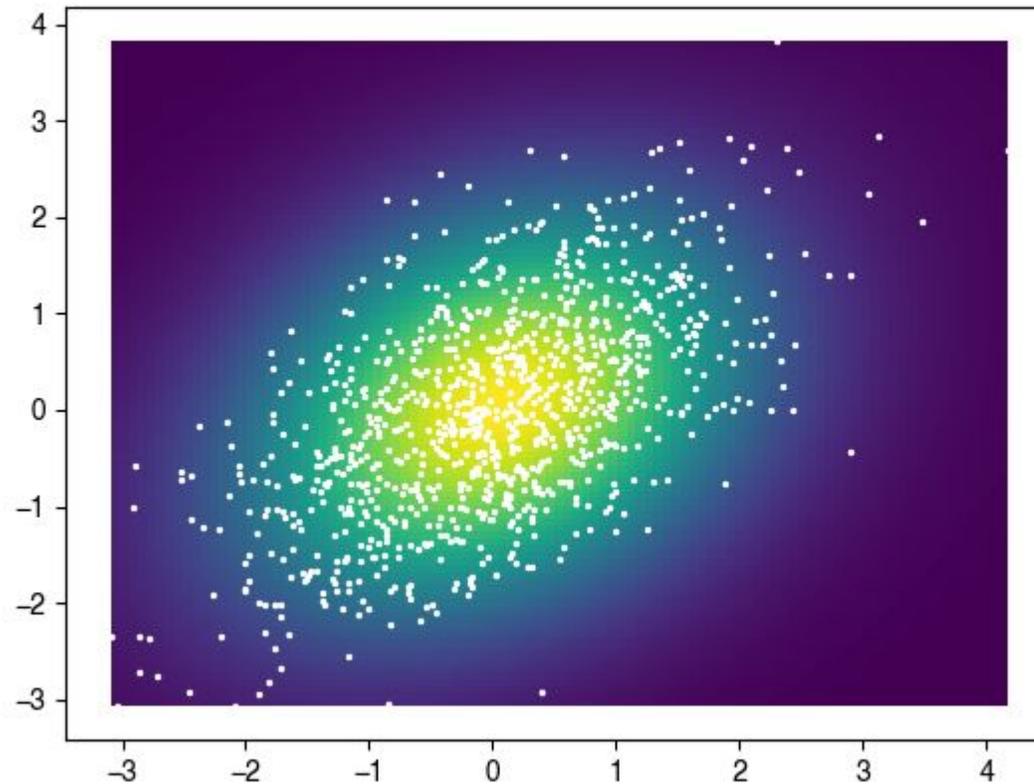
1. You observe a *limited set of samples* from some unknown distribution
2. How do we *estimate the original distribution* it came from?
3. We can add a small *kernel* around each sample
4. And *evaluate their sums* in a dense grid

Kernel Density Estimation!



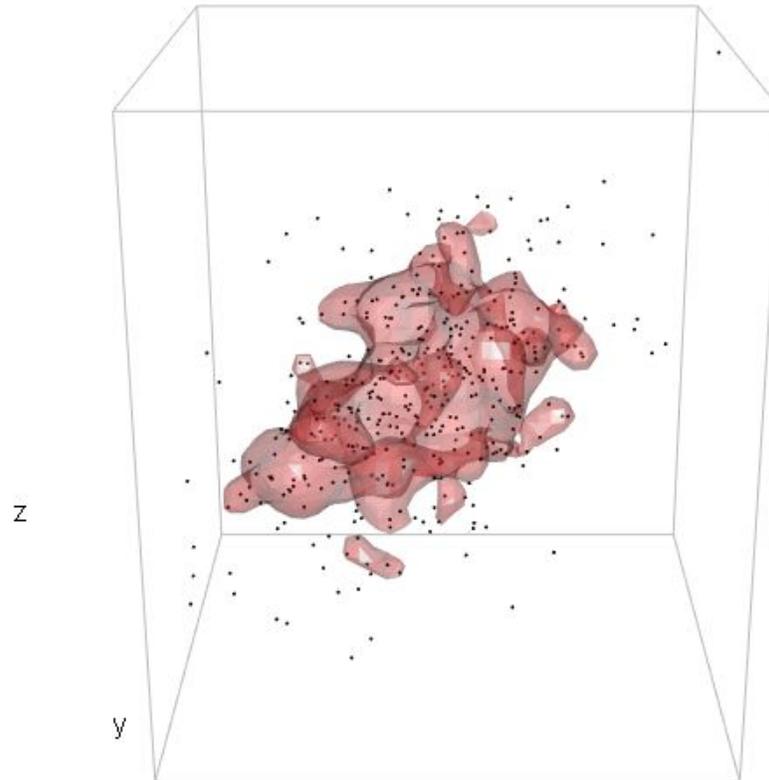
More Dimensions!

1. The basic principle is easily extended to **2D**



More Dimensions!

1. The basic principle is easily extended to **2D**
2. and **3D**



Slow...

1. Naive algorithm is **quite slow**:
 - a. For every sample
 - i. For every evaluated grid cell
 1. Accumulate the value of the kernel



1. Naive algorithm is **quite slow**
2. But we can **parallelize**, across:
 - a. Each sample
 - b. Each grid cell a sample touches



Slow...

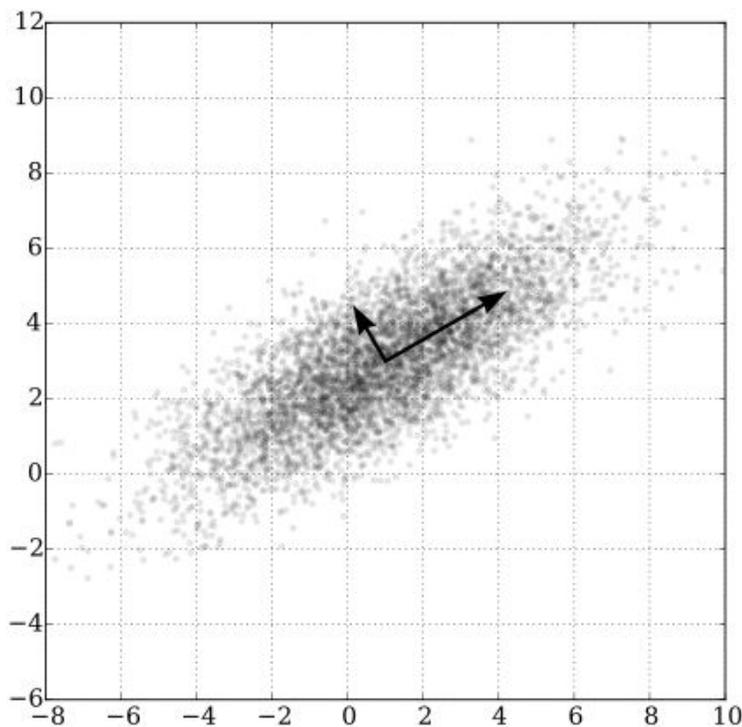
1. Naive algorithm is **quite slow**
2. But we can **parallelize**, across:
 - a. Each sample
 - b. Each grid cell a sample touches
3. And we can precompute exactly which cells will be affected*

**Lopez-Novoa et al. 2016*

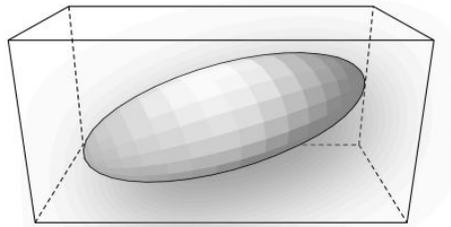


How to figure out the 3D ellipsoid kernel?

1. Covariance matrix of the data encodes the “ellipsoid” shape of the data
2. Ellipsoid Matrix with the shape of the data is the ***inverse of the covariance***



1. For each sample, calculate its 3D bounding box



Chop & Crop

1. For each sample, calculate its 3D bounding box
2. **Chop the ellipsoid kernel** along the Z axis

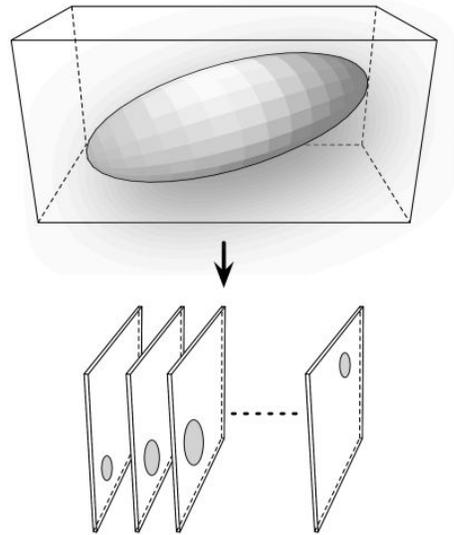


Fig. 2 Chopping a 3D bounding box into 2D slices



Chop & Crop

1. For each sample, calculate its 3D bounding box
2. **Chop the ellipsoid kernel** along the Z axis
3. **Crop each 2D slice** so that only the cells inside the kernel are considered

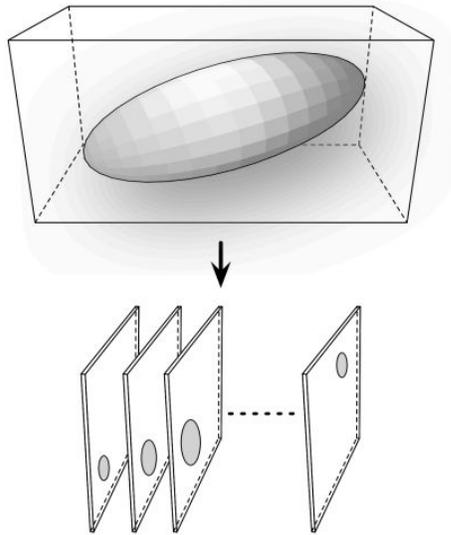


Fig. 2 Chopping a 3D bounding box into 2D slices

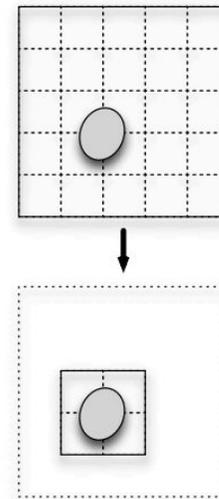


Fig. 3 Cropping a 2D slice to obtain a minimum-size bounding rectangle



Chop & Crop

1. For each sample, calculate its 3D bounding box
2. **Chop the ellipsoid kernel** along the Z axis
3. **Crop each 2D slice** so that only the cells inside the kernel are considered
4. Launch a parallel thread **for each cell of each cropped slice for each sample**

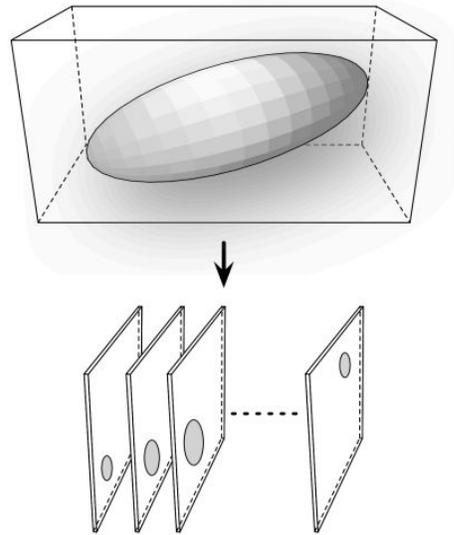


Fig. 2 Chopping a 3D bounding box into 2D slices

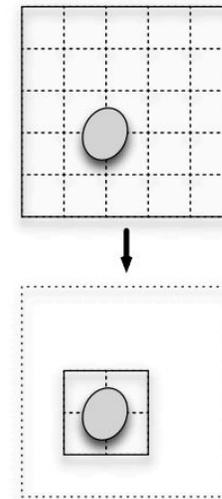
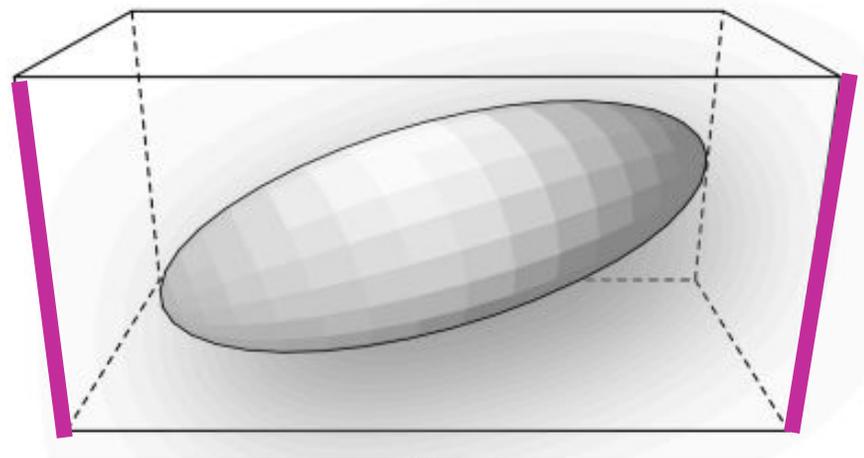


Fig. 3 Cropping a 2D slice to obtain a minimum-size bounding rectangle

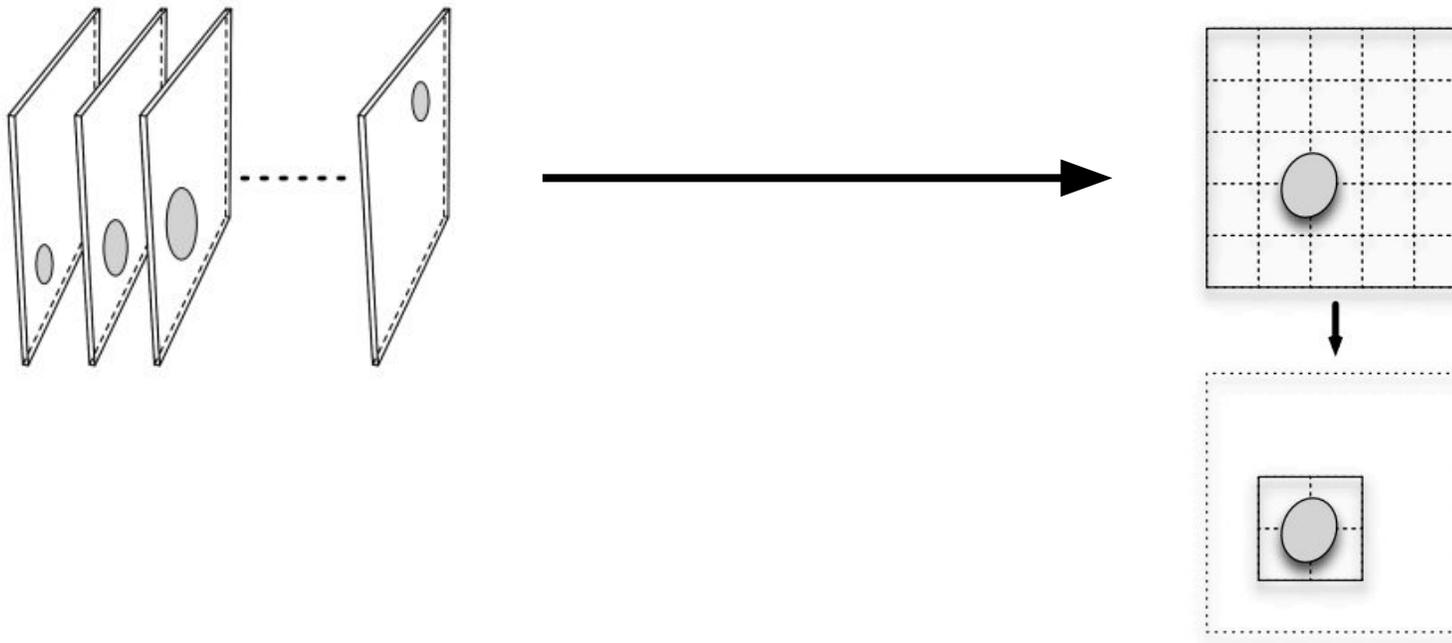


1. CUDA *thread considers a single sample*
2. Stores the *ellipsoid kernel boundaries* along the **Z** dimension
3. No race conditions



Implementation - *Crop Kernel*

1. CUDA *thread* considers a *single slice of a single sample*
2. Stores the *X and Y boundaries* within that slice
3. No race conditions



Implementation - *Crop Kernel*

1. CUDA *thread considers a single slice of a single sample*
2. Stores the *X and Y boundaries* within that slice
3. No race conditions

Mathematical Interlude



Implementation - Crop Kernel

1. CUDA thread considers a single slice of a single sample
2. Stores the X and Y boundaries within that slice
3. No race conditions

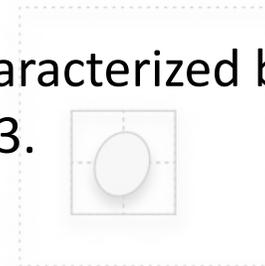
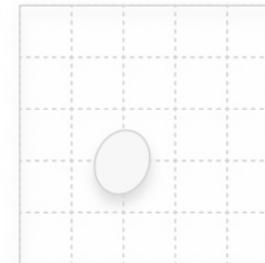
Mathematical Interlude

The general form of a rotated, translated ellipse is a quadratic equation: $Ax^2 + Bxy + Cy^2 + Dx + Ey + F < 1$

This can be represented with matrices:

$$\begin{bmatrix} x & y & 1 \end{bmatrix} * \begin{bmatrix} A & B/2 & D/2 \\ B/2 & C & E/2 \\ D/2 & E/2 & F \end{bmatrix} * \begin{bmatrix} x \\ y \\ 1 \end{bmatrix}$$

...or $x^T A x < 1$ for short



This all generalizes to higher dimensions. Our ellipsoid kernel is characterized by a 4x4 matrix, except that because it has no translation, we can get away with 3x3.



Implementation - *Crop Kernel*

1. CUDA *thread considers a single slice of a single sample*
2. Stores the *X and Y boundaries* within that slice
3. No race conditions

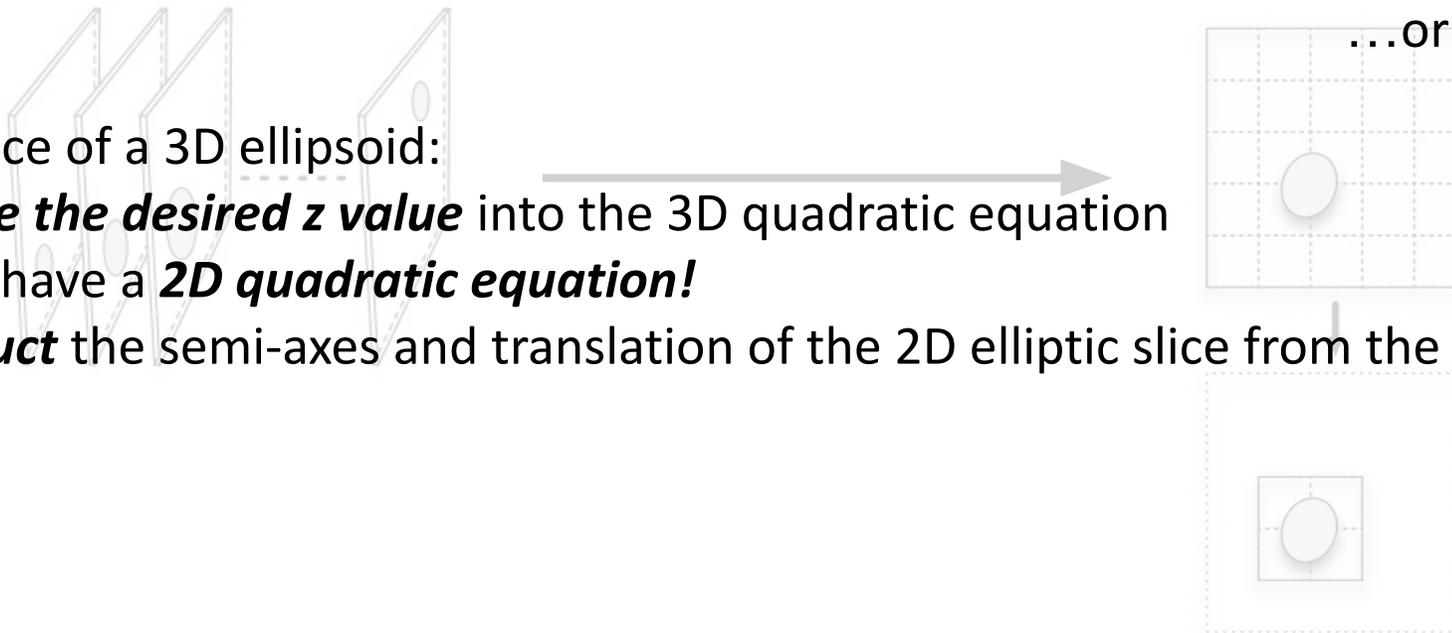
Mathematical Interlude

The general form of a rotated, translated ellipse is a quadratic equation: $Ax^2 + Bxy + Cy^2 + Dx + Ey + F < 1$

...or $x^T A x < 1$ for short

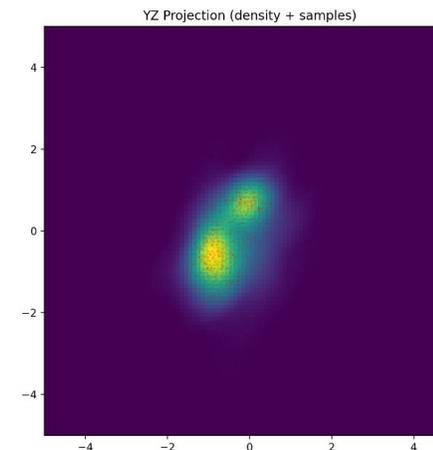
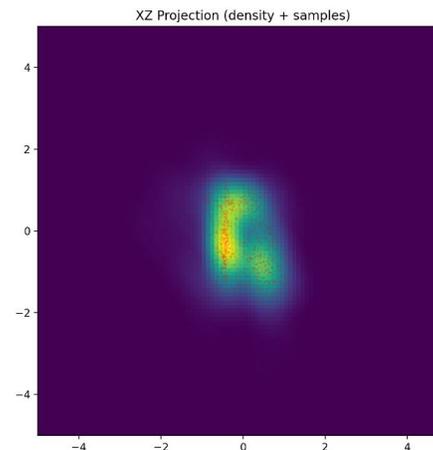
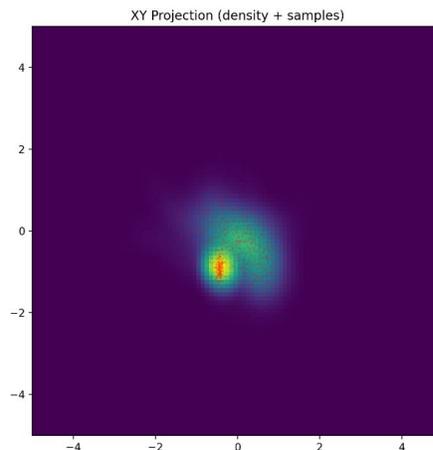
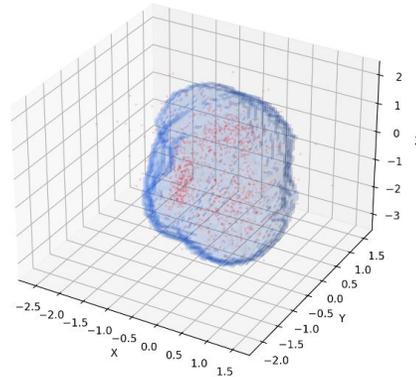
To get a 2D slice of a 3D ellipsoid:

1. **Substitute the desired z value** into the 3D quadratic equation
2. Now you have a **2D quadratic equation!**
3. **Reconstruct** the semi-axes and translation of the 2D elliptic slice from the quadratic coefficients



Implementation - *KDE Kernel*

1. CUDA *thread considers a single cell affected by a single sample*
2. Checks the value of the ellipsoid kernel
3. And *atomically accumulates the result* in the final grid



1. In the original paper, densities were accumulated on the CPU, which was a serious ***bottleneck***

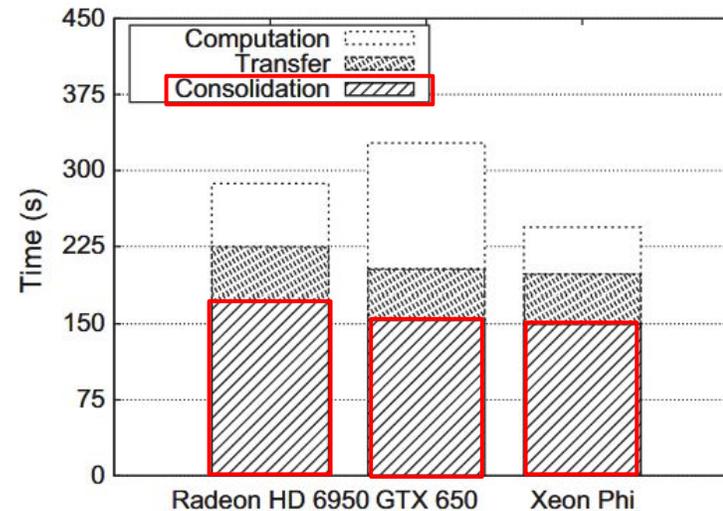


Fig. 8 Dissected execution time of OpenCL S-KDE. 1M dataset and 194M evaluation points



1. In the original paper, densities were accumulated on the CPU, which was a serious ***bottleneck***
2. By doing everything on the GPU, ***we can visualize the results directly***



1. In the original paper, densities were accumulated on the CPU, which was a serious ***bottleneck***
2. By doing everything on the GPU, ***we can visualize the results directly***
3. We can also visualize ***the algorithm itself!***

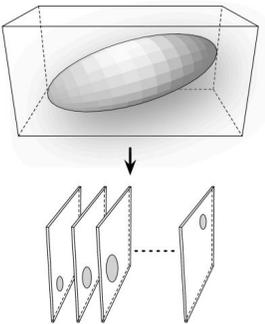


Fig. 2 Chopping a 3D bounding box into 2D slices

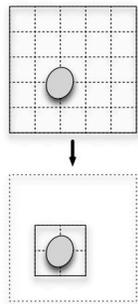


Fig. 3 Cropping a 2D slice to obtain a minimum-size bounding rectangle



1. In the original paper, densities were accumulated on the CPU, which was a serious ***bottleneck***
2. By doing everything on the GPU, ***we can visualize the results directly***
3. We can also visualize ***the algorithm itself!***

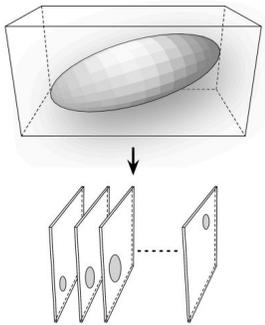


Fig. 2 Chopping a 3D bounding box into 2D slices

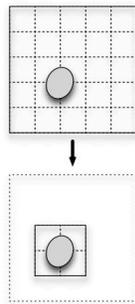
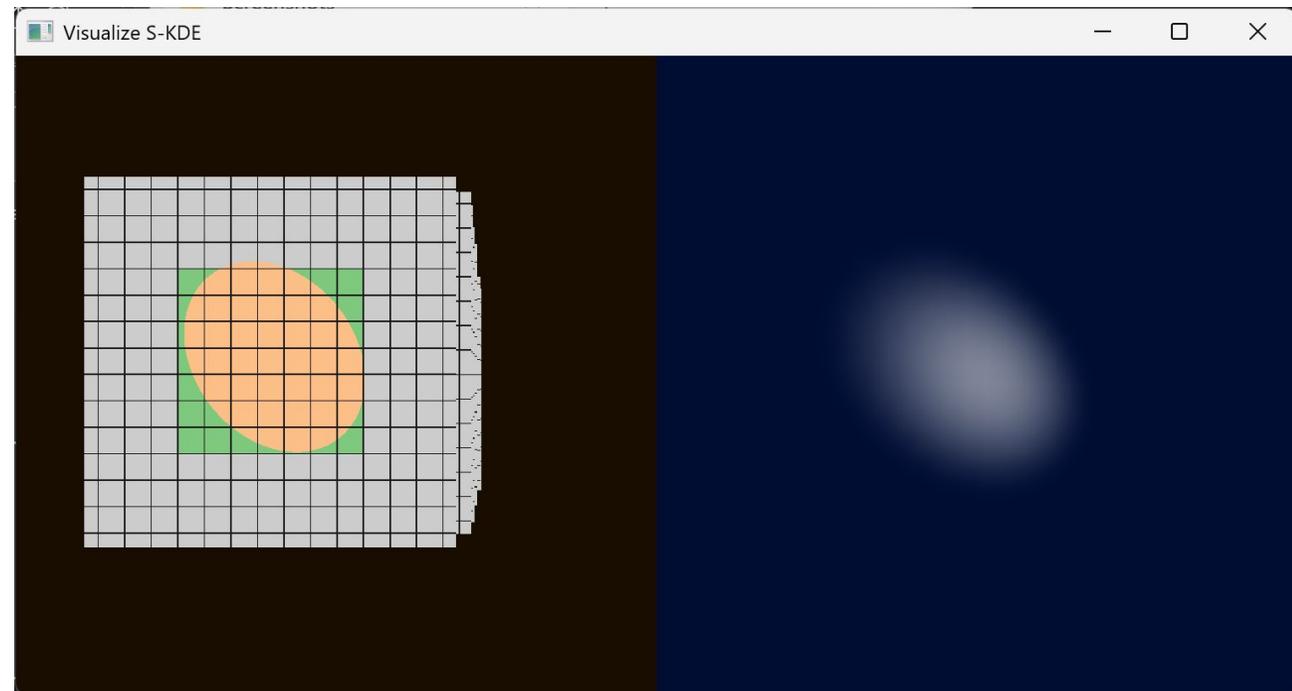


Fig. 3 Cropping a 2D slice to obtain a minimum-size bounding rectangle



Implementation - *Visualization*

1. In the original paper, densities were accumulated on the CPU, which was a serious ***bottleneck***
2. By doing everything on the GPU, ***we can visualize the results directly***
3. We can also visualize ***the algorithm itself!***

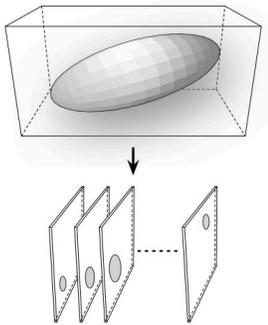


Fig. 2 Chopping a 3D bounding box into 2D slices

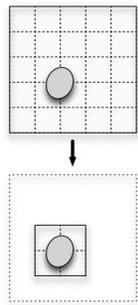
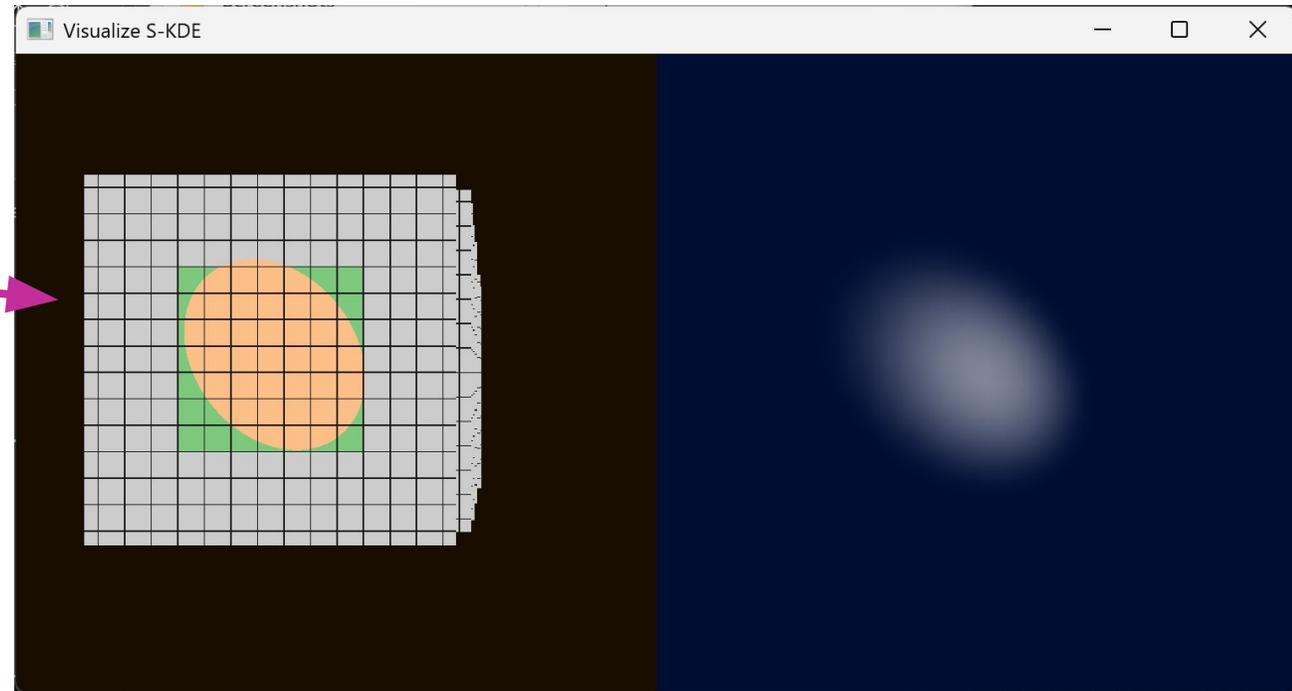


Fig. 3 Cropping a 2D slice to obtain a minimum-size bounding rectangle



Implementation - *Visualization*

1. In the original paper, densities were accumulated on the CPU, which was a serious ***bottleneck***
2. By doing everything on the GPU, ***we can visualize the results directly***
3. We can also visualize ***the algorithm itself!***

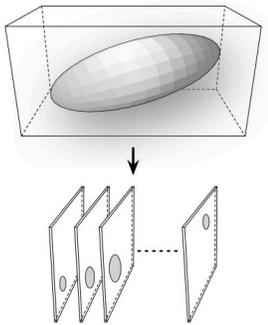


Fig. 2 Chopping a 3D bounding box into 2D slices

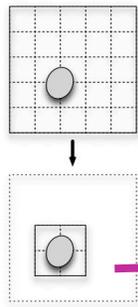
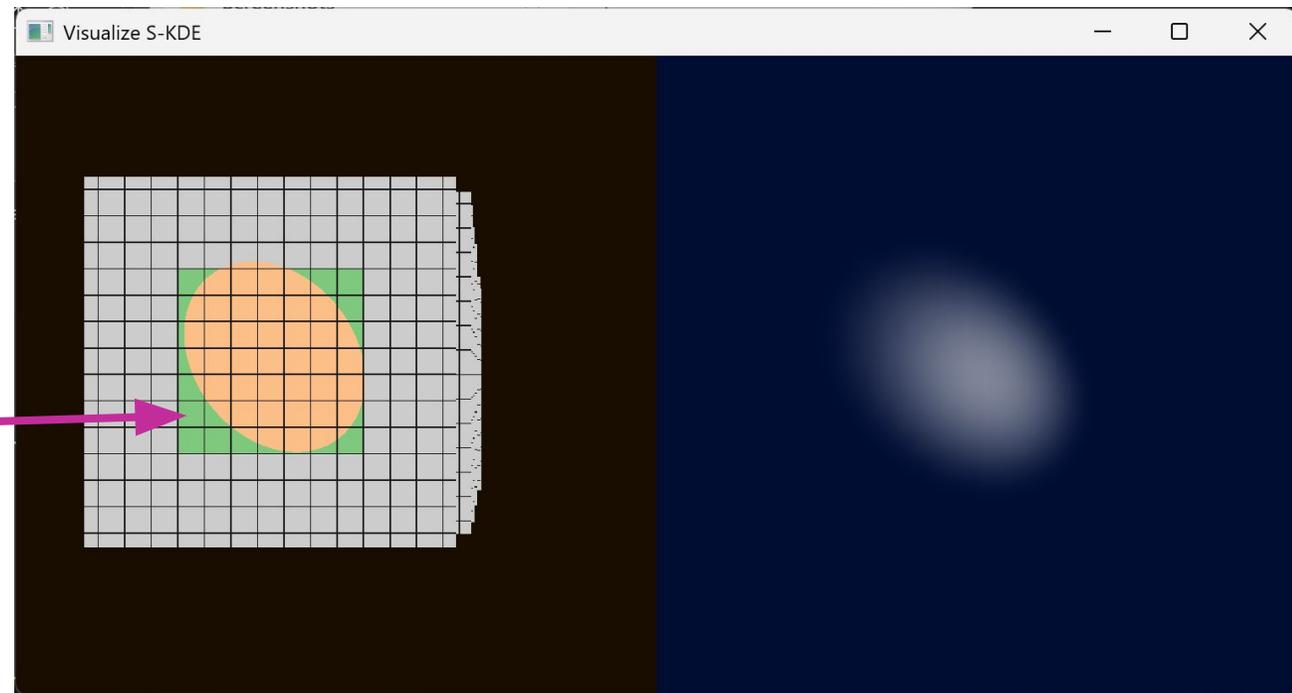


Fig. 3 Cropping a 2D slice to obtain a minimum-size bounding rectangle



1. In the original paper, densities were accumulated on the CPU, which was a serious ***bottleneck***
2. By doing everything on the GPU, ***we can visualize the results directly***
3. We can also visualize ***the algorithm itself!***

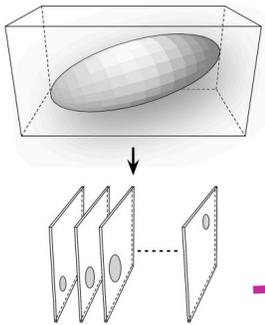


Fig. 2 Chopping a 3D bounding box into 2D slices

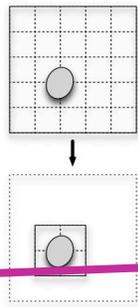
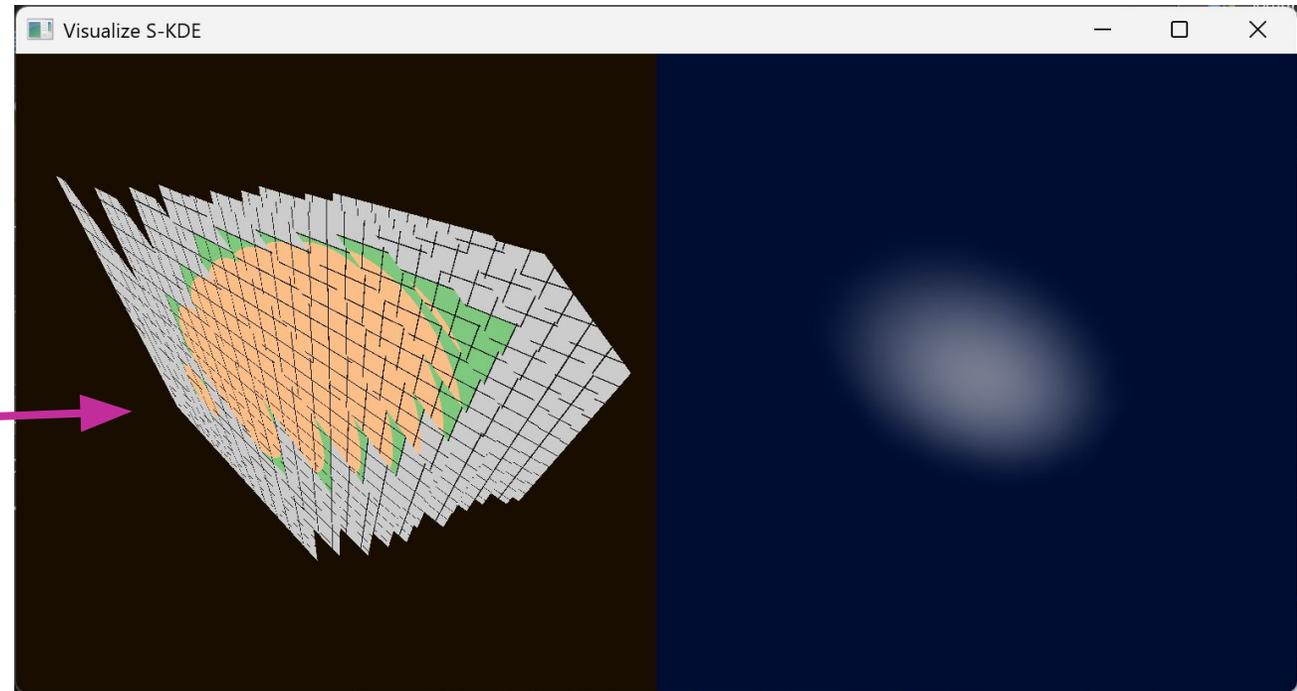


Fig. 3 Cropping a 2D slice to obtain a minimum-size bounding rectangle



1. In the original paper, densities were accumulated on the CPU, which was a serious ***bottleneck***
2. By doing everything on the GPU, ***we can visualize the results directly***
3. We can also visualize ***the algorithm itself!***

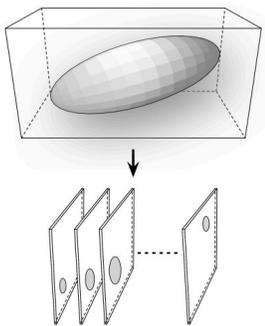


Fig. 2 Chopping a 3D bounding box into 2D slices

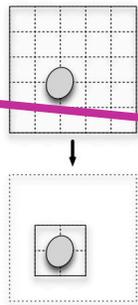
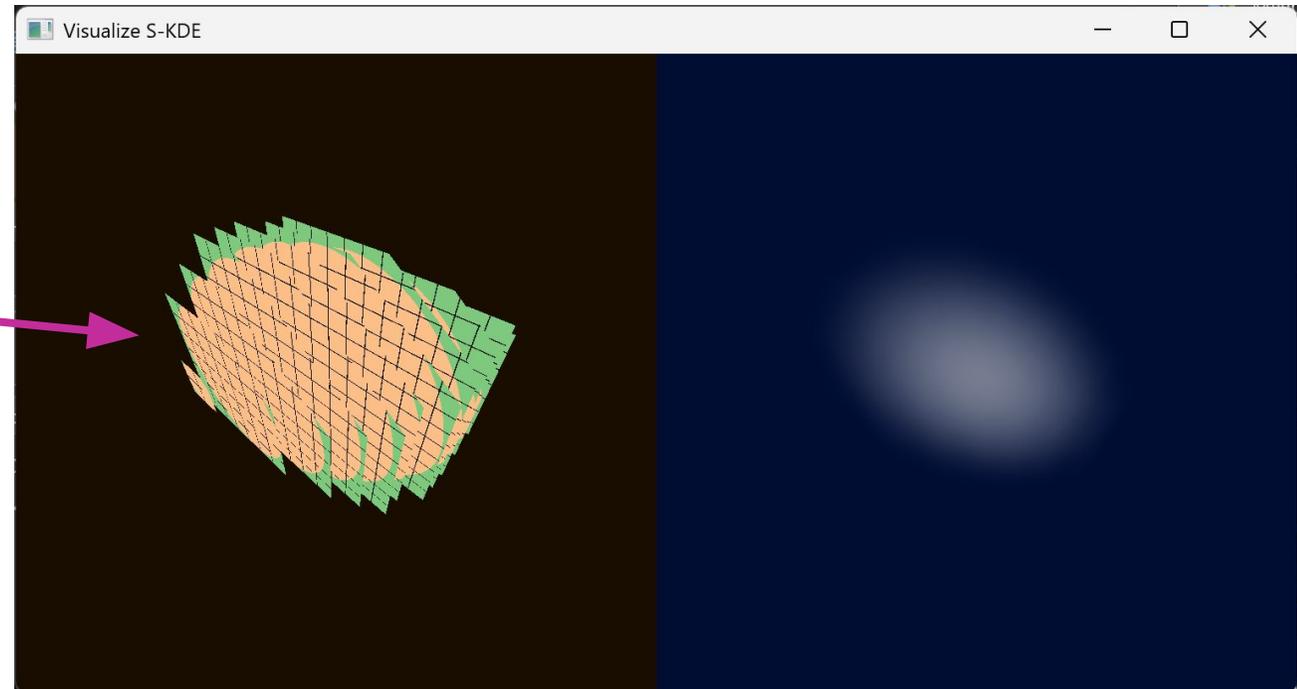
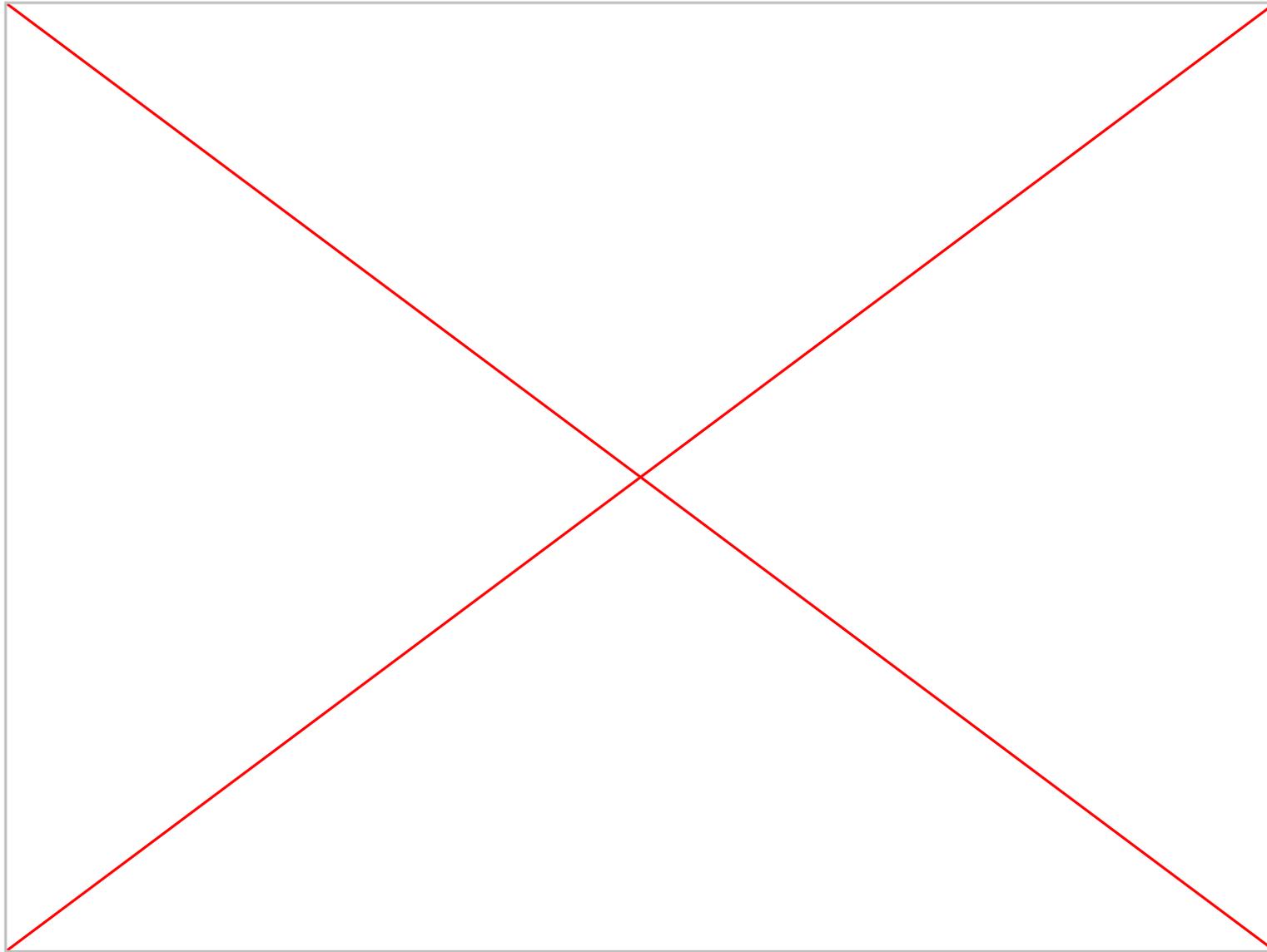


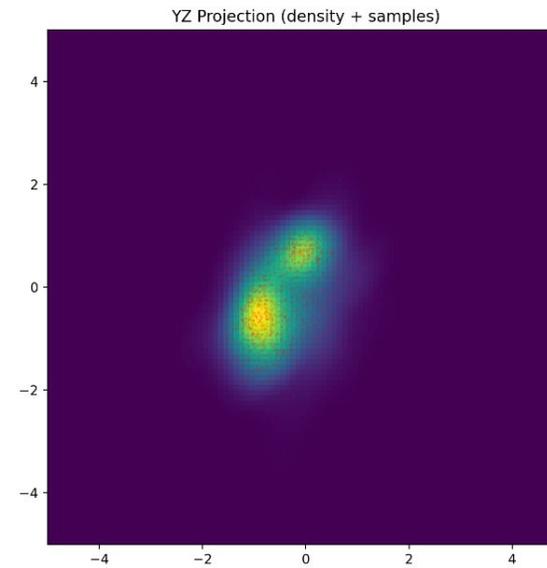
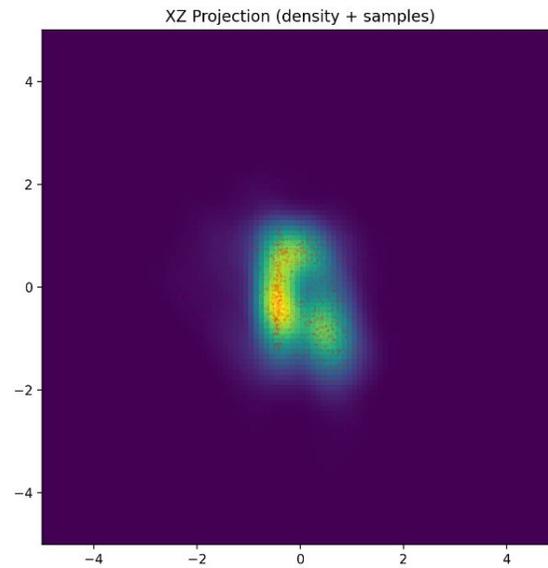
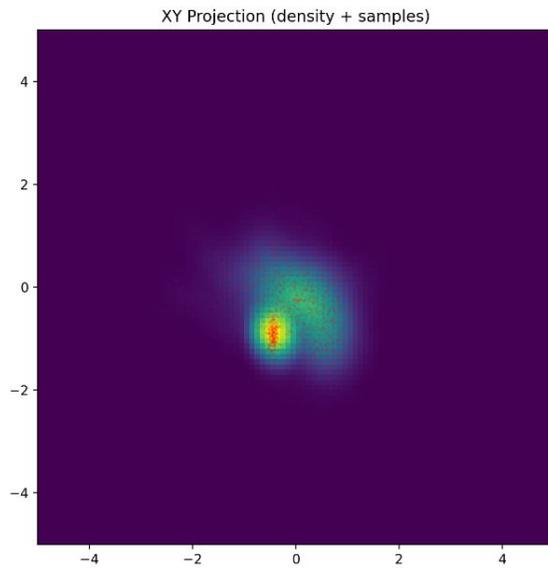
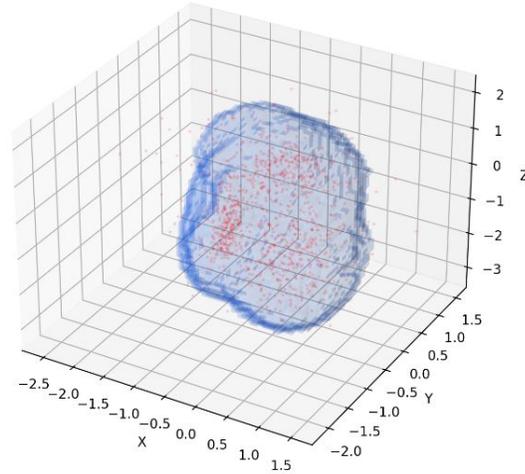
Fig. 3 Cropping a 2D slice to obtain a minimum-size bounding rectangle



Implementation - *Visualization*



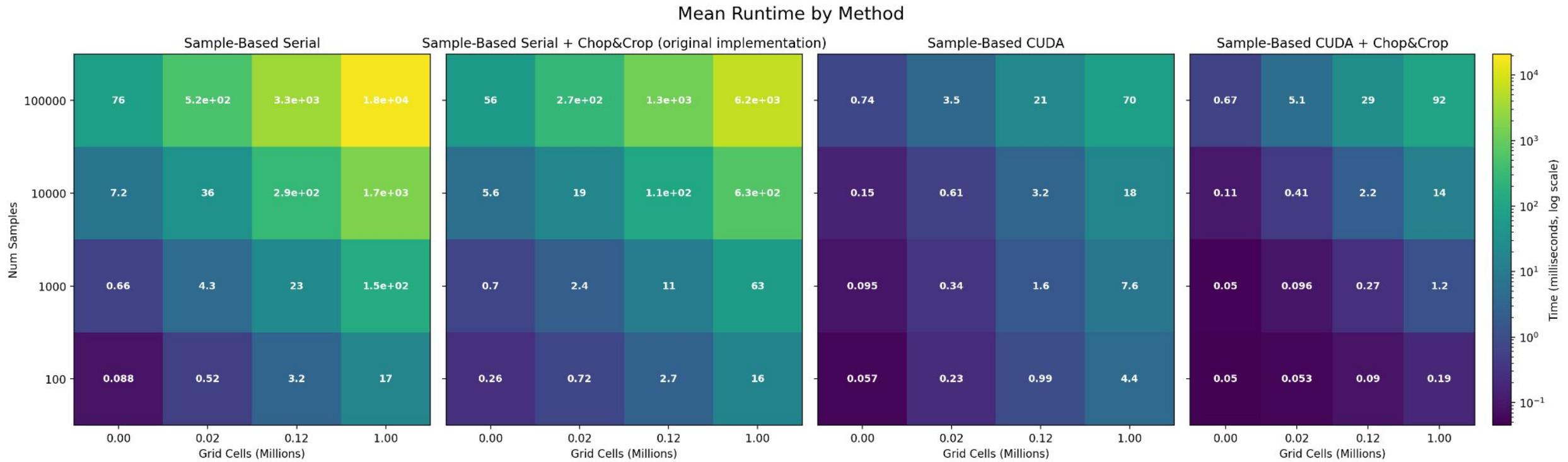
Example Result



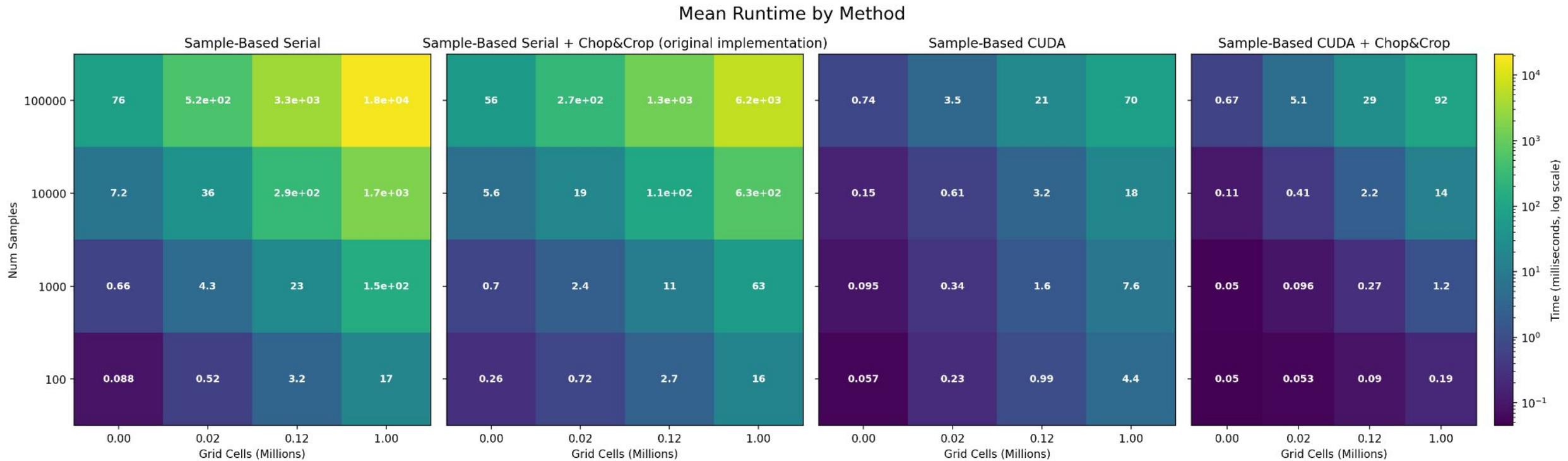
1. Compare timing against:
 - a. CPU, serial
 - b. CPU, serial + Chop&Crop (Original Reference Impl.)
 - c. CUDA, parallel sample-wise
 - d. CUDA, parallel sample-and-cell-wise + Chop&Crop
2. Generate samples from a ***set of 3D Gaussians***
3. 3 repetitions per testing condition
4. 3 warmup runs for CUDA



Experiments and Profiles



Experiments and Profiles



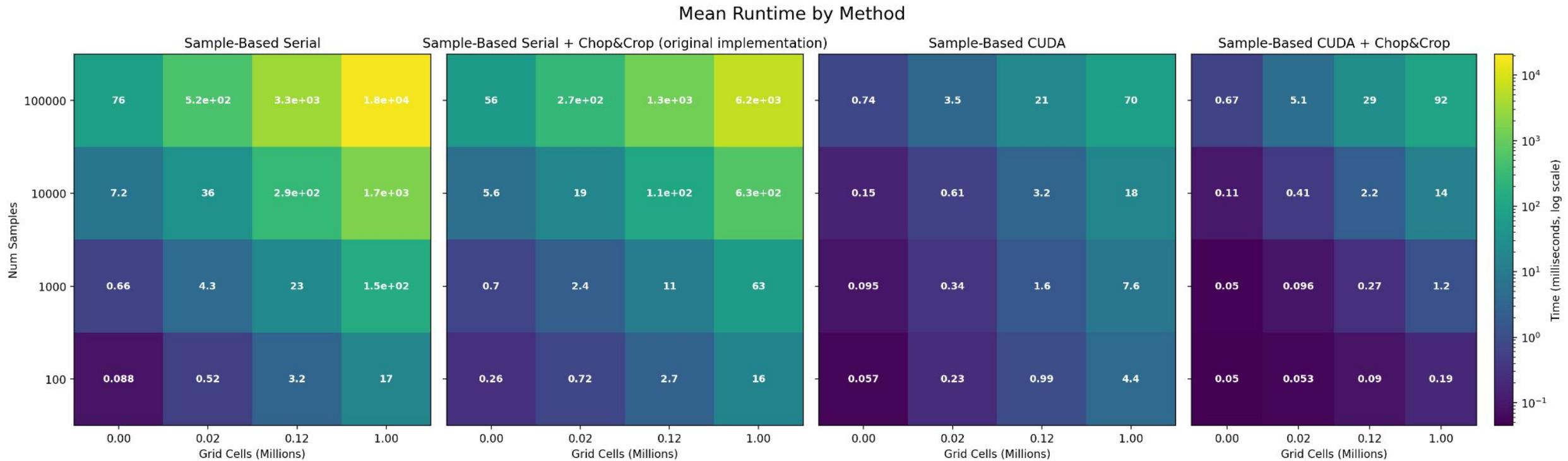
The more the cells, the slower we get



The more the samples, the slower we get



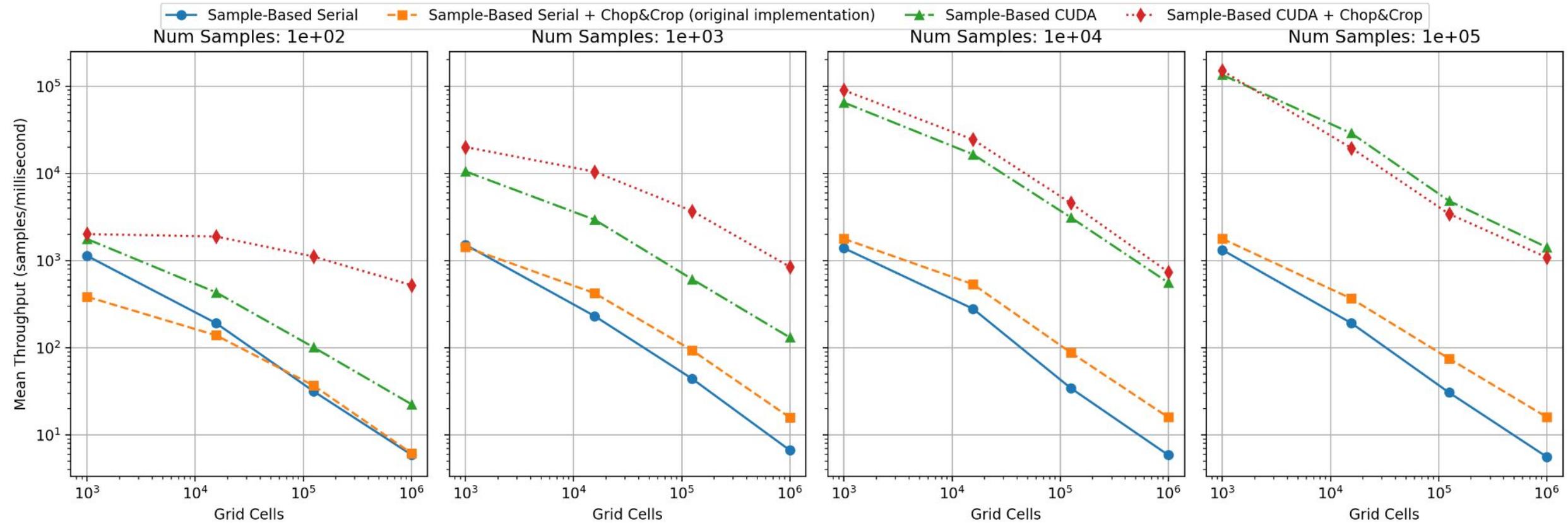
Experiments and Profiles



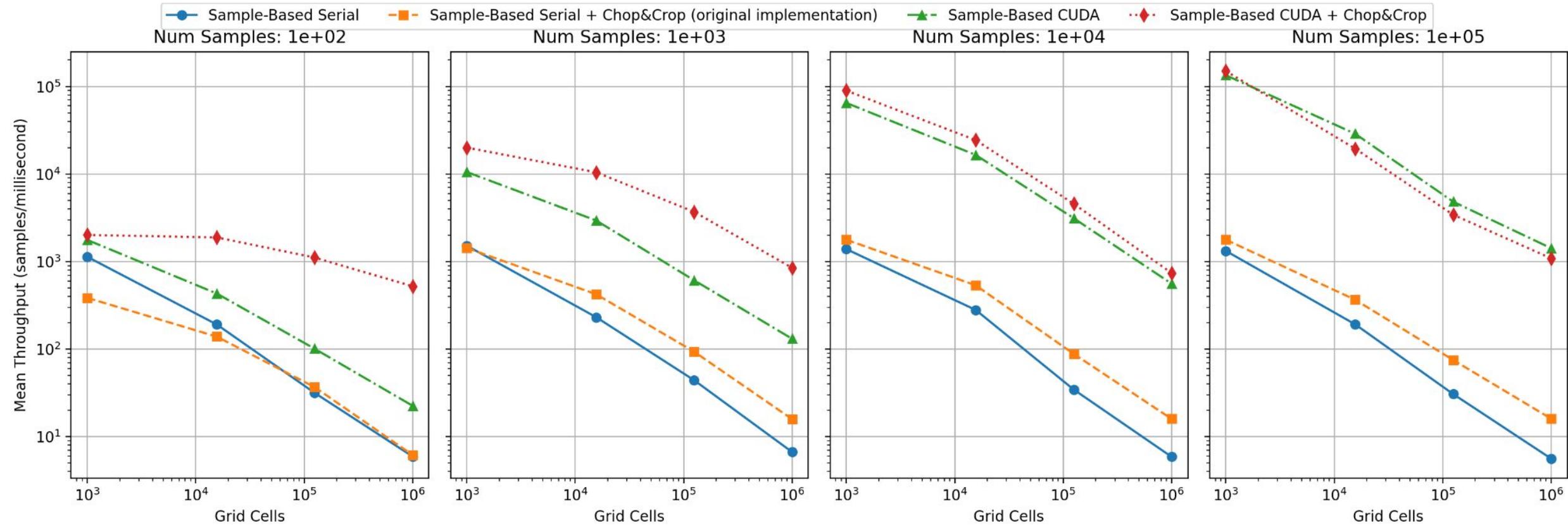
Both Chop&Crop and parallelism result in speedups



Experiments and Profiles



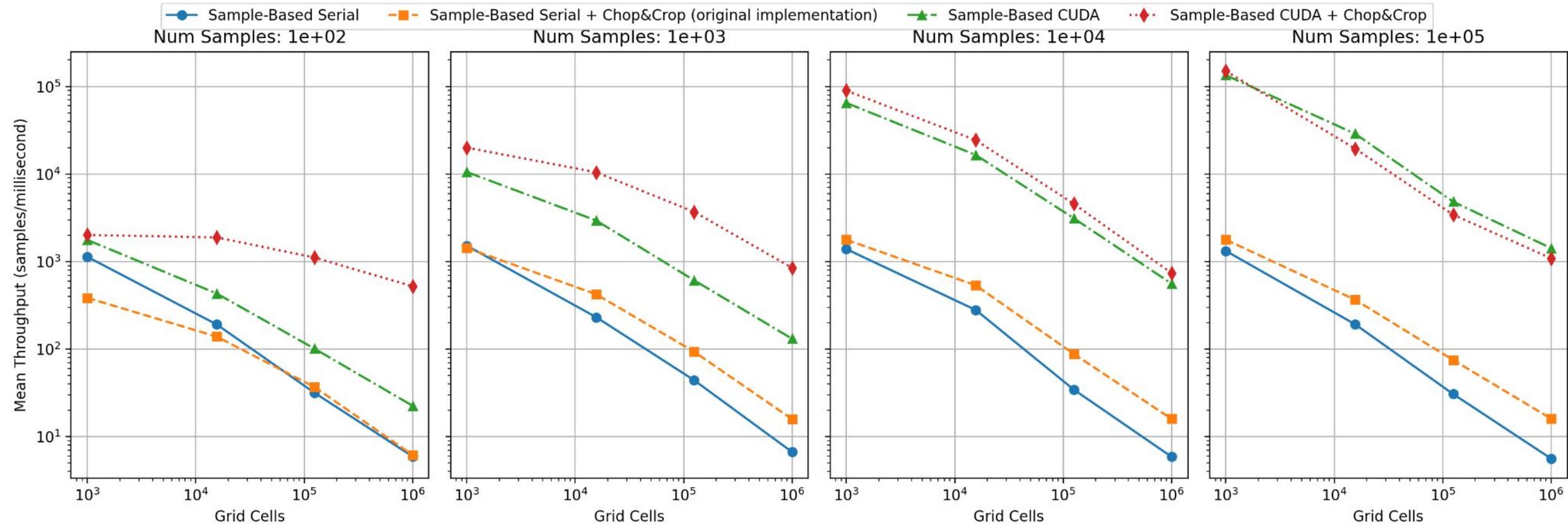
Experiments and Profiles



Chop&Crop improvements decay for large number of samples



Experiments and Profiles



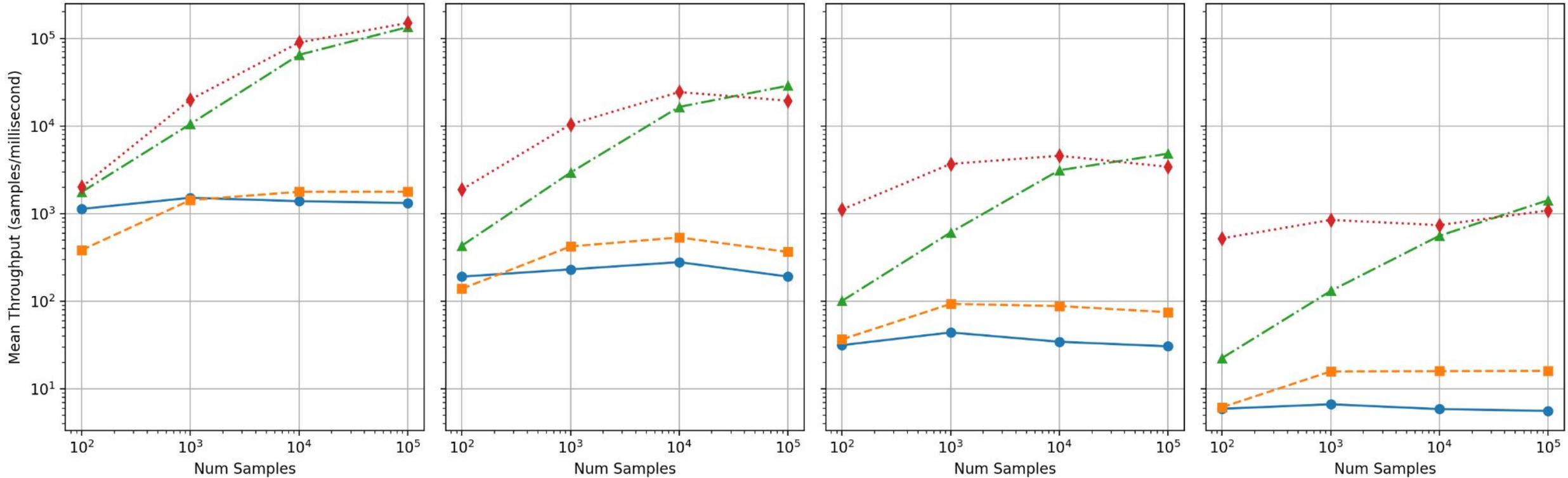
All methods suffer similar slowdown with the number of grid cells



Experiments and Profiles

Legend:
● Sample-Based Serial
■ Sample-Based Serial + Chop&Crop (original implementation)
▲ Sample-Based CUDA
◆ Sample-Based CUDA + Chop&Crop

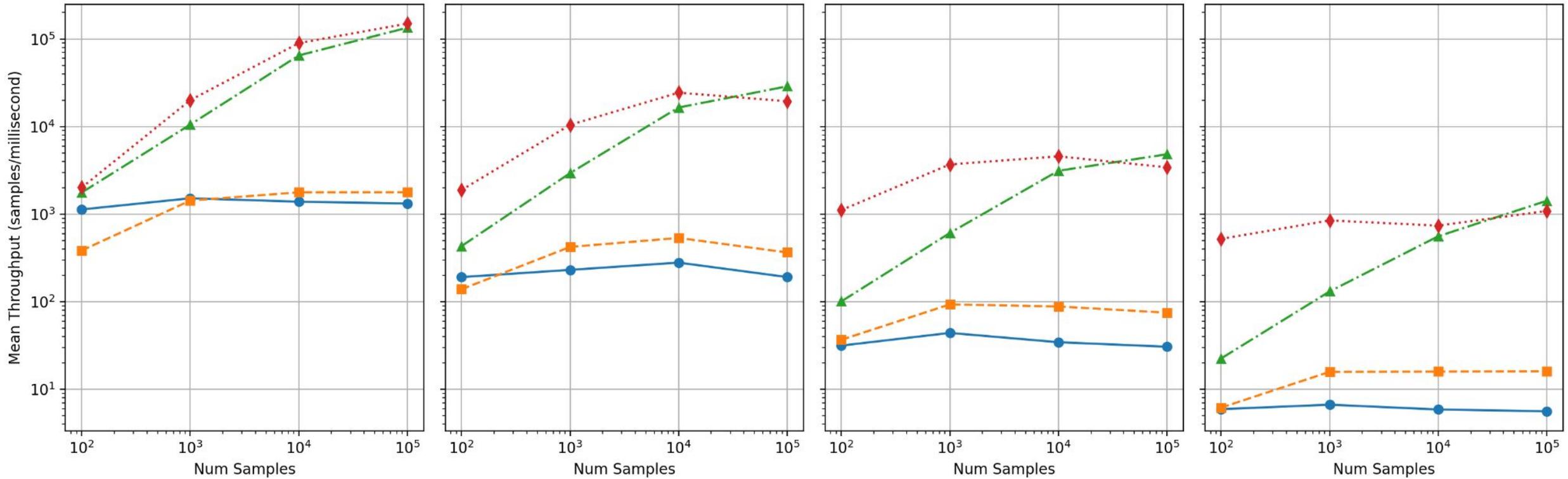
Grid Cells: $1e+03$ Grid Cells: $1.56e+04$ Grid Cells: $1.25e+05$ Grid Cells: $1e+06$



Experiments and Profiles

Legend:
● Sample-Based Serial
■ Sample-Based Serial + Chop&Crop (original implementation)
▲ Sample-Based CUDA
◆ Sample-Based CUDA + Chop&Crop

Grid Cells: $1e+03$ Grid Cells: $1.56e+04$ Grid Cells: $1.25e+05$ Grid Cells: $1e+06$

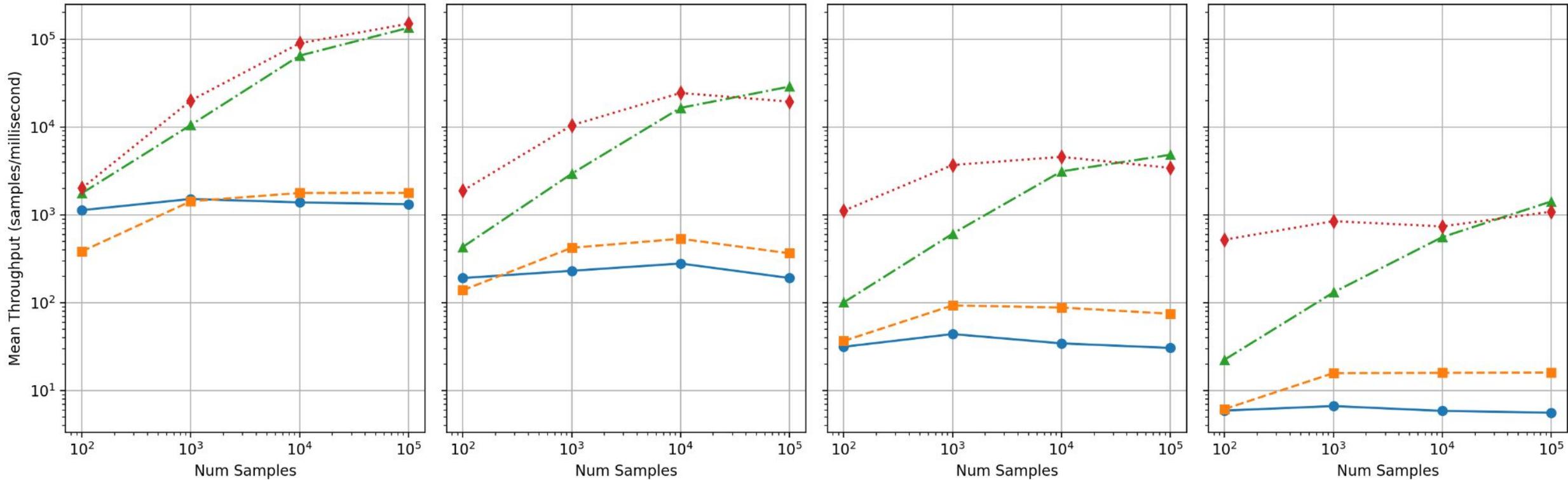


Serial Methods maintain its speed regardless of the dataset size



Experiments and Profiles

Legend: Sample-Based Serial (blue solid line with circles), Sample-Based Serial + Chop&Crop (original implementation) (orange dashed line with squares), Sample-Based CUDA (green dash-dot line with triangles), Sample-Based CUDA + Chop&Crop (red dotted line with diamonds).
Grid Cells: $1e+03$ Grid Cells: $1.56e+04$ Grid Cells: $1.25e+05$ Grid Cells: $1e+06$



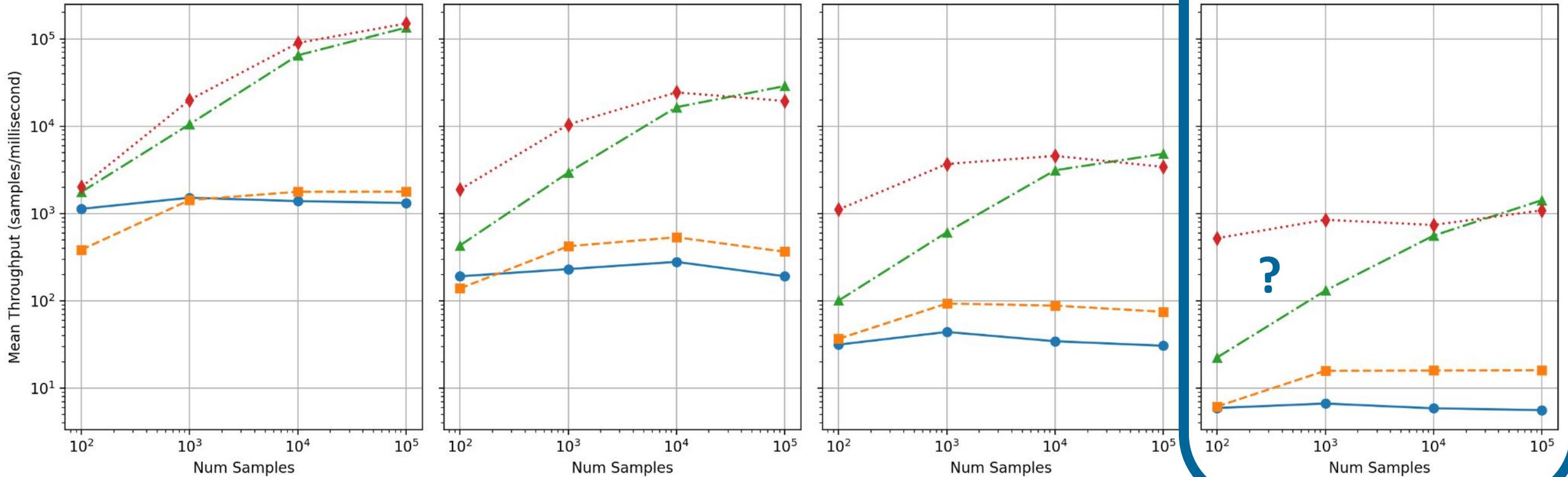
Parallel methods shine when given lots of samples



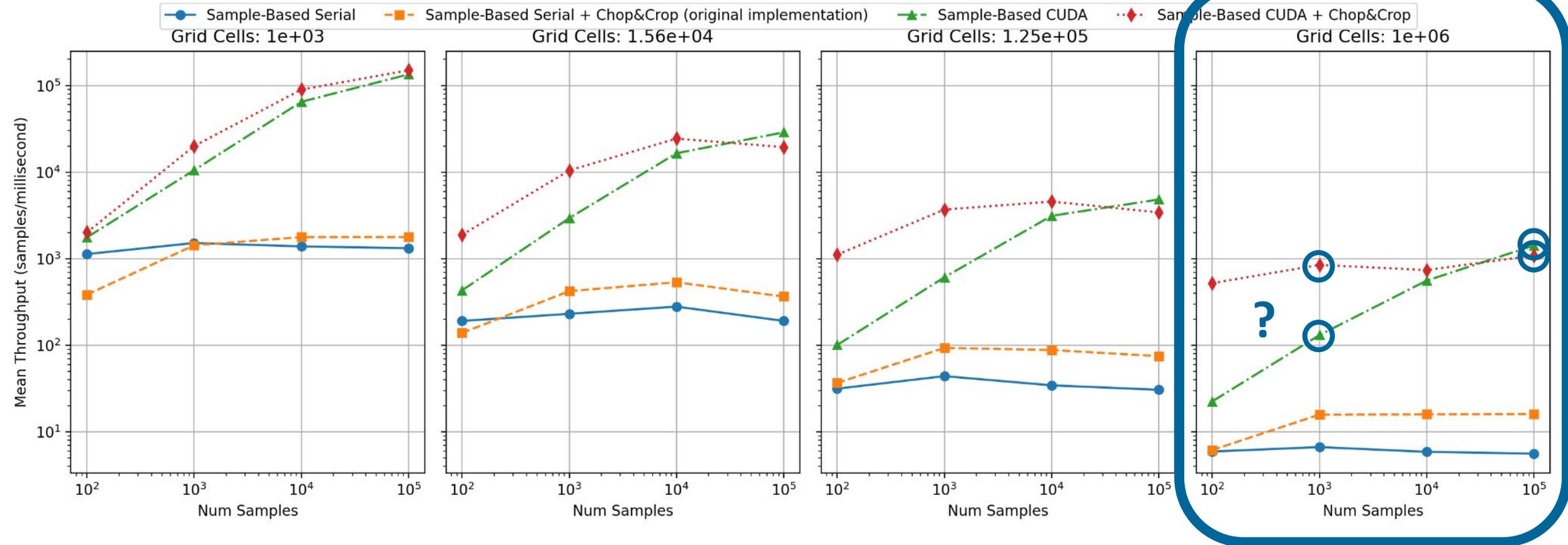
Experiments and Profiles

Legend:
● Sample-Based Serial
■ Sample-Based Serial + Chop&Crop (original implementation)
▲ Sample-Based CUDA
◆ Sample-Based CUDA + Chop&Crop

Grid Cells: $1e+03$ Grid Cells: $1.56e+04$ Grid Cells: $1.25e+05$ Grid Cells: $1e+06$



Experiments and Profiles



Sample-wise CUDA + Chop&Crop

Function Name	Demangled Name	Duration (1.01009e+09)
chopKernel	..chopKernel..Slice...	0.02
cropKernel	..cropKernel..Ellipse...	1.87
writeDensitiesKernel	..writeDensitiesKer...	251.62



L2 Cache Hit Rate (%)

	Naive CUDA	CUDA SKDE (<i>Final Kernel</i>)
1K samples	82.40	97.91
1M samples	98.89	98.67



Slowest Code Lines

	Naive CUDA	CUDA SKDE (<i>Final Kernel</i>)
1K samples	98% Location sm_60_atomic_functions.hpp:79 (0x500c5a6b0 in kde_kernel) ↗ sm_60_atomic_functions.hpp:79 (0x500c5a3f0 in kde_kernel) ↗ sm_60_atomic_functions.hpp:79 (0x500c5a970 in kde_kernel) ↗ sm_60_atomic_functions.hpp:79 (0x500c5aab0 in kde_kernel) ↗	~100% Location cuda_skde.cuh:169 (0x500c61bc0 in writeDensitiesKernel) ↗
1M samples	96% sm_60_atomic_functions.hpp:79 (0x500c5a3f0 in kde_kernel) ↗ sm_60_atomic_functions.hpp:79 (0x500c5a6b0 in kde_kernel) ↗ sm_60_atomic_functions.hpp:79 (0x500c5a970 in kde_kernel) ↗ sm_60_atomic_functions.hpp:79 (0x500c5aab0 in kde_kernel) ↗	~100% Location cuda_skde.cuh:169 (0x500c61bc0 in writeDensitiesKernel) ↗

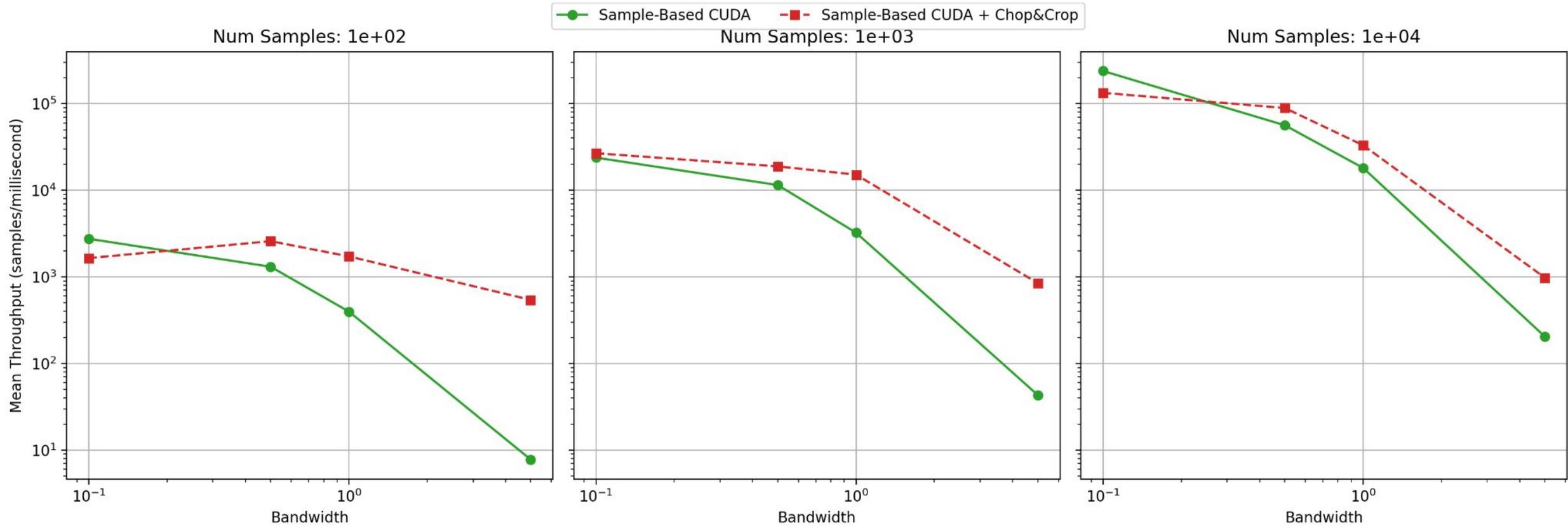


Slowest Code Lines

	Naive CUDA	CUDA SKDE (<i>Final Kernel</i>)
1K samples	98%	<p><code>atomicAdd(&pde[evaluationId], density);</code></p> <p>Location cuda_skde.cuh:169 (0x500c61bc0 in writeDensitiesKernel)</p>
1M samples	96%	<p>~100%</p> <p>Location cuda_skde.cuh:169 (0x500c61bc0 in writeDensitiesKernel)</p>



Experiments and Profiles



Compute Throughput (%)

	Naive CUDA	CUDA SKDE (<i>Final Kernel</i>)
1K samples	9.58	85.39
1M samples	85.04	85.98



Achieved Occupancy (Theoretical) (%)

	Naive CUDA	CUDA SKDE (<i>Final Kernel</i>)
1K samples	16.59 (67)	62.56 (100)
1M samples	59.02 (67)	66.84 (100)



Achieved Occupancy (Theoretical) (%)

	Naive CUDA	CUDA SKDE (<i>Final Kernel</i>)
1K samples	16.59 (67)	62.56 (100)
1M samples	59.02 (67)	66.84 (100)

Difference likely due to thread-per-sample vs thread-per-sample-cell



1. Faster, parallel, smarter density accumulation

```
atomicAdd(&pde[evaluationId], density);
```



1. Faster, parallel, smarter density accumulation
2. Improved, coalesced data access for both samples and evaluation points

Uncoalesced Global Accesses
Est. Speedup: 45.46%



Lopez-Novoa, Unai and Mendiburu, Alexander and Miguel-Alonso, Jose

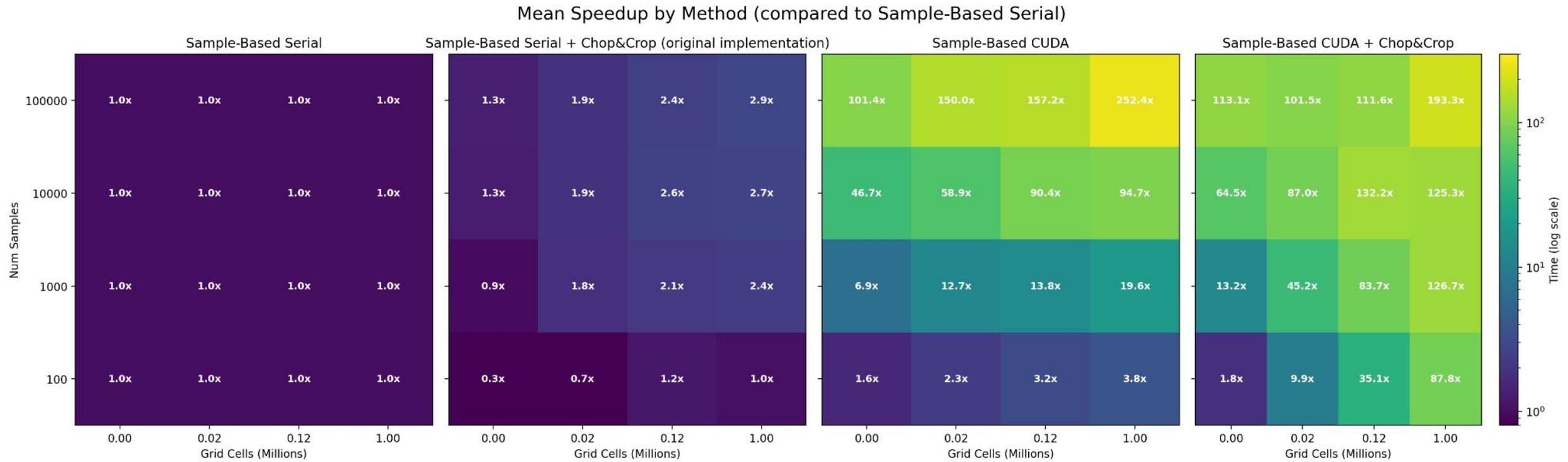
Kernel density estimation in accelerators: Implementation and performance evaluation

The Journal of Supercomputing Vol. 72, N. 2

2016

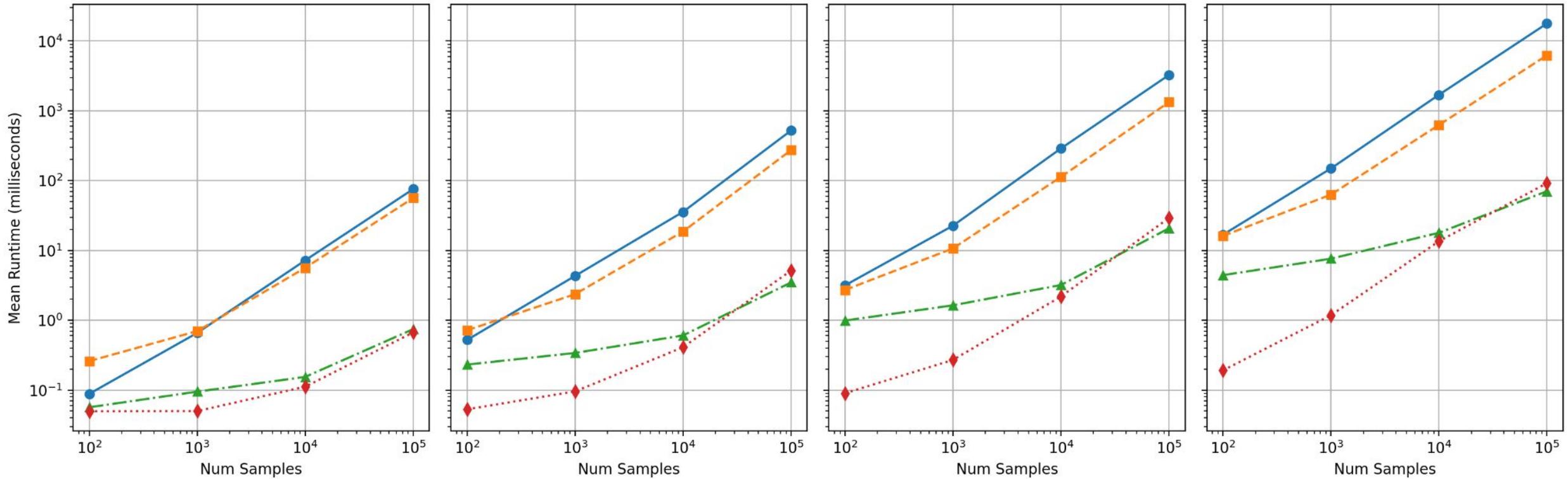


Extra Plots



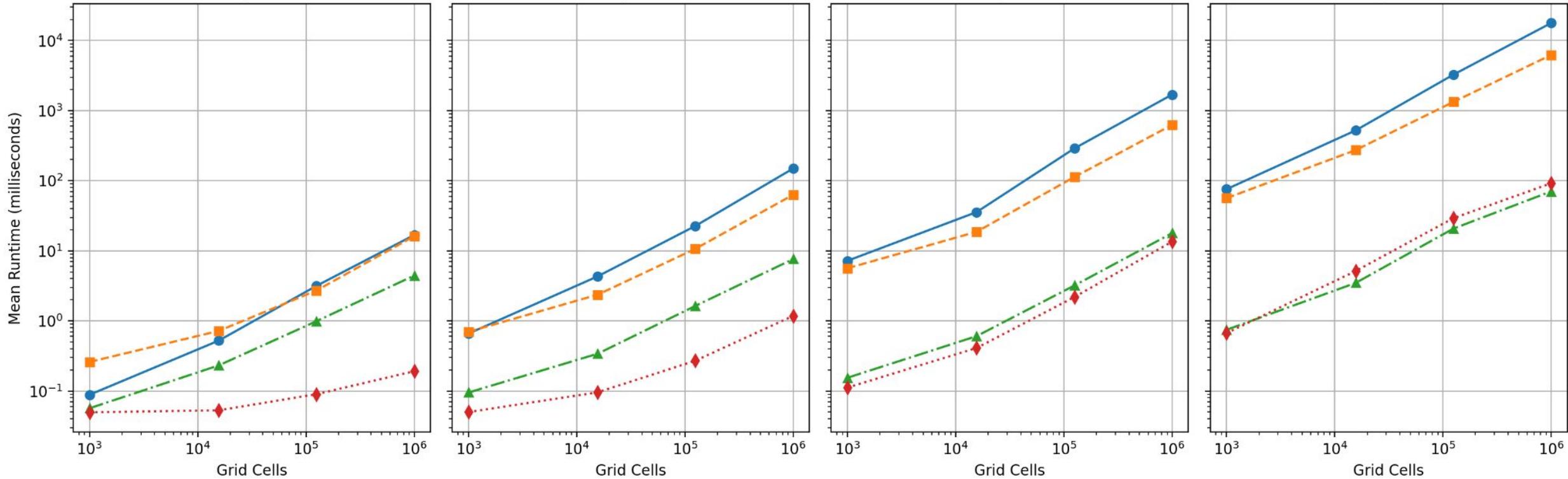
Extra Plots

Legend:
● Sample-Based Serial
■ Sample-Based Serial + Chop&Crop (original implementation)
▲ Sample-Based CUDA
◆ Sample-Based CUDA + Chop&Crop
Grid Cells: $1e+03$ Grid Cells: $1.56e+04$ Grid Cells: $1.25e+05$ Grid Cells: $1e+06$

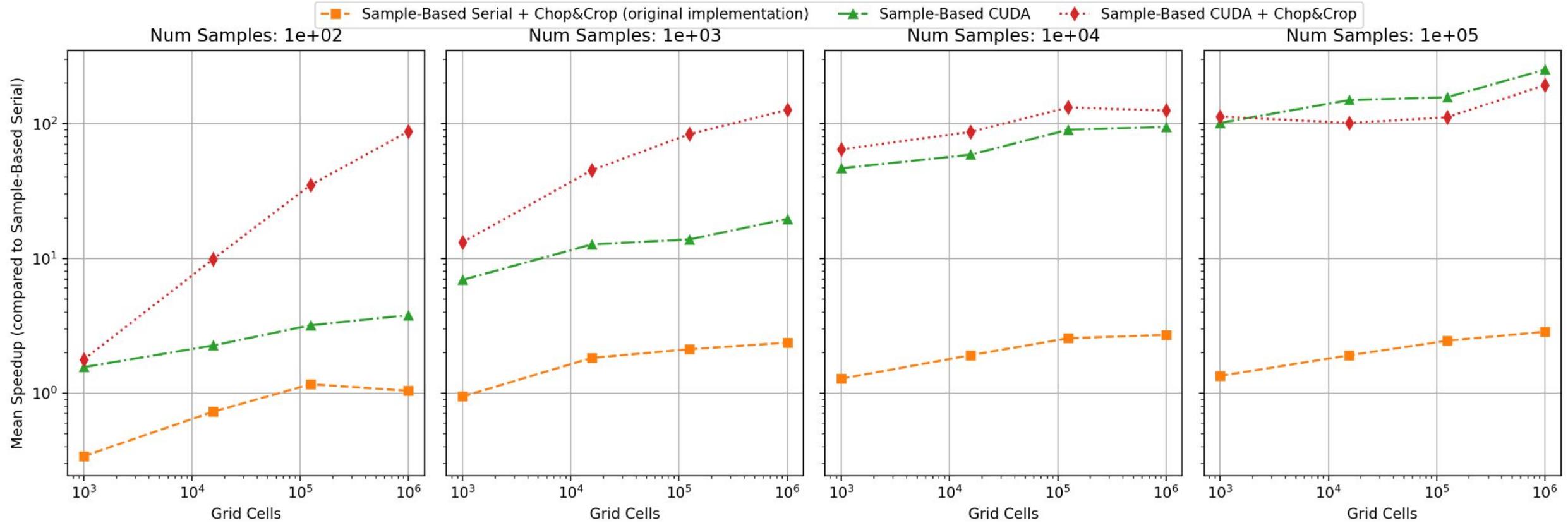


Extra Plots

Legend:
● Sample-Based Serial
■ Sample-Based Serial + Chop&Crop (original implementation)
▲ Sample-Based CUDA
◆ Sample-Based CUDA + Chop&Crop
Num Samples: 1e+02
Num Samples: 1e+03
Num Samples: 1e+04
Num Samples: 1e+05



Extra Plots



Extra Plots

