# SAGE: Self-Tuning Approximation for Graphics Engines

**Mehrzad Samadi**[1], Janghaeng Lee[1], D. Anoushe Jamshidi[1], Amir Hormati[2], and Scott Mahlke[1]
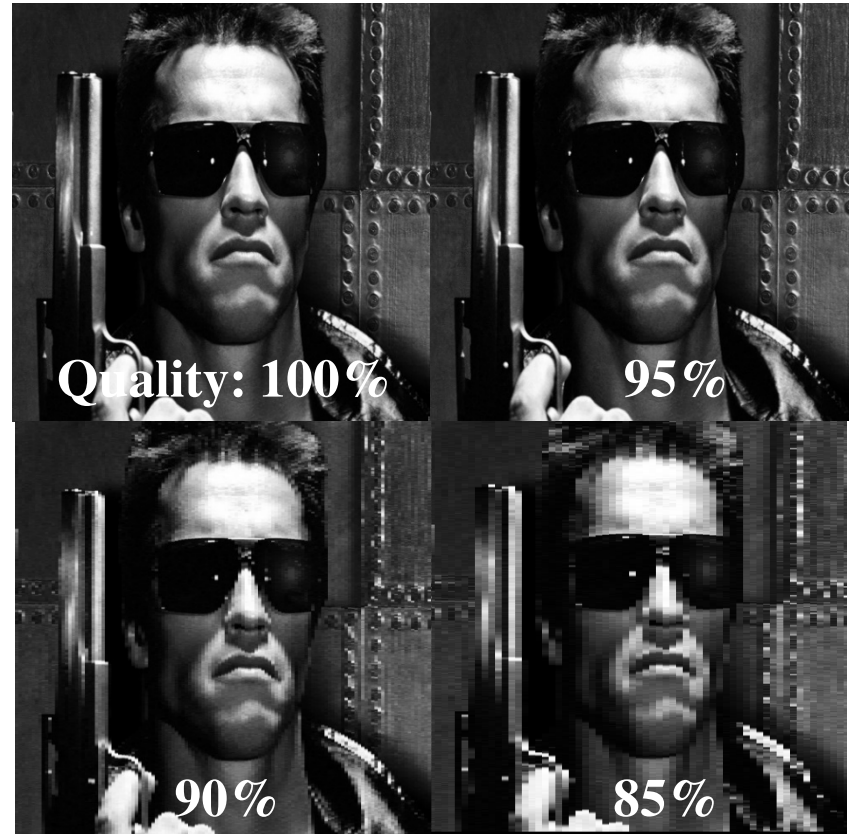
University of Michigan[1], Google Inc.[2]

December 2013

University of Michigan
Electrical Engineering and Computer Science

CCC

Compilers creating custom processors

# Approximate Computing

- Different domains:
  - Machine Learning
  - Image Processing
  - Video Processing
  - Physical Simulation
  - …



Quality: 100%    95%
90%    85%

Less work ➡ Higher performance
Lower power consumption

# Ubiquitous Graphics Processing Units

- Wide range of devices



Super Computers        Servers        Desktops        Cell Phones
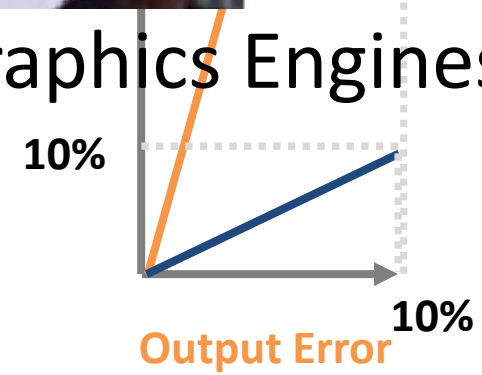
- Mostly regular applications
- Works on large data sets
  Good opportunity for automatic approximation
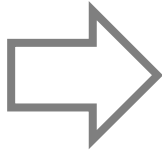
# SAGE Framework



- Sir

  - 
  - 

- Self-Tuning Approximation on Graphics Engines

  - Write the program once
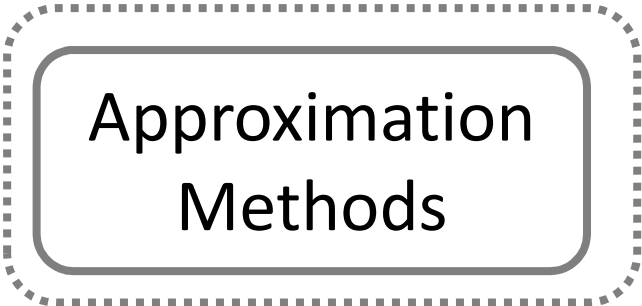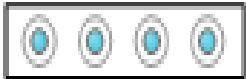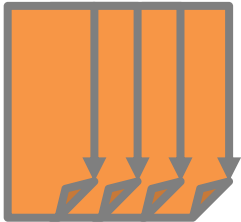  - Automatic approximation
  - Self-tuning dynamically

10%

10%

**Output Error**
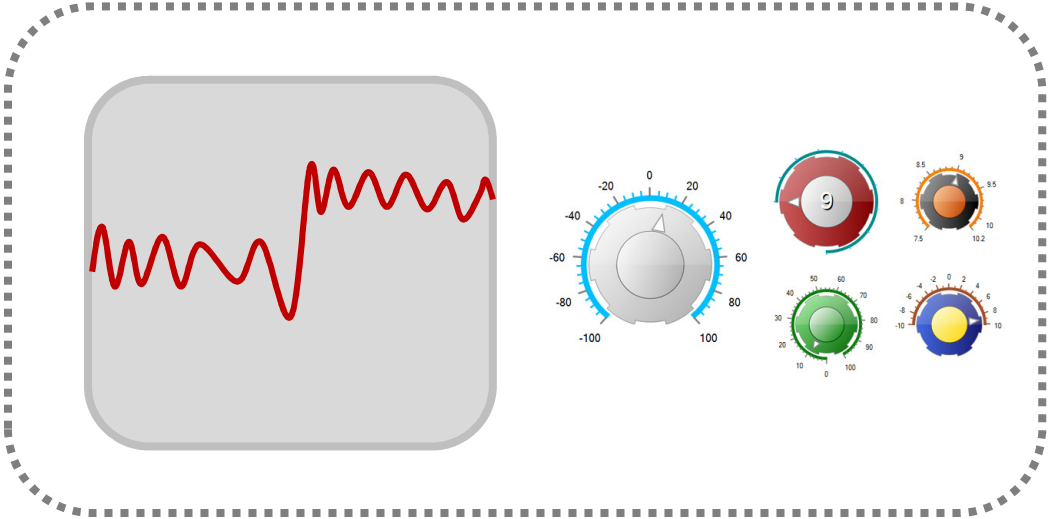
4

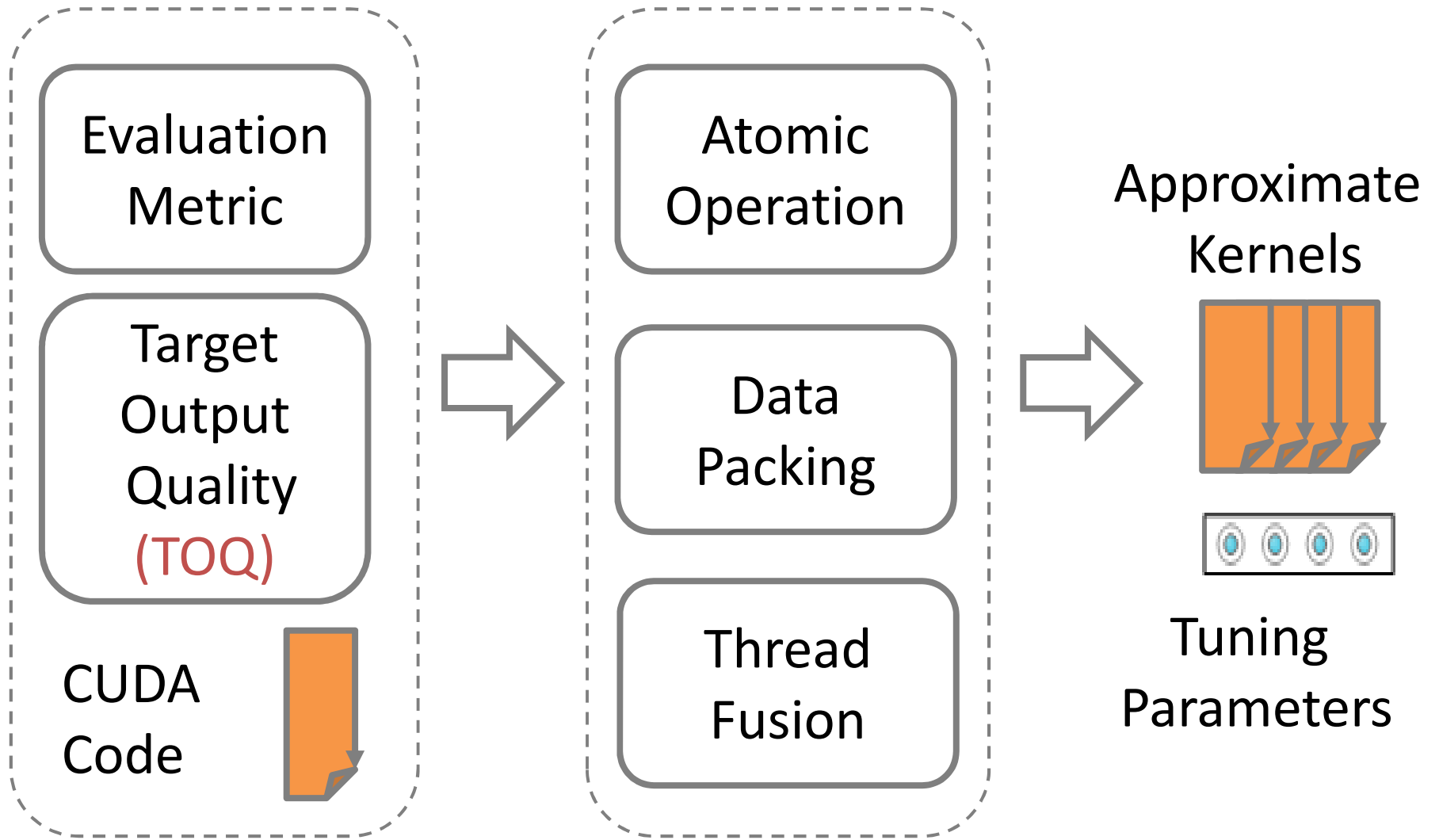# Overview

# SAGE Framework

Input Program

Static Compiler

Approximation Methods

Approximate Kernels

Runtime system

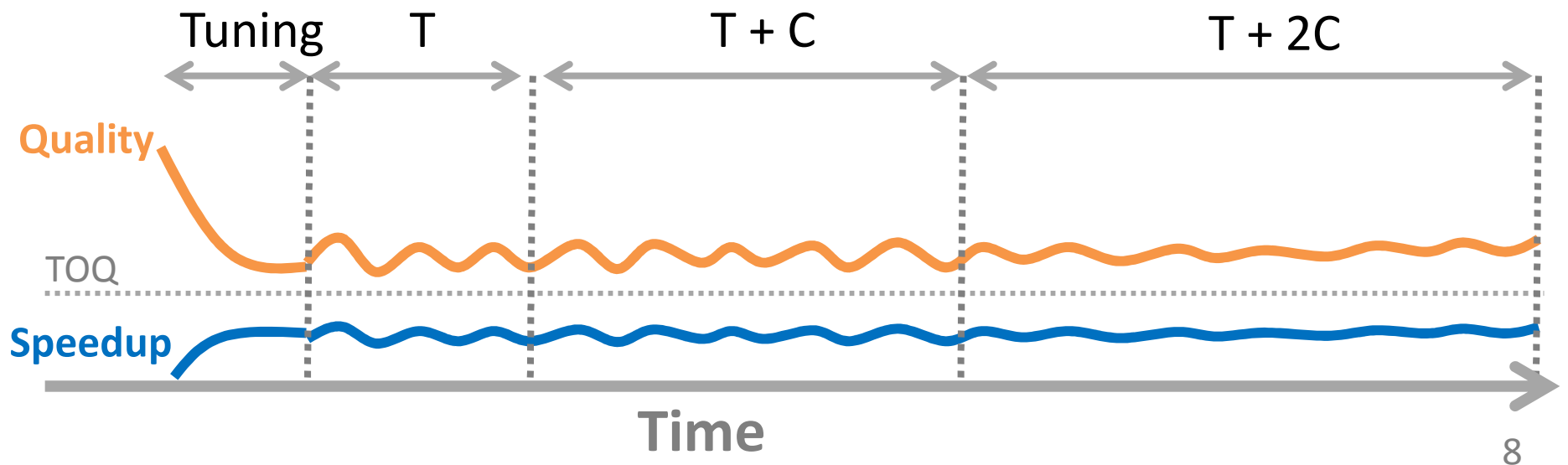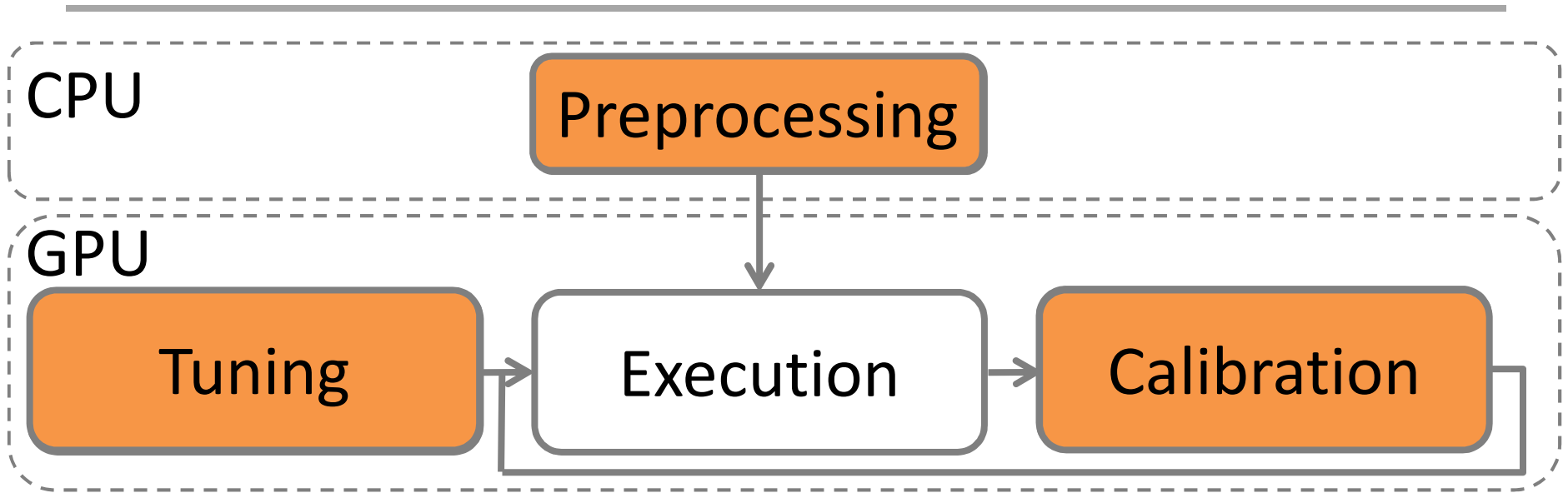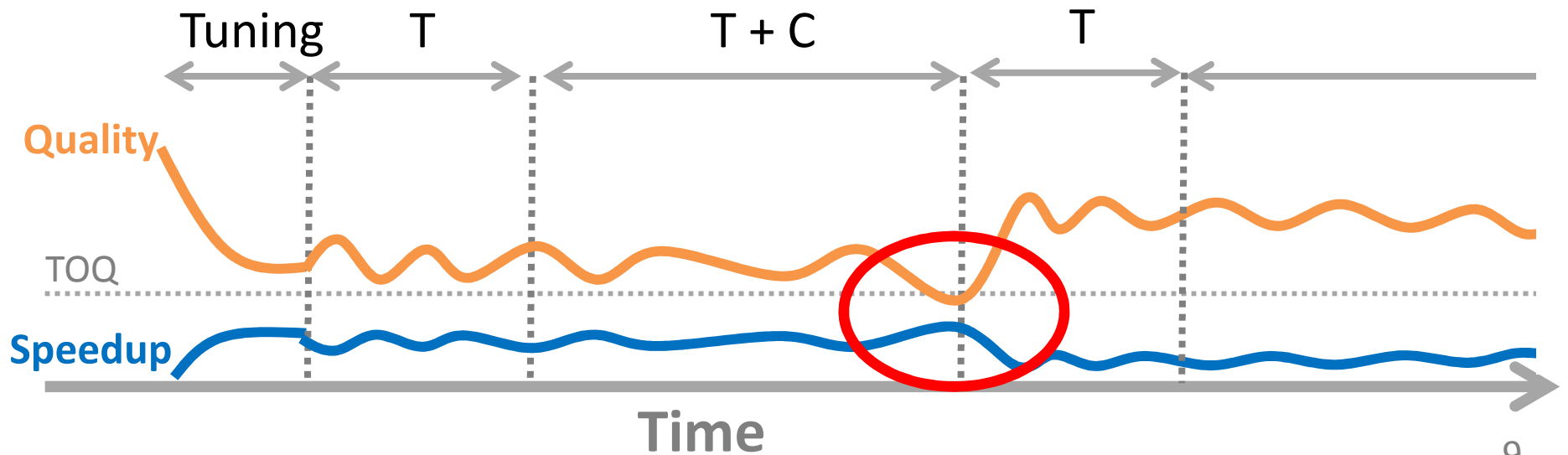Tuning Parameters

# Static Compilation

# Runtime System

CPU

Preprocessing

GPU

Tuning → Execution → Calibration

Tuning | T | T + C | T + 2C

Quality

TOQ

Speedup
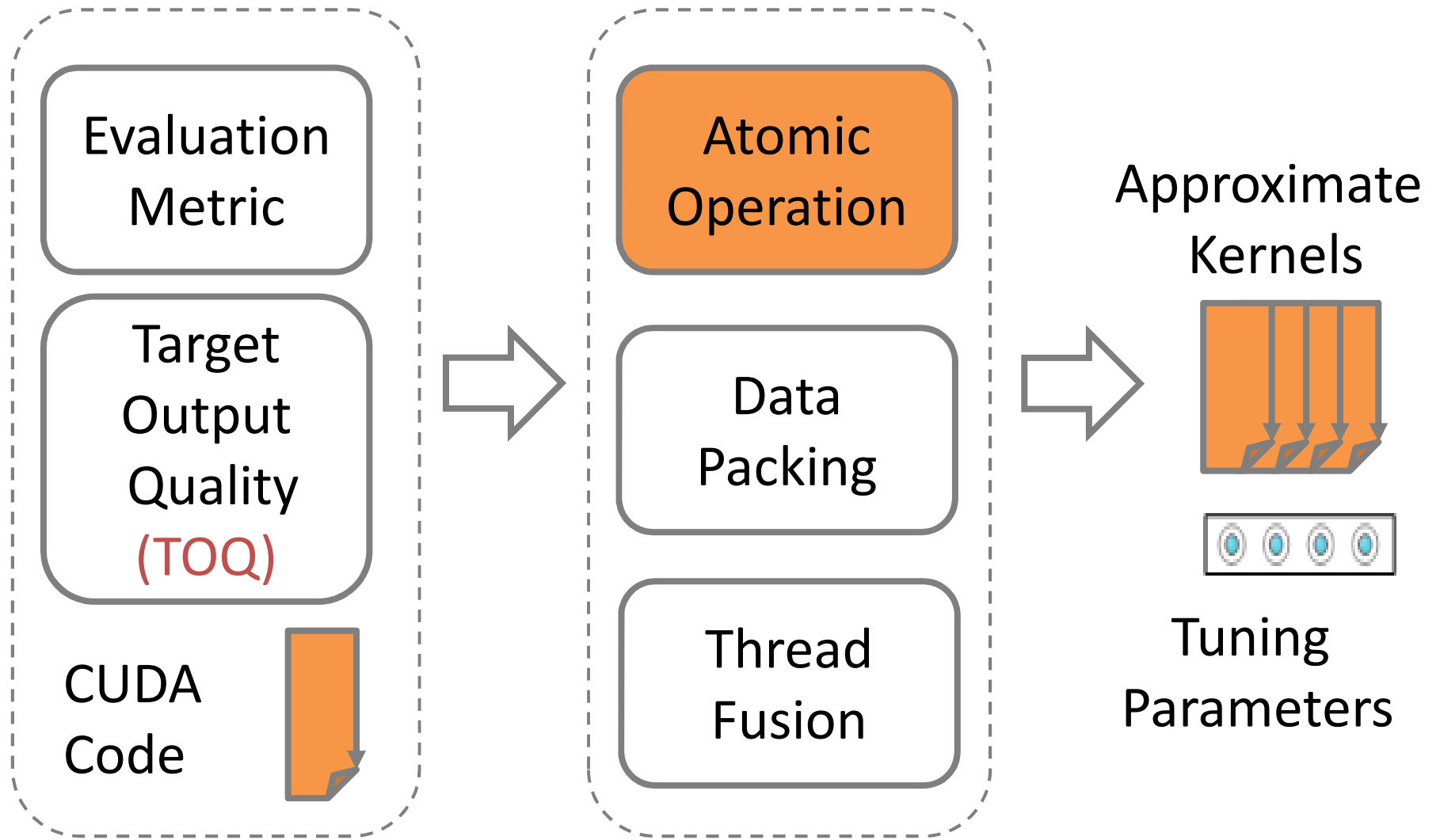
**Time**

# Runtime System

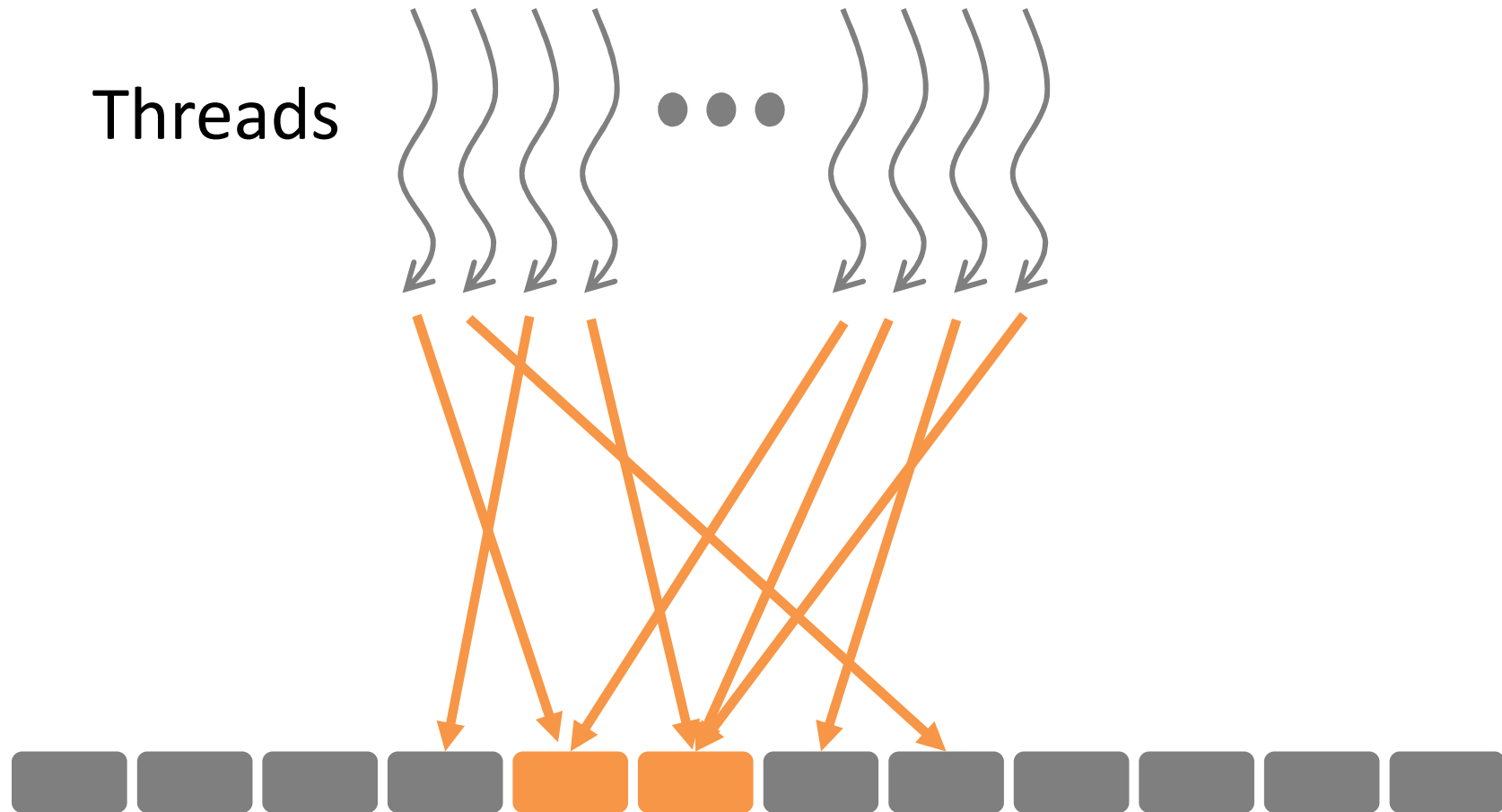# Approximation Methods

# Approximation Methods
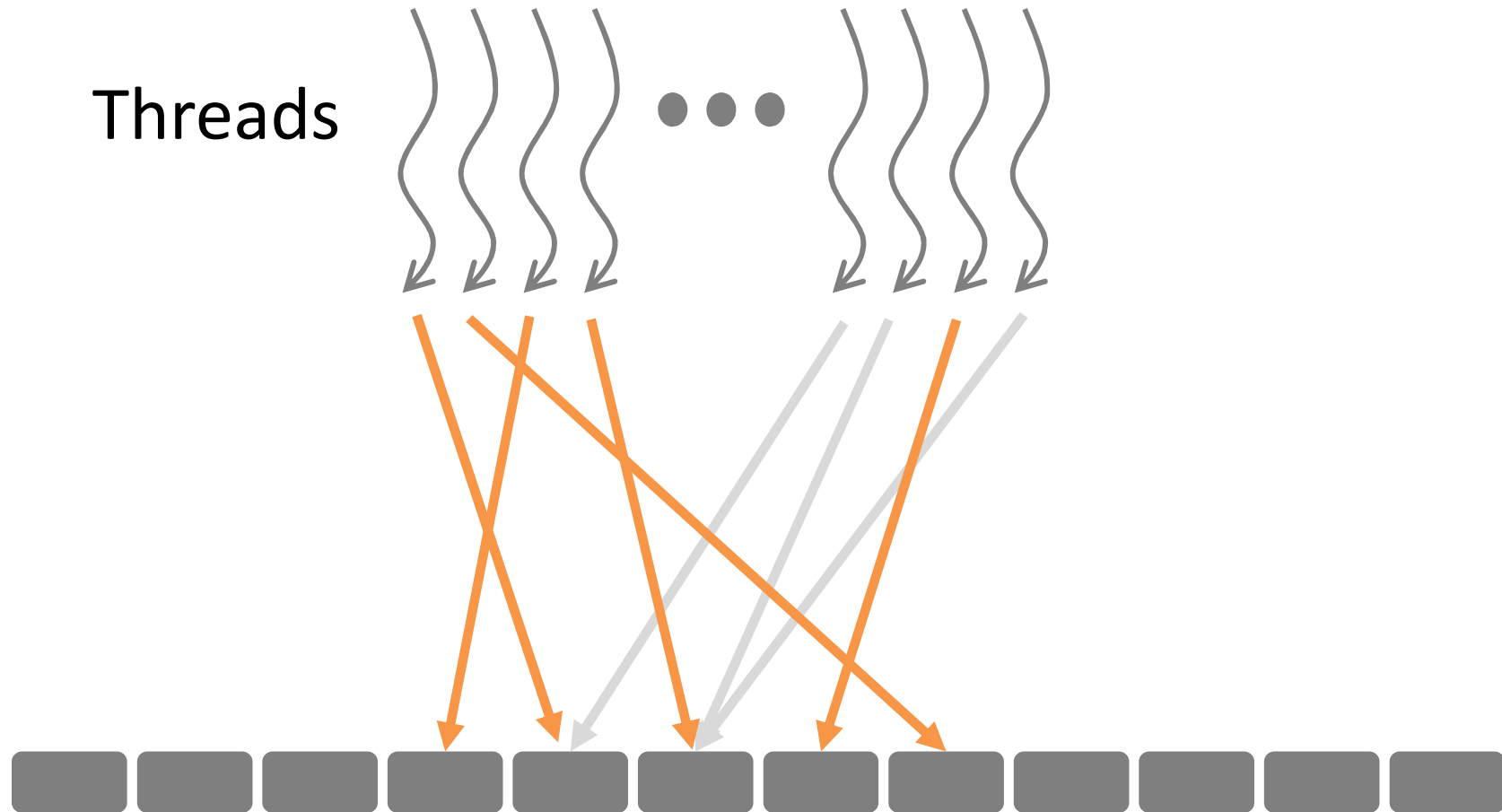
# Atomic Operations

- Atomic operations update a memory location such that the update appears to happen atomically

```
// Compute histogram of colors in an image
__global__ void histogram(int n, int* color, int* bucket)
int tid = threadIdx.x + blockDim.x * blockIdx.x;
int nThreads = gridDim.x * blockDim.x;
for ( int i = tid ; tid < n; tid += nThreads)
    int c = colors[i];
    atomicAdd(&bucket[c], 1);
```

# Atomic Operations



Threads

# Atomic Operations

Threads

# Atomic Operations

Threads

# Atomic Operations

Threads

# Atomic Add

Slowdown



Conflicts per Warp

# Atomic Operation Tuning

Iterations

```
// Compute histogram of colors in an image
__global__ void histogram(int n, int* color, int* bucket)
int tid = threadIdx.x + blockDim.x * blockIdx.x;
int nThreads = gridDim.x * blockDim.x;
for ( int i = tid ; tid < n; tid += nThreads)
    int c = colors[i];
    atomicAdd(&bucket[c], 1);
```
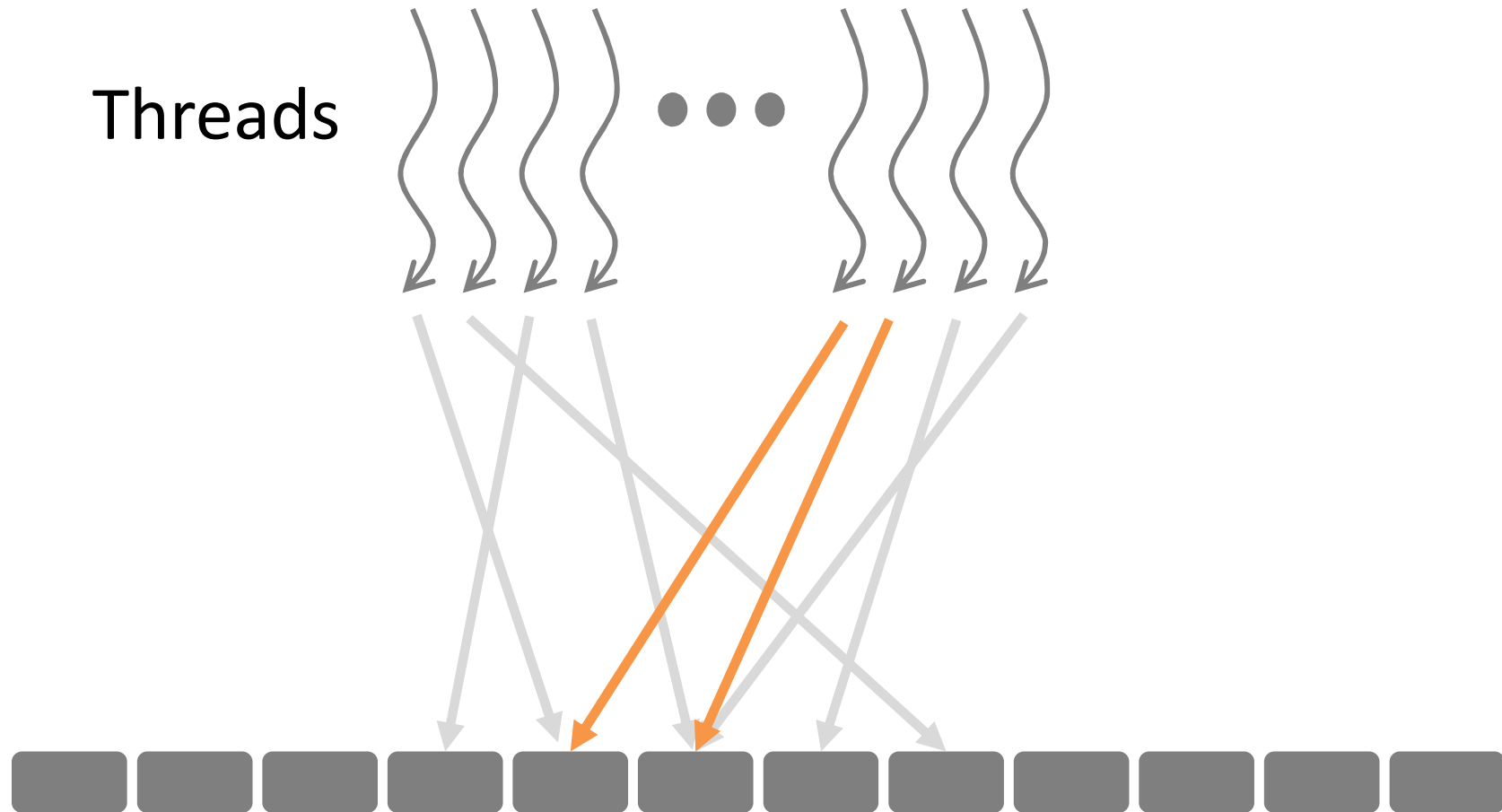
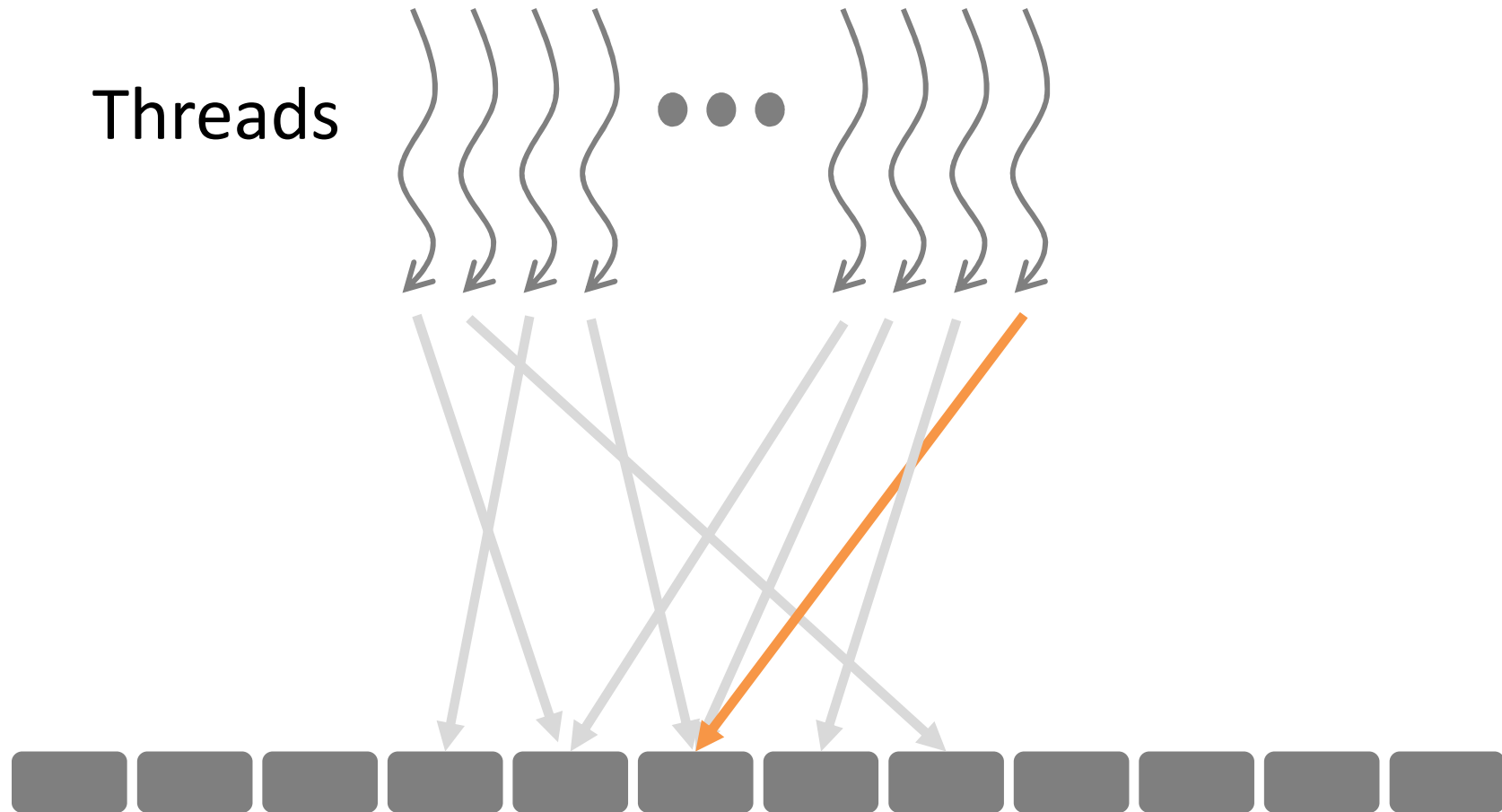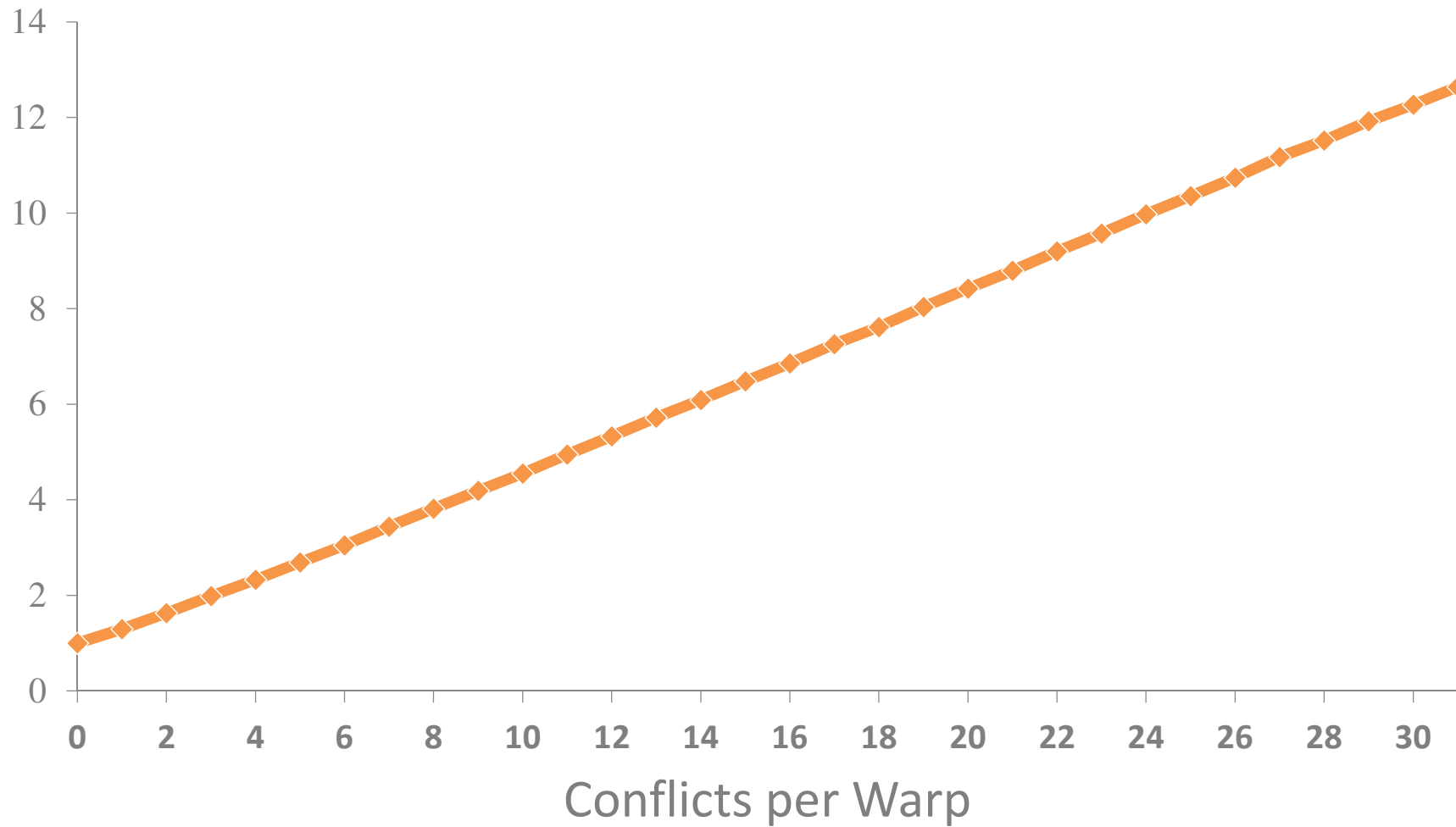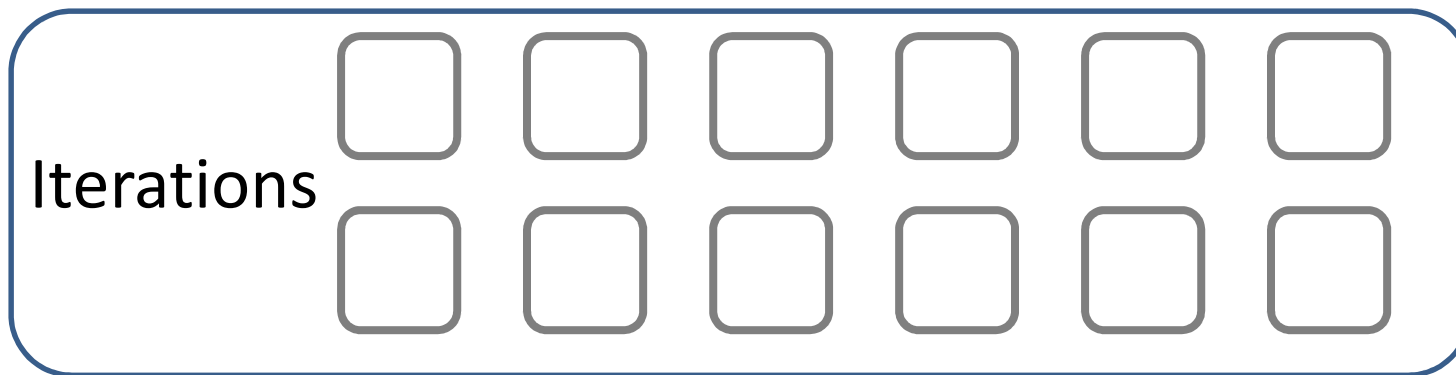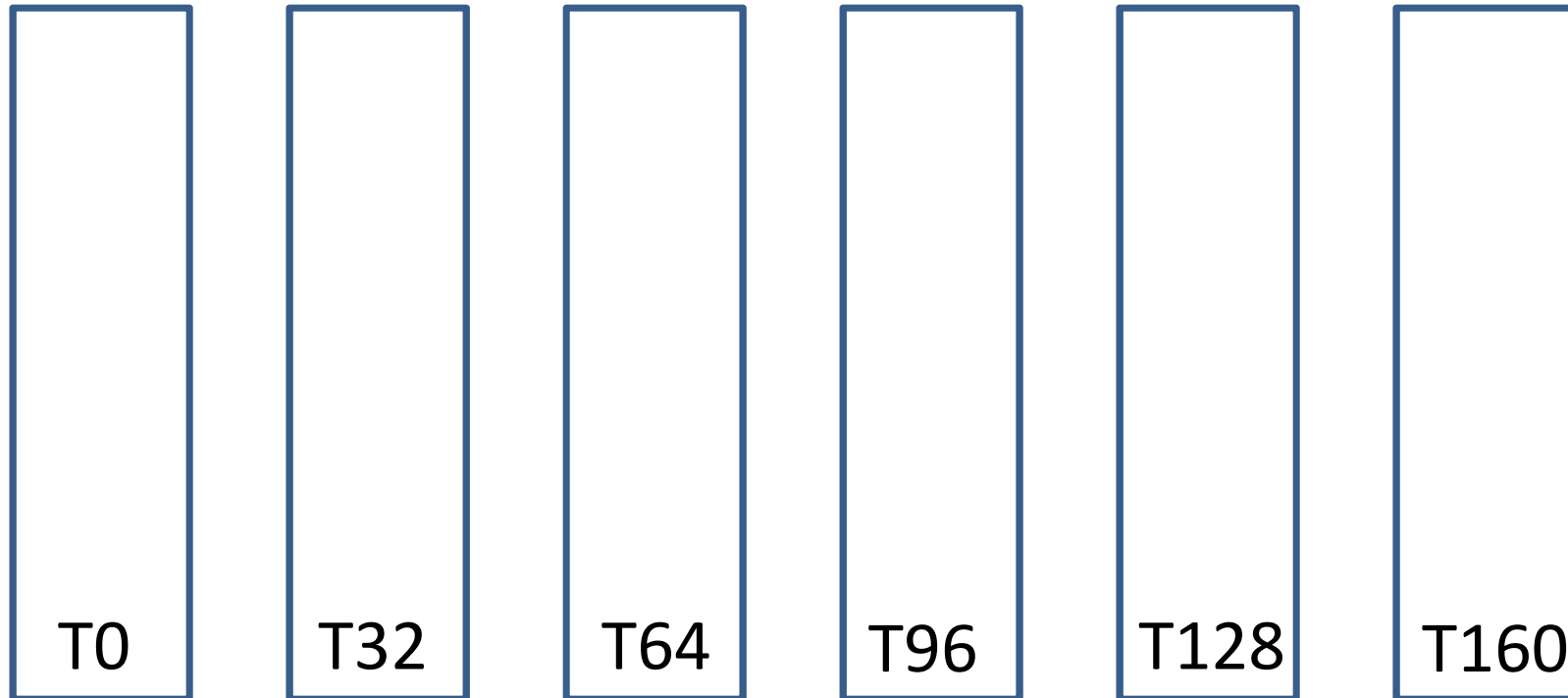T0    T32    T64    T96    T128    T160

# Atomic Operation Tuning

- SAGE skips one iteration per thread
- Iterations To improve the performance, it drops the iteration with the maximum number of conflicts

T0     T32     T64     T96     T128     T160

# Atomic Operation Tuning
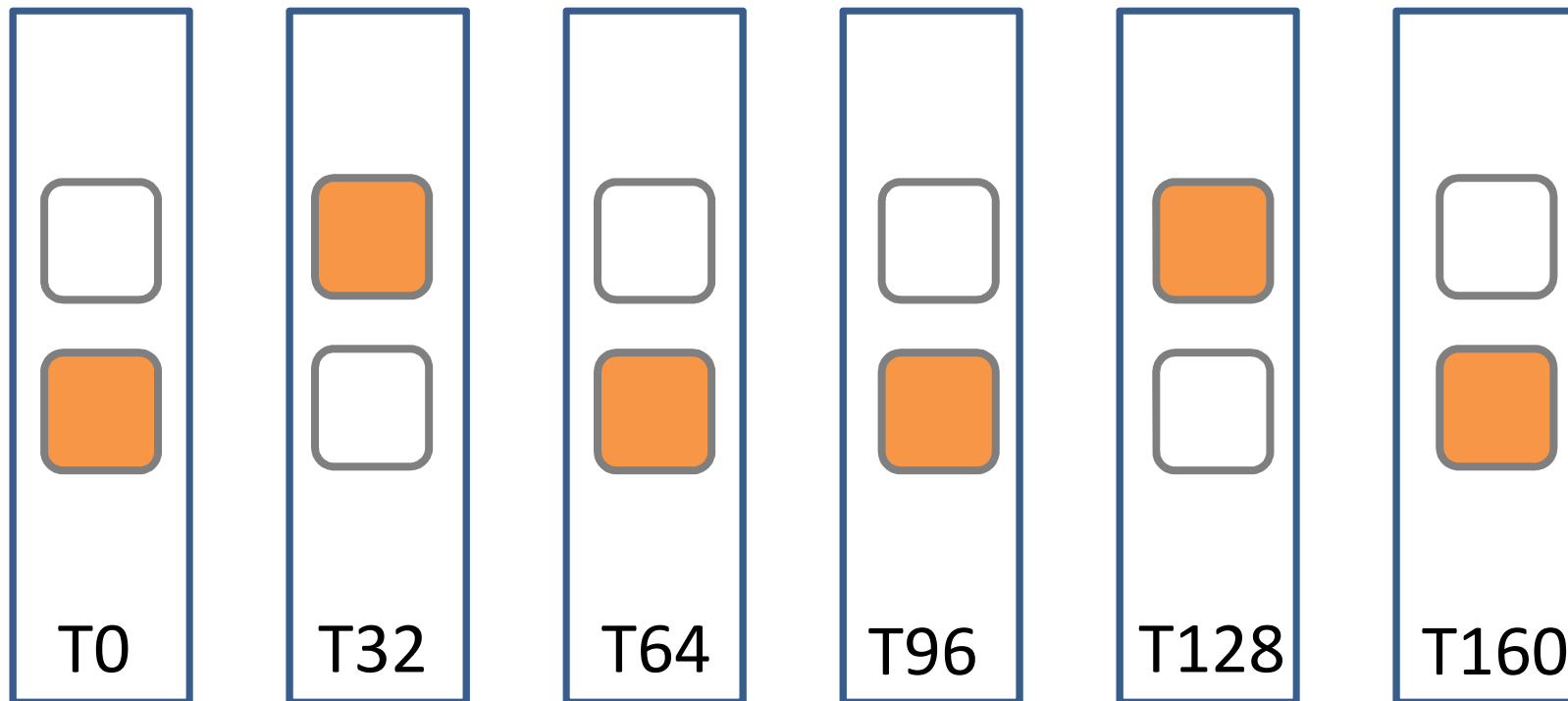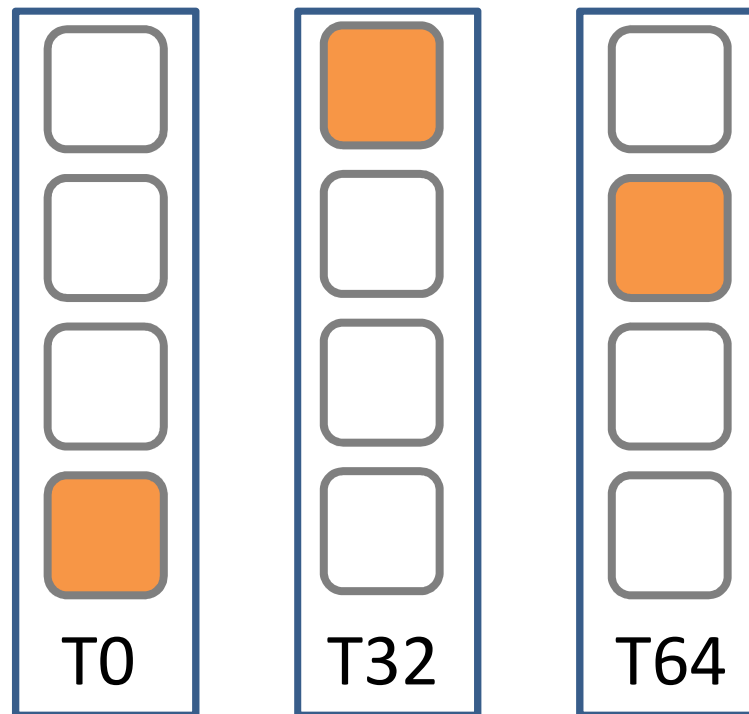
- SAGE skips one iteration per thread
- To improve the performance, it drops the
iteration with the maximum number of conflicts
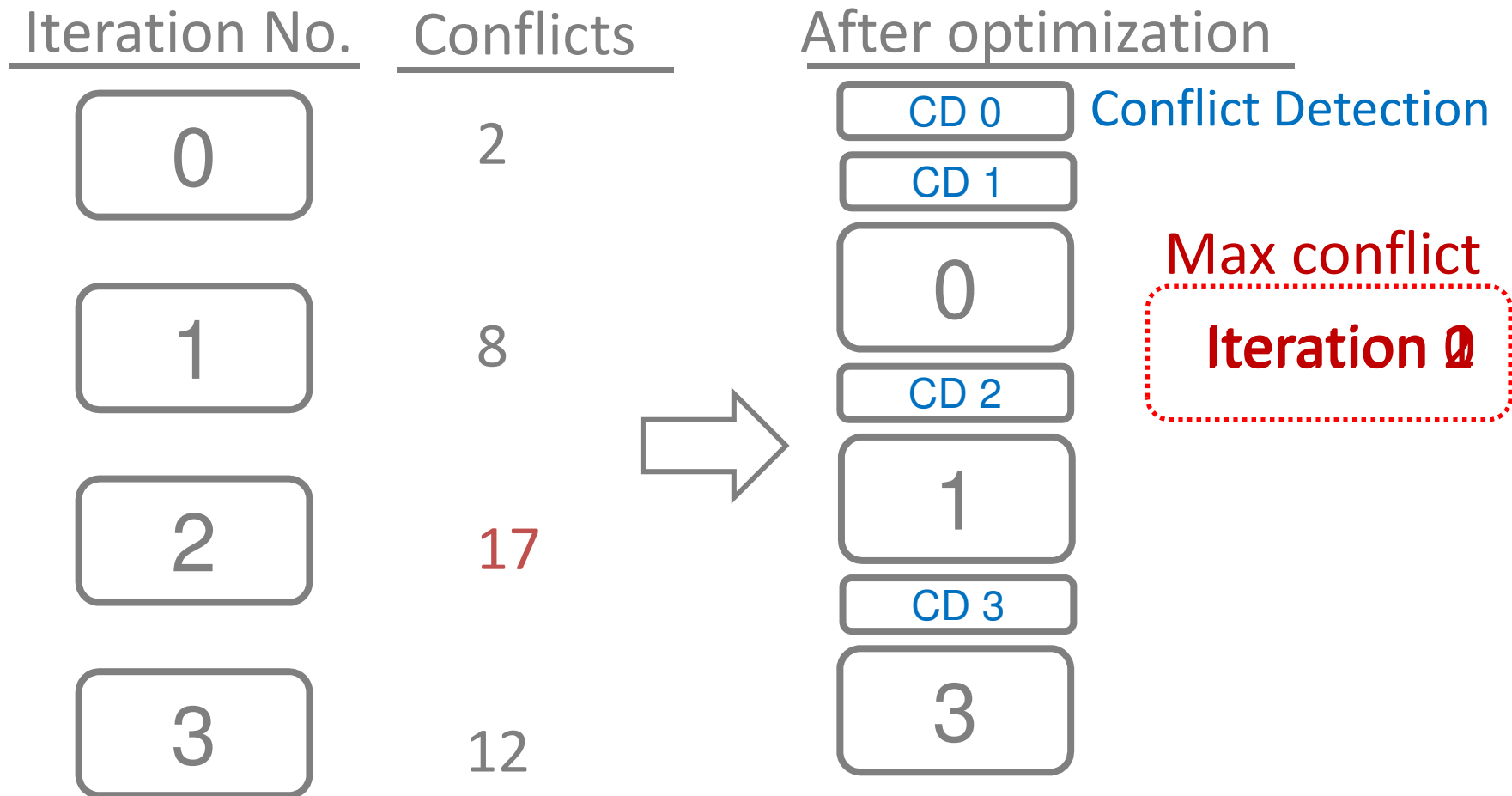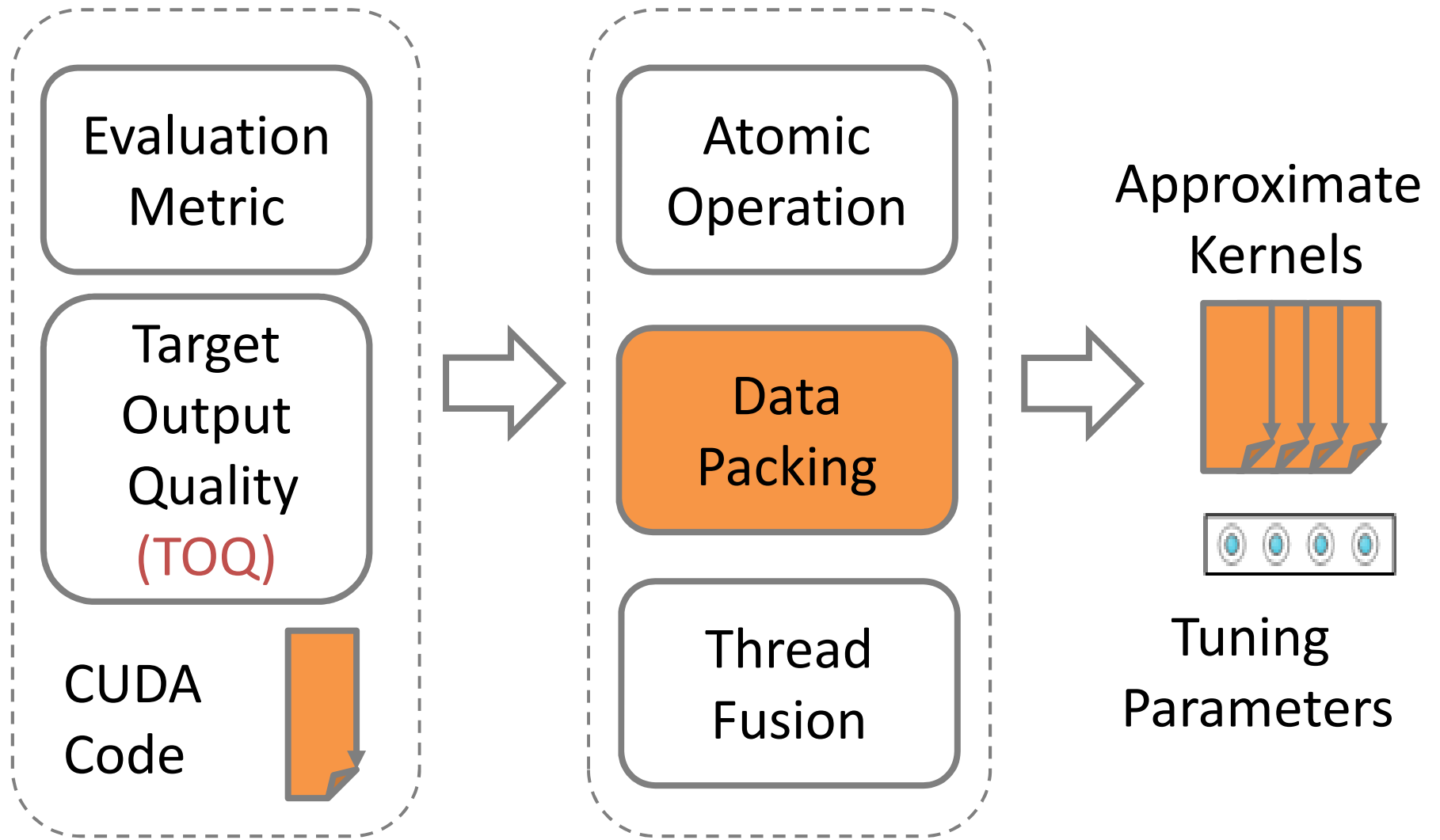
It drops 50% of iterations

T0    T32    T64    T96    T128    T160

# Atomic Operation Tuning

Drop rate goes down to 25%



T0    T32    T64

# Dropping One Iteration

Iteration No.

Conflicts

After optimization

| 0 | 2 |
| 1 | 8 |
| 2 | 17 |
| 3 | 12 |

CD 0 — Conflict Detection

CD 1

0

CD 2

1

CD 3

3

Max conflict

Iteration 0

# Approximation Methods
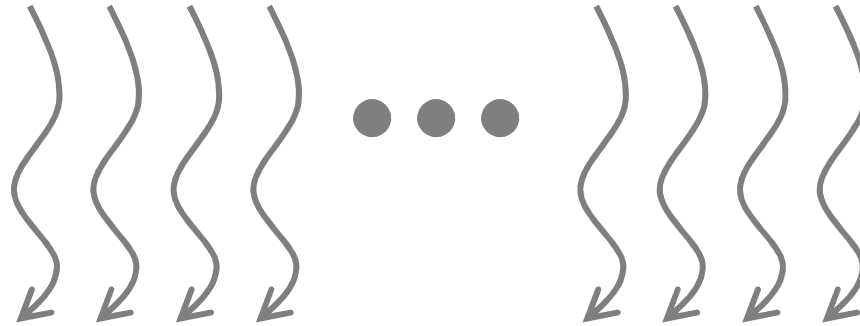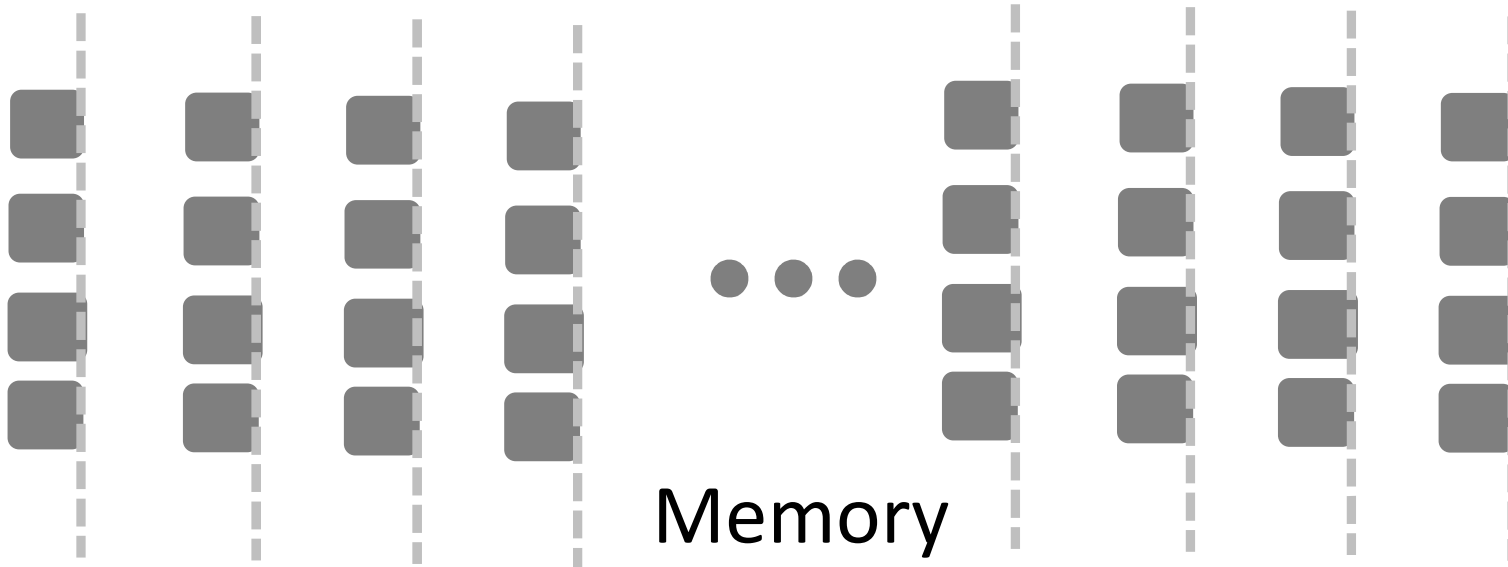
# Data Packing

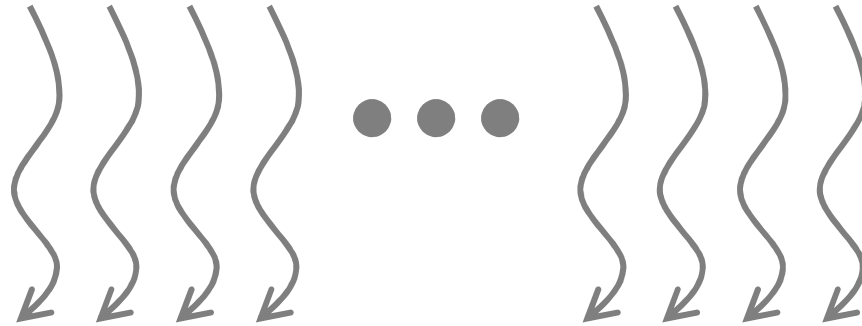Threads

Memory

# Data Packing

Threads

Memory

# Data Packing

Threads

Memory

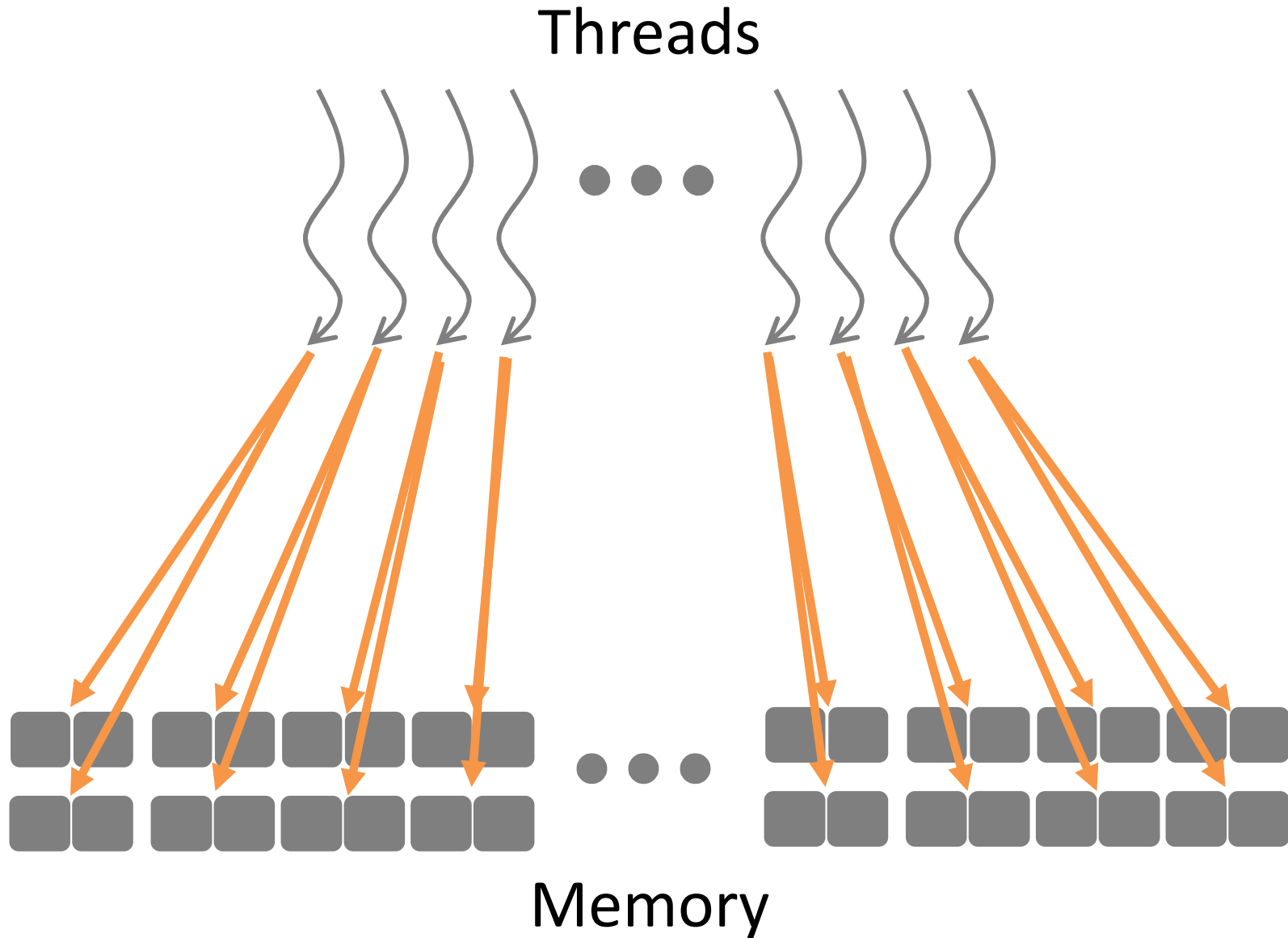# Data Packing

Threads



Memory

# Quantization

- Preprocessing finds min and max of the input sets and packs the data



- During execution, each thread unpacks the data and transforms the quantization level to data by applying a linear transformation
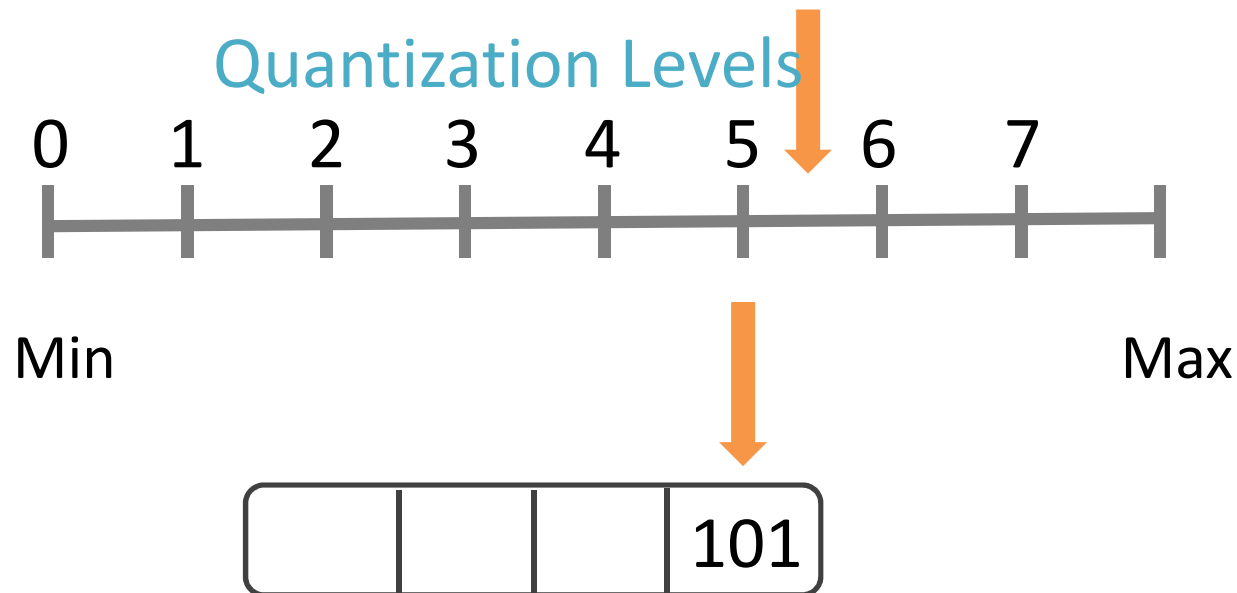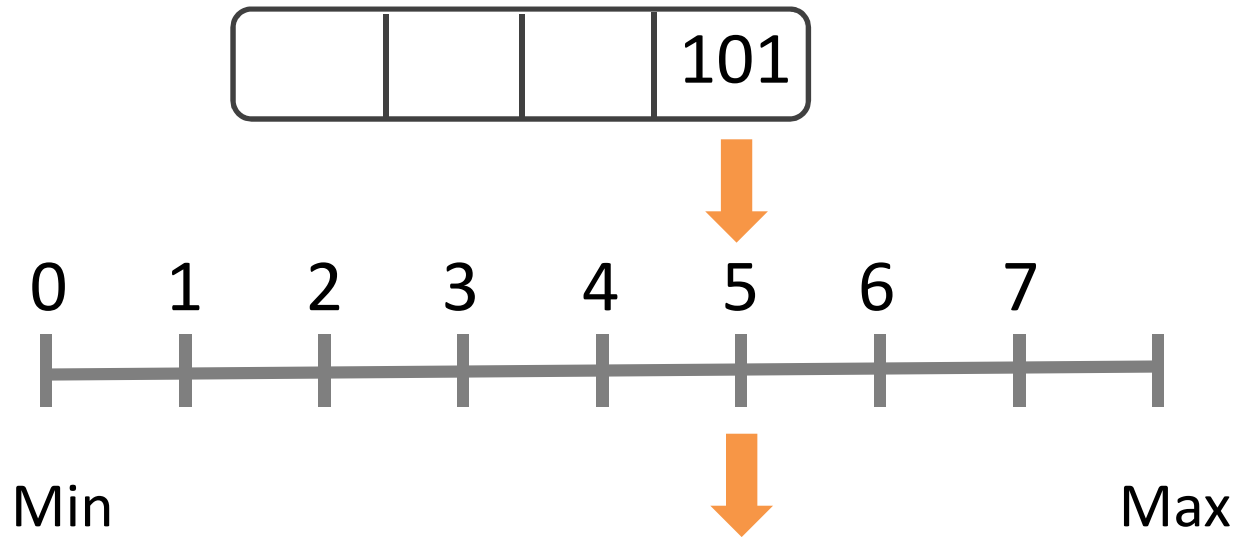
# Quantization

- Preprocessing finds max and min of the input sets and packs the data



- During execution, each thread unpacks the data and transforms the quantization level to data by applying a linear transformation
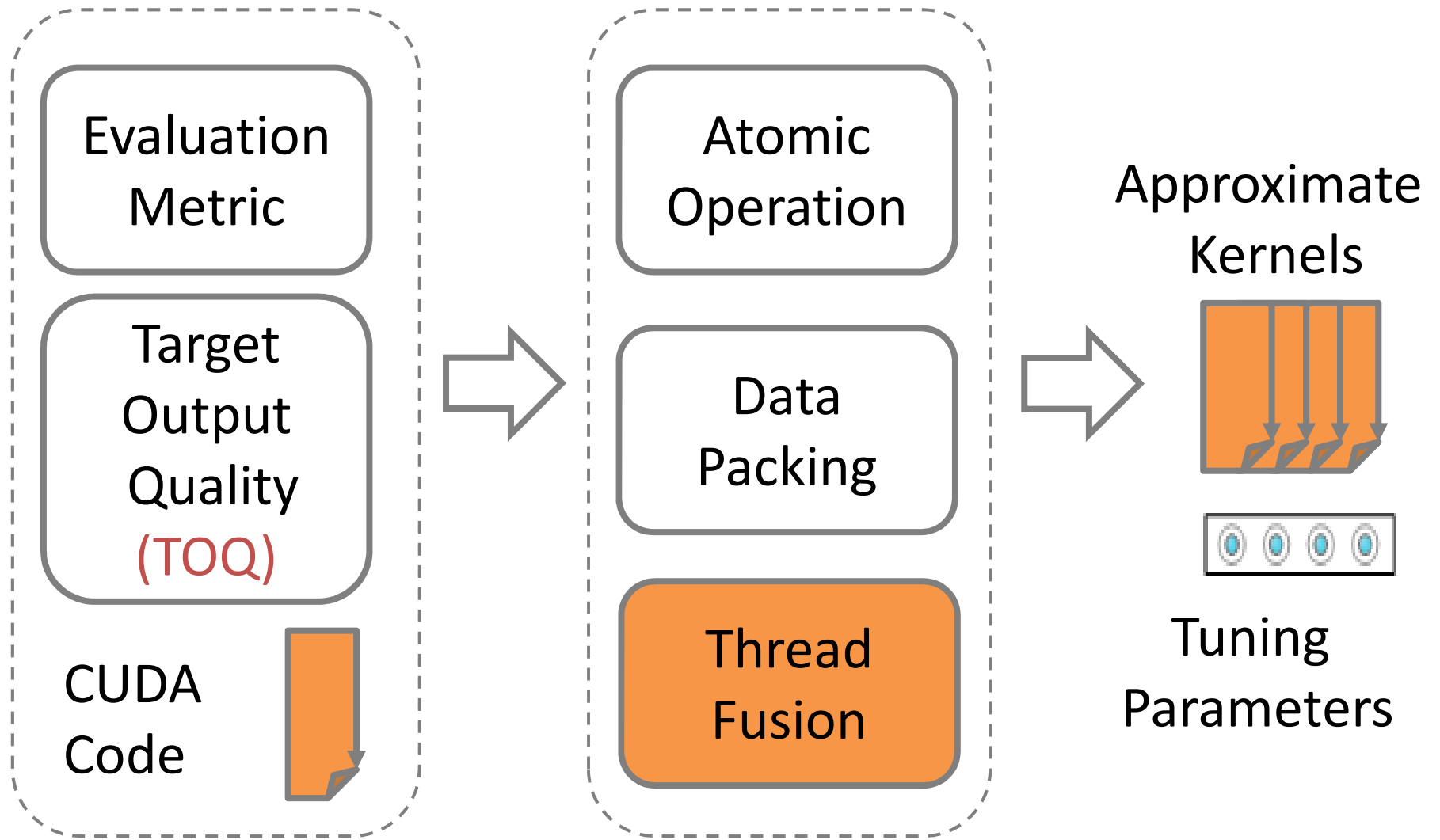
# Approximation Methods

Evaluation Metric

Target Output Quality (TOQ)

CUDA Code

Atomic Operation

Data Packing

Thread Fusion

Approximate Kernels

Tuning Parameters

# Thread Fusion

Memory

Threads

Memory

# Thread Fusion



Memory

Threads

Memory

# Thread Fusion

Computation
Output writing

Computation

Output writing

T0 T1

# Thread Fusion

# Thread Fusion



Reducing the number of threads per block results in poor resource utilization

T0      T127      T0      T127

Block 0      Block 1

# Block Fusion



T0  T1     T254 T255

Block 0 & 1 fused

# Runtime

# How to Compute Output Quality?

Approximate Version    Accurate Version    Evaluation Metric

- High overhead
- Tuning should find a good enough approximate version as fast as possible
- Greedy algorithm

# Tuning

**TOQ = 90%**  $\boxed{\text{K(0,0)}}$ Quality = 100%
Speedup = 1x

Tuning parameter of the
First optimization

Tuning parameter of the
Second optimization

**K(x,y)**

# Tuning

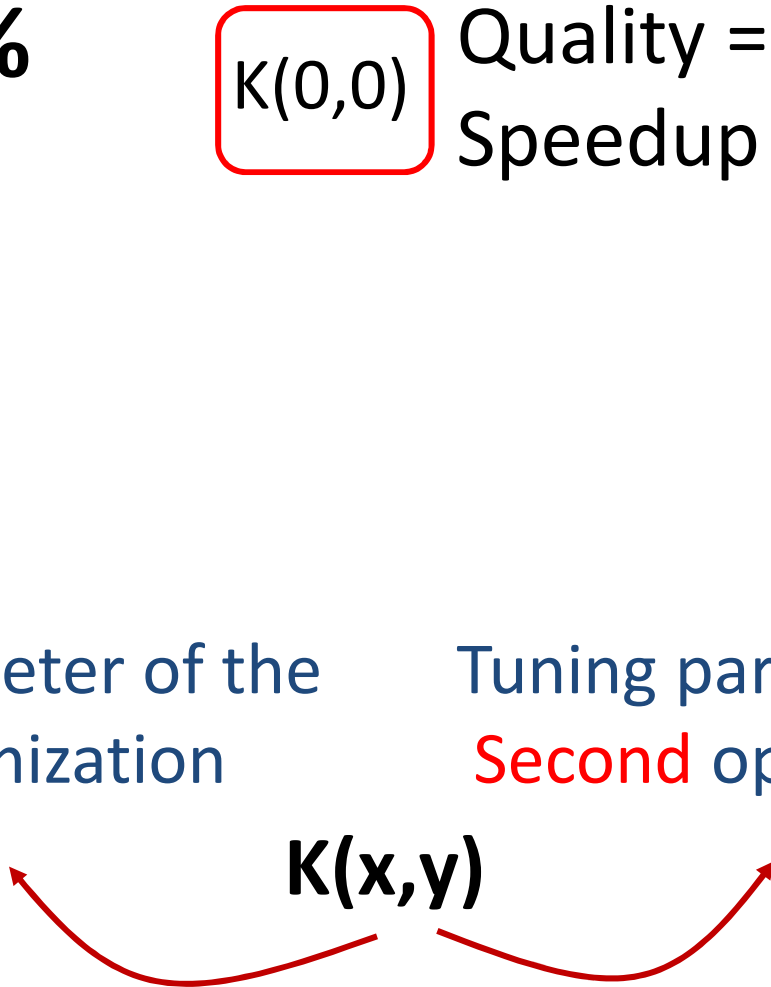**TOQ = 90%**

K(0,0)

Quality = 100%
Speedup = 1x

94%
1.15X K(1,0)

K(0,1) 96%
1.5X

# Tuning

**TOQ = 90%**

K(0,0)

Quality = 100%
Speedup = 1x

94%
1.15X  K(1,0)

K(0,1)  96%
1.5X

94%
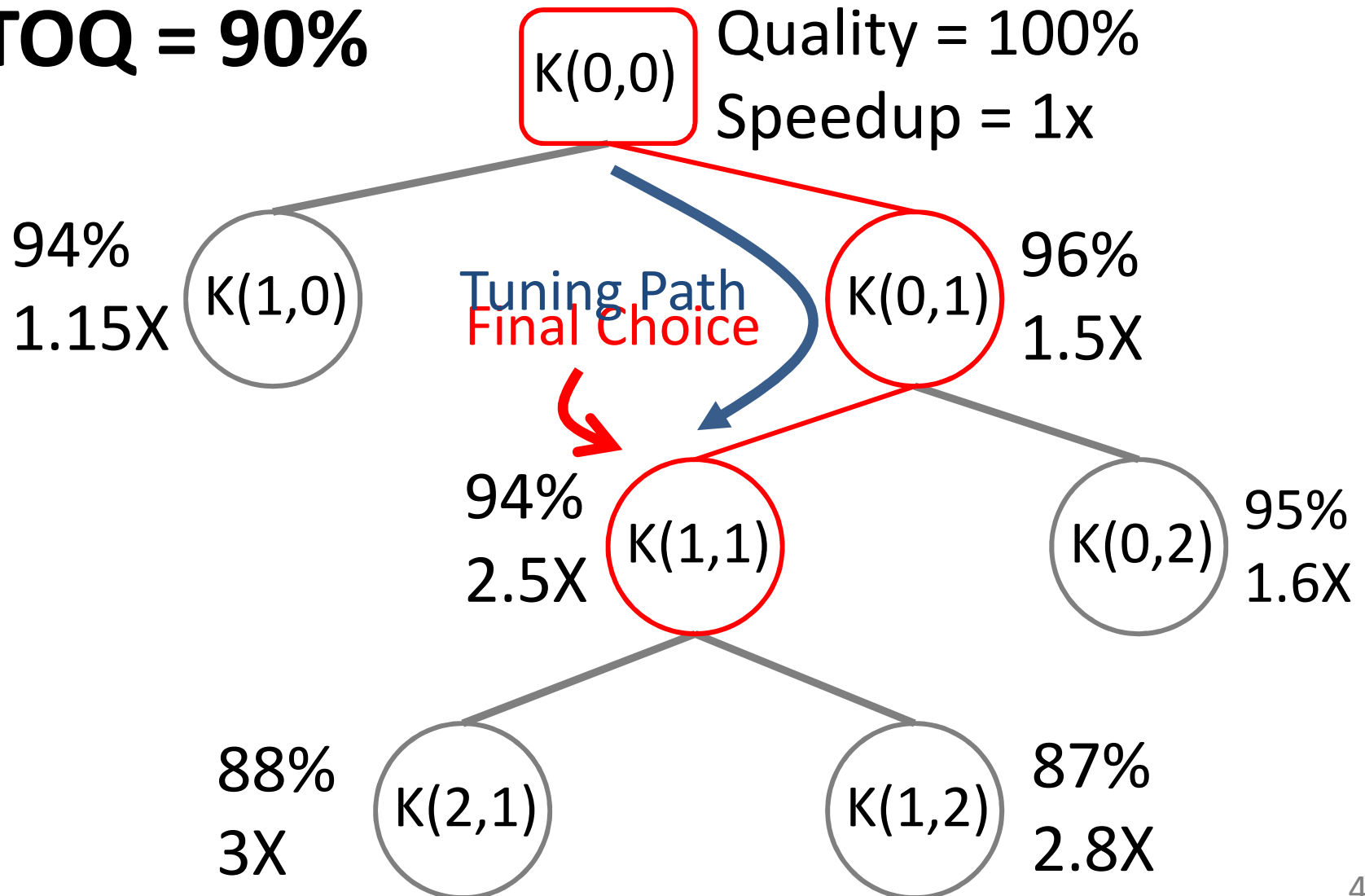2.5X  K(1,1)

K(0,2)  95%
1.6X

# Tuning

**TOQ = 90%**



K(0,0)  Quality = 100%  Speedup = 1x

94%  1.15X  K(1,0)

Tuning Path

Final Choice

K(0,1)  96%  1.5X

94%  2.5X  K(1,1)

K(0,2)  95%  1.6X

88%  3X  K(2,1)
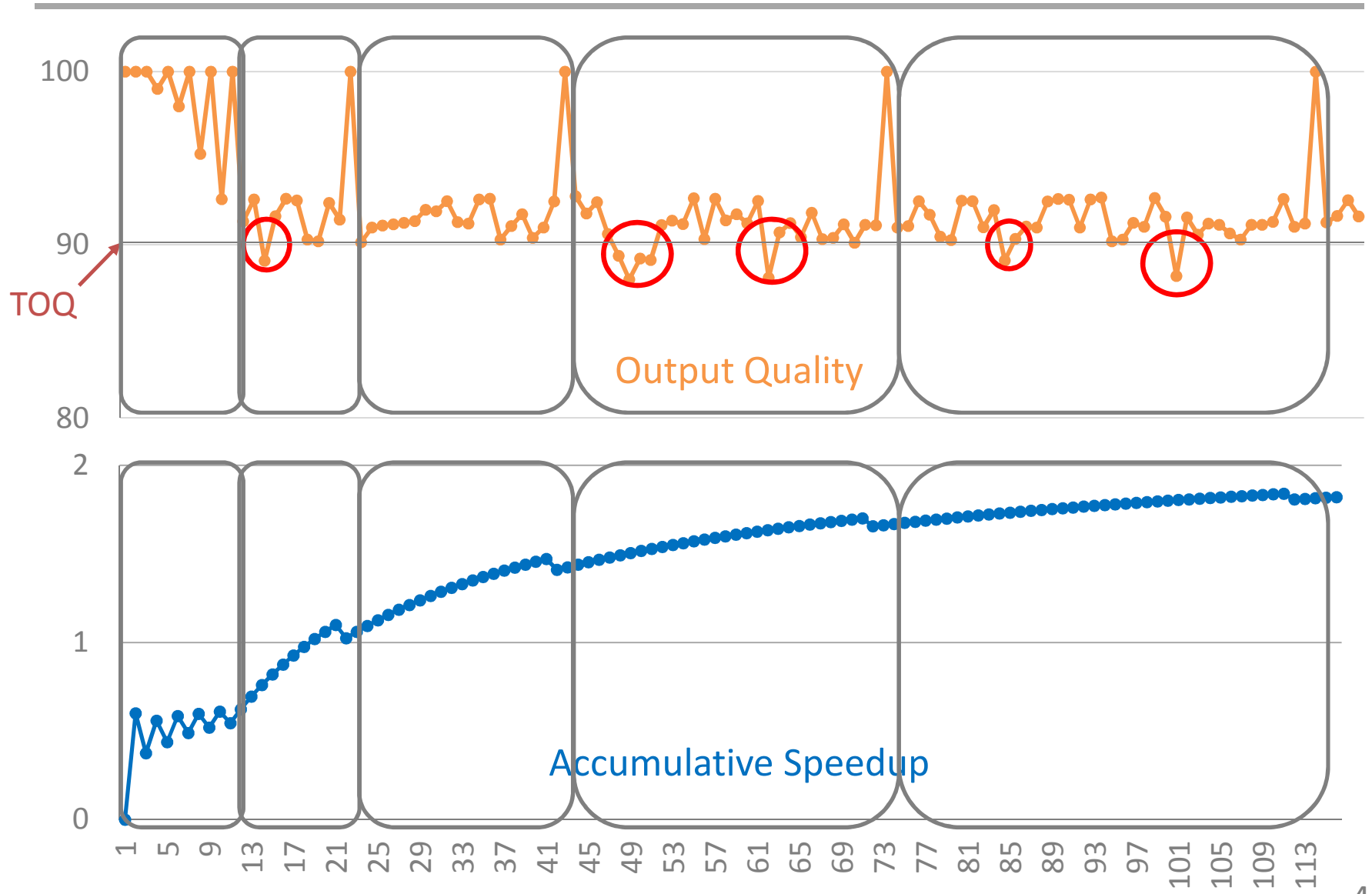
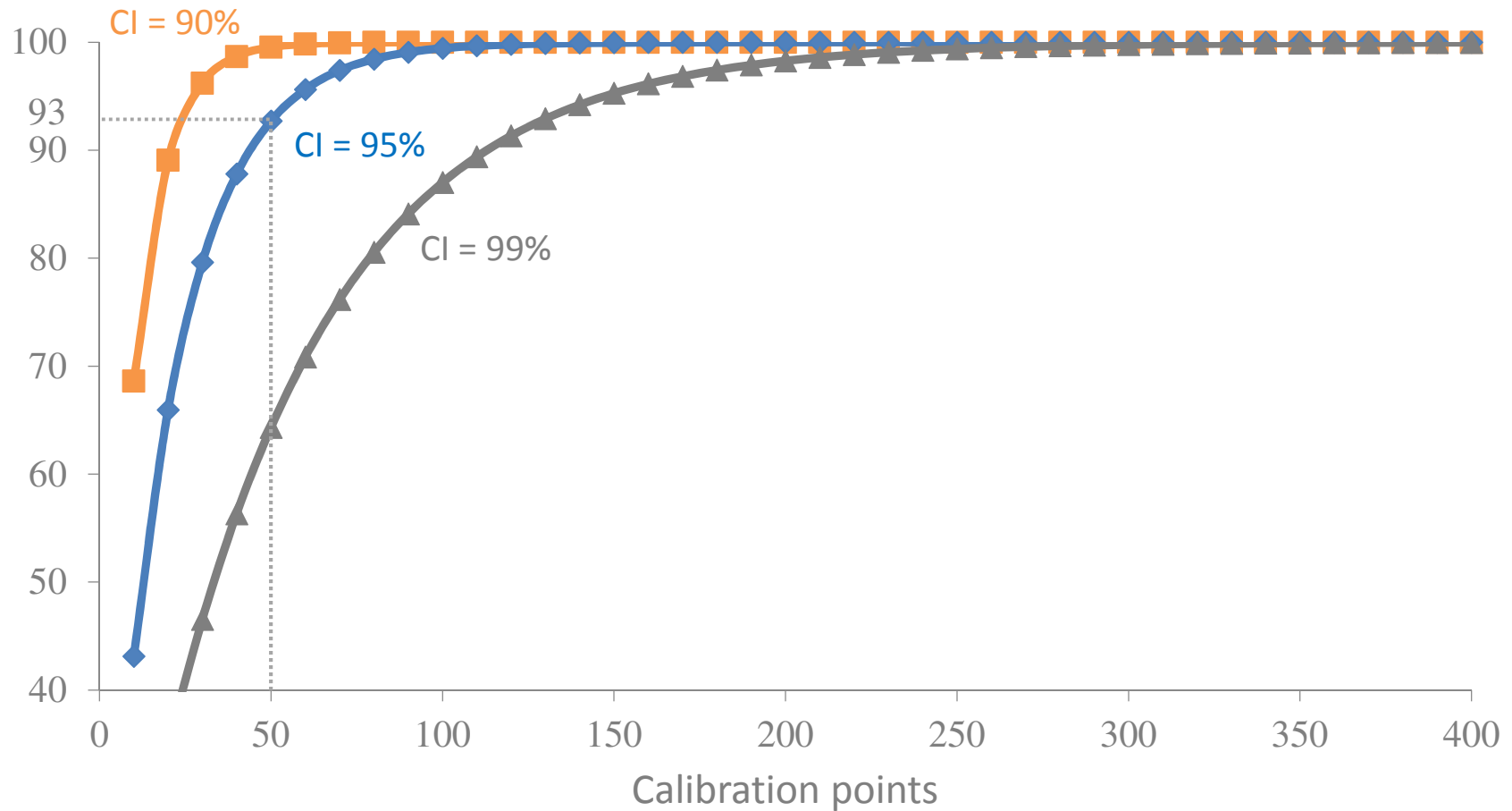K(1,2)  87%  2.8X

# Evaluation

# Experimental Setup

- Backend of Cetus compiler

- GPU
  - NVIDIA GTX 560
    - 2GB GDDR 5

- CPU
  - Intel Core i7

- Benchmarks
  - Image processing
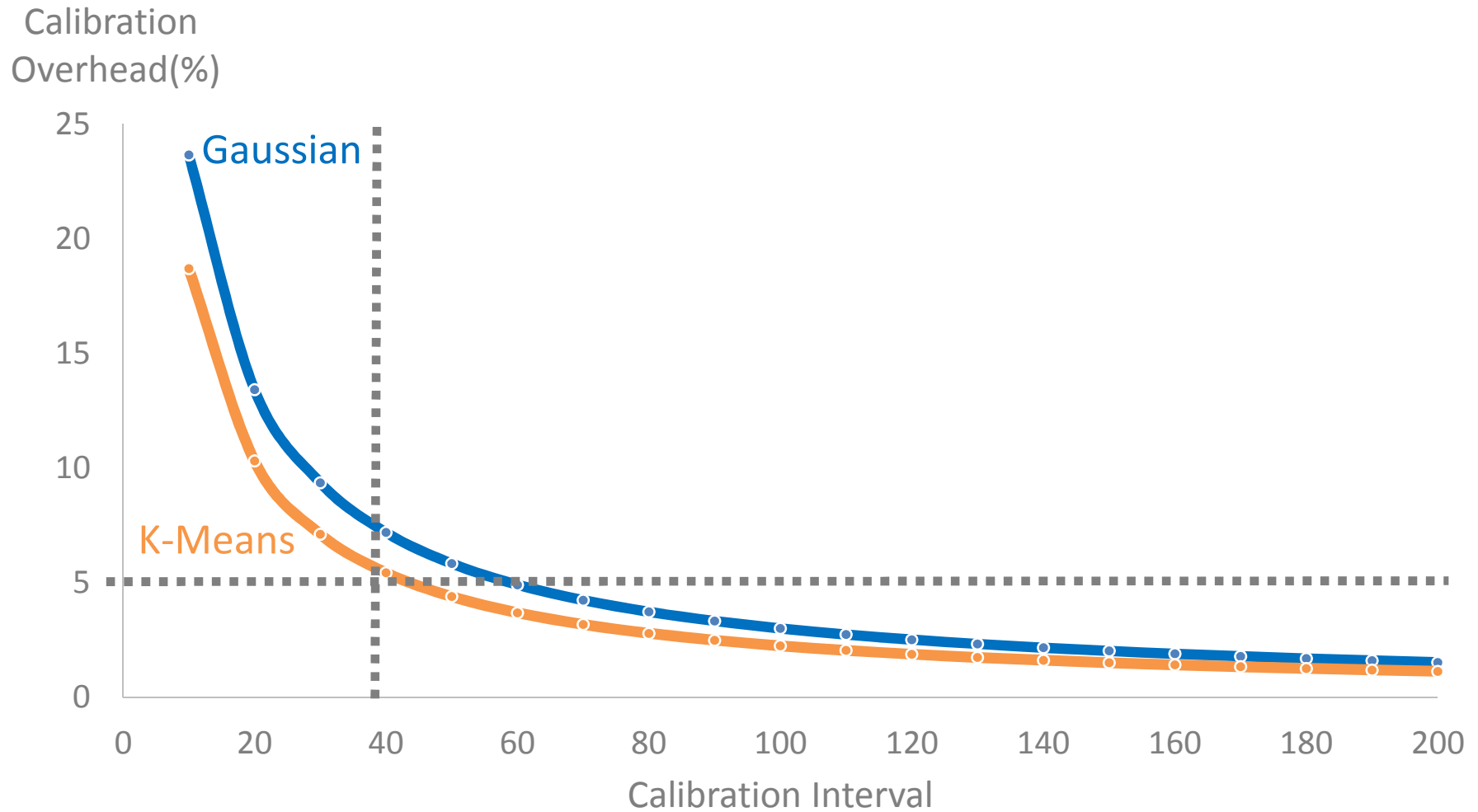  - Machine Learning

# K-Means

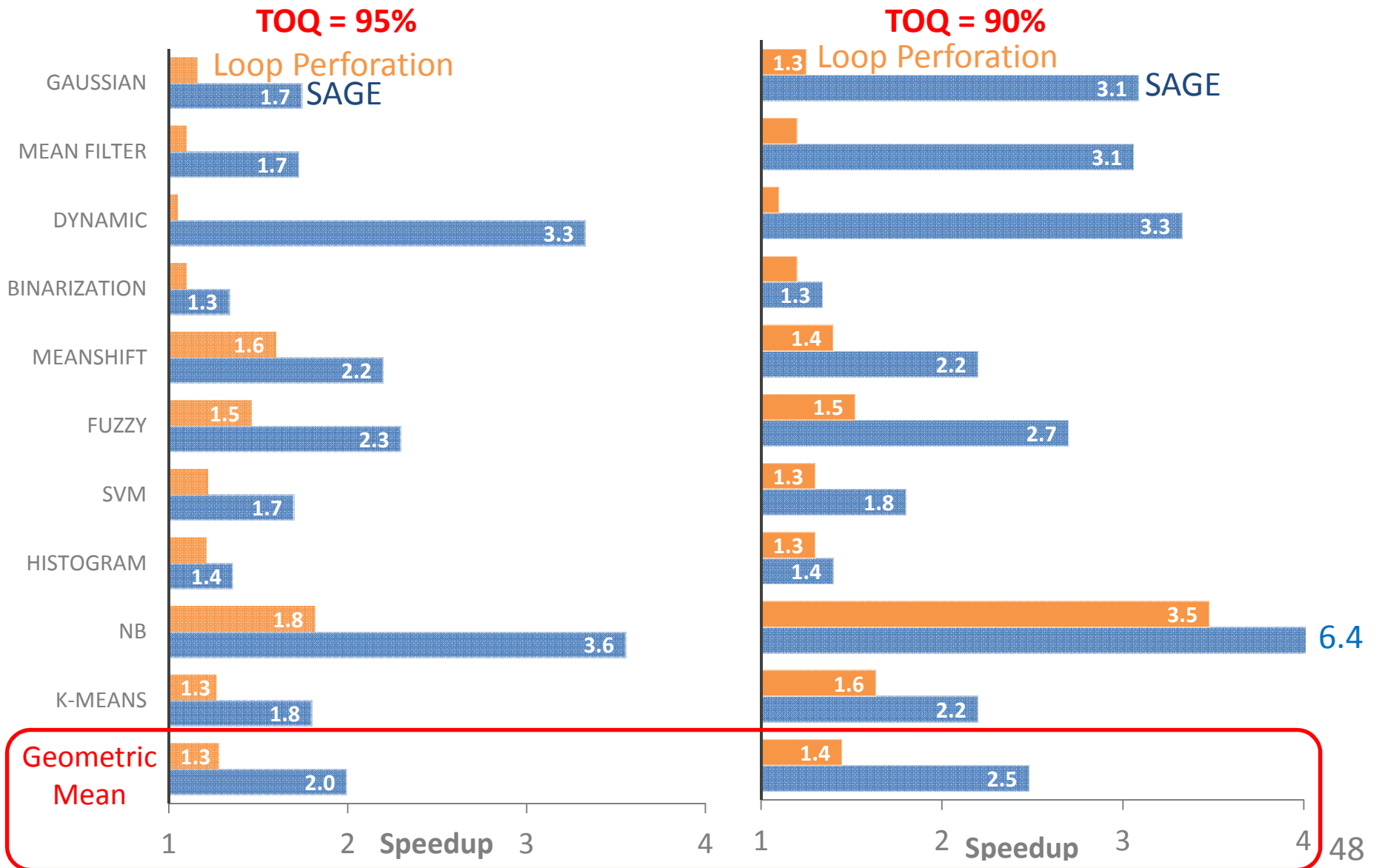# Confidence



After checking 50 samples, we will be 93% confident that 95% of the outputs satisfy the TOQ threshold

Beta   Uniform   Binomial

*posterior ∝ prior × likelihood*

46

# Calibration Overhead

# Performance



**TOQ = 95%**

**TOQ = 90%**

# Conclusion

- Automatic approximation is possible

- SAGE automatically generates approximate kernels with different parameters

- Runtime system uses tuning parameters to control the output quality during execution

- 2.5x speedup with less than 10% quality loss compared to the accurate execution

# SAGE: Self-Tuning Approximation for Graphics Engines

**Mehrzad Samadi**[1], Janghaeng Lee[1], D. Anoushe Jamshidi[1], Amir Hormati[2], and Scott Mahlke[1]
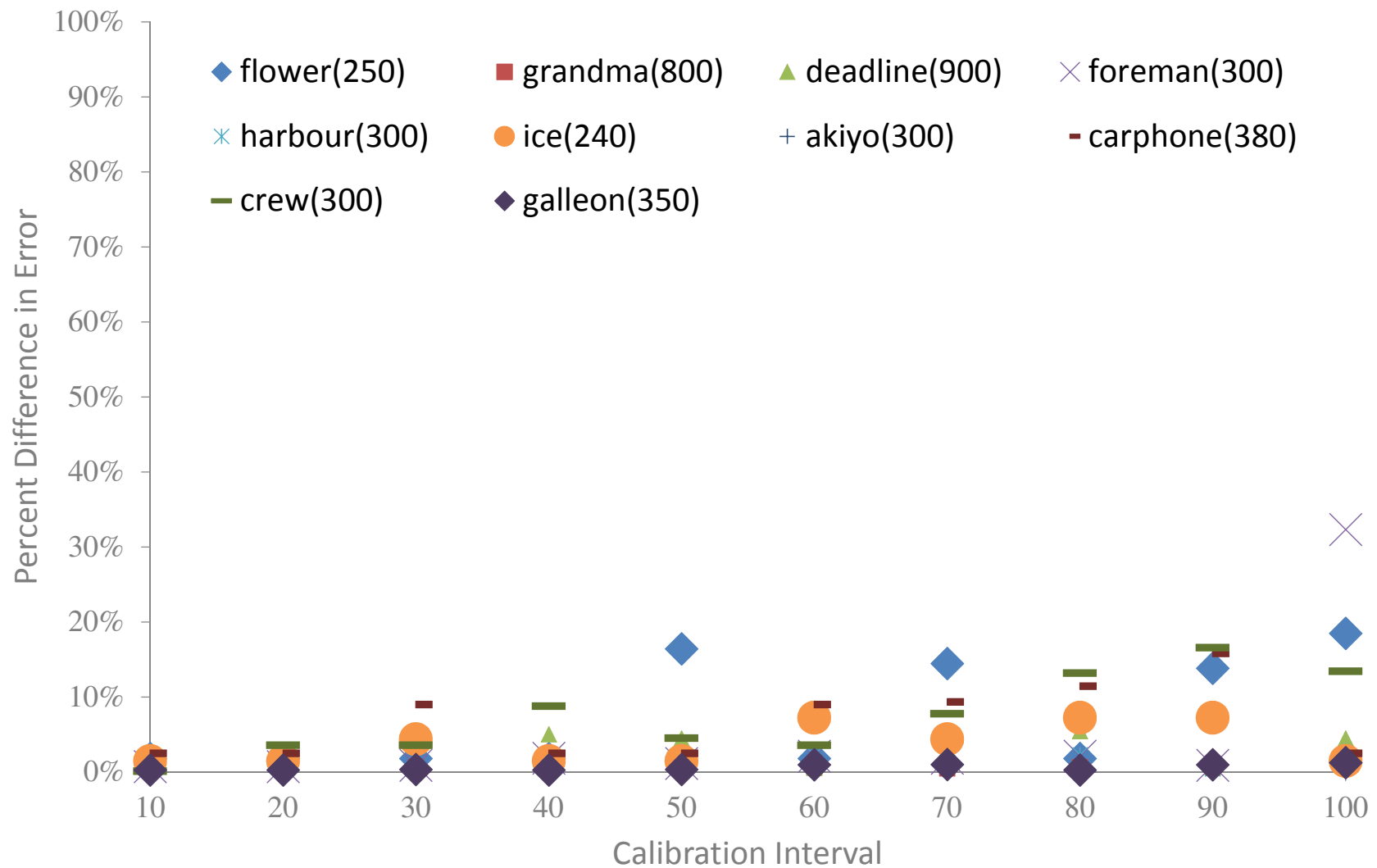
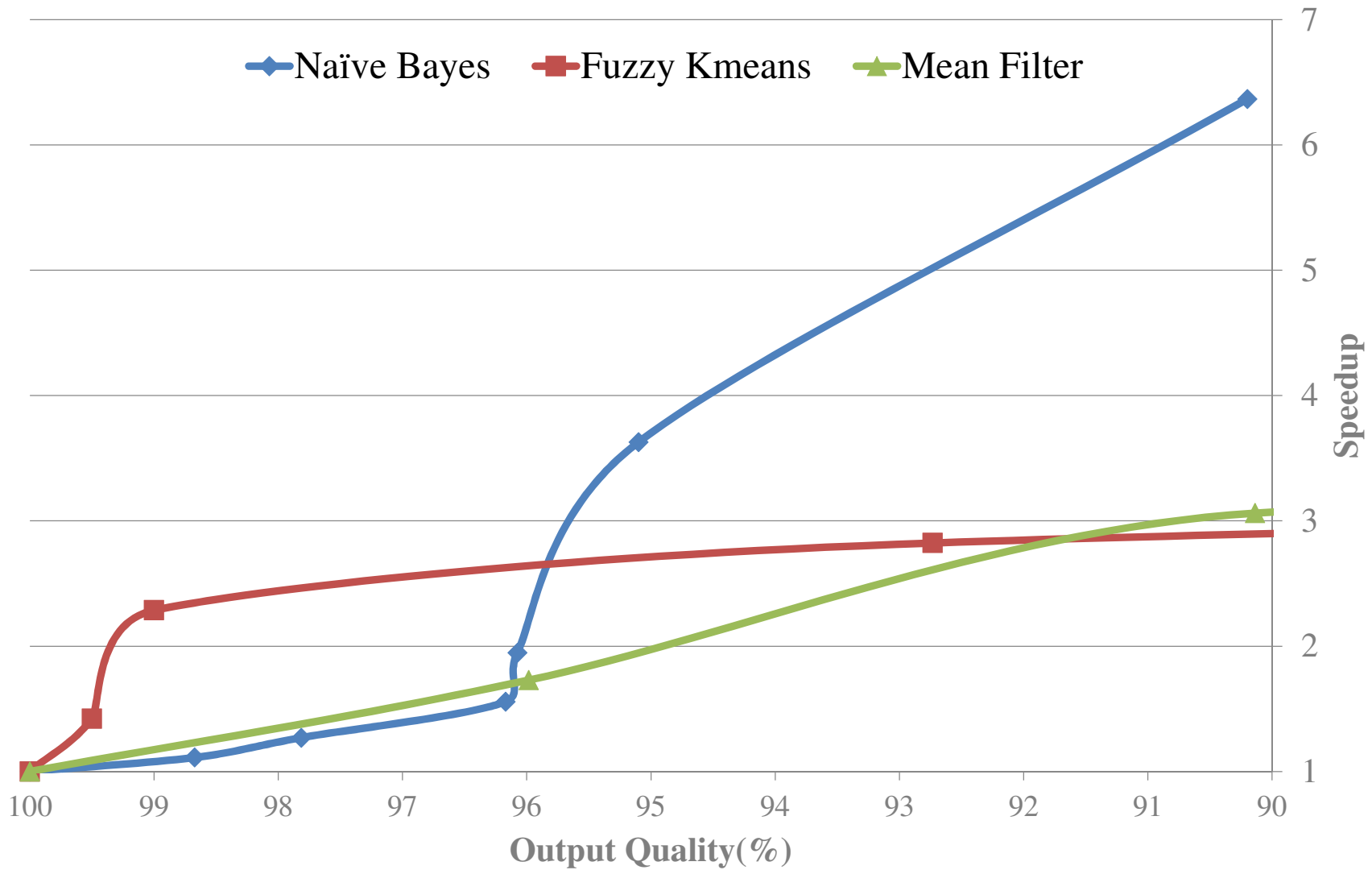University of Michigan[1], Google Inc.[2]

December 2013

# Backup Slides

# What Does Calibration Miss?

# SAGE Can Control The Output Quality

# Distribution of Errors