

Autotuning OpenCL Workgroup Size for Stencil Patterns

Chris Cummins



THE UNIVERSITY *of* EDINBURGH
informatics

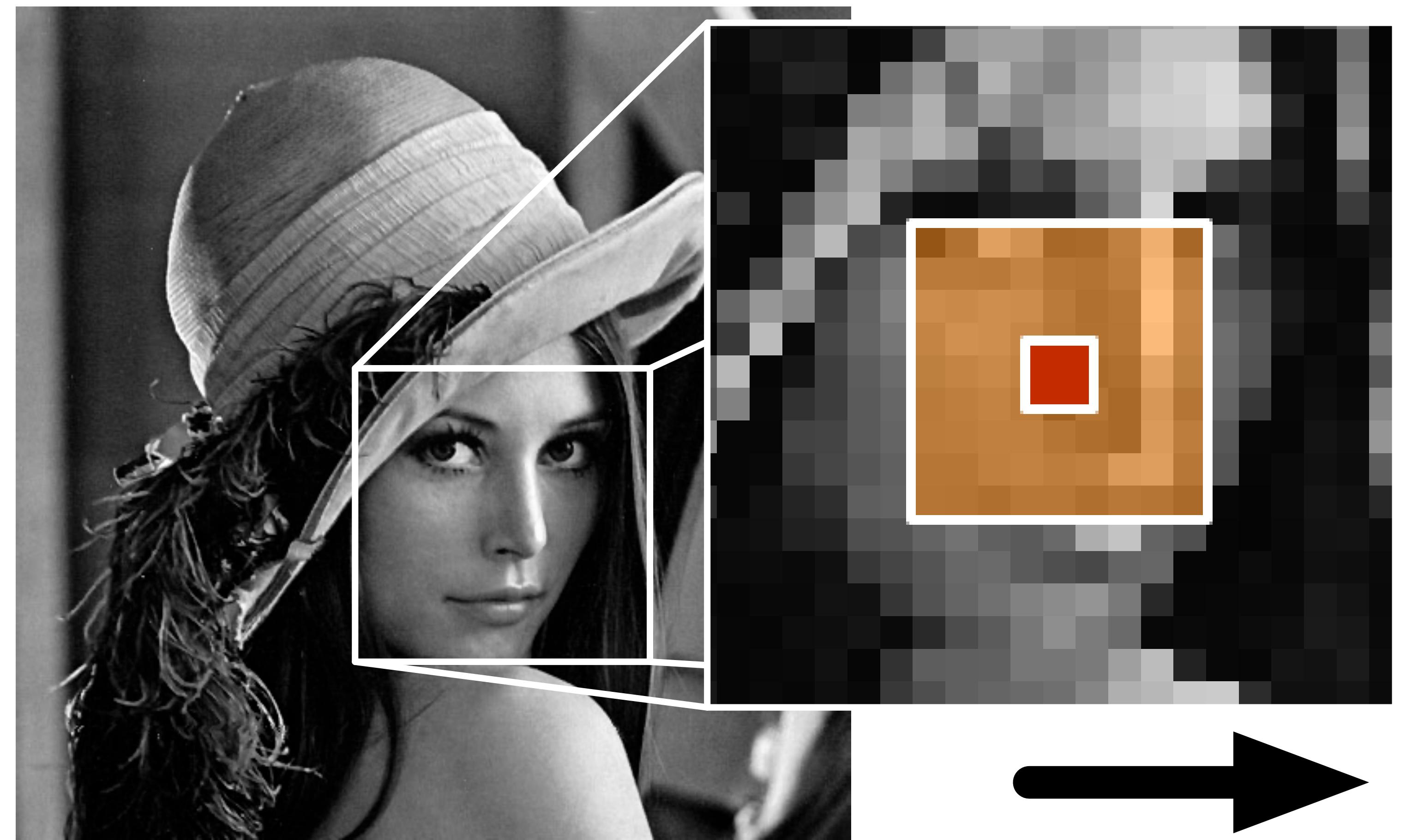
**EPSRC Centre for Doctoral Training in
Pervasive Parallelism**

EPSRC
Engineering and Physical Sciences
Research Council

<http://chriscummins.cc>

**Stencils &
Workgroup
size**

**Stencils &
Workgroup
size**



input

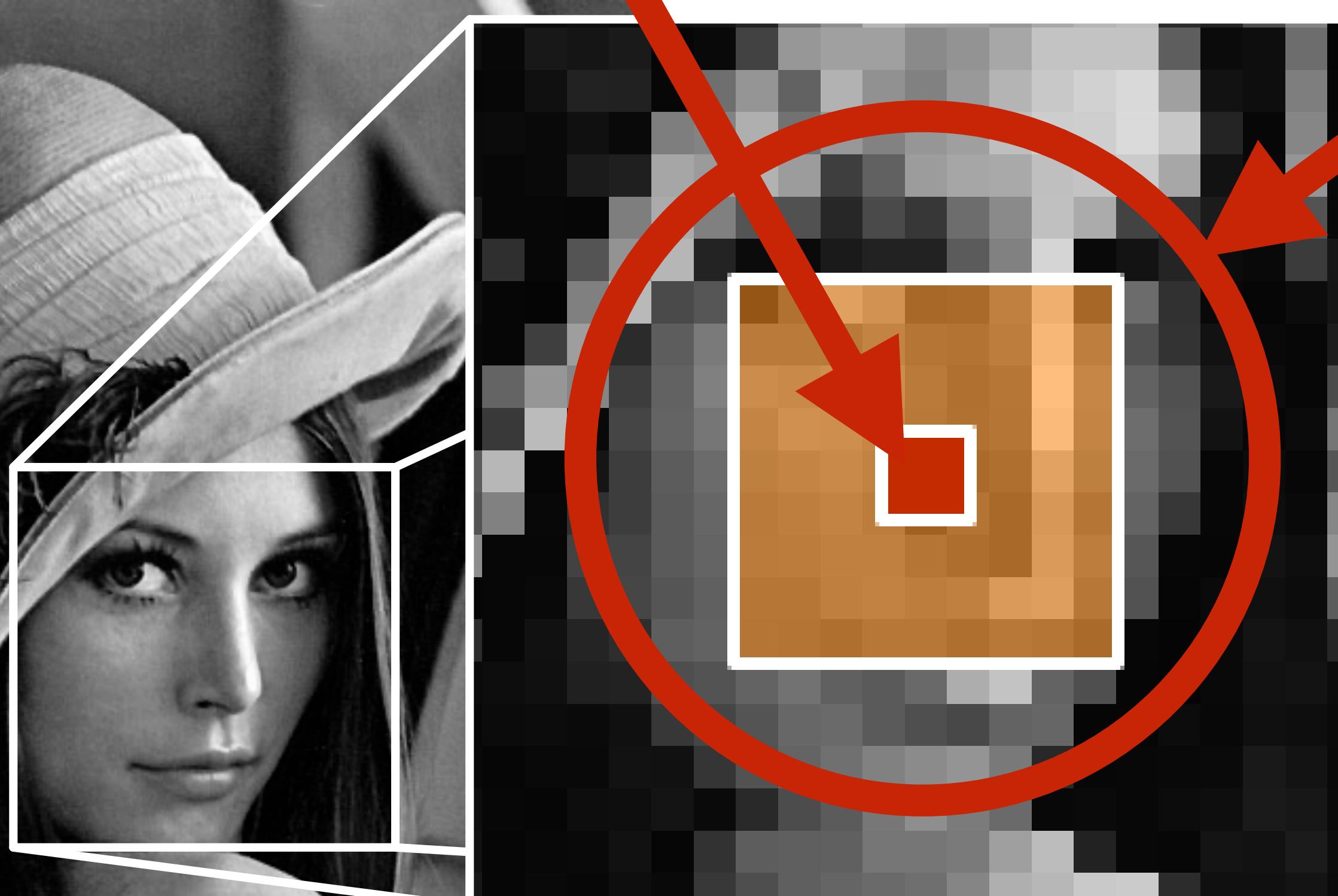
stencil



output

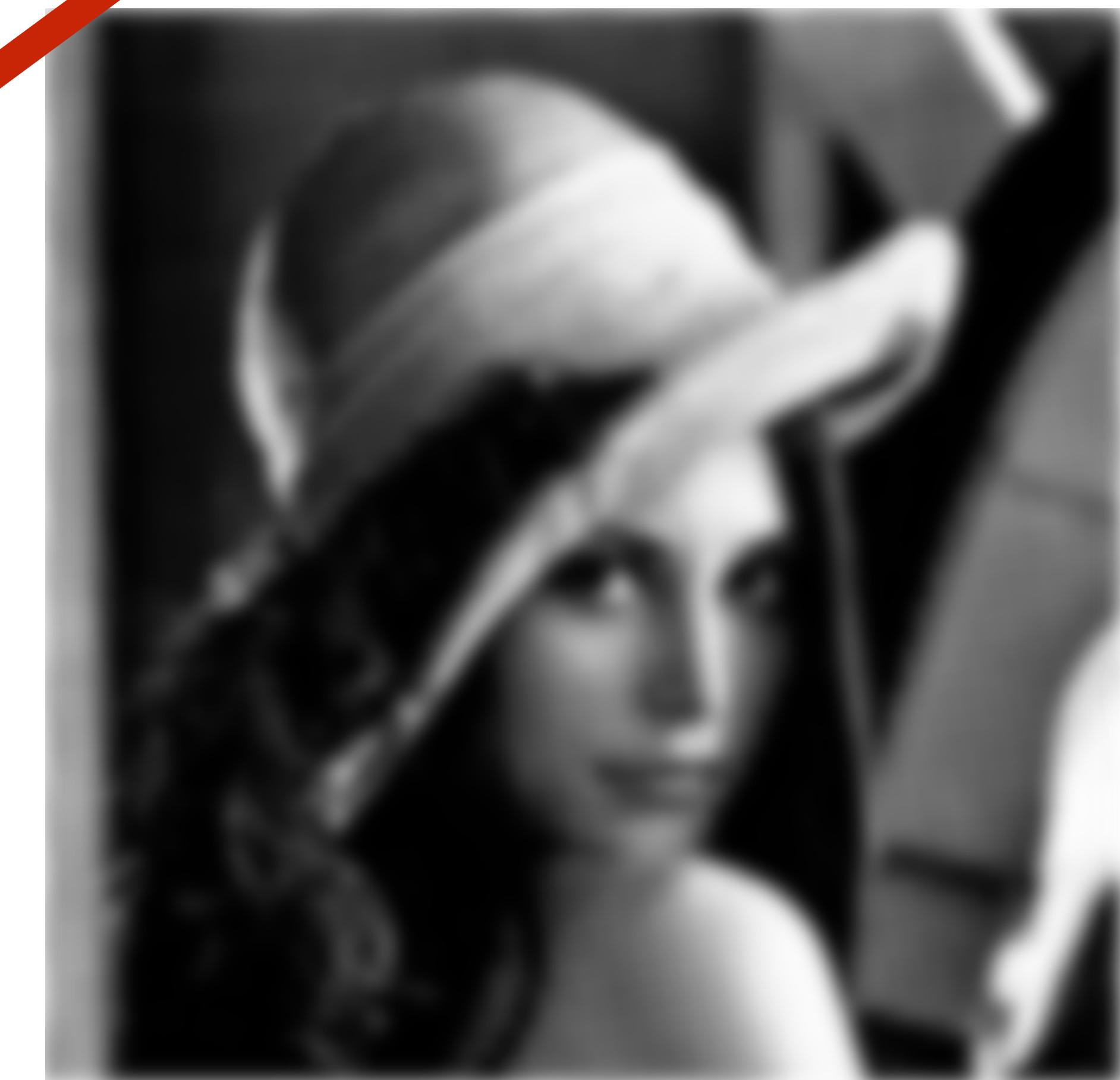
element

border region



input

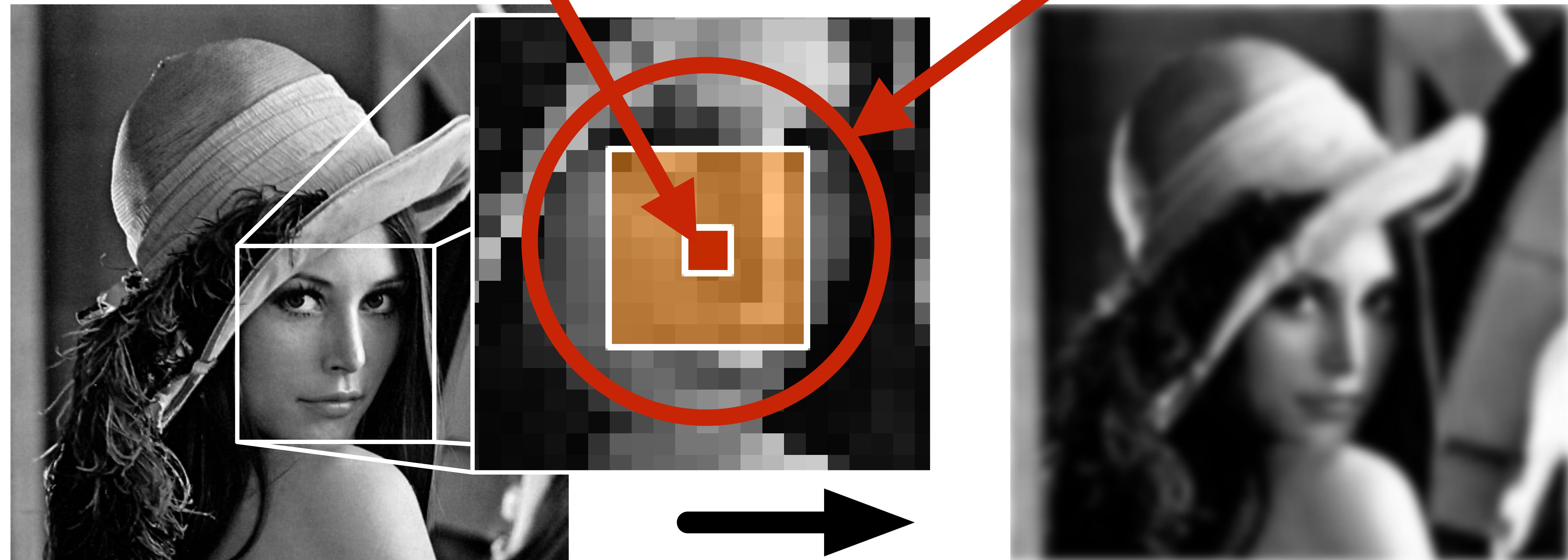
stencil



output

10^6 elements

10^6 border regions



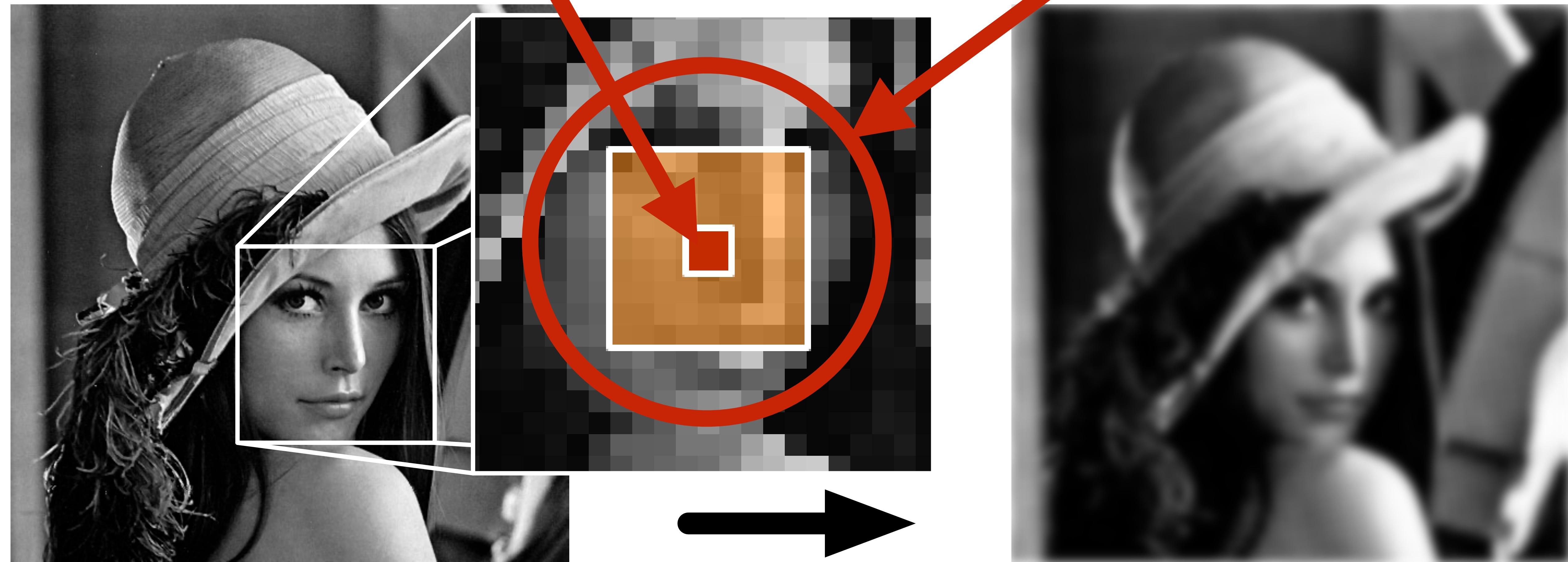
input

stencil

output

10^6 elements

10^6 border regions



input

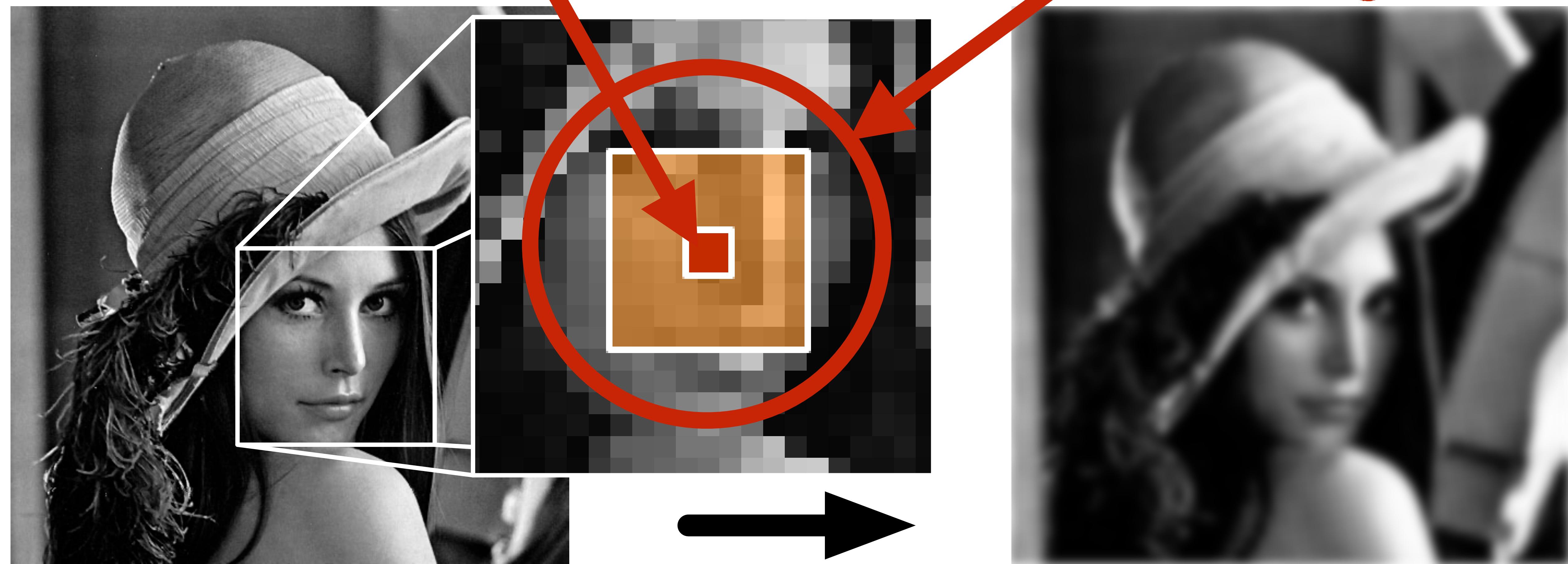
stencil

output

Multiple independent computations

10^6 elements

10^6 border regions

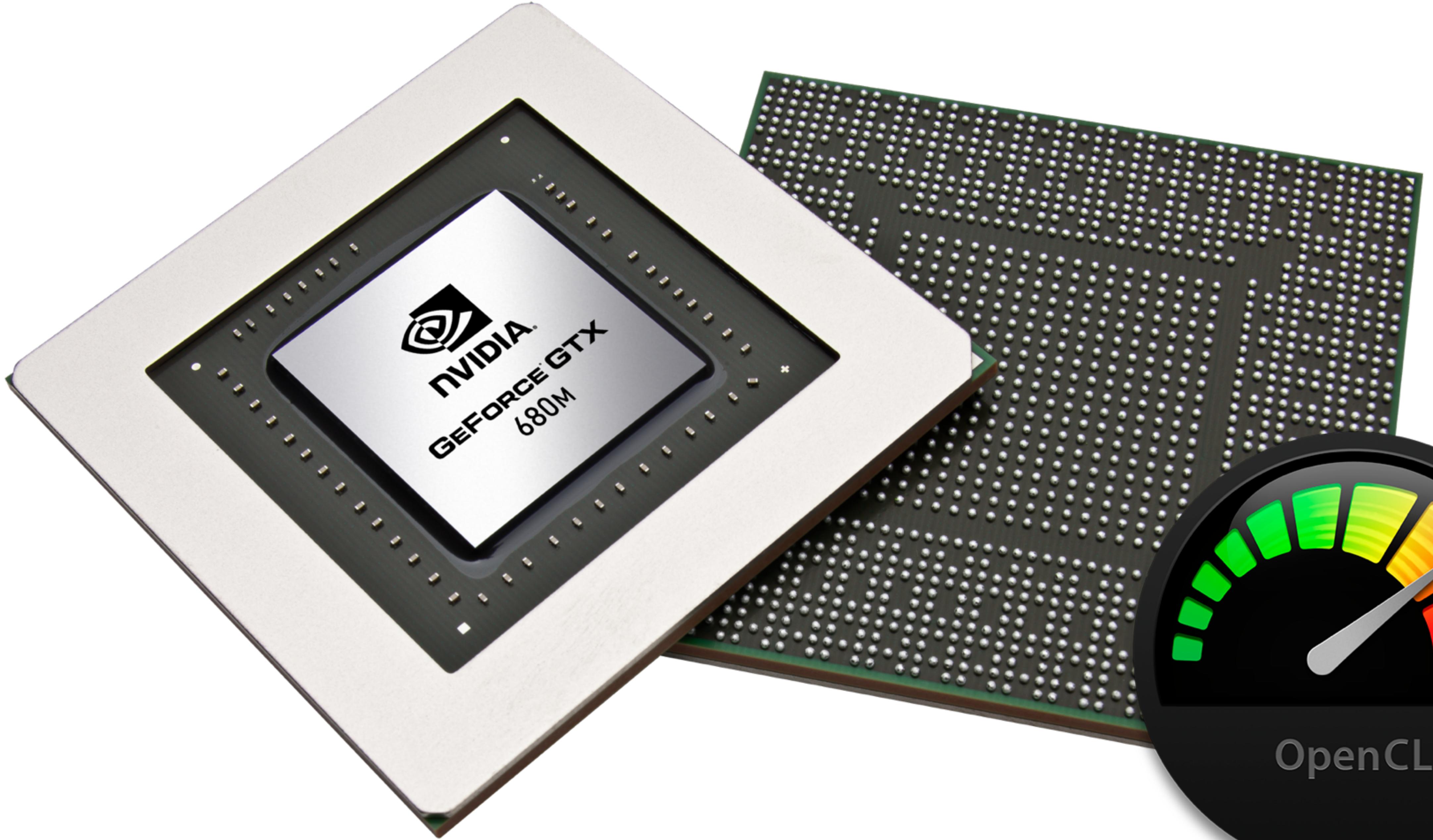


input

stencil

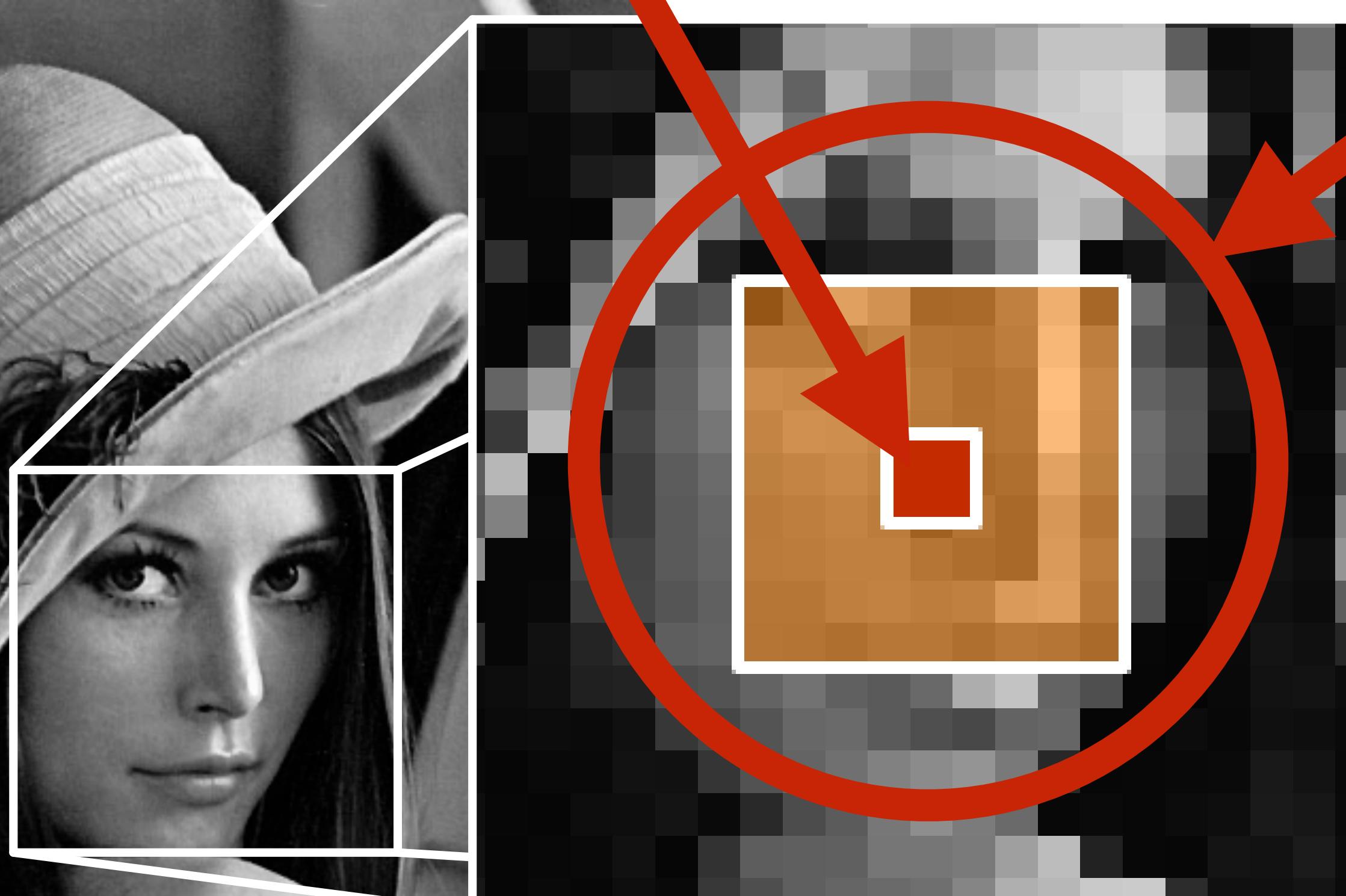
output

Multiple (overlapping) memory accesses



element

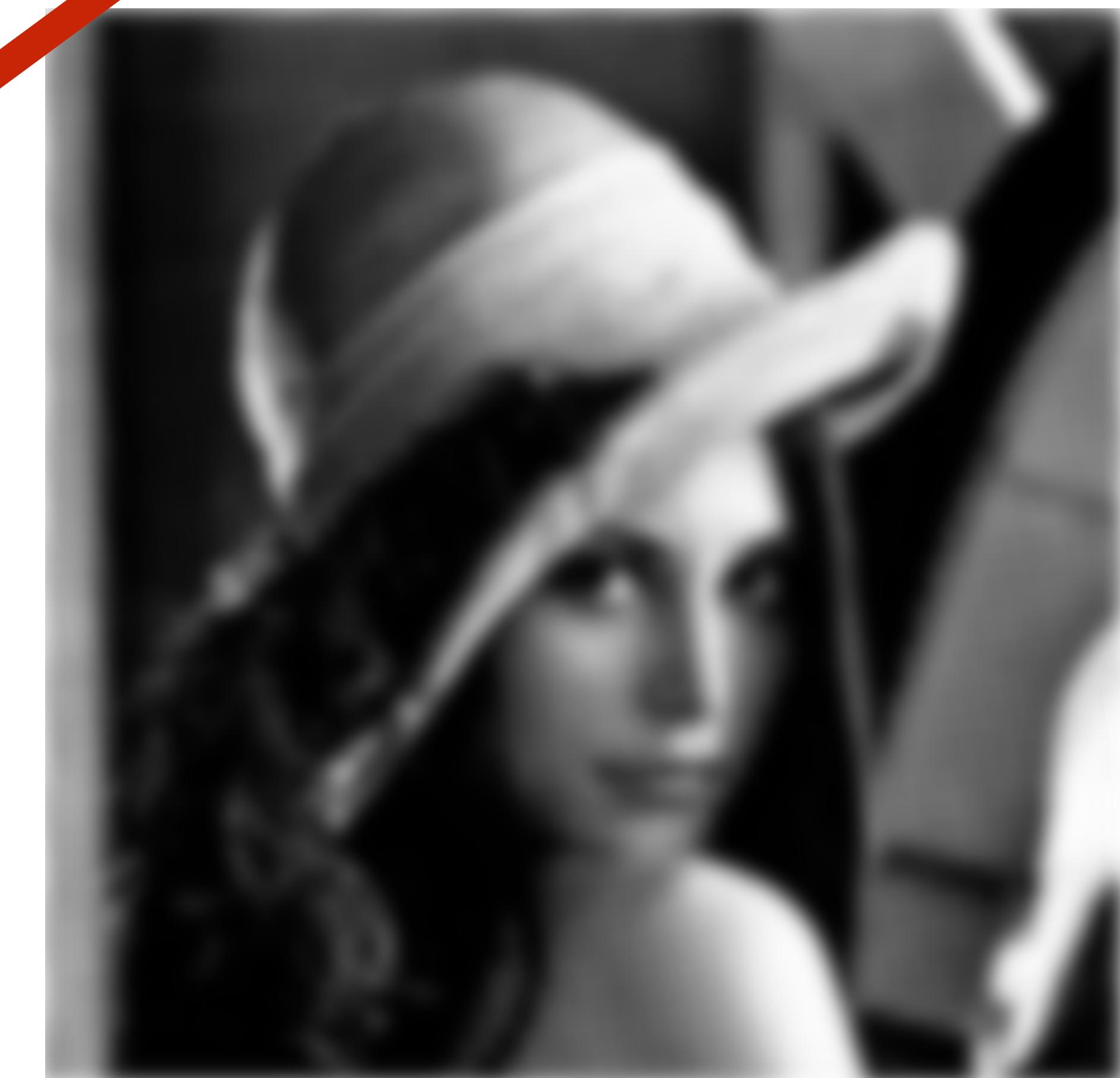
border region



input

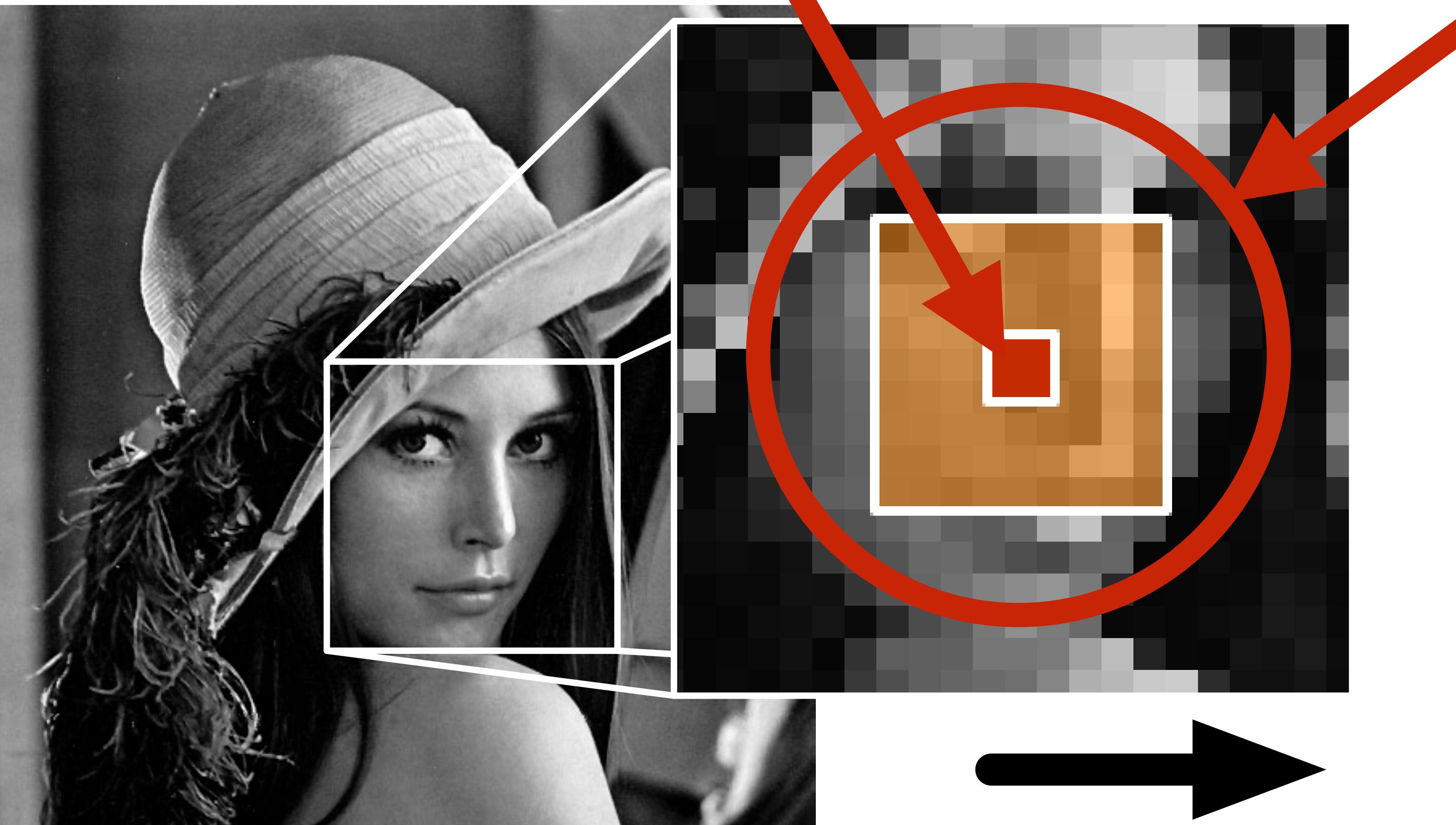
stencil

output



element

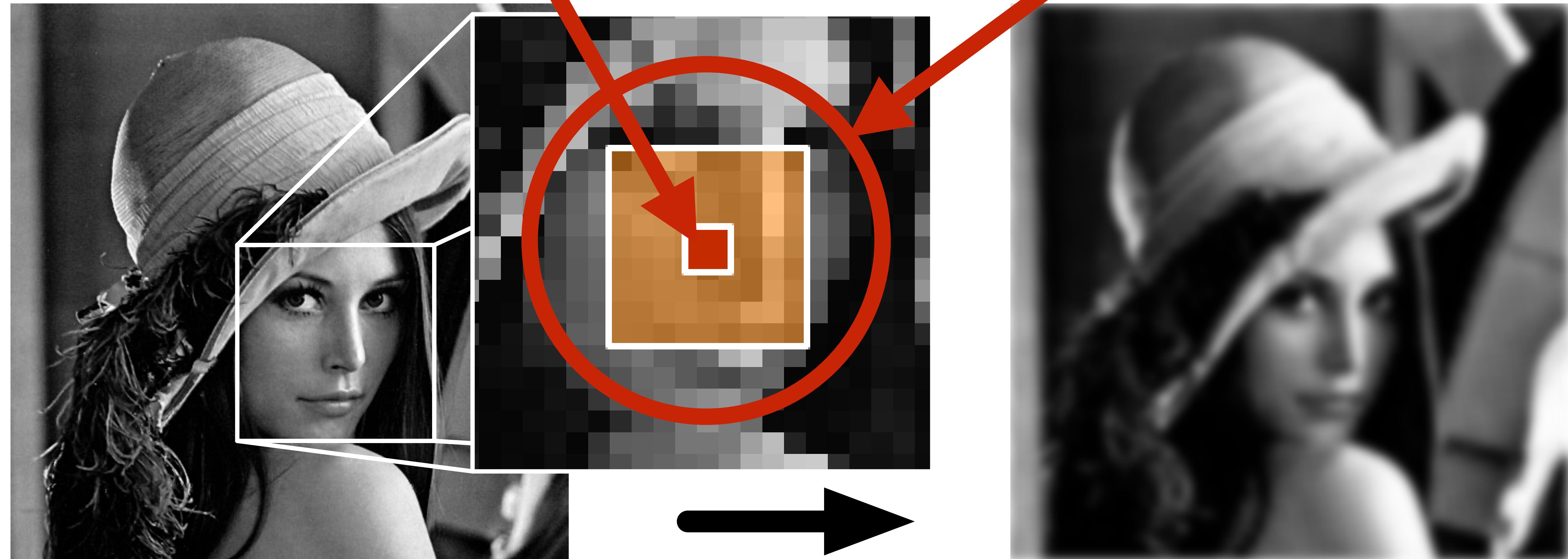
border region



input

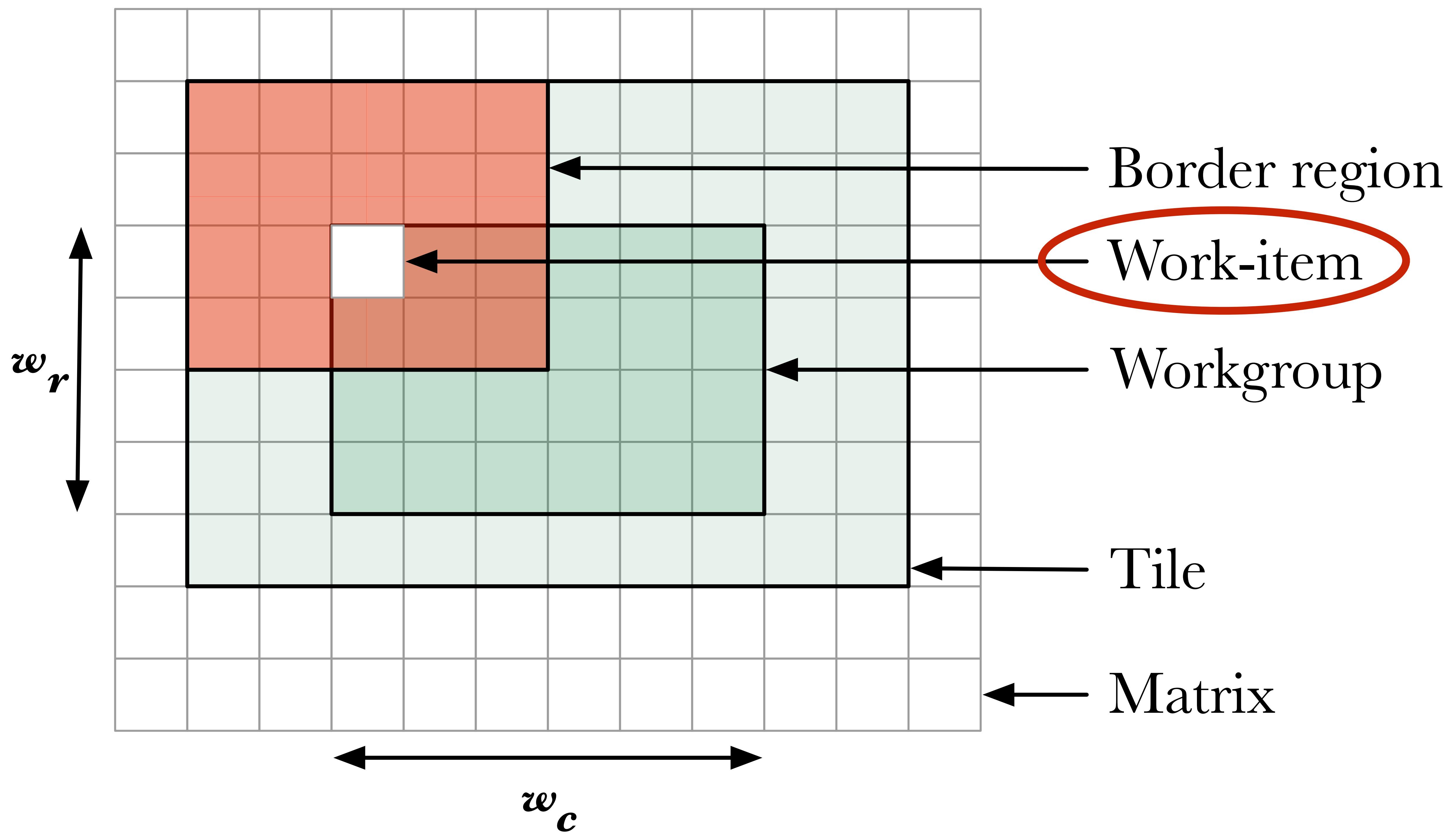
stencil *kernel* output

~~-element work-item border region~~



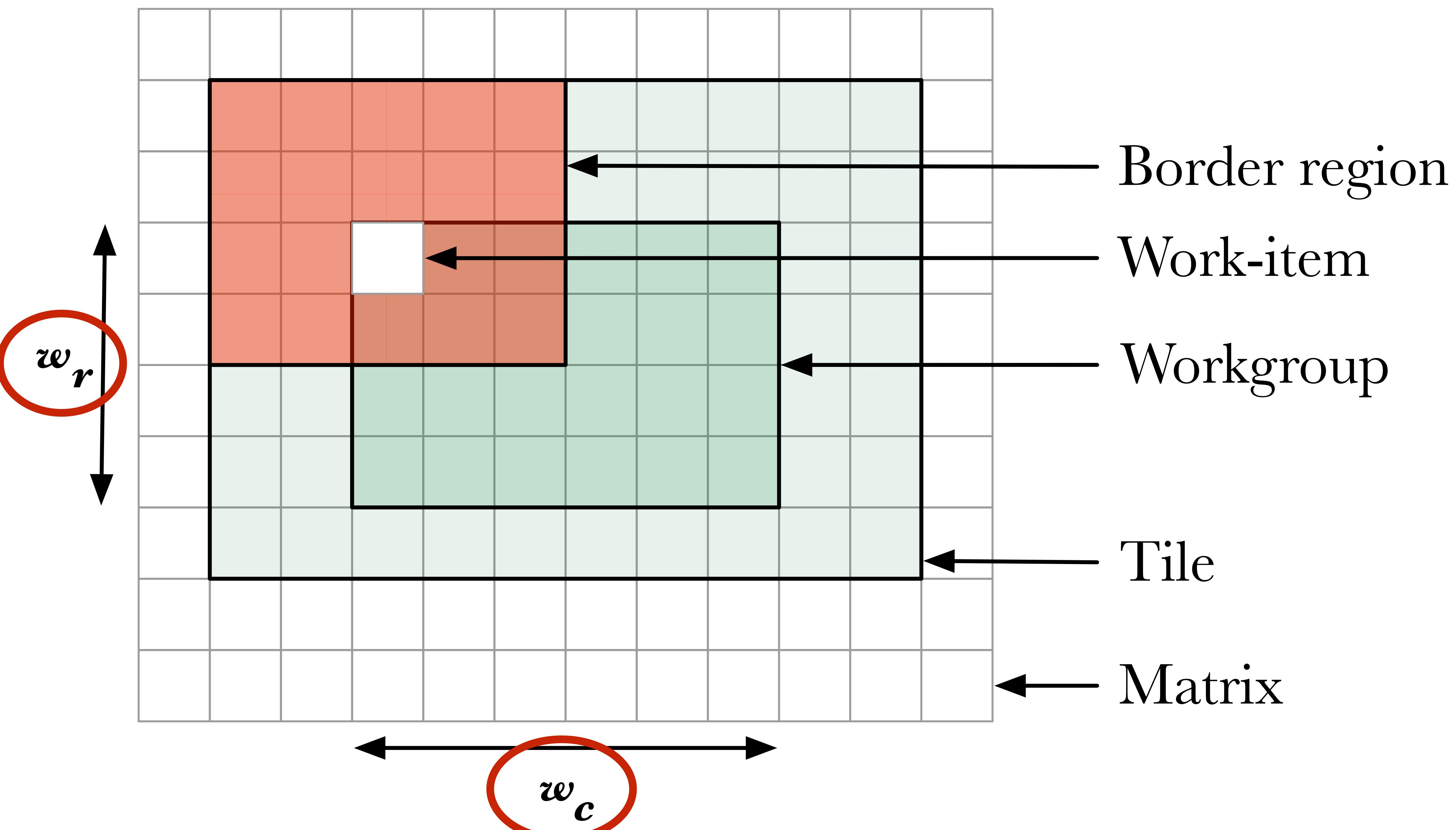
input

stencil ~~kernel~~ output



**Stencils &
Workgroup
size**

**Stencils &
Workgroup
size**



Workgroup size affects
mapping to SIMD hardware.
device occupancy.
local memory utilisation.

**Pop
Quiz!**

What is the best workgroup size for ...

Gaussian blur, 512px x 512px,
floats, on:

1. AMD HD7990?
2. Nvidia GTX Titan?
3. Intel i7-3820?

What is the best workgroup size for ...

Gaussian blur, 512px x 512px,
floats, on:

1. AMD HD7990? 64×4
2. Nvidia GTX Titan? 96×4
3. Intel i7-3820? 40×24

What is the best workgroup size for ...

Nvidia GTX 590, 4096 x 4096 elements
running:

1. Sobel edge detection?
2. Heat equation?
3. Game of life?

What is the best workgroup size for ...

Nvidia GTX 590, 4096 x 4096 elements
running:

1. Sobel edge detection? 256×2
2. Heat equation? 128×2
3. Game of life? 32×6

What is the best workgroup size for ...

1. Intel i5-2430, game of life,
 4096×4096 ?
2. Nvidia GTX 690, threshold,
 512×512 ?
3. Intel i7-3820, NMS, 512×512 ?

What is the best workgroup size for ...

1. Intel i5-2430, game of life,
 $4096 \times 4096?$ *196×20*
2. Nvidia GTX 690, threshold,
 $512 \times 512?$ *32×4*
3. Intel i7-3820, NMS, $512 \times 512?$
 88×8

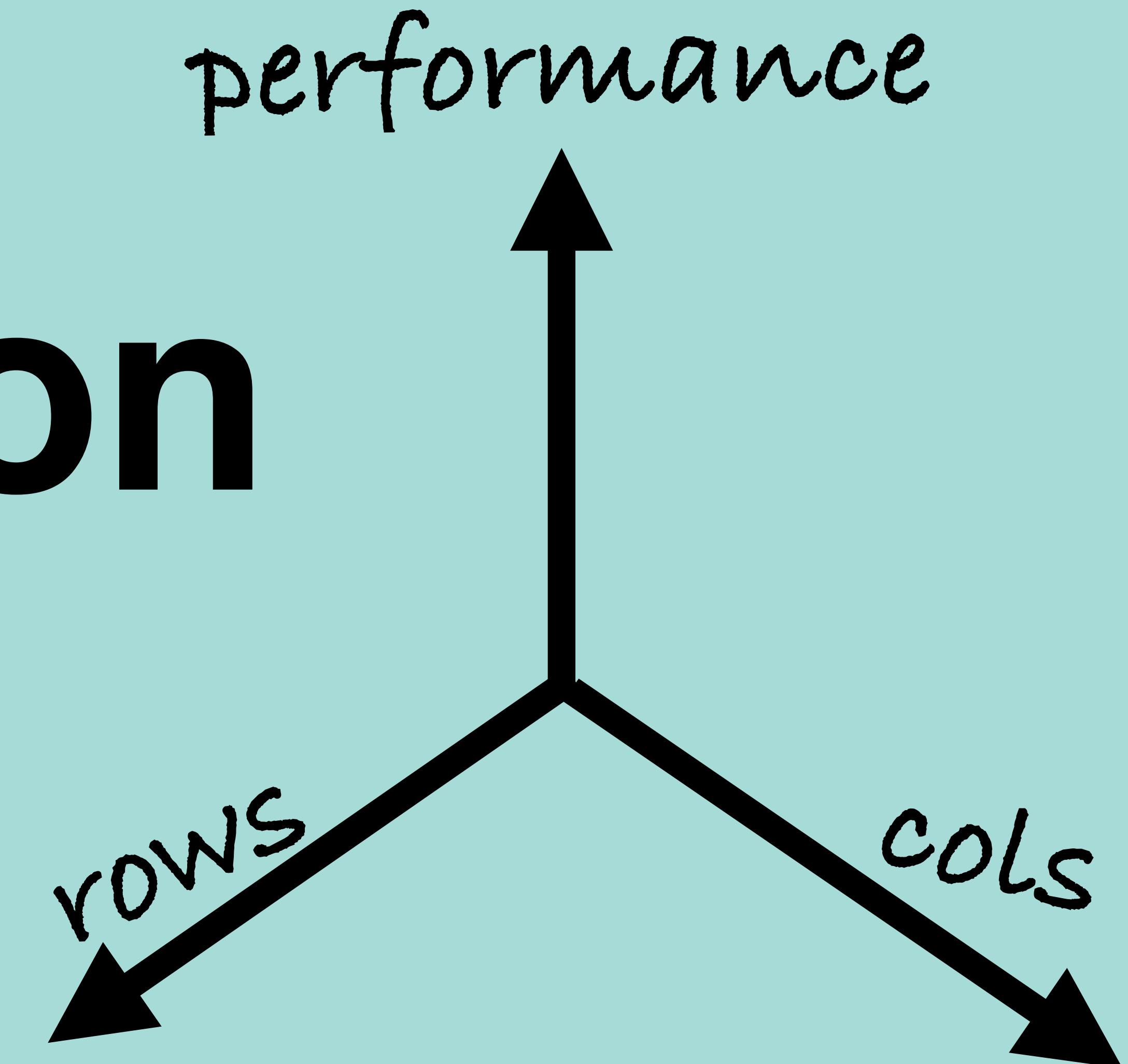
**One size
does not
fit all!**

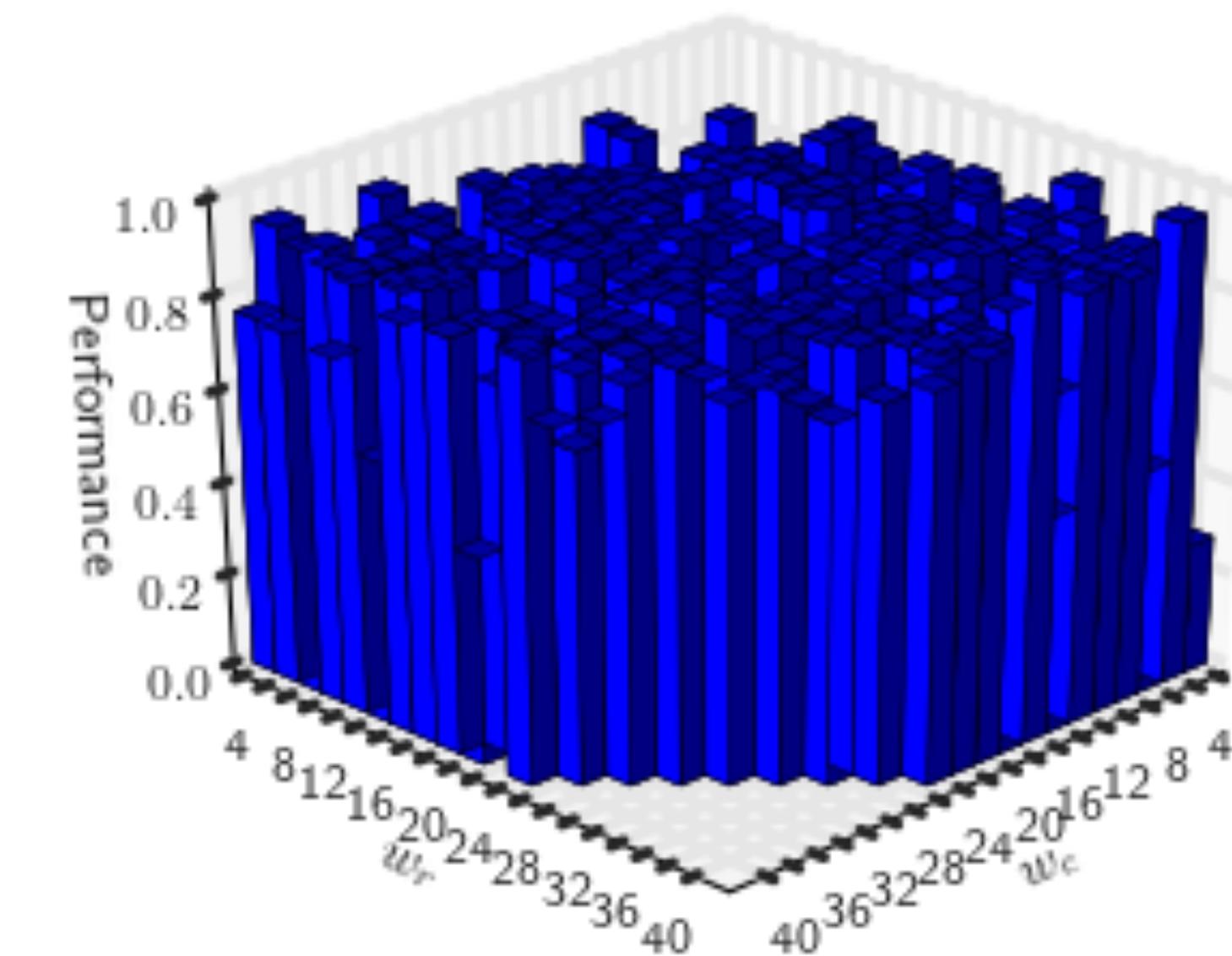
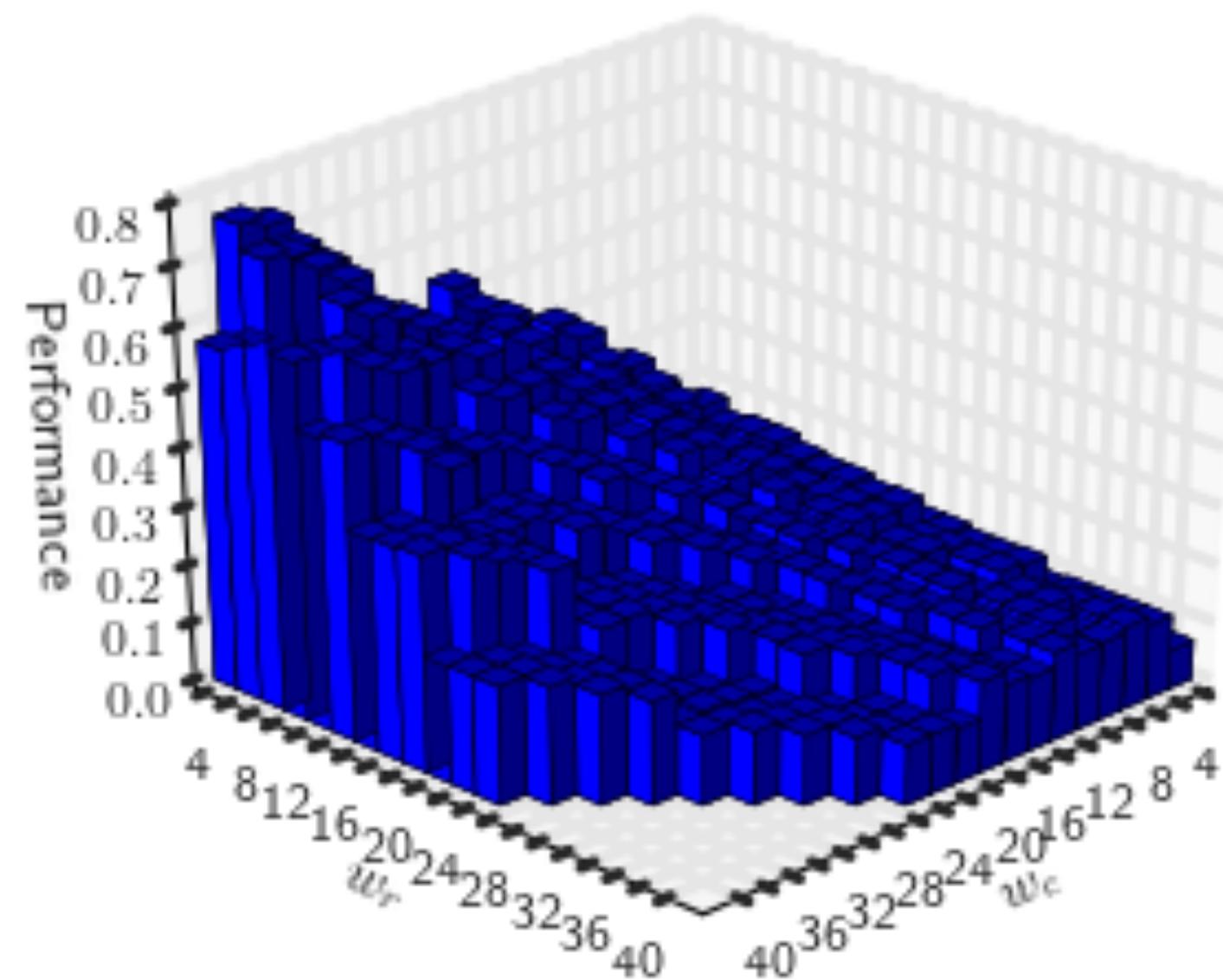
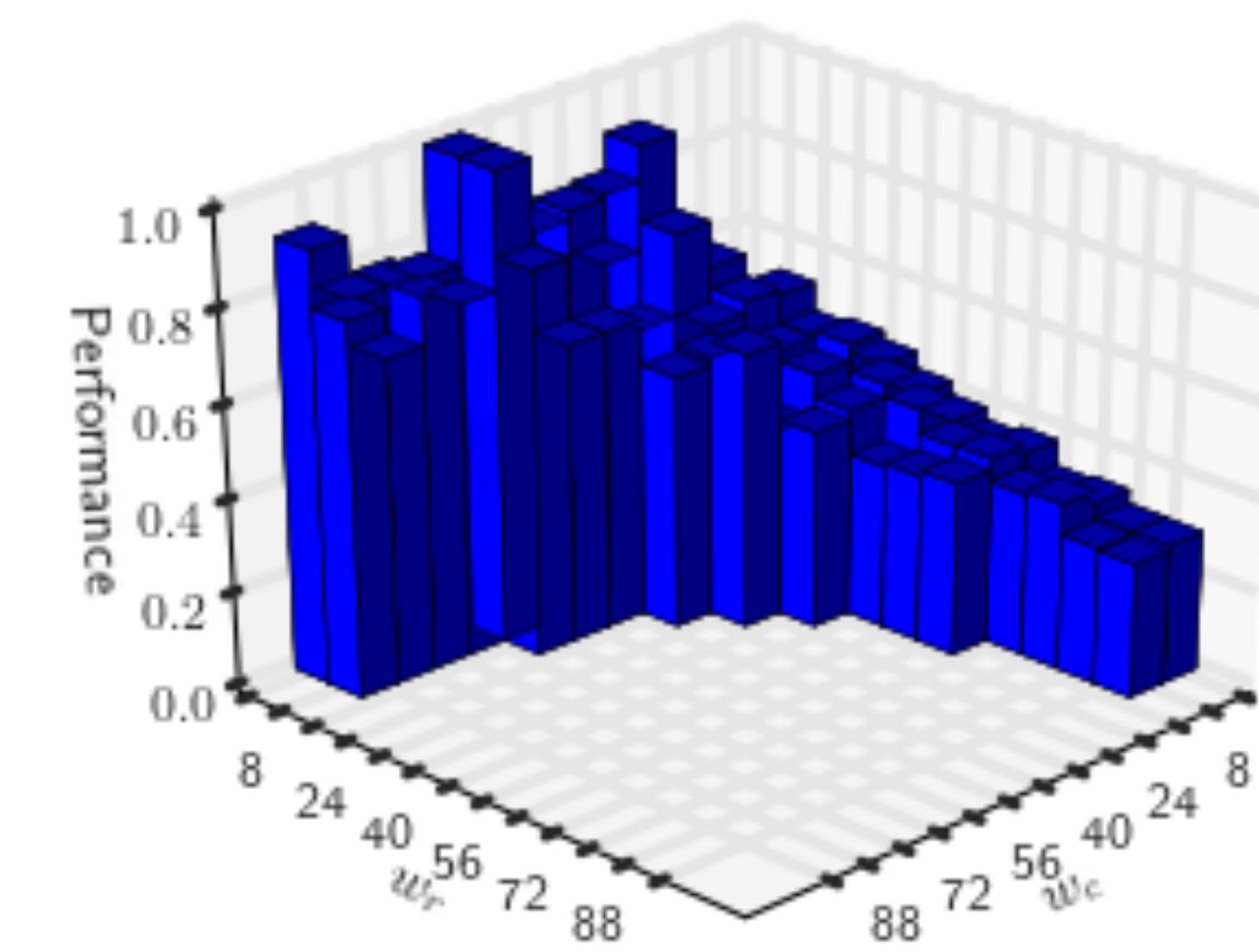
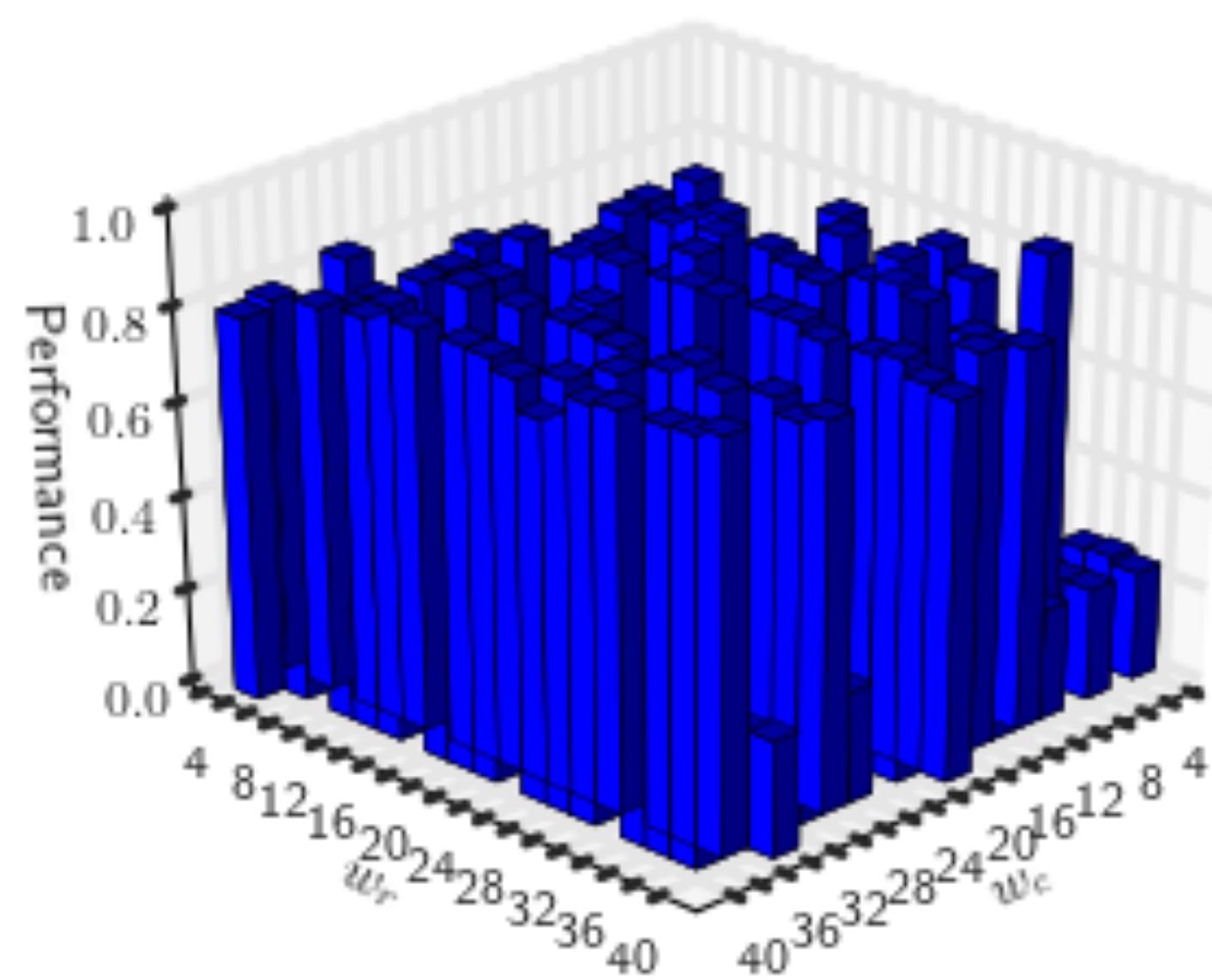
Choosing workgroup size depends on:

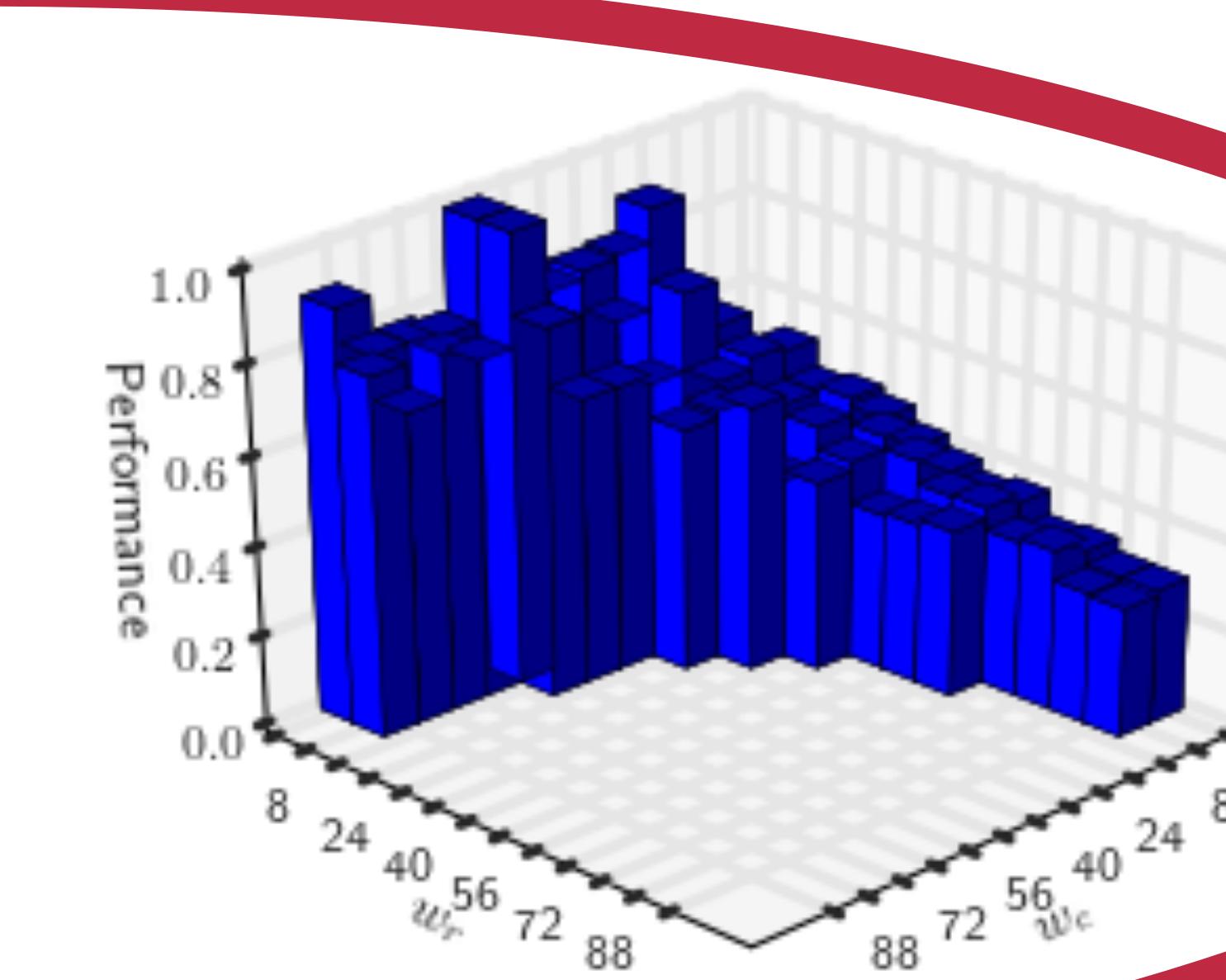
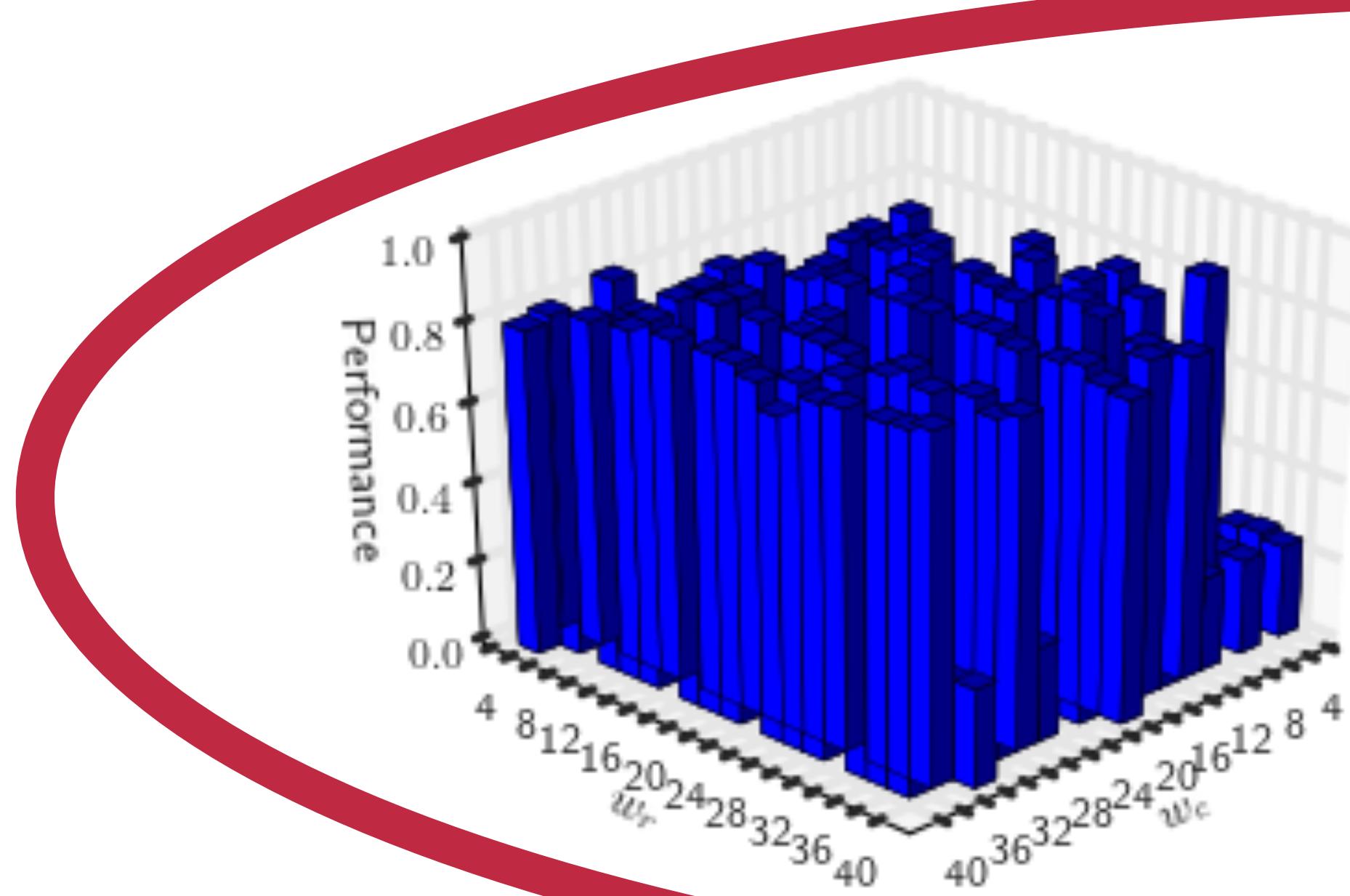
1. Device
2. Program
3. Dataset

Optimisation

space



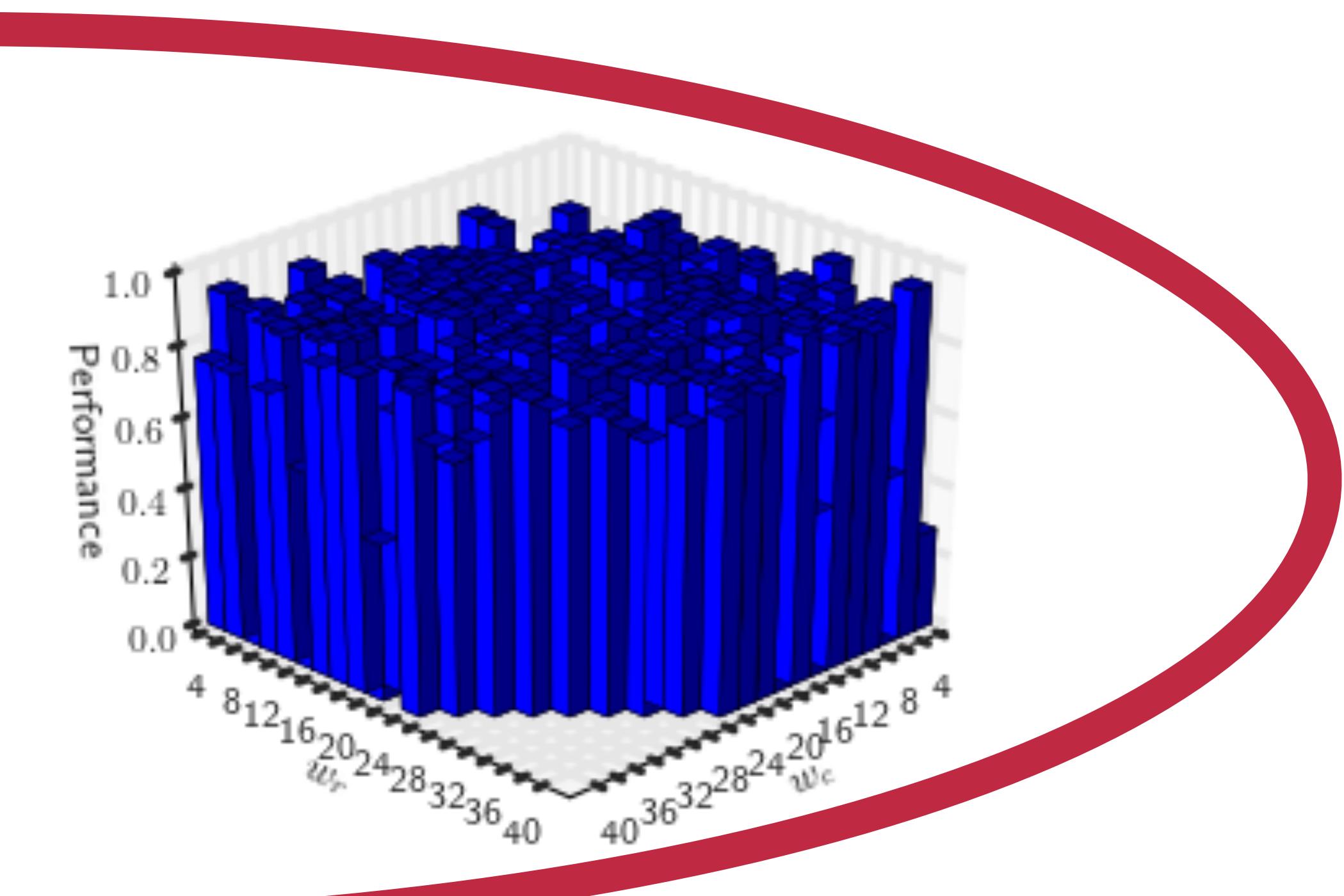
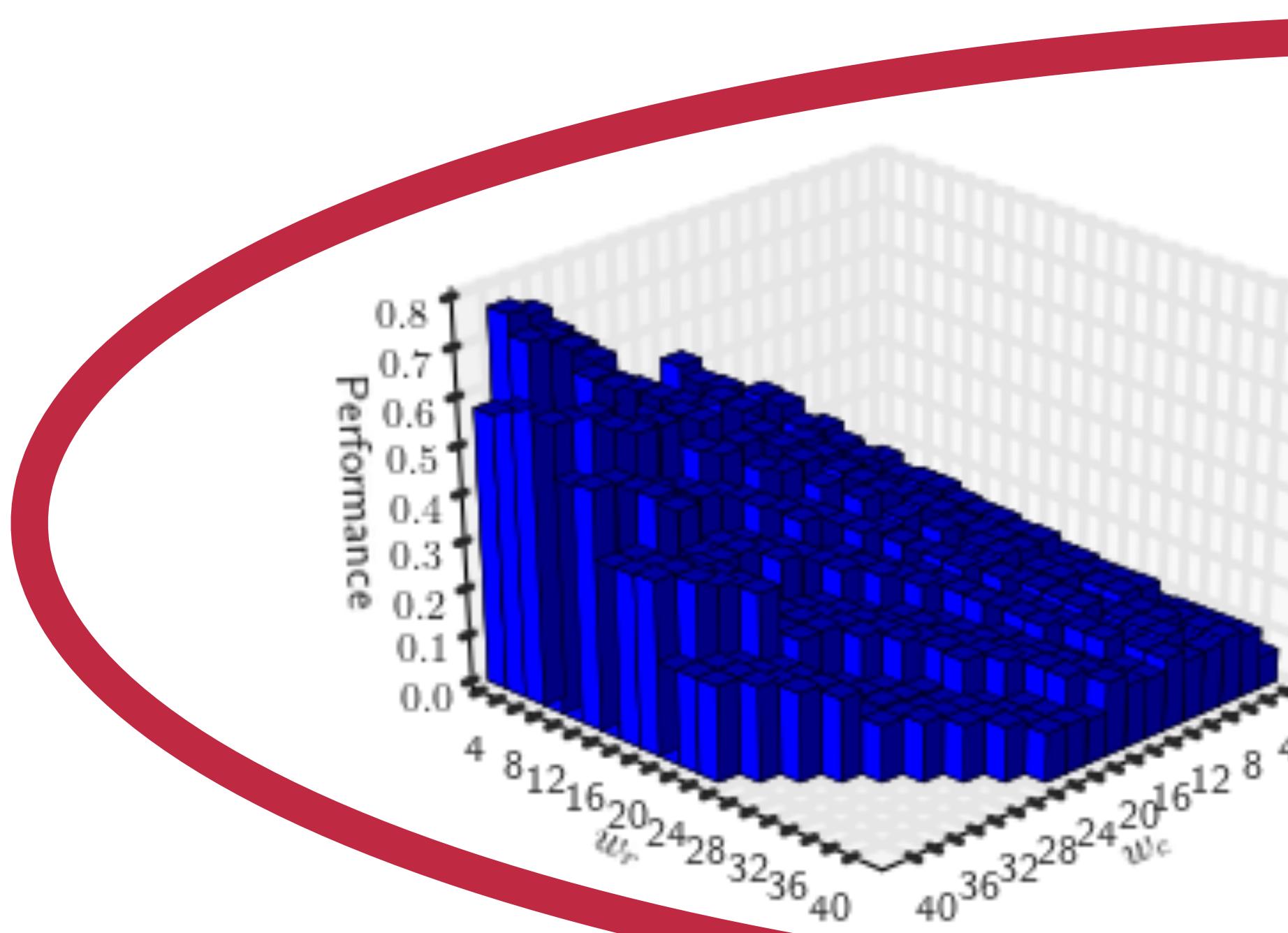


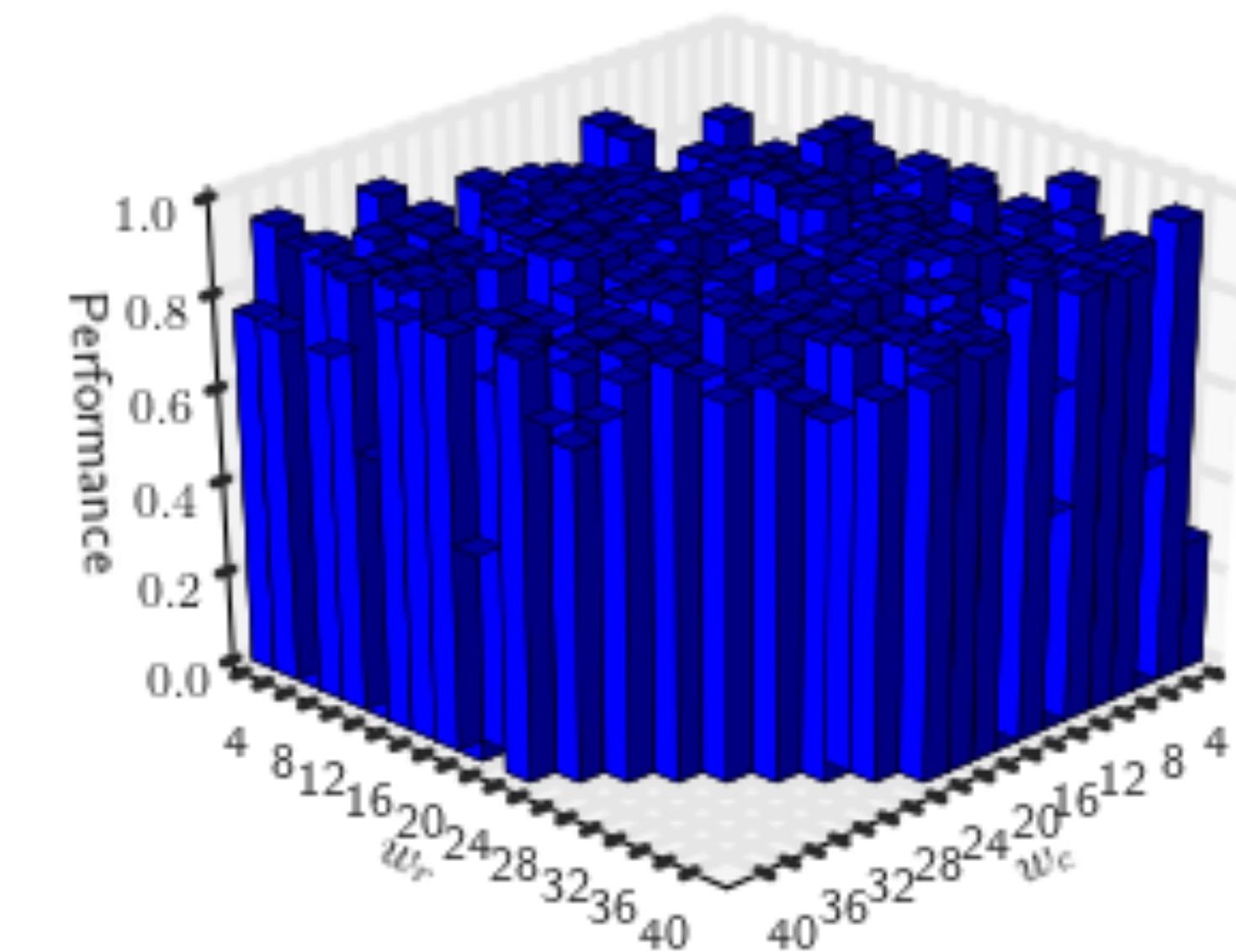
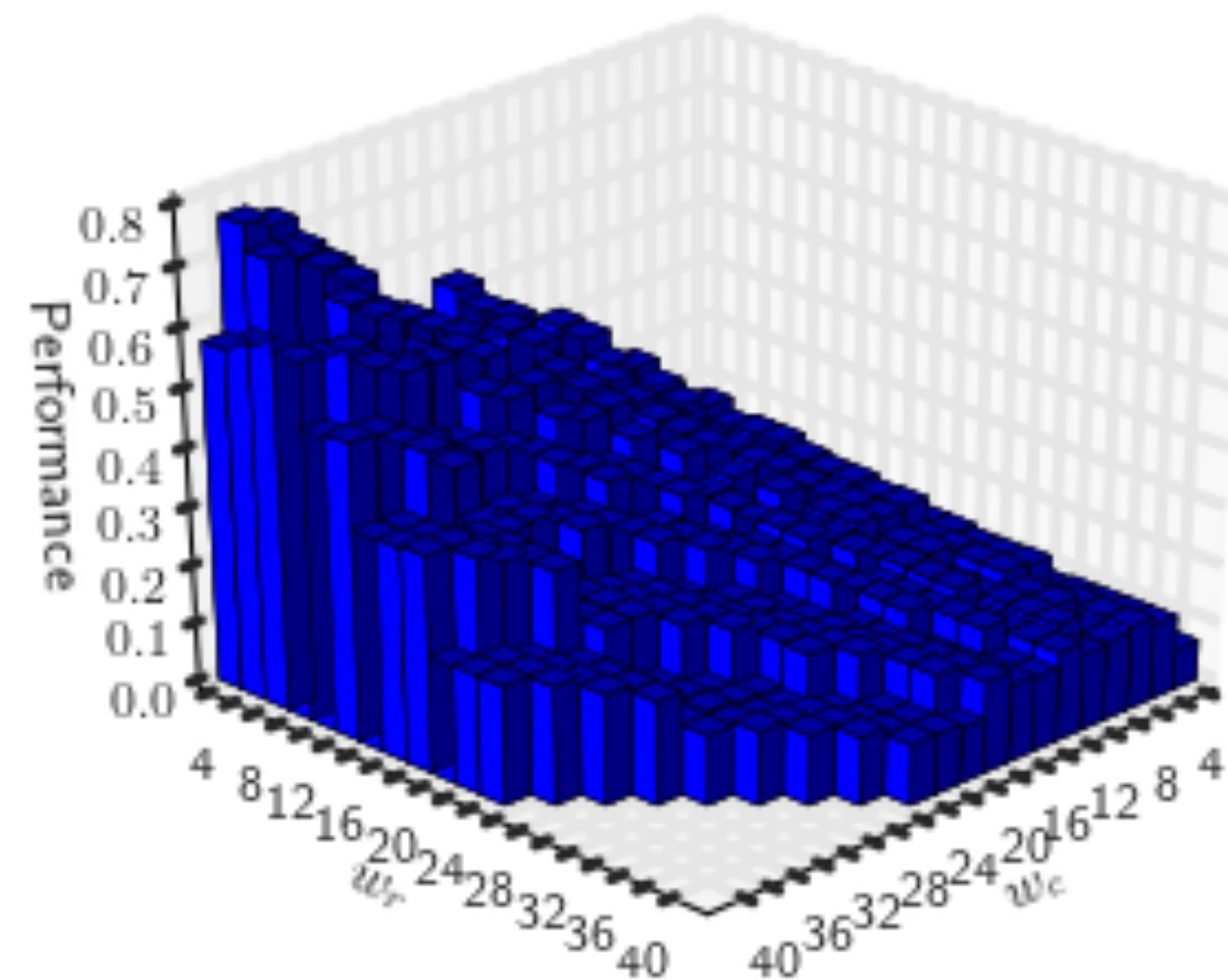
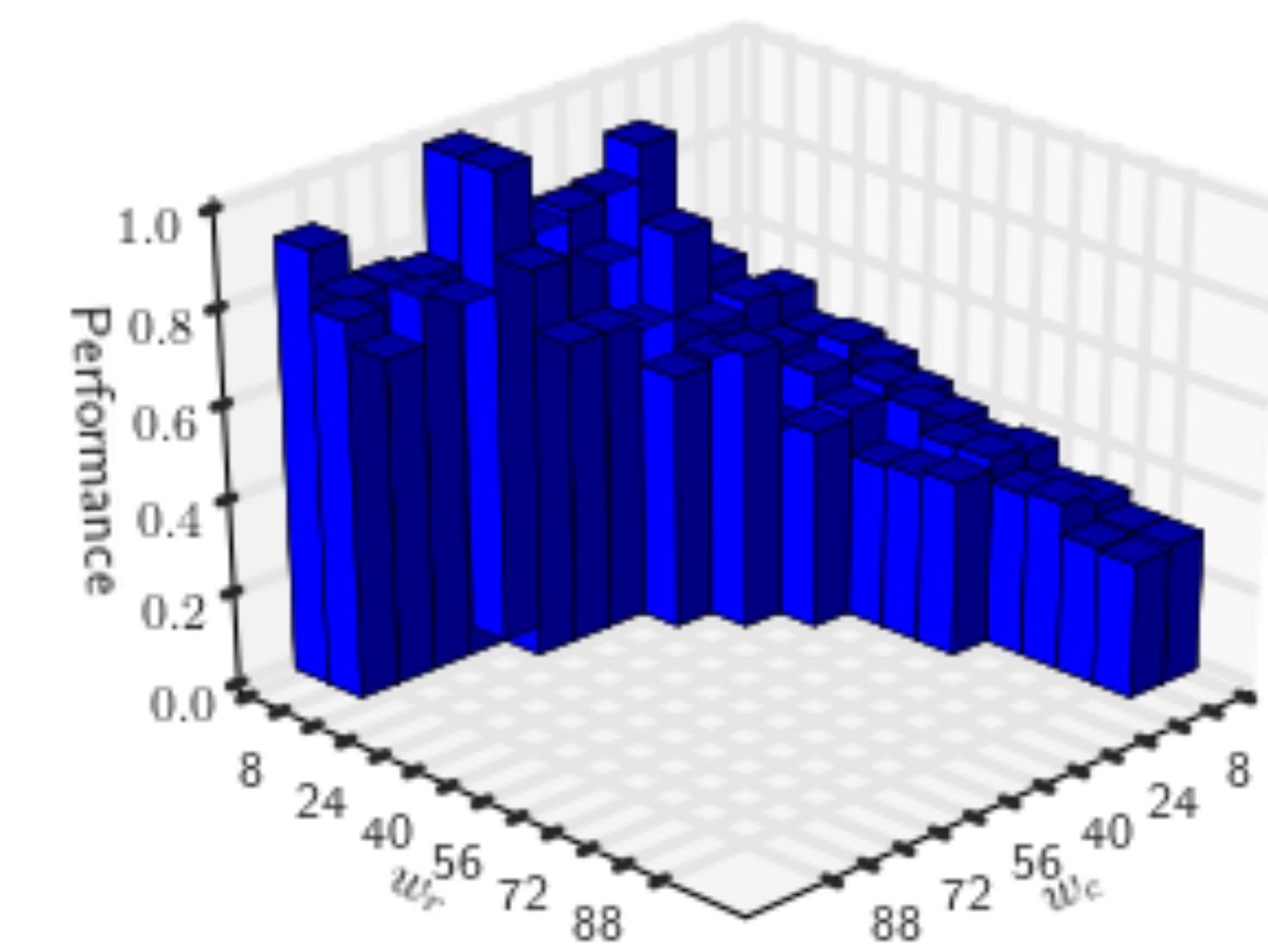
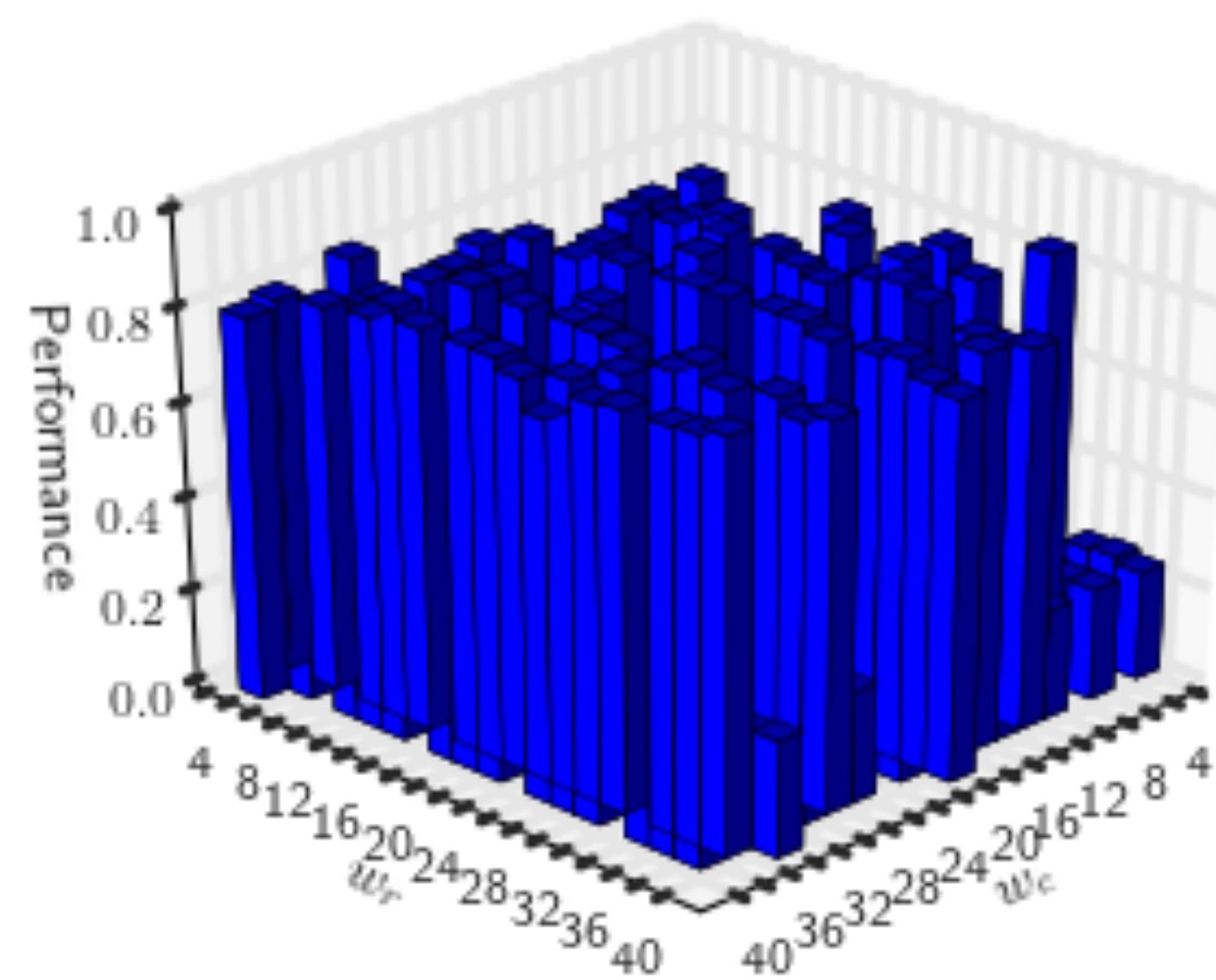


Same stencil!
Different device!

The figure consists of two 3D bar charts side-by-side, enclosed within a large red oval. Overlaid on the bottom chart is the text "Same stencil!" in red, stylized font. Overlaid on the bottom chart is also the text "Different device!" in red, stylized font.

Same device!
Different stencil!





Workgroup Size + Stencils

1. Non-linear, non-continuous
2. Device, program, dataset
3. Not all values are legal

Autotuning

**Set a workgroup size
Execute and time program**

Set a workgroup size

Execute and time program

... (continue until done / bored)

Pick the best one you tried

Set a workgroup size

Execute and time program

... (continue until done / bored)

Pick the best one you tried

*(íterative
compilación)*

BAD!

Takes a loooong time

BAD!

Takes a loooong time

BAD!

Must be repeated for every new “x”



Let's improve

Set a workgroup size

Execute and time program

... (continue until done / bored)

Pick the best one you tried

Set a workgroup size

Execute and time program

... (continue until done / bored)

Pick the best one you tried

1 data point

Collect **data points**

Extract “features”

Train machine learning classifier

Extract “features”

Input to classifier

GOOD!

Can make *predictions* on unseen “X”

device
dataset
program

GOOD!

Can make *predictions* on unseen “X”



device
dataset
program

GOOD!

Many unanswered questions ...

Questions:

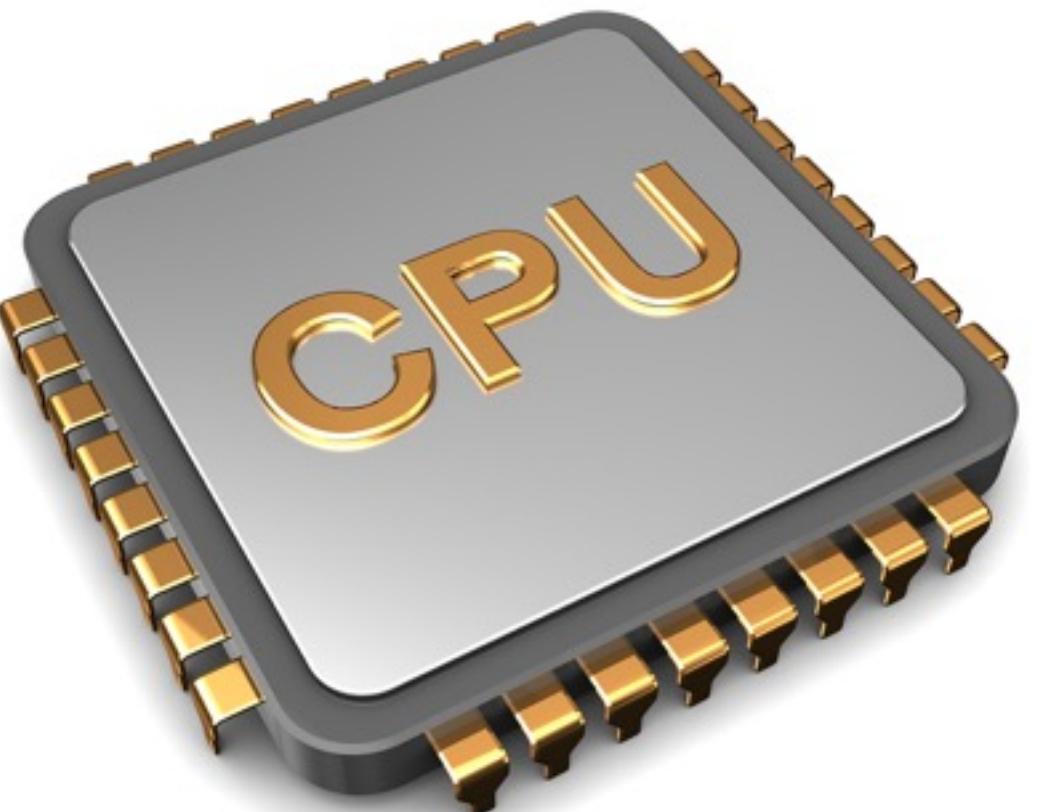
- 1. What features do we need?**
- 2. What programs do we train on?**
- 3. How do we make predictions?**

Questions:

1. What features do we need?
2. What programs do we train on?
3. How do we make predictions?

- 1. Device**
- 2. Kernel**
- 3. Dataset**

- 1. Device**
- 2. Kernel**
- 3. Dataset**



or



How many compute units?

How much memory?

Cache size?

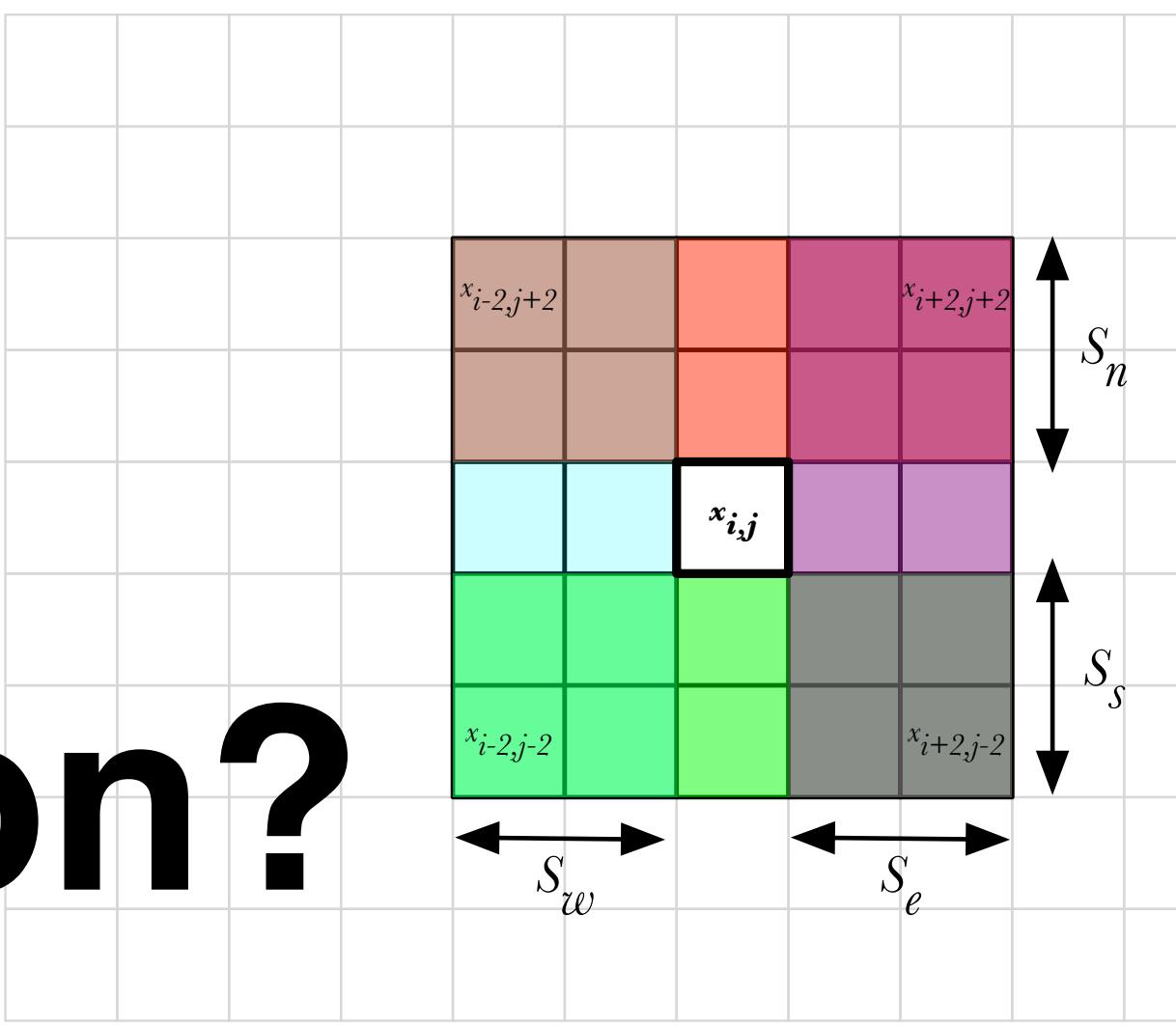
etc.



- 1. Device**
- 2. Kernel**
- 3. Dataset**

- 1. Device**
- 2. Kernel**
- 3. Dataset**

1. Device
2. Kernel
3. Dataset



How big is border region?

What shape is it?

How many instructions?

What type of instructions?

etc.



1. Device
2. Kernel
3. Dataset

- 1. Device**
- 2. Kernel**
- 3. Dataset**

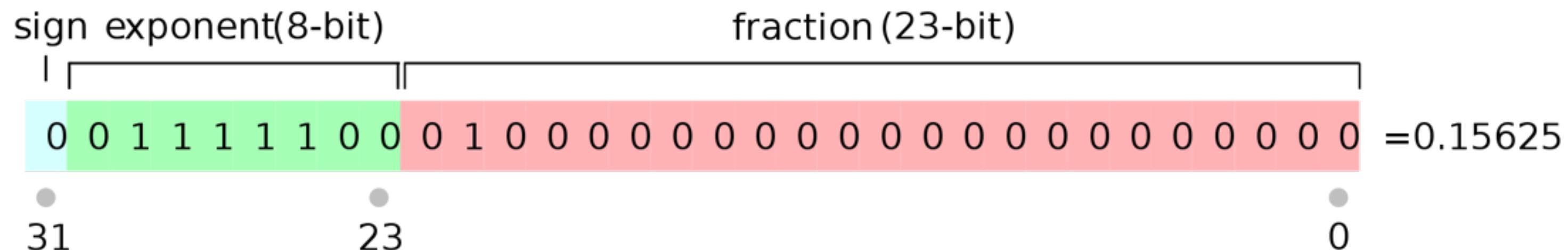
1. Device
2. Kernel
3. Dataset

$$\mathbf{A} = \begin{bmatrix} A_{11} & A_{12} & \cdots & A_{1n} \\ A_{21} & & & A_{2n} \\ \vdots & & & \vdots \\ A_{n1} & A_{n2} & \cdots & A_{nn} \end{bmatrix}$$

How big is the data?

What type is the input?

What type is the output?



- 1. Device**
- 2. Kernel**
- 3. Dataset**

- 1. Device**
- 2. Kernel**
- 3. Dataset**

Questions:

1. What features do we need?
2. What programs do we train on?
3. How do we make predictions?

Questions:

1. What features do we need? ✓
2. What programs do we train on?
3. How do we make predictions?

1. Learn by example
2. Learn by exploration

Use benchmark programs

Hope that they are representative

1. Learn by example

2. Learn by exploration

- 1. Learn by example**
- 2. Learn by exploration**

1. Learn by example
 2. Learn by exploration
- Create own benchmarks
Explore (the huge!) program space

Questions:

1. What features do we need? ✓
2. What programs do we train on?
3. How do we make predictions?

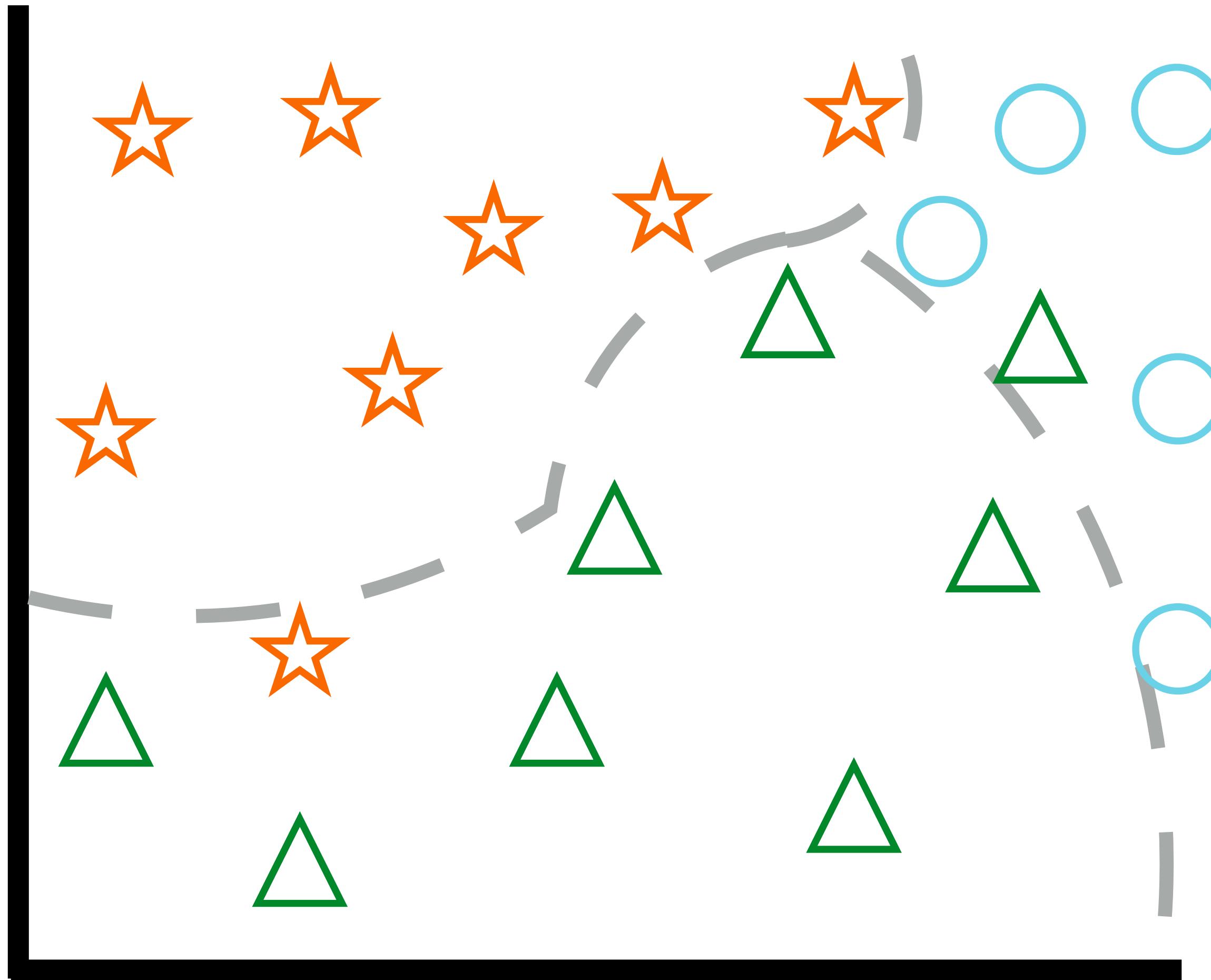
Questions:

1. What features do we need? ✓
2. What programs do we train on? ✓
3. How do we make predictions?

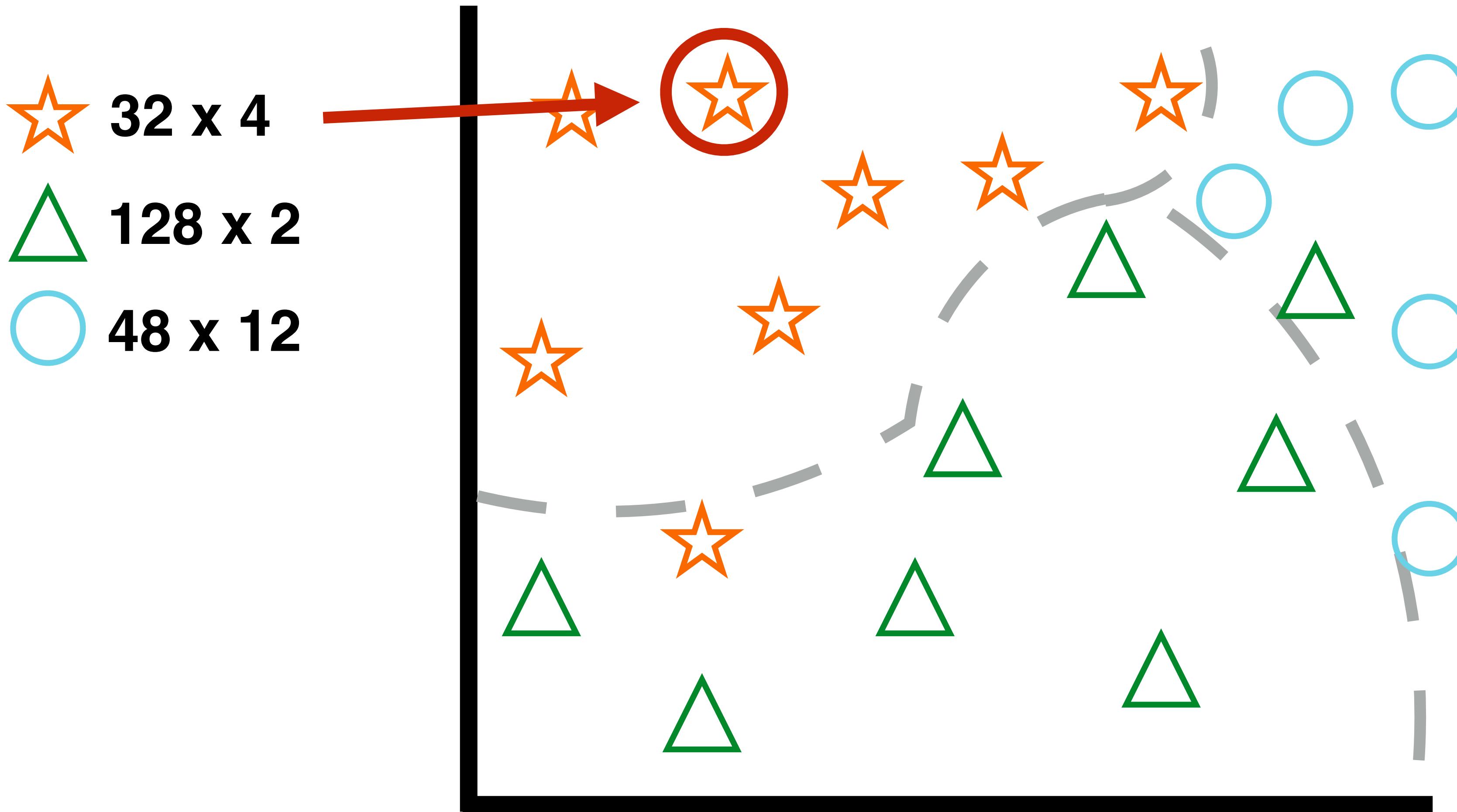
- 1. Classifier**
- 2. Runtime Regressor**
- 3. Speedup Regressor**

- 1. Classifier**
- 2. Runtime Regressor**
- 3. Speedup Regressor**

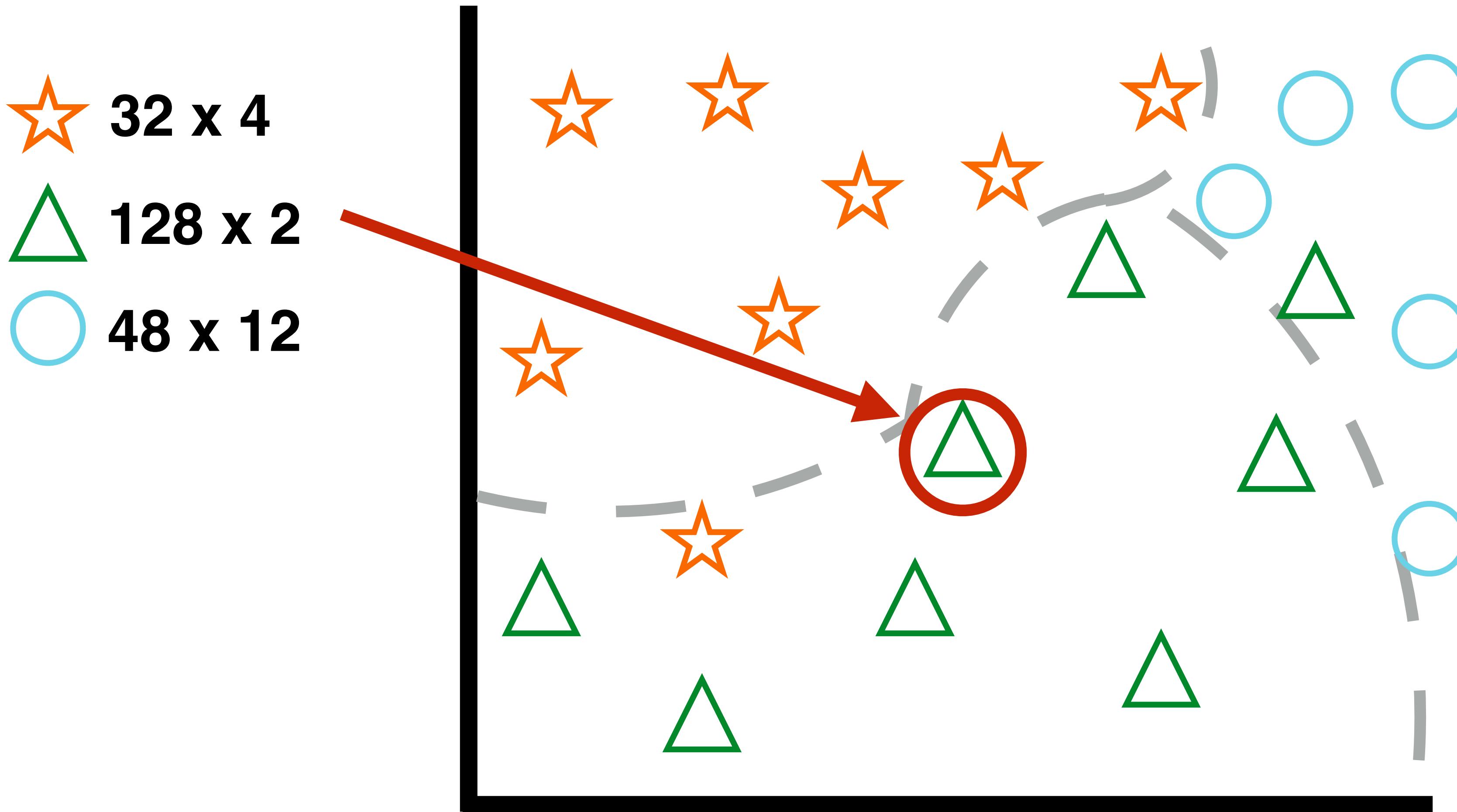
- ★ 32 x 4
- △ 128 x 2
- 48 x 12



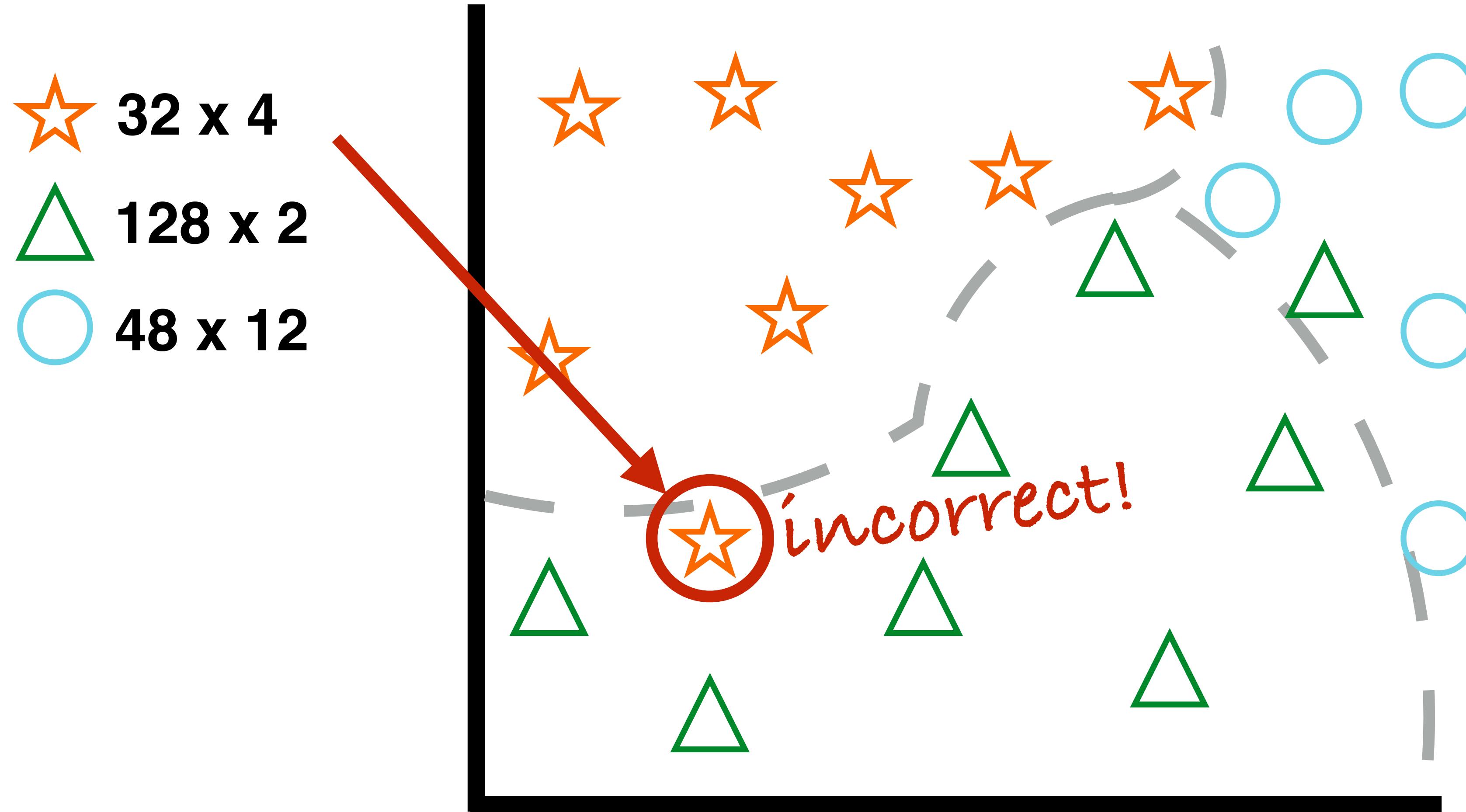
Predict category (*optimal workgroup size*) for scenario



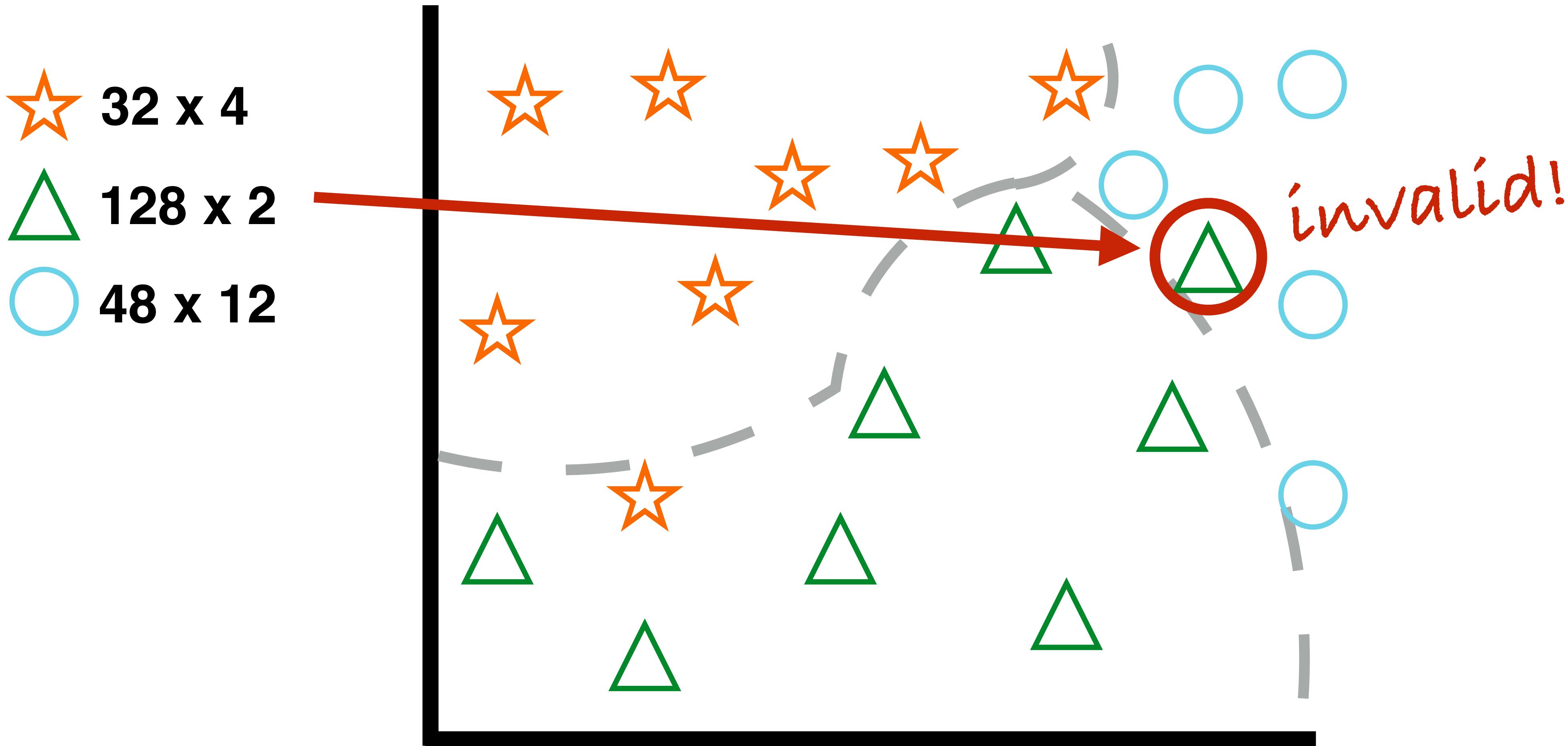
Predict category (*optimal workgroup size*) for scenario



Predict category (*optimal workgroup size*) for scenario



Predict category (*optimal workgroup size*) for scenario



Predict category (*optimal workgroup size*) for scenario

**What do we do if the
classifier predicts an
illegal parameter value?**

Fallback Handlers

- 1. Baseline**
- 2. Random**
- 3. Nearest Neighbour**

Fallback Handlers

1. Baseline
2. Random
3. Nearest Neighbour

“pick something we
know is safe”

Fallback Handlers

1. Baseline
2. Random
3. Nearest Neighbour

“pick a random
value”

Fallback Handlers

1. Baseline
2. Random
3. Nearest Neighbour

“pick the closest value we think will work”

- 1. Classifier**
- 2. Runtime Regressor**
- 3. Speedup Regressor**

- 1. Classifier**
- 2. Runtime Regressor**
- 3. Speedup Regressor**

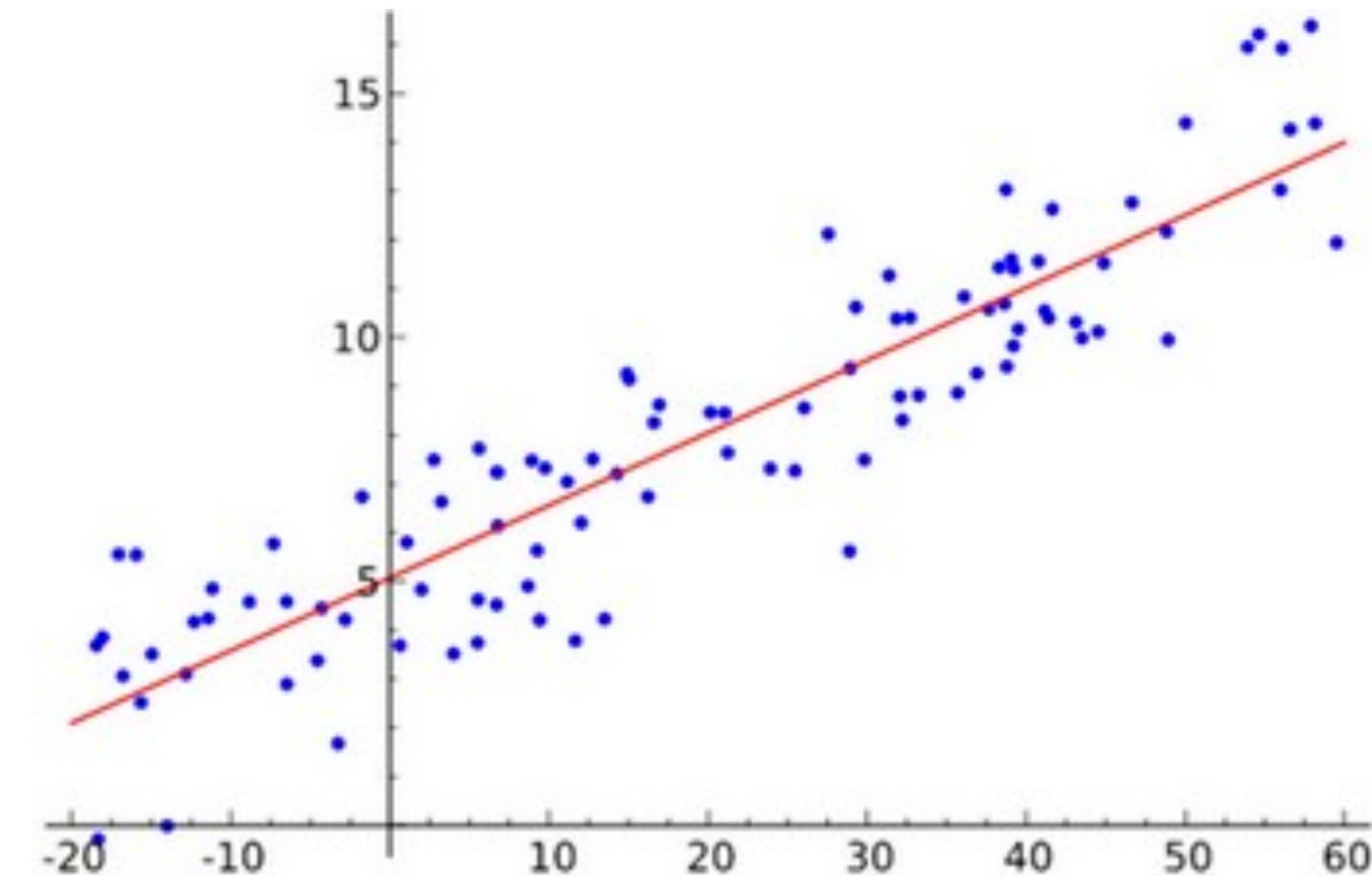
1. Classifier

2. Runtime Regressor

3. Speedup Regressor

**Predict *runtime* of
program for
workgroup size**

**Search for *lowest*
*runtime***



1. Classifier

2. Runtime Regressor

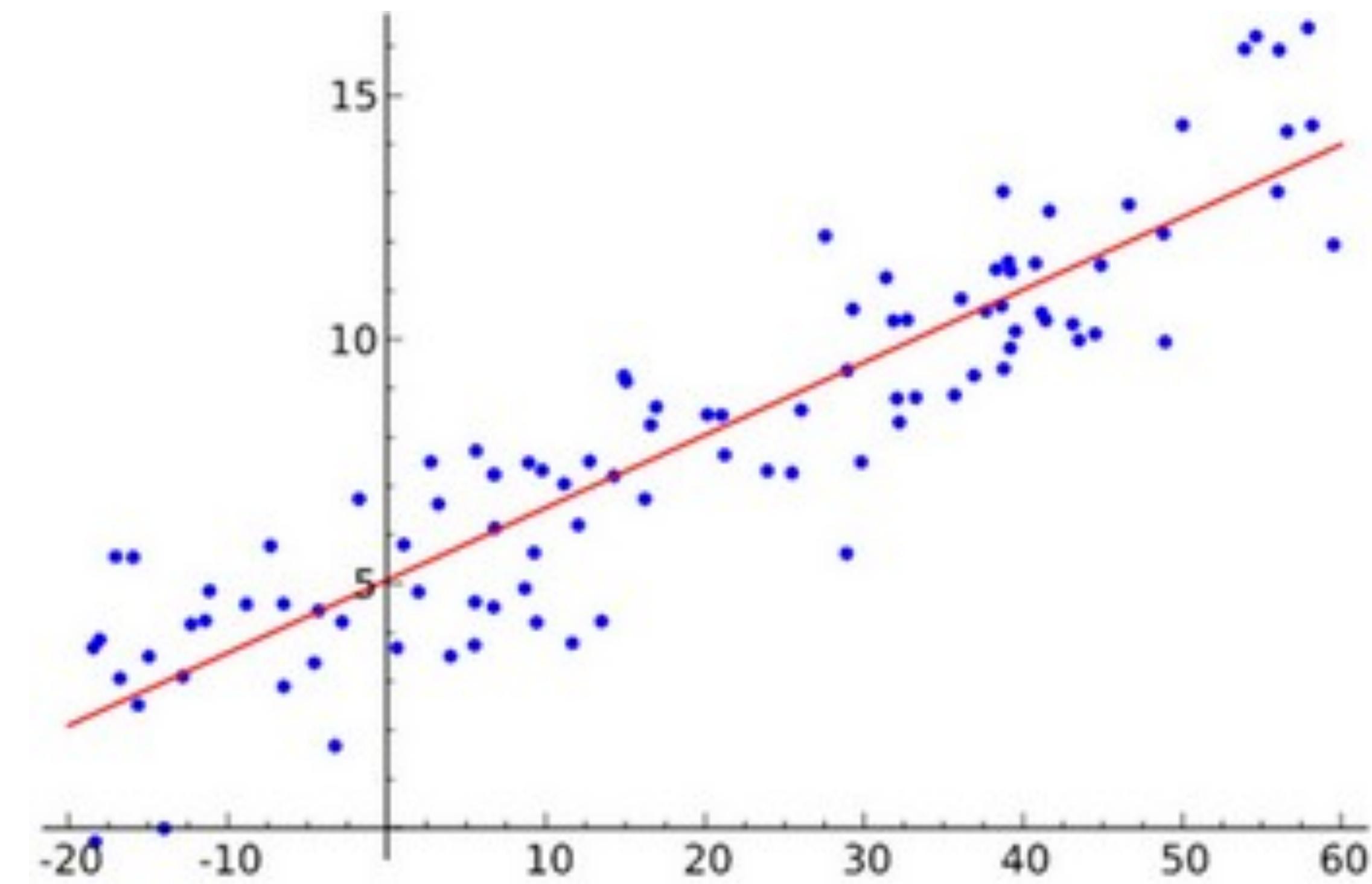
3. Speedup Regressor

- 1. Classifier**
- 2. Runtime Regressor**
- 3. Speedup Regressor**

1. Classifier
2. Runtime Regressor
3. Speedup Regressor

**Predict speedup of
workgroup size A
over B for program**

**Search for *highest*
speedup**



1. Classifier
2. Runtime Regressor
3. Speedup Regressor

- 1. Classifier**
- 2. Runtime Regressor**
- 3. Speedup Regressor**

Questions:

1. What features do we need? ✓
2. What programs do we train on? ✓
3. How do we make predictions?

Questions:

1. What features do we need? ✓
2. What programs do we train on? ✓
3. How do we make predictions? ✓

Experiment

Implementation

Modified SkelCL stencil pattern

Python server process for autotuning

5 classifiers, random forest regressor

Experimental Setup

6 stencil benchmarks + synthetic.

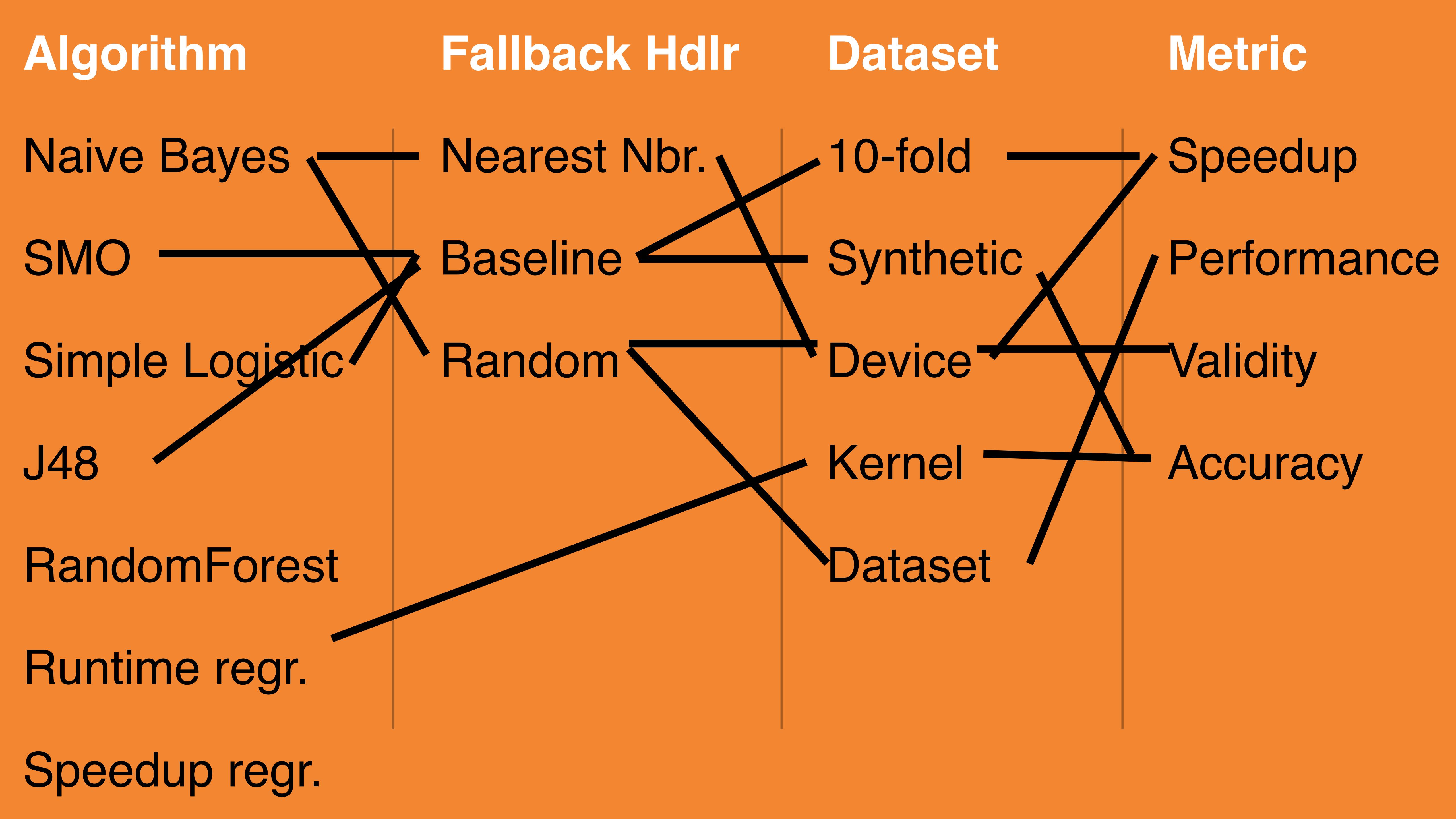
7 different GPUs & CPUs.

4 dataset sizes.

**Exhaustive search of workgroup size
space for each**

Results

Algorithm	Fallback Hdlr	Dataset	Metric
Naive Bayes	Nearest Nbr.	10-fold	Speedup
SMO	Baseline	Synthetic	Performance
Simple Logistic	Random	Device	Validity
J48		Kernel	Accuracy
RandomForest		Dataset	
Runtime regr.			
Speedup regr.			

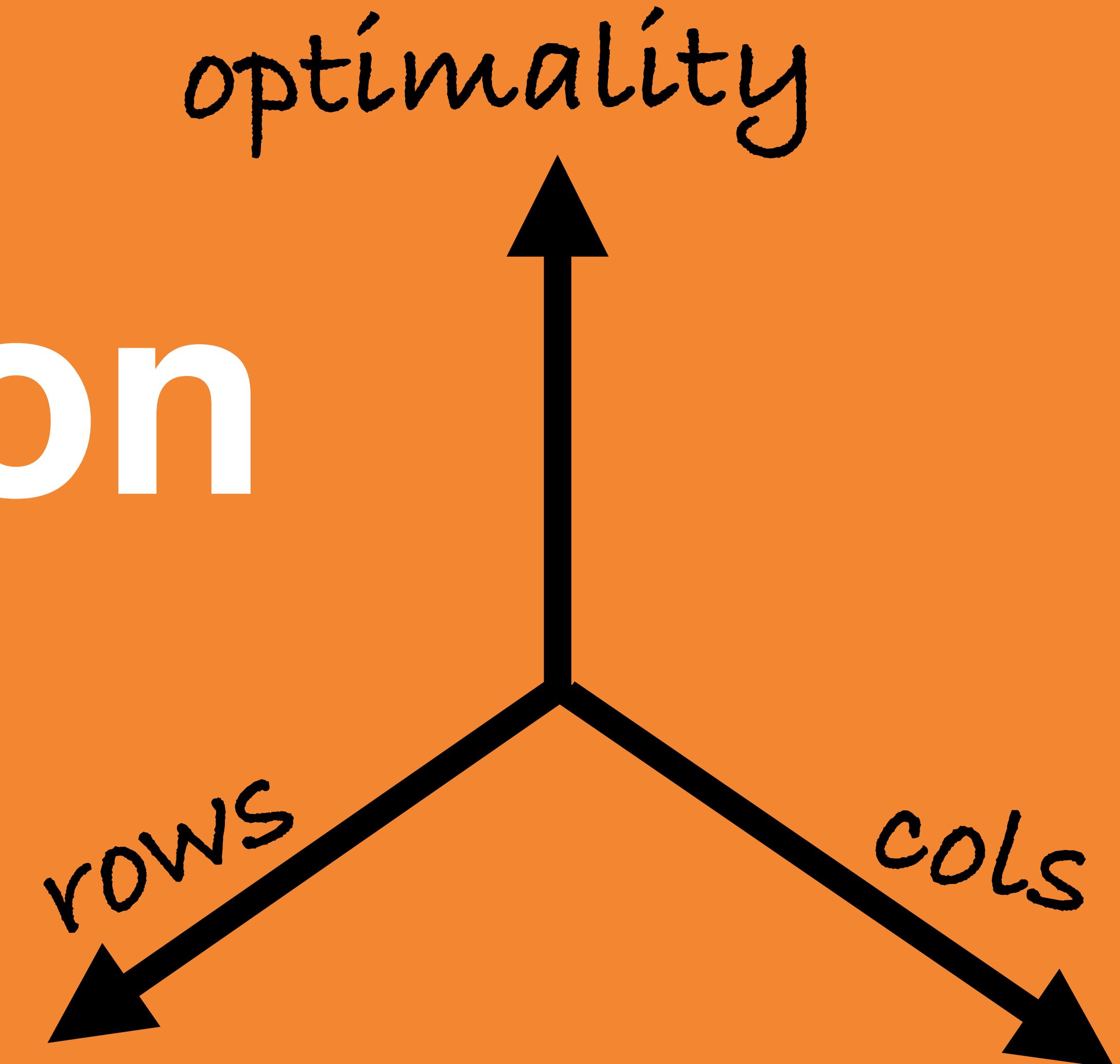


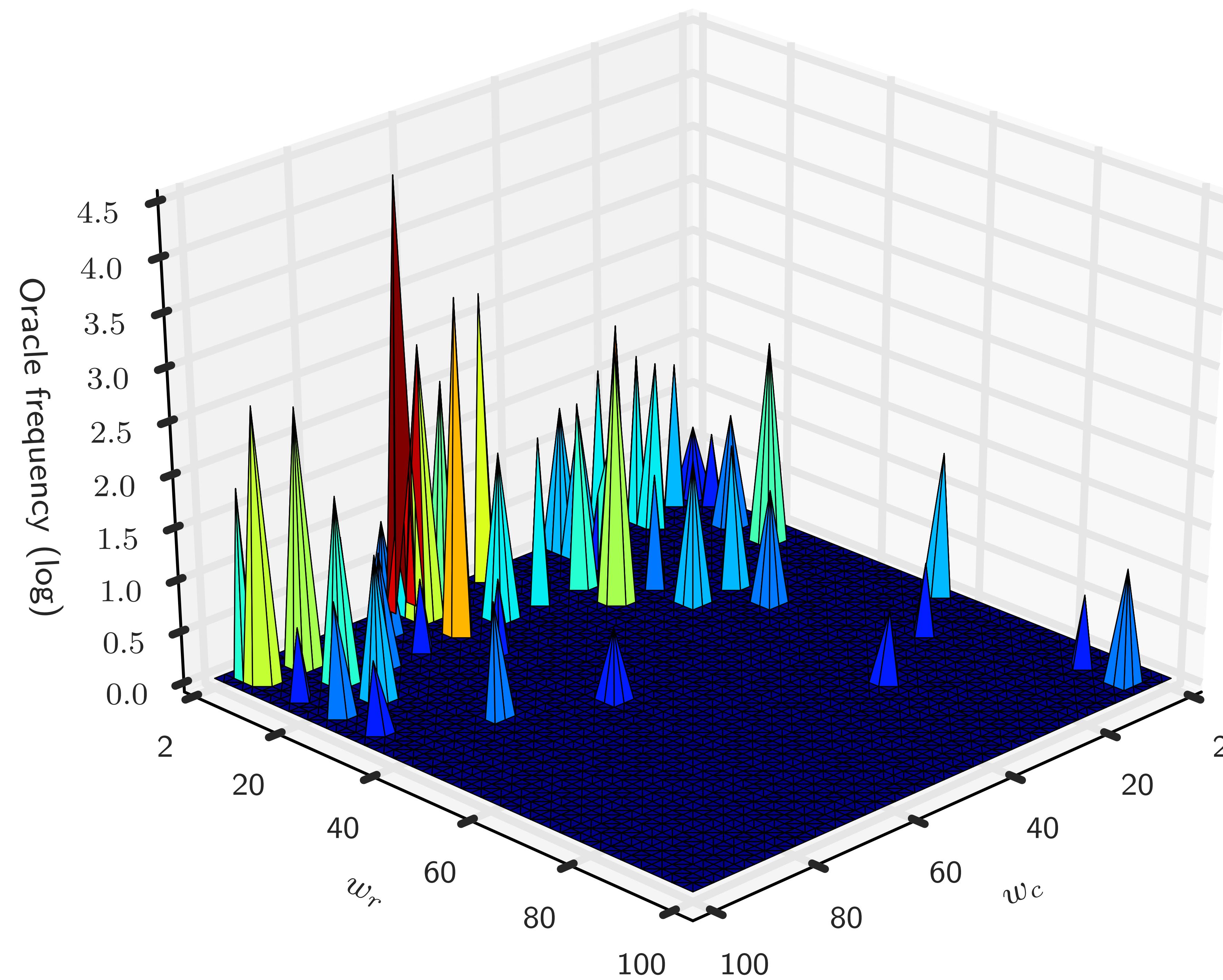
(selected)

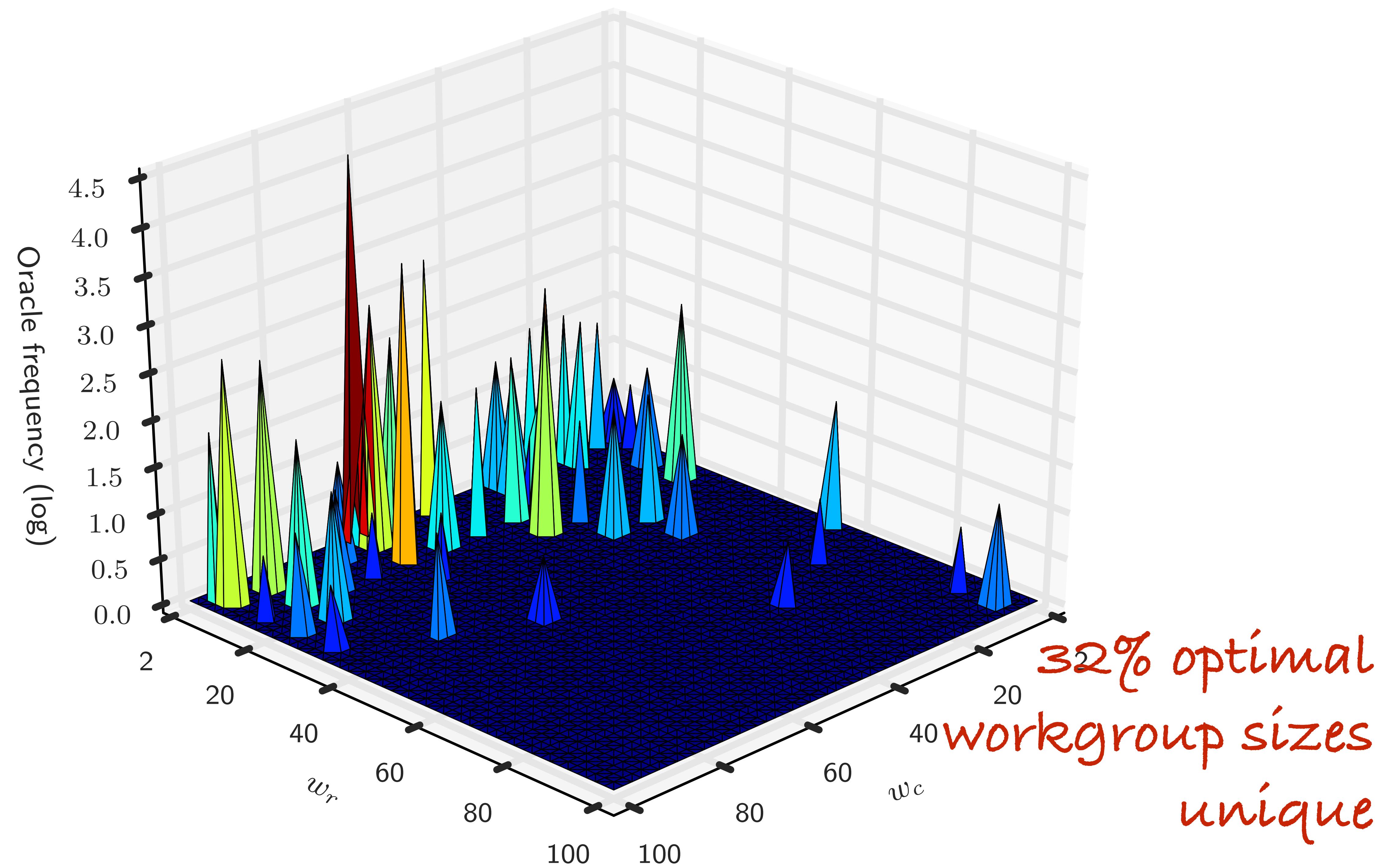
Results

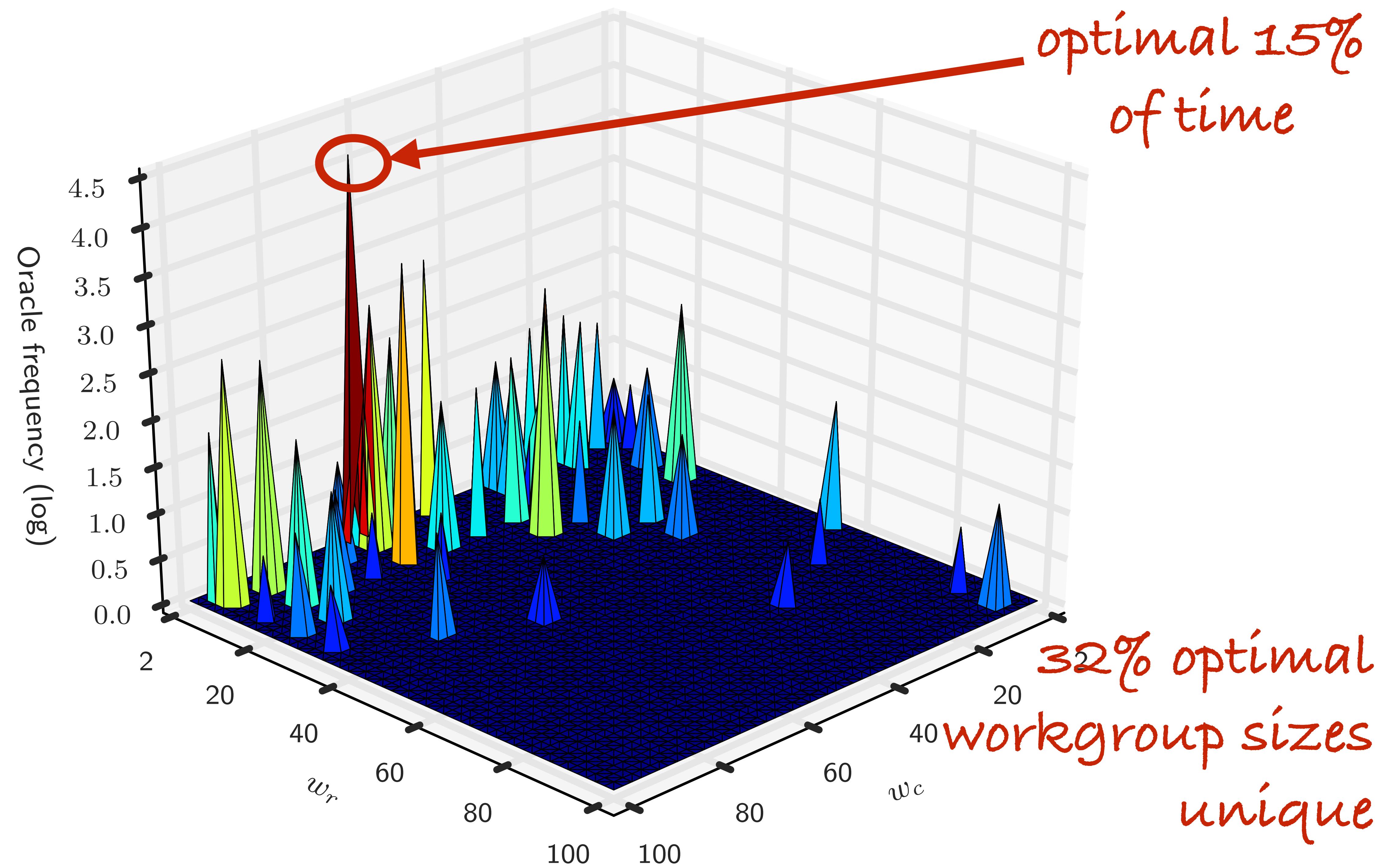
Optimisation

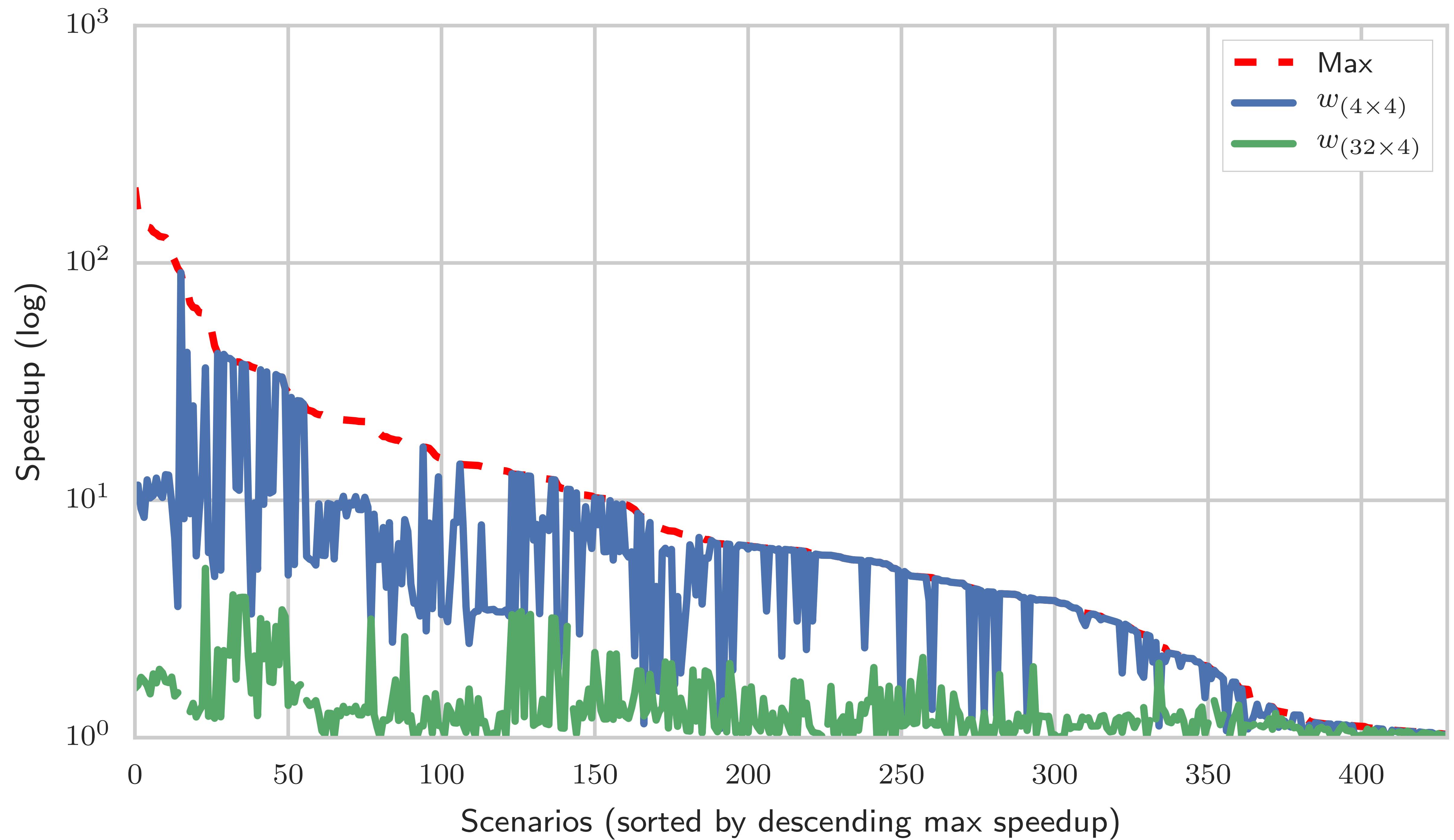
space

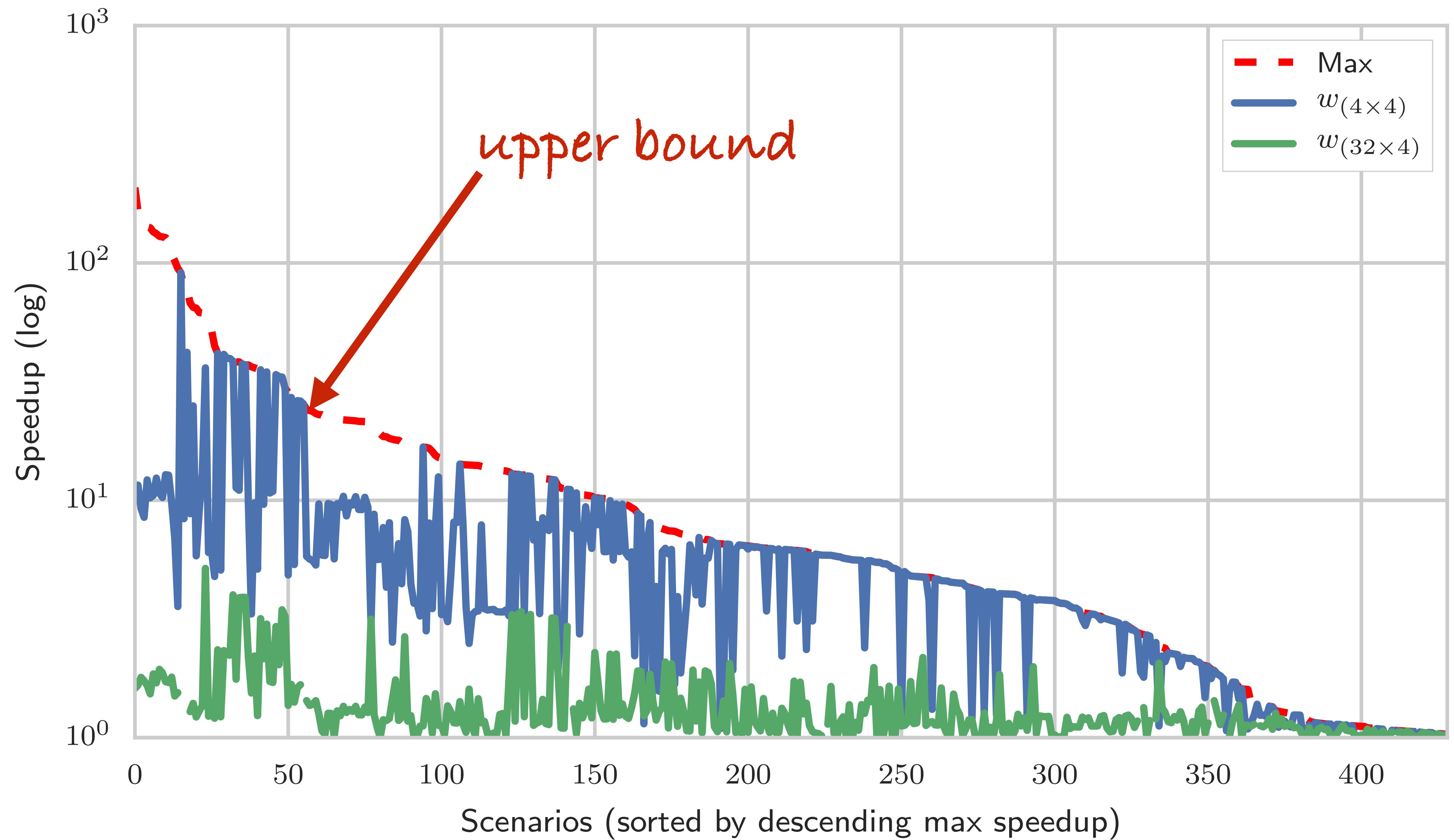


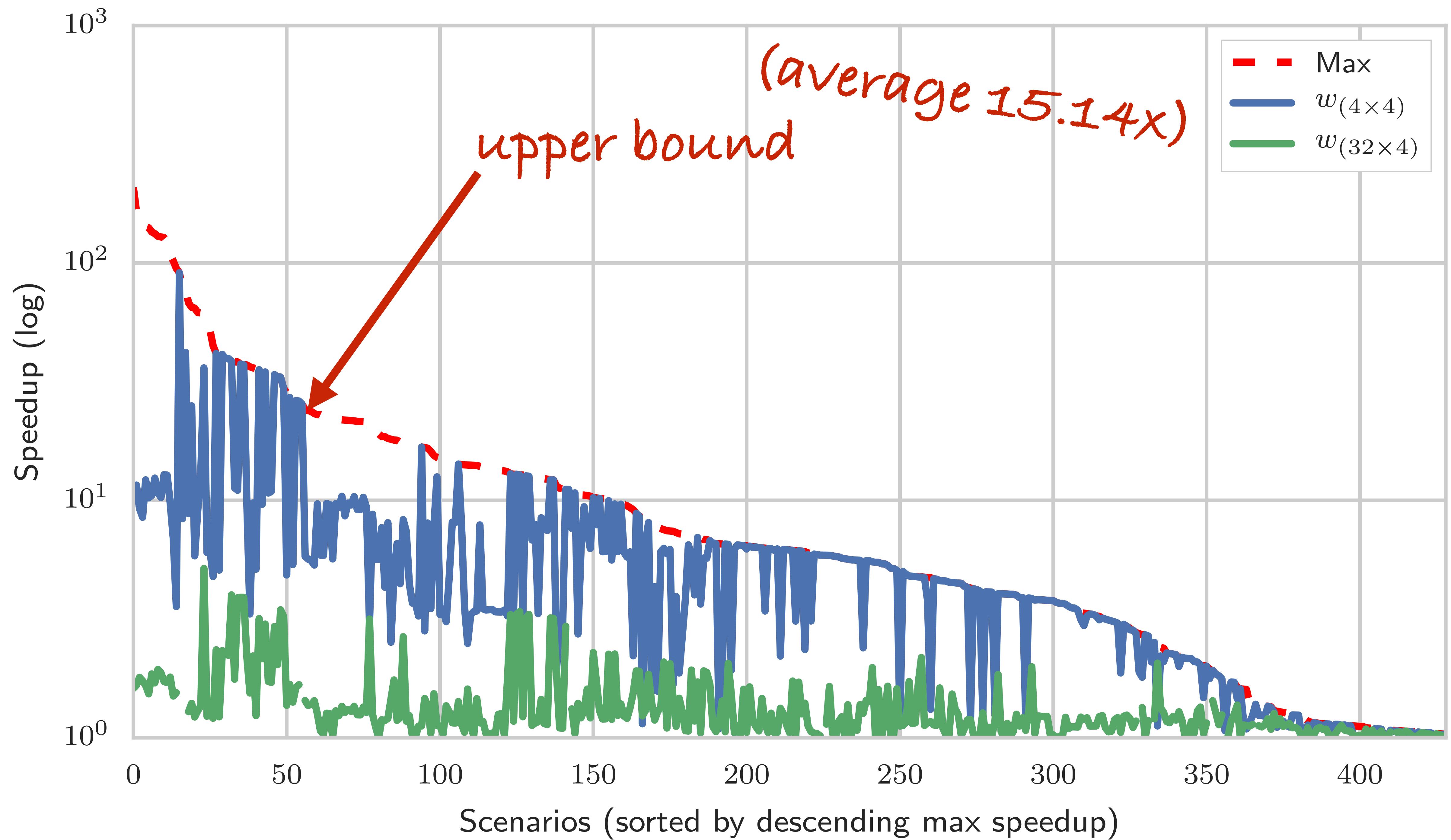


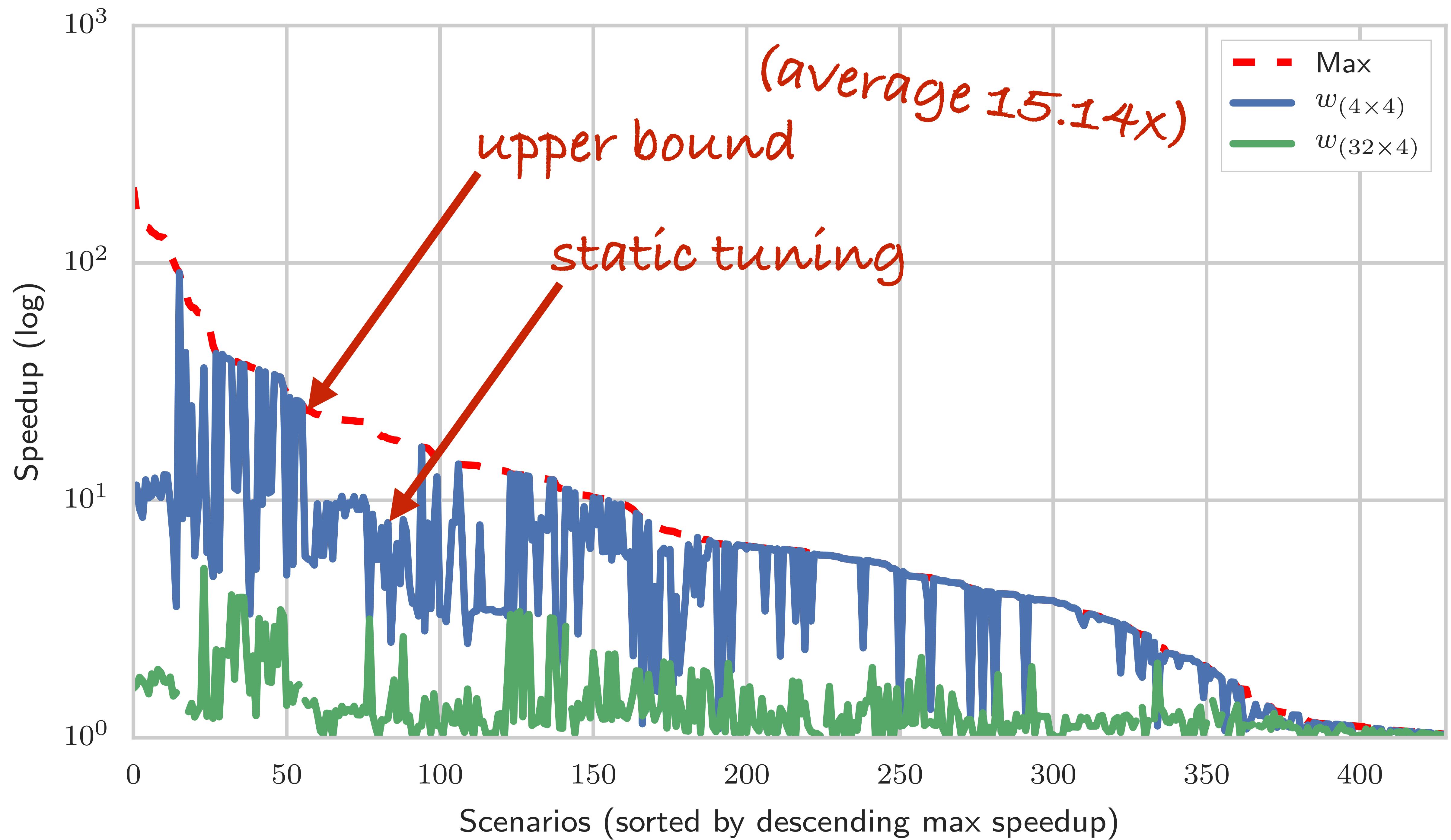


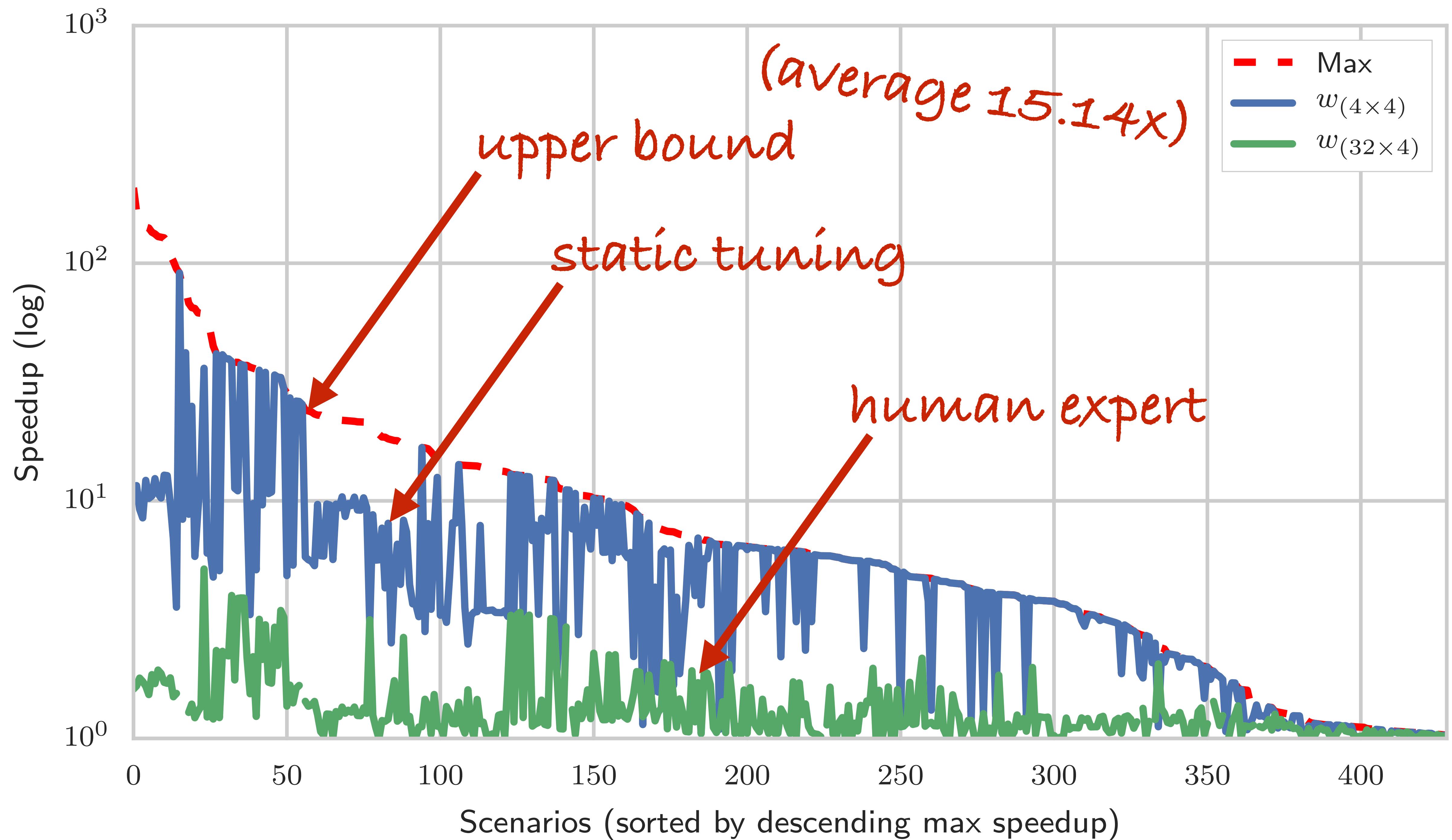




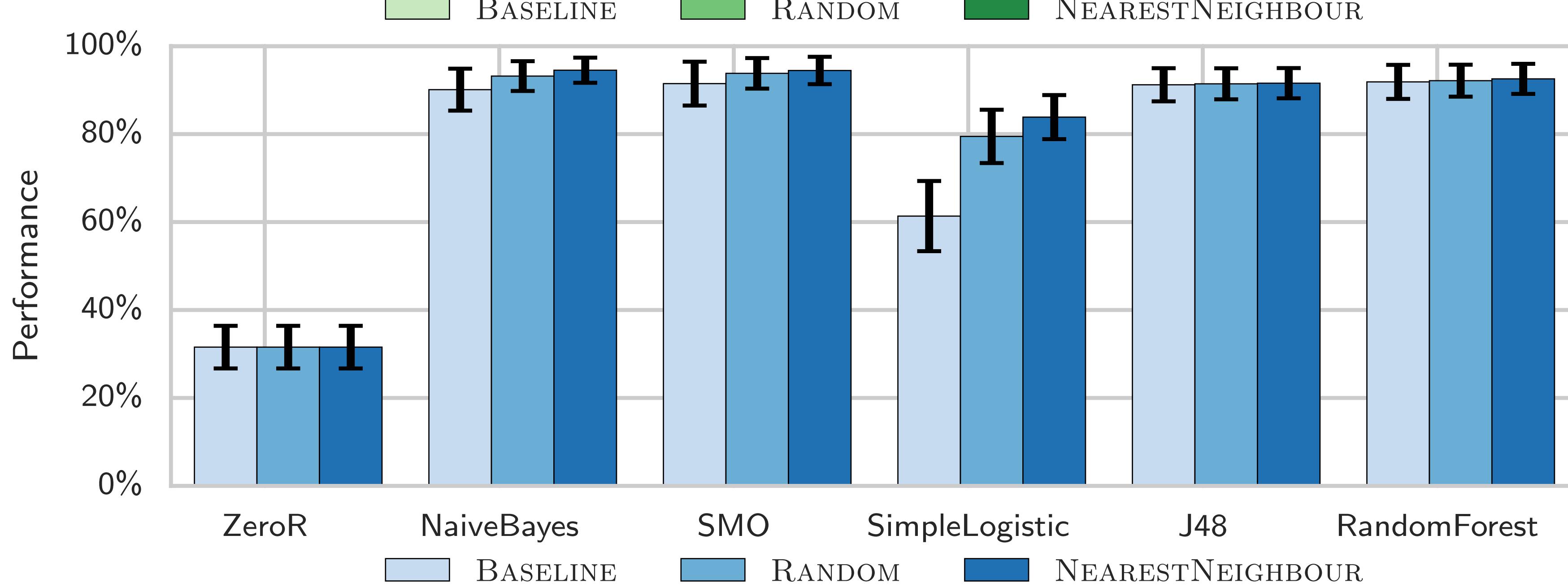
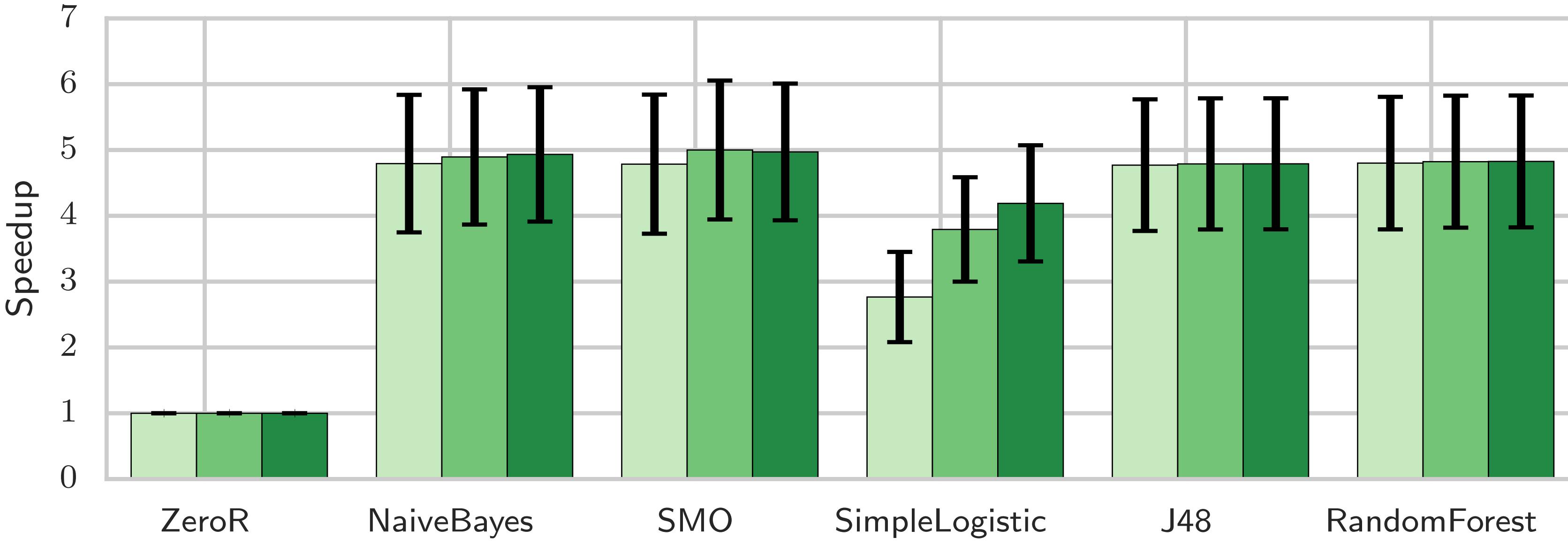


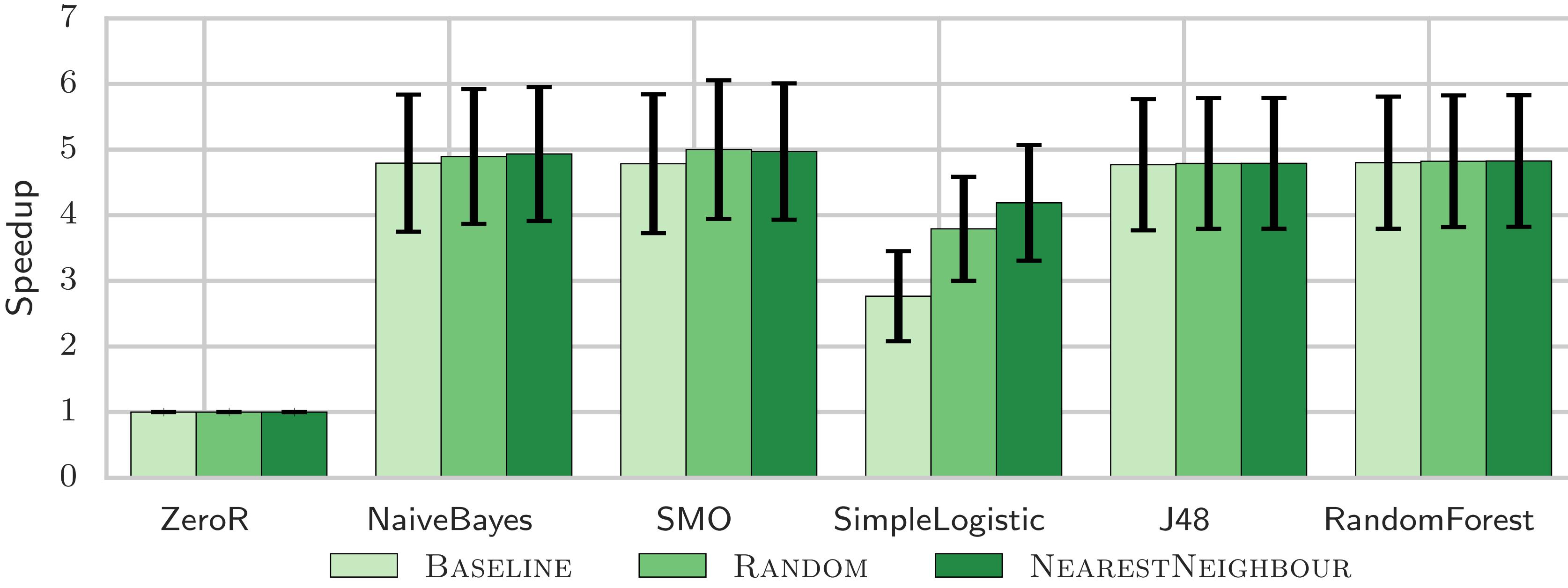




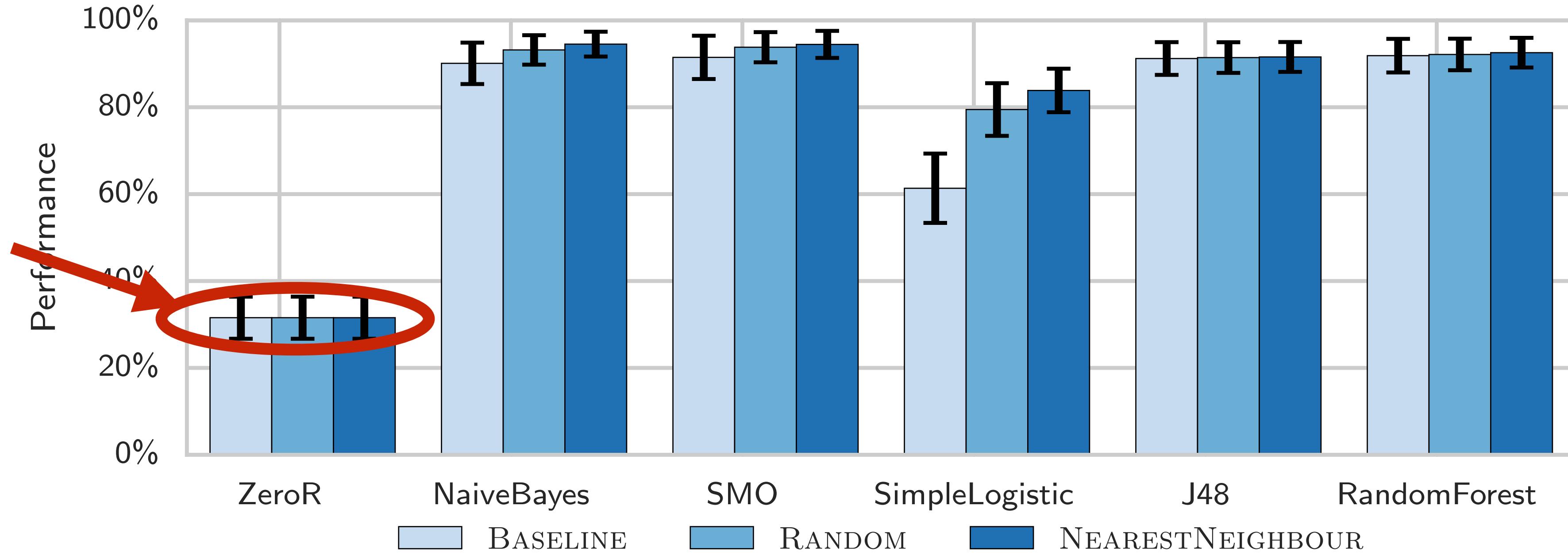


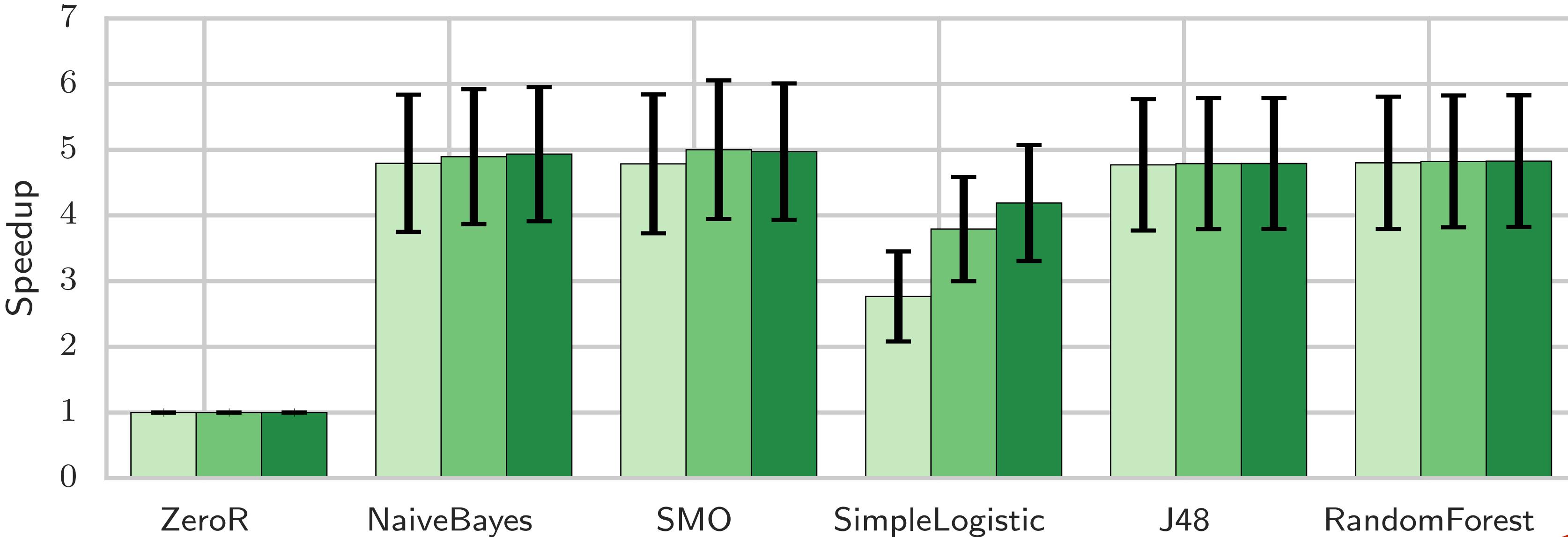
Autotuning Classification



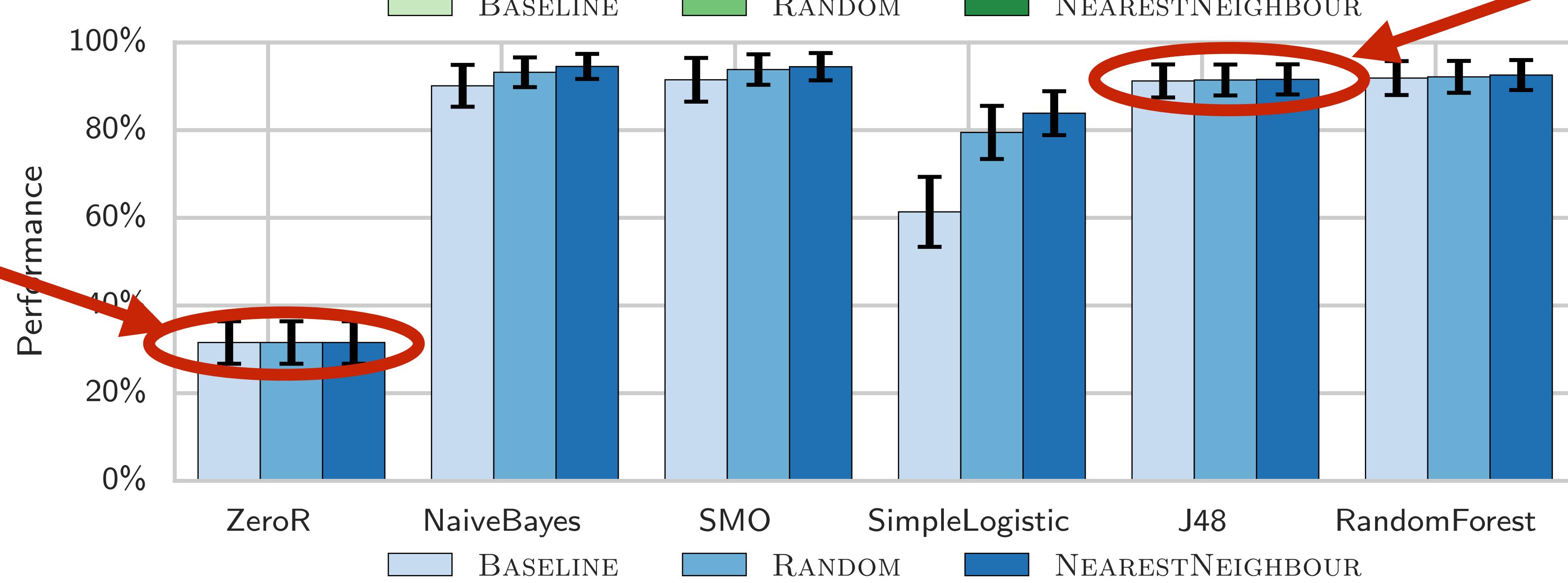


*26%
optimal*



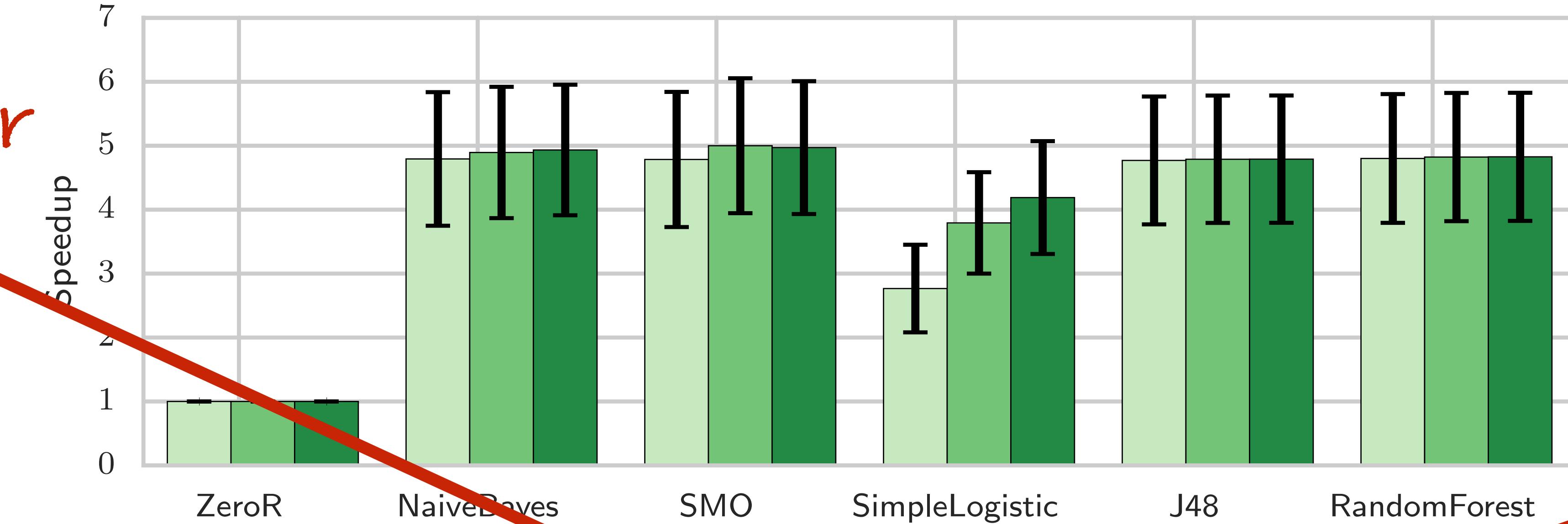


90% optimal



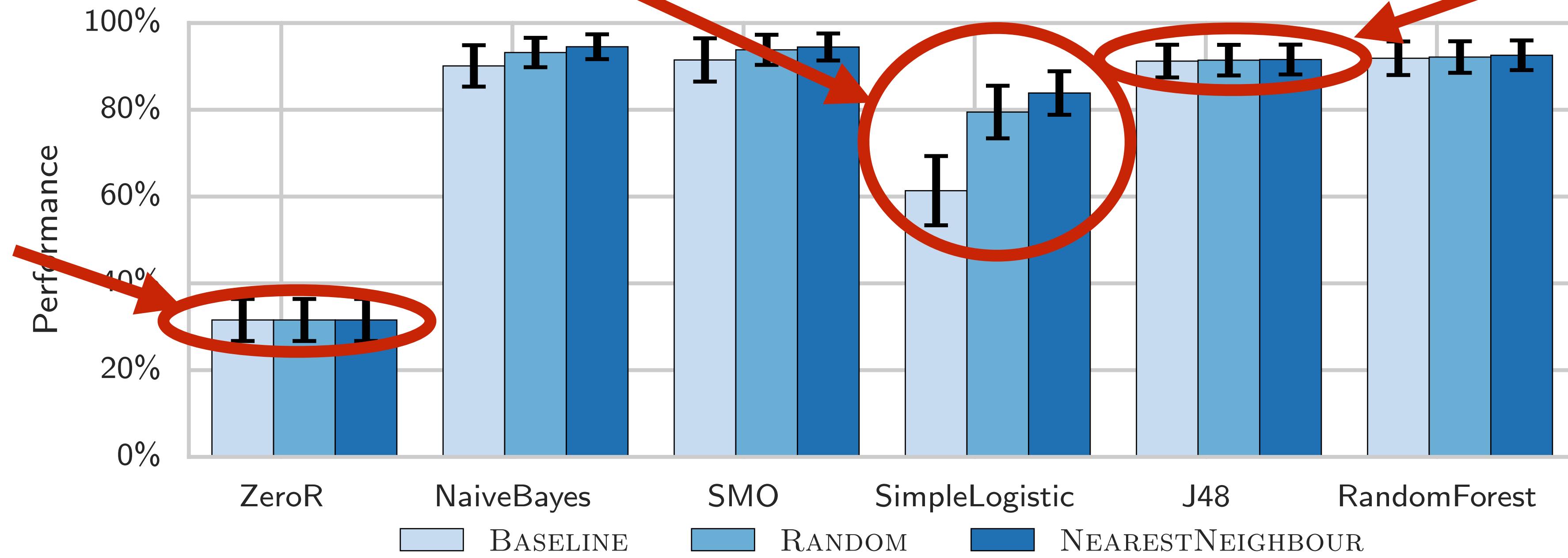
26% optimal

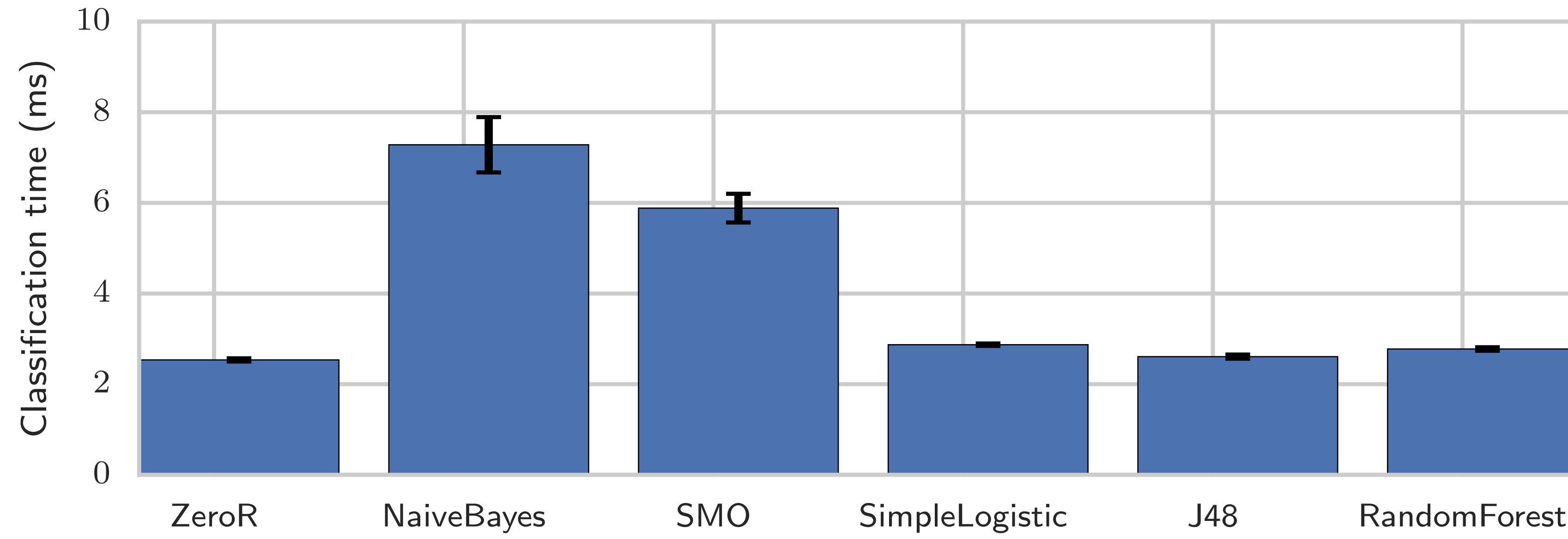
Nearest
neighbour
best



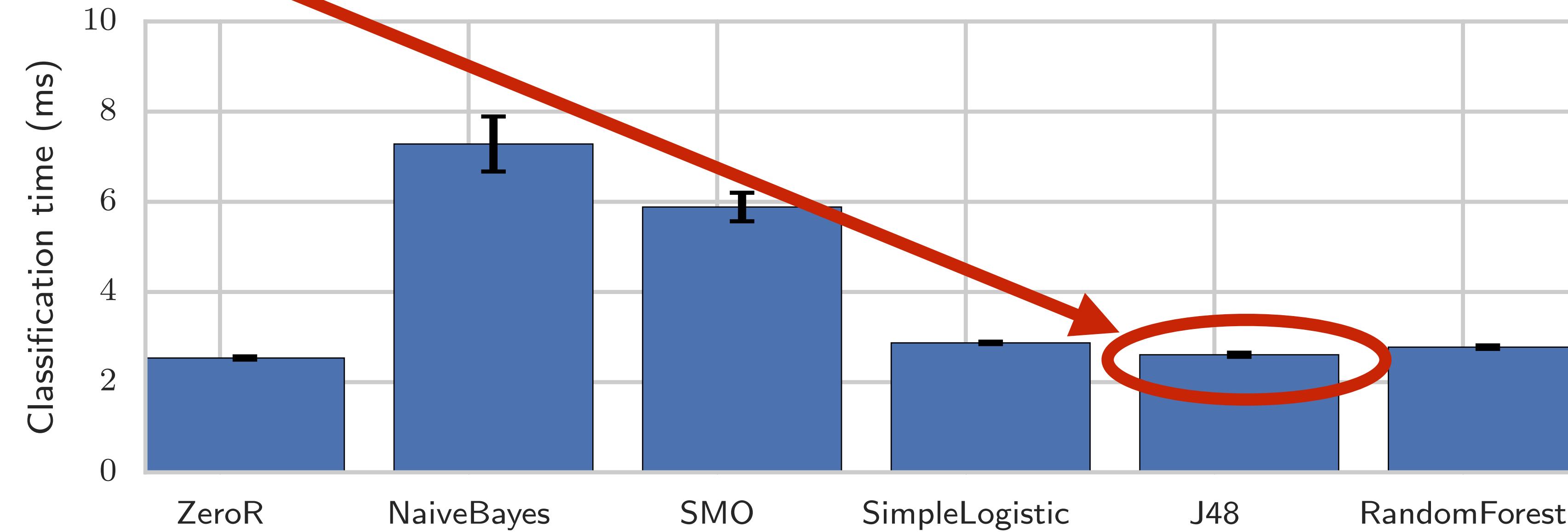
90%
optimal

26%
optimal



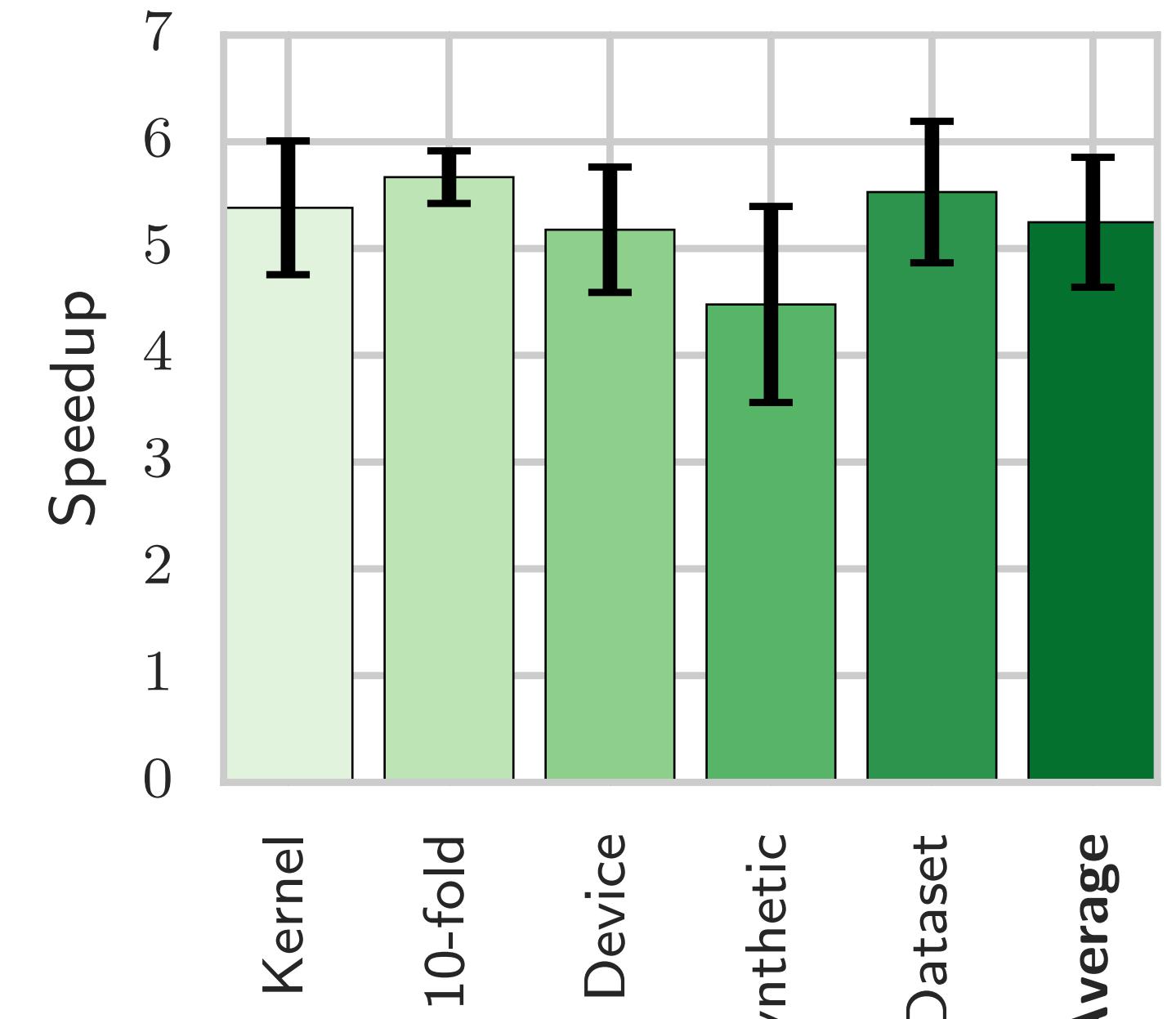
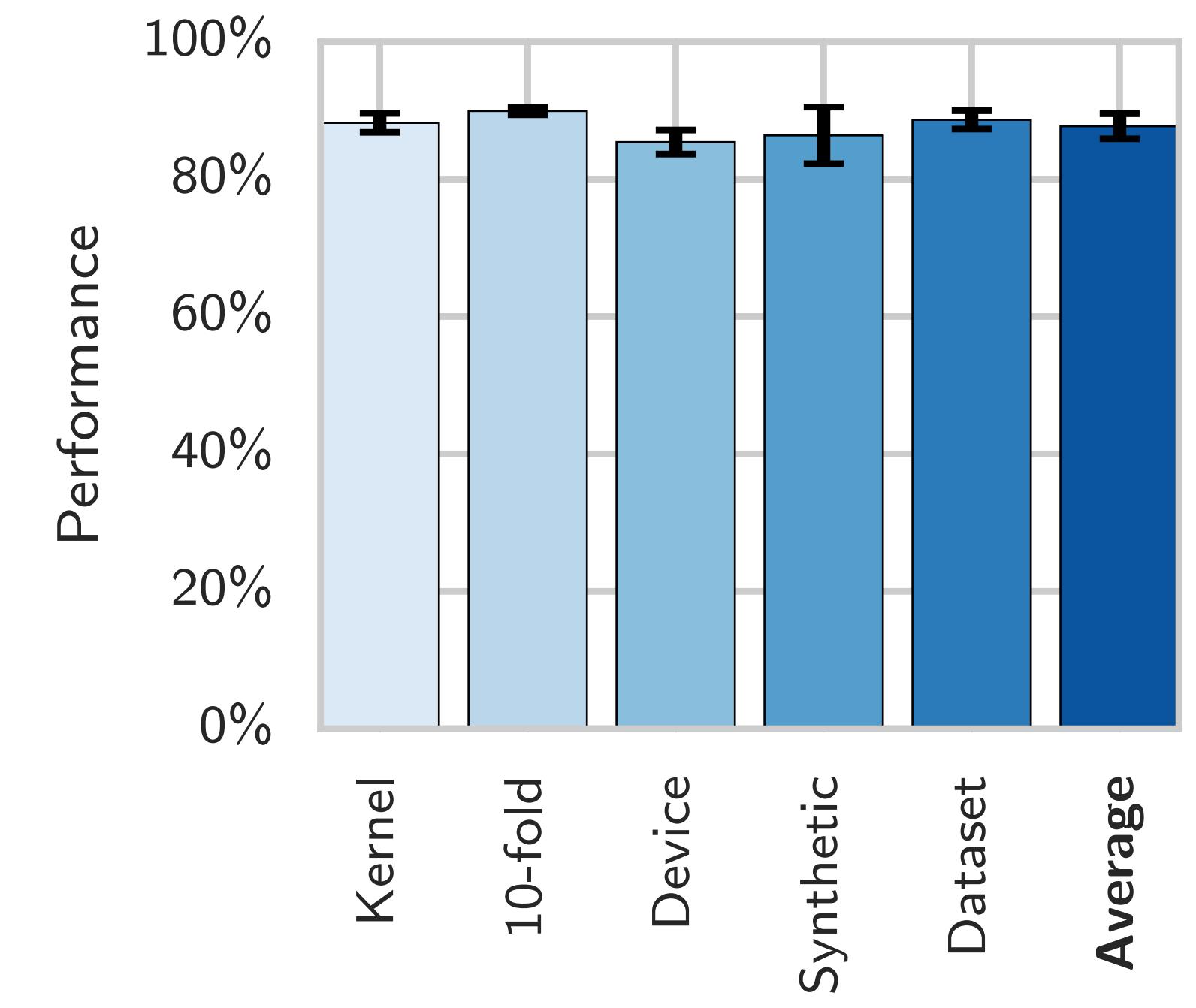
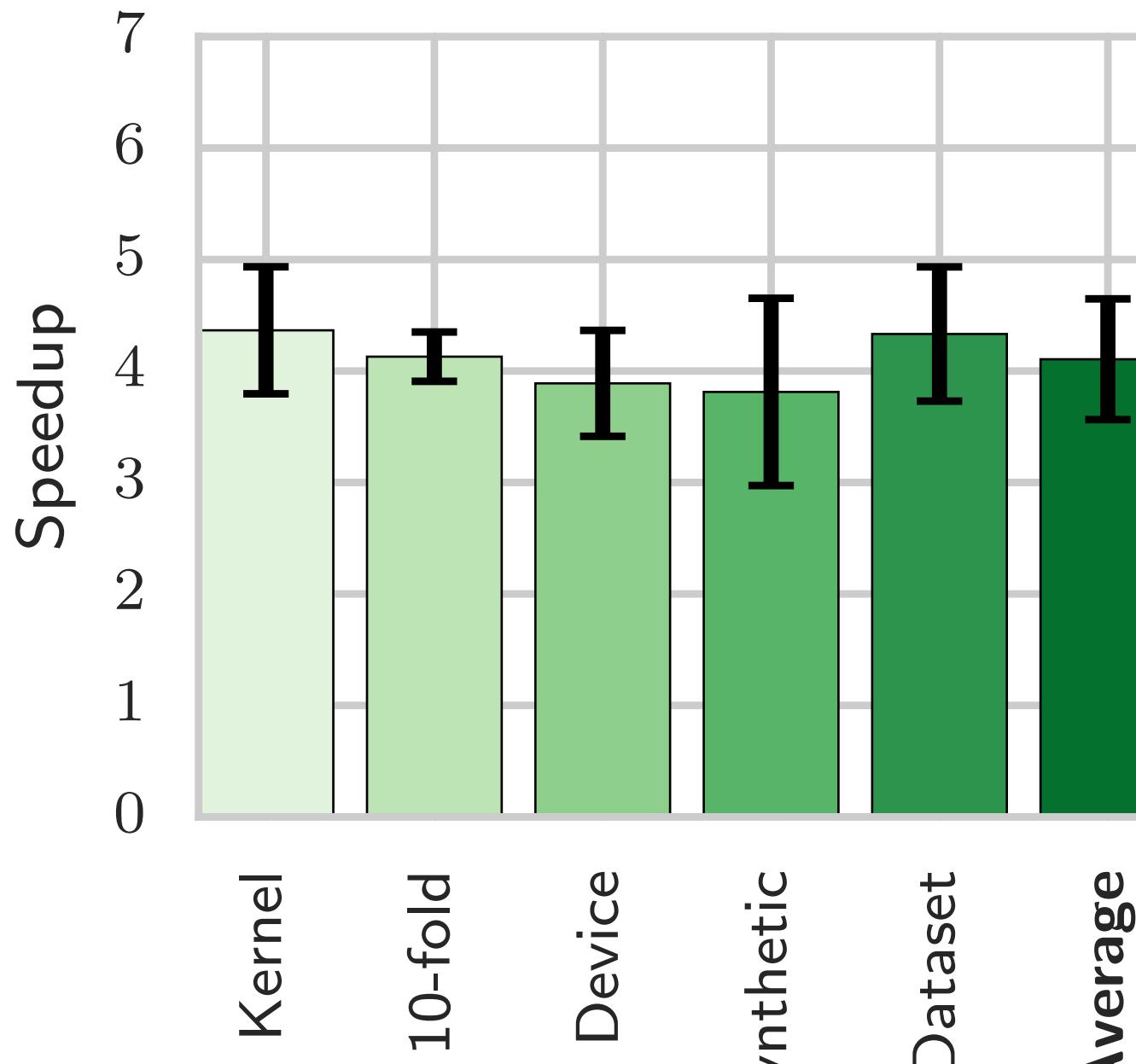
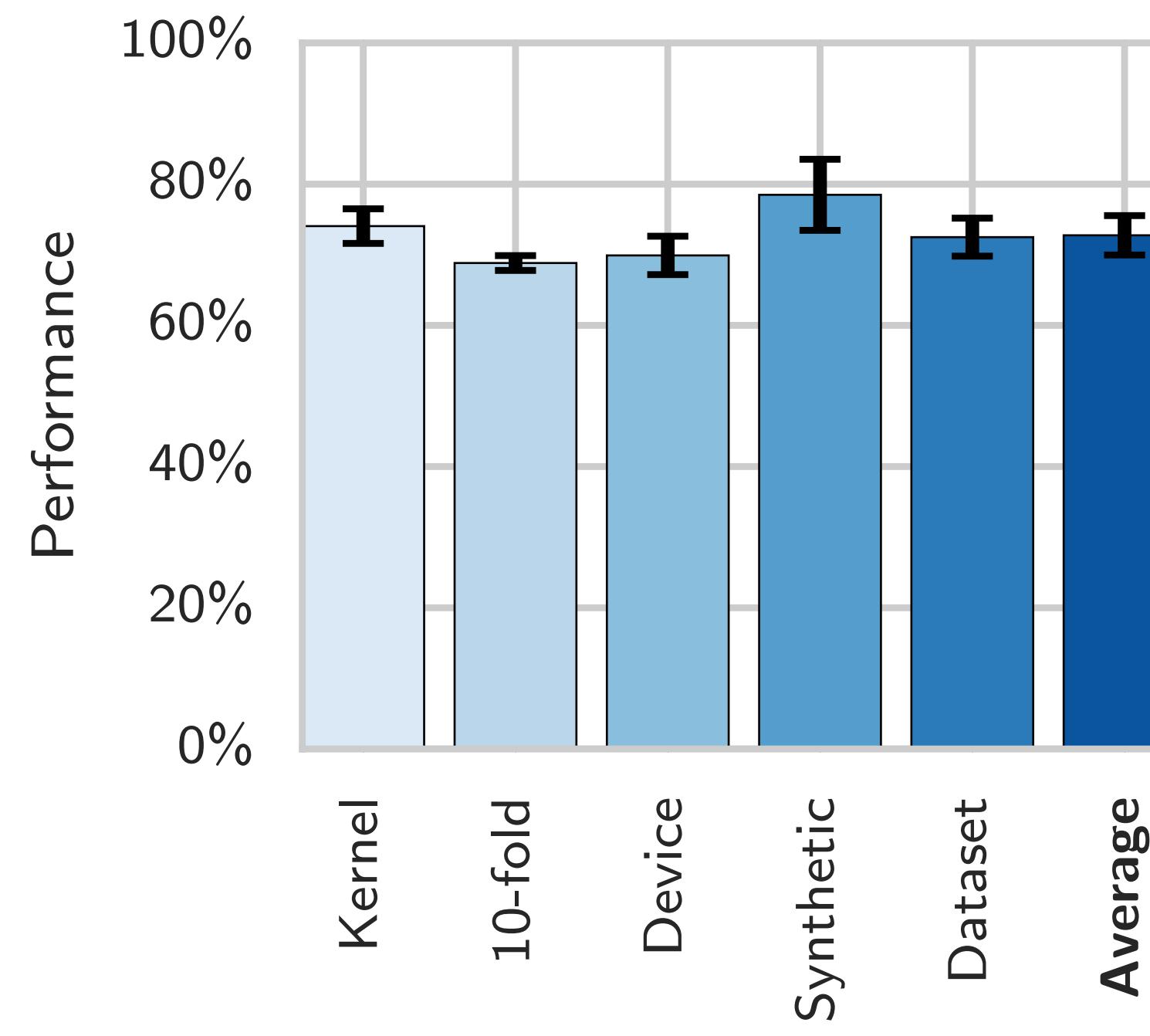


2.5ms RTT



Autotuning Regression

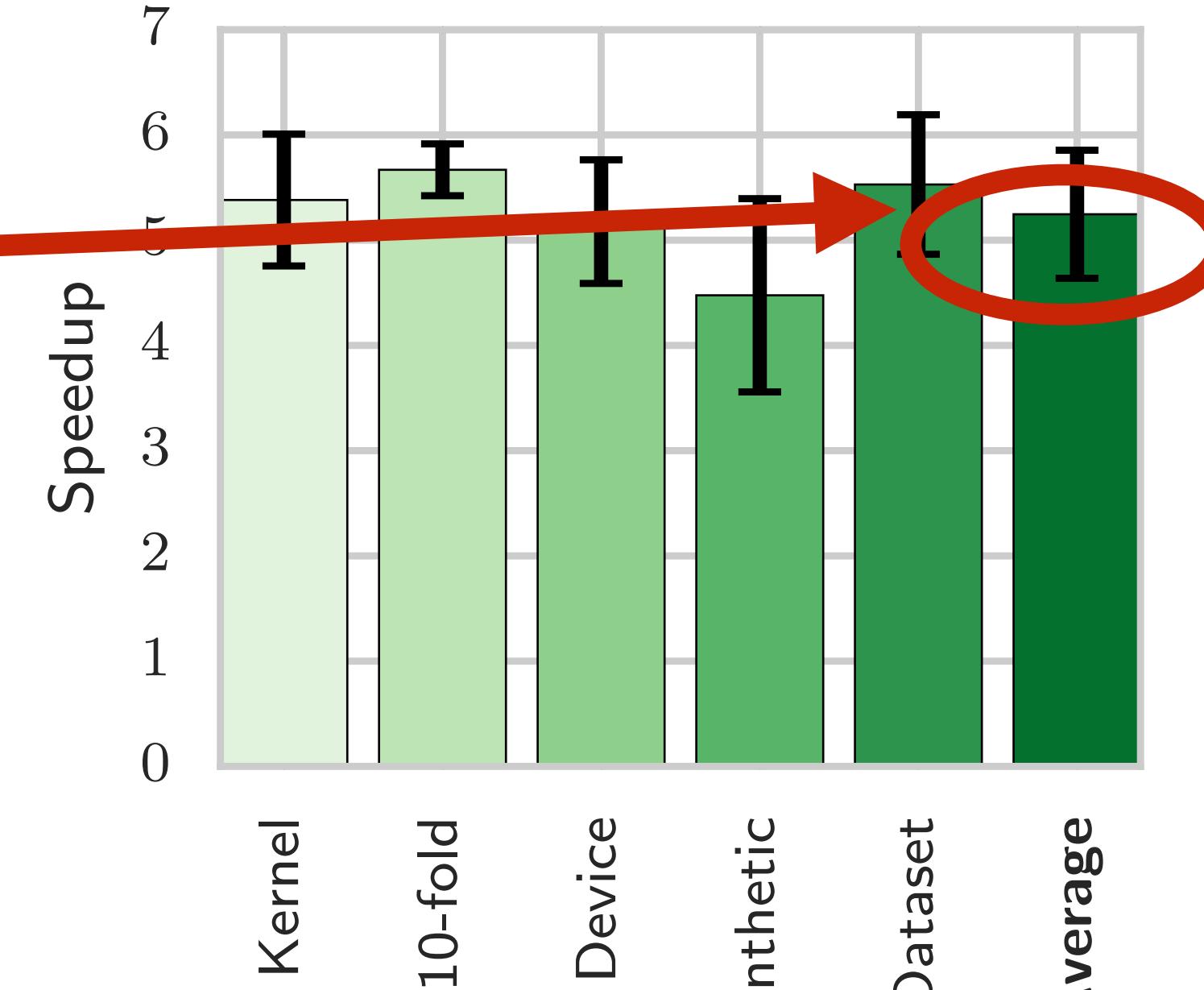
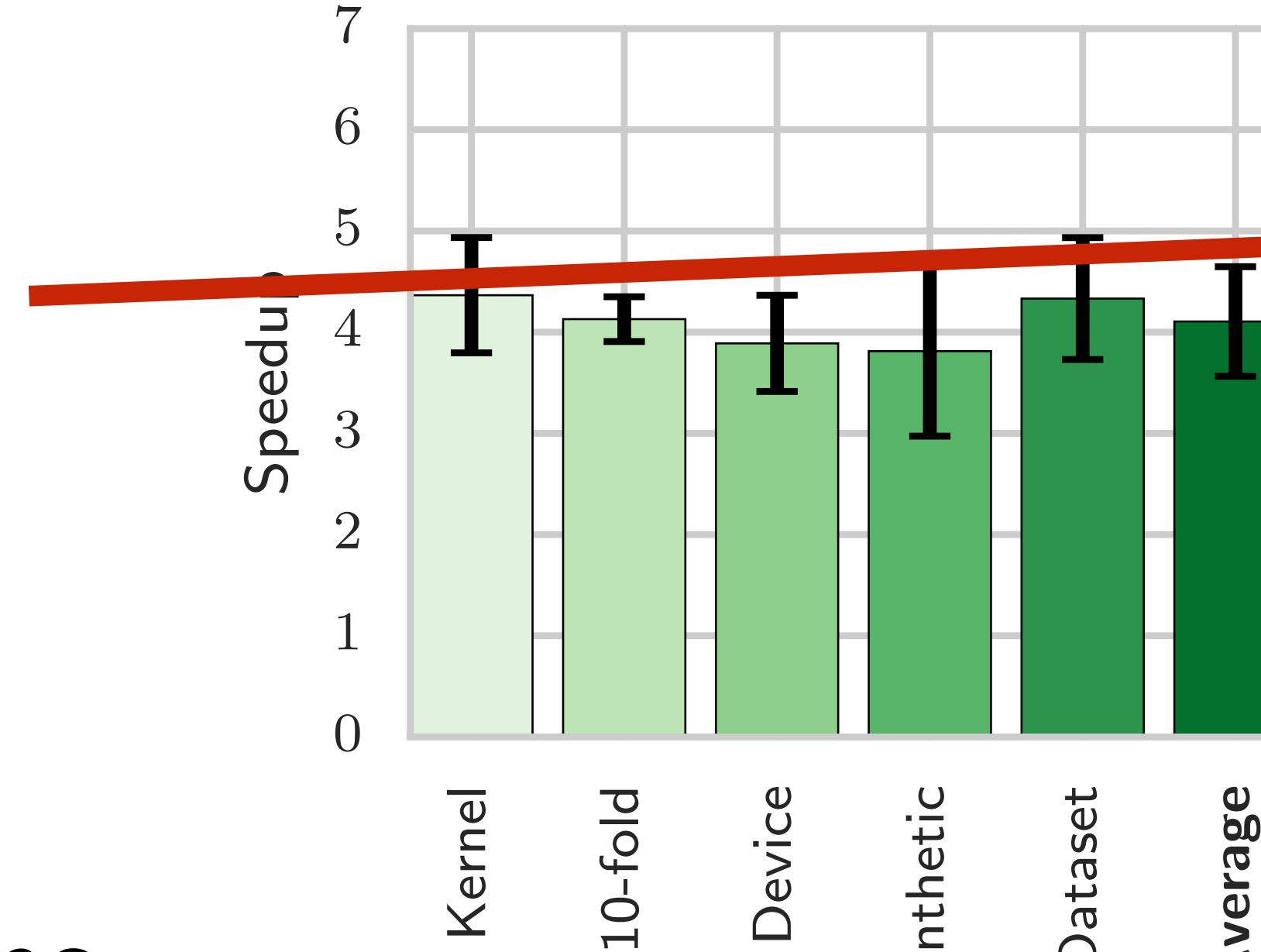
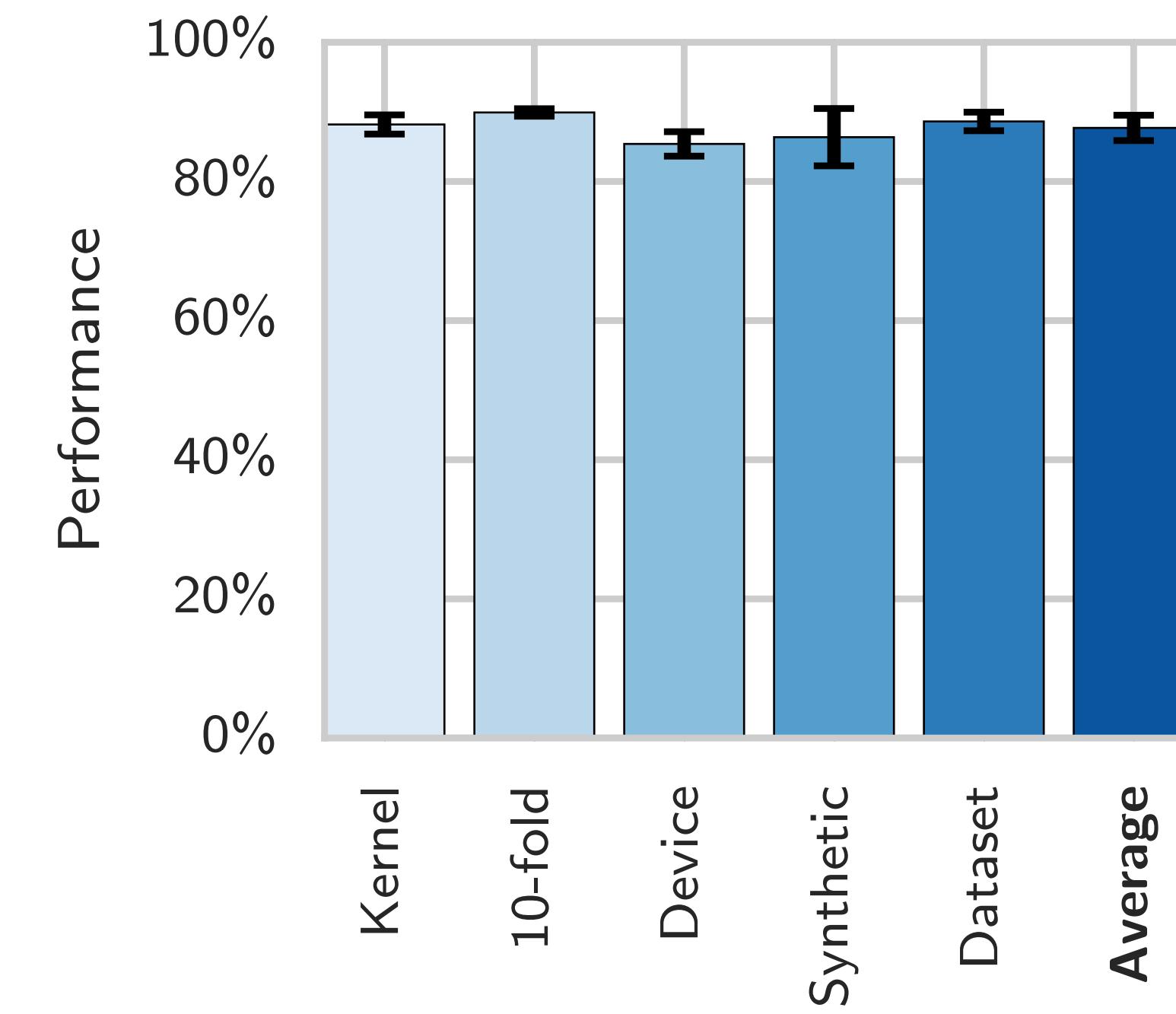
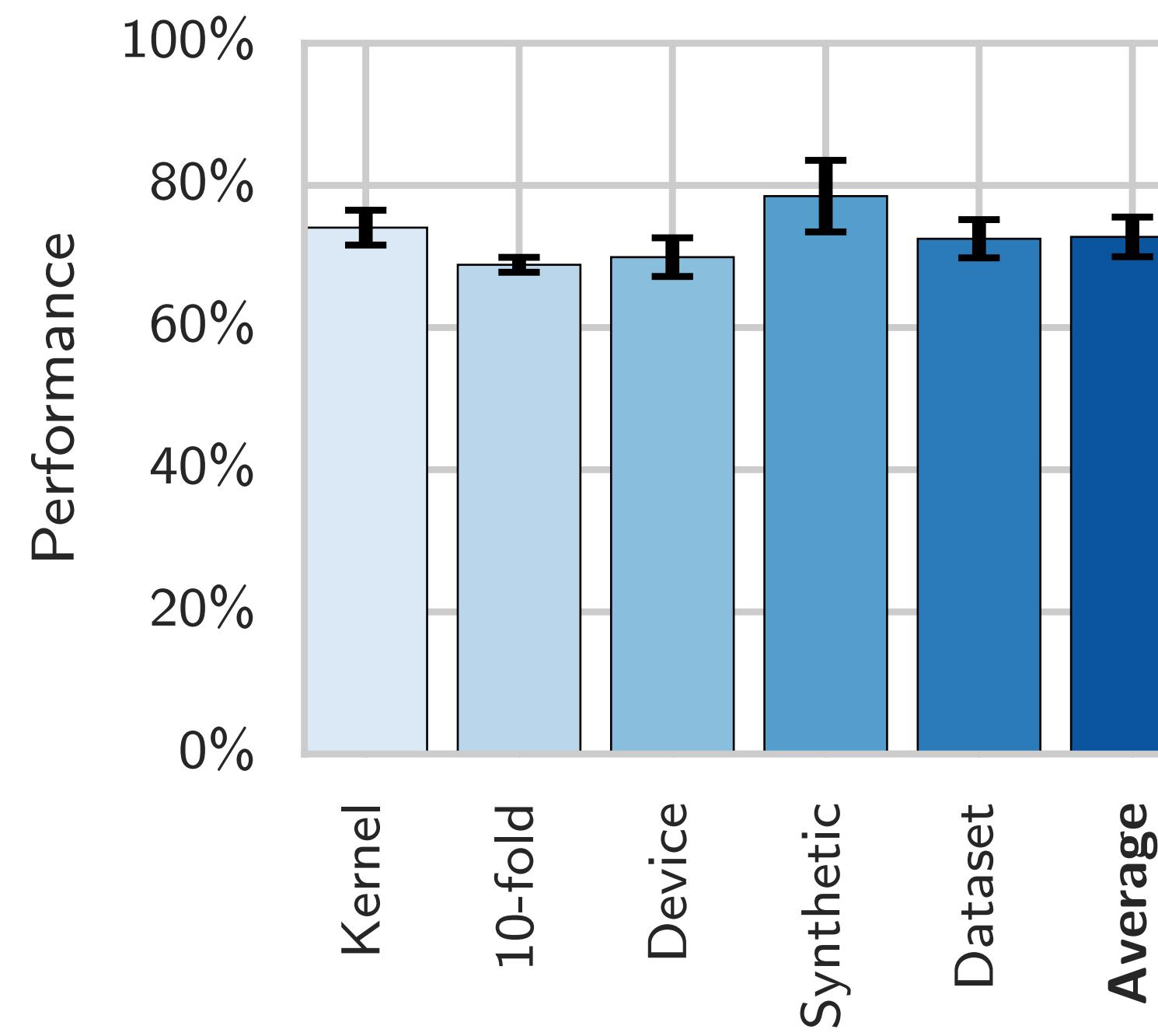
Runtime regression



Speedup regression

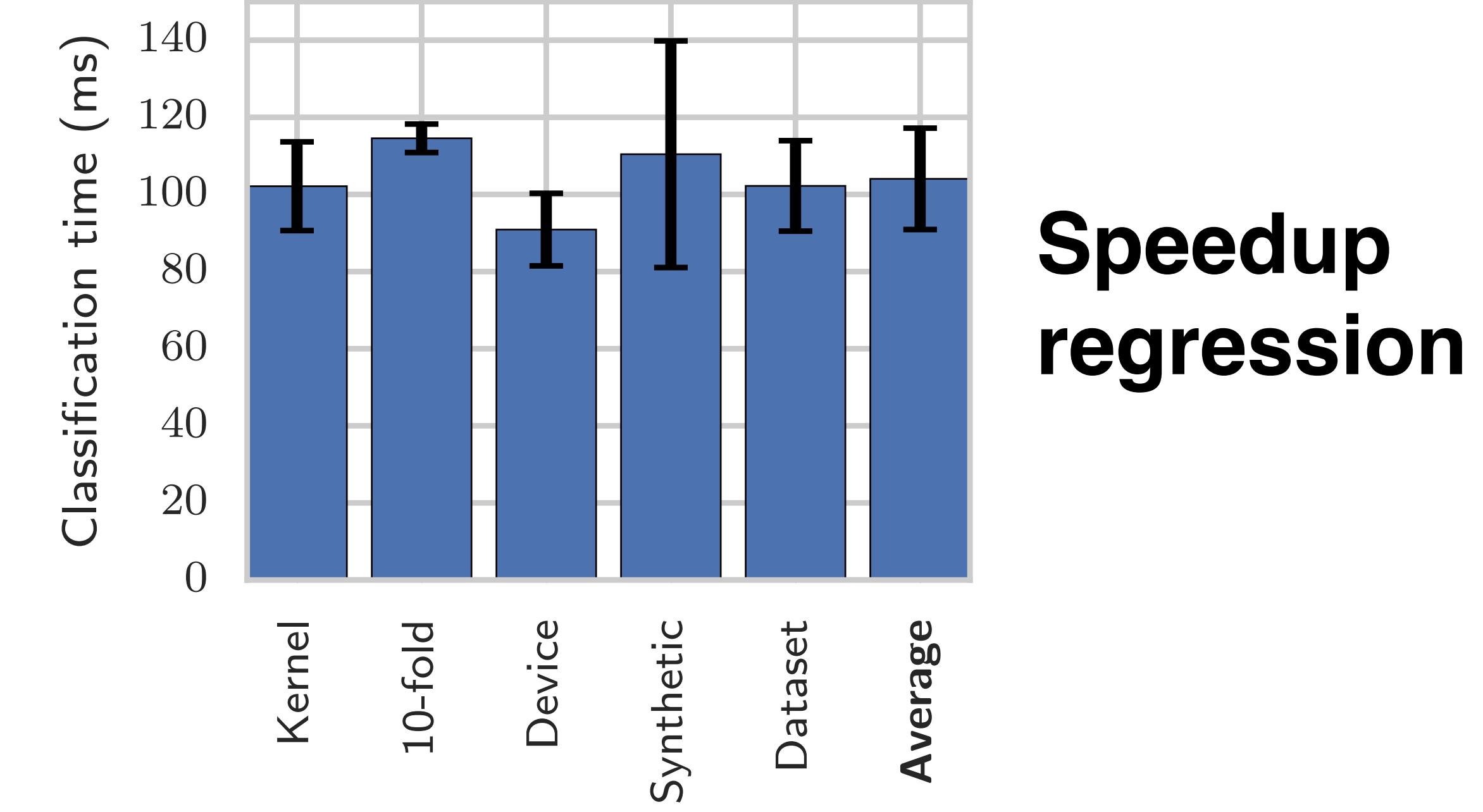
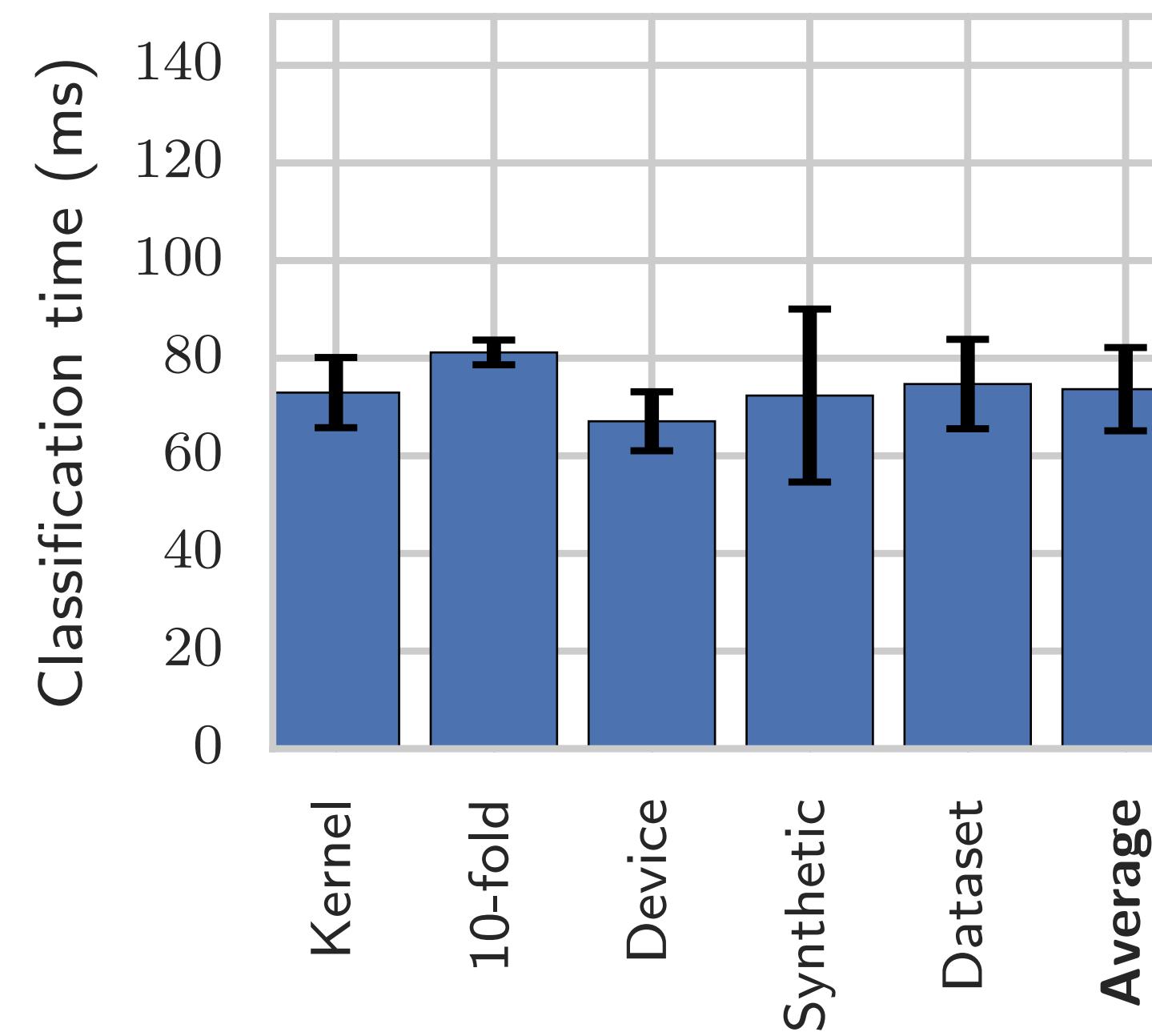
Highest speedup

Runtime regression



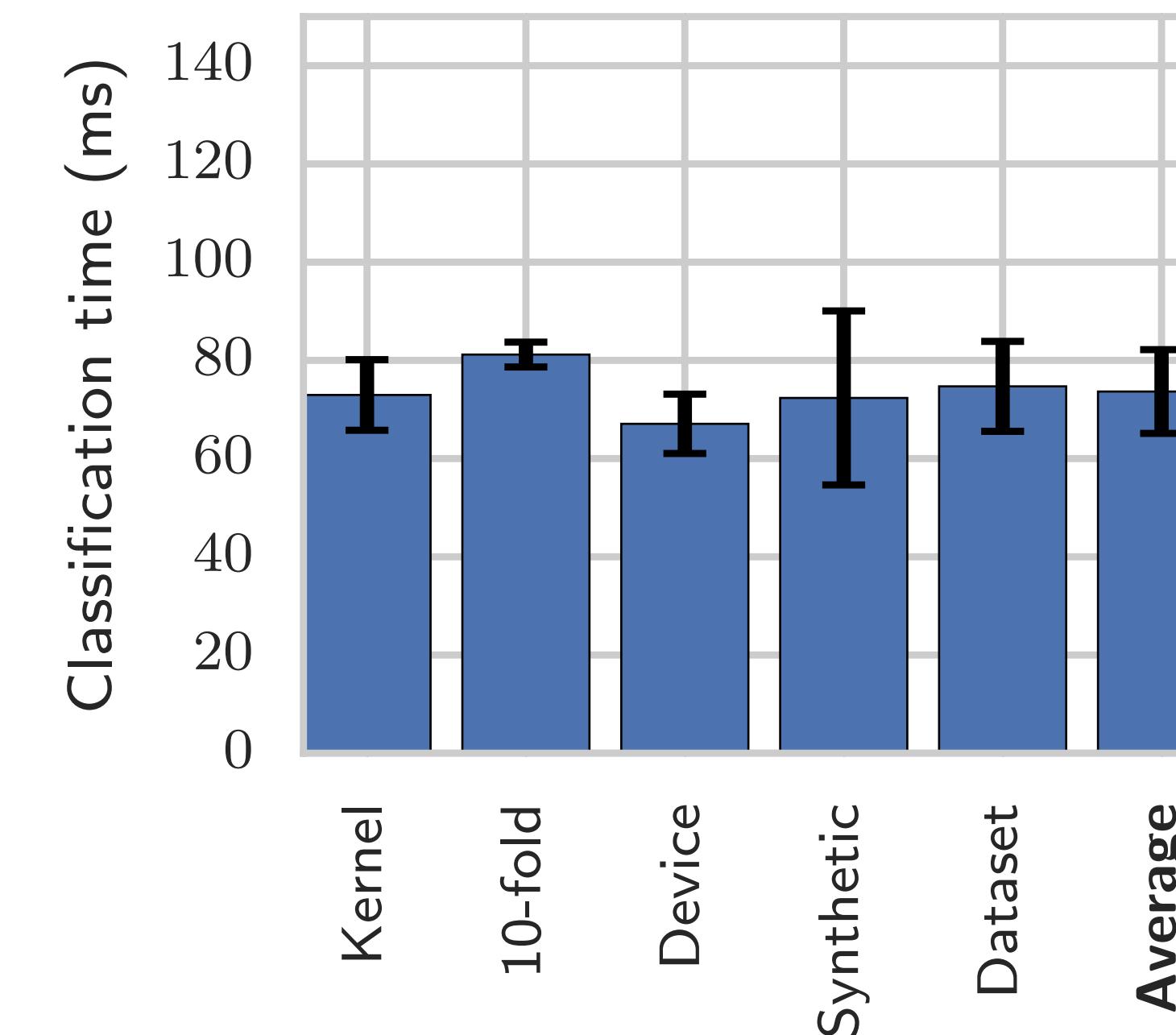
Speedup regression

Runtime regression

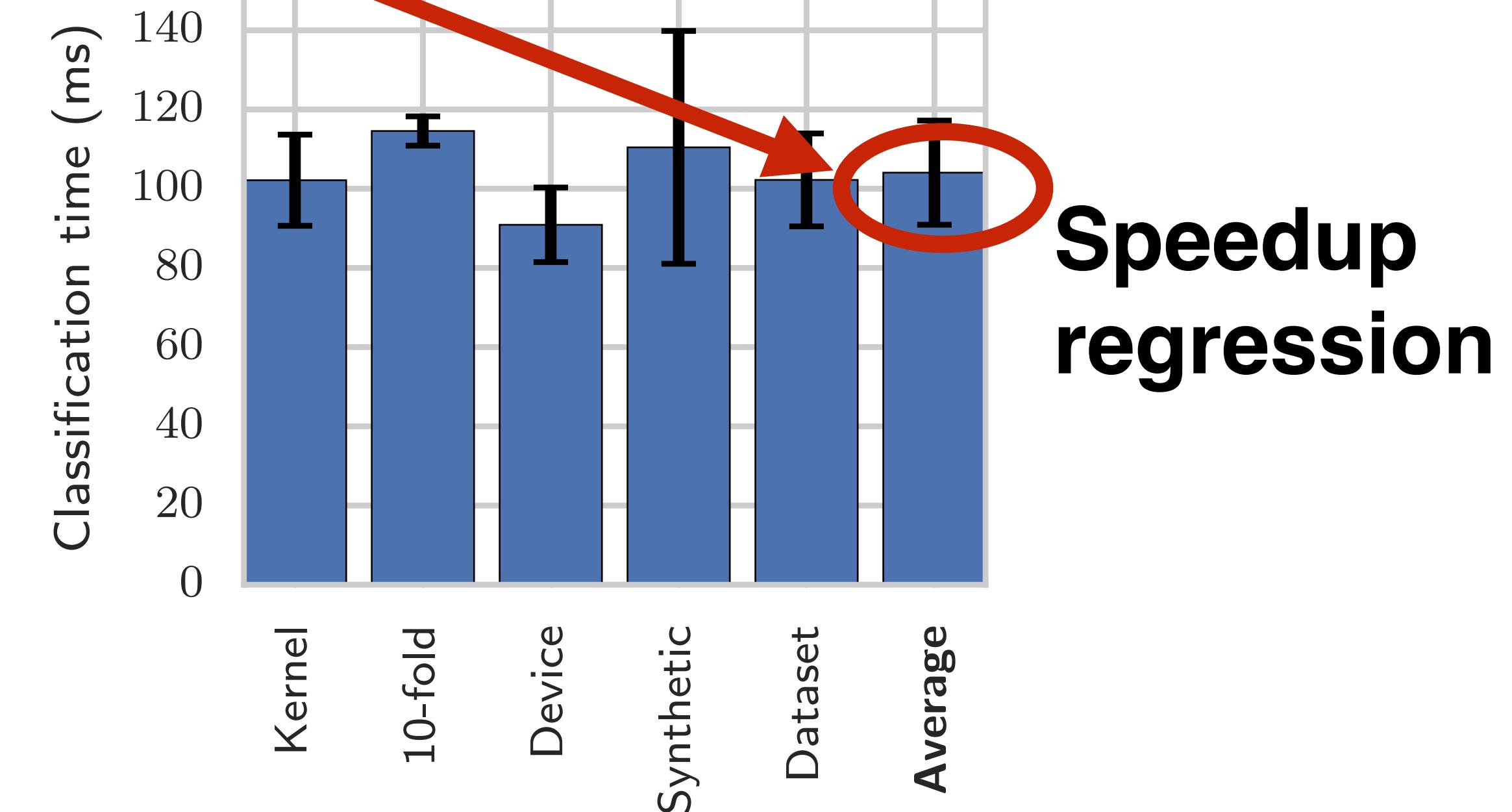


Speedup regression

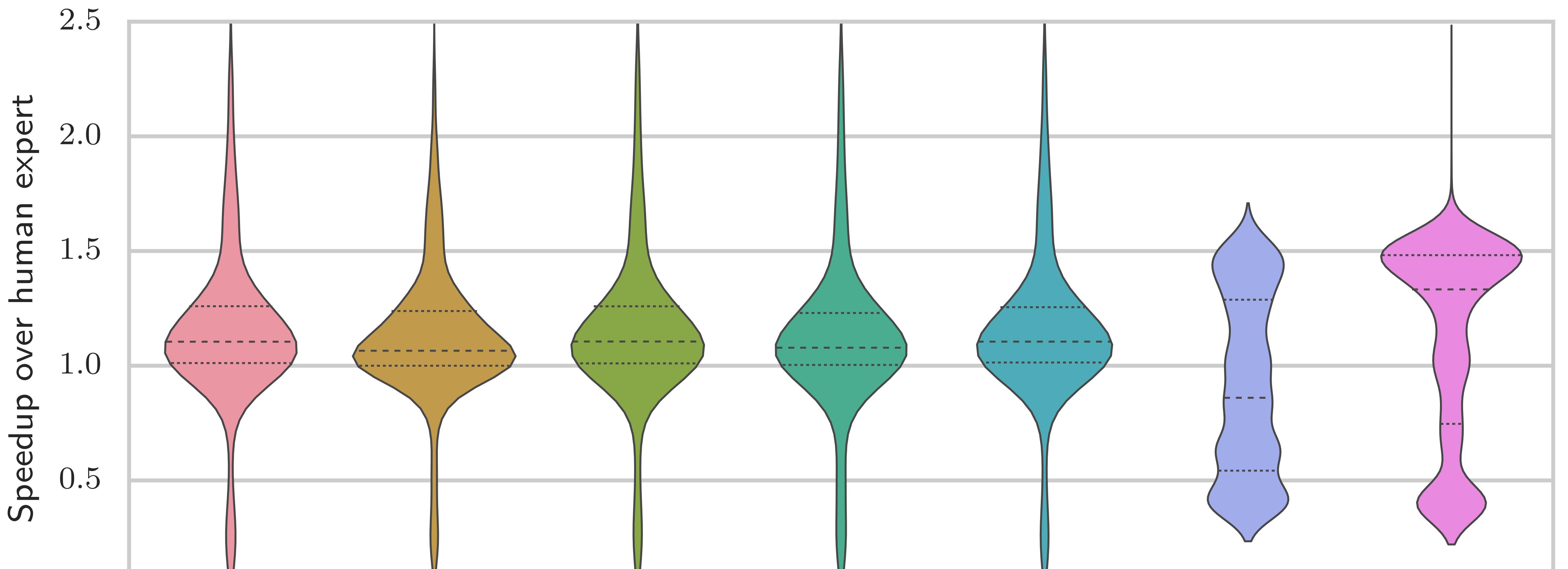
Runtime regression



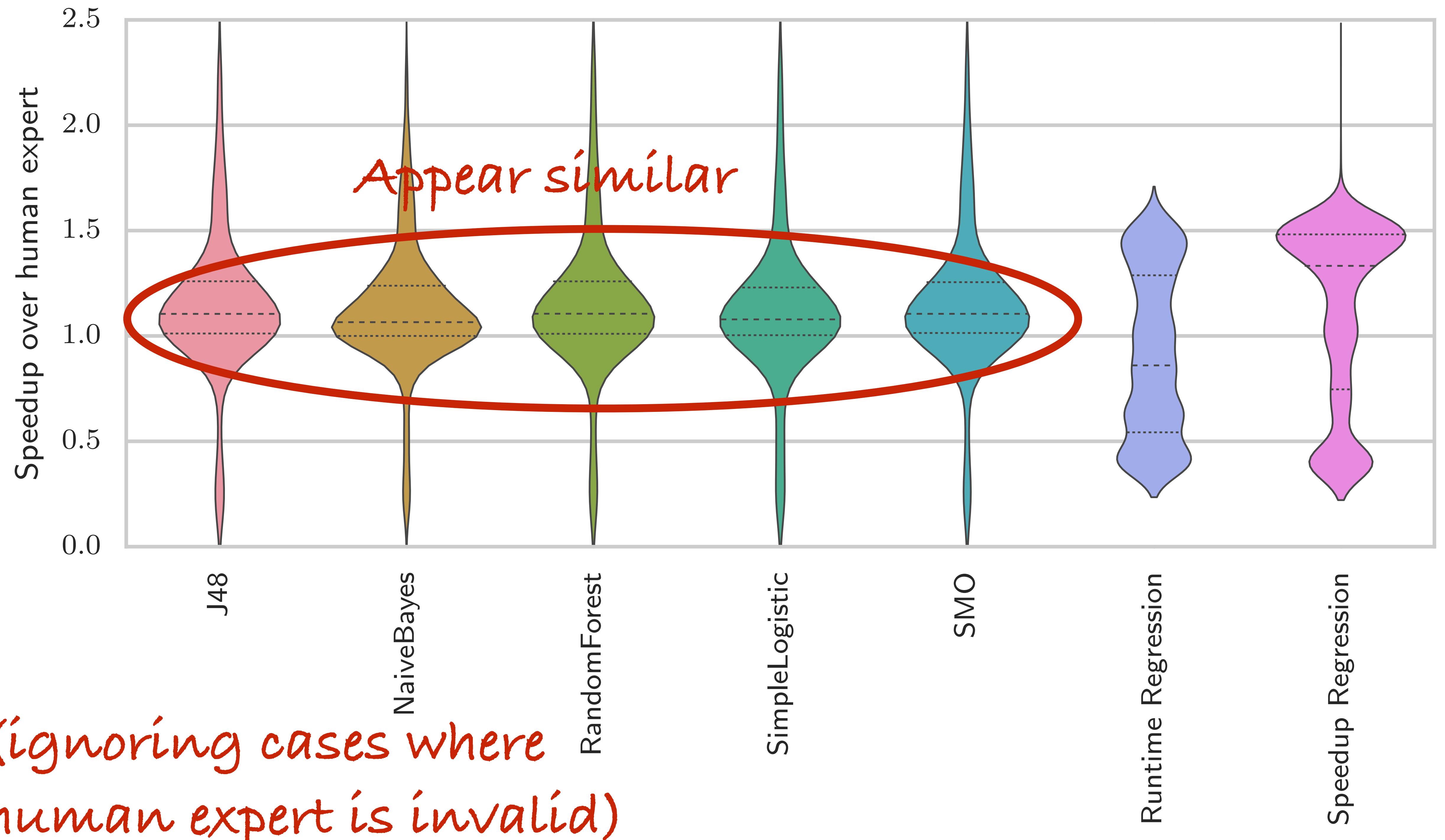
40x slower
than j48

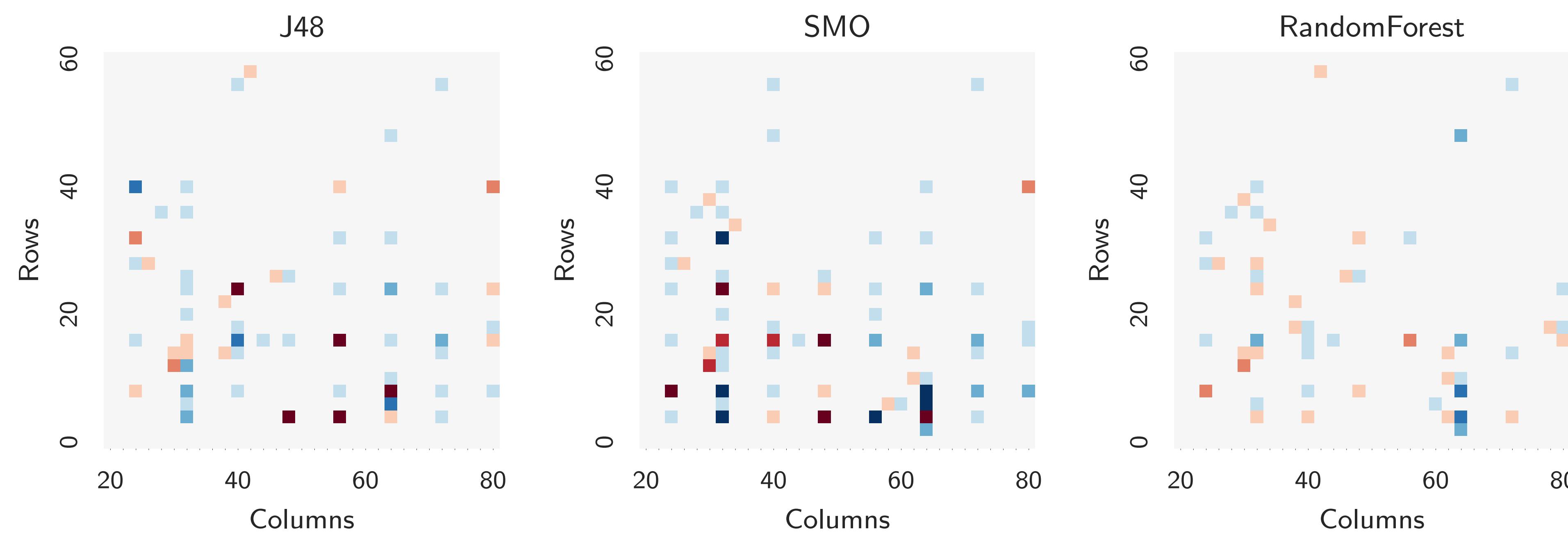


Speedup regression

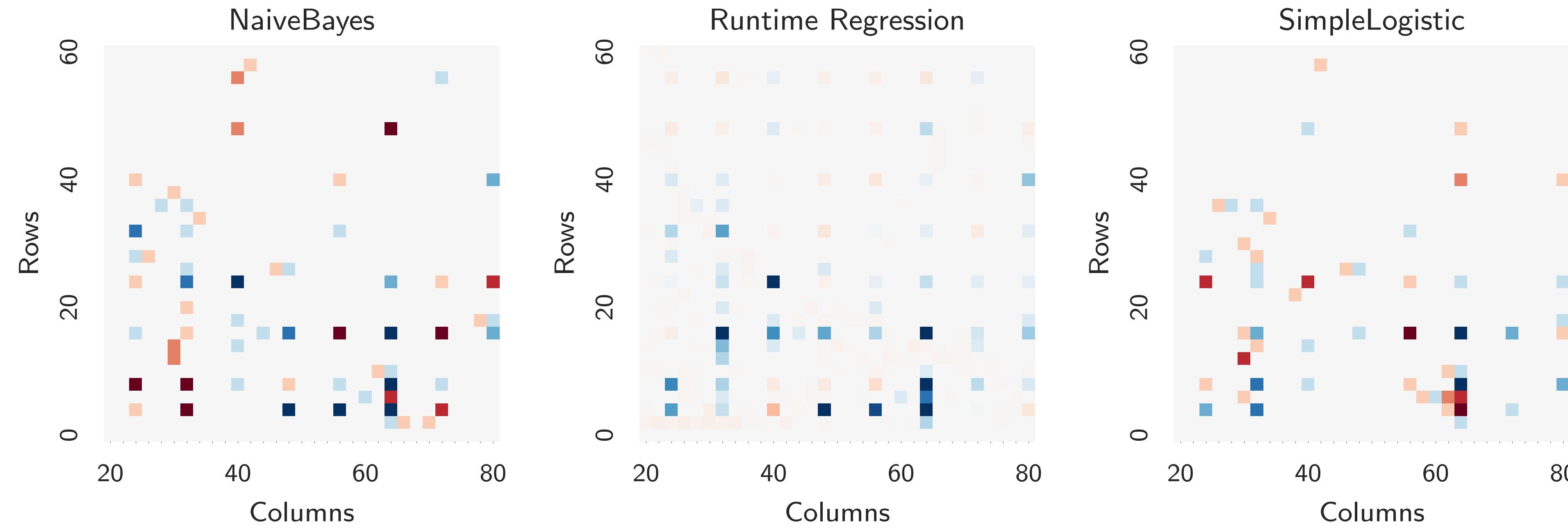


(ignoring cases where
human expert is invalid)





very different prediction characteristics



Results

1. Average 15x speedup best/worst
2. 32% of optimal params unique
3. Classification achieves 90% of optimal performance
4. Speedup regression has highest median perf, but takes longer

Conclusions

Workgroup size is critical for stencil performance

Setting workgroup size depends on device, kernel, dataset

Static tuning fails to exploit available performance

We present *three* methodologies for autotuning OpenCL workgroup size

Trade-offs between prediction cost and training cost

Achieving average 1.22x speedup over *human expert*, with increased reliability

Details in the paper!

Autotuning OpenCL Workgroup Size for Stencil Patterns

Chris Cummins Pavlos Petoumenos Michel Steuwer Hugh Leather

University of Edinburgh

c.cummins@ed.ac.uk, ppetoume@inf.ed.ac.uk, michel.steuwer@ed.ac.uk, hleather@inf.ed.ac.uk

Abstract

Selecting an appropriate workgroup size is critical for the performance of OpenCL kernels, and requires knowledge of the underlying hardware, the data being operated on, and the implementation of the kernel. This makes portable performance of OpenCL programs a challenging goal, since simple heuristics and statically chosen values fail to exploit the available performance. To address this, we propose the use of machine learning-enabled autotuning to automatically predict workgroup sizes for stencil patterns on CPUs and multi-GPUs.

We present three methodologies for predicting workgroup sizes. The first, using classifiers to select the optimal workgroup size. The second and third proposed methodologies employ the novel use of regressors for performing classification by predicting the runtime of kernels and the relative performance of different workgroup sizes, respectively. We evaluate the effectiveness of each technique in an empirical study of 429 combinations of architecture, kernel, and dataset, comparing an average of 629 different workgroup sizes for each. We find that autotuning provides a median $3.79 \times$ speedup over the best possible fixed workgroup size, achieving 94% of the maximum performance.

1. Introduction

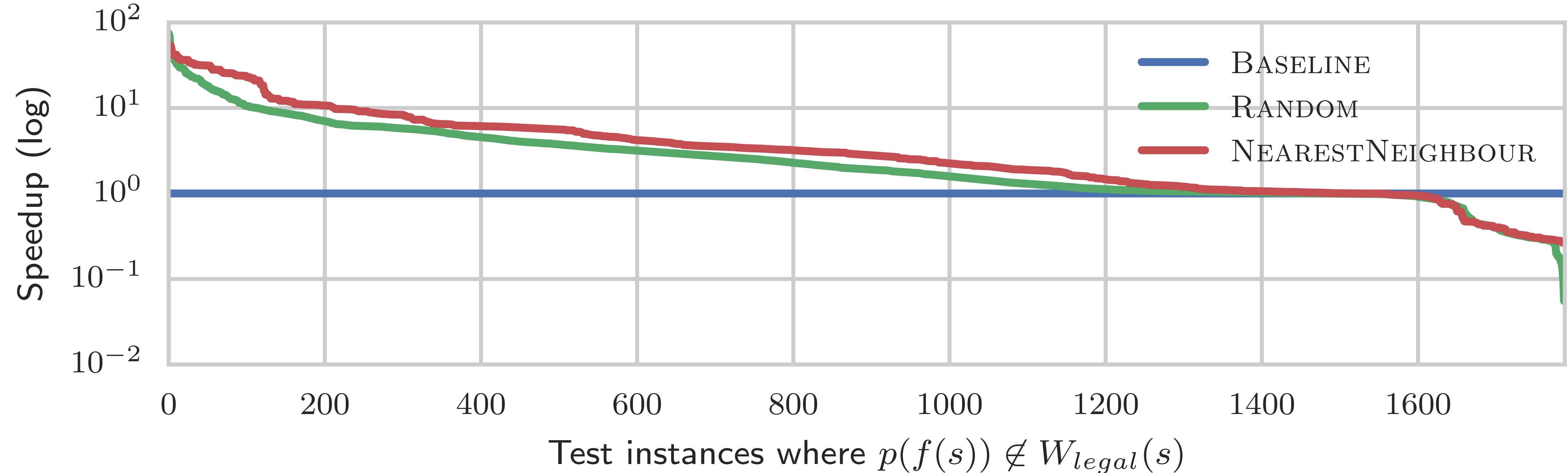
Stencil codes have a variety of computationally demanding uses from fluid dynamics to quantum mechanics. Efficient, tuned stencil implementations are highly sought after, with early work in 2003 by Bolz et al. demonstrating the capability of GPUs for massively parallel stencil operations [1]. Since then, the introduction of the OpenCL standard has introduced greater programmability of heterogeneous devices by providing a vendor-independent layer of abstraction for data parallel programming of CPUs, GPUs, DSPs, and other devices [2]. However, achieving portable performance of OpenCL programs is a hard task — OpenCL kernels are sensitive to properties of the underlying hardware, to the implementation, and even to the *dataset* that is operated upon. This forces developers to laboriously hand-tune performance on a case-by-case basis, since simple heuristics fail to exploit the available performance.

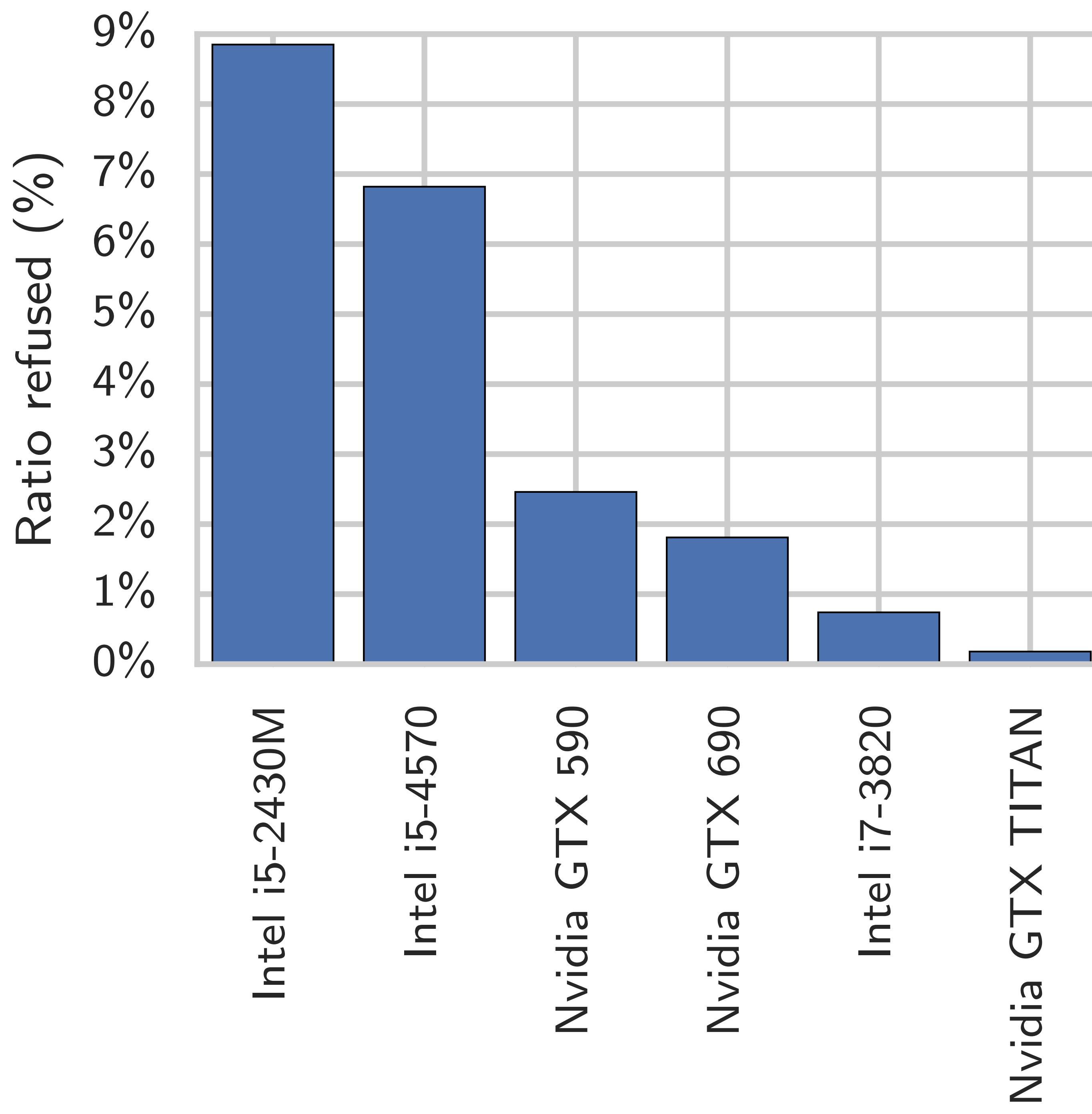


Autotuning OpenCL Workgroup Size for Stencil Patterns

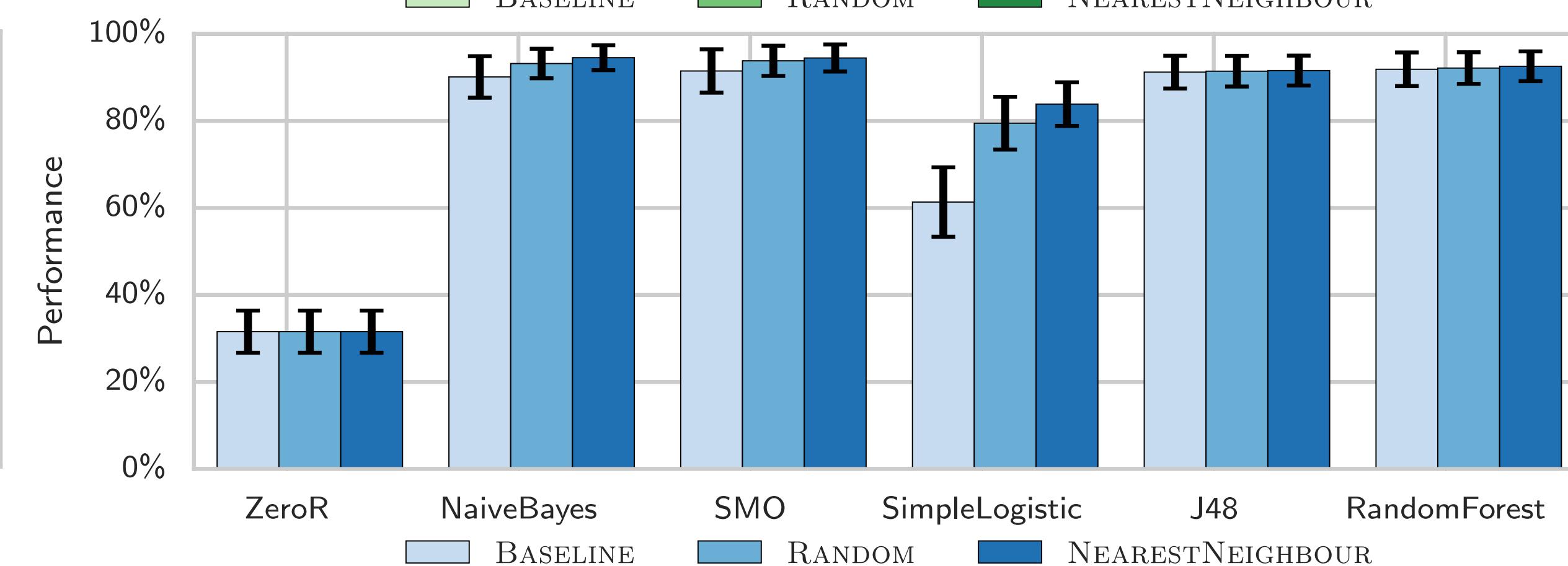
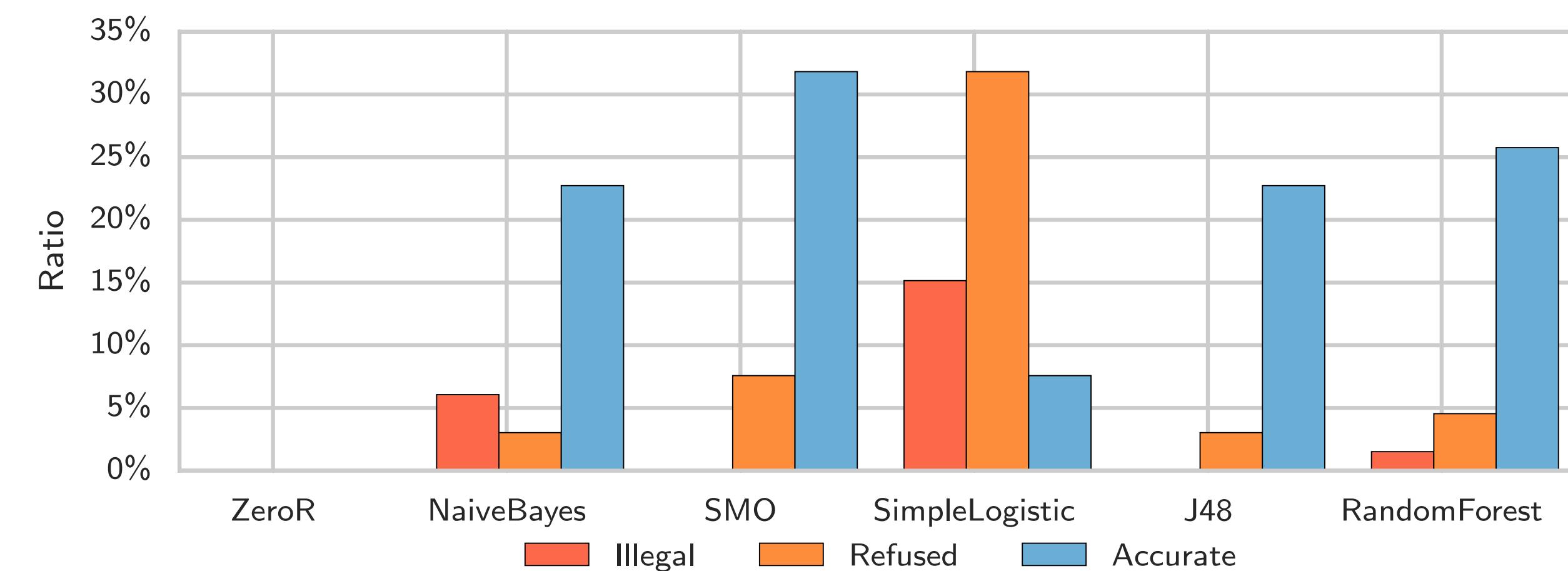
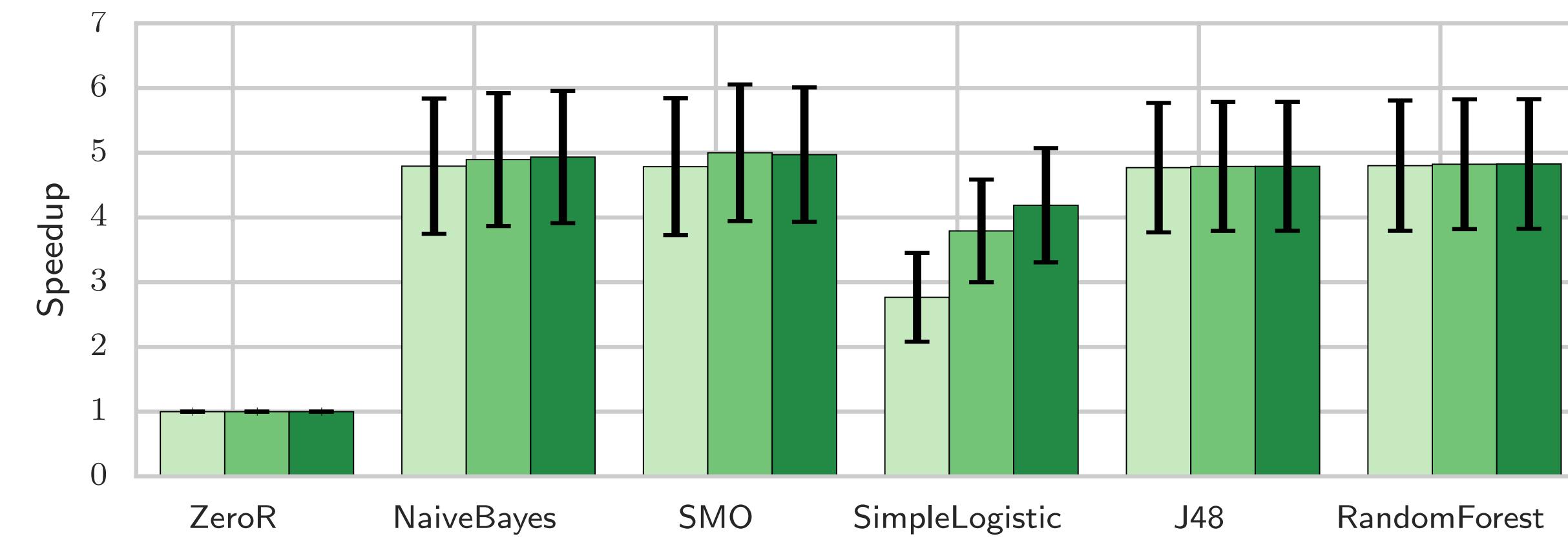
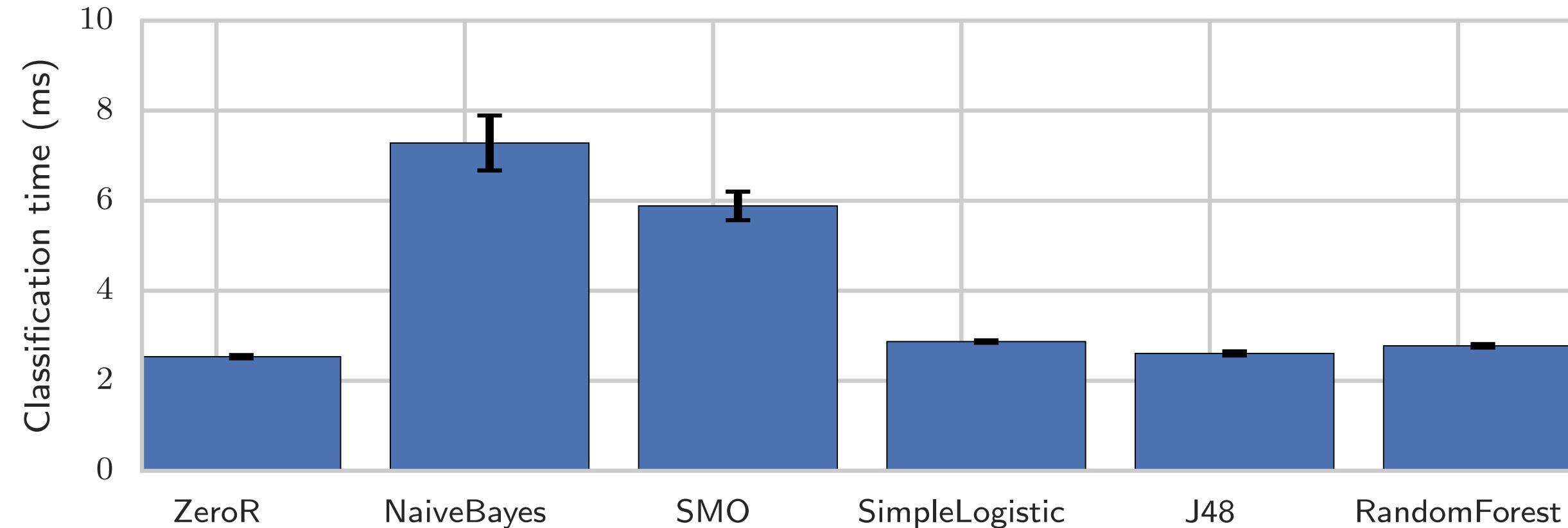
<http://chriscummins.cc>

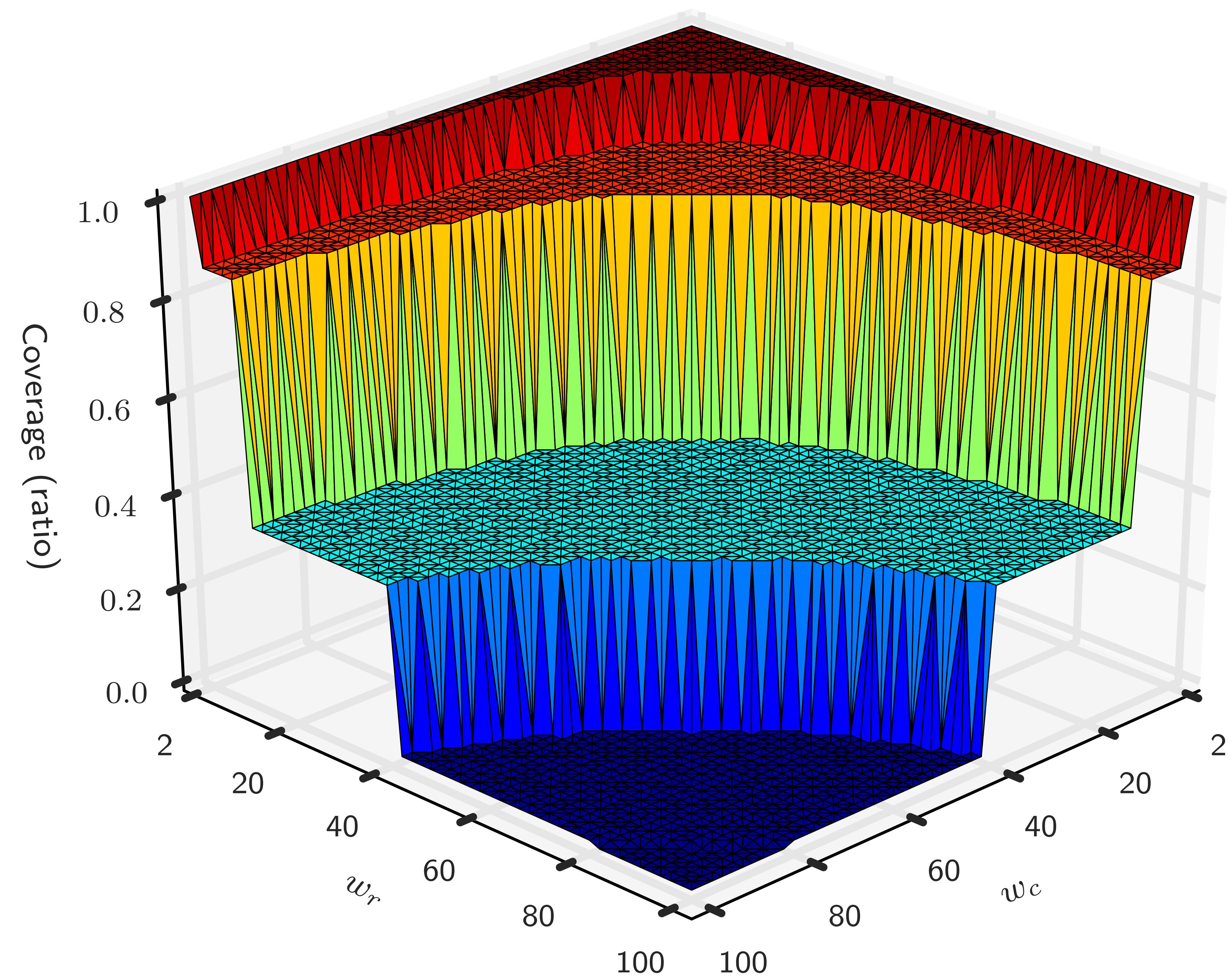
Fallback Handlers

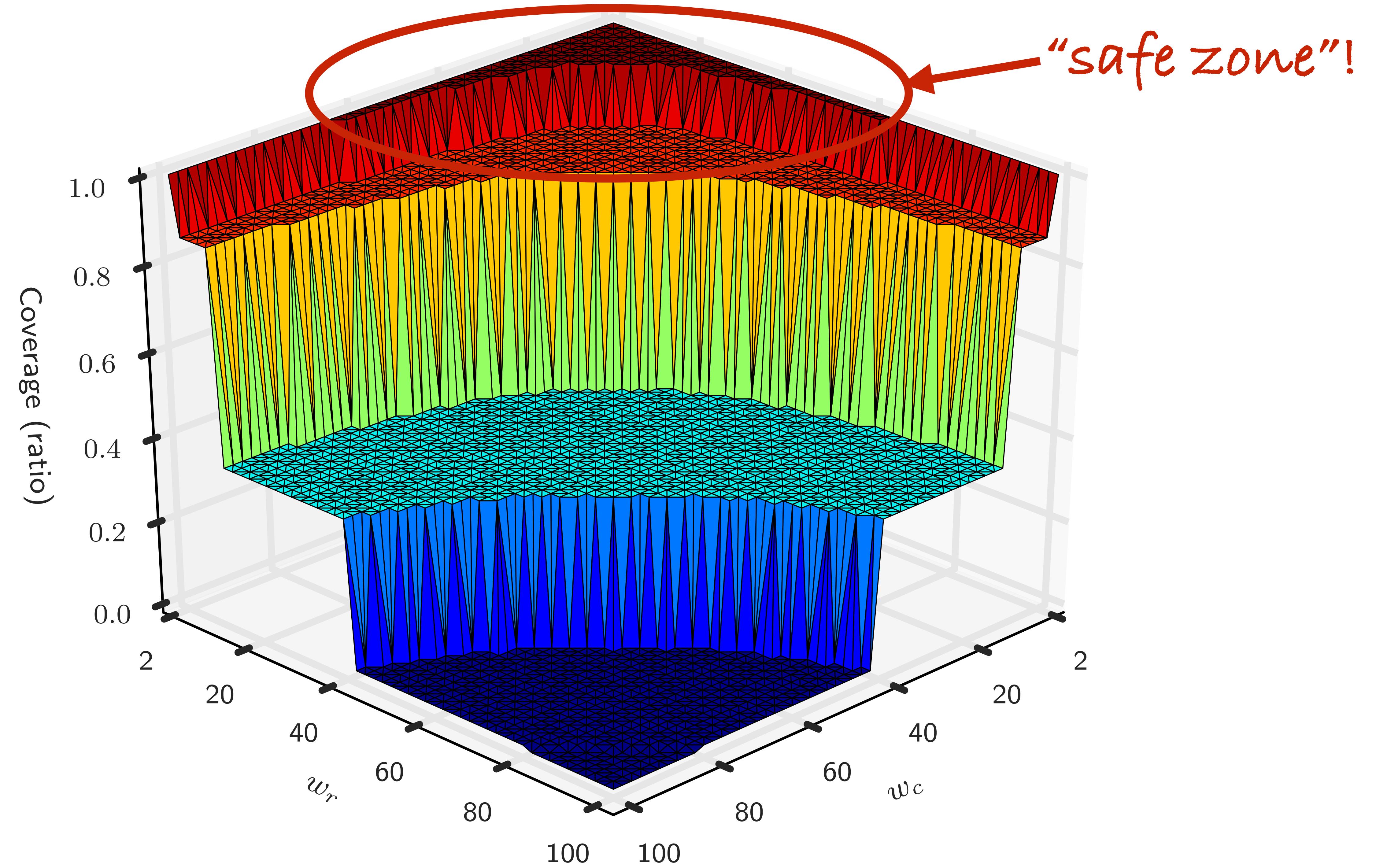


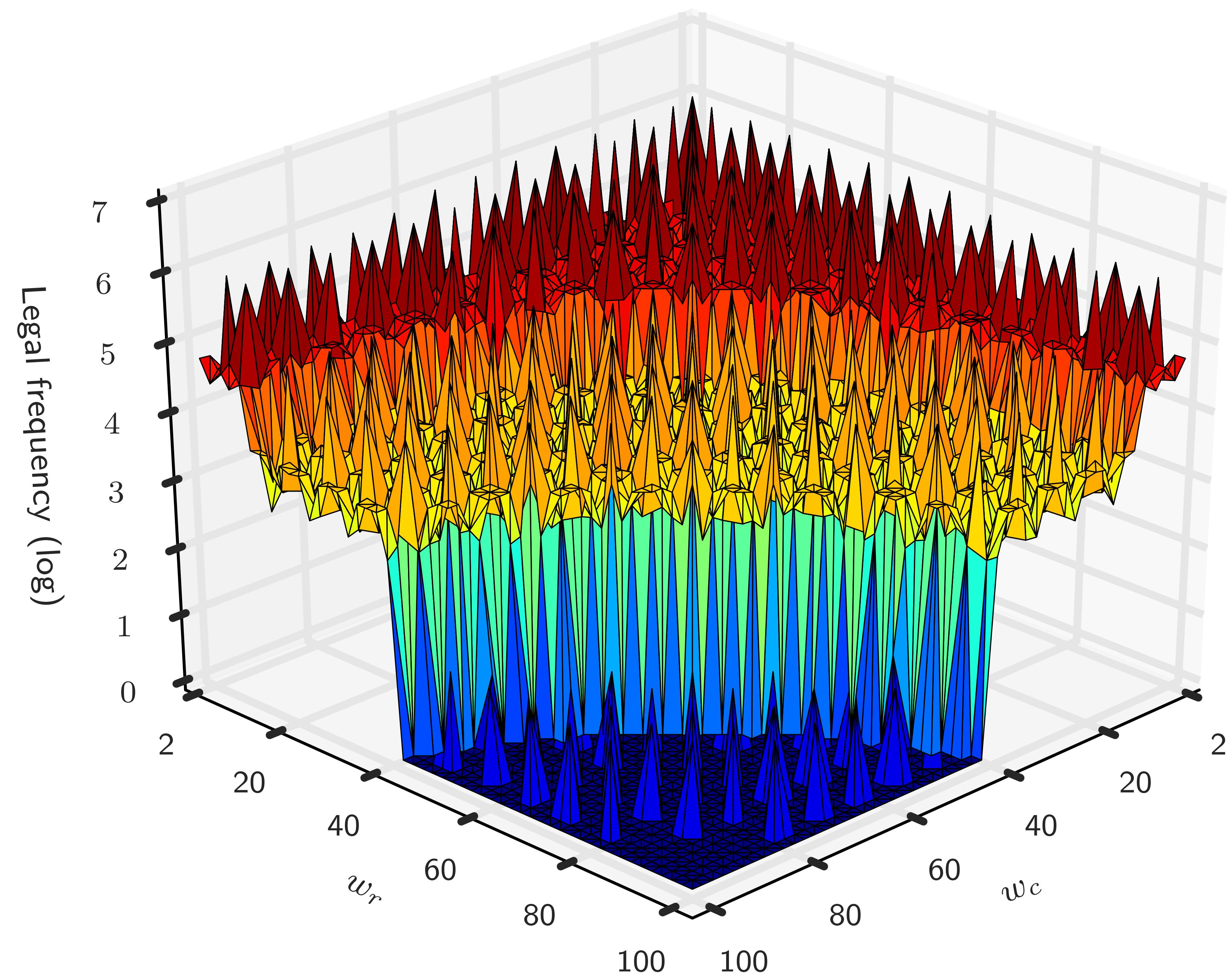


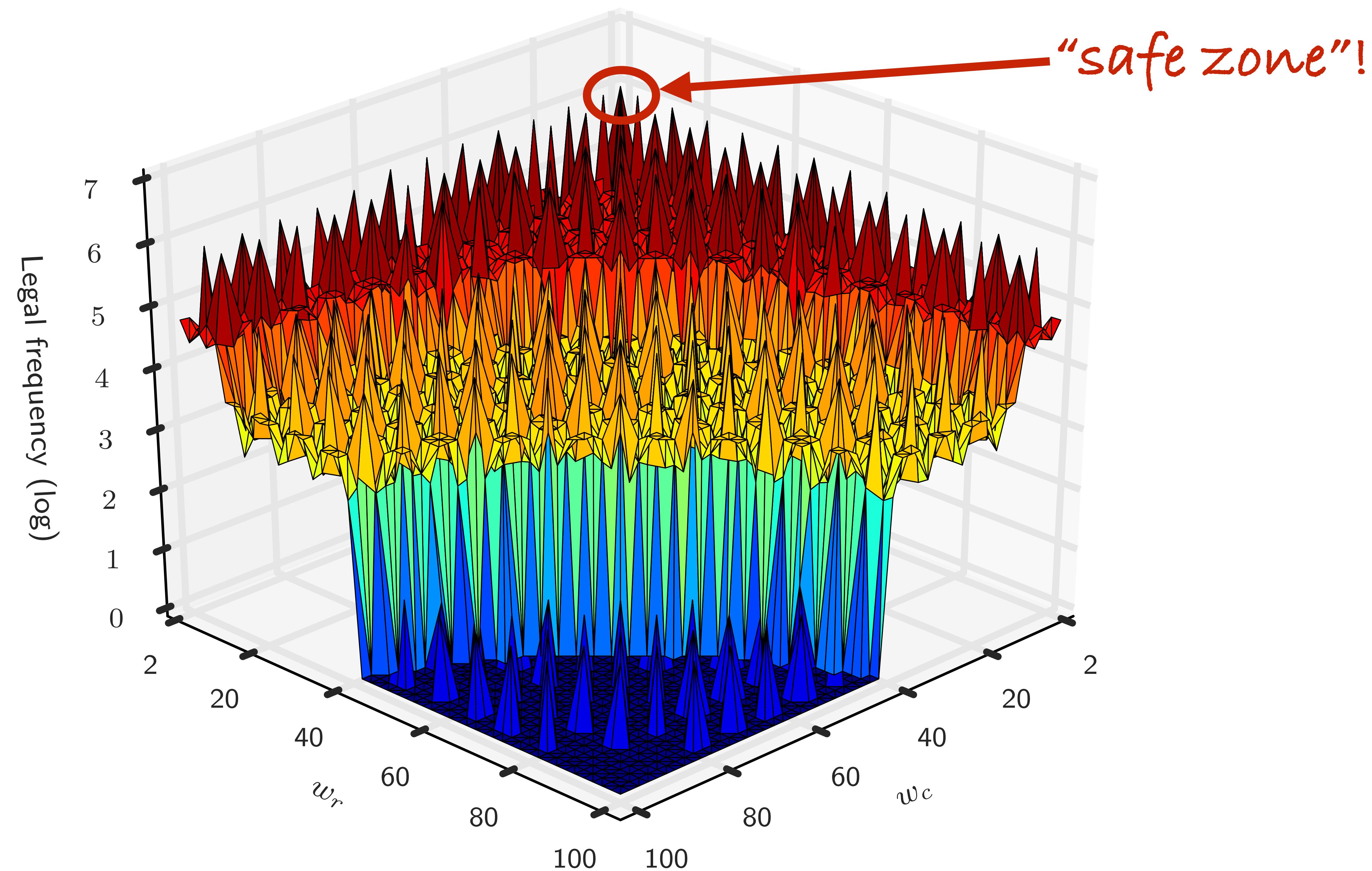
Classifier performance











Palette