

Towards *Collaborative Performance Tuning* of Algorithmic Skeletons

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>

Uncontentious statement:

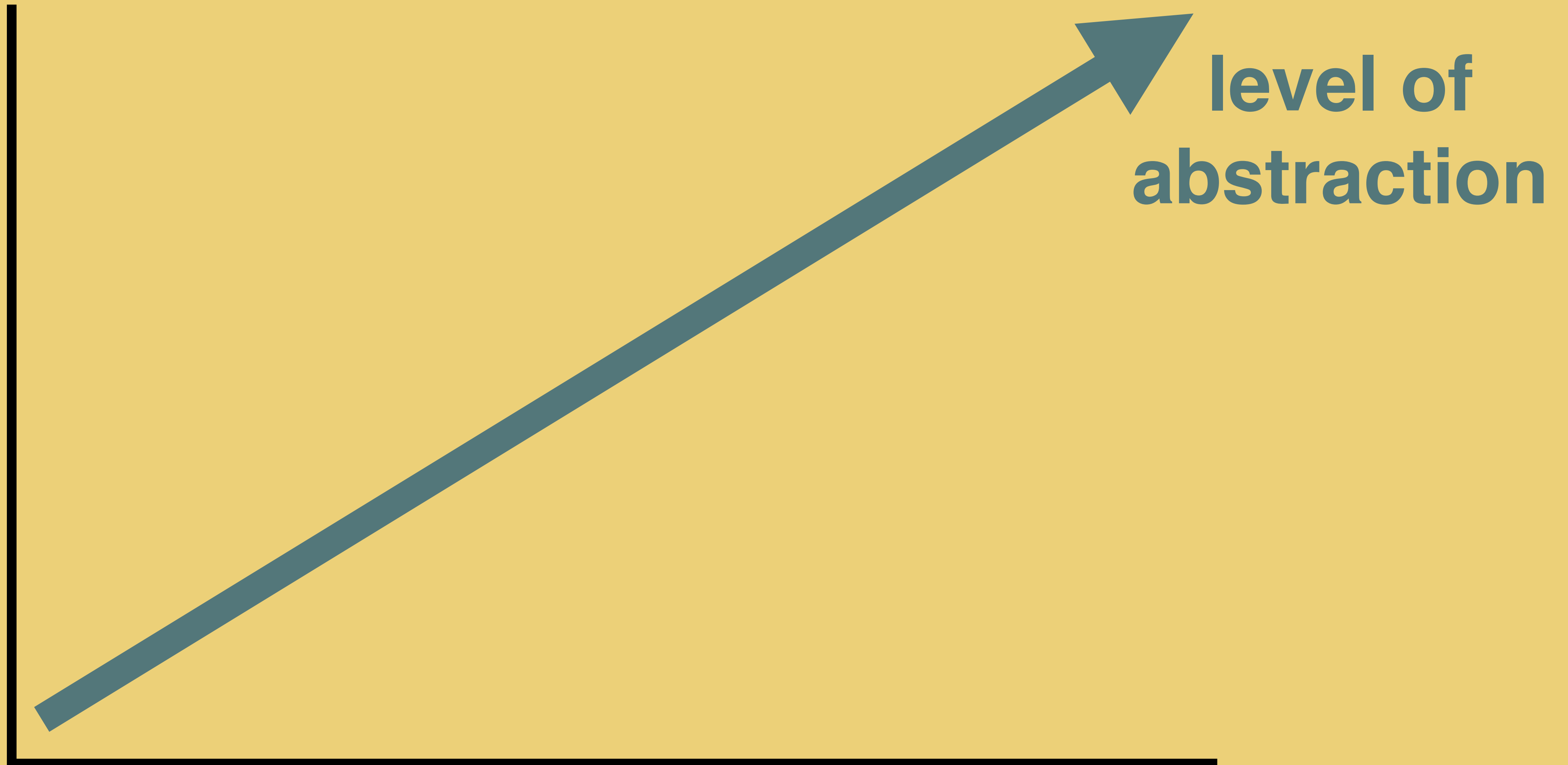
**High level
programming is
great!**

**So why do application
programmers resort
to writing *low level*
code?**

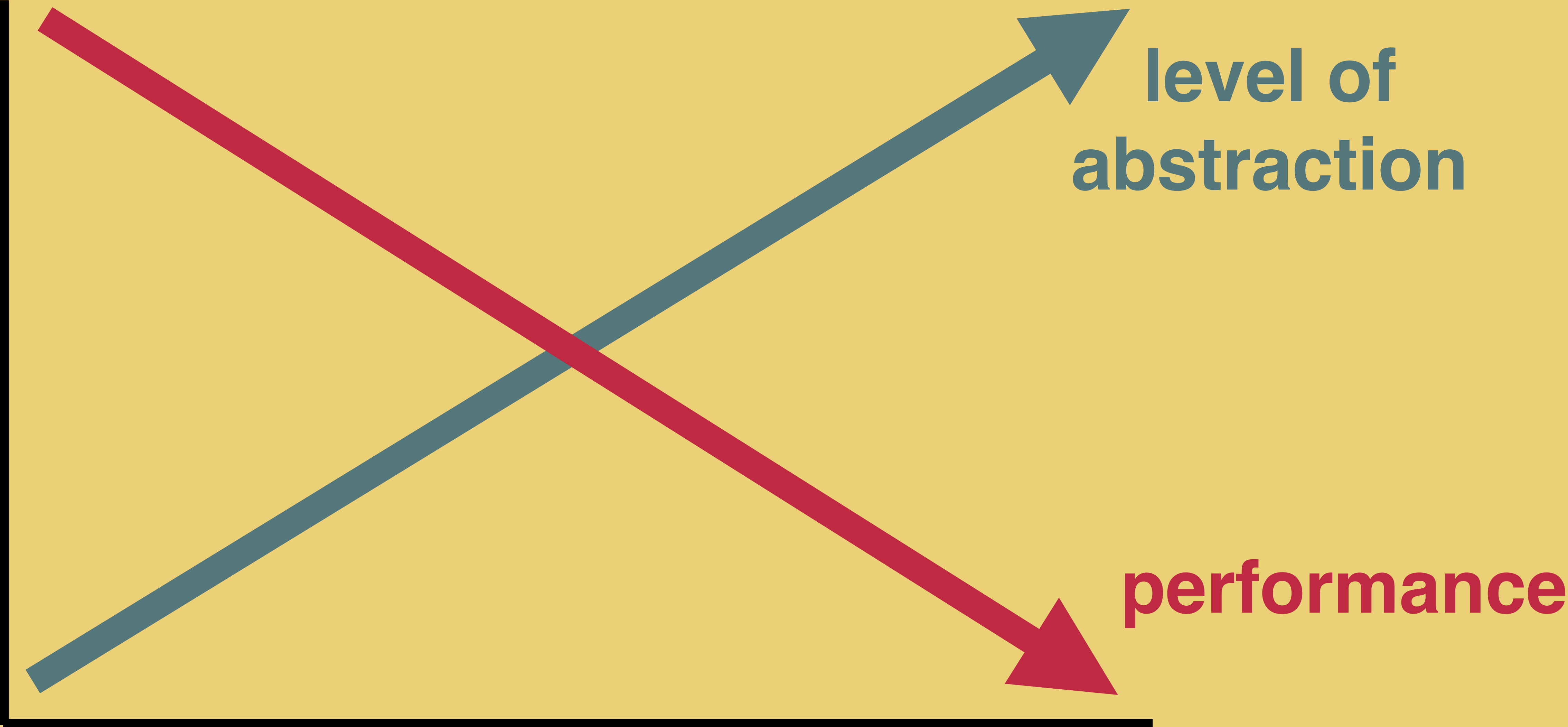
So why do application
programmers resort
to writing *low level*
code?

Performance!

common perception is ...



common perception is ...





**programmers
who care about
*performance***

**programmers
who care about
*abstractions***

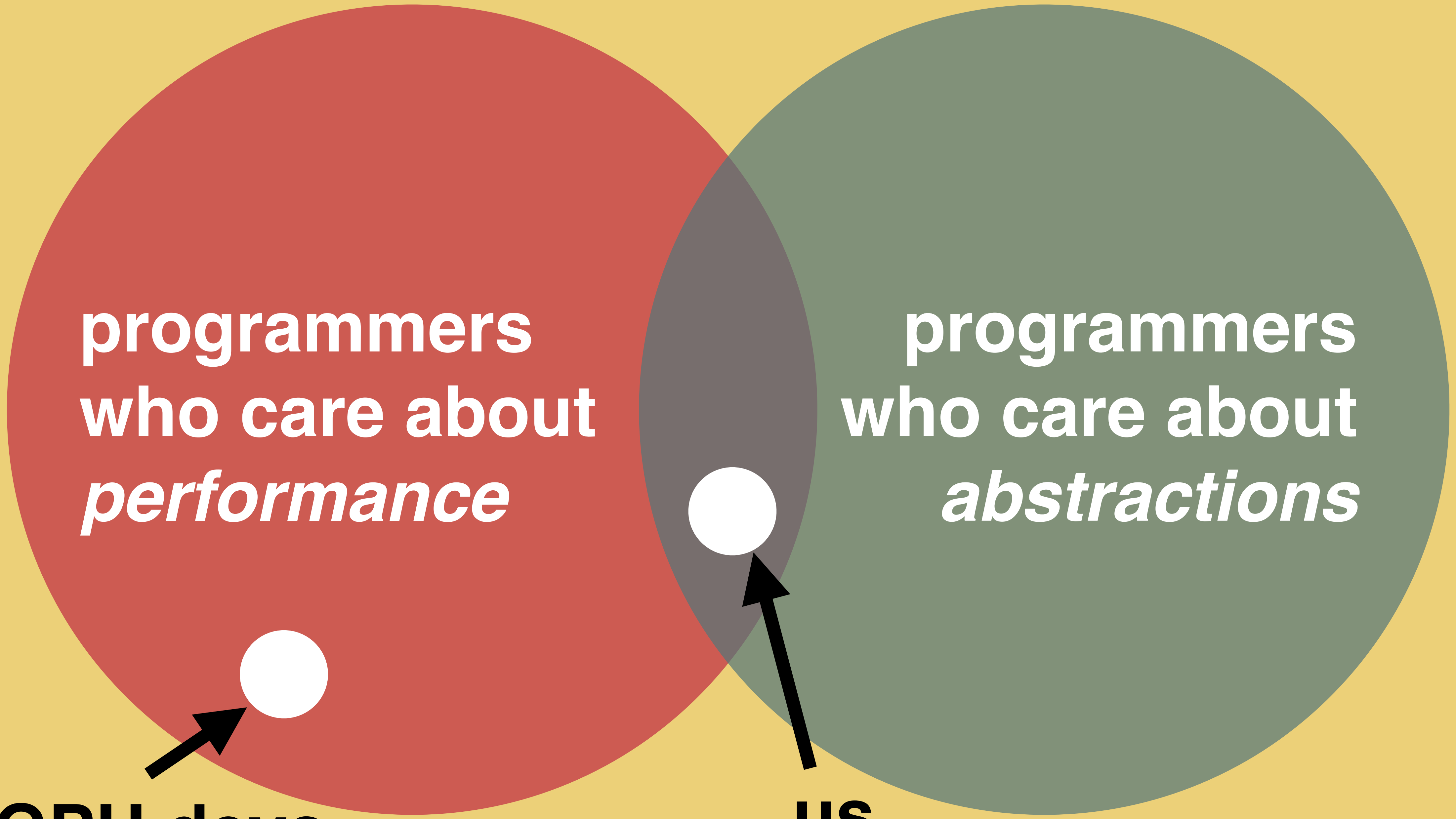


programmers
who care about
performance

programmers
who care about
abstractions

GPGPU devs





A Venn diagram with two overlapping circles on a yellow background. The left circle is red and contains the text 'programmers who care about performance'. The right circle is green and contains the text 'programmers who care about abstractions'. The intersection of the two circles is shaded dark grey and contains a white dot. An arrow points from the text 'GPGPU devs' at the bottom left to this white dot. Another arrow points from the text 'us' at the bottom right to the same white dot.

programmers
who care about
performance

programmers
who care about
abstractions

GPGPU devs

us

**How do we break
the illusion?**

High level code
needs to be at least
competitive with
low level

High level code
needs to be at least
competitive with
low level (but faster would be nice)

Reasons for low level:

**Reasons for low level:
Domain-specific
optimisations**

Reasons for low level:

**Domain-specific
optimisations**

Parameter tuning

Reasons for low level:

**Domain-specific
optimisations**

Parameter tuning

Parameter tuning for Algorithmic Skeletons

Parameter tuning for ~~Algorithmic~~ *Stencil* ~~Skeletons~~ *skeletons*

OpenCL workgroup size

~~Parameter tuning~~

~~for Algorithmic~~

~~Skeletons~~

stencil
skeletons

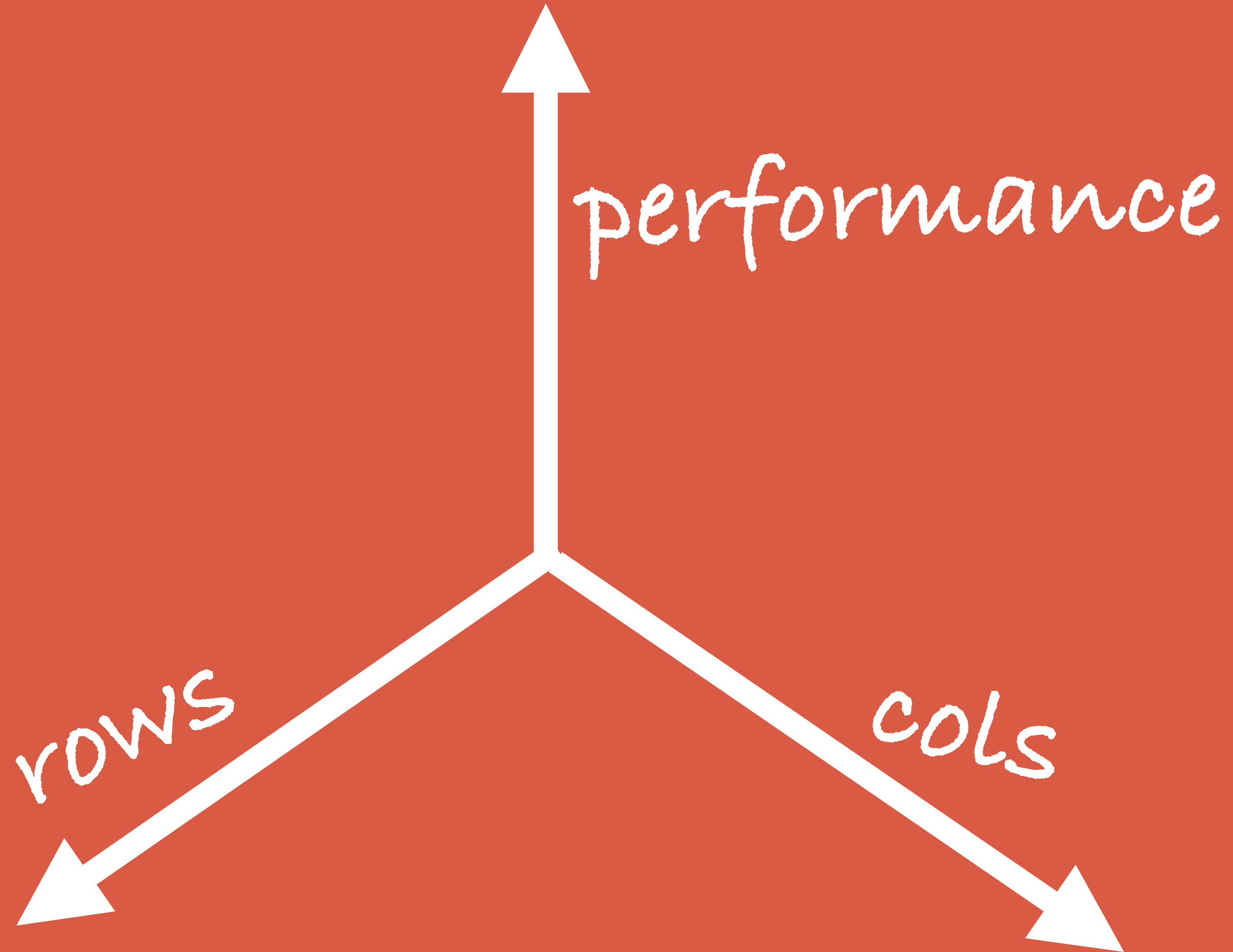
OpenCL workgroup size:

OpenCL workgroup size:
Controls decomposition of
threads.

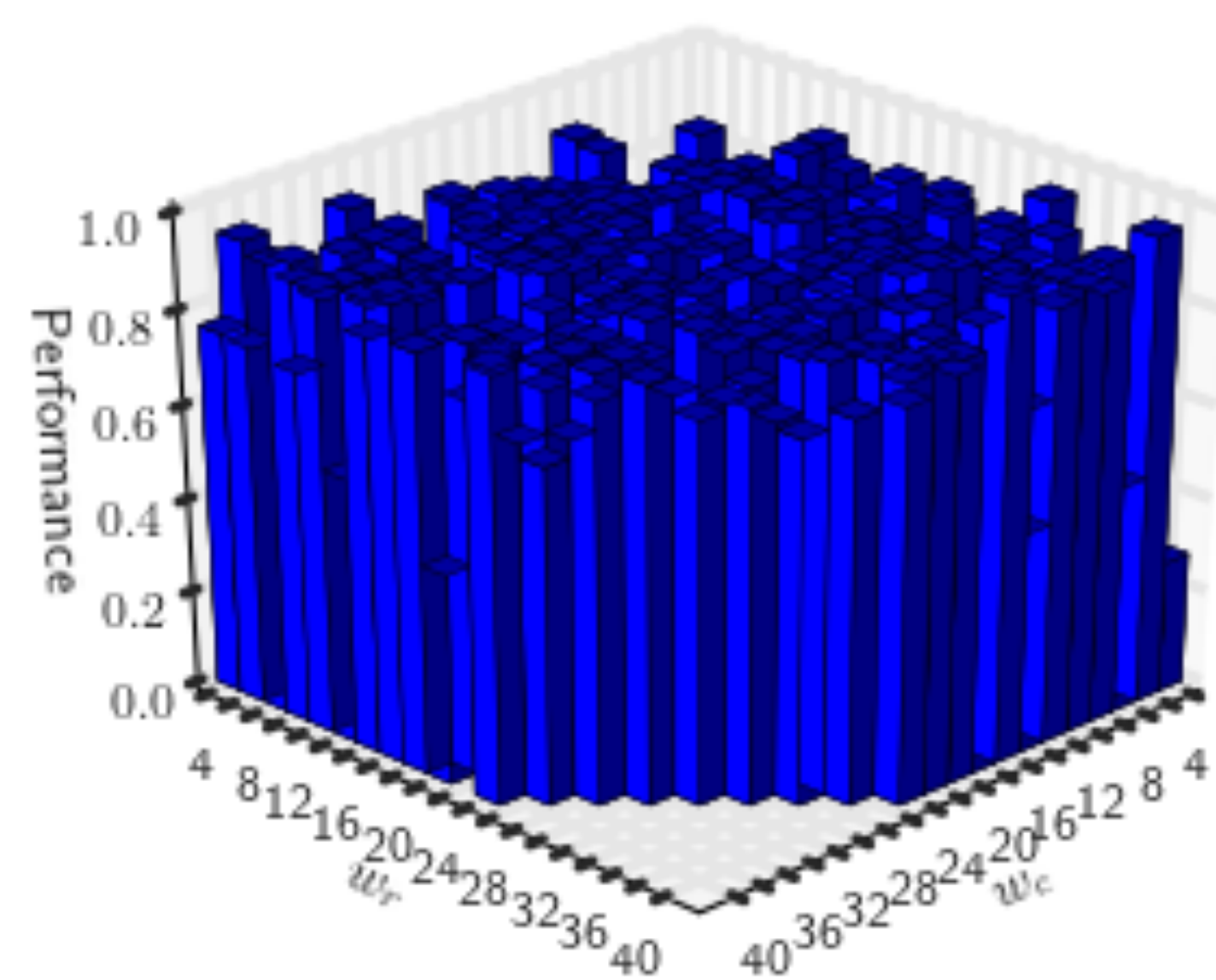
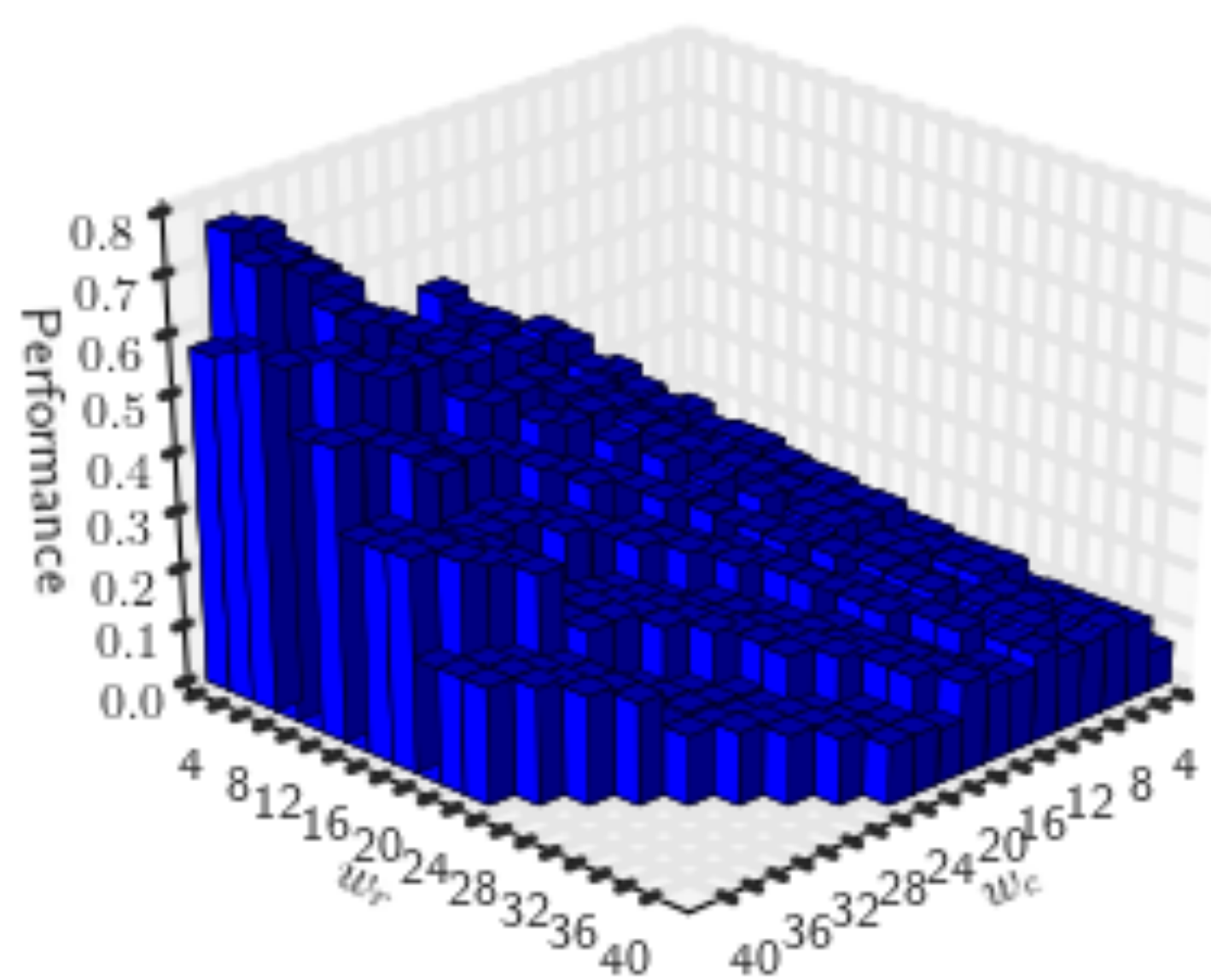
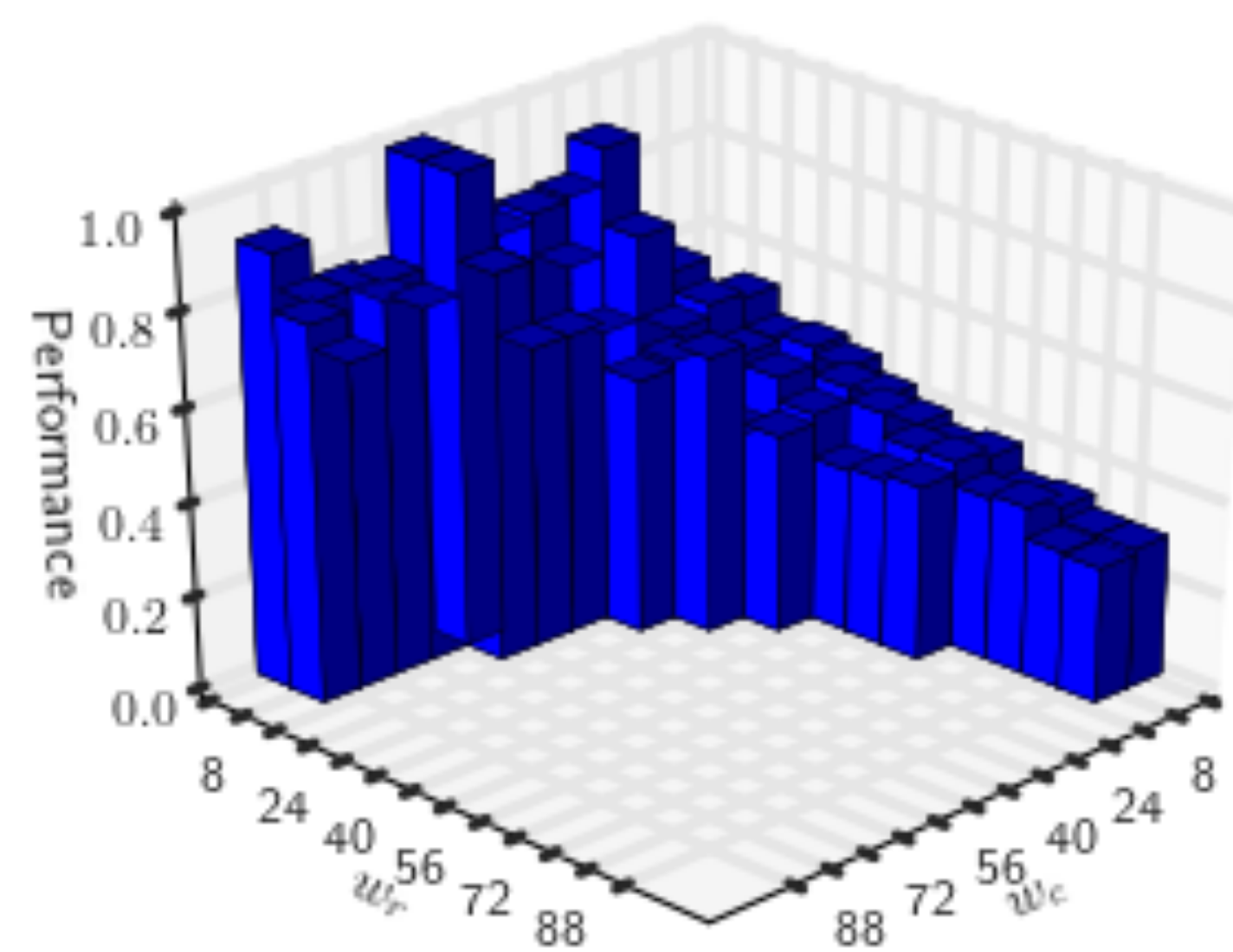
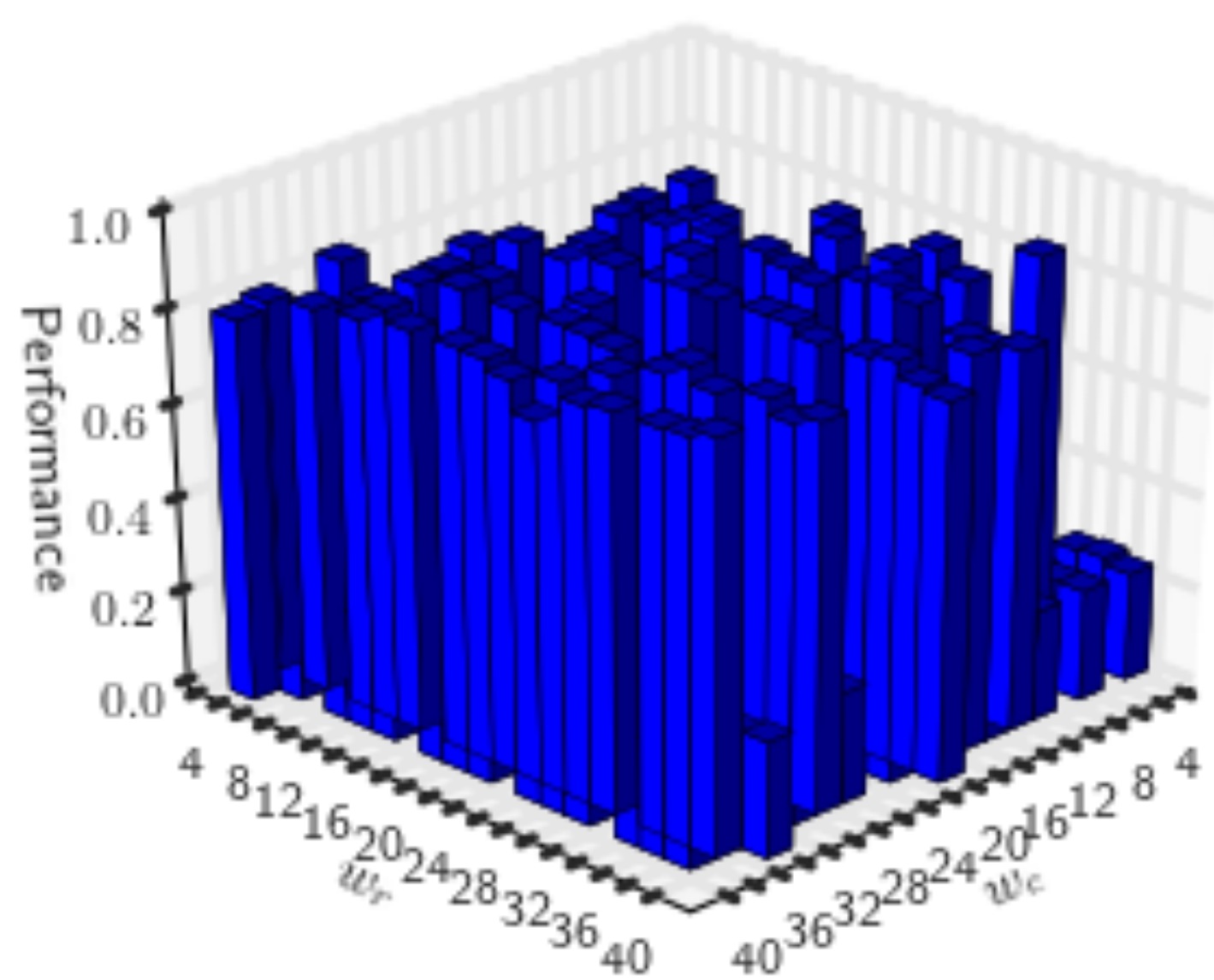
OpenCL workgroup size:
Controls decomposition of
threads.
Is a 2D parameter (rows x cols).

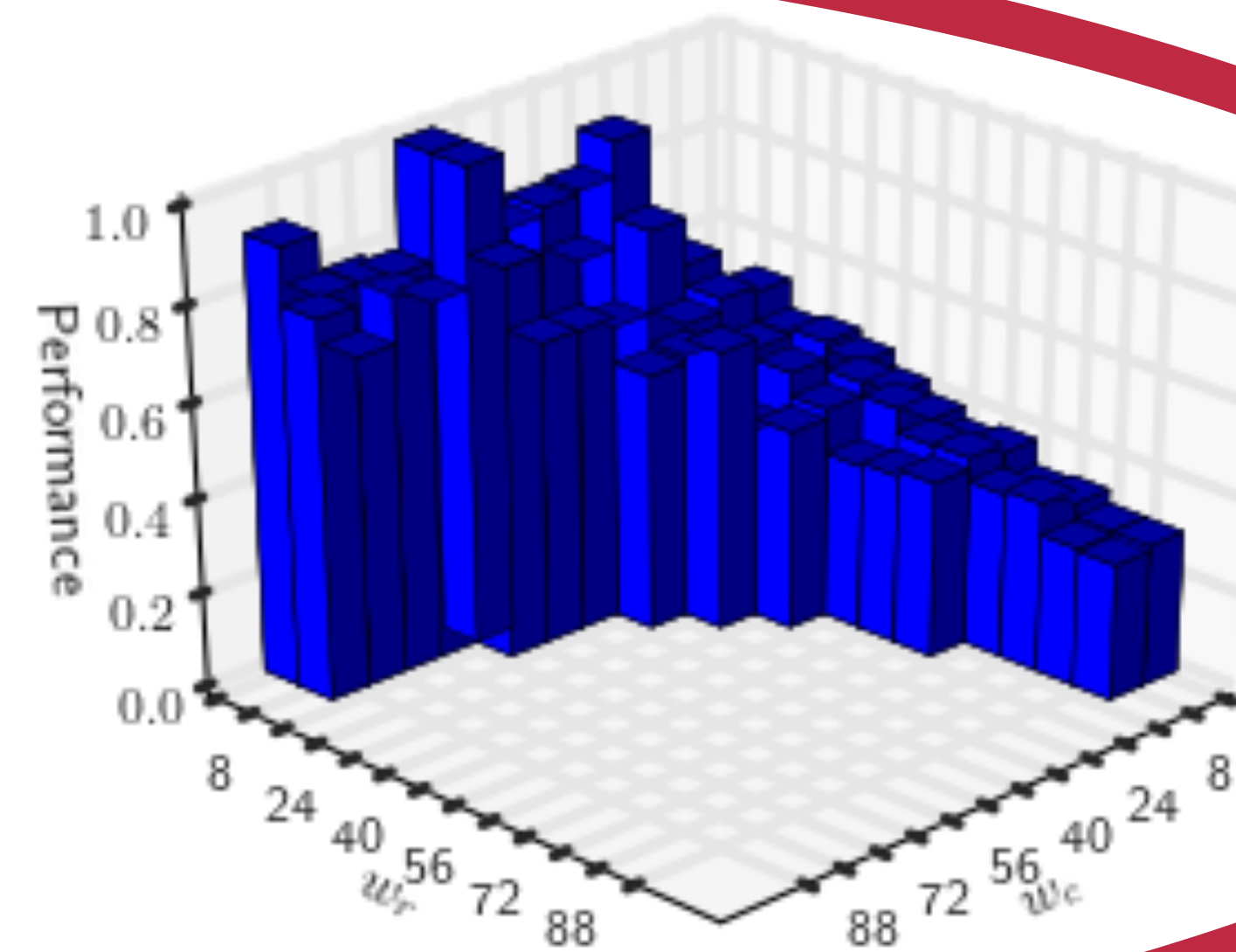
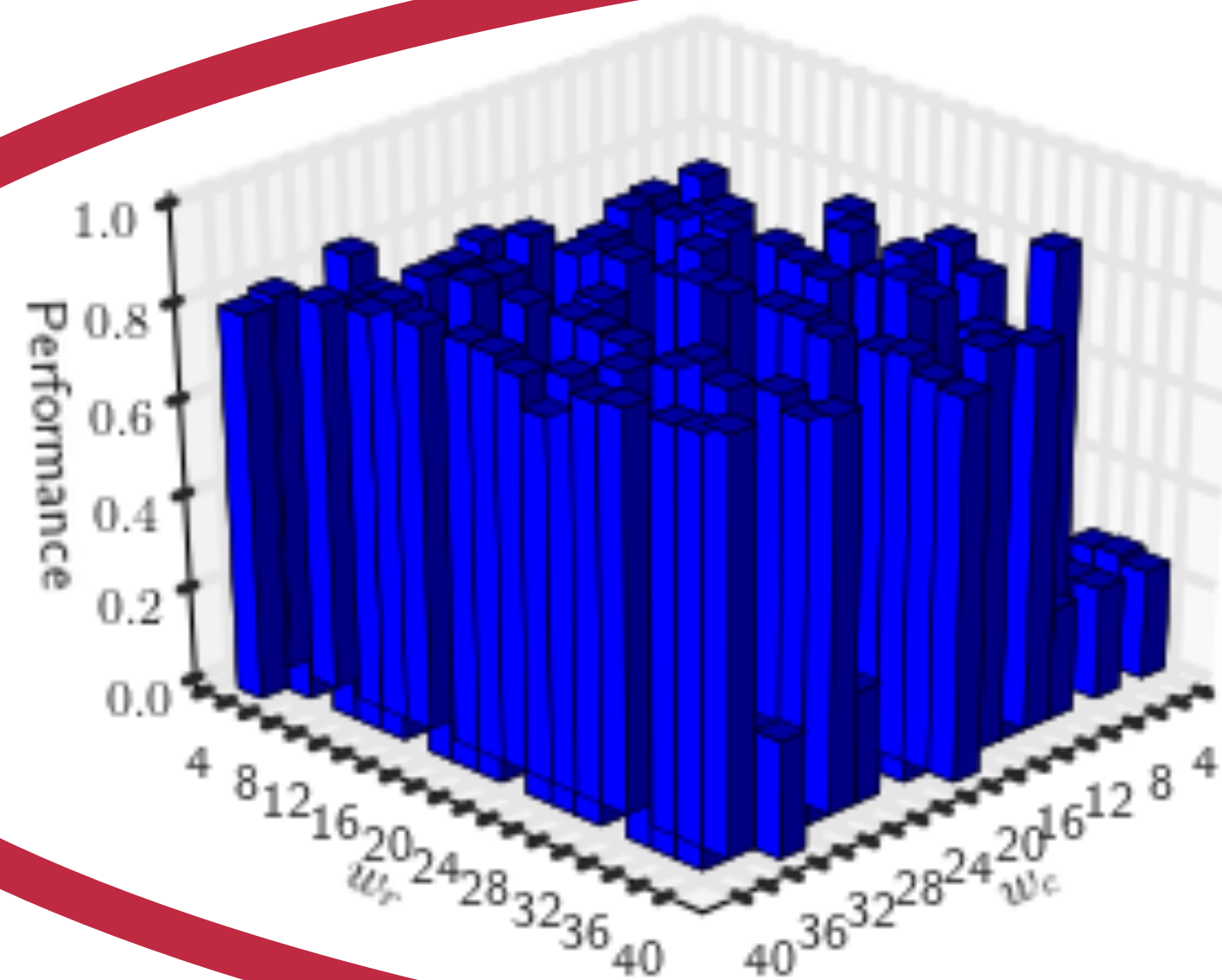
OpenCL workgroup size:
Controls decomposition of
threads.
Is a 2D parameter (rows x cols).
Critical to performance.

OpenCL workgroup size:

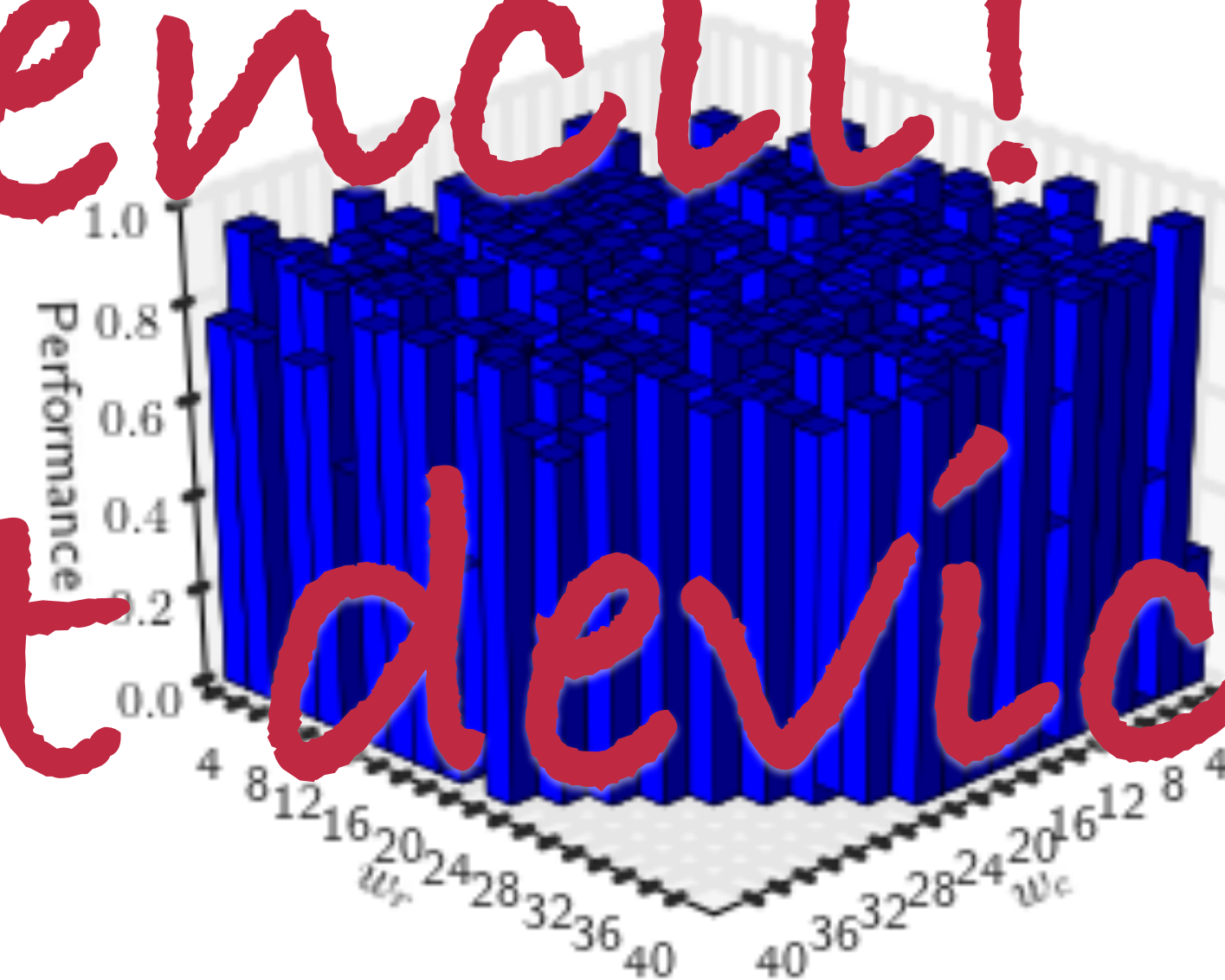
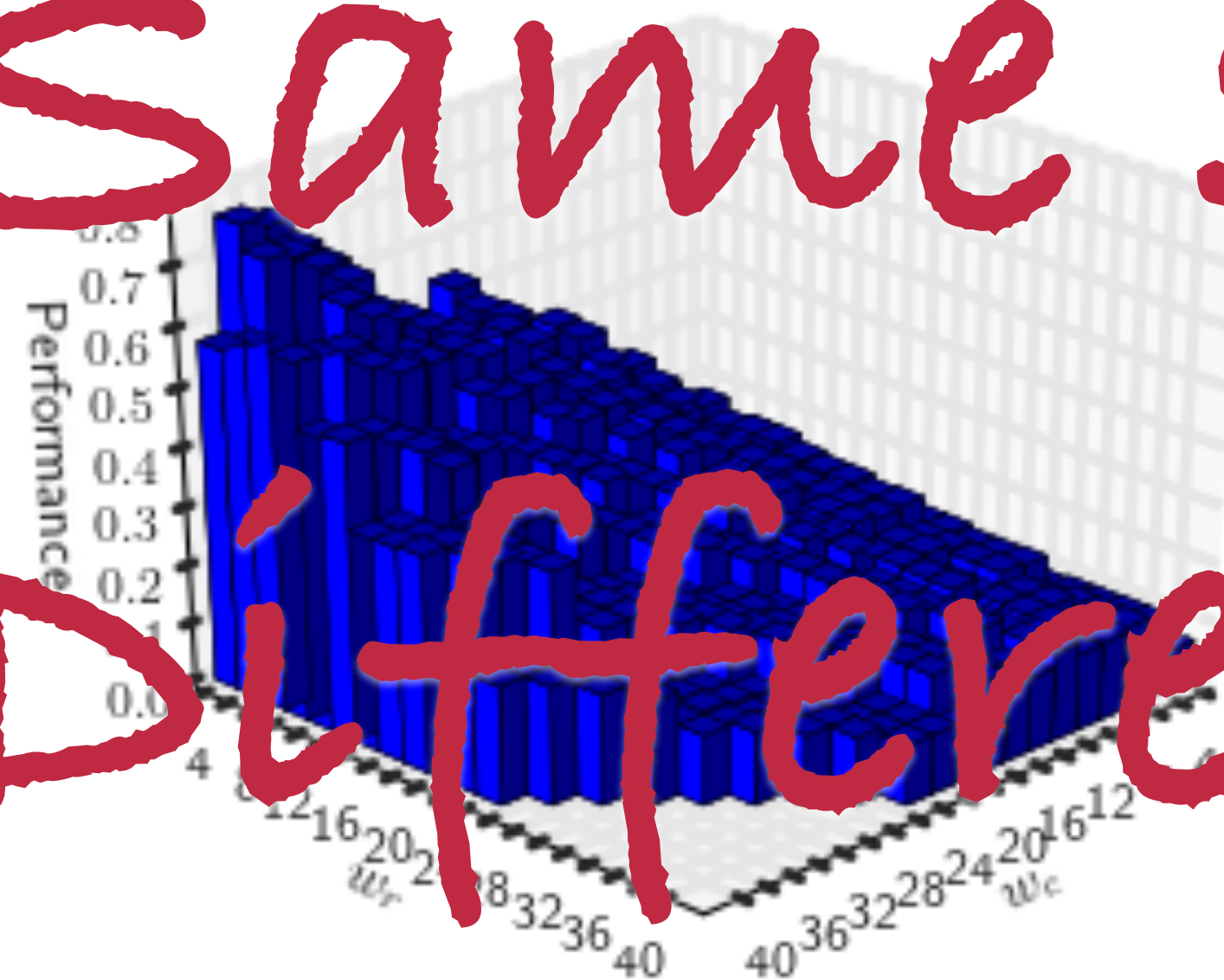


Examples

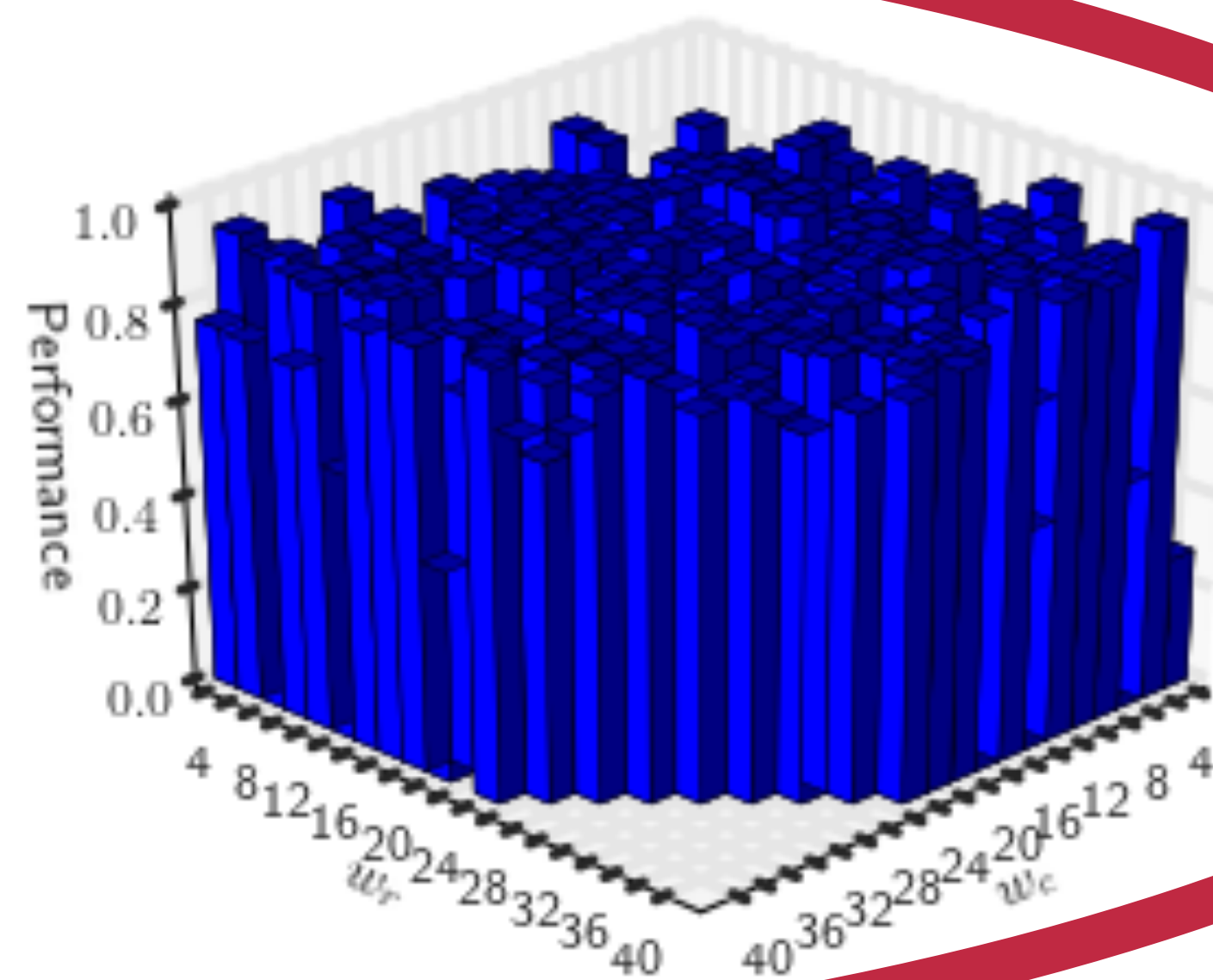
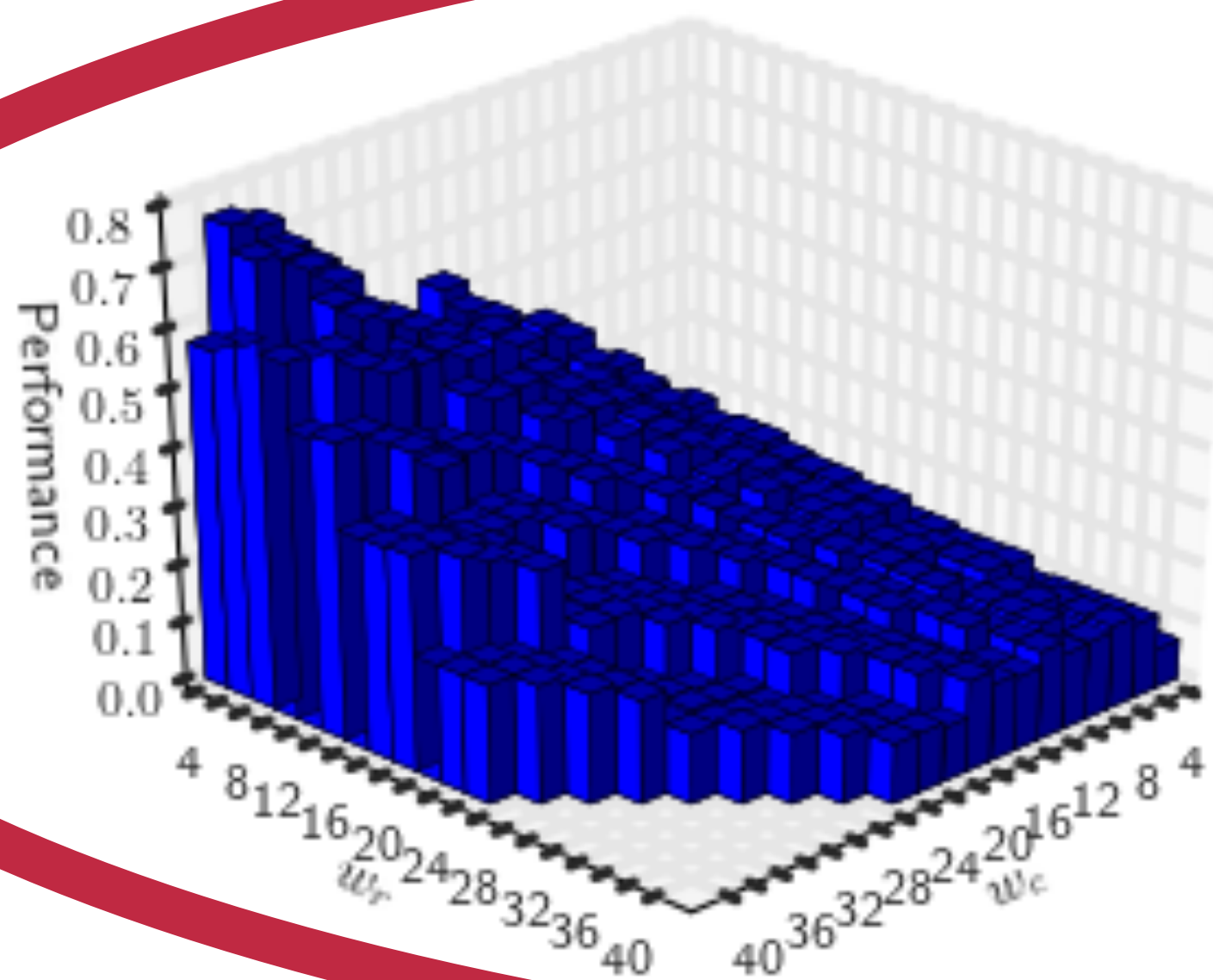


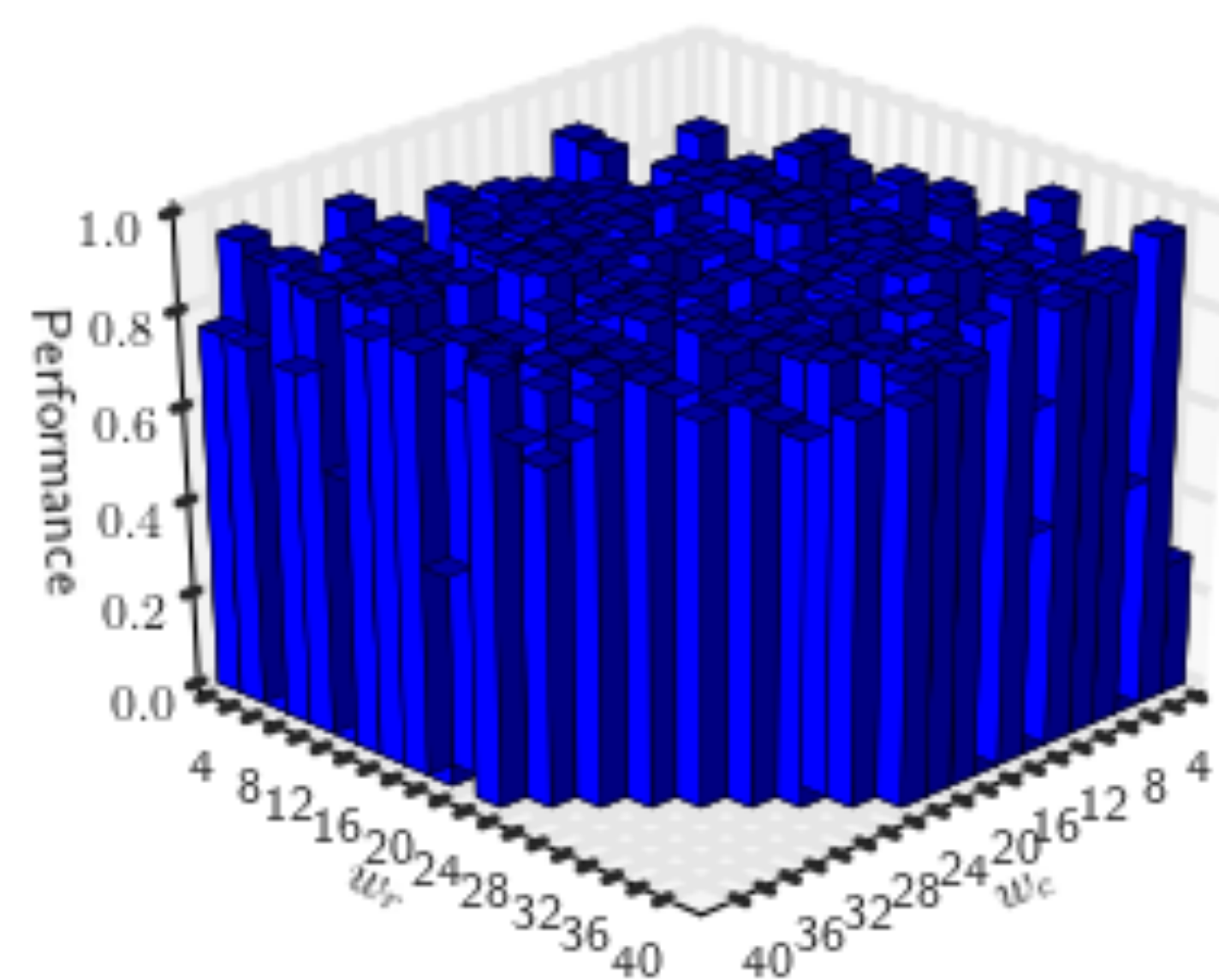
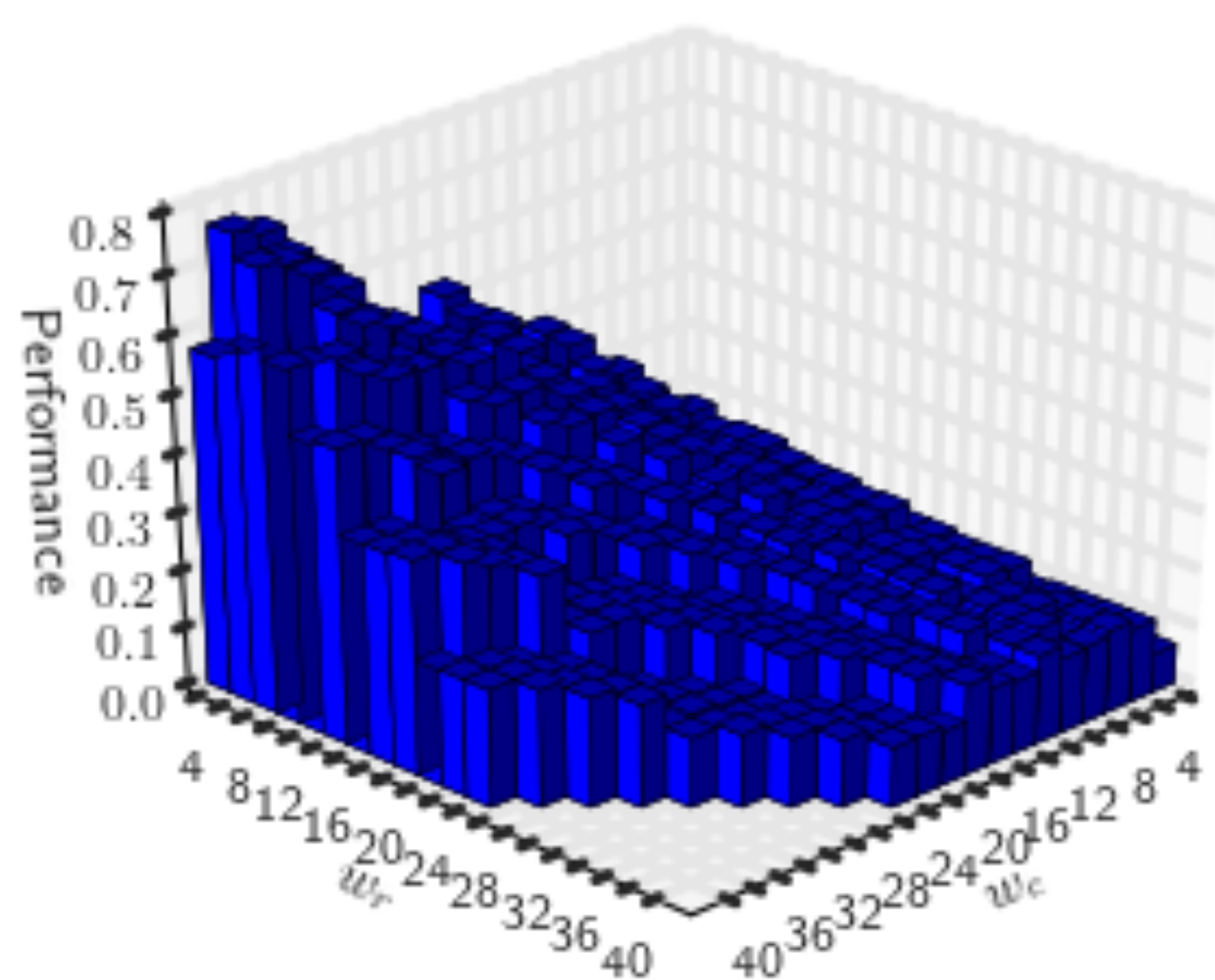
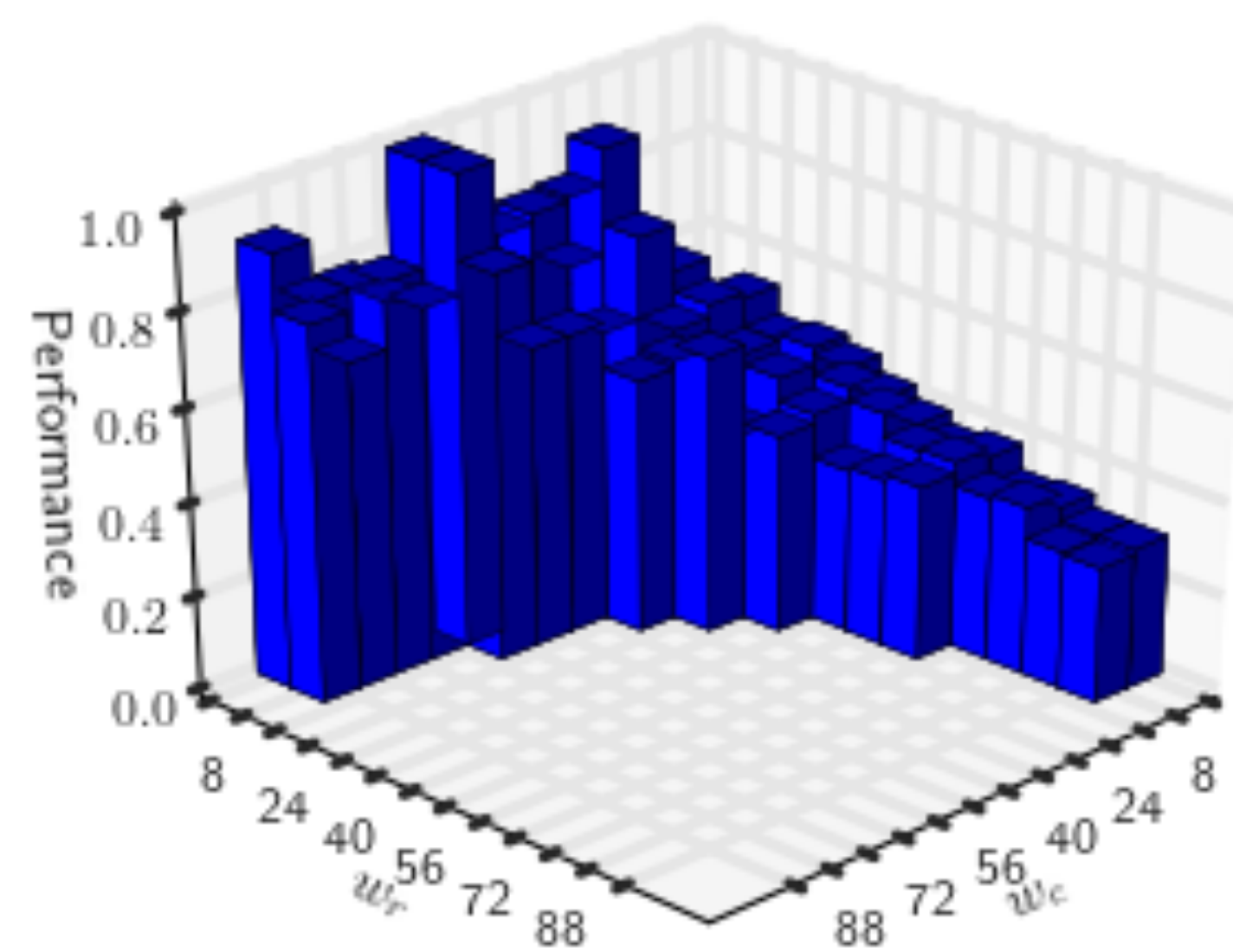
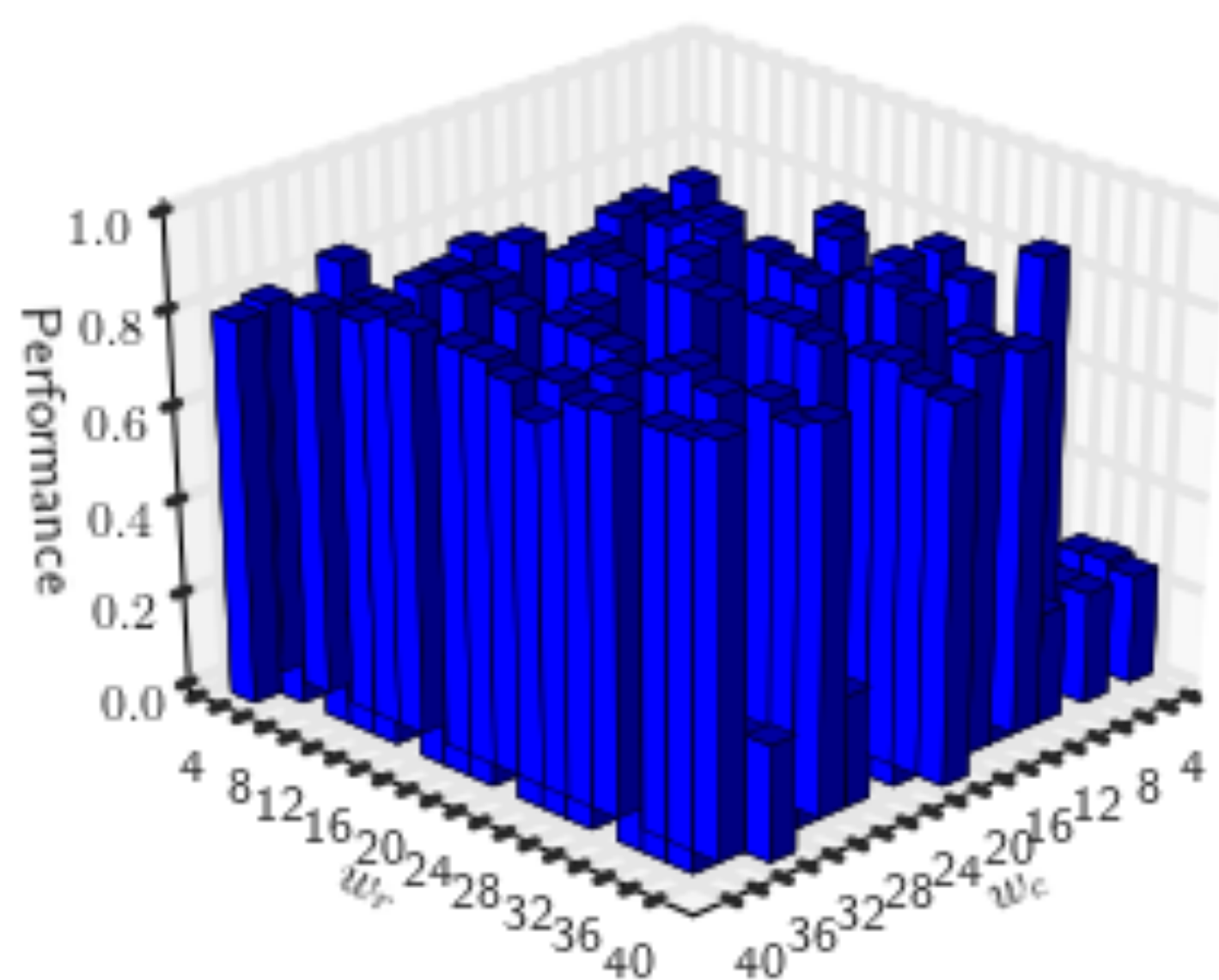


Same stencil!
Different device!



Same device!
Different stencil!





choosing workgroup size
depends on:

1. device
2. program
3. dataset

**Let's automate
this!**

Approach 1

Set a workgroup size

Execute and time program

Set a workgroup size

Execute and time program

Set a workgroup size

Execute and time program

Set a workgroup size

Execute and time program

Set a workgroup size

Execute and time program

Set a workgroup size

Execute and time program

Set a workgroup size

Execute and time program

Set a workgroup size

Execute and time program

Set a workgroup size

Execute and time program

Set a workgroup size

Execute and time program

Set a workgroup size

Execute and time program

Set a workgroup size

Execute and time program

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

Set a workgroup size

Execute and time program

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

*(iterative
compilation)*

BAD!

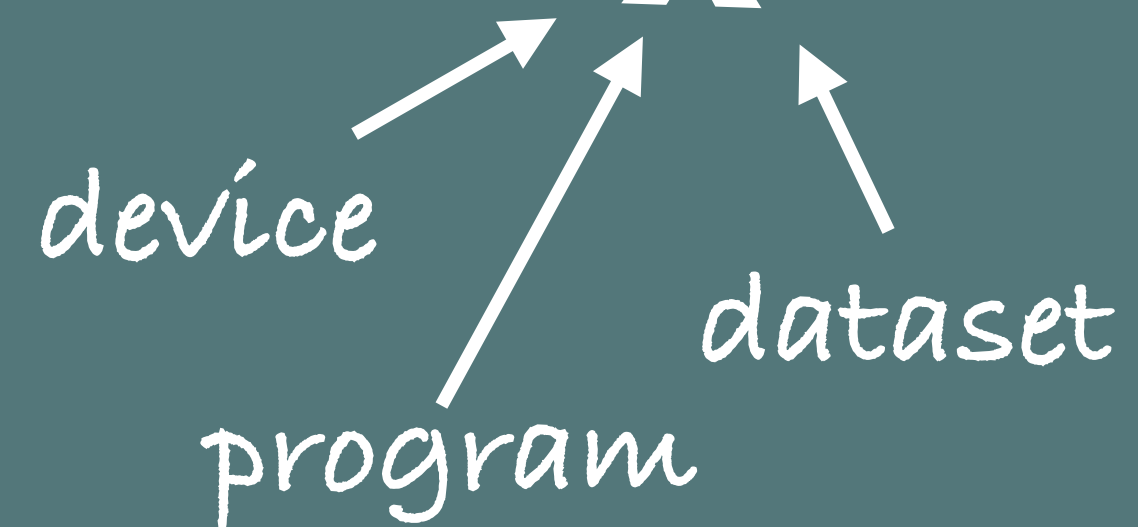
Takes a looooong time

BAD!

Takes a looooong time

BAD!

Must be repeated for every new “x”



Approach 2

Set a workgroup size

Execute and time program

Set a workgroup size

Execute and time program

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

Set a workgroup size

Execute and time program

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

1 data point



Collect **data points**

Extract “features”

Train machine learning classifier

Extract “features”

Input to classifier

BETTER!

Can make *predictions* on unseen “x”



```
graph LR; device --> x["x"]; program --> x; dataset --> x;
```

A diagram with three labels: 'device', 'program', and 'dataset'. Arrows point from each label to a large 'x' in quotes. 'device' is at the top left, 'program' is at the bottom left, and 'dataset' is at the bottom right.

BETTER!

Can make *predictions* on unseen “x”



A diagram with three handwritten labels: 'device', 'program', and 'dataset'. Three arrows originate from these labels and point towards a large 'x' that is part of the text 'unseen “x”'. The 'device' arrow points from the top-left, the 'program' arrow from the bottom-left, and the 'dataset' arrow from the bottom-right.

BETTER!

Still takes a looooong time

Can make *predictions* on unseen “x”



A diagram with three labels: 'device', 'program', and 'dataset'. Arrows point from each label to a large 'x' in quotes. 'device' has an arrow pointing up and to the right. 'program' has an arrow pointing up and to the left. 'dataset' has an arrow pointing up and to the right.

BETTER!

Still takes a loooooong time

Requires a lot of code

Our wish list:

- 1. Reduce training costs**
- 2. Reduce implementation costs**
- 3. Minimise runtime overheads**

Our Approach ...

OmniTune

1. Allows *collaborative* performance tuning

Reduce training costs ✓

2. Provides re-usable *implementations*

Reduce implementation costs ✓

3. Provides lightweight runtime interface

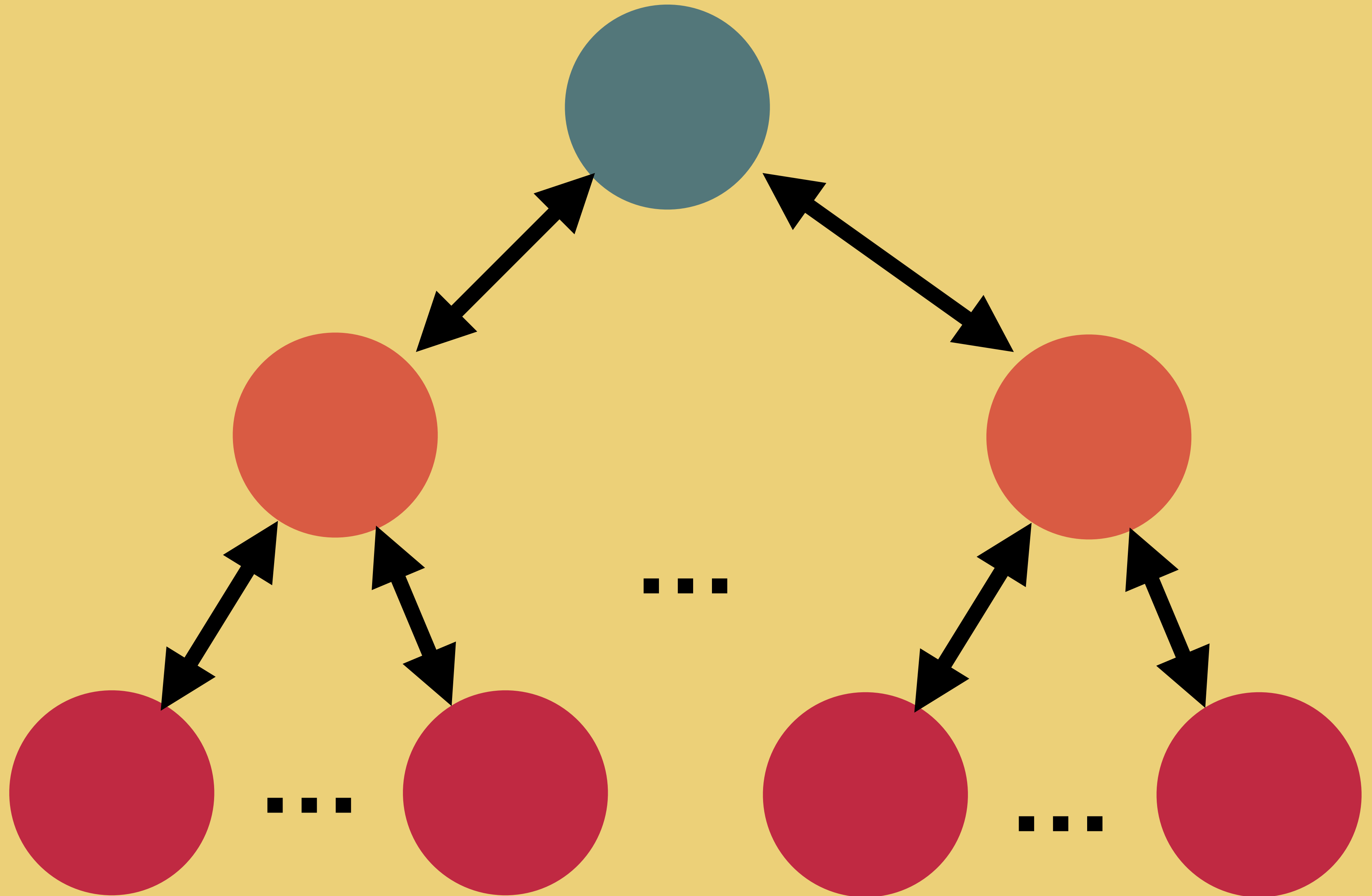
Minimise runtime overheads ✓

How does it work?

Remote

Servers

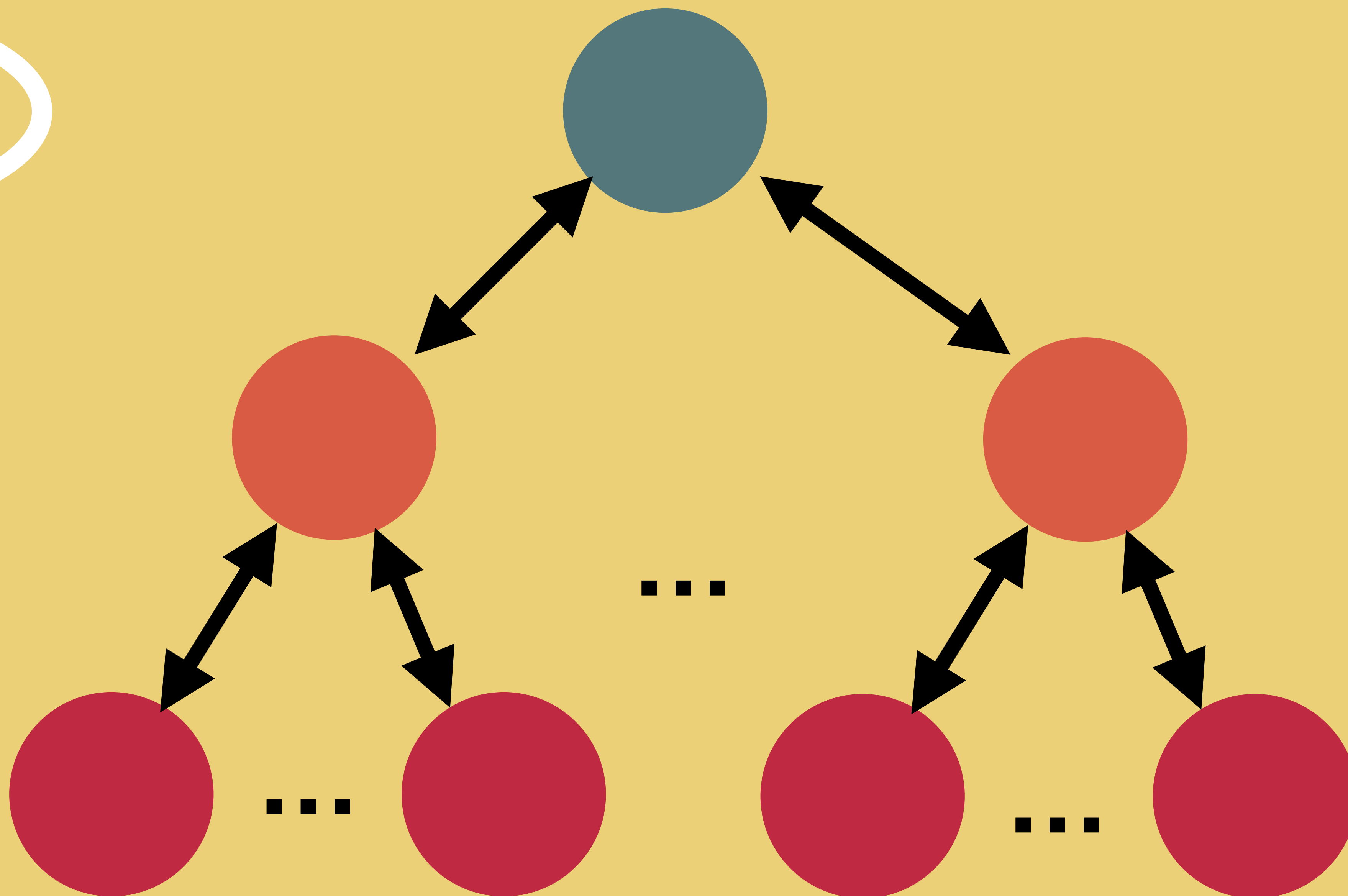
Clients



Remote

Servers

Clients



Remote

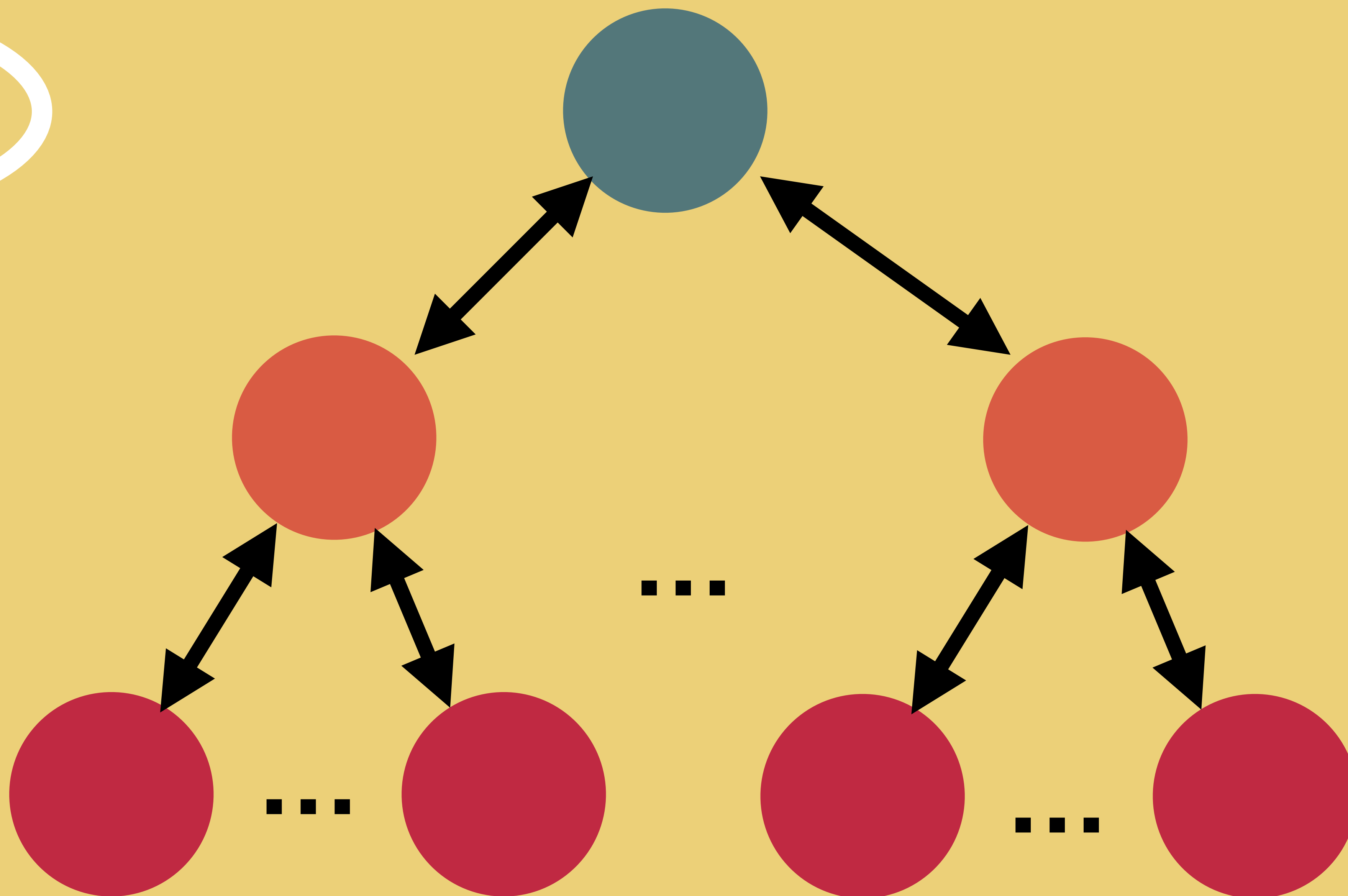
Book-keeper

**Manages and
stores training
data**

Remote

Servers

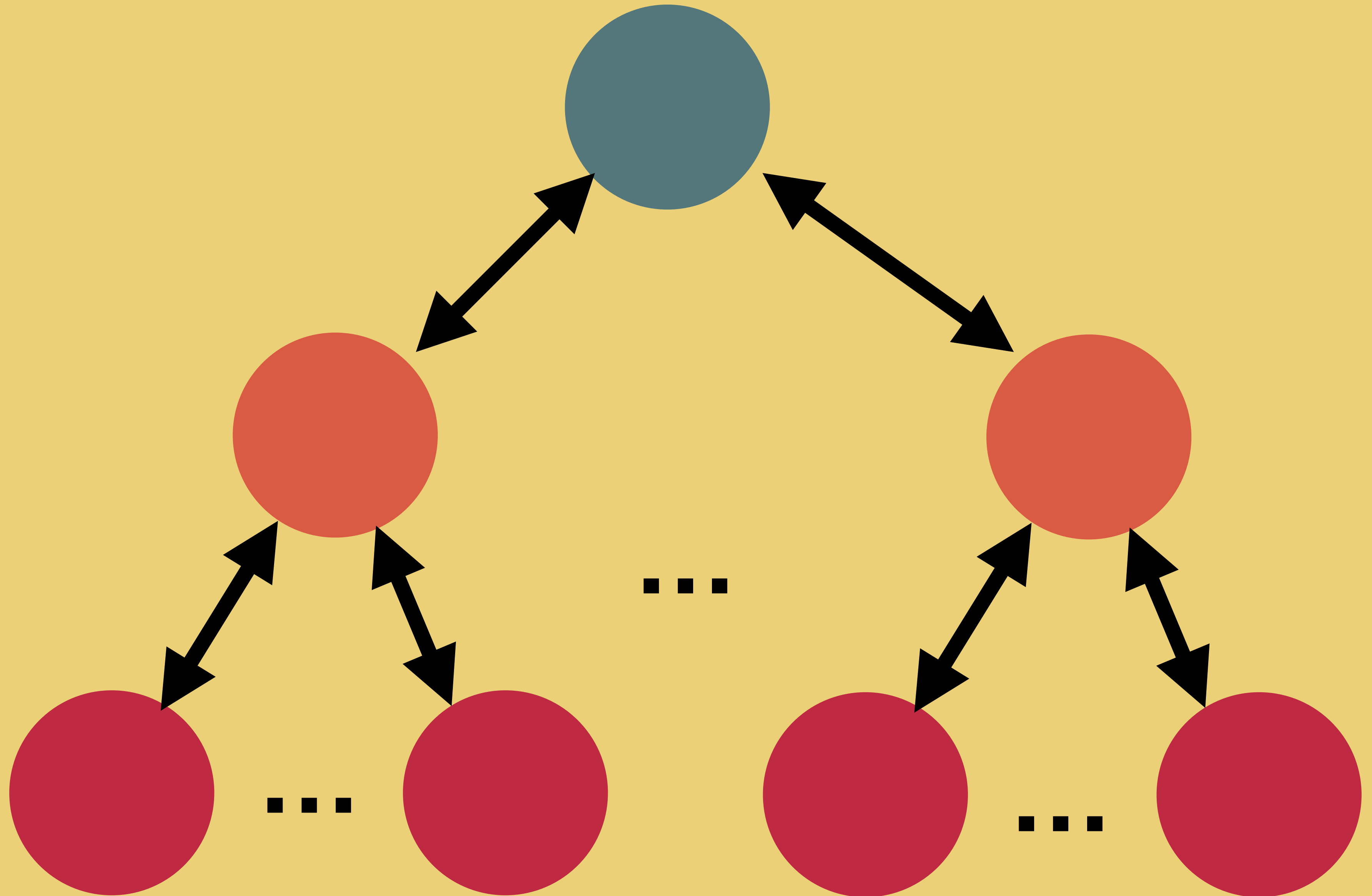
Clients



Remote

Servers

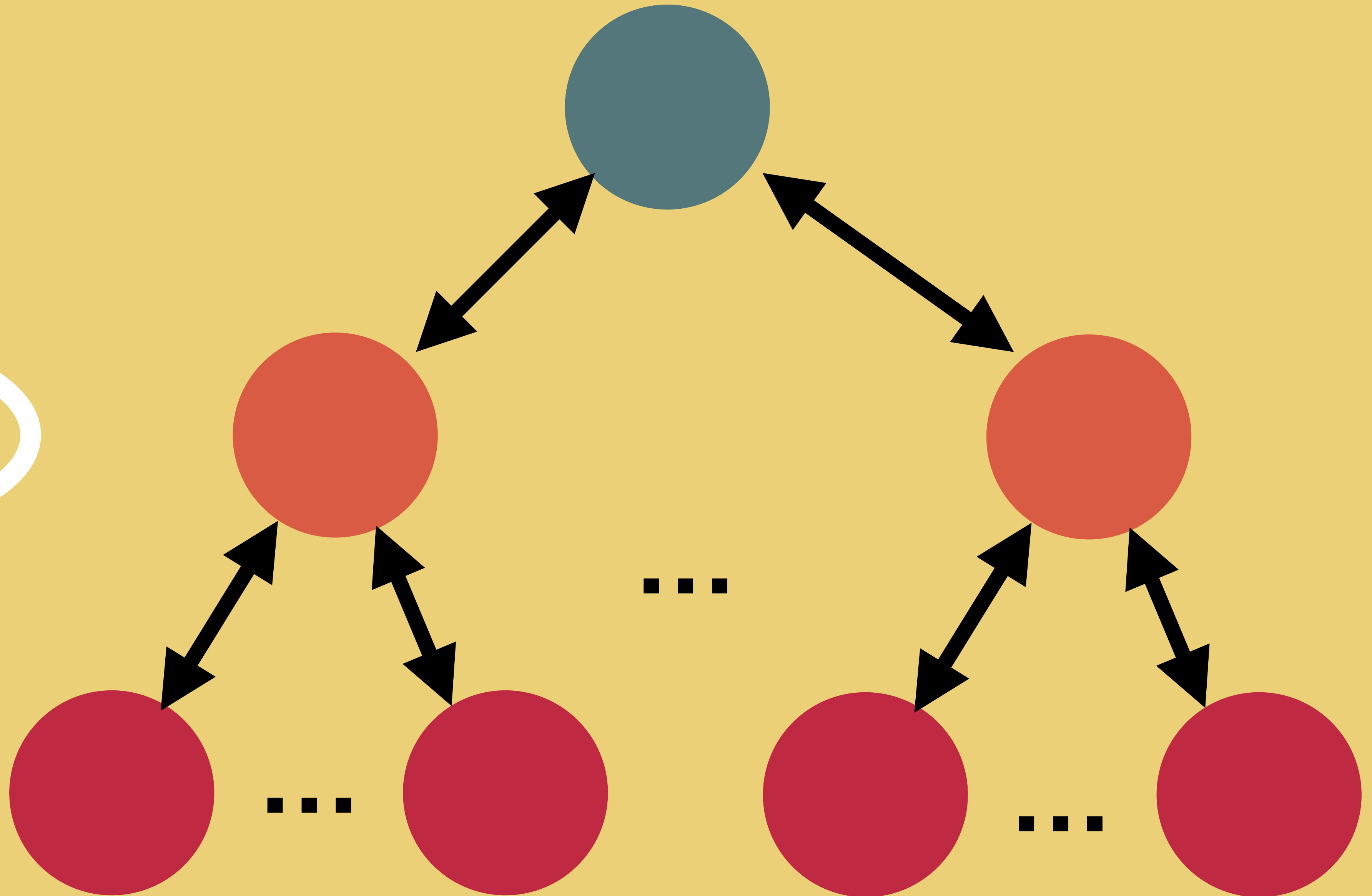
Clients



Remote

Servers

Clients



**Autotuning
engine**



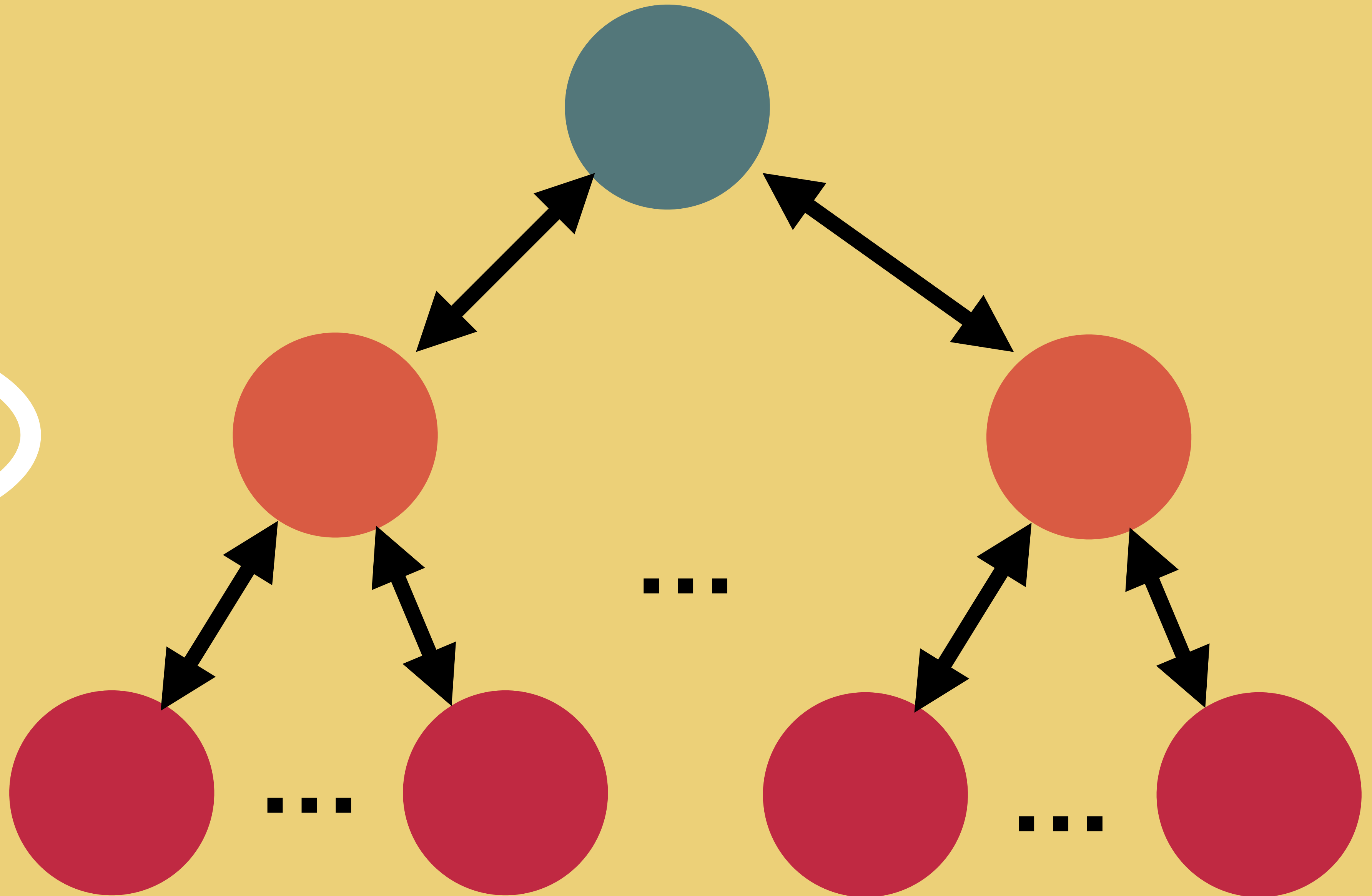
Servers

**Performs
machine learning**

Remote

Servers

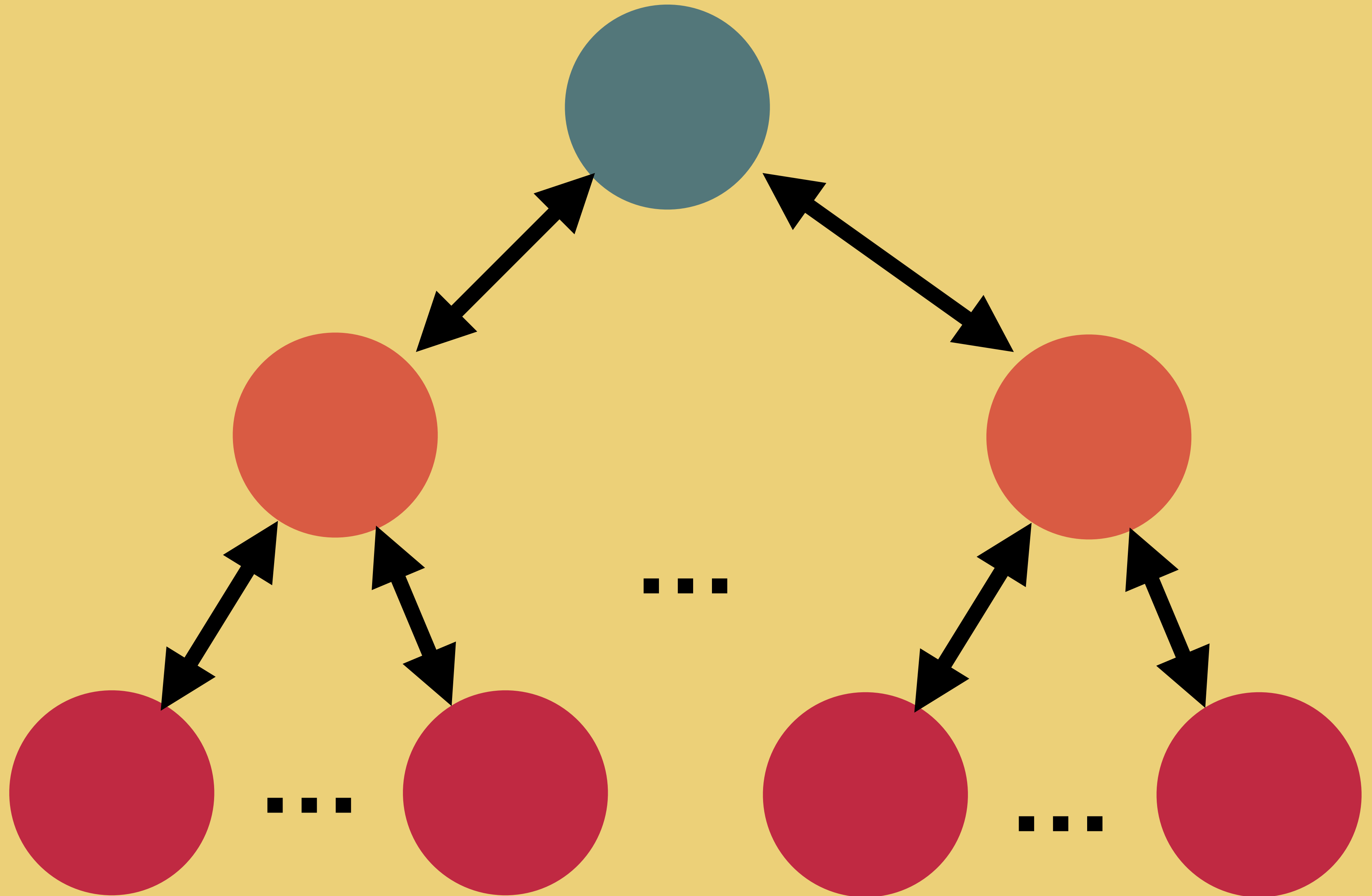
Clients



Remote

Servers

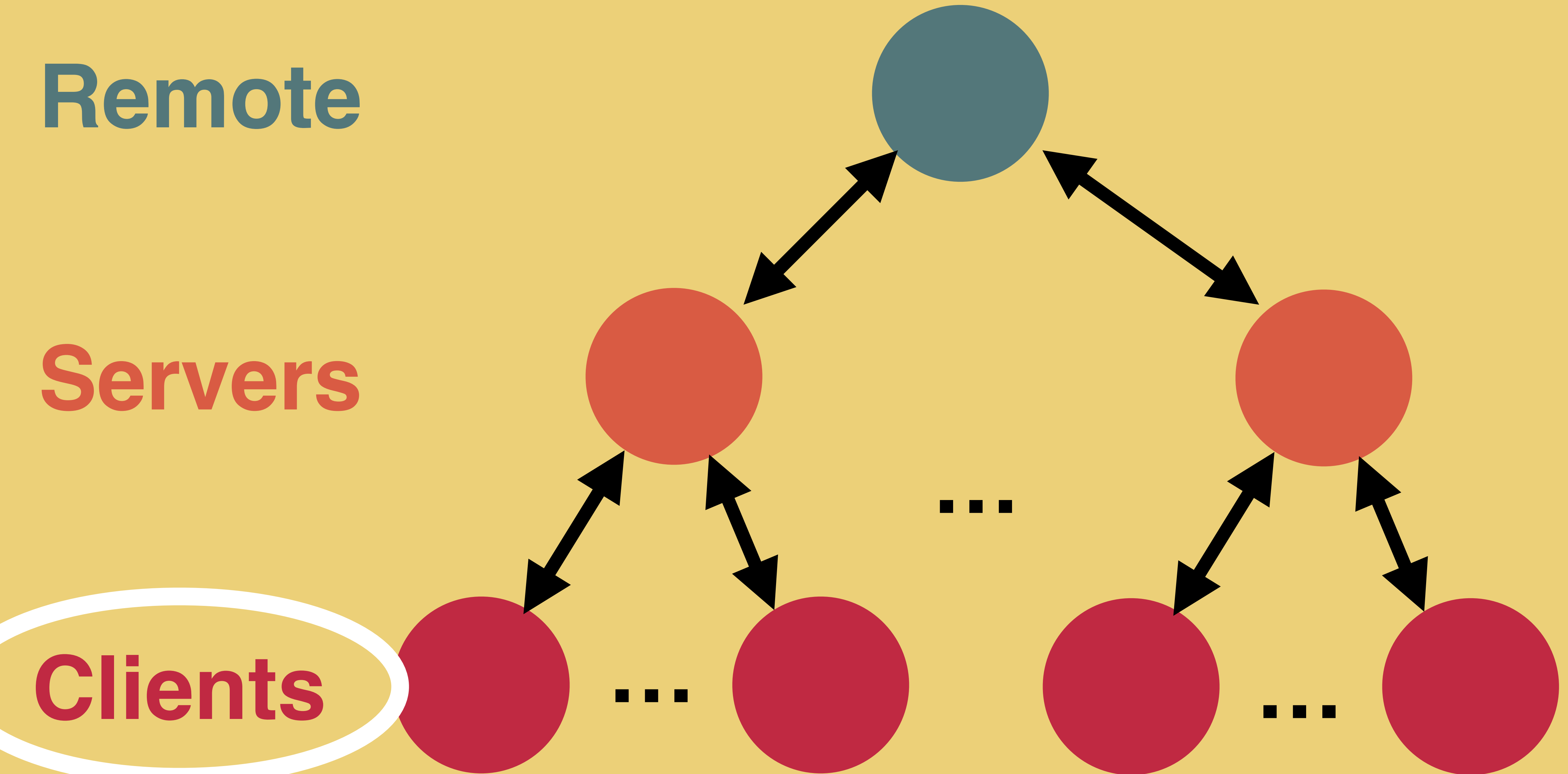
Clients



Remote

Servers

Clients



**Target
applications**

**Programs we
want to tune**

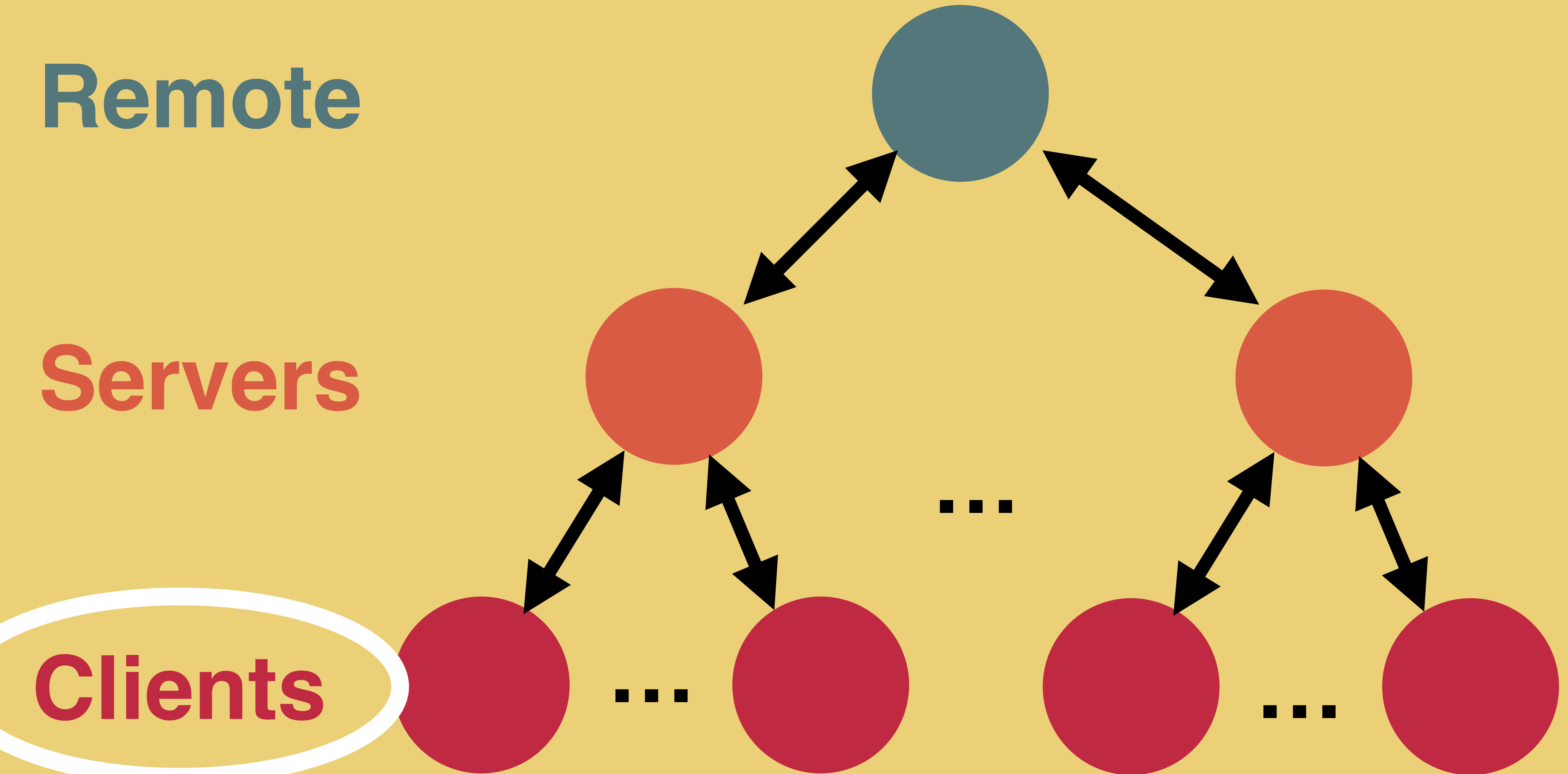


Clients

Remote

Servers

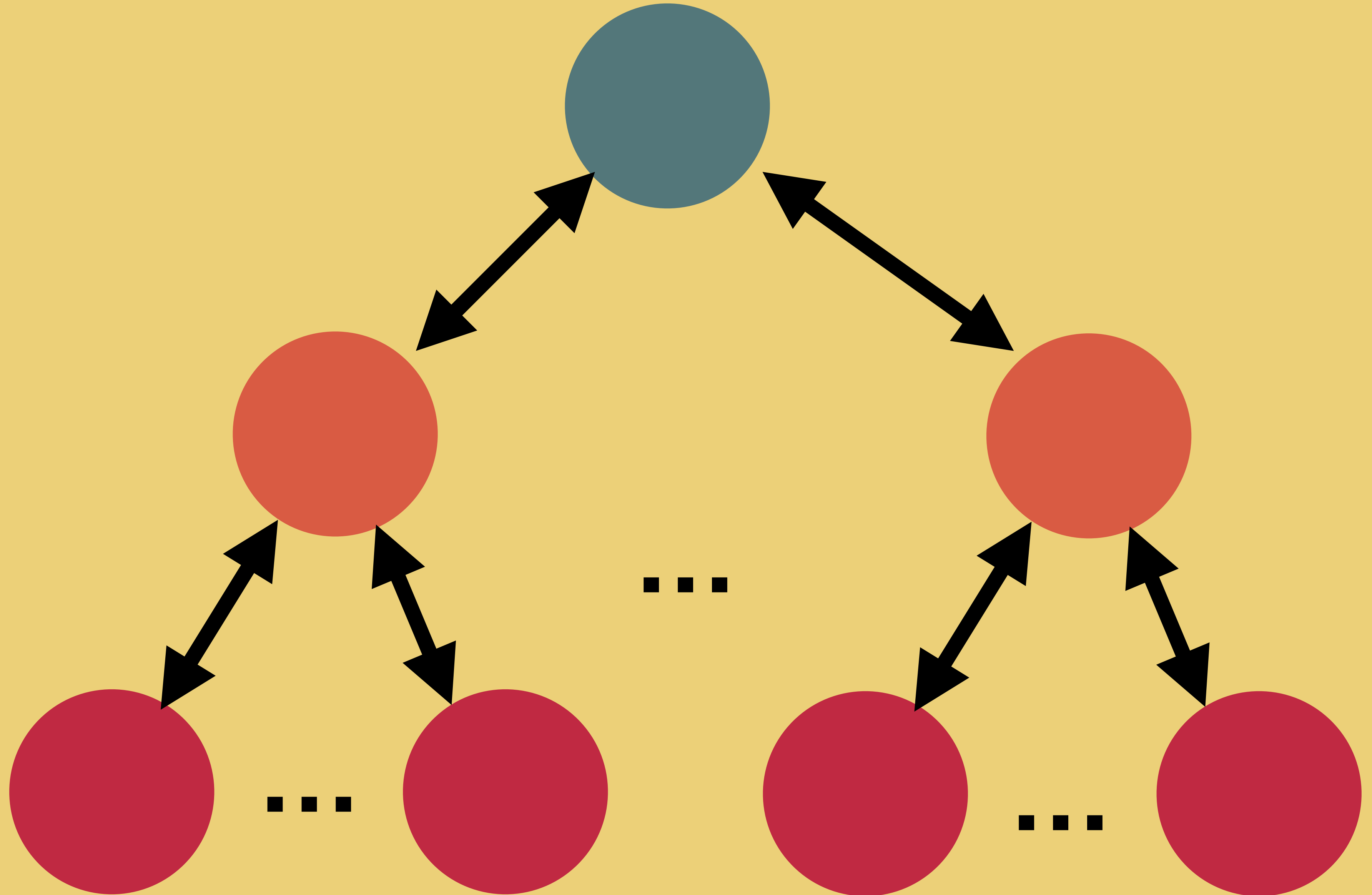
Clients



Remote

Servers

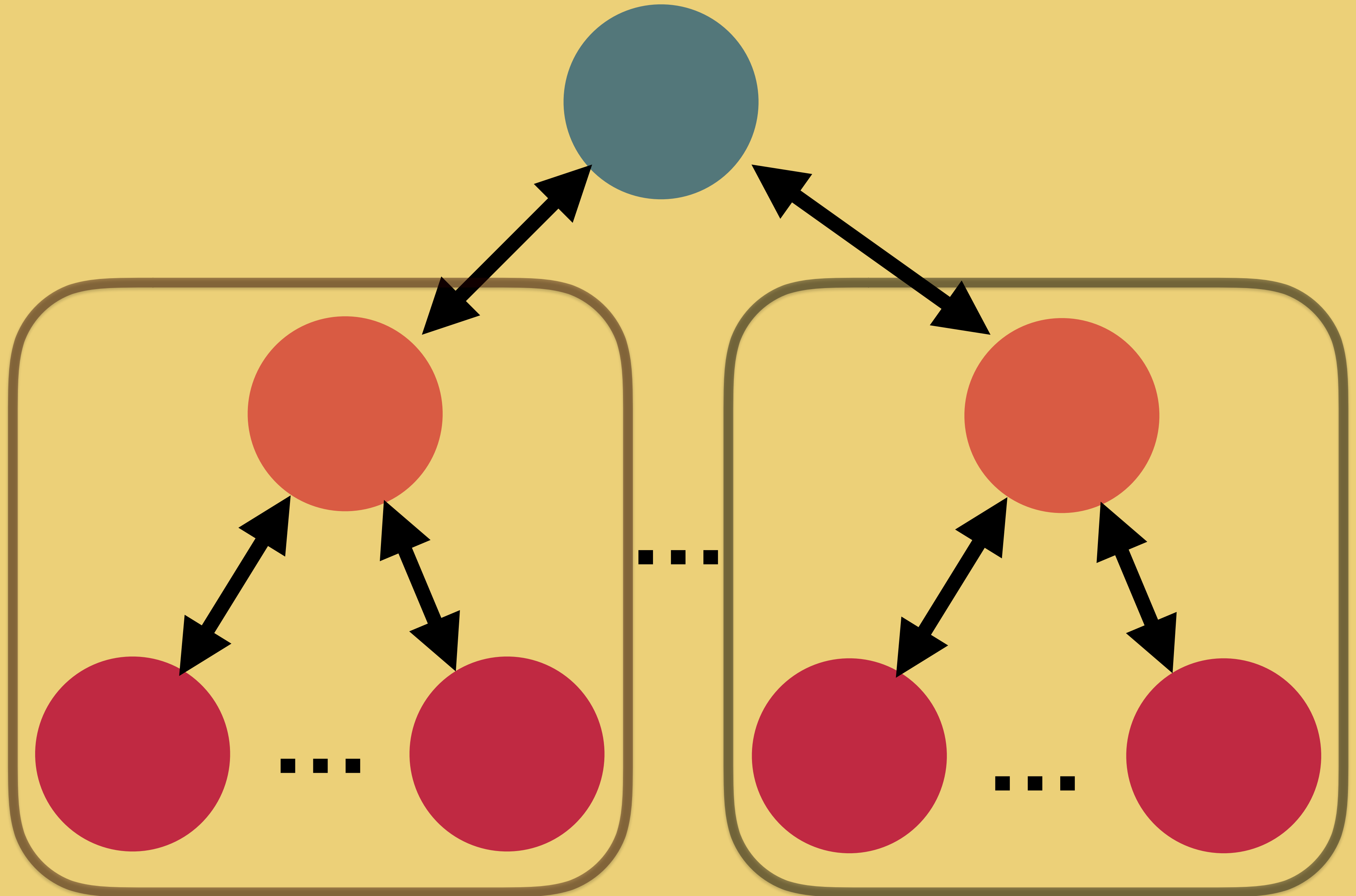
Clients



Remote

Servers

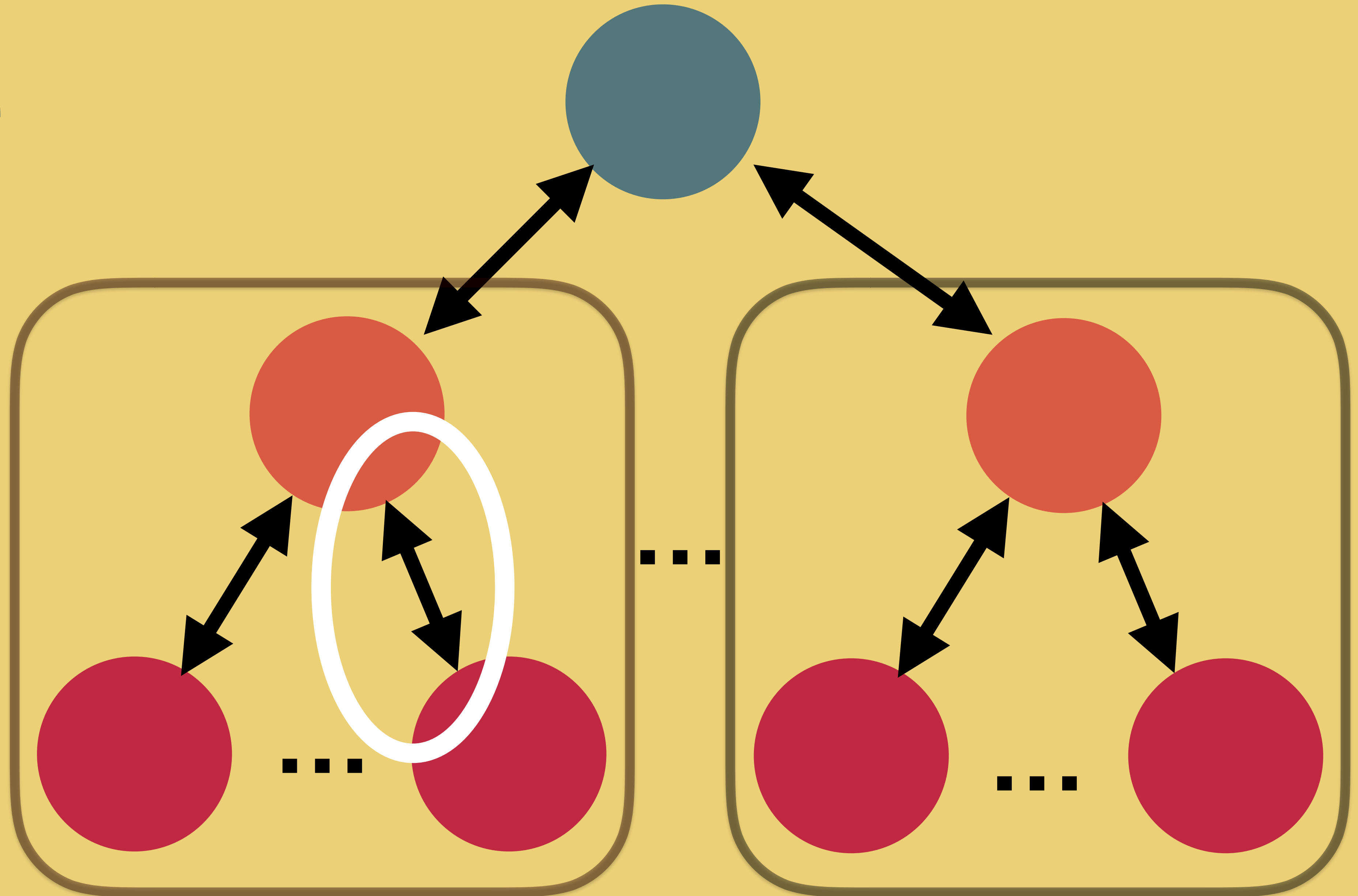
Clients



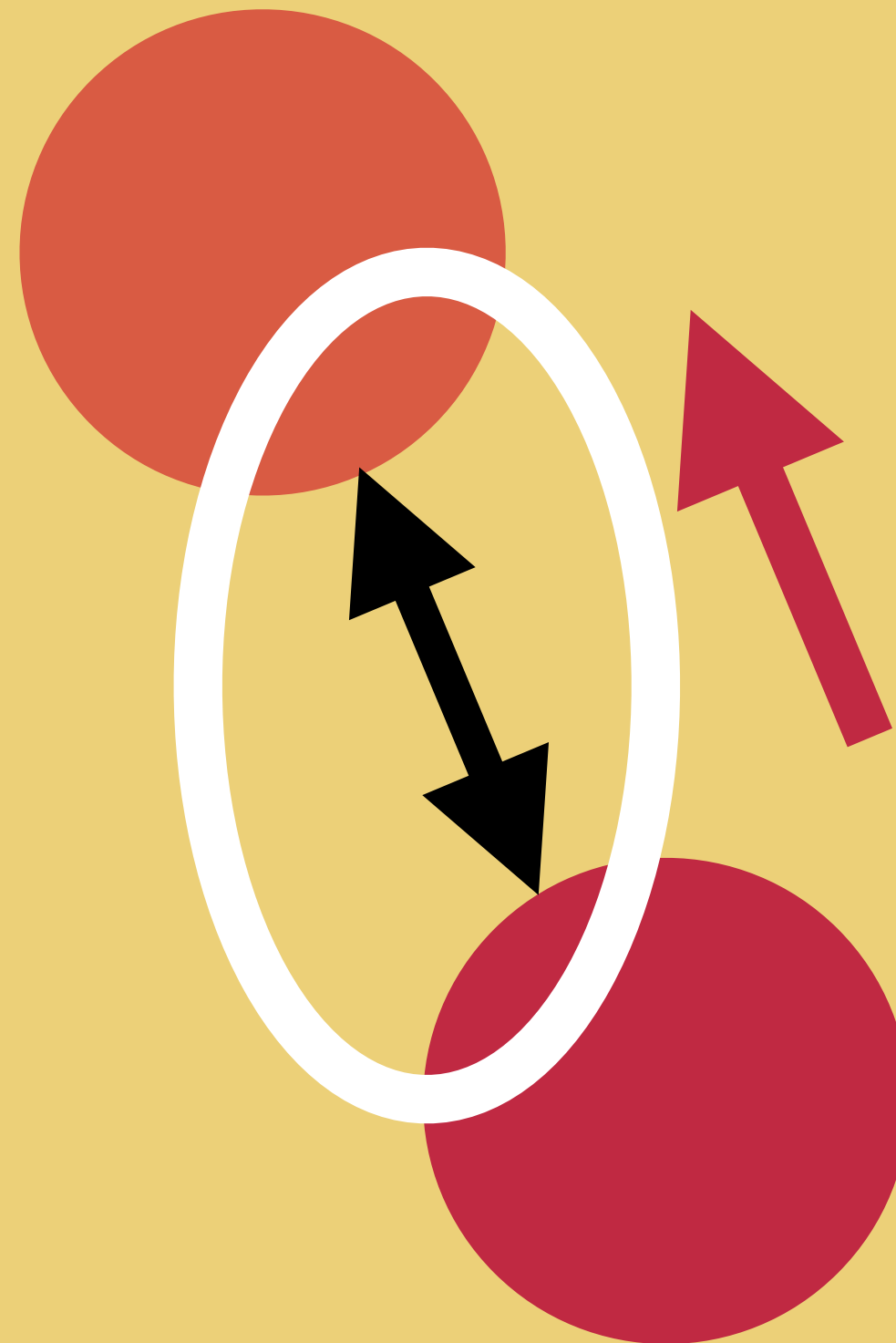
Remote

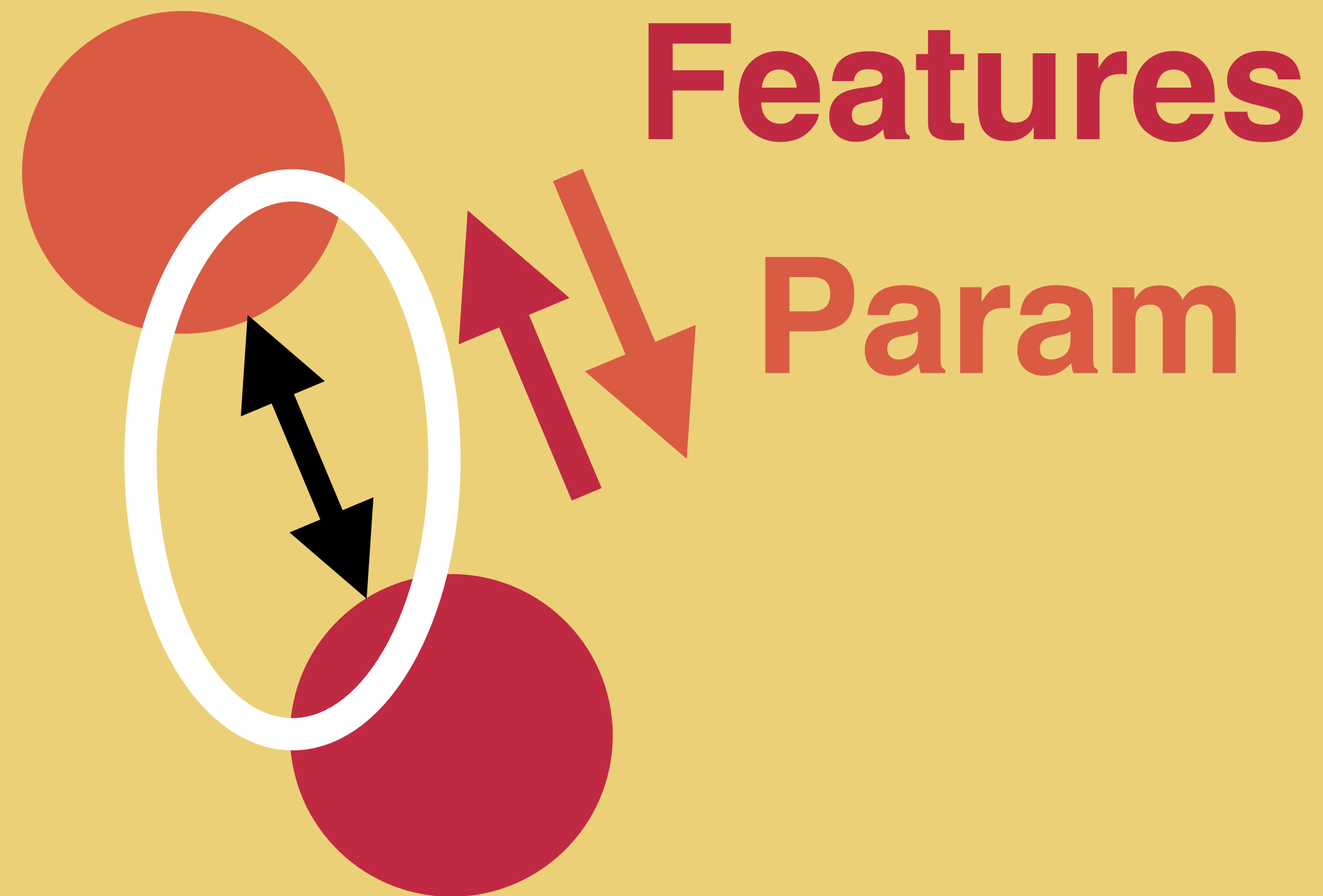
Servers

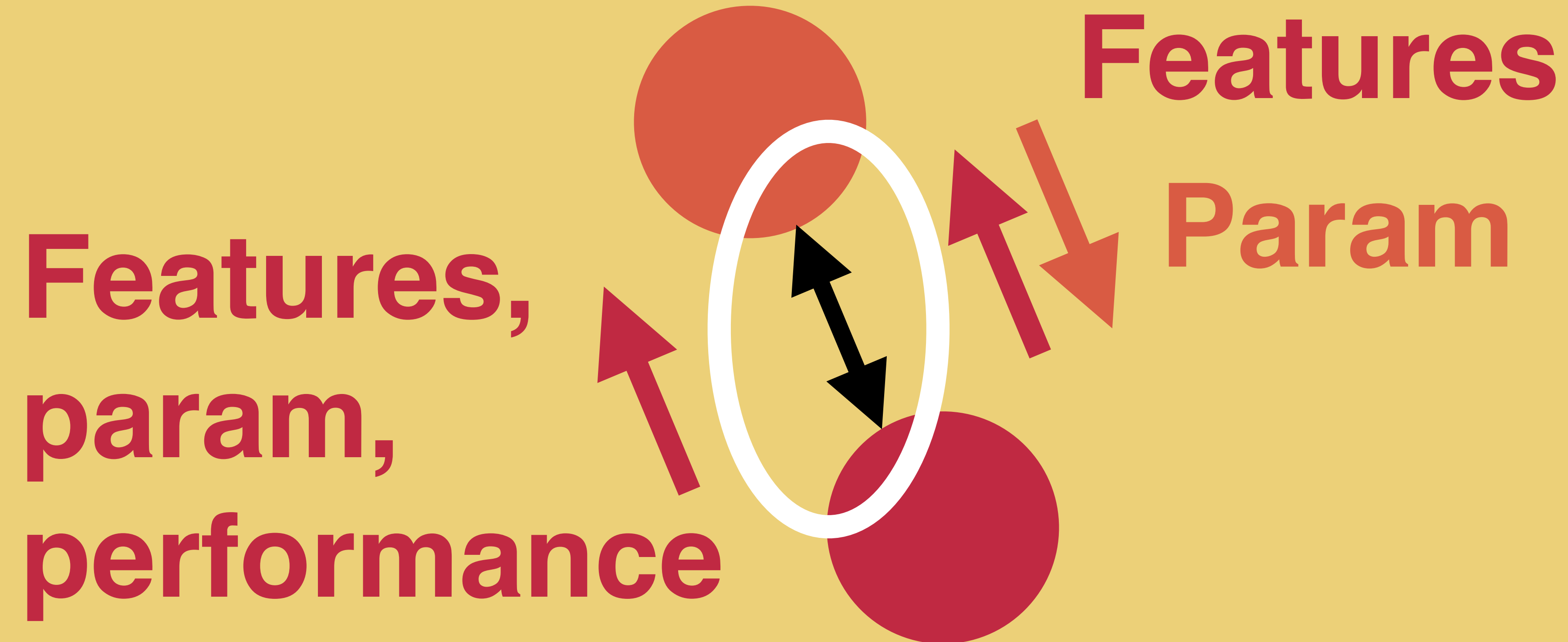
Clients



Features



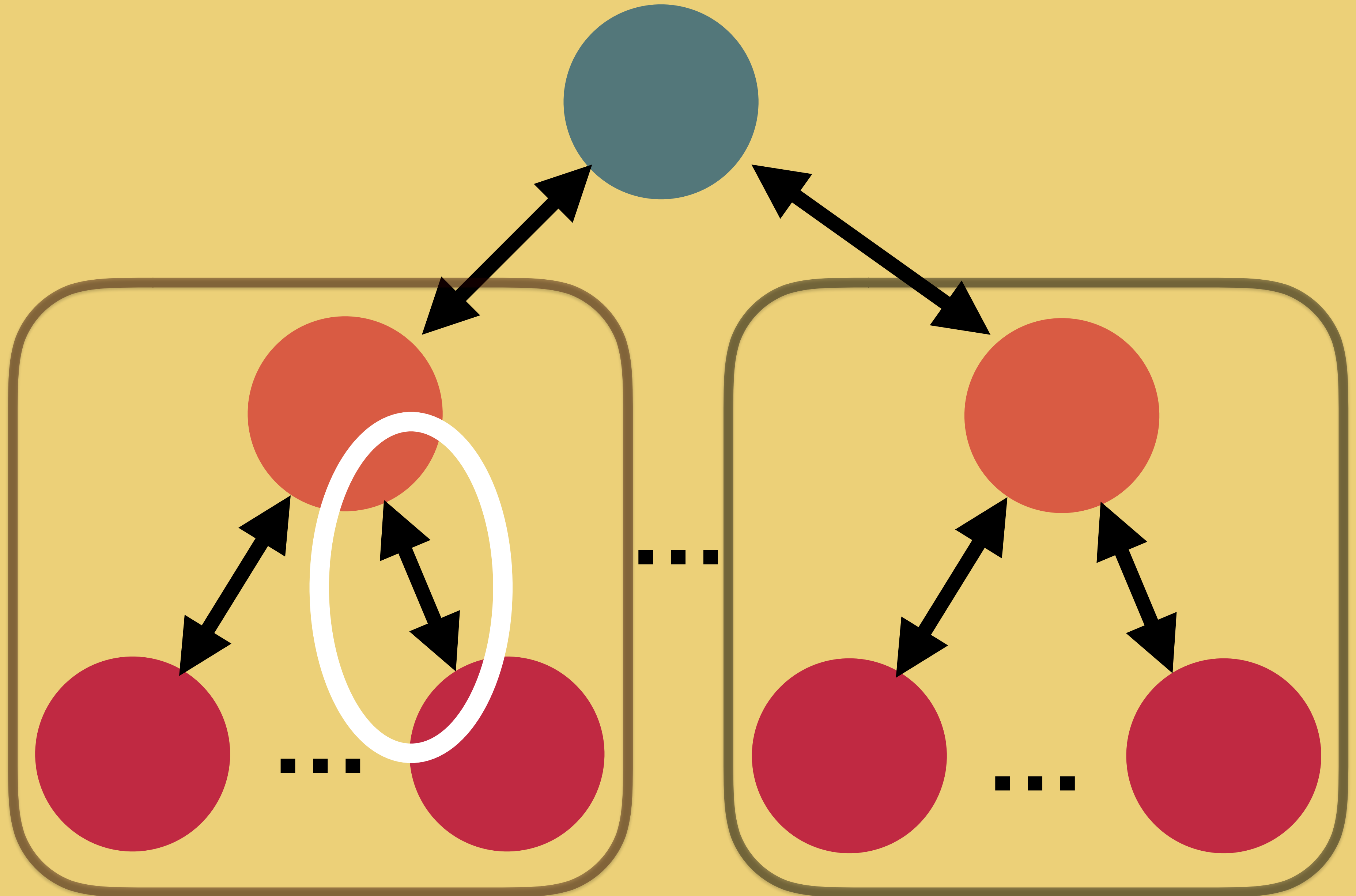




Remote

Servers

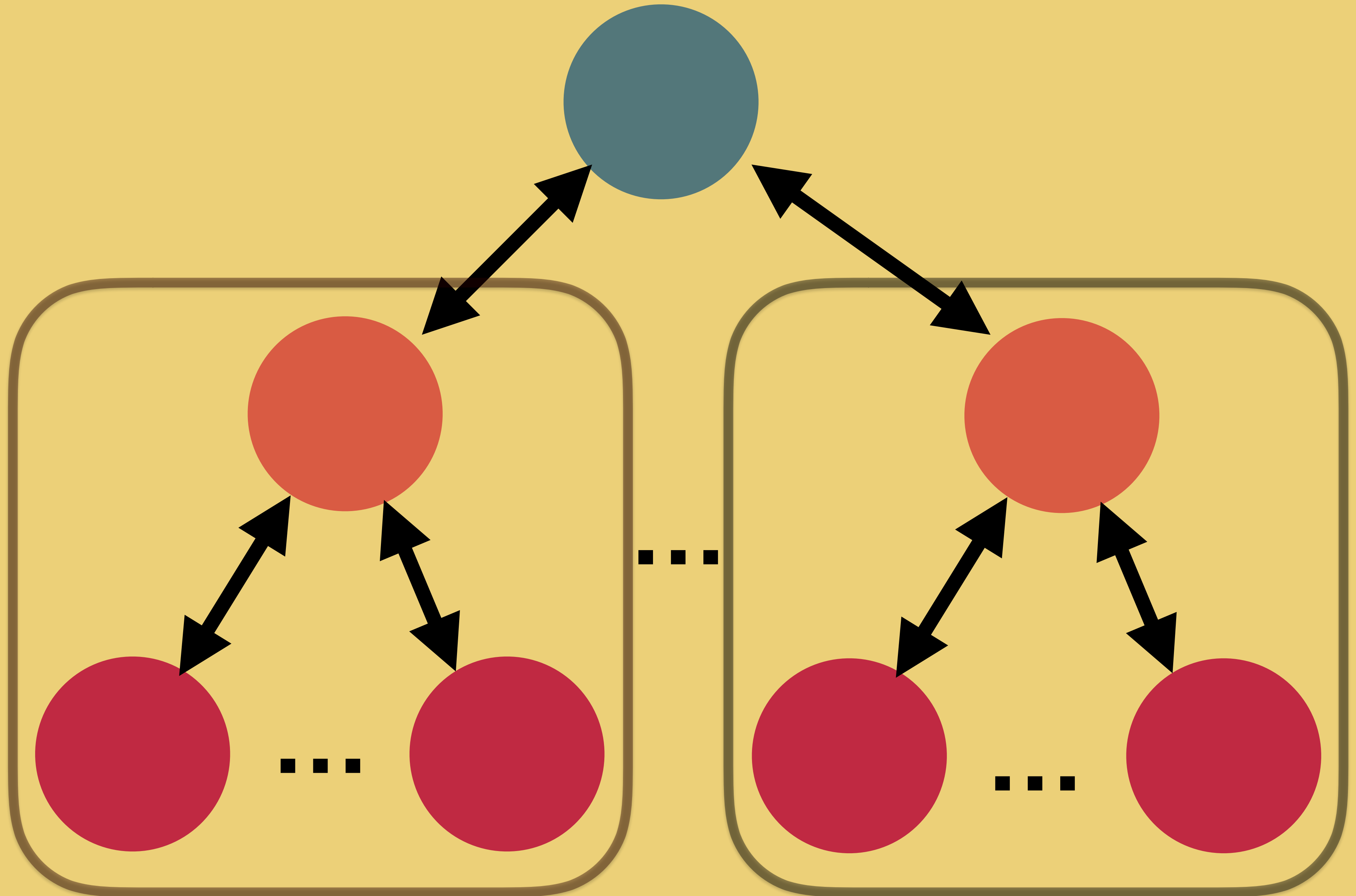
Clients



Remote

Servers

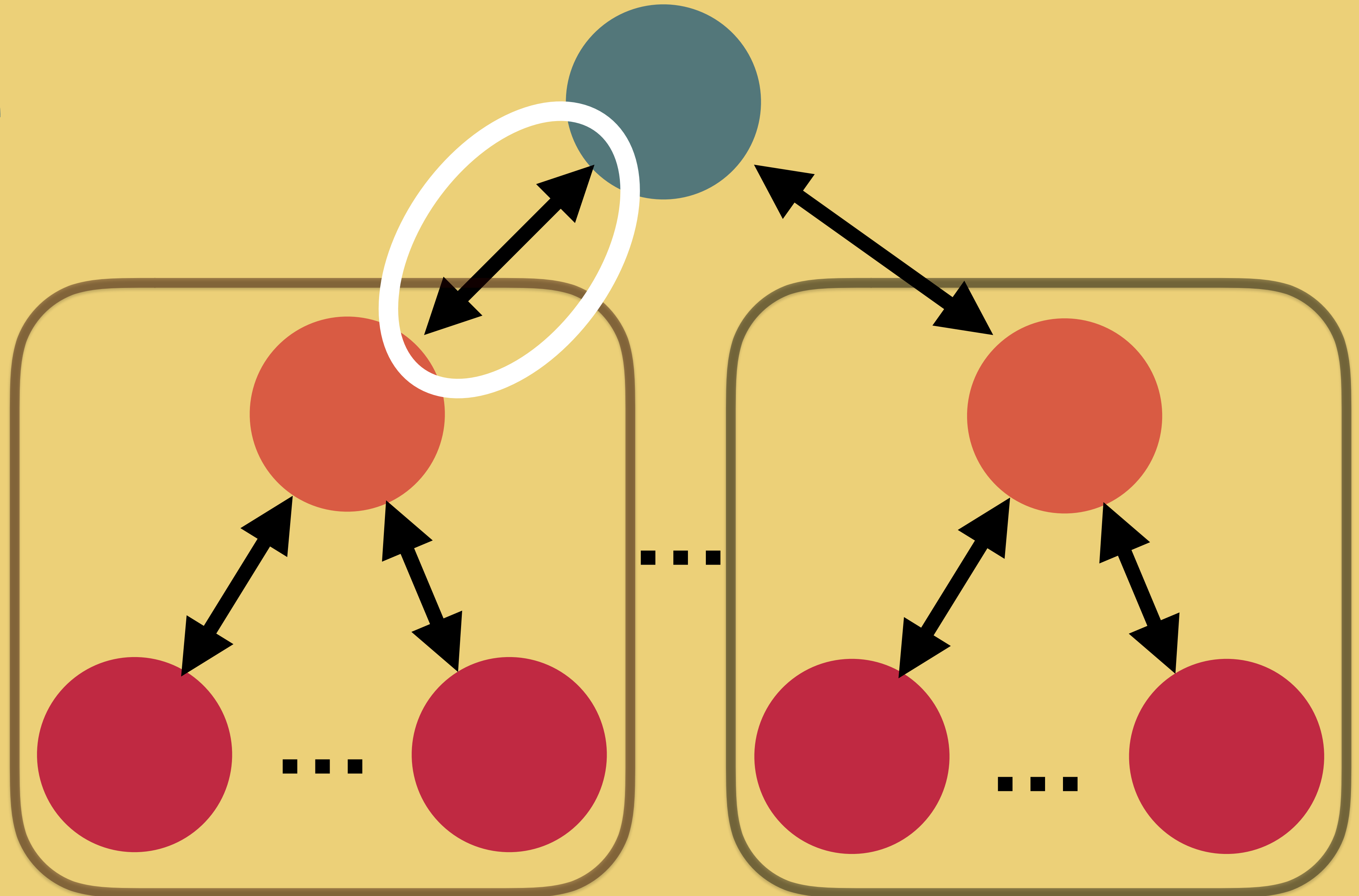
Clients



Remote

Servers

Clients



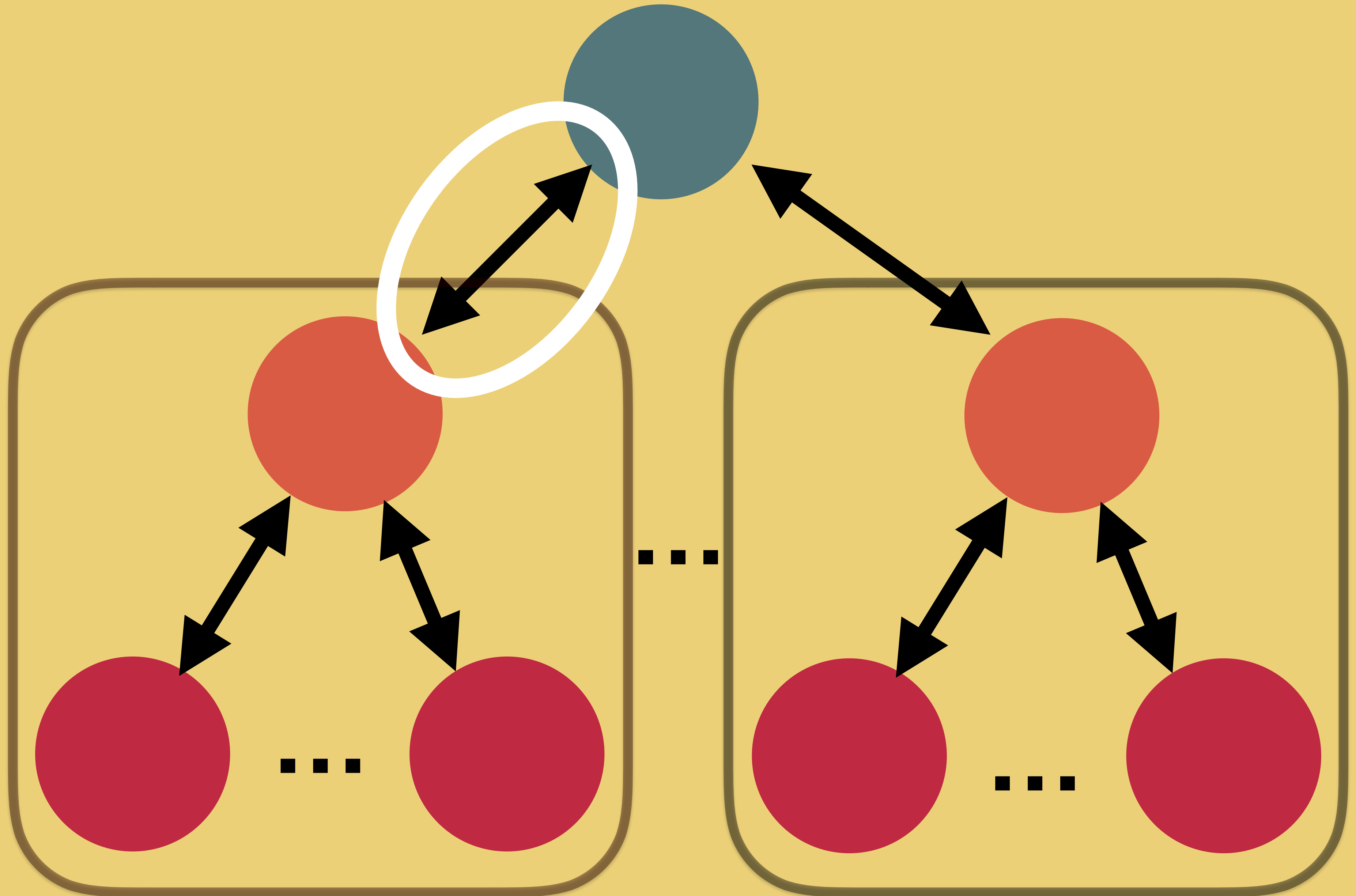




Remote

Servers

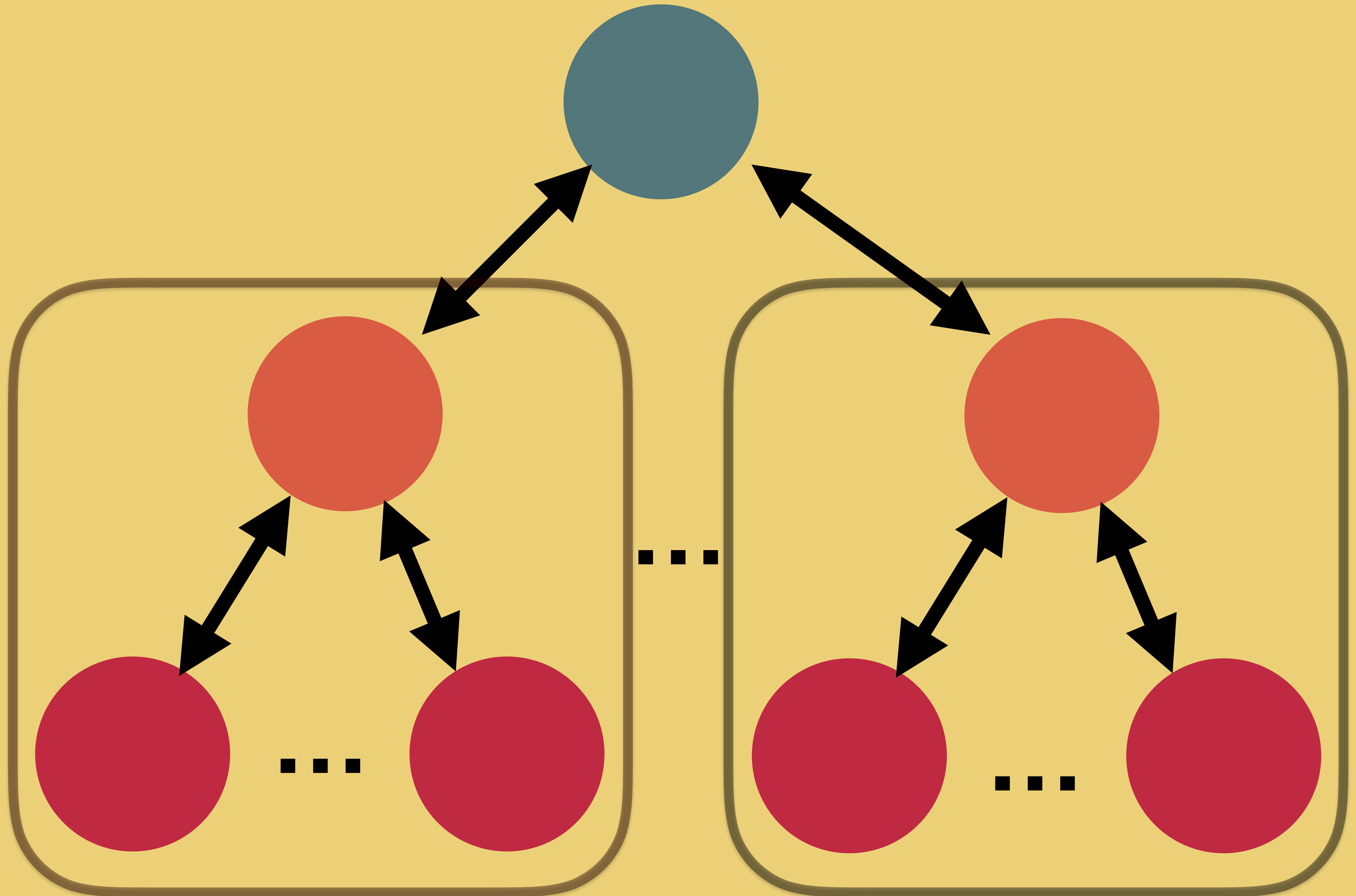
Clients



Remote

Servers

Clients



Demonstration

Implementation:

Remote: AWS instance + MySQL

Server: standalone system daemon,
decision tree classifier

Client: modified SkelCL stencil
pattern

Remote

TCP/IP

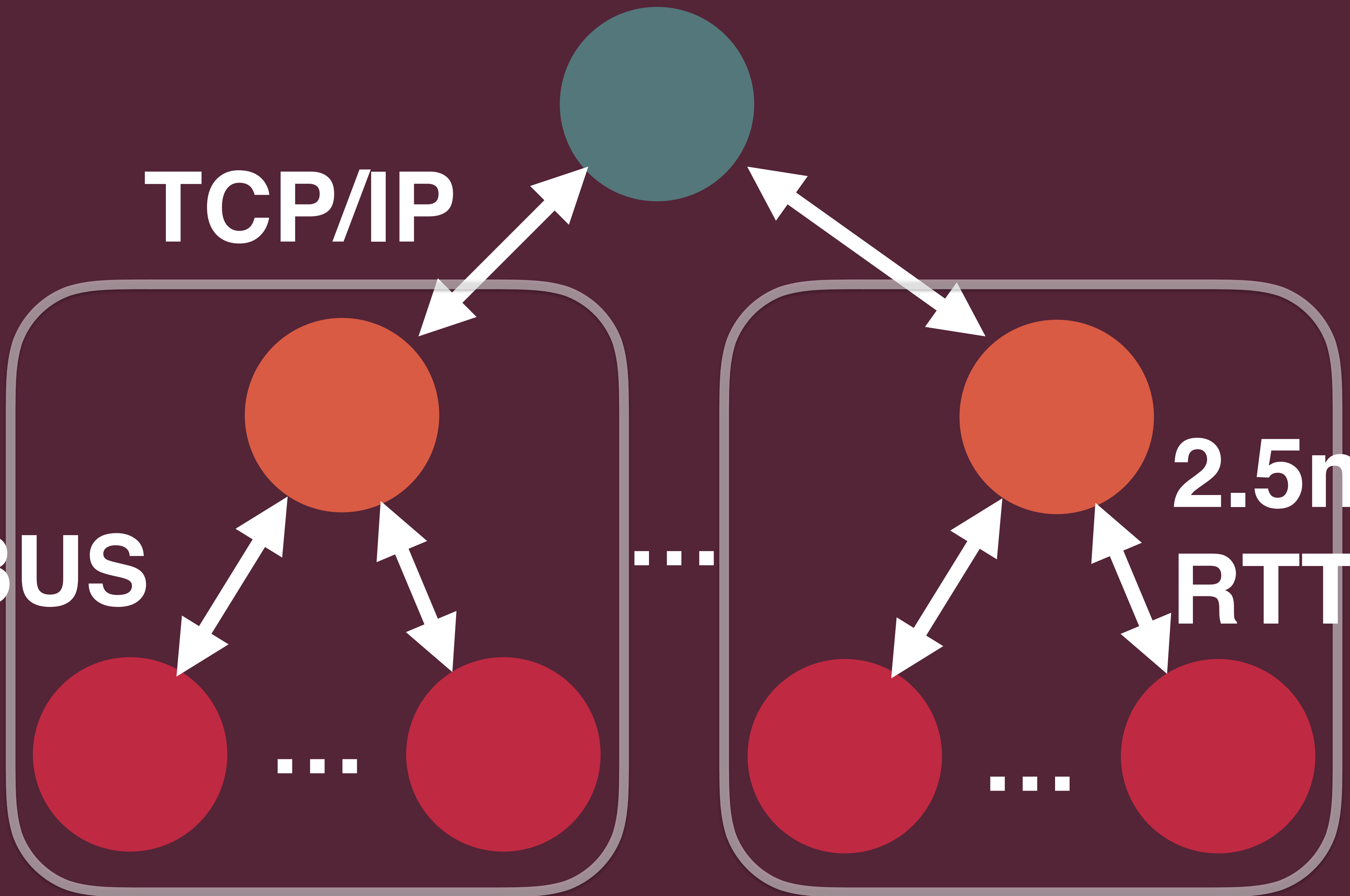
Servers

DBus

SkelCL

2.5ms

RTT



Experimental Setup:

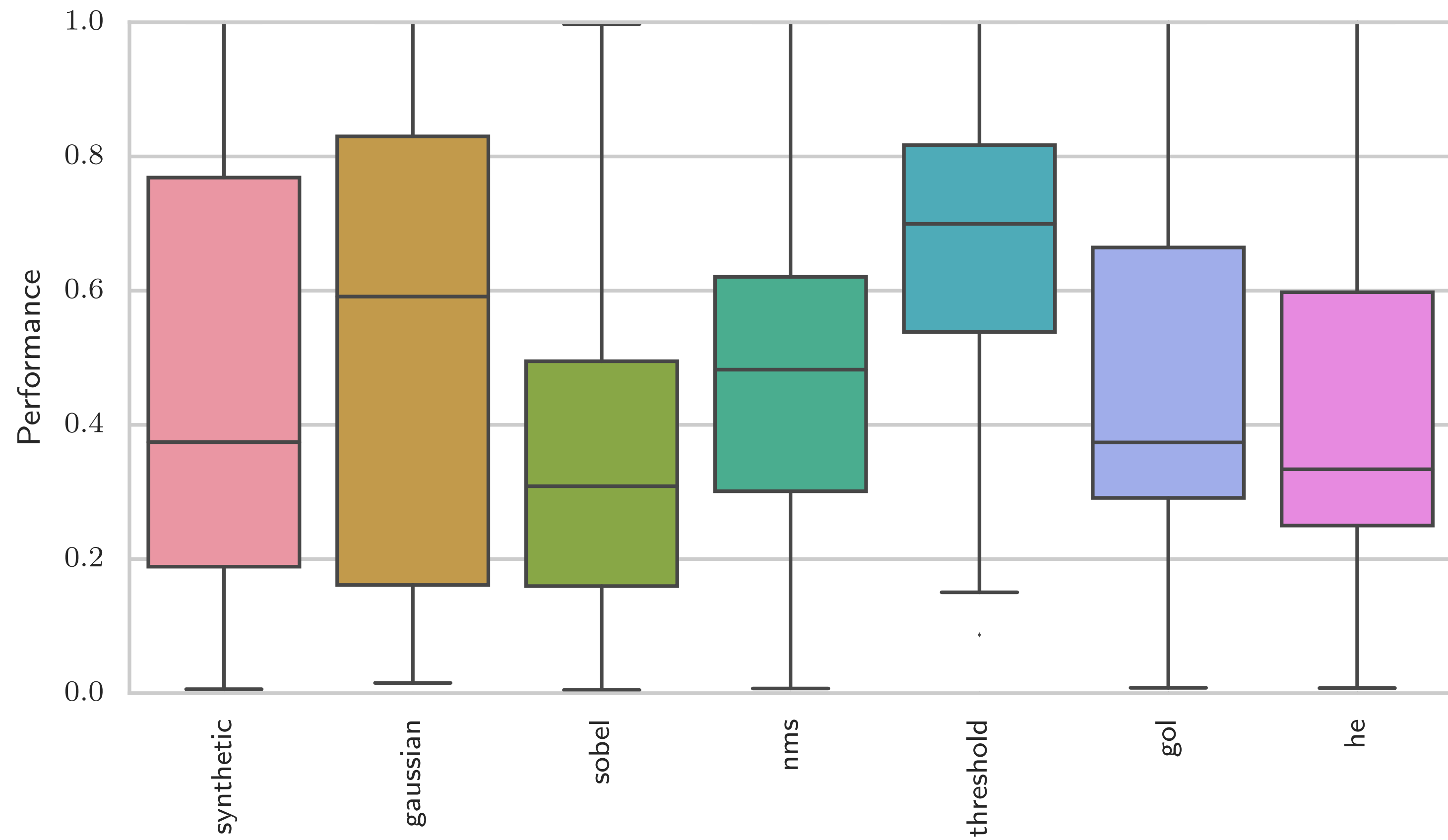
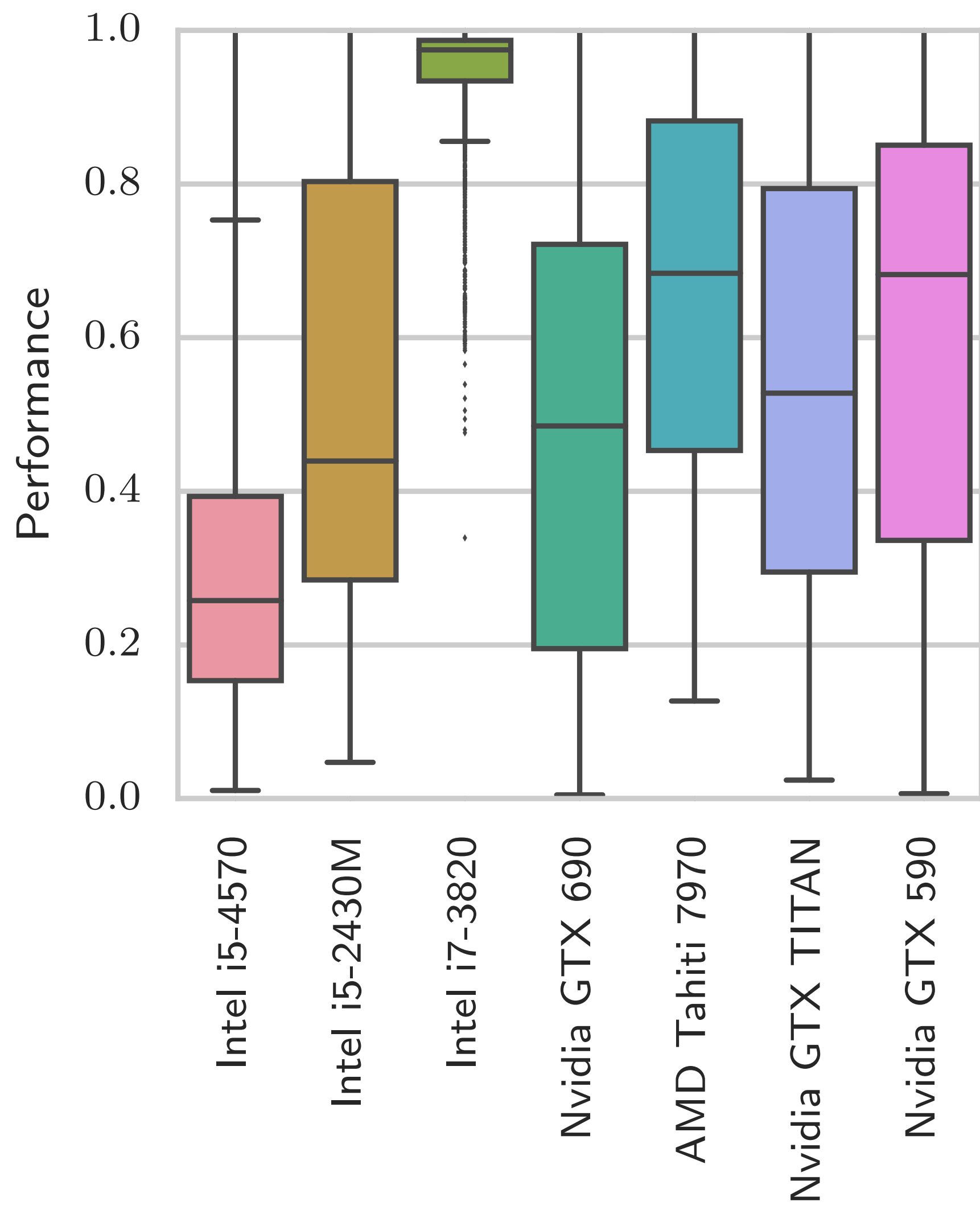
6 stencil benchmarks + synthetic.

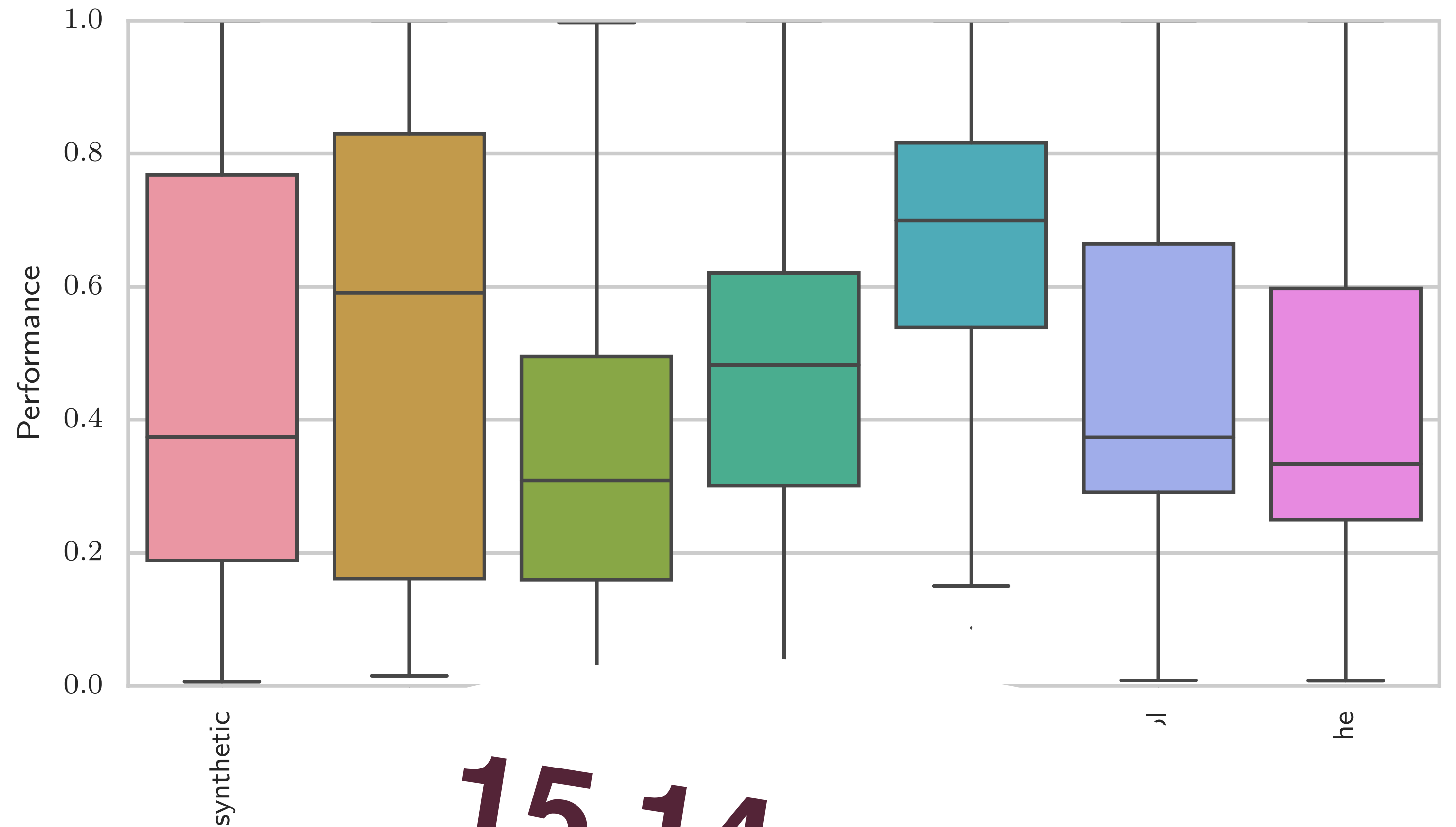
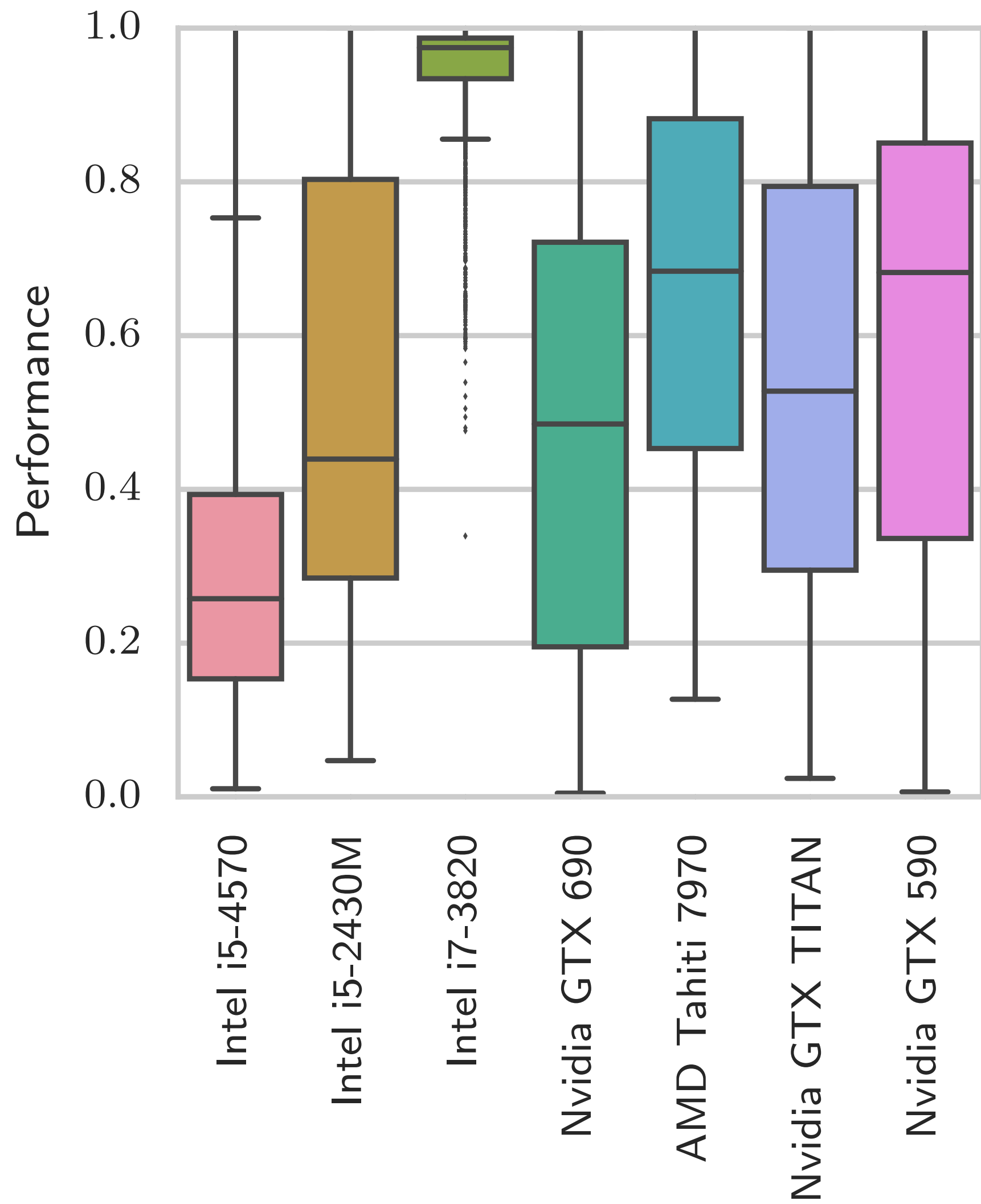
7 different GPUs & CPUs.

4 dataset sizes.

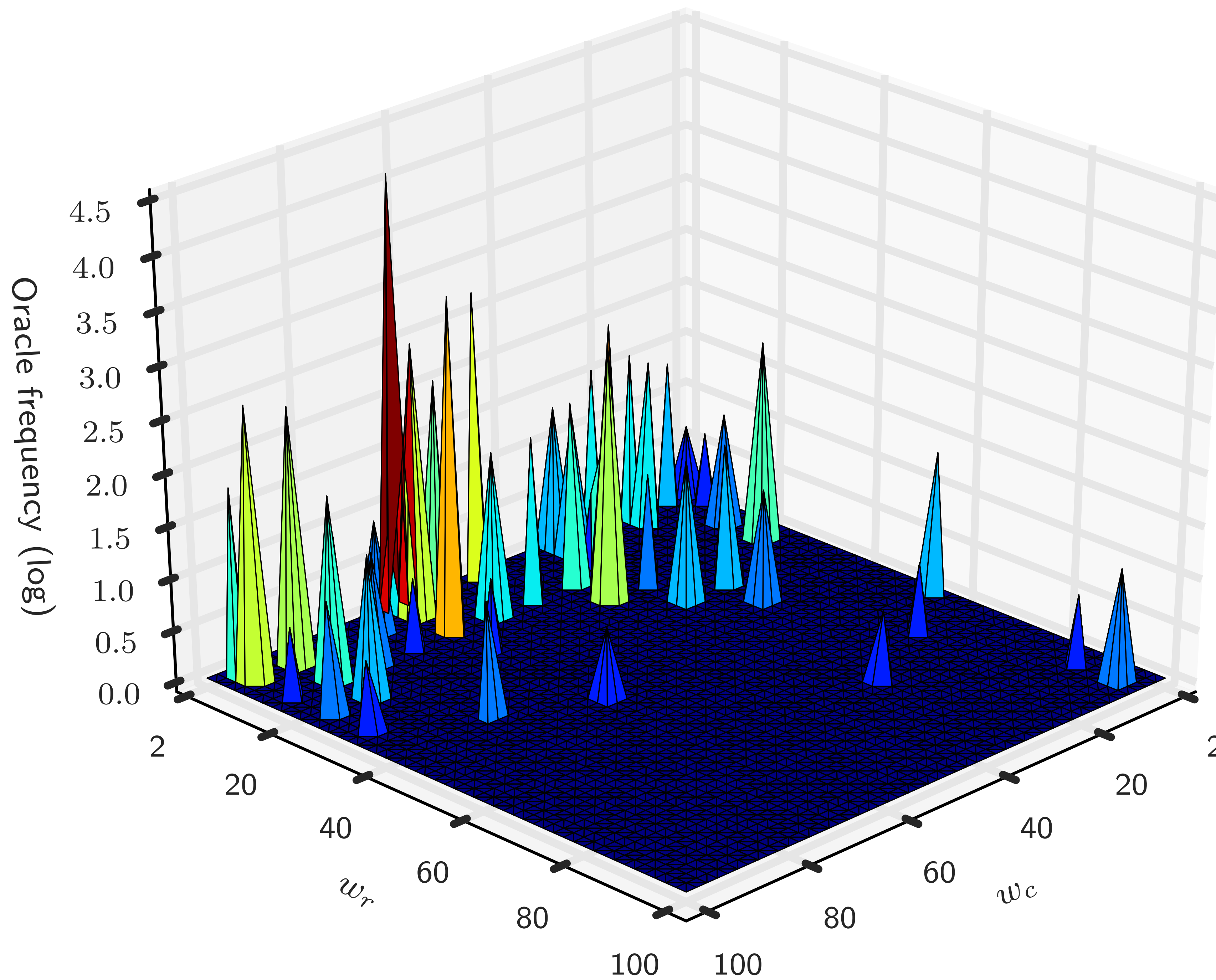
Exhaustive search of workgroup size
space for each

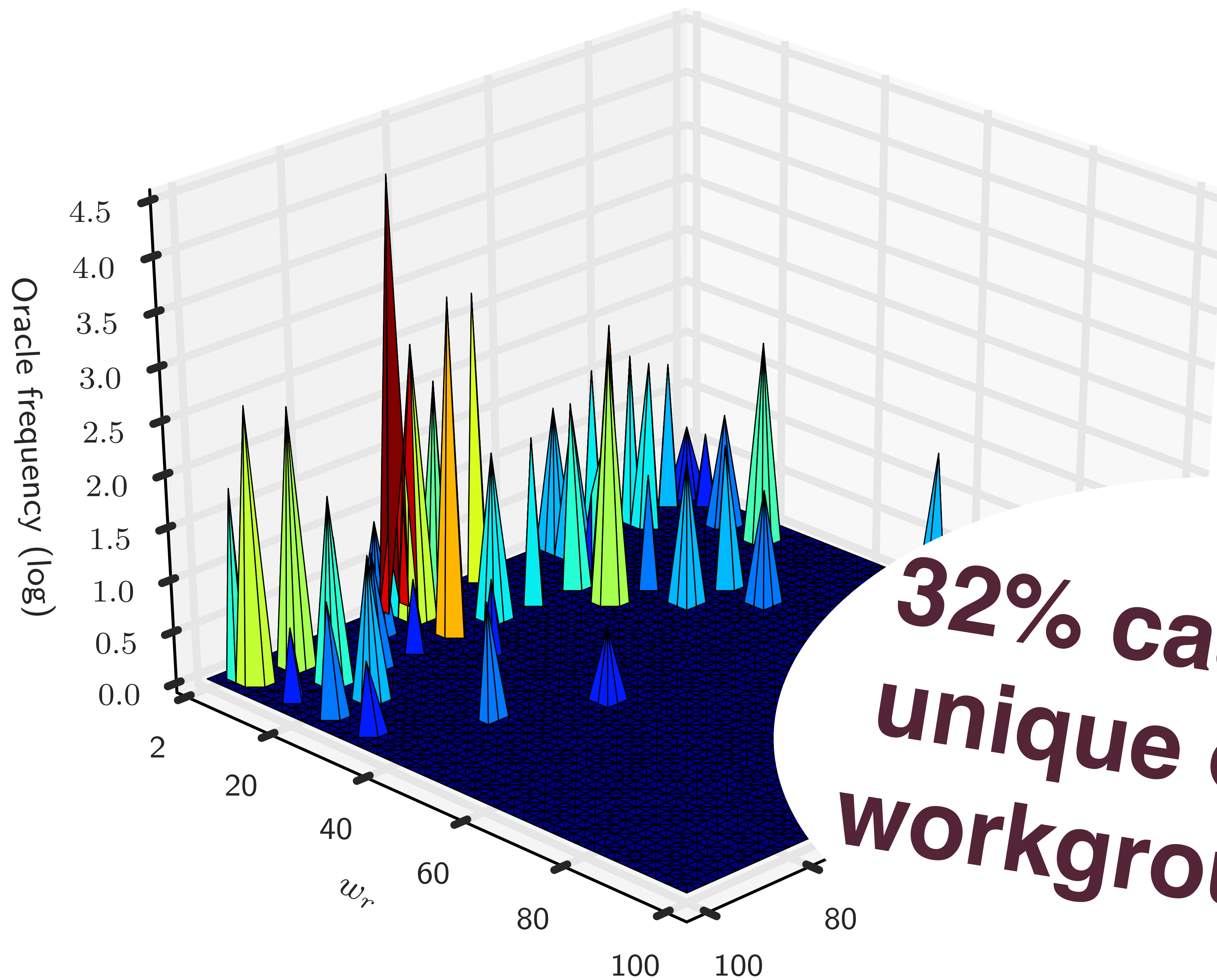
Results





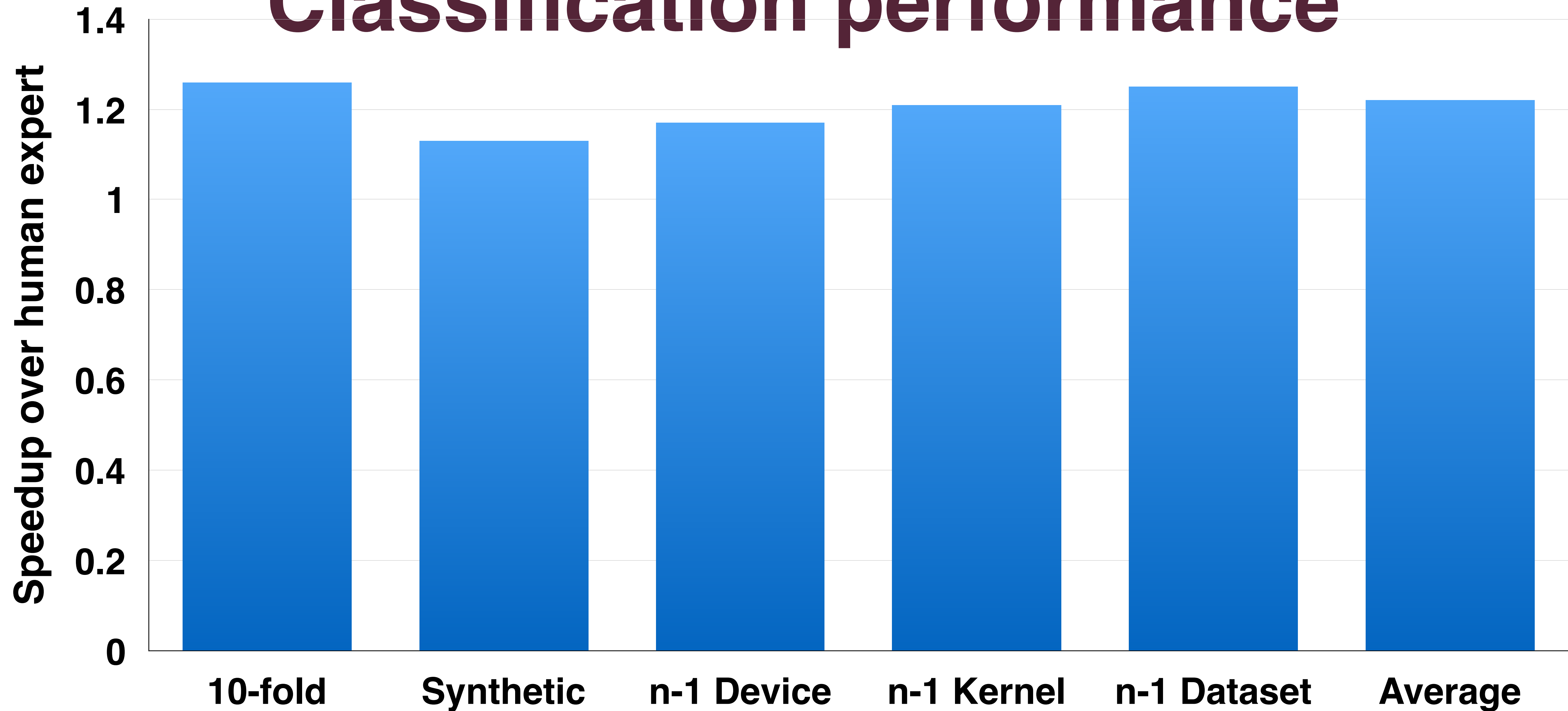
**15.14x speedup
best over worst**





**32% cases has
unique optimal
workgroup size**

Classification performance



Classification performance

Speedup over human expert

1.4
1.2
1
0.8
0.6
0.4
0.2
0

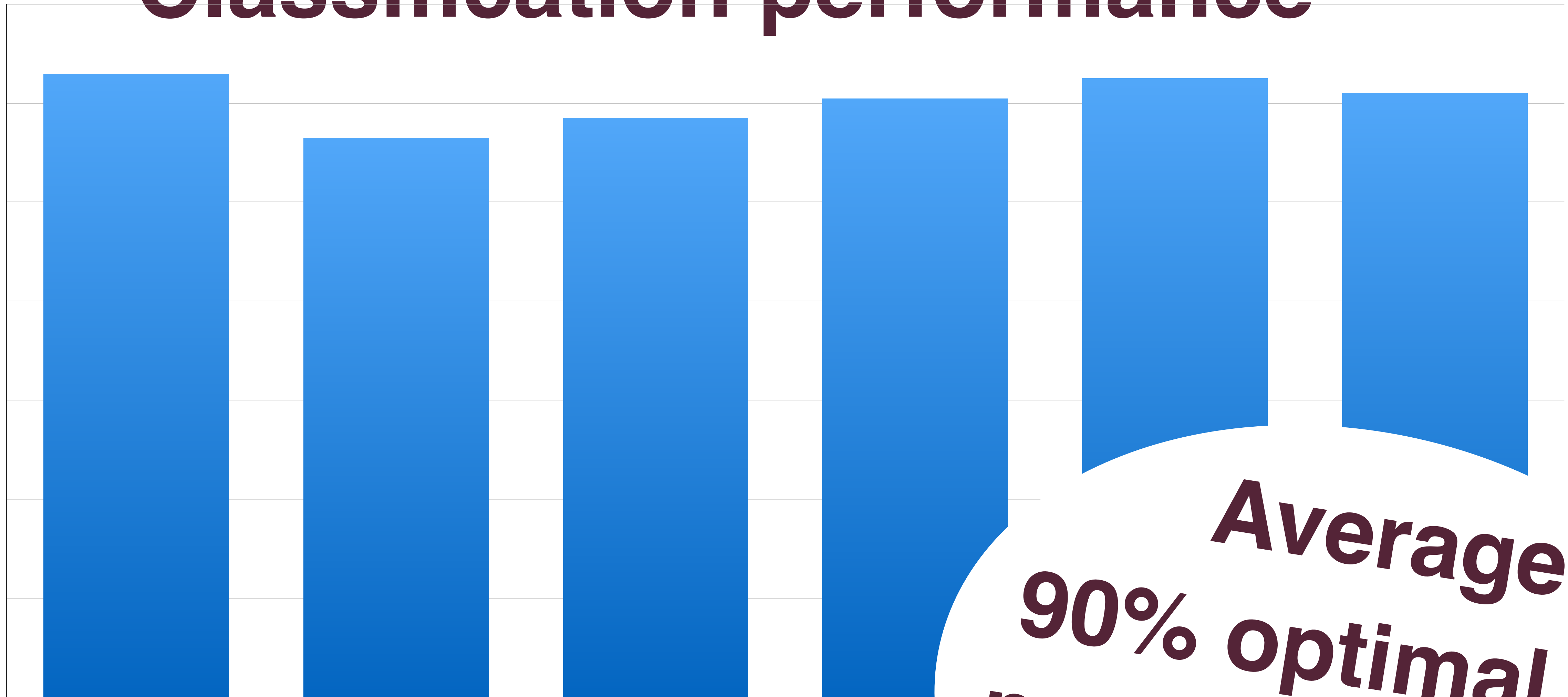
10-fold

Synthetic

n-1 Device

n-1 Ke

*Average
90% optimal
performance*



Conclusions

High level GPU code must compete
with low level on *performance*

That means *automating* the kind of
tuning which is typical of low level

**We present a framework for doing this
using machine learning**

Demonstrated using SkelCL stencils

**Achieves average 1.22x speedup over
*human expert***

Details in the paper!

Towards Collaborative Performance Tuning of Algorithmic Skeletons

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

The physical limitations of microprocessor design have forced the industry towards increasingly heterogeneous designs to extract performance. This trend has not been matched with adequate software tools, leading to a growing disparity between the availability of parallelism and the ability for application developers to exploit it.

Algorithmic skeletons simplify parallel programming by providing high-level, reusable patterns of computation. Achieving performant skeleton implementations is a difficult task; skeleton authors must attempt to anticipate and tune for a wide range of architectures and use cases. This results in implementations that target the general case and cannot provide the performance advantages that are gained from providing low level optimization parameters. Autotuning combined with machine learning offers promising performance benefits in these situations, but the high cost of autotuning and lack of available tools limits the practicality of performing autotuning at the level of the skeleton library can overcome these issues.

In this work, we present OmniTune — an extensible and distributed framework for dynamic autotuning of optimization parameters at runtime. OmniTune uses a client-server model with a flexible API to support machine learning enabled autotuning. Training data is shared across a network of cooperating systems, using a collective approach to performance tuning.

We demonstrate the practicality of OmniTune in a case study using the algorithmic skeleton library SkelCL. By automatically tuning the workgroup size of OpenCL Stencil skeleton kernels, we show that that static tuning across a range of GPUs and programs can achieve only 26% of the optimal performance, while OmniTune achieves 92% of this maximum, equating to an average 5.65× speedup. OmniTune achieves this without introducing a significant runtime overhead, and enables portable, cross-device and cross-program tuning.

1. Introduction

General purpose programming with GPUs has been shown to provide huge parallel throughput, but poses a significant programming challenge, requiring application developers to master an unfamiliar programming model (such as provided by CUDA or OpenCL) and architecture (SIMD with a multi-level memory hierarchy). As a result, GPGPU programming is often considered beyond the realm of everyday development. If steps are not taken to increase the accessibility of such parallelism, the gap between potential and utilized performance will continue to widen as hardware core counts increases.

Algorithmic skeletons offer a solution to this this *programmability challenge* by raising the level of abstraction. This simplifies parallel programming, allowing developers to focus on solving problems rather than coordinating parallel resources. Skeleton frameworks provide robust parallel implementations of common patterns of computation which developers parameterise with their application-specific code. This greatly reduces the challenge of parallel programming, allowing users to structure their problem-solving logic sequentially, while offloading the cognitive cost of parallel coordination to the skeleton library author. The rising number of skeleton frameworks supporting graphics hardware illustrates the demand for high level abstractions for GPGPU programming [1, 2]. The challenge is in maintaining portable performance across the breadth of devices in the rapidly developing GPU and heterogeneous architecture landscape.

1.1 The Performance Portability Challenge

There are many factors — or *parameters* — which influence the behavior of parallel programs. For example, setting the number of threads to launch for a particular algorithm. The performance of parallel programs is sensitive to the values of these parameters, and when tuning to maximize performance, one size does not fit all. The suitability of parameter values depends on the program implementation, the target hardware, and the dataset that is operated upon. Iterative compilation and autotuning have been shown to help in these cases by automating the process of tuning parameter values to match individual execution environments [3]. However, there have been few attempts to develop general mechanisms for these techniques, and the time taken to develop ad-hoc autotuning solutions and gather performance data is often prohibitively expensive.

We believe that by embedding autotuning at the skeletal level, it is possible to achieve performance with algorithmic skeletons that is competitive with — and in some cases, exceeds — that of hand tuned parallel implementations which traditionally came at the cost of many man hours of work from expert programmers to develop.

Incorporating autotuning into algorithmic skeleton libraries has two key benefits: first, it minimizes development effort by requiring only a modification to the skeleton implementation rather than to every user program; and second, by targeting a library, it enables a broader and more substantive range of performance data to be gathered than with ad-hoc tuning of individual programs.

2016/1/6

2016/1/6

2016/1/6

Towards *Collaborative* *Performance Tuning* of Algorithmic Skeletons

<http://chriscummins.cc>