



Autotuning and Algorithmic Skeletons

Uncontentious statement:

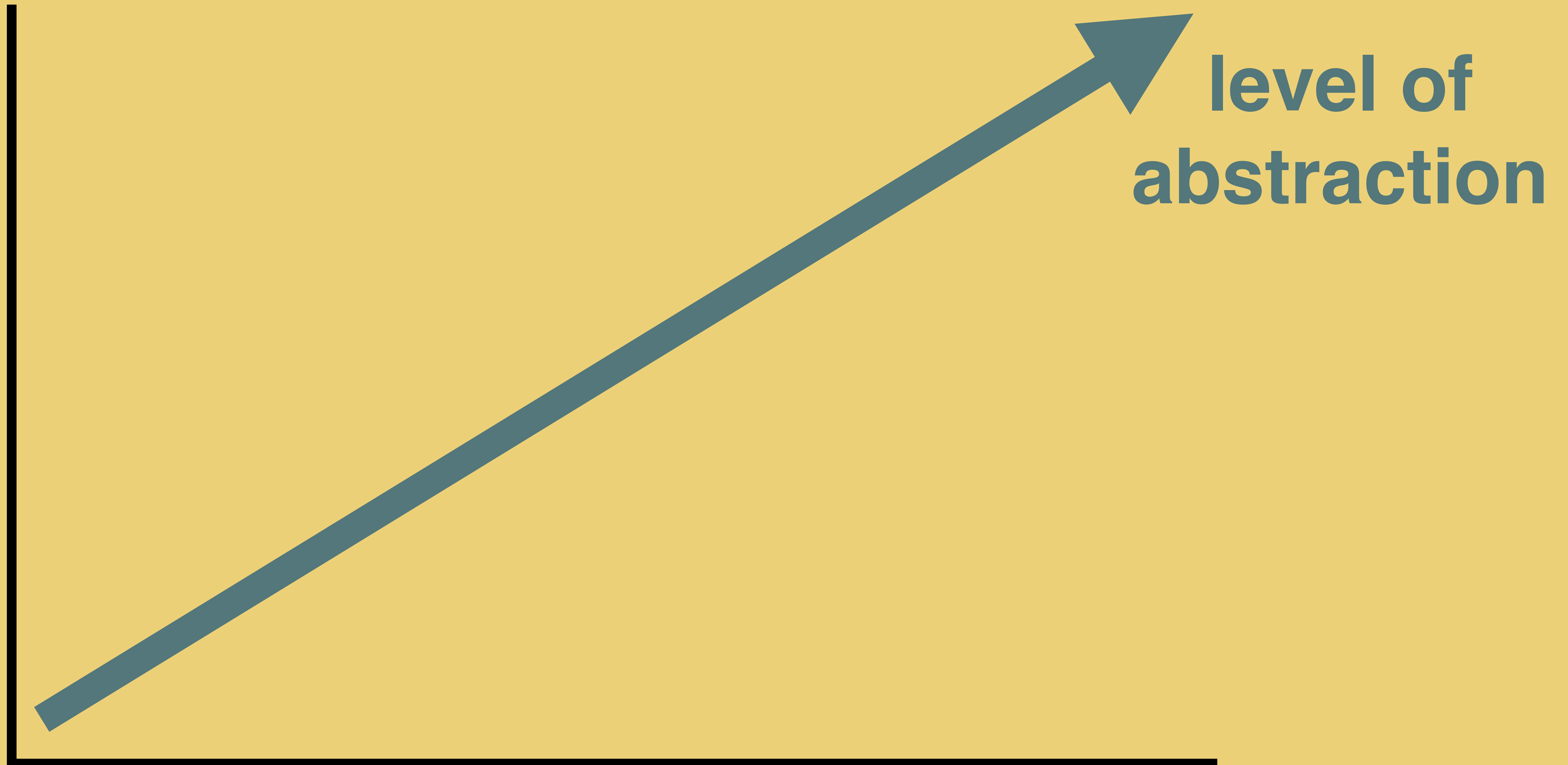
**High level
programming
is great!**

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

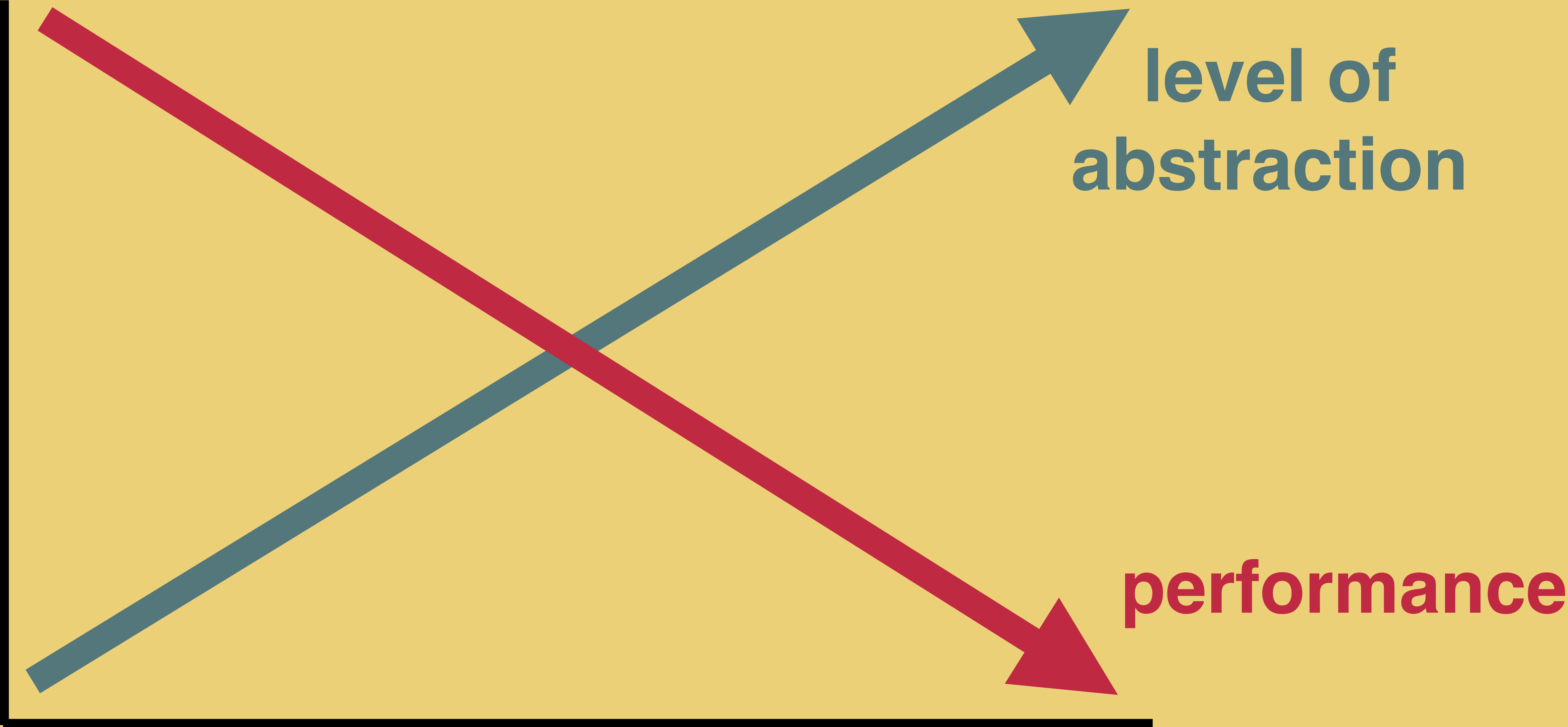
So why do application
programmers resort
to writing *low level*
code?

Performance!

common perception is ...



common perception is ...





**programmers
who care about
*performance***

**programmers
who care about
*abstractions***



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 label 'us' at the bottom to this white dot.

**programmers
who care about
*performance***

**programmers
who care about
*abstractions***

us



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 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

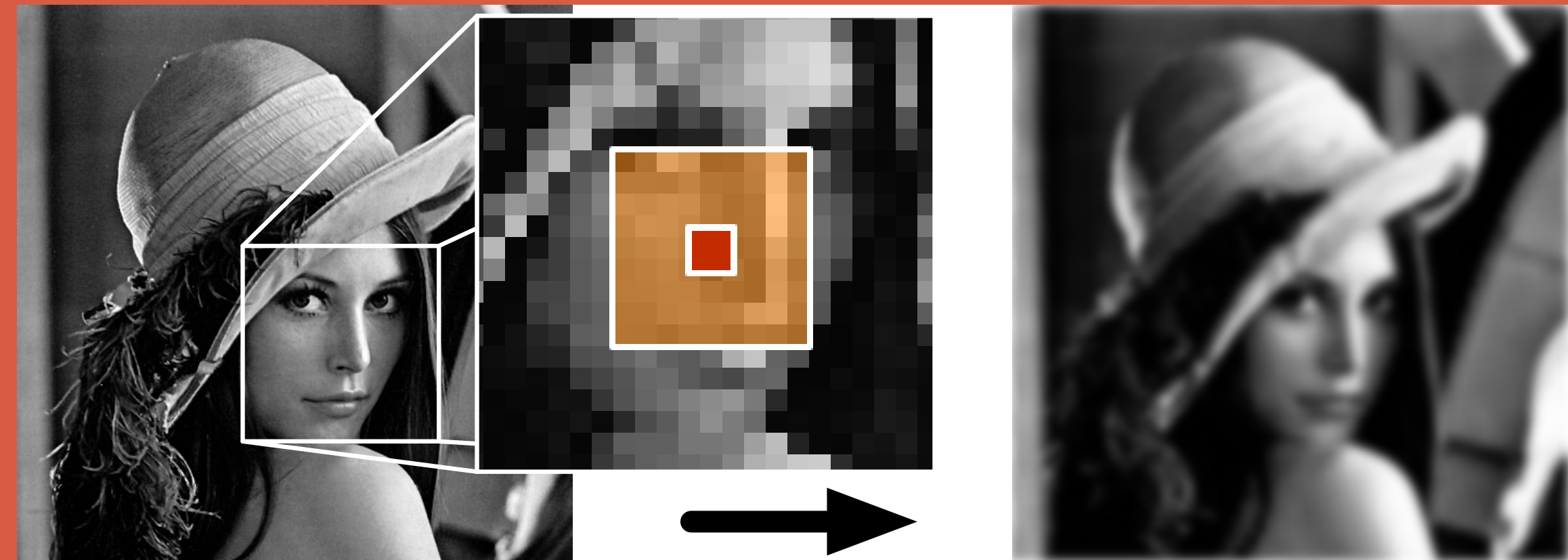
Reasons for low level:

**Domain-specific
optimisations**

Parameter tuning

Parameter tuning for Algorithmic Skeletons

Parameter tuning for ~~Algorithmic~~ ~~Skeletons~~ *stencils*



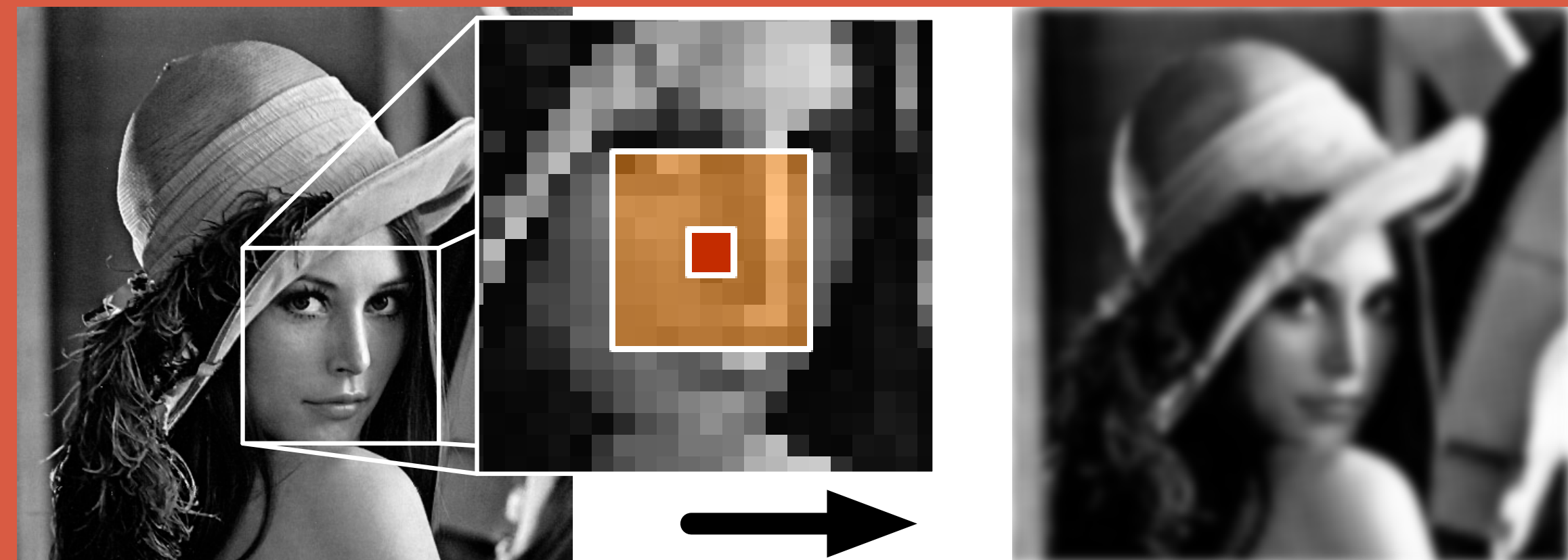
OpenCL workgroup size

~~Parameter tuning~~

~~for Algorithmic~~

~~Skeletons~~

stencils



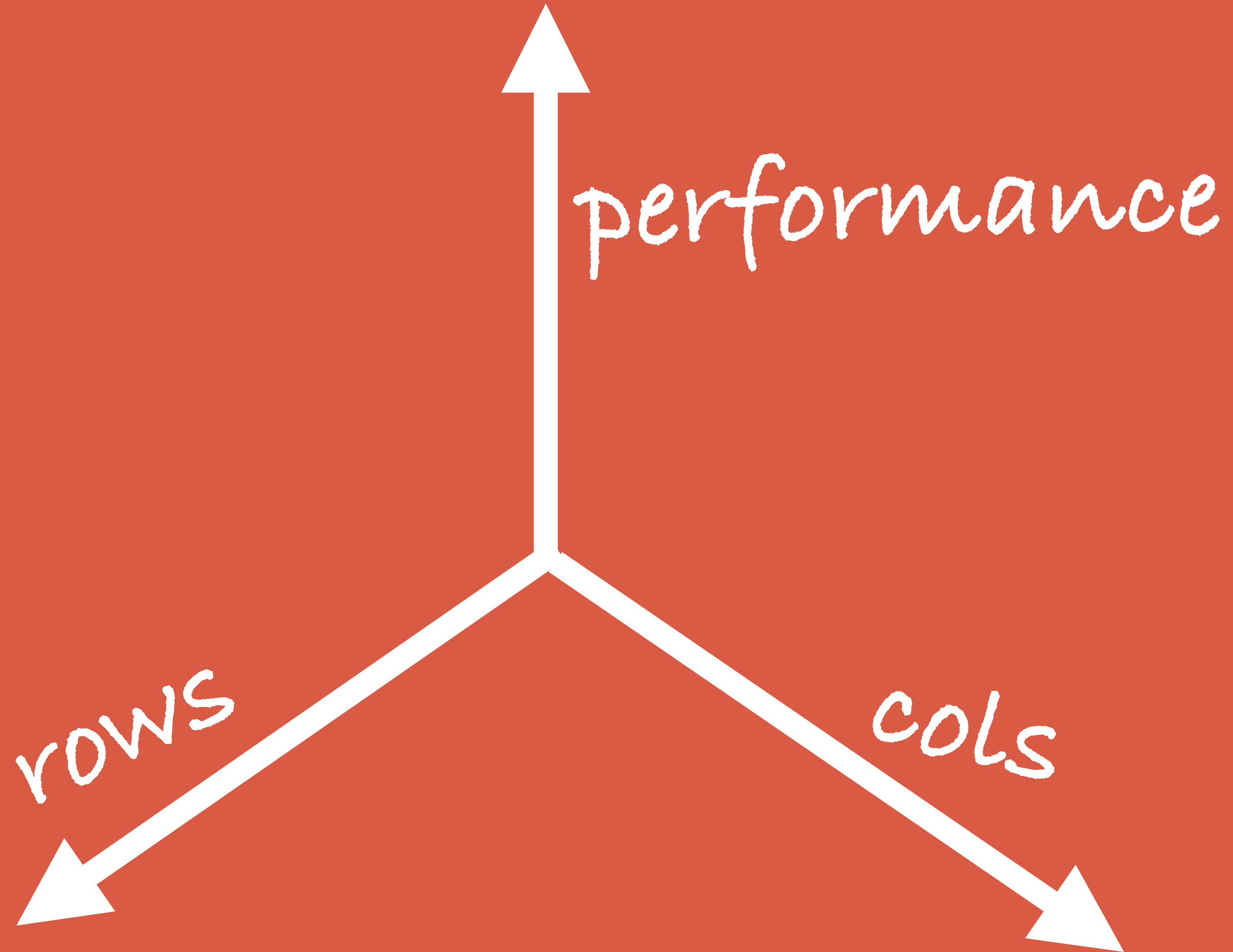
OpenCL workgroup size:

OpenCL workgroup size:
Controls composition of
hardware threads.

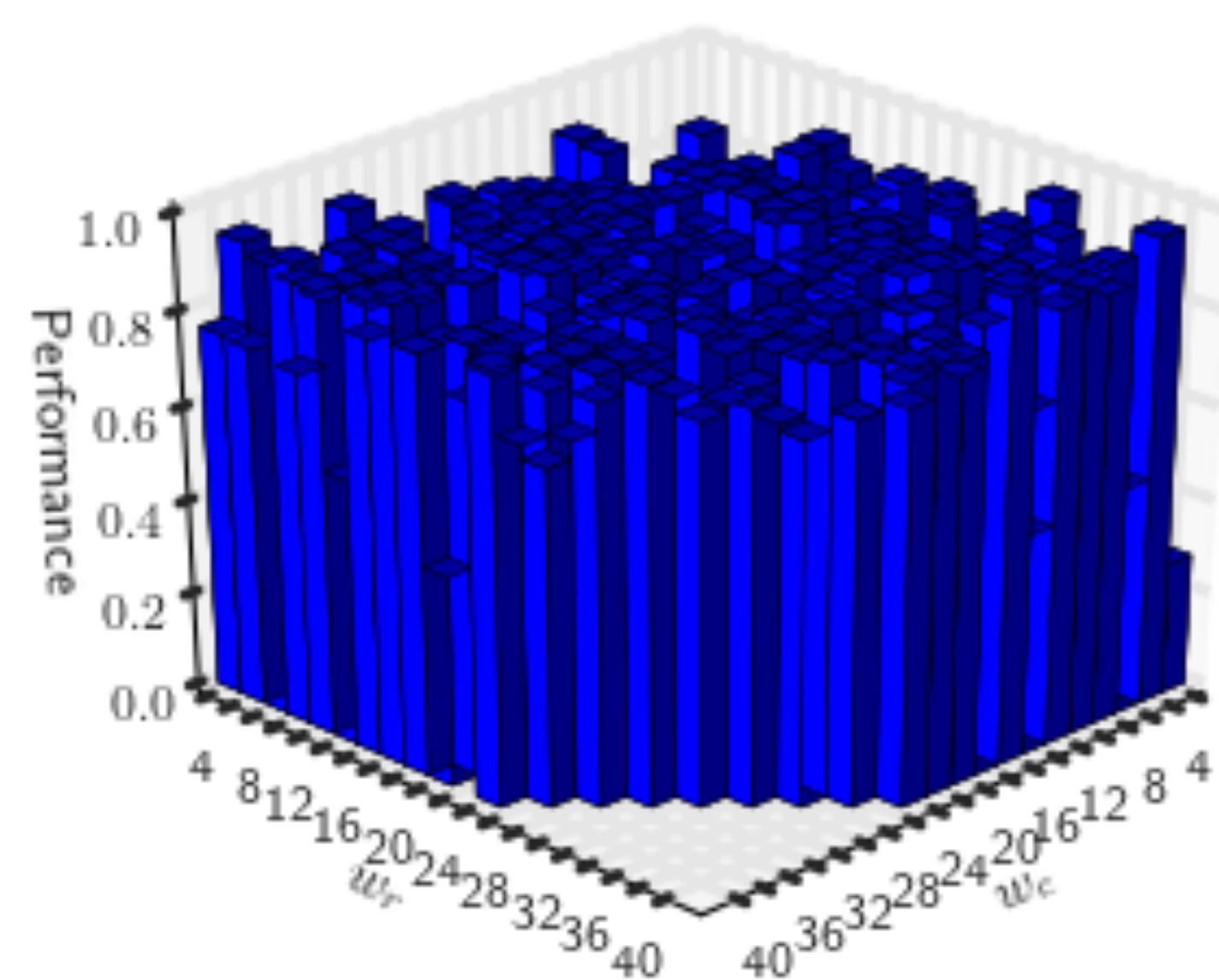
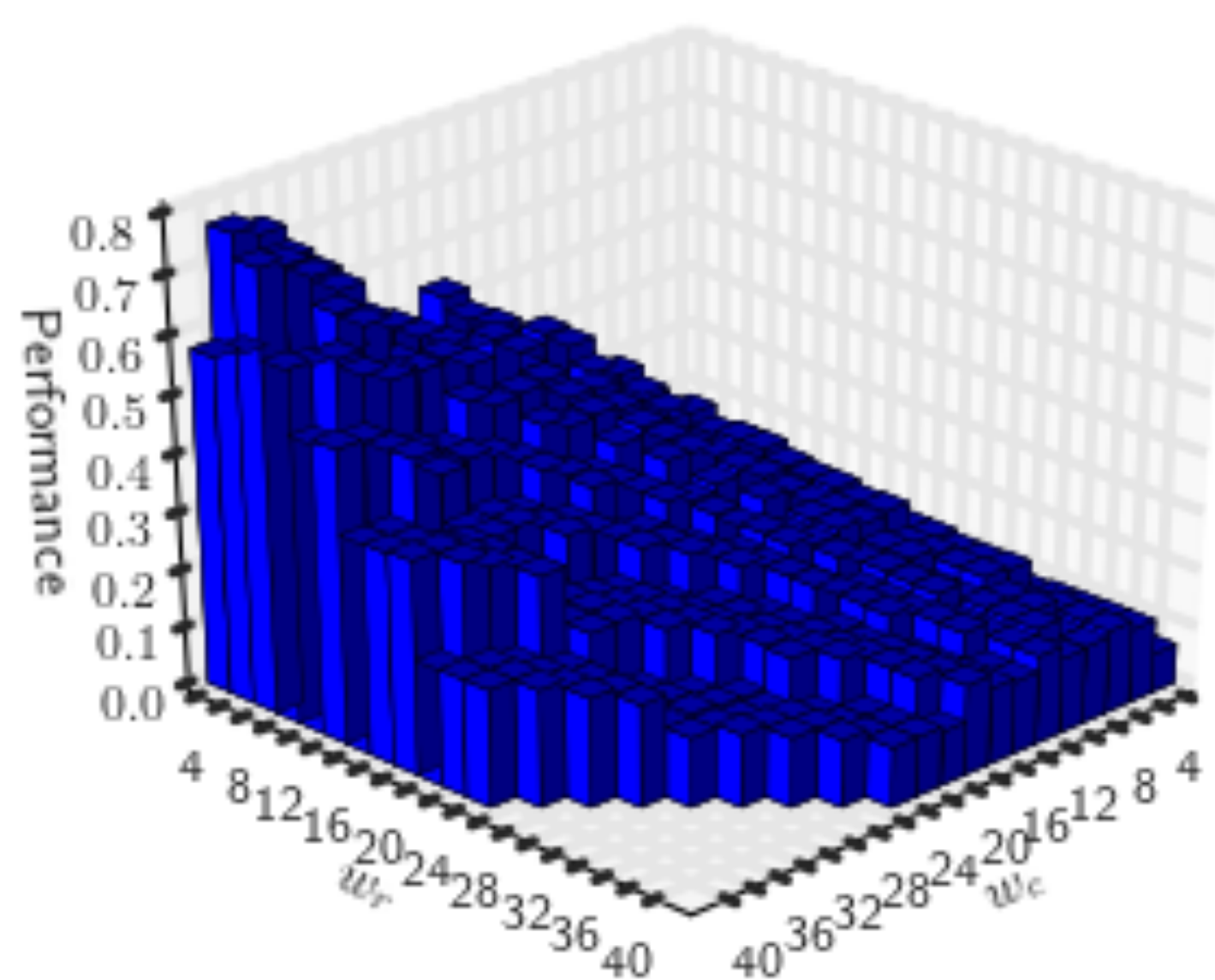
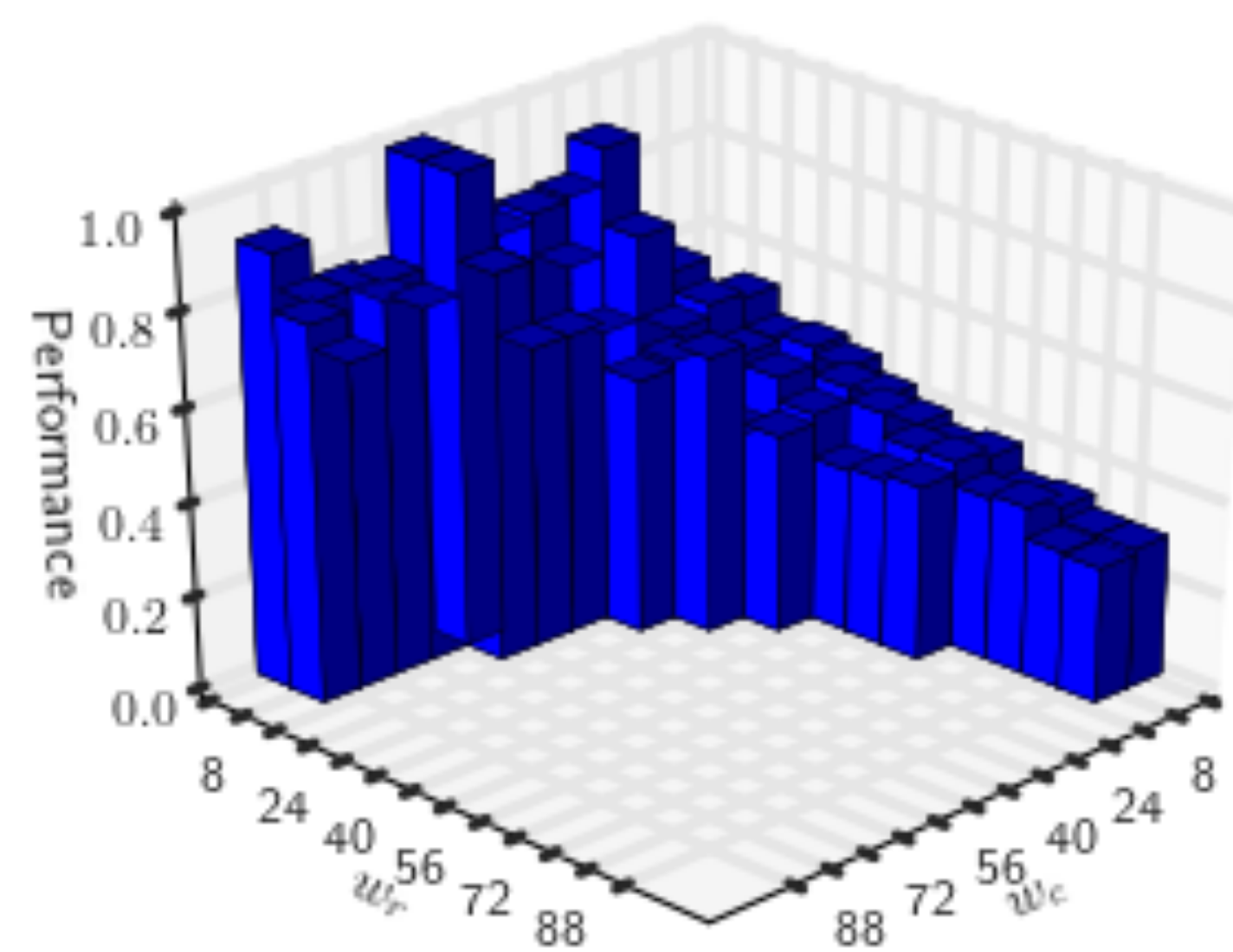
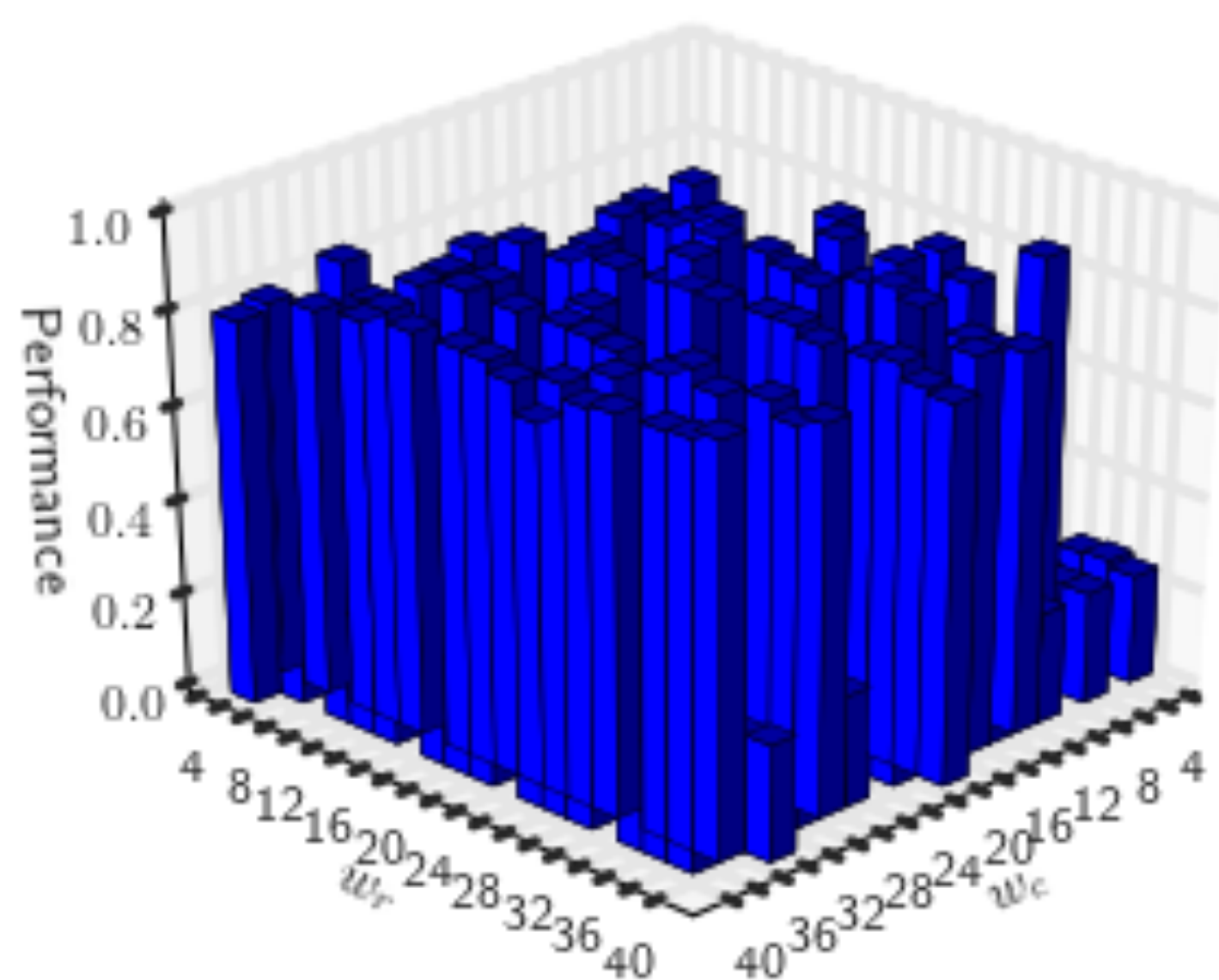
OpenCL workgroup size:
Controls composition of
hardware threads.
Is a 2D parameter (rows x cols).

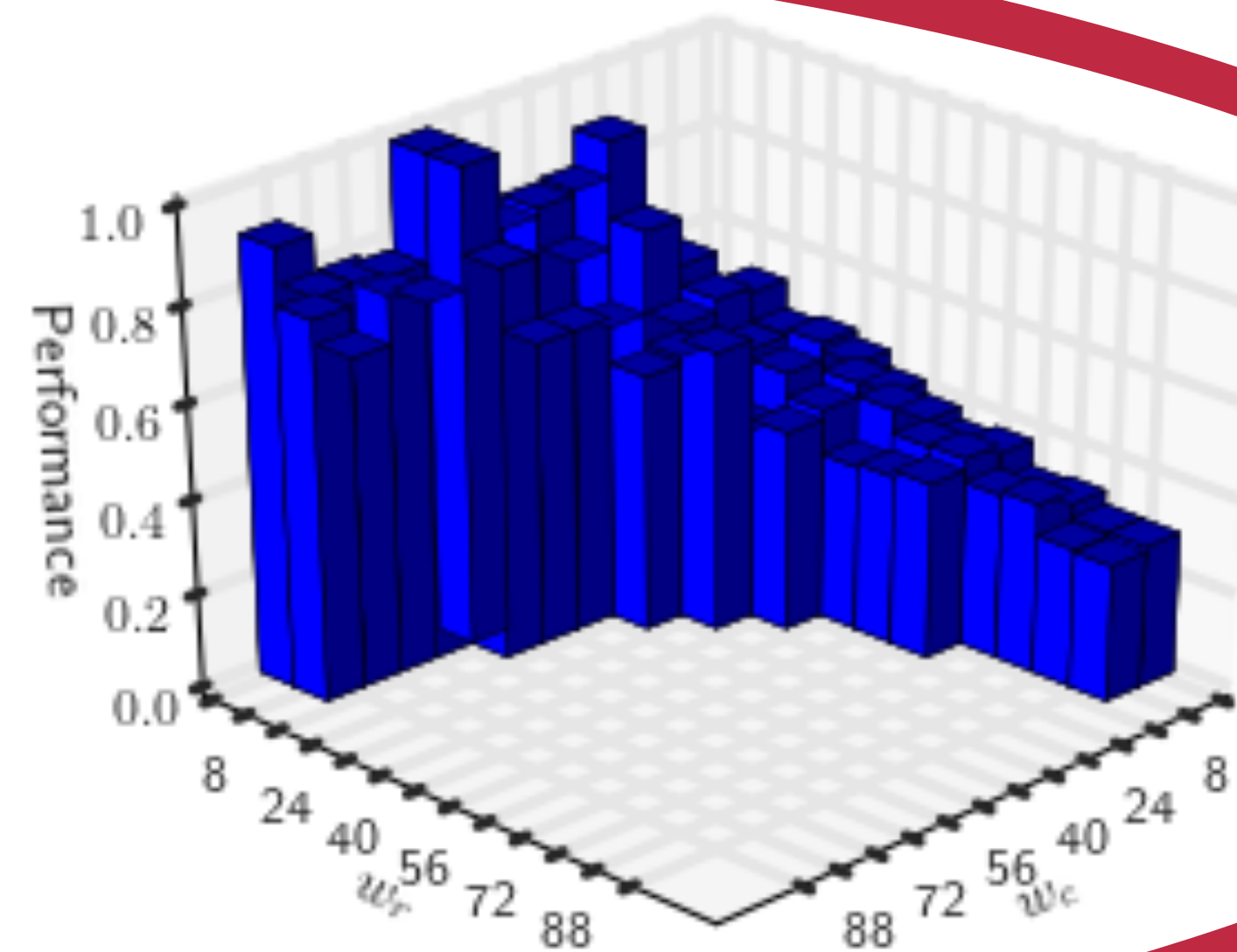
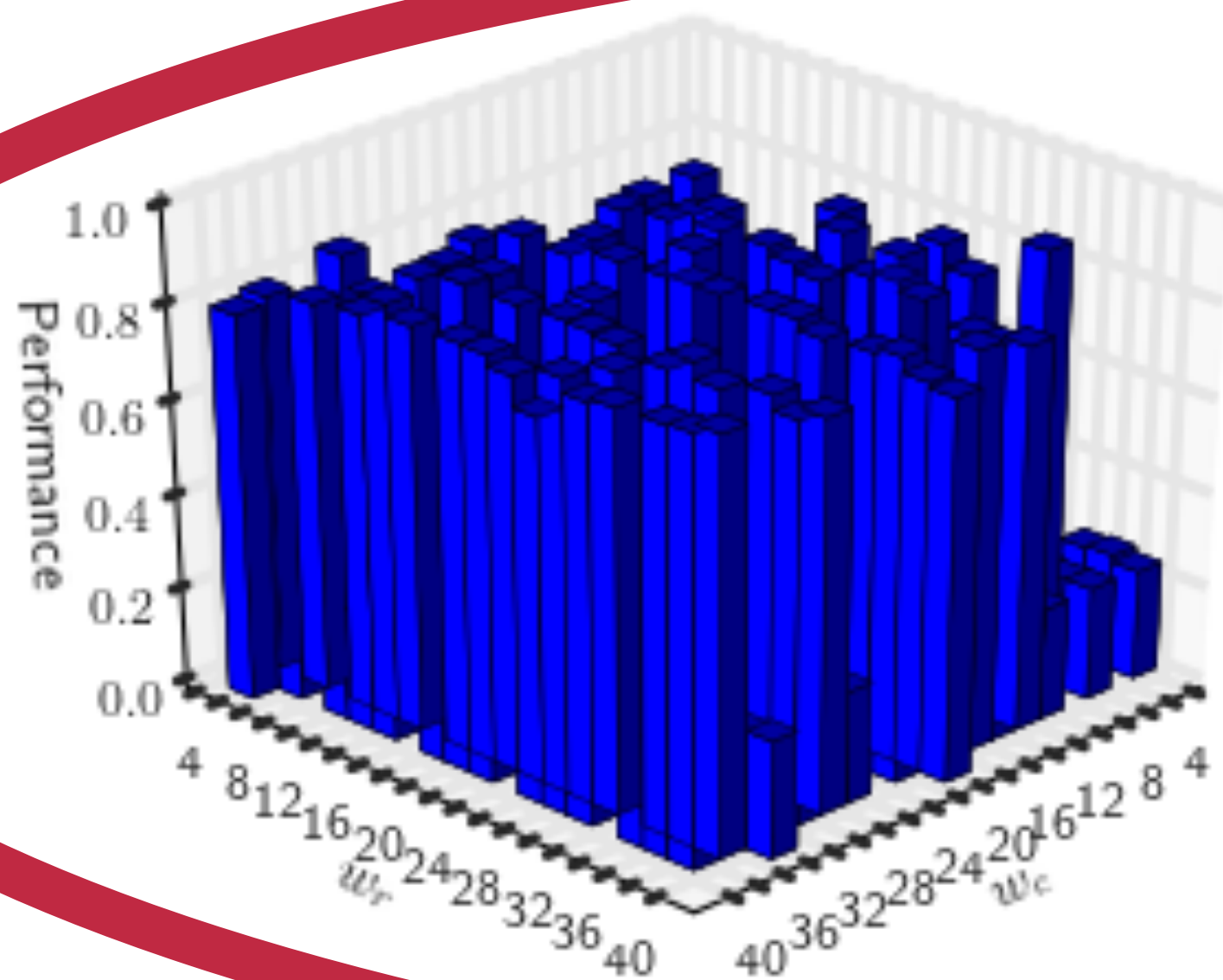
OpenCL workgroup size:
Controls composition of
hardware threads.
Is a 2D parameter (rows x cols).
Critical to performance.

OpenCL workgroup size:

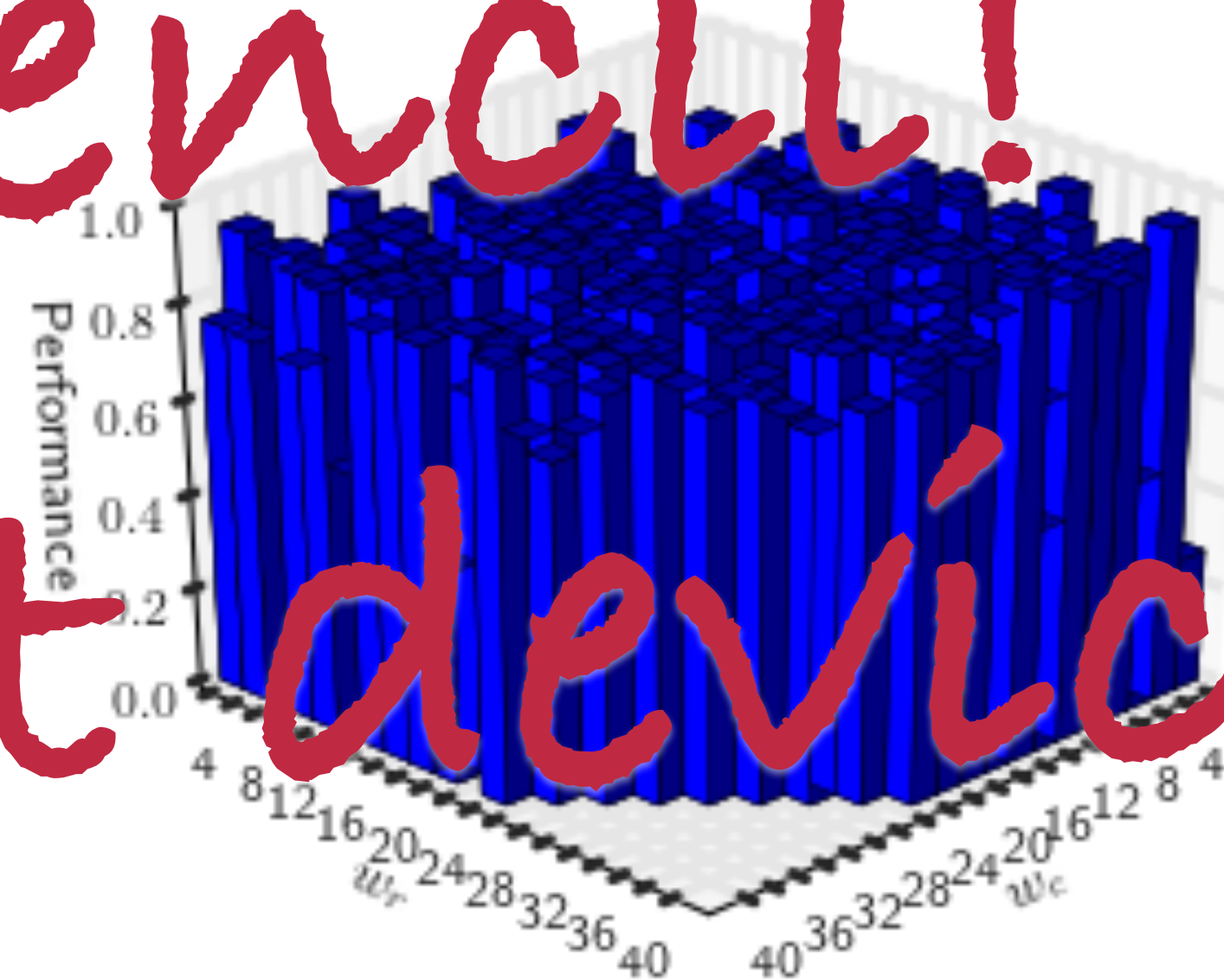
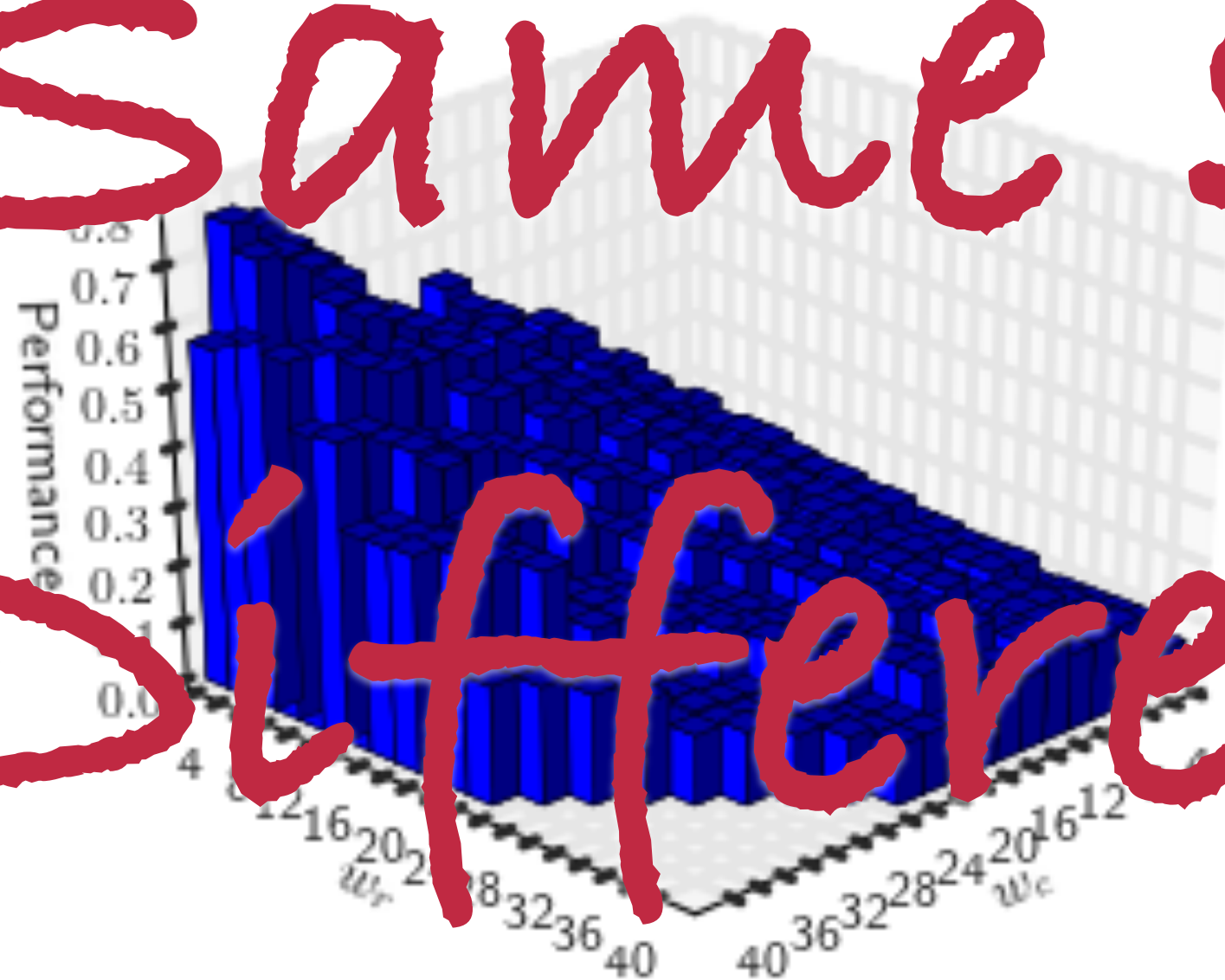


Examples

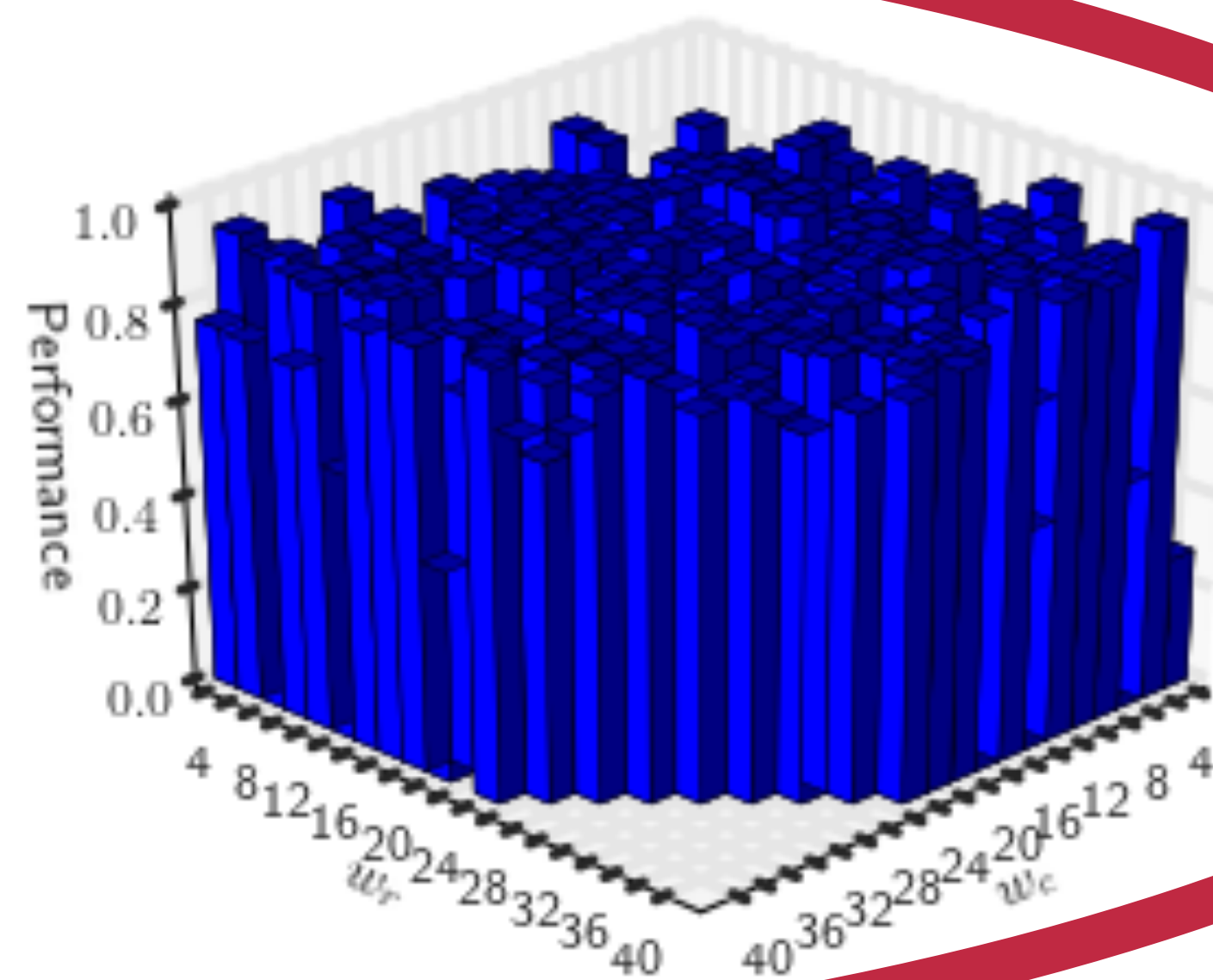
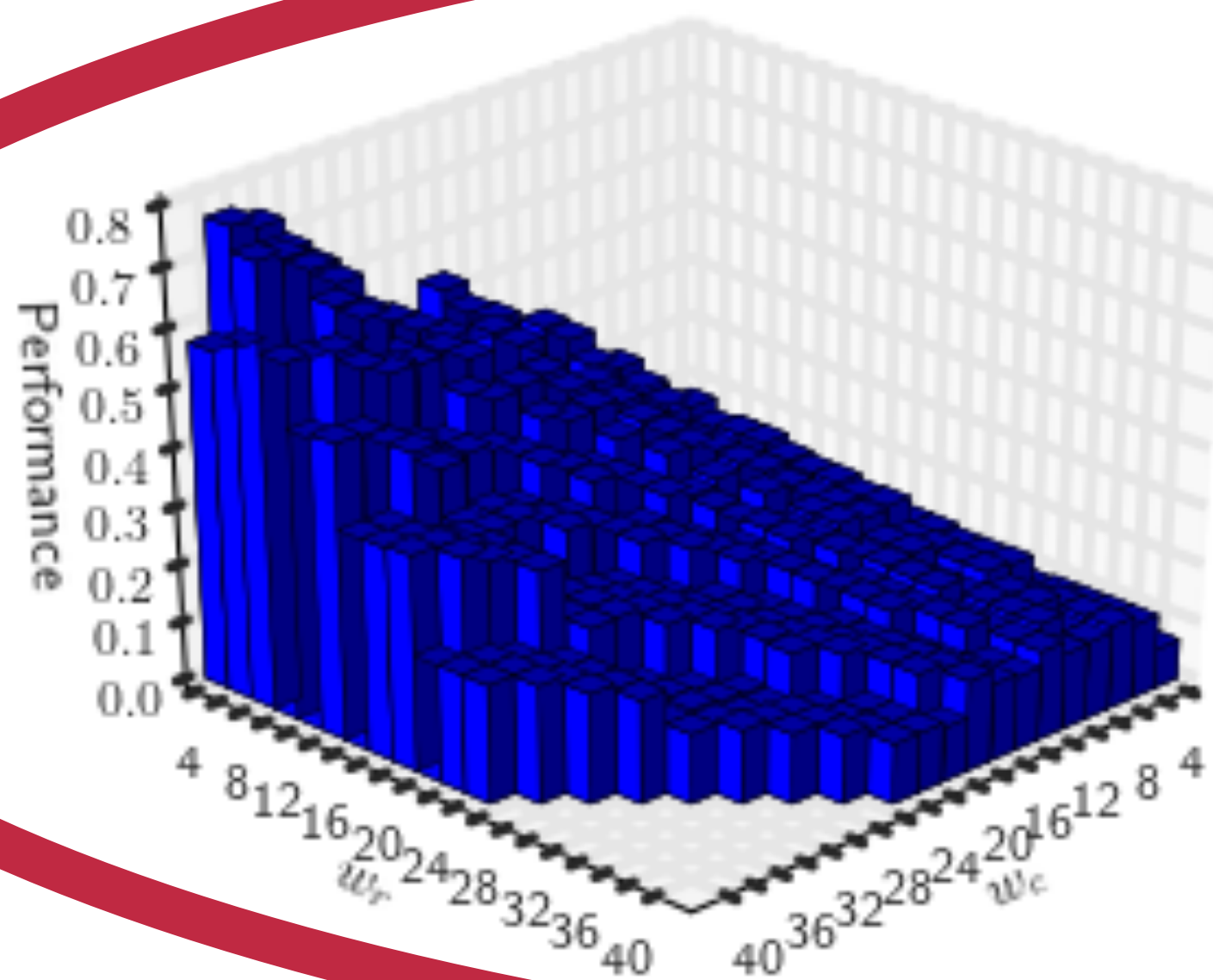


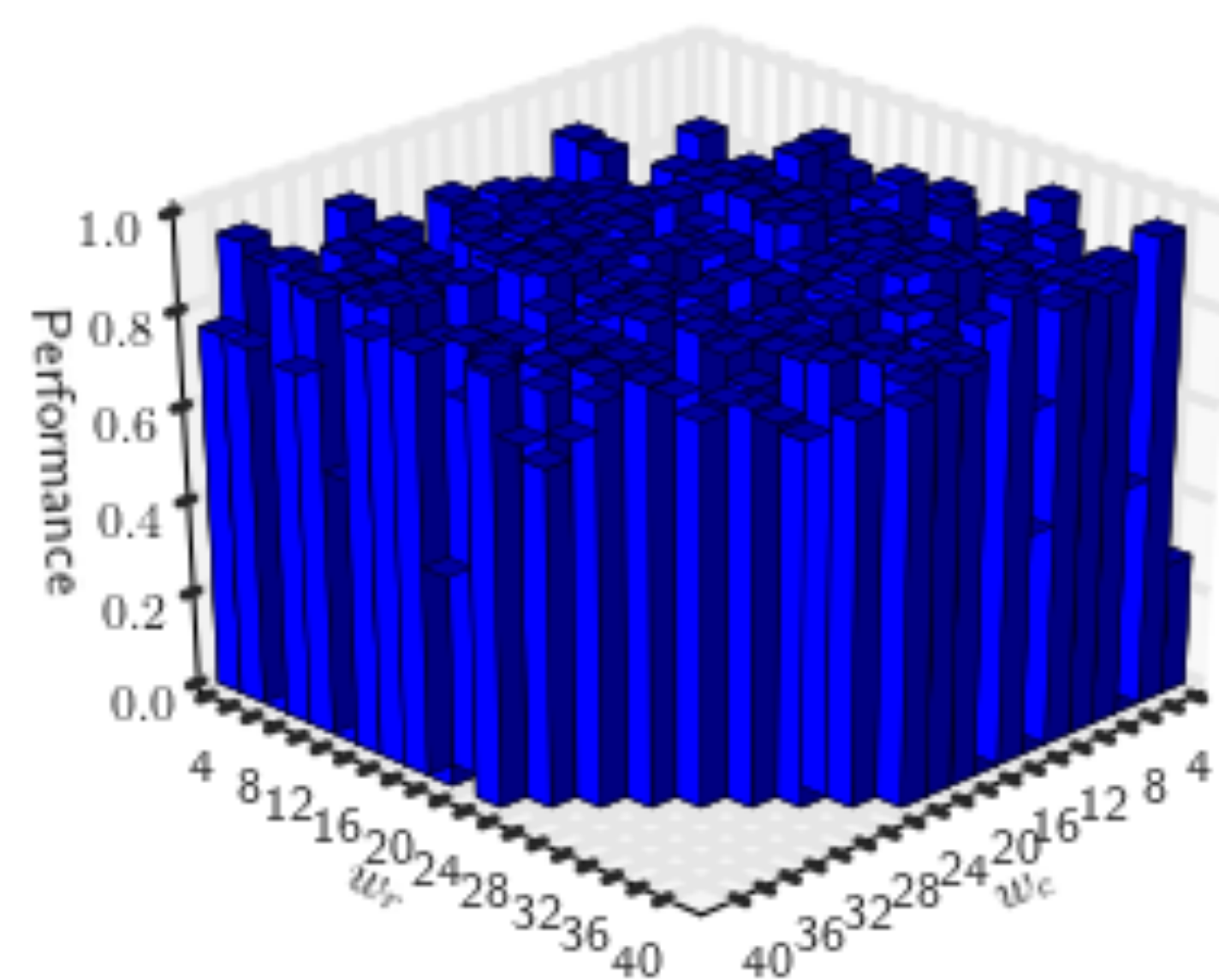
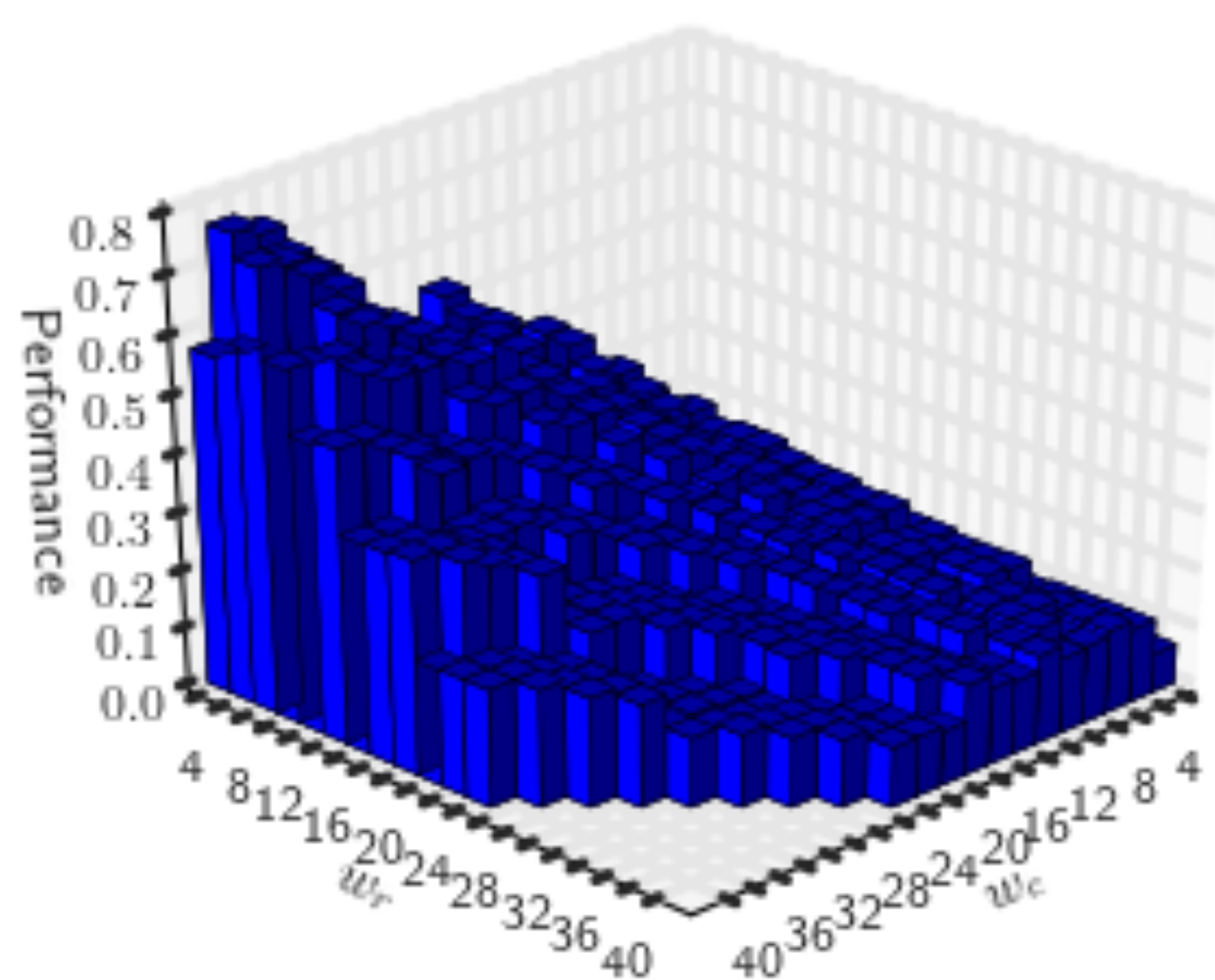
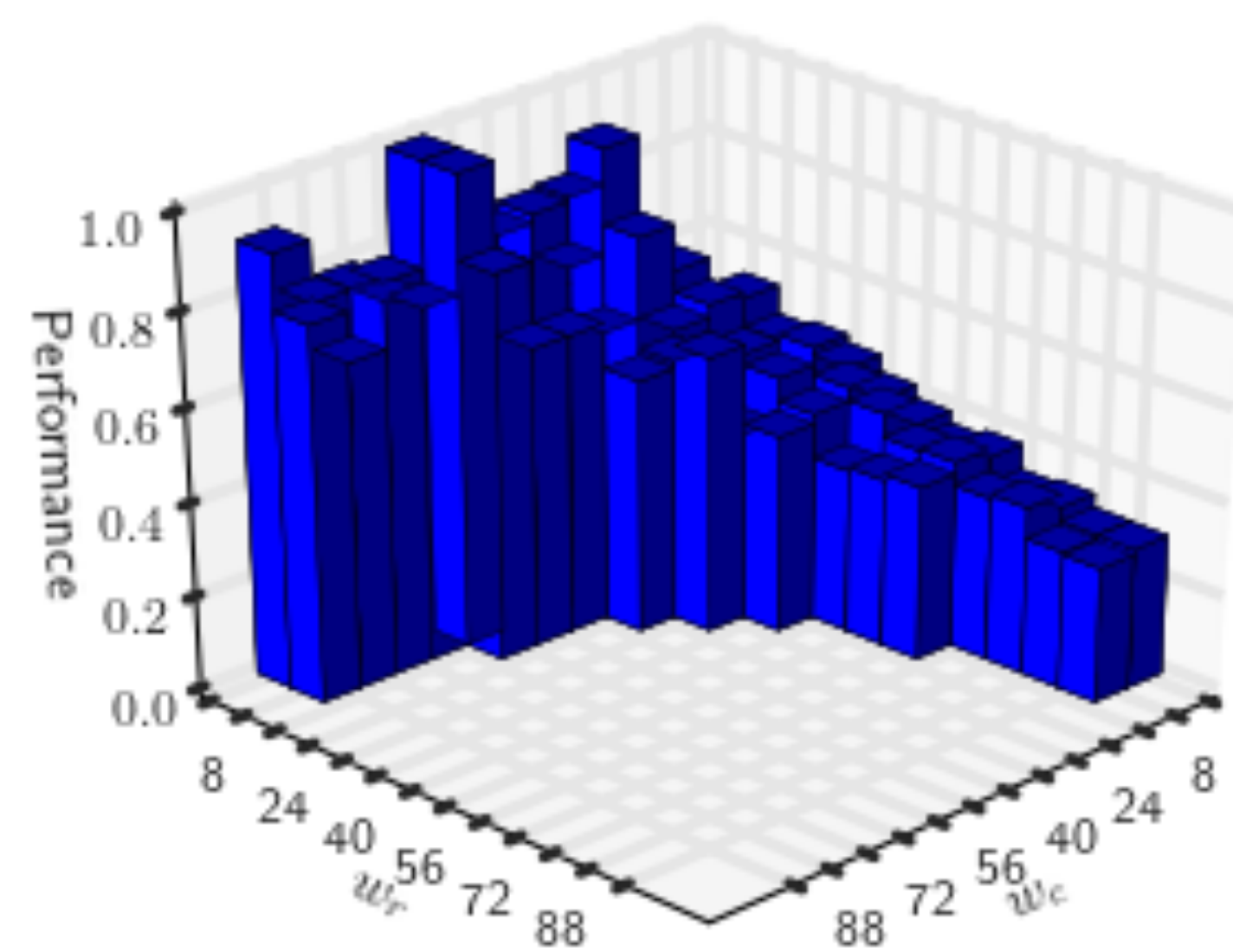
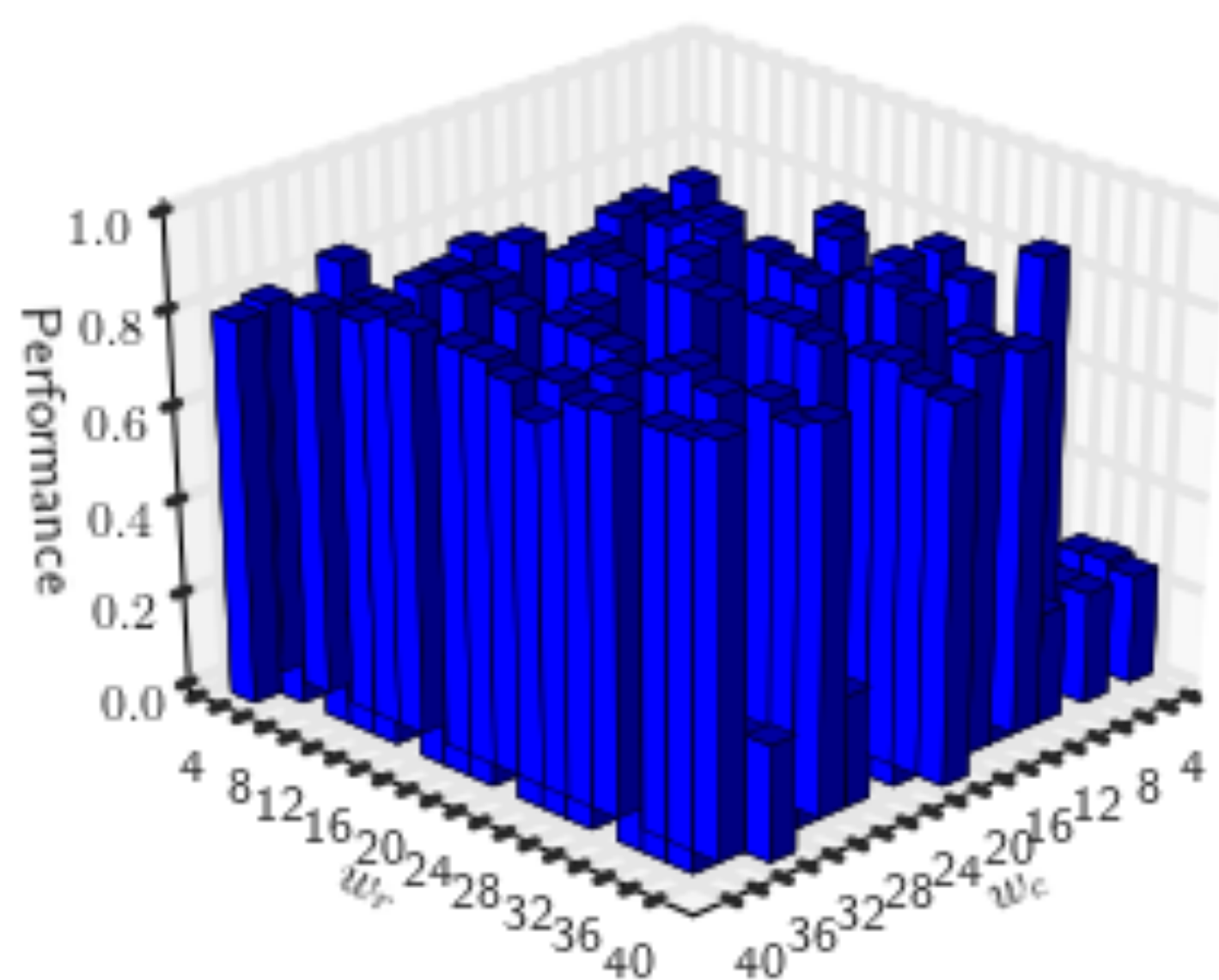


Same stencil!
Different device!



Same device!
Different stencil!





choosing workgroup size
depends on:

1. device
2. program
3. dataset

Picking good
parameters

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

*a job for the
interns*

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

*a job for the
interns*

*/ iterative
compilation*

BAD!

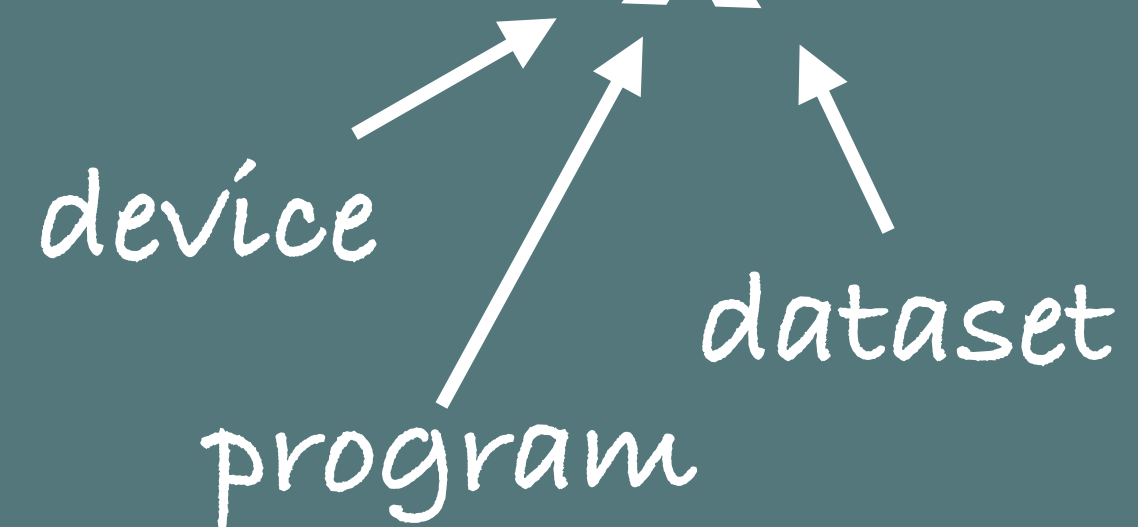
Takes a looooong time

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”



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

device

program

dataset

BETTER!

Still takes a looooong time

Can make *predictions* on unseen “x”



A diagram with three labels: 'device', 'program', and 'dataset'. Arrows point from 'device' and 'program' to a central point, and an arrow points from 'dataset' to a point slightly to the right. These two points then have arrows pointing to the 'x' in the text '“x”'.

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**

A new approach ...

OmniTune

1. Allows *collaborative* performance tuning

Reduce training costs ✓

2. Provides re-usable *implementations*

Reduce implementation costs ✓

3. Provides lightweight runtime interface

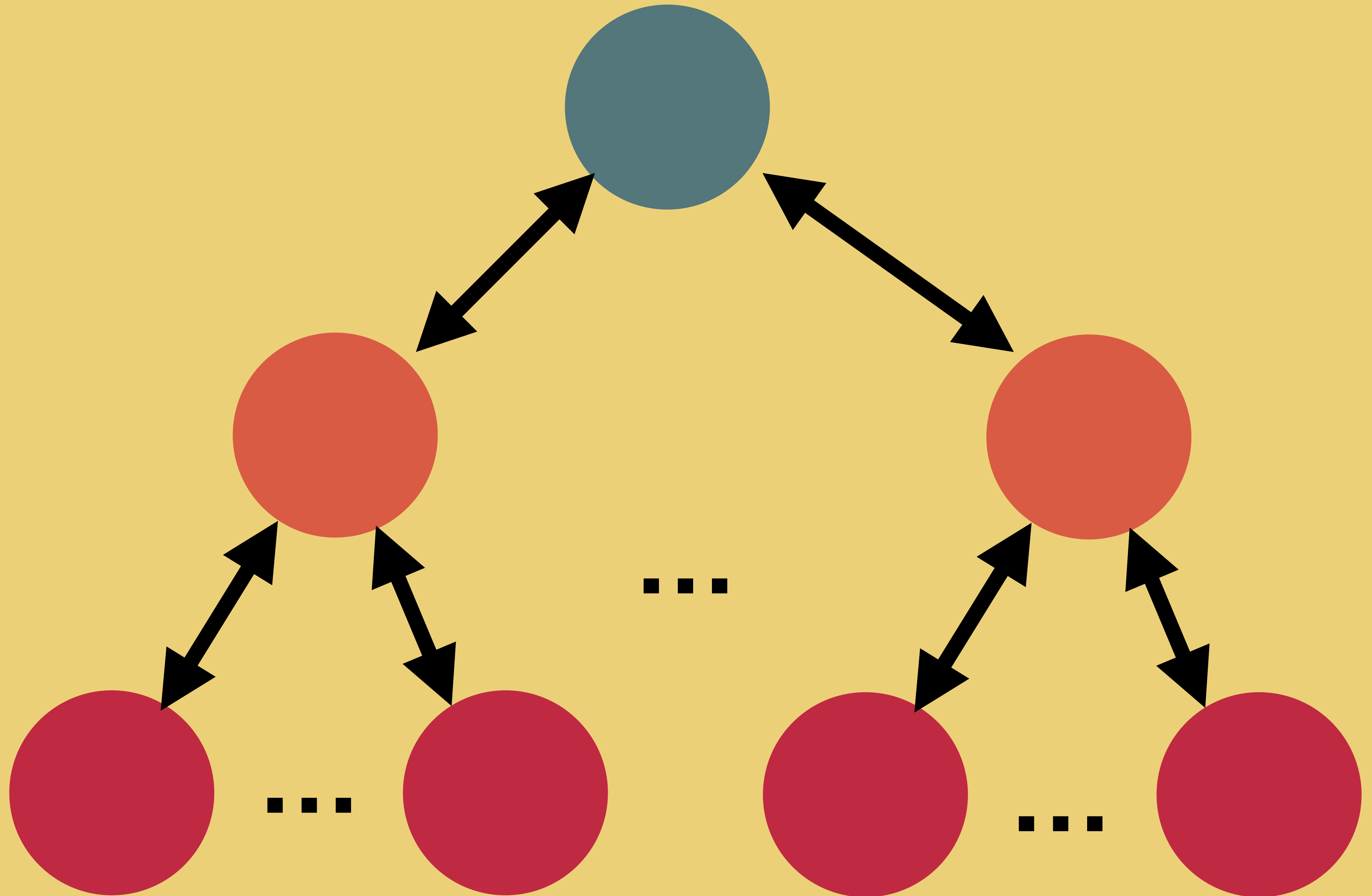
Minimise runtime overheads ✓

How does it work?

Remote

Servers

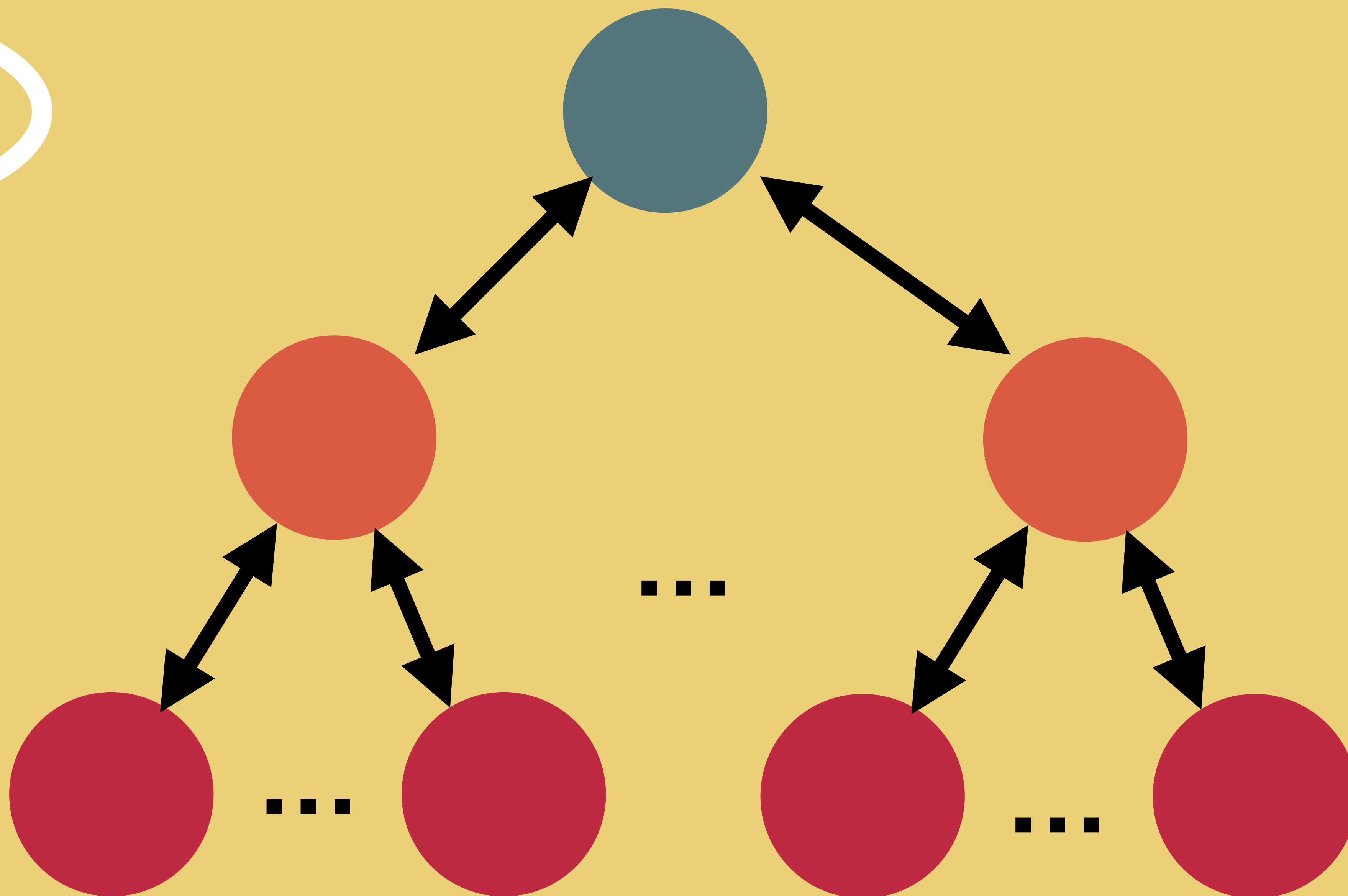
Clients



Remote

Servers

Clients



Remote

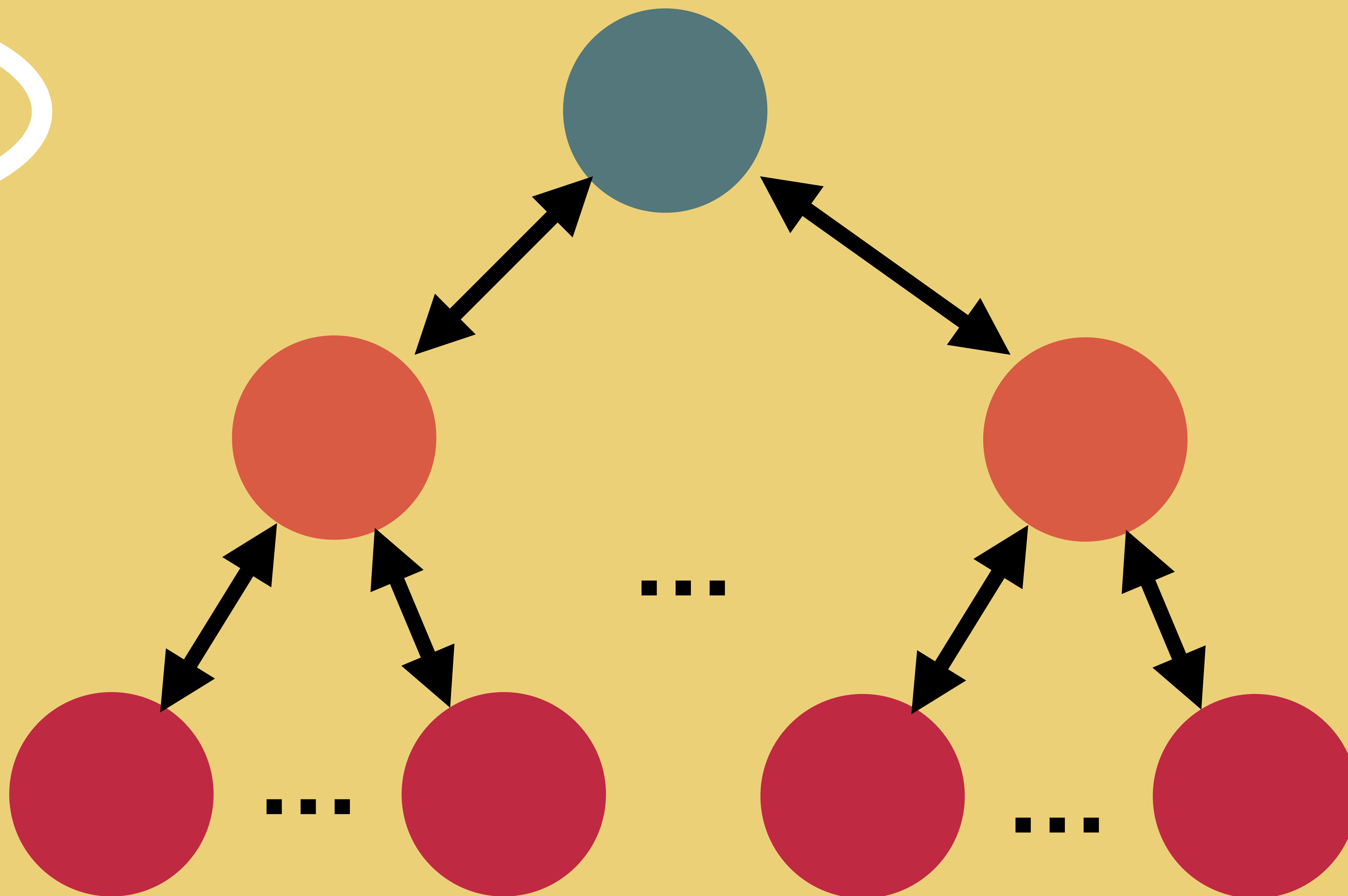
Book-keeper

**Manages and
stores training
data**

Remote

Servers

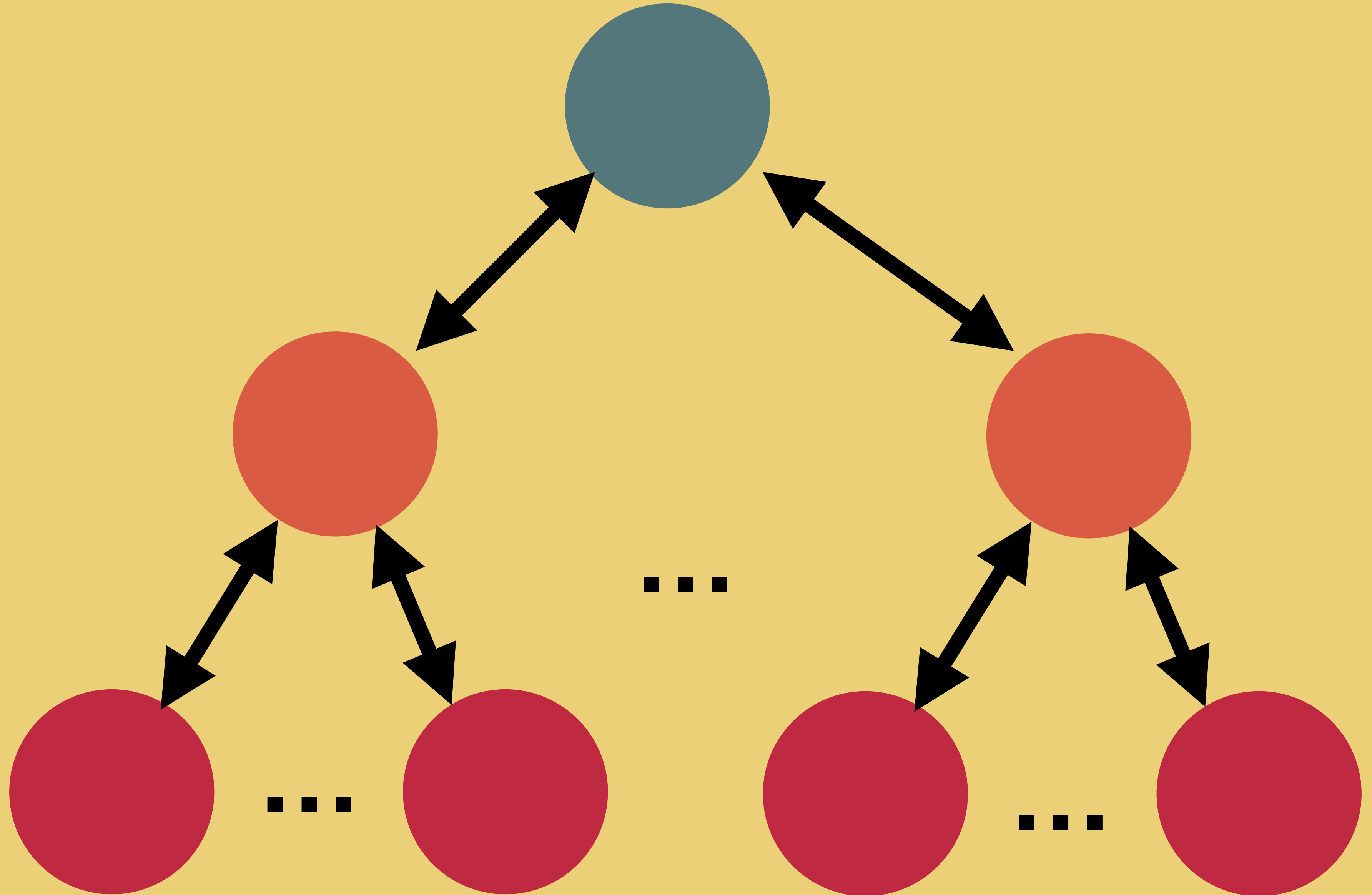
Clients



Remote

Servers

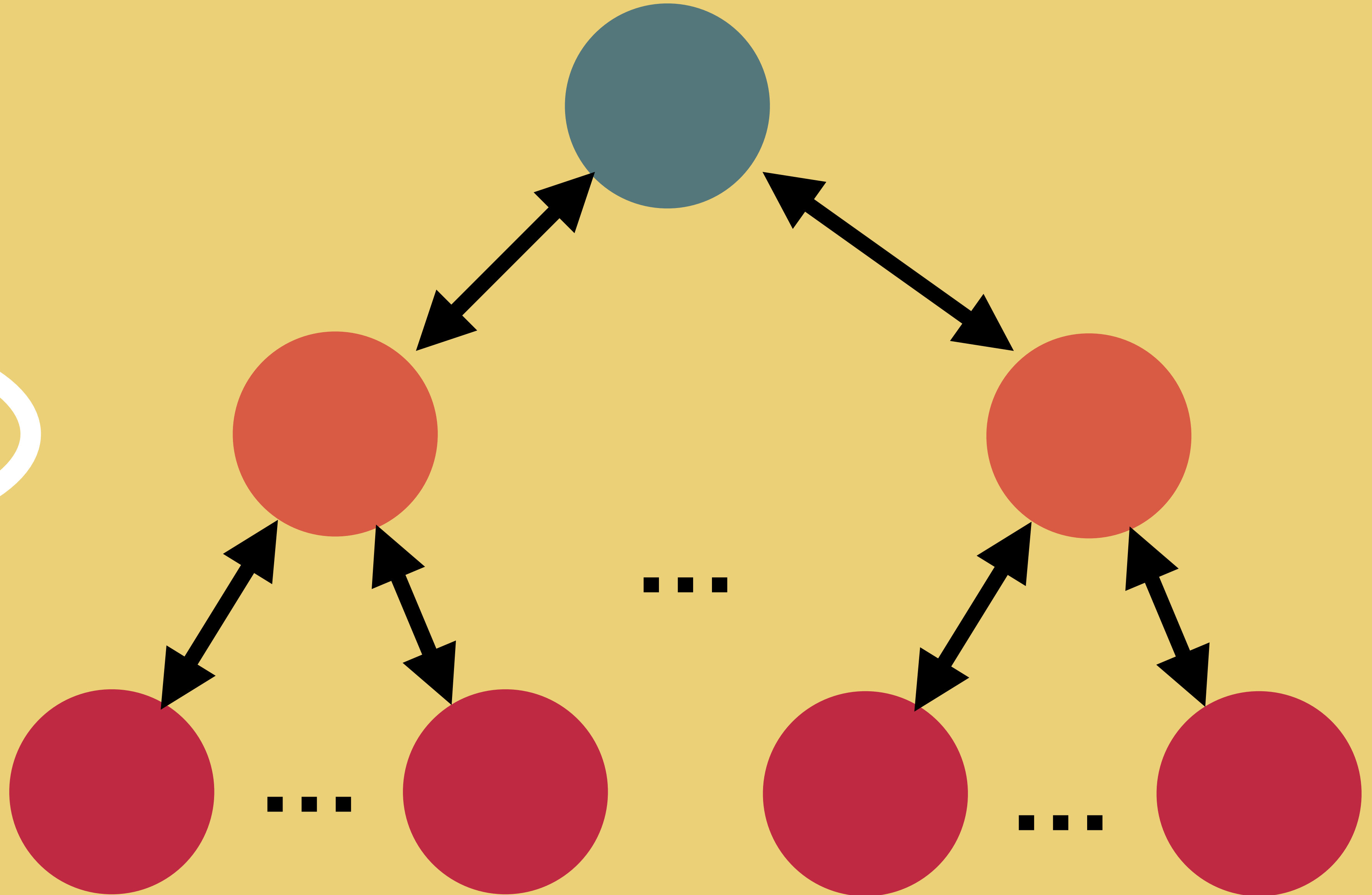
Clients



Remote

Servers

Clients



**Autotuning
engine**



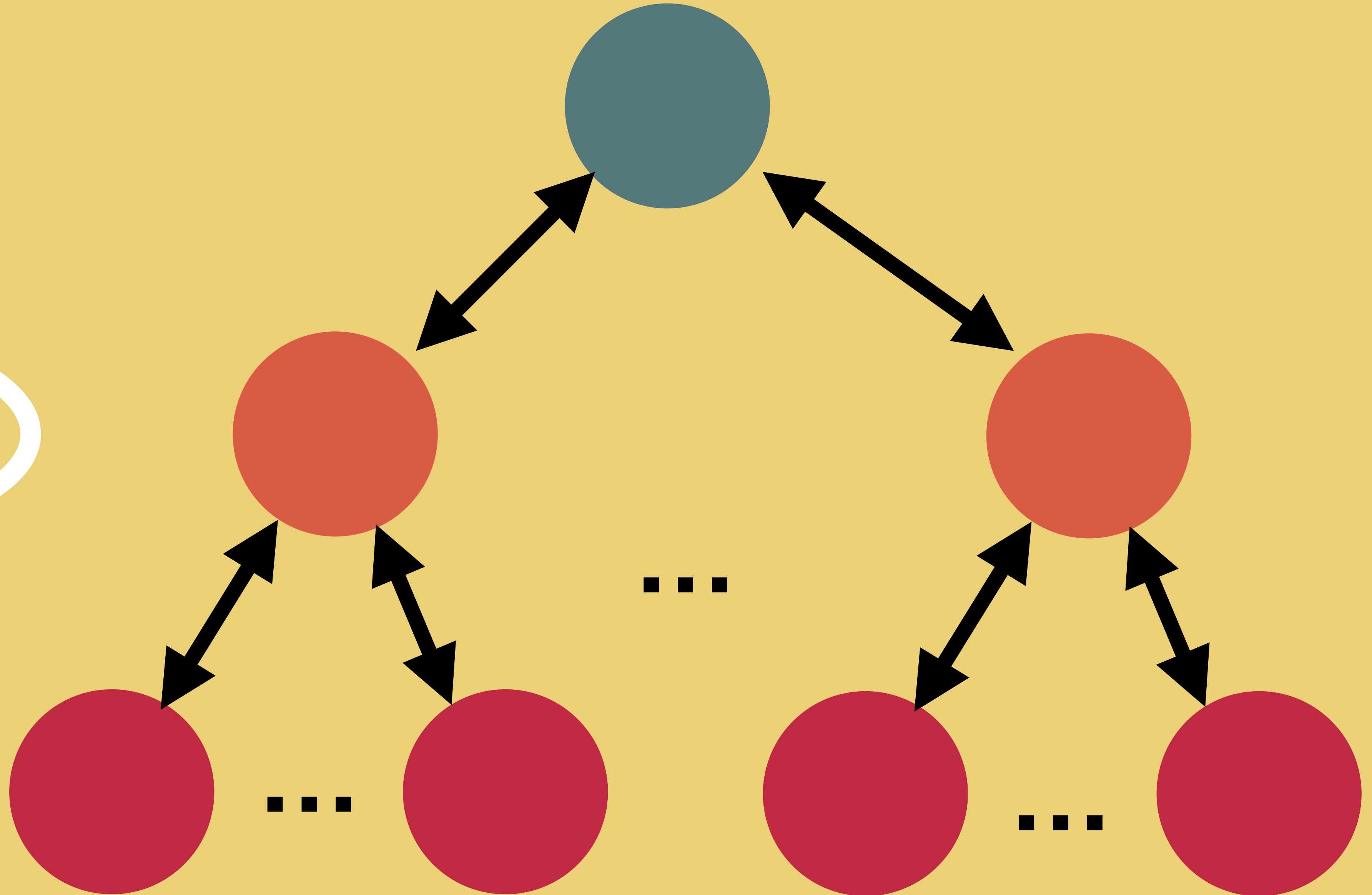
Servers

**Performs
machine learning**

Remote

Servers

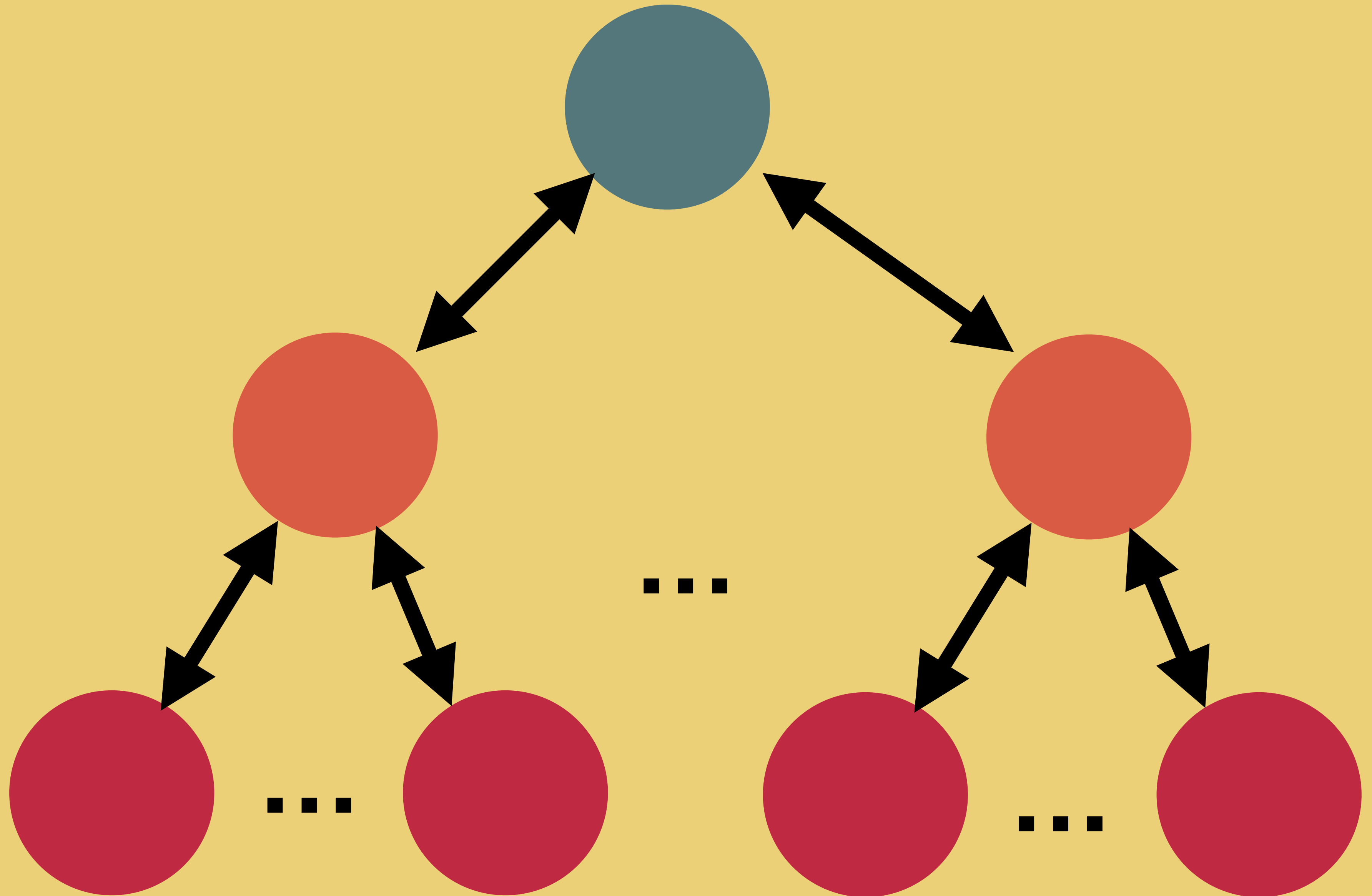
Clients



Remote

Servers

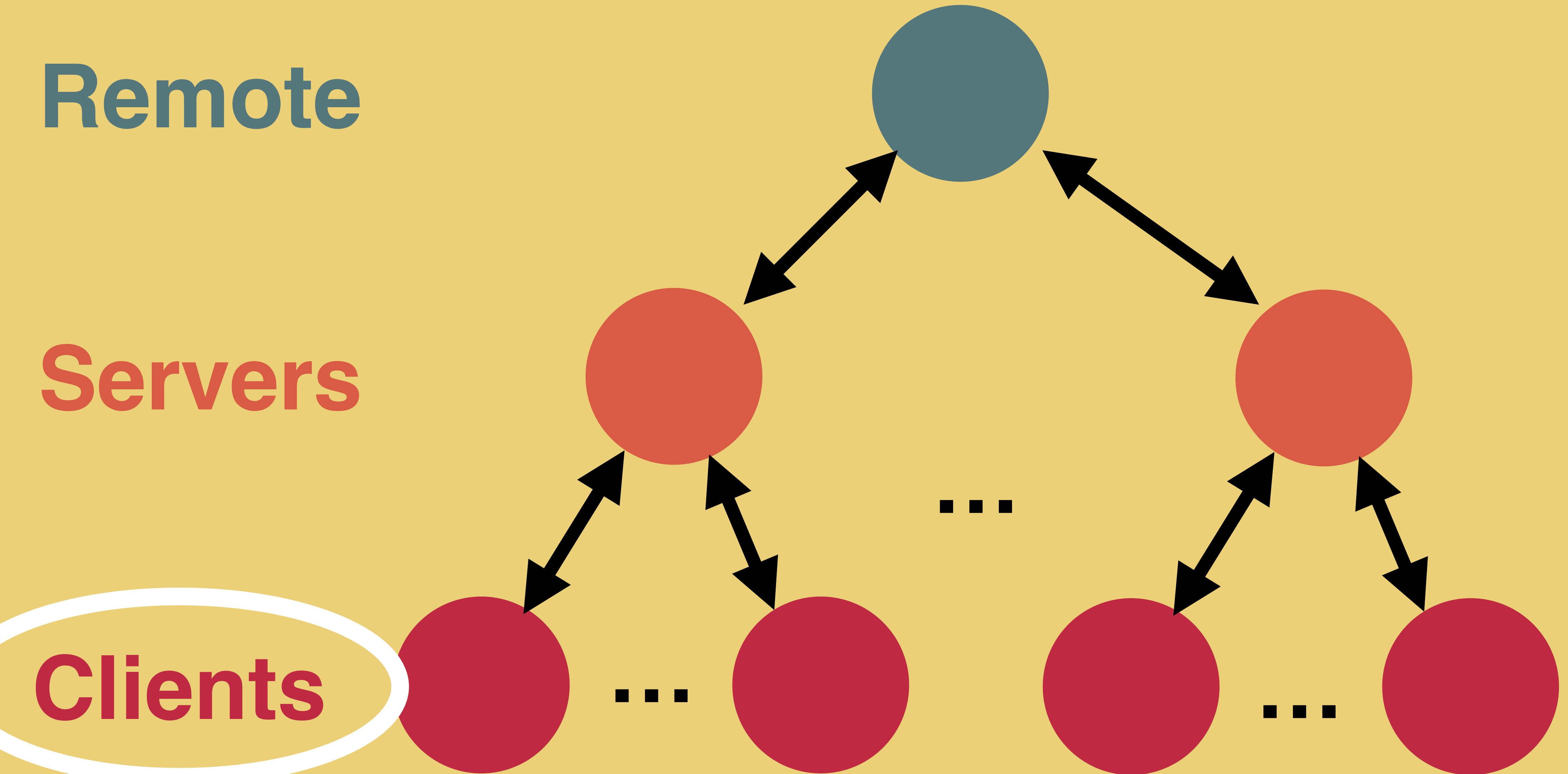
Clients



Remote

Servers

Clients



**Target
applications**

**Programs we
want to tune**

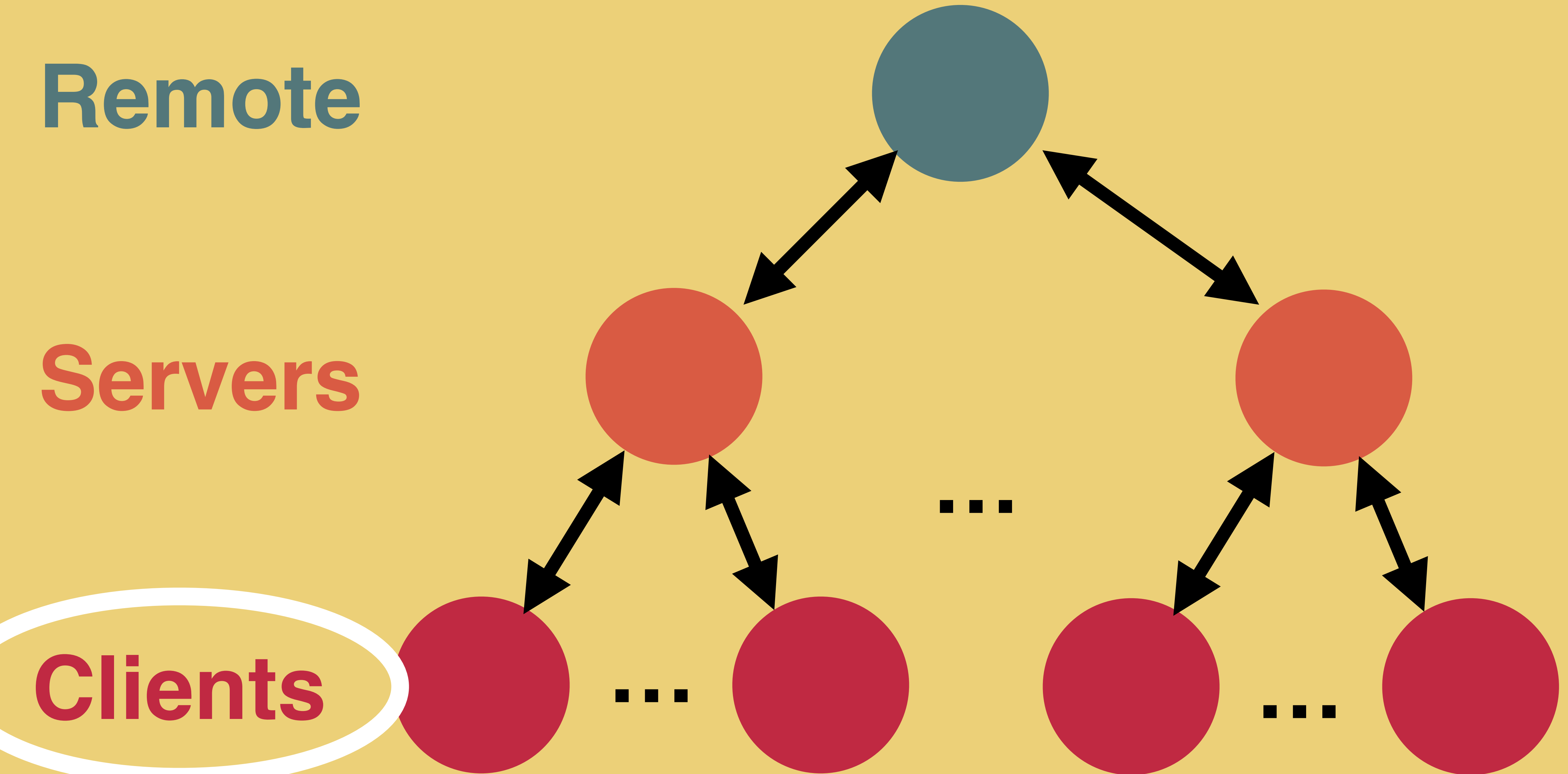


Clients

Remote

Servers

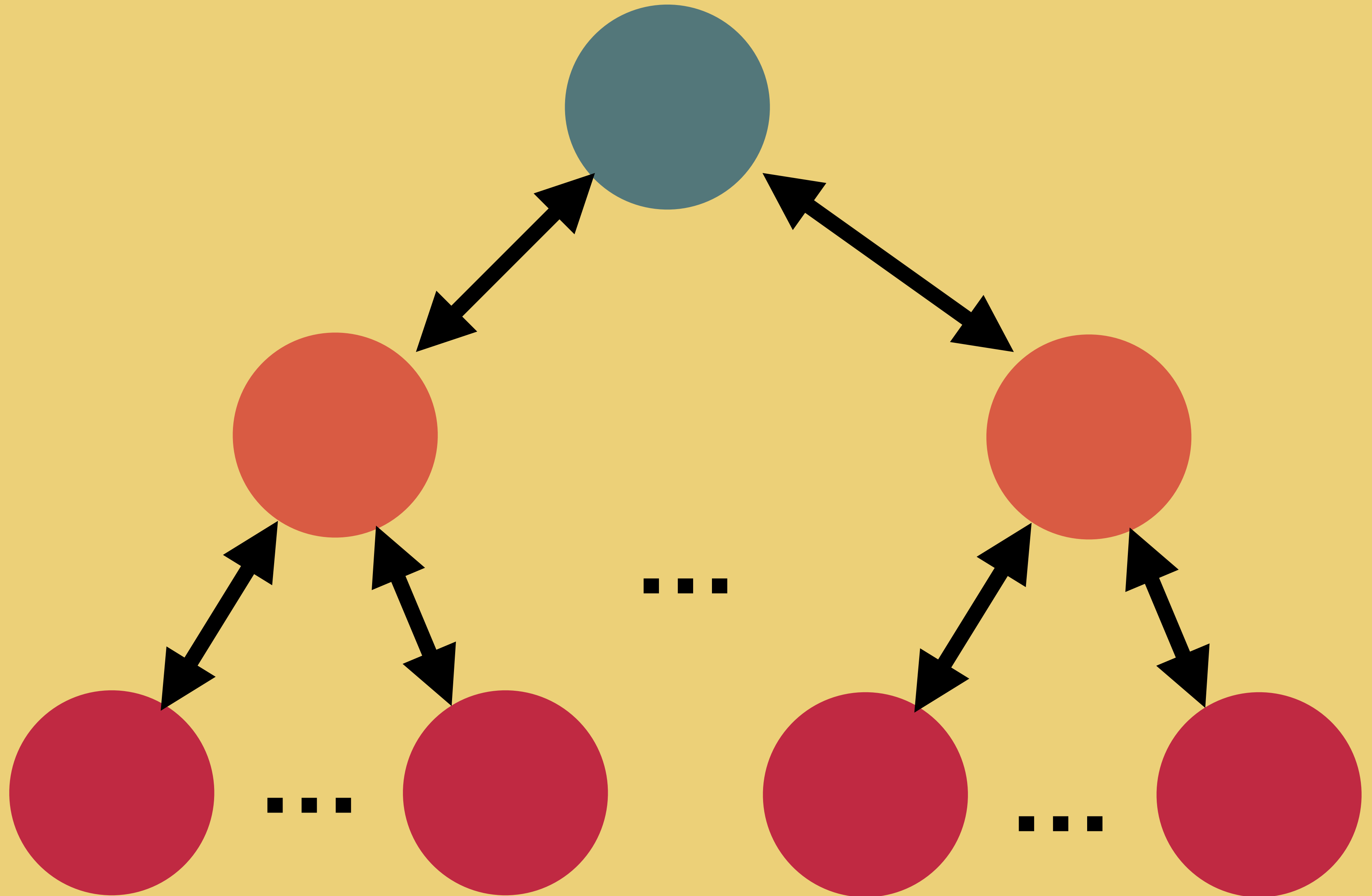
Clients



Remote

Servers

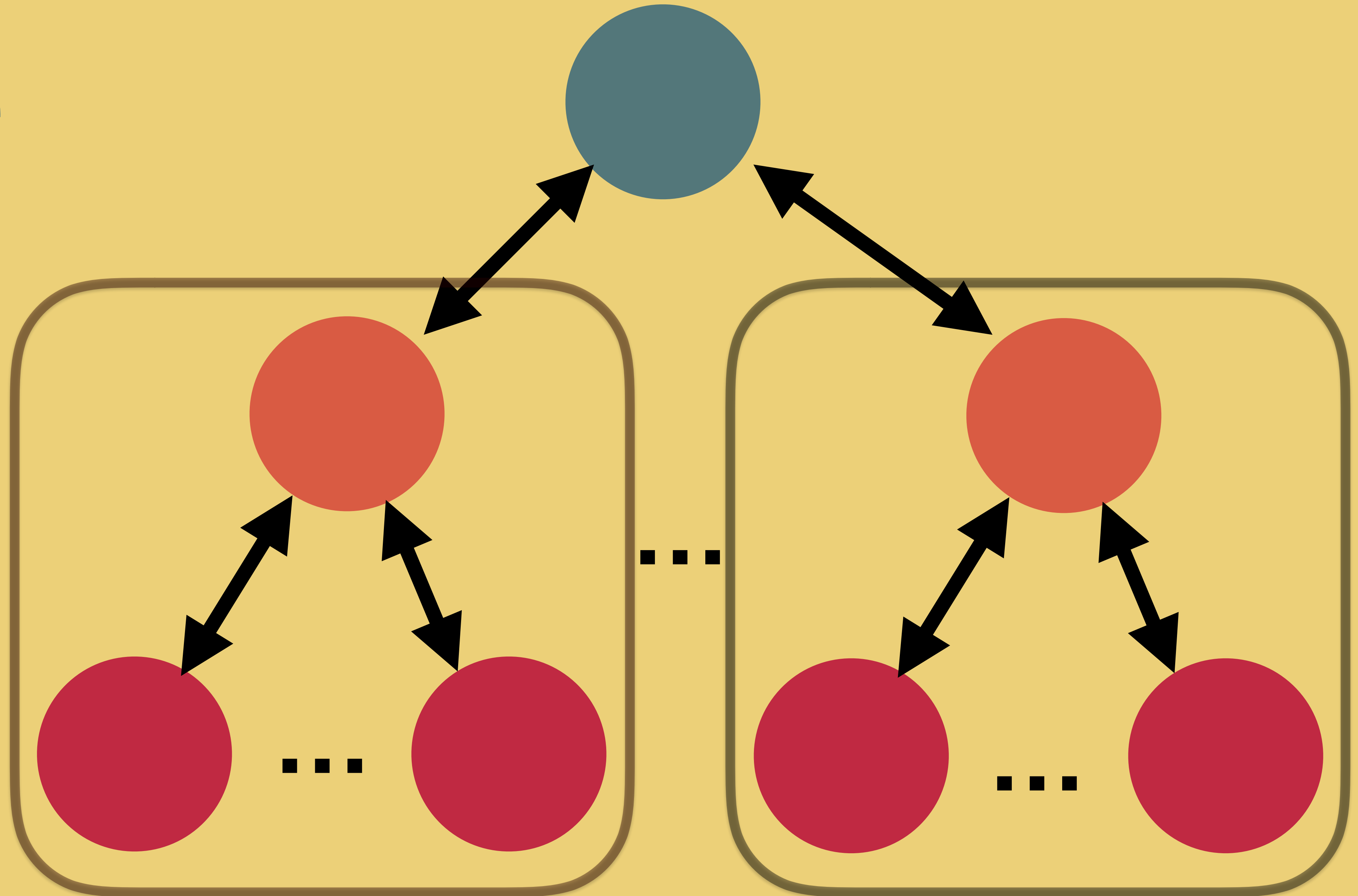
Clients



Remote

Servers

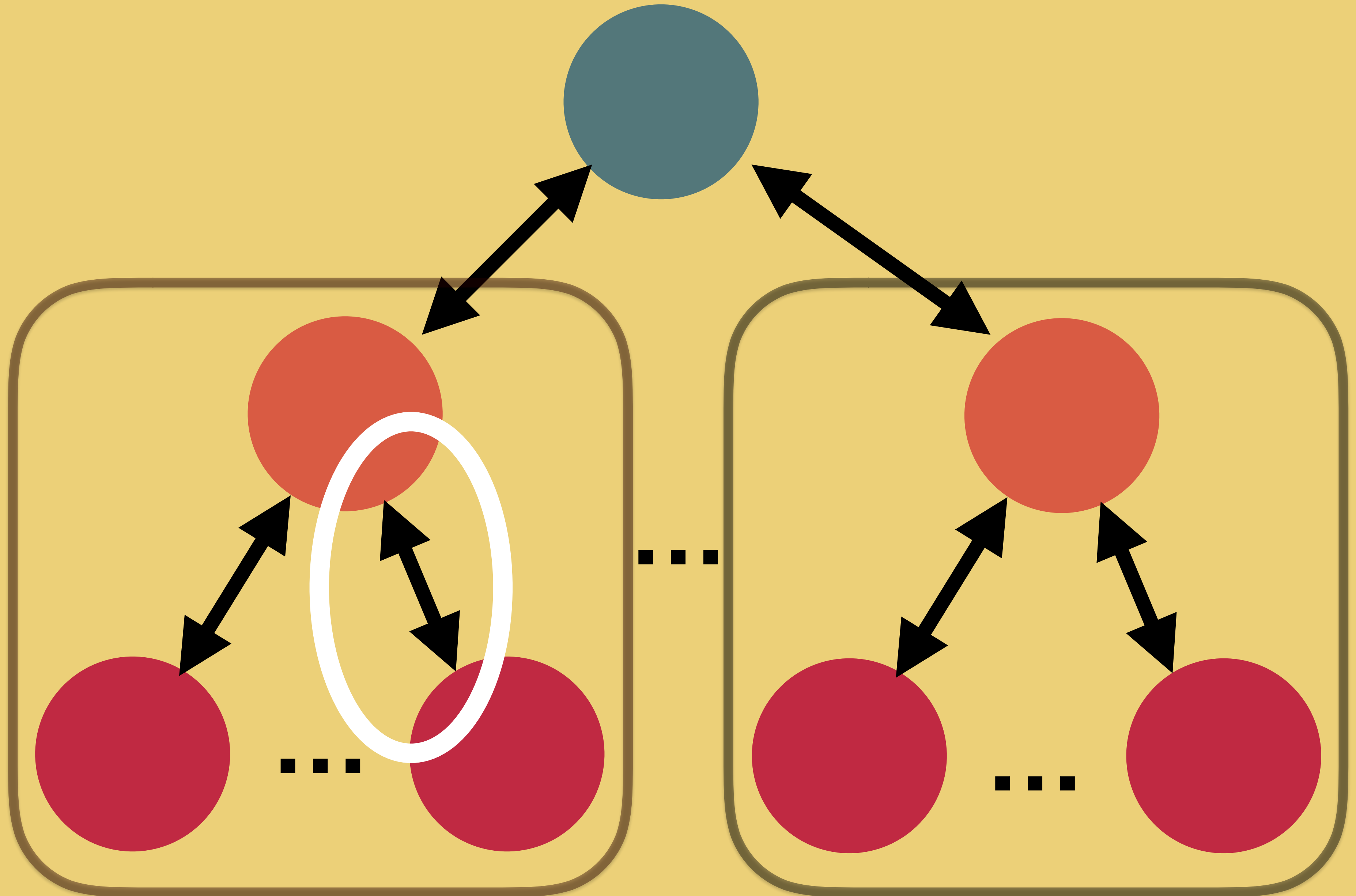
Clients

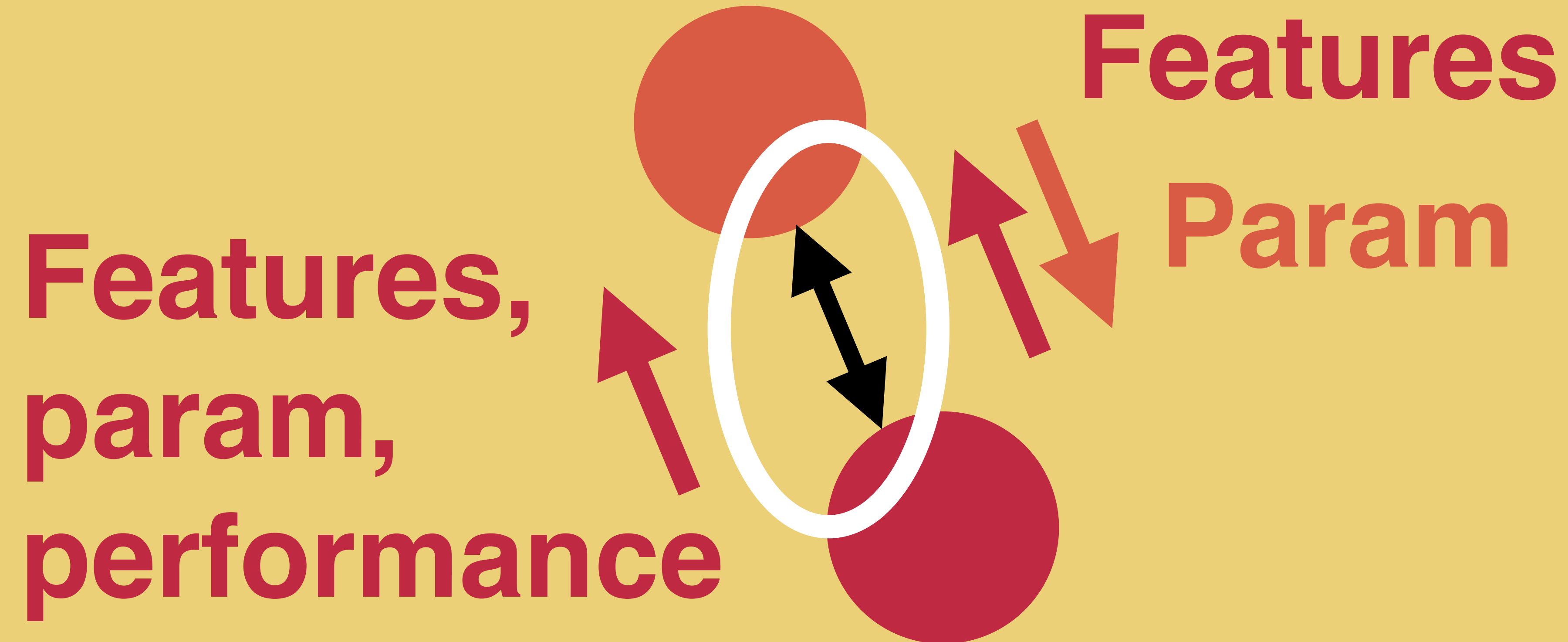


Remote

Servers

Clients

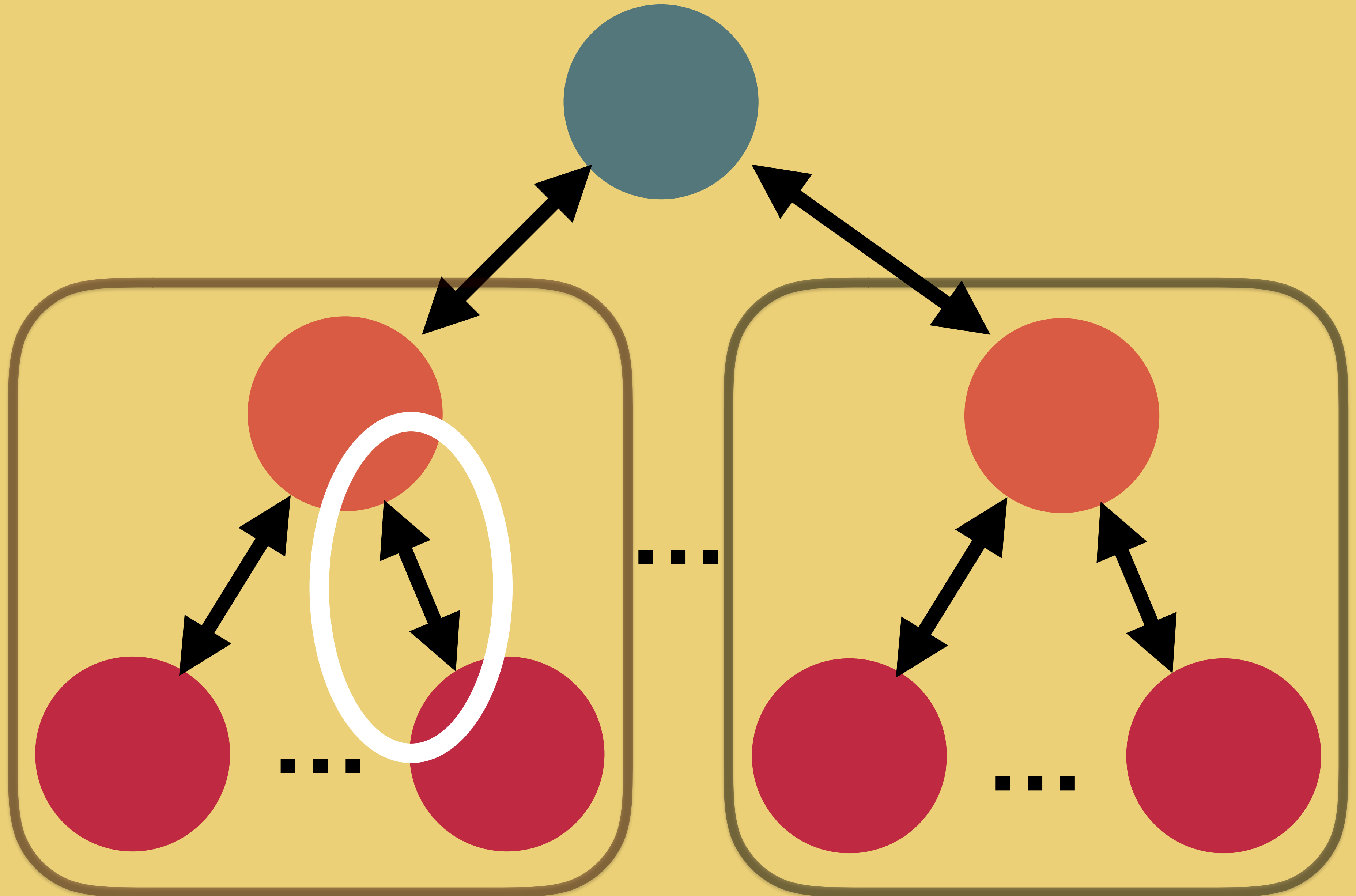




Remote

Servers

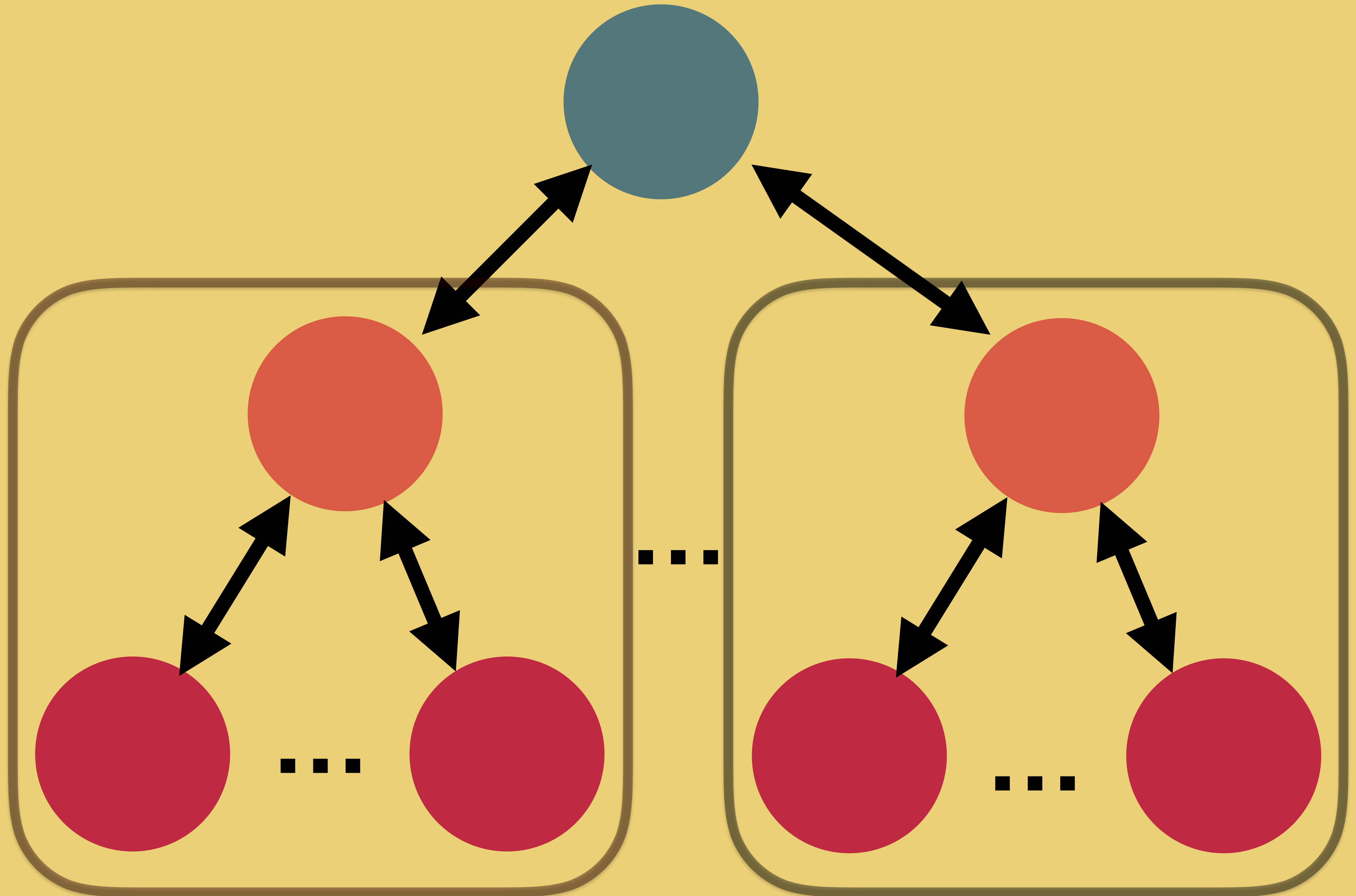
Clients



Remote

Servers

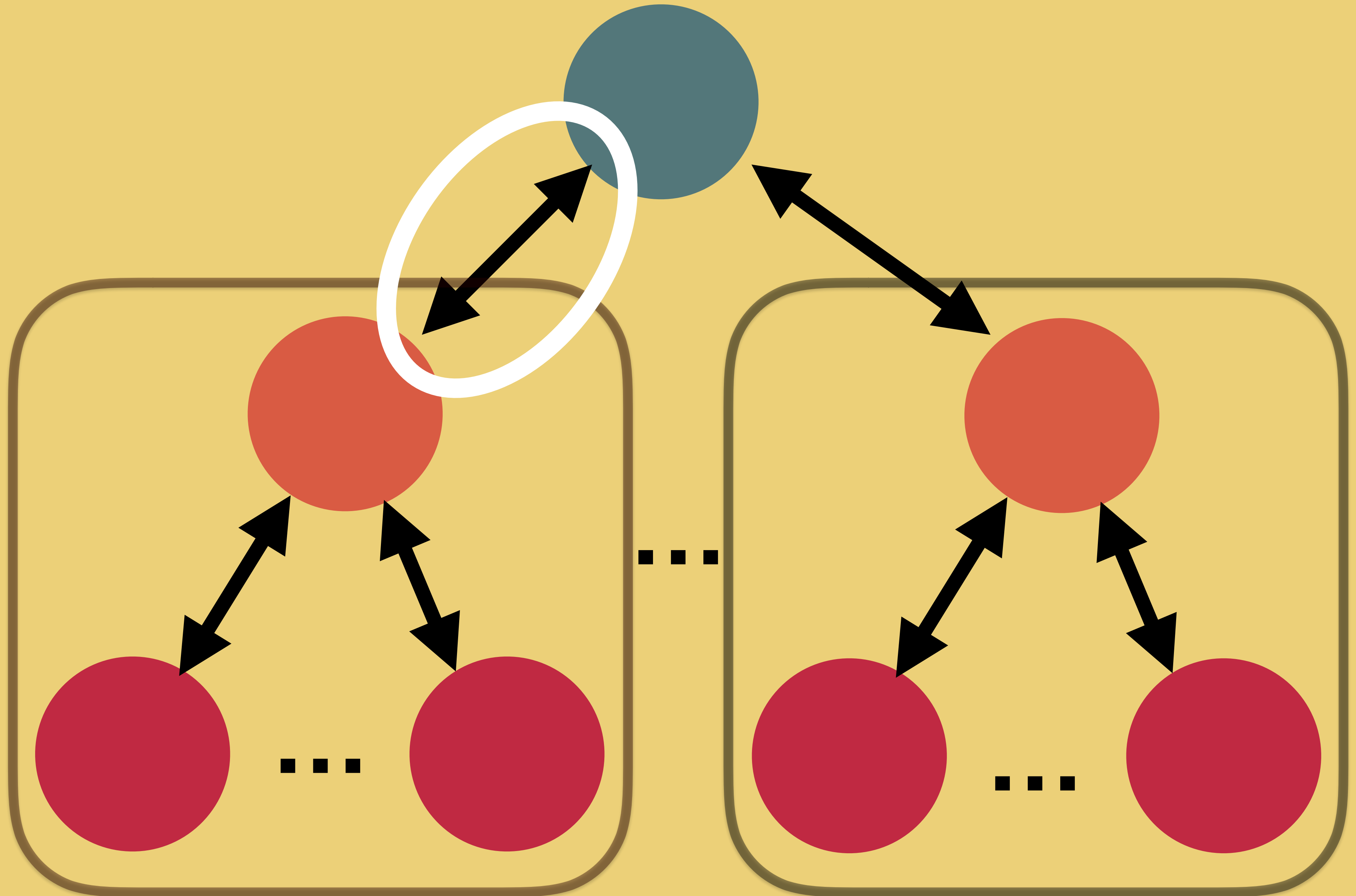
Clients



Remote

Servers

Clients



New training data

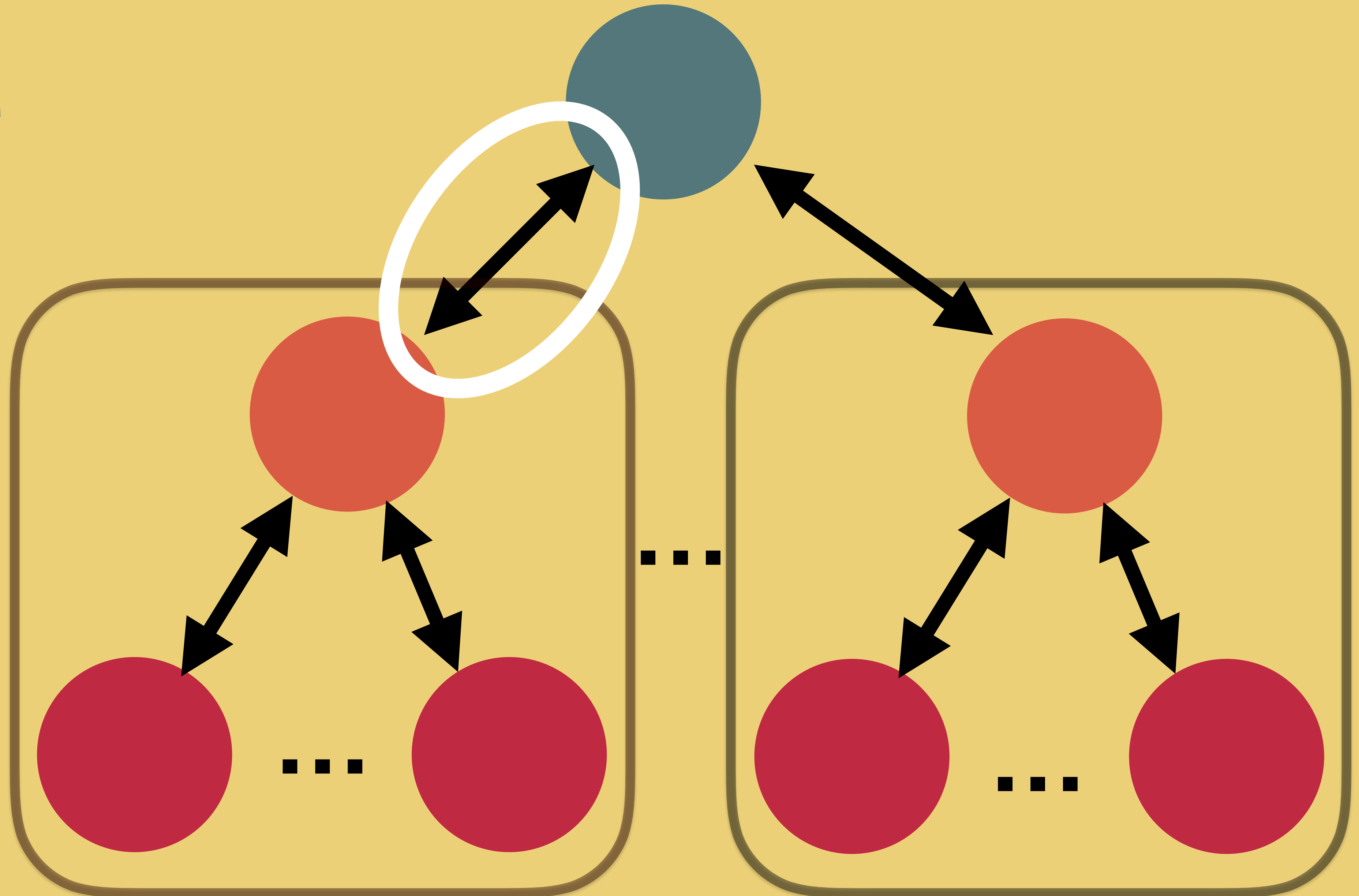


New training data

Remote

Servers

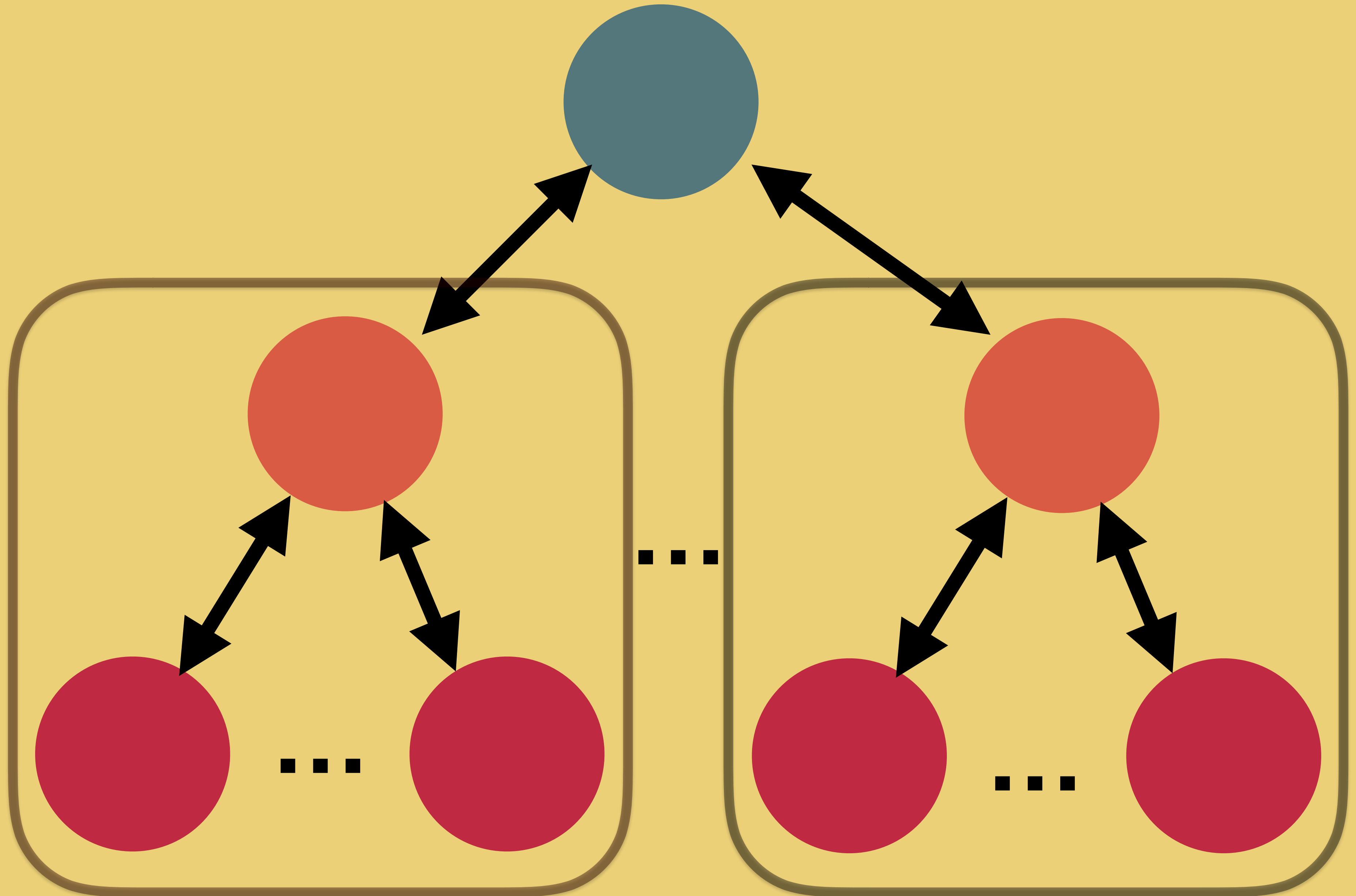
Clients



Remote

Servers

Clients

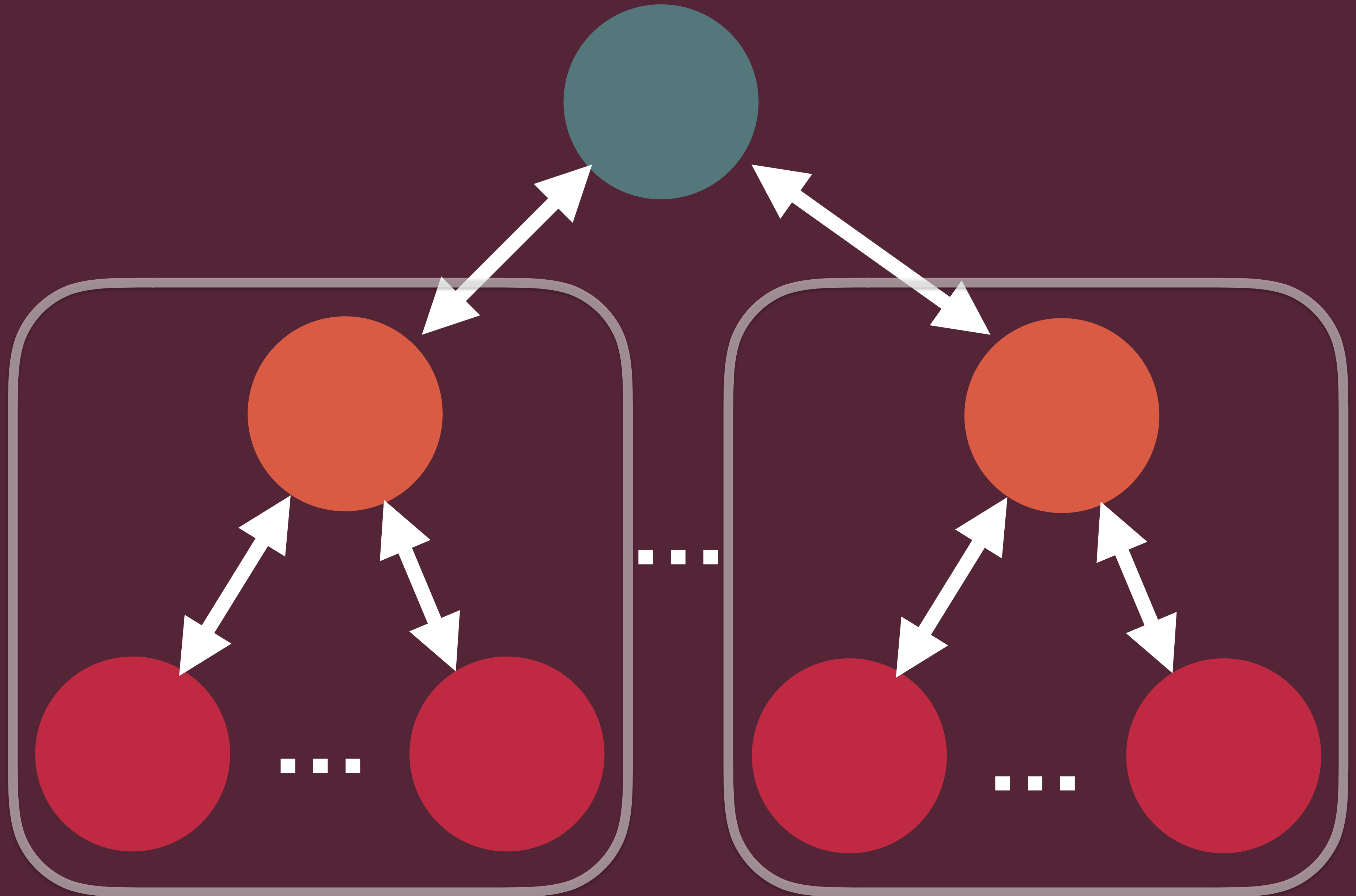


Implementation

Remote

Servers

Clients



AWS + MySQL

TCP/IP

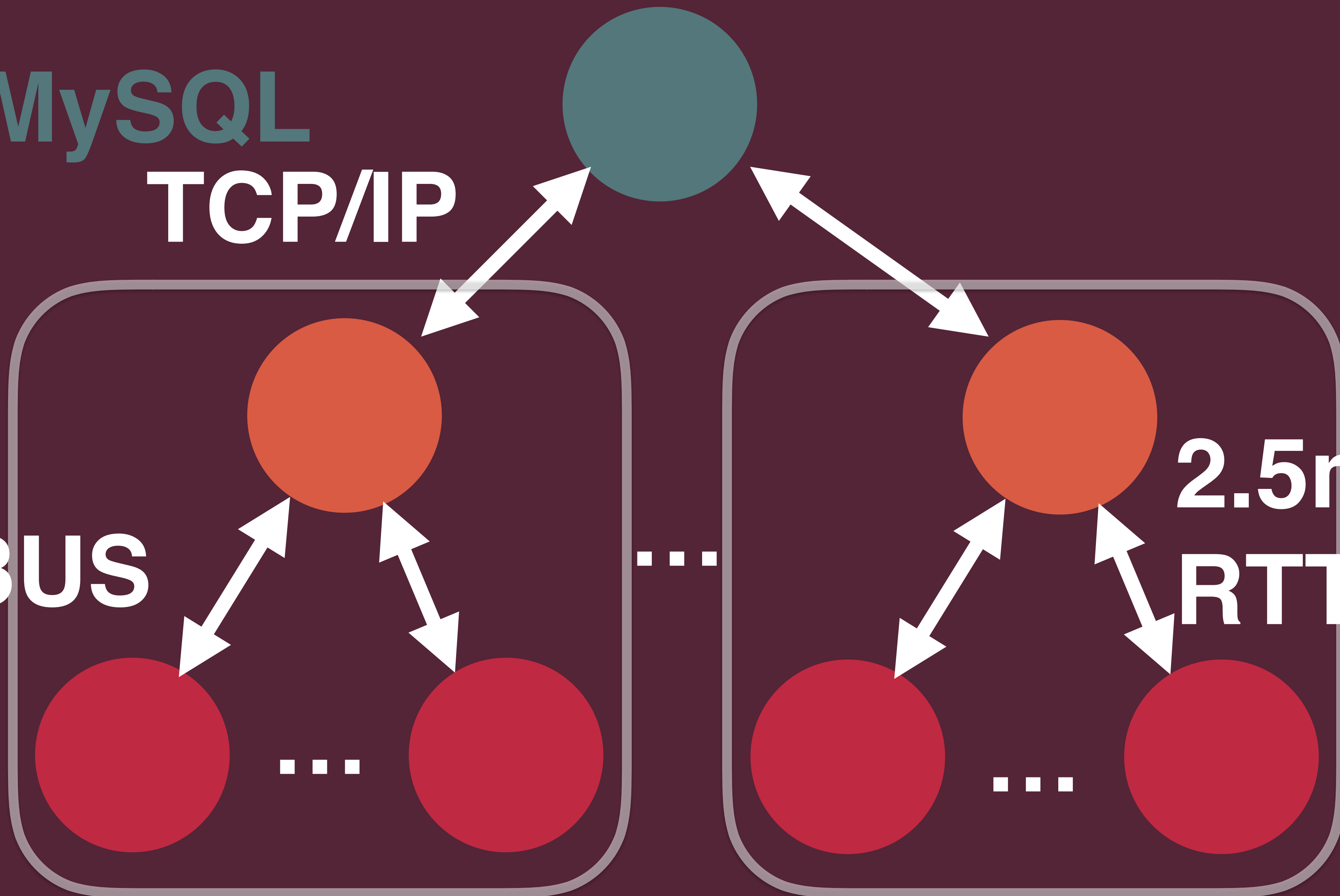
Python

DBus

SkelCL

2.5ms

RTT



Experiment

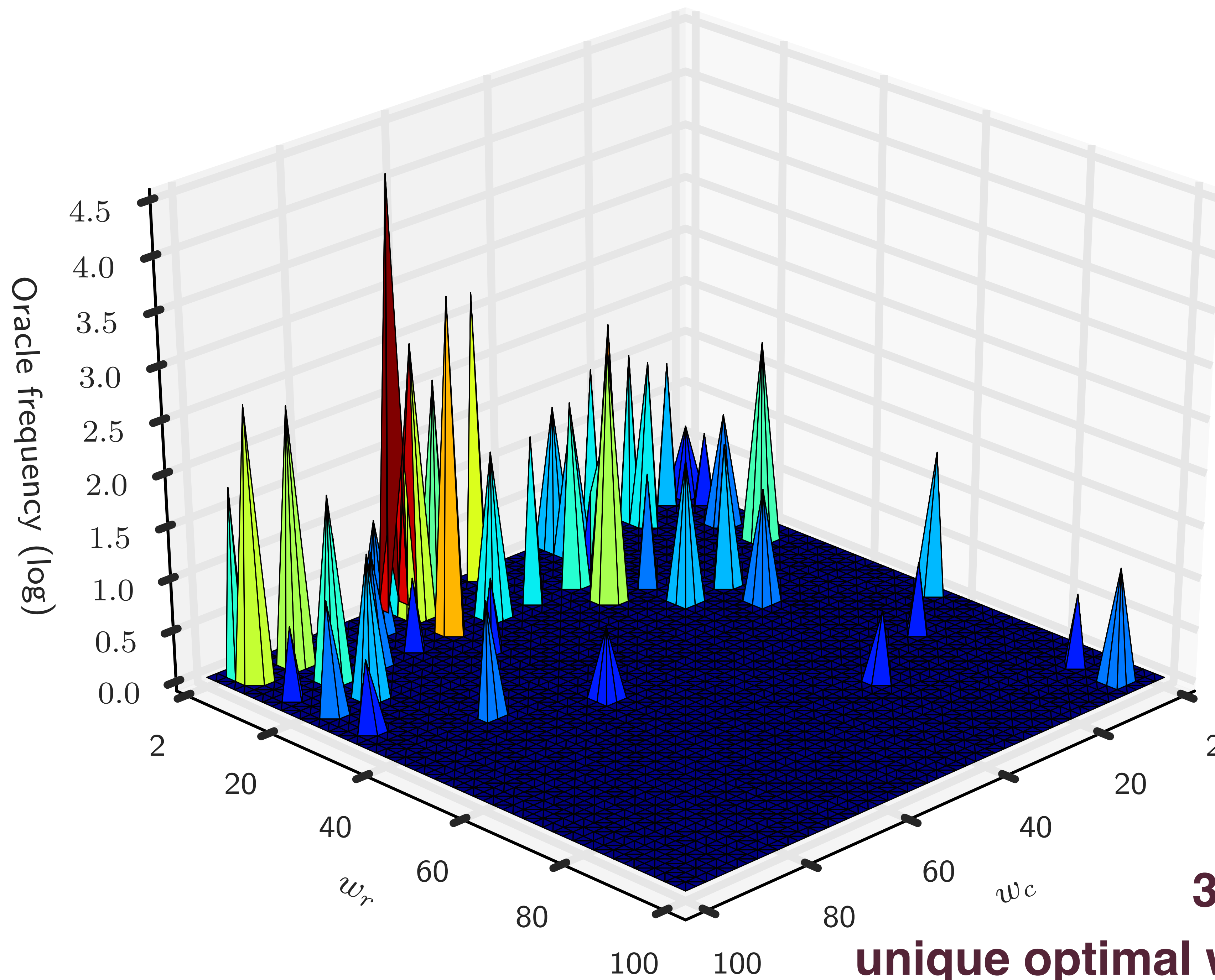
Setup:

6 stencil benchmarks + synthetic.

7 different GPUs & CPUs.

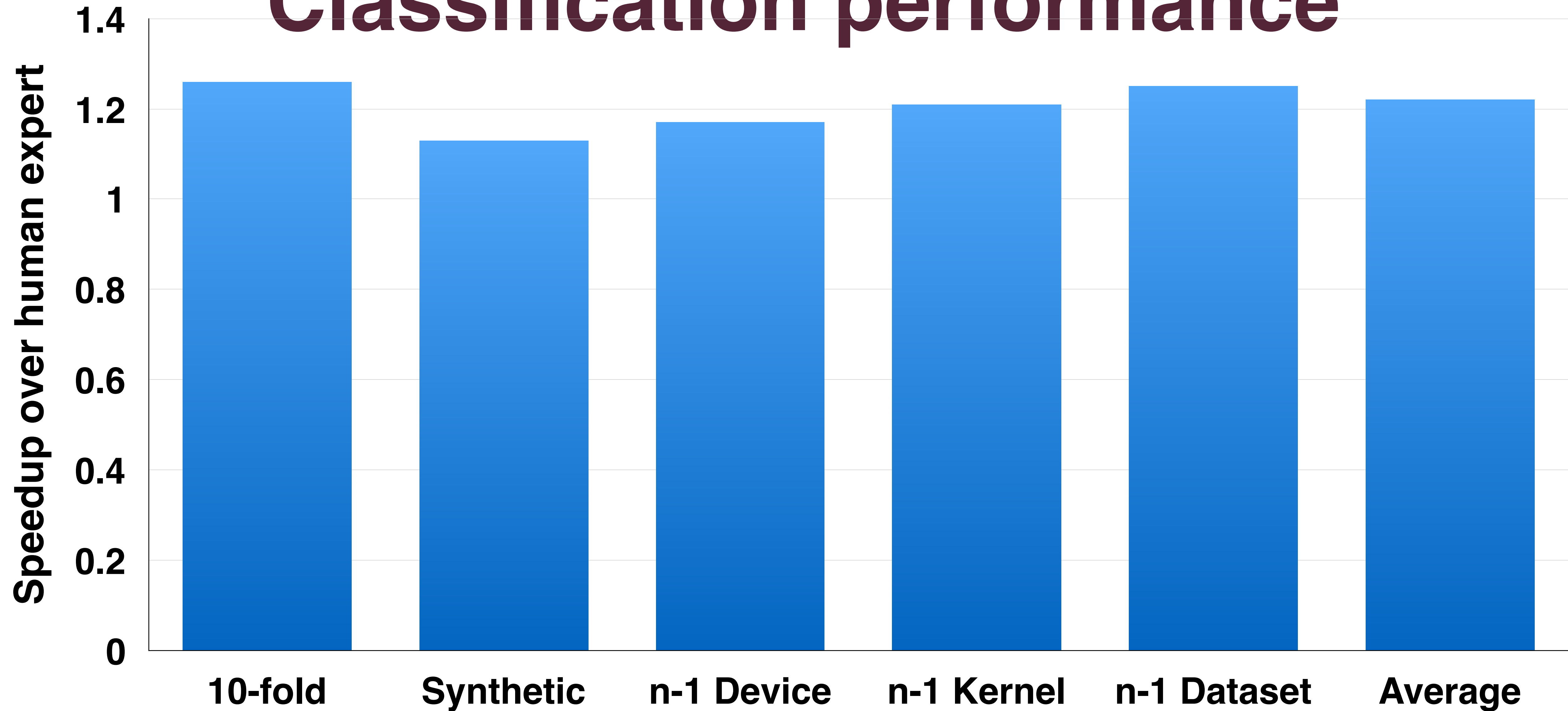
4 dataset sizes.

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

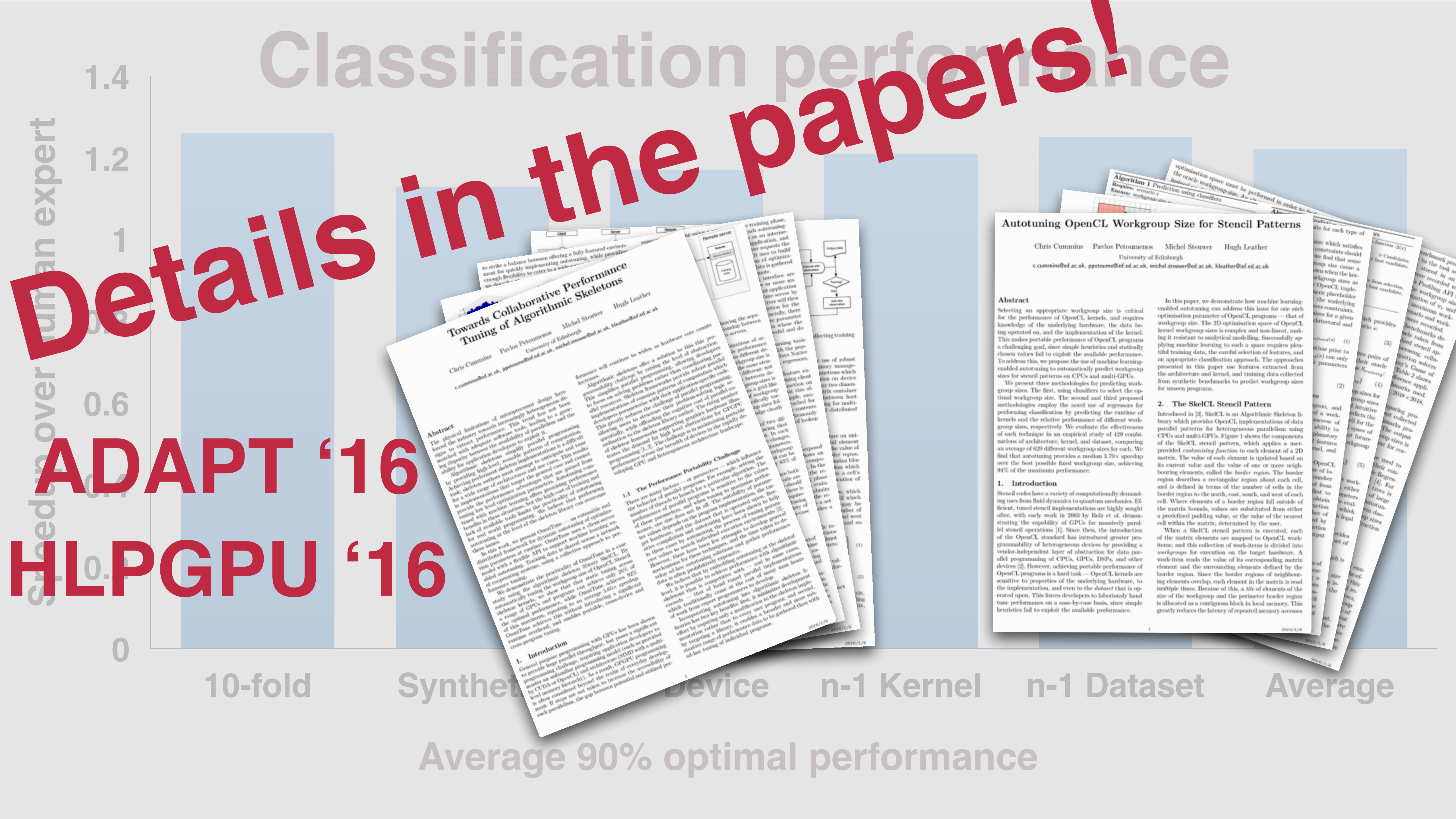


**32% cases have
unique optimal workgroup size**

Classification performance



Average 90% optimal performance



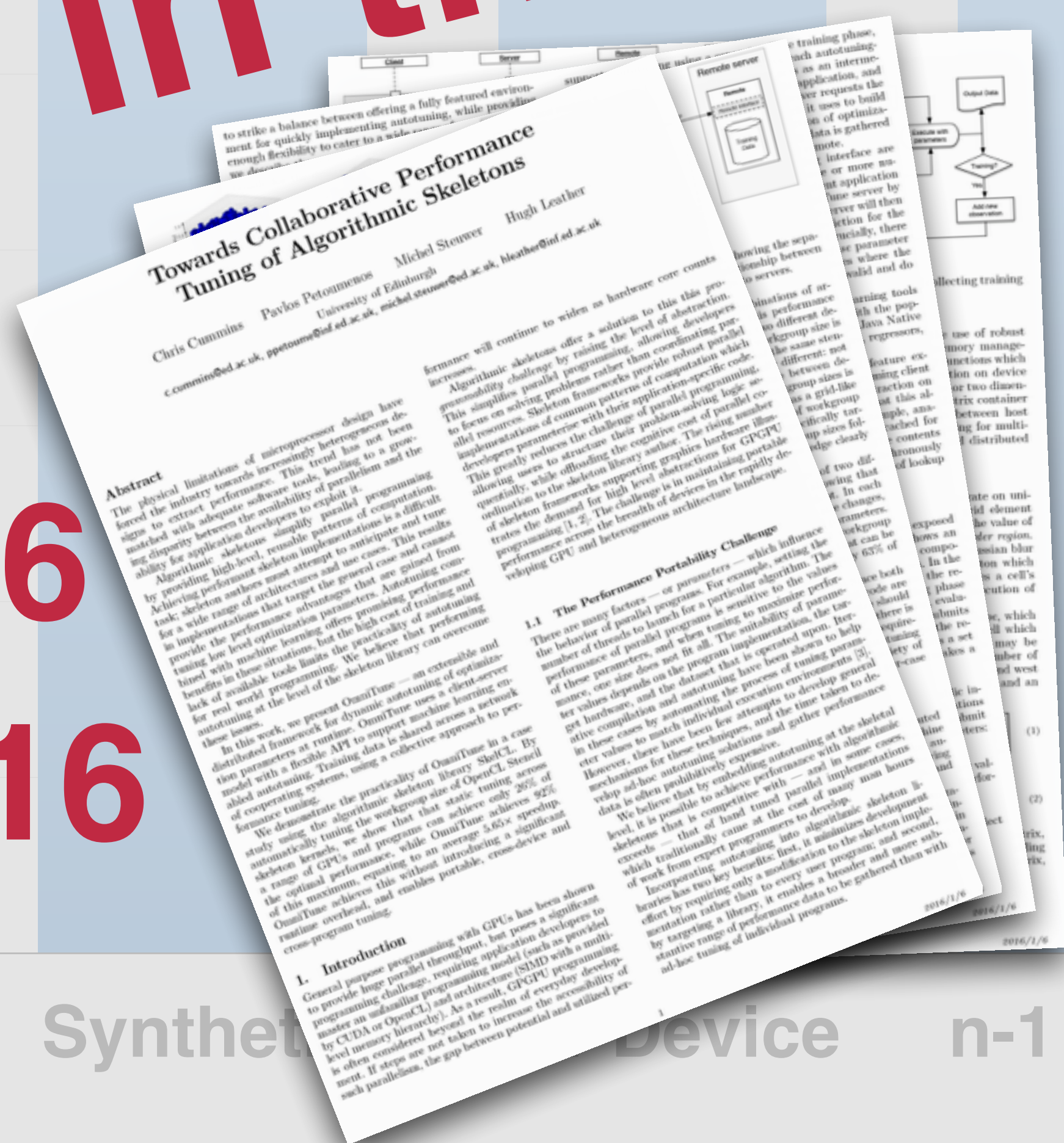
Classification performance

Human expert
Speedup over
CPU

ADAPT '16 HLPGPU '16

10-fold Synthetic Device n-1 Kernel n-1 Dataset Average

Average 90% optimal performance



Conclusions

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

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

**I designed a framework for doing this
using machine learning**

Demonstrated using SkelCL stencils

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



Autotuning and Algorithmic Skeletons