



ceph days
LONDON 2026



Cuttle out the middleman

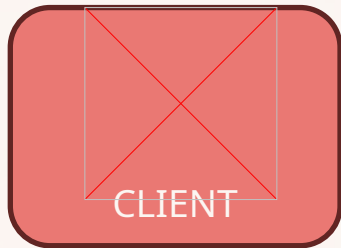
EC Direct Reads

Alex Ainscow

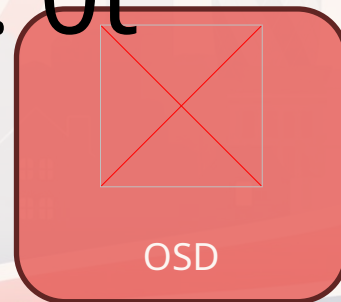
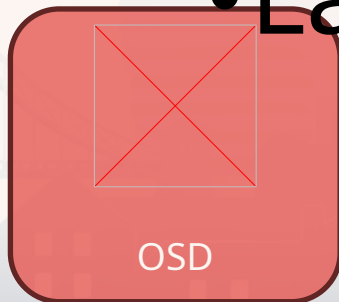
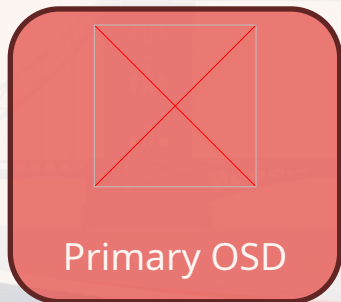




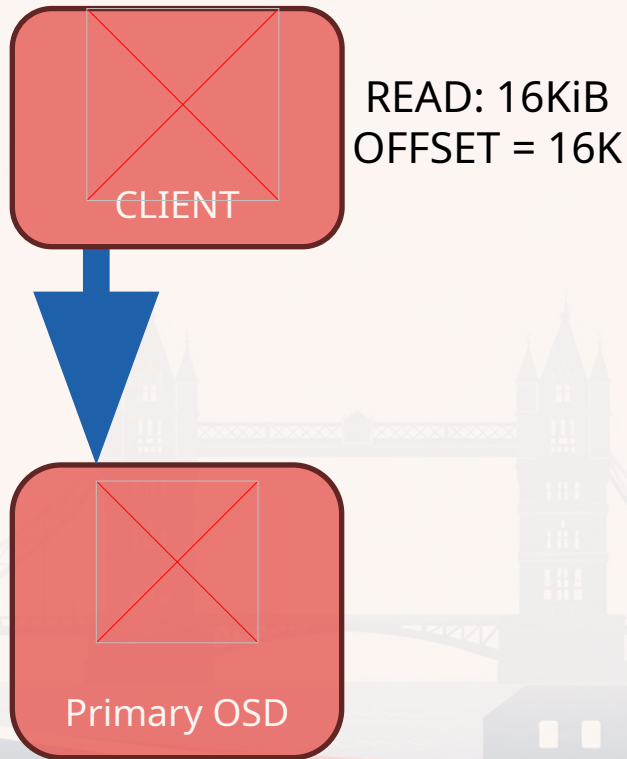
Small Reads



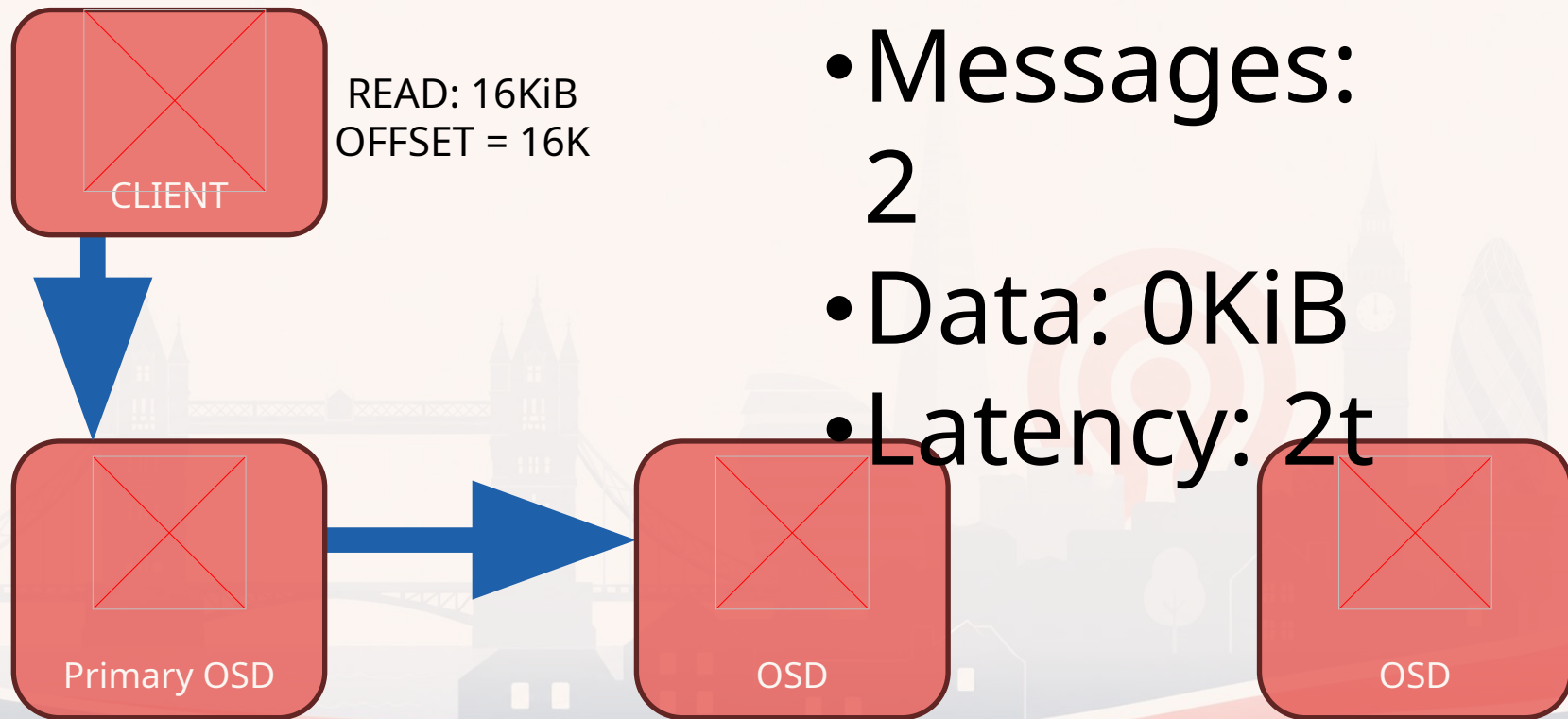
READ: 16KiB
OFFSET = 16K



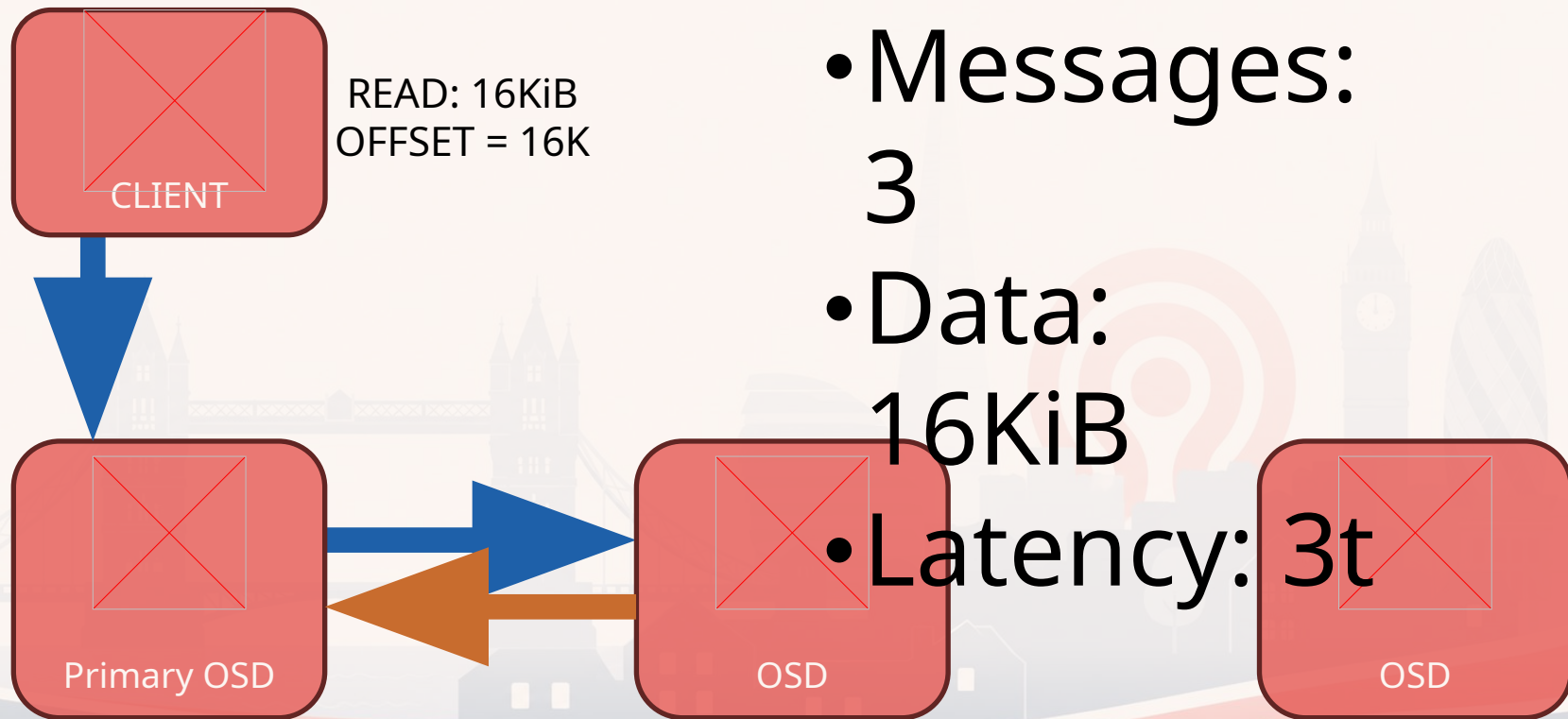
- Messages: 0
- Data: 0KiB
- Latency: 0t

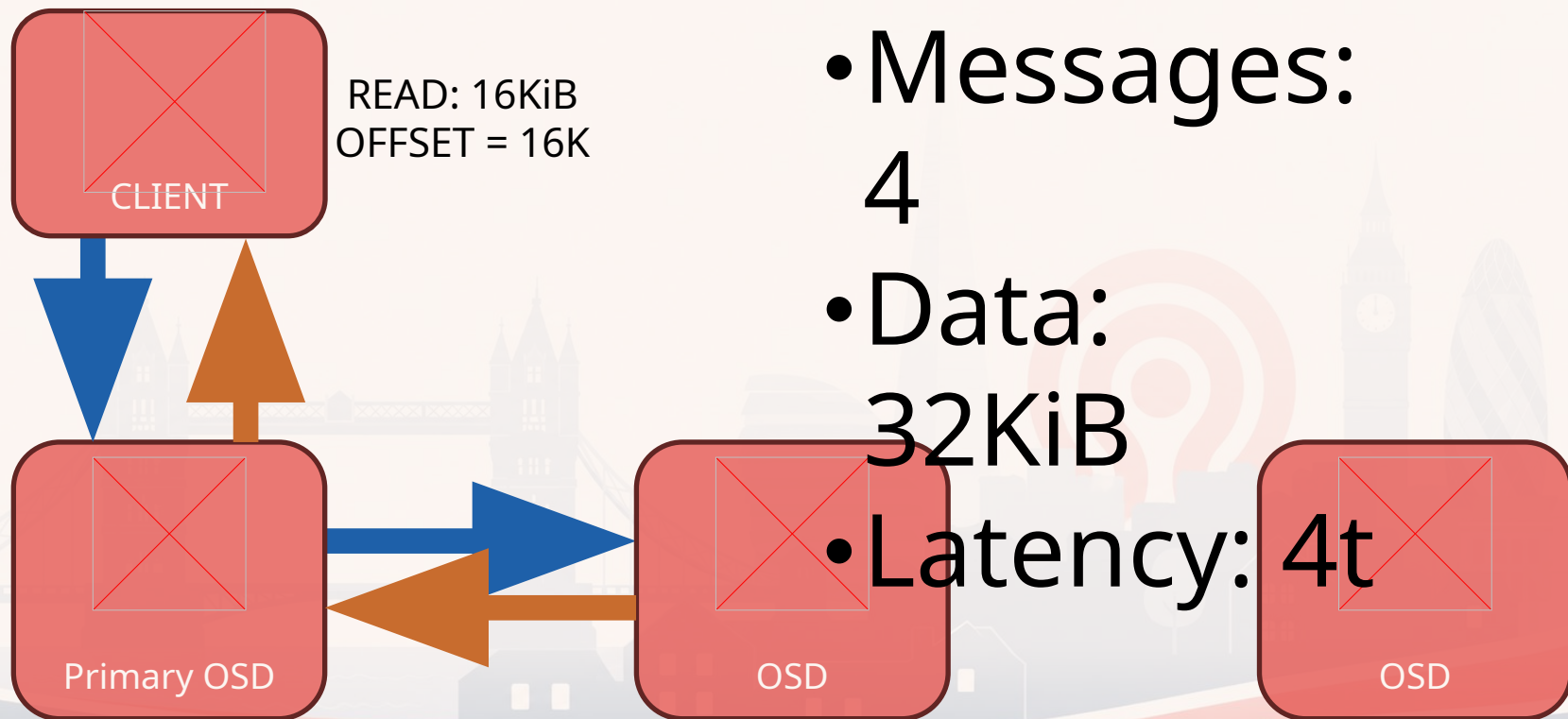


- Messages: 1
- Data: 0KiB
- Latency: 1t

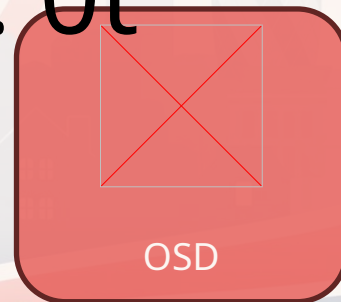
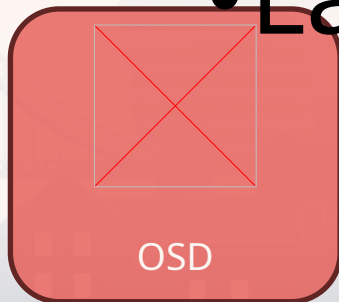
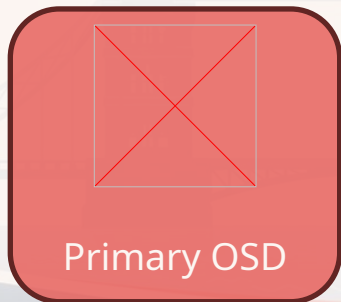
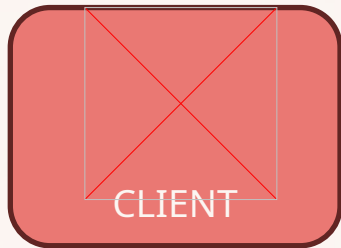


- Messages: 2
- Data: 0KiB
- Latency: 2t



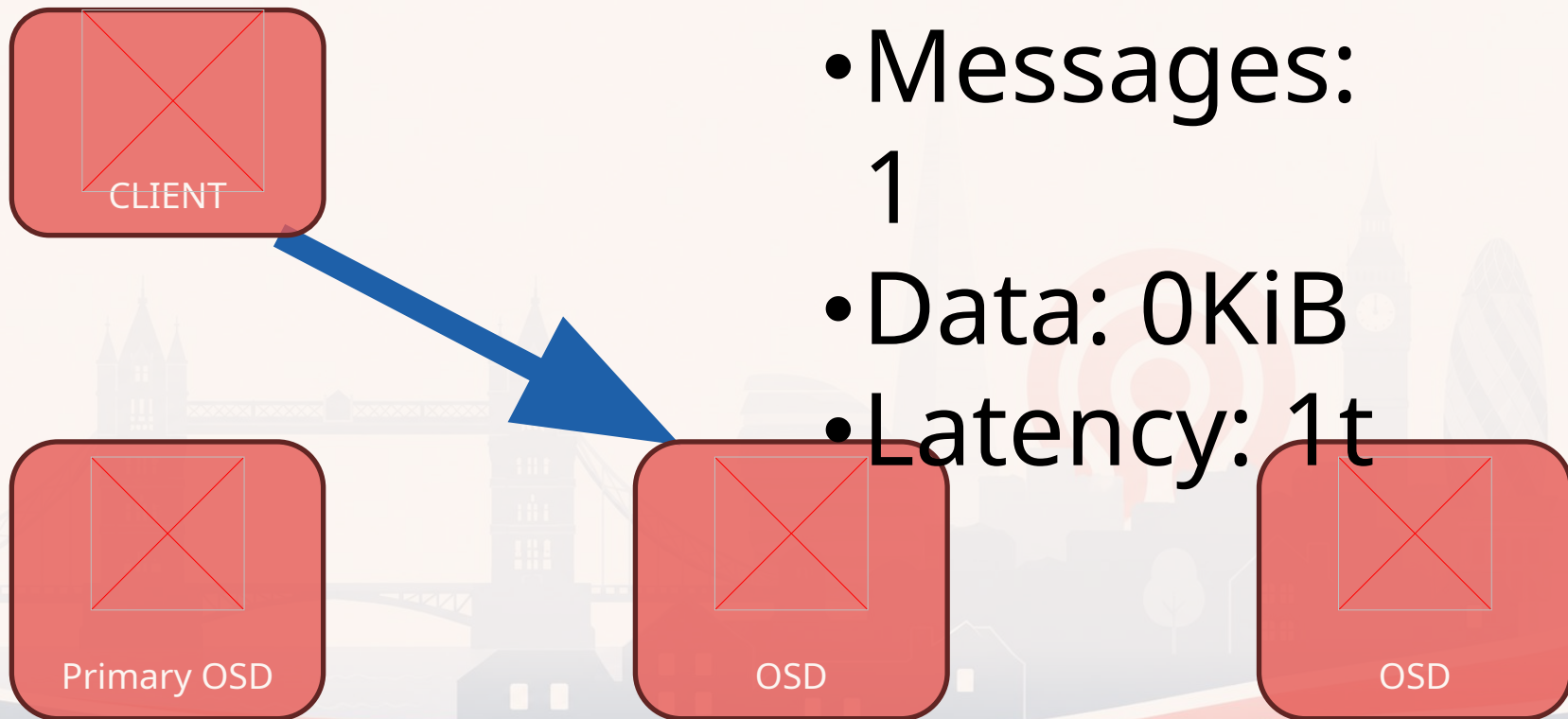


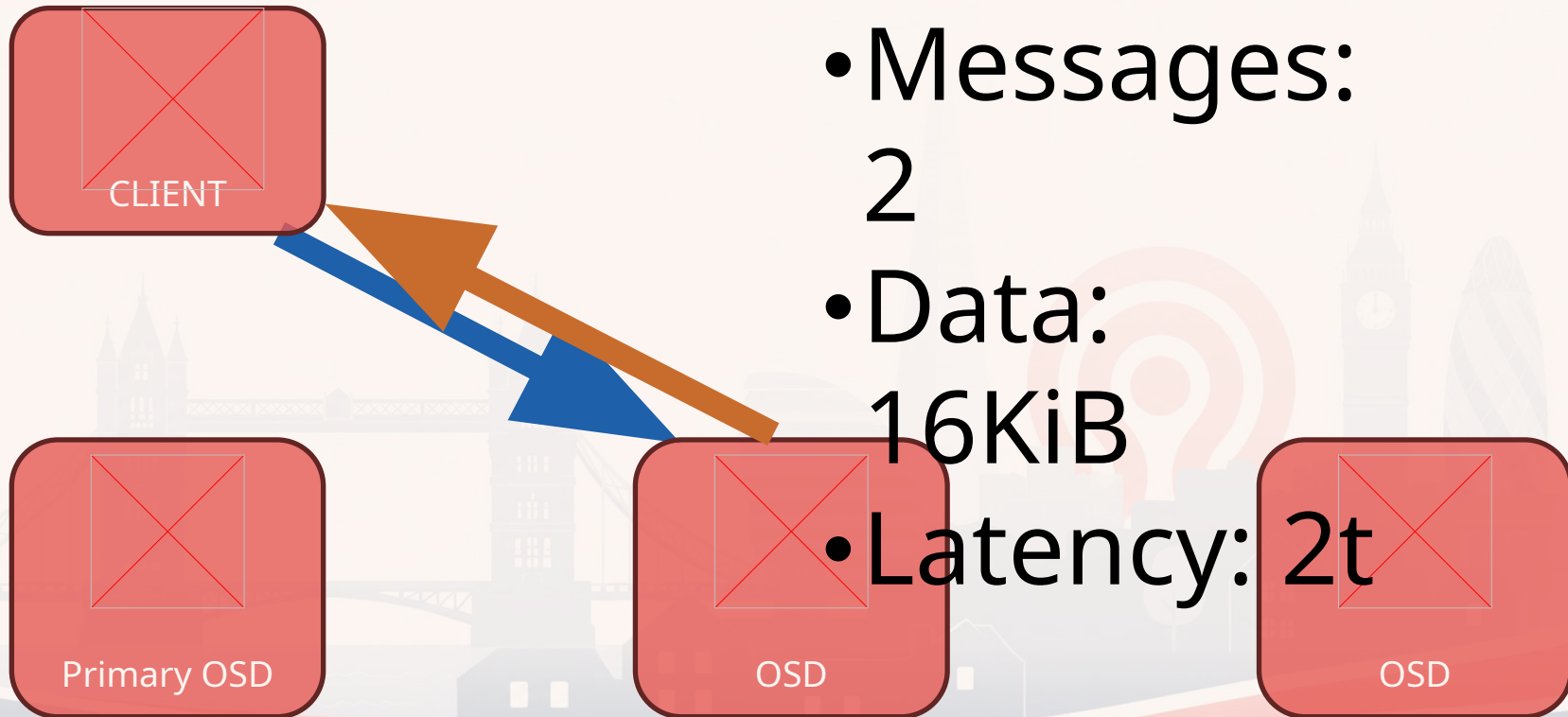
- Messages: 4
- Data: 32KiB
- Latency: 4t



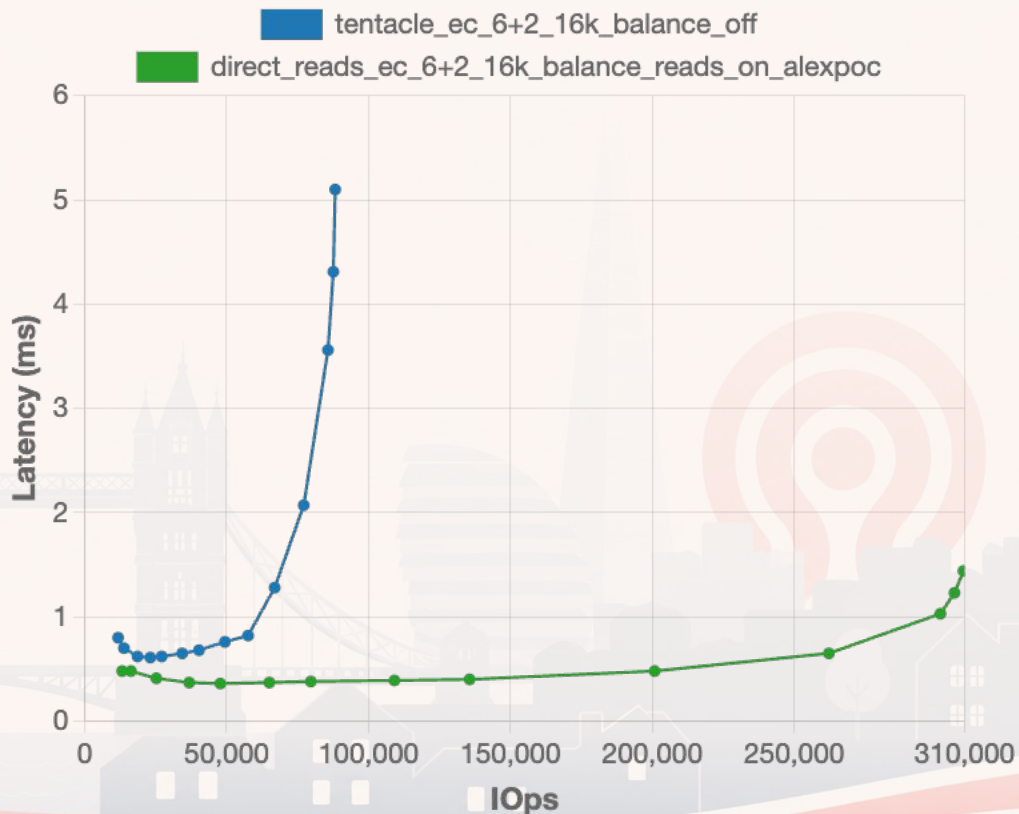
- Messages: 0
- Data: 0KiB
- Latency: 0t

Direct Read Path

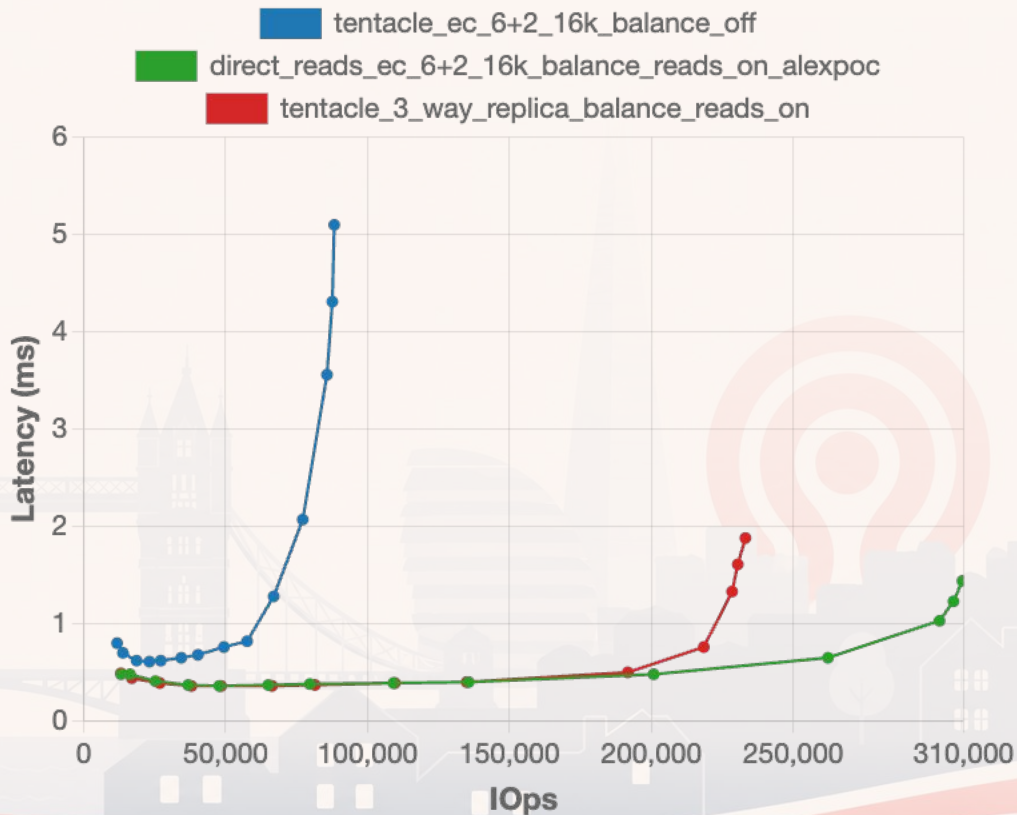




Performance – 8k reads

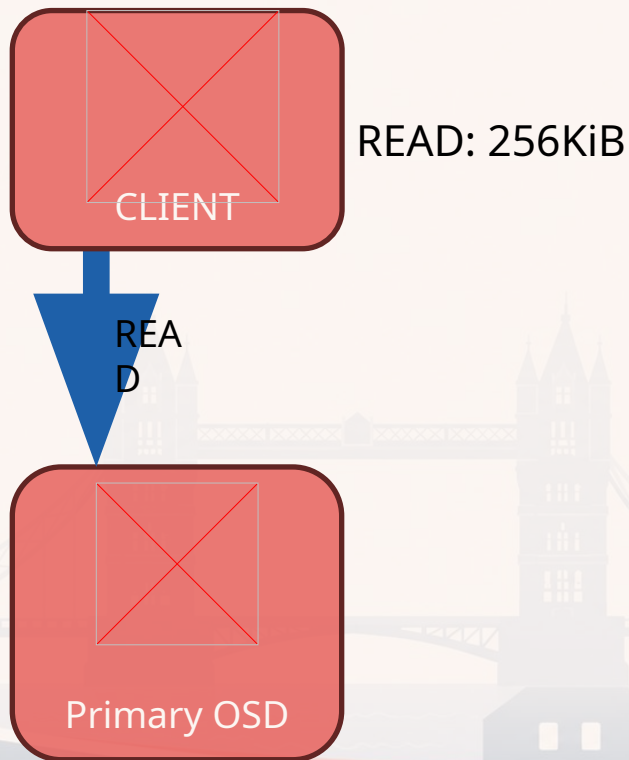


Performance – 8k reads

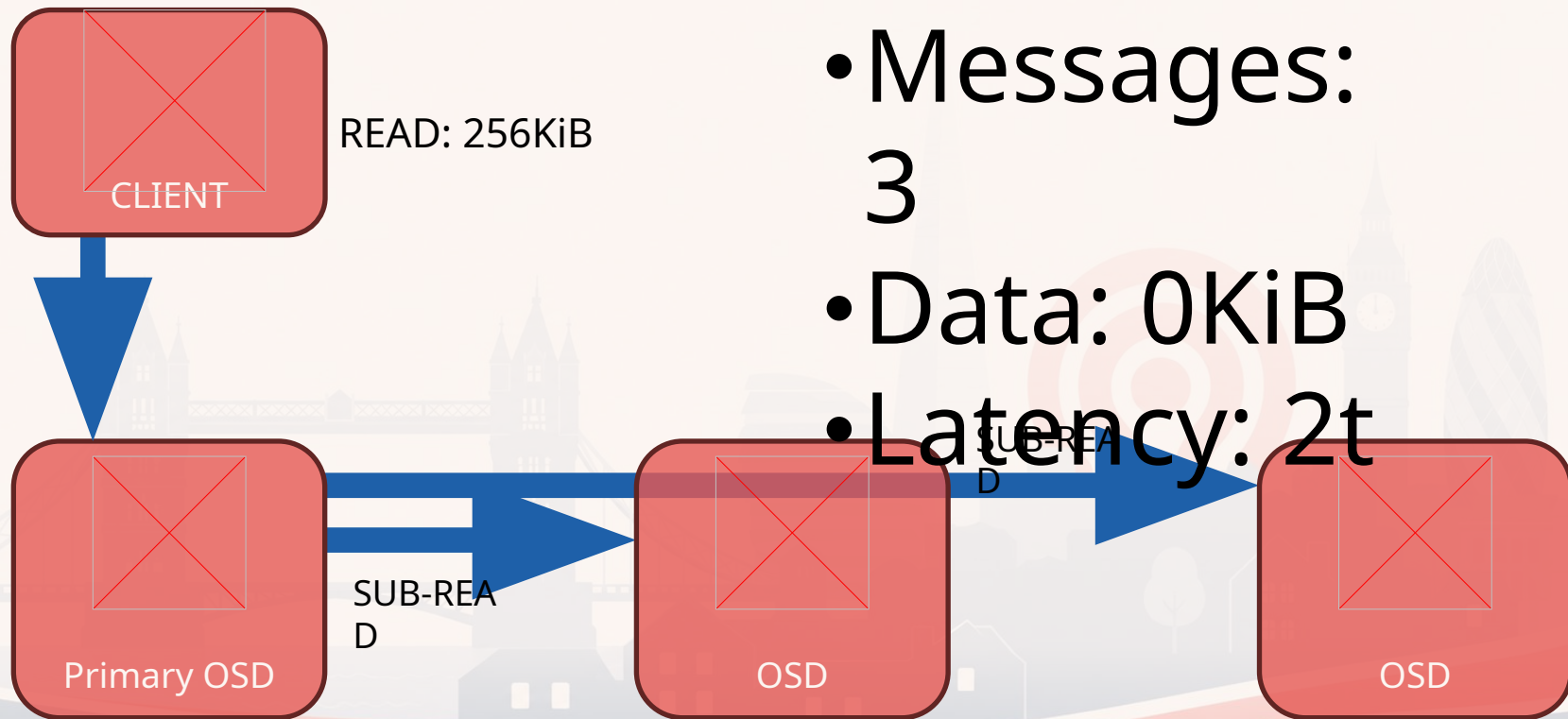




Large Reads

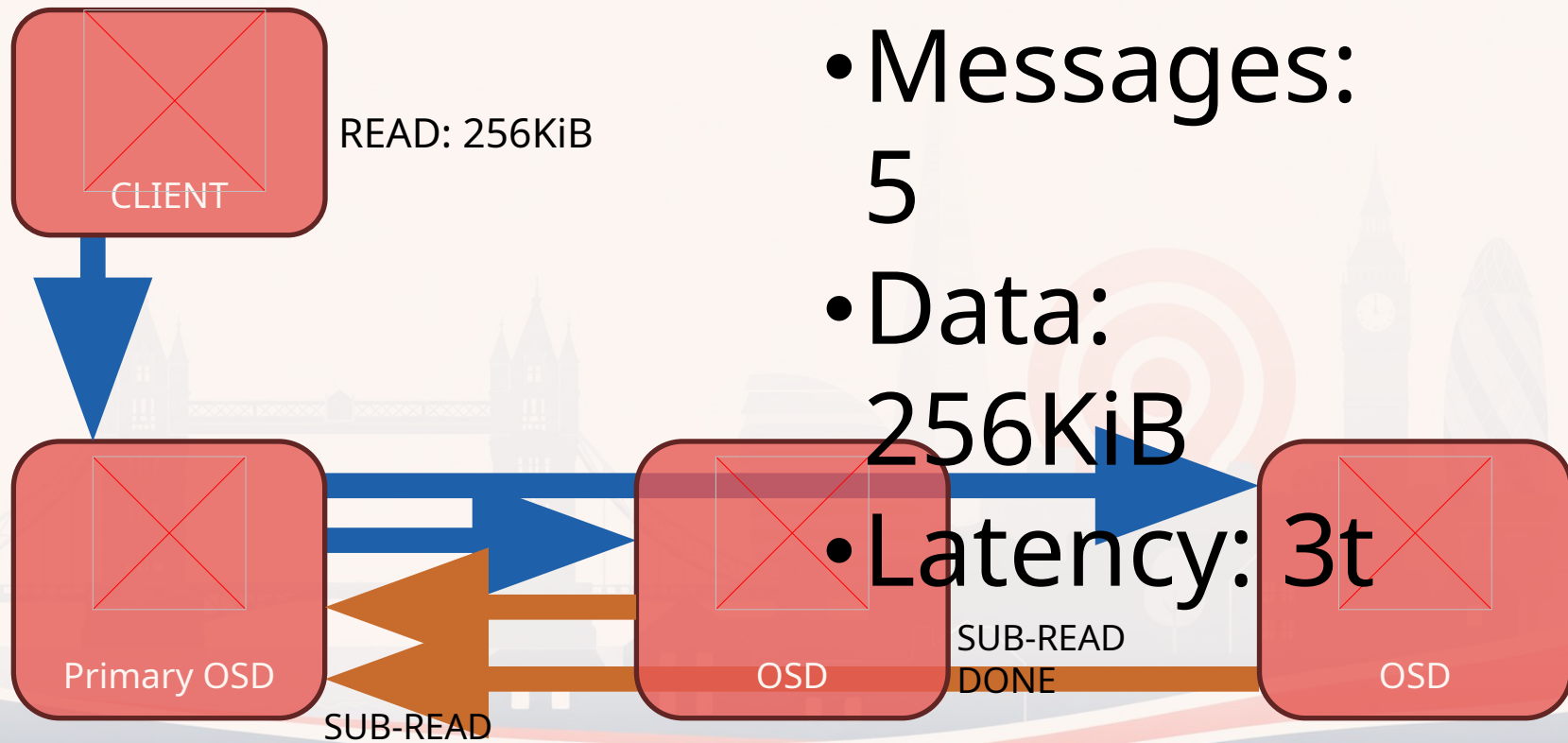


- Messages: 1
- Data: 0KiB
- Latency: 1t

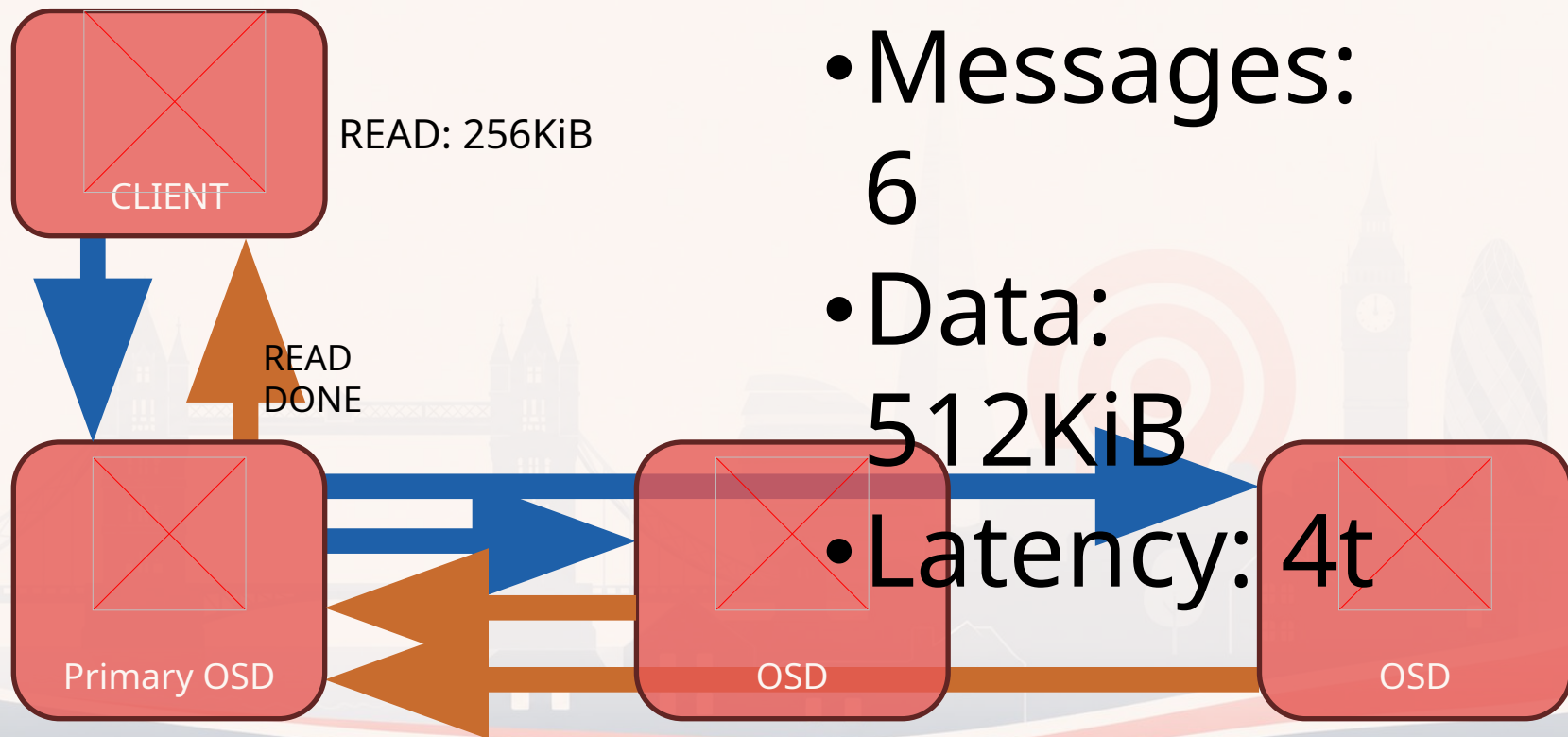


- Messages: 3
- Data: 0KiB
- Latency: 2t

As-is Read Path

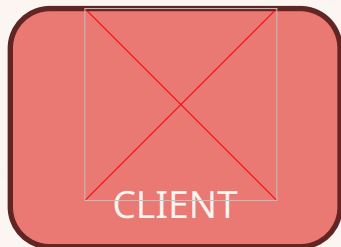


As-is Read Path

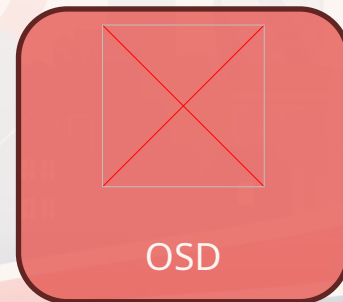
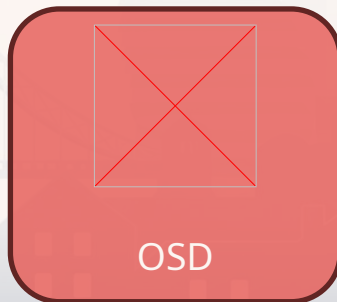
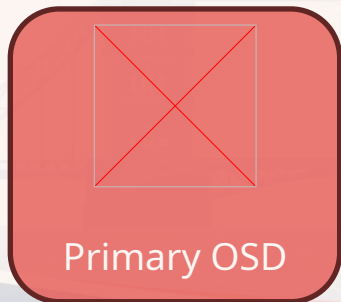


- Messages: 6
- Data: 512KiB
- Latency: 4t

Direct Read Path



- Message
- s:
- Data:
- Latency:



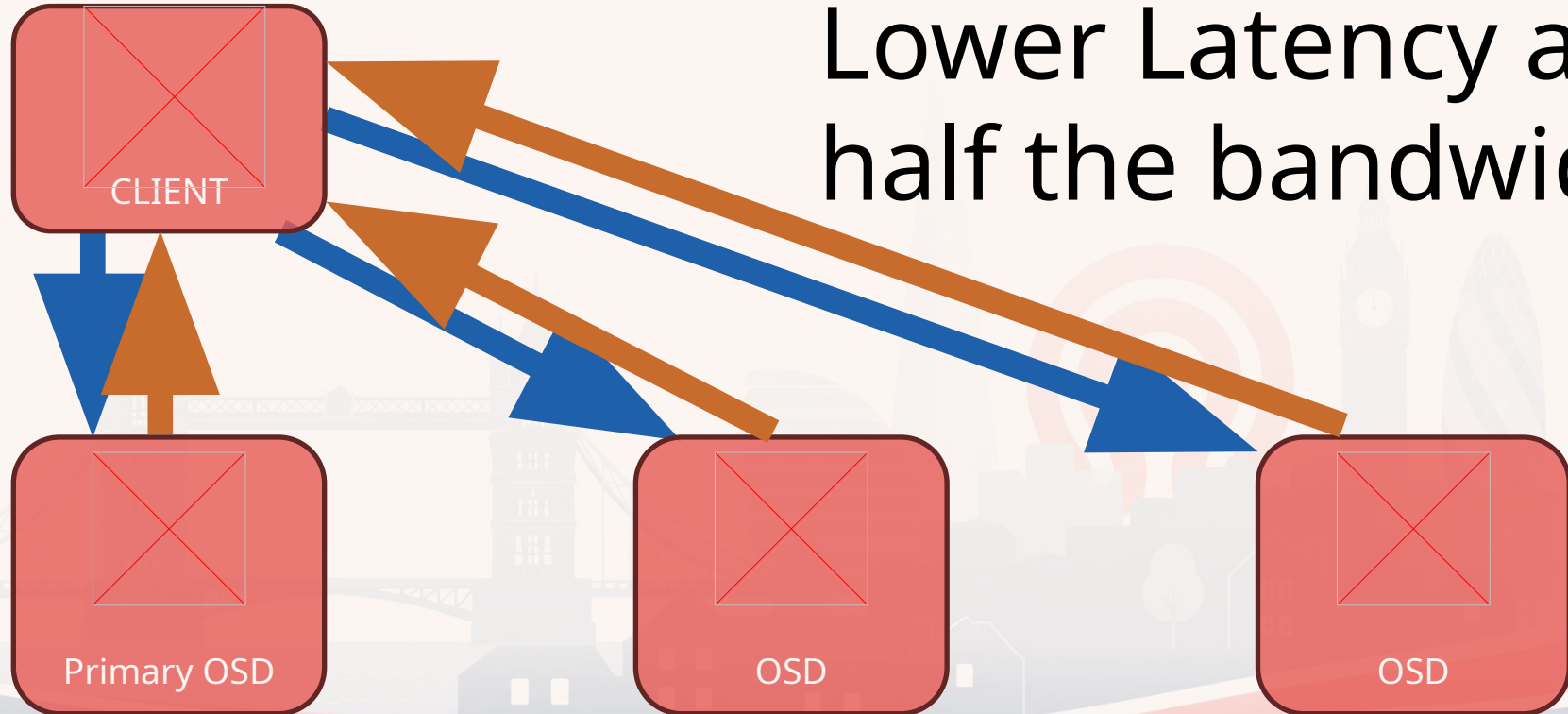
Direct Read Path



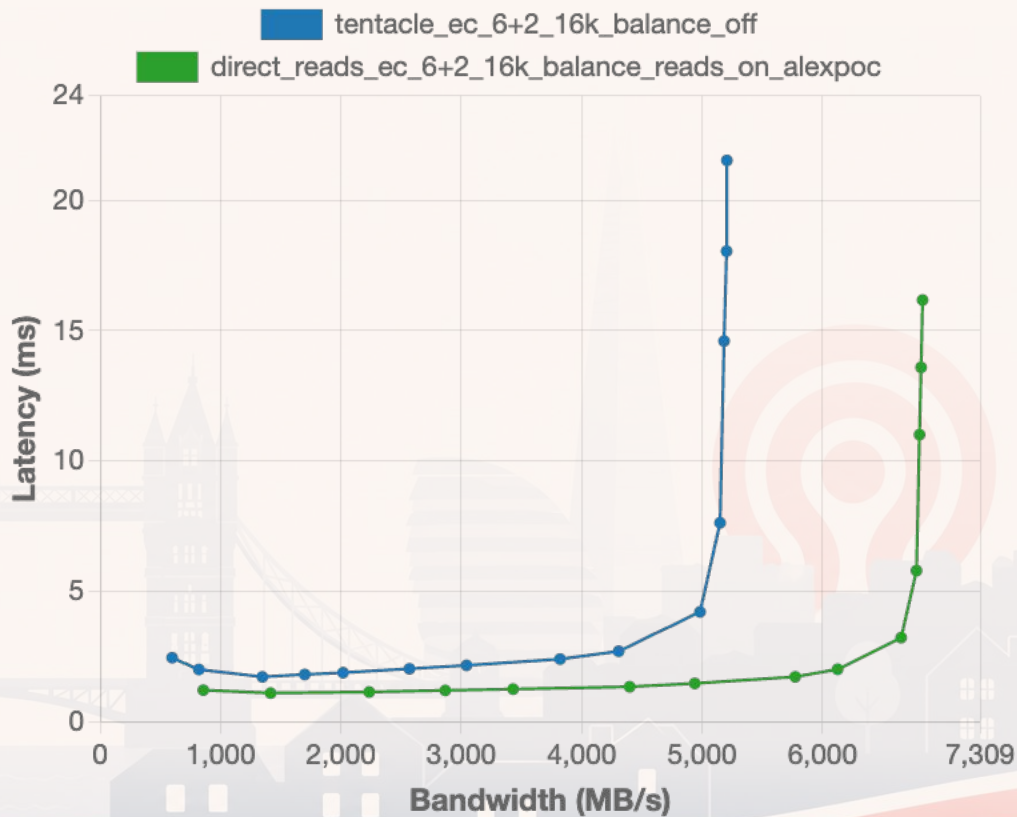
- Messages: 3
- Data: 0KiB
- Latency: 1t

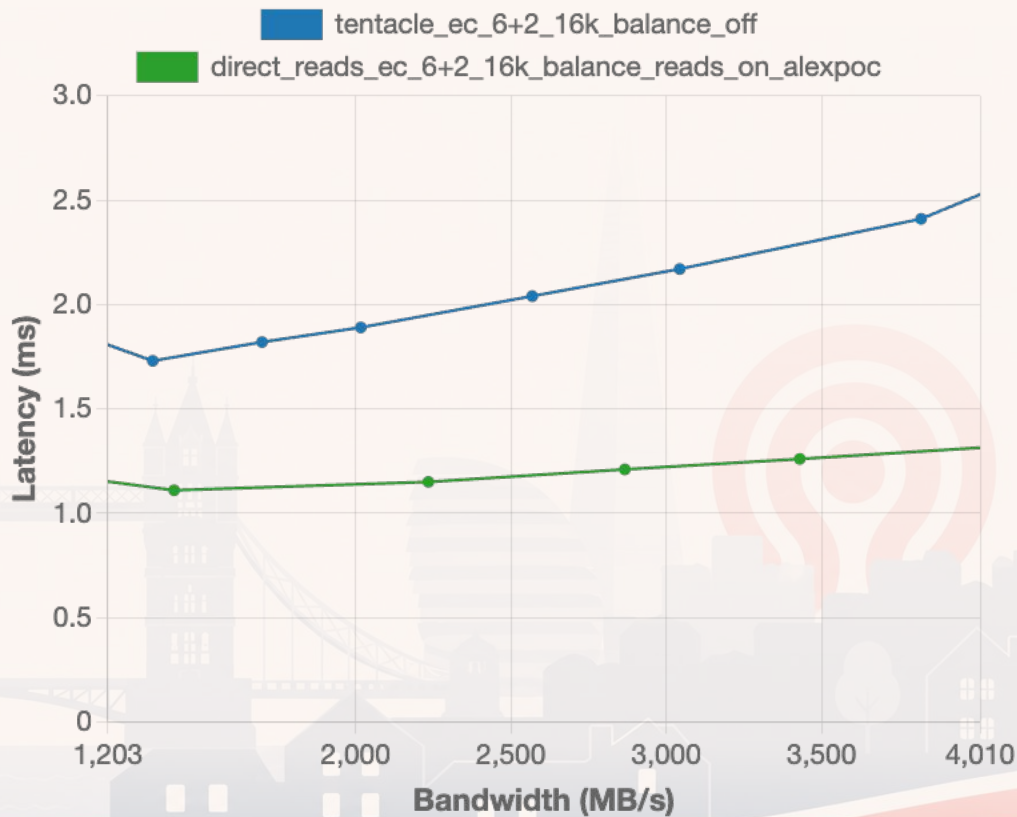
Direct Read Path



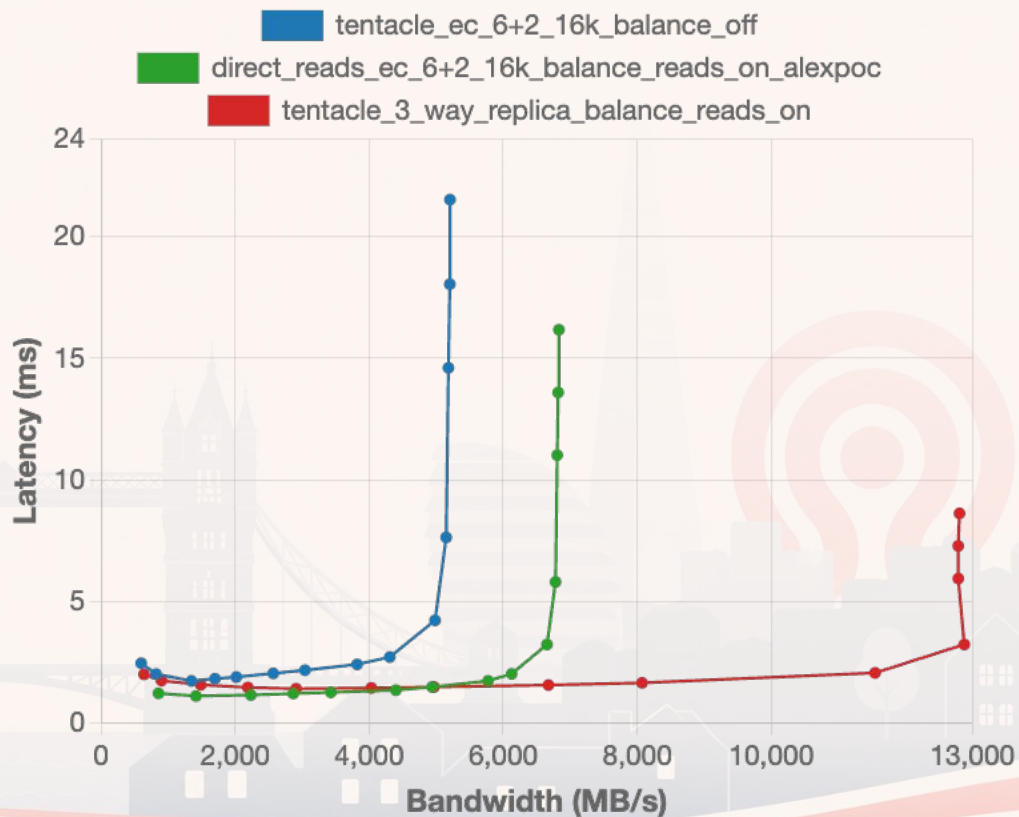


Lower Latency at
half the bandwidth

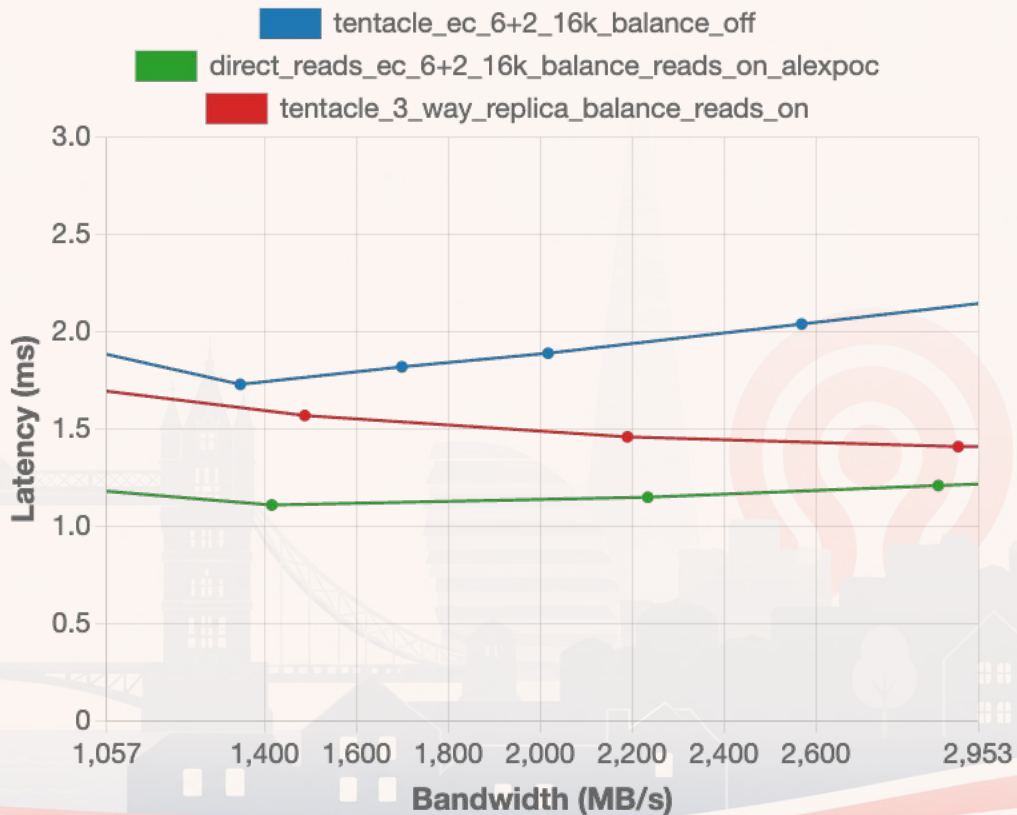




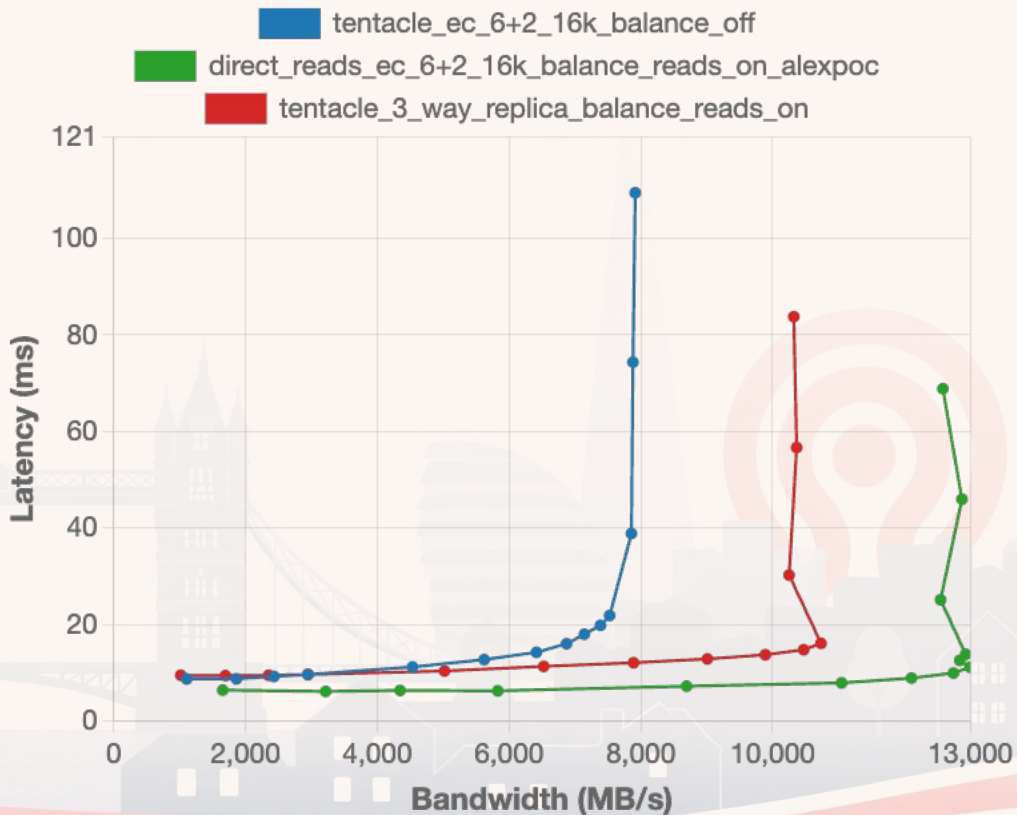
Performance vs Replica



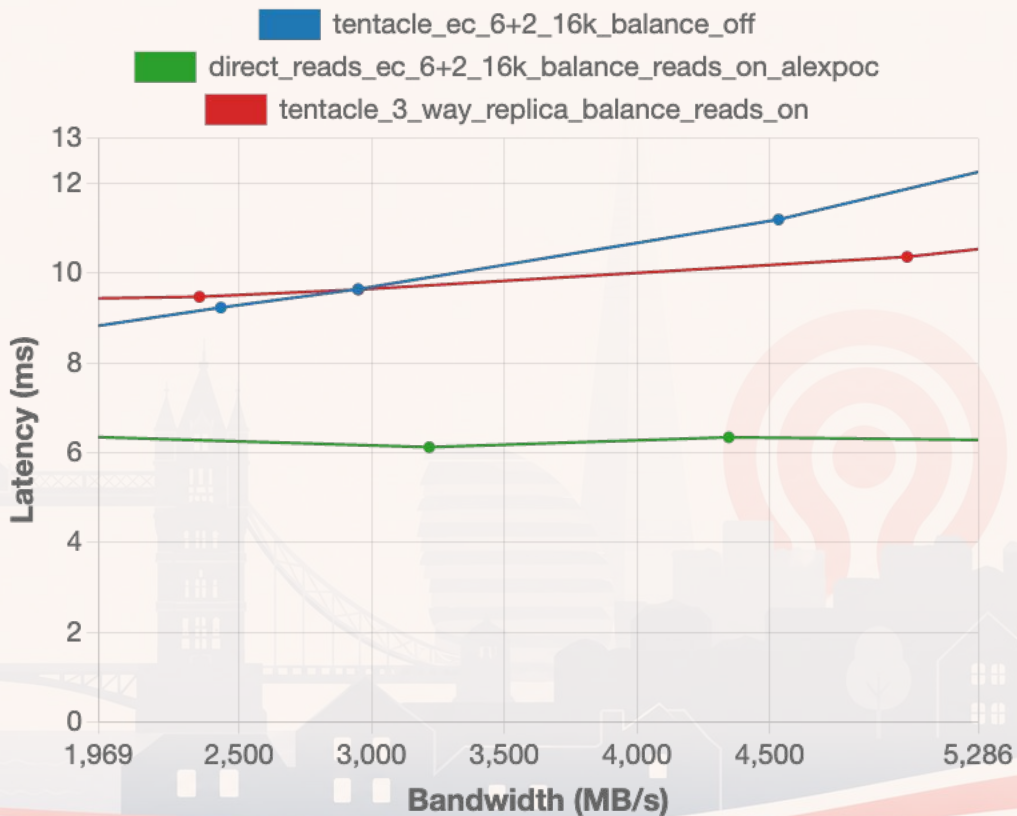
Performance vs Replica



Performance with 4MiB reads!



Performance with 4MiB reads!

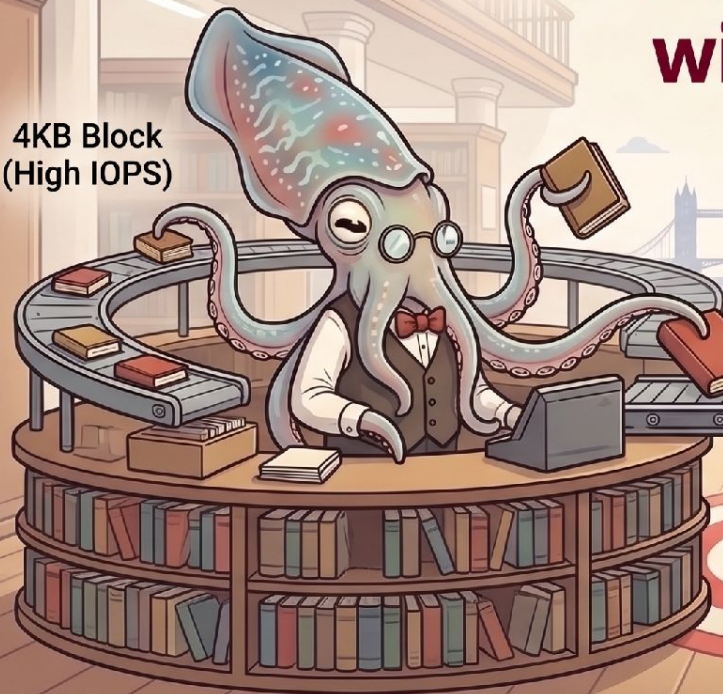




ceph days
LONDON 2026

Performance Variation with IO Size

**4KB Block
(High IOPS)**



**128KB Block
(Balanced)**

**1MB Block
(Low IOPS)**

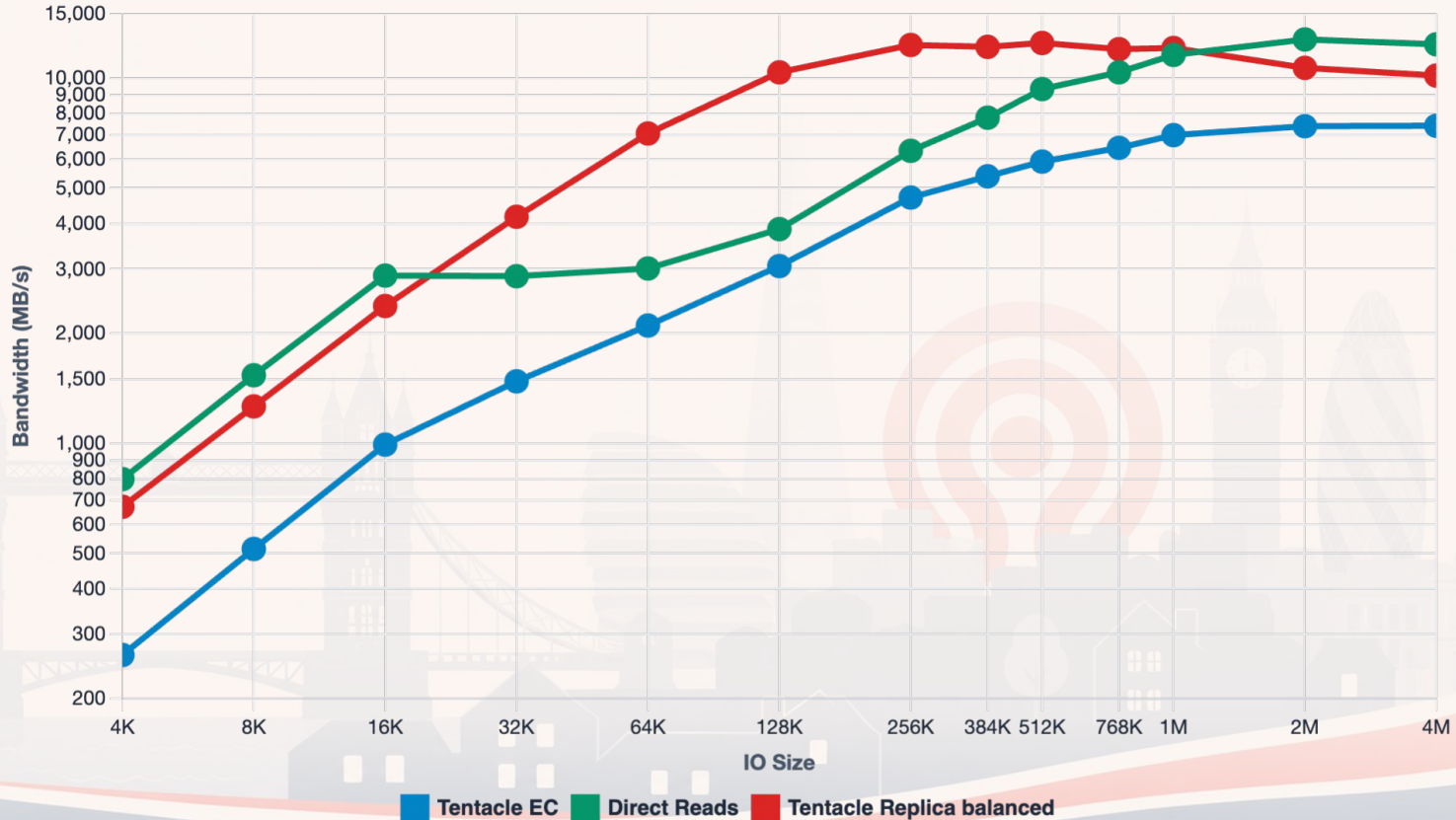


**Data Archive:
Performance
Studies**

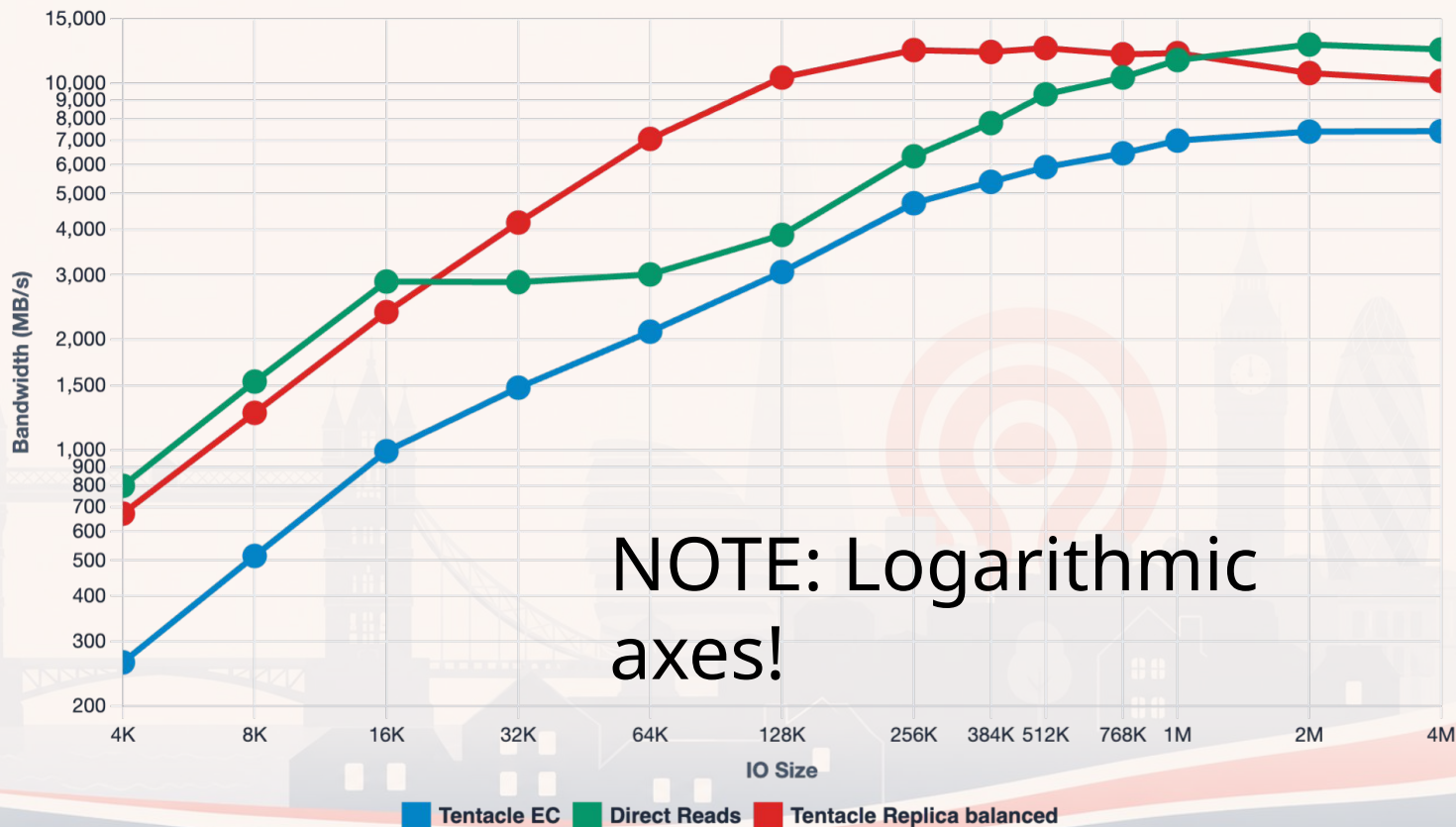
**IO Size
Management**



Bandwidth variations vs IO size



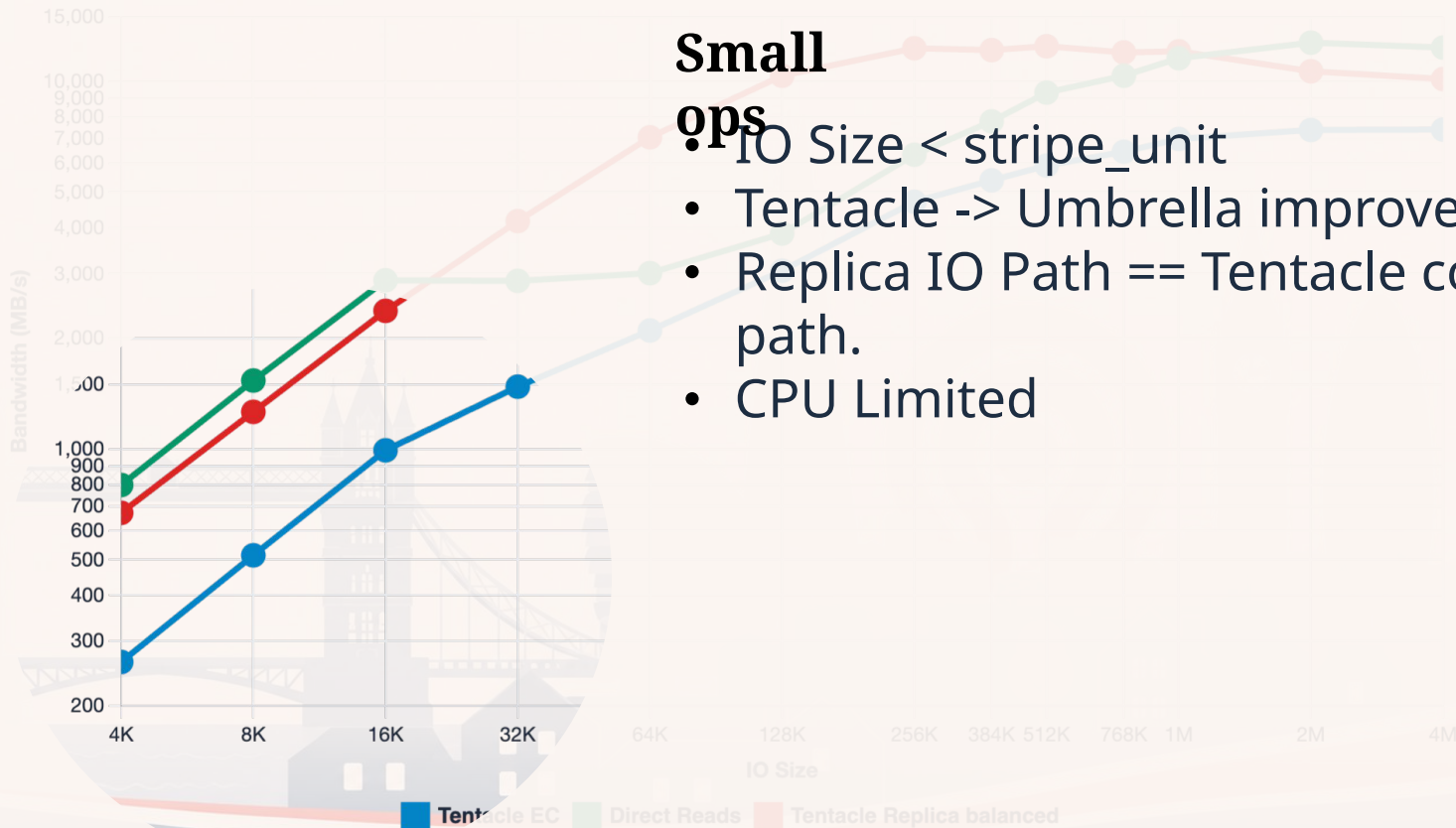
Bandwidth variations vs IO size



NOTE: Logarithmic axes!



Bandwidth variations vs IO size



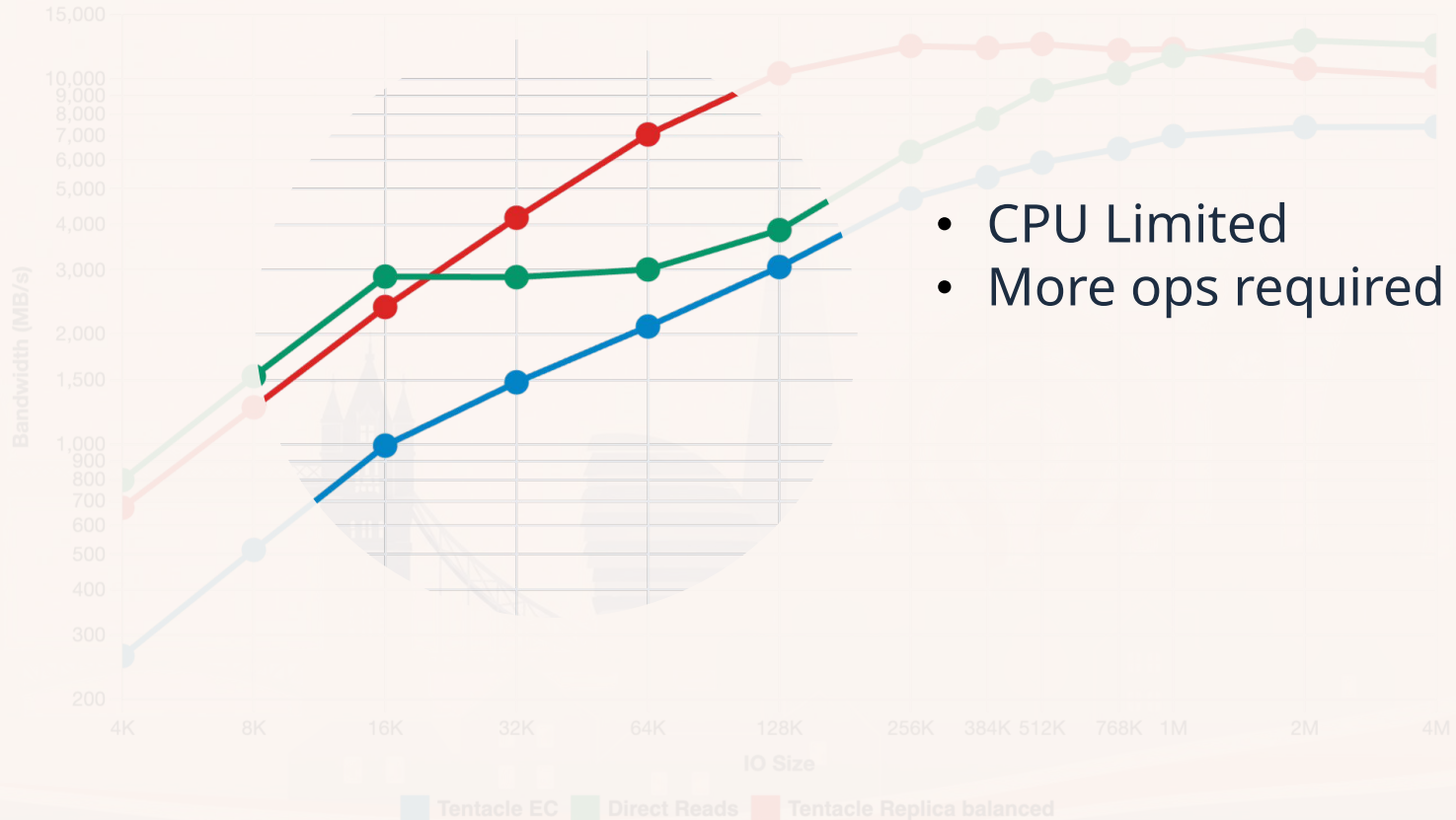
Small

ops

- IO Size < stripe_unit
- Tentacle -> Umbrella improved ?
- Replica IO Path == Tentacle code path.
- CPU Limited

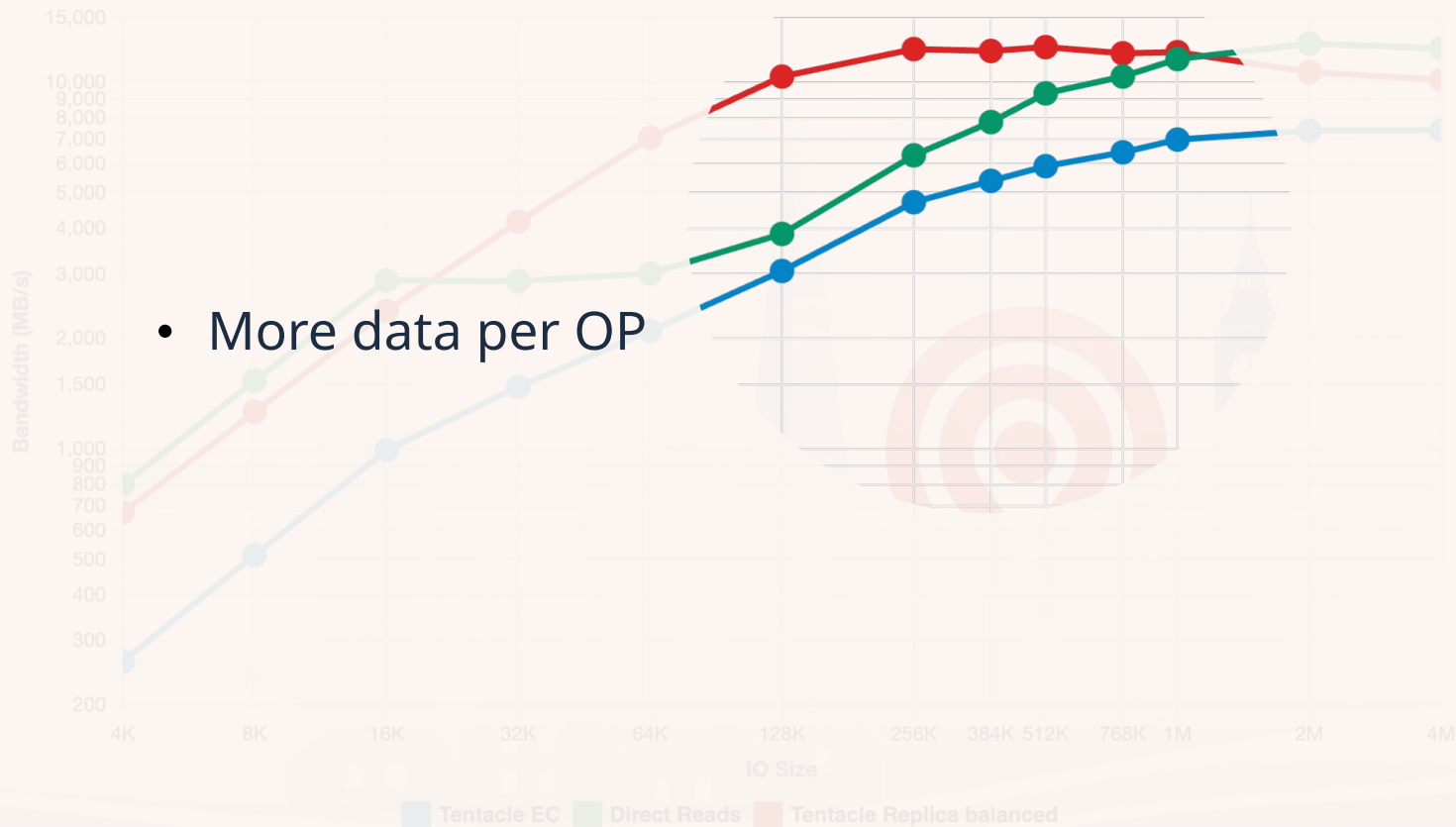


Bandwidth variations vs IO size



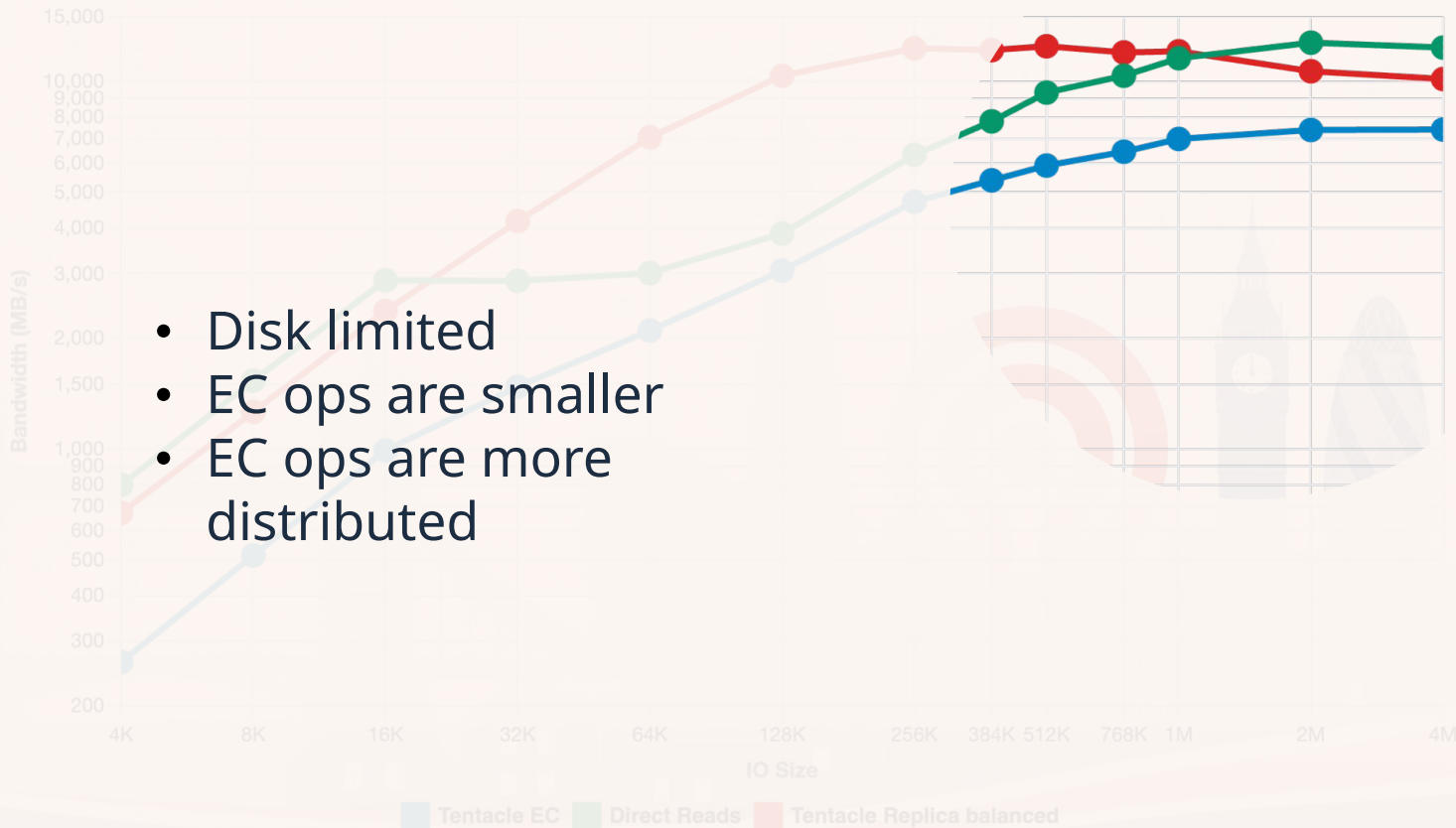


Bandwidth variations vs IO size



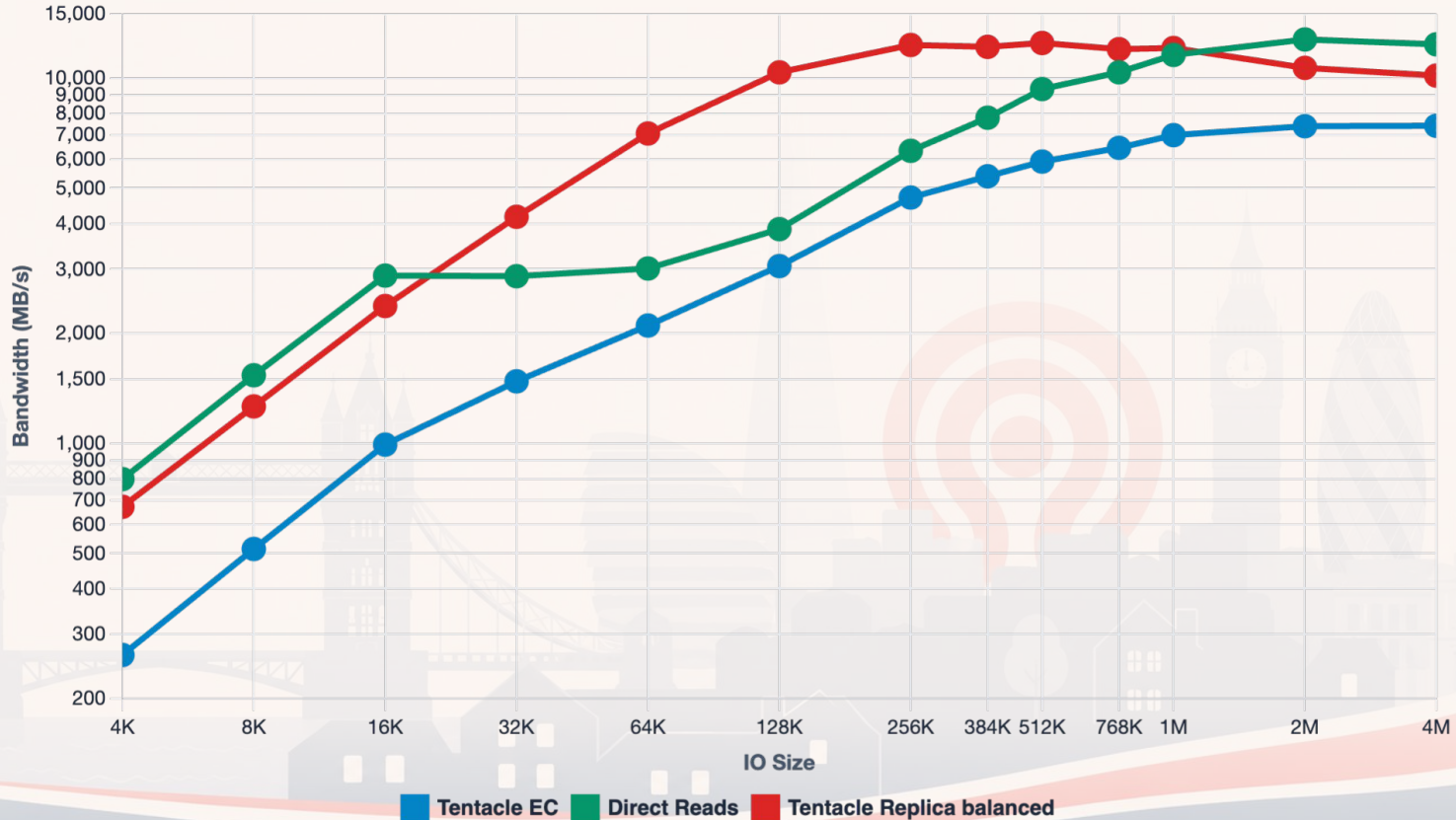


Bandwidth variations vs IO size

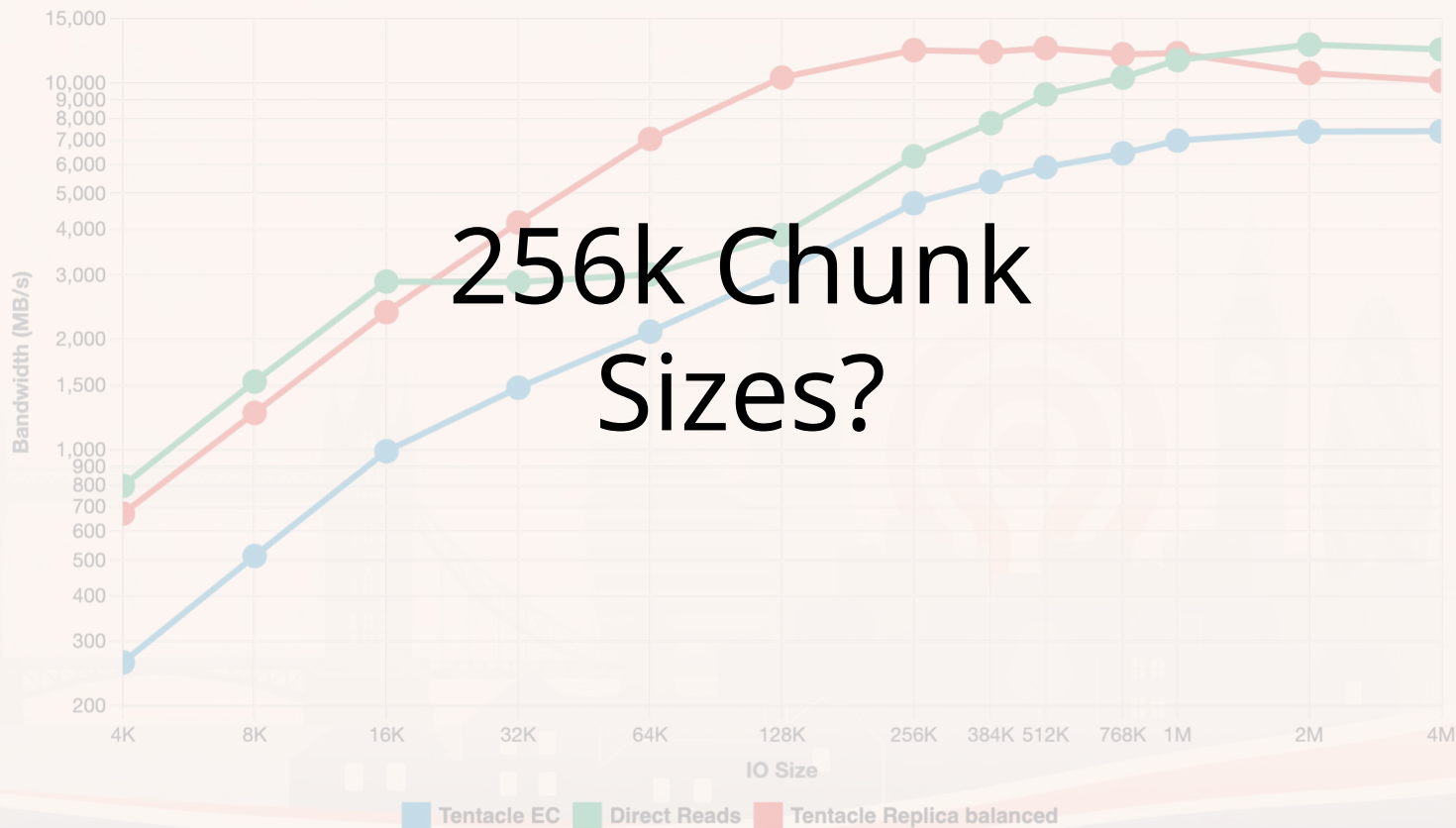




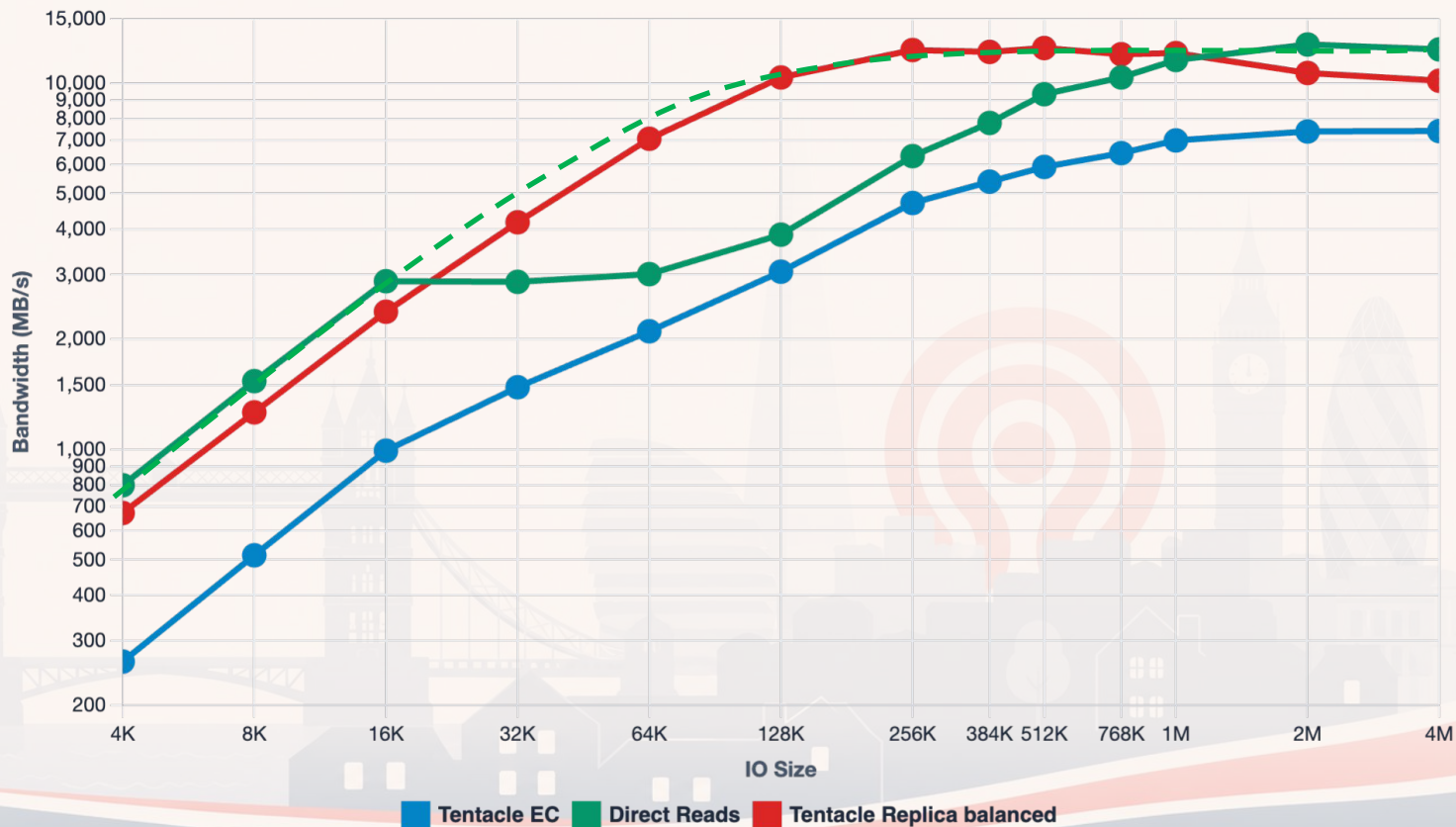
Bandwidth variations vs IO size

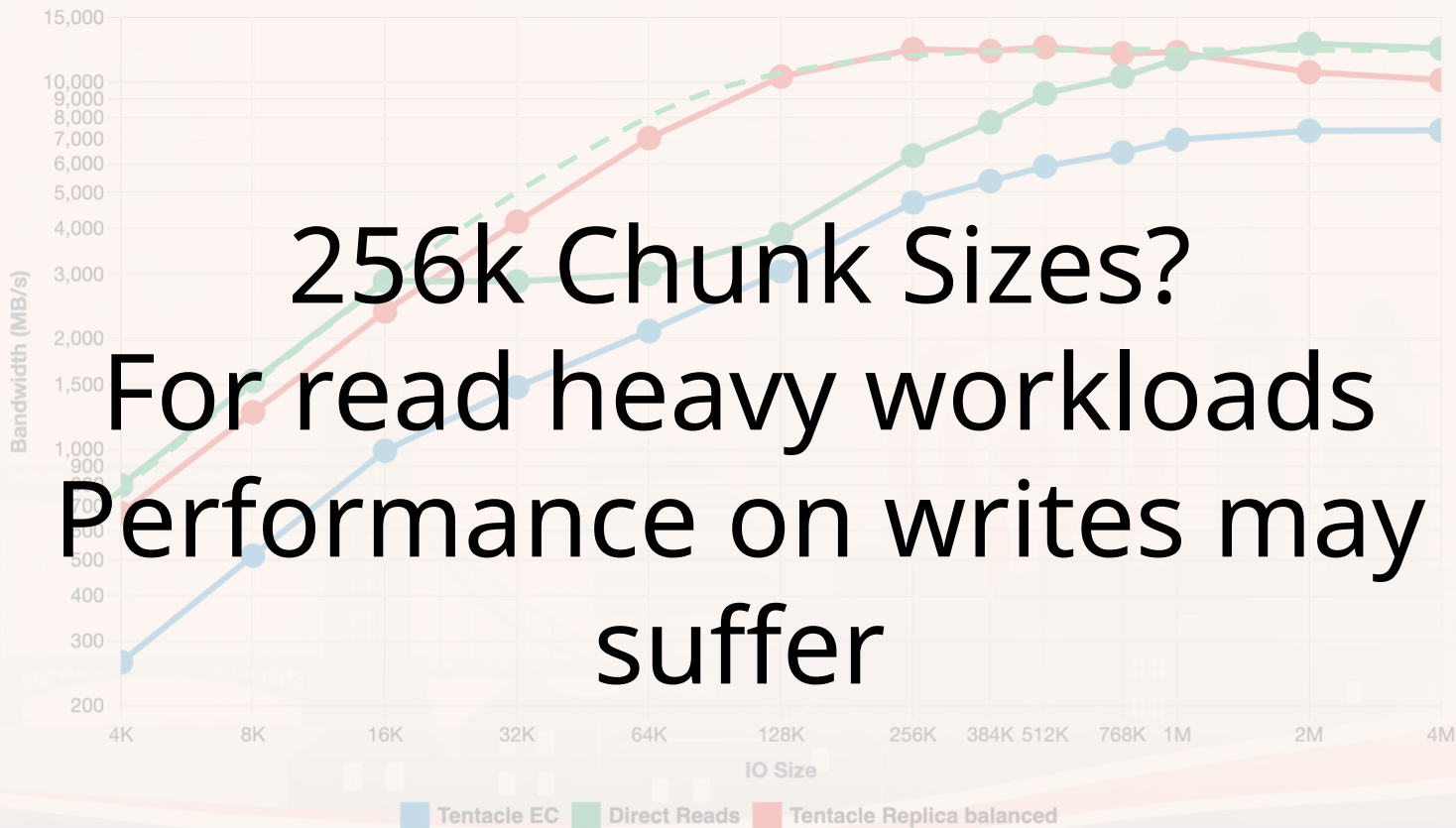


Bandwidth variations vs IO size



Bandwidth variations vs IO size







Failed OSDs

- Client Detects and re-routes ops to primary
- Redrives to primary may occur during transition.

- Uncommitted writes can cause reads to the same object to be re-driven
- Reads can over-take outstanding writes: All Ceph clients expect and handle this cleanly
- Reads will never be torn (i.e. contain part of a write)



ceph days
LONDON 2026



Can we do this for REPLICA too?

(for large reads)



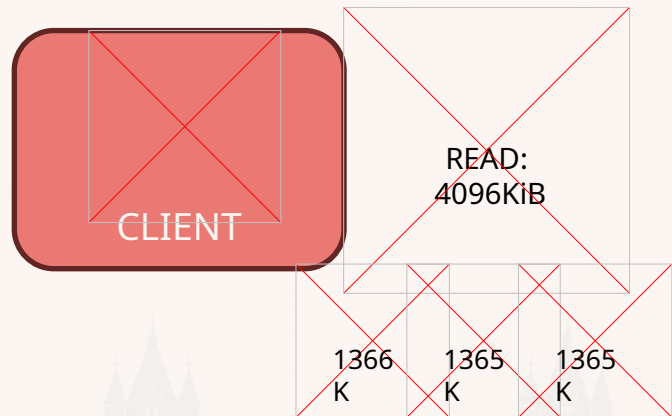
ceph days
LONDON 2026



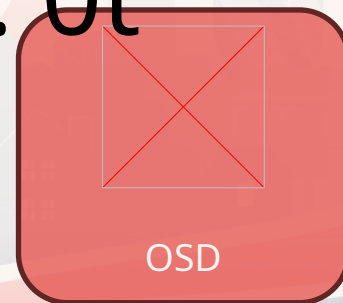
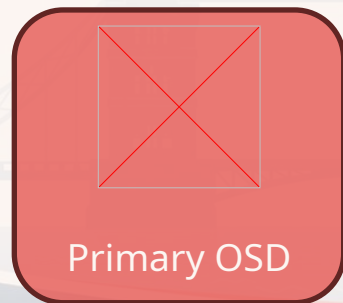
Yes!

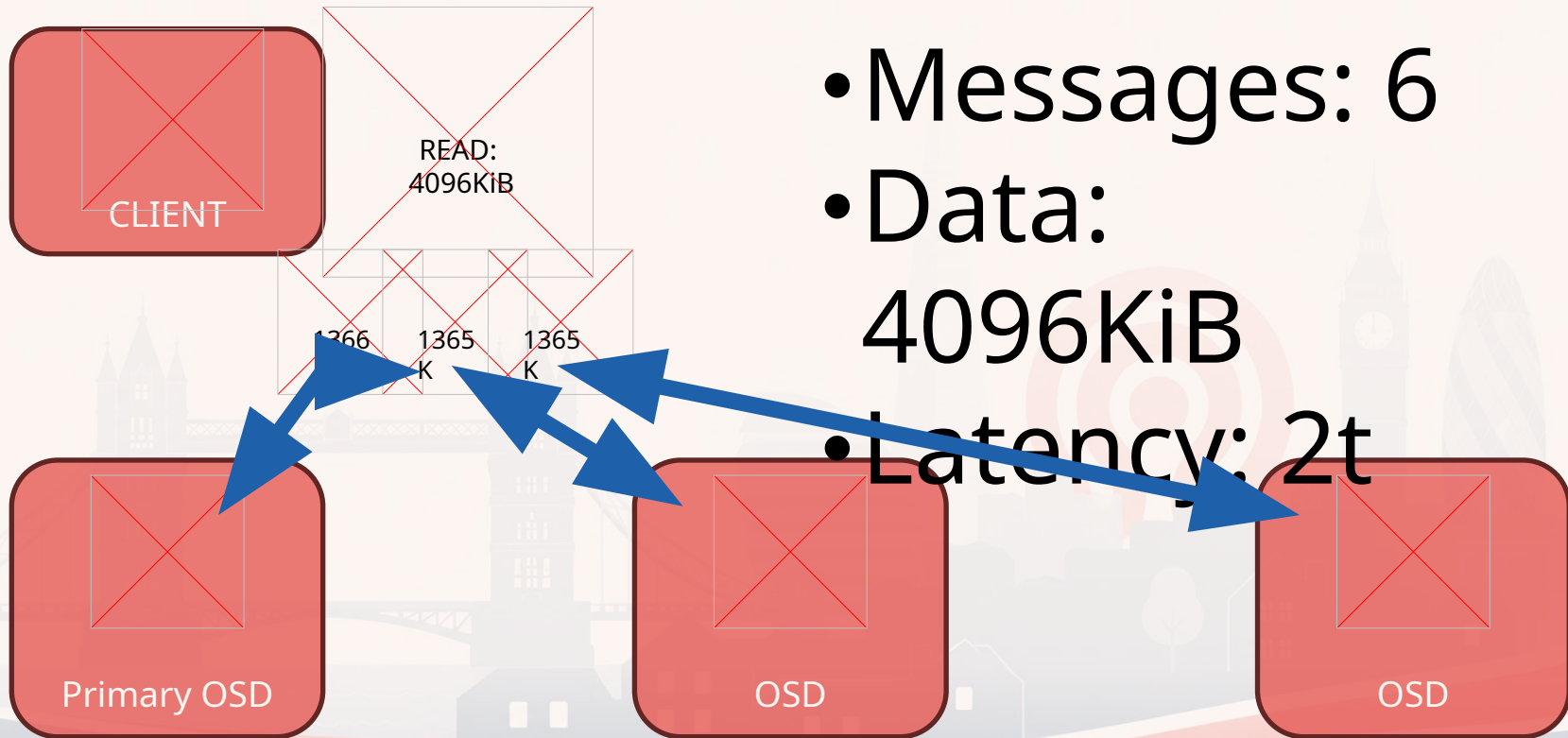
(for large reads)





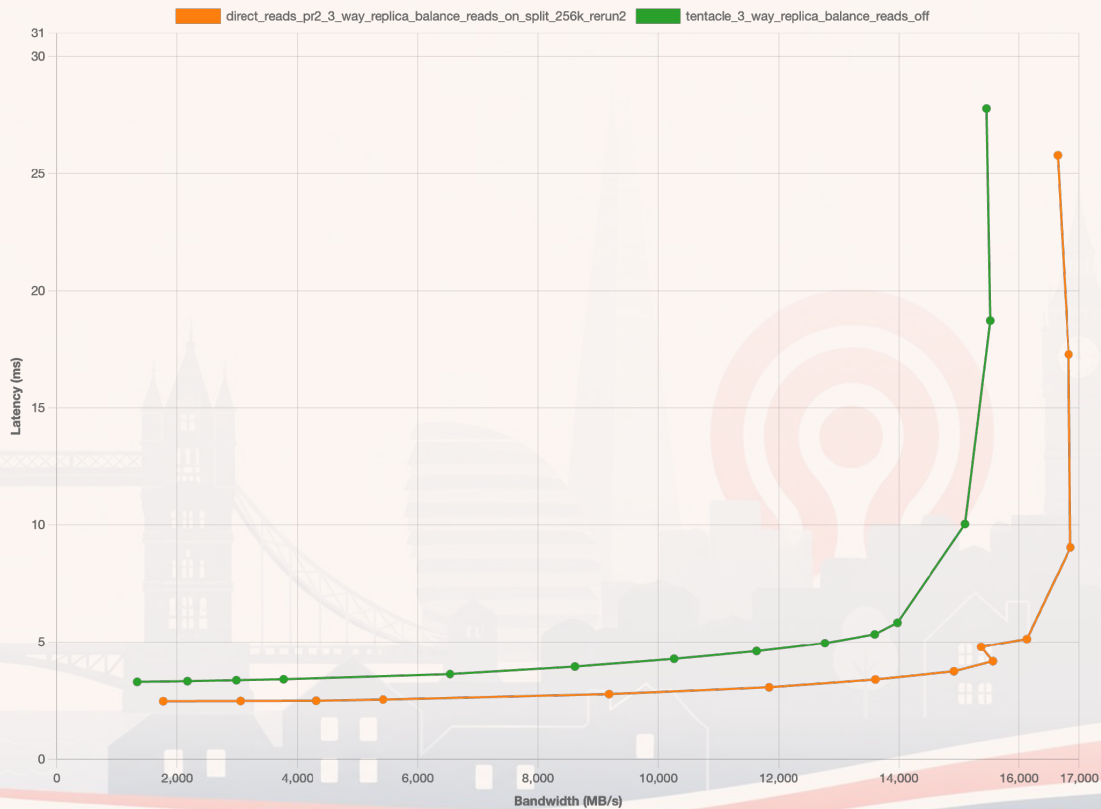
- Messages: 0
- Data: 0KiB
- Latency: 0t



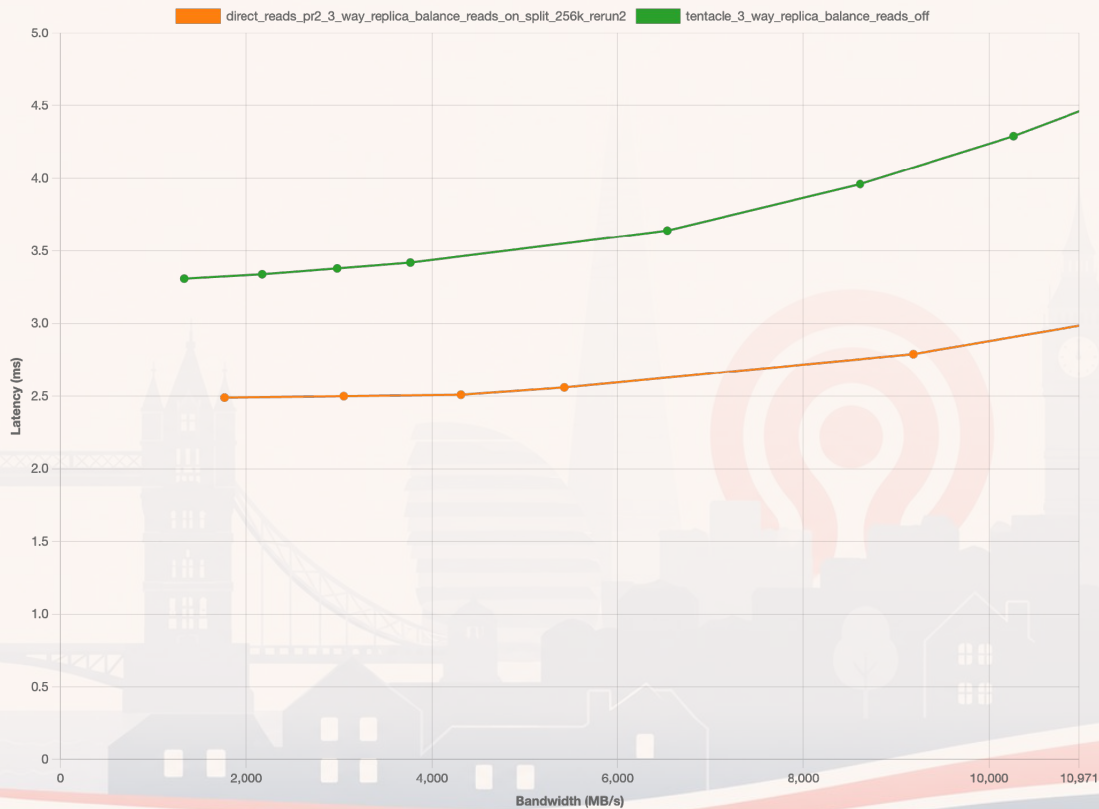


- Messages: 6
- Data: 4096KiB
- Latency: 2t

Split Replica Performance (2MiB)



Split Replica Performance (2MiB)



Turning Direct Reads ON!



This is a config-option:

```
[client]
rbd_read_from_replica_policy = balance
```

<https://docs.ceph.com/en/latest/rbd/rbd-config-ref>

This is a mount option:

```
read_from_replica=balance
```

<https://docs.ceph.com/en/latest/man/8/mount.ceph/>

Balanced reads can be turned on globally using the config option:

```
rados_replica_read_policy = balance
```

NOTE: May not work well with unusual/unsupported features.

Enable balanced reads

AND

```
[global]  
osd_min_split_replica_read_size = 262144
```

Any Questions?

