

Survey: Istio users' cross-cluster routing today

Gangmuk Lim (UIUC), Aditya Prerepa (UIUC), Brighten Godfrey (UIUC & Broadcom), Radhika Mittal (UIUC)



To learn more about SLATE project, Visit our project website!

<https://servicelayernetworking.github.io/slate>

Intro

Who are we?

- We are a research group in the Department of Computer Science at the University of Illinois at Urbana-Champaign, working on improved platforms for multi-cluster microservice deployments.

What is this survey?

- The purpose of this survey is to understand common multi-cluster deployments in practice. The results of this survey will only be used for research purposes. It includes 18 questions regarding multi-cluster/cross-cluster routing.

Where was the survey distributed?

- The total number of responses is 31. Six of them were excluded since they do not run multi-cluster and have less than 10 nodes. The respondents of the survey were from a variety of internet businesses at varying scales, from 2 clusters and a few nodes to over 50 clusters and thousands of nodes.

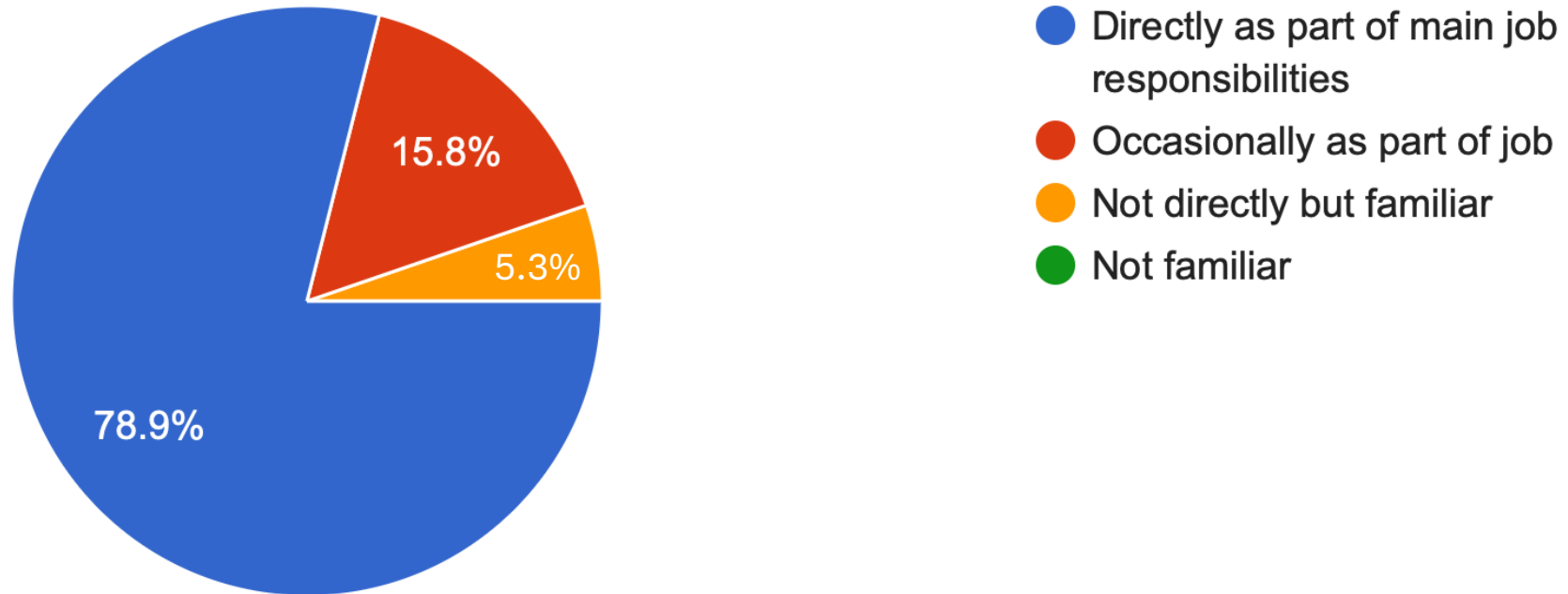
Potential bias of the result

- This survey was conducted in Istio community. We specified this survey was about multi-cluster & cross-cluster routing. Hence, the population could be biased to the people who are familiar with multi-cluster deployment and who actually have experience to deploying multiple clusters.

Q1

Do you manage Kubernetes cluster(s) as part of your job? (Throughout this survey, "clusters" refers to Kubernetes clusters.)

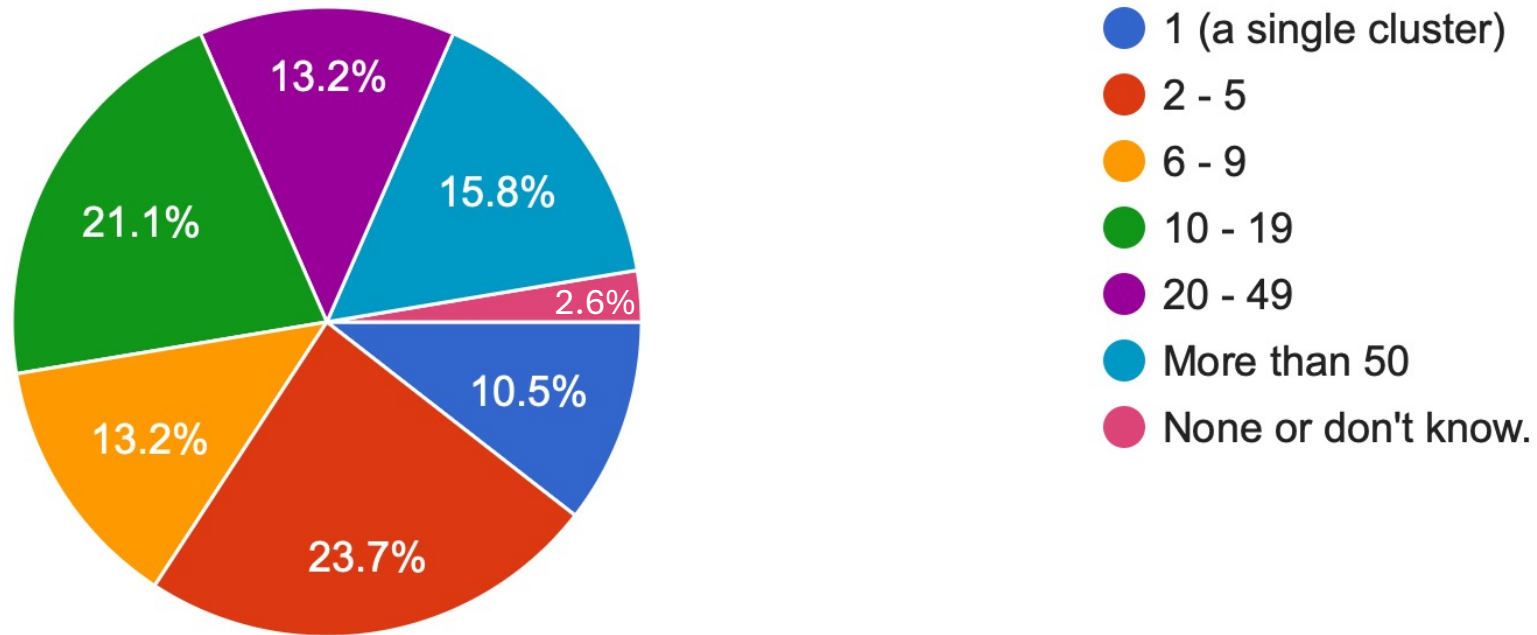
38 responses



Q2

Roughly how many production clusters do you have?

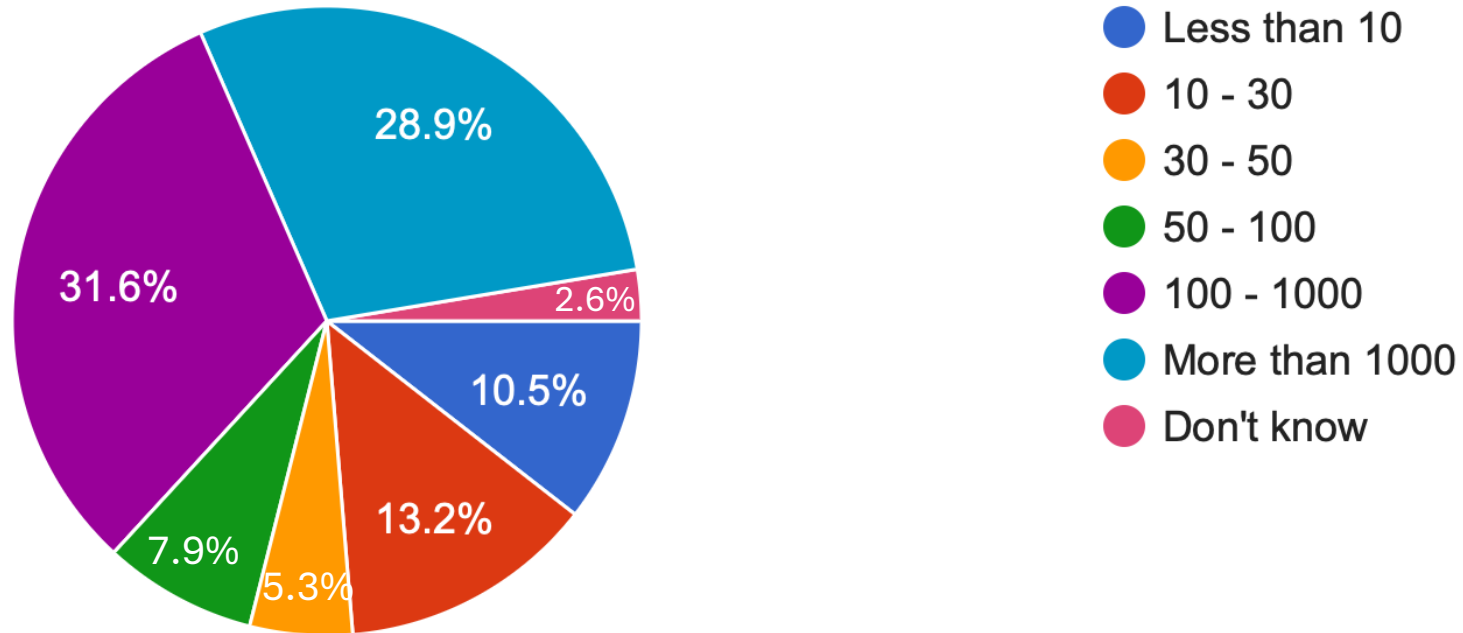
38 responses



Q3

Across all of your clusters together, roughly how many nodes (VMs or physical machines) are there?

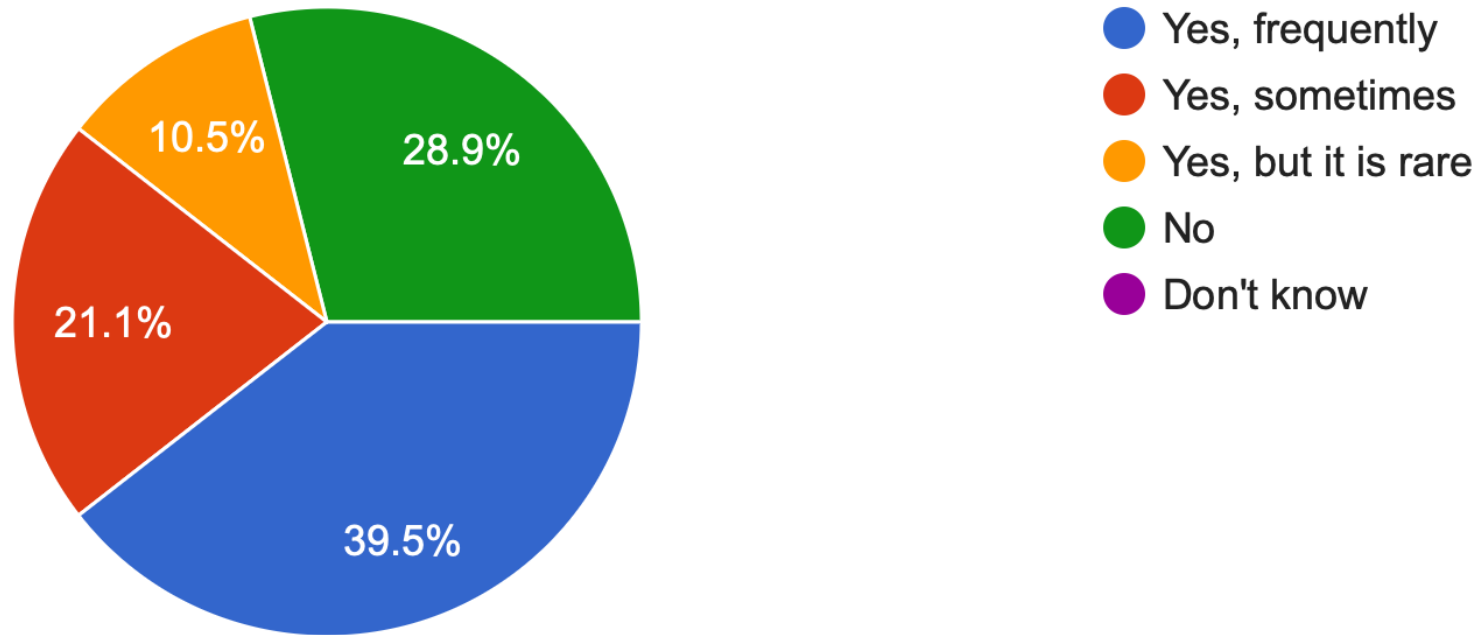
38 responses



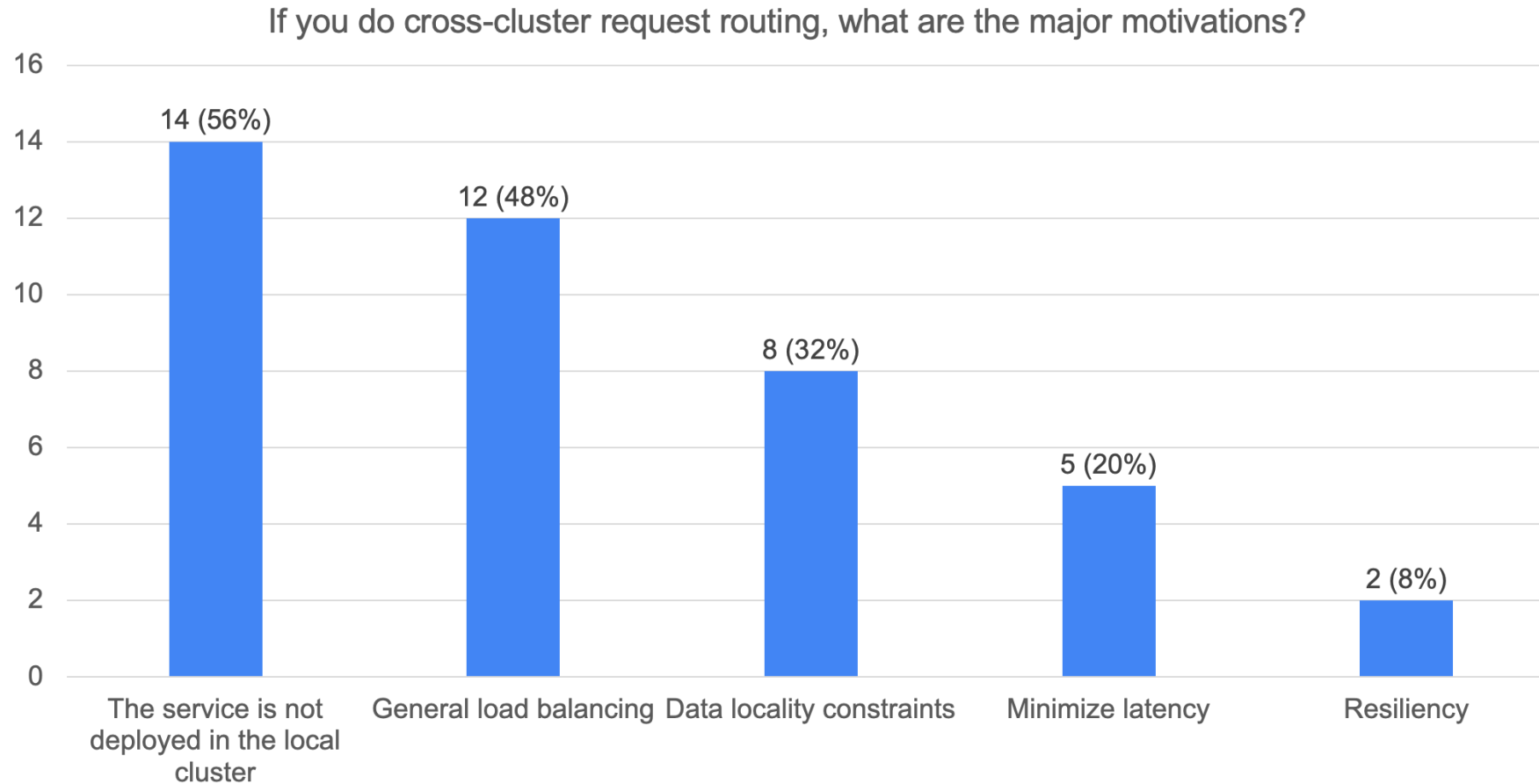
Q4

Does your deployment do cross-cluster request routing?

38 responses



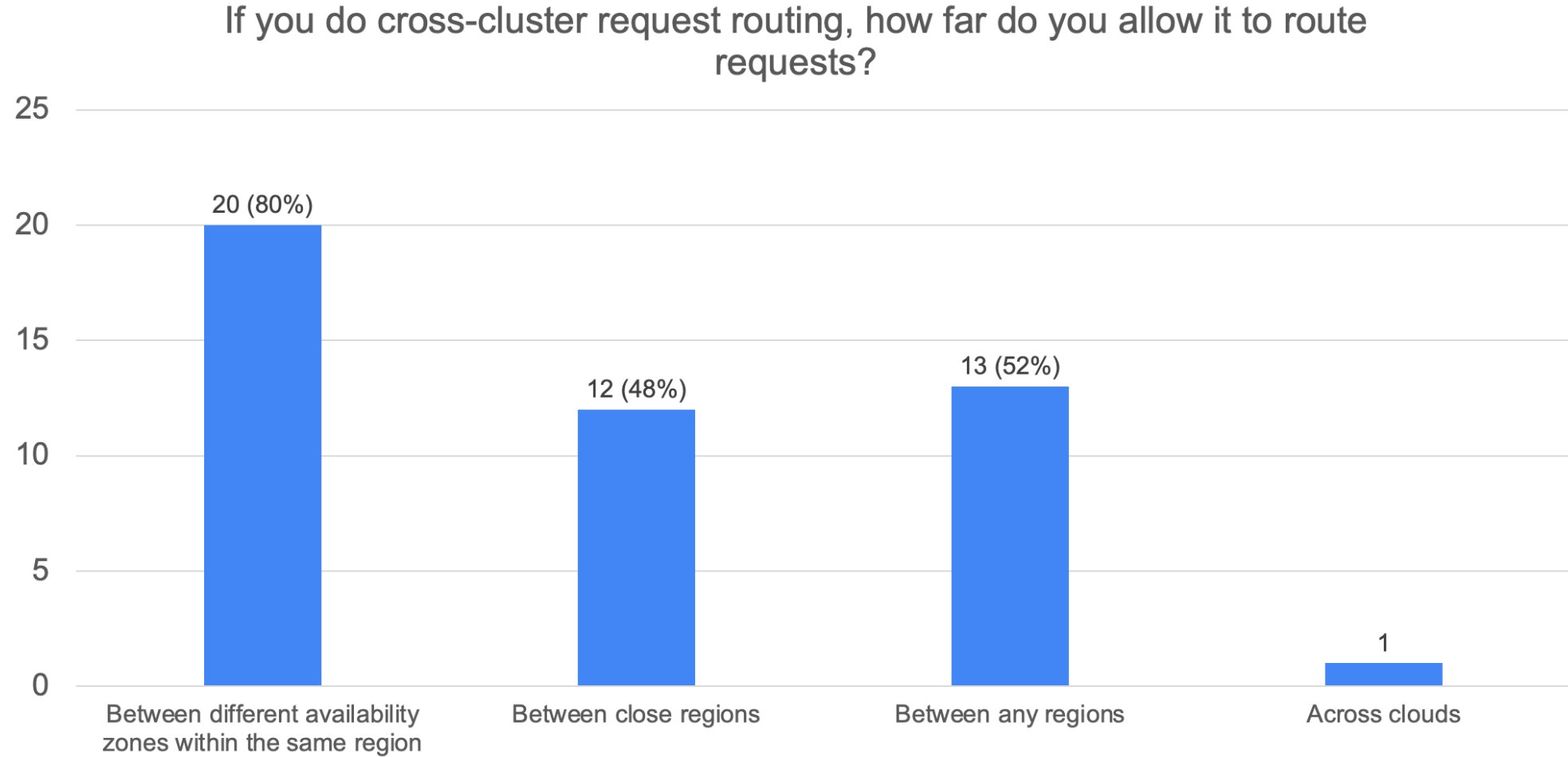
Q5



Other answers

- Migration to new clusters on a regular basis (1)
- Implementing Edge Gateway (cross region load balancing, etc) (2)
- Yet to do (1)
- Other reasons (2)

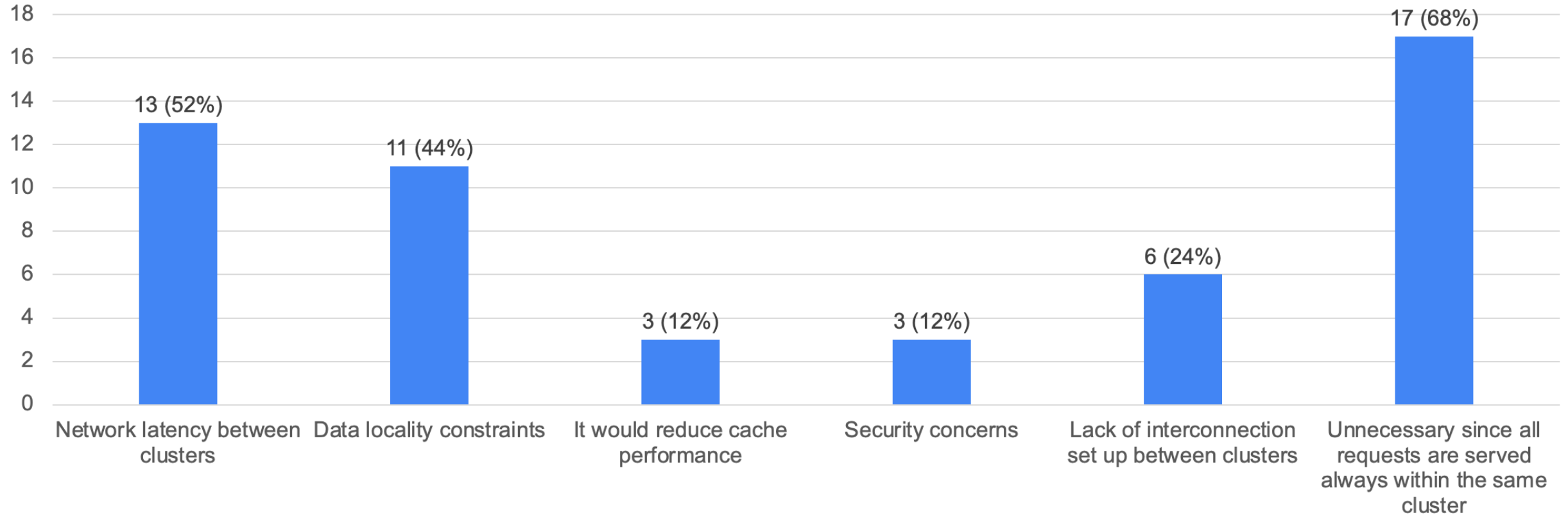
Q6



close region: e.g., us-east-region 1 <-> us-east-region 2
any region: e.g., us-east-region <-> us-west-region

Q7

For your services that do not use cross-cluster request routing, what are the major reasons? (You can choose more than one.)



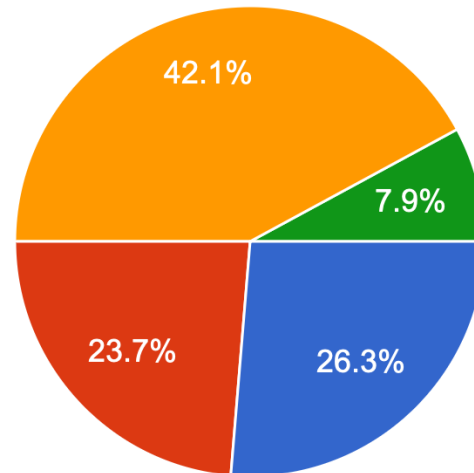
Other reasons

- Reliability
- All services use cross-cluster request routing
- Complexity of setup and knowledge gap of application development teams
- We are exploring the ways to do it

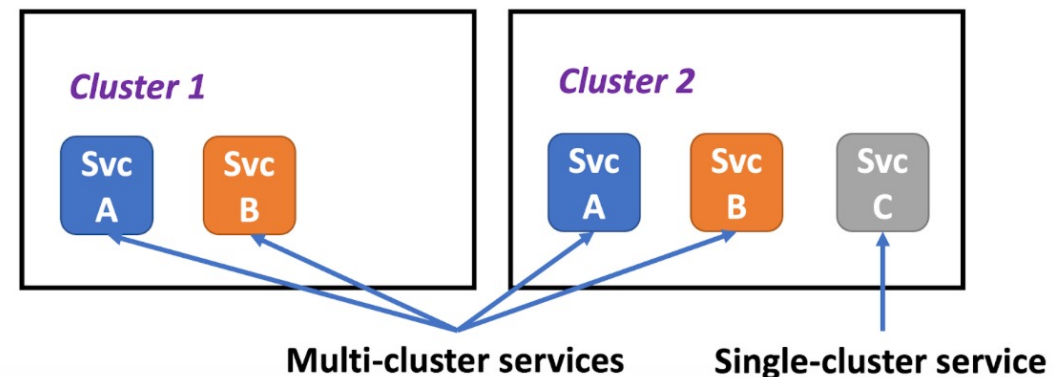
Q8

Do you deploy multi-cluster services? We define a multi-cluster service as a single service that has production replicas running in more than one clust...example of multi-cluster services is shown below.

38 responses

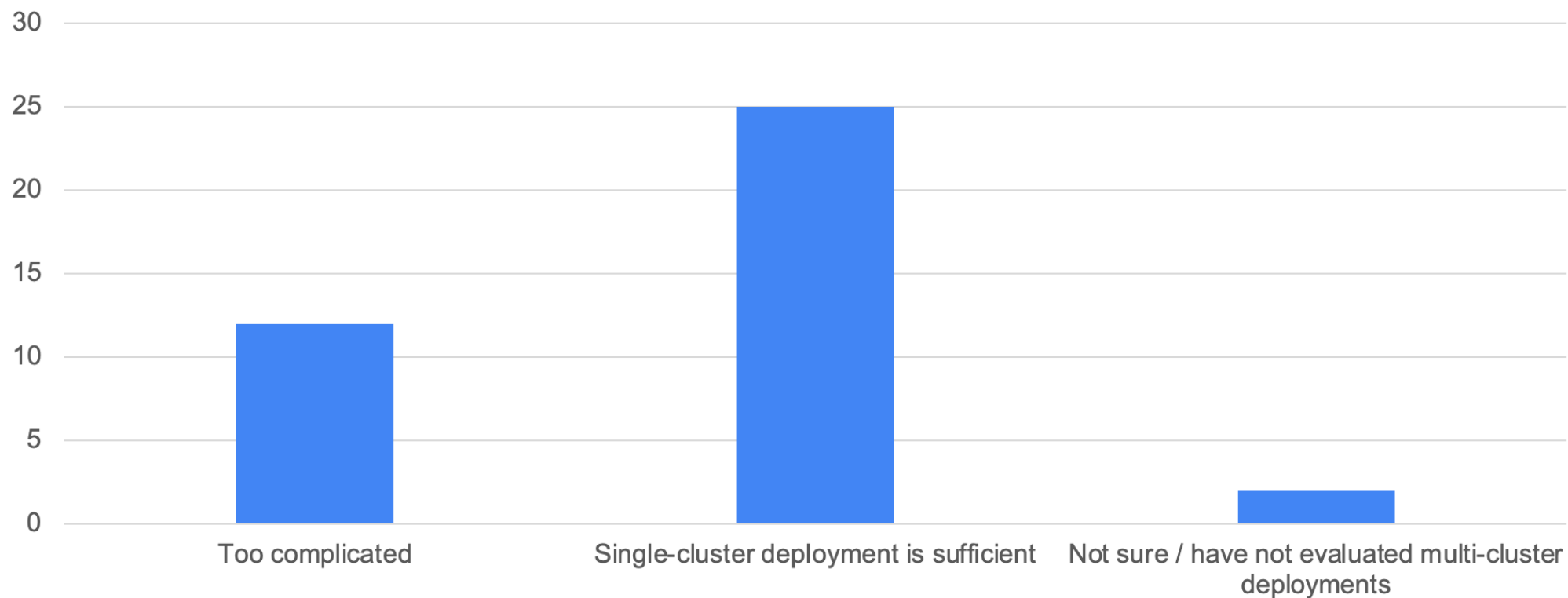


- Yes, most or all of our services are deployed in multiple clusters
- Yes, some of our services are deployed in multiple clusters
- Not now, but we expect multi-cluster deployments in the future
- No, and we don't expect multi-cluster deployments in the future
- Don't know



Q9

For your services which are **not** deployed in multiple clusters, what are the major reasons? (You can choose more than one.)



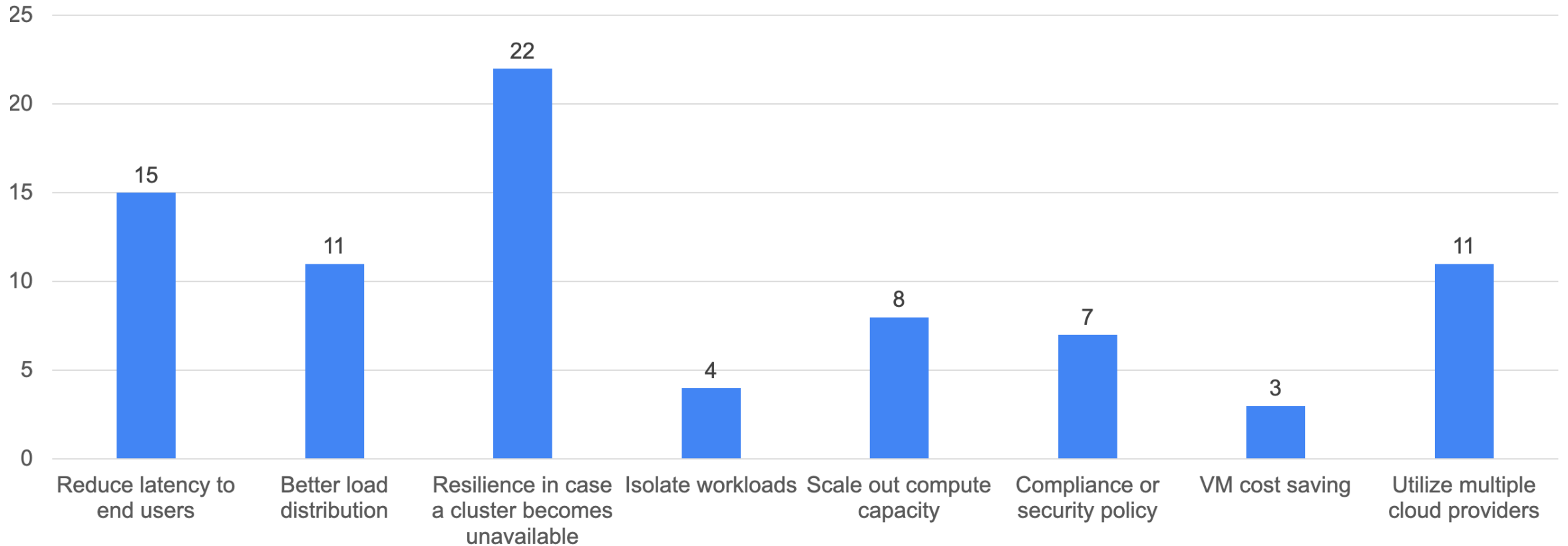
Other answers

- geolocation dns is sufficient for now
- IaC Architecture does not support multiple cluster service setup
- Long distances between datacenters plus customers data concentrated in a single cluster.
- Intentionally segregated for security or usage reasons
- some services are ok to be deployed in a single AZ
- Service Mesh (Istio) implementation is in an early stage yet

If you don't deploy multi-cluster services and don't plan to deploy them in the future, you can stop the survey [here](#) and submit the response.

Q10

For your services which are or will be deployed multi-cluster, what are the major reasons? (You can choose more than one.)



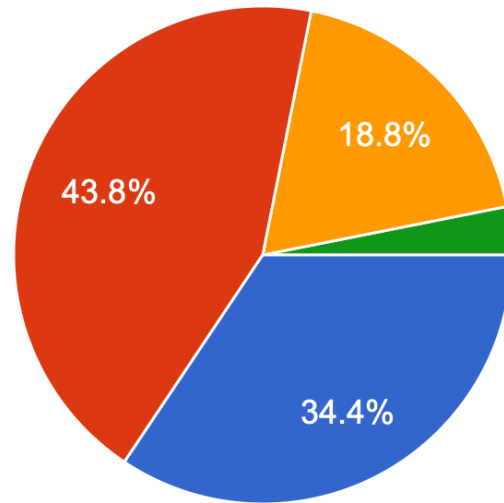
Other answers

- Migration to new clusters regularly
- Consolidating load balancers
- Cutting costs
- Safer Kubernetes upgrades - move away from in place upgrade to A/B cluster upgrade.

Q11

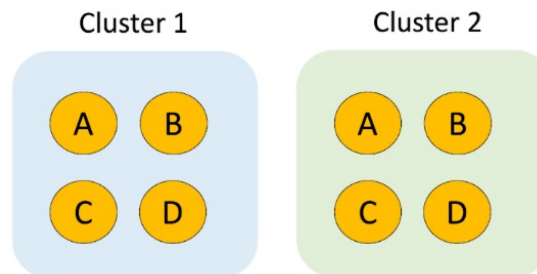
How do or will you deploy your multi-cluster services?

32 responses

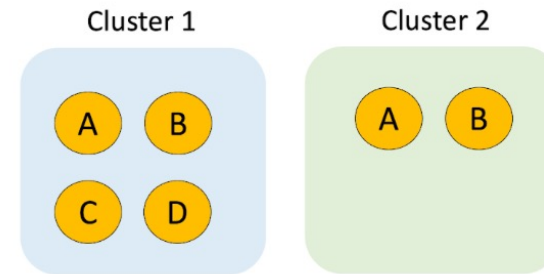


- Complete replication: Each cluster contains all of the application's microservices
- Partial replication: Some clusters contain only some of the application's microservices (and therefore, some requests may need to be routed across clusters)
- Both
- Don't know

Complete replication



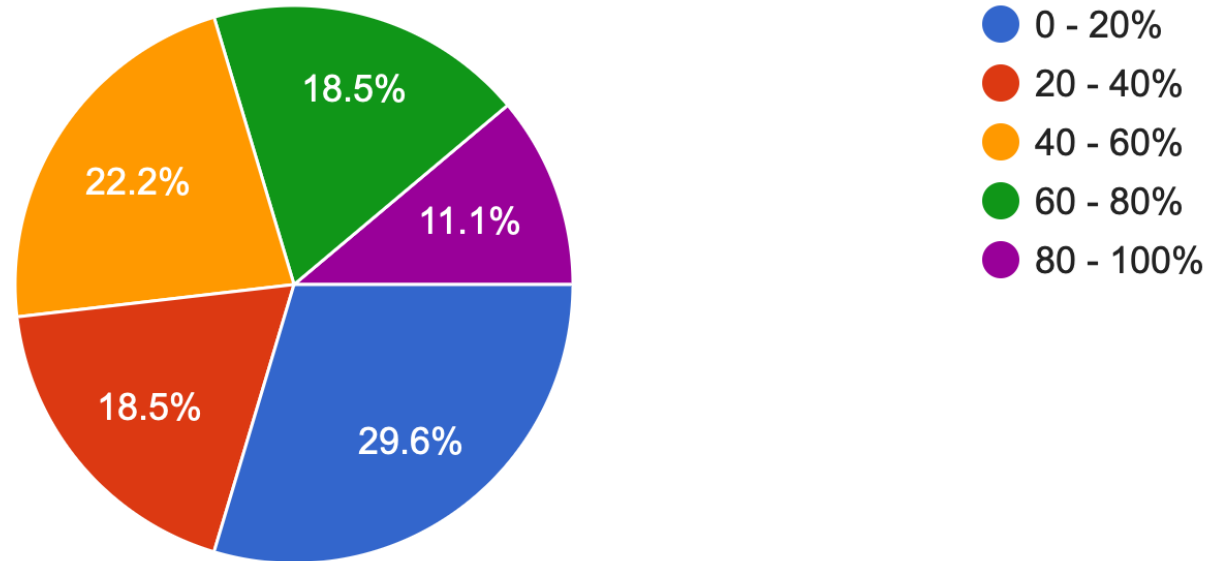
Partial replication



Q12

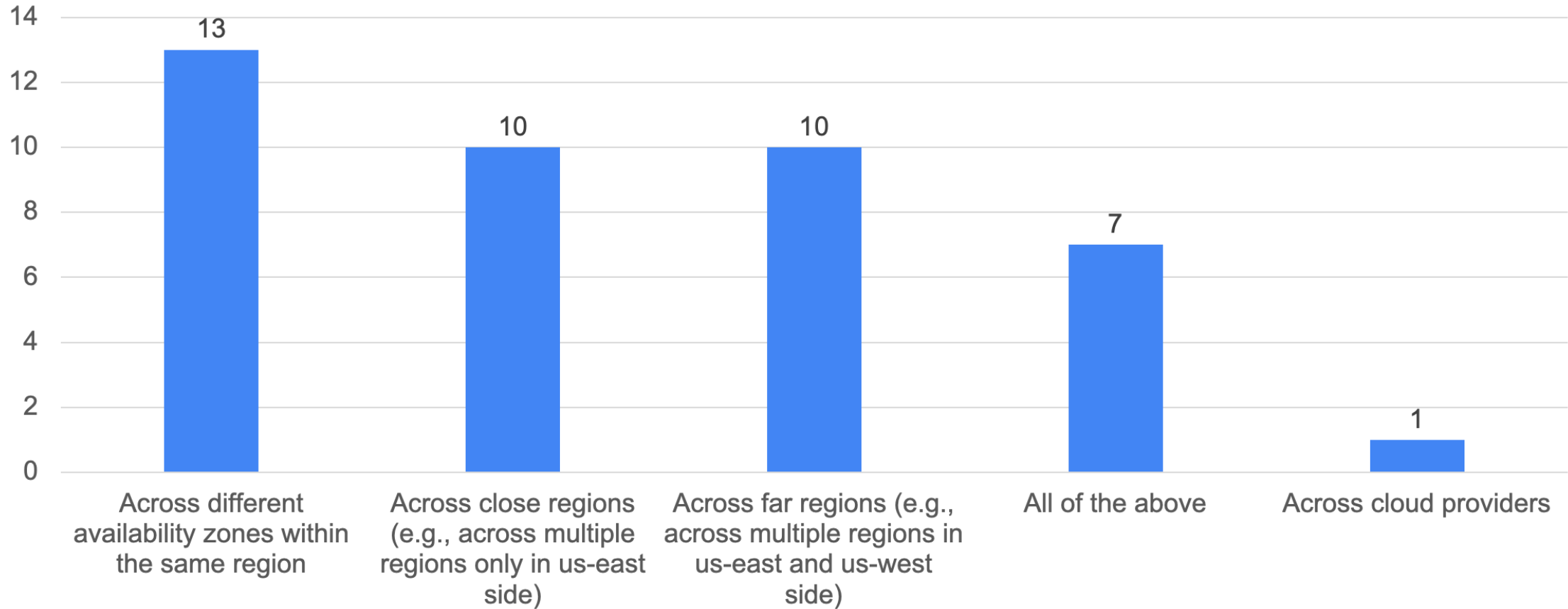
How much percentage of your services are multi-cluster services?

27 responses



Q13

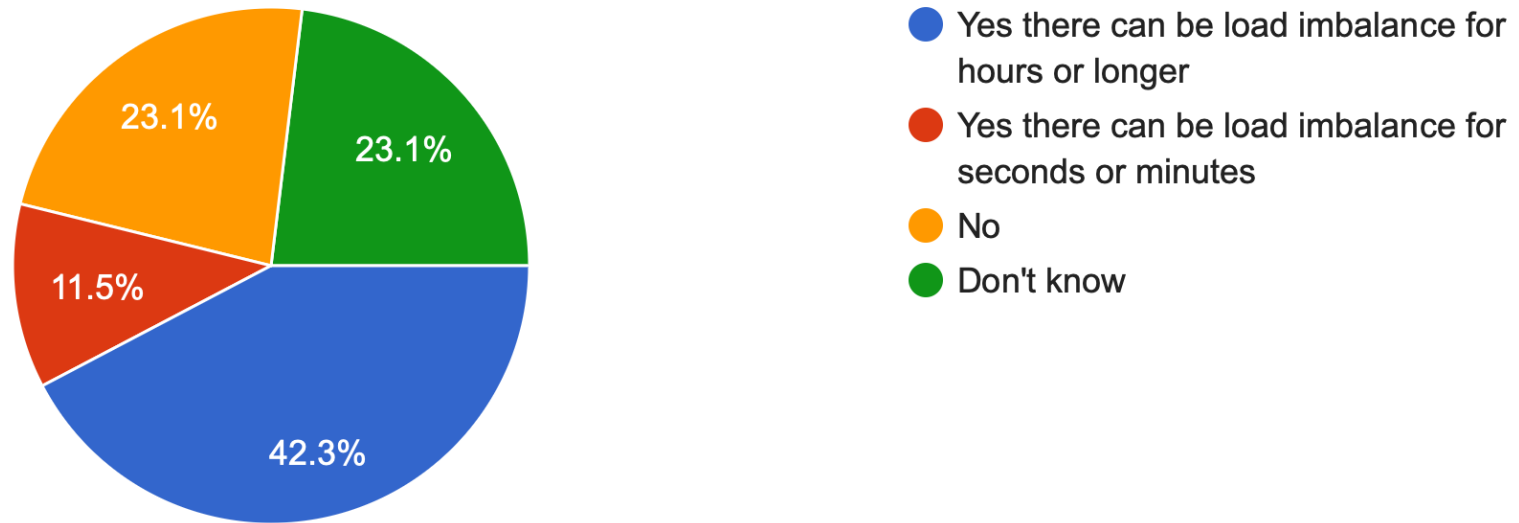
Where do or will you deploy your multi-cluster services?



Q14

Is there considerable imbalance in load between clusters in your multi-cluster services? (you can ignore this question if you have not deployed multi-cluster services yet)

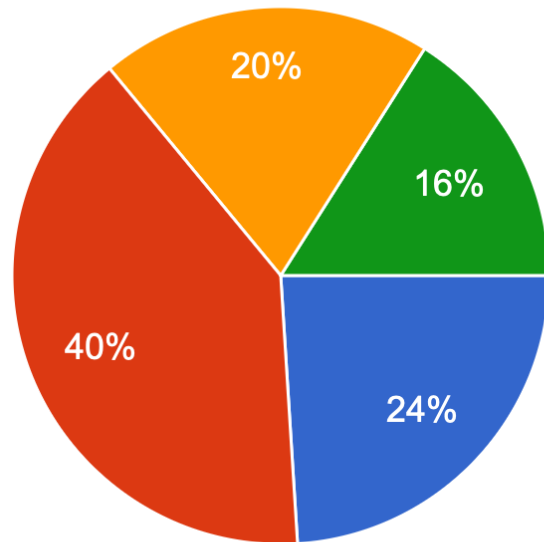
26 responses



Q15

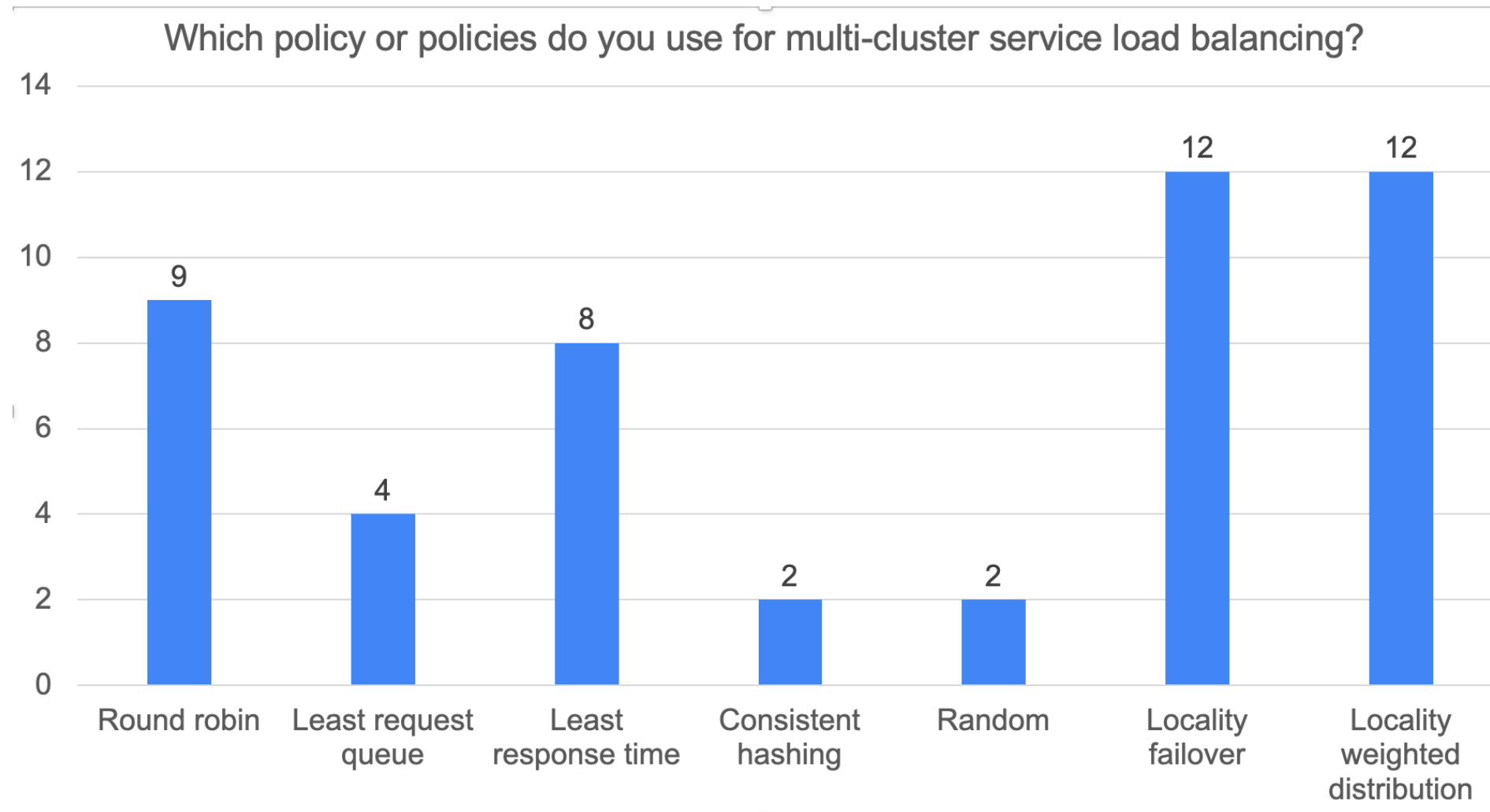
Is there difference in cost to serve requests between clusters in your multi-cluster services? e.g., VM price difference, network bandwidth cost difference... you have not deployed multi-cluster services yet)

25 responses



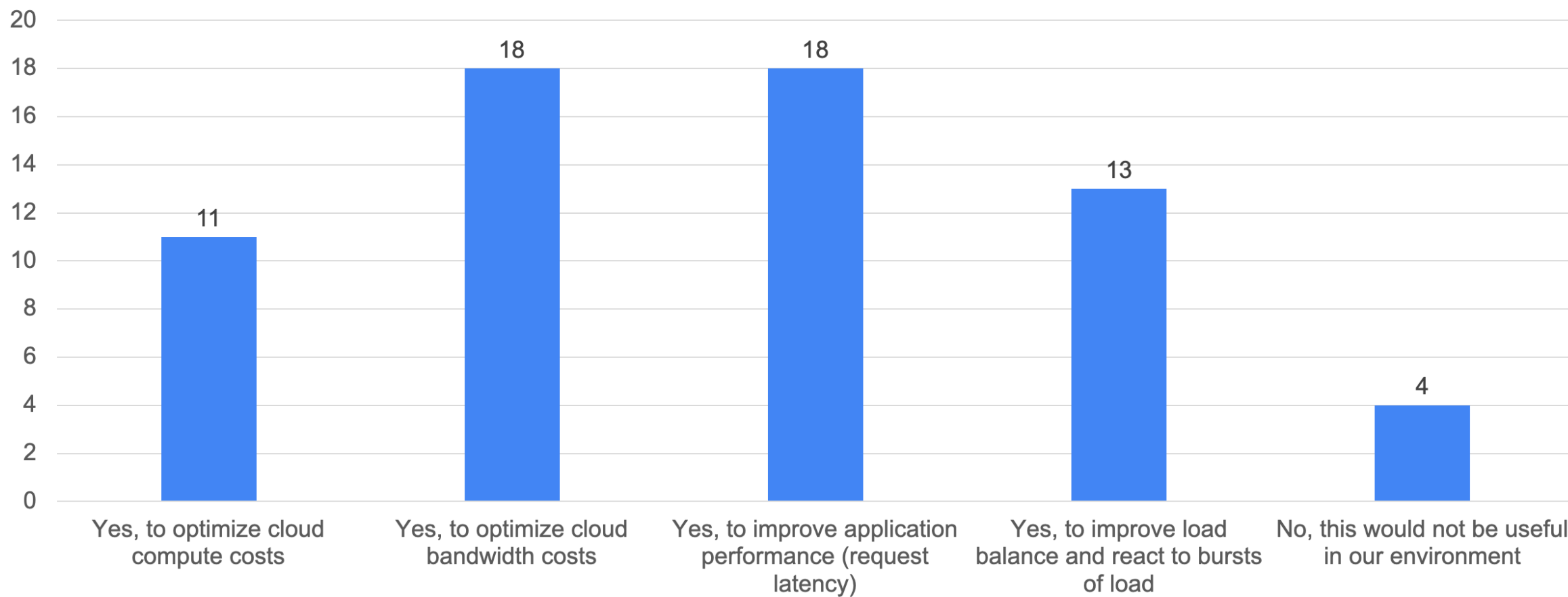
- Yes, in terms of compute cost (e.g., VM price), serving request in one cluster can cost less than other clusters.
- Yes, in terms of network bandwidth cost, serving request in one cluster can cost less than other clusters.
- No
- Don't know

Q16

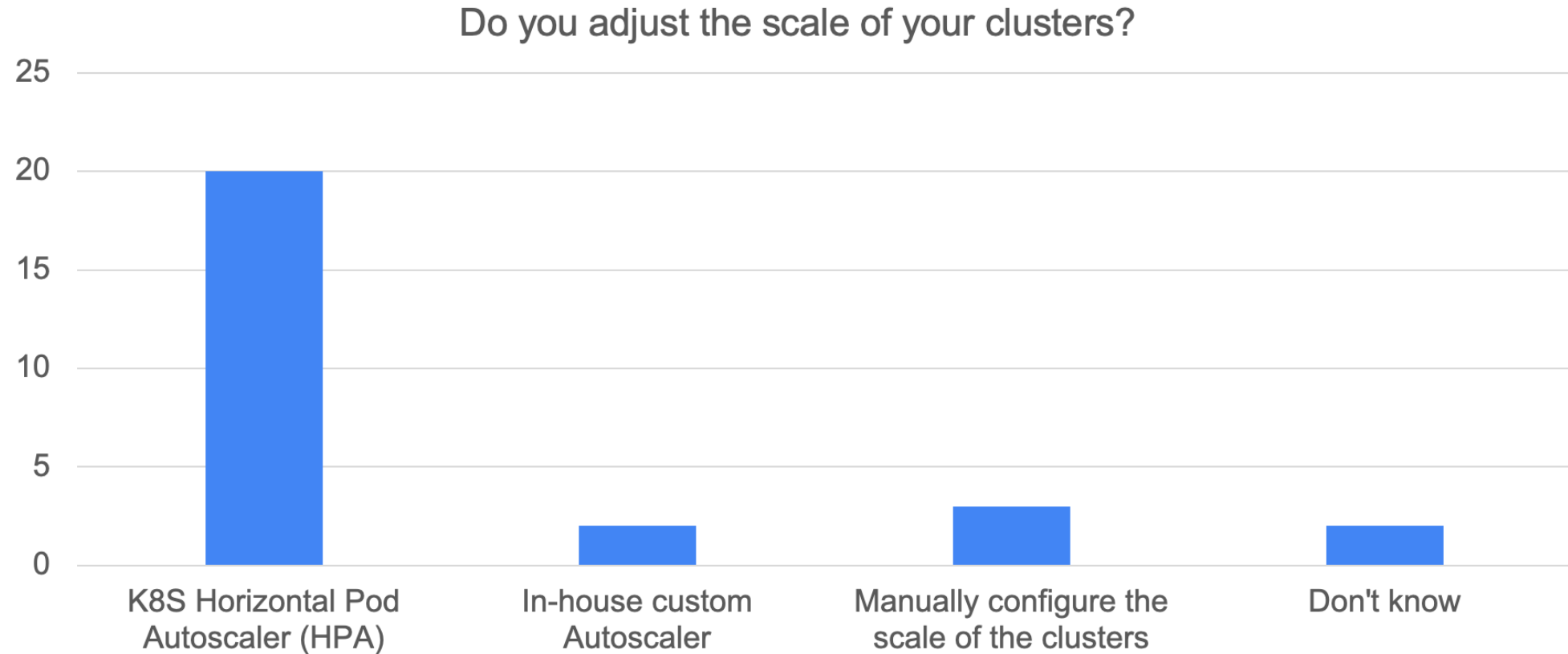


Q17

Suppose your cluster infrastructure systems could automatically optimize the **cross-cluster routing among multi-cluster services**. Would this be useful in your environment? (Check all that apply.)



Q18



Other answers

- Node autoscaler (2)
- Karpenter (2)
- KEDA (2)
- Spot.io (1)