

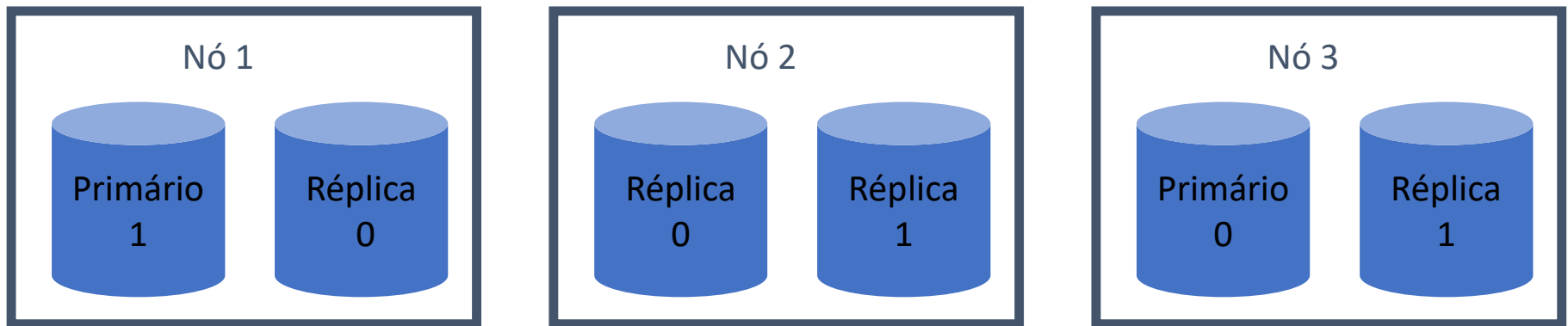
# Um índice é dividido em shards.

Cada shard pode estar em um diferente nó do cluster

Cada shard contém um índice do Lucene

# Shards primários e réplicas

Este índice tem dois shard primários e duas réplicas  
A aplicação deve alternar as requisições entre eles



Requisições de leitura são direcionadas ao shard primário, então replicadas  
Requisições de leitura são direcionadas ao shard primário ou qualquer réplica

## Quantos Shards eu preciso?

- Você pode adicionar mais Shards a qualquer momento sem reindexação
- Mas você não pode usa-los ilimitado, você pode inicialmente fazer 1000 em um mesmo nó: overhead
- Você quer super alocar, mas não com exagero
- Considere escalar em fases, assim você tem tempo de reindexar antes da próxima etapa



# Não existe reposta mágica

---

- Tipo de dados
- Balanço entre leitura e escrita
- Complexidade do schema
- Indexação
- Análises

# Como fazer?

- Comece com um único servidor usando o mesmo hardware de Produção, com um shard e sem replicação
- Importe com documentos reais e execute consultas reais
- "Force" até seu limite: agora você sabe a capacidade de um único shard

Lembre-se:  
replicas  
podem ser  
adicionadas



Aplicações com muita leitura  
podem adicionar mais replicas  
sem reindexação



Porém observe que isso só  
ajuda se as novas replicas  
estiverem em outro hardware