

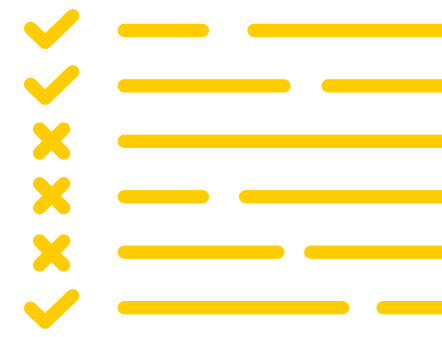
Yandex

Big Data

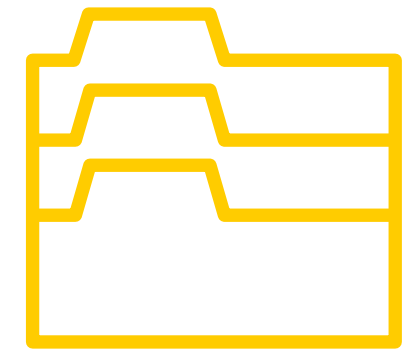
Scaling Distributed File System



Databases



Transaction
logs



Document
stores

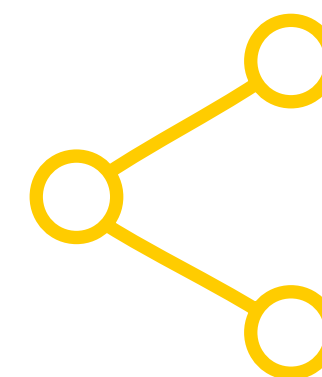
Where does Big Data come from?



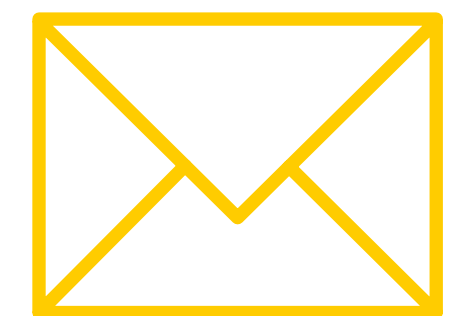
Instant
messages



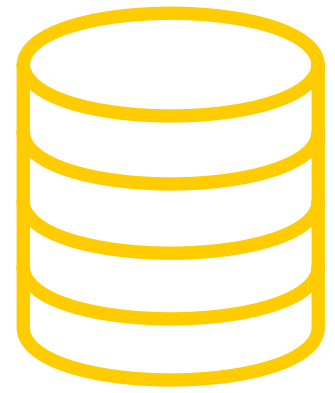
Videos



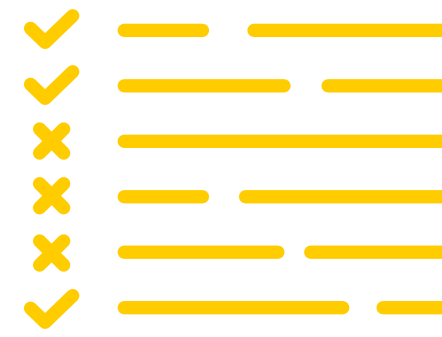
Social media



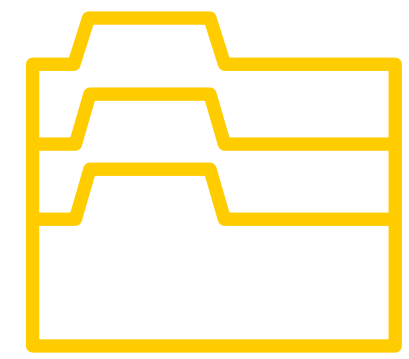
Emails



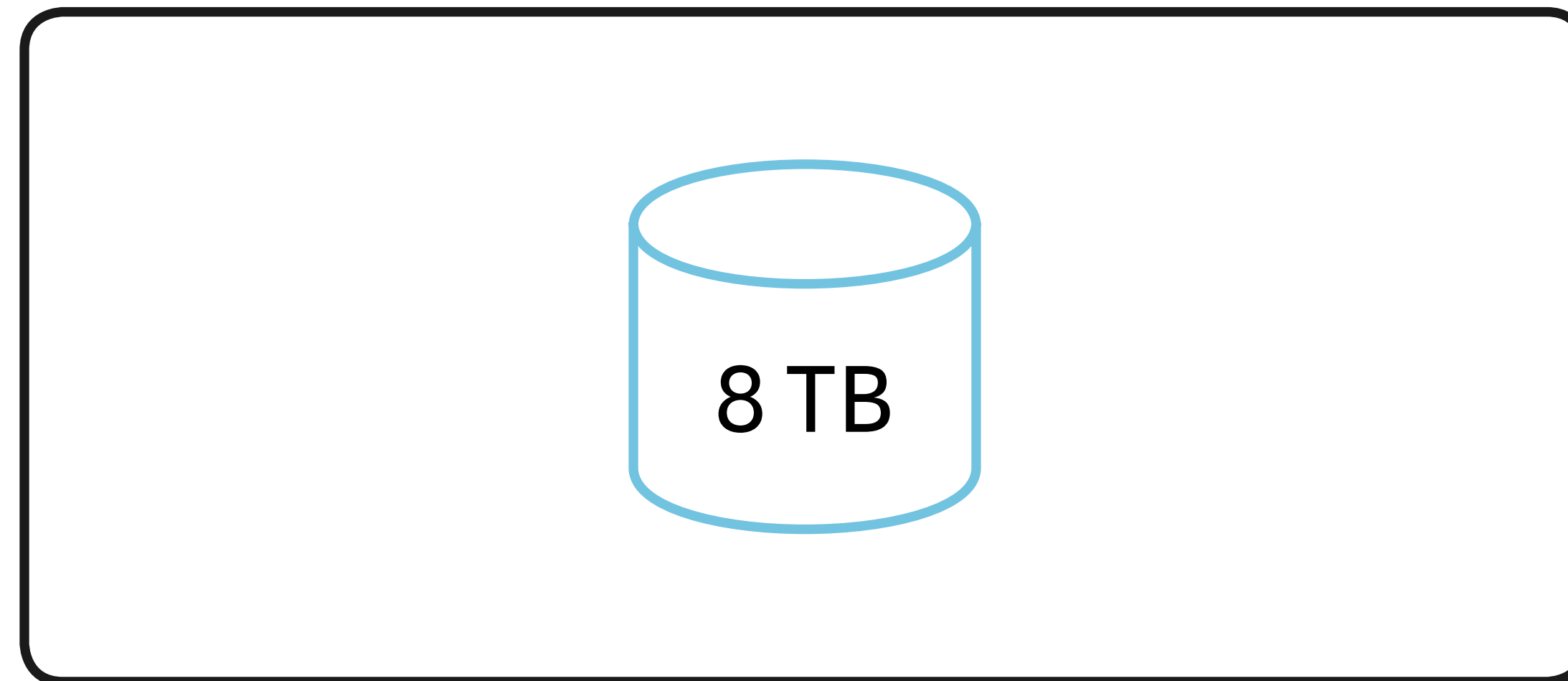
Databases



Transaction
logs



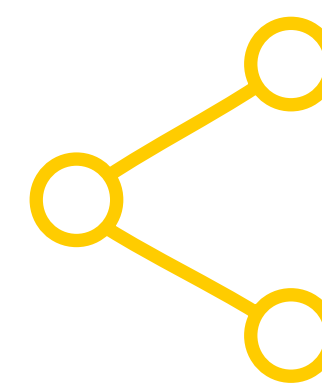
Document
stores



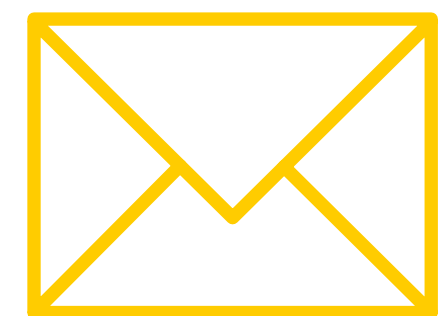
Instant
messages



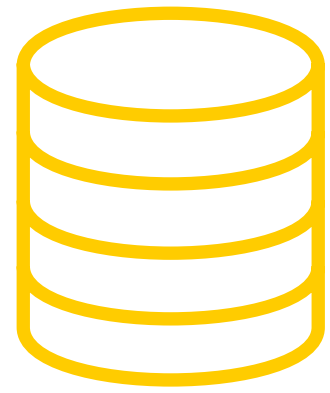
Videos



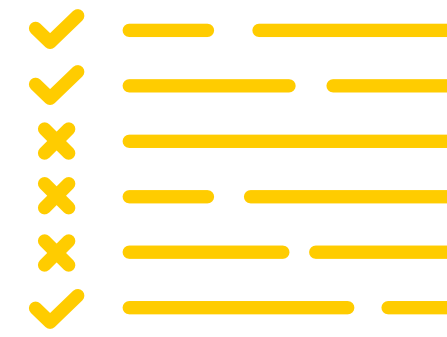
Social media



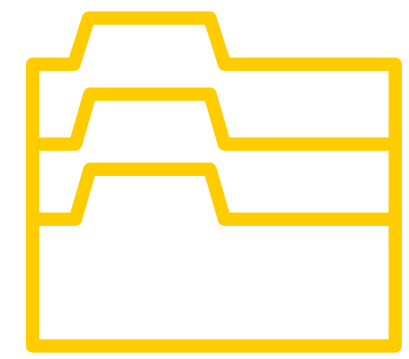
Emails



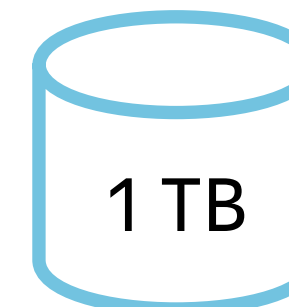
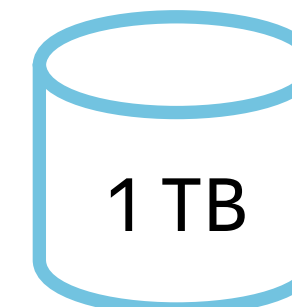
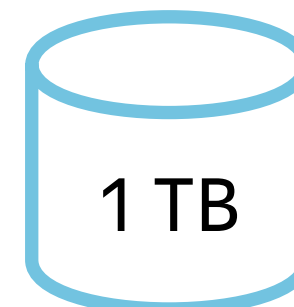
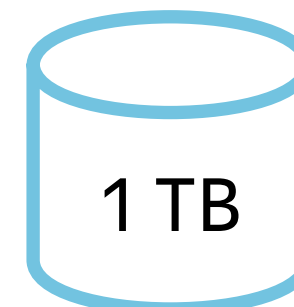
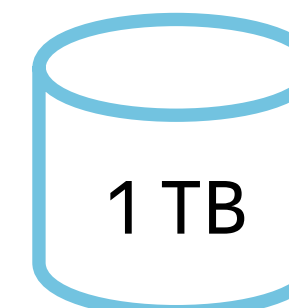
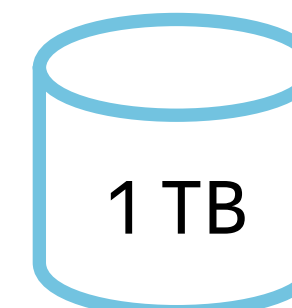
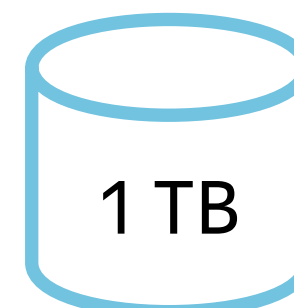
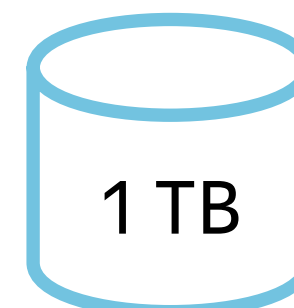
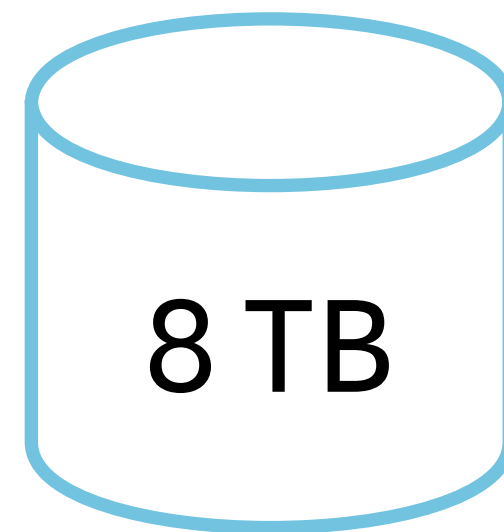
Databases



Transaction
logs



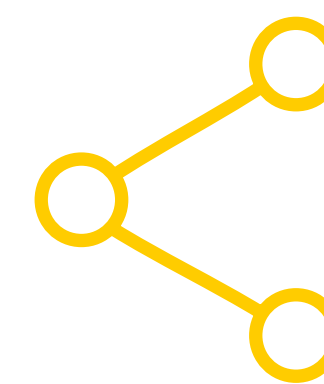
Document
stores



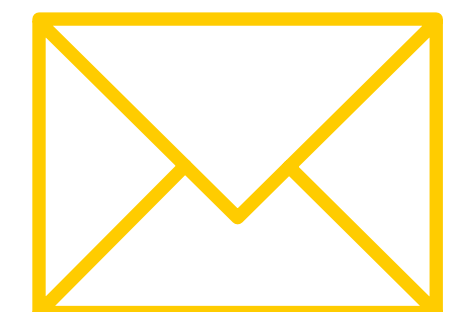
Instant
messages



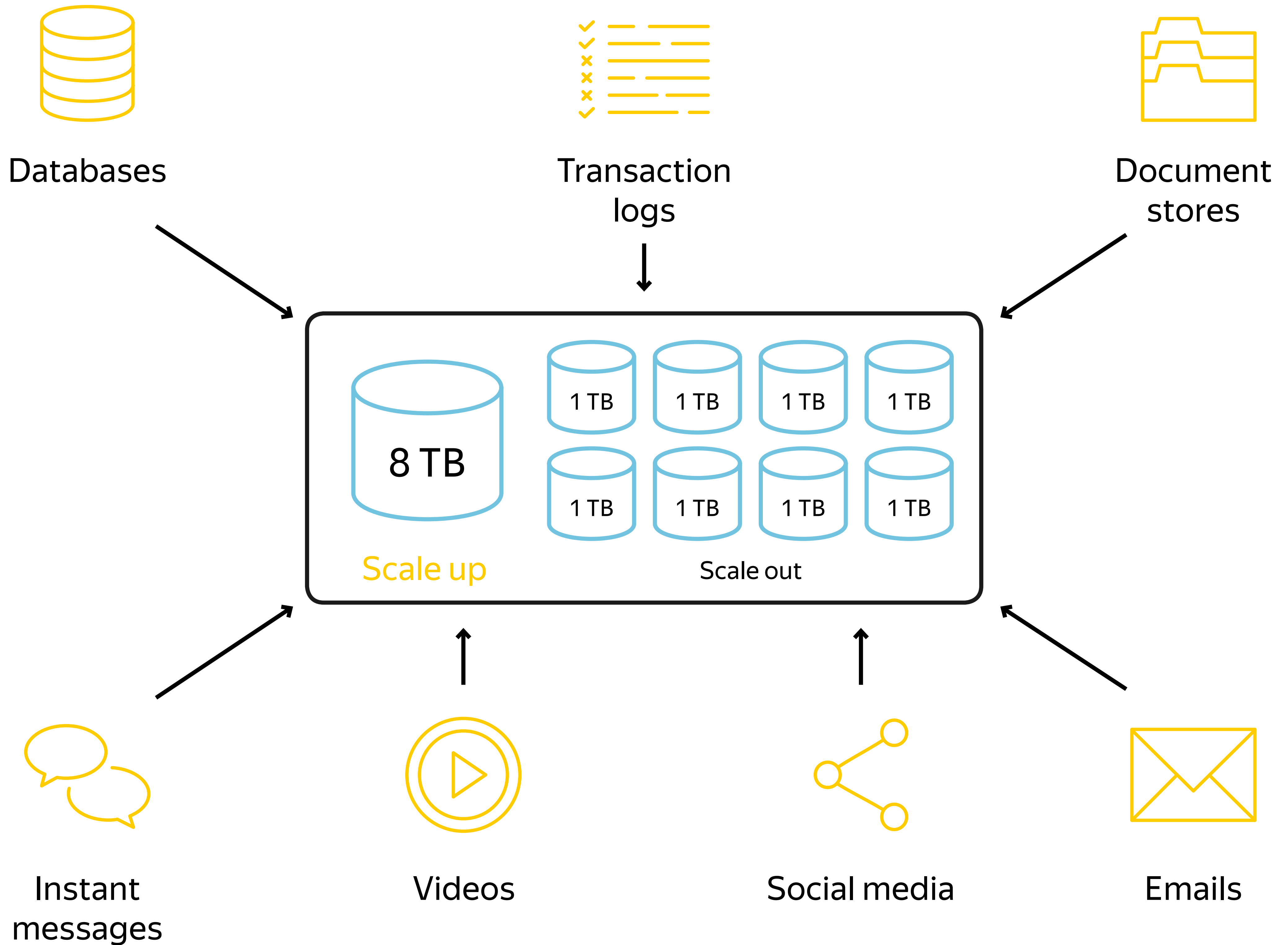
Videos

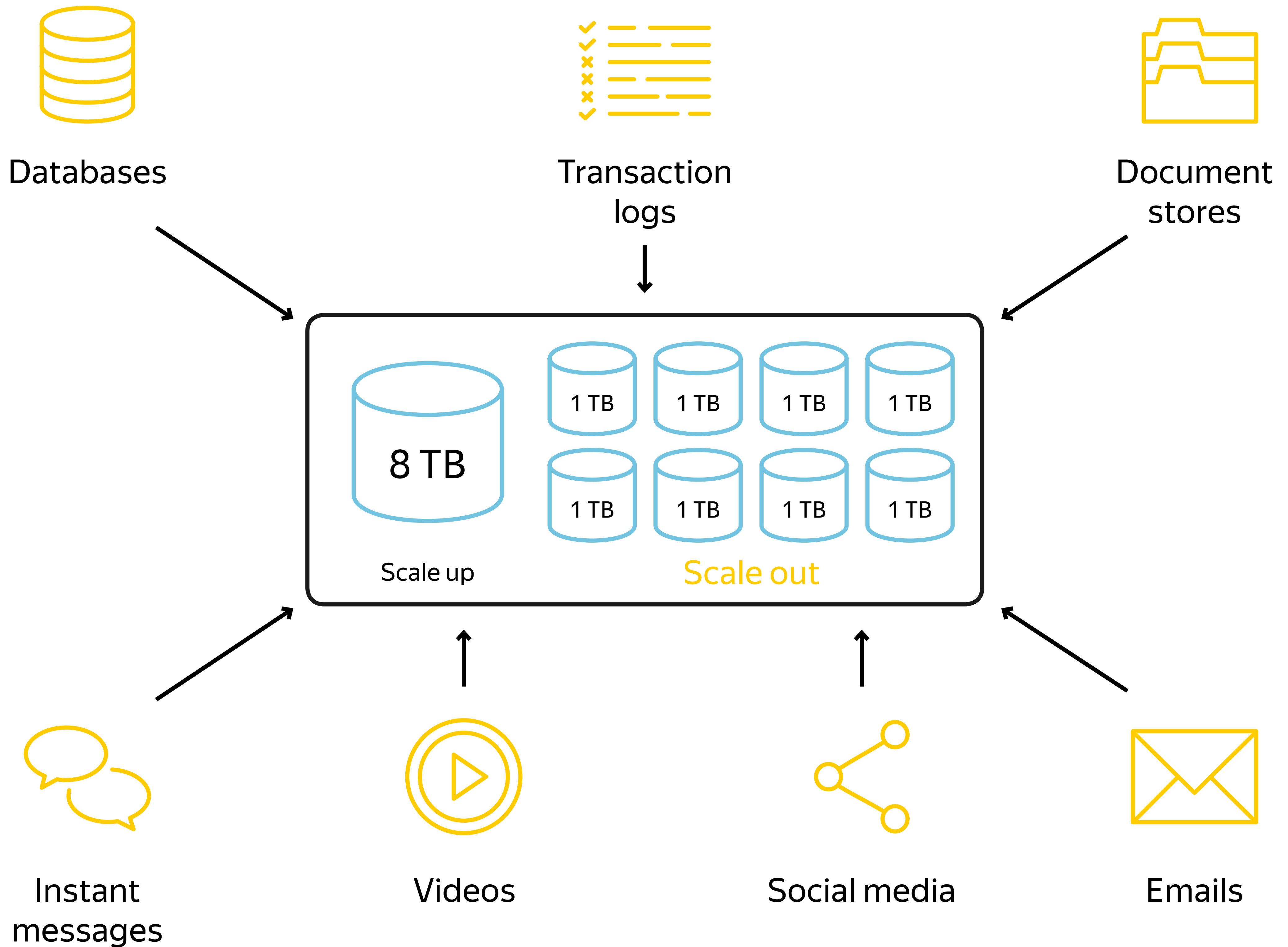


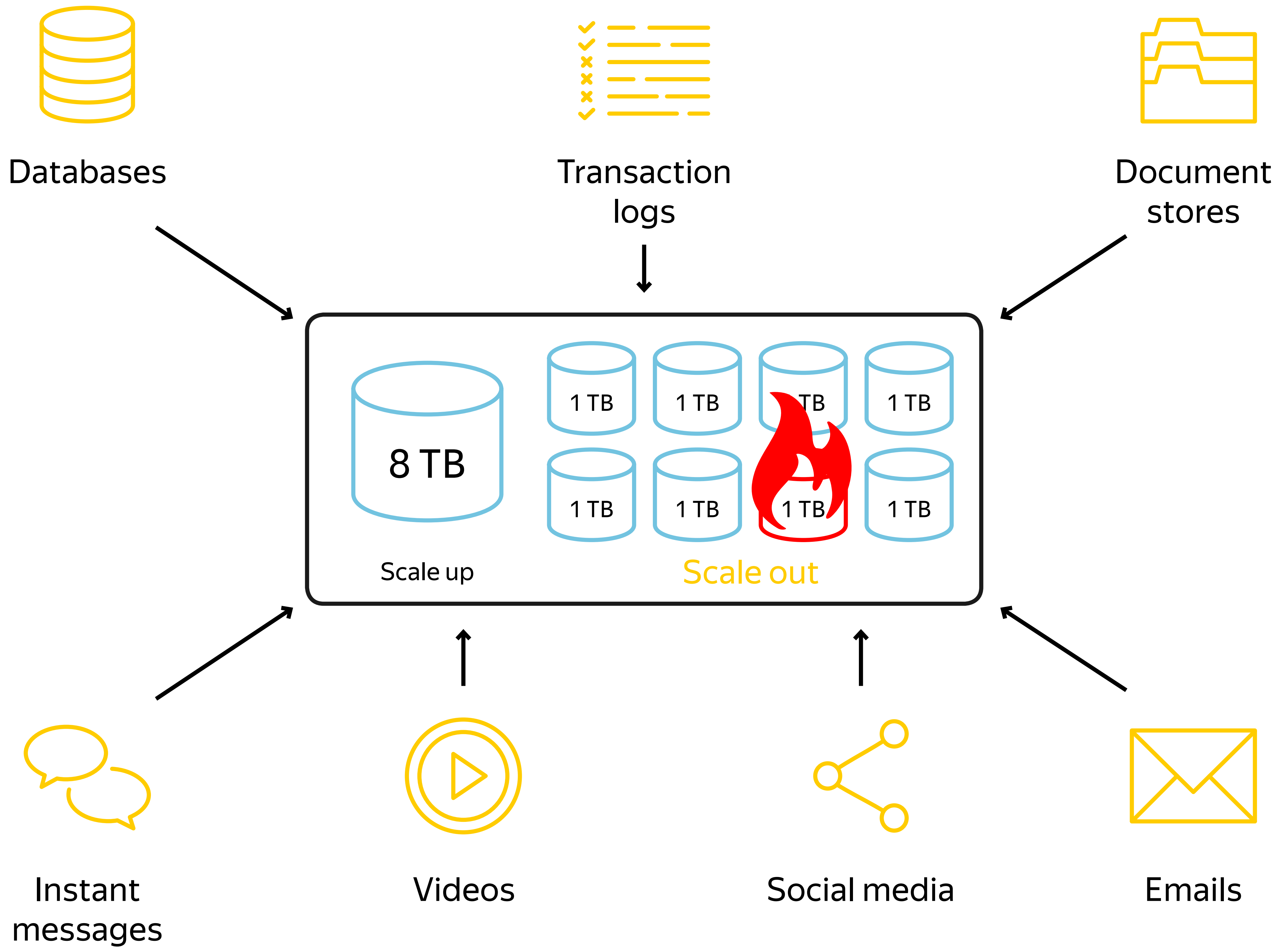
Social media



Emails







The Google File System

Sanjay Ghemawat, Howard Gobioff, and Shun-Tak Leung

Google*

ABSTRACT

We have designed and implemented the Google File System, a scalable distributed file system for large distributed data-intensive applications. It provides fault tolerance while running on inexpensive commodity hardware, and it delivers high aggregate performance to a large number of clients.

While sharing many of the same goals as previous distributed file systems, our design has been driven by observations of our application workloads and technological environment, both current and anticipated, that reflect a marked departure from some earlier file system assumptions. This has led us to reexamine traditional choices and explore radically different design points.

The file system has successfully met our storage needs. It is widely deployed within Google as the storage platform for the generation and processing of data used by our service as well as research and development efforts that require large data sets. The largest cluster to date provides hundreds of terabytes of storage across thousands of disks on over a thousand machines, and it is concurrently accessed by hundreds of clients.

In this paper, we present file system interface extensions designed to support distributed applications, discuss many aspects of our design, and report measurements from both microbenchmarks and real world use.

1. INTRODUCTION

We have designed and implemented the Google File System (GFS) to meet the rapidly growing demands of Google's data processing needs. GFS shares many of the same goals as previous distributed file systems such as performance, scalability, reliability, and availability. However, its design has been driven by key observations of our application workloads and technological environment, both current and anticipated, that reflect a marked departure from some earlier file system design assumptions. We have reexamined traditional choices and explored radically different points in the design space.

First, component failures are the norm rather than the exception. The file system consists of hundreds or even thousands of storage machines built from inexpensive commodity parts and is accessed by a comparable number of client machines. The quantity and quality of the components virtually guarantee that some are not functional at any given time and some will not recover from their current failures. We have seen problems caused by application bugs, hardware errors, and the failures of disks, memory, connectors, networking, and power supplies. Therefore, constant monitoring, error detection, fault tolerance, and automatic recovery must be integral to the system.

The Google File System, Symposium on Operating Systems Principles (SOSP, 2003)

GFS Key Components

GFS Key Components

- › components failures are a norm (\rightarrow replication)

GFS Key Components

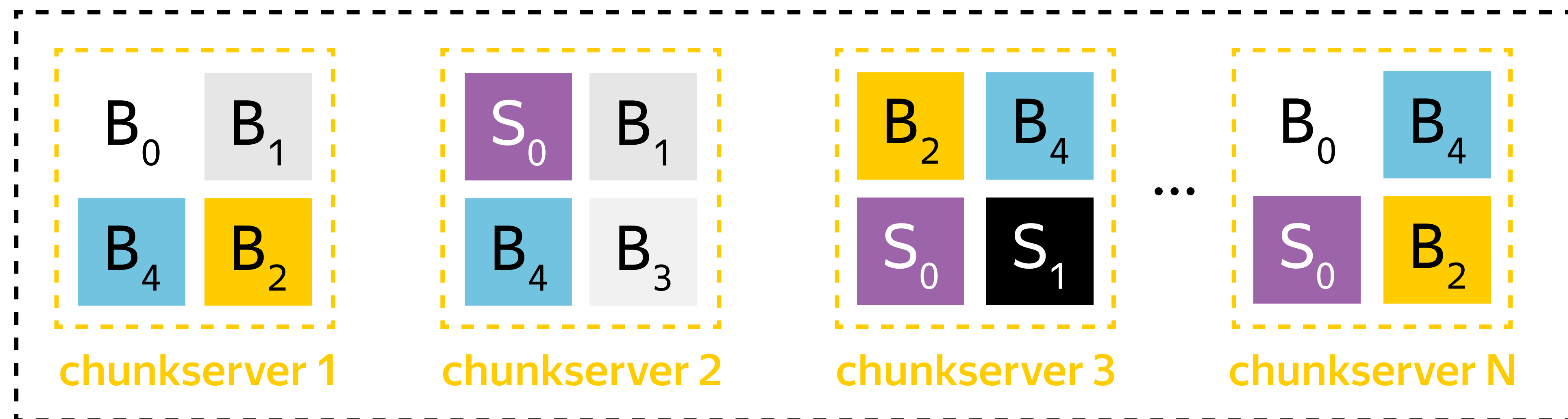
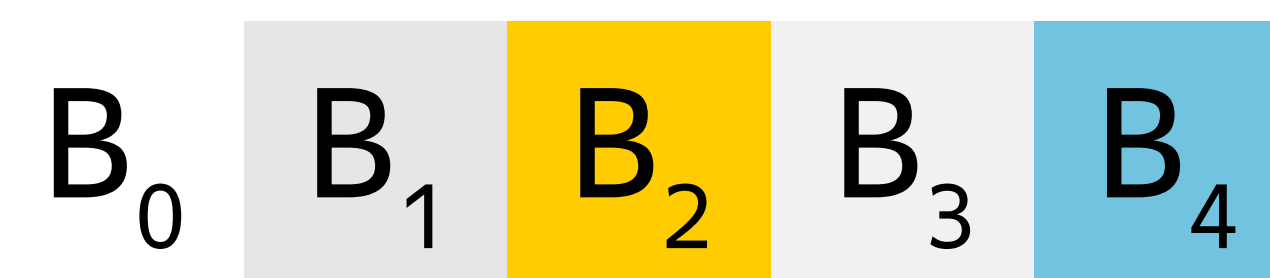
- › components failures are a norm (→ replication)
- › even space utilisation

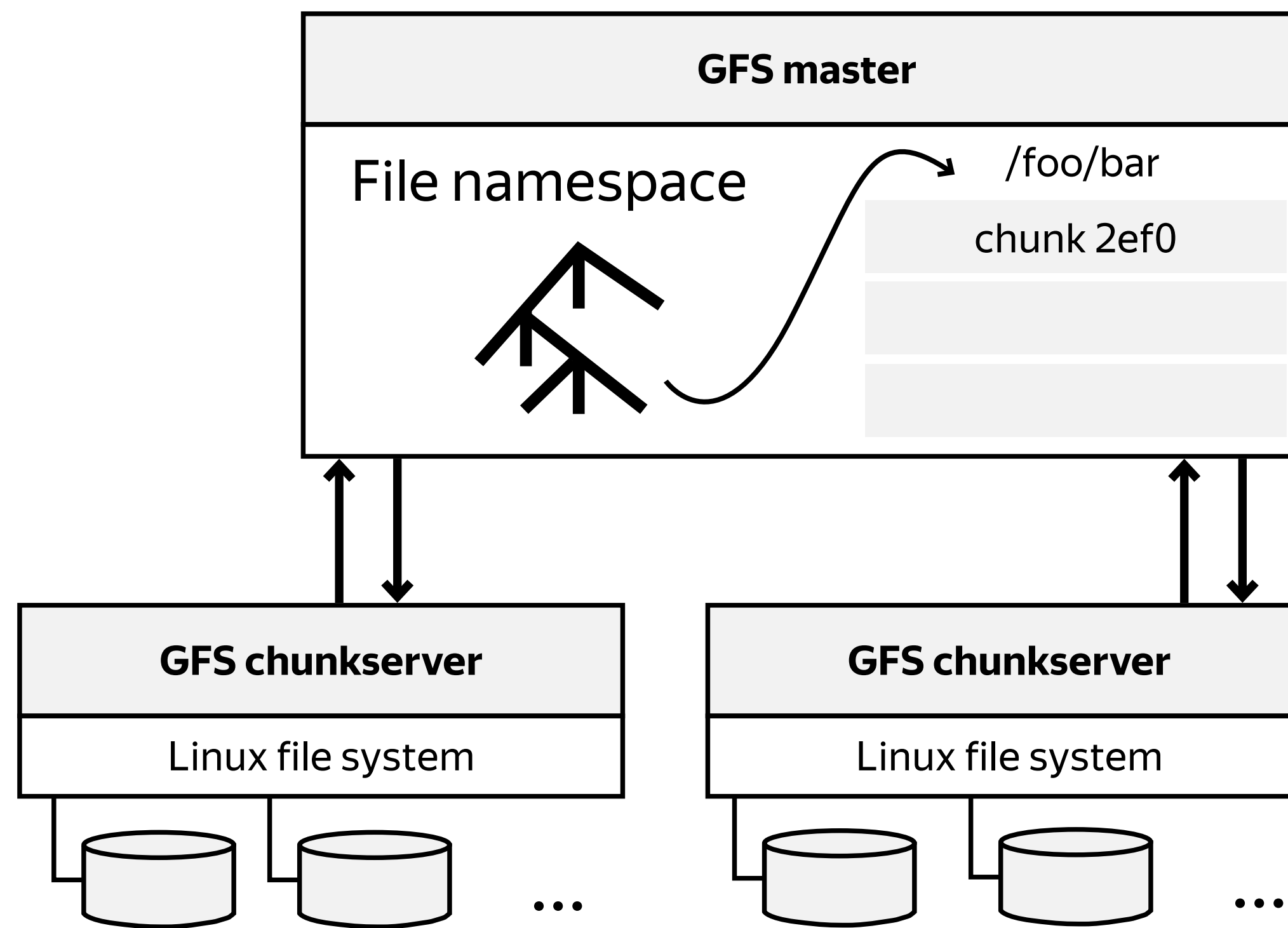
GFS Key Components

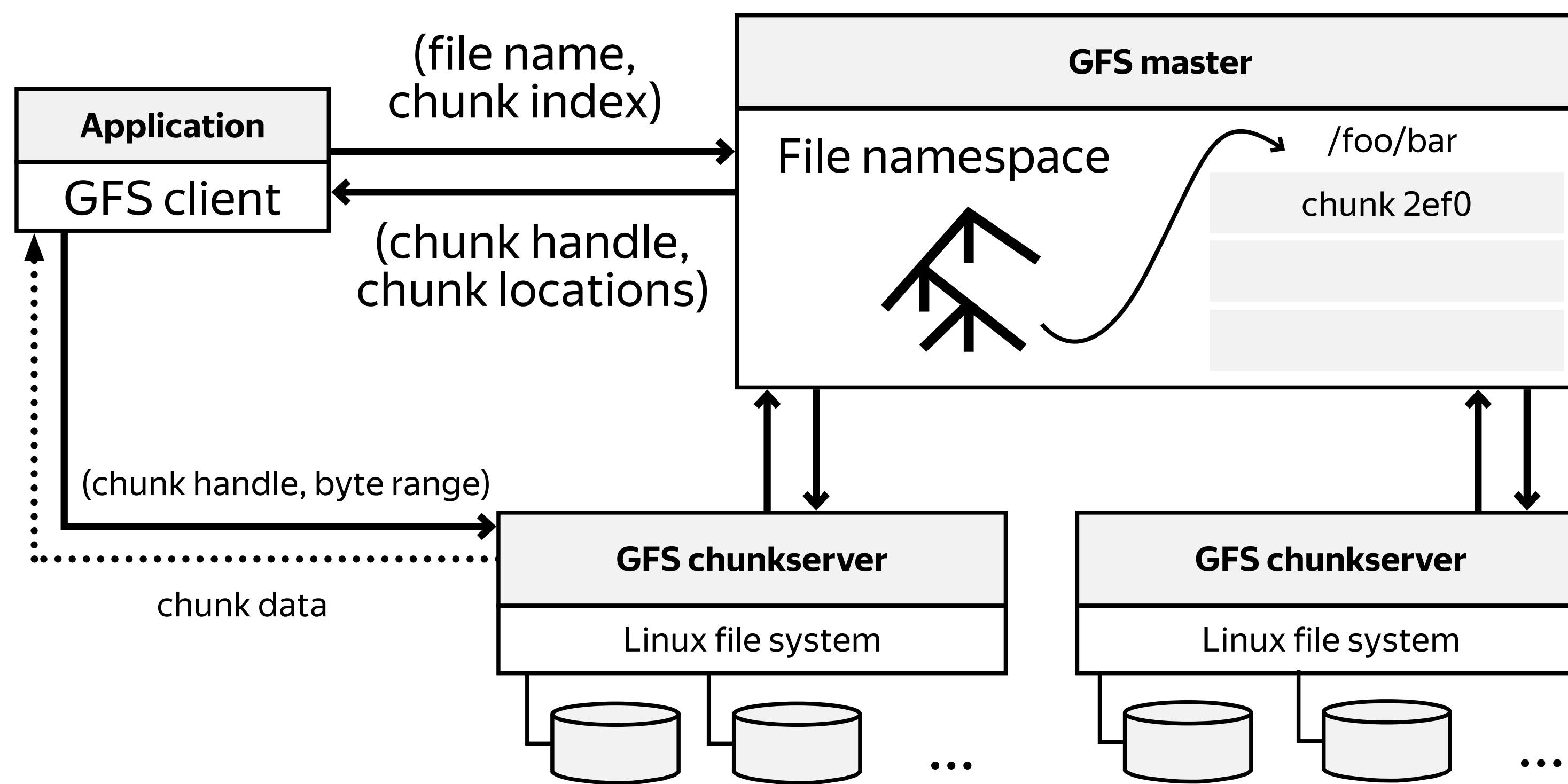
- › components failures are a norm (→ replication)
- › even space utilisation
- › write-once-read-many

Replication

S.txt + B.txt:

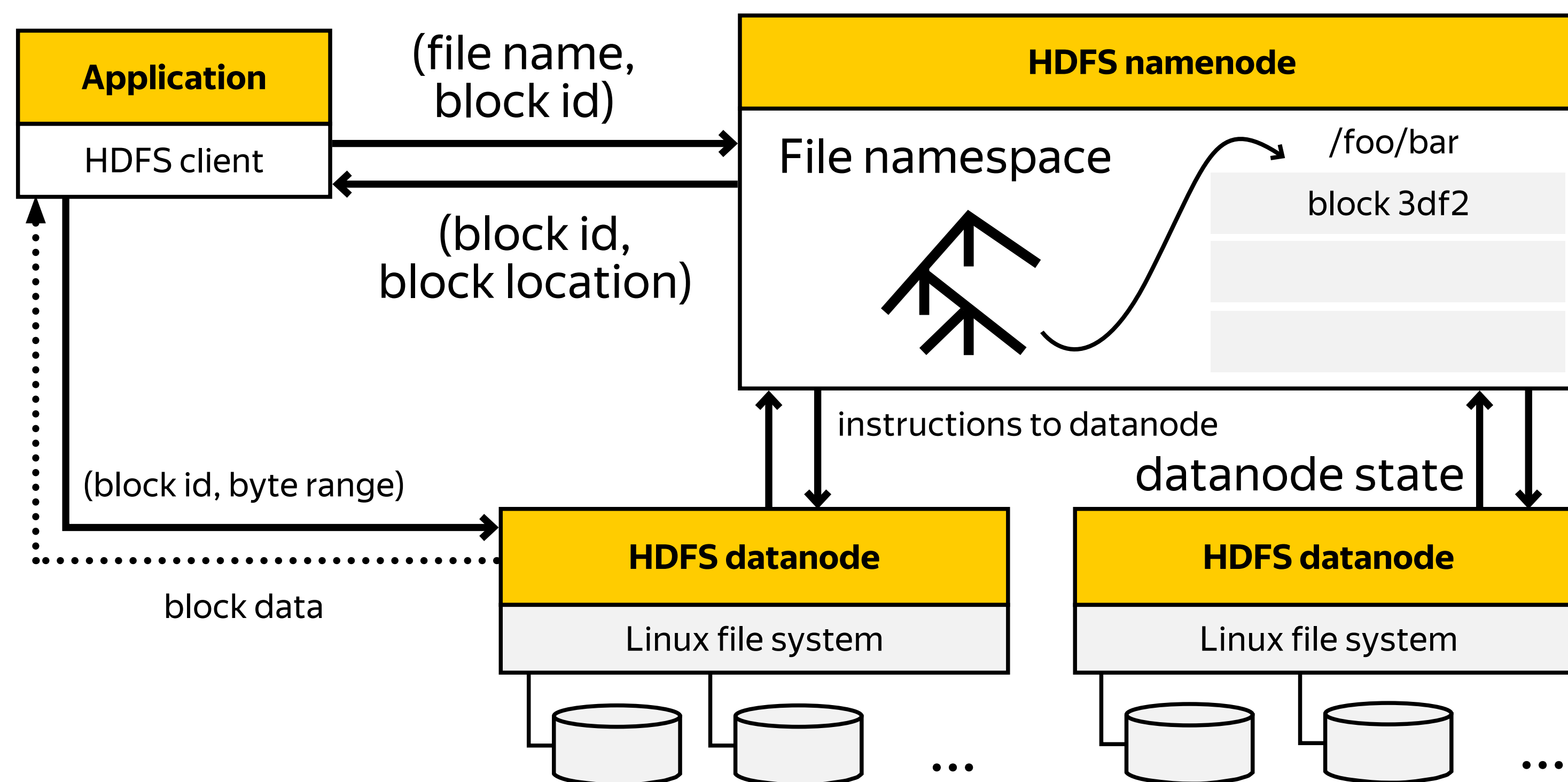




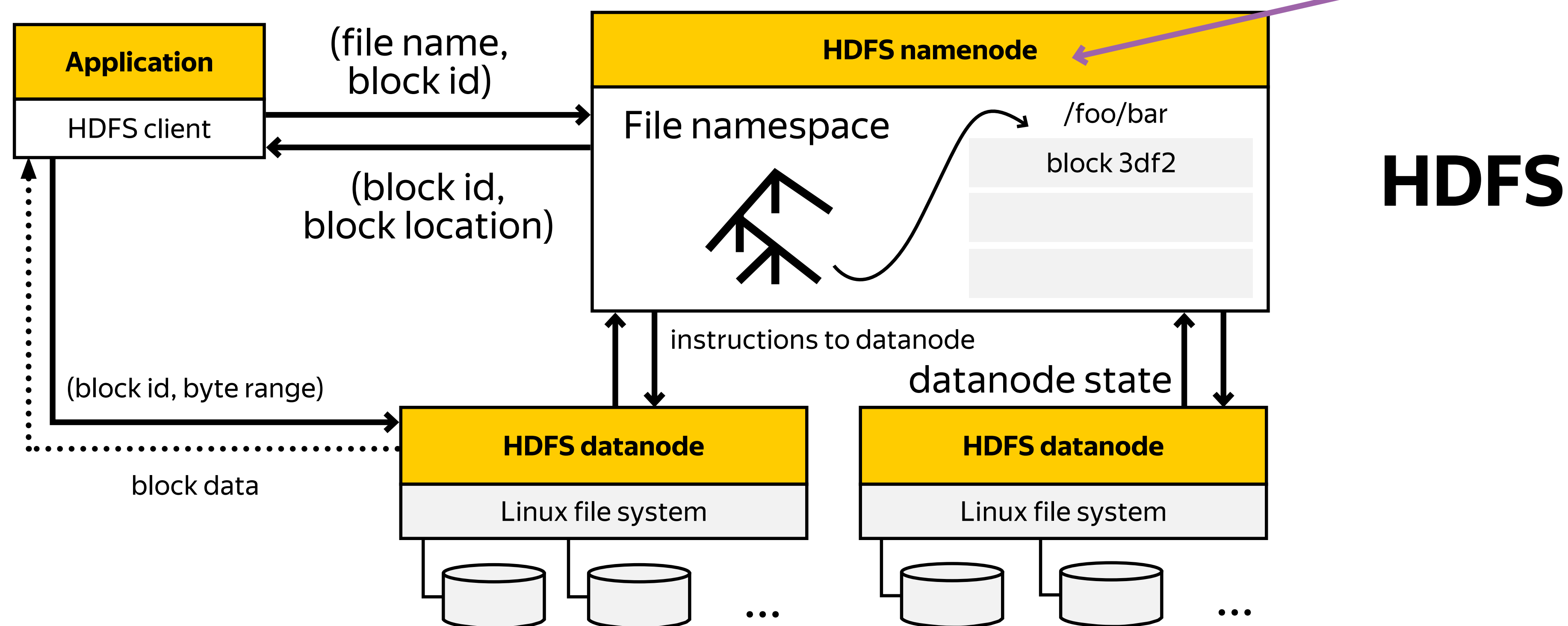
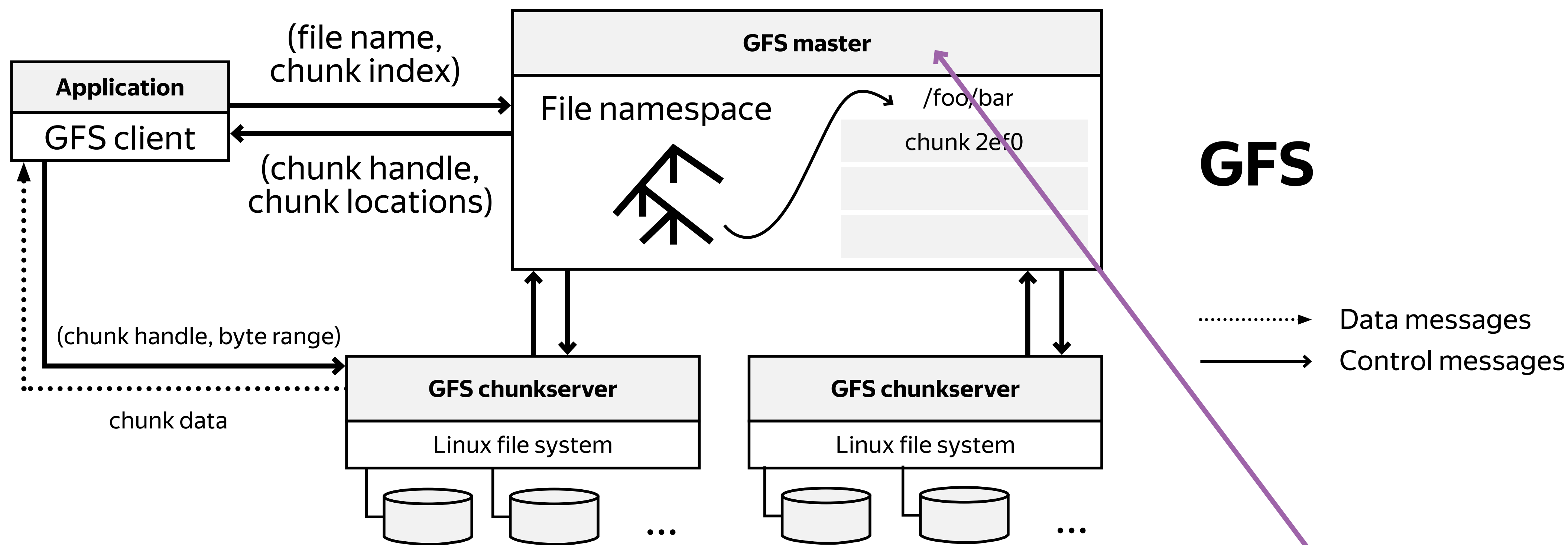


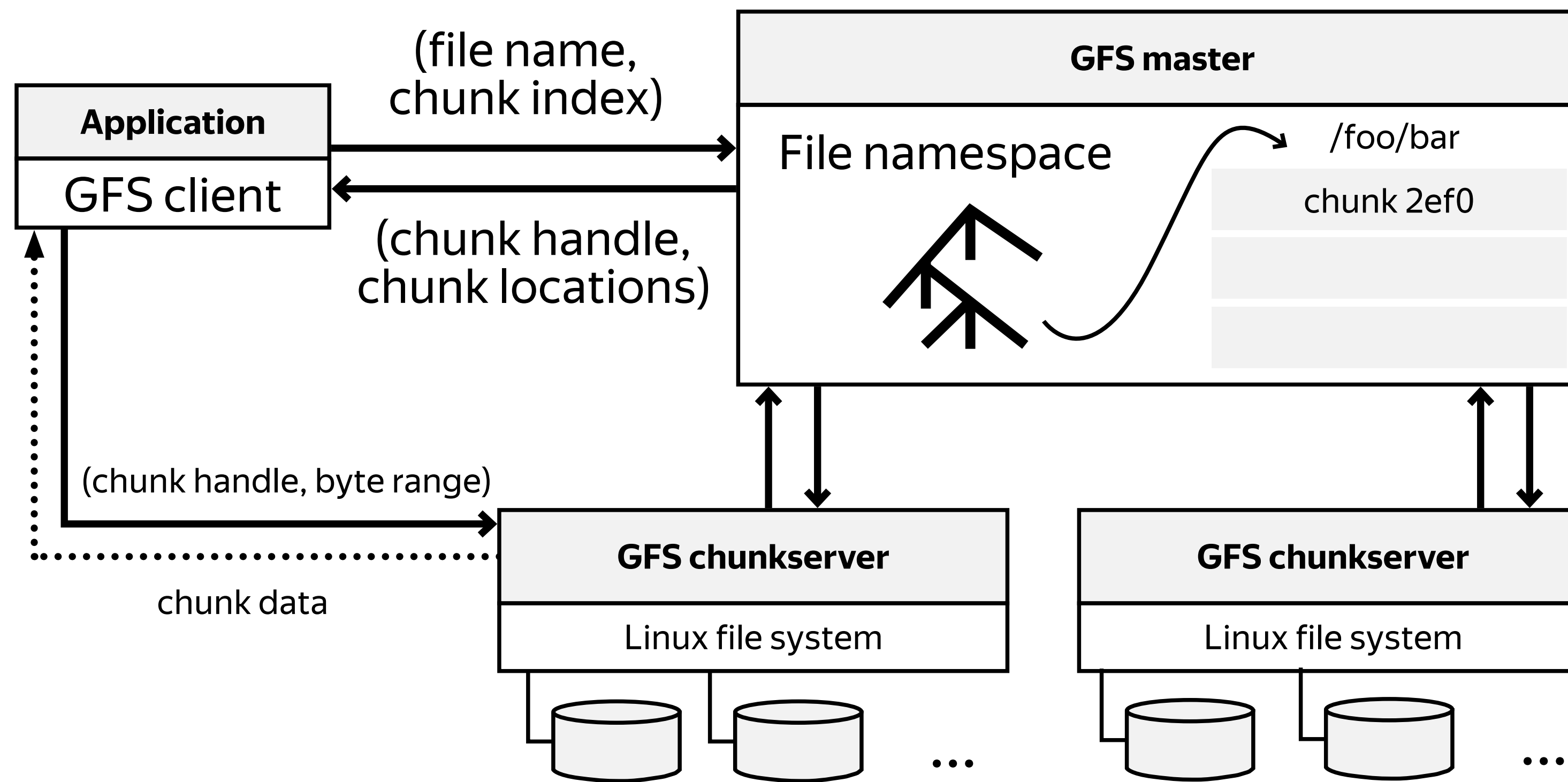
GFS

.....> Data messages
——> Control messages



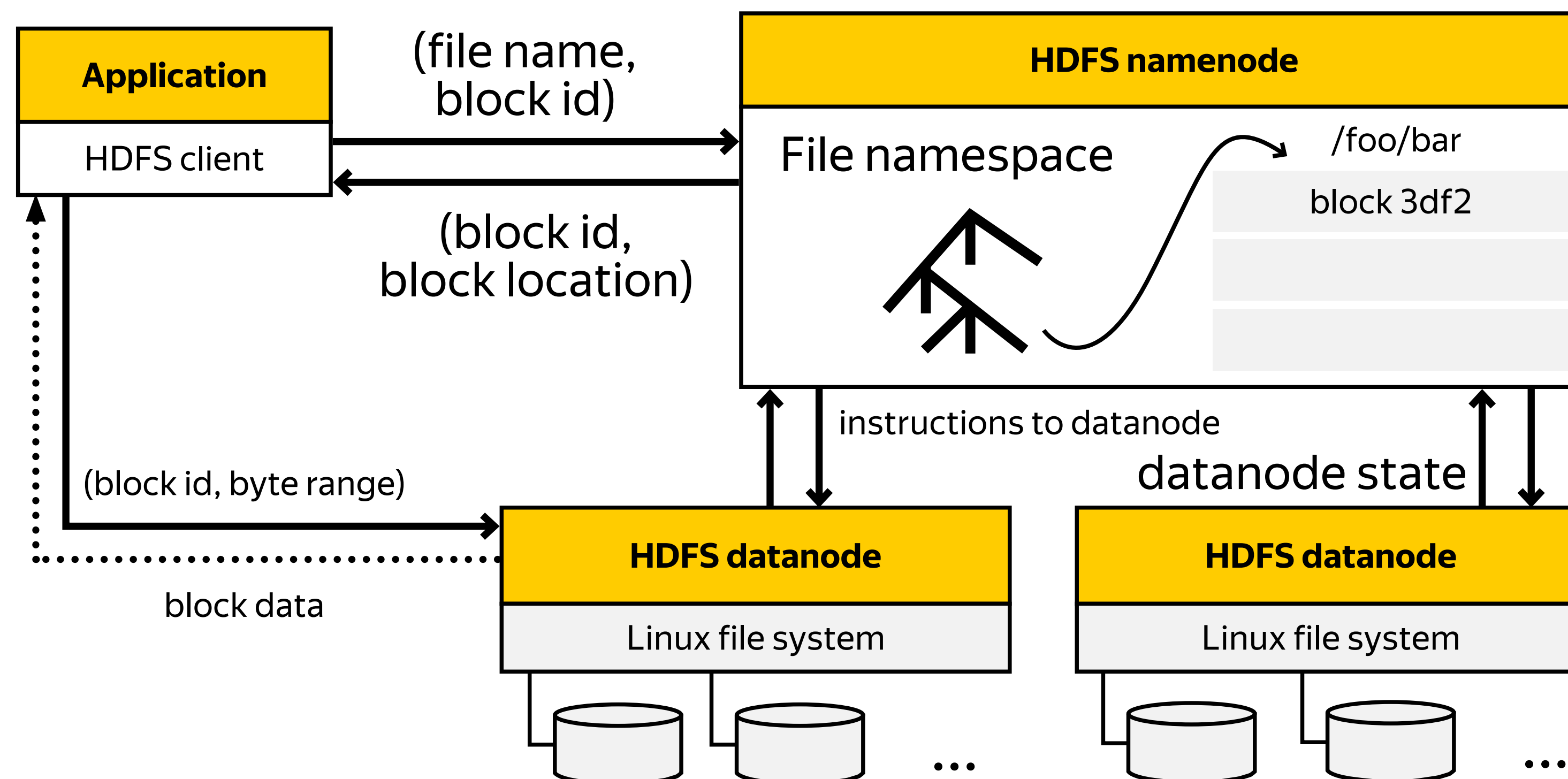
HDFS



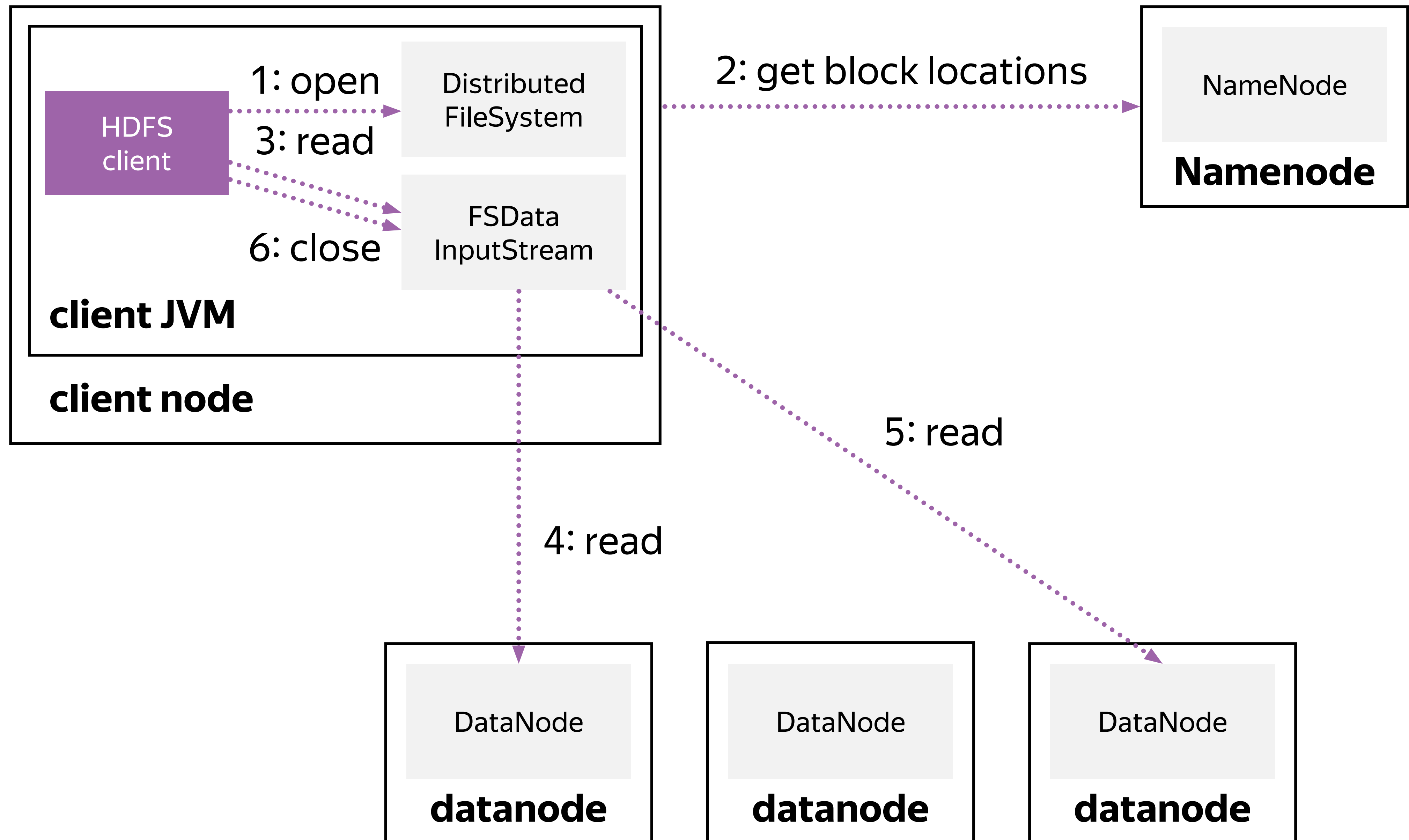


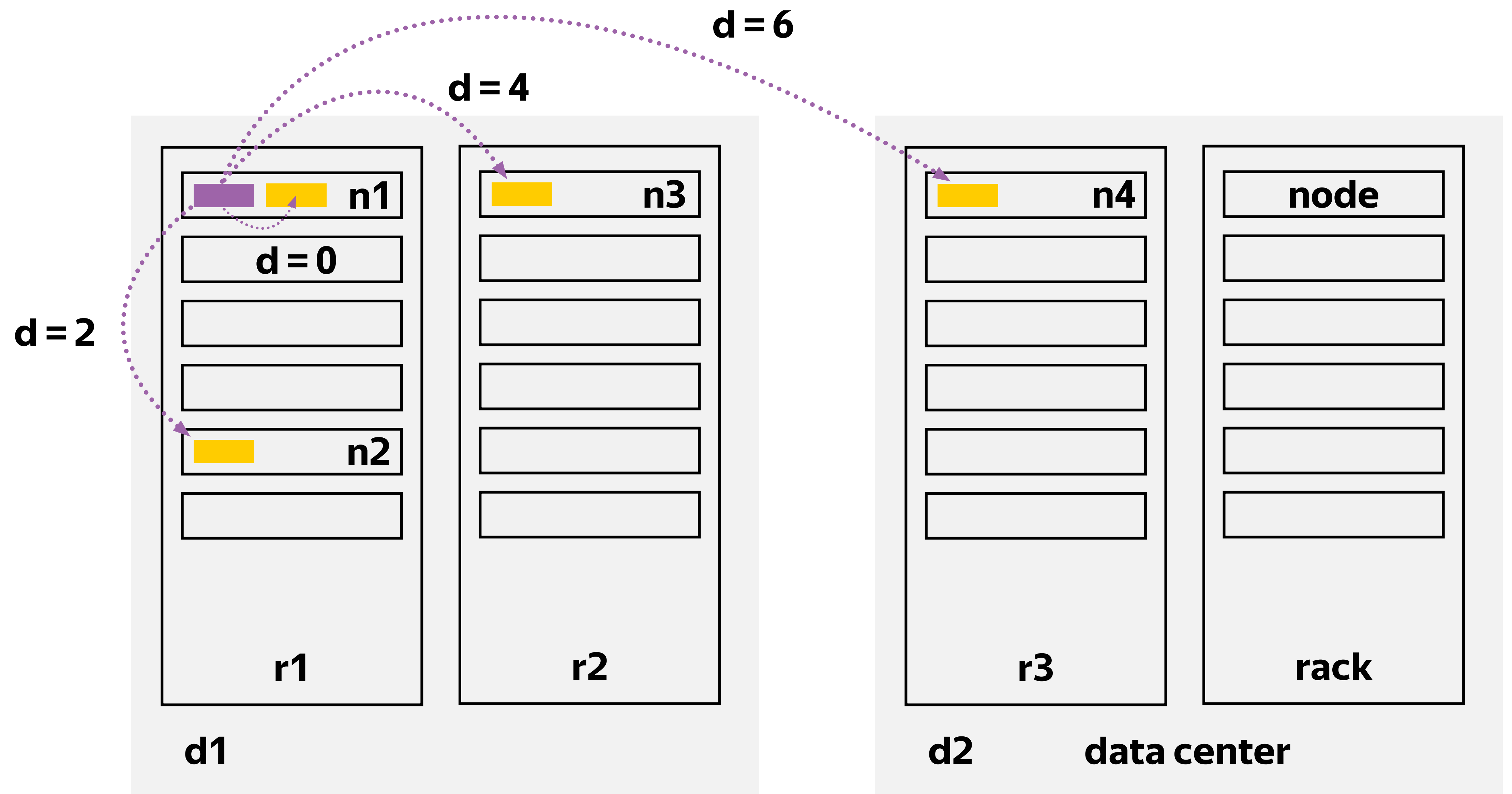
GFS, C++

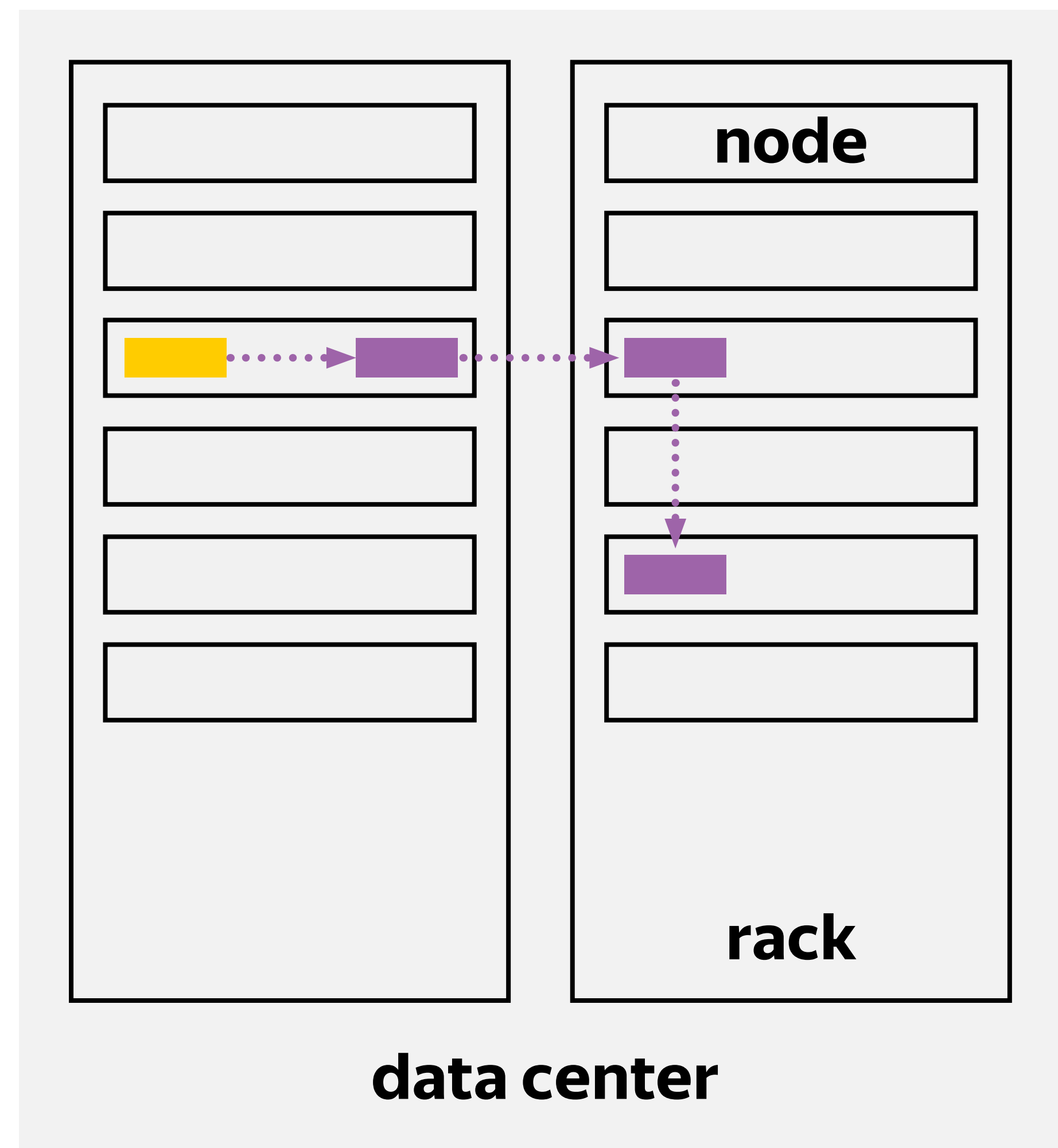
.....➔ Data messages
 ———➔ Control messages

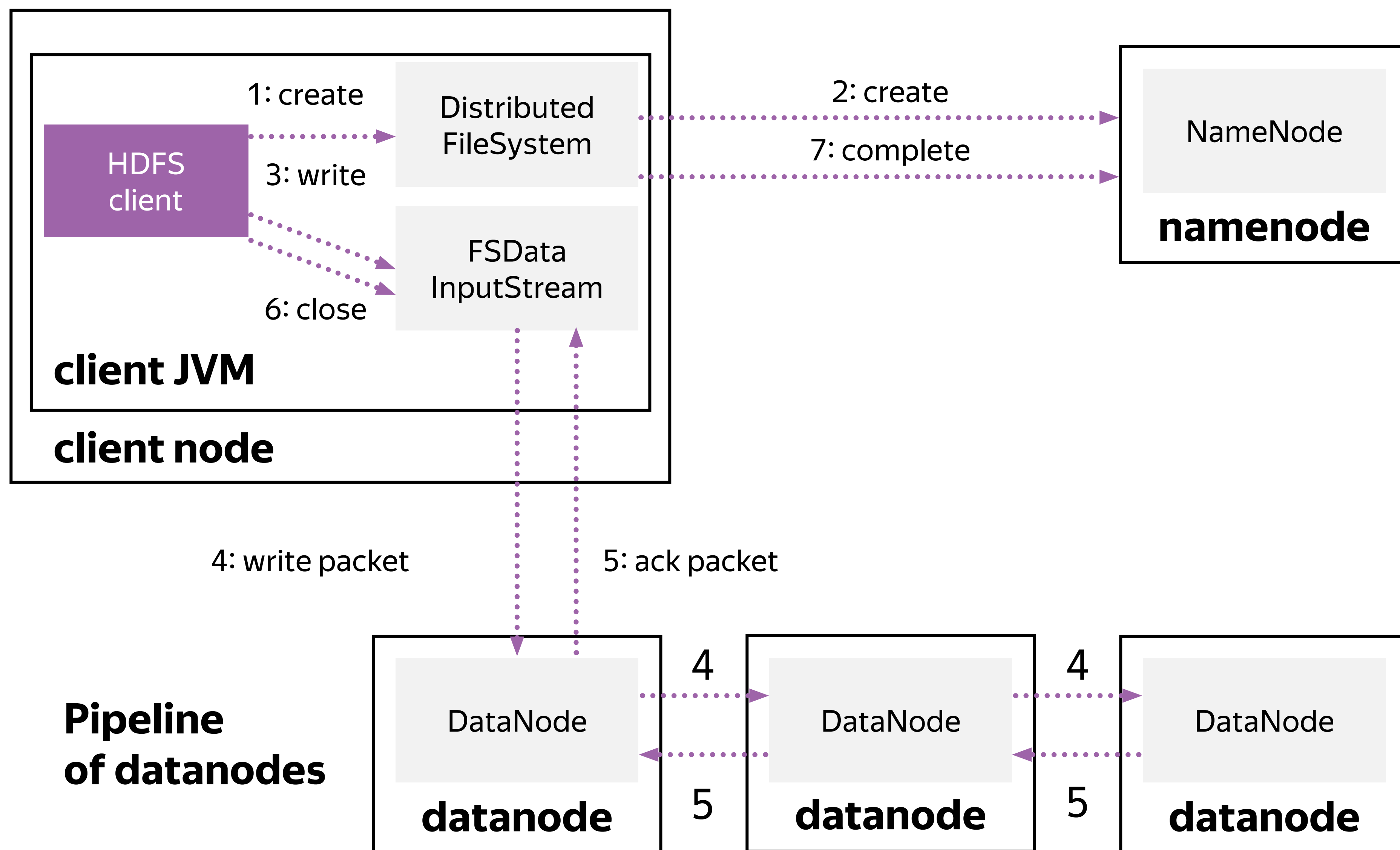


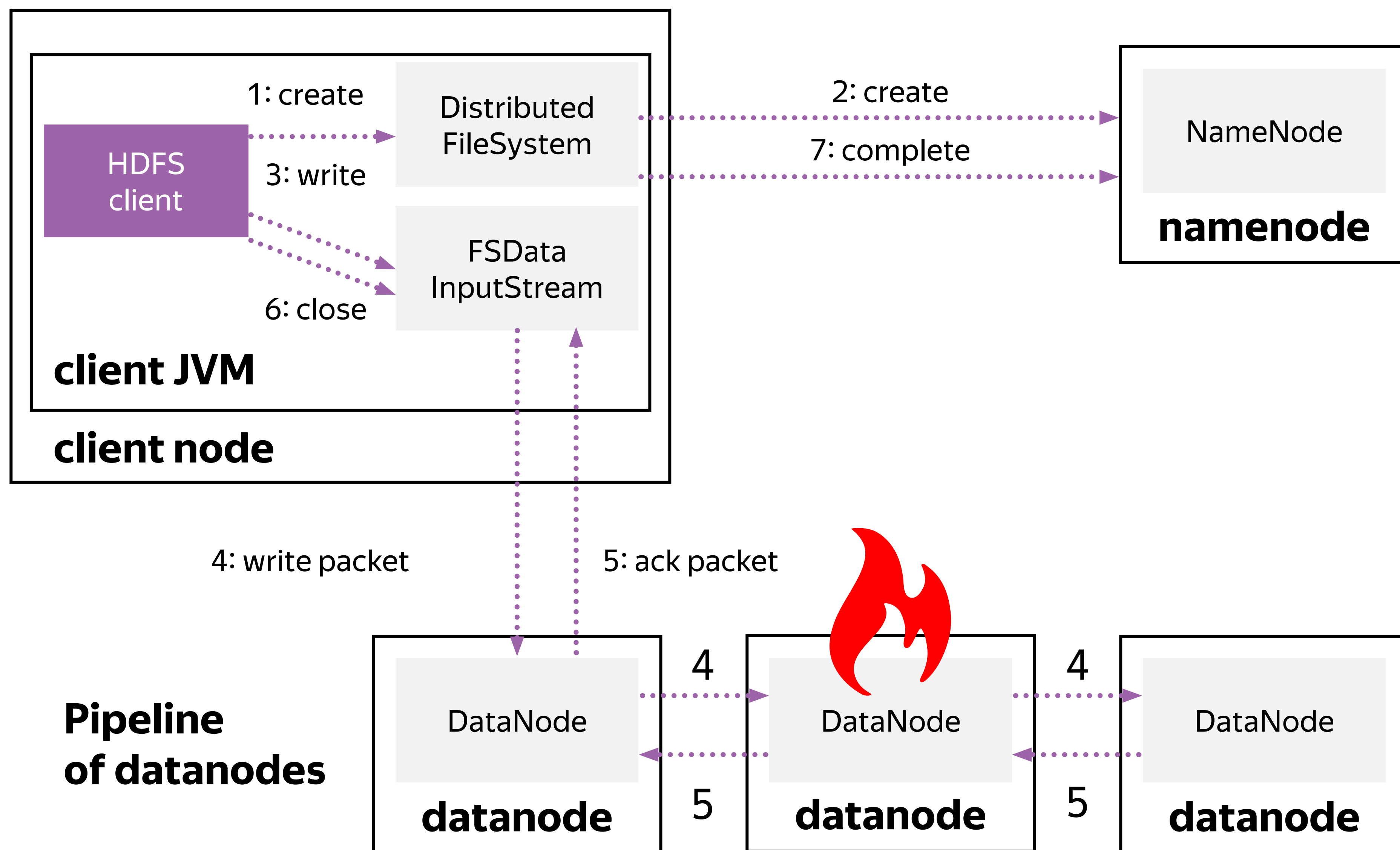
HDFS, Java

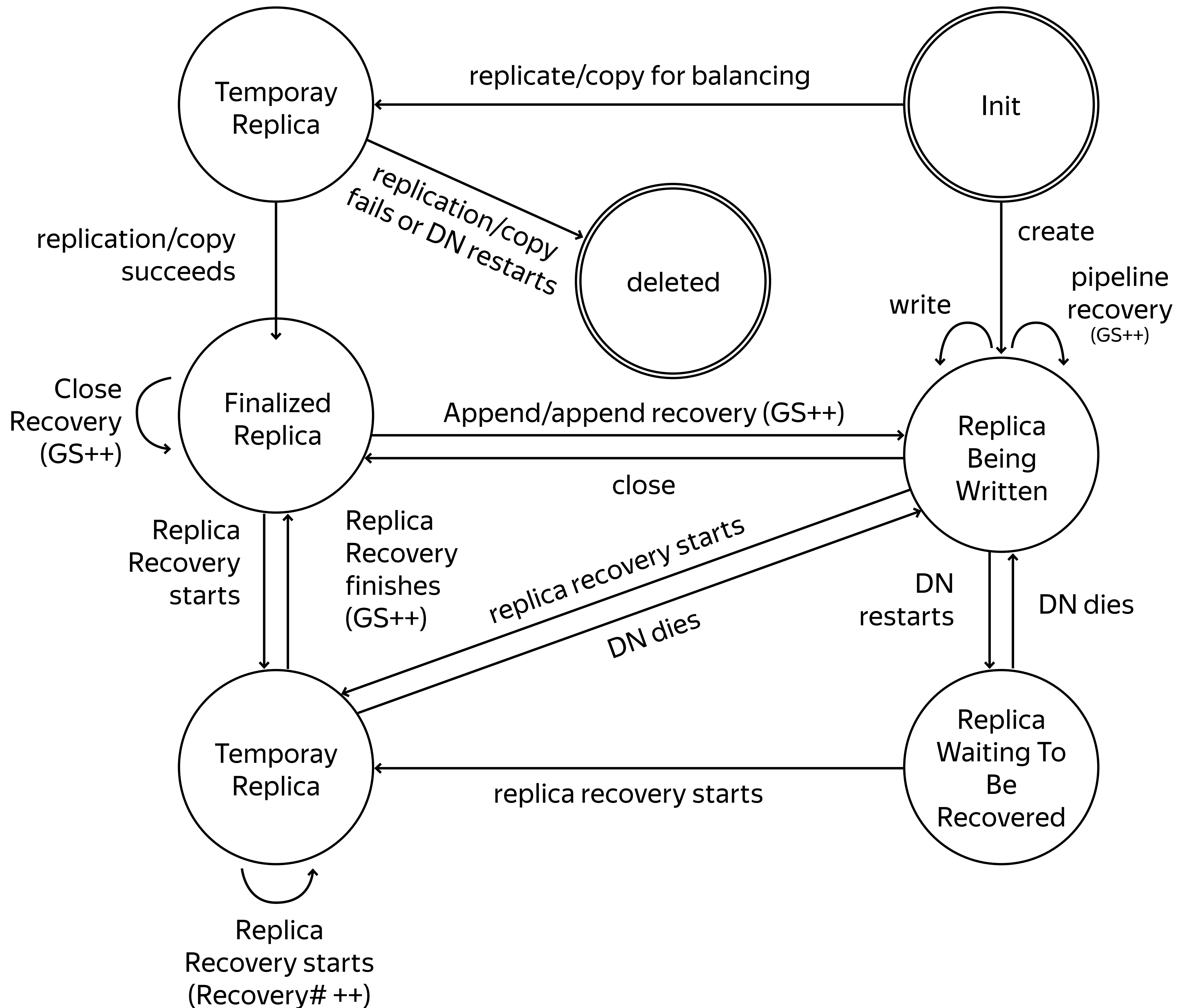












Summary

Summary

- › you can **explain** what vertical and horizontal scaling is

Summary

- › you can **explain** what vertical and horizontal scaling is
- › you can **list** server roles in HDFS

Summary

- › you can **explain** what vertical and horizontal scaling is
- › you can **list** server roles in HDFS
- › you can **explain** how topology affects replica placement

Summary

- › you can **explain** what vertical and horizontal scaling is
- › you can **list** server roles in HDFS
- › you can **explain** how topology affects replica placement
- › you can **explain** what chunk / block size is used for

Summary

- › you can **explain** what vertical and horizontal scaling is
- › you can **list** server roles in HDFS
- › you can **explain** how topology affects replica placement
- › you can **explain** what chunk / block size is used for
- › you can explain in detail how HDFS client reads and writes data

BigDATAteam