



## **Resilient and Fast Persistent Container Storage Leveraging Linux's Storage Functionalities**

Philipp Reisner, CEO LINBIT



# Leading Open Source OS based SDS



## COMPANY OVERVIEW

- Developer of DRBD and LINSTOR
- 100% founder owned
- Offices in Europe and US
- Team of 30+ highly experienced Linux experts
- Exclusivity Japan: SIOS

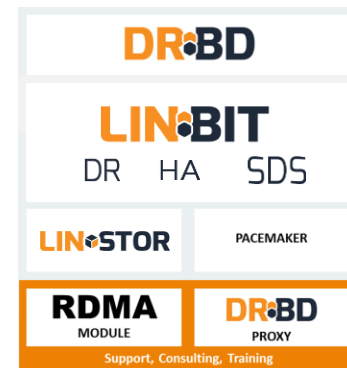


## REFERENCES



## PRODUCT OVERVIEW

- Leading Open Source Block Storage (included in Linux Kernel (v2.6.33))
- Open Source DRBD supported by proprietary LINBIT products / services
- OpenStack with DRBD Cinder driver
- Kubernetes Driver
- 6 x faster than CEPH
- Install base of >2 million



## SOLUTIONS

### DRBD Software Defined Storage (SDS)

New solution (introduced 2016)

Perfectly suited for SSD/NVMe high performance storage

### DRBD High Availability (HA), DRBD Disaster Recovery (DR)

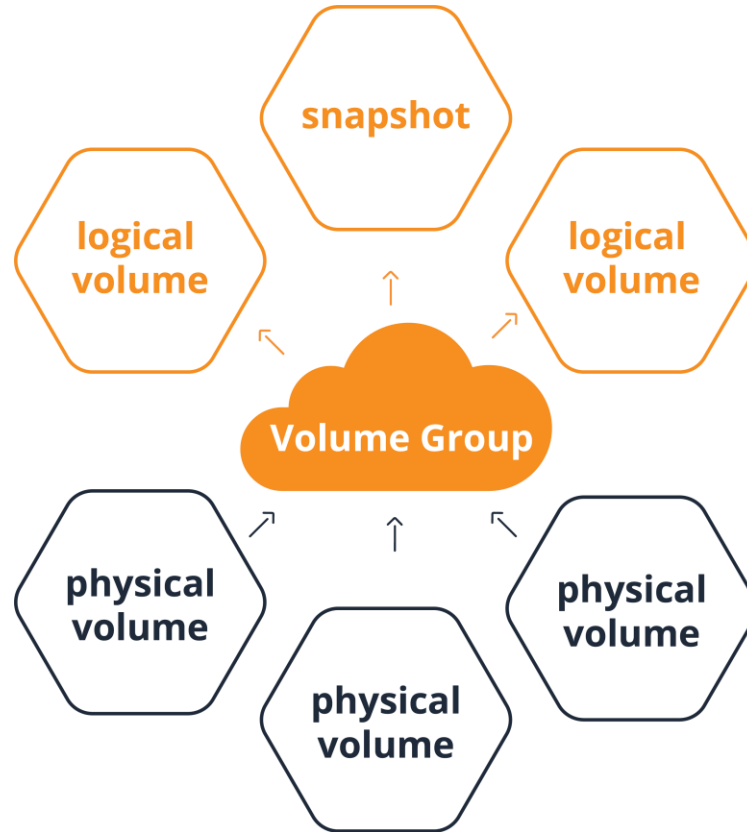
Market leading solutions since 2001, over 600 customers

Ideally suited to power HA and DR in OEM appliances (Cisco, IBM, Oracle)

# Linux Storage Gems

**LVM, RAID, SSD cache tiers, deduplication, targets & initiators**

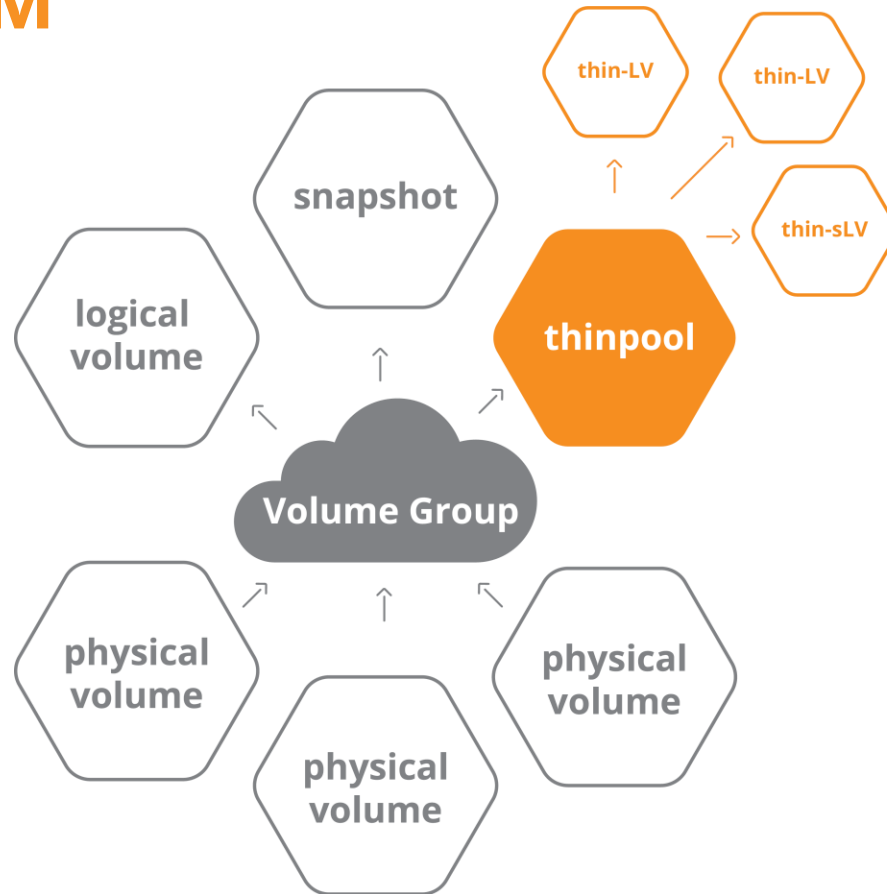
# Linux's LVM



# Linux's LVM

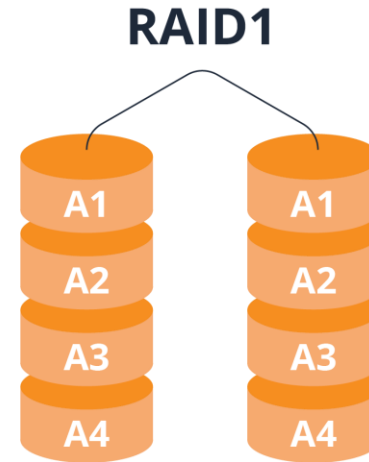
- based on device mapper
- original objects
  - PVs, VGs, LVs, snapshots
  - LVs can scatter over PVs in multiple segments
- thinlv
  - thinpools = LVs
  - thin LVs live in thinpools
  - multiple snapshots became efficient!

# Linux's LVM



# Linux's RAID

- original MD code
  - `mdadm` command
  - Raid Levels: 0,1,4,5,6,10
- Now available in LVM as well
  - device mapper interface for MD code
  - do not call it 'dmraid'; that is software for hardware fake-raid
  - `lvcreate --type raid6 --size 100G VG_name`



# SSD cache for HDD

- dm-cache
  - device mapper module
  - accessible via LVM tools
- bcache
  - generic Linux block device
  - slightly ahead in the performance game



# Linux's DeDupe

- Virtual Data Optimizer (VDO) since RHEL 7.5
  - Red hat acquired Permabit and is GPLing VDO
- Linux upstreaming is in preparation
- in-line data deduplication
- kernel part is a device mapper module
- indexing service runs in user-space
- async or synchronous writeback
- Recommended to be used below LVM

# Linux's targets & initiators

- Open-ISCSI initiator
- letd, STGT, SCST
  - mostly historical
- **LIO**
  - iSCSI, iSER, SRP, FC, FCoE
  - SCSI pass through, block IO, file IO, user-specific-IO
- NVMe-OF
  - target & initiator



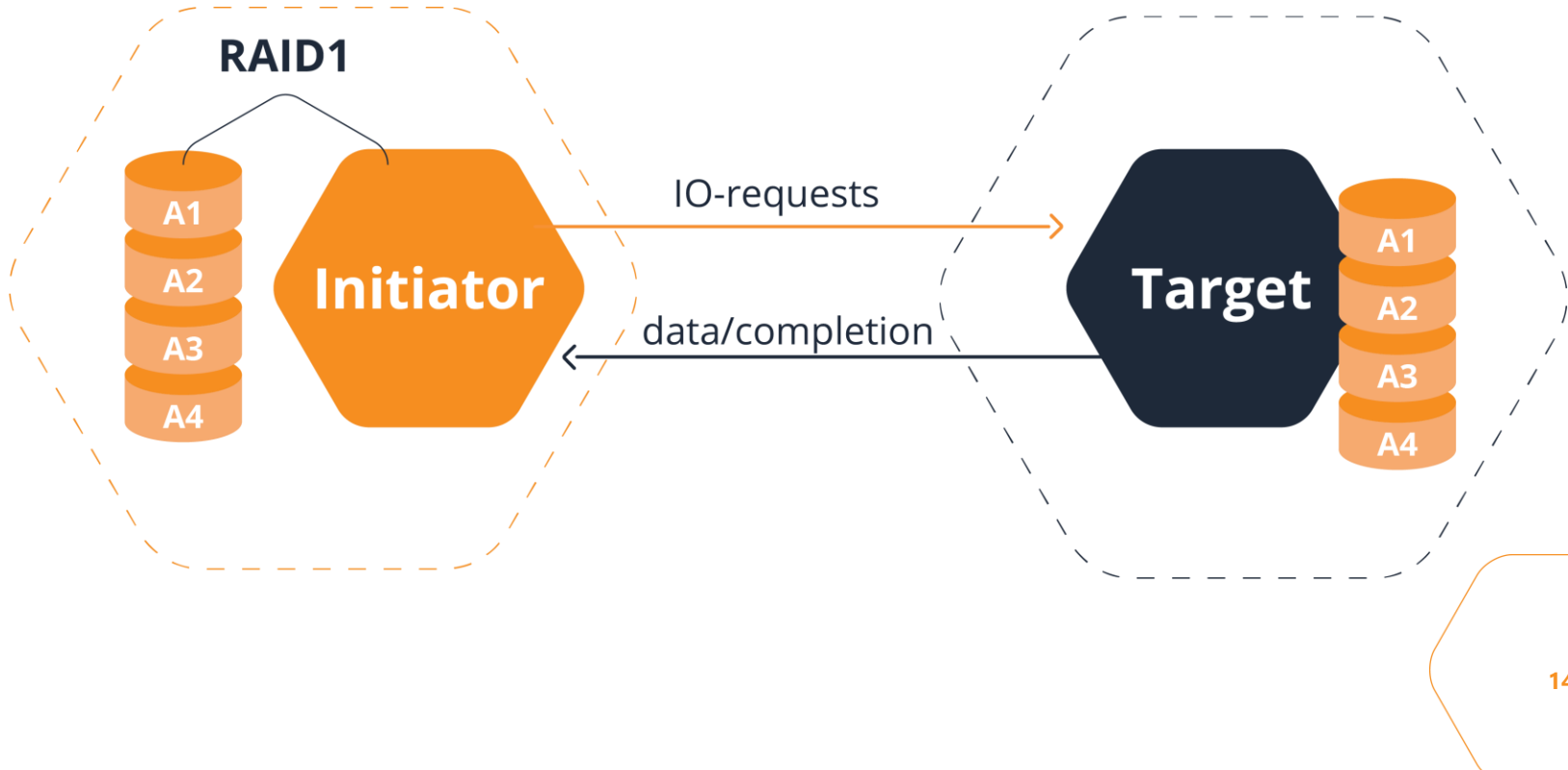
# ZFS on Linux

- Ubuntu eco-system only
- has its own
  - logic volume manager (zVols)
  - thin provisioning
  - RAID (RAIDz)
  - caching for SSDs (ZIL, SLOG)
  - and a file system!

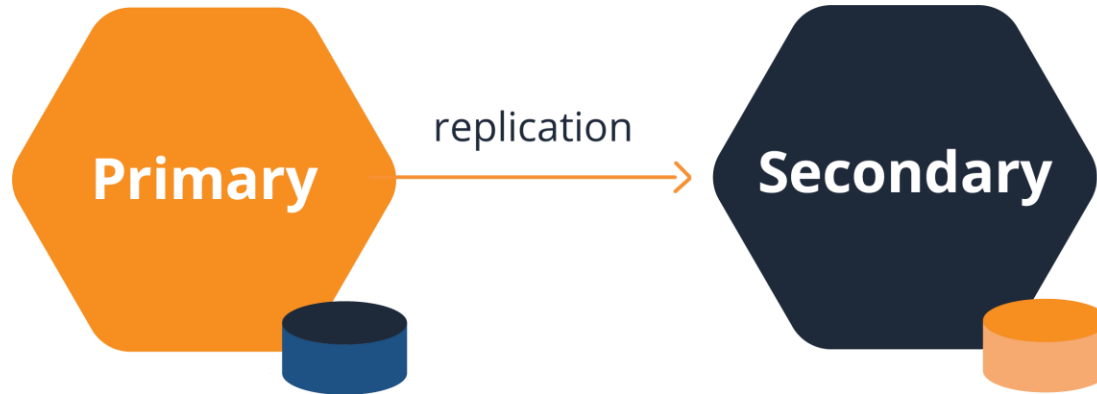


**Put in simplest form**

# DRBD – think of it as ...

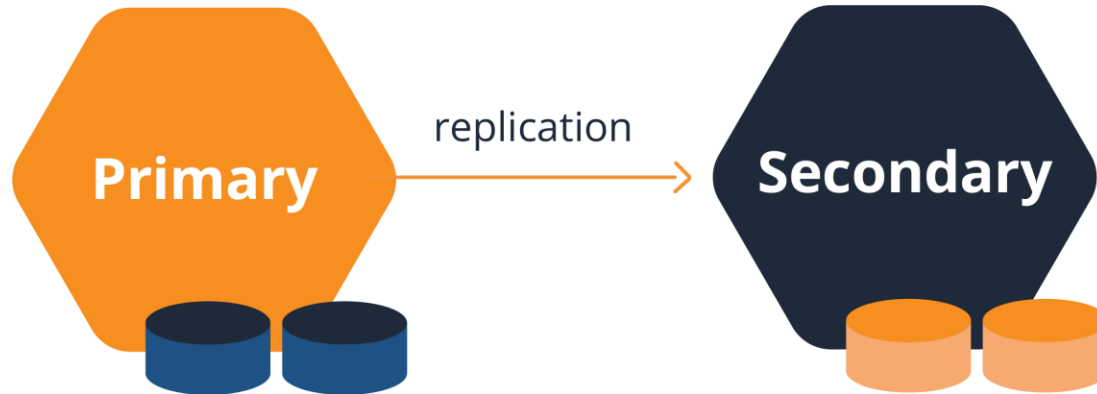


# DRBD Roles: Primary & Secondary



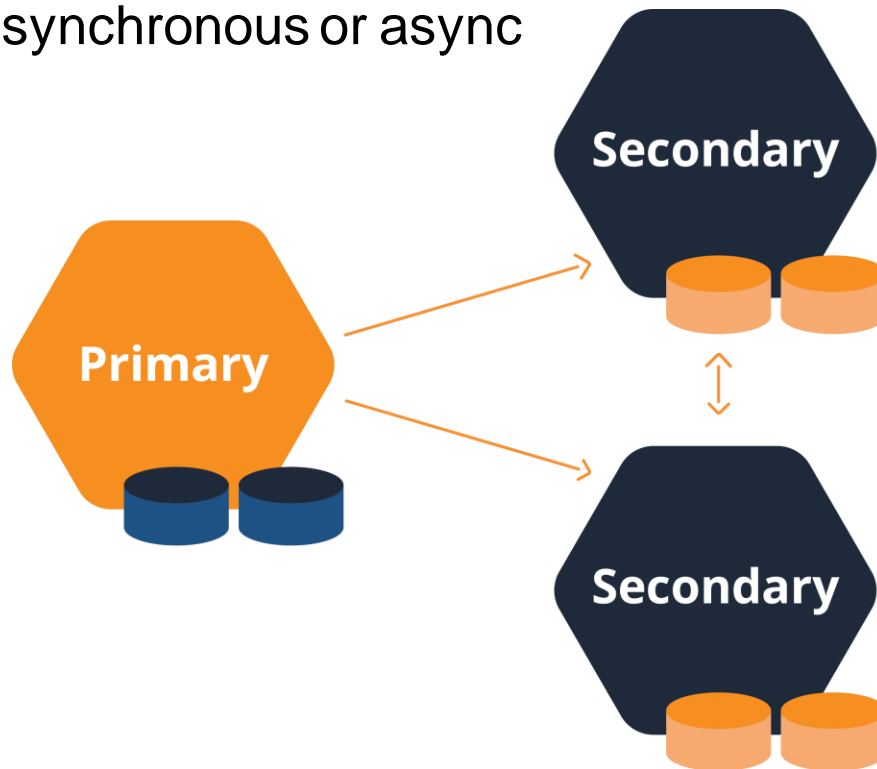
# DRBD – multiple Volumes

- consistency group



# DRBD – up to 32 replicas

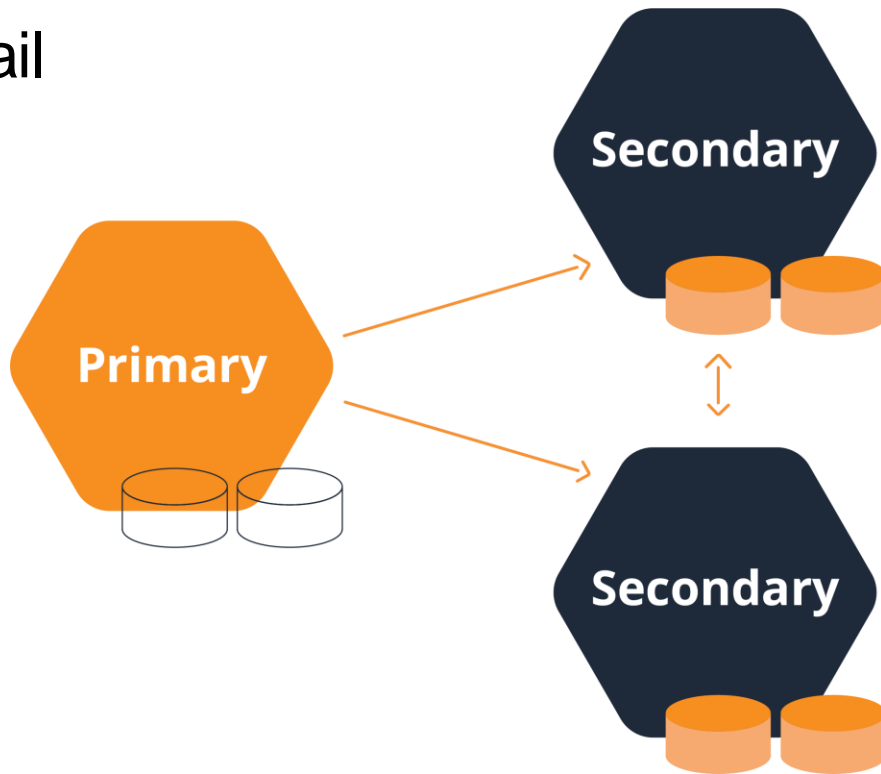
- each may be synchronous or async





# DRBD – Diskless nodes

- intentional diskless (no change tracking bitmap)
- disks can fail



## DRBD - more about

- a node knows the version of the data it exposes
- automatic partial resync after connection outage
- checksum-based verify & resync
- split brain detection & resolution policies
- fencing
- quorum
- multiple resources per node possible (1000s)
- dual Primary for live migration of VMs only!

- Recent optimizations
  - meta-data on PMEM/NVDIMMS
  - Improved, fine-grained locking for parallel workloads
- ROADMAP
  - Eurostars grant: DRBD4Cloud
    - erasure coding (2020)
  - Long distance replication
    - send data once over long distance to multiple replicas



**The combination is more than the sum of its parts**

# LINSTOR - goals



- storage build from generic (x86) nodes
- for SDS consumers (K8s, OpenStack, OpenNebula)
- building on existing Linux storage components
- multiple tenants possible
- deployment architectures
  - distinct storage nodes
  - hyperconverged with hypervisors / container hosts
- LVM, thin LVM or ZFS for volume management (stratis later)
- **Open Source, GPL**

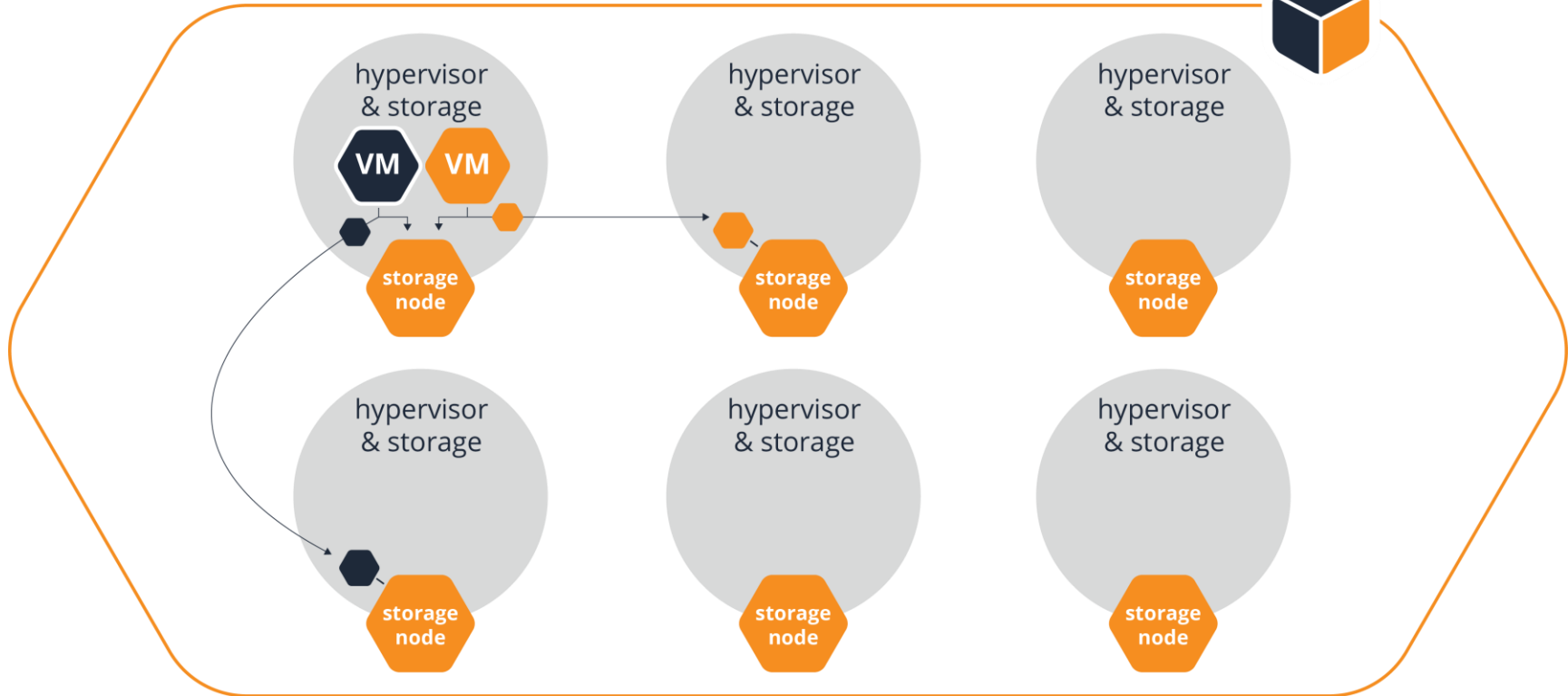


**Examples**

# LINSTOR - Hyperconverged

LINBIT

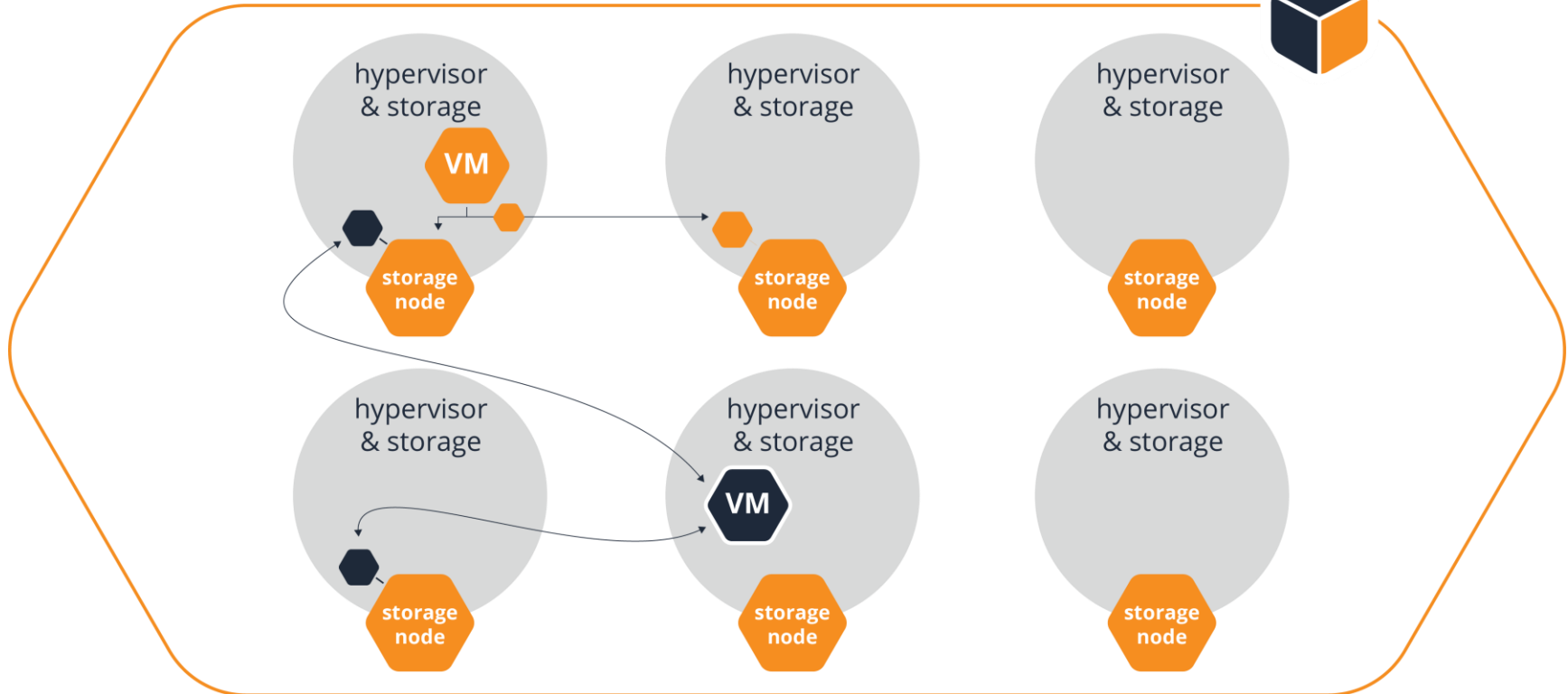
LINSTOR



# LINSTOR - VM migrated

LINBIT

LINSTOR

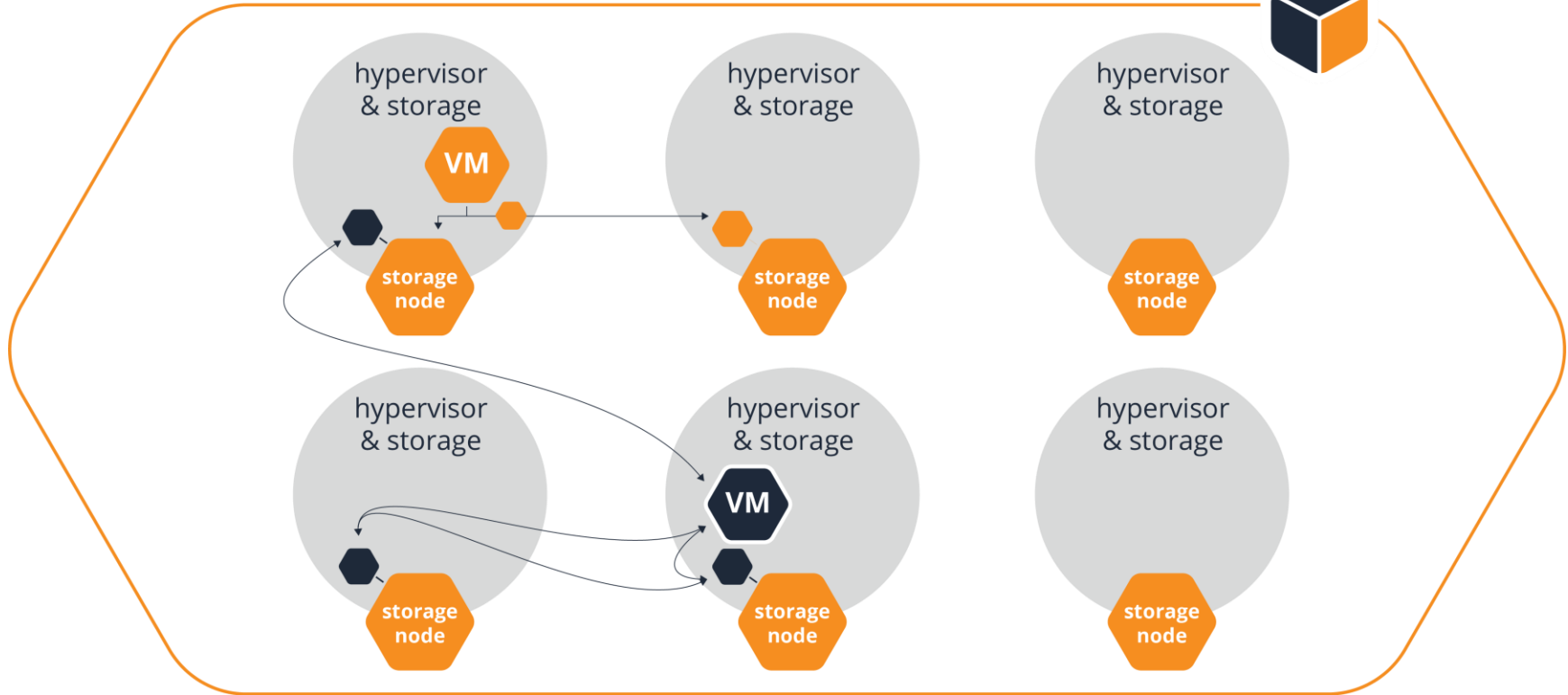




# LINSTOR - add local replica

LINBIT

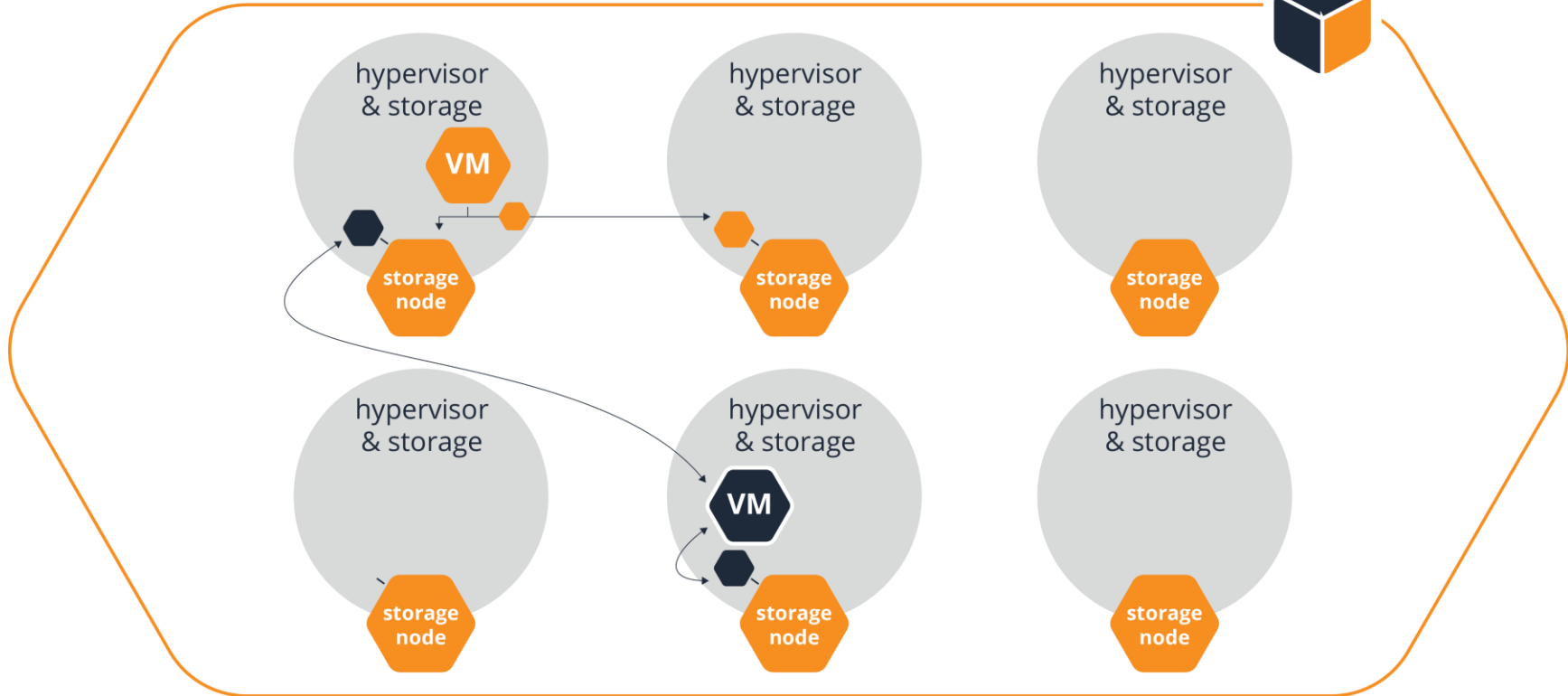
LINSTOR



# LINSTOR - remove 3<sup>rd</sup> copy

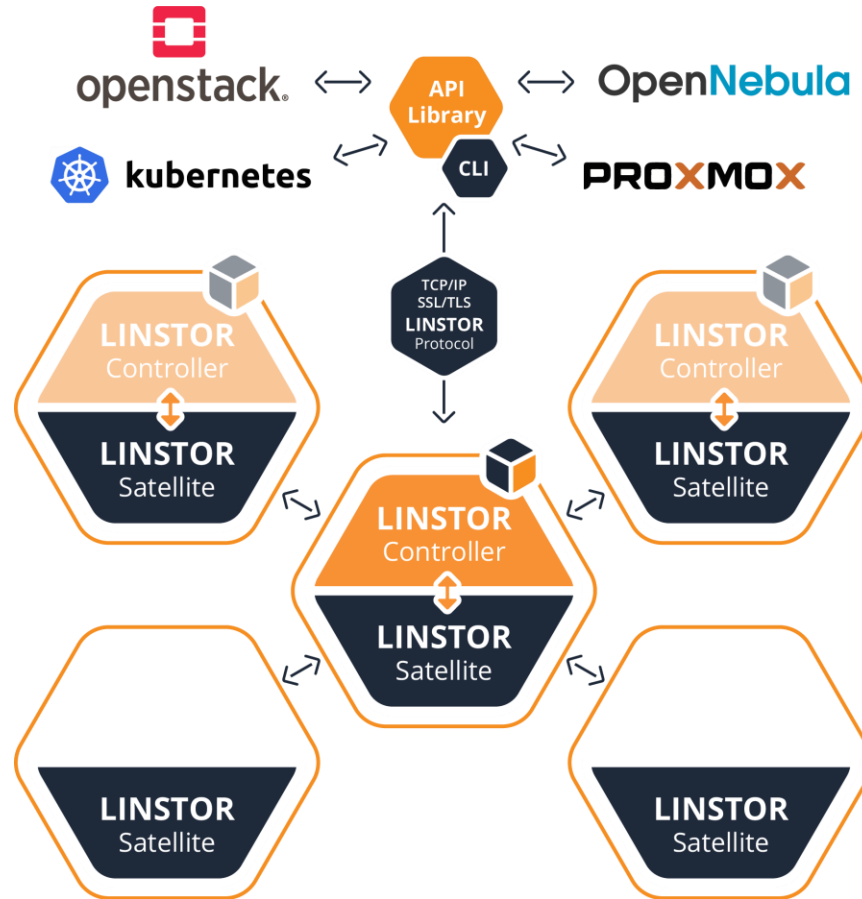
LINBIT

LINSTOR





**Architecture and functions**



# LINSTOR data placement

- arbitrary tags on nodes
  - require placement on equal/different/named tag values
- prohibit placements with named existing volumes
  - different failure domains for related volumes

## Example policy

3 way redundant, where two copies are in the same rack but in different fire compartments (synchronous) and a 3<sup>rd</sup> replica in a different site (asynchronous)

## Example tags

rack = number  
room = number  
site = city

# LINSTOR network path selection

- a storage pool may preferred a NIC
  - express NUMA relation of NVMe devices and NICs
- DRBD's multi pathing supported
  - load balancing with the RDMA transport
  - fail-over only with the TCP transport

# LINSTOR connectors



- Kubernetes
  - FlexVolume & External Provisioner
  - CSI (Docker Swarm, Mesos)



- OpenStack/Cinder
  - since Stein release (April 2019)



- OpenNebula



- Proxmox VE



- XenServer / XCP-ng

# Piraeus Datastore



- Publicly available containers of all components
- Deployment by single YAML-file
- Joint effort of LINBIT & DaoCloud

<https://piraeus.io>

<https://github.com/piraeusdatastore>





# Case study - intel



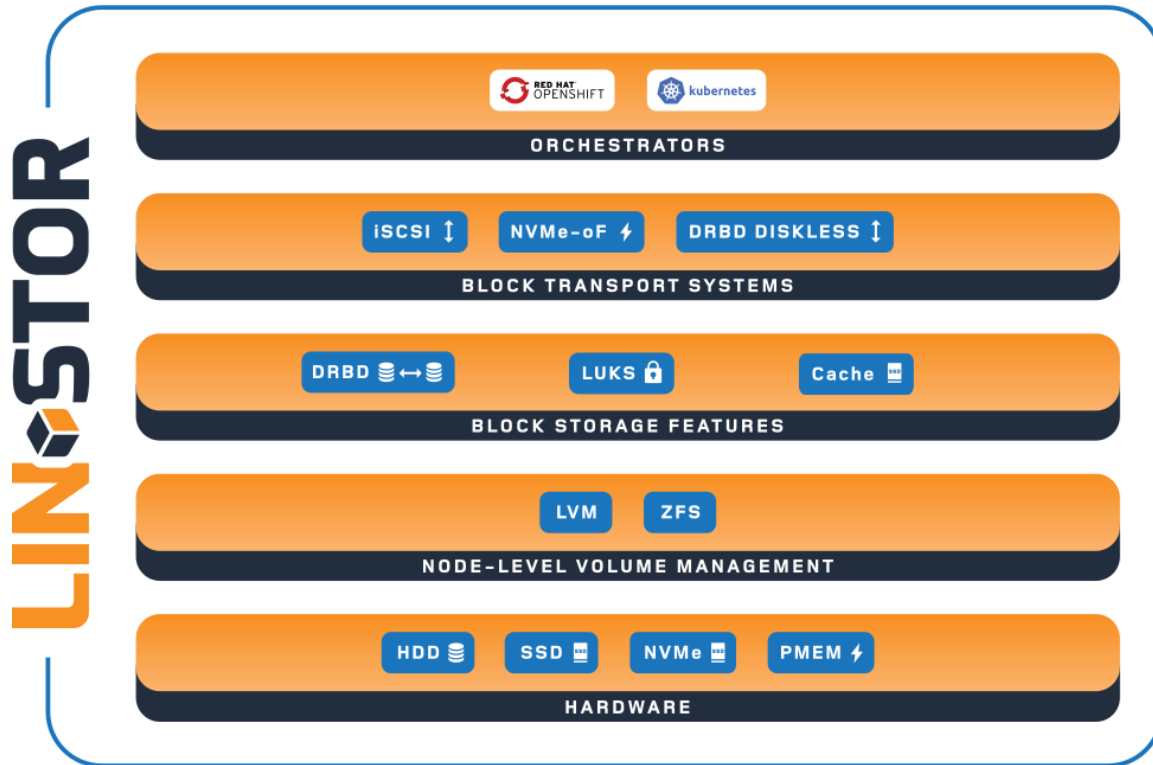
Intel® Rack Scale Design (Intel® **RSD**) is an industry-wide architecture for disaggregated, composable infrastructure that fundamentally changes the way a data center is built, managed, and expanded over time.

## LINBIT working together with Intel

LINSTOR is a storage orchestration technology that brings storage from generic Linux servers and SNIA Swordfish enabled targets to containerized workloads as persistent storage. LINBIT is working with Intel to develop a Data Management Platform that includes a storage backend based on LINBIT's software. LINBIT adds support for the SNIA Swordfish API and NVMe-oF to LINSTOR.

# Summary

## LINUX BLOCK STORAGE MANAGEMENT FOR CONTAINERS





# Thank you

<https://www.linbit.com>

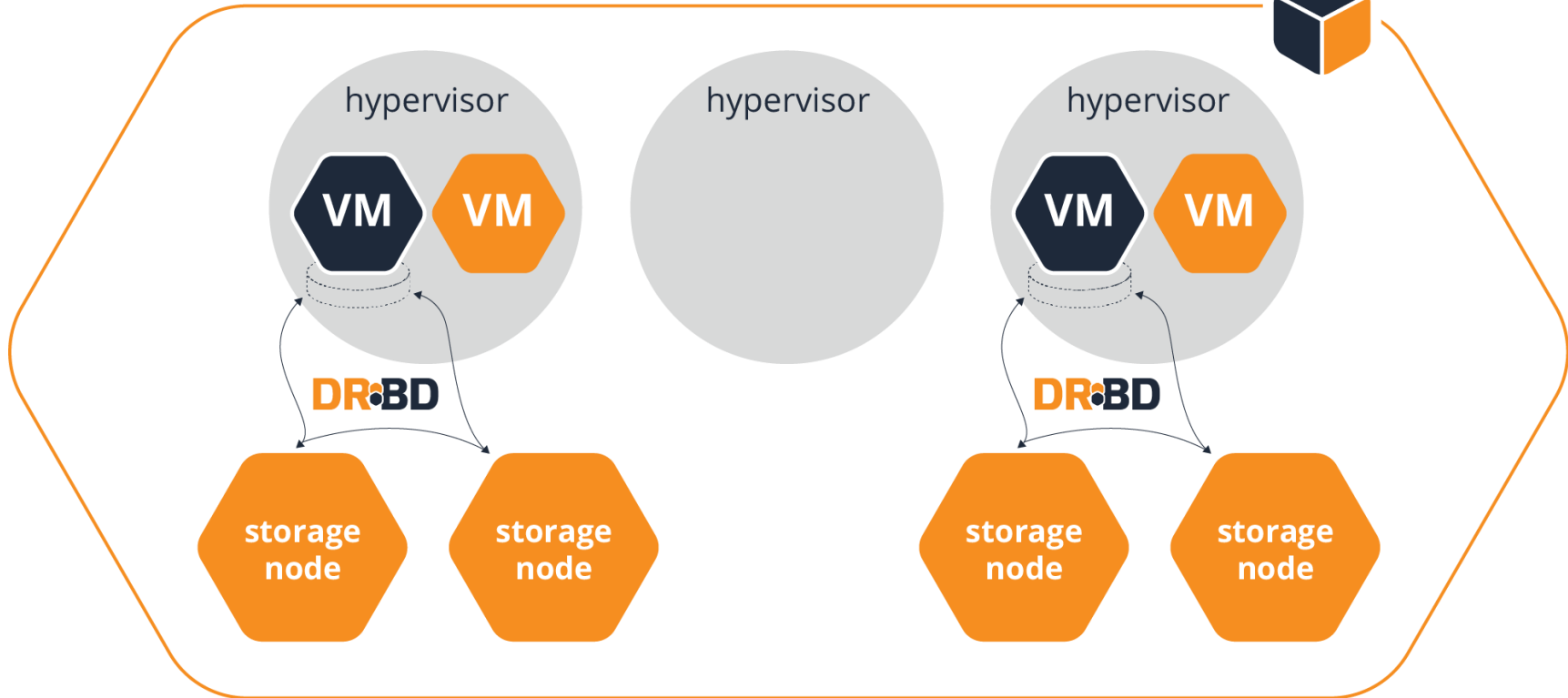


## Appendix Slides: Example Disaggregated Architecture

# LINSTOR – disaggregated stack

LINBIT

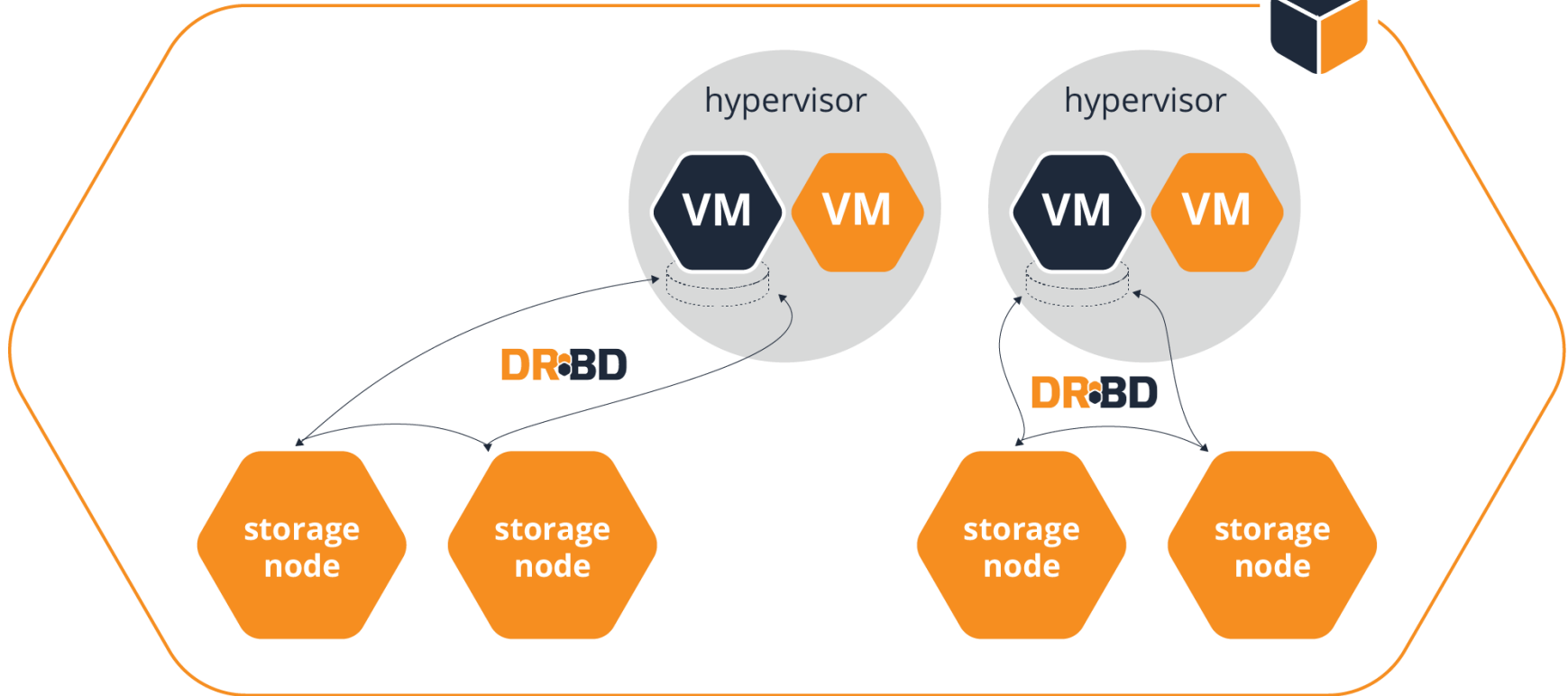
LINSTOR



# LINSTOR / failed Hypervisor

LINBIT

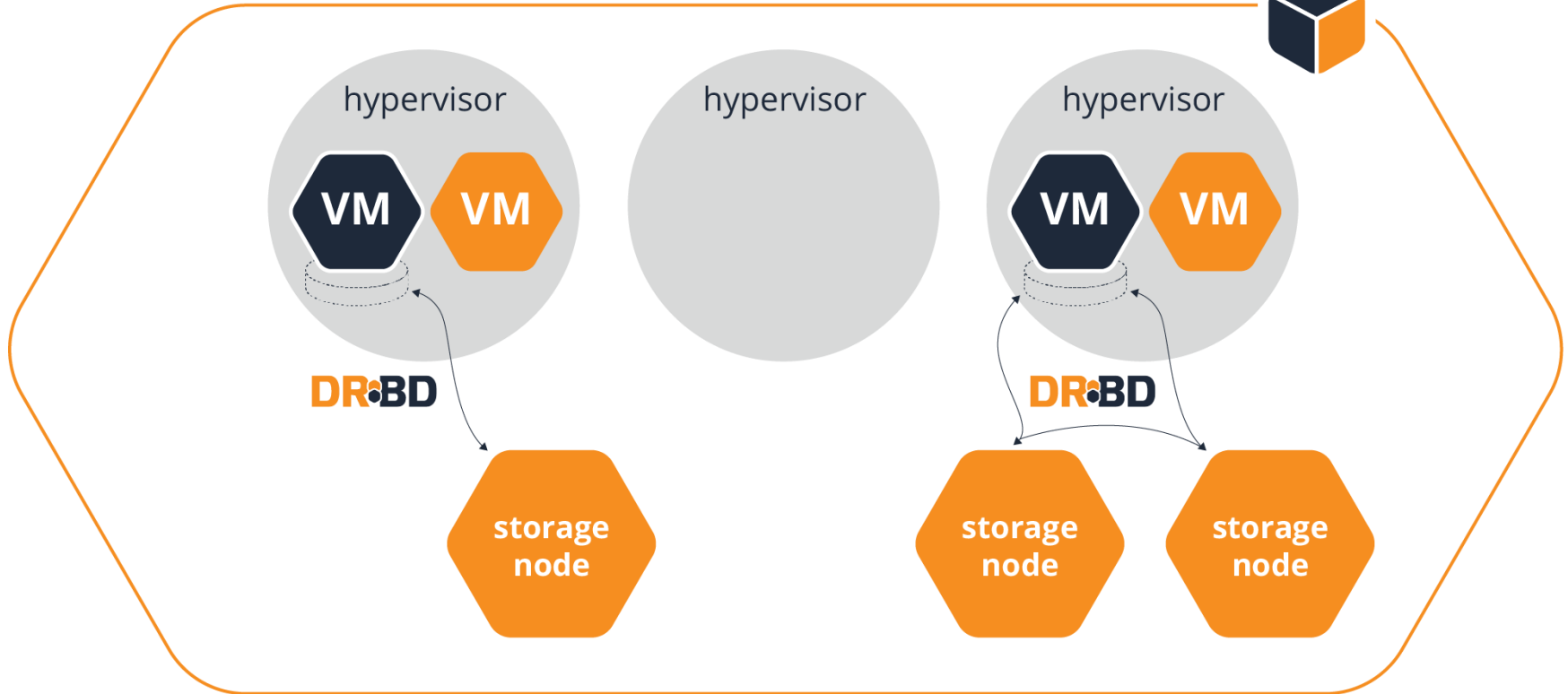
LINSTOR



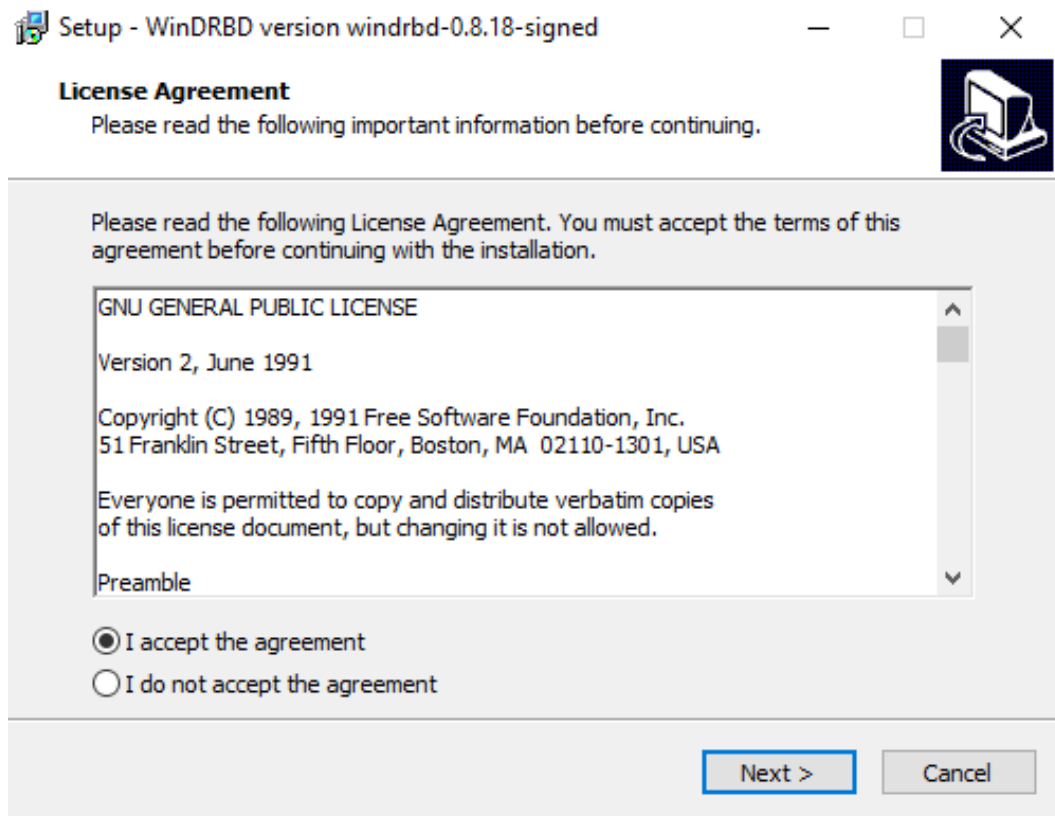
# LINSTOR / failed storage node

LINBIT

LINSTOR



## WinDRBD





- in public beta
  - <https://www.linbit.com/en/drbd-community/drbd-download/>
- Windows 7sp1, Windows 10, Windows Server 2016
- wire protocol compatible to Linux version
- driver tracks Linux version with one day release offset
- WinDRBD user level tools are merged into upstream