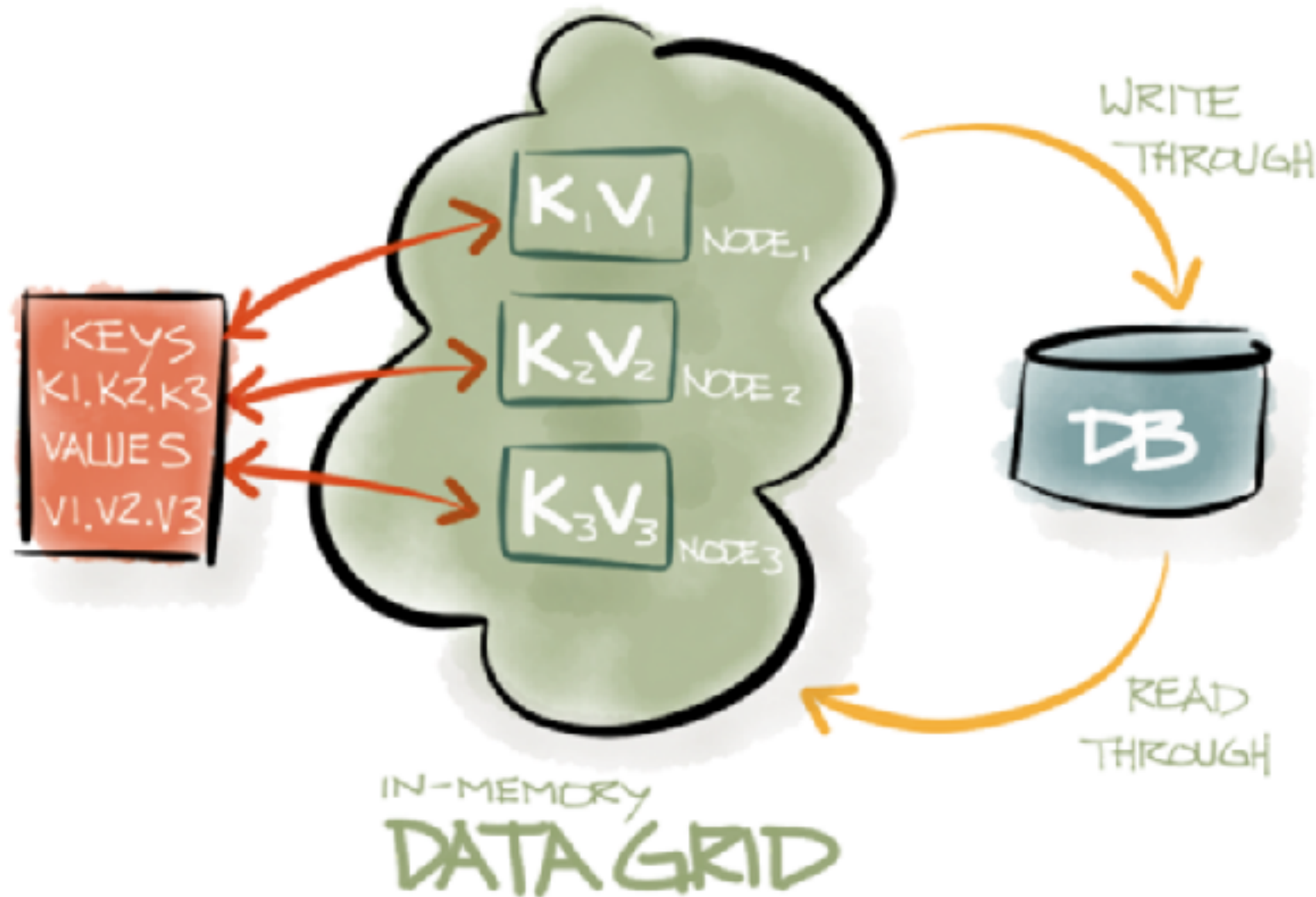# Agenda

- Features that in-memory data grids lack
- Apache Ignite way: durability through page memory architecture
- Durability: use cases and solutions
  - Storage management use cases
  - Data backups use cases
- Durability: performance tricks

GridGain

# Agenda

- **Features that in-memory data grids lack**
- Apache Ignite way: durability through page memory architecture
- Durability: use cases and solutions
  - Storage management use cases
  - Data backups use cases
- Durability: performance tricks

# In-memory Data Grid

Good, but

GridGain

Do you need to access all your data at in-memory speed?

GridGain

# In-memory Data Grid

Good, but

- Storing all data in RAM is expensive
  RAM ~8$ per GB, SSD ~0.2$ per GB

GridGain

Sooner or later, cluster will require maintenance

GridGain

# In-memory Data Grid

Good, but

- Storing all data in RAM is expensive
  RAM ~8$ per GB, SSD ~0.2$ per GB
- Cluster maintenance is complicated
  Grid restart requires data reloading

9

GridGain

Anything that can go wrong will go wrong

# In-memory Data Grid

Good, but

- Storing all data in RAM is expensive
  RAM ~8$ per GB, SSD ~0.2$ per GB
- Cluster maintenance is complicated
  Grid restart requires data reloading
- Disaster protection
  Data backups would be handy

GridGain

# Agenda

- Features that in-memory data grids lack
- **Apache Ignite way: durability through page memory architecture**
- Durability: use cases and solutions
  - Storage management use cases
  - Data backups use cases
- Durability: performance tricks

GridGain
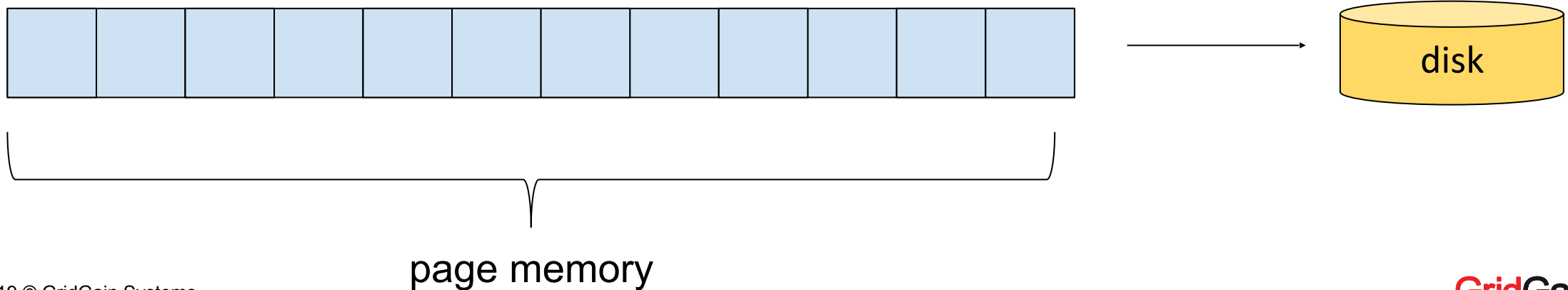
# How to gain in-memory speed and durability?

- Apache Ignite: transparent page memory architecture

GridGain

- Pages are always on disk, optionally in RAM



disk

page memory

GridGain

# Transparent Page Memory Architecture

- Pages are always on disk, optionally in RAM
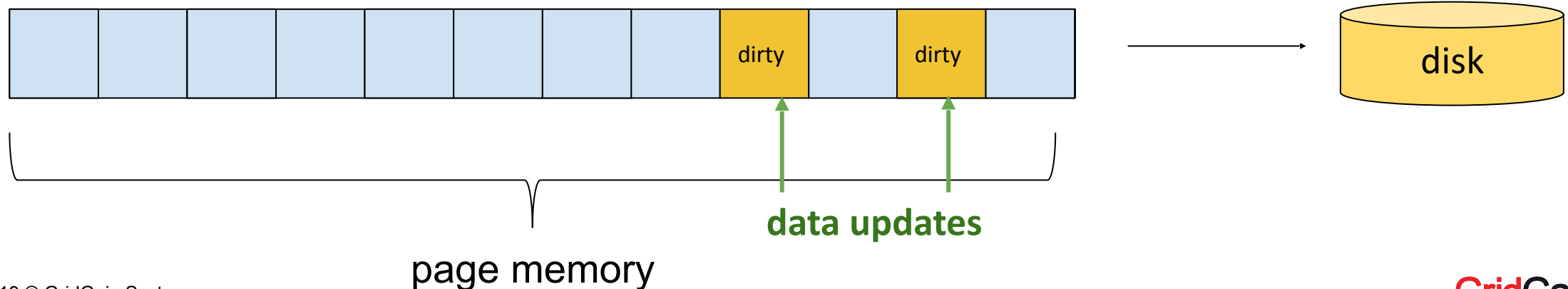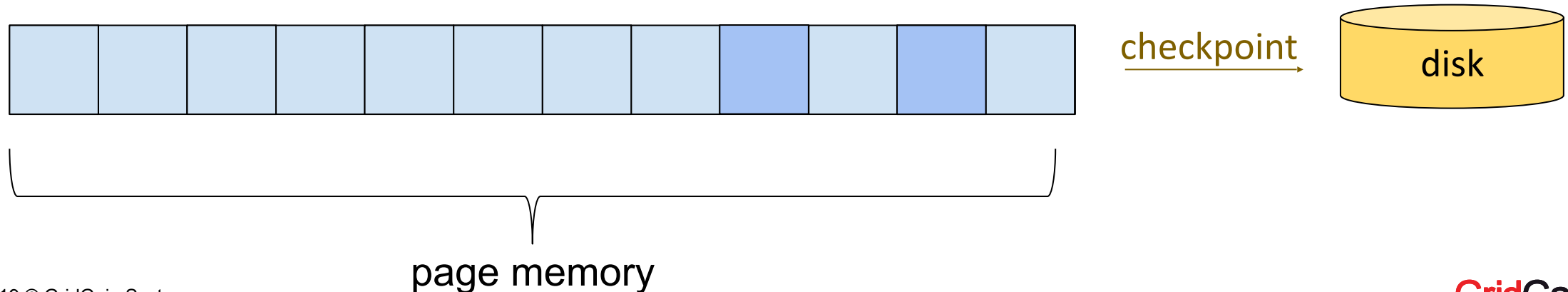- Dirty pages are accumulated in RAM
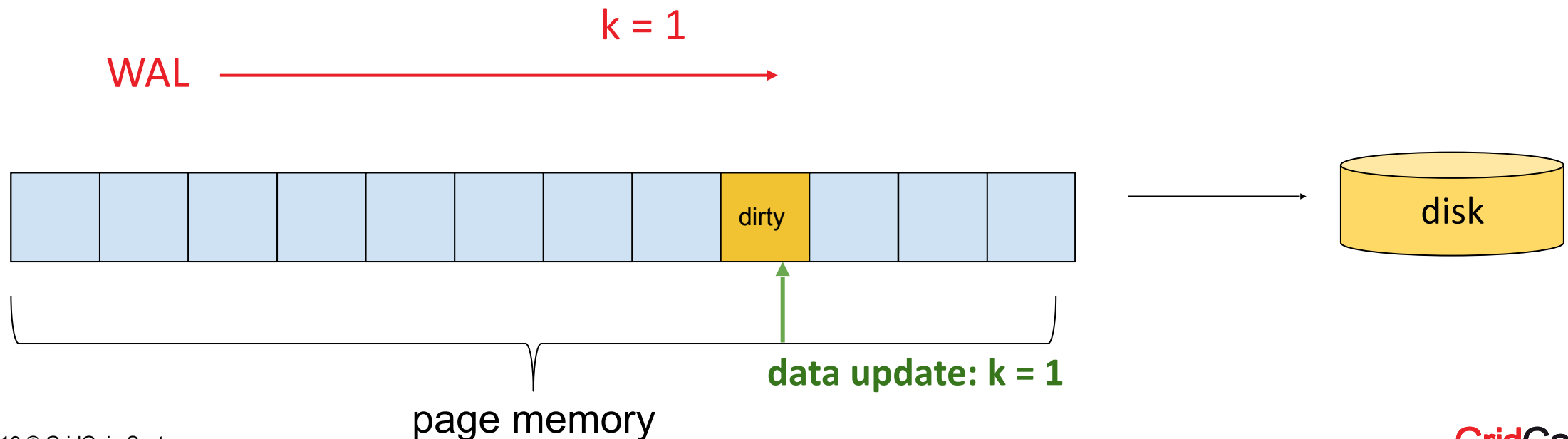
# Transparent Page Memory Architecture

- Pages are always on disk, optionally in RAM
- Dirty pages are accumulated in RAM
- Checkpoint: batch of dirty pages is written to disk



checkpoint

disk

page memory

- Pages are always on disk, optionally in RAM
- Dirty pages are accumulated in RAM
- Checkpoint: batch of dirty pages is written to disk
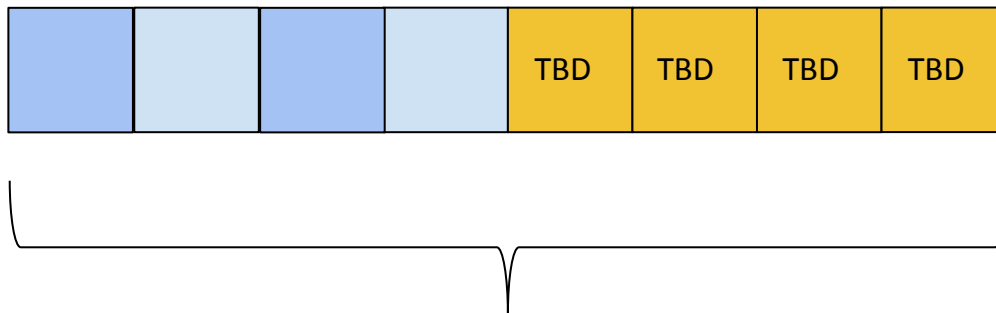- WAL: updates between checkpoints are logged

$k = 1$

WAL

dirty

disk

data update: $k = 1$

page memory

GridGain
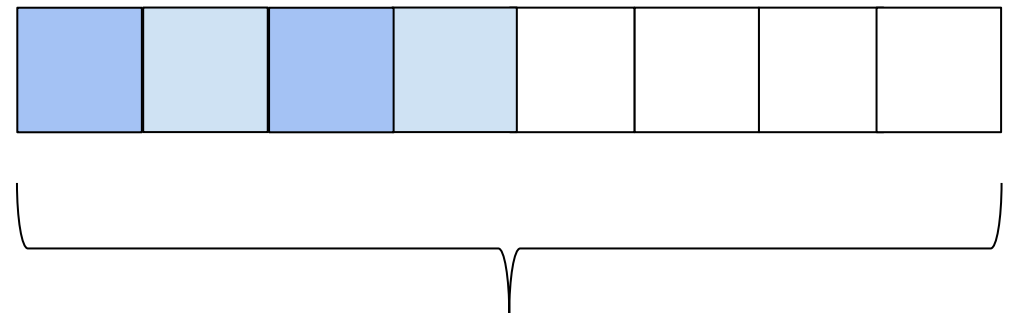
- Scan disk storage, copy pages to snapshot



snapshot

disk storage

snapshot storage

- Scan disk storage, copy pages to snapshot
- Next checkpoint is going to update yet not written page?

snapshot →

| | | | | TBD | TBD | TBD | TBD |

disk storage

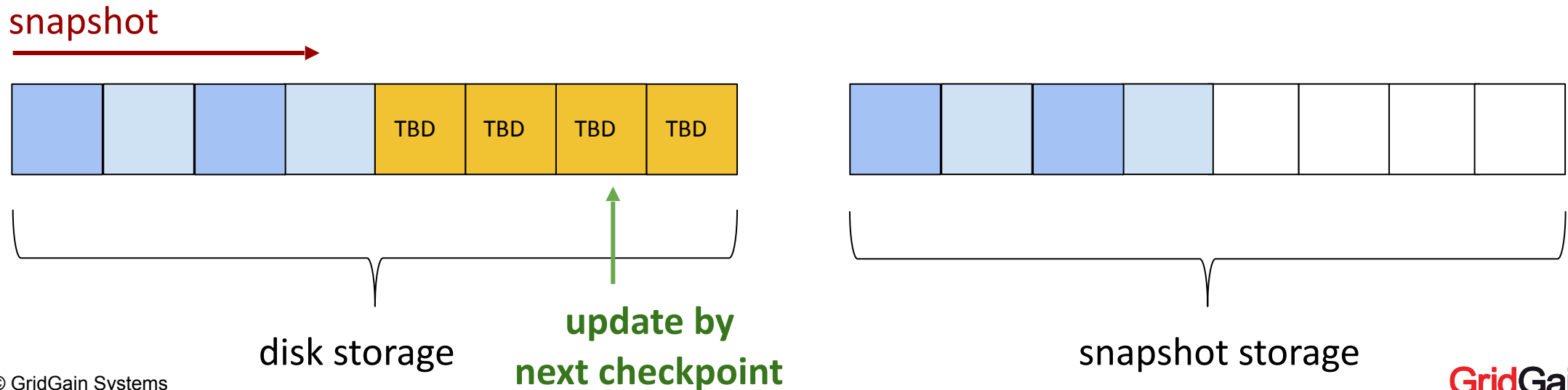update by
next checkpoint

| | | | | | | | |

snapshot storage

GridGain

# Snapshot under load: copy-on-write

- Scan disk storage, copy pages to snapshot
- Next checkpoint is going to update yet not written page?
- Let it write page to snapshot first!



cooperative write to snapshot

snapshot

| | | | | TBD | TBD | new state | TBD |

disk storage

update by next checkpoint

| | | | | | old state | |

snapshot storage

GridGain

# Agenda

- Features that in-memory data grids lack
- Apache Ignite way: durability through page memory architecture
- Durability: use cases and solutions
  - Storage management use cases
  - Data backups use cases
- Durability: performance tricks

GridGain

# From theory to practice: Data Storage Configuration

- Use cases:
  - Limit RAM usage

# From theory to practice: Data Storage Configuration

- Use cases:
  - Limit RAM usage
  - Different RAM limitations for different caches

- Use cases:
  - Limit RAM usage
  - Different RAM limitations for different caches
  - Fast cluster restart and cheaper data storing

# From theory to practice: Data Storage Configuration

- Use cases:
  - Limit RAM usage
  - Different RAM limitations for different caches
  - Fast cluster restart and cheaper data storing
  - Hot and cold data

- Default: in-memory mode

# Use case: limit node RAM consumption

- Default: in-memory mode
- Overall RAM usage limit is configurable

GridGain

- Default: in-memory mode
- Overall RAM usage limit is configurable
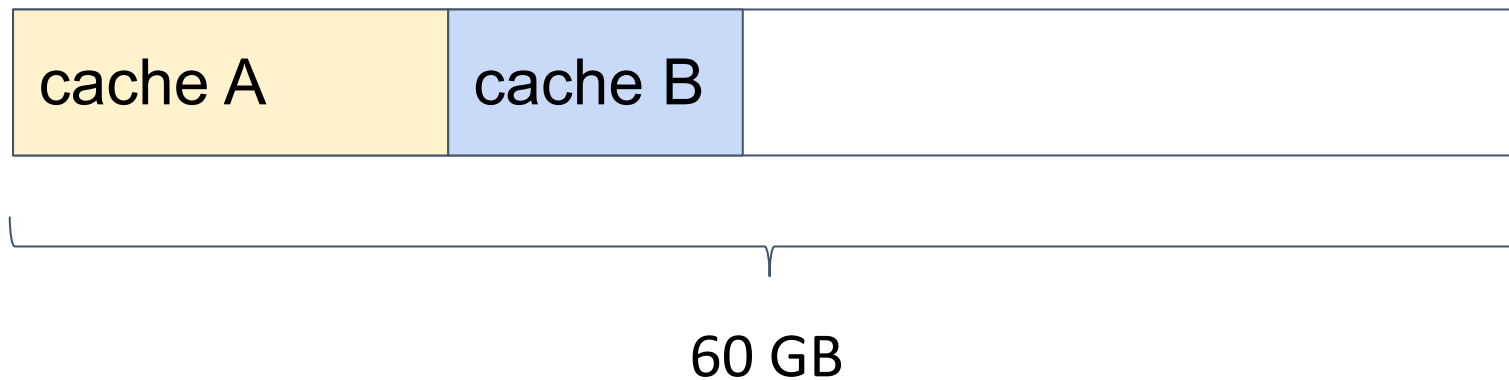- Available RAM allocated by caches on demand

GridGain

- Default: in-memory mode
- Overall RAM usage limit is configurable
- Available RAM allocated by caches on demand

```
new DataStorageConfiguration()
    .setDefaultDataRegionConfiguration(
        new DataRegionConfiguration().setMaxSize(60L * 1024 * 1024 * 1024));
```

- Default: in-memory mode
- Overall RAM usage limit is configurable
- Available RAM allocated by caches on demand

2019 © GridGain Systems

- Several "data regions"

# Use case: limit RAM consumption for specific cache

- Several "data regions"
- Each region has its own limit

GridGain

# Use case: limit RAM consumption for specific cache

- Several "data regions"
- Each region has its own limit
- Optional eviction mode: old data above the limit is removed

GridGain

- Several "data regions"
- Each region has its own limit
- Optional eviction mode: old data above the limit is removed

```
new DataStorageConfiguration()
    .setDefaultDataRegionConfiguration(
        new DataRegionConfiguration().setMaxSize(45L * 1024 * 1024 * 1024))
    .setDataRegionConfigurations(
        new DataRegionConfiguration().setName("region-with-eviction")
            .setMaxSize(15L * 1024 * 1024 * 1024)
            .setPageEvictionMode(DataPageEvictionMode.RANDOM_LRU));
```
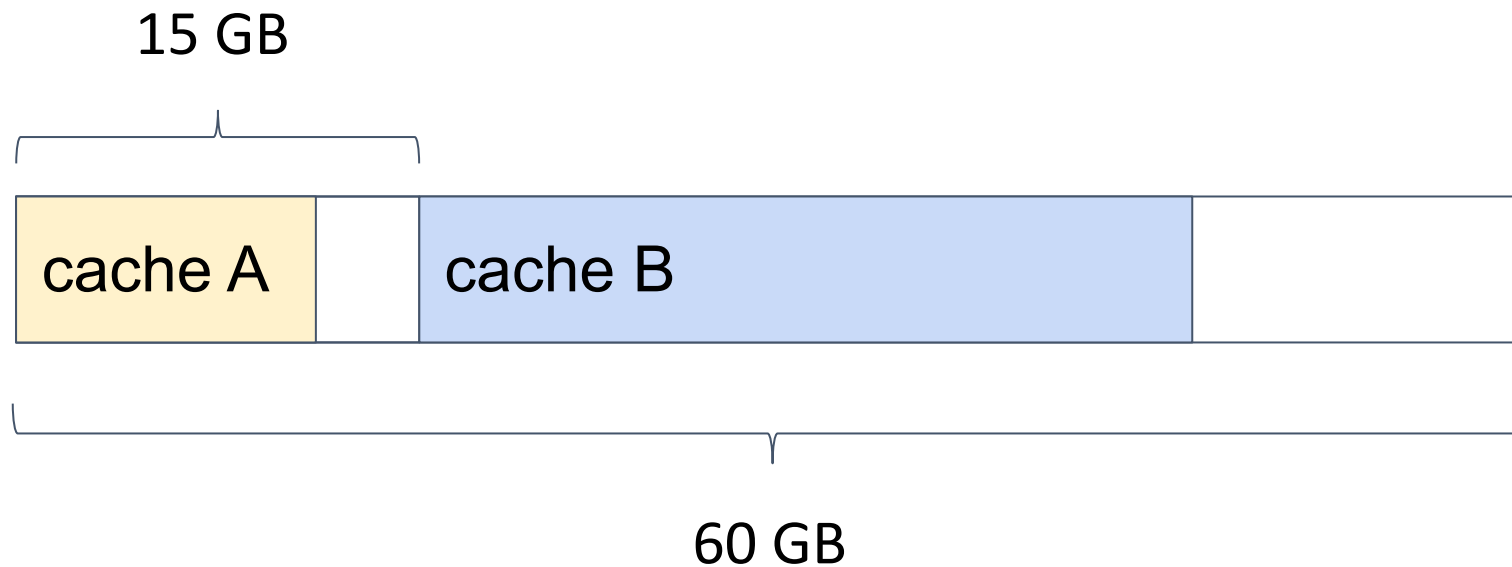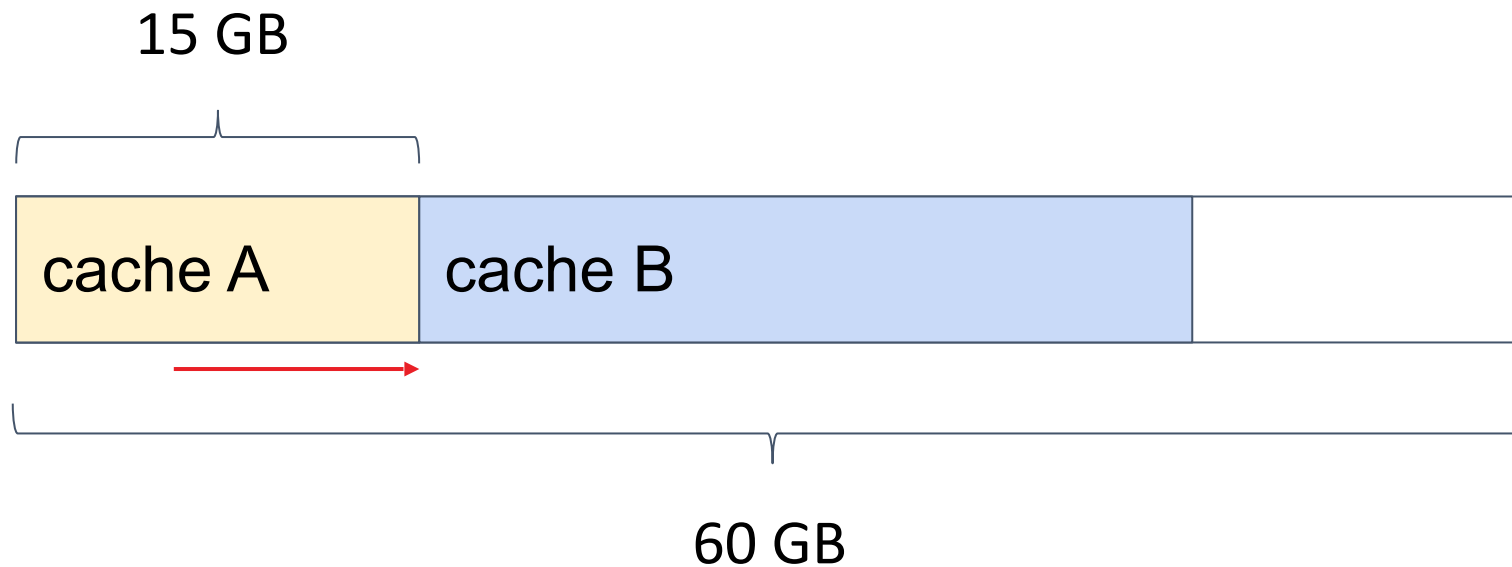
GridGain

# Use case: limit RAM consumption for specific cache

- Several "data regions"
- Each region has its own limit
- Optional eviction mode: old data above the limit is removed
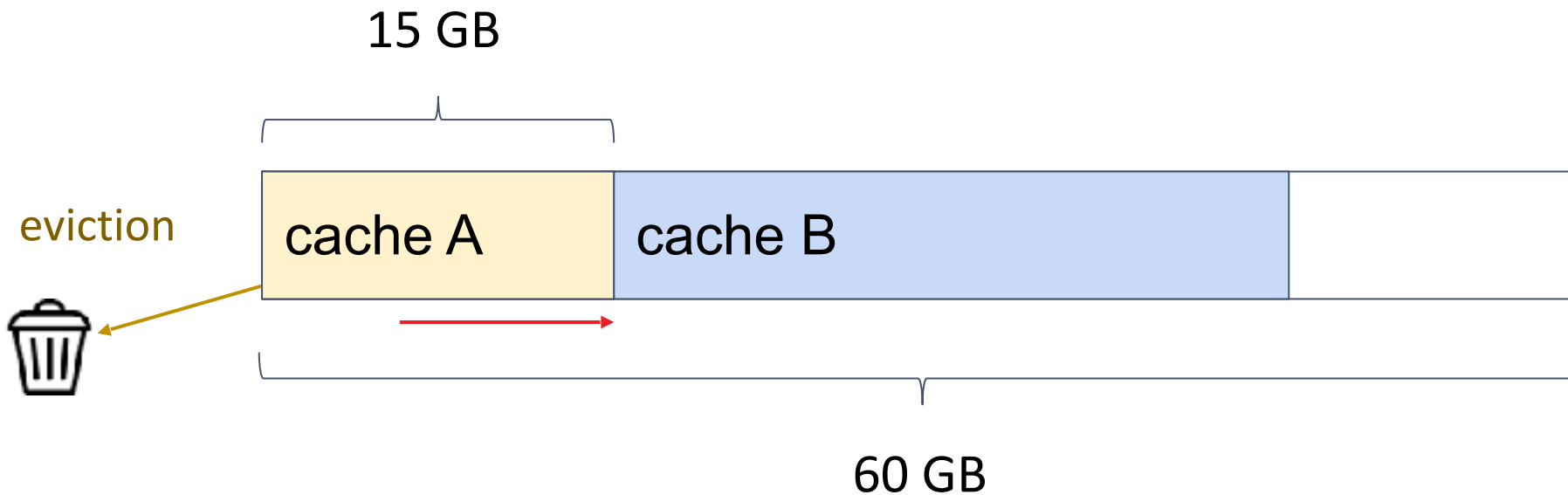
15 GB

| cache A | | cache B | |
|---------|---|---------|---|

60 GB

- Several "data regions"
- Each region has its own limit
- Optional eviction mode: old data above the limit is removed

15 GB

| cache A | cache B | |

60 GB

GridGain

- Several "data regions"
- Each region has its own limit
- Optional eviction mode: old data above the limit is removed

15 GB

eviction

| cache A | cache B | |

60 GB

GridGain

- Persistent mode: all pages on disk, subset of pages in RAM

GridGain

- Persistent mode: all pages on disk, subset of pages in RAM
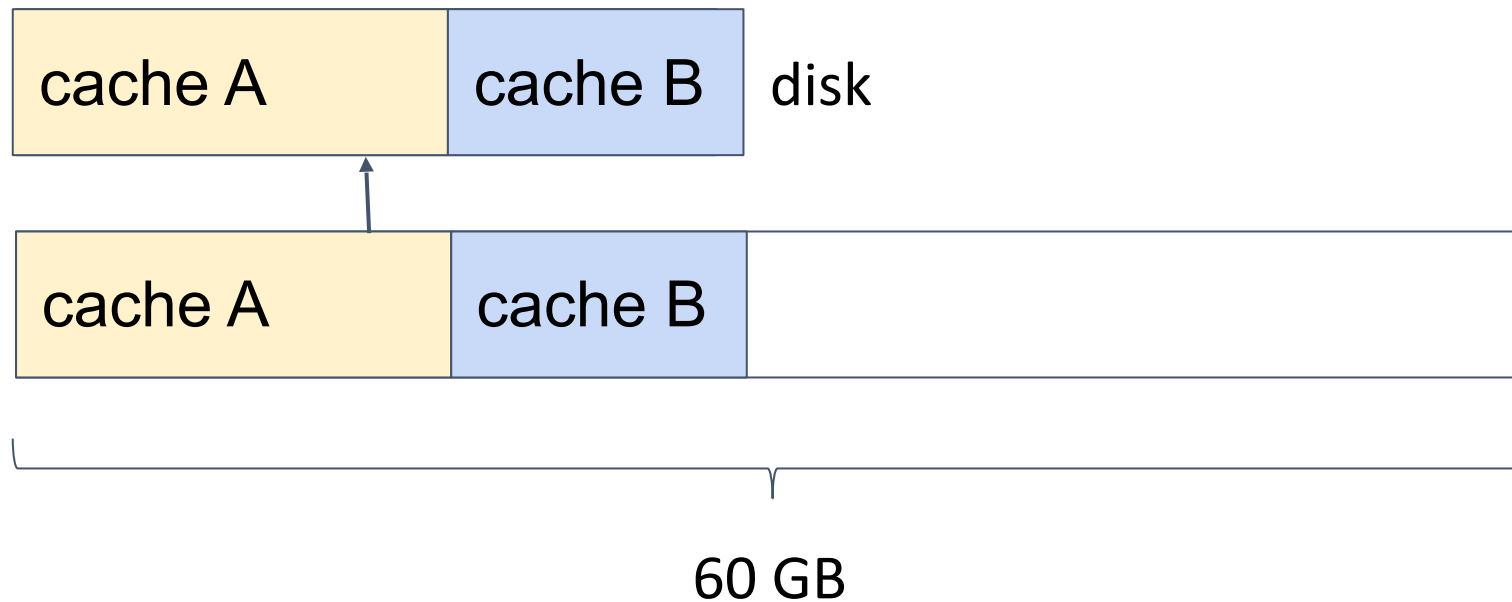- Cold pages replaced to disk on demand

- Persistent mode: all pages on disk, subset of pages in RAM
- Cold pages replaced to disk on demand

```java
new DataStorageConfiguration()
    .setDefaultDataRegionConfiguration(
        new DataRegionConfiguration().setMaxSize(60L * 1024 * 1024 * 1024)
            .setPersistenceEnabled(true));
```

- Persistent mode: all pages on disk, subset of pages in RAM
- Cold pages replaced to disk on demand

| cache A | cache B | disk |
|---------|---------|------|

| cache A | cache B | |
|---------|---------|--|

60 GB

- Persistent mode: all pages on disk, subset of pages in RAM
- Cold pages replaced to disk on demand

- Persistent mode: all pages on disk, subset of pages in RAM
- Cold pages replaced to disk on demand

2019 © GridGain Systems

# Use case: hot and cold data

- Small memory region for big cold dataset

# Use case: hot and cold data

- Small memory region for big cold dataset
- Large memory region for small hot dataset

GridGain

# Use case: hot and cold data

- Small memory region for big cold dataset
- Large memory region for small hot dataset

```
new DataStorageConfiguration()
    .setDefaultDataRegionConfiguration(
        new DataRegionConfiguration().setMaxSize(45L * 1024 * 1024 * 1024)
            .setPersistenceEnabled(true))
    .setDataRegionConfigurations(new DataRegionConfiguration().setName("cold")
        .setMaxSize(15L * 1024 * 1024 * 1024)
        .setPersistenceEnabled(true)));
```

GridGain

- Small memory region for big cold dataset
- Large memory region for small hot dataset

# Agenda

- Features that in-memory data grids lack
- Apache Ignite way: durability through page memory architecture
- Durability: use cases and solutions
  - Storage management use cases
  - Data backups use cases
- Durability: performance tricks

GridGain

- Use cases:
    - Disaster protection

# From theory to practice: Data Snapshots

- Use cases:
  - Disaster protection
  - Optimization: snapshots of non-volatile data

GridGain

# From theory to practice: Data Snapshots

- Use cases:
  - Disaster protection
  - Optimization: snapshots of non-volatile data
  - When local snapshot is not enough: remote snapshot catalog

GridGain

- Snapshot create
  - Background process
  - Current state of disk store copied to snapshot directory

2019 © GridGain Systems

- Snapshot create
  - Background process
  - Current state of disk store copied to snapshot directory
- Snapshot restore
  - Disk storage is replaced by previously saved state

GridGain

# Use case: regular snapshots of non-volatile data

GridGain

# Use case: regular snapshots of non-volatile data

- Incremental snapshot create
  - Only changed pages are written

GridGain

- Incremental snapshot create
  - Only changed pages are written
- Special page type to track changes

| idx=0 | Meta page |
|-------|-----------|
| idx=1 | Tracking page   0101010100  **1**110001001 |
| idx=2 | Regular page |
| idx=3 | Regular page |
| idx=4 | Regular page |

← **data update**

2019 © GridGain Systems

- Complete disaster (local snapshots are lost as well)

GridGain

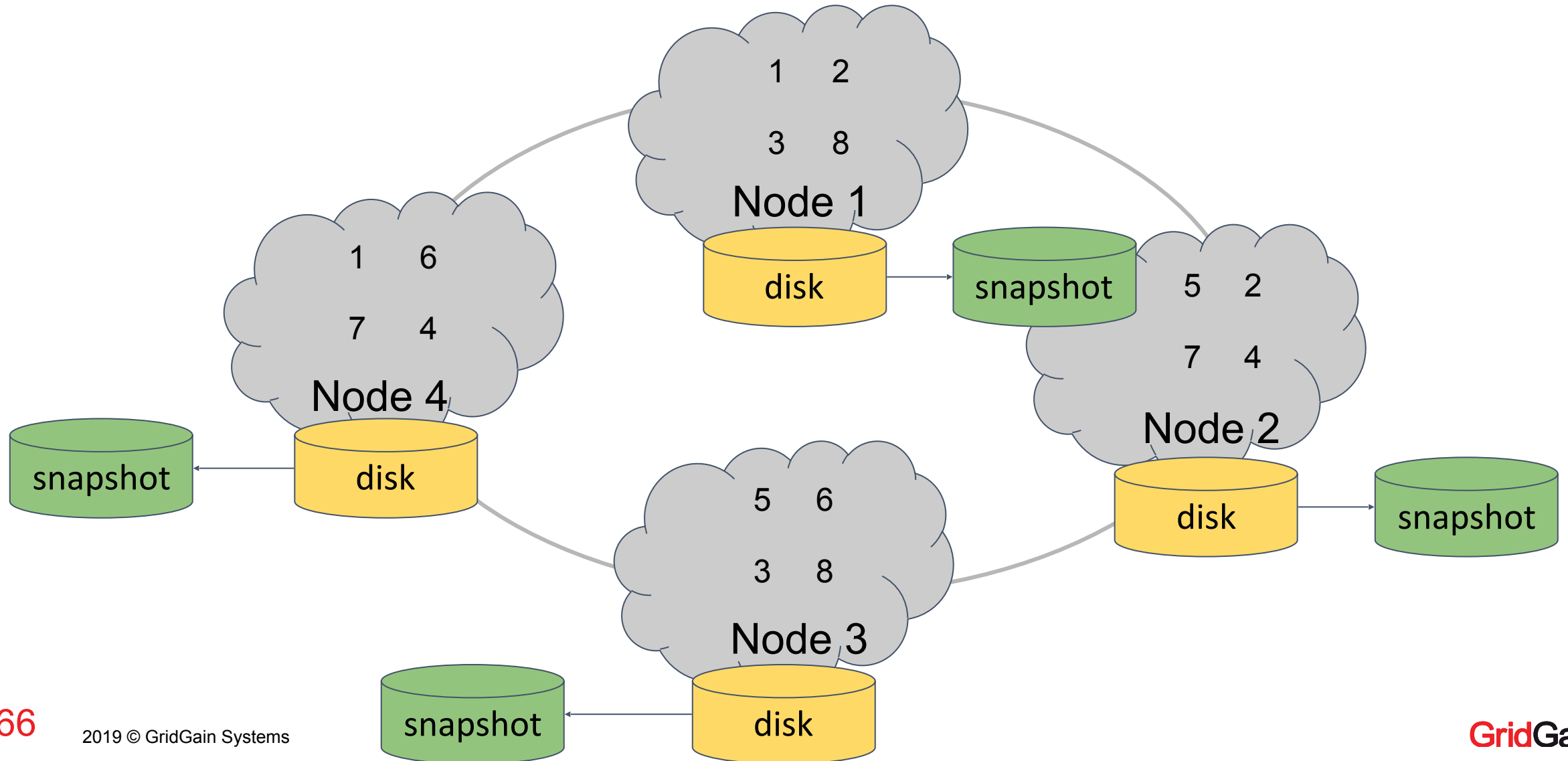- Complete disaster (local snapshots are lost as well)
- Daily snapshot catalog

- Complete disaster (local snapshots are lost as well)
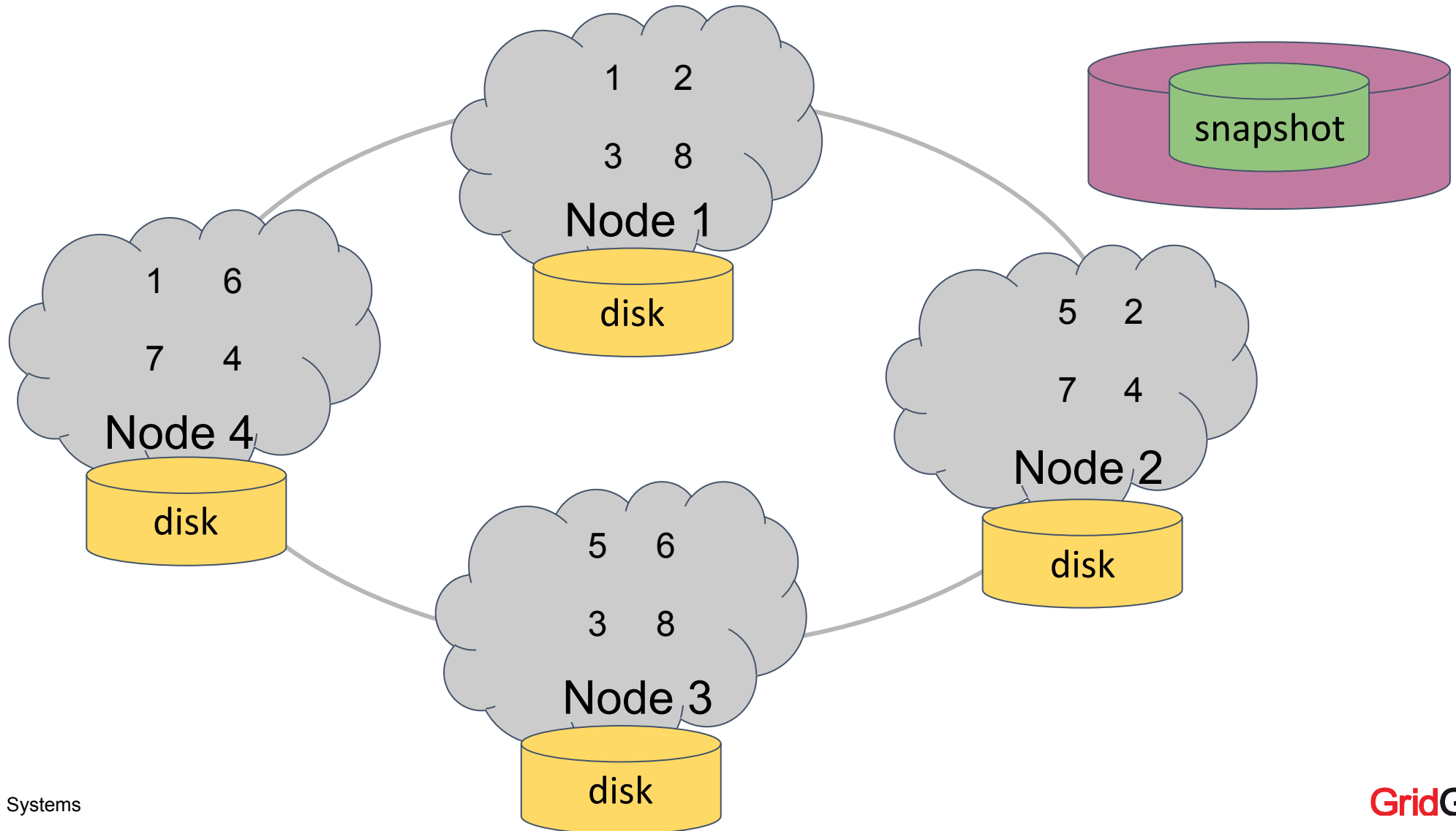- Daily snapshot catalog
- Restore after topology change

GridGain

- Snapshot move to shared folder
  - Data from the whole cluster is moved to reliable network storage

- Snapshot move to shared folder
  - Data from the whole cluster is moved to reliable network storage
- Snapshot restore from shared folder
  - Even if topology was changed, all data partitions will be found

GridGain

# Agenda

- Features that in-memory data grids lack
- Apache Ignite way: durability through page memory architecture
- Durability: use cases and solutions
  - Storage management use cases
  - Data backups use cases
- Durability: performance tricks

GridGain

# Memory / disk ratio affects performance directly
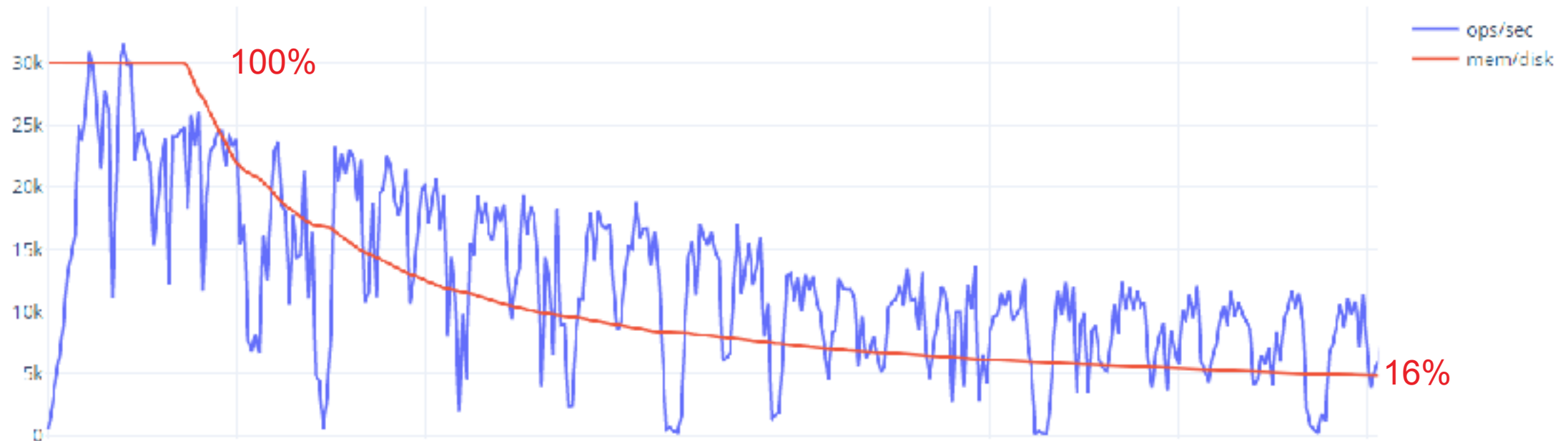
- Every page absent in RAM will require synchronous read

GridGain

# Memory / disk ratio affects performance directly

- Every page absent in RAM will require synchronous read
- Latency grows along with share of "disk only" pages

2019 © GridGain Systems

- Peak load throughput can be higher than disk throughput

- Peak load throughput can be higher than disk throughput

- Peak load throughput can be higher than disk throughput

```
dataStorageCfg.setWriteThrottlingEnabled(true);
```

# Overprovision your SSD

- SSDs are designed to be easily adopted by OS
  - "write K bytes to 0xFF…" interface like HDD

GridGain

- SSDs are designed to be easily adopted by OS
  - "write K bytes to 0xFF…" interface like HDD
- But actually SSD is a complex computer itself

2019 © GridGain Systems

# Overprovision your SSD

- SSD has pages and blocks (64/128/256 pages)

GridGain

# Overprovision your SSD

- SSD has pages and blocks (64/128/256 pages)
- Data is written in page granularity, erased in block granularity

GridGain

# Overprovision your SSD

- SSD has pages and blocks (64/128/256 pages)
- Data is written in page granularity, erased in block granularity
- Block erase requires shifting useful data to another block

GridGain

# Overprovision your SSD

- SSD has pages and blocks (64/128/256 pages)
- Data is written in page granularity, erased in block granularity
- Block erase requires shifting useful data to another block
- Shifting is easier when more free blocks are available

GridGain

# Overprovision your SSD



Random Writes (4KB sustained)

2019 ©

GridGain

# Overprovision your SSD

- Beware: SSD performance issues do not appear immediately

GridGain

# Overprovision your SSD

- Beware: SSD performance issues do not appear immediately



GC activity is no longer able to keep up in background

- Every data update is written twice in persistent mode

# Disable WAL on initial data load

- Every data update is written twice in persistent mode
  - To journal (write-ahead log), synchronously

GridGain

# Disable WAL on initial data load

- Every data update is written twice in persistent mode
  - To journal (write-ahead log), synchronously
  - To disk storage (on checkpoint), asynchronously

GridGain

# Disable WAL on initial data load

- Every data update is written twice in persistent mode
  - To journal (write-ahead log), synchronously
  - To disk storage (on checkpoint), asynchronously
- WAL can be disabled on purpose

# Disable WAL on initial data load

- Every data update is written twice in persistent mode
    - To journal (write-ahead log), synchronously
    - To disk storage (on checkpoint), asynchronously
- WAL can be disabled on purpose
    - Crash recovery is not guaranteed

GridGain

# Disable WAL on initial data load

- Every data update is written twice in persistent mode
  - To journal (write-ahead log), synchronously
  - To disk storage (on checkpoint), asynchronously
- WAL can be disabled on purpose
  - Crash recovery is not guaranteed
  - At least 2x load throughput boost

GridGain

# Disable WAL on initial data load

- Every data update is written twice in persistent mode
  - To journal (write-ahead log), synchronously
  - To disk storage (on checkpoint), asynchronously
- WAL can be disabled on purpose
  - Crash recovery is not guaranteed
  - At least 2x load throughput boost
- `igniteCluster.disableWal(cacheToLoad);`

GridGain

- Persistent Ignite node has four disk write activities

GridGain

# Consider using separate devices

- Persistent Ignite node has four disk write activities
    - Checkpointing

2019 © GridGain Systems

GridGain

# Consider using separate devices

- Persistent Ignite node has four disk write activities
    - Checkpointing
    - Writing WAL

GridGain

# Consider using separate devices

- Persistent Ignite node has four disk write activities
  - Checkpointing
  - Writing WAL
  - Transferring old WAL segments to WAL archive dir

# Consider using separate devices

- Persistent Ignite node has four disk write activities
  - Checkpointing
  - Writing WAL
  - Transferring old WAL segments to WAL archive dir
  - Data snapshotting

GridGain

# Consider using separate devices

- Persistent Ignite node has four disk write activities
  - Checkpointing
  - Writing WAL
  - Transferring old WAL segments to WAL archive dir
  - Data snapshotting
- Separate path can be configured for each
  - `dataStorageCfg.setStoragePath(…);`
  - `dataStorageCfg.setWalPath(…);`
  - `dataStorageCfg.setWalArchivePath(…);`
  - `snapshotCfg.setSnapshotsPath(…);`

GridGain

# Performance tips: summary

- Plan memory / disk ratio for your performance requirements
- Use throttling for smooth throughput
- Overprovision your SSD
- Disable WAL on initial data load
- Split disk activities on separate storage devices

GridGain

**Thanks for your attention!**
**Questions?**

GridGain