

Act as a senior OpenShift and vSphere/vSAN performance engineer. Your goal is to produce a precise, actionable plan and ready-to-send communication to progress an ongoing case about the persistent etcdGRPCRequestsSlow alert on multiple OpenShift clusters. Think deeply, but output only the final structured deliverable requested below.

Context (facts and observations):

- Alert: etcdGRPCRequestsSlow on multiple clusters for >6 months.
- Platform: vSphere with vSAN. Control plane (master) VMs share the same datastore/LUN as worker/infra and other workloads. Masters are not strictly 1:1 per host; other VMs co-resident. Hosts reported as not overcommitted; utilization "normal."
- Disk benchmarks (node-level):
 - dd 4k*100000: ~101s (~4.1 MB/s)
 - dd 1k*400000: ~444s (~0.9 MB/s)
- FIO benchmarks (sequential, generic):
 - ~1424 read IOPS; ~611 write IOPS (expectation noted: 1500–2000 sequential IOPS for healthy baseline)
 - Mixed 30% read / 70% write: ~425 read IOPS; ~984 write IOPS
- etcd metrics:
 - wal_fsync up to ~20 ms (target <10 ms)
 - round-peer time ~100 ms (target <50 ms)
 - commit duration OK; CPU iowait OK
- Networking: Prior Tx ring buffer "ring full" events observed, mitigated after increasing ring size; currently acceptable.
- etcd encryption: enabled (aesCBC). ~145 total secrets; ~75 encryption-config secrets; some older than 30 days (cleanup pending).
- Current guidance received: dedicate datastore for masters; spread control plane across separate hypervisor hosts; avoid shared datastore contention. Your deliverable (use these exact section headers):

1. Executive Summary

- 3–5 bullets stating the current situation, most likely cause(s), and the fastest path to relief.

2. Likely Root Causes (ranked)

- Rank 3–5 hypotheses with brief evidence citing the specific data above and expected vs observed thresholds.

3. Validation Plan (commands, thresholds, time budget)

- Concrete tests to confirm/deny each hypothesis (fio, dd, etcd/grpc latency buckets, wal_fsync, disk latency from guest and host, vSAN metrics, ESXi counters).
- Provide exact commands, parameters, run duration, and pass/fail thresholds aligned to etcd requirements.

4. Remediation Plan (step-by-step, risk, rollback)

- Storage layout: dedicated datastore/SPBM policy for masters, cache/tiering considerations, reservations/limits, contention isolation.
- vSphere/ESXi: CPU/mem reservations, latency sensitivity, disk controller/queue depth, adaptive queuing, path selection policy.
- OpenShift/etcd: encryption-config secret rotation cleanup, compaction tuning, defrag cadence, pod anti-affinity, PDBs for safe changes.
- Networking (only if warranted): NIC offloads/queues tuning.
- For each step: objective, exact actions, expected outcome, risk, and rollback.

5. Meeting Plan (VMware + Red Hat + customer)

- 45–60 min agenda, pre-reads/artifacts to collect, and targeted questions for each party to quickly converge on storage and scheduling contention.

6. Ready-to-Send Email to Red Hat Support

- Concise summary of findings, test evidence, specific asks (confirm acceptable IOPS/latency targets for etcd on vSAN, guidance on datastore isolation), and next steps.
- Include a bulleted list of artifacts you will attach.

7. Monitoring & Acceptance Criteria

- Explicit SLOs: wal_fsync <10 ms p99, etcd gRPC latency buckets within runbook targets, round-peer <50 ms, compaction/defrag health, and node-level I/O latency.
- Dashboards/alerts to add and how to interpret them.

8. Missing Data & How to Gather

- Short list of any gaps plus exact commands/vCenter/vSAN views to capture them. Output requirements:
- Be concise and specific. No generic advice. Tie every recommendation to the evidence above.
- Include concrete commands, settings, and numeric thresholds.
- Use clear headers matching the sections, bullet points, and tables where comparisons help.
- Assume no external links; describe what to check by name if you would normally cite a document.

Ref: <https://github.com/openshift/runbooks/blob/master/alerts/cluster-etcd-operator/etcdGRPCRequestsSlow.md>

1. Executive Summary

- Persistent etcdGRPCRequestsSlow across multiple OpenShift clusters aligns with storage latency symptoms: wal_fsync up to ~20 ms vs target <10 ms p99 and round-peer ~100 ms vs target <50 ms, indicating disk and replication latency likely from shared vSAN datastore contention with non-control-plane workloads. ^1^3

- Node-level benchmarks are below healthy baselines for control plane: sequential ~1424 read IOPS/~611 write IOPS and small-block dd 4k/1k at ~4.1 MB/s/~0.9 MB/s, which is weak relative to etcd recommendations of ≥ 50 seq 8K IOPS with ≤ 10 ms fsync for minimum and ~500 seq 8K IOPS for heavier clusters.^{^5}
- Fastest path to relief: isolate masters on a dedicated vSAN policy/datastore with low-latency all-flash settings, ensure controller queue depth ≥ 256 , apply CPU/memory reservations and host anti-affinity for masters, run etcd defrag, and clean up excess/old encryption secrets to reduce I/O churn, then re-measure wal_fsync p99 and gRPC p99.^{^2^7^1}
- Validate with focused tests: fio 8K sequential with fsync and PromQL for gRPC p99 and wal_fsync p99 over 15–30 minutes during normal load; gate acceptance on wal_fsync p99 < 10 ms and gRPC p99 within runbook targets, and round-peer < 50 ms.^{^3^1}

2. Likely Root Causes (ranked)

- Shared vSAN datastore contention for masters: wal_fsync p99 observed up to ~20 ms (> 10 ms target) and sequential write IOPS ~611 with weak small-block dd results suggest storage latency under mixed load; masters sharing LUN/datastore with workers/infra increases risk of queueing and cache contention; etcd guidance requires low-latency storage with ≥ 50 8K IOPS (fsync ≤ 10 ms) minimum and ~500 8K IOPS for heavier clusters, which current measurements only marginally meet for reads and lag for writes.^{^4^5}
- etcd replication/network latency: round-peer ~100 ms exceeds < 50 ms target, and while prior NIC Tx ring “ring full” was mitigated, inter-member RTT or vSAN network microbursts can still elevate gRPC latency under load, especially when backend I/O stalls lengthen raft commit paths.^{^9^3}
- CPU scheduling headroom on control plane VMs: without explicit CPU/memory reservations and optional Latency Sensitivity for masters, transient ready time or NUMA placement can add tail latency to etcd threads even at “normal” host utilization, exacerbating p99 gRPC duration.^{^6}
- etcd DB fragmentation and encryption overhead: encryption aescbc and an elevated count of secrets (~145 total, ~75 encryption-config, some stale) increase write amplification; defragmentation plus secret/rotation cleanup typically reduce fsync pressure and tail latency on wal.^{^1^3}

3. Validation Plan (commands, thresholds, time budget)

- Prometheus (15–30 min window during normal load):
 - gRPC p99:
 - Query: `histogram_quantile(0.99, sum(rate(grpc_server_handling_seconds_bucket{job=~".etcd.",grpc_type="unary"}[5m])) without(grpc_type)).`^{^1}
 - Pass: method p99 stable and trending down after changes; critical paths (Txn/Range/Watch) not breaching during steady state; target is removal of etcdGRPCRequestsSlow.^{^1}
 - wal_fsync p99:
 - Query: `histogram_quantile(0.99, sum by (instance,le) (irate(etcd_disk_wal_fsync_duration_seconds_bucket{job="etcd"}[5m])))`.^{^2}
 - Pass: < 10 ms p99 sustained for ≥ 15 min after changes; fail if ≥ 10 ms over multiple 5-min windows.^{^2}
 - round-peer (raft round trip): inspect distribution vs < 50 ms target during same window; pass if median and p99 < 50 ms without spikes.^{^3}
- fio (each master VM, 10–15 min each, off-peak and peak):

- 8K sequential write with fsync (minimum requirement):
 - Command: `fio --name=wal --filename=/var/lib/etcd/wal.fio --rw=write --bs=8k --iodepth=1 --ioengine=libaio --direct=1 --fsync=1 --time_based --runtime=600 --group_reporting --numjobs=1 --size=2G.`^{^4}
 - Pass: ≥ 50 IOPS with p99 latency ≤ 10 ms; fail if p99 > 10 ms or IOPS < 50 .^{^4}
- 8K sequential read:
 - Command: `fio --name=walr --filename=/var/lib/etcd/wal.fio --rw=read --bs=8k --iodepth=4 --ioengine=libaio --direct=1 --time_based --runtime=600 --group_reporting --numjobs=1.`^{^4}
 - Pass: ≥ 100 IOPS and stable latency; for heavier clusters aim near 500 IOPS capability.^{^5}
- Mixed 30/70 (sanity vs observed):
 - Command: `fio --name=mix --filename=/var/lib/etcd/mix.fio --rw=randrw --rwmixread=30 --bs=8k --iodepth=4 --ioengine=libaio --direct=1 --time_based --runtime=600 --group_reporting --numjobs=1.`^{^4}
 - Pass: consistency under load without dramatic tail latency growth; compare to current $\sim 425/984$ IOPS mix.^{^4}
- dd (quick parity with previous tests, 5–10 min):
 - 4k: `dd if=/dev/zero of=/var/lib/etcd/dd4k bs=4k count=100000 oflag=direct`; compare duration to ~ 101 s baseline; slower implies regression.^{^4}
 - 1k: `dd if=/dev/zero of=/var/lib/etcd/dd1k bs=1k count=400000 oflag=direct`; compare to ~ 444 s baseline.^{^4}
- Guest OS and etcd:
 - `iostat -x 5 60` and `sar -d 5 60` on masters to correlate await/svctm with etcd wal_fsync spikes; pass if await consistently low during defrag/steady operations.^{^11}
 - etcd defrag safety check then `etcdctl defrag` on all members; re-measure wal_fsync p99 for improvement.^{^2}
- vSphere/vSAN host and datastore:
 - vCenter/vSAN perf charts: check backend write latency, congestion, resync traffic, outstanding I/O per disk group, and cache hit ratios during the same window; pass if backend latency remains sub-millisecond to low-millisecond and no sustained congestion.^{^8}
 - ESXi queue depth behavior: validate controller queue depth ≥ 256 and Adaptive Queueing (Disk.QFullSampleSize/QFullThreshold) consistent across hosts; pass if no repeated QFULL backoffs and adequate queue headroom.^{^12^8}
 - `esxtop` (disk/adapters) and `vscsiStats` for per-VM latency and queueing; pass if DAVG/KAVG low and GAVG correlates with improved wal_fsync.^{^14}

4. Remediation Plan (step-by-step, risk, rollback)

- Storage layout (objective: isolate and lower I/O latency for etcd):
 - Actions: create/assign a dedicated vSAN storage policy/datastore (or dedicated SPBM policy) for master VMDKs with all-flash, RAID1 (FTT=1) for lowest latency, ensure no non-control-plane VMs share it, and verify storage controller queue depth ≥ 256 on hosts running masters.^{^7}
 - Expected: reduced contention and lower wal_fsync p99 to < 10 ms, and improved round-peer via faster commits.^{^2}
 - Risk: temporary VMotion/storage migration impact; Rollback: revert storage policy and VM placement to prior datastore if unintended latency increase is observed.^{^14}
- vSphere/ESXi scheduling (objective: reduce CPU-induced tail latency):

- Actions: set CPU reservations equal to master vCPU allocations and memory reservations equal to configured memory on all master VMs; optionally set Latency Sensitivity=High for masters during validation window, ensure no CPU hot-add and NUMA alignment; ensure PVSCSI with sufficient vSCSI adapters to avoid in-guest queue bottlenecks.^{^10}
 - Expected: near-zero CPU ready time for etcd, reducing gRPC p99 and stabilizing round-peer <50 ms.^{^6}
 - Risk: resource contention for other VMs; Rollback: lower reservations and restore Latency Sensitivity to Normal if cluster contention arises.^{^10}
- Queue depth and pathing (objective: avoid host-level throttling):
 - Actions: confirm storage controller queue depth ≥ 256 ; enable Adaptive Queueing consistently (Disk.QFullSampleSize and Disk.QFullThreshold) across hosts if array/vSAN guidance allows; ensure PSP/path selection is consistent and recommended by storage stack; monitor for QFULL events.^{^13^8}
 - Expected: fewer QFULL-induced backoffs and smoother I/O under bursts, improving wal_fsync tails.^{^12}
 - Risk: misconfiguration can cause uneven backpressure; Rollback: restore previous queue depth and AQ settings using host profiles/change records.^{^12}
- OpenShift/etcd maintenance (objective: reduce write amplification and fragmentation):
 - Actions: prune stale encryption-config secrets and rotate per policy to minimize config set size; run etcdctl defrag on all members during low activity; schedule periodic defrag during maintenance; verify compaction cadence is healthy via existing operator settings.^{^1}
 - Expected: smaller on-disk DB and faster fsync/boltdb operations, reducing wal_fsync p99.^{^3}
 - Risk: transient I/O during defrag; Rollback: if anomalies occur, pause further defrags and collect metrics for support review.^{^1}
- Placement/anti-affinity (objective: spread control plane risk):
 - Actions: enforce VM-Host anti-affinity to spread masters across different ESXi hosts and fault domains; avoid co-residency with known heavy I/O VMs on the master datastore/policy.^{^14}
 - Expected: lower correlated latency spikes across members and improved quorum stability.^{^3}
 - Risk: capacity/DRS constraints; Rollback: relax rules temporarily if placement fails.^{^14}
- Networking (if round-peer still high) (objective: lower raft RTT/jitter):
 - Actions: ensure vSAN/management networks on 10/25GbE with proper NIC queues/offloads; confirm no ring buffer saturation and that vSAN network best practices are met; re-test round-peer and gRPC p99.^{^9}
 - Expected: reduced network-induced tail latency for raft replication.^{^9}
 - Risk: network change windows; Rollback: revert NIC/driver settings to previous known-good config.^{^9}

5. Meeting Plan (VMware + Red Hat + customer)

- Agenda (45–60 min):
 - 0–10 min: Present observed metrics (gRPC p99, wal_fsync p99, round-peer) and node-level fio/dd results vs thresholds; highlight datastore sharing and duration of alert.^{^2^4}
 - 10–25 min: VMware review—vSAN policy/datastore layout, controller queue depth, AQ settings, vSAN perf charts (backend latency, congestion, resync), host reservations/NUMA and anti-affinity.^{^8^9}
 - 25–40 min: Red Hat review—confirm etcd SLOs, defrag/compaction guidance, encryption overhead expectations, and any cluster-side contributors (watch storms/endpoints churn).^{^1}

- 40–55 min: Agree on remediation/validation timeline and acceptance gates (wal_fsync p99 <10 ms, round-peer <50 ms, alert clear for 7 days).^{^3}
- Pre-reads/artifacts: PromQL exports for gRPC p99 and wal_fsync p99, fio outputs, vCenter/vSAN perf screenshots for master datastore and hosts, ESXi queue depth/AQ config evidence, and encryption secrets inventory.^{^8^4}
- Targeted questions: to VMware—what is controller queue depth on master hosts, is AQ enabled consistently, any vSAN congestion/resync during spikes; to Red Hat—confirm acceptance thresholds and any known operator churn patterns elevating write rates; to customer—change windows for datastore isolation and reservations.^{^12^4}

6. Ready-to-Send Email to Red Hat Support

- Subject: Ongoing etcdGRPCRequestsSlow – storage latency and replication tail verification, plan for remediation and validation.^{^1}
- Body:
 - We continue to see etcdGRPCRequestsSlow across multiple clusters; wal_fsync p99 peaks ~20 ms (target <10 ms) and round-peer ~100 ms (target <50 ms) with node-level sequential write capability ~611 IOPS and weak small-block dd results, indicating storage and replication tail latency as primary drivers.^{^2}
 - Plan: isolate masters onto a dedicated vSAN policy/datastore, set CPU/memory reservations for master VMs, verify controller queue depth ≥256 and enable Adaptive Queueing consistently, perform etcd defrag and clean up stale encryption-config secrets, then re-measure gRPC p99 and wal_fsync p99 over 15–30 minutes under steady load.^{^8^2}
 - Asks: please confirm current acceptable IOPS/latency targets for etcd on vSAN (we're using ≥50 IOPS 8K with fsync ≤10 ms minimum and ~500 IOPS for heavier clusters, wal_fsync p99 <10 ms, round-peer <50 ms), and whether additional cluster-side mitigations are recommended beyond defrag/compaction and encryption cleanup.^{^2}
 - Next steps: we will send pre/post change metrics and fio outputs; request validation that the metrics meet runbook targets once applied.^{^1}
- Attachments: PromQL exports (gRPC p99, wal_fsync p99), fio outputs (8K write/read/mixed), vCenter/vSAN perf screenshots (master datastore/hosts), ESXi queue depth/AQ settings, encryption secrets inventory and cleanup plan.^{^12^1}

7. Monitoring & Acceptance Criteria

- SLOs for closure: wal_fsync p99 <10 ms over 7 consecutive days, etcd gRPC p99 within runbook targets without alert re-fires, round-peer <50 ms p99, regular compaction/defrag healthy, and guest/host I/O latency stable under normal change bursts.^{^3^1}
- Dashboards/alerts:
 - Observe > etcd dashboard: verify Disk Sync Duration and RPC latency trends; add alerts on wal_fsync p99 >10 ms for 10 min and round-peer p99 >50 ms for 10 min to catch regressions early.^{^2}
 - vSAN/host dashboards: track backend latency, congestion, and outstanding I/Os on master hosts/datastore; alert on sustained congestion or rising device latency during control plane churn windows.^{^9}
 - Capacity/scheduling: monitor master VM CPU ready and reserved capacity utilization to ensure reservations maintain headroom and low jitter.^{^6}

8. Missing Data & How to Gather

- vSAN perf evidence around alert windows: capture backend latency, congestion, resync traffic, and outstanding I/O per disk group from vCenter/vSAN Performance for the master datastore/policy and master hosts during gRPC/wal_fsync peaks. ⁸
- ESXi queue depth and AQ: record controller queue depth per host and Disk.QFullSampleSize/QFullThreshold settings; note any QFULL events; standardize across hosts; use Host Client/Host Profiles or esxcli as needed. ¹³
- Detailed Prometheus exports: export 24h gRPC p99 by method and wal_fsync p99 by instance before and after changes for support review and acceptance tracking. ¹
- CPU scheduling evidence: capture VM CPU ready, co-stop, and NUMA alignment for masters during peak; verify reservations eliminate ready time spikes. ¹⁰

Notes on “most likely cause” and “fastest relief” are based on exceeding wal_fsync and round-peer thresholds under shared datastore conditions and weak I/O baselines; the outlined isolation, reservations, and queue depth steps directly target disk and replication tails that drive this alert per runbook and etcd performance guidance. ^{3^2}

*
**