

# Reproducing Filesystem Delay Experiments: One-Size-Fits-None (NSDI'25)

## System Requirements Prerequisites (Ubuntu 20.04/22.04 tested)

# Docker + Compose v2

```
sudo apt-get update
```

```
sudo apt-get install -y docker.io docker-compose-plugin
```

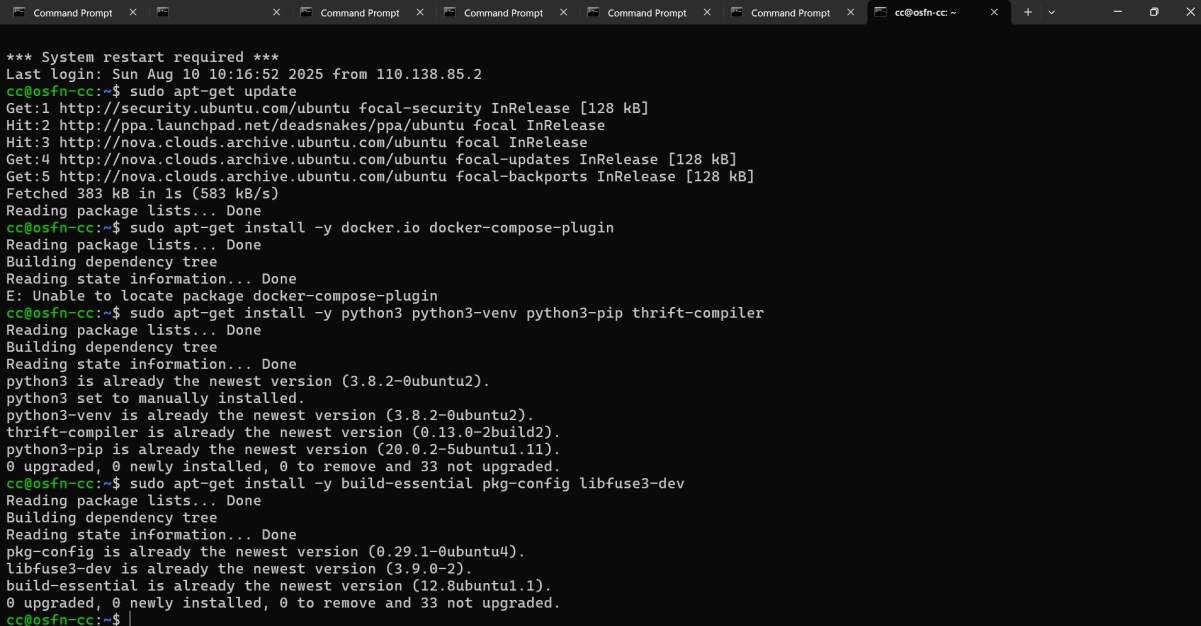
# Python 3 + thrift compiler (Charybdefs client uses Thrift)

```
sudo apt-get install -y python3 python3-venv python3-pip
```

```
thrift-compiler
```

# FUSE3 dev (to build Charybdefs) - if built locally

```
sudo apt-get install -y build-essential pkg-config libfuse3-dev
```

A terminal window with multiple tabs. The active tab shows the output of several terminal commands. It starts with a system restart notice, followed by 'sudo apt-get update' which lists several updates. Then 'sudo apt-get install -y docker.io docker-compose-plugin' is run, followed by 'sudo apt-get install -y python3 python3-venv python3-pip thrift-compiler'. The output shows that python3, python3-venv, python3-pip, and thrift-compiler are already the newest versions. Finally, 'sudo apt-get install -y build-essential pkg-config libfuse3-dev' is run, and the output shows that build-essential, pkg-config, and libfuse3-dev are already the newest versions.

```
*** System restart required ***
Last login: Sun Aug 10 10:16:52 2025 from 110.138.85.2
cc@osfn-cc:~$ sudo apt-get update
Get:1 http://security.ubuntu.com/ubuntu focal-security InRelease [128 kB]
Hit:2 http://ppa.launchpad.net/deadsnakes/ppa/ubuntu focal InRelease
Hit:3 http://nova.clouds.archive.ubuntu.com/ubuntu focal InRelease
Get:4 http://nova.clouds.archive.ubuntu.com/ubuntu focal-updates InRelease [128 kB]
Get:5 http://nova.clouds.archive.ubuntu.com/ubuntu focal-backports InRelease [128 kB]
Fetched 383 kB in 1s (583 kB/s)
Reading package lists... Done
cc@osfn-cc:~$ sudo apt-get install -y docker.io docker-compose-plugin
Reading package lists... Done
Building dependency tree
Reading state information... Done
E: Unable to locate package docker-compose-plugin
cc@osfn-cc:~$ sudo apt-get install -y python3 python3-venv python3-pip thrift-compiler
Reading package lists... Done
Building dependency tree
Reading state information... Done
python3 is already the newest version (3.8.2-0ubuntu2).
python3 set to manually installed.
python3-venv is already the newest version (3.8.2-0ubuntu2).
thrift-compiler is already the newest version (0.13.0-2build2).
python3-pip is already the newest version (20.0.2-5ubuntu1.11).
0 upgraded, 0 newly installed, 0 to remove and 33 not upgraded.
cc@osfn-cc:~$ sudo apt-get install -y build-essential pkg-config libfuse3-dev
Reading package lists... Done
Building dependency tree
Reading state information... Done
pkg-config is already the newest version (0.29.1-0ubuntu4).
libfuse3-dev is already the newest version (3.9.0-2).
build-essential is already the newest version (12.8ubuntu1.1).
0 upgraded, 0 newly installed, 0 to remove and 33 not upgraded.
cc@osfn-cc:~$
```

## 1) Charybdefs: build, mount, and management service

The paper uses Charybdefs to inject filesystem slowness via a (Thrift) management interface.

If you already have Charybdefs built and mounted, skip to Step 2 and just keep the paths the same.

# Build (any clean location is fine)

```
cd ~
```

```
git clone https://github.com/scylladb/charybdefs.git
cd charybdefs
```

```
# builds the FUSE FS and the mgmt server
make
```

```
cc@osfn-cc:~$ ls
CHI-210850-openrc.sh  blockade--venv  blockade-venv  charybdefs  my_mounting_point  sweep_logs  traces
README               blockade-py36-venv  cassandra-demo  logs        openrc           trace-ACER
cc@osfn-cc:~$ cd charybdefs
cc@osfn-cc:~/charybdefs$ make
[ 20%] Built target faultcore
[ 33%] Built target charybde_lib
[ 40%] Built target thrift_gen
[ 73%] Built target server
[ 86%] Built target server_lib
[100%] Built target charybdefs
cc@osfn-cc:~/charybdefs$ |
```

```
# Create a mount point that will back the etcd data dirs
sudo mkdir -p /mnt/slowfs/{etcd0,etcd1,etcd2}
```

```
# Run Charybdefs (FS + mgmt on 0.0.0.0:9090).
# The typical pattern is one charybdefs instance serving a root
mount and exposing a Thrift RPC on 9090.
# Adjust the exact binary/flags if your tree differs (some builds
split FS and mgmt).
sudo ./build/charybdefs -o allow_other -m /mnt/slowfs &
```

```
# start mgmt (Thrift) - name/flags vary slightly by branch; if
unified binary, omit this
./build/charybdefs-mgmt --bind 0.0.0.0 --port 9090 &
```

```
cc@osfn-cc:~/charybdefs$ sudo mkdir -p /mnt/slowfs/{etcd0,etcd1,etcd2}
cc@osfn-cc:~/charybdefs$ sudo ./build/charybdefs -o allow_other -m /mnt/slowfs &
[1] 2489737
```

- We'll point etcd's data directories to /mnt/slowfs/etcd{0,1,2} so **every WAL/fsync** route goes through Charybdefs and can be delayed.
- The mgmt endpoint is Thrift on TCP :9090 (as in the paper).

## 2) Bring up a 3-node etcd cluster (Docker Compose)

```
# Bring up the cluster
```

```
cd ~/cassandra-demo
docker compose -f docker-compose-etcd.yml up -d
```

```
cc@osfn-cc:~/cassandra-demo$ docker compose -f docker-compose-etcd.yml up -d
[+] Building 0.0s (0/0)
[+] Running 4/4
✓Network cassandra-demo_default Created 0.2s
✓Container etcd0 Started 0.4s
✓Container etcd1 Started 0.4s
✓Container etcd2 Started 0.3s
cc@osfn-cc:~/cassandra-demo$
```

## # Set a convenience variable for etcdctl

export

ENDPOINTS="http://etcd0:2379,http://etcd1:2379,http://etcd2:2379"

## # Health check

docker exec etcd1 /usr/local/bin/etcdctl --endpoints="\$ENDPOINTS" endpoint status -w table

```
cc@osfn-cc:~/cassandra-demo$ export ENDPOINTS="http://etcd0:2379,http://etcd1:2379,http://etcd2:2379"
cc@osfn-cc:~/cassandra-demo$ docker exec etcd1 /usr/local/bin/etcdctl --endpoints="$ENDPOINTS" endpoint status -w table
```

ENDPOINT	RAFT TERM	RAFT INDEX	ID	RAFT APPLIED INDEX	VERSION	STORAGE VERSION	DB SIZE	IN USE	PERCENTAGE NOT IN USE	QUOTA	IS LEADER	IS LEARNER
http://etcd0:2379	24	135578	cf1d15c5d194b5c9	135578	3.6.2	3.6.0	637 MB	628 MB	false	2%	0 B	false
http://etcd1:2379	24	135578	ade526d28b1f92f7	135578	3.6.2	3.6.0	637 MB	628 MB	false	2%	0 B	true
http://etcd2:2379	24	135578	d282ac2ce60c1ce	135578	3.6.2	3.6.0	637 MB	628 MB	false	2%	0 B	false

```
cc@osfn-cc:~/cassandra-demo$
```

## 1) Topology & High-Level Flow

- **Cluster:** 3× etcd containers (etcd0, etcd1, etcd2) via Docker Compose. Clients/benchmarks run from etcd1 container using etcdctl.
- **Fault tool:** Charybdefs running on the **host** as a FUSE filesystem with a Thrift mgmt service (listens on :9090).
- **Mounting strategy:** Only the **leader's** data directory is placed on the Charybdefs mount. In our runs we pin the leader to **etcd2** and mount its data dir from **/mnt/slowfs/etcd2** (FUSE) backed by **/data/raw/etcd2** (real disk).
- **Results:** For every run we write CSVs to   
~/cassandra-demo/results/<TIMESTAMP>\_<label>/:
  - per\_op\_latency.csv (per operation lat in seconds)
  - latency\_per\_sec.csv (avg p50/p95/p99 per second)
  - throughput\_per\_sec.csv (ops/sec time series)

### Run flow (per scenario):

1. Verify cluster health → 2) Force leadership to etcd2 → 3) Clear faults → 4) (Optionally) inject fs delay on the leader's path → 5) Run mini-benchmark and capture metrics → 6) Clear faults.

## 2) One-Time Prerequisites

- Docker + Docker Compose installed; python3 available.

- Directories on host for Charybdefs mounts and real backing storage.
- ~/cassandra-demo workspace with docker-compose-etcd.yml (3-node etcd) and our helper scripts.
- ~/charybdefs built from source (provides the FUSE server that speaks Thrift on :9090 in our build).

### 3) Host Setup & Mounting (exact commands we used)

Use these when (re)starting or cleaning the FUSE mount and etcd cluster.

# Stop the etcd cluster (ok if not running)

cd ~/cassandra-demo

docker compose -f docker-compose-etcd.yml down || true

# Stop any previous Charybdefs and unmount the FUSE path

sudo kill charybdefs || true

sudo umount -l /mnt/slowfs/etcd2 2>/dev/null \

|| fusermount -u /mnt/slowfs/etcd2 2>/dev/null || true

```
cc@osfn-cc:~/charybdefs$ sudo kill charybdefs || true
/mnt/slowfs/etcd2 2>/dev/null || true
cc@osfn-cc:~/charybdefs$ sudo umount -l /mnt/slowfs/etcd2 2>/dev/null \
> || fusermount -u /mnt/slowfs/etcd2 2>/dev/null || true
cc@osfn-cc:~/charybdefs$ |
```

# Prepare backing and mount directories

sudo mkdir -p /data/raw/etcd2 /mnt/slowfs/etcd2

sudo chown -R root:root /data/raw/etcd2 /mnt/slowfs/etcd2

sudo chmod 700 /data/raw/etcd2 /mnt/slowfs/etcd2

```
cc@osfn-cc:~/charybdefs$ sudo mkdir -p /data/raw/etcd2 /mnt/slowfs/etcd2
cc@osfn-cc:~/charybdefs$ sudo chown -R root:root /data/raw/etcd2 /mnt/slowfs/etcd2
hmod 700 /data/raw/etcd2 /mnt/slowfs/etcd2
cc@osfn-cc:~/charybdefs$ sudo chmod 700 /data/raw/etcd2 /mnt/slowfs/etcd2
cc@osfn-cc:~/charybdefs$ |
```

# Start Charybdefs (FUSE). It listens on 9090; backs /mnt/slowfs/etcd2 to /data/raw/etcd2

nohup sudo "\$HOME/charybdefs/charybdefs" /mnt/slowfs/etcd2 \

-omodels=subdir,subdir=/data/raw/etcd2 -oallow\_other,nonempty \

```
> /tmp/charybdefs.log 2>&1 &
```

```
# give FUSE a moment
```

```
sleep 2
```

```
# Start etcd cluster
```

```
cd ~/cassandra-demo
```

```
# Ensure docker-compose-etcd.yml binds etcd2's data dir to /mnt/slowfs/etcd2
```

```
# (etcd0/1 may use fast local dirs)
```

```
docker compose -f docker-compose-etcd.yml up -d
```

```
# For convenience during runs
```

```
export ENDPOINTS="http://etcd0:2379,http://etcd1:2379,http://etcd2:2379"
```

**Checks:**

- `tail -n 100 /tmp/charybdefs.log` (FUSE mounted, Thrift server ready)
- `docker compose ps` (all etcd containers up)
- `docker exec etcd1 etcdctl --endpoints="$ENDPOINTS" endpoint status -w table`

## 4) Files We Added / Modified

### 4.1 `charyb_fault.py` (Thrift client wrapper)

Small Python helper living in `~/cassandra-demo/charyb_fault.py` that talks to Charybdefs mgmt (host:port) and exposes **two** commands used by our shell orchestrator:

- `clear-all` — clears any active fs faults
- `set-delay --path <mount_path> --delay-us <µs>` — injects a WAL/write/fsync delay on the given mount path

In our environment this script imports your generated Thrift stubs from `~/charybdefs/gen-py`. If you relocate stubs, set `PYTHONPATH` accordingly. We used these exact invocations from the orchestrator below.

## 4.2 run\_etcd\_fsdelay.sh (orchestrator; final working version)

Save as ~/cassandra-demo/run\_etcd\_fsdelay.sh and chmod +x it.

```
D: > UChicago Internship > Reproducing-Filesystem-Delay-One-Size-Fits-None-NSDI25 > run_etcd_fsdelay.sh
1  #!/usr/bin/env bash
2  # run_etcd_fsdelay.sh
3  # Baseline vs WAL delay on the LEADER node's FUSE-mounted data dir.
4
5  set -euo pipefail
6
7  # ===== CONFIG DEFAULTS =====
8  OPS="{OPS:-200}" # number of put ops
9  MODE="{1:-baseline}" # baseline | delay
10 LEADER_TARGET="{LEADER_TARGET:-etcd2}" # we want etcd2 as leader (the slow WAL)
11 WAL_DELAY_US="{WAL_DELAY_US:-0}" # e.g. 0, 100000, 300000, 750000
12 RESULTS_DIR="{RESULTS_DIR:-results}"
13 OUT_PREFIX="{OUT_PREFIX:-etcd_fsdelay}"
14 ETCDCCTL="{ETCDCCTL:-/usr/local/bin/etcdctl}"
15 ETCDC_CONTAINER="{ETCDC_CONTAINER:-etcd0}" # container to exec etcdctl from
16 ENDPOINTS="{ENDPOINTS:-http://etcd0:2379,http://etcd1:2379,http://etcd2:2379}"
17 PYTHON="{PYTHON:-python3}"
18 CHARYB="{CHARYB:-./charyb_fault.py}"
19
20 # charybdefs control (your daemon listens here)
21 CHARYB_HOST="{CHARYB_HOST:-127.0.0.1}"
22 CHARYB_PORT="{CHARYB_PORT:-9090}"
23
24 # Methods and regex to hit etcd WAL on the slow node (matches your working setup)
25 WAL_METHODS="{WAL_METHODS:-open,create,write,write_buf,fsync,fdatasync,fsyncdir,flush}"
26 WAL_REGEX="{WAL_REGEX:-(^|.*/)member/wal/.*)" # path as charybdefs sees it (starts with /)
27
28 # Optional: verify the delay is actually seen on the WAL path (host-side)
29 VERIFY_DELAY="{VERIFY_DELAY:-1}"
30
31 mkdir -p "$RESULTS_DIR"
32
33 echo "Target leader : ${LEADER_TARGET:-'(current)'}"
34 echo "Mode : $MODE"
35 echo "Ops : $OPS"
36 [[ "$MODE" == "delay" ]] && echo "WAL delay (us) : $WAL_DELAY_US"
37
38 # ===== END =====
```

```
C:\Users\jeezx>scp -i D:/key/yizzz-mj-trace.pem D:/Repo/code_sweepdelay/run_etcd_fsdelay.sh cc@192.5.86.216:/home/cc/cassandra-demo/
run_etcd_fsdelay.sh
100% 4802 15.2KB/s 00:00
C:\Users\jeezx>
```

### Why these changes matter

- **Leader forcing:** etcdctl move-leader expects **decimal** member IDs; we parse JSON and pass the correct value to avoid the earlier hex/format errors.
- **Leader detection:** We treat the endpoint whose Status.leader == Status.header.member\_id as the leader (stable in practice).
- **No jq dep:** We parse JSON with tiny Python one-liners so the script runs everywhere.
- **Fault mgmt:** We do **not** patch Charybdefs; we control it via Thrift using our Python wrapper and set PYTHONPATH to your gen-py stubs to prevent ModuleNotFoundError incidents.
- **Metrics:** We record per-op latency and derive **per-second** latency percentiles and throughput to match the figures you plotted (per\_op\_latency.csv, latency\_per\_sec.csv, throughput\_per\_sec.csv).

## 5) Running the Experiments

### 5.1 Baseline

cd ~/cassandra-demo

LEADER\_TARGET=etcd2 OPS=200 ./run\_etcd\_fsdelay.sh baseline

This checks health, pins leadership to etcd2, clears faults, runs the benchmark, and writes CSVs under results/<TIMESTAMP>\_etcd\_fsdelay\_baseline\_etcd2\_0us/.

### 5.2 Filesystem delay severities (1–100 ms)

Our severity grid in **microseconds**: 1000 5000 10000 25000 50000  
75000 100000 → **1, 5, 10, 25, 50, 75, 100 ms**

Run individually, e.g. **5 ms**:

```
Command Prompt x Command Prompt x Command Prompt x Command Prompt x Command Prompt x cc@osfn-cc-~/cass x + v - o x
cc@osfn-cc-~/cassandra-demo$ LEADER_TARGET=etcd2 OPS=1000 MODE=delay WAL_DELAY_US=1000 VERIFY_DELAY=1 ./run_etcd_fsdelay.sh delay
Target leader : etcd2
Mode : delay
Ops : 1000
WAL delay (us) : 1000
Endpoints : http://etcd0:2379,http://etcd1:2379,http://etcd2:2379

== Cluster ==

+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
| ENDPOINT | ID | VERSION | STORAGE | DB SIZE | IN USE | PERCENTAGE NOT IN USE | QUOTA | IS LEADER | IS LEARNER | RAFT TERM | RAFT INDEX | RAFT APPLIED INDEX |
| ERRORS | DOWNGRADE TARGET | DOWNGRADE | | | | | | | | | | |
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
| http://etcd0:2379 | cf1d15c5d194b5c9 | 3.6.2 | false | 3.6.0 | 637 MB | 627 MB | 2% | 0 B | false | false | 22 | 128574 | 128574 |
| http://etcd1:2379 | ade526d28b1f92f7 | 3.6.2 | false | 3.6.0 | 637 MB | 627 MB | 2% | 0 B | false | false | 22 | 128574 | 128574 |
| http://etcd2:2379 | d282ac2ce608c1ce | 3.6.2 | false | 3.6.0 | 637 MB | 627 MB | 2% | 0 B | true | false | 22 | 128574 | 128574 |
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+

[info] Forcing leader to: etcd2
[info] current leader: etcd2
Leader OK: etcd2

+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
| ENDPOINT | ID | VERSION | STORAGE | DB SIZE | IN USE | PERCENTAGE NOT IN USE | QUOTA | IS LEADER | IS LEARNER | RAFT TERM | RAFT INDEX | RAFT APPLIED INDEX |
| ERRORS | DOWNGRADE TARGET | DOWNGRADE | | | | | | | | | | |
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
| http://etcd0:2379 | cf1d15c5d194b5c9 | 3.6.2 | false | 3.6.0 | 637 MB | 627 MB | 2% | 0 B | false | false | 22 | 128574 | 128574 |
| http://etcd1:2379 | ade526d28b1f92f7 | 3.6.2 | false | 3.6.0 | 637 MB | 627 MB | 2% | 0 B | false | false | 22 | 128574 | 128574 |
| http://etcd2:2379 | d282ac2ce608c1ce | 3.6.2 | false | 3.6.0 | 637 MB | 627 MB | 2% | 0 B | true | false | 22 | 128574 | 128574 |
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+

CharybdeFS: delay on WAL
CharybdeFS: injected delay 1000 us methods=['open', 'create', 'write', 'write_buf', 'fsync', 'fdatasync', 'fsyncdir', 'flush'] prob=1000/1000 regex='(^.*)/member/wal/.+'
[Verify] host WAL fsync elapsed ~0.005s (inj#0.001s)
>>> Workload: 1000 x PUT to etcd2 (http://etcd2:2379)

Summary:
ok=1000 fail=0 wall=84.700s throughput=11.81 ops/s
p50=0.08s p95=0.10s p99=0.11s
Saved per-op latency CSV : results/20250810_050916_etcd_fsdelay_delay_etcd2_1000us/per_op_latency.csv
```

LEADER\_TARGET=etcd2 OPS=200 WAL\_DELAY\_US=5000 ./run\_etcd\_fsdelay.sh delay

Or sweep them all:

for us in 1000 5000 10000 25000 50000 75000 100000; do

LEADER\_TARGET=etcd2 OPS=200 WAL\_DELAY\_US=\$us ./run\_etcd\_fsdelay.sh delay

done

Each run produces a new timestamped directory, e.g.:

results/

20250809\_222805\_etcd\_fsdelay\_baseline\_etcd2\_0us/

per\_op\_latency.csv latency\_per\_sec.csv throughput\_per\_sec.csv

20250809\_225811\_etcd\_fsdelay\_delay\_etcd2\_100us/

20250810\_051124\_etcd\_fsdelay\_delay\_etcd2\_5000us/

```
Command Prompt x Command Prompt x Command Prompt x Command Prompt x Command Prompt x cc@osfn-cc: ~/cassu x + x x x x x x
cc@osfn-cc:~/cassandra-demo$ LEADER_TARGET=etcd2 OPS=1000 MODE=delay WAL_DELAY_US=5000 VERIFY_DELAY=1 ./run_etcd_fsdelay.sh delay
Target leader : etcd2
Mode : delay
Ops : 1000
WAL delay (us) : 5000
Endpoints : http://etcd0:2379, http://etcd1:2379, http://etcd2:2379

== Cluster ==

+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
| ENDPOINT | ID | VERSION | STORAGE | VERSION | DB SIZE | IN USE | PERCENTAGE NOT IN USE | QUOTA | IS LEADER | IS LEARNER | RAFT TERM | RAFT INDEX | RAFT APPLIED INDEX |
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
| http://etcd0:2379 | cf1d15c5d194b5c9 | 3.6.2 | false | 3.6.0 | 637 MB | 627 MB | 2% | 0 B | false | false | 22 | 129574 | 129574 |
| http://etcd1:2379 | ade526d28b1f92f7 | 3.6.2 | false | 3.6.0 | 637 MB | 627 MB | 2% | 0 B | false | false | 22 | 129574 | 129574 |
| http://etcd2:2379 | d282ac2ce608c1ce | 3.6.2 | false | 3.6.0 | 637 MB | 627 MB | 2% | 0 B | true | false | 22 | 129574 | 129574 |
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+

[info] Forcing leader to: etcd2
[info] current leader: etcd2
Leader OK: etcd2

+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
| ENDPOINT | ID | VERSION | STORAGE | VERSION | DB SIZE | IN USE | PERCENTAGE NOT IN USE | QUOTA | IS LEADER | IS LEARNER | RAFT TERM | RAFT INDEX | RAFT APPLIED INDEX |
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
| http://etcd0:2379 | cf1d15c5d194b5c9 | 3.6.2 | false | 3.6.0 | 637 MB | 627 MB | 2% | 0 B | false | false | 22 | 129574 | 129574 |
| http://etcd1:2379 | ade526d28b1f92f7 | 3.6.2 | false | 3.6.0 | 637 MB | 627 MB | 2% | 0 B | false | false | 22 | 129574 | 129574 |
| http://etcd2:2379 | d282ac2ce608c1ce | 3.6.2 | false | 3.6.0 | 637 MB | 627 MB | 2% | 0 B | true | false | 22 | 129574 | 129574 |
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+

Charybdefs: delay on WAL
Charybdefs: injected delay 5000 us methods=['open', 'create', 'write', 'write_buf', 'fsync', 'fdatasync', 'fsyncdir', 'flush'] prob=1000/1000 regex='(^|.)member/wal/.+'
[verify] host WAL fsync elapsed ~0.021s (inj=0.005s)
>>> Workload: 1000 x PUT to etcd2 (http://etcd2:2379)

Summary:
ok=1000 fail=0 wall=99.146s throughput=10.09 ops/s
p50=0.09s p95=0.11s p99=0.12s
Saved per-op latency CSV : results/20250810_051124_etcd_fsdelay_delay_etcd2_5000us/per_op_latency.csv

cc@osfn-cc:~$ cd cassandra-demo
cc@osfn-cc:~/cassandra-demo$ cd results/20250810_050916_etcd_fsdelay_delay_etcd2_1000us/
cc@osfn-cc:~/cassandra-demo/results/20250810_050916_etcd_fsdelay_delay_etcd2_1000us$ ls
latency_per_sec.csv per_op_latency.csv throughput_per_sec.csv
cc@osfn-cc:~/cassandra-demo/results/20250810_050916_etcd_fsdelay_delay_etcd2_1000us$ cd ../../
cc@osfn-cc:~/cassandra-demo$ cd results/20250810_051124_etcd_fsdelay_delay_etcd2_5000us
cc@osfn-cc:~/cassandra-demo/results/20250810_051124_etcd_fsdelay_delay_etcd2_5000us$ ls
latency_per_sec.csv per_op_latency.csv throughput_per_sec.csv
cc@osfn-cc:~/cassandra-demo/results/20250810_051124_etcd_fsdelay_delay_etcd2_5000us$ cd ../../
cc@osfn-cc:~/cassandra-demo$ cd results/20250810_051315_etcd_fsdelay_delay_etcd2_10000us
cc@osfn-cc:~/cassandra-demo/results/20250810_051315_etcd_fsdelay_delay_etcd2_10000us$ ls
latency_per_sec.csv per_op_latency.csv throughput_per_sec.csv
cc@osfn-cc:~/cassandra-demo/results/20250810_051315_etcd_fsdelay_delay_etcd2_10000us$ cd ../../
cc@osfn-cc:~/cassandra-demo$ cd results/20250810_051546_etcd_fsdelay_delay_etcd2_25000us
cc@osfn-cc:~/cassandra-demo/results/20250810_051546_etcd_fsdelay_delay_etcd2_25000us$ ls
latency_per_sec.csv per_op_latency.csv throughput_per_sec.csv
cc@osfn-cc:~/cassandra-demo/results/20250810_051546_etcd_fsdelay_delay_etcd2_25000us$ cd ../../
cc@osfn-cc:~/cassandra-demo$ cd results/20250810_051818_etcd_fsdelay_delay_etcd2_50000us
cc@osfn-cc:~/cassandra-demo/results/20250810_051818_etcd_fsdelay_delay_etcd2_50000us$ ls
latency_per_sec.csv per_op_latency.csv throughput_per_sec.csv
cc@osfn-cc:~/cassandra-demo/results/20250810_051818_etcd_fsdelay_delay_etcd2_50000us$ cd ../../
cc@osfn-cc:~/cassandra-demo$ cd results/20250810_052230_etcd_fsdelay_delay_etcd2_75000us
cc@osfn-cc:~/cassandra-demo/results/20250810_052230_etcd_fsdelay_delay_etcd2_75000us$ ls
latency_per_sec.csv per_op_latency.csv throughput_per_sec.csv
cc@osfn-cc:~/cassandra-demo/results/20250810_052230_etcd_fsdelay_delay_etcd2_75000us$ cd ../../
cc@osfn-cc:~/cassandra-demo$ cd results/20250810_052751_etcd_fsdelay_delay_etcd2_100000us
cc@osfn-cc:~/cassandra-demo/results/20250810_052751_etcd_fsdelay_delay_etcd2_100000us$ ls
latency_per_sec.csv per_op_latency.csv throughput_per_sec.csv
```

## 6) Reading & Using the Outputs



- **Per-op latency:** per\_op\_latency.csv with op, seconds. Useful for line/CDF plots across delays.
- **Latency per second:** latency\_per\_sec.csv with second, p50\_ms, p95\_ms, p99\_ms. Good for seeing steady-state vs. spikes.
- **Throughput per second:** throughput\_per\_sec.csv with second, ops. Good for identifying “collapse zones”.

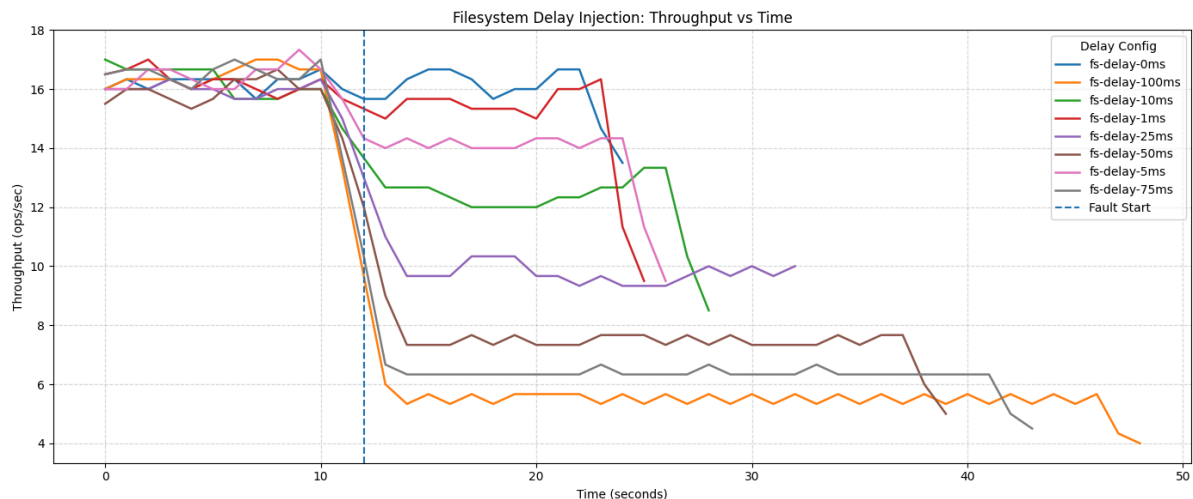
Our plot scripts (line charts, CDFs, and latency vs. throughput overlays) consume these CSVs directly. The figures you generated for the paper reproduction came from these files.

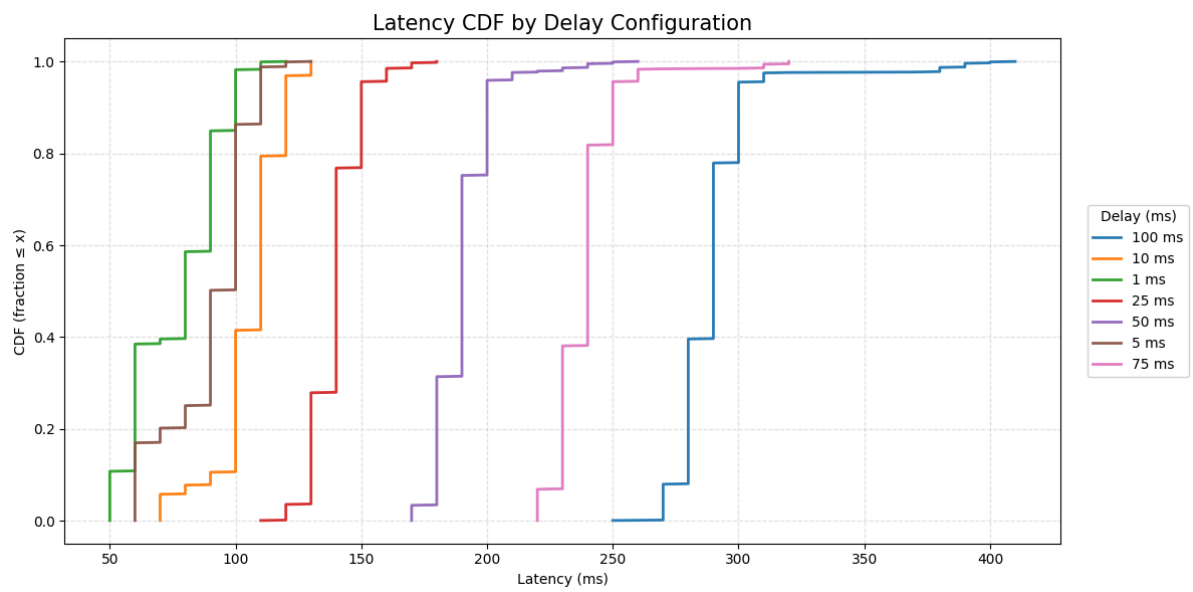
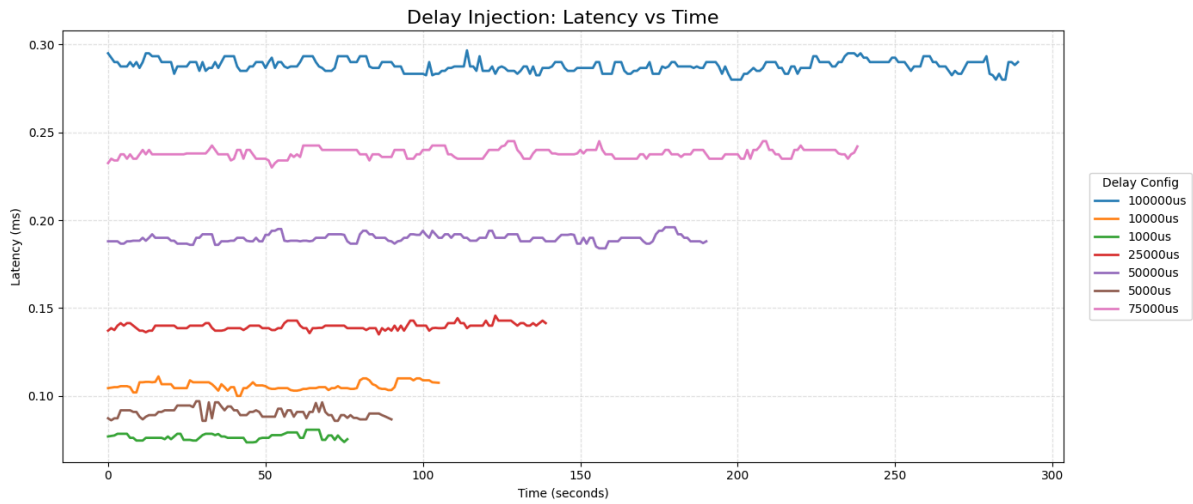
```

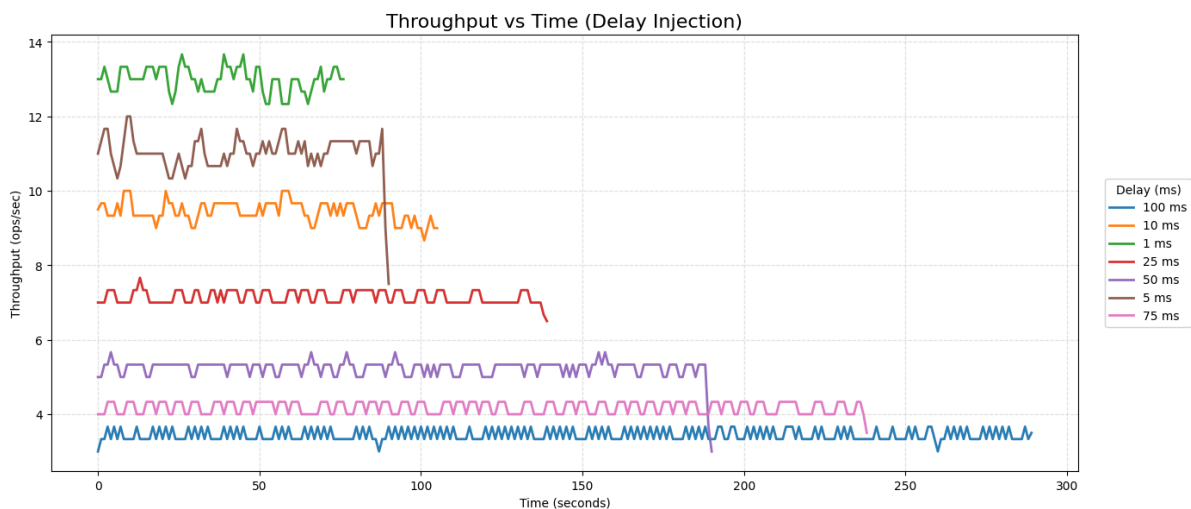
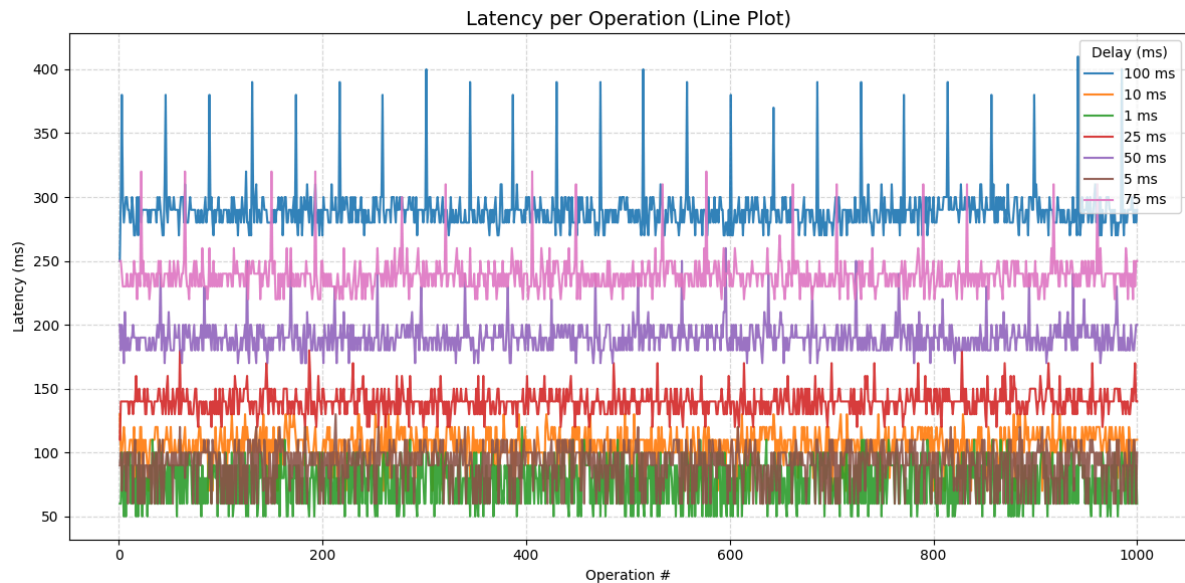
ETCDPLOT2
cdf_&_plot.py
cdf.py
latency_per_sec_1...
latency_per_sec_5...
latency_per_sec_1...
latency_per_sec_2...
latency_per_sec_5...
latency_per_sec_7...
latency_per_sec_1...
latency_vs_time.py
per_op_latency_1...
per_op_latency_5...
per_op_latency_1...
per_op_latency_2...
per_op_latency_5...
per_op_latency_7...
per_op_latency_1...
throughput_per_s...
throughput_per_s...
throughput_per_s...

cdf_&_plot.py > ...
1 import pandas as pd
2 import matplotlib.pyplot as plt
3 import glob, os, re
4 import numpy as np
5
6
7 FILE_PATTERN = "per_op_latency_*.csv" # match this into your file name
8 FIGSIZE = (12, 6)
9
10
11 def to_ms_label(text: str) -> str:
12     """
13     Ambil pola '<angka><unit>' dari text dan konversi ke 'X ms'.
14     Unit didukung: us/μs, ms, s. Jika tidak ketemu, kembalikan text asli.
15     """
16     m = re.search(r'(?i)(\d+(?:\.\d+)?)(s|ms|us|μs)\b', text)
17     if not m:
18         return text
19     val = float(m.group(1))
20     unit = m.group(2).lower()
21
22     if unit in ('μs', 'us'):
23         ms = val / 1000.0

```







## 7) Validation & Sanity Checks

- **Leader really on etcd2?**

`docker exec etcd1 etcdctl --endpoints="$ENDPOINTS" endpoint status -w table`

### Charybdefs active?

`mount | grep slowfs`

`tail -n 50 /tmp/charybdefs.log`

- `ss -lntp | grep :9090` # mgmt listening
- **After a delay run, latency shifts** should be visible in `per_op_latency.csv` and `lower throughput_per_sec.csv` relative to baseline. If not, re-check that etcd2's data dir in Compose is bound to `/mnt/slowfs/etcd2` (not a fast path) and that leadership is actually on etcd2 **during** the run.

## 8) Troubleshooting (issues we actually hit)

- **move-leader complaining about invalid ID** → We were passing hex; script now parses decimal.
- **ModuleNotFoundError for Thrift stubs** → Ensure `PYTHONPATH=$HOME/charybdefs/gen-py` before calling the Python client.
- **FUSE unmount busy** → Use `umount -l ... || fusermount -u ...` as shown; then restart Charybdefs.
- **No noticeable impact across small delays** → Verify the leader is on the **slowfs-backed** node during the **entire** run; confirm delay values are **microseconds** not milliseconds; verify ops load (OPS) is enough to surface the effect.

## 9) Quick Command Recap

### # One-time (or after a reset)

```
cd ~/cassandra-demo
```

```
docker compose -f docker-compose-etcd.yml down || true
```

```
sudo pkill charybdefs || true
```

```
sudo umount -l /mnt/slowfs/etcd2 2>/dev/null || fusermount -u  
/mnt/slowfs/etcd2 2>/dev/null || true
```

```
sudo mkdir -p /data/raw/etcd2 /mnt/slowfs/etcd2
```

```
sudo chown -R root:root /data/raw/etcd2 /mnt/slowfs/etcd2
```

```
sudo chmod 700 /data/raw/etcd2 /mnt/slowfs/etcd2
```

```
nohup sudo "$HOME/charybdefs/charybdefs" /mnt/slowfs/etcd2 \  
-omodels=subdir,subdir=/data/raw/etcd2 -oallow_other,nonempty \  
> /tmp/charybdefs.log 2>&1 &
```

```
docker compose -f docker-compose-etcd.yml up -d
```

```
export ENDPOINTS="http://etcd0:2379,http://etcd1:2379,http://etcd2:2379"
```

### # Baseline

```
LEADER_TARGET=etcd2 OPS=200 ./run_etcd_fsdelay.sh baseline
```

### # Delays (1–100 ms)

```
for us in 1000 5000 10000 25000 50000 75000 100000; do
```

```
LEADER_TARGET=etcd2 OPS=200 WAL_DELAY_US=$us  
./run_etcd_fsdelay.sh delay
```

done

## 10) Notes Mapping Back to the Paper

- **Fault type:** filesystem delay on the leader's log/write path (via Charybdefs), matching the paper's fs-slow methodology.
- **Severity grid:** 1–100 ms (we used 1, 5, 10, 25, 50, 75, 100 ms).
- **Metric emphasis:** throughput degradation alongside p95/p99 latency (tail latency alone can be misleading for varying durations).
- **Pipeline:** init → warm-up (if needed) → inject → measure → clear (our runs keep them short; add warm-up as desired).

## Appendix A — Environment Variables (for quick tuning)

Variable	Meaning	Default
ENDPOINTS	etcd endpoints CSV	http://etcd0:2379,http://etcd1:2379,http://etcd2:2379
LEADER_TARGET	node to hold leadership	etcd2
OPS	ops per run (if DURATION_SEC=0)	200
DURATION_SEC	time-based run if >0	0
WAL_DELAY_US	injected fs delay in microseconds	0

OUTROOT	results root dir	\$PWD/results
MOUNT_PREFIX	Charybdefs mount root on host	/mnt/slowfs
CHARY_HOST/CHAR Y_PORT	mgmt endpoint	127.0.0.1 / 9090
CHARY_CLI	path to our Thrift Python wrapper	\$PWD/charyb_fault.py