

## Scaling Validation Report v3

### Key Findings:

- 8 replicas held 56.14 req/s at 400 users (0 failures) after ramp/timeout tuning; +13.
- The earlier debug run logged 25.95% failures; rerun now completes 10,111 calls with 0 fails.
- CPU snapshots show 95-154% utilisation per inference container (avg ~130%), indicating some contention.
- 200-user / 8-replica sweep remains warm-up limited (20.44 req/s, 0 fails) and is treated as a failure.

### Artifacts:

- reports/scaling\_validation\_v3/scaling\_metrics.csv (and .json)
- reports/scaling\_validation\_v3/scaling\_metadata.csv (and .json)
- reports/scaling\_validation\_v3/rps\_vs\_replicas.[png|pdf] (plus median/p95/failure charts)
- locust/results/multi\_replica\_v1/replicas8\_u400/ (clean rerun + CPU stats)
- locust/results/multi\_replica\_v1/replicas8\_u400\_obsolete/ (superseded artifacts)

### Next Steps:

- Optional: capture the 800-user / 8-replica lane with -r 30 and refresh visuals.