# Analysis and Design of Scalable Programming Auto grader

MTP Final Thesis Presentation

Guide:
Prof Varsha Apte

Submitted by:
Chaitanya Varma
203050026

# Table of Contents

# Need for Scalability in Programming auto-graders

# Typical programming auto grader architecture



API Layer     Service Layer     DB Layer

User

Request

Response

Web Server

Compilation

Execution

Comparing outputs

processing time depends upon the user code and cannot be profiled

**Scalability**: Ability of the system to scale its performance with increase in resources

At peak usage, system need to scale its performance

Not possible to estimate resource requirement in prior

# Background about Evalpro

# Evalpro Design



Processes evaluation tasks asynchronously

To avoid

Timeouts

Reloads

Load inflation

Internal Web server

Load balancer

Users

Evaluation request

HAProxy

WSGI

Compilation

Execution

Postgres Db

Celery task queue

Compare outputs

Nginx

Redis

Serves static content

Broker, enques evaluation tasks to celery queue

Indicates container

Managed by Docker

File system

Media Folder
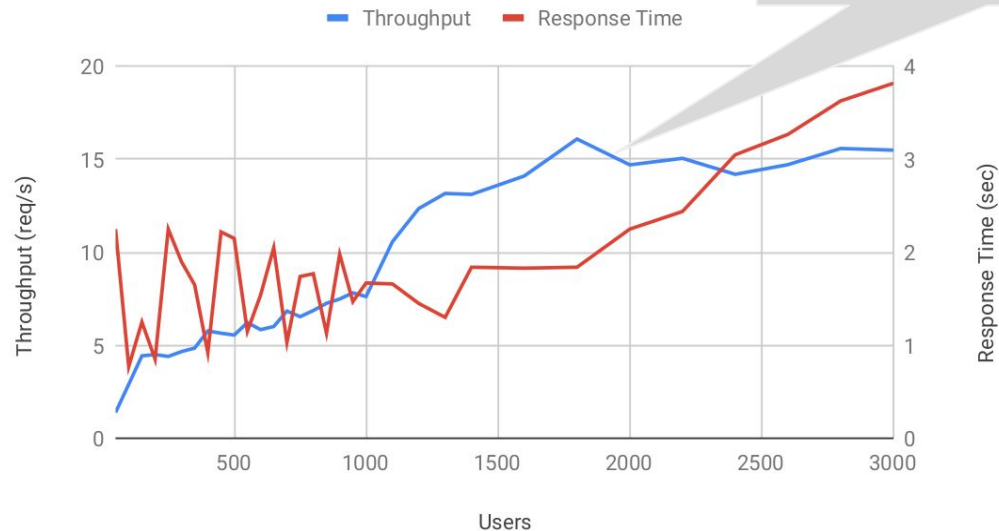
User uploaded files are stored

# Prior work of Evalpro scalability

# Prior work (MTP 2019) and its limitations

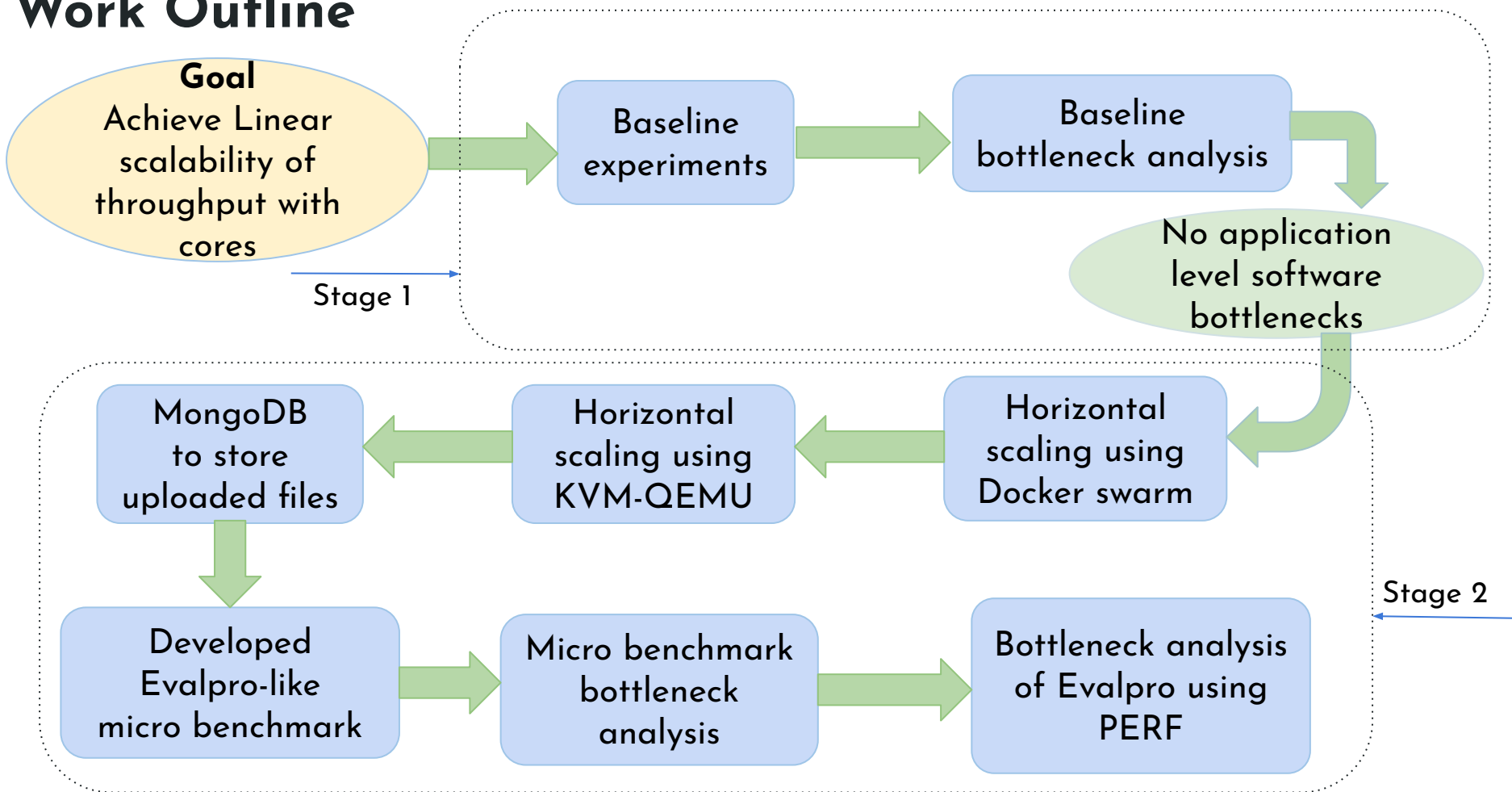40 CPU core server

Throughput and Response Time



- Throughput flattened
- CPU utilization not exceeded 5%

Scalability issue

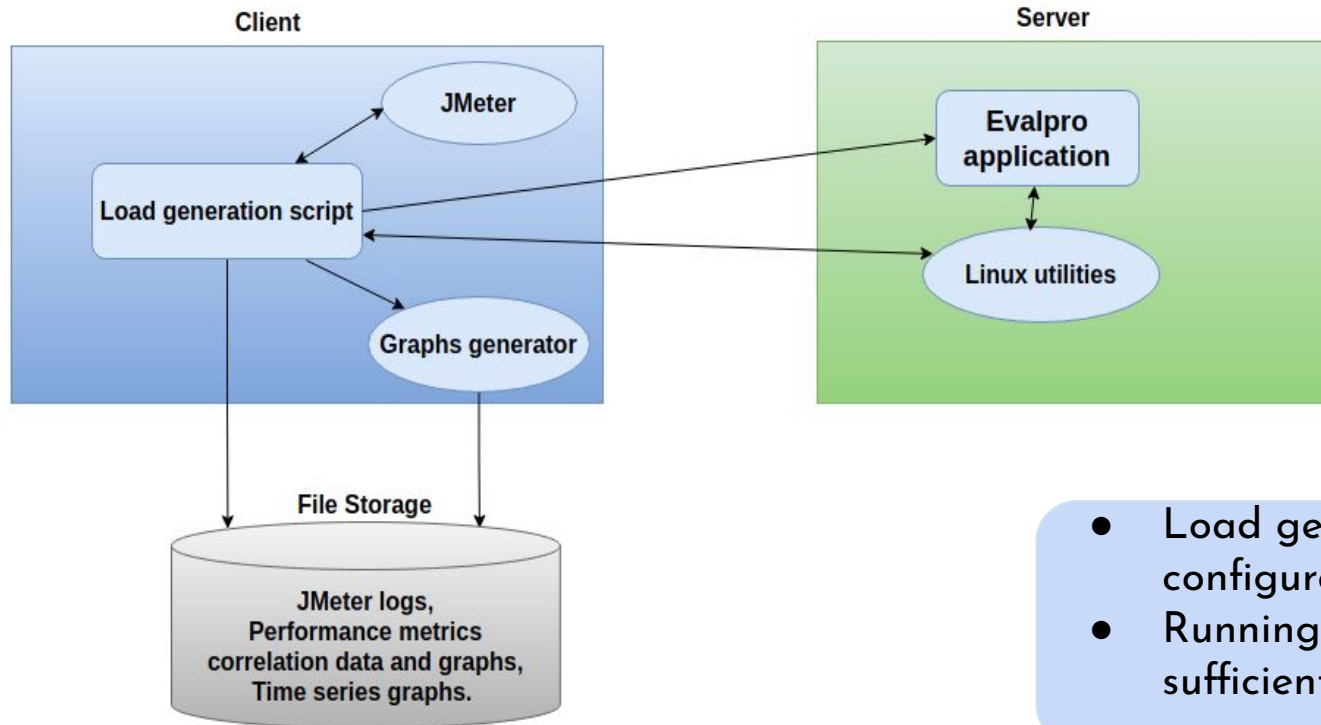Increasing  WSGI+Celery replicas increased the throughput and CPU utilization

- Reason for improvement in scalability not justified
- Scalability issue with single replica, configurations affecting the performance not studied
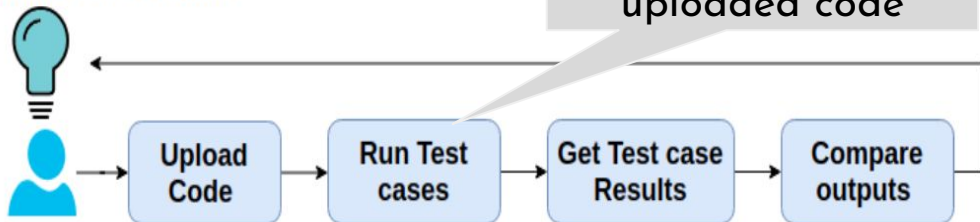
*"The data and figures in this slide were taken from Anshul MTP thesis 2019"*

# Work Outline

# MTP stage 1

# Load generation and performance measurement Infrastructure

# Baseline Experiment setup

Includes compilation and execution of uploaded code

16 CPU cores

Evalpro instance

**Think time= 15 seconds**

Upload Code → Run Test cases → Get Test case Results → Compare outputs

Realistic user session

HAProxy ↔ WSGI

Celery task queue

Postgres DB

Nginx

Redis

Media Folder
File system

☐ *Indicates container*

⬭ *Managed by Docker*

| Type | CPU | Cores | Memory | L3 cache | L2 cache | L1 cache |
|------|-----|-------|--------|----------|----------|----------|
| Server | Intel$^R$ Xeon$^R$ CPU E5-2650 v2 @ 2.60GHz | 16 | 16GB | 20MB | 256KB | 32KB |
| Client | AMD Opteron$^{TM}$ Processor 6212 | 16 | 16GB | 6MB | 2MB | 64KB |

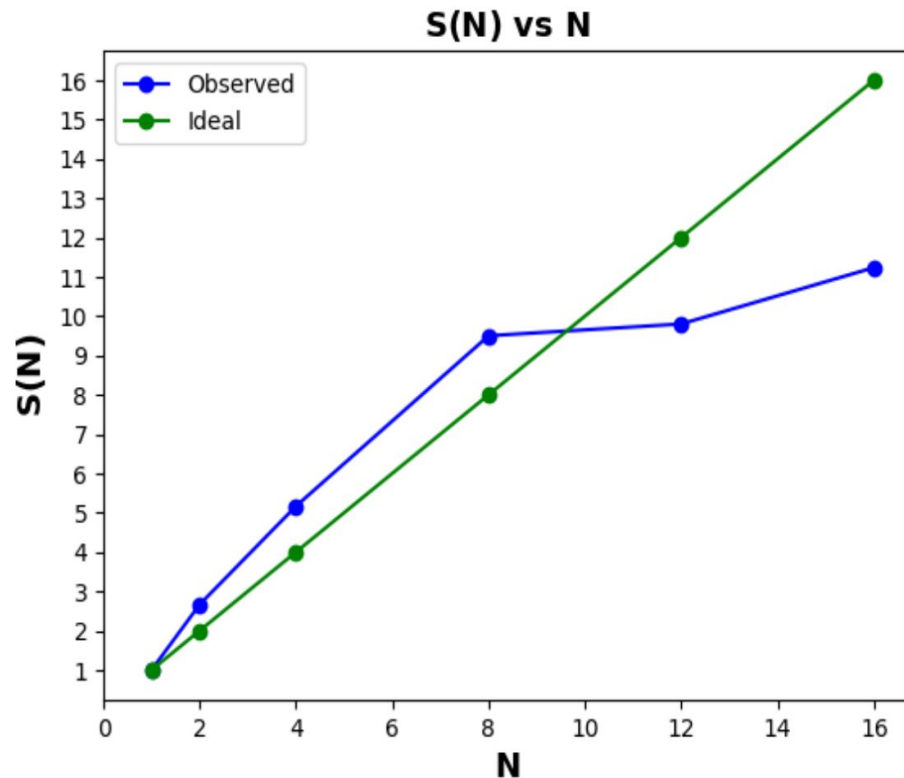Hardware specifications for baseline experiments

# Baseline 16

To avoid obvious bottlenecks, Postgres DB connections set to 10000 Celery threads set to number of cores

$Throughput_{max}(N)$ : Maximum throughput with N CPU cores
$S(N)$ : Throughput scalability factor with N CPU cores

$S(N) = Throughput_{max}(N) \div Throughput_{max}(1)$
Ideal $S(N) = N$



S(N) vs N

| $Throughput_{max}(1)$ | $Throughput_{max}(16)$ | Ideal $Throughput_{max}(16)$ | $S(16)$ |
|---|---|---|---|
| 0.3 req/sec | 3.37 req/sec | 4.8 req/sec | 11.23 |

Linear scaling of throughput not achieved

# Baseline bottleneck analysis



**16 CPU cores**

Increased gunicorn workers

Setting celery CPU affinities

Increased prefetch count

Increased celery threads

Made celery queue transient

Disabled logging

Increased open connections to redis

HAProxy

WSGI

Postgres

Celery task queue

Nginx

Redis

*Indicates container*

*Managed by Docker*

None of these experiments improved throughput scalability

No application level software bottlenecks

All cores 90-100 % utilized

# MTP Stage-2

# Horizontal scaling with Docker swarm

Load balances request load across multiple WSGI + Celery replicas



Indicates container

Managed by Docker

Docker swarm didn't improve throughput scalability

|  | Baseline | Using Docker swarm |
|---|---|---|
| $Throughput_{max}(1)$ (req/sec) | 0.3 | 0.3 |
| $Throughput_{max}(16)$ (req/sec) | 3.37 | 3.4 |
| Ideal $Throughput_{max}(16)$ (req/sec) | 4.8 | 4.8 |
| $S(16)$ | 11.23 | 11.3 |

# Horizontal scaling with KVM-QEMU

# Completely Isolated VM setup



|  | Baseline | Isolated VM setup |
|---|---|---|
| $Throughput_{max}(1)$ (req/sec) | 0.3 | 0.3 |
| $Throughput_{max}(16)$ (req/sec) | 3.37 | 4.45 |
| Ideal $Throughput_{max}(16)$ (req/sec) | 4.8 | 4.8 |
| $S(16)$ | 11.23 | 14.5 |

Isolated setup improved throughput scalability

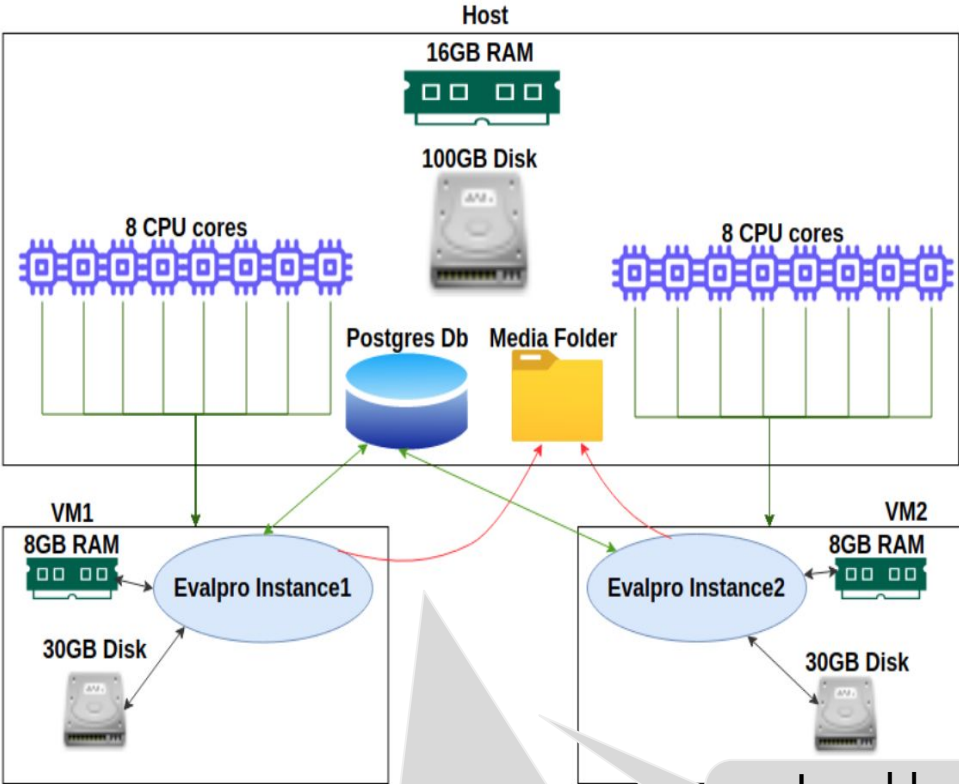# Completely Isolated VM setup limitations



Each VM need to maintain separate user related data

Load balancer need have user to VM mapping

Migrating user data with new VM addition is very complex and not feasible.

Not the desired solution

# User data and files sharing VM setup



|  | Baseline | Data and file sharing VM setup |
|---|---|---|
| $Throughput_{max}(1)$ (req/sec) | 0.3 | 0.3 |
| $Throughput_{max}(16)$ (req/sec) | 3.37 | 2.78 |
| Ideal $Throughput_{max}(16)$ (req/sec) | 4.8 | 4.8 |
| $S(16)$ | 11.23 | 9.62 |

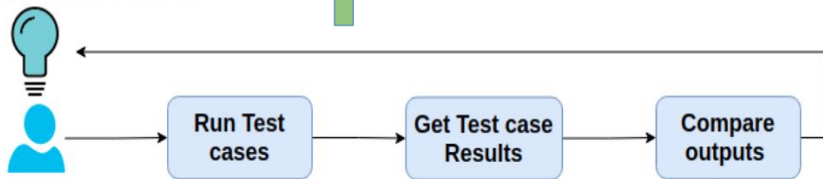Adding new VM doesn't require migration of user data

Load balancer need not maintain user to VM mapping

Throughput scalability decreased with this setup
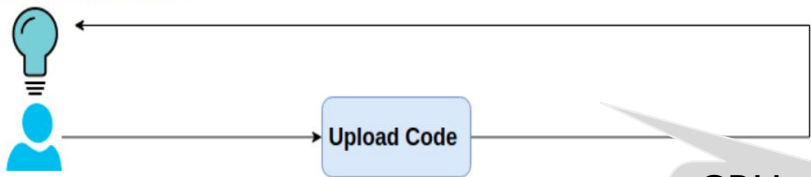
# Reason for reduction in throughput scalability

Performed experiments with different user sessions
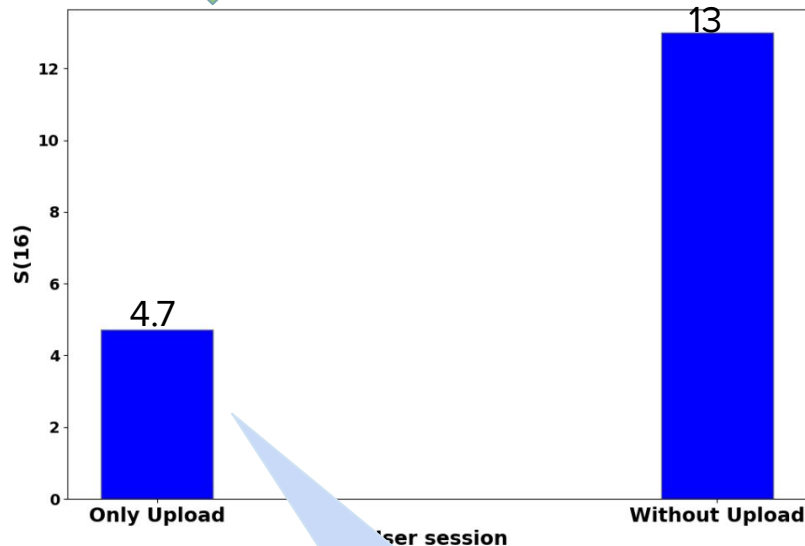
Think time= 8 seconds

Run Test cases → Get Test case Results → Compare outputs

User session: Without Upload

Think time= 15 seconds

Upload Code

User session: Only Upload

CPU utilization never exceeded 35%
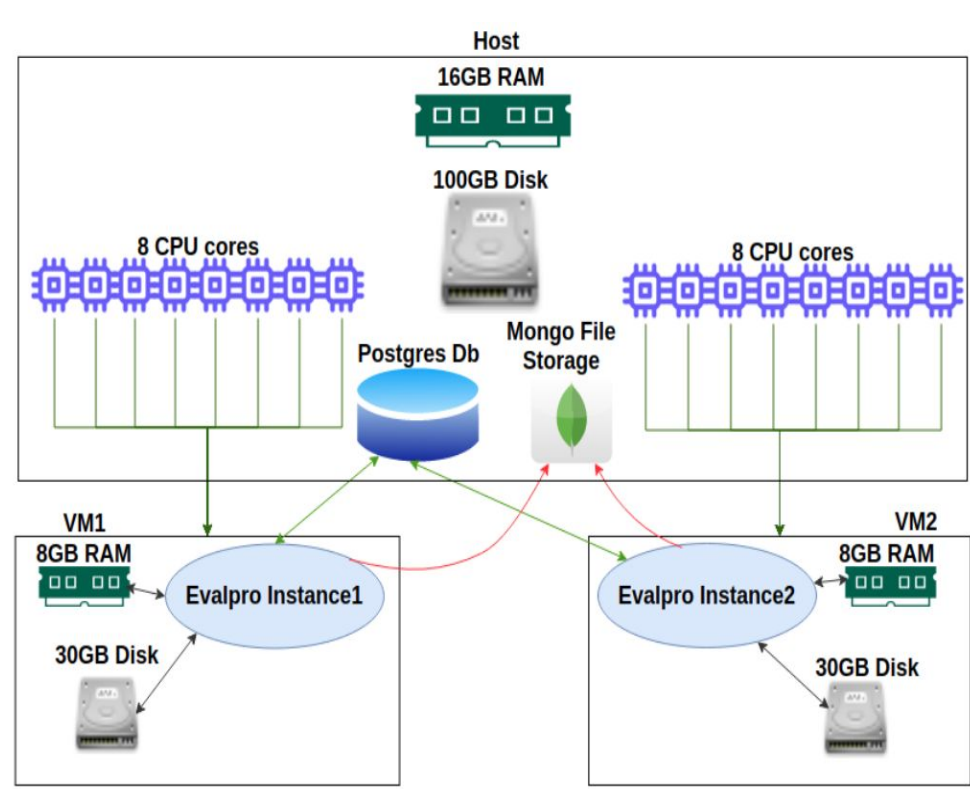
Uploading file by VMs in host is reducing scalability

S(16)

| | 13 |
| Only Upload | Without Upload |
| 4.7 | |

User session

# Using MongoDB for File Storage

# Evalpro updated architecture



Stores user uploaded files

File data is stored as stream of bytes

Users

HAProxy

WSGI

Compilation

Execution

Postgres Db

Celery task queue

Compare outputs

Nginx

Redis

MongoDB

Indicates container

Managed by Docker

# User data and files sharing VM setup with MongoDB



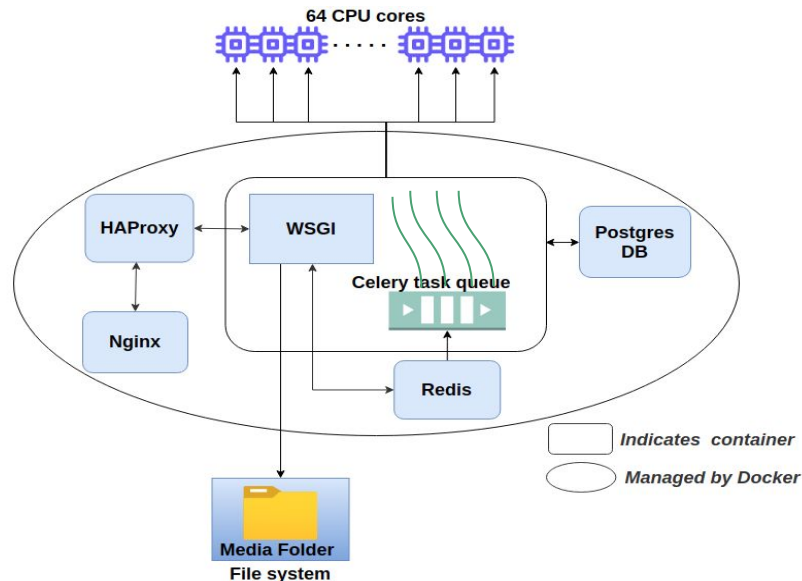| | Baseline | MongoDB for file storage |
|---|---|---|
| $Throughput_{max}(1)$ (req/sec) | 0.3 | 0.29 |
| $Throughput_{max}(16)$ (req/sec) | 3.37 | 3.5 |
| Ideal $Throughput_{max}(16)$ (req/sec) | 4.8 | 4.64 |
| $S(16)$ | 11.23 | 12 |

**Hypothesis**:
With large number of CPU cores throughput scalability will further increase

Throughput scalability slightly increased

# Experiment setup using 64 CPU cores

| Type | CPU | Cores | Memory | L3 Cache | L2 cache | L1 cache |
|------|-----|-------|--------|----------|----------|----------|
| Server | Intel[R] Xeon[R] CPU E5-2683 v4 @ 2.10GHz | 64 | 128GB | 80MB | 8MB | 1MB |
| Client | AMD Opteron[TM] Processor 6212 | 16 | 16GB | 6MB | 2MB | 64KB |
| Client | AMD Opteron[TM] Processor 6278 | 16 | 16GB | 6MB | 2MB | 64KB |
| Client | Intel[R] Xeon[R] CPU E5-2650 v2 @ 2.60GHz | 16 | 16GB | 20MB | 256KB | 32KB |

# Baseline-64



| | |
|---|---|
| $Throughput_{max}(1)$ | 0.66 req/sec |
| $Throughput_{max}(16)$ | 8.75 req/sec |
| $Throughput_{max}(64)$ | 19.56 req/sec |
| Ideal $Throughput_{max}(64)$ | 42 req/sec |
| $S(64)$ | 29.6 |

All CPU cores 90-100% utilized

Linear scaling of throughput not achieved

# User data and files sharing VM setup with MongoDB - 64



| | Baseline-64 | MongoDB for file storage - 64 |
|---|---|---|
| $Throughput_{max}(1)$ (req/sec) | 0.66 | 0.47 |
| $Throughput_{max}(64)$ (req/sec) | 19.56 | 11.27 |
| Ideal $Throughput_{max}(64)$ (req/sec) | 42 | 30 |
| $S(64)$ | 29.6 | 24 |

Our hypothesis is wrong

Throughput scalability decreased

# Completely Isolated VM setup - 64



|  | Baseline-64 | Isolated setup- 64 |
|---|---|---|
| $Throughput_{max}(1)$ (req/sec) | 0.66 | 0.66 |
| $Throughput_{max}(64)$ (req/sec) | 19.56 | 15.03 |
| Ideal $Throughput_{max}(64)$ (req/sec) | 42 | 42 |
|  | 29.6 | 23 |

With higher number of cores, even this setup didn't improve throughput scalability
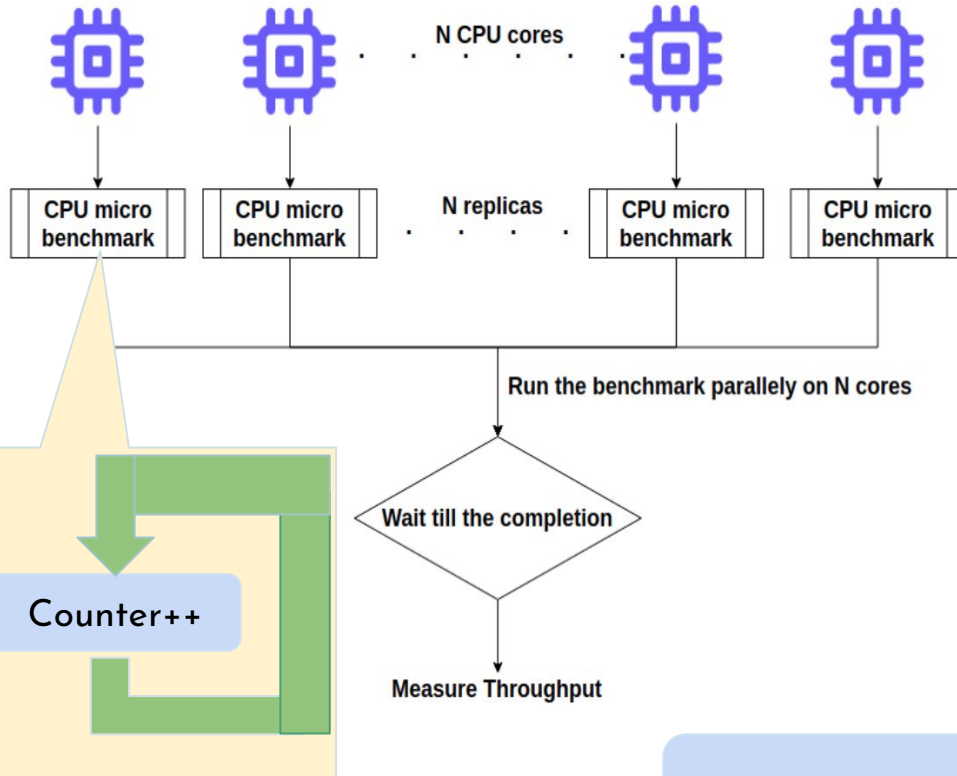
Throughput scalability decreased

# Micro benchmark Experiments

# Need for micro benchmarks

Simple, lightweight and portable → Can be easily profiled and tuned (less lines of code) → Isolating reason for scalability limitation

Establishing best scalability of workload on given platform → Two micro benchmarks

**CPU micro benchmark**
Raw CPU workload
Increments counter for large iterations in tight loop

**Evalpro micro benchmark**
Evalpro-like workload
Evaluation task run in tight loop

# Experiments on CPU micro benchmark



N CPU cores

N replicas

CPU micro benchmark

CPU micro benchmark

CPU micro benchmark

CPU micro benchmark

Run the benchmark parallely on N cores

Wait till the completion

Measure Throughput

Counter++

S(N) vs N

Observed
Ideal

S(64) = 48

Throughput scalability is close to linear

# Experiments on Evalpro micro benchmark

N CPU cores

N replicas

| Evalpro micro benchmark | Evalpro micro benchmark | . . . | Evalpro micro benchmark | Evalpro micro benchmark |

**Hard disk**

**Run the benchmark parallely on N cores**

Compilation

Execution

Compare o/ps

Wait till the completion

**Measure Throughput**

### S(N) vs N

- Evalpro Application
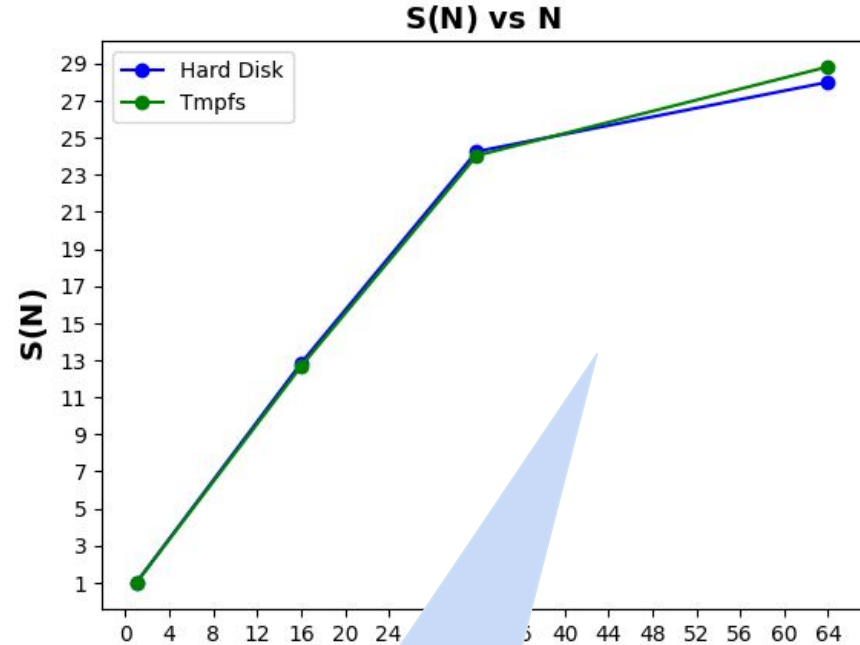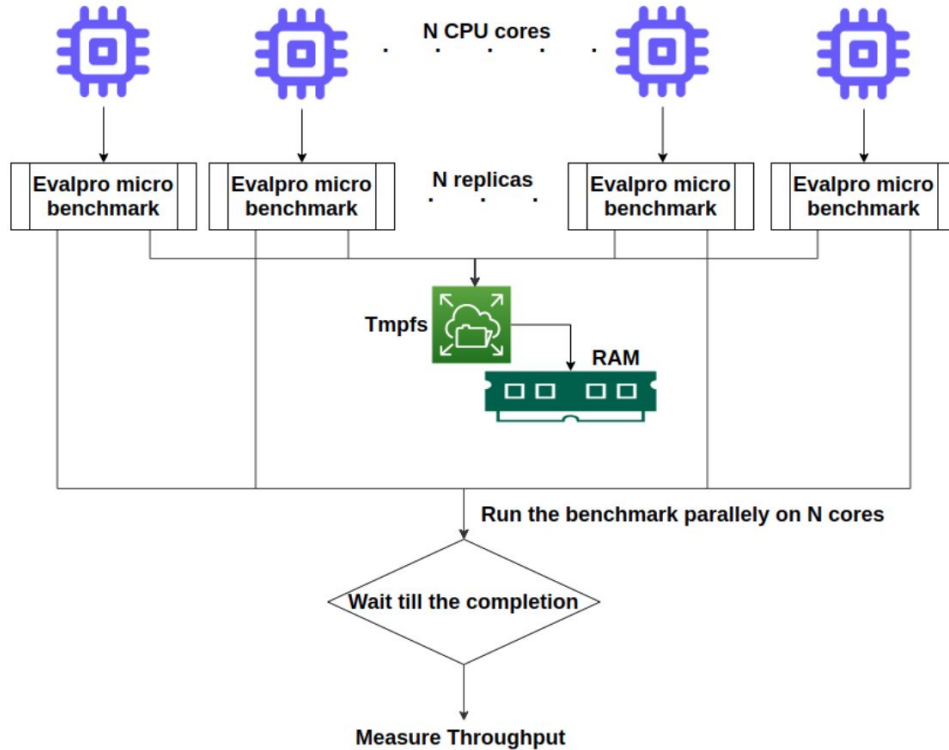- Evalpro micro benchmark

S(64) = 28

Throughput scalability is almost same as Evalpro application

# Micro benchmark Bottleneck Analysis

# Using Tmpfs in place of Hard disk



N CPU cores

Evalpro micro benchmark    Evalpro micro benchmark    N replicas    Evalpro micro benchmark    Evalpro micro benchmark

Tmpfs

RAM

Run the benchmark parallely on N cores

Wait till the completion

Measure Throughput

**S(N) vs N**

Hard Disk
Tmpfs

Using Tmpfs didn't improve throughput scalability
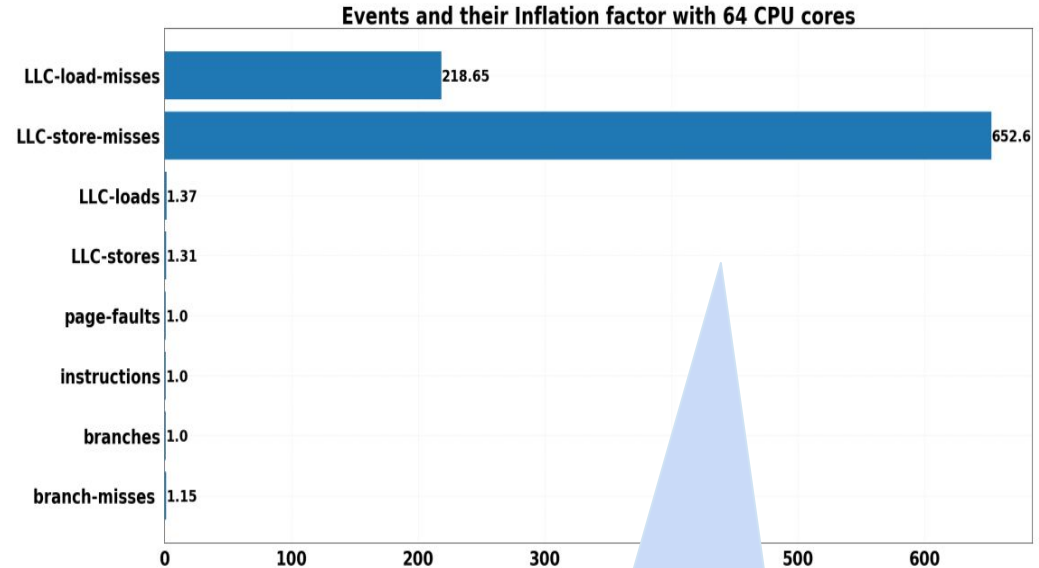
# PERF analysis on Evalpro micro benchmark

PERF profiles different software, hardware events

$Count_{event}(N)$ : Occurrence count of event with N cores
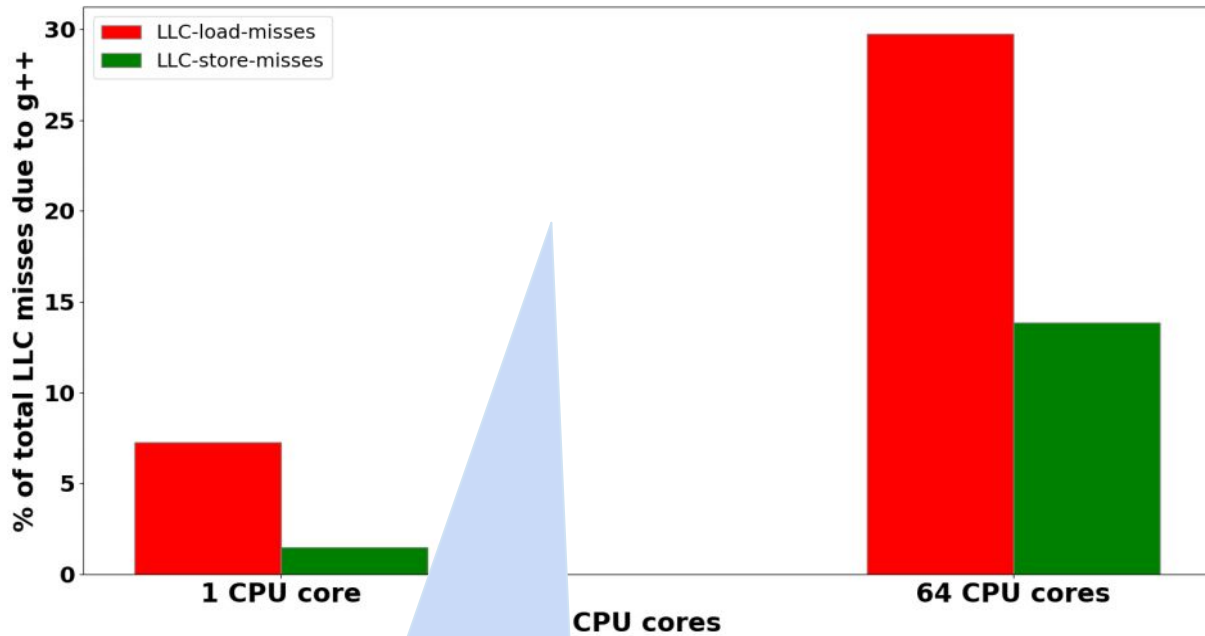$Inflation\_factor_{event}(N)$ : Inflation factor of event with N cores

$Inflation\_factor_{event}(N) =$
$Count_{event}(N) \div ( N \times Count_{event}(1) )$

Ideal $Inflation\_factor_{event}(N) = 1$

**Events and their Inflation factor with 64 CPU cores**

| Event | Inflation factor |
|---|---|
| LLC-load-misses | 218.65 |
| LLC-store-misses | 652.6 |
| LLC-loads | 1.37 |
| LLC-stores | 1.31 |
| page-faults | 1.0 |
| instructions | 1.0 |
| branches | 1.0 |
| branch-misses | 1.15 |

(x-axis: 0, 100, 200, 300, 500, 600)

LLC load and store misses, disproportionately inflated with 64 CPU cores

# PERF analysis at Program level



Hypothesis
CPU cache
is reason for
scalability limitation
of Evalpro application

LLC load and store misses increased drastically by g++ when number of cores is high

CPU cache is the reason for scalability limitation of Evalpro micro benchmark
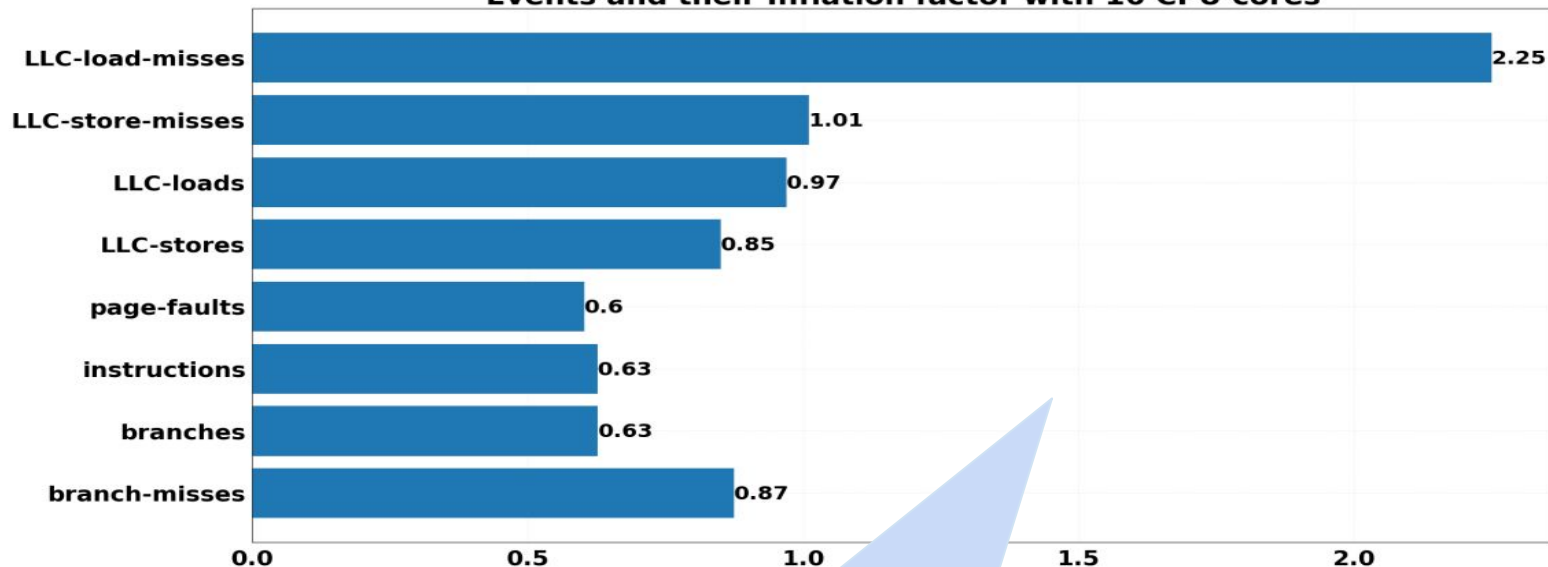
# PERF analysis on Evalpro application

| 16 core | Baseline -16 | 20 MB CPU cache Intel$^R$ Xeon$^R$ CPU E5-2650 v2 @ 2.60GHz |

## Events and their Inflation factor with 16 CPU cores

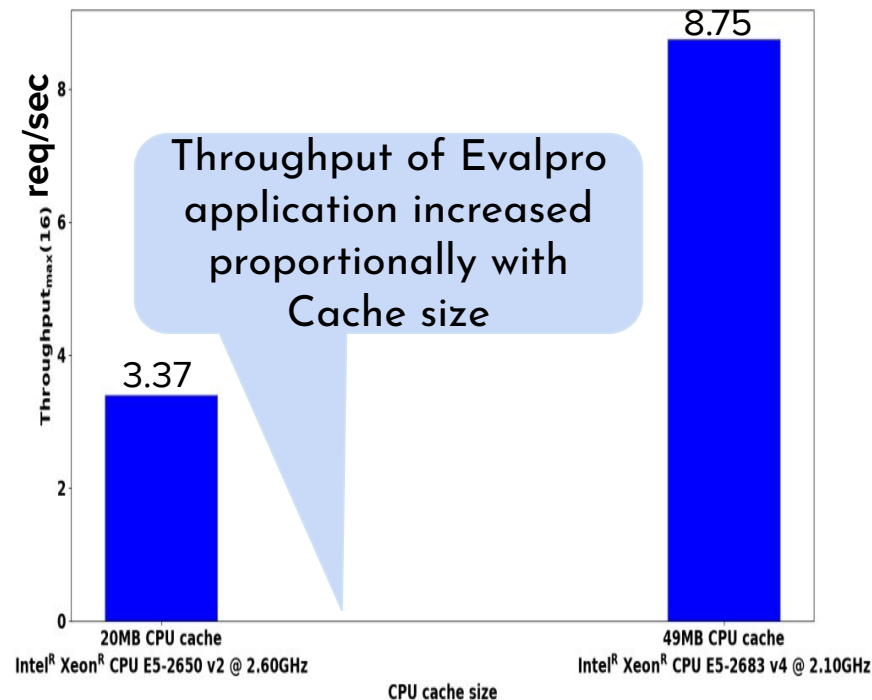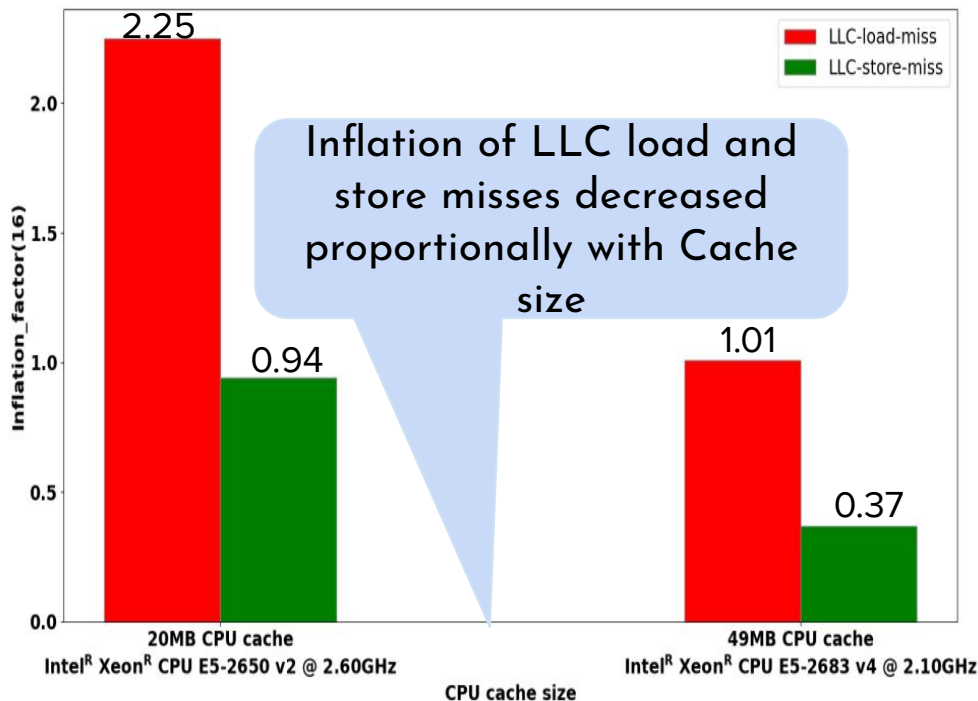| Event | Inflation factor |
|---|---|
| LLC-load-misses | 2.25 |
| LLC-store-misses | 1.01 |
| LLC-loads | 0.97 |
| LLC-stores | 0.85 |
| page-faults | 0.6 |
| instructions | 0.63 |
| branches | 0.63 |
| branch-misses | 0.87 |

LLC load misses inflated
LLC store misses slightly inflated

# Cache size affect on Evalpro application



Baseline - 16
16 cores
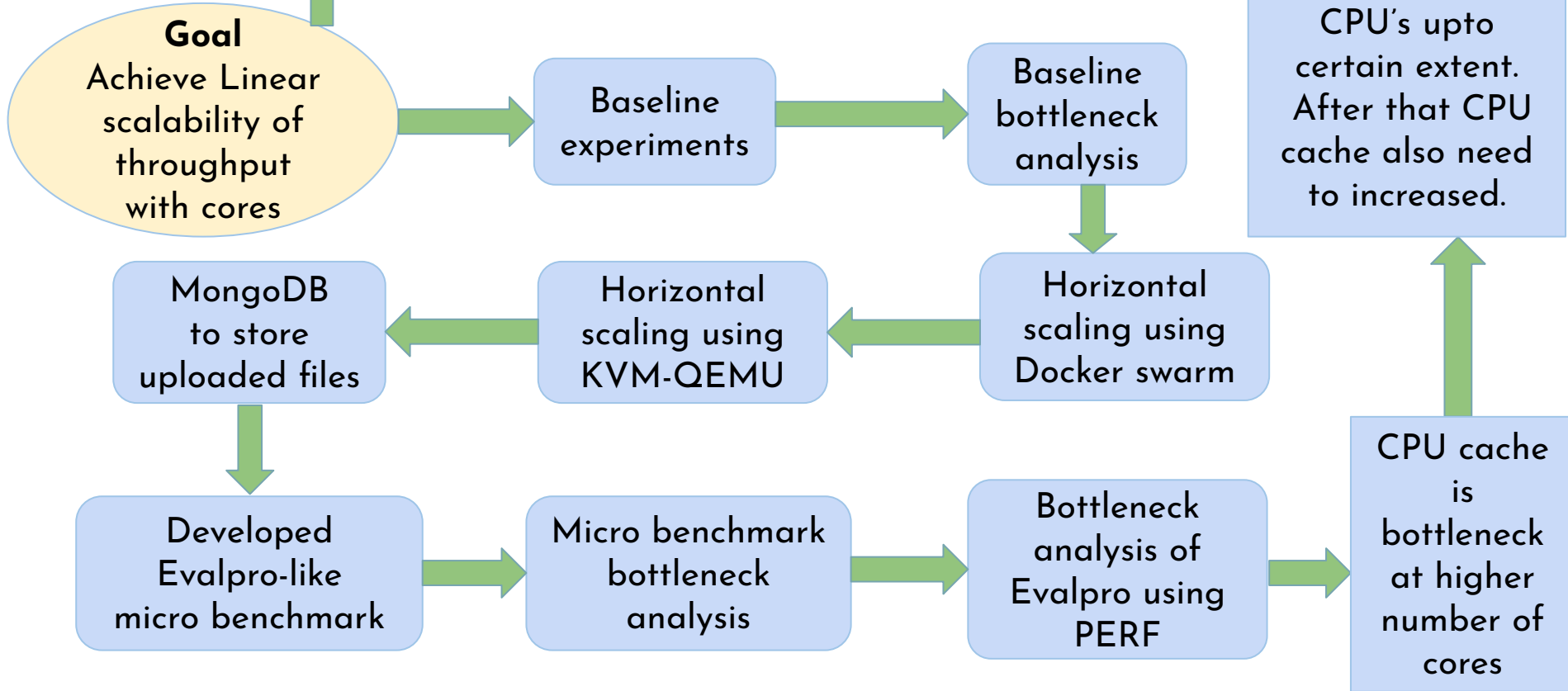20 MB cache
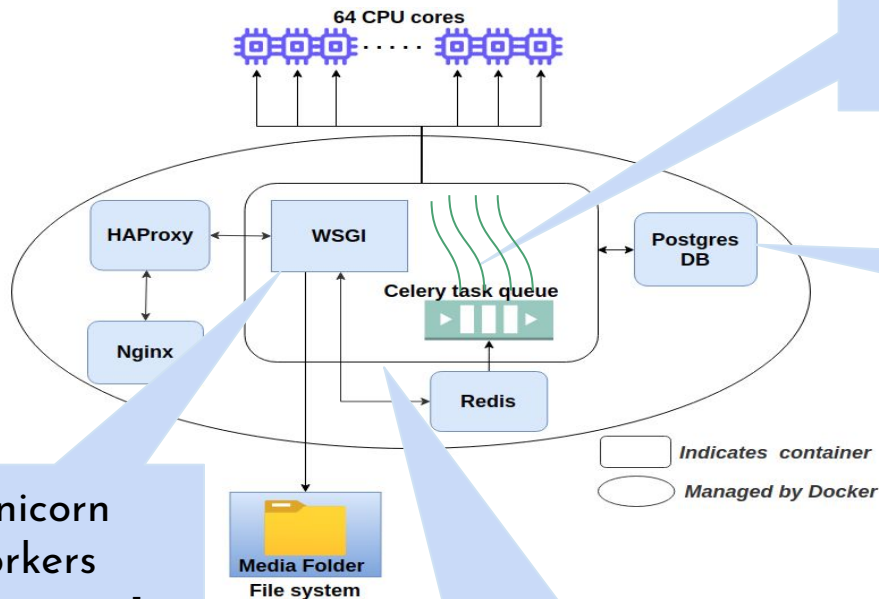
Baseline -64
only 16 cores enabled
89 MB cache

Inflation of LLC load and store misses decreased proportionally with Cache size

Throughput of Evalpro application increased proportionally with Cache size

# Cache size affect summary

Our hypothesis
is correct
i.e CPU Cache
Is the bottleneck

|  | Baseline-16 | Baseline-64 | % of increase or decrease |
|---|---|---|---|
| CPU cache size for 16 CPU cores | 20 MB | 49 MB | + 145% |
| $Inflation\_factor_{LLC\text{-load-misses}}(16)$ | 2.25 | 0.94 | - 140% |
| $Inflation\_factor_{LLC\text{-store-misses}}(16)$ | 1.01 | 0.37 | - 170% |
| $Throughput_{max}(16)$ | 3.37 req/sec | 8.75 req/sec | + 160% |

# Conclusion

**Goal**
Achieve Linear scalability of throughput with cores

Baseline experiments → Baseline bottleneck analysis

Linear scaling of throughput with CPU's upto certain extent. After that CPU cache also need to increased.

MongoDB to store uploaded files ← Horizontal scaling using KVM-QEMU ← Horizontal scaling using Docker swarm

Developed Evalpro-like micro benchmark → Micro benchmark bottleneck analysis → Bottleneck analysis of Evalpro using PERF → CPU cache is bottleneck at higher number of cores

# Final Recommendation for Scaling EvalPro

**64 CPU cores**

celery threads at least number of cores

postgres connections to large number i.e around 100000

**HAProxy**

**WSGI**

**Celery task queue**

**Postgres DB**

**Nginx**

**Redis**

*Indicates container*

*Managed by Docker*

**Media Folder**
**File system**

Gunicorn workers = 2*ncores+1

No need for multiple WSGI+Celery replicas with above configurations

Migrate Evalpro to cloud

Pay for use

Easily scale performance linearly with CPU cores and cache size

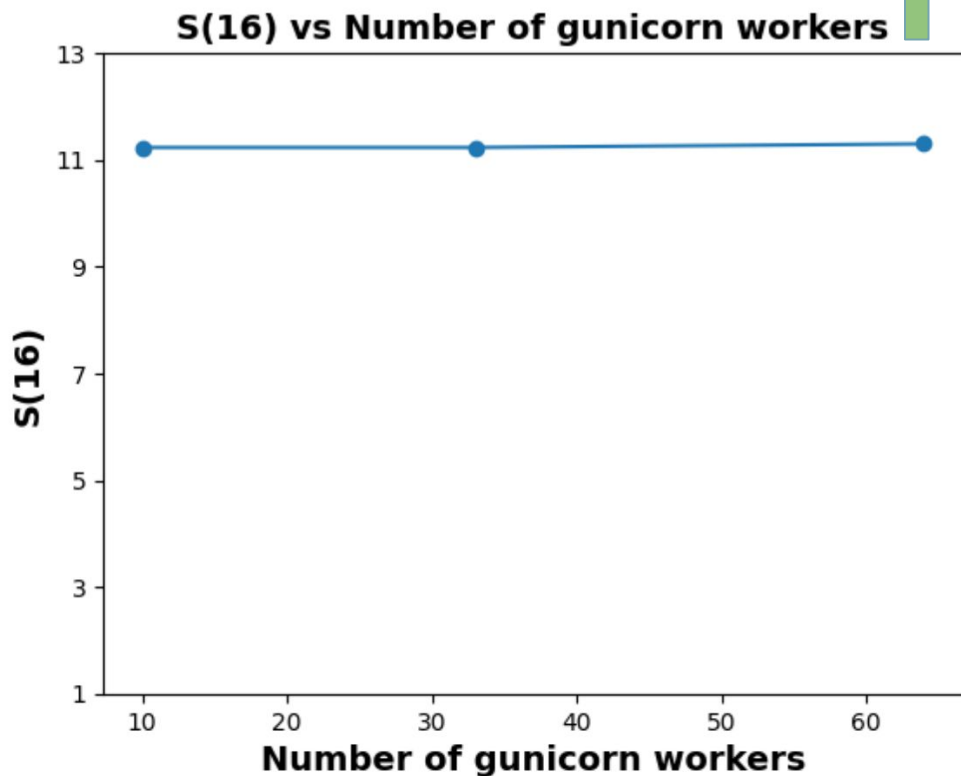# Future work

Scalability of compilation workloads

Developing a browser plugin to do compilation at client side

Improving single core throughput of Evalpro

# THANK YOU

# Backup Slides

# Increasing gunicorn workers



S(16) vs Number of gunicorn workers

Increasing gunicorn workers didn't improve throughput scalability

Gunicorn workers are WSGI worker threads

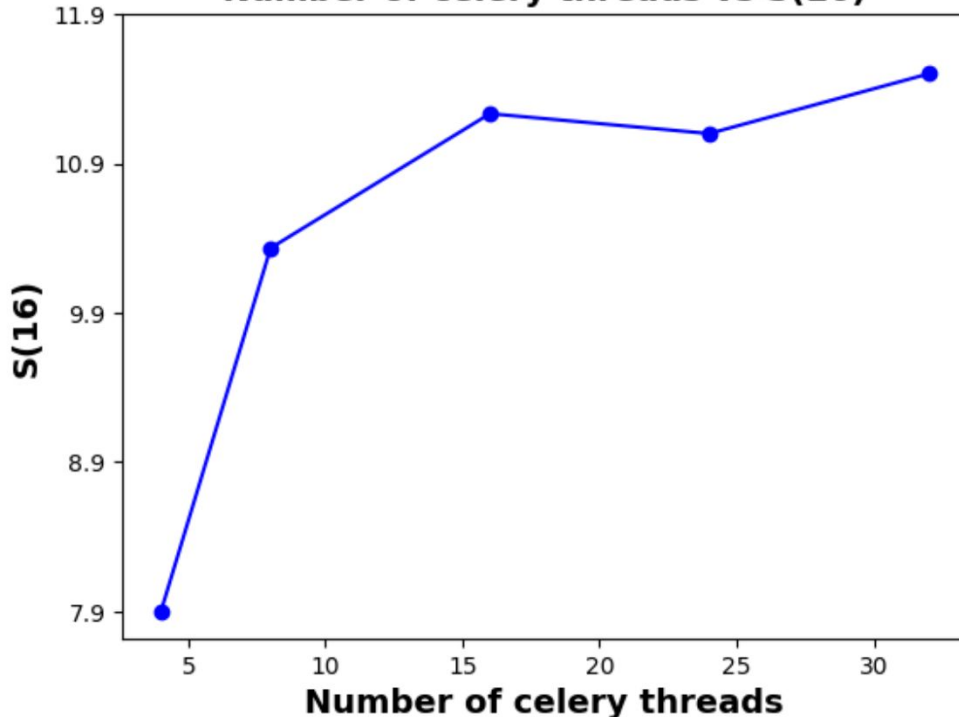Baseline gunicorn workers = 10

Gunicorn workers are not bottleneck

Gunicorn documentation says,
Ideal value = 2 x ncores +1
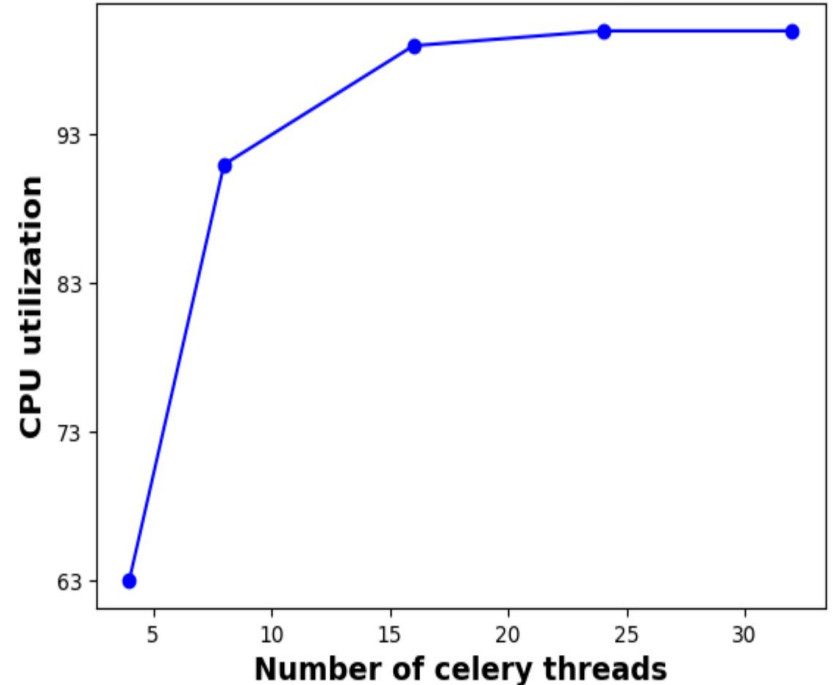
# Tuning celery threads

Auto grading tasks processed asynchronously

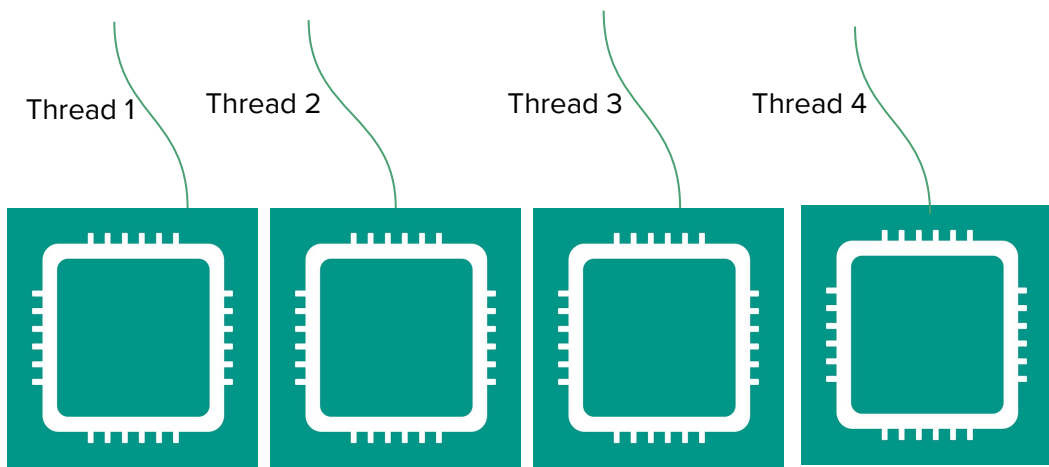Celery threads are bottleneck when less than number of CPU cores



Number of celery threads vs S(16)



CPU utilization vs No of celery threads
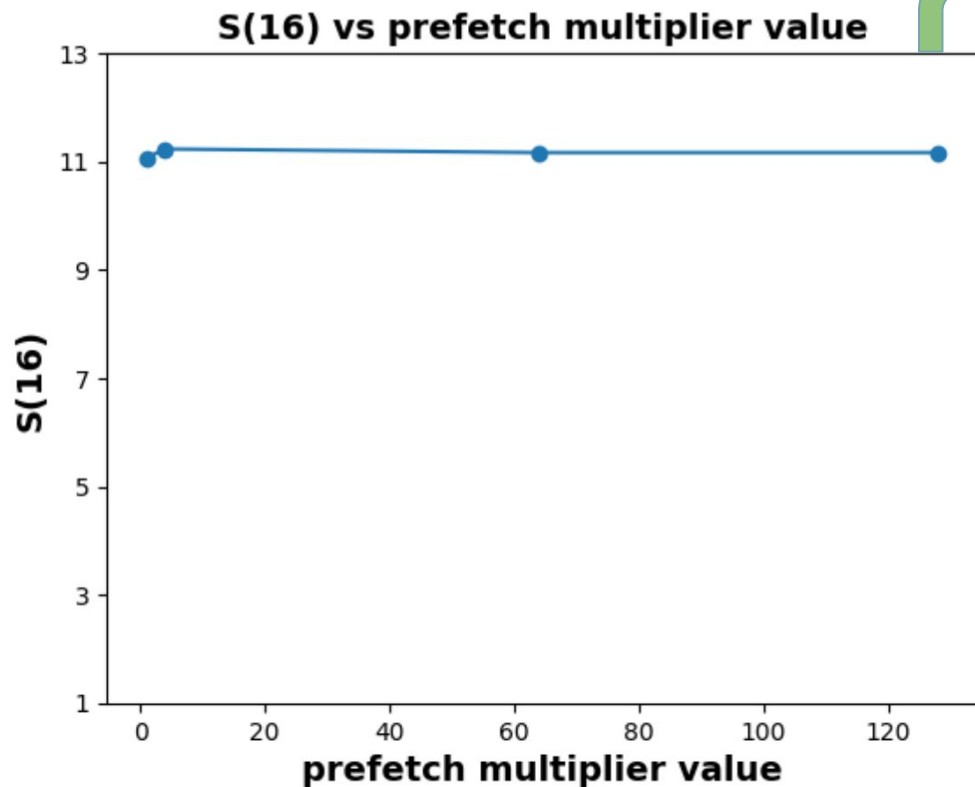
# Setting CPU affinities for celery

| Type | Throughput (req/sec) | S(16) |
|------|------|------|
| Without affinity | 3.37 | 11.23 |
| With affinity | 3.32 | 11 |

Thread 1    Thread 2    Thread 3    Thread 4

· · · · 16 cores

Setting CPU affinities didn't improve throughput scalability

# Increasing celery prefetch multiplier



S(16) vs prefetch multiplier value

Increasing prefetch multiplier didn't improve throughput scalability

No of tasks celery picks and keeps in memory at a time

Baseline prefetch multiplier value = 4

Celery Prefetch multiplier is not bottleneck

# Increasing celery broker pool limit

# Using transient celery queue

Celery threads

queue(by default persistent)

Redis

Stores task metadata

File

Bypass this step

| Type | Throughput (req/sec) | S(16) |
|---|---|---|
| Persistent | 3.37 | 11.23 |
| Transient | 3.3 | 11 |

Making celery queue transient
didn't improve
throughput scalability

# Disabling writing to log files

| Type | Throughput (req/sec) | S(16) |
|---|---|---|
| Persistent | 3.37 | 11.23 |
| Transient | 3.3 | 11 |

Request → Gunicorn

Request info

Gunicorn log

Task → Celery

Task Execution info

Celery log

Disable logging

Disable logging

Disabling logs didn't improve throughput scalability