# PROJECT: Diagnosing the latency–throughput tradeoff in vLLM serving under concurrent loads

- **Experiment 1:** Baseline setup and workload harness for mixed prompt-length serving
- **Experiment 1.1:** Fix max-tokens and context-budget failures by adding dynamic `max_tokens` to prevent request rejections
- **Experiment 2:** BAD (unfair) scheduling characterization with concurrency sweeps to expose head-of-line blocking and tail-latency spikes
- **Experiment 3:** Tail latency and throughput variation with concurrency sweep by holding workload and output length constant to isolate scheduling effects
- **Experiment 4:** Sweep of concurrency and `max-num-seqs` under a deliberately bad configuration to map scaling limits, queueing onset, and short request degradation from long-prefill contention
- **Experiment 5:** Find a strong FAIR scheduling configuration by first sweeping `max-num-batched-tokens` at lower concurrency (c=16), then validating at c=32, and finally running a `max-num-batched-tokens × max-num-seqs` grid at c=32 to pick a stable throughput vs tail-latency sweet spot

## Installing libraries

```
1 !pip -q install datasets transformers
2 !pip -q install tqdm
3
4 !pip -q install "jedi>=0.16"
5 !pip -q install vllm
```

## Importing Libraries

```
1 from google.colab import drive
2 import os
3
4 import json, random, statistics
5 from pathlib import Path
6 from datasets import load_dataset
7 from transformers import AutoTokenizer
8
9 import itertools
10
11 import numpy as np
12 import pandas as pd
13 from tqdm.notebook import tqdm
14
15 import random
16
17 import vllm
18 import torch
19 from vllm import LLM, SamplingParams
20 import requests, time
21
22 import concurrent.futures as cf
23 import pandas as pd
24
25 import shutil
26 import math
27 import subprocess
28
29
```

## Mounting the Drive

```
1 # Mounting Google Drive
2 drive.mount("/content/drive")
```
```
Mounted at /content/drive
```

```
1 # Write locally first for faster write, then copy to Drive at the end
2 OUT_DIR = "/content/workloads"
```

```
 3 os.makedirs(OUT_DIR, exist_ok=True)
 4
 5 # Where to persist outputs in Drive
 6 DRIVE_OUT_DIR = "/content/drive/MyDrive/llm_loadtest_workloads"
 7 os.makedirs(DRIVE_OUT_DIR, exist_ok=True)
 8
 9 print("OUT_DIR:", OUT_DIR)
10 print("DRIVE_OUT_DIR:", DRIVE_OUT_DIR)
11
```

```
OUT_DIR: /content/workloads
DRIVE_OUT_DIR: /content/drive/MyDrive/llm_loadtest_workloads
```

```
 1 MODEL = "Qwen/Qwen2.5-3B-Instruct"
 2 tok = AutoTokenizer.from_pretrained(MODEL, use_fast=True)
 3
 4 def tok_len(text: str) -> int:
 5     return len(tok.encode(text, add_special_tokens=False))
 6
 7 def in_range(n, r):
 8     return r[0] <= n <= r[1]
 9
10 random.seed(42)
11
12 SHORT_RANGE = (8, 24)        # NQ-friendly short prompts
13 MED_RANGE   = (128, 256)     # UltraChat-friendly medium prompts
14 LONG_RANGE  = (2048, 4096)   # long prompts via concatenation
15
16
17 # How many prompts you want in each bucket
18 WANT_SHORT  = 800
19 WANT_MEDIUM = 400
20 WANT_LONG   = 600
21
22 # Requests per workload file
23 TOTAL_REQS_PER_FILE = 400
24
25 print("Tokenizer loaded for:", MODEL)
26
```

```
tokenizer_config.json:     7.30k/? [00:00<00:00, 100kB/s]

vocab.json:     2.78M/? [00:00<00:00, 22.0MB/s]

merges.txt:     1.67M/? [00:00<00:00, 12.0MB/s]

tokenizer.json:     7.03M/? [00:00<00:00, 24.7MB/s]

Tokenizer loaded for: Qwen/Qwen2.5-3B-Instruct
```

## Function to enable Dataset Streaming

```
 1 def try_load_streaming(dataset_name, config_split_candidates):
 2     last_err = None
 3
 4     for cfg, split in config_split_candidates:
 5         try:
 6             if cfg is None:
 7                 ds = load_dataset(dataset_name, split=split, streaming=True)
 8
 9             else:
10                 ds = load_dataset(dataset_name, cfg, split=split, streaming=True)
11             print(f"Loaded {dataset_name} (config={cfg}, split={split})")
12
13             return ds
14
15         except Exception as e:
16             last_err = e
17             print(f"Failed {dataset_name} (config={cfg}, split={split}): {type(e).__name__}")
18
19     raise RuntimeError(f"Could not load {dataset_name}. Last error: {last_err}")
20
21
```

```
 1 # streaming both the datasets
 2 ultra_stream = try_load_streaming(
 3     "HuggingFaceH4/ultrachat_200k",
 4     config_split_candidates=[
 5         (None, "train_sft"),
 6         ("default", "train_sft"),
```

```
 7          ("default", "train"),
 8      ],
 9  )
10
11  nq_stream = try_load_streaming(
12      "natural_questions",
13      config_split_candidates=[
14          (None, "train"),
15          ("default", "train"),
16      ],
17  )
18
```

```
README.md:       3.90k/? [00:00<00:00, 123kB/s]

Loaded HuggingFaceH4/ultrachat_200k (config=None, split=train_sft)

README.md:       13.7k/? [00:00<00:00, 239kB/s]

Resolving data files: 100%                              287/287 [00:00<00:00, 33.69it/s]

Resolving data files: 100%                              287/287 [00:00<00:00, 3896.64it/s]

Loaded natural_questions (config=None, split=train)
```

## Check if the Dataset loaded properly

```
1 ex = next(iter(nq_stream))
2 print("Keys:", ex.keys())
3 print("Sample:", {k: ex[k] for k in list(ex.keys())[:5]})
```

```
Keys: dict_keys(['id', 'document', 'question', 'long_answer_candidates', 'annotations'])
Sample: {'id': '4549465242785278785', 'document': {'html': '<!DOCTYPE html>\n<HTML class="client-js ve-not-availa
```

```
1 ex = next(iter(ultra_stream))
2 print("Keys:", ex.keys())
3 print("Sample:", {k: ex[k] for k in list(ex.keys())[:5]})
```

```
Keys: dict_keys(['prompt', 'prompt_id', 'messages'])
Sample: {'prompt': "These instructions apply to section-based themes (Responsive 6.0+, Retina 4.0+, Parallax 3.0+
```

## Functions to extract prompt from Dataset

```python
 1 def ultrachat_user_text(ex):
 2     prompt = ex.get("prompt")
 3     if isinstance(prompt, str) and prompt.strip():
 4         return prompt.strip()
 5
 6     # This is a fallback: scan the "messages" list for the first user turn
 7     msgs = ex.get("messages")
 8     if isinstance(msgs, list):
 9         for m in msgs:
10             if not isinstance(m, dict):
11                 continue
12
13             role = (m.get("role") or "").lower()
14             content = m.get("content") or ""
15
16             if role == "user" and isinstance(content, str) and content.strip():
17                 return content.strip()
18
19     return None
20
21
22 def nq_question_text(ex):
23     q = ex.get("question")
24     if isinstance(q, dict):
25         t = q.get("text")
26         return t.strip() if isinstance(t, str) and t.strip() else None
27     if isinstance(q, str):
28         return q.strip() if q.strip() else None
29     return None
30
```

## Create Small and Medum Buckets

```python
1  from tqdm import tqdm
2
3  def stream_bucket(
4      stream,
5      text_fn,
6      want_n,
7      token_range,
8      max_read,
9      desc,
10     print_every=500,
11     also_keep_pool=True,
12 ):
13     """
14     Stream from a HF streaming dataset, extract text, collect:
15       – bucket: texts in token_range until want_n
16       – pool: all extracted texts (optional)
17     """
18     bucket = []
19     pool = [] if also_keep_pool else None
20
21     for i, ex in tqdm(enumerate(stream), total=max_read, desc=desc):
22         if i >= max_read:
23             break
24
25         s = text_fn(ex)
26         if not s:
27             continue
28
29         if also_keep_pool:
30             pool.append(s)
31
32         n = tok_len(s)
33         if len(bucket) < want_n and in_range(n, token_range):
34             bucket.append(s)
35
36         if (i % print_every) == 0:
37             tqdm.write(f"read={i:,}, {desc}_bucket={len(bucket)}")
38
39         if len(bucket) >= want_n:
40             break
41
42     return bucket, pool
43
```

```python
1  MAX_NQ_READ = 50000
2
3  short_bucket, nq_pool = stream_bucket(
4      stream=nq_stream,
5      text_fn=nq_question_text,
6      want_n=WANT_SHORT,
7      token_range=SHORT_RANGE,
8      max_read=MAX_NQ_READ,
9      desc="short",
10     print_every=200,    # keep your original cadence
11     also_keep_pool=True,
12 )
13
14 print("short_bucket:", len(short_bucket))
15
```

```
short:    0%|          | 2/50000 [00:18<104:09:36,  7.50s/it]read=0, short_bucket=1
short:    0%|          | 208/50000 [00:30<26:42, 31.07it/s]read=200, short_bucket=201
short:    1%|          | 402/50000 [00:42<2:21:13,  5.85it/s]read=400, short_bucket=401
short:    1%|          | 604/50000 [00:58<35:59, 22.88it/s]read=600, short_bucket=601
short:    2%||         | 799/50000 [01:08<1:10:09, 11.69it/s]
short_bucket: 800
```

```python
1  MAX_ULTRA_READ = 200000
2
3  medium_bucket, ultra_pool = stream_bucket(
4      stream=ultra_stream,
5      text_fn=ultrachat_user_text,
6      want_n=WANT_MEDIUM,
7      token_range=MED_RANGE,
8      max_read=MAX_ULTRA_READ,
9      desc="medium",
10     print_every=500,    # keep your original cadence
11     also_keep_pool=True,
12 )
13
14 print("medium_bucket:", len(medium_bucket))
```

```
medium:    0%|              | 14/200000 [00:03<9:37:24,  5.77it/s] read=0, medium_bucket=0
medium:    0%|              | 528/200000 [00:04<10:42, 310.50it/s]read=500, medium_bucket=51
medium:    1%|              | 1040/200000 [00:06<09:46, 339.40it/s]read=1,000, medium_bucket=106
medium:    1%|              | 1537/200000 [00:08<12:04, 274.11it/s]read=1,500, medium_bucket=158
medium:    1%|              | 2042/200000 [00:09<09:52, 334.10it/s]read=2,000, medium_bucket=198
medium:    1%||             | 2543/200000 [00:11<10:07, 325.01it/s]read=2,500, medium_bucket=252
medium:    2%||             | 3004/200000 [00:13<23:16, 141.10it/s]read=3,000, medium_bucket=305
medium:    2%||             | 3614/200000 [00:19<07:18, 448.23it/s]read=3,500, medium_bucket=355
medium:    2%||             | 3990/200000 [00:20<17:07, 190.85it/s]medium_bucket: 400
```

## Create Large Bucket

```python
1 def build_long_prompt(target_tokens, rng, snippets_pool, max_unique_before_reset=2000):
2     parts = []
3     used = set()
4
5     while True:
6         idx = rng.randrange(len(snippets_pool))
7         if idx in used:
8             continue
9         used.add(idx)
10
11        parts.append(snippets_pool[idx].replace("\n", " ").strip())
12        text = "\n\n".join(parts)
13
14        if tok_len(text) >= target_tokens:
15            break
16
17        if len(used) > min(max_unique_before_reset, len(snippets_pool) - 1):
18            used.clear()
19
20    ids = tok.encode(text, add_special_tokens=False)[:target_tokens]
21    return tok.decode(ids)
22
23 def build_long_bucket(
24     want_n,
25     token_range,
26     snippets_pool,
27     seed=123,
28     desc="long",
29     print_every=10,
30 ):
31     assert len(snippets_pool) > 0, "snippets_pool must not be empty"
32
33     rng = random.Random(seed)
34     bucket = []
35
36     for j in tqdm(range(want_n), desc=f"Building {desc}"):
37         target = rng.randint(token_range[0], token_range[1])
38         lp = build_long_prompt(target, rng, snippets_pool)
39         bucket.append(lp)
40
41         if ((j + 1) % print_every) == 0:
42             tqdm.write(f"built={j+1}/{want_n}, last_len={tok_len(lp)} tokens")
43
44     return bucket
45
```

```python
1 snippets_pool = short_bucket + medium_bucket
2 long_bucket = build_long_bucket(
3     want_n=WANT_LONG,
4     token_range=LONG_RANGE,
5     snippets_pool=snippets_pool,
6     seed=123,
7     desc="long",
8     print_every=10,
9 )
10
11 print("long_bucket:", len(long_bucket))
12
```

Show hidden output

## Create Mixed Buckets

```python
1
2 from tqdm.auto import tqdm
3
4 TOTAL_REQS_PER_FILE = 400
5 SEED = 999
6
7 OVERWRITE_LOCAL = False
8 OVERWRITE_DRIVE = False
9
10 ENABLE_LONG_HEAVY = True
11 HEAVY_LONG_RANGE = (3500, 4096)
12 WANT_LONG_HEAVY = 600
13 USE_LONG_HEAVY_FOR = {"mix_60_40.jsonl", "mix_50_50.jsonl"}
14
15 MIX_SPECS = {
16     "short_only.jsonl":    {"short": 1.0},
17     "medium_only.jsonl":   {"medium": 1.0},
18     "long_only.jsonl":     {"long": 1.0},
19
20     "mix_90_10.jsonl":     {"short": 0.9, "long": 0.1},
21     "mix_70_30.jsonl":     {"short": 0.7, "long": 0.3},
22     "mix_70_20_10.jsonl":  {"short": 0.7, "medium": 0.2, "long": 0.1},
23
24     # these use long heavy
25     "mix_60_40.jsonl":     {"short": 0.6, "long": 0.4},
26     "mix_50_50.jsonl":     {"short": 0.5, "long": 0.5},
27 }
28
29 assert "tok" in globals(), "Tokenizer 'tok' not found."
30 assert "short_bucket" in globals() and len(short_bucket) > 0, "short_bucket missing/empty."
31 assert "medium_bucket" in globals() and len(medium_bucket) > 0, "medium_bucket missing/empty."
32 assert "long_bucket" in globals() and len(long_bucket) > 0, "long_bucket missing/empty."
33
34
35 def tok_len(text: str) -> int:
36     return len(tok.encode(text, add_special_tokens=False))
37
38
39 def build_long_prompt(target_tokens, rng, snippets_pool):
40     parts = []
41     used = set()
42     while True:
43         idx = rng.randrange(len(snippets_pool))
44         if idx in used:
45             continue
46         used.add(idx)
47
48         parts.append(snippets_pool[idx].replace("\n", " ").strip())
49         text = "\n\n".join(parts)
50
51         if tok_len(text) >= target_tokens:
52             break
53
54         if len(used) > min(2000, len(snippets_pool) - 1):
55             used.clear()
56
57     ids = tok.encode(text, add_special_tokens=False)[:target_tokens]
58     return tok.decode(ids)
59
60
61 def build_long_bucket(want_n, token_range, snippets_pool, seed=1234, desc="long_bucket"):
62     rng = random.Random(seed)
63     bucket = []
64
65     for _ in tqdm(range(want_n), desc=f"Building {desc} {token_range}"):
66         target = rng.randint(token_range[0], token_range[1])
67         bucket.append(build_long_prompt(target, rng, snippets_pool))
68
69     return bucket
70
71
72 long_heavy_bucket = None
73 if ENABLE_LONG_HEAVY and len(USE_LONG_HEAVY_FOR) > 0:
74     snippets_pool = short_bucket + medium_bucket
75     long_heavy_bucket = build_long_bucket(
76         want_n=WANT_LONG_HEAVY,
77         token_range=HEAVY_LONG_RANGE,
78         snippets_pool=snippets_pool,
79         seed=1234,
80         desc="long_heavy_bucket",
81     )
82
```

```python
 83     print(
 84         "long_heavy_bucket:",
 85         len(long_heavy_bucket),
 86         "example_len_tokens:",
 87         tok_len(long_heavy_bucket[0]),
 88     )
 89
 90 def write_jsonl(path, prompts, max_tokens=128, temperature=0.0):
 91     with open(path, "w", encoding="utf-8") as f:
 92         for p in prompts:
 93             row = {
 94                 "messages": [{"role": "user", "content": p}],
 95                 "max_tokens": max_tokens,
 96                 "temperature": temperature,
 97             }
 98             f.write(json.dumps(row) + "\n")
 99
100
101 def sample_prompts(bucket, k, rng):
102     if k <= 0:
103         return []
104     if len(bucket) < k:
105         raise ValueError(f"Bucket too small: need {k}, have {len(bucket)}")
106     return rng.sample(bucket, k)
107
108
109 def normalize_ratios(ratios: dict) -> dict:
110     r = dict(ratios)
111     s = float(r.get("short", 0.0))
112     m = float(r.get("medium", 0.0))
113     l = float(r.get("long", 0.0))
114     total = s + m + l
115
116     if total > 1.000001:
117         raise ValueError(f"Ratios sum > 1.0: {ratios}")
118
119     if "long" not in r:
120         r["long"] = 1.0 - (s + m)
121
122     return r
123
124 def make_mix(total, ratios, rng, short_b, medium_b, long_b):
125     ratios = normalize_ratios(ratios)
126     k_short  = int(total * ratios.get("short", 0))
127     k_medium = int(total * ratios.get("medium", 0))
128     k_long   = total - k_short - k_medium
129
130     prompts = []
131     prompts += sample_prompts(short_b,  k_short,  rng)
132     prompts += sample_prompts(medium_b, k_medium, rng)
133     prompts += sample_prompts(long_b,   k_long,   rng)
134
135     rng.shuffle(prompts)
136
137     return prompts
138
139 def should_write(path, overwrite=False):
140     return overwrite or (not os.path.exists(path))
141
142
143 def safe_copy(local_path, drive_path, overwrite=False):
144     if not overwrite and os.path.exists(drive_path):
145         return False
146
147     shutil.copy(local_path, drive_path)
148
149     return True
150
151 rng = random.Random(SEED)
152
153 written = []
154 skipped = []
155
156 for fn, ratios in MIX_SPECS.items():
157     local_path = os.path.join(OUT_DIR, fn)
158     drive_path = os.path.join(DRIVE_OUT_DIR, fn)
159
160     if (not OVERWRITE_LOCAL) and os.path.exists(local_path) and (not OVERWRITE_DRIVE) and os.path.exists(dr
161         skipped.append(fn)
162         continue
163
164     use_heavy = (fn in USE_LONG_HEAVY_FOR) and (long_heavy_bucket is not None)
```

```
165        long_source = long_heavy_bucket if use_heavy else long_bucket
166
167     prompts = make_mix(
168         total=TOTAL_REQS_PER_FILE,
169         ratios=ratios,
170         rng=rng,
171         short_b=short_bucket,
172         medium_b=medium_bucket,
173         long_b=long_source,
174     )
175
176     if should_write(local_path, overwrite=OVERWRITE_LOCAL):
177         write_jsonl(local_path, prompts)
178     else:
179         skipped.append(fn)
180         continue
181
182     copied = safe_copy(local_path, drive_path, overwrite=OVERWRITE_DRIVE)
183     if not copied:
184         print("Exists on Drive, skipping copy:", drive_path)
185
186     written.append(fn)
187
188 print("Done.")
189 print("Written:", written)
190 print("Skipped:", skipped)
191 print("Local dir:", OUT_DIR)
192 print("Drive dir:", DRIVE_OUT_DIR)
193
```

```
Building long_heavy_bucket (3500, 4096): 100%                          600/600 [07:49<00:00, 1.57it/s]
long_heavy_bucket: 600 example_len_tokens: 3951
Exists on Drive, skipping copy: /content/drive/MyDrive/llm_loadtest_workloads/short_only.jsonl
Exists on Drive, skipping copy: /content/drive/MyDrive/llm_loadtest_workloads/medium_only.jsonl
Exists on Drive, skipping copy: /content/drive/MyDrive/llm_loadtest_workloads/long_only.jsonl
Exists on Drive, skipping copy: /content/drive/MyDrive/llm_loadtest_workloads/mix_90_10.jsonl
Exists on Drive, skipping copy: /content/drive/MyDrive/llm_loadtest_workloads/mix_70_30.jsonl
Exists on Drive, skipping copy: /content/drive/MyDrive/llm_loadtest_workloads/mix_70_20_10.jsonl
Exists on Drive, skipping copy: /content/drive/MyDrive/llm_loadtest_workloads/mix_60_40.jsonl
Exists on Drive, skipping copy: /content/drive/MyDrive/llm_loadtest_workloads/mix_50_50.jsonl
Done.
Written: ['short_only.jsonl', 'medium_only.jsonl', 'long_only.jsonl', 'mix_90_10.jsonl', 'mix_70_30.jsonl', 'mix_
Skipped: []
Local dir: /content/workloads
Drive dir: /content/drive/MyDrive/llm_loadtest_workloads
```

```
 1 WORKLOAD_DIR = "/content/workloads"
 2 FILES = sorted([f for f in os.listdir(WORKLOAD_DIR) if f.endswith(".jsonl")])
 3
 4
 5 def tok_len(text: str) -> int:
 6     return len(tok.encode(text, add_special_tokens=False))
 7
 8 def read_jsonl(path: str):
 9     rows = []
10
11     with open(path, "r", encoding="utf-8") as f:
12         for line in f:
13             if not line.strip():
14                 continue
15
16             rows.append(json.loads(line))
17
18     return rows
19
20 def extract_prompt(row: dict) -> str:
21     msgs = row.get("messages", [])
22
23     if isinstance(msgs, list) and len(msgs) > 0 and isinstance(msgs[0], dict):
24         return (msgs[0].get("content") or "").strip()
25
26     return ""
27
28 def stats(arr):
29     arr = np.array(arr, dtype=float)
30
31     if len(arr) == 0:
32         return {}
33
34     return {
35         "n": int(len(arr)),
36         "min": float(np.min(arr)),
```

```python
37            "p50": float(np.percentile(arr, 50)),
38            "p90": float(np.percentile(arr, 90)),
39            "p95": float(np.percentile(arr, 95)),
40            "p99": float(np.percentile(arr, 99)),
41            "max": float(np.max(arr)),
42            "mean": float(np.mean(arr)),
43        }
44
45  results = []
46  per_file_details = {}
47
48  for fn in tqdm(FILES, desc="Verifying workload JSONLs"):
49      path = os.path.join(WORKLOAD_DIR, fn)
50      rows = read_jsonl(path)
51
52      prompts = [extract_prompt(r) for r in rows]
53      prompts = [p for p in prompts if p]  # drop empties
54
55      tok_lens = [tok_len(p) for p in prompts]
56      char_lens = [len(p) for p in prompts]
57
58      s_tok = stats(tok_lens)
59      s_chr = stats(char_lens)
60
61      results.append({
62          "file": fn,
63          "n": s_tok.get("n", 0),
64          "tok_min": s_tok.get("min"),
65          "tok_p50": s_tok.get("p50"),
66          "tok_p90": s_tok.get("p90"),
67          "tok_p95": s_tok.get("p95"),
68          "tok_p99": s_tok.get("p99"),
69          "tok_max": s_tok.get("max"),
70          "tok_mean": s_tok.get("mean"),
71          "char_p50": s_chr.get("p50"),
72          "char_p99": s_chr.get("p99"),
73      })
74
75      # store some examples for sanity checks
76      if len(tok_lens) > 0:
77          idx_sorted = np.argsort(tok_lens)
78          smallest = idx_sorted[:2].tolist()
79          largest = idx_sorted[-2:].tolist()
80          per_file_details[fn] = {
81              "smallest_examples": [(tok_lens[i], prompts[i][:200]) for i in smallest],
82              "largest_examples": [(tok_lens[i], prompts[i][:200]) for i in largest],
83          }
84
85  df_stats = pd.DataFrame(results).sort_values("file")
86  display(df_stats)
87
88  print("\n--- Sanity samples (first 200 chars) ---")
89  for fn, det in per_file_details.items():
90      print(f"\n### {fn}")
91      print("Smallest:")
92      for L, s in det["smallest_examples"]:
93          print(f"  tok={L} | {s!r}")
94      print("Largest:")
95      for L, s in det["largest_examples"]:
96          print(f"  tok={L} | {s!r}")
97
```

Verifying workload JSONLs: 100%                                    8/8 [00:20<00:00,  2.03s/it]

| | file | n | tok_min | tok_p50 | tok_p90 | tok_p95 | tok_p99 | tok_max | tok_mean | char_p50 | char_p99 |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | long_only.jsonl | 400 | 2051.0 | 3068.5 | 3931.1 | 4028.10 | 4083.01 | 4092.0 | 3061.0350 | 14489.5 | 19827.18 |
| 1 | medium_only.jsonl | 400 | 128.0 | 169.0 | 231.1 | 240.05 | 253.00 | 256.0 | 175.9325 | 819.0 | 1378.01 |
| 2 | mix_50_50.jsonl | 400 | 8.0 | 1761.0 | 4018.1 | 4062.10 | 4090.02 | 4096.0 | 1918.2400 | 7905.0 | 19973.59 |
| 3 | mix_60_40.jsonl | 400 | 8.0 | 12.5 | 3983.1 | 4042.10 | 4089.01 | 4095.0 | 1532.5800 | 56.0 | 19942.21 |
| 4 | mix_70_20_10.jsonl | 400 | 8.0 | 11.0 | 432.9 | 3189.05 | 3898.31 | 4090.0 | 357.9650 | 49.5 | 18275.65 |
| 5 | mix_70_30.jsonl | 400 | 8.0 | 11.0 | 3525.5 | 3852.10 | 4070.13 | 4092.0 | 939.5800 | 50.0 | 19346.25 |
| 6 | mix_90_10.jsonl | 400 | 8.0 | 10.0 | 239.1 | 3332.80 | 3887.60 | 4090.0 | 334.2300 | 46.0 | 18965.18 |
| 7 | short_only.jsonl | 400 | 8.0 | 10.0 | 14.0 | 15.00 | 19.01 | 23.0 | 10.5925 | 45.0 | 88.01 |

--- Sanity samples (first 200 chars) ---

### long_only.jsonl
Smallest:
  tok=2051 | 'who wrote the song come thou fount of every blessing\n\nCan you summarize the setting and atmospher
  tok=2058 | 'Write a fast-paced, plot-driven thriller novel that chronicles the story of a young woman abducted
Largest:
  tok=4090 | "elvis presley daddy don't you walk so fast\n\nhow old do you have to be to get a tattoo in utah\n\n
  tok=4092 | 'when was the last time a total eclipse happened\n\nWrite a detailed report that analyzes and summar

### medium_only.jsonl
Smallest:
  tok=128 | 'Great way to get your message out there, "Happy Birthday", "Will You Marry Me" and More!\nLet them k
  tok=128 | 'Implement a queue data structure in Java that can enqueue and dequeue elements. The queue should hav
Largest:
  tok=255 | 'such as. The grüffelo) or poems for children. Do you have any tips? Recommendations? Or would you li
  tok=256 | '1. Roast and powder coriander seeds, cumin seeds and red chillies.\n2. Grind together ginger-garlic.

### mix_50_50.jsonl
Smallest:
  tok=8 | 'most commonly used punishment in the united states'
  tok=8 | 'i m in the love with the coco'
Largest:
  tok=4095 | 'who does the voice of shaggy on scooby doo\n\ntom and jerry movie the fast and the furry\n\nIs it p
  tok=4096 | 'deeper than the holler what is a holler\n\nwho played tj on head of the class\n\nis the song fight

### mix_60_40.jsonl
Smallest:
  tok=8 | 'does president have to be born in usa'
  tok=8 | 'who wants to marry a prince tv show'
Largest:
  tok=4095 | 'Can you provide a brief summary of the central problem facing Dell management according to the HBS
  tok=4095 | 'who does the voice of shaggy on scooby doo\n\ntom and jerry movie the fast and the furry\n\nIs it p

### mix_70_20_10.jsonl
Smallest:
  tok=8 | 'properties of red black tree in data structure'
  tok=8 | 'what are the parts of the finger called'
Largest:
  tok=4084 | 'Write a free verse poem that captures the beauty and essence of each season, drawing inspiration fr
  tok=4090 | "elvis presley daddy don't you walk so fast\n\nhow old do you have to be to get a tattoo in utah\n\n

### mix_70_30.jsonl
Smallest:
  tok=8 | 'what number book is mark in the bible'
  tok=8 | 'the legend of the monkey king tv series'
Largest:
  tok=4090 | "elvis presley daddy don't you walk so fast\n\nhow old do you have to be to get a tattoo in utah\n\n
  tok=4092 | 'when was the last time a total eclipse happened\n\nWrite a detailed report that analyzes and summar

### mix_90_10.jsonl
Smallest:
  tok=8 | 'why do ice cubes get smaller over time'
  tok=8 | 'early blues musicians used which combination of instruments'
Largest:
  tok=4032 | "what is the meaning of yanny and laurel\n\nearly blues musicians used which combination of instrume
  tok=4090 | "elvis presley daddy don't you walk so fast\n\nhow old do you have to be to get a tattoo in utah\n\n

### short_only.jsonl
Smallest:
  tok=8 | 'who sang take that look off your face'
  tok=8 | 'where is the us military base in japan'
Largest:
  tok=21 | 'nasa conducted a poll with 9 names to choose the name of the rover which landed on mars recently'
  tok=23 | 'erupted on august 3 1965 as the result of a routine arrest of a drunk driver'

## Setting up vLLM server and Dataset loading

```
1 print("vLLM:", vllm.__version__)
2 print("torch:", torch.__version__, "cuda:", torch.version.cuda)
3 print("GPU:", torch.cuda.get_device_name(0))
```

```
vLLM: 0.13.0
torch: 2.9.0+cu126 cuda: 12.6
GPU: NVIDIA L4
```

```
1 # Use this cell, when restarting the Server again in the same runtime
2 # it clears out the previous instances prevents issues in server setup
3
4 !pkill -f "vllm.entrypoints.openai.api_server" || true
5 !pkill -f "EngineCore" || true
6 !pkill -f "vllm" || true
7
```

```
^C
^C
^C
```

```
 1 MODEL_ID = "Qwen/Qwen2.5-3B-Instruct"
 2 PORT = 8000
 3
 4 MAX_MODEL_LEN = 4096
 5 GPU_UTIL = 0.70
 6 DTYPE = "half"              # fp16 weights/compute
 7
 8 # Launch server
 9 # run it in the background and log to a file.
10 cmd = f"""
11 nohup python3 -m vllm.entrypoints.openai.api_server \
12   --model {MODEL_ID} \
13   --host 127.0.0.1 \
14   --port {PORT} \
15   --max-model-len {MAX_MODEL_LEN} \
16   --gpu-memory-utilization {GPU_UTIL} \
17   --dtype {DTYPE} \
18   --enforce-eager \
19   --disable-log-stats \
20   > vllm_server.log 2>&1 &
21 """
22
23 print("Launching vLLM server...")
24 !bash -lc "$cmd"
25
26 print("Server process launched. Showing last 30 log lines:")
27 !tail -n 30 vllm_server.log
28
```

```
Launching vLLM server...
Server process launched. Showing last 30 log lines:
```

## Function to check the successful startup of the server

```
 1 import time, requests
 2
 3 BASE = "http://127.0.0.1:8000"
 4
 5 def wait_for_server(timeout_s=600):
 6     t0 = time.time()
 7     while time.time() - t0 < timeout_s:
 8         try:
 9             r = requests.get(f"{BASE}/v1/models", timeout=2)
10             if r.status_code == 200:
11                 return True, r.json()
12         except Exception:
13             pass
14         time.sleep(2)
15     return False, None
16
17 ok, payload = wait_for_server()
18 print("Ready:", ok)
19 if ok:
20     print("Models:", [m["id"] for m in payload.get("data", [])])
21 else:
```

```
22    print("Not ready yet. Check logs:")
23    !tail -n 60 vllm_server.log
24
```

```
Ready: True
Models: ['Qwen/Qwen2.5-3B-Instruct']
```

```
 1 # Server take almost a minute to start,
 2 # keep checking the logs to confirm successful startup
 3 # before running any experiments
 4
 5 !tail -n 200 vllm_server.log
 6
```

Show hidden output

## Sanity check after Server start

```
 1 import requests, json
 2
 3 BASE="http://127.0.0.1:8000"
 4
 5 payload = {
 6   "model": "Qwen/Qwen2.5-3B-Instruct",
 7   "messages": [{"role":"user","content":"Say hi in one sentence."}],
 8   "max_tokens": 32,
 9   "temperature": 0.0
10 }
11
12 r = requests.post(f"{BASE}/v1/chat/completions", json=payload, timeout=120)
13 print(r.status_code)
14 print(r.json()["choices"][0]["message"]["content"])
15
```

```
200
Hi there! How can I assist you today?
```

```
 1 # check the GPU consumption once the vLLM server is up
 2 ! nvidia-smi
```

```
Sun Jan 11 19:55:23 2026
+-----------------------------------------------------------------------------------------+
| NVIDIA-SMI 550.54.15              Driver Version: 550.54.15      CUDA Version: 12.4      |
|-----------------------------------------+------------------------+----------------------+
| GPU  Name                 Persistence-M | Bus-Id          Disp.A | Volatile Uncorr. ECC |
| Fan  Temp   Perf          Pwr:Usage/Cap |           Memory-Usage | GPU-Util  Compute M. |
|                                         |                        |               MIG M. |
|=========================================+========================+======================|
|   0  NVIDIA L4                      Off |   00000000:00:03.0 Off |                    0 |
| N/A   56C    P0               30W /  72W |   16225MiB /  23034MiB |      0%      Default |
|                                         |                        |                  N/A |
+-----------------------------------------+------------------------+----------------------+

+-----------------------------------------------------------------------------------------+
| Processes:                                                                              |
|  GPU   GI   CI        PID   Type   Process name                              GPU Memory |
|        ID   ID                                                               Usage      |
|=========================================================================================|
+-----------------------------------------------------------------------------------------+
```

## Load the dataset and run tests

```
 1 from google.colab import drive
 2 drive.mount("/content/drive")
 3
```

```
Drive already mounted at /content/drive; to attempt to forcibly remount, call drive.mount("/content/drive", force
```

## Load the dataset from Drive

```
 1 import json, os, random
 2 from pathlib import Path
 3
 4 PROMPT_DIR = Path("/content/drive/MyDrive/llm_loadtest_workloads")
 5
 6 FILES = [
 7     "short only isonl"
```

```
 7       short_only.jsonl",
 8       "medium_only.jsonl",
 9       "long_only.jsonl",
10       "mix_90_10.jsonl",
11       "mix_70_30.jsonl",
12       "mix_70_20_10.jsonl",
13       "mix_60_40.jsonl",
14       "mix_50_50.jsonl"
15 ]
16
17 def load_jsonl_prompts(path: Path, limit=None, seed=123):
18     prompts = []
19     with path.open("r", encoding="utf-8") as f:
20         for line in f:
21             obj = json.loads(line)
22
23             if "prompt" in obj:
24                 prompts.append(obj["prompt"])
25
26             elif "messages" in obj:
27                 users = [m["content"] for m in obj["messages"] if m.get("role") == "user"]
28
29                 if users:
30                     prompts.append(users[-1])
31
32             if limit and len(prompts) >= limit:
33                 break
34
35     random.Random(seed).shuffle(prompts)
36     return prompts
37
38 prompt_sets = {}
39 for fn in FILES:
40     p = PROMPT_DIR / fn
41     if p.exists():
42         prompt_sets[fn] = load_jsonl_prompts(p)
43         print(f"Loaded {fn}: {len(prompt_sets[fn])} prompts")
44     else:
45         print(f"Missing: {p}")
46
```

```
Loaded short_only.jsonl: 400 prompts
Loaded medium_only.jsonl: 400 prompts
Loaded long_only.jsonl: 400 prompts
Loaded mix_90_10.jsonl: 400 prompts
Loaded mix_70_30.jsonl: 400 prompts
Loaded mix_70_20_10.jsonl: 400 prompts
Loaded mix_60_40.jsonl: 400 prompts
Loaded mix_50_50.jsonl: 400 prompts
```

## ⌄ Checking if the dataset is properly read

```
 1 from transformers import AutoTokenizer
 2 import numpy as np, random
 3
 4 TOK_MODEL = "Qwen/Qwen2.5-3B-Instruct"
 5 tok = AutoTokenizer.from_pretrained(TOK_MODEL, use_fast=True)
 6
 7 def tok_len(s: str) -> int:
 8     return len(tok.encode(s, add_special_tokens=False))
 9
10 def summarize_lengths(name, prompts, n=200):
11     sample = random.sample(prompts, min(n, len(prompts)))
12     lens = np.array([tok_len(p) for p in sample])
13     print(f"{name:18s}  n={len(sample):4d}  min={lens.min():4d}  p50={int(np.median(lens)):4d}  p90={int(np.
14
15 for fn, prompts in prompt_sets.items():
16     summarize_lengths(fn, prompts)
17
```

```
short_only.jsonl     n= 200  min=   8  p50=  10  p90=  13  p99=  16  max=  16
medium_only.jsonl    n= 200  min= 128  p50= 166  p90= 230  p99= 246  max= 255
long_only.jsonl      n= 200  min=2051  p50=3196  p90=3949  p99=4062  max=4093
mix_90_10.jsonl      n= 200  min=   8  p50=  10  p90=2319  p99=3944  max=4002
mix_70_30.jsonl      n= 200  min=   8  p50=  11  p90=3312  p99=3924  max=4031
mix_70_20_10.jsonl   n= 200  min=   8  p50=  11  p90= 231  p99=3660  max=3833
mix_60_40.jsonl      n= 200  min=   8  p50=  12  p90=3970  p99=4082  max=4090
mix_50_50.jsonl      n= 200  min=   8  p50=3510  p90=4002  p99=4084  max=4096
```

## Experiment 1

In this experiment, we stress-test vLLM serving under mixed prompt-length workloads to reproduce two real-world issues:

(1) context-budget failures where long prompts combined with a fixed max_tokens cause requests to error out (HTTP 400 max token limit), and

(2) degraded serving performance as the fraction of long prompts increases. We run controlled workloads (short-only and mixed ratios like 90/10 and 70/30), then compare success rate, throughput (RPS), and latency metrics (p50/p90/p99) to quantify how long prompts amplify tail latency and reduce effective throughput.

This sets up the next experiments where we mitigate the failure mode (dynamic max_tokens / context budgeting) and run experiments to see how can we configure our servers to perform better under mixed loads.

## Single-Request Client Helper (Latency + Token Usage)

```python
1  import time, requests
2
3  BASE = "http://127.0.0.1:8000"
4  MODEL_ID = "Qwen/Qwen2.5-3B-Instruct"
5
6  def one_request(prompt, max_tokens=128, temperature=0.0, timeout=120):
7      t0 = time.perf_counter()
8      payload = {
9          "model": MODEL_ID,
10         "messages": [{"role": "user", "content": prompt}],
11         "max_tokens": max_tokens,
12         "temperature": temperature,
13     }
14
15     r = requests.post(f"{BASE}/v1/chat/completions", json=payload, timeout=timeout)
16
17     latency = time.perf_counter() - t0
18
19     # Hard fail with useful error
20     if r.status_code != 200:
21         txt = r.text[:500]
22         raise RuntimeError(f"HTTP {r.status_code}: {txt}")
23
24     j = r.json()
25     usage = j.get("usage", {})
26     return {
27         "latency_s": latency,
28         "prompt_tokens": usage.get("prompt_tokens"),
29         "completion_tokens": usage.get("completion_tokens"),
30         "total_tokens": usage.get("total_tokens"),
31         "status_code": r.status_code,
32     }
33
```

## Some checks before starting the experiment

```python
1  # Is the server process actually running?
2  !ps aux | grep -E "vllm\.entrypoints\.openai\.api_server|api_server" | grep -v grep
3
4  # Is anything listening on port 8000, and which process owns it?
5  !ss -ltnp | grep ":8000" || true
6
7  # What did the server log say most recently?
8  !tail -n 80 vllm_server.log || true
9
10
11 print(requests.get("http://127.0.0.1:8000/v1/models").status_code)
12 print(requests.get("http://127.0.0.1:8000/v1/models").json())
```

Show hidden output

## Server Health Check + Verify vLLM OpenAI Endpoint Is Working

```python
1  import requests, json
2
3  BASE = "http://127.0.0.1:8000"
4  MODEL_ID = "Qwen/Qwen2.5-3B-Instruct"
```

```
 5
 6 payload = {
 7     "model": MODEL_ID,
 8     "messages": [{"role": "user", "content": "Say hello in one sentence."}],
 9     "max_tokens": 32,
10     "temperature": 0.0,
11 }
12
13 r = requests.post(f"{BASE}/v1/chat/completions", json=payload, timeout=120)
14 print("status:", r.status_code)
15 print("text:", r.text[:400])
16
```

```
status: 200
text: {"id":"chatcmpl-bc2f4e8695ca898b","object":"chat.completion","created":1768163326,"model":"Qwen/Qwen2.5-3B-
```

## Single-Request Sanity Check (Real Workload Prompt)

```
1 test_prompt = prompt_sets["short_only.jsonl"][0]
2 print("prompt tokens (rough chars):", len(test_prompt))
3
4 out = one_request(test_prompt, max_tokens=32)
5 print(out)
6
```

```
prompt tokens (rough chars): 49
{'latency_s': 0.8831302409998898, 'prompt_tokens': 39, 'completion_tokens': 32, 'total_tokens': 71, 'status_code'
```

## Client-Side Load Generator (ThreadPool) + Latency/RPS Summary

```
 1 def run_load(prompts, concurrency=8, total_requests=200, max_tokens=128, desc="load"):
 2     rng = random.Random(123)
 3     picked = [rng.choice(prompts) for _ in range(total_requests)]
 4
 5     results = []
 6     t_start = time.perf_counter()
 7
 8     def task(i):
 9         p = picked[i]
10         try:
11             res = one_request(p, max_tokens=max_tokens)
12             res["ok"] = True
13         except Exception as e:
14             res = {"ok": False, "error": str(e), "latency_s": None,
15                     "prompt_tokens": None, "completion_tokens": None, "total_tokens": None}
16         return res
17
18     ok_ct = 0
19     fail_ct = 0
20
21     with cf.ThreadPoolExecutor(max_workers=concurrency) as ex:
22         futs = [ex.submit(task, i) for i in range(total_requests)]
23
24         pbar = tqdm(total=total_requests, desc=desc)
25         for f in cf.as_completed(futs):
26             r = f.result()
27             results.append(r)
28             if r.get("ok"):
29                 ok_ct += 1
30             else:
31                 fail_ct += 1
32             pbar.update(1)
33             pbar.set_postfix(ok=ok_ct, fail=fail_ct)
34         pbar.close()
35
36     t_end = time.perf_counter()
37
38     df = pd.DataFrame(results)
39     ok_df = df[df["ok"] == True].copy()
40
41     summary = {
42         "concurrency": concurrency,
43         "total_requests": total_requests,
44         "ok": int(ok_df.shape[0]),
45         "fail": int(df.shape[0] - ok_df.shape[0]),
46         "wall_time_s": t_end - t_start,
47         "rps": ok_df.shape[0] / (t_end - t_start) if (t_end - t_start) > 0 else None,
```

```
48        "p50_s": float(ok_df["latency_s"].quantile(0.50)) if len(ok_df) else None,
49        "p90_s": float(ok_df["latency_s"].quantile(0.90)) if len(ok_df) else None,
50        "p99_s": float(ok_df["latency_s"].quantile(0.99)) if len(ok_df) else None,
51        "mean_s": float(ok_df["latency_s"].mean()) if len(ok_df) else None,
52    }
53    return df, summary
54
```

## Reproducing Context-Length Failures Under Load (Max Token Limit)

```
1  workloads = [
2      ("short_only", "short_only.jsonl"),
3      ("mix_90_10", "mix_90_10.jsonl"),
4      ("mix_70_30", "mix_70_30.jsonl"),
5  ]
6
7  all_summaries = []
8  all_dfs = {}
9
10 for name, fn in tqdm(workloads, desc="Workloads"):
11     print("\nRunning:", name)
12     df, summ = run_load(
13         prompt_sets[fn],
14         concurrency=8,
15         total_requests=200,
16         max_tokens=128,
17         desc=f"{name} (c=8)"
18     )
19     all_dfs[name] = df
20     all_summaries.append({"workload": name, **summ})
21     print(summ)
22
23 summary_df = pd.DataFrame(all_summaries)
24 summary_df
25
```

```
Workloads: 100%                                          3/3 [04:25<00:00, 90.89s/it]

Running: short_only
short_only (c=8): 100%                        200/200 [01:19<00:00, 2.34it/s, fail=0, ok=200]
{'concurrency': 8, 'total_requests': 200, 'ok': 200, 'fail': 0, 'wall_time_s': 79.18587238500186, 'rps': 2.525703

Running: mix_90_10
mix_90_10 (c=8): 100%                         200/200 [01:26<00:00, 2.30it/s, fail=2, ok=198]
{'concurrency': 8, 'total_requests': 200, 'ok': 198, 'fail': 2, 'wall_time_s': 86.622785784999, 'rps': 2.28577271

Running: mix_70_30
mix_70_30 (c=8): 100%                         200/200 [01:39<00:00, 2.96it/s, fail=5, ok=195]
{'concurrency': 8, 'total_requests': 200, 'ok': 195, 'fail': 5, 'wall_time_s': 99.59230680399924, 'rps': 1.957982
```

|   | workload | concurrency | total_requests | ok | fail | wall_time_s | rps | p50_s | p90_s | p99_s | mean_s |
|---|----------|-------------|----------------|-----|------|-------------|----------|----------|----------|----------|----------|
| 0 | short_only | 8 | 200 | 200 | 0 | 79.185872 | 2.525703 | 3.553791 | 3.576180 | 3.586765 | 3.104981 |
| 1 | mix_90_10 | 8 | 200 | 198 | 2 | 86.622786 | 2.285773 | 3.615653 | 4.137271 | 4.528919 | 3.433756 |
| 2 | mix_70_30 | 8 | 200 | 195 | 5 | 99.592307 | 1.957983 | 4.101179 | 4.894588 | 5.376055 | 4.054638 |

```
1  for name, df in all_dfs.items():
2      print("\n===", name, "===")
3      print("rows:", len(df))
4      print("cols:", list(df.columns))
5
6      if "ok" in df.columns:
7          fails = df[df["ok"] == False]
8      elif "status" in df.columns:
9          fails = df[df["status"] != 200]
10     else:
11         print("I don't see ok/status columns — tell me df.columns and I'll adapt.")
12         continue
13
14     print("fail count:", len(fails))
15     if len(fails):
16         cols_to_show = [c for c in ["status", "error", "latency_s", "prompt_len", "prompt_tokens"] if c in fa
17         display(fails[cols_to_show].head(10))
18         if "error" in fails.columns:
19             print("\nTop error reasons:")
20             display(fails["error"].value_counts().head(10))
```

```
     21
```

```
=== short_only ===
rows: 200
cols: ['latency_s', 'prompt_tokens', 'completion_tokens', 'total_tokens', 'status_code', 'ok']
fail count: 0

=== mix_90_10 ===
rows: 200
cols: ['latency_s', 'prompt_tokens', 'completion_tokens', 'total_tokens', 'status_code', 'ok', 'error']
fail count: 2
```

|  | error | latency_s | prompt_tokens |
|---|---|---|---|
| 159 | HTTP 400: {"error":{"message":"'max_tokens' or... | NaN | NaN |
| 168 | HTTP 400: {"error":{"message":"'max_tokens' or... | NaN | NaN |

```
Top error reasons:
```

|  | count |
|---|---|
| **error** |  |
| HTTP 400: {"error":{"message":"'max_tokens' or 'max_completion_tokens' is too large: 128. This model's maximum context length is 4096 tokens and your request has 4030 input tokens (128 > 4096 - 4030). None","type":"BadRequestError","param":null,"code":400}} | 1 |
| HTTP 400: {"error":{"message":"'max_tokens' or 'max_completion_tokens' is too large: 128. This model's maximum context length is 4096 tokens and your request has 4040 input tokens (128 > 4096 - 4040). None","type":"BadRequestError","param":null,"code":400}} | 1 |

**dtype:** int64

```
=== mix_70_30 ===
rows: 200
cols: ['latency_s', 'prompt_tokens', 'completion_tokens', 'total_tokens', 'status_code', 'ok', 'error']
fail count: 5
```

|  | error | latency_s | prompt_tokens |
|---|---|---|---|
| 85 | HTTP 400: {"error":{"message":"'max_tokens' or... | NaN | NaN |
| 90 | HTTP 400: {"error":{"message":"'max_tokens' or... | NaN | NaN |
| 118 | HTTP 400: {"error":{"message":"'max_tokens' or... | NaN | NaN |
| 178 | HTTP 400: {"error":{"message":"'max_tokens' or... | NaN | NaN |
| 189 | HTTP 400: {"error":{"message":"'max_tokens' or... | NaN | NaN |

```
Top error reasons:
```

|  | count |
|---|---|
| **error** |  |
| HTTP 400: {"error":{"message":"'max_tokens' or 'max_completion_tokens' is too large: 128. This model's maximum context length is 4096 tokens and your request has 4040 input tokens (128 > 4096 - 4040). None","type":"BadRequestError","param":null,"code":400}} | 2 |
| HTTP 400: {"error":{"message":"'max_tokens' or 'max_completion_tokens' is too large: 128. This model's maximum context length is 4096 tokens and your request has 3984 input tokens (128 > 4096 - 3984). None","type":"BadRequestError","param":null,"code":400}} | 1 |
| HTTP 400: {"error":{"message":"'max_tokens' or 'max_completion_tokens' is too large: 128. This model's maximum context length is 4096 tokens and your request has 4015 input tokens (128 > 4096 - 4015). None","type":"BadRequestError","param":null,"code":400}} | 1 |

## Conclusion

In the mixed workloads (mix_90_10 and mix_70_30), we start seeing HTTP 400 failures even though the same code works fine for short-only prompts. The error message shows why: the model has a fixed maximum context length of 4096 tokens, and some of our long prompts already consume ~3980–4060 input tokens.

When we blindly request a fixed max_tokens=128 for completion, the request becomes impossible because the remaining context budget (4096 - input_tokens) is smaller than 128. vLLM therefore rejects those requests with a BadRequest error (e.g., "128 > 4096 – 4040").

This confirms a real serving issue under mixed prompt lengths: static max_tokens can cause hard failures whenever inputs approach the context window, motivating the next experiment where we dynamically cap max_tokens based on the available token budget per request.

```
1 import matplotlib.pyplot as plt
2 import numpy as np
3
```

```
 4
 5 labels = summary_df["workload"].tolist()
 6
 7 # p50/p90/p99
 8 p50 = summary_df["p50_s"].values
 9 p90 = summary_df["p90_s"].values
10 p99 = summary_df["p99_s"].values
11
12 x = np.arange(len(labels))
13 w = 0.25
14
15 plt.figure()
16 plt.bar(x - w, p50, width=w, label="p50")
17 plt.bar(x,     p90, width=w, label="p90")
18 plt.bar(x + w, p99, width=w, label="p99")
19 plt.xticks(x, labels, rotation=20)
20 plt.ylabel("Latency (s)")
21 plt.title("Latency percentiles by workload")
22 plt.legend()
23 plt.show()
24
25 # RPS
26 plt.figure()
27 plt.bar(labels, summary_df["rps"].values)
28 plt.ylabel("Requests/sec")
29 plt.title("Throughput (RPS) by workload")
30 plt.xticks(rotation=20)
31 plt.show()
32
```





Plot Interpretation

These two plots summarize the impact of prompt-length mix on both latency tail and throughput at a fixed concurrency (c=8, 200 total requests each).

As we move from short_only -> mix_90_10 -> mix_70_30, the proportion of long prompts increases, which pushes up p90/p99 latency and reduces overall RPS. This is the classic "mixed-length serving" effect.

Long prompts spend much more time in prefill and consume more KV-cache and batching budget, so they slow down batch completion and increase queueing for everyone, which shows up most clearly in tail latency.

The throughput chart reinforces this: more long requests means fewer requests completed per second. Also note that the latency percentiles here are computed over successful (ok) requests.

The few failed requests from long prompts are not included in the percentile bars, and those failures are exactly what we address next with the max_tokens safety and dynamic max_tokens fix.

```python
1  for name, df in all_dfs.items():
2      if "latency_s" in df.columns:
3          vals = df[df["latency_s"].notna()]["latency_s"].values
4          plt.figure()
5          plt.hist(vals, bins=30)
6          plt.title(f"Latency histogram: {name}")
7          plt.xlabel("seconds")
8          plt.ylabel("count")
9          plt.show()
10
```

**Latency histogram: short_only**

**Latency histogram: mix_90_10**

**Latency histogram: mix_70_30**

## Plot Interpretation

These histograms show the full latency distribution for each workload, not just the percentile summaries. The short_only workload is tightly clustered, which indicates stable, predictable latency when all prompts are small.

As we introduce longer prompts in mix_90_10 and mix_70_30, the distribution shifts to the right and spreads out, meaning higher average latency and greater variability. You can also see a heavier "right tail" in the mixed workloads, which is exactly what drives the p90/p99 increases we saw earlier.
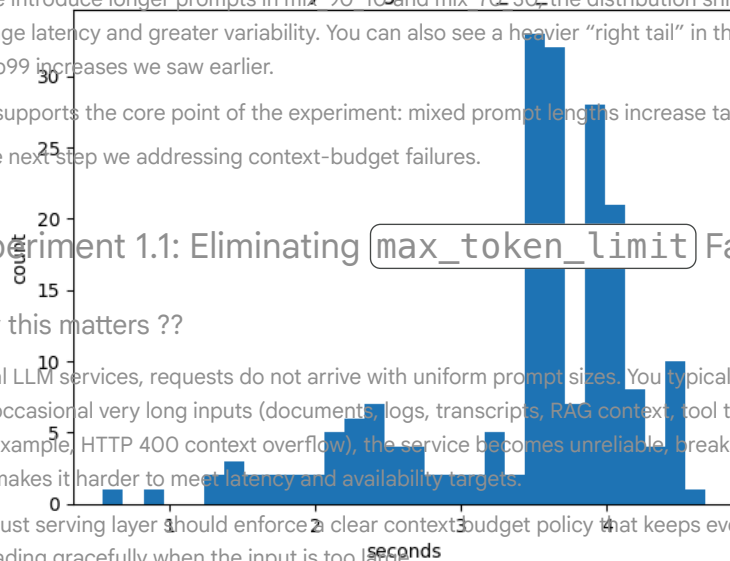
This supports the core point of the experiment: mixed prompt lengths increase tail latency and make serving behavior less predictable.

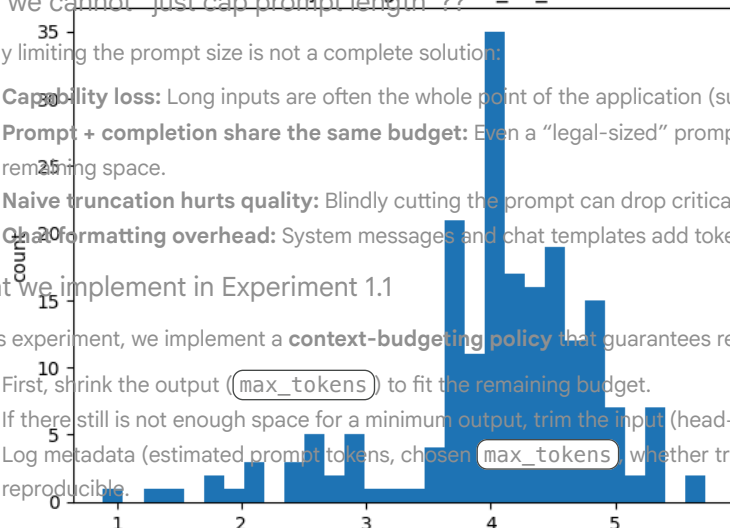In the next step we addressing context-budget failures.

## Experiment 1.1: Eliminating `max_token_limit` Failures with Context-Budgeting

### Why this matters ??

In real LLM services, requests do not arrive with uniform prompt sizes. You typically see a mix of short questions, medium chat turns, and occasional very long inputs (documents, logs, transcripts, RAG context, tool traces). If the system sometimes returns hard failures (for example, HTTP 400 context overflow), the service becomes unreliable, breaks downstream pipelines (agents, batch jobs, retries), and makes it harder to meet latency and availability targets.

A robust serving layer should enforce a clear context-budget policy that keeps every request within the model's context window, while degrading gracefully when the input is too large.

### Why we cannot "just cap prompt length" ??

Simply limiting the prompt size is not a complete solution:

- **Capability loss:** Long inputs are often the whole point of the application (summarization, document Q&A, log diagnosis).
- **Prompt + completion share the same budget:** Even a "legal-sized" prompt can overflow if `max_tokens` is too large for the remaining space.
- **Naive truncation hurts quality:** Blindly cutting the prompt can drop critical context and silently degrade answer quality.
- **Chat formatting overhead:** System messages and chat templates add tokens that a simple character cap does not account for.

### What we implement in Experiment 1.1

In this experiment, we implement a **context-budgeting policy** that guarantees requests fit the model's context window:

- First, shrink the output (`max_tokens`) to fit the remaining budget.
- If there still is not enough space for a minimum output, trim the input (head+tail) to create room.
- Log metadata (estimated prompt tokens, chosen `max_tokens`, whether trimming occurred) so the behavior is measurable and reproducible.

```
1 # clean out previous instances of vLLM if you want to start a new server
2
3 !pkill -9 -f "vllm.entrypoints.openai.api_server" || true
4 !pkill -9 -f "EngineCore" || true
5 !pkill -9 -f "vllm" || true
6 !sleep 2
7 !nvidia-smi
^C
^C
^C
```

```
Sun Jan 11 22:01:00 2026
+-----------------------------------------------------------------------------------------+
| NVIDIA-SMI 550.54.15              Driver Version: 550.54.15      CUDA Version: 12.4      |
|-----------------------------------------+------------------------+----------------------+
| GPU  Name                 Persistence-M | Bus-Id          Disp.A | Volatile Uncorr. ECC |
| Fan  Temp   Perf          Pwr:Usage/Cap |           Memory-Usage | GPU-Util  Compute M. |
|                                         |                        |               MIG M. |
|=========================================+========================+======================|
|   0  NVIDIA L4                      Off |   00000000:00:03.0 Off |                    0 |
| N/A   54C    P8             18W /   72W |     3MiB /  23034MiB |      0%      Default |
|                                         |                        |                  N/A |
+-----------------------------------------+------------------------+----------------------+

+-----------------------------------------------------------------------------------------+
| Processes:                                                                              |
|  GPU   GI   CI        PID   Type   Process name                              GPU Memory |
|        ID   ID                                                               Usage      |
|=========================================================================================|
|  No running processes found                                                             |
+-----------------------------------------------------------------------------------------+
```

```
 1 !lsof -i :8000 || true
 2
 3
```

```
COMMAND  PID USER   FD   TYPE DEVICE SIZE/OFF NODE NAME
python3 1144 root   79u  IPv4 606432      0t0  TCP localhost:49490->localhost:8000 (CLOSE_WAIT)
```

```
 1
 2 MODEL_ID = "Qwen/Qwen2.5-3B-Instruct"
 3 PORT = 8000
 4
 5 MAX_MODEL_LEN = 4096
 6 GPU_UTIL = 0.70
 7 DTYPE = "half"
 8
 9 cmd = f"""
10 nohup python -m vllm.entrypoints.openai.api_server \
11   --model {MODEL_ID} \
12   --host 127.0.0.1 \
13   --port {PORT} \
14   --max-model-len {MAX_MODEL_LEN} \
15   --gpu-memory-utilization {GPU_UTIL} \
16   --dtype {DTYPE} \
17   --disable-log-stats \
18   > vllm_server.log 2>&1 &
19 """
20
21 print("Launching vLLM server...")
22 !bash -lc "$cmd"
23 !tail -n 30 vllm_server.log
24
```

```
Launching vLLM server...
```

## Function to check the successful startup of the server.

```
 1 BASE = "http://127.0.0.1:8000"
 2
 3 def wait_for_server(timeout_s=600):
 4     t0 = time.time()
 5     while time.time() - t0 < timeout_s:
 6         try:
 7             r = requests.get(f"{BASE}/v1/models", timeout=2)
 8             if r.status_code == 200:
 9                 return True, r.json()
10         except Exception:
11             pass
12         time.sleep(2)
13     return False, None
14
15 ok, payload = wait_for_server()
16 print("Ready:", ok)
17 if ok:
18     print("Models:", [m["id"] for m in payload.get("data", [])])
19 else:
20     print("Not ready yet. Check logs:")
21     !tail -n 60 vllm_server.log
22
```

```
Ready: True
Models: ['Qwen/Qwen2.5-3B-Instruct']
```

## ˅  Dynamic Output Budgeting and Prompt Trimming to Avoid Context Overflow

```
 1 SAFETY_MARGIN = 16          # keep small headroom to avoid borderline failures
 2 DESIRED_MAX_TOKENS = 128    # same as Experiment 1
 3 MIN_OUTPUT_TOKENS = 32      # minimum output you want to guarantee
 4
 5 def count_prompt_tokens(messages):
 6     """
 7     Counts tokens for the full chat-formatted prompt that will be sent to vLLM.
 8     Qwen tokenizer supports apply_chat_template().
 9     """
10     ids = tok.apply_chat_template(
11         messages,
12         tokenize=True,
13         add_generation_prompt=True,
14         return_tensors=None
15     )
16     return len(ids)
17
18 def head_tail_trim_text_by_tokens(text, target_tokens, head_frac=0.35):
19     """
20     Keep some head + some tail tokens, drop middle.
21     """
22     ids = tok.encode(text, add_special_tokens=False)
23     if len(ids) <= target_tokens:
24         return text, False
25
26     head_tokens = int(target_tokens * head_frac)
27     tail_tokens = target_tokens - head_tokens
28
29     trimmed_ids = ids[:head_tokens] + ids[-tail_tokens:]
30     return tok.decode(trimmed_ids), True
31
32 def build_payload_with_budget(prompt_text,
33                               desired_max_tokens=DESIRED_MAX_TOKENS,
34                               min_output_tokens=MIN_OUTPUT_TOKENS,
35                               safety_margin=SAFETY_MARGIN,
36                               temperature=0.0):
37     """
38     Returns (payload, meta).
39     - Ensures no context overflow.
40     - If needed, trims user prompt (head+tail) to make room for min output tokens.
41     """
42
43     messages = [{"role": "user", "content": prompt_text}]
44     ptoks = count_prompt_tokens(messages)
45
46     # First attempt: just shrink output tokens if needed
47     max_out = min(desired_max_tokens, MAX_MODEL_LEN - safety_margin - ptoks)
48
49     trimmed = False
50     original_ptoks = ptoks
51     original_prompt = prompt_text
52
53     # If not enough room even for min output -> trim input to create room
54     if max_out < min_output_tokens:
55         target_prompt_budget = MAX_MODEL_LEN - safety_margin - min_output_tokens
56         if target_prompt_budget <= 0:
57             raise RuntimeError("Model context too small after safety margin.")
58
59         # We need to reduce prompt tokens from ptoks down to target_prompt_budget.
60         # We'll trim only the USER CONTENT tokens.
61         #
62         # Use a simple search on how many user-content tokens we can keep.
63         # (Because chat template adds overhead tokens.)
64         lo, hi = 0, len(tok.encode(prompt_text, add_special_tokens=False))
65         best_text = ""
66         best_ptoks = None
67
68         while lo <= hi:
69             mid = (lo + hi) // 2
70             candidate_text, _ = head_tail_trim_text_by_tokens(prompt_text, mid)
71             candidate_msgs = [{"role": "user", "content": candidate_text}]
72             candidate_ptoks = count_prompt_tokens(candidate_msgs)
73
74             if candidate_ptoks <= target_prompt_budget:
75                 best_text = candidate_text
```

```
 76                    best_ptoks = candidate_ptoks
 77                    lo = mid + 1
 78                else:
 79                    hi = mid - 1
 80
 81          if best_ptoks is None:
 82              raise RuntimeError(
 83                  f"Prompt too long to fit even after trimming. "
 84                  f"prompt_tokens={ptoks}, budget={target_prompt_budget}"
 85              )
 86
 87          prompt_text = best_text
 88          trimmed = True
 89          messages = [{"role": "user", "content": prompt_text}]
 90          ptoks = best_ptoks
 91
 92          max_out = min(desired_max_tokens, MAX_MODEL_LEN - safety_margin - ptoks)
 93
 94      # Final guard
 95      if max_out <= 0:
 96          raise RuntimeError(
 97              f"No room for output tokens. prompt_tokens={ptoks}, "
 98              f"max_len={MAX_MODEL_LEN}, safety_margin={safety_margin}"
 99          )
100
101      payload = {
102          "model": MODEL_ID,
103          "messages": messages,
104          "max_tokens": int(max_out),
105          "temperature": temperature,
106      }
107
108      meta = {
109          "prompt_tokens_est": ptoks,
110          "max_tokens_used": int(max_out),
111          "trimmed": trimmed,
112          "prompt_tokens_est_original": original_ptoks,
113          "original_prompt_len_chars": len(original_prompt),
114          "final_prompt_len_chars": len(prompt_text),
115      }
116      return payload, meta
117
```

## Request Execution with Dynamic Token Budgeting

```
 1
 2
 3 def one_request(prompt, desired_max_tokens=128, temperature=0.0, timeout=120):
 4     t0 = time.perf_counter()
 5
 6     payload, meta = build_payload_with_budget(
 7         prompt_text=prompt,
 8         desired_max_tokens=desired_max_tokens,
 9         min_output_tokens=MIN_OUTPUT_TOKENS,
10         safety_margin=SAFETY_MARGIN,
11         temperature=temperature,
12     )
13
14     r = requests.post(f"{BASE}/v1/chat/completions", json=payload, timeout=timeout)
15     latency = time.perf_counter() - t0
16
17     if r.status_code != 200:
18         raise RuntimeError(f"HTTP {r.status_code}: {r.text[:800]}")
19
20     j = r.json()
21     usage = j.get("usage", {})
22
23     return {
24         "latency_s": latency,
25         "prompt_tokens": usage.get("prompt_tokens"),
26         "completion_tokens": usage.get("completion_tokens"),
27         "total_tokens": usage.get("total_tokens"),
28         "status_code": r.status_code,
29         # extra debugging info (useful for your writeup)
30         **meta
31     }
32
```

```
1  # smoke test of one_request function
2
3  test_prompt = "Explain in 2 sentences what chunked prefill does in vLLM."
4  out = one_request(test_prompt, desired_max_tokens=128, temperature=0.0, timeout=120)
5  out
6
```

```
{'latency_s': 1.3368042790007166,
 'prompt_tokens': 46,
 'completion_tokens': 49,
 'total_tokens': 95,
 'status_code': 200,
 'prompt_tokens_est': 46,
 'max_tokens_used': 128,
 'trimmed': False,
 'prompt_tokens_est_original': 46,
 'original_prompt_len_chars': 57,
 'final_prompt_len_chars': 57}
```

## Load Test Harness Using Budget-Aware Requests (Dynamic max_tokens + Optional Prompt Trimming)

```
1  def run_load(prompts, concurrency=8, total_requests=200, desired_max_tokens=128, desc="load"):
2      rng = random.Random(123)
3      picked = [rng.choice(prompts) for _ in range(total_requests)]
4
5      results = []
6      t_start = time.perf_counter()
7
8      def task(i):
9          p = picked[i]
10         try:
11             res = one_request(p, desired_max_tokens=desired_max_tokens, temperature=0.0, timeout=180)
12             res["ok"] = True
13         except Exception as e:
14             res = {"ok": False, "error": str(e), "latency_s": None,
15                    "prompt_tokens": None, "completion_tokens": None, "total_tokens": None}
16         return res
17
18     ok_ct = 0
19     fail_ct = 0
20
21     with cf.ThreadPoolExecutor(max_workers=concurrency) as ex:
22         futs = [ex.submit(task, i) for i in range(total_requests)]
23         pbar = tqdm(total=total_requests, desc=desc)
24         for f in cf.as_completed(futs):
25             r = f.result()
26             results.append(r)
27             if r.get("ok"):
28                 ok_ct += 1
29             else:
30                 fail_ct += 1
31             pbar.update(1)
32             pbar.set_postfix(ok=ok_ct, fail=fail_ct)
33         pbar.close()
34
35     t_end = time.perf_counter()
36
37     df = pd.DataFrame(results)
38     ok_df = df[df["ok"] == True].copy()
39
40     summary = {
41         "concurrency": concurrency,
42         "total_requests": total_requests,
43         "ok": int(ok_df.shape[0]),
44         "fail": int(df.shape[0] - ok_df.shape[0]),
45         "wall_time_s": t_end - t_start,
46         "rps": ok_df.shape[0] / (t_end - t_start) if (t_end - t_start) > 0 else None,
47         "p50_s": float(ok_df["latency_s"].quantile(0.50)) if len(ok_df) else None,
48         "p90_s": float(ok_df["latency_s"].quantile(0.90)) if len(ok_df) else None,
49         "p99_s": float(ok_df["latency_s"].quantile(0.99)) if len(ok_df) else None,
50         "mean_s": float(ok_df["latency_s"].mean()) if len(ok_df) else None,
51     }
52     return df, summary
53
54
55
```

## Run Safe (Budget-Aware) Load Tests Across Workloads

```
1 workloads = [
2     ("short_only", "short_only.jsonl"),
3     ("mix_90_10", "mix_90_10.jsonl"),
4     ("mix_70_30", "mix_70_30.jsonl"),
5 ]
6
7
8 all_summaries = []
9 all_dfs = {}
10
11 for name, fn in tqdm(workloads, desc="Workloads (safe)"):
12     print("\nRunning:", name)
13     df, summ = run_load(
14         prompt_sets[fn],
15         concurrency=8,
16         total_requests=200,
17         desired_max_tokens=128,
18         desc=f"{name} (c=8, safe)"
19     )
20     all_dfs[name] = df
21     all_summaries.append({"workload": name, **summ})
22     print(summ)
23
24 summary_df_safe = pd.DataFrame(all_summaries)
25 summary_df_safe
```

```
Workloads (safe): 100%                                    3/3 [04:30<00:00, 92.91s/it]

Running: short_only
short_only (c=8, safe): 100%                    200/200 [01:19<00:00, 2.30it/s, fail=0, ok=200]
{'concurrency': 8, 'total_requests': 200, 'ok': 200, 'fail': 0, 'wall_time_s': 79.25956869500078, 'rps': 2.523354

Running: mix_90_10
mix_90_10 (c=8, safe): 100%                     200/200 [01:28<00:00, 2.34it/s, fail=0, ok=200]
{'concurrency': 8, 'total_requests': 200, 'ok': 200, 'fail': 0, 'wall_time_s': 88.17387503600003, 'rps': 2.268245

Running: mix_70_30
mix_70_30 (c=8, safe): 100%                     200/200 [01:42<00:00, 3.26it/s, fail=0, ok=200]
{'concurrency': 8, 'total_requests': 200, 'ok': 200, 'fail': 0, 'wall_time_s': 102.89451057299993, 'rps': 1.94373
```

|   | workload | concurrency | total_requests | ok | fail | wall_time_s | rps | p50_s | p90_s | p99_s | mean_s |
|---|----------|-------------|----------------|-----|------|-------------|-----|-------|-------|-------|--------|
| 0 | short_only | 8 | 200 | 200 | 0 | 79.259569 | 2.523355 | 3.556058 | 3.573445 | 3.581660 | 3.106944 |
| 1 | mix_90_10 | 8 | 200 | 200 | 0 | 88.173875 | 2.268246 | 3.640731 | 4.242185 | 4.532153 | 3.469050 |
| 2 | mix_70_30 | 8 | 200 | 200 | 0 | 102.894511 | 1.943738 | 4.175408 | 4.893263 | 5.393245 | 4.073671 |

## Conclusion

In Experiment 1.1 (the budget-aware request builder), the key improvement is reliability: all three workloads now complete with 0 failures (200/200 OK), including mixes that previously hit HTTP 400 context overflow when long prompts left insufficient room for the requested completion (max_tokens).

This confirms the fix works as intended, because each request now dynamically adjusts max_tokens to fit the remaining context window and trims the prompt only when necessary to guarantee a minimum output budget.

Performance trends remain consistent with Experiment 1: as the proportion of long prompts increases (short_only -> mix_90_10 -> mix_70_30), latency rises and throughput (RPS) drops due to heavier prefill and decoding work, but crucially this happens without any request-level hard failures, making the system robust under mixed prompt lengths.

## Experiment 2: Diagnosing Tail Latency with a "Bad" vLLM Batching Config

```
1 !pkill -f "vllm.entrypoints.openai.api_server" || true
2 !pkill -f "VLLM::EngineCore" || true
3 !pkill -f "EngineCore" || true
4 time.sleep(2)
5
6 !rm -f vllm_bad.log
7 print("killed + log removed")
8
```

```
^C
^C
^C
killed + log removed
```

## Creating a Bad Config for LLM Server

To replicate a **"chunked prefill effectively OFF"** situation

**Why this is a "bad" config:** We set `--max-num-batched-tokens 8192` (very large) while leaving `--max-num-seqs`
unconstrained (defaults can be high). This encourages vLLM to pack big prefills into a batch, which can let long prompts dominate
compute and delay short requests.

**Chunked prefill caveat:** Newer vLLM versions have chunked prefill **enabled by default**, but its *practical* benefit depends on admission
and batching knobs. If `max-num-batched-tokens` is too permissive (and/or `max-num-seqs` allows over-admission), long prefills
can still monopolize batch capacity, effectively **diminishing chunked prefill's ability to interleave work and protect short requests**.

**Head-of-line blocking risk:** With mixed workloads (50/50 short and long), long-prefill requests can occupy most of the batch token
budget, causing short prompts to wait behind long prefills, inflating **p90/p99** latency.

```
 1 !nohup python3 -m vllm.entrypoints.openai.api_server \
 2   --model Qwen/Qwen2.5-3B-Instruct \
 3   --host 127.0.0.1 --port 8000 \
 4   --dtype half \
 5   --max-model-len 4096 \
 6   --gpu-memory-utilization 0.70 \
 7   --disable-log-stats \
 8   --enforce-eager \
 9   --max-num-batched-tokens 8192 \
10   --long-prefill-token-threshold 1024 \
11   > vllm_bad.log 2>&1 </dev/null &
12
```

```
 1 BASE = "http://127.0.0.1:8000"
 2
 3 def wait_for_server(timeout_s=600, poll_s=2):
 4     t0 = time.time()
 5     last_err = None
 6     while time.time() - t0 < timeout_s:
 7         try:
 8             r = requests.get(f"{BASE}/v1/models", timeout=2)
 9             if r.status_code == 200:
10                 return True, r.json()
11             last_err = f"HTTP {r.status_code}: {r.text[:200]}"
12         except Exception as e:
13             last_err = f"{type(e).__name__}: {e}"
14         time.sleep(poll_s)
15     return False, last_err
16
17 ok, models_or_err = wait_for_server()
18 print("Ready:", ok)
19
20 if ok:
21     print("Models:", [m.get("id") for m in models_or_err.get("data", [])])
22 else:
23     print("Not ready yet. Last error:", models_or_err)
24     print("Check logs:")
25     !tail -n 80 vllm_bad.log
26
```

```
Ready: True
Models: ['Qwen/Qwen2.5-3B-Instruct']
```

```
 1 !tail -n 60 vllm_bad.log
```

Show hidden output

## Load Tokenizer

```python
1  import requests
2  from transformers import AutoTokenizer
3
4  BASE = "http://127.0.0.1:8000"
5  MODEL_ID = "Qwen/Qwen2.5-3B-Instruct"
6
7  tok = AutoTokenizer.from_pretrained(MODEL_ID, use_fast=True)
8
9  models = requests.get(f"{BASE}/v1/models").json()
10 model_info = next(m for m in models["data"] if m["id"] == MODEL_ID)
11 MAX_MODEL_LEN = int(model_info.get("max_model_len", 4096))
12
13 print("tok loaded")
14 print("MAX_MODEL_LEN from server:", MAX_MODEL_LEN)
15
```

```
tok loaded
MAX_MODEL_LEN from server: 4096
```

## Run Load Test on Bad Server Config

```python
1  import pandas as pd
2
3  WL_FILE = "mix_50_50.jsonl"
4  CONCURRENCY = 32
5  TOTAL_REQ = 200
6  DESIRED_MAX_TOKENS = 16
7
8
9  df_bad, summ_bad = run_load(
10     prompt_sets[WL_FILE],
11     concurrency=CONCURRENCY,
12     total_requests=TOTAL_REQ,
13     desired_max_tokens=DESIRED_MAX_TOKENS,
14     desc=f"BAD_unfair | {WL_FILE} | c={CONCURRENCY}"
15 )
16
17
18 print("BAD summary:", summ_bad)
19 df_bad.head()
20
```

```
BAD_unfair | mix_50_50.jsonl | c=32: 100%                              200/200 [00:42<00:00,  9.61it/s, fail=0, ok=200]
BAD summary: {'concurrency': 32, 'total_requests': 200, 'ok': 200, 'fail': 0, 'wall_time_s': 42.19113462899986, '
```

|   | latency_s | prompt_tokens | completion_tokens | total_tokens | status_code | prompt_tokens_est | max_tokens_used | trim |
|---|-----------|---------------|-------------------|--------------|-------------|-------------------|-----------------|------|
| 0 | 8.678932  | 37            | 16                | 53           | 200         | 37                | 16              |      |
| 1 | 8.703615  | 40            | 16                | 56           | 200         | 40                | 16              |      |
| 2 | 8.683890  | 42            | 16                | 58           | 200         | 42                | 16              |      |
| 3 | 8.700286  | 37            | 16                | 53           | 200         | 37                | 16              |      |
| 4 | 8.697547  | 37            | 16                | 53           | 200         | 37                | 16              |      |

## Compute Overall and Short-Prompt Tail Latency Metrics (p50/p90/p99)

```python
1  def tail_metrics(df, label, short_cutoff_tokens=64):
2      ok = df[df["ok"] == True].copy()
3
4      # overall
5      overall = {
6          "label": label,
7          "subset": "overall",
8          "n": len(ok),
9          "p50": float(ok["latency_s"].quantile(0.50)),
10         "p90": float(ok["latency_s"].quantile(0.90)),
11         "p99": float(ok["latency_s"].quantile(0.99)),
12         "mean": float(ok["latency_s"].mean()),
13     }
14
15     # short-only subset (inside the mixed workload)
```

```
16    short = ok[ok["prompt_tokens"] <= short_cutoff_tokens].copy()
17    if len(short) == 0:
18        short_stats = {"label": label, "subset": f"short<= {short_cutoff_tokens}", "n": 0,
19                       "p50": None, "p90": None, "p99": None, "mean": None}
20    else:
21        short_stats = {
22            "label": label,
23            "subset": f"short<= {short_cutoff_tokens}",
24            "n": len(short),
25            "p50": float(short["latency_s"].quantile(0.50)),
26            "p90": float(short["latency_s"].quantile(0.90)),
27            "p99": float(short["latency_s"].quantile(0.99)),
28            "mean": float(short["latency_s"].mean()),
29        }
30
31    return overall, short_stats
32
```

```
1 overall_bad, short_bad = tail_metrics(df_bad, "BAD_unfair", short_cutoff_tokens=64)
2 overall_bad, short_bad
3
```

```
({'label': 'BAD_unfair',
  'subset': 'overall',
  'n': 200,
  'p50': 6.181310366500156,
  'p90': 10.99196271639994,
  'p99': 11.637151490029915,
  'mean': 6.493085269189996},
 {'label': 'BAD_unfair',
  'subset': 'short<= 64',
  'n': 108,
  'p50': 5.480035559000271,
  'p90': 8.717383960199822,
  'p99': 10.618033962269847,
  'mean': 5.465330376453692})
```

## ∨ Expermient 3: Tail latency and Throughput variation with Concurrency Sweep

This sweep isolates scheduling effects by holding workload and max batched token length constant, then measuring how tail latency and throughput evolve with concurrency.

```
1 from google.colab import drive
2 drive.mount("/content/drive")
3
4 import os, json, shutil, time
5 import pandas as pd
6
7 WL_FILE = "mix_50_50.jsonl"
8 TOTAL_REQ = 200
9 DESIRED_MAX_TOKENS = 16
10 CONCURRENCY_LIST = [1, 2, 4, 8, 16, 32, 48, 64]
11
12 SHORT_CUTOFF = 64
13 LONG_CUTOFF  = 2048
14
15 LABEL = "BAD_unfair"
16
17 RUN_DIR = "/content/drive/MyDrive/vllm_sweeps/exp2_bad_unfair"
18 os.makedirs(RUN_DIR, exist_ok=True)
19
20 CKPT_CSV      = os.path.join(RUN_DIR, "concurrency_checkpoint.csv")
21 MANIFEST_JSON = os.path.join(RUN_DIR, "concurrency_manifest.json")
22
23 SAVE_RAW_DF = False
24 RAW_DIR = os.path.join(RUN_DIR, "raw_dfs")
25 if SAVE_RAW_DF:
26    os.makedirs(RAW_DIR, exist_ok=True)
27
28 manifest = {
29    "label": LABEL,
30    "workload_file": WL_FILE,
31    "total_requests": TOTAL_REQ,
32    "desired_max_tokens": DESIRED_MAX_TOKENS,
33    "concurrency_list": CONCURRENCY_LIST,
34    "short_cutoff": SHORT_CUTOFF,
35    "long_cutoff": LONG_CUTOFF,
36 }
37
38 if os.path.exists(MANIFEST_JSON):
```

```python
39     try:
40         old = json.loads(open(MANIFEST_JSON, "r").read())
41         if old != manifest:
42             print("Manifest differs from existing checkpoint. You may be mixing runs.")
43             print("Existing:", old)
44             print("New:", manifest)
45     except Exception:
46         print("Could not read existing manifest. Proceeding.")
47 else:
48     with open(MANIFEST_JSON, "w") as f:
49         json.dump(manifest, f, indent=2)
50
51 def summarize_run(df, summ, label, workload, concurrency):
52     ok = df[df["ok"] == True].copy()
53
54     row = {
55         "label": label,
56         "workload": workload,
57         "concurrency": int(concurrency),
58
59         "ok": summ.get("ok"),
60         "fail": summ.get("fail"),
61         "rps": summ.get("rps"),
62         "wall_time_s": summ.get("wall_time_s"),
63         "p50_s": summ.get("p50_s"),
64         "p90_s": summ.get("p90_s"),
65         "p99_s": summ.get("p99_s"),
66         "mean_s": summ.get("mean_s"),
67     }
68
69     short = ok[ok["prompt_tokens"] <= SHORT_CUTOFF]
70     row["short_n"] = int(len(short))
71     row["short_p50_s"] = float(short["latency_s"].quantile(0.50)) if len(short) else None
72     row["short_p90_s"] = float(short["latency_s"].quantile(0.90)) if len(short) else None
73     row["short_p99_s"] = float(short["latency_s"].quantile(0.99)) if len(short) else None
74
75     longg = ok[ok["prompt_tokens"] >= LONG_CUTOFF]
76     row["long_n"] = int(len(longg))
77     row["long_p50_s"] = float(longg["latency_s"].quantile(0.50)) if len(longg) else None
78     row["long_p90_s"] = float(longg["latency_s"].quantile(0.90)) if len(longg) else None
79     row["long_p99_s"] = float(longg["latency_s"].quantile(0.99)) if len(longg) else None
80
81     return row
82
83 def load_ckpt_df(path):
84     if not os.path.exists(path):
85         return pd.DataFrame()
86     df = pd.read_csv(path)
87     if "concurrency" in df.columns:
88         df["concurrency"] = pd.to_numeric(df["concurrency"], errors="coerce").astype("Int64")
89     return df
90
91 def done_concurrencies(df_ckpt):
92     if df_ckpt.empty:
93         return set()
94     good = df_ckpt[df_ckpt["concurrency"].notna()].copy()
95     if "status" in good.columns:
96         good = good[good["status"].isna()].copy()
97     return set(int(x) for x in good["concurrency"].dropna().tolist())
98
99 def append_ckpt_row(row, path):
100     df_row = pd.DataFrame([row])
101     header = not os.path.exists(path)
102     df_row.to_csv(path, mode="a", header=header, index=False)
103
104 def maybe_save_raw_df(df, concurrency):
105     if not SAVE_RAW_DF:
106         return
107     outp = os.path.join(RAW_DIR, f"raw_c{int(concurrency)}.parquet")
108     try:
109         df.to_parquet(outp, index=False)
110     except Exception:
111         outp = os.path.join(RAW_DIR, f"raw_c{int(concurrency)}.csv")
112         df.to_csv(outp, index=False)
113
114 ckpt_df = load_ckpt_df(CKPT_CSV)
115 done = done_concurrencies(ckpt_df)
116 print("Loaded checkpoint rows:", len(ckpt_df))
117 print("Completed concurrencies:", sorted(done))
118
119 for c in CONCURRENCY_LIST:
120     c = int(c)
```

```python
121        if c in done:
122            print("Skipping:", c)
123            continue
124
125        print("Running:", LABEL, WL_FILE, "c=", c)
126
127        try:
128            df, summ = run_load(
129                prompt_sets[WL_FILE],
130                concurrency=c,
131                total_requests=TOTAL_REQ,
132                desired_max_tokens=DESIRED_MAX_TOKENS,
133                desc=f"{LABEL} | {WL_FILE} | c={c}",
134            )
135
136            row = summarize_run(df, summ, label=LABEL, workload=WL_FILE, concurrency=c)
137            append_ckpt_row(row, CKPT_CSV)
138            maybe_save_raw_df(df, concurrency=c)
139
140            ckpt_df = load_ckpt_df(CKPT_CSV)
141            done = done_concurrencies(ckpt_df)
142
143        except Exception as e:
144            fail_row = {
145                "label": LABEL,
146                "workload": WL_FILE,
147                "concurrency": c,
148                "ok": None,
149                "fail": None,
150                "rps": None,
151                "wall_time_s": None,
152                "p50_s": None,
153                "p90_s": None,
154                "p99_s": None,
155                "mean_s": None,
156                "short_n": None,
157                "short_p50_s": None,
158                "short_p90_s": None,
159                "short_p99_s": None,
160                "long_n": None,
161                "long_p50_s": None,
162                "long_p90_s": None,
163                "long_p99_s": None,
164                "status": "run_failed",
165                "error": type(e).__name__,
166            }
167            append_ckpt_row(fail_row, CKPT_CSV)
168            print("run_load failed:", type(e).__name__, str(e)[:200])
169
170 summary_df = load_ckpt_df(CKPT_CSV).copy()
171 if "status" in summary_df.columns:
172     summary_df = summary_df[summary_df["status"].isna()].copy()
173
174 summary_df = summary_df.sort_values("concurrency").reset_index(drop=True)
175
176 print("Concurrency sweep done/resumed. Summary DF loaded from Drive:")
177 display(summary_df)
178
179 FINAL_CSV = os.path.join(RUN_DIR, "summary_df_final.csv")
180 summary_df.to_csv(FINAL_CSV, index=False)
181 print("Final snapshot saved:", FINAL_CSV)
182
```

```
Drive already mounted at /content/drive; to attempt to forcibly remount, call drive.mount("/content/drive", force
Loaded checkpoint rows: 0
Completed concurrencies: []
Running: BAD_unfair mix_50_50.jsonl c= 1
BAD_unfair I mix_50_50.jsonl I c=1: 100%                              200/200 [02:32<00:00,  1.22it/s, fail=0, ok=200]
Running: BAD_unfair mix_50_50.jsonl c= 2
BAD_unfair I mix_50_50.jsonl I c=2: 100%                              200/200 [01:25<00:00,  2.07it/s, fail=0, ok=200]
Running: BAD_unfair mix_50_50.jsonl c= 4
BAD_unfair I mix_50_50.jsonl I c=4: 100%                              200/200 [00:55<00:00,  3.03it/s, fail=0, ok=200]
Running: BAD_unfair mix_50_50.jsonl c= 8
BAD_unfair I mix_50_50.jsonl I c=8: 100%                              200/200 [00:40<00:00,  6.62it/s, fail=0, ok=200]
Running: BAD_unfair mix_50_50.jsonl c= 16
BAD_unfair I mix_50_50.jsonl I c=16: 100%                             200/200 [00:33<00:00,  7.47it/s, fail=0, ok=200]
Running: BAD_unfair mix_50_50.jsonl c= 32
BAD_unfair I mix_50_50.jsonl I c=32: 100%                             200/200 [00:33<00:00, 13.21it/s, fail=0, ok=200]
Running: BAD_unfair mix_50_50.jsonl c= 48
BAD_unfair I mix_50_50.jsonl I c=48: 100%                             200/200 [00:32<00:00, 10.68it/s, fail=0, ok=200]
Running: BAD_unfair mix_50_50.jsonl c= 64
BAD_unfair I mix_50_50.jsonl I c=64: 100%                             200/200 [00:32<00:00, 35.51it/s, fail=0, ok=200]
Concurrency sweep done/resumed. Summary DF loaded from Drive:
```

|   | label | workload | concurrency | ok | fail | rps | wall_time_s | p50_s | p90_s | p99_s | mean_s | short |
|---|-------|----------|-------------|-----|------|-----|-------------|-------|-------|-------|--------|-------|
| 0 | BAD_unfair | mix_50_50.jsonl | 1 | 200 | 0 | 1.310728 | 152.586979 | 0.455211 | 1.264692 | 1.302439 | 0.762749 | |
| 1 | BAD_unfair | mix_50_50.jsonl | 2 | 200 | 0 | 2.350512 | 85.087858 | 0.753871 | 1.356037 | 1.788700 | 0.848931 | |
| 2 | BAD_unfair | mix_50_50.jsonl | 4 | 200 | 0 | 3.588961 | 55.726434 | 0.974289 | 1.860058 | 2.692315 | 1.103148 | |
| 3 | BAD_unfair | mix_50_50.jsonl | 8 | 200 | 0 | 4.897939 | 40.833506 | 1.472894 | 2.478667 | 3.694609 | 1.622580 | |
| 4 | BAD_unfair | mix_50_50.jsonl | 16 | 200 | 0 | 5.902047 | 33.886548 | 2.537444 | 4.032396 | 6.204649 | 2.649116 | |
| 5 | BAD_unfair | mix_50_50.jsonl | 32 | 200 | 0 | 5.912537 | 33.826429 | 5.829220 | 7.589303 | 9.536481 | 5.360095 | |
| 6 | BAD_unfair | mix_50_50.jsonl | 48 | 200 | 0 | 6.029853 | 33.168302 | 7.492978 | 11.804981 | 13.755477 | 7.446991 | |
| 7 | BAD_unfair | mix_50_50.jsonl | 64 | 200 | 0 | 5.817353 | 34.379898 | 9.910197 | 16.760385 | 21.667766 | 10.589799 | |

```
Final snapshot saved: /content/drive/MyDrive/vllm_sweeps/exp2_bad_unfair/summary_df_final.csv
```

Next steps: ( Generate code with `summary_df` ) ( New interactive sheet )

## Plot some meaningful charts from Sweep results

```python
1  df = summary_df.copy()
2
3  num_cols = [
4      "concurrency","rps","wall_time_s","p50_s","p90_s","p99_s","mean_s",
5      "short_p50_s","short_p90_s","short_p99_s",
6      "long_p50_s","long_p90_s","long_p99_s",
7      "short_n","long_n"
8  ]
9  for c in num_cols:
10     if c in df.columns:
11         df[c] = pd.to_numeric(df[c], errors="coerce")
12
13 df = df.sort_values("concurrency").reset_index(drop=True)
14
15 # output directory
16 PLOT_DIR = "plots_bad_sweep"
17 os.makedirs(PLOT_DIR, exist_ok=True)
18
19
20 def savefig(name):
21     path = os.path.join(PLOT_DIR, name)
22     plt.tight_layout()
23     plt.savefig(path, dpi=200, bbox_inches="tight")
24     print("Saved:", path)
25
26
27 # Throughput vs concurrency
28 plt.figure()
29 plt.plot(df["concurrency"], df["rps"], marker="o")
30 plt.xlabel("Client concurrency")
31 plt.ylabel("Throughput (requests/sec)")
32 plt.title("BAD_unfair: Throughput vs Concurrency")
33 plt.grid(True, alpha=0.3)
34 savefig("bad_throughput_vs_concurrency.png")
```
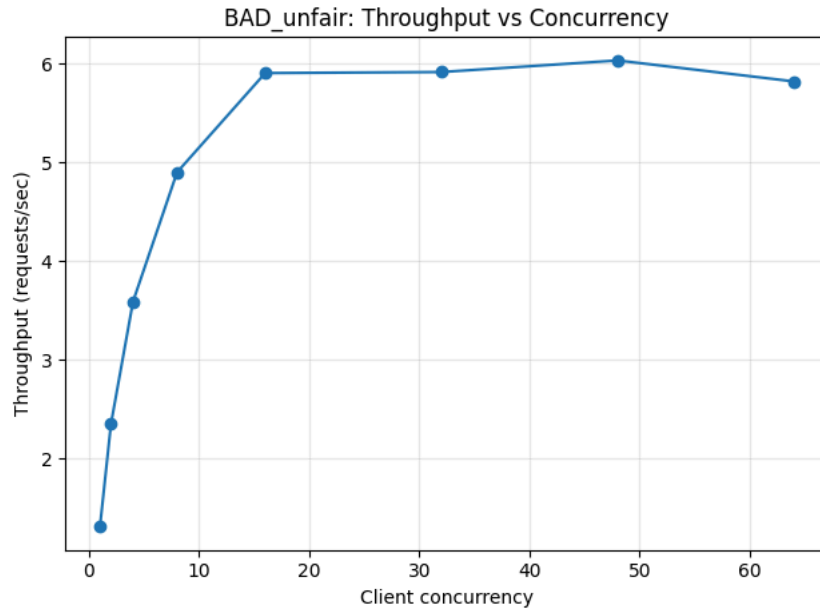
```python
35 plt.show()
36
37 # Overall tail latency vs concurrency (p50 / p90 / p99)
38 for metric, title in [
39     ("p50_s", "BAD_unfair: Overall p50 vs Concurrency"),
40     ("p90_s", "BAD_unfair: Overall p90 vs Concurrency"),
41     ("p99_s", "BAD_unfair: Overall p99 vs Concurrency"),
42 ]:
43     plt.figure()
44     plt.plot(df["concurrency"], df[metric], marker="o")
45     plt.xlabel("Client concurrency")
46     plt.ylabel(f"Latency (s) — {metric.replace('_s','')}")
47     plt.title(title)
48     plt.grid(True, alpha=0.3)
49     savefig(f"bad_{metric}_vs_concurrency.png")
50     plt.show()
51
52 # Short vs Long p99 comparison (same plot)
53 plt.figure()
54 plt.plot(df["concurrency"], df["short_p99_s"], marker="o", label="short p99")
55 plt.plot(df["concurrency"], df["long_p99_s"], marker="o", label="long p99")
56 plt.xlabel("Client concurrency")
57 plt.ylabel("Latency (s)")
58 plt.title("BAD_unfair: Short vs Long p99 vs Concurrency")
59 plt.grid(True, alpha=0.3)
60 plt.legend()
61 savefig("bad_short_long_p99_vs_concurrency.png")
62 plt.show()
63
64 # "Knee" visualization: p99 vs RPS (tradeoff curve)
65 plt.figure()
66 plt.plot(df["rps"], df["p99_s"], marker="o")
67 for _, r in df.iterrows():
68     plt.text(r["rps"], r["p99_s"], str(int(r["concurrency"])), fontsize=9)
69
70 plt.xlabel("Throughput (requests/sec)")
71 plt.ylabel("Overall p99 latency (s)")
72 plt.title("BAD_unfair: p99 vs Throughput (labels = concurrency)")
73 plt.grid(True, alpha=0.3)
74 savefig("bad_p99_vs_rps.png")
75 plt.show()
76
77 print("Done. Plot folder:", PLOT_DIR)
78 df[["concurrency","rps","p50_s","p90_s","p99_s","short_p99_s","long_p99_s"]]
79
```

Saved: plots_bad_sweep/bad_throughput_vs_concurrency.png

## BAD_unfair: Throughput vs Concurrency



Saved: plots_bad_sweep/bad_p50_s_vs_concurrency.png

## BAD_unfair: Overall p50 vs Concurrency



Saved: plots_bad_sweep/bad_p90_s_vs_concurrency.png

## BAD_unfair: Overall p90 vs Concurrency



Saved: plots_bad_sweep/bad_p99_s_vs_concurrency.png

## BAD_unfair: Overall p99 vs Concurrency

Saved: plots_bad_sweep/bad_short_long_p99_vs_concurrency.png



Saved: plots_bad_sweep/bad_p99_vs_rps.png



Done. Plot folder: plots_bad_sweep

| | concurrency | rps | p50_s | p90_s | p99_s | short_p99_s | long_p99_s |
|---|---|---|---|---|---|---|---|
| 0 | 1 | 1.310728 | 0.455211 | 1.264692 | 1.302439 | 0.459093 | 1.303359 |
| 1 | 2 | 2.350512 | 0.753871 | 1.356037 | 1.788700 | 0.975944 | 1.797350 |
| 2 | 4 | 3.588961 | 0.974289 | 1.860058 | 2.692315 | 1.870385 | 2.711042 |
| 3 | 8 | 4.897939 | 1.472894 | 2.478667 | 3.694609 | 2.762490 | 3.715869 |
| 4 | 16 | 5.902047 | 2.537444 | 4.032396 | 6.204649 | 4.027770 | 6.210381 |

| 5 | 32 | 5.912537 | 5.829220 | 7.589303 | 9.536481 | 7.235541 | 10.189636 |
| 6 | 48 | 6.029853 | 7.492978 | 11.804981 | 13.755477 | 11.278168 | 13.872306 |
| 7 | 64 | 5.817353 | 9.910197 | 16.760385 | 21.667766 | 15.790760 | 21.937133 |

## Conclusion

- **Throughput scales then saturates:** RPS rises from **~1.31 (c=1)** to **~5.90 (c=16)** and then stays around **~5.8–6.0 RPS** through **c=64**, showing the server hits a throughput ceiling.
- **Clear knee and tail blow-up:** overall **p99** grows from **~1.30s (c=1)** to **~6.20s (c=16)**, then jumps to **~9.54s (c=32)** and reaches **~21.67s (c=64)**.
- **Short requests get dragged into the tail:** short **p99** inflates from **~0.46s (c=1)** to **~4.03s (c=16)** and **~15.79s (c=64)**, which indicates strong queueing interference from long prompts.
- **At high concurrency, short and long tails converge:** by **c=64**, short **p99 ~15.79s** and long **p99 ~21.94s**, meaning even short prompts suffer large delays under this bad/unfair configuration.

## ⌄ Saving plots to Drive

```
1 LOCAL_PLOT_DIR = "plots_bad_sweep"
2 DRIVE_PLOT_DIR = "/content/drive/MyDrive/vllm_sweeps/exp2_bad_unfair/plots_bad_sweep"
3
4 if os.path.exists(DRIVE_PLOT_DIR):
5     shutil.rmtree(DRIVE_PLOT_DIR)
6 shutil.copytree(LOCAL_PLOT_DIR, DRIVE_PLOT_DIR)
7
8 print("Copied plots to:", DRIVE_PLOT_DIR)
```

```
Copied plots to: /content/drive/MyDrive/vllm_sweeps/exp2_bad_unfair/plots_bad_sweep
```

## ⌄ Experiment 4: Sweep of concurrency and `max-num-seqs`

In this sweep, we stress-test a deliberately "bad" vLLM server configuration under a mixed workload (mix_50_50) to quantify how tail latency behaves as system load increases. By sweeping both client concurrency and the server's admission capacity (max_num_seqs), we expose where throughput stops scaling, where queueing begins to dominate, and how long-prefill requests can degrade the latency of short requests in a shared serving setup. The goal is to establish a clear baseline of failure modes before introducing improved scheduling/configurations in the next steps.

```
 1 BASE = "http://127.0.0.1:8000"
 2 MODEL_ID = "Qwen/Qwen2.5-3B-Instruct"
 3
 4
 5 CHECKPOINT_CSV = "bad_grid_checkpoint.csv"
 6 MANIFEST_JSON  = "bad_grid_manifest.json"
 7
 8
 9 def kill_server():
10     os.system('pkill -f "vllm.entrypoints.openai.api_server" || true')
11     os.system('pkill -f "VLLM::EngineCore" || true')
12     time.sleep(1.5)
13
14 def tail_log(log_path, n=60):
15     if os.path.exists(log_path):
16         print(f"\n--- tail {log_path} (last {n}) ---")
17         os.system(f"tail -n {n} {log_path}")
18     else:
19         print(f"(log missing: {log_path})")
20
21 def wait_for_ready(base=BASE, log_path="vllm_bad.log", timeout_s=420):
22     t0 = time.time()
23     tries = 0
24     while True:
25         tries += 1
26         try:
27             r = requests.get(f"{base}/v1/models", timeout=2)
28             if r.status_code == 200:
29                 data = r.json().get("data", [])
30                 ids = [m.get("id") for m in data]
31                 print(f" Server ready. Models: {ids}")
32                 return True
33             else:
34                 print(f"[{tries}] /v1/models status={r.status_code}")
35         except Exception as e:
36             print(f"[{tries}] not ready yet: {type(e).__name__}")
37
38         alive = os.system('ps -ef | grep -E "vllm.entrypoints.openai.api_server" | grep -v grep > /dev/null
39         if not alive:
```

```
40                print(" vLLM process is NOT running anymore.")
41                tail_log(log_path, n=140)
42                return False
43
44          if time.time() - t0 > timeout_s:
45                print(" Timed out waiting for server.")
46                tail_log(log_path, n=140)
47                return False
48
49          if tries % 5 == 0:
50                tail_log(log_path, n=30)
51
52          time.sleep(2)
53
54 def warmup_request():
55     payload = {
56         "model": MODEL_ID,
57         "messages": [{"role": "user", "content": "hi"}],
58         "max_tokens": 8,
59         "temperature": 0.0,
60     }
61     try:
62         _ = requests.post(f"{BASE}/v1/chat/completions", json=payload, timeout=30)
63     except Exception:
64         pass
65
66 def start_bad_server(max_num_seqs, log_path="vllm_bad.log"):
67     kill_server()
68     try:
69         os.remove(log_path)
70     except FileNotFoundError:
71         pass
72
73     cmd = f"""
74 nohup python3 -m vllm.entrypoints.openai.api_server \
75   --model {MODEL_ID} \
76   --host 127.0.0.1 --port 8000 \
77   --dtype half \
78   --max-model-len 4096 \
79   --gpu-memory-utilization 0.70 \
80   --disable-log-stats \
81   --enforce-eager \
82   --max-num-batched-tokens 8192 \
83   --long-prefill-token-threshold 1024 \
84   --max-num-seqs {int(max_num_seqs)} \
85   > {log_path} 2>&1 </dev/null &
86 """
87     os.system(cmd)
88     time.sleep(1.0)
89
90
91
92 def _safe_quantile(series, q):
93     if series is None or len(series) == 0:
94         return None
95     return float(series.quantile(q))
96
97 def _safe_mean(series):
98     if series is None or len(series) == 0:
99         return None
100    return float(series.mean())
101
102 def subset_metrics(df, short_cutoff_tokens=64, long_cutoff_tokens=2048):
103    ok = df[df["ok"] == True].copy()
104    lat = ok["latency_s"]
105
106    out = {
107        "overall_n": int(len(ok)),
108        "overall_p50": _safe_quantile(lat, 0.50),
109        "overall_p90": _safe_quantile(lat, 0.90),
110        "overall_p99": _safe_quantile(lat, 0.99),
111        "overall_mean": _safe_mean(lat),
112    }
113
114    short = ok[ok["prompt_tokens"] <= short_cutoff_tokens]
115    out.update({
116        "short_n": int(len(short)),
117        "short_p50": _safe_quantile(short["latency_s"], 0.50),
118        "short_p90": _safe_quantile(short["latency_s"], 0.90),
119        "short_p99": _safe_quantile(short["latency_s"], 0.99),
120        "short_mean": _safe_mean(short["latency_s"]),
121    })
```

```
122
123     long = ok[ok["prompt_tokens"] >= long_cutoff_tokens]
124     out.update({
125         "long_n": int(len(long)),
126         "long_p50": _safe_quantile(long["latency_s"], 0.50),
127         "long_p90": _safe_quantile(long["latency_s"], 0.90),
128         "long_p99": _safe_quantile(long["latency_s"], 0.99),
129         "long_mean": _safe_mean(long["latency_s"]),
130     })
131     return out
132
133 def row_from_run(label, workload, concurrency, max_num_seqs, df, summ):
134     m = subset_metrics(df, short_cutoff_tokens=64, long_cutoff_tokens=2048)
135     return {
136         "label": label,
137         "workload": workload,
138         "concurrency": int(concurrency),
139         "max_num_seqs": int(max_num_seqs),
140
141         "status": "ok",
142         "ok": int(summ.get("ok")) if summ.get("ok") is not None else int((df["ok"] == True).sum()),
143         "fail": int(summ.get("fail")) if summ.get("fail") is not None else int((df["ok"] != True).sum()),
144         "rps": float(summ.get("rps")) if summ.get("rps") is not None else None,
145         "wall_time_s": float(summ.get("wall_time_s")) if summ.get("wall_time_s") is not None else None,
146         "summ_p50_s": float(summ.get("p50_s")) if summ.get("p50_s") is not None else None,
147         "summ_p90_s": float(summ.get("p90_s")) if summ.get("p90_s") is not None else None,
148         "summ_p99_s": float(summ.get("p99_s")) if summ.get("p99_s") is not None else None,
149         "summ_mean_s": float(summ.get("mean_s")) if summ.get("mean_s") is not None else None,
150     }
151
152
153 def load_checkpoint_df(path=CHECKPOINT_CSV):
154     if os.path.exists(path):
155         df = pd.read_csv(path)
156         # normalize types
157         if "max_num_seqs" in df.columns:
158             df["max_num_seqs"] = df["max_num_seqs"].astype(int)
159         if "concurrency" in df.columns:
160             df["concurrency"] = df["concurrency"].astype(int)
161         return df
162     return pd.DataFrame()
163
164 def completed_pairs(df_ckpt):
165     # Only treat status=="ok" as completed (you can also include server_failed if you want)
166     if df_ckpt.empty:
167         return set()
168     done = df_ckpt[df_ckpt["status"] == "ok"][["max_num_seqs", "concurrency"]].dropna()
169     return set((int(r.max_num_seqs), int(r.concurrency)) for r in done.itertuples(index=False))
170
171 def append_row_checkpoint(row, path=CHECKPOINT_CSV):
172     df_row = pd.DataFrame([row])
173     header = not os.path.exists(path)
174     df_row.to_csv(path, mode="a", header=header, index=False)
175
176 def write_manifest(payload, path=MANIFEST_JSON):
177     # don't overwrite if already exists unless identical
178     if os.path.exists(path):
179         try:
180             old = json.loads(open(path, "r").read())
181             if old != payload:
182                 print(" Manifest differs from existing one. You may be resuming with different sweep params
183                 print("Existing:", old)
184                 print("New:", payload)
185         except Exception:
186             print(" Could not read existing manifest. Proceeding.")
187         return
188     with open(path, "w") as f:
189         json.dump(payload, f, indent=2)
190
191
192
193 # sweep runner
194
195 WL_FILE = "mix_50_50.jsonl"
196 TOTAL_REQ = 200
197 DESIRED_MAX_TOKENS = 16
198
199 concurrency_list = [1, 2, 4, 8, 16, 32, 48, 64]
200 max_num_seqs_list = [16, 32, 48, 64, 80, 96]
201
202 # save manifest to know what this checkpoint corresponds to
203 write_manifest({
```

```
204        "workload": WL_FILE,
205        "total_requests": TOTAL_REQ,
206        "desired_max_tokens": DESIRED_MAX_TOKENS,
207        "concurrency_list": concurrency_list,
208        "max_num_seqs_list": max_num_seqs_list,
209        "model_id": MODEL_ID,
210        "server_flags": {
211            "max_num_batched_tokens": 8192,
212            "long_prefill_token_threshold": 1024,
213            "gpu_memory_utilization": 0.70,
214            "max_model_len": 4096,
215            "dtype": "half",
216            "enforce_eager": True
217        }
218 })
219
220 ckpt_df = load_checkpoint_df(CHECKPOINT_CSV)
221 done = completed_pairs(ckpt_df)
222 print(f" Loaded checkpoint rows: {len(ckpt_df)} | completed runs: {len(done)}")
223
224 for max_num_seqs in max_num_seqs_list:
225     # If all concurrencies for this max_num_seqs are done, skip server restart entirely
226     all_done_for_this = all((max_num_seqs, c) in done for c in concurrency_list)
227     if all_done_for_this:
228         print(f"\n Skipping max_num_seqs={max_num_seqs} (all concurrencies already done)")
229         continue
230
231     print("\n====================================================")
232     print(f"Starting BAD server with max-num-seqs={max_num_seqs}")
233     print("====================================================")
234     start_bad_server(max_num_seqs=max_num_seqs, log_path="vllm_bad.log")
235
236     ok = wait_for_ready(BASE, log_path="vllm_bad.log", timeout_s=420)
237     if not ok:
238         # record server_failed for missing concurrencies only
239         for c in concurrency_list:
240             if (max_num_seqs, c) in done:
241                 continue
242             row = {
243                 "label": "BAD_unfair",
244                 "workload": WL_FILE,
245                 "concurrency": int(c),
246                 "max_num_seqs": int(max_num_seqs),
247                 "status": "server_failed",
248             }
249             append_row_checkpoint(row, CHECKPOINT_CSV)
250         # refresh checkpoint view
251         ckpt_df = load_checkpoint_df(CHECKPOINT_CSV)
252         done = completed_pairs(ckpt_df)
253         continue
254
255     warmup_request()
256
257     for c in concurrency_list:
258         if (max_num_seqs, c) in done:
259             print(f"  Skip (already done): max_seqs={max_num_seqs}, c={c}")
260             continue
261
262         print(f"\n--- RUN: BAD_unfair | {WL_FILE} | max_seqs={max_num_seqs} | c={c} ---")
263
264         try:
265             df_bad, summ_bad = run_load(
266                 prompt_sets[WL_FILE],
267                 concurrency=c,
268                 total_requests=TOTAL_REQ,
269                 desired_max_tokens=DESIRED_MAX_TOKENS,
270                 desc=f"BAD_unfair | {WL_FILE} | max_seqs={max_num_seqs} | c={c}",
271             )
272
273             row = row_from_run(
274                 label="BAD_unfair",
275                 workload=WL_FILE,
276                 concurrency=c,
277                 max_num_seqs=max_num_seqs,
278                 df=df_bad,
279                 summ=summ_bad,
280             )
281             append_row_checkpoint(row, CHECKPOINT_CSV)
282
283         except Exception as e:
284             # record run_failed so you can retry later
285             row = {
```

```
286                        "label": "BAD_unfair",
287                        "workload": WL_FILE,
288                        "concurrency": int(c),
289                        "max_num_seqs": int(max_num_seqs),
290                        "status": "run_failed",
291                        "error": type(e).__name__,
292                    }
293                    append_row_checkpoint(row, CHECKPOINT_CSV)
294                    print(" run_load failed:", type(e).__name__)
295                    tail_log("vllm_bad.log", n=80)
296
297            # refresh done set so mid-run restarts still work cleanly
298            ckpt_df = load_checkpoint_df(CHECKPOINT_CSV)
299            done = completed_pairs(ckpt_df)
300
301 kill_server()
302
303 summary_df = load_checkpoint_df(CHECKPOINT_CSV)
304 print("\n BAD (concurrency x max-num-seqs) sweep (resume-capable) done.")
305 summary_df
306
```

```
✅ Loaded checkpoint rows: 0 | completed runs: 0

==================================================
Starting BAD server with max-num-seqs=16
==================================================
[1] not ready yet: ConnectionError
[2] not ready yet: ConnectionError
[3] not ready yet: ConnectionError
[4] not ready yet: ConnectionError
[5] not ready yet: ConnectionError

--- tail vllm_bad.log (last 30) ---
[6] not ready yet: ConnectionError
[7] not ready yet: ConnectionError
[8] not ready yet: ConnectionError
[9] not ready yet: ConnectionError
[10] not ready yet: ConnectionError

--- tail vllm_bad.log (last 30) ---
[11] not ready yet: ConnectionError
[12] not ready yet: ConnectionError
[13] not ready yet: ConnectionError
[14] not ready yet: ConnectionError
[15] not ready yet: ConnectionError

--- tail vllm_bad.log (last 30) ---
[16] not ready yet: ConnectionError
[17] not ready yet: ConnectionError
[18] not ready yet: ConnectionError
[19] not ready yet: ConnectionError
[20] not ready yet: ConnectionError

--- tail vllm_bad.log (last 30) ---
[21] not ready yet: ConnectionError
[22] not ready yet: ConnectionError
[23] not ready yet: ConnectionError
[24] not ready yet: ConnectionError
✅ Server ready. Models: ['Qwen/Qwen2.5-3B-Instruct']

--- RUN: BAD_unfair | mix_50_50.jsonl | max_seqs=16 | c=1 ---
```
BAD_unfair | mix_50_50.jsonl | max_seqs=16 | c=1: 100%                                                       200/200 [02:30<00:00,  1.21it/s, fail=0, ok=200]

```
--- RUN: BAD_unfair | mix_50_50.jsonl | max_seqs=16 | c=2 ---
```
BAD_unfair | mix_50_50.jsonl | max_seqs=16 | c=2: 100%                                                       200/200 [01:13<00:00,  2.78it/s, fail=0, ok=200]

```
--- RUN: BAD_unfair | mix_50_50.jsonl | max_seqs=16 | c=4 ---
```
BAD_unfair | mix_50_50.jsonl | max_seqs=16 | c=4: 100%                                                       200/200 [00:42<00:00,  4.44it/s, fail=0, ok=200]

```
--- RUN: BAD_unfair | mix_50_50.jsonl | max_seqs=16 | c=8 ---
```
BAD_unfair | mix_50_50.jsonl | max_seqs=16 | c=8: 100%                                                       200/200 [00:26<00:00,  9.00it/s, fail=0, ok=200]

```
--- RUN: BAD_unfair | mix_50_50.jsonl | max_seqs=16 | c=16 ---
```
BAD_unfair | mix_50_50.jsonl | max_seqs=16 | c=16: 100%                                                      200/200 [00:18<00:00, 10.34it/s, fail=0, ok=200]

```
--- RUN: BAD_unfair | mix_50_50.jsonl | max_seqs=16 | c=32 ---
```
BAD_unfair | mix_50_50.jsonl | max_seqs=16 | c=32: 100%                                                      200/200 [00:15<00:00, 19.97it/s, fail=0, ok=200]

```
--- RUN: BAD_unfair | mix_50_50.jsonl | max_seqs=16 | c=48 ---
```
BAD_unfair | mix_50_50.jsonl | max_seqs=16 | c=48: 100%                                                      200/200 [00:14<00:00, 19.86it/s, fail=0, ok=200]

```
--- RUN: BAD_unfair | mix_50_50.jsonl | max_seqs=16 | c=64 ---
```
BAD_unfair | mix_50_50.jsonl | max_seqs=16 | c=64: 100%                                                      200/200 [00:13<00:00, 19.43it/s, fail=0, ok=200]

```
==================================================
Starting BAD server with max-num-seqs=32
==================================================
[1] not ready yet: ConnectionError
[2] not ready yet: ConnectionError
[3] not ready yet: ConnectionError
[4] not ready yet: ConnectionError
[5] not ready yet: ConnectionError

--- tail vllm_bad.log (last 30) ---
[6] not ready yet: ConnectionError
[7] not ready yet: ConnectionError
[8] not ready yet: ConnectionError
[9] not ready yet: ConnectionError
[10] not ready yet: ConnectionError

--- tail vllm_bad.log (last 30) ---
[11] not ready yet: ConnectionError
[12] not ready yet: ConnectionError
[13] not ready yet: ConnectionError
[14] not ready yet: ConnectionError
```

```
[14] not ready yet: ConnectionError
[15] not ready yet: ConnectionError

--- tail vllm_bad.log (last 30) ---
[16] not ready yet: ConnectionError
[17] not ready yet: ConnectionError
[18] not ready yet: ConnectionError
[19] not ready yet: ConnectionError
[20] not ready yet: ConnectionError

--- tail vllm_bad.log (last 30) ---
[21] not ready yet: ConnectionError
[22] not ready yet: ConnectionError
[23] not ready yet: ConnectionError
[24] not ready yet: ConnectionError
[25] not ready yet: ConnectionError

--- tail vllm_bad.log (last 30) ---
✅ Server ready. Models: ['Qwen/Qwen2.5-3B-Instruct']

--- RUN: BAD_unfair | mix_50_50.jsonl | max_seqs=32 | c=1 ---
BAD_unfair | mix_50_50.jsonl | max_seqs=32 | c=1: 100%          200/200 [02:31<00:00,  1.21it/s, fail=0, ok=200]

--- RUN: BAD_unfair | mix_50_50.jsonl | max_seqs=32 | c=2 ---
BAD_unfair | mix_50_50.jsonl | max_seqs=32 | c=2: 100%          200/200 [01:12<00:00,  2.06it/s, fail=0, ok=200]

--- RUN: BAD_unfair | mix_50_50.jsonl | max_seqs=32 | c=4 ---
BAD_unfair | mix_50_50.jsonl | max_seqs=32 | c=4: 100%          200/200 [00:42<00:00,  4.55it/s, fail=0, ok=200]

--- RUN: BAD_unfair | mix_50_50.jsonl | max_seqs=32 | c=8 ---
BAD_unfair | mix_50_50.jsonl | max_seqs=32 | c=8: 100%          200/200 [00:27<00:00, 10.03it/s, fail=0, ok=200]

--- RUN: BAD_unfair | mix_50_50.jsonl | max_seqs=32 | c=16 ---
BAD_unfair | mix_50_50.jsonl | max_seqs=32 | c=16: 100%         200/200 [00:18<00:00, 10.30it/s, fail=0, ok=200]

--- RUN: BAD_unfair | mix_50_50.jsonl | max_seqs=32 | c=32 ---
BAD_unfair | mix_50_50.jsonl | max_seqs=32 | c=32: 100%         200/200 [00:12<00:00, 30.93it/s, fail=0, ok=200]

--- RUN: BAD_unfair | mix_50_50.jsonl | max_seqs=32 | c=48 ---
BAD_unfair | mix_50_50.jsonl | max_seqs=32 | c=48: 100%         200/200 [00:10<00:00, 22.96it/s, fail=0, ok=200]

--- RUN: BAD_unfair | mix_50_50.jsonl | max_seqs=32 | c=64 ---
BAD_unfair | mix_50_50.jsonl | max_seqs=32 | c=64: 100%         200/200 [00:10<00:00, 30.86it/s, fail=0, ok=200]

==================================================
Starting BAD server with max-num-seqs=48
==================================================
[1] not ready yet: ConnectionError
[2] not ready yet: ConnectionError
[3] not ready yet: ConnectionError
[4] not ready yet: ConnectionError
[5] not ready yet: ConnectionError

--- tail vllm_bad.log (last 30) ---
[6] not ready yet: ConnectionError
[7] not ready yet: ConnectionError
[8] not ready yet: ConnectionError
[9] not ready yet: ConnectionError
[10] not ready yet: ConnectionError

--- tail vllm_bad.log (last 30) ---
[11] not ready yet: ConnectionError
[12] not ready yet: ConnectionError
[13] not ready yet: ConnectionError
[14] not ready yet: ConnectionError
[15] not ready yet: ConnectionError

--- tail vllm_bad.log (last 30) ---
[16] not ready yet: ConnectionError
[17] not ready yet: ConnectionError
[18] not ready yet: ConnectionError
[19] not ready yet: ConnectionError
[20] not ready yet: ConnectionError

--- tail vllm_bad.log (last 30) ---
[21] not ready yet: ConnectionError
[22] not ready yet: ConnectionError
[23] not ready yet: ConnectionError
[24] not ready yet: ConnectionError
✅ Server ready. Models: ['Qwen/Qwen2.5-3B-Instruct']

--- RUN: BAD_unfair | mix_50_50.jsonl | max_seqs=48 | c=1 ---
BAD_unfair | mix_50_50.jsonl | max_seqs=48 | c=1: 100%          200/200 [02:31<00:00,  1.22it/s, fail=0, ok=200]
```

```
--- RUN: BAD_unfair | mix_50_50.jsonl | max_seqs=48 | c=2 ---
BAD_unfair | mix_50_50.jsonl | max_seqs=48 | c=2: 100%          200/200 [01:12<00:00, 2.77it/s, fail=0, ok=200]

--- RUN: BAD_unfair | mix_50_50.jsonl | max_seqs=48 | c=4 ---
BAD_unfair | mix_50_50.jsonl | max_seqs=48 | c=4: 100%          200/200 [00:42<00:00, 5.14it/s, fail=0, ok=200]

--- RUN: BAD_unfair | mix_50_50.jsonl | max_seqs=48 | c=8 ---
BAD_unfair | mix_50_50.jsonl | max_seqs=48 | c=8: 100%          200/200 [00:27<00:00, 8.21it/s, fail=0, ok=200]

--- RUN: BAD_unfair | mix_50_50.jsonl | max_seqs=48 | c=16 ---
BAD_unfair | mix_50_50.jsonl | max_seqs=48 | c=16: 100%         200/200 [00:17<00:00, 9.75it/s, fail=0, ok=200]

--- RUN: BAD_unfair | mix_50_50.jsonl | max_seqs=48 | c=32 ---
BAD_unfair | mix_50_50.jsonl | max_seqs=48 | c=32: 100%         200/200 [00:12<00:00, 21.39it/s, fail=0, ok=200]

--- RUN: BAD_unfair | mix_50_50.jsonl | max_seqs=48 | c=48 ---
BAD_unfair | mix_50_50.jsonl | max_seqs=48 | c=48: 100%         200/200 [00:10<00:00, 17.80it/s, fail=0, ok=200]

--- RUN: BAD_unfair | mix_50_50.jsonl | max_seqs=48 | c=64 ---
BAD_unfair | mix_50_50.jsonl | max_seqs=48 | c=64: 100%         200/200 [00:09<00:00, 38.46it/s, fail=0, ok=200]

==================================================
Starting BAD server with max-num-seqs=64
==================================================
[1] not ready yet: ConnectionError
[2] not ready yet: ConnectionError
[3] not ready yet: ConnectionError
[4] not ready yet: ConnectionError
[5] not ready yet: ConnectionError

--- tail vllm_bad.log (last 30) ---
[6] not ready yet: ConnectionError
[7] not ready yet: ConnectionError
[8] not ready yet: ConnectionError
[9] not ready yet: ConnectionError
[10] not ready yet: ConnectionError

--- tail vllm_bad.log (last 30) ---
[11] not ready yet: ConnectionError
[12] not ready yet: ConnectionError
[13] not ready yet: ConnectionError
[14] not ready yet: ConnectionError
[15] not ready yet: ConnectionError

--- tail vllm_bad.log (last 30) ---
[16] not ready yet: ConnectionError
[17] not ready yet: ConnectionError
[18] not ready yet: ConnectionError
[19] not ready yet: ConnectionError
[20] not ready yet: ConnectionError

--- tail vllm_bad.log (last 30) ---
[21] not ready yet: ConnectionError
[22] not ready yet: ConnectionError
[23] not ready yet: ConnectionError
[24] not ready yet: ConnectionError
✅ Server ready. Models: ['Qwen/Qwen2.5-3B-Instruct']

--- RUN: BAD_unfair | mix_50_50.jsonl | max_seqs=64 | c=1 ---
BAD_unfair | mix_50_50.jsonl | max_seqs=64 | c=1: 100%          200/200 [02:31<00:00, 1.21it/s, fail=0, ok=200]

--- RUN: BAD_unfair | mix_50_50.jsonl | max_seqs=64 | c=2 ---
BAD_unfair | mix_50_50.jsonl | max_seqs=64 | c=2: 100%          200/200 [01:13<00:00, 2.77it/s, fail=0, ok=200]

--- RUN: BAD_unfair | mix_50_50.jsonl | max_seqs=64 | c=4 ---
BAD_unfair | mix_50_50.jsonl | max_seqs=64 | c=4: 100%          200/200 [00:42<00:00, 4.33it/s, fail=0, ok=200]

--- RUN: BAD_unfair | mix_50_50.jsonl | max_seqs=64 | c=8 ---
BAD_unfair | mix_50_50.jsonl | max_seqs=64 | c=8: 100%          200/200 [00:26<00:00, 8.99it/s, fail=0, ok=200]

--- RUN: BAD_unfair | mix_50_50.jsonl | max_seqs=64 | c=16 ---
BAD_unfair | mix_50_50.jsonl | max_seqs=64 | c=16: 100%         200/200 [00:18<00:00, 12.46it/s, fail=0, ok=200]

--- RUN: BAD_unfair | mix_50_50.jsonl | max_seqs=64 | c=32 ---
BAD_unfair | mix_50_50.jsonl | max_seqs=64 | c=32: 100%         200/200 [00:12<00:00, 27.29it/s, fail=0, ok=200]

--- RUN: BAD_unfair | mix_50_50.jsonl | max_seqs=64 | c=48 ---
BAD_unfair | mix_50_50.jsonl | max_seqs=64 | c=48: 100%         200/200 [00:09<00:00, 41.25it/s, fail=0, ok=200]
```

```
--- RUN: BAD_unfair | mix_50_50.jsonl | max_seqs=64 | c=64 ---

BAD_unfair | mix_50_50.jsonl | max_seqs=64 | c=64: 100%                    200/200 [00:08<00:00, 35.42it/s, fail=0, ok=200]


=================================================
Starting BAD server with max-num-seqs=80
=================================================
[1] not ready yet: ConnectionError
[2] not ready yet: ConnectionError
[3] not ready yet: ConnectionError
[4] not ready yet: ConnectionError
[5] not ready yet: ConnectionError

--- tail vllm_bad.log (last 30) ---
[6] not ready yet: ConnectionError
[7] not ready yet: ConnectionError
[8] not ready yet: ConnectionError
[9] not ready yet: ConnectionError
[10] not ready yet: ConnectionError

--- tail vllm_bad.log (last 30) ---
[11] not ready yet: ConnectionError
[12] not ready yet: ConnectionError
[13] not ready yet: ConnectionError
[14] not ready yet: ConnectionError
[15] not ready yet: ConnectionError

--- tail vllm_bad.log (last 30) ---
[16] not ready yet: ConnectionError
[17] not ready yet: ConnectionError
[18] not ready yet: ConnectionError
[19] not ready yet: ConnectionError
[20] not ready yet: ConnectionError

--- tail vllm_bad.log (last 30) ---
[21] not ready yet: ConnectionError
[22] not ready yet: ConnectionError
[23] not ready yet: ConnectionError
[24] not ready yet: ConnectionError
[25] not ready yet: ConnectionError

--- tail vllm_bad.log (last 30) ---
✅ Server ready. Models: ['Qwen/Qwen2.5-3B-Instruct']

--- RUN: BAD_unfair | mix_50_50.jsonl | max_seqs=80 | c=1 ---

BAD_unfair | mix_50_50.jsonl | max_seqs=80 | c=1: 100%                     200/200 [02:31<00:00,  1.21it/s, fail=0, ok=200]


--- RUN: BAD_unfair | mix_50_50.jsonl | max_seqs=80 | c=2 ---

BAD_unfair | mix_50_50.jsonl | max_seqs=80 | c=2: 100%                     200/200 [01:13<00:00,  2.78it/s, fail=0, ok=200]


--- RUN: BAD_unfair | mix_50_50.jsonl | max_seqs=80 | c=4 ---

BAD_unfair | mix_50_50.jsonl | max_seqs=80 | c=4: 100%                     200/200 [00:42<00:00,  4.63it/s, fail=0, ok=200]


--- RUN: BAD_unfair | mix_50_50.jsonl | max_seqs=80 | c=8 ---

BAD_unfair | mix_50_50.jsonl | max_seqs=80 | c=8: 100%                     200/200 [00:27<00:00,  8.01it/s, fail=0, ok=200]


--- RUN: BAD_unfair | mix_50_50.jsonl | max_seqs=80 | c=16 ---

BAD_unfair | mix_50_50.jsonl | max_seqs=80 | c=16: 100%                    200/200 [00:18<00:00,  9.05it/s, fail=0, ok=200]


--- RUN: BAD_unfair | mix_50_50.jsonl | max_seqs=80 | c=32 ---

BAD_unfair | mix_50_50.jsonl | max_seqs=80 | c=32: 100%                    200/200 [00:13<00:00, 17.77it/s, fail=0, ok=200]


--- RUN: BAD_unfair | mix_50_50.jsonl | max_seqs=80 | c=48 ---

BAD_unfair | mix_50_50.jsonl | max_seqs=80 | c=48: 100%                    200/200 [00:10<00:00, 20.46it/s, fail=0, ok=200]


--- RUN: BAD_unfair | mix_50_50.jsonl | max_seqs=80 | c=64 ---

BAD_unfair | mix_50_50.jsonl | max_seqs=80 | c=64: 100%                    200/200 [00:08<00:00, 22.60it/s, fail=0, ok=200]


=================================================
Starting BAD server with max-num-seqs=96
=================================================
[1] not ready yet: ConnectionError
[2] not ready yet: ConnectionError
[3] not ready yet: ConnectionError
[4] not ready yet: ConnectionError
[5] not ready yet: ConnectionError

--- tail vllm_bad.log (last 30) ---
[6] not ready yet: ConnectionError
[7] not ready yet: ConnectionError
[8] not ready yet: ConnectionError
[9] not ready yet: ConnectionError
[10] not ready yet: ConnectionError
```

```
--- tail vllm_bad.log (last 30) ---
[11] not ready yet: ConnectionError
[12] not ready yet: ConnectionError
[13] not ready yet: ConnectionError
[14] not ready yet: ConnectionError
[15] not ready yet: ConnectionError

--- tail vllm_bad.log (last 30) ---
[16] not ready yet: ConnectionError
[17] not ready yet: ConnectionError
[18] not ready yet: ConnectionError
[19] not ready yet: ConnectionError
[20] not ready yet: ConnectionError

--- tail vllm_bad.log (last 30) ---
[21] not ready yet: ConnectionError
[22] not ready yet: ConnectionError
[23] not ready yet: ConnectionError
[24] not ready yet: ConnectionError
```
✅ Server ready. Models: ['Qwen/Qwen2.5-3B-Instruct']

--- RUN: BAD_unfair | mix_50_50.jsonl | max_seqs=96 | c=1 ---

BAD_unfair l mix_50_50.jsonl l max_seqs=96 l c=1: 100%          200/200 [02:30<00:00, 1.22it/s, fail=0, ok=200]

--- RUN: BAD_unfair | mix_50_50.jsonl | max_seqs=96 | c=2 ---

BAD_unfair l mix_50_50.jsonl l max_seqs=96 l c=2: 100%          200/200 [01:23<00:00, 2.05it/s, fail=0, ok=200]

--- RUN: BAD_unfair | mix_50_50.jsonl | max_seqs=96 | c=4 ---

BAD_unfair l mix_50_50.jsonl l max_seqs=96 l c=4: 100%          200/200 [00:52<00:00, 2.91it/s, fail=0, ok=200]

--- RUN: BAD_unfair | mix_50_50.jsonl | max_seqs=96 | c=8 ---

BAD_unfair l mix_50_50.jsonl l max_seqs=96 l c=8: 100%          200/200 [00:36<00:00, 6.51it/s, fail=0, ok=200]

--- RUN: BAD_unfair | mix_50_50.jsonl | max_seqs=96 | c=16 ---

BAD_unfair l mix_50_50.jsonl l max_seqs=96 l c=16: 100%          200/200 [00:23<00:00, 7.26it/s, fail=0, ok=200]

--- RUN: BAD_unfair | mix_50_50.jsonl | max_seqs=96 | c=32 ---

BAD_unfair l mix_50_50.jsonl l max_seqs=96 l c=32: 100%          200/200 [00:15<00:00, 12.58it/s, fail=0, ok=200]

--- RUN: BAD_unfair | mix_50_50.jsonl | max_seqs=96 | c=48 ---

BAD_unfair l mix_50_50.jsonl l max_seqs=96 l c=48: 100%          200/200 [00:14<00:00, 32.98it/s, fail=0, ok=200]

--- RUN: BAD_unfair | mix_50_50.jsonl | max_seqs=96 | c=64 ---

BAD_unfair l mix_50_50.jsonl l max_seqs=96 l c=64: 100%          200/200 [00:11<00:00, 21.06it/s, fail=0, ok=200]

✅ BAD (concurrency x max-num-seqs) sweep (resume-capable) done.

| | label | workload | concurrency | max_num_seqs | status | ok | fail | rps | wall_time_s | summ_p50_s | ... | sh |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | BAD_unfair | mix_50_50.jsonl | 1 | 16 | ok | 200 | 0 | 1.327445 | 150.665335 | 0.456500 | ... | |
| 1 | BAD_unfair | mix_50_50.jsonl | 2 | 16 | ok | 200 | 0 | 2.735775 | 73.105440 | 0.729429 | ... | |
| 2 | BAD_unfair | mix_50_50.jsonl | 4 | 16 | ok | 200 | 0 | 4.670095 | 42.825685 | 0.773686 | ... | |
| 3 | BAD_unfair | mix_50_50.jsonl | 8 | 16 | ok | 200 | 0 | 7.587554 | 26.358956 | 1.061564 | ... | |
| 4 | BAD_unfair | mix_50_50.jsonl | 16 | 16 | ok | 200 | 0 | 11.004723 | 18.174014 | 1.386675 | ... | |
| 5 | BAD_unfair | mix_50_50.jsonl | 32 | 16 | ok | 200 | 0 | 13.177642 | 15.177222 | 2.401322 | ... | |
| 6 | BAD_unfair | mix_50_50.jsonl | 48 | 16 | ok | 200 | 0 | 13.587360 | 14.719563 | 3.413464 | ... | |
| 7 | BAD_unfair | mix_50_50.jsonl | 64 | 16 | ok | 200 | 0 | 13.562172 | 14.746900 | 4.487614 | ... | |
| 8 | BAD_unfair | mix_50_50.jsonl | 1 | 32 | ok | 200 | 0 | 1.319701 | 151.549532 | 0.458507 | ... | |
| 9 | BAD_unfair | mix_50_50.jsonl | 2 | 32 | ok | 200 | 0 | 2.747514 | 72.793065 | 0.732680 | ... | |
| 10 | BAD_unfair | mix_50_50.jsonl | 4 | 32 | ok | 200 | 0 | 4.677001 | 42.762445 | 0.775419 | ... | |
| 11 | BAD_unfair | mix_50_50.jsonl | 8 | 32 | ok | 200 | 0 | 7.324900 | 27.304127 | 1.088049 | ... | |
| 12 | BAD_unfair | mix_50_50.jsonl | 16 | 32 | ok | 200 | 0 | 10.693156 | 18.703553 | 1.404330 | ... | |
| 13 | BAD_unfair | mix_50_50.jsonl | 32 | 32 | ok | 200 | 0 | 15.372674 | 13.010098 | 2.089596 | ... | |
| 14 | BAD_unfair | mix_50_50.jsonl | 48 | 32 | ok | 200 | 0 | 17.569705 | 11.383230 | 2.574769 | ... | |
| 15 | BAD_unfair | mix_50_50.jsonl | 64 | 32 | ok | 200 | 0 | 15.976706 | 12.518225 | 3.737646 | ... | |
| 16 | BAD_unfair | mix_50_50.jsonl | 1 | 48 | ok | 200 | 0 | 1.322014 | 151.284371 | 0.458063 | ... | |
| 17 | BAD_unfair | mix_50_50.jsonl | 2 | 48 | ok | 200 | 0 | 2.743948 | 72.887680 | 0.731331 | ... | |
| 18 | BAD_unfair | mix_50_50.jsonl | 4 | 48 | ok | 200 | 0 | 4.716243 | 42.406635 | 0.777884 | ... | |

## Conclusion

| 19 | BAD_unfair | mix_50_50.jsonl | 8 | 48 | ok | 200 | 0 | 7.384087 | 27.085270 | 1.056648 | ... |
| 20 | BAD_unfair | mix_50_50.jsonl | 16 | 48 | ok | 200 | 0 | 11.097723 | 18.021716 | 1.366025 | ... |
| 21 | BAD_unfair | mix_50_50.jsonl | 32 | 48 | ok | 200 | 0 | 15.919387 | 12.563298 | 2.105477 | ... |
| 22 | BAD_unfair | mix_50_50.jsonl | 48 | 48 | ok | 200 | 0 | 17.584097 | 11.373914 | 2.540315 | ... |
| 23 | BAD_unfair | mix_50_50.jsonl | 64 | 48 | ok | 200 | 0 | 18.785692 | 10.646400 | 2.921785 | ... |
| 24 | BAD_unfair | mix_50_50.jsonl | 1 | 64 | ok | 200 | 0 | 1.319597 | 151.561409 | 0.457572 | ... |
| 25 | BAD_unfair | mix_50_50.jsonl | 2 | 64 | ok | 200 | 0 | 2.721553 | 73.487456 | 0.733294 | ... |
| 26 | BAD_unfair | mix_50_50.jsonl | 4 | 64 | ok | 200 | 0 | 4.658324 | 42.933894 | 0.793960 | ... |
| 27 | BAD_unfair | mix_50_50.jsonl | 8 | 64 | ok | 200 | 0 | 7.484882 | 26.900224 | 1.056196 | ... |
| 29 | BAD_unfair | mix_50_50.jsonl | 32 | 64 | ok | 200 | 0 | 15.559964 | 12.853500 | 1.932129 | ... |

- **Throughput scaling:** RPS increases with client concurrency, but the gains flatten after roughly $c=32$ for most $max\_num\_seqs$ values (diminishing returns region).

- **Tail latency inflation:** As concurrency rises, **p99 grows much faster than p50**, showing queueing and head-of-line effects dominating at high load.

- **Short vs long separation:** In the mixed workload, **short requests stay much faster than long requests**, but short p99 still inflates as concurrency increases because shorts get stuck behind long prefill work.

- **Impact of max_num_seqs:** Larger $max\_num\_seqs$ generally allows higher throughput at mid/high concurrency, but it can also **worsen tail latency** if the server admits too many concurrent sequences and amplifies contention/queueing under this "bad" configuration.

```python
1  OUT_DIR = "/content/drive/MyDrive/vllm_sweeps/bad"
2  os.makedirs(OUT_DIR, exist_ok=True)
3
4  # Copy any existing on-disk artifacts to Drive (if they exist)
5  for f in ["bad_grid_checkpoint.csv", "bad_grid_manifest.json", "vllm_bad.log"]:
6      if os.path.exists(f):
7          shutil.copy(f, f"{OUT_DIR}/{f}")
8          print("✅ Copied:", f"{OUT_DIR}/{f}")
9      else:
10         print("⚠️ Missing locally:", f)
11
12 # save the final dataframe we have in memory (summary_df)
13 try:
14     summary_df_path = f"{OUT_DIR}/bad_summary_df.csv"
15     summary_df.to_csv(summary_df_path, index=False)
16     print("✅ Saved summary_df:", summary_df_path)
17 except NameError:
18     print("⚠️ summary_df not found in memory. If your df has a different name, save that instead.")
19
```

```
Drive already mounted at /content/drive; to attempt to forcibly remount, call drive.mount("/content/drive", force
✅ Copied: /content/drive/MyDrive/vllm_sweeps/bad/bad_grid_checkpoint.csv
✅ Copied: /content/drive/MyDrive/vllm_sweeps/bad/bad_grid_manifest.json
✅ Copied: /content/drive/MyDrive/vllm_sweeps/bad/vllm_bad.log
✅ Saved summary_df: /content/drive/MyDrive/vllm_sweeps/bad/bad_summary_df.csv
```

| 42 | BAD_unfair | mix_50_50.jsonl | 4 | 96 | ok | 200 | 0 | 5.773505 | 32.984304 | 0.903605 | ... |
| 43 | BAD_unfair | mix_50_50.jsonl | 8 | 96 | ok | 200 | 0 | 5.551586 | 36.025742 | 1.418485 | ... |
| 44 | BAD_unfair | mix_50_50.jsonl | 16 | 96 | ok | 200 | 0 | 8.342541 | 23.973511 | 1.792479 | ... |
| 45 | BAD_unfair | mix_50_50.jsonl | 32 | 96 | ok | 200 | 0 | 12.902239 | 15.501185 | 2.657199 | ... |
| 46 | BAD_unfair | mix_50_50.jsonl | 48 | 96 | ok | 200 | 0 | 13.514812 | 14.798578 | 3.581036 | ... |
| 47 | BAD_unfair | mix_50_50.jsonl | 64 | 96 | ok | 200 | 0 | 14.880155 | 13.440721 | 3.983440 | ... |

48 rows × 28 columns

## Plots for Bidirectional Sweep in Bad Config

```python
1  df = summary_df.copy()
2
3  # Keep only good rows
4  df = df[df["status"] == "ok"].copy()
5
6  # Ensure numeric
7  for col in ["concurrency", "max_num_seqs", "rps", "short_p99", "long_p99", "overall_p99", "wall_time_s"]:
8      if col in df.columns:
9          df[col] = pd.to_numeric(df[col], errors="coerce")
10
11 # Sort for consistent plotting
12 df = df.sort_values(["max_num_seqs", "concurrency"]).reset_index(drop=True)
13
14 concurrency_list = sorted(df["concurrency"].unique().tolist())
15 max_num_seqs_list = sorted(df["max_num_seqs"].unique().tolist())
16
17 print("Concurrencies:", concurrency_list)
18 print("max_num_seqs:", max_num_seqs_list)
19 print("Rows:", len(df))
20
```

```
Concurrencies: [1, 2, 4, 8, 16, 32, 48, 64]
max_num_seqs: [16, 32, 48, 64, 80, 96]
Rows: 48
```

```python
1  def plot_heatmap(pivot_df, title, xlabel="concurrency", ylabel="max_num_seqs", fmt="{:.2f}"):
2      """
3      pivot_df: index=max_num_seqs, columns=concurrency, values=metric
4      """
5      data = pivot_df.values.astype(float)
6
7      plt.figure(figsize=(10, 5.5))
8      im = plt.imshow(data, aspect="auto", interpolation="nearest")
```

```
 9    plt.colorbar(im)
10
11    plt.title(title)
12    plt.xlabel(xlabel)
13    plt.ylabel(ylabel)
14
15    plt.xticks(ticks=np.arange(pivot_df.shape[1]), labels=pivot_df.columns.tolist())
16    plt.yticks(ticks=np.arange(pivot_df.shape[0]), labels=pivot_df.index.tolist())
17
18
19    for i in range(pivot_df.shape[0]):
20        for j in range(pivot_df.shape[1]):
21            val = data[i, j]
22            if np.isfinite(val):
23                plt.text(j, i, fmt.format(val), ha="center", va="center", fontsize=8)
24
25    plt.tight_layout()
26    plt.show()
27
28 # Heatmap: rps
29 p_rps = df.pivot_table(index="max_num_seqs", columns="concurrency", values="rps", aggfunc="mean")
30 plot_heatmap(p_rps, "BAD grid: Throughput (rps)")
31
32 # Heatmap: short_p99
33 p_sp99 = df.pivot_table(index="max_num_seqs", columns="concurrency", values="short_p99", aggfunc="mean")
34 plot_heatmap(p_sp99, "BAD grid: short p99 (s)")
35
36 # Heatmap: long_p99
37 p_lp99 = df.pivot_table(index="max_num_seqs", columns="concurrency", values="long_p99", aggfunc="mean")
38 plot_heatmap(p_lp99, "BAD grid: long p99 (s)")
39
```

## BAD grid: Throughput (rps)

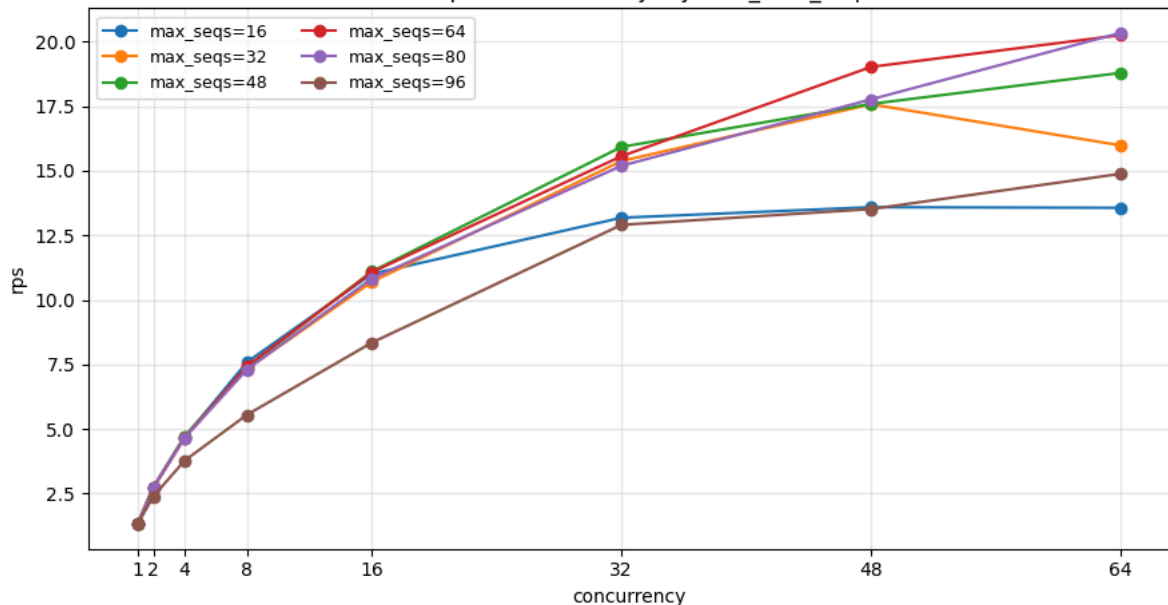| max_num_seqs | | | | | | | |
|---|---|---|---|---|---|---|---|
| 16 | 1.33 | 2.74 | 4.67 | 7.59 | 11.00 | 13.18 | 13.59 | 13.56 |
| 32 | 1.32 | 2.71 | 4.66 | 7.52 | 11.10 | 15.37 | 17.57 | 15.98 |
| 48 | 1.32 | 2.72 | 4.72 | 7.29 | 11.10 | 15.92 | 17.58 | 18.79 |
| 64 | 1.32 | 2.72 | 4.66 | 7.43 | 11.07 | 15.56 | 19.02 | 20.25 |
| 80 | 1.33 | 2.72 | 4.64 | 7.50 | 10.82 | 15.19 | 17.76 | 20.34 |
| 96 | 1.33 | 2.73 | 5.57 | 8.34 | 12.90 | 13.51 | 14.88 |

## Heatmap Interpretaion

- **Throughput scales with concurrency until ~32–64**, but the gain after ~32 is smaller than the latency cost. You can see the plateau: rps improves from low concurrency to 32, then flattens for many max_num_seqs settings.

- **max_num_seqs acts like a capacity knob for concurrency.** At high concurrency (48–64), small max_num_seqs (like 16) cannot keep up, so both short and long p99 inflate. Mid-range values (around 64–80) sustain higher rps while keeping p99 lower.

- **Too large max_num_seqs can hurt.** The row at max_num_seqs = 96 shows lower throughput and much higher short and long p99 at high concurrency, consistent with extra batching/queueing overhead and heavier contention when you allow too many in-flight sequences.

- **Short requests still suffer under mixed workloads when scheduling is unfair.** Even though "short p99" is lower than long p99, it rises sharply with concurrency (especially at c=48–64), showing head-of-line blocking where long prefills delay short requests.

- **Practical "knee" region:** max_num_seqs around **64–80** with concurrency around **16–32** gives a better tradeoff. It keeps rps high without the steep tail-latency blow-up seen at the extreme high concurrency settings.

```python
1  plt.figure(figsize=(9, 5))
2  for m in max_num_seqs_list:
3      sub = df[df["max_num_seqs"] == m].sort_values("concurrency")
4      plt.plot(sub["concurrency"], sub["rps"], marker="o", label=f"max_seqs={m}")
5
6  plt.title("BAD: rps vs concurrency (by max_num_seqs)")
7  plt.xlabel("concurrency")
8  plt.ylabel("rps")
9  plt.xticks(concurrency_list)
10 plt.grid(True, alpha=0.3)
11 plt.legend(ncol=2, fontsize=9)
12 plt.tight_layout()
13 plt.show()
14
```



BAD: rps vs concurrency (by max_num_seqs)

| max_num_seqs | | | | | | | |
|---|---|---|---|---|---|---|---|
| 48 | 1.30 | 1.32 | 1.55 | 1.89 | 2.49 | 3.51 | 4.91 | 6.06 |
| 64 | 1.31 | 1.32 | 1.62 | 2.06 | 2.06 | 3.20 | 3.77 | 4.68 |
| 80 | 1.31 | 1.31 | 1.56 | 1.95 | 2.15 | 3.49 | 4.43 | 4.86 |
| 96 | | | | | | | 5.97 | 6.44 |

## Interpretation: BAD configuration throughput (RPS) vs concurrency, grouped by max_num_seqs
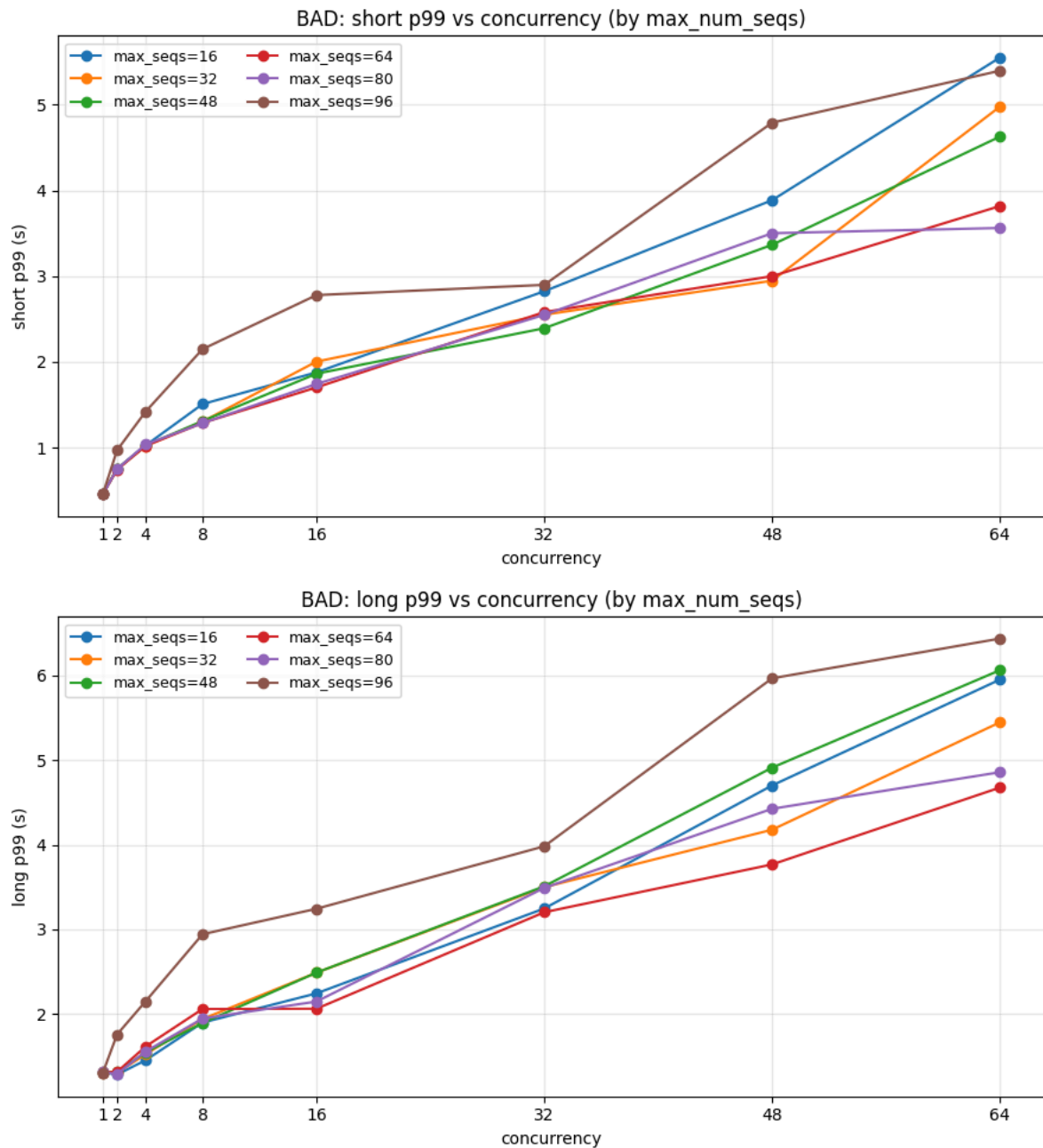
- **RPS increases with concurrency for all settings up to ~16–32**, which is the "scaling region" where the GPU is getting better utilized as we add more in-flight work.

- **max_num_seqs controls how far you can scale before flattening.**
  - With **max_num_seqs=16**, the curve **plateaus early** (around concurrency 32) and stays near ~13–14 rps, meaning the server cannot admit enough parallel sequences to benefit from more client concurrency.

- With **max_num_seqs=48/64/80**, the curves keep improving to higher concurrencies and reach the best throughput at the right edge (48–64).

- **The best throughput in this sweep occurs at high concurrency (48–64) with max_num_seqs around 64–80** (top curves near ~20 rps). That suggests this range is closer to the capacity point of this setup.

- **max_num_seqs=96 underperforms at low and mid concurrency and only partially catches up later.** This is consistent with "too many allowed in-flight sequences" adding overhead (more queueing, scheduling contention, memory pressure), so the system wastes efficiency instead of converting concurrency into throughput.

- **Diminishing returns are visible beyond ~48 concurrency**, even for the best max_num_seqs values. The slope flattens, meaning you are approaching saturation, and extra concurrency mostly increases queueing rather than throughput.

```python
# short p99
plt.figure(figsize=(9, 5))
for m in max_num_seqs_list:
    sub = df[df["max_num_seqs"] == m].sort_values("concurrency")
    plt.plot(sub["concurrency"], sub["short_p99"], marker="o", label=f"max_seqs={m}")

plt.title("BAD: short p99 vs concurrency (by max_num_seqs)")
plt.xlabel("concurrency")
plt.ylabel("short p99 (s)")
plt.xticks(concurrency_list)
plt.grid(True, alpha=0.3)
plt.legend(ncol=2, fontsize=9)
plt.tight_layout()
plt.show()

# long p99
plt.figure(figsize=(9, 5))
for m in max_num_seqs_list:
    sub = df[df["max_num_seqs"] == m].sort_values("concurrency")
    plt.plot(sub["concurrency"], sub["long_p99"], marker="o", label=f"max_seqs={m}")

plt.title("BAD: long p99 vs concurrency (by max_num_seqs)")
plt.xlabel("concurrency")
plt.ylabel("long p99 (s)")
plt.xticks(concurrency_list)
plt.grid(True, alpha=0.3)
plt.legend(ncol=2, fontsize=9)
plt.tight_layout()
plt.show()
```

BAD: short p99 vs concurrency (by max_num_seqs)

BAD: long p99 vs concurrency (by max_num_seqs)

## Interpretation: BAD short p99 vs concurrency and BAD long p99 vs concurrency

- For both short and long requests, **p99 increases monotonically with concurrency**. This is classic queueing behavior: higher offered load increases worst-case wait time.
- **Short p99 inflates sharply at high concurrency** (c=48, c=64), showing that short requests lose their "fast path" when mixed with long prompts under unfair scheduling.
- **Long p99 also climbs steadily**, and it typically climbs faster than short p99 because long prompts contribute more prefill work and hold resources longer.
- **max_num_seqs=64 (and sometimes 80) provides a better tradeoff** than 16/32 (too restrictive) and 96 (too aggressive): higher throughput without blowing up p99 as badly as the overcommitted setting.

## Experiment 5: Creating Fair Scheduling to improve the latency and RPS (or the trade-off)

## Fair sanity check: lower max-num-batched-tokens (btok=768) at c=32

Before running the full FAIR sweeps, I ran a single controlled trial by reducing `max-num-batched-tokens` from the BAD baseline (8192) down to **768**, while keeping everything else identical and evaluating at the **same client concurrency (c=32)** used for the BAD runs.

- **Purpose:** verify that `max-num-batched-tokens` is a meaningful lever under mixed 50/50 traffic, and that lowering it can reduce the "long-prefill domination" effect that inflates tail latency.
- **Why 768:** it's a mid-range value expected to encourage more frequent scheduling turns (smaller batches), which can improve responsiveness for short requests in mixed workloads.
- **Why c=32:** match the BAD configuration's stress level so the comparison is apples-to-apples and any improvement is not just because the system is under-loaded.

```
1 !pkill -f "vllm.entrypoints.openai.api_server" || true
2 !pkill -f "VLLM::EngineCore" || true
3 !rm -f vllm_fair.log
4 !sleep 1
5
```

```
1
2 !nohup python3 -m vllm.entrypoints.openai.api_server \
3   --model Qwen/Qwen2.5-3B-Instruct \
4   --host 127.0.0.1 --port 8000 \
5   --dtype half \
6   --max-model-len 4096 \
7   --gpu-memory-utilization 0.70 \
8   --disable-log-stats \
9   --enforce-eager \
10  --max-num-batched-tokens 768 \
11  --long-prefill-token-threshold 1024 \
12  > vllm_fair.log 2>&1 </dev/null &
```
```
^C
^C
```

```
1 !tail -n 40 vllm_fair.log
2
```

Show hidden output

```
1 df_fair_50_50, summ_fair_50_50 = run_load(
2     prompt_sets["mix_50_50.jsonl"],
3     concurrency=32,
4     total_requests=200,
5     desired_max_tokens=16,
6     desc="FAIR_protected | mix_50_50 | c=16"
7 )
8
9 print("FAIR summary:", summ_fair_50_50)
10 df_fair_50_50.head()
11
```

FAIR_protected | mix_50_50 | c=16: 100%                    200/200 [00:41<00:00, 11.12it/s, fail=0, ok=200]

FAIR summary: {'concurrency': 32, 'total_requests': 200, 'ok': 200, 'fail': 0, 'wall_time_s': 41.268839522000235,

| | latency_s | prompt_tokens | completion_tokens | total_tokens | status_code | prompt_tokens_est | max_tokens_used | trim |
|---|---|---|---|---|---|---|---|---|
| 0 | 3.036116 | 37 | 16 | 53 | 200 | 37 | 16 | |
| 1 | 3.129861 | 37 | 16 | 53 | 200 | 37 | 16 | |
| 2 | 3.074753 | 45 | 16 | 61 | 200 | 45 | 16 | |
| 3 | 3.092510 | 37 | 16 | 53 | 200 | 37 | 16 | |
| 4 | 3.141762 | 38 | 16 | 54 | 200 | 38 | 16 | |

```
1
2 def latency_stats(df, label, mask, subset_name):
3     x = df.loc[mask & (df["ok"] == True), "latency_s"].dropna().values
4     if len(x) == 0:
5         return {"label": label, "subset": subset_name, "n": 0, "p50": None, "p90": None, "p99": None, "mean"
6     return {
7         "label": label,
8         "subset": subset_name,
9         "n": int(len(x)),
10        "p50": float(np.quantile(x, 0.50)),
11        "p90": float(np.quantile(x, 0.90)),
12        "p99": float(np.quantile(x, 0.99)),
13        "mean": float(np.mean(x)),
14    }
15
16 def compare_short_long(df, label, short_cut=64, long_cut=2048):
17     return [
```

```
18        latency_stats(df, label, mask=df["prompt_tokens"] >= 0, subset_name="overall"),
19        latency_stats(df, label, mask=df["prompt_tokens"] <= short_cut, subset_name=f"short<= {short_cut}"),
20        latency_stats(df, label, mask=df["prompt_tokens"] >= long_cut, subset_name=f"long>= {long_cut}"),
21    ]
22
23 bad_stats  = compare_short_long(df_bad,   "BAD_unfair", short_cut=64, long_cut=2048)
24 fair_stats = compare_short_long(df_fair_50_50, "FAIR_768", short_cut=64, long_cut=2048)
25
26 bad_stats, fair_stats
27
```

```
([{'label': 'BAD_unfair',
   'subset': 'overall',
   'n': 200,
   'p50': 5.534205468999517,
   'p90': 10.461404818999926,
   'p99': 11.045583048669595,
   'mean': 6.076663721505029},
  {'label': 'BAD_unfair',
   'subset': 'short<= 64',
   'n': 108,
   'p50': 4.6954104774999905,
   'p90': 8.427761869799998,
   'p99': 10.15342957642024,
   'mean': 5.170086654425977},
  {'label': 'BAD_unfair',
   'subset': 'long>= 2048',
   'n': 92,
   'p50': 7.038852010499795,
   'p90': 10.96015420079957,
   'p99': 11.058848744559482,
   'mean': 7.140906365467394}],
 [{'label': 'FAIR_768',
   'subset': 'overall',
   'n': 200,
   'p50': 6.04282219300012,
   'p90': 9.18175588419981,
   'p99': 11.05694531869025,
   'mean': 6.266493499729964},
  {'label': 'FAIR_768',
   'subset': 'short<= 64',
   'n': 108,
   'p50': 5.629765399999542,
   'p90': 8.433177871999579,
   'p99': 10.081671103819572,
   'mean': 5.578912106185119},
  {'label': 'FAIR_768',
   'subset': 'long>= 2048',
   'n': 92,
   'p50': 7.1285672229996635,
   'p90': 9.995110139300506,
   'p99': 11.105451201459774,
   'mean': 7.073654266065215}])
```

## Conclusion: btok = 768 and c = 32

- **Overall:** FAIR_768 slightly **hurts median** (p50: 6.04s vs 5.53s) but **improves p90** (9.18s vs 10.46s). **p99 is basically unchanged** (~11.06s in both).
- **Short requests (≤64 tokens):** FAIR_768 **slows p50** (5.63s vs 4.70s) but gives a **small p99 improvement** (10.08s vs 10.15s). Net: better tail for short, worse typical latency.
- **Long requests (≥2048 tokens):** FAIR_768 **improves p90** (10.00s vs 10.96s) while **p99 stays ~the same** (11.11s vs 11.06s). This suggests btok=768 reduces "most long requests" queueing, but the worst-case is still dominated by contention.
- **Implication:** Lowering `max-num-batched-tokens` to 768 shifts the distribution (especially **p90**) but doesn't meaningfully reduce the **worst-case p99** at this concurrency—so a sweep is needed to find the btok "sweet spot" that improves tails without sacrificing too much throughput.

## ˅ `max-num-batched-token` Sweep @ c = 16

```
 1 BASE = "http://127.0.0.1:8000"
 2 MODEL_ID = "Qwen/Qwen2.5-3B-Instruct"
 3
 4 def kill_server():
 5     os.system('pkill -f "vllm.entrypoints.openai.api_server" || true')
 6     os.system('pkill -f "VLLM::EngineCore" || true')
 7     time.sleep(1.5)
 8
 9 def start_fair_server(max_btok, log_path="vllm_fair.log"):
10     # clean + restart
11     kill_server()
```

```
12      try:
13          os.remove(log_path)
14      except FileNotFoundError:
15          pass
16
17      cmd = f"""
18  nohup python3 -m vllm.entrypoints.openai.api_server \
19    --model {MODEL_ID} \
20    --host 127.0.0.1 --port 8000 \
21    --dtype half \
22    --max-model-len 4096 \
23    --gpu-memory-utilization 0.70 \
24    --disable-log-stats \
25    --enforce-eager \
26    --max-num-batched-tokens {int(max_btok)} \
27    --long-prefill-token-threshold 1024 \
28    > {log_path} 2>&1 </dev/null &
29  """
30      os.system(cmd)
31      time.sleep(1.0)
32
33  def tail_log(log_path, n=60):
34      if os.path.exists(log_path):
35          print(f"\n--- tail {log_path} (last {n}) ---")
36          os.system(f"tail -n {n} {log_path}")
37      else:
38          print(f"(log missing: {log_path})")
39
40  def wait_for_ready(base=BASE, log_path="vllm_fair.log", timeout_s=300):
41      t0 = time.time()
42      tries = 0
43      while True:
44          tries += 1
45          try:
46              r = requests.get(f"{base}/v1/models", timeout=2)
47              if r.status_code == 200:
48                  data = r.json().get("data", [])
49                  ids = [m.get("id") for m in data]
50                  print(f"✅ Server ready. Models: {ids}")
51                  return True
52              else:
53                  print(f"[{tries}] /v1/models status={r.status_code}")
54          except Exception as e:
55              print(f"[{tries}] not ready yet: {type(e).__name__}")
56
57          # if process died, fail fast
58          alive = os.system('ps -ef | grep -E "vllm.entrypoints.openai.api_server" | grep -v grep > /dev/null
59          if not alive:
60              print("❌ vLLM process is NOT running anymore.")
61              tail_log(log_path, n=120)
62              return False
63
64          if time.time() - t0 > timeout_s:
65              print("❌ Timed out waiting for server.")
66              tail_log(log_path, n=120)
67              return False
68
69          if tries % 5 == 0:
70              tail_log(log_path, n=30)
71
72          time.sleep(2)
73
74  def warmup_request():
75      # forces one small request to populate caches
76      payload = {
77          "model": MODEL_ID,
78          "messages": [{"role": "user", "content": "hi"}],
79          "max_tokens": 8,
80          "temperature": 0.0
81      }
82      try:
83          _ = requests.post(f"{BASE}/v1/chat/completions", json=payload, timeout=30)
84      except Exception:
85          pass
86
87
88
89  #  Sweep runner
90
91  WL_FILE = "mix_50_50.jsonl"
92  CONCURRENCY = 16
93  TOTAL_REQ = 400
```

```python
 94 DESIRED_MAX_TOKENS = 16
 95
 96 btok_list = [256, 384, 512, 768, 1024, 1536, 2048]  # edit as you like
 97 rows = []
 98 fair_runs = {}
 99
100 for btok in btok_list:
101     print(f"\n===============================")
102     print(f"Starting FAIR server: max-num-batched-tokens={btok}")
103     print(f"===============================")
104
105     start_fair_server(btok, log_path="vllm_fair.log")
106
107     ok = wait_for_ready(BASE, log_path="vllm_fair.log", timeout_s=420)
108     if not ok:
109         rows.append({
110             "label": "FAIR",
111             "workload": WL_FILE,
112             "concurrency": CONCURRENCY,
113             "max_num_batched_tokens": btok,
114             "status": "server_failed"
115         })
116         continue
117
118     warmup_request()
119
120     df_fair, summ_fair = run_load(
121         prompt_sets[WL_FILE],
122         concurrency=CONCURRENCY,
123         total_requests=TOTAL_REQ,
124         desired_max_tokens=DESIRED_MAX_TOKENS,
125         desc=f"FAIR | {WL_FILE} | c={CONCURRENCY} | btok={btok}"
126     )
127
128     fair_runs[btok] = (df_fair, summ_fair)
129
130     m = subset_metrics(df_fair)
131
132     rows.append({
133         "label": "FAIR",
134         "workload": WL_FILE,
135         "concurrency": CONCURRENCY,
136         "max_num_batched_tokens": btok,
137         "status": "ok",
138         "ok_count": int((df_fair["ok"] == True).sum()),
139         "fail_count": int((df_fair["ok"] != True).sum()),
140         "p50_s": float(df_fair[df_fair["ok"]==True]["latency_s"].quantile(0.50)),
141         "p90_s": float(df_fair[df_fair["ok"]==True]["latency_s"].quantile(0.90)),
142         "p99_s": float(df_fair[df_fair["ok"]==True]["latency_s"].quantile(0.99)),
143         "mean_s": float(df_fair[df_fair["ok"]==True]["latency_s"].mean()),
144         **m
145     })
146
147
148 kill_server()
149
150 fair_summary_df = pd.DataFrame(rows)
151 print("\n✅ FAIR btok sweep done.")
152 fair_summary_df
153
```