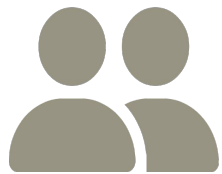


How Cookpad Leverages Triton Inference Server To Boost Their Model Serving





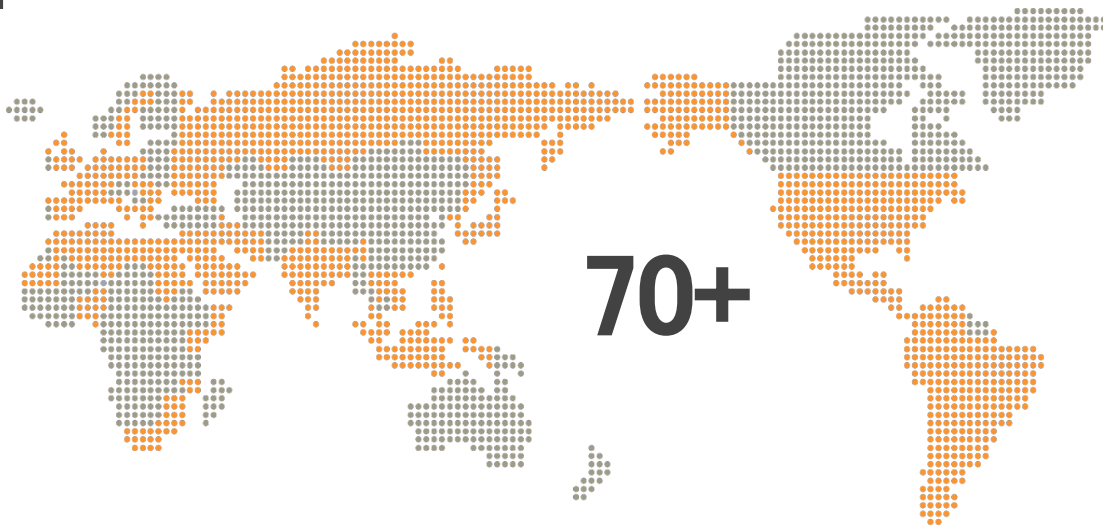
100m +



6m +



30 +



70+

← In season now



Spinach



Potato



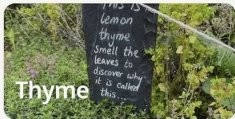
Purple Sprouting Broccoli



Rhubarb



Rosemary



Thyme



Pollock



Wild Garlic



Spring Onion



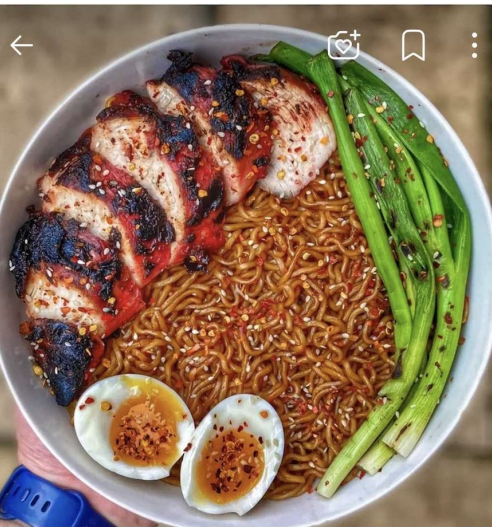
Nettle



Spring Greens



Crab



Char Siu Chicken w/ Hong Kong Style Noodles



Craig Stokes

@whatcraigcooks

📍 Manchester, England, United King...

A Chinese inspired chicken with Hong Kong style noodles, great for any evening

Origin: China

🕒 40 minutes

📷 1 cooksnap

🔍 Search recipes, ingredients or tips

Today's popular searches



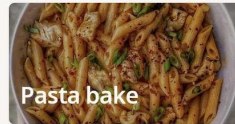
Side dishes



Ukraine



Hong kong



Pasta bake



Spicy food



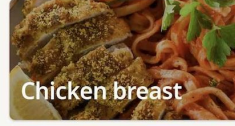
Fish recipes



Cookies



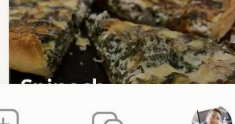
Puff pastry



Chicken breast



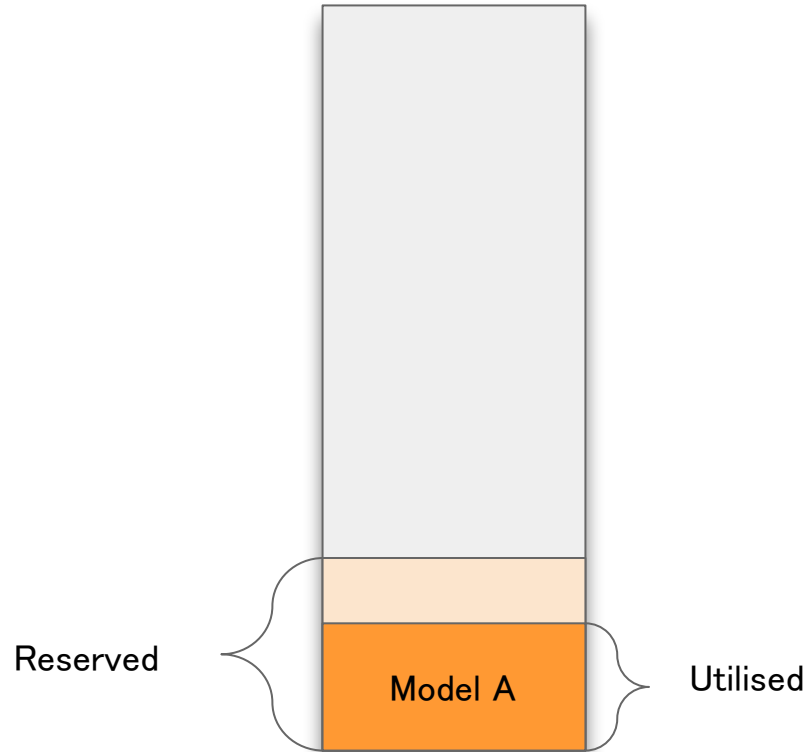
Cauliflower



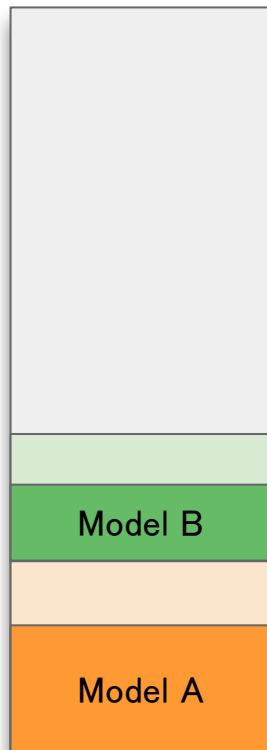
Improving your model deployment with Triton Inference Server

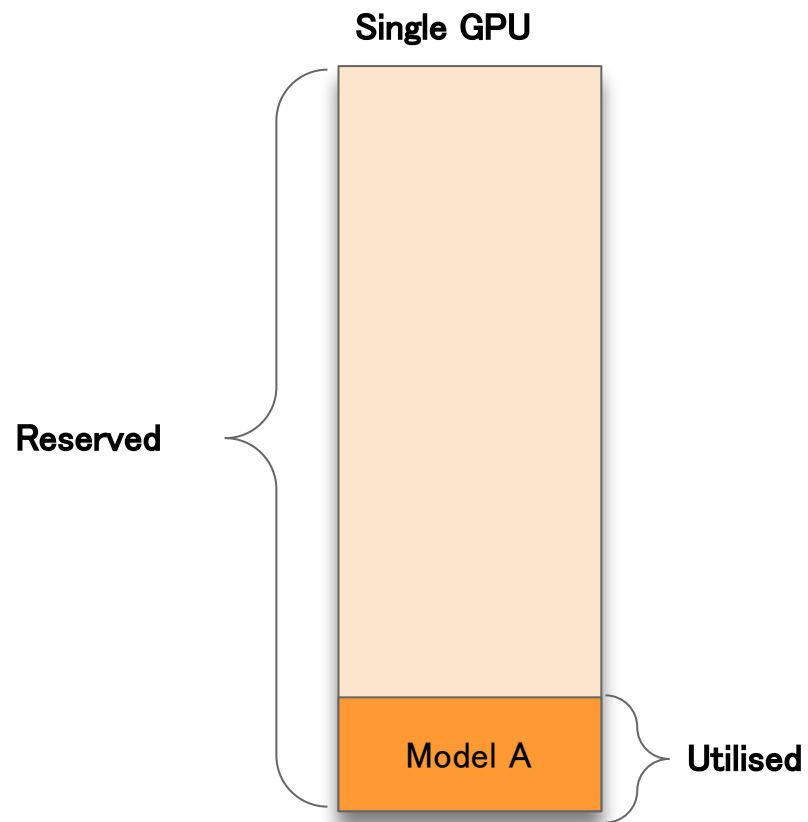


CPU / Memory

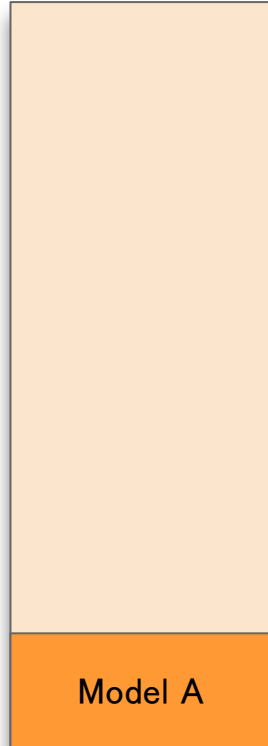


CPU / Memory





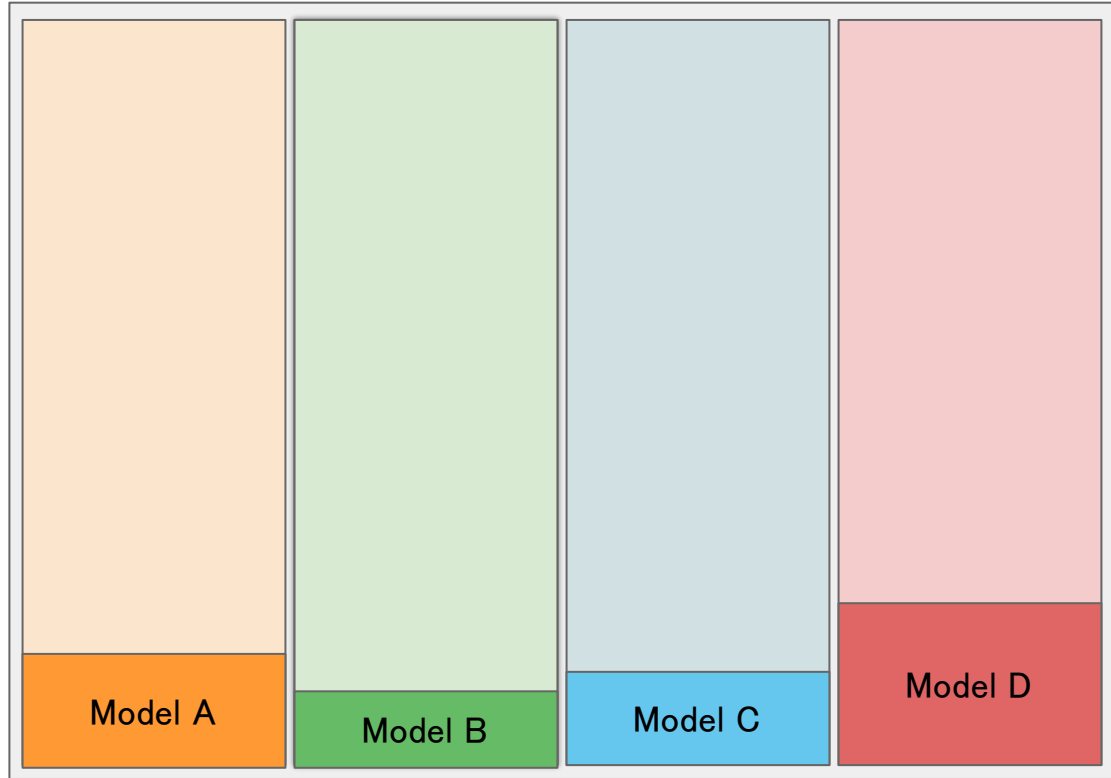
Single GPU



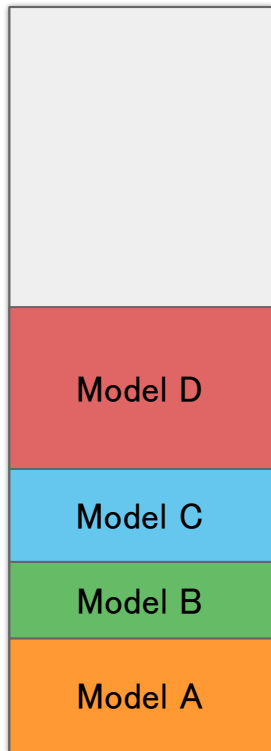
Single GPU



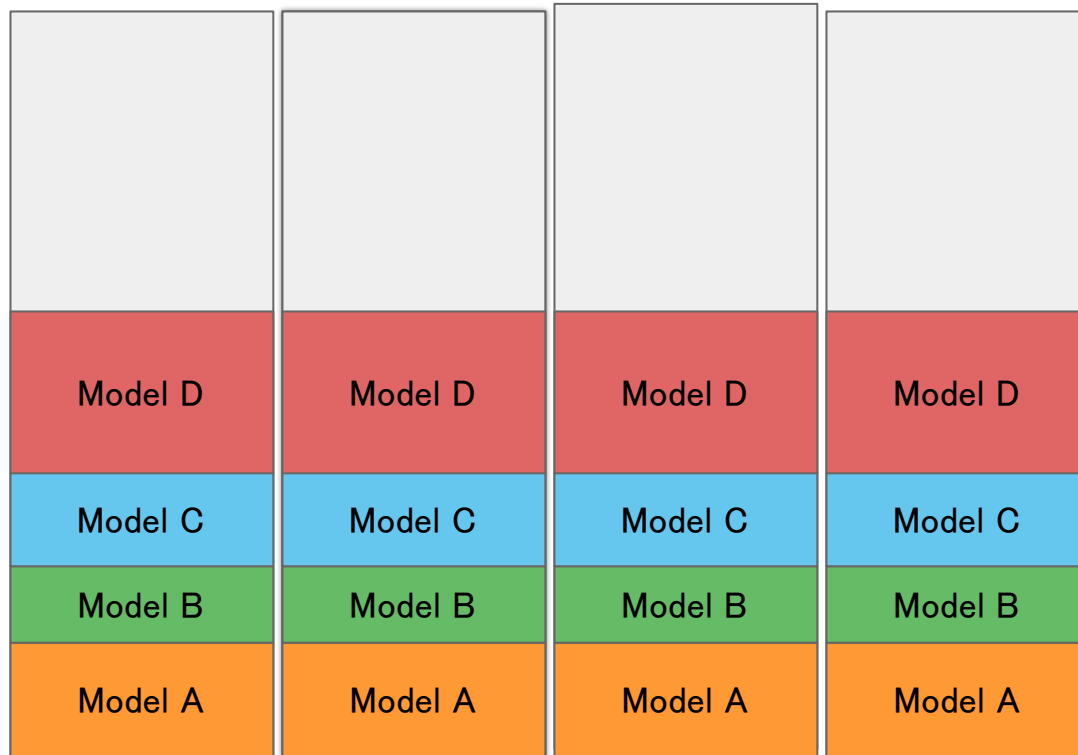
Multiple GPUs



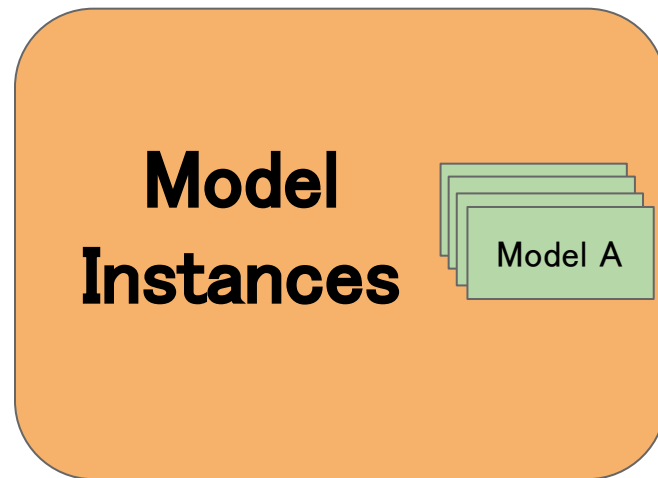
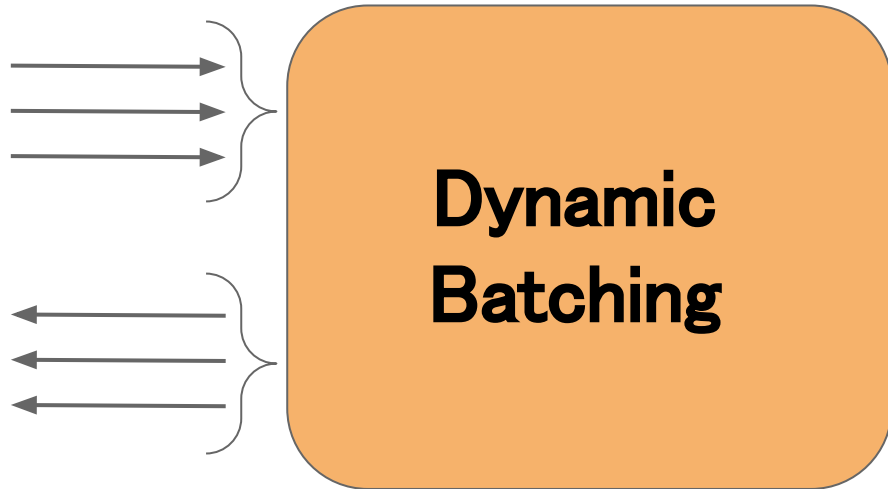
Triton on single GPU

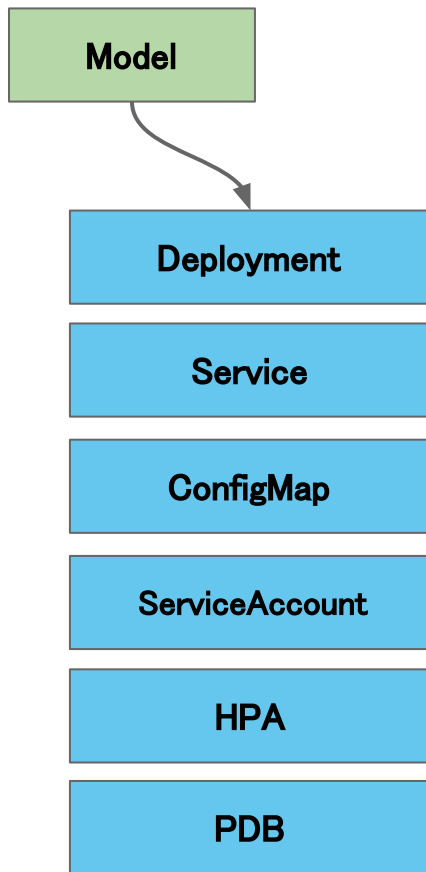


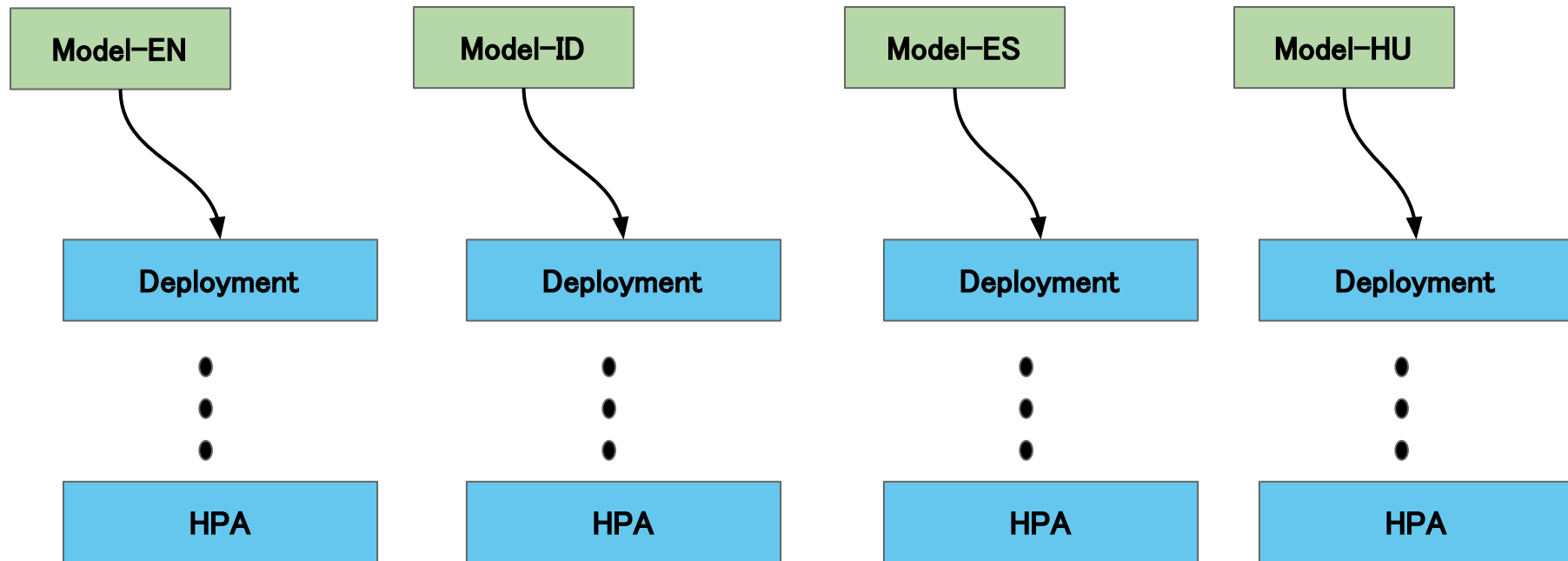
Triton Multiple GPUs











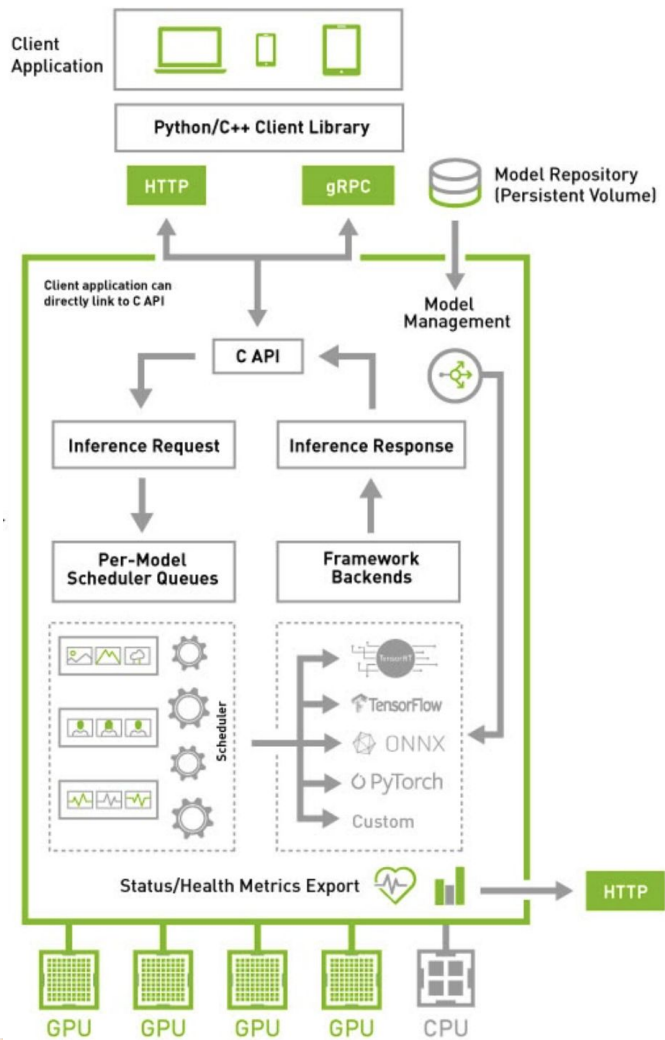
```
local_model_dir: /models
```

```
model_files:
```

- text-embeddings/en/model.tar.gz
- text-embeddings/es/model.tar.gz
- text-embeddings/id/model.tar.gz
- text-embeddings/hu/model.tar.gz
- image-embeddings/model.tar.gz
- categorisations/optimized/model.tar.gz
- onnx_xlm_roberta_base.tar.gz

```
s3_bucket: <bucket_name>
```

```
s3_models_dir: deployment-artifacts
```

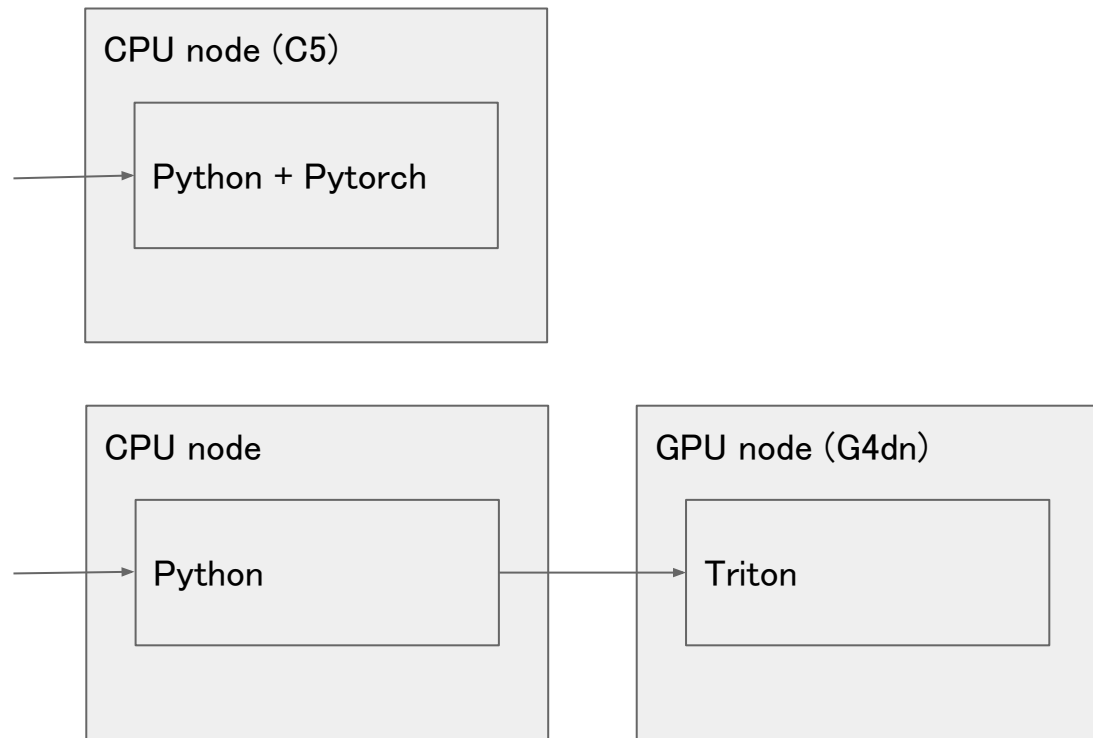




[3, 244, 244]



[300]



Model



Deployment

Service

ConfigMap

ServiceAccount

HPA

PDB

config.yaml

✓ image_embedding

✓ 20190320110924

└─ model.plan

≡ config.pbtxt

≡ config.pbtxt ×

image_embedding > ≡ config.pbtxt

```
1 platform: "tensorrt_plan"
2 max_batch_size: 1
3 input [
4     {
5         name: "tensor"
6         data_type: TYPE_FP32
7         dims: [ 3, 224, 224 ]
8     }
9 ]
10 output [
11     {
12         name: "1930"
13         data_type: TYPE_FP32
14         dims: [ 300 ]
15     }
16 ]
17
```

! config.yaml ×

! config.yaml

```
1  local_model_dir: /models
2  model_files:
3  |  - image_embedding.tar.gz
4  s3_bucket: ember.us-east-1
5  s3_models_dir: packages
6
```


| | | | | | |
|-------------|-----|-----|------|-----|-----|
| Processing: | 257 | 272 | 14.1 | 269 | 337 |
| Waiting: | 254 | 272 | 14.1 | 269 | 337 |
| Total: | 257 | 272 | 14.1 | 269 | 337 |

Percentage of the requests served within a certain time (ms)

| | |
|------|-----------------------|
| 50% | 269 |
| 66% | 273 |
| 75% | 277 |
| 80% | 279 |
| 90% | 286 |
| 95% | 297 |
| 98% | 327 |
| 99% | 337 |
| 100% | 337 (longest request) |

root@ubuntu: /#

| | | | | | |
|-------------|----|----|-----|----|----|
| Processing: | 17 | 18 | 1.0 | 18 | 22 |
| Waiting: | 17 | 18 | 0.9 | 18 | 22 |
| Total: | 17 | 19 | 1.0 | 19 | 23 |

Percentage of the requests served within a certain time (ms)

| | |
|------|----------------------|
| 50% | 19 |
| 66% | 19 |
| 75% | 19 |
| 80% | 20 |
| 90% | 20 |
| 95% | 21 |
| 98% | 23 |
| 99% | 23 |
| 100% | 23 (longest request) |

root@ubuntu: /#


```
Processing: 269 538 43.0 538 631
Waiting:    268 537 42.7 538 631
Total:      269 538 43.0 538 631
```

Percentage of the requests served within a certain time (ms)

```
50%    538
66%    549
75%    552
80%    553
90%    568
95%    576
98%    631
99%    631
100%   631 (longest request)
```

root@ubuntu: /#

```
Total:      18 20 1.6 19 29
```

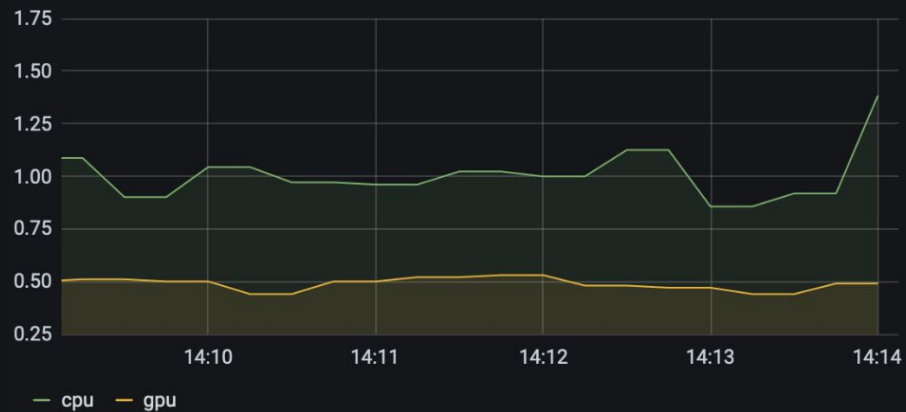
ERROR: The median and mean for the initial connection time are more than twice the standard deviation apart. These results are NOT reliable.

Percentage of the requests served within a certain time (ms)

```
50%    19
66%    20
75%    20
80%    20
90%    22
95%    23
98%    27
99%    29
100%   29 (longest request)
```

root@ubuntu: /#

CPU/GPU



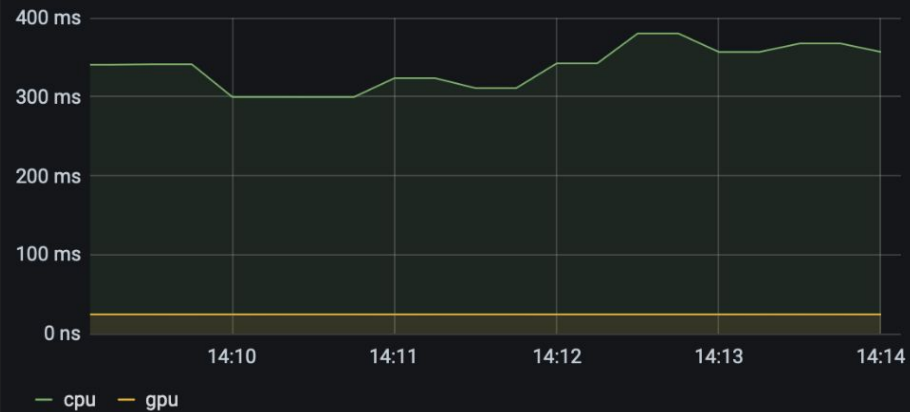
Memory Usage



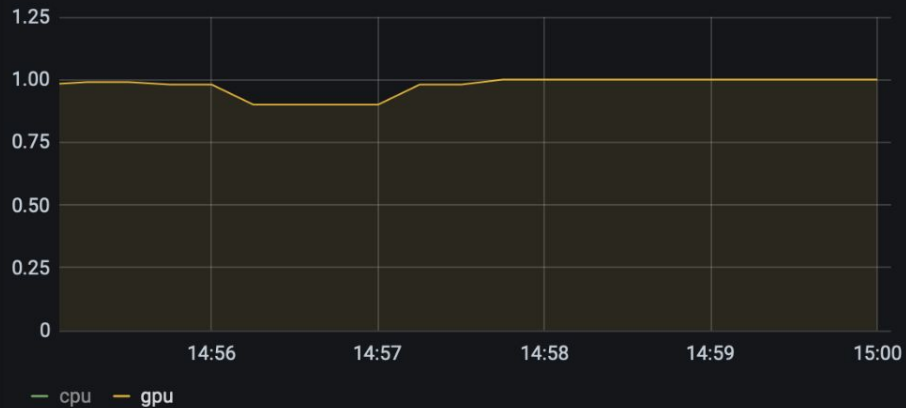
Request Rate



P95 Latency



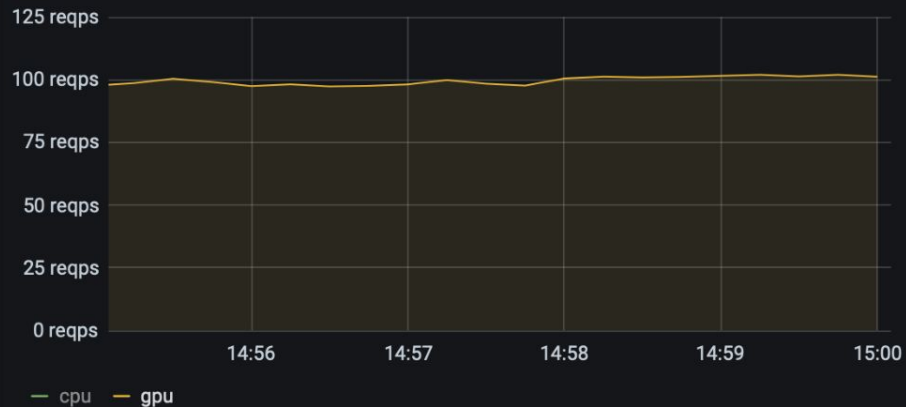
CPU/GPU



Memory Usage



Request Rate



P95 Latency



≡ config.pbtxt ✕

models > image_embedding > ≡ config.pbtxt

```
1 platform: "tensorrt_plan"
2 max_batch_size: 8
3 > input [...]
9 ]
10 > output [...]
16 ]
17 ∨ instance_group [
18 ∨ {
19   · count: 2
20   · kind: KIND_GPU
21   · gpus: [ 0 ]
22 }
23 ]
```

≡ config.pbtxt ✕

models > image_embedding > ≡ config.pbtxt

```
1 platform: "tensorrt_plan"
2 max_batch_size: 8
3 > input [...]
9 ]
10 > output [...]
16 ]
17 ∨ dynamic_batching {
18   · preferred_batch_size: [ 4, 8 ],
19   · max_queue_delay_microseconds: 100
20 }
```

Q & A

We are hiring

<https://careers.cookpad.com>

