

# Synthèse des résultats d'analyse de Graphana et k6

Bryan Joya

JOYB74080104

## 1. Analyse de la latence avant et après le Load Balancer

Avant le load balancer:

Lancement du test de charge Load-test.js :

```
THRESHOLDS

http_req_duration
X 'p(95)<500' p(95)=7.01s

http_req_failed
✓ 'rate<0.01' rate=0.00%

TOTAL RESULTS

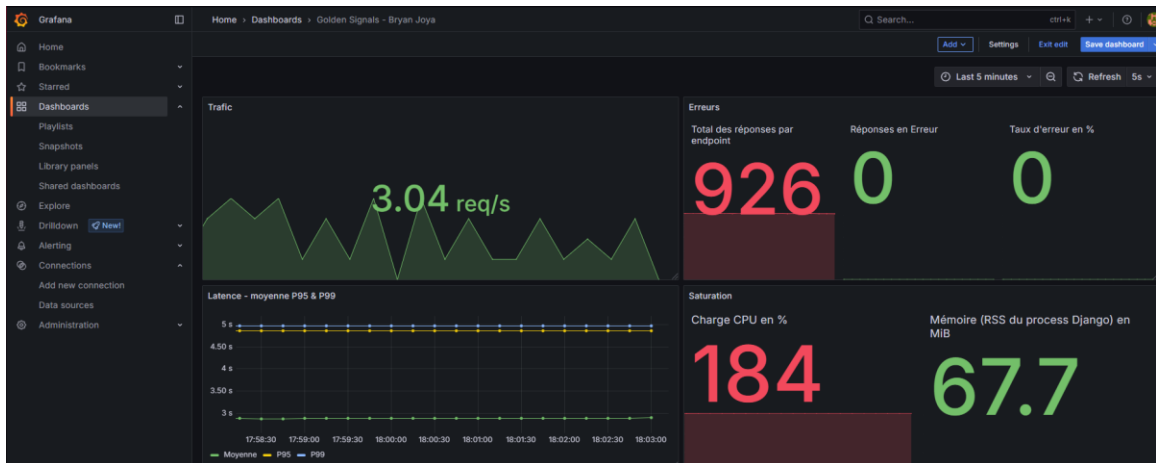
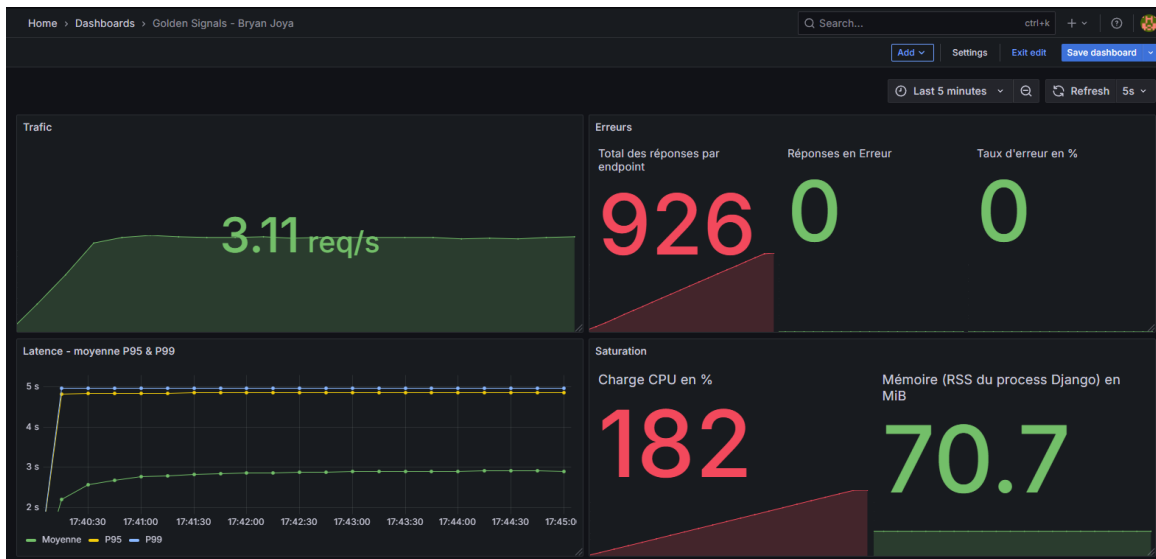
checks_total.....: 921      2.980015/s
checks_succeeded.....: 100.00% 921 out of 921
checks_failed.....: 0.00%  0 out of 921

✓ stock 200
✓ rapport 200
✓ update 200

HTTP
http_req_duration.....: avg=6.29s min=1.54s med=6.4s max=7.96s p(90)=6.85s
{ expected_response:true }.....: avg=6.29s min=1.54s med=6.4s max=7.96s p(90)=6.85s
http_req_failed.....: 0.00%  0 out of 921
http_reqs.....: 921      2.980015/s

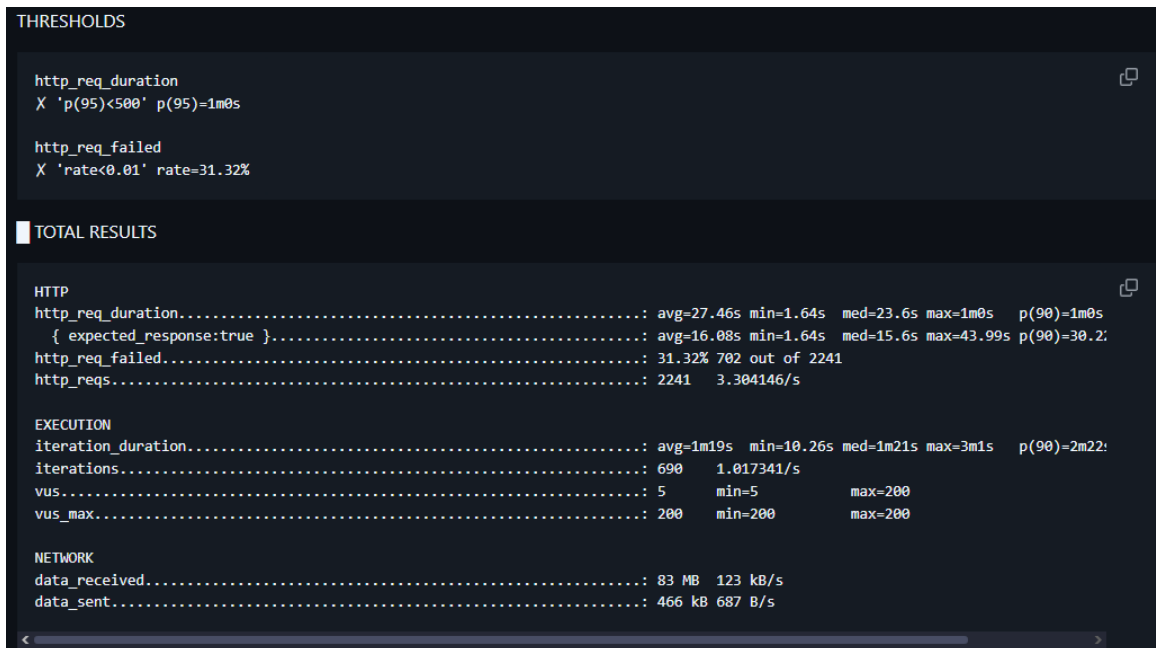
EXECUTION
iteration_duration.....: avg=19.87s min=9.4s med=20.11s max=21.42s p(90)=20.8s
iterations.....: 307      0.993338/s
vus.....: 2      min=2      max=20
vus_max.....: 20      min=20      max=20

NETWORK
data_received.....: 2.4 MB 7.7 kB/s
data_sent.....: 183 kB 592 B/s
```



Avant son déploiement, les tests de charge indiquaient une latence élevée, avec une durée moyenne de requête à 6,29 secondes et un pic à 7,96 secondes. Toutefois, aucun échec de requête n'a été détecté. Le nombre total de requêtes (926) et la vitesse d'exécution (3.11req par secondes) indiquent une activité soutenue. On peut observer que 184% du CPU est utilisé par la machine.

Lancement du test de charge Load-test-2.js :



Sous de 200 utilisateurs concurrents, l'application ne satisfait ni le critère de latence (p(95) < 500 ms) ni celui du taux d'échec (rate < 1 %) : en effet, 31 % des requêtes sont interrompues et la latence du 95<sup>e</sup> centile atteint systématiquement 60 s, valeur qui correspond au timeout par défaut de k6. Ces timeouts massifs indiquent que, sous cette charge, le serveur ne peut pas traiter les requêtes dans les délais et que beaucoup d'itérations restent inachevées.

## Analyse des points faibles et recommandations d'optimisation

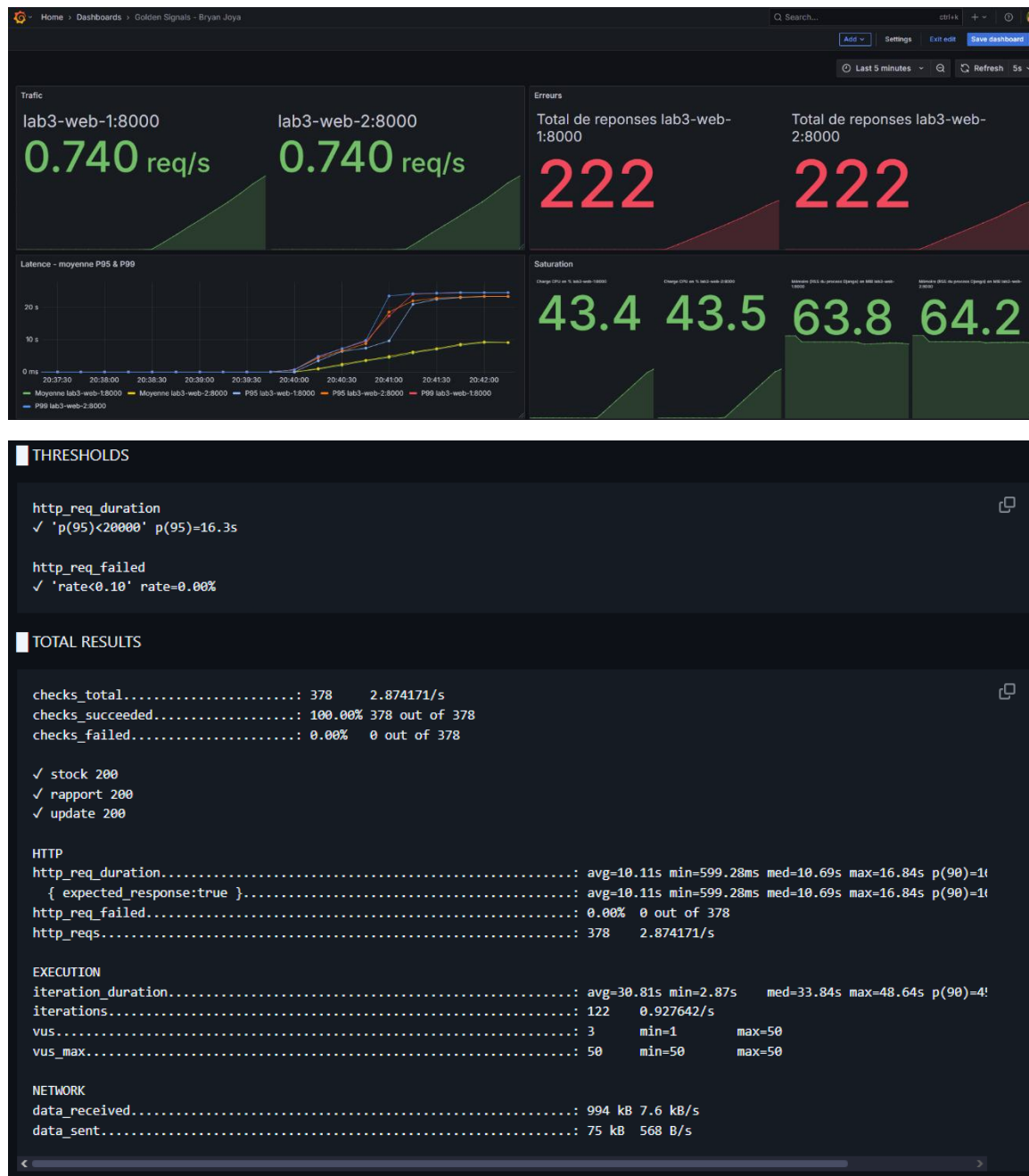
- Pool de connexions saturé
  - Nombre de connexions ouvertes atteint rapidement le maximum configuré → file d'attente, latences accrues.
  - Pas de timeout ni de « recycle » sur les connexions inactives.
  - Absence ou insuffisance d'index
  - Requêtes de lecture sur les tables volumineuses (produits, stocks, rapports) déclenchent des parcours séquentiels (table scan).
- Requêtes SQL sous-optimales
- Absence de mise en cache
  - Appels répétés au même endpoint /rapport/ ou à la même ressource stock/produit sans cache → requêtes identiques refaites à chaque VU.
  - Aucun mécanisme de cache HTTP ni cache applicatif (Redis, in-memory).

### Recommandations d'amélioration

- Pool de connexions
  - Activer un « connection lifetime » pour fermer les connexions longues
- Mise en cache
  - Côté API, ajouter un cache en mémoire ou Redis pour les réponses fréquentes.
  - Activer le cache HTTP pour permettre aux clients et aux proxy de conserver les résultats.
  - Mettre en place un cache de requêtes SQL pour éviter la recompilation répétée.

## Après le load balancer CONFIGURATION ROUND ROBIN:

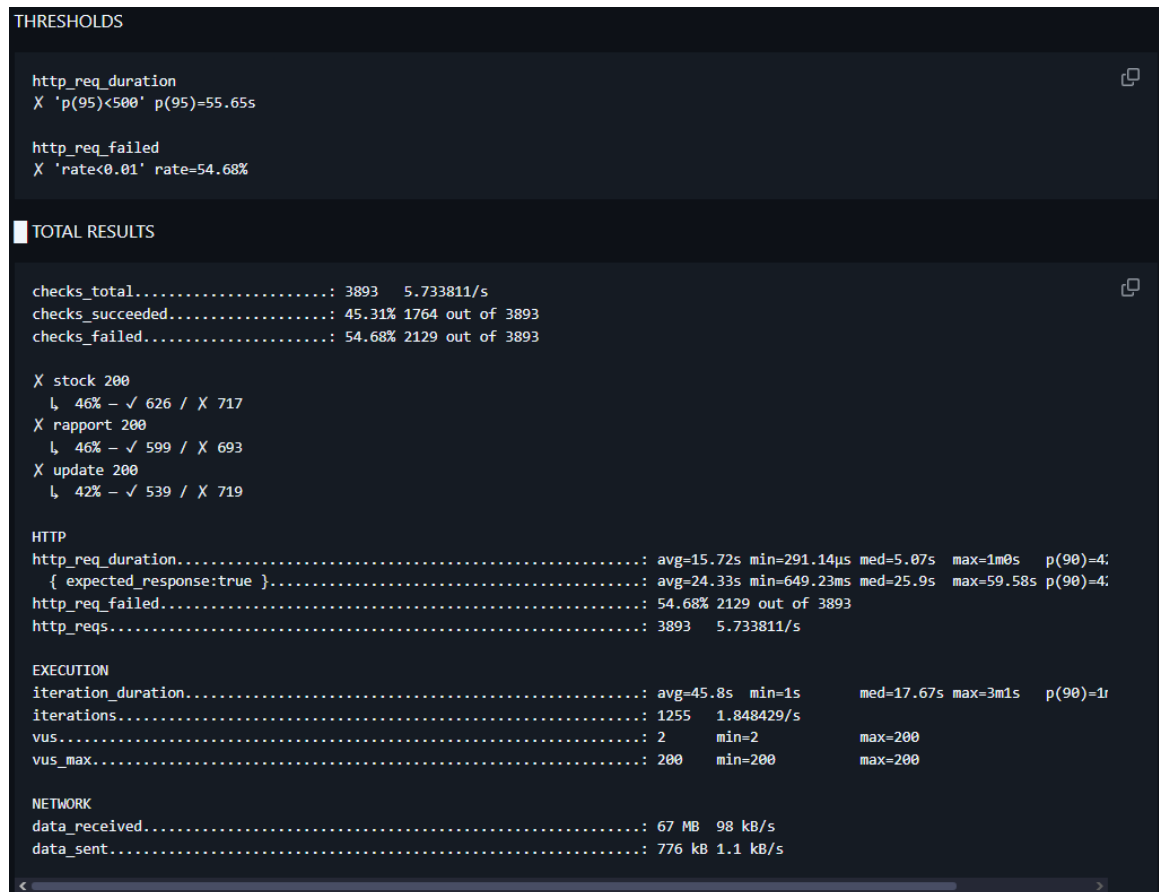
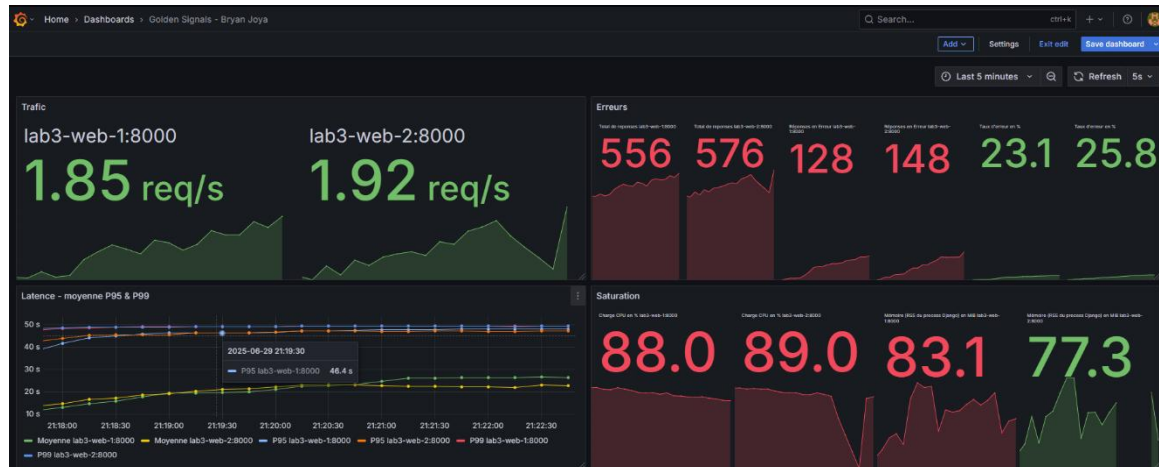
### Lancement du load-test.js :



On observe une répartition équilibrée du trafic entre les deux serveurs (0,740 requêtes/sec chacun), ce qui montre que la charge est bien distribuée. Le taux d'erreurs reste nul, avec 378 vérifications réussies et aucune échouée. La latence moyenne des requêtes HTTP est de 10,11 secondes, avec une valeur maximale à 16,84 secondes, ce qui respecte les seuils définis (p95 < 20s). Les indicateurs de saturation et d'utilisation des ressources sont dans des plages acceptables, reflétant une bonne gestion du flux réseau et des performances

optimisées. Enfin, aucun échec n'a été détecté et toutes les opérations (stock, rapport, update) ont retourné un code 200, signalant leur succès.

## Lancement du load-test-2.js :



On peut voir un taux d'erreur d'environ 25% par instance avec une latence moyenne de 24 secondes.

## 2. Évolution du trafic avec plusieurs instances : requêtes par seconde CONFIGURATION ROUND ROBIN

2 INSTANCES - LOAD-TEST.JS



Avec 2 instances, le système atteint environ **1,5 requêtes par seconde**, chacune gérant la moitié de la charge.

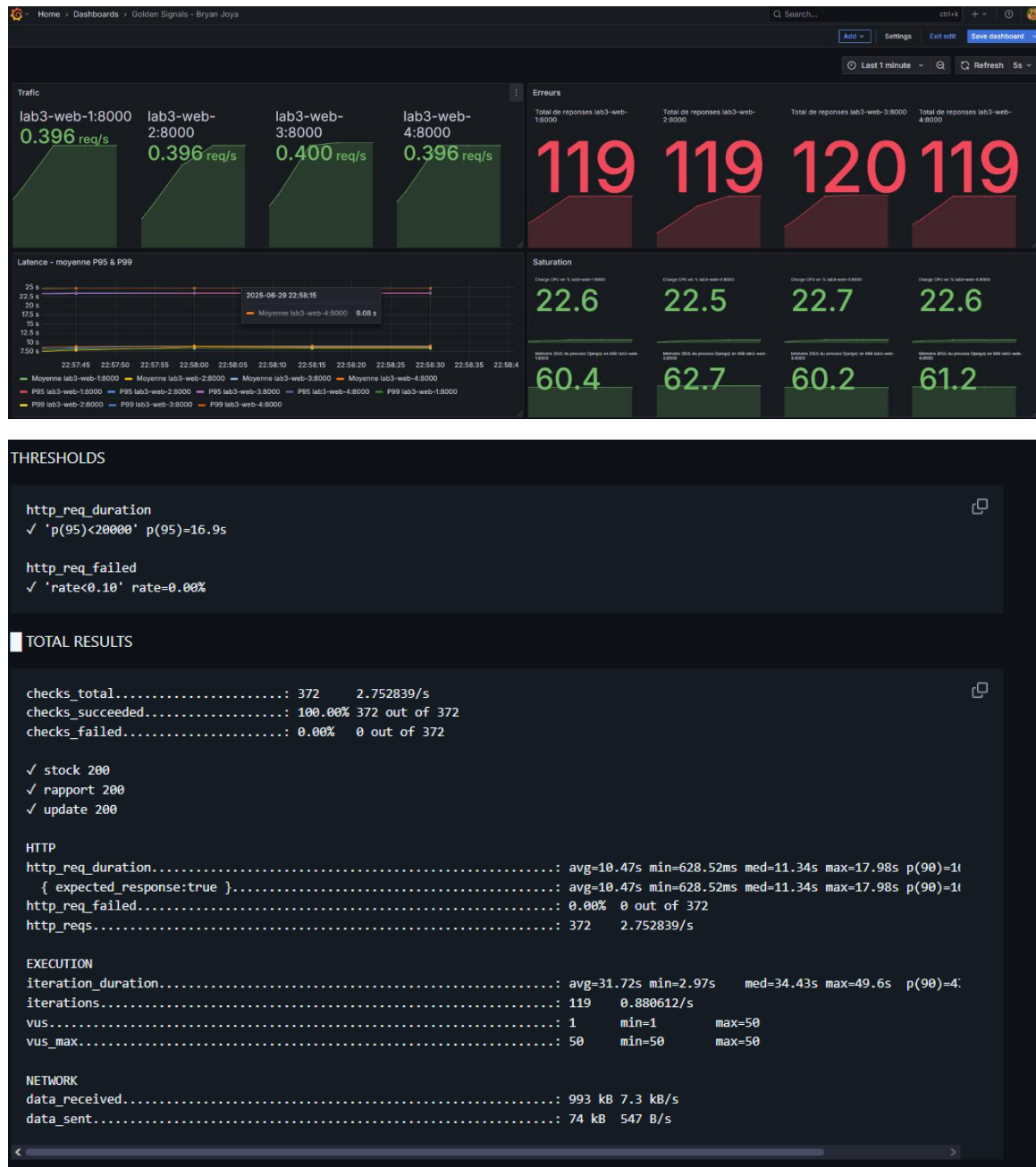
### 3 INSTANCES - LOAD-TEST.JS



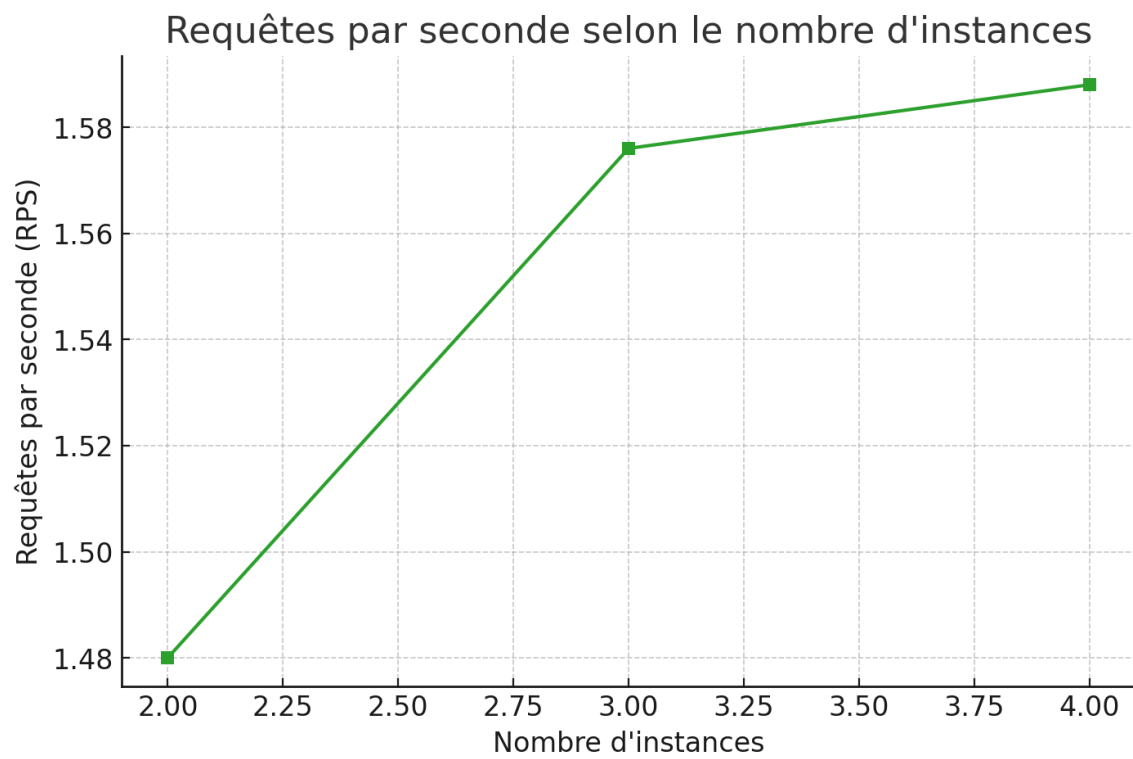
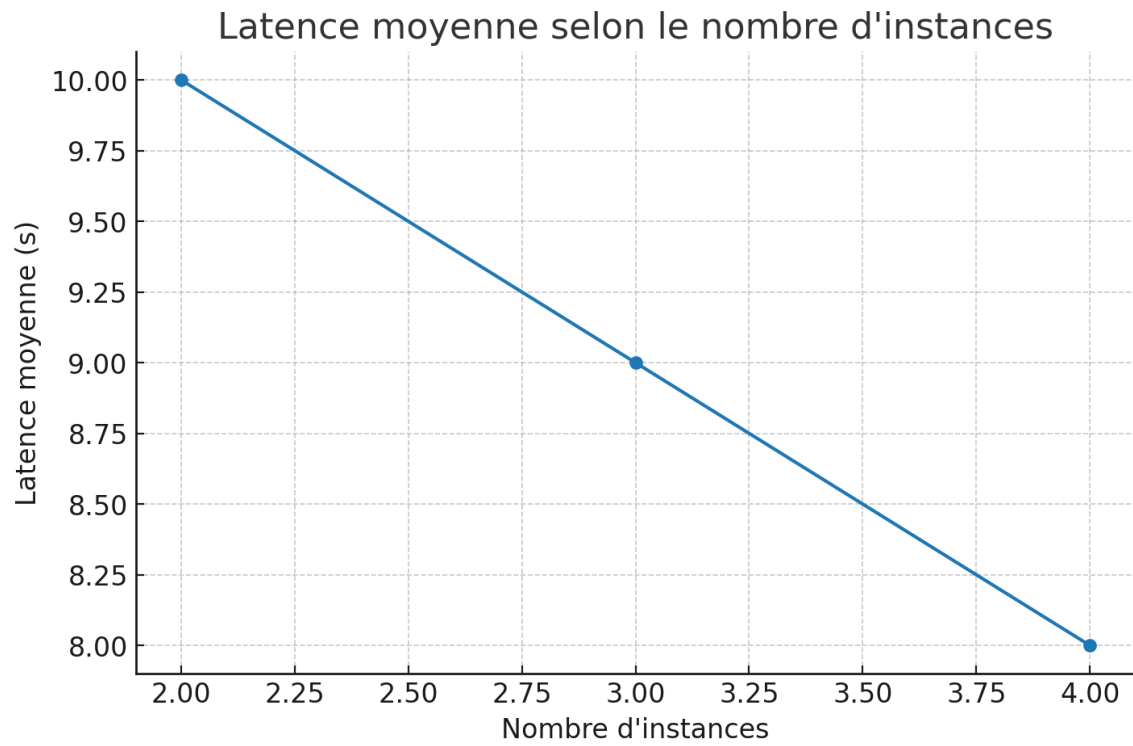
Avec 3 instances, la charge est encore mieux répartie, chaque serveur prenant environ un tiers des requêtes. On observe une légère amélioration des performances globales avec plus de réponses traitées. Mais le temps de réponse été le temps d'itération augmente légèrement.

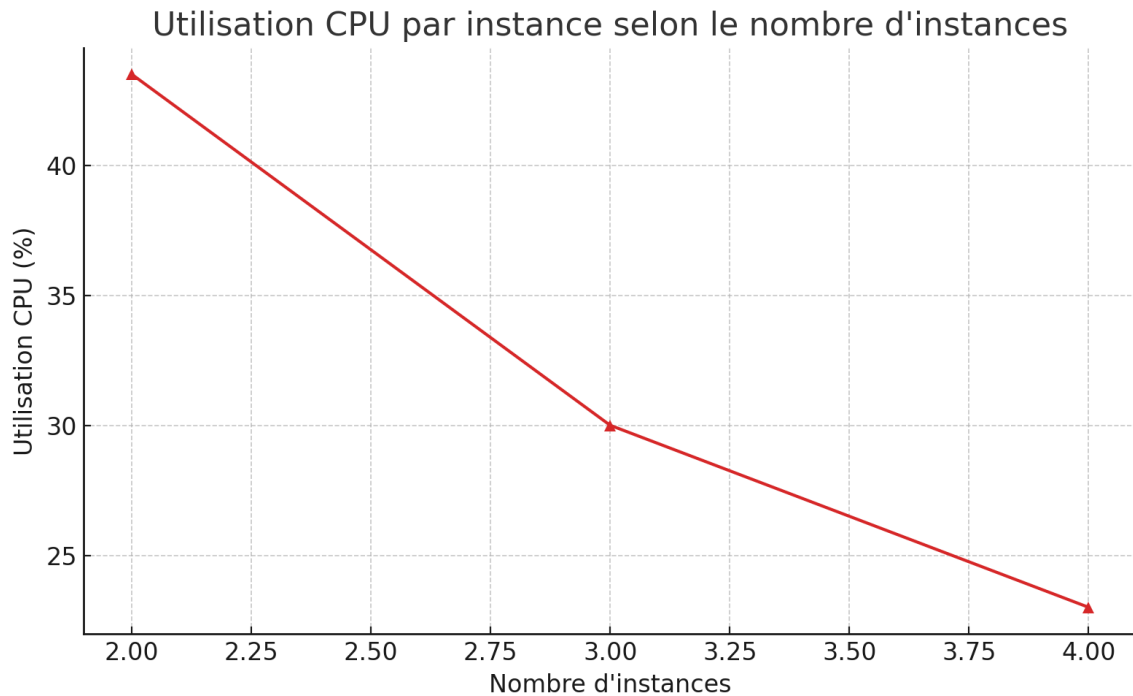


## 4 INSTANCES - LOAD-TEST.JS



**Avec 4 instances**, le système devient plus stable et les ressources CPU sont moins saturées. Il y a toujours une augmentation de la durée des itérations et un gain léger en réponses totales traitées.





Avec 2 instances, le système traite environ **1.48 requêtes par seconde**, totalisant **444 réponses**, avec une **latence moyenne de 10 secondes** et une **utilisation CPU de 43,5 % par instance**. En passant à 3 instances, les performances s'améliorent : **1.576 requêtes par seconde**, **474 réponses** au total, une **latence réduite à 9 secondes** et un **CPU à 30 %**. Avec 4 instances, on atteint **1.588 requêtes par seconde**, **477 réponses**, une **latence moyenne de 8 secondes** et seulement **23 % d'utilisation CPU**. Ces résultats montrent que l'ajout d'instances permet une meilleure répartition de la charge, une réduction de la latence et une baisse de la pression sur chaque instance.

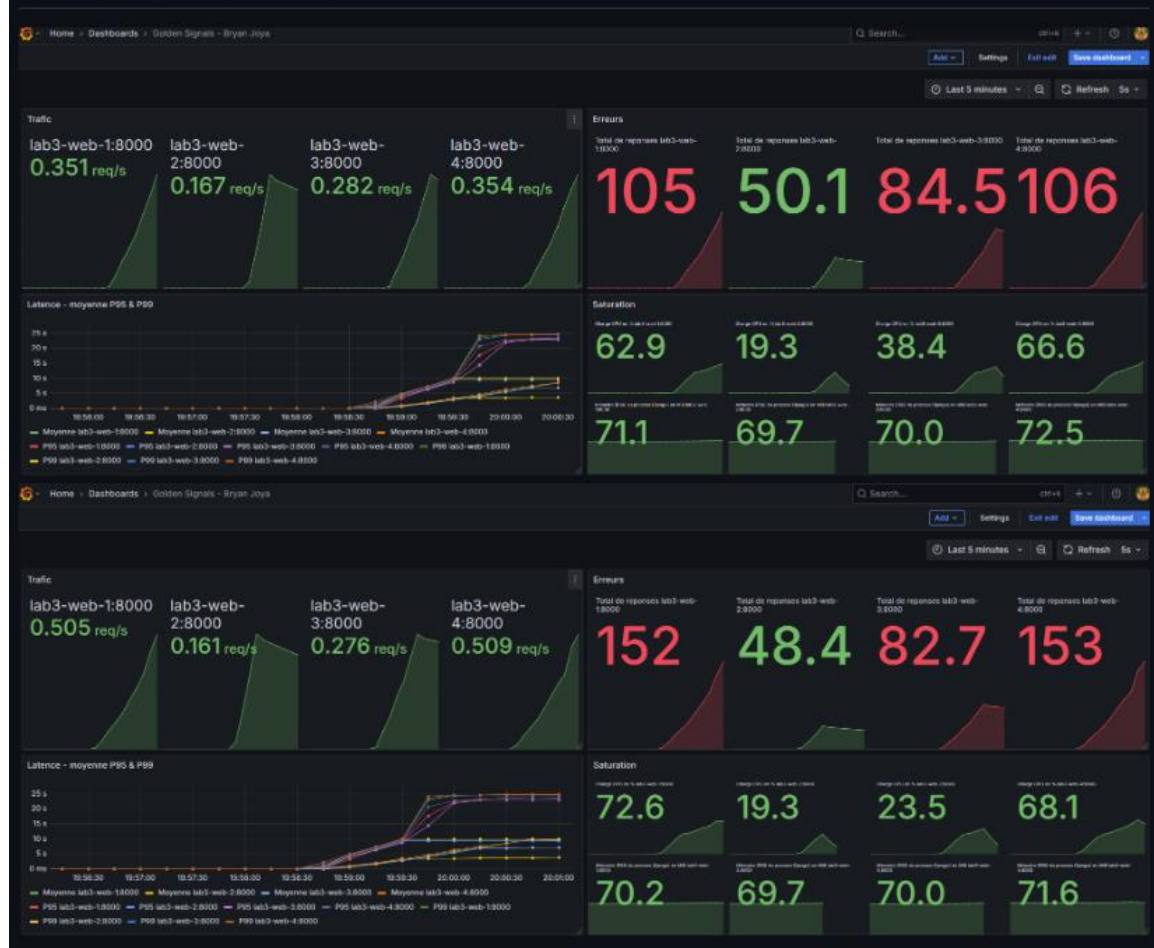
### 3. Test de tolérance aux pannes

LOAD-TEST.JS

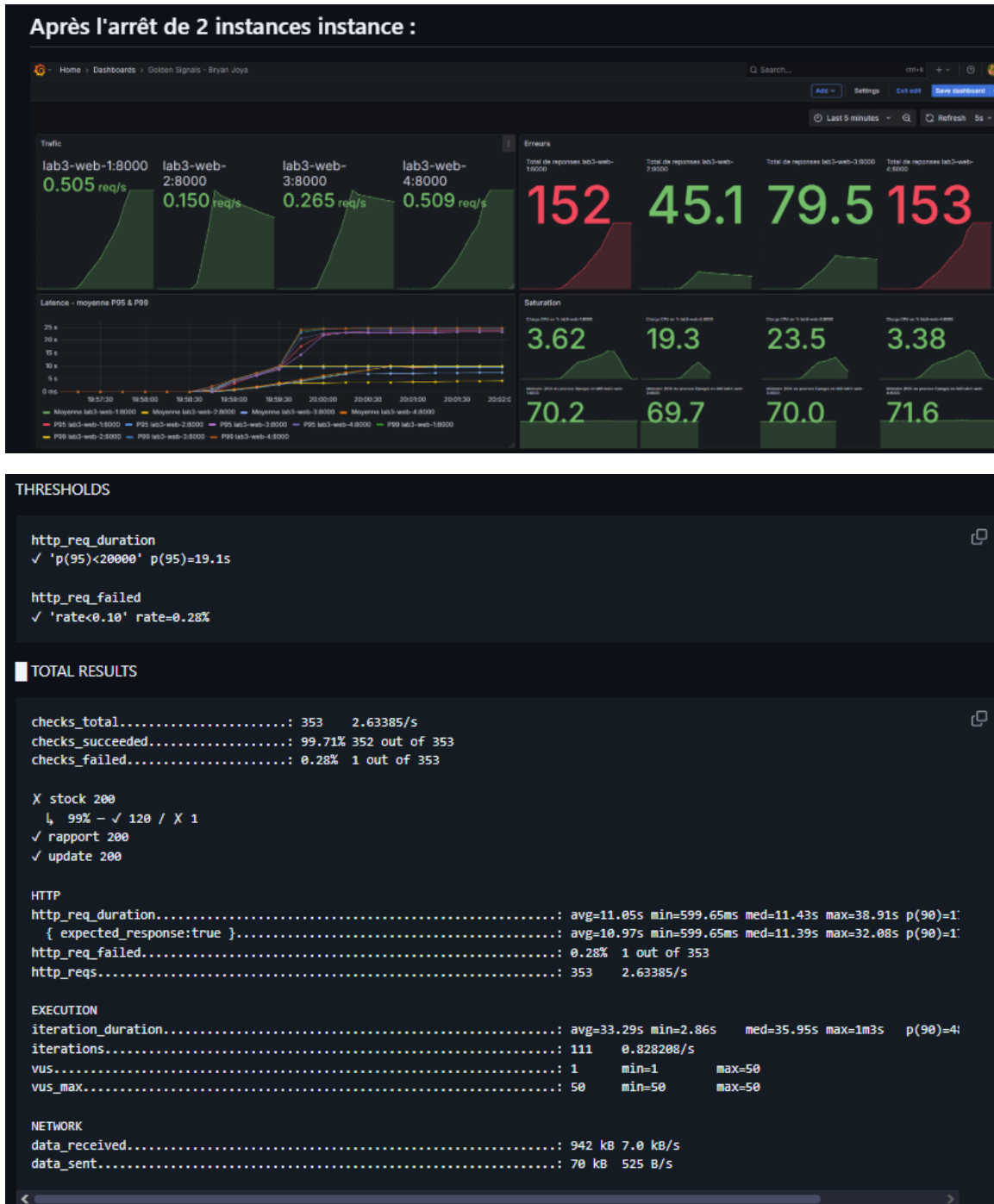
#### Avec 4 instances fonctionnelles



#### Après l'arrêt d'une instance :



## LOAD-TEST.JS

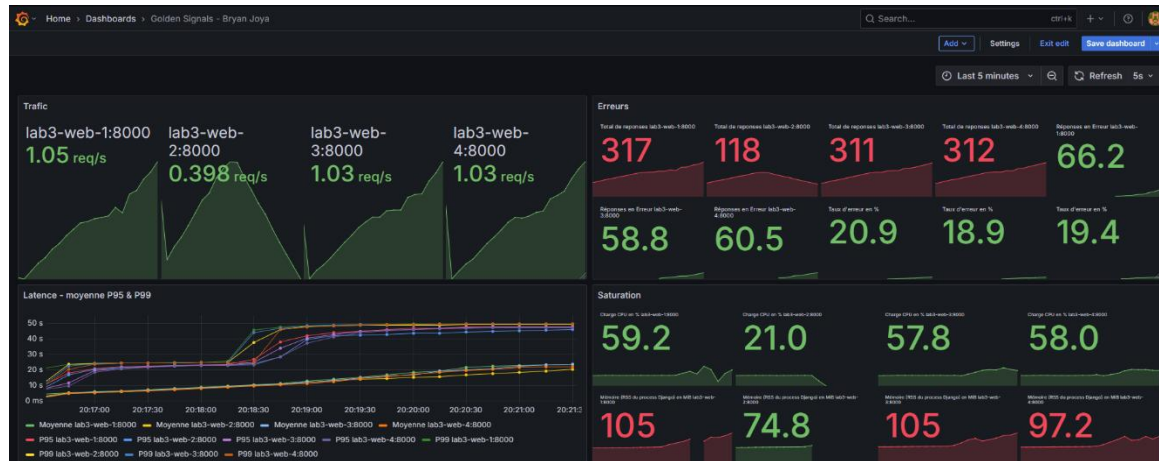


On peut observer la répartition des requêtes après l'arrêt du premier et du deuxième conteneur. La panne a engendré 1 erreur de visualisation de stock.

## LOAD-TEST-2.JS

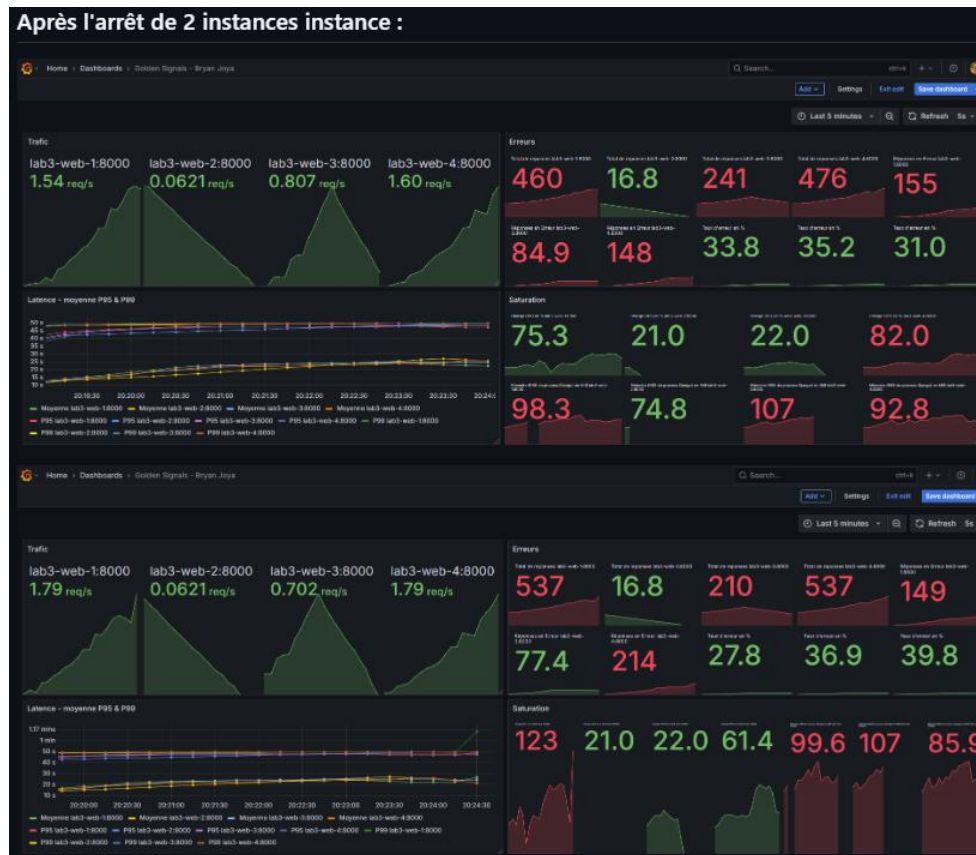


## ARRÊT D'UNE INSTANCE





## LOAD-TEST-2.JS



```
THRESHOLDS

http_req_duration
X 'p(95)<500' p(95)=1m0s

http_req_failed
X 'rate<0.01' rate=43.10%

TOTAL RESULTS

checks_total.....: 2960 4.369191/s
checks_succeeded.....: 56.89% 1684 out of 2960
checks_failed.....: 43.10% 1276 out of 2960

X stock 200
  ↳ 58% - ✓ 689 / X 429
X rapport 200
  ↳ 56% - ✓ 560 / X 423
X update 200
  ↳ 54% - ✓ 515 / X 424

HTTP
http_req_duration.....: avg=20.63s min=332.94µs med=20.13s max=1m0s p(90)=4i
{ expected_response:true }.....: avg=22.52s min=635.52ms med=22.84s max=59.96s p(90)=4i
http_req_failed.....: 43.10% 1276 out of 2960
http_reqs.....: 2960 4.369191/s

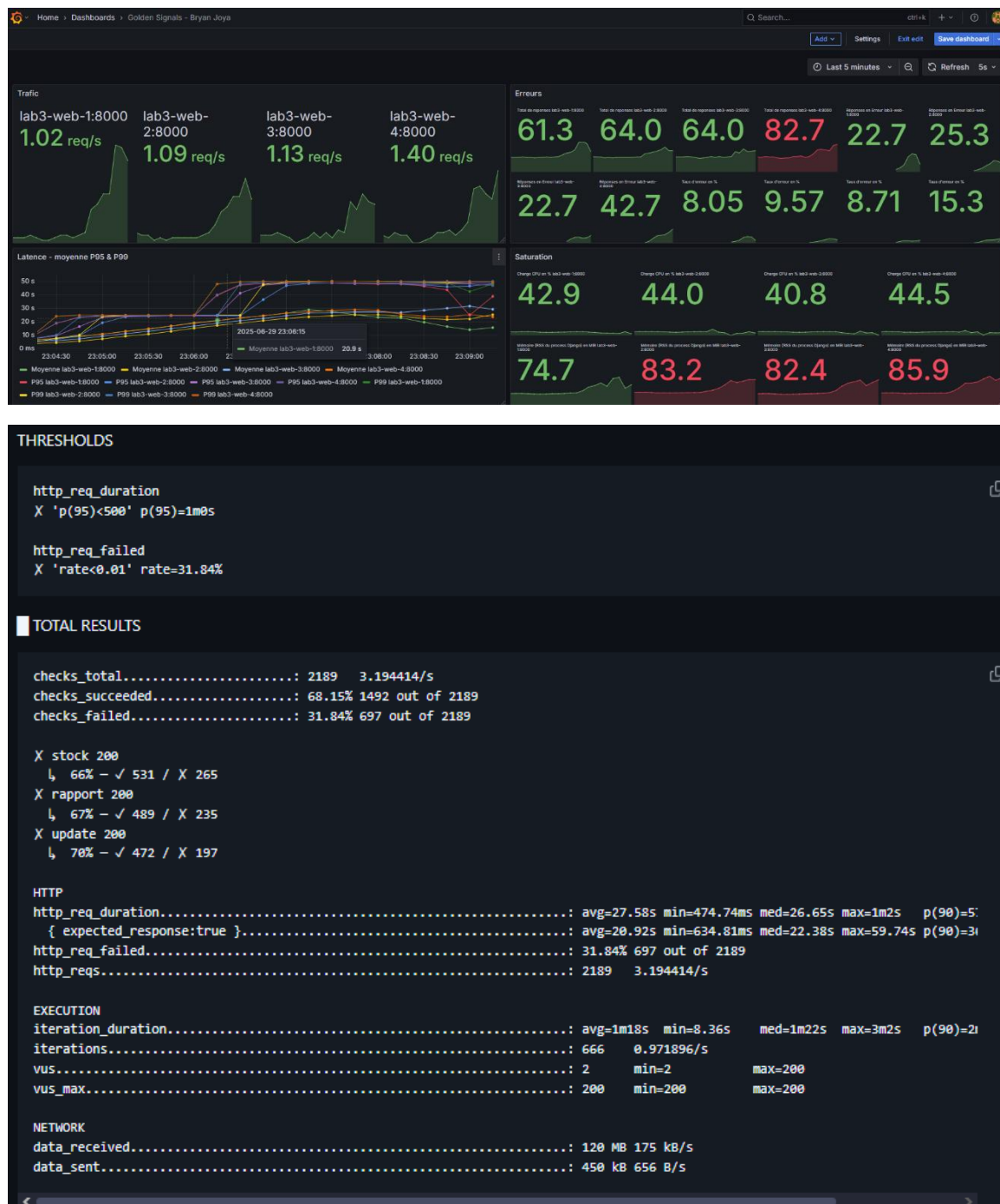
EXECUTION
iteration_duration.....: avg=59.22s min=1s med=1m3s max=3m1s p(90)=2i
iterations.....: 937 1.383085/s
vus.....: 4 min=4 max=200
vus_max.....: 200 min=200 max=200

NETWORK
data_received.....: 115 MB 169 kB/s
data_sent.....: 595 KB 878 B/s
```

Les pics dans les graphes démontrent l'arrêt des deux conteneurs. On peut observer la répartition de charge vers les deux autres conteneurs opérationnels.

## 4. Analyse des stratégies de Load Balancing

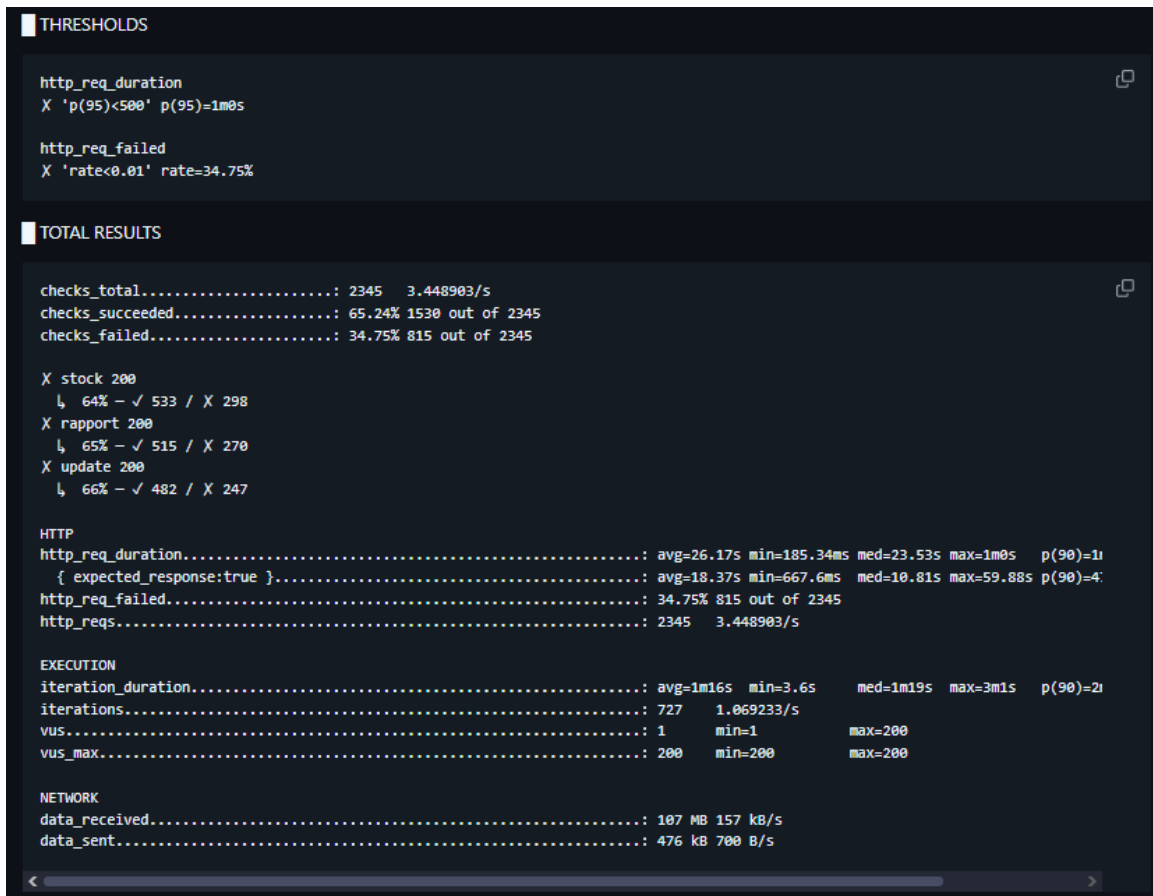
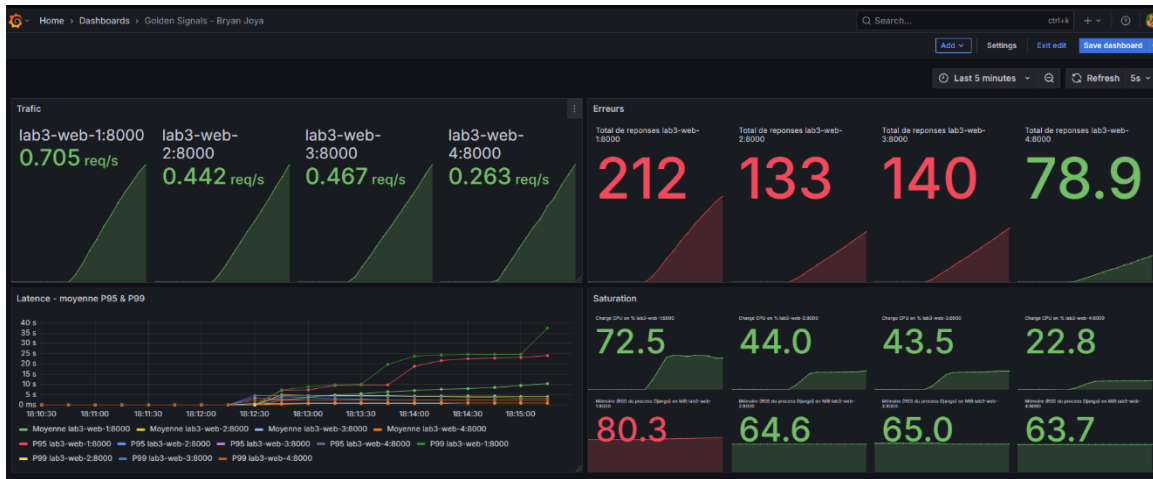
### ROUND ROBIN



Répartition équitable et simple, efficace pour une charge uniforme.

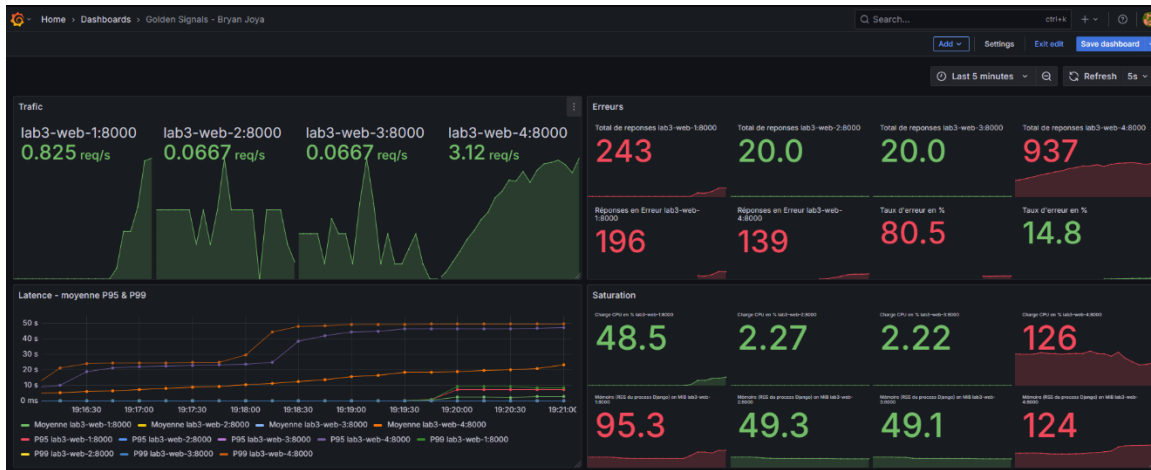


## WEIGHTED ROUND ROBIN : 4 – 2 – 2 – 1



**Weighted Round Robin** : permet de privilégier certains serveurs (poids 4 – 2 – 2 – 1), utile si les instances ont des puissances inégales. Cela se reflète dans le nombre de requêtes par secondes traitées par les instances (0.7)-(0.4)-(0.4)-(0.25). Un taux d'erreur plus élevé que celui du round robin est présent.

## IP HASH : EXÉCUTÉ EN LOCAL DONC UNE SEULE IP



### THRESHOLDS

```
http_req_duration
X 'p(95)<500' p(95)=1m0s

http_req_failed
X 'rate<0.01' rate=33.12%
```

### TOTAL RESULTS

```
checks_total.....: 2427 3.548472/s
checks_succeeded.....: 66.87% 1623 out of 2427
checks_failed.....: 33.12% 804 out of 2427

X stock 200
  ↳ 64% - ✓ 567 / X 310
X rapport 200
  ↳ 68% - ✓ 553 / X 252
X update 200
  ↳ 67% - ✓ 503 / X 242

HTTP
http_req_duration.....: avg=25.49s min=132.31ms med=22.46s max=1m0s p(90)=1i
{ expected_response:true }.....: avg=22.47s min=757.51ms med=21.09s max=59.84s p(90)=4i
http_req_failed.....: 33.12% 804 out of 2427
http_reqs.....: 2427 3.548472/s

EXECUTION
iteration_duration.....: avg=1m13s min=1.46s med=1m14s max=3m1s p(90)=2i
iterations.....: 742 1.084865/s
vus.....: 1 min=1 max=200
vus_max.....: 200 min=200 max=200

NETWORK
data_received.....: 124 MB 181 kB/s
data_sent.....: 496 kB 725 B/s
```

**IP Hash** : peu concluant ici, car les tests sont en local (une seule IP), donc tout est dirigé vers la même instance.

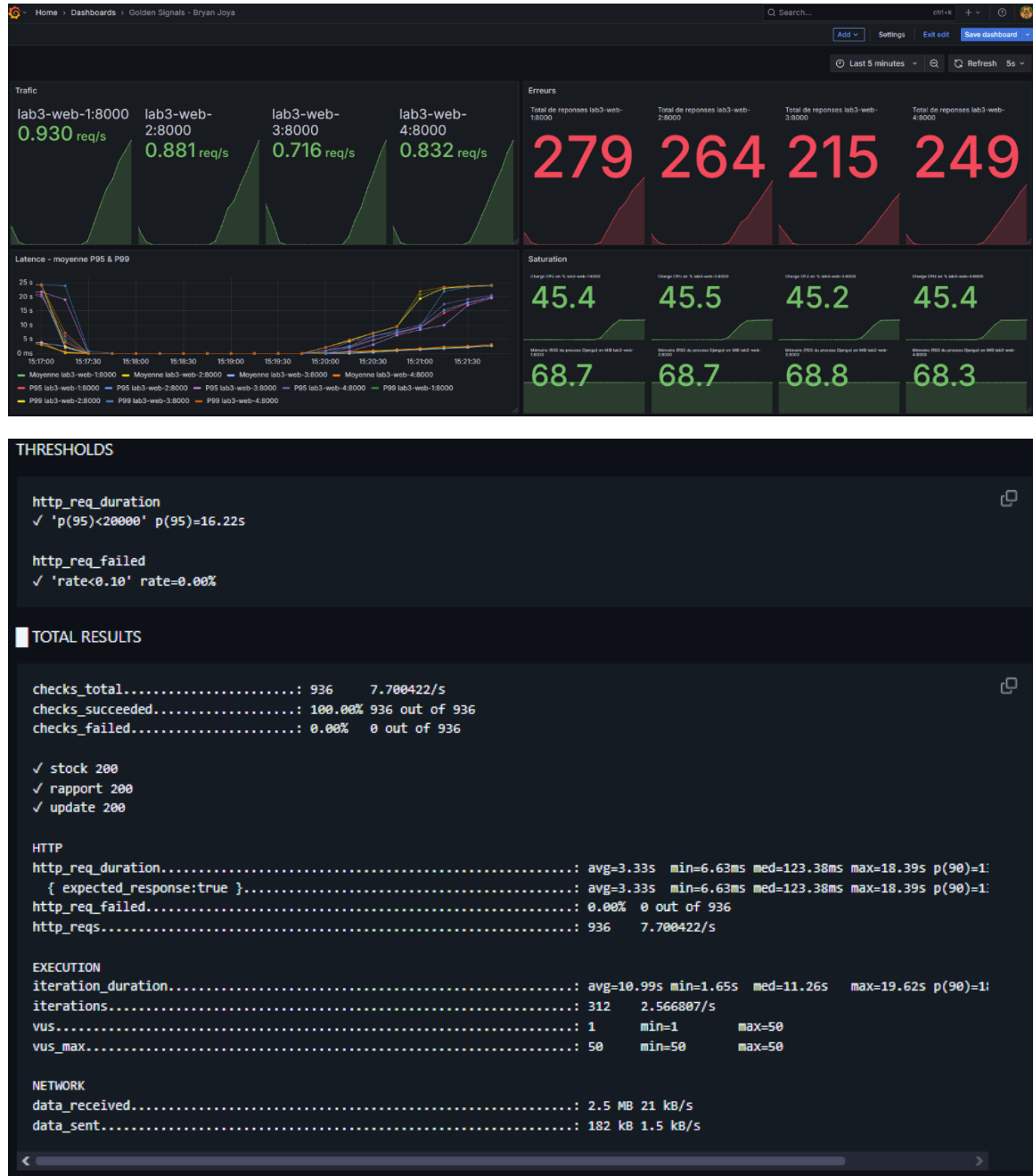
## LEAST CONNECTIONS



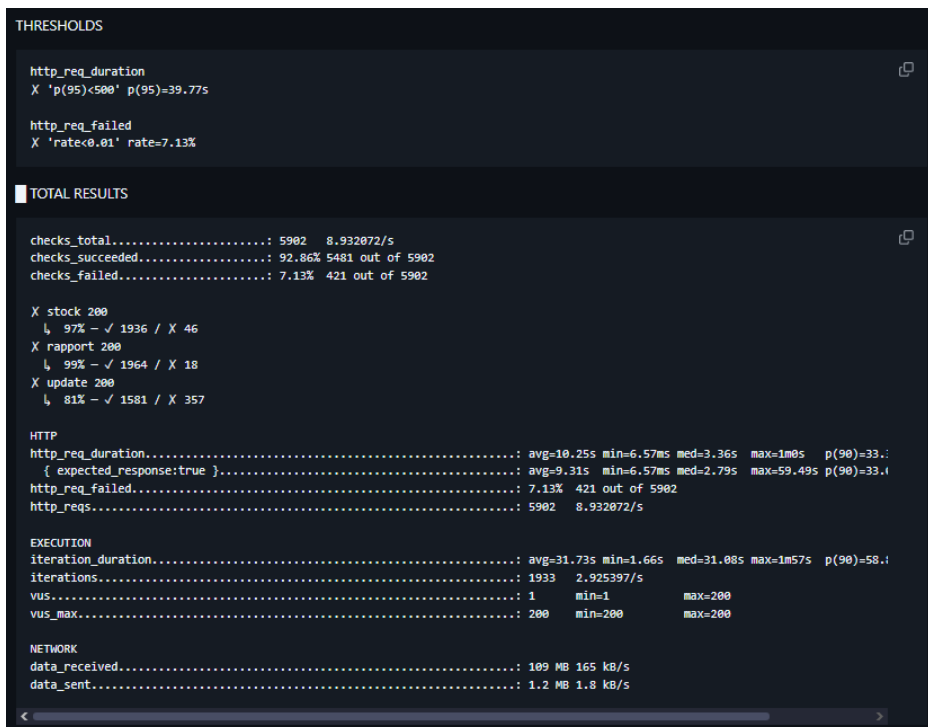
**Least Connections** : la stratégie la plus intelligente dans le cas de requêtes longues. Elle dirige les nouvelles connexions vers les serveurs les moins occupés, améliorant l'efficacité sous forte charge.

## 5. Performances après mise en cache : CONFIGURATION LEAST CONN

LOAD-TEST.JS



## LOAD-TEST-2.JS



Les résultats après avoir fait la mise en cache de STOCK et RAPPORT ne sont pas similaires. On peut assister à une augmentation **drastique** des performances du système :

Environ 30% plus de réussite par endpoint.

Les Checks Total ont plus que doublés (2078 à 5902).

Le nombre de requêtes par secondes à triplé : (3.0 req par secondes à 8.92 req par secondes).

## Représentation du stress test avant et après cache :

Comparons le Load-test-2.js (stress test) avec configuration Least connections

### Graphana:

Métrique	Avant cache	Après cache
Requêtes par secondes	Entre 0.9 et 1.1 par instance	Entre 2.3 et 2.5 par instance
Latence moyenne	Environs 24 secondes	Environs 8 secondes
CPU	25% par instance	43% par instance
Réponses totales traitées	285+326+290+313 = 1214	709 + 697+ 766 + 753 = 2925
Erreurs	33.6% + 41.2% +30.8% + 41.9%	9.6%+9.22%+7.92%+12.2%

### K6

Métrique	Avant cache	Après cache
Checks total	2078 et 3.01rqs/s	5902 et 8.93rqs/s
Checks succeeded	65.49% de 2078 = 1361	92.86% de 5902 = 5481
Iterations_duration	Moyenne de 1min19s	Moyenne de 31 secondes