

# Pruebas de Carga

En este documento detallamos las pruebas a realizar, el análisis de los resultados y la tecnología empleada.

## Estructuración del plan de pruebas

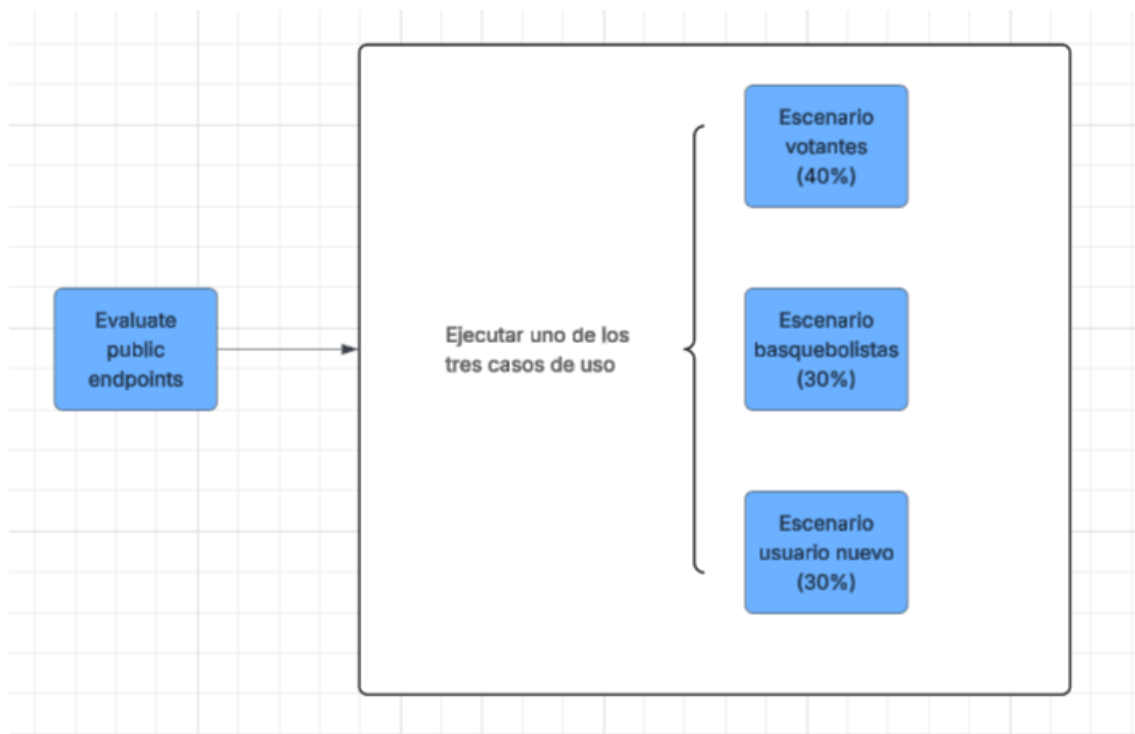
### 1. Evaluación básica a nivel de endpoints

Evaluar de forma aislada los tiempos de respuesta de cada endpoint.

### 2. Escenarios realistas de uso de la aplicación

Se evalúan varios escenarios que representan posibles casos de uso común para los usuarios de la aplicación.

- **Votantes:** un usuario inicia sesión, mira el ranking, ve algunos videos y vota
- **Uso de basquetbolistas:** un usuario inicia sesión, sube un video y va a la lista de video subidos a mirar el estado de sus videos.
- **Usuarios nuevos:** un usuario crea una cuenta, inicia sesión y mira algunos videos



## Carga progresiva

La prueba de carga se configuró con un esquema de incremento progresivo de usuarios virtuales (ramping). Durante seis etapas consecutivas de un minuto cada una, se agregan 100 usuarios por etapa, iniciando con 100 usuarios concurrentes y aumentando linealmente hasta alcanzar un máximo de 600 usuarios al sexto minuto. Este enfoque permite observar cómo responde el sistema bajo un crecimiento sostenido de la demanda, evaluando tanto su capacidad de escalamiento como la estabilidad del rendimiento a medida que aumenta la concurrencia.

## Resultados

```
CUSTOM
auth_response_time.....: avg=7543.502423 min=210 med=7279 max=16630 p(90)=14685.8 p(95)=15517.2
my_videos_response_time.....: avg=7469.833333 min=107 med=7179.5 max=16595 p(90)=14609.9 p(95)=15377
public_videos_response_time.....: avg=7312.603017 min=125 med=7141 max=16442 p(90)=14017 p(95)=15350
rankings_response_time.....: avg=7522.838771 min=99 med=7272 max=16347 p(90)=14671 p(95)=15401.5
scenario_count.....: 1782 4.569061/s
users_created.....: 502 1.287132/s
video_download_response_time....: avg=8655.497161 min=155 med=8467 max=19042 p(90)=16888.8 p(95)=17569.4
video_upload_response_time.....: avg=13656.696154 min=1012 med=11998.5 max=36144 p(90)=24891.8 p(95)=30874.45
video_upload_success_rate.....: 0.00% 0 out of 260
vote_response_time.....: avg=7838.762346 min=325 med=7639 max=16965 p(90)=15600 p(95)=16021.65

HTTP
http_req_duration.....: avg=7.72s min=98.22ms med=7.42s max=36.14s p(90)=14.85s p(95)=15.66s
{ expected_response:true }...: avg=7.47s min=98.22ms med=7.22s max=16.96s p(90)=14.62s p(95)=15.38s
http_req_failed.....: 14.70% 2244 out of 15255
http_reqs.....: 15255 39.113929/s

EXECUTION
iteration_duration.....: avg=54.18s min=3.38s med=53.78s max=1m58s p(90)=1m36s p(95)=1m43s
iterations.....: 1627 4.17164/s
vus.....: 494 min=1 max=599
vus_max.....: 600 min=600 max=600

NETWORK
data_received.....: 95 MB 243 kB/s
data_sent.....: 2.6 MB 6.7 kB/s

running (6m30.0s), 000/600 VUs, 1627 complete and 494 interrupted iterations
stress ✓ [=====] 493/600 VUs 6m0s
```

## Análisis de los Resultados

Durante la evaluación básica a nivel de endpoints se identificó que la mayoría de las operaciones críticas presentaron tiempos de respuesta elevados bajo condiciones de carga progresiva. Los endpoints de autenticación (/api/auth/login), videos públicos (/api/public/videos), rankings (/api/public/rankings) y videos de usuario (/api/videos) tuvieron tiempos de respuesta promedio entre 7 y 8 segundos, con valores en el percentil 95

superiores a los 15 segundos. Esto evidencia que, conforme aumentó la concurrencia, el sistema no pudo mantener latencias aceptables, lo que impacta directamente la experiencia del usuario.

En los escenarios realistas de uso de la aplicación, se observaron los siguientes comportamientos:

- **Votantes:** los usuarios lograron iniciar sesión y consultar rankings y videos, pero con latencias superiores a lo esperado ( $\approx 7-8$  segundos en promedio), lo cual afecta la fluidez de la interacción. Aunque el proceso de votación se ejecutó, en muchos casos la operación presentó demoras notables.
- **Basquetbolistas:** este escenario fue el más crítico. Las cargas de video tuvo tiempos promedio de 13.6 segundos y un pico máximo de 36 segundos. Este resultado evidencia que el sistema presenta debilidades para soportar operaciones pesadas de subida de archivos bajo alta concurrencia con la instancia única.
- **Usuarios nuevos:** los procesos de creación de cuenta y exploración de videos sí se realizaron, pero nuevamente con tiempos de respuesta muy altos, comparables a los del resto de endpoints (7-8 segundos en promedio).

En términos globales, la prueba generó 15,255 solicitudes HTTP, de las cuales un 14.7% fallaron (2,244 errores). Asimismo, se registraron 494 iteraciones interrumpidas, lo cual indica que múltiples usuarios virtuales no pudieron completar su flujo debido a fallas o tiempos de espera excesivos.

El sistema alcanzó su punto de estrés a partir de  $\sim 493$  usuarios concurrentes, momento en el cual la degradación del servicio se hizo evidente: los tiempos de respuesta aumentaron abruptamente, las tasas de error crecieron y varios escenarios dejaron de completarse correctamente. Esto significa que la plataforma es capaz de sostener un volumen de usuarios concurrentes cercano a los 500, pero colapsa al acercarse al objetivo de 600.

En conclusión, los resultados muestran que el sistema no logró sostener la carga de 600 usuarios concurrentes definida en la prueba. Las altas latencias y el porcentaje significativo de reflejan limitaciones tanto en el backend como posiblemente en la infraestructura de soporte. Antes de escalar a escenarios mayor uso, será necesario realizar optimizaciones en el procesamiento de peticiones, mejorar la eficiencia de operaciones críticas como la carga de videos y evaluar la capacidad de los recursos de infraestructura disponibles.

## Tecnología Empleada

Para la ejecución de las pruebas se utilizó k6, que es una herramienta de pruebas de carga y estrés orientada a desarrolladores, diseñada para evaluar el rendimiento y la resiliencia de aplicaciones y servicios web. Permite definir escenarios de prueba mediante scripts en JavaScript, lo que facilita la creación de flujos realistas que simulan el comportamiento de los usuarios finales. Una de sus principales ventajas es su bajo consumo de recursos y su capacidad para ejecutar pruebas de alta concurrencia de manera eficiente, ya sea en entornos locales o en la nube. Además, proporciona métricas detalladas como latencias, tasas de error y percentiles de tiempo de respuesta, lo que la convierte en una opción sólida para detectar cuellos de botella y validar la capacidad de la infraestructura antes de enfrentar escenarios de uso real.