



상위 1% 월급쟁이를 위한

쿠버네티스 어나더 클래스

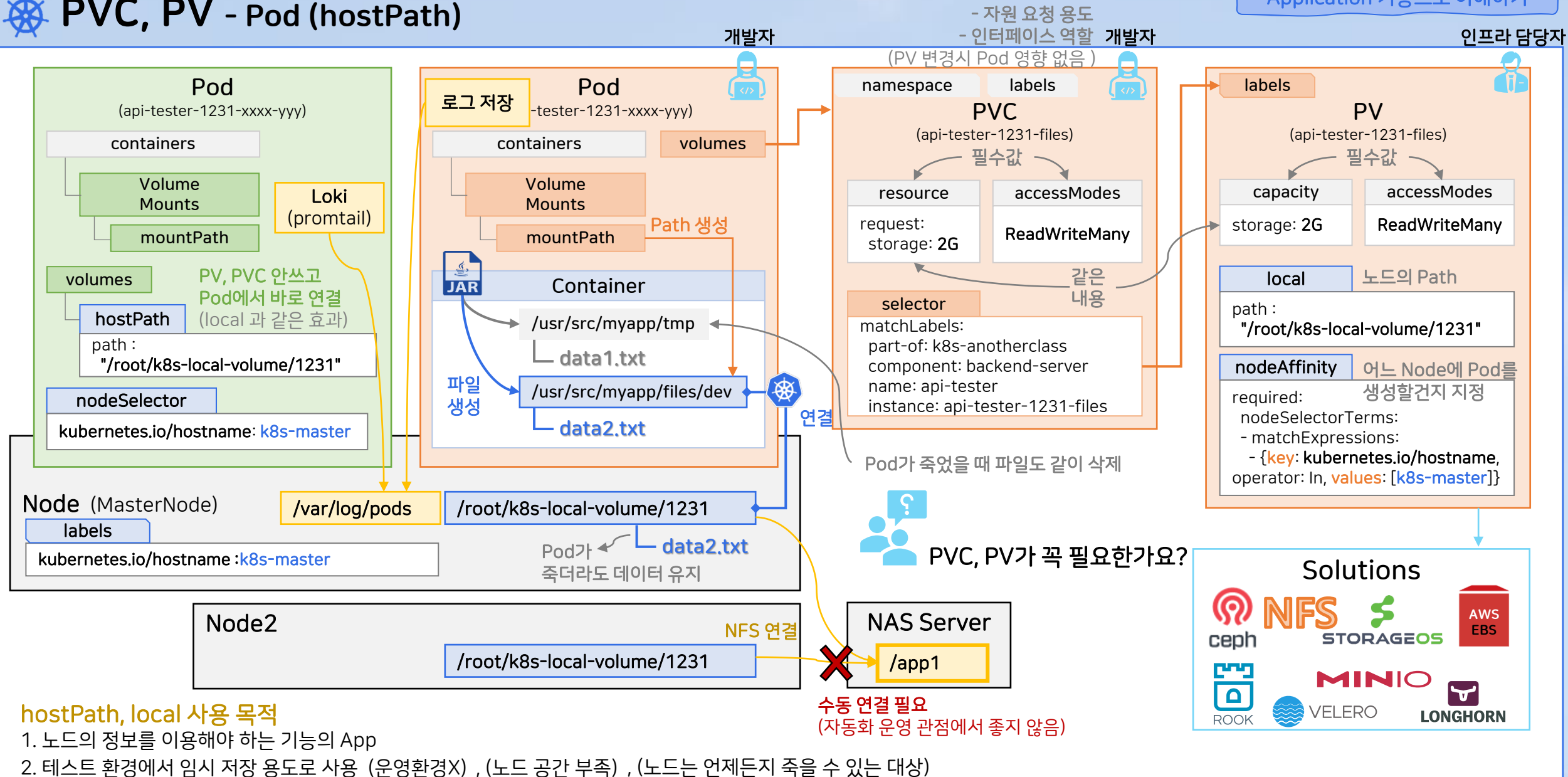
Application 기능을 이해하기

PVC, PV / Deployment / HPA / Service



PVC, PV - Pod (hostPath)

Application 기능으로 이해하기

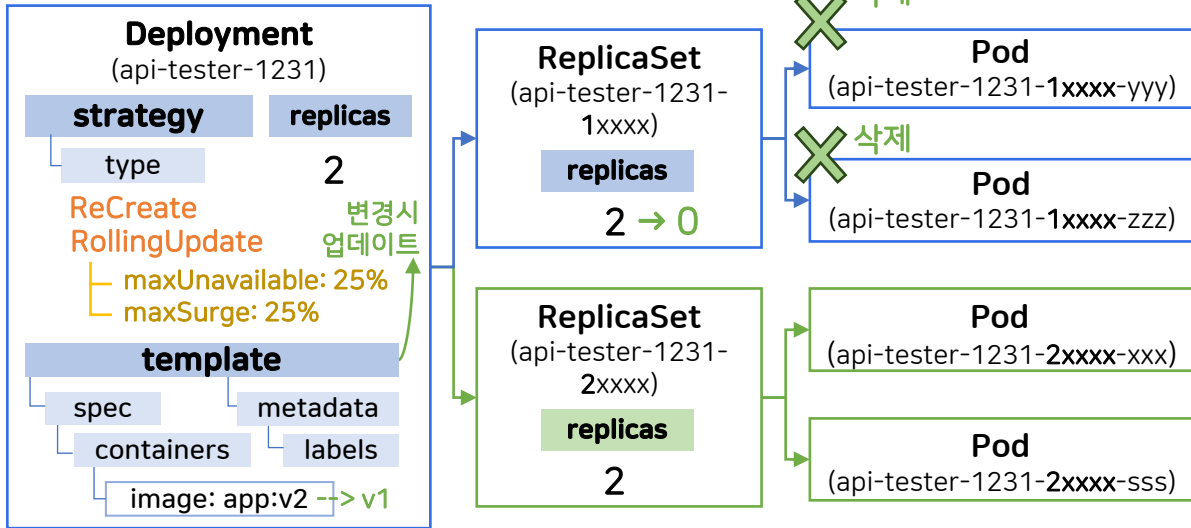




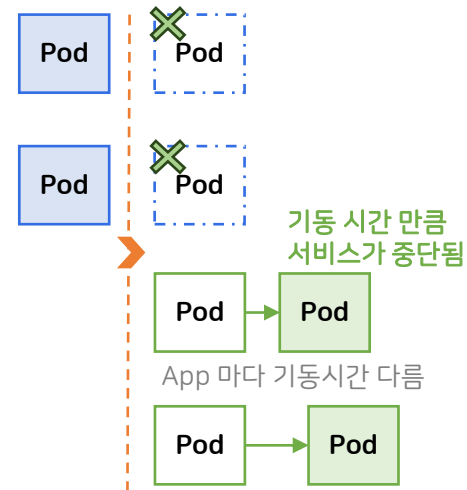
Deployment - update

Application 기능으로 이해하기

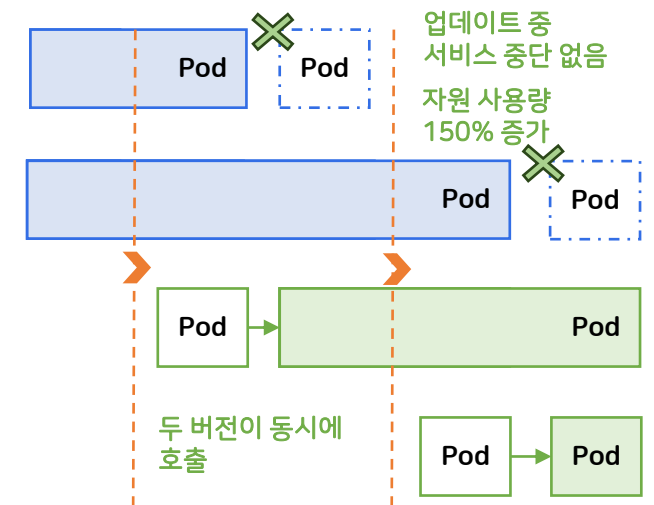
Deployment Update



Recreate



RollingUpdate



RollingUpdate - maxUnavailable, maxSurge

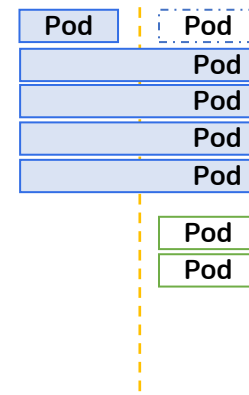
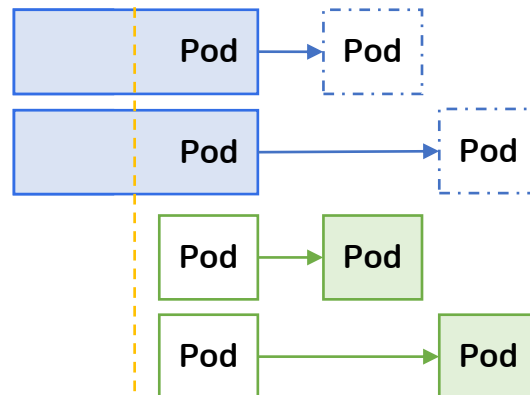
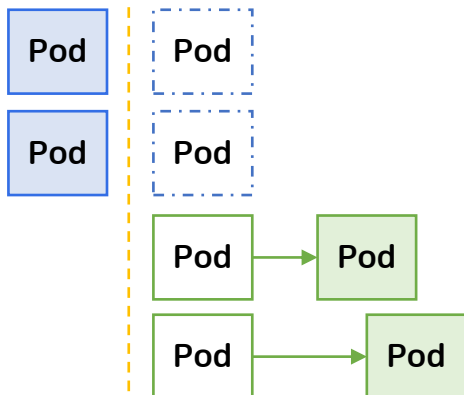
maxUnavailable: 100%
maxSurge: 100%

Recreate와 같은 효과

maxUnavailable: 0%
maxSurge: 100%

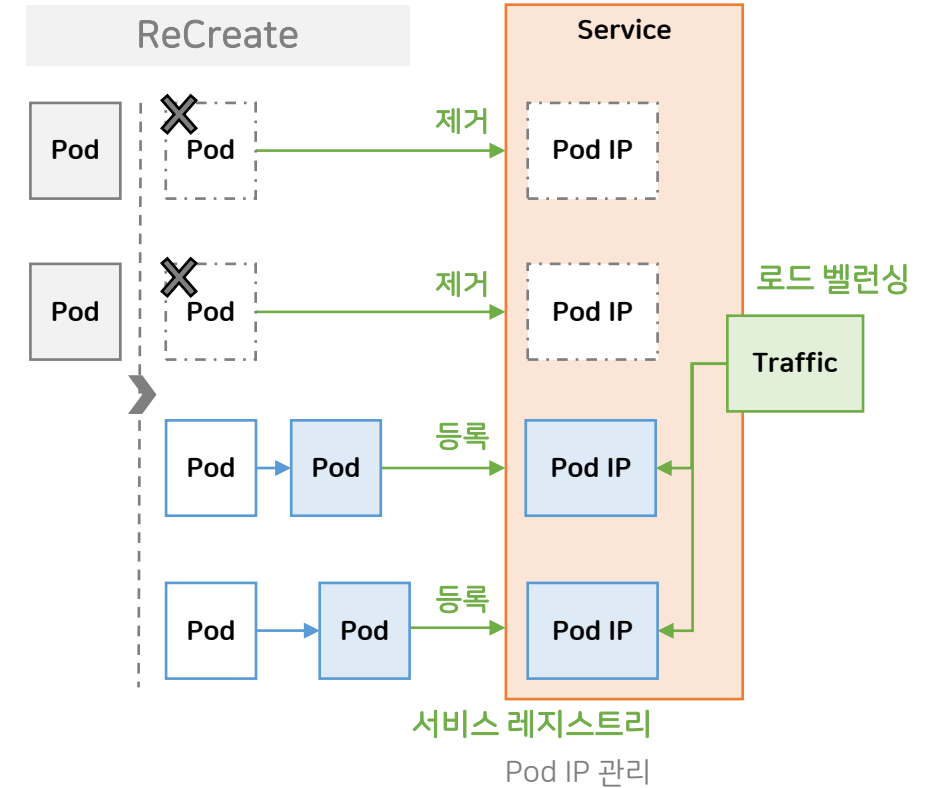
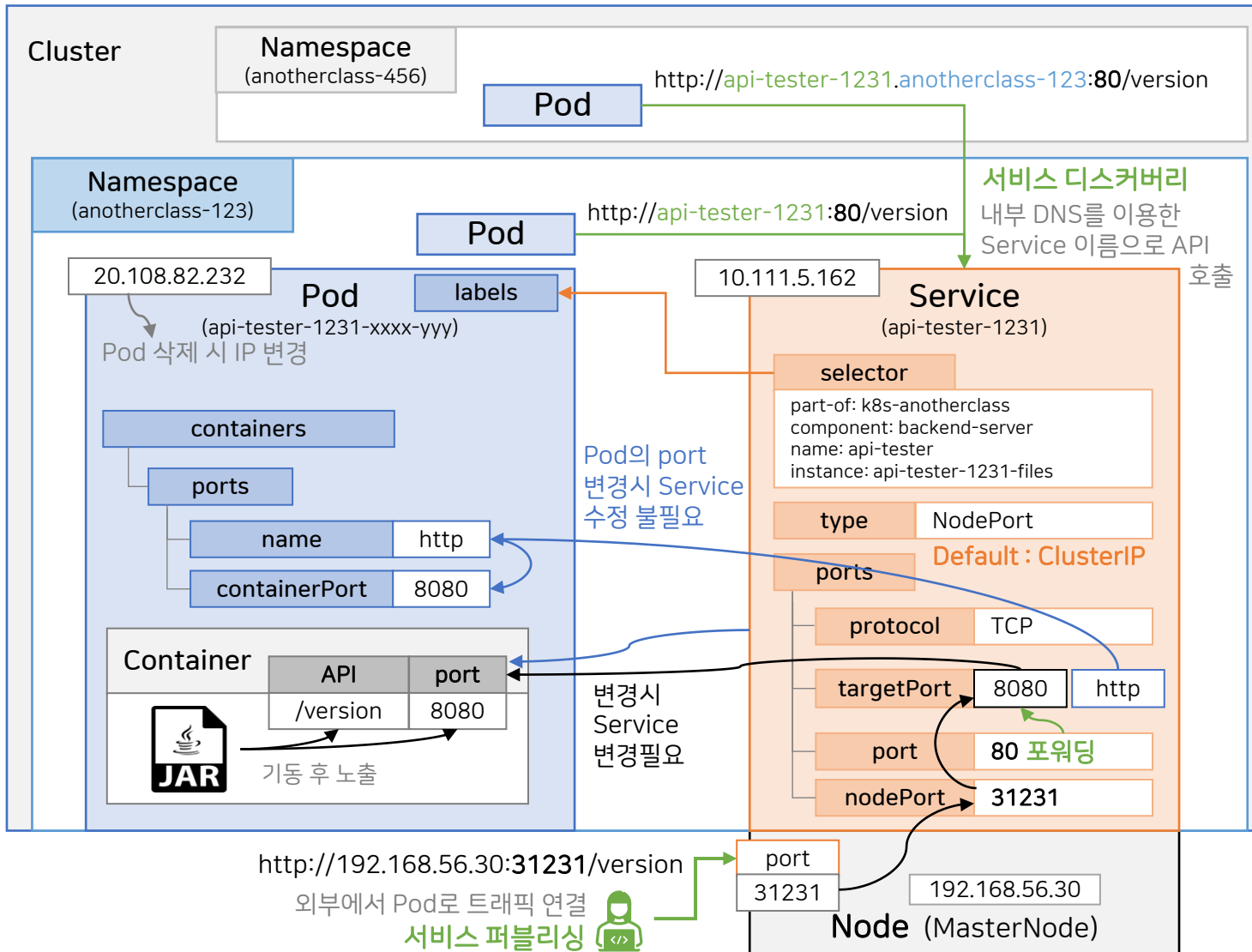
Blue/Green에 가까운 효과

maxUnavailable: 25%
maxSurge: 25%

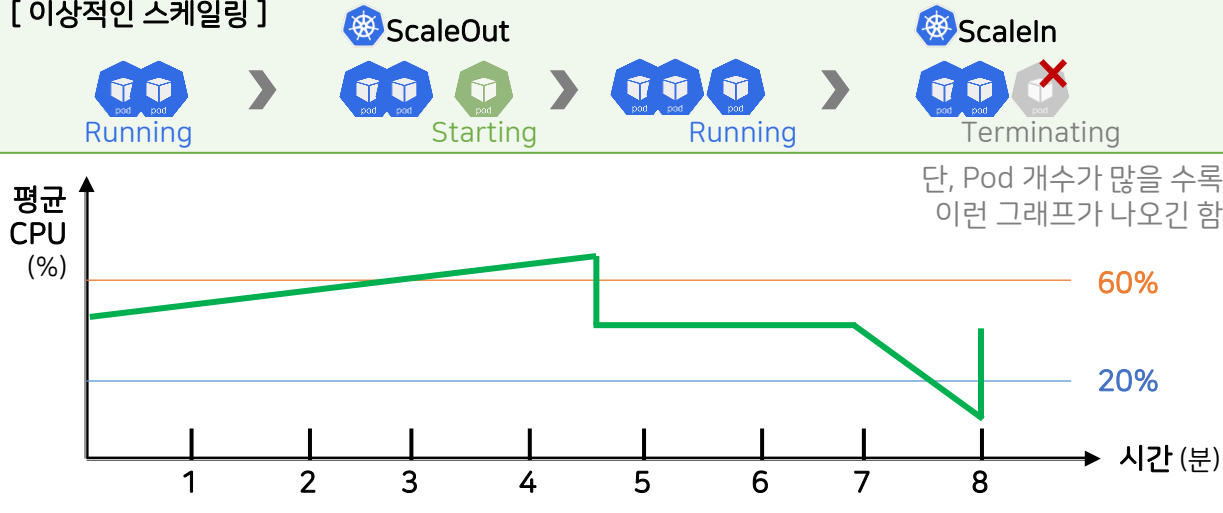


[Blue/Green]

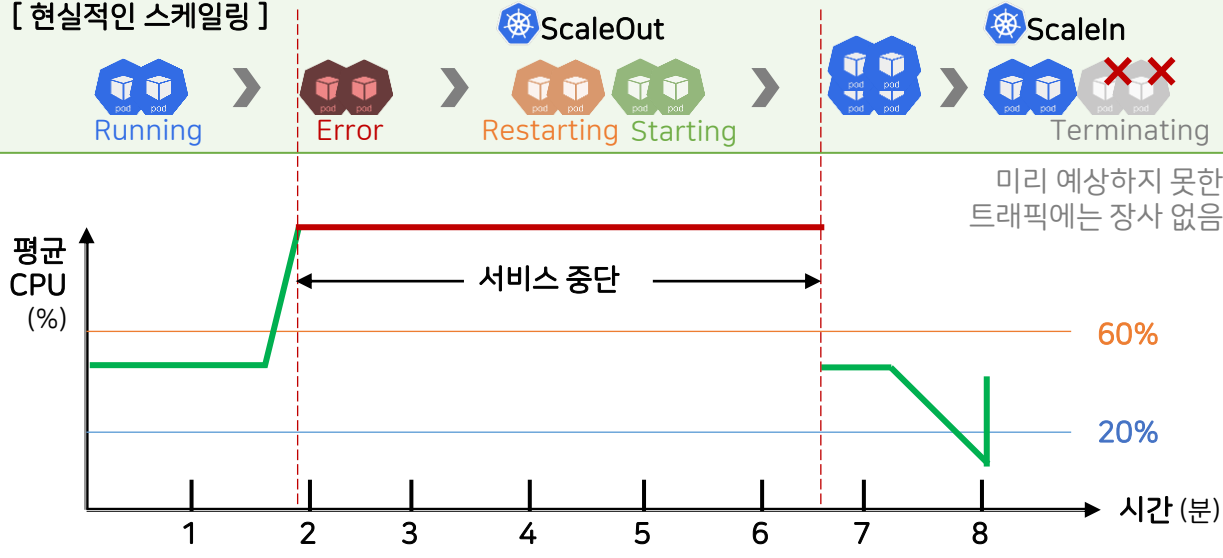
- 업데이트 동안 두 버전이 동시에 호출되지 않음
- 자원 사용량 200% 증가
- 별도 배포 솔루션에서 제공



[이상적인 스케일링]



[현실적인 스케일링]



HPA (api-tester-1231-default)

```

scaleTargetRef
  apiVersion: apps/v1
  kind: Deployment
  name: api-tester-1231
  minReplicas: 2
  maxReplicas: 4
  metrics
    - type: Resource
      resource:
        name: cpu
        target:
          type: Utilization
          averageUtilization: 60
  
```

수치를 결정하는데 영향도가 큼

현재 Pod 수 * (평균 CPU / HPA CPU) → 변경될 Pod 수

$2 * ((150+1)/2) / 60 = 3 (2.51)$

behavior

scaleUp: stabilizationWindowSeconds: 120

scaleDown: stabilizationWindowSeconds: 600

polices:

- type: Pods
- value: 1
- periodSeconds: 60

1분에 1개씩 Pod 제거

Deployment (api-tester-1231)

```

replicas: 2
  
```

ReplicaSet

Pod -Xmx 50M

resource	container
limits:	200%
cpu: 200m	150m
memory: 200Mi	70m
requests:	40Mi
cpu: 100m	
memory: 100Mi	
100% 계산을 위한 기준 값	
CPU	Memory

부하에 따른 증감 대상 아님

2분 동안 60%이상 유지시 scaleOut

부하가 감소해도 10분동안은 Pod 수 개수