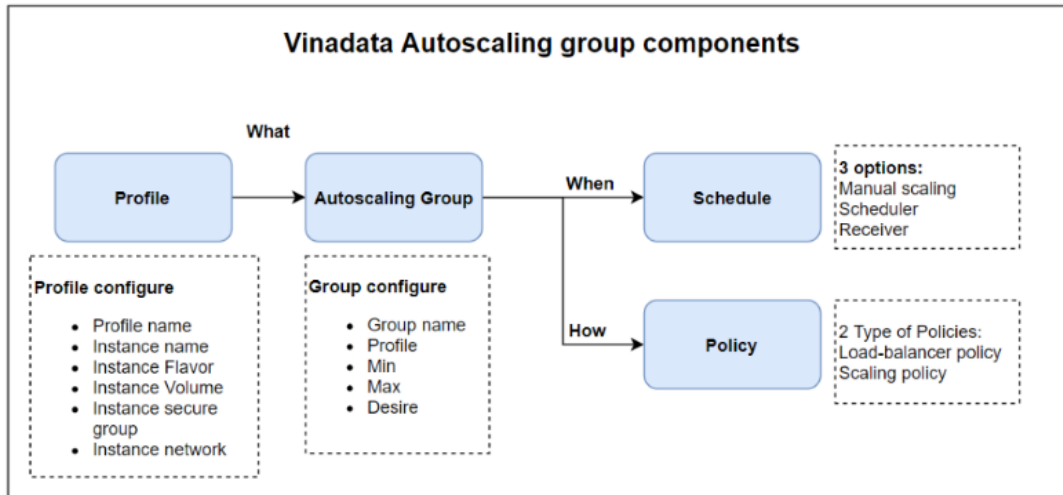


Dch V Auto Scaling HCM 02

Auto-scaling là dịch vụ tự động tăng giảm (scale) máy chủ (VM) theo yêu cầu do user tạo ra. Auto-scaling mở rộng tính sẵn sàng cao của hệ thống mà vẫn tiết kiệm chi phí

Các thành phần của Autoscaling



Profile: Là thiết lập cấu hình mẫu của instance bao gồm Flavor, Volume, Secure group, network. Auto scaling group dựa trên thiết lập của profile để scale Instance có cấu hình tương ứng.

Một profile có thể gắn nhiều scaling group, tuy nhiên mỗi scaling group chỉ có thể gắn với 1 profile duy nhất.

Cấu hình mẫu của profile sau khi tạo thành công không thể edit. Nếu có nhu cầu thay đổi thì bạn cần tạo profile mới.

Auto Scaling group: Một auto scaling group sẽ gồm nhiều instances có cùng một loại cấu hình như nhau từ profile. Vì vậy mỗi auto scaling group chỉ có thể gắn với 1 profile và không thể có nhiều hơn 1 profile cho mỗi group.

Auto scaling group sẽ tự động tăng giảm vì thiết lập "desired capacity" và duy trì số lượng này cho dù có 1 instance bị lỗi (un-healthy). Các tính năng "healthy check" & "Recover node" sẽ kiểm tra, terminate các instance bị lỗi thì tạo instance mới, giúp cho số lượng instance luôn mở rộng.

Schedule: Component này sẽ xác định thời điểm nào auto-scale sẽ thực hiện, Có 3 cách thiết lập là:

- Manual: Bằng cách thay đổi số "desired capacity" từ giao diện cấu hình Auto scaling group, số lượng instance sẽ tăng thêm / giảm xuống tùy theo nhu cầu.
- Scheduler: thiết lập lịch cụ thể tùy theo nhu cầu. Ví dụ thiết lập lịch tăng giảm instance khi giá cao hơn và thiết lập lịch giảm số instance xuống giá thấp hơn để tiết kiệm chi phí.
- Receiver: Cho phép bạn scale (in/out) một auto-scaling group thông qua http url. Webhook receiver thường sử dụng khi bạn có sẵn một monitor tool.

Policy: Khi gắn vào 1 auto-scaling group, policy sẽ quyết định cách thức scale của group. Policy gồm 2 loại là:

- Load balancer policy: Policy này giúp các instance thoát ra khỏi auto-scaling để thêm vào pool của Load balancer.
- Scaling policy: Policy này sẽ quyết định các tính chất gồm: Scaling out (Tăng thêm), scaling in (giảm xuống), mức tăng / giảm bao nhiêu instance, cooldown giữa các lần scale

Một policy có thể gắn vào nhiều auto-scaling group khác nhau (trừ load balancer policy), giúp bạn tiết kiệm thời gian khi có nhiều auto scaling group. Một auto-scaling group chỉ có thể gắn nhiều hơn 1 policy, tuy nhiên không thể gắn cùng một loại policy (ví dụ: một group chỉ có 1 scaling out hoặc 1 scaling in policy).