

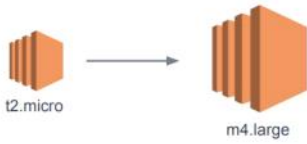
EC2 Auto Scaling



Auto Scaling

What is Scaling?

Vertical Scaling



Scale Up / Down

Horizontal Scaling



Scale Out / In

Scaling konusu bizim her yerde karşımıza çıkabilecek bir konu. Sadece AWS e has bir şey değil. O yüzden bu konuyu iyi anlamak lazım.

2 çeşit ölçeklendirme vardır: Vertical ve Horizontal scaling

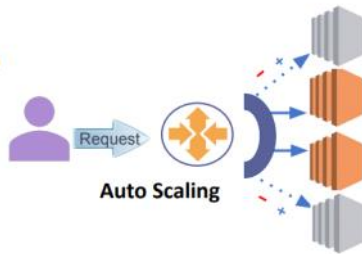
Vertical scalingte bir bina var mesela üzerine kat çıkmak gibi. Dikey bir büyüme var. T2.micro az geldi mesela diyelim ki o zaman m4.large ile büyütüyoruz. Kapladığı alan aynı ama özelliği farklı mesela memoryi arttırıyoruz ya da cpusunu arttırıyoruz gibi. Literatürde buna scale up veya scale down diyoruz. Burada quality/nitelik arttırıyoruz.

Horizontalda ise bir instance yetmediyse mesela onun kapasitesini arttırmak yerine sayısını arttırıyoruz. 1 instance değilse mesela 3 instance run ediyoruz gibi. Burada nicelik arttırıyoruz. Bunu da literatürde scale out, scale in diye kullanıyoruz.

Aws çoğunlukla horizontal scaling yapıyor ama launch template ettiğimizde hem horizontal hem vertical scaling yapabiliyoruz.

Auto Scaling

What is Auto Scaling?

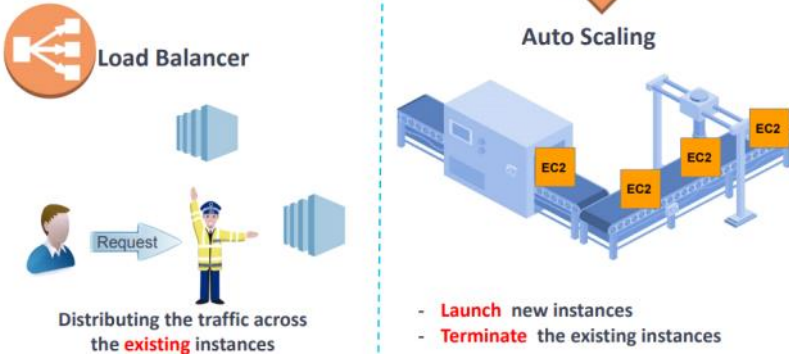


Auto scaling ise bizim bir yükümüz var diyelim. Bu yüke yetecek kadar isin içine instance katmamız demek scaling yapmak demek.

- Amazon EC2 Auto Scaling is a component that helps you ensure that you have the **correct number of Amazon EC2 instances** available to **handle the load** for your application.
- Auto Scaling **adds or removes instances** to keep your system steady state.
- You can **automate the increasing or decreasing** of virtual machines depending on your policy.

Auto Scaling

Auto Scaling vs Load Balancer



Load balancer yoktan var etme gucu yok. Elinde kac instance varsa onu yoneter. Mesela 3 instanceta birisi bozuldu diyelim diger ikisine gonderir. Ikisi bozulduysa yuku kalan instancea gonderir ama o da bozulursa shuttle down yapar ve tum yuku hepsine gonderir hangisi neyi alirsa artik mantigiyla. Ama auto scalingte ihtiyac hasil oldu diyelim hemen verilen algoritmalar dahilinde instance ayaga kildirir ya da instance azaltir. Auto scaling target grouplara ihtiyac olan instancelari sagliyor. Load balancerda bu instancelara dengeli yuk dagilimlarini yapiyor. Auto scaling uretme ve terminate etmek yetkisi var. Load balancer ise bu auto scalingin urettigi veya terminate ettigi instancelar arasinda yuk dagilimini sagliyor.

Auto Scaling

Features of Auto Scaling

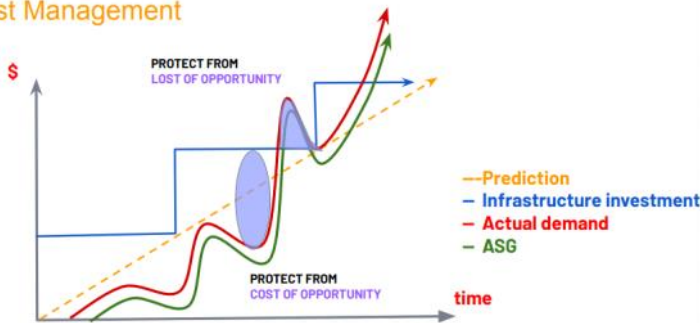


- Auto Scaling Policy
- Launch Configurations or Launch Templates.
- Fault tolerance.
- Compatible with Elastic Load Balancer
- Better Cost Management.

Bu yonetme sekli policylerle oluyor. Bizim belirledigimiz policylere gore yani bizim belirledigimiz kriterlere gore auto scaling uretime geciyor belirli zamanlarda ya da terminate ediyor instancelari kullanim azligina gore. Ama bunlarin hepsini policylerde biz belirliyoruz cunku gorunmez bi makine var ortada onu yonlendirmek icin bir kurallar butunune ihtiyacimiz var. Auto scalingin bir instance ayaga kaldirabilmesi icin launch configuration ve launch template ihtiyaci var ki hangi kistaslara gore bir instance kaldiracagini bilsin- Her zaman instance dengesini tutturmamabiliyoruz yani bazen cok bazen eksik instance ayaga kaldirtabiliyoruz auto scalingde. Bunun dengersiz oldugu durumlarda zararlar olabilir. ELB ile isbirlikci calisir.

Auto Scaling

Better Cost Management

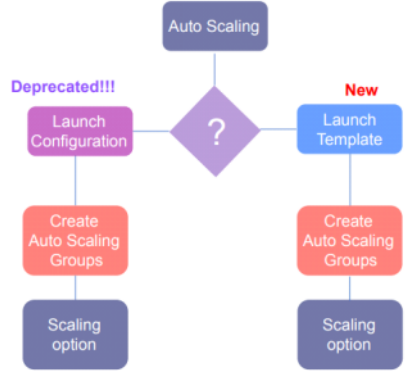


Bu bizim tahmini yatırım ve zaman çizelgemiz. Sarı olan çizgi bizim tahminimiz. Mavi olan kısımda dikeye gectigi zamanlar instance ihtiyacinin oldugu zamanlar. Yukari dogru gittigi vakit yeni instancelar eklemis demek. Ongorulebilir belirli araliklarla infrastructura yatırım yapmis oluyoruz. Ama bazen bizim yaptigimiz yatirimla actual demand yani gercek ihtiyaclar uyusmaya biliyo kirmizi çizgideki gibi. Sarı bizim tahmini egri. Mavi bizim yatirimimiz. Kirmizi ise gerceklesen egri. Yesil olan bizim auto scaling groupumuz. Bizim mavi çizgimizden ziyade kendi iyaptigi ongoruye gore gercege en yakin olan instance araligini belirleyip instance uretiyor ya da terminate ediyor. Yani yeterli yatırım yapamama ya da fazla gereksiz yatırım yapma olayini cozmus oluyor talep edilmeye veya edilmeme durumu dahilinde ve policy ve kriterlere gore.

Auto Scaling

Auto Scaling Creating Process

- First, you need to select either the **Launch Template** or the **Launch Configuration** option and create it.
- Then, create an **Auto Scaling Group**.
- Finally, Finish Creating Auto Scaling



Auto scaling yaratmanın iki yolu var: Launch configuration (deprecated olmuş bu) ve Launch template

Bir launch template yaratıyoruz instance yaratacağı zaman neye göre yaratacağını bilsin diye. Daha sonra bir auto scaling grubu oluşturun ve bittikten sonra scaling option ekliyoruz.

Auto Scaling

Launch Configuration Deprecation Schedule

The announcement from AWS

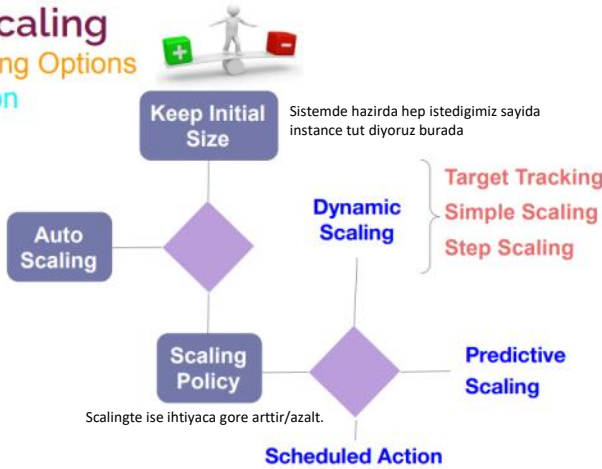
- March 31, 2023 - New accounts created after this date will not be able to create new launch configurations via the console. **API and CLI access will remain available to support customers** with automation use cases.
- December 31, 2023 - New accounts created after this date **will not be able to create new launch configurations**.

Launch configuration artık tedavülden kalkmış.

Auto Scaling

Auto Scaling Options

Next session



Sistemde hazırda hep istediğimiz sayıda instance tutuyoruz burada

Scalingte ise ihtiyaca göre arttır/azalt.

Auto scalingi iki şekilde kullanacağız yani iki turunu da kullanacağız: keep initial size ve scaling policy. Mesela biz auto scalingten her zaman 3 sağlıklı instance oluşturmamızı isteriz. Çünkü birisinin başına bir is gelmesi dahilinde load balancer sadece iki instance arasında yük dağıtır yük aynı olmasına rağmen. O yüzden auto scalingin her zaman hazırda üç instance tutmasını istiyoruz diyelim ki. O hemen bir sağlıklı instance daha üretir 3 istediğimiz için biz hep en az.

Scaling policyde ise şunu diyoruz:trafik arttığında arttır azaldığında instanceları azalt. Burada da 3 şekilde scaling yapılıyor: scheduled, predictive, dynamic.

Scheduledta zaman belli. Mesela her gün saat 03.00 da benim ekstra instance ihtiyacım var sen arttır. Ve akşam 5te de kaldır instanceları ihtiyac yok diyorum mesela. Benim ayarladığım saate göre ben ongoruyorum ve bildiriyorum.

Predictivede ise kendi el atıyo mevzuyla. Diyoki mben sistemi en az bi 48 saat inceleyeyim, bakayım bir hangi saatlerde ihtiyac fazla ya da az. Ona göre kendi kendine o aralıklarda ec2 arttırıp azaltmaya başlıyor. Bir analiz var burada.

Dynamicte ise anlık olarak response var. Anlık olarak ihtiyac arttı mesela anında ec2 kaldırıyor ya da ihtiyac bitti diyelim hemen terminate ediyor. Burada bir netlik yok. Urunun ne kadar satılacağı bilinmiyordur mesela bu kullanılır. Burada sürekli bir gözlem var. Hep dinamik. Burada da 3 mesele var: Target trackingte mesela diyoruz ki bir pcpu belirliyoruz mesela. Bunu hep %40in altında tutuyoruz. O da ona göre anlık reaksiyon gösteriyor.yani %30-35 olduğunda aws harekete geçiyo hemen aktif bir şekilde gözlem yapıyor ve reaksiyon gösteriyor.

Simple scalingte, eğer cpu %40i geçerse harekete geç diyoruz yani %40 olana kadar hareket etme bi dur diyoruz. Targetta gecmeden önce reaksiyon basılıyordu. Burada ise gectikten sonra basılıyor.

Stepte ise simplein bir üst hali. Yanimesela diyoruz ki %60-80 arasında 1 instance ayaga kaldırı ama %80 geçerse 3 instance ya da %90 geçerse ne varsa kullan diyoruz mesela bu şekilde steplendiriyoruz.

SINAVLARDA SIMPLE VE TARGETİ AYNI ANDA VERMEZLER AMA AWS TARGET TAVSİYE EDER. ORAN VS VARSAN TARGETTİR

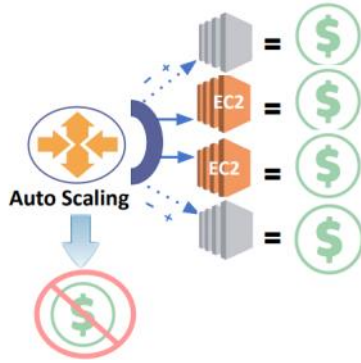
SPEŞİFİK BİR SAAT İSE SCHEDULED SECİYORUZ.

ORAN VERİRSE TARGET, SAAT VERİRSE SCHEDULED

CLARUSWAY

► Auto Scaling

Pricing for Amazon EC2 Auto Scaling



CLARUSWAY
WAY TO REINVENT YOURSELF



Auto scaling ücretsiz. Arkada bunu yapan makineyi biz görmüyoruz.

Ama ürettiği nihai çıktılarından ücret alıyor.

Yani auto scaling gece ucte arttirip azaltmaktan değil, kullandigi kadar instance parasi alıyor.

Load balancer ve cloudwatch icinde para aliyor belli bir miktardan sonra.

Cloudwatch auto scalingin onemli bir parcasi. Buradan alarm kuruyor arttirmasi gerektiği zamanları anlamasi için.