
使用随机分区进行工作负载隔离

Colm MacCárthaigh



使用随机分区进行工作负载隔离

Copyright © 2019 Amazon Web Services, Inc. and/or its affiliates. All rights reserved.

如今，Amazon Route 53 上托管了许多全球最大的企业和最受欢迎的网站，但在上线之初可没有如此名气。

采用 DNS 托管

AWS 开始提供服务后不久，AWS 客户就明确表示他们希望能够在自己“根”域上使用 Amazon Simple Storage Service (S3)、Amazon CloudFront 和 Elastic Load Balancing 服务。这里的“根”域指的是诸如“amazon.com”之类的域名，而不仅仅是“www.amazon.com”之类的域名。

这听上去似乎很简单。但是，由于 DNS 协议的设计方案要追溯到 20 世纪 80 年代，因此这项任务比看起来要难。DNS 具有一项称为 CNAME 的功能，该功能允许域的拥有者将域的一部分卸载到其他提供商进行托管，但不能在根域或者顶级域上使用。为了满足客户的需求，我们必须实际托管客户的域。当托管客户的域时，我们可以返回 Amazon S3、Amazon CloudFront 或者 Elastic Load Balancing 的当前 IP 地址集。但是这些服务都在不断扩展并添加 IP 地址，因此客户也无法轻松地自己的域配置中进行硬编码。

托管 DNS 并非易事。如果 DNS 存在问题，那么整个企业都可能会连接中断。在确定了需求之后，我们便立即着手按照 Amazon 的典型工作方式来解决问题。我们组建了一支小型工程师团队，然后开始工作。

应对 DDoS 攻击

对于任何 DNS 提供商来说，如何应对分布式拒绝服务 (DDoS) 攻击是大家公认的最大的挑战。DNS 建立在 UDP 协议的基础上，这意味着在许多缺乏监管的互联网中充斥着大量欺诈性的 DNS 请求。但 DNS 同时也是至关重要的基础设施，这样一来，DNS 便成为了易受攻击的目标，吸引了许多想要勒索企业的不法分子、出于各种原因企图制造可引发停机的恶意“引导程序”，以及一些偶然误入歧途的麻烦制造者，这些人似乎没有意识到他们的所作所为是需要承担法律责任的严重犯罪行为。无论出于什么原因，每天都会发生数千次针对域的 DDOS 攻击。

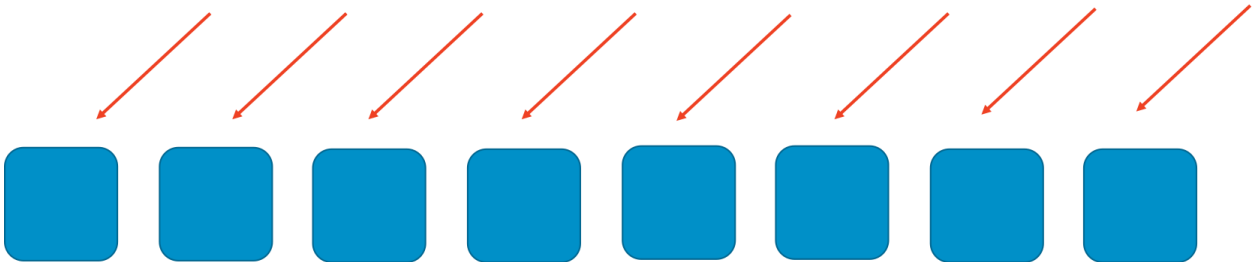
减少这类攻击的一种方法是使用高服务器容量。尽管拥有良好的容量基准很重要，但是这种方法并不能真正实现扩展。提供商新增一台服务器需要花费数千美元，而如果攻击者使用感染性极强的僵尸网络，就可以极少的成本制造出大量假客户端。对于提供商而言，提高服务器容量并不可行。

在我们构建 Amazon Route 53 时，DNS 防御的最新技术是采用专用网络设备，这些设备可以通过各种方式对流量进行高速“清理”。我们在 Amazon 上将许多此类设备用于我们现有的内部 DNS 服务，另外我们还与硬件供应商讨论了其他可用工具。我们发现要完全覆盖每个 Route 53 域所需购买的设备需要花费数千万美元，并且设备的交付、安装和投入运行时间将会使我们的原有计划延长数月。这与我们计划要求的快速、节省的原则背道而驰，因此我们从未认真考虑过使用这些设备。我们需要找到一种方法，将所花费的资源全部用于保护实际遭受攻击的域。我们从“需求是发明之母”这句古老原则受到了启发。我们的需求是使用适量的资源快速构建世界一流的 100% 正常运行时间 DNS 服务。由此我们发明了随机分区。

什么是随机分区？

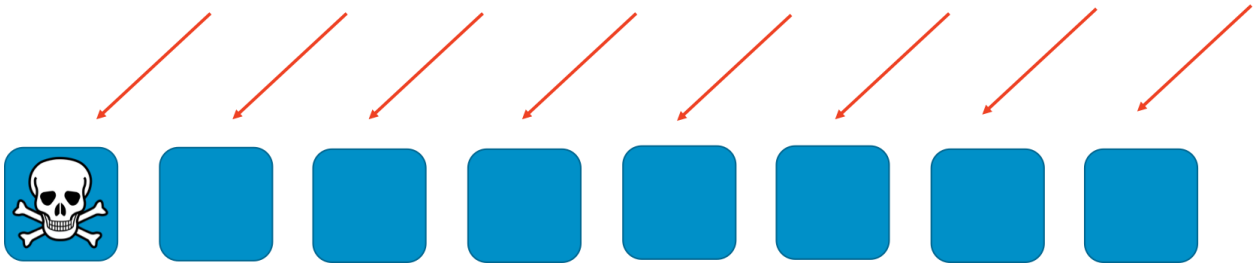
随机分区简单而不失强大功能。它的强大远超我们最初的设想。我们反复地使用随机分区，如今它已成为一种核心模式，使 AWS 可以提供具有高性价比的多租户服务，为每个客户提供单租户体验。

要了解随机分区的工作原理，首先请思考如何通过普通分区提升系统的可扩展性和弹性。想象一下一个由八个工作单元组成的可水平扩展的系统或服务。下图展示了工作单元及其要求。工作单元可以是服务器、队列或数据库，换言之，可以是任意构成系统的“基本单位”。



未分区时，一组工作单元负责处理所有工作。每个工作单元必须能够处理任何请求。这对于提高效率 and 冗余非常有效。如果一个工作单元发生故障，其他七个工作单元可以接手其工作，因此系统中所需的备用容量相对较小。但是，如果通过特定类型的请求或大量请求（例如 DDoS 攻击）触发工作单元故障，就会出现大问题。以下两个图片展示了这种攻击的发展进程。

使用随机分区进行工作负载隔离

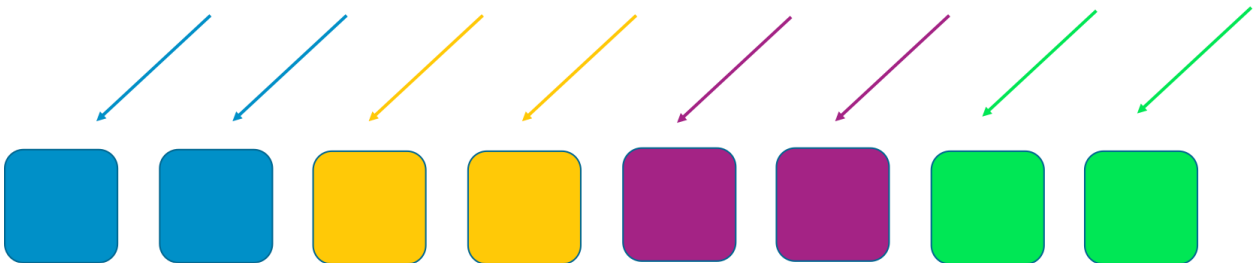


这种问题首先会影响第一个工作单元并使其发生故障，并且随着其余工作单元的接管，还会继续逐一影响其他工作单元。很快，所有工作单元以及整个服务就会瘫痪。



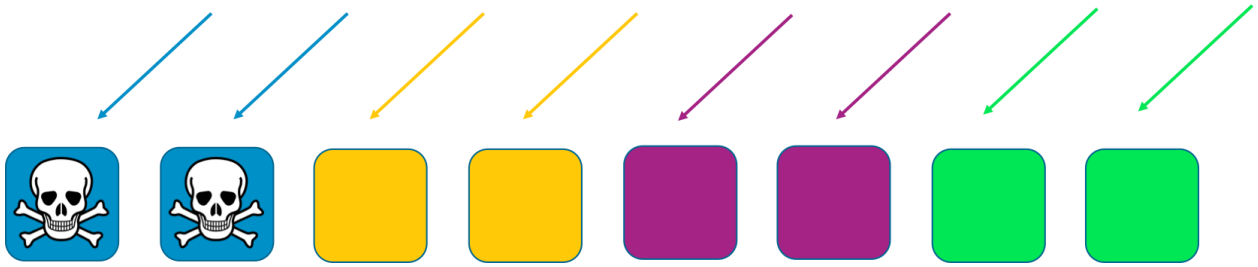
这种故障将会影响“所有人和所有事物”。整个服务会发生中断。所有客户都会受到影响。正如我们在可用性工程中所说的那样：这不是最佳选择。

通过常规分区，我们可以做得更好。如果将一个小组分为 4 个工作单元分区，虽然牺牲了一定的效率，但可以改善上述问题造成的影响。以下两个图片展示了分区限制 DDoS 攻击造成的影响的原理。



在此示例中，每个分区包含两个工作单元。我们在各个分区之间分配资源（例如客户域）。这种方式仍有冗余，但是因为每个分区只有两个工作单元，所以我们必须在系统中保留更多的未用容量用于处理任何故障。相应的，受影响的范围也会大大缩小。

使用随机分区进行工作负载隔离

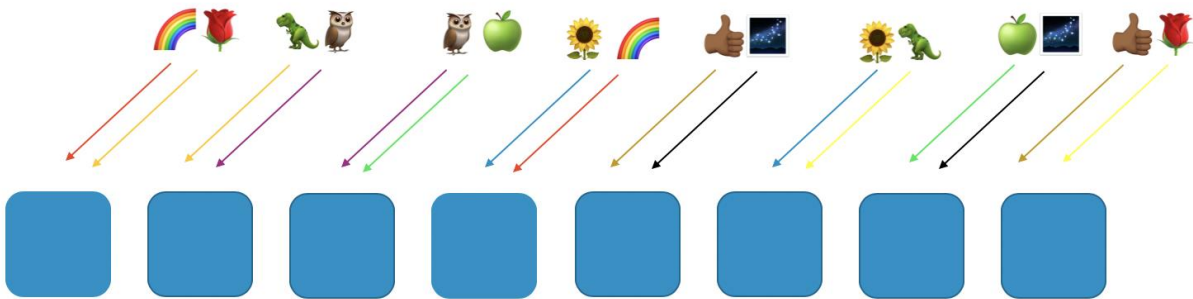


在这个分区系统中，分区数量越多，受影响的范围越小。示例中有四个分区，如果某个客户遇到问题，则该客户所在的托管分区以及该分区上的所有其他客户都可能受到影响。但是，该分区仅占整体服务的四分之一。25% 的服务受影响总比 100% 受影响要好得多。通过随机分区，我们还可以做得更好。

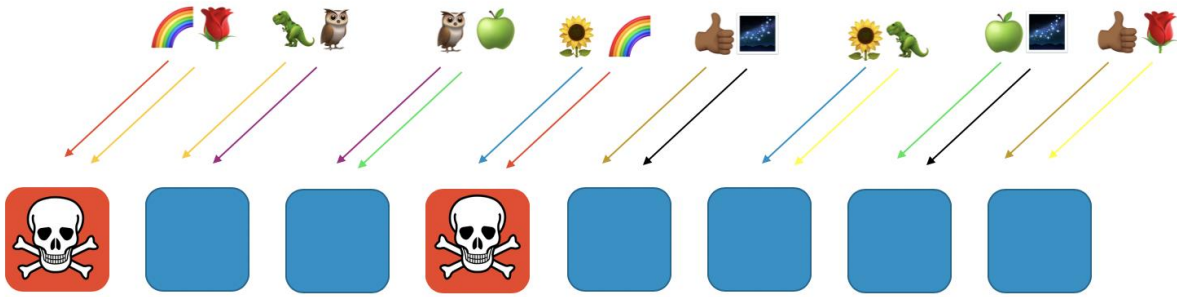
使用随机分区，我们创建了几个虚拟分区，每个虚拟分区包含两个工作单元，然后将客户、资源或任何我们想隔离的东西分配给其中一个虚拟分区。

下图展示了一个随机分区布局示例，其中有八个工作单元和八个客户，每个客户被分配给两个工作单元。通常，我们的客户数量要远多于工作单元数量。但是为了便于说明，我们在此缩小规模进行展示。我们将重点关注两个客户：“彩虹”客户和“玫瑰”客户。

在示例中，我们将“彩虹”客户分配给 1 号和 4 号工作单元。这两个工作单元共同构成了该客户的随机分区。其他客户也将被分配到由两个不同工作单元组成的虚拟分区。例如，“玫瑰”顾客被分配给 1 号和 8 号工作单元。



如果分配给 1 号和 4 号工作单元的“彩虹”客户遇到问题（例如有害请求或过量请求），则该问题会影响该虚拟分区，但不会完全影响其他任何随机分区。实际上，最多只会会有一个其他随机分区工作单元受此影响。如果请求者支持容错并且可以解决此问题（例如重试），那么就可以继续为其余分区上的客户或资源提供不间断的服务，如下图所示。



换句话说，虽然所有服务于“彩虹”的工作单元都可能遇到问题或受到攻击，但其他工作单元完全没有受到影响。对于客户而言，这意味着即使“玫瑰”客户和“向日葵”客户都与“彩虹”客户共用一个工作单元，这二者也都不会受到影响。“玫瑰”客户可以从 8 号工作单元获得服务，“向日葵”客户可以从 6 号工作单元获得服务，如下图所示。



当发生问题时，我们虽然还是会损失全部服务的四分之一，但是这种客户或资源的分配方式意味着随机分区模式下的受影响范围要小得多。如果有 8 个工作单元，就有 28 种由两个工作单元组成的唯一组合，这意味着有 28 种可能的随机分区。如果我们有数百个或更多的客户，并且将每个客户分配给一个随机分区，那么此类问题造成的影响范围就只占整体的 1/28。这比常规分区的防御效果提升了 7 倍。

好消息是，您拥有的工作单元和客户越多，这一数字还会呈指数式增长。从这些角度来说，大多数情况下扩展会愈发困难，但是随机分区却变得越来越有效。实际上，如果有足够的工作单元，随机分区的数量可能会比客户数量更多，并且每个客户都可以相互隔离。

Amazon Route 53 与随机分区

随机分区的这些优势对 Amazon Route 53 有什么好处？在 Route 53 中，我们决定将容量分配到总共 2048 个虚拟名称服务器中。这些服务器是虚拟的，因为它们与托管 Route 53 的物理服务器之间不存在对应关系。我们可以调整这些服务器来帮助管理容量。然后，我们将每个客户域分配给一个包含四个虚拟名称服务器的随机分区。按这种方式，将会有 7300 亿个可能的随机分区。随机分区的数量如此之多，让我们足以为每个域分配一个唯一的随机分区。实际上，我们可以更进一步，确保没有任何客户域与任何其他客户域共享两个以上的虚拟名称服务器。

测试结果让人惊喜。如果某个客户域成为了 DDoS 攻击的目标，则分配给该域四个虚拟名称服务器的流量将出现激增，但其他客户的域则不受影响。我们不甘让受攻击的客户蒙受损失。随机分区意味着我们可以识别受到攻击的客户，并将其隔离到特殊的专用攻击容量中。此外，我们还开发了 AWS Shield 流量清理器专有层。即使在发生这些事件时，随机分区也能够有效确保全体 Route 53 客户获得无缝体验。

结论

我们已经在我们的许多其他系统中嵌入了随机分区。此外，我们还会提出一些改进措施，例如递归随机分区，将项目分多层进行分区，从而隔离了客户的客户。随机分区适用范围广，是分配现有资源的明智之选。随机分区通常也无需额外付费，经济实惠但效果显著。

如果您想亲身体验随机分区，请查看我们的开源 [Route 53 Infima](#) 库。该库中包含可用于分配或安排资源的随机分区的几种不同实施方案。